

Question 1.2: Computing the probabilities

(4 Points)

Name: Yixi Rao
UID: U6826541

Can you compute probabilities of the class being in program A and program B given the observation, and from there reason about which program is more likely?

Hint: You may need to write out expressions involving the number of students in this class, say n , which is fixed, but not provided in this question. Do the probability of the observed class being in program A depend on class size n ?

You may want to consider the ratio of likelihoods.

let X be a random variable indicating the program I came in. So. $X = \{A, B\}$

let O be a random variable indicating the percentage of boys who are in the program I came in.
So $O = \{0\%, 1\%, \dots, 100\%\}$

$$\begin{aligned} P(X=A | O=55\%) &= \frac{P(O=55\% | X=A) \cdot P(X=A)}{P(O=55\%)} = \frac{P(O=55\% | X=A) \cdot P(X=A)}{\sum_n P(X=n, O=55\%)} = \frac{P(O=55\% | X=A) \cdot P(X=A)}{\sum_n P(X=n | O=55\%) P(X=n)} \\ &= \frac{P(O=55\% | X=A) \cdot P(X=A)}{P(O=55\% | X=A) P(X=A) + P(O=55\% | X=B) P(X=B)} \end{aligned}$$

Suppose that there are two classes in Program A and B

define $x\%$ be the percentage of boys in class 1 in Program A, and $y\%$ for class 2 in Program A

e.g.

$$\text{So we have } 50x + 50y = 2n \cdot 0.65 \Rightarrow x+y = 1.3, \text{ We can conclude that } \begin{cases} 0.3 \leq x \leq 1 \\ 0.3 \leq y \leq 1 \end{cases}$$

so the X or O has 70 possible percentage's values, so $P(O=55\% | X=A) = \frac{1}{70}$

do the same trick to calculate $P(O=55\% | X=B)$. We can get $P(O=55\% | X=B) = \frac{1}{70}$

And $P(X=A) = P(X=B) = \frac{1}{2}$ Since I enter in a classroom in random

$$\text{So } P(X=A | O=55\%) = \frac{1}{70} \quad P(X=B | O=55\%) = \frac{1}{70}$$

$$P(X=A | O=55\%) > P(X=B | O=55\%)$$

Question 1.3: Well-reasoned hunch

(1 Points)

Without computing the posterior probabilities for program A or B, could you make a reasoned argument, based on properties of the binomial distribution, that program B is more likely?

define X_A is a random variable indicating the number of boys sampled from program A in n trials

$$\text{so } P(X_A=k) = C_n^k 0.65^k 0.35^{n-k}$$

define X_B is a random variable indicating the number of boys sampled from program B in n trials

$$\text{so } P(X_B=k) = C_n^k 0.45^k 0.55^{n-k}$$

$$\left\{ \begin{array}{l} D(X_A) = n \cdot 0.65 \cdot 0.35 = 0.225n \\ D(X_B) = n \cdot 0.45 \cdot 0.55 = 0.2475n \end{array} \right. \Rightarrow D(X_B) > D(X_A)$$

The variance of binomial for program B is larger than the variance of binomial for program A, so Program B has more variability in the number of boys, so it is likely that B has as many as 55% boys.

Question 2.1: Verifying valid distribution

(5 Points)

Show that $q(x | \eta)$ given by Eqs. (2.1) and (2.2) is a valid probability density function.

$$q(x | \eta) = \exp(\eta^T \cdot u(x) - \psi(\eta))$$

$$\psi(x) = \log \int \exp(\eta^T \cdot u(x)) dx$$

$$\begin{aligned} q(x | \eta) &= \frac{1}{\exp(\psi(x))} \cdot \exp(\eta^T \cdot u(x)) = \frac{1}{\exp(\log \int \exp(\eta^T \cdot u(x)) dx)} \cdot \exp(\eta^T \cdot u(x)) \\ &= \frac{\exp(\eta^T \cdot u(x))}{\int \exp(\eta^T \cdot u(x)) dx} \end{aligned}$$

To prove $q(x | \eta)$ is a valid probability density function, we need to prove

① $q(x | \eta)$ is nonnegative for each value of the random variable

$$\textcircled{2} \quad \int_x q(x | \eta) dx = 1$$

$$\textcircled{1}: q(x | \eta) = \frac{\exp(\eta^T \cdot u(x))}{\int \exp(\eta^T \cdot u(x)) dx}, \text{ and we know } \forall x, \exp(x) > 0$$

so $q(x | \eta)$ must > 0

$$\textcircled{2}: \int_x q(x | \eta) dx = \int_x \frac{\exp(\eta^T \cdot u(x))}{\int \exp(\eta^T \cdot u(x)) dx} = \frac{\int_x \exp(\eta^T \cdot u(x)) dx}{\int \exp(\eta^T \cdot u(x)) dx} = 1$$

