

1)  
RHS

$$\arg \min_{\theta \in \Theta} E_{\hat{p}(n)} [D_{KL}(\hat{p}(y|n) \parallel p_{\theta}(y|n))]$$

$$= \arg \min_{\theta \in \Theta} E_{\hat{p}(n)} [\log \hat{p}(y|n)] - E_{\hat{p}(n)} [\log p_{\theta}(y|n)]$$

Since first term does not depend on  $\theta$ , it is equivalent to

$$\arg \max_{\theta \in \Theta} E_{\hat{p}(n)} [\log p_{\theta}(y|n)]$$

LHS

$$\arg \max_{\theta \in \Theta} E_{\hat{p}(n,y)} [\log p_{\theta}(y|n)]$$

$$= \arg \max_{\theta \in \Theta} \sum_n \hat{p}(n,y) [\log p_{\theta}(y|n)]$$

$$= \arg \max_{\theta \in \Theta} \sum_n \hat{p}(n) \hat{p}(y|n) [\log p_{\theta}(y|n)]$$

Since  $\hat{p}(y|n)$  is not dependent on  $\theta$ , it can be taken out, So, the above expression is equivalent to

$$= \arg \max_{\theta \in \Theta} \sum_n \hat{p}(n) [\log p_{\theta}(y|n)]$$

$$= \arg \max_{\theta \in \Theta} E_{\hat{p}(n)} [\log p_{\theta}(y|n)]$$

which is what our RHS is equivalent to. Hence, both are equivalent problems.

2) LHS

$$P_0(y|x) = \frac{P_0(x,y)}{P_0(x)}$$

$$= \frac{\pi_y \exp\left(-\frac{1}{2\sigma^2} (x-\mu_y)^T (x-\mu_y)\right) \cdot Z^{-1}(\sigma)}{\sum_i^K \pi_i \exp\left(-\frac{1}{2\sigma^2} (x-\mu_i)^T (x-\mu_i)\right) \cdot Z^{-1}(\sigma)}$$

Since,  $Z^{-1}(\sigma)$  in the denominator is not dependent of  $i$ , it can be taken out of summation and cancel with  $Z^{-1}(\sigma)$  in numerator.

$$= \frac{\pi_y \exp\left(-\frac{1}{2\sigma^2} (x-\mu_y)^T (x-\mu_y)\right)}{\sum_i^K \pi_i \exp\left(-\frac{1}{2\sigma^2} (x-\mu_i)^T (x-\mu_i)\right)} \quad \text{--- (1)}$$

Simplifying the numerator of (1)

$$\begin{aligned} & -\frac{1}{2\sigma^2} (x-\mu_y)^T (x-\mu_y) \\ &= -\frac{1}{2\sigma^2} (x^T - \mu_y^T) (x - \mu_y) \end{aligned}$$

$$= -\frac{1}{2\sigma^2} \left( x^T (x - \mu_y) + (\mu_y^T x - \mu_y^T \mu_y) \right)$$

$$= x^T \frac{(x - \mu_y)}{-2\sigma^2} + \left( \mu_y^T x - \mu_y^T \mu_y \right)$$

In this  $\frac{(x - \mu_y)}{-2\sigma^2} \in \mathbb{R}^n$

$$\left( \frac{\mu_y^T x - \mu_y^T \mu_y}{-2\sigma^2} \right) \in \mathbb{R}_{++}$$

$$= x^T g^y + t^y$$

where  $g^y = \frac{(x - \mu_y)}{-2\sigma^2}$

$$t^y = \frac{\mu_y^T x - \mu_y^T \mu_y}{-2\sigma^2}$$

So (1) can be written as

$$= \frac{\pi_y \exp(n^T s^y + t^y)}{\sum_i^k \pi_i \exp(n^T s^i + t^i)}$$

$$= \frac{\exp(\log \pi_y) \exp(n^T s^y + t^y)}{\sum_i^k \exp(\log \pi_i) \exp(n^T s^i + t^i)}$$

$$= \frac{\exp(n^T s^y + (t^y + \log \pi_y))}{\sum_{i=1}^k \exp(n^T s^i + (t^i + \log \pi_i))}$$

which is equivalent to RHS.

