



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

AY: 2023-24

| | | | |
|---------------------|--------|---------------------|----------------------|
| Class: | TE | Semester: | VI |
| Course Code: | CSL604 | Course Name: | Machine Learning Lab |

| | |
|---------------------------------|--|
| Name of Student: | Parth Manoj Raut |
| Roll No.: | 44 |
| Experiment No.: | 5 |
| Title of the Experiment: | Implementation of Expectation-Maximization algorithm |
| Date of Performance: | |
| Date of Submission: | |

Evaluation

| Performance Indicator | Max. Marks | Marks Obtained |
|------------------------------------|------------|----------------|
| Performance | 5 | |
| Understanding | 5 | |
| Journal work and timely submission | 10 | |
| Total | 20 | |

| Performance Indicator | Exceed Expectations (EE) | Meet Expectations (ME) | Below Expectations (BE) |
|------------------------------------|--------------------------|------------------------|-------------------------|
| Performance | 4-5 | 2-3 | 1 |
| Understanding | 4-5 | 2-3 | 1 |
| Journal work and timely submission | 8-10 | 5-8 | 1-4 |

Checked by

Name of Faculty : Mr Raunak Joshi

Signature :

Date :



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

Aim: Implementation of Expectation-Maximization Algorithm.

Objective: Able to interpret the how Expectation Maximization Algorithm is applied in unsupervised clustering algorithm to handle the non observable features.

Theory:

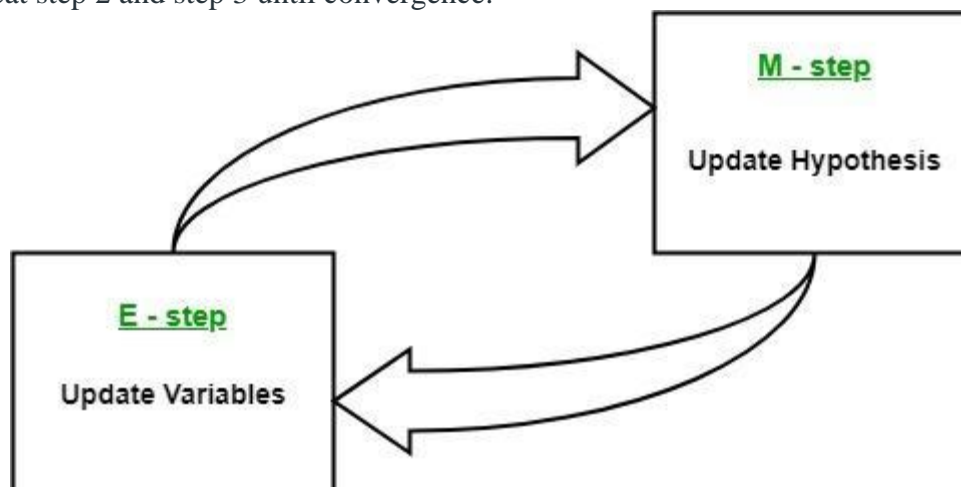
In the real-world applications of machine learning, it is very common that there are many relevant features available for learning but only a small subset of them are observable. So, for the variables which are sometimes observable and sometimes not, then we can use the instances when that variable is visible is observed for the purpose of learning and then predict its value in the instances when it is not observable.

On the other hand, **Expectation-Maximization algorithm** can be used for the latent variables (variables that are not directly observable and are actually inferred from the values of the other observed variables) too in order to predict their values with the condition that the general form of probability distribution governing those latent variables is known to us. This algorithm is actually at the base of many unsupervised clustering algorithms in the field of machine learning.

It was explained, proposed and given its name in a paper published in 1977 by Arthur Dempster, Nan Laird, and Donald Rubin. It is used to find the *local maximum likelihood parameters* of a statistical model in the cases where latent variables are involved and the data is missing or incomplete.

Algorithm:

1. Given a set of incomplete data, consider a set of starting parameters.
2. **Expectation step (E – step):** Using the observed available data of the dataset, estimate (guess) the values of the missing data.
3. **Maximization step (M – step):** Complete data generated after the expectation (E) step is used in order to update the parameters.
4. Repeat step 2 and step 3 until convergence.





