

Q1) To prove,

$$\arg \max_{\theta \in \Theta} E_{\hat{p}(z, y)} [\log p_{\theta}(y|z)] = \arg \min_{\theta \in \Theta} E_{\hat{p}(z)} [D_{KL}(\hat{p}(y|z) || p_{\theta}(y|z))]$$

Proof:

KL divergence for each z :

$$\begin{aligned} D_{KL}(\hat{p}(y|z) || p_{\theta}(y|z)) &= \sum_y \hat{p}(y|z) \log \frac{\hat{p}(y|z)}{p_{\theta}(y|z)} \\ &= \sum_y \hat{p}(y|z) (\log \hat{p}(y|z) - \log p_{\theta}(y|z)) \end{aligned}$$

Taking expectation over z ,

$$\begin{aligned} E_{z \sim \hat{p}(z)} [D_{KL}(\hat{p}(y|z) || p_{\theta}(y|z))] &= \sum_z \hat{p}(z) D_{KL}(\hat{p}(y|z) || p_{\theta}(y|z)) \\ &= \sum_z \hat{p}(z) \sum_y \hat{p}(y|z) [\log \hat{p}(y|z) - \log p_{\theta}(y|z)] \\ &= \sum_{z, y} \hat{p}(z, y) [\log \hat{p}(y|z) - \log p_{\theta}(y|z)] \\ &= E_{\hat{p}(z, y)} [\log \hat{p}(y|z)] - E_{\hat{p}(z, y)} [\log p_{\theta}(y|z)] \end{aligned}$$

$E_{\hat{p}(z, y)} [\log \hat{p}(y|z)]$ is independent of θ & can be treated as a constant, C .

$$\therefore E_{z \sim \hat{p}(z)} [D_{KL}(\hat{p}(y|z) || p_{\theta}(y|z))] = C - E_{\hat{p}(z, y)} [\log p_{\theta}(y|z)]$$

As C is a constant, minimizing the KL divergence over θ is equivalent to maximizing the expected log-likelihood.

$$\arg \min_{\theta} (C - E_{\hat{p}(z,y)} [\log p_{\theta}(y|z)]) = \arg \max_{\theta} E_{\hat{p}(z,y)} [\log p_{\theta}(y|z)]$$

$$\therefore \arg \max_{\theta \in \Theta} E_{\hat{p}(z,y)} [\log p_{\theta}(y|z)] = \arg \min_{\theta \in \Theta} E_{\hat{p}(z)} [D_{KL}(\hat{p}(y|z) || p_{\theta}(y|z))]$$

$$Q2) p_{\theta}(y) = \pi, \sum_{y=1}^k \pi_y = 1, p_{\theta}(z|y) = \mathcal{N}(z | \mu_y, \sigma^2 I)$$

$$p_r(y|z) = \frac{\exp(\tilde{x}^T w_y + b_y)}{\sum_{i=1}^n \exp(\tilde{x}^T w_i + b_i)}$$

To prove - For any choice of Θ , there exists γ such that

$$p_{\theta}(y|z) = p_{\gamma}(y|z)$$

Proof:

$$p_{\theta}(z, y) = \pi_y \mathcal{N}(z | \mu_y, \sigma^2 I)$$

$$\mathcal{N}(z | \mu_y, \sigma^2 I) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \|z - \mu_y\|^2\right)$$

$$p_{\theta}(y|z) = \frac{p_{\theta}(z, y)}{p_{\theta}(z)} = \frac{\pi_y \mathcal{N}(z | \mu_y, \sigma^2 I)}{\sum_{i=1}^k \pi_i \mathcal{N}(z | \mu_i, \sigma^2 I)}$$

$$p_{\theta}(y|z) = \frac{\pi_y \exp(-1/2\sigma^2 \|z - \mu_y\|^2)}{\sum_{i=1}^k \pi_i \exp(-1/2\sigma^2 \|z - \mu_i\|^2)}$$

