

Exploratory Analysis

Cool Beans Programming

2023-04-22

Loading libraries and data

```
library(ISLR2)
college <- read.csv("https://www.statlearning.com/s/College.csv")
rownames(college) <- college[, 1]
View(college)
```

Variables

Private: public/private indicator

Apps: number of applications received

Accept: number of applications accepted

Enroll: number of new students enrolled

Top10perc: new students from the top 10% of high school class

Top25perc: new students from the top 25% of high school class

F.Undergrad: number of full-time undergraduates

P.Undergrad: number of part-time undergraduates

Outstate: out-of-state tuition

Room.Board: room and board costs

Books: estimated book costs

Personal: estimated personal spending

PhD: percent of faculty with Ph.D.'s

Terminal: percent of faculty with terminal degree

S.F. Ratio: student/faculty ratio

perc.alumni: percent of alumni who donate

Expend: instructional expenditure per student

Grad.Rate: graduation rate

Exploring the data

```
college <- college[, -1]
head(college)
```

```
##               Private Apps Accept Enroll Top10perc Top25perc
## Abilene Christian University    Yes 1660   1232    721      23      52
## Adelphi University              Yes 2186   1924    512      16      29
## Adrian College                  Yes 1428   1097    336      22      50
## Agnes Scott College              Yes  417    349    137      60      89
## Alaska Pacific University        Yes  193    146     55      16      44
## Albertson College                Yes  587    479    158      38      62
##               F.Undergrad P.Undergrad Outstate Room.Board Books
## Abilene Christian University    2885      537    7440      3300    450
## Adelphi University              2683     1227   12280      6450    750
## Adrian College                  1036       99   11250      3750    400
## Agnes Scott College              510       63   12960      5450    450
## Alaska Pacific University        249      869    7560      4120    800
## Albertson College                678       41   13500      3335    500
##               Personal PhD Terminal S.F.Ratio perc.alumni Expend
## Abilene Christian University    2200   70      78     18.1      12    7041
## Adelphi University              1500   29      30     12.2      16   10527
## Adrian College                  1165   53      66     12.9      30    8735
## Agnes Scott College              875   92      97      7.7      37   19016
## Alaska Pacific University        1500   76      72     11.9       2   10922
## Albertson College                675   67      73      9.4      11    9727
##               Grad.Rate
## Abilene Christian University     60
## Adelphi University               56
## Adrian College                   54
## Agnes Scott College               59
## Alaska Pacific University         15
## Albertson College                 55
```

