

# Modern Data Mining, HW 3

Group Member 1

Group Member 2

Group Member 3

11:59 pm, 03/17, 2024

## Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>PartI: Model Building</b>  | <b>2</b> |
| 1.1      | Objectives . . . . .  | 2        |
| 1.2      | Review materials . . . . .  | 2        |
| 1.3      | Case study: COVID19 . . . . .                                       | 2        |
| <b>2</b> | <b>Part II: Logistic Regression</b>                                 | <b>2</b> |
| 2.1      | Objectives . . . . .  | 3        |
| 2.2      | R Markdown / Knitr tips . . . . .                                   | 3        |
| 2.3      | Review . . . . .  | 4        |
| 2.4      | Framingham heart disease study . . . . .                            | 4        |
| 2.4.1    | Identify risk factors . . . . .                                     | 4        |
| 2.4.1.1  | Understand the likelihood function . . . . .                        | 4        |
| 2.4.1.2  | Identify important risk factors for <b>Heart.Disease.</b> . . . . . | 4        |
| 2.4.1.3  | Model building . . . . .  | 5        |
| 2.4.2    | Classification analysis . . . . .                                   | 5        |
| 2.4.2.1  | ROC/FDR . . . . .   | 5        |
| 2.4.2.2  | Cost function/ Bayes Rule . . . . .                                 | 5        |

# 1 PartI: Model Building

Multiple regression is one of the most popular methods used in statistics as well as in machine learning. We use linear models as a working model for its simplicity and interpretability. It is important that we use domain knowledge as much as we could to determine the form of the response as well as the function format for the factors. Then, when we have many possible features to be included in the working model it is inevitable that we need to choose a best possible model with a sensible criterion. Regularizations such as LASSO are introduced. Be aware that if a model selection is done formally or informally, the inferences obtained with the final `lm()` fit may not be valid. Some adjustment will be needed. This last step is beyond the scope of this class. Check the research line that Linda and collaborators have been working on.

The main job in this part is a rather involved case study about devastating covid19 pandemic. Please read through the case study first. This project is for sure a great one listed in your CV.

For covid case study, the major time and effort would be needed in EDA portion.

## 1.1 Objectives

- Model building process
- Methods
  - Model selection
    - \* LASSO (L1 penalty)
    - \* A quick backward elimination
- Understand the criteria
  - Testing Errors
  - K fold Cross Validation
  - LASSO
- Packages
  - `lm()`, `Anova`
  - `regsubsets()`
  - `glmnet()` & `cv.glmnet()`

## 1.2 Review materials

- Study lecture: Regularization
- Study lecture: Multiple regression

Review the code and concepts covered during lectures: multiple regression, penalized regression through elastic net (only LSASSO).

**Important Notice:** The focus of this part is Covid case study.

## 1.3 Case study: COVID19

See a separate file `covid_case_study_2024.Rmd` for details.

- Start the EDA as earlier as possible.
- Please check previous midterms where we used the same dataset.

# 2 Part II: Logistic Regression

Logistic regression is used for modeling categorical response variables. The simplest scenario is how to identify risk factors of heart disease? In this case the response takes a possible value of YES or NO. Logit

link function is used to connect the probability of one being a heart disease with other potential risk factors such as **blood pressure**, **cholesterol level**, **weight**. Maximum likelihood function is used to estimate unknown parameters. Inference is made based on the properties of MLE. We use AIC to help nailing down a useful final model. Predictions in categorical response case is also termed as **Classification** problems. One immediately application of logistic regression is to provide a simple yet powerful classification boundaries. Various metrics/criteria are proposed to evaluate the quality of a classification rule such as **False Positive**, **FDR** or **Mis-Classification Errors**.

LASSO with logistic regression is a powerful tool to get dimension reduction. We will not use it here in this work.

