

# CSCI 720 Big Data Analytics HW03 Results

---

Student: Guo, Zizhun

Email: [zg2808@cs.rit.edu](mailto:zg2808@cs.rit.edu)

Phone: 585-284-0464

Submission: Feb/15/2020

Due Date: Feb/16/2020 11:59 PM

---

## 1. Did your mentor program use pandas to help you along?

Yes.

```
1 | # location: HW_03_Guo_Zizhun_Mentor.py
2 |
3 | import pandas as pd # data analysis and manipulation package
4 |
```

## 2. Did your mentor program create your training program?

Yes.

```
1 | # location: HW_03_Guo_Zizhun_Mentor.py
2 |
3 | def write_file(best_attribute, is_positive_correlated):
4 |     """ writes trained program based on selected best attribute
5 |     Paras:
6 |         @best_attribute: a string represents the name of best attribute
7 |         @is_positive_correlated: a boolean determines the classification rule
8 |     Return:
9 |         void
10 |    """
```

```

1 | # initialize a string to contains codes for dumping in the trained program
2 |     lines = ""
3 |         ..
4 |         .. implementations ..
5 |         ..
6 |     # create a file with given name "filepath"
7 |     f = open('HW03_Guo_Zizhun_Trained.py', "w")
8 |     # write string to the filepath
9 |     f.write(lines)
10 |    f.close()

```

### 3. Did the training program use pandas?

Yes.

```

1 |     # location: HW_03_Guo_Zizhun_Trained.py
2 |
3 |     import pandas as pd
4 |

```

### 4. Which attribute is most strongly cross-correlated with the target variable?

Attribute	Correlation with Target variable
Bread	-0.012
Vitamins	-0.448
Vegetable	0.016
Milk	-0.035
Banana	-0.070
PeanutButter	0.582
Chocolate	-0.030
Citrus	-0.057
Cookie	0.050
IceCream	0.013
Soda	-0.061
Fruit	-0.005

**PeanutButter** has the highest cross correlation absolute value, which means it is the most strongly cross-correlated attribute with the target variable.

## 5. Which is best feature out of the above items for your One Rule?

PeanutButter.

## 6. What structure did your One Rule classifier have?

The if-statement condition sets the structure based on whether the **value of cross-correlation** is **positive** or **negative**.

**Sets structure for One Rule classifier in trained program:**

```
1 | # location: HW_03_Guo_Zizhun_Mentor.py
2 | if is_positive_correlated:
3 |
4 |     lines += f"\n\tfor val in data:"
5 |     lines += f"\n\t\tif val > 0:"
6 |     lines += f"\n\t\t\ttprint('1')"
7 |     lines += f"\n\t\telse:"
8 |     lines += f"\n\t\t\ttprint('0')"
9 | else:
10 |     lines += f"\n\tfor val in data:"
11 |     lines += f"\n\t\tif val == 0:"
12 |     lines += f"\n\t\t\ttprint('1')"
13 |     lines += f"\n\t\telse:"
14 |     lines += f"\n\t\t\ttprint('0')"
15 |
```

### Q1: What was the if-else rule you got?

**Struture below:** (based on the given training dataset)

```
1 | # location: HW_03_Guo_Zizhun_Trained.py
2 |
3 | if attribute > 0:
4 |     print('1')
5 | else:
6 |     print('0')
7 |
```

## 7. Run the original training data back through your classifier.

**Q1: What was the accuracy of your resulting classifier, on the training data?**

The Trained Program would print the result as the homework asked, but the accuracy can be captured by **frequency table** created in Mentor Program:

```
1 | # shell console
2 |
3 | PeanutButter  Sickness
4 | 0             0         399
5 |               1         108
6 | 1             0         101
7 |               1         392
8 | dtype: int64
9 |
```

$numberOfCorrectness = numberOf(PeanutButter : 0, Sickness : 0) +$

$numberOf(PeanutButter : 1, Sickness : 1) = 399 + 392 = 791$

$numberOfTotal = 1000$

$Accuracy = numberOfCorrectness / numberOfTotal = 791 / 1000 = 0.79$

The **Accuracy** is **0.79**.

