

# Chapter 1

## Introduction

Brain ageing depends on many complex factors related to both genetics and postnatal environment. It is associated with cellular damage, changes in morphological and functional brain connectivity, changes in gray and white matter structure, frontal lobe thinning and other factors that result in cognitive decline, decreased memory capacity and a variety of neurological and neurodegenerative disorders such as Alzheimer’s and Parkinson’s diseases [2, 8, 4]. With the increasing global population lifespan, age-related brain disorders have become more common, motivating the research in biomarkers to represent an individual’s brain health status, reveal common risk factors, give insight into the mechanisms of brain ageing, potentially delaying the onset of the disease, slowing it down and reducing its prevalence. With the increased availability of brain imaging data through the longitudinal population-wide studies (such as the UK Biobank), and increasing computing power (enabling the analysis of rich and otherwise computationally intractable datasets), brain ageing is becoming a popular research direction.

One biomarker for estimating the overall brain health is the so-called *brain age*, defined as the *apparent* age of an individual’s brain compared to the typical population.

and derived from the neuroimaging analysis of the patient, most commonly magnetic resonance imaging (MRI) data. which is a *prediction* of a person’s chronological (actual) age .

Recent studies have linked the deviations between the brain age estimate and the true chronological age (the *brain age gap*) to the occurrence of various neurological conditions such as multiple sclerosis and dementia [6], presumably because the overestimate of the chronological age indicates the accelerating ageing and higher cellular damage accumulation of the brain.

Population graphs, where nodes represent the neuroimaging data and edges represent some similarity metric between patients is a promising approach to analysing those datasets, as it allows to leverage both the individual subject data as well as the outcomes of similar patients, just like in diagnosing patients clinicians would look at the examples of other patients who had those particular symptoms and which disease they were diagnosed with.

This dissertation proposes applying the (novel?) population graph paradigm to the brain age estimation task.

## 1.1 Related work

### 1.1.1 Population graphs

Semisupervised population graphs in Parisot et al. for classifying healthy patients and patients having Mild Cognitive Impairment/Alzheimer’s or autism spectrum disorder. These achieve state-of-the-art accuracies around 70% even for a binary classification task which indicates that brain conditions are generally a complex task. Due to the different nature of the task and different evaluation metrics, this performance cannot be directly compared to the predictive power of the brain age regression task.

### 1.1.2 Brain age prediction from neuroimaging data

Most of current machine learning-based methods for brain age gap prediction work on per-brain basis and do not consider pairwise similarities between patients, not taking into account the population as a whole.

Graphs are static so it makes sense that most practical approaches would not use this as their model, because it prevents using applying it for new patients

Kaufmann et al. [6] uses gradient boosting based techniques (XGBoost) [1] for brain age gap prediction from structural magnetic resonance imaging (MRI) data. The study uses a significantly larger dataset of 45,000 people and presents separate models for female and male brain age gap prediction, without considering any pairwise similarities between individuals.

Another framework for brain age gap prediction is based on Gaussian Processes regression, using raw T1-weighted MRI scans, segmenting them and using principal components analysis (PCA) for dimensionality reduction [2].

# Chapter 2

## Preparation

This section presents the mathematical definition of the population graphs, the neuroimaging data those graphs will incorporate, and the neural network architectures that will be used to train the population graphs.

### 2.1 Brain age estimation

To estimate the brain age gap the difference between the chronological and the brain age must be known. The common technique is to develop a machine learning method to estimate the chronological age from the brain imaging features. The estimated age is then considered to be the brain age although it has been trained to estimate the chronological age and the actual value of the brain age is never known. The reason is why it works... [8]

We represent the brain age  $y_b$  as the sum of the known chronological age  $y_c$  and the unknown brain age gap  $\varepsilon_g$ :

$$y_b = y_c + \varepsilon_g. \quad (2.1)$$

On the other hand, the machine learning model estimates the brain age  $y_b$  from some function of the brain imaging features  $X$  and a prediction error  $\varepsilon_e$ :

$$y_b = f(X) + \varepsilon_e. \quad (2.2)$$

Expressing  $y_c$  as  $y_b - \varepsilon_g$  and substituting the previous result, we get the estimate of chronological age as

$$y_c = f(X) + \varepsilon_g - \varepsilon_e, \quad (2.3)$$

Since the machine learning models estimate chronological age as a function of brain imaging features (with some estimation error  $\varepsilon'_e$ )

$$y_c = f(X) + \varepsilon'_e, \quad (2.4)$$

the error of the proposed machine learning model  $\varepsilon'_e$  will contain in itself both the brain age gap  $\varepsilon_g$  (which we are interested in) and the brain age prediction error  $\varepsilon_e$ .

