# MAXIMUM LIKELIHOOD ESTMATION VS PARZEN WINDOW ESTIMATION

CS662 Project 2

*Xiaobo Zhang*

## Section 1. Introduction

The goal of this project is to compare two density function estimation: maximum likelihood estimation and Parzen window estimation, under the size of training samples, variance of data and computation speed to determine which estimation works, which estimation not works and which is more advantage than the other. Section 2 will introduce maximum likelihood estimation and Parzen window estimation and how they are used for classification. Section 3 will experiment these two estimations with the size of training samples. Section 4 will experiment these two estimations with variance of data. Section 5 will experiment the computation speed of two estimations. Section 6 will conclude the result of this project.

## Section 2. Introduction of maximum likelihood estimation and Parzen window estimation

According the Bayes decision rule:

$$Given\ x\ \epsilon\ X$$

$$find\ \underset{w_i \in \{w_1, w_2, \dots w_n\}}{\operatorname{argmax}}\ \rho(w_i|x)Prob(w_i)$$

$$where\ \rho(w_i|x)\ is\ the\ distribution\ of\ X\ under\ the\ class\ w_i$$

, it uses probability density function $\rho(w_i|x)$ and priori probability $Prob(w_i)$ to calculate the posterior probability and determine which class the point(samples) belongs to. However, the probability density function is not given in real word, so we have to find some ways to make estimated parameters approximate to true parameters. In this project, there are two methods for estimation: parametric estimation and non-parametric estimation. Both methods will through observed data to construct an estimation.

For parametric estimation, the project uses maximum likelihood estimation to estimate the $\hat{\mu}$ and $\hat{\Sigma}$ as the parameters of density function for Bayes decision rule classification. The estimated $\hat{\mu}$ and $\hat{\Sigma}$ are:

$$Given\ trainnig\ data\ \{x_1, x_2, x_3, \ldots, x_n\} \in X$$

$$\hat{\mu} = \frac{1}{N} \sum_{x_i \in X} x_i \quad and \quad \hat{\Sigma} = \frac{1}{N} \sum_{x_i \in X} (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

Using estimated $\hat{\mu}$ and $\hat{\Sigma}$ as parameters in discrimination function g(x) for classification:

$$g(x) = g_1(x) - g_2(x)$$

$$= -\frac{1}{2}(x - \hat{\mu}_2)^T \hat{\Sigma}_2^{-1}(x - \hat{\mu}_2) - \frac{1}{2}(x - \hat{\mu}_1)^T \hat{\Sigma}_1^{-1}(x - \hat{\mu}_1)$$

$$+ \frac{1}{2} \ln \frac{|\hat{\Sigma}_2|}{|\hat{\Sigma}_1|} + \ln \frac{Prob(\hat{w}_1)}{Prob(\hat{w}_2)}$$

It is same as project 1, choose class 1 when g(x) > 0, and choose class 2 when g(x) < 0. The only difference is changing the true parameter with estimated parameter.

For non-parametric estimation, this project uses Parzen window estimation to estimate the density of the test point (samples) p(x) for Bayes decision rule classification. Assume there is a point $x_0$ in the region $R_n$, and the estimated $\rho(x_0|w_i)$ will be:

$$\rho(x_0|w_i) = \frac{K_n}{N * V_n}$$

*where $K_n$ is number of points in $R_n$ belong to class i and $V_n$ is the volumn of $R_n$*

The $K_n$ is calculated by a window function $\varphi$:

$$\varphi(u) \begin{cases} 1 & if\ |u_1| < \frac{1}{2}\ and\ if\ |u_2| < \frac{1}{2} \ldots\ldots\ and\ |u_n| < \frac{1}{2} \\ 0\ else \end{cases}$$