\*Reference:

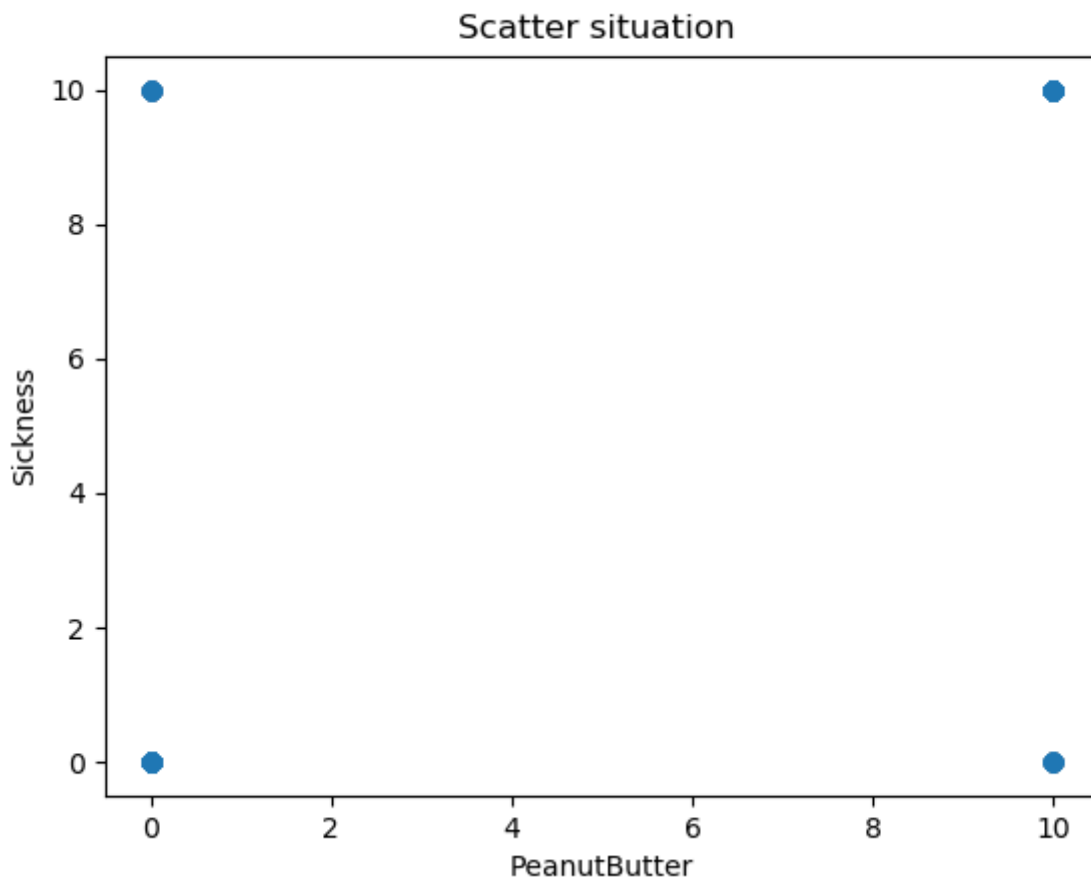
Another way to define Accuracy:  $ACC = (TP + TN) / (TP + TN + FP + FN)$

## 8. Conclusion: (2)

### Q1: What else did you learn along the way here?

I also learned using scatter plot to present the correlation situation.

