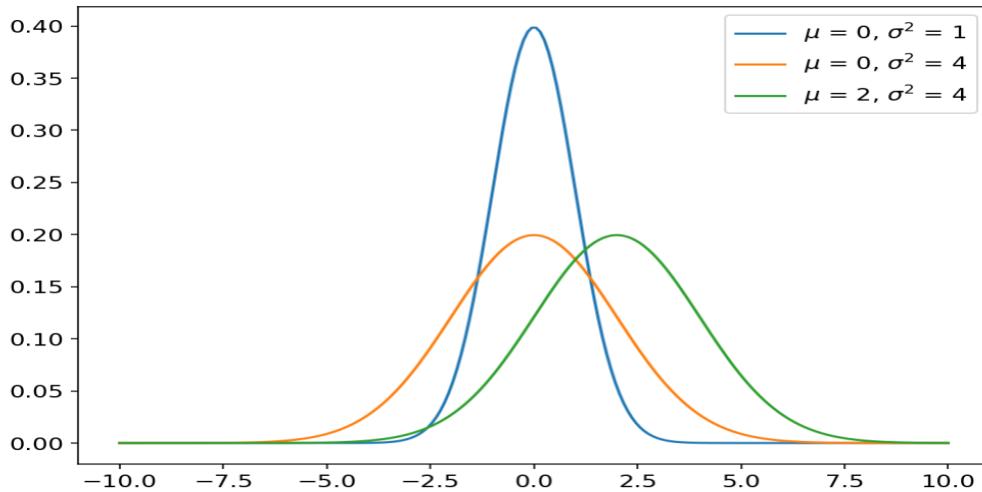


Name: Ashmitha Dale Pais

SID: 923069586

1.1

Suppose lets say $P(X)$ takes a form of a weighted sum of K different gaussian components.

