

# Deterministic and Probabilistic Weather Forecasting: An Analysis towards Big Data Samples Using the Google Search Random Surfer

Amarasingha Arachchige Mihiri Chathurika<sup>1</sup>, Bhupendra Nath Tiwari<sup>1,2,\*</sup>, Chandra Kishore<sup>1</sup>

<sup>1</sup>University of Information Science and Technology "St. Paul the Apostle", Republic of North Macedonia

<sup>2</sup>INFN-Laboratori Nazionali di Frascati, Italy

*Received August 1, 2019; Revised September 3, 2019; Accepted September 9, 2019*

Copyright©2019 by authors, all rights reserved. Authors agree that this article remains permanently open access under the terms of the Creative Commons Attribution License 4.0 International License

**Abstract** The concept of big data has become one of the most important topics in the field of information science and engineering. In this paper, we offer modeling of data and its stability and forecasting by considering anti-symmetric traceless and symmetric models for atmospheric pressure variations. The data sample has been collected every 10 minutes for several years during 2009-2016 at the Weather Station, Max Planck Institute for Biogeochemistry, Jena, Germany. Subsequently, we extend the proposed model with a probabilistic transformation matrix by considering the Google search random surfer matrix with a small damping factor  $p$  ( $0 < p < 1$ ). Following the Principal Component Analysis (PCA), our study plays a vital role in big data samples and their stability analysis. A comparative discussion is provided for the above transformation matrix and its probabilistic counterpart. Finally, predictions are made towards feature selection, PCA and data compression sensing in the light of big data.

**Keywords** Weather Forecasting, Big Data Analytics, Sampling Theory, Principal Component Analysis, Google Search Random Surfer

---

## 1. Introduction

Forecasting is an important problem concerning the weather and its modeling [1]. Following the same, in this paper, we provide modeling of data samples in the light of their stability and forecasting. Our focus is towards the study of big data samples. First of all, we consider that a big data sample is large and it has complex data structures. The concepts of big data sets comply with its main qualities

such as the volume, velocity, variety and veracity [2]. On the one hand, a big data sample can have a high dimensionality [3]. Such challenges are considered through our modeling techniques, see [2] for an introduction towards big data.

Big data configurations hold great promises for discovering subtle patterns in population sciences and heterogeneities that are usually not possible in a small scale data size. Thus, a large sample size that involves a vast volume and high dimensionality of big data requires unique computational and statistical challenges including the scalability, storage bottleneck, noise accumulation, spurious correlation, incidental endogeneity and measurement errors [4]. Furthermore, the modeling of data and its stability are important in predicting the future behavior of configuration. In this direction, we have studied the Richardson Integration over fluctuations of its step sizes for arbitrary real valued integratable functions [5].

The present investigation supports data modeling that arises as the process of creating data samples towards the optimal information system designing and applications. Our focus includes development of a model that determines possible atmospheric pressure variations. In this concern, Barometric pressure and atmospheric pressure fluctuations are among our future studies. We also focus on anti-symmetric traceless and symmetric models and their relation to convex analysis techniques [6]. Here, the convex analysis plays an important role in the light of optimization theory. Namely, the convex analysis simplifies the optimization with respect to constrains of the concept of convexity into the problem. This is also named as convex optimization as it is optimizing the convex function. For a convex function, all local minima is same as global minima. It is because the minimization method is based on the analysis of convexity properties of a distance

function [7].

Existence of the global minimum via the convex function techniques are defined by a convex set [8]. This is characterized by directions of the flow of the objective function. For determining its global minima, a nonconvex function can be convexified [8]. A convex set is connected in its nonempty relative interior and it has feasible direction at any point. A real-valued convex function is continuous and it has good differentiability properties. Herewith, closed convex cones arise as the self-dual object with respect to the polarity [9]. In the light of convex analysis, lower semi-continuous functions can equally be treated by self-dual objects with respect to the conjugacy. Here, linear problems are often solved by convex methods and vice-versa [10,11]. Concerning the polyhedral convex sets, the duality concept plays a vital role.