$$\Pr(y|x) = \frac{\exp(n^T w_y + b_y)}{\sum_i \exp(n^T w_i + b_i)}$$

3(a)

$$P(x_1, x_2, \dots, x_n) = P(x_1) P(x_2|x_1) P(x_3|x_1, x_2) \dots P(x_n|x_1, \dots, x_{n-1})$$

-  $P(x_1)$  requires  $k_1 - 1$  parameters

-  $P(x_2|x_1)$  requires  $k_2 \cdot (k_1 - 1)$  parameters

-  $P(x_3|x_1, x_2)$  requires  $k_3 \cdot k_2 \cdot (k_1 - 1)$

⋮

-  $P(x_n|x_1, \dots, x_{n-1})$  requires  $(k_1 - 1) \cdot k_2 \dots k_{n-1}$  parameters

So, total number of params required are

$$= (k_1 - 1)(1 + k_2 + k_2 \cdot k_3 + \dots + k_2 \dots k_{n-1})$$

(b) If each variable  $X_i$  is independent of others, then

$$P(x_1, x_2, \dots, x_n) = P(x_1) P(x_2) \dots P(x_n)$$

- For  $P(x_1)$ , we require  $k_1 - 1$  parameters
  - For  $P(x_n)$ , we require  $k_n - 1$  parameters
- So, total we require  $\sum_{i=1}^n k_i - 1$  parameters

(c)

$$P(x_1, \dots, x_n) = P(x_1) \cdot P(x_2 | x_1) \cdot P(x_3 | x_1, x_2) \cdot \dots \cdot P(x_{m+1} | x_1, x_2, \dots, x_m) \cdot P(x_{m+2} | x_2, \dots, x_{m+1}) \cdot \dots \cdot P(x_n | x_{n-m}, \dots, x_{n-1})$$

As discussed earlier, for first "m" terms, the number of parameters required are

$$\text{total}_{\text{first-m}} = (k_1 - 1)(1 + k_2 + k_2 k_3 + \dots + k_2 \dots k_{m-1})$$

For each term 'i' in the above equation (There will be "n-m" such terms)

$$P_i = (k_{i-m} - 1)(k_{i-m+1}) \dots (k_{i-1}) \\ = (k_{i-m} - 1) \prod_{j=i-m+1}^{i-1} k_j$$

So, total numbers of params for last n-m terms are

$$\text{total}_{\text{last-i-m}} = \sum_{t=m+1}^n P_m$$



$$\text{So, overall total} = \text{total}_{\text{first-m}} + \text{total}_{\text{last-m}}$$

4) Let's take the case where  $n=2$  i.e. there are only 2 variables  $X_1, X_2$ .

where

$$P_f(X) = N(0, 1)$$

$$P_f(x_2 | x_1) = N(\mu_\theta(x_1), \Sigma_\theta(x_1))$$

$$P_f(x_2) = \int P_\theta(x_2, x_1) dx_1$$

Essentially we need to find all the possible values of  $x_1$  and for each get  $P(x_2 | x_1)$  to find  $P(x_2)$ . That is a mixture of infinite gaussians which is not a gaussian. Whereas,

$$P_r(x_2) = N(0, 1)$$

An intuitive example of this is large language models which are used to generate english sentences. A model trained in reverse order may not be able to used for QA use cases because in that "Q" comes before "A" and "A" is dependent on "Q" which is not modelled in reverse order.

$S(n)$

$$\begin{aligned} E(A) &= E \left[ \frac{1}{K} \sum_{i=1}^K p(x|z^{(i)}) \right] \\ &= \frac{1}{n} \sum_{j=1}^n \left( \frac{1}{K} \sum_{i=1}^K p(x|z^{(i)}) \right) \\ &= \frac{1}{K} \sum_{i=1}^K \left( \frac{1}{n} \sum_{j=1}^n p(x|z^{(i)}) \right) \end{aligned}$$

Let's say "n" is a very large number,  
then

$$\begin{aligned} &= \frac{1}{K} \sum_{i=1}^K \left( \sum p(z) p(x|z^{(i)}) \right) \\ &= \frac{1}{K} \sum_{i=1}^K p(x) \\ &= p(x) \end{aligned}$$