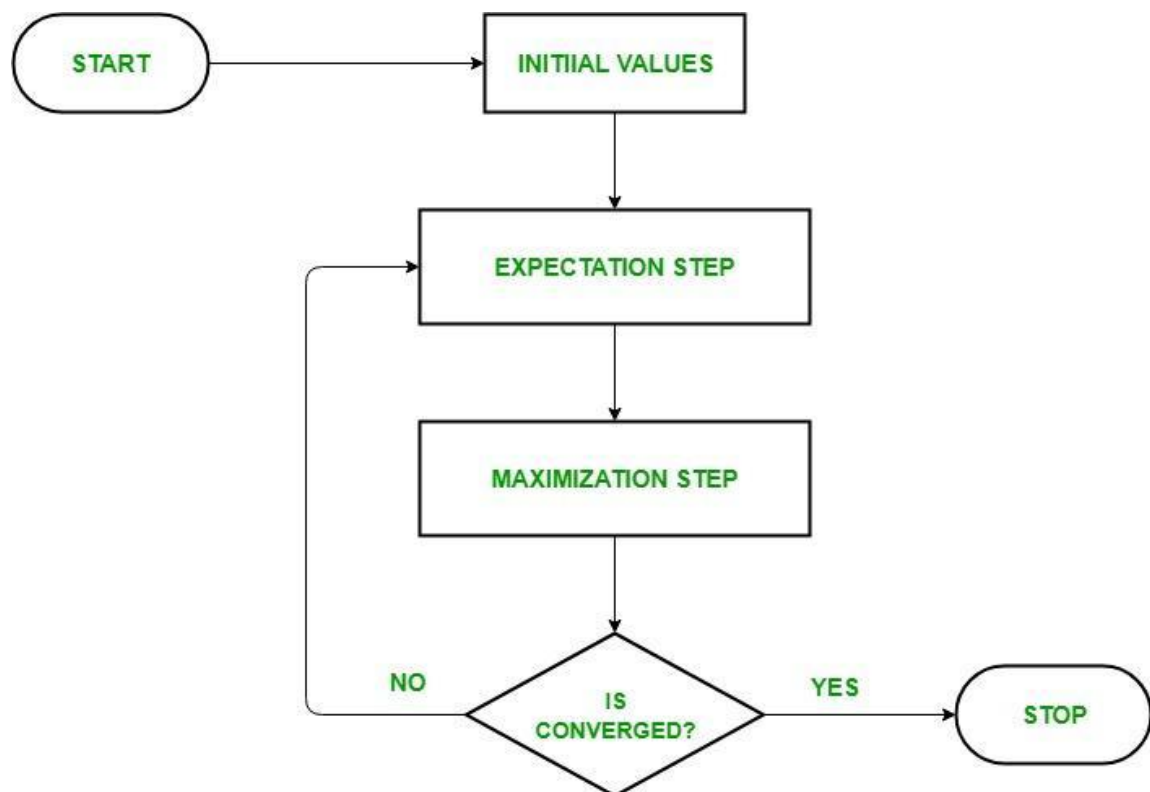
Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

The essence of Expectation-Maximization algorithm is to use the available observed data of the dataset to estimate the missing data and then using that data to update the values of the parameters.

- Initially, a set of initial values of the parameters are considered. A set of incomplete observed data is given to the system with the assumption that the observed data comes from a specific model.
- The next step is known as “Expectation” – step or *E-step*. In this step, we use the observed data in order to estimate or guess the values of the missing or incomplete data. It is basically used to update the variables.
- The next step is known as “Maximization”-step or *M-step*. In this step, we use the complete data generated in the preceding “Expectation” – step in order to update the values of the parameters. It is basically used to update the hypothesis.
- Now, in the fourth step, it is checked whether the values are converging or not, if yes, then stop otherwise repeat *step-2* and *step-3* i.e. “Expectation” – step and “Maximization” – step until the convergence occurs.

Flow chart for EM algorithm –



Usage of EM algorithm –

- It can be used to fill the missing data in a sample.



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

- It can be used as the basis of unsupervised learning of clusters.
- It can be used for the purpose of estimating the parameters of Hidden Markov Model (HMM).
- It can be used for discovering the values of latent variables.

Advantages of EM algorithm –

- It is always guaranteed that likelihood will increase with each iteration.
- The E-step and M-step are often pretty easy for many problems in terms of implementation.
- Solutions to the M-steps often exist in the closed form.

