KULLIYYAH OF INFORMATION & COMMUNICATION TECHNOLOGY

SEMESTER 3, 2021/2022

INFO 4313     DATA MINING
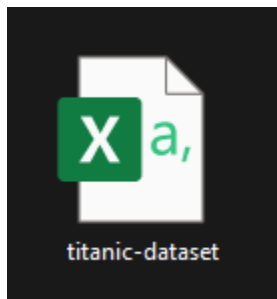
SECTION 01

*"FINAL ASSESSMENT"*

PREPARED BY:
JOYADDAR MD JOBAYER 1731833

1. I have downloaded the excel file as given in part A which is Titanic.csv.


titanic-dataset

2.

The titanic dataset contains 891 data. After downloading the dataset, I did some pre-processing with the dataset. Firstly, I used a substitute formula in excel to clean the Name column. After that, I have deleted to Colum from the dataset to run with Weka which are cabin and ticket. Also, I filled up with 0 for missing data or empty cell. Lastly, I have changed with 1 for yes and 0 for No in survive Colum. I used IF formula to change the value from 0,1 to yes, no.

| Passenger | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22 | 1 | 0 | A/5 21171 | 7.25 | | S |
| 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Thayer) | female | 38 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26 | 0 | 0 | STON/O2. | 7.925 | | S |
| 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35 | 1 | 0 | 113803 | 53.1 | C123 | S |
| 5 | 0 | 3 | Allen, Mr. William Henry | male | 35 | 0 | 0 | 373450 | 8.05 | | S |
| 6 | 0 | 3 | Moran, Mr. James | male | | 0 | 0 | 330877 | 8.4583 | | Q |
| 7 | 0 | 1 | McCarthy, Mr. Timothy J | male | 54 | 0 | 0 | 17463 | 51.8625 | E46 | S |
| 8 | 0 | 3 | Palsson, Master. Gosta Leonard | male | 2 | 3 | 1 | 349909 | 21.075 | | S |
| 9 | 1 | 3 | Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg) | female | 27 | 0 | 2 | 347742 | 11.1333 | | S |
| 10 | 1 | 2 | Nasser, Mrs. Nicholas (Adele Achem) | female | 14 | 1 | 0 | 237736 | 30.0708 | | C |
| 11 | 1 | 3 | Sandstrom, Miss. Marguerite Rut | female | 4 | 1 | 1 | PP 9549 | 16.7 | G6 | S |
| 12 | 1 | 1 | Bonnell, Miss. Elizabeth | female | 58 | 0 | 0 | 113783 | 26.55 | C103 | S |
| 13 | 0 | 3 | Saundercock, Mr. William Henry | male | 20 | 0 | 0 | A/5. 2151 | 8.05 | | S |
| 14 | 0 | 3 | Andersson, Mr. Anders Johan | male | 39 | 1 | 5 | 347082 | 31.275 | | S |
| 15 | 0 | 3 | Vestrom, Miss. Hulda Amanda Adolfina | female | 14 | 0 | 0 | 350406 | 7.8542 | | S |
| 16 | 1 | 2 | Hewlett, Mrs. (Mary D Kingcome) | female | 55 | 0 | 0 | 248706 | 16 | | S |
| 17 | 0 | 3 | Rice, Master. Eugene | male | 2 | 4 | 1 | 382652 | 29.125 | | Q |
| 18 | 1 | 2 | Williams, Mr. Charles Eugene | male | | 0 | 0 | 244373 | 13 | | S |
| 19 | 0 | 3 | Vander Planke, Mrs. Julius (Emelia Maria Vandemoortele) | female | 31 | 1 | 0 | 345763 | 18 | | S |
| 20 | 1 | 3 | Masselmani, Mrs. Fatima | female | | 0 | 0 | 2649 | 7.225 | | C |
| 21 | 0 | 2 | Fynney, Mr. Joseph J | male | 35 | 0 | 0 | 239865 | 26 | | S |
| 22 | 1 | 2 | Beesley, Mr. Lawrence | male | 34 | 0 | 0 | 248698 | 13 | D56 | S |
| 23 | 1 | 3 | McGowan, Miss. Anna "Annie" | female | 15 | 0 | 0 | 330923 | 8.0292 | | Q |
| 24 | 1 | 1 | Sloper, Mr. William Thompson | male | 28 | 0 | 0 | 113788 | 35.5 | A6 | S |
| 25 | 0 | 3 | Palsson, Miss. Torborg Danira | female | 8 | 3 | 1 | 349909 | 21.075 | | S |
| 26 | 1 | 3 | Asplund, Mrs. Carl Oscar (Selma Augusta Emilia Johansson) | female | 38 | 1 | 5 | 347077 | 31.3875 | | S |
| 27 | 0 | 3 | Emir, Mr. Farred Chehab | male | | 0 | 0 | 2631 | 7.225 | | C |
| 28 | 0 | 1 | Fortune, Mr. Charles Alexander | male | 19 | 3 | 2 | 19950 | 263 | C23 C25 C2 | S |
| 29 | 1 | 3 | O'Dwyer, Miss. Ellen "Nellie" | female | | 0 | 0 | 330959 | 7.8792 | | Q |
| 30 | 0 | 3 | Todoroff, Mr. Lalio | male | | 0 | 0 | 349216 | 7.8958 | | S |
| 31 | 0 | 1 | Uruchurtu, Don. Manuel E | male | 40 | 0 | 0 | PC 17601 | 27.7208 | | C |
| 32 | 1 | 1 | Spencer, Mrs. William Augustus (Marie Eugenie) | female | | 1 | 0 | PC 17569 | 146.5208 | B78 | C |
| 33 | 1 | 3 | Glynn, Miss. Mary Agatha | female | | 0 | 0 | 335677 | 7.75 | | Q |
| 34 | 0 | 2 | Wheadon, Mr. Edward H | male | 66 | 0 | 0 | C.A. 24579 | 10.5 | | S |

