

Problem1.1

1. Use X.shape to get the number of features and the number of data points. Report both numbers, mentioning which number is which. (5 points)

```
In [256]: 1 import numpy as np
          2 import matplotlib.pyplot as plt
          3 nych = np.genfromtxt ("nyc_housing.txt", delimiter = None)
          4 Y = nych[:, -1]
          5 X = nych[:, 0:-1]
```

```
In [257]: 1 X.shape
```

```
Out[257]: (300, 3)
```

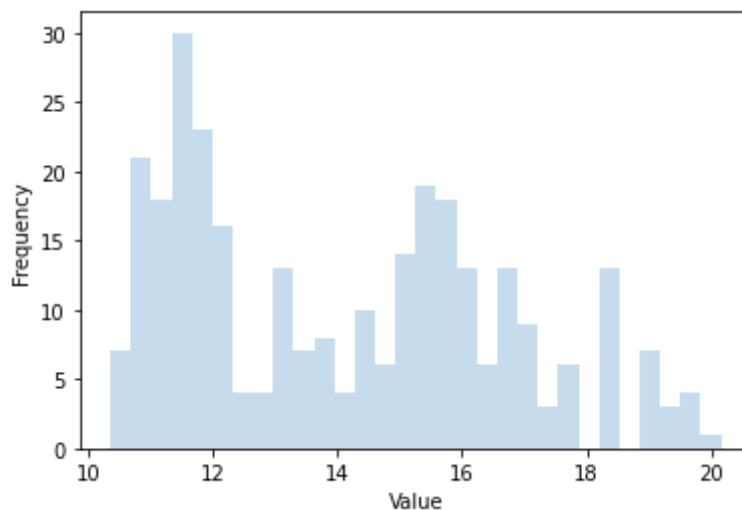
300 = Number of data points (row), 3 = Number of features (col)

Problem1.2

For each feature, plot a histogram (plt.hist) of the data values. (5 points)

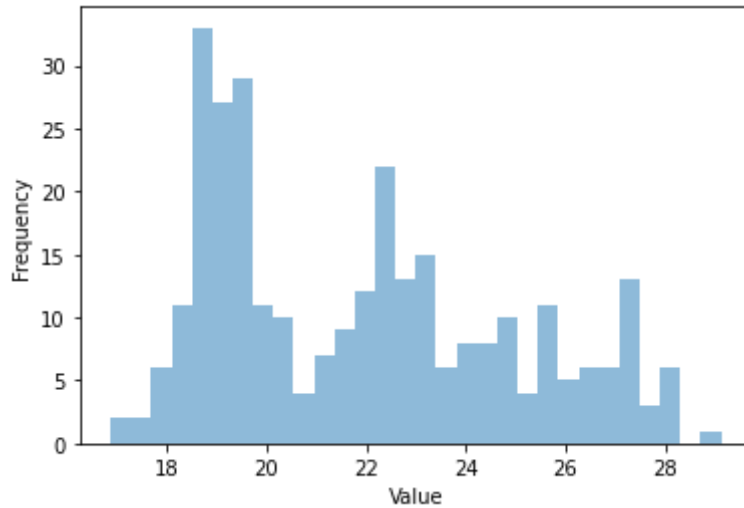
```
In [89]: 1 #First histogram
          2 #feature1 = [i[0] for i in X]
          3
          4 feature1 = nych[:,0]
          5 plt.hist(feature1, bins = 30, alpha = 0.25)
          6 plt.xlabel('Value')
          7 plt.ylabel('Frequency')
```

```
Out[89]: Text(0, 0.5, 'Frequency')
```



