



GRAPH LEARNING METHODS FOR CLIMATE MODELS WITH GRAPH SPARSIFICATION

Justin Clarke

Supervisor: Associate Professor Yanan Fan

School of Mathematics and Statistics
UNSW Sydney

November 2023

SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS OF THE DEGREE OF
BACHELOR OF SCIENCE WITH HONOURS

Plagiarism statement

I declare that this thesis is my own work, except where acknowledged, and has not been submitted for academic credit elsewhere.

I acknowledge that the assessor of this thesis may, for the purpose of assessing it:

- Reproduce it and provide a copy to another member of the University; and/or,
- Communicate a copy of it to a plagiarism checking service (which may then retain a copy of it on its database for the purpose of future plagiarism checking).

I certify that I have read and understood the University Rules in respect of Student Academic Misconduct, and am aware of any potential plagiarism penalties which may apply.

By signing this declaration I am agreeing to the statements and conditions above.

Signed: Justin Clarke

Date: 17/11/2023

Acknowledgements

Abstract

Graph sparification techniques for graph neural networks have traditionally been used to accelerate training and inference on real-world graphs which have billions of paramaters. There are also many different climate models which use complex mathematical models to model the interactions between energy and matter over the world. Many of these models share components and the structure of these relationships is not easily found due to the complexity of these climate models. The space of all possible graphs grows exponentially with the number of nodes and as such any correlation or causality is difficult to find. In this thesis, I attempt to quantify these relationships with graph sparsification techniques. (Talk more about climate?)

Contents

Chapter 1	Introduction	3
1.1	Climate Models	3
1.2	Ensemble Models	4
1.3	Machine Learning	5
1.4	Sparsification of Deep Learning	5
1.5	Motivation	6
1.6	Outline	6
Chapter 2	Methods and Related Techniques	8
2.1	Deep Learning	8
2.2	Graph Neural Networks	9
2.3	Sparsification of Graph Neural Networks	10
2.4	Variational Autoencoder	13
Chapter 3	Framework	14
3.1	Climate Models	14
3.2	Dataset	14
3.3	Exploratory Data Analysis	15
3.4	Problem Formulation	17
3.5	Implementation	18
Chapter 4	Results	19
4.1	Model verification	19
4.2	Model results	19
Chapter 5	Discussion	20
Chapter 6	Appendix	21
References		22

List of Figures

2.1	Examples of (a) matrix convolution with a 3×3 kernel and (b) a graph where a convolution is much more difficult to define.	9
2.2	Illustration of the Unified Graph Sparsification.	13
3.1	Visualisation of Climate Change in CMIP6 Models	15
3.2	Plot of model aggregated temperatures at time slices from 1850-2100	16
3.3	Correlation Heatmap of CMIP6 Models	16
3.4	Partial Correlation Heatmap of CMIP6 Models	17

List of Tables

1.1	Shared Socioeconomic Pathways Descriptions	4
1.2	Representative Concentration Pathways	4
3.1	CMIP6 Models used in this thesis. Model number is assigned for node numbering purposes.	15

CHAPTER 1

Introduction

1.1 Climate Models

Climate models are based on well-known scientific processes and attempt to simulate the movement of fluids and energy throughout a system. The simplest of these have existed since the 1950's with the very first computers modelling simple variables on small two-dimensional climates [9]. Modern Global Climate Models (GCMs) have been expanded to covering the globe in three dimensions. Most of these model the complex interactions between various physical, chemical, biological and geological processes. These models are constantly being updated with new data as many different groups and institutions implement higher spatial and temporal resolutions. This has in part been driven by developments in computational techniques and the vast amount of data available worldwide to train these models. [38] A common approach to studying the complex system of Earth's climate has been through sophisticated mathematical modelling with a range of temporal and spatial data [27] and at the centre of most of these approaches are the Navier-Stokes equations which are used describe the movements of liquids and gases in our oceans and atmosphere. [39]

The primary focus of these models has been forecasting time series data, as climate science aims to address a fundamental question: What impact have contemporary human greenhouse gas emissions had on the Earth and its future? Understanding this is important as it can inform us on the potential harms to society and the environment and guide the population and policy makers who are implementing change. The Coupled Model Intercomparison Project is a collaborative project between many meteorological institutions that aims to improve our understanding of climate change. The sixth iteration of these models or CMIP6 are the premier models for this task but while it was expected to contain around 100 models made by 49 separate groups, delays have caused only 40 have been published so far. However, the results from these 40 models so far indicates a far greater sensitivity to increases in greenhouse gases when compared to the previous generation of CMIP5 models. [16]

The basis behind the Intergovernmental Panel on Climate Change and its 2021 IPCC sixth assessment from the Coupled Model Intercomparison Project (CMIP6) were the Shared Socioeconomic Pathways (SSP). [29] These SSP's represent a broad set of possible changes in population, economic and technological growth, and urbanisation that would influence future changes to the climate. [43] These are directly

related to Representative Concentration Pathways (RCP) introduced by the previous CMIP5 which are categorisations based on the estimated future concentrations of greenhouse gases in the atmosphere. [15]. Tables 1.1 and 1.2 detail the levels of these SSPs and RCPs.

SSP	Description
SSP1	Sustainability: The world shifts gradually, but pervasively, toward a more sustainable path, emphasizing more inclusive development that respects perceived environmental boundaries.
SSP2	Middle of the road: The world follows a path in which social, economic, and technological trends do not shift markedly from historical patterns.
SSP3	Regional rivalry: A resurgent nationalism, concerns about competitiveness and security, and regional conflicts push countries to increasingly focus on domestic or, at most, regional issues.
SSP4	Inequality: Highly unequal investments in human capital, combined with increasing disparities in economic opportunity and political power, lead to increasing inequalities and stratification both across and within countries.
SSP5	Fossil-fueled development: This world places increasing faith in competitive markets, innovation and participatory societies to produce rapid technological progress and development of human capital as the path to sustainable development. Global markets are increasingly integrated.

Table 1.1: Shared Socioeconomic Pathways Descriptions

RCP	Temperature Increase (2081–2100)	Sea Level Rise (2081–2100)
2.6	1.0° C	0.4m
4.5	1.8° C	0.47m
6.0	2.2° C	0.48m
8.5	3.7° C	0.63m

Table 1.2: Representative Concentration Pathways

Scenarios are named based on the conjunction of their SSP level and RCP values. For example, on the lower end, SSP126 assumes an increasingly sustainable world where consumption is oriented towards minimising material resource and energy usage while SSP585 assumes a worst case scenario where fossil fuel usage and an energy-intensive lifestyle intensifies. The main output of CMIP6 models are the “scenario runs” which predict various outcomes in temperature, precipitation, air pressure and solar radiation given a certain SSP over time from 1850–2100.

