# HW5 ISYE 6501

*Kunle Lawal, Anubhav Rana, Mihir Tulpule, Ali Mujtaba Lakdawala*

## Question 8.1

*Describe a situation or problem from your job, everyday life, current events, etc., for which a linear regression model would be appropriate. List some (up to 5) predictors that you might use.*

In the future (after graduating from the MSA program) we might want to buy a house we would like to predict the price of a house given certain characteristics. To solve this problem a linear regression could work well.

We could use the following as predictors to determine the value of a house to us personally, with a numerical rating instead from 1-10 rather than a price. This model could be tailored to determine the estimated value to us as a user rather than the market price of the house.

Predictors that could be used in such a model are:

1. Age of house
2. Location
3. Number of schools in a 5 mile radius.
4. Median income of population
5. Number of bedrooms

## Question 8.2

*Using crime data from http://www.statsci.org/data/general/uscrime.txt (file uscrime.txt, description at http://www.statsci.org/data/general/uscrime.html ), use regression (a useful R function is lm or glm) to predict the observed crime rate in a city with the following data:*

M = 14.0
So = 0
Ed = 10.0
Po1 = 12.0
Po2 = 15.5
LF = 0.640
M.F = 94.0
Pop = 150
NW = 1.1
U1 = 0.120
U2 = 3.6
Wealth = 3200
Ineq = 20.1
Prob = 0.04
Time = 39.0

*Show your model (factors used and their coefficients), the software output, and the quality of fit.*