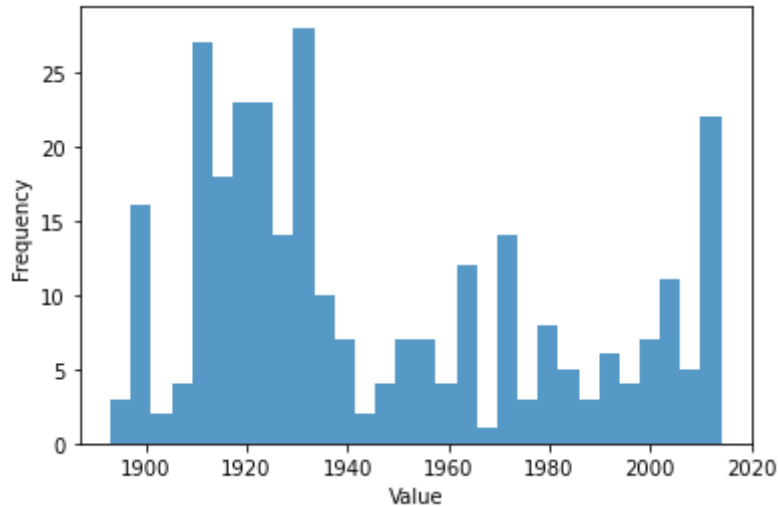
```
In [88]: 1 #Second Histogram
2 feature2 = [i[1] for i in X]
3 plt.hist(feature2, bins = 30, alpha = 0.5)
4 plt.xlabel('Value')
5 plt.ylabel('Frequency')
```

Out[88]: Text(0, 0.5, 'Frequency')



```
In [90]: 1 #Thrid Histogram
2 feature3 = [i[2] for i in X]
3 plt.hist(feature3, bins = 30, alpha = 0.75)
4 plt.xlabel('Value')
5 plt.ylabel('Frequency')
```

Out[90]: Text(0, 0.5, 'Frequency')



Problem1.3

```
In [263]: 1 a = X[:, 0:1]
          2
          3 print("feature 1 \n mean = " , np.mean(a) , "and standard deivation = "
          4
          5
          6 b = X[:, 1:2]
          7 print("feature 2 \n mean =", np.mean(b), " and standard deivation = ",
          8
          9 c = X[:, 2:3]
         10 print("feature 3 \n mean =", np.mean(c), " and standard deivation =", n
         11
         12
```

```
feature 1
mean = 14.118392438424483 and standard deivation = 2.569090284260317

feature 2
mean = 21.907116176170856 and standard deivation = 2.9785784999947165

feature 3
mean = 1946.3533333333332 and standard deivation = 35.39889577687731
```

