

CSC411 Assignment 2

1. (a) WTP: $H(X) = \sum_x p(x) \log_2 \left(\frac{1}{p(x)} \right) \geq 0$

• Since X is a discrete random variable, probability mass function $p(x) \in [0, 1]$

• Case 1: When $p(x) = 0$, by the condition in the question:

$$p(x) \log_2 \left(\frac{1}{p(x)} \right) = 0$$

• Case 2: When $p(x) \in (0, 1] \Leftrightarrow \frac{1}{p(x)} \in [1, +\infty) \Leftrightarrow \log_2 \left(\frac{1}{p(x)} \right) > 0$
 $\Leftrightarrow p(x) \log_2 \left(\frac{1}{p(x)} \right) > 0$

• Therefore, for each possible value of $x \in X$, we have:

$$p(x) \log_2 \left(\frac{1}{p(x)} \right) \geq 0$$

$$\Leftrightarrow \sum_x p(x) \log_2 \left(\frac{1}{p(x)} \right) \geq 0$$

$\Leftrightarrow H(X) \geq 0$, i.e.: entropy $H(X)$ is non-negative.

(b) WTP: $KL(p \parallel q) = \sum_x p(x) \log_2 \frac{p(x)}{q(x)} \geq 0$.

• By appendix: $\log(x)$ is concave on the set of positive real number

→ In order to use Jensen's inequality, which states:

if $\phi(t)$ is a convex function of t , then:

$$\phi(E[t]) \leq E[\phi(t)]$$

→ We need to transfer concave function $\log(x)$ to convex one.

By the defⁿ of concave function, we have:

f is concave if $-f$ is convex

→ Therefore $-\log(x)$ is convex function on the set of positive real number.

$$\begin{aligned} KL(p \parallel q) &= \sum_x p(x) \log_2 \frac{p(x)}{q(x)} \\ &= \sum_x p(x) \left(-\log_2 \frac{q(x)}{p(x)} \right) && \# \text{ transfer to convex function.} \\ &= E \left[-\log_2 \frac{q(x)}{p(x)} \right] && \# \phi(t) = -\log_2 t; t = \frac{q(x)}{p(x)} \\ &\geq -\log_2 \left(\sum_x p(x) \frac{q(x)}{p(x)} \right) && \# \text{ By Jensen's inequality.} \\ &= -\log_2 \left(\sum_x q(x) \right) = -\log_2 1 = 0 \end{aligned}$$

• Therefore, $KL(p \parallel q) \geq 0$, i.e.: $KL(p \parallel q)$ is non-negative.

(c) WTS: $I(Y; X) = KL(p(x, y) \parallel p(x)p(y))$

where $p(x) = \sum_y p(x, y)$ is the marginal distribution of X .

$$\begin{aligned} I(Y; X) &= H(Y) - H(Y|X) && \# \text{ By defⁿ of } I(Y; X). \\ &= \left(-\sum_y p(y) \log p(y) \right) - \left(-\sum_x \sum_y p(x, y) \log p(y|x) \right) && \# \text{ By defⁿ of } H(X) \\ &= \left(-\sum_y \left(\sum_x p(x, y) \right) \log p(y) \right) + \sum_x \sum_y p(x, y) \log p(y|x) && \# \text{ By marginal distribution of } Y \\ &= \sum_x \sum_y p(x, y) \log \frac{1}{p(y)} + \sum_x \sum_y p(x, y) \log p(y|x) \\ &= \sum_x \sum_y p(x, y) \left(\log \frac{1}{p(y)} + \log p(y|x) \right) \\ &= \sum_x \sum_y p(x, y) \cdot \log \frac{p(x, y)}{p(x)p(y)} && (*) \end{aligned}$$

$$\begin{aligned} KL(p(x, y) \parallel p(x)p(y)) &= \sum_{x, y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} && \# \text{ By defⁿ of } KL \\ &= \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)} && (***) \end{aligned}$$

• By (*) and (**), we conclude that:

$$I(Y; X) = KL(p(x, y) \parallel p(x)p(y))$$

2. WTP: $\forall x, t \quad L(\bar{h}(x), t) \leq \frac{1}{m} \sum_{i=1}^m L(h_i(x), t)$

where $\bar{h}(x) = \frac{1}{m} \sum_{i=1}^m h_i(x)$

- $\bar{h}(x) = \frac{1}{m} \sum_{i=1}^m h_i(x) = E[h_i(x)]$ # expectation of $h_i(x)$ (x)
- $L(\bar{h}(x), t) = \frac{1}{2} (\bar{h}(x) - t)^2$ # By defⁿ of squared error loss function
- $= \frac{1}{2} (E[h_i(x)] - t)^2$ # By (2), let $\phi(s) = \frac{1}{2}(s-t)^2, s = h_i(x)$
- $\leq E[\frac{1}{2} (h_i(x) - t)^2]$ # By Jensen's inequality,
- $= \frac{1}{2} (E[(h_i(x) - t)^2])$ # By the property of expected value.
- $= \frac{1}{2} \cdot \frac{1}{m} \sum_{i=1}^m (h_i(x) - t)^2$
- $= \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_i(x) - t)^2$
- $= \frac{1}{m} \sum_{i=1}^m L(h_i(x), t)$ # By defⁿ of squared error loss function