Question 2.2: A Bayesian example (Part 1)

(13 Points)

Suppose that we have a likelihood and prior distribution given by,

$$\begin{aligned} > x | \mu \sim \mathcal{N}(\mu, \sigma^2) & \quad l(x | \eta) = \exp(\eta^T u(x) - \psi(\eta)), \\ > \mu \sim \text{EXP}(u, \eta), & \quad \psi(\eta) = \log \int \exp(\eta^T u(x)) dx \end{aligned}$$

where \mathcal{N} is a 1-dimensional Gaussian distribution with mean μ and standard deviation σ . Show that $\mu | x \sim \text{EXP}(\hat{u}, \hat{\eta})$, for some \hat{u} and $\hat{\eta}$.

Make sure to clearly specify the function \hat{u} and parameters $\hat{\eta}$.

Hint: The integral for $q(\mu)$ does not need to be simplified / solved.

$$\text{want to find } \hat{u}, \hat{\eta}, \text{ to satisfy } P(M | x, \hat{\eta}) = \text{Exp}(\hat{u}, \hat{\eta}) = \frac{\exp(\hat{\eta}^T \cdot \hat{u} | M)}{\int \exp(\hat{\eta}^T \cdot \hat{u} | M) du}$$



$$P(u|x) = \frac{P(x|u)P(u)}{P(x)} = \frac{\overbrace{P(x|u)P(u)}^{\int P(x|u)P(u) du}}{\int \overbrace{P(x|u)P(u)}^{N(u, \sigma^2)} \cdot \frac{\exp(\eta^T u(u))}{\int \exp(\eta^T u(u)) du} du} = \frac{N(u, \sigma^2) \cdot \frac{\exp(\eta^T u(u))}{\int \exp(\eta^T u(u)) du} du}{\int N(u, \sigma^2) \cdot \frac{\exp(\eta^T u(u))}{\int \exp(\eta^T u(u)) du} du} \quad \textcircled{1}$$

find $\hat{\eta}, \hat{u}$ to match $N(u, \sigma^2) \cdot \frac{\exp(\eta^T u(u))}{\int \exp(\eta^T u(u)) du} = \exp(\hat{\eta}^T \hat{u}(u))$

$$\textcircled{1} = \frac{1}{N^{2\pi\sigma^2}} \cdot \exp\left(-\frac{(x-u)^2}{2\sigma^2}\right) \cdot \frac{\exp(\eta^T u(u))}{\int \exp(\eta^T u(u)) du}$$

$$\textcircled{1} = \exp\left(\log(2\pi\sigma^2)^{-\frac{1}{2}}\right) \cdot \exp\left(-\frac{1}{2\sigma^2}(x^2 - 2ux + u^2)\right) \cdot \frac{\exp(\eta^T u(u))}{\int \exp(\eta^T u(u)) du}$$

$$\textcircled{1} = \exp\left(-\frac{1}{2}\log(2\pi\sigma^2)\right) \cdot \exp\left(-\frac{x^2}{2\sigma^2} + \frac{ux}{\sigma^2} - \frac{u^2}{2\sigma^2}\right) \cdot \frac{\exp(\eta^T u(u))}{\int \exp(\eta^T u(u)) du}$$

$$\textcircled{1} = \exp\left(-\frac{1}{2}\log(2\pi\sigma^2) - \frac{x^2}{2\sigma^2} + \frac{ux}{\sigma^2} - \frac{u^2}{2\sigma^2} + \eta^T u(u)\right) \cdot \frac{1}{\int \exp(\eta^T u(u)) du}$$

$$\textcircled{1} = \exp\left(\frac{ux}{\sigma^2} - \frac{u^2}{2\sigma^2} + \eta^T u(u) - \frac{1}{2}\log(2\pi\sigma^2) - \frac{x^2}{2\sigma^2}\right) \cdot \exp\left(-\log\left(\int \exp(\eta^T u(u)) du\right)\right)$$

$$\textcircled{1} = \exp\left(\frac{ux}{\sigma^2} - \frac{u^2}{2\sigma^2} + \eta^T u(u) - \frac{1}{2}\log(2\pi\sigma^2) - \frac{x^2}{2\sigma^2} - \log\left(\int \exp(\eta^T u(u)) du\right)\right)$$

