HIGH DIMENSIONAL STATISTICS LECTURE I: MOTIVATION FOR THE COURSE

Msc programme in Mathematics and Statistics University of Helsinki Matti Pirinen

"HIGH DIMENSIONAL", WHAT IS IT AND WHY BOTHER?

Each row is a gene (~ 10^4) and each column a cancer patient (~ 10^3)

Red/blue colors represent gene expression levels

I.Are different subtypes of cancer different in gene expression?

2. Can we tell from gene expression which subtype a patient has?

3. Could we choose treatment based on the gene expression levels?



From: "Integrated genomic and molecular characterization of cervical cancer." Nature 543.

"HIGH DIMENSIONAL", WHAT IS IT AND WHY BOTHER?

UK Biobank releases brain Images of 100,000s participants

I. How brain activity changes between rest and tasks?

2. Can we tell from measured brain activity what is the context of the individual ?

3. What are the statistical associations between brain activity patterns and 1000s of measured congnitive, behavioral, lifestyle or genetic variables?



Nature Neuroscience 2015 19, p.1523–1536

"HIGH DIMENSIONAL", WHAT IS IT AND WHY BOTHER?

- Examples of ongoing data explosion
 - Life sciences with high-throughput technologies, e.g. genomics, metabolomics
 - Physics and engineering, e.g. CERN, astronomy, robotics, sensors
 - Humanities and digitalization, e.g. libraries of historical texts
 - Internet, e.g. images, videos, sound, text, social media
- Working definition: High Dimensional (HD) data set X has a lot of observations (n x p is large)
 - *n* units (or samples) as rows of X (e.g. individuals $n \sim 10^5$)
 - *p* variables (or features) as columns of X (e.g. genetic variants $p \sim 10^6$)
 - often HD means that "p >> n is large", but more general definition for us is " $n \ge p$ is large"
- Unprecedented potential for new knowledge, but a need for suitable methods

STATISTICS AND MACHINE LEARNING

- Methods we will consider have (some of) following properties
 - Role of the variables is interpretable in the fitted model
 - Simple and complete description in terms of a probability model
 - Conceptually straightforward quantification of uncertainty of parameters and predictions (although not always easy to compute in HDs)
- Our methods are instances of **statistical learning** subset of machine learning (ML) methods
 - ML has also powerful methods for prediction that are more of "black boxes" and are defined by algorithms rather than by probability models
 - Deep learning, random forests etc.
- In modern data science, we should know all types of learning methods
 - In this course, we will focus on statistical learning

EXCELLENT BOOKS AVAILABLE ONLINE

Springer Texts in Statistics

Gareth James Daniela Witten Trevor Hastie Robert Tibshirani

An Introduction to Statistical Learning

with Applications in R

- Shorthand: ISL
- Very easy to read
- Thorough examples in R
- Little maths
- Excellent book to get intuition behind the concepts and models
- Video lectures of chapters available
- Most recent edition June 2023

https://www.statlearning.com/

Description Springer

- Shorthand: ESL
- Comprehensive collection of methods
- Mathematical descriptions included



https://hastie.su.domains/ElemStatLearn/

CONTENTS

- Weeks I-2: Large-scale inference, i.e., what are the statistical ideas and measures used when we carry out thousands of tests/comparisons simultaneously
- Weeks 3-5: Regression with a large numbers of predictors, variable selection, bias-variance trade-off, bootstrap
- Weeks 6-7: Dimension reduction
- Week 7: Summary
- Exam

PASSING THE COURSE

- 6 sets of exercises to be returned through Moodle area "High Dimensional Statistics 2024"
 - Half of the exercise points needed to pass (i.e. 30 out of 60)
 - Computer exercises to be done with R, preferably R Markdown
 - Tutored exercise class with a teacher on Mondays 12.15-14.00 B222
 - Exercise session videos available after each week's deadline (Tue 10.00)
- Exam 3hrs is to be done remotely with a computer (vote for time)
- Lectures on Tuesdays & Thursdays 10.15 -12.00 in B222
- Course material on the course home page (link in Moodle)

EXAMPLE FROM GENOMICS

- Which genetic variants are **associated** with cardiovascular disease (CVD), #I cause of death in the Western world? (weeks I-2)
- Which genetic variants are **causal** for CVD? (weeks 3-5)
- How can we best **predict** genetic risk for CVD? (weeks 3-5)
- How can we **visualize** and extract the **main structure** of very highdimensional genetic data in just a few dimensions? (weeks 6-7)

GENOME-DISEASE ASSOCIATION STUDY

- Collect 10,000s of cases (individuals with the disease) and controls (individuals from the general population who do not have the disease)
- Genotype everyone in 1,000,000s of genomic positions
- Do a statistical test at each position to see whether genotype distributions are different between cases and controls



WHICH VARIANTS ARE INTERESTING?



- Each variant is tested for a statistical difference between cases and controls
- Millions of tests, how to do inference?
- In Fig., what kind of threshold measures are
 - P-value = 5e-8
 - FDR = 5%
- SNP = name of the variant EAF = effect allele frequency OR = odds ratio

Nelson et al. Nat Genetics 2017

WHICH VARIANTS ARE CAUSAL (AND NOT JUST ASSOCIATED)?



- Variants physically near each other are highly correlated and show similar effect sizes / P-values
- Which one(s) of them is/are truly driving the signal and which are just passengers?
- We need to analyze them jointly, which becomes a HD regression problem

Benner et al.AJHG 2017

WHICH VARIANTS ARE CAUSAL (NOT JUST ASSOCIATED)?



Benner et al.AJHG 2017

19q13/APOE association with LDL cholesterol

PREDICTION MODEL

Accurate Genomic Prediction of Human Height

Louis Lello,* Steven G. Avery,* Laurent Tellier,*^{,†,‡} Ana I. Vazquez,[§] Gustavo de los Campos,^{§,**} and Stephen D. H. Hsu^{*,†,1}

Genetics, Vol. 210, 477–497 October 2018

- Starts with 650,000 genetic variants and 420,000 individuals with height measurements
- Uses the LASSO method for building the predictive model (same method that we will look at next weeks)

IDENTIFYING RELEVANT VARIANTS



• About 20,000 variants are identified by LASSO and each with its effect size will be used in predicting the height of a new test individual

TESTING THE PREDICTOR



Predictor has correlation of 0.61 with actual height, that is, it explains about 37% (=0.61²) of variation of height in **the test sample**.

Most individuals are within 4cm of their predicted height.



PC1

- From genotypes (10³ x 10⁶) to pairwise covariances (10³ x 10³) to first 2 principal components (2 x 10³)
- Reduction is of order 10⁵ and the main structure is not only preserved but has also become more easily visible