

Name: Hootan Hosseinzadeganbushehri

XXXXXXXXXXXX

XXXXXXXXXXXX

UCI Email: hootanh@uci.edu

# Problem 1 - Part 1:

```
In [161]: import numpy as np
import matplotlib.pyplot as plt

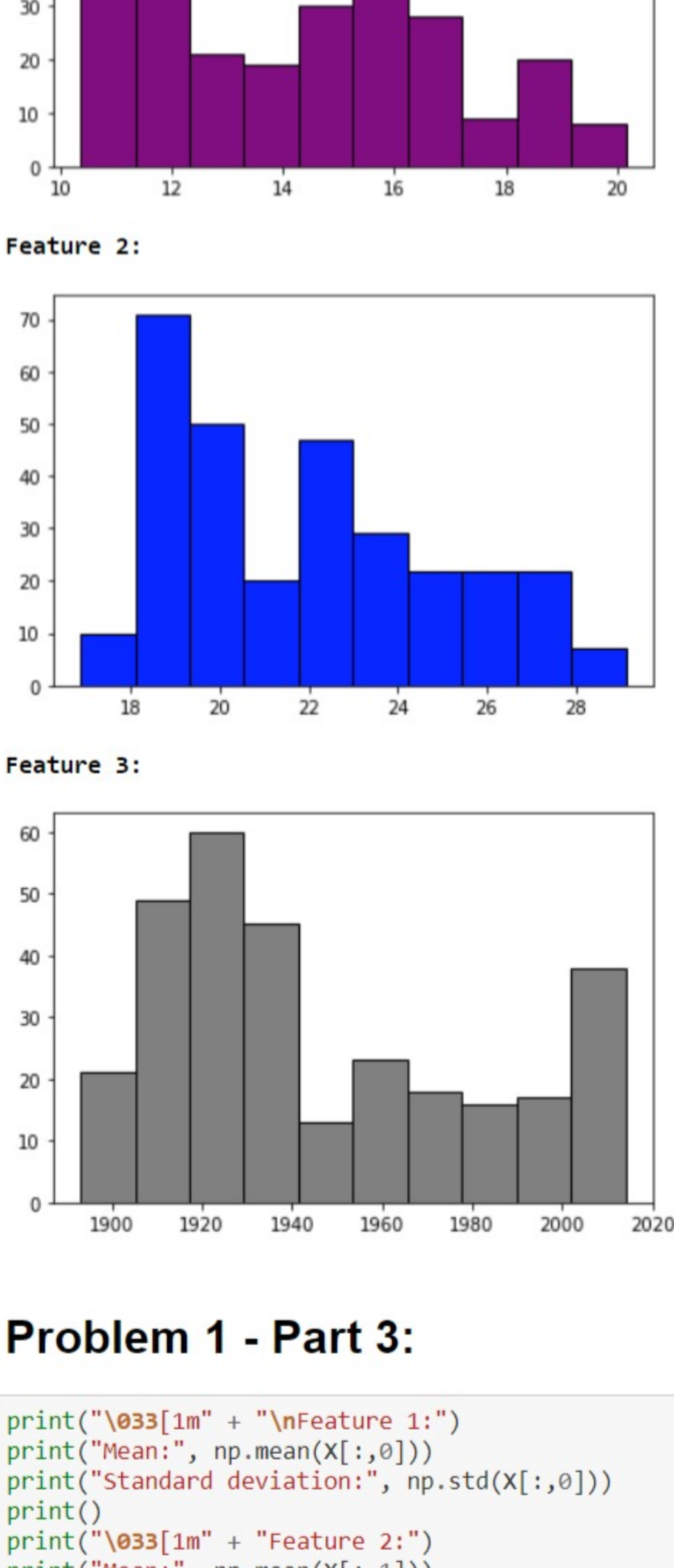
nych = np.genfromtxt("data/nycc_housing.txt", delimiter = None)
Y = nych[:,1]
X = nych[:,0:-1]

print("\033[1m" + "\nnumber of data points:", X.shape[0], "\033[1m" + "\nnumber of features:", X.shape[1])
```

