

HONOURS PROJECT IN PHYSICS

A Study of Northeast Monsoon Forecast Errors

Author:

CHAN Man Yau (A0116571J)

Supervisors:

Professor LIM Hock

Dr. SUN Xiang Ming

Abstract

The forecast error covariance matrix, \mathbf{B} , is essential for data assimilation. However, \mathbf{B} is too big to be directly computed and stored. To overcome this problem, it is common to diagonalize \mathbf{B} through a two step process. The first step transforms the forecast errors, \mathbf{x}^b , into a set of errors that possess zero cross-covariances, resulting in a block diagonal \mathbf{B} containing only spatial auto-covariances. The second step then removes the spatial auto-covariances, resulting in a diagonal \mathbf{B} . This project focused on the first step and utilized linear relationships to decorrelate the stream function (ψ) and velocity potential (ϕ) forecast errors for the Singapore Variable resolution model, during the Maritime Continent Northeast monsoon season in 2015. However, the linear regression coefficient matrix relating ϕ and ψ is too big to be computed and must thus be modelled.

We will show that the two commonly used regression coefficient matrix models in the Weather Research and Forecasting Data Assimilation system do not perform well in the region and season considered. The first model assumes that ψ and ϕ are disconnected across model levels and are independent of the horizontal positions of ψ and ϕ considered. The second model disposes of the first assumption, but keeps the second assumption. Examinations of the columns of \mathbf{B} (estimated via the National Meteorology Centre method) revealed that cross-covariances depend on the horizontal positions of the ψ and ϕ , meaning that the linear regression coefficients depend heavily on the horizontal positions considered. This breaks down the common assumption of the two models. We also considered the efficacy of the models in removing cross-covariances, and found them to be unsurprisingly bad at cross-covariance removal.

Examinations of the columns of \mathbf{B} also revealed that the values of forecast error variables surrounding a point can be taken as samples of the forecast error variables at that point (neighbourhood equivalence assumption). Following this assumption, the forecast model grid is broken into subdomains and the regression coefficients are calculated between pairs of subdomains. This new model was found to be much better at removing cross-covariances because its assumption holds.

Acknowledgements

This project would not be possible without the guidance and mentorship of Professor Lim Hock from the National University of Singapore, Faculty of Science, Department of Physics, and Dr. Sun Xiangming from the Meteorological Services Singapore, Centre for Climate Research Singapore, Weather Modelling and Prediction. Their immense patience and understanding as I floundered through the theory, created buggy routines and occasionally overtaxing servers, are nothing short of admirable. I would also like to thank Dr. Huang Xiang-Yu and the staff from the Centre for Climate Research Singapore for their encouragement during this project and for letting me use their Athena computing cluster.

Aside from that, I would like to thank Soong Yun Ting and Liu Xiao Rui for being there when I raved and frothed at the mouth (metaphorically) while debugging routines and working out theoretical matters. Thank you Ng Wei Jie for letting me bounce ideas around. I would also like to thank Dr. Siu Zhuo Bin and John Ouyang for proofreading this thesis.

Lastly and most importantly, I would like to thank my parents, whose patience and love helped me to get through this project in one piece.

Contents

1	Introduction	1
1.1	Data assimilation background	1
1.2	Background error covariance matrix	3
1.2.1	Overwhelming size and unknown true state	3
1.2.2	Specificity to region and season	5
1.2.3	Dependence on source of \mathbf{X}^b	7
1.3	Overview	7
2	Materials and methods	15
2.1	Overcoming the \mathbf{X}^t requirement with the NMC method	15
2.2	SINGV model setup	17
2.3	GEN_BE processing	19
2.3.1	Pre-GEN_BE processing	19
2.3.2	GEN_BE	20
3	Sanity checking the SINGV outputs	23
3.1	Time-averaged thermodynamic profile is as expected	25
3.2	900 hPa monsoon flow is as expected	27
3.3	Cold surges manifested in the simulations	28
3.4	Interesting vortex features	28
3.5	Summary	30

4	Error covariance features and regression coefficient models	31
4.1	Physical meaning of ϕ and ψ	32
4.2	Columns of \mathbf{B}	34
4.2.1	Interpreting columns of \mathbf{B}	34
4.2.2	Computing columns of \mathbf{B} from NMC	35
4.3	Features observed in columns of \mathbf{B}	36
4.3.1	Deep convection	36
4.3.2	Vortex patterns	38
4.3.3	Horizontally variant and anisotropic cross-covariances	39
4.4	Regression coefficient models	42
4.4.1	Conventional regression coefficient models and their flawed assumptions	42
4.4.2	New regression coefficient matrix model	42
4.5	Summary	47
5	Results from regression coefficient models	49
5.1	Performance	49
5.2	Variance occupation	52
5.3	Summary	53
6	Conclusions and future work	55
A	The roles of the background error covariance matrix in DA	61
A.1	Solutions to Cost Function	61
A.1.1	Best linear unbiased estimate solution	61
A.1.2	Single and double-observation BLUE solutions	64
A.2	Roles of \mathbf{B} in data assimilation	65
A.3	Summary of derivation	67

B	Details about the first step of CVT	69
C	Additional plots	73

Chapter 1

Introduction

1.1 Data assimilation background

Need for data assimilation To obtain a good numerical weather forecast, the initial conditions must accurately represent the true state of the pre-forecast atmosphere (Caron and Fillion, 2010). Unfortunately, it is difficult in practice to obtain good initial conditions. Observations alone are insufficient to generate a full set of initial conditions — there are typically far fewer observations than model grid points (Bouttier and Courtier, 2002).¹ To make matters worse, while educated guesses (*e.g.*, result of a previous forecast) can easily generate a full set of initial conditions, they are susceptible to error. The solution is to combine the observations and educated guesses to produce a set of initial conditions that has the greatest probability of being accurate (Lorenc, 1986; Bannister, 2008a). This combination process is known as data assimilation, or DA.

DA problem To be exact, data assimilation seeks to determine a model state (state vector, \mathbf{X}) that has the maximum probability density of being true, given a set of observations (observation vector, \mathbf{Y}), and in the light of an educated guess (background state, \mathbf{X}^b). This probability density can be expressed using Bayes' theorem (Lorenc, 1986; Kalnay, 2003):

$$p(\mathbf{X}|\mathbf{Y}) \propto p(\mathbf{Y}|\mathbf{X})p(\mathbf{X}). \quad (1.1)$$

Contribution from \mathbf{X}^b The prior probability density, $p(\mathbf{X})$, represents the probability that \mathbf{X} is an accurate representation of the atmosphere, before considering the evidence.

¹Forecast models typically run on at least $O(10^6)$ grid points and $O(10)$ variables. *I.e.*, an unrealistic $O(10^7)$ observations are needed for initialization.

During this phase, \mathbf{X}^b is typically taken to be the best guess at what the true state \mathbf{X}^t looks like. As such, $p(\mathbf{X})$ is typically modelled by

$$p(\mathbf{X}) \propto \exp \left\{ -\frac{1}{2} (\mathbf{X} - \mathbf{X}^b)^\top \mathbf{B}^{-1} (\mathbf{X} - \mathbf{X}^b) \right\},$$

where \mathbf{B} is the background error covariance matrix (Lorenc, 1986). \mathbf{B} is defined as

$$\mathbf{B} \equiv \frac{1}{N_T} \sum_{t=1}^{N_T} (\mathbf{X}^b(t) - \mathbf{X}^t(t)) (\mathbf{X}^b(t) - \mathbf{X}^t(t))^\top \quad (1.2)$$

N_T is the number of samples considered, \mathbf{X}^t is the true state of the atmosphere and t is an index representing the sample number. For the ease of notation, let us define the background state error vector,

$$\mathbf{x}^b(t) \equiv \mathbf{X}^b(t) - \mathbf{X}^t(t). \quad (1.3)$$

This also means that

$$\mathbf{B} = \frac{1}{N_T} \sum_{t=1}^{N_T} \mathbf{x}^b(t) [\mathbf{x}^b(t)]^\top$$

Contribution from \mathbf{Y} The other term of Eqn (1.1), $p(\mathbf{Y}|\mathbf{X})$, is typically modelled by a similar Gaussian distribution, except that the reference (or the “mean”) is taken to be \mathbf{X} (Lorenc, 1986). Since \mathbf{X} exists in a different basis from \mathbf{Y} , an operator, H , that maps \mathbf{X} to \mathbf{Y} is needed. Taken together, the model is

$$p(\mathbf{Y}|\mathbf{X}) \propto \exp \left\{ -\frac{1}{2} [\mathbf{Y} - H(\mathbf{X})]^\top \mathbf{R}^{-1} [\mathbf{Y} - H(\mathbf{X})] \right\}.$$

The observation error covariance matrix, \mathbf{R} , is defined as (Bouttier and Courtier, 2002),

$$\mathbf{R} \equiv \frac{1}{N_T} \sum_{t=1}^{N_T} [\mathbf{Y}(t) - H(\mathbf{X}^t(t))] [\mathbf{Y}(t) - H(\mathbf{X}^t(t))]^\top.$$

Full DA problem Bringing these two probability densities into Eqn (1.1) yields:

$$p(\mathbf{X}|\mathbf{Y}) \propto \exp \left\{ -\frac{1}{2} (\mathbf{X} - \mathbf{X}^b)^\top \mathbf{B}^{-1} (\mathbf{X} - \mathbf{X}^b) - \frac{1}{2} [\mathbf{Y} - H(\mathbf{X})]^\top \mathbf{R}^{-1} [\mathbf{Y} - H(\mathbf{X})] \right\}.$$

DA cost function It is immediately apparent that maximizing $P(\mathbf{X}|\mathbf{Y})$ is equivalent to minimizing

$$J(\mathbf{X}) = (\mathbf{X} - \mathbf{X}^b)^\top \mathbf{B}^{-1} (\mathbf{X} - \mathbf{X}^b) + [\mathbf{Y} - H(\mathbf{X})]^\top \mathbf{R}^{-1} [\mathbf{Y} - H(\mathbf{X})]. \quad (1.4)$$

$J(\mathbf{X})$ is known in the literature as the cost function (Kalnay, 2003; Bannister, 2008a). In short, data assimilation determines the best initial conditions (\mathbf{X}^a) by minimizing Eqn (1.4), given a set of observations and an educated guess.²

1.2 Background error covariance matrix

Roles of \mathbf{B} in DA The background error covariance matrix (\mathbf{B}), as defined in Eqn (1.2), plays several important roles in data assimilation. First of all, \mathbf{B} weighs the contributions from \mathbf{Y} and \mathbf{X}^b to \mathbf{X}^a , according to their variances. The greater the relative variance of one component, the smaller its contribution to \mathbf{X}^a . Aside from that, \mathbf{B} is also responsible for spreading the difference in the observation and the background state throughout the model grid and across different variables. This means that \mathbf{B} allows a single observation to affect the entire set of initial conditions. Lastly, when multiple observations are assimilated, \mathbf{B} causes the observations to interact at locations/variables different from the observation location and variable (Bouttier and Courtier, 2002; Kalnay, 2003; Bannister, 2008a). Clearly, an accurate \mathbf{B} is essential for obtaining good initial conditions via data assimilation.³

1.2.1 Overwhelming size and unknown true state

Difficulties in obtaining \mathbf{B} Unfortunately, it is impossible to compute \mathbf{B} directly from the definition in Eqn (1.2) due to two problems. First of all the definition in Eqn (1.2) requires the unobtainable \mathbf{X}^t . This is typically handled by the National Meteorology Center method (Parrish and Derber, 1992), which will be covered in Chapter 2. Furthermore, since \mathbf{B} is an $N \times N$ matrix (N being the number of grid points times variables in the model), \mathbf{B} typically contains $O(10^{14})$ elements. In other words, computations involving \mathbf{B} are severely constrained by the impractical amount of computational memory required ($\sim O(10^2)$ terabytes). These mean direct computations of \mathbf{B} are clearly impractical (Bannister, 2008b).

²A solution to Eqn (1.4) can be derived under the assumption that the interpolation operator H can be linearized for small perturbations in \mathbf{x} . See Appendix A.1 for the details.

³All three roles of \mathbf{B} can be derived from Eqn (1.4). See Appendix A.2 for the details.

Diagonalizing \mathbf{B} by eigenvectors Transforming \mathbf{x}^b into a vector of mutually uncorrelated elements immediately diagonalizes \mathbf{B} (Lanczos, 1957; Bannister, 2008b). Since there are only N non-zero elements in the diagonalized \mathbf{B} , \mathbf{B} will not have computational memory problems. In principle, the ideal way of doing so is through the use of the eigenvectors of \mathbf{B} . When \mathbf{B} is transformed to the basis of its eigenvectors, it is diagonal. The transformation operator for \mathbf{x}^b then is simply an $N \times N$ matrix containing \mathbf{B} 's eigenvectors as its rows.

Replace eigenvector transform with control variable transform Unfortunately, it is difficult to compute the eigenvectors if \mathbf{B} cannot be computed in the first place. As such, the eigenvector transform method is replaced with a two-step transformation process (Parrish and Derber, 1992; Wu et al., 2002; Bannister, 2008b). Suppose that \mathbf{x}^b contains the background error of L variables. The first step transforms \mathbf{x}^b into a vector, $\mathbf{x}^{b'}$, where the cross-variable covariances (cross-covariance for short) are zero. The second step then removes the auto-covariances (covariance of the same variable, but between any pair of model grid points) of $\mathbf{x}^{b'}$, resulting in the vector $\mathbf{x}^{b''}$. Since $\mathbf{x}^{b''}$ will have zero cross-covariances and zero-autocovariances, its covariance matrix,

$$\langle \mathbf{x}^{b''} \mathbf{x}^{b''\top} \rangle = \frac{1}{N_T} \sum_{t=1}^{N_T} \mathbf{x}^{b''}(t) \mathbf{x}^{b''}(t)^\top,$$

will be diagonal. This two-step transformation is typically called the control variable transform (CVT)⁴.

Removing cross-covariances with least squares linear regression residues A good choice of variables for $\mathbf{x}^{b'}$ is the residues from performing least square linear regressions between pairs of \mathbf{x}^b variables (Daley, 1991). To see the transformation, suppose that \mathbf{x}^b is arranged into blocks containing only one variable each, or,

$$\mathbf{x}^{b\top} = \left[\mathbf{x}^{b_1\top}, \mathbf{x}^{b_2\top}, \dots, \mathbf{x}^{b_l\top}, \dots, \mathbf{x}^{b_{L_1}\top}, \mathbf{x}^{b_L\top} \right],$$

and $\mathbf{x}^{b'}$ is set up in a similar fashion,

$$\mathbf{x}^{b'\top} \equiv \left[\mathbf{x}^{b'_1\top}, \mathbf{x}^{b'_2\top}, \dots, \mathbf{x}^{b'_l\top}, \dots, \mathbf{x}^{b'_{L_1}\top}, \mathbf{x}^{b'_L\top} \right].$$

⁴A detailed description of the CVT method and the underlying mathematics of the first transformation can be found in Appendix B

The first transformation step is defined as

$$\mathbf{x}^{b'_l}(t) = \mathbf{x}^{b_l}(t) - \sum_{\ell=1}^{l-1} \boldsymbol{\alpha}_{l,\ell} \mathbf{x}^{b'_\ell}(t), \quad (1.5)$$

where,

$$\boldsymbol{\alpha}_{l,\ell} \equiv \left\langle \mathbf{x}^{b_l} \left(\mathbf{x}^{b'_\ell} \right)^\top \right\rangle \left\langle \mathbf{x}^{b'_\ell} \left(\mathbf{x}^{b'_\ell} \right)^\top \right\rangle^{-1}, \quad \ell = 1, 2, \dots, l-2, l-1. \quad (1.6)$$

In other words, $\mathbf{x}^{b'_l}(t)$ is the residue from performing a least square linear regression of \mathbf{x}^{b_l} against $\mathbf{x}^{b'_1}, \mathbf{x}^{b'_2}, \dots, \mathbf{x}^{b'_{l-1}}, \mathbf{x}^{b'_l}$. This transformation guarantees that $\left\langle \mathbf{x}^{b'_l} \left(\mathbf{x}^{b'_\ell} \right)^\top \right\rangle$ is zero for all $l \neq \ell$.⁵

Regression coefficient matrices must be modelled The CVT transform alone does not solve the overwhelming size problem – for L variables, there are $(L^2 - L)/2$ regression coefficient matrices, each with dimensions of $N/L \times N/L$. Furthermore, the auto-covariance matrices must be invertible to determine $\boldsymbol{\alpha}_{l,\ell}$. Unfortunately, there are typically insufficient samples of \mathbf{x}^b to ensure that the $N/L \times N/L$ auto-covariance matrices are invertible (Bannister, 2008b). These severely restrict our ability to determine $\boldsymbol{\alpha}_{l,\ell}$. However, it is possible to model $\boldsymbol{\alpha}_{l,\ell}$ by making some assumptions about how the regression coefficients vary spatially in the model grid. This modelling is the subject of this project.

1.2.2 Specificity to region and season

Regional and seasonal dependence The roles of \mathbf{B} described earlier clearly imply that \mathbf{B} contains information about how background errors are organized and connected. That in turn depends on the meteorological features of the region considered. These features also exhibit seasonal variations. As such, \mathbf{B} changes with the region and season of study. These two dependences must be considered in any study of \mathbf{B} (Bouttier and Courtier, 2002; Kalnay, 2003; Bannister, 2008b).

Paucity of Maritime Continent \mathbf{B} literature While there is a wealth of literature concerning \mathbf{B} in the mid-latitudes and on the global scale, scant attention has been given to \mathbf{B} at the Maritime Continent. However, there is ample motivation to study \mathbf{B} in this region. First of all, there is socioeconomic value in such a study: Indonesia, Malaysia and Singapore have a combined population of 85 million (Badan Perencanaan

⁵See the proof for Eqn (B.1) in Appendix B for the origin of this zero property.

Pembangunan Nasional, 2013). Furthermore, two of the most commonly used regression coefficient models explicitly assume that the regression coefficients do not change with the horizontal position of any pair of points considered (Chen et al., 2013; Descombes et al., 2015). In other words, the models assume that the regression coefficients are horizontally invariant and isotropic. Given the complex topography of the Maritime Continent (Chang et al., 2016), these models may not hold, necessitating a study of \mathbf{B} to test them and possibly formulate a new model. Unfortunately, to date, we have only been able to find a single publication concerning \mathbf{B} in this specific region (Chen et al., 2013). There is clearly a need for a study of \mathbf{B} in the Maritime Continent.

Paucity of Northeast monsoon \mathbf{B} literature To make matters worse, the above-mentioned singular study was only performed during the Southwest monsoon season (Chen et al., 2013). The climate of equatorial Southeast Asia has two monsoon seasons (Chang et al., 2005; Fong and Ng, 2012): the Southwest monsoon (July to September), the Northeast monsoon (December to March)⁶. In other words, a study \mathbf{B} in the Northeast monsoon season (NEM) of the Maritime Continent is warranted.

Northeast monsoon season The Northeast monsoon season (NEM) begins in December when the latitude of maximal solar heating is in the Maritime Continent-northern Australian region. This generates a pressure trough south of the Equator. The cooling of the northern Eurasian landmass prompts the formation of a surface, cold-core, high pressure system over Siberia (Ding, 1994). The pressure gradient drives the general near-surface northeasterly monsoon flow that characterizes the Northeast monsoon.

Cold surges Southward bursts of intense cold air have also been observed to emanate from this system when it moves towards the coastline of China and the western Pacific (Wang, 2006). These bursts, also known as cold surges, arrive at the Maritime Continent from the northeast, and amplify the near-surface northeasterlies (Fong and Ng, 2012) due to the orientation of the region’s topography (Chang et al., 2005). The topography of the Maritime Continent also acts to channel the cold surge towards the Equator (Chang et al., 2005). These cold surges are typically moistened by their trajectories over the South China Sea, and are associated with enhanced convection in the Maritime Continent (see Figure 6 of Chang et al., 2005).

⁶The two monsoon seasons are separated by intermonsoon seasons.

1.2.3 Dependence on source of \mathbf{X}^b

B varies with source of \mathbf{X}^b B is also dependent upon the source of educated guesses used for data assimilation. Consider the scenario where \mathbf{X}^b comes from the region's climatology and another scenario where \mathbf{X}^b is the outcome of a previous numerical forecast. Intuitively, the two \mathbf{X}^b will exhibit different errors, and thus have different B 's. Clearly, these dependences must be accounted for when studying B .

SINGV \mathbf{X}^b source Given this dependence, a source of \mathbf{X}^b must be selected in this project. The Singapore Variable resolution model (SINGV) is a forecasting model developed collaboratively by the Centre for Climate Research Singapore (Meteorological Services Singapore) and the UK MetOffice to enhance the forecasting capabilities of Singapore. To further enhance Singapore's forecasting capabilities, a study of B , using SINGV as its source, is needed.

1.3 Overview

Goal The goal of this project is to determine an appropriate model for the regression coefficient matrices involved in the CVT, over the Maritime Continent, during the Northeast monsoon season, for the Singapore Variable resolution model.

Wind field B As a first approach to this topic, we will only consider 2 forecast error variables: the velocity potential forecast error, ϕ , and the stream function forecast error, ψ , of the horizontal wind components. Their corresponding forecast variables are denoted by the corresponding capitalized Greek letters, Φ and Ψ . The two variables, and their errors, are related to the horizontal wind field (Bijlsma et al., 1986), and its error, via

$$\mathbf{U}_H = \nabla_H \Phi - \nabla_H \times (\Psi \hat{\mathbf{z}}), \quad (1.7)$$

$$\mathbf{u}_H = \nabla_H \phi - \nabla_H \times (\psi \hat{\mathbf{z}}), \quad (1.8)$$

where

$$\nabla_H \equiv \hat{\mathbf{x}} \frac{\partial}{\partial x} + \hat{\mathbf{y}} \frac{\partial}{\partial y}.$$

The gradient term and curl terms can be viewed as the diverging part of the wind (wind error) and rotating part of the wind (wind error) respectively.

CVT formulation As such, the state vector of forecast errors at time t is

$$\mathbf{x}^b_t = [\boldsymbol{\psi}_t^\top, \boldsymbol{\phi}_t^\top]^\top, \quad (1.9)$$

where $\boldsymbol{\psi}_t$ is a vector of $N/2$ elements containing the forecast errors of ψ at all positions on the forecast model grid at time t , and likewise for $\boldsymbol{\phi}$. The corresponding control variable transform for this study is

$$\begin{cases} \boldsymbol{\psi} = \boldsymbol{\psi}' \\ \boldsymbol{\phi}' = \boldsymbol{\phi} - \boldsymbol{\alpha}_{\phi,\psi}\boldsymbol{\psi}' \end{cases} \quad (1.10)$$

Two commonly used $\boldsymbol{\alpha}_{\phi,\psi}$ models A total of three ways to model $\boldsymbol{\alpha}_{\phi,\psi}$ will be considered in this project. The first model assumes that ψ and ϕ are only related when they are on the same model level, and that their relationship is independent of the horizontal positions considered. The second model removes the same-level restriction, allowing for inter-model level relationships. However, the horizontal invariance assumption applies to the second model as well. The algorithms for these models are illustrated in Figures 1.1 and 1.2.

New $\boldsymbol{\alpha}_{\phi,\psi}$ model The third model will follow a completely different assumption: the evaluations of ϕ and ψ in the vicinity of a selected model point are equivalent to samples of the evaluations of ϕ and ψ at the said model point. The resulting third model can be thought of as a coarse-resolution version of \mathbf{B} , allowing both horizontal and vertical variations in the relationship between ψ and ϕ . The way which regression coefficients are computed and employed in this model are illustrated in Figures 1.3 and 1.4, respectively.

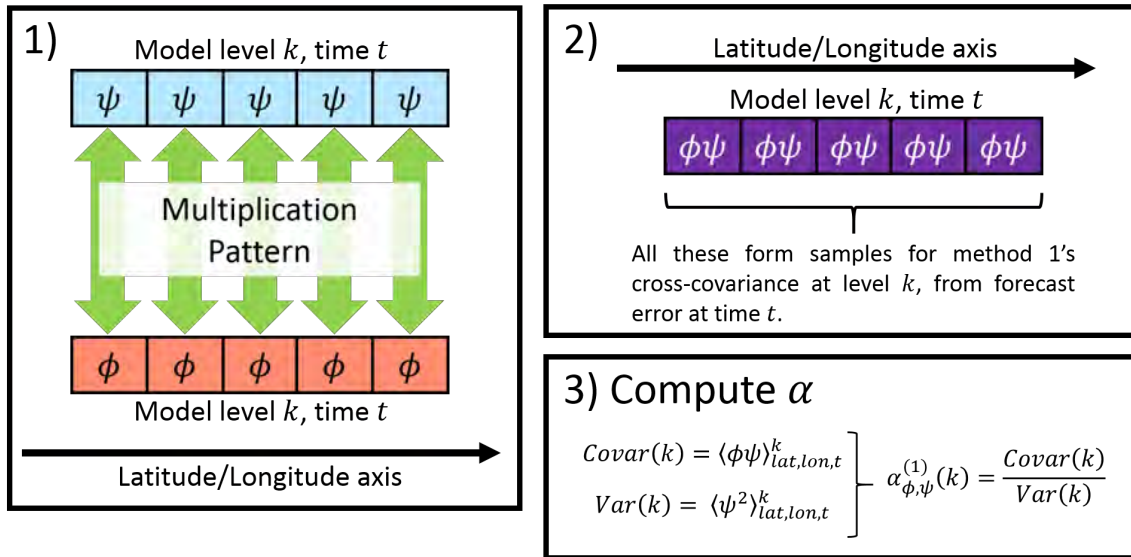
Overview – Setup and checks We will begin with explaining how \mathbf{x}^b can be determined by using pairs of forecasts, via the National Meteorology Center (NMC) method, in Chapter 2. Following that, Chapter 2 will detail the process of SINGV data generation and subsequent processing to produce \mathbf{x}^b . The SINGV forecast outputs will then be checked to ensure that it is capturing the general characteristics of the Northeast monsoon in Chapter 3.