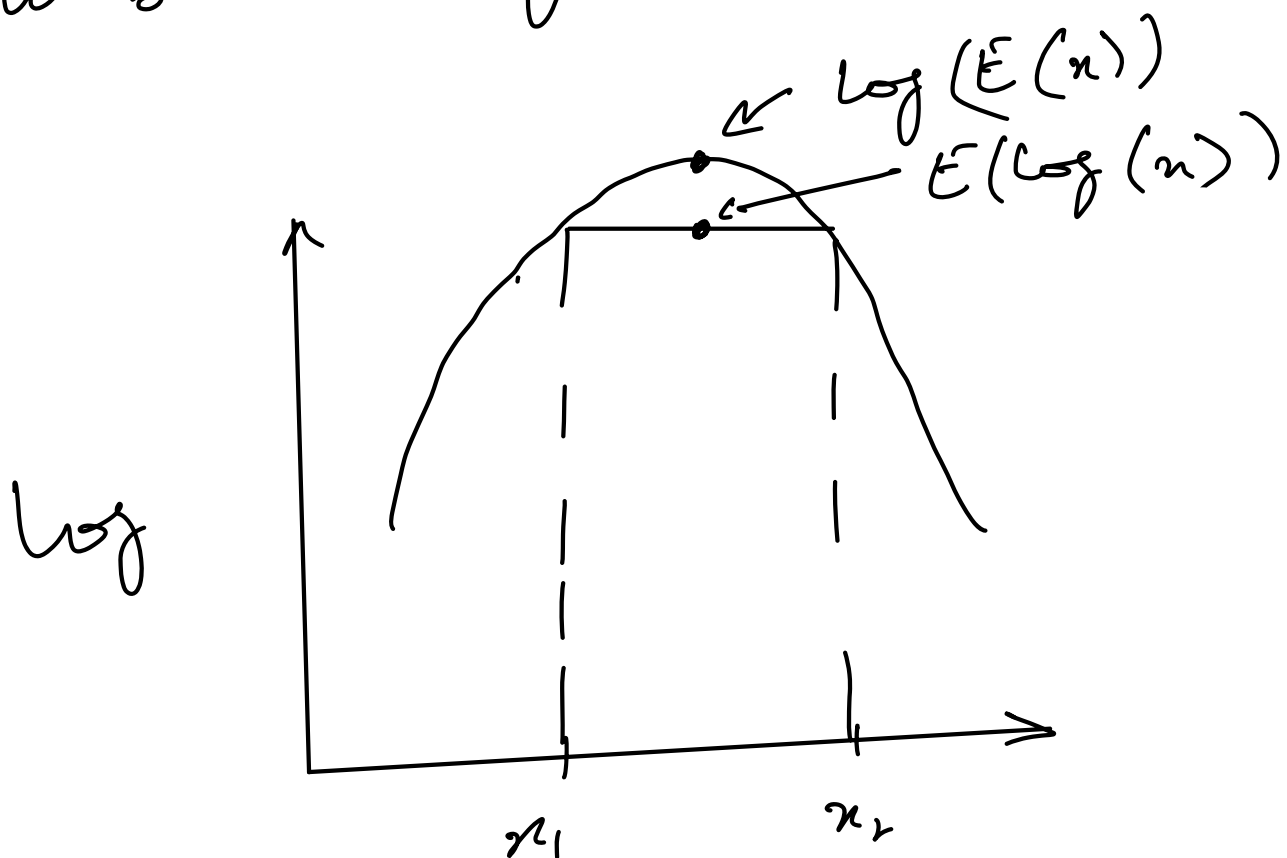
Hence, A is an unbiased  
estimator of  $p(x)$ .

5(b)

$$E(\log A) \leq \log E(A) \quad (\text{Jensen's Eq.})$$

$$\leq \log p(x) \quad \because E(A) = p(x) \text{ as per 5(a)}$$

So, by definition of unbiased estimator,  $\log A$  is not an unbiased estimator. For that, LHS had to be strictly equal to RHS



$$6a) \quad n = 16 \quad \left( \lceil \log_2 50257 \rceil \right) \left( \text{ceil}(\log_2(50257)) \right)$$

6b) Weight matrix of fully connected layer will increase from  $(768 \times 50257)$  to  $(768 \times 60000)$ .

No. of bias terms will increase from 50275 to 60000

So the increase will be

$$768 (60000 - 50257) + (60000 - 50257) = 7,492,367$$