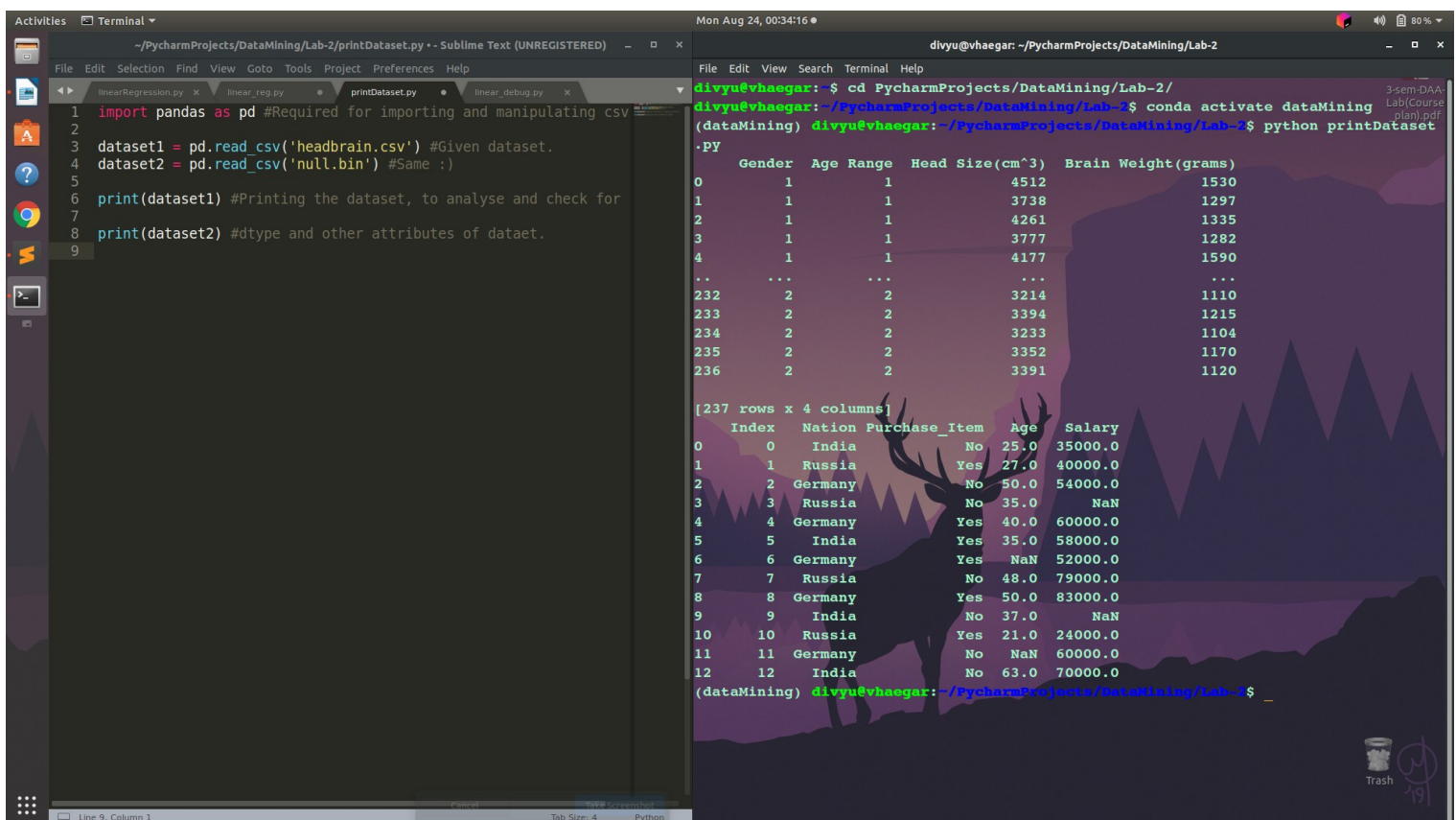


Data-Mining and predictions by machine~

~Given datasets

- > **headbrain** -> Datasets of shape (237 x 4), the columns of significant usecase being `head size` and `brain wieght`. Provides information about brain wieght(in g) w.r.t head size(in cm³).
- > **age_salary** -> Dataset of size (13 X 4), the columns of significant usecase being `head size` and `brain wieght`. Provides information about salary of an individual w.r.t his/her age.

