

# Statistical Machine Learning: Exercise 2



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

Group A: Kexin Wang (2540047), Paul Philipp Seitz (2337506)

Summer Term 2022

## Task 1: Optimization (25 Points)

### 1a) Numerical Optimization (15 Points)

Derivative of the function:

$$\begin{aligned}\frac{\partial f}{\partial x_j} &= \sum_{i=1}^N 200(x_i - x_{i-1}^2)(\epsilon_{i,j} - 2x_{i-1}\epsilon_{i-1,j}) - 2(1 - x_{i-1})\epsilon_{i-1,j} \\ &= 200(x_j - x_{j-1}^2) - 400x_j(x_{j+1} - x_j^2) - 2(1 - x_j)\end{aligned}$$

but the first and the last x are exception:

$$\frac{\partial f}{\partial x_0} = -400x_0(x_1 - x_0^2) - 2(1 - x_0)$$

$$\frac{\partial f}{\partial x_{N-1}} = 200(x_{N-1} - x_{N-2}^2)$$

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def rosenbrock(x):
5     sum = 0
6     for i in range(len(x)-1):
7         sum += 100*(x[i+1] - x[i]**2)**2 + (x[i] - 1)**2
8
9     return sum
10
11 def rosenbrock_gradient(x, n):
12     grad = np.zeros(n)
13
14     for i in range(n-1):
15         grad[i] = 2*(x[i]-1) - 400*x[i]*(x[i+1]-x[i]**2)
16
17         if i != 0:
18             grad[i] += 200*(x[i]-x[i-1]**2)
19
20     grad[n-1] = 200*(x[n-1]-x[n-2]**2)
21     return grad
22
23
24 def loss(x):
25     return np.linalg.norm(rosenbrock(x))
26
```

---