It is out of scope of this dissertation to prove that the brain age gap component  $\varepsilon_g$  is larger than the error term; however, the keen reader is referred to Niu et al. [8] where this is verified through correlation of the brain age gaps to the cognitive behaviour scores.

## 2.2 Neuroimaging dataset

### 2.2.1 United Kingdom Biobank

The United Kingdom Biobank (UK Biobank) [9] is a continuous, large, population-wide study of over 500,000 participants containing a wide range of phenotypic and genetic data that is used by the researchers to analyse the risk factors and development of various health conditions.

Of particular relevance to this dissertation are the UK Biobank participants with neuroimaging data records that have been denoised, motion-corrected and otherwise processed for further analysis, a total of 17,550 participants. The data has been initially preprocessed with the standard UK Biobank pipelines,<sup>1</sup> and further denoised, parcellated and kindly provided by Dr Richard Bethlehem of the Department of Psychiatry. The Department of Psychiatry pipelines have been co-authored with Dr Rafael Romero-Garcia and Dr Lisa Ronan.<sup>2</sup> The details on the neuroimaging dataset are described below.

### 2.2.2 Parcellation

A *parcellation* or *atlas* refers to the way the brain is split into meaningful regions for further analysis. Whether two voxels of a brain belong to the same parcel may depend on their proximity, empirical evidence of that the voxels are responsible for the same function and so forth, and can be used to compare the locations in two different brains. When the brain is imaged there is a choice whether to warp the image of the brain to the fixed atlas or whether to warp the atlas to match the variable brain images. The former makes it easier to process a dataset of many images and find the matching regions of two brains faster, but the latter remains more faithful to the unique structure of the individual patient's brain.

---

<sup>1</sup>[https://biobank.ctsu.ox.ac.uk/crystal/crystal/docs/brain\\_mri.pdf](https://biobank.ctsu.ox.ac.uk/crystal/crystal/docs/brain_mri.pdf)

<sup>2</sup><https://github.com/ucam-department-of-psychiatry/UKB>

Both functional and structural datasets use one of the most common parcellations by Glasser et al. [5], which divides the brain into 360 cortical regions and 16 subcortical regions.

### 2.2.3 Structural data

Structural brain imaging data encapsulates brain features related to its structure, such as cortical thickness, white

Combined T1-weighted and T2-weighted FLAIR images which emphasise different aspects of the MRI scan and therefore might help to extract different structural features. For the dataset used in this dissertation, the combination of the two types of images to derive the structural features using the HCP Freesurfer pipeline.<sup>3</sup>

### 2.2.4 Euler indices

Euler index<sup>4</sup> is a quality control metric and corresponds to the number of times the Freesurfer brain reconstruction software failed to connect two slices of an MRI image. The higher the Euler index, the worse is the quality of the scan. Euler indices might be used to remove the subjects with low-quality scans to avoid them affecting the analysis [6]. Otherwise they can be used as a covariate in a machine learning model (as brain similarity metric or a node feature) to correct for any bias in prediction that might be related to scan quality.

### 2.2.5 Functional data

The resting state functional MRI (rs-fMRI) is the representation of the brain activity over time. In MRI scanner this is measured by the changes in blood oxygenation as neural activity regulates the oxygen demand, resulting in blood oxygenation level dependent (BOLD) time-series measured at each voxel of the brain.

*TODO figure of BOLD timeseries*

We are interested in estimating which parts of the brain are connected to each other, which we do by making use of the assumption that parts of the brain that have related function would also have similar activity patterns. As a consequence, we would expect higher correlation of the corresponding BOLD time-series. For time-series  $T_1$  and  $T_2$ , *Pearson's correlation* (denoted as  $r$ ) is computed as

$$r(T_1, T_2) = \frac{\text{cov}(T_1, T_2)}{\sigma_{T_1} \sigma_{T_2}} \quad (2.5)$$

where  $\text{cov}(\cdot, \cdot)$  denotes covariance and  $\sigma$  stands for standard deviation.

