

CS772: Assignment 2

Gurpreet Singh — 150259

Question 1

We can write $P(\theta)$ as a convex combination of K priors as follows

$$P(\theta) = \sum_{k=1}^K P(z = k)P(\theta|z = k) \quad (1)$$

Using Chain Rule of Probability, we can write the posterior in the following manner

$$P(\theta|\mathcal{D}) = \sum_{k=1}^K P(\theta|\mathcal{D}, z = k)P(z = k|\mathcal{D}) \quad (2)$$

Now, using Bayes rule, we can write

$$P(\theta|\mathcal{D}, z = k) = \frac{P(\mathcal{D}|\theta, z = k)P(\theta|z = k)}{P(\mathcal{D}|z = k)} \quad (3)$$

Note: The denominator is independent of θ .

$$P(\mathcal{D}|z = k) = \int_{\theta} P(\mathcal{D}|\theta, z = k)P(\theta|z = k) \quad (4)$$

Now, since we are also given that $P(\theta|z = k)$ is conjugate to $P(\mathcal{D}|\theta|z = k)$, hence $P(\theta|\mathcal{D}, z = k)$ is tractable and can be computed. Also

$$\sum_{k=1}^K P(z = k|\mathcal{D}) = 1$$

Therefore, using equation 2 we can say that $P(\theta|\mathcal{D})$ can be written as a convex combination of conjugate distributions. Precisely

$$P(\theta|\mathcal{D}) = \sum_{k=1}^K \frac{P(\mathcal{D}|\theta, z = k)P(\theta|z = k)P(z = k)}{P(\mathcal{D}|z = k)} \quad (5)$$

Question 2

Part A

We know $\mathbf{x} \in \mathbb{R}^D$ and $y \in \mathbb{R}$. This suggests that for the joint distribution of (x, y) , we should use Gaussian Distribution. Therefore

$$(\mathbf{x}, y) \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

This distribution can be justified as we are essentially modelling y as a function of \mathbf{x} alongside modelling \mathbf{x} as a part of the generative model. This can be later seen when we compute the MLE estimate of \mathbf{x} and Σ .

Part B

Using the Chain Rule of Probability, we can say

$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}, y)}{P(\mathbf{x})}$$

Also, we can write

$$P(\mathbf{x}) = \int_{y \in \mathbb{R}} P(\mathbf{x}, y) dy$$

Using properties of Gaussian distributions, we can easily marginalize over y to get $P(\mathbf{x}, y)$

$$P(y|\mathbf{x}) = \mathcal{N}(\mu', \Sigma')$$

where μ' and Σ' are defined as follows

$$\begin{aligned}\mu' &= \boldsymbol{\mu}_y + \Sigma_{yx} \Sigma_{xx}^{-1} (\mathbf{x} - \boldsymbol{\mu}_x) \\ \Sigma' &= \Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}\end{aligned}$$

Note: Here $\Sigma' \in \mathbb{R}$

Part C

For MLE estimate, we need to maximize $P(\mathcal{D})$. Assuming all points to iid, we can write

$$P(\mathcal{D}) = \prod_{n=1}^N p(\mathbf{x}^n, y^n)$$

We can also minimize the NLL of $P(\mathcal{D})$

$$\begin{aligned}NLL &= \sum_{n=1}^N -\log(P(\mathbf{x}^n, y^n)) \\ &= \sum_{n=1}^N \frac{1}{2} (\mathbf{p}^n - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{p}^n - \boldsymbol{\mu})\end{aligned}$$

where

$$\mathbf{p}^n = \begin{bmatrix} \mathbf{x}^n \\ y^n \end{bmatrix}$$

Since this is the NLL of a multivariate gaussian distribution, we can directly write the MLE estimates of μ' and Σ'

$$\begin{aligned} \boldsymbol{\mu} &= \frac{1}{N} \sum_{n=1}^N \mathbf{p}^n \\ \Rightarrow \mu_x &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}^n \quad \mu_y = \frac{1}{N} \sum_{n=1}^N y^n \\ \Sigma &= \frac{1}{N} \sum_{n=1}^N (\mathbf{p}^n - \boldsymbol{\mu})(\mathbf{p}^n - \boldsymbol{\mu})^T \\ \Rightarrow \Sigma &= \frac{1}{N} \sum_{n=1}^N \begin{pmatrix} (\mathbf{x}^n - \boldsymbol{\mu}_x)(\mathbf{x}^n - \boldsymbol{\mu}_x)^T & (\mathbf{x}^n - \boldsymbol{\mu}_x)(y^n - \mu_y) \\ (y^n - \mu_y)(\mathbf{x}^n - \boldsymbol{\mu}_x)^T & (y^n - \mu_y)(y^n - \mu_y) \end{pmatrix} \end{aligned}$$