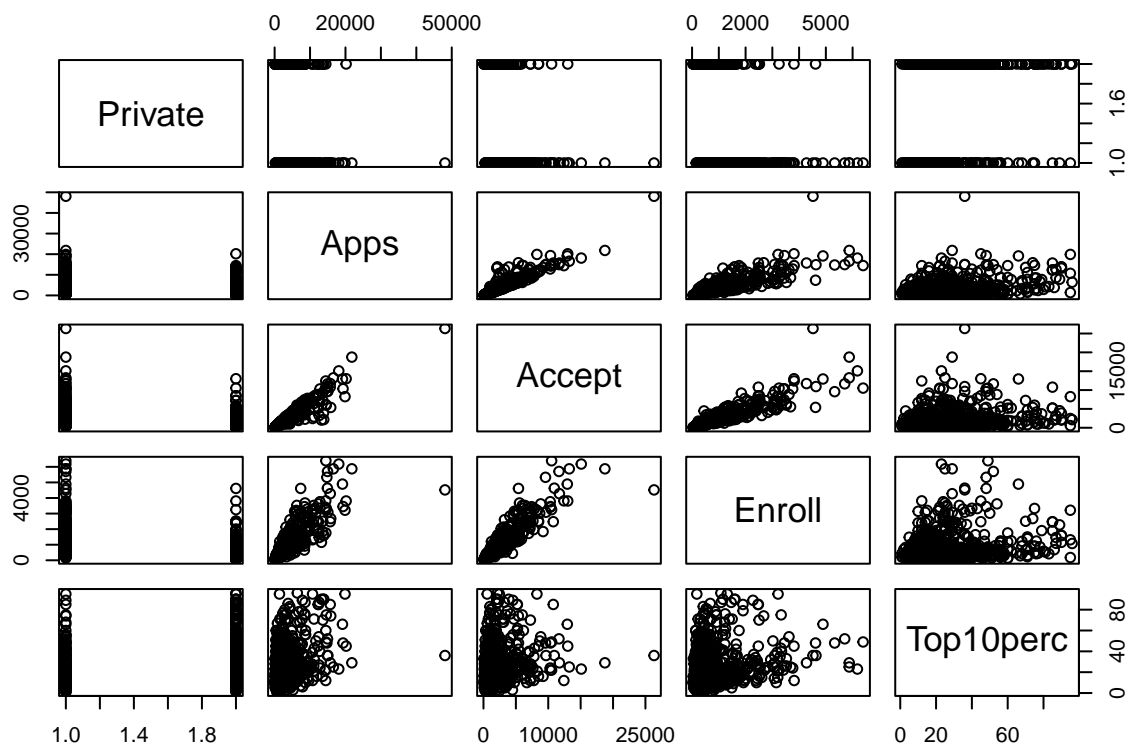
```
summary(college)
```

```
##      Private      Apps      Accept      Enroll
## Length:777      Min.   : 81      Min.   : 72      Min.   : 35
## Class :character 1st Qu.: 776      1st Qu.: 604      1st Qu.: 242
## Mode  :character Median : 1558      Median : 1110      Median : 434
##               Mean   : 3002      Mean   : 2019      Mean   : 780
##               3rd Qu.: 3624      3rd Qu.: 2424      3rd Qu.: 902
##               Max.   :48094      Max.   :26330      Max.   :6392
##      Top10perc      Top25perc      F.Undergrad      P.Undergrad
## Min.   : 1.00      Min.   : 9.0      Min.   : 139      Min.   : 1.0
## 1st Qu.:15.00      1st Qu.: 41.0      1st Qu.: 992      1st Qu.: 95.0
## Median :23.00      Median : 54.0      Median : 1707      Median : 353.0
## Mean   :27.56      Mean   : 55.8      Mean   : 3700      Mean   : 855.3
## 3rd Qu.:35.00      3rd Qu.: 69.0      3rd Qu.: 4005      3rd Qu.: 967.0
## Max.   :96.00      Max.   :100.0      Max.   :31643      Max.   :21836.0
##      Outstate      Room.Board      Books      Personal
## Min.   : 2340      Min.   :1780      Min.   : 96.0      Min.   : 250
```

```
## 1st Qu.: 7320    1st Qu.:3597    1st Qu.: 470.0    1st Qu.: 850
## Median : 9990    Median :4200    Median : 500.0    Median :1200
## Mean   :10441    Mean   :4358    Mean   : 549.4    Mean   :1341
## 3rd Qu.:12925    3rd Qu.:5050    3rd Qu.: 600.0    3rd Qu.:1700
## Max.   :21700    Max.   :8124    Max.   :2340.0    Max.   :6800
##      PhD      Terminal      S.F.Ratio      perc.alumni
## Min.    : 8.00    Min.    : 24.0    Min.    : 2.50    Min.    : 0.00
## 1st Qu.: 62.00    1st Qu.: 71.0    1st Qu.:11.50    1st Qu.:13.00
## Median : 75.00    Median : 82.0    Median :13.60    Median :21.00
## Mean   : 72.66    Mean   : 79.7    Mean   :14.09    Mean   :22.74
## 3rd Qu.: 85.00    3rd Qu.: 92.0    3rd Qu.:16.50    3rd Qu.:31.00
## Max.   :103.00    Max.   :100.0    Max.   :39.80    Max.   :64.00
##      Expend      Grad.Rate
## Min.    : 3186    Min.    : 10.00
## 1st Qu.: 6751    1st Qu.: 53.00
## Median : 8377    Median : 65.00
## Mean   : 9660    Mean   : 65.46
## 3rd Qu.:10830    3rd Qu.: 78.00
## Max.   :56233    Max.   :118.00
```