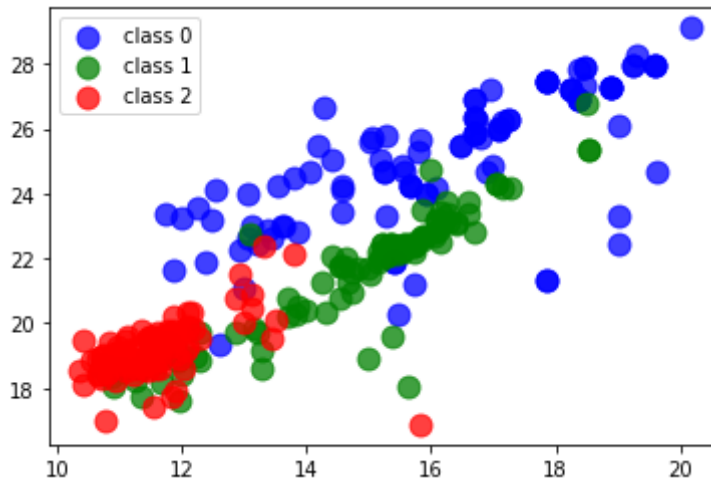
Problem1.4

For each pair of features (1,2), (1,3), and (2,3), plot a scatterplot (see `plt.plot` or `plt.scatter`) of the feature values, colored according to their target value (class). (For example, plot all data points with $y = 0$ as blue, $y = 1$ as green, and $y = 2$ as red.) (5 points)

```

In [93]: 1 colors = ['blue', 'green', 'red']
          2
          3 for i, c in enumerate(np.unique(nych[:, -1])):
          4     mask = np.where(nych[:, -1] == c)[0] # Finding the right points
          5     plt.scatter(nych[mask, 0], nych[mask, 1], s=120, c=colors[i], alpha
          6
          7 plt.legend()
          8 plt.show()

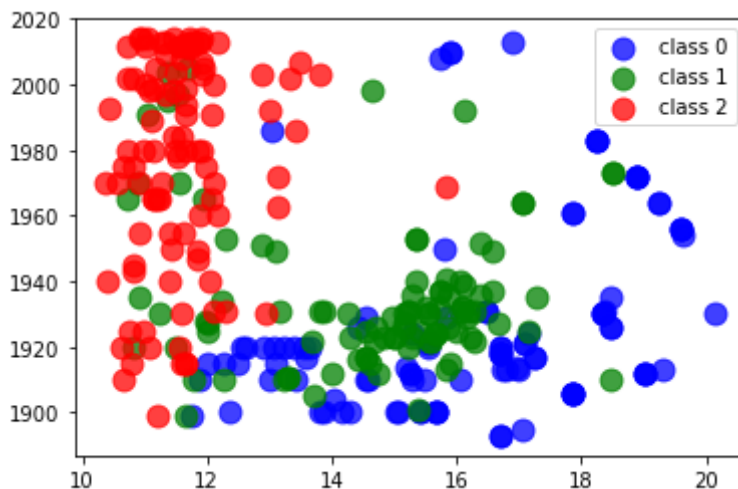
```



```

In [94]: 1 colors = ['blue', 'green', 'red']
          2
          3 for i, c in enumerate(np.unique(nych[:, -1])):
          4     mask = np.where(nych[:, -1] == c)[0] # Finding the right points
          5     plt.scatter(nych[mask, 0], nych[mask, 2], s=120, c=colors[i], alpha
          6
          7 plt.legend()
          8 plt.show()

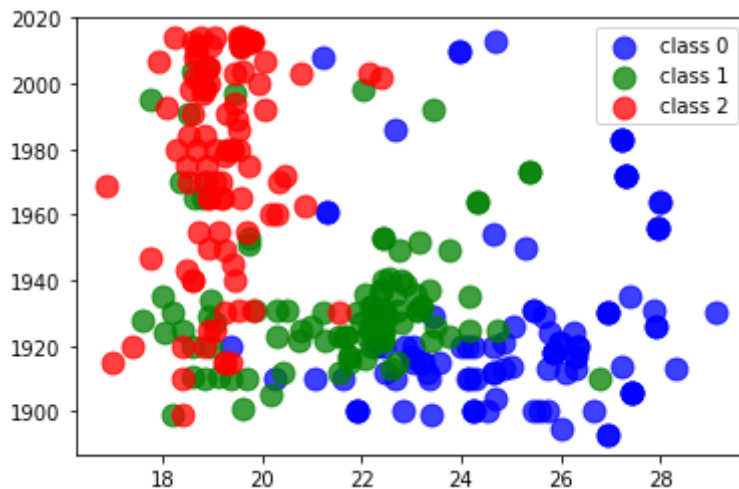
```



```

In [95]: 1 colors = ['blue', 'green', 'red']
          2
          3 for i, c in enumerate(np.unique(nych[:, -1])):
          4     mask = np.where(nych[:, -1] == c)[0] # Finding the right points
          5     plt.scatter(nych[mask, 1], nych[mask, 2], s=120, c=colors[i], alpha
          6
          7 plt.legend()
          8 plt.show()

```



Problem 2: k-nearest-neighbor predictions (25 points)

