

HAP 618 – Fall xxx

Final Project Paper

Breast Cancer

Background and problem description

Breast cancer is the 2nd most common type of cancer in the United States, with estimated new cases of 255,000 and estimated deaths of 41,000 for both male and female every year according to cdc.gov. These numbers can be reduced with several different methods. One of the many different methods is identifying if the tumor type is benign or malignant.

There are two types of tumor and they are benign and malignant. Benign is non-cancerous, it is a growth that is not cancer. It does not invade nearby tissue or spread to other parts of the body (cancer.gov). However, malignant is cancerous and it can invade and destroy nearby tissue and spread to other parts of the body (cancer.gov).

The problem that this project is identifying is, if there is a correlation between the tumor type and the nine attributes (independent variables) in this project. Understanding the correlation between the tumor type and the different attributes could potentially help identify the type of tumor more accurately. But also, analyze the data to understand what attributes correlate the most or least with the tumor type. According to New England Journal of Medicine, experts believe as many as 31 percent of all breast cancer cases are misdiagnosed, and many as 90,00 women are misdiagnosed with invasive breast cancer (The New York Times). Those are huge percentage and number of women's misdiagnosed.

The problem could be accomplished by using Unix system, Python, and HTML. The Unix system will mainly serve as the platform to storing all the data, Python will do the main job of aggregating the data, and HTML will display the graph and information nicely on the web. The dataset I will be using is from UCI and it is called *Breast Cancer Wisconsin (Original)* dataset.

Design of the solution:

For the project the codes that were used were; Python, Unix, and HTML. All three of these codes intertwine within each other to bring the project together. All three of these codes could have been executed on either server or on a client's computer, but for this project majority of the code is on the Unix server. There are few Python codes that are on my personal device. The project is not interactive, it consists of static webpages, except the bar charts refreshes with updated data. It was mainly designed for understanding the breast cancer data, analyzing the outcomes, correlation between tumor type and the attributes for breast cancer, and prediction model.

Python was used to aggregate the data to create bar charts and to create prediction model. The libraries that were used in the Python program are pandas, matplotlib.pyplot, and sklearn (sci-kit). The pandas library was used to import the breastcancer.csv file. The matplotlib.pyplot was used to create the bar charts and adding small detailed features such as; x-label, y-label, title, legend, assigning colors to specific class, and the x-tick marks. The graph is available for the users to see and analyze. Lastly, sklearn was used to create a prediction model by getting prediction probability on test data, and getting prediction accuracy score.

The Unix system was used to store the data and code (breast cancer data, html code, and python code). It is also responsible for running all the codes as it provided the support of displaying the result on the internet browser and creating the bar charts. The breastcancer.csv file was stored in the user's home of the Unix system. All the HTML codes were saved in /public_html and all Python codes were saved in /cgi-bin. Python code was used to read the breastcancer.csv file. Once the file was opened, it was used to aggregate the data, create bar charts, and save the bar chart as image.png in /public_html this gave the

image an URL. HTML code was used to create a web page and used the image URL to display the graph accordingly. In whole, Unix system took all the different parts of the code and displayed it as one.

HTML code, as stated before was used mainly to display the information and graph accordingly. HTML code could have been executed from the server or on a client computer, but in this project the code was executed from the server.

Implementation:

