

# Assignment 2

## Assignment 2

### Biomedical Data Science

**Due on Thursday 18th March 2020, 5:00pm**

The assignment is marked out of 100 points, and will contribute to 30% of your final mark. Please knit this document in PDF format and submit using the gradescope link on Learn. If you can't knit to PDF directly, knit it to word and you should be able to either convert to PDF or print it and scan to PDF using a scanning app on your phone. If you have any code that doesn't run you won't be able to knit the document so comment it as you might still get some grades for partial code. Clear and reusable code will be rewarded so pay attention to indentation, choice of variable identifiers, comments, error checking, etc. An initial code chunk is provided after each subquestion but create as many chunks as you feel is necessary to make a clear report. Add plain text explanations in between the chunks as and when required and any comments necessary within code chunks to make it easier to follow your code/reasoning.

### Problem 1 (27 points)

File `wdbc2.csv` (available from the accompanying zip folder on Learn) refers to a study of breast cancer where the outcome of interest is the type of the tumour (benign or malignant, recorded in column "diagnosis"). The study collected 30 imaging biomarkers on 569 patients.

#### Problem 1.a (7 points)

Using package `caret`, create a data partition so that the training set contains 70% of the observations (set the random seed to 984065 beforehand). Fit both a ridge regression model and a lasso model which uses cross-validation on the training set to diagnose the type of tumour from the 30 biomarkers. Then use a plot to help identify the penalty parameter  $\lambda$  that maximizes the AUC. Note: There is no need to use the `prepare.glmnet()` function from lab 4, using `as.matrix()` with the required columns is sufficient.

```
# Enter code here.
```

#### Problem 1.b (2 points)

Create a data table that for each value of 'lambda.min' and 'lambda.1se' for each model fitted in problem 1.a reports: \* the corresponding AUC, \* the corresponding model size. Use 3 significant digits for floating point values and comment on these results. Hint: The AUC values are stored in the field called 'cvm'.

```
# Enter code here.
```

#### Problem 1.c (7 points)

Perform both backward (we'll later refer to this as model B) and forward (model S) stepwise selection on the same training set derived in problem 1.a. Report the variables selected and their standardized regression coefficients in decreasing order of the absolute value of their standardized regression coefficient. Discuss the results and how the different variables entering or leaving the model influenced the final result.

```
# Enter code here.
```

### Problem 1.d (3 points)

Compare the goodness of fit of model B and model S in an appropriate way.

```
# Enter code here.
```

### Problem 1.e (2 points)

Compute the training AUC for model B and model S.

```
# Enter code here.
```

### Problem 1.f (6 points)

Use the four models to predict the outcome for the observations in the test set (use the lambda at 1 standard error for the penalised models). Plot the ROC curves of these models (on the sameplot, using different colours) and report their test AUCs. Compare the training AUCs obtained in problems 1.b and 1.e with the test AUCs and discuss the fit of the different models.

```
# Enter code here.
```

## Problem 2 (40 points)

File GDM.raw.txt (available from the accompanying zip folder on Learn) contains 176 SNPs to be studied for association with incidence of gestational diabetes (a form of diabetes that is specific to pregnant women). SNP names are given in the form “rs1234\_X” where “rs1234” is the official identifier (rsID), and “X” (one of A, C, G, T) is the reference allele.

### Problem 2.a (3 points)

Read file GDM.raw.txt into a data table named gdm.dt. Impute missing values in gdm.dt according to SNP-wise median allele count.

```
# Enter code here.
```

### Problem 2.b (8 points)

Write function `univ.glm.test <- function(x, y, order = FALSE)` where `x` is a data table of SNPs, `y` is a binary outcome vector, and `order` is a boolean. The function should fit a logistic regression model for each SNP in `x`, and return a data table containing SNP names, regression coefficients, odds ratios, standard errors and p-values. If `order` is set to `TRUE`, the output data table should be ordered by increasing p-value.

```
# Enter code here.
```

