

AN ADVANCED ARTIFICIAL INTELLIGENCE SYSTEM FOR INVESTIGATING  
THE TROPICAL CYCLONE RAPID INTENSIFICATION

by

Yijun Wei  
A Dissertation  
Submitted to the  
Graduate Faculty  
of  
George Mason University  
in Partial Fulfillment of  
The Requirements for the Degree  
of  
Doctor of Philosophy  
Computational Sciences and Informatics

Committee:

_____	Dr. Jason Kinser, Committee Chair
_____	Dr. Ruixin Yang, Committee Member
_____	Dr. Igor Griva, Committee Member
_____	Dr. Olga Gkountouna, Committee Member
_____	Dr. Jason Kinser, Department Chairperson
_____	Dr. Donna M. Fox, Associate Dean, Office of Student Affairs & Special Programs, College of Science
_____	Dr. Fernando R. Miralles-Wilhelm, Dean, College of Science
Date: _____	Fall Semester 2020 George Mason University Fairfax, VA

AN ADVANCED Artificial Intelligence System for Investigating the Tropical Cyclone  
Rapid Intensification

A Dissertation submitted in partial fulfillment of the requirements for the degree of  
Doctor of Philosophy at George Mason University

by

Yijun Wei  
Master of Art  
University of Michigan – Ann Arbor, 2013  
Master of Science  
University of Michigan – Ann Arbor, 2013

Director: Ruixin Yang, Associate Professor  
Geography and Geoinformation Science

Fall Semester 2020  
George Mason University  
Fairfax, VA

Copyright 2020 Yijun Wei  
All Rights Reserved

## **DEDICATION**

I dedicate this dissertation to my parents Wu Gao and Yong Wei, my advisor Dr. Ruixin Yang and my dear friends.

## **ACKNOWLEDGEMENTS**

I would like to thank many friends, relatives, and supporters who have made this happen. Dr. Ruixin Yang, advised me and gave me great help for conducting research. My parents, Wu Gao and Yong Wei, supported me spiritually. Dr. Kinser, Dr. Griva, and Dr. Gkountouna were of invaluable help. Finally, thanks to all the people who help me along this long journey.

# TABLE OF CONTENTS

List of Tables .....	vii
List of Figures .....	x
Abstract .....	xii
Chapter 1 Introduction .....	1
1.1 Background .....	1
1.1.1 Dynamical models .....	2
1.1.2 Statistical Model .....	3
1.1.3 Statistical-Dynamical Model. ....	5
1.2 The problem .....	8
1.2.1 SHIPS-RII model (KD03).....	8
1.2.2 Revised RII model (KDK10) .....	10
1.2.3 Enhanced RII model (KRD15) .....	11
1.2.4 Systematic machine learning and data mining models .....	12
1.3 Proposed approach .....	14
1.3.1 COR-SHIPS model .....	15
1.3.2 LLE-SHIPS model and DL-SHIPS model.....	16
Chapter 2 Data .....	19
2.1 SHIPS Developmental Data.....	19
2.2 ECMWF ERA-Interim reanalysis data .....	28
2.3 NHC best track data .....	31
Chapter 3 Data filter.....	34
3.1 SHIPS data filter .....	34
3.1.1 ASCII text to attribute-relation table .....	34
3.1.2 Preprocessing of the SHIPS data in attribute-relation table.....	37
3.1.3 Removal of highly correlated variables .....	42
3.2 ERA-Interim data filters .....	44
3.2.1 Local Linear Embedding (LLE) for filtering near core ERA-Interim data .....	45
3.2.2 Deep learning (ERA-Interim data filter for DL-SHIPS model).....	50
Chapter 4 GMM-SMOTE Sampler.....	64
Chapter 5 XGBoost classifier and Hyperparameter tuning process .....	73
5.1 XGBoost classifier .....	73
5.2 Hyperparameter tuning process .....	78
Chapter 6 Result.....	80
6.1 COR-SHIPS model .....	81
6.1.1 Hyperparameter tuning for model selection .....	81
6.1.2 COR-SHIPS result on test data.....	90
6.1.3 Feature importance.....	91
6.2 LLE-SHIPS model .....	94
6.2.1 Hyperparameter tuning for model selection .....	94

6.2.2 LLE-SHIPS result on test data .....	101
6.2.3 Feature importance.....	102
6.3 DL-SHIPS model .....	111
6.3.1 Hyperparameters tuning and result .....	111
6.3.2 Model result on test data .....	123
6.3.3 Feature importance.....	124
6.4 Model performance comparison .....	131
6.4.1 Model performance in Yang (2016) and Kaplan et al. (2015).....	132
6.4.2 Model comparison .....	134
6.5 Feature importance.....	138
6.5.1 Feature importance comparison between COR-SHIPS model, LLE-SHIPS model, and DL-SHIPS model .....	138
6.5.2 Feature importance comparison between previous studies and this study .....	142
Chapter 7 Conclusion and discussion .....	146
Appendix 1 Principal component analysis.....	150
1.1 Principal component analysis .....	150
1.2 Kernel PCA.....	152
Appendix 2 Additional Tables .....	156
References .....	185
Biography.....	202

# LIST OF TABLES

Table

	Page
Table 1.1: Candidate variables and their abbreviations included in DeMaria and Kaplan (1994).....	5
Table 1.2: Candidate variables and their abbreviations included in KD03. ....	8
Table 2.1: List of one-time variables (SHIPS 2018a), explaining the details of each variable, and the values for each corresponding time column (Adopted from SHIPS (2018c)).....	21
Table 2.2: List of the special notations (SHIPS 2018a), explaining the details of each variable (Adopted from SHIPS (2018c)). ....	25
Table 2.3: Temporal and spatial coverage of the ERA-interim pressure level data and its available pressure levels and variables. ....	29
Table 2.4: Variable names, abbreviations, units, and description for the 14 variables in the ERA-Interim pressure level dataset. ....	29
Table 2.5: A TC record in NHC best track data. ....	32
Table 3.1: One row of the attribute-relation table converted from original SHIPS data showed in Figure 2.1.....	36
Table 3.2: Variables with missing value and the missing percentage in SHIPS data. Variables are sorted according to the percentage. Variables without missing values are not listed.....	39
Table 3.3: Variables with higher than 50% single values in the SHIPS Data, the single values, and the percentages.....	40
Table 3.4: Correlation matrix for pairs among BD06, BD12, and BD18, and the highly correlated group lists leading with each variable based on a 0.8 correlation threshold....	43
Table 3.5: Hyperparameters for the LLE and their searching range defined by the Min(imum) and Max(imum).....	50
Table 4.1: Hyperparameters and their searching space in GMM-SMOTE sampling process.....	72
Table 5.1: Hyperparameters, their searching space defined by the minimum and maximum, and the initial values in GMM-SMOTE sampling process and XGBoost classifier. ....	77
Table 5.2: Confusion matrix. ....	78
Table 6.1: Kappa scores of the 5 best 10-fold cross-validation results and their means for different correlation thresholds. ....	82
Table 6.2: The number of minority and total instances, and the Imbalance Ratio (with population RI ratio at 5.1%) for the 6 clusters.....	84
Table 6.3: Top performed hyperparameter sets, the corresponding cross-validation kappa scores, and specific values of the tuned hyperparameters. ....	87



Table 6.4: The descending value ranking of individual hyperparameter among the top 5 performed cases, and the corresponding conservativeness ranking scores in parentheses. ....	88
Table 6.5: Confusion matrix values of our model after (before) hyperparameter tuning. ....	90
Table 6.6: Performance comparisons. MB and MA denote the models before and after the hyperparameter tuning. ....	91
Table 6.7: Features of top ten importance, their importance scores, and feature description from SHIPS (2018c) in the COR-SHIPS model. ....	92
Table 6.8: The performance for models with different sets of values of the hyperparameters, no_dimension and no_neighbors.....	95
Table 6.9: The number of minority and total cases, and the Imbalance Rate (with population RI ratio at 5.1%) for the 5 clusters generated by GMM. ....	97
Table 6.10: Top 5 performed hyperparameter sets, the corresponding cross-validation kappa scores, and specific values of the tuned hyperparameters.....	99
Table 6.11: The descending value ranking of individual hyperparameter among the top 5 performed cases, and the corresponding conservativeness ranking scores in parentheses. ....	99
Table 6.12: Confusion matrix values after (before) hyperparameter tuning with the test data. ....	101
Table 6.13: Performance comparisons. MB and MA denote the models before and after the hyperparameters in GMM-SMOTE and XGBoost are tuned. ....	102
Table 6.14: Variables of top ten importance, their importance scores, and feature description from SHIPS (2018c) in LLE-SHIPS model. ....	103
Table 6.15: ERA-Interim variable group with top 5 importance scores, calculated from the second step. ....	107
Table 6.16: 14 variables, their summed importance score, non-zero features extracted from each variable network, and the corresponding missing variables due to all zeros. ....	114
Table 6.17: Highly correlated variable groups .....	114
Table 6.18: Dimensions of the compressed features of auto-encoder after tuning based on the summed importance score described in Table 8 for each of the 14 variables. ....	116
Table 6.19: Number of minority, total cases, and the IIR (with population RI ratio at 5.1%) for the 3 clusters generated by GMM.....	120
Table 6.20: Top performed hyperparameter sets, the corresponding cross-validation kappa scores, and specific values of the tuned hyperparameters.....	121
Table 6.21: The descending value ranking of individual hyperparameter among the top 5 performed cases, and the corresponding conservativeness ranking scores in parentheses. ....	122
Table 6.22: Confusion matrix values after (before) hyperparameter tuning with the test data. ....	123
Table 6.23: Performance comparisons. MB and MA denote the models before and after the hyperparameters in GMM-SMOTE and XGBoost are tuned. ....	124
Table 6.24: Variable importance in DL-SHIPS model. ....	126

Table 6.25: Summed variable importance score, the number of non-zero, non-correlated features, the feature-wise averaged importance score, and its ranking for each ERA-Interim variable. ....	127
Table 6.26: Performance comparison between our models, and Y16 and KRD15. ....	135
Table 6.27: Performance comparison between 3 models developed in this study, ‘X’ in the table indicates that not available value. ....	138
Table 6.28: Top 36 most important variables in COR-SHIPS model, LLE-SHIPS model, and DL-SHIPS model. ....	139

# LIST OF FIGURES

Figure	Page
Figure 1.1: The Artificial Intelligence (AI) system structure designed in this study. ....	15
Figure 1.2: COR-SHIPS model structure.....	16
Figure 1.3: LLE-SHIPS and DL-SHIPS model structure. ....	17
Figure 2.1: An example of data block of original SHIPS ASCII text file with 141 lines. ....	27
Figure 3.1: The 33*33 grid boxes centered at the grid box consisting of the center of a TC, denoted as the black dot, and the blue area presents the near core grids.....	45
Figure 3.2: Demonstration of the convolution operation.....	54
Figure 3.3: Convolution operation for input with multiple channels.....	55
Figure 3.4: Max pooling example. ....	56
Figure 3.5: Dimension changes of the ERA data through the 3D CNN auto-encoder layers. ....	57
Figure 3.6: The deconvolution operation, reverse of operations shown in Figure 3.3. ....	60
Figure 3.7: An unpooling example. ....	61
Figure 3.8: The CNN on the right that first runs through the input to the output (from bottom to top), and the position of the Max Pooling pixel is saved as a switch that will be used later for the unpooling operation on the left .....	62
Figure 3.9: Combined deep learning filters for the 14 variables in ERA-Interim data. ...	63
Figure 5.1. CART sample. ....	74
Figure 6.1: BIC ( $10^5$ ) for GMM with different number of clusters. ....	84
Figure 6.2: Variation of Cross-validation kappa scores over Bayesian Optimization iteration numbers. ....	86
Figure 6.3: (a) Precision and POD score vs. decision threshold, and (b) Kappa score vs. decision threshold. ....	89
Figure 6.4: BIC ( $10^6$ ) for GMM with a different number of clusters in LLE-SHIPS model.....	96
Figure 6.5: Variation of Cross-validation kappa scores over Bayesian Optimization iteration numbers for LLE-SHIPS model. ....	98
Figure 6.6: (a) Precision and POD score vs. decision threshold. (b) Kappa score vs. decision threshold in LLE-SHIPS model.....	100
Figure 6.7: Network training loss over iterations for pv, z, t, q, w, vo, d, u, v, r, o3, clwc, ciwc, cc from left to right and from top to bottom.....	116
Figure 6.8: Structure for adjusted auto-encoder network. ....	118
Figure 6.9: BIC ( $10^6$ ) for GMM with different number of clusters. ....	119
Figure 6.10: Variation of Cross-validation kappa scores over Bayesian Optimization iteration numbers. ....	121
Figure 6.11: (a) Precision and POD score vs. decision threshold. (b) Kappa score vs. decision threshold .....	123

Figure 6.12: 3 channels, 64 feature maps for the first layer (dimension: 3 (channel) * 64 (feature map) * 30 (feature map dimension) * 30 (feature map dimension)) of the network that is immediate after the input layer (dimension: 37 (pressure level) * 4 (-18h, -12h, -6h, and 0h) * 33 (vertical grid) * 33 (horizontal grid)) for variable q with its network structure displayed in Figure 6.8a. ....	128
Figure 6.13: Same as Figure 6.12 but for a RI instance (a) RI in channel 1. (b) RI in channel 2. (c) RI in channel 3. ....	129
Figure 6.14: Same as Figure 6.12 but for variable vo with its network structure in Figure 6.8a in a non-RI instance: (a) non-RI in channel 1. (b) non-RI in channel 2. (c) non-RI in channel 3. ....	130
Figure 6.15: Same as Figure 6.12 but for variable vo with its network structure in Figure 6.8a in 3 channels in a RI instance: (a) RI in channel 1. (b) RI in channel 2. (c) RI in channel 3. ....	130
Figure 6.16: Same as Figure 6.12 but for variable u with its network structure in Figure 6.8a in a non-RI instance: (a) non-RI in channel 1. (b) non-RI in channel 2. (c) non-RI in channel 3. ....	131
Figure 6.17: Same as Figure 6.12 but for variable u with its network structure in Figure 6.8a in 3 channels in a RI instance: (a) RI in channel 1. (b) RI in channel 2. (c) RI in channel 3. ....	131
Figure 6.18: Kappa, POD, and FAR for (a) C4.5 decision tree. (b) ADTree. Data are from Y16. ....	133
Figure 6.19: Different model's performance regarding Peirce's skill score (PSS) based on data from KRD15. ....	134
Figure 6.20: Model performance comparison: Model's test kappa, FAR, and POD score in the best model in Yang (2016), SHIPS model, LLE-SHIPS model, and DL-SHIPS model. ....	135
Figure 6.21: Model's test PSS, FAR, and POD score in KRD15 (Kaplan et al., 2015), COR-SHIPS model, LLE-SHIPS model, and DL-SHIPS model. ....	137

# ABSTRACT

## AN ADVANCED ARTIFICIAL INTELLIGENCE SYSTEM FOR INVESTIGATING THE TROPICAL CYCLONE RAPID INTENSIFICATION

Yijun Wei, Ph.D.

George Mason University, 2020

Dissertation Director: Dr. Ruixin Yang

Tropical cyclones (TCs) can cause heavy casualties due to storm surge, high wind gusts, heavy rainfall and flooding, and landslides, so predicting TC is important. There are mainly two elements of TC forecasting: tracking prediction and intensity prediction. So far, it is found that tracking prediction is more mature than the intensity prediction. Various models are developed for TC intensity prediction and can be simple enough to run for a few seconds or complex enough to run for a couple of hours on a supercomputer. Although with so many models are developed, the intensity prediction accuracy is still very low, and one primary reason is the existence of Rapid Intensification (RI).

Currently, most RI prediction studies are conducted based on a subset of the SHIPS database using a relatively simple model structure. However, variables (features) in the SHIPS database are built upon expert knowledge in TC intensity studies, and the variable values are derived from gridded model outputs or satellite observations. Are there any more

21 important variables in TC intensity predictions but not identified in the SHIPS dataset? In  
22 this study, two AI-based techniques are used to extract new features from a widely used  
23 comprehensive gridded reanalysis data set. The original SHIPS data, and the newly derived  
24 features are used as inputs to an artificial intelligence (AI) for the RI prediction.

25 This study first constructs a complicated artificial intelligence (AI) system, the COR-  
26 SHIPS model, based on the complete SHIPS dataset that handles feature engineering and  
27 selection, imbalance, prediction, and hyper parameter-tuning simultaneously. The COR-  
28 SHIPS model is derived to improve the performance of the current researches in RI  
29 prediction and to identify other essential SHIPS variables that are ignored by previous  
30 studies with variable importance scores. COR-SHIPS is also used as the baseline model in  
31 the dissertation.

32 To distill new variables from vast amounts of gridded data, two models, with a similar  
33 structure to the COR-SHIPS model but with an additional data filters, are designed in the  
34 dissertation to identify new features related to TC intensity changes in general and RI in  
35 particular. Here, we adopt the Local linear embedding (LLE) and deep learning (DL)  
36 techniques respectively to filter the near center and large-scale spatial data of the European  
37 Centre for Medium-Range Weather Forecasts (ECMWF) ERA-Interim reanalysis data, one  
38 of the best reanalysis products at the moment, for identifying new variables related to RI,  
39 and term the corresponding LLE-SHIPS model and DL-SHIPS model, respectively.

40 The result of the three models outperforms most of the earlier studies by at least  
41 approximately 30%, 60%, and 75%, respectively. In addition to the well-known SHIPS  
42 database, we specify the 400 and 450 hPa wind speeds, identify 1000 hPa potential vorticity

43 and vertical pressure speed, and differentiate humidity southeast, vorticity north, and  
44 eastward wind north to the TC centers that could help the prediction and understanding the  
45 occurrence of RI.  
46  
47

# CHAPTER 1 INTRODUCTION

## 1.1 Background

Tropical cyclones (TCs) can cause heavy casualties due to storm surge, high wind gusts, heavy rainfall and flooding, and landslides (Pacific Disaster Center. n.d.). On May 2, 2008, Cyclone Nargis sent a storm surge in Myanmar and killed at least 138,000 people (Enz et al. 2009). Prediction of the behavior of TCs can minimize deaths and losses. Therefore, skillful TC prediction is significant to mitigate risk by timely planning and preparation.

The first known TC forecast was conducted by Lt. Col. William Reid of the Corps of Royal Engineers in the western hemisphere in 1846, and barometric pressure was used as the basis for Reid's approach (Reid 1846). Most forecasts before 1900 were obtained by direct observation at weather stations through the telegraph. Significant changes were made in data collection since 1900, where radiosondes (1930), aircraft (1943), coast weather radar (1950), and weather satellite (1960) were introduced (Sheets 1990).

There are mainly two elements of TC forecasting: track forecasting and intensity forecasting. These two predictions are critical in disaster prevention, but the development of these two presents a difference. So far, it is found that tracking prediction is more mature than intensity prediction. DeMaria et al. (2007) examined the National Hurricane Center (NHC) and Joint Typhoon Warning Center operational TC intensity forecasts for the three major northern hemisphere TC basins (Atlantic, eastern North Pacific, and



68 western North Pacific) for the past two decades. The intensity forecasts were compared to  
69 the track forecasts for the same data sample. The performance of the two forecasts was  
70 comparable at 12 h, but the track forecasts were 2 to 5 times more skillful by 72 h, with  
71 the largest ratio in the western Pacific. As lead time increases, tracking prediction became  
72 more skillful than intensity prediction. Cangialosi and Franklin (2017) indicated that in  
73 the Atlantic and Pacific, the skill of track prediction is at least 3-7 times larger than that  
74 of intensity in 12, 24, 36 hours, and 10-40 times more skillful by 48, 72, 96, and 120  
75 hours in 2016. Since the tracking forecasting is relatively accurate, and the intensity  
76 forecasting is with low skills, recent research on TC prediction mainly focuses on  
77 intensity forecasting on a time scale from 12 hours to 120 hours (Cangialosi and Franklin  
78 2017).

79 Various models were developed for TC intensity prediction and can be simple enough  
80 to run for a few seconds or complex enough to run for a couple of hours on a  
81 supercomputer. Based on the mechanism, these models can be characterized as the  
82 dynamical model, the statistical model, and the statistical-dynamical model.

### 83 **1.1.1 Dynamical models**

84 Dynamical models, also known as the numerical models, consider complex  
85 physical processes and are used on the supercomputer to solve the ordinary and partial  
86 differential equations in physics. One of the most critical and influential models is the  
87 Geophysical Fluid Dynamics Laboratory (GFDL) model (Kurihara et al. 1998), which  
88 was used for a research purpose during 1973 and 1980. Encouraged by the performance  
89 of the GFDL model, the research model (GFDL) was converted to a comprehensive

90 prediction system that started in the mid-1980s. The process took about 15 years, and  
91 GFDL model became operational in 1995. The interpolated form of GFDL model (GFDI)  
92 became available for intensity in 1996, and a U.S. Navy's version of GFDI was added in  
93 1999 (Rennick 1999). Another widely used dynamical model is the hurricane weather  
94 research and forecast (H-WRF) model, which became operational in 2007 (Miller 2007).

### 95 **1.1.2 Statistical Model**

96 A statistical model does not include the physics of the atmosphere but instead is  
97 based on the relationship between specific information and behavior of TCs. The first  
98 statistical model was developed in 1972 to help generate TC track forecasts, and the  
99 model was named Climatology and Persistence (CLIPER5). In 1979, the Statistical  
100 Hurricane Intensity Forecast (SHIFOR) model, which consisted of climatology and  
101 persistence variables, began operational for TC intensity prediction (Jarvinen and  
102 Neumann 1979). A 5-day SHIFOR version (SHIFOR5) was implemented in 2001 (Knaff  
103 et al. 2003). Decay-SHIFOR5 is a form of SHIFOR5 that includes a weakening  
104 component when TCs move inland, and Decay-SHIFOR5 modifies intensity over land  
105 using CLIPER track and climatological decay rate (Rhome 2007).

106 A well-known statistical model is the Statistical Hurricane Intensity Prediction  
107 Scheme (SHIPS) developed by DeMaria and Kaplan (1994). In SHIPS, different multiple  
108 regression models with persistence, synoptic, and climatological variables were derived  
109 to predict TC intensity changes in 12, 24, 36, 48, 72 hours for the Atlantic basin. The  
110 storm intensity, i.e., the dependent variable, is measured by the maximum 1-min  
111 sustained surface wind, and an independent variable (predictor) is considered significant

112 if the probability that the regression coefficient is different from 0 with exceeds 95%  
113 confidence.

114       The candidate variables were displayed in Table 1.1. Among those, JDATE,  
115 VMX, DVMX, LAT, LONG, USM, VSM, and CSM were climatological and persistence  
116 variables, and POT, SHR, DSHR, REFC, PEFC, SIZE, and DTL were synoptic variables.  
117 POT took the effect of sea surface temperature (SST) into account since SST is closely  
118 related to TC intensification (Merrill 1987). SHR and DSHR were used to evaluate the  
119 vertical shear, and plenty of studies have shown vertical shear of the horizontal wind has  
120 a negative influence on TC intensification (e.g., Gray 1967; Merrill 1988). REFC and  
121 PEFC are included to account for positive interactions between the TC and synoptic-scale  
122 systems. SIZE was included as a measure of the extent of the outer circulation of the TC.  
123 Although all landfall cases were eliminated from the data, the proximity to land might  
124 still have a modifying influence on the storm intensity, and DTL was involved in the  
125 model. A simple backward-stepping procedure was conducted for the variable selection  
126 and POT, SHR, DVMX, REFC, PRFC, JDATE, LONG, DTL, SIZE, and DSHR were  
127 selected as the variables for the multiple regression model. The model was tested using  
128 the Jackknife procedure, and the result indicated that the intensity errors were 10% - 15%  
129 smaller than the errors from a model that used only climatology and persistence  
130 (SHIFOR5). However, the forecast only explained about 50% of the variability of the  
131 observed intensity change, which indicates that a statistical model with large-scale  
132 variables is not able to explain all types of storm process effects adequately (DeMaria and  
133 Kaplan 1994).

134

135 **Table 1.1: Candidate variables and their abbreviations included in DeMaria and**  
 136 **Kaplan (1994).**

Variable	Abbreviation
Absolute value of Julian date – 253	JDATE
Initial storm intensity	VMX
Intensity change during previous 12 h	DVMX
Initial storm latitude	LAT
Initial storm longitude	LONG
Eastward component of storm motion vector	USM
Northward component of storm motion vector	VSM
Magnitude of storm motion vector	CSM
Maximum possible intensity - initial intensity	POT
Magnitude of 850-200-mb vertical shear	SHR
Time tendency of vertical shear magnitude	DSHR
The 200-mb relative eddy angular momentum flux convergence	REFC
The 200-mb planetary eddy angular momentum flux convergence	PEFC
The 850-mb relative angular momentum	SIZE
Distance to nearest major landmass	DTL

137

### 138 **1.1.3 Statistical-Dynamical Model.**

139 Statistical-dynamical model blends the statistical model and the dynamical model  
 140 (NHC Track and Intensity Models 2017). In other words, the statistical-dynamical model  
 141 employs variables derived from the dynamical models. Although SHIPS proposed in  
 142 DeMaria and Kaplan (1994) was regarded as a statistical or statistical-synoptic model. In  
 143 1997, SHIPS was converted to a statistical-dynamical model by using large-scale

144 variables in Global Forecast System (GFS) (DeMaria et al. 2005). Therefore, the SHIPS  
145 after 1997 is regarded as a statistical-dynamical model.

146 DeMaria et al. (2005) described modifications to the NHC operational SHIPS  
147 intensity model from 1997 to 2003, including an additional method to account for the use  
148 of variables from the dynamical model in 1997, the storm decay over land in 2000, the  
149 extension of the forecasts from 3 to 5 days in 2001, and the use of the GFS, a global  
150 numerical computer model run by National Oceanic and Atmospheric Administration  
151 (NOAA) in 2001. The study showed that SHIPS performs well in predicting 72 h  
152 intensity in the Atlantic, and at 48 and 72 h in the east Pacific. The inclusion of the  
153 effects of the decay over land beginning in 2000 reduces the short period Atlantic  
154 intensity error but not for 72 h forecasting. An experimental version of SHIPS consisted  
155 of satellite variables during the 2002 and 2003 seasons significantly improved skill in the  
156 east Pacific forecasts by up to 7% at 12–72 h, and 3.5% through 72 h in Atlantic forecasts  
157 (DeMaria et al. 2005).

158 The latest version of SHIPS, which has an inland decay component, was known  
159 as Decay-SHIPS (DSHIPS). The DSHIPS typically provides more accurate TC intensity  
160 forecasts when TCs encounter or interact with the land. Over open water with no land  
161 interactions, the intensity forecasts from DSHIPS and SHIPS are identical (Rhome 2007).

162 As SHIPS model is mainly used in Atlantic and northeast Pacific, a similar model  
163 known as the Statistical Typhoon Intensity Prediction Scheme (STIPS) was developed for  
164 the northwest Pacific Ocean and Southern Hemisphere by Knaff et al. (2005).

165           The values of SHIPS variables, available openly, are considered as the SHIPS  
166 database. Every year, instances from the previous year and new variables from other  
167 sources may be added to the database, while some old variables may be removed. The  
168 development of the SHIPS database was described in DeMaria and Kaplan (1994, 1999),  
169 DeMaria et al. (2005), and Kaplan et al. (2010, hereafter, KDK10). The most recent  
170 version of the database is the SHIPS Developmental Data, a complete dataset with known  
171 different types of the parameter related to TC intensity changes (SHIPS 2018a). The  
172 SHIPS database was used by many TC intensity related types of research. One such  
173 example is the logistic growth equation model (LGEM) (DeMaria 2009), which was also  
174 a type of statistical-dynamical TC intensity model and used the same input as SHIPS but  
175 in the framework of a simplified dynamical prediction system. LGEM estimated the only  
176 parameter in LGEM - population growth rate, which is proportional to the maximum  
177 sustained wind, using four free parameters. These four parameters were the time-  
178 dependent growth rate, maximum potential intensity (MPI), and two constants that  
179 determine how quickly the intensity relaxes toward the MPI, i.e., vertical shear (S) and a  
180 convective instability parameter (C). LGEM was found to explain observed intensity  
181 variations better than SHIPS (DeMaria and Kaplan 1994). LGEM-MR, a version of  
182 LGEM, where the remaining parameters are determined by a multiple regression method  
183 using a subset of the SHIPS database, came to work in real-time from 2006. The average  
184 skill of LGEM-MR forecasts is up to 17% better than those from the SHIPS model  
185 (DeMaria 2009).

186 Since a large number of intensity forecasting models came into available, the  
187 consensus prediction model (ICON) was developed by Sampson et al. (2008). A three  
188 models consensus of DSHIPS, GFDL, and GFNI (The interpolated Navy version of  
189 GFDL hurricane model) were found outperform almost all the single intensity forecasting  
190 model in the Atlantic basin.

## 191 **1.2 The problem**

192 Although the statistical-dynamical models have been used since early 1990, the  
193 prediction accuracy is still not high. One primary reason is the existence of Rapid  
194 Intensification (RI) (Kaplan and DeMaria 2003; DeMaria et al. 2005; Yang et al. 2007).

### 195 **1.2.1 SHIPS-RII model (KD03)**

196 RI was defined in Kaplan and DeMaria (2003, hereafter, KD03), as the maximum  
197 sustained surface wind speed increase of 30 kt or more over a 24-h period, and KD03  
198 derived the initial version of the Statistical Hurricane Intensity Prediction Scheme Rapid  
199 Intensification Index (SHIPS-RII) model for RI prediction for the Atlantic basin. A two-  
200 sided t-test was utilized in KD03 to determine if the 16 different variables, listed in Table  
201 1.2, display significant differences in RI instances and non-RI instances.

202

203 **Table 1.2: Candidate variables and their abbreviations included in KD03.**

Variable	Abbreviation
Maximum sustained surface wind speed	VMAX
Latitude	LAT
Longitude	LON
Storm speed	SPD
Intensity change during the previous 12 h	DVMX
Storm motion	USTM
The absolute value of (Julian date - 253)	JDAY

Sea surface temperature	SST
Maximum possible intensity - initial intensity	POT
850–200-hPa vertical shear	SHR
200-hPa wind	U200
200-hPa temperature	T200
850–700-hPa relative humidity	RHLO
850-hPa relative vorticity	Z850
200-hPa relative eddy angular momentum flux	REFC
Steering layer	SLYR

204

205 The result indicated that

- 206       • 11 of 16 variables except for VMAX, SPD, JDAY, T200, and Z850 show
- 207               significantly difference at the 95% significance level;
- 208       • And among them, 10 variables except for LON show significantly difference at
- 209               the 99% significance level;
- 210       • 7 of them, i.e., LAT, DMAX, SST, RHLO, POT, SHR, and U200, show
- 211               significantly difference at the 99.9% significance level

212 in RI instances and non-RI instances.

213       Variables that are significant at 95%, 99%, and 99.9% confidence level were used to

214 evaluate the composite probability of RI, respectively, and variables at 99.9% level were

215 found to have the highest probability. However, LAT and U200 were further removed

216 because they were found to be highly correlated to SHR, and SST, and POT, respectively,

217 and highly correlated variables do not give us much additional information. Therefore,

218 DVMX, SHR, SST, POT, and RHLO were remained to achieve the highest composite

219 probability of RI.



220 KD03 also found that RI instances tended to occur farther south and west than the  
221 non-RI instances. In addition, the RI instances had a more westerly component of motion  
222 and were intensifying more during the preceding 12h than the non-RI instances.  
223 Furthermore, the RI instances appeared farther from their maximum potential intensity  
224 and in regions of warmer water, higher lower-tropospheric relative humidity, lower  
225 vertical shear, and more easterly upper-tropospheric flow than the non-RI instances.  
226 Interestingly, RI was more likely to occur for systems that are in an environment where  
227 forcing from upper-level troughs or cold lows was weaker than the average of all.

#### 228 **1.2.2 Revised RII model (KDK10)**

229 To employ a more sophisticated statistical method, compared to KD03, KDK10 used  
230 four more variables, one large-scale variable, and three satellite-derived variables,

- 231 • 200-hPa divergence from the 0–1000-km radius (D200),
- 232 • Percent area from 50 to 200 km covered by  $\leq -30^{\circ}\text{C}$  infrared (IR) imagery  
233 cloud-top brightness temperatures (PX30),
- 234 • Standard deviation of 50–200-km IR cloud-top brightness temperatures (SDBT),  
235 and
- 236 • Ocean heat content (OHC),

237 to conduct a linear discriminant analysis for RI prediction both in the Atlantic and in the  
238 eastern North Pacific based on TCs happened during 1995 and 2006. KDK10 evaluated  
239 the performance of the model in terms of probability of detection (POD) and false alarm  
240 ratios (FAR), and the model was found better than any other operational RI prediction  
241 models at that time. Meanwhile, D200, SHRD, and the PER were found to be the most

important variables in the Atlantic basin RI prediction. In contrast, PER, SDBT, and POT were found to be the most important variables in the eastern North Pacific basin.

### **1.2.3 Enhanced RII model (KRD15)**

To improve the usefulness of the revised RII model, Kaplan et al. (2015, hereafter, KRD15) reevaluated the variables for RI with 20-55 knots intensity changes in 12 to 48 hours (seven combinations) for both the Atlantic basin and the eastern North Pacific and selected ten variables (replaced two and added two in comparison to those in KDK10). They then used the linear discriminant analysis technique to develop an enhanced SHIPS-RII. The enhanced SHIPS-RII model, along with the logistic regression model and the Bayesian classification models by Rozoff and Kossin (2011), were fed into a probabilistic model, and resulted in a better RI prediction.

Although KD03, KDK10, and KRD15 achieved certain prediction skills for the RI prediction, the test method used in those studies for variable selection is a one-by-one t-test, which is a trial-and-error process on individual factors. There are possibilities that a single variable may be insignificantly correlated to response, but multiple variables together may have a significant impact on the response's prediction skill (Trevor et al. 2009). Furthermore, only a few variables (usually less than 20) are selected for these studies, and many useful variables may be neglected. Therefore, more systematic methods are needed to conduct an exhaustive search for the most influential factors contributing to RI in a given set of factors. Efforts were made by Yang et al. (2008), and Yang et al. (2011), which employed the association rule for feature selection among the variables identified by KD03.

#### 264    **1.2.4 Systematic machine learning and data mining models**

265           Association rule is an unsupervised and automatic data exploration method to  
266   explore multiple-to-one associations for discovering interesting relationships hidden in  
267   large databases (Yang et al. 2007). The strength of association rule can be measured  
268   concerning its support and confidence. Support determines how often a rule applies to a  
269   given data set, while confidence determines how frequently the rule happens (Tan 2015).  
270   Yang et al. (2007) adopted association rules with the 11 independent variables being  
271   discretized into two value ranges (High-Low) for each of the variables to predict RI. A  
272   three variables association rule (47.6% confidence, and 1.3% support) mined out has a  
273   higher RI probability than that with five variables (41% confidence, 0.7% support)  
274   identified by KD03. Yang et al. (2011) used association rule with more variables from  
275   KD03 database for the period 1997-2003. The result showed that the association rule  
276   reaches the support of 5.5% with an accuracy of no less than 70%. However, there are a  
277   large number of RI instances that do not follow the rule; a more generalized approach  
278   should be used.

279           Furthermore, Yang (2016, hereafter, Y16) employed WEKA (Holmes et al. 1994),  
280   a machine learning toolbox, to conduct an exhaustive and systematic examination for  
281   classification-based RI prediction with various models, subset features, and cost values  
282   for imbalance handling. Y16 split the entire dataset into a training dataset for model  
283   fitting and a test dataset for model evaluation. Although the performance of the best  
284   model in Y16 achieved a decent training result, the performance on test data was not as  
285   good. Apparently, the commonly known overfitting caused an accuracy discrepancy.

286           One way to improve the performance of Y16 is to tweak the so-called  
287 hyperparameters, the set of model parameters that do not change over the training process  
288 because in Y16, only the default hyperparameter setting is used. The other way is to  
289 improve the cost-effective approach used in Y16 that handles the highly imbalanced RI  
290 and non-RI instances.

291           So far, most RI prediction studies, including those introduced above, are  
292 conducted based on SHIPS database. However, variables in the SHIPS database are built  
293 upon human expertise in defining a relevant event based on hard and subjective  
294 thresholds. There are possibilities that those expert engineered variables in SHIPS may  
295 not be comprehensive enough, or some useful information may be ignored by the experts  
296 since the mechanism of TC intensification, and RI process is still not fully understood.

297           Therefore, other data sources should be employed in addition to SHIPS data to  
298 enhance the performance of the model. As one of the best reanalysis products at the  
299 moment, European Centre for Medium-Range Weather Forecasts (ECMWF) ERA-  
300 Interim reanalysis data can be a candidate. A large number of researches regarding TCs  
301 are conducted based on ECMWF ERA-Interim reanalysis data.

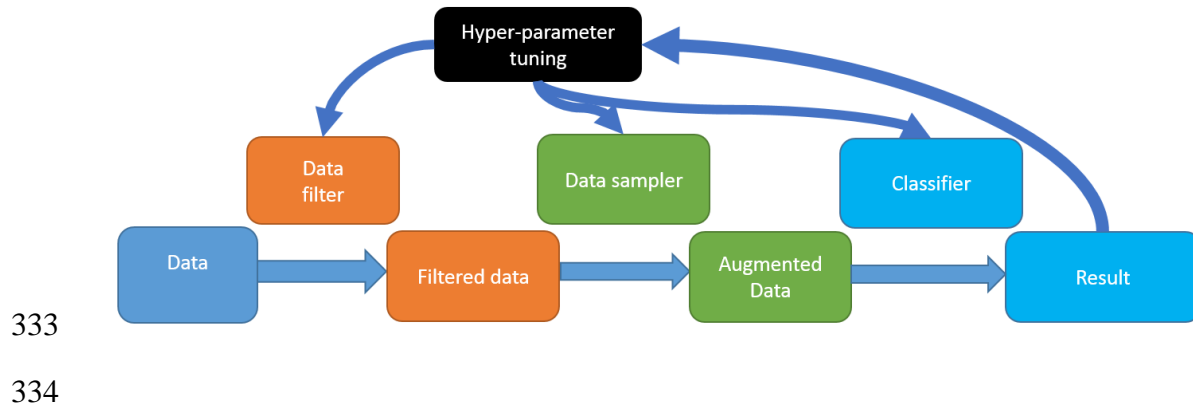
302           For example, in Wang et al. (2015), the relationship between the vertical wind  
303 shear (VWS) and the intensity change is analyzed based on ECMWF ERA-Interim  
304 reanalysis data. VWS was found negatively correlated with the intensity change, and  
305 furthermore, instead of commonly used shear between 200 and 850 hPa, the shear  
306 between 300 and 1000 hPa displays a higher negative correlation with the TC intensity  
307 change. Wang et al. (2015) also indicated that the probability for TC intensifies the

intensity and suffers RI increases when the VWS is falling and sea surface temperature (SST) is increasing. Other researches with the ERA-Interim data on TC related topics include Qian et al. (2016), Wang (2018), Astier et al. (2015), and Ferrara et al. (2017).

### **1.3 Proposed approach**

To improve the performance of previous RI prediction researches and to identify the new additional essential variables from the prediction, this study constructs a well-tailored artificial intelligence (AI) system that uses a data filter to process the input data into attribute-relation format, adopts a customized data sampler for overcoming the imbalance, employs a very powerful state-of-the-art classifier, and tweaks the hyperparameters for optimal results. The structure of the proposed AI system is displayed in Figure 1.1. The input data could be a single data source, as well as multiple data sources, and for a single data source input, the input data will be processed to attribute-relation format by the data filter to be fed into the data sampler. If there are multiple input data sources, each data source will be processed into a separate attribute-relation table and are concatenated with each other before feeding into the data sampler. The data sampler will upsample RI (minority) instances and downsample non-RI (majority) instances accordingly, leading to a balanced augmented data set. The classifier will then classify the balanced data into RI or non-RI instances. Although the hyperparameter tuning is displayed as one component in Figure 1.1, the process could take place in multiple steps, either independently for one component, or in several components for the whole AI system.

329 Based on the AI system in Figure 1.1, three models, the COR-SHIPS model, the  
 330 LLE-SHIPS model, and the DL-SHIPS model, are developed, and some details are  
 331 elaborated below.  
 332



335 **Figure 1.1: The Artificial Intelligence (AI) system structure designed in this study.**  
 336 **One data filter is displayed in the Figure to process one data source, but if there are**  
 337 **multiple input data sources, multiple data filters will be used, and each input data**  
 338 **will be processed separately from each data filter into a separate attribute-relation**  
 339 **table. All of the data filters' output is concatenated together before feeding into the**  
 340 **data sampler.**

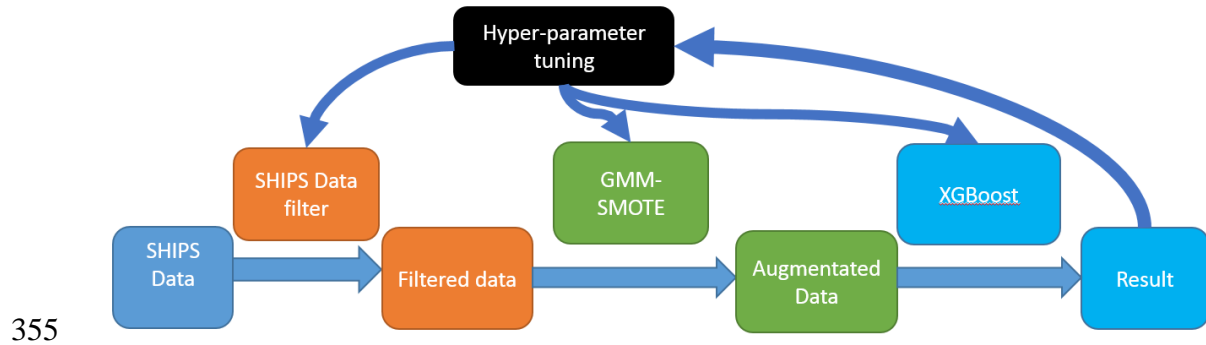
### 341 1.3.1 COR-SHIPS model

342

343 COR-SHIPS model employs only SHIPS developmental data and is the continued  
 344 work of Y16. Comparing with Y16, the COR-SHIPS model adopts a different data filter,  
 345 upsamples RI instances, employs a more powerful classifier, and tunes their  
 346 hyperparameters to improve the performance. Figure 1.2 displays the structure of the  
 347 COR-SHIPS model, which consists of four components, SHIPS data filter, GMM-  
 348 SMOTE data sampler, XGBoost classifier, and hyperparameter tuning component. The

SHIPS data filter is first used to convert the SHIPS developmental case-based data blocks to commonly used attribute-relation format and to filter the variables to generate a reduced variable set as the input for the data sampler. The GMM-SMOTE data sampler, XGBoost classifier, and hyperparameter tuning component are working the same as the corresponding component in Figure 1.1.

354



**Figure 1.2: COR-SHIPS model structure.**

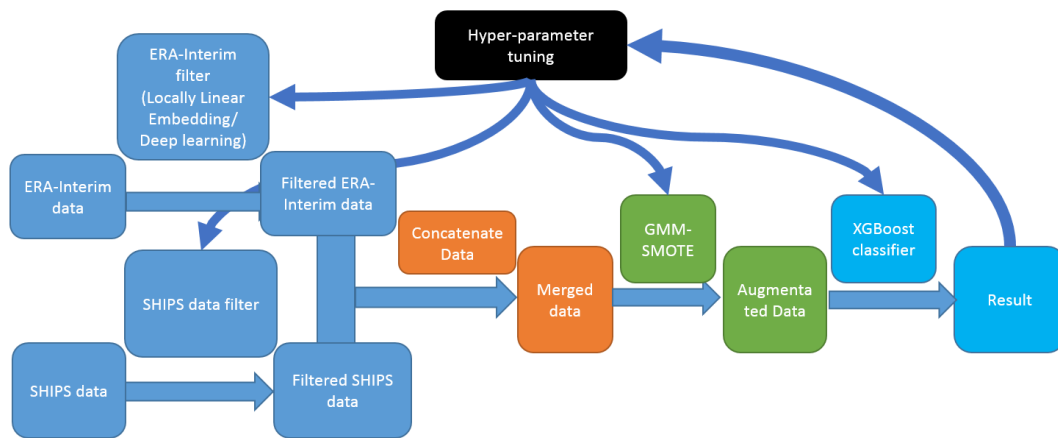
### 1.3.2 LLE-SHIPS model and DL-SHIPS model

However, similar to Y16, the COR-SHIPS model still only employs the SHIPS dataset, which is largely based on expert experiences. Since the mechanism of the TC is unknown, the knowledge from the domain scientist may not be comprehensive, which may result in some important variables not being included in the SHIPS dataset. Therefore, an additional input data source, comprehensive ECMWF ERA-Interim reanalysis data, is used, and different ERA-Interim data filters are used to process the data into the attribute-relation format. Two models, the LLE-SHIPS model, and the DL-SHIPS model, are proposed based on the structure that is described in Figure 1.3 with

almost the same structure as that of the COR-SHIPS model in Figure 1.2 except two differences.

Firstly, in addition to using the SHIPS data filter, LLE-SHIPS model and DL-SHIPS model employ other ERA-Interim data filters to process ERA-Interim data. The ERA-Interim data filter in the LLE-SHIPS model is used to only value near the TC center (near center information) that is derived from the reanalysis data. Due to the limitation on LLE implementation, the DL-SHIPS model will be used to derive additional large-scale variables, up to 1,200 km from the TC centers.

Secondly, the hyperparameter tuning component of the LLE-SHIPS model and DL-SHIPS model tune of ERA-Interim data filter independently from the other components, which are similar to the COR-SHIPS model.



**Figure 1.3: LLE-SHIPS and DL-SHIPS model structure.**



382           In sum, this study unprecedently designs an AI system that automates data  
383 filtering, data augmentation, classification, and hyperparameter tuning for TC intensity  
384 prediction to improve the RI prediction performance and identify new variables that are  
385 critical in the prediction process. This is also one of the few attempts to explore how the  
386 state of art machine learning models perform in TC RI prediction.

387           The outline of the remainder of this dissertation is constructed as follows. The  
388 datasets, including but not limit to ERA-Interim and SHIPS for this study, are introduced  
389 in Chapter 2. Chapter 3 specifies data filters, and Chapter 4 describes the data sampler.  
390 The classifier and hyperparameter tuning components are discussed in Chapter 5. Chapter  
391 6 delivers the result and discusses variable importance. Chapter 7 concludes the study and  
392 discusses future research.

## CHAPTER 2 DATA

The variables used in this study originated from three datasets for the 1982 to 2017 period. The first data is SHIPS Developmental Data, the most complete dataset with known different types of the parameters related to TC intensity changes (SHIPS 2018a). Another data set is the ECMWF ERA-Interim pressure level reanalysis data, “a reanalysis of the global atmosphere covering the data-rich period since 1979 available every 6 hours over 37 different vertical pressure levels” (Berrisford et al. 2011). The last is the National Hurricane Center (NHC) best track data, which has “a comma-delimited, text format with six-hourly information on the location, maximum winds, central pressure, and (beginning in 2004) size of all known TCs and subtropical cyclones” (Landsea and Franklin 2013). To prepare for further analysis in this study, National Hurricane Center (NHC) best track data is first used to locate the center of the TC, and related spatial subsets around TC centers is cropped from ECMWF ERA-Interim pressure level reanalysis data. Then ERA-Interim variables are processed through the ERA-Interim data filter and then concatenated with the SHIPS variables processed from the SHIPS data filter.

### **2.1 SHIPS Developmental Data**

At the moment, SHIPS Developmental Data (SHIPS 2018a), collected in ASCII text file (SHIPS 2018b), is the most complete dataset with known different types of the parameters related to TC intensity changes. Every year, instances from the previous year, and possibly new variables from other sources are added to the SHIPS Developmental Data while some old variables may be removed. The 2018 version of the SHIPS

414 Developmental Data used in this study had TC instances from 1982 to 2017 in the  
415 Atlantic basin.

416 SHIPS data consists of synoptic, climatological, persistence, geographical, satellite,  
417 experimental variables. Some of the variable examples are horizontal wind difference  
418 between 850 and 200 hPa (SHRD)<sup>1</sup>, relative humidity between 850 and 700 hPa  
419 (RHLO), the maximum potential intensity from Kerry Emanuel equation (VMPI),  
420 previous 12-hour change intensity (BD12), the 200 hPa zonal wind (U200), the 200 hPa  
421 zonal wind with radius 0-200 km (U20C), and Reynolds SST (RSST) (DeMaria and  
422 Kaplan 1994; DeMaria and Kaplan 1999; DeMaria et al. 2005).

423 The values of these variables are derived from multiple types of data sources and are  
424 accumulated based on multiyear data processing. It is almost impossible to describe all of  
425 them in detail, so we describe some data source examples below (Readers are referred to  
426 the SHIPS data description file for more information (SHIPS 2018c)).

427 TPW (Total Precipitable Water), which is the volume of water vapor in a column that  
428 from the earth surface to the atmosphere, is a meteorological parameter used for heavy  
429 precipitation prediction. TPW is created by two satellite instruments, AMSU (Advanced  
430 Microwave Sounding Unit) on three NOAA satellites, and SSM/I (Special Sensor  
431 Microwave Imager) on three DMSP (Defense Meteorological Satellite Program) satellites  
432 (Kidder and Jones 2007), using a blending algorithm. In this study, 40 variables are TPW  
433 variables in different spatial scale related to TC center, including 21 variables in MTPW

---

<sup>1</sup> Original abbreviations used in SHIPS are used here. Readers are referred to SHIPS documentation (SHIPS 2018c) for details.

(explained in Table 2.1 and also will be discussed below), and 19 variables, PW01 to PW19, in different temporal scale of MTPW.

Infrared (IR) imagery is produced by sensing the electromagnetic radiations emitted or reflected from a target surface. IR imagery obtained from GEOS east and GEOS west has high temporal and spatial resolution and also is well known for its correlation with TC rapid intensification (Knaff et al. 2008). Variables in SHIPS created from GEOS infrared (IR) imagery are: IR00, IRXX, IRM3, and IRM1 are in different spatial and temporal scales. More details are explained in Table 2.1.