All the programs that were used in this project are important, being said that, without one it is not a complete project. However, if there is one important program that must be identified it would be the Python codes for the bar charts. The bar charts are used to analyze the data which is the main topic of the project. The following is one of the codes for bar chart.

```

110 import cgi
111 import pandas as pd
112 import matplotlib.pyplot as plt
113
114 mydata = pd.read_csv('/home/[redacted]/breastcancer.csv')
115
116 #####
117 class_2 = mydata[mydata['Class'] == 2]
118 class_4 = mydata[mydata['Class'] == 4]
119 #####
120 c2_uni_of_cellshape = class_2['Uniformity_of_Cell_Shape'].value_counts().sort_index()
121 cnt_uni_of_cellshape_c2 = c2_uni_of_cellshape.size
122 c4_uni_of_cellshape = class_4['Uniformity_of_Cell_Shape'].value_counts().sort_index()
123 cnt_uni_of_cellshape_c4 = c4_uni_of_cellshape.size
124 #####
125 def uni_of_cellshape():
126     uni_of_cellshape = mydata['Uniformity_of_Cell_Shape'].value_counts().sort_index()
127     cnt_uni_of_cellshape = uni_of_cellshape.size
128     ct_labels = ['1','2','3','4','5','6','7','8','9','10']
129     plt.bar(range(cnt_uni_of_cellshape_c2), list(c2_uni_of_cellshape), color = 'orange', label = 'Benign')
130     plt.bar(range(cnt_uni_of_cellshape_c4), list(c4_uni_of_cellshape), label = 'Malignant')
131     plt.xticks(range(cnt_uni_of_cellshape), ct_labels)
132     plt.title('Number of Instances by Uniformity of Cell Shape and Class')
133     plt.xlabel('Uniformity of Cell Shape')
134     plt.ylabel('Number of Instances')
135     plt.legend()
136     plt.savefig('/home/[redacted]/public_html/uni_of_cellshape.png')
137
138 uni_of_cellshape()

```

The following is the step-by-step on how the Python code works:

Line 110: Imports CGI

Line 111: Import pandas – used to import the breastcancer.csv

Line 112: Import matplotlib.pyplot – used to create the bar chart

Line 114: Reading the breastcancer.csv using pandas library

Line 117: Selecting the data only where Class = 2

Line 118: Selecting the data only where Class = 4

Line 120: Selecting only data from the column “Uniformity_of_Cell_Shape” and counting the values and sorting it by the index for data in class = 2.

Line 121: Counting the number of distinct values

Line 122: Selecting only data from the column “Uniformity_of_Cell_Shape” and counting the values and sorting it by the index for data in class = 4

Line 123: Counting the number of distinct values

Line 125: Defining a function “uni_of_cellshape”

Line 126 to Line 136: Is creating the bar chart including; xticks, title, x-label, y-label, legend, and saving the result as a image.png

Line 138: Printing the function “uni_of_cellshape”

Other Python codes that were implemented in this project is the prediction model codes. These Python codes were used to create a prediction model using a data mining method. This was an interesting perspective of the data, and further helps understand the significance in attribute type and result correlation to tumor type.

The HTML code was utilized to display the information and graphs and charts on the web, it was also used to make the web page look nice and organized. The HTML code was implemented towards bringing all the pieces of the project together into one web page. Also, how the images were saved as .png in html folder allowed flexibility to display the image on the web.

As it has been previously mentioned before, the implementation of the Unix system was to store the codes and run the codes. The Unix system also was used to store Python result for bar charts, which was in the format of .png, into HTML file. This easily made it possible to display the bar chart on the HTML web page with correct .png HTML URL.

Testing:

My project consists of four HTML pages. Majority of the codes are stored within Unix system and is running from the Unix system. Only codes that are not in the Unix system is prediction model codes.

Breast Cancer Dataset the highlighted “breastcancer.csv”:

a) Unix:

```
[~]@students ~]$ ls
anotherfile      assignment4-a.file.txt  breastcancer.csv  fall17  hap618  okay  python
assignformdata.txt  bfile                  class.png         fall618  homedir.txt  public_html  python.py
[~]@students ~]$
```

b) Content of the file breastcancer.csv:

```
1369821,10,10,10,10,5,10,10,10,7,4
1371026,5,10,10,10,4,10,5,6,3,4
1371920,5,1,1,1,2,1,3,2,1,2
466906,1,1,1,1,2,1,1,1,1,2
466906,1,1,1,1,2,1,1,1,1,2
534555,1,1,1,1,2,1,1,1,1,2
536708,1,1,1,1,2,1,1,1,1,2
566346,3,1,1,1,2,1,2,3,1,2
603148,4,1,1,1,2,1,1,1,1,2
654546,1,1,1,1,2,1,1,1,8,2
654546,1,1,1,3,2,1,1,1,1,2
695091,5,10,10,5,4,5,4,4,1,4
714039,3,1,1,1,2,1,1,1,1,2
763235,3,1,1,1,2,1,2,1,2,2
776715,3,1,1,1,3,2,1,1,1,2
841769,2,1,1,1,2,1,1,1,1,2
888820,5,10,10,3,7,3,8,10,2,4
897471,4,8,6,4,3,4,10,6,1,4
897471,4,8,8,5,4,5,10,4,1,4
[~]@students ~]$ cat breastcancer.csv
```