*Note that because there are only 47 data points and 15 predictors, you'll probably notice some overfitting. We'll see ways of dealing with this sort of problem later in the course.*

```
setwd("/Users/alimujtaba/Google Drive/isye6501modelling/isye6501homeworks/hw5")
require(data.table)
```

```
## Loading required package: data.table
```
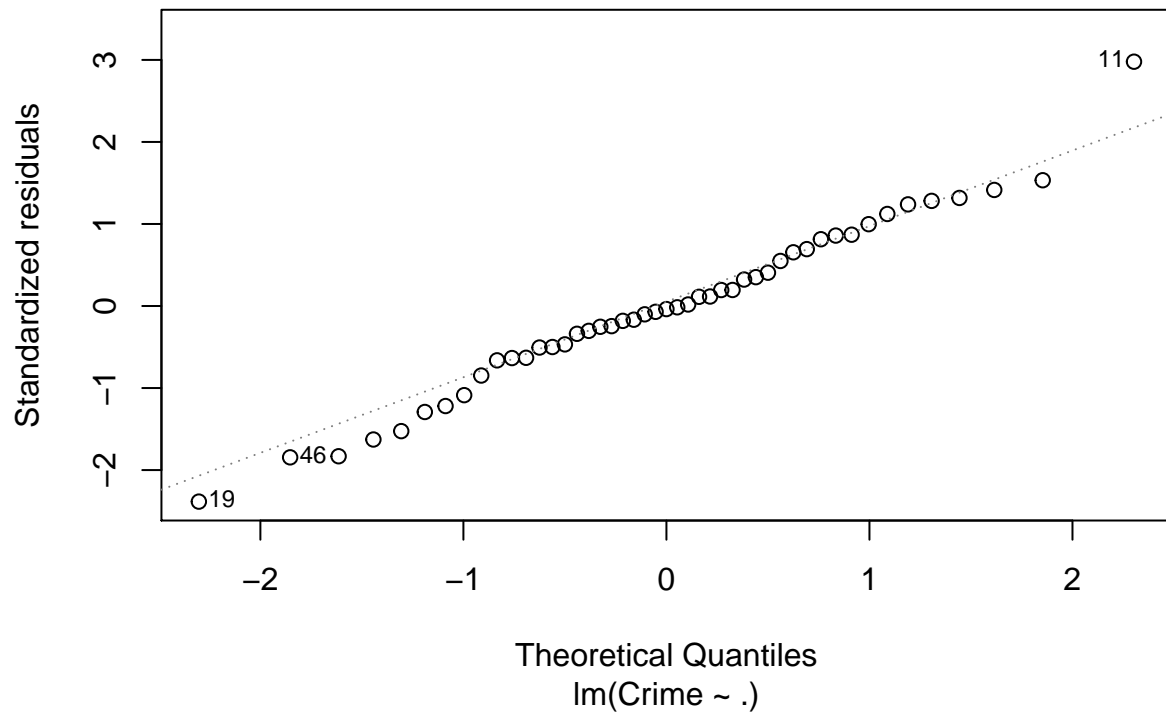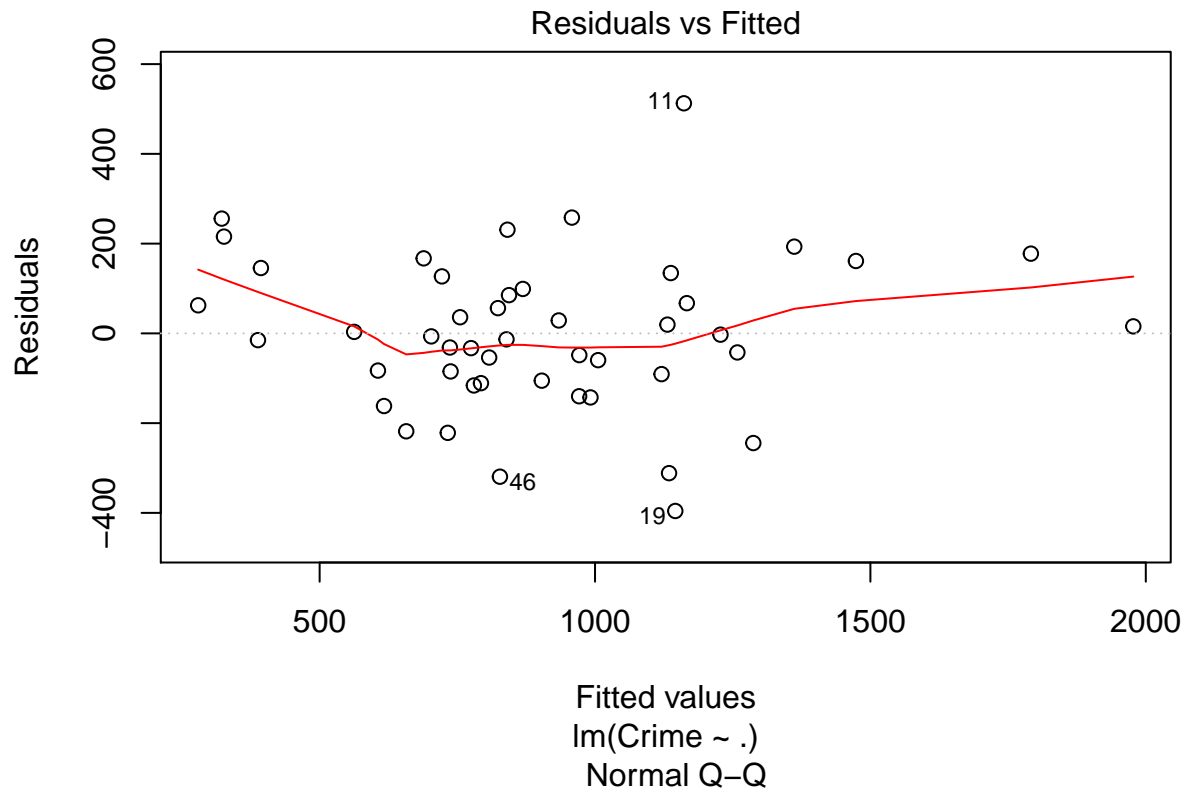
```r
crime_data <- read.table("uscrime.txt", header = TRUE)

crime_data
```
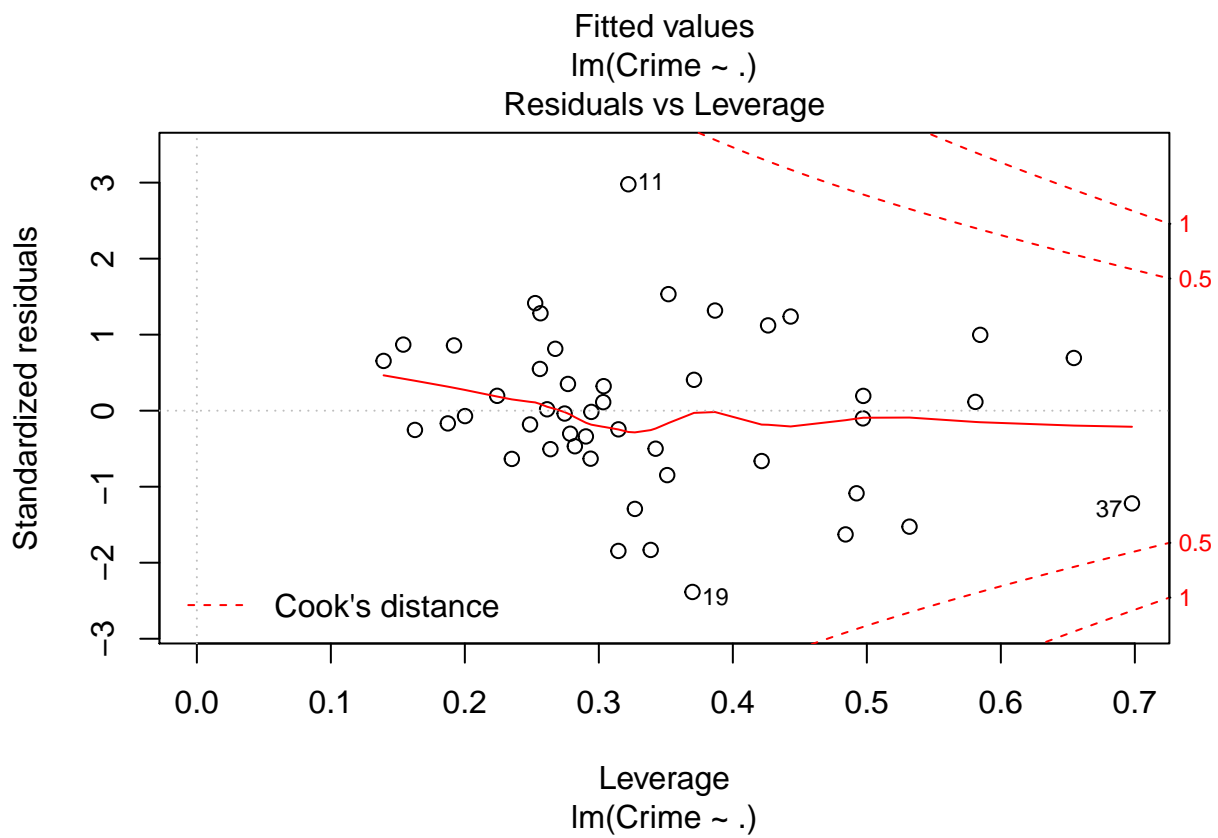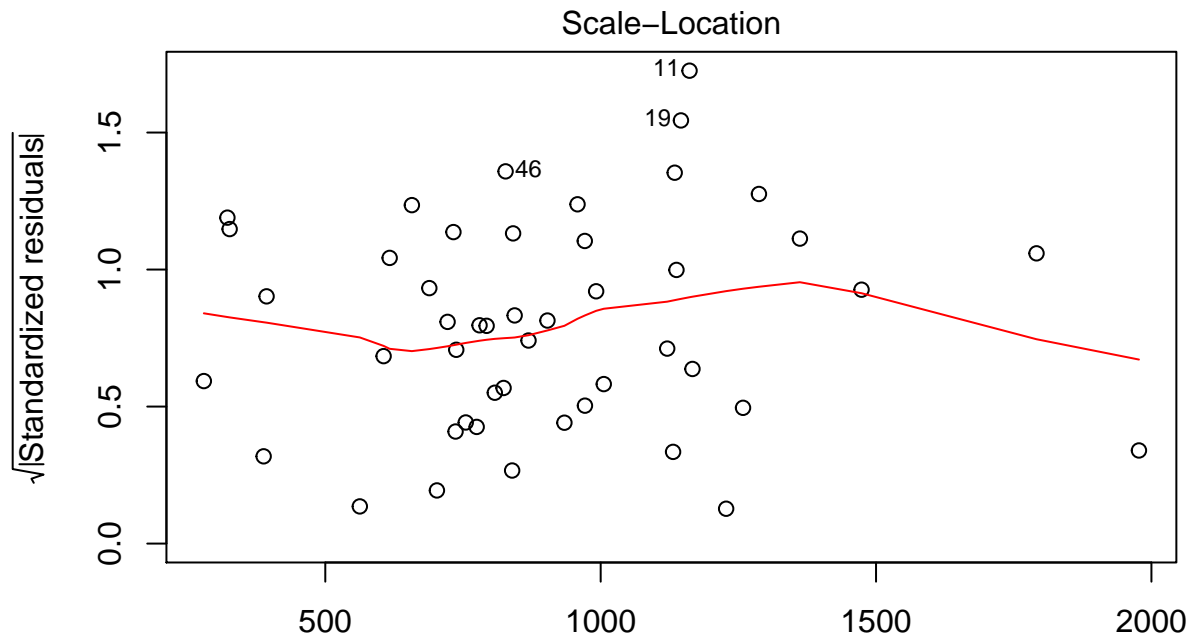
```
##        M So   Ed  Po1  Po2    LF   M.F Pop   NW    U1  U2 Wealth Ineq
## 1  15.1  1  9.1  5.8  5.6 0.510  95.0  33 30.1 0.108 4.1   3940 26.1
## 2  14.3  0 11.3 10.3  9.5 0.583 101.2  13 10.2 0.096 3.6   5570 19.4
## 3  14.2  1  8.9  4.5  4.4 0.533  96.9  18 21.9 0.094 3.3   3180 25.0
## 4  13.6  0 12.1 14.9 14.1 0.577  99.4 157  8.0 0.102 3.9   6730 16.7
## 5  14.1  0 12.1 10.9 10.1 0.591  98.5  18  3.0 0.091 2.0   5780 17.4
## 6  12.1  0 11.0 11.8 11.5 0.547  96.4  25  4.4 0.084 2.9   6890 12.6
## 7  12.7  1 11.1  8.2  7.9 0.519  98.2   4 13.9 0.097 3.8   6200 16.8
## 8  13.1  1 10.9 11.5 10.9 0.542  96.9  50 17.9 0.079 3.5   4720 20.6
## 9  15.7  1  9.0  6.5  6.2 0.553  95.5  39 28.6 0.081 2.8   4210 23.9
## 10 14.0  0 11.8  7.1  6.8 0.632 102.9   7  1.5 0.100 2.4   5260 17.4
## 11 12.4  