

CS 199
Final Project - MOOC Data Analysis

May 16, 2014

Sam Laane <laane2@illinois.edu>
José Vicente Ruiz <ruizcep2@illinois.edu>
Nick Jeffrey <njeffre2@illinois.edu>

Contents

| | | |
|----------|-------------------------------|----------|
| 1 | Introduction | 2 |
| 2 | Predicting exam scores | 2 |
| 2.1 | Implementation | 2 |
| 3 | Finding correlations | 2 |
| 3.1 | Implementation | 2 |
| 3.1.1 | Results | 4 |
| 3.2 | Conclusions | 6 |
| 3.2.1 | Correlations | 6 |

1 Introduction

This assignment consisted of analysing the provided *MOOC data* and determining whether we could predict final exam scores or determine whether bloom taxonomy levels actually provided a meaningful differentiation between questions.

For that purpose, a dataset was provided by *Magical Data Fairies*. The programming languages used have been R and Python, and this report have been typeset using \LaTeX .

2 Predicting exam scores

2.1 Implementation

???

3 Finding correlations

3.1 Implementation

```

1  # Define functions.
2  printf <- function(...) invisible(print(sprintf(...)))
3  library(ggplot2)
4
5  # Read the data.
6  toRm = c("Q4Q07_corr", "Q4Q13_corr")
7  dataWithNames<- read.csv('combinednodropna.csv', h=T) # Headers = True.
8
9  #remove user id's
10 data <- dataWithNames[-1]
11
12 # remove values with std of zero
13
14 toRm = c()
15 #check for std of 0
16 for (d in names(data)){
17   s = sd(data[[d]], na.rm = TRUE)
18   if (s == 0){
19     print(d)
20     toRm = c(toRm, d)
21   }
22 }
23 data <- data[,!(names(data)) %in% toRm]
24
25 #quizData <- data[,grep(sprintf("^Q%dQ[0-9][0-9]_corr$",4),names(data))
26   ]
27 #testData <- data[,-(1:(20 * 8-2))]
28 #ncol(quizData)
29 #View(quizData)
30
31 units = c(0,1,2,3,4,5,6,7)
32 #units = c(4,5,6)
33 #quizData <- data[,1:(20*8-2)]

