# Thesis Artifact

Multivariate Normative Models Using Variational Auto-Encoders: A Study on Covariate Embedding and Robustness to Site-Variance using Gen R Data

Remy Duijsens

*Delft University of Technology*

Delft, The Netherlands

r.w.j.p.duijsens@student.tudelft.nl

RESEARCH MOTIVATION

Representation learning concerns itself with learning meaningful patterns in data to perform predetermined tasks. These tasks vary widely and include regression, classification, generation, clustering, dimensionality reduction, and anomaly detection. Different domains and applications have embraced representation learning techniques to enhance the quality of their models. Often, a vast amount of complex data that has previously been underutilized is available. Deep learning is a specific type of representation learning that uses neural networks with many layers to learn representations of data. These deep models have shown to be an extremely powerful and successful tool in finding insightful patterns when data and features grow exponentially.

One important domain that has adopted representation learning in its analyses is predictive (clinical) psychiatry. A recent application in this field is the creation of normative models, an alternative to case-control studies. In normative modeling, it is essential to map individuals to a reference group. The idea of such a reference group comes from biological differences between individuals. Nuisance factors can explain these differences, introducing variability to individual data measurements. In the specific settings of biological brain age estimation, the data in question is related to brain imaging (e.g., raw MRI brain images or brain volumes and surface areas). Here, important nuisance factors are age and sex. This is because these factors provide information about the underlying biological differences in brain characteristics and can help improve the model. The reference groups can thus be partitioned using these factors. The factors intended to be embedded into models are referred to as covariates. The importance of effective covariate embedding techniques in deep learning models becomes apparent when shifting from univariate normative models (concerning a single variable) to a multivariate approach that can benefit significantly from the presence and coherence of a large number of (brain) features.

Examining current practices in deep learning modeling reveals that these nuisance factors are often not accounted for, even though their contribution can be vital to understanding data variability and improving model quality. Here, a lesson can be taken from normative modeling. Let us apply this insight to a popular toy and benchmark problem, such as handwritten digit recognition, using the MNIST dataset as an example. Here, the data consists of pixel values, which are the only features available. However, just like the brain images mentioned previously, this data originates from the real world and is influenced by nuisance factors that can contribute to the learning process. For example, consider the difference between left and right-handed people and the differences in handwriting between males and females.

The general question then becomes how these covariates can be embedded to be most valuable to our models. These covariates must not necessarily be of the same type and modality as the data, and even between the covariates themselves, there might be differences (e.g., continuous vs. categorical, tabular vs. image, bounded vs. unbounded). Thus, a natural extension of this question is how different data modalities can be embedded and what methods are most effective for specific types of covariates.

In contrast to these nuisance factors, noise factors are present in all real-world data. These factors can include random noise and variations due to many unknown parameters. In normative modeling of brain image data, noise factors include the measurement design, such as the use of different measurement hardware and device configurations (systematic noise), as well as biological factors not explained directly by the data and the selected covariates. Here, the question is similar to embedding the nuisance factors, yet the effects are negated: how can these noise factors be separated from the learned representation? Since there is no direct access to the factors generating the noise, proxies must be found that might explain parts of the noise or correlate with it. Covariates, and thus different covariate embedding techniques, could be used to separate the noise signals by finding the correlation between the hidden noise factors and the suspected correlated covariate (e.g., site when trying to separate noise from systematic site variance). Once separated, the noise signals could be disentangled into distinct latent factors or removed entirely from the learned representation. The model's robustness can be studied to evaluate its handling of noise.

The developed multivariate normative models should be examined to determine whether these new methods improve current univariate modeling efforts. Here, the comparison

is made regarding model characteristics (e.g., performance, robustness, complexity) and clinical performance, in which the biological age and the brain age gap are estimated. In addition, there is interest in metrics that can be used to compare the different modeling techniques. Primarily, the focus is on metrics that can be used to evaluate the performance of different covariate embedding techniques (e.g., loss-based functions) and metrics that evaluate the influence of noise on the model (e.g., robustness).

In conclusion, the necessity of this work is motivated based on 1) improving the understanding of the role of nuisance factors in deep learning and how covariate embedding techniques can improve model performance, 2) how noise factors can be addressed to improve robustness to noise, and 3) how multivariate normative modeling compares to existing univariate efforts—both in terms of the model characteristics and their performance in clinical applications. The corresponding hypotheses to be tested include whether various covariate embedding techniques improve model performance, whether addressing noise factors enhances robustness, and whether multivariate normative models offer better characteristics over univariate ones. An experimental platform will be built to validate these hypotheses through a series of experiments. The contribution of this work will provide new insights into these modeling techniques and advance the understanding of multivariate normative models. Furthermore, the findings regarding covariate embedding techniques and noise separation extend naturally to other domains and applications in deep learning and provide valuable insights into the field of AI.