Number of data points: 300  
Number of features: 3

# Problem 1 - Part 2: II

```
In [162]: print("\033[1m" + "\nfeature 1:")
plt.hist(X[:,0], color = "purple", edgecolor = "black")
plt.show()
print()
print("\033[1m" + "feature 2:")
plt.hist(X[:,1], color = "blue", edgecolor = "black")
plt.show()
print()
print("\033[1m" + "feature 3:")
plt.hist(X[:,2], color = "gray", edgecolor = "black")
plt.show()
```



# Problem 1 - Part 3:

```
In [163]: print("\033[1m" + "\nfeature 1 and 2:")
print("Mean:", np.mean(X[:,0]))
print("Standard deviation:", np.std(X[:,0]))
print()
print("\033[1m" + "feature 2:")
print("Mean:", np.mean(X[:,1]))
print("Standard deviation:", np.std(X[:,1]))
print()
print("\033[1m" + "feature 3:")
print("Mean:", np.mean(X[:,2]))
print("Standard deviation:", np.std(X[:,2]))
```

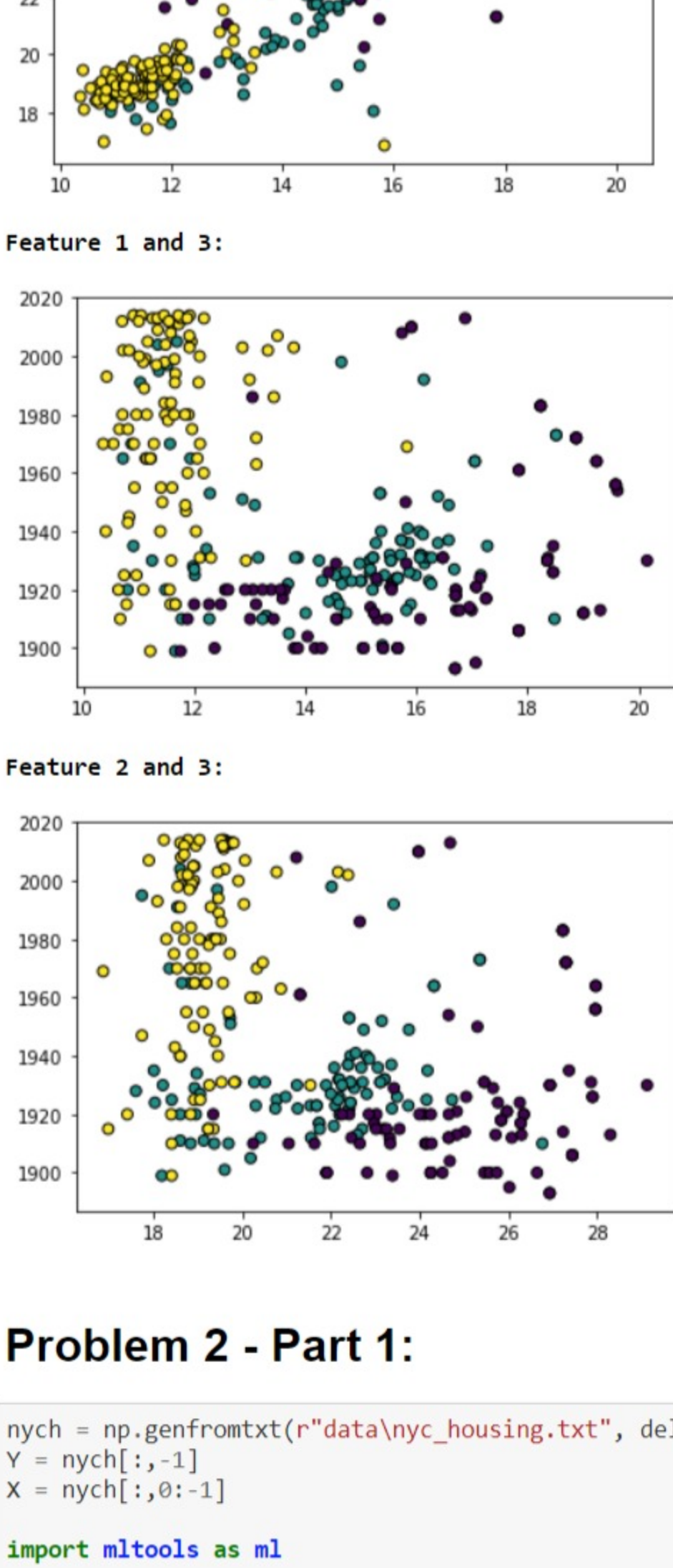
Feature 1:  
Mean: 14.118392438424893  
Standard deviation: 2.569998284260317