Overview – modelling $\boldsymbol{\alpha}_{\phi\psi}$ Upon establishing the sanity of the SINGV model, we will then examine some of the forecast error covariance features estimated by the NMC method in Chapter 4⁷. The negative implications of the observed features for the two commonly used models (shown in Figures 1.1 and 1.2) will then be elaborated. In response, Chapter

⁷Essentially, we will be looking at columns of a \mathbf{B} estimated by the NMC.

4 will formulate new way of modelling regression coefficient matrices (the third model illustrated in Figures 1.3 and 1.4). Chapter 5 will then compare the performance of two commonly used models against that of the new model, and show that the new model far outperforms the former two. We will then conclude with some comments and directions for future work in Chapter 6.

Method 1: $\alpha_{\phi,\psi}^{(1)}$ computation



Method 1: using $\alpha_{\phi,\psi}^{(1)}$ for control variable transform

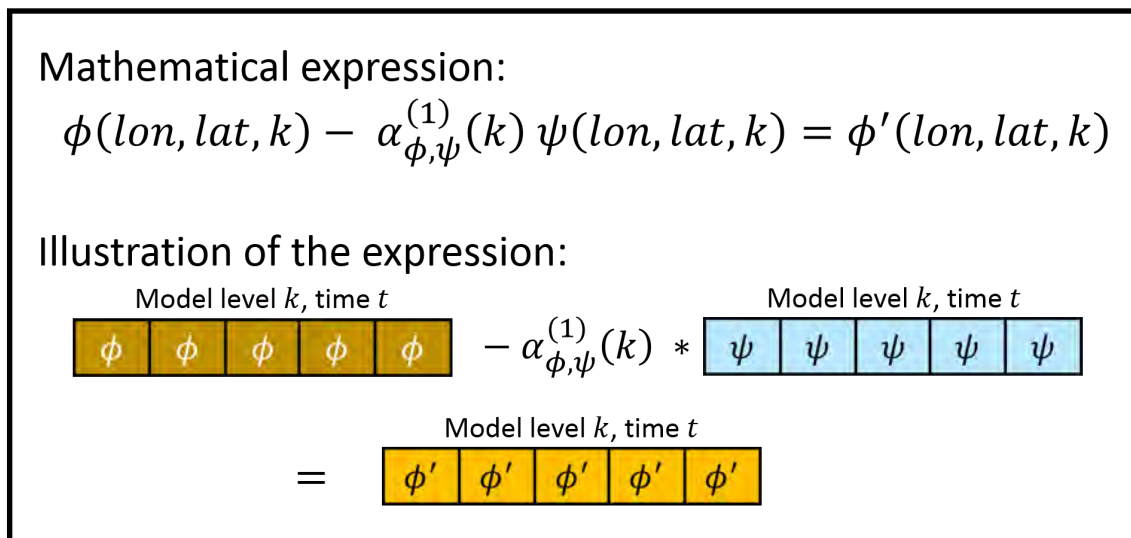
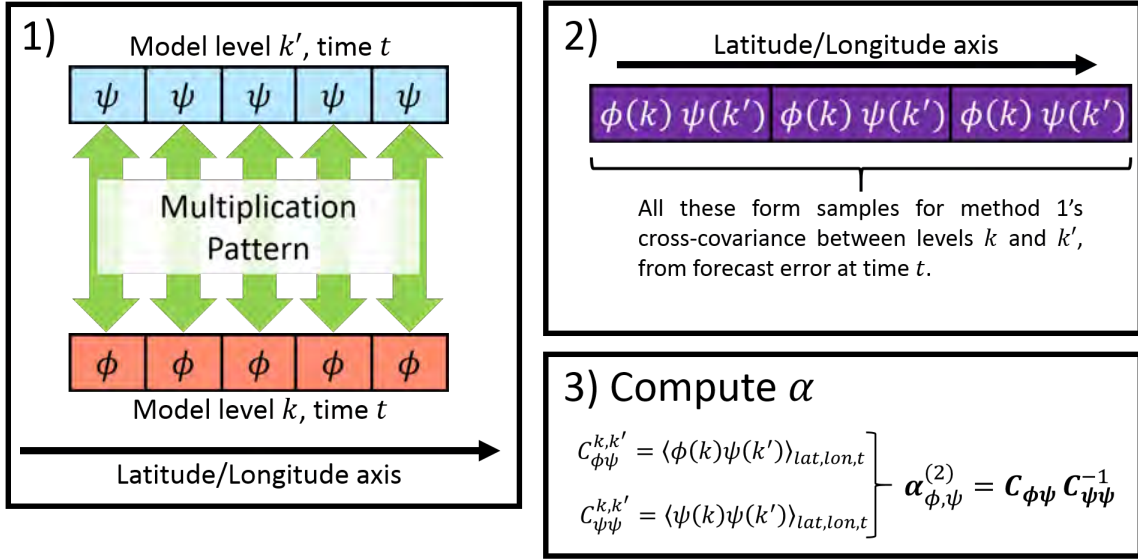


Figure 1.1: Illustration of how the first model computes regression coefficients (top) and how it is utilized for the control variable transform (bottom). Each colored box represents a model grid point with the indicated error variable indicated. The panels are to be read in sequence of the labelled numbers. This model assumes that ψ and ϕ are disconnected across model levels and their relationship on the same model level is invariant of horizontal position. In other words, there are only N_Z regression coefficients, one for each model level. Note that lon , lat and k refer to the longitude, latitude and model level of the model grid. The subscripts appended to the averages indicate the dimensions of averaging.

Method 2: $\alpha_{\phi,\psi}^{(2)}$ computation



Method 2: using $\alpha_{\phi,\psi}^{(2)}$ for control variable transform

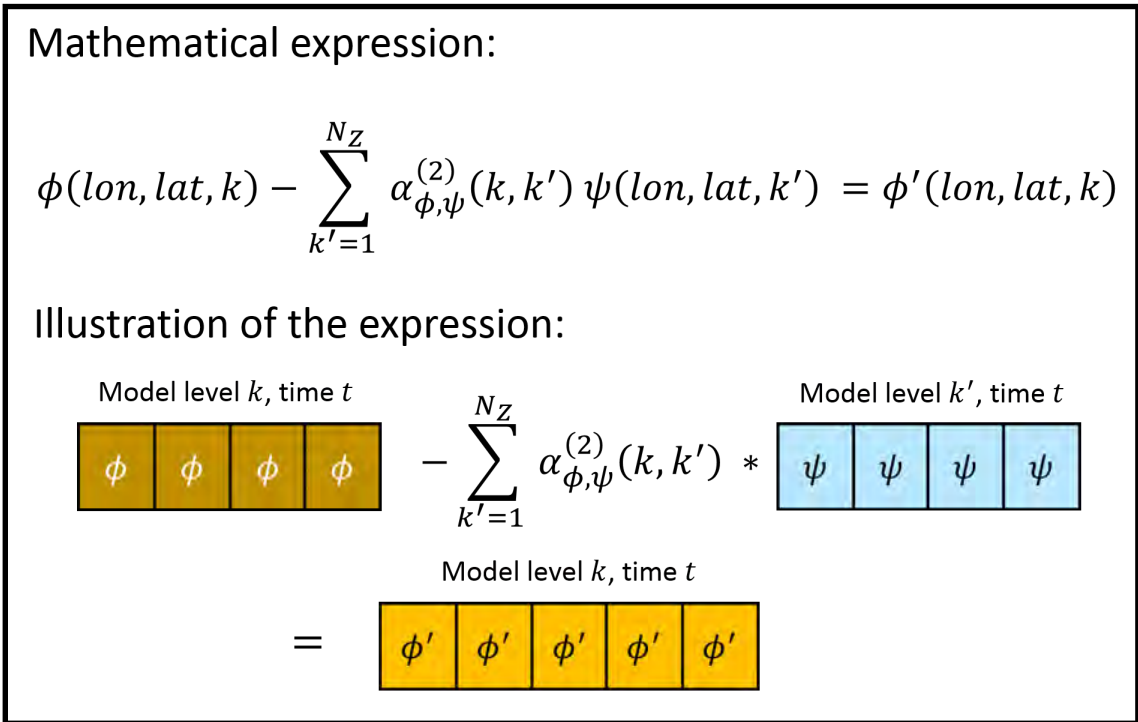


Figure 1.2: Illustration of how the second model computes regression coefficients (top) and how it is utilized for the control variable transform (bottom). Each colored box represents a model grid point with the indicated error variable indicated. The panels are to be read in sequence of the labelled numbers. This model assumes that the relationship between ψ and ϕ is invariant of horizontal position. In other words, there are only N_Z regression coefficients, one for each model level. Note that lon , lat and k refer to the longitude, latitude and model level of the model grid. The subscripts appended to the averages indicate the dimensions of averaging.

Method 3: $\alpha_{\phi,\psi}^{(3)}$ computation

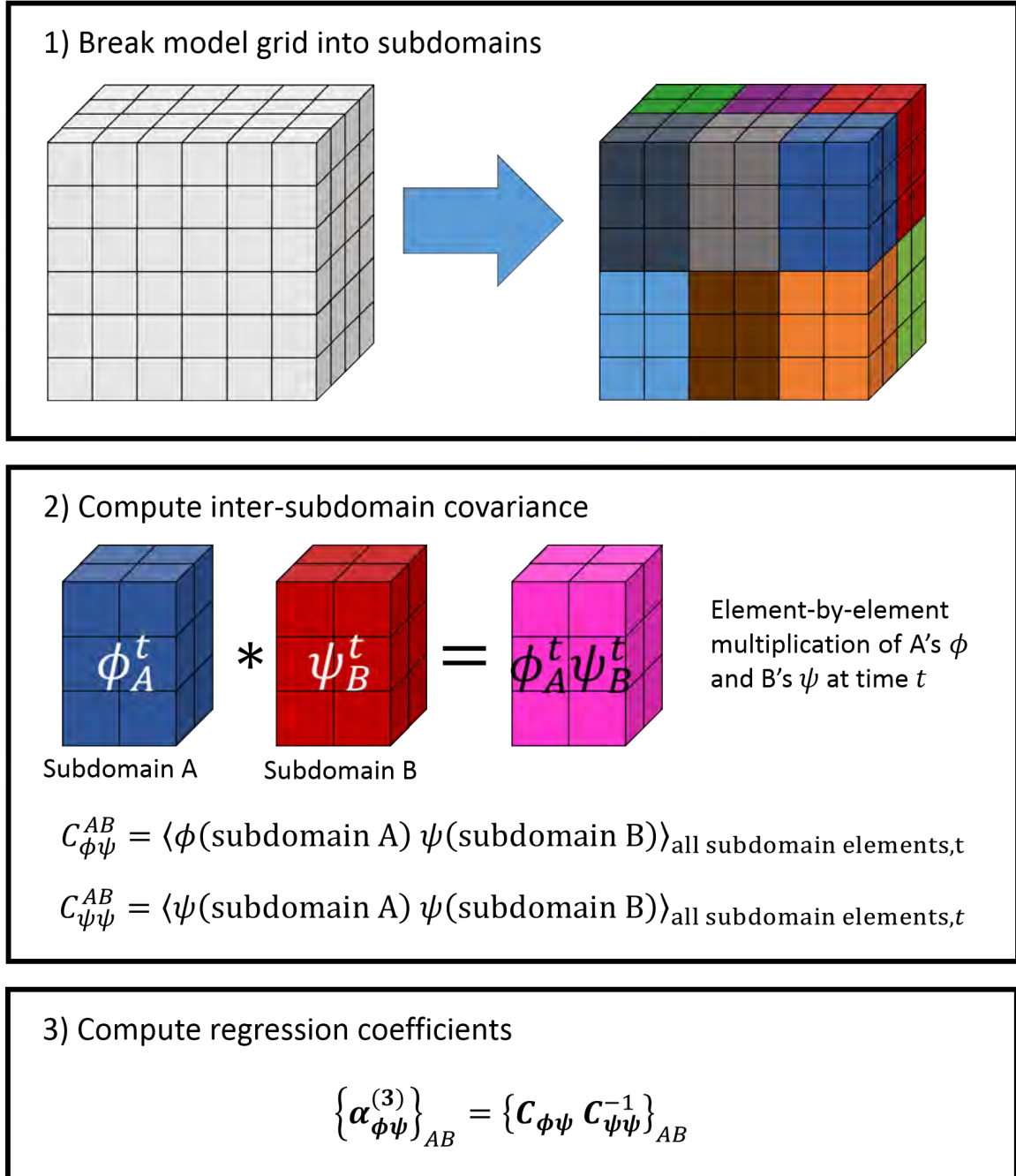


Figure 1.3: Illustration of how the third model computes regression coefficients. The panels are to be read in sequence of the labelled numbers. Each colored box represents a model grid point with the indicated error variable indicated. The panels are to be read in sequence of the labelled numbers. This model assumes that within a subdomain, all values of ψ and ϕ can be treated as samples of the center of the subdomain. The subscripts appended to the averages indicate the dimensions of averaging.

Method 3: using $\alpha_{\phi,\psi}^{(3)}$ for control variable transform

Mathematical expression:

$$\phi(\text{subdomain}_A, t) - \sum_B^{\text{all subdomains}} \alpha_{\phi,\psi}^{(3)}(A, B) \psi(\text{subdomain}_B, t) = \phi'(\text{subdomain}_A, t)$$

Illustration of the expression:

The diagram illustrates the mathematical expression using 3D cubes. On the left, a blue cube represents Subdomain A with the variable ϕ_A^t . This is followed by a minus sign and a summation over all subdomains B, with the coefficient $\alpha_{\phi,\psi}^{(3)}(A, B)$. This is multiplied by a red cube representing Subdomain B with the variable ψ_B^t . The result is a green cube representing Subdomain A with the transformed variable ϕ'_A^t .

Figure 1.4: Illustration of how the third model is used in the control variable transform. The panels are to be read in sequence of the labelled numbers. Each colored box represents a model grid point with the indicated error variable indicated. The panels are to be read in sequence of the labelled numbers. This model assumes that within a subdomain, all values of ψ and ϕ can be treated as samples of the center of the subdomain.

Chapter 2

Materials and methods

Overview In this chapter, we will explain how \mathbf{x}^b can be estimated using the National Meteorology Center method, in spite of the unknown \mathbf{X}^t . Following that, the setup of the Singapore Variable resolution model (SINGV) will be discussed. During that discussion, the exact manner which the National Meteorology Centre method is applied will also be elaborated. To obtain the velocity potential and stream functions, we utilized the GEN_BE package created by the National Corporation for Atmospheric Research and the UK Met Office.

2.1 Overcoming the \mathbf{X}^t requirement with the NMC method

Purpose The National Meteorology Centre (NMC) method mitigates the requirement of the unobtainable \mathbf{X}^t in any computation of \mathbf{B} (Parrish and Derber, 1992). In essence, if two different \mathbf{X}^b from the same source can be obtained for the same time, t , \mathbf{B} can be approximated by

$$\mathbf{B} \approx \frac{1}{2N_T} \sum_{t=1}^{N_T} (\mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t} - \boldsymbol{\mu}) (\mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t} - \boldsymbol{\mu})^\top, \quad (2.1)$$

where,

$$\boldsymbol{\mu} \equiv \frac{1}{N_T} \sum_{t=1}^{N_T} (\mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t}).$$

Assumptions of NMC method The NMC method operates under two assumptions. First of all, it is assumed that the auto-covariance of the two background errors are equivalent on average. In other words,

$$\sum_{t=1}^{N_T} \left(\mathbf{x}^b_{(1),t} \mathbf{x}^{b\top}_{(1),t} \right) \approx \sum_{t=1}^{N_T} \left(\mathbf{x}^b_{(2),t} \mathbf{x}^{b\top}_{(2),t} \right),$$

where $\mathbf{x}^b_{(1),t}$ and $\mathbf{x}^b_{(2),t}$ are the errors of $\mathbf{X}^b_{(1),t}$ and $\mathbf{X}^b_{(2),t}$ respectively. The second assumption is that $\mathbf{x}^b_{(1),t}$ and $\mathbf{x}^b_{(2),t}$ are mutually uncorrelated, or,

$$\sum_{t=1}^{N_T} \left(\mathbf{x}^b_{(1),t} \mathbf{x}^{b\top}_{(2),t} \right) = \sum_{t=1}^{N_T} \left(\mathbf{x}^b_{(2),t} \mathbf{x}^{b\top}_{(1),t} \right) = \mathbf{0}.$$

NMC derivation The NMC method in Eqn (2.1) can be shown to be consistent with the definition of \mathbf{B} in Eqn (1.2) under these assumptions. Starting from Eqn (2.1), and drawing on the fact that $\mathbf{x}^b_{(1),t} \equiv \mathbf{X}^b_{(1),t} - \mathbf{X}^t_t$ and likewise for $\mathbf{x}^b_{(2),t}$,

$$\begin{aligned} & \frac{1}{N_T} \sum_{t=1}^{N-1} \left\{ \mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t} \right\} \left\{ \mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t} \right\}^\top \\ &= \frac{1}{N_T} \sum_{t=1}^{N-1} \left\{ (\mathbf{x}^b_{(1),t} + \mathbf{X}^t_t - \mathbf{x}^b_{(2),t} - \mathbf{X}^t_t) (\mathbf{x}^b_{(1),t} + \mathbf{X}^t_t - \mathbf{x}^b_{(2),t} - \mathbf{X}^t_t)^\top \right\} \\ &= \frac{1}{N_T} \sum_{t=1}^{N-1} \left\{ (\mathbf{x}^b_{(1),t} - \mathbf{x}^b_{(2),t}) (\mathbf{x}^b_{(1),t} - \mathbf{x}^b_{(2),t})^\top \right\} \\ &= \frac{1}{N_T} \sum_{t=1}^{N-1} \left\{ \mathbf{x}^b_{(1),t} \mathbf{x}^{b\top}_{(1),t} - \mathbf{x}^b_{(1),t} \mathbf{x}^{b\top}_{(2),t} - \mathbf{x}^b_{(2),t} \mathbf{x}^{b\top}_{(1),t} + \mathbf{x}^b_{(2),t} \mathbf{x}^{b\top}_{(2),t} \right\} \\ &\approx \frac{1}{N_T} \sum_{t=1}^{N-1} \left\{ \mathbf{x}^b_{(1),t} \mathbf{x}^{b\top}_{(1),t} + \mathbf{x}^b_{(2),t} \mathbf{x}^{b\top}_{(2),t} \right\} \approx 2 \sum_{t=1}^{N-1} \left\{ \mathbf{x}^b_{(1),t} \mathbf{x}^{b\top}_{(1),t} \right\} \approx 2\mathbf{B}. \end{aligned}$$

Modification due to bias There is a hidden pitfall in using the NMC method: the average of $\mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t}$ may not be zero. As the general definition of a covariance matrix calls for mean removal before computing the covariance, the NMC method must be modified. The NMC method employed in this project is thus (Descombes et al., 2015):

$$\mathbf{B} \approx \frac{1}{2N_T} \sum_{t=1}^{N_T} (\mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t} - \boldsymbol{\mu}) (\mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t} - \boldsymbol{\mu})^\top,$$

where,

$$\boldsymbol{\mu} \equiv \frac{1}{N_T} \sum_{t=1}^{N_T} (\mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t}).$$

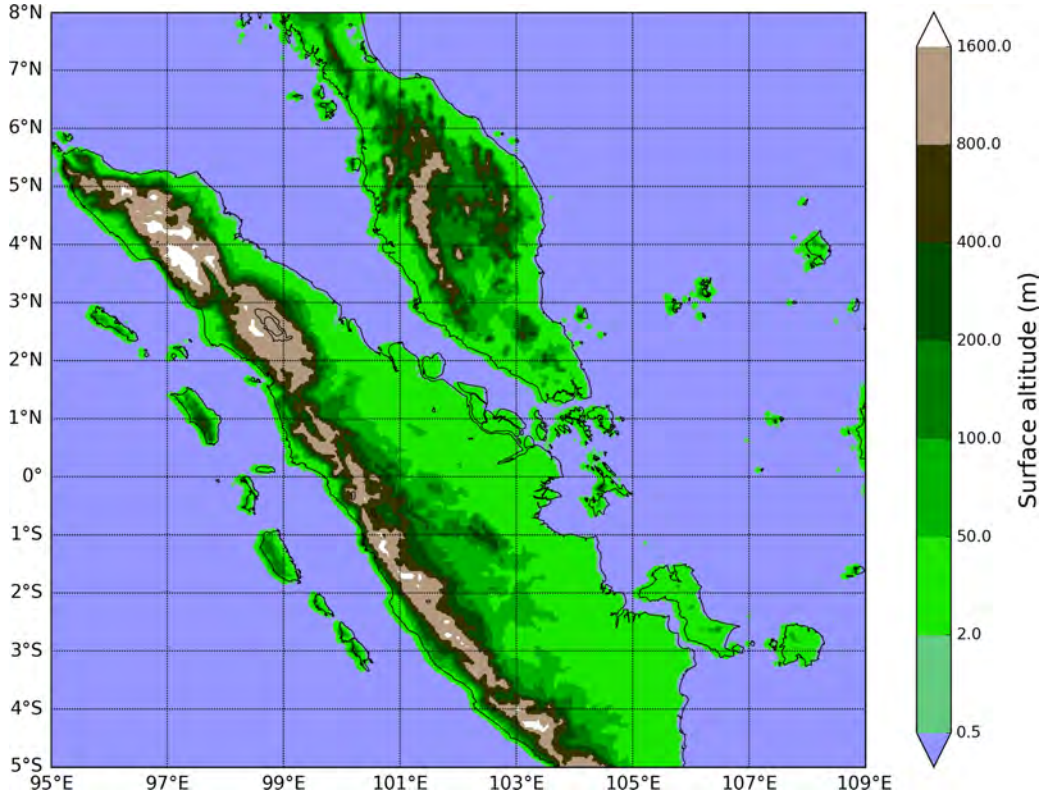


Figure 2.1: Plot of the surface heights of the SINGV simulation domain.

This also means that the background error at time t , \mathbf{x}^b_t , can be viewed as:

$$\mathbf{x}^b_t \approx \frac{1}{\sqrt{2}} (\mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t} - \langle \mathbf{X}^b_{(1),t} - \mathbf{X}^b_{(2),t} \rangle).$$

Forecast perturbation-based NMC Since the source of \mathbf{X}^b in this project is the Singapore Variable resolution model forecasts (SINGV), the NMC requires pairs of forecasts valid at the same time to form \mathbf{B} (Parrish and Derber, 1992; Bannister, 2008b). If the time of data assimilation is T hours, then SINGV forecasts were initiated 12 hours and 6 hours before T , and their outputs at T are used as as the two \mathbf{X}^b at time T .

2.2 SINGV model setup

Background and computational resource This study employed the Singapore Variable Resolution Model, version 2.1 (SINGV), which is modified from the UK MetOffice’s (UKMO) Unified Model (UM). The SINGV v2.1 is produced by a multi-year collaboration project between the UK MetOffice and the Center for Climate Research Singapore to build a tropical convective-scale NWP for Singapore and the surrounding region. The model was run using the Center for Climate Research Singapore’s High Performance Com-

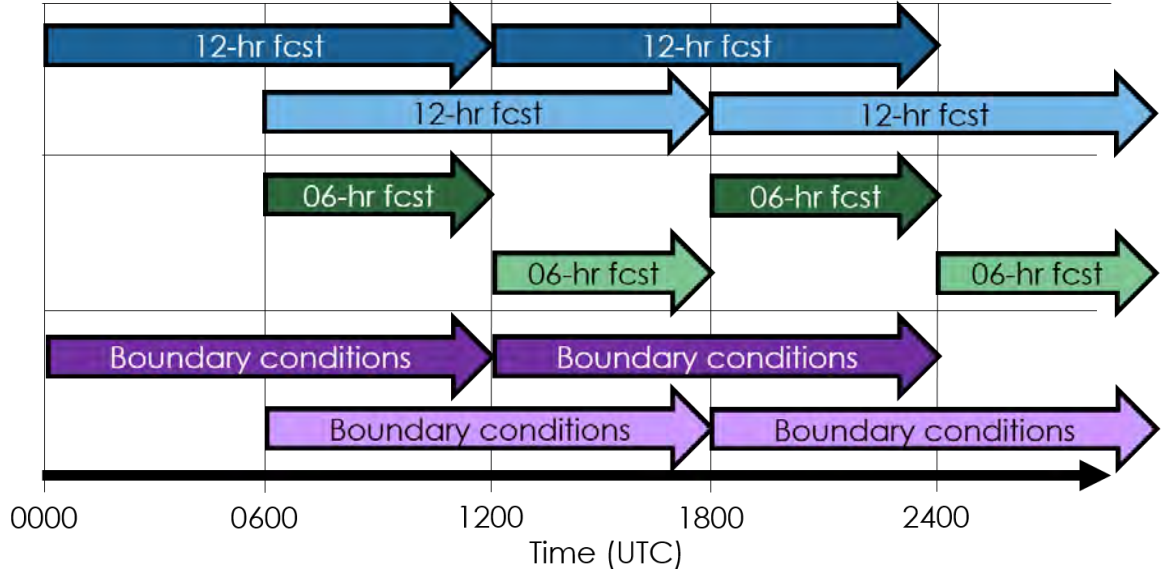


Figure 2.2: Illustration of how pairs of forecasts are set up, to be utilized by the NMC method. See text for the complementing description.

puting Cluster, Athena. Forecasts from the SINGV model were obtained over 15 days of the Northeast monsoon of 2015, from 5th December 2015, 0000 UTC, to 19th December 2015, 1800 UTC.