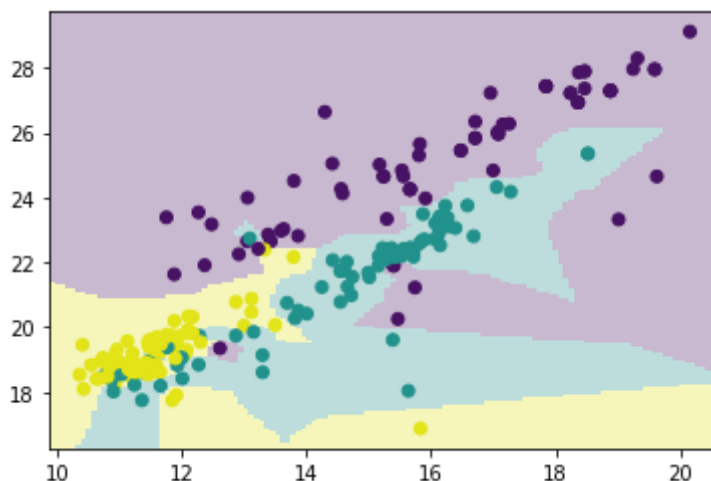
Problem 2.1

```

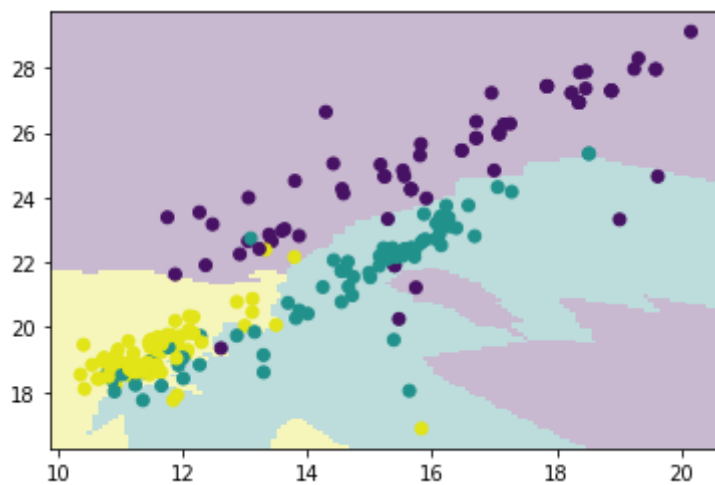
In [155]: 1 nych = np.genfromtxt("nyc_housing.txt",delimiter=None) # load the data
2 Y = nych[:, -1]
3 X = nych[:, 0:-1]
4 # Note: indexing with ":" indicates all values (in this case, all rows)
5 # indexing with a value ("0", "1", "-1", etc.) extracts only that value
6 # indexing rows/columns with a range ("1:-1") extracts any row/column
7
8 import mltools as ml
9 # We'll use some data manipulation routines in the provided class code
10 # Make sure the "mltools" directory is in a directory on your Python path
11 # export PYTHONPATH=$\${PYTHONPATH}:/path/to/parent/dir
12 # or add it to your path inside Python:
13 # import sys
14 # sys.path.append('/path/to/parent/dir/');
15
16 np.random.seed(0) # set the random number seed
17 X,Y = ml.shuffleData(X,Y); # shuffle data randomly
18 # (This is a good idea in case your data are ordered in some systematic way)
19
20 Xtr,Xva,Ytr,Yva = ml.splitData(X,Y, 0.75); # split data into 75/25 train/test
21
22 K=[1,5,10,50]
23 knn = ml.knn.knnClassify()
24 for i,k in enumerate(K):
25     print("K = "+str(k) )
26
27     knn.train(Xtr, Ytr, k)
28
29     Yva = knn.predict ( Xva)
30     knn = ml.knn.knnClassify()
31     knn.train(Xtr[:, :2], Ytr, k)
32
33     ml.plotClassify2D(knn, Xtr[:, :2], Ytr)
34     plt.show()