$$\textcircled{1} = \exp\left(\begin{bmatrix} u & u^2 & 1 \end{bmatrix} \begin{bmatrix} \frac{x}{\sigma^2} \\ -\frac{1}{2\sigma^2} \\ \eta^T u(u) - \frac{1}{2}\log(2\pi\sigma^2) - \frac{x^2}{2\sigma^2} - \log\left(\int \exp(\eta^T u(u)) du\right) \end{bmatrix}\right)$$

$$\text{so } \hat{u}(u) = \begin{bmatrix} u \\ u^2 \\ 1 \end{bmatrix}$$

$$\text{and } \hat{\eta} = \begin{bmatrix} \frac{x}{\sigma^2} \\ -\frac{1}{2\sigma^2} \\ \eta^T u(u) - \frac{1}{2}\log(2\pi\sigma^2) - \frac{x^2}{2\sigma^2} - \log\left(\int \exp(\eta^T u(u)) du\right) \end{bmatrix}$$

Question 2.3: A Bayesian example (Part 2)

Suppose that we have the same likelihood as per Question 2.2 but with a Gaussian prior distribution

$$x \mid \mu \sim \mathcal{N}(\mu, \sigma^2)$$

$$\mu \sim \mathcal{N}(\mu_0, \sigma_0^2).$$

First, find the parameters of an exponential family such that $\text{EXP}(u, \eta) = \mathcal{N}(\mu_0, \sigma_0^2)$ (in distribution) with $u = (x, x^2)$.

Then use the result of Question 2.2 to find the parameters of $\mu \mid x$.

$$\begin{aligned} P(\mu \mid \mu_0, \sigma_0^2) &= \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left(-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}\right) = \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left(-\frac{1}{2\sigma_0^2}(\mu^2 - 2\mu_0\mu + \mu_0^2)\right) \\ &= \exp\left(\log(2\pi\sigma_0^2)^{-\frac{1}{2}}\right) \cdot \exp\left(-\frac{1}{2\sigma_0^2}[-2\mu_0, 1] \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} - \frac{\mu_0^2}{2\sigma_0^2}\right) \\ &= \exp\left(-\frac{1}{2} \log(2\pi\sigma_0^2)\right) \cdot \exp\left(-\frac{1}{2\sigma_0^2}[-2\mu_0, 1] \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} - \frac{\mu_0^2}{2\sigma_0^2}\right) \\ &= \exp\left([- \frac{\mu_0}{\sigma_0^2}, -\frac{1}{2\sigma_0^2}] \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} - \frac{\mu_0^2}{2\sigma_0^2} - \frac{1}{2} \log(2\pi\sigma_0^2)\right) \\ &= \exp\left([- \frac{\mu_0}{\sigma_0^2}, -\frac{1}{2\sigma_0^2}] \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} - \left(\frac{\mu_0^2}{2\sigma_0^2} + \frac{1}{2} \log(2\pi\sigma_0^2)\right)\right) \end{aligned}$$

$$\text{so } u = \begin{bmatrix} \frac{\mu_0}{\sigma_0^2} \\ -\frac{1}{2\sigma_0^2} \end{bmatrix} \quad u(u) = \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix}$$

$$\hat{\eta} = \begin{bmatrix} \frac{x}{\sigma^2} \\ -\frac{1}{2\sigma^2} \\ \eta^T u(u) - \frac{1}{2} \log(2\pi\sigma^2) - \frac{x^2}{2\sigma^2} - \log(\int \exp(u^T \cdot u(u)) du) \end{bmatrix}$$

$$= \begin{bmatrix} \frac{x}{\sigma^2} \\ \frac{1}{\sigma^2} \\ -\frac{1}{2\sigma^2} \\ \frac{\mu_0\mu}{\sigma_0^2} - \frac{\mu^2}{2\sigma_0^2} - \frac{1}{2} \log(2\pi\sigma^2) - \frac{x^2}{2\sigma^2} - \log(\int \exp(\frac{\mu_0\mu}{\sigma_0^2} - \frac{\mu^2}{2\sigma_0^2}) du) \end{bmatrix}$$