The default image would be looks like this:



**Image 1: The 2D Sickness based on Peanutbutter**

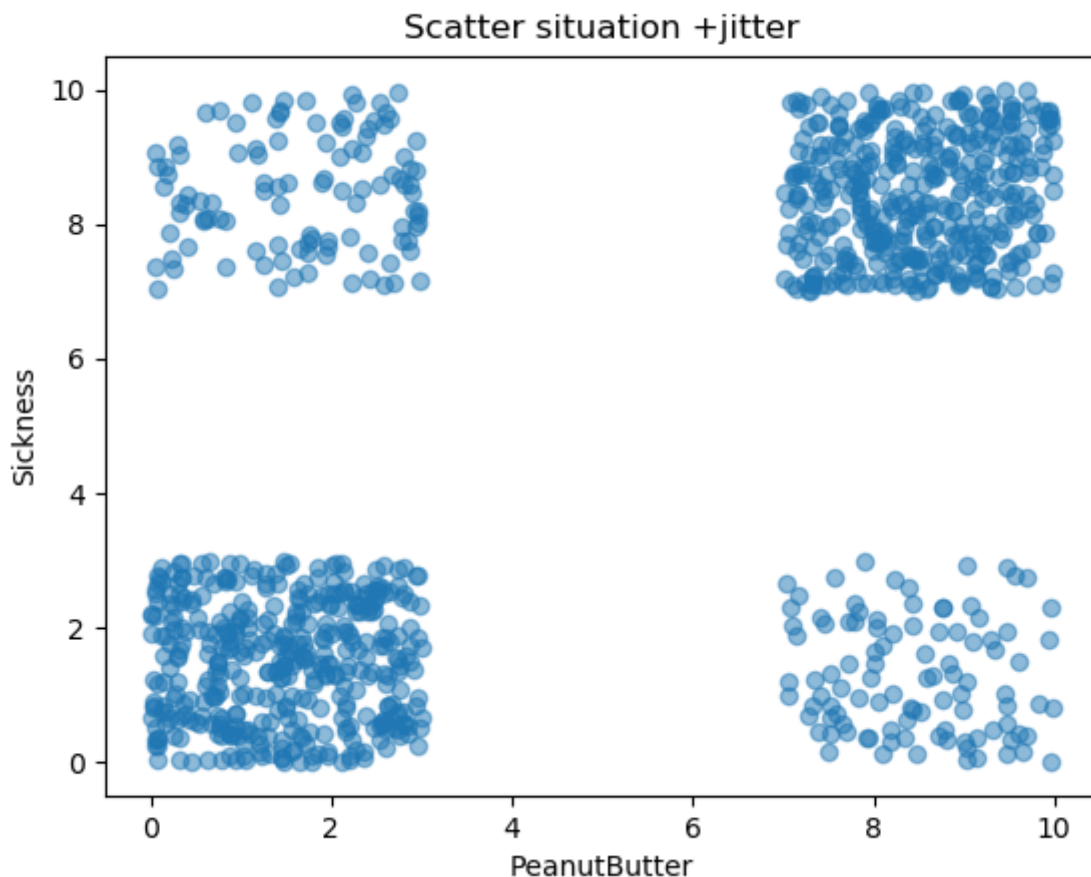
Each dot actually overlaps so many points. In order to see how condense each points are, we need to add jitter to scatter it up a little bit. So just choose a randomly scattered seed generated by Numpy package.

```

1      # location: HW_03_Guo_Zizhun_Mentor/ditribution_plot_with_jitter.py
2
3      x = df_filtered[best_attribute] \
4          * (1 - scatter_fraction_rate) \
5          * scatter_scale \
6          + np.random.ranf(size) \
7          * scatter_fraction_rate \
8          * scatter_scale
9
10     ...
11     ...
12     plt.scatter(x, y, alpha=0.5)      # plot the image alpha = 0.5 indicates the level of
13                                     # Area with darker color have more points overlapped

```

After adding the **Jitter**:



**Image 2: The 2D Sickness based on Peanutbutter with Jitter**

Here **Jitter** as the **Noise** added enabling easy observation on how the points are settled. As can be seen, points from **lower left** and **upper right** are largely placed, which means the two x-variable and y-variable are **strongly positive cross-correlated**.

**This tendency implies:** The y-variable tends to have the same value of x\_variable, since in our case, we only have two value options (0 and 1).

## Q2: What can you conclude?

The tendency for attribute values to be changed as the target variables is defined by the **sign** of **cross-correlation**. This conclusion can be used to determine the One Rule of trained programs that if the cross-correlation(CC) is less than 0, the prediction should be **opposite** as the value of attribute, whereas if CC is greater than 0, the prediction should be **same** result of the attribute value.