**Table 2.1: List of one-time variables (SHIPS 2018a), explaining the details of each variable, and the values for each corresponding time column (Adopted from SHIPS (2018c)).**

Variable Name	Variable description
HIST	Storm history variable. The number of 6 hour periods the storm max wind has been above 20, 25, ..., 120 kt
IR00	Variables from GOES data (not time dependent). The 20 values in this record are as follows: TIME = 0: Time (hhmm) of the GOES image, relative to this case TIME = 6: Average GOES ch 4 brightness temp (deg C *10), r=0-200 km TIME = 12: Standard deviation of GOES BT (deg C*10), r=0-200 km TIME = 18: Same as 2) for r=100-300 km TIME = 24: Same as 3) for r=100-300 km TIME = 30: Percent area r=50-200 km of GOES ch 4 BT < -10 C TIME = 36: Same as 6 for BT < -20 C TIME = 42: Same as 6 for BT < -30 C TIME = 48: Same as 6 for BT < -40 C TIME = 54: Same as 6 for BT < -50 C TIME = 60: Same as 6 for BT < -60 C TIME = 66: max BT from 0 to 30 km radius (deg C*10) TIME = 72: avg BT from 0 to 30 km radius (deg C*10) TIME = 78: radius of max BT (km) TIME = 84: min BT from 20 to 120 km radius (deg C*10)

	<p>TIME = 92: avg BT from 20 to 120 km radius (deg C*10)</p> <p>TIME = 98: radius of min BT (km)</p> <p>TIME = 102 to 120: Variables need for storm size estimation</p>
IRXX	Same as IR00 above, but generated from other variables (not satellite data). These should only be used to fill in for missing IR00 if needed
IRM1	Same as IR00 but at 1.5 hours before initial time
IRM3	Same as IR00 but at three hours before initial time
PSLV	<p>Pressure of the center of mass (hPa) of the layer where storm motion best matches environmental flow (t=0 only). Also, the information used to calculate the steering layer pressure.</p> <p>All fields are valid at TIME = 0, and those in the TIME = 6 to TIME = 102 columns include the following:</p> <p>TIME = 6 column: The observed zonal storm motion component (m/s *10)</p> <p>TIME = 12 column: The observed meridional storm motion component (m/s *10)</p> <p>TIME=18, TIME=24 columns: Same as t=6, 12 hr columns but for the 1000 to 100 hPa mass weighted deep layer environmental wind (m/s *10) t=30, t=36 columns: Same as t=6,12 columns but for the optimally weighted deep layer mean flow (m/s *10)</p> <p>TIME=42 column: The parameter alpha that controls the constraint on the weights from being not too “far” from the deep layer mean weights (non-dimensional, *100)</p> <p>TIME=48 to TIME=102 columns: The optimal vertical weights for p=100, 150, 200, 250, 300, 400, 500, 700, 850 and 1000 hPa (non-dimensional *1000)</p>

MTPW	<p>Total Precipitable Water (TPW) variables at t=0 from the GFS analysis. The 21 values in this record are as follows:</p> <p>TIME = 0: 0-200 km average TPW (mm * 10)</p> <p>TIME = 6: 0-200 km TPW standard deviation (mm * 10)</p> <p>TIME = 12: 200-400 km average TPW (mm * 10)</p> <p>TIME = 18: 200-400 km TPW standard deviation (mm * 10)</p> <p>TIME = 24: 400-600 km average TPW (mm * 10)</p> <p>TIME = 30: 400-600 km TPW standard deviation (mm * 10)</p> <p>TIME = 36: 600-800 km average TPW (mm * 10)</p> <p>TIME = 42: 600-800 km TPW standard deviation (mm * 10)</p> <p>TIME = 48: 800-1000 km average TPW (mm * 10)</p> <p>TIME = 54: 800-1000 km TPW standard deviation (mm * 10)</p> <p>TIME = 60: 0-400 km average TPW (mm * 10)</p> <p>TIME = 66: 0-400 km TPW standard deviation (mm * 10)</p> <p>TIME = 72: 0-600 km average TPW (mm * 10)</p> <p>TIME = 78: 0-600 km TPW standard deviation (mm * 10)</p> <p>TIME = 84: 0-800 km average TPW (mm * 10)</p> <p>TIME = 90: 0-800 km TPW standard deviation (mm * 10)</p> <p>TIME = 96: 0-1000 km average TPW (mm * 10)</p> <p>TIME = 102: 0-1000 km TPW standard deviation (mm * 10)</p> <p>TIME = 108: %TPW less than 45 mm, r=0 to 500 km in 90 deg azimuthal quadrant centered on up-shear direction</p> <p>TIME = 114: 0-500 km averaged TPW (mm * 10) in 90 deg up-shear quadrant</p> <p>TIME = 120: 0-500 km average TPW (mm * 10)</p>
------	--

446

447       The NCODA system is “an oceanographic version of the multivariate optimum  
448 interpolation (MVOI) technique widely used in operational atmospheric forecasting  
449 systems. The ocean analysis variables in NCODA are temperature, salinity, geopotential  
450 (dynamic height), and velocity” (Cummings 2005). Related variables are NSST (SST  
451 from the NCODA analysis), NTMX (Max ocean temperature in the NCODA vertical  
452 profile), NDFR (Depth of the lowest model level in the NCODA analysis), NTFR (Ocean  
453 temperature at the lowest level in the NCODA analysis), NOHC (Ocean heat content  
454 from the NCODA analysis relative to the 26 degree C isotherm), NO20 (Same as NOHC

455 with respect to the 20 degree C isotherm), and XNST-XO20 (Climatological values of the  
456 NCODA variables with relate to the depth of the 30, 28, ..., 16 deg C isotherms).

457 The original format for SHIPS data is an ASCII text file, which consists of a large  
458 number of blocks, and each block involves variable names and their values from the  
459 current time up to 120 hours in a 6-hour interval. Some satellite and count variables only  
460 have one time only, and a few other variables are with values up to 12 hours before the  
461 current time.

462 Figure 2.1 displays an example of one data block, which has 141 lines with the line  
463 name at the end of each row with two special notation HEAD/LAST for the start/end of  
464 the block. TIME indicates the relative hour to the current time, and 9999 is filled when  
465 the value of a variable is not available. Based on the contents of the data block, all the  
466 lines can be divided into three categories:

- 467 • Special notations: HEAD, LAST, and TIME. The detailed information of HEAD,  
468 LAST, and TIME are displayed in Table 2.2.
- 469 • One-time variables: HIST, IRXX, IR00, IRM1, IRM3, PSLV, and MTPW: Those  
470 variables have only one-time values such as satellite and count variables with  
471 each TC instance, and the values on a particular column are with different  
472 meanings rather than time-dependent values. In Figure 2.1, 0 to 120 hours are  
473 corresponding to 21 values. The detailed information is displayed in Table 2.2,  
474 which also indicates that within each line, each element (column) only presents a  
475 different variable rather than providing values of variables at 6 hour interval. For  
476 example, that MTPW at TIME 6 column is 44 implies 0-200 km Total

477           Precipitable Water standard deviation (mm \* 10) is 44 based on Figure 2.1 and  
478           Table 2.1.

- 479           • Time-dependent variables: all other lines (variable names): Each line provides  
480           values of the corresponding variable at 6 hour interval at most from past 12 hours  
481           (-12 in TIME line) up to 120 hours (120 in TIME line) in the future as displayed  
482           in Figure 2.1. For example, VMAX at TIME 6 is 25 implies that 6 hours later  
483           from the moment of the block the VMAX is 25 (kt) based on Figure 2.1.

484

485   **Table 2.2: List of the special notations (SHIPS 2018a), explaining the details of each**  
486   **variable (Adopted from SHIPS (2018c)).**

Variable Name	Variable description
HEAD	Header line (1 <sup>st</sup> four letters of storm name, 2-digit year, month, day, and UTC time, maximum winds, lat, lon, minimum sea level pressure, and ATCF ID number (e.g., AL011982) at t=0 of the current case)
TIME	Time away from current
LAST	The last line for this case

487

488

489

490

491

492

493

494



(a)

[illegible]

496

497

498

499

500



501

(b)

```

580 572 578 562 533 541 513 499 489 502 513 501 494 478 477 452 444 9999 9999 9999 9999 Pw03
68 61 68 95 102 61 52 73 69 47 76 77 72 78 69 57 9999 9999 9999 9999 Pw04
536 521 534 529 494 493 487 493 483 489 486 491 477 459 457 439 429 9999 9999 9999 9999 Pw05
98 91 88 87 94 78 75 83 82 69 86 78 75 72 91 91 95 9999 9999 9999 9999 Pw06
477 477 507 512 480 483 486 481 480 485 489 483 454 436 441 420 414 9999 9999 9999 9999 Pw07
133 115 103 94 87 87 88 90 91 71 77 73 91 103 96 106 111 9999 9999 9999 9999 Pw08
465 459 485 484 475 490 516 479 471 476 481 454 431 421 424 402 389 9999 9999 9999 9999 Pw09
140 138 131 113 97 80 72 84 97 73 72 81 101 104 111 121 117 9999 9999 9999 9999 Pw10
589 583 585 571 553 552 520 504 492 505 516 504 500 486 483 458 454 9999 9999 9999 9999 Pw11
66 59 63 90 103 60 53 72 63 43 70 71 74 69 74 64 57 9999 9999 9999 9999 Pw12
558 547 556 548 520 519 501 498 487 496 499 496 487 471 469 447 440 9999 9999 9999 9999 Pw13
90 85 83 91 102 77 68 79 74 60 80 75 75 73 85 81 82 9999 9999 9999 9999 Pw14
522 516 534 532 503 503 495 491 484 491 495 490 472 456 457 435 429 9999 9999 9999 9999 Pw15
119 106 95 94 98 84 78 84 82 65 79 74 84 89 91 94 97 9999 9999 9999 9999 Pw16
502 496 516 514 493 499 503 486 479 485 490 477 457 443 445 423 414 9999 9999 9999 9999 Pw17
129 122 113 104 99 82 76 84 88 69 77 79 93 96 100 106 106 9999 9999 9999 9999 Pw18
113 83 116 439 758 328 494 700 695 418 778 746 847 848 942 885 860 9999 9999 9999 9999 Pw19
510 526 549 461 411 476 448 442 431 446 413 418 411 408 389 377 376 9999 9999 9999 9999 Pw20
575 564 570 555 534 535 512 503 488 501 505 499 491 479 475 452 446 9999 9999 9999 9999 Pw21
0 -105 200 -122 253 38 30 23 18 14 13 -54 -94 13 -173 -103 77 67 112 115 9999 IRXX
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 IR00
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 IRM1
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 IRM3
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 PC00
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 PCM1
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 PCM3
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 RD20
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 RD26
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 RHCN
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 NSST
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 NTMX
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 NDTX
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 NDML
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 ND30
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 ND28
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 ND26
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 ND24
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 ND22
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 ND20
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 ND18
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 ND16
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 NDFR
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 NTRF
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 NOHC
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 NO20
270 281 284 284 284 281 278 279 280 278 278 277 276 276 276 278 9999 9999 9999 9999 XNST
271 281 285 284 284 281 278 279 280 279 278 278 277 276 276 276 278 9999 9999 9999 9999 XTMT
0 0 2 5 4 1 1 1 1 1 1 1 1 1 1 1 9999 9999 9999 9999 XDTX
15 15 25 38 34 18 12 14 19 14 14 13 13 12 11 11 11 9999 9999 9999 9999 XDML
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 XD30
9999 10 21 29 27 13 9999 9999 11 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 XD28
21 37 62 87 78 40 23 27 38 28 26 24 23 20 19 20 9999 9999 9999 9999 XD26
50 65 99 128 119 67 44 29 65 50 47 45 43 40 38 39 41 9999 9999 9999 9999 XD24
86 97 133 166 156 96 66 72 94 73 70 67 65 61 64 64 65 9999 9999 9999 9999 XD22
123 131 171 214 201 126 93 100 125 101 98 95 92 89 90 81 9999 9999 9999 9999 XD20
166 170 222 273 261 164 122 132 164 133 129 126 123 120 122 116 9999 9999 9999 9999 XD18
240 233 287 347 333 215 161 172 213 175 169 166 163 162 161 140 9999 9999 9999 9999 XD16
1183 1929 2941 2958 3000 2963 2662 3000 3000 3000 3000 3000 3000 2635 1021 387 91 9999 9999 9999 9999 XDPR
57 46 44 43 44 44 45 45 44 44 45 45 45 47 90 143 205 9999 9999 9999 9999 XTFR
7 20 39 56 50 24 13 15 21 15 14 13 12 11 10 11 12 9999 9999 9999 9999 XOHC
177 223 326 421 392 226 151 167 219 169 161 154 147 138 136 140 142 9999 9999 9999 9999 XO20
266 280 283 281 281 278 274 275 278 275 275 274 274 272 272 272 273 9999 9999 9999 9999 XDST
9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 9999 LAST

```

502

503

504

505

506

**Figure 2.1: An example of data block of original SHIPS ASCII text file with 141 lines. a) the first 70 lines of the block. b) the rest lines with the first line duplicated with that in a).**

507

Only special notations and one-time variables are discussed in detail above

508

because all the time-dependent variable values are organized in the same format.

509

## 510 **2.2 ECMWF ERA-Interim reanalysis data**

511 ECMWF was founded in 1973 and is a research institute sponsored by several  
512 countries to produce numerical weather prediction to its member countries. ERA-Interim  
513 reanalysis data are generated by ECMWF every 6 hours and global wise covered with a  
514 horizontal resolution approximately 80km from its forecast numerical model to "improve  
515 on various technical aspects of reanalysis such as data selection, quality control, bias  
516 correction, and performance monitoring, each of which can have a major impact on the  
517 quality of the reanalysis products." (Dee et al. 2011). In short, the ERA-Interim reanalysis  
518 data is derived from the assimilating atmospheric model and can be regarded as the  
519 observed data. ERA-Interim reanalysis data has five data products, model level dataset,  
520 potential temperature dataset, potential velocity dataset, pressure level dataset, and  
521 surface dataset. Among them, the pressure level dataset is the most frequently used in TC  
522 researches. For example, Wang et al. (2015) adopted ERA-Interim pressure level data  
523 product to evaluate how vertical wind shear is influencing TC intensity change; Li et al.  
524 (2017) derived vorticity analysis from identifying TCs' track in Northwest Pacific Ocean  
525 Region based on ERA-Interim pressure level data product.

526 As the primary focus of this study is the TC in the Atlantic basin, ERA-Interim  
527 pressure level dataset, stored in a netCDF format, in the Atlantic basin is used. The  
528 Temporal and spatial coverage, the pressure levels, as well as the variable abbreviations  
529 of ERA-Interim pressure level dataset, are listed in Table 2.3, and the variables are  
530 explained in Table 2.4.

531

532

533 **Table 2.3: Temporal and spatial coverage of the ERA-interim pressure level data**  
 534 **and its available pressure levels and variables.**

Temporal Coverage	Four times daily, midnight, 6 am, noon, and 6 pm, January 1979 to two months delay of the moment
Spatial coverage	Global grid, 0.75 degree resolution
Pressure levels (37 to 1)	1, 2, 3, 5, 7, 10, 20, 30, 50, 70, 100, 125, 150, 175, 200, 225, 250, 300, 350, 400, 450, 500, 550, 600, 650, 700, 750, 775, 800, 825, 850, 875, 900, 925, 950, 975, 1000 hPa
Variables (short names explained in Table 2.4)	cc, ciwc, clwc, d, z, w, o3, pv, r, q, t, u, v, and vo

535

536 **Table 2.4: Variable names, abbreviations, units, and description for the 14 variables**  
 537 **in the ERA-Interim pressure level dataset.**

Variable Name	Abbreviation	Units	Description
Fraction of cloud cover	cc	percentage	Horizontal fraction of the grid box covered by cloud
Cloud ice water content	ciwc	kg kg <sup>-1</sup>	Grid-box mean specific cloud ice water content (mass of condensate / mass of moist air)
Cloud liquid water content	clwc	kg kg <sup>-1</sup>	Grid-box mean specific cloud liquid water content (mass of condensate / mass of moist air)
Divergence	d	s <sup>-1</sup>	Relative divergence
Geopotential	z	m <sup>2</sup> s <sup>-2</sup>	At the surface: orography
Vertical velocity	w	Pa s <sup>-1</sup>	Pressure vertical velocity dp/dt. In the model equations it is usually denoted by the Greek letter omega
Ozone mass mixing ratio	o3	kg kg <sup>-1</sup>	Mass mixing ratio of Ozone
Potential vorticity	pv	K m <sup>2</sup> kg <sup>-1</sup> s <sup>-1</sup>	The ability of air to rotate in the atmosphere. “conservation equation directly ties together the dynamics and the heating” (Molinari 1989)

Relative humidity	r	percentage	Relative humidity is defined with respect to saturation of the mixed phase, i.e. with respect to saturation over ice below -23°C and with respect to saturation over water above 0°C. In the regime in between a quadratic interpolation is applied
Specific humidity	q	kg kg <sup>-1</sup>	Grid box mean (mass of water vapour / mass of moist air)
Temperature	t	K	Temperature
U component of wind	u	m s <sup>-1</sup>	West to east flow (eastward wind)
V component of wind	v	m s <sup>-1</sup>	South to north flow (northward wind)
Vorticity (relative)	vo	s <sup>-1</sup>	Measure of the rotation of air in the horizontal

538

539           As shown in Table 2.3, there are 14 variables and 37 pressure levels in the  
540 pressure level dataset, and some of these variables are the same as those used in SHIPS  
541 database, for example, vo in 850-hpa, t in 200-hPa, and v in 850-hpa to 200-hpa (Kaplan  
542 and DeMaria 2003).

543 Table 2.4 displays the 14 variables with their explanations. Fraction of cloud cover (cc),  
544 cloud ice water content (ciwc), and cloud liquid water content (clwc) are similar variables  
545 since cc presents the proportional of a grid box covered either by liquid cloud or ice  
546 cloud, while ciwc presents in each grid box the mass of cloud ice particles per kilogram  
547 of the total mass of dry air, water vapor, cloud liquid, cloud ice, rain, and falling snow,  
548 and clwc is almost the same as ciwc except that cloud ice particles is replaced by cloud  
549 liquid water droplets. Divergence (d) implies the horizontal divergence rate of the  
550 velocity that the air is spreading out from a point, and d is positive when the air is

551 spreading out and negative vice versa. Geopotential ( $z$ ) indicates the amount of work to  
552 lift a unit of air from the mean sea level to a certain point against the gravity. Vertical  
553 velocity ( $w$ ) indicates the speed of air moving upward or downward. Ozone mass mixing  
554 ratio ( $o_3$ ) indicates “the mass of ozone per kilogram of air.” Potential vorticity ( $pv$ )  
555 indicates the ability of air to rotate in the atmosphere and usually is used to look for  
556 places where wind storm likely to occur. Relative humidity ( $r$ ) indicates water vapor  
557 pressure as a percentage of a fixed vapor pressure value when water vapor becomes  
558 liquid water or ice. Specific humidity ( $q$ ) is the same as  $ciwc$ , but the mass of cloud ice  
559 particles per kilogram is replaced with the mass of water vapor per kilogram.  
560 Temperature ( $t$ ) indicates the temperature in the atmosphere. U component of wind ( $u$ )  
561 and V component of wind ( $v$ ) indicates “the horizontal speed of air moving towards the  
562 east, in meters per second” and “the horizontal speed of air moving towards the north, in  
563 meters per second,” respectively.  $u$  and  $v$  can be combined to calculate the speed and  
564 direction of the wind. The negative sign of  $u$  and  $v$  shows that air moving to the west and  
565 south. Relative vorticity ( $vo$ ) is a clockwise (positive) or counter clockwise (negative) air  
566 spin.

### 567 **2.3 NHC best track data**

568 To determine which part of the ERA-interim data should be used, the center of  
569 TCs need to be located. The NHC best track (HURDAT2) data, available every 6 hours  
570 (midnight (UTC 0), 6 am (UTC 600), noon (UTC 1200), and 6 pm (UTC 1800)),  
571 including the time, longitude, latitude, maximum sustained wind speed of the TCs will be  
572 used. Table 2.5 is a concise version from the original NHC best track data and shows a

TC cataloged as AL011982, indicating the TC is the 1<sup>st</sup> TC in 1982 that occurred in the Atlantic Ocean and had seventeen records named ALBERTO. The first record indicates the TC was recorded for June 2<sup>nd</sup>, 1982 at 1200 UTC. The TC was centered at 21.7°N and 87.1°W with an intensity of 20 knots and minimum pressure 1005 millibars. RI occurs if and only if sustained wind speed in the next 24 hours increases 30 knots or more. Therefore, TC AL011982 undergoes RI at 1200UTC and 1800 UTC on June 2<sup>nd</sup> 1981, and at 0 UTC and 600 UTC on June 3<sup>rd</sup> 1981.

**Table 2.5: A TC record in NHC best track data.**

AL011982			ALBERTO	17		
Date	Time	System status	Latitude	Longitude	Maximum sustained wind speed (in knots)	Minimum Pressure (in millibars)
19820602	1200	TD	21.7N	87.1W	20	1005
19820602	1800	TD	22.2N	86.5W	25	1004
19820603	0	TD	22.6N	85.8W	30	1003
19820603	600	TS	22.8N	85.0W	40	1001
19820603	1200	TS	23.2N	84.2W	50	995
19820603	1800	HU	24.0N	83.6W	75	985
19820604	0	HU	24.8N	83.4W	65	992
19820604	600	TS	24.9N	84.1W	55	998
19820604	1200	TS	24.9N	84.8W	45	1002
19820604	1800	TS	25.0N	84.2W	40	1005
19820605	0	TD	25.1N	84.1W	30	1007
19820605	600	TD	25.2N	84.0W	25	1008
19820605	1200	TD	25.3N	83.9W	25	1009
19820605	1800	TD	25.4N	83.6W	25	1010
19820606	0	TD	25.5N	83.3W	25	1010
19820606	600	TD	25.5N	83.0W	25	1010
19820606	1200	TD	25.5N	82.6W	20	1010

583           NHC best track data is used to identify TC's center, and its related information  
584   from the entire ERA-Interim reanalysis data. The two datasets, with the same temporal  
585   resolution, will be processed through the ERA-Interim data filter, and details will be  
586   discussed in Chapter 3.



## CHAPTER 3 DATA FILTER

As indicated in Figure 1.1, the COR-SHIPS model, LLE-SHIPS model, and DL-SHIPS model, although differ in certain parts, share the same logical structure. They share three of the four components, data sampler, classifier, and hyperparameter tuning process and differ only in the data filter. COR-SHIPS model employs only the SHIPS data filter, while the LLE-SHIPS model and DL-SHIPS model adopt extra ERA-Interim data with filters based on local linear embedding (LLE) and deep learning (DL), respectively, in addition to the SHIPS data filter.

### **3.1 SHIPS data filter**

#### **3.1.1 ASCII text to attribute-relation table**

For the RI classification analysis by the XGBoost classifier, the input data model is attribute-relation tables commonly used for relational databases, and therefore, each SHIPS instance block should be transformed into one entry in an attribute-relation table. One sample block of the SHIPS ASCII text file is displayed in Figure 2.1, which will be converted to one entry in an attribute-relation table, as shown in Figure 3.1. During the conversion, special notation, one-time variables, and time-dependent variables are handled differently and are described below in detail. The number of variables for each category is also listed for tracking purposes only.

1. Special notations (3 lines in total: HEAD, LAST, and TIME): The HEAD line has ten elements: TC name (NAME); two-digit year (YEAR), month (MONTH), date (DATE); UTC; maximum surface wind; center latitude; center longitude; the

608 minimum sea level pressure; and the ATCF ID number (ATCF). Because the  
609 maximum surface wind; center latitude; center longitude; the minimum sea level  
610 pressure are included in the time dependent variables, only the rest six variables are  
611 extracted from the HEAD line: NAME, YEAR, MONTH, DATE, UTC, and ATCF.  
612 No information is retrieved from the TIME and LAST lines (result in 6 variables).

613 2. One-time variables (7 lines: HIST, IRXX, IR00, IRM1, IRM3, PSLV, and MTPW):  
614 As indicated in section 2.1, since these seven lines contain values for many one-time  
615 variables, and index is added to denote and to distinguish the corresponding variables.  
616 Therefore, HIST\_0, HIST\_1, ..., HIST\_20 are created using values in the HIST line.  
617 Similarly, IRXX\_1, ..., IRXX\_20, IR00\_1, ..., IR00\_20, IRM1\_1, ..., IR M1\_20,  
618 IRM3\_1, ..., IRM3\_20 (when TIME = 0, IRXX, IR00, IRM1, and IRM3 present the  
619 relative time of GOES image relative to the instance, not related to the problem in the  
620 study, and therefore, are removed), PSLV\_1, ..., PSLV\_18 (values in 108, 114, and  
621 120 hour are filled with 9999), MTPW\_0, ..., MTPW\_20 are created using values in  
622 the corresponding lines (result in 140 variables).

623 3. Time dependent variables (the remaining 131 lines<sup>2</sup>): The current value of a time  
624 dependent parameter ( $t=0$ ) is associated with the corresponding TC instance. The  
625 values for other times (from previous 12 hours to future 120 hours) may be used to  
626 derive other variables or simply ignored (result in 131 variables, total 277 variables).  
627

---

<sup>2</sup> Three lines, PC00, PCM1, PCM3 were mistreated as time-dependent initially. See details in the main text.

628  
629  
630

**Table 3.1: One row of the attribute-relation table converted from original SHIPS data showed in Figure 2.1. Number 1 to 277 corresponds to the 1<sup>st</sup> to the 277<sup>th</sup> columns in the attribute-relation table is added for notation only here.**

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
NAME	YEAR	MONT H	DATE	TIME	ATCF	VMA X	MSLP	TYPE	DELV	INCV	LAT	LON	CSST	CD20	CD26
ALBE	1982	6	2	12	AL011 982	20	1005	1	0	9999	217	871	274	150	47
17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32
COHC	DTL	OAGE	NAGE	RSST	DSST	DSTA	U200	U20C	V20C	E000	EPOS	ENEG	EPSS	ENSS	RHLO
24	21	0	0	280	274	276	224	238	99	3528	113	8	51	16	70
33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48
RHMD	RHHI	Z850	D200	REFC	PEFC	T000	R000	Z000	TLAT	TLON	TWAC	TWXC	G150	G200	G250
57	50	7	64	4	-1	259	84	-27	207	873	67	100	-5	3	13
49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64
V000	V850	V500	V300	TGRD	TADV	PENC	SHDC	SDDC	SHGC	DIVC	T150	T200	T250	SHRD	SHTD
44	64	67	21	10	2	103	264	80	310	49	-665	-526	-408	251	94
65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80
SHRS	SHIS	SHRG	PENV	VMPH	VVAV	VMFX	VVAC	HE07	HE05	O500	O700	CFLX	PW01	PW02	PW03
90	137	296	95	123	1288	860	1382	0	-29	-82	-60	139	618	44	580
81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96
PW04	PW05	PW06	PW07	PW08	PW09	PW10	PW11	PW12	PW13	PW14	PW15	PW16	PW17	PW18	PW19
68	536	98	477	133	465	140	589	66	558	90	522	119	502	129	113
97	98	99	100	101	102	103	104	105	106	107	108	109	110	111	112
PW20	PW21	PC00	PCM1	PCM3	RD20	RD26	RHCN	NSST	NTMX	NDTX	NDML	ND30	ND28	ND26	ND24
510	975	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999
113	114	115	116	117	118	119	120	121	122	123	124	125	126	127	128
ND22	ND20	ND18	ND16	NDFR	NTFR	NOHC	NO20	XNST	XTMX	XDTX	XDML	XD30	XD28	XD26	XD24
9999	9999	9999	9999	9999	9999	9999	9999	270	271	0	15	9999	9999	21	50
129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144
XD22	XD20	XD18	XD16	XDFR	XTFR	XOHC	XO20	XDST	HIST_1	HIST_2	HIST_3	HIST_4	HIST_5	HIST_6	HIST_7
86	123	166	240	1183	57	7	177	266	1	0	0	0	0	0	0
145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160
HIST_8	HIST_9	HIST_1 0	HIST_1 1	HIST_1 2	HIST_1 13	HIST_1 14	HIST_1 15	HIST_1 16	HIST_1 8	HIST_1 9	HIST_1 0	HIST_2 0	HIST_2 1	PSLV_1	PSLV_2
0	0	0	0	0	0	0	0	0	0	0	0	0	0	548	31
161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176
PSLV_3	PSLV_4	PSLV_5	PSLV_6	PSLV_7	PSLV_8	PSLV_9	PSLV_10	PSLV_11	PSLV_12	PSLV_13	PSLV_14	PSLV_15	PSLV_16	PSLV_17	PSLV_18
23	24	13	28	16	40	23	64	82	77	89	92	85	208	200	79
177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192
MTPW_1	MTPW_2	MTPW_3	MTPW_4	MTPW_5	MTPW_6	MTPW_7	MTPW_8	MTPW_9	MTPW_10	MTPW_11	MTPW_12	MTPW_13	MTPW_14	MTPW_15	MTPW_16
618	44	580	68	536	98	477	133	465	140	589	66	558	90	522	119
193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208
MTPW_17	MTPW_18	MTPW_19	MTPW_20	MTPW_21	IRXX_1	IRXX_2	IRXX_3	IRXX_4	IRXX_5	IRXX_6	IRXX_7	IRXX_8	IRXX_9	IRXX_10	IRXX_11
502	129	113	510	575	0	-105	200	-122	253	38	30	23	18	14	13
209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224
IRXX_12	IRXX_13	IRXX_14	IRXX_15	IRXX_16	IRXX_17	IRXX_18	IRXX_19	IRXX_20	IR00_1	IR00_2	IR00_3	IR00_4	IR00_5	IR00_6	IR00_7
-54	-94	13	-173	-103	77	67	112	115	9999	9999	9999	9999	9999	9999	9999
225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240
IR00_8	IR00_9	IR00_10	IR00_11	IR00_12	IR00_13	IR00_14	IR00_15	IR00_16	IR00_17	IR00_18	IR00_19	IR00_20	IRM1_1	IRM1_2	IRM1_3
9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999
241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256
IRM1_4	IRM1_5	IRM1_6	IRM1_7	IRM1_8	IRM1_9	IRM1_10	IRM1_11	IRM1_12	IRM1_13	IRM1_14	IRM1_15	IRM1_16	IRM1_17	IRM1_18	IRM1_19
9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999
257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272
IRM1_20	IRM3_1	IRM3_2	IRM3_3	IRM3_4	IRM3_5	IRM3_6	IRM3_7	IRM3_8	IRM3_9	IRM3_10	IRM3_11	IRM3_12	IRM3_13	IRM3_14	IRM3_15
9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999
273	274	275	276	277											
IRM3_16	IRM3_17	IRM3_18	IRM3_19	IRM3_20											
9999	9999	9999	9999	9999											

631  
632

633 In SHIPS, three lines, PC00, PCM1, PCM3, give the first nine principal components  
634 (PCs) for IR imageries at the current time, 1.5 hours before, and 3 hours before. Initially,  
635 we misinterpret the variables as time-dependent variables and keep the current values,  
636 PC00\_1, PCM1\_1, and PCM3\_1, only. Although it is not a correct treatment for those  
637 values, fortunately, the kept values are for the 2<sup>nd</sup> principal components, which are the  
638 only important components identified by previous studies (KRD15).

639 With those information extractions and format conversion, each block (a TC  
640 instance) is converted to a row similar to the form in Figure 3.1, with a total of 277  
641 variables. In Figure 3.1, the first row is added to denote the column series numbers from  
642 1 to 277 for notation purposes only. The second row indicates the names of the variables  
643 or headers, and the third row gives the values of the variables for one instance. And all  
644 values for all other instances are stacked together to constitute an attribute-relation table  
645 with all TC instances for this study.

### 646 **3.1.2 Preprocessing of the SHIPS data in attribute-relation table**

647 The raw attribute-relation table obtained above is hard to be used directly due to  
648 variable natures, i.e., irrelevant variables, heavy missing values (9999 in original SHIPS  
649 data), scaling issue, and the inter-correlation between variables. In addition, some  
650 potential variables are not available directly from the simple conversion procedure. As a  
651 result, a preprocessing is performed on the raw attribute-relation table.

#### 652 **3.1.2.1 Adding additional variables**

653 Based on previous studies (e.g., DeMaria et al. 2005), intensity change is critical  
654 for rapid intensification prediction. Therefore, the previous 6-hour intensity change

655 (BD06) is calculated as subtracting the current intensity by the intensity 6 hours before,  
656 and BD06 for the first instance of each TC is set as missing. Previous 12-hour intensity  
657 change (BD12) and 18-hour intensity change (BD18) are calculated similarly. The first  
658 two (three) instances of each TC for BD12 (BD18) are set as missing. Since BD06 and  
659 BD12 contain the information of DELV and INCV, the latter two are removed (3  
660 variables added and 2 removed resulting in 1 more variable, so 278 variables remained in  
661 total).

662 To include temporal variation of RI, annual Julian date is created to combine the  
663 information of MONTH and DATE, while MONTH, DATE, and UTC are removed. In  
664 addition, TYPE (storm type) should not have any influence of RI prediction and hence is  
665 also removed (1 variable added and 4 removed leading to 3 variables less, so 275  
666 variables remained in total).

### 667 ***3.1.2.2 Variable removal***

668 Some variables, such as the first four letters of storm name (NAME), are unique  
669 information used for tracing back the specific TC, and are unrelated to the RI prediction,  
670 which should be removed. ATCF and YEAR are also such variables and hence are  
671 removed. In addition, PSLV1 to PSLV18 represent storm motion information, where  
672 PSLV1 to PSLV8 are the storm motion component, and PSLV\_9 to PSLV\_18 indicates  
673 the optimal vertical weights for various pressures levels. Therefore, PSLV\_9 to PSLV\_18  
674 are not related to TC rapid intensification and hence is removed (13 variables removed,  
675 so 262 variables remained in total).

IRXX should only be used when IR00 values are missing. Therefore, missing values in IR00\_1, ..., IR00\_20, are replaced with the corresponding values in IRXX\_1, ..., IRXX\_20 respectively and then IRXX\_1, ..., IRXX\_20 are removed (20 variables removed, so 242 variables remained in total).

Missing values in a variable do not provide much information for that variable. The more missing values in a variable, the less information it has. Table 3.1 displays the missing value percentage in the SHIPS data (SHIPS 2018b) for variables with at least one missing filling. XD30 contains more than 95% missing values, while the next nine variables have more than 50% missing values. Other variables have less than 45% missing percentages. Since variables with more than 50% are not expected to give much information in the RI prediction, XD30, NDML, ND30, ND28, ND24, ND22, ND18, ND16, NO20, and XD28 are removed. After this removal, the remaining missing values are coded as NA; a notation can be easily handled later for sampling and classification (10 variables removed, so 232 variables remained in total).

690

691

**Table 3.2: Variables with missing value and the missing percentage in SHIPS data. Variables are sorted according to the percentage. Variables without missing values are not listed.**

Variable	Percentage	Variable	Percentage	Variable	Percentage
XD30	98.08%	RHCN	34.18%	XDML	5.05%
NDML	61.98%	PCM1	21.86%	XDTX	4.39%
ND30	61.84%	XD26	20.46%	BD06	4.30%
ND28	61.84%	IRM1_2 to IRM_20	14.12%	XDFR	2.41%
ND24	61.84%			CD26	2.10%
ND22	61.84%	PC00	14.06%	MSLP	1.08%

ND18	61.84%	IRM1_1	13.84%	TWAC	0.88%
ND16	61.84%	IRM3_2 to IRM2_20	12.93%	TWXC	0.88%
NO20	61.84%			DIVC	0.88%
XD28	51.78%	BD18	12.88%	COHC	0.48%
RD20	41.41%	XD24	12.87%	XNST	0.36%
RD26	41.41%	PCM3	12.84%	XTMX	0.36%
NDFR	39.57%	IRM3_1	12.57%	XTFR	0.36%
NTFR	37.62%	XD22	10.32%	XOHC	0.36%
NSST	37.31%	XD20	9.23%	XO20	0.36%
NTMX	37.31%	IR00_1 to IR00_20	9.12%	XDST	0.36%
NDTX	37.31%			DSST	0.03%
ND26	37.31%	XD18	8.96%	PSLV_2 to PSLV_8	0.01%
ND20	37.31%	XD16	8.65%		
NOHC	37.31%	BD12	8.60%		

695

696 In addition to variables with high missing value percentages, there are variables  
697 with a very high percentage of a single value, and those variables are of less value in the  
698 RI prediction. Table 3.2 shows variables with its largest percentages (greater than 50) of  
699 single values, and PSLV\_8, the constraint on the weight, has more than 99.999% of  
700 instances with a single value, 40, is removed by assuming a threshold of 90% (1 variable  
701 removed, so 231 variables remained in total).

702

703 **Table 3.3: Variables with higher than 50% single values in the SHIPS Data, the**  
704 **single values, and the percentages.**

Variable	Value	Percentage
PSLV_8	40	99.99%
HIST_21	0	89.76%
HIST_20	0	88.10%
HIST_19	0	86.39%

HIST_18	0	83.84%
HIST_17	0	80.59%
HIST_16	0	78.77%
HIST_15	0	75.79%
HIST_14	0	73.81%
HIST_13	0	70.94%
HIST_12	0	66.30%
IRM3_1	-245	63.46%
PCM3	-245	63.27%
HIST_11	0	61.81%
IR00_1	15	58.89%
PC00	15	58.75%
HIST_10	0	56.41%
HIST_9	0	50.62%

705

706           The above components, data conversion from ASCII blocks to the attribute-  
707 relation table, irrelevant variable removal, missing value, and single value handling, new  
708 attribute creators together construct the SHIPS data filter. Through this filter, the original  
709 ASCII based block SHIPS dataset (SHIPS, 2018c) is filtered into an attribute-relation  
710 table for all TC instances with one TC instance as one row with 231 attributes (columns).

### 711 ***3.1.2.3 Rescale variables between 0 and 1.***

712           To make the data internally consistent, numerical values for all variables are  
713 rescaled except for missing values (NA). All values are rescaled to between 0 and 1 using

$$714 \quad \frac{Value_{variable} - Min_{variable}}{Max_{variable} - Min_{variable}} \quad (3.1)$$



715 where  $\text{Value}_{\text{variable}}$  represents the value of a particular variable that needs to be  
716 standardized,  $\text{Max}_{\text{variable}}$  and  $\text{Min}_{\text{variable}}$  represent the maximum and minimum values  
717 of the particular variable across all instances.

### 718 **3.1.3 Removal of highly correlated variables**

719 Highly correlated input variables could influence the accuracy of the variable  
720 importance evaluation. Therefore, among highly correlated variables in the SHIPS  
721 dataset, only one variable should be kept while others being removed.

722 The definition of “highly correlated” depends on a predefined correlation threshold,  
723 which is related to the number of variables to be removed (or kept). This correlation  
724 threshold is one of the so-called hyperparameters.

725 As the first step, pairwise correlations of all the variables are calculated and  
726 compared with the correlation threshold. For each variable, its highly correlated variables  
727 (correlation higher than the threshold) are identified and placed in a group started with  
728 that variable. Therefore, there is a group that started with each of the 231 variables. All  
729 groups are sorted in descending order based on the number of variables they have, and if  
730 the length of the two groups is the same, the sorting is alphabetically based on the leading  
731 variable names. Then starting from the first group, all following groups starting with any  
732 current group members are eliminated, and then member variables are also removed from  
733 all other remaining groups to guarantee that each variable will appear only once. This  
734 process continues until the last group is reached. After that, the first variables in all  
735 remained groups are selected as filtered variables.

For example, Table 3.3 displays a correlation matrix between BD06, BD12, and BD18, and the selection threshold is defined as 0.8. BD12 has two correlated variables, BD06 and BD18, while BD06 and BD18 have only one correlated variable, BD12 each. Therefore BD12's group with length two is placed before BD06's, and BD18's of length one in the sorted list (not shown). In the removal process, groups started with BD06, and BD18 are removed because both are in the group starting BD12. If there were more groups (not this case), BD06 and BD18 should also be removed in all appearance. For this extremely simplified example, only variable BD12 is kept.

**Table 3.4: Correlation matrix for pairs among BD06, BD12, and BD18, and the highly correlated group lists leading with each variable based on a 0.8 correlation threshold.**

	<b>BD06</b>	<b>BD12</b>	<b>BD18</b>	<b>Highly correlated variable (<math>\geq 0.8</math>)</b>
<b>BD06</b>	1.00	0.86	0.75	BD06, BD12
<b>BD12</b>	0.86	1.00	0.92	BD12, BD06, BD18
<b>BD18</b>	0.75	0.92	1.00	BD18, BD12

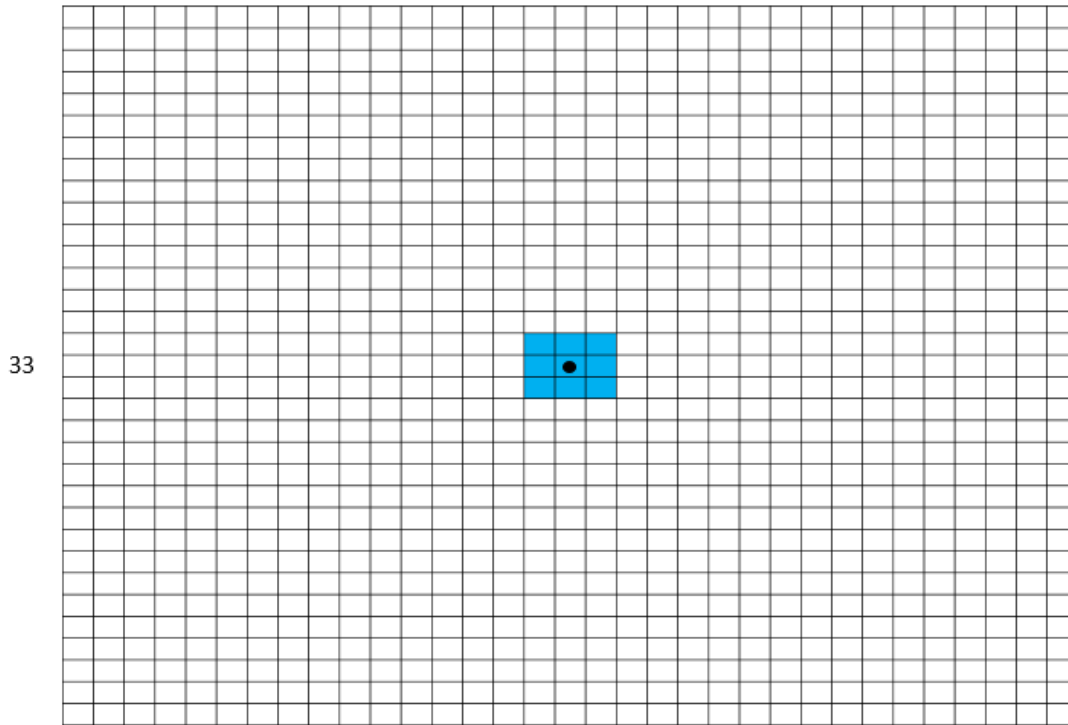
The final result from the above procedure should be sensitive to the to-be-tuned hyperparameter correlation threshold. If the threshold is too low, important variables could be removed, and the accuracy of the model is reduced; if the threshold is too high, the variable importance will be evaluated inaccurately because the variable importance score is largely influenced by highly correlated variables (multicollinearity). In this study,

754 0.7, 0.8, 0.9, and 0.95 are predetermined as the searching space to tune (find the best  
755 value for) this hyperparameter. The details will be discussed in Chapter 6.

### 756 **3.2 ERA-Interim data filters**

757 The SHIPS data filter converts ASCII SHIPS data from the instance block format to  
758 an attribute-relation table and is the only data filter used in the COR-SHIPS model. To  
759 identify new features beyond SHIPS and to improve the RI prediction performance, the  
760 ERA-Interim data set is used together with SHIPS data in this study. As for the SHIPS  
761 data, this gridded pressure level data should also be filtered into the attribute-relation  
762 table format.

763 ERA data is filtered in two different ways to inspect the near core information and  
764 large-scale effects, respectively. As indicated in Figure 3.2, the average of the blue 3\*3  
765 grid boxes is used to present the near center information, and the local linear embedding  
766 (LLE) model is employed to filter the near center values, leading to the LLE-SHIPS  
767 model. For large-scale features, 33\*33 grid boxes around the TC centers are filtered with  
768 a deep learning (DL) model, and the result is the DL-SHIPS model. Comparing with the  
769 LLE data filter, which can only process small scale near core information, DL can extract  
770 more information in the large-scale, including the smaller core information processed by  
771 the LLE model, but also has the risk of overfitting with the complicated structure.  
772 Moreover, the structure of the DL model is very complicated; therefore, it is much more  
773 difficult to evaluate the feature importance for the DL-SHIPS model than for the LLE-  
774 SHIPS model.



**Figure 3.1: The 33\*33 grid boxes centered at the grid box consisting of the center of a TC, denoted as the black dot, and the blue area presents the near core grids.**

### 3.2.1 Local Linear Embedding (LLE) for filtering near core ERA-Interim data

The nine grid boxes around a TC center together with an approximated 240km\*240km size are considered as the near center area, and values for 14 variables at 37 pressure levels are averaged over the nine boxes to have 37\*14 representative values.

As well known, the RI status is not only based on the current moment but also that of the last 18 hours. Therefore, for each instance, we include data from the previous 18 hours to current, and therefore have  $4*37*14 = 2072$  features (variables). The features are labeled in a “time\_variable\_level” format. 18 hours before, 12 hours before, 6 hours

787 before, and at present are represented as NT18, NT12, NT06, and NT00. Level 1 (1000  
788 hPa) to level 37 (1 hPa) are represented as l1, ..., l37. For example, NT18\_pv\_l37  
789 represents the horizontally averaged value of pv (potential vorticity) at pressure level 37  
790 (1000 hPa) at 18 hours before the TC instance time.

791 For a classification problem, 2072 is a very high dimension number, and overfitting  
792 cannot be avoided if all of the 2072 features are used in a model. Therefore, we need to  
793 reduce the feature dimension, which can be achieved using the principal component  
794 analysis (PCA), designed to reduce the feature dimensions while keeping as much  
795 statistical information as possible. Pearson (1901) first came up with the idea of PCA,  
796 and later, this idea was independently proposed by Hotelling (1933). PCA is applied in a  
797 large number of areas such as exploratory analysis (Li and Ralph 2019), dimension  
798 reduction approach (Ron 2000; Labib and Vemuri 2006), and Geostationary Operational  
799 Environmental Satellite (GOES-East and GOES-West) infrared imagery variables  
800 creation (Kaplan et al. 2015).

801 However, the principal components are constructed with a linear combination of the  
802 original features; hence nonlinear structure between these features is missed. To break  
803 this limit, kernel PCA, which uses a non-linear kernel to transfer the original feature  
804 space to a kernel Hilbert space and, therefore, to account for the nonlinear structure  
805 between features, was proposed (Schölkopf et al. 1998; Yang et al. 2006). Local linear  
806 embedding (LLE) is a type of kernel PCA and was first introduced in Roweis and Saul  
807 (2000) for dimension reduction. In traditional dimension reduction approaches such as  
808 the regular kernel PCA or multidimensional scaling (MDS), when the new reduced space

809 is searched, the geometry distance between different observations are not preserved. That  
 810 is, far away observations in the original feature space may be mapped to their  
 811 neighborhood in the new reduced feature space. By contrast, LLE preserves the global  
 812 geometry structure from locally linear fits in the new space. In other words, low  
 813 dimensional representation of the high dimensional data is discovered, where these local  
 814 relationships are best preserved (Géron 2017).

815 LLE is used to handle the near core ERA data in this study and leads to the LLE-  
 816 SHIPS model. As indicated above, each instance has  $4 \times 37 \times 14 = 2072$  features. Those  
 817 features are rescaled to numbers between 0 and 1 using  $(value - \min) / (max - \min)$   
 818 where value presents the raw value while min and max imply the minimum and  
 819 maximum value within 18 hours over all 37 levels for each of the 14 variables, similar to  
 820 rescaling process in Equation (3.1).

821 There are two steps in LLE itself. The first step is to evaluate “how each training  
 822 instance linearly relates to its closest neighbors,” and then in the second step, “looking for  
 823 a low dimensional representation of all training instances, where these local relationships  
 824 are best preserved” (Géron 2017).

825 Mathematically, we can elaborate the first step of LLE details with the following  
 826 equations:

$$\hat{W} \sim \underset{W}{argmin} \sum_{i=1}^m ||x^{(i)} - \sum_{j=1}^k w_{i,j} x^{(j)} ||^2 \quad (3.2)$$

827 subject to

$$\begin{cases} w_{i,j} = 0 & \text{if } x^{(j)} \text{ is not one of the } k \text{ nearest neighbor of } x^{(i)}, \\ \sum_{j=1}^m w_{i,j} = 1. \end{cases} \quad (3.3)$$

In the above equations, there are  $m$  instances in the entire dataset  $(x^{(i)}, i = 1, \dots, m)$ , where each instance  $x^{(i)}$  ( $i = 1, \dots, m$ ) is a vector of dimension 2072.  $k$  is a to-be-tuned integer defining the neighborhood size. Then  $w_{i,j}$ ,  $j = 1, \dots, k$  are the weights for the  $k$  nearest neighbors of  $x^{(i)}$ , and  $w_{i,j}$  is summed up to 1 over all neighbors and set as 0 when  $x^{(j)}$  is not the neighbor of  $x^{(i)}$ .  $w_{i,j}$  are trained to minimize the sum of the square distance between  $x^{(i)}$  and its weighted neighbors' sum,  $\sum_{j=1}^k w_{i,j} x^{(j)}$ .  $\widehat{W}$  is the solution of the weight matrix  $W$  (the matrix form of  $w_{i,j}$ ), that satisfies Equation (3.2).

In the second step, after the trained weights are calculated, instances in the entire dataset are mapped to a  $d$ -dimensional space ( $d < 2072$ , which is undefined, and its tuning will be discussed later together with  $k$ ) while preserving the relationship between instances as much as possible.  $z^{(i)}$  is the image of  $x^{(i)}$  in the  $d$ -dimension space,  $i = 1, \dots, m$ . The weight  $\widehat{W}$  derived from step 1 is fixed, and the sum of the squared distance between  $z^{(i)}$  and its weighted neighbors should be minimized to look for  $z^{(j)}$  ( $j = 1, \dots, m$ ). That is,

$$\hat{Z} = \underset{Z}{\operatorname{argmin}} \sum_{i=1}^m ||z^{(i)} - \sum_{j=1}^k w_{i,j} z^{(j)} ||^2 \quad (3.4)$$

subject to

$$\begin{cases} \sum_{i=1}^m z^{(i,j)} = 0 & j=1, \dots, d \\ \frac{1}{m} Z' Z = I_d \end{cases} \quad (3.5)$$

where each the  $z^{(i,j)}$  ( $i = 1, \dots, m; j = 1, \dots, d$ ) represents  $i$ -th instances in  $j$ -th dimension, and  $Z$  is the matrix form of  $z_{i,j}$ .  $z^{(i)}$  is summed to 0 over all  $d$  dimensions, and the covariance matrix of  $Z$  be the ( $d$ -dimensional)  $I_d$ , where  $I_d$  indicates the identity matrix with  $d \times d$  dimension (Roweis and Saul 2000).  $\hat{Z}$  is the solution of the weight matrix  $Z$  that satisfies Equation (3.4).

Based on Ginsburg et al. (2016), LLE can be derived as a kernel PCA with kernel  $K = \lambda_{max} I - (1 - \hat{W})(1 - \hat{W}^T)$ . More technical details of PCA and Kernel PCA are given in Appendix 1.

In this work, LLE is used to reduce the original space with 2072 dimensions to the new  $d$ -dimensional space while preserving the maximal global geometry structure. How to define  $d$  is very subjective; if  $d$  is too large, the reduction of dimension is light, and the possibility of overfitting is not reduced much. And if  $d$  is too small, some important information about the original space may be lost in the new reduced space. Therefore,  $d$  is the hyperparameter of LLE that needs to be determined. Since SHIPS data filter outputs approximately 72 variables, 10 and 90 are defined as the lower bound and the upper bound to search for  $d$ .

Another hyperparameter that describes how much geometry information should be kept is the number of the nearest neighbors identified for each observation,  $k$ . If  $k$  is too low, less geometry information is kept, and the new variables lose too much information



from the original variables. If  $k$  is too high, the computational cost is too high, and overfitting cannot be avoided. Therefore, to keep the balance, 5 and 15 are defined as the lower bound and the upper bound for searching  $k$ .

Although  $k$  and  $d$  are independent of each other, the number of dimension ( $d$ ) is usually larger than the number of neighbors ( $k$ ). The search range for  $k$  and  $d$  are summarized in Table 3.4, and the tuning details will be discussed in Chapter 6.

**Table 3.5: Hyperparameters for the LLE and their searching range defined by the Min(imum) and Max(imum).**

Hyperparameter	Explanation	Min	Max
no_neighbors ( $k$ )	The number of the nearest neighbors	5	15
no_dimension ( $d$ )	The dimensions in the reduced space	10	90

### 3.2.2 Deep learning (ERA-Interim data filter for DL-SHIPS model)

The large-scale range area with 33\*33 grid cells cannot be processed by LLE because even a supercomputer cannot handle the computational cost. Correspondingly, an alternative data filter based on deep learning (DL), a well-known technique for its capacity to handle a large amount of data, is used to process information in the large-scale range.

DL is a kind of Artificial Neural Network (ANN) model, which is designed to solve learning tasks by imitating the human biological neural network. The first functional ANN like model was proposed by Hodgkin and Huxley (1952), who had used non-linear

885 features and multiple layers to develop a model. However, the ANN model was  
886 inefficient until 1985, when backpropagation was first employed in ANN (Holyoak  
887 1987). ANN became popular after 2006 when Hinton (2007) proposed the concept of  
888 "deep learning," an architecture with many more layers than ANN. Hinton (2007) also  
889 proved that backpropagation works efficiently in multilayer ANN learning.

890       Although deep learning becomes very popular and a variety of implementations were  
891 developed since then, one significant breakthrough of deep learning was Alexnet deep  
892 learning model (Krizhevsky et al. 2012), which won the first prize and achieved exciting  
893 accuracy in ImageNet 2012 challenge and marked the start of the broad implementation  
894 of deep learning. Alexnet was the first end to end deep objective classification learning  
895 system and achieved 15.3% top-5 classification error rate for the ImageNet 2012  
896 challenge. Later work such as VGG (Simonyan and Zisserman 2014) and GoogleNet  
897 (Szegedy et al. 2015), which reach 7.3% and 6.7% top-5 error rate separately, are derived  
898 from Alexnet.

899       All of those works are based on Convolutional Neural Network (CNN), one of the  
900 significant components in deep learning that extracts features, i.e., variables, directly  
901 from pixel-based images. CNN is an ANN-based network that is mainly used for  
902 processing natural images with 3 RGB channels, and it significantly outperforms all other  
903 data mining techniques (Krizhevsky and Hinton 2012; Simonyan and Zisserman 2014;  
904 Szegedy et al. 2015). To be specific, CNN can be viewed as a 2D version of ANN, where  
905 the one dimensional hidden layer is replaced by multiple 2D layers.

906 In addition to the astonishing accuracy in image object classification, CNN is  
907 successfully applied in other areas like text classification (Karpathy et al. 2014; Lai et al.  
908 2015), sentiment analysis (Santos and Gatti 2014), and extreme weather prediction (Liu  
909 et al. 2016; Racah et al. 2016).