---

<sup>3</sup><https://www.ncbi.nlm.nih.gov/pubmed/23668970>

<sup>4</sup><https://www.ncbi.nlm.nih.gov/pubmed/29278774>

This (or other correlation types) to derive the *functional connectivity matrix* that stores the pairwise correlations between the different voxels as the overall representation of functional brain connectivity. For time-series  $T_1, \dots, T_N$ ,

$$\text{fcm}(T_1, \dots, T_N) = \begin{bmatrix} r(T_1, T_1) & \cdots & r(T_1, T_N) \\ \vdots & \ddots & \vdots \\ r(T_N, T_1) & \cdots & r(T_N, T_N) \end{bmatrix}, \quad (2.6)$$

of which (due to symmetry and the non-informative diagonal) only the flattened lower triangle is usually used for the machine learning implementations.

### 2.2.6 Phenotype data

For the population graph construction, the neuroimaging features are associated with the individual subjects (nodes). The similarity metric is defined by phenotype data, which is all important but not directly neuroimaging-related data. Some examples of phenotype data that is related to the brain include the patient's sex, mental health, other potential health issues, full-time education, bipolar disorder status etc. Indeed the metric that is being predicted is correlated with the brain tissue age indicative of various neurological and neurodegenerative diseases.

TODO a table of the actual phenotypes used when the results are ready

## 2.3 Population graphs

We connect the multi-modal MRI imaging (structural and/or functional), quality control and phenotype data of the set of patients  $S$  into a sparse *population graph*  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ , where  $\mathcal{V}$  is the set of graph vertices (with one vertex uniquely representing one patient), and  $\mathcal{E}$  is the set of edges (representing the *similarity* of patients).

Each vertex  $v \in \mathcal{V}$  is the vector containing the individual subject's neuroimaging data, whether structural, functional, or both. The edge  $(v, w) \in \mathcal{E}$  connects patients  $v$  and  $w$  based on phenotypic similarity that is defined by some *similarity metric* as described below.

### 2.3.1 Similarity metrics

The topology of the graph is determined by a similarity metric that uses the non-imaging (phenotypic) information of the subjects to create connections between the nodes containing brain imaging data. Defining a good similarity metric is important to correct for the confounding effects on the feature vectors (for example, subject's sex has an impact on the brain volume) as well as to cluster patients into the most informative neighbourhoods where the predicted variable (brain age) explains the most variance in node features.

Similarity metrics are defined using some *similarity function*  $\text{sim}(\cdot, \cdot)$  which takes two subjects and returns the similarity score between them (the higher the score, the more similar the subjects):

$$\text{sim}(S_v, S_w) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}[M_i(v) = M_i(w)]. \quad (2.7)$$

Here  $\{M_1, \dots, M_n\}$  is a set of phenotypic metrics that are used in computing subject similarity (such as sex, mental health, fluid intelligence score etc.) and  $\mathbf{1}[\cdot]$  is an indicator function, in this case returning a non-zero value when values for a given phenotypic feature  $M_i$  match for two subjects  $S_v$  and  $S_w$ .

TODO *Qualitative (categorical) metrics can be defined as Kronecker delta  $\delta$  and quantitative (numerical) metrics have an indicator whether the distance between values  $M_i(v)$  and  $M_i(w)$  is within some  $\epsilon$ .*

To avoid memory issues when  $|\mathcal{E}| \sim O(|\mathcal{V}|^2)$  and minimise the size of the neighbourhood to only the highly similar subjects, a *similarity threshold*  $\theta$  is defined such that

$$(S_v, S_w) \in \mathcal{E} \iff \text{sim}(S_v, S_w) \geq \theta. \quad (2.8)$$

### 2.3.2 Training task

[7]: The task is to predict the label of the nodes, where the labels are visible only for the training and validation but not the test nodes. Semi-supervised learning allows the label information and parameter weights to spread to node neighbourhoods to the similar nodes as defined by the similarity metric,

A population graph is trained in a *semi-supervised* manner. This means that the entire dataset (all nodes and all labels) is included in the graph. At each training step (*epoch*), the parameters for all nodes are updated but only based on the feedback from training and validation nodes. The final predictive power of the model is evaluated based on its performance on previously hidden test node labels.

