

(Fall 2021)

Shubham Shrivastava

email: shubhams@stanford.edu

P1 Data distribution:  $\hat{p}(x, y)$   
 inputs:  $x \in X$   
 outputs:  $y \in Y$

KL divergence is defined as:

$$D_{KL}(p(x) \| q(x)) = E_{x \sim p(x)} [\log p(x) - \log q(x)]$$

We have to show that:

$$\begin{aligned} & \arg \max_{\theta \in \Theta} E_{\hat{p}(x, y)} [\log \hat{p}_\theta(y|x)] \\ &= \arg \min_{\theta \in \Theta} E_{\hat{p}(x)} [D_{KL}(\hat{p}(y|x) \| p_\theta(y|x))] \end{aligned}$$

i

$\Rightarrow$  Let's expand the KL term:

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

ii

$$\begin{aligned} \Rightarrow & \arg \min_{\theta \in \Theta} E_{\hat{p}(x)} [D_{KL}(\hat{p}(y|x) \| p_\theta(y|x))] = \\ & \arg \min_{\theta \in \Theta} E_{\hat{p}(x)} \left[ E_{y \sim \hat{p}(y|x)} [\log \hat{p}(y|x) - \log p_\theta(y|x)] \right] \end{aligned}$$

$$\begin{aligned}
 &= \arg \min_{\theta \in \Theta} E_{\hat{P}(x)}^1 \left[ \log \hat{p}(y|x) \right] \\
 &\quad - \arg \min_{\theta \in \Theta} E_{\hat{P}(x)}^1 \left[ \log p_\theta(y|x) \right] \\
 &\quad \text{---} \\
 &\quad E_{\hat{P}(x)}^1 \cdot \hat{p}(y|x) = E_{\hat{P}(x,y)}^1 \\
 &= \arg \min_{\theta \in \Theta} E_{\hat{P}(x,y)}^1 \left[ \log \hat{p}(y|x) \right] \\
 &\quad - \arg \min_{\theta \in \Theta} E_{\hat{P}(x,y)}^1 \left[ \log p_\theta(y|x) \right] - \text{(iii)}
 \end{aligned}$$

The first term in above equation is independent of  $\theta$  and hence does not contribute to the optimization problem — this term can be replaced by a constant 'C'. Also, utilizing the properties in equation (iv), we can re-write (iii) as (iv)

$$\max_{\theta} f(\theta) \cong \min_{\theta} [-f(\theta)] - \text{(iv)}$$

So eq<sup>n</sup> (iii) can be re-written as

$$\boxed{\arg \max_{\theta \in \Theta} E_{\hat{P}(x,y)}^1 \left[ \log p_\theta(y|x) \right]} - \text{(iv)}$$

P2 — The generative process is given as:

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

From Bayes rule:

$$p_{\theta}(y|x) = \frac{p_{\theta}(x|y) p_{\theta}(y)}{p_{\theta}(x)} \quad \text{--- (i)}$$

where,  $p_{\theta}(y)$  is given to be  $\pi_y$ ;  $\sum_{j=1}^k \pi_j = 1$

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

,  $p_{\theta}(x)$  is a mixture of  $k$  gaussians

expanding the numerator in eq<sup>n</sup>=i

$$p_{\theta}(x|y) p_{\theta}(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu_y)^T(x-\mu_y)}{2\sigma^2}\right) \cdot \pi_y \quad \text{--- (ii)}$$

Also,  $p_{\theta}(x)$  can be represented as:

$$p_{\theta}(x) = \sum_{i=1}^k \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu_i)^T(x-\mu_i)}{2\sigma^2}\right) \quad \text{--- (iii)}$$

Substituting (ii) and (iii) in eq<sup>n</sup>(i)

$$p_\theta(y|x) = \frac{\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu_y)^T(x-\mu_y)}{2\sigma^2}\right) \pi_y}{\sum_{j=1}^K \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu_j)^T(x-\mu_j)}{2\sigma^2}\right)}$$

$$= \frac{\cancel{\frac{1}{\sqrt{2\pi\sigma^2}}} \exp\left(-\frac{(x-\mu_y)^T(x-\mu_y)}{2\sigma^2}\right) \pi_y}{\cancel{\frac{1}{\sqrt{2\pi\sigma^2}}} \sum_{j=1}^K \exp\left(-\frac{(x-\mu_j)^T(x-\mu_j)}{2\sigma^2}\right)}$$

$$= \frac{\exp\left(\frac{-1}{2\sigma^2}(x^T x - 2x^T \mu_y + \mu_y^T \mu_y)\right) \pi_y}{\sum_{j=1}^K \exp\left(\frac{-1}{2\sigma^2}(x^T x - 2x^T \mu_j + \mu_j^T \mu_j)\right)}$$