Before pre-processing the titanic dataset

| | A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|---|
| | PassengerId | Name | Sex | Age | SibSp | Parch | Fare | Embarked | Pclass | Survived |
| 1 | 1 | Braund  Mr. Owen Harris | male | 22 | 1 | 0 | 7.25 | S | 3 | No |
| 2 | 2 | Cumings  Mrs. John Bradley  Florence Briggs Thayer | female | 38 | 1 | 0 | 71.2833 | C | 1 | Yes |
| 3 | 3 | Heikkinen  Miss. Laina | female | 26 | 0 | 0 | 7.925 | S | 3 | Yes |
| 4 | 4 | Futrelle  Mrs. Jacques Heath  Lily May Peel | female | 35 | 1 | 0 | 53.1 | S | 1 | Yes |
| 5 | 5 | Allen  Mr. William Henry | male | 35 | 0 | 0 | 8.05 | S | 3 | No |
| 6 | 6 | Moran  Mr. James | male | 0 | 0 | 0 | 8.4583 | Q | 3 | No |
| 7 | 7 | McCarthy  Mr. Timothy J | male | 54 | 0 | 0 | 51.8625 | S | 1 | No |
| 8 | 8 | Palsson  Master. Gosta Leonard | male | 2 | 3 | 1 | 21.075 | S | 3 | No |
| 9 | 9 | Johnson  Mrs. Oscar W  Elisabeth Vilhelmina Berg | female | 27 | 0 | 2 | 11.1333 | S | 3 | Yes |
| 10 | 10 | Nasser  Mrs. Nicholas  Adele Achem | female | 14 | 1 | 0 | 30.0708 | C | 2 | Yes |
| 11 | 11 | Sandstrom  Miss. Marguerite Rut | female | 4 | 1 | 1 | 16.7 | S | 3 | Yes |
| 12 | 12 | Bonnell  Miss. Elizabeth | female | 58 | 0 | 0 | 26.55 | S | 1 | Yes |
| 13 | 13 | Saundercock  Mr. William Henry | male | 20 | 0 | 0 | 8.05 | S | 3 | No |
| 14 | 14 | Andersson  Mr. Anders Johan | male | 39 | 1 | 5 | 31.275 | S | 3 | No |
| 15 | 15 | Vestrom  Miss. Hulda Amanda Adolfina | female | 14 | 0 | 0 | 7.8542 | S | 3 | No |
| 16 | 16 | Hewlett  Mrs. Mary D Kingcome | female | 55 | 0 | 0 | 16 | S | 2 | Yes |
| 17 | 17 | Rice  Master. Eugene | male | 2 | 4 | 1 | 29.125 | Q | 3 | No |
| 18 | 18 | Williams  Mr. Charles Eugene | male | 0 | 0 | 0 | 13 | S | 2 | Yes |
| 19 | 19 | Vander Planke  Mrs. Julius  Emelia Maria Vandemoortele | female | 31 | 1 | 0 | 18 | S | 3 | No |
| 20 | 20 | Masselmani  Mrs. Fatima | female | 0 | 0 | 0 | 7.225 | C | 3 | Yes |
| 21 | 21 | Fynney  Mr. Joseph J | male | 35 | 0 | 0 | 26 | S | 2 | No |
| 22 | 22 | Beesley  Mr. Lawrence | male | 34 | 0 | 0 | 13 | S | 2 | Yes |
| 23 | 23 | McGowan  Miss. Anna  Annie | female | 15 | 0 | 0 | 8.0292 | Q | 3 | Yes |
| 24 | 24 | Sloper  Mr. William Thompson | male | 28 | 0 | 0 | 35.5 | S | 1 | Yes |
| 25 | 25 | Palsson  Miss. Torborg Danira | female | 8 | 3 | 1 | 21.075 | S | 3 | No |
| 26 | 26 | Asplund  Mrs. Carl Oscar  Selma Augusta Emilia Johansson | female | 38 | 1 | 5 | 31.3875 | S | 3 | Yes |
| 27 | 27 | Emir  Mr. Farred Chehab | male | 0 | 0 | 0 | 7.225 | C | 3 | No |
| 28 | 28 | Fortune  Mr. Charles Alexander | male | 19 | 3 | 2 | 263 | S | 1 | No |
| 29 | 29 | O Dwyer  Miss. Ellen  Nellie | female | 0 | 0 | 0 | 7.8792 | Q | 3 | Yes |
| 30 | 30 | Todoroff  Mr. Lalio | male | 0 | 0 | 0 | 7.8958 | S | 3 | No |
| 31 | 31 | Uruchurtu  Don. Manuel E | male | 40 | 0 | 0 | 27.7208 | C | 1 | No |
| 32 | 32 | Spencer  Mrs. William Augustus  Marie Eugenie | female | 0 | 1 | 0 | 146.5208 | C | 1 | Yes |
| 33 | 33 | Glynn  Miss. Mary Agatha | female | 0 | 0 | 0 | 7.75 | Q | 3 | Yes |
| 34 | 34 | Wheadon  Mr. Edward H | male | 66 | 0 | 0 | 10.5 | S | 2 | No |