We notice that $\hat{\eta}$ still has some μ , so we re-factorise it

$$\hat{\eta}^T \cdot u(u) = \begin{bmatrix} \frac{x}{\sigma^2} \\ -\frac{1}{2\sigma^2} \\ \frac{M_0 \cdot M}{\sigma^2} - \frac{M^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) - \frac{x^2}{2\sigma^2} - \log(\int \exp(M_0 u - \frac{M^2}{2\sigma^2}) du) \end{bmatrix} \begin{bmatrix} u^3 \\ u^2 \\ 1 \end{bmatrix}$$

$$= \frac{M \cdot x}{\sigma^2} - \frac{M^2}{2\sigma^2} + \frac{M_0 \cdot M}{\sigma^2} - \frac{M^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) - \frac{x^2}{2\sigma^2} - \log(\int \exp(\eta^T u u) du)$$

from the expression:

$$P(M | M_0, \eta_0) = \exp \left(\left[\frac{M_0}{\sigma^2} - \frac{1}{2\sigma^2} \right] \begin{bmatrix} u \\ u^2 \end{bmatrix} - \left(\frac{M_0^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2) \right) \right)$$

We know $\frac{M_0^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)$ is $\psi(x)$

$$\text{so } \log \int \exp(\eta^T u(u)) du = \frac{M_0^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)$$

$$\hat{\eta}^T u(u) = -\frac{M^2}{2\sigma^2} + \frac{M_0 \cdot M}{\sigma^2} - \frac{M_0^2}{2\sigma^2} - \frac{x^2}{2\sigma^2} + \frac{M \cdot x}{\sigma^2} - \frac{M^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2} \log(2\pi\sigma^2)$$

$$\hat{\eta}^T u(u) = M^2 \left(-\frac{1}{2\sigma^2} - \frac{1}{2\sigma^2} \right) + M \left(\frac{M_0}{\sigma^2} + \frac{x}{\sigma^2} \right) - \frac{M_0^2}{2\sigma^2} - \frac{x^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2} \log(2\pi\sigma^2)$$

$$\text{so } \hat{h}_u = \begin{bmatrix} -\frac{1}{2\sigma^2} - \frac{1}{2\sigma^2} \\ \frac{M_0}{\sigma^2} + \frac{x}{\sigma^2} \\ -\frac{M_0^2}{2\sigma^2} - \frac{x^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2} \log(2\pi\sigma^2) \end{bmatrix}$$

$$u(u) = \begin{bmatrix} u^3 \\ u^2 \\ 1 \end{bmatrix}$$

Question 2.5: KL-Divergence for exponential families

(15 Points)

Show that the KL-divergence of distributions $\text{EXP}(\boldsymbol{u}, \eta_1)$ and $\text{EXP}(\boldsymbol{u}, \eta_2)$ within the same exponential family is given by:

$$D_{\text{KL}}[\eta_1 : \eta_2] = \psi(\eta_2) - \psi(\eta_1) - \lambda_1^\top (\eta_2 - \eta_1).$$

Hint: Use Eq. (2.7).

derive Eq. (2.7).

$$\varphi(\eta) = \log \int \exp(\eta^\top \boldsymbol{u}(x)) dx$$

$$\exp(\varphi(\eta)) = \int \exp(\eta^\top \boldsymbol{u}(x)) dx$$

$$\exp(\varphi(\eta)) \cdot \varphi'(\eta) = \frac{\partial}{\partial \eta} \left(\int \exp(\eta^\top \boldsymbol{u}(x)) dx \right)$$

$$= \int \exp(\eta^\top \boldsymbol{u}(x)) \cdot \boldsymbol{u}(x) dx$$

$$\varphi'(\eta) = \frac{1}{\exp(\varphi(\eta))} \int \exp(\eta^\top \boldsymbol{u}(x)) \cdot \boldsymbol{u}(x) dx$$

$$\varphi'(\eta) = \int \exp(\eta^\top \boldsymbol{u}(x) - \varphi(\eta)) \cdot \boldsymbol{u}(x) dx$$

$$\Rightarrow \Delta \varphi(\eta) = \varphi(\eta) = \int \exp(\eta^\top \boldsymbol{u}(x) - \varphi(\eta)) \cdot \boldsymbol{u}(x) dx = \int q(x|\eta) \cdot \boldsymbol{u}(x) dx = \mathbb{E}_{x \sim \text{Exp}(\boldsymbol{u}, \eta)} [\boldsymbol{u}(x)] = \lambda \quad (2.7.1)$$

now:

$$D_{\text{KL}}[q(x|\eta_1) : q(x|\eta_2)] = \int q(x|\eta_1) \log \frac{q(x|\eta_1)}{q(x|\eta_2)} dx$$

$$= \int \exp(\eta_1^\top \boldsymbol{u}(x) - \varphi(\eta_1)) \cdot \log \frac{\exp(\eta_1^\top \boldsymbol{u}(x))}{\exp(\varphi(\eta_1))} \cdot \frac{\exp(\varphi(\eta_2))}{\exp(\eta_2^\top \boldsymbol{u}(x))} dx$$