3. WTS: $err'_t = \frac{\sum_{i=1}^N w_i I\{h_t(x^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^N w_i} = \frac{1}{2}$

• By the condition in the question we have:

$t^{(i)} \in \{-1, 1\}$ and $h_t(x^{(i)}) \in \{-1, 1\}$

• Therefore, we can define:

$\begin{cases} t^{(i)} h_t(x^{(i)}) = 1, & i \in E^c \\ t^{(i)} h_t(x^{(i)}) = -1, & i \in E \end{cases}$ # correct classify
misclassify.

• [Case 1]: When $i \in E^c$, by defⁿ of w_i , we have:

$\begin{aligned} w_i &\leftarrow w_i \exp(-dt) \\ &\leftarrow w_i \exp(-\frac{1}{2} \log \frac{1-err_t}{err_t}) \quad \# \text{ By defⁿ of } d_t \\ &\leftarrow w_i \exp(\log \sqrt{\frac{err_t}{1-err_t}}) \\ &\leftarrow w_i \sqrt{\frac{err_t}{1-err_t}} \end{aligned}$

• [Case 2]: When $i \in E$, by defⁿ of w_i , we have:

$\begin{aligned} w_i &\leftarrow w_i \exp(dt) \\ &\leftarrow w_i \exp(\frac{1}{2} \log \frac{1-err_t}{err_t}) \\ &\leftarrow w_i \exp(\log \sqrt{\frac{1-err_t}{err_t}}) \\ &\leftarrow w_i \sqrt{\frac{1-err_t}{err_t}} \end{aligned}$

• Therefore, $err'_t = \frac{\sum_{i=1}^N w_i I\{h_t(x^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^N w_i}$

$$\begin{aligned} &= \frac{\sum_{i \in E} w_i \sqrt{\frac{1-err_t}{err_t}}}{\sum_{i \in E} w_i \sqrt{\frac{1-err_t}{err_t}} + \sum_{i \in E^c} w_i \sqrt{\frac{err_t}{1-err_t}}} \\ &= \frac{(\frac{1-err_t}{err_t})^{1/2} \sum_{i \in E} w_i}{(\frac{1-err_t}{err_t})^{1/2} \sum_{i \in E} w_i + (\frac{err_t}{1-err_t})^{1/2} \sum_{i \in E^c} w_i} \quad (*) \\ &= \frac{\sum_{i \in E} w_i}{\sum_{i \in E} w_i + (\frac{err_t}{1-err_t}) \sum_{i \in E^c} w_i} \quad (**) \end{aligned}$$

• We also know that:

$err_t = \frac{\sum_{i \in E} w_i}{\sum_{i \in E} w_i + \sum_{i \in E^c} w_i}$

$\Rightarrow \sum_{i \in E} w_i = err_t (\sum_{i \in E} w_i + \sum_{i \in E^c} w_i)$

$\Rightarrow \sum_{i \in E^c} w_i = (\frac{1-err_t}{err_t}) \sum_{i \in E} w_i \quad (***)$

• Plug (***) into (**), we have:

$err'_t = \frac{\sum_{i \in E} w_i}{\sum_{i \in E} w_i + (\frac{err_t}{1-err_t}) (\frac{1-err_t}{err_t}) \sum_{i \in E} w_i} = \frac{\sum_{i \in E} w_i}{\sum_{i \in E} w_i + \sum_{i \in E} w_i} = \frac{1}{2}$

Interpretation:

In the class, we know that a single weak classifier is not capable of making the training error very small. It only performs slightly better than chance, i.e.: the error of classifier h according to the given weights $\{w^{(1)}, \dots, w^{(N)}\}$ is at most $\frac{1}{2} - \gamma$ for some $\gamma > 0$.

After the t^{th} iteration, reweight all the weights and apply the old weak learner, now we get the error rate proven above: $\text{err}_t' = \frac{1}{2}$, which implies that the error rate is now at its maximum and therefore, we cannot learn anything new with the old weak learner. This result forces us to use a new weak learner in the next iteration in order to decrease the error rate.