Feature 2:  
Mean: 21.507116176170856  
Standard deviation: 2.9785784999947165

Feature 3:  
Mean: 1946.3533333333332  
Standard deviation: 35.398895776877931

# Problem 1 - Part 4:

```
In [164]: print("\033[1m" + "\nfeature 1 and 2:")
plt.scatter(X[:,0], X[:,1], c = Y, edgecolors = "black")
plt.show()
print()
print("\033[1m" + "feature 1 and 3:")
plt.scatter(X[:,0], X[:,2], c = Y, edgecolors = "black")
plt.show()
print()
print("\033[1m" + "feature 2 and 3:")
plt.scatter(X[:,1], X[:,2], c = Y, edgecolors = "black")
plt.show()
```



# Problem 2 - Part 1:

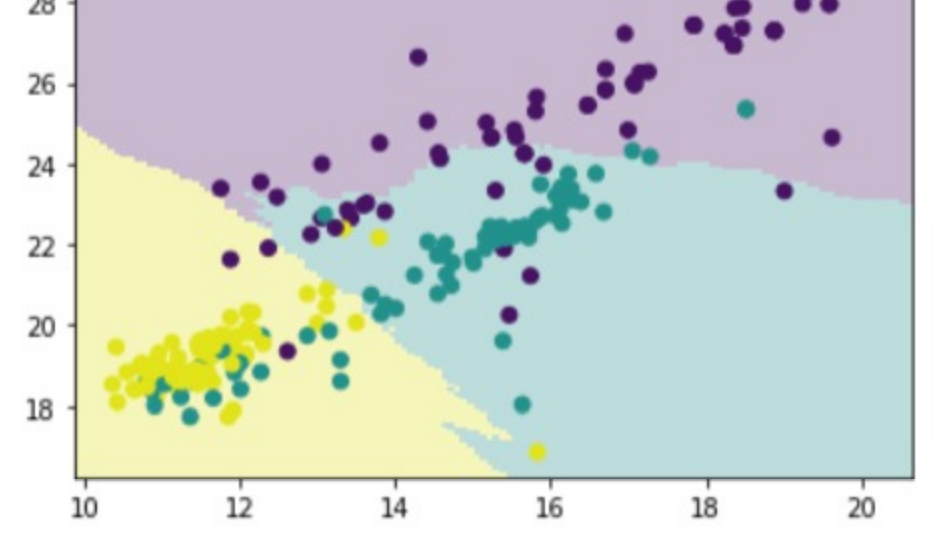