1.2 Ensemble Models

Ensemble modelling is a process where multiple models are used in combination for a task and are more performant when the base models are diverse and independent [26]. Simple averages of multiple models are regularly shown to have better correlation with historical observations than any individual model alone. [22] [10]

CMIP6 is known as an ‘ensemble of opportunity’ [25], where the makeup the ensemble is determined by the ability of each base model to contribute. [1] This is the main benefit of ensemble methods as models can be weighted based on how accurate they are certain predictions. However, as research has become far more interconnected in the modern era, many aspects such as expertise, code and literature are often shared between groups. As such, many of the models that contribute to CMIP6 are highly likely to be dependent of each other. [1]. The degree of dependence between these models is difficult to ascertain as this would require a qualitative investigation into the personel, code and references between each component of CMIP6. Various novel approaches such as stochastic Markov chains [27] have been used to provide a more optimal ensemble mean which may account for this dependence. However, one would expect there to be some ground truth graph structure that links all models together through some dependence.

1.3 Machine Learning

The size of climate data has lead us to explore more scalable techniques to model this data. There are many different kinds of machine learning but the advent of the deep learning paradigm in the past decade has become the gold-standard of machine learning [4]. This has led to countless advancements in many practical applications such as, audio processing [3], visual data processing [49] and natural language processing [54] to name a few examples. The name deep learning refers to the ability of deep learning models to extract high-level, abstract features from raw data by using many layers of simple, computer understandable representations. Computers perform well on complex logical tasks such as arithmetic but often struggle with more simplistic tasks that are harder to quantify such as visual recognition and language. The ability of deep learning models to quantify these simple tasks have allowed artificial intelligence to apporach near human level understanding. [14] The first neural networks developed in the late 1950’s sought to simulate how human brain learned and operated. [44] The next development in neural networks also came from neuroscientific principles [18] with Convolutional Neural Networks (CNN) which could train models to be equivariant to translations in data and process data with grid-like structure. In recent years, Graph Neural Networks (GNN) have become the premier method of processing data with non-cartesian structure as standard convolutions on a graph structure much harder to define. The main feature of GNNs is the message passing framework, where information from features on each node is passed to neighbouring nodes then aggregated and embedded. This is then propagated through a neural network structure to perform a range of tasks on the entire graph, individual nodes or edges. Much of this data exists in the world in applications such as chemical analysis [53], social networks analysis [42], link prediction [56] and unstructured data processing [37].

1.4 Sparsification of Deep Learning

An estimated 80 to 90 percent of the worlds’ 79 zetabytes of data is unstructured and graphs make up a significant proportion of this data [20]. Traditional neural networks are not able to extract meaninful insights from this unstructured data without significant cleaning. The modern age has also produced many advancements in computing such as Massively Parallel Processing (MPP) [36] and big data which has led

to an exponential growth in the size and complexity of these deep learning models. The well-known Generative-Pretrained Transformer 3 (GPT-3) model by OpenAI commonly used for ChatGPT had variants with 175 billion parameters which required 800 gigabytes to store. [40] Although Deep Neural Networks (DNN) tend to generalise well even when overparameterised [6], this level of overparameterisation makes inference and prediction highly costly when the same performance could be achieved on a far more simple model. To address this, the concept of the Lottery Ticket Hypothesis (LTH) [12] was introduced which explored the possibility of simplifying redundant models by trainable sparse subnetworks whilst still training to full accuracy. Chen et.al. [7] extended the LTH to Graph Neural Networks by co-optimising graph and neural networks weights and zeroing out edges with the lowest magnitude. This reduced computational costs by over 85% depending on the size of the graph whilst maintaining a strong baseline accuracy.

1.5 Motivation

The goal of this thesis is to investigate whether these graph sparsification techniques can be used to determine some dependence structure within a graph of models which are all attempting to model the same scenario in the climate. Graph structure and dependence learning is already possible with unsupervised methods such as Variational Graph Autoencoders (VGAE) [55] but to our best knowledge, graph sparsification has not been used before as a method for graph structure learning or inferring dependence. Existing methods for determining multiple correlation such as partial correlation and multiple correlation coefficient R^2 are linear methods. Graph sparsification is primarily used for simplifying graphs which have been overparameterised and grown too large but by removing these edges, the edges that remain should theoretically hold some relation in dependence to the graph which is the motivation for this thesis.

1.6 Outline

The thesis is structured as follows. In Chapter 2, we will review the studies relating to neural networks, the extensions towards graph neural networks and the lottery ticket hypothesis (LTH). Current methods along with the benefits and shortcomings will be discussed and the terminology and definitions for will also be outlined in this section.

In Chapter 3, we will perform an exploratory data analysis (EDA) and provide a high-level overview of the dataset. This will provide a more in-depth understanding of the scenarios and CMIP6 models. There will also be a more comprehensive analysis of an individual climate model? (pick an example?) the climate modelling process and ensemble weighting method? We then introduce our problem formulation with by formalising the regression problem we will be using to perform the sparsification algorithm.

In Chapter 4, we visualise the results of the sparsification algorithm and we compare it to various correlation and partial correlation matrices. We may also look at

Mutual Information Criterion (MIC) and compare with Variational Graph Autoencoder?

Finally in Chapter 5, we summarise our findings and present some directions for possible future research.

CHAPTER 2

Methods and Related Techniques

2.1 Deep Learning

The advent of deep learning provided algorithms that could automatically extract higher-level features from raw data. [8] The multi-layer perceptron [45] is formulated using linked layers of nodes which transforms a set of inputs into an output. The single layer version of this model can be represented as

$$f(x) = \sigma(\Theta^\top X) \text{ where } \Theta = \begin{bmatrix} b \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \text{ and } X = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad (2.1.1)$$

Where $\sigma(\cdot)$ is some activation function such as ReLu, hyperbolic tangent or logistic function, X is the data, Θ is the learned parameters and b represents a bias term. These parameters are set to some initial values and are iteratively updated in a back-propagation training process,

$$\theta^{t+1} = \theta^t - \eta \frac{\partial E(X, \theta^t)}{\partial \theta} \quad (2.1.2)$$

Where $E(\cdot)$ is some loss function and η is the learning rate.

The next advancement in the deep learning space came with the Convolutional Neural Network (CNN) which was a regularised MLP that could handle data with data with structure and multiple dimensions far better than the traditional MLP due to its use of weight sharing, sampling and local receptive fields. [14] Suppose we have an image or some other kind of data in two-dimensional matrix form. Let $\mathbf{X} \in \mathbb{R}^{H \times W}$ be the input image and $\mathbf{W} \in \mathbb{R}^{h \times w}$ be the kernel or filter. By performing a convolution, we are effectively 'sliding' our weight matrix kernel over our input and the resulting feature map $\mathbf{Z} = \mathbf{X} * \mathbf{W}$,

$$Z_{i,j} = \sum_{u=0}^{h-1} \sum_{v=0}^{w-1} x_{i+u,j+v} w_{u,v} \quad (2.1.3)$$

The novelty of the convolutional layer compared to a linear layer is that the kernel is shared across all locations of the input and therefore if a pattern in the input moves, the corresponding output will also follow this movement. This provides shift equivariance which is something that early MLP's failed to achieve. [34]