```

27 def gradient_descent(x0, alpha, max_iter):
28     x = x0
29     history = [loss(x)]
30
31     for i in range(max_iter):
32         x = x - alpha*rosenbrock_gradient(x, len(x))
33         current_loss = loss(x)
34         history.append(current_loss)
35
36     return x, history
37
38 def optimize(alpha):
39     np.random.seed(7431)
40     x0 = np.random.rand(20)
41     max_iter = 10000
42     x, history = gradient_descent(x0, alpha, max_iter)
43
44     # Plot history
45     plt.figure(figsize=[8,6])
46     plt.plot(history, linewidth=3.0)
47     plt.legend(['Loss with alpha=%.3f' % alpha], fontsize=18)
48     plt.title('GD Learning Curve', fontsize=16)
49     plt.show()
50
51 # Higher learning rate leads to divergence
52 # Lower learning rate leads to slower progress and therefor higher loss after 10000 steps
53 optimize(1e-3)

```

---

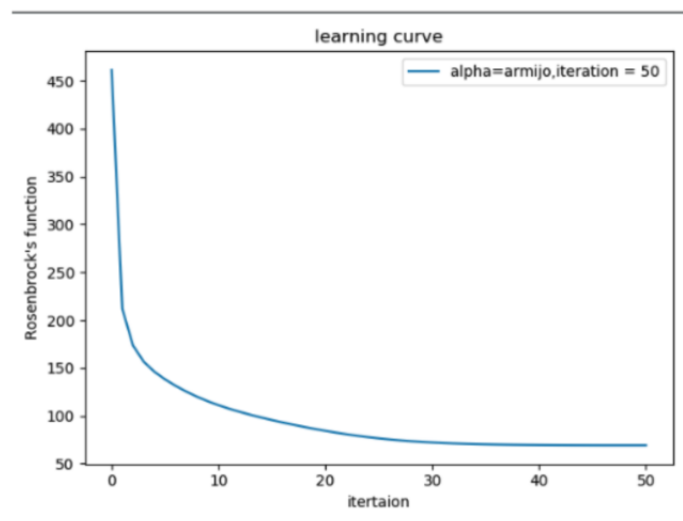


Figure 1: gradient descent algorithm after 10,000 steps on Rosenbrock's function

---

## 1b) Gradient Descent Variants (10 Points)

---

### 1).A. Batch Gradient Descent, BGD<sup>[3]</sup> :

They process all samples simultaneously in a large batch, i.e. use all samples for gradient update at each iteration.

**Advantages:**

- During training, a fixed learning rate is used and there is no need to worry about the appearance of learning rate recession.
- The direction determined by the full data set is more representative of the total sample and thus more accurate towards the direction where the extremes are located. When the objective function is convex, it must converge to the global minimum, and if the objective function is non-convex it converges to the local minimum.
- It has an unbiased estimate of the gradient. The more samples there are, the lower the standard deviation.
- One iteration is computed for all samples, and at this time, parallelism is achieved by using vectorization for the operation.

#### **Disadvantages:**

- Although vectorization is used in the computation process, it still takes a lot of time to traverse all the samples, especially when the data set is large (millions or even hundreds of millions), it is a bit difficult to complete.
- Since each update occurs after traversing the entire sample, it is only after all of them are completed that some examples may be redundant and not very useful for parameter updating. In terms of the number of iterations, the number of BGD iterations is relatively small.

In terms of the number of iterations, the number of BGD iterations is relatively small.

#### **B. Stochastic Gradient Descent, SGD<sup>[1]</sup> :**

Stochastic Gradient Descent uses a sample for parameter update at each iteration (mini-batch size = 1).

#### **Advantages:**

- Noise is added to the learning process, which improves the generalization error.
- Since the loss function is not on all training data, but on a randomly optimized loss function on one training data in each iteration, the parameters are updated much faster in each round.

#### **Disadvantages**

- Does not converge and fluctuates around the minimum value.
- Cannot use vectorization calculations in one sample, and the learning process becomes slow.
- A single sample does not represent the trend of the whole sample.
- When local minima or saddle points are encountered, SGD gets stuck at a gradient of 0.

#### **C. Mini-batch Gradient Descent, MBGD:**

MBGD is in between the above two, using more than one but not all training samples. In each step of the algorithm, we randomly draw a small batch of samples from a training set of  $m$  samples (which have been disordered). Usually uses a batch of 32 or 64 samples.

#### **Advantages**

- The computation is faster than Batch Gradient Descent because the update can be performed by traversing only some samples.
- Random sample selection helps to avoid repeating redundant samples and samples that contribute less to parameter updates.
- Using one batch at a time can greatly reduce the number of iterations required for convergence, and at the same time can make the converged result closer to the effect of gradient descent.

#### **Disadvantages**

- During the iteration, the learning process fluctuates because of the presence of noise. Therefore, it hovers in the region of the minimum and does not converge.
- The learning process will have more oscillations. To get closer to the minimum, the learning rate decay needs to be increased to reduce the learning rate and avoid excessive oscillations.

- choice of batch size may cause some problems.

If the sample size is relatively small, the batch gradient descent algorithm is used. If the sample is too large, or online algorithm, use stochastic gradient descent algorithm. In the practical general case, use the small batch gradient descent algorithm.

## 2). Momentum<sup>[2]</sup> :

Momentum methods are designed to accelerate learning (speed up gradient descent), especially when dealing with high curvature, small but consistent gradients, or gradients with noise. The Momentum algorithm accumulates a moving average of the previous exponential decay of the gradient and continues to move in that direction.

The idea of Momentum is simply to update a portion more of the updates from the previous iteration to smooth out the gradient for this iteration. From a physical point of view, it is like a small ball rolling down is affected by its own historical momentum, that is why it is called Momentum algorithm.

This has the direct effect of making the turnaround less dramatic when the gradient is falling, so that it can rush more smoothly and quickly to the local minimum. For example, in the SGD algorithm with momentum, when the general value of the momentum parameter is changed, it may be 2, 10, or 100 times faster than the SGD algorithm. And it does improve the convergence rate of some algorithms, for example, in the case of convex batch gradients, Nesterov momentum improves the error convergence rate. But not all algorithms are applicable, which means momentum is not always useful. for instance, in the case of stochastic gradient, Nesterov momentum does not improve the convergence rate.

---

## Task 2: Density Estimation (35 Points)

---

### 2a) Gaussian Maximization Likelihood Estimate (10 Points)

---

$$L(x_1, \dots, x_n | \mu, \Sigma) = \prod_{i=1}^n p(x_i | \mu, \Sigma) = \prod_{i=1}^n \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}$$

Using Log-Likelihood to determine ML estimates for the mean and covariance of the multivariate Gaussian distribution.

$$\tilde{L}(x_1, \dots, x_n | \mu, \Sigma) = \ln L(x_1, \dots, x_n | \mu, \Sigma) = \sum_{i=1}^n \ln p(x_i | \mu, \Sigma)$$

$$\ln p(x_i | \mu, \Sigma) = -\frac{1}{2} (D \ln(2\pi) + \ln(|\Sigma|) + (x_i - \mu)^T \Sigma^{-1} (x_i - \mu))$$

$$\tilde{L}(x_1, \dots, x_n | \mu, \Sigma) = -\frac{n}{2} (D \ln(2\pi) + \ln(|\Sigma|)) - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

We begin with  $\mu_{MLE}$ :

$$\frac{\partial \tilde{L}}{\partial \mu} = -\frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \mu} [(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)] \stackrel{(86)}{=} \sum_{i=1}^n \Sigma^{-1} (x_i - \mu) = 0$$

$$\sum_{i=1}^n \Sigma^{-1} x_i = n \Sigma^{-1} \mu$$

$$\mu_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

Now for  $\Sigma_{MLE}$ :

$$\frac{\partial \tilde{L}}{\partial \Sigma} = -\frac{n}{2} \frac{\partial}{\partial \Sigma} \ln(|\Sigma|) - \frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \Sigma} [(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)]$$

$$\begin{aligned}
-\frac{n}{2} \frac{\partial}{\partial \Sigma} \ln(|\Sigma|) &= -\frac{n}{2} \frac{1}{|\Sigma|} \frac{\partial}{\partial \Sigma} |\Sigma| \stackrel{(49)}{=} -\frac{n}{2} \frac{1}{|\Sigma|} |\Sigma| (\Sigma^{-1})^T \stackrel{(\Sigma^{-1} \text{ symmetric})}{=} -\frac{n}{2} \Sigma^{-1} \\
\frac{\partial}{\partial \Sigma} [(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)] &\stackrel{(61)}{=} -\Sigma^{-T} (x_i - \mu) (x_i - \mu)^T \Sigma^{-T} \stackrel{(\Sigma^{-1} \text{ symmetric})}{=} -(x_i - \mu) (x_i - \mu)^T \Sigma^{-1} \Sigma^{-1} \\
\frac{\partial \tilde{L}}{\partial \Sigma} &= -\frac{n}{2} \Sigma^{-1} + \frac{1}{2} \left( \sum_{i=1}^n (x_i - \mu) (x_i - \mu)^T \right) \Sigma^{-1} \Sigma^{-1} = 0 \\
\left( \sum_{i=1}^n (x_i - \mu) (x_i - \mu)^T \right) \Sigma^{-1} \Sigma^{-1} &= n \Sigma^{-1} \\
\Sigma_{MLE} &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu) (x_i - \mu)^T
\end{aligned}$$

---

## 2b) Prior Probabilities (2 Points)

---

The prior probability  $p(A) = \frac{\text{count}(A)}{\text{count}(A) + \text{count}(B)}$  of each class A from the dataset is computed in the python code below.

```

1 def get_priors(d1,d2):
2     #counting the number of class data point than normalizing
3     #d1,d2 are the data load from txt files.
4     p1 = 1. * d1.shape[0] / (d1.shape[0] + d2.shape[0])
5     p2 = 1. * d2.shape[0] / (d1.shape[0] + d2.shape[0])
6     return p1,p2

```

The results are shown as the following figure. The prior possibility of dataset 1 is the  $p(C_1) = 0.24$  and the prior possibility of dataset 2 is  $p(C_2) = 0.76$

```

0.23847695390781562
0.7615230460921844

```

Figure 2

---

## 2c) Biased ML Estimate (5 Points)

---

The bias of an estimator is the difference between this estimator's expected value and the true value of the parameter being estimated. We construct some estimator  $\hat{\theta}$  that maps observed data to values that we hope are close to  $\theta$ . The bias of  $\hat{\theta}$  relative to  $\theta$  is defined as

$$\text{Bias}(\hat{\theta}, \theta) = \text{Bias}_{\theta}[\hat{\theta}] = E_{x|\theta}[\hat{\theta}] - \theta = E_{x|\theta}[\hat{\theta} - \theta].$$

Where  $E_{x|\theta}$  denotes expected value over the distribution  $P(x|\theta)$

```

1 def estimate_parameters(data):
2     covar = np.zeros((data.shape[1], data.shape[1]))
3
4     N = data.shape[0]
5     mu = np.sum(data, axis=0) / N
6     for xi in data:
7         diff = xi - mu
8         covar += np.outer(diff, diff.T)
9
10    covar_unbiased = covar / (N - 1)
11    covar /= N
12
13    return mu, covar, covar_unbiased

```

---

```

14 mu1, covar1, covar_unb1 = estimate_parameters(data_1.values)
15 mu2, covar2, covar_unb2 = estimate_parameters(data_2.values)
16 print("mean1=", mu1)
17 print("biasedcovariance1=", covar1)
18 print("unbiasedcovariance1=", covar_unb1)
19 print("mean2=", mu2)
20 print("biasedcovariance2=", covar2)
21 print("unbiasedcovariance2=", covar_unb2)

```

---

After computing them, the results are given as the following figure.