```

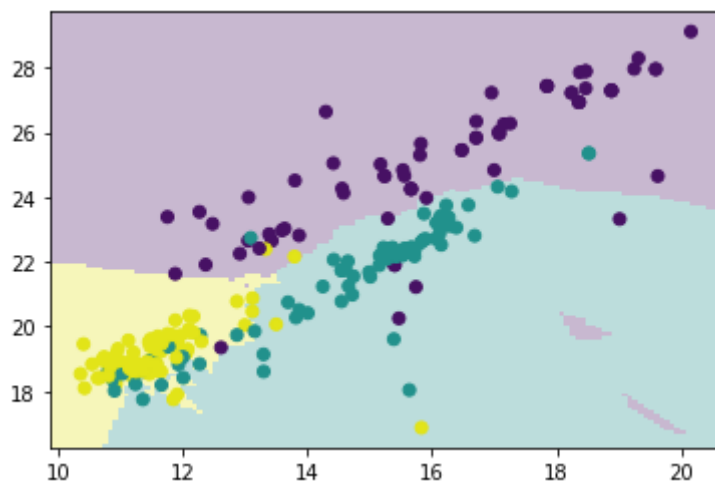
K = 1



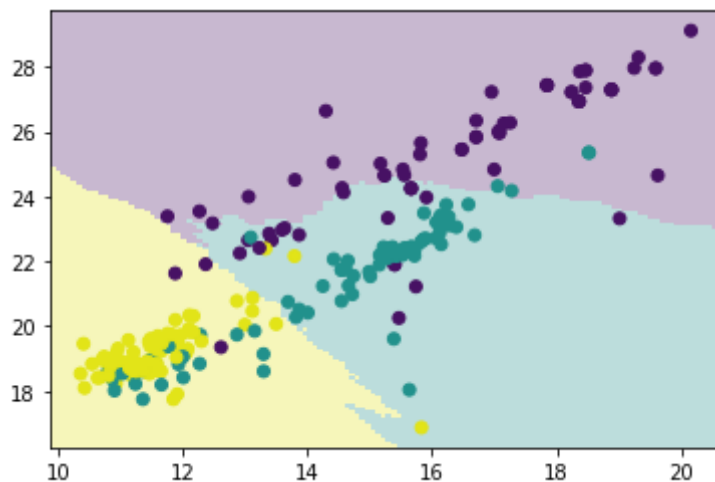
K = 5



K = 10



K = 50



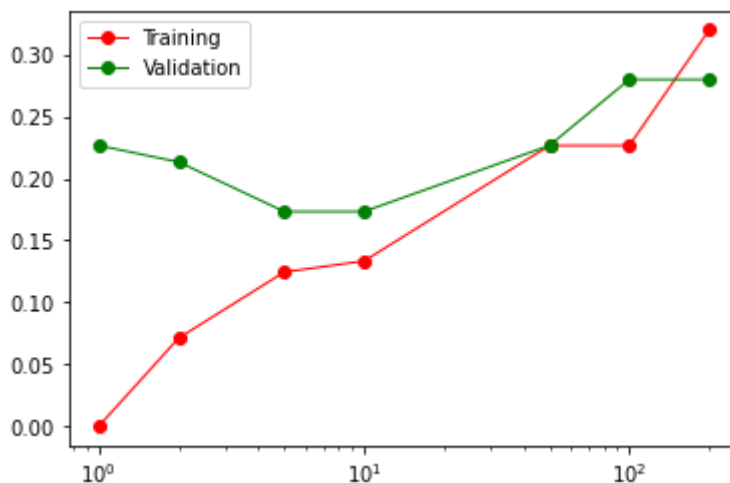
Problem 2.2

Again using only the first two features, compute the error rate (number of misclassifications) on both the training and validation data as a function of $K = [1, 2, 5, 10, 50, 100, 200]$. You can do this most easily with a for-loop:

```

In [217]: 1 #with training error in red and validation error in green.
2 Y = nych[:, -1] #only the first two features
3 X = nych[:, 0:2] #only the first two features
4 np.random.seed(0)
5 X,Y = ml.shuffleData(X, Y); # shuffle data randomly
6 Xtr,Xva,Ytr,Yva = ml.splitData(X, Y, 0.75) # split data into 75/
7
8 K = [1,2,5,10,50,100,200]
9
10 errTrain = [None] * len(K)
11 errVal = [None] * len(K)
12 for i, k in enumerate(K):
13     learner = ml.knn.knnClassify() # create the object and train it
14     learner.train(Xtr, Ytr, k) # where K is an integer, e.g. 1 for
15     errTrain[i] = learner.err(Xtr, Ytr)
16     errVal[i] = learner.err(Xva, Yva)
17
18 plt.semilogx(K, errTrain, color='red', lw=1, marker='o', label='Tra
19 plt.semilogx(K, errVal, color='green', lw=1, marker='o', label='Val
20 # Adding a legend to the plot that will use the labels from the 'label'
21 plt.legend()
22
23 # Controlling the axis.
24 ax.semilogx(.8, 200)
25 ax.semilogx(0, 1)
26 plt.show()
27 # pratie one