2.2 Graph Neural Networks

When it comes to data in a graph-like structure, standard CNN's cannot be applied due to the non-euclidean nature of a graph. In an image or a matrix, our kernel is generally a $n \times n$ matrix which can be applied to the entirety of the data. In graphs this is not always possible due to the fact that any number of nodes can be connected by any number of edges. [46]

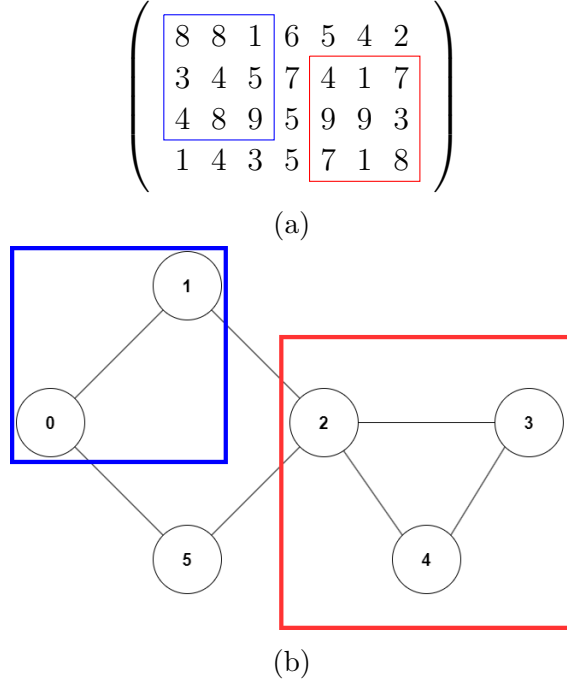


Figure 2.1: Examples of (a) matrix convolution with a 3×3 kernel and (b) a graph where a convolution is much more difficult to define.

The Graph Neural Network (GNN) was developed for this purpose and they can be broadly categorised into gating and attention based methods [50] and spectral or spatial methods within Graph Convolutional Network (GCN) research. [24]

We define a graph as $\mathcal{G} = (\mathcal{V}, \mathbf{A})$, where \mathcal{V} represents a set of vertices which contains a list of nodes $\{v_1, \dots, v_n\}$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ the adjacency matrix which contains information on the graph topology. If an edge is also weighted from node j to i , it is noted as e_{ji} . If an edge exists between two node v_i and v_j , then $\mathbf{A}_{ij} = 1$ else, $\mathbf{A}_{ij} = 0$. We also define the degree matrix as $\mathbf{D} = \sum_j A_{ij}$ where each entry on the diagonal is equal to the row sum of the adjacency matrix \mathbf{A} . Each node has a p -dimensional feature vector $x_i \in \mathbb{R}^p$ which describes some information about the node in the graph. By combining all n feature vectors from all nodes, we have a feature matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$. The graph also has a response $Y \in \mathbb{R}^p$ which is for a graph-level task but a node level task would simply have $Y \in \mathbb{R}^{n \times p}$. The convolution operation from [24] can be expressed as

$$\mathbf{X}' = \tilde{D}^{-1/2}(A + I)\tilde{D}^{-1/2}\mathbf{X}\Theta \quad (2.2.1)$$

or on a node-level, the operation known as the message passing framework can be thought of as

$$x'_i = \Theta^\top \sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{e_{ji}}{\sqrt{\hat{g}_j \hat{g}_i}} x_j. \quad (2.2.2)$$

where e_{ij} is the edge weight between nodes i and j and $\mathcal{N}(i)$ are the neighbouring nodes connected by an edge to node i . $\hat{g}_i = \sum_{j \in \mathcal{N}(i) \cup \{i\}} e_{ji}$ is the row sum of the edge weights. This effectively takes the node embeddings from each neighbouring node and aggregates it by a permutation invariant function such as the sum. Over many layers this passes the signal from all nodes around the graph to each other. [46] The final learning function for a two layer network is then just

$$f(\mathbf{A}, \mathbf{X}) = \sigma_2(\hat{\mathbf{A}}_2 \sigma_1(\hat{\mathbf{A}}_1 \mathbf{X} \Theta^{(0)}) \Theta^{(1)}) \quad (2.2.3)$$

where $\sigma_1(\cdot)$ and $\sigma_2(\cdot)$ are an activation function such as ReLU, and $\hat{A} = \tilde{D}^{-1/2}(A + I_N)\tilde{D}^{-1/2}$ is the symmetrically normalised adjacency matrix. The propagation rule is then as follows:

$$H^{(l+1)} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} \Theta^{(l)}) \quad (2.2.4)$$

Here, $\tilde{A} = A + I_N$ is the adjacency matrix of the undirected graph \mathcal{G} with self-connections from the identity matrix I_N . $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$ is the row sum of \tilde{A} and $\Theta^{(l)}$ is a layer trainable weight matrix. $H^{(l)} \in \mathbb{R}^{N \times D}$ is the matrix of activations in the l th layer with $H^0 = \mathbf{X}$ or the input.

2.3 Sparsification of Graph Neural Networks

Hornik [17] showed that in a single layer MLP and provided enough hidden neuron units, the neural network could model any smooth truth function as for each added neuron, the decision space can be segmented to added linear regions to conform to any response. Many experiments have also found that deep networks with many layers perform better than shallow ones [33] [41] as the usage of many layers allows deeper layers to leverage features produced by earlier layers. The novelty of neural networks is that they do not tend to be affected greatly by overparameterisation [6], and as they generally improve with more neurons and layers these models have grown exponentially in size with some models using billions of parameters and most of these models having more parameters than training observations which has made both inference and prediction incredibly costly. (Reference some math from the number of linear decision regions and some graphs?)

The most basic approach to inducing model sparsity is through an l_1 penalty also known as LASSO in the loss function. [34] Unlike the l_2 regularisation which penalises large magnitude weights, the l_1 regularisation minimises and zeros weights which encourages sparsity. On a linear regression, this is done with a MAP estimation with a Laplace prior and is equivalent to optimising

$$\mathcal{L}(\Theta) = \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\Theta\|_2^2 + \lambda \|\Theta\|_1 \quad (2.3.1)$$

where $\lambda > 0$ is some tunable sparsity parameter. This is a form of the bridge estimator [11] where $\alpha = 1$ and Θ is a solution to the objective function

$$\operatorname{argmax}_{\Theta} \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\Theta\|_2^2 + \lambda \sum_{j=1}^p |\theta_j|^\alpha \quad (2.3.2)$$

where $0 < \alpha < 1$. This bridge penalty induces the exponential power prior distribution on Θ in the form

$$\pi(\Theta|\lambda, \alpha) \propto \prod_{j=1}^p \frac{\lambda^{1/\alpha}}{2\Gamma(1 + 1/\alpha)} \exp(-\lambda|\theta_j|^\alpha) \quad (2.3.3)$$

