程序代写代談認知編程辅导Clustering



In-class Exercises

K-Medians

Problem 1: Consider a modified version of the K-means objective, where we use L_1 distance instead.

$$WeC_k,\underline{a},\underline{t},\underline{c}$$

This variation of the algorithm is called K-medians. Derive the Lloyd's algorithm for this model. Assignment Project Exam Help 1. Updating the cluster assignments z_{ik} is the same as for the K-means algorithm:

Email:
$$tutorcs@163.com$$

2. The updates for μ_k 's should solve

QQ:
$$749389176_{\mu_k} = \underset{\mu_k}{\text{arg min}} \sum_{i=1}^{476} z_{ik} ||x_i - \mu_k||_1$$

The objective for least single-controld to can be rewritten as

$$egin{aligned} \mathcal{J}(oldsymbol{X},oldsymbol{Z},oldsymbol{\mu}_k) &= \sum_{i=1}^N oldsymbol{z}_{ik}||oldsymbol{x}_i - oldsymbol{\mu}_k||_1 \ &= \sum_{i=1}^N oldsymbol{z}_{ik} \sum_{d=1}^D |oldsymbol{x}_{id} - oldsymbol{\mu}_{kd}| \end{aligned}$$

Clearly, this is a convex function of μ_k , as it is a sum of piecewise linear functions. We can actually solve for each μ_{kd} separately, as they do not interact in the objective, by finding the roots of the derivatives.

Observe, that

$$\frac{\partial}{\partial \boldsymbol{\mu}_{kd}} |\boldsymbol{x}_{id} - \boldsymbol{\mu}_{kd}| = \begin{cases} \frac{\partial}{\partial \boldsymbol{\mu}_{kd}} (\boldsymbol{\mu}_{kd} - \boldsymbol{x}_{id}) = 1 & \text{if } \boldsymbol{\mu}_{kd} > \boldsymbol{x}_{id} \\ \frac{\partial}{\partial \boldsymbol{\mu}_{kd}} (\boldsymbol{x}_{id} - \boldsymbol{\mu}_{kd}) = -1 & \text{if } \boldsymbol{\mu}_{kd} < \boldsymbol{x}_{id} \\ 0 & \text{if } \boldsymbol{\mu}_{kd} = \boldsymbol{x}_{id}. \end{cases}$$

(Note: actually the absolute value function 15 Not differentiable 1 undefined. A rigorous treatment of this problem would require us to use subgradients (see https://web.stanford.edu/class/ee364b/lectures/subgradients_notes.pdf), but just fices for our purpose.) "pretending" that

Hence, the deriv ective is

$$m{z}_{ik} |m{x}_{id} - m{\mu}_{kd}|$$

The first sum represents "number of points x_i assigned to class k, such that $x_{id} < \mu_{kd}$ ". Each of these sums represents the humber of points in day (C. that are located to the left (right) of the given value of μ_{kd} . Because we want to set the gradient to zero, we are looking for such a μ_{kd} , that along the axis d exactly $N_k/2$ points are to left of it, and another $N_k/2$ points are to the right (where A = Single Like). This is exact the definition of a median Help Therefore, the optimal update is given as

Email: tutores@163.com

Gaussian Mixture Model Problem 2: Derive the Stepupdate for the Gaussian mixture model.

In the E-step we have to evaluate the posterior distribution over the latent variables given the current parameters, i.e. $\gamma_t(\mathbf{Z})$ Because GMMs assume that the latent variables are independent, $\gamma_t(\mathbf{Z}) = \prod_{i=1}^N \gamma_t(\mathbf{z}_i)$ and it is enough to derive the E-step for a single data point. The update rule follows directly from Bayes' theorem.

