

ECON-UB 251 Econometrics I Assignment 1

Kevin Song

Theory

1.1

- Derive the OLS estimator of β_1

1.1 $Y_i = \beta_1 X_i + u_i$

D) The OLS regression can be modeled as the sum of squared deviations from the regression line:

$$TSS(b_0, b_1) = \sum_{i=1}^n [Y_i - b_0 - b_1 X_i]^2$$

Since our model has no b_0 , we only have to take the partial derivative of b_1 .

$$\frac{d TSS(b_0, b_1)}{d b_1} = -2 \sum_{i=1}^n X_i (Y_i - b_1 X_i)$$
$$= -2 \left(\sum_{i=1}^n X_i Y_i - b_1 \sum_{i=1}^n X_i^2 \right) = 0$$

3)

$$\sum_{i=1}^n X_i Y_i - b_1 \sum_{i=1}^n X_i^2 = 0$$
$$b_1 \sum_{i=1}^n X_i^2 = \sum_{i=1}^n X_i Y_i$$
$$b_1 = \frac{\sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n X_i^2}$$

1.1b This is called regression through the origin because without β_0 , the regression will go through $(0,0)$. This can be problematic if the data does not actually fit well with this forced fit.

1.1c Yes, imagine this graph:

The regression through the origin can not get an accurate β_1 without an intercept β_0 .

1.2

- What happens to the OLS estimators if we multiply the independent or dependent variables by 100?

1.2

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \cdot \bar{x}$$

If we multiply the dependent by 100: $100y_i$, then
 $\hat{\beta}_1$ will multiply by 100 because $(y_i - \bar{y}) \cdot 100$ in the numerator.
 $\hat{\beta}_0$ will also multiply by 100 because $100\bar{y} - 100\hat{\beta}_1 \cdot \bar{x}$
equals $100(\bar{y} - \hat{\beta}_1 \bar{x}) = \hat{\beta}_0$

If we multiply the independent by 100: $100x_i$, then
 $\hat{\beta}_1$ will divide by 100 because the $(x_i - \bar{x})$ in the denom
is squared, so $100^2 = 10000$. The 100 in the numerator
offsets this partially: $\frac{100}{10000} = \frac{1}{100}\hat{\beta}_1$.
 $\hat{\beta}_0$ will not change because $\bar{y} - \frac{\hat{\beta}_1}{100} \cdot 100\bar{x}$
equals $\bar{y} - \hat{\beta}_1 \bar{x}$ still.

$$1.2 \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\begin{aligned} y_i &= 100y_i' \\ 100 &\in 100y_i' \end{aligned}$$

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i' - \bar{y}') \cdot 100}{\sum_{i=1}^n (x_i - \bar{x})^2} = [100\hat{\beta}_1]$$

$$\begin{aligned} \hat{\beta}_0 &= 100\bar{y}' - 100\hat{\beta}_1 \cdot \bar{x} \\ &= 100(\bar{y}' - \hat{\beta}_1 \cdot \bar{x}) = [100\hat{\beta}_0] \end{aligned}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (100x_i - 100\bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (100x_i - 100\bar{x})^2}$$

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{100^2 \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{100 \sum_{i=1}^n (x_i - \bar{x})^2} = \boxed{\frac{1}{100}\hat{\beta}_1}$$

$$\begin{aligned} \hat{\beta}_0 &= \bar{y} - \frac{1}{100}\hat{\beta}_1 \cdot 100\bar{x} \\ &= \bar{y} - \hat{\beta}_1 \bar{x} = \boxed{\hat{\beta}_0} \end{aligned}$$

1.3

3. What happens to the OLS estimators if we standardize the variables?

1.3 When β_0 and β_1 become standardized so that

\bar{x} and $\bar{y} = 0$, and σ_x and σ_y become 1:

$$\hat{\beta}_1 = \text{corr}(x, y) \frac{\sigma_y}{\sigma_x} = \text{corr}(x, y) \cdot \frac{1}{1}$$

$$\hat{\beta}_1 = \text{corr}(x, y)$$

$$\hat{\beta}_0 = \bar{y}^* - \hat{\beta}_1 \bar{x}^* = 0 - \hat{\beta}_1 \cdot 0 = 0$$

The intercept becomes 0 so the regression will go through the origin and the slope simply becomes the correlation of x and y .