Disadvantages of EM algorithm –

- It has slow convergence.
- It makes convergence to the local optima only.
- It requires both the probabilities, forward and backward (numerical optimization requires only forward probability).

Implementation:

```
import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style("white")
%matplotlib inline
#for matrix math
import numpy as np
#for normalization + probability density function computation
from scipy import stats
#for data preprocessing
import pandas as pd
from math import sqrt, log, exp, pi
from random import uniform
print("import done")

import done

[2]: random_seed=36788765
np.random.seed(random_seed)

Mean1 = 2.0 # Input parameter, mean of first normal probability distribution
Standard_dev1 = 4.0 #@param {type:"number"}
Mean2 = 9.0 # Input parameter, mean of second normal probability distribution
Standard_dev2 = 2.0 #@param {type:"number"}

# generate data
y1 = np.random.normal(Mean1, Standard_dev1, 1000)
y2 = np.random.normal(Mean2, Standard_dev2, 500)

# For data visualisation calculate left and right of the graph
Min_graph = min(data)
Max_graph = max(data)
x = np.linspace(Min_graph, Max_graph, 2000) # to plot the data

print('Input Gaussian {}:  $\mu = \{1.2\}$ ,  $\sigma = \{1.2\}$ '.format("1", Mean1, Standard_dev1))
print('Input Gaussian {}:  $\mu = \{1.2\}$ ,  $\sigma = \{1.2\}$ '.format("2", Mean2, Standard_dev2))
sns.distplot(data, bins=20, kde=False)

Input Gaussian 1:  $\mu = 2.0$ ,  $\sigma = 4.0$ 
Input Gaussian 2:  $\mu = 9.0$ ,  $\sigma = 2.0$ 
<ipython-input-2-8d5c25c80112>:21: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).

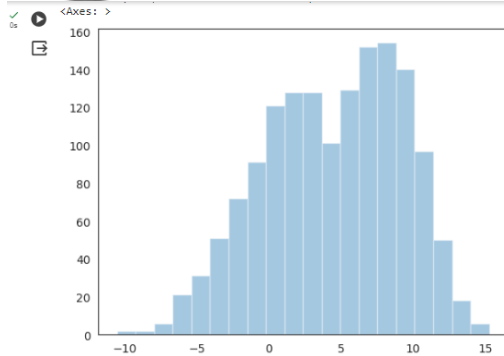
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

sns.distplot(data, bins=20, kde=False)
```



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science



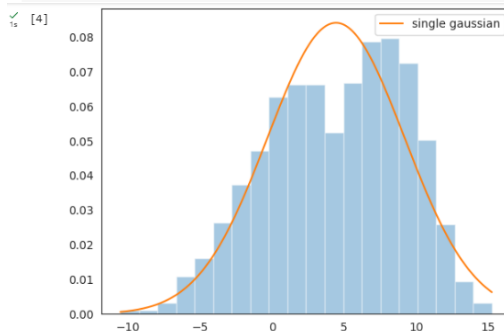
```
[3] class Gaussian:
    """Model univariate Gaussian"""
    def __init__(self, mu, sigma):
        #mean and standard deviation
        self.mu = mu
        self.sigma = sigma

    #probability density function
    def pdf(self, datum):
        """Probability of a data point given the current parameters"""
        u = (datum - self.mu) / abs(self.sigma)
        y = (1 / (sqrt(2 * pi) * abs(self.sigma))) * exp(-u * u / 2)
        return y

    def __repr__(self):
        return 'Gaussian({0:4.6}, {1:4.6})'.format(self.mu, self.sigma)
print("done")

done
```

```
[4] best_single = Gaussian(np.mean(data), np.std(data))
print('Best single Gaussian:  $\mu = {:.2}, \sigma = {:.2}$ '.format(best_single.mu, best_single.sigma))
#fit a single gaussian curve to the data
g_single = stats.norm(best_single.mu, best_single.sigma).pdf(x)
sns.distplot(data, bins=20, kde=False, norm_hist=True)
plt.plot(x, g_single, label='single gaussian')
plt.legend()
```



Expectation Maximization with Gaussian Mixture Model

```
Best single Gaussian:  $\mu = 4.4, \sigma = 4.8$ 
<ipython-input-4-0e8931e2d352>:5: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

sns.distplot(data, bins=20, kde=False, norm_hist=True)
<matplotlib.legend.Legend at 0x7a4dc8920940>
```