After pre-processing the titanic dataset

3. I have generated 3 models such as Naive Bayes, Logistic and J48. But I have chosen Naïve Bayes models to predict which passengers have high likely to survive.



Figure: Naïve Bayes model

Figure: Logistic model

Figure: J48 model

4. In Naïve Bayes model, after using the 10-fold cross validation calculating the accuracy and confusion matrix which can be produced to support that passengers have high likely to survive from the dataset.

**NAÏVE BAYES**

=== Confusion Matrix ===

```
  a   b   <-- classified as
 474 75 |  a = No
 115 227 |  b = Yes
```

From the confusion matrix we can say that,

TP = 474

FP = 75

TN = 227

FN = 115

Accuracy = TP + TN / TP + FP + FN + TN

= 474 + 227 / 474 + 75 + 115 + 227

= 701 / 891

= 0.79

According to the calculation, Naïve Bayes has 0.79 or 79% accuracy of survivors.

5.

```
  ▼        0        (normalized) Cabin=D28
  +        0.7797   (normalized) Cabin=E17
  +       -0.5098   (normalized) Cabin=A24
  +        0        (normalized) Cabin=C50
  +        0.01     (normalized) Cabin=B42
  +        0.4898   (normalized) Cabin=C148
  +        0.11     (normalized) Embarked=S
  +        0.3199   (normalized) Embarked=C
  +        0.06     (normalized) Embarked=Q
  +        0.45
```

```
=== Re-evaluation on test set ===

User supplied test set
Relation:    predict
Instances:   unknown (yet). Reading incrementally
Attributes:  12

=== Predictions on user test set ===

     inst#    actual  predicted error prediction
       1       1:?        1:No      1
       2       1:?        1:No      1
       3       1:?        1:No      1
       4       1:?        1:No      1
       5       1:?        1:No      1
       6       1:?        1:No      1
       7       1:?        1:No      1
       8       1:?        1:No      1
       9       1:?        1:No      1

=== Summary ===

Total Number of Instances                  0
Ignored Class Unknown Instances            9
```