Model grid setup The simulation domain spans over a geographical region from 94.543°E to 109.2445°E and 5.74°S to 8.2325°N, with a horizontal spatial resolution of 4.5 km¹. The model employed a total of 80 static, terrain-following, hybrid height model levels. The orography of the simulation domain has been plotted out in Figure 2.1.

Employing the NMC method As mentioned earlier, the NMC method can be used to sidestep the \mathbf{X}^t requirement in computing \mathbf{B} . Pairs of \mathbf{X}^b valid at the same times are needed for that approximation. In the case of SINGV, pairs of 12-hour and 6-hour forecasts were used. In other words, the NMC estimate of \mathbf{x}^b_t is

$$\mathbf{x}^b_t \approx \frac{1}{\sqrt{2}} (\mathbf{X}^b_{12\text{-hr},t} - \mathbf{X}^b_{06\text{-hr},t} - \boldsymbol{\mu}), \quad \boldsymbol{\mu} \equiv \frac{1}{N_T} \sum_{t=1}^{N_T} (\mathbf{X}^b_{12\text{-hr},t} - \mathbf{X}^b_{06\text{-hr},t}) \quad (2.2)$$

We will use Figure 2.2 to illustrate the how the forecast pairs were set up. Suppose that a 12-hour forecast was initiated at 0000 UTC. Then to obtain the second forecast for the NMC, a 06-hour forecast was initiated at 0600 UTC. The initial conditions used for both forecasts were taken from the UM Global analysis data, released at 0000 UTC and 0600 UTC respectively. Both forecasts utilized the same set of hourly boundaries obtained

¹Horizontal grid arrangement: 364 (longitudinal) by 346 (latitudinal).

from the 36-hour UM Global that was released at 0000 UTC. The pair of forecasts at 1200 UTC forms the first forecast pair for the NMC. This process is repeated at 0600 UTC, 1200 UTC, 1800 UTC, and so forth, every 6 hours. Note that this method of matching boundary conditions is known to enhance the accuracy of the NMC method (Berre et al., 2006; Bannister, 2008a).

Initial and boundary conditions The boundary and initial conditions for SINGV were taken from the UM Global 36-hour global forecasts. These global forecasts are performed daily at 00, 06, 12 and 18 hours UTC and are released to the Center for Climate Research Singapore upon completion. The SINGV system then performs a vertical linear and horizontal bilinear interpolation to pass the UM Global data into the aforementioned SINGV grid.

2.3 GEN_BE processing

2.3.1 Pre-GEN_BE processing

Need for pre-GEN_BE processing To use the GEN_BE package, the forecasts must be converted from the UK Met Office PP files into the Weather Research and Forecasting model NetCDF output files. Aside from the file format difference, the SINGV model uses a height-based, terrain following, vertical coordinate system, whereas the WRF model uses a pressure-based, terrain following, vertical coordinate system (η levels). The η level value of a point with pressure P^{WRF} with a model top pressure P^{top} and model surface pressure P^{sfc} is:

$$\eta = \frac{P^{\text{WRF}} - P^{\text{t}}}{P^{\text{sfc}} - P^{\text{t}}}$$

Hybrid height to η -level transform The WRF Preprocessing system is capable of transforming fields from hybrid height coordinates to η levels via vertical linear interpolation (metgrid step). To maximize interpolation accuracy, η levels that minimize the difference between η level pressure and hybrid height pressure were used. *I.e.*,

$$\frac{\partial}{\partial \eta_k} (P_{ijk}^{\text{SINGV}} - P_{ijk}^{\text{WRF}})^2 = 0 \quad (2.3)$$

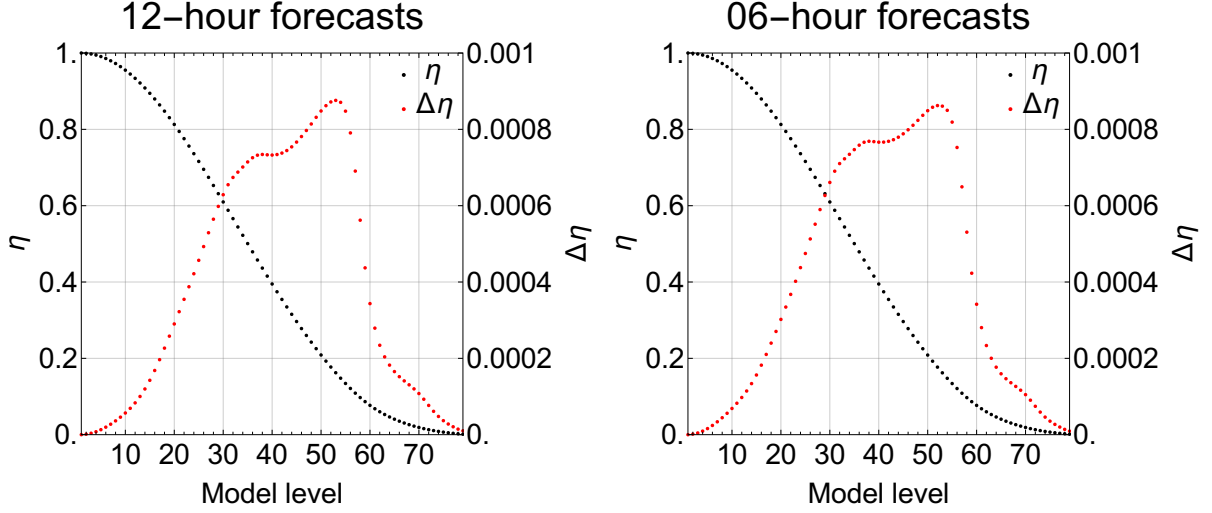


Figure 2.3: Plots of the average and standard deviation of η with respect to each model height, for the 12-hour forecast (left) and the 6-hour forecast (right).

Since $\eta = \frac{P^{\text{WRF}} - P^{\text{top}}}{P^{\text{sfc}} - P^{\text{top}}}$ implies $P_{ijk}^{\text{WRF}} = \eta_k (P_{ij}^{\text{sfc}} - P_{ij}^{\text{top}}) + P_{ij}^{\text{top}}$, the set of η values that solves Eqn (2.3) is thus found to be:

$$\eta_k = \frac{P_{ij}^{\text{top}} P_{ij}^{\text{top}} - P_{ij}^{\text{top}} P_{ijk}^{\text{SINGV}} - P_{ij}^{\text{sfc}} P_{ij}^{\text{top}} + P_{ij}^{\text{sfc}} P_{ijk}^{\text{SINGV}}}{P_{ij}^{\text{top}} P_{ij}^{\text{top}} - 2P_{ij}^{\text{top}} P_{ij}^{\text{sfc}} + P_{ij}^{\text{sfc}} P_{ij}^{\text{sfc}}} \quad (2.4)$$

η -levels are similar across 54 forecasts Surprisingly, despite the fact that Eqn (2.4) does not guarantee that the η values of each forecast set will be different, the η values ultimately turned out to be virtually identical. The average η value, and the standard deviation of η , across model levels, for both the 12-hour forecast (unstaggered) and the 6-hour forecast (staggered) are as shown in Figure 2.3.

ARW-WRF conversion After the η level and WPS procedure, the WRF model was run for 0 time steps to produce the SINGV data in the appropriate NetCDF format with the WRF model variables.

2.3.2 GEN_BE

Uses of GEN_BE The GEN_BE package was utilized on the WRF-processed forecast pairs to produce the velocity potential (Φ) and stream function (Ψ) for each pair. GEN_BE then follows Eqn (2.2) to produce the NMC-estimated forecast error.

Φ and Ψ determination As mentioned in the introduction, Φ and Ψ are defined for the horizontal wind field. In other words, for the k -th model level, Φ and Ψ are related to the horizontal wind field at level k , \mathbf{U}_k , by (Descombes et al., 2015)

$$\nabla_{\mathbf{H}} \cdot \mathbf{U}_k(x, y) = \nabla_H^2 \Phi(x, y, k) \quad , \quad \{\nabla_{\mathbf{H}} \times \mathbf{U}_k(x, y)\} \cdot \hat{\mathbf{z}} = \nabla_H^2 \Psi(x, y, k) \quad (2.5)$$

where

$$\mathbf{U}_k(x, y) \equiv U_{x,k}(x, y) \hat{\mathbf{x}} + U_{y,k}(x, y) \hat{\mathbf{y}} \quad , \quad \nabla_{\mathbf{H}} \equiv \hat{\mathbf{x}} \frac{\partial}{\partial x} + \hat{\mathbf{y}} \frac{\partial}{\partial y} \quad ,$$

and $\hat{\mathbf{x}}$, $\hat{\mathbf{y}}$ and $\hat{\mathbf{z}}$ are the usual right-handed Cartesian unit vectors. The GEN_BE package solves Eqn (2.5) using a discrete sine transform, with zero boundary conditions. It then utilizes Eqn (2.2) to determine the forecast errors at time t , \mathbf{x}^b_t .

Removal surface level and boundaries The GEN_BE removes the first model level during its processing. Also, to prevent the boundary conditions of SINGV coming from the UM from interfering with the analysis, 18 points are cropped from both the north and south boundaries of the model grid and 17 points are cropped from both the east and west boundaries. This turns the model grid from a $364 \times 346 \times 80$ cuboid lattices of points (in longitude, latitude, vertical order) to $330 \times 310 \times 79$.

Chapter 3

Sanity checking the SINGV outputs

Purpose of chapter Before examining the features of the forecast error, it is important to ensure that the SINGV outputs display the general features of the Northeast monsoon. This acts as a sanity check to ensure that SINGV is behaving as expected. As such, we will now examine the synoptic features of the SINGV data set.

Overview This chapter will start with examining and confirming that the time-averaged thermodynamic and horizontal wind fields of the SINGV forecasts are as expected of the Northeast monsoon season. Afterwards, the forecasts will be examined for cold surges (important feature in the Northeast monsoon) as a further confirmation that the model is capturing the key Northeast monsoon features. We will then comment on an intriguing, short-lived, low-level, counter-clockwise circulation feature observed in the forecasts.

Hybrid height model levels will be used Before discussing the features of the SINGV dataset, note that the data in this section will be plotted based on the hybrid height model levels used in the SINGV model. However, the model level index will not be stated: instead, the rough atmospheric pressure level that corresponds to that model level will be stated. The correspondence between the model level and atmospheric pressure, and between model level and geopotential height, are as plotted in Figures 3.1 and 3.2.

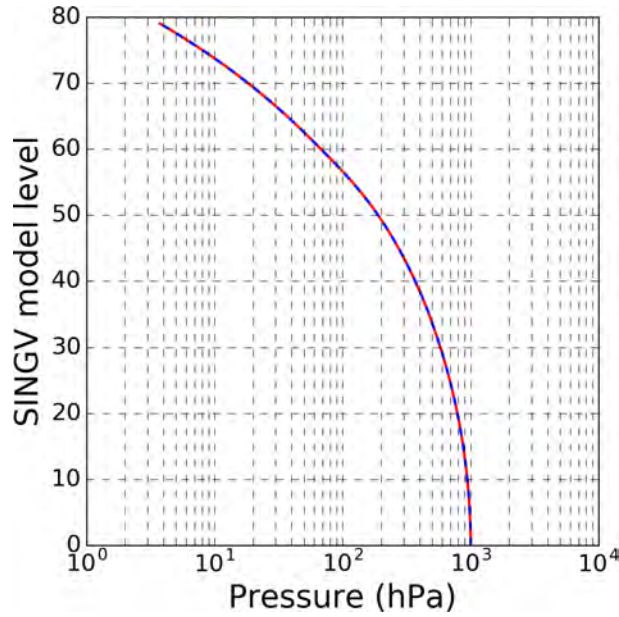


Figure 3.1: Plots of temporally and horizontally averaged pressure with respect to SINGV model levels. The solid red and dashed blue lines respectively indicate the averaged pressure profile of the 12-hour and 06-hour forecasts. Note that we have actually shaded one standard deviation on the left and right of pressure curve (see Figure 3.2 for example). However, the standard deviations are so small that they cannot be distinguished from the average curves. Also, both pressure profiles are virtually identical.

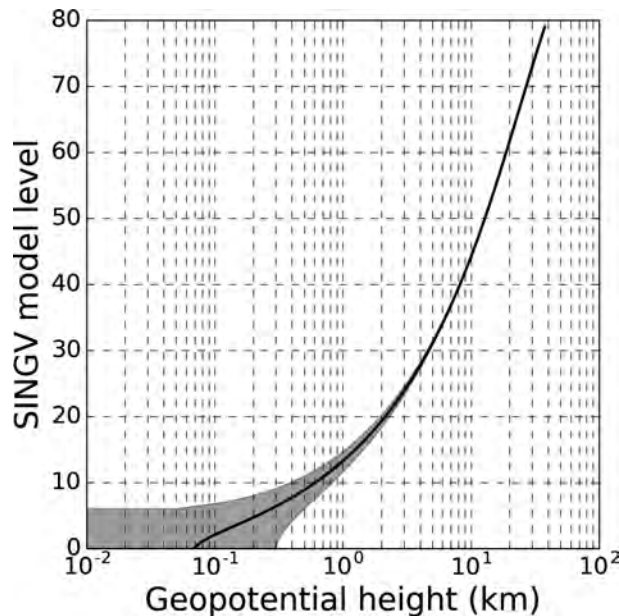


Figure 3.2: Plot of horizontally averaged geopotential height with respect to SINGV model levels. The shaded grey area indicates one standard deviation on either side of the averaged curve, per model level. Note that only one geopotential height curve is needed as the geopotential height of each model level remains the same across all forecasts in this study.

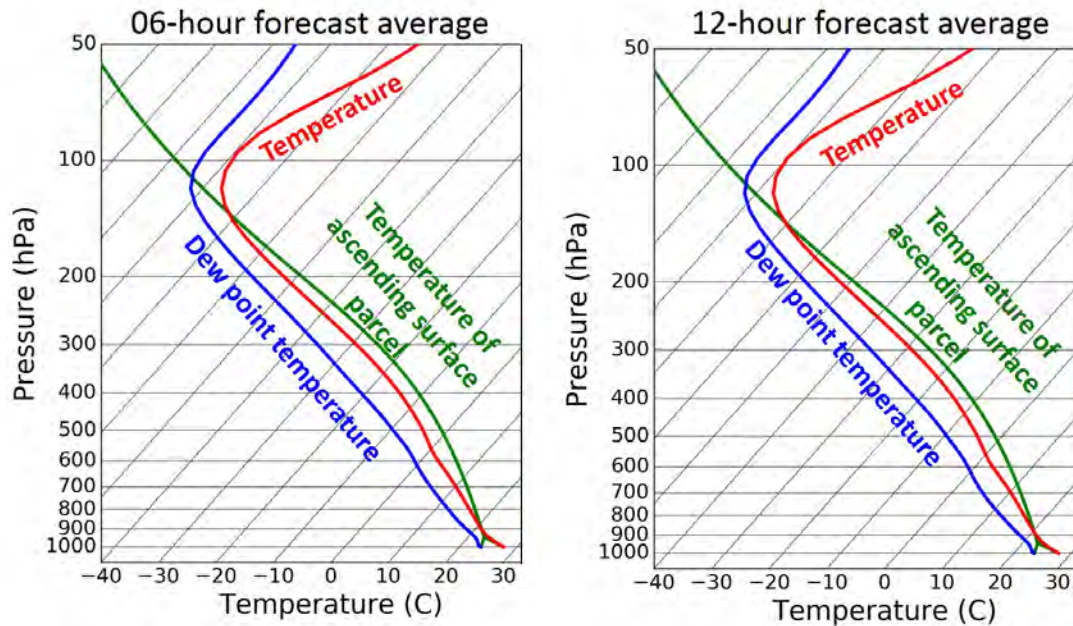


Figure 3.3: SkewT-logP plots of the two forecast sets (left: 06-hour forecasts, right: 12-hour forecasts), obtained by averaging pressure, specific humidity and temperature for every model level. The average is performed across the horizontal domain of each model level, across all forecasts with the same length. Temperature and dewpoint temperature profiles are as indicated by the solid red and blue lines respectively. The solid green curve indicates the profile of a surface air parcel as it ascends adiabatically. Note that the base of the green curve actually coincides with the base of the red curve. The green “leg” jutting out from the base of the blue curve traces the isotherm of the surface dew point temperature as surface parcel ascends. The place where the “leg” and the green curve meet is where the relative humidity of the surface parcel has reached 100%.

3.1 Time-averaged thermodynamic profile is as expected

Tropopause pressure level is correct SkewT-logP diagrams are useful for examining the likelihood of convection and the type of convection in the atmosphere. In Figure 3.3, skewT-logP¹ diagrams of both forecast sets are shown. The plotted profiles are constructed by taking the horizontal and temporal average of temperature, specific humidity and pressure for each model level. It is clear that the temperature minima occurs slightly above the 100 hPa level, suggesting that the tropopause resides roughly at that level. This is as expected for an equatorial region.

¹These skewT-logP plots are produced with a modified version of the `SkewT` package in Python, released by Thomas Chubb, Department of Mathematical Sciences, Monash University, copyrighted 2014.

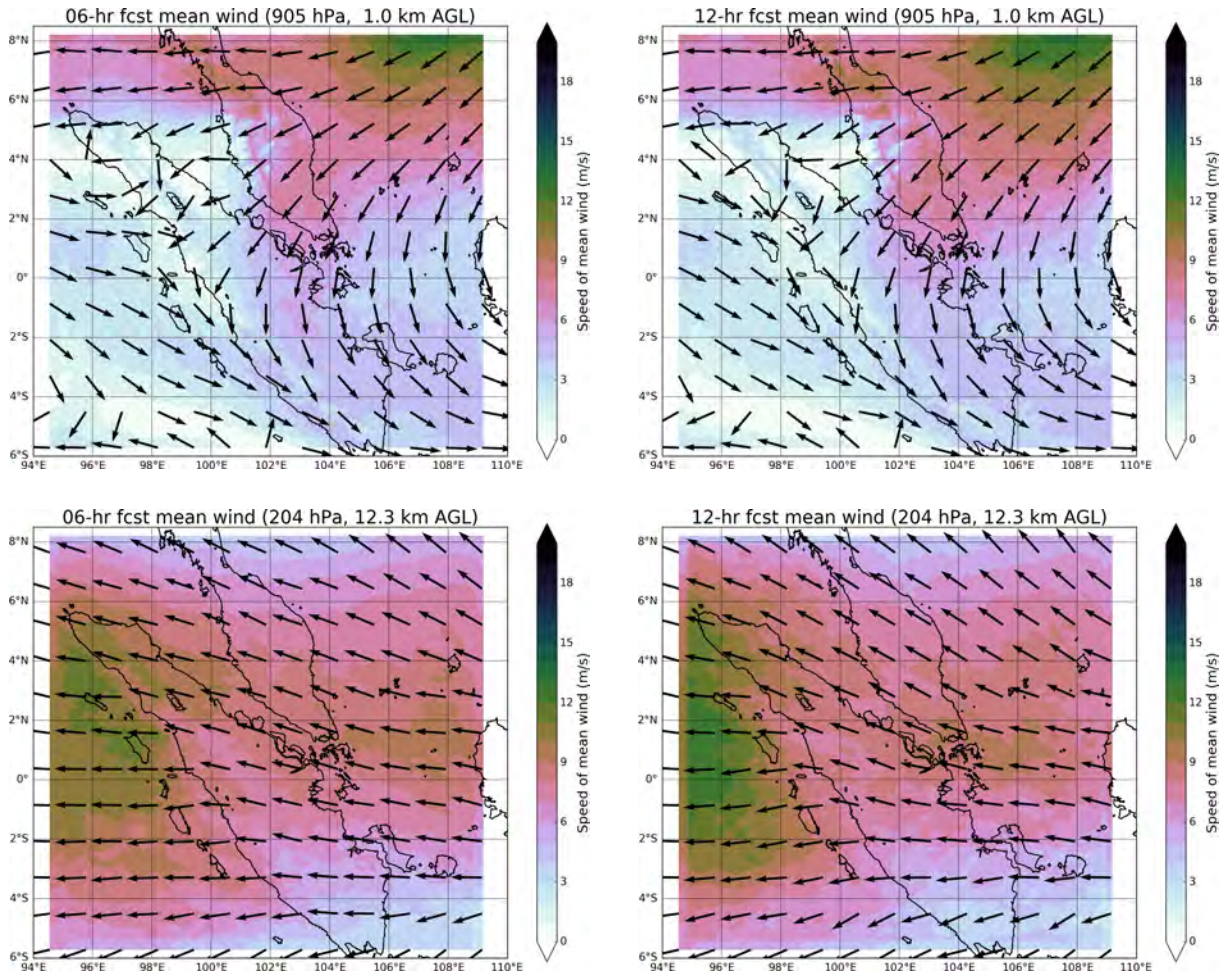


Figure 3.4: Contour plots of the time-averaged wind of the 06-hour forecasts (left column) and the 12-hour forecasts (right column), overlaid with arrows indicating the direction of wind flow, at the SINGV model levels with average pressures of 900 hPa (top) and 200 hPa (bottom). Note that the arrows only indicate the time-averaged wind direction, meaning that all arrows have the same lengths.

Atmosphere is unstable Figure 3.3 also shows that the atmosphere for both forecast sets are unstable and prime for deep convection. While the temperature of a surface air parcel undergoing dry adiabatic ascension is cooler than the environmental temperature for the bottom 100 hPa of the atmosphere, the difference is rather slight. In other words, only a small amount of energy is needed for the parcel to overcome the barrier to adiabatic ascension in the first 100 hPa of the atmosphere. Once the parcel passes the 900 hPa level (which is the level of free convection), the parcel ascends via a moist adiabat and remains continuously warmer than the atmosphere, allowing it to freely ascend, until it reaches the 130 hPa level. Clearly, the atmosphere is rather unstable. The closeness of this level to the tropopause indicates that the atmosphere is prime for deep convection. This is expected for the Maritime Continent during the Northeast monsoon (Chang et al., 2005).