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

```
class GaussianMixture_self:
    "Model mixture of two univariate Gaussians and their EM estimation"

    def __init__(self, data, mu_min=min(data), mu_max=max(data), sigma_min=1, sigma_max=1, mix=.5):
        self.data = data
        #todo the Algorithm would be numerical enhanced by normalizing the data first, next do all the EM steps and do the de-normalising at the end

        #init with multiple gaussians
        self.one = Gaussian(uniform(mu_min, mu_max),
                             uniform(sigma_min, sigma_max))
        self.two = Gaussian(uniform(mu_min, mu_max),
                             uniform(sigma_min, sigma_max))

        #as well as how much to mix them
        self.mix = mix

    def Estep(self):
        "Perform an E(steimation)-step, assign each point to gaussian 1 or 2 with a percentage"
        # compute weights
        self.loglike = 0. # = log(p = 1)
        for datum in self.data:
            # unnormalized weights
            wp1 = self.one.pdf(datum) * self.mix
            wp2 = self.two.pdf(datum) * (1. - self.mix)
            # compute denominator
            den = wp1 + wp2
            # normalize
            wp1 /= den
            wp2 /= den # wp1+wp2= 1, it either belongs to gaussian 1 or gaussian 2
            # add into loglike
            self.loglike += log(den) #freshening up self.loglike in the process
            # yield weight tuple
            yield (wp1, wp2)
```

```
    def Mstep(self, weights):
        "Perform an M(aximization)-step"
        # compute denominators
        (left, right) = zip(*weights)
        one_den = sum(left)
        two_den = sum(right)

        # compute new means
        self.one.mu = sum(w * d for (w, d) in zip(left, data)) / one_den
        self.two.mu = sum(w * d for (w, d) in zip(right, data)) / two_den

        # compute new sigmas
        self.one.sigma = sqrt(sum(w * ((d - self.one.mu) ** 2)
                                   for (w, d) in zip(left, data)) / one_den)
        self.two.sigma = sqrt(sum(w * ((d - self.two.mu) ** 2)
                                   for (w, d) in zip(right, data)) / two_den)

        # compute new mix
        self.mix = one_den / len(data)
```

```
    def iterate(self, N=1, verbose=False):
        "Perform N iterations, then compute log-likelihood"
        for i in range(1, N+1):
            self.Mstep(self.Estep()) #The heart of the algorithm, perform E-step and next M-step
            if verbose:
                print('{0:2} {1}'.format(i, self))
            self.Estep() # to freshen up self.loglike

    def pdf(self, x):
        return (self.mix)*self.one.pdf(x) + (1-self.mix)*self.two.pdf(x)

    def __repr__(self):
        return 'GaussianMixture({0}, {1}, mix={2.03})'.format(self.one,
                                                                self.two,
                                                                self.mix)

    def __str__(self):
        return 'Mixture: {0}, {1}, mix={2.03}'.format(self.one,
                                                        self.two,
                                                        self.mix)

print("done")
```

done

```
[6] n_iterations = 20
best_mix = None
best_loglike = float('-inf')
mix = GaussianMixture_self(data)
for _ in range(n_iterations):
    try:
```