$$= \int \exp(\eta_1^\top \boldsymbol{u}(x) - \varphi(\eta_1)) \cdot \log \exp(\eta_1^\top \boldsymbol{u}(x) + \varphi(\eta_2) - \varphi(\eta_1) - \eta_2^\top \boldsymbol{u}(x)) dx$$

$$= \int \exp(\eta_1^\top \boldsymbol{u}(x) - \varphi(\eta_1)) \cdot (\underbrace{\eta_1^\top \boldsymbol{u}(x) + \varphi(\eta_2)}_{\eta_1^\top \boldsymbol{u}(x) + \varphi(\eta_2) - \varphi(\eta_1) - \eta_2^\top \boldsymbol{u}(x)}) dx$$

$$= \int \exp(\eta_1^\top \boldsymbol{u}(x) - \varphi(\eta_1)) \cdot (\eta_1^\top \boldsymbol{u}(x) + \eta_2^\top \boldsymbol{u}(x)) dx + \int \exp(\eta_1^\top \boldsymbol{u}(x) - \varphi(\eta_1)) \cdot \varphi(\eta_2) dx - \dots - \int \varphi(\eta_1) \cdot \exp(\eta_1^\top \boldsymbol{u}(x) - \varphi(\eta_1)) dx$$

according to Eq (2.7.1), and $\int q(x|\eta) dx = 1$. We get

$$= (\eta_1 - \eta_2)^\top \lambda_1 + \varphi(\eta_2) - \varphi(\eta_1)$$

$$= \psi(\eta_2) - \psi(\eta_1) - \lambda_1^\top (\eta_2 - \eta_1). \text{ proved.}$$

Question 2.6: Pythagorean Theorem for exponential families

(10 Points)

Let $\text{EXP}(\boldsymbol{\eta}_1)$, $\text{EXP}(\boldsymbol{\eta}_2)$, and $\text{EXP}(\boldsymbol{\eta}_3)$ be distributions within the same exponential family. Furthermore, let

$$a^2 = D_{\text{KL}}[\boldsymbol{\eta}_1 : \boldsymbol{\eta}_2]; \quad b^2 = D_{\text{KL}}[\boldsymbol{\eta}_2 : \boldsymbol{\eta}_3]; \quad c^2 = D_{\text{KL}}[\boldsymbol{\eta}_1 : \boldsymbol{\eta}_3].$$

Prove that Fig. 1 holds. That is, show that $a^2 + b^2 = c^2$ iff the difference in natural parameters ($\boldsymbol{\eta}_2 - \boldsymbol{\eta}_3$) is perpendicular to the difference in expectation parameters ($\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2$).

We need to prove $a^2 + b^2 = c^2 \Leftrightarrow (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T(\boldsymbol{\eta}_2 - \boldsymbol{\eta}_3) = 0$ in two directions

$$\textcircled{1} \quad a^2 + b^2 = c^2 \Rightarrow (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T(\boldsymbol{\eta}_2 - \boldsymbol{\eta}_3) = 0$$

$$\boldsymbol{a}^2 + \boldsymbol{b}^2 = \boldsymbol{c}^2$$

$$D_{\text{KL}}[\boldsymbol{\eta}_1 : \boldsymbol{\eta}_2] + D_{\text{KL}}[\boldsymbol{\eta}_2 : \boldsymbol{\eta}_3] = D_{\text{KL}}[\boldsymbol{\eta}_1 : \boldsymbol{\eta}_3] \quad (\text{use Eq in Q2.5})$$

$$\psi(\boldsymbol{\eta}_2) - \psi(\boldsymbol{\eta}_1) - \boldsymbol{\lambda}_1^T(\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1) + \psi(\boldsymbol{\eta}_3) - \psi(\boldsymbol{\eta}_2) - \boldsymbol{\lambda}_2^T(\boldsymbol{\eta}_3 - \boldsymbol{\eta}_2) = \psi(\boldsymbol{\eta}_3) - \psi(\boldsymbol{\eta}_1) - \boldsymbol{\lambda}_1^T(\boldsymbol{\eta}_3 - \boldsymbol{\eta}_1)$$

$$\psi(\boldsymbol{\eta}_2) - \boldsymbol{\lambda}_1^T(\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1) - \psi(\boldsymbol{\eta}_2) - \boldsymbol{\lambda}_2^T(\boldsymbol{\eta}_3 - \boldsymbol{\eta}_2) = -\boldsymbol{\lambda}_1^T(\boldsymbol{\eta}_3 - \boldsymbol{\eta}_1)$$