Glance of the given datasets



```
import pandas as pd #Required for importing and manipulating csv
dataset1 = pd.read_csv('headbrain.csv') #Given dataset.
dataset2 = pd.read_csv('null.bin') #Same :)
print(dataset1) #Printing the dataset, to analyse and check for
print(dataset2) #dtype and other attributes of dataet.
```

```
divyu@vhaegar: ~/PycharmProjects/DataMining/Lab-2/
divyu@vhaegar:~/PycharmProjects/DataMining/Lab-2$ conda activate dataMining
(dataMining) divyu@vhaegar:~/PycharmProjects/DataMining/Lab-2$ python printDataset
.PY
  Gender  Age  Range  Head Size(cm^3)  Brain Weight(grams)
0      1     1      1           4512             1530
1      1     1      1           3738             1297
2      1     1      1           4261             1335
3      1     1      1           3777             1282
4      1     1      1           4177             1590
...     ...     ...           ...             ...
232     2     2      2           3214             1110
233     2     2      2           3394             1215
234     2     2      2           3233             1104
235     2     2      2           3352             1170
236     2     2      2           3391             1120

[237 rows x 4 columns]
  Index  Nation  Purchase_Item  Age  Salary
0      0   India             No  25.0  35000.0
1      1  Russia             Yes  27.0  40000.0
2      2  Germany             No  50.0  54000.0
3      3  Russia             No  35.0   NaN
4      4  Germany             Yes  40.0  60000.0
5      5   India             Yes  35.0  58000.0
6      6  Germany             Yes  NaN  52000.0
7      7  Russia             No  48.0  79000.0
8      8  Germany             Yes  50.0  83000.0
9      9   India             No  37.0   NaN
10     10  Russia             Yes  21.0  24000.0
11     11  Germany             No  NaN  60000.0
12     12   India             No  63.0  70000.0
(dataMining) divyu@vhaegar:~/PycharmProjects/DataMining/Lab-2$
```

Which ML-model is best applicable?

Since, data-points seem to have linear-dependance i.e can be represented as linear-algebraic combination of the known-variable, Single variable linear regression model will be best applicable.

I will be calculating Linear Regression using Least Square method, the required parameters for the same are,

- > x (known variable ~ headsize / age.
- > y (unknown variable ~ brain weight / salary.
- > $x - \text{mean}(x)$
- > $y - \text{mean}(y)$
- > $[x - \text{mean}(x)]^2$ {Taking the square is necessary because, $x - \text{mean}(x)$ can be negative at instances, since it's deviation from mean, sum of all data-points would be 0, hence squaring the term is necessary}
- > $[y - \text{mean}(y) * x - \text{mean}(x)]$

The goal is to find the equation of the best-fit line that describes the relation from known to unknown variables with least error.

Coefficient (or slope of line) :

$$\frac{\sum [x - \text{mean}(x)]^2}{\sum \{[y - \text{mean}(y)] * [x - \text{mean}(x)]\}}$$

Intercept (or y-intercept of line):