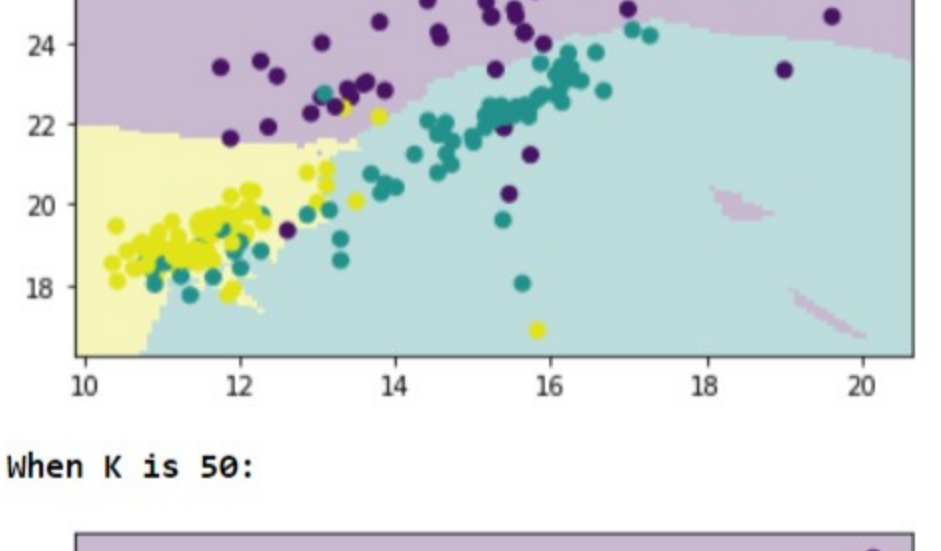
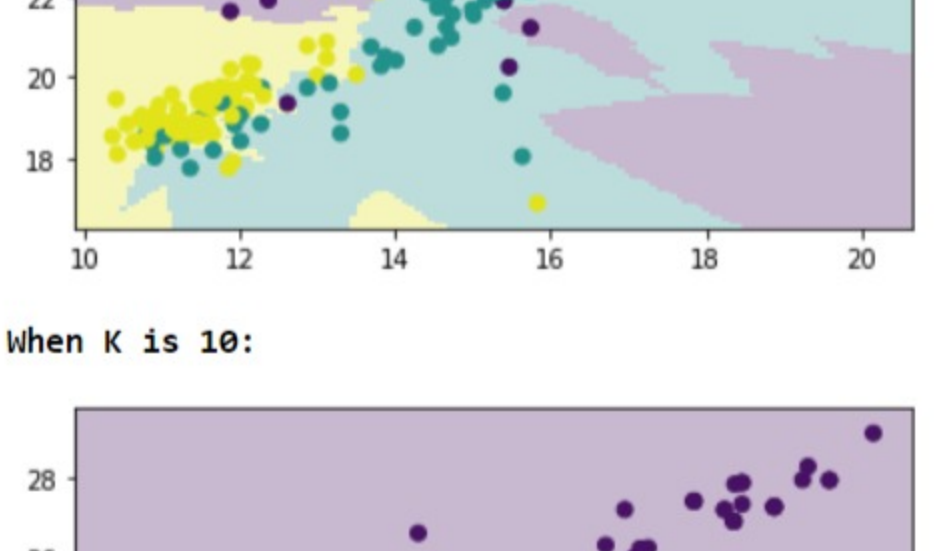
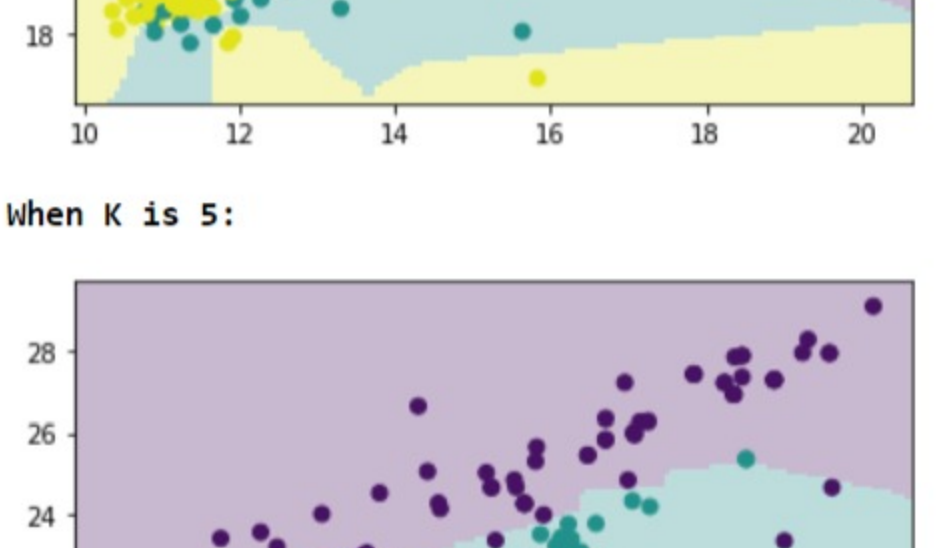
```
In [165]: nych = np.genfromtxt("data/nycc_housing.txt", delimiter = None)
Y = nych[:,1]
X = nych[:,0:-1]

import mltools as ml

np.random.seed(0)
X, Y = ml.shuffleData(X, Y)

Xtr, Xva, Ytr, Yva = ml.splitData(X, Y, 0.75)
k = [1, 5, 10, 50]

for i in range(len(k)):
    print("\033[1m" + "when K is " + str(k[i]) + ":")
    knn = ml.knn.knnClassify()
    knn.train(Xtr[:, 0:-2], Ytr, k[i])
    ml.plotClassify2D(knn, Xtr[:, 0:-2], Ytr)
    plt.show()
```