$$\begin{split} \gamma_t(\boldsymbol{z}_i = k) &= \mathrm{p}(\boldsymbol{z}_i = k \mid \boldsymbol{x}_i, \boldsymbol{\pi}^{(t)}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}) \\ &= \frac{\mathrm{p}(\boldsymbol{x}_i \mid \boldsymbol{z}_i = k, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}) \; \mathrm{p}(\boldsymbol{z}_i = k \mid \boldsymbol{\pi}^{(t)})}{\mathrm{p}(\boldsymbol{x}_i \mid \boldsymbol{\pi}^{(t)}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)})} \\ &= \frac{\mathrm{p}(\boldsymbol{x}_i \mid \boldsymbol{z}_i = k, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}) \; \mathrm{p}(\boldsymbol{z}_i = k \mid \boldsymbol{\pi}^{(t)})}{\sum_{j=1}^K \mathrm{p}(\boldsymbol{x}_i \mid \boldsymbol{z}_i = j, \boldsymbol{\pi}^{(t)}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}) \; \mathrm{p}(\boldsymbol{z}_i = j \mid \boldsymbol{\pi}^{(t)})} \\ &= \frac{\boldsymbol{\pi}_k^{(t)} \mathcal{N} \left(\boldsymbol{x}_i \mid \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)}\right)}{\sum_{j=1}^K \boldsymbol{\pi}_j^{(t)} \mathcal{N} \left(\boldsymbol{x}_i \mid \boldsymbol{\mu}_j^{(t)}, \boldsymbol{\Sigma}_j^{(t)}\right)} \end{split}$$

Problem 3: Derive the M-step update for the Gaussian mixture model.

In the M-step we maximize delication of the M-st we plug in the definition of the expected value and expand, we get

$$\mathcal{L} = \sum_{i=1}^{N} \sum_{k=1}^{K} \gamma_{i} \qquad \qquad = k \mid \pi, \mu, \Sigma$$

$$= \sum_{i=1}^{N} \sum_{k=1}^{K} \gamma_{i} \qquad \qquad = k, \mu, \Sigma) \text{ p}(\mathbf{z}_{i} = k \mid \pi)$$

$$= \sum_{i=1}^{N} \sum_{k=1}^{K} \gamma_{i} \qquad \qquad = k, \mu, \Sigma) + \sum_{i=1}^{N} \sum_{k=1}^{K} \gamma_{t}(\mathbf{z}_{i} = k) \log \mathbf{p}(\mathbf{z}_{i} = k \mid \pi)$$

$$\mathcal{L}_{\mathbf{z}}$$

where \mathcal{L}_z only depends on μ and Σ . To find the optimal π , we need to maximize \mathcal{L}_z with respect to π . Since π has several constraints placed on it, we will have to solve the following convex optimization problem.

Assignment Project Exam Help subject to $\sum \pi_k - 1 = 0$

Before we formulate the Empain, we supported the support of the su

$$\begin{aligned} \mathcal{L}_{\boldsymbol{z}} = \sum_{k=1}^{N} \sum_{k=1}^{K} \gamma_{t}(\boldsymbol{z}_{i} = k) \, \log p(\boldsymbol{z}_{i} = k \mid \boldsymbol{\pi}) = \sum_{k=1}^{K} N_{k} \log \pi_{k} \\ \text{where } N_{k} = \sum_{i=1}^{N} \gamma_{t}(\boldsymbol{z}_{i} = k) \text{ is the size of the k-th cluster. The Lagrangian is given by} \end{aligned}$$

https://filtores:
$$(O_{k=1}^K \pi_k)$$

and it has its maximum in π at

$$\frac{\partial f}{\partial \boldsymbol{\pi}_k} = \frac{N_k}{\boldsymbol{\pi}_k} - \lambda \stackrel{!}{=} 0 \Leftrightarrow \boldsymbol{\pi}_k = \frac{N_k}{\lambda}$$

because f is concave as a function of π . This gives us the dual function as

$$g(\lambda) = \max_{\boldsymbol{\pi}} f(\boldsymbol{\pi}, \lambda) = f\left(\left(\frac{N_1}{\lambda}, \dots, \frac{N_K}{\lambda}\right), \lambda\right) = \sum_{k=1}^K N_k \log \frac{N_k}{\lambda} + \lambda - N.$$

When f is concave, the dual is convex and we find the minimum of g at

$$\frac{\partial g}{\partial \lambda} = \sum_{k=1}^K N_k \frac{\lambda}{N_k} \left(-\frac{N_k}{\lambda^2} \right) + 1 = 1 - \frac{N}{\lambda} \stackrel{!}{=} 0 \Leftrightarrow \lambda = N.$$

This means that the M-step for π is $\pi_k^{(t+1)} = \frac{N_k}{N}$.