Labels are available for a small subset of nodes and are spread across neighbourhoods through a regularisation term such as Laplacian [7], predicting the labels for the remaining nodes

TODO *an illustration of the graph with marked training and validation nodes with visible labels and how the training happens on those followed by how the graph is evaluated on the test nodes where the parameters were updated but the labels never seen*

## 2.4 Graph convolutional networks

TODO Graph definition through degree and adjacency matrices, the weights of the edges being binary (either no edge or a weight-1 edge). Graph Laplacian. Diagonalisation of

Laplacian matrix because it's positive semi-definite. Spatial vs spectral graph domains. Graph Fourier transform. Fourier transform as transformation to the eigenbasis obtained through graph Laplacian diagonalisation. Note that the basis depends completely on the graph structure and slight permutation can change the eigenbasis. Graph convolution with a filter  $g$  as multiplication of signal and filter in the Fourier domain. Parametric and non-parametric convolution filters using this  $g_\theta$ ; what does  $\theta$  mean in relation to the Fourier coefficients; what exactly Fourier coefficients refer to. Eigenvalues as graph frequencies and learning as a low-pass filter of frequencies. Approximation of the Laplacian diagonalisation with Chebyshev polynomials (ChebNet) with the convolution defined for  $k$ -th order polynomials corresponding to  $k$  hops across neighbourhood. [3] Kipf and Welling simplifying ChebNet hops as stacking single-hop layers and using renormalisation trick for the node to use the most information from itself. [7] Averaging representations of neighbour features and smoothing labels as a consequence. [11] (this citation also very good for follow through for what the training is doing for all layers in general rather than a single feature vector, i.e. the training flow); can also mention [12]; convolution on regular vs irregular grids (i.e. intro to simple convolution before carrying on?); [3]: "convolution is linear operator that diagonalises in the fourier domain represented by the eigenvectors of the Laplacian operator"

[10]: graph definition in the Fourier domain based on the eigendecomposition of the graph Laplacian, filters further approximated by Chebyshev expansion of the graph Laplacian avoiding the expensive eigendecomposition operation (ChebNet); further improvements by filter restriction to local neighbourhoods and stacking them if necessary

### 2.4.1 Graph Fourier transform

## 2.5 Graph attention networks

[10] More based on the non-spectral (i.e. spatial) approaches where an operator is defined to work on the neighbourhoods of different sizes, maintaining the weight sharing property (usually done by defining a separate operator on the neighbourhoods of those sizes).

*Self-attention* is also added in to the concept where different parts of the node's neighbourhood are considered with different importance weights, getting a representation of the rest of the neighbourhood.

### 2.5.1 Graph attentional layer

For the layer of an  $N$ -node graph with  $F$  input and  $F'$  output features,

input node features  $\{\mathbf{h}_1, \dots, \mathbf{h}_N\}, \mathbf{h}_i \in \mathbb{R}^F$

output node features  $\{\mathbf{h}'_1, \dots, \mathbf{h}'_N\}, \mathbf{h}'_i \in \mathbb{R}^{F'}$

linear transformation with weight matrix  $\mathbf{W} \in \mathbb{R}^{F' \times F}$

self attention  $a : \mathbb{R}^{F'} \times \mathbb{R}^{F'} \rightarrow \mathbb{R}$



attention coefficients *can already include neighbourhoods*  $e_{ij} = a(\mathbf{W}\mathbf{h}_i, \mathbf{W}\mathbf{h}_j)$

weight vector  $\mathbf{a} \in \mathbb{R}^{2F'}$

$$\alpha_{ij} = \text{softmax}_j(\mathbf{a}^T[\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_j]) \quad (2.9)$$

$$= \frac{\exp(\sigma_1(\mathbf{a}^T[\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_j]))}{\sum_{k \in \mathcal{N}_i} \exp(\sigma_1(\mathbf{a}^T[\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_k]))} \quad (2.10)$$

where  $\alpha_{ij}$  is the attention coefficient for an edge  $i \rightsquigarrow j$  (corresponding to the importance of features in node  $j$  to the features in node  $i$ ), normalised across  $i$ 's neighbourhood  $\mathcal{N}_i$  (defined as e.g. all nodes one hop away);  $\sigma_1$  is a non-linearity.  $\parallel$  means concatenation

The coefficients  $\alpha_{ij}$  and the weight matrix are used to compute the output features:

$$\mathbf{h}'_i = \sigma\left(\sum_{k \in \mathcal{N}_i} \alpha_{ik} \mathbf{W}\mathbf{h}_k\right) \quad (2.11)$$

### 2.5.2 Multiple attention

The above attention mechanism can be repeated several times to stabilise the performance, where one independent application of attention is called an attention head. The outputs of the independent attention heads are concatenated together until the last layer when they are averaged into a single output. For  $K$  attention heads, the results of (2.11) are concatenated:

$$\mathbf{h}'_i = \parallel_{k=1}^K \sigma\left(\sum_{k \in \mathcal{N}_i} \alpha_{ik}^k \mathbf{W}^k \mathbf{h}_k\right) \quad (2.12)$$