If the window function is hypercube $R^n$, then $K_n$ and $V_n$ will be:

$$K_n = \sum_{x_i \in \text{ all training data}} \varphi(\frac{x_i - x_0}{h})$$
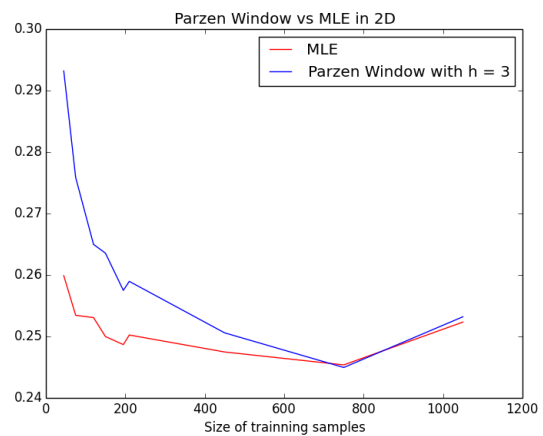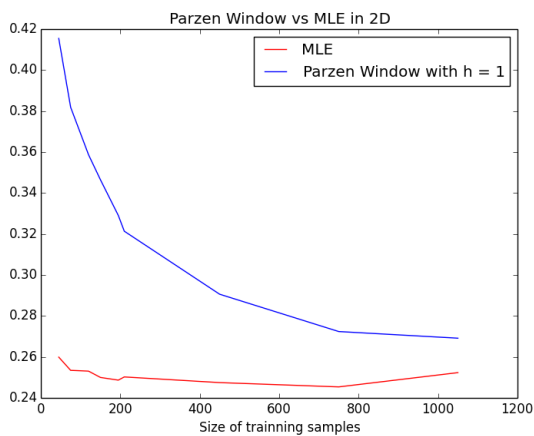
$$V_n = (h)^n$$

Using estimated density function $\rho(x_0|w_i)$ and prior probability $Prob(w_i)$ as parameters for Bayes decision classification.
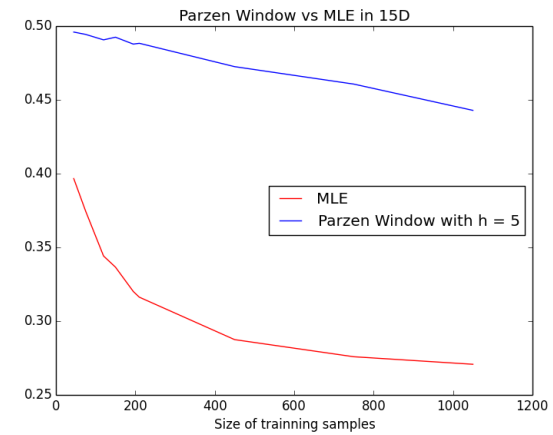
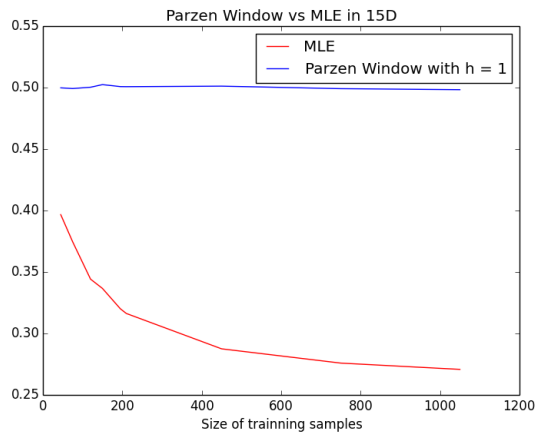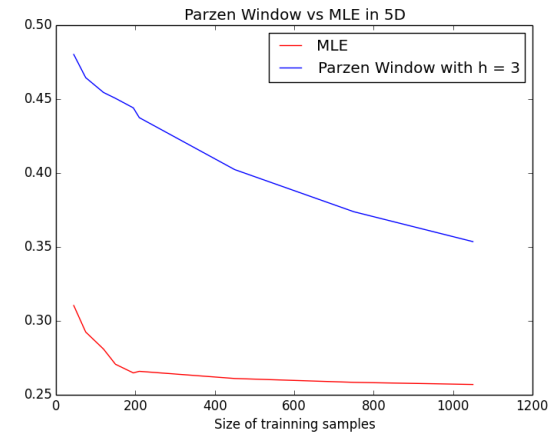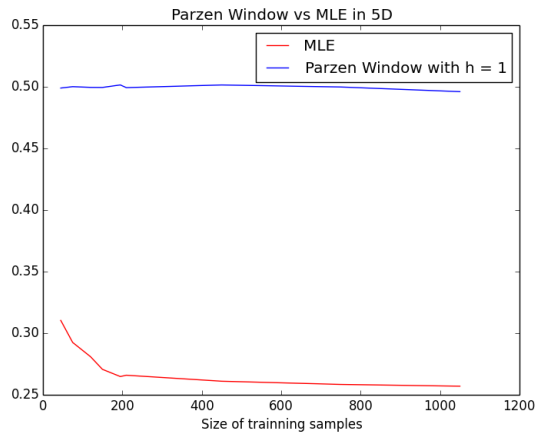## Section 3. Size of training samples