### Problem 2.c (5 points)

Using function `univ.glm.test()`, run an association study for all the SNPs in `gdm.dt` against having gestational diabetes (column “pheno”). For the SNP that is most strongly associated to increased risk of gestational diabetes and the one with most significant protective effect, report the summary statistics from the GWAS as well as the 95% and 99% confidence intervals on the odds ratio.

```
# Enter code here.
```

### Problem 2.d (4points)

Merge your GWAS results with the table of gene names provided in file GDM.annot.txt (available from the accompanying zip folder on Learn). For SNPs that have p-value  $< 10^{-4}$  (hit SNPs) report SNP name, effect allele, chromosome number and corresponding gene name. Separately, report for each ‘hit SNP’ the names of the genes that are within a 1Mb window from the SNP position on the chromosome. Note: That’s genes that fall within  $\pm 1,000,000$  positions using the ‘pos’ column in the dataset.

```
# Enter code here.
```

### Problem 2.e (8 points)

Build a weighted genetic risk score that includes all SNPs with p-value  $< 10^{-4}$ , a score with all SNPs with p-value  $< 10^{-3}$ , and a score that only includes SNPs on the FTO gene (hint: ensure that the ordering of SNPs is respected). Add the three scores as columns to the `gdm.dt` data table. Fit the three scores in separate logistic regression models to test their association with gestational diabetes, and for each report odds ratio, 95% confidence interval and p-value.

```
# Enter code here.
```

### Problem 2.f (4 points)

File `GDM.test.txt` (available from the accompanying zip folder on Learn) contains genotypes of another 40 pregnant women with and without gestational diabetes (assume that the reference allele is the same one that was specified in file `GDM.raw.txt`). Read the file into variable `gdm.test`. For the set of patients in `gdm.test`, compute the three genetic risk scores as defined in problem 2.e using the same set of SNPs and corresponding weights. Add the three scores as columns to `gdm.test` (hint: use the same columnnames as before).

```
# Enter code here.
```

### Problem 2.g (4 points)

Use the logistic regression models fitted in problem 2.e to predict the outcome of patients in `gdm.test`. Compute the test log-likelihood for the predicted probabilities from the three genetic risk score models.

```
# Enter code here.
```

### Problem 2.h (4points)

File `GDM.study2.txt` (available from the accompanying zip folder on Learn) contains the summary statistics from a different study on the same set of SNPs. Perform a meta-analysis with the results obtained in problem 2.c (hint: remember that the effect alleles should correspond) and produce a summary of the meta-analysis results for the set of SNPs with meta-analysis p-value  $< 10^{-4}$  sorted by increasing p-value.

```
# Enter code here.
```

## Problem 3 (33 points)

File `nki.csv` (available from the accompanying zip folder on Learn) contains data for 144 breast cancer patients. The dataset contains a binary outcome variable (“Event”, indicating the insurgence of further complications after operation), covariates describing the tumour and the age of the patient, and gene expressions for 70 genes found to be prognostic of survival.

### Problem 3.a (6 points)

Compute the matrix of correlations between the gene expression variables, and display it so that a block structure is highlighted. Discuss what you observe. Write some code to identify the unique pairs of (distinct) variables that have correlation coefficient greater than 0.80 in absolute value and report their correlation coefficients.

```
# Enter code here.
```

### Problem 3.b (8 points)

Run PCA (only over the columns containing gene expressions), in order to derive a patient-wise summary of all gene expressions (dimensionality reduction). Decide which components to keep and justify your decision. Test

if those principal components are associated with the outcome in unadjusted logistic regression models and in models adjusted for age, estrogen receptor and grade. Justify the difference in results between unadjusted and adjusted models.

```
# Enter code here.
```

### **Problem 3.c (8 points)**

Use plots to compare with the correlation structure observed in problem 2.a and to examine how well the dataset may explain your outcome. Discuss your findings and suggest any further steps if needed.

```
# Enter code here.
```

### **Problem 3.d (11 points)**

Based on the models we examined in the labs, fit an appropriate model with the aim to provide the most accurate prognosis you can for patients. Discuss and justify your decisions.

```
# Enter code here.
```