$$\psi(\boldsymbol{\eta}_2) - \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_2 + \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_1 - \psi(\boldsymbol{\eta}_2) - \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_3 + \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_2 = -\boldsymbol{\lambda}_1^T\boldsymbol{\eta}_3 + \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_1$$

$$\psi(\boldsymbol{\eta}_2) - \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_2 - \psi(\boldsymbol{\eta}_2) - \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_3 + \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_2 = -\boldsymbol{\lambda}_1^T\boldsymbol{\eta}_3$$

$$-\boldsymbol{\lambda}_1^T\boldsymbol{\eta}_2 - \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_3 + \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_2 = -\boldsymbol{\lambda}_1^T\boldsymbol{\eta}_3$$

$$-\boldsymbol{\lambda}_1^T\boldsymbol{\eta}_3 + \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_2 + \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_3 - \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_2 = 0$$

$$(\boldsymbol{\lambda}_1^T - \boldsymbol{\lambda}_2^T) \cdot (\boldsymbol{\eta}_2 - \boldsymbol{\eta}_3) = 0 \Rightarrow (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T(\boldsymbol{\eta}_2 - \boldsymbol{\eta}_3) = 0 \Rightarrow \text{proved}$$

$$\textcircled{2} \quad (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T(\boldsymbol{\eta}_2 - \boldsymbol{\eta}_3) = 0 \Rightarrow a^2 + b^2 = c^2$$

$$(\boldsymbol{\lambda}_1^T - \boldsymbol{\lambda}_2^T) \cdot (\boldsymbol{\eta}_2 - \boldsymbol{\eta}_3) = \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_2 - \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_3 - \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_2 + \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_3 = 0$$

$$\text{So } \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_2 - \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_2 - \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_3 = -\boldsymbol{\lambda}_1^T\boldsymbol{\eta}_3$$

$$\left. \begin{aligned} \boldsymbol{a}^2 + \boldsymbol{b}^2 &= \psi(\boldsymbol{\eta}_2) - \psi(\boldsymbol{\eta}_1) - \boldsymbol{\lambda}_1^T(\boldsymbol{\eta}_2 - \boldsymbol{\eta}_1) + \psi(\boldsymbol{\eta}_3) - \psi(\boldsymbol{\eta}_2) - \boldsymbol{\lambda}_2^T(\boldsymbol{\eta}_3 - \boldsymbol{\eta}_2) \\ &= \psi(\boldsymbol{\eta}_2) - \psi(\boldsymbol{\eta}_1) - \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_2 + \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_1 + \psi(\boldsymbol{\eta}_3) - \psi(\boldsymbol{\eta}_2) - \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_3 + \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_2 \end{aligned} \right\}$$

$$\xrightarrow{\text{Replace}} \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_2 - \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_2 - \boldsymbol{\lambda}_2^T\boldsymbol{\eta}_3 + \psi(\boldsymbol{\eta}_2) - \psi(\boldsymbol{\eta}_1) + \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_1 + \psi(\boldsymbol{\eta}_3) - \psi(\boldsymbol{\eta}_2)$$

$$= -\boldsymbol{\lambda}_1^T\boldsymbol{\eta}_3 + \psi(\boldsymbol{\eta}_2) - \psi(\boldsymbol{\eta}_1) + \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_1 + \psi(\boldsymbol{\eta}_3) - \psi(\boldsymbol{\eta}_2)$$

$$= -\boldsymbol{\lambda}_1^T\boldsymbol{\eta}_3 - \psi(\boldsymbol{\eta}_1) + \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_1 + \psi(\boldsymbol{\eta}_3) = -\psi(\boldsymbol{\eta}_1) - \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_3 + \boldsymbol{\lambda}_1^T\boldsymbol{\eta}_1 = D_{\text{KL}}[\boldsymbol{\eta}_1 : \boldsymbol{\eta}_3] = c^2 \text{ proved.}$$