This experiment has several assumptions:

1. The data is multivariate Gaussian distribution.

2. The size of data is 3000.

3. $\Sigma_1 = \Sigma_2$.

4. The mean distance of two class is 3.

5. The variance of each dimension is less than 10.

6. The priori probability P(w1) = P(w2) = 0.5

The experiment will compare the error rate for each method by different size of training samples and different dimension.

Given above graph, the maximum likelihood estimation has lower error rate than Parzen window estimation when the size of training samples is small. With the size of training samples increasing, the error rate of Parzen window estimation decrease faster than maximum likelihood estimation. It is because

more and more training samples in the space makes the estimated density approximate to true density. When the size of training samples is really large, the error rate of two estimation methods will converge at some point.

There is an interesting thing is when the dimension increase, the error rate of Parzen window estimation always stays at 0.5 given the height of window h = 1. It perhaps the volume of the window is too small to hold enough points. If increasing the window size, the error rate will decrease with the size of training samples increase as usual. Choosing the right window size will influence the error rate a lot.

## Section 4. Variance of data

For this experiment, there are several assumptions:

1. The data is multivariate Gaussian distribution.

2. The size of data is 3000.

3. $\Sigma_1 = \Sigma_2$.

4. The mean distance of two class is 3.

5. The priori probability P(w1) = P(w2) = 0.5

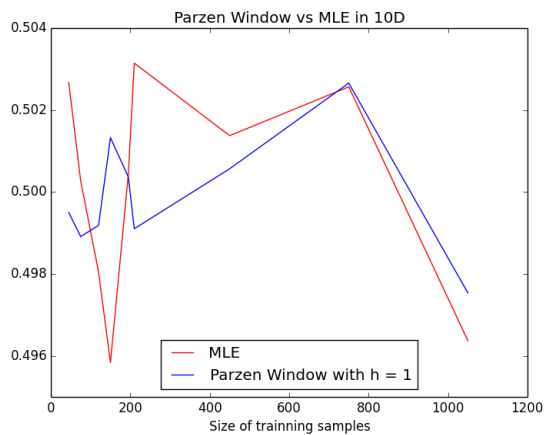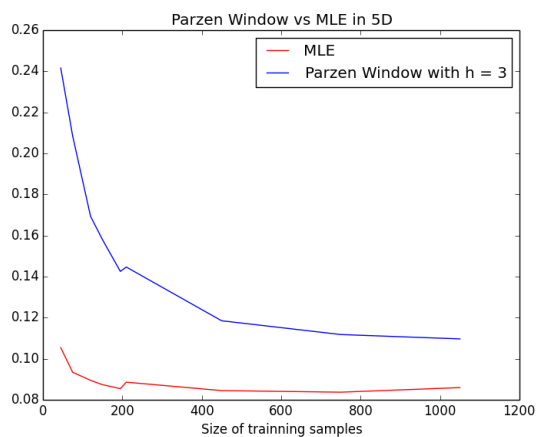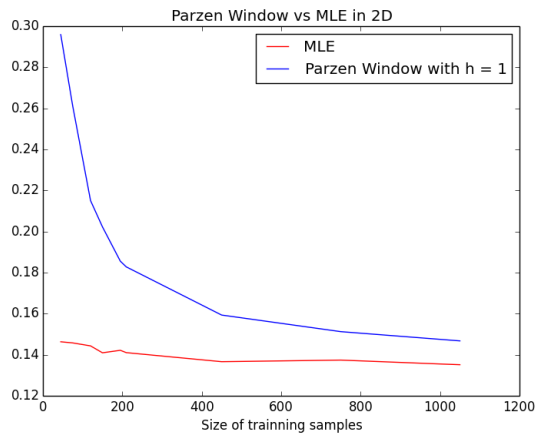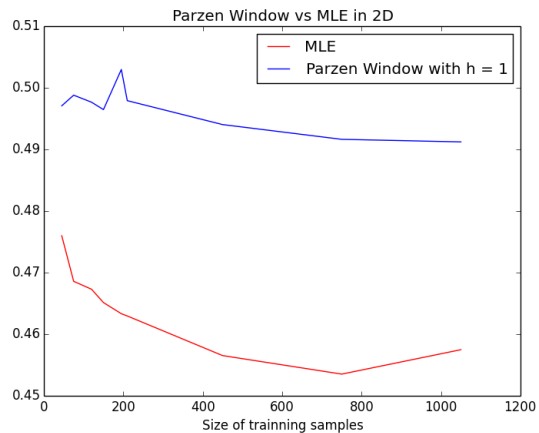The experiment will compare if estimation methods work with different variance of data in 2D, 5D, 10D, 15D.
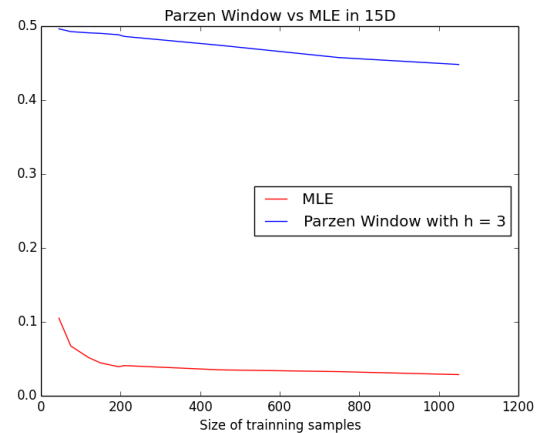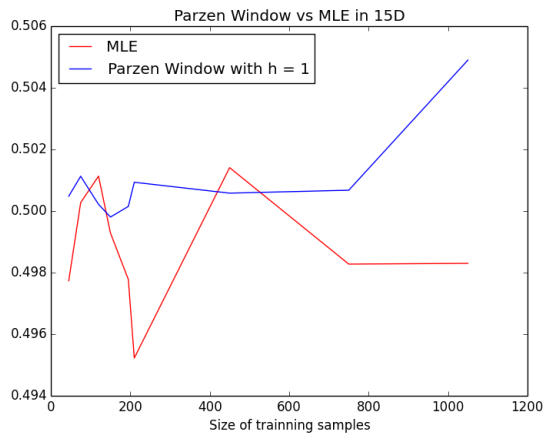
Case 1: variance between 100 and 1000 on each dimension(diagonal)

Case 2: variance between 0 and 5 on each dimension(diagonal)

Case 1                                    Case 2

When the variance of data is large, the error rate of both estimation methods is very high and doesn't decrease a lot with the size of training samples increasing. For maximum likelihood estimation, it needs more training samples to make estimated parameters $\hat{\mu}$ and $\hat{\Sigma}$ close to true parameters. For Parzen window estimation, the window needs more points in the volume to let estimated density function close to true density function.