This is easily extended to neural networks by applying this penalty to the weights in the layers of the network. [31] Despite this benefit, modern GPUs are optimised for dense matrix multiplication and as such there aren't many computational benefits from regularisation if certain weights across the network are zero. Methods that encourage *group* sparsity are able to prune whole nodes and layers out of our model result in *block sparse* weight matrices which provide much more substantial computational savings. [47] [52] [32] [30]

Dropout [48] is another strategy traditionally used to prevent overfitting but has been formalised as structured shrinkage prior. [35] At each training iteration, we turn off all outputs from each neuron with probability p . In this way, each neuron is trained to learn well individually even if some neurons are missing and prevents units from overfitting and depending on each other. We can view this as estimating a noise version of the weights [34],

$$\theta_{lji} = w_{lji} \epsilon_{li} \text{ where } \epsilon_{li} \sim \text{Ber}(1 - p) \quad (2.3.4)$$

When testing, this dropout is generally turned off but if used at test time, this becomes an ensemble of networks each with sparse graphs [13] [21] and becomes Monte Carlo dropout. This can be expressed as

$$p(\mathbf{Y}|\mathbf{X}, \mathcal{D}) \cong \frac{1}{N} \sum_{n=1}^N p(\mathbf{Y}|\mathbf{X}, \hat{W}\epsilon^n + \hat{\mathbf{b}}) \quad (2.3.5)$$

The Lottery Ticket Hypothesis (LTH) [12] explored the possibility of simplifying redundant models by trainable sparse subnetworks whilst still training to full accuracy. Chen et. al. [7] extended the LTH to Graph Neural Networks by co-simplifying both the adjacency matrix of the graph and the weights in the network of the model. For a semi-supervised classification task, the objective function is:

$$\mathcal{L}(\mathcal{G}, \Theta) = -\frac{1}{|\mathcal{V}_{\text{label}}|} \sum_{v_i \in \mathcal{V}_{\text{label}}} y_i \log(z_i), \quad (2.3.6)$$

where \mathcal{L} is the cross-entropy error of all samples and y_i is the label vector of node v_i . The Unified GNN Sparsification (UGS) framework then introduced two masks m_g and m_θ with the same shape as the adjacency matrix \mathbf{A} and the weights matrix Θ , which gives the following objective function:

$$\mathcal{L}_{\text{UGS}} = \mathcal{L}(\{m_g \odot A, \mathbf{X}\}, m_\theta \odot \Theta) + \gamma_1 \|m_g\|_1 + \gamma_2 \|m_\theta\|_1, \quad (2.3.7)$$

where \odot is the element-wise product, γ_1 and γ_2 are hyperparameters to control the shrinkage of m_g and m_θ . After training, the lowest magnitude elements in m_g and m_θ are set to zero with respect to some set values of p_g and p_θ . These sparse masks are then applied which prune \mathbf{A} and Θ . The algorithm is then as listed in Algorithm 1 and 2

Algorithm 1 Unified GNN Sparsification [7]

Input: Graph $\mathcal{G} = \{A, \mathbf{X}\}$, GNN $f(\mathcal{G}, \Theta_0)$, weight initialisation Θ_0 , masks $m_g^0 = A$ and $m_\theta^0 = 1 \in \mathbb{R}^{\|\Theta_0\|_0}$, step size η , λ_g and λ_θ , preset ratios p_g and p_θ

Output: Sparse masks m_g and m_θ

- 1: **for** iteration $i = 0, 1, 2, \dots, N - 1$ **do**
 - 2: Forward $f(\cdot, m_\theta^i \odot \Theta_i)$ with $\mathcal{G} = \{m_g^i \odot A, \mathbf{X}\}$ ▷ Computes 2.3.7
 - 3: Backpropagate to update $\Theta_{i+1} \leftarrow \Theta_i - \eta \nabla_{\Theta_i} \mathcal{L}_{\text{UGS}}$
 - 4: Update $m_g^{i+1} \leftarrow m_g^i - \eta \nabla_{m_g^i} \mathcal{L}_{\text{UGS}}$
 - 5: Update $m_\theta^{i+1} \leftarrow m_\theta^i - \eta \nabla_{m_\theta^i} \mathcal{L}_{\text{UGS}}$
 - 6: **end for**
 - 7: Set p_g of the lowest magnitude values in m_g^N to 0 and others to 1, then obtain m_g . ▷ Set at 5%
 - 8: Set p_θ of the lowest magnitude values in m_θ^N to 0 and others to 1, then obtain m_θ . ▷ Set at 20%
-

Algorithm 2 Iterative UGS to find GLT [7]

Input: Graph $\mathcal{G} = \{A, \mathbf{X}\}$, GNN $f(\mathcal{G}, \Theta_0)$, weight initialisation Θ_0 , initial masks $m_g^0 = A$ and $m_\theta^0 = 1 \in \mathbb{R}^{\|\Theta_0\|_0}$, predefined sparsity levels s_g and s_θ .

Output: Graph Lottery Ticket (GLT) $f(\{m_g \odot A, \mathbf{X}\}, m_\theta \odot \Theta_0)$

- 1: **while** $1 - \frac{\|m_g\|_0}{\|A\|_0} < s_g$ **and** $1 - \frac{\|m_\theta\|_0}{\|\Theta_0\|_0} < s_\theta$ **do**
 - 2: Sparsify GNN $f(\cdot, m_\theta^i \odot \Theta_i)$ with $\mathcal{G} = \{m_g^i \odot A, \mathbf{X}\}$ according to Algorithm 1
 - 3: Update m_g and m_θ accordingly
 - 4: Rewind GNN weights to Θ_0
 - 5: Rewind masks, $m_g = m_g \odot A$
 - 6: **end while**
-

By iteratively training and pruning both the graph to a specified level of sparsity, we gain computational benefit as training progresses and a sparse graph and neural network by the end of the training process. An illustration of the process is given in Figure 2.2