According to the screenshot above, I used Function SGD to get the prediction:

inst# actual predicted error prediction

1 1:? 1:No 1

2 1:? 1:No 1

3 1:? 1:No 1

4 1:? 1:No 1

5 1:? 1:No 1

6 1:? 1:No 1

7 1:? 1:No 1

8 1:? 1:No 1

9 1:? 1:No 1

PART 2:

LINEAR REGRESSION:

**Predicting the value of weight when the height are 80,82,60 and 58**

Let, Height =X and Weight = Y

| X | Y | X - mean(x) | Y - mean(y) | X - mean(x) * Y - mean(y) | X - mean(x)^2 |
|---|---|---|---|---|---|
| 63 | 127 | -6.3 | -31.8 | 200.34 | 39.69 |
| 64 | 121 | -5.3 | -37.8 | 200.34 | 28.09 |
| 66 | 142 | -3.3 | -16.8 | 55.44 | 10.89 |
| 69 | 157 | -0.3 | -1.8 | 0.54 | 0.09 |
| 69 | 162 | -0.3 | 3.2 | -0.96 | 0.09 |
| 71 | 156 | 1.7 | -2.8 | -4.76 | 2.89 |
| 71 | 169 | 1.7 | 10.2 | 17.34 | 2.89 |
| 72 | 165 | 2.7 | 6.2 | 16.74 | 7.29 |
| 73 | 181 | 3.7 | 22.2 | 82.14 | 13.69 |
| 75 | 208 | 5.7 | 49.2 | 280.44 | 32.49 |
| | | | | 847.6 | 138.1 |

| | |
|---|---|
| Mean (x) | 69.3 |
| Mean (y) | 158.8 |

y=mx+c

| | |
|---|---|
| m | 6.138 |
| c | -266.534 |

y= 6.138x + (-266.534)

| | E | F |
|---|---|---|
| 52 | Mean (y) | 158.8 |
| 53 | | |
| 54 | y=mx+c | |
| 55 | m | 6.138 |
| 56 | c | -266.534 |
| 57 | | |
| 58 | y= 6.138x + (-266.534) | |

| | E | G/H | I |
|---|---|---|---|
| 61 | **When x=80** | **So the predict values are,** | |
| 62 | y= 6.138 (80) + (-266.534) | **X** | **Py** |
| 63 | y = 224.506 | 80 | 224.506 |
| 64 | **When x=82** | 82 | 236.782 |
| 65 | y= 6.138 (82) + (-266.534) | 60 | 101.746 |
| 66 | y = 236.782 | 58 | 89.47 |
| 67 | **When x= 60** | | |
| 68 | y= 6.138 (60) + (-266.534) | | |
| 69 | y = 101.746 | | |
| 70 | **When x= 58** | | |
| 71 | y= 6.138 (58) + (-266.534) | | |
| 72 | y = 89.47 | | |
| 73 | | | |
| 74 | | | |

Cell reference: O62

| | A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|---|
| 72 | | | | | y = 89.47 | | | | | |
| 73 | | | | | | | | | | |
| 74 | | | | | | | | | | |
| 75 | | | | | | | | | | |
| 76 | Finding the R^2 , MSE and MAE | | | | | | | | | |
| 77 | | | | | | | | | | |
| 78 | | | | | | | | | | |
| 79 | | | | | | | | | | |
| 80 | | R^2 | | | | | | | | |
| 81 | | | | | | | | | | |
| 82 | | | X | Y | Py | Py - Y | (Py-Y)^2 | (Y-mean(y))^2 | | |
| 83 | | | 63 | 127 | 120.16 | -6.84 | 46.7856 | 1011.24 | | |
| 84 | | | 64 | 121 | 126.298 | 5.298 | 28.068804 | 1428.84 | | |
| 85 | | | 66 | 142 | 138.574 | -3.426 | 11.737476 | 282.24 | | |
| 86 | | | 69 | 157 | 156.988 | -0.012 | 0.000144 | 3.24 | | |
| 87 | | | 69 | 162 | 156.988 | -5.012 | 25.120144 | 10.24 | | |
| 88 | | | 71 | 156 | 169.264 | 13.264 | 175.933696 | 7.84 | | |
| 89 | | | 71 | 169 | 169.264 | 0.264 | 0.069696 | 104.04 | | |
| 90 | | | 72 | 165 | 175.402 | 10.402 | 108.201604 | 38.44 | | |
| 91 | | | 73 | 181 | 181.54 | 0.54 | 0.2916 | 492.84 | | |
| 92 | | | 75 | 208 | 193.816 | -14.184 | 201.185856 | 2420.64 | | |
| 93 | | | | | | | 597.39462 | 5799.6 | | |
| 94 | | | | | | | | | | |
| 95 | | | | | | | | | | |
| 96 | | | | | R^2 | 0.897 | | | | |