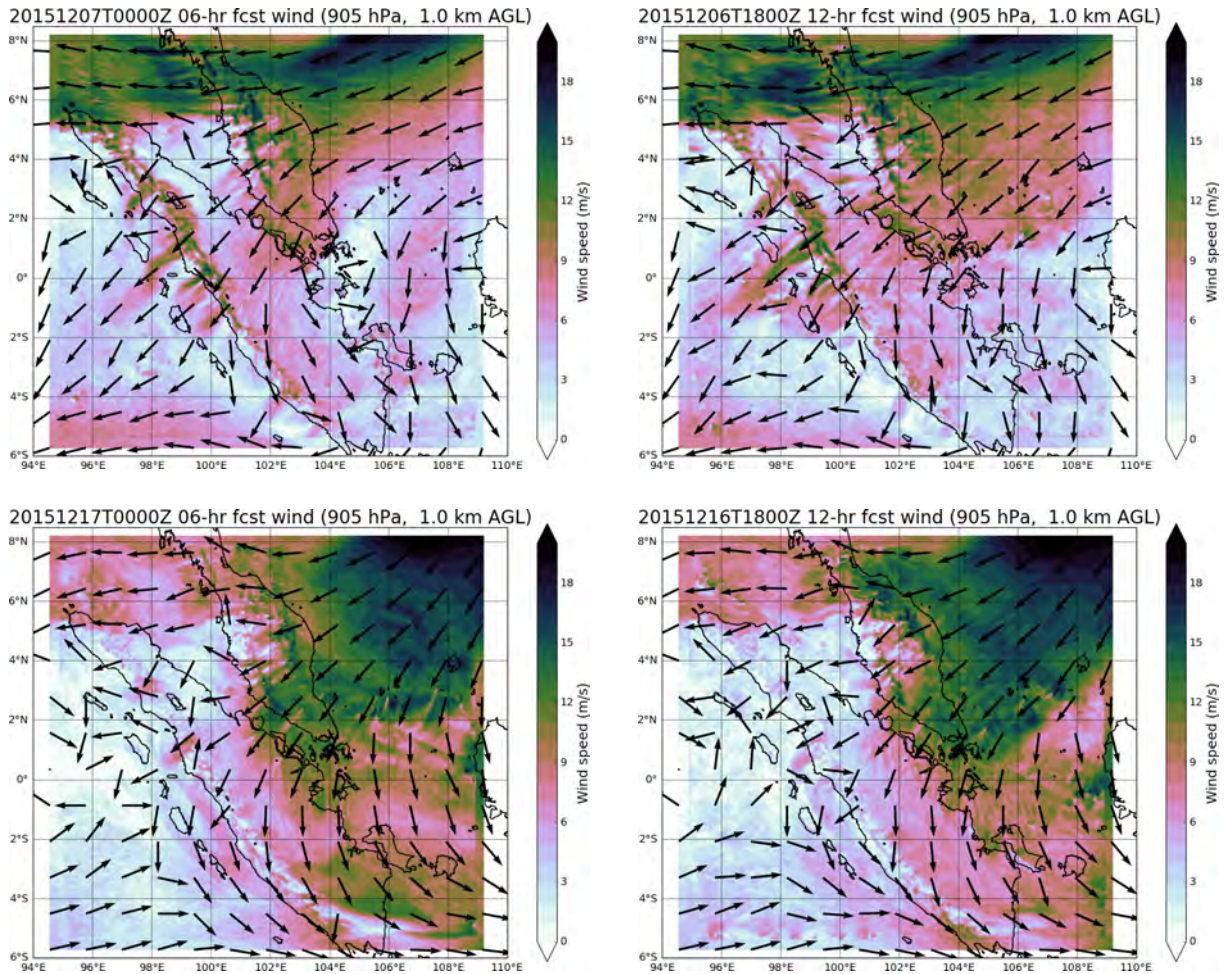


Figure 3.5: Same type of plots as Figure 3.4, but on the SINGV level corresponding to the 900 hPa. The data shown are snapshots of the cold surges at 0600 UTC, 7th December (top), and 0600 UTC, 17th December (bottom), observed in the 06-hour forecasts (left) and 12-hour forecasts (right).

3.2 900 hPa monsoon flow is as expected

Observed monsoon flow features Aside from that, on average, the 06-hour forecasts and the 12-hour forecast both show a low-level northeasterly wind entering the domain from the northeastern corner of the domain (Figure 3.4, 900 hPa plots). This wind flow is deflected leftwards as it approaches and crosses the Equator, resulting in a northwesterly low-level wind exiting the domain at the southeastern corner. This monsoon flow pattern is as expected from the monsoon’s historical climatology (Chang et al., 2005; Fong and Ng, 2012).²

²The monsoon flow pattern is explained by the paragraph labelled “Northeast monsoon season” in the introduction.

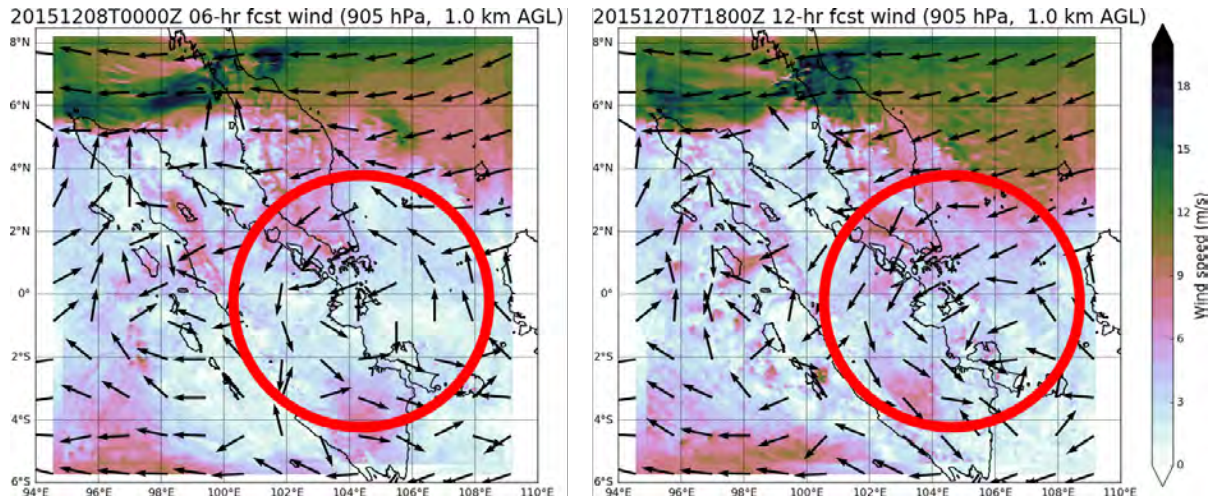


Figure 3.6: Same type of plots as Figure 3.4, but on the SINGV level corresponding to the 900 hPa. The data shows the Borneo vortex that appeared in the 06-hour forecast (left) and 12-hour forecast (right) valid on 0600 UTC, 8th December 2015. The vortex shows up as a closed, counterclockwise wind direction pattern (red circles) in both forecasts.

3.3 Cold surges manifested in the simulations

Cold surges observed Another important feature observed would be cold surges – bursts of low-level cold air that occasionally emanate from the Siberian surface high pressure system. The cold surge is a characteristic of the Northeast monsoon (Fong and Ng, 2012).³ Two clear cold surges were observed in the forecasts. The cold surges appeared in the 06-hour and 12-hour forecasts valid from 5th December, 1200 UTC, to 9th December, 1200 UTC (Figure 3.5), and from 15th December, 0000 UTC, to the end of the simulations. Clearly, the SINGV forecasts are able to resolve this key feature of the Northeast monsoon season.

3.4 Interesting vortex features

Vortex-like feature observed We also observed a rather interesting large-scale counterclockwise circulation (several hundred kilometres across) over equatorial South China Sea (Figure 3.6). The circulation shows up as a low-level closed, counterclockwise, circulation loop on the forecasts, valid from 0000UTC, 8th December to 0600 UTC, 9th December (Figure 3.6). Interestingly, the Australian Bureau of Meteorology’s Gradient Wind Analysis on the 9th of December 2015 also shows a similar circulation at the place where our

³A description and explanation of the cold surge can be found in the introduction, in the paragraph labelled “cold surges”.

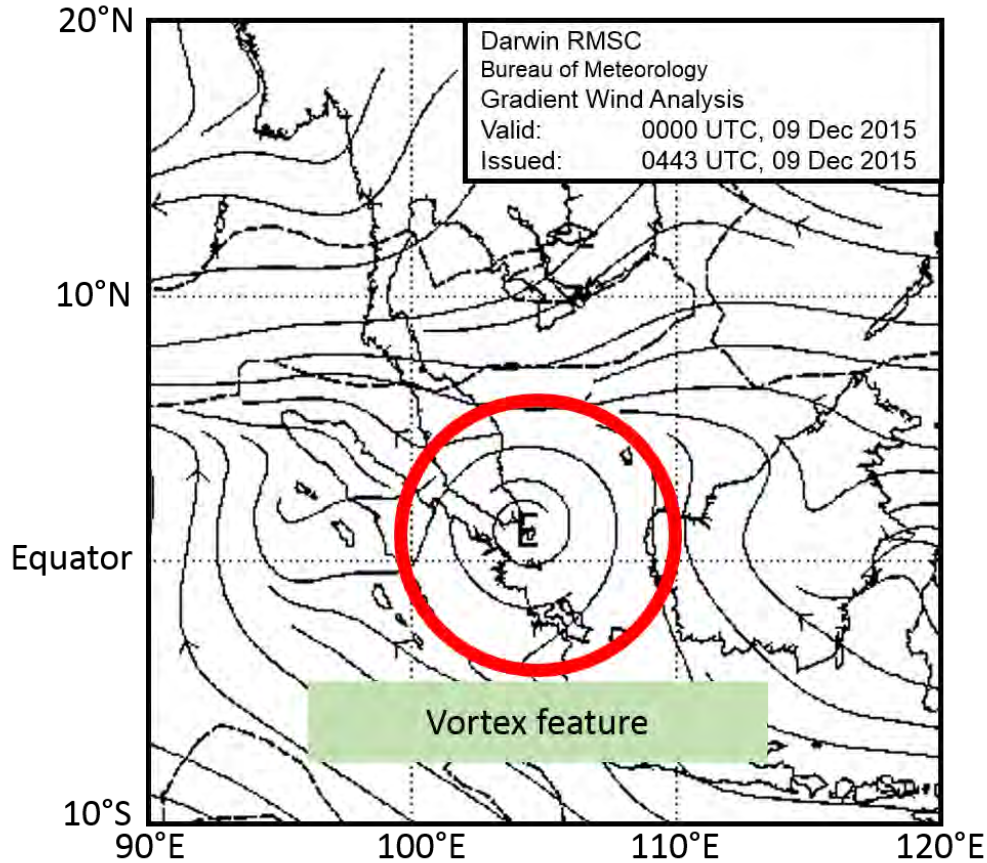


Figure 3.7: Streamlines plotted by the Australian Bureau of Meteorology (Regional Specialized Meteorological Center) on 9th December 2015, 0000 UTC. These streamlines show the wind flow at 1 km above the sea level. Note the vortex feature indicated in the red circle. This plot is modified from the analysis charts stored in the Australian Bureau of Meteorology Analysis Chart Archive (Australian Government, 2016).

circulation appears in the forecast (Figure 3.7). In other words, this circulation pattern may not be a spurious artefact of the model.

Uncertainty over identity of circulation At the time of writing, we are uncertain if the circulation is a Borneo vortex⁴, which is a characteristic of the Northeast monsoon. While the scale and morphology of the circulation pattern matches, its position is too far west: Borneo vortices usually occur in the region between the western tip of Borneo and off the east coast of the Malay Peninsula (Chang et al., 2005).

Circulation unlikely to affect error covariances significantly Even if the circulation ultimately is not the characteristic Borneo vortex of the Northeast monsoon, the

⁴Low-level, converging, cyclonic vortices that tend to occur over the western part of Borneo (near the Equator, and east of the model domain), and are typically accompanied with deep convection.

short duration of this circulation means that it is unlikely to have much of an impact on the overall error covariance of the region. The circulation clearly vanished from the forecast after 0600 UTC, 10th December. In other words, the circulation only persisted for 12% of the 15 days considered. Furthermore, the circulation itself wandered from the southeastern part of the domain to the northeastern part of the domain, diluting the influence of any errors relating to the circulation over a large trajectory. As such, we can safely ignore the influence of the circulation when analysing the forecast error covariances.

3.5 Summary

In summary, we have confirmed that the SINGV forecasts are behaving as they should in the presence of the Northeast monsoon: the identified tropopause is in the correct pressure level, the monsoon flow is as expected and the characteristic cold surges are observed. The unusual circulation observed in the forecasts is also unlikely to have a great influence on the covariances of the forecast error. We can thus infer that the forecast error covariances estimated from these forecasts to be representative of the Maritime Continent's Northeast monsoon.

Chapter 4

Error covariance features and regression coefficient models

Using NMC to examine error covariance features While it is impossible to determine \mathbf{B} due to memory problems, it is possible to determine several columns of \mathbf{B} . In this chapter, we will utilize the NMC method to determine several columns of the \mathbf{B} matrix. Each column of \mathbf{B} contains the covariance between the forecast error of a particular variable, at a particular point on the model grid, against the forecast error of both variables, across the entire model grid. Such columns can be used to diagnose whether a regression coefficient model's assumptions holds.

Overview We will begin with explaining the physical meaning of the two error variables we are studying in this project. Following that, we will explain how the single-observation solution to the data assimilation cost function can be used to interpret physical meaning from columns of \mathbf{B} . Several columns of \mathbf{B} will then be displayed and interpreted. Using these columns, we will overturn the horizontal invariance assumption utilized in two commonly used models of the control variable transform regression coefficients¹. A new assumption will then be proposed from examining the columns of \mathbf{B} and a new model will be formulated based on that².

Components of horizontal wind error

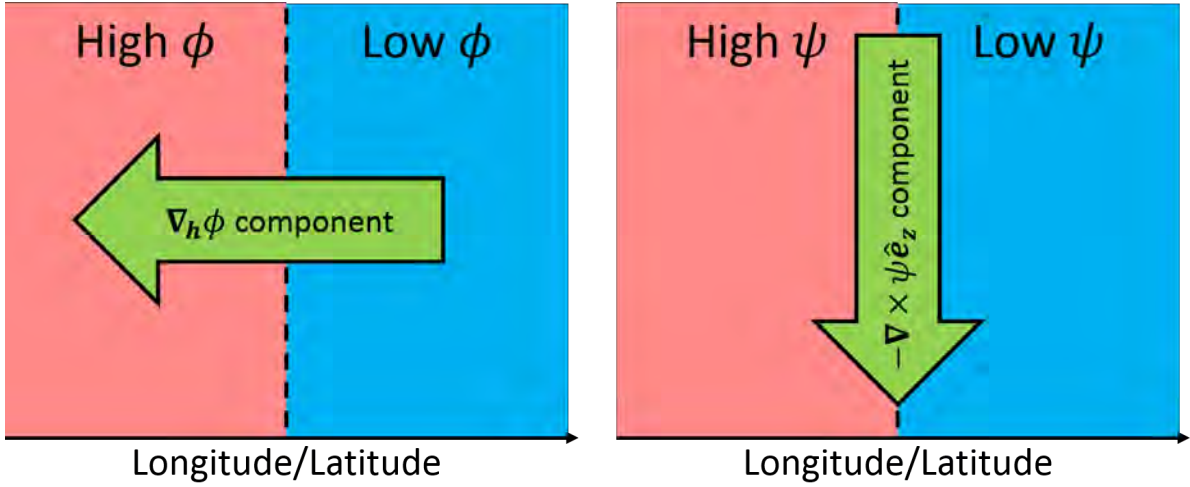


Figure 4.1: Diagrams showing the directions of the horizontal wind error vector field, induced by ϕ (left) and ψ (right). Note that the z -axis is pointing out of the page.

4.1 Physical meaning of ϕ and ψ

Meaning of ϕ As mentioned in the introduction, the error in the horizontal wind field can be written as

$$\mathbf{u}_k = \nabla\phi(x, y, k) - \nabla \times (\psi(x, y, k) \hat{\mathbf{z}}). \quad (4.1)$$

The general appearances of both parts of the wind field's forecast error can be directly inferred from the contour lines of ϕ and ψ . For the case of the divergent part of the wind ($\nabla\phi$), consider the fact that the gradient of a scalar is perpendicular to the scalar's contour, and points in the direction of maximum increase (Figure 4.1, left panel). In other words, a positive hill of ϕ corresponds to a local flow convergence that is centered around the maximum and *vice versa* for a valley of ϕ (see Figure 4.2 for illustrations).

Meaning of ψ It is slightly more difficult to interpret ψ . Defining $\mathbf{u}_\psi \equiv -\nabla_{\mathbf{H}} \times (\psi \hat{\mathbf{z}})$, consider the following integral over an area bounded by a ψ contour line on a fixed model level:

$$\int_S \nabla_{\mathbf{H}}^2 \psi dA = \int_S \nabla_{\mathbf{H}} \times \mathbf{u}_\psi \cdot \hat{\mathbf{z}} dA.$$

¹These are schematically illustrated in Figures 1.1 and 1.2.

²The computation of the regression coefficients under the new model is schematically illustrated in Figure 1.3. Figure 1.4 shows how this model is utilized in the control variable transform.

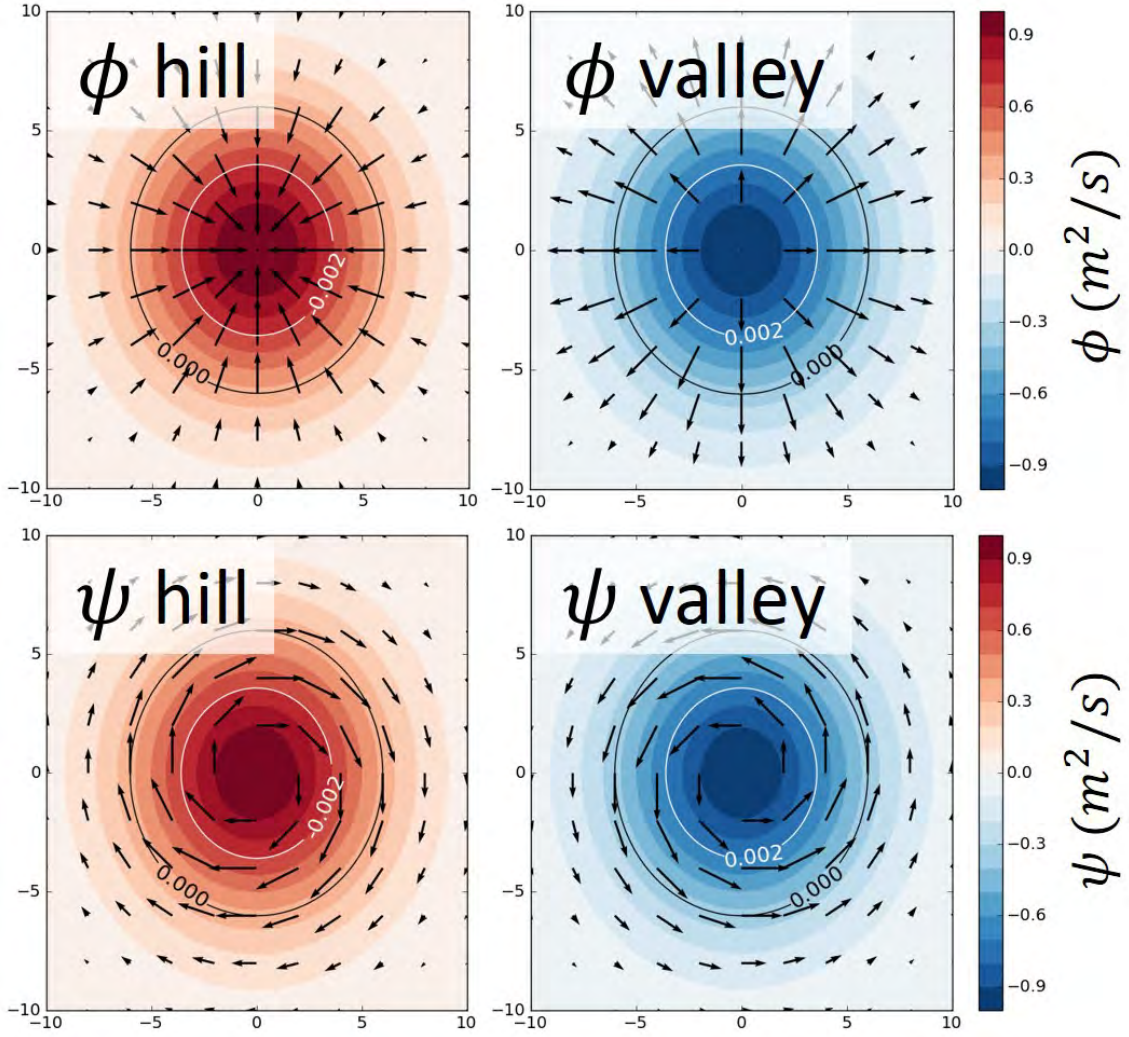


Figure 4.2: Contour plots showing examples of hills (left column) and valleys (right column) of ϕ (top row) and ψ (bottom row), overlaid with vector arrows indicating the respective components of wind implied by these two variables. The unfilled contours indicate horizontal divergence (top row) and the z -component of vorticity (bottom row).

Applying the 2D form of the divergence theorem and Green's theorem yields,

$$\oint_{\partial S} \nabla_H \psi \cdot d\mathbf{l}_\perp = \oint_{\partial S} \mathbf{u}_\psi \cdot d\mathbf{l}_\parallel.$$

where $d\mathbf{l}_\perp$ is a vector pointing out of area S and perpendicular to the contour, and $d\mathbf{l}_\parallel$ is a length vector element along the contour, oriented such that S lies on the left of $d\mathbf{l}_\parallel$. Since the contour is on a fixed model level, then

$$d\mathbf{l}_\parallel \times \hat{\mathbf{z}} = d\mathbf{l}_\perp.$$

Applying this relation and considering $(\nabla_{\mathbf{H}}\psi) \cdot (d\mathbf{l}_{\parallel} \times \hat{\mathbf{z}}) = (d\mathbf{l}_{\parallel}) \cdot (\hat{\mathbf{z}} \times \nabla_{\mathbf{H}}\psi)$ yields,

$$\oint_{\partial S} (d\mathbf{l}_{\parallel}) \cdot (\hat{\mathbf{z}} \times \nabla_{\mathbf{H}}\psi) = \oint_{\partial S} \mathbf{u}_{\psi} \cdot d\mathbf{l}_{\parallel},$$

or,

$$(\hat{\mathbf{z}} \times \nabla_{\mathbf{H}}\psi) \cdot d\mathbf{l}_{\parallel} = \mathbf{u}_{\psi} \cdot d\mathbf{l}_{\parallel}.$$

Since the gradient is horizontal and perpendicular to the contour line, $(\hat{\mathbf{z}} \times \nabla_{\mathbf{H}}\psi)$ yields a vector is either antiparallel or parallel to $d\mathbf{l}_{\parallel}$. The right panel of Figure 4.1 illustrates how \mathbf{u}_{ψ} can be inferred from ψ fields using this idea. For the case of a ψ hill bounded by the contour line, $(\hat{\mathbf{z}} \times \nabla_{\mathbf{H}}\psi)$ is antiparallel to $d\mathbf{l}_{\parallel}$. In other words, cyclonic circulation will be represented by a ψ valley (see Figure 4.2 for illustrations).

4.2 Columns of \mathbf{B}

4.2.1 Interpreting columns of \mathbf{B}

Single-observation solution interprets \mathbf{B} To understand the physical meaning of a column of \mathbf{B} , we need a way to relate the information in the column of error covariances to the result of data assimilation itself. Interpreting \mathbf{B} through the lens of the single-observation solution to the data assimilation cost function in Eqn (1.4) is one intuitive way.

Single-observation solution conditions and statement Supposing that the observation corresponds to the n -th element on the model state vector, \mathbf{X} , then the ℓ -th element of the analysis state vector, \mathbf{X}^a , can be shown to be (Bannister, 2008a)

$$X_{\ell}^a - X_{\ell}^b = B_{\ell n} \frac{Y - X_n^b}{\sigma + B_{nn}}. \quad (4.2)$$

This is the single-observation solution³, and it operates under three conditions. First of all, only one observation, Y , is assimilated. Secondly, this observation must be on the same location as one of the forecast model’s grid points. Lastly, the observation must be an observation of one of the variables in the model. Henceforth, $\mathbf{X}^a - \mathbf{X}^b$ will be called the “analysis increment”.

³The derivation of this solution is available in Appendix A.

An analysis increment is proportional to a column of \mathbf{B} As can be seen from Eqn (4.2), the analysis increment is proportional to the n -th column of \mathbf{B} , with $\frac{Y-X_n^b}{\sigma+B_{nn}}$ as the proportionality constant. In other words, we can simply infer how \mathbf{X}^b will be modified to form \mathbf{X}^a from a column of \mathbf{B} . We will use this framework to determine the physical implications of the information contained in \mathbf{B} .

4.2.2 Computing columns of \mathbf{B} from NMC

General form of \mathbf{B} As mentioned earlier, we can determine several columns of \mathbf{B} via the NMC method, without suffering from computer memory problems. We will explain this procedure in detail. For that, we need the form of \mathbf{B} . It is simply

$$\mathbf{x}^b \equiv \begin{bmatrix} \boldsymbol{\psi} \\ \boldsymbol{\phi} \end{bmatrix} \implies \mathbf{B} = \begin{bmatrix} \langle \boldsymbol{\psi} \boldsymbol{\psi}^\top \rangle_t & \langle \boldsymbol{\psi} \boldsymbol{\phi}^\top \rangle_t \\ \langle \boldsymbol{\phi} \boldsymbol{\psi}^\top \rangle_t & \langle \boldsymbol{\phi} \boldsymbol{\phi}^\top \rangle_t \end{bmatrix},$$

where $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$ are vectors containing the respective forecast errors ($N/2$ elements each) and the averaging is done over time.

Meaning of a column of \mathbf{B} The m -th column in the left half of \mathbf{B} ($m < N/2$) contains 2 sets of information. It contains the covariance between $\boldsymbol{\psi}$ at the position corresponding to m -th element on the $\boldsymbol{\psi}$, versus $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$ everywhere on the model grid. Likewise, the m' -th column on the right half of \mathbf{B} ($m' > N/2$) contains the covariance between $\boldsymbol{\phi}$ at the $m' - N/2$ position on $\boldsymbol{\phi}$ vector, versus $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$ everywhere on the model grid.