```

mean1= [-0.70623122 -0.81921721]
biasedcovariance1= [[9.05734193 2.68490698]
 [2.68490698 3.60344808]]
unbiasedcovariance1= [[9.09555856 2.6962357 ]
 [2.6962357 3.61865251]]
mean2= [3.98460051 3.98500772]
biasedcovariance2= [[4.18045643 0.02797194]
 [0.02797194 2.75628915]]
unbiasedcovariance2= [[4.18596428 0.02800879]
 [0.02800879 2.75992063]]

```

Figure 3

---

## 2d) Class Density (5 Points)

---

```

1 def multivariate_gaussian(data, mu, covar):
2     """
3     return likelihood for all given samples
4     """
5     out = np.empty(data.shape[0])
6     # denominator
7     y = np.sqrt((2 * math.pi) * np.linalg.det(covar))
8
9     # compute for each datapoint
10    for i, x in enumerate(data):
11        diff = x - mu
12        out[i] = np.exp(-0.5 * diff.T.dot(np.linalg.inv(covar)).dot(diff)) / y
13
14    return out
15
16 # Visualization of class density
17 def visualize_class_density(mu, covar, data):
18     steps = 100
19
20     x_data = data[:, 0]
21     y_data = data[:, 1]
22
23     x_min = x_data.min()
24     x_max = x_data.max()
25     y_min = y_data.min()
26     y_max = y_data.max()
27
28     x = np.arange(x_min - 1, x_max + 1, (x_max - x_min + 2) / steps)
29     y = np.arange(y_min - 1, y_max + 1, (y_max - y_min + 2) / steps)
30
31     Y, X = np.meshgrid(y, x)

```

```

32 Z = np.empty((steps, steps))
33
34 for i in range(n_models):
35     for j in range(steps):
36         # construct vector with same x and all possible y to cover the plot space
37         points = np.append(X[j], Y[j]).reshape(2, x.shape[0]).T
38         Z[j] = multivariate_gaussian(points, mu[i], covar[i])
39     c_plot = plt.contour(X, Y, Z)
40     plt.clabel(c_plot, inline=1, fontsize=10)
41
42 mus = np.append(mu1, mu2).reshape(n_models, mu1.shape[0])
43 sigmas = np.append(covar_unb1, covar_unb2).reshape(n_models, covar_unb1.shape[0],
44     covar_unb1.shape[1])
44 data = np.append(data_1, data_2, axis=0)
45 priors = np.append(prior1, prior2).reshape(n_models, 1)
46
47 plt.figure(0)
48 visualize_class_density(mus, sigmas, data)
49
50 # plot the samples
51 plt.plot(data_1.values[:, 0], data_1.values[:, 1], 'co', zorder=1, color="pink",
52     label="Density1")
53 plt.plot(data_2.values[:, 0], data_2.values[:, 1], 'cx', zorder=1, label="Density2")
54 plt.xlabel('x')
55 plt.ylabel('y')
56 plt.title("2d) Class Density: Estimate of Unbiased Gaussian Distributions")
57 plt.legend()

```

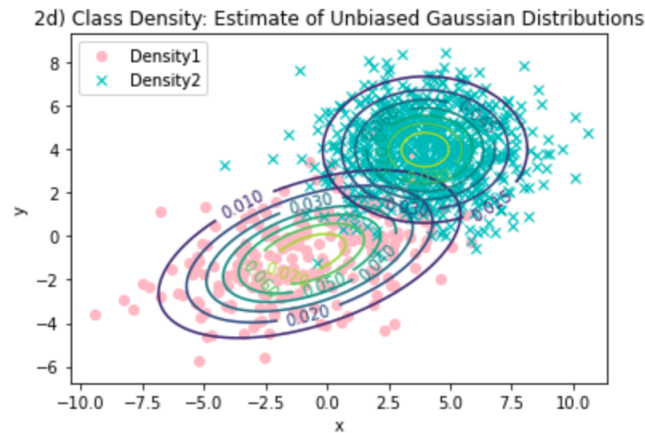


Figure 4

## 2e) Posterior (8 Points)

```

1 def visualize_posterior(mu, covar, data, prior):
2     steps = 100
3
4     # x_data = data[:, 0]
5     # y_data = data[:, 1]
6
7     x_data = np.arange(-15, 15, 30/998)
8     y_data = data[:, 1]
9

```