Which is averaged in the last layer:

$$\mathbf{h}'_i = \sigma\left(\frac{1}{K} \sum_{k=1}^K \sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \mathbf{h}_j\right). \quad (2.13)$$

TODO Talk about the initial features, *diagram of the attention heads I guess etc.*

## 2.6 Requirements analysis

Tasks to be implemented (according to proposal: work to be done, success criteria, possible extensions), their relative importance (priority) and difficulty. Provide the order in which the tasks should be carried out to show good planning skills and account for the changes in proposal where the preprocessing pipeline turned out to be more important than the neural network implementation.

## 2.7 Software engineering practice

Implementing a flexible preprocessing pipeline which could be customised in the future for a variety of machine learning tasks even outside graph neural networks (a package).

Modular structure encapsulating specific task and having well defined documentations of the others.

Description of software engineering techniques: planning out and executing the project based on requirements analysis, setting tasks, and smoothly meeting the success criteria.

Code reuse (of open source well tested libraries), follow documentation and follow the PEP-8 style guide (or whatever PyCharm encourages).

Incremental development.

Modular structure: e.g. data processing, graph construction, graph neural network modules, robustness evaluation framework. Figure out where validation and cross validation sections should be (while training, separately etc.)

Diagram of the pipelines and module interaction (like in google design docs)

## 2.8 Choice of tools

PyTorch, PyTorch geometric extension, graph spectral filters/convolutions, message passing, time-series preprocessing into correlation matrices, IDEs, backup strategies

## 2.9 Starting point

- dataset, preprocessed by Dr Richard A.I. Bethlehem
- PyTorch, PyTorch geometric implementing GCN and GAT APIs and the graph API
- no previous experience with graph neural networks or the mathematics behind it
- no previous experience with PyTorch; limited experience with machine learning frameworks (basics of TensorFlow), no experience with neuroimaging data

# Chapter 3

## Implementation

Could be split into *preprocessing, parameter tuning, evaluation framework, software engineering techniques, repository overview*

### 3.0.1 Precomputation and preprocessing of connectivity matrices

Tangent works better than correlation or partial correlation.

### 3.0.2 Structural and functional data extraction

### 3.0.3 Graph construction pipeline

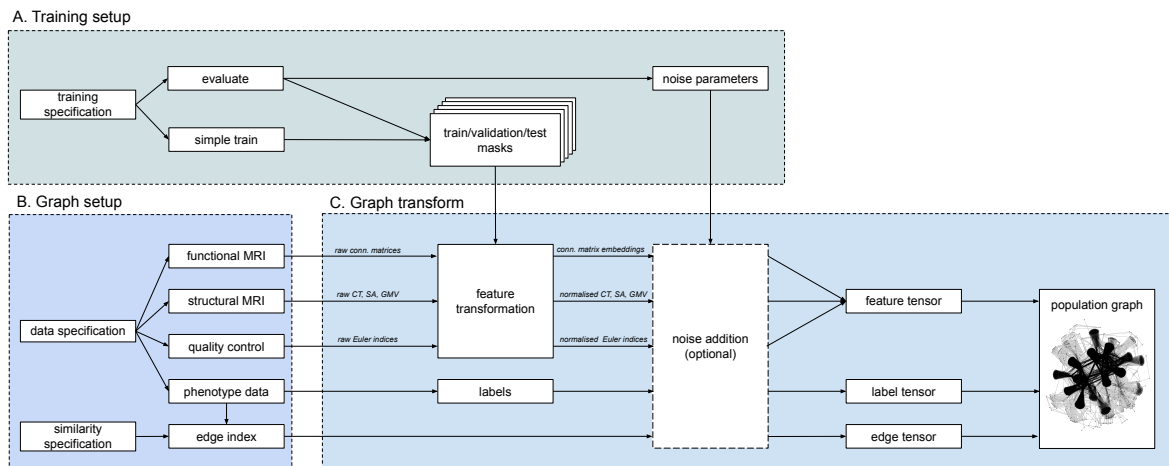


Figure 3.1: Graph pipeline.

### 3.0.4 Train, test, validation split

## 3.1 Non-graph baselines

I have played around with *xgboost*, ElasticNet and simple multilayer perceptrons to get a good idea of the minimum baselines my new architectures need to reach to be considered more effective.