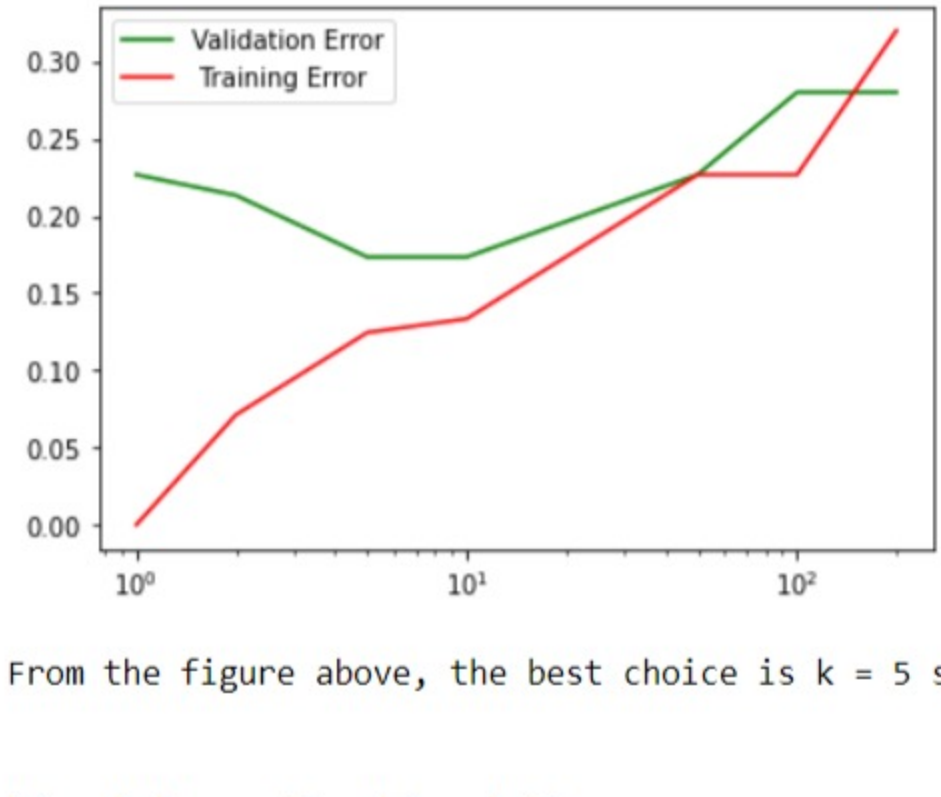


# Problem 2 - Part 2:

```
In [166]: K = [1, 2, 5, 10, 50, 100, 200]
errTrain = [None]*len(K)
errValid = [None]*len(K)

for i, k in enumerate(K):
    learner = ml.knn.knnClassify(Xtr[:,0:-2], Ytr, k)
    errTrain[i] = learner.err(Xtr[:,2], Ytr)
    errValid[i] = learner.err(Xva[:,2], Yva)

plt.semilogx(K, errValid, "green", label = "Validation Error")
plt.semilogx(K, errTrain, "red", label = "Training Error")
plt.legend()
print()
print("\033[1m" + "Error rate on both the training and validation data:")
print("From the figure above, the best choice is k =", K[np.argmin(errValid)], "so I would recommend", K[np.argmin(errValid)], "for K.")
```



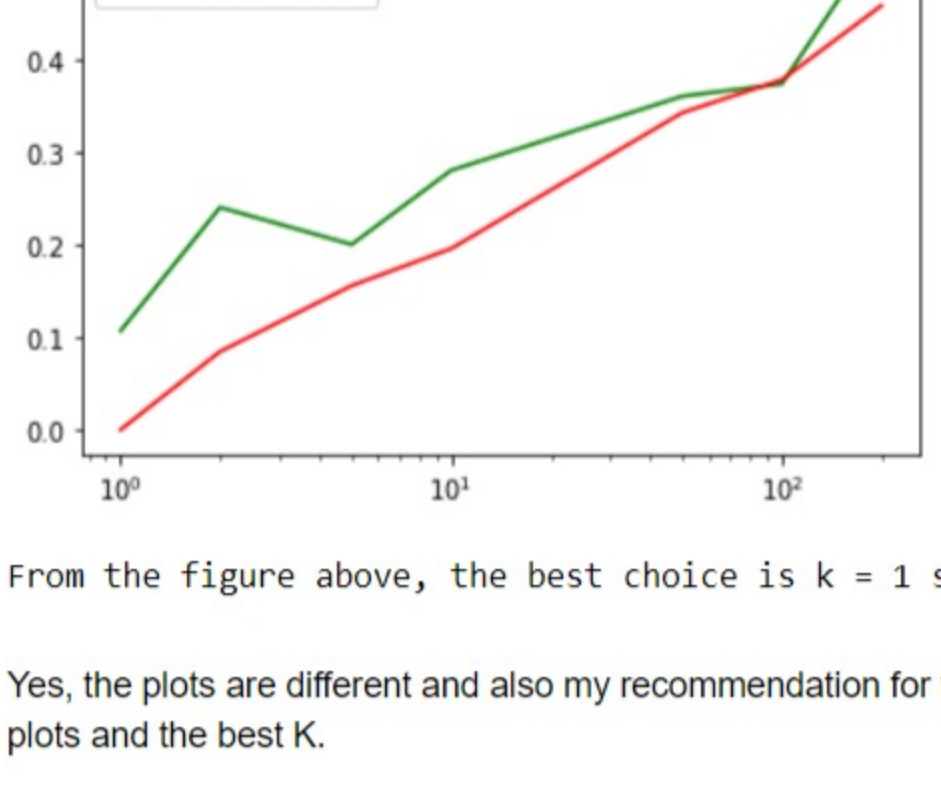
From the figure above, the best choice is k = 5 so I would recommend 5 for K.

# Problem 2 - Part 3:

```
In [167]: K = [1, 2, 5, 10, 50, 100, 200]
errTrain = [None]*len(K)
errValid = [None]*len(K)

for i, k in enumerate(K):
    learner = ml.knn.knnClassify(Xtr, Ytr, k)
    errTrain[i] = learner.err(Xtr[:,2], Ytr)
    errValid[i] = learner.err(Xva, Yva)

plt.semilogx(K, errValid, "green", label = "Validation Error")
plt.semilogx(K, errTrain, "red", label = "Training Error")
plt.legend()
print()
print("\033[1m" + "Error rate on both the training and validation data:")
print("From the figure above, the best choice is k =", K[np.argmin(errValid)], "so I would recommend", K[np.argmin(errValid)], "for K.")
```



From the figure above, the best choice is k = 1 so I would recommend 1 for K.

