# ReadMe

1.) Read the whole file.

2.) Total 3 types of distributions were used to model data distribution within the dataset. 1. Histograms, 2. Uni-axial Gaussian Model, and 3. Uni-axial Gaussian Mixture models.

3.) File 'Results' has all parameters values with distribution graphs.

4.) Datasets used. 1. Iris Data set, 2. Indian Liver Patient Dataset

5.) Note: My major is civil engineering and this is the first time I used python. So files are hard coded specifically, categorical distribution files. Gaussian model and Gaussian mixture model files can be used for other datasets by changing the input variables as explained below.

**Categorical Files: File name starts with 'Cat_'.**

Step 1: Import Libraries and read excel file. Save column data in variables.

```
1    import numpy as np
2    import pandas as pd
3    import matplotlib.pyplot as plt
4
5    data = pd.read_excel('Indian Liver Patient Dataset.xlsx')
6    data = data.to_numpy()
7    classdata = data[:,10]
```

Step 2: Create variables corresponding to data type. Loop through all data and count the no. of occurrence of that data type.

```
# Create variables corresponding to class type in data
a = 0 # for Male
b = 0 # for Female

for x in classdata:
    if x == 'Male':
        a = a+1
    else:
        b = b+1
```

Step 3:Sum the total no. of data/events occurrences. Find probability by dividing data type to total no. of data/event occurrence.

```
# Total weight to find probability/Normalization
sum = a + b

# Probaility of getting Male,Female
pa = a/sum #Male
pb = b/sum #Female
```

Step 4: Print the probability of data type occurrence and plot the data to visualize distributions.

```python
# Print the probility of data type
print('Probaility of Male =', pa)
print('Probaility of Female =', pb)

# Plotting categorical Data to visualize
xaxis = ['Male','Female']
yaxis = [pa,pb]
plt.title('Categorical Distribution of Data')
plt.xlabel('Gender Type')
plt.ylabel('Probability')
plt.bar(xaxis,yaxis)
plt.show()
```

**Gaussian Files: File name 'Gaussian.py'**

Step 1: Import Libraries and read excel file. Save column data in variable. To run this file enter excel file name containing data in red rectangle box. Type the no. of column in green box to find Gaussian distribution of that data.

```python
1    import numpy as np
2    import pandas as pd
3    import matplotlib.pyplot as plt
4
5    data = pd.read_excel('Indian Liver Patient Dataset.xlsx')
6    data = data.to_numpy()
7    sepal_length = data[:,9]
```

Step 2: Create bins of the data to display the frequency, find out mean and standard deviation of the data and print these parameters.

```python
10    # bins and histogram
11    count, bins, ignored = plt.hist(sepal_length,15, density=True)
12    print(count)
13    print(bins)
14
15    #Calculate mean and std
16    mean = np.mean(sepal_length)
17    sd = np.std(sepal_length)
18    print('Mean of this data =',mean)
19    print('Standard Deviation of this data =',sd)
```

Step 3: Create function to find out the Gaussian probability density of the data point and apply function to data.

```python
19   # Creating a Function
20   def normal_dist(bins, mean, sd):
21       #prob_density = (np.pi*sd) * np.exp(-0.5*((data_age-mean)/sd)**2)
22       prob_density = 1/(sd*np.sqrt(2*np.pi)) * np.exp(-(bins.astype(float) - mean)**2/(2*sd**2))
23       return prob_density
24
25   #Apply function to data.
26   pdf = normal_dist(bins,mean,sd)
```

Step 4: Plot the results.

```python
28   #Plotting the results
29   plt.plot(bins,pdf, color = 'red')
30   plt.xlabel(xlabel)
31   plt.ylabel('Probability Density')
32   plt.title(Title)
33   plt.show()
```

**Gaussian Mixture Files: File name 'GaussianMixture.py' (Code gives different parameters each time)**

Step 1: Import Libraries and read excel file. Save column data in variable. To run this file enter excel file name containing data in red rectangle box. Type the no. of column in green box to find Gaussian distribution of that data.

```python
1    import numpy as np
2    import pandas as pd
3    import math
4    import matplotlib.pyplot as plt
5
6    # Read the data
7    data = pd.read_excel('iris.xlsx')
8    data = data.to_numpy()
9    sepal_length = data[:,2]
```

Step 2: Input no. of classes, and iteration no.

```python
11   # Input: no. of classes, no. of iterations
12   k = 2
13   iter = 150
```

Step 3: Divide the data into bins

```python
15   # Calculate Mean and Standard Deviation
16   count, bins, ignored = plt.hist(sepal_length, 20, density=True)
17
```