Each row indicates a different college or university. Most of the variables are quantitative while only the private/public indicator variable is categorical. Next, we convert this indicator to a factor variable and create a scatter plot matrix of the first 5 variables.

```
# convert Private to factor variable
college$Private <-as.factor(college$Private)
pairs(college[,1:5])
```



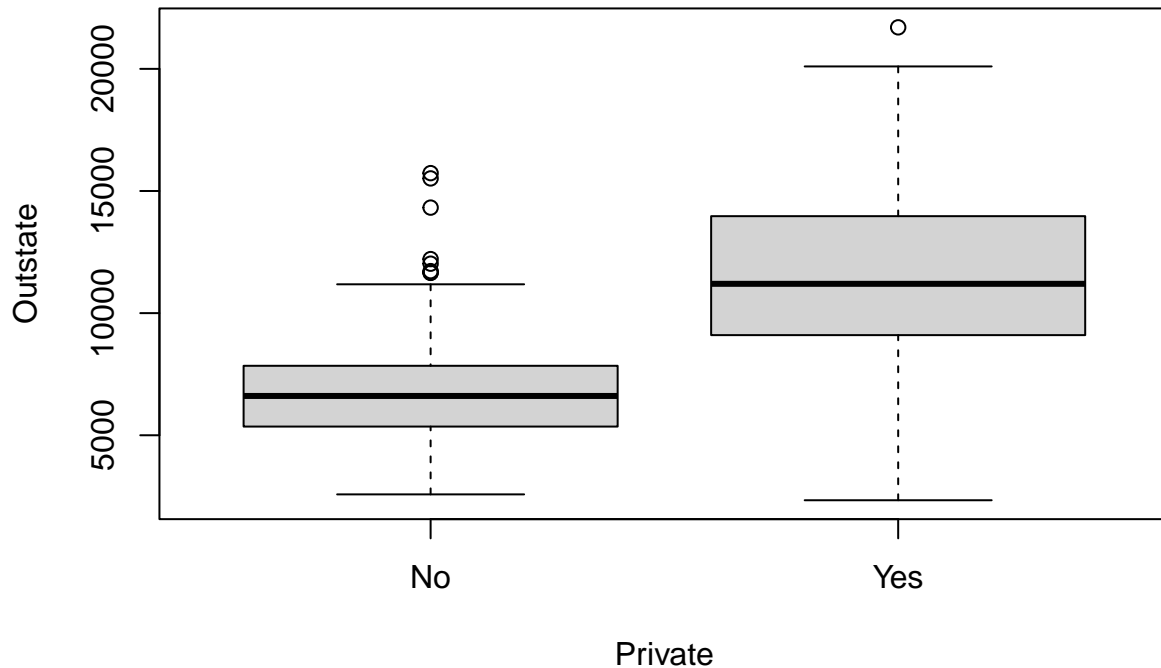
From the graph, Apps and Accept appear to have a strong linear relationship with maybe one outlier. Accept and Enroll have a linear relationship as well.

Graphing the data

Side-by-side box plots of out of state tuition by private and public schools indicates that the median out of state tuition for private schools is larger than that of out of state public students, while public out of state students have a larger distribution.

```
boxplot(college$Outstate ~ college$Private, xlab= "Private", ylab="Outstate",
        main="Out of State Tuition vs Private")
```

Out of State Tuition vs Private



Feature Engineering

We now create a new qualitative variable called Elite by binning the Top10perc variable. The schools will be divided into two groups based on if the proportion of students coming from the top 10% of their high school classes exceeds 50%.

```
Elite <-rep("No", nrow(college))
Elite[college$Top10perc > 50] <- "Yes"
Elite <-as.factor(Elite)
college <-data.frame(college,Elite)
```