However, the error rate of both methods is much smaller when the variance of data is small. It probably because small variance data are more concrete, and the estimator can tell the entire distribution of the data without too many training samples and output parameters approximate to true parameters. The small variance data works for both estimation methods and it performs better on maximum likelihood estimation than Parzen window estimation when the size of training samples is small. The error rate of Parzen window estimation is not as low as the maximum likelihood estimation when the size of training samples is small, but the error rate decrease when the size of training samples increasing. It is because there are enough points (training samples) to tell entire distribution of data.

## Section 5. Computation speed

In this experiment, it has several assumptions:

1. The data is multivariate Gaussian distribution.

2. The size of data is 3000.

3.The probability of training samples is 0.5 of entire data.

4. $\Sigma_1 = \Sigma_2$.

5. The mean distance between 2 class is 6.

6. The priori probability P(w1) = P(w2) = 0.5

The experiment will compare the computation time of each method in 2D, 5D, 10D and 25D.

In each dimension, we can find that the computation cost of Parzen window estimation is higher than maximum likelihood estimation obviously. It is because each test point(sample) has to compute distance with all the training samples to get the density p(x). The computation cost will grow as the size of training samples increasing, but maximum likelihood estimation just uses estimated parameters $\hat{\mu}$ and $\hat{\Sigma}$, and discrimination function directly for classification. The size of training samples doesn't have strong influence to maximum likelihood estimation on computation cost, but the Parzen window estimation does.



For each estimation method, we can find that the computation time will increase with the size of data increasing. In addition, the computation cost of each estimation method will increase with dimension of the data increasing. According the graph, maximum likelihood estimation method grows linearly with the size of data increase, but for Parzen window estimation, it grows exponentially. So this experiment can conclude that the size of data highly influences the computation speed, especially for Parzen window estimation.

## Section 6. Conclusion

This project experimented performance of two estimation methods: maximum likelihood estimation and Parzen window estimation, through the size of training samples, variance of data and computation speed. We can found that both estimation methods work good with small variance data, not good for large variance data. The maximum likelihood estimation works better than Parzen window estimation when the size of training samples is small, because little training samples cannot help Parzen window estimation figure out the distribution of data very well. Parzen window estimation will works well when distribution of data is clear, which means there are many training samples on the space. The maximum likelihood estimation is much faster than Parzen window estimation on computation speed. It is because Parzen window needs to scan all the training samples for just one test samples, but maximum likelihood estimation doesn't have to scan all training samples. Parzen window estimation needs a right window size h to make the classification good, the size of h shouldn't be too small, and shouldn't be too large either.

## Reference:

1. Professor's lecture notes

2. Kernel density estimation via the Parzen-Rosenblatt window method

http://sebastianraschka.com/Articles/2014_kernel_density_est.html

**mle.py**

```python
1. import numpy as np
2.
3. def estimate_mu(data):
4.     data =  np.array(data)
5.     return np.mean(data,axis=0)
6.
7. def estimate_covariance(data):
8.     data =  np.array(data)
9.     return np.cov(data.T)
```

