

# CSC411 Fall 2018: Homework 2

Niloufar Afsariardchi

March 9, 2020

## 1 Problem 1

a) We know  $0 \leq p(x) \leq 1$ , therefore  $1 \leq (1/p(x))$ , hence  $0 \leq \log_2(1/p(x))$ . Since  $p(x)$  is non-negative and  $\log_2(1/p(x))$  is also non-negative, the multiplication of two term is also non-negative, we therefore have  $0 \leq p(x) \log_2(1/p(x))$ . Since for each  $x$ , the term  $p(x) \log_2(1/p(x))$  is non-negative, the summation over all possible  $x$  is also non-negative and  $0 \leq \sum_x p(x) \log_2(1/p(x))$

b) Let's define  $\phi(x) = -\log(x)$ , note that since  $\log(x)$  is a concave function,  $-\log(x)$  is convex. Let's define a random variable  $U(X) = \frac{q(X)}{p(X)}$ , then from Jensen's inequality we know:

$$\phi(E_X[U(X)]) \leq E_X[\phi(U(X))] \quad (1)$$

$$-\log_2 \left( \sum_x p(x) \frac{q(X)}{p(X)} \right) \leq - \sum_x p(x) \log_2 \left( \frac{q(X)}{p(X)} \right) \quad (2)$$

$$-\log_2 \left( \sum_x q(x) \right) \leq KL(p||q) \quad (3)$$

$$-\log_2 1 \leq KL(p||q) \quad (4)$$

$$0 \leq KL(p||q) \quad (5)$$

where in equation (1) we used Jensen's inequality for convex function  $\phi$ , equation (2) is from the definition of the expectation for discrete distribution, and in equation (3) we know the sum of the probabilities over all outcomes is 1.

c) From definition we know:

$$KL(p(x, y)||p(x)p(y)) = \sum_{x,y} p(x, y) \log_2 \frac{p(x, y)}{p(x)p(y)} \quad (6)$$

$$= \sum_{x,y} p(x, y) \log_2 p(x, y) - \sum_{x,y} p(x, y) \log_2 p(x) - \quad (7)$$

$$- \sum_{x,y} p(x, y) \log_2 p(y) \quad (8)$$

$$= -H(X, Y) + H(X) + H(Y) \quad (9)$$

$$= H(Y) - H(Y|X) \quad (10)$$

$$= H(X) - H(X|Y) \quad (11)$$

$$= IG(X; Y) \quad (12)$$

where equation (10) is the result of the chain rule discussed in the class.

## 2 Problem 2

a) Define convex function  $\phi(y) = 0.5(y - t)^2$  and random variable  $\hat{y} \in \hat{Y} = \{y_1, y_2, \dots, y_m\}$ . From Jansen's inequality we have:

$$\phi(E_{\hat{Y}}[\hat{Y}]) \leq E_{\hat{Y}}[\phi(U(\hat{y}))] \quad (13)$$

$$0.5\left(t - \sum_{y_i \in \hat{Y}} p(y_i) y_i\right)^2 \leq \sum_{y_i \in \hat{Y}} p(y_i) 0.5(t - y_i)^2 \quad (14)$$

$$0.5\left(t - \frac{1}{m} \sum_{y_i \in \hat{Y}} y_i\right)^2 \leq \frac{1}{m} \sum_{y_i \in \hat{Y}} 0.5(t - y_i)^2 \quad (15)$$

$$L(t, h(\bar{x})) \leq \frac{1}{m} \sum_{i=1}^m L(t, h_i(x)) \quad (16)$$

note that I assume uniform distribution for random variable  $\hat{Y}$ . There for  $p(y_i) = 1/m$  in equation (15).

b) Let's define the set  $E = \{i : h_t(x^{(i)}) \neq t^{(i)}\}$ . Therefore,

$$\text{err}'_t = \frac{\sum_{i \in E} w'_i}{\sum_{i=1}^N w'_i} \quad (17)$$

$$= \frac{\sum_{i \in E} w_i \exp(-\alpha_t h(x^{(i)}) t^{(i)})}{\sum_{i \in E} w_i \exp(-\alpha_t h(x^{(i)}) t^{(i)}) + \sum_{i \in E^c} w_i \exp(-\alpha_t h(x^{(i)}) t^{(i)})} \quad (18)$$

$$= \frac{\sum_{i \in E} w_i \exp(\alpha_t)}{\sum_{i \in E} w_i \exp(\alpha_t) + \sum_{i \in E^c} w_i \exp(-\alpha_t)} \quad (19)$$

$$= \frac{\sum_{i \in E} w_i \exp(0.5 \log \frac{1 - \text{err}_t}{\text{err}_t})}{\sum_{i \in E} w_i \exp(0.5 \log \frac{1 - \text{err}_t}{\text{err}_t}) + \sum_{i \in E^c} w_i \exp(-0.5 \log \frac{1 - \text{err}_t}{\text{err}_t})} \quad (20)$$

$$= \frac{\sum_{i \in E} w_i \left(\frac{1 - \text{err}_t}{\text{err}_t}\right)^{0.5}}{\sum_{i \in E} w_i \left(\frac{1 - \text{err}_t}{\text{err}_t}\right)^{0.5} + \sum_{i \in E^c} w_i \left(\frac{1 - \text{err}_t}{\text{err}_t}\right)^{-0.5}} \quad (21)$$

$$= \frac{\sum_{i \in E} w_i}{\sum_{i \in E} w_i + \sum_{i \in E^c} w_i \left(\frac{\text{err}_t}{1 - \text{err}_t}\right)} \quad (22)$$

$$= \frac{1}{1 + \frac{\sum_{i \in E^c} w_i}{\sum_{i \in E} w_i} \left(\frac{\text{err}_t}{1 - \text{err}_t}\right)} \quad (23)$$

$$= \frac{1}{1 + \left(\frac{\sum_{i \in E^c} w_i}{\sum_{i \in E} w_i} + 1 - 1\right) \left(\frac{\text{err}_t}{1 - \text{err}_t}\right)} \quad (24)$$

$$= \frac{1}{1 + \left(\frac{\sum_{i=1}^N w_i}{\sum_{i \in E} w_i} - 1\right) \left(\frac{\text{err}_t}{1 - \text{err}_t}\right)} \quad (25)$$

$$= \frac{1}{1 + \left(\frac{1}{\text{err}_t} - 1\right) \left(\frac{\text{err}_t}{1 - \text{err}_t}\right)} \quad (26)$$

$$= \frac{1}{2} \quad (27)$$

equation (19) is the result of  $h(x^{(i)})t^{(i)} = -1$  for false and  $h(x^{(i)})t^{(i)} = 1$  for correct estimation.

The interpretation of this result is that using new weights the old weak learner is as good as a random guess. Therefore, with the new weights we need to update the weak learner, and we should do a better job at estimation than the previous round so that the error drop below  $1/2$ . This highlights the need to increase the complexity of our estimator at each iteration (e.g., by having a new weak learner added to the ensemble of the weak learners from previous iterations), so that bias decreases over iteration.