

MAS202M Final Project

Kári Hlynsson

University of Iceland

April 2024



Overview

- Part A: Data cleaning / clustering (30%)
 - Removal of duplicated samples
 - Treatment of missing values (NAs)
 - Hierarchical clustering, K-means and PCA
- Part B: Prediction / inference (70%)
 - EHP1: LASSO and Random Forest
 - temp: SVMs and KNN



Part A



Data cleaning

- Raw data is 1772×8 , reduced to 1672×5
- The first column, `index`, was removed
- Covariates `col3` and `col6` were identical except at points where the other was `NA`. The missing observations in `col3` were derived from `col6`; the latter was removed
- 97 observations were found to be duplicates and were removed from the data set
- 3 observations were duplicates but contained one `NA` field, and were not initially detected. These were also removed
- Following removal of duplicated, the `id` column held no relevant information and was subsequently removed

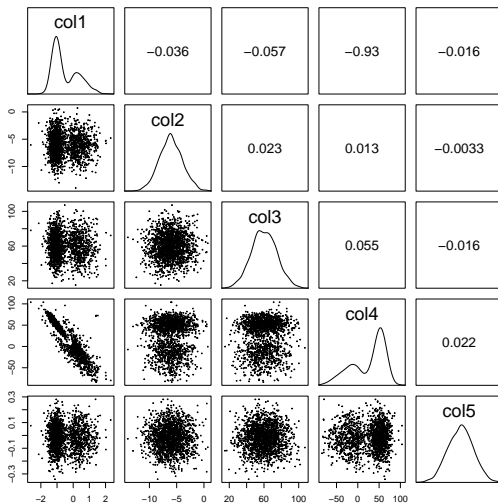


Data cleaning

- The raw data contained 20 missing values
- Means of respective covariates was used to impute these values



Clustering

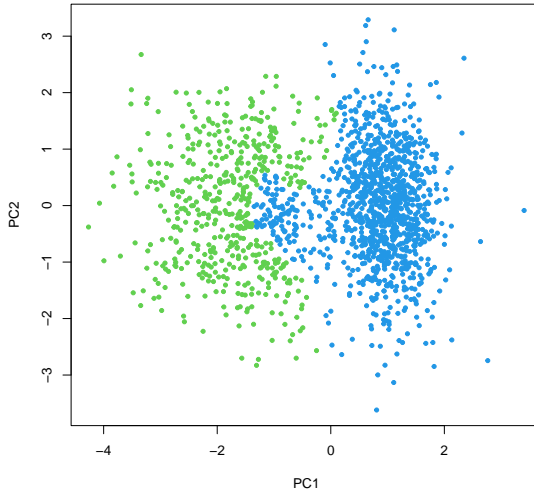


Clustering

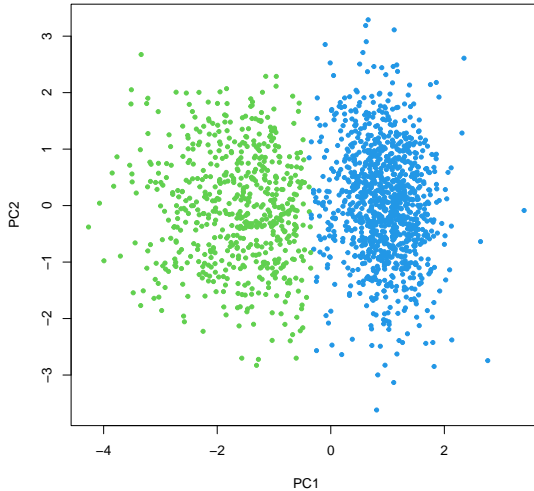
- To cluster the data, hierarchical clustering and K-means were considered
 - Each of the two methods was run on *raw* and *scaled* data, as well as on the first two *principal components* of the data
 - Clustering accuracy was mainly assessed through visual confirmation
- Optimal clustering results were obtained for K-means clustering on PCA data