2.1

```
library(readr)
sample_orig_2012 <- read_delim("/cloud/project/sample_orig_2012.txt", "|", escape_double = FALSE, col_na

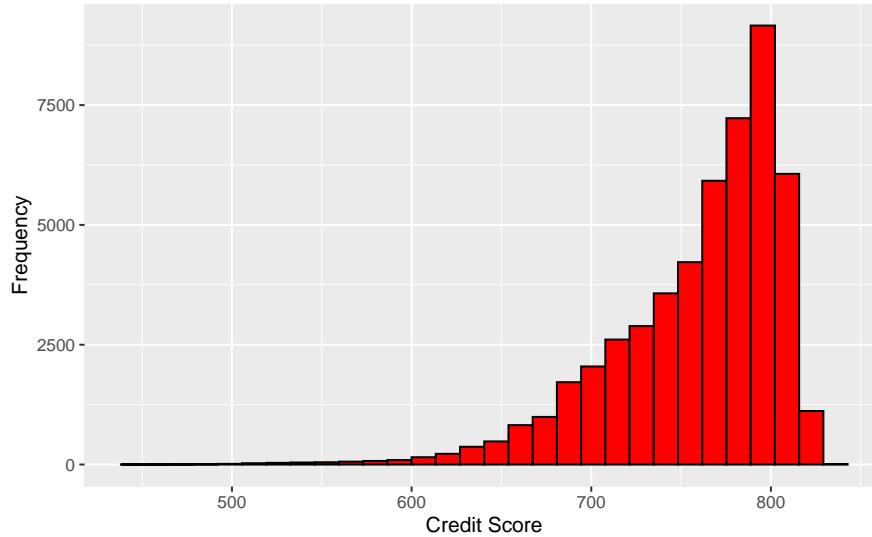
orig.names <- c("Credit_Score", "First_Payment_ate", "First_Time_Homebuyer", "Maturity_Date",
"MSA", "Mortgage_Insurance_Percentage", "Number_Units", "Occupancy_Status", "CLTV",
"DTI", "UPB", "LTV", "Interest_Rate", "Channel", "Prepayment_Penalty",
"Amortization_Type", "State", "Property_Type", "Postal_Code", "Sequence_Number",
"Purpose", "Loan_Term", "Number_Borrowers", "Seller_Name", "Servicer_Name",
"Super_Conforming", "Pre-HARP_Loan", "Program_Indicator", "HARP_Indicator",
"Valuation_Method", "Interest_Only")
colnames(sample_orig_2012) <- orig.names
```

3.1

```
library(ggplot2)
library(dplyr)
Credit_Score <- sample_orig_2012$Credit_Score
Interest_Rate <- sample_orig_2012$Interest_Rate
sample_orig_2012 = filter(sample_orig_2012, sample_orig_2012$Credit_Score != 9999)
ggplot(sample_orig_2012, aes(x = Credit_Score)) +
  geom_histogram(fill = "red", color = "black") +
  labs(title = "Distribution of Credit Score",
       subtitle = "2012 US",
```

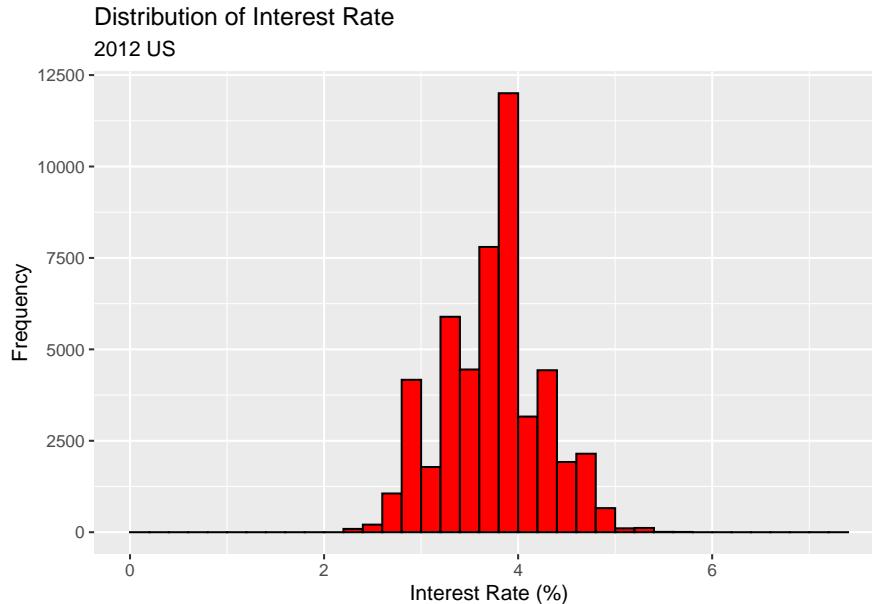
```
x = "Credit Score",
y = "Frequency")
```

Distribution of Credit Score
2012 US



There are no outliers, nor missing values, but the data does skew to the left, with data points gradually going farther and farther left from the mode.

```
library(ggplot2)
ggplot(sample_orig_2012, aes(x = Interest_Rate)) +
  geom_histogram(breaks = seq(0, 7.5, 0.2), fill = "red", color = "black") +
  labs(title = "Distribution of Interest Rate",
       subtitle = "2012 US",
       x = "Interest Rate (%)",
       y = "Frequency")
```