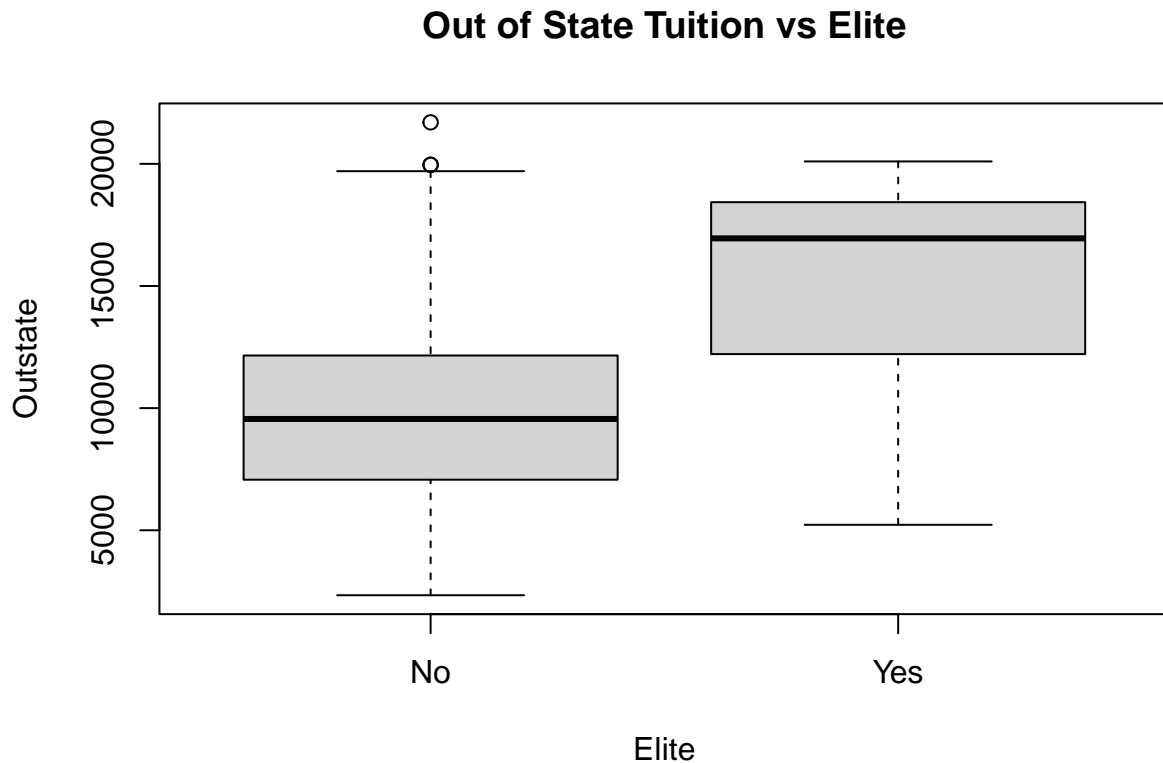
The majority of students are not categorized as elite; only 11.1% of students fall into this category.

```
summary(college)
```

```
## Private      Apps      Accept      Enroll      Top10perc
## No :212   Min.    :   81   Min.    :   72   Min.    :   35   Min.    : 1.00
## Yes:565   1st Qu.:  776   1st Qu.:  604   1st Qu.:  242   1st Qu.:15.00
##          Median : 1558   Median : 1110   Median :  434   Median :23.00
##          Mean   : 3002   Mean   : 2019   Mean   :  780   Mean   :27.56
##          3rd Qu.: 3624   3rd Qu.: 2424   3rd Qu.:  902   3rd Qu.:35.00
##          Max.   :48094   Max.   :26330   Max.   :6392   Max.   :96.00
## Top25perc  F.Undergrad  P.Undergrad      Outstate
## Min.    :   9.0   Min.    :  139   Min.    :   1.0   Min.    : 2340
## 1st Qu.:  41.0   1st Qu.:  992   1st Qu.:  95.0   1st Qu.: 7320
## Median :  54.0   Median : 1707   Median : 353.0   Median : 9990
```

```
## Mean : 55.8 Mean : 3700 Mean : 855.3 Mean :10441
## 3rd Qu.: 69.0 3rd Qu.: 4005 3rd Qu.: 967.0 3rd Qu.:12925
## Max. :100.0 Max. :31643 Max. :21836.0 Max. :21700
## Room.Board Books Personal PhD
## Min. :1780 Min. : 96.0 Min. : 250 Min. : 8.00
## 1st Qu.:3597 1st Qu.: 470.0 1st Qu.: 850 1st Qu.: 62.00
## Median :4200 Median : 500.0 Median :1200 Median : 75.00
## Mean :4358 Mean : 549.4 Mean :1341 Mean : 72.66
## 3rd Qu.:5050 3rd Qu.: 600.0 3rd Qu.:1700 3rd Qu.: 85.00
## Max. :8124 Max. :2340.0 Max. :6800 Max. :103.00
## Terminal S.F.Ratio perc.alumni Expend
## Min. : 24.0 Min. : 2.50 Min. : 0.00 Min. : 3186
## 1st Qu.: 71.0 1st Qu.:11.50 1st Qu.:13.00 1st Qu.: 6751
## Median : 82.0 Median :13.60 Median :21.00 Median : 8377
## Mean : 79.7 Mean :14.09 Mean :22.74 Mean : 9660
## 3rd Qu.: 92.0 3rd Qu.:16.50 3rd Qu.:31.00 3rd Qu.:10830
## Max. :100.0 Max. :39.80 Max. :64.00 Max. :56233
## Grad.Rate Elite
## Min. : 10.00 No :699
## 1st Qu.: 53.00 Yes: 78
## Median : 65.00
## Mean : 65.46
## 3rd Qu.: 78.00
## Max. :118.00
```