$$= \frac{\pi_y \cdot \exp\left(\frac{-x^T x}{2\sigma^2}\right) / \exp\left(\frac{1}{2\sigma^2}(-2x^T \mu_y + \mu_y^T \mu_y)\right)}{\sum_{j=1}^K \left[ \exp\left(\frac{-x^T x}{2\sigma^2}\right) / \exp\left(\frac{1}{2\sigma^2}(-2x^T \mu_j + \mu_j^T \mu_j)\right) \right]}$$

$$= \frac{\pi_y \cdot \exp\left(\frac{-x^T x}{2\sigma^2}\right) / \exp\left(\frac{1}{2\sigma^2}(-2x^T \mu_y + \mu_y^T \mu_y)\right)}{\sum_{j=1}^K \left[ \exp\left(\frac{-x^T x}{2\sigma^2}\right) / \exp\left(\frac{1}{2\sigma^2}(-2x^T \mu_j + \mu_j^T \mu_j)\right) \right]}$$

$$\begin{aligned}
&= \frac{\pi_j \cdot \exp\left(\frac{-x^T x}{2\sigma^2}\right) / \exp\left(\frac{1}{2\sigma^2}(-2x^T \mu_j + \mu_j^T \mu_j)\right)}{\sum_{i=1}^K \left[ \exp\left(\frac{-x^T x}{2\sigma^2}\right) / \exp\left(\frac{1}{2\sigma^2}(-2x^T \mu_i + \mu_i^T \mu_i)\right) \right]} \\
&= \frac{\pi_j \cdot \exp\left(\frac{-x^T x}{2\sigma^2}\right) / \exp\left(\frac{1}{2\sigma^2}(-2x^T \mu_j + \mu_j^T \mu_j)\right)}{K \cdot \exp\left(\frac{-x^T x}{2\sigma^2}\right) / \sum_{i=1}^K \exp\left(\frac{1}{2\sigma^2}(-2x^T \mu_i + \mu_i^T \mu_i)\right)} \\
&= \frac{\exp\left(\frac{-1}{2\sigma^2}(2x^T \mu_j - \mu_j^T \mu_j)\right) \pi_j}{\sum_{i=1}^K \left[ \exp\left(\frac{-1}{2\sigma^2}(2x^T \mu_i - \mu_i^T \mu_i)\right) \right] + 1}
\end{aligned}$$

$\pi_j$  can also be written as  $\exp(\log_e \pi_j)$ ,  
and 1 can be rewritten as  $\exp(\log_e \pi_i)$  where,  $\pi_i = 1$

$$\Rightarrow p_{\theta}(y|x) = \frac{\exp\left(\frac{-1}{2\sigma^2}(2x^T\mu_y - \mu_y^T\mu_y + \log_e \pi_y)\right)}{\sum_{j=1}^K \left[\exp\left(\frac{-1}{2\sigma^2}(2x^T\mu_j - \mu_j^T\mu_j + \log_e \pi_j)\right)\right]}$$

$$= \frac{\exp\left(x^T\left(\frac{-2\mu_y}{2\sigma^2}\right) + \left(\frac{\mu_y^T\mu_y + \log_e \pi_y}{2\sigma^2}\right)\right)}{\sum_{j=1}^K \exp\left(x^T\left(\frac{-2\mu_j}{2\sigma^2}\right) + \left(\frac{\mu_j^T\mu_j + \log_e \pi_j}{2\sigma^2}\right)\right)}$$

(iv)

If we let  $w_y = \frac{-2\mu_y}{2\sigma^2}$ ;  $b_y = \frac{\mu_y^T\mu_y + \log_e \pi_y}{2\sigma^2}$

$$w_i = \frac{-2\mu_i}{2\sigma^2}; b_i = \frac{\mu_i^T\mu_i + \log_e \pi_i}{2\sigma^2}$$

then equation (iv) becomes

$$p_{\theta}(y|x) = \frac{\exp(x^T w_j + b_j)}{\sum_{i=1}^K \exp(x^T w_i + b_i)} = p_y(y|x)$$

P<sub>3</sub>

① Without any conditional independence assumption, the joint distribution can be expressed using the chain rule:

$$\text{i.e. } p(x_1, x_2, \dots, x_n) = p(x_1) \underbrace{p(x_2 | x_1)}_{\text{Number of parameters} \rightarrow (k_1 - 1)} \underbrace{\dots}_{(k_2 - 1)} \underbrace{p(x_n | x_1 \dots x_{n-1})}_{(k_n - 1)}$$

So, total number of parameters:

$$(k_1 - 1) + k_1(k_2 - 1) + k_1 k_2 (k_3 - 1) + \dots + k_1 k_2 \dots k_{n-1} (k_n - 1)$$
$$= (k_1 - 1) + \sum_{i=2}^n (k_i - 1) \prod_{j=1}^{i-1} k_j$$