Question 3.1: Deriving the EMM expectation step

Derive the expectation we will maximise in EMM EM. In particular, show that

$$\sum_Z p(Z | X, \theta^{\text{old}}) \log p(X, Z | \theta) = \sum_{n=1}^N \sum_{k=1}^K \gamma^{\text{old}}(z_{nk}) (\log \pi_k + \log q(x_n | \eta_k)) \quad (3.5)$$

where

$$\gamma^{\text{old}}(z_{nk}) = \frac{\pi_k^{\text{old}} q(x_n | \eta_k^{\text{old}})}{\sum_j \pi_j^{\text{old}} q(x_n | \eta_j^{\text{old}})}. \quad (3.6)$$

let $\theta^{\text{old}} = \theta^t$

~~$\sum_Z P(Z | X, \theta^t) \log P(X, Z | \theta)$~~

$$= \sum_Z \log \prod_{n=1}^N P(X_n, Z_n | \theta) \prod_{n=1}^N P(Z_n | X_n, \theta^t)$$

$$= \sum_Z \sum_{n=1}^N \log P(X_n, Z_n | \theta) \cdot \prod_{n=1}^N P(Z_n | X_n, \theta^t) = \sum_Z \sum_{n=1}^N \log P(X_n | Z_n, \theta) \cdot P(Z_n | \theta) \cdot \prod_{n=1}^N P(Z_n | X_n, \theta^t)$$

$$= \sum_Z \sum_{n=1}^N \log \prod_{k=1}^K (q(X_n | \eta_k))^{\gamma_{nk}} \cdot \prod_{k=1}^K \pi_k^{\gamma_{nk}} \cdot P(Z | X, \theta^t)$$

$$= \sum_Z \sum_{n=1}^N \log \prod_{k=1}^K (q(X_n | \eta_k) \cdot \pi_k)^{\gamma_{nk}} \cdot P(Z | X, \theta^t)$$

$$= \sum_Z \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \cdot \log (q(X_n | \eta_k) \cdot \pi_k) \cdot P(Z | X, \theta^t)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \sum_Z \gamma_{nk} \cdot P(Z | X, \theta^t) \cdot \log (q(X_n | \eta_k) \cdot \pi_k)$$

$$= \sum_{n=1}^N \sum_{k=1}^K P(Z_{nk}=1 | X_n, \theta^t) \log (q(X_n | \eta_k) \cdot \pi_k) \quad (\text{using Eq. 3.1})$$

$$= \sum_{n=1}^N \sum_{k=1}^K \frac{P(Z_{nk}=1, X_n | \theta^t)}{P(X_n | \theta^t)} \log (q(X_n | \eta_k) \cdot \pi_k)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \frac{P(Z_{nk}=1, X_n | \theta^t)}{\sum_{k=1}^K P(Z_{nk}=1, X_n | \theta^t)} \log (q(X_n | \eta_k) \cdot \pi_k)$$

Now derive $P(Z_{nk}=1, X_n | \theta^t)$

$$P(X_n, Z_{nk}=1 | \theta^t) = P(Z_{nk}=1 | \theta^t) P(X_n | Z_{nk}=1, \theta^t)$$

$$= P(X_n | Z_{nk}=1, \theta^t) \cdot P(Z_{nk}=1 | \theta^t) \quad (\text{then using Eq. 3.3})$$

$$= q(X_n | \eta_k^t) \cdot \pi_k^t \quad (\text{using Eq. 3.4})$$

$$\text{So } \sum_{n=1}^N \sum_{k=1}^K \frac{P(Z_{nk}=1, X_n | \theta^t)}{\sum_{k=1}^K P(Z_{nk}=1, X_n | \theta^t)} \log (q(X_n | \eta_k^t) \cdot \pi_k^t) = \sum_{n=1}^N \sum_{k=1}^K \frac{q(X_n | \eta_k^t) \cdot \pi_k^t}{\sum_j \pi_j^t q(X_n | \eta_j^t)} (\log \pi_k^t + \log q(X_n | \eta_k^t)) \\ = \sum_{n=1}^N \sum_{k=1}^K \pi_k^t \gamma_{nk} (\log \pi_k^t + \log q(X_n | \eta_k^t))$$

Proved.

Question 3.2: Deriving the EMM maximisation step

(7 Points)

Derive the maximisation updates for EMM EM for mixing parameters π and expectation parameters Λ . That is, assuming that Eq. (3.5) is concave in π and H , show that

$$\pi_k = \frac{N_k}{N} \quad (3.9)$$

$$\text{and} \quad \lambda_k = \frac{1}{N_k} \sum_{n=1}^N (u(x_n) \cdot \gamma^{\text{old}}(z_{nk})), \quad (3.10)$$

where $N_k = \sum_{n=1}^N \gamma^{\text{old}}(z_{nk})$.