Yes, the plots are different and also my recommendation for the best K is different as well; therefore, a change in the number of the features can change both plots and the best K.

# Problem 3 - Part 1:

$p(y=-1) = 6/10 = 3/5$

$p(y=1) = 4/10 = 2/5$

$p(x=1|y=-1) = 3/6 = 1/2$

$p(x=2|y=-1) = 3/6 = 1/2$

$p(x=3|y=-1) = 4/6 = 2/3$

$p(x=4|y=-1) = 5/6$

$p(x=5|y=-1) = 2/6 = 1/3$

$p(x=1|y=1) = 3/4$

$p(x=2|y=1) = 0/4 = 0$

$p(x=3|y=1) = 3/4$

$p(x=4|y=1) = 2/4 = 1/2$

$p(x=5|y=1) = 1/4$

$p(x=1|y=-1) = 3/6 = 1/2$

$p(x=2|y=-1) = 1/6$

$p(x=3|y=-1) = 2/6 = 1/3$

$p(x=4|y=-1) = 1/6$

$p(x=5|y=-1) = 4/6 = 2/3$

$p(x=1|y=1) = 1/4$

$p(x=2|y=1) = 4/4 = 1$

$p(x=3|y=1) = 1/4$

$p(x=4|y=1) = 2/4 = 1/2$

$p(x=5|y=1) = 3/4$

# Problem 3 - Part 2:

$p(00000|y=1) = (4/10)(1/4)(1/4)(1/4)(2/4)(3/4) = 0.009375$

$p(00000|y=-1) = (6/10)(3/6)(1/6)(2/6)(1/6)(4/6) = 0.001851852$

==>

$p(00000|y=1) > p(00000|y=-1) \Rightarrow$  Class  $y = 1$  is bigger than class  $y = -1$

==> Therefore, the email must be read.

$p(11010|y=1) = (4/10)(3/4)(0)(3/4)(2/4)(1/4) = 0$

$p(11010|y=-1) = (6/10)(3/6)(5/6)(2/6)(5/6)(4/6) = 0.046296296$

==>

$p(11010|y=1) < p(11010|y=-1) \Rightarrow$  Class  $y = -1$  is bigger than class  $y = 1$

==> Therefore, the email must be discarded.

# Problem 3 - Part 3:

$p(y=1|00000) =$

$= p(00000|y=1) p(y=1) / p(y=1) p(00000|y=1) + p(00000|y=-1) p(y=-1)$

$= 0.8351$

$\Rightarrow p(y=1|00000) = 0.8351$

$p(y=1|11010) =$

$= p(11010|y=1) p(y=1) / p(y=1) p(11010|y=1) + p(11010|y=-1) p(y=-1)$

$= 0$

$\Rightarrow p(y=1|11010) = 0$

# Problem 3 - Part 4:

Since here we have many parameters, joint Bayes classifier is going to take a lot of time and can be very slow as well. Also, joint Bayes classifier can increase the complexity in the computation. However, it is more efficient for us to use naive Bayes classifier instead because the data features are independent and can be done by using fewer parameters. In addition, naive Bayes classifier can simplify the calculations for the probabilities and increase the speed of the process which is going to be faster than joint Bayes classifier. Therefore, we should not use a joint Bayes classifier for these data.

# Problem 3 - Part 5:

We do not need to re-train the model because all the features are independent from each other; therefore, it is not needed to re-train the model even after deleting a feature. Furthermore, we simply just exclude the first feature(x1) meaning that we do not use  $p(x1=1|y=1)$  which is the probability associated with the first feature in the calculation and there is no change in the rest of the probabilities. Therefore, we do not have to re-train the model or re-calculate parameters after excluding the probability for the first feature(x1).