0 10.5 12.1 11.6 0.580  96.6 101 10.6 0.077 3.5   6570 17.0
## 12 13.4  0 10.8  7.5  7.1 0.595  97.2  47  5.9 0.083 3.1   5800 17.2
## 13 12.8  0 11.3  6.7  6.0 0.624  97.2  28  1.0 0.077 2.5   5070 20.6
## 14 13.5  0 11.7  6.2  6.1 0.595  98.6  22  4.6 0.077 2.7   5290 19.0
## 15 15.2  1  8.7  5.7  5.3 0.530  98.6  30  7.2 0.092 4.3   4050 26.4
## 16 14.2  1  8.8  8.1  7.7 0.497  95.6  33 32.1 0.116 4.7   4270 24.7
## 17 14.3  0 11.0  6.6  6.3 0.537  97.7  10  0.6 0.114 3.5   4870 16.6
## 18 13.5  1 10.4 12.3 11.5 0.537  97.8  31 17.0 0.089 3.4   6310 16.5
## 19 13.0  0 11.6 12.8 12.8 0.536  93.4  51  2.4 0.078 3.4   6270 13.5
## 20 12.5  0 10.8 11.3 10.5 0.567  98.5  78  9.4 0.130 5.8   6260 16.6
## 21 12.6  0 10.8  7.4  6.7 0.602  98.4  34  1.2 0.102 3.3   5570 19.5
## 22 15.7  1  8.9  4.7  4.4 0.512  96.2  22 42.3 0.097 3.4   2880 27.6
## 23 13.2  0  9.6  8.7  8.3 0.564  95.3  43  9.2 0.083 3.2   5130 22.7
## 24 13.1  0 11.6  7.8  7.3 0.574 103.8   7  3.6 0.142 4.2   5400 17.6
## 25 13.0  0 11.6  6.3  5.7 0.641  98.4  14  2.6 0.070 2.1   4860 19.6
## 26 13.1  0 12.1 16.0 14.3 0.631 107.1   3  7.7 0.102 4.1   6740 15.2
## 27 13.5  0 10.9  6.9  7.1 0.540  96.5   6  0.4 0.080 2.2   5640 13.9
## 28 15.2  0 11.2  8.2  7.6 0.571 101.8  10  7.9 0.103 2.8   5370 21.5
## 29 11.9  0 10.7 16.6 15.7 0.521  93.8 168  8.9 0.092 3.6   6370 15.4
## 30 16.6  1  8.9  5.8  5.4 