

ITCT Midterm

- R07944007 網媒碩一 林良翰

1.

a. Chain Rule for Relative Entropy

- By the definition of relative entropy

$$\begin{aligned} D(p(x, y) \parallel q(x, y)) &= \sum_{x \in X} \sum_{y \in Y} \left[p(x, y) \log \frac{p(x, y)}{q(x, y)} \right] \\ &= \sum_{x \in X} \sum_{y \in Y} [p(x, y) \log p(x, y) - p(x, y) \log q(x, y)] \\ &= \sum_{x \in X} \sum_{y \in Y} [p(x, y) \log (p(x) p(y|x)) - p(x, y) \log (q(x) q(y|x))] \\ &= \sum_{x \in X} \sum_{y \in Y} \left[p(x, y) \log \frac{p(x)}{q(x)} + p(x, y) \log \frac{p(y|x)}{q(y|x)} \right] \\ &= \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} + \sum_{x \in X} p(x) \sum_{y \in Y} p(y|x) \log \frac{p(y|x)}{q(y|x)} \\ &= D(p(x) \parallel q(x)) + D(p(y|x) \parallel q(y|x)) \end{aligned}$$

b. Jensen's Inequality

- By the definition of the concave function, $p_1 + p_2 = 1$

$$p_1 f(x_1) + p_2 f(x_2) \leq f(p_1 x_1 + p_2 x_2)$$

- Suppose the theorem $E[f(X)] \leq f(E[X])$ is true for distributions with $K - 1$ mass points when f is a concave function

$$\text{Let } p'_i = \frac{p_i}{1 - p_k}$$

$$\begin{aligned}
f(E[X]) &= f\left(\sum_{i=1}^k p_i x_i\right) \\
&= f\left(\sum_{i=1}^k (1 - p_k) p'_i x_i\right) \\
&= f\left((1 - p_k) \sum_{i=1}^k p'_i x_i\right) \\
&= f\left((1 - p_k) p'_k x_k + (1 - p_k) \sum_{i=1}^{k-1} p'_i x_i\right) \\
&= f\left(p_k x_k + (1 - p_k) \sum_{i=1}^{k-1} p'_i x_i\right) \\
&\geq p_k f(x_k) + (1 - p_k) f\left(\sum_{i=1}^{k-1} p'_i x_i\right) \\
&\geq p_k f(x_k) + (1 - p_k) \sum_{i=1}^{k-1} p'_i f(x_i) \\
&= \sum_{i=1}^k p_i f(x_i)
\end{aligned}$$

c. Convexity of Mutual Information

- We want to prove $I(X; Y)$ is a convex function of $p(Y|X)$ for fixed $p(X)$

It is complicated to find approve directly from the definition of $I(X; Y)$

- We introduce an auxiliary random variable \tilde{Y} with a mixing distribution

$$p(\tilde{y}|x) = \lambda p_1(y|x) + (1 - \lambda) p_2(y|x)$$

- To prove convexity, we need to prove

$$I(X; \tilde{Y}) \leq \lambda I_{p_1}(X; Y) + (1 - \lambda) I_{p_2}(X; Y)$$

- Since $I(X; \tilde{Y}) = D(p(x, \tilde{y}) || p(x) p(\tilde{y}))$

$$\begin{aligned}
p(\tilde{y}) &= \sum_{x \in X} p(x) p(\tilde{y}|x) \\
&= \sum_{x \in X} p(x) \lambda p_1(y|x) + p(x) (1 - \lambda) p_2(y|x) \\
&\stackrel{p(x) \text{ fixed}}{=} \lambda p_1(y) + (1 - \lambda) p_2(y)
\end{aligned}$$

$$\begin{aligned}
p(x, \tilde{y}) &= p(x) p(\tilde{y}|x) \\
&= p(x) \lambda p_1(y|x) + p(x) (1 - \lambda) p_2(y|x) \\
&= \lambda p_1(x, y) + (1 - \lambda) p_2(x, y)
\end{aligned}$$