910 Liu et al. (2016) built an Alexnet alike CNN model to classify three extreme types of  
911 weather, TCs, atmospheric rivers, and weather fronts based on the CAM5.1 historical  
912 run, ERA-interim reanalysis, 20-century reanalysis, and NCEP-NCAR reanalysis data.  
913 The overall accuracy achieves more than 88%, and the TC detection rate reaches 98%.

914 Although regular CNN achieves excellent accuracy in tasks like image classification,  
915 CNN cannot handle problems with temporal information involved. Tran et al. (2015)  
916 proposed a 3D CNN aiming at handling video analysis problem by adding another  
917 temporal dimension on to CNN.

918 The large-scale ERA-Interim dataset consists of 14 variables of 33\*33 gridded data  
919 with a temporal coverage from the previous 18 hours to the current time and 37 pressure  
920 levels. A 3D CNN can be used to extract features from each individual variable in such  
921 an arrangement. The 37 pressure levels are viewed as the 37 channels similar to RGB  
922 channels of video, the gridded data as images, and the temporal coverage as image  
923 sequence of a video.

924 Another important structure of deep learning is the auto-encoder network, which is  
925 “a type of ANN that is trained to attempt to copy its input to its output. Internally, it has a  
926 hidden layer that describes a code used to represent the input. The network may be

927 viewed as consisting of two parts: an encoder represents a feature extracted process and a  
928 decoder that produces an input reconstruction” (Wei et al. 2018). Auto-encoder is used  
929 for dimension reduction when the original data space dimension is too large and is also  
930 used for classification and prediction (Gogna and Majumdar 2019).

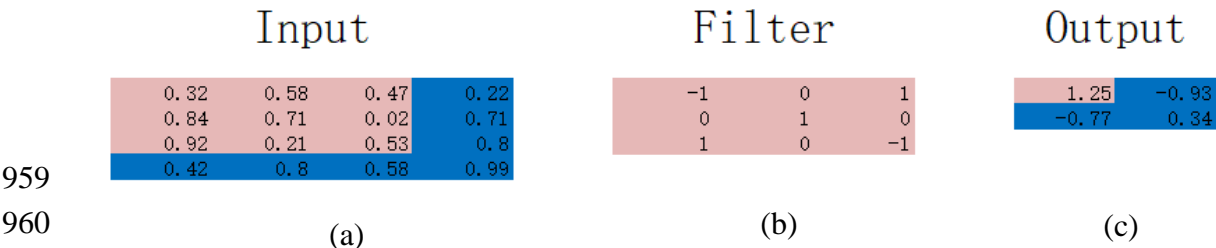
931 Racah et al. (2016) proposed an auto-encoder CNN architecture for a semi-  
932 supervised classification on the extreme weather. Since there are a large number of  
933 unlabeled extreme weather images, and to expand the training dataset, Racah et al. (2016)  
934 employed a bounding box technique to recognize the location of extreme weather, and  
935 the classification is based on those data. Although the classification performance in  
936 Racah et al. (2016) still needs improvement, it reveals that there are many promises to  
937 consider deep learning techniques in the weather community.

938 And in this work, the standard deep learning model is used to filter the large-scale  
939 ECMWF ERA-Interim reanalysis data, the data associated with all grids in Figure 3.2.  
940 Each instance has  $4(-18h, -12h, -6h, 0h) \times 37(\text{pressure level}) \times 33(\text{grid vertical}$   
941  $\text{dimension}) \times 33(\text{grid horizontal dimension})$  dimensions (values) (instead of  $4 \times 37 \times 14$  in  
942 LLE-SHIPS model), and the values are scaled to between 0 and 1 again, as did for LLE.

943 In a 2D convolutional layer, the same learnable filter is applied to each group of  
944 nearby pixels to extract features. The filter is defined as a  $p \times q$  ( $p, q$  are integers) size  
945 rectangle that can be convolved through the entire input array with the  $m \times n$  dimension.  
946 The dot product is computed between the filter weights and the input, and producing an  
947  $(m-p+1) \times (n-q+1)$  output array after scanning assuming a stride of 1. Figure 3.3 displays

948 an example of the convolution operation. A 3\*3 filter is convolved through a 4\*4 array  
 949 and output a 2\*2 array with values calculated by the dot product of the sliding filter and  
 950 the original data value. If the input array has more than one channel, as in a natural image  
 951 with RGB channels, there will be the 3rd dimension (depth) added to the previous two-  
 952 dimension filter, and the output array will still be two-dimension with value summing  
 953 over the depth dimensions. Figure 3.4 shows a multi-channel example with a 4\*4 image  
 954 with 3 channels (Figure 3.4a). A three-dimension filter (Figure 3.4b) are designed and  
 955 each is applied to the corresponding channel, and the result will be 3\*2\*2 outputs (Figure  
 956 3.4c). Then these 3 outputs will be simply summed up together, leading to one 2\*2 output  
 957 (Figure 3.4d).

958



959  
 960  
 961  
 962  
 963  
 964  
 965  
 966  
 967  
 968  
 969  
 970

**Figure 3.2: Demonstration of the convolution operation. (a) a 4\*4 array, (b) a 3\*3 filter and its weights, and (c) the resulting output array.**

Input	Filter	Output	Final output																																	
<table><tr><td>0.32</td><td>0.58</td><td>0.47</td><td>0.22</td></tr><tr><td>0.84</td><td>0.71</td><td>0.02</td><td>0.71</td></tr><tr><td>0.92</td><td>0.21</td><td>0.53</td><td>0.8</td></tr><tr><td>0.42</td><td>0.8</td><td>0.58</td><td>0.99</td></tr></table>	0.32	0.58	0.47	0.22	0.84	0.71	0.02	0.71	0.92	0.21	0.53	0.8	0.42	0.8	0.58	0.99	<table><tr><td>-1</td><td>0</td><td>1</td></tr><tr><td>0</td><td>1</td><td>0</td></tr><tr><td>1</td><td>0</td><td>-1</td></tr></table>	-1	0	1	0	1	0	1	0	-1	<table><tr><td>1.25</td><td>-0.93</td></tr><tr><td>-0.77</td><td>0.34</td></tr></table>	1.25	-0.93	-0.77	0.34	<table><tr><td>5.34</td><td>3.88</td></tr><tr><td>3.45</td><td>7.27</td></tr></table>	5.34	3.88	3.45	7.27
0.32	0.58	0.47	0.22																																	
0.84	0.71	0.02	0.71																																	
0.92	0.21	0.53	0.8																																	
0.42	0.8	0.58	0.99																																	
-1	0	1																																		
0	1	0																																		
1	0	-1																																		
1.25	-0.93																																			
-0.77	0.34																																			
5.34	3.88																																			
3.45	7.27																																			
<table><tr><td>0.11</td><td>0.58</td><td>0.47</td><td>0.95</td></tr><tr><td>0.76</td><td>0.45</td><td>0.18</td><td>0.71</td></tr><tr><td>0.92</td><td>0.21</td><td>0.53</td><td>0.8</td></tr><tr><td>0.78</td><td>0.8</td><td>0.58</td><td>0.99</td></tr></table>	0.11	0.58	0.47	0.95	0.76	0.45	0.18	0.71	0.92	0.21	0.53	0.8	0.78	0.8	0.58	0.99	<table><tr><td>-1</td><td>0</td><td>2</td></tr><tr><td>3</td><td>1</td><td>0</td></tr><tr><td>1</td><td>0</td><td>-1</td></tr></table>	-1	0	2	3	1	0	1	0	-1	<table><tr><td>2.66</td><td>2.3</td></tr><tr><td>2.2</td><td>3.8</td></tr></table>	2.66	2.3	2.2	3.8					
0.11	0.58	0.47	0.95																																	
0.76	0.45	0.18	0.71																																	
0.92	0.21	0.53	0.8																																	
0.78	0.8	0.58	0.99																																	
-1	0	2																																		
3	1	0																																		
1	0	-1																																		
2.66	2.3																																			
2.2	3.8																																			
<table><tr><td>0.18</td><td>0.41</td><td>0.09</td><td>0.18</td></tr><tr><td>0.63</td><td>0.07</td><td>0.37</td><td>0.41</td></tr><tr><td>0</td><td>0.19</td><td>0.04</td><td>0.6</td></tr><tr><td>0.07</td><td>0.01</td><td>0.2</td><td>0.4</td></tr></table>	0.18	0.41	0.09	0.18	0.63	0.07	0.37	0.41	0	0.19	0.04	0.6	0.07	0.01	0.2	0.4	<table><tr><td>0</td><td>1</td><td>1</td></tr><tr><td>2</td><td>2</td><td>0</td></tr><tr><td>4</td><td>0</td><td>0</td></tr></table>	0	1	1	2	2	0	4	0	0	<table><tr><td>1.43</td><td>2.51</td></tr><tr><td>2.02</td><td>3.13</td></tr></table>	1.43	2.51	2.02	3.13					
0.18	0.41	0.09	0.18																																	
0.63	0.07	0.37	0.41																																	
0	0.19	0.04	0.6																																	
0.07	0.01	0.2	0.4																																	
0	1	1																																		
2	2	0																																		
4	0	0																																		
1.43	2.51																																			
2.02	3.13																																			
(a)	(b)	(c)	(d)																																	

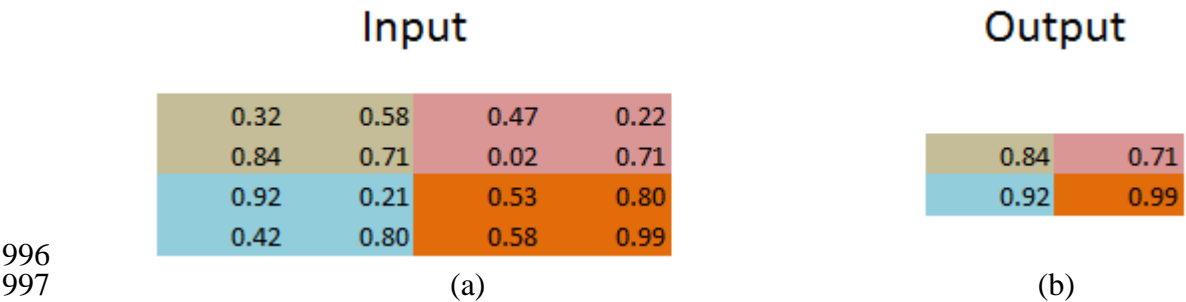
**Figure 3.3: Convolution operation for input with multiple channels. (a) 3\*4\*4 arrays, (b) a three-dimension filter, (c) the output arrays after filtering the 3 channels one by one, and (d) the final output array with values being the sum of values on the depth dimension.**

When the input is 3D arrays, the 3D filter and its convolve operation are the same to that of 2D except that an additional dimension is added. To be specific, a 3D  $p*q*r$  filter is used to extract 3D information from the  $m*n*o$  input, and result in an  $(m-p+1)*(n-q+1)*(o-r+1)$  output array assuming stride of 1 for all dimensions. Multiple filters may be applied to the same data to extract different levels of information.

The above described convolution procedure only extracts linear information, and for obtaining nonlinear information, an activation layer is introduced after each convolutional layer. Rectified Linear Units (ReLU) is the most commonly used activation function that maps negative values to 0, and keeps the positive values, respectively. This function will not affect the size of the data arrays.

A pooling layer is usually applied after the convolution and activation transformation to reduce the input's dimension in order to avoid overfitting, and unlike convolution,

there's no overlap in pooling operations for each pooling layer. Max pooling is the most widely used pooling method. Figure 3.5 displays the 2D maximum pooling; the maximum value is taken from each block (2\*2) of the original image (4\*4) and generate a 2\*2 pooled image. Similarly, 3D max-pooling layer uses pooling operation in 3D space, where all 3 dimensions are reduced simultaneously.



**Figure 3.4: Max pooling example. A 4\*4 image is sampled by a 2\*2 max pooling. (a) the original image, and (b) the pooled image.**

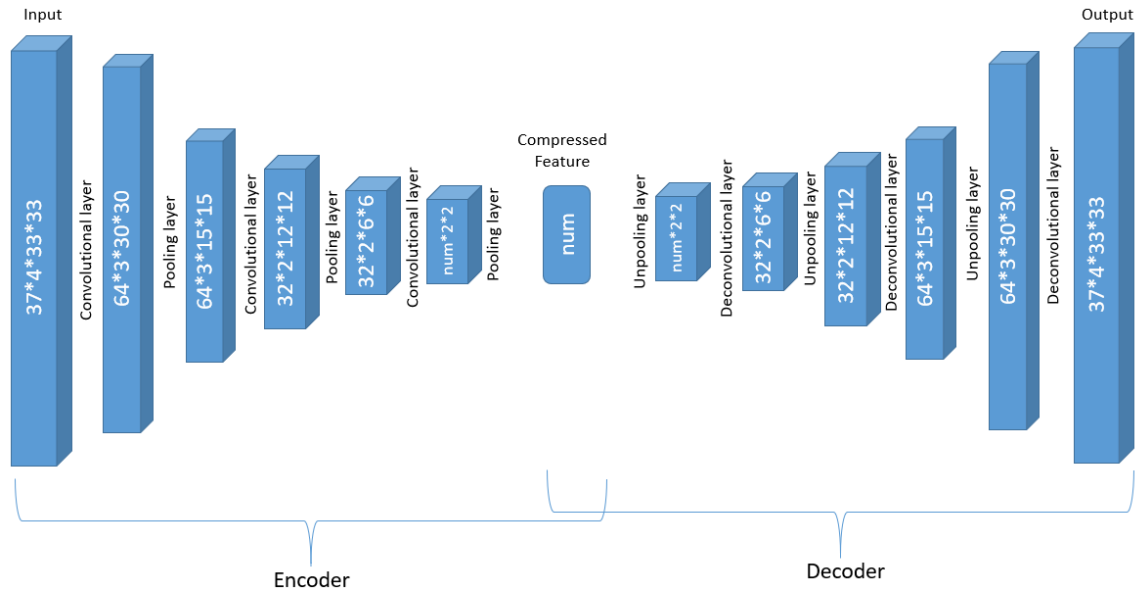
There are various types of deep learning models, and the most appropriate model for converting the gridded data into features for mining purposes is the auto-encoder network. Each auto-encoder network is composed of multiple deep learning layers, which is divided into two parts: an encoder represents a feature extraction process from the input and a decoder that reconstructs the input.

With the 14\*4\*37\*33\*33 dimensional ERA-Interim data, a more efficient auto-encoder network is a 3D Conv-auto-encoder. That is an auto-encoder with a group of 3D convolutional, activation (ReLU), and pooling layers. In each 3D convolutional layer,

there are multiple 3D convolutional filters with learnable weights with an additional channel dimension on the input channels. Moreover, the 14 variables in ERA-Interim data are treated differently than usual spatial or temporal dimension, and therefore, 14 different 3D Conv-auto-encoders are adopted to handle the ERA-Interim data.

To be specific, the input of the encoder are observations with dimension of  $37*4*33*33$ , with pressure level (37) as its channel. There are 14 such auto-encoder networks.

The network working on a single variable is elaborated below in detail, and the dimension changes of the data are displayed in Figure 3.6.



**Figure 3.5: Dimension changes of the ERA data through the 3D CNN auto-encoder layers.**



- 1023       • The first convolution layer is with 64 different  $37(\text{channel}) \times 2 \times 4 \times 4$  filters and  
1024       converts the  $37 \times 4 \times 33 \times 33$  array for one variable to 64  $3 \times 30 \times 30$  arrays. In other  
1025       words, a  $37 \times 2 \times 4 \times 4$  filter is applied and the results are summed up in the channel  
1026       dimension (37), and therefore the vertical pressure layer dimension number is  
1027       reduced to 1. This procedure is repeated 64 times with different convolution  
1028       weights. Therefore, after the first convolution layer, the original  $37 \times 4 \times 33 \times 33$   
1029       array becomes 64  $3 \times 30 \times 30$  arrays. The activation applications after each filter in  
1030       the convolution layer do not change the array size and therefore are not shown in  
1031       Figure 3.6.
- 1032       • A  $1 \times 2 \times 2$  pooling layer converts the 64 arrays with dimension  $3 \times 30 \times 30$  to 64  
1033       arrays with dimension  $3 \times 15 \times 15$ .
- 1034       • The second convolution layer has 32 different filters with dimensions  $64 \times 2 \times 4 \times 4$ ,  
1035       and in this layer, the new dimension due to 64 different filters in the previous  
1036       layer is considered as “channels” and the filtered arrays will be summed over the  
1037       channel dimension. As a result, each of the 32 filters converts the  $64 \times 3 \times 15 \times 15$   
1038       array to 1 array with reduced dimensions  $2 \times 12 \times 12$  with the same operation as that  
1039       of the first convolution layer, and finally there are 32 such arrays.
- 1040       • The same  $1 \times 2 \times 2$  pooling layer is applied to the 32  $2 \times 12 \times 12$  arrays and that results  
1041       in 32  $2 \times 6 \times 6$  arrays.
- 1042       • Similar to the previous two convolutional layers, the third convolution layer has  
1043       num different convolution filters  $32 \times 2 \times 5 \times 5$ , and the dimension of 32 is treated as

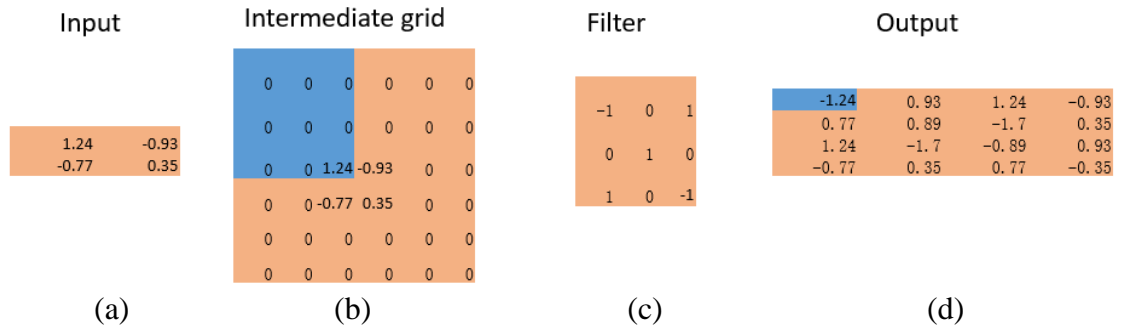
1044 channel again. The result after this filtering process is num arrays of dimension  
1045  $1*2*2$ , where the num is a to be determined hyperparameter.

1046 • The  $1*2*2$  pooling layers will finally compress the arrays into num scalar  
1047 features.

1048 The decoder is the reverse of the encoder by using the deconvolutional and unpooling  
1049 layers in DeConvNet network (Zeiler et al. 2010, 2011, 2013) to reconstruct the  
1050 convolutional networks, i.e., reverse the convolution and the pooling operations. In  
1051 deconvolutional, 0s are padded to the neighbor of the input (output of the corresponding  
1052 convolution operation) to generate an intermediate grid, and a learnable filter is used to  
1053 convolve through the intermediate grid to generate the output, which is the reconstruction  
1054 of the convolutional input. The learnable filter in deconvolution is updated in the same  
1055 way as the learnable filter in convolution.

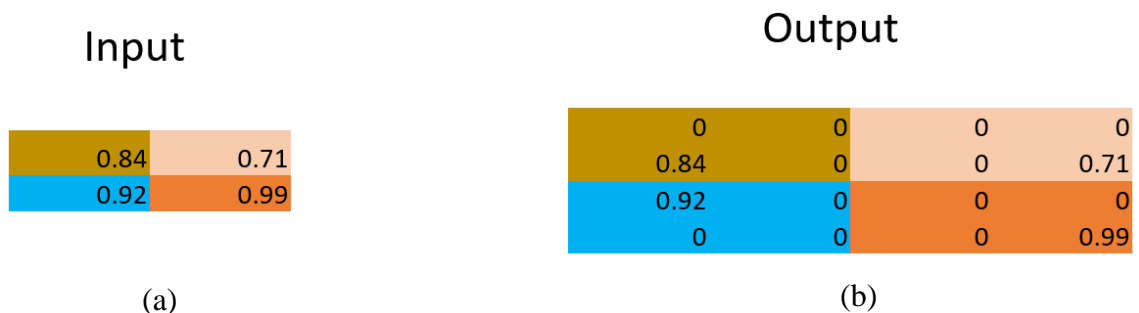
1056 An example is displayed in Figure 3.7, which is the deconvolution operation to  
1057 reverse the convolution operation described in Figure 3.3. To reverse the  $2*2$  image to  
1058 the  $4*4$  image, the original  $2*2$  image (a) is padded with 2 rows and 2 columns 0 around  
1059 each pixel to generate the intermediate grid (b), which is then convolved through the  
1060 filter (c), and result in the output (d), which is the upsampled result with regard to (a).

1061



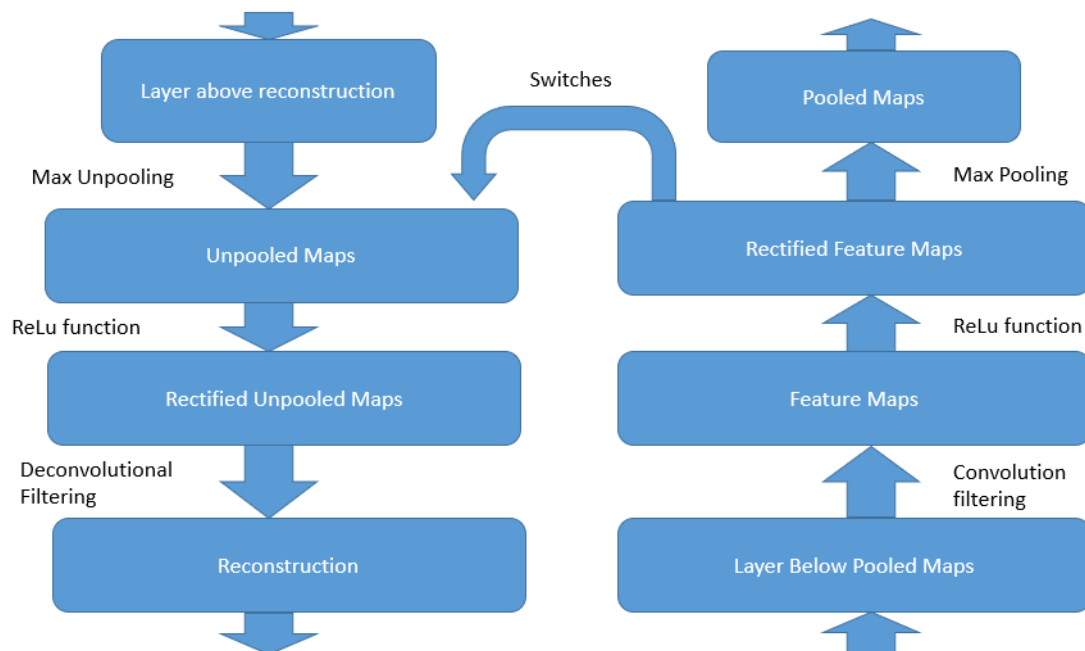
**Figure 3.6: The deconvolution operation, reverse of operations shown in Figure 3.3. (a) a 2\*2 array, (b) padded 0 to (a), (c) the filter, and (d) resulting output array after filtering.**

To reconstruct the Max pooling operation, the location of the feature map that has the maximum value (location of the passed value through Max pooling) is recorded in a switch during the corresponding Max pooling operation. Then the input of the Unpooling is upsampled where the maximum value is put to the saved position in the switch, and 0 is put into everywhere else. An example is displayed in Figure 3.8, which is used to reconstruct the output of Figure 3.5. ReLu function in the convolutional network is the same as in the DeConvNet network.



**Figure 3.7: An unpooling example.** A 2\*2 image is upsampled by a unpooling process to a 4\*4 image. The position of the valid value (nonzero) in each 2\*2 sub-image is based on the location of maximum value during pooling (see Figure 3.5 for details). (a) the original image, and (b) the unpooled image.

The structure of DeConvNet network is displayed on the left in Figure 3.9, which is used to reconstruct the corresponding CNN.



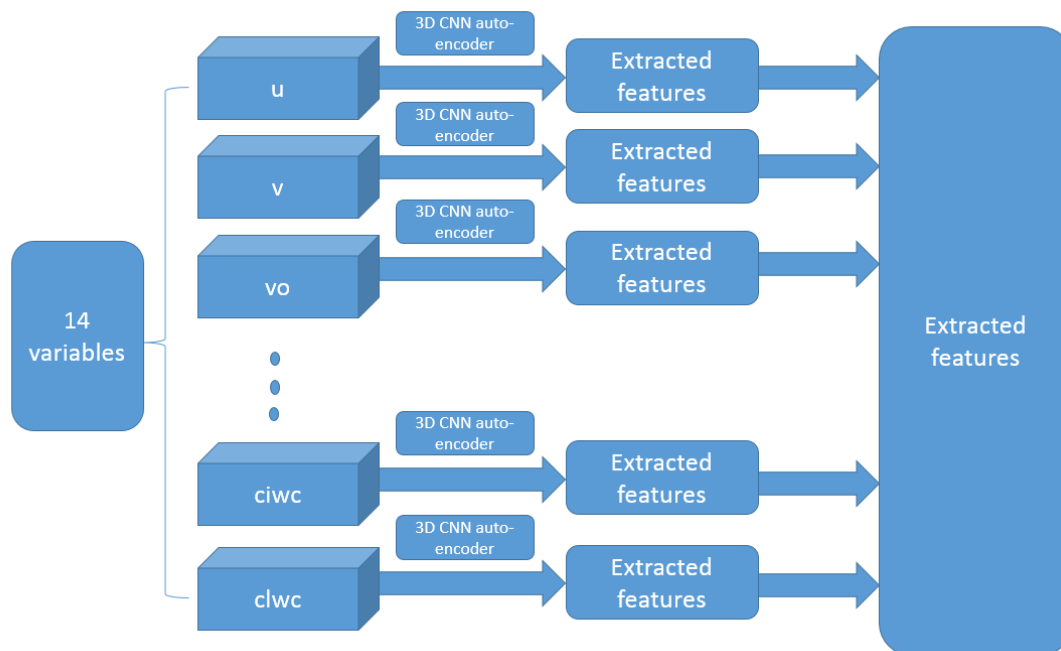
1088 **Figure 3.8: The CNN on the right that first runs through the input to the output**  
1089 **(from bottom to top), and the position of the Max Pooling pixel is saved as a switch**  
1090 **that will be used later for the unpooling operation on the left. On the left, the**  
1091 **structure of the corresponding DeConvNet network consisting of one unpooling**  
1092 **layer, one ReLu function (the same as in CNN), and one deconvolution layer to**  
1093 **reconstruct the CNN on the right based on its output (top to bottom) (Zeiler et al.**  
1094 **2011).**  
1095

1096 The network is trained through the backpropagation, where the mean square error  
1097 (Trevor et al. 2009) is used as the loss function, and Adam optimizer (Ruder 2016) is  
1098 used as the optimizer to update the filter weights through backpropagation.

1099 14 separate networks with the same structure displayed in Figure 3.6 are trained  
1100 separately for 14 different ERA-Interim variables as shown in Figure 3.10. The  
1101 compressed features from each of the networks are merged with filtered SHIPS variables  
1102 and together used as the input of the GMM-SMOTE.

1103

1104



**Figure 3.9: Combined deep learning filters for the 14 variables in ERA-Interim data.**

1110

## CHAPTER 4 GMM-SMOTE SAMPLER

1111           In a binary classification (prediction) problem such as the RI vs. non-RI  
1112 prediction, we use the machine learning (ML) model to look for a decision boundary in  
1113 the feature space (Friedl and Brodley 1997) to separate RI and non-RI instances, and the  
1114 prediction is made for any new instance based on its location in the feature space.

1115           The decision boundary in the feature space will be highly skewed if the data are  
1116 with highly imbalanced class samples. Unfortunately, the RI and non-RI instances are  
1117 highly imbalanced because only about 5% of TC instances undergo RI process.  
1118 Generally, three types of techniques are used to handle the imbalanced data problem,  
1119 algorithm level approach, cost-sensitive approach, and resampling data approach (Last et  
1120 al. 2017).

1121           The algorithm level approach aims at modifying ML algorithm that applies to  
1122 regular balanced data to cope with imbalanced data to correct the skewed decision  
1123 boundary. Such techniques, including changing the decision threshold and training a  
1124 separate model (Anyfantis et al. 2006; Chawla et al. 2004; Galar et al. 2012). Cost-  
1125 sensitive approach assigns different costs for incorrectly classifying different classes in  
1126 the ML model to correct the skewed decision boundary. In other words, cost-sensitive  
1127 approach assigns a lower cost for misclassifying majority class and a higher cost for  
1128 misclassifying minority class (Galar et al. 2012; Castro et al. 2013; López et al. 2015).  
1129 Instead of making modifications of the ML algorithm structure, the resampling data

1130 approach resamples the imbalanced dataset to create a balanced one to decrease the effect  
1131 of the skewed distribution in the ML model's learning process (Krawczyk et al. 2014).

1132         At present, resampling is the most widely used approach to overcome the  
1133 imbalanced data problem. The resampling approach falls into three categories:  
1134 upsampling, downsampling, and the hybrid method (Last et al. 2017). The upsampling  
1135 approach upsamples the data by duplicating observations in minority class until the  
1136 number of observations in minority class matches that of the majority class, and the  
1137 downsampling removes additional observations in the majority class (Japkowicz 2000).  
1138 However, simply upsampling or downsampling does not significantly improve the  
1139 minority class prediction accuracy because they do not fortify the decision boundary. The  
1140 hybrid method combines both of them by generating new instances different from  
1141 existing ones for the minority class and removing the majority class instances  
1142 simultaneously. As an approach that combines upsampling and downsampling  
1143 approaches to improve their drawbacks, the hybrid method should be used. Among all the  
1144 hybrid resampling approaches, Synthetic Minority Over-sampling Technique (SMOTE)  
1145 has been widely employed by researchers and scientists to solve the real-world problem  
1146 and academy problem due to its simplicity and its advantages to random sampling  
1147 (Shaiba and Hahsler 2016). SMOTE was proposed by Chawla et al. (2002) to handle the  
1148 imbalanced dataset, which upsamples minority classes by constructing "synthetic"  
1149 examples rather than upsampling with replications and outperforms upsampling and  
1150 downsampling alone (Akbari et al. 2004; Batista et al. 2004; Liu et al. 2006).



1151           The decision boundary that separates RI and non-RI instances in SMOTE may not  
1152 be enforced since instances in minority class far from the decision boundary have the  
1153 same probability of being selected as those closed ones. In addition, SMOTE may further  
1154 amplify the noise present in the data (Nguyen et al. 2011). When all the instances in the  
1155 minority class have an equal probability of being selected, those noise observations may  
1156 be amplified and hence decrease the accuracy of the model (Bunkhumpornpat et al.  
1157 2009).

1158           To decrease the influence of the noise observations, Han et al. (2005) proposed  
1159 two approaches, SMOTE1 and SMOTE2, to improve SMOTE by splitting minority class  
1160 instances into three groups, i.e., noise, safe, and danger using a k-nearest neighbor  
1161 approach. An instance is regarded as noise if all its neighbors are belonging to the  
1162 majority class and as safe if more than half of the neighbors are belonging to the minority  
1163 class. Otherwise, that instance is regarded as danger. Noise instances are useless because  
1164 they do not provide information about the minority class, and so do safe instances, in that  
1165 no matter what classification model is used, they are less likely to be misclassified.  
1166 Therefore, augmenting danger instances could be the most efficient way to increase  
1167 classification accuracy. Therefore, approaches proposed by Han et al. (2005) only  
1168 augment danger minority instances.

1169           To avoid the influence of the extreme values in the sampling process, and to make  
1170 the resample process more efficient (Jo and Japkowicz 2004), Song et al. (2016)  
1171 proposed a bi-directional sampling approach, where minority and majority classes are  
1172 separately clustered using K-Means. The majority class is downsampled by only selecting

1173 instances near the cluster center, and the minority instances are upsampled by SMOTE  
1174 using instances in the same cluster. The downsampling and upsampling processes are  
1175 replicated multiple times until instances in majority class and minority class are balanced.  
1176 Last et al. (2017) proposed K-means SMOTE by first clustering the entire population into  
1177 different clusters. Then only clusters with more than a certain ratio of minority class  
1178 instances are selected, and each cluster is assigned with a weight equals to the number of  
1179 minority class elements divided by the sum of their distance to the center of that cluster.  
1180 Then minority instances are augmented the number of times proportional to their weight.

1181         However, K-means is not working efficiently on complex geometrical shaped  
1182 data, especially in a high dimensional space. Furthermore, K-means SMOTE does not  
1183 handle missing values in attributes. Finally, clustering is an unsupervised approach that  
1184 the selection of the number of clusters is very subjective. If too few clusters are specified,  
1185 underfitting may occur - to cluster apparently different instances into one cluster hence  
1186 unable to identify the difference. If too many clusters are specified, overfitting may occur  
1187 - instances that have a similar property could be clustered into different clusters.

1188         To better fit SMOTE based approach to the high-dimensional data, here, the  
1189 Gaussian Mixture Model (GMM) with a weighted Euclidean distance is used for  
1190 clustering (Friedl and Brodley 1997). As a type of model-based clustering approach,  
1191 GMM has been used in a large number of areas, such as speech recognition (Reynolds et  
1192 al. 2000), and feature extraction (Torres-Carrasquillo et al. 2002).

1193         However, the number of clusters ( $M$ ) in GMM should be defined before fitting the  
1194 model to the data. Bayesian Information Criterion (BIC) (Volinsky and Raftery 2000) is a

1195 statistic calculated for each clustered dataset based on the likelihood to identify enough  
 1196 clustering information while avoiding overfitting as much as possible. BIC can be used to  
 1197 help select the best number of clusters, which is defined as  $n\_cluster$ , a to be tuned  
 1198 hyperparameter in GMM. The search space for  $n\_cluster$  is defined as 2 to 10, which  
 1199 implies BIC will be calculated for each  $n\_cluster$  starting from 2 iteratively, and the  
 1200 process stops when BIC stops decreasing for two continuous iterations or  $n\_cluster$   
 1201 equals 10 is reached.

1202 GMM is an unsupervised approach that assumes each observation in the  
 1203 population can be represented as a weighted sum of several (the number of pre-  
 1204 determined clusters) Gaussian distributions, and the weights are summed to 1. Each  
 1205 cluster corresponds to one Gaussian distribution, and observation will be assigned to the  
 1206 cluster with the highest weight (Fraley and Raftery 1998), which is unlike clustering  
 1207 approaches such as K-means, assigning each observation to a different cluster directly.

1208 To be specific, assume  $x$  is an observation in population  $X$  with  $D$  feature  
 1209 dimensions,  $M$  is the number of Gaussian distribution (clusters) within  $X$ ,  $w_i$  is the mixed  
 1210 weight for each component and sum to 1, and  $g(x|u_i, \Sigma_i)$  is the  $i$ -th Gaussian distribution  
 1211 with mean  $u_i$  and standard deviation  $\sigma_i$ .  $u_i$  and  $\sigma_i$  is calculated by expectation-  
 1212 maximization (E-M) algorithm, which is a likelihood based approach that starts from  
 1213 some initial estimates and stops until convergence arrives (Dempster et al. 1977).

$$p(x|\lambda) = \sum_{i=1}^M w_i g(x|u_i, \Sigma_i) \quad (4.1)$$

$$g(x|u_i, \sigma_i) = \frac{1}{2\pi^{1/2}|\sigma_i|^{1/2}} \exp \left\{ -\frac{1}{2} (x - u_i)^2 \sigma_i^{-1} \right\} \quad (4.2)$$

1214 The GMM clustering and the SMOTE sampling processes is combined in this  
 1215 study as the data sampler to resample the unbalanced RI dataset, and the data sampler  
 1216 employs the BIC to determine the number of clusters for GMM from a predefined range  
 1217 (searching space) between 2 and 10; then GMM is used to cluster all the instances. Every  
 1218 cluster is defined as safe, noise, or dangerous based on the instance imbalance rate (IIR),  
 1219 i.e., the ratio of the minority instances in the cluster divided by the ratio for the entire  
 1220 population. For example, if there are 3% minority instances in a cluster and 5% minority  
 1221 instance in the population,  $IIR = 3\%/5\% = 0.6$  for this cluster. When a cluster is  
 1222 composed of mainly majority instances, the classification of all instances would be  
 1223 majority class no matter the actual instance is the majority or not. Those clusters cannot  
 1224 make any contribution to improve the classification accuracy and are termed as noise  
 1225 clusters. Similarly, a cluster is defined as safe when its minority instances are dominant in  
 1226 the cluster and are less possible to be misclassified. In this study, 0.2 (5) of IIR value is  
 1227 set as the threshold, and any clusters with  $IIR \leq 0.2$  ( $\geq 5$ ) are termed noise (safe).  
 1228 Otherwise ( $0.2 < IIR < 5$ ), the classification of the minority is difficult, and the cluster is  
 1229 termed dangerous. Similar to Last et al. (2017), instances can also be identified as safe,  
 1230 noise, or danger based on the number of minority instances in their  $m\_neighbors$   
 1231 neighbors, but slightly different here, an instance is termed as noise (safe) if none (more  
 1232 than half) of its neighbors is in the minority class; otherwise, as danger. Only danger  
 1233 instances in dangerous clusters are upsampled with SMOTE following

$$\begin{aligned}
 1234 \quad & u_{new} = u_c + w * (u_n - u_c), \quad w \sim U(0,1) \\
 1235 \quad & \text{if } u_n \text{ is minority class; otherwise } w \sim U(0,0.5)
 \end{aligned} \tag{4.3}$$

1236 where  $u_c$  presents the selected minority class instance that needs to be augmented;  $u_n$   
1237 presents the randomly selected neighbor of  $u_c$  using  $k$  nearest neighbor in the same  
1238 cluster with  $k = k\_neighbors$ , another to-be-determined hyperparameter;  $U(0,1)$  and  
1239  $U(0,0.5)$  present random numbers with a uniform distribution between 0 and 1, and 0 and  
1240 0.5, respectively. Because the number of majority instances is much larger than that of  
1241 the minority instance, the increase of  $m\_neighbors$  will result in the increase of the  
1242 number of majority neighbors for an instance, leading to the increased possibility of the  
1243 instance is classified as noise; hence the instance is less likely to be upsampled.  
1244 Therefore, the number of instances that are upsampled are fewer, and the variety of the  
1245 upsampled instances is decreased. Similarly, smaller  $k\_neighbors$  will lead to a smaller  
1246 variety because fewer neighbors will be selected for the upsampling. Smaller (larger)  
1247 variety represents less (more) coverage in the feature space and would more likely result  
1248 in an underfitted (overfitted) model, or a conservative (aggressive) model. Therefore,  
1249 large (small)  $m\_neighbors$  and small (large)  $k\_neighbors$  will lead to conservative  
1250 (aggressive) models.

1251 Each danger instance in dangerous clusters is augmented  $N_a$  times using equation  
1252 (4.3) where  $N_a$  is defined as the integer part of

$$1253 \quad 0.75 * \frac{\text{The number of majority instance in the population}}{\text{The total number of minority instances need to be augmented}} - 1,$$

1254 which makes the final number of minority instances are approximately 75% of the  
1255 majority instance number, assuming most of the minority instances will be augmented.  
1256 Then, as the final step of the resampling, 25% instances in the majority class are

1257 randomly removed (downsampled) to make the majority class and minority class have  
1258 similar numbers of instances (Song et al. 2016; Last et al. 2017). Clusters with the  
1259 number of instances less than the maximum of `m_neighbors` and `k_neighbors` will be  
1260 removed.

1261       The GMM-SMOTE approach described above will be used to augment the output  
1262 data from the data filter of the COR-SHIPS model, LLE-SHIPS model, and DL-SHIPS  
1263 model. There are 3 hyperparameters that need to be tuned, as listed in Table 4.1. To avoid  
1264 overfitting and underfitting, i.e., not to select too large or too small `k_neighbors`, and  
1265 `m_neighbors`, the search space is defined as 3 to 14, and 3 to 10, respectively.  
1266 Furthermore, `n_cluster` should guarantee that each cluster has at least `k_neighbors+1`  
1267 instances. Therefore, the search space for `n_cluster` is defined as 2 to 10. A cluster will be  
1268 removed if the number of its total instances is less than `k_neighbors` or `m_neighbors`.

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278 **Table 4.1: Hyperparameters and their searching space in GMM-SMOTE sampling**  
 1279 **process.**

Hyperparameter	Component	Explanation	Minimum	Maximum	Initial value
n_cluster	GMM-SMOTE	The number of clusters in the Gaussian Mixture Model function	2	10	1
m_neighbors	GMM-SMOTE	The number of nearest neighbors used to determine if a minority sample is in danger	3	10	10
k_neighbors	GMM-SMOTE	The number of nearest neighbors used to construct synthetic samples	3	14	5

1280

## CHAPTER 5 XGBOOST CLASSIFIER AND HYPERPARAMETER TUNING PROCESS

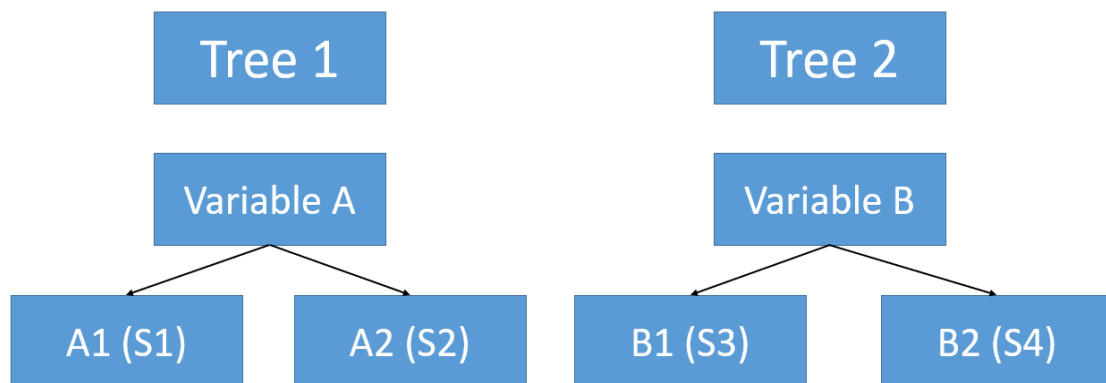
### 5.1 XGBoost classifier

The classification and regression tree model (CART) (Friedl and Brodley 1997) is one of the most popular classifications and prediction models in the machine learning community that is capable of capturing the linear and nonlinear relationship between predictors. However, a single CART usually is a weak classifier, only slightly better than random guessing. Therefore, a boosting approach is commonly used, which trains a set of weak classifiers to enhance the classification performance, with new classifiers being trained to correct the mistake made by the previous classifiers. The performance of the boosting model with a large set of weak classifiers usually outperforms the single strong classifier.

State of the art boosting tree model is the Gradient Boosting Tree (GBT) (Trevor et al. 2009), which trains a series of weak CART classifiers iteratively. Each classifier is constructed on the remaining errors of previous classifiers, and new classifiers are trained to correct the mistake made by the previous classifiers. The output of each classifier is a leaf level score, as shown in Figure 5.1, as an example for RI prediction. An instance with attributes A and B lands at leaf A1 of Tree 1 based on the A value, and leaf B1 of Tree 2 based on the B value. The raw classification score for being a RI instance is a weighted sum,  $w_1 * S1 + w_2 * S3$ , with  $w_1$  and  $w_2$  being the weights for Tree 1 and Tree 2, respectively. The raw score is rescaled to a value between 0 and 1 using a sigmoid



function (Trevor et al. 2009). In the binary classification in this study, where 0 and 1 are used to encode non-RI and RI classes, if the rescaled score is less than a decision threshold, which is a to-be-tuned hyperparameter, and preselected as 0.5, the instance is predicted as non-RI. Otherwise, the instance is predicted as RI. A greedy algorithm like GBT may generate too many weak classifiers fitting in the residuals that the total model can be easily overfitted, and a regularized distributed GBT, i.e., XGBoost, is designed to control the overfitting (Chen and Guestrin 2016).



**Figure 5.1. CART sample. Two regression trees 1 and 2, split based on values of variable A and variable B, respectively. A1, A2, B1 and B2 are the leaf names and S1, S2, S3, and S4 are the corresponding classification scores.**

XGBoost generates weaker classifiers iteratively by minimizing an objective function, consisting a loss function and a regularization function. The loss function is based on the errors between the predicted classes and the ground truth classes for all instances. The regularization function is a function of classifiers, and its purpose is to control the overfitting of the final classification. In short, the strategy in XGBoost is to

1320 have the best prediction (minimum loss function) with the regularized complexity of the  
1321 tree structure (to avoid overfitting).

1322         The regularization constraints can be roughly divided into two categories. The  
1323 first category is on the overall structure outside individual classifier, which includes  
1324 shrinkage ratio (shrinkage), the number of classifiers (n\_estimator), subsample ratio  
1325 (subsample), and features ratio (colsample). The second category is on the individual  
1326 CART (classifier) level, which includes L1 regularization (reg\_alpha) and L2  
1327 regularization (reg\_lambda), minimum loss reduction required to make a split (Split  
1328 criteria, aka, gamma) and the minimum sum of instance weight in a split  
1329 (min\_child\_weight), and the maximum depth of the CART (max\_depth).

1330         In the first category, since the subsequent classifiers are iteratively fitted into the  
1331 remained error from the previous classifiers, the subsequent classifiers contribute less and  
1332 less as boosted trees go deeper. Therefore, similar to gradient descent algorithm with  
1333 decreasing steps for better approximation (Trevor et al. 2009), we decrease the  
1334 contribution of even weaker classifiers with a rate of shrinkage, a hyperparameter within  
1335 (0,1]. Empirically XGBoost was found to perform best with the shrinkage around 0.1,  
1336 and the search space is defined as 0 to 0.3 here. While more classifiers would result in  
1337 better accuracy, too many classifiers will result in overfitting. Therefore, the number of  
1338 classifiers, n\_estimator, is limited in [100, 2000] and will be searched in that range to  
1339 avoid underfitting and overfitting simultaneously. In addition to the classifier number,  
1340 overfitting can also be reduced by using only a reduced set of datasets and variables  
1341 (Trevor et al. 2009). Subsample ratio (subsample) and features ratio (colsample) are used

1342 to control the sizes of randomly selected reduced datasets and feature sets, and the search  
1343 spaces are defined as 0.5 (50% instances) to 1 and 0.4 (40% features) to 1, respectively.  
1344 In total, four constrains: shrinkage, n\_estimator, subsample, and colsample are adopted as  
1345 hyperparameters for overall constrains in the classification process.

1346 In the second category, the concern is the same, to avoid overfitting by similar  
1347 strategies but on individual trees. Although the number of features is reduced by  
1348 colsample, that process is random. The number of the features are further controlled by  
1349 L1 regularization and L2 regularization based on their importance, where L1  
1350 regularization (reg\_alpha) and L2 regularization (reg\_lambda) are similar to Ridge and  
1351 Lasso (Friedl and Brodley 1997) in linear regression but apply on CART to reduce the  
1352 impact of less-predictive features. The search spaces of L1 and L2 regularization are  
1353 specified as 0 to 20, and 0.1 to 20, respectively. Another tree-level constraint is to limit  
1354 tree growth. This can be achieved by setting the minimum loss reduction required to  
1355 make a split (gamma) and a minimum sum of instance weight in a split  
1356 (min\_child\_weight), and 0 to 10, and 0.5 to 5 are defined as their search space,  
1357 respectively. Finally, the maximum depth (max\_depth) is used to control the depth of a  
1358 CART and is searched in the space from 3 to 10.

1359 The last to-be-tuned hyperparameter is the decision threshold. If the summed  
1360 score output from the XGboost is above the threshold, an instance is classified as RI.  
1361 Tentatively, the decision threshold is preselected as 0.5, and will be tuned.

1362 Details of all hyperparameters are specified in Table 5.1. For a short note, lower  
1363 m\_neighbors, k\_neighbors, shrinkage, n\_estimators, subsample, colsample, max\_depth,

and higher `reg_alpha`, `reg_lambda`, `gamma`, `min_child_weight` lead to a more conservative model. `n_cluster` and decision threshold will not influence the conservativeness of the model. More details of XGBoost will be elaborated in Chapter 6.

**Table 5.1: Hyperparameters, their searching space defined by the minimum and maximum, and the initial values in GMM-SMOTE sampling process and XGBoost classifier.**

Hyperparameter	Component	Explanation	Min	Max	Initial value
<code>shrinkage</code>	XGBoost	Shrinkage ratio for each feature	0	0.3	0.1
<code>n_estimator</code>	XGBoost	The number of CART to grow	100	2000	100
<code>subsample</code>	XGBoost	Subsample ratio of the training instances	0.5	1	1
<code>colsample</code>	XGBoost	Subsample ratio of columns for creating each classifier	0.4	1	1
<code>reg_alpha</code>	XGBoost	L1 regularization term on weights	0	20	0
<code>reg_lambda</code>	XGBoost	L2 regularization term on weights	0.1	20	1
<code>gamma</code>	XGBoost	Minimum loss reduction required to make a further partition on a leaf node of the CART	0	10	0
<code>min_child_weight</code>	XGBoost	Minimum sum of instance weight in a split	0.5	5	1
<code>max_depth</code>	XGBoost	Max depth of each CART model in XGBoost	3	10	3
<code>decision threshold</code>	XGBoost	Decision threshold on the XGBoost classifier output	0	1	0.5

## 1372 **5.2 Hyperparameter tuning process**

1373 Hyperparameter tuning is based on pre-defined measures of classification  
1374 performance, and all performance measures are derived from the elements of the so-  
1375 called confusion matrix, as shown in Table 5.2. The commonly used accuracy,  
1376  $(TP+TN)/(TP+FP+FN+TN)$ , is not a good measure for the unbalanced RI cases. Instead,  
1377 Probability Of Detection (POD), False Alarm Ratio (FAR), Peirce's Skill Score (PSS),  
1378 and kappa scores are often used in RI prediction evaluations (Wilks 2011; Yang 2016;  
1379 Kaplan et al. 2015).

1380

1381 **Table 5.2: Confusion matrix.**

	Predicted positive	Predicted negative
Actual positive	Truth positive (TP)	False negative (FN)
Actual negative	False positive (FP)	Truth negative (TN)

1382

1383 POD (aka recall) is defined as  $TP / (TP + FN)$ , or the ratio of correct positive  
1384 prediction cases to all the positive cases. FAR is defined as  $FP / (TP + FP)$  measuring the  
1385 false positive prediction ratio compared to all the positive predictions. PSS, defined as  
1386  $TP / (TP + FN) - FP / (FP + TN)$ , can be interpreted as the sum of the class level accuracy.  
1387 And the kappa score, defined as  $2 * (TP * TN - FN * FP) / [(TP + FN)(FN + TN) -$   
1388  $(TP + FP)(FP + TN)]$ , conceptually the same as Brier Skill Score (BSS), measuring the  
1389 relative improvement of the prediction against the prediction based on samples without

1390 any models (Wilks 2011). In this study, a single metric, the kappa score, is mainly used  
1391 for the hyperparameters tuning. Other measures are used mostly to report the  
1392 performance of the prediction and to compare with previous studies.

1393 To find the optimal values of the hyperparameter set (maximizing the kappa score  
1394 in this study), grid search is performed in most previous works but the grid search is too  
1395 time-consuming, especially with a large number of hyperparameters in the model.  
1396 Bayesian optimization (BO) is then designed to reduce the time consumption, which uses  
1397 an iteration procedure to search the global optimum. In reality, it is difficult for BO to  
1398 find the global optimum, and instead, the BO will converge to local optima and diverge  
1399 from them during the global optimum searching process. Therefore, the iteration should  
1400 stop at a pre-defined maximum iteration number to avoid almost never-end global  
1401 optimum search (Snoek et al. 2012; Shahriari et al. 2015). In this study, BO is used with  
1402 pre-defined ranges for most of the searched hyperparameters, and the ranges are  
1403 independent of each other.

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416

1417

## CHAPTER 6 RESULT

1418           In data mining and machine learning, a dataset is usually divided into training  
1419 dataset, validation dataset, and test dataset. The training dataset is used to train the model,  
1420 the validation dataset is adopted for hyperparameter tuning, and the test dataset is  
1421 employed to provide the final assessment of the model (James et al. 2013).

1422           In this study, the whole dataset is divided into 90% for training-validation and  
1423 10% for test. To further divide the training-validation dataset into the training dataset and  
1424 validation dataset, the entire training-validation dataset is split into ten mutually exclusive  
1425 equal-sized subsets. One of the subsets is reserved as the validation dataset, and the rest  
1426 nine subsets are used for training. Each of the subsets is used in turns for validation  
1427 dataset once. This is defined as 10-fold cross-validation (Kohavi 1995). The model's  
1428 performance (kappa score) is evaluated by the mean performance (mean kappa score) on  
1429 the 10 validation datasets.