**parzen.py**

```python
1. #parzen.py
2. #function of parzen window
3. import math
4. import numpy as np
5.
6. def window_function(sample,test_data,h):
7.     #print sample
8.     #print "x",test_data
9.     for x,x_i in zip(sample,test_data):
10.
11.             if abs(x-x_i) > float(h)/2:
12.                 return 0
13.         return 1
14.
15.     def estimate_single_px(samples,test_data,h):
16.         inRegion = 0
17.         outRegion = 0
18.
19.         for sample in samples:
20.             if window_function(sample,test_data,h) == 1:
21.                 inRegion  += 1
22.             else:
23.                 outRegion += 1
24.
25.         return float(inRegion)/((inRegion+outRegion)*(h**len(test_data)))
26.
27.     def get_px(samples,test_datas,h):
28.         px =[]
29.
30.         for x in test_datas:
31.             x =np.array(x)
32.
```

```python
33.              val = estimate_single_px(samples,x,h)
34.              px.append(val)
35.          return px
```

plotting.py
```python
1. import matplotlib.pyplot as plt
2.
3. def plot(title, parzen, mle,xtitle,ytitle,x_attr,h):
4.     plt.clf()
5.     plt.title(title)
6.     plt.xlabel(xtitle)
7.     plt.plot(x_attr,mle,label="MLE",color = 'r')
8.     plt.plot(x_attr,parzen,label="Parzen Window with "+str(h),col
   or = 'b')
9.     plt.legend(loc='best')
10.        #plt.ylim(0,0.8)
11.        plt.show()
```

testing_mle.py
```python
1. import mle
2. import numpy as np
3. import random
4. import parzen
5. from numpy import matrix
6. from mpl_toolkits.mplot3d import Axes3D
7. import matplotlib.pyplot as plt
8. from matplotlib.mlab import bivariate_normal
9. import math
10.    from numpy.linalg import linalg
11.    import time
12.
13.
14.    def compute_gx(mu_1,mu_2,cov_1,cov_2,w1,w2,x):
15.
16.        mu_1 =matrix(mu_1)
17.        mu_2 = matrix(mu_2)
18.        cov_1 = matrix(cov_1)
19.        cov_2 = matrix(cov_2)
20.
21.
```

```
22.         const = 0.5 * math.log(linalg.det(cov_2)/linalg.det(cov_
    1)) + math.log(float(w1)/w2)
23.
24.         #print "const",const
25.         #print "xx",(0.5* (x-mu_2) * cov_2.I * (x-
    mu_2).T ) - (0.5* (x-mu_1) * cov_1.I * (x-mu_1).T )
26.         gx = (0.5 * (x-mu_2) * cov_2.I * (x-
    mu_2).T ) - (0.5* (x-mu_1) * cov_1.I * (x-mu_1).T )
27.         #print gx
28.         v = gx.tolist()
29.         return v[0][0]+const
30.
31.
32.
33.     def testing(train_samples_1,test_samples_1,train_samples_2,t
    est_samples_2):
34.
35.         now = time.time()
36.
37.
38.         estimate_mu_1 = mle.estimate_mu(train_samples_1)
39.         estimate_cov_1 = mle.estimate_covariance(train_samples_1
    )
40.
41.         estimate_mu_2 = mle.estimate_mu(train_samples_2)
42.         estimate_cov_2 = mle.estimate_covariance(train_samples_2
    )
43.
44.         w1 = len(train_samples_1)/float(len(train_samples_1)+len
    (train_samples_2))
45.         w2 = len(train_samples_2)/float(len(train_samples_1)+len
    (train_samples_2))
46.
47.         error = 0
48.         for d in test_samples_1:
49.             val = compute_gx(estimate_mu_1 ,estimate_mu_2,estima
    te_cov_1,estimate_cov_2,w1,w2,d)
50.             if val > 0:
51.                 pass
52.             else:
53.                 error += 1
54.
55.         for d in test_samples_2:
56.             val = compute_gx(estimate_mu_1 ,estimate_mu_2,estima
    te_cov_1,estimate_cov_2,w1,w2,d)
57.             if val > 0:
```

```
58.                    error += 1
59.              else:
60.                    pass
61.
62.          res =  float(error)/float(len(test_samples_1)+len(test_s
   amples_2))
63.
64.          return res
65.          #errors.append(res)
66.          #times.append(time.time()-now)
67.
68.          #print errors
69.          #print "error rate:",sum(errors)/len(errors)
70.          #print "average_time:",sum(times)/len(times)
71.          #return sum(errors)/len(errors)
72.          #return sum(times)/len(times)
```