- Finally, by the convexity of $D(p||q)$

$$\begin{aligned}
D(p(x, \tilde{y}) || p(x) p(\tilde{y})) &= D(\lambda p_1(x, y) + (1 - \lambda) p_2(x, y) || p(x) \lambda p_1(y) + p(x) (1 - \lambda) p_2(y)) \\
&\leq \lambda D(p_1(x, y) || p(x) p_1(y)) + (1 - \lambda) D(p_2(x, y) || p(x) p_2(y)) \\
&= \lambda I_{p_1}(X; Y) + (1 - \lambda) I_{p_2}(X; Y)
\end{aligned}$$

$$\Rightarrow I(X; \tilde{Y}) \leq \lambda I_{p_1}(X; Y) + (1 - \lambda) I_{p_2}(X; Y)$$

d. Fano's Inequality

- Define a random variable E with following properties

$$E = \begin{cases} 0, & X = \hat{X} \\ 1, & X \neq \hat{X} \end{cases}$$

- By the chain rule for entropies

$$\begin{aligned} H(E, X | \hat{X}) &= H(X | \hat{X}) + H(E | X, \hat{X}) \\ &= H(E | \hat{X}) + H(X | E, \hat{X}) \end{aligned}$$

- By the data processing inequality

$$H(X | \hat{X}) \geq H(X | Y)$$

- E is the function of X, \hat{X}

$$H(E | X, \hat{X}) = 0$$

- By the property of conditional entropy

$$H(E | \hat{X}) \leq H(E) = H(P_e)$$

- If \hat{X} is known and $E = 0$, which means the prediction is correct

$$H(X | E = 0, \hat{X}) = 0$$

If \hat{X} is known and $E = 1$, which means the prediction is not correct, but we can assure that $X \neq \hat{X} \Rightarrow X$ has $|X| - 1$ possible ground truth

$$H(X | E = 1, \hat{X}) = \log(|X| - 1)$$

Therefore

$$\begin{aligned} H(X | E, \hat{X}) &= P(E = 0) H(X | E = 0, \hat{X}) + P(E = 1) H(X | E = 1, \hat{X}) \\ &= (1 - P_e) 0 + P_e \log(|X| - 1) \end{aligned}$$

- Combining all the inequations above

$$H(X | Y) \leq H(P_e) + P_e \log(|X| - 1)$$

Weaken the boundary

$$H(X | Y) \leq 1 + P_e \log |X|$$

$$P_e \geq \frac{H(X|Y)-1}{\log |X|}$$

a. Uniform distribution of finite discrete sources gives Maximal Entropy

- Let p be the probability function of random variable X
- Let U be uniform distributed random variable, $u(x) = \frac{1}{|X|}$
- By the fact that relative entropy is always larger or equals to 0

$$\begin{aligned}
 0 &\leq D(p||u) \\
 &= \sum_{x \in X} p(x) \log \frac{p(x)}{u(x)} \\
 &= \sum_{x \in X} p(x) \log |X| + \sum_{x \in X} p(x) \log p(x) \\
 &= \log |X| - H(X)
 \end{aligned}$$

- $\Rightarrow H(X) \leq \log |X|$

If $H(X)$ reaches maximum

$$\begin{aligned}
 H(X) &= \log |X| \\
 - \sum_{x \in X} p(x) \log p(x) &= \sum_{x \in X} p(x) \log |X| \\
 \Rightarrow p(x) &= \frac{1}{|X|}
 \end{aligned}$$

- p must be a uniform distribution to maximize it's entropy

b. Normal distribution of infinite continuous sources gives Maximal Differential Entropy

- Let X be a random variable with a probability density function f , the differential entropy of X is defined as

$$h(X) = \int_{-\infty}^{\infty} f(x) \log f(x) dx$$

- Let $g(x)$ be a Gaussian PDF with mean μ and variance σ^2 and $f(x)$ and arbitrary PDF with same variance

$$h(g) = \frac{1}{2} \log(2\pi e \sigma^2)$$

- Apply the same method \rightarrow the relative entropy is always larger or equals to 0

$$\begin{aligned}
 0 &\leq D(f||g) \\
 &= \int_{-\infty}^{\infty} f(x) \log \frac{f(x)}{g(x)} dx \\
 &= -h(f) - \int_{-\infty}^{\infty} f(x) \log(g(x)) dx
 \end{aligned}$$

Note that

$$\begin{aligned}
\int_{-\infty}^{\infty} f(x) \log(g(x)) dx &= \int_{-\infty}^{\infty} f(x) \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}\right) dx \\
&= \int_{-\infty}^{\infty} f(x) \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) dx + \log(e) \int_{-\infty}^{\infty} f(x) \left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \\
&= -\frac{1}{2} \log(2\pi\sigma^2) - \log(e) \frac{\sigma^2}{2\sigma^2} \\
&= -\frac{1}{2} \log(2\pi e \sigma^2) \\
&= -h(g)
\end{aligned}$$

- Therefore

$$h(f) \leq h(g)$$

and the equality holds when $f = g$

3.

a.