**MSE**

| | X | Y | Py | Py - Y | (Py-Y)^2 |
|---|---|---|---|---|---|
| | 63 | 127 | 120.16 | -6.84 | 46.7856 |
| | 64 | 121 | 126.298 | 5.298 | 28.068804 |
| | 66 | 142 | 138.574 | -3.426 | 11.737476 |
| | 69 | 157 | 156.988 | -0.012 | 0.000144 |
| | 69 | 162 | 156.988 | -5.012 | 25.120144 |
| | 71 | 156 | 169.264 | 13.264 | 175.933696 |
| | 71 | 169 | 169.264 | 0.264 | 0.069696 |
| | 72 | 165 | 175.402 | 10.402 | 108.201604 |
| | 73 | 181 | 181.54 | 0.54 | 0.2916 |
| | 75 | 208 | 193.816 | -14.184 | 201.185856 |
| | | | | | 597.39462 |
| | | | | | |
| | | | MSE | 59.739462 | |

**MAE**

| X | Y | Py | Py - Y | \|Py - Y\| |
|----|-----|---------|---------|----------|
| 63 | 127 | 120.16  | -6.84   | 6.84     |
| 64 | 121 | 126.298 | 5.298   | 5.298    |
| 66 | 142 | 138.574 | -3.426  | 3.426    |
| 69 | 157 | 156.988 | -0.012  | 0.012    |
| 69 | 162 | 156.988 | -5.012  | 5.012    |
| 71 | 156 | 169.264 | 13.264  | 13.264   |
| 71 | 169 | 169.264 | 0.264   | 0.264    |
| 72 | 165 | 175.402 | 10.402  | 10.402   |
| 73 | 181 | 181.54  | 0.54    | 0.54     |
| 75 | 208 | 193.816 | -14.184 | 14.184   |
|    |     |         |         | 59.242   |

| MAE | 5.9242 |
|-----|--------|

*** CTRL WITH ~ TO SHOW ALL PROCESS

Formula bar — C35: Let, Height =X and Weight = Y

**R^2**

| | X | Y | Py | Py - Y | (Py-Y)^2 | (Y-mean(y))^2 |
|---|---|---|---|---|---|---|
| 83 | 63 | 127 | =6.138*(63)+(-266.534) | =E83-D83 | =POWER(F83,2) | =POWER(H38,2) |
| 84 | 64 | 121 | =6.138*(64)+(-266.534) | =E84-D84 | =POWER(F84,2) | =POWER(H39,2) |
| 85 | 66 | 142 | =6.138*(66)+(-266.534) | =E85-D85 | =POWER(F85,2) | =POWER(H40,2) |
| 86 | 69 | 157 | =6.138*(69)+(-266.534) | =E86-D86 | =POWER(F86,2) | =POWER(H41,2) |
| 87 | 69 | 162 | =6.138*(69)+(-266.534) | =E87-D87 | =POWER(F87,2) | =POWER(H42,2) |
| 88 | 71 | 156 | =6.138*(71)+(-266.534) | =E88-D88 | =POWER(F88,2) | =POWER(H43,2) |
| 89 | 71 | 169 | =6.138*(71)+(-266.534) | =E89-D89 | =POWER(F89,2) | =POWER(H44,2) |
| 90 | 72 | 165 | =6.138*(72)+(-266.534) | =E90-D90 | =POWER(F90,2) | =POWER(H45,2) |
| 91 | 73 | 181 | =6.138*(73)+(-266.534) | =E91-D91 | =POWER(F91,2) | =POWER(H46,2) |
| 92 | 75 | 208 | =6.138*(75)+(-266.534) | =E92-D92 | =POWER(F92,2) | =POWER(H47,2) |
| 93 | | | | | =SUM(G83:G92) | =SUM(H83:H92) |

| | | | Py | Py - Y | | |
|---|---|---|---|---|---|---|
| 96 | | | R^2 | =1-(G93/H93) | | |