Question 3

For a one dimensional GMM, $x \in R$ and $z \in \{1, 2 \dots K\}$. Hence, we model $x|z = k$ as a univariate gaussian, and z as a multinoulli.

$$\begin{aligned} x|z = k &\sim \mathcal{N}(\mu_k, \sigma_k^2) \\ z &\sim \text{Multinoulli}(\boldsymbol{\pi}) \end{aligned}$$

Assuming that z is observed, we can write the joint distribution, x, z as follows

$$P(x, z) = \frac{1}{\sqrt{2\pi}\sigma_z} \boldsymbol{\pi}_z e^{-\frac{1}{2\sigma_z^2}(x-\mu_z)^2}$$

This can be written in the exponential form as follows

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2\sigma_z^2}(x-\mu_z)^2 - \log \sigma_z + \log \boldsymbol{\pi}_z}$$

Looking at only the part in the exponent

$$\begin{aligned} \text{exp} &= -\frac{1}{2\sigma_z^2}(x-\mu_z)^2 - \log \sigma_z + \log \boldsymbol{\pi}_z \\ &= -\left[\frac{x^2}{2\sigma_z^2} - \frac{x\mu_z}{\sigma_z^2} + \frac{\mu_z^2}{2\sigma_z^2} + \log \sigma_z - \log \boldsymbol{\pi}_z \right] \end{aligned}$$

If we represent z in the one-hot representation, i.e. $\mathbf{z}_k = \mathbb{I}[z = k]$

$$\begin{aligned} \text{exp} &= -\sum_{k=1}^K \mathbf{z}_k \left(\frac{x^2}{2\sigma_k^2} - \frac{x\mu_k}{\sigma_k^2} + \frac{\mu_k^2}{2\sigma_k^2} + \log \sigma_k - \log \boldsymbol{\pi}_k \right) \\ &= \boldsymbol{\theta}^T \boldsymbol{\phi}(x, z) - A(\boldsymbol{\theta}) \end{aligned}$$

where

$$\boldsymbol{\theta} = \begin{bmatrix} -\frac{1}{2\sigma_k^2} \\ \frac{\mu_k}{\sigma_k^2} \\ \log \boldsymbol{\pi}_k \end{bmatrix}_{k=1}^K \quad \boldsymbol{\phi}(x, z) = \begin{bmatrix} \mathbf{z}_k x^2 \\ \mathbf{z}_k x \\ \mathbf{z}_k \end{bmatrix}_{k=1}^K$$

Since we can write the joint distribution in the exponential form, we can say that x, z is an exponential distribution if z is observed.

Question 4

Since each data point is a binary vector of D dimensions (assuming the features to be independent), and we have K components, we can model $\mathbf{x}|z$ as a Multivariate Bernoulli Distribution, and z as a Multinoulli Distribution.

$$P(\mathbf{x}|z = k) = \prod_{d=1}^D \rho_{k_d}^{x_d} (1 - \rho_{k_d})^{(1-x_d)}$$

$$P(z) = \prod_{k=1}^K \pi_k^{\mathbb{I}[z=k]}$$

Similar to the GMM case, we compute the Complete Log Likelihood ($\log P(\mathbf{X}, z)$) it is not easy to compute the MAP estimate directly. It is easier to represent z in a one hot representation, such that $\mathbf{z}_k = \mathbb{I}[z = k]$

$$P(\mathbf{x}, z) = \prod_{k=1}^K \prod_{d=1}^D \left(\rho_{k_d}^{x_d} (1 - \rho_{k_d})^{(1-x_d)} \pi_k \right)^{z_k}$$

$$P(\mathbf{X}, \mathbf{z}) = \prod_{n=1}^N P(\mathbf{x}^n, \mathbf{z}^n)$$

Therefore, CLL can be written as

$$CLL = \sum_{n=1}^N \sum_{k=1}^K z_k^n \log \left(\rho_{k_d}^{x_d^n} (1 - \rho_{k_d})^{(1-x_d^n)} \right)$$

Here, we are assuming that z or \mathbf{z} is observed, however, that is not the case. It is also hard to simultaneously model both \mathbf{x} and \mathbf{z} . Therefore, we alternately model them, using EM algorithm.