0.521  97.3  46 25.4 0.072 2.6   3960 23.7
## 31 14.0  0  9.3  5.5  5.4 0.535 104.5   6  2.0 0.135 4.0   4530 20.0
## 32 12.5  0 10.9  9.0  8.1 0.586  96.4  97  8.2 0.105 4.3   6170 16.3
## 33 14.7  1 10.4  6.3  6.4 0.560  97.2  23  9.5 0.076 2.4   4620 23.3
## 34 12.6  0 11.8  9.7  9.7 0.542  99.0  18  2.1 0.102 3.5   5890 16.6
## 35 12.3  0 10.2  9.7  8.7 0.526  94.8 113  7.6 0.124 5.0   5720 15.8
## 36 15.0  0 10.0 10.9  9.8 0.531  96.4   9  2.4 0.087 3.8   5590 15.3
## 37 17.7  1  8.7  5.8  5.6 0.638  97.4  24 34.9 0.076 2.8   3820 25.4
## 38 13.3  0 10.4  5.1  4.7 0.599 102.4   7  4.0 0.099 2.7   4250 22.5
## 39 14.9  1  8.8  6.1  5.4 0.515  95.3  36 16.5 0.086 3.5   3950 25.1
## 40 14.5  1 10.4  8.2  7.4 0.560  98.1  96 12.6 0.088 3.1   4880 22.8
## 41 14.8  0 12.2  7.2  6.6 0.601  99.8   9  1.9 0.084 2.0   5900 14.4
## 42 14.1  0 10.9  5.6  5.4 0.523  96.8   4  0.2 0.107 3.7   4890 17.0
## 43 16.2  1  9.9  7.5  7.0 0.522  99.6  40 20.8 0.073 2.7   4960 22.4
## 44 13.6  0 12.1  9.5  9.6 0.574 101.2  29  3.6 0.111 3.7   6220 16.2
## 45 13.9  1  8.8  4.6  4.1 0.480  96.8  19  4.9 0.135 5.3   4570 24.9
## 46 12.6  0 10.4 10.6  9.7 0.599  98.9  40  2.4 0.078 2.5   5930 17.1
## 47 13.0  0 12.1  9.0  9.1 0.623 104.9   3  2.2 0.113 4.0   5880 16.0
##        Prob    Time Crime
## 1  0.084602 26.2011   791
```

```