1430           When the study is first conducted with the SHIPS data, there were 11,317  
1431 instances (cases) from 1982 to 2016, and 571 (approximately 5%) were under rapid  
1432 intensification (RI). A random stratified sampling on RI and non-RI cases was drawn  
1433 with a similar proportion, and that resulted in 10,185 instances (including 523 RI cases,  
1434 5.1%) in the training-validation set, and 1,132 instances (including 48 RI cases, 4.2%) in  
1435 the test dataset.

1436           After the training was done, however, 465 instances in 2017 and the last  
1437 tropical cyclone in 2016 are added to the SHIPS developmental database, and all of these

1438 instances are added to the test dataset. The test dataset proportion is ended up with 1,597  
1439 (14.1%) instances in total with 95 RI instances (5.9%).

1440 Therefore, in this study, COR-SHIPS, LLE-SHIPS, and DL-SHIPS models are  
1441 trained using the training dataset initially. Then their hyperparameter tuning process will  
1442 be derived on the validation dataset following steps introduced in Chapter 6, and their  
1443 performance will be evaluated in the test dataset, which will be compared with previous  
1444 works in Y16 and KRD15. Finally, their variable importance will be evaluated and  
1445 discussed.

1446 All algorithms, including those for data processing, data visualization, and data  
1447 mining and machine learning in this study are performed with R (version 3.5.1), python  
1448 base (version 3.7.0), python multiprocessing package (version 2.5), scikit-learn package  
1449 (version 1.9.2), XGBoost package (version 0.83), and pyspark package (Spark API  
1450 (version 2.21). The entire process is implemented on Amazon Web service (72 cpus (3.0  
1451 GHz Intel Xeon Platinum processors), 144 G memory), and a local machine (8 i7 cores,  
1452 32G RAM, and GTX 1080).

## 1453 **6.1 COR-SHIPS model**

### 1454 **6.1.1 Hyperparameter tuning for model selection**

#### 1455 ***6.1.1.1 Hyperparameters tuning for SHIPS data filter***

1456 The structure of the COR-SHIPS model is displayed in Figure 1.1, and the  
1457 correlation threshold in the SHIPS data filter is tuned based on the trial-and-error with 4  
1458 values, 0.7, 0.8, 0.9, and 0.95. For each of the correlation threshold, variables are first  
1459 filtered as discussed before. Then Bayesian Optimization with 40 iterations is used to



tune hyperparameters in Table 4.1 and Table 5.1 with no clustering and the preset 0.5 classification decision threshold. With each iteration, a 10 cross-validation kappa score is recorded with a corresponding set of hyperparameter values, and the top 5 out of 40 kappa scores for each threshold are listed in Table 6.1. Threshold 0.95, 0.9, 0.8, 0.7 reach mean kappa scores of 0.401, 0.409, 0.411, and 0.343, respectively. This indicates the model performs the best at 0.8 and 0.9 among the four given numbers.

**Table 6.1: Kappa scores of the 5 best 10-fold cross-validation results and their means for different correlation thresholds. The name of the 1<sup>st</sup> to 5<sup>th</sup> indicates the 1<sup>st</sup> to 5<sup>th</sup> best kappa scores. The “number variables selected” is the number of variables kept after highly correlated variables removal.**

Correlation Threshold	5 <sup>th</sup>	4 <sup>th</sup>	3 <sup>rd</sup>	2 <sup>nd</sup>	1 <sup>st</sup>	Mean	Number variables selected
0.7	0.314	0.321	0.334	0.352	0.392	0.343	56
0.8	0.404	0.407	0.411	0.417	0.418	0.411	72
0.9	0.402	0.403	0.405	0.415	0.419	0.409	99
0.95	0.387	0.397	0.401	0.409	0.411	0.401	136

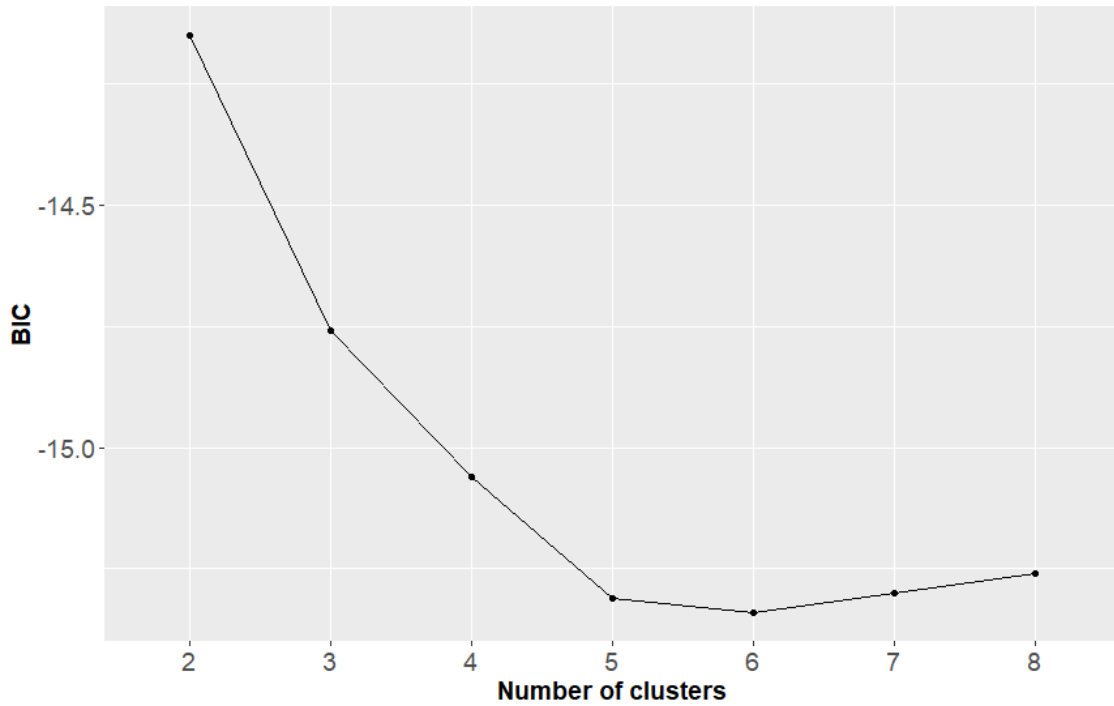
Table 6.1 also shows there are 136, 99, 72, and 56 variables left for threshold 0.95, 0.9, 0.8, and 0.7, respectively after filtering. Approximately 30 variables are reduced when the threshold is changed from 0.95 to 0.9, 0.9 to 0.8, and 0.8 to 0.7. Having less number of variables results in a lower possibility of overfitting, and the score of threshold 0.8 is higher than that of 0.9. Therefore, 0.8 is selected as the correlation threshold for removing highly correlated attributes in the SHIPS data. After the removal,

1478 there are 72 groups of highly correlated variables, as listed in Table A1, and the first  
1479 variable in each group will be selected to form the 72 selected variables.

#### 1480 ***6.1.1.2 The number of clusters selected in GMM-SMOTE***

1481 After the correlation threshold is determined, the hyperparameters for GMM-  
1482 SMOTE and XGBoost still need to be tuned for the best results. In GMM tuning for the  
1483 “optimal” cluster numbers, BIC is chosen as the selection criterion, and the BIC values  
1484 with the different number of clusters are displayed in Figure 6.1. The BIC values  
1485 decrease with the cluster number first and then increase. The BIC values decrease with  
1486 the cluster number for small cluster numbers, but stop decreasing at n\_cluster=6, and  
1487 increases when the cluster numbers increase to 7 and 8. Therefore, n\_cluster is selected as  
1488 6, associated with the lowest BIC value.

1489



**Figure 6.1: BIC ( $10^5$ ) for GMM with different number of clusters.**

The six clusters with the numbers of minority (RI) and total instances, and the IIR in each are displayed in Table 6.2. As we defined danger clusters with 0.2-5 IIR range, Clusters 3 and 5 could be excluded due to too few minority cases in the following augmentation, which removed a total of 2,401 instances with 17 RI cases (0.71%).

**Table 6.2: The number of minority and total instances, and the Imbalance Ratio (with population RI ratio at 5.1%) for the 6 clusters.**

Cluster	1	2	3	4	5	6	Total
Number of the minority instance	84	69	12	235	5	118	523
Number of the total instance	2275	1481	1255	2390	1146	1638	10185
Imbalance Rate	0.724	0.914	0.187	1.928	0.086	1.413	1

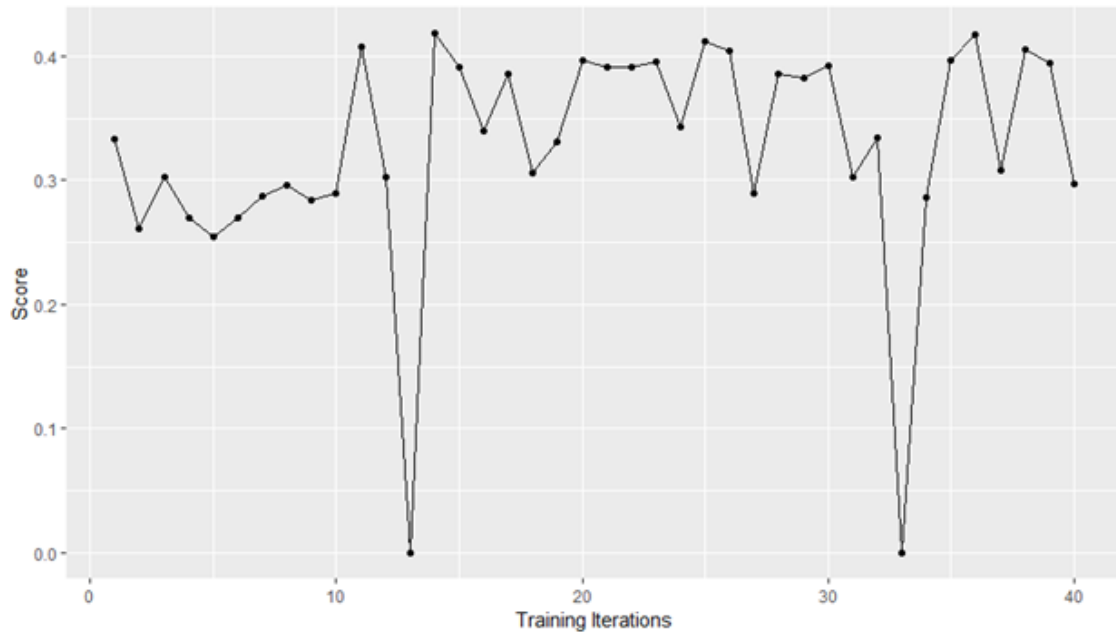
1500

### 1501 ***6.1.1.3 Hyperparameters tuning for GMM-SMOTE and XGBoost***

1502         While the correlation threshold in the SHIPS data filter and the cluster number are  
1503         tuned separately, all other eleven hyperparameters except for the decision threshold are  
1504         tuned together by Bayesian Optimization (BO). Those parameters are listed in Table 4.1  
1505         and Table 5.1 with the pre-defined searching spaces and initial values.

1506         Figure 6.2 shows the change of the 10-fold cross-validation kappa scores, the  
1507         tuning criterion, over a total 40 BO iterations. The kappa score is fluctuating over the  
1508         iterations. For example, BO process helps the model reach a local optimum at iteration 8,  
1509         and diverges from the local optimum to look for the global optimum and reach another  
1510         local optimum after iteration 11. This process continues until the global optimum is  
1511         found, which is barely possible, or the maximum iteration, 40 preset in this case, is  
1512         reached. Since the trend with the iteration is unpredictable, hyperparameter sets with the  
1513         best 5 kappa scores are selected, and the scores and hyperparameter values are displayed  
1514         in Table 6.3. The different hyperparameter set is named as MX, where X is defined by the  
1515         iteration number. For example, M11 implies the parameter set is selected after the 11<sup>th</sup>  
1516         iteration.

1517



**Figure 6.2: Variation of Cross-validation kappa scores over Bayesian Optimization iteration numbers.**

Based on Table 6.3, the performance of M14 (0.418) and M36 (kappa=0.417) are better than M25 (0.411), M38 (0.405), and M11(0.407). As indicated in previous sections, lower `k_neighbors`, `shrinkage`, `n_estimators`, `subsample`, `colsample_bytree`, `max_depth`, and higher `m_neighbors`, `reg_alpha`, `reg_lambda`, `gamma`, `min_child_weight` will lead to a more conservative prediction model. Therefore, we will analyze the values of hyperparameters to find a balanced (not too conservative and not too aggressive) model based on the ranking of those values.

1533 **Table 6.3: Top performed hyperparameter sets, the corresponding cross-validation**  
1534 **kappa scores, and specific values of the tuned hyperparameters. The numbers after**  
1535 **“M” denoting the iteration numbers.**

Name	M38	M11	M25	M36	M14
Kappa score	0.405	0.407	0.411	0.417	0.418
m_neighbors	4	3	3	3	5
k_neighbors	6	11	9	10	10
shrinkage	0.29	0.23	0.3	0.3	0.21
n_estimators	2000	572	2000	376	1510
subsample	0.75	0.5	0.5	0.5	0.67
colsample	0.99	0.78	0.99	0.9	0.99
reg_alpha	0.5	1.34	0.5	0.5	0.5
reg_lambda	20	20	20	18.91	20
gamma	0	0	0	0	0
min_child_weight	0.5	0.5	2.12	1.26	0.91
max_depth	7	8	7	7	10

1536

1537 To look for the balanced model, a system is designed to score the  
1538 conservativeness of the hyperparameter set. The scores are based on a 1 (the least  
1539 conservative) to 5 (the most conservative) scale associated with the hyperparameter value  
1540 ranks, as listed in Table 6.4. For hyperparameters favoring smaller values for  
1541 conservativeness (k\_neighbors, shrinkage, n\_estimators, subsample, colsample,  
1542 max\_depth), the scores are the same as the descending parameter value ranks. When a tie  
1543 appears, the tied values will have the same rank (score), and the next rank value depends  
1544 on how many values tie together. For example, in k\_neighbors, M11 has the largest  
1545 value, 11, hence M11 is scored 1. M36 and M14 have the second largest value, 10, hence  
1546 they are scored 2. The next largest value, 9, is ranked the 4<sup>th</sup> (instead of the 3<sup>rd</sup>) in M25,  
1547 and is scored 4. For other hyperparameters (m\_neighbors, reg\_alpha, reg\_lambda,  
1548 gamma, and min\_child\_weight), the conservativeness ranking scores are opposite to the

1549 descending value ranks. After the ranking scores are assigned to all of the 11  
1550 hyperparameters in the five local optimal cases, the scores are summed up for the five  
1551 cases (Table 6.4). Since our goal is to choose a model neither conservative nor  
1552 aggressive, the parameter set M36 with the middle conservativeness ranking score is  
1553 chosen for following implementation and discussion.

1554

1555 **Table 6.4: The descending value ranking of individual hyperparameter among the**  
1556 **top 5 performed cases, and the corresponding conservativeness ranking scores in**  
1557 **parentheses. The variables with normal font are those favoring smaller values for**  
1558 **conservativeness, and those italicized favoring larger values.**

1559

Name	M38	M11	M25	M36	M14
<i>m_neighbors</i>	2 (4)	3 (1)	3 (1)	3 (1)	1 (5)
k_neighbors	5 (5)	1 (1)	4 (4)	2 (2)	2 (2)
shrinkage	3 (3)	4 (4)	1 (1)	1 (1)	5 (5)
n_estimators	1 (1)	4 (4)	1 (1)	5 (5)	3 (3)
subsample	1 (1)	3 (3)	3 (3)	3 (3)	2 (2)
colsample	1 (1)	3 (3)	1 (1)	2 (2)	1 (1)
<i>reg_alpha</i>	2 (1)	1 (5)	2 (1)	2 (1)	2 (1)
<i>reg_lambda</i>	1 (2)	1 (2)	1 (2)	5 (1)	1 (2)
<i>gamma</i>	1 (1)	1 (1)	1 (1)	1 (1)	1 (1)
<i>min_child_weight</i>	5 (1)	5 (1)	1 (5)	2 (4)	3 (3)
max_depth	3 (3)	2 (2)	3 (3)	3 (3)	1 (1)
<b>Total score</b>	<b>23</b>	<b>27</b>	<b>23</b>	<b>24</b>	<b>26</b>

1560

#### 1561 **6.1.1.4 Hyperparameters tuning for XGBoost decision threshold**

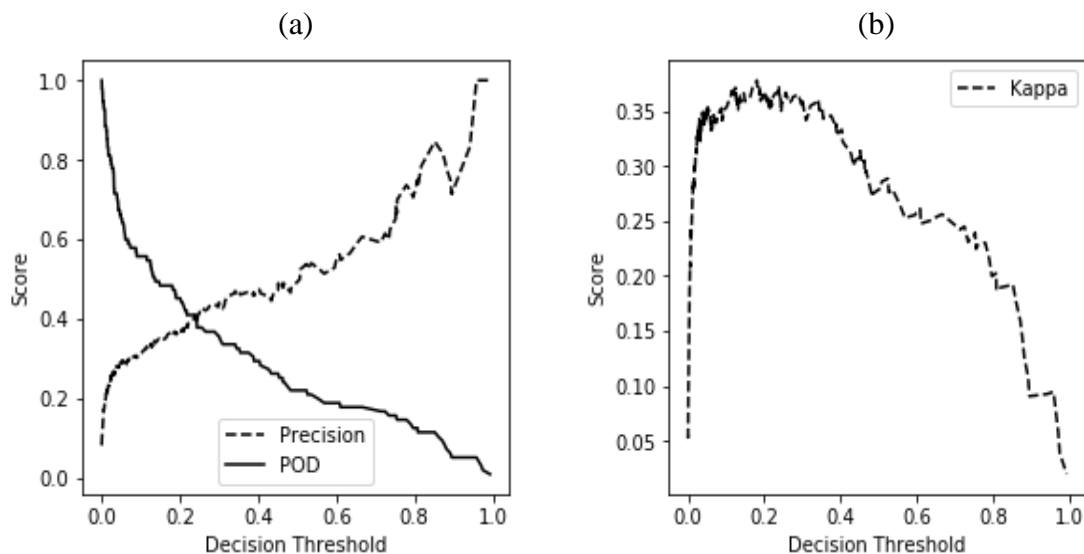
1562 The last tuned hyperparameter is the decision threshold on the XGBoost classifier  
1563 output, which was set as 0.5 before tuning. To tune this hyperparameter, we use a graphic  
1564 method based on the relative values of POD and FAR as well as the kappa score. Since

1565 POD and FAR are monotonically decreasing function with the decision threshold, we  
 1566 instead use precision (1-FAR) for identifying a threshold that balances the POD and  
 1567 FAR.

1568 Figure 6.3a displays the variations of precision and POD as functions of the  
 1569 decision threshold. The precision and POD curves cross each other around 0.2 of the  
 1570 threshold value, a relatively balanced point for POD and FAR. At the same point, the  
 1571 kappa scores shown in Figure 6.3b close to the highest value, 0.35. As a result, 0.2 is  
 1572 selected as the decision threshold, which is expected to balance the POD and FAR, and  
 1573 therefore could minimize the overfitting effect in the final classification results.

1574

1575



1576

1577

1578 **Figure 6.3: (a) Precision and POD score vs. decision threshold, and (b) Kappa score**  
 1579 **vs. decision threshold.**

1580



### 1581 6.1.2 COR-SHIPS result on test data

1582 Unlike the traditional classifier such as decision tree (Yang et al. 2016), modern  
 1583 classifier such as the XGBoost gives almost perfect classification on the training-  
 1584 validation data, i.e.,  $POD \approx 1$  and  $FAR \approx 0$ . Due to this fact, the performance measures on  
 1585 the training results are of little value, and therefore, the evaluation of the prediction is on  
 1586 the test data only.

1587 The confusion matrix for the testing data, before hyperparameter tuning (MB),  
 1588 and after hyperparameter tuning (MA) is displayed in Table 6.5. The hyperparameter set  
 1589 of MB is displayed in the last column of Table 4.1, and Table 5.1 (initial values set by the  
 1590 software) with the decision threshold as 0.5, and the MA is with the hyperparameter set  
 1591 of M36 (Table 6.3) with the 0.2 decision threshold and 6 clusters. MA's TN (1,438) and  
 1592 FN (56) are slightly smaller than 1,447 and 62 of MB, while MA's TP (39) and FP (64)  
 1593 are larger than 33 and 55 of MB. Smaller TP and larger FN in MB implies that MB is  
 1594 worse at correctly predicting RI instances, and vice versa. It seems that there is a trade-  
 1595 off between correctly predicting RI and non-RI, i.e., if we want to better predict the RI  
 1596 instance, we should sacrifice the power we predict for non-RI instances, and vice versa.

1597

1598 **Table 6.5: Confusion matrix values of our model after (before) hyperparameter**  
 1599 **tuning**

	Predicted RI	Predicted non-RI	Actual
Actual RI	39 (33)	56 (62)	95
Actual non-RI	64 (55)	1438 (1447)	1502

Predicted	103 (88)	1494 (1509)	
-----------	----------	-------------	--

1600

1601 Kappa, PSS, POD, and FAR are used for the model evaluation, and their values  
1602 for MB and MA are elaborated in Table 6.6. The POD and FAR values for MB and MA  
1603 cases demonstrated the importance of hyperparameter tuning. After tuning, POD  
1604 increases 26.1% from 0.326 to 0.411, while FAR increase only from 0.617 to 0.621,  
1605 almost nill (0.6%). That is, the benefit of higher correct RI prediction is much higher than  
1606 the cost of false alarm with the hyperparameter tuning. The overall statistics PSS and  
1607 kappa score also increased from 0.293 to 0.368 (25.6%) and from 0.315 to 0.354  
1608 (12.4%), respectively, confirming the significant improvement on RI prediction with the  
1609 hyperparameter tuning procedure.

1610

1611 **Table 6.6: Performance comparisons. MB and MA denote the models before and**  
1612 **after the hyperparameter tuning.**

Model	Kappa	PSS	POD	FAR
MB	0.315	0.293	0.326	0.617
MA	0.354	0.368	0.411	0.621
Improvement MB	12.4%	25.6%	26.1%	+0.6%

1613

### 1614 6.1.3 Feature importance

1615 Generally, the variable (feature) importance is used to leverage the variable  
1616 contribution and is defined as a quantitative score. The higher the score is, the more the  
1617 variable contributes and the more useful that variable is for classifying RI. The classifier

used in this study, XGBoost, provides the scaled importance scores with the sum of all scores being one.

Table 6.7 displays the variables with the top 10 importance scores and their definition (SHIPS 2018c). The scores of the full 72 variables are given in Table A4. The past 12-hour intensity change, BD12, has the largest importance score, 0.0362, which almost doubles the importance score of the second important variable. Because BD18 and BD06 are highly correlated with BD12 (see Table A1), we can safely assume that they are as important as BD12. The second most important variable is DTL, the distance from a TC to the nearest major land. The importance of DTL is slightly higher than the third to seventh most important variables, CFLX, SHRD, G150, jd, and VAMX, which are related to dry air, vertical wind shear magnitude at 850-200 hPa, the temperature perturbation at 150 hPa, Julian day, and the current TC intensity. The eighth to ninth variables are IRM1\_5 and PW08. The tenth most important variable is VMPI, the Maximum potential intensity, which ranked higher in other RI studies.

**Table 6.7: Features of top ten importance, their importance scores, and feature description from SHIPS (2018c) in the COR-SHIPS model.**

Variable	Importance	Description
BD12	0.0362	The past 12 hour intensity change
DTL	0.0217	The distance to nearest major land
CFLX	0.0207	Dry air predictor based on the difference in surface moisture flux between air with the observed (GFS) RH value, and with RH of air mixed from 500 hPa to the surface

SHDC	0.0206	850-200 hPa shear magnitude (kt *10) vs time (200-800 km) but with vortex removed and averaged from 0-500 km relative to 850 hPa vortex center
G150	0.0205	Temperature perturbation at 150 hPa due to the symmetric vortex calculated from the gradient thermal wind. Averaged from r=200 to 800 km centered on input lat/lon (not always the model/analysis vortex position) (deg C*10)
jd	0.0204	Julian date
VMAX	0.0201	Maximum Surface Wind
IRM1_5	0.0199	Predictors from GOES data (not time dependent) for r=100-300 km but at 1.5 hours before initial time
PW08	0.0191	Time dependent 600-800 km TPW standard deviation (mm * 10)
VMPI	0.0190	Maximum potential intensity from Kerry Emanuel equation (kt)

1635

1636        It is interesting to notice that, IRM1\_5, the standard deviation (STD) of GOES (Knaff  
1637 et al. 2008) BT (brightness temperature) in 100-300 km radius 1.5 hours before the initial  
1638 time, is more important than the average BT value itself (IRM1\_2). The phenomenon  
1639 plausibly says that the non-uniform BT distribution around TC center plays a more  
1640 important role than the uniform BT level for the RI. The same thing takes place with  
1641 PW08, the 600-800 km TPW (Total Precipitable Water) standard deviation from the GFS  
1642 analysis (Berger 2016), which is more important than the corresponding TPW value,  
1643 PW07 represented by the highly correlated RHMD (Table A1). This finding is consistent  
1644 with the relationship between TC intensity and the symmetry of the TC structure. Asif

et al. (2020) used the STD and other statistics of brightness temperature in centric bands to establish a relationship with TC intensity, and those statistics play a similar role of the variance of the deviation angle described by Piñeros et al. (2011) and Ritchie et al. (2012).

## **6.2 LLE-SHIPS model**

### **6.2.1 Hyperparameter tuning for model selection**

#### ***6.2.1.1 Hyperparameters tuning for data filters***

LLE-SHIPS model is trained with both the SHIPS data and the near core ERA gridded data. The SHIPS data will inherit the filtered data for the COR-SHIPS model, and the ERA data will be filtered with LLE. Two new hyperparameters will present with the LLE filter, the number of the nearest neighbors for each observation (no\_neighbors), and the number of dimensions in the reduced space (no\_dimension). As we did in the tuning process for the COR-SHIPS model, BO with 40 iterations is used to tune the two new hyperparameters with no clustering and the preset 0.5 classification decision threshold. In addition, a 10-fold cross-validation kappa score is recorded with a corresponding set of hyperparameter values in each iteration, and the 5 hyperparameter sets with the best 10-fold cross-validation kappa scores are listed in Table 6.8. The search range for no\_neighbors and no\_dimension are pre-defined to 5 to 15 and 10 to 90, respectively.

Based on Table 6.8, we can find that the best kappa score is 0.297 achieved at no\_dimension being 90, and no\_neighbors equaling 15. Therefore, 90 and 15 are selected as final parameter values for no\_dimension and no\_neighbors. With the above filter

1667 setting, the original  $4 \times 37 \times 14 = 2072$  features are filtered into 90 new variables, named as  
 1668 lle1 to lle90. Based on the property of LLE, lle1 to lle90 are independent of each other,  
 1669 i.e., the correlation between any of them is 0. In addition, after calculation, the absolute  
 1670 correlation between new lle variables and the SHIPS variables are less than 0.8; hence no  
 1671 additional variables will be removed in this phase.

1672

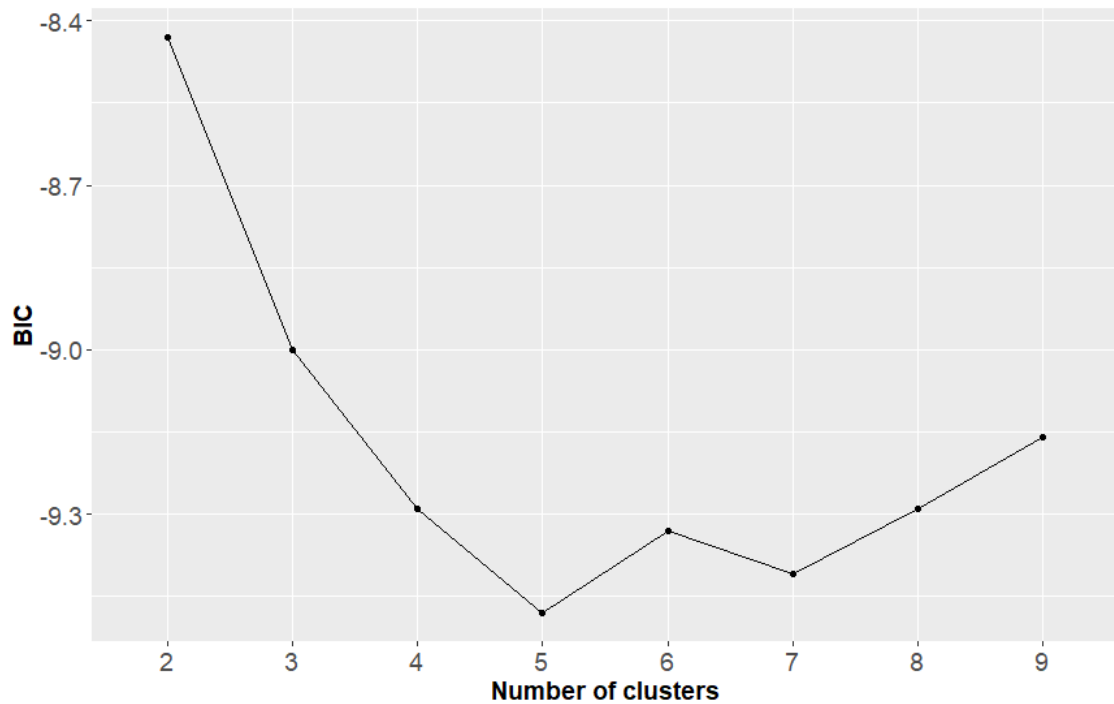
1673 **Table 6.8: The performance for models with different sets of values of the**  
 1674 **hyperparameters, no\_dimension and no\_neighbors.**

Kappa Score	Rank	no_dimension	no_neighbors
0.297	1	90	15
0.295	2	89	15
0.290	3	87	13
0.287	4	90	14
0.282	5	80	15

1675

#### 1676 **6.2.1.2 The number of clusters selected in GMM-SMOTE**

1677 After the hyperparameters in data filters are tuned, the hyperparameters for  
 1678 GMM-SMOTE and XGBoost still need to be tuned for the best results. Similar to the  
 1679 COR-SHIPS model, the BIC values with the different number of clusters are displayed in  
 1680 Figure 6.4, and the BIC values decrease with the cluster number first and then increase.  
 1681 BIC stops decreasing for the next two iterations first at cluster 7, and the minimum BIC  
 1682 is reached at 5 before reaching at 7. Therefore, n\_cluster is selected as 5.



**Figure 6.4: BIC ( $10^6$ ) for GMM with a different number of clusters in LLE-SHIPS model.**

The clustering result is displayed in Table 6.9 with the numbers of minority (RI) and total instances, and the IIR in each cluster. As we defined danger clusters with a 0.2-5 IIR range, all the clusters are included in the following augmentation. Although there is no instance removed, the synthetic instances are created only using instances in the same cluster, which decreases the possibility of outlier creations.

1695 **Table 6.9: The number of minority and total cases, and the Imbalance Rate (with**  
1696 **population RI ratio at 5.1%) for the 5 clusters generated by GMM.**

Cluster	1	2	3	4	5	Total
Number of the minority instance	196	26	112	39	150	523
Number of the total instance	3048	1858	1078	1608	2593	10185
Imbalance Rate	1.286	0.280	2.078	0.485	1.157	1

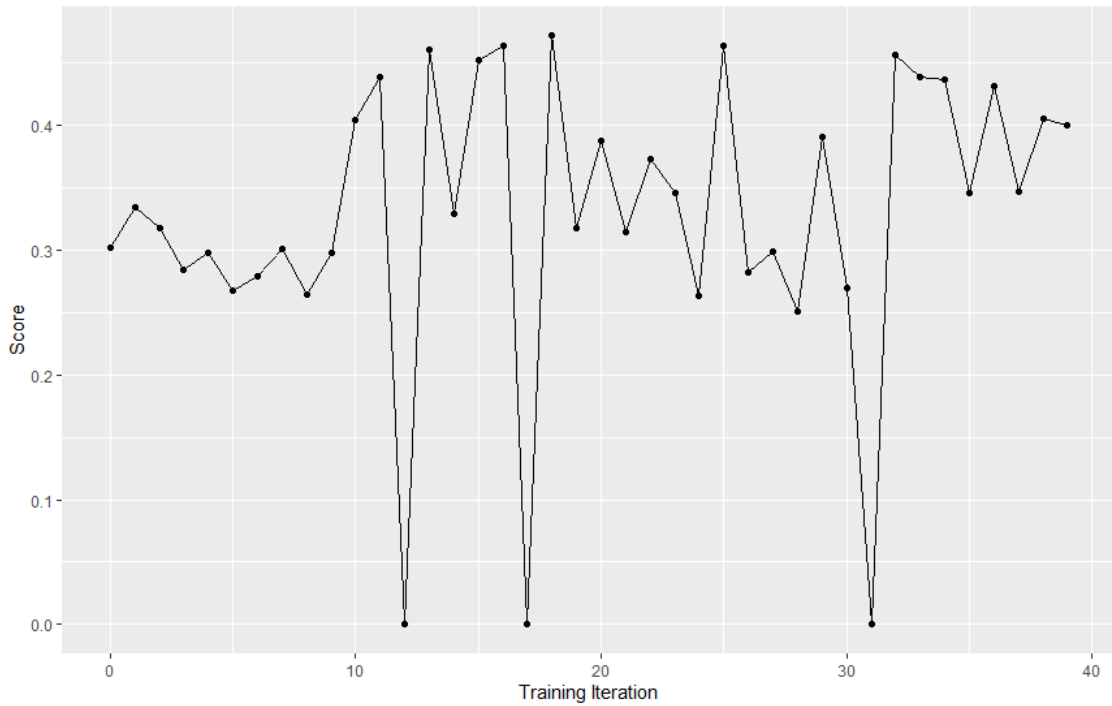
1697

### 1698 *6.2.1.3 Hyperparameters tuning for GMM-SMOTE and XGBoost*

1699 While the hyperparameter tuning is conducted separately for the two data filters,  
1700 the hyperparameters in the GMM-SMOTE and XGBoost in Table 4.1 and Table 5.1 with  
1701 their pre-defined searching space are tuned together by BO as for the COR-SHIPS model.  
1702 Figure 6.5 shows the change of the 10-fold cross-validation kappa scores, the tuning  
1703 criterion, over a total 40 BO iterations.

1704





**Figure 6.5: Variation of Cross-validation kappa scores over Bayesian Optimization iteration numbers for LLE-SHIPS model.**

As shown in Table 6.10, the top 5 performed hyperparameter sets are M18 (kappa=0.472), M25 (0.464), M16 (0.464), M13 (0.461), and M32 (0.456). To find a balanced model, the same score system described for the COR-SHIPS model is used, and the value ranks and their corresponding conservativeness ranking scores are listed in Table 6.11. The total conservativeness scores are also calculated, which are 25, 20, 30, 31, and 23, respectively for M18, M25, M16, M13, and M32, and the parameter set M18 with the middle conservativeness ranking score is chosen for the following implementation and discussion. Coincidentally, M18 is also the set associated with the highest kappa score.

1719 **Table 6.10: Top 5 performed hyperparameter sets, the corresponding cross-**  
1720 **validation kappa scores, and specific values of the tuned hyperparameters. The**  
1721 **numbers after “M” denoting the iteration numbers.**

Name	M18	M25	M16	M13	M32
Kappa score	0.472	0.464	0.464	0.461	0.456
n_cluster	5	5	5	5	5
m_neighbors	8	8	3	7	7
k_neighbors	12	14	3	14	14
shrinkage	0.21	0.16	0.23	0.16	0.16
n_estimators	731	2000	2000	1286	1819
subsample	0.85	0.9	0.64	0.77	0.72
colsample_bytree	0.8	0.99	0.4	0.4	0.99
reg_alpha	0.1	0.1	0.1	0.1	0.1
reg_lambda	0.5	0.5	20	17.93	12.2
gamma	0	0	0	0	0
min_child_weight	1.4	0.5	0.5	3.48	0.5
max_depth	10	10	3	7	8

1722

1723 **Table 6.11: The descending value ranking of individual hyperparameter among the**  
1724 **top 5 performed cases, and the corresponding conservativeness ranking scores in**  
1725 **parentheses. The variables with normal font are those favoring smaller values for**  
1726 **conservativeness, and those *italicized* favoring larger values.**

Name	M18	M25	M16	M13	M32
<i>m_neighbors</i>	<i>1 (4)</i>	<i>1 (4)</i>	<i>5 (1)</i>	<i>3 (2)</i>	<i>3 (2)</i>
k_neighbors	4 (4)	1 (1)	5 (5)	1 (1)	1 (1)
shrinkage	2 (2)	3 (3)	1 (1)	3 (3)	3 (3)
n_estimators	5 (5)	1 (1)	1 (1)	4 (4)	3 (3)
subsample	2 (2)	1 (1)	5 (5)	3 (3)	4 (4)
colsample_bytree	3 (3)	1 (1)	4 (4)	4 (4)	1 (1)
<i>reg_alpha</i>	<i>1 (1)</i>	<i>1 (1)</i>	<i>1 (1)</i>	<i>1 (1)</i>	<i>1 (1)</i>
<i>reg_lambda</i>	<i>4 (1)</i>	<i>4 (1)</i>	<i>1 (5)</i>	<i>2 (4)</i>	<i>3 (3)</i>
<i>gamma</i>	<i>1 (1)</i>	<i>1 (1)</i>	<i>1 (1)</i>	<i>1 (1)</i>	<i>1 (1)</i>
<i>min_child_weight</i>	<i>2 (4)</i>	<i>3 (1)</i>	<i>3 (1)</i>	<i>1 (5)</i>	<i>3 (1)</i>
max_depth	1 (1)	1 (1)	5 (5)	4 (4)	3 (3)
Total score	28	16	30	32	23

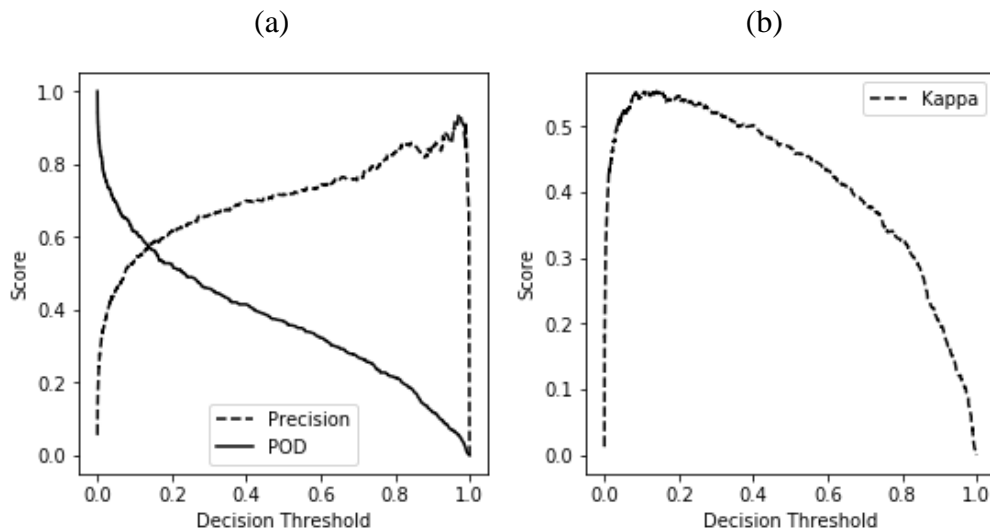
1727

#### 1728 6.2.1.4 Hyperparameter tuning for the decision threshold

1729 The last tuned hyperparameter is the decision threshold on the XGBoost classifier  
1730 output. Initially, the decision threshold is set as 0.5. Figure 6.6 (b) displays that the kappa  
1731 score approximates the highest value, 0.52 when the decision threshold reaches 0.15.  
1732 Figure 6.6 (a) displays variations of precision (1-FAR) and POD variations as functions  
1733 of the decision threshold from 10-fold cross-validation in the training/validation data.  
1734 Decision threshold of 0.15 also close to the intersection of the POD and precision score, a  
1735 relatively balanced point for POD and FAR. As a result, 0.15 is selected as the decision  
1736 threshold.

1737

1738



1739

1740 **Figure 6.6: (a) Precision and POD score vs. decision threshold. (b) Kappa score vs.**  
1741 **decision threshold in LLE-SHIPS model.**

1742

## 1743 6.2.2 LLE-SHIPS result on test data

1744 As discussed in the COR-SHIPS model, the evaluation of the prediction  
 1745 performance is on test data only. The test confusion matrix for the model, before  
 1746 hyperparameter tuning (MB), and after hyperparameter tuning (MA) is displayed in Table  
 1747 6.12. Learning from the Table, MA's TN (1,469) and FN (55) is slightly smaller than  
 1748 1,491 and 75 of MB, while MA's TP (40) and FP (33) is significant larger than 20 and 11  
 1749 of MB.

1750

1751 **Table 6.12: Confusion matrix values after (before) hyperparameter tuning with the**  
 1752 **test data.**

	Predicted RI	Predicted non-RI	Actual
Actual RI	40 (20)	55 (75)	95
Actual non-RI	33 (11)	1469 (1491)	1502
Predicted	73 (88)	1524 (1509)	

1753

1754 Kappa, PSS, POD, and FAR are used for the model evaluation, and their values  
 1755 for MB and MA are elaborated in Table 6.13. The POD and FAR values for MB and MA  
 1756 cases demonstrated the importance of hyperparameter tuning. After tuning, POD  
 1757 increases 99.5% from 0.211 to 0.421, while FAR decreases from 0.645 to 0.563, 12.7%.  
 1758 The overall statistics PSS and kappa score also increased from 0.203 to 0.399 (96.6%)  
 1759 and from 0.297 to 0.454 (52.9%), respectively, confirming the significant improvement  
 1760 on RI prediction with the hyperparameter tuning procedure, and apparently, the model

1761 was overfitted before tuning process with so many features. Furthermore, the  
 1762 improvement is almost 4 times than that of the COR-SHIPS model (25.6% and 12.4% in  
 1763 PSS and kappa improvement), indicating the hyperparameter tuning is more efficient in  
 1764 the more complicated LLE-SHIPS model.  
 1765

1766 **Table 6.13: Performance comparisons. MB and MA denote the models before and**  
 1767 **after the hyperparameters in GMM-SMOTE and XGBoost are tuned.**

Model	Kappa	PSS	POD	FAR
MB	0.297	0.203	0.211	0.645
MA	0.454	0.399	0.421	0.563
Improvement MB	52.9%	96.6%	99.5%	-12.7%

1768

### 1769 6.2.3 Feature importance

1770 We learned from section 6.1.3 that how XGBoost evaluates the variable (feature)  
 1771 importance score for COR-SHIPS model. In LLE-SHIPS model, the same approach can  
 1772 be used to calculate the importance score for the SHIPS variables and lle1, ..., lle90, but  
 1773 we cannot directly calculate the importance score for the original ERA-Interim variables  
 1774 in their feature space; hence they cannot be linked to the original ERA variables. So  
 1775 instead, we try to relate the IS of all the lle1, ..., lle90 to individual ERA parameter  
 1776 groups (based on correlation). Therefore, the feature importance evaluation in the LLE-  
 1777 SHIPS model is divided into 2-steps. First, the XGBoost is used to evaluate the  
 1778 importance score for SHIPS variables and lle1 to lle90 in the same way as in COR-  
 1779 SHIPS model, and then a feature permutation approach is used to evaluate the importance

score for the original ERA-Interim feature space separately based on the importance score generated from the first step.

### 6.2.3.1 Variable importance in XGBoost

Table 6.14 displays the 10 most important variables among the 162 selected not highly correlated variables (90 LLE variables and 72 SHIPS variables) and their definition. We need to notice that none of the LLE variables are among the top 10 variables. The reason might be there are too many variables generated from the LLE, and having so many variables reduces the importance score for each of them. The assumption is confirmed by summing the importance score for lle1 to lle90, which is 0.4288, only a bit smaller than that of the SHIPS variables (0.5712).

**Table 6.14: Variables of top ten importance, their importance scores, and feature description from SHIPS (2018c) in LLE-SHIPS model.**

Variable	Importance	Description
BD12	0.018769	The past 12 hour intensity change
VMAX	0.016706	Maximum Surface Wind
DTL	0.013821	The distance to nearest major land
SHRD	0.012952	850-200 hPa shear magnitude
TWXC	0.01151	Maximum 850 hPa symmetric tangential wind at 850 hPa from NCEP analysis
G150	0.011339	Temperature perturbation at 150 hPa due to the symmetric vortex calculated from the gradient thermal wind. Averaged from r=200 to 800 km centered on input lat/lon (not always the model/analysis vortex position) (deg C*10)
VMPI	0.011262	Maximum potential intensity from Kerry Emanuel equation (kt)
REFC	0.011262	Relative eddy momentum flux convergence

TGRD	0.01106	The magnitude of the temperature gradient between 850 and 700 hPa averaged from 0 to 500 km estimated from the geostrophic thermal wind
IRM1_5	0.010672	Predictors from GOES data (not time dependent) for r=100-300 km but at 1.5 hours before initial time

1793

1794           The past 12 hour intensity change, BD12, has the largest importance score,  
1795 0.018769, which slightly better than the importance score of the second important  
1796 variable, VMAX, the Maximum Surface Wind, and its highly correlated variable, the  
1797 Minimum Sea Level Pressure. The importance score between BD12 and VMAX is very  
1798 similar. The rest top 10 variables are DTL, SHRD, TWXC, G150, VMPI, REFC, TGRD,  
1799 and IRM1\_5. The highly correlated variable groups with the important variables and the  
1800 importance scores of all of the 162 variables can also be found in Table A2.

#### 1801 **6.2.3.1 Group importance in LLE**

1802           Molnar (2019) described a feature permutation approach to evaluate the  
1803 importance of features on training dataset for nonlinear models where the importance  
1804 score cannot be derived easily. We assume that for the feature space in any given dataset  
1805 is  $X$ ,  $f(X)$  is the predicted value by the classifier  $f$ , and  $y$  is the ground truth. We further  
1806 assume the loss of the classifier is  $L(y, f(X))$ . Then for each feature in the feature space  $X$ ,  
1807 permute its value to 0 for all the observations while keeping other features unchanged  
1808 (represented as  $X^{permute}$ ). Finally, the difference between the loss of the permuted  
1809 feature space ( $X^{permute}$ ), and the original loss is calculated for each feature, and the  
1810 difference is used as its importance.

1811 Although feature permutation is an efficient approach to evaluate the feature  
1812 importance for different models, especially for a black-box model such as LLE, Molnar  
1813 (2019) also indicates that the permuted feature importance could be biased by the  
1814 highly correlated features. For example, if we evaluate the importance score for each of  
1815 the 2,072 variables, the result, i.e., the importance score is not accurate due to the  
1816 existence of the highly correlated variables because they could influence each other.  
1817 Similar to the removal of highly correlated variables in the SHIPS data filter, pairwise  
1818 correlations of all the features are calculated and compared with the correlation threshold  
1819 0.8. This process results in 135 groups, and the details of the Group are elaborated in  
1820 Table A2. Then an importance score is calculated for each Group, and details will be  
1821 elaborated later in the section.

1822 The group-level importance score is calculated specifically as:

1823  $f$ : trained model;  $X$ : original feature space;  $y$ : ground truth;  $L(y, f(X))$ : loss between the  
1824 ground truth and the predicted value by the classifier.

- 1825 1. Calculate  $Importance_{LLE\_Total}$  as the sum of the importance score of lle1 to lle90  
1826 derived from XGBoost, here is 0.4288.
- 1827 2. Calculate the original model error  $L(y, f)$ .
- 1828 3. for each group:
  - 1829 a) Generate feature matrix  $X^{permute}$  by setting features in that group to 0,  
1830 which breaks the corresponding correlation between all the features
  - 1831 b) Calculated error  $L(y, f(X^{permute}))$
  - 1832 c) Estimate the importance for the group  $imp = L(y, f(X^{permute})) - L(y, f)$



1833 d) Associate the score to the group

1834 e) Negative importance is set to 0

1835 4. Group importance score is rescaled as attributing the total important scores by

1836 LLE variables based on the ratio of loss of a particular group to the total loss (sum

1837 of a group losses), the specific calculation is:

$$1838 \quad Imp_{group} = Importance_{LLE\_Total} * imp_{group} / \sum_i \text{for all the groups } imp_i \quad (6.1)$$

1839 5. Sort groups by their number of features.

1840 Based on the Algorithm, the importance score for each Group is calculated, and

1841 groups with the top five important scores are list in Table 6.15. Intuitively, turning much

1842 more variables to 0 could reduce the model's performance more than turning much fewer

1843 variables to 0 because changing more variables are likely to alternate the model's

1844 performance more. However, based on Table 6.15, we find that the top 5 does not contain

1845 too many groups with a large number of variables, which indicates that these groups with

1846 few variables play a more important role in RI prediction than other groups, especially

1847 the groups with significant more variables.

1848

1849

1850

1851

1852 **Table 6.15: ERA-Interim variable group with top 5 importance scores, calculated**  
1853 **from the second step. The group number (Group), the number of variables in the**  
1854 **group (Group size), and the importance score (Importance).**

Group	Group size	Importance
49	5	0.023614
88	1	0.021988
1	309	0.019687
29	11	0.019662
3	148	0.017280

1855

1856           Group 49 (G49) has the highest importance score (IS), 0.024, and it has 5  
1857 variables, NT12\_v\_118, NT06\_v\_118, NT00\_v\_117, NT00\_v\_118, NT06\_v\_117, which  
1858 indicates that the northward wind speed on level 17 (450 hPa) at 6 hours before, and at  
1859 present, together with level 18 (400 hPa) at 12 hours before, 6 hours before, and at  
1860 present are important in RI prediction. We can find that the middle level's (400 hPa and  
1861 450 hPa) northward wind plays a significant role than that in the lower level, i.e., 1000  
1862 hPa, and higher level (1 hPa). The reason could be when the RI starts to occur, the  
1863 northward wind speed in 400 and 450 hPa change faster than that of other levels. We can  
1864 also find that both 6 hours before and the present northward wind speed are important at  
1865 400 and 450 hPa, which indicates that the northward wind speed in 400 and 450 hPa start  
1866 to change immediately before the occurrence of RI, and 18 hours before are too long to  
1867 influence the occurrence of RI.

1868           Wang et al. (2015) found that “In the active (inactive) season, the low-level (deep  
1869 layer) shear is more negatively correlated with the TC intensity change than the deep-  
1870 layer (low level) shear.” Our study identifies that importance of the northward wind

1871 speed in the 400 and 450 hPa for the RI prediction, that could contribute to the 400 and  
1872 450 hPa VWS (vertical wind shear), which recognizes the importance of the mid-layer  
1873 shear with regard to the intensity change in addition to the finding of Wang et al. (2015).

1874         The second most important group, the G88 with an IS of 0.0220, only has one  
1875 variable, NT18\_pv\_11, the potential vorticity at 18 hours before on the first level (1000  
1876 hPa). The importance score for NT18\_pv\_11 is even 17% higher than that for the most  
1877 important variable in Table 6.13, BD12 with a 0.0188 score, which is also the highest  
1878 importance for a single variable. This result demonstrated that the machine learning  
1879 method could identify important features, which may not be in the commonly used data  
1880 set, such as the SHIPS database. However, the role of pv in RI was identified by others  
1881 already (e.g., Martinez et al. 2019; Tsujino and Kuo 2020). Tsujino and Kuo (2020)  
1882 detailed the changed of pv during the RI of Supertyphoon Haiyan (2013) with numerical  
1883 simulation. They emphasized the pv increasing around 3-5km height at the beginning  
1884 stage of the RI. Carefully checking their results (Fig. 2b&c), one can find the pv actually  
1885 increases simultaneously around the sea level in 20-40 km range from the center, which  
1886 is the same as what we identified here by the NT18\_pv\_11.

1887         All other level 1 pv (3 of them) are grouped in G63 with importance scores (IS)  
1888 (0.010746). All level 2 in G55 with 4 members and IS 0.004744. All other pv are in G4  
1889 with 140 members but IS being only 0.006264. Those numbers demonstrated that only  
1890 lower layer pv affects the RI process.

1891           The third most important Group is G1, which has 309 features in the Group, and  
1892   with IS 0.0197. Since all types of ERA-Interim variables are included in the Group, it is  
1893   difficult to trace back which variable is more important. This implies the SHIPS dataset is  
1894   very useful because it removes a lot of highly correlated variables and only extract  
1895   important information from these variables.