NMC \mathbf{B} column computation method From Eqn (2.2), we can compute the vectors of the NMC-estimated forecast errors at time index t via

$$\boldsymbol{\psi}_t \approx \frac{1}{\sqrt{2}} (\boldsymbol{\psi}_{12\text{-hr fcst},t} - \boldsymbol{\psi}_{06\text{-hr fcst},t} - \langle \boldsymbol{\psi}_{12\text{-hr fcst},t} - \boldsymbol{\psi}_{06\text{-hr fcst},t} \rangle_t),$$

and likewise for $\boldsymbol{\phi}$. As such, to compute the column of \mathbf{B} corresponding to $\boldsymbol{\psi}$ at model index (i, j, k) (longitudinal, latitudinal, vertical), we simply use

$$\begin{aligned} \langle \boldsymbol{\psi} \boldsymbol{\psi} (i, j, k) \rangle_t &\approx \frac{1}{N_T} \sum_{t=1}^{N_T} \boldsymbol{\psi}_t (i, j, k) * \boldsymbol{\psi}_t, \\ \langle \boldsymbol{\phi} \boldsymbol{\psi} (i, j, k) \rangle_t &\approx \frac{1}{N_T} \sum_{t=1}^{N_T} \boldsymbol{\psi}_t (i, j, k) * \boldsymbol{\phi}_t. \end{aligned}$$

A similar method can be utilized to compute the column of \mathbf{B} corresponding to ϕ at (i, j, k) . When interpreting the columns of \mathbf{B} , we will view the covariances on the left hand side at all displayed locations on the model grid.

4.3 Features observed in columns of \mathbf{B}

4.3.1 Deep convection

Recap: forecasts show deep convection tendencies With the single-observation solution framework to interpret columns of \mathbf{B} and the NMC column estimation laid out, we are ready to examine the physical meaning implied by the columns of \mathbf{B} . We will begin with deep convection. As per the skewT-logP discussion in Chapter 3, the forecasts of the Northeast monsoon display a tendency for deep convection. As such, we would expect to see the forecast error covariances to reflect that.

Auto-covariances of ϕ at the surface level To confirm that, 9 columns of the $\langle \phi\phi^\top \rangle$ sub-matrix of \mathbf{B} were computed. The 9 columns are such that their ϕ reference points are all located on the lowest model level (indicated by the 9 green dots in Figure 4.3). The covariances between ϕ at each of these 9 locations (green dots), and ϕ in a 450 km by 450 km square centred on each location, are plotted out in Figure 4.3.

Mesoscale surface convergence with positive ϕ observation increment Let us suppose that we are assimilating an observation of ϕ that is bigger than the corresponding background state's ϕ (*i.e.*, positive ϕ observation increment). Following the discussion on interpreting the columns of \mathbf{B} through Eqn (4.2), Figure 4.3 implies that assimilating that observation at any green dot causes a hill of ϕ ($O(10^2)$ km diameter) to appear in the analysis increment, centred on the said green dot. Physically, this means that the assimilation process enhances mesoscale convergence ($O(10^2)$ km wide) on the surface level.

Deep atmosphere mesoscale divergence with positive ϕ observation increment Furthermore, when we consider the vertical structure of these ϕ - ϕ covariances, we notice that the positive observation increment at the surface triggers a ϕ valley aloft. This is demonstrated by Figure 4.4, which shows the vertical structure of the said covariances, from the three green dots at 1°N in Figure 4.3. Clearly, assimilating the observation

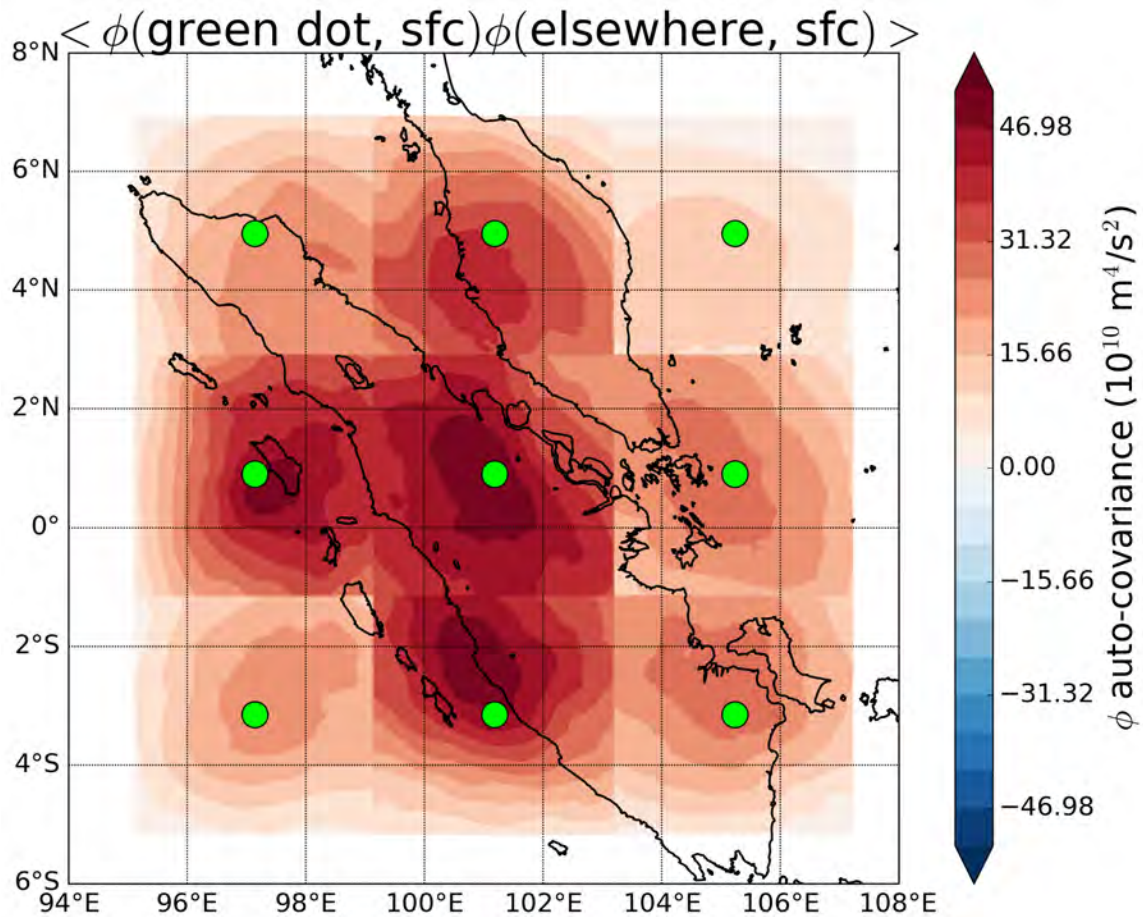


Figure 4.3: The contour plot shows the covariance between ϕ at the green dots and ϕ in squares centred on the green dots, on the surface level. The dimensions of each square is $450 \text{ km} \times 450 \text{ km}$, or roughly $4^\circ \times 4^\circ$. The landmasses in the domain are outlined in thick, solid black lines.

enhances divergence aloft ($O(10^2)$ km across). Furthermore, the valley feature occurs in the vicinity of the 50-th model level, which corresponds to the 110 hPa layer, according to Figure 3.1. In other words, the data assimilation triggers deep atmosphere mesoscale divergence.

Deep convection error feature Taken together, the mesoscale surface convergence and ~ 100 hPa mesoscale divergence pattern implies that the error covariances are organized in a deep convection pattern. Furthermore, given the ubiquity of this pattern, it seems that using a positive ϕ observation increment anywhere on the surface enhances convection. In other words, the ϕ - ϕ covariances reflect the tendency for deep convection (as discussed in Chapter 3).

Foreshadow – Covariances \approx variances near green dots An interesting observation can be made from Figures 4.3 and 4.4. Notice that the covariances tend to vary slowly

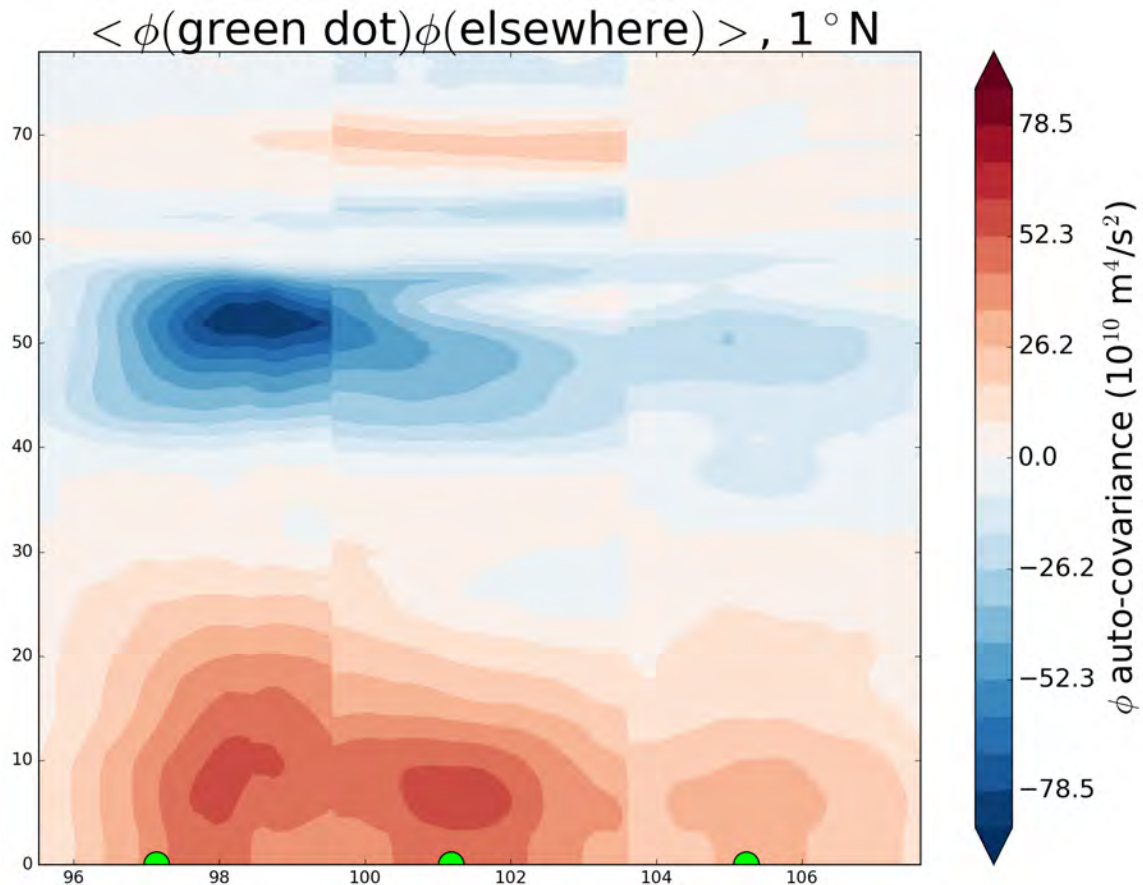


Figure 4.4: The contour plot shows the covariances of ϕ from the green dots along the 1°N latitude in Figure 4.3, and ϕ everywhere on a vertical slice of ϕ , along the 1°N latitude. This plot is separated into 3 vertical strips, each centred on a the green dot, spanning 4° longitude in the horizontal and all 79 model levels in the vertical. The horizontal axis is longitude and the vertical axis is model level.

immediately around the green dots. In other words, the covariances immediately around the green dots are roughly equal to the variance of ϕ at the green dots. This intriguing property will be used later when formulating our own regression coefficient model.

4.3.2 Vortex patterns

Surface circulations from ψ observation increment When we computed the covariance between ψ at the same 9 points, and ψ everywhere on the surface of the model, hills similar to those of ϕ - ϕ surface covariances appear (Figure 4.5). However, the physical meaning for the analysed wind field is completely different from that of ϕ . Assimilating a positive ψ observation increment generates positive ψ hills in the analysis increment. Following the interpretation in Figure 4.2, these ψ hills mean that this assimilation enhances clockwise (or anticyclonic) circulation surrounding the green dots.

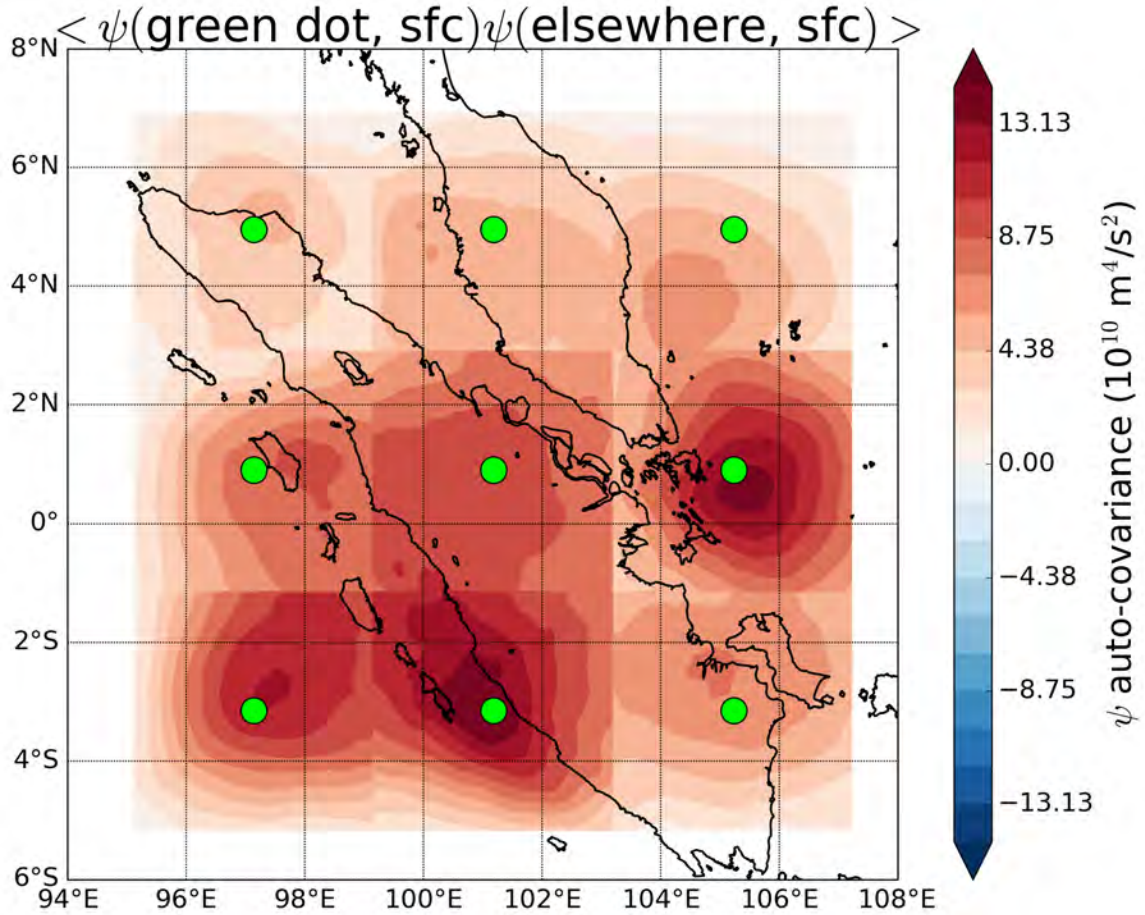


Figure 4.5: The contour plot shows the covariance between ψ at the green dots and ψ in squares centred on the green dots, on the surface level. The dimensions of each square is $450 \text{ km} \times 450 \text{ km}$, or roughly $4^\circ \times 4^\circ$. The landmasses in the domain are outlined in thick, solid black lines.

Foreshadow – Covariances \approx variances near green dots At the time of writing, we are unable to determine the physical origin of these patterns for ψ , both on the surface (Figure 4.5) and in the vertical (Figure C.1 in Appendix C). However, like the case of ϕ - ϕ covariances, we notice that the covariances tend to be similar to variances immediately around the green dots. This pattern will be utilized later to formulate our own regression coefficient model.

4.3.3 Horizontally variant and anisotropic cross-covariances

ϕ - ψ covariance columns are horizontally variant and anisotropic We will now turn our attention to the cross-covariance between ϕ and ψ . The covariances between ϕ at the same 9 reference positions and ψ elsewhere in 450 km by 450 km squares centred on the reference positions are computed and plotted in Figure 4.6. It is immediately apparent that unlike the ϕ - ϕ and ψ - ψ columns of \mathbf{B} , these ϕ - ψ columns show vastly

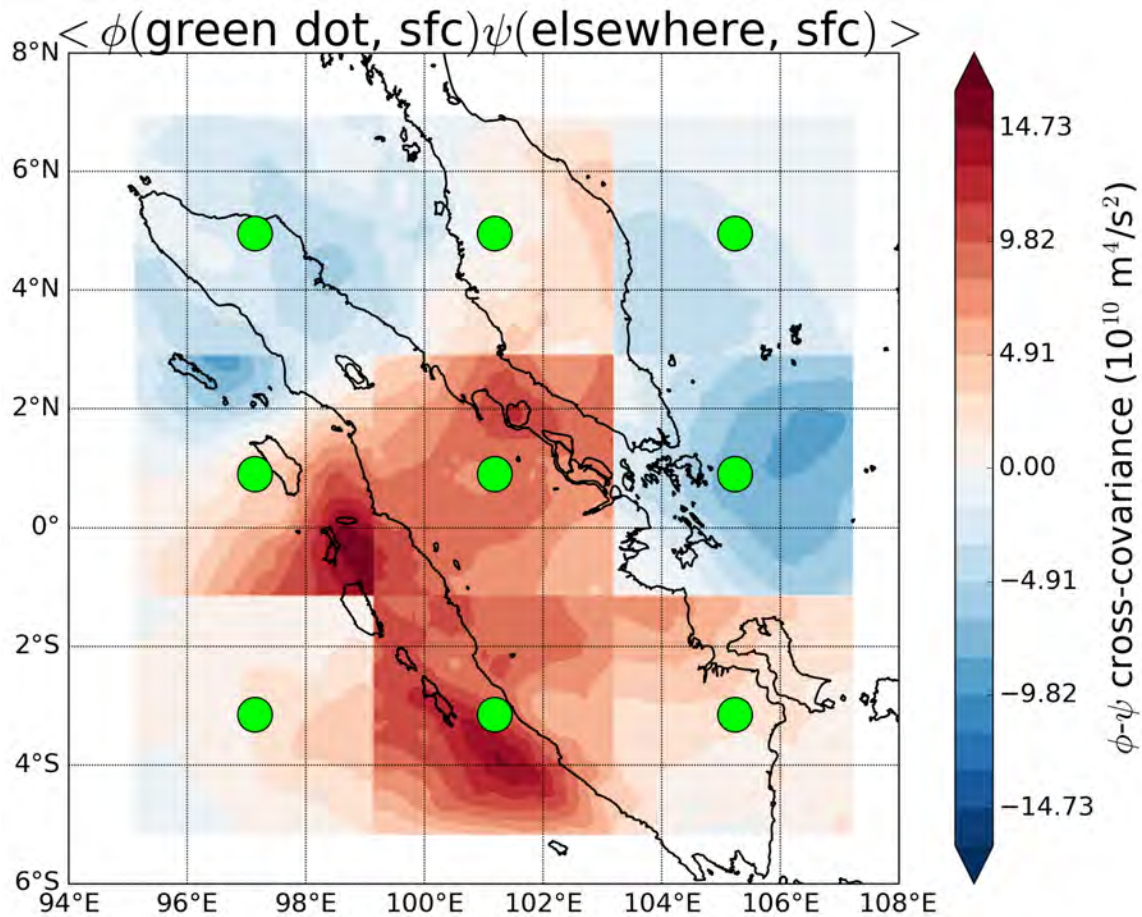


Figure 4.6: The contour plot shows the covariance between ϕ at the green dots and ψ in squares centred on the green dots, on the surface level. The dimensions of each square is $450 \text{ km} \times 450 \text{ km}$, or roughly $4^\circ \times 4^\circ$. The landmasses in the domain are outlined in thick, solid black lines.

different error patterns from column to column. Furthermore, the covariances do not vary isotropically around the green dots. In other words, the cross-covariance columns exhibit huge variations between horizontal locations and are severely anisotropic.⁴

Horizontal variations and anisotropy may be due to complicated topographic interactions

The variations and anisotropy may be due to complicated topographic interactions. As an illustration, we will examine the cross-covariances connected to the green dot at 1°N , 97°E of Figure 4.6. Figure 4.7 display the cross-covariances between ϕ at the green dot, and everywhere else in the model grid, on the surface level.⁵ When an assimilated observation of ϕ at the green dot (1°N , 97°E) increases ϕ at that location, the analysis increment will contain a ϕ hill centered on the green dot (see Figure 4.3) and a ψ dipole, with the negative patch to the northwest and the positive patch to the southeast

⁴Similar observations are also made when we consider the 9 columns corresponding to the covariance between ψ at the same 9 positions and ϕ everywhere on the surface (Figure C.2 in Appendix C).

⁵Essentially, an expanded version of the square in Figure 4.6 that is centred at 1°N , 97°E .

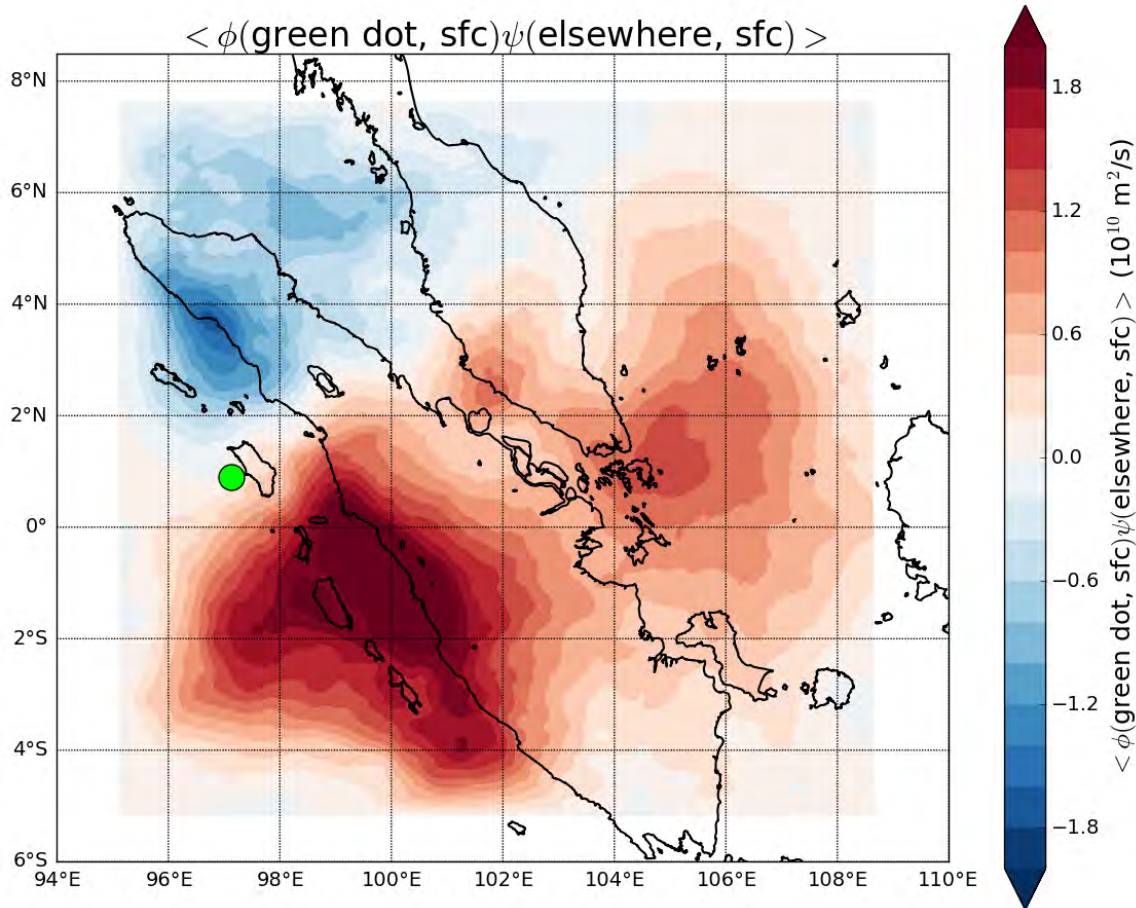


Figure 4.7: The contour plot shows the covariance between ϕ at the green dot, surface level, and ψ everywhere else on the surface level of the model grid. The landmasses in the domain are outlined in thick, solid black lines.

(Figure 4.7). This analysis increment can be interpreted as a southwesterly enhancement of the wind flowing through the green dot⁶. The southwesterly impinges on the island beside the green dot, resulting in surface convergence due to topographic obstruction, giving rise to the ϕ hill. Downstream of the green dot, this enhanced wind encounters the pass of the mountain range that runs along Sumatra (see Figure 2.1). As such, the wind partially diverges to the northwest and southeast and the remainder passes through the mountain pass. The divergence explains the intense ψ dipole feature on the windward side of the mountain range, while the jet resulting from the passage through the mountain pass explains the weaker dipole seen downstream of the pass. In the other locations of the domain, the topography is completely different, meaning that different interactions, and thus cross-covariance patterns, are to be expected.⁷ In other words, topography may be part of the reason behind the dramatic horizontal variations and anisotropies observed.