To find the M-step rule Land Cyc Sed tox Long CS 编程辅导

$$\mathcal{L}_{\boldsymbol{x}} = \sum_{i=1}^{N} \sum_{k=1}^{K} \sum_{i=1}^{N} \sum_{k=1}^{N} \sum_{k=1}^{N} \sum_{k=1}^{N} (\boldsymbol{x}_{i} - \boldsymbol{\mu}_{k}) \sum_{k=1}^{N} (\boldsymbol{x}_{i} - \boldsymbol{\mu}_{k}) + D \log (2\pi) + \log \det \boldsymbol{\Sigma}_{k}).$$

where D is the feature $m{\psi}$ ake the derivative with respect to μ_k

$$\frac{\partial \mathcal{L}_{\boldsymbol{x}}}{\partial \boldsymbol{\mu}_{k}} = -\frac{1}{2} \sum_{i=1}^{N} \gamma_{t}(\boldsymbol{z}_{i} = k) \left((-1) \cdot \left(\boldsymbol{\Sigma}_{k}^{-1} + \boldsymbol{\Sigma}_{k}^{-T} \right) (\boldsymbol{x}_{i} - \boldsymbol{\mu}_{k}) \right) = \sum_{i=1}^{N} \gamma_{t}(\boldsymbol{z}_{i} = k) \left(\boldsymbol{\Sigma}_{k}^{-1} (\boldsymbol{x}_{i} - \boldsymbol{\mu}_{k}) \right)$$

and then find its root WeChat: cstutorcs

$$\frac{\partial \mathcal{L}_{\boldsymbol{x}}}{\partial \boldsymbol{\mu}_{k}} = 0 \Leftrightarrow \sum_{i=1}^{N} \gamma_{t}(\boldsymbol{z}_{i} = k) \boldsymbol{\Sigma}_{k}^{-1} \boldsymbol{x}_{i} = N_{k} \boldsymbol{\Sigma}_{k}^{-1} \boldsymbol{\mu}_{k} \Leftrightarrow \boldsymbol{\mu}_{k} = \frac{1}{N_{k}} \sum_{i=1}^{N} \gamma_{t}(\boldsymbol{z}_{i} = k) \boldsymbol{x}_{i}$$

which gives us the update Signment Project Exam Help

$$\mu_k^{(t+1)} = \frac{1}{N_k} \sum_{i=1}^N \gamma_t(z_i = k) x_i.$$
Email: tutorcs@163.com

It remains to find the M-step for Σ . Again we proceed by taking the derivative with respect to Σ_k

$$\frac{\partial \mathcal{L}_{x}}{\partial \boldsymbol{\Sigma}_{k}} = \mathbf{O} \sum_{i=1}^{N} \gamma_{t} (\boldsymbol{z}_{i} - \boldsymbol{\mu}_{k})^{T} \mathbf{\Sigma}_{k}^{-T} + \boldsymbol{\Sigma}_{k}^{-T} \Big]$$

$$= -\frac{1}{2} \left(N_{k} I_{D} - \sum_{i=1}^{N} \gamma_{t} (\boldsymbol{z}_{i} = k) \left[\boldsymbol{\Sigma}_{k}^{-T} (\boldsymbol{x}_{i} - \boldsymbol{\mu}_{k}) (\boldsymbol{x}_{i} - \boldsymbol{\mu}_{k})^{T} \right] \right) \boldsymbol{\Sigma}_{k}^{-T}$$

$$\mathbf{https:} / \mathbf{tutorcs.com}$$

where I_D is the D-dimensional identity matrix. We finish by finding its root

$$\frac{\partial \mathcal{L}_{\boldsymbol{x}}}{\partial \boldsymbol{\Sigma}_k} = 0 \Leftrightarrow N_k I_D = \boldsymbol{\Sigma}_k^{-T} \sum_{i=1}^N \gamma_t(\boldsymbol{z}_i = k) (\boldsymbol{x}_i - \boldsymbol{\mu}_k) (\boldsymbol{x}_i - \boldsymbol{\mu}_k)^T$$

which produces the final update rule

$$\Sigma_k^{(t+1)} = \frac{1}{N_k} \sum_{i=1}^N \gamma_t(z_i = k) (x_i - \mu_k) (x_i - \mu_k)^T.$$