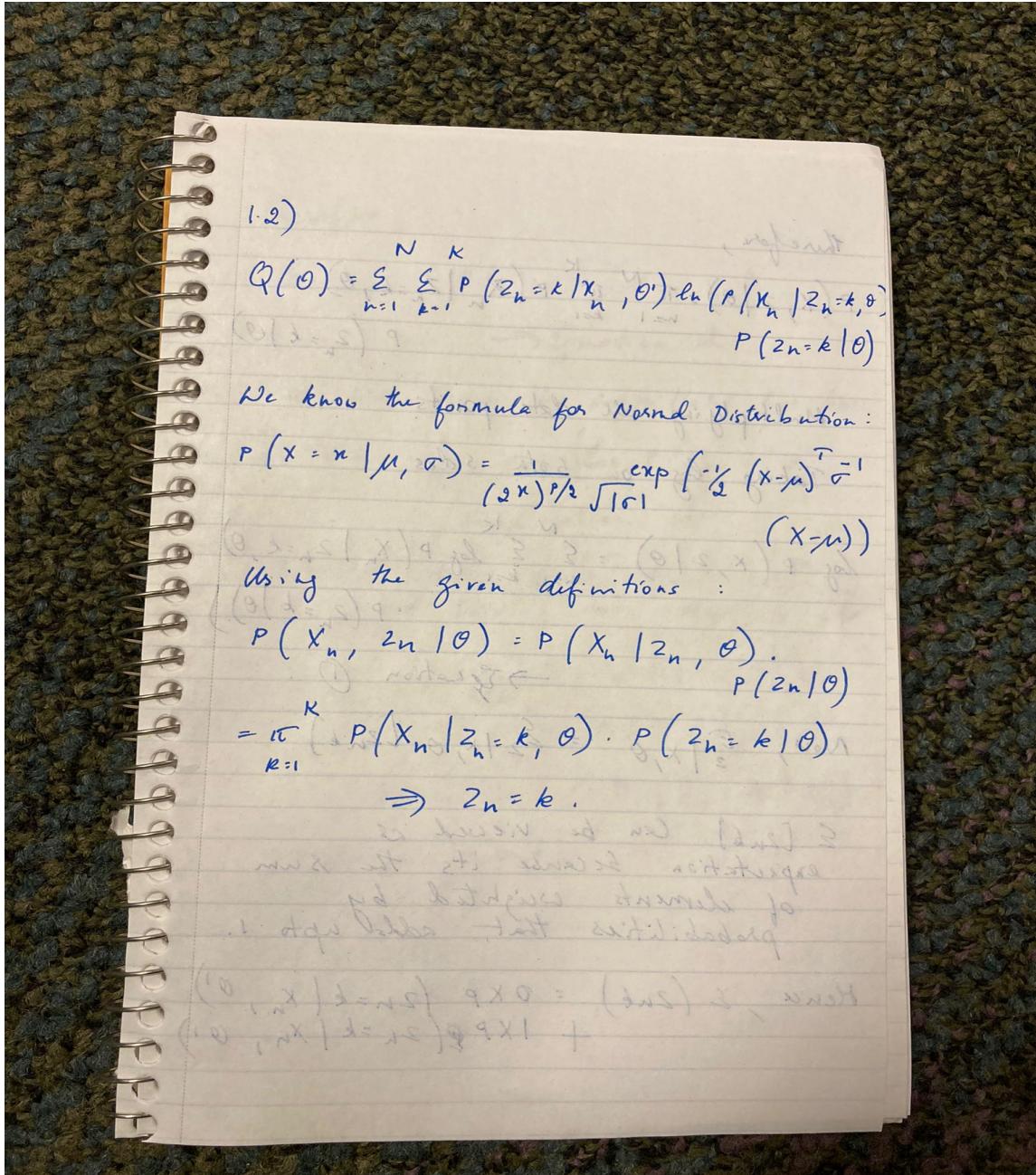


$$P(X_n) = \Pi_1 N(x | \mu_1, \sigma_1^2) + \Pi_2 N(x | \mu_2, \sigma_2^2) + \Pi_3 N(x | \mu_3, \sigma_3^2)$$

$$\text{Then, } P(X_n) = \sum_{k=1}^K \Pi_k N(x_n | \mu_k, \sigma_k^2)$$

$$\text{Where } \Pi_k = \text{mixing weights of } \sum_{k=1}^K \Pi_k = 1$$

1.2



therefore,

$$P(X, Z | \theta) = \prod_{n=1}^N \prod_{k=1}^K P(X_n | Z_n=k, \theta) \cdot P(Z_n=k | \theta)$$

Multiplying all data points

Taking log on both sides

$$\log P(X, Z | \theta) = \sum_{n=1}^N \sum_{k=1}^K \log P(X_n | Z_n=k, \theta) \cdot P(Z_n=k | \theta)$$

\rightarrow Equation ①.

$$\text{Now, } E_Z | X, \theta = E_{Z|X, \theta}[Z_{nk}]$$

$E[Z_{nk}]$ can be viewed as expectation because its the sum of elements weighted by probabilities that adds upto 1.

$$\text{Hence, } E(Z_{nk}) = 0 \times P(Z_n=k | X_n, \theta) + 1 \times P(Z_n=k | X_n, \theta)$$

Therefore,

$$\varepsilon_2 | x, \theta = p(z_n=k | X_n, \theta) \quad)$$

→ Equation 2

Now, we know the auxiliary function

$$Q(\theta) = \varepsilon_2 | x, \theta \cdot [\log p(x, z | \theta)]$$

→ Equation 3

Putting values from equation ①
and ② into equation ③

$$\begin{aligned} \Rightarrow Q(\theta) &= \sum_{n=1}^N \sum_{k=1}^K p(z_n=k | x_n, \theta) \\ &\quad \cdot \ln(p(x_n | z_n=k, \theta) \\ &\quad \cdot p(z_n=k | \theta)) \end{aligned}$$

→ Hence proved.

1.3

1.3)

$$P(z_n=k | x_n, \theta) = \frac{N(x_n | \mu_k, \sigma_k^2) \pi_k}{\sum_{k=1}^K N(x_n | \mu_k, \sigma_k^2) \pi_k}$$

We know that

$$\begin{aligned} P(x_n, z_n | \theta) &= P(x_n | z_n, \theta) \cdot P(z_n | \theta) \\ &= \prod_{k=1}^K (P(x_n | z_{nk}=1, \theta) P(z_{nk}=1 | \theta)) \end{aligned}$$

Here, if $z_{nk}=1$, then that means
 x_n came from the k th Gaussian,
so this is a normal distribution
which is :

$$P(x_n | z_{nk}=1, \theta) = N(x_n | \mu_k, \sigma_x^2)$$

$$\text{And, } P(z_{nk}=1 | \theta) = \pi_k$$

$$\text{So, } P(x_n, z_n | \theta) = \prod_{k=1}^K (N(x_n | \mu_k, \sigma_x^2) \pi_k)^{z_{nk}}$$

Therefore,

$$P(X, Z | \theta) = \prod_{n=1}^N \prod_{k=1}^K P(x_n | \mu_k, \sigma_k^2)^{z_{nk}}$$

$$\log P(X, Z | \theta) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \log [P(x_n | \mu_k, \sigma_k^2)^{z_{nk}}]$$