⁶Refer to Figure 4.1 to see how the ψ part of the wind field can be interpreted from ψ contours

⁷Also, the atmospheric forcing can be different. In the eastern half of the domain, the error cross-covariances are likely influenced by the monsoon flow.

4.4 Regression coefficient models

4.4.1 Conventional regression coefficient models and their flawed assumptions

Flawed conventional regression coefficient models These horizontal variations and anisotropies are contrary to the assumptions of the two commonly used regression coefficient models (see Figures 1.1 and 1.2). The first model assumes that there are no horizontal variations in the regression coefficient and that there are zero covariances between variables on different model levels (Chen et al., 2013). The second model relaxes the second assumption, but keeps the assumption of horizontal invariance (Descombes et al., 2015). On top of that, these models assume that the relationship between variables on two model grid points are isotropic (Bannister, 2008b). This manifests in the lack of any horizontal directionality in these definitions of regression coefficients. However, as observed earlier, the cross-covariances are hardly horizontally invariant or isotropic. *I.e.*, the regression coefficients are likely to be horizontally variant and anisotropic! We can already foresee that these conventional models will have bad performance in our study.

Conventional models will still be used Despite the breakdown of the assumptions, these two models are still worth testing. First of all, the performance of these two models can be used to benchmark the performance of our own model (which will be introduced later). Furthermore, these two models result in very small matrices of α , and can thus be easily computed. Lastly, due to the limited number of forecast pairs, any computation of regression coefficient matrices must contend with the problem of undersampling. If there are insufficient samples of forecast errors, the forecast error auto-covariance matrices are definitely singular. In other words, inverting these matrices to determine the regression coefficients is out of the question. However, these two models circumvent this undersampling problem by considering all the points in each model level as a sample. This typically overcomes the problem of undersampling.

4.4.2 New regression coefficient matrix model

Requirements A better model than the two conventional models would be one that uses more realistic assumptions while preserving the benefit of computational ease and circumvents undersampling. Here, we will deduce a more suitable assumption from examining the earlier plots of the columns of \mathbf{B} .

Assumption: points near reference point are samples of a specified point The assumption that we have decided to use is: the valuations of ϕ and ψ in the immediate vicinity of a specified point are equivalent to samples of ϕ and ψ at the specified point. For the ease of writing, this assumption will be called the **neighbourhood equivalence assumption**. We will justify this assumption mathematically and from examining the earlier auto-covariance plots.

Mathematical justification of neighbourhood equivalence assumption Consider the Taylor expansion of ψ as a function of position and time. The expansion about a position \mathbf{r} , up to the linear order of \mathbf{h} , is simply,

$$\psi(\mathbf{r} + \mathbf{h}, t) \approx \psi(\mathbf{r}, t) + \mathbf{h} \cdot \nabla \psi(\mathbf{r}, t).$$

For the ease of discussion, the coordinates are Cartesian and ∇ is the 3D Cartesian gradient operator. We will also set the temporal mean of ψ to zero. Clearly, the temporal auto-covariance and variance of ψ in this small region are

$$\begin{aligned} \text{Var} \left\{ \psi(\mathbf{r}) \right\} &= \langle \psi(\mathbf{r})^2 \rangle, \\ \text{Covar} \left\{ \psi(\mathbf{r}), \psi(\mathbf{r} + \mathbf{h}) \right\} &\approx \langle \psi(\mathbf{r})^2 \rangle + \mathbf{h} \cdot \langle \psi(\mathbf{r}) \nabla \psi(\mathbf{r}) \rangle. \end{aligned}$$

In other words, in the vicinity of \mathbf{r} , the valuation of the auto-covariance can be very close to the variance at \mathbf{r} ! This implies that if the valuations of ψ on two nearby points, they can be treated as two samples of the same quantity.⁸

Justification from columns of \mathbf{B} : horizontal view The auto-covariance columns of \mathbf{B} also support this notion. As mentioned in the earlier parts of this chapter, we noticed that the auto-covariances tend to be similar to the variances of the green dots immediately around the green dots. In fact, when we divide the auto-covariances of ϕ with the variance of ϕ , and likewise for ψ , in Figure 4.8, the resulting ratio tends to be within 0.2 of 1 within the 0.5° of the green dot. In other words, we can roughly consider ϕ and ψ within a $1^\circ \times 1^\circ$ square centred at a point to be samples of ϕ and ψ at that point.

Justification from columns of \mathbf{B} : vertical view The same argument also applies in the vertical. From Figure 4.9, we infer that within roughly 5 model levels of a reference

⁸We acknowledge that there are scenarios where covariance of two variables can be equal to their variances, even though the two variables are clearly different. However, since we are dealing with auto-covariances, it seems intuitive to consider them the same quantity.

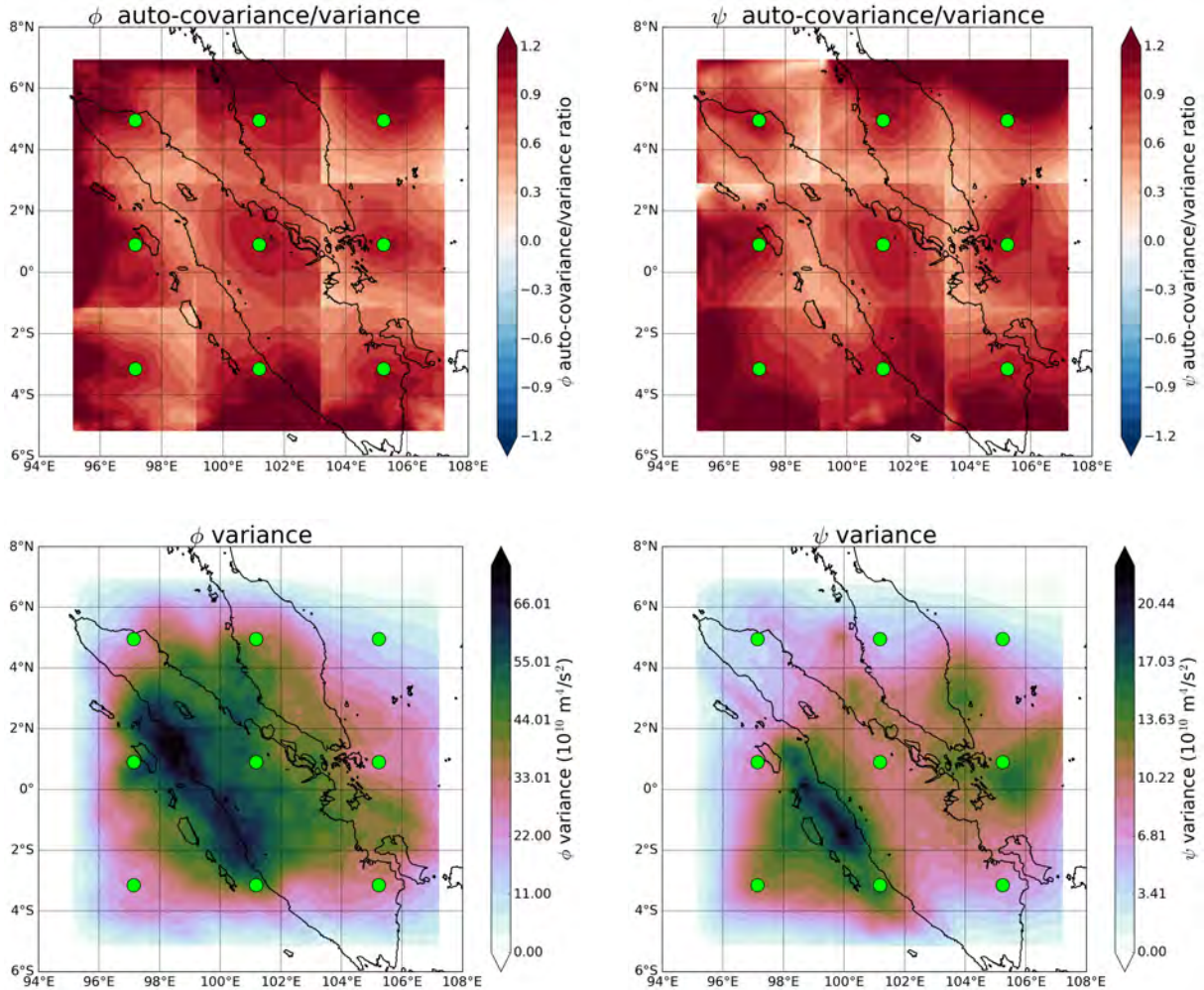


Figure 4.8: The horizontal contour plots show the ratio between the auto-covariances in Figures 4.3 (top left, for ϕ) and 4.5 (top right, for ψ), against the local variances. The local variances have been plotted in the corresponding columns of the bottom row.

point, the neighbourhood equivalence assumption applies.⁹

Neighbourhood equivalence assumption allows affordable α computation The neighbourhood equivalence assumption causes the regression coefficient matrix computation to be affordable. Since the equivalence appears to hold for a set of points within a subdomains with sides 1°-by-1°-by-5 model levels, we can split the entire model grid into such subdomains. The regression coefficients only need to be calculated between each pair of such subdomains. This drastically reduces the size of the regression coefficient matrix, making it affordable to compute.

Neighbourhood equivalence assumption prevents undersampling Furthermore, the assumption prevents undersampling from happening. Under the assumption, all points

⁹However, due to computational speed limitations, we will consider groups of 10 model levels later.

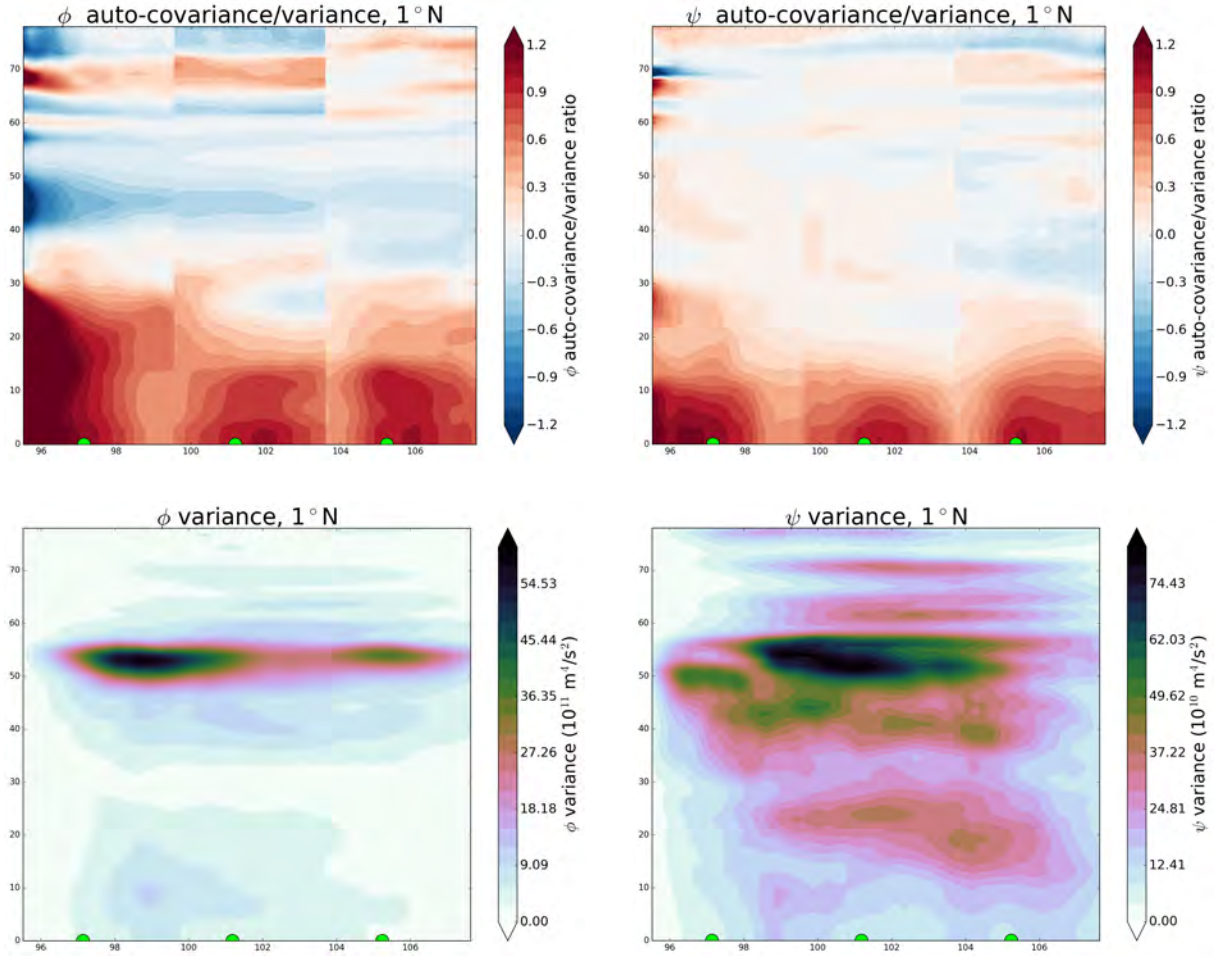


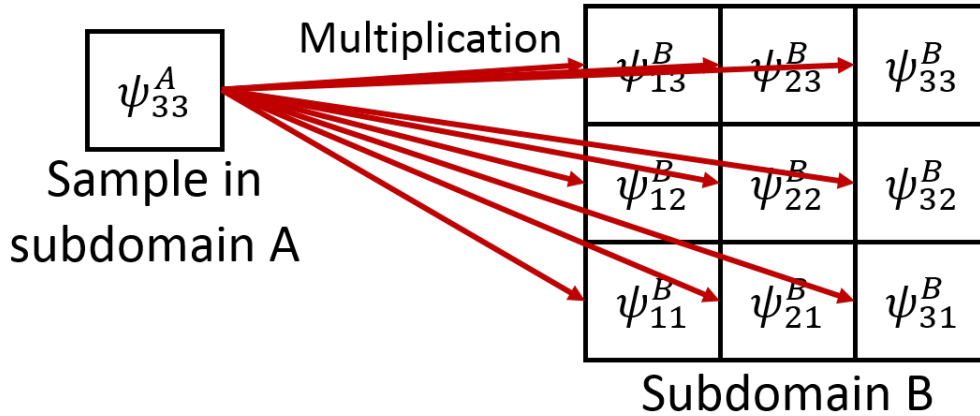
Figure 4.9: The vertical contour plots show the ratio between the auto-covariances in Figures 4.4 (top left, for ϕ) and C.1 (top right, for ψ), against the local variances. The local variances have been plotted in the corresponding columns of the bottom row.

within a subdomain are samples. This dramatically increases the number of samples by a factor of the number of points in the subdomain. Furthermore, the size of the auto-covariance matrices needed to compute the regression coefficients is reduced by the same factor, reducing the number of samples needed to prevent singular matrices. Clearly, formulating a model of regression coefficients under this assumption is feasible.

Need for sampling pattern of new model However, the assumption alone does not tell us the exact way to compute the reduced auto-covariance matrices. Suppose that we want to determine the auto-covariance of ψ between subdomain A and subdomain B. By the definition of covariance, we need to multiply the mean-removed samples of ψ in both subdomains to compute the covariance.

Point-by-point multiplication sampling We must avoid multiplying each mean-removed sample in subdomain A to all the mean-removed samples in subdomain B (Figure

Forbidden covariance sampling method



Permitted covariance sampling method

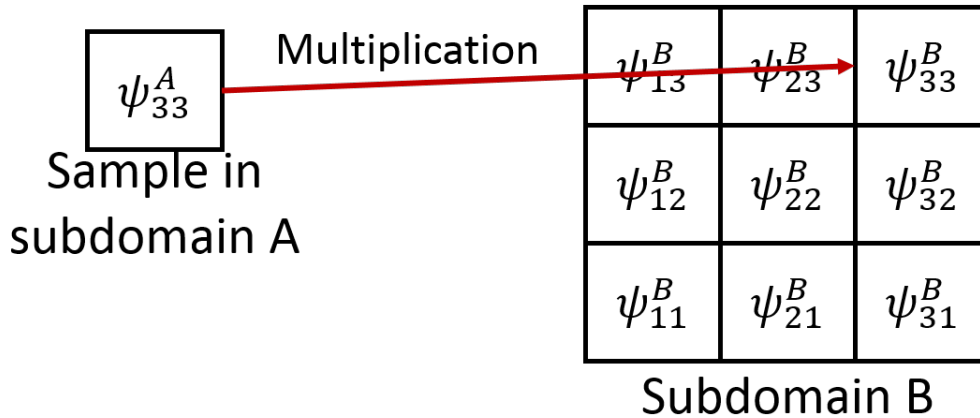


Figure 4.10: Two ways to use the samples in a pair of subdomains (A and B) to compute the multiples needed to generate auto-covariance of ψ between A and B.

4.10, top row). Mathematically, such an approach is exactly the same as taking the mean of all the values in A and B and use that as the sample. This approach kills off the increase in available samples. Instead, we should create the multiplication samples as per the second row in Figure 4.10. In other words, there is a one-to-one pairing of samples across the two subdomains. That will utilize all the samples, without causing undersampling.

Formulation of new model Following this sampling method, we can define the mathematical formulation of the new model. First, we will consider the (i, j, k) -th subdomain (counting from the bottommost-southmost-westmost subdomain) on the model grid to be the n -th subdomain. For the ease of computation, we will use subdomains that have the size of 1° by 1° by 10 model levels (to be precise, 20 by 20 by 10 model grid points). This

means that the model grid is split into a 3D array of $16 \times 15 \times 7$ subdomains. Thus, we can define the mapping from (i, j, k) to n ,

$$n(i, j, k) = i + j * 16 + k * 16 * 15. \quad (4.3)$$

Then, we can define the full mathematical formulation of the model (see Figure 1.3 for an intuitive schematic illustration of the model).

$$\begin{aligned} \psi'_{n;t}{}^{(3)} &= \psi_{n;t}{}^{(3)}, \quad \phi'_{n;t}{}^{(3)} = \phi_{b;t}{}^{(3)} - \sum_{i'=1}^{16} \sum_{j'=1}^{15} \sum_{k'=1}^7 \alpha_{\phi, \psi'}^{(3)}(n, n') \psi'_{n';t}{}^{(3)} \\ C_{\psi', \psi'; n, n'}^{(3)} &\equiv \sum_{t=1}^{N_T} \frac{\psi'_{n;t}{}^{(3)} \cdot \psi'_{n';t}{}^{(3)}}{20 * 20 * 10 * N_T}, \quad C_{\phi, \psi'; n, n'}^{(3)} \equiv \sum_{t=1}^{N_T} \frac{\phi_{n;t}{}^{(3)} \cdot \psi'_{n';t}{}^{(3)}}{20 * 20 * 10 * N_T} \\ \alpha_{\phi, \psi'}^{(3)}(n, n') &\equiv \left\{ C_{\phi, \psi'}^{(3)} \left(C_{\psi', \psi'}^{(3)} \right)^{-1} \right\}_{nn'}. \end{aligned} \quad (4.4)$$

4.5 Summary

In summary, we have shown that the columns of the auto-covariance submatrices of \mathbf{B} implicates deep convection and that the cross-covariance submatrices are not horizontally invariant and isotropic. Following this variance and anisotropic idea, we argued that the assumption behind two commonly used regression coefficient models (Figures 1.1 and 1.2) do not hold. A new assumption was then formulated from realizing that within the vicinity of a specified point, auto-covariances tend to be approximately equal to variances. A new model was then formulated based on this realization.

Chapter 5

Results from regression coefficient models

Overview The three regression coefficient models are now applied onto the outputs of the GEN_BE system. The first model assumes that the regression coefficients are horizontally invariant and isotropic, and that there is no relationship between different model levels (Figure 1.1). The second model also assumes horizontal invariance, but does not assume that error variables are disconnected across different model levels (Figure 1.2). The third method simply follows the neighbourhood equivalence assumption (Figures 1.3 and 1.4). The performance of the modelled regression coefficients are evaluated by their ability to remove cross-covariances after applying the first step of the CVT. Another result of interest would be the amount of ϕ variance that the linear ϕ - ψ relationship can account for. If most of the variance of ϕ can be explained by the ϕ - ψ relationship, we can consider ϕ' to be completely decoupled from all other meteorological error variables in future works.

5.1 Performance

Basis of metric If a perfect set of regression coefficients can be produced, then the cross-covariance between variables in $\mathbf{x}^{b'}$ is zero¹. We can utilize this to determine the performance of the modelled regression coefficient matrices.

¹See the proof for Eqn (B.1) in Appendix B for the details.

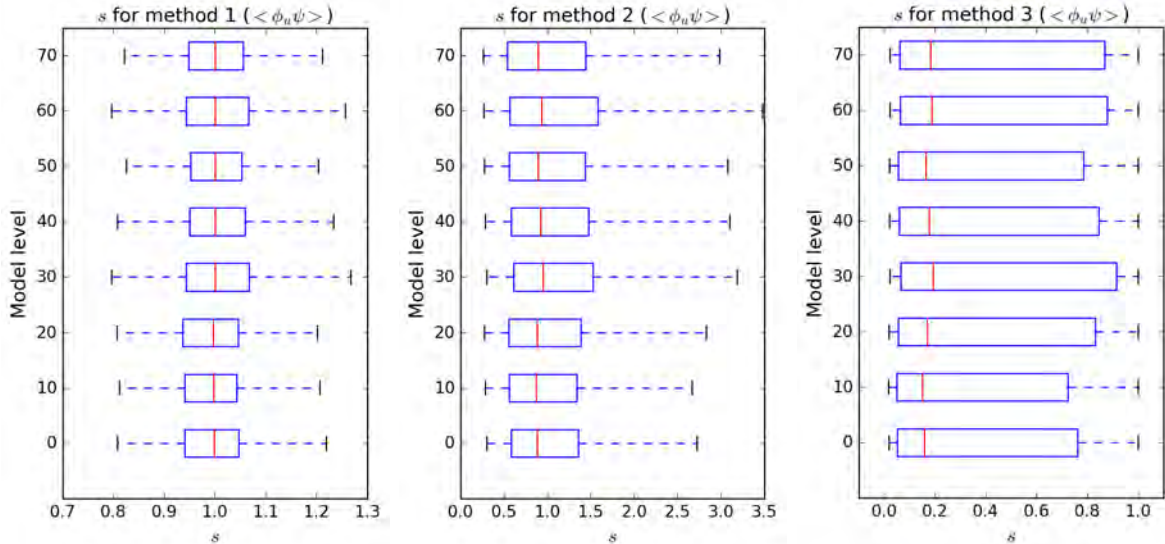


Figure 5.1: Box-and-whiskers plots of $s(\mathbf{r}, \mathbf{r}')$. These metrics were evaluated at the 72 \mathbf{r}' locations specified in the text, for the first model (left), the second model (middle) and the third model (right). The boxes bound the interquartile range, the red vertical line indicates the median and the whiskers bound the 10th and 90th percentiles of the $9 * 160 * 150$ evaluations at each of the 8 level.

Description of s We used a metric to evaluate the performance of the regression coefficient models pairs of points in the model. The mathematical formulation is

$$s(\mathbf{r}, \mathbf{r}') \equiv \left\| \sum_{t=1}^{N_T} \phi'(\mathbf{r}, t) \psi'(\mathbf{r}', t) \right\| / \left\{ \sum_{t=1}^{N_T} \phi(\mathbf{r}, t) \psi(\mathbf{r}', t) \right\}, \quad (5.1)$$

where \mathbf{r} and \mathbf{r}' are position vectors. This metric computes the ratio of the residual cross-covariance after the first step of CVT, against the original cross-covariance. The ideal value of s is thus zero. Naturally, the best model would be one that turns s to the smallest absolute value, when employed in the CVT.

s computer memory constraint Due to computational memory constraints, only 72 \mathbf{r}' were sampled. Likewise, \mathbf{r} was sampled every two grid points in the horizontal, and every grid point in the vertical. The 72 \mathbf{r}' positions are split across 8 model levels², 9 per model level. The horizontal positions of each set of 9 points are the same as the green dots in Figure 4.8.

Plotting the s data The s valuations from these 72 test locations are grouped according to the 8 model levels of the test locations and then plotted in Figure 5.1. In other words, each box and whisker plot in Figure 5.1 contains $9 * 160 * 150$ evaluations of s .

²Model level indices are $k = 0, 10, 20, \dots, 60, 70$.

Third model outperforms the first model According to Figure 5.1, 50% of the s evaluations for third model (below the median) have values less than 0.2. In other words, after applying the CVT, 50% of all tested points have cross-covariances with magnitudes that are less than 20% of the original cross-covariances. Furthermore, median to 90% percentile evaluations of s for the third model indicate that 40% of the new cross-covariances range from 20% to 100% of the original cross-covariances. In contrast, 40% of the s evaluations for the first model (10th percentile to median) indicate that the absolute value of the new cross-covariances are more than 80% of the original. To make matters worse, 40% of the s evaluations for the first model (median to 90th percentile) indicate that applying the first model resulted in larger cross-covariances! Clearly, the third model is far more adept at removing cross-covariances than the first model.

