# Regression and Classification Exercises

*Soma S Dhavala*

```
knitr::opts_chunk$set(echo = FALSE)
set.seed(123)

# UCI ML datasets and many simulated datasets available
require(mlbench)
```

```
## Loading required package: mlbench
```

```
path=getwd()

# no of records
n = 100
# no of feature
p = 10

giveColNames <- function(p)
{
  cnames <- c('target',paste(rep("feature_",p),formatC(1:p,width=floor(1+log10(p)),flag=0,format="d"),s
}

writeData <- function(X,y,fname)
{
  p <- ncol(X)
  cnames <- c('target',paste(rep("feature_",p),formatC(1:p,width=floor(1+log10(p)),flag=0,format="d"),s
  # regression problem
  df <- as.data.frame(cbind(y,X))
  colnames(df) <- cnames
  write.csv(df,paste(fname,".regr.csv",sep=""),row.names = F)

  yhat <- y-median(y,na.rm = T)
  prob <- 1/(1+exp(-yhat))
  y <-rbinom(n,1,prob)
  df$target <- y
  write.csv(df,paste(fname,".class.csv",sep=""),row.names = F)
}
```

## 1. Single Feature

**Data**

Generate 100 records with 1 features

**Comments**

Should be straight forward

## 2. Multiple Features

### 2a. independent features

**Data**

Generate 100 records with 10 indepdent features.

**Comments**

- last two features are not important (with lasso, it should produce exact zero)

### 2b correlated features

**Data**

Generate 100 records with 10 correlated features.

**Comments**

- last two features are not important (with lasso, it should produce exact zero)
- gradient descent would be unstable
- variable selection is not consistent

### 2c features of different scale

**Data**

Generate 100 records with 10 indepdent features. Each feature is on a different scale and different mean

**Comments**

- last two features are not important (with lasso, it should produce exact zero)
- gradient descent would be unstable

### 2d correlated features and with different scale

**Data**

Generate 100 records with 10 correlated features. Each feature is on a different scale and different mean

**Comments**

- last two features are not important (with lasso, it should produce exact zero)
- gradient descent would be unstable
- variable selection is not consistent

### 2e correlated features and with different scale, mising data and outliers.

**Data**

Generate 100 records with 10 correlated features. Each feature is on a different scale and different mean

**Comments**

- last two features are not important (with lasso, it should produce exact zero)
- gradient descent would be unstable
- variable selection is not consistent
- regression/classfication are noisy

```
## Warning in rbinom(n, 1, prob): NAs produced
```

# 3 Non-Linear regression

## 3a Friedman-1 benchmark dataset

**Data**

Generarte data from

$y = 10\sin(\pi x_1 x_2) + 20(x_3 - 0.5)^2 + 10x_4 + 5x_5 + e$

It has 100 records and 10 features and only five are used

**Comments**

- can fit linear regression with additional features
- non-parametric method is better in the absence of additional info
- only few features are useful

## 3b Friedman-2 benchmark dataset

**Data**

Generarte data from

$y = (x_1^2 + (x_2 x_3 - (1/x_2 x_4))^2)^{0.5} + e$

It has 100 records and 4 features

**Comments**

- non-parametric method is better in the absence of additional info
- linear models will be poor fit

## 4 Ozone Data Set

**Data**

Leo Breiman, Department of Statistics, UC Berkeley. Data used in Leo Breiman and Jerome H. Friedman (1985), Estimating optimal transformations for multiple regression and correlation, JASA, 80, pp. 580-598.

**Comments**

- predict V4 (non-parametric method is better in the absence of additional info
- linear models will be poor fit

**Exercises**

1. Is a simple linear regression model better choice? Explain in your words what is the functional relationship between the target and the predictor? Can it still be called a linear model?

   - DataSets: 1
   - Miconception: Meaning of Linearity
   - Concepts: run simple linear regression and log-linear model, understand the blackbox, implement simple gradient dsescent and compare model with libraries

2. Is a multiple linear regression model better choice? Explain in your words what is the functional relationship between the target and the predictor?

   - DataSets: >1;
   - Miconception: Meaning of Linearity
   - Concepts: Model Selection, Idea of Baseline Model

3. Comment on the numerical stability of the model fit?

   - DataSets: 2c-2e;
   - Miconception: Ml is black-box approach
   - Concepts: Dataset Standarization, Collinearity, Robust regression, Missing Value treatment

4. Is the model explaing the data? Is your model a good model?

   - DataSets: All;
   - Miconception: ML is a black-box approach, I've THE best model
   - Concepts: Model assessment, explainability vs predictive power

5. Is it necessary to preprocess the data? If yes, what sort of data preparation is needed?

   - DataSets: >1;
   - Miconception: I will be given nice, clean data, all that I need to do is just call a function.
   - Concepts: Data cleaning, transformations, check residuals, Iterate between input-model-output-validate
   - Methods: Best subset selection (forward, backward, stagewise), lasso, LARS

6. Provide diagnostic plots and critique the model fit

   - DataSets: All;
   - Miconception:
   - Concepts: Residual plots, Generalization Error, Test and Train errors, Model fit statistics such as AIC, BIC

   - Techniques/Methods: Cross-Validation