1896           The fourth most important group is G29, which has the IS, 0.020, and consists of  
1897   11 variables, i.e., NT18\_u\_11, NT12\_u\_11, NT06\_u\_11, NT00\_t\_110, NT00\_u\_116,  
1898   NT00\_u\_117, NT06\_u\_117, NT12\_u\_117, NT18\_u\_118, NT12\_u\_118. We can find that  
1899   most of the variables in the group is u, the eastward wind speed. Similar with G49 but a  
1900   slightly different, the eastward wind speed at level 17 (450 hPa), and 18 (400 hPa) play  
1901   an important role in RI prediction. Other than 400 and 450 hPa eastward wind speed, the  
1902   eastward wind at 1000 hPa at 6, 12, and 18 hours before RI also plays an important role.  
1903   As discussed above, Wang et al. (2015) found that “low-level shear between 850 (or 700)  
1904   and 1000 hPa is more negatively correlated with TC intensity change than any deep-layer  
1905   shear during the active typhoon season,” which matches our findings that eastward wind  
1906   speed, related to the VWS, at 1000 hPa are significant in RI prediction. Additionally, we  
1907   also recognize that the mid-level (400 and 450 hPa) eastward wind speed (VWS) are  
1908   important to TC intensity change. One exception variable in the Group is the temperature  
1909   (t) at 775 hPa, NT00\_t\_110, although highly correlated with u in terms of value, possibly  
1910   misplaced in the group because there’s only one t variable in the Group.

1911           The fifth most important Group, G3, has IS 0.017 and 148 variables, consists of w  
1912   (the pressure vertical velocity) at all levels over 18 hours before, 12 hours before, 6 hours

1913 before, and at present. This indicates the vertical pressure speed plays an important role  
1914 in RI prediction, which includes SHIPS variable O500 (highly correlated with O700 in  
1915 Table A1), ranked 72 in Table A5. This indicates that other than 500 and 700 hPa, the  
1916 pressure vertical velocity ( $w$ ) at other pressure levels is also important, and more  
1917 researches should be done to figure out more details.

1918 In sum, two out of the top five important groups, G45, and G29, contain eastward  
1919 and northward wind speed variables, especially at 400, 450, and 1000 hPa, which  
1920 indicates that wind velocity, hence the VWS at 400, 450, and 1000 hPa pressure level  
1921 plays a significant role in RI prediction, not only matches what have been found in Wang  
1922 et al. (2015) but also identifies the importance of the mid-level vertical shear to intensity  
1923 change. Another group, G3, only contains the pressure vertical velocity, indicates that  
1924 vertical pressure speed is critical in RI prediction. Other than the O500, and O700  
1925 included in the SHIPS database, it is necessary to dig out other significant pressure levels  
1926 for the pressure vertical velocity. One variable, 18 hour before potential vorticity at 1000  
1927 hPa, is more significant than BD12, and more researches need to be done for this specific  
1928 variable.

1929 Here we derive the group level importance score for ERA-Interim variables.  
1930 Although because the AI system is consisted of too many components that the score is  
1931 not 100% accurate, the system is still able to identify useful features in addition to SHIPS  
1932 database.

1933 More details of the XGBoost scores can be found in Table A5, group importance  
1934 scores can be found in Table A3.

## 1935 **6.3 DL-SHIPS model**

### 1936 **6.3.1 Hyperparameters tuning and result**

#### 1937 ***6.3.1.1 Hyperparameters tuning for data filter***

1938           The structure of the DL-SHIPS model is almost the same as that of the LLE-  
1939 SHIPS model except that the ERA-Interim data is filtered with a CNN-based autoencoder  
1940 network, as introduced in section 3.2.3. Similar to the LLE-SHIPS model, in the DL-  
1941 SHIPS model, the correlation threshold (0.8) in the SHIPS data filter is inherited directly  
1942 from that of the COR-SHIPS model. The difference is that in the DL-SHIPS model,  
1943 unlike in the LLE-SHIPS model where all variables in ERA-Interim reanalysis data are  
1944 treated together with only one dimension reducing model for the feature extraction, one  
1945 autoencoder network is trained to extract information from each individual of 14 ERA-  
1946 Interim variables. In other words, there are 14 different autoencoder networks in total for  
1947 the 14 variables. Other components in the LLE-SHIPS model and the DL-SHIPS model  
1948 are the same.

1949           In the tuning process of the ERA-Interim data filter in the DL-SHIPS model, all  
1950 14 ERA-Interim original variables are initialed with the same 3D CNN auto-encoder  
1951 structure as described in Figure 3.6. The to-be-determined hyperparameter, the dimension  
1952 of the compressed feature (num), is pre-determined as 8, therefore, 8 new variables are  
1953 generated from each network, labeled as variable+order in the compressed feature layer  
1954 ('1' to '8'). For example, v1 to v8 are new variables derived from the trained 3D CNN  
1955 auto-encoder for variable v. A three-step tuning process for “num,” which is similar with  
1956 the tuning process of SHIPS data filter. In the first step, training a separated 3D CNN

1957 auto-encoder for each of the 14 variables for 200 epochs. In the second step, the output  
 1958 from SHIPS data filter (72 variables) and the output of the ERA data filter (the trained  
 1959 DL models), which is all the 14 variables from variable1 to variable8 ( $14 \times 8 = 112$ ) are  
 1960 concatenated to form the input to the sampler, where zero (zero values for all instances)  
 1961 and highly correlated variables are removed. Then the BO with 40 iterations is used to  
 1962 tune hyperparameters in for GMM-SMOTE in Table 4.1 and XGBoost related  
 1963 hyperparameters listed in Table 5.1 with no clustering, and the preset 0.5 classification  
 1964 decision threshold. Finally, the hyperparameter set with the highest 10-fold cross-  
 1965 validation is selected; instead of kappa score, here, we use the importance score for  
 1966 hyperparameter tuning instead of the kappa score used in the SHIPS data filter tuning.  
 1967 The importance score for each input variable is derived. In the third step, the importance  
 1968 for each of the 14 variables is calculated as the sum of the variable1 to variable8, and the  
 1969 num is determined by the summed importance score that will be discussed later.

1970 Figure 6.7 shows the training losses (mean square error) of 14 auto-encoder  
 1971 networks change over iterations that are trained for 200 epochs, respectively, which  
 1972 indicates that all the networks are converged after 100 epochs. 8 new variables are  
 1973 engineered from each ERA-interim variable first, because of the characteristic of the  
 1974 auto-encoder, features with no information, i.e., zero feature, could be created. The zero  
 1975 features are listed in the last column of Table 6.16, and they are removed. A correlation  
 1976 check is conducted among the 72 SHIPS variables and the newly derived non-zero  
 1977 variables, and highly correlated ( $>0.8$ ) variables should be removed. There are two  
 1978 correlations, i.e., the correlation between ERA variables vs. SHIPS variables, and the

correlation between only ERA variables, as highly correlated SHIPS variables are already removed in section 3.1.3. Correlation between ERA variables vs. SHIPS variables is all less than 0.8 so no variable will be removed in this phase. However, when we check the correlation between ERA variables, we find highly correlated ( $>0.8$ ) variables are existed, and should be removed. The same procedure in section 3.1.3 in SHIPS data filter is used again, and highly correlated variables are sorted in Table 6.17 (only variables with highly correlated ones are displayed). Therefore, d5, r4, clwc4, w6 are kept as well as those not highly correlated features, while d2, d3, cc8, q7, cc3, and w1 are removed. Then all remained features are concatenated together and fit into the GMM-SMOTE sampler and XGBoost classier. The hyperparameters of the model is tuned by BO based on no clustering, decision threshold set at 0.5, and other hyperparameters as listed in the final column of Table 4.1 and Table 5.1. The summed importance score for all 14 variables are displayed in Table 6.16, where variables are sorted based on their summed importance scores. Based on Table 6.16, we can find that with the decreasing of the summed importance score, the number of kept (removed) features are roughly decreasing (increasing). Therefore, we are tuning the hyperparameter, dimension of the compressed feature (num), based on the summed importance score for each variable, and the details are displayed in Table 6.16.



1999 **Table 6.16: 14 variables, their summed importance score, non-zero features**  
2000 **extracted from each variable network, and the corresponding missing variables due**  
2001 **to all zeros.**

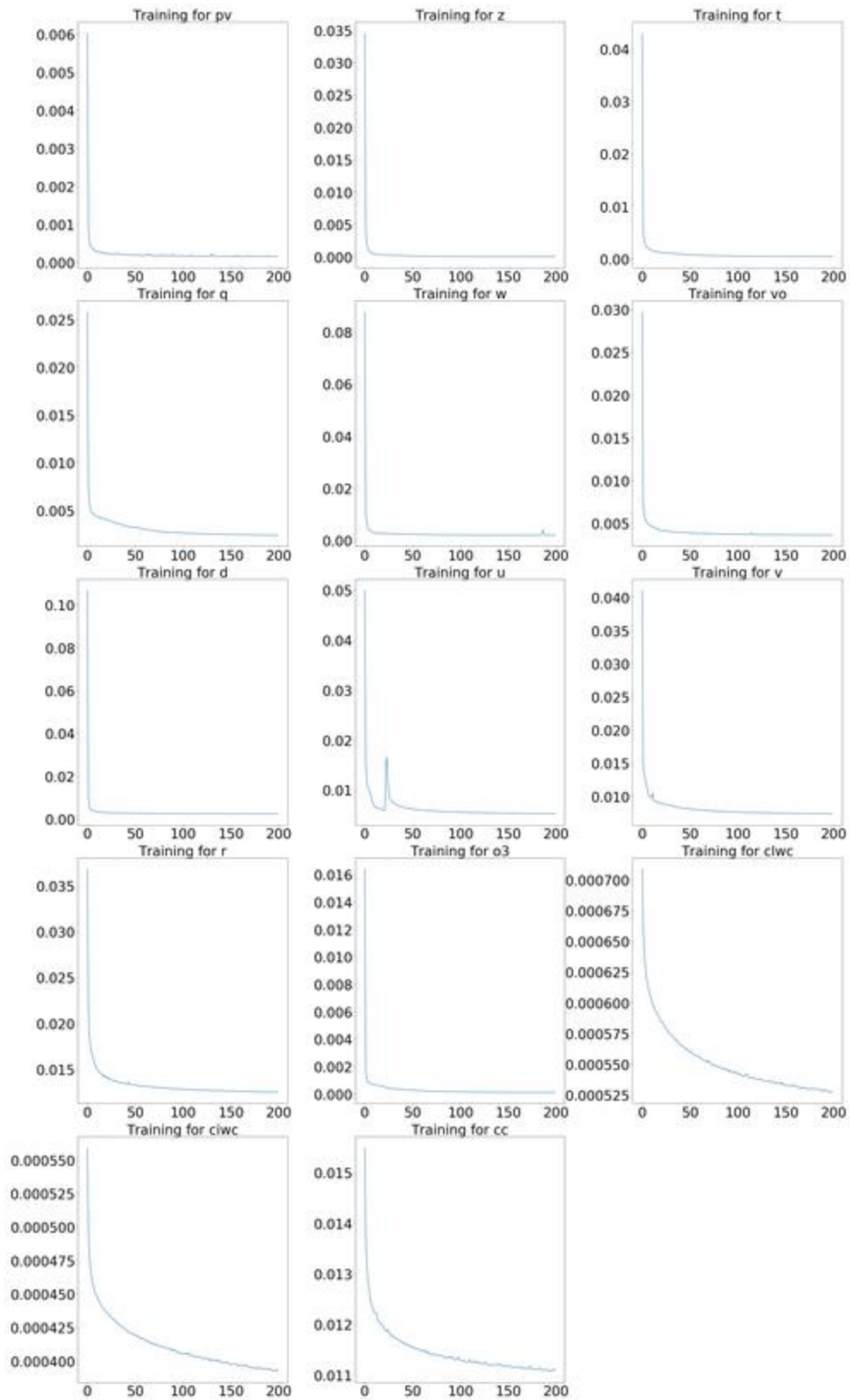
Variable	Summed Importance Score	The number of kept features	Corresponding missing features
q	0.062	6	q3, q8
r	0.055	7	r4
u	0.055	8	
v	0.056	7	v7
pv	0.050	6	pv7, pv8
vo	0.049	5	vo1, vo6, vo7
w	0.042	6	w3, w4
d	0.039	5	d1, d4, d6
t	0.020	4	t1, t2, t5, t6
z	0.017	3	z2, z3, z5, z7, z8
o3	0.019	3	o31, o32, o33, o34, o36
clwc	0.014	2	clwc2, clwc4, clwc5, clwc6, clwc7, clwc8
cc	0.011	2	cc1, cc2, cc4, cc5, cc6, cc7
ciwc	0.008	1	ciwc2, ciwc3, ciwc4, ciwc5, ciwc6, ciwc7, ciwc8

2002

2003 **Table 6.17: Highly correlated variable groups. Only groups with more than one**  
2004 **variable is displayed. “Variable” column indicates the selected variable, and its**  
2005 **highly correlated (>0.8) variables are displayed in “Highly correlated variables.”**

Variables	Highly correlated variables
d5	d2, d3
r4	cc8, q7
clwc4	cc3
w6	w1

2006



2007

**Figure 6.7: Network training loss over iterations for pv, z, t, q, w, vo, d, u, v, r, o3, clwc, ciwc, cc from left to right and from top to bottom. The 14 graphs represent the training loss change for each variable respectively. The y-axis represents the training loss, and decreasing from top (the maximum loss value of the variable) to bottom (0). The x axis represents the iterations, and increasing from left (0) to right (200).**

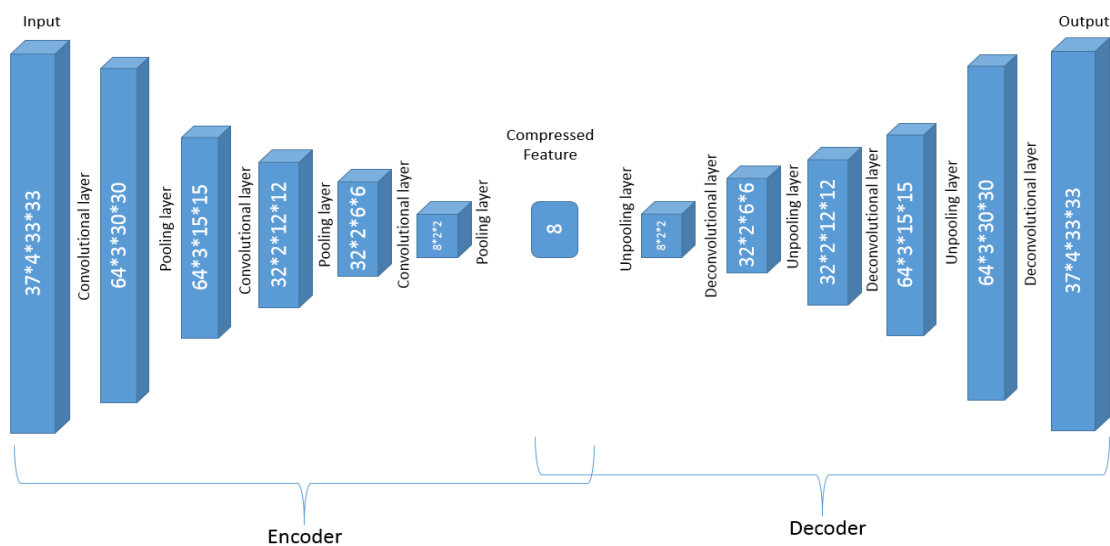
These variables are categorized into 3 classes in Table 6.18 based on their summed importance score described in Table 6.16, and their new structures are displayed in Figure 6.8. The dimension of the compressed feature (num) equal to 8, 4, and 2 in Table 6.18 corresponding to (a), (b), and (c), respectively in Figure 6.8. Therefore, the auto-encoder network for each variable is retrained with the new structure, i.e., Figure 6.8 (a) for variables pv, q, r, u, v, and vo, Figure 6.8 (b) for variables, w, d, t, Figure 6.8 (c) for z, o3, cc, ciwc, and clwc.

**Table 6.18: Dimensions of the compressed features of auto-encoder after tuning based on the summed importance score described in Table 8 for each of the 14 variables.**

Importance sum	Dimension of the compressed feature (num)
Less than 0.02	2
0.02 to 0.045	4
Above 0.045	8

2028

(a)



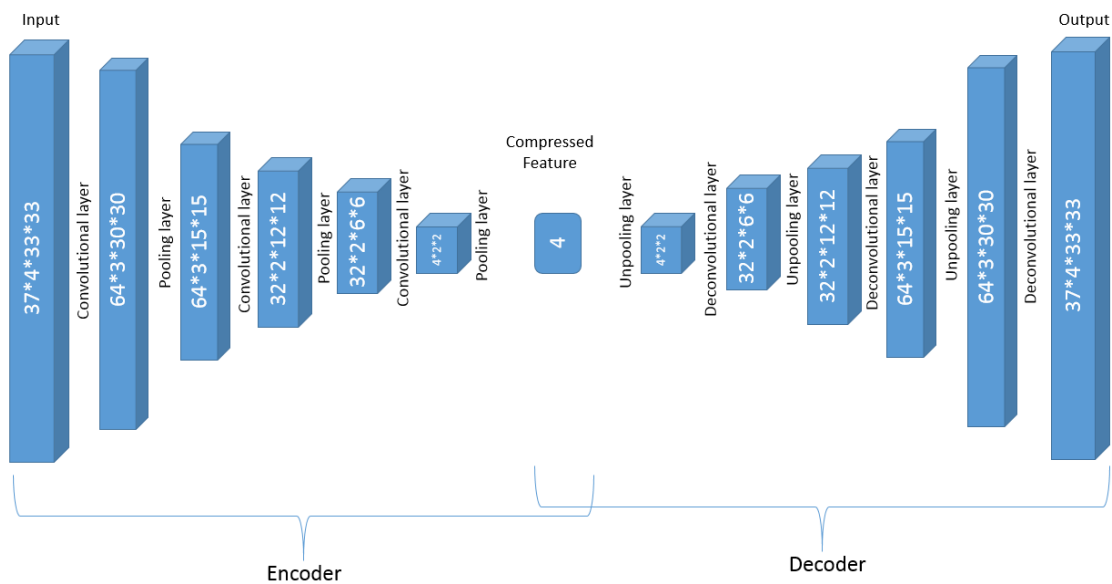
2029

2030

2031

2032

(b)

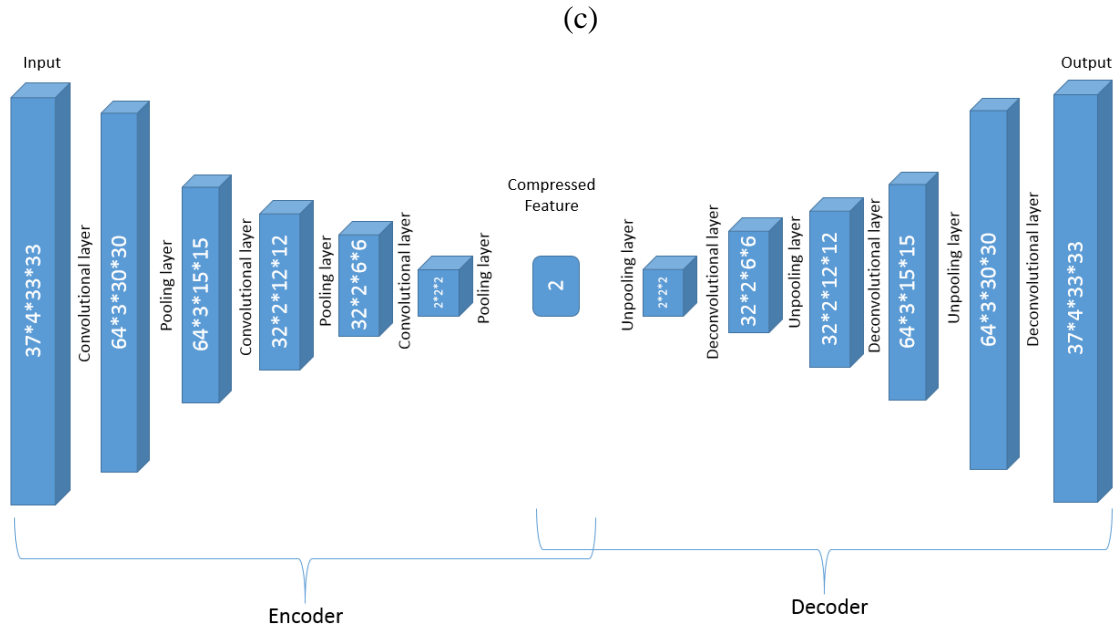


2033

2034

2035

2036  
2037  
2038



2039  
2040  
2041

**Figure 6.8: Structure for adjusted auto-encoder network.**

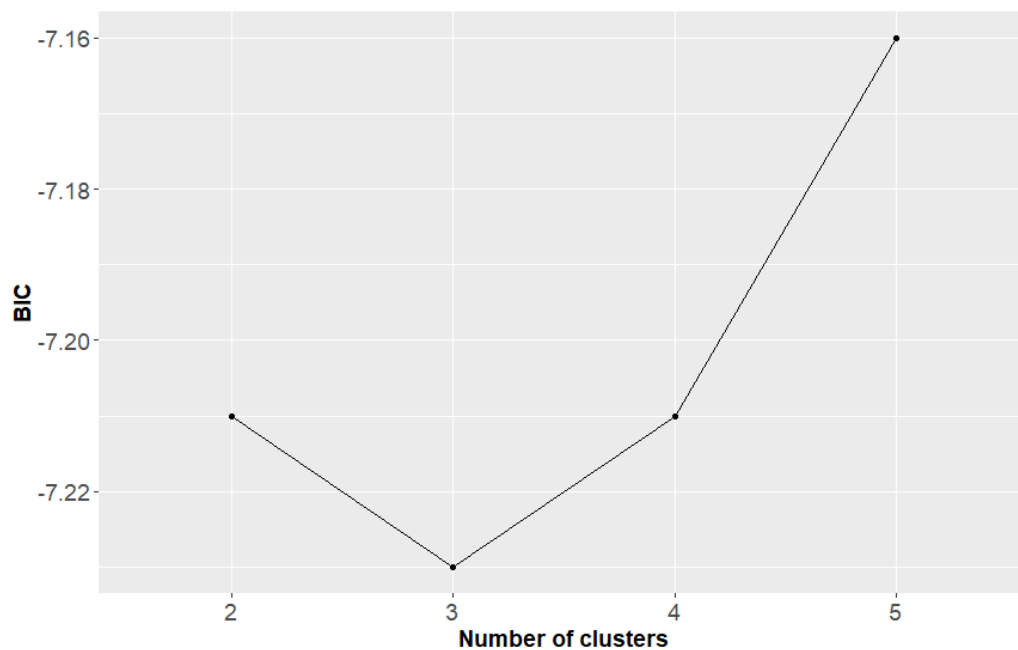
2042           After the new structure for DL Interim filter is derived, we retrain each network  
2043 for 200 times, and the training loss is very similar with those shown in Figure 6.7 so their  
2044 training loss graph is not displayed. We also find that all the networks are converged in  
2045 100 iterations. As same as being processed above, pv2, pv4, pv5, q1, u4, u7, v1, v2, vo2,  
2046 vo7, w4, d1, d2, d4, t2, t3, t4, z1, and o32 are zero features, i.e., contain all zeros, are  
2047 removed. cc2 is highly correlated ( $>0.8$ ) with cc1, and w2 and w3 are highly correlated  
2048 with w1. Hence clwc1, w2, and w3 are removed. The remained 48 ( $6*8+3*4+5*2-22$ )  
2049 features are concatenated with the filtered SHIPS variables, and are used as the input to  
2050 the GMM-SMOTE.

2051 **6.3.1.2 The number of clusters selected in GMM-SMOTE**

2052 After the hyperparameters in data filters are tuned, the hyperparameters for  
2053 GMM-SMOTE and XGBoost still need to be tuned for the best results. Similar to the  
2054 LLE-SHIPS model, the BIC values with the different number of clusters are displayed in  
2055 Figure 6.9, n\_cluster is selected as 3 with the smallest BIC value.

2056 The clustering result is displayed in Table 6.19 with the numbers of minority (RI)  
2057 and total instances, and the IIR in each cluster. As we defined danger clusters with 0.2-5  
2058 IR range, Clusters 1, 2, and 3 are all included in the following augmentation.

2059



2060

2061 **Figure 6.9: BIC (10<sup>6</sup>) for GMM with different number of clusters.**

2062

2063 **Table 6.19: Number of minority, total cases, and the IIR (with population RI ratio**  
2064 **at 5.1%) for the 3 clusters generated by GMM.**

Cluster	1	2	3	Total
Number of the minority instance	209	36	258	523
Number of the total instance	3222	2645	4318	10185
Imbalance Rate	1.297	0.272	1.195	1

2065

2066

### 2067 *6.3.1.3 Hyperparameters tuning for GMM-SMOTE and XGBoost*

2068 Similar to LLE-SHIPS model, Figure 6.10 shows the 10-fold cross-validation

2069 kappa scores on the training-validation dataset change over a total 40 BO iterations.

2070 Similarly, since the trend with the iteration is unpredictable, hyperparameter sets with the

2071 best 5 kappa scores are selected, and their performance and hyperparameter values are

2072 displayed in Table 6.20 with the same MX notation.

2073 The top 5 performed hyperparameter sets are M23 (0.516), M25 (0.506), M27

2074 (0.502), M31 (0.501), and M21 (0.498). Similarly, the ranking for conservativeness

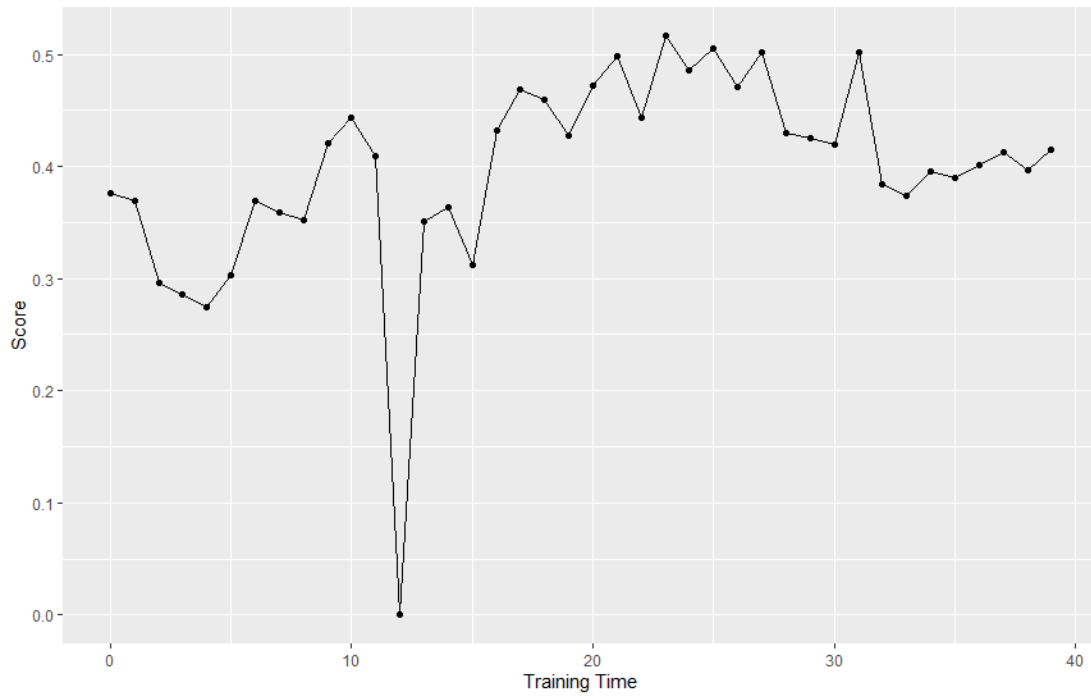
2075 among the five groups for individual hyperparameters are listed in Table 6.21. The total

2076 conservativeness scores are also calculated, which are 17, 25, 23, 32, and 32 respectively

2077 for M21, M31, M27, M25, and M23. Since our goal is to choose a model neither

2078 conservative nor aggressive, the parameter set M31 with the middle conservativeness

2079 ranking score is chosen for following implementation and discussion.



**Figure 6.10: Variation of Cross-validation kappa scores over Bayesian Optimization iteration numbers.**

**Table 6.20: Top performed hyperparameter sets, the corresponding cross-validation kappa scores, and specific values of the tuned hyperparameters. The numbers after “M” denoting the iteration numbers.**

Name	M21	M31	M27	M25	M23
Kappa score	0.498	0.501	0.502	0.506	0.516
m_neighbors	10	10	10	10	9
k_neighbors	14	9	14	7	8
shrinkage	0.23	0.19	0.15	0.12	0.15
n_estimators	2000	2000	2000	1088	1603
subsample	0.50	0.50	0.50	1.00	0.88
colsample_bytree	1.00	1.00	1.00	0.50	0.82
reg_alpha	0.50	0.50	0.50	0.50	0.82
reg_lambda	20.00	20.00	0.50	20.00	18.38
gamma	0	0	0	0	0
min_child_weight	0.50	0.50	0.50	0.50	0.76
max_depth	6	3	3	5	3



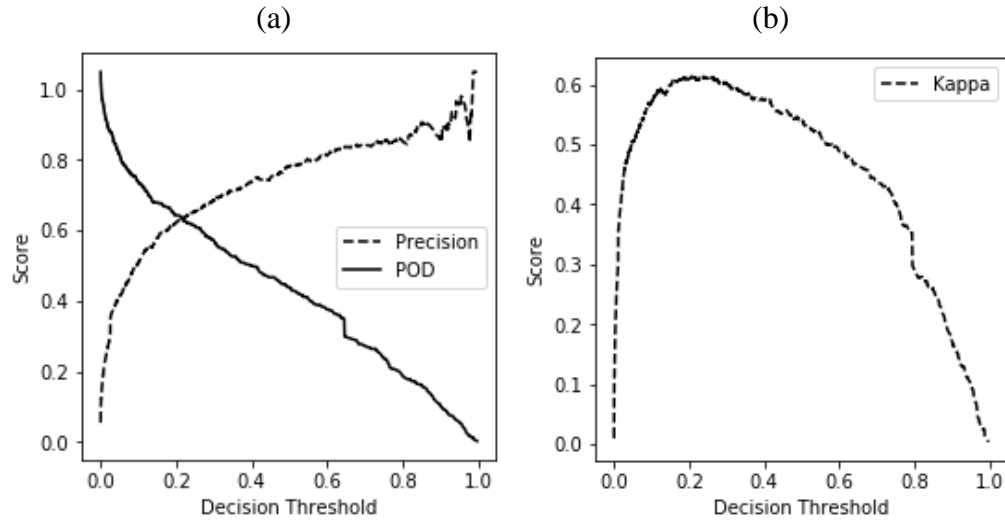
**Table 6.21: The descending value ranking of individual hyperparameter among the top 5 performed cases, and the corresponding conservativeness ranking scores in parentheses. The parameters with normal font are those favoring smaller values for conservativeness, and those *Italicized* favoring larger values.**

Name	M21	M31	M27	M25	M23
<i>m_neighbors</i>	<i>1 (2)</i>	<i>1 (2)</i>	<i>1 (2)</i>	<i>1 (2)</i>	<i>5 (1)</i>
k_neighbors	1 (1)	3 (3)	1 (1)	5 (5)	4 (4)
shrinkage	1 (1)	2 (2)	3 (3)	5 (5)	3 (3)
n_estimators	1 (1)	1 (1)	1 (1)	5 (5)	4 (4)
subsample	3 (3)	3 (3)	3 (3)	1 (1)	2 (2)
colsample_bytree	1 (1)	1 (1)	1 (1)	5 (5)	4 (4)
<i>reg_alpha</i>	<i>2 (1)</i>	<i>2 (1)</i>	<i>2 (1)</i>	<i>2 (1)</i>	<i>1 (5)</i>
<i>reg_lambda</i>	<i>1 (3)</i>	<i>1 (3)</i>	<i>5 (1)</i>	<i>1 (3)</i>	<i>4 (2)</i>
<i>gamma</i>	<i>1 (1)</i>	<i>1 (1)</i>	<i>1 (1)</i>	<i>1 (1)</i>	<i>1 (1)</i>
<i>min_child_weight</i>	<i>2 (1)</i>	<i>2 (1)</i>	<i>2 (1)</i>	<i>2 (1)</i>	<i>1 (5)</i>
max_depth	1 (1)	5 (5)	5 (5)	2 (2)	5 (5)
Total score	16	23	20	31	36

#### 6.3.1.4 Hyperparameters tuning for XGBoost

Similarly, to tune the decision threshold, Figure 6.11 (a) displays variations of precision and POD variations as functions of the decision threshold from 10-fold cross-validation in the training/validation data. The precision and POD curves cross each other around 0.2 of the threshold value, a relatively balanced point for POD and FAR. At the same point, the kappa scores shown in Figure 6.11 (b) is closer to the highest value, 0.61. As a result, 0.2 is selected as the decision threshold as before.

2104



2105

2106

2107

2108

**Figure 6.11: (a) Precision and POD score vs. decision threshold. (b) Kappa score vs. decision threshold**

2109

### 6.3.2 Model result on test data

2110

2111

2112

2113

2114

2115

Similar to LLE-SHIPS model, the evaluation of the prediction for DL-SHIPS model is on the test data only. The test confusion matrix for the model, before hyperparameter tuning (MB), and after hyperparameter tuning (MA) is displayed in Table 6.22. The result indicates that the hyperparameter tuning procedure does help the model performance.

2116

2117

**Table 6.22: Confusion matrix values after (before) hyperparameter tuning with the test data.**

	Predicted RI	Predicted non-RI	Actual
Actual RI	48 (29)	47 (66)	95
Actual non-RI	37 (31)	1465 (1471)	1502
Total Predicted	85 (60)	1512 (1537)	

2118

2119 Kappa, PSS, POD, and FAR are used for the model evaluation, and their values  
2120 for MB and MA are elaborated in Table 6.23. The POD and FAR values for MB and MA  
2121 cases demonstrated the importance of hyperparameter tuning. After tuning, POD  
2122 increases 65.6% from 0.305 to 0.505, while FAR decreases from 0.517 to 0.435, 15.9%.  
2123 The overall statistics PSS and kappa score also increased from 0.285 to 0.481 (68.8%)  
2124 and from 0.344 to 0.506 (47.1%), respectively, confirming the significant improvement  
2125 on RI prediction with the hyperparameter tuning procedure, and apparently, the model  
2126 was overfitted before tuning process with so many variables.  
2127

2128 **Table 6.23: Performance comparisons. MB and MA denote the models before and**  
2129 **after the hyperparameters in GMM-SMOTE and XGBoost are tuned.**

Model	Kappa	PSS	POD	FAR
MB	0.344	0.285	0.305	0.517
MA	0.506	0.481	0.505	0.435
Improvement MB	47.1%	68.8%	65.6%	-15.9%

2130

### 2131 6.3.3 Feature importance

2132 Similar to LLE-SHIPS model, the importance score could be derived from XGBoost  
2133 for the output of the data filters. However, in the DL-SHIPS model, since there are even  
2134 significantly more variables (each grid in each variable could be regarded as a feature)  
2135 than that of the input for the LLE data filter, it is even more computationally expensive  
2136 and impossible to implement the same feature importance evaluation approach

(permutation) as in LLE-SHIPS model, i.e., tracing back the importance of each ERA-Interim feature is almost impossible for DL-SHIPS model. Therefore, although we can evaluate the importance of the output of data filters, how to evaluate the contribution from each of the original ERA-Interim variables is a notoriously difficult task for deep learning networks, a.k.a, autoencoder network. Here we are roughly evaluating the importance of the ERA-Interim variables by calculating their summed importance score derived from the XGBoost classifier for each of the 14 ERA-Interim variables, as well as the averaged individual score for parameters associated with each variable. And for the variables with higher importance score, the feature level information from the individual 3D auto-encoder are traced back based on the feature map (Zeiler 2014), where the extracted information, for example, the geometric location, is visualized.

Table 6.24 displays the 10 most important variables among the 120 selected variables, including not highly correlated 72 SHIPS variables, and 48 variables extracted from the DL ERA-interim data filter. We can find that among the top 10, there are six SHIPS variables and four DL variables, and the total importance score for DL variables is 0.4119, while that of SHIPS variables is 0.5881. Therefore, the average score per SHIPS variable/ERA variable is 0.0082/0.0086. The fact that the average score for the SHIPS variable is less than that of ERA variables indicates that the ERA-Interim data filter has a similar importance score comparing to that of SHIPS variables; hence DL ERA-interim data filter is working efficiently.

2158 **Table 6.24: Variable importance in DL-SHIPS model.**

Variable	Importance	Description
BD12	0.019747	The past 12 hour intensity change
VMAX	0.017600	Maximum Surface Wind
SHRD	0.014810	850-200 hPa shear magnitude
DTL	0.014381	The distance to nearest major land
IRM1_5	0.013737	Predictors from GOES data (not time dependent) for r=100-300 km but at 1.5 hours before initial time
o31	0.013308	3 <sup>rd</sup> variable in o3
G150	0.013093	Temperature perturbation at 150 hPa due to the symmetric vortex calculated from the gradient thermal wind. Averaged from r=200 to 800 km centered on input lat/lon (not always the model/analysis vortex position) (deg C*10)
q7	0.013093	7 <sup>th</sup> variable in q
u3	0.012878	3 <sup>rd</sup> variable in u
q4	0.012878	4 <sup>th</sup> variable in q

2159

2160 Table 6.24 also indicates that BD12 has the largest importance score, 0.0197, and the  
 2161 second most important variable is VMAX. The third and fourth most important variables,  
 2162 SHRD and DTL. The fifth to tenth variables are IRM1\_5, o31, G150, q7, u3, and q4, and  
 2163 four of them are derived from the DL-interim data filter. o31 is the first variable extracted  
 2164 from o3's network, while q7, u3, and q4 are the seventh, third, and fourth variables of q,  
 2165 u, and q. The importance scores of all of the 120 variables can also be found in Table A6.

2166 Since the 3D auto-encoder model structure is different over 14 ERA-Interim  
 2167 variables, the summed importance score, which is the sum over all the output from the  
 2168 same network, for example, the summed importance score for r is the sum of the  
 2169 importance score over r1, r2, ..., r8, for each of the 14 variables, as well as the averaged

importance score on the non-zero, non-highly-correlated features. For example, the averaged importance score for  $r$  is the summed importance score for  $r$  divided by 8, are described in Table 6.25. Based on Table 6.25, we can find that  $q$ ,  $vo$ , and  $u$  are scored in the top 5 in terms of both summed score and the average score. Therefore, below we are looking at feature maps from the first layer of the networks for  $q$ ,  $vo$ , and  $u$  between an example RI and non-RI instances to roughly estimate what plays a more significant role to distinguish between RI and non-RI instances.

2177

**Table 6.25: Summed variable importance score, the number of non-zero, non-correlated features, the feature-wise averaged importance score, and its ranking for each ERA-Interim variable.**

Variable	The number of features	Summed Importance Score	Importance score rank	Average importance score	Average importance score rank
$q$	7	0.062	1	0.010843	3
$vo$	6	0.055	2	0.010583	4
$u$	6	0.055	3	0.00975	5
$v$	6	0.056	4	0.008483	7
$pv$	5	0.050	5	0.00882	6
$r$	8	0.049	6	0.004838	14
$ciwc$	2	0.042	7	0.0072	10
$o3$	1	0.039	8	0.0133	1
$cc$	2	0.020	9	0.006	12
$d$	1	0.017	10	0.0118	2
$t$	1	0.019	11	0.0082	8
$z$	1	0.014	12	0.0077	9

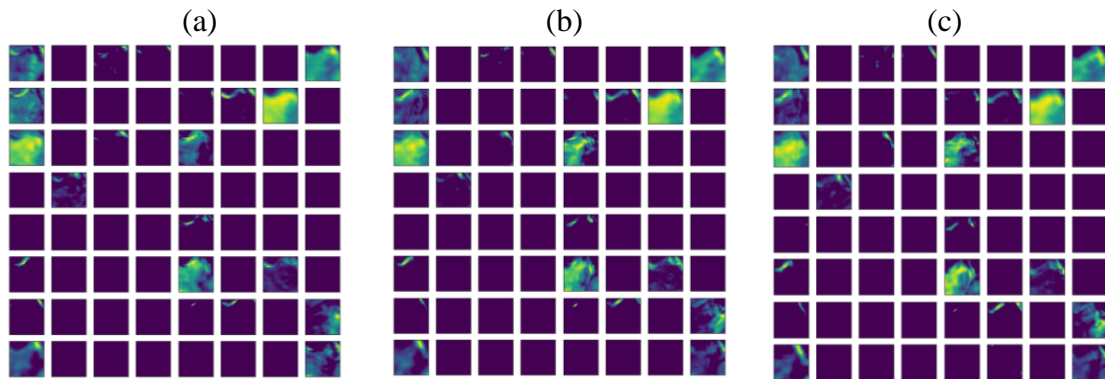
w	1	0.011	13	0.0071	11
clwc	1	0.008	14	0.0058	13

2181

2182        Figure 6.12 and Figure 6.13 represent examples of non-RI and RI instances features  
2183        extracted by the 3D autoencoder in all 3 channels in terms of variable relative humidity  
2184        (q), and based on the figures we can find that the extracted feature maps are sparse, with  
2185        only 22/64, and 24/64 non empty feature maps, where the deep blue feature maps  
2186        represents all the pixels have value 0, for non-RI and RI instances respectively. With the  
2187        limited available information, we can find the non-RI instances are extracting features  
2188        from the northeast (upper left) of the center or the whole domain, while RI instances are  
2189        extracting features mainly on the southeast (bottom right) of the center. So we can  
2190        conclude that the relative humidity (q) in northeast of the center is more important in RI  
2191        occurrence.

2192

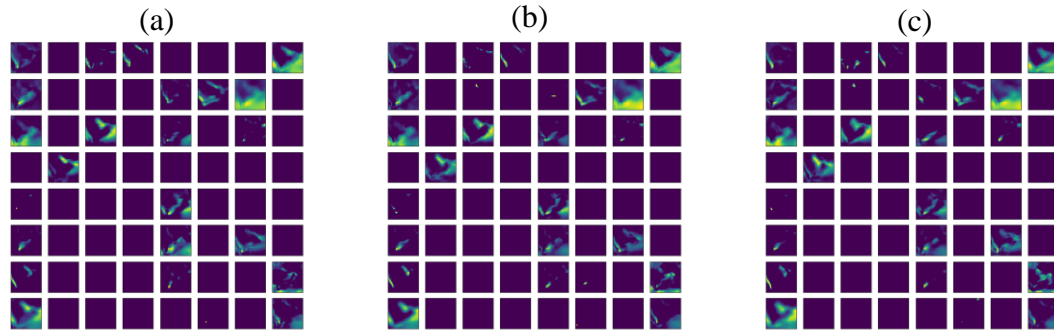
2193



2194

2195 **Figure 6.12: 3 channels, 64 feature maps for the first layer (dimension: 3 (channel)**  
2196 **\* 64 (feature map) \* 30 (feature map dimension) \* 30 (feature map dimension)) of**  
2197 **the network that is immediate after the input layer (dimension: 37 (pressure level) \***

2198 4 (-18h, -12h, -6h, and 0h) \* 33 (vertical grid) \* 33 (horizontal grid)) for variable q  
 2199 with its network structure displayed in Figure 6.8a. This is an example for a non-RI  
 2200 instance, and (a) Non-RI in channel 1. (b) Non-RI in channel 2. (c) Non-RI in  
 2201 channel 3, and the sequence of the channel does not matter. In each channel, there  
 2202 are 64 (8 in the row and 8 in the column) feature maps, and each feature map has 30  
 2203 (pixels) \* 30 (pixels) dimension. Deep blue implies the value in that pixel is 0, and  
 2204 the brighter the color is, the high the value in that pixel.  
 2205  
 2206

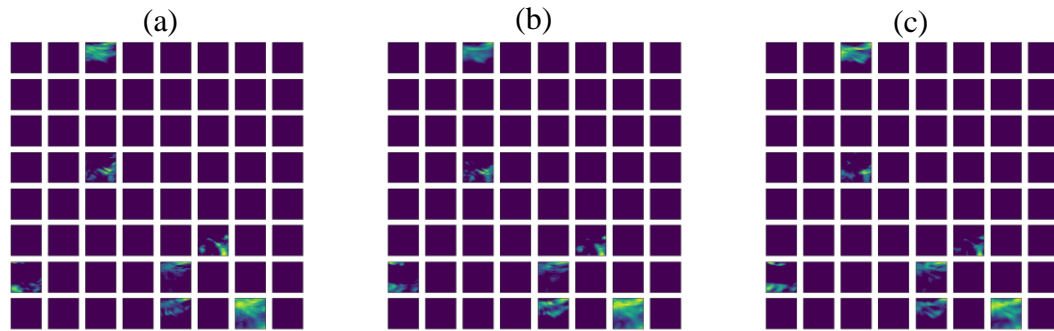


2207  
 2208 **Figure 6.13: Same as Figure 6.12 but for a RI instance (a) RI in channel 1. (b) RI**  
 2209 **in channel 2. (c) RI in channel 3.**  
 2210

2211 Figure 6.14 describes the examples of non-RI instances and Figure 6.15 describes the  
 2212 RI features extracted by the 3D autoencoder in all 3 channels in terms of variable relative  
 2213 vorticity (vo), and we can find that the feature map is even more sparse comparing to  
 2214 variable relative humidity (q), with only approximately 7/64, and 7/64 non empty feature  
 2215 maps for both situations. Among them, 5/7 feature maps are towards north (top; of the  
 2216 center) and the rest 2 feature maps are towards south (bottom; of the center). In  
 2217 comparison, RI instance indicates that 5 are towards south (lower) and 2 are towards  
 2218 north (of the center). So we can conclude that relative vorticity (vo) in the south of the  
 2219 center are more important in RI occurrence.  
 2220

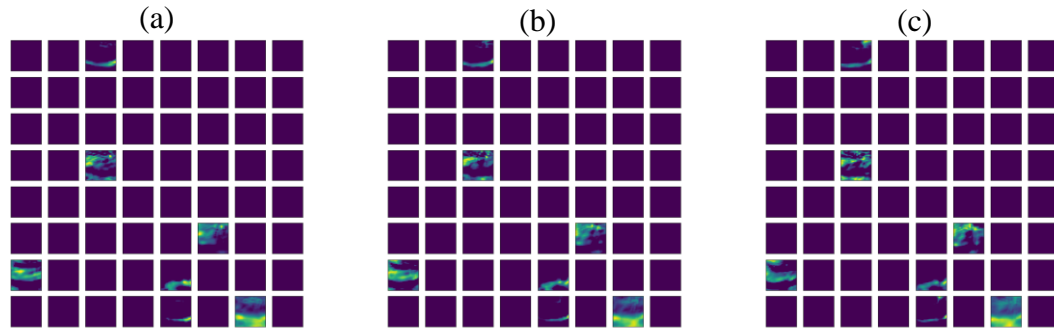


2221



2222

2223 **Figure 6.14: Same as Figure 6.12 but for variable vo with its network structure in**  
2224 **Figure 6.8a in a non-RI instance: (a) non-RI in channel 1. (b) non-RI in channel 2.**  
2225 **(c) non-RI in channel 3.**



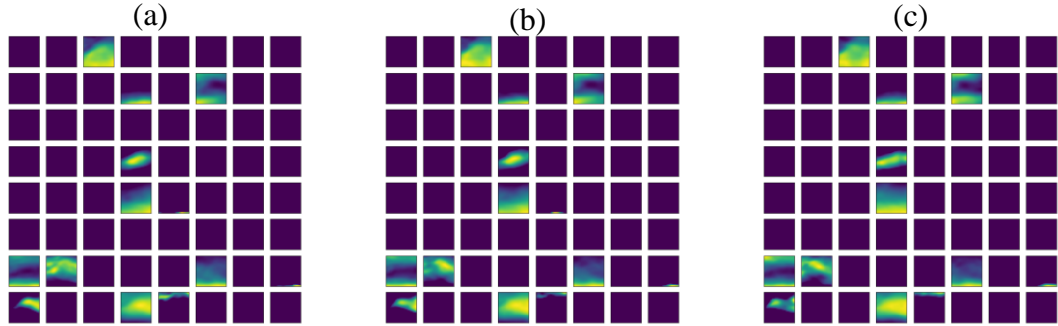
2226

2227 **Figure 6.15: Same as Figure 6.12 but for variable vo with its network structure in**  
2228 **Figure 6.8a in 3 channels in a RI instance: (a) RI in channel 1. (b) RI in channel 2.**  
2229 **(c) RI in channel 3.**

2230

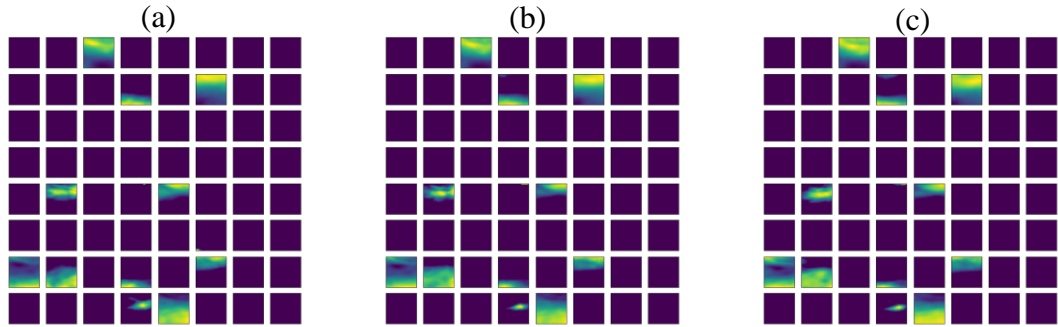
2231 Similarly, based on Figure 6.16 and 6.17, 11/64 and 11/64 are non-empty feature  
2232 maps for non-RI and RI instances in variable eastward wind (u), non-RI instance has 6, 1,  
2233 and 4 feature maps toward south, north, and east of the center. In comparison, RI instance  
2234 has 6, 4, and 1 feature maps concentrating north, south, and east of the center, which  
2235 indicated that eastward wind (u) in the north of the center is more possible to result in RI.

2236



2237

2238 **Figure 6.16:** Same as Figure 6.12 but for variable u with its network structure in  
2239 **Figure 6.8a** in a non-RI instance: (a) non-RI in channel 1. (b) non-RI in channel 2.  
2240 (c) non-RI in channel 3.



2241

2242 **Figure 6.17:** Same as Figure 6.12 but for variable u with its network structure in  
2243 **Figure 6.8a** in 3 channels in a RI instance: (a) RI in channel 1. (b) RI in channel 2.  
2244 (c) RI in channel 3.

2245

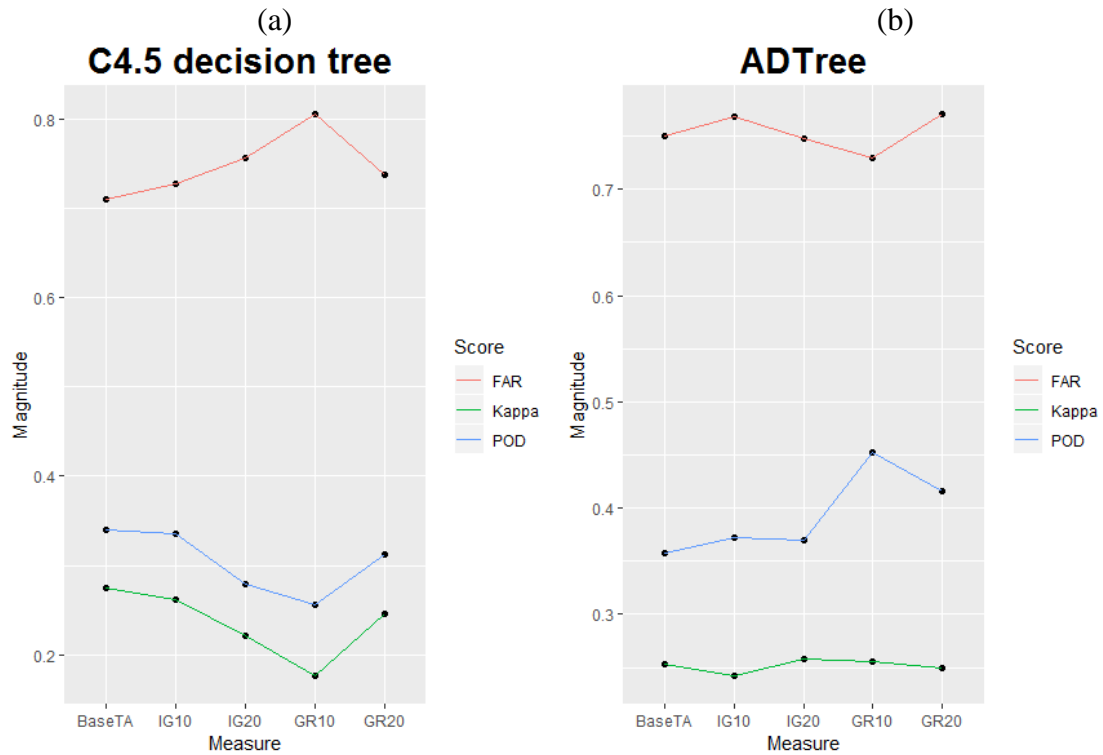
## 2246 **6.4 Model performance comparison**

2247 Two works, i.e., the best model in Y16 and KRD15, which outperforms almost all  
2248 of the other works in the RI prediction, are used to compare with the performance of the  
2249 COR-SHIPS model, LLE-SHIPS model, and DL-SHIPS model.