Navigation:

- a) Code: The following codes are identical in all three of my separate html web page.

```
<div class="navbar">
  <h2>Breast Cancer Data Anlysis</h2>
  <a href="http://students.hi.gmu.edu/[redacted]fp_about_the_data" class="button">About the Data</a>
  <a href="http://students.hi.gmu.edu/[redacted]fp_overview_of_the_data" class="button">Overview of the Data</a>
  <a href="http://students.hi.gmu.edu/[redacted]fp_visualization_of_data" class="button">Visualization of the Data</a>
  <a href="http://students.hi.gmu.edu/[redacted]fp_data_mining" class="button">Prediction Model</a>
<br>
<span style="font-size:10px">HAP 618 Final Project by [redacted]/span>
</div>
```

- b) Result (output):

Breast Cancer Data Anlysis



Bar Graph:

- a) Python Code: Used the following code to create graph from the data. This python code format/layout was used to generate all the bar charts. Within the Unix system, the Python code for each graph had to be individually and separately saved onto different .cgi files. Otherwise, if more than one graph was in one .cgi file, the graph would overlap against each other. Therefore, I created total of 10 .cgi files (Number of Instances by Class + 9 different graphs for each attributes) for the bar charts.

```
110 import cgi
111 import pandas as pd
112 import matplotlib.pyplot as plt
113
114 mydata = pd.read_csv('/home/[redacted]breastcancer.csv')
115
116 #####
117 class_2 = mydata[mydata['Class'] == 2]
118 class_4 = mydata[mydata['Class'] == 4]
119 #####
120 c2_uni_of_cellshape = class_2['Uniformity_of_Cell_Shape'].value_counts().sort_index()
121 cnt_uni_of_cellshape_c2 = c2_uni_of_cellshape.size
122 c4_uni_of_cellshape = class_4['Uniformity_of_Cell_Shape'].value_counts().sort_index()
123 cnt_uni_of_cellshape_c4 = c4_uni_of_cellshape.size
124 #####
125 def uni_of_cellshape():
126     uni_of_cellshape = mydata['Uniformity_of_Cell_Shape'].value_counts().sort_index()
127     cnt_uni_of_cellshape = uni_of_cellshape.size
128     ct_labels = ['1','2','3','4','5','6','7','8','9','10']
129     plt.bar(range(cnt_uni_of_cellshape_c2), list(c2_uni_of_cellshape), color = 'orange', label = 'Benign')
130     plt.bar(range(cnt_uni_of_cellshape_c4), list(c4_uni_of_cellshape), label = 'Malignant')
131     plt.xticks(range(cnt_uni_of_cellshape), ct_labels)
132     plt.title('Number of Instances by Uniformity of Cell Shape and Class')
133     plt.xlabel('Uniformity of Cell Shape')
134     plt.ylabel('Number of Instances')
135     plt.legend()
136     plt.savefig('/home/[redacted]public_html/uni_of_cellshape.png')
137
138 uni_of_cellshape()
```

About the Data:

a) HTML Code: The following is part of the code, which consists of the code for the table.

```
183 <div class="main">
184 About the Dataset:
185 <br><table style="table-layout: fixed; border: 0.5px solid black;">
186 <tr>
187   <td><font size="2"><b>Data Set:</b></td><td><font size="2">Breast Cancer Wisconsin (Original)</td>
188 </tr>
189 <tr>
190   <td><font size="2"><b>Number of Instances:</b></td><td><font size="2">699</td>
191 </tr>
192 <tr>
193   <td><font size="2"><b>Number of Attributes:</b></td><td><font size="2">10</td>
194 </td>
195 </table>
196 <br>
197 <table style="table-layout:fixed; border: 0.5px solid black;">
198 Attribute Information:
199 <tr>
200   <td><font size="2"><b>Sample code number:</b></td><td><font size="2">Id number</td>
201 </tr>
202 <tr>
203   <td><font size="2"><b>Clump Thickness:</b></td><td><font size="2">1-10</td>
204 </tr>
205 <tr>
206   <td><font size="2"><b>Uniformity of Cell Size:</b></td><td><font size="2">1-10</td>
207 </tr>
208 <tr>
209   <td><font size="2"><b>Uniformity of Cell Shape:</b></td><td><font size="2">1-10</td>
210 </tr>
211 <tr>
212   <td><font size="2"><b>Marginal Adhesion:</b></td><td><font size="2">1-10</td>
213 </tr>
214 <tr>
215   <td><font size="2"><b>Single Epithelial Cell Size:</b></td><td><font size="2">1-10</td>
216 </tr>
217 <tr>
218   <td><font size="2"><b>Bare Nuclei:</b></td><td><font size="2">1-10</td>
219 </tr>
```