## 2  0.029599 25.2999  1635
## 3  0.083401 24.3006   578
## 4  0.015801 29.9012  1969
## 5  0.041399 21.2998  1234
## 6  0.034201 20.9995   682
## 7  0.042100 20.6993   963
## 8  0.040099 24.5988  1555
## 9  0.071697 29.4001   856
## 10 0.044498 19.5994   705
## 11 0.016201 41.6000  1674
## 12 0.031201 34.2984   849
## 13 0.045302 36.2993   511
## 14 0.053200 21.5010   664
## 15 0.069100 22.7008   798
## 16 0.052099 26.0991   946
## 17 0.076299 19.1002   539
## 18 0.119804 18.1996   929
## 19 0.019099 24.9008   750
## 20 0.034801 26.4010  1225
## 21 0.022800 37.5998   742
## 22 0.089502 37.0994   439
## 23 0.030700 25.1989  1216
## 24 0.041598 17.6000   968
## 25 0.069197 21.9003   523
## 26 0.041698 22.1005  1993
## 27 0.036099 28.4999   342
## 28 0.038201 25.8006  1216
## 29 0.023400 36.7009  1043
## 30 0.075298 28.3011   696
## 31 0.041999 21.7998   373
## 32 0.042698 30.9014   754
## 33 0.049499 25.5005  1072
## 34 0.040799 21.6997   923
## 35 0.020700 37.4011   653
## 36 0.006900 44.0004  1272
## 37 0.045198 31.6995   831
## 38 0.053998 16.6999   566
## 39 0.047099 27.3004   826
## 40 0.038801 29.3004  1151
## 41 0.025100 30.0001   880
## 42 0.088904 12.1996   542
## 43 0.054902 31.9989   823
## 44 0.028100 30.0001  1030
## 45 0.056202 32.5996   455
## 46 0.046598 16.6999   508
## 47 0.052802 16.0997   849
```

```r
model <- lm(Crime ~ ., data = crime_data)
# Plotting the model and the summary
plot(model)
```

