

# Linear Regression

The following data\* represents the number of units (Units) being replaced in a computer and the corresponding time taken to repair the computer in minutes (Minutes)

ITTP 39km Railer | 6 AM - 6 PM

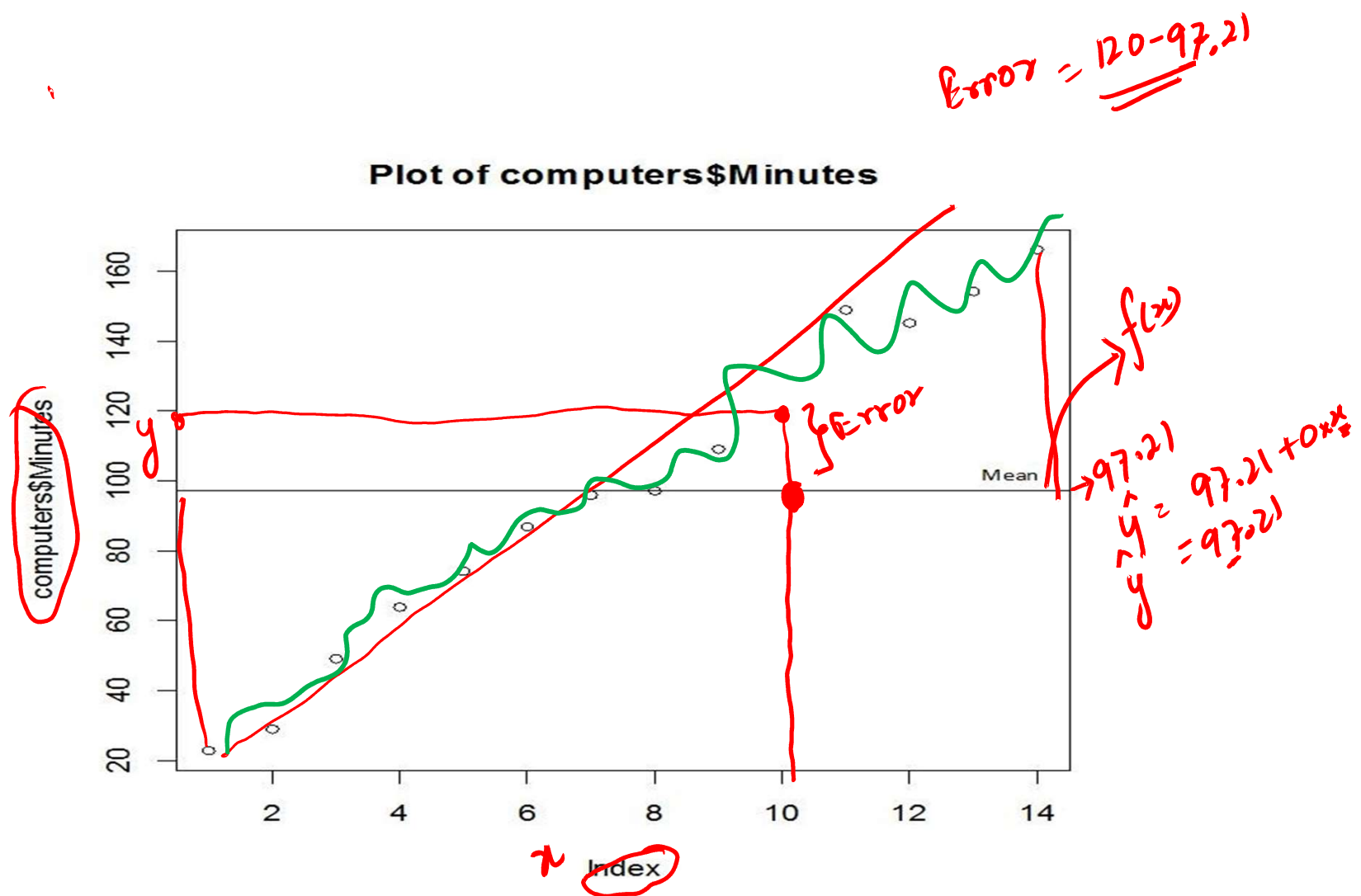
| computers | Units | Minutes |
|-----------|-------|---------|
| 1         | 1     | 23      |
| 2         | 2     | 29      |
| 3         | 3     | 49      |
| 4         | 4     | 64      |
| 5         | 4     | 74      |
| 6         | 5     | 87      |
| 7         | 6     | 96      |
| 8         | 6     | 97      |
| 9         | 7     | 109     |
| 10        | 8     | 119     |
| 11        | 9     | 149     |
| 12        | 9     | 145     |
| 13        | 10    | 154     |
| 14        | 10    | 166     |

If we were to estimate the typical time taken by the computer shop to repair a computer, which of the following would be appropriate?

- a) Arithmetic mean = 97.21 minutes
- b) Median = 96.50 minutes

Answer: Either Mean or Median can be used when trying to predict the expected value of a single variable.