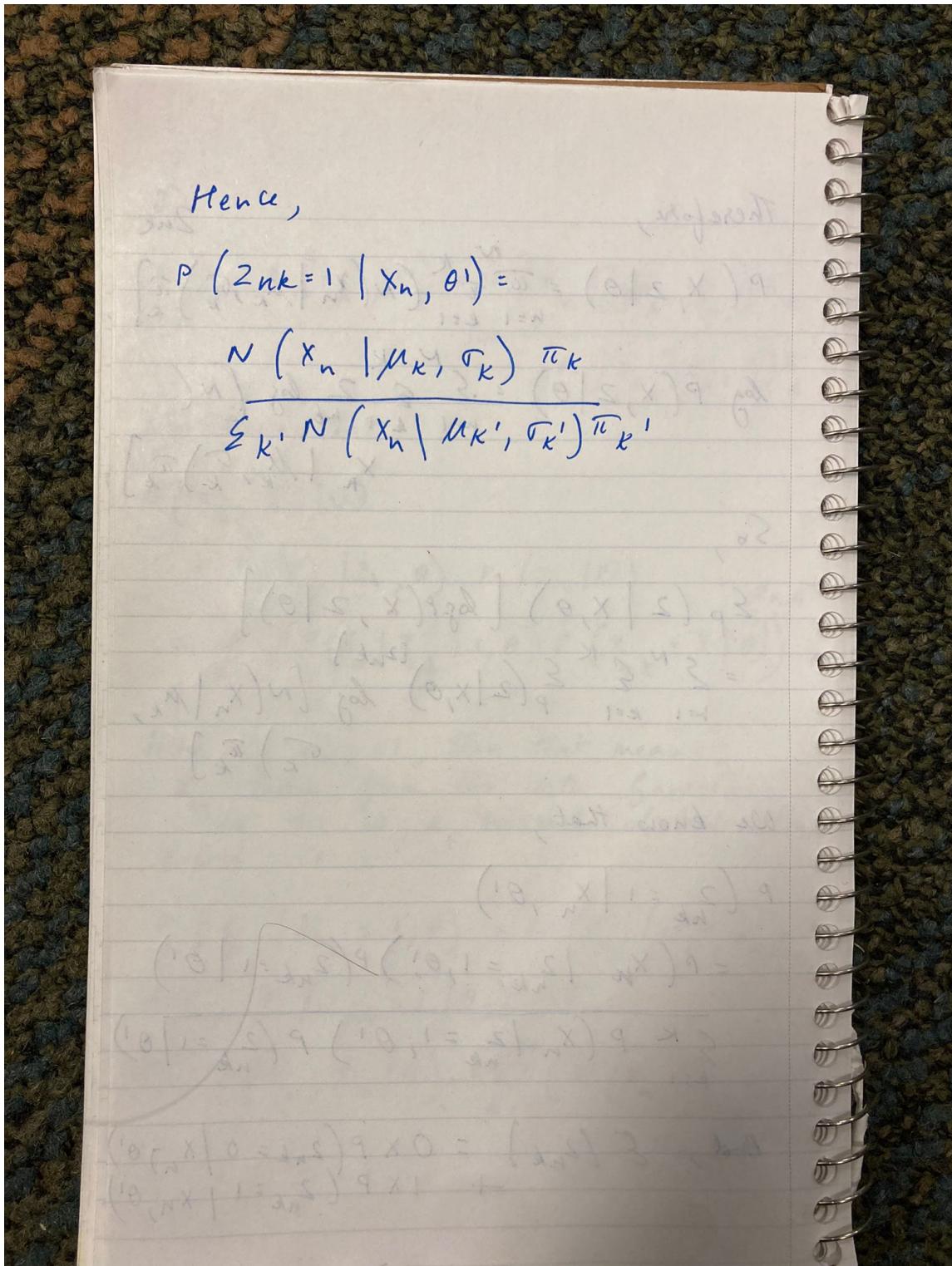
So,

$$\begin{aligned} & E_p(z | X, \theta) [\log P(X, Z | \theta)] \\ &= \sum_{n=1}^N \sum_{k=1}^K \sum_p P(z | X, \theta)^{[z_{nk}]} \log [P(x_n | \mu_k, \sigma_k^2)^{z_{nk}}] \end{aligned}$$

We know that,

$$\begin{aligned} & P(z_{nk}=1 | X_n, \theta') \\ &= P(x_n | z_{nk}=1, \theta') P(z_{nk}=1 | \theta') \\ &\quad \sum_{k=1}^K P(x_n | z_{nk}=1, \theta') P(z_{nk}=1 | \theta') \end{aligned}$$

$$\text{And, } E(z_{nk}) = 0 \times P(z_{nk}=0 | X_n, \theta') + 1 \times P(z_{nk}=1 | X_n, \theta')$$