Expanding the quadratic term:

$$\|x - \mu_y\|^2 = (x - \mu_y)^T (x - \mu_y) = x^T x - 2x^T \mu_y + \mu_y^T \mu_y$$

$$p_\theta(y|x) = \frac{\pi_y \exp\left(-\frac{x^T x}{2\sigma^2} + \frac{x^T \mu_y}{\sigma^2} - \frac{1}{2\sigma^2} \mu_y^T \mu_y\right)}{\sum_{i=1}^K \pi_i \exp\left(-\frac{x^T x}{2\sigma^2} + \frac{x^T \mu_i}{\sigma^2} - \frac{\mu_i^T \mu_i}{2\sigma^2}\right)}$$

$$p_\theta(y|x) = \frac{\pi_y \exp\left(\frac{x^T \mu_y}{\sigma^2} - \frac{\mu_y^T \mu_y}{2\sigma^2}\right) \cdot \exp\left(-\frac{x^T x}{2\sigma^2}\right)}{\exp\left(-\frac{x^T x}{2\sigma^2}\right) \sum_{i=1}^K \pi_i \exp\left(\frac{x^T \mu_i}{\sigma^2} - \frac{\mu_i^T \mu_i}{2\sigma^2}\right)}$$

$$p_\theta(y|x) = \frac{\pi_y \exp\left(\frac{x^T \mu_y}{\sigma^2} - \frac{\mu_y^T \mu_y}{2\sigma^2}\right)}{\sum_{i=1}^K \pi_i \exp\left(\frac{x^T \mu_i}{\sigma^2} - \frac{\mu_i^T \mu_i}{2\sigma^2}\right)}$$

we need to express the posterior in the form:

$$p_\gamma(y|x) = \frac{\exp(x^T w_y + b_y)}{\sum_{i=1}^K \exp(x^T w_i + b_i)}$$

to match the exponents of $p_\theta(y|x)$ with $p_\gamma(y|x)$,
let $w_y = \frac{\mu_y}{\sigma^2}$, $b_y = -\frac{\mu_y^T \mu_y}{2\sigma^2} + \log \pi_y$

Thus $p_\theta(y|x)$ can be written as

$$p_\theta(y|x) = \frac{\exp(x^T w_y + b_y)}{\sum_{i=1}^K \exp(x^T w_i + b_i)} = p_\gamma(y|x)$$

This is standard form of multi-class logistic regression

Q3)

n discrete random variables $\{X_i\}_{i=1}^n$ each having K_i different outcomes

A) Total no. of parameters needed to express the joint distribution without conditional independence

$$\begin{aligned}\text{Total number of joint outcomes} &= K_1 \times K_2 \times \dots \times K_n \\ &= \prod_{i=1}^n K_i\end{aligned}$$

But since probabilities sum to 1, we would need 1 parameter less to express the joint.

$$\text{Hence, total parameters} = \prod_{i=1}^n K_i - 1.$$

B) To use only $\sum_{i=1}^n (K_i - 1)$ parameters, the parameters for each variable need to be specified independently.

This is achieved when the random variables are mutually independent.

C) For $i \leq m$, No conditional independence. X_i depends on all predecessors

For $i > m$, X_i is conditionally independent of all its ancestors except most recent m ones

For a Bayesian network, no. of independent parameters:

(no. of values for X_{i-1}) \times (No. of parent config)

→ For $i \leq m$,

No. of possible parent config: $\prod_{j=1}^{i-1} K_j \quad \forall i=2, 3, \dots, m$

For each config, K_{i-1} free parameters.

\therefore Total parameters = $(K_{i-1}) \prod_{j=1}^{i-1} K_j$

Total for $i \leq m$, we get $\rightarrow \sum_{i=1}^m (K_{i-1}) \prod_{j=1}^{i-1} K_j$

→ For $i > m$,

Total no. of parent config: $\prod_{j=i-m}^{i-1} K_j$

Total no. of free parameters = $(K_{i-1}) \prod_{j=i-m}^{i-1} K_j$

Total for $i > m$, we get $\rightarrow \sum_{i=m+1}^n (K_{i-1}) \prod_{j=i-m}^{i-1} K_j$

Total parameters:

$$\sum_{i=1}^m (K_{i-1}) \prod_{j=1}^{i-1} K_j + \sum_{i=m+1}^n (K_{i-1}) \prod_{j=i-m}^{i-1} K_j$$

$$5) P(z) = \int_{\mathbf{z}} P(z, \mathbf{z}) d\mathbf{z}$$

$$A) A(z^{(1)}, \dots, z^{(K)}) = \frac{1}{K} \sum_{i=1}^K P(z | z^{(i)}), \quad z^{(i)} \sim p(z)$$