$f(x) = y = \hat{y} = \beta_0 + \beta_1 x$   
 output  $\rightarrow$  Predicted output



[Data: 23 29 49 64 74 87 96 97 109 119 149 145 154 166, Mean: 97.21]

# Regression Analysis

Regression analysis is a statistical process for estimating the relationships among variables. It includes many techniques for modeling and analyzing several variables, when the focus is on the relationship between a dependent variable and one or more independent variables (or 'predictors').

Regression analysis focuses on building a model that can be used to predict the value of the dependent variable based on the predictor variables.

The regression model is represented using a mathematical model of form  $y = f(X)$ , where  $y$  is the dependent variable and  $X$  is the set of predictor variables ( $x_1, x_2 \dots x_n$ ).

# Regression Analysis

*dependent* ↘

↙ features → Parameters  
Independent Variable

$$f(x) = \beta_0 + \beta_1 x + \epsilon$$

↑  
Error

Linear form:  $f(X) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$  ↙

Non-linear form:  $f(X) = \beta_0 + \beta_1 x_1^{p_1} + \beta_2 x_2^{p_2} + \dots + \beta_n x_n^{p_n} + \epsilon$  ↙

# Regression Analysis

**Some commonly used types of linear forms are:**

Simple linear form – Here there is one predictor and one dependent variable :  $f(X) =$

$$\beta_0 + \beta_1 x_1 + \epsilon$$

Multiple linear form – Here there are multiple predictor variables and one dependent variable:

$$f(X) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

**Some commonly used types of non-linear forms are:**

Polynomial form:  $f(X) = \beta_0 + \beta_1 x_1^{p_1} + \beta_2 x_2^{p_2} + \dots + \beta_n x_n^{p_n} + \epsilon$  ✓

Quadratic form:  $f(X) = \beta_0 + \beta_1 x_1^2 + \beta_2 x_1 + \epsilon$  ✓

Logistic form:  $f(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}} + \epsilon$  ✓  
↑ classification

Where  $\beta_0, \beta_1, \beta_2 \dots \beta_n$  are said to be the regression coefficients and  $\epsilon$  accounts for the error in prediction. The regression coefficients and the error in prediction are real numbers.

# Regression Analysis

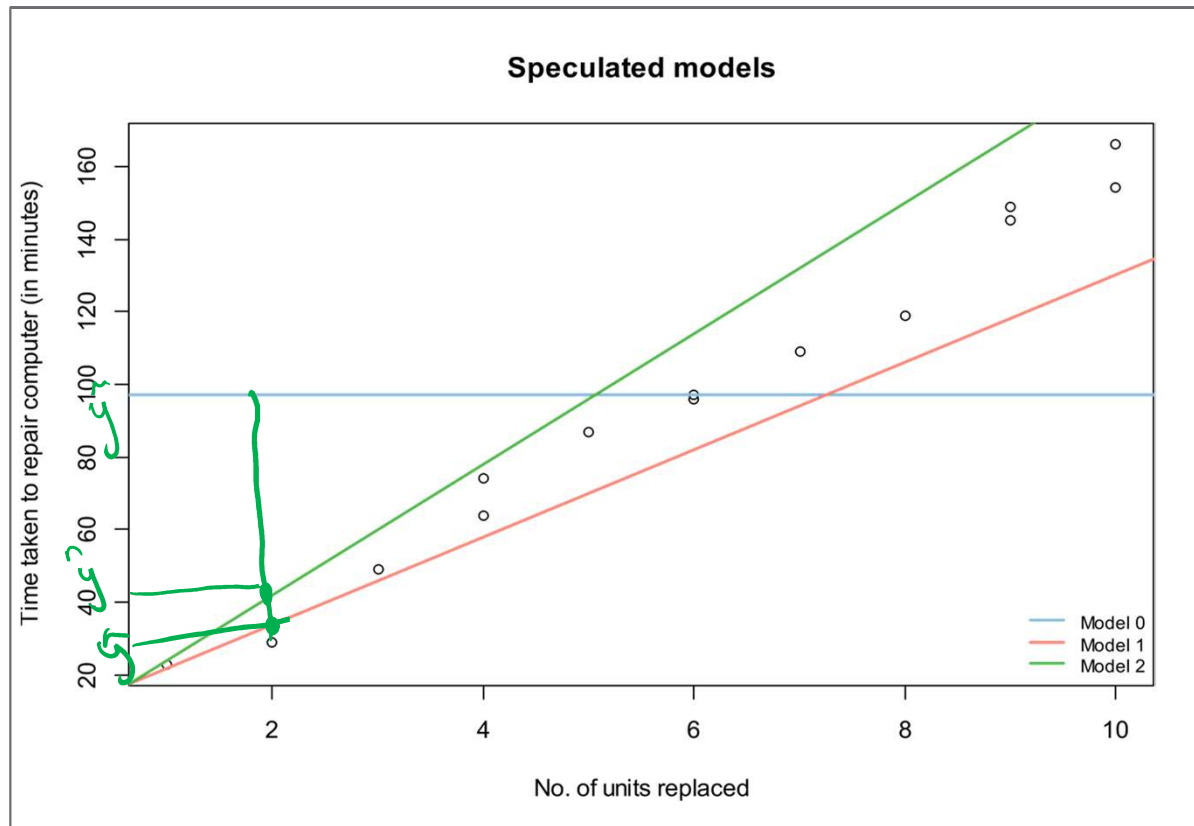
Time taken to repair the computer =  $\beta_0 + \beta_1$  \* No. of units replaced