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

```
#train!
mix.iterate(verbose=True)
if mix.loglike > best_loglike:
    best_loglike = mix.loglike
    best_mix = mix

except (ZeroDivisionError, ValueError, RuntimeWarning): # Catch division errors from bad starts, and just throw them out...
    print("One Less")
    pass
```

```
1 Mixture: Gaussian(3.19248, 4.4423), Gaussian(8.68074, 2.98154), mix=0.771)
1 Mixture: Gaussian(3.12652, 4.59545), Gaussian(8.65989, 2.56288), mix=0.76)
1 Mixture: Gaussian(3.02684, 4.48971), Gaussian(8.72741, 2.33123), mix=0.75)
1 Mixture: Gaussian(2.92835, 4.44966), Gaussian(8.81862, 2.18843), mix=0.742)
1 Mixture: Gaussian(2.84224, 4.4038), Gaussian(8.89978, 2.0987), mix=0.735)
1 Mixture: Gaussian(2.77093, 4.3601), Gaussian(8.96487, 2.04307), mix=0.729)
1 Mixture: Gaussian(2.71313, 4.32176), Gaussian(9.01402, 2.0096), mix=0.724)
1 Mixture: Gaussian(2.66624, 4.2895), Gaussian(9.04941, 1.99044), mix=0.721)
1 Mixture: Gaussian(2.62767, 4.26288), Gaussian(9.07368, 1.98042), mix=0.717)
1 Mixture: Gaussian(2.59522, 4.24097), Gaussian(9.08939, 1.97624), mix=0.714)
1 Mixture: Gaussian(2.56723, 4.2228), Gaussian(9.0987, 1.97576), mix=0.712)
1 Mixture: Gaussian(2.54248, 4.20748), Gaussian(9.10337, 1.97763), mix=0.709)
1 Mixture: Gaussian(2.5201, 4.1943), Gaussian(9.10474, 1.98008), mix=0.707)
1 Mixture: Gaussian(2.49947, 4.18271), Gaussian(9.1038, 1.98526), mix=0.705)
1 Mixture: Gaussian(2.48015, 4.17227), Gaussian(9.10128, 1.99009), mix=0.703)
1 Mixture: Gaussian(2.46186, 4.1627), Gaussian(9.09768, 1.99524), mix=0.7)
1 Mixture: Gaussian(2.44437, 4.15377), Gaussian(9.09337, 2.00054), mix=0.698)
1 Mixture: Gaussian(2.42754, 4.14533), Gaussian(9.0886, 2.0059), mix=0.696)
1 Mixture: Gaussian(2.41127, 4.13728), Gaussian(9.08354, 2.01124), mix=0.694)
1 Mixture: Gaussian(2.39549, 4.12953), Gaussian(9.07831, 2.01651), mix=0.693)
```

✓ Finding best Mixture Gaussian Model

```
n_iterations = 300
n_random_restarts = 4
best_mix = None
best_loglike = float('-inf')
print('Computing best model with random restarts...\n')
for _ in range(n_random_restarts):
    mix = GaussianMixture_self(data)
    for _ in range(n_iterations):
        try:
            mix.iterate()
            if mix.loglike > best_loglike:
                best_loglike = mix.loglike
                best_mix = mix
        except (ZeroDivisionError, ValueError, RuntimeWarning): # Catch division errors from bad starts, and just throw them out...
            pass
    #print('Best Gaussian Mixture :  $\mu = \{:.2\}$ ,  $\sigma = \{:.2\}$  with  $\mu = \{:.2\}$ ,  $\sigma = \{:.2\}$ '.format(best_mix.one.mu, best_mix.one.sigma, best_mix.two.mu, best_mix.two.sigma))

print('Input Gaussian (1):  $\mu = \{:.2\}$ ,  $\sigma = \{:.2\}$ '.format("1", Mean1, Standard_dev1))
print('Input Gaussian (2):  $\mu = \{:.2\}$ ,  $\sigma = \{:.2\}$ '.format("2", Mean2, Standard_dev2))
print('Gaussian (1):  $\mu = \{:.2\}$ ,  $\sigma = \{:.2\}$ , weight =  $\{:.2\}$ '.format("1", best_mix.one.mu, best_mix.one.sigma, best_mix.mix))
print('Gaussian (2):  $\mu = \{:.2\}$ ,  $\sigma = \{:.2\}$ , weight =  $\{:.2\}$ '.format("2", best_mix.two.mu, best_mix.two.sigma, (1-best_mix.mix)))
#Show mixture
sns.distplot(data, bins=20, kde=False, norm_hist=True)
g_both = [best_mix.pdf(e) for e in x]
plt.plot(x, g_both, label='gaussian mixture')
```

```
g_left = [best_mix.one.pdf(e) * best_mix.mix for e in x]
plt.plot(x, g_left, label='gaussian one')
g_right = [best_mix.two.pdf(e) * (1-best_mix.mix) for e in x]
plt.plot(x, g_right, label='gaussian two')
plt.legend()
```

Computing best model with random restarts...