Residuals vs Fitted

lm(Crime ~ .)

Normal Q–Q

Theoretical Quantiles
lm(Crime ~ .)

Scale–Location
lm(Crime ~ .)



Residuals vs Leverage
lm(Crime ~ .)

```r
summary(model)
```

```
##
## Call:
## lm(formula = Crime ~ ., data = crime_data)
##
```

```
## Residuals:
##     Min      1Q  Median      3Q     Max
## -395.74  -98.09   -6.69  112.99  512.67
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept) -5.984e+03  1.628e+03  -3.675 0.000893 ***
## M            8.783e+01  4.171e+01   2.106 0.043443 *
## So          -3.803e+00  1.488e+02  -0.026 0.979765
## Ed           1.883e+02  6.209e+01   3.033 0.004861 **
## Po1          1.928e+02  1.061e+02   1.817 0.078892 .
## Po2         -1.094e+02  1.175e+02  -0.931 0.358830
## LF          -6.638e+02  1.470e+03  -0.452 0.654654
## M.F          1.741e+01  2.035e+01   0.855 0.398995
## Pop         -7.330e-01  1.290e+00  -0.568 0.573845
## NW           4.204e+00  6.481e+00   0.649 0.521279
## U1          -5.827e+03  4.210e+03  -1.384 0.176238
## U2           1.678e+02  8.234e+01   2.038 0.050161 .
## Wealth       9.617e-02  1.037e-01   0.928 0.360754
## Ineq         7.067e+01  2.272e+01   3.111 0.003983 **
## Prob        -4.855e+03  2.272e+03  -2.137 0.040627 *
## Time        -3.479e+00  7.165e+00  -0.486 0.630708
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 209.1 on 31 degrees of freedom
## Multiple R-squared:  0.8031, Adjusted R-squared:  0.7078
## F-statistic: 8.429 on 15 and 31 DF,  p-value: 3.539e-07
```

The model quality is excellent shown by the multiple R-squared value of 0.8031 which shows that the model is very strong at predicting crime. The most significant predictors are:
1. M (percentage of males (14-24) in population 2. Ed = Mean years of schooling for those above 25 3. Ineq = Income inequality 4. Prob = Probability of imprisonment.

These conclusions are on a benchmark of alpha = 0.05. This does not mean that the rest of the predictors are insignificant. It just means that in the presence of the above predictors they do not add significant value to the model.