## 3.2 Graph convolutional network

Describe the architecture in BrainGCN

## 3.3 Graph attention network

## 3.4 Robustness evaluation framework

## 3.5 Repository overview

Figure out how to split framework modules.

Could have a base BrainGNN class which can then be *extended* with BrainGAT, BrainGCN.

# Chapter 4

## Evaluation

### 4.1 Overview of the predictive power for the two graph neural network approaches

Discuss Pearson's  $r$ , coefficient of determination  $r^2$  and other common performance metrics.

### 4.2 Hyperparameters

Effects of hyperparameter tuning

### 4.3 Unit testing

### 4.4 Other quantitative and qualitative results

Statistical significance of the results compared to the baseline (e.g. Pearson's  $r$  method in Python seems to return some  $p$ -value of it..?)

### 4.5 Robustness of the graph neural networks to noisy and missing data

Add noise to the graphs and measure the rate of drop in predictive power.

### 4.6 Comparison against existing benchmarks

Compare to the Kaufmann et al.'s *xgboost* approach [6] ( $r \sim 0.93$ ); and the other package that was cited in the same paper.

Possibly compare to other non-graph (relatively baseline) (neural network) architectures, e.g. ElasticNet, MLP,...

## 4.7 Interpretation of the model behaviour

One of the advantages of the graphs that they *should* be interpretable I guess.

# Chapter 5

## Conclusion

### 5.1 Successes and failures

### 5.2 The project in hindsight (lessons learnt)

Lessons learnt: graph *representation* is more important than the framework used. Good representations of the dataset can help guide the learning algorithm in the right way as it gets the intuitions faster just because of the way the data was represented in the first place. It's not good if it captures bias, but in this case representation made the difference between the model not learning anything and the model getting great results.

Was a mistake analysing functional data first without analysing the other methods in detail. Functional imaging data by itself gave very high dimensionality which either could not be learnt by the network because of the low number of examples or other factors. Most of the literature uses just the structural data for age prediction, and indeed this turned out to be more effective. Also makes sense intuitively as structural features would be related to the signs of brain atrophy while it is not clear the pattern of how resting state brain activity would change with ageing brain.

Mention Niu et al. 2019 raising the issue that there is systematic bias in brain age gap prediction but not many studies use this knowledge to correct for it.

### 5.3 Possible continuations of the project

- Include DNA methylation data as it is widely used in other studies and is claimed to improve the predictive power of the model [2].





# Bibliography

- [1] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794. ACM, 2016.
- [2] James Cole, Stuart Ritchie, Mark Bastin, Maria Valdés Hernández, et al. Brain age predicts mortality. *Molecular Psychiatry*, 23(5):1385–1392, 2018.
- [3] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems*, pages 3844–3852, 2016.
- [4] Katja Franke and Christian Gaser. Ten years of BrainAGE as a neuroimaging biomarker of brain aging: What insights have we gained? *Frontiers in Neurology*, 10:789, 2019.
- [5] Matthew F Glasser, Timothy S Coalson, Emma C Robinson, Carl D Hacker, John Harwell, Essa Yacoub, Kamil Ugurbil, Jesper Andersson, Christian F Beckmann, Mark Jenkinson, et al. A multi-modal parcellation of human cerebral cortex. *Nature*, 536(7615):171, 2016.
- [6] Tobias Kaufmann, Dennis van der Meer, Nhat Trung Doan, Emanuel Schwarz, et al. Common brain disorders are associated with heritable patterns of apparent aging of the brain. *Nature Neuroscience*, 22(10):1617–1623, 2019.
- [7] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- [8] Xin Niu, Fengqing Zhang, John Kounios, and Hualou Liang. Improved prediction of brain age using multimodal neuroimaging data. *Human Brain Mapping*, 2019.
- [9] Cathie Sudlow, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul Downey, Paul Elliott, Jane Green, Martin Landray, et al. UK Biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLOS Medicine*, 12(3):e1001779, 2015.
- [10] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph Attention Networks. *International Conference on Learning Representations*, 2018.

- [11] Felix Wu, Tianyi Zhang, Amauri Holanda de Souza Jr, Christopher Fifty, Tao Yu, and Kilian Q Weinberger. Simplifying graph convolutional networks. *arXiv preprint arXiv:1902.07153*, 2019.
- [12] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S Yu. A comprehensive survey on graph neural networks. *arXiv preprint arXiv:1901.00596*, 2019.