On the other hand, the probabilistic inference is the task of deriving the probability of one or more random variables taking a specific value or set of values. Probabilistic inference uses stochastic models to classify the problem by minimizing the loss function [12]. The problem is to determine its efficiency in making predictions. Hereby, one uses a distance function in the light of linear difference model. Further, fuzzy logic provides a viable principle of approximate reasoning with a limiting case. It has many-valued logic in which the truth values of variables may be any real number between 0 and 1 inclusively. This defines the notion of fuzziness. Fuzzy logic could thus be used to predict chaotic behavior [13].

Further, in the realm of non-convex theories, there are gradient descent and Quasi-Newton method and others that are used to minimize the non-convex errors of continuous functions in a high dimensional space. A trust region approaches are used to determine the second order methods to non-convex problems. In such methods, in order to remove a negative curvature, the Hessian techniques are defined as a damped configuration. This is realized by adding a constant to the diagonal of its Hessian matrix. This is equivalent to adding a constant to each of its eigenvalues [14]. Hereby, from the outset of this paper, the convex analysis and convex programming are investigated on the same footing of the optimization theory.

In the light of optimization problems, finding the best solution from all possible solutions is among the main concerns of the present research. Depending on the number and the type of variables of the problem, optimization can be either continuous or discrete. In the case of a discrete optimization, the approach is for an integer, permutation or graph from a countable set. In a continuous problem, the approaches utilized include free optimization problems. Constrained optimization problems use the Lagrange techniques and other multimodal optimizations such as the gradient method, standard genetic algorithms, particle swarm and artificial bee colony methods [15, 16].

In the light of forecasting and data analysis, Google search random surfers are used as a matrix for various

probabilistic transformations. Namely, Google matrix is among one of the most accurate methods that is utilized for predictions. It was created as the result of dangling and disconnected nodes [17]. Here, the nodes refer to the considered websites that are visited during the search. In practice, a Google surfer matrix generates the best optimal prediction by using a large matrix of a considered data with reference to their connectivity as outlined in the next section. Based on the worldwide web server searching, weather is a very diverse phenomenon in its characters and predictions. It spreads among huge uncertainties where data points may not be connected to each other. Following the same, we propose a model for exploiting the Google surfer matrix towards the forecasting of weather [18].

In this research, instead of web sites we take atmospheric pressure variations to make predictions. We have downloaded a data sample as taken from Kaggle data science server [19]. The data has been collected every 10 minutes for several years during 2009-2016 at the Weather Station, Max Planck Institute for Biogeochemistry, Jena, Germany. We consider the difference of the atmospheric pressure to create the transformation matrix in parallel with the Google surfer matrix. For a better prediction, we define a probabilistic transformation model using transformation matrix with a damping factor. With the help of PCA, we provide the prediction of atmospheric pressure by obtaining the maximum and minimum eigenvalues of the probabilistic transformation model.

The rest of the paper is presented as follows. In section 2, we offer an overview of the model with convex optimization, Google surfer matrix and stochastic optimization. In section 3, we proposed two models based on anti-symmetric traceless and symmetric matrices. In section 4, we generalize models in the light of optimization theory. In the section 5, we provide verification of the models with the evaluations of the results. In section 6, we discuss our conclusions and future directions for further research and developments.

## 2. Review of the Model

In this section, we present an overview of the proposed model. Let's consider a generic form of continuous problem defined as the objective function that is to be optimized where we can apply convex programming and Google surfer matrix.

### 2.1. Convex Optimization

To study the convex optimization techniques, we introduce a cost vector  $C$  with its transpose

$$C^T = (c_1 + c_2 + c_3 + \dots + c_n).$$

Therefore, the underlying objective function is defined as

$$C^T X = c_1 x_1 + c_2 x_2 + \dots + c_n x_n.$$

Here, the variables  $x_1, x_2, x_n$  are variables with respect to which the objective function  $C^T X$  is sought to be optimized. As a minimization problem, the above objective function reads as

$$\min_X (C^T X) \text{ such that } AX \leq b.$$

In this setting, we can optimize the nodes through the following notations

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{pmatrix}, X = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix}, \text{ and } b = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{pmatrix}.$$

When we have equality, we solve the problem by applying Gauss elimination method and get  $b = A^{-1}X$ . In the case of  $b \geq AX$ , we use the Fourier-Motzkin elimination techniques [20] to find the optimal results through Google surfer matrix models.