- Show that a suffix code is uniquely decodable
- Suppose we have a code X satisfying prefix condition.
- If X doesn't satisfy unique decodable, then we can find a code $A = (a_1, a_2, \dots, a_N)$ from X such that this code can be constructed from other codes.
- We can extract a code $B = (b_1, b_2, \dots, b_K)$, $K \leq N$ from A , then

$$a_1 = b_1$$

$$a_2 = b_2$$

...

$$a_K = b_K$$

Thus B becomes the prefix of A , and this contradicts with the fact that A satisfies prefix condition $\Rightarrow B$ not exists, A is uniquely decodable.

- Therefore, when a code satisfies suffix condition, it must also be uniquely decodable.

b.

- Show that the minimum average length over all codes satisfying the suffix condition is the same as the average length of the Huffman code for that random variable.
- If the average code length of the shortest code which its suffix code A is shorter than Huffman code, then we reverse this code A' . The suffix code of A' will become prefix code of A , thus we get a set of prefix code which has average code length shorter than Huffman code. This violates the fact that Huffman code is the optimal prefix code. Therefore, the suffix code which has average code length shorter than Huffman code does not exist.

- If the average code length of the shortest code which its suffix code is longer than Huffman code, where the Huffman code is prefix code, we can reverse all the codes in Huffman code, and these will become suffix code. These suffix code has average code length shorter than the suffix code of the code with shortest average code length we found in the beginning \Rightarrow contradict.
- By the two cases above, we can conclude that the suffix code of the code with shortest average length has the same average code length of Huffman code.

4.

a.

- Step 1 - always divide the largest node (largest probability)
- Step 2 - divide until having 2^l nodes

b.

- $P_A = 0.7, P_B = 0.2, P_C = 0.1$
- Example input - *ABCCBAAAC*
- The code tree

$$\left\{ \begin{array}{l} A - 0.7 \\ B - 0.2 \\ C - 0.1 \end{array} \right\} \left\{ \begin{array}{l} AA - 0.49 \\ AB - 0.14 \\ AC - 0.07 \end{array} \right\} \left\{ \begin{array}{l} AAA - 0.343 \\ AAB - 0.098 \\ AAC - 0.049 \end{array} \right\} \left\{ \begin{array}{l} AAAA \\ AAAB \\ AAAC \end{array} \right\}$$

c.

- "variable source symbol"-to-"fixed length codeword"
 - Pros - certain strings of symbols are frequently repeated, strings can be assigned code words that represent the "entire string of symbols".
 - Cons - usually efficient only for long files or messages.
- Huffman "variable codeword length"-to-"fixed source symbol"
 - Pros - asymptotically approach the source entropy for long messages, the receiver does not require prior knowledge of the coding table constructed by the transmitter, allow the receiver to construct its own decoding table "on the fly", the information required to do this is transmitted early in the coded messages.
 - Cons - initially "expand" rather than "compress" the data, less and less information need be sent to aid the receiver as time progresses.

5.

a.

- 這堂課到目前為止給了我最大的收穫是重新構建了我在機率方面的知識，尤其是 Entropy 這一部分。其實在大二學機率的時候，看到 KL-Divergence, Mutual Information,... 之類的名詞，往往都只能背公式，而無法深入去理解它背後的含義，但是上過這門課之後，我對他們的定義至少有完整的基礎概念，而且還多學到了很多不等式，其中的 Jensen's Inequality 我覺得是我印象最深刻的，因為他的公式非常直覺，而且證明起來也不會用到太難的數學基礎。

b.

- 這個學期中有介紹 JPEG 編碼技術，因此我也想到如果可以的話，希望也能聽到關於 PNG 方面的技術，因為他還能夠記錄透明度，我個人就比較好奇如果是透明度的話，那該如何去壓縮？然後他和 JPEG 比較的話，又有哪些優缺點？