```



Based on these plots, I would like to recommend value of 10 (between $10^{(0.5)}$ and $10^{(1)}$). Because this value is the lowest

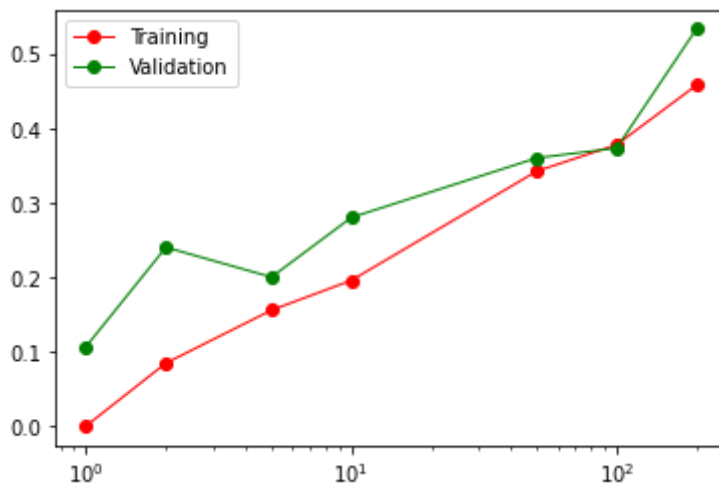
Problem2.3

Again using only the first two features, compute the error rate (number of misclassifications) on both the training and validation data as a function of $K = [1, 2, 5, 10, 50, 100, 200]$. You can do this most easily with a for-loop:


```

In [215]: 1 Y = nych[:, -1]
2 X = nych[:, 0:-1]
3 np.random.seed(0)
4 X,Y = ml.shuffleData(X, Y); # shuffle data randomly
5 Xtr,Xva,Ytr,Yva = ml.splitData(X, Y, 0.75) # split data into 75/25
6
7 K = [1,2,5,10,50,100,200]
8
9 errTrain = [None] * len(K)
10 errVal = [None] * len(K)
11 for i, k in enumerate(K):
12     learner = ml.knn.knnClassify() # create the object and train
13     learner.train(Xtr, Ytr, k) # where K is an integer, e.g. 1, 2, 5, 10, 50, 100, 200
14     errTrain[i] = learner.err(Xtr, Ytr) # to calculate the training error
15     errVal[i] = learner.err(Xva, Yva) # to calculate the validation error
16
17 plt.semilogx(K, errTrain, color='red', lw=1, marker='o', label='Training')
18 plt.semilogx(K, errVal, color='green', lw=1, marker='o', label='Validation')
19
20
21 plt.legend() # Adding a legend to the plot that will use the labels from the legend
22
23
24 ax.semilogx(.8, 200) # Controlling the axis.
25 ax.semilogx(0, 1)
26 plt.show()
27
28 # practice one

```



Based on these plots, I would like to recommend value of 5. Because this value is the lowest

Problem3 Naïve Bayes Classifiers (35 points)

Problem3.1

Compute all the probabilities necessary for a naïve Bayes classifier, i.e., the class probability $p(y)$ and all the individual feature probabilities $p(x_i | y)$, for each class y and feature x_i . (7 points)

```
In [137]: 1 from IPython.display import Image
          2 Image( filename = 'page/178-hw1-9.jpg')
```