## 2.2. Google Surfer Matrix

In order to illustrate the Google surfer matrix, let's consider an example of four nodes as in Figure 1. By assigning weights to each edge as below, we get the linear algebraic equations as

$$\begin{aligned} x_1 &= x_2 + \frac{1}{2}x_3, \\ x_2 &= \frac{1}{4}x_1 + \frac{1}{3}x_4, \\ x_3 &= \frac{1}{4}x_1 + \frac{1}{3}x_2 + \frac{1}{3}x_4, \\ x_4 &= \frac{1}{4}x_1 + \frac{1}{4}x_2 + \frac{1}{2}x_3. \end{aligned}$$

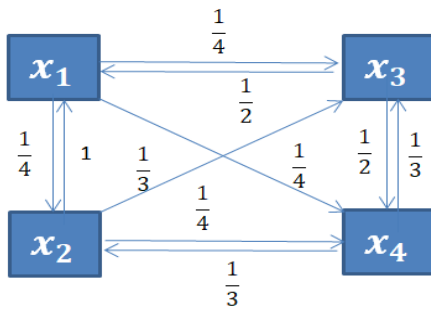


Figure 1. Four nodes problem

From the above four equations, we get the following transformation

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} 0 & 1 & 1/2 & 0 \\ 1/4 & 0 & 0 & 1/3 \\ 1/4 & 1/3 & 0 & 1/3 \\ 1/4 & 1/4 & 1/2 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix}.$$

In the matrix notation, it can be expressed as

$$X' = AX,$$

where  $X'$  is the output,  $A$  is transformation matrix and  $X$  is the input.

In this system, the relevant results are found after surfing for a very long time, where there is no way of reducing/removing the outliers [2].

## 2.3. Stochastic Optimization

In order to remove above disadvantages of outliers, we analyze transformation matrix  $A$  with probability  $p$ , where  $p \in [0, 1]$ . Here,  $p$  is called damping factor or bargaining factor. Instead of the transformation matrix  $A$ , we define the objective function with a probabilistic transformation matrix  $M$ . In this case, in the light of the Google surfer matrix,  $M$  is also known as the Page Rank matrix [18] that is defined as

$$M = (1 - p)A + pB,$$

where the matrix  $A$  is the same as before in the foregoing subsection and the matrix  $B$  reads as

$$B = \frac{1}{n} \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}.$$

The purpose of having  $B$  as the matrix with its all elements unity is to avoid outliers [2]. Here,  $n$  is the number of nodes or the size of the transfer matrix.

For example, for the system as considered above the matrix  $B$  can be defined as

$$B = \frac{1}{4} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}.$$

In order to make reliable predictions, we need to find the eigenvalue of  $M$  as a function of  $p$ . In special cases such as for a small  $p$  ( $p \approx 0$ ), the eigenvalue of  $M$  is equal to eigenvalue of  $A$ . It's called stochastic PCA. In such a transfer matrix, the eigenvalues will transform the matrix into linear system and tell about the stability of the data. Hereby, we obtain the minimum and maximum eigenvalues of the transfer matrix  $M$ .

## 3. Generalization of the Model

In this section, we provide the generalization of the symmetric and anti-symmetric traceless models in the light of the optimization theory as below. With the optimization formulation as in section 2, we may redefine the problem as

$$y = AX - b.$$

Here, the norm of  $y$  reads as

$$\|y\|^2 = y_1^2 + y_2^2 + \dots + y_n^2.$$

In this setting, our proposed optimization reads as the minimization of distance

$$\min_X \|y\|^2 = \min_X \|AX - b\|^2.$$

Hereby, the above equation generates a linear regression model.