In this exercise we have used the following matrix calculus rules which you can look up in the matrix cookbook.

$$\frac{\partial \boldsymbol{a}^T \boldsymbol{X} \boldsymbol{a}}{\partial \boldsymbol{a}} = \left(\boldsymbol{X} + \boldsymbol{X}^T \right) \boldsymbol{a}^T \qquad \frac{\partial \boldsymbol{a}^T \boldsymbol{X}^{-1} \boldsymbol{b}}{\partial \boldsymbol{X}} = -\boldsymbol{X}^{-T} \boldsymbol{b} \boldsymbol{a}^T \boldsymbol{X}^{-T} \qquad \frac{\partial \log |\det \boldsymbol{X}|}{\partial \boldsymbol{X}} = \boldsymbol{X}^{-T} \boldsymbol{a}^T \boldsymbol{a}^T$$

Expectation Maximizat程 Apprite 写代故 CS编程辅导 Problem 4: Consider a mixture model where the components are given by independent Bernoulli vari-

ables. This is useful when modelling, e.g., binary images, where each of the D dimensions of the image xplack or white. More formally, we have corresponds to a differen

$$=\prod_{d=1}^D oldsymbol{ heta}_{kd}^{x_d} (1-oldsymbol{ heta}_{kd})^{1-x_d}.$$

ave a product of independent Bernoullis, where $oldsymbol{ heta}_{kd}$ denotes

That is, for a given mixtile the Bernoulli parameter

 \bullet $\theta = \{\theta_{kd} \mid k = 1, \dots, K, d = 1, \dots, D\}$ of a mixture of Derive the EM algorithm Bernoullis.

 \mathbf{l} ixel d.

Assume here for simplicity, that the distribution of components p(z) is uniform: $p(z) = \prod_{k=1}^K \pi_k^{z_k} = \prod_{k=1}^K \left(\frac{1}{K}\right)^{z_k}$. WeChat: cstutorcs

Due to the uniform prior on
$$z_i$$
, the $p(z_i)$ cancel and the responsibilities compute as
$$\underbrace{Assignment}_{\gamma_t(z_i=k)} \underbrace{ProjectExam}_{p(x_i\mid z_i=k,\theta)\cdot p(z_i=l)} = \underbrace{\frac{p(x_i\mid z_i=k,\theta)}{\sum_{l=1}^K p(x_i\mid z_i=l,\theta)}}_{\frac{K}{l=1}} \underbrace{Help}_{p(x_i\mid z_i=l,\theta)}$$

which constitues the Estermail: tutorcs@163.com

It remains to derive the M-step. Similiar to mixture of Gaussians:

$$\sum_{\boldsymbol{z} \sim \gamma_{t}(\boldsymbol{z})}^{\mathbb{E}} [\log p(\boldsymbol{X}, \boldsymbol{z})] = \sum_{i=1}^{N} \sum_{k=1}^{K} 9 \boldsymbol{z} + 9 \boldsymbol{z} +$$

The constant C collects all terms independent of θ and hence irrelevant for further optimization.

We now need to take derivatives with respect to θ .

$$\begin{split} \frac{\partial \mathcal{L}_i}{\partial \boldsymbol{\theta}_{k',d'}} &= \sum_{k=1}^K \gamma_t(\boldsymbol{z}_i = k) \sum_{d=1}^D \left(\boldsymbol{x}_{id} \frac{\partial \log \boldsymbol{\theta}_{kd}}{\partial \boldsymbol{\theta}_{k',d'}} + (1 - \boldsymbol{x}_{id}) \frac{\partial \log (1 - \boldsymbol{\theta}_{kd})}{\partial \boldsymbol{\theta}_{k',d'}} \right) \\ &= \gamma_t(\boldsymbol{z}_i = k) \left(\frac{\boldsymbol{x}_{id}}{\boldsymbol{\theta}_{k',d'}} - \frac{1 - \boldsymbol{x}_{id}}{1 - \boldsymbol{\theta}_{k',d'}} \right) \end{split}$$

We observe that the θ_{kd} do not interact, so their optimal values are independent from each other and we can handle them individually.