## 9. BONUS: (+1)

**Q1: In addition to the best feature, can you find another feature that might help classify sickness?**

Yes. **Vitamin** can be this feature, since

## Q2: What makes you say this?

Because it has the second highest cross-correlation, which is **-0.448**.

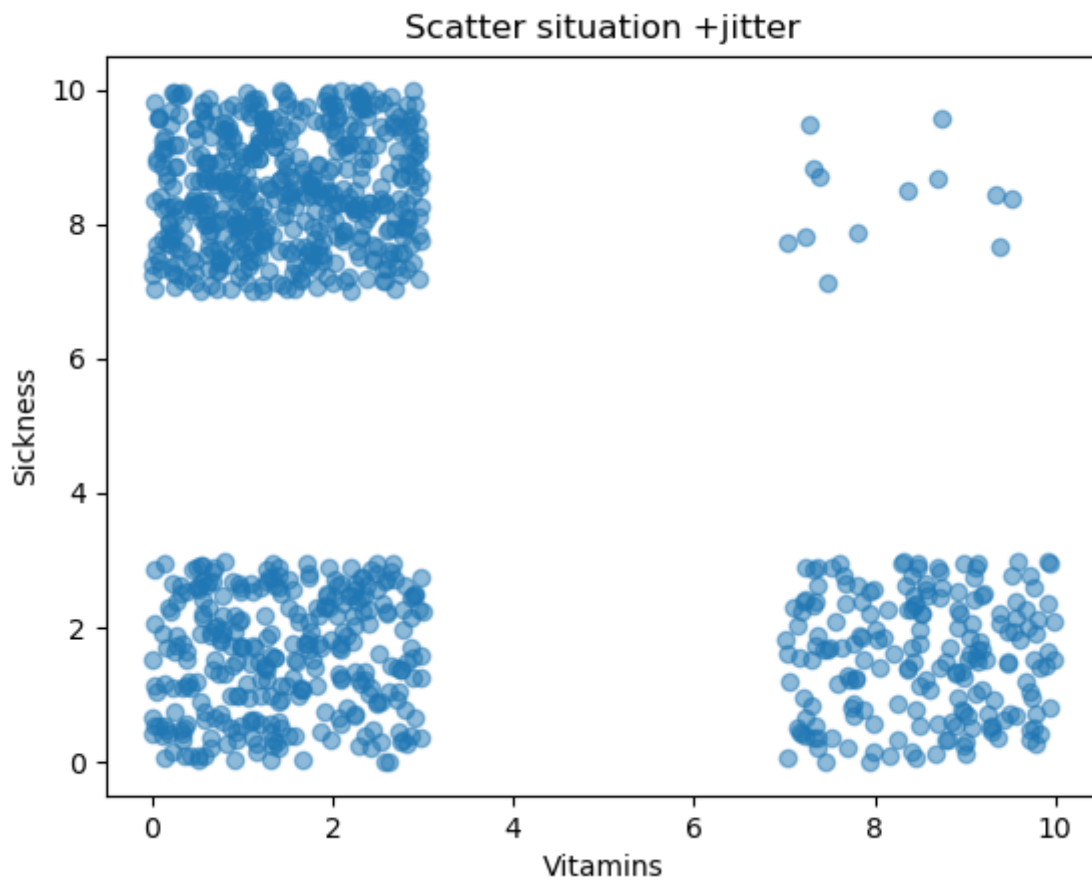


Image 3: The 2D Sickness based on Vitamin with Jitter

## Q3: How did you test this hypothesis?

I can print the frequency table to calculate the Accuracy.

Vitamins	Sickness	#
0	0	305
	1	487
1	0	195
	1	13
		total: 1000

$$\text{Accuracy} = (487+195)/1000 = \mathbf{0.682}$$

\*Since the correlation is less than 0, it needs to sum up the # of {0, 1} and {1, 0} pairs.