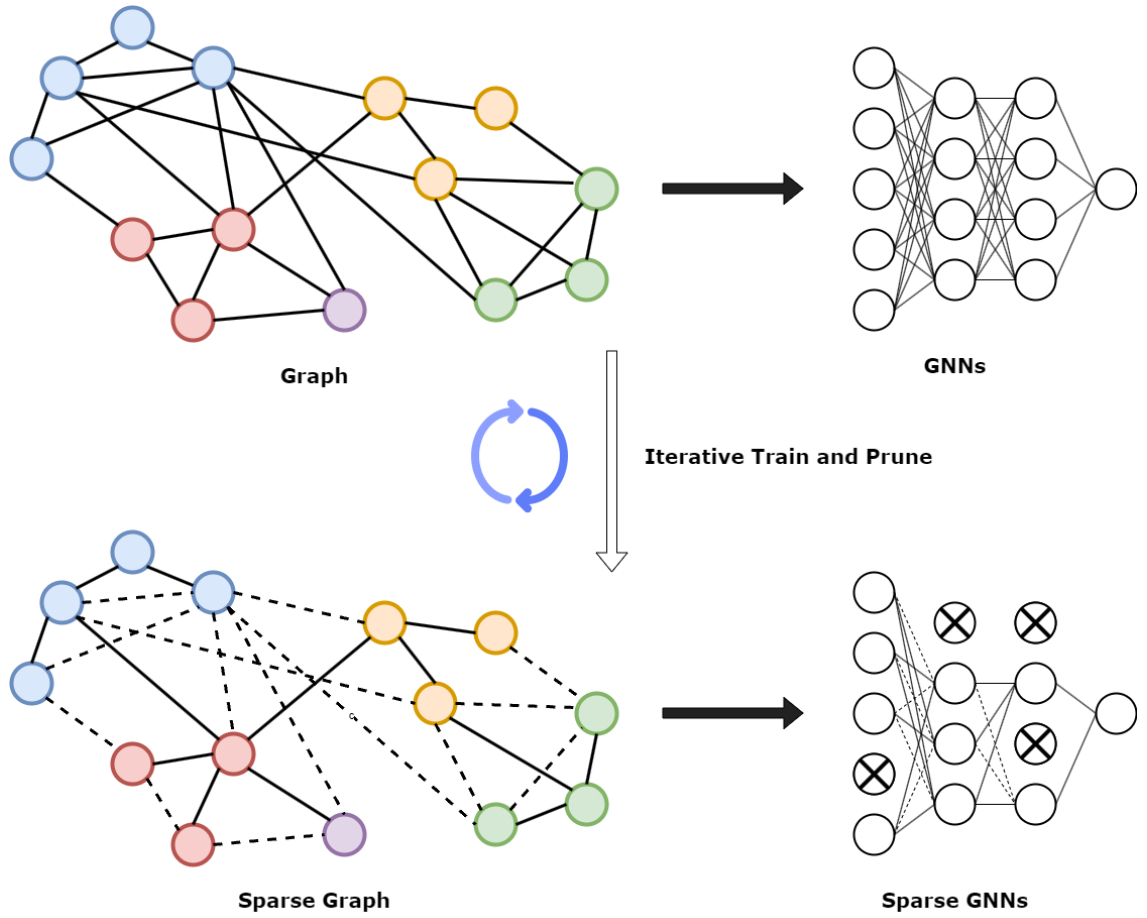


Figure 2.2: Illustration of the Unified Graph Sparsification.

2.4 Variational Autoencoder

This section should be done with more time if the original sparsification section is completed. Show how this is an alternative in graph discovery.

CHAPTER 3

Framework

3.1 Climate Models

Talk about the training time of one model ACCESS CM2 and go into detail about it.

An example of one of the models used in the CMIP6 Project is the Australian Community Climate and Earth System Simulator (ACCESS). [5] The core components that make up ACCESS include the UK Met Office Unified Model (UM) for the atmosphere [51], the CSIRO Community Atmosphere Biosphere Land Exchange (CABLE) model for land interactions [28], the United States' National Oceanic and Atmospheric Administration (NOAA)/ Geophysical Fluid Dynamics Laboratory (GFDL) Modular Ocean Model (MOM) for the ocean [2] and the Los Alamos National Laboratory CICE model for sea ice. [19] These components are shared between many of the CMIP6 models but there are also many locally produced components that may affect the dependence structure.

ACCESS was trained using the National Computational Infrastructure (NCI) with up to 900 cores and took up to 2-3 weeks for 100 model years.

(Maybe add something about needing a more scalable method to determine this kind of dependence?)

3.2 Dataset

The datasets used are from the CMIP6 scenario runs made available on the KNMI Climate Explorer website <https://climexp.knmi.nl/>. The KNMI is part of the World Meteorological Organization (WMO) and the scenario runs include monthly predictions for temperature, min temperature, max temperature precipitation, radiation and pressure all over the globe in a 192×144 grid between 1850–2100 for 40 different models from different institutions likely with some clusters of dependence. Due to the time and computational limitations, this was filtered to just temperature during the 1960–1980 period in just Australia or latitudes -44° to -12° and longitudes 288° to 336° . We chose filter the data down to these boundaries also since the goal of this thesis is to produce the optimal sparsification of the graph rather than producing the more performant model. The 1960–1980 period has the most stable data between the models and we also have historical truth observations for this time period. We chose scenario 245 as it is the middle of the road scenario and as such we should not see the greatest variance which should allow the sparsification algorithm to detect dependence more easily. An enumerated table of the models used are listed in Table 3.1

Model Number	Model Name	Model Number	Model Name
1	ACCESS-CM2	21	GFDL-CM4
2	ACCESS-ESM1-5	22	GFDL-ESM4
3	AWI-CM-1-1-MR	23	GISS-E2-1-G p1
4	BCC-CSM2-MR	24	GISS-E2-1-G p3
5	CAMS-CSM1-0	25	HadGEM3-GC31-LL f3
6	CanESM5 p1	26	HadGEM3-GC31-MM f3
7	CanESM5 p2	27	INM-CM4-8
8	CanESM5-CanOE p2	28	INM-CM5-0
9	CESM2	29	IPSL-CM6A-LR
10	CESM2-WACCM	30	KACE-1-0-G
11	CIESM	31	MCM-UA-1-0
12	CMCC-CM2-SR5	32	MIROC6
13	CNRM-CM6-1 f2	33	MIROC-ES2L f2
14	CNRM-CM6-1-HR f2	34	MPI-ESM1-2-HR
15	CNRM-ESM2-1 f2	35	MPI-ESM1-2-LR
16	EC-Earth3	36	MRI-ESM2-0
17	EC-Earth3-Veg	37	NESM3
18	FGOALS-f3-L	38	NorESM2-LM
19	FGOALS-g3	39	NorESM2-MM
20	FIO-ESM-2-0	40	UKESM1-0-LL f2

Table 3.1: CMIP6 Models used in this thesis. Model number is assigned for node numbering purposes.

3.3 Exploratory Data Analysis

Scenario 245 temperature at sea level shows a sigmoidal shaped rise. Properties are summarised in Figure 3.1

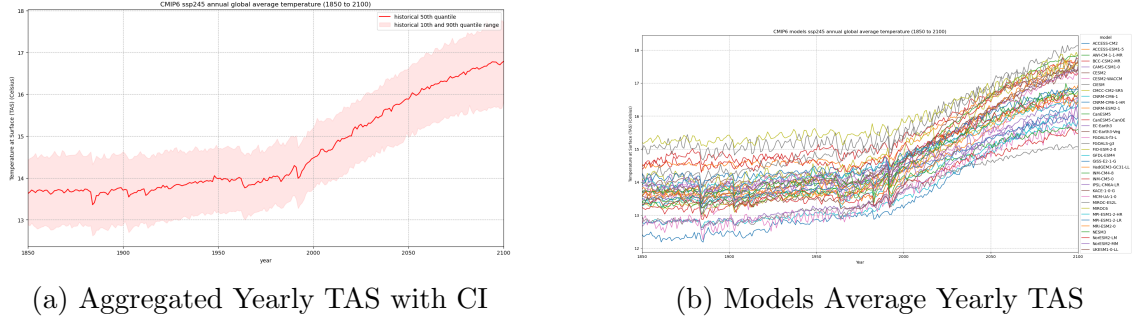


Figure 3.1: Visualisation of Climate Change in CMIP6 Models

A spatial plot of the aggregated averages of the models over specific time slices from 1850–2100 in Figure 3.2 show clear variations especially at the poles and extremities.