```

10 x_min = x_data.min()
11 x_max = x_data.max()
12 y_min = y_data.min()
13 y_max = y_data.max()
14
15 x = np.arange(x_min - 1, x_max + 1, (x_max - x_min + 2) / steps)
16 y = np.arange(y_min - 1, y_max + 1, (y_max - y_min + 2) / steps)
17
18 Y, X = np.meshgrid(y, x)
19 Z = np.empty((steps, steps))
20
21 for i in range(n_models):
22     for j in range(steps):
23         # construct vector with same x and all possible y to cover the plot space
24         points = np.append(X[j], Y[j]).reshape(2, x.shape[0]).T
25         Z[j] = multivariate_gaussian(points, mu[i], covar[i]) * prior[i]
26     c_plot = plt.contour(X, Y, Z)
27     plt.clabel(c_plot, inline=1, fontsize=10)
28
29 plt.figure(1)
30 visualize_posterior(mus, sigmas, data, priors)
31 plt.plot(data_1.values[:, 0], data_1.values[:, 1], 'co', zorder=1, color="grey",
32         label="Density1")
33 plt.plot(data_2.values[:, 0], data_2.values[:, 1], 'cx', zorder=1, color="yellow",
34         label="Density2")
35 plt.xlabel('x')
36 plt.ylabel('y')
37 plt.title("2e) Posterior: Posterior Distributions and Decision Boundary")
38 plt.legend()
39 plt.show()

```

The below plot shows the posterior distribution of each class  $P(C_i|x)$  and the decision boundary.

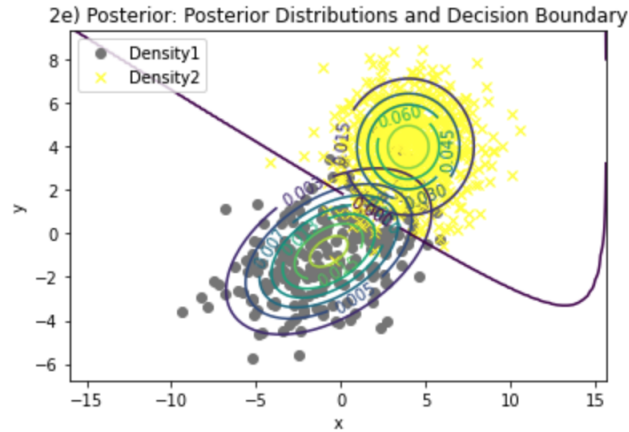


Figure 5

As shown in the figure, the decision boundary is not linear because the covariance of Gaussian is not identical.



---

### Task 3: Non-parametric Density Estimation (20 Points)

---

#### 3a) Histogram (4 Points)

---

```
1 import numpy as np
2 from matplotlib import pyplot as plt
3
4 train_data= np.loadtxt('nonParamTrain.txt')
5 test_data= np.loadtxt('nonParamTest.txt')
6
7 binwidth = 0.02
8 n,bins,patches = plt.hist(train_data,bins=np.arange(min(train_data),max(train_data),binwidth))
```

---

for binwidth=0.02, the figure is shown as the following figure.

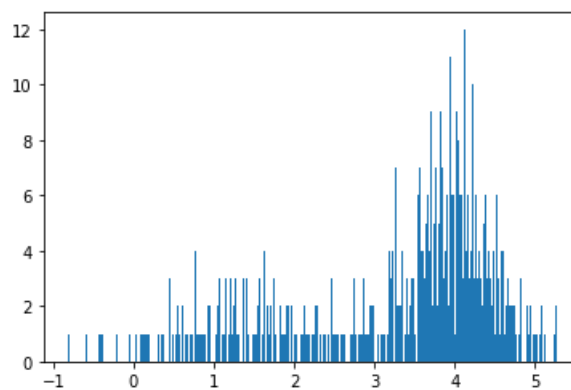


Figure 6

---

```
1 binwidth = 0.5
2 n,bins,patches = plt.hist(train_data,bins=np.arange(min(train_data),max(train_data),binwidth))
```

---

for binwidth=0.5, the figure is shown as the following figure.

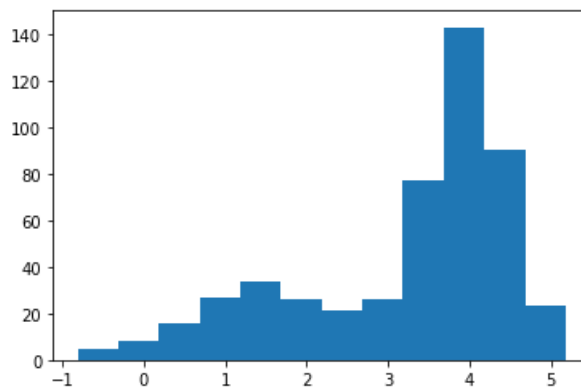


Figure 7

---

```
1 binwidth = 2.0
2 n,bins,patches = plt.hist(train_data,bins=np.arange(min(train_data),max(train_data),binwidth))
```

---

for binwidth=2.0, the figure is shown as the following figure.

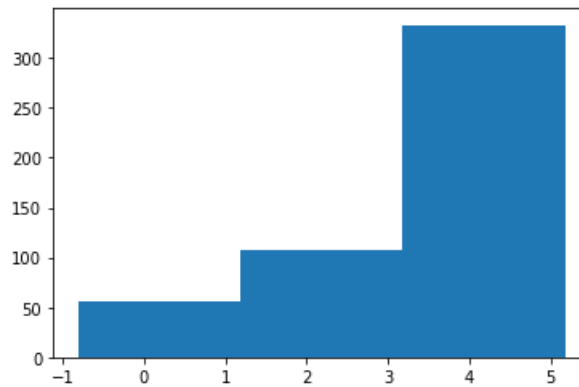


Figure 8

The bin size of 0.5 performs best. The 0.02 seems to be overfitting the data as it is not smooth enough. In contrast, 2.0 might be too smooth and could thus capture too little data and be too vague. However, we cannot easily determine which bin size results in the highest quality of estimation without knowing the underlying model.

---

### 3b) Kernel Density Estimate (6 Points)

---

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 train_data = np.loadtxt('nonParamTrain.txt')
5 test_data = np.loadtxt('nonParamTest.txt')
6
7 x = np.linspace(-4,8,1000)
8
9 def kde(x,h,data):
10     distance = (np.repeat([data],np.size(x),0).T - np.repeat([x],np.size(data),0))**2
11     p = np.sum(np.exp(-distance/2/h**2),0)
12     p = p*(1/np.sqrt(2*np.pi)/h/np.size(data))
13     return p
14
15 p1x = kde(x,0.03,train_data)
16 p2x = kde(x,0.2,train_data)
17 p3x = kde(x,0.8,train_data)
18
19 fig, ax = plt.subplots()
20 ax.plot(x,p1x)
21 ax.plot(x,p2x)
22 ax.plot(x,p3x)
23 plt.legend(['h = 0.03', 'h =0.2', 'h=0.8'])
24 plt.xlabel('x')
25 plt.ylabel('p(x)')
```

---

From the above figures, it can be seen that  $\sigma = 0.2$  performs best since it has low fluctuations with an acceptable log-likelihood.

When  $\sigma = 0.8$  the density estimation is too much smoothed and does not represent the characteristics of the original distribution.

For  $\sigma = 0.03$  even though the log-likelihood is most desirable, the estimate fluctuates strongly which results in the graph being too noisy (overfitting).