Further, the optimal feature selection and data compressed sensing are done via the stochastic transformation matrix

$$M = (1 - p)A + pB$$

$$= (1 - p) \left[ A + \frac{p}{1-p} B \right].$$

Define a new parameter

$$\lambda = \frac{p}{1-p}; 0 \leq p \leq 1,$$

whereby we get the following probabilistic transformation matrix

$$M = k(A + \lambda B), \text{ where } k = 1 - p.$$

The analysis is perform by computing the transformation vector

$$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}.$$

Herewith, the concerned optimization problem can be redefined as

$$MX = kAX + \lambda BX.$$

In terms of the norms, this simplifies as follows

$$\|MX\|^2 = k\|AX\|^2 + \lambda\|BX\|^2,$$

$$\min\|MX\|^2 = \min_X (k\|AX\|^2 + \lambda\|BX\|^2).$$

This is termed as the LASSO method [21] Herewith, it is observed that a probabilistic model is same as the LASSO model in the light of feature selection [21, 22]. Such directions are performed in the light of our proposal as below.

## 4. Proposed Models

In this section, we provide two models to analyze the data sample. This allows to finding the spreading through the maximum and minimum eigenvalues of the transformation matrix. These models support forecasting and dynamical behavior of the sample. Predictions are made towards the atmospheric pressure determination.

### 4.1. Model 1: Anti-symmetric Configuration

Using atmospheric pressure variations, we build predictions based on an anti-symmetric traceless matrix as the following linear difference model. Considering the matrix elements

$$t_{ij} = t_i - t_j.$$

The transformation matrix  $A$  reads as

$$A = (t_i - t_j)_{n \times n}.$$

Following the same, we have the following transformation matrix

$$A = \begin{pmatrix} 0 & a_{12} & \cdots & a_{1n} \\ -a_{12} & 0 & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ -a_{1n} & -a_{2n} & \cdots & 0 \end{pmatrix}.$$

### 4.2. Model 2: Symmetric Configuration

Similarly, for the case of a positive definite matrix  $A$ , we define the elements  $A$  as the absolute difference  $a_{ij} = |a_i - a_j|$ . Hereby, it follows that the transformation matrix reads as the following symmetric traceless matrix

$$A = \begin{pmatrix} 0 & a_{12} & \cdots & a_{1n} \\ a_{12} & 0 & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & \cdots & 0 \end{pmatrix}.$$

In both above models we observe that there are no autocorrelations, that is, we have  $a_{ii} = 0$ . Consequently, our model predicts correlations between two distinct observation points.

## 5. Verification of the Models

From the data sample as in [19], we have made three sub samples as Data Set 1, Data Set 2 and Data Set 3 that contain the consequent first thousand, second thousand and third thousand from the main data sample.

After creating a matrix  $A$ , with Python libraries named Pandas and Numpy, we have obtained the maximum and minimum eigenvalue of the stochastic transfer matrix  $M$  with  $p = 0.33$  for the above three dataset.

Here, in order to reduce the memory size, we convert the data type of the atmospheric pressure measured in the unit of pmbar to np.float32 unit as for both the proposed anti-symmetric and symmetric models.

### 5.1. Model 1: Anti-symmetric Configuration

For the case of the anti-symmetric traceless matrix as above in section 4.1, using Python, we find the maximum and minimum of the eigenvalues of  $M$  as below:

Data Set 1:

$$(max, min) = (165 + 3226.962j, -1.0236668e - 05 - 4.797735e - 06j).$$

Data Set 2:

$$(max, min) = (165 + 4310.9688j, -5.8782107e - 06 - 2.7488336e - 06j).$$

Data Set 3:

$$(max, min) = (164.99998 + 7031.0015j, -5.1989286e - 06 - 3.7833154e - 06j).$$

In this framework, we see that the maximum eigenvalue is quite large in comparison to minimum eigenvalue of  $M$ . The imaginary components show that there is a damping in the chosen model. The real component of the maximum eigenvalue is nearly same for all the data sets. The minimum eigenvalues have different orders.