testing_parzen.py

```
1. import numpy as np
2. import random
3. import parzen
4. import time
5. from numpy import matrix
6.
7.
8. def classify(x,train_samples_1,train_samples_2,w1,w2,h):
9.     px_1 = parzen.estimate_single_px(train_samples_1,x,h)
10.        px_2 = parzen.estimate_single_px(train_samples_2,x,h)
11.        if px_1*w1 > px_2*w2:
12.            return 1
13.        else:
14.            return 2
15.
16.
17.     def testing(train_samples_1,test_samples_1,train_samples_2,t
   est_samples_2,h):
18.        errors = []
19.        times = []
20.
21.        now = time.time()
22.
23.        w1 = len(train_samples_1)/float(len(train_samples_1)+len
   (train_samples_2))
24.        w2 = len(train_samples_2)/float(len(train_samples_1)+len
   (train_samples_2))
25.        error = 0
```

```python
26.          for x in test_samples_1:
27.              label = classify(x,train_samples_1,train_samples_2,w
      1,w2,h)
28.              if label == 1:
29.                  pass
30.              else:
31.                  error +=1
32.
33.          for x in test_samples_2:
34.              label = classify(x,train_samples_1,train_samples_2,w
      1,w2,h)
35.              if label == 2:
36.                  pass
37.              else:
38.                  error +=1
39.
40.
41.          res = error/float(len(test_samples_1)+len(test_samples_2
      ))
42.          return res
43.          #errors.append(res)
44.          #times.append(time.time()-now)
45.
46.          #print errors
47.          #print "error rate:",sum(errors)/len(errors)
48.          #print "average_time:",sum(times)/len(times)
49.          #return sum(errors)/len(errors)
50.          #return sum(times)/len(times)
```