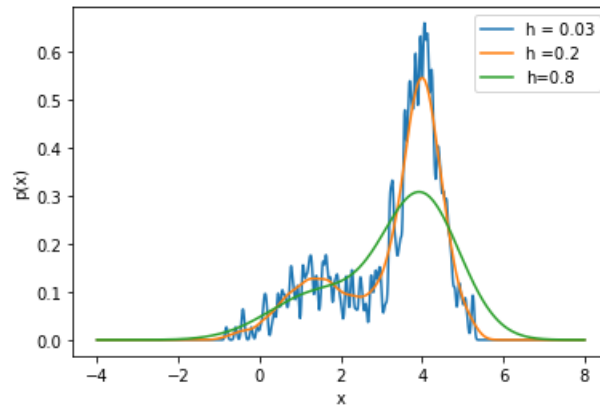


Figure 9: KDE

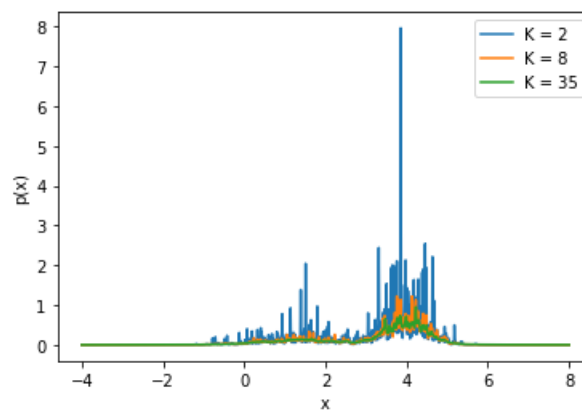


Figure 10: KDE\_Test

### 3c) K-Nearest Neighbors (6 Points)

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 train_data = np.loadtxt('nonParamTrain.txt')
5 test_data = np.loadtxt('nonParamTest.txt')
6
7 x = np.linspace(-4,8,1000)
8
9 def knn(x,K,data):
10     distance = abs(np.repeat([data],np.size(x),0).T - np.repeat([x],np.size(data),0))
11     distance.sort(axis=0)
12     V = []
13     p = []
14     for i in range(np.size(x)):
15         V.append(2*distance[K-1,i])
16         p.append(K/(np.size(data)*V[i]))
17     return np.array(p)
18
19 p1x = knn(x,2,train_data)
20 p2x = knn(x,8,train_data)
21 p3x = knn(x,35,train_data)
22

```

---

```

23 fig, ax = plt.subplots()
24 ax.plot(x,p1x)
25 ax.plot(x,p2x)
26 ax.plot(x,p3x)
27 plt.legend(['K = 2', 'K = 8', 'K = 35'])
28 plt.xlabel('x')
29 plt.ylabel('p(x)')

```

---

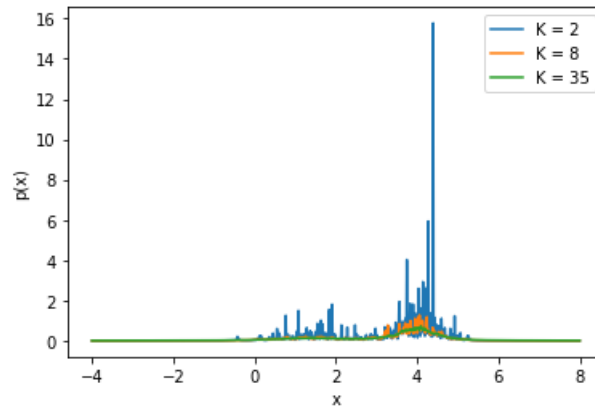


Figure 11: KNN

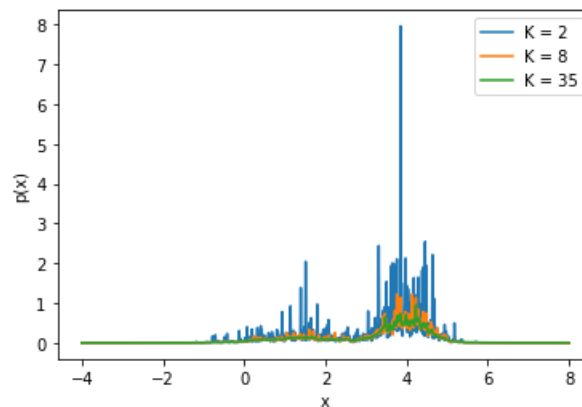


Figure 12: KNN\_Test

When  $K = 2$  we get noisy results and overfit, when  $K = 8$  we still get a bit noisy results and overfit and when  $K = 35$  the results are better.

---

### 3d) Comparison of the Non-Parametric Methods (4 Points)

---

The purpose of training machine learning is to be able to better predict new unknown samples, which means the generalization capability plays a very significant role. If the machine learning algorithm learns only for the training datasets, what it actually learns may be noise. And if the model obtained using the training datasets can get good results on the training datasets, but poor results in the face of the test datasets, then it means that the model generalization capability is very weak, for this case, we generally called **overfitting**.

What we are really looking for is the model with the best generalization capability. So measuring the model generalization capability is what makes it more meaningful to divide the datasets additional to the test datasets.

The log-likelihood of the testing data is estimated using the KDE estimators and the KNN estimator. We need to test

them on a different data set to avoid the problem of overfitting the model.  
Log-likelihoods of the estimators for both training and testing sets are given in the table below:

KDE Estimators			KNN Estimators		
	Training set	Testing Set		Training set	Testing Set
h=0.03	2425.84	-inf	k=2	-1256.04	-1672.38
h=0.2	2383.64	-7853.36	k=8	-1127.73	-1544.07
h=0.8	2305.52	853.49	k=35	-949.19	-1356.53

Table 1

The test set is used to predict the performance of the model on previously unseen data. We would choose the KNN estimator with k=2 because it has the best log-likelihood on the testing set with least degree of overfitting.

---

#### Task 4: Expectation Maximization (20 Points)

---

##### 4a) Gaussian Mixture Update Rules (2 Points)

---

The model parameters are the means, the covariance and the responsibility of each mixture component.

At  $E$  - step, we compute the posterior distribution, also called responsibility, for each mixture component and for all data point:

$$\alpha_{nj} = p(j|x_n) = \frac{\pi_j N(X_n|\mu_j, \sigma_j)}{\sum_{i=1}^N \pi_i N(X_n|\mu_i, \sigma_i)}$$

At  $M$  - step, we compute the model parameters:

$$\mu_j = \frac{\sum_{i=1}^N \alpha_{nj} x_n}{\sum_{i=1}^N \alpha_{nj}}, \sigma_j^2 = \frac{\sum_{i=1}^N \alpha_{nj} (x_n - \mu_j)(x_n - \mu_j)^T}{\sum_{i=1}^N \alpha_{nj}}, \pi_j = \frac{\sum_{i=1}^N \alpha_{nj}}{N}.$$

---

##### 4b) EM (18 Points)

---

The Expectation-Maximization algorithm for Gaussian Mixture Models is implemented and initialized uniformly. The multivariate normal probability density function is computed with the help of **multivariate\_normal** function from *scipy.stats*. The code is given below performs this operation.

---