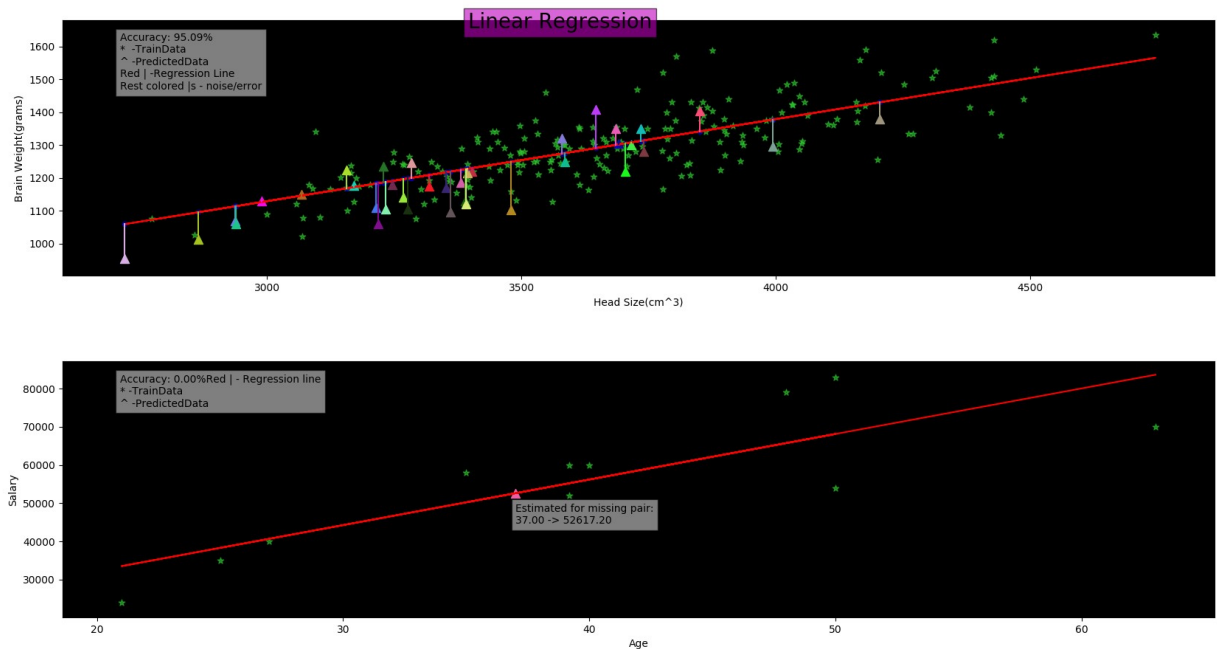
$$\text{mean}(y) - \text{slop} * \text{mean}(x)$$