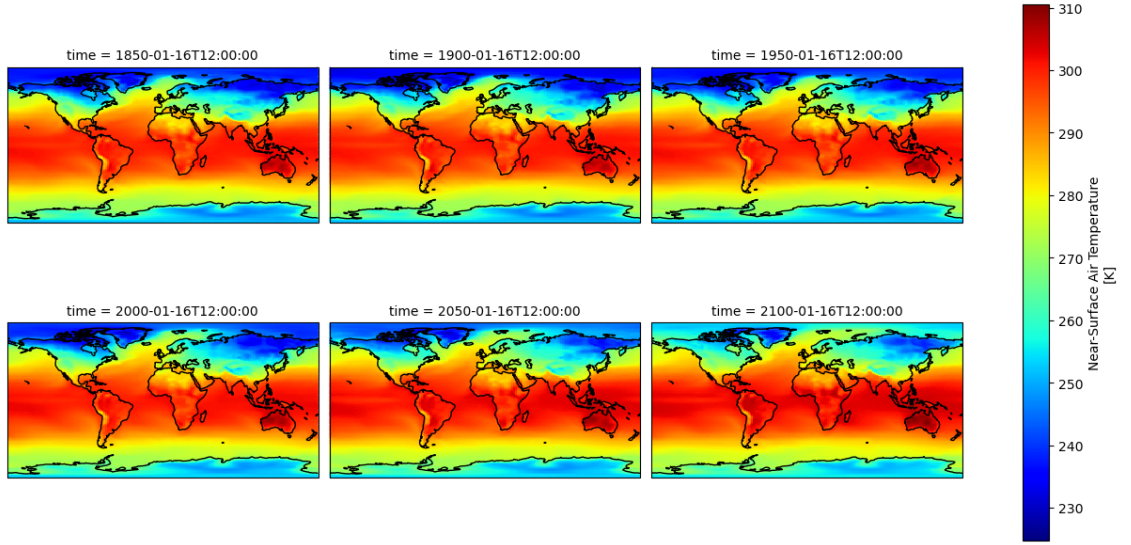


Figure 3.2: Plot of model aggregated temperatures at time slices from 1850-2100

Looking at a plot of all the models, there is a clear correlation between all the models and the correlation heatmap affirms this as the correlation between each ranges from 0.96–1.

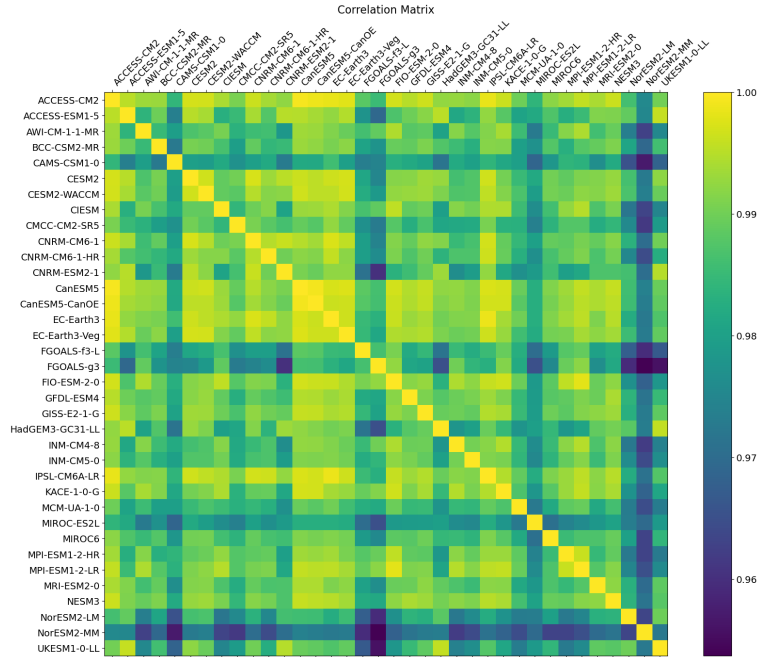


Figure 3.3: Correlation Heatmap of CMIP6 Models

Suppose we have a set of random variables $\mathbf{V} = \mathbf{X}_1, \dots, \mathbf{X}_n$. Partial correlation measures the correlation between two variables after removing the effect of all other random variables. [23] For example, with three variables the partial correlation of \mathbf{X}_i and \mathbf{X}_j given \mathbf{X}_k is,

$$r_{ij|k} = \frac{r_{ij} - r_{ik}r_{jk}}{\sqrt{1 - r_{ik}^2}\sqrt{1 - r_{jk}^2}} \quad (3.3.1)$$

where $r_{ij} = \text{corr}(x_i, x_j)$. Like correlation, this value ranges from 1 to -1 to indicate a strong positive or negative relationship respectively. We see in the partial correlation plot in Figure 3.4 that most partial correlation value are near zero but there are some notable exceptions. The diagonal is expectedly equal to one and the models produced by the same institution such as the CanESM5 (models 5, 6 and 7) or EC models (models 16 and 17) have very high partial correlation. The heatmap shows some unexpected dependences between unrelated models and even some high negative correlations between UKESM1 (model 40) and NESM3 (model 37) for example. The high values of correlation and low values of partial correlation suggests that while the relations may not be causal, there is a high degree of dependence between the models.

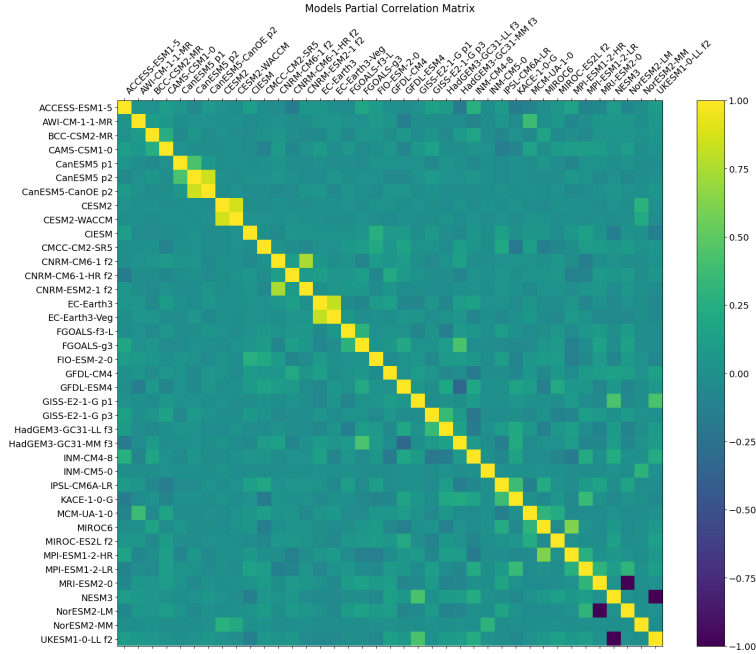


Figure 3.4: Partial Correlation Heatmap of CMIP6 Models

3.4 Problem Formulation

Maybe add this to the background section and put more of the regression, diagrams of the process and shrinkage prior stuff here that it more specific to this thesis