```

1 import math
2
3 import numpy as np
4 import pandas as pd
5 from matplotlib import pyplot as plt
6
7 training_data = pd.read_csv("gmm.txt", sep=" ")
8 k = 4
9 n_iter = 30
10
11
12 def main():
13     # init
14     covar = np.array(
15         [np.identity(training_data.shape[1]) for _ in range(k)])
16     mu = np.random.uniform(training_data.min().min(), training_data.max().max(),

```

```

17         size=(k, training_data.shape[1]))
18     pi = np.random.uniform(size=(k,))
19
20     log_likelihood = np.empty((n_iter,))
21
22     for i in range(n_iter):
23         alpha = e(x=training_data.values, mu=mu, covar=covar, pi=pi)
24         mu, covar, pi = m(x=training_data.values, alpha=alpha)
25
26         # plot at given steps
27         if i + 1 in [1, 3, 5, 10, 30]:
28             plt.figure(i)
29             visualize(mu, covar, training_data.values, i)
30             log_likelihood[i] = calculate_log_likelihood(training_data.values, mu, covar, pi)
31
32             print("Finished iteration {}".format(i))
33
34     plt.figure(n_iter + 1)
35     plt.plot(log_likelihood)
36     plt.xlabel("Iterations")
37     plt.ylabel("Log-Likelihood")
38     plt.title("Log-Likelihood Progress")
39     plt.show()
40
41
42     def e(x, mu, covar, pi):
43         alpha = np.empty((k, x.shape[0]))
44         for i in range(k):
45             alpha[i] = pi[i] * multivariate_gaussian(x, mu[i], covar[i])
46
47         # sum over all models per data point
48         denominator = np.sum(alpha, axis=0)
49
50         return alpha / denominator
51
52
53     def m(x, alpha):
54         N = np.sum(alpha, axis=1) # sum over all data points per model
55
56         mu = np.zeros((k, x.shape[1]))
57         covar = np.zeros((k, x.shape[1], x.shape[1]))
58
59         for i in range(k):
60             # update mu
61             for j, val in enumerate(x):
62                 mu[i] += (alpha[i, j] * val)
63
64             mu[i] /= N[i]
65
66             # update covariance
67             for j, val in enumerate(x):
68                 diff = val - mu[i]
69                 covar[i] += alpha[i, j] * np.outer(diff, diff.T)
70
71             covar[i] /= N[i]
72
73         # update pi
74         pi = N / x.shape[0]
75
76         return mu, covar, pi

```