Clustering – Hierarchical



Clustering – K-means



Part B



Quantitative variable

- We will consider **EHPB1**
- LASSO and random forest regression



LASSO

- High dimensional setting, $n \ll p$, motivates use of methods that *reduce number of covariates* used in modelling
- LASSO was chosen for this task in the linear setting. Other potential candidates include ridge, PCR, PLS, AIC step procedures, etc.
- Assumes a *linear relationship* between response and covariates which may be overly stringent
- Optimise over cost hyperparameter, λ :

$$\beta_{\text{LASSO}} = \underset{\beta}{\operatorname{argmin}} \left[\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \left| \sum_{j=1}^p \beta_j \right| \right]$$



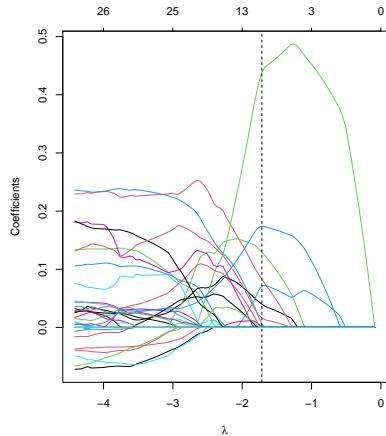
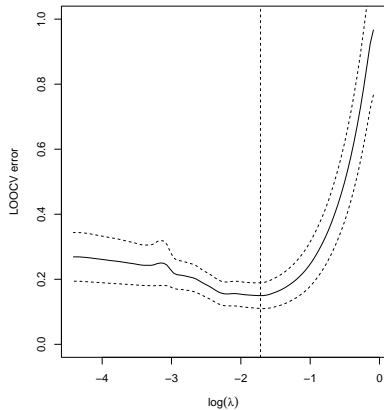
LASSO

- Optimal value of λ determined using 75% train-test split using LOOCV. This gave $\lambda \approx 0.1797$
- Test set MSE = 0.1080
- 8 predictors were non-zero, specifying the model

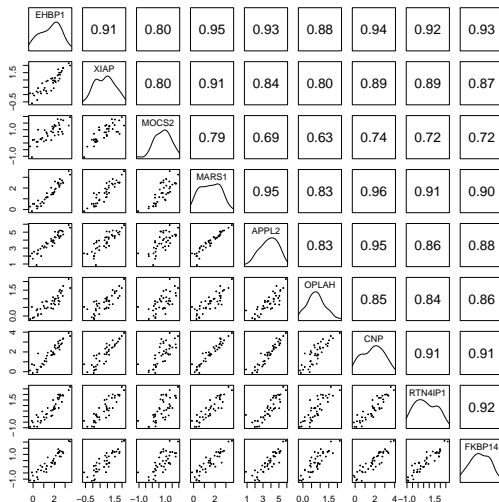
$$\begin{aligned} \text{EHBP1} = & 0.5504 + 0.0403 \times \text{XIAP} + 0.0826 \times \text{MOCS2} \\ & + 0.4383 \times \text{MARS1} + 0.006 \times \text{APPL2} \\ & + 0.1296 \times \text{OPLAH} + 0.071 \times \text{CNP} \\ & + 0.0164 \times \text{RTN4IP1} + 0.1734 \times \text{FKBP14} \end{aligned}$$



