

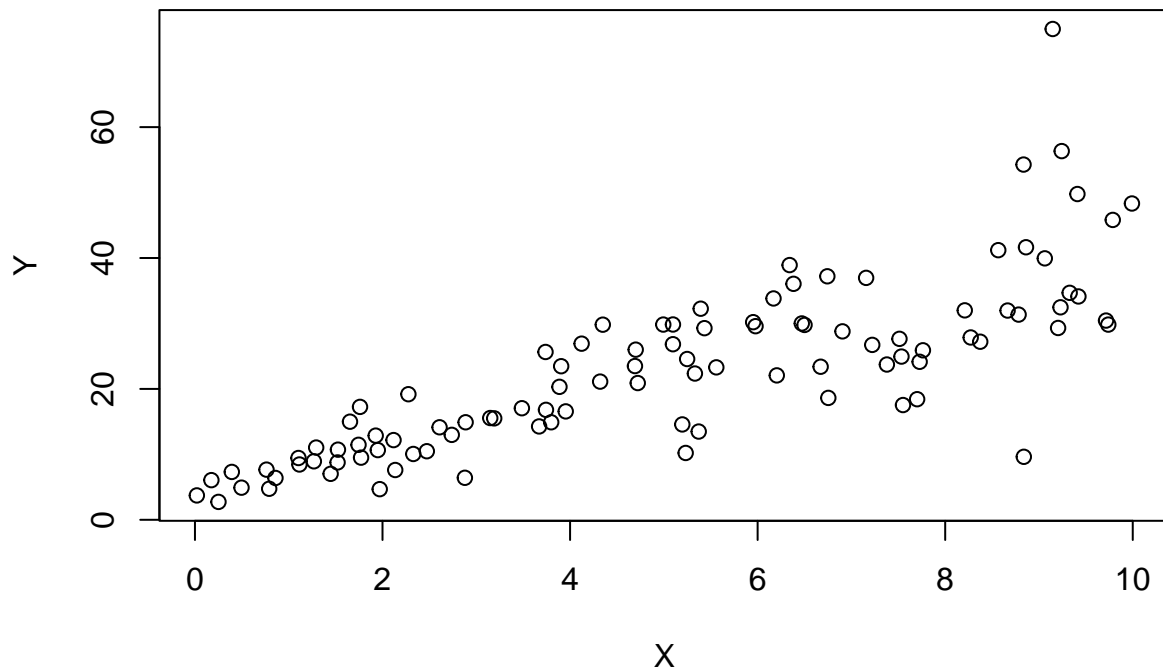
Regression Assignment

Venkata Naga Siddartha Gutha

2022-11-13

1. a)

```
set.seed(2017)
X=runif(100)*10
Y=X*4+3.45
Y=rnorm(100)*0.29*Y+Y
plot(Y~X)
```



from the graph it can be seen that value of y increases with increase in x. This indicates that there is some relation between X and Y. So, yes they can be put into a linear model to explain y based on x

1. b)

```
model<-lm(Y~X)
model
```

```
##
## Call:
## lm(formula = Y ~ X)
##
## Coefficients:
## (Intercept)          X
##      4.465      3.611
```

```
summary(model)
```

```
##
## Call:
## lm(formula = Y ~ X)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -26.755  -3.846  -0.387   4.318  37.503
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   4.4655     1.5537   2.874  0.00497 **
## X              3.6108     0.2666  13.542 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.756 on 98 degrees of freedom
## Multiple R-squared:  0.6517, Adjusted R-squared:  0.6482
## F-statistic: 183.4 on 1 and 98 DF,  p-value: < 2.2e-16
```

Equation: $Y = 4.465 + (3.611) \cdot X$

Accuracy of a model is determined by r Squared value and in this case it is 0.6517, this means that the model explains 65.17% variability of the response variable i.e Y.

1. c)

In the above case Correlation Coefficient is equal to the Coefficient of Determination (R^2) as the regression is based on single variable. Therefore Correlation Coefficient = $R^2 = 0.6517$

2. a)

```
head(mtcars)
```

```
##           mpg  cyl  disp  hp drat   wt  qsec vs  am  gear  carb
## Mazda RX4    21.0   6  160  110 3.90 2.620 16.46 0   1    4     4
## Mazda RX4 Wag 21.0   6  160  110 3.90 2.875 17.02 0   1    4     4
## Datsun 710    22.8   4  108   93 3.85 2.320 18.61 1   1    4     1
## Hornet 4 Drive 21.4   6  258  110 3.08 3.215 19.44 1   0    3     1
## Hornet Sportabout 18.7   8  360  175 3.15 3.440 17.02 0   0    3     2
## Valiant      18.1   6  225  105 2.76 3.460 20.22 1   0    3     1
```

```
lm_james<-lm(hp~wt,data = mtcars)
summary(lm_james)
```

```
##
## Call:
## lm(formula = hp ~ wt, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -83.430 -33.596 -13.587   7.913 172.030
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -1.821     32.325  -0.056   0.955
## wt             46.160      9.625   4.796 4.15e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 52.44 on 30 degrees of freedom
## Multiple R-squared:  0.4339, Adjusted R-squared:  0.4151
## F-statistic:    23 on 1 and 30 DF,  p-value: 4.146e-05
```

```
lm_chris<-lm(hp~mpg,data = mtcars)
summary(lm_chris)
```

```
##
## Call:
## lm(formula = hp ~ mpg, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -59.26 -28.93 -13.45  25.65 143.36
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   324.08     27.43  11.813 8.25e-13 ***
## mpg           -8.83      1.31  -6.742 1.79e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 43.95 on 30 degrees of freedom
## Multiple R-squared:  0.6024, Adjusted R-squared:  0.5892
## F-statistic: 45.46 on 1 and 30 DF,  p-value: 1.788e-07
```