Model 0: Time taken to repair the computer = 97.21 (*mean*)

Model 1: Time taken to repair the computer = 10 + 12 \* No. of units replaced

Model 2: Time taken to repair the computer = 6 + 18 \* No. of units replaced

# Regression Analysis





# Regression Analysis

model0

|    | Units replaced | Observed time taken | Expected value | Expected - observed value | Error<br>( $\hat{y} - y$ ) |
|----|----------------|---------------------|----------------|---------------------------|----------------------------|
| 7  |                |                     |                |                           |                            |
| 8  |                |                     |                |                           |                            |
| 9  | 1              | 23                  | 97.21429       | 74.2142857                |                            |
| 10 | 2              | 29                  | 97.21429       | 68.2142857                |                            |
| 11 | 3              | 49                  | 97.21429       | 48.2142857                |                            |
| 12 | 4              | 64                  | 97.21429       | 33.2142857                |                            |
| 13 | 4              | 74                  | 97.21429       | 23.2142857                |                            |
| 14 | 5              | 87                  | 97.21429       | 10.2142857                |                            |
| 15 | 6              | 96                  | 97.21429       | 1.2142857                 |                            |
| 16 | 6              | 97                  | 97.21429       | 0.2142857                 |                            |
| 17 | 7              | 109                 | 97.21429       | -11.7857143               |                            |
| 18 | 8              | 119                 | 97.21429       | -21.7857143               |                            |
| 19 | 9              | 149                 | 97.21429       | -51.7857143               |                            |
| 20 | 9              | 145                 | 97.21429       | -47.7857143               |                            |
| 21 | 10             | 154                 | 97.21429       | -56.7857143               |                            |
| 22 | 10             | 166                 | 97.21429       | -68.7857143               |                            |

$$\begin{aligned} \text{Error} &= \text{Sum}(\text{minutes mean} - \text{computer\$minutes}) \\ &= -8.25e-14 \end{aligned}$$

$$\begin{aligned} \text{Error} &= \text{Sum}((\text{minutes mean} - \text{computer\$minutes})^2) \\ &= 27768.36 \end{aligned}$$

Sum of Squared Errors

# Regression Analysis

model1

| Units replaced | Observed time | Expected time | Expected - observed value |
|----------------|---------------|---------------|---------------------------|
| 1              | 23            | 22            | -1                        |
| 2              | 29            | 34            | 5                         |
| 3              | 49            | 46            | -3                        |
| 4              | 64            | 58            | -6                        |
| 4              | 74            | 58            | -16                       |
| 5              | 87            | 70            | -17                       |
| 6              | 96            | 82            | -14                       |
| 6              | 97            | 82            | -15                       |
| 7              | 109           | 94            | -15                       |
| 8              | 119           | 106           | -13                       |
| 9              | 149           | 118           | -31                       |
| 9              | 145           | 118           | -27                       |
| 10             | 154           | 130           | -24                       |
| 10             | 166           | 130           | -36                       |

$$\text{Error} = \text{Sum}((\text{minutesmean} - \text{computer\$minutes})^2)$$

= 4993

# Regression Analysis

model2

| Units replaced | observed time | Expected time | Expected - observed value |
|----------------|---------------|---------------|---------------------------|
| 1              | 23            | 24            | 1                         |
| 2              | 29            | 42            | 13                        |
| 3              | 49            | 60            | 11                        |
| 4              | 64            | 78            | 14                        |
| 4              | 74            | 78            | 4                         |
| 5              | 87            | 96            | 9                         |
| 6              | 96            | 114           | 18                        |
| 6              | 97            | 114           | 17                        |
| 7              | 109           | 132           | 23                        |
| 8              | 119           | 150           | 31                        |
| 9              | 149           | 168           | 19                        |
| 9              | 145           | 168           | 23                        |
| 10             | 154           | 186           | 32                        |
| 10             | 166           | 186           | 20                        |

$\epsilon$

$$\begin{aligned}\text{Error} &= \text{Sum}((\text{minutesmean} - \text{computer\$minutes})^2) \\ &= 5001\end{aligned}$$

# Regression Analysis

The best fit model is obtained by solving the linear regression model  $\hat{y} = \beta_0 + \beta_1 x + \varepsilon$

To determine the coefficients  $\beta_0$  and  $\beta_1$  such that the error (as shown below) is minimum:

$$\text{Error} = \sum (\hat{y} - y)^2$$

$$\text{Error} = \sum (\beta_0 + \beta_1 x - y)^2$$

$$\frac{\partial \text{Error}}{\partial \beta_0} =$$

$$\frac{\partial \text{Error}}{\partial \beta_1} =$$

Differentiating the equation  $\sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$  w.r.t  $b_0$  and equating to 0 we get

no. of Rows  $\rightarrow$   $b_1 = \frac{[ \sum_{i=1}^n (x_i y_i) ] - n \bar{x} \bar{y}}{[ \sum_{i=1}^n x_i^2 ] - n \bar{x}^2}$   $\leftarrow$  mean bar

Similarly, differentiating the equation  $\sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$  w.r.t  $b_1$  and equating to 0 we get

$$b_0 = \bar{y} - b_1 \bar{x}$$

Here,  $b_0$  and  $b_1$  are the regression coefficients for the sample,  $\bar{x}$  and  $\bar{y}$  are the sample means of the variables X and Y respectively.