1.4)

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

$$\text{where } r_{nk} = \frac{N(x_n | \mu_k, \sigma_k) \pi_k}{\sum_{k'} N(x_n | \mu_{k'}, \sigma_{k'}) \pi_{k'}}$$

We know from the previous proofs that,

$$E(z_{nk}) = \frac{N(x_n | \mu_k, \sigma_k) \pi_k}{\sum_{k'} N(x_n | \mu_{k'}, \sigma_{k'}) \pi_{k'}}$$

And,

$$E(\log p(x, z | \theta)) = \sum_{n=1}^N \sum_{k=1}^K E(z_{nk}) [\log N(x_n | \mu_k, \sigma_k) + \log \pi_k]$$

Taking the derivative of μ_k

$$\frac{\partial}{\partial \mu_k} \log N(x_n | \mu_k, \sigma_k)$$

$$= \frac{\partial}{\partial \mu_k} \left(-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right)$$

$$= \frac{\partial}{\partial \mu_k} \left(-\frac{1}{2} (\mu_k^T \Sigma_k^{-1} \mu_k - 2 x_n^T \Sigma_k^{-1} \mu_k) \right)$$

We know that,

$$\frac{\partial}{\partial x} (x^T A x) = (A^T + A)x$$

And, $\frac{\partial}{\partial x} (A^T x) = A$

So, $\frac{1}{2} (2 \Sigma_k^{-1} \mu_k - 2 \Sigma_k^{-1} x_n)$

$$\Rightarrow \Sigma_k^{-1} x_n = \Sigma_k^{-1} \mu_k$$

So, $\frac{\partial}{\partial \mu_k} \Sigma (by p(x, z | \theta))$

$$\Rightarrow \sum_{n=1}^N \Sigma [2_{nk}] \left(\Sigma_k^{-1} x_n - \Sigma_k^{-1} \mu_k \right) = 0$$

$$\Rightarrow \sum_{n=1}^N \Sigma [2_{nk}] x_n = \sum_{n=1}^N \Sigma [2_{nk}] \mu_k$$

Therefore,

$$M_K = \frac{\sum_{n=1}^N r_{nk} X_n}{\sum_{n=1}^N r_{nk}}$$

→ Here, $\sum_{n=1}^N r_{nk}$ is taken to evaluate my answer so I replaced it by the variable r_{nk} in the final answer.

$$\lambda = (x^T A) \frac{b}{r_{nk}}$$

$$\left(\frac{x^T b - \lambda}{N} \right) \frac{1}{r_{nk}}$$

$$\rightarrow \lambda^2 = x^T b$$

$$((0.1x_1 + 0.2x_2) 3 - b) \frac{1}{r_{nk}}$$

$$0 = (4.3 - x_1^2 - x_2^2) [0.1x_1 + 0.2x_2] \frac{1}{r_{nk}}$$

$$x_1 (0.1x_1 + 0.2x_2) \frac{1}{r_{nk}} = \sqrt{b^2 - 4x_1^2 - 4x_2^2} \frac{1}{r_{nk}}$$

2

The process can be divided into estimation and maximization steps:

Estimation Step (E-step):

1. Initially, we set our model parameters: the mean (μ_k), covariance matrix (Σ_k), and mixing coefficients (π_k).
2. For each data point, we compute the posterior probabilities of the data points belonging to each centroid using the current parameter values. These probabilities are often represented by the latent variables y_k .
3. Finally, we estimate the values of the latent variables y_k based on the current parameter values.

Maximization Step:

1. Here, we update the parameter values: the mean (μ_k), covariance matrix (Σ_k), and mixing coefficients (π_k) using the estimated latent variable y_k .
2. We update the mean of the cluster point (μ_k) by computing the weighted average of data points using the corresponding latent variable probabilities.
3. The covariance matrix (Σ_k) is updated by calculating the weighted average of the squared differences between the data points and the mean, using the corresponding latent variable probabilities.
4. The mixing coefficients (π_k) are updated by averaging the latent variable probabilities for each component.

We repeat the E-step and M-step until convergence:

1. We iterate between the estimation step and maximization step until the change in the log-likelihood or the parameters falls below a predefined threshold or until a maximum number of iterations is reached.
2. Essentially, in the estimation step, we update the latent variables based on the current parameter values.
3. Conversely, in the maximization step, we update the parameter values using the estimated latent variables.
4. This iterative process continues until our model converges.