

# Problem Set 2

API222: Big Data & Machine Learning

Brendan Hellweg

10/14/2021

## Problem 1

1 LDA and Logistic regression both have linear decision boundaries and will perform similarly poorly when applied to nonlinear decision boundaries. QDA and KNN both have nonlinear decision boundaries. Depending on the type of nonlinearity of the Bayes Decision Boundary, either QDA or KNN May perform better. KNN will perform better if the assumption of normalcy does not hold, while QDA performs better when the assumption of normalcy holds.

## Problem 2

**Problem 2A** See the table below.

##	X1	X2	X3	Y	distance
## 1	0	3	2	Blue	3.61
## 2	0	2	3	Blue	3.61
## 3	0	4	1	Orange	4.12
## 4	1	3	1	Green	3.32
## 5	1	5	0	Orange	5.10
## 6	-1	4	2	Orange	4.58
## 7	0	4	2	Orange	4.47

**Problem 2B** The prediction with  $K = 1$  is **Green**. The decision criteria for KNN is to select the  $K$  data points closest to the test point and assign a value based on the observed value. In this case, the one nearest neighbor (where the Euclidian distance is minimized) has property  $Y = \text{Green}$ .

**Problem 2C** The prediction with  $K = 3$  is **Blue**. For the three nearest neighbors (where the Euclidian distance is minimized), two have the property  $Y = \text{Blue}$  and one has property  $Y = \text{Green}$ . Different KNN tests will have different decision criteria for inconsistent values, but typically the value with the most or majority observations in the set is what's selected.

**Problem 2D** If the Bayes decision boundary is highly nonlinear, we'd want a small  $K$  so we can be sensitive to nonlinearities.

## Problem 3

When  $\lambda = 0$ , the Lasso model becomes equivalent to the residual sum of squares, therefore becoming identical to the OLS model. When  $\lambda$  approaches infinity, all coefficients approach zero and the model becomes a constant function.

## Problem 4

While Lasso uses  $\|\beta\|_1$  with the sums of absolute values of  $\beta_j$ , Ridge uses  $\|\beta\|_2$  with the sums of squares of  $\beta_j$ . Using  $\|\beta\|_1$  forces some of the coefficient estimates to reach exactly zero, meaning that Lasso produces a sparse and therefore more interpretable model. Ridge, on the other hand, produces a dense and therefore harder to interpret model.

If all I care about is the predictive ability and I don't care at all about the level of interpretability, I would use Cross Validation to compare the two potential models (after also using it to identify the optimal values of  $\lambda$  to design each model). Cross Validation works by estimating the test error rate by holding aside subsets of the training observations in the fitting process and then applying the model to those saved subsets. The metric to compare is the test MSE.

## Problem 5

We cannot draw this conclusion, as we do not know where along the ROC curve we want to minimize the number of false positives. Model A could start steeper and level off sooner than Model B, meaning that for early parts of the ROC curve it outperforms Model B. Conversely, Model B could dominate throughout the ROC curve. Therefore, we cannot make a conclusion from the information provided.

## Data Questions

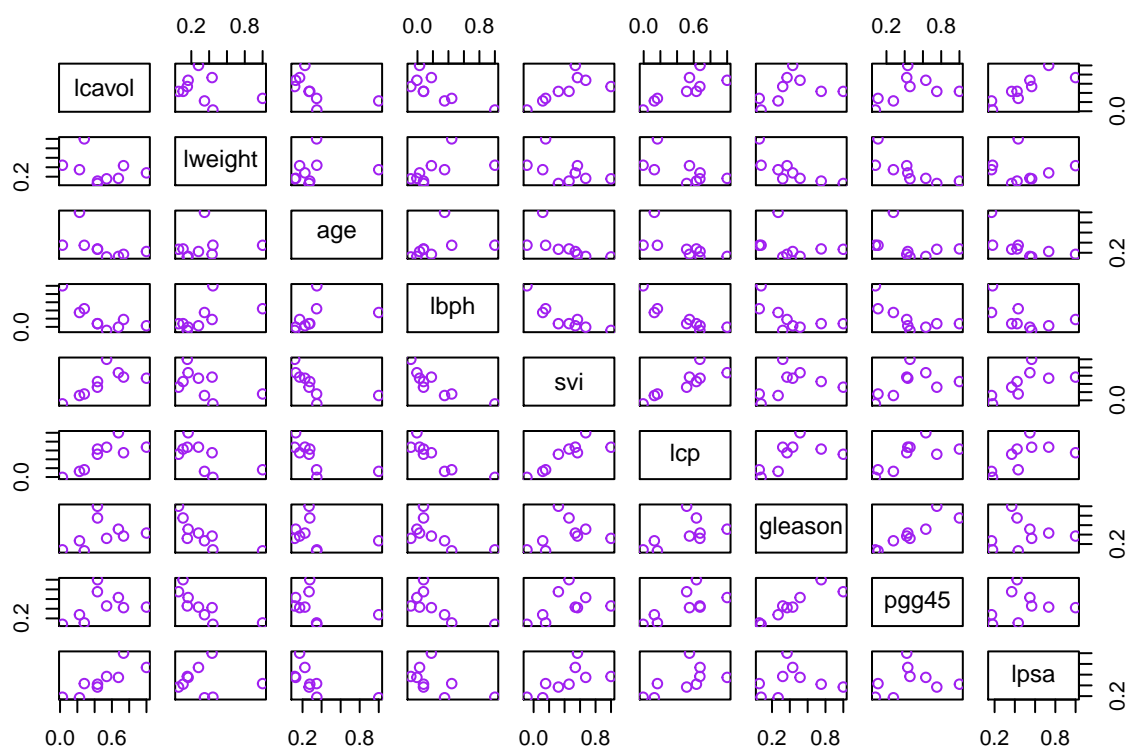
### Problem 1

```
c("The dimensions are: ",  
  nrow(prostate), " by ", ncol(prostate),".",  
  " The final column represents whether the entry belongs in the test or the treatment set.") %>%  
str_c(.,collapse = '')
```