② If the joint distribution are assumed to be independent then they can be represented as

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

Number of parameters needed to represent

$$x_1 \rightarrow (k_1 - 1), x_2 \rightarrow (k_2 - 1), \dots, x_n \rightarrow (k_n - 1)$$

So, total number of parameters:

$$\sum_{i=1}^n (k_i - 1)$$

③ The joint distribution can be expressed as following using chain rule:

$$p(x_1, x_2, x_3, \dots, x_n) = \underbrace{p(x_1, x_2, x_3, \dots, x_m)}_{\text{no conditional independence assumption}} \cdot p(x_{m+1}, x_{m+2}, \dots, x_n | x_1, x_2, \dots, x_m)$$

For the first term, number of parameters needed to specify the joint distribution is:

$$(k, -1) + \sum_{i=2}^m (k_i, -1) \prod_{j=1}^{i-1} k_j$$

Let's expand the second term:

$$\begin{aligned} p(x_{m+1}, x_{m+2}, \dots, x_n | x_1, x_2, \dots, x_m) &= p(x_{m+1} | x_1, \dots, x_m) \\ &\quad p(x_{m+2} | x_1, \dots, x_{m+1}) \\ &\quad \vdots \\ &\quad p(x_n | x_1, \dots, x_{n-1}) \end{aligned}$$

Given the assumption that for  $j > m$ , the random variable  $X_j$  is conditional independent of all ancestors given the previous  $m$  ancestors, above eq<sup>n</sup> can be re-written as:

$$\begin{aligned}
 p(x_{m+1}, x_{m+2}, \dots, x_n | x_1, x_2, \dots, x_m) &= p(x_{m+1} | x_1, \dots, x_m) \\
 &\quad p(x_{m+2} | x_2, \dots, x_{m+1}) \\
 &\quad \vdots \\
 &\quad p(x_n | x_{n-m}, \dots, x_{n-1})
 \end{aligned}$$

Number of parameters needed to specify the above distribution:

$$\begin{aligned}
 & (k_{m+1}-1) k_1 k_2 \dots k_m + (k_{m+2}-1) k_2 k_3 \dots k_{m+1} \\
 & + \dots + (k_n-1) k_{n-m} k_{n-m+1} \dots k_{n-1} \\
 & = \sum_{i=m+1}^n (k_i-1) \prod_{j=i-m}^{i-1} k_j
 \end{aligned}$$

Hence, total number of parameters needed to specify the full joint:

$$(k_1-1) + \sum_{i=2}^m (k_i-1) \prod_{j=1}^{i-1} k_j + \sum_{i=m+1}^n (k_i-1) \prod_{j=i-m}^{i-1} k_j$$

P4

Let's consider the case where  $n=2$

$$p_f(x_1, x_2) = p_f(x_1) p_f(x_2 | x_1)$$

where,  $p_f(x_i) = \mathcal{N}(x_i | \mu_i(0), \sigma_i^2(0))$

Since,  $x_{\leq i} = \begin{cases} [x_1, \dots, x_{i-1}]^T & \text{if } i \geq 1 \\ 0 & \text{if } j < 1 \end{cases}$

It is given that  $R^0 = \{0\} \Rightarrow \mu_i(0) = c_1 \text{ & } \sigma_i^2(0) = c_2$   
for some constants  $c_1 \text{ & } c_2$

$$\Rightarrow p_f(x_1) = \mathcal{N}(x_1 | c_1, c_2^2)$$

$$p_f(x_2 | x_1) = \mathcal{N}(x_2 | \mu_2(x_1), \sigma_2^2(x_1))$$

Using Bayes rule:  $p_f(x_1 | x_2) = \frac{p_f(x_2 | x_1)}{p_f(x_2)}$

where,  $p_f(x_2) = \int_{x_1=-\infty}^{\infty} p_f(x_1) p_f(x_2 | x_1) dx_1$

$$= \int_{x_1=-\infty}^{\infty} \mathcal{N}(x_1 | c_1, c_2^2) \cdot \mathcal{N}(x_2 | \mu_2(x_1), \sigma_2^2(x_1)) dx_1$$

multi-modal gaussian

(1)

Let's look at the reverse function:

$$p_r(x_1, x_2) = p_r(x_1 | x_2) p_r(x_2)$$

where,  $p_r(x_2) = \mathcal{N}(x_2 | \underbrace{c_3, c_4^2}_{\text{constants}}) \quad (ii)$

It is clear from equations (i) and (ii) that  $p_f(x_2)$  is a multivariate gaussian whereas  $p_r(x_2)$  is a univariate gaussian. Hence, for some choice of  $(\mu_i, \sigma_i)$ , the neural network  $\{\hat{\mu}_i, \hat{\sigma}_i\}$  cannot represent the same distribution.