test.py

```python
1. import testing_mle
2. import testing_parzen
3. import plotting
4. from matplotlib.mlab import bivariate_normal
5. import matplotlib.pyplot as plt
6. import random
7. from numpy import matrix
8. import numpy as np
9.
10.      def create_variance_matrix(d):
11.          cov = []
12.          for i in range(d):
13.              vec = []
14.              for j in range(d):
15.                  if i == j:
16.                      vec.append(random.randint(1, 10))
```

```python
                        #vec.append(1)
                    else:
                        vec.append(0)
                cov.append(vec)

        for i in range(d):
            for j in range(i+1,d,1):
                if i != j:
                    #cov[i][j] = random.randint(1, 3)
                    cov[i][j] = round(random.random(),2)
                    cov[j][i] = cov[i][j]
        return matrix(cov)

    def generate_random_point(dimesion,distance):
        p1 = []
        p2 = []
        for t in range(dimesion):
            x = random.random()*10-10
            p1.append(x)
            p2.append(x)

        l = len(p1)
        i = int(random.random()*l)
        p2[i] += distance
        return p1,p2

    def generate_normal_data(mean,cov,size):
        data = np.random.multivariate_normal(mean, cov, size)
        return data

    def generate_possion_data(size,dimension):
        data = np.random.poisson(10,size=(size,dimension))
        return data

    def generate_uniform_data(start,end,size,dimesion):
        data = np.random.uniform(start,end,size=(size,dimesion))

        return data

    def sampling(data_1,data_2,prob):
        train_1 = []
        test_1 = []
        train_2 = []
        test_2 = []
        for d in data_1:
            if random.random() < prob:
```

```python
                    train_1.append(d)
            else:
                    test_1.append(d)

        for d in data_2:
            if random.random() < prob:
                    train_2.append(d)
            else:
                    test_2.append(d)
        return train_1,test_1,train_2,test_2


    #testing methods by computation time
    dimensions = [2,5,10,25]
    sizes = [100,200,500,1000,2000]

    for d in dimensions:
        parzen = []
        mle = []

        for size in sizes:
            print d,size
            parzen.append(testing_parzen.testing(d,size,0.5,1))

            mle.append(testing_mle.testing(d,size,0.5))
        print "d"
        title = "Parzen Window vs MLE in "+str(d)+"D"
        title = str(title);
        xtitle = "Size of data"
        ytitle = "Computation time(sec)"
        x_attr = sizes

        plotting.plot(title, parzen, mle,xtitle,ytitle,x_attr)


    colors = ['r', 'b', 'g', 'k', 'm']
    i = 0
    for d in dimensions:
        parzen = []
        mle = []

        for size in sizes:
            print d,size
            parzen.append(testing_parzen.testing(d,size,0.5,1))

            #mle.append(testing_mle.testing(d,size,0.5))
```

```python
106.        plt.plot(sizes,parzen,label=str(d)+"D",color=colors[i])

107.        #plt.plot(sizes,mle,label=str(d)+"D",color=colors[i])
108.
109.        i+= 1
110.
111.    title = "Parzen Window"
112.    #title = "MLE"
113.    plt.title(title)
114.    xtitle = "Size of data"
115.    ytitle = "Computation time(sec)"
116.    plt.xlabel(xtitle)
117.    plt.ylabel(ytitle)
118.    plt.legend(loc='best')
119.
120.    plt.show()
121.
122.
123.
124.    #testing by trainning sampling size and variance of data
125.    dimensions = [10,15,25]
126.    size = 1500
127.    probs = [0.03,0.05,0.08,0.1,0.13,0.14,0.3,0.5,0.7]
128.    for d in dimensions:
129.
130.        #generate parameter
131.        mean_1,mean_2 = generate_random_point(d,3)
132.        cov_1 = create_variance_matrix(d)
133.        cov_2 = cov_1
134.
135.        #generate data
136.        data_1 = generate_normal_data(mean_1,cov_1,size)
137.        #data_1 = generate_uniform_data(-10,10,size,d)
138.        data_2 = generate_normal_data(mean_2,cov_2,size)
139.        #data_2 = generate_uniform_data(-10,10,size,d)
140.
141.        parzen_1 = []
142.        parzen_3 = []
143.
144.        mle = []
145.
146.        for prob in probs:
147.            print d,size*prob
148.            parzen_errors_1 = []
149.            parzen_errors_3 = []
150.
```

```python
151.            mle_errors = []
152.            for t in range(10):
153.                train_samples_1,test_samples_1,train_samples_2,t
    est_samples_2 = sampling(data_1,data_2,prob)
154.                parzen_errors_1.append(testing_parzen.testing(tr
    ain_samples_1,test_samples_1,train_samples_2,test_samples_2,1))
155.                parzen_errors_3.append(testing_parzen.testing(tr
    ain_samples_1,test_samples_1,train_samples_2,test_samples_2,3))
156.                mle_errors.append(testing_mle.testing(train_samp
    les_1,test_samples_1,train_samples_2,test_samples_2))
157.            parzen_1.append(sum(parzen_errors_1)/len(parzen_erro
    rs_1))
158.            parzen_3.append(sum(parzen_errors_3)/len(parzen_erro
    rs_3))
159.
160.            mle.append(sum(mle_errors)/len(mle_errors))
161.        print parzen_1
162.        print parzen_3
163.        print mle
164.
165.        title = "Parzen Window vs MLE in "+str(d)+"D"
166.        title = str(title);
167.        xtitle = "Size of trainning samples"
168.        ytitle = "error rate"
169.        x_attr = [x*size for x in probs]
170.
171.        plotting.plot(title, parzen_1, mle,xtitle,ytitle,x_attr,
    "h = 1")
172.        plotting.plot(title, parzen_3, mle,xtitle,ytitle,x_attr,
    "h = 3")
```