Thus show that we have the update equation for the natural parameters H given by

$$\text{let all } x^{\text{old}} = x^t \quad \eta_k = \nabla \varphi \left(\frac{1}{N_k} \sum_{n=1}^N (u(x_n) \cdot \gamma^{\text{old}}(z_{nk})) \right). \quad (3.11)$$

$$(3.9) \quad \pi_k^{t+1} = \arg \max_{\pi_k} \sum_{n=1}^N \sum_{k=1}^K (\log \pi_k + \log q(x_n | \eta_k)) \cdot r^t(z_{nk}), \text{ s.t. } \sum_{k=1}^K \pi_k = 1$$

We observe that this formula: $\sum_{n=1}^N \sum_{k=1}^K (\log \pi_k + \log q(x_n | \eta_k)) \cdot \frac{q(x_n | \eta_k) \cdot \pi_k^t}{\sum_j \pi_j^t q(x_n | \eta_j^t)}$, only $\sum_{n=1}^N \sum_{k=1}^K (\log \pi_k \cdot r^t(z_{nk}))$ contains π_k , so

$$\text{define } \mathcal{L}(\pi_k, \lambda) = \sum_{n=1}^N \sum_{k=1}^K (\log \pi_k \cdot r^t(z_{nk})) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \sum_{n=1}^N \frac{1}{\pi_k} \cdot r^t(z_{nk}) + \lambda$$

$$\text{let } \frac{\partial \mathcal{L}}{\partial \pi_k} = 0 \text{ so}$$

$$\sum_{n=1}^N \frac{1}{\pi_k} \cdot r^t(z_{nk}) + \lambda = 0 \quad (\text{times } \pi_k \text{ on both sides})$$

$$\sum_{n=1}^N r^t(z_{nk}) + \pi_k \lambda = 0 \quad \text{This is true for all } k \in \{1, 2, 3, \dots, K\}, \text{ so we can summarise it as}$$

$$\sum_{k=1}^K \sum_{n=1}^N r^t(z_{nk}) + \pi_k \lambda = 0$$

$$\sum_{k=1}^K \sum_{n=1}^N r^t(z_{nk}) + \sum_{k=1}^K \pi_k \lambda = 0 \quad (\text{We know } \sum_{k=1}^K \pi_k = 1)$$

$$\sum_{n=1}^N \sum_{k=1}^K r^t(z_{nk}) + \lambda = 0 \quad (\text{We know } \sum_{k=1}^K r^t(z_{nk}) = 1)$$

$$\Rightarrow N + \lambda = 0 \Rightarrow \lambda = -N$$

$$\text{This equation can be } \sum_{n=1}^N r^t(z_{nk}) + \pi_k \lambda = 0 \equiv \sum_{n=1}^N r^t(z_{nk}) - N \pi_k = 0 \Rightarrow \pi_k^{t+1} = \frac{1}{N} \sum_{n=1}^N r^t(z_{nk}) = \frac{N_k}{N}. \text{ proved.}$$

(3.10)

$$\eta_k^{t+1} = \arg \max_{\eta_k} \sum_{n=1}^N \sum_{k=1}^K (\log \pi_k + \log q(x_n | \eta_k)) \cdot r^t(z_{nk}) \text{ only } \sum_{n=1}^N \sum_{k=1}^K \log(q(x_n | \eta_k)) \cdot r^t(z_{nk}) \text{ contains } \eta_k$$

$$\text{So define } f(\eta_k) = \sum_{n=1}^N \sum_{k=1}^K \log(q(x_n | \eta_k)) \cdot r^t(z_{nk}) = \sum_{n=1}^N \sum_{k=1}^K (\eta_k^T u(x_n) - \psi(\eta_k)) \cdot r^t(z_{nk})$$

$$\frac{\partial f(\eta_k)}{\partial \eta_k} = \sum_{n=1}^N (u(x_n) - \nabla \psi(\eta_k)) r^t(z_{nk}) = \sum_{n=1}^N (u(x_n) - \lambda_k) r^t(z_{nk})$$

$$\text{let } \frac{\partial f(\eta_k)}{\partial \eta_k} = 0 = \sum_{n=1}^N u(x_n) \cdot r^t(z_{nk}) - \sum_{n=1}^N r^t(z_{nk}) \cdot \lambda_k$$

$$\sum_{n=1}^N u(x_n) r^t(z_{nk}) = \lambda_k \cdot \sum_{n=1}^N r^t(z_{nk}) = \lambda_k N_k$$

$$\therefore \lambda_k^{t+1} = \frac{1}{N_k} \sum_{n=1}^N u(x_n) r^t(z_{nk}). \text{ proved.}$$

(3.11)

We can use the Eq (3.8) without proving, so we can use the proved Eq (3.10) to show that

$$\eta_k = \nabla \varphi(\lambda_k) = \nabla \varphi\left(\frac{1}{N_k} \sum_{n=1}^N (u(x_n) r^t(z_{nk}))\right). \text{ proved.}$$