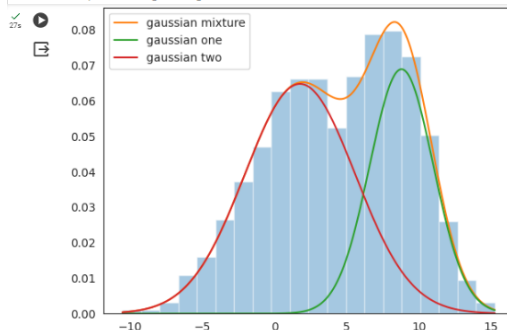
```
Input Gaussian 1:  $\mu = 2.0$ ,  $\sigma = 4.0$ 
Input Gaussian 2:  $\mu = 9.0$ ,  $\sigma = 2.0$ 
Gaussian 1:  $\mu = 8.8$ ,  $\sigma = 2.2$ , weight = 0.38
Gaussian 2:  $\mu = 1.8$ ,  $\sigma = 3.8$ , weight = 0.62
<ipython-input-7-8906b55554ce>:23: UserWarning:
```

'distplot' is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mvaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(data, bins=20, kde=False, norm_hist=True)
<matplotlib.legend.Legend at 0x7a4d8a9f3850>
```





Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

Using Sklearn

```
[8] from sklearn.mixture import GaussianMixture
gmm = GaussianMixture(n_components = 2, tol=0.000001)
gmm.fit(np.expand_dims(data, 1)) # Parameters: array-like, shape (n_samples, n_features), 1 dimension dataset so 1 feature
Gaussian_nr = 1
print('Input Gaussian {}:  $\mu = \{:.2\}$ ,  $\sigma = \{:.2\}$ '.format("1", Mean1, Standard_dev1))
print('Input Gaussian {}:  $\mu = \{:.2\}$ ,  $\sigma = \{:.2\}$ '.format("2", Mean2, Standard_dev2))
for mu, sd, p in zip(gmm.means_.flatten(), np.sqrt(gmm.covariances_.flatten()), gmm.weights_):
    print('Gaussian {}:  $\mu = \{:.2\}$ ,  $\sigma = \{:.2\}$ , weight =  $\{:.2\}$ '.format(Gaussian_nr, mu, sd, p))
    g_s = stats.norm(mu, sd).pdf(x) * p
    plt.plot(x, g_s, label='gaussian sklearn')
    Gaussian_nr += 1
sns.distplot(data, bins=20, kde=False, norm_hist=True)
gmm_sum = np.exp([gmm.score_samples(e.reshape(-1, 1)) for e in x]) #gmm gives log probability, hence the exp() function
plt.plot(x, gmm_sum, label='gaussian mixture')
plt.legend()
```

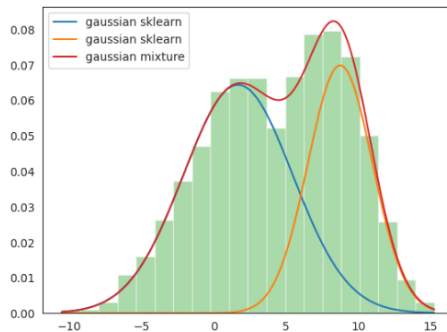
Input Gaussian 1: $\mu = 2.0$, $\sigma = 4.0$
Input Gaussian 2: $\mu = 9.0$, $\sigma = 2.0$
Gaussian 1: $\mu = 1.7$, $\sigma = 3.8$, weight = 0.61
Gaussian 2: $\mu = 8.8$, $\sigma = 2.2$, weight = 0.39
<ipython-input-8-b6a32aee3662>:12: UserWarning:

'distplot' is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mvaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(data, bins=20, kde=False, norm_hist=True)
<matplotlib.legend.Legend at 0x7a4d87b708b0>
```



Conclusion:

1) What is the use of Expectation-Maximization Algorithm?

The Expectation-Maximization (EM) algorithm is used in cases where data contains latent or unobserved variables, such as clustering and density estimation tasks. EM iteratively estimates the parameters of probabilistic models while maximizing the likelihood function. It alternates between two steps: the E-step calculates the anticipated values of latent variables based on current parameter estimations, and the M-step changes the parameters to maximize the expected log-likelihood. EM is very useful for dealing with incomplete data or when direct parameter estimate is impossible owing to lacking information.

2) How it can be used to handle missing data?

The Expectation-Maximization (EM) technique is ideal for addressing missing data since it iteratively estimates missing values while updating model parameters. The E-step calculates the expected values of missing data using the observed data and current parameter estimates. The M-step maximizes the likelihood function by updating the parameters with observed and estimated missing data. This iterative approach continues until convergence, effectively imputing missing values and allowing for robust parameter estimate in the presence of partial data.