$$\frac{\partial \mathbb{E}_{\boldsymbol{z} \sim \gamma_t(\boldsymbol{z})}[\log p(\boldsymbol{X}, \boldsymbol{z} \mid \boldsymbol{\theta})]}{\partial \boldsymbol{\theta}_{kd}} = \sum_{i=1}^{N} \frac{\partial \mathcal{L}_i}{\partial \boldsymbol{\theta}_{kd}} = \sum_{i=1}^{N} \gamma_t(\boldsymbol{z}_i = k) \left(\frac{\boldsymbol{x}_{id}}{\boldsymbol{\theta}_{kd}} - \frac{1 - \boldsymbol{x}_{id}}{1 - \boldsymbol{\theta}_{kd}}\right)$$

By finding the roots $\frac{\partial \mathcal{L}(\mathcal{F}_{p}(t,t_{\theta}))}{\partial \theta_{kd}} = 0$, we obtain the optimal update in a similar fashion as in the standard Bernoulli MLE:



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

https://tutorcs.com

程序代写代做 CS编程辅导

Homework

Gaussian Mixture Mode

Problem 5: Consider a



Derive the expected value

Hint: it is helpful to remember the identity $\text{Cov}[\boldsymbol{x}] = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^T] - \mathbb{E}[\boldsymbol{x}]\mathbb{E}[\boldsymbol{x}]^T$.

For $\mathbb{E}[x]$ we use the law vice ted expetiations tutorcs

$$\begin{array}{l} \mathbb{E}[\boldsymbol{z}] = \mathbb{E}\left[\mathbb{E}[\boldsymbol{z}\mid\boldsymbol{z}]\right] = \sum_{k=1}^{K} \pi_{k} \, \mathbb{E}[\boldsymbol{z}\mid\boldsymbol{z}=k] = \sum_{k=1}^{K} \pi_{k} \mu_{k} \\ Assignment Project Exam Help \end{array}$$

For covariance, we first compute $\mathbb{E}[xx^T]$ again using the law of iterated expectations

ExEmail: tutores@163.com

$$\mathbf{Q} \mathbf{Q}_{k}^{\mathbf{z}=1} \mathbf{\pi}_{k} \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^{T} \mid \boldsymbol{z} = k] \\
\mathbf{Q} \mathbf{Q}_{k}^{\mathbf{z}=1} \mathbf{749389476} \\
= \sum_{k=1}^{K} \boldsymbol{\pi}_{k} \left(\operatorname{Cov}[\boldsymbol{x} \mid \boldsymbol{z} = k] + \mathbb{E}[\boldsymbol{x} \mid \boldsymbol{z} = k] \mathbb{E}[\boldsymbol{x} \mid \boldsymbol{z} = k]^{T} \right)$$

https://etutorcs.com

and thus

$$\operatorname{Cov}[\boldsymbol{x}] = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^T] - \mathbb{E}[\boldsymbol{x}] \, \mathbb{E}[\boldsymbol{x}]^T = \sum_{k=1}^K \boldsymbol{\pi}_k (\boldsymbol{\Sigma}_k + \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T) - \sum_{k=1}^K \sum_{j=1}^K \boldsymbol{\pi}_k \boldsymbol{\pi}_j \boldsymbol{\mu}_k \boldsymbol{\mu}_j^T$$

Problem 6: Consider a mixture of K isotropic Gaussians, all with the same known covariances $\Sigma_k = \sigma^2 I$. Derive the EM algorithm for the case when $\sigma^2 \to 0$, and show that it's equivalent to Lloyd's algorithm for K-means.

We consider a GMM willident call it otropic opvirities. Consider a GMM willident call it otropic opvirities. Consider a GMM will be a considerable and a considerable a considerable

$$\begin{array}{c|c}
\mathbf{z}_{i} \mid \mathbf{z}_{ik} = 1, \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \, \mathbf{p}(\mathbf{z}_{ik} = 1 \mid \boldsymbol{\pi}_{k}) \\
\int \mathbf{p}(\mathbf{x}_{i} \mid \mathbf{z}) \, \mathbf{p}(\mathbf{z}) \, \mathrm{d}\mathbf{z}
\end{array} \tag{1}$$

$$\pi_k \mathcal{N}(\mathbf{x}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\
\stackrel{K}{\underset{l=1}{}} \pi_l \mathcal{N}(\mathbf{x}_i \mid \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)$$
(2)