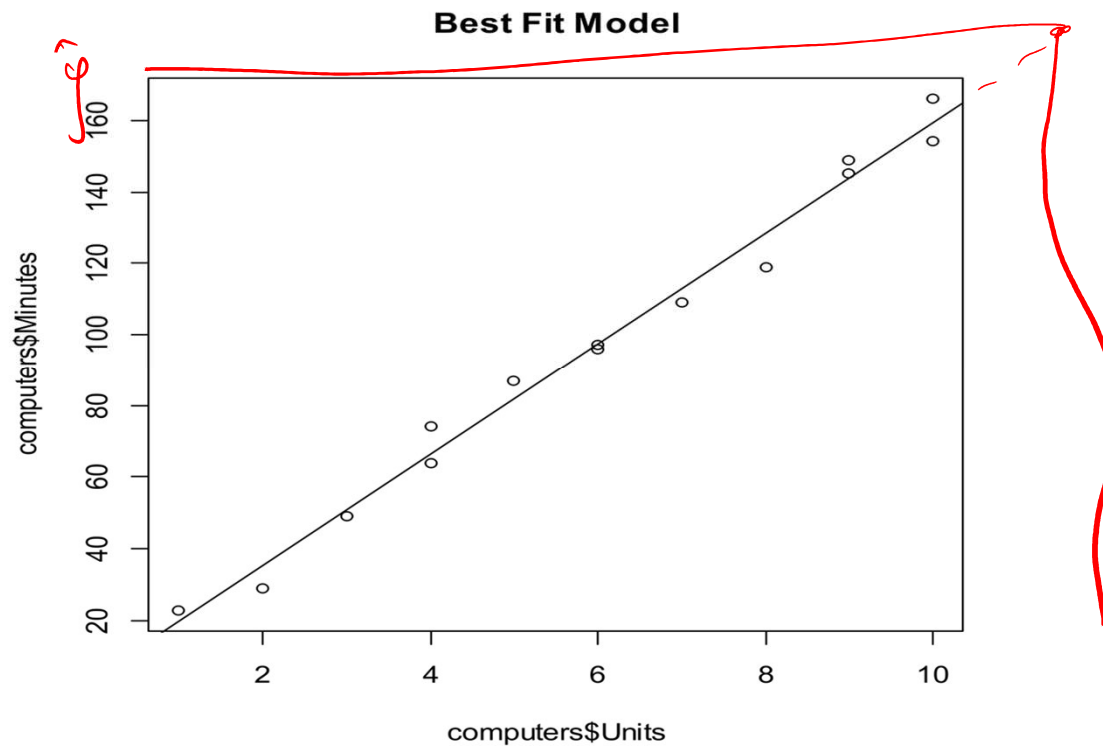
## Best fit Model

Using the previously obtained formula,  $b_0$  and  $b_1$  for the given sample dataset is determined as 4.162 (approximately) and 15.509 (approximately) respectively as shown.



```
1  #Computing values for b0 and b1
2  x <- computers$Units
3  y <- computers$Minutes
4  xiyi <- x*y
5  n <- nrow(computers)
6  xmean <- mean(computers$Units)
7  ymean <- mean(computers$Minutes)
8  numerator <- sum(xiyi)-n*xmean*ymean
9  denominator <- sum(x^2)-n*(xmean^2)
10 b1 <- numerator/denominator
11 b0 <- ymean-b1*xmean
12 #values of b0 and b1
13 b0
14 [1] 4.161654
15 b1
16 [1] 15.50877
```

The plot of this model along with the given data is as shown below



```
1 #Plot of Best fit Model: minutes = 4.161654 + 15.50877*units to be replaced
2 plot(computers$Units,computers$Minutes, main="Model 1")
3 abline(b0,b1)
```

## Best fit Model

The following code snippet shows the number of units replaced, observed time taken, expected time taken (based on the model) and the difference between predicted and observed values for the best fit model. The total sum of squared errors for the best fit model is 348.8484.

```
1 #Best fit model Units, Observed time, Expected time, Difference between observed and expected time
2 best_fit <- data.frame(matrix(data=c(computers$Units,computers$Minutes,(b0+b1*computers$Units),
3                                     ((b0+b1*computers$Units)-computers$Minutes)),ncol=4))
4 colnames(best_fit) <- c("Units replaced", "Observed time", "Expected time",
5                         "Expected - observed value")
6 best_fit
7 Units replaced Observed time Expected time Expected - observed value
8              1           23      19.67043      -3.3295739
9              2           29      35.17920       6.1791980
10             3           49      50.68797       1.6879699
11             4           64      66.19674       2.1967419
12             4           74      66.19674      -7.8032581
13             5           87      81.70551      -5.2944862
14             6           96      97.21429       1.2142857
15             6           97      97.21429       0.2142857
16             7          109     112.72306       3.7230576
17             8          119     128.23183       9.2318296
18             9          149     143.74060      -5.2593985
19             9          145     143.74060      -1.2593985
20            10          154     159.24937       5.2493734
21            10          166     159.24937      -6.7506266
```

```
1 #Sum of squared errors for Best fit model
2 sum((((b0+b1*computers$Units)-computers$Minutes))^2)
3 [1] 348.8484 →
```