b) Result (output): http://students.hi.gmu.edu/~xxx/fp_about_the_data

About the Dataset:	
Data Set:	Breast Cancer Wisconsin (Original)
Number of Instances:	699
Number of Attributes:	10

Attribute Information:	
Sample code number:	Id number
Clump Thickness:	1-10
Uniformity of Cell Size:	1-10
Uniformity of Cell Shape:	1-10
Marginal Adhesion:	1-10
Single Epithelial Cell Size:	1-10
Bare Nuclei:	1-10
Bland Chromatin:	1-10
Normal Nucleoli:	1-10
Mitoses:	1-10
Class:	2 for benign 4 for malignant

Citation:

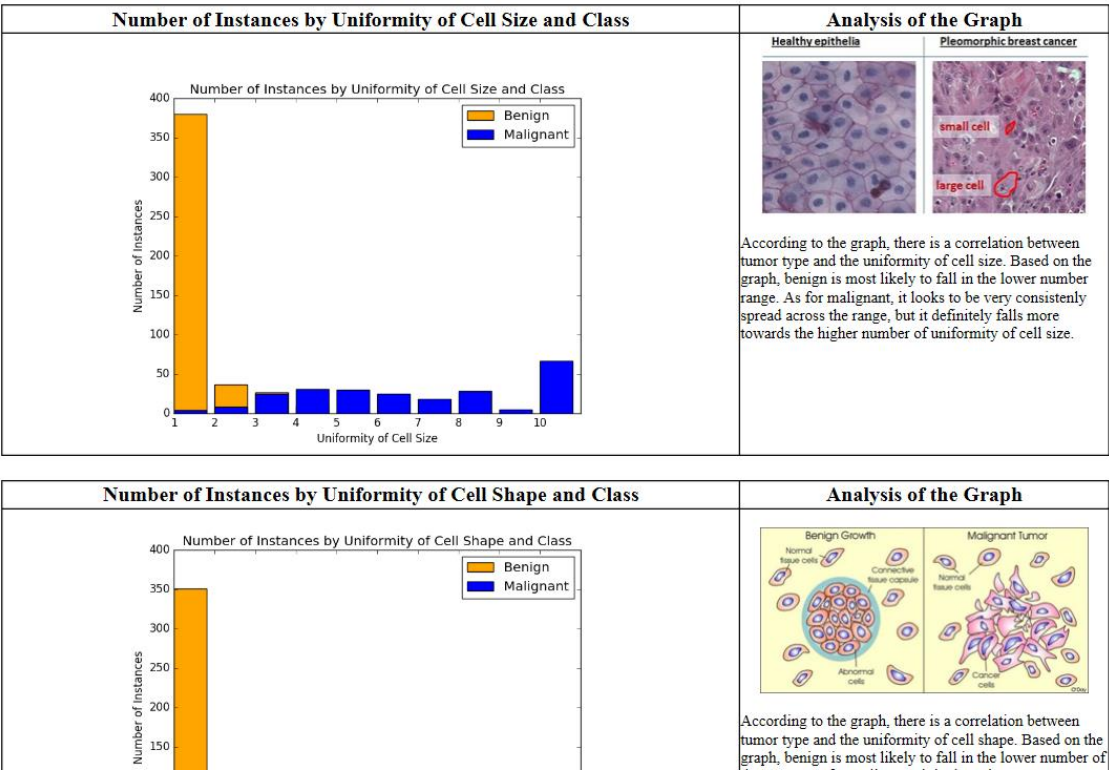
1. O. L. Mangasarian and W. H. Wolberg: "Cancer diagnosis via linear programming", SIAM News, Volume 23, Number 5, September 1990, pp 1 & 18.
2. William H. Wolberg and O.L. Mangasarian: "Multisurface method of pattern separation for medical diagnosis applied to breast cytology", Proceedings of the National Academy of Sciences, U.S.A., Volume 87, December 1990, pp 9193-9196.
3. O. L. Mangasarian, R. Setiono, and W.H. Wolberg: "Pattern recognition via linear programming: Theory and application to medical diagnosis", in: "Large-scale numerical optimization", Thomas F. Coleman and Yuying Li, editors, SIAM Publications, Philadelphia 1990, pp 22-30.
4. K. P. Bennett & O. L. Mangasarian: "Robust linear programming discrimination of two linearly inseparable sets", Optimization Methods and Software 1, 1992, 23-34 (Gordon & Breach Science Publishers).