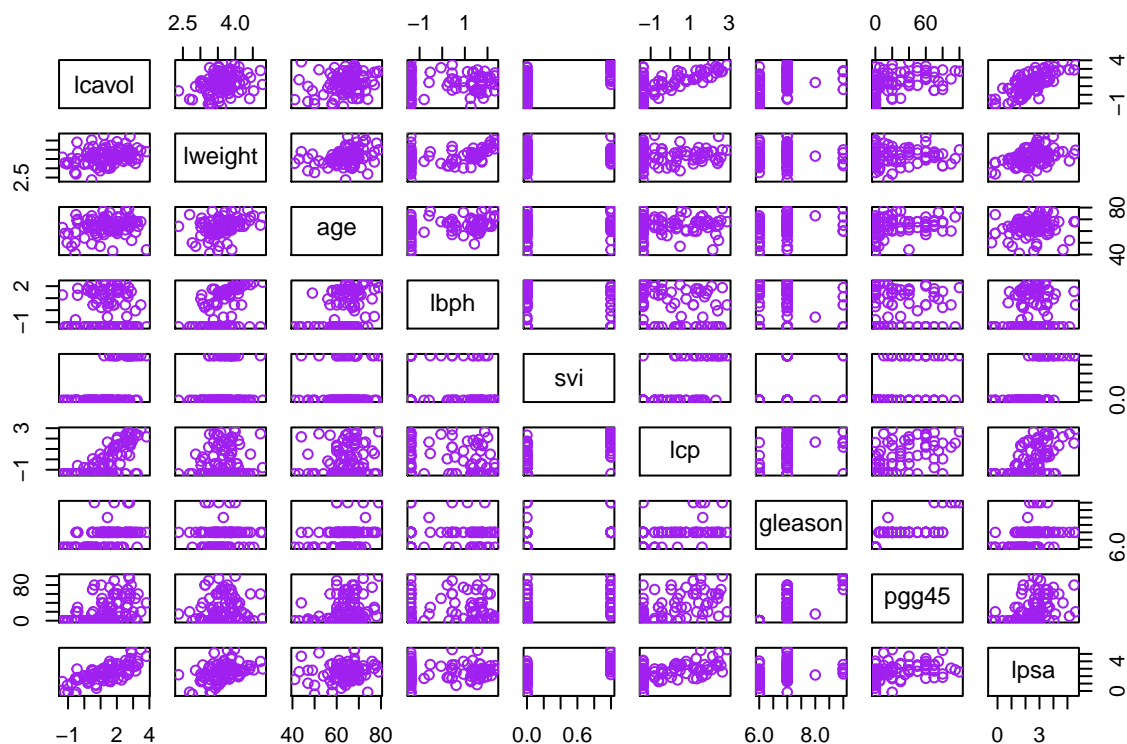
```
## [1] "The dimensions are: 97 by 10. The final column represents whether the entry belongs in the test
```

### Problem 2

```
pairs(cor(prostate[1:9]), col = "purple")
```



```
pairs((prostate[1:9]), col = "purple", rowlatop = T)
```



The Gleason variable has a small number of possible values in this dataset, so the results will be constrained to a small number of bands, which are vertical because the other variable has a strong degree of variation along a given value of gleason.