```

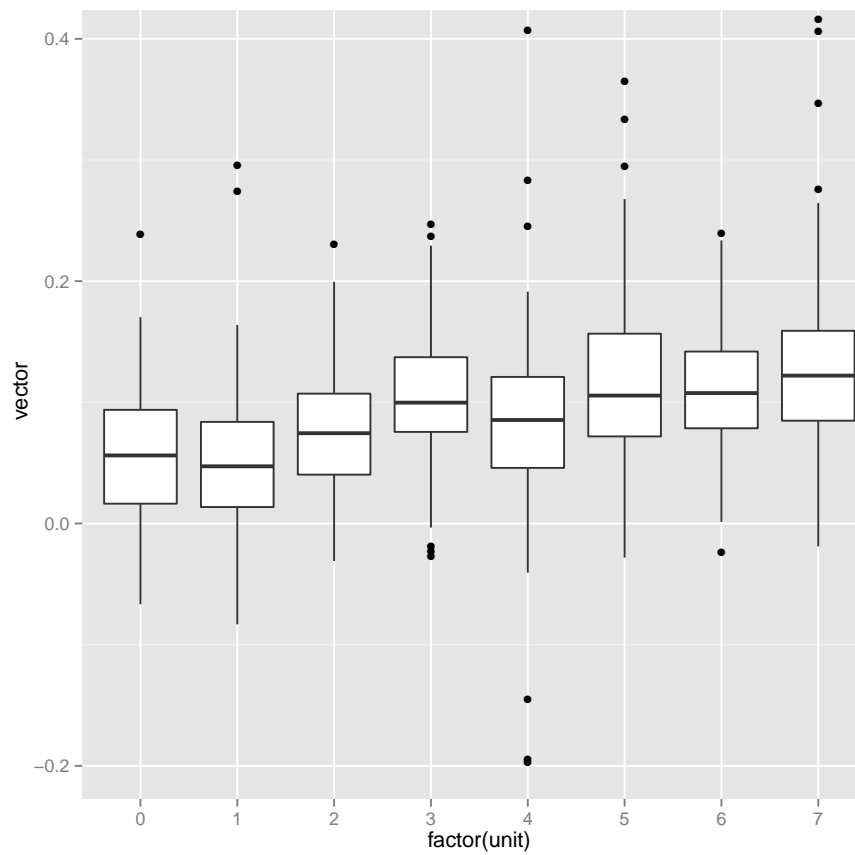
```

33 #testData <- data[,-(1:(20 * 8-2))]
34
35 #calculate Correlation matrix
36
37 correlationMatrix = cor(data, use = "na.or.complete")
38
39
40 #find how correlations compare.
41 correlationSums <- (apply(correlationMatrix, 2, sum )) -1
42 correlationMeans <- correlationSums / (ncol(data) - 1 )
43 #View(correlationMeans)
44 #print(mean(correlationMeans))
45
46
47 #now lets compare the quiz to the test
48 quizCol <- correlationMatrix[,1:(20 * 8-2)]
49 testRow <- quizCol[-(1:(20 * 8-2)), ]
50 #View(testRow)
51 df <- data.frame(vector = numeric(), unit = numeric())
52 dff <- data.frame(vector = numeric(), unit = numeric())
53 linedata <- data.frame(unit = numeric(), same = numeric(), others =
    numeric(), diff = numeric())
54 units = c(0,1,2,3,4,5,6,7)
55 for (i in units){
56     questionsInUnit <- testRow[,grep(sprintf("^Q%dQ[0-9][0-9]_corr$",i)
    ,names(data))]
57     t <- questionsInUnit[(1:5) + 5 * i,]
58     f <- questionsInUnit[-((1:5) + 5 * i),]
59     corInUnit <- mean(t)
60     corOutUnit <- mean(f)
61     printf("mean correlation of week %d quiz to week %d test questions:
    %f", i, i, corInUnit)
62     printf("max correlation of week %d quiz to week %d test questions:
    %f", i, i, max(t))
63     printf("mean correlation of week %d quiz to other test questions: %
    f", i, corOutUnit)
64     printf("max correlation of week %d quiz to other test questions: %
    f", i, max(f))
65     printf("the diffence of correlation: %f", corInUnit - corOutUnit)
66     #View(t)
67     #Sys.sleep(1200)
68     df <- rbind(df, data.frame(vector= as.vector(t), unit=i))
69     dff <- rbind(dff, data.frame(vector= as.vector(f), unit=i))
70     linedata <- rbind(linedata, data.frame(unit = i, same = mean(t),
    other = mean(f), diff = corInUnit - corOutUnit))
71 }
72 ggplot(dff, aes(y=vector, x=factor(unit))) + geom_boxplot()
73 ggsave('unrelatedcorrelationbox.pdf')

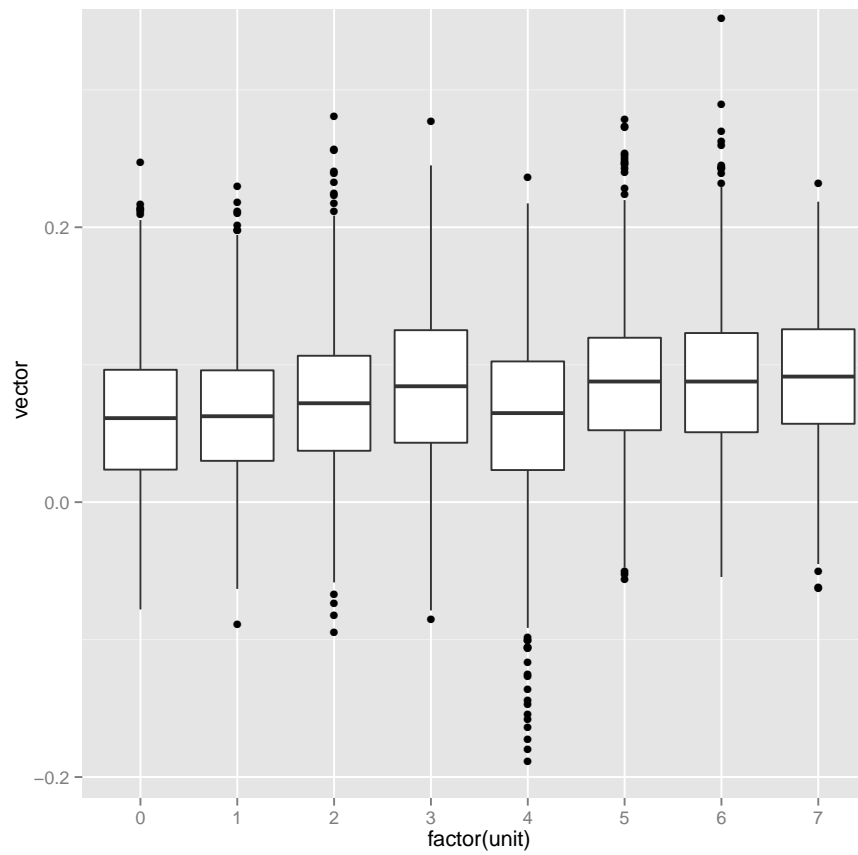
```