The goal in this thesis is to find some graph dependence structure between these models and this must be done through a multivariate regression task. We take the model outputs from our $n = 40$ models at spatially and temporally uniform times across each model and let the node embedding for node i in our GNN be $x_i \in \mathbb{R}^p$ where the size of p is

$$p = \text{no. of latitude points} \times \text{no. of longitude points} \times \text{no. of time periods}.$$

We begin with a fully connected graph with each node being connected to every other node by an undirected edge. By performing the sparsification algorithm, we would expect the redundant edges in this full connected graph to be set to zero. Finally, our Y is the historical observations from the 1960–1980 period and the final regression problem can be formulated as

$$f : \mathcal{G} \times X \rightarrow Y \quad (3.4.1)$$

where f denotes the learning function, \mathcal{G} the graph, $X \in \mathbb{R}^{n \times p}$ denotes the time series input and $Y \in \mathbb{R}^p$ the regression target. The shape of this Y makes this a graph-level regression task. For this task, we minimise the Mean Squared Error

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)$$

During the training process we perform a similar process to the Algorithm 1 to find the graph lottery ticket but instead of a dropout-like process of setting the lowest magnitude weights to zero, we use a shrinkage prior in the form of setting

$$m_g^{i+1} = \min(|m_g^i - \eta \nabla_{m_g^i} \mathcal{L}|, J_n) \quad (3.4.2)$$

where $J_n \in 1^{n \times n}$ represents the matrix of all-ones. This artificially restricts the edge weight $e_{ij} \in [0, 1]$.

Need to describe the math behind graph sparsification. More about shrinkage see Xiongwen.

3.5 Implementation

Python was used with Pytorch Geometric (PYG) to implement the models computationally. Pytorch is the the basis behind PYG and many of the numerical calculations are done with the help of Numpy. The data was cleaned and formatted with Pandas and Xarray from netcdf (.nc) data files. The models were training and stored on my personal computer with an AMD Ryzen 5 3600 (3.6 GHZ) 6 core 12 thread CPU, 32GB RAM and a NVIDIA RTX3080 (12GB) GPU.

CHAPTER 4

Results

4.1 Model verification

If the VGAE section is completed, we can compare the sparsified graph with the VGAE produced graph to determine how good graph sparsification is when used for graph discovery and thereby correlation in a graph structure.

4.2 Model results

Is there some way we can test the model's results depending on how sparse we make the graph etc. Research required to find some quantitative measure for this.

Some figures of the NN structure would also be helpful for this. Need to use nx or some other graph representation tool in python for this.

CHAPTER 5

Discussion

CHAPTER 6

Appendix

References

- [1] G. Abramowitz and C. H. Bishop. Climate model dependence and the ensemble dependence transformation of cmip projections. *Journal of Climate*, 28(6):2332 – 2348, 2015.
- [2] Alistair Adcroft, Whit Anderson, V Balaji, Chris Blanton, Mitchell Bushuk, Carolina O Dufour, John P Dunne, Stephen M Griffies, Robert Hallberg, Matthew J Harrison, et al. The gfdl global ocean and sea ice model om4. 0: Model description and simulation features. *Journal of Advances in Modeling Earth Systems*, 11(10):3167–3211, 2019.
- [3] Ahsan Adeel, Mandar Gogate, and Amir Hussain. Contextual deep learning-based audio-visual switching for speech enhancement in real-world environments. *Information Fusion*, 59:163–170, 2020.
- [4] Laith Alzubaidi, Jinglan Zhang, Amjad J Humaidi, Ayad Al-Dujaili, Ye Duan, Omran Al-Shamma, José Santamaría, Mohammed A Fadhel, Muthana Al-Amidie, and Laith Farhan. Review of deep learning: Concepts, cnn architectures, challenges, applications, future directions. *Journal of big Data*, 8:1–74, 2021.
- [5] Daohua Bi, Martin Dix, Simon Marsland, Siobhan O’farrell, Arnold Sullivan, Roger Bodman, Rachel Law, Ian Harman, Jhan Srbinovsky, Harun A Rashid, et al. Configuration and spin-up of access-cm2, the new generation australian community climate and earth system simulator coupled model. *Journal of Southern Hemisphere Earth Systems Science*, 70(1):225–251, 2020.
- [6] Yuan Cao and Quanquan Gu. Generalization error bounds of gradient descent for learning over-parameterized deep relu networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):3349–3356, Apr. 2020.
- [7] Tianlong Chen, Yongduo Sui, Xuxi Chen, Aston Zhang, and Zhangyang Wang. A unified lottery ticket hypothesis for graph neural networks, 2021.
- [8] Li Deng and Dong Yu. *Deep Learning: Methods and Applications*. Now Foundations and Trends, 2014.
- [9] Paul N Edwards. History of climate modeling. *Wiley Interdisciplinary Reviews: Climate Change*, 2(1):128–139, 2011.
- [10] Jinming Feng, Dong-Kyou Lee, Congbin Fu, Jianping Tang, Yasuo Sato, Hisashi Kato, John L McGregor, and Kazuo Mabuchi. Comparison of four ensemble methods combining regional climate simulations over asia. *Meteorology and Atmospheric Physics*, 111:41–53, 2011.
- [11] LLDiko E Frank and Jerome H Friedman. A statistical view of some chemometrics regression tools. *Technometrics*, 35(2):109–135, 1993.
- [12] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M. Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis, 2020.