# Step 4: Initialize some variables in E-step and M-step

```python
18    # Create Pi,mu, and sigma
19    mu = np.random.randint(5, size=(k,1))
20    sigma = np.random.randint(5, size=(k,1))
21    pi = np.random.dirichlet(np.ones(k),size=1).transpose()
22
23    Gamma = np.zeros((len(bins),k))
24    pdf = np.zeros((len(bins),k))
25    xGamma = np.zeros((len(bins),k))
26    xsigma = np.zeros((len(bins),k))
27
```

$$mu = \begin{bmatrix} \ \end{bmatrix}_{k \times 1} \quad (mean)$$

K = no. of classes
1 → one Variable

$$sigma = \begin{bmatrix} \ \end{bmatrix}_{k \times 1} \quad (variance)$$

K = no. of classes
1 → one Variable

$$pi = \begin{bmatrix} \ \end{bmatrix}_{k \times 1} \quad (weightage)$$

K = no. of classes
1 → one Variable

2. **E step**. Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}. \qquad (9.23)$$

$$Gamma = \begin{bmatrix} \ \end{bmatrix}_{n \times k}$$

n → no of data points
k → no of classes

$$PDF = \begin{bmatrix} \ \end{bmatrix}_{n \times k}$$
↓
Probability
Density
Function

n → no of data points
k → no of classes

Step 4: Create function to find out the Gaussian probability density of the data point.

```python
28    # Creating a Function for calculating gaussian distribution probability
29  ∨ def normal_dist(data, mean, sd):
30        prob_density = 1/(sd*np.sqrt(2*np.pi)) * np.exp(-(data.astype(float) - mean)**2/(2*sd**2))
31        return prob_density
32
```

Step 5: E-Step

```python
33      # E STEP CALCULATION
34  ∨ def estep(mu,sigma,pi):
35          # Calculate PDF
36  ∨      for x in range(k):
37              pdf[:,x] = normal_dist(bins,mu[x],sigma[x])
38
39          # Calculate pipdf2
40          pipdf2 = pdf.copy()
41  ∨      for x in range(k):
42              pipdf2[:,x] = pi[x] * pipdf2[:,x]
43          pipdf2 = pipdf2.sum(axis=1)
44
45          # Calculate pipdf
46          pipdf = pdf.copy()
47  ∨      for x in range(k):
48              pipdf[:,x] = pi[x] * pdf[:,x]
49
50          # Calculate Gamma
51  ∨      for x in range(k):
52              Gamma[:,x] = pipdf[:,x] / pipdf2[:]
53
54          return Gamma
```

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}.$$

<u>Step 5: M-Step</u>

```python
56    # M STEP CALCULATION
57    def mstep(resposibility):
58        xGamma = np.zeros((len(bins),k))
59        xsigma = np.zeros((len(bins),k))
60        # Calculate NK
61        NK = Gamma.sum(axis=0)
62
63        # Calculate new mu
64        for x in range(k):
65            xGamma[:,x] = bins * Gamma[:,x]
66        xGamma = xGamma.sum(axis=0)
67        munew = xGamma / NK
68
69        # Calculate Sigma
70        for x in range(k):
71            xsigma[:,x] = Gamma[:,x] * (bins-munew[x])* (bins-munew[x])
72        xsigma = xsigma.sum(axis=0)
73        sigmanew = np.sqrt(xsigma / NK)
74
75        # Calculate pinew
76        pinew = NK / len(bins)
77
78        return munew,sigmanew,pinew
```

3. **M step**. Re-estimate the parameters using the current responsibilities

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk}) \mathbf{x}_n \tag{9.24}$$

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk}) \left(\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}\right) \left(\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}\right)^{\text{T}} \tag{9.25}$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} \tag{9.26}$$

where

$$N_k = \sum_{n=1}^{N} \gamma(z_{nk}). \tag{9.27}$$

Step 6: Iterations. (input value)

```
80    for x in range(iter):
81        responsibility = estep(mu,sigma,pi)
82        mu,sigma,pi = mstep(responsibility)
83
```

Step 7: Calculate final PDFs, print parameters, and plot the distribution. PDF need to be multiplied by weightages.

```
84    pdf1 = normal_dist(bins,mu[0],math.sqrt(sigma[0]))
85    pdf2 = normal_dist(bins,mu[1],math.sqrt(sigma[1]))
86
87    print('mu =', mu)
88    print('sigma =', sigma)
89    print('pi =', pi)
90
91
92    #Plotting the results
93    plt.plot(bins,pi[0]*pdf1, color = 'red')
94    plt.plot(bins,pi[1]*pdf2, color = 'red')
95    plt.xlabel('Age')
96    plt.ylabel('Probability Density')
97    plt.show()
```

## Files on Github

| | |
|---|---|
| 📄 /Iris Data Set/bezdekIris.data | ✕ |
| 📄 /Iris Data Set/Cat_Iris types.py | ✕ |
| 📄 /Iris Data Set/class_type.png | ✕ |
| 📄 /Iris Data Set/Gaussian.py | ✕ |

These files are for
Iris data set

| | |
|---|---|
| 📄 /Indian Liver Patient Dataset/Alkphos.png | ✕ |
| 📄 /Indian Liver Patient Dataset/Alkphos.py | ✕ |

These files are for
Indian Liver Patient
Dataset