# Problem 4:

```
In [171]: nych = np.genfromtxt("data/nycc_housing.txt", delimiter = None)
Y = nych[:,1]
X = nych[:,0:-1]

import mltools as ml

np.random.seed(0)
X, Y = ml.shuffleData(X, Y)

Xtr, Xva, Ytr, Yva = ml.splitData(X, Y, 0.75)

learner = ml.bayes.gaussClassify(Xtr[:,0:-2], Ytr, k)
errTrain = learner.err(Xtr[:,2], Ytr)
errValid = learner.err(Xva[:,2], Yva)

Xtr = Xtr[:,2:]
Xtr_class0 = Xtr[Ytr==0]
Xtr_class1 = Xtr[Ytr==1]

data = np.array([Xtr_class0,Y])

print(Xtr.shape)

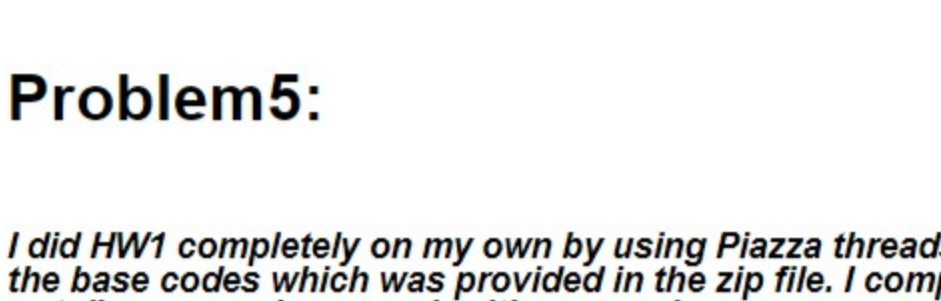
print(np.mean(Xtr, axis = 0))

bc = ml.bayes.gaussClassify(Xtr[:, 0:-2], Ytr);
ml.plotClassify2D(bc, Xtr[:, 0:-2], Ytr);
Xtr_class0.transpose()

print(errTrain)
print(errValid)
print(X[Y==0].mean(axis=0))
print(X[Y==1].mean(axis=1))
b=np.mean(Xtr_class0)
b=np.cov(Xtr_class0)

ml.plotGauss2D(a, b)

(225, 2)
[14.02464449 21.7962618 ]
0.32
0.28
[ 16.14898626 25.07251957 1926.94 ]
[51.09655676 665.25450304 655.68381395 661.67278942 655.03270718
654.03664401 659.32921024 645.33296736 649.8053967 658.18405506
656.79896487 650.45539046 653.15506344 655.12973532 654.64290185
655.16339453 650.14932592 652.14840455 675.98900382 659.26977693
656.59408079 656.22190887 657.0954563 668.45947557 653.32652193
666.80485228 647.63740881 658.81938777 657.0954563 678.23388977
661.20386612 651.09655676 673.52646224 652.68764132 649.43251925
648.80957793 666.41368164 655.49807526 647.0515745 655.89236648
661.61235066 654.45248877 656.57714004 672.29174325 656.76956238
653.88246328 663.58040357 642.95272036 658.48430021 677.98685989
654.03664401 653.44343742 652.63851768 655.93273989 654.67057218
674.70448874 651.2029516 655.12174625 651.8102234 654.19691872
658.83142933 664.78935731 654.72573185 663.84694944 653.32652193
652.53857149 653.12727469 646.30151612 678.53196973 663.58940357
655.42899253 660.45947557 659.62113611 655.93273989 654.30420021
651.75057779 647.64347447 656.22190887 655.24041756 652.53111556
658.98284052 655.06859945 658.37383229 657.7841416 659.80317764]
```



# Problem5:

I did HW1 completely on my own by using Piazza threads before I start to do my homework. In addition, I used the lecture notes and read through the best codes which was provided in the zip file. I completely followed the academic honesty guidelines which is on our canvas website and I did not discuss my homework with anyone in-person.