3.2

```
# Calculate Mean  
mean(Credit_Score)
```

```
[1] 759.1127
```

The mean is not within the modal bar. This is because of the left skew lowering the mean compared to the median and mode.

```
# Calculate Standard Deviation  
sd(Credit_Score)
```

```
[1] 47.31753
```

Credit Score's standard deviation is about 47 units, which is quite a large spread. This number is large also in part due to the skew of the graph.

```
# Calculate Skewness  
library(moments)  
skewness(Credit_Score)
```

```
[1] -1.364729
```

The skewness is smaller than -1, meaning there is a high left skew, which affects the mean and std.

```
# Calculate Kurtosis  
kurtosis(Credit_Score)
```

```
[1] 5.429348
```

Kurtosis is much greater than 3, meaning that there are high amounts of data in the tails. We can observe this in how there are data points below 500 in a graph with mean = 759 and std = 47. This is more than 5 standard deviations which in a normal distribution would be practically nonexistent.

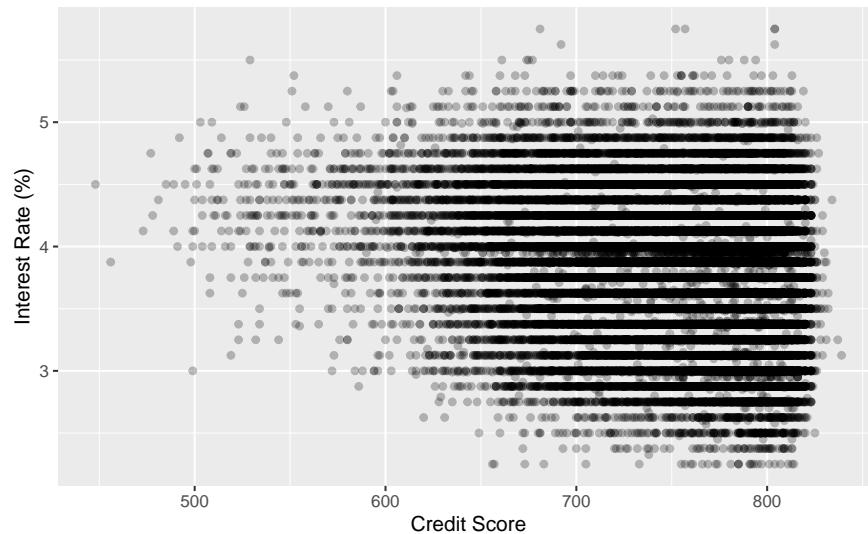
3.3

Scatter plot of the credit score and interest rate variables:

```
library(ggplot2)  
ggplot(sample_orig_2012, aes(x = sample_orig_2012$Credit_Score, y = sample_orig_2012$Interest_Rate)) +  
  geom_point(color = "black", alpha = 0.25) +  
  labs(title = "Scatterplot of Credit Score and Interest Rate",  
       subtitle = "2012 US",  
       x = "Credit Score",  
       y = "Interest Rate (%)")
```

Scatterplot of Credit Score and Interest Rate

2012 US



In the scatterplot, it does not seem like interest rate has much correlation with credit score. The regression appears to be slightly negative, but the two variables are probably independent from each other.

4.1

$$\hat{\beta}_1 = \frac{\frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2}$$

$$\hat{\beta}_1 = -2.338e-03$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$$\hat{\beta}_2 = 5.532$$

4.2

Because $\hat{\beta}_1 = -2.338e-03$, there is a negative linear relationship between Credit Scores and Interest rate; as Credit Score increases by 1 the Interest rate decreases by $2.338e-03$. The sign is consistent, because higher credit score means that person is more consistently paying their card, meaning lower risk for the company and thus lower interest rates.

4.3

```
regression <- lm(Interest_Rate ~ Credit_Score)
summary(regression)
```

Call:

```
lm(formula = Interest_Rate ~ Credit_Score)
```

Residuals:

Min	1Q	Median	3Q	Max
-1.74798	-0.33196	0.04071	0.31536	2.09810

```

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 5.532e+00 3.587e-02 154.20 <2e-16 ***
Credit_Score -2.338e-03 4.717e-05 -49.58 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Residual standard error: 0.499 on 49998 degrees of freedom
 Multiple R-squared: 0.04685, Adjusted R-squared: 0.04684
 F-statistic: 2458 on 1 and 49998 DF, p-value: < 2.2e-16

- R^2 is the correlation of determination, 0.04685 suggests very little of the variance in interest rate is explained by this credit score in this model. This supports the null hypothesis that interest rate is independent from credit score.
- SER is the average distance that the observed values fall from the regression. SER = 0.499 with a small 4.6% R^2 means that the error of the residuals are very large.