## 2.1 Objectives

- Understand the model
  - logit function
    - \* interpretation
  - Likelihood function
- Methods
  - Maximum likelihood estimators
    - \* Z-intervals/tests
    - \* Chi-squared likelihood ratio tests
- Metrics/criteria
  - Sensitivity/False Positive
  - True Positive Prediction/FDR
  - Misclassification Error/Weighted MCE
  - Residual deviance
  - Training/Testing errors
- R functions/Packages
  - `glm()`, `Anova`
  - `pROC`
- Data needed
  - `Framingham.dat`

## 2.2 R Markdown / Knitr tips

You should think of this R Markdown file as generating a polished report, one that you would be happy to show other people (or your boss). There shouldn't be any extraneous output; all graphs and code run should clearly have a reason to be run. That means that any output in the final file should have explanations.

A few tips:

- Keep each chunk to only output one thing! In R, if you're not doing an assignment (with the `<-` operator), it's probably going to print something.
- If you don't want to print the R code you wrote (but want to run it, and want to show the results), use a chunk declaration like this: `{r, echo=F}`. Notice this is set as a global option.
- If you don't want to show the results of the R code or the original code, use a chunk declaration like: `{r, include=F}`
- If you don't want to show the results, but show the original code, use a chunk declaration like: `{r, results='hide'}`.
- If you don't want to run the R code in a chunk at all use `{r, eval = F}`.
- We show a few examples of these options in the below example code.
- For more details about these R Markdown options, see the [documentation](#).
- Delete the instructions and this R Markdown section, since they're not part of your overall report.

## 2.3 Review

Review the code and concepts covered in

- Module Logistic Regressions/Classification

## 2.4 Framingham heart disease study

We will continue to use the Framingham Data (`Framingham.dat`) so that you are already familiar with the data and the variables. All the results are obtained through training data.

Liz is a patient with the following readings: `AGE=50`, `GENDER=FEMALE`, `SBP=110`, `DBP=80`, `CHOL=180`, `FRW=105`, `CIG=0`. We would be interested to predict Liz's outcome in heart disease.

To keep our answers consistent, use a subset of the data, and exclude anyone with a missing entry. For your convenience, we've loaded it here together with a brief summary about the data.

We note that this dataset contains 311 people diagnosed with heart disease and 1095 without heart disease.

```
0    1
1095 311
```

After a quick cleaning up here is a summary about the data:

Lastly we would like to show five observations randomly chosen.

|     | HD | AGE | SEX    | SBP | DBP | CHOL | FRW | CIG |
|-----|----|-----|--------|-----|-----|------|-----|-----|
| 643 | 1  | 61  | MALE   | 140 | 68  | 248  | 104 | 20  |
| 11  | 0  | 45  | MALE   | 110 | 88  | 183  | 90  | 0   |
| 576 | 1  | 58  | MALE   | 150 | 95  | 296  | 100 | 15  |
| 560 | 1  | 59  | MALE   | 260 | 130 | 246  | 111 | 20  |
| 702 | 0  | 45  | FEMALE | 122 | 74  | 178  | 88  | 5   |

### 2.4.1 Identify risk factors

**2.4.1.1 Understand the likelihood function** Conceptual questions to understand the building blocks of logistic regression. All the codes in this part should be hidden. We will use a small subset to run a logistic regression of HD vs. SBP.

- Take a random subsample of size 5 from `hd_data_f` which only includes HD and SBP. Also set `set.seed(471)`. List the five observations neatly below. No code should be shown here.
- Write down the likelihood function using the five observations above.
- Find the MLE based on this subset using `glm()`. Report the estimated logit function of SBP and the probability of HD=1. Briefly explain how the MLE are obtained based on ii. above.
- Evaluate the probability of Liz having heart disease.

**2.4.1.2 Identify important risk factors for Heart.Disease.** We focus on understanding the elements of basic inference method in this part. Let us start a fit with just one factor, SBP, and call it `fit1`. We then add one variable to this at a time from among the rest of the variables. For example

- Which single variable would be the most important to add? Add it to your model, and call the new fit `fit2`.