**MSE**

| | X | Y | Py | Py - Y | (Py-Y)^2 |
|---|---|---|---|---|---|
| 102 | 63 | 127 | =6.138*(63)+(-266.534) | =E102-D102 | =POWER(F102,2) |
| 103 | 64 | 121 | =6.138*(64)+(-266.534) | =E103-D103 | =POWER(F103,2) |
| 104 | 66 | 142 | =6.138*(66)+(-266.534) | =E104-D104 | =POWER(F104,2) |
| 105 | 69 | 157 | =6.138*(69)+(-266.534) | =E105-D105 | =POWER(F105,2) |
| 106 | 69 | 162 | =6.138*(69)+(-266.534) | =E106-D106 | =POWER(F106,2) |
| 107 | 71 | 156 | =6.138*(71)+(-266.534) | =E107-D107 | =POWER(F107,2) |
| 108 | 71 | 169 | =6.138*(71)+(-266.534) | =E108-D108 | =POWER(F108,2) |
| 109 | 72 | 165 | =6.138*(72)+(-266.534) | =E109-D109 | =POWER(F109,2) |
| 110 | 73 | 181 | =6.138*(73)+(-266.534) | =E110-D110 | =POWER(F110,2) |
| 111 | 75 | 208 | =6.138*(75)+(-266.534) | =E111-D111 | =POWER(F111,2) |
| 112 | | | | | =SUM(G102:G111) |

| | | | Py | Py - Y | |
|---|---|---|---|---|---|
| 115 | | | MSE | =G112/10 | |

C35 | Let, Height =X and Weight = Y

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 99 | | MSE | | | | | | |
| 100 | | | | | | | | |
| 101 | | | X | Y | Py | Py - Y | (Py-Y)^2 | |
| 102 | | | 63 | 127 | =6.138*(63)+(-266.534) | =E102-D102 | =POWER(F102,2) | |
| 103 | | | 64 | 121 | =6.138*(64)+(-266.534) | =E103-D103 | =POWER(F103,2) | |
| 104 | | | 66 | 142 | =6.138*(66)+(-266.534) | =E104-D104 | =POWER(F104,2) | |
| 105 | | | 69 | 157 | =6.138*(69)+(-266.534) | =E105-D105 | =POWER(F105,2) | |
| 106 | | | 69 | 162 | =6.138*(69)+(-266.534) | =E106-D106 | =POWER(F106,2) | |
| 107 | | | 71 | 156 | =6.138*(71)+(-266.534) | =E107-D107 | =POWER(F107,2) | |
| 108 | | | 71 | 169 | =6.138*(71)+(-266.534) | =E108-D108 | =POWER(F108,2) | |
| 109 | | | 72 | 165 | =6.138*(72)+(-266.534) | =E109-D109 | =POWER(F109,2) | |
| 110 | | | 73 | 181 | =6.138*(73)+(-266.534) | =E110-D110 | =POWER(F110,2) | |
| 111 | | | 75 | 208 | =6.138*(75)+(-266.534) | =E111-D111 | =POWER(F111,2) | |
| 112 | | | | | | | =SUM(G102:G111) | |
| 113 | | | | | | | | |
| 114 | | | | | | | | |
| 115 | | | | | MSE | =G112/10 | | |
| 116 | | | | | | | | |
| 117 | | | | | | | | |
| 118 | | MAE | | | | | | |
| 119 | | | | | | | | |
| 120 | | | X | Y | Py | Py - Y | \|Py - Y\| | |
| 121 | | | 63 | 127 | =6.138*(63)+(-266.534) | =E121-D121 | 6.84 | |
| 122 | | | 64 | 121 | =6.138*(64)+(-266.534) | =E122-D122 | 5.298 | |
| 123 | | | 66 | 142 | =6.138*(66)+(-266.534) | =E123-D123 | 3.426 | |
| 124 | | | 69 | 157 | =6.138*(69)+(-266.534) | =E124-D124 | 0.012 | |
| 125 | | | 69 | 162 | =6.138*(69)+(-266.534) | =E125-D125 | 5.012 | |
| 126 | | | 71 | 156 | =6.138*(71)+(-266.534) | =E126-D126 | 13.264 | |
| 127 | | | 71 | 169 | =6.138*(71)+(-266.534) | =E127-D127 | 0.264 | |
| 128 | | | 72 | 165 | =6.138*(72)+(-266.534) | =E128-D128 | 10.402 | |
| 129 | | | 73 | 181 | =6.138*(73)+(-266.534) | =E129-D129 | 0.54 | |
| 130 | | | 75 | 208 | =6.138*(75)+(-266.534) | =E130-D130 | 14.184 | |
| 131 | | | | | | | =SUM(G121:G130) | |
| 132 | | | | | | | | |
| 133 | | | | | | | | |
| 134 | | | | | | | | |
| 135 | | | | | MAE | =G131/10 | | |