Tutor CS
$$k \exp\left(\frac{-||x_i - \mu_k||^2}{2\sigma^2}\right)$$

$$l \pi_l \exp\left(\frac{-||x_i - \mu_l||^2}{2\sigma^2}\right)$$
(3)

$$= \frac{1}{\sum_{l} \frac{\pi_{l}}{\pi_{k}} \exp\left(\frac{-||\boldsymbol{x}_{i} - \boldsymbol{\mu}_{l}||^{2} + ||\boldsymbol{x}_{i} - \boldsymbol{\mu}_{k}||^{2}}{2\sigma^{2}}\right)}$$
(4)

If μ_k denotes the center that is closest to x_i , then

Assignment Project Exam Help

for all l, with equality if and only if k = l. For $\sigma \to 0$, the denominator of Equation 4 converges to 1: If k = l, the argument of $\exp(\cdot)$ is exactly zero, while for $k \neq l$ we are exponentiating increasingly negative arguments.

If μ_k denotes a center that is not closest to x_i , there is at least one $l \neq k$ for which

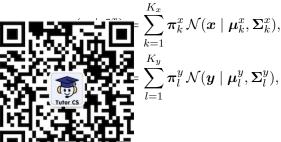
$$000 \cdot 74292589476 \infty \quad \text{as } \sigma \to 0.$$

Consequently, the denominator of Equation 4 diverges to ∞ .

This means that the respectivities descent to a hard one hat assignment of the data point x_i to the component closest to x_i . This coincides with step 1 of Lloyd's algorithm.

Inserting one-hot responsibilities into the general GMM M-step immediately yields step 2 in Lloyd's algorithm. Notice that we do not learn covariances, they are assumed fixed. Moreover, we don't have to worry about π_k s, because they are irrelevant as the term π_l/π_k always gets overshadowed by the $\exp(\cdot)$ next to it.

We can conclude that Lloyd's algorithm for K-Means is a special case of the more general EM algorithm for GMM.



and the random variable

- a) Describe a generat \mathbf{z} f drawing samples) for \mathbf{z} .
- b) Explain in a few sentences why $p(z \mid \theta^x, \theta^y)$ is again a mixture of Gaussians.
- c) State the probability density in the probability of the constant $\mathbf{C}_{\mathbf{S}}^{\mathbf{U}}$ of $\mathbf{C}_{\mathbf{S}}^{\mathbf{U}}$ of $\mathbf{C}_{\mathbf{S}}^{\mathbf{U}}$
 - a) Draw a sample x from $p(x \mid \theta^x)$ with the usual GMM sampling method and the same for y from $p(y \mid \theta^y)$. No significant part of the same for y from $p(y \mid \theta^y)$.
 - b) Let \boldsymbol{x} be drawn from the component k of $p(\boldsymbol{x} \mid \boldsymbol{\theta}^x)$ and \boldsymbol{y} be drawn from the component l of $p(\boldsymbol{y} \mid \boldsymbol{\theta}^y)$. Then \boldsymbol{z} is the sum of two normally distributed random variables $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}_k^x, \boldsymbol{\Sigma}_k^x)$ and $\boldsymbol{y} \sim \mathcal{N}(\boldsymbol{\mu}_l^y, \boldsymbol{\Sigma}_l^y)$. There are $K_x \cdot K_y$ such possible (k, l) combinations, each having probability $\boldsymbol{\pi}_k^x \boldsymbol{\pi}_l^y$ respectively.

That is, $p(\boldsymbol{z} \mid \boldsymbol{\theta}^x, \boldsymbol{\theta}^y)$ is a mixture of $K_x K_y$ Gaussians.

c) It follows from the argument in 4) 9 as the 9 challift density function of z is

$$\begin{aligned} & \text{https://tettercs.com} \\ & \text{https://tettercs.com} \end{aligned}$$

Problem 8: Download the notebook exercise_12_clustering.ipynb from Moodle. Fill in the missing code and run the notebook. Convert the evaluated notebook to PDF and append it to your other solutions before uploading.