$m_0 = -27$   
 $m_1 = 4997$   
 $m_2 = 5001$   
 $m_3 = 348$



# Best Fit Model

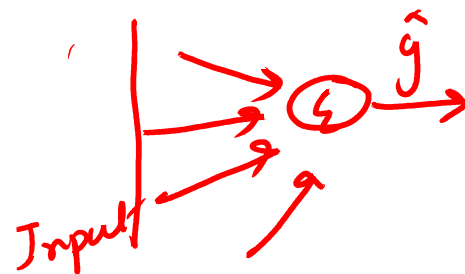
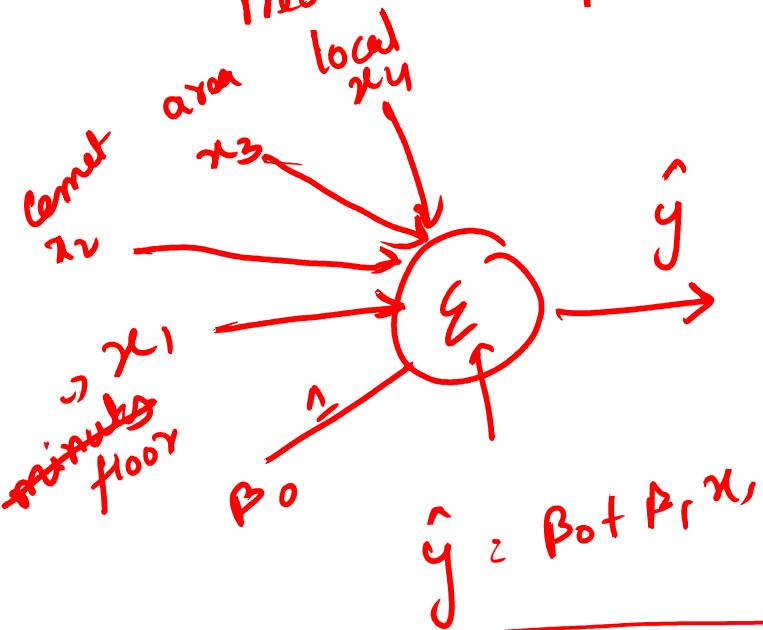
The best fit linear regression model can also be obtained in R using the `lm()` command as shown below:

Syntax for the `lm()` function : `lm(dependent variable ~ predictor variable)`

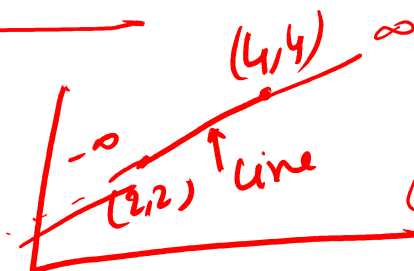
```
1 lm.model <- lm(computers$Minutes~computers$Units)
2 lm.model
3
4 Call:
5 lm(formula = computers$Minutes ~ computers$Units)
6
7 Coefficients:
8 (Intercept)  computers$Units
9  4.162        15.509
```

The value of coefficients  $b_0$  and  $b_1$  indicate that it takes approximately 15.509 minutes to replace a unit and a fixed time of approximately 4.162 minutes to understand a given repair.

Predict Price of a flat



$\Sigma_{i=1}^n$



length?  
 $(-\infty, \infty)$



# Logistic Regression

# Iris Dataset

SL, S.W, P.L, P.W



flower

multidimensional  
flat

| Sepal Length | Sepal Width | Petal Length | Petal Width | Class           |
|--------------|-------------|--------------|-------------|-----------------|
| 5.1          | 3.5         | 1.4          | 0.2         | Iris-setosa     |
| 4.9          | 3           | 1.4          | 0.2         | Iris-setosa     |
| 5.6          | 2.5         | 3.9          | 1.1         | Iris-versicolor |
| 5.9          | 3.2         | 4.8          | 1.8         | Iris-versicolor |
| 6.3          | 2.9         | 5.6          | 1.8         | Iris-virginica  |
| 6.5          | 3           | 5.8          | 2.2         | Iris-virginica  |