Third model outperforms the second model Similarly, the third model outperforms the second model. The 10th percentile to median evaluations of s for the second model indicates that applying the second model can only remove less than 80% of the cross-covariances for 40% of the points. In contrast, for the same range, the third model removes more than 80%. To make matters worse for the second model, the median to 90th percentile range of s evaluations indicate that applying the second model amplifies the cross-covariances! For the same percentile range, the third model is still able to constrain the absolute values of the new cross-covariances to be smaller than that of the original cross-covariances. The third model is clearly more suited to removing cross-covariances than the second model.

First and second model are unsuitable for our region The first model is clearly unsuitable for our region as it only modifies the cross-covariances by 20%, and it does amplify cross-covariances for 50% of the evaluated cases. While the second model is able to remove between more than 40% of the cross-covariances for 25% of the points tested (below the 25th percentile in Figure 5.1), the second model actually amplifies cross-covariances for roughly half the tested points (above the medians in Figure 5.1). Furthermore, 25% of the tested points (upwards of the 75th percentile) indicate that the second model actually amplifies the cross-covariances by more than 30%, and 10% of the tested points (upward of the 90th percentile) indicate that the second model actually increased the cross-covariance by more than 150%. These results mean that the first and second model are highly unsuitable for use in the CVT for our region of interest.

Violation of assumptions is the reason for first and second models' performance As per the discussion in the preceding chapter, it is hardly surprising that the

performance of the first and second models is terrible in our region. In the preceding chapter, we observed from Figure 4.6 that cross-covariances vary tremendously with the horizontal position considered and show great anisotropy around said position. As explained, this violates the horizontal invariance and isotropy assumption that underlies the first and second model. In other words, the bad performance by the two models are as expected.

5.2 Variance occupation

An additional boon for third model Of course, if the third model can also account for most of the variance of ϕ , future work utilizing this model can also neglect relationships between ϕ' and other meteorological variables. To examine this, we defined the following fraction

$$f(\mathbf{r}) \equiv \left\{ \sum_{t=1}^{N_t} \phi(\mathbf{r}) \phi_b(\mathbf{r}) \right\} / \left\{ \sum_{t=1}^{N_t} \phi(\mathbf{r}) \phi(\mathbf{r}) \right\}, \quad (5.2)$$

where

$$\phi_b \equiv \alpha_{\phi, \psi} \psi'.$$

In essence, if f is close to 1, then the linear relationship between ϕ and ψ is able to account for most of the variance in ϕ . The f evaluated for all three methods, across the model grid, are aggregated by model levels and displayed as box-and-whisker plots in Figure 5.2.

Third model explains most of ϕ variance Two observations can be made from Figure 5.2. First of all, from the horizontal scales of Figure 5.2, it is clear that utilizing the third model results in the greatest f valuations (which are very close to 1). In other words, when the third model is used, most of the variance in ϕ can be explained by ψ . In contrast, the other two methods only explain a small portion of the variance. As such, if the third model is implemented, subsequent studies can simply consider error variables apart from ψ to be largely unrelated to ϕ .

Fall-off in third model's f at upper level The second observation is that in the upper levels, there is a notable fall in the f valuations of the third model. This is most likely because the subdomains only cover the bottom 70 levels of the model and neglect the remaining 9 levels.

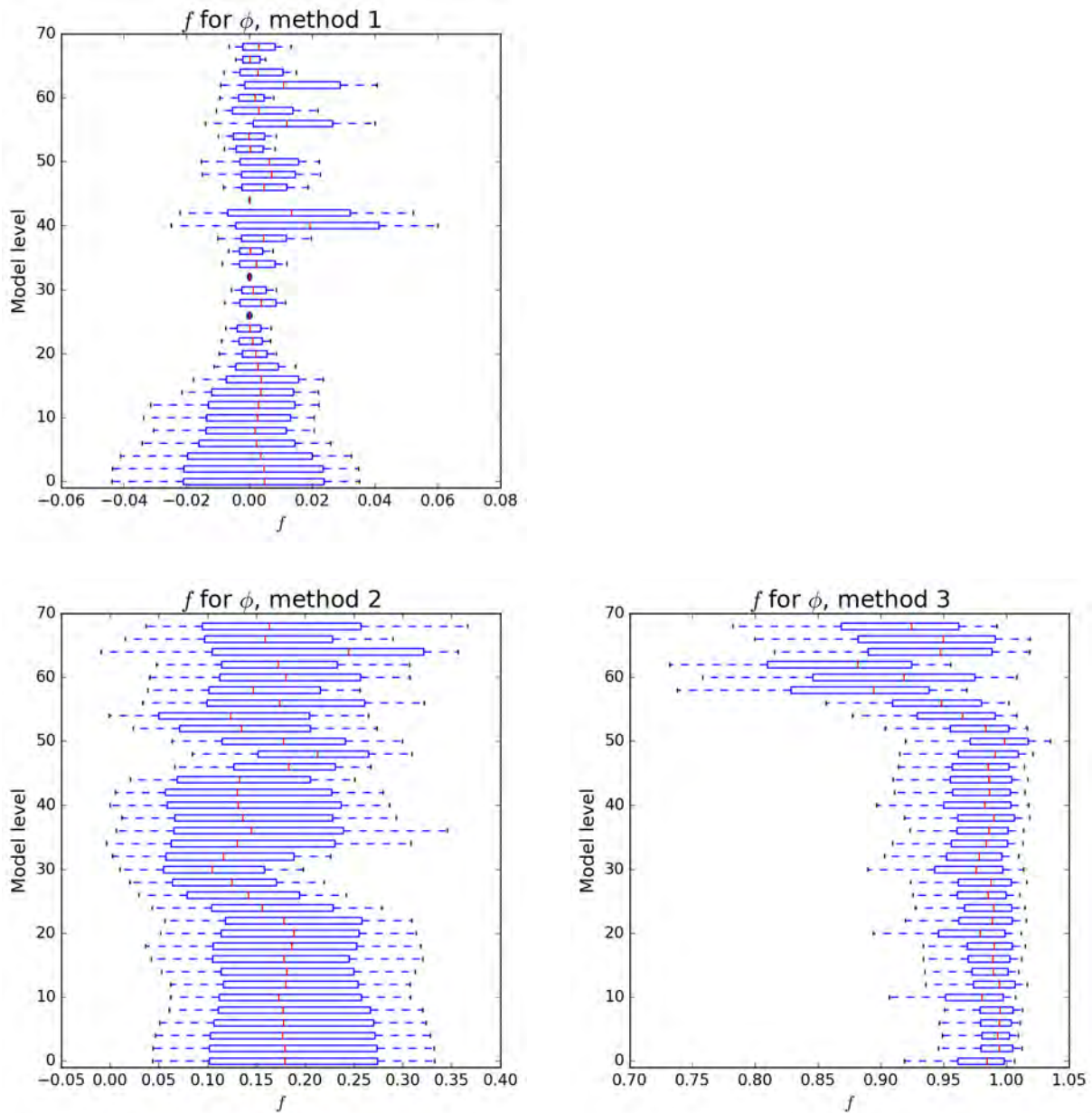


Figure 5.2: Box-and-whiskers diagrams of f obtained by the three methods of modelling regression coefficients. The boxes bound the interquartile range, the red vertical line indicates the median and the whiskers bound the 10th and 90th percentiles of the $9 * 160 * 150$ evaluations at each of the 8 level.

5.3 Summary

In summary, because the assumptions underlying the two commonly used models do not hold in reality, while the third model’s assumption is based off the features observed in \mathbf{B} , the third model shows the greatest ability to remove cross-covariances. Furthermore, the third model has the additional benefit of accounting for most of ϕ ’s variance, meaning that ϕ ’ can be essentially considered to be decoupled from all other error variables introduced by future work.

Chapter 6

Conclusions and future work

Commonly used models are inadequate The goal of this project is to determine an appropriate model for the regression coefficient matrices involved in the CVT, over the Maritime Continent, during the Northeast monsoon season, using forecasts from the Singapore Variable resolution model. After sanity checking SINGV forecasts (5th December 2015 to 19th December 2015), we examined several columns of \mathbf{B} and realized that the cross-covariance between ϕ and ψ exhibits large variations depending on the horizontal location of the reference points. This means that any regression coefficient model based on the assumption that regression coefficients are invariant in the horizontal are unlikely to work well in the Maritime Continent, during the Northeast monsoon. Later tests confirmed their inadequate performance in this scenario.

Neighbourhood equivalence assumption and the third model is the best Examinations of the columns of \mathbf{B} also revealed that there is a tendency for the auto-covariances between evaluations of an error variable at two nearby points to be approximately equal to their variances. A third model was formulated based on this assumption (neighbourhood equivalence assumption), and then tested. When utilized in a CVT, the third model removed most of the cross-covariances ϕ and ψ , making it an adequate model for our scenario, fulfilling our goal.

Third model's affordable computational cost It must be said here that the third model has a much higher computational cost than the other two. This is to be expected: the regression coefficient matrix generated by the third model is much bigger than the first and second model. However, this cost is still affordable! We were able generate the regression coefficients within 6 hours of computation in a laptop with an Intel i7 quad-core and 16 gigabytes of Random Access Memory, without using parallel computing

(only 1 core was employed). Furthermore, when we ran a CVT based on the third model, using the same laptop, without parallel programming, the first step of the CVT took only 3 minutes. Also, the algorithm that utilizes the third model for the CVT is actually an embarrassingly parallel algorithm (many repetitive matrix multiplications). It seems likely that if the model was rewritten for parallel computing, it would be affordable for data assimilation in operational forecasting.

Future work: increase number of subdomains in third model Even though our new model has much better performance than the other two, it is not optimal yet. In Chapter 4, we mentioned that the neighbourhood equivalence assumption appears to hold for 5 model levels. However, we utilized subdomains with 10 model levels each due to computational time limitations. Future work on this model can further enhance the performance of the model by using half the number of model levels per subdomain.

Future work: apply neighbourhood equivalence assumption to CVT's second step Aside from subdomain size optimization, another avenue of future work would be to apply the neighbourhood equivalence assumption to the second step of the CVT: the removal of auto-covariances in $\mathbf{x}^{b'}$. We can safely say from Figure 4.5 that the assumption should hold fairly well when dealing with the auto-covariance of ψ as it is unchanged by the first step of the CVT. However, work is needed to determine if the assumption does hold for ϕ' . If that is true, it may be feasible to model the auto-covariance of $\mathbf{x}^{b'}$ using the eigenvectors of the auto-covariance matrix of the subdomains. Future work can then test this new model against other pre-existing models of auto-covariances.

Future work: generality of neighbourhood equivalence assumption and third model A third avenue of future work would be to determine whether the neighbourhood equivalence assumption can be utilized on other commonly used variables in data assimilation (*e.g.*, air temperature and relative humidity). If so, then it may be possible to utilize the third model when removing their cross-covariances through the CVT.

Outlook for the future Advancements in computer technology have worked wonders for operational forecasting. In a mere three decades, the resolution of global circulation forecast models has improved from several hundreds of kilometres to the current less than 20 km resolution (the European Centre for Medium Range Weather Forecasts currently uses a 9 km horizontal resolution in their 36-hour global forecasts). In fact, at the time of writing, numerical weather prediction and limited area models with a horizontal resolution of less than 5 km are fairly common. It would not be surprising if sub-convective scale

(sub-kilometre scale) forecast models become common by the end of the next decade. However, it is difficult to expect the density of observation to match up to the resolution of forecast models. Clearly, \mathbf{B} 's capacity for spreading the influence of observations from their sparse locations throughout a model's initial conditions during data assimilation will prove to be important in this future. The art of data assimilation and studies on modelling \mathbf{B} will remain invaluable for a long time.

Bibliography

- Australian Government, Bureau of Meteorology. 2016. *Analysis chart archive*.
- Badan Perencanaan Pembangunan Nasional. 2013. *Indonesian population projection*.
- Bannister, R. N. 2008a. *A review of forecast error covariance statistics in atmospheric variational data assimilation. I: Characteristics and measurements of forecast error covariances*, Quarterly Journal of the Royal Meteorological Society **134**, no. 637, 1951–1970, available at [arXiv:0801.1618v2](https://arxiv.org/abs/0801.1618v2).
- . 2008b. *A review of forecast error covariance statistics in atmospheric variational data assimilation. II: Modelling the forecast error covariance statistics*, Quarterly Journal of the Royal Meteorological Society **134**, no. 637, 1971–1996.
- Berre, L., S. Ecaterina Stefanescu, and M. Belo Pereira. 2006. *The representation of the analysis effect in three error simulation techniques*, Tellus A **58**, no. 2, 196–209.
- Bijlsma, S. J., L. M. Hafkenscheid, and Peter Lynch. 1986. *Computation of the streamfunction and velocity potential and reconstruction of the wind field*, Monthly Weather Review **114**, no. 8, 1547–1551.
- Bouttier, F. and P. Courtier. 2002. *Data assimilation concepts and methods*, Meteorological training course lecture series.
- Caron, J-F. and L. Fillion. 2010. *An examination of background error correlations between mass and rotational wind over precipitation regions*, Monthly Weather Review **138**, no. 2, 563–578.
- Chang, C-P., P. a. Harr, and H-J. Chen. 2005. *Synoptic disturbances over the equatorial South China Sea and western Maritime Continent during boreal winter*, Monthly Weather Review **133**, no. 3, 489–503.
- Chang, C-P., M-M. Lu, and H. Lim. 2016. *Monsoon convection in the Maritime Continent: interaction of large-scale motion and complex terrain*, Meteorological Monographs - Multiscale Convection-Coupled Systems in the Tropics: A Tribute to Dr. Michio Yanai **56**, 6.1–6.29.
- Chen, Y., S. R. H. Rizvi, X. Y. Huang, J. Min, and X. Zhang. 2013. *Balance characteristics of multivariate background error covariances and their impact on analyses and forecasts in tropical and Arctic regions*, Meteorology and Atmospheric Physics **121**, no. 1-2, 79–98.
- Daley, R. 1991. *Atmospheric data analysis*, Vol. 15, Cambridge University Press.
- Descombes, G., T. Auligné, F. Vandenberghe, D. M. Barker, and J. Barré. 2015. *Generalized background error covariance matrix model (GEN-BE v2.0)*, Geoscientific Model Development **8**, no. 3, 669–696.
- Ding, Y. 1994. *The summer monsoon in East Asia*, Monsoons over china, pp. 1–90.
- Fong, M. and L. K. Ng. 2012. *The weather and climate of Singapore*, Meteorological Service Singapore, Singapore.

Kalnay, E. 2003. *Atmospheric modelling, data assimilation and predictability*, Vol. 129, Cambridge University Press.

Lanczos, C. 1957. *Applied analysis*, Dover Publications.

Lorenc, A. C. 1986. *Analysis methods for numerical weather prediction*, Quart J. R. Mrt. Soc. **112**, 1177–1194.

Parrish, D. F. and J. C. Derber. 1992. *The National Meteorological Center's Spectral Statistical-Interpolation Analysis System*, Vol. 120.

Wang, B. 2006. *The Asian monsoon*, Springer.

Wu, W-S., R. J. Purser, and D. F. Parrish. 2002. *Three-dimensional variational analysis with spatially inhomogeneous covariances*, Monthly Weather Review **130**, no. 12, 2905–2916.

Appendix A

The roles of the background error covariance matrix in DA

Overview Here, the roles of \mathbf{B} in data assimilation (mentioned in Chapter 1) will be shown through the single-observation and double-observation assimilation solutions that minimize Eqn (1.4). As such, this chapter will be divided into two sections: the derivation of the solutions to Eqn (1.4) and the application of said solutions to demonstrate said roles of \mathbf{B} . Afterwards, this chapter will conclude with a brief summary. Note that the single-observation solution and several roles of \mathbf{B} in data assimilation will be referenced in later chapters of this document to interpret columns of \mathbf{B} .

A.1 Solutions to Cost Function

A.1.1 Best linear unbiased estimate solution

Starting point – cost function As mentioned in Chapter 1, goal of data assimilation is to determine a model state that minimizes Eqn (1.4). In the field’s parlance, this model state is called the analysis state, \mathbf{X}^a . \mathbf{X}^a will be derived in this section. For convenience, the cost function is rewritten in tensor notation,

$$J(\mathbf{X}) = X_i B_{ij}^{-1} X_j + X_i^b B_{ij}^{-1} X_j^b - X_i B_{ij}^{-1} X_j^b - X_i^b B_{ij}^{-1} X_j \\ + Y_m R_{mn}^{-1} Y_n + H(\mathbf{X})_m R_{mn}^{-1} H(\mathbf{X})_n - Y_m R_{mn}^{-1} H(\mathbf{X})_n - H(\mathbf{X})_m R_{mn}^{-1} Y_n.$$

Consider the gradient of the cost function:

$$\begin{aligned} \frac{\partial J(\mathbf{X})}{\partial x_l} &= B_{lj}^{-1} X_j + X_i B_{il}^{-1} - B_{lj}^{-1} X_j^b - X_i^b B_{il}^{-1} \\ &+ \frac{\partial}{\partial x_l} \left\{ H(\mathbf{X})_m R_{mn}^{-1} H(\mathbf{X})_n - Y_m R_{mn}^{-1} H(\mathbf{X})_n - H(\mathbf{X})_m R_{mn}^{-1} Y_n \right\}. \end{aligned}$$

Since $\mathbf{X} = \mathbf{X}^a$ minimizes the cost function, then,

$$\begin{aligned} \left. \frac{\partial J(\mathbf{X})}{\partial x_l} \right|_{\mathbf{X}=\mathbf{X}^a} &= 0 = B_{lj}^{-1} X_j^a + X_i^a B_{il}^{-1} - B_{lj}^{-1} X_j^b - X_i^b B_{il}^{-1} \\ &+ \left. \frac{\partial H(\mathbf{X})_m}{\partial x_l} \right|_{\mathbf{X}=\mathbf{X}^a} R_{mn}^{-1} H(\mathbf{X}^a)_n + H(\mathbf{X}^a)_m \left. \frac{\partial H(\mathbf{X})_n}{\partial x_l} \right|_{\mathbf{X}=\mathbf{X}^a} \\ &- Y_m R_{mn}^{-1} \left. \frac{\partial H(\mathbf{X})_n}{\partial x_l} \right|_{\mathbf{X}=\mathbf{X}^a} - \left. \frac{\partial H(\mathbf{X})_m}{\partial x_l} \right|_{\mathbf{X}=\mathbf{X}^a} R_{mn}^{-1} Y_n. \end{aligned}$$

Symmetric \mathbf{B}^{-1} and \mathbf{R}^{-1} Before going any further, it is useful to note that \mathbf{B} and \mathbf{R} are, by definition, symmetric and positive definite matrices. In other words, \mathbf{B} and \mathbf{R} are invertible and their inverses are symmetric. Applying these ideas condenses the cost function differential:

$$\begin{aligned} \left. \frac{\partial J(\mathbf{X})}{\partial x_l} \right|_{\mathbf{X}=\mathbf{X}^a} &= 0 = 2 B_{lj}^{-1} X_j^a - 2 B_{lj}^{-1} X_j^b + 2 \left\{ \left. \frac{\partial H(\mathbf{X})_m}{\partial x_l} \right|_{\mathbf{X}=\mathbf{X}^a} R_{mn}^{-1} H(\mathbf{X}^a)_n \right\} \\ &- 2 \left\{ Y_m R_{mn}^{-1} \left. \frac{\partial H(\mathbf{X})_n}{\partial x_l} \right|_{\mathbf{X}=\mathbf{X}^a} \right\}. \end{aligned}$$

Gradient of H Defining

$$\mathbf{H} \equiv \left. \nabla H(\mathbf{X}) \right|_{\mathbf{X}=\mathbf{X}^a} \quad (\text{A.1})$$

simplifies the problem even further:

$$B_{lj}^{-1} X_j^a - B_{lj}^{-1} X_j^b + H_{lm} R_{mn}^{-1} H(\mathbf{X}^a)_n - Y_m R_{mn}^{-1} H_{nl} = 0.$$

Note that the gradient operator in Eqn (A.1) is in model space. In other words, for a N -element state vector and a M -element observation vector, \mathbf{H} is a $M \times N$ matrix. Thus,

$$\begin{aligned} B_{lj}^{-1} X_j^a - B_{lj}^{-1} X_j^b &= + Y_m R_{mn}^{-1} H_{nl} - H_{lm} R_{mn}^{-1} H(\mathbf{X}^a)_n \\ \implies \mathbf{B}^{-1} (\mathbf{X}^a - \mathbf{X}^b) &= \mathbf{H}^\top \mathbf{R}^{-1} \{\mathbf{Y} - H(\mathbf{X}^a)\} \end{aligned}$$

Tangent linear hypothesis To proceed any further, it is assumed that \mathbf{X}^a and \mathbf{X}^b are close enough for H to be linear. In other words, for $\mathbf{X}^a + \delta\mathbf{X} = \mathbf{X}^b$,

$$H(\mathbf{X}^a + \delta\mathbf{X} - \delta\mathbf{X}) \approx H(\mathbf{X}^a + \delta\mathbf{X}) - \mathbf{H}\delta\mathbf{X} = H(\mathbf{X}^b) + \mathbf{H}\mathbf{X}^a - \mathbf{H}\mathbf{X}^b. \quad (\text{A.2})$$

This assumption is called the ‘‘tangent linear hypothesis’’ in literature. Thus,

$$\begin{aligned} \mathbf{B}^{-1}(\mathbf{X}^a - \mathbf{X}^b) &= \mathbf{H}^\top \mathbf{R}^{-1} \{ \mathbf{Y} - H(\mathbf{X}^a + \delta\mathbf{X} - \delta\mathbf{X}) \} \\ \implies \mathbf{B}^{-1}(\mathbf{X}^a - \mathbf{X}^b) &= \mathbf{H}^\top \mathbf{R}^{-1} \{ \mathbf{Y} - H(\mathbf{X}^b) - \mathbf{H}\mathbf{X}^a + \mathbf{H}\mathbf{X}^b \} \\ \implies (\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H}) (\mathbf{X}^a - \mathbf{X}^b) &= \mathbf{H}^\top \mathbf{R}^{-1} \{ \mathbf{Y} - H(\mathbf{X}^b) \} \end{aligned}$$

Assume invertibility of $(\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})$ To proceed any further, it is necessary to assume that $(\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})$ is invertible. Thus:

$$(\mathbf{X}^a - \mathbf{X}^b) = (\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{R}^{-1} [\mathbf{Y} - H(\mathbf{X}^b)]$$

Tidying up In the literature, \mathbf{K} is usually written as $\mathbf{B}\mathbf{H}^\top (\mathbf{H}\mathbf{B}\mathbf{H}^\top + \mathbf{R})^{-1}$. For completeness, the two forms of \mathbf{K} can be shown to be equivalent. Suppose that:

$$\mathbf{B}\mathbf{H}^\top (\mathbf{H}\mathbf{B}\mathbf{H}^\top + \mathbf{R})^{-1} = (\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{R}^{-1}$$

Then,

$$\begin{aligned} \mathbf{B}\mathbf{H}^\top &= (\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{R}^{-1} (\mathbf{H}\mathbf{B}\mathbf{H}^\top + \mathbf{R}) \\ &= (\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})^{-1} (\mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H}\mathbf{B} + \mathbf{I}) \mathbf{H}^\top \\ &= (\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})^{-1} (\mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H} + \mathbf{B}^{-1}) \mathbf{B}\mathbf{H}^\top \\ &= \mathbf{B}\mathbf{H}^\top \end{aligned}$$

Thus, the two forms of \mathbf{K} are consistent with each other. For consistency with the literature, the following form of the solution will be used for the rest of this document¹:

$$\mathbf{X}^a = \mathbf{X}^b + \mathbf{B}\mathbf{H}^\top (\mathbf{H}\mathbf{B}\mathbf{H}^\top + \mathbf{R})^{-1} [\mathbf{Y} - H(\mathbf{X}^b)]. \quad (\text{A.3})$$

¹An alternative derivation of the solution, based on minimizing the error of the analysis state, can be found in the appendix.

A.1.2 Single and double-observation BLUE solutions

Coincidence assumption The roles of \mathbf{B} in data assimilation can be directly observed from considering the solution in Eqn (A.3) for the case where only one observation is assimilated, and for the case where two observations are assimilated. To simplify the problem, it will be assumed that the observation(s) is located at a model grid point(s), and concerns one (two) of the variables in \mathbf{B} .

H under coincidence assumption Under this collocation assumption, the operator H is effectively the matrix \mathbf{H} . Let the normalized basis vector for that model point and position in model basis be $\hat{\mathbf{X}}_n$ (*i.e.*, the n -th element on \mathbf{X}) and likewise, the observation basis vector be $\hat{\mathbf{Y}}_m$. As such, \mathbf{H} for the single and double-observation cases are,

$$\mathbf{H} = 1\hat{\mathbf{Y}}\hat{\mathbf{X}}_n^\top, \text{ and, } \mathbf{H} = 1\hat{\mathbf{Y}}_1\hat{\mathbf{X}}_{n_1}^\top + 1\hat{\mathbf{Y}}_2\hat{\mathbf{X}}_{n_2}^\top, \quad (\text{A.4})$$

respectively.