Visualization of the Data:

a) HTML Code:

```
2 <style>
3
4 table {
5   border-collapse: collapse;
6   width: 65%;
7 }
8 td, th {
9   border: 1px solid black;
10 }
11 .textNode{
12   text-align:left;
13   padding:0;
14   margin:0;
15   vertical-align:top;
16   font-size: 12px;
17   word-wrap: break-word;
18 }
19 .button{
20   background-color: #555555; /* Green */
21   border: none;
22   color: white;
23   padding: 10px 30px;
24   text-align: center;
25   text-decoration: none;
26   display: inline-block;
27   font-size: 14px;
28   margin: 4px 2px;
29   cursor: pointer;
30 }
31 .navbar {
32   overflow: hidden;
33   background-color: #FFFFFF;
34   position: fixed;
35   top: 0;
36   width: 100%;
37 }
38 .main {
39   padding: 100px;
40   margin-top: 30px;
41 }
42 </style>
43
44 <body>
45 <center>
46
47 <div class="navbar">
48   <h2>Breast Cancer Data Analysis</h2>
49   <a href="http://students.hi.gmu.edu/~xxx/fp_draft4" class="button">About the Data</a>
50   <a href="http://students.hi.gmu.edu/~xxx/fp_draft3" class="button">Overview of the Data</a>
51   <a href="http://students.hi.gmu.edu/~xxx/fp_draft2" class="button">Visualization of Data</a>
52 </div>
53 <span style="font-size:10px">HAP 618 Final Project by <span>
54 </div>
55
56 <div class="main">
57 <table style="table-layout: fixed;">
58 <tr>
59 <th colspan="3"> Number of Instance by Class </th>
60 </tr>
61 <tr>
62 <td colspan="3"><b>Benign</b></td>
63 <td colspan="3"><b>Malignant</b></td>
64 <td colspan="3"><b>Analysis on Graph</b></td>
65 </tr>
66 <tr>
67 <td colspan="3"><b>Benign</b></td>
68 <td colspan="3"><b>Malignant</b></td>
69 <td colspan="3"><b>Analysis on Graph</b></td>
70 </tr>
71 <tr>
72 <td colspan="3"></td>
73 <td colspan="3"></td>
74 <td colspan="3"></td>
75 </tr>
76 </table>
77 <div class="textnode">
78 <table style="table-layout: fixed;">
79 <tr>
80 <th colspan="2"> Number of Instances by Clump Thickness and Class</th>
81 <th colspan="2"><b>Benign</b></th>
82 <th colspan="2"><b>Malignant</b></th>
83 </tr>
84 <tr>
85 <td colspan="2"></td>
86 <td colspan="2"></td>
87 <td colspan="2"></td>
88 <td colspan="2"></td>
89 </tr>
90 </table>
91 </div>
92 </body>
93 </html>
```

b) Result (output): http://students.hi.gmu.edu/~xxx/fp_visualization_of_data



Overview of the data:

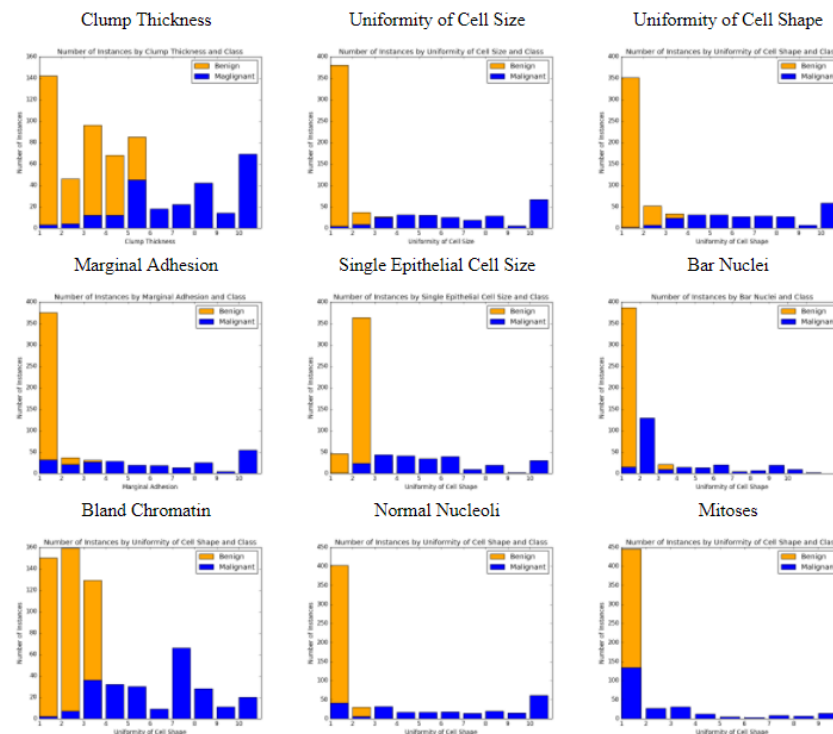
a) HTML Code: Arranged the graphs using the tables html code.

```

141 <div class="main">
142 <table style="table-layout: fixed;">
143 <tr>
144 <td><font size="2"><center>Clump Thickness</center></td>
145 <td><font size="2"><center>Uniformity of Cell Size</center></td>
146 <td><font size="2"><center>Uniformity of Cell Shape</center></td>
147 </tr>
148 <tr>
149 <td>
150 <td>
151 <td>
152 </tr>
153 <tr>
154 <td><font size="2"><center>Marginal Adhesion</center></td>
155 <td><font size="2"><center>Single Epithelial Cell Size</center></td>
156 <td><font size="2"><center>Bar Nuclei</center></td>
157 </tr>
158 <tr>
159 <td>
160 <td>
161 <td>
162 </tr>
163 <tr>
164 <td><font size="2"><center>Bland Chromatin</center></td>
165 <td><font size="2"><center>Normal Nucleoli</center></td>
166 <td><font size="2"><center>Mitoses</center></td>
167 </tr>
168 <tr>
169 <td>
170 <td>
171 <td>
172 </tr>
173 </table>
174 </center>
175 </div>
176 </body>
177 </html>

```

b) Result (output): http://students.hi.gmu.edu/~xxx/fp_overview_of_the_data



Prediction Model

a) Python Code

```
from sklearn.metrics import accuracy_score
from sklearn.naive_bayes import GaussianNB
import matplotlib.pyplot as plt

size = y_test.size

# Initialize our classifier
gnb = GaussianNB()
# Train our classifier
model = gnb.fit(x_train, y_train)
# Make predictions
preds = gnb.predict(x_test)
preds2 = gnb.predict_proba(x_test)
print(preds2)
print(accuracy_score(y_test, preds))

##### Plot GaussianNB

x = [i[0] for i in preds2]
y = [i[1] for i in preds2]

plt.scatter(x,y, label = 'Prediction Probability')
plt.title('GaussianNB Prediction Probability')
plt.xlabel('Probability')
plt.ylabel('Probability')

plt.scatter(x,y, label = 'Prediction Probability')
plt.scatter(range(size), y_test, alpha = 0.4, label = 'Actual Result')
plt.text(160,0.35, 'Accuracy Score: 0.9619')
plt.rcParams["figure.figsize"] = (10,5)
plt.legend()
plt.title('GaussianNB')
plt.ylabel('Prediction Probability')
plt.xlabel('Number of Cases')
plt.xticks()
plt.show()
```