Note: I am using z to represent the integer value of z where as \mathbf{z} represents the one-hot representation

Expectation Step

We first start with an initial value of all the parameters, namely $\{\rho_k\}_{k=1}^K$ and π . Then, we find the expected values of \mathbf{z}_k^n using these parameters. As discussed in class, it is most optimal to calculate the expectation wrt the posterior distribution of the latent variable.

$$\begin{aligned} \mathbb{E}[z_k^n] &= 0 \cdot P(z^n \neq k | \mathbf{x}^n, \{\rho_k\}_{k=1}^K, \pi) + 1 \cdot P(z^n = k | \mathbf{x}^n, \{\rho_k\}_{k=1}^K, \pi) \\ &= P(z^n = k | \mathbf{x}^n, \{\rho_k\}_{k=1}^K, \pi) \\ &= \frac{\pi_k P(\mathbf{x} | z = k, \rho_k)}{\sum_{k'=1}^K \pi_{k'} P(\mathbf{x} | z = k', \rho_{k'})} \end{aligned}$$

Now using these values, we will find out the best MAP estimates for the parameters mentioned above.

Maximization Step

In order to find the MAP estimate, we require the posterior distributions of the parameters.

MAP estimate of π

Since π defines a parameter for a Multinoulli distribution, it is only logical to assume Dirichlet Distribution, i.e $\pi \sim \text{Dir}(\gamma)$.

$$P(\pi) = \frac{1}{B(\gamma)} \prod_{k=1}^K \pi_k^{\gamma_k}$$

$$P(\pi|\mathbf{X}, \mathbf{z}) \propto P(\mathbf{X}, \mathbf{z}|\pi)P(\pi)$$

$$\begin{aligned} \pi_{MAP} &= \arg \max_{\pi} \log P(\pi|\mathbf{X}, \mathbf{z}) \\ \text{s.t. } &\sum_{k=1}^K \pi_k = 1 \end{aligned}$$

This reduces to the following optimization function

$$\pi_{MAP} = \arg \max_{\pi, \lambda} \sum_{n=1}^N \sum_{k=1}^K z_k^n \log \pi_k + \sum_{k=1}^K \pi_k (\gamma_k - 1) + \lambda (1 - \sum_{k=1}^K \pi_k)$$

where λ is the Lagrange Multiplier. We need to replace z_k^n with its expected value, as calculated in the Expectation Step. After solving, we finally get

$$\pi_{MAP} = \left\{ \frac{\sum_{n=1}^N \mathbb{E}[z_k^n] + \gamma_k - 1}{N + \sum_{k'=1}^K \gamma_{k'} - K} \right\}_{k=1}^K$$

MAP estimate of $\{\rho\}_{k=1}^K$

Since we are using ρ_k to model multivariate independent binary data, hence it makes sense to take ρ_{kd} as a Beta distribution.

$$P(\rho_{kd}) = \rho_{kd}^{(\alpha_{kd}-1)} (1 - \rho_{kd})^{(\beta_{kd}-1)}$$

$$P(\rho_{kd}|\mathbf{X}, \mathbf{z}) \propto P(\mathbf{X}, \mathbf{z}|\rho_{kd})P(\rho_{kd})$$

$$\rho_{kd,MAP} = \arg \max_{\rho_{kd}} P(\rho_{kd}|\mathbf{X}, \mathbf{z})$$

The above optimization can be easily solved by differentiating the objective function partially wrt ρ_{kd}

$$\pi_{MAP} = \frac{\sum_{n=1}^N \mathbb{E}[z_k^n] x_d^n + \gamma_k - 1}{\sum_{n=1}^N \mathbb{E}[z_k^n] + \alpha_{kd} + \beta_{kd} - 2}$$

We can alternate between these steps until convergence.

Question 5

We have the following distributions

$$\begin{aligned}\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta &\sim \prod_{n=1}^N \mathcal{N}(y^n|\mathbf{x}^n, \mathbf{w}, \beta^{-1}) \\ \mathbf{w}|\lambda &\sim \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbb{I}_D) \\ \mathbf{w}|\mathbf{y}, \mathbf{X}, \lambda, \beta &\sim \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \Sigma)\end{aligned}$$

where $\boldsymbol{\mu}$ and Σ are defined as follows

$$\Sigma = (\beta \mathbf{X}^T \mathbf{X} + \lambda \mathbb{I}_D)^{-1} \quad \boldsymbol{\mu} = (\mathbf{X}^T \mathbf{X} + \frac{\lambda}{\beta} \mathbb{I}_D)^{-1} \mathbf{X}^T \mathbf{y}$$