To show the above expression is an unbiased

estimator, we need to show,

$$E[A(z^{(1)}, \dots, z^{(K)})] = p(z)$$

proof:

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

$$\text{we know, } E_{z \sim p(z)}[p(z|z)] = \int p(z|z)p(z)dz = p(z)$$

$$\begin{aligned} \therefore E[A(z^{(1)}, \dots, z^{(K)})] &= \frac{1}{K} \sum_{i=1}^K p(z) \\ &= \frac{1}{K} \cdot K \cdot p(z) \\ &= p(z) \end{aligned}$$

Thus the Monte Carlo estimator A provides an unbiased estimate of $p(z)$.

B) To find if $\log A$ is an unbiased estimator or not,

$$\text{we need to check if } E[\log A] = \log p(z)$$

Bring Jensen's inequality for concave function

$$f(E[X]) \geq E[f(X)]$$

$$\log(E[A]) \geq E[\log A]$$

we know $E[\log A] = p(z)$

$\therefore \log(p(z)) > E[\log A]$, which shows
 $\log A$ is not an unbiased estimator.

6)

a) Minimum bits to represent 50257 tokens

$$2^n > 50257$$

$$2^{15} = 32768, 2^{16} = 65536$$

Thus minimal $n = 16$ bits

b) Increase in no. of parameters when expanding
from 50257 to 60000

$$\text{Original parameters} = 50257 \times 768$$

$$\text{New parameters} = 60000 \times 768$$

$$\text{Difference} = (60000 - 50257) \times 768 = 7482624$$

$$\begin{aligned} \text{Total difference} &= \text{contribution from embedding} + \\ &\quad \text{Fully connected layer} \\ &= 7482624 \times 2 = 14974848 \end{aligned}$$

4) Setup with example, where $n=2$

$$P_f(z_1, z_2) = P_f(z_1) P_f(z_2 | z_1)$$

- Marginal for $z_1 \rightarrow P_f(z_1) = N(z_1 | 0, 1)$
- conditional for z_2 given z_1 :

$$p_f(z_2|z_1) = \mathcal{N}(z_2 | \mu_2(z_1), \epsilon),$$

$$\text{where } \mu_2(z_1) = \begin{cases} 0 & \text{if } z_1 \leq 0 \\ 1 & \text{if } z_1 > 0 \end{cases}$$

$$p_f(z_2) = \int p_f(z_1) p_f(z_2|z_1) dz_1$$

$$\because z_1 \sim \mathcal{N}(0, 1):$$

$$\rightarrow \text{Probability that } z_1 \leq 0 \approx 0.5$$

$$\rightarrow \text{Probability that } z_1 > 0 \approx 0.5$$

$$\text{For } z_1 \leq 0,$$

$$\mu_2(z_1) = 0 \text{ \& } p_f(z_1 \leq 0) = 0.5$$

$$\therefore p_f(z_2) = 0.5 \mathcal{N}(z_2 | 0, \epsilon)$$

$$\text{For } z_1 > 0$$

$$\mu_2(z_1) = 1 \text{ \& } p_f(z_1 > 0) = 0.5$$

$$\therefore p_f(z_2) = 0.5 \mathcal{N}(z_2 | 1, \epsilon)$$

The combined marginal:

$$p_f(z_2) \approx 0.5 \mathcal{N}(z_2 | 0, \epsilon) + 0.5 \mathcal{N}(z_2 | 1, \epsilon)$$

For the reverse autoregressive model,

$$p_g(z_1, z_2) = p_g(z_1|z_2) p_g(z_2)$$

$$p_g(z_2) = \mathcal{N}(z_2 | \mu^2(0), \sigma^2(0))$$

- In the reverse process, z_2 being the first variable has no previous variable to condition on.
 - $p_{\pi}(z_2)$ is hence just defined by a single gaussian & is always Unimodal, defined by a single peak.
 - This is in contrast to $p_f(z_2)$ which is a bimodal distribution being a mixture 2 Gaussian with no overlap
 - Since $p_f(z_2) \neq p_{\pi}(z_2)$, the forward & reverse auto regressive models do not cover the same hypothesis space.
-