150 Records



Iris Versicolor



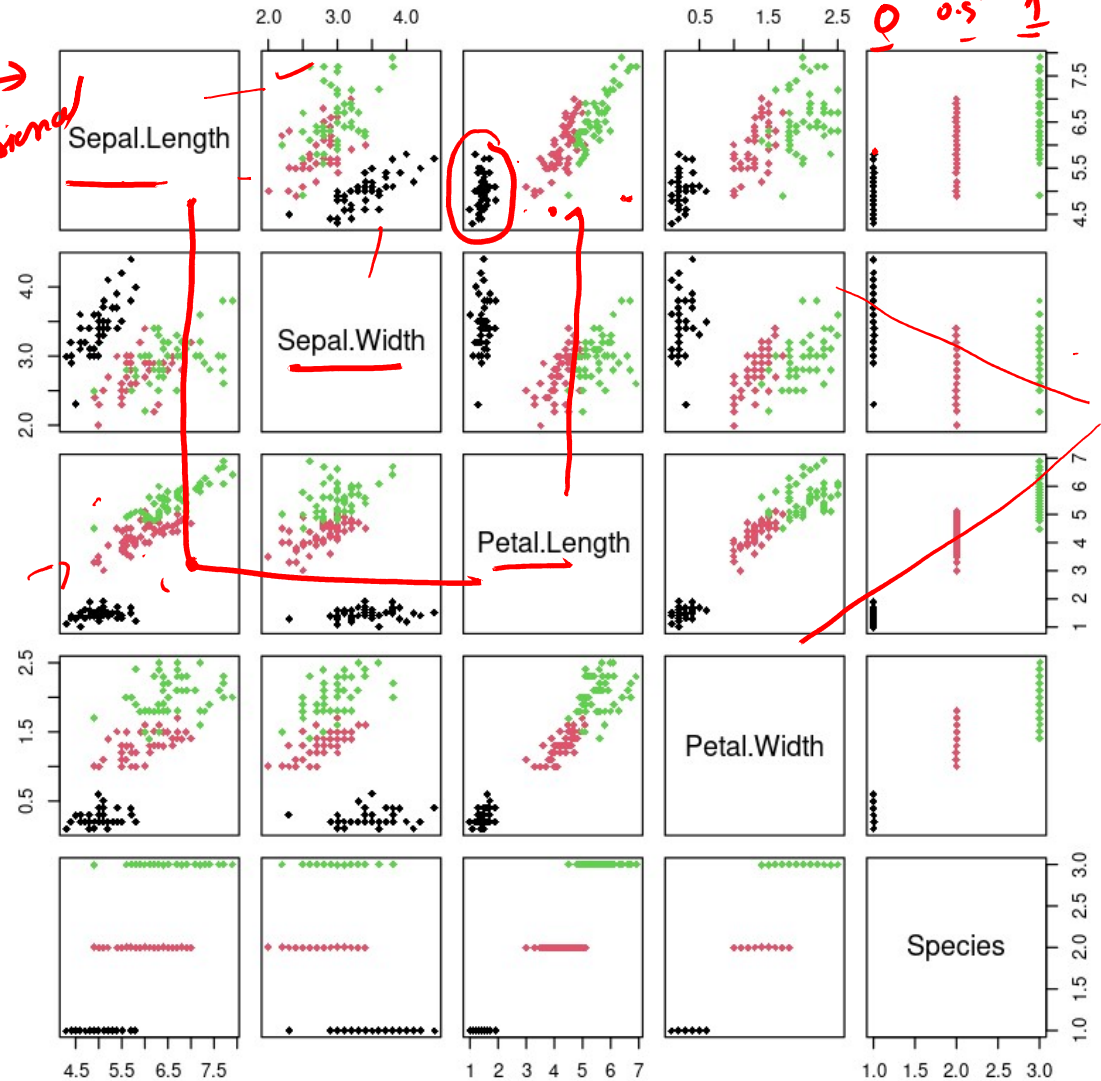
Iris Setosa



Iris Virginica

Catzen

0.5



# Learning Function

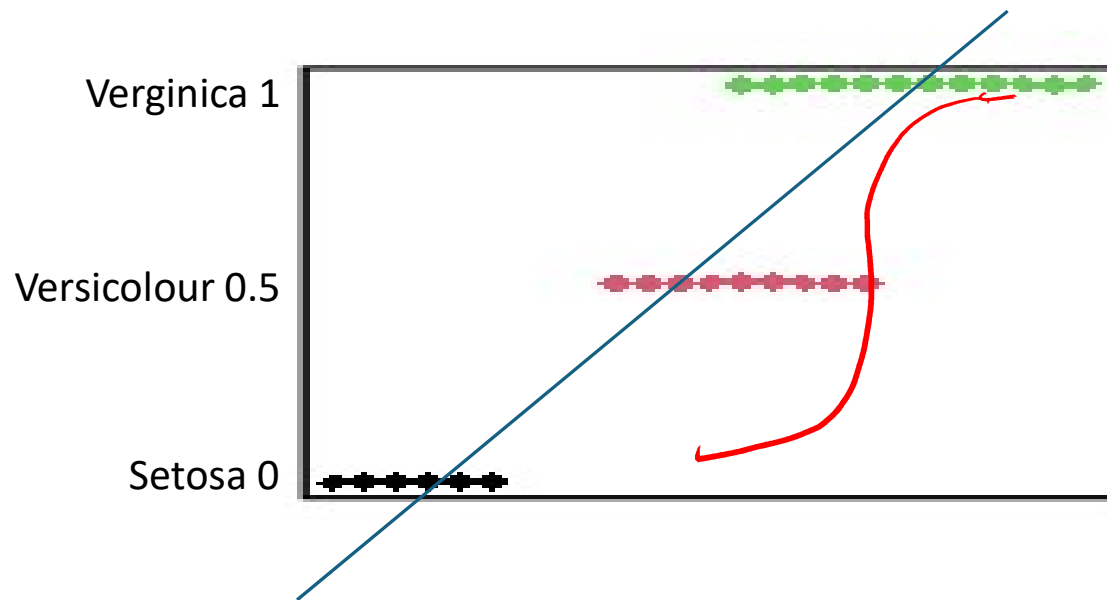
Multiple Linear Regression

$$F(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4$$

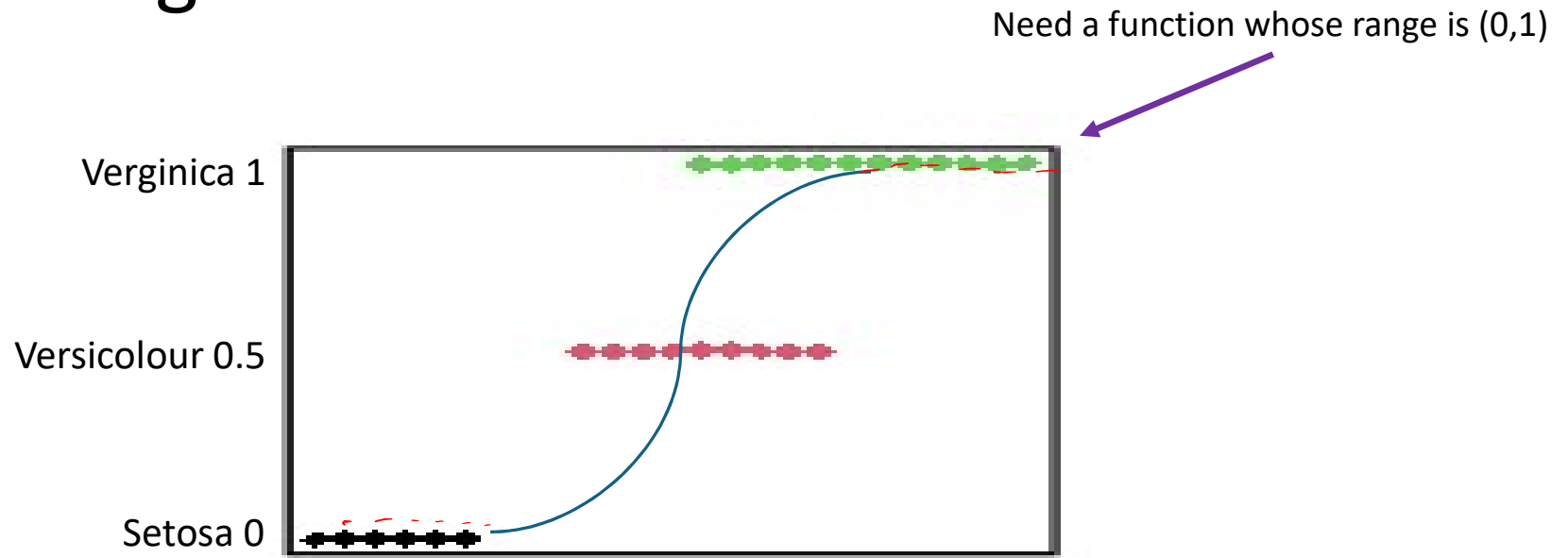
Simple Linear Regression

$$F(x) = \beta_0 + \beta_1 x_1$$

Range:  $(-\infty, \infty)$



# Learning Function



# Odds

**Odds** of my winning to losing is 2 is to 8.

- **Odds** – Something happening to something not happening

$$\text{Odds} = \frac{\text{Something happening}}{\text{Something not happening}} = \frac{\text{Win}}{\text{Loose}}$$

- **Probability** – Something happening to all the set of events

$$P_{\text{win}} = \frac{\text{Something happening}}{(\text{Something happening}) + (\text{Something not happening})} = \frac{\text{Win}}{\text{Win} + \text{Loose}}$$

$$P_{\text{Loose}} = \frac{\text{Something not happening}}{(\text{Something happening}) + (\text{Something not happening})} = \frac{\text{Loose}}{\text{Win} + \text{Loose}}$$

$$\text{Odds} = \frac{\text{Something happening}}{\text{Something not happening}} = \frac{\text{Win}}{\text{Loose}} = \frac{\text{Win} / (\text{Win} + \text{Loose})}{\text{Loose} / (\text{Win} + \text{Loose})} = \frac{P_{\text{win}}}{P_{\text{Loose}}}$$