3.1.1 Results

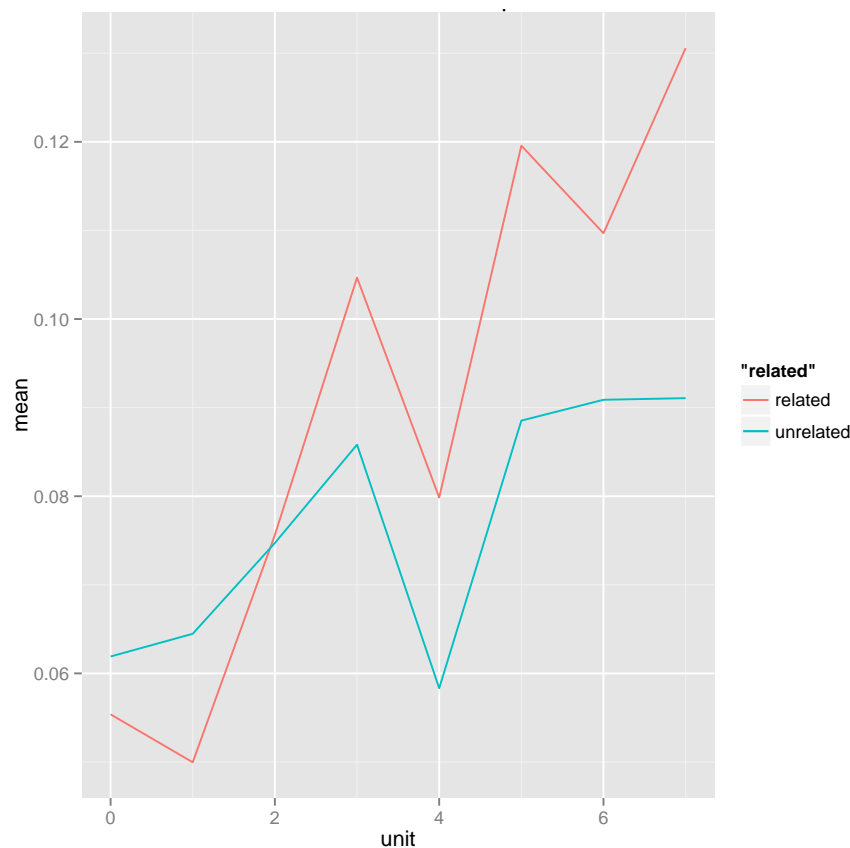
This figure shows a box and whisker plot of correlation for every question on a quiz to every related question on the final.



This figure shows a Box and Whisker plot for every question on a quiz to the unrelated questions on the final.



This figure shows the mean correlations of quiz questions on related and unrelated questions on the final side by side.



3.2 Conclusions

3.2.1 Correlations

Judging from the graphs provided, the correlation of quiz scores with the final had a few interesting trends. One was an outcome we expected, which is that quizzes later in the course more strongly correlate with their questions on the final, most likely due to the simple fact that the material is fresher in the student's minds. This trend can be seen both in the box plot and in the line graph.

Another trend can be seen in the box plots as well: That there are many outliers, especially later on in the course. The outliers that have much higher correlation are easy to explain because one would expect that questions which are testing the same material would have stronger correlations than those testing different material in the same subject. However, there are also outliers in the other direction, even going so far as to have a few questions that are negatively correlated with exam scores. We suspect that this could be due to those being trick questions, that students only studied because they happened to guess wrongly on the quiz, but without having access to the questions themselves, this can at best only be a guess as to why these trends are appearing.