### 5.2. Model 2: Symmetric Configurations

For the case of the symmetric traceless matrix as in above section 4.2, we obtain the maximum and minimum of the eigenvalues of  $M$  as

Data Set 1:

$$(max, min) = (4223.138, -2213.86).$$

Data Set 2:

$$(max, min) = (5163.6606, -3661.3755).$$

Data Set 3:

$$(max, min) = (8637.3, -5250.997).$$

We observe that both the maximum and minimum eigenvalues are of different orders. There is a large difference between the minimum and maximum eigenvalues of  $M$ . The absence of the imaginary components  $j$  shows that there is no damping in this model.

## 6. Conclusions

In this paper, we have discussed the optimization of a continuous problem in the light of atmospheric pressure variations. Our model uses Google search matrix as it is one of the most efficient and accurate algorithm. We have considered transfer matrix and optimization of big data samples. Hereby, notice that a large matrix calculation requires applications of fast computers, high memory capacity and advanced mathematical formulations, the Google search matrix can be used in forecasting of the optimal value of the weather.

In particular, for both anti-symmetric and symmetric traceless matrix of the probabilistic transformation models, we have respectively complex and real numbers as the minimum and maximum eigenvalue of the transfer matrix. In the case of anti-symmetric traceless transfer matrix, we get an oscillatory system with damping. It has the effect of preventing its oscillations by reducing its amplitude. Hence, the anti-symmetric traceless configuration is not stable.

On the other hand, configurations with a symmetric traceless transfer matrix are stable. In sort, anti-symmetric traceless configurations are rather useful towards the prediction of the atmospheric pressure and related

forecasting of the weather. Related numeric analysis and development of algorithm [23] are left open for future research and investigations.

## Appendix: Eigenvalues and Eigenvectors of a Matrix

In this appendix, we introduce the eigenvalues of a matrix and the norm of an eigenvector. Namely, let  $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$  be a 2x2 matrix, then its eigenvalue  $\lambda$  is found by solving the equation

$$AX = \lambda X,$$

where  $X$  is a two dimensional vector. That is, we solve the equation

$$\text{def}(A - \lambda I) = 0$$

$$\text{def} \left( \begin{pmatrix} a & b \\ c & d \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right) = 0$$

$$\text{def} \begin{pmatrix} a - \lambda & b \\ c & d - \lambda \end{pmatrix} = 0$$

$$(a - \lambda)(d - \lambda) - bc = 0$$

$$\lambda^2 - (a + d)\lambda + (ad - bc) = 0.$$

As per the above equations, trace is the sum of the diagonal elements, that is, we have  $\text{tr}(A) = a + d$  and its determinant  $\text{det}(A)$  is equal to the quantity  $ad - bc$ . In other words, it follows that we have the following determinant

$$\|A\| = ad - bc.$$

Therefore, the above characteristic equation simplifies as

$$\lambda^2 - \text{tr}(A)\lambda + \|A\| = 0.$$

The solution of above the equation is given by

$$\lambda = \frac{\text{tr}(A) \pm \sqrt{\text{tr}(A)^2 + 4\|A\|}}{2} = \lambda_1, \lambda_2.$$

Similarly, the given eigenvector

$$X = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix},$$

with its transpose

$$X^T = (x_1 \quad x_2),$$

the respective norm of  $X$  is defined as

$$X^2 = X^T X = (x_1 \quad x_2) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}.$$

In other words, the above norm reads as the summation

$$\|X\|^2 = x_1^2 + x_2^2.$$

The same is used in probabilistic optimization models, e.g. the LASSO model.

## Acknowledgement

B.N.T. would like to thank the Yukawa Institute for Theoretical Physics at Kyoto University. Discussions during the workshop YITP-T-18-04 "New Frontiers in String Theory 2018" were useful towards the initiation of this work.