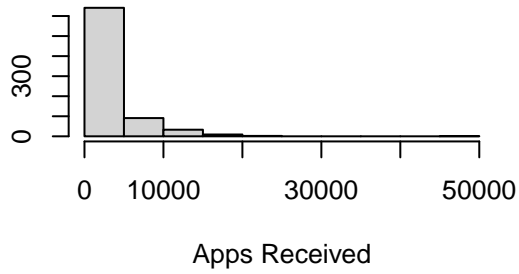
```
boxplot(college$Outstate ~ college$Elite, xlab="Elite",ylab="Outstate",
        main="Out of State Tuition vs Elite")
```



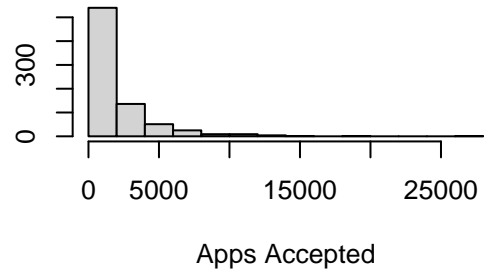
The distribution of non-elite students paying out of state tuition greater than that of elite, out of state students. The median out of state tuition price is greater for elite students than for non-elite students.

```
par(mfrow=c(2,2))
hist(college$Apps, xlab="Apps Received", ylab="", main="Dist. of Applications Received")
hist(college$Accept, xlab="Apps Accepted", ylab="", main="Dist. of Applications Accepted")
hist(college$Room.Board, xlab="Room and Board Cost", ylab="", main="Dist. of Room & Board Costs")
hist(college$Books, xlab="Estimated Book Cost", ylab="", main="Dist. of Textbook Costs")
```

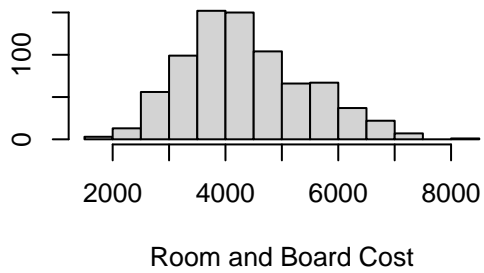
Dist. of Applications Received



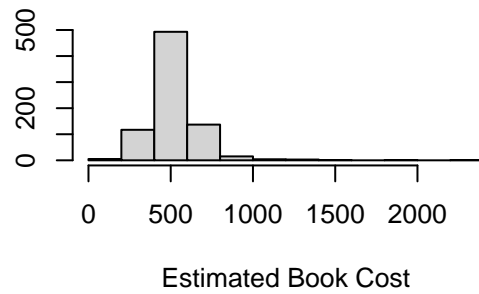
Dist. of Applications Accepted



Dist. of Room & Board Costs



Dist. of Textbook Costs



Histograms are plotted for 4 different variables. The distribution of applications received and applications accepted is right skewed while room & board costs and textbook costs are more normally distributed.