- [13] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [14] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
- [15] Thomas Harrisson. Explainer: How ‘shared socioeconomic pathways’ explore future climate change, Apr 2018.
- [16] Thomas Harrisson. Cmp6: The next generation of climate models explained, Oct 2021.
- [17] Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):251–257, 1991.
- [18] D. H. Hubel and T. N. Wiesel. Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex. *The Journal of Physiology*, 160(1):106–154, 1962.
- [19] Elizabeth C Hunke, William H Lipscomb, Adrian K Turner, Nicole Jeffery, and Scott Elliott. Cice: the los alamos sea ice model documentation and software user’s manual version 4.1 la-cc-06-012. *T-3 Fluid Dynamics Group, Los Alamos National Laboratory*, 675:500, 2010.
- [20] W.H. Inmon and A. Nesavich. *Tapping into Unstructured Data: Integrating Unstructured Data and Textual Analytics into Business Intelligence*. Pearson Education, 2007.
- [21] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in neural information processing systems*, 30, 2017.
- [22] Viatcheslav V Kharin and Francis W Zwiers. Climate predictions with multi-model ensembles. *Journal of Climate*, 15(7):793–799, 2002.
- [23] Seongho Kim. ppcor: an r package for a fast calculation to semi-partial correlation coefficients. *Communications for statistical applications and methods*, 22(6):665, 2015.
- [24] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks, 2017.
- [25] Reto Knutti, Gab Abramowitz, M. Collins, Veronika Eyring, Peter Gleckler, Bruce Hewitson, and Linda Mearns. Good practice guidance paper on assessing and combining multi model climate projections. pages 1–15, 01 2010.
- [26] Vijay Kotu and Bala Deshpande. Chapter 2 - data science process. In Vijay Kotu and Bala Deshpande, editors, *Data Science (Second Edition)*, pages 19–37. Morgan Kaufmann, second edition edition, 2019.
- [27] Max Kulinich. *A Markov chain method for weighting climate model ensembles and uncertainty estimation on spatially explicit data*. PhD thesis, UNSW, 2022.
- [28] Rachel Law, Michael Raupach, Gabriel Abramowitz, Imtiaz Dharssi, Vanessa Haverd, Andrew Pitman, Luigi Renzullo, Albert Van Dijk, and Yingping Wang. *The Community Atmosphere Biosphere Land Exchange (CABLE) model Roadmap for 2012-2017*. Centre for Australian Weather and Climate Research, 2012.
- [29] June-Yi Lee, Jochem Marotzke, Govindasamy Bala, Long Cao, Susanna Corti, John P Dunne, Francois Engelbrecht, Erich Fischer, John C Fyfe, Christopher Jones, et al. Future global climate: scenario-based projections and near-term

- information. In *Climate change 2021: The physical science basis. Contribution of working group I to the sixth assessment report of the intergovernmental panel on climate change*, pages 553–672. Cambridge University Press, 2021.
- [30] Christos Louizos, Karen Ullrich, and Max Welling. Bayesian compression for deep learning. *Advances in neural information processing systems*, 30, 2017.
 - [31] Rongrong Ma, Jianyu Miao, Lingfeng Niu, and Peng Zhang. Transformed l1 regularization for learning sparse deep neural networks. *Neural Networks*, 119:286–298, 2019.
 - [32] Dmitry Molchanov, Arsenii Ashukha, and Dmitry Vetrov. Variational dropout sparsifies deep neural networks. In *International Conference on Machine Learning*, pages 2498–2507. PMLR, 2017.
 - [33] Guido F Montufar, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio. On the number of linear regions of deep neural networks. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
 - [34] Kevin P. Murphy. *Probabilistic Machine Learning: An introduction*. MIT Press, 2022.
 - [35] Eric Nalisnick, José Miguel Hernández-Lobato, and Padhraic Smyth. Dropout as a structured shrinkage prior. In *International Conference on Machine Learning*, pages 4712–4722. PMLR, 2019.
 - [36] Tomas Nordström and Bertil Svensson. Using and designing massively parallel computers for artificial neural networks. *Journal of parallel and distributed computing*, 14(3):260–285, 1992.
 - [37] Sara Nouri Golmaei. *Improving the Performance of Clinical Prediction Tasks by Using Structured and Unstructured Data Combined with a Patient Network*. PhD thesis, 2021.
 - [38] Jonathan Overpeck, Gerald Meehl, Sandrine Bony, and David Easterling. Climate data challenges in the 21st century. *Science (New York, N.Y.)*, 331:700–2, 02 2011.
 - [39] Tim Palmer and Phoebe Williams. Stochastic physics and climate models. *Philosophical transactions. Series A, Mathematical, physical, and engineering sciences*, 366:2421–7, 04 2008.
 - [40] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
 - [41] Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 2847–2854. PMLR, 06–11 Aug 2017.
 - [42] Bhavtosh Rath, Aadesh Salecha, and Jaideep Srivastava. Detecting fake news spreaders in social networks using inductive representation learning. In *2020 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, pages 182–189. IEEE, 2020.
 - [43] Keywan Riahi, Detlef P. van Vuuren, Elmar Kriegler, Jae Edmonds, Brian C. O’Neill, Shinichiro Fujimori, Nico Bauer, Katherine Calvin, Rob Dellink, Oliver Fricko, Wolfgang Lutz, Alexander Popp, Jesus Crespo Cuaresma, Samir

- KC, Marian Leimbach, Leiwen Jiang, Tom Kram, Shilpa Rao, Johannes Emmerling, Kristie Ebi, Tomoko Hasegawa, Petr Havlik, Florian Humpenöder, Lara Aleluia Da Silva, Steve Smith, Elke Stehfest, Valentina Bosetti, Jiyong Eom, David Gernaat, Toshihiko Masui, Joeri Rogelj, Jessica Strefler, Laurent Drouet, Volker Krey, Gunnar Luderer, Mathijs Harmsen, Kiyoshi Takahashi, Lavinia Baumstark, Jonathan C. Doelman, Mikiko Kainuma, Zbigniew Klimont, Giacomo Marangoni, Hermann Lotze-Campen, Michael Obersteiner, Andrzej Tabeau, and Massimo Tavoni. The shared socioeconomic pathways and their energy, land use, and greenhouse gas emissions implications: An overview. *Global Environmental Change*, 42:153–168, 2017.
- [44] F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6):386–408, 1958.
 - [45] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
 - [46] Benjamin Sanchez-Lengeling, Emily Reif, Adam Pearce, and Alexander B. Wiltschko. A gentle introduction to graph neural networks. *Distill*, 2021. <https://distill.pub/2021/gnn-intro>.
 - [47] Simone Scardapane, Danilo Comminiello, Amir Hussain, and Aurelio Uncini. Group sparse regularization for deep neural networks. *Neurocomputing*, 241:81–89, 2017.
 - [48] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
 - [49] Haiman Tian, Shu-Ching Chen, and Mei-Ling Shyu. Evolutionary programming based deep learning feature selection and network construction for visual data classification. *Information systems frontiers*, 22:1053–1066, 2020.
 - [50] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks, 2018.
 - [51] David Walters, Anthony J Baran, Ian Boutle, Malcolm Brooks, Paul Earnshaw, John Edwards, Kalli Furtado, Peter Hill, Adrian Lock, James Manners, et al. The met office unified model global atmosphere 7.0/7.1 and jules global land 7.0 configurations. *Geoscientific Model Development*, 12(5):1909–1963, 2019.
 - [52] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. *Advances in neural information processing systems*, 29, 2016.
 - [53] Katherine Xu and Janice Lan. Chemistry insights for large pretrained GNNs. In *NeurIPS 2022 AI for Science: Progress and Promises*, 2022.
 - [54] Tom Young, Devamanyu Hazarika, Soujanya Poria, and Erik Cambria. Recent trends in deep learning based natural language processing. *ieee Computational intelligence magazine*, 13(3):55–75, 2018.
 - [55] Yue Yu, Jie Chen, Tian Gao, and Mo Yu. DAG-GNN: DAG structure learning with graph neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7154–7163. PMLR, 09–15 Jun 2019.
 - [56] Muhan Zhang, Pan Li, Yinglong Xia, Kai Wang, and Long Jin. Revisiting graph neural networks for link prediction. 2020.