PS ①  $A$  is given as:

$$A(z^{(1)}, \dots, z^{(k)}) = \frac{1}{k} \sum_{j=1}^k p(x|z^{(j)}) \text{, where } z^{(i)} \sim p(z)$$

Let's compute  $E[A]$

$$E[A(z^{(1)}, \dots, z^{(k)})] = \frac{1}{k} \sum_{i=1}^k E_{z^{(i)}} p(x|z^{(i)})$$

$\underbrace{\phantom{E_{z^{(i)}}} p(z)}$

$$= p(z) p(x|z)$$

$$= \int_z p(x,z) dz$$

$$= p(x)$$

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

② From Jensen's inequality

$$E[\log y] \leq \log E[y]$$

$$\Rightarrow E[\log(A(z^{(1)}, \dots, z^{(k)}))] \leq \log E[A(z^{(1)}, \dots, z^{(k)})]$$

$$\Rightarrow E[\log(A(z^{(1)}, \dots, z^{(k)}))] \leq \log b(x) \xrightarrow{\text{from previous solution}}$$

An estimator  $\hat{\theta}$  is an unbiased estimator of  $\theta$  iff  $E[\hat{\theta}] = \theta$ . Here we have :  $E[\log A] \leq \log b(x)$   
Hence,  $\log A$  is not an unbiased estimator of  $\log(b(x))$

**P6** ① 50257 tokens represented by n-bits

$(a_1, a_2, a_3, \dots, a_n)$ , where  $a_i \in \{0, 1\}$ ,  $\forall i = 1, 2, \dots, n$

for minimal n,  $2^n \geq 50,257$

$$\Rightarrow n = 16 \text{-bits}$$

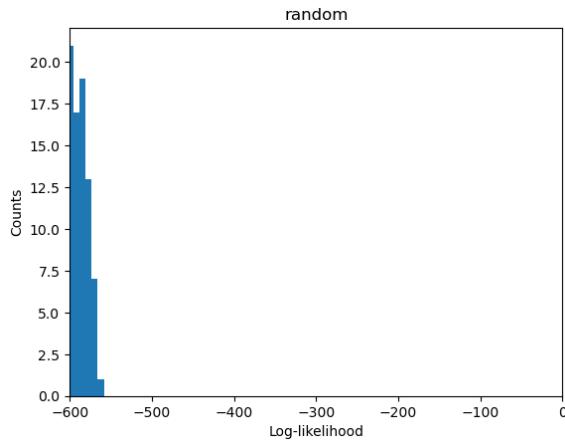
② Number of increased dimension:  $60000 - 50257$   
 $= 9,743$

For fully-connected layers, number of additional parameters:  $768 \times 9,743 + 9743 = 7,492,367$

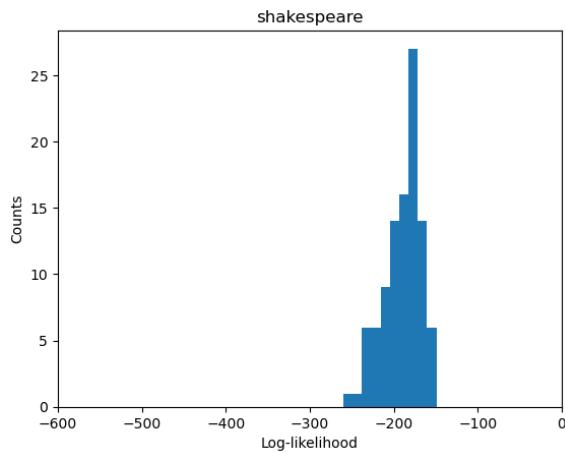
So, total increase in number of parameters:

$$7,492,367$$

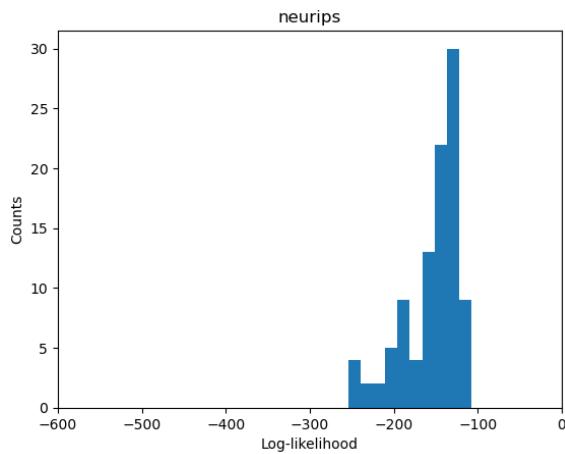
## 1 Question 6.4



(a) Histogram of the log-likelihoods of random strings



(b) Histogram of the log-likelihoods of string snippets from Shakespeare's work



(c) Histogram of the log-likelihoods of string snippets from NeurIPS 2015

Figure 1: Histogram of the log-likelihoods of strings.