Code for the Single Variable Linear Regression model

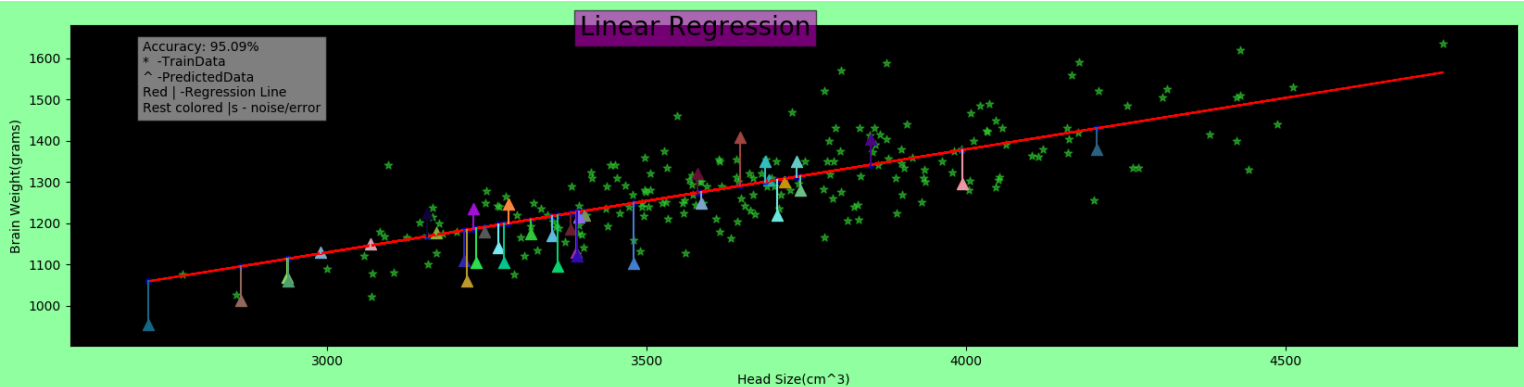
```
1 # Linear Regression Model.
2 # Code by,Dibyasom Puhan          24th August, 2020.
3 #
4 # Linear regression?
5 #-----
6 # Linear regression is an effective ML-model to predict unknown variable when the unknown variable varies with some
7 # constant proportionality with the known variable(s), whether positive or negative.
8 #
9 # Using this particular piece of code,
10 # in main() call the class linear_reg with args <fileName>, <Known Variable>, <Unknown variable>,
11 # <Title for visualization>, <Estimation of missing data points, T/F ?>
12 # Enjoy all necessary insights about data, with beautiful visualizations.
13
14 import numpy as np
15 import pandas as pd
16 import matplotlib.pyplot as plt
17
18 class linear_reg:
19
20     def __init__(self, fileName, Xcolumn, Ycolumn, title, estReq=False): #Initializes all the necessary variables to keep track of model.
21         self.estimationRequired = estReq
22         self.title = title
23         self.dataset = pd.read_csv(fileName)
24         self.Xcolumn = Xcolumn
25         self.Ycolumn = Ycolumn
26         self.yPred = []
27         self.bestFit = {}
28         self.X = np.array([])
29         self.Y = np.array([])
30         self.X_test = np.array([])
31         self.Y_test = np.array([])
32         self.yPred_all = np.array([])
33         self.acc = 0
34
35     def preprocess(self): #Substitutes missing data in Independent column, making it ready to be used by model
36         self.dataset[self.Xcolumn].fillna(value=self.dataset[self.Xcolumn].mean(), inplace=True)
37
38     def split(self): #Splits the dataset for training and testing purpose
39         self.X = np.array(self.dataset[self.Xcolumn][:200])
40         self.Y = np.array(self.dataset[self.Ycolumn][:200])
41
42         self.X_test = np.array(self.dataset[self.Xcolumn][200:])
43         self.Y_test = np.array(self.dataset[self.Ycolumn][200:])
44
45     def newMatrix_noNaN(self): #Mimics the dataset but with no missing data.
46         data = self.dataset.dropna()
47         self.Y = np.array(data['Salary'])
48         self.X = np.array(data['Age'])
49         self.X_test = np.array([x for x in self.dataset['Age'] if x not in self.X])
50
51     def apply_linear_regression(self): #Applies linear regression and fetches the best-fit eq parameters in dictionary.
52         x_mean, y_mean = np.mean(self.X), np.mean(self.Y)
53         x_x, y_y = self.X-x_mean, self.Y-y_mean
54         x_ssd = np.array(np.square(x_x))
55         yDiff_xDiff = np.array(np.multiply(x_x, y_y))
56         coefficient = np.sum(yDiff_xDiff) / np.sum(x_ssd)
57         intercept = y_mean - coefficient*x_mean
58         self.bestFit = {'coefficient': coefficient, 'intercept': intercept}
59
60     def generate_yPred(self): #Generates predicted-Y using best-fit eq generated via linear regression.
61         self.yPred = [(self.bestFit['coefficient']*x+self.bestFit['intercept']) for x in self.X_test]
62         self.yPred_all = [(self.bestFit['coefficient']*x+self.bestFit['intercept']) for x in self.X]
63         print(self.yPred)
64
65     def calculate_accuracy(self): #Calculates accuracy of prediction of the model.
66         y_error = [abs(y_test-y_pred)/y_test for y_test,y_pred in zip(self.Y_test, self.yPred)]
67         self.acc = (100-sum(y_error)/len(y_error)*100)
68
69     def trainAndTest(self): #Calls all necessary functions in proper heirarchy for smooth control flow and convenience, hehe:)
70         self.split()
71         self.apply_linear_regression()
72         self.generate_yPred()
73         self.calculate_accuracy()
74
75     def estimateMissing(self): #Calls all necessary functions in defined heirarchy to predict and fill missing values in dataset too.
76         self.preprocess()
77         self.newMatrix_noNaN()
78         self.apply_linear_regression()
79         self.generate_yPred()
```

```
79
80 def visualize(packet): #Visualizes the linear model, with regression line, predicted data-points, given data-points, and noise cloud.
81     fig, axs = plt.subplots(2)
82
83     for plot_id, self in enumerate(packet):
84         keyContent = 'Accuracy: {:.2f}%'.format(self.acc)+self.title
85         fig.patch.set_facecolor('xkcd:mint green')
86         axs[plot_id].text(0.05, 0.95, keyContent, transform=axs[plot_id].transAxes, fontsize=10, bbox=dict(facecolor='w', alpha=0.5), verticalalign=
87         axs[plot_id].set_facecolor('k')
88         axs[plot_id].plot(self.X test, self.yPred, color='red')
89         axs[plot_id].plot(self.X, self.yPred all, color='red')
90         axs[plot_id].scatter(self.X, self.Y, c='#32CD32', s=40, marker="*", alpha = '0.6')
91         if self.estimateRequired:
92             for x,y in zip(self.X test, self.yPred):
93                 rgb = (np.random.random(), np.random.random(), np.random.random())
94                 axs[plot_id].scatter(x, y, c=rgb, s=75, marker="^", alpha=1)
95                 axs[plot_id].text(x, y-8000, 'Estimated for missing pair:\n{:.2f} -> {:.2f}'.format(x, y), fontsize=10, bbox=dict(facecolor='w', alp
96             else:
97                 for x,y in zip(self.X test, self.yPred):
98                     rgb = (np.random.random(), np.random.random(), np.random.random())
99                     axs[plot_id].scatter(x, _y, c='b', s=20, marker="o", alpha=0.5)
100                     axs[plot_id].scatter(x, y, c=rgb, s=75, marker="^", alpha=1)
101                     axs[plot_id].plot([x, x], [_y, y], color=rgb)
102
103     for ax,self in zip(axs.flat, packet):
104         ax.set(xlabel= self.Xcolumn, ylabel= self.Ycolumn)
105
106     fig.suptitle('Linear Regression', fontsize=20, bbox=dict(facecolor='m', alpha=0.6))
107     fig.tight_layout()
108     plt.show()
109
110 def main(): #Main funtion, basically the linear-regression applicable dataset should be called as memeber of class linear-reg for model to initiate.
111     HeadBrain = linear_reg('headbrain.csv', 'Head Size(cm^3)', 'Brain Weight(grams)', '\n* -TrainData\n^ -PredictedData\nRed | -Regression Line\n
112     HeadBrain.trainAndTest()
113     AgeSalary = linear_reg('null.bin', 'Age', 'Salary', 'Red | - Regression line\n* -TrainData\n^ -PredictedData', True)
114     AgeSalary.estimateMissing()
115     visualize([HeadBrain, AgeSalary]) #Calls visualizer to plot every insight gained from dataset.
116
117 if __name__ == "__main__":
118     main()
119
```