The model plots also show that there are very few outliers. QQ plot shows normality is good. The residuals vs. fitted plot shows that there is good variance in the model. No Heteroskedasticity.

```r
# Getting a sense of the distribution of the data points.
summary(crime_data)
```

```
##        M               So               Ed              Po1
##  Min.   :11.90   Min.   :0.0000   Min.   : 8.70   Min.   : 4.50
##  1st Qu.:13.00   1st Qu.:0.0000   1st Qu.: 9.75   1st Qu.: 6.25
##  Median :13.60   Median :0.0000   Median :10.80   Median : 7.80
##  Mean   :13.86   Mean   :0.3404   Mean   :10.56   Mean   : 8.50
##  3rd Qu.:14.60   3rd Qu.:1.0000   3rd Qu.:11.45   3rd Qu.:10.45
##  Max.   :17.70   Max.   :1.0000   Max.   :12.20   Max.   :16.60
##       Po2              LF              M.F             Pop
##  Min.   : 4.100   Min.   :0.4800   Min.   : 93.40   Min.   : 3.00
##  1st Qu.: 5.850   1st Qu.:0.5305   1st Qu.: 96.45   1st Qu.: 10.00
##  Median : 7.300   Median :0.5600   Median : 97.70   Median : 25.00
##  Mean   : 8.023   Mean   :0.5612   Mean   : 98.30   Mean   : 36.62
```

```
## 3rd Qu.: 9.700    3rd Qu.:0.5930    3rd Qu.: 99.20    3rd Qu.: 41.50
## Max.   :15.700    Max.   :0.6410    Max.   :107.10    Max.   :168.00
##        NW                U1                U2             Wealth
## Min.   : 0.20    Min.   :0.07000    Min.   :2.000    Min.   :2880
## 1st Qu.: 2.40    1st Qu.:0.08050    1st Qu.:2.750    1st Qu.:4595
## Median : 7.60    Median :0.09200    Median :3.400    Median :5370
## Mean   :10.11    Mean   :0.09547    Mean   :3.398    Mean   :5254
## 3rd Qu.:13.25    3rd Qu.:0.10400    3rd Qu.:3.850    3rd Qu.:5915
## Max.   :42.30    Max.   :0.14200    Max.   :5.800    Max.   :6890
##        Ineq              Prob              Time             Crime
## Min.   :12.60    Min.   :0.00690    Min.   :12.20    Min.   : 342.0
## 1st Qu.:16.55    1st Qu.:0.03270    1st Qu.:21.60    1st Qu.: 658.5
## Median :17.60    Median :0.04210    Median :25.80    Median : 831.0
## Mean   :19.40    Mean   :0.04709    Mean   :26.60    Mean   : 905.1
## 3rd Qu.:22.75    3rd Qu.:0.05445    3rd Qu.:30.45    3rd Qu.:1057.5
## Max.   :27.60    Max.   :0.11980    Max.   :44.00    Max.   :1993.0
```

```r
new_data_point = data.frame(M = 14.0, So = 0,  Ed = 10.0,
Po1 = 12.0,
Po2 = 15.5,
LF = 0.640,
M.F = 94.0,
Pop = 150,
NW = 1.1,
U1 = 0.120,
U2 = 3.6,
Wealth = 3200,
Ineq = 20.1,
Prob = 0.04,
Time = 39.0)

predict <- predict.lm(model, new_data_point, interval = "prediction")
predict
```

```
##        fit       lwr      upr
## 1 155.4349 -1370.845 1681.715
```

The final prediction is below the min of crime and has a very large prediction interval. This is caused by the fact that a lot of the inputs tend towards the min and max of their respective ranges making it hard to be confident about the prediction.


**Additonal Analysis:**

Using the most significant predictors from our initial model we created a new linear model with the predcitors Mo, Ed, Prob and Ineq.

```r
model2 <- lm(Crime ~ M + Ed + Prob + Ineq, data = crime_data)
# Plotting the model and the summary
#plot(model)

summary(model2)
```

```
##
## Call:
```

```
## lm(formula = Crime ~ M + Ed + Prob + Ineq, data = crime_data)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -532.97 -254.03  -55.72  137.80  960.21
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1339.35    1247.01  -1.074  0.28893
## M              35.97      53.39   0.674  0.50417
## Ed            148.61      71.92   2.066  0.04499 *
## Prob        -7331.92    2560.27  -2.864  0.00651 **
## Ineq           26.87      22.77   1.180  0.24458
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 347.5 on 42 degrees of freedom
## Multiple R-squared:  0.2629, Adjusted R-squared:  0.1927
## F-statistic: 3.745 on 4 and 42 DF,  p-value: 0.01077
```

```
predict2 <- predict.lm(model2, new_data_point, interval = "prediction")
predict2
```

```
##        fit      lwr      upr
## 1 897.2307 184.0633 1610.398
```

While those are the most significant predictors, by themselves they result in a poor model also resulting in a widely different prediction. From this we can see that all the other predictors while not significant contribute additional information to the construction of a good model.