GenNet Framework: Interpretable Neural Networks for Predicting Phenotype from Genotype

Arno van Hilten¹, Steven A. Kushner², Hieab H.H. Adams^{1,3}, Wiro J. Niessen^{1,4}, Gennady V. Roshchupkin^{1,5}

¹ Department of Radiology and Nuclear Medicine, Erasmus MC, Rotterdam, The Netherlands; ² Department of Psychiatry, Erasmus MC, Rotterdam, The Netherlands; ³ Department of Clinical Genetics, Erasmus MC, Rotterdam, Netherlands. ; ⁴ Department of Imaging Physics, Delft University of Technology, Delft, The Netherlands; ⁵ Department of Epidemiology, Erasmus MC, Rotterdam, The Netherlands

Introduction

Neural networks are currently the state of the art in many areas of scientific research and applications but are seldom applied in population genomics and **predictive genetics** due to the computational burden and lack of **interpretability**. Here, we propose a novel **open-source** deep learning framework, **GenNet**, for interpretable predictions of phenotypes from genotypes. In this framework, public prior **biological knowledge** (e.g. DNA and pathway annotations) is used to define a sparse, memory-efficient networks.

Results

Erasmus Medical Center Rotterdam

BIGE

The framework is validated in linear and non-linear simulations (see figure 2) and applied to the Rotterdam study (exonic variants) the UK Biobank and the Schizophrenia WES Data (see table 1). Identifying commonly associated genes for identifying hair and eye color such as OCA2 and HERC2, validating the interpretability of the network.

Trait	Dataset (type)	Subjects &	Phenotype	— Heritability	AUC	AUC	GenNet:
		Class I	Class II		Lasso	GenNet	top 3 genes

Methods

In the framework, **prior knowledge** is used to create groups of connected nodes to reduce the number of learnable parameters in comparison to a fully connected neural network. For example in the first layer, where **millions** of single nucleotide polymorphisms (SNPs) inputs are only connected to their corresponding genes, creating meaningful and interpretable connections while significantly **reducing** the total number of **parameters**.



Eye color	Rotterdam genotype array	4041 Blue	Other	0.80-0.98	0.68	0.75	HERC2, OCA2, LAMC1
Hair color	UK Biobank	1648	1656	0.70-0.97	0.78	0.83	MC1R*. OCA2. TC2N
	WES data	Blond	Red		0110		, , , , , , , , , , , , , , , , , , ,
	UK Biobank	1672	1664	0 70 0 07	0 70	0 00	MC1D* OCAD 7CCUCA
	WES data	Dark brown	Red	0.70-0.97	0.79	U.00	$MCIR^{+}$, $UCA2$, $ZCCHC4$
	UK Biobank	4352	4343	0.70-0.97	0.64	0.75	OCA2, TC2N, EXOC2
	WES data	Blond	Dark brown				
Mala baldnaca	UK Biobank	3454	3454	0.60-0.70	0.57	0.57	NGEF, NKRD18B,
wale baldness	WES data	No balding	Severe balding				SYNJ2
Dinalar	UK Biobank	343	347	0.73-0.93	0.59	0.60	LINC00266-1, CSMD1,
Dibolat	WES data	Cases	Controls				TCERG1L
Cohizophronia	Sweden	4969	6245	0.80-0.85	0.65	0.74	7NIE772 DCNIT DVCE
Schizophrenia	WES data	Cases	Controls				ZINFIIS, PCINI, DISF

Table 1: Summary of the experiments and results in this study for the simplest network in our framework that contains the input SNPs, the gene layer and the output layer. *MC1R was not present in gene annotations but identified by linkage disequilibrium.



Figure 1: Overview of neural networks in the **GenNet framework**. Different types of prior biological knowledge can be used to **define** the **connections**. In the shown network gene annotations were used to define the connections to the first layer and pathway annotations to group the genes from the first layer. Aside from gene and pathway annotations, our framework provides layers built from exon annotations, chromosome annotations and cell and tissue type expressions.

GenNet is an open-source framework written in **Keras and Tensorflow**, code and tutorials are available on https://github. com/arnovanhilten/GenNet/.

Conclusion

- We developed a framework with memory-efficient and interpretable neural networks by using publicly available **biological** knowledge.
- These networks obtain good predictive performance even while only using **exomic variants**.
- We anticipate this approach as having the potential for uncovering **novel insights** into the genetic architecture of complex diseases

Figure 2: A) Non-linear simulation showing the basic principle of the network, thickness of the connections represents the learned weight (causal in red). Run online https://tinyurl.com/y8hh8rul. B) Simulations with synthetic data showing the performance of GenNet expressed in the area under the curve for increasing levels of heritability and training set size (C). In black the theoretical maximum of the AUC versus heritability. **D)** Manhattan plot of the importance of the genes according to the network for distinguishing between schizophrenia cases and controls. E) This manhattan plot is a cross section of the trained network between the gene layer and the outcome.