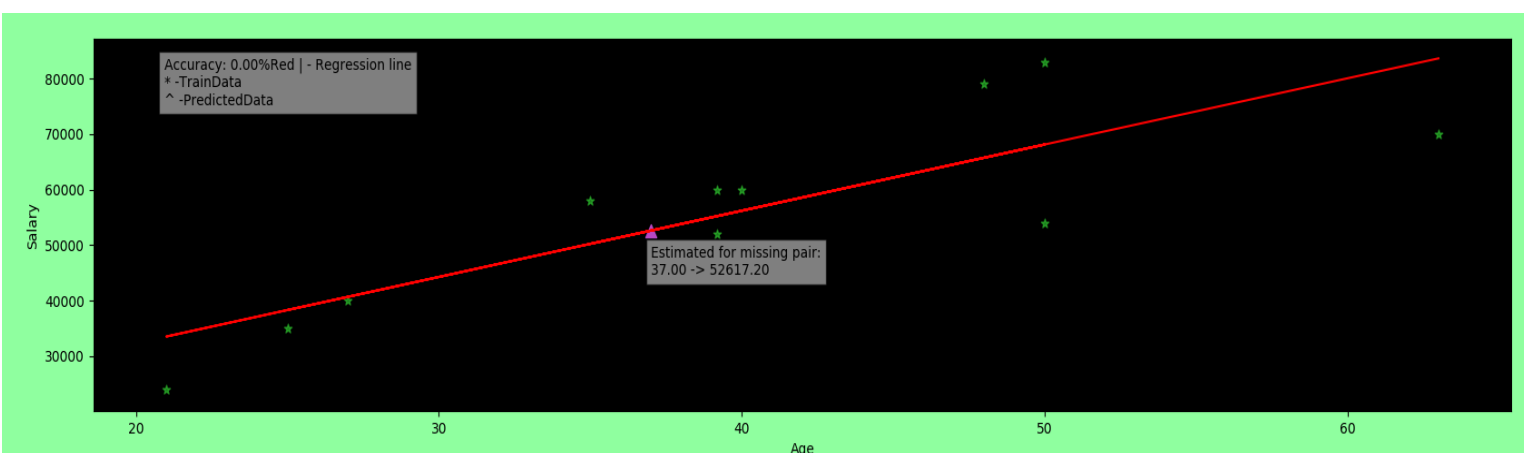
Output of the model with visualization, in next page >>



Zooming in | Predicts Brain-Weight on the basis of Head-Size.
| **COLORFUL** data-points are predicted data.
| Lines connecting them represent the **error**.
| **Green** data points are plots of data in data-set.



Zooming into age-salary



Dibyasom Puhan
SAP id: 500076104

Date: 24th August, 2020
Roll no.: R177219075