Out[137]:

1. Compute all the probabilities necessary for a naïve Bayes classifier, i.e., the class probability $p(y)$ and all the individual feature probabilities $p(x_i|y)$, for each class y and feature x_i . (7 points)

$$P(y=1) = \frac{4}{10} = \frac{2}{5}$$

$$\begin{aligned} P(x_1=1 | y=1) &= \frac{3}{4}, & P(x_1=0 | y=1) &= \frac{1}{4} \\ P(x_2=1 | y=1) &= \frac{0}{4} = 0, & P(x_2=0 | y=1) &= \frac{4}{4} = 1 \\ P(x_3=1 | y=1) &= \frac{3}{4}, & P(x_3=0 | y=1) &= \frac{1}{4} \\ P(x_4=1 | y=1) &= \frac{2}{4} = \frac{1}{2}, & P(x_4=0 | y=1) &= \frac{2}{4} = \frac{1}{2} \\ P(x_5=1 | y=1) &= \frac{1}{4}, & P(x_5=0 | y=1) &= \frac{3}{4} \end{aligned}$$

$$P(y=-1) = \frac{6}{10} = \frac{3}{5}$$

$$\begin{aligned} P(x_1=1 | y=-1) &= \frac{3}{6} = \frac{1}{2}, & P(x_1=0 | y=-1) &= \frac{3}{6} = \frac{1}{2} \\ P(x_2=1 | y=-1) &= \frac{5}{6}, & P(x_2=0 | y=-1) &= \frac{1}{6} \\ P(x_3=1 | y=-1) &= \frac{4}{6} = \frac{2}{3}, & P(x_3=0 | y=-1) &= \frac{2}{6} = \frac{1}{3} \\ P(x_4=1 | y=-1) &= \frac{5}{6}, & P(x_4=0 | y=-1) &= \frac{1}{6} = \frac{1}{6} \\ P(x_5=1 | y=-1) &= \frac{2}{6} = \frac{1}{3}, & P(x_5=0 | y=-1) &= \frac{4}{6} = \frac{2}{3} \end{aligned}$$

Problem3.2

```
In [139]: 1 from IPython.display import Image
          2 Image( filename = 'page/178-hw1-10.jpg')
```

Out[139]:

2. Which class would be **predicted** for $\underline{x} = (0\ 0\ 0\ 0\ 0)$? What about for $\underline{x} = (1\ 1\ 0\ 1\ 0)$? (7 points)

$$\underline{x} = (0\ 0\ 0\ 0\ 0)$$

$$P(x|y=1), \frac{1}{4} \times 1 \times \frac{1}{4} \times \frac{2}{4} \times \frac{2}{4} = \frac{6}{4^4} = \frac{6}{256} = \frac{3}{128}$$

$$P(x|y=-1), \frac{3}{6} \times \frac{1}{6} \times \frac{2}{6} \times \frac{1}{6} \times \frac{4}{6} = \frac{24}{6^5} = \frac{24}{7776} = \frac{12}{3888} = \frac{6}{1944} = \frac{3}{972} = \frac{1}{324}$$

$$P(y=1)P(x|y=1) = \frac{4}{10} \times \frac{3}{128} = 0.009375$$

$$P(y=-1)P(x|y=-1) = \frac{6}{10} \times \frac{1}{324} = 0.00185185185 = 0.001852$$

$$P(y=1)P(x|y=1) > P(y=-1)P(x|y=-1)$$

\therefore Predicted as "Read".