$$= \frac{P_{\text{win}}}{1 - P_{\text{win}}}$$

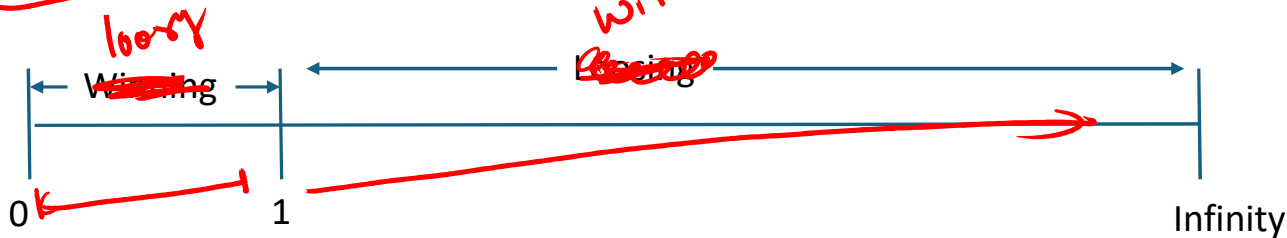
# Odds

## ~~Winning~~

- $\text{Odd}(1:2) = \frac{1}{2} = 0.5$
- $\text{Odd}(1:4) = \frac{1}{4} = 0.25$
- $\text{Odd}(1:8) = \frac{1}{8} = 0.125$
- $\text{Odd}(1:16) = \frac{1}{16} = 0.0625$
- $\text{Odd}(1:32) = \frac{1}{32} = 0.03125$

## ~~Loosing~~

- $\text{Odd}(2:1) = \frac{2}{1} = 2$
- $\text{Odd}(4:1) = \frac{4}{1} = 4$
- $\text{Odd}(8:1) = \frac{8}{1} = 8$
- $\text{Odd}(16:1) = \frac{16}{1} = 16$
- $\text{Odd}(32:1) = \frac{32}{1} = 32$



6000:1



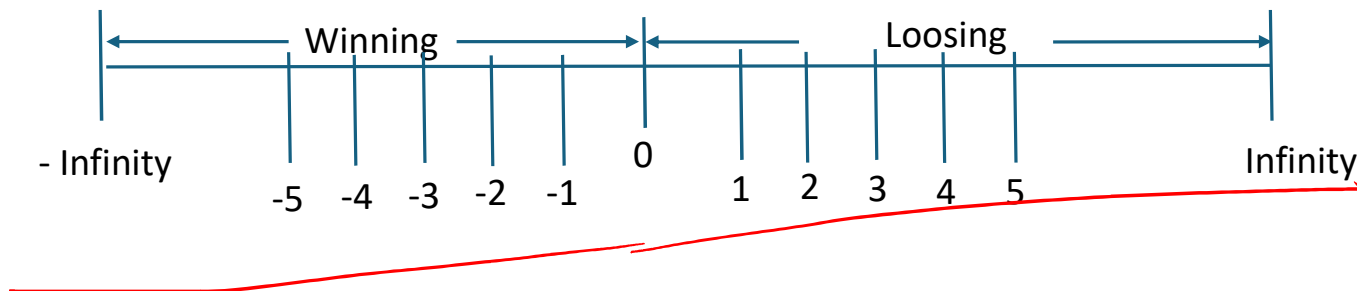
# Log of Odds

## Winning

- $\text{Odd}(1:2) = \frac{1}{2} = 0.5 = \log(0.5) = -1$
- $\text{Odd}(1:4) = \frac{1}{4} = 0.25 = \log(0.25) = -2$
- $\text{Odd}(1:8) = \frac{1}{8} = 0.125 = \log(0.125) = -3$
- $\text{Odd}(1:16) = \frac{1}{16} = 0.0625 = \log(0.0625) = -4$
- $\text{Odd}(1:32) = \frac{1}{32} = 0.03125 = \log(0.03125) = -5$

## Loosing

- $\text{Odd}(2:1) = \frac{2}{1} = 2 = \log(2) = 1$
- $\text{Odd}(4:1) = \frac{4}{1} = 4 = \log(4) = 2$
- $\text{Odd}(8:1) = \frac{8}{1} = 8 = \log(8) = 3$
- $\text{Odd}(16:1) = \frac{16}{1} = 16 = \log(16) = 4$
- $\text{Odd}(32:1) = \frac{32}{1} = 32 = \log(32) = 5$



# Log of Odds

$$\underline{Odds} = \frac{P_{win}}{P_{Loose}}$$

$$\log(Odds) = \log\left(\frac{P_{win}}{P_{Loose}}\right)$$

$$\log(Odds) = \log\left(\frac{P_{win}}{1 - P_{win}}\right)$$

$$\log\left(\frac{P_{win}}{1 - P_{win}}\right) = \underline{\beta_0} + \beta_1 x$$

$$\frac{P_{win}}{1 - P_{win}} = e^{\beta_0 + \beta_1 x}$$

$$P_{win} = (1 - P_{win})e^{\beta_0 + \beta_1 x}$$

$$P_{win} = 1e^{\beta_0 + \beta_1 x} - P_{win}e^{\beta_0 + \beta_1 x}$$

$$P_{win} + P_{win}e^{\beta_0 + \beta_1 x} = 1e^{\beta_0 + \beta_1 x}$$

$$P_{win}(1 + e^{\beta_0 + \beta_1 x}) = 1e^{\beta_0 + \beta_1 x}$$

$$P_{win} = \frac{e^{\beta_0 + \beta_1 x}}{(1 + e^{\beta_0 + \beta_1 x})}$$

Divide numerator and Denominator by  $e^{\beta_0 + \beta_1 x}$

$$P_{win} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

logistic  
function

Logistic Function/Sigmoid Function