```

77
78
79 def multivariate_gaussian(data, mu, covar):
80
81     out = np.empty(data.shape[0])
82     denominator = np.sqrt((2 * math.pi) * np.linalg.det(covar))
83     covar_inv = np.linalg.inv(covar)
84
85     # compute for each datapoint
86     for i, x in enumerate(data):
87         diff = x - mu
88         out[i] = np.exp(-0.5 * diff.T.dot(covar_inv).dot(diff)) / denominator
89
90     return out
91
92
93 def visualize(mu, covar, data, iteration):
94
95     steps = 100
96
97     x_data = data[:, 0]
98     y_data = data[:, 1]
99
100     x_min = x_data.min()
101     x_max = x_data.max()
102     y_min = y_data.min()
103     y_max = y_data.max()
104
105     x = np.arange(x_min - 1, x_max + 1, (x_max - x_min + 2) / steps)
106     y = np.arange(y_min - 1, y_max + 1, (y_max - y_min + 2) / steps)
107
108     Y, X = np.meshgrid(y, x)
109     Z = np.empty((steps, steps))
110
111     for i in range(k):
112         for j in range(steps):
113             # construct vector with same x and all possible y to cover the plot space
114             points = np.append(X[j], Y[j]).reshape(2, x.shape[0]).T
115             Z[j] = multivariate_gaussian(points, mu[i], covar[i])
116             plt.contour(X, Y, Z, 1)
117
118     # plot the samples
119     plt.plot(x_data, y_data, 'co', zorder=1)
120     plt.xlabel('x')
121     plt.ylabel('y')
122     plt.title("Mixtures after {} steps".format(iteration + 1))
123
124
125 def calculate_log_likelihood(x, mu, covar, pi):