We will pick up the variable either with highest  $|z|$  value, or smallest  $p$  value. Report the summary of your `fit2` Note: One way to keep your output neat, we will suggest you using `xtable`. And here is the summary report looks like.

- Is the residual deviance of `fit2` always smaller than that of `fit1`? Why or why not?

|             | Estimate | Std. Error | z value | Pr(> z ) |
|-------------|----------|------------|---------|----------|
| (Intercept) | -6.4863  | 0.7918     | -8.19   | 0.0000   |
| SBP         | 0.0143   | 0.0022     | 6.38    | 0.0000   |
| AGE         | 0.0577   | 0.0142     | 4.06    | 0.0000   |

- iii. Perform both the Wald test and the Likelihood ratio tests (Chi-Squared) to see if the added variable is significant at the .01 level. What are the p-values from each test? Are they the same?

**2.4.1.3 Model building** Start with all variables. Our goal is to fit a well-fitting model, that is still small and easy to interpret (parsimonious).

- i. Use backward selection method. Only keep variables whose coefficients are significantly different from 0 at .05 level. Kick out the variable with the largest p-value first, and then re-fit the model to see if there are other variables you want to kick out.
- ii. Use AIC as the criterion for model selection. Find a model with small AIC through exhaustive search. Does exhaustive search guarantee that the p-values for all the remaining variables are less than .05? Is our final model here the same as the model from backwards elimination?
- iii. Use the model chosen from part ii. as the final model. Write a brief summary to describe important factors relating to Heart Diseases (i.e. the relationships between those variables in the model and heart disease). Give a definition of “important factors”.
- iv. What is the probability that Liz will have heart disease, according to our final model?

## 2.4.2 Classification analysis

### 2.4.2.1 ROC/FDR

- i. Display the ROC curve using `fit1`. Explain what ROC reports and how to use the graph. Specify the classifier such that the False Positive rate is less than .1 and the True Positive rate is as high as possible.
- ii. Overlay two ROC curves: one from `fit1`, the other from `fit2`. Does one curve always contain the other curve? Is the AUC of one curve always larger than the AUC of the other one? Why or why not?
- iii. Estimate the Positive Prediction Values and Negative Prediction Values for `fit1` and `fit2` using .5 as a threshold. Which model is more desirable if we prioritize the Positive Prediction values?
- iv. For `fit1`: overlay two curves, but put the threshold over the probability function as the x-axis and positive prediction values and the negative prediction values as the y-axis. Overlay the same plot for `fit2`. Which model would you choose if the set of positive and negative prediction values are the concerns? If you can find an R package to do so, you may use it directly.

**2.4.2.2 Cost function/ Bayes Rule** Bayes rules with risk ratio  $\frac{a_{10}}{a_{01}} = 10$  or  $\frac{a_{10}}{a_{01}} = 1$ . Use your final model obtained from Part 1 to build a class of linear classifiers.

- i. Write down the linear boundary for the Bayes classifier if the risk ratio of  $a_{10}/a_{01} = 10$ .
- ii. What is your estimated weighted misclassification error for this given risk ratio?
- iii. How would you classify Liz under this classifier?
- iv. Bayes rule gives us the best rule if we can estimate the probability of HD-1 accurately. In practice we use logistic regression as our working model. How well does the Bayes rule work in practice? We hope to show in this example it works pretty well.

Now, draw two estimated curves where  $x = \text{threshold}$ , and  $y = \text{misclassification errors}$ , corresponding to the thresholding rule given in x-axis.

- v. Use weighted misclassification error, and set  $a_{10}/a_{01} = 10$ . How well does the Bayes rule classifier perform?
- vi. Use weighted misclassification error, and set  $a_{10}/a_{01} = 1$ . How well does the Bayes rule classifier perform?