LASSO



LASSO

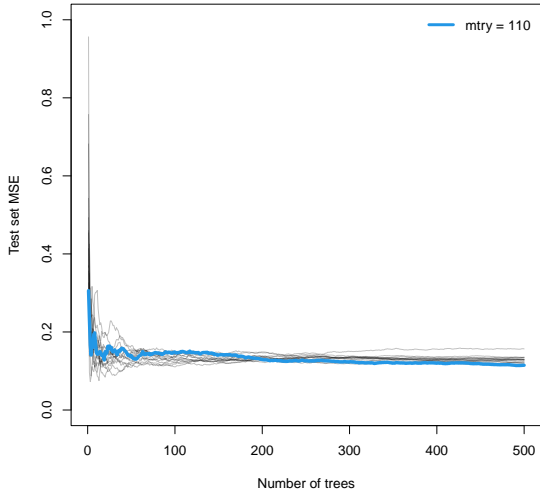


Random forest

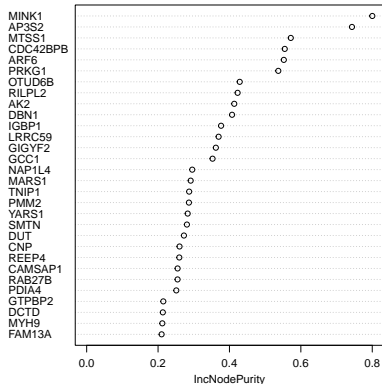
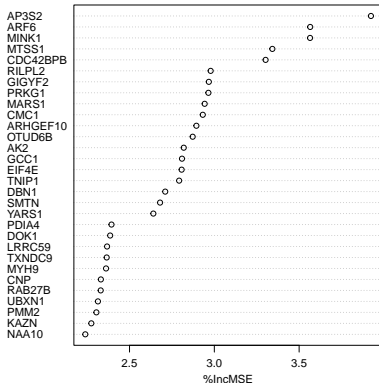
- Random forest regressor was tuned for optimal `mtry` with `ntree` = 500, yielding `mtry` = 110
- Works well in the *high-dimensional setting* $n \ll p$ due to restricting decision trees to a *subset of predictors*
- More *candidate causal predictors* uncovered – favourable in "breadth-first" scanning, e.g. for biomarkers
- Test set MSE of 0.1147, marginally unfavourable to LASSO.



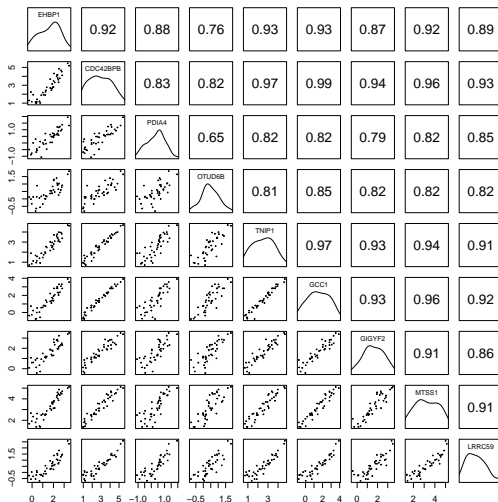
Random forest



Random forest



Random forest

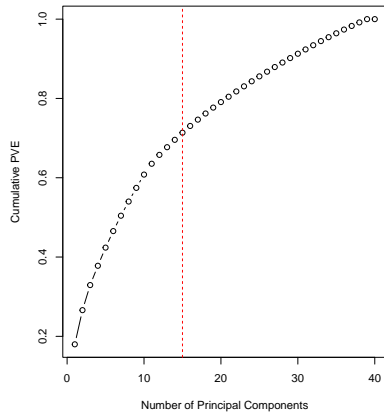
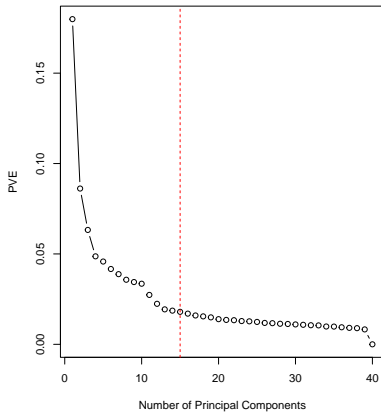


Categorical variable

- We will study the **temp** variable
- SVM and KNN
- Fitted on data transformed using PCA to reduce computational complexity (15 components, 71.4% var.)



Categorical variable



SVM

- Optimised over three kernels; linear (SVC), polynomial and radial and respective parameters
 - Linear: `cost`
 - Polynomial: `cost`, `degree`
 - Radial: `cost`, `gamma`
- Linear (`cost` = 0.0785) and radial kernels gave (`cost` = 54.5559, `gamma` = 0.0034) best performance with 70% accuracy
- Polynomial kernel (`cost` = 6.1585, `degree` = 1) lagged behind with 65% accuracy



KNN

- Strong out-of-the-box classifier, *nonparametric*
- Tends to suffer in high-dimensional settings
- Optimised for the number of nearest observations to inspect when classifying, K
- $K = 9$ gave the optimal accuracy of 75%



KNN