Let, Height =X and Weight = Y

| | | | | | | |
|---|---|---|---|---|---|---|
| 99 | | **MSE** | | | | |
| 100 | | | | | | |
| 101 | | | X | Y | Py | Py - Y | (Py-Y)^2 |
| 102 | | | 63 | 127 | =6.138*(63)+(-266.534) | =E102-D102 | =POWER(F102,2) |
| 103 | | | 64 | 121 | =6.138*(64)+(-266.534) | =E103-D103 | =POWER(F103,2) |
| 104 | | | 66 | 142 | =6.138*(66)+(-266.534) | =E104-D104 | =POWER(F104,2) |
| 105 | | | 69 | 157 | =6.138*(69)+(-266.534) | =E105-D105 | =POWER(F105,2) |
| 106 | | | 69 | 162 | =6.138*(69)+(-266.534) | =E106-D106 | =POWER(F106,2) |
| 107 | | | 71 | 156 | =6.138*(71)+(-266.534) | =E107-D107 | =POWER(F107,2) |
| 108 | | | 71 | 169 | =6.138*(71)+(-266.534) | =E108-D108 | =POWER(F108,2) |
| 109 | | | 72 | 165 | =6.138*(72)+(-266.534) | =E109-D109 | =POWER(F109,2) |
| 110 | | | 73 | 181 | =6.138*(73)+(-266.534) | =E110-D110 | =POWER(F110,2) |
| 111 | | | 75 | 208 | =6.138*(75)+(-266.534) | =E111-D111 | =POWER(F111,2) |
| 112 | | | | | | | =SUM(G102:G111) |
| 113 | | | | | | | |
| 114 | | | | | | | |
| 115 | | | | | | MSE | =G112/10 |
| 116 | | | | | | | |
| 117 | | | | | | | |
| 118 | | **MAE** | | | | | |
| 119 | | | | | | | |
| 120 | | | X | Y | Py | Py - Y | \|Py - Y\| |
| 121 | | | 63 | 127 | =6.138*(63)+(-266.534) | =E121-D121 | 6.84 |
| 122 | | | 64 | 121 | =6.138*(64)+(-266.534) | =E122-D122 | 5.298 |
| 123 | | | 66 | 142 | =6.138*(66)+(-266.534) | =E123-D123 | 3.426 |
| 124 | | | 69 | 157 | =6.138*(69)+(-266.534) | =E124-D124 | 0.012 |
| 125 | | | 69 | 162 | =6.138*(69)+(-266.534) | =E125-D125 | 5.012 |
| 126 | | | 71 | 156 | =6.138*(71)+(-266.534) | =E126-D126 | 13.264 |
| 127 | | | 71 | 169 | =6.138*(71)+(-266.534) | =E127-D127 | 0.264 |
| 128 | | | 72 | 165 | =6.138*(72)+(-266.534) | =E128-D128 | 10.402 |
| 129 | | | 73 | 181 | =6.138*(73)+(-266.534) | =E129-D129 | 0.54 |
| 130 | | | 75 | 208 | =6.138*(75)+(-266.534) | =E130-D130 | 14.184 |
| 131 | | | | | | | =SUM(G121:G130) |
| 132 | | | | | | | |
| 133 | | | | | | | |
| 134 | | | | | | | |
| 135 | | | | | | MAE | =G131/10 |