#### 2250    **6.4.1 Model performance in Yang (2016) and Kaplan et al. (2015)**

2251       In Y16, experiments are conducted using C4.5 decision tree (Quinlan, 1993),  
2252   alternating decision tree (ADTree; Freund and Mason, 1999), random forest (Breiman,  
2253   2001), classification and regression tree (CART; Breiman et al. 2017), logistic model tree  
2254   (LMT; Landwehr et al. 2005), the repeated incremental pruning to produce error  
2255   reduction (RIPPER; Cohen 1995), function-based classification such as support vector  
2256   machines with sequential minimal optimization (SMO; Platt 1999), naïve Bayes scheme  
2257   (Tan et al. 2015), and the decision tree with naïve Bayes classifiers at the leaves  
2258   (NBTree; Kohavi 1996) with cost ratio 4.6 to predict RI. Figure 6.18 shows the best  
2259   performed two classifiers C4.5 decision tree and ADTree, where measures in Figure 6.18  
2260   is defined as the different groups of variables selected by different variable selection  
2261   criteria. The performances of ADTree over all measures are more robust since its kappa  
2262   scores over different dataset is more stable than that of C4.5 decision tree, although the  
2263   best test kappa (27.5%) is achieved by C4.5 decision tree.

2264



2265

2266

2267

2268

**Figure 6.18: Kappa, POD, and FAR for (a) C4.5 decision tree. (b) ADTree. Data are from Y16.**

2269

2270

2271

2272

2273

2274

2275

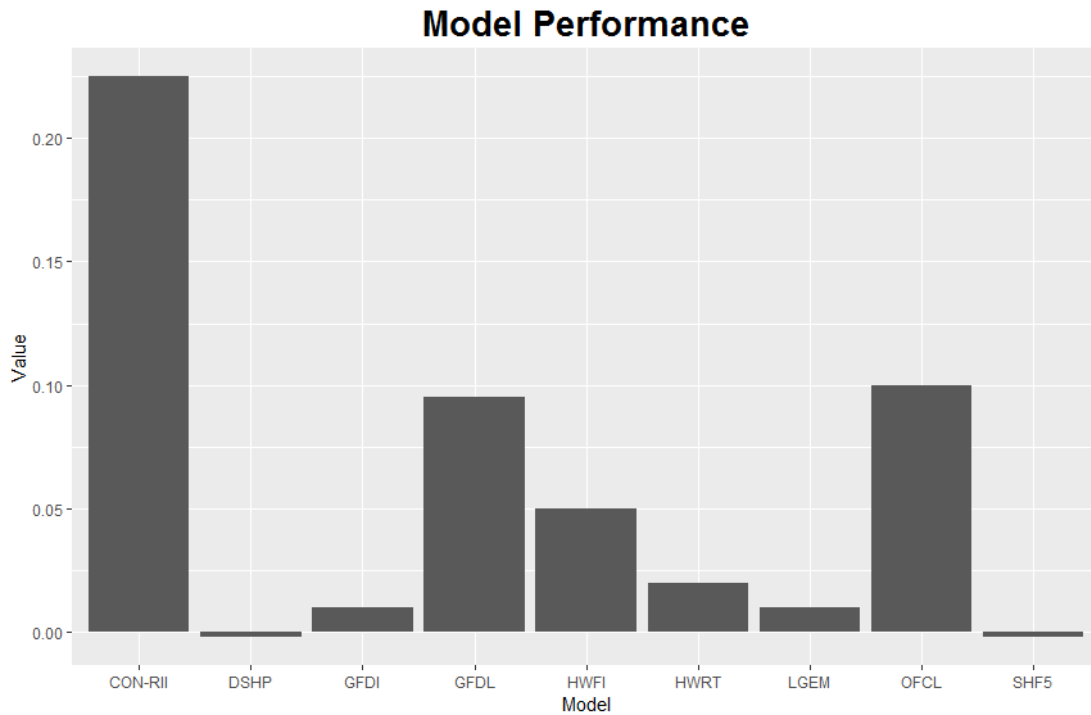
2276

2277

In KRD15, CON-R11 outperforms all other operational models, i.e., the 5-day SHIFOR model (SHF5; Knaff et al. 2003), the decay version of the SHIPS model (DSHP; DeMaria et al. 2005), the logistic growth equation model (LGEM; DeMaria 2009), the Geophysical Fluid Dynamical Laboratory (GFDL) hurricane prediction model early (GFDI) and late (GFDL) versions (Kurihara et al. 1998) and the Hurricane Weather Research and Forecasting Model early (HWFI) and late (HWRF) versions (Tallapragada et al. 2014), and the NHC official forecast (OFCL) in RI prediction at the threshold 30 knots in the 24-hour lead-time in terms of Peirce's skill score (PSS) with approximately 0.225 for TC cases in 2008–13. The result details of all the models are presented in

2278 Figure 6.19, the CON-RII with PSS 0.225 performs little over two times than the second  
 2279 best-performed model (OFCL).

2280



2281

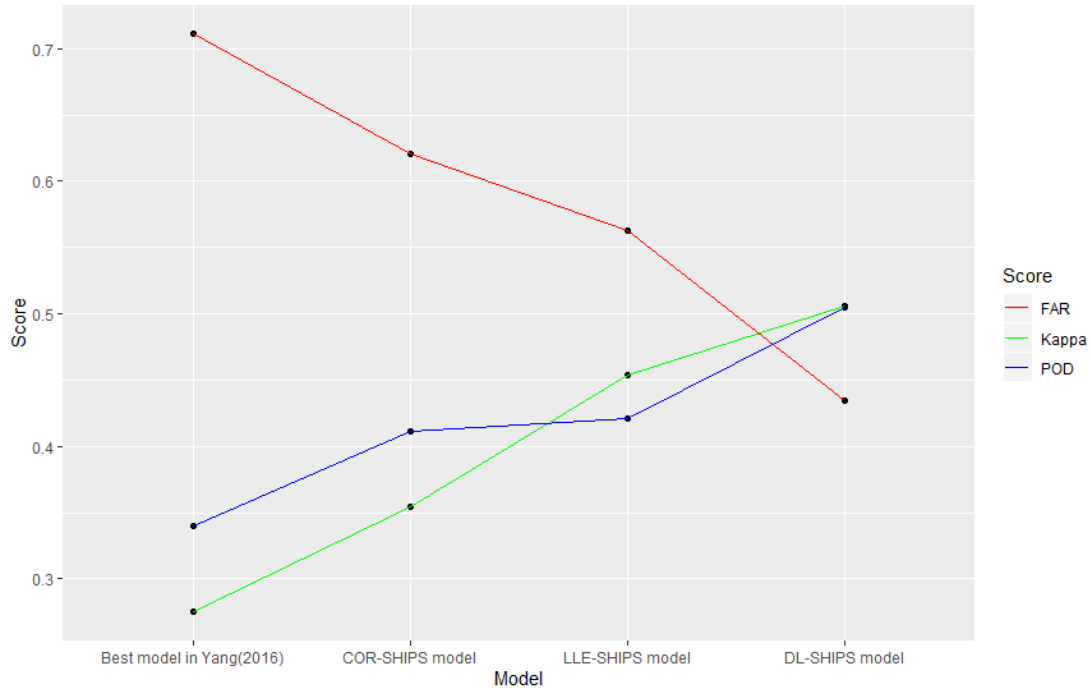
2282 **Figure 6.19: Different model's performance regarding Peirce's skill score (PSS)**  
 2283 **based on data from KRD15.**

2284

## 2285 6.4.2 Model comparison

2286 Figure 6.20 displays the model performance comparison between the best model in  
 2287 Yang (2016) (a.k.a. Y16) and the three newly developed models in this study for kappa  
 2288 score, POD, and FAR. The performance of the COR-SHIPS model, LLE-SHIPS model,  
 2289 and DL-SHIPS model is significantly better than that of the best model in Yang (2016).

Details are listed in Table 6.26, and we can find that the performance (kappa score, POD, and FAR) improvement by using the entire SHIPS database with regard to Y16 is medium with 28.7%, 20.9%, and -12.7%, respectively. If we use ERA-Interim data in addition to the SHIPS database, we achieve significant improvement by at least 65.1%, 23.8%, and -20.8% in terms of kappa score, POD, and FAR (the smaller, the better).



**Figure 6.20: Model performance comparison: Model's test kappa, FAR, and POD score in the best model in Yang (2016), SHIPS model, LLE-SHIPS model, and DL-SHIPS model.**

**Table 6.26: Performance comparison between our models, and Y16 and KRD15.**

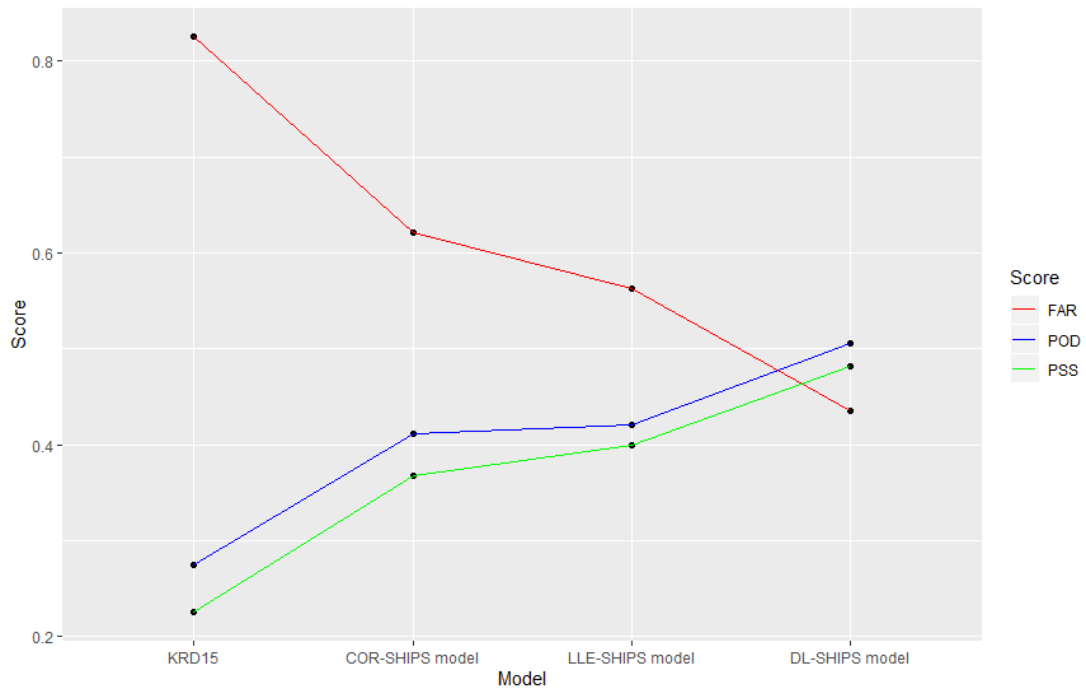
Model	Kappa	PSS	POD	FAR	Improvement Yang			Improvement KRD15		
					Kappa	POD	FAR	PSS	POD	FAR
COR-SHIPS	0.354	0.368	0.411	0.621	28.7%	20.9%	-12.7%	63.6%	49.5%	-24.7%

LLE-SHIPS	0.454	0.399	0.421	0.563	65.1%	23.8%	-20.8%	77.3%	53.1%	-31.8%
DL-SHIPS	0.506	0.481	0.505	0.435	84.0%	48.5%	-38.8%	114.0%	83.6%	-47.3%
Y16	0.275	NA	0.340	0.711						
KRD15	NA	0.225	0.275	0.825						

2303

2304        Figure 6.21 displays the model performance comparison among CON-RII in Kaplan  
2305 et al. (2015) (a.k.a. KRD15), COR-SHIPS model, LLE-SHIPS model, and DL-SHIPS  
2306 model in terms of PSS, POD, and FAR. The performance of the SHIPS model, LLE-  
2307 SHIPS model, and DL-SHIPS model is significantly better than that of KRD15 as at least  
2308 63.6%,49.5%, 24.7% improvement in PSS score, POD, and FAR, and more details are  
2309 elaborated in Table 6.26. The performance improvement in KRD15 for our 3 models are  
2310 more significant than that in Y16, especially we are considering that the performance of  
2311 the model in KRD15 is evaluated in the training dataset, and the performance in Y16 and  
2312 our models are evaluated based on the test dataset.

2313



**Figure 6.21: Model's test PSS, FAR, and POD score in KRD15 (Kaplan et al., 2015), COR-SHIPS model, LLE-SHIPS model, and DL-SHIPS model.**

The performance comparison between our 3 models – COR-SHIPS, LLE-SHIPS, and DL-SHIPS are elaborated in Table 6.27. We can find a moderate improvement was made from the COR-SHIPS model to the LLE-SHIPS model in terms of kappa, PSS, POD, and FAR of 28.2%, 8.4%, 2.4%, and 9.3%, which indicates that most of the near center information has already been explored/extracted by variables in SHIPS database. In comparison, a significant change is made from the COR-SHIPS model to the DL-SHIPS model in terms of kappa, PSS, POD, and FAR of 42.9%, 30.7%, 22.9%, and 30.0%, which indicates that the DL-SHIPS model catches large-scale information not only as averages (from SHIPS database) but also the variations (from DL model). With



almost half (48 vs. 90) variables generated from ERA-Interim data filter, the DL-SHIPS model still outperforms the LLE-SHIPS model significantly, which not only indicates that DL ERA-interim data filter is extracting more efficient information, but also shows that using large-scale ERA variables provides much more information than incorporating near core ERA variables only in RI prediction.

**Table 6.27: Performance comparison between 3 models developed in this study, ‘X’ in the table indicates that not available value.**

Model	Kappa	PSS	POD	FAR	Improvement COR-SHIPS				Improvement LLE-SHIPS			
					Kappa	PSS	POD	FAR	Kappa	PSS	POD	FAR
COR-SHIPS	0.354	0.368	0.411	0.621	X	X	X	X	X	X	X	X
LLE-SHIPS	0.454	0.399	0.421	0.563	28.2%	8.4%	2.4%	-9.3%	X	X	X	X
DL-SHIPS	0.506	0.481	0.505	0.435	42.9%	30.7%	22.9%	-30.0%	11.5%	20.6%	20.0%	-22.7%

## **6.5 Feature importance**

### **6.5.1 Feature importance comparison between COR-SHIPS model, LLE-SHIPS model, and DL-SHIPS model**

Table 6.28 lists the 36 most important variables for each of the three models. Learning from Table 6.28, 7 out of top 10 in the LLE-SHIPS model is in the top 10 of the COR-SHIPS model, and in all the SHIPS variables identified in top 10 in the DL-SHIPS model are included in that of the COR-SHIPS model and LLE-SHIPS model. This indicates that all three models identify similar important variables.

2344

2345 **Table 6.28: Top 36 most important variables in COR-SHIPS model, LLE-SHIPS**  
 2346 **model, and DL-SHIPS model.**

Rank	COR-SHIPS model variable	COR-SHIPS model variable Importance Score	LLE-SHIPS model variable	LLE-SHIPS model variable Importance Score	DL-SHIPS model variable	DL-SHIPS model variable Importance Score
1	BD12	0.0362	BD12	0.0188	BD12	0.0197
2	DTL	0.0217	VMAX	0.0167	VMAX	0.0176
3	CFLX	0.0207	DTL	0.0138	SHRD	0.0148
4	SHDC	0.0206	SHRD	0.0130	DTL	0.0144
5	G150	0.0205	TWXC	0.0115	IRM1_5	0.0137
6	jd	0.0204	G150	0.0113	o31	0.0133
7	VMAX	0.0199	VMPI	0.0113	G150	0.0131
8	IRM1_5	0.0199	REFC	0.0113	q7	0.0131
9	PW08	0.0191	TGRD	0.0111	u3	0.0129
10	VMPI	0.019	IRM1_5	0.0107	q4	0.0129
11	SHTD	0.0187	IR00_12	0.0107	G200	0.0129
12	IR00_12	0.0183	V300	0.0105	vo3	0.0127
13	HE07	0.018	VVAC	0.0105	REFC	0.0124
14	MTPW_2	0.0177	G200	0.0103	vo5	0.0122
15	XD18	0.0177	PEFC	0.0096	vo8	0.0120
16	SHTS	0.0175	MTPW_2	0.0096	PEFC	0.0120
17	PW14	0.0173	XDTX	0.0095	d3	0.0118
18	TWXC	0.0172	PSLV_1	0.0095	CFLX	0.0116
19	R000	0.0168	T150	0.0095	PSLV_3	0.0116
20	V300	0.0167	CFLX	0.0094	T150	0.0114
21	OAGE	0.0165	HIST_2	0.0094	jd	0.0114
22	PSLV_1	0.0162	HE07	0.0092	R000	0.0114
23	Z850	0.0161	SHTS	0.0091	TWXC	0.0112
24	SHRS	0.0161	PSLV_3	0.0089	u8	0.0112
25	SDDC	0.0157	SHTD	0.0085	PW08	0.0112
26	VVAC	0.0156	G250	0.0085	q3	0.0112
27	PSLV_5	0.0156	CD26	0.0085	XDTX	0.0112
28	TGRD	0.0154	lle84	0.0082	CD26	0.0109
29	T150	0.0153	EPSS	0.0082	q8	0.0109
30	CD26	0.0153	R000	0.0077	pv3	0.0107
31	TADV	0.0152	SDDC	0.0076	v4	0.0107
32	V850	0.0151	IRM3_19	0.0075	r1	0.0105
33	PSLV_4	0.0148	RD26	0.0075	u1	0.0101
34	Z000	0.0145	PW08	0.0074	q5	0.0099
35	REFC	0.0145	SHRS	0.0074	IR00_12	0.0099
36	RD26	0.0142	NDTX	0.0074	vo4	0.0099

2347

2348       The 6 common variables are BD12, VMAX, SHRD, DTL, IRM1\_5, and G150.  
2349       Although there is no LLE engineered variables in the top 10 in the LLE-SHIPS model,  
2350       the importance score for the most important variable, BD12, is only approximately 11%  
2351       higher than the second important variable VMAX. In comparison, the importance score  
2352       for BD12 in the COR-SHIPS model has 50% higher importance score than the second  
2353       most important variable, DTL, which indicates that the performance of the LLE-SHIPS  
2354       model does not heavily rely on one individual variable. In addition, the performance of  
2355       the LLE-SHIPS model has approximately 28% higher kappa score than that of the COR-  
2356       SHIPS model, which indicates that the ERA-interim data filter efficiently extracts  
2357       important near center features that help the RI prediction.

2358       DL-SHIPS model has 120 variables, which is 42 (33%) variables less than that of the  
2359       LLE-SHIPS model, and the performance of the DL-SHIPS model (0.506 kappa value) is  
2360       approximately 11% better than that of the LLE-SHIPS model (0.454). The fact that with  
2361       much fewer variables, the DL-SHIPS model is performing much better than the LLE-  
2362       SHIPS model, indicates that DL ERA-Interim data filter extracts large-scale features that  
2363       is more representative of RI than that of LLE ERA-Interim data filter, i.e., near center  
2364       feature. However, unlike the LLE-SHIPS model, the DL-SHIPS model extracts features  
2365       from each variable separately, and interaction between different variables are ignored. If  
2366       we also extract the interaction between terms in the large-scale dataset, we can get better  
2367       performance.

2368       Since we have 72 variables in the COR-SHIPS model, we further compare the top  
2369       36, i.e., 50% number of variables in the COR-SHIPS model, for the COR-SHIPS model,

2370 LLE-SHIPS model, and DL-SHIPS model that described in Table 6.28. LLE-SHIPS  
2371 model has 34 SHIPS variables, and 2 ERA-Interim variables, and 25 of 34 variables  
2372 (73.5%) are overlapped with that of the COR-SHIPS model, which also support the fact  
2373 that LLE ERA-interim data filter extracts important near center features that help the RI  
2374 prediction, although the new features seem not as efficient since there are only 2 in the  
2375 top 36. However, with 162 variables in total, the LLE-SHIPS model has the total  
2376 importance score 0.259 for the overlap variables, and only with 72 variables, COR-  
2377 SHIPS model has 0.456 for the same overlap variables, which almost double that of the  
2378 LLE-SHIPS model. The reason is there are significantly more variables in the LLE-  
2379 SHIPS model, and it is not surprising that the importance score for the overlapped  
2380 variables is significantly different in the COR-SHIPS model, and LLE-SHIPS model.

2381 In contrast, the DL-SHIPS model has 19 SHIPS variables, and 17 ERA-Interim  
2382 variables. Among all these variables, 14 of 19 SHIPS variables (73.7%) are overlapped  
2383 with that of the COR-SHIPS model, and the summed importance scores are 0.276, and  
2384 0.180 respectively for the COR-SHIPS model and DL-SHIPS model, which indicates  
2385 DL-SHIPS model relies less on SHIPS variables. 18 of 19 SHIPS variables (94.7%) are  
2386 overlapped with that of the LLE-SHIPS model, and the summed importance scores are  
2387 0.199, and 0.231 respectively for the LLE-SHIPS model and DL-SHIPS model, almost  
2388 same.

2389 With significantly less ERA-Interim variables and almost the same number of SHIPS  
2390 variables, we can conclude that the DL ERA-interim data filter is efficient at either  
2391 improving the prediction accuracy, or extracting new variables at the large-scale. There

2392 might be two reasons for this. The first reason might be large-scale features are more  
2393 efficient than the near center features in RI prediction, which matches domain scientist  
2394 experience, because most of the SHIPS variables are large scale variables. The second  
2395 reason might be compared with the LLE data filter, the DL data filter has multiple  
2396 convolutional layers, which distilled the all levels of large-scale information  
2397 comprehensively that may be ignored by domain scientists for a long time.

2398

## 2399 **6.5.2 Feature importance comparison between previous studies and this study**

2400 A two-side t-test is used for variable selection in KD03, and the RII model was  
2401 built based on the five variables, DVMX (Intensity change during the previous 12 h),  
2402 SST, POT (Maximum potential intensity (MPI) – maximum sustained surface wind  
2403 speed), SHR (850-200-hPa vertical shear averaged from  $r = 200$ -800 km), and RHLO  
2404 (850-700-hPa relative humidity averaged from  $r=200$ -800 km), which are found  
2405 significant in a 99.9% level and with the highest individual RI prediction power. In the  
2406 first 10 importance variables identified by the COR-SHIPS model, BD12 (ranked 1<sup>st</sup>),  
2407 SHRD (4<sup>th</sup>), VMPI (10<sup>th</sup>), and VMAX (7<sup>th</sup>) (POT = VMPI - VMAX) are consistent with  
2408 the selected variables in KD03. The missed variables in the top ten compared with the top  
2409 five in KD03 are SST and RHLO. SST is highly correlated with the selected E000, which  
2410 is listed 52<sup>th</sup> in the importance ranks. RHLO is highly correlated with RHMD, which is  
2411 listed 57<sup>th</sup> in the importance ranks (Table A1 and Table A4).

2412 Compared with variables selected by KD03, in KDK10, SST is removed and 4  
2413 additional variables, D200 (time averaged 200-hPa divergence within a 1000-km radius),

2414 OHC (time averaged oceanic heat content), SDBT (STD of GOES-IR BT ( $t = 0h$ ) within  
2415 a 50–200-km radius), and PX30 (the percentage area from 50- to 200-km radius covered  
2416 by IR cloud-top BT of  $-30^{\circ}\text{C}$  or colder), are added. Among the four new variables, D200  
2417 is ranked 44<sup>th</sup> with a 0.0131 importance score (Table A4) in the COR-SHIPS model. The  
2418 OHC related parameters include COHC, NOHC, and RHCN, and among them, the  
2419 highest importance score is achieved by COHC, which is highly correlated with CD26  
2420 ranked 30th with a score of 0.0153. The PX30 is corresponding to IR00\_8, which is  
2421 highly correlated with IRM1\_16 ranked 50<sup>th</sup> with a 0.0119 score value. The only caught  
2422 new KRD10 variable in our top ten in the COR-SHIPS model is the SDBT by IRM1\_5  
2423 (ranked 8<sup>th</sup>), representing GOES BT STD within the 100-300 km around the TC centers  
2424 but 1.5 hours before the current time.

2425 KRD15 replaced RHLO with TPW (Percentage of an area with  $\text{TPW} < 45 \text{ mm}$   
2426 within a 500-km radius and  $\pm 45^{\circ}$  of the upshear SHIPS wind direction ( $t = 0h$ )), and  
2427 PX30 with PC2 (second principal component of GOES-IR imagery within a 440-km  
2428 radius ( $t = 0h$ )), and added 2 new variables, ICDA (Inner-core dry-air predictor (time  
2429 avg)), and VMX0 (Max sustained wind ( $t = 0h$ )), comparing with variable used in  
2430 KDK10. Among the four new variables, VMX0 is consistent with VMAX, ranked 7<sup>th</sup> in  
2431 the COR-SHIPS model importance list. ICDA is not directly included in SHIPS data, but  
2432 the related parameter found is CFLX, the dry air predictor except for a factor of VMX0 in  
2433 KRD15, and CFLX is ranked the 3<sup>rd</sup> in the top 10 parameter list. The definition of TPW  
2434 is the same as MTPW\_19 in the SHIPS, which ranked only 37<sup>th</sup> with an importance score  
2435 of 0.014. The PC2 equivalent parameter in SHIPS is PC00, which ranked only the 70<sup>th</sup>.

2436 In summary, variables used by KD03, KDK10, and KRD15 for RI prediction are  
2437 mostly caught up with our top 10 variables in the COR-SHIPS model. The missed  
2438 variable in KD03 is RHLO and SST, and RHLO was actually replaced by TPW later  
2439 (KRD15), and TPW is ranked 37<sup>th</sup> in our list, much more important than the RHLO via  
2440 the highly correlated RHMD at the 57<sup>th</sup> place. Among the 4 newly added parameters in  
2441 KDK10, three, OHC, D200, and PX30, are outside the top 10 list. There are several  
2442 variables in SHIPS representing the OHC. The most important one is found to be  
2443 climatological OHC via the highly correlated parameter CD26 at the 30<sup>th</sup> rank. KRD15  
2444 mentioned that OHC works well only when the other two variables, POT and ICDA are  
2445 included in a model. D200 was introduced to SHIPS in 1998 (DeMaria and Kaplan  
2446 1999), but it was eliminated in 2001 and added back in 2002 (DeMaria et al. 2005).  
2447 DeMaria et al. (2005) also found that the role of this divergence in TC intensity  
2448 forecasting is sensitive to the data sources. Therefore, it is not very unusual if this model  
2449 did not rank this predictor high. The last parameter not in the top 10 list, PX30, was  
2450 replaced by PC2 in KRD15. Actually, PC2 is ranked 70<sup>th</sup> in this study, and it is hard to  
2451 interpret the result. It is very unfortunate that the GOES-IR principal components were  
2452 mistreated initially in this work, and we missed the opportunity to rank the importance of  
2453 other PCs among the first nine PCs. The other missed parameter in KRD15 is the TPW,  
2454 ranked only 37<sup>th</sup>. It is plausible that the humidity effects are also reflected in the 3<sup>rd</sup>  
2455 ranked parameter CFLX.

2456 Comparing to the COR-SHIPS model, the LLE-SHIPS model not only employs  
2457 SHIPS parameters but also ERA-Interim near center parameters for predicting RI, and we

2458 divide the ERA-Interim parameters into different groups with highly correlated  
2459 parameters to evaluate the group importance instead of evaluating the importance for  
2460 each parameter. Based on the top 5 important groups, we can find that wind speed,  
2461 especially at 400, 450, and 1000 hPa play a significant role in RI prediction. Another  
2462 important information we get, potential vorticity at 1000 hPa, although possibly the  
2463 importance score is over-estimated, is more important than the most important SHIPS  
2464 variable, BD12. Finally, vertical pressure speed is also found important in RI prediction.

2465 In addition to the LLE-SHIPS model, the DL-SHIPS model, which adopts the large-  
2466 scale ERA-Interim information, further improves the performance of COR-SHIPS  
2467 significantly. Relative humidity (q), relative vorticity (vo), and eastward wind (u) are  
2468 found to be top 3 important variables, but evaluating the contribution from each of the  
2469 original ERA-Interim parameters is a notoriously difficult task for deep learning  
2470 networks. So we roughly evaluating their 3D auto-encoder first layer weights, and find  
2471 out that RI instances tend to have higher weights in southeast humidity (q), north relative  
2472 vorticity (vo), and north eastward wind (u), where the direction is with regard to the TC  
2473 center  
2474



2475

## CHAPTER 7 CONCLUSION AND DISCUSSION

2476       To improve RI prediction with modern techniques, this study constructs a well-  
2477 tailored artificial intelligence (AI) system that goes back to the SHIPS database, the most  
2478 complete dataset with parameters known to be related to TC intensity changes, as well as  
2479 the ERA-Interim dataset, the best reanalysis product at the moment, to extract  
2480 information from a more complete set of variables. This system consists of four major  
2481 components, data filters to remove variables unrelated to RI, reduce variables among  
2482 highly correlated variables, screen out variables with high missing value rates, and  
2483 engineer/extract a reduced set of variables from the high dimensional variable space; a  
2484 customized sampler to upsample the minority (RI) instances and to downsample majority  
2485 instances simultaneously by a GMM-SMOTE sampler; a very powerful state-of-the-art  
2486 classifier, the XGBoost, to classify instances into RI and non-RI and to evaluate variable  
2487 importance based on the information gain; a hyperparameter tuning procedure tweaking  
2488 hyperparameters appearing in all of the three above components, within pre-defined value  
2489 ranges.

2490       Based on the AI system shown in Figure 1.1, three models, the COR-SHIPS  
2491 model, the LLE-SHIPS model, and the DL-SHIPS model, are developed. The COR-  
2492 SHIPS model only employs SHIPS data and is the continued work of Y16. Comparing  
2493 with Y16, the COR-SHIPS model adopts a different data filter, oversamples RI instances,  
2494 employs a more powerful classifier, and tunes their hyper-parameters to improve the  
2495 performance.

2496           However, the COR-SHIPS model still only employs the SHIPS dataset, which is  
2497 largely based on expert experiences. Since the mechanism of the tropical cyclone RI is  
2498 unknown, the knowledge from the domain scientist may not be comprehensive, which  
2499 indicates some important variables may not be included in the SHIPS dataset. Therefore,  
2500 ECMWF ERA-Interim reanalysis data are used to improve the performance of the COR-  
2501 SHIPS model. Two automatic feature extraction approaches, local linear embedding  
2502 (LLE) and deep learning (DL), are used to extract features from near center data (small  
2503 scale) and large-scale data respectively to create the LLE-SHIPS model and DL-SHIPS  
2504 model.

2505           The entire dataset is split into training/validation and test set, where the former is  
2506 used to fit our model and to tune the hyperparameters, and the latter is used for the  
2507 performance evaluation and comparison. The performance of our model on the test data  
2508 shows improvement in the RI prediction with hyperparameter tuning.

2509           It is found that our three model outperforms Y16 by 20.9%, 23.8%, and 48.5%,  
2510 and KRD15 by 49.5%, 53.1%, and 83.6% on POD, while reducing the FAR by 12.7%,  
2511 20.8%, and 38.8% comparing with Y16, and 24.7%, 31.8%, and 47.3% with KRD15  
2512 respectively. Our model also improves the kappa score of 28.7%, 65.1%, and 84.0% vs.  
2513 Y16 and the PSS 63.6%, 77.3%, and 114% against KRD15. With the difficulties in RI  
2514 prediction and the slow improvement rates in previous studies (KD03, KDK10, KRD15,  
2515 Y16), we believe the improvement by this work is substantial. The significant  
2516 improvement made by the three models also challenges the mainstream point of selecting  
2517 only a few variables fitting in the simple model for the prediction, i.e., involving more

2518 variables in the complicated model with high penalty terms is better than a simple model  
2519 with few variables.

2520 The variable importance is also evaluated, and BD12, the past 12-hour intensity  
2521 change, is found to have the largest importance score and contributes more than other  
2522 variables in all three models, and the common 6 variables in top 10 are BD12, VMAX,  
2523 SHRD, DTL, IRM1\_5, and G150.

2524 Previous important variables for RI prediction are determined by the significance test  
2525 in KD03, KDK10, and KRD15, and most of them are consistent with top 10 variables in  
2526 our three models with some exceptions. The variables in the top 10 list but not considered  
2527 in other RI studies may be helpful for future RI studies, especially for DLT, jd, G150,  
2528 PW08, TWXC, REFC, and TGRD. The additional significant variables identified by the  
2529 ERA-Interim data filter are the wind speed, especially at 400 and 450 hPa, potential  
2530 vorticity at 1000 hPa, vertical pressure in the near center and southeast humidity (q),  
2531 north relative vorticity (vo), and north eastward wind (u) in the large-scale and those  
2532 would help understand the mechanism of the TC intensification.

2533 COR-SHIPS, LLE-SHIPS, and DL-SHIPS model designed in this study performed  
2534 significantly better than most of the previous works such as KRD15, and Y16. Although  
2535 we are able to evaluate the importance of the output of different data filters, and we can  
2536 somehow even trace back the importance to the feature level, the feature level importance  
2537 is not accurate. Accurately tracing back the feature level importance is related to  
2538 interpretability or explainability of the complicated machine learning model, which is still  
2539 a challenging problem in the AI field because with so many non-linear transforms

2540 happening in the machine learning model structure, no one could easily tell what's going  
2541 on (Molnar 2019). There are some attempts for the model interpretability mentioned  
2542 above, such as Lime (Local Interpretable Model-Agnostic Explanations) (Ribeiro et al.  
2543 2016), that a simple/explainable model is used to approximate the underline model,  
2544 SHAP (Shaley Additive Explanations) (Lundberg and Lee 2017) that a game theory-  
2545 based approach is used to evaluate the feature level importance score, but these are all the  
2546 workaround solutions, and more researches need to be done for more accurately  
2547 describing the importance of the original variables.

## 2548 APPENDIX 1 PRINCIPAL COMPONENT ANALYSIS

### 2549 1.1 Principal component analysis

2550 Principal component analysis (PCA) is a dimension reduction approach to identify  
2551 new features (principal component, i.e., PC) that contain as much statistical information  
2552 of the original features as possible, and PCs that are not correlated to each other.  
2553 Statistical information is represented by the variance of the original features. The  
2554 correlation between different PCs equal to 0, and PCs are sorted by their variance.

2555 PCA is elaborated in math format as below:

2556 Assume there are  $m$  observations in the entire dataset w.r.t. each observation  
2557  $x^{(i)}$   $i = 1, \dots, m$ , where each observation  $x^{(i)}$  ( $i = 1, \dots, m$ ) is a multi-dimension  
2558 vector.

2559 The data  $X = [x^{(1)}, \dots, x^{(m)}]$  is centered through  $\hat{X} = X - [\frac{1}{m} \sum_{i=1}^m x_i, \dots, \frac{1}{m} \sum_{i=1}^m x_i]$ ,

2560 where  $\hat{X}$  represents that each column of  $X$  is subtracted by  $\frac{1}{m} \sum_{i=1}^m x_i$ . Assume  $A =$

2561  $[A_1, \dots, A_m]$  is the projection matrix that

$$2562 \quad Y = A\hat{X}$$

2563 Where  $A_i^T A_i = 1$  for  $i = 1, \dots, m$ ,  $Y = [Y_1, \dots, Y_m]$ , and  $Y_i Y_j = 0$  for  $i, j = 1, \dots, m; i \neq$   
2564  $j$ .

2565

2566 (the largest variance for each  $Y_i$  while  $Y_i Y_j = 0$  for  $i, j = 1, \dots, m; i \neq j$ ). Therefore, since  
2567 the variance of  $Y_i$  that  $Y_i Y_j = 0$  for  $i, j = 1, \dots, m; i \neq j$  should be maximized, hence the  
2568 trace of the covariance matrix of  $A\hat{X}$  should be maximized, which results in (1).

2569

$$\begin{aligned} A &= \underset{A}{\operatorname{argmax}} \operatorname{trace}(A^T S A) & (1) \\ \text{s. t.} \quad & A^T A = I \text{ and } S = \frac{1}{m} \hat{X} \hat{X}^T \end{aligned}$$

2570

2571 By using Lagrangian multiplier and taking the derivative on (1), we get

$$2572 \quad SA = \lambda A$$

2573 where  $\lambda$  is Lagrangian multiplier.

2574 Based on eigen-decomposition described in Stoer and Bulirsch (2013),  $\lambda$  is the  
2575 diagonal eigenvalue matrix of  $S$ , where eigenvalues of  $S$  are located in the diagonal of  $\lambda$ ,  
2576 and sorted decreasingly from left upper corner to the right lower corner.  $A$  is the

2577 corresponding eigenvector matrix of  $S$ , where  $A_i$  for  $i = 1, \dots, m$  is the eigenvector that  
2578 corresponds to  $i$ -th largest eigenvalue of  $S$ .  $A_i$  for  $i = 1, \dots, m$  represent the first  $m$  PCs.

2579 In applications, the first few PCs are chosen with the largest contribution to the total  
2580 variance (variation).

2581 More details can be found in Friedman et al. (2001).

## 2582 **1.2 Kernel PCA**

2583 Kernel PCA maps the original data  $X$  to a kernel Hilbert space through a  
2584 transformation  $\phi$  to perform the PCA, and the kernel space is unknown. Similarly,  
2585 assume there are  $m$  observation in the entire dataset w.r.t. each observation  $x^{(i)}$   $i =$   
2586  $1, \dots, m$ . Therefore, data  $X = [x^{(1)}, \dots, x^{(m)}]$ , and  $A = [A_1, \dots, A_m]$  is the projection  
2587 matrix that

$$2588 \quad Y = A\phi(X)$$

$$2589 \quad Y = [Y_1, \dots, Y_m], \text{ and } Y_i Y_j = 0 \text{ for } i, j = 1, \dots, m; i \neq j$$

2590 Similar to PCA, the problem can be rewritten as (2)

2591

2592

$$A = \underset{A}{\operatorname{argmax}} \operatorname{trace}(A^T C A) \quad (2)$$

$$s. t. \quad A^T A = I \text{ and } C = \frac{1}{m} \widehat{\phi(X)} \widehat{\phi(X)}^T$$

2593

2594 where  $\widehat{\phi(X)} = \phi(X) - [\frac{1}{m} \sum_{i=1}^m \phi(x_i), \dots, \frac{1}{m} \sum_{i=1}^m \phi(x_i)]$

2595 Equation (2) could be solved as same as Equation (1), hence convert to

2596

$$CA = \lambda A \quad (3)$$

2597 (3) can be transformed to vector form

2598

$$CA_k = \lambda_k A_k \quad (4)$$

2599 since  $C = \frac{1}{m} \widehat{\phi(X)} \widehat{\phi(X)}^T$ , (4) becomes



$$\frac{1}{m} \sum_{i=1}^m \widehat{\phi(x_i)} [\widehat{\phi(x_i)}^T A_k] = \lambda_k A_k \quad (5)$$

2600

2601 If both sides of the (5) is divided by  $\lambda_k$ ,  $A_k$  can be rewritten as

$$A_k = \sum_{i=1}^m t_{ki} \phi(x_i) \quad (6)$$

2602

$$\frac{1}{m} \sum_{i=1}^m \widehat{\phi(x_i)} \widehat{\phi(x_i)}^T \sum_{j=1}^m t_{kj} \phi(x_j) = \lambda_k \sum_{i=1}^m t_{ki} \phi(x_i) \quad (7)$$

2603 Define  $k(x_i, x_j) = \widehat{\phi(x_i)}^T \widehat{\phi(x_j)}$  for  $i, j = 1, \dots, m$ .  $K$  is a  $m * m$  dimensional matrix

2604 that  $K(i, j) = k(x_i, x_j)$  for  $i, j = 1, \dots, m$ .

$$\frac{1}{m} \sum_{i=1}^m k(x_i, x_i) \sum_{j=1}^m t_{kj} k(x_i, x_j) = \lambda_k \sum_{i=1}^m t_{ki} k(x_i, x_i) \quad (8)$$

2605 Then (8) becomes

$$K^2 t_k = \lambda_k m K t_k \quad (9)$$

2606

2607 where  $t_k = [t_{k1}, \dots, t_{km}]^T$

2608 (9) is divided by  $K$

$$K t_k = \lambda_k m t_k \quad (10)$$

2609  $Y_k$  is then calculated as:

$$Y_k = \widehat{\phi(X)}^T A_k = \sum_{i=1}^m t_{ki} k(x, x_i) \quad (11)$$

2610 (Schölkopf et al. 1997)

2611 More details can be found in Schölkopf et al. (1998).

2612

2613

2614

## APPENDIX 2 ADDITIONAL TABLES

2616 **Table A1: The highly correlated parameter group list.**

2617	[IRM1_16, IR00_10, IR00_13, IR00_15, IR00_16, IR00_2, IR00_4, IR00_6, IR00_7,
2618	IR00_8, IR00_9, IRM1_10, IRM1_12, IRM1_13, IRM1_15, IRM1_2, IRM1_4, IRM1_6,
2619	IRM1_7, IRM1_8, IRM1_9, IRM3_10, IRM3_12, IRM3_13, IRM3_15, IRM3_16,
2620	IRM3_2, IRM3_4, IRM3_6, IRM3_7, IRM3_8, IRM3_9]
2621	[E000, CSST, DSST, DSTA, ENEG, ENSS, EPOS, MTPW_10, MTPW_12, MTPW_14,
2622	MTPW_16, MTPW_20, MTPW_2, MTPW_4, MTPW_6, MTPW_8, PW03, PW05,
2623	PW07, PW09, PW11, PW13, PW15, PW17, PW21, RSST, T000, XDST, XNST,
2624	XTMX]
2625	[HIST_8, HIST_9, HIST_10, HIST_11, HIST_12, HIST_13, HIST_14, HIST_2,
2626	HIST_3, HIST_4, HIST_5, HIST_6, HIST_7]
2627	[PW14, MTPW_11, MTPW_13, MTPW_15, MTPW_3, MTPW_5, PW04, PW06,
2628	PW12, PW16]
2629	[CD26, CD20, COHC, ND26, XD20, XD22, XD24, XD26, XO20, XOHC]
2630	[MTPW_19, MTPW_0, MTPW_18, PW01, PW19, PW20], [HIST_15, HIST_16,
2631	HIST_17, HIST_18, HIST_19, HIST_20]
2632	[IRM3_19, IRM1_18, IRM1_19, IRM1_20, IRM3_18, IRM3_20]
2633	[PW08, MTPW_9, MTPW_17, PW10, PW18]
2634	[PSLV_4, PSLV_2, PSLV_6, U200, U20C]
2635	[IRM1_5, IR00_5, IRM1_3, IRM3_3, IRM3_5]
2636	[SHRD, SHDC, SHGC, SHRG]
2637	[V850, TWAC, V000, V500]
2638	[IRM3_11, IR00_11, IRM1_11]
2639	[RHMD, RHHI, RHLO]
2640	[IR00_20, IR00_18, IR00_19]
2641	[DTL, LON, TLON]
2642	[PSLV_3, PSLV_5, PSLV_7]
2643	[PENV, PENC, Z000]
2644	[VVAC, VMFX, VVAV]
2645	[BD12, BD06, BD18]
2646	[T250, T200]

2647	[NSST, NTMX]
2648	[HIST_1, HIST_0]
2649	[NOHC, RHCN]
2650	[RD26, RD20]
2651	[XD18, XD16]
2652	[D200, DIVC]
2653	[PC00, IR00_1]
2654	[OAGE, NAGE]
2655	[NTFR, XTFR]
2656	[PCM1, IRM1_1]
2657	[HE07, HE05]
2658	[PEFC, V20C]
2659	[PCM3, IRM3_1]
2660	[VMAX, MSLP]
2661	[TLAT, LAT]
2662	[MTPW_1, PW02]
2663	[O500, O700]
2664	[NDFR, XDFR]
2665	[IR00_12]
2666	[VMPI]
2667	[IR00_3]
2668	[ND20]
2669	[EPSS]
2670	[TWXC]
2671	[G150]
2672	[SHTD]
2673	[NDTX]
2674	[XDML]
2675	[Z850]
2676	[CFLX]
2677	[XDTX]

2678 [SHTS]  
 2679 [SDDC]  
 2680 [jd]  
 2681 [SHRS]  
 2682 [IR00\_14]  
 2683 [IR00\_17]  
 2684 [G250]  
 2685 [G200]  
 2686 [REFC]  
 2687 [PSLV\_1]  
 2688 [V300]  
 2689 [IRM1\_17]  
 2690 [T150]  
 2691 [TGRD]  
 2692 [TADV]  
 2693 [IRM3\_14]  
 2694 [R000]  
 2695 [IRM3\_17]  
 2696 [IRM1\_14]  
 2697  
 2698

2699 **Table A2: Detail of the Group with the original score from LLE data filter.**

2700 Group1, 0.07131729492630356, [NT18\_cc\_l11, NT18\_cc\_l12, NT18\_cc\_l13,  
 2701 NT18\_cc\_l14, NT18\_cc\_l15, NT18\_ciwc\_l12, NT18\_ciwc\_l13, NT18\_ciwc\_l14,  
 2702 NT18\_ciwc\_l15, NT18\_ciwc\_l16, NT18\_o3\_l1, NT18\_o3\_l10, NT18\_o3\_l11,  
 2703 NT18\_o3\_l14, NT18\_o3\_l15, NT18\_o3\_l16, NT18\_o3\_l17, NT18\_o3\_l12, NT18\_o3\_l19,  
 2704 NT18\_q\_l1, NT18\_q\_l2, NT18\_q\_l3, NT18\_q\_l37, NT18\_q\_l4, NT18\_q\_l5,  
 2705 NT18\_q\_l6, NT18\_r\_l11, NT18\_r\_l12, NT18\_r\_l14, NT18\_r\_l15, NT18\_r\_l16,  
 2706 NT18\_t\_l10, NT18\_t\_l14, NT18\_t\_l15, NT18\_t\_l16, NT18\_t\_l17, NT18\_t\_l18,  
 2707 NT18\_t\_l19, NT18\_t\_l20, NT18\_t\_l21, NT18\_t\_l22, NT18\_t\_l5, NT18\_t\_l6, NT18\_t\_l7,  
 2708 NT18\_t\_l8, NT18\_t\_l9, NT18\_u\_l11, NT18\_u\_l12, NT18\_u\_l13, NT18\_u\_l14,  
 2709 NT18\_u\_l15, NT18\_u\_l16, NT18\_u\_l25, NT18\_u\_l26, NT18\_u\_l27, NT18\_u\_l28,  
 2710 NT18\_u\_l29, NT18\_u\_l30, NT18\_u\_l31, NT18\_u\_l32, NT18\_u\_l33, NT18\_u\_l34,  
 2711 NT18\_u\_l35, NT18\_u\_l36, NT18\_u\_l37, NT12\_cc\_l11, NT12\_cc\_l12, NT12\_cc\_l13,  
 2712 NT12\_cc\_l14, NT12\_cc\_l15, NT12\_ciwc\_l12, NT12\_ciwc\_l13, NT12\_ciwc\_l14,  
 2713 NT12\_ciwc\_l15, NT12\_ciwc\_l16, NT12\_o3\_l1, NT12\_o3\_l10, NT12\_o3\_l11,

2714 NT12\_o3\_l14, NT12\_o3\_l15, NT12\_o3\_l16, NT12\_o3\_l17, NT12\_o3\_l2, NT12\_o3\_l5,  
 2715 NT12\_o3\_l9, NT12\_q\_l1, NT12\_q\_l13, NT12\_q\_l14, NT12\_q\_l2, NT12\_q\_l3,  
 2716 NT12\_q\_l37, NT12\_q\_l4, NT12\_q\_l5, NT12\_q\_l6, NT12\_r\_l11, NT12\_r\_l12,  
 2717 NT12\_r\_l13, NT12\_r\_l14, NT12\_r\_l15, NT12\_r\_l16, NT12\_t\_l10, NT12\_t\_l11,  
 2718 NT12\_t\_l14, NT12\_t\_l15, NT12\_t\_l16, NT12\_t\_l17, NT12\_t\_l18, NT12\_t\_l19,  
 2719 NT12\_t\_l20, NT12\_t\_l21, NT12\_t\_l22, NT12\_t\_l5, NT12\_t\_l6, NT12\_t\_l7, NT12\_t\_l8,  
 2720 NT12\_t\_l9, NT12\_u\_l11, NT12\_u\_l12, NT12\_u\_l13, NT12\_u\_l14, NT12\_u\_l15,  
 2721 NT12\_u\_l16, NT12\_u\_l25, NT12\_u\_l26, NT12\_u\_l27, NT12\_u\_l28, NT12\_u\_l29,  
 2722 NT12\_u\_l30, NT12\_u\_l31, NT12\_u\_l32, NT12\_u\_l33, NT12\_u\_l34, NT12\_u\_l35,  
 2723 NT12\_u\_l36, NT12\_u\_l37, NT12\_z\_l15, NT12\_z\_l16, NT12\_z\_l17, NT12\_z\_l18,  
 2724 NT06\_cc\_l11, NT06\_cc\_l12, NT06\_cc\_l13, NT06\_cc\_l14, NT06\_cc\_l15,  
 2725 NT06\_ciwc\_l12, NT06\_ciwc\_l13, NT06\_ciwc\_l14, NT06\_ciwc\_l15, NT06\_ciwc\_l16,  
 2726 NT06\_o3\_l1, NT06\_o3\_l10, NT06\_o3\_l11, NT06\_o3\_l12, NT06\_o3\_l13, NT06\_o3\_l14,  
 2727 NT06\_o3\_l15, NT06\_o3\_l16, NT06\_o3\_l17, NT06\_o3\_l2, NT06\_o3\_l5, NT06\_o3\_l9,  
 2728 NT06\_q\_l1, NT06\_q\_l13, NT06\_q\_l14, NT06\_q\_l15, NT06\_q\_l2, NT06\_q\_l3,  
 2729 NT06\_q\_l4, NT06\_q\_l5, NT06\_q\_l6, NT06\_r\_l11, NT06\_r\_l12, NT06\_r\_l13,  
 2730 NT06\_r\_l14, NT06\_r\_l15, NT06\_r\_l16, NT06\_t\_l10, NT06\_t\_l11, NT06\_t\_l14,  
 2731 NT06\_t\_l15, NT06\_t\_l16, NT06\_t\_l17, NT06\_t\_l18, NT06\_t\_l19, NT06\_t\_l20,  
 2732 NT06\_t\_l21, NT06\_t\_l22, NT06\_t\_l28, NT06\_t\_l29, NT06\_t\_l30, NT06\_t\_l31,  
 2733 NT06\_t\_l32, NT06\_t\_l5, NT06\_t\_l6, NT06\_t\_l7, NT06\_t\_l8, NT06\_t\_l9, NT06\_u\_l11,  
 2734 NT06\_u\_l12, NT06\_u\_l13, NT06\_u\_l14, NT06\_u\_l15, NT06\_u\_l16, NT06\_u\_l25,  
 2735 NT06\_u\_l26, NT06\_u\_l27, NT06\_u\_l28, NT06\_u\_l29, NT06\_u\_l30, NT06\_u\_l31,  
 2736 NT06\_u\_l32, NT06\_u\_l33, NT06\_u\_l34, NT06\_u\_l35, NT06\_u\_l36, NT06\_u\_l37,  
 2737 NT06\_z\_l14, NT06\_z\_l15, NT06\_z\_l16, NT06\_z\_l17, NT06\_z\_l18, NT06\_z\_l19,  
 2738 NT00\_cc\_l11, NT00\_cc\_l12, NT00\_cc\_l13, NT00\_cc\_l14, NT00\_cc\_l15,  
 2739 NT00\_ciwc\_l11, NT00\_ciwc\_l12, NT00\_ciwc\_l13, NT00\_ciwc\_l14, NT00\_ciwc\_l15,  
 2740 NT00\_ciwc\_l16, NT00\_o3\_l1, NT00\_o3\_l10, NT00\_o3\_l11, NT00\_o3\_l12,  
 2741 NT00\_o3\_l13, NT00\_o3\_l14, NT00\_o3\_l15, NT00\_o3\_l16, NT00\_o3\_l17, NT00\_o3\_l2,  
 2742 NT00\_o3\_l5, NT00\_o3\_l9, NT00\_q\_l1, NT00\_q\_l13, NT00\_q\_l14, NT00\_q\_l15,  
 2743 NT00\_q\_l2, NT00\_q\_l3, NT00\_q\_l4, NT00\_q\_l5, NT00\_q\_l6, NT00\_q\_l7, NT00\_r\_l11,  
 2744 NT00\_r\_l12, NT00\_r\_l13, NT00\_r\_l14, NT00\_r\_l15, NT00\_r\_l16, NT00\_t\_l11,  
 2745 NT00\_t\_l14, NT00\_t\_l15, NT00\_t\_l16, NT00\_t\_l17, NT00\_t\_l18, NT00\_t\_l19,  
 2746 NT00\_t\_l20, NT00\_t\_l21, NT00\_t\_l22, NT00\_t\_l27, NT00\_t\_l28, NT00\_t\_l29,  
 2747 NT00\_t\_l30, NT00\_t\_l31, NT00\_t\_l32, NT00\_t\_l33, NT00\_t\_l34, NT00\_t\_l5,  
 2748 NT00\_t\_l6, NT00\_t\_l7, NT00\_t\_l8, NT00\_t\_l9, NT00\_u\_l11, NT00\_u\_l12,  
 2749 NT00\_u\_l13, NT00\_u\_l14, NT00\_u\_l15, NT00\_u\_l26, NT00\_u\_l27, NT00\_u\_l28,  
 2750 NT00\_u\_l29, NT00\_u\_l30, NT00\_u\_l31, NT00\_u\_l32, NT00\_u\_l33, NT00\_u\_l34,  
 2751 NT00\_u\_l35, NT00\_u\_l36, NT00\_u\_l37, NT00\_z\_l13, NT00\_z\_l14, NT00\_z\_l15,  
 2752 NT00\_z\_l16, NT00\_z\_l17, NT00\_z\_l18, NT00\_z\_l19, NT18\_r\_l13]  
 2753 Group2, -0.040873502972860964, [NT18\_cc\_l16, NT18\_cc\_l17, NT18\_cc\_l18,  
 2754 NT18\_cc\_l19, NT18\_cc\_l20, NT18\_cc\_l21, NT18\_cc\_l22, NT18\_cc\_l23, NT18\_cc\_l27,  
 2755 NT18\_cc\_l28, NT18\_cc\_l29, NT18\_cc\_l30, NT18\_cc\_l31, NT18\_cc\_l32, NT18\_cc\_l33,  
 2756 NT18\_cc\_l34, NT18\_ciwc\_l18, NT18\_ciwc\_l19, NT18\_ciwc\_l20, NT18\_ciwc\_l21,  
 2757 NT18\_clwc\_l21, NT18\_clwc\_l22, NT18\_clwc\_l24, NT18\_clwc\_l25, NT18\_clwc\_l26,