Single-observation solution For the single-observation case, it is clear that \mathbf{R} consists of a single element (σ^2). Applying the single-observation \mathbf{R} and \mathbf{H} into Eqn (A.3) yields the single-observation data assimilation equation

$$X_l^a - X_l^b = B_{ln} \frac{Y - X_n^b}{B_{nn} + \sigma^2}. \quad (\text{A.5})$$

Note that subscripts refer to the corresponding elements in the relevant vector/matrix.

Double-observation solution For convenience, suppose that the two observations for the double-observation cases are uncorrelated. Then, $\mathbf{R} = \sigma_1^2 \hat{\mathbf{Y}}_1 \hat{\mathbf{Y}}_1^\top + \sigma_2^2 \hat{\mathbf{Y}}_2 \hat{\mathbf{Y}}_2^\top$. It can be shown that the corresponding BLUE solution is

$$\begin{aligned} X_l^a - X_l^b &= \frac{B_{ln_1}}{\gamma} (y_1 - X_{n_1}^b) + \frac{B_{ln_2}}{\gamma} (y_2 - X_{n_2}^b) \\ &\quad + \frac{B_{ln_1}B_{n_2n_2} - B_{n_1n_2}B_{ln_2}}{\gamma} (Y_1 - X_{n_1}^b) + \frac{B_{ln_2}B_{n_1n_1} - B_{ln_1}B_{n_1n_2}}{\gamma} (Y_2 - X_{n_2}^b) \end{aligned} \quad (\text{A.6})$$

where $\gamma \equiv B_{n_1n_1}B_{n_2n_2} + B_{n_1n_1}\sigma_2^2 + B_{n_2n_2}\sigma_1^2 + \sigma_1^2\sigma_2^2 - B_{n_1n_2}^2$.

Both solutions match It can also be shown that the double-observation solution can be reduced to the single-observation solution when either σ_1^2 or σ_2^2 is set to infinity. This makes sense as when there is infinite observation error, the observation's impact on data assimilation should be zero.

A.2 Roles of \mathbf{B} in data assimilation

\mathbf{B} weighs \mathbf{Y} and \mathbf{X}^b locally The first role of \mathbf{B} in data assimilation can be illustrated using Eqn (A.5). Expressing and rearranging Eqn (A.5) at the n -th element of the state vector gives,

$$X_n^a = \frac{B_{nn}}{B_{nn} + \sigma^2} Y + \frac{\sigma^2}{B_{nn} + \sigma^2} X_n^b. \quad (\text{A.7})$$

Clearly, \mathbf{B} serves to weigh the contributions of the observation and background state, according to the relative confidence for the two pieces of information.

\mathbf{B} spreads $\mathbf{Y} - H(\mathbf{X}^b)$ globally Aside from local weighing, \mathbf{B} also spreads the difference between the observation and background state, throughout the model grid and across model variables. This can be inferred by considering Eqn (A.5) and the general form of \mathbf{B} for L variables:

$$\mathbf{B} = \begin{bmatrix} \langle \mathbf{x}^b_1 \mathbf{x}^b_1^\top \rangle & \langle \mathbf{x}^b_1 \mathbf{x}^b_2^\top \rangle & \cdots & \langle \mathbf{x}^b_1 \mathbf{x}^b_l^\top \rangle & \cdots & \langle \mathbf{x}^b_1 \mathbf{x}^b_{L-1}^\top \rangle & \langle \mathbf{x}^b_1 \mathbf{x}^b_L^\top \rangle \\ \langle \mathbf{x}^b_2 \mathbf{x}^b_1^\top \rangle & \langle \mathbf{x}^b_2 \mathbf{x}^b_2^\top \rangle & \cdots & \langle \mathbf{x}^b_2 \mathbf{x}^b_l^\top \rangle & \cdots & \langle \mathbf{x}^b_2 \mathbf{x}^b_{L-1}^\top \rangle & \langle \mathbf{x}^b_2 \mathbf{x}^b_L^\top \rangle \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ \langle \mathbf{x}^b_l \mathbf{x}^b_1^\top \rangle & \langle \mathbf{x}^b_l \mathbf{x}^b_2^\top \rangle & \cdots & \langle \mathbf{x}^b_l \mathbf{x}^b_l^\top \rangle & \cdots & \langle \mathbf{x}^b_l \mathbf{x}^b_{L-1}^\top \rangle & \langle \mathbf{x}^b_l \mathbf{x}^b_L^\top \rangle \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ \langle \mathbf{x}^b_{L-1} \mathbf{x}^b_1^\top \rangle & \langle \mathbf{x}^b_{L-1} \mathbf{x}^b_2^\top \rangle & \cdots & \langle \mathbf{x}^b_{L-1} \mathbf{x}^b_l^\top \rangle & \cdots & \langle \mathbf{x}^b_{L-1} \mathbf{x}^b_{L-1}^\top \rangle & \langle \mathbf{x}^b_{L-1} \mathbf{x}^b_L^\top \rangle \\ \langle \mathbf{x}^b_L \mathbf{x}^b_1^\top \rangle & \langle \mathbf{x}^b_L \mathbf{x}^b_2^\top \rangle & \cdots & \langle \mathbf{x}^b_L \mathbf{x}^b_l^\top \rangle & \cdots & \langle \mathbf{x}^b_L \mathbf{x}^b_{L-1}^\top \rangle & \langle \mathbf{x}^b_L \mathbf{x}^b_L^\top \rangle \end{bmatrix}. \quad (\text{A.8})$$

\mathbf{X}^a is the sum of \mathbf{X}^b and a multiple of a column of \mathbf{B} . Clearly, any column of \mathbf{B} cuts through an auto-covariance matrix, and $L - 1$ cross-covariance matrices. The auto-covariance matrix's column spreads the observation-background difference throughout the model grid, for the assimilated variable. At the same time, the columns of the $L - 1$ cross-

covariance matrices spreads the observation-background difference to different variables, across the entire model grid. In other words, \mathbf{B} spreads the influence of $\mathbf{Y} - H(\mathbf{X}^b)$ throughout the entire set of initial conditions.

\mathbf{B} overlaps the influences of different observations When multiple observations are assimilated, \mathbf{B} allows their respective influences to overlap, as though two single-observation solutions were used. This can be inferred from the first line of the double-observation solution, Eqn (A.6):

$$\frac{B_{ln_1}}{\gamma} (y_1 - X_{n_1}^b) + \frac{B_{ln_2}}{\gamma} (y_2 - X_{n_2}^b).$$

Notice that for a N -element state vector, for all $l \in [1, N]$, $l \in \mathbb{Z}$, these terms do not cancel out or vanish. As such, the two observation-background differences are included in every point on the grid and across all variables, weighed by the corresponding columns of \mathbf{B} . In other words, \mathbf{B} causes the influence of observations to overlap in the model grid.

\mathbf{B} causes observation interactions Aside from the overlapping effect, \mathbf{B} also causes the observation-background differences to interaction at grid points and/or variables that differ from the observation location and/or variables. The second line of the double-observation BLUE solution in Eqn (A.6) alludes to this effect:

$$\frac{B_{ln_1}B_{n_2n_2} - B_{n_1n_2}B_{ln_2}}{\gamma} (y_1 - X_{n_1}^b) + \frac{B_{ln_2}B_{n_1n_1} - B_{ln_1}B_{n_1n_2}}{\gamma} (y_2 - X_{n_2}^b)$$

In other words, at points/variables different from the observations, an additional interaction between the two observation-background differences pop up. \mathbf{B} supplies the interaction itself.

\mathbf{B} limits observation interactions It is also interesting to note that the interacting term from one observation vanishes at the location/variable of the other observation. Specifically, $B_{ln_1}B_{n_2n_2} - B_{n_1n_2}B_{ln_2} = 0$ when $l = n_2$, and, $B_{ln_2}B_{n_1n_1} - B_{ln_1}B_{n_1n_2} = 0$ when $l = n_1$. This means that at the position/variable of one observation, the only influence from the other observation comes from the overlapping part. In other words, \mathbf{B} limits observation interactions to places/variables away from observation points/variables.

A.3 Summary of derivation

In summary, a general solution to the cost function of Eqn (1.4) have been derived under the tangent linear hypothesis and the assumption that $(\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})$ is invertible. Two specific solutions (single-observation and double-observation) were then derived from the general solution, and then utilized to demonstrate that \mathbf{B} has the following roles in data assimilation:

1. \mathbf{B} weighs \mathbf{Y} and \mathbf{X}^b locally,
2. \mathbf{B} spreads $\mathbf{Y} - H(\mathbf{X}^b)$ globally,
3. \mathbf{B} overlaps influences of different observations,
4. \mathbf{B} causes observation interactions, and,
5. \mathbf{B} limits observation interactions.

In the later chapters, the first two roles of \mathbf{B} and the single-observation solutions will be referenced when interpreting the results of this project's computations.

Appendix B

Details about the first step of CVT

Purpose of CVT As described in the introduction, the goal of the CVT method is to transform the inputted state vector of background error variables (\mathbf{x}^b) into a vector of control variables ($\mathbf{x}^{b''}$) possessing zero covariances. When this happens, \mathbf{B} is effectively transformed into an easily managed diagonal matrix.

Diagonalisation inspiration The CVT model actually emulates the effect of applying \mathbf{B} 's eigenvector matrices to itself, *i.e.*, linear algebra's diagonalisation. Consider:

$$\mathbf{B} = \mathbf{V} \mathbf{S} \mathbf{V}^\top$$

where \mathbf{V} and \mathbf{S} are matrices of the same dimensions as \mathbf{B} , the columns of \mathbf{V} contain the eigenvectors of \mathbf{B} and \mathbf{S} is a diagonal matrix of the eigenvalues. In other words, if one defines $\mathbf{x}^{b''} \equiv \mathbf{V}^\top \mathbf{x}^b$, *i.e.*, $\mathbf{V} \mathbf{x}^{b''} = \mathbf{x}^b$, then,

$$\langle \mathbf{x}^{b''} \mathbf{x}^{b''\top} \rangle = \langle \mathbf{V}^\top \mathbf{x}^b \mathbf{x}^{b\top} \mathbf{V} \rangle = \mathbf{V}^\top \langle \mathbf{x}^b \mathbf{x}^{b\top} \rangle \mathbf{V} = \mathbf{V}^\top \mathbf{V} \mathbf{S} \mathbf{V}^\top \mathbf{V} = \mathbf{S}.$$

In other words, the best CVT operator is none other than \mathbf{V}^\top . However, given that \mathbf{B} must be fully determined before any computation of \mathbf{V} , this optimal operator cannot be used.

Formulation of \mathbf{U}_p Instead, the CVT method seeks to emulate the diagonalizing effect of \mathbf{V} through the two stages described in the introduction. We will focus on the cross-

covariance removal stage, whose operator is \mathbf{U}_p ¹. The transformation is simply

$$\mathbf{x}^b = \mathbf{U}_p \mathbf{x}^{b'},$$

where \mathbf{x}^b and $\mathbf{x}^{b'}$ are as defined in Chapter 1. It should be apparant that

$$\mathbf{U}_p = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \boldsymbol{\alpha}_{1,2} & \mathbf{I} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \boldsymbol{\alpha}_{1,L-1} & \boldsymbol{\alpha}_{2,L-1} & \cdots & \mathbf{I} & \mathbf{0} \\ \boldsymbol{\alpha}_{1,L} & \boldsymbol{\alpha}_{2,L} & \cdots & \boldsymbol{\alpha}_{L-1,L} & \mathbf{I} \end{bmatrix}.$$

This also means that

$$\mathbf{B} = \mathbf{U}_p \langle \mathbf{x}^{b'} \mathbf{x}^{b'\top} \rangle \mathbf{U}_p^\top.$$

Decorrelated $\mathbf{x}^{b'}$ If $\mathbf{x}^{b'}$ has zero cross-covariances (over time), then $\langle \mathbf{x}^{b'} \mathbf{x}^{b'\top} \rangle$ is a block diagonal matrix. We will prove this zero cross-covariance property,

$$\langle \mathbf{x}^{b'_l} \mathbf{x}^{b'_m\top} \rangle = 0 \quad \forall l \neq m. \quad (\text{B.1})$$

First layer of proof for Eqn (B.1) The proof for Eqn (B.1) will be done inductively. First, consider the scenario where there are only 2 background error variables, then $\mathbf{x}^{b'_1} = \mathbf{x}^{b_1}$, $\mathbf{x}^{b'_2} = \mathbf{x}^{b_2} - \boldsymbol{\alpha}_{2,1} \mathbf{x}^{b'_1}$, and, $\boldsymbol{\alpha}_{2,1} \equiv \langle \mathbf{x}^{b_2} \mathbf{x}^{b'_1\top} \rangle \langle \mathbf{x}^{b'_1} \mathbf{x}^{b'_1\top} \rangle^{-1}$. Thus,

$$\begin{aligned} \langle \mathbf{x}^{b'_2} \mathbf{x}^{b'_1\top} \rangle &= \langle \mathbf{x}^{b_2} \mathbf{x}^{b'_1\top} \rangle - \langle \boldsymbol{\alpha}_{2,1} \mathbf{x}^{b'_1} \mathbf{x}^{b'_1\top} \rangle \\ &= \langle \mathbf{x}^{b_2} \mathbf{x}^{b'_1\top} \rangle - \langle \mathbf{x}^{b_2} \mathbf{x}^{b'_1\top} \rangle \langle \mathbf{x}^{b'_1} \mathbf{x}^{b'_1\top} \rangle^{-1} \langle \mathbf{x}^{b'_1} \mathbf{x}^{b'_1\top} \rangle = \mathbf{0}. \end{aligned}$$

Since $\langle \mathbf{x}^{b'_1} \mathbf{x}^{b'_2\top} \rangle = \langle \mathbf{x}^{b'_2} \mathbf{x}^{b'_1\top} \rangle^\top$, then $\langle \mathbf{x}^{b'_1} \mathbf{x}^{b'_2\top} \rangle = \mathbf{0}$. As such, Eqn (B.1) holds for the 2 background error variable case.

Inductive proof for Eqn (B.1) Now, assume that Eqn (B.1) holds for a set of n transformed background error variables and an additional transformed background error is added. Under this supposition, for Eqn (B.1) to hold for $n+1$ transformed background error variables, we just need to show that

$$\langle \mathbf{x}^{b'_{n+1}} \mathbf{x}^{b'_m\top} \rangle = 0 \quad \forall m < n+1.$$

¹In the field, \mathbf{U}_p is often called the physical balance operator.

By the definitions in Chapter 1,

$$\begin{aligned}
\langle \mathbf{x}_{n+1}^{b'} \mathbf{x}_m^{b'\top} \rangle &= \langle \mathbf{x}_{n+1}^b \mathbf{x}_m^{b'\top} \rangle - \sum_{p=1}^n \alpha_{n+1,p} \langle \mathbf{x}_p^{b'} \mathbf{x}_m^{b'\top} \rangle \\
&= \langle \mathbf{x}_{n+1}^b \mathbf{x}_m^{b'\top} \rangle - \left(\sum_{p=1}^m \alpha_{n+1,p} \langle \mathbf{x}_p^{b'} \mathbf{x}_m^{b'\top} \rangle + \alpha_{n+1,m} \langle \mathbf{x}_m^{b'} \mathbf{x}_m^{b'\top} \rangle + \sum_{p=m+1}^n \alpha_{n+1,p} \langle \mathbf{x}_p^{b'} \mathbf{x}_m^{b'\top} \rangle \right) \\
&= \langle \mathbf{x}_{n+1}^b \mathbf{x}_m^{b'\top} \rangle - \alpha_{n+1,m} \langle \mathbf{x}_m^{b'} \mathbf{x}_m^{b'\top} \rangle - \left(\sum_{p=1}^m \alpha_{n+1,p} \langle \mathbf{x}_p^{b'} \mathbf{x}_m^{b'\top} \rangle + \sum_{p=m+1}^n \alpha_{n+1,p} \langle \mathbf{x}_p^{b'} \mathbf{x}_m^{b'\top} \rangle \right) \\
&= \langle \mathbf{x}_{n+1}^b \mathbf{x}_m^{b'\top} \rangle - \langle \mathbf{x}_{n+1}^b \mathbf{x}_m^{b'\top} \rangle - \left(\sum_{p=1}^m \alpha_{n+1,p} \langle \mathbf{x}_p^{b'} \mathbf{x}_m^{b'\top} \rangle + \sum_{p=m+1}^n \alpha_{n+1,p} \langle \mathbf{x}_p^{b'} \mathbf{x}_m^{b'\top} \rangle \right) \\
&= - \sum_{p=1}^m \alpha_{n+1,p} \langle \mathbf{x}_p^{b'} \mathbf{x}_m^{b'\top} \rangle - \sum_{p=m+1}^n \alpha_{n+1,p} \langle \mathbf{x}_p^{b'} \mathbf{x}_m^{b'\top} \rangle.
\end{aligned}$$

In the two remaining summations, it is clear that $p \leq n$. Since Eqn (B.1) is supposed to hold for the first n variables, then, $\langle \mathbf{x}_p^{b'} \mathbf{x}_m^{b'\top} \rangle = \mathbf{0}$ for $p \leq n$. This clearly means that

$$\langle \mathbf{x}_{n+1}^{b'} \mathbf{x}_m^{b'\top} \rangle = \mathbf{0}. \tag{B.2}$$

Notice that Eqn (B.2) rests on the supposition that Eqn (B.1) holds for the first n transformed background error variables, where n is an unspecified integer. In other words, Eqn (B.1) should also hold for n transformed background error variables if it holds for $n-1$ transformed background error variables. If we repeat this backtracking continuously, we will eventually arrive at the 2 transformed background error variables scenario. Eqn (B.1) has been previously proven to hold for that particular scenario. In other words, by induction, Eqn (B.1) holds for any number of transformed background error variables.

Appendix C

Additional plots

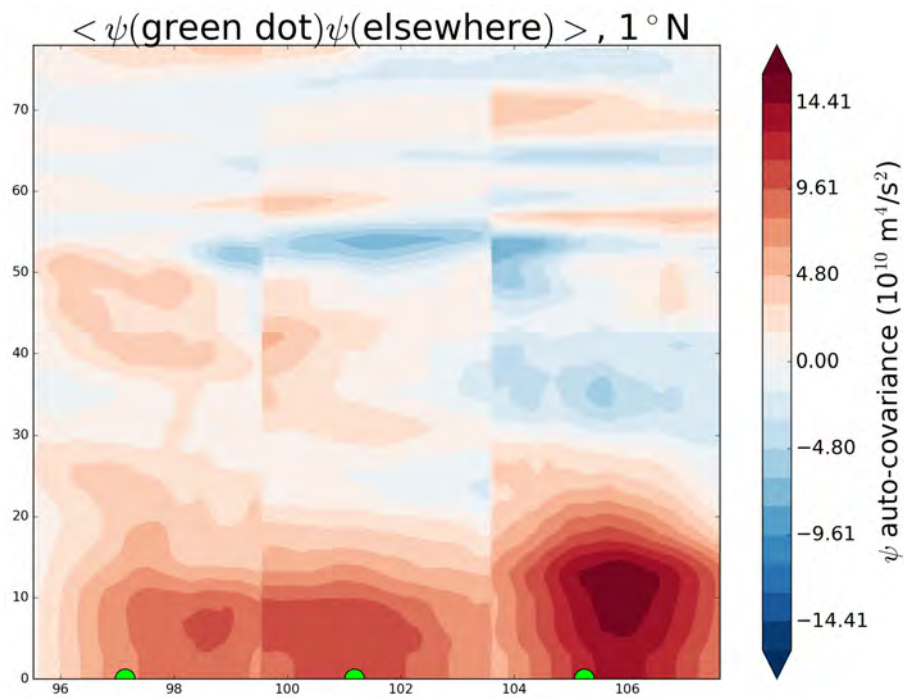


Figure C.1: The contour plot shows the covariance of ψ at the green dots on the surface, at 1°N , and ψ everywhere on a vertical slice of ϕ , along the 1°N latitude. This plot is separated into 3 vertical strips, each centred on a the green dot, spanning 2° longitude in the horizontal and all 79 model levels in the vertical. The horizontal axis is longitude and the vertical axis is model level. This plot is referenced to in the discussion near Figure 4.5

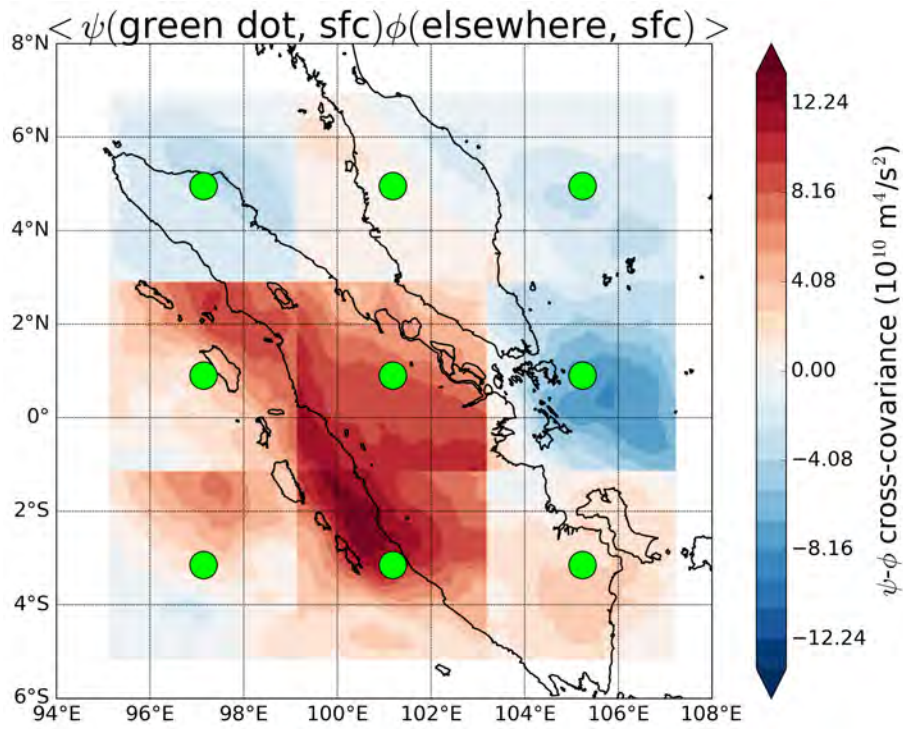


Figure C.2: The contour plot shows the covariance between ψ at the green dots and ϕ in squares centred on the green dots, on the surface level. The dimensions of each square is $450 \text{ km} \times 450 \text{ km}$, or roughly $4^\circ \times 4^\circ$. The landmasses in the domain are outlined in thick, solid black lines.

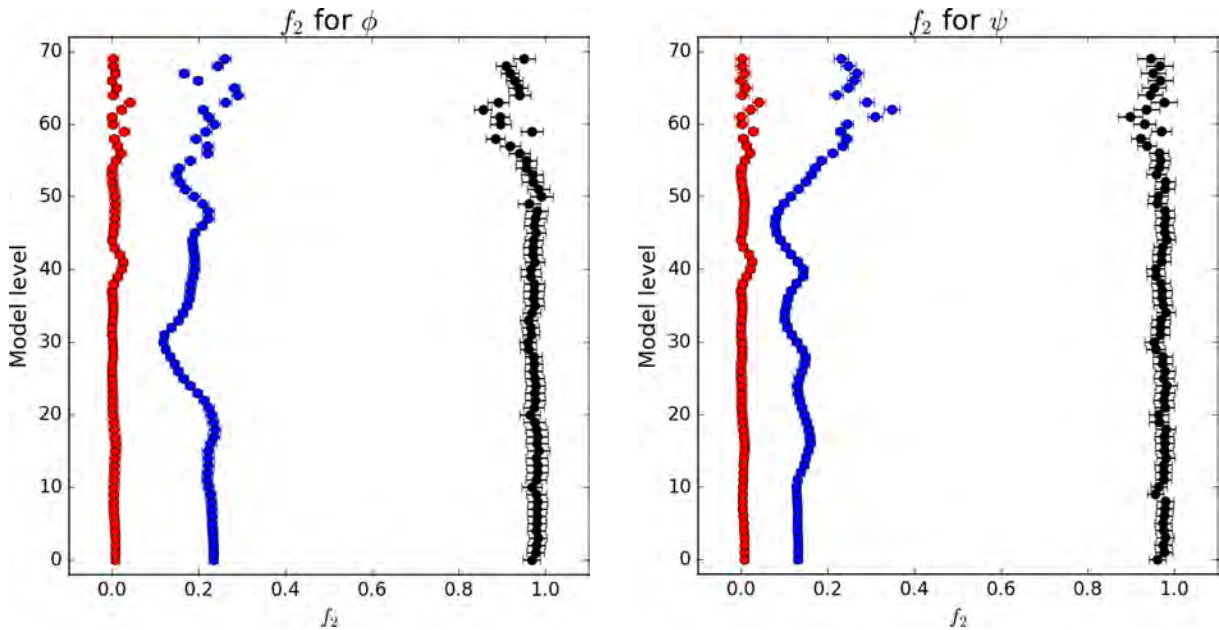


Figure C.3: Contribution fractions (Chen et al., 2013) for the three ways of modelling regression coefficients. The red data points are for the first model, the blue for the second model and the black for the third model.