Chris is correct as the model made on Chris opinion has an accuracy of 0.6024 which is higher than accuracy of 0.4339 of the model made on James opinion

2. b)

```
model_hp<-lm(hp~cyl+mpg,data = mtcars)
summary(model_hp)
```

```
##
## Call:
## lm(formula = hp ~ cyl + mpg, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -53.72 -22.18 -10.13  14.47 130.73
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   54.067     86.093   0.628  0.53492
## cyl           23.979      7.346   3.264  0.00281 **
## mpg          -2.775      2.177  -1.275  0.21253
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 38.22 on 29 degrees of freedom
## Multiple R-squared:  0.7093, Adjusted R-squared:  0.6892
## F-statistic: 35.37 on 2 and 29 DF,  p-value: 1.663e-08
```

```
predict(model_hp, data.frame(cyl=4, mpg=22))
```

```
##           1
## 88.93618
```

From the model it is estimated that the Horse Power of a car with 4 calendar and mpg of 22 is 88.93618.

3. a)

```
library(mlbench)
```

```
## Warning: package 'mlbench' was built under R version 4.2.2
```

```
data("BostonHousing")
head(BostonHousing)
```

```
##      crim zn indus chas    nox    rm  age    dis rad tax ptratio    b lstat
## 1 0.00632 18  2.31    0 0.538 6.575 65.2 4.0900    1 296    15.3 396.90  4.98
## 2 0.02731  0  7.07    0 0.469 6.421 78.9 4.9671    2 242    17.8 396.90  9.14
## 3 0.02729  0  7.07    0 0.469 7.185 61.1 4.9671    2 242    17.8 392.83  4.03
## 4 0.03237  0  2.18    0 0.458 6.998 45.8 6.0622    3 222    18.7 394.63  2.94
## 5 0.06905  0  2.18    0 0.458 7.147 54.2 6.0622    3 222    18.7 396.90  5.33
## 6 0.02985  0  2.18    0 0.458 6.430 58.7 6.0622    3 222    18.7 394.12  5.21
##   medv
## 1 24.0
## 2 21.6
## 3 34.7
## 4 33.4
## 5 36.2
## 6 28.7
```

```
Boston_model<-lm(medv~crim+zn+prratio+chas,data = BostonHousing)
summary(Boston_model)
```

```
##
## Call:
## lm(formula = medv ~ crim + zn + prratio + chas, data = BostonHousing)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -18.282  -4.505  -0.986   2.650  32.656
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  49.91868    3.23497   15.431 < 2e-16 ***
## crim        -0.26018    0.04015   -6.480 2.20e-10 ***
## zn           0.07073    0.01548    4.570 6.14e-06 ***
## prratio     -1.49367    0.17144   -8.712 < 2e-16 ***
## chas1        4.58393    1.31108    3.496 0.000514 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.388 on 501 degrees of freedom
## Multiple R-squared:  0.3599, Adjusted R-squared:  0.3547
## F-statistic: 70.41 on 4 and 501 DF,  p-value: < 2.2e-16
```

From the summary of the model it can be seen r squared value is 0.3588. This implies that accuracy of the model is low.

3. b) i.

When two houses are identical in all aspects then the one that bounds Charles River is more expensive by \$4583.9. This is because coefficient of Charles variable is 4.5839 and price is measured in 1000's of USD.

3. b) ii.

Coefficient of pupil-teacher ratio from the model is -1.493. This indicates that for every unit increase of pupil-teacher ratio the price of house decreases by \$1493. So, the price of house with pupil-teacher ratio 15, 18 is reduced by 22395 and 26874 USD respectively. Therefore the house with pupil_teacher ratio 15 is USD 4479 more expensive than the house with pupil_teacher ratio 18.

3. c)

The variables crime rate (crim), proportion of residential land zoned for lots over 25,000 sq.ft (zn), the local pupil-teacher ratio (prratio) and whether the tract bounds Charles River(chas) that are used in the model are all statistically important as their p-values lies between 0 and 0.001.

3. d)

```
anova(Boston_model)
```

```
## Analysis of Variance Table
##
## Response: medv
##           Df Sum Sq Mean Sq F value    Pr(>F)
## crim       1  6440.8   6440.8 118.007 < 2.2e-16 ***
## zn         1  3554.3   3554.3  65.122 5.253e-15 ***
## ptratio    1  4709.5   4709.5  86.287 < 2.2e-16 ***
## chas       1   667.2    667.2  12.224 0.0005137 ***
## Residuals 501 27344.5     54.6
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

From the anova analysis the order of importance of the four variables is as follows:

- 1.crime crate (crim)
- 2.the local pupil-teacher ratio (ptratio)
- 3.proportion of residential land zoned for lots over 25,000 sq.ft (zn)
- 4.the tract bounds Charles River(chas)