126     likelihood = np.empty((k, x.shape[0]))
127     for i in range(k):
128         likelihood[i] = pi[i] * multivariate_gaussian(x, mu[i], covar[i])
129
130     return np.sum(np.log(np.sum(likelihood, axis=0)))
131
132
133 if __name__ == '__main__':
134     main()

```

Plots at different iteration  $t_i \in 1, 3, 5, 10, 30$  are generated for showing the data and the mixture component. The

log-likelihood for every iteration is plotted as well.

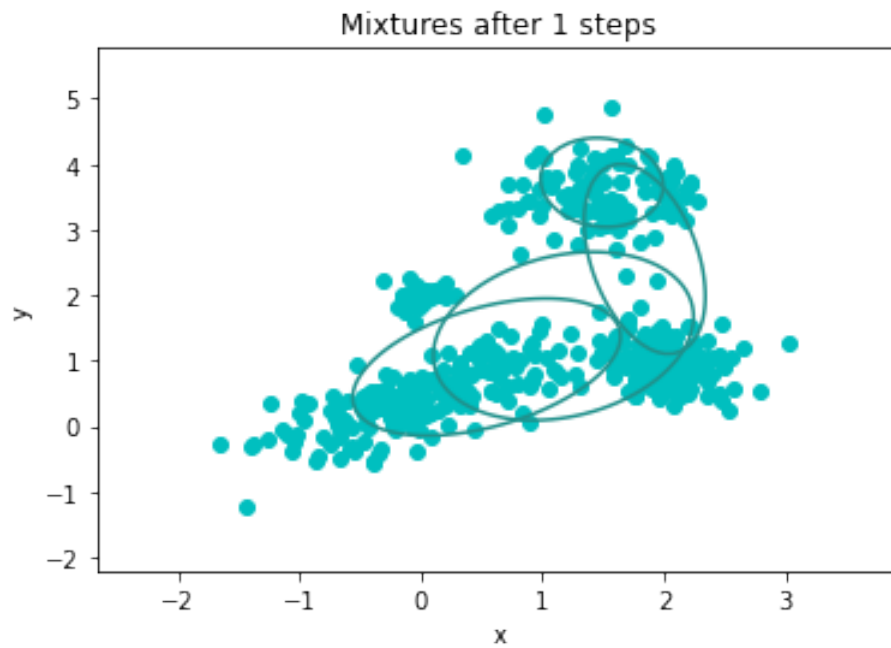


Figure 13

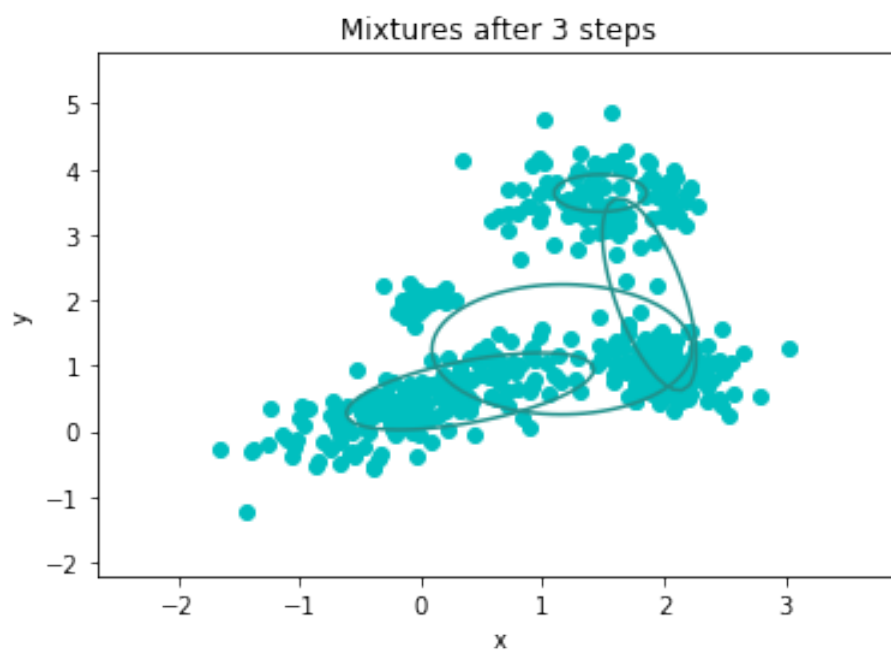


Figure 14

It is obvious that the algorithm has already to cluster all data points in the 30th iteration. Besides that, according to the log-likelihood, there is no improvement in the estimation after the 8th iteration so that the stable state has already achieved.



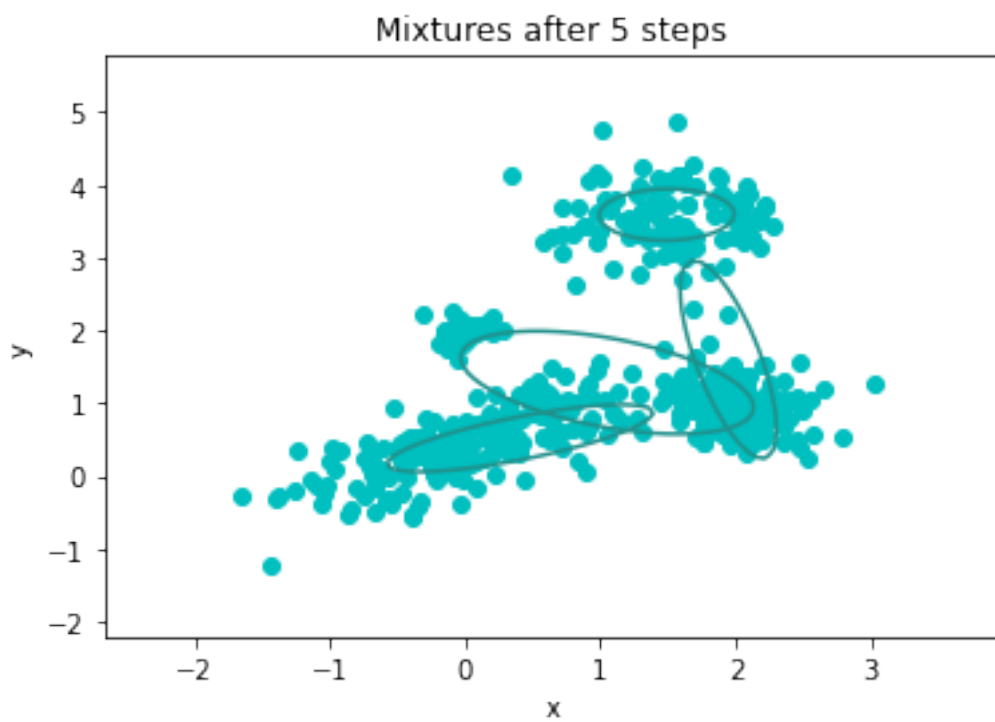


Figure 15

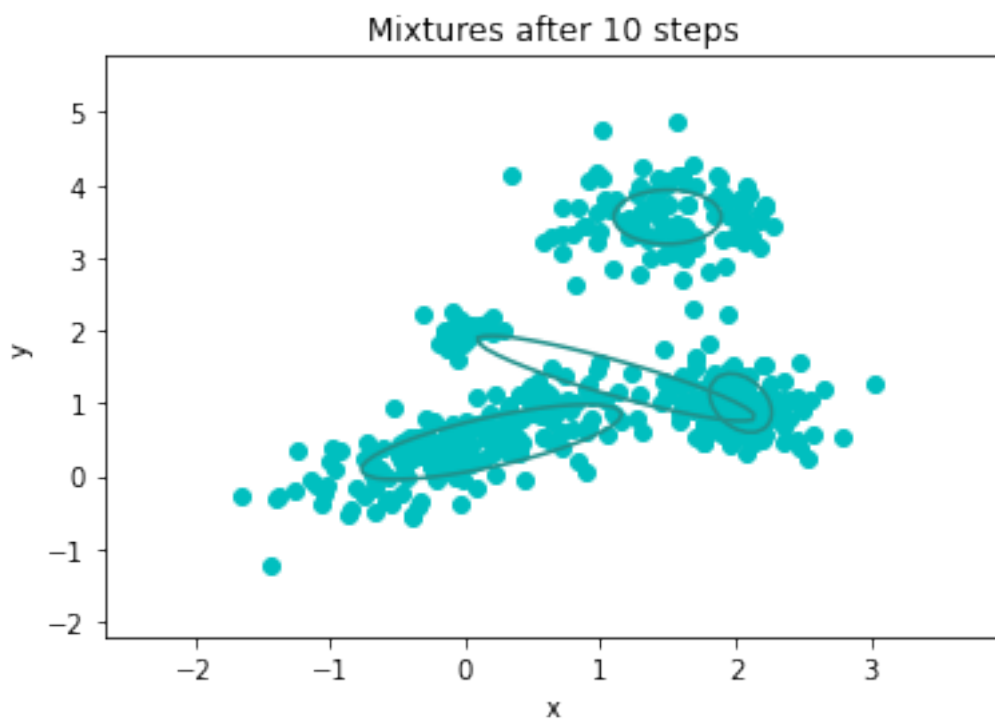


Figure 16

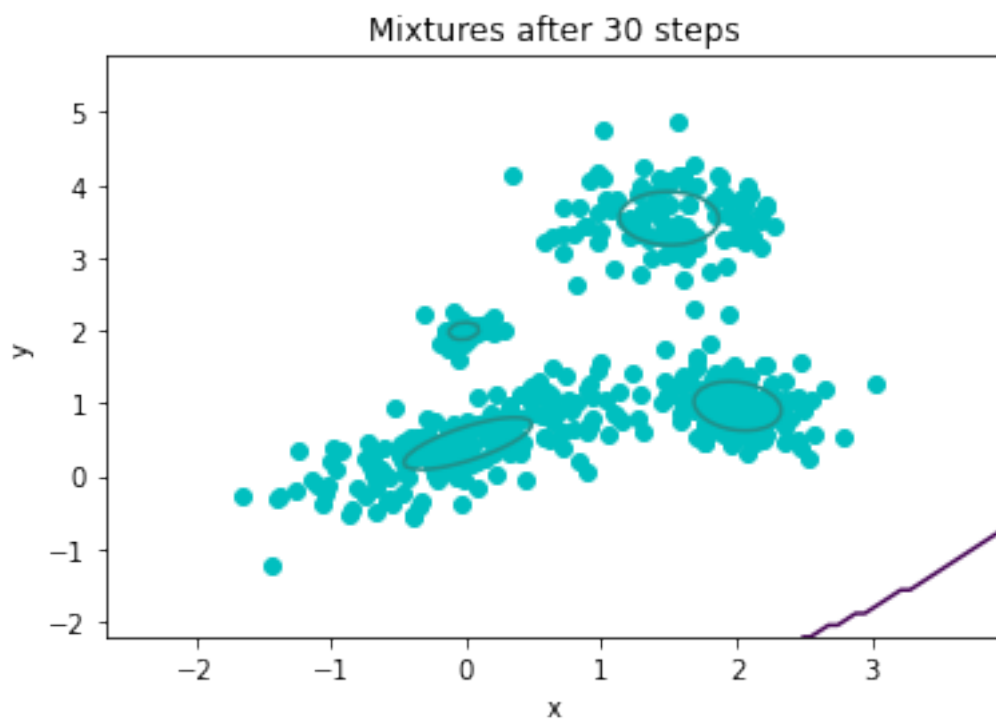


Figure 17

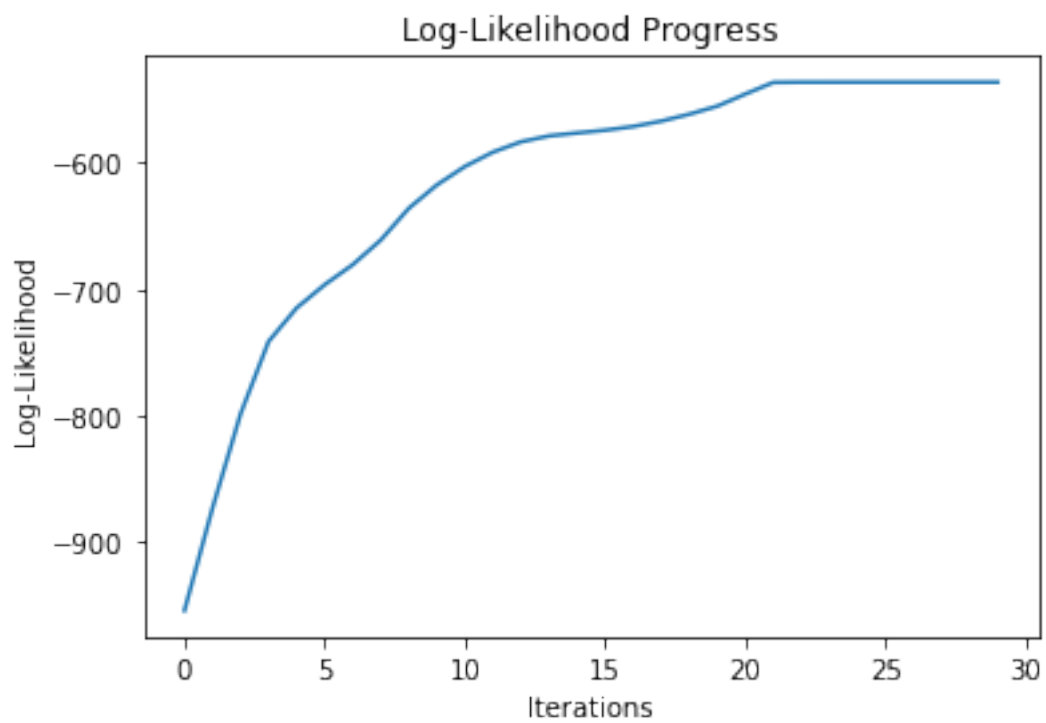


Figure 18

---

## References

---

- [1] Bhavesh Bhatt. *Stochastic vs Batch vs Mini-Batch Gradient Descent*. URL: <https://www.youtube.com/watch?v=Ne3hjpP7KSI>.
- [2] G-kdom. *SGD with Momentum, SGD with Nesterov Momentum*. URL: <https://zhuanlan.zhihu.com/p/73264637>.
- [3] dProgrammer lopez. *Mini-batch, Stochastic Batch Gradient Descent | Tutorial*. URL: <https://www.youtube.com/watch?v=EAZs9qnR6mM>.