### Problem 3

```
seatbelts2 <- drop_na(USSeatBelts)
nrow(seatbelts) - nrow(seatbelts2)
```

```
## [1] 209
```

### Problem 4

```
cats <- seatbelts2 %>%
  select(where(is.factor))
names(cats)
```

```
## [1] "state" "year" "speed65" "speed70" "drinkage" "alcohol" "enforce"
```

```
str_c(names(cats[1]),": ",nrow(distinct(cats[1])),collapse = '')
```

```
## [1] "state: 51"
```

```
str_c(names(cats[2]),": ",nrow(distinct(cats[2])),collapse = '')
```

```
## [1] "year: 15"
```

```
str_c(names(cats[3]),": ",nrow(distinct(cats[3])),collapse = '')
```

```
## [1] "speed65: 2"
```

```
str_c(names(cats[4]),": ",nrow(distinct(cats[4])),collapse = '')
```

```
## [1] "speed70: 2"
```

```
str_c(names(cats[5]),": ",nrow(distinct(cats[5])),collapse = '')
```

```
## [1] "drinkage: 2"
```

```
str_c(names(cats[6]),": ",nrow(distinct(cats[6])),collapse = '')
```

```
## [1] "alcohol: 2"
```

```
str_c(names(cats[7]),": ",nrow(distinct(cats[7])),collapse = '')
```

```
## [1] "enforce: 3"
```

## Problem 5

```
seatbelts2 %>%  
  slice_max(.,order_by = .$fatalities,n = 1)
```

```
##   state year miles fatalities seatbelt speed65 speed70 drinkage alcohol income  
## 1    SC 1985 26677 0.03564868   0.128      no      no      no      no 11206  
##      age enforce  
## 1 33.42841      no
```

```
seatbelts2 %>%  
  slice_min(.,order_by = .$fatalities,n = 1)
```

```
##   state year miles fatalities seatbelt speed65 speed70 drinkage alcohol income  
## 1    MA 1996 49956 0.008327328   0.54     yes      no     yes      no 29591  
##      age  enforce  
## 1 37.25541 secondary
```

## Problem 6

```

seatbelts3 <- seatbelts2 %>%
  select(-c(1:2)) %>%
  dummy_cols(.,select_columns = c('speed65',
                                   'speed70',
                                   'drinkage',
                                   'alcohol',
                                   'enforce')) %>%

  select(-c(4:7,10))

seatbelt_lm <- lm(fatalities~.,seatbelts3)
summary(seatbelt_lm)

```

```

##
## Call:
## lm(formula = fatalities ~ ., data = seatbelts3)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.0097759 -0.0022532 -0.0002876  0.0020290  0.0140265
##
## Coefficients: (5 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.828e-02  3.889e-03   9.841 < 2e-16 ***
## miles        -1.491e-09  3.353e-09  -0.445  0.656810
## seatbelt      -7.528e-04  1.696e-03  -0.444  0.657267
## income       -7.959e-07  4.869e-08 -16.345 < 2e-16 ***
## age          -5.546e-05  1.136e-04  -0.488  0.625587
## speed65_no    -1.593e-04  4.623e-04  -0.345  0.730528
## speed65_yes           NA          NA      NA      NA
## speed70_no    -2.232e-03  5.352e-04  -4.171  3.54e-05 ***
## speed70_yes           NA          NA      NA      NA
## drinkage_no     1.174e-03  9.092e-04   1.291  0.197160
## drinkage_yes           NA          NA      NA      NA
## alcohol_no     1.824e-03  4.803e-04   3.798  0.000162 ***
## alcohol_yes           NA          NA      NA      NA
## enforce_no     -9.388e-04  5.194e-04  -1.807  0.071252 .
## enforce_primary 9.109e-04  5.273e-04   1.727  0.084680 .
## enforce_secondary NA          NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.003532 on 545 degrees of freedom
## Multiple R-squared:  0.5156, Adjusted R-squared:  0.5067
## F-statistic: 58.01 on 10 and 545 DF,  p-value: < 2.2e-16

```

```

str_c("Adjusted R^2: ",round(100*seatbelt_lm[["residuals"]][["98"]],2))

```

```

## [1] "Adjusted R^2: 0.51"

```

## **Setup for Regression & KNN Section**

For the next section, I will filter the dataset to remove rows with NA values and create treatment and test bins.

### **Problem 7**

### **Problem 8**

### **Problem 9**

Note: This value has a random element due to the process of generating MSE\_TE of a KNN model. Over several instances running this code block, I found values ranging from approx. 450 and 650.