2758 NT18\_clwc\_l27, NT18\_clwc\_l28, NT18\_clwc\_l29, NT18\_clwc\_l30, NT18\_clwc\_l31,  
 2759 NT18\_q\_l15, NT18\_q\_l16, NT18\_q\_l17, NT18\_q\_l18, NT18\_q\_l19, NT18\_q\_l20,  
 2760 NT18\_q\_l21, NT18\_q\_l22, NT18\_q\_l23, NT18\_q\_l28, NT18\_q\_l29, NT18\_q\_l30,  
 2761 NT18\_q\_l31, NT18\_q\_l32, NT18\_q\_l33, NT18\_q\_l34, NT18\_r\_l17, NT18\_r\_l18,  
 2762 NT18\_r\_l19, NT18\_r\_l20, NT18\_r\_l21, NT18\_r\_l22, NT18\_r\_l23, NT18\_r\_l24,  
 2763 NT18\_r\_l25, NT18\_r\_l26, NT18\_r\_l27, NT18\_r\_l28, NT18\_r\_l29, NT18\_r\_l30,  
 2764 NT18\_r\_l31, NT18\_r\_l32, NT18\_r\_l33, NT18\_r\_l34, NT18\_r\_l35, NT18\_r\_l36,  
 2765 NT12\_cc\_l16, NT12\_cc\_l17, NT12\_cc\_l18, NT12\_cc\_l19, NT12\_cc\_l20, NT12\_cc\_l21,  
 2766 NT12\_cc\_l22, NT12\_cc\_l23, NT12\_cc\_l27, NT12\_cc\_l28, NT12\_cc\_l29, NT12\_cc\_l30,  
 2767 NT12\_cc\_l31, NT12\_cc\_l32, NT12\_cc\_l33, NT12\_ciwc\_l17, NT12\_ciwc\_l18,  
 2768 NT12\_ciwc\_l19, NT12\_ciwc\_l20, NT12\_ciwc\_l21, NT12\_ciwc\_l22, NT12\_clwc\_l19,  
 2769 NT12\_clwc\_l20, NT12\_clwc\_l21, NT12\_clwc\_l22, NT12\_clwc\_l23, NT12\_clwc\_l24,  
 2770 NT12\_clwc\_l25, NT12\_clwc\_l26, NT12\_clwc\_l27, NT12\_clwc\_l28, NT12\_clwc\_l29,  
 2771 NT12\_clwc\_l30, NT12\_clwc\_l31, NT12\_q\_l16, NT12\_q\_l17, NT12\_q\_l18,  
 2772 NT12\_q\_l19, NT12\_q\_l20, NT12\_q\_l21, NT12\_q\_l22, NT12\_q\_l23, NT12\_q\_l27,  
 2773 NT12\_q\_l28, NT12\_q\_l29, NT12\_q\_l30, NT12\_q\_l31, NT12\_q\_l32, NT12\_q\_l33,  
 2774 NT12\_r\_l17, NT12\_r\_l18, NT12\_r\_l19, NT12\_r\_l20, NT12\_r\_l22, NT12\_r\_l23,  
 2775 NT12\_r\_l24, NT12\_r\_l25, NT12\_r\_l26, NT12\_r\_l27, NT12\_r\_l28, NT12\_r\_l29,  
 2776 NT12\_r\_l30, NT12\_r\_l31, NT12\_r\_l32, NT12\_r\_l33, NT12\_r\_l34, NT12\_r\_l35,  
 2777 NT12\_r\_l36, NT12\_r\_l37, NT06\_cc\_l16, NT06\_cc\_l17, NT06\_cc\_l18, NT06\_cc\_l19,  
 2778 NT06\_cc\_l20, NT06\_cc\_l21, NT06\_cc\_l22, NT06\_cc\_l23, NT06\_cc\_l28, NT06\_cc\_l29,  
 2779 NT06\_cc\_l30, NT06\_cc\_l31, NT06\_cc\_l32, NT06\_cc\_l33, NT06\_ciwc\_l18,  
 2780 NT06\_ciwc\_l19, NT06\_ciwc\_l20, NT06\_ciwc\_l21, NT06\_ciwc\_l22, NT06\_clwc\_l20,  
 2781 NT06\_clwc\_l21, NT06\_clwc\_l22, NT06\_clwc\_l23, NT06\_clwc\_l24, NT06\_clwc\_l25,  
 2782 NT06\_clwc\_l26, NT06\_clwc\_l27, NT06\_clwc\_l28, NT06\_clwc\_l29, NT06\_clwc\_l30,  
 2783 NT06\_clwc\_l31, NT06\_q\_l16, NT06\_q\_l17, NT06\_q\_l18, NT06\_q\_l19, NT06\_q\_l20,  
 2784 NT06\_q\_l21, NT06\_q\_l22, NT06\_q\_l23, NT06\_q\_l24, NT06\_q\_l25, NT06\_q\_l26,  
 2785 NT06\_q\_l27, NT06\_q\_l28, NT06\_q\_l29, NT06\_q\_l30, NT06\_q\_l31, NT06\_q\_l32,  
 2786 NT06\_r\_l17, NT06\_r\_l18, NT06\_r\_l19, NT06\_r\_l20, NT06\_r\_l21, NT06\_r\_l22,  
 2787 NT06\_r\_l23, NT06\_r\_l24, NT06\_r\_l25, NT06\_r\_l26, NT06\_r\_l27, NT06\_r\_l28,  
 2788 NT06\_r\_l29, NT06\_r\_l30, NT06\_r\_l31, NT06\_r\_l32, NT06\_r\_l33, NT06\_r\_l34,  
 2789 NT06\_r\_l35, NT06\_r\_l36, NT00\_cc\_l16, NT00\_cc\_l17, NT00\_cc\_l18, NT00\_cc\_l19,  
 2790 NT00\_cc\_l20, NT00\_cc\_l21, NT00\_cc\_l22, NT00\_cc\_l31, NT00\_cc\_l32,  
 2791 NT00\_ciwc\_l18, NT00\_ciwc\_l19, NT00\_ciwc\_l20, NT00\_ciwc\_l21, NT00\_clwc\_l21,  
 2792 NT00\_clwc\_l22, NT00\_clwc\_l23, NT00\_clwc\_l24, NT00\_clwc\_l25, NT00\_clwc\_l26,  
 2793 NT00\_clwc\_l27, NT00\_clwc\_l28, NT00\_q\_l17, NT00\_q\_l18, NT00\_q\_l19,  
 2794 NT00\_q\_l20, NT00\_q\_l21, NT00\_q\_l22, NT00\_q\_l23, NT00\_q\_l24, NT00\_q\_l25,  
 2795 NT00\_q\_l26, NT00\_q\_l27, NT00\_r\_l17, NT00\_r\_l18, NT00\_r\_l19, NT00\_r\_l20,  
 2796 NT00\_r\_l21, NT00\_r\_l22, NT00\_r\_l23, NT00\_r\_l24, NT00\_r\_l25, NT00\_r\_l26,  
 2797 NT00\_r\_l27, NT00\_r\_l28, NT00\_r\_l29, NT00\_r\_l30, NT00\_r\_l31, NT00\_r\_l32,  
 2798 NT00\_r\_l33, NT00\_r\_l34, NT00\_r\_l35, NT12\_r\_l21]  
 2799 Group3, 0.06259891129094408, [NT18\_w\_l1, NT18\_w\_l10, NT18\_w\_l11,  
 2800 NT18\_w\_l12, NT18\_w\_l13, NT18\_w\_l14, NT18\_w\_l15, NT18\_w\_l16, NT18\_w\_l17,  
 2801 NT18\_w\_l18, NT18\_w\_l19, NT18\_w\_l2, NT18\_w\_l20, NT18\_w\_l21, NT18\_w\_l22,

2802 NT18\_w\_l23, NT18\_w\_l24, NT18\_w\_l25, NT18\_w\_l26, NT18\_w\_l27, NT18\_w\_l28,  
 2803 NT18\_w\_l29, NT18\_w\_l3, NT18\_w\_l30, NT18\_w\_l31, NT18\_w\_l32, NT18\_w\_l33,  
 2804 NT18\_w\_l34, NT18\_w\_l35, NT18\_w\_l36, NT18\_w\_l37, NT18\_w\_l4, NT18\_w\_l5,  
 2805 NT18\_w\_l6, NT18\_w\_l7, NT18\_w\_l8, NT18\_w\_l9, NT12\_w\_l1, NT12\_w\_l10,  
 2806 NT12\_w\_l11, NT12\_w\_l12, NT12\_w\_l13, NT12\_w\_l14, NT12\_w\_l15, NT12\_w\_l16,  
 2807 NT12\_w\_l17, NT12\_w\_l18, NT12\_w\_l19, NT12\_w\_l2, NT12\_w\_l20, NT12\_w\_l21,  
 2808 NT12\_w\_l22, NT12\_w\_l23, NT12\_w\_l24, NT12\_w\_l25, NT12\_w\_l26, NT12\_w\_l27,  
 2809 NT12\_w\_l28, NT12\_w\_l29, NT12\_w\_l3, NT12\_w\_l30, NT12\_w\_l31, NT12\_w\_l32,  
 2810 NT12\_w\_l33, NT12\_w\_l34, NT12\_w\_l35, NT12\_w\_l36, NT12\_w\_l37, NT12\_w\_l4,  
 2811 NT12\_w\_l5, NT12\_w\_l6, NT12\_w\_l7, NT12\_w\_l8, NT12\_w\_l9, NT06\_w\_l1,  
 2812 NT06\_w\_l10, NT06\_w\_l11, NT06\_w\_l12, NT06\_w\_l13, NT06\_w\_l14, NT06\_w\_l15,  
 2813 NT06\_w\_l16, NT06\_w\_l17, NT06\_w\_l18, NT06\_w\_l19, NT06\_w\_l2, NT06\_w\_l20,  
 2814 NT06\_w\_l21, NT06\_w\_l22, NT06\_w\_l23, NT06\_w\_l24, NT06\_w\_l25, NT06\_w\_l26,  
 2815 NT06\_w\_l27, NT06\_w\_l28, NT06\_w\_l29, NT06\_w\_l3, NT06\_w\_l30, NT06\_w\_l31,  
 2816 NT06\_w\_l32, NT06\_w\_l33, NT06\_w\_l34, NT06\_w\_l35, NT06\_w\_l36, NT06\_w\_l37,  
 2817 NT06\_w\_l4, NT06\_w\_l5, NT06\_w\_l6, NT06\_w\_l7, NT06\_w\_l8, NT06\_w\_l9,  
 2818 NT00\_w\_l1, NT00\_w\_l10, NT00\_w\_l11, NT00\_w\_l12, NT00\_w\_l13, NT00\_w\_l14,  
 2819 NT00\_w\_l15, NT00\_w\_l16, NT00\_w\_l17, NT00\_w\_l18, NT00\_w\_l19, NT00\_w\_l2,  
 2820 NT00\_w\_l20, NT00\_w\_l21, NT00\_w\_l22, NT00\_w\_l23, NT00\_w\_l24, NT00\_w\_l25,  
 2821 NT00\_w\_l26, NT00\_w\_l27, NT00\_w\_l28, NT00\_w\_l29, NT00\_w\_l3, NT00\_w\_l30,  
 2822 NT00\_w\_l31, NT00\_w\_l32, NT00\_w\_l33, NT00\_w\_l34, NT00\_w\_l35, NT00\_w\_l36,  
 2823 NT00\_w\_l37, NT00\_w\_l4, NT00\_w\_l5, NT00\_w\_l6, NT00\_w\_l8, NT00\_w\_l9,  
 2824 NT00\_w\_l7]  
 2825 Group4, 0.02269292645906973, [NT18\_pv\_l10, NT18\_pv\_l11, NT18\_pv\_l12,  
 2826 NT18\_pv\_l13, NT18\_pv\_l14, NT18\_pv\_l15, NT18\_pv\_l16, NT18\_pv\_l17,  
 2827 NT18\_pv\_l18, NT18\_pv\_l19, NT18\_pv\_l20, NT18\_pv\_l21, NT18\_pv\_l22,  
 2828 NT18\_pv\_l23, NT18\_pv\_l24, NT18\_pv\_l25, NT18\_pv\_l26, NT18\_pv\_l27,  
 2829 NT18\_pv\_l28, NT18\_pv\_l29, NT18\_pv\_l3, NT18\_pv\_l30, NT18\_pv\_l31, NT18\_pv\_l32,  
 2830 NT18\_pv\_l33, NT18\_pv\_l34, NT18\_pv\_l35, NT18\_pv\_l36, NT18\_pv\_l37, NT18\_pv\_l4,  
 2831 NT18\_pv\_l6, NT18\_pv\_l7, NT18\_pv\_l8, NT18\_pv\_l9, NT12\_pv\_l10, NT12\_pv\_l11,  
 2832 NT12\_pv\_l12, NT12\_pv\_l13, NT12\_pv\_l14, NT12\_pv\_l15, NT12\_pv\_l16,  
 2833 NT12\_pv\_l17, NT12\_pv\_l18, NT12\_pv\_l19, NT12\_pv\_l20, NT12\_pv\_l21,  
 2834 NT12\_pv\_l22, NT12\_pv\_l23, NT12\_pv\_l24, NT12\_pv\_l25, NT12\_pv\_l26,  
 2835 NT12\_pv\_l27, NT12\_pv\_l28, NT12\_pv\_l29, NT12\_pv\_l3, NT12\_pv\_l30, NT12\_pv\_l31,  
 2836 NT12\_pv\_l32, NT12\_pv\_l33, NT12\_pv\_l34, NT12\_pv\_l35, NT12\_pv\_l36,  
 2837 NT12\_pv\_l37, NT12\_pv\_l4, NT12\_pv\_l5, NT12\_pv\_l6, NT12\_pv\_l7, NT12\_pv\_l8,  
 2838 NT12\_pv\_l9, NT06\_pv\_l10, NT06\_pv\_l11, NT06\_pv\_l12, NT06\_pv\_l13, NT06\_pv\_l14,  
 2839 NT06\_pv\_l15, NT06\_pv\_l16, NT06\_pv\_l17, NT06\_pv\_l18, NT06\_pv\_l19,  
 2840 NT06\_pv\_l20, NT06\_pv\_l21, NT06\_pv\_l22, NT06\_pv\_l23, NT06\_pv\_l24,  
 2841 NT06\_pv\_l25, NT06\_pv\_l26, NT06\_pv\_l27, NT06\_pv\_l28, NT06\_pv\_l29, NT06\_pv\_l3,  
 2842 NT06\_pv\_l30, NT06\_pv\_l31, NT06\_pv\_l32, NT06\_pv\_l33, NT06\_pv\_l34,  
 2843 NT06\_pv\_l35, NT06\_pv\_l36, NT06\_pv\_l37, NT06\_pv\_l4, NT06\_pv\_l5, NT06\_pv\_l6,  
 2844 NT06\_pv\_l7, NT06\_pv\_l8, NT06\_pv\_l9, NT00\_pv\_l10, NT00\_pv\_l11, NT00\_pv\_l12,  
 2845 NT00\_pv\_l13, NT00\_pv\_l14, NT00\_pv\_l15, NT00\_pv\_l16, NT00\_pv\_l17,



2846 NT00\_pv\_118, NT00\_pv\_119, NT00\_pv\_120, NT00\_pv\_121, NT00\_pv\_122,  
 2847 NT00\_pv\_123, NT00\_pv\_124, NT00\_pv\_125, NT00\_pv\_126, NT00\_pv\_127,  
 2848 NT00\_pv\_128, NT00\_pv\_129, NT00\_pv\_13, NT00\_pv\_130, NT00\_pv\_131, NT00\_pv\_132,  
 2849 NT00\_pv\_133, NT00\_pv\_134, NT00\_pv\_135, NT00\_pv\_136, NT00\_pv\_137, NT00\_pv\_14,  
 2850 NT00\_pv\_15, NT00\_pv\_16, NT00\_pv\_17, NT00\_pv\_18, NT00\_pv\_19, NT18\_pv\_15]  
 2851 Group5, -0.042245901124683405, [NT18\_d\_11, NT18\_d\_110, NT18\_d\_111,  
 2852 NT18\_d\_112, NT18\_d\_113, NT18\_d\_114, NT18\_d\_115, NT18\_d\_116, NT18\_d\_117,  
 2853 NT18\_d\_118, NT18\_d\_119, NT18\_d\_12, NT18\_d\_120, NT18\_d\_121, NT18\_d\_122,  
 2854 NT18\_d\_123, NT18\_d\_124, NT18\_d\_125, NT18\_d\_126, NT18\_d\_128, NT18\_d\_129,  
 2855 NT18\_d\_13, NT18\_d\_130, NT18\_d\_131, NT18\_d\_14, NT18\_d\_15, NT18\_d\_16,  
 2856 NT18\_d\_17, NT18\_d\_18, NT18\_d\_19, NT12\_d\_11, NT12\_d\_110, NT12\_d\_111,  
 2857 NT12\_d\_112, NT12\_d\_113, NT12\_d\_114, NT12\_d\_115, NT12\_d\_116, NT12\_d\_117,  
 2858 NT12\_d\_118, NT12\_d\_119, NT12\_d\_12, NT12\_d\_120, NT12\_d\_121, NT12\_d\_122,  
 2859 NT12\_d\_123, NT12\_d\_124, NT12\_d\_125, NT12\_d\_126, NT12\_d\_127, NT12\_d\_128,  
 2860 NT12\_d\_129, NT12\_d\_13, NT12\_d\_130, NT12\_d\_131, NT12\_d\_132, NT12\_d\_14,  
 2861 NT12\_d\_15, NT12\_d\_16, NT12\_d\_17, NT12\_d\_18, NT12\_d\_19, NT06\_d\_11,  
 2862 NT06\_d\_110, NT06\_d\_111, NT06\_d\_112, NT06\_d\_113, NT06\_d\_114, NT06\_d\_115,  
 2863 NT06\_d\_116, NT06\_d\_117, NT06\_d\_118, NT06\_d\_119, NT06\_d\_12, NT06\_d\_120,  
 2864 NT06\_d\_121, NT06\_d\_122, NT06\_d\_123, NT06\_d\_124, NT06\_d\_125, NT06\_d\_126,  
 2865 NT06\_d\_127, NT06\_d\_128, NT06\_d\_129, NT06\_d\_13, NT06\_d\_130, NT06\_d\_131,  
 2866 NT06\_d\_132, NT06\_d\_14, NT06\_d\_15, NT06\_d\_16, NT06\_d\_17, NT06\_d\_18,  
 2867 NT06\_d\_19, NT00\_d\_11, NT00\_d\_110, NT00\_d\_111, NT00\_d\_112, NT00\_d\_113,  
 2868 NT00\_d\_114, NT00\_d\_115, NT00\_d\_116, NT00\_d\_117, NT00\_d\_118, NT00\_d\_119,  
 2869 NT00\_d\_12, NT00\_d\_120, NT00\_d\_121, NT00\_d\_122, NT00\_d\_123, NT00\_d\_124,  
 2870 NT00\_d\_125, NT00\_d\_126, NT00\_d\_127, NT00\_d\_128, NT00\_d\_129, NT00\_d\_13,  
 2871 NT00\_d\_130, NT00\_d\_131, NT00\_d\_132, NT00\_d\_14, NT00\_d\_15, NT00\_d\_16,  
 2872 NT00\_d\_17, NT00\_d\_18, NT00\_d\_19, NT18\_d\_127]  
 2873 Group6, -0.02193022436520109, [NT18\_vo\_11, NT18\_vo\_110, NT18\_vo\_12,  
 2874 NT18\_vo\_120, NT18\_vo\_121, NT18\_vo\_122, NT18\_vo\_123, NT18\_vo\_124,  
 2875 NT18\_vo\_125, NT18\_vo\_126, NT18\_vo\_127, NT18\_vo\_128, NT18\_vo\_129, NT18\_vo\_13,  
 2876 NT18\_vo\_130, NT18\_vo\_131, NT18\_vo\_132, NT18\_vo\_133, NT18\_vo\_134,  
 2877 NT18\_vo\_135, NT18\_vo\_136, NT18\_vo\_137, NT18\_vo\_14, NT18\_vo\_15, NT18\_vo\_16,  
 2878 NT18\_vo\_17, NT18\_vo\_18, NT18\_vo\_19, NT12\_vo\_11, NT12\_vo\_110, NT12\_vo\_12,  
 2879 NT12\_vo\_121, NT12\_vo\_122, NT12\_vo\_123, NT12\_vo\_124, NT12\_vo\_125,  
 2880 NT12\_vo\_126, NT12\_vo\_127, NT12\_vo\_128, NT12\_vo\_129, NT12\_vo\_13, NT12\_vo\_130,  
 2881 NT12\_vo\_131, NT12\_vo\_132, NT12\_vo\_133, NT12\_vo\_134, NT12\_vo\_135,  
 2882 NT12\_vo\_136, NT12\_vo\_137, NT12\_vo\_14, NT12\_vo\_15, NT12\_vo\_16, NT12\_vo\_17,  
 2883 NT12\_vo\_18, NT12\_vo\_19, NT06\_vo\_11, NT06\_vo\_110, NT06\_vo\_12, NT06\_vo\_121,  
 2884 NT06\_vo\_122, NT06\_vo\_123, NT06\_vo\_124, NT06\_vo\_125, NT06\_vo\_126,  
 2885 NT06\_vo\_127, NT06\_vo\_128, NT06\_vo\_129, NT06\_vo\_13, NT06\_vo\_130, NT06\_vo\_131,  
 2886 NT06\_vo\_132, NT06\_vo\_133, NT06\_vo\_134, NT06\_vo\_135, NT06\_vo\_136,  
 2887 NT06\_vo\_137, NT06\_vo\_14, NT06\_vo\_15, NT06\_vo\_16, NT06\_vo\_17, NT06\_vo\_19,  
 2888 NT00\_vo\_11, NT00\_vo\_110, NT00\_vo\_12, NT00\_vo\_121, NT00\_vo\_122, NT00\_vo\_123,  
 2889 NT00\_vo\_124, NT00\_vo\_125, NT00\_vo\_126, NT00\_vo\_127, NT00\_vo\_128,

2890 NT00\_vo\_l29, NT00\_vo\_l3, NT00\_vo\_l30, NT00\_vo\_l31, NT00\_vo\_l32, NT00\_vo\_l33,  
 2891 NT00\_vo\_l34, NT00\_vo\_l35, NT00\_vo\_l36, NT00\_vo\_l37, NT00\_vo\_l4, NT00\_vo\_l5,  
 2892 NT00\_vo\_l6, NT00\_vo\_l7, NT00\_vo\_l8, NT00\_vo\_l9, NT06\_vo\_l8]  
 2893 Group7, -0.053063711383979584, [NT18\_clwc\_l32, NT18\_clwc\_l33, NT18\_q\_l35,  
 2894 NT18\_q\_l36, NT18\_r\_l37, NT18\_t\_l1, NT18\_t\_l12, NT18\_t\_l13, NT18\_t\_l2,  
 2895 NT18\_t\_l23, NT18\_t\_l24, NT18\_t\_l25, NT18\_t\_l26, NT18\_t\_l27, NT18\_t\_l28,  
 2896 NT18\_t\_l29, NT18\_t\_l3, NT18\_t\_l30, NT18\_t\_l31, NT18\_t\_l32, NT18\_t\_l33,  
 2897 NT18\_t\_l34, NT18\_t\_l35, NT18\_t\_l36, NT18\_t\_l37, NT18\_t\_l4, NT12\_clwc\_l32,  
 2898 NT12\_clwc\_l33, NT12\_q\_l15, NT12\_q\_l35, NT12\_q\_l36, NT12\_t\_l1, NT12\_t\_l13,  
 2899 NT12\_t\_l2, NT12\_t\_l23, NT12\_t\_l24, NT12\_t\_l25, NT12\_t\_l26, NT12\_t\_l27,  
 2900 NT12\_t\_l28, NT12\_t\_l29, NT12\_t\_l3, NT12\_t\_l30, NT12\_t\_l31, NT12\_t\_l32,  
 2901 NT12\_t\_l33, NT12\_t\_l34, NT12\_t\_l35, NT12\_t\_l36, NT12\_t\_l37, NT12\_t\_l4,  
 2902 NT06\_clwc\_l32, NT06\_clwc\_l33, NT06\_r\_l37, NT06\_t\_l1, NT06\_t\_l13, NT06\_t\_l2,  
 2903 NT06\_t\_l23, NT06\_t\_l24, NT06\_t\_l25, NT06\_t\_l26, NT06\_t\_l27, NT06\_t\_l3,  
 2904 NT06\_t\_l33, NT06\_t\_l34, NT06\_t\_l35, NT06\_t\_l36, NT06\_t\_l37, NT06\_t\_l4,  
 2905 NT00\_clwc\_l32, NT00\_clwc\_l33, NT00\_q\_l16, NT00\_r\_l36, NT00\_r\_l37, NT00\_t\_l1,  
 2906 NT00\_t\_l13, NT00\_t\_l2, NT00\_t\_l23, NT00\_t\_l24, NT00\_t\_l26, NT00\_t\_l3,  
 2907 NT00\_t\_l35, NT00\_t\_l36, NT00\_t\_l37, NT00\_t\_l4, NT00\_t\_l25]  
 2908 Group8, 0.0, [NT18\_clwc\_l1, NT18\_clwc\_l10, NT18\_clwc\_l11, NT18\_clwc\_l12,  
 2909 NT18\_clwc\_l13, NT18\_clwc\_l14, NT18\_clwc\_l15, NT18\_clwc\_l16, NT18\_clwc\_l17,  
 2910 NT18\_clwc\_l2, NT18\_clwc\_l3, NT18\_clwc\_l4, NT18\_clwc\_l5, NT18\_clwc\_l6,  
 2911 NT18\_clwc\_l7, NT18\_clwc\_l8, NT18\_clwc\_l9, NT12\_clwc\_l1, NT12\_clwc\_l10,  
 2912 NT12\_clwc\_l11, NT12\_clwc\_l12, NT12\_clwc\_l13, NT12\_clwc\_l14, NT12\_clwc\_l15,  
 2913 NT12\_clwc\_l16, NT12\_clwc\_l17, NT12\_clwc\_l2, NT12\_clwc\_l3, NT12\_clwc\_l4,  
 2914 NT12\_clwc\_l5, NT12\_clwc\_l6, NT12\_clwc\_l7, NT12\_clwc\_l8, NT12\_clwc\_l9,  
 2915 NT06\_clwc\_l1, NT06\_clwc\_l10, NT06\_clwc\_l11, NT06\_clwc\_l12, NT06\_clwc\_l13,  
 2916 NT06\_clwc\_l14, NT06\_clwc\_l15, NT06\_clwc\_l16, NT06\_clwc\_l17, NT06\_clwc\_l2,  
 2917 NT06\_clwc\_l3, NT06\_clwc\_l4, NT06\_clwc\_l5, NT06\_clwc\_l6, NT06\_clwc\_l7,  
 2918 NT06\_clwc\_l8, NT06\_clwc\_l9, NT00\_clwc\_l1, NT00\_clwc\_l10, NT00\_clwc\_l11,  
 2919 NT00\_clwc\_l12, NT00\_clwc\_l13, NT00\_clwc\_l14, NT00\_clwc\_l15, NT00\_clwc\_l16,  
 2920 NT00\_clwc\_l17, NT00\_clwc\_l3, NT00\_clwc\_l4, NT00\_clwc\_l5, NT00\_clwc\_l6,  
 2921 NT00\_clwc\_l7, NT00\_clwc\_l8, NT00\_clwc\_l9, NT00\_clwc\_l2]  
 2922 Group9, -0.026883331336958416, [NT18\_ciwc\_l24, NT18\_o3\_l12, NT18\_o3\_l13,  
 2923 NT18\_o3\_l15, NT18\_r\_l10, NT18\_t\_l11, NT18\_z\_l11, NT18\_z\_l12, NT18\_z\_l13,  
 2924 NT18\_z\_l14, NT18\_z\_l15, NT18\_z\_l16, NT18\_z\_l18, NT18\_z\_l19, NT18\_z\_l20,  
 2925 NT18\_z\_l21, NT18\_z\_l22, NT18\_z\_l23, NT18\_z\_l24, NT18\_z\_l25, NT12\_o3\_l12,  
 2926 NT12\_o3\_l13, NT12\_r\_l10, NT12\_z\_l11, NT12\_z\_l12, NT12\_z\_l13, NT12\_z\_l14,  
 2927 NT12\_z\_l19, NT12\_z\_l20, NT12\_z\_l21, NT12\_z\_l22, NT12\_z\_l23, NT12\_z\_l24,  
 2928 NT12\_z\_l25, NT06\_r\_l10, NT06\_u\_l24, NT06\_z\_l10, NT06\_z\_l11, NT06\_z\_l12,  
 2929 NT06\_z\_l13, NT06\_z\_l20, NT06\_z\_l21, NT06\_z\_l22, NT06\_z\_l23, NT06\_z\_l24,  
 2930 NT06\_z\_l25, NT00\_r\_l10, NT00\_u\_l24, NT00\_u\_l25, NT00\_z\_l10, NT00\_z\_l11,  
 2931 NT00\_z\_l12, NT00\_z\_l20, NT00\_z\_l21, NT00\_z\_l22, NT00\_z\_l23, NT00\_z\_l24,  
 2932 NT00\_z\_l25, NT18\_z\_l17]

2933 Group10, 0.006458219563131418, [NT18\_o3\_l25, NT18\_o3\_l26, NT18\_o3\_l27,  
 2934 NT18\_o3\_l28, NT18\_o3\_l29, NT18\_o3\_l30, NT18\_o3\_l31, NT18\_o3\_l32,  
 2935 NT18\_o3\_l33, NT18\_o3\_l34, NT18\_o3\_l35, NT12\_o3\_l24, NT12\_o3\_l25,  
 2936 NT12\_o3\_l26, NT12\_o3\_l27, NT12\_o3\_l28, NT12\_o3\_l29, NT12\_o3\_l30,  
 2937 NT12\_o3\_l31, NT12\_o3\_l32, NT12\_o3\_l33, NT12\_o3\_l34, NT12\_o3\_l35, NT12\_z\_l37,  
 2938 NT06\_o3\_l24, NT06\_o3\_l25, NT06\_o3\_l26, NT06\_o3\_l27, NT06\_o3\_l28,  
 2939 NT06\_o3\_l29, NT06\_o3\_l30, NT06\_o3\_l31, NT06\_o3\_l33, NT06\_o3\_l34,  
 2940 NT06\_o3\_l35, NT06\_z\_l37, NT00\_o3\_l24, NT00\_o3\_l25, NT00\_o3\_l26, NT00\_o3\_l27,  
 2941 NT00\_o3\_l28, NT00\_o3\_l29, NT00\_o3\_l30, NT00\_o3\_l31, NT00\_o3\_l32,  
 2942 NT00\_o3\_l33, NT00\_o3\_l34, NT00\_o3\_l35, NT06\_o3\_l32]  
 2943 Group11, -0.04429699024103351, [NT18\_v\_l27, NT18\_v\_l28, NT18\_v\_l29,  
 2944 NT18\_v\_l30, NT18\_v\_l31, NT18\_v\_l32, NT18\_v\_l33, NT18\_v\_l34, NT18\_v\_l35,  
 2945 NT18\_v\_l36, NT18\_v\_l37, NT12\_v\_l27, NT12\_v\_l28, NT12\_v\_l29, NT12\_v\_l30,  
 2946 NT12\_v\_l31, NT12\_v\_l32, NT12\_v\_l33, NT12\_v\_l34, NT12\_v\_l35, NT12\_v\_l36,  
 2947 NT12\_v\_l37, NT06\_v\_l27, NT06\_v\_l28, NT06\_v\_l29, NT06\_v\_l30, NT06\_v\_l32,  
 2948 NT06\_v\_l33, NT06\_v\_l34, NT06\_v\_l35, NT06\_v\_l36, NT06\_v\_l37, NT00\_v\_l27,  
 2949 NT00\_v\_l28, NT00\_v\_l29, NT00\_v\_l30, NT00\_v\_l31, NT00\_v\_l32, NT00\_v\_l33,  
 2950 NT00\_v\_l34, NT00\_v\_l35, NT00\_v\_l36, NT00\_v\_l37, NT06\_v\_l31]  
 2951 Group12, 0.016470284903957744, [NT18\_vo\_l11, NT18\_vo\_l12, NT18\_vo\_l13,  
 2952 NT18\_vo\_l14, NT18\_vo\_l15, NT18\_vo\_l16, NT18\_vo\_l17, NT18\_vo\_l18,  
 2953 NT18\_vo\_l19, NT12\_vo\_l11, NT12\_vo\_l12, NT12\_vo\_l13, NT12\_vo\_l14,  
 2954 NT12\_vo\_l15, NT12\_vo\_l16, NT12\_vo\_l17, NT12\_vo\_l18, NT12\_vo\_l19,  
 2955 NT12\_vo\_l20, NT06\_vo\_l12, NT06\_vo\_l13, NT06\_vo\_l14, NT06\_vo\_l15,  
 2956 NT06\_vo\_l16, NT06\_vo\_l17, NT06\_vo\_l18, NT06\_vo\_l19, NT06\_vo\_l20,  
 2957 NT00\_vo\_l11, NT00\_vo\_l12, NT00\_vo\_l13, NT00\_vo\_l14, NT00\_vo\_l15,  
 2958 NT00\_vo\_l16, NT00\_vo\_l17, NT00\_vo\_l18, NT00\_vo\_l19, NT00\_vo\_l20,  
 2959 NT06\_vo\_l11]  
 2960 Group13, 0.023626937292063666, [NT18\_cc\_l25, NT18\_cc\_l26, NT12\_cc\_l24,  
 2961 NT12\_cc\_l25, NT12\_cc\_l26, NT12\_cc\_l34, NT12\_cc\_l35, NT12\_cc\_l36, NT12\_cc\_l37,  
 2962 NT06\_cc\_l24, NT06\_cc\_l25, NT06\_cc\_l26, NT06\_cc\_l34, NT06\_cc\_l35, NT06\_cc\_l36,  
 2963 NT06\_cc\_l37, NT06\_clwc\_l34, NT06\_clwc\_l35, NT06\_clwc\_l36, NT06\_clwc\_l37,  
 2964 NT00\_cc\_l23, NT00\_cc\_l24, NT00\_cc\_l25, NT00\_cc\_l26, NT00\_cc\_l27, NT00\_cc\_l28,  
 2965 NT00\_cc\_l29, NT00\_cc\_l30, NT00\_cc\_l33, NT00\_cc\_l34, NT00\_cc\_l35, NT00\_cc\_l36,  
 2966 NT00\_cc\_l37, NT06\_cc\_l27]  
 2967 Group14, -0.07471531491655992, [NT18\_v\_l10, NT18\_v\_l12, NT18\_v\_l13, NT18\_v\_l14,  
 2968 NT18\_v\_l15, NT18\_v\_l16, NT18\_v\_l17, NT18\_v\_l18, NT18\_v\_l19, NT12\_v\_l10,  
 2969 NT12\_v\_l13, NT12\_v\_l14, NT12\_v\_l15, NT12\_v\_l16, NT12\_v\_l17, NT12\_v\_l18, NT12\_v\_l19,  
 2970 NT06\_v\_l10, NT06\_v\_l13, NT06\_v\_l14, NT06\_v\_l15, NT06\_v\_l16, NT06\_v\_l17,  
 2971 NT06\_v\_l18, NT06\_v\_l19, NT00\_v\_l10, NT00\_v\_l12, NT00\_v\_l13, NT00\_v\_l14,  
 2972 NT00\_v\_l15, NT00\_v\_l17, NT00\_v\_l18, NT00\_v\_l19, NT00\_v\_l16]  
 2973 Group15, 0.0, [NT18\_ciwc\_l1, NT18\_ciwc\_l2, NT18\_ciwc\_l3, NT18\_ciwc\_l4,  
 2974 NT18\_ciwc\_l5, NT18\_ciwc\_l6, NT18\_ciwc\_l7, NT18\_ciwc\_l8, NT18\_ciwc\_l9,  
 2975 NT12\_ciwc\_l1, NT12\_ciwc\_l2, NT12\_ciwc\_l3, NT12\_ciwc\_l4, NT12\_ciwc\_l5,  
 2976 NT12\_ciwc\_l7, NT12\_ciwc\_l8, NT06\_ciwc\_l1, NT06\_ciwc\_l2, NT06\_ciwc\_l3,

2977 NT06\_ciwc\_l4, NT06\_ciwc\_l5, NT06\_ciwc\_l6, NT06\_ciwc\_l7, NT06\_ciwc\_l8,  
 2978 NT00\_ciwc\_l1, NT00\_ciwc\_l2, NT00\_ciwc\_l3, NT00\_ciwc\_l4, NT00\_ciwc\_l5,  
 2979 NT00\_ciwc\_l6, NT00\_ciwc\_l7, NT00\_ciwc\_l8, NT12\_ciwc\_l6]  
 2980 Group16, 0.025737968406863487, [NT18\_r\_l1, NT18\_r\_l2, NT18\_r\_l3, NT18\_r\_l4,  
 2981 NT18\_r\_l5, NT18\_r\_l6, NT18\_r\_l7, NT18\_r\_l8, NT12\_r\_l1, NT12\_r\_l2, NT12\_r\_l3,  
 2982 NT12\_r\_l4, NT12\_r\_l5, NT12\_r\_l6, NT12\_r\_l7, NT12\_r\_l8, NT06\_r\_l1, NT06\_r\_l2,  
 2983 NT06\_r\_l3, NT06\_r\_l4, NT06\_r\_l5, NT06\_r\_l6, NT06\_r\_l7, NT06\_r\_l8, NT00\_r\_l1,  
 2984 NT00\_r\_l2, NT00\_r\_l3, NT00\_r\_l4, NT00\_r\_l6, NT00\_r\_l7, NT00\_r\_l8, NT00\_r\_l5]  
 2985 Group17, 0.0032523694495293354, [NT18\_z\_l26, NT18\_z\_l27, NT18\_z\_l28,  
 2986 NT18\_z\_l29, NT18\_z\_l30, NT18\_z\_l31, NT18\_z\_l32, NT18\_z\_l33, NT12\_z\_l26,  
 2987 NT12\_z\_l27, NT12\_z\_l28, NT12\_z\_l29, NT12\_z\_l30, NT12\_z\_l31, NT12\_z\_l32,  
 2988 NT12\_z\_l33, NT06\_z\_l26, NT06\_z\_l27, NT06\_z\_l28, NT06\_z\_l29, NT06\_z\_l30,  
 2989 NT06\_z\_l31, NT06\_z\_l32, NT00\_z\_l27, NT00\_z\_l28, NT00\_z\_l29, NT00\_z\_l30,  
 2990 NT00\_z\_l31, NT00\_z\_l32, NT00\_z\_l26]  
 2991 Group18, -0.012730578721368735, [NT18\_v\_l11, NT18\_v\_l12, NT18\_v\_l13,  
 2992 NT18\_v\_l14, NT18\_v\_l15, NT18\_v\_l16, NT18\_v\_l17, NT12\_v\_l11, NT12\_v\_l12,  
 2993 NT12\_v\_l13, NT12\_v\_l15, NT12\_v\_l16, NT12\_v\_l17, NT06\_v\_l11, NT06\_v\_l12,  
 2994 NT06\_v\_l13, NT06\_v\_l14, NT06\_v\_l15, NT06\_v\_l16, NT00\_v\_l12, NT00\_v\_l13,  
 2995 NT00\_v\_l14, NT00\_v\_l15, NT00\_v\_l16, NT12\_v\_l14]  
 2996 Group19, -0.0462370812127475, [NT18\_v\_l21, NT18\_v\_l22, NT18\_v\_l23,  
 2997 NT18\_v\_l24, NT18\_v\_l25, NT18\_v\_l26, NT12\_v\_l21, NT12\_v\_l22, NT12\_v\_l23,  
 2998 NT12\_v\_l24, NT12\_v\_l26, NT06\_v\_l21, NT06\_v\_l22, NT06\_v\_l23, NT06\_v\_l24,  
 2999 NT06\_v\_l25, NT06\_v\_l26, NT00\_v\_l21, NT00\_v\_l22, NT00\_v\_l23, NT00\_v\_l24,  
 3000 NT00\_v\_l25, NT00\_v\_l26, NT12\_v\_l25]  
 3001 Group20, -0.03416642562085348, [NT18\_u\_l20, NT18\_u\_l21, NT18\_u\_l22,  
 3002 NT18\_u\_l23, NT18\_u\_l24, NT18\_u\_l9, NT12\_u\_l20, NT12\_u\_l21, NT12\_u\_l22,  
 3003 NT12\_u\_l23, NT12\_u\_l24, NT12\_u\_l9, NT06\_u\_l20, NT06\_u\_l21, NT06\_u\_l22,  
 3004 NT06\_u\_l23, NT06\_u\_l9, NT00\_u\_l20, NT00\_u\_l21, NT00\_u\_l22, NT00\_u\_l9,  
 3005 NT00\_u\_l23]  
 3006 Group21, -0.015980995488767524, [NT18\_o3\_l4, NT18\_z\_l10, NT18\_z\_l8, NT18\_z\_l9,  
 3007 NT12\_o3\_l4, NT12\_z\_l8, NT12\_z\_l9, NT06\_o3\_l4, NT06\_z\_l7, NT06\_z\_l8,  
 3008 NT06\_z\_l9, NT00\_o3\_l4, NT00\_z\_l5, NT00\_z\_l6, NT00\_z\_l7, NT00\_z\_l8, NT00\_z\_l9,  
 3009 NT12\_z\_l10]  
 3010 Group22, -0.022620312860484826, [NT18\_z\_l3, NT18\_z\_l4, NT18\_z\_l5, NT18\_z\_l6,  
 3011 NT18\_z\_l7, NT12\_z\_l3, NT12\_z\_l4, NT12\_z\_l5, NT12\_z\_l6, NT12\_z\_l7, NT06\_z\_l3,  
 3012 NT06\_z\_l4, NT06\_z\_l5, NT06\_z\_l6, NT00\_z\_l2, NT00\_z\_l3, NT00\_z\_l4]  
 3013 Group23, -0.009525048312305717, [NT18\_ciwc\_l11, NT18\_ciwc\_l17, NT12\_ciwc\_l11,  
 3014 NT12\_q\_l34, NT06\_ciwc\_l11, NT06\_q\_l34, NT06\_q\_l35, NT00\_ciwc\_l17,  
 3015 NT00\_clwc\_l29, NT00\_clwc\_l30, NT00\_clwc\_l31, NT00\_q\_l34, NT00\_q\_l35,  
 3016 NT06\_ciwc\_l17]  
 3017 Group24, 0.016470284903957744, [NT18\_o3\_l37, NT12\_o3\_l37, NT06\_o3\_l37,  
 3018 NT06\_q\_l36, NT00\_o3\_l37, NT00\_q\_l28, NT00\_q\_l29, NT00\_q\_l30, NT00\_q\_l31,  
 3019 NT00\_q\_l32, NT00\_q\_l33, NT00\_q\_l36, NT06\_q\_l33]

3020 Group25, 0.003052627038463651, [NT06\_ciwc\_l25, NT06\_ciwc\_l26, NT06\_ciwc\_l27,  
 3021 NT06\_ciwc\_l29, NT06\_ciwc\_l30, NT06\_ciwc\_l31, NT06\_ciwc\_l32, NT06\_ciwc\_l33,  
 3022 NT06\_ciwc\_l34, NT06\_ciwc\_l35, NT06\_ciwc\_l36, NT06\_ciwc\_l37, NT06\_ciwc\_l28]  
 3023 Group26, 0.025737968406863487, [NT18\_ciwc\_l25, NT18\_ciwc\_l26, NT18\_ciwc\_l27,  
 3024 NT18\_ciwc\_l29, NT18\_ciwc\_l30, NT18\_ciwc\_l31, NT18\_ciwc\_l32, NT18\_ciwc\_l33,  
 3025 NT18\_ciwc\_l34, NT18\_ciwc\_l35, NT18\_ciwc\_l36, NT18\_ciwc\_l37, NT18\_ciwc\_l28]  
 3026 Group27, 0.031109314730230153, [NT18\_o3\_l18, NT18\_o3\_l19, NT18\_o3\_l20,  
 3027 NT12\_o3\_l18, NT12\_o3\_l19, NT12\_o3\_l20, NT06\_o3\_l18, NT06\_o3\_l19,  
 3028 NT06\_o3\_l20, NT00\_o3\_l19, NT00\_o3\_l20, NT00\_o3\_l18]  
 3029 Group28, 0.0032523694495293354, [NT18\_u\_l4, NT18\_u\_l5, NT18\_u\_l6, NT12\_u\_l4,  
 3030 NT12\_u\_l5, NT12\_u\_l6, NT06\_u\_l4, NT06\_u\_l6, NT00\_u\_l4, NT00\_u\_l5, NT00\_u\_l6,  
 3031 NT06\_u\_l5]  
 3032 Group29, 0.0712300559164073, [NT18\_u\_l1, NT18\_u\_l18, NT12\_u\_l1, NT12\_u\_l17,  
 3033 NT12\_u\_l18, NT06\_u\_l1, NT06\_u\_l17, NT00\_t\_l10, NT00\_u\_l16, NT00\_u\_l17,  
 3034 NT18\_u\_l17]  
 3035 Group30, 0.04046369350173762, [NT18\_u\_l10, NT18\_u\_l7, NT12\_u\_l10, NT12\_u\_l7,  
 3036 NT12\_u\_l8, NT06\_u\_l10, NT06\_u\_l7, NT06\_u\_l8, NT00\_u\_l7, NT00\_u\_l8,  
 3037 NT18\_u\_l8]  
 3038 Group31, 0.026777194098555834, [NT12\_ciwc\_l25, NT12\_ciwc\_l26, NT12\_ciwc\_l27,  
 3039 NT12\_ciwc\_l28, NT12\_ciwc\_l29, NT12\_ciwc\_l30, NT12\_ciwc\_l31, NT06\_ciwc\_l24,  
 3040 NT12\_ciwc\_l24]  
 3041 Group32, -0.0032999087322507226, [NT18\_d\_l33, NT18\_d\_l34, NT18\_d\_l35,  
 3042 NT18\_d\_l36, NT18\_d\_l37, NT12\_d\_l33, NT12\_d\_l34, NT00\_d\_l33, NT18\_d\_l32]  
 3043 Group33, -0.03088094284403675, [NT18\_u\_l3, NT12\_u\_l2, NT12\_u\_l3, NT06\_u\_l2,  
 3044 NT06\_u\_l3, NT00\_u\_l1, NT00\_u\_l2, NT00\_u\_l3, NT18\_u\_l2]  
 3045 Group34, -0.015980995488767524, [NT12\_d\_l35, NT06\_d\_l34, NT06\_d\_l35,  
 3046 NT06\_d\_l36, NT06\_d\_l37, NT00\_d\_l34, NT00\_d\_l35, NT06\_d\_l33]  
 3047 Group35, 0.016470284903957744, [NT18\_v\_l19, NT18\_v\_l20, NT12\_v\_l19,  
 3048 NT06\_v\_l19, NT06\_v\_l20, NT00\_v\_l19, NT00\_v\_l20, NT12\_v\_l20]  
 3049 Group36, -0.015980995488767524, [NT18\_o3\_l7, NT12\_o3\_l7, NT12\_o3\_l8,  
 3050 NT06\_o3\_l7, NT06\_o3\_l8, NT00\_o3\_l7, NT00\_o3\_l8, NT18\_o3\_l8]  
 3051 Group37, -0.009776487332806783, [NT18\_o3\_l22, NT12\_o3\_l21, NT12\_o3\_l22,  
 3052 NT06\_o3\_l21, NT06\_o3\_l22, NT00\_o3\_l21, NT00\_o3\_l22, NT18\_o3\_l21]  
 3053 Group38, 0.012734297741761935, [NT00\_ciwc\_l30, NT00\_ciwc\_l31, NT00\_ciwc\_l33,  
 3054 NT00\_ciwc\_l34, NT00\_ciwc\_l35, NT00\_ciwc\_l36, NT00\_ciwc\_l37, NT00\_ciwc\_l32]  
 3055 Group39, -0.01927724806311193, [NT18\_cc\_l35, NT18\_cc\_l36, NT18\_cc\_l37,  
 3056 NT18\_clwc\_l23, NT18\_clwc\_l36, NT18\_clwc\_l37, NT18\_cc\_l24]  
 3057 Group40, 0.029881404183944138, [NT18\_q\_l24, NT18\_q\_l25, NT18\_q\_l26,  
 3058 NT18\_q\_l27, NT12\_q\_l24, NT12\_q\_l25, NT12\_q\_l26]  
 3059 Group41, 0.02269292645906973, [NT18\_q\_l7, NT18\_q\_l8, NT12\_q\_l7, NT12\_q\_l8,  
 3060 NT06\_q\_l8, NT00\_q\_l8, NT06\_q\_l7]  
 3061 Group42, -0.0032999087322507226, [NT00\_ciwc\_l23, NT00\_ciwc\_l25,  
 3062 NT00\_ciwc\_l26, NT00\_ciwc\_l27, NT00\_ciwc\_l28, NT00\_ciwc\_l29, NT00\_ciwc\_l24]

