# Modelling human complex traits with regression and neural-network based methods

**Marton Kelemen**

Wellcome Sanger Institute
University of Cambridge

This dissertation is submitted for the degree of
*Doctor of Philosophy*

Darwin College                                   November 2020

# Declaration

This thesis is the result of my own work and includes nothing which is the outcome of work done in collaboration except as declared in the Preface and specified in the text. It is not substantially the same as any that I have submitted, or, is being concurrently submitted for a degree or diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text. I further state that no substantial part of my thesis has already been submitted, or, is being concurrently submitted for any such degree, diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text. It does not exceed the prescribed word limit for the relevant Degree Committee. This dissertation contains fewer than the limit of 60,000 words set by the Biology Degree Committee.

<div align="right">

Marton Kelemen

November 2020

</div>

# Modelling human complex traits with regression and neural-network based methods

## Marton Kelemen

Identifying how epistasis, non-linear genetic effects, contribute to phenotypic variance in humans has been an enduring challenge. So far neither the computational resources that could accommodate higher-order interactions at scale nor the large-scale population cohorts with adequate statistical power were available up until recently. With the advent of graphics processing unit computing farms and neural-network based methods, together with large biobank-scale data sets, such as the UK Biobank which offers a sample size of ~500K, this has been changing. These developments offer opportunities for the development of novel approaches that could provide insights into the genetic underpinnings of complex disease risk and trait variation.

After reviewing the necessary background material, this work consists of three research chapters. The organising theme of these is the building of genotype-phenotype maps, which grow from the simple additive, through the two-way interactions, up to higher-order interactions in the last chapter.

I begin by covering the common quality control steps and basic additive association analyses I carried out that explored the information boundaries of my data which serves as the foundation for the rest of my work. I managed to recover primary association signals described in the literature for my cohorts confirming the validity of my data processing steps. I also describe a novel method that exploits shared genetic effects to improve risk prediction for related traits. Relative to baselines, this improved squared correlations between observed and predicted sub-phenotypes by ~25% and ~19% for ulcerative colitis and Crohn's disease, respectively.

Building on the previously prepared data sets, I searched for two-way interactions using standard statistical methods belonging to the regression framework. In the UK Biobank cohort I pursued a hypothesis-free approach to consider interactions both within and between the genomic domains of SNP, transcription and protein derived predictors. For the much smaller inflammatory bowel disease studies, I followed a hypothesis driven strategy to reduce search space which only considered haplotype-specific interactions between biologically plausible loci to increase power. I found that the results from both of these approaches were consistent with the null hypothesis of no significant contribution to phenotypic variance from non-linear genetic effects.

Parallel to my search for epistasis using regression based models, I also considered the neural-network framework to find indirect evidence for non-linear effects contributing to

phenotypic variance. I confirmed via a large-scale simulation study the potential of neural-networks to be able to identify interactions at a higher accuracy than standard regression based methods. In the real datasets, I searched for individual epistatic interactions using both experimental approaches from the literature, together with methods that I developed for this purpose. However, I was unable to find convincing evidence for statistical interactions contributing to complex trait variance.

In summary, I found that despite the large cohorts I had access to and the modern non-linear methods I deployed, evidence for non-linear genetic effects contributing to complex human trait variance remained elusive.

# Acknowledgements

This thesis is dedicated to my mother, Judith Dimitrova, and my father, Lajos Kelemen.

# Table of contents

# List of figures

# List of tables

# Nomenclature

**Acronyms / Abbreviations**

**AUC**    Area under the curve

**BMI**    Body mass index

**CD**    Crohn's disease

**CNN**    Convolutional neural-network

**FDR**    False discovery rate

**FIS**    Fluid intelligence score

**FNN**    Fully-connected neural-network

**GPU**    Graphics processing unit

**GWAS**    Genome-wide association study

**HLA**    Human leukocyte antigen

**HMM**    Hidden Markov model

**HWE**    Hardy–Weinberg equilibrium

**IBD**    Inflammatory bowel disease

**KRR**    Kernel-ridge regression

**LD**    Linkage disequilibrium

**LMM**    Linear mixed-effects model

**MAF**    Minor allele frequency

**NN** Neural-network

**OLS** Ordinary least squares

**PRS** Polygenic risk score

**QC** Quality control

**REML** Restricted maximum likelihood

**RKHS** Reproducing kernel Hilbert space

**ROC** Receiver operating characteristic

**RR** Ridge regression

**SNP** Single-nucleotide polymorphism

**T1D** Type 1 diabetes

**TF** Transcription factor

**TWAS** Transcriptome-wide association study

**UC** Ulcerative colitis

**UKBB** UK Biobank

**WGS** Whole-genome sequencing