## REFERENCES

- [1] Showalter, A. K. (1953). A stability index for thunderstorm forecasting. *Bulletin of the American Meteorological Society*, 34(6), 250-252.
- [2] Bucchianico, D. A. (2017). The mathematics behind big data. Data Science Center Eindhoven, 4TU AMI SRO Big Data Meeting, Big Data: Mathematics in Action.
- [3] Wu, X., Zhu, X., Wu, G. Q., & Ding, W. (2013). Data mining with big data. *IEEE transactions on knowledge and data engineering*, 26(1), 97-107.
- [4] Mathur, A., Sihag, A., Bagaria, E. G., & Rajawat, S. (2014). A new perspective to data processing: Big Data. In 2014 International Conference on Computing for Sustainable Global Development (INDIACom) (pp. 110-114). IEEE.
- [5] Tiwari, B. N., & Chaturika, A. A. M. (2019). Optimization of the richardson integration over fluctuations of its step sizes. *Cogent Mathematics & Statistics*, (just-accepted), 1643438.
- [6] Xie, D., Hu, X., & Zhang, L. (2003). The solvability conditions for inverse eigenproblem of symmetric and anti - persymmetric matrices and its approximation. *Numerical Linear Algebra with Applications*, 10(3), 223-234.
- [7] De Leeuw, J. (2011). Applications of convex analysis to multidimensional scaling.
- [8] Strongin, R. G., & Sergeyev, Y. D. (2013). *Global optimization with non-convex constraints: Sequential and parallel algorithms* (Vol. 45). Springer Science & Business Media.
- [9] Suneja, S. K., Aggarwal, S., & Davar, S. (2002). Multiobjective symmetric duality involving cones. *European Journal of Operational Research*, 141(3), 471-479.
- [10] (2012). *Lecture: Convex Analysis and Optimization*. Massachusetts Institute of Technology.
- [11] Bertsekas, D. P. (2009). *Convex optimization theory* (pp. 157-226). Belmont: Athena Scientific.
- [12] Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- [13] Martinez-Soto, R., Castillo, O., Aguilar, L. T., & Melin, P. (2010). Fuzzy logic controllers optimization using genetic algorithms and particle swarm optimization. In *Mexican International Conference on Artificial Intelligence* (pp. 475-486). Springer, Berlin, Heidelberg.
- [14] Dauphin, Y. N., Pascanu, R., Gulcehre, C., Cho, K., Ganguli, S., & Bengio, Y. (2014). Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. *Neural information processing systems* (pp. 2933-2941).
- [15] Shen, J. *Introduction to Optimization Theory*. Lecture Notes, School of Economics, Shandong University.
- [16] Ong, Y. S., Nair, P. B., & Keane, A. J. (2003). Evolutionary optimization of computationally expensive problems via surrogate modeling. *AIAA journal*, 41(4), 687-696.
- [17] Paparo, G. D., & Martin-Delgado, M. A. (2012). Google in a quantum network. *Scientific reports*, 2, 444.
- [18] Tanase, R., & Radu, R. (2016). Lecture# 3: PageRank algorithm—the mathematics of Google Search. Online: [www.math.cornell.edu/mec/Winter2009/RalucaRemus/Lecture3/lecture3.html](http://www.math.cornell.edu/mec/Winter2009/RalucaRemus/Lecture3/lecture3.html).
- [19] Kris. (2017). Weather Archive Jena. Online: <https://www.kaggle.com/pankrzysiu/weather-archive-jena>.
- [20] Zhen, J., Den Hertog, D., & Sim, M. (2018). Adjustable robust optimization via Fourier–Motzkin elimination. *Operations Research*, 66(4), 1086-1100.
- [21] Wright, S. (2015). *Optimization Techniques for Learning and Data Analysis*. University of Wisconsin-Madison, IPAM Summer School.
- [22] Zhao, P., & Yu, B. (2006). On model selection consistency of Lasso. *Journal of Machine learning research*, 7(Nov), 2541-2563.
- [23] Kalnay, E., Kanamitsu, M., & Baker, W. E. (1990). Global numerical weather prediction at the National Meteorological Center. *Bulletin of the American Meteorological Society*, 71(10), 1410-1428.