Using these distributions, we can find the CLL (Complete Log Likelihood) i.e. $(\log P(\mathbf{y}, \mathbf{w}|\mathbf{X}, \beta, \lambda))$

$$\begin{aligned}P(\mathbf{y}, \mathbf{w}|\mathbf{X}, \beta, \lambda) &= P(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)P(\mathbf{w}|\lambda) \\ \log P(\mathbf{y}, \mathbf{w}|\mathbf{X}, \beta, \lambda) &= \log P(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) + \log P(\mathbf{w}|\lambda) \\ \implies CLL &= \frac{1}{2} \left[N \log \beta + D \log \lambda - \beta(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) \right. \\ &\quad \left. \lambda \mathbf{w}^T \mathbf{w} - (N + D) \log 2\pi \right]\end{aligned}$$

We have computed the CLL assuming we have observed the value of \mathbf{w} , however, that is not the case. Therefore, we will use EM Algorithm to alternatively model the parameters and the hyperparameters.

Expectation Step

We first start with an initial value of all the hyper parameters, namely β and λ . Then, we find the expected values of \mathbf{z}_k^n using these parameters. As discussed in class, it is most optimal to calculate the expectation wrt the posterior distribution of the latent variable. Here, the latent variable is \mathbf{w} .

$$\mathbb{E}[\mathbf{w}] = \int_{\mathbf{w} \in \mathbb{R}^D} \mathbf{w} P(\mathbf{w}|\mathbf{y}, \mathbf{X}, \lambda, \beta) d\mathbf{w}$$

Since the posterior of \mathbf{w} is a normal distribution, we can directly write the expectation of \mathbf{w}

$$= \boldsymbol{\mu}$$

Since we are to maximize the expectation of CLL, we also need to compute $\mathbb{E}[CLL]$

$$\begin{aligned}\mathbb{E}[CLL] &= \mathbb{E} \left[\frac{1}{2} \left[N \log \beta + D \log \lambda - \beta(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) - \lambda \mathbf{w}^T \mathbf{w} - (N + D) \log 2\pi \right] \right] \\ &= \frac{1}{2} \left[N \log \beta + D \log \lambda - \beta(\mathbf{y}^T \mathbf{y} - \mathbb{E}[\mathbf{w}^T] \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbb{E}[\mathbf{w}] + \mathbb{E}[\mathbf{w}^T (\mathbf{X}^T \mathbf{X} + \frac{\lambda}{\beta} \mathbb{I}_D) \mathbf{w}]) \right. \\ &\quad \left. - (N + D) \log 2\pi \right]\end{aligned}$$

We can compute the expectations as follows

$$\begin{aligned}\mathbb{E}[\mathbf{w}^T] &= \mathbb{E}[\mathbf{w}]^T \\ &= \boldsymbol{\mu}\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\mathbf{w}\mathbf{w}^T] &= Cov(\mathbf{w}) + \mathbb{E}[\mathbf{w}]\mathbb{E}[\mathbf{w}]^T \\ &= \Sigma + \boldsymbol{\mu}\boldsymbol{\mu}^T\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\mathbf{w}^T R \mathbf{w}] &= Tr(R\mathbb{E}[\mathbf{w}\mathbf{w}^T]) \\ &= Tr(R(\Sigma + \boldsymbol{\mu}\boldsymbol{\mu}^T))\end{aligned}$$

Hence, we can write the expected value of CLL

$$\mathbb{E}[CLL] = \frac{1}{2} \left[N \log \beta + D \log \lambda - \beta (\mathbf{y}^T \mathbf{y} - \boldsymbol{\mu}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \boldsymbol{\mu} + Tr((\mathbf{X}^T \mathbf{X} + \frac{\lambda}{\beta} \mathbb{I}_D) \boldsymbol{\mu} \boldsymbol{\mu}^T)) - (N + D) \log 2\pi \right]$$

Maximization Step

MLE Estimate of β

We simply require to maximize the Expectation of CLL wrt β

$$\begin{aligned}\beta_{MLE} &= \arg \max_{\beta} \mathbb{E}[CLL] \\ \beta_{MLE}^{-1} &= \frac{1}{N} (\mathbf{y}^T \mathbf{y} - 2\boldsymbol{\mu}^T \mathbf{X}^T \mathbf{y} + Tr((\mathbf{X}^T \mathbf{X}) \boldsymbol{\mu} \boldsymbol{\mu}^T))\end{aligned}$$

MLE Estimate of λ

We simply require to maximize the Expectation of CLL wrt λ

$$\begin{aligned}\lambda_{MLE} &= \arg \max_{\lambda} \mathbb{E}[CLL] \\ \lambda_{MLE}^{-1} &= \frac{Tr(\boldsymbol{\mu} \boldsymbol{\mu}^T)}{D}\end{aligned}$$