3063 Group43, -0.022620312860484826, [NT18\_q\_l12, NT18\_q\_l13, NT12\_q\_l12,  
 3064 NT06\_q\_l12, NT00\_q\_l12, NT18\_q\_l14]  
 3065 Group44, 0.010064393684136919, [NT12\_u\_l19, NT06\_u\_l18, NT06\_u\_l19,  
 3066 NT00\_u\_l10, NT00\_u\_l19, NT18\_u\_l19]  
 3067 Group45, -0.025199275642583396, [NT12\_ciwc\_l32, NT12\_ciwc\_l33, NT12\_ciwc\_l34,  
 3068 NT12\_ciwc\_l36, NT12\_ciwc\_l37, NT12\_ciwc\_l35]  
 3069 Group46, 0.04046369350173762, [NT12\_z\_l36, NT06\_z\_l35, NT00\_z\_l35,  
 3070 NT00\_z\_l36, NT00\_z\_l37, NT06\_z\_l36]  
 3071 Group47, 0.03294056980791549, [NT18\_clwc\_l19, NT18\_clwc\_l20, NT00\_clwc\_l19,  
 3072 NT00\_clwc\_l20, NT06\_clwc\_l19]  
 3073 Group48, 0.04348999451297464, [NT18\_o3\_l23, NT12\_o3\_l23, NT06\_o3\_l23,  
 3074 NT00\_o3\_l23, NT18\_o3\_l24]  
 3075 Group49, 0.08554519157484475, [NT12\_v\_l18, NT06\_v\_l18, NT00\_v\_l17,  
 3076 NT00\_v\_l18, NT06\_v\_l17]  
 3077 Group50, 0.026777194098555834, [NT12\_cc\_l10, NT12\_ciwc\_l10, NT06\_ciwc\_l10,  
 3078 NT00\_cc\_l10, NT06\_cc\_l10]  
 3079 Group51, 0.013290747725846708, [NT18\_clwc\_l34, NT12\_clwc\_l35, NT12\_clwc\_l36,  
 3080 NT12\_clwc\_l34]  
 3081 Group52, -0.019890565190348153, [NT12\_z\_l34, NT06\_z\_l34, NT00\_z\_l33,  
 3082 NT06\_z\_l33]  
 3083 Group53, 0.02269292645906973, [NT18\_o3\_l3, NT06\_o3\_l3, NT00\_o3\_l3,  
 3084 NT12\_o3\_l3]  
 3085 Group54, 0.010064393684136919, [NT18\_o3\_l36, NT12\_o3\_l36, NT00\_o3\_l36,  
 3086 NT06\_o3\_l36]  
 3087 Group55, 0.01718412356692567, [NT12\_pv\_l2, NT06\_pv\_l2, NT00\_pv\_l2,  
 3088 NT18\_pv\_l2]  
 3089 Group56, 0.026777194098555834, [NT18\_q\_l11, NT12\_q\_l11, NT00\_q\_l11,  
 3090 NT06\_q\_l11]  
 3091 Group57, 0.0032523694495293354, [NT18\_r\_l9, NT12\_r\_l9, NT00\_r\_l9, NT06\_r\_l9]  
 3092 Group58, -0.02851378448120112, [NT12\_q\_l9, NT06\_q\_l9, NT00\_q\_l9, NT18\_q\_l9]  
 3093 Group59, -0.03528302096554892, [NT18\_q\_l10, NT06\_q\_l10, NT00\_q\_l10,  
 3094 NT12\_q\_l10]  
 3095 Group60, -0.051288793867985794, [NT00\_clwc\_l35, NT00\_clwc\_l36, NT00\_clwc\_l34]  
 3096 Group61, -0.01927724806311193, [NT18\_z\_l2, NT06\_z\_l2, NT12\_z\_l2]  
 3097 Group62, 0.025737968406863487, [NT12\_o3\_l6, NT00\_o3\_l6, NT06\_o3\_l6]  
 3098 Group63, 0.03892763335500793, [NT12\_pv\_l1, NT00\_pv\_l1, NT06\_pv\_l1]  
 3099 Group64, 0.006458219563131418, [NT18\_clwc\_l18, NT06\_clwc\_l18, NT12\_clwc\_l18]  
 3100 Group65, 0.05435473046696537, [NT00\_q\_l37, NT06\_q\_l37]  
 3101 Group66, 0.026777194098555834, [NT12\_d\_l37, NT12\_d\_l36]  
 3102 Group67, -0.013099278001989068, [NT00\_d\_l37, NT00\_d\_l36]  
 3103 Group68, 0.00961854131774742, [NT18\_z\_l35, NT18\_z\_l34]  
 3104 Group69, 0.02269292645906973, [NT18\_z\_l36, NT18\_z\_l37]  
 3105 Group70, 0.006149179017971962, [NT06\_t\_l12, NT12\_t\_l12]  
 3106 Group71, -0.009776487332806783, [NT18\_ciwc\_l10, NT18\_cc\_l10]

3107 Group72, -0.03187470567985029, [NT00\_cc\_l9, NT00\_ciwc\_l9]  
 3108 Group73, 0.03739216600762729, [NT06\_cc\_l9, NT06\_ciwc\_l9]  
 3109 Group74, -0.013099278001989068, [NT12\_cc\_l9, NT12\_ciwc\_l9]  
 3110 Group75, -0.022620312860484826, [NT00\_ciwc\_l22]  
 3111 Group76, 0.01960401690104774, [NT18\_clwc\_l35]  
 3112 Group77, 0.006458219563131418, [NT18\_ciwc\_l22]  
 3113 Group78, 0.00961854131774742, [NT12\_clwc\_l37]  
 3114 Group79, -0.06863045489336894, [NT00\_clwc\_l37]  
 3115 Group80, 0.05516234025094624, [NT00\_v\_l11]  
 3116 Group81, 0.04046369350173762, [NT00\_u\_l18]  
 3117 Group82, 0.031109314730230153, [NT18\_v\_l18]  
 3118 Group83, -0.016470284903957633, [NT00\_z\_l34]  
 3119 Group84, 0.03294056980791549, [NT12\_z\_l35]  
 3120 Group85, 0.034274390568905466, [NT12\_v\_l2]  
 3121 Group86, -0.013099278001989068, [NT06\_v\_l2]  
 3122 Group87, -0.03528302096554892, [NT00\_clwc\_l18]  
 3123 Group88, 0.07965337355823643, [NT18\_pv\_l1]  
 3124 Group89, 0.03294056980791549, [NT00\_t\_l12]  
 3125 Group90, 0.026777194098555834, [NT00\_ciwc\_l10]  
 3126 Group91, 0.02269292645906973, [NT06\_ciwc\_l23]  
 3127 Group92, -0.0065008857713687584, [NT06\_v\_l1]  
 3128 Group93, -0.05485916703871707, [NT18\_z\_l1]  
 3129 Group94, -0.047768415870933056, [NT18\_v\_l1]  
 3130 Group95, -0.047768415870933056, [NT18\_ciwc\_l23]  
 3131 Group96, 0.016470284903957744, [NT12\_ciwc\_l23]  
 3132 Group97, 0.0, [NT18\_o3\_l6]  
 3133 Group98, -0.042245901124683405, [NT00\_v\_l1]  
 3134 Group99, -0.053063711383979584, [NT12\_v\_l1]  
 3135 Group100, 0.01960401690104774, [NT12\_z\_l1]  
 3136 Group101, -0.03187470567985029, [NT00\_z\_l1]  
 3137 Group102, 0.03294056980791549, [NT06\_z\_l1]  
 3138 Group103, 0.0, [NT18\_cc\_l2]  
 3139 Group104, 0.0, [NT06\_cc\_l1]  
 3140 Group105, 0.0, [NT12\_cc\_l7]  
 3141 Group106, 0.0, [NT06\_cc\_l8]  
 3142 Group107, 0.0, [NT00\_cc\_l6]  
 3143 Group108, 0.0, [NT18\_cc\_l8]  
 3144 Group109, 0.0, [NT00\_cc\_l3]  
 3145 Group110, 0.0, [NT00\_cc\_l4]  
 3146 Group111, 0.0, [NT12\_cc\_l1]  
 3147 Group112, 0.0, [NT12\_cc\_l5]  
 3148 Group113, 0.0, [NT12\_cc\_l8]  
 3149 Group114, 0.0, [NT12\_cc\_l2]  
 3150 Group115, 0.0, [NT18\_cc\_l4]

3151 Group116, 0.0, [NT00\_cc\_11]  
 3152 Group117, 0.0, [NT00\_cc\_17]  
 3153 Group118, 0.0, [NT18\_cc\_15]  
 3154 Group119, 0.0, [NT00\_cc\_18]  
 3155 Group120, 0.0, [NT06\_cc\_14]  
 3156 Group121, 0.0, [NT06\_cc\_13]  
 3157 Group122, 0.0, [NT00\_cc\_15]  
 3158 Group123, 0.0, [NT06\_cc\_17]  
 3159 Group124, 0.0, [NT06\_cc\_15]  
 3160 Group125, 0.0, [NT18\_cc\_16]  
 3161 Group126, 0.0, [NT18\_cc\_19]  
 3162 Group127, 0.0, [NT12\_cc\_13]  
 3163 Group128, 0.0, [NT18\_cc\_13]  
 3164 Group129, 0.0, [NT06\_cc\_16]  
 3165 Group130, 0.0, [NT06\_cc\_12]  
 3166 Group131, 0.0, [NT18\_cc\_11]  
 3167 Group132, 0.0, [NT18\_cc\_17]  
 3168 Group133, 0.0, [NT00\_cc\_12]  
 3169 Group134, 0.0, [NT12\_cc\_14]  
 3170 Group135, 0.0, [NT12\_cc\_16]

3171

3172

3173 **Table A3: Group level importance from LLE data filter with the group number, the**  
 3174 **number of variables in each group, the importance score, and the rank.**

Group	Num	score	Rank
49	5	0.023614	1
88	1	0.021988	2
1	309	0.019687	3
29	11	0.019662	4
3	148	0.01728	5
80	1	0.015227	6
65	2	0.015004	7
48	5	0.012005	8
81	1	0.01117	9
46	6	0.01117	10
30	11	0.01117	11
63	3	0.010746	12
73	2	0.010322	13
85	1	0.009461	14
47	5	0.009093	15



84	1	0.009093	16
102	1	0.009093	17
89	1	0.009093	18
82	1	0.008587	19
27	12	0.008587	20
40	7	0.008249	21
90	1	0.007392	22
56	4	0.007392	23
50	5	0.007392	24
66	2	0.007392	25
31	9	0.007392	26
62	3	0.007105	27
16	32	0.007105	28
26	13	0.007105	29
13	34	0.006522	30
91	1	0.006264	31
41	7	0.006264	32
53	4	0.006264	33
69	2	0.006264	34
4	140	0.006264	35
100	1	0.005412	36
76	1	0.005412	37
55	4	0.004744	38
35	8	0.004546	39
96	1	0.004546	40
12	39	0.004546	41
24	13	0.004546	42
51	4	0.003669	43
38	8	0.003515	44
54	4	0.002778	45
44	6	0.002778	46
78	1	0.002655	47
68	2	0.002655	48
77	1	0.001783	49
10	49	0.001783	50
64	3	0.001783	51

70	2	0.001697	52
17	30	0.000898	53
57	4	0.000898	54
28	12	0.000898	55
25	13	0.000843	56
135	1	0	57
115	1	0	58
114	1	0	59
113	1	0	60
112	1	0	61
15	33	0	62
111	1	0	63
117	1	0	64
110	1	0	65
109	1	0	66
108	1	0	67
107	1	0	68
106	1	0	69
8	68	0	70
105	1	0	71
104	1	0	72
116	1	0	73
120	1	0	74
118	1	0	75
119	1	0	76
134	1	0	77
133	1	0	78
132	1	0	79
131	1	0	80
130	1	0	81
129	1	0	82
128	1	0	83
127	1	0	84
126	1	0	85
125	1	0	86
124	1	0	87

123	1	0	88
122	1	0	89
121	1	0	90
97	1	0	91
103	1	0	92
32	9	0	93
42	7	0	94
92	1	0	95
23	14	0	96
71	2	0	97
37	8	0	98
18	25	0	99
74	2	0	100
67	2	0	101
86	1	0	102
21	18	0	103
34	8	0	104
36	8	0	105
83	1	0	106
61	3	0	107
39	7	0	108
52	4	0	109
6	109	0	110
75	1	0	111
22	17	0	112
43	6	0	113
45	6	0	114
9	59	0	115
58	4	0	116
33	9	0	117
72	2	0	118
101	1	0	119
20	22	0	120
87	1	0	121
59	4	0	122
2	255	0	123

5	127	0	124
98	1	0	125
11	44	0	126
19	24	0	127
95	1	0	128
94	1	0	129
60	3	0	130
99	1	0	131
7	86	0	132
93	1	0	133
79	1	0	134
14	34	0	135

3175

3176

3177

**Table A4: COR-SHIPS model feature importance and its ranking.**

Importance score	Variable	Ranking
0.0362	BD12	1
0.0217	DTL	2
0.0207	CFLX	3
0.0206	SHRD	4
0.0205	G150	5
0.0204	jd	6
0.0199	VMAX	7
0.0199	IRM1_5	8
0.0191	PW08	9
0.019	VMPI	10
0.0187	SHTD	11
0.0183	IR00_12	12
0.018	HE07	13
0.0177	MTPW_1	14
0.0177	XD18	15
0.0175	SHTS	16
0.0173	PW14	17
0.0172	TWXC	18
0.0168	R000	19
0.0167	V300	20

0.0165	OAGE	21
0.0162	PSLV_1	22
0.0161	Z850	23
0.0161	SHRS	24
0.0157	SDDC	25
0.0156	VVAC	26
0.0156	PCM3	27
0.0154	TGRD	28
0.0153	T150	29
0.0153	CD26	30
0.0152	TADV	31
0.0151	V850	32
0.0148	PSLV_4	33
0.0145	PSLV_3	34
0.0145	REFC	35
0.0142	RD26	36
0.014	MTPW_19	37
0.0138	ND20	38
0.0138	XDML	39
0.0137	PENV	40
0.0137	EPSS	41
0.0136	G200	42
0.0134	IR00_3	43
0.0131	D200	44
0.013	NTFR	45
0.0124	T250	46
0.0124	O500	47
0.0124	IR00_20	48
0.012	NSST	49
0.0119	IRM1_16	50
0.0119	TLAT	51
0.0115	E000	52
0.0112	IRM3_17	53
0.0112	IRM3_11	54
0.0109	HIST_1	55
0.0109	G250	56

0.0109	RHMD	57
0.0106	NDFR	58
0.0104	IR00_17	59
0.0099	IRM1_17	60
0.0098	NOHC	61
0.0089	PEFC	62
0.0083	IR00_14	63
0.0075	IRM3_14	64
0.0063	PCM1	65
0.0057	IRM1_14	66
0.0052	NDTX	67
0.0043	HIST_8	68
0.0041	XDTX	69
0.0031	PC00	70
0.0023	IRM3_19	71
0.0019	HIST_15	72

3178

3179

3180

**Table A5: LLE-SHIPS model feature importance score and its ranking.**

Importance score	Variable	Ranking
0.0188	BD12	1
0.0167	VMAX	2
0.0138	DTL	3
0.013	SHRD	4
0.0115	TWXC	5
0.0113	G150	6
0.0113	VMPI	7
0.0113	REFC	8
0.0111	TGRD	9
0.0107	IRM1_5	10
0.0107	IR00_12	11
0.0105	V300	12
0.0105	VVAC	13
0.0103	G200	14
0.0096	PEFC	15
0.0096	MTPW_1	16

0.0095	XDTX	17
0.0095	PSLV_1	18
0.0095	T150	19
0.0094	CFLX	20
0.0094	HIST_2	21
0.0092	HE07	22
0.0091	SHTS	23
0.0089	PSLV_3	24
0.0085	SHTD	25
0.0085	G250	26
0.0085	CD26	27
0.0082	Ile84	28
0.0082	EPSS	29
0.0077	R000	30
0.0076	SDDC	31
0.0075	IRM3_19	32
0.0075	RD26	33
0.0074	PW08	34
0.0074	SHRS	35
0.0074	NDTX	36
0.0073	Ile75	37
0.0073	jd	38
0.0073	TADV	39
0.0072	NDFR	40
0.0072	E000	41
0.0072	HIST_9	42
0.0072	PSLV_4	43
0.0071	MTPW_19	44
0.007	IRM1_16	45
0.007	PW14	46
0.007	OAGE	47
0.007	Ile78	48
0.0068	XD18	49
0.0068	Ile1	50
0.0066	Ile49	51
0.0066	Ile24	52

0.0065	PCM1	53
0.0065	ND20	54
0.0064	Ile2	55
0.0064	IR00_20	56
0.0063	NSST	57
0.0062	Z850	58
0.0062	NTFR	59
0.0061	IR00_14	60
0.0061	NOHC	61
0.0061	Ile89	62
0.006	IRM3_14	63
0.006	Ile71	64
0.006	Ile3	65
0.006	Ile52	66
0.0059	Ile19	67
0.0059	Ile51	68
0.0059	Ile72	69
0.0059	IR00_17	70
0.0059	Ile4	71
0.0058	O500	72
0.0058	Ile53	73
0.0058	V850	74
0.0058	TLAT	75
0.0057	Ile66	76
0.0057	Ile30	77
0.0056	RHMD	78
0.0056	Ile60	79
0.0056	Ile16	80
0.0055	IR00_3	81
0.0055	Ile76	82
0.0055	Ile57	83
0.0054	Ile69	84
0.0054	Ile54	85
0.0054	PC00	86
0.0053	Ile8	87
0.0052	Ile77	88



0.0052	IRM3_11	89
0.0052	IRM1_14	90
0.0052	Ile33	91
0.0051	Ile10	92
0.0051	IRM1_17	93
0.0051	Ile55	94
0.0051	Ile56	95
0.0051	T250	96
0.005	Ile26	97
0.005	Ile81	98
0.005	Ile73	99
0.0049	XDML	100
0.0049	Ile13	101
0.0049	D200	102
0.0048	Ile17	103
0.0048	Ile37	104
0.0048	Ile21	105
0.0048	Ile20	106
0.0048	Ile48	107
0.0048	Ile9	108
0.0047	Ile7	109
0.0047	Ile39	110
0.0047	Ile29	111
0.0046	Ile11	112
0.0046	IRM3_17	113
0.0046	Ile44	114
0.0046	Ile12	115
0.0046	Ile42	116
0.0045	Ile27	117
0.0045	Ile74	118
0.0045	Ile22	119
0.0044	Ile62	120
0.0044	Ile82	121
0.0044	PCM3	122
0.0044	Ile47	123
0.0044	Ile32	124

0.0043	Ile43	125
0.0043	Ile80	126
0.0043	Ile83	127
0.0043	Ile31	128
0.0042	Ile6	129
0.0042	Ile15	130
0.0042	Ile70	131
0.0042	Ile63	132
0.0042	Ile38	133
0.0042	Ile46	134
0.0042	Ile23	135
0.0041	Ile34	136
0.0041	PENV	137
0.0041	Ile58	138
0.004	Ile36	139
0.0039	Ile59	140
0.0039	Ile68	141
0.0039	Ile65	142
0.0039	Ile25	143
0.0039	Ile28	144
0.0038	Ile45	145
0.0038	Ile14	146
0.0038	Ile64	147
0.0038	Ile61	148
0.0037	Ile41	149
0.0037	Ile85	150
0.0037	Ile90	151
0.0036	Ile88	152
0.0036	Ile79	153
0.0035	Ile87	154
0.0035	Ile5	155
0.0035	Ile35	156
0.0034	Ile18	157
0.0034	Ile67	158
0.0034	Ile86	159
0.0032	Ile50	160

0.003	lle40	161
0.0018	HIST_16	162

3181

3182

3183

**Table A6: DL-SHIPS model feature importance score, and its ranking.**

Importance score	Variable	Ranking
0.0197	BD12	1
0.0176	VMAX	2
0.0148	SHRD	3
0.0144	DTL	4
0.0137	IRM1_5	5
0.0133	o31	6
0.0131	G150	7
0.0131	q7	8
0.0129	u3	9
0.0129	q4	10
0.0129	G200	11
0.0127	vo3	12
0.0124	REFC	13
0.0122	vo5	14
0.012	vo8	15
0.012	PEFC	16
0.0118	d3	17
0.0116	CFLX	18
0.0116	PSLV_3	19
0.0114	T150	20
0.0114	jd	21
0.0114	R000	22
0.0112	TWXC	23
0.0112	u8	24
0.0112	PW08	25
0.0112	q3	26
0.0112	XDTX	27
0.0109	CD26	28
0.0109	q8	29
0.0107	pv3	30

0.0107	v4	31
0.0105	r1	32
0.0101	u1	33
0.0099	q5	34
0.0099	IR00_12	35
0.0099	vo4	36
0.0097	HE07	37
0.0097	u6	38
0.0097	q2	39
0.0094	r6	40
0.0094	vo6	41
0.0094	MTPW_1	42
0.0092	u2	43
0.0092	r4	44
0.0092	pv7	45
0.0092	pv6	46
0.009	PSLV_1	47
0.009	TADV	48
0.0088	v8	49
0.0088	HIST_2	50
0.0088	VMPI	51
0.0088	V300	52
0.0088	SHRS	53
0.0086	VVAC	54
0.0086	MTPW_19	55
0.0086	v5	56
0.0082	t1	57
0.0082	RD26	58
0.0082	SDDC	59
0.0082	q6	60
0.0082	O500	61
0.0082	v7	62
0.0082	IRM3_11	63
0.0079	E000	64
0.0079	PW14	65
0.0077	z2	66

0.0077	G250	67
0.0075	pv1	68
0.0075	cc1	69
0.0075	XDML	70
0.0075	pv8	71
0.0073	vo1	72
0.0073	ciwc1	73
0.0073	v3	74
0.0073	SHTS	75
0.0073	v6	76
0.0071	ciwc2	77
0.0071	w1	78
0.0071	IRM3_19	79
0.0071	IR00_17	80
0.0069	Z850	81
0.0069	SHTD	82
0.0069	NOHC	83
0.0067	OAGE	84
0.0064	XD18	85
0.0064	IR00_3	86
0.0064	IRM1_16	87
0.0064	PSLV_4	88
0.0062	NTFR	89
0.0062	HIST_9	90
0.006	ND20	91
0.006	IR00_14	92
0.0058	IRM3_17	93
0.0058	EPSS	94
0.0058	clwc2	95
0.0058	D200	96
0.0058	V850	97
0.0056	PC00	98
0.0056	r8	99
0.0054	u5	100
0.0054	NDFR	101
0.0052	PCM1	102

0.0052	NSST	103
0.0052	PENV	104
0.0052	TGRD	105
0.0047	IRM3_14	106
0.0047	IR00_20	107
0.0047	T250	108
0.0047	RHMD	109
0.0047	IRM1_14	110
0.0047	IRM1_17	111
0.0045	cc2	112
0.0034	NDTX	113
0.0034	r7	114
0.0032	TLAT	115
0.0024	PCM3	116
0.0015	HIST_16	117
0.0006	r2	118
0	r3	119
0	r5	120

3184  
 3185  
 3186  
 3187  
 3188  
 3189  
 3190  
 3191  
 3192  
 3193  
 3194  
 3195  
 3196  
 3197  
 3198  
 3199  
 3200  
 3201  
 3202  
 3203  
 3204  
 3205

3206  
3207  
3208  
3209  
3210

## REFERENCES

- 3211 Asif, A., M. Dawood, B. Jan, J. Khurshid, M. DeMaria and F. Minhas, 2020: PHURIE:  
3212 hurricane intensity estimation from infrared satellite imagery using machine  
3213 learning. *Neural Comput & Applic*, **32**, 4821–4834.
- 3214 Akbani, R., S. Kwek, and N. Japkowicz, 2004: Applying support vector machines to  
3215 imbalanced datasets. In *European conference on machine learning*. Springer,  
3216 Berlin, Heidelberg, 39-50.
- 3217 Anyfantis, D., M. Karagiannopoulos, S. Kotsiantis, and P. Pintelas, 2007: Robustness of  
3218 learning techniques in handling class noise in imbalanced datasets. In *IFIP*  
3219 *International Conference on Artificial Intelligence Applications and Innovations*.  
3220 Springer, Boston, MA, 21-28.
- 3221 Astier, N. , M. Plu, and C. Claud, 2015: Associations between tropical cyclone activity in  
3222 the Southwest Indian Ocean and El Niño Southern Oscillation. *Atmos. Sci. Lett.*,  
3223 **16**, 506-511.
- 3224 Batista, G. E., R. C. Prati, and M. C. Monard, 2004: A study of the behavior of several  
3225 methods for balancing machine learning training data. *ACM SIGKDD explorations*  
3226 *newsletter*, 6(1), 20-29.



3227 Berrisford, P., D. Dee, P. Poli, R. Brugge, K. Fielding, M. Fuentes, P. Kallberg, S.  
 3228 Kobayashi, S. Uppala, and A. Simmons, 2011: The ERA-Interim archive, version  
 3229 2.0.  
 3230 Berger, E, 2016: The European forecast model already kicking America's butt just  
 3231 improved. *Ars Technica*, 11 March 2016. Retrieved 16 August 2016.  
 3232 Breiman, L., 2001: Random forests. *Machine learning*, **45**(1), 5-32.  
 3233 Breiman, L., 2017: Classification and regression trees. Routledge.  
 3234 Bunkhumpornpat, C., K. Sinapiromsaran, and C. Lursinsap, 2009: Safe-level-smote:  
 3235 Safe-level-synthetic minority over-sampling technique for handling the class  
 3236 imbalanced problem. In *Pacific-Asia conference on knowledge discovery and data*  
 3237 *mining*. Springer, Berlin, Heidelberg, 475-482.  
 3238 Szegedy, C., W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V.  
 3239 Vanhoucke, and A. Rabinovich, 2015: Going deeper with convolutions. In  
 3240 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*  
 3241 (CVPR), 1–9.  
 3242 Cangialosi, J. P., and J. L. Franklin, 2017: 2016 National Hurricane Center Forecast  
 3243 Verification Report. National Hurricane Center.  
 3244 Castro, C. L., and A. P. Braga, 2013: Novel cost-sensitive approach to improve the  
 3245 multilayer perceptron performance on imbalanced data. *IEEE transactions on*  
 3246 *neural networks and learning systems*, **24**(6), 888-899

3247 Chawla, N. V., K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, 2002: SMOTE:  
3248 synthetic minority over-sampling technique. *Journal of artificial intelligence*  
3249 *research*, **16**, 321-357.

3250 Chawla, N. V., N. Japkowicz, and A. Kotcz, 2004: Special issue on learning from  
3251 imbalanced data sets. *ACM SIGKDD explorations newsletter*, **6**(1): 1–6.

3252 Chen, T., and C. Guestrin, 2016: Xgboost: A scalable tree boosting system. In  
3253 *Proceedings of the 22nd acm sigkdd international conference on knowledge*  
3254 *discovery and data mining*. ACM, 785-794.

3255 Cohen, W. W., 1995: Fast effective rule induction. In *Machine Learning Proceedings*  
3256 *1995*, 115-123.

3257 Cummings, J. A., 2005: Operational multivariate ocean data assimilation. *Quarterly*  
3258 *Journal of the Royal Meteorological Society*, **131**(613), 3583-3604.

3259 Dee, D. P., S. M. Uppala, A. J. Simmons, P. Berrisford, P. Poli, S. Kobayashi, U. Andrae,  
3260 M.A. Balmaseda, G. Balsamo, P. Bauer, P. Bechtold, A. C. M. Beljaars, L. van de  
3261 Berg, J. Bidlot, N. Bormann, C. Delsol, R. Dragani, M. Fuentes, A. J. Geer, L.  
3262 Haimberger, S. B. Healy, H. Hersbach, E.V. Hólm, L. Isaksen, P. Kållberg, M.  
3263 Köhler, M. Matricardi, A. P. McNally, B. M. Monge-Sanz, J. -J. Morcrette, B. -K.  
3264 Park, C. Peubey, P. de Rosnay, C. Tavalato, J. -N. Thépaut, and F. Vitart, 2011:  
3265 The ERA-Interim reanalysis: Configuration and performance of the data  
3266 assimilation system. *Quarterly Journal of the royal meteorological society*, **137**,  
3267 553-597.

3268 DeMaria, M., 2009: A simplified dynamical system for tropical cyclone intensity  
 3269 prediction. *Monthly Weather Review*, **137**(1), 68-82.

3270 DeMaria, M., and J. Kaplan, 1994: A statistical hurricane intensity prediction scheme  
 3271 (SHIPS) for the Atlantic basin. *Weather and Forecasting*, **9**(2), 209-220.

3272 DeMaria, M., and J. Kaplan, 1999: An Updated Statistical Hurricane Intensity Prediction  
 3273 Scheme (SHIPS) for the Atlantic and Eastern North Pacific Basins Mark. *Weather*  
 3274 *and Forecasting*, **14**(3), 326–337.

3275 DeMaria, M., J. A. Knaff, and C. Sampson, 2007: Evaluation of long-term trends in  
 3276 tropical cyclone intensity forecasts. *Meteorology and Atmospheric Physics*, **97**(1-  
 3277 4), 19.

3278 DeMaria, M., M. Mainelli, L. K. Shay, J. A. Knaff, and J. Kaplan, 2005: Further  
 3279 improvements to the statistical hurricane intensity prediction scheme  
 3280 (SHIPS). *Weather and Forecasting*, **20**(4), 531-543.

3281 Dempster, A. P., N. M. Laird, and D. B. Rubin, 1977: Maximum likelihood from  
 3282 incomplete data via the EM algorithm. *Journal of the royal statistical society. Series*  
 3283 *B (methodological)*, **39**(1), 1-22

3284 Enz, R., P. Zimmerli, and S. Schwarz, 2009: Natural catastrophes and man-made  
 3285 disasters in 2008: North America and Asia suffer heavy losses. *National*  
 3286 *Emergency Training Center*.

3287 Ferrara, M., F. Groff, Z. Moon, K. Keshavamurthy, S. M. Robeson, and C. Kieu, 2017:  
 3288 Large-scale control of the lower stratosphere on variability of tropical cyclone  
 3289 intensity. *Geophysical Research Letters*, **44**(9), 4313-4323.

3290 Fraley, C., and A. E. Raftery, 1998: How many clusters? Which clustering method?  
3291 Answers via model-based cluster analysis. *The computer journal*, **41**(8), 578-588.

3292 Freund, Y., and L. Mason, 1999: The alternating decision tree learning algorithm. In *icml*,  
3293 **99**, 124-133.

3294 Friedl, M. A., and C. E. Brodley, 1997: Decision tree classification of land cover from  
3295 remotely sensed data. *Remote Sensing of Environment*, **61**(3), 399-409.

3296 Friedman, J., T. Hastie, and R. Tibshirani, 2001: The elements of statistical learning. *New*  
3297 *York: Springer series in statistics*, **1**, 337-387.

3298 Galar, M., A. Fernandez, E. Barrenechea, H. Bustince, and F. Herrera, 2012: A review on  
3299 ensembles for the class imbalance problem: bagging-, boosting-, and hybrid-based  
3300 approaches. *IEEE Transactions on Systems, Man, and Cybernetics, Part C*  
3301 *(Applications and Reviews)*, **42**(4), 463-484

3302 Géron, A., 2017: Hands-on machine learning with Scikit-Learn and TensorFlow:  
3303 concepts, tools, and techniques to build intelligent systems. *O'Reilly Media, Inc.*

3304 Ginsburg, S. B., G. Lee, S. Ali, and A. Madabhushi, 2016: Feature importance in  
3305 nonlinear embeddings (FINE): applications in digital pathology. *IEEE transactions*  
3306 *on medical imaging*, **35**(1), 76-88.

3307 Gogna, A., and A. Majumdar, 2019: Discriminative Autoencoder for Feature Extraction:  
3308 Application to Character Recognition. *Neural Processing Letters*, **49**(3), 1723-  
3309 1735.

3310 Gray, W. M., and L. R. Brody, 1967: Global view of the origin of tropical disturbances  
3311 and storms. *Colorado State University, Department of Atmospheric Science*.

3312 Han, H., W. Y. Wang, and B. H. Mao, 2005: Borderline-SMOTE: a new over-sampling  
 3313 method in imbalanced data sets learning. In *International Conference on Intelligent*  
 3314 *Computing*. Springer, Berlin, Heidelberg, 878-887.

3315 Hinton, G. E., N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, 2012:  
 3316 Improving neural networks by preventing co-adaptation of feature detectors. *arXiv*  
 3317 *preprint arXiv:1207.0580*.

3318 Hodgkin, A. L., and A. F. Huxley, 1952: A quantitative description of membrane current  
 3319 and its application to conduction and excitation in nerve. *The Journal of*  
 3320 *physiology*, **117**(4), 500-544.

3321 Holmes, G., A. Donkin, and I. H. Witten, 1994: Weka: A machine learning workbench.  
 3322 In *Proceedings of ANZIS'94-Australian New Zealand Intelligent Information*  
 3323 *Systems Conference*. IEEE. 357-361.

3324 Holyoak, K. J., 1987: Parallel distributed processing: explorations in the microstructure  
 3325 of cognition. *Science*, **236**, 992-997.

3326 Hotelling, H., 1933: Analysis of a complex of statistical variables into principal  
 3327 components. *Journal of educational psychology*, **24**(6), 417.

3328 James, G., D. Witten, T. Hastie, and R. Tibshirani, 2013: *An Introduction to Statistical*  
 3329 *Learning*. Springer.

3330 Japkowicz, N., 2000: The class imbalance problem: Significance and strategies.  
 3331 In *Proceeding of the Int'l Conf. on Artificial Intelligence*.

3332 Jarvinen, B. R., and C. J. Neumann, 1979: Statistical forecasts of tropical cyclone  
 3333 intensity for the North Atlantic basin.

3334 Jo, T., and N. Japkowicz, 2004: Class imbalances versus small disjuncts. *ACM Sigkdd*  
3335 *Explorations Newsletter*, **6**(1), 40-49

3336 Kaplan, J., and M. DeMaria, 2003: Large-scale characteristics of rapidly intensifying  
3337 tropical cyclones in the North Atlantic basin. *Weather and Forecasting*, **18**(6),  
3338 1093-1108.

3339 Kaplan, J., M. DeMaria, and J. A. Knaff, 2010: A revised tropical cyclone rapid  
3340 intensification index for the Atlantic and eastern North Pacific basins. *Weather and*  
3341 *Forecasting*, **25**(1), 220-241.

3342 Kaplan, J., C. M. Rozoff, M. DeMaria, C. R. Sampson, J. P. Kossin, C. S. Velden, J. J.  
3343 Cione, J. P. Dunion, J. A. Knaff, J. A. Zhang, J. F. Dostalek, J. D. Hawkins, T. F.  
3344 Lee, and J. E. Solbrig, 2015: Evaluating environmental impacts on tropical cyclone  
3345 rapid intensification predictability utilizing statistical models. *Weather and*  
3346 *Forecasting*, **30**(5), 1374-1396.

3347 Karpathy, A., G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, 2014:  
3348 Large-scale video classification with convolutional neural networks.  
3349 In *Proceedings of the IEEE conference on Computer Vision and Pattern*  
3350 *Recognition*, 1725-1732.

3351 Kidder, S. Q., and A. S. Jones, 2007: A blended satellite total precipitable water product  
3352 for operational forecasting. *Journal of Atmospheric and Oceanic*  
3353 *Technology*, **24**(1), 74-81.

3354 Knaff, J. A., DeMaria, M., Sampson, C. R., & Gross, J. M. (2003). Statistical, 5-day  
 3355 tropical cyclone intensity forecasts derived from climatology and  
 3356 persistence. *Weather and Forecasting*, **18**(1), 80-92.

3357 Knaff, J. A., C. R. Sampson, and M. DeMaria, 2005: An operational statistical typhoon  
 3358 intensity prediction scheme for the western North Pacific. *Weather and*  
 3359 *Forecasting*, **20**(4), 688-699.

3360 Knaff, J. A., T. A. Cram, A. B. Schumacher, J. P. Kossin, and M. DeMaria,  
 3361 2008: Objective identification of annular hurricanes. *Weather and*  
 3362 *Forecasting*, **23**(1), 17–88.

3363 Kohavi, R., 1995: A study of cross-validation and bootstrap for accuracy estimation and  
 3364 model selection. In *IJCAI*, **14**(2), 1137-1145.

3365 Kohavi, R., 1996: Scaling up the accuracy of Naive-Bayes classifiers: a decision-tree  
 3366 hybrid. In *KDD*, **96**, 202-207.

3367 Krawczyk, B., M. Woźniak, and G. Schaefer, 2014: Cost-sensitive decision tree  
 3368 ensembles for effective imbalanced classification. *Applied Soft Computing*, **14**,  
 3369 554-562.

3370 Krizhevsky, A., I. Sutskever, and G. E. Hinton, 2012: Imagenet classification with deep  
 3371 convolutional neural networks. In *Advances in neural information processing*  
 3372 *systems*, 1097-1105.

3373 Kurihara, Y., R. E. Tuleya, and M. A. Bender, 1998: The GFDL hurricane prediction  
 3374 system and its performance in the 1995 hurricane season. *Monthly weather*  
 3375 *review*, **126**(5), 1306-1322.

- 3376 Labib, K., and V. R. Vemuri, 2006: An application of principal component analysis to the  
3377 detection and visualization of computer network attacks. In *Annales des*  
3378 *télécommunications*, Springer-Verlag, **61**(1-2), 218-234.
- 3379 Lai, S., L. Xu, K. Liu, and J. Zhao, 2015: Recurrent Convolutional Neural Networks for  
3380 Text Classification. In *AAAI*, **333**, 2267-2273.
- 3381 Landsea, C. W. and J. L. Franklin, 2013: Atlantic Hurricane Database Uncertainty and  
3382 Presentation of a New Database Format. *Monthly Weather Review*, **141**, 3576-  
3383 3592.
- 3384 Landwehr, N., M. Hall, and E. Frank, 2005: Logistic model trees. *Machine*  
3385 *learning*, **59**(1-2), 161-205.
- 3386 Last, F., G. Douzas, and F. Bacao, 2017: Oversampling for Imbalanced Learning Based  
3387 on K-Means and SMOTE. *arXiv preprint arXiv:1711.00837*.
- 3388 Li, H., and P. Ralph, 2019: Local PCA shows how the effect of population structure  
3389 differs along the genome. *Genetics*, **211**(1), 289-304.
- 3390 Li, J., S. Pan, Y. Chen, and Y. Pan, 2017: Assessment of tropical cyclones in ECMWF  
3391 reanalysis data over Northwest Pacific Ocean. In *The 27th International Ocean and*  
3392 *Polar Engineering Conference*. International Society of Offshore and Polar  
3393 Engineers
- 3394 Liu, Y., A. An, and X. Huang, 2006: Boosting prediction accuracy on imbalanced  
3395 datasets with SVM ensembles. In *Pacific-Asia Conference on Knowledge*  
3396 *Discovery and Data Mining*, Springer, Berlin, Heidelberg, 107-118.



3397 Liu, Y., E. Racah, J. Correa, A. Khosrowshahi, D. Lavers, K. Kunkel, M. Wehner, and  
 3398 W. Collins, 2016: Application of deep convolutional neural networks for detecting  
 3399 extreme weather in climate datasets, *arXiv preprint arXiv:1605.01156*.  
 3400 López, V., S. del Río, J. M. Benítez, and F. Herrera, 2015: Cost-sensitive linguistic fuzzy  
 3401 rule based classification systems under the MapReduce framework for imbalanced  
 3402 big data. *Fuzzy Sets and Systems*, **258**, 5-38.  
 3403 Lundberg, S. M., and S. I. Lee, 2017: A unified approach to interpreting model  
 3404 predictions. In *Advances in neural information processing systems*, 4765-4774.  
 3405 Martinez, J., M. M. Bell, R. F. Rogers, and J. D. Doyle, 2019: Axisymmetric Potential  
 3406 Vorticity Evolution of Hurricane Patricia (2015). *J. Atmos. Sci.*, **76**, 2043–  
 3407 2063, <https://doi.org/10.1175/JAS-D-18-0373.1>.  
 3408 Merrill, R. T., 1987: An experiment in the statistical prediction of tropical cyclone  
 3409 intensity change. In *17th Conf. Hurricanes and Tropical Meteorology*, Hurricanes  
 3410 and Tropical Meteorology, Miami, FL. *Amer. Meteor. Soc.*, 302-304.  
 3411 Merrill, R. T. (1988). Environmental influences on hurricane intensification. *Journal of*  
 3412 *the Atmospheric Sciences*, **45**(11), 1678-1687.  
 3413 Miller, D., 2007: NEW ADVANCED HURRICANE MODEL AIDS NOAA  
 3414 FORECASTERS. <http://www.noaanews.noaa.gov/stories2007/s2885.html>.  
 3415 [Accessed 29 May 2017](#)  
 3416 Molinari, J., and D. Vollaro, 1989: External influences on hurricane intensity. Part I:  
 3417 Outflow layer eddy angular momentum fluxes. *Journal of the Atmospheric*  
 3418 *Sciences*, **46**(8), 1093-1105.

- 3419 Molnar, C., 2018: Interpretable machine learning: A guide for making black box models  
3420 explainable. *Christoph Molnar, Leanpub*.
- 3421 Nguyen, H. M., E. W. Cooper, and K. Kamei, 2011: Borderline over-sampling for  
3422 imbalanced data classification. *International Journal of Knowledge Engineering  
3423 and Soft Data Paradigms*, **3**(1), 4-21.
- 3424 NHC (National Hurricane Center), 2019: *NHC Track and Intensity Models*. Web site:  
3425 <http://www.nhc.noaa.gov/modelsummary.shtml>. Last accessed on August 18, 2020.
- 3426 Pacific Disaster Center. (n.d.): Tropical Cyclones. Retrieved November 20, 2017, from  
3427 <http://www.pdc.org/resources/natural-hazards/tropical-cyclones/>.
- 3428 Pearson, K., 1901: LIII. On lines and planes of closest fit to systems of points in  
3429 space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of  
3430 Science*, **2**(11), 559-572.
- 3431 Piñeros, M. F., E. A. Ritchie, and J. S. Tyo, 2011: Estimating tropical cyclone intensity  
3432 from infrared image data. *Weather and Forecasting*, **26**(5), 690–698.
- 3433 Platt, J. C., 1999: 12 fast training of support vector machines using sequential minimal  
3434 optimization. *Advances in kernel methods*, 185-208.
- 3435 Qian, Y., C. Liang, S. Peng, S. Chen, and S. Wang, 2016: [A Horizontal Index for the  
3436 Influence of Upper-Level Environmental Flow on Tropical Cyclone Intensity](#).  
3437 *Weather and Forecasting*, **31**, 237–253, <https://doi.org/10.1175/WAF-D-15-0091.1>
- 3438 Quinlan, J. R., 2014: C4. 5: programs for machine learning. *Elsevier*.

3439 Racah, E., C. Beckham, T. Maharaj, and C. Pal, 2016: Semi-Supervised Detection of  
 3440 Extreme Weather Events in Large Climate Datasets. *arXiv preprint*  
 3441 *arXiv:1612.02095*.  
 3442 Reid, W., 1846: An Attempt to Develop the Law of Storms by Means of Facts, Arranged  
 3443 According to Place and Time: And Hence to Point Out a Cause for the Variable  
 3444 Winds, with the View to Practical Use in Navigation. *Illustrated by Charts and*  
 3445 *Wood Cuts*. John Weale.  
 3446 Rennick, M. A., 1999: Performance of the Navy's tropical cyclone prediction model in  
 3447 the western North Pacific basin during 1996. *Weather and Forecasting*, **14**(3), 297-  
 3448 305.  
 3449 Reynolds, D. A., T. F. Quatieri, and R. B. Dunn, 2000: Speaker verification using  
 3450 adapted Gaussian mixture models. *Digital signal processing*, **10** (1-3), 19-41.  
 3451 Rhome, J. R., 2007: Technical summary of the National Hurricane Center track and  
 3452 intensity models.  
 3453 Ribeiro, M. T., S. Singh, and C. Guestrin, 2016: " Why should I trust you?" Explaining  
 3454 the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD*  
 3455 *international conference on knowledge discovery and data mining*, 1135-1144.  
 3456 Ritchie EA, G. Valliere-Kelley, M. F. Piñeros, J. S. Tyo, 2012: Tropical cyclone intensity  
 3457 estimation in the North Atlantic Basin using an improved deviation angle variance  
 3458 technique. *Weather and Forecasting*, **27**(5): 1264–1277.  
 3459 Ron, U., 2000: A practical guide to swap curve construction. Ottawa: Bank of Canada,  
 3460 **2000** (17).

3461 Roweis, S. T., and L. K. Saul, 2000: Nonlinear dimensionality reduction by locally linear  
 3462 embedding. *science*, **290**(5500), 2323-2326.

3463 Rozoff, C. M., and J. P. Kossin, 2011: New probabilistic forecast models for the  
 3464 prediction of tropical cyclone rapid intensification. *Weather and*  
 3465 *Forecasting*, **26**(5), 677-689.

3466 Ruder, S. (2016). An overview of gradient descent optimization algorithms. *arXiv*  
 3467 *preprint arXiv:1609.04747*.

3468 Sampson, C. R., J. L. Franklin, J. A. Knaff, and M. DeMaria, 2008: Experiments with a  
 3469 simple tropical cyclone intensity consensus. *Weather and Forecasting*, **23**(2), 304-  
 3470 312.

3471 Santos, d. C., M. Gatti, 2014: Deep Convolutional Neural Networks for Sentiment  
 3472 Analysis of Short Texts. In *Proceedings of COLING 2014, the 25th International*  
 3473 *Conference on Computational Linguistics: Technical Papers*, 69-78.

3474 Schölkopf, B., A. Smola, and K. R. Müller, 1997: Kernel principal component analysis.  
 3475 In *International Conference on Artificial Neural Networks*, Springer, Berlin,  
 3476 Heidelberg, 583-588.

3477 Schölkopf, B., A. Smola, and K. R. Müller, 1998: Nonlinear component analysis as a  
 3478 kernel eigenvalue problem. *Neural computation*, **10**(5), 1299-1319.

3479 Shahriari, B., K. Swersky, Z. Wang, R. P. Adams, and N. De Freitas, 2015: Taking the  
 3480 human out of the loop: A review of Bayesian optimization. *Proceedings of the*  
 3481 *IEEE*, **104**(1), 148-175.

3482 Shaiba, H. and M. Hahsler, 2016: Research Article Applying Machine Learning Methods  
 3483 for Predicting Tropical Cyclone Rapid Intensification Events. *Research Journal of*  
 3484 *Applied Sciences, Engineering and Technology* **13**(8), 638-651.  
 3485 Sheets, R. C., 1990: The National Hurricane Center—past, present, and future. *Weather*  
 3486 *and Forecasting*, **5**(2), 185-232.  
 3487 SHIPS, 2018a: SHIPS statistical tropical cyclone intensity forecast technique  
 3488 development, developmental data. [Available online  
 3489 at [http://rammb.cira.colostate.edu/research/tropical\\_cyclones/ships/developmental\\_data.as](http://rammb.cira.colostate.edu/research/tropical_cyclones/ships/developmental_data.asp)  
 3490 [p.](http://rammb.cira.colostate.edu/research/tropical_cyclones/ships/developmental_data.asp)]  
 3491 SHIPS, 2018b:  
 3492 [http://rammb.cira.colostate.edu/research/tropical\\_cyclones/ships/docs/AL/lsdiaga\\_1](http://rammb.cira.colostate.edu/research/tropical_cyclones/ships/docs/AL/lsdiaga_1982_2017_sat_ts.dat)  
 3493 [982\\_2017\\_sat\\_ts.dat](http://rammb.cira.colostate.edu/research/tropical_cyclones/ships/docs/AL/lsdiaga_1982_2017_sat_ts.dat) (a link to the 2018 version of the SHIPS developmental data.  
 3494 [last accessed on February 3, 2020.]  
 3495 SHIPS, 2018c:  
 3496 [http://rammb.cira.colostate.edu/research/tropical\\_cyclones/ships/docs/ships\\_predict](http://rammb.cira.colostate.edu/research/tropical_cyclones/ships/docs/ships_predict_or_file_2018.doc)  
 3497 [or\\_file\\_2018.doc](http://rammb.cira.colostate.edu/research/tropical_cyclones/ships/docs/ships_predict_or_file_2018.doc) (a link to the 2018 version of the SHIPS developmental data  
 3498 variables. [last accessed on February 3, 2020.]  
 3499 Simonyan, K., and A. Zisserman, 2014: Very deep convolutional networks for large-scale  
 3500 image recognition. *arXiv preprint arXiv:1409.1556*.  
 3501 Snoek, J., H. Larochelle, and R. P. Adams, 2012: Practical Bayesian optimization of  
 3502 machine learning algorithms. In *Advances in Neural Information Processing*  
 3503 *Systems*, 2951-2959.

3504 Song, J., X. Huang, S. Qin, and Q. Song, 2016: A bi-directional sampling based on K-  
 3505 means method for imbalance text classification. In *2016 IEEE/ACIS 15th*  
 3506 *International Conference on Computer and Information Science (ICIS)*, 1-5, IEEE.  
 3507 Stoer, J., and R. Bulirsch, 2013: Introduction to numerical analysis. *Courier Corporation*.  
 3508 Tallapragada, V., 2014: Performance of the 2013 NCEP Operational HWRF and Plans  
 3509 for 2014 Hurricane Season. Retrieved from  
 3510 [https://www.ofcm.noaa.gov/meetings/TCORF/ihc14/presentations/Session4/s04-](https://www.ofcm.noaa.gov/meetings/TCORF/ihc14/presentations/Session4/s04-02tallapragada.pdf)  
 3511 [02tallapragada.pdf](https://www.ofcm.noaa.gov/meetings/TCORF/ihc14/presentations/Session4/s04-02tallapragada.pdf)  
 3512 Tan, P.-N., M. Steinbach, and V. Kumar, 2015: Introduction to data mining. *Dorling*  
 3513 *Kindersley: Pearson*.  
 3514 Torres-Carrasquillo, P. A., E. Singer, M. A. Kohler, R. J. Greene, D. A. Reynolds, and J.  
 3515 R. Deller Jr, 2002: Approaches to language identification using Gaussian mixture  
 3516 models and shifted delta cepstral features. In *Seventh International Conference on*  
 3517 *Spoken Language Processing*  
 3518 Tran, D., L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, 2015: Learning  
 3519 spatiotemporal features with 3d convolutional networks. In *Proceedings of the*  
 3520 *IEEE international conference on computer vision*, 4489-4497.  
 3521 Trevor, H., T. Robert, and J. Friedman, 2009: *The Elements of Statistical Learning: Data*  
 3522 *Mining, Inference, and Prediction*. Springer Science & Business Media.  
 3523 Tsujino, S., and H. C. Kuo, 2020: Potential Vorticity Mixing and Rapid Intensification in  
 3524 the Numerically Simulated Supertyphoon Haiyan (2013). *Journal of the*  
 3525 *Atmospheric Sciences*, **77**(6), 2067-2090.

- 3526 Volinsky, C. T., and A. E. Raftery, 2000: Bayesian information criterion for censored  
3527 survival models. *Biometrics*, **56**(1), 256-262.
- 3528 Wang, Y., Y. Rao, Z. M. Tan, and D. Schönemann, 2015: A statistical analysis of the  
3529 effects of vertical wind shear on tropical cyclone intensity change over the western  
3530 North Pacific. *Monthly Weather Review*, **143**(9), 3434-3453.
- 3531 Wang, Z., 2018: What is the key feature of convection leading up to tropical cyclone  
3532 formation?. *Journal of the Atmospheric Sciences*, **75**(5), 1609-1629.
- 3533 Wei, Y., L. Sartore, J. Abernethy, D. Miller, K. Toppin, and M. Hyman, 2018: Deep  
3534 Learning for Data Imputation and Calibration Weighting. In *JSM Proceedings*,  
3535 Statistical Computing Section. Alexandria, VA: American Statistical Association,  
3536 1121-1131.
- 3537 Wilks, D. S., 2011: *Statistical Methods in the Atmospheric Sciences*. 3<sup>rd</sup> ed. Elsevier, 676.
- 3538 Yang, R., 2016: A Systematic Classification Investigation of Rapid Intensification of  
3539 Atlantic Tropical Cyclones with the SHIPS Database. *Weather and*  
3540 *Forecasting*, **31**(2), 495-513.
- 3541 Yang, R., J. Tan, and M. Kafatos, 2006: A Pattern Selection Algorithm in Kernel PCA  
3542 Applications. In *proceedings of the first International Conference on Software and*  
3543 *Data Technologies*, Setubal, Portugal, 195-202.
- 3544 Yang, R., J. Tang, and M. Kafatos, 2007: Improved associated conditions in rapid  
3545 intensifications of tropical cyclones. *Geophysical Research Letters*, **34**(20).

3546 Yang, R., J. Tang, and D. Sun, 2011: Association rule data mining applications for  
 3547 Atlantic tropical cyclone intensity changes. *Weather and Forecasting*, **26**(3), 337-  
 3548 353.  
 3549 Zeiler, M. D., D. Krishnan, G. W. Taylor, and R. Fergus, 2010: Deconvolutional  
 3550 networks. In *2010 IEEE Computer Society Conference on computer vision and*  
 3551 *pattern recognition, IEEE*, 2528-2535.  
 3552 Zeiler, M. D., G. W. Taylor, and R. Fergus, 2011: Adaptive deconvolutional networks for  
 3553 mid and high level feature learning. In *2011 International Conference on Computer*  
 3554 *Vision, IEEE*, 2018-2025.  
 3555 Zeiler, Matthew D, and R. Fergus, 2014: Visualizing and understanding convolutional  
 3556 neural networks. In *Proceedings of the 13th European Conference Computer*  
 3557 *Vision and Pattern Recognition, Zurich, Switzerland*, 6-12.  
 3558  
 3559  
 3560  
 3561  
 3562  
 3563  
 3564  
 3565  
 3566  
 3567  
 3568  
 3569  
 3570  
 3571  
 3572  
 3573  
 3574



3575

## **BIOGRAPHY**

3576 Yijun Wei graduated from Nanshan high School, Mianyang, Sichuan, China, in 2006. He  
3577 received his Bachelor of Arts from Sichuan University in 2011. He received his Master of  
3578 Arts in Math and Master of Science in Statistics from University of Michigan – Ann  
3579 Arbor in 2013.