$$\underline{x} = (1\ 1\ 0\ 1\ 0) \quad \begin{matrix} x_1=1 \& y=1 \\ \uparrow \\ x_2=0, y=1 \\ \uparrow \end{matrix}$$

$$P(x|y=1) = \frac{3}{4} \times \frac{0}{4} \times \frac{1}{4} \times \frac{2}{4} \times \frac{3}{4} = \frac{0}{1024}$$

$$P(x|y=-1) = \frac{3}{6} \times \frac{5}{6} \times \frac{2}{6} \times \frac{5}{6} \times \frac{4}{6} = \frac{600}{7776}$$

$$P(y=1)P(x|y=1) = \frac{4}{10} \times \frac{0}{1024}$$

$$P(y=-1)P(x|y=-1) = \frac{6}{10} \times \frac{600}{7776} = 0.04629629629$$

$$\approx 0.04630$$

$$P(x|y=1) < P(x|y=-1)$$

\therefore Predict "Discard"

Problem3.3

```
In [146]: 1 from IPython.display import Image
          2 Image( filename = 'page/178-hw1-11.jpg')
```

Out[146]:

3. Compute the posterior probability that $y = +1$ given the observation $\underline{x} = (0\ 0\ 0\ 0\ 0)$. Also compute the posterior probability that $y = +1$ given the observation $\underline{x} = (1\ 1\ 0\ 1\ 0)$. (7 points)

$$\underline{x} = (0\ 0\ 0\ 0\ 0)$$

$$P(y=1|\underline{x}) = \frac{P(y=1) P(\underline{x}|y=1)}{P(y=1) \cdot P(\underline{x}|y=1) + P(y=-1) P(\underline{x}|y=-1)}$$

$$= \frac{\frac{4}{10} \times \frac{24}{1024}}{\frac{4}{10} \times \frac{24}{1024} + \frac{6}{10} \times \frac{24}{7776}} \approx 0.8351$$

$$\underline{x} = (1\ 1\ 0\ 1\ 0)$$

$$P(y=1|\underline{x}) = \frac{P(y=1) P(\underline{x}|y=1)}{P(y=1) \cdot P(\underline{x}|y=1) + P(y=-1) P(\underline{x}|y=-1)}$$

$$= \frac{\frac{4}{10} \times 0}{\frac{4}{10} \times 0 + \frac{6}{10} \times \frac{600}{7776}} = 0.00$$

Problem3.4


```
In [143]: 1 from IPython.display import Image
          2 Image( filename = 'page/178-hw1-12.jpg' )
```

Out[143]:

4. Why should we probably not use a “joint” Bayes classifier (using the joint probability of the features x , as opposed to the conditional independencies assumed by naïve Bayes) for these data? (7 points)

We can think about two chararistst of "Joint Bayes". One the depedent and antoher is total number.

First, case of "Joint Bayes" need to be dependent. But this problems (email) event is not dependent not independent.

Second, for the "Joint Bayes", we need $32(2^5)$ parameters. But, in this case, we just have 10 parameters.

Problem3.5

```
In [145]: 1 from IPython.display import Image
          2 Image( filename = 'page/178-hw1-13.jpg' )
```

Out[145]:

5. Suppose that before we make our predictions, we lose access to my address book, so that we cannot tell whether the email author is known. Do we need to re-train the model to classify based solely on the other four features? If so, how? If not, what changes about how our trained parameters are used? Hint: what parameters do I need for a naïve Bayes model over only features x_2, \dots, x_5 ? Do I need to re-calculate any new parameter values in our new setting? What, if anything, changes about the parameters or the way they are used? (7 points)

No, we don't have to re-train the model. Because this is independency event, it will not change the probabilities.

For example, probability of $P(x_1 = 1 | y = 1)$ will not be changed even if $P(x_n = 1 | y = 1)$ is changed.

Problem 4: Gaussian Bayes Classifiers (15 points)

Now, using the NYC Housing data, we will explore a classifier based on Bayes rule. Again, we'll use only the first two features of NYC Housing, shuffled and split in to training and validation sets as before.

1. Splitting your training data by class, compute the empirical mean vector and covariance matrix of the data in each class. (You can use mean and cov for this.) (5 points)

In []:

1

Problem 5: Statement of Collaboration (5 points)

- Piazza question @116 "How to use err() "
- Pual Sung : talk about basic principle of "k-nearest-neighbor predictions"
- Paul Sung : discuss about probability for naive bayes