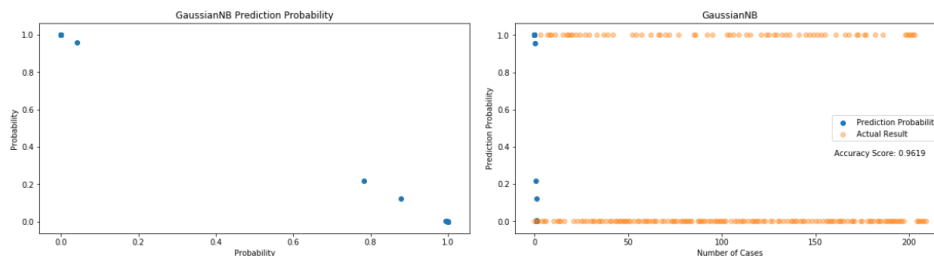
b) HTML Code:

```
<div class="main">
<center><b><span style="font-size:26px">Prediction Model</span></b></center>
<table style="table-layout: fixed;">
<tr>
<td></td>
<td></td>
</tr>
<br>
<tr>
<td></td>
<td></td>
</tr>
<br>
<tr>
<td></td>
<td></td>
</tr>
</table>
</div>
```

- The images were saved on my personal laptop and uploaded on the web.

c) Result (output): http://students.hi.gmu.edu/~xxx/fp_data_mining

Prediction Model



Conclusions & Future Work (2-3 paragraphs or more)

Based on the bar chart results there is a significant correlation between the two tumor types and attribute results. The pattern that is seemingly repetitive in relation to tumor type, for almost every attribute, was interesting to analyze. Tumor type benign attribute results was relatively on the lower end of the number range, however for malignant attribute results were opposite of benign where malignant results were relatively on the higher end of the number range. In the project prediction model was also created. Based on the accuracy score for the prediction, it was roughly ≈ 0.96 for all three (GaussianNB, KNeighbor Classifier N=30, and Support Vector Machine) methods used in the project. This result shows that attribute result is reliable source to identifying if the tumor is benign or malignant

This project has potential growth and room for further development in the field of healthcare. As stated early on, breast cancer is becoming a nationwide problem and more or less an epidemic. It is definitely preventable and treatable with correct diagnosis and treatments. Considering the data used for this project was fairly old (between for this project is fairly old (between 1989 to 1991), it would be interesting to analyze with a recent dataset with same attributes and format.

Breast cancer is a growing problem, and there is no “right” solution to tackling this issue. However, decreasing misdiagnosis and increase in correct diagnosis can significantly help an individual go towards the correct path of treatment, as necessary.

References:

1. Dataset Citation: [https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(original\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(original))
 - a. 1. O. L. Mangasarian and W. H. Wolberg: "Cancer diagnosis via linear programming", SIAM News, Volume 23, Number 5, September 1990, pp 1 & 18.
 - b. 2. William H. Wolberg and O.L. Mangasarian: "Multisurface method of pattern separation for medical diagnosis applied to breast cytology", Proceedings of the National Academy of Sciences, U.S.A., Volume 87, December 1990, pp 9193-9196.
 - c. 3. O. L. Mangasarian, R. Setiono, and W.H. Wolberg: "Pattern recognition via linear programming: Theory and application to medical diagnosis", in: "Large-scale numerical optimization", Thomas F. Coleman and Yuying Li, editors, SIAM Publications, Philadelphia 1990, pp 22-30.
 - d. 4. K. P. Bennett & O. L. Mangasarian: "Robust linear programming discrimination of two linearly inseparable sets", Optimization Methods and Software 1, 1992, 23-34 (Gordon & Breach Science Publishers).
2. HTML. w3schools.com. <https://www.w3schools.com/html/>. Accessed November 22, 2017.
3. What is Breast Cancer? Breastcancer.org. <https://www.cancer.org/cancer/breast-cancer/about/what-is-breast-cancer.html>. Published September 21, 2017. Accessed November 22, 2017.
4. Breast Cancer Statistics. cdc.gov. <https://www.cdc.gov/cancer/breast/statistics/index.htm>. Published June 7, 2017. Accessed November 22, 2017.
5. New England Journal of Medicine
6. www.preventcancer.org
7. www.google.com for the images
8. <http://scikit-learn.org/stable/index.html>
9. <https://stackoverflow.com/questions/>
10. <https://pandas.pydata.org/pandas-docs/stable/>
11. Class lecture notes on python, unix, and html.