

## CSC311 HOMEWORK 1

Q1.

(a)

$$E(Z) = E(|X - Y|^2) = E(X^2 + Y^2 - 2XY)$$

$$= E(X^2) + E(Y^2) - 2E(XY) \text{ by linearity of expectation}$$

$$= \text{Var}(X) + E(X)^2 + \text{Var}(Y) + E(Y)^2 - 2E(XY) \text{ since } \text{Var}(X) = E(X^2) - E(X)^2$$

$$= \text{Var}(X) + E(X)^2 + \text{Var}(Y) + E(Y)^2 - 2E(X)E(Y) \text{ since } X, Y \text{ are independent}$$

Since X and Y are two independent univariate random variables sampled uniformly from the unit interval [0, 1], the expectation of X and Y is  $E(X) = E(Y) = \frac{1}{2}$ ,  $\text{Var}(X) = \text{Var}(Y) = \frac{1}{12}$

$$E(X^2) = E(Y^2) = \int_0^1 x^2 dx = \frac{1}{3}$$

$$E(X^3) = E(Y^3) = \int_0^1 x^3 dx = \frac{1}{4}$$

$$E(X^4) = E(Y^4) = \int_0^1 x^4 dx = \frac{1}{5}$$

$$E(Z) = \text{Var}(X) + E(X)^2 + \text{Var}(Y) + E(Y)^2 - 2E(X)E(Y)$$

$$= \frac{1}{12} + \frac{1}{4} + \frac{1}{12} + \frac{1}{4} - 2 \cdot \frac{1}{4}$$

$$= \frac{1}{6}$$

$$\text{Var}(Z) = E(Z^2) - E(Z)^2$$

$$= E((X - Y)^4) - \frac{1}{6^2}$$

$$= E(X^4 - 4X^3Y + 6X^2Y^2 - 4XY^3 + Y^4) - \frac{1}{36}$$

by linearity of expectation and the fact that X and Y are independent

$$= E(X^4) - 4E(X^3)E(Y) + 6E(X^2)E(Y^2) - 4E(X)E(Y^3) + E(Y^4) - \frac{1}{36}$$

$$= \frac{1}{5} - 4 \cdot \frac{1}{4} \cdot \frac{1}{2} + 6 \cdot \frac{1}{3} \cdot \frac{1}{3} - 4 \cdot \frac{1}{2} \cdot \frac{1}{4} + \frac{1}{5} - \frac{1}{36}$$

$$= \frac{7}{180}$$

Therefore, the expectation of Z is  $\frac{1}{6}$  and the variance of Z is  $\frac{7}{180}$ .

(b)

$$E(R) = E(Z_1 + Z_2 + \dots + Z_d)$$

$$= E(Z_1) + E(Z_2) + \dots + E(Z_d)$$

$$= E(|X_1 - Y_1|^2) + E(|X_2 - Y_2|^2) + \dots + E(|X_d - Y_d|^2)$$

$$\begin{aligned}
\text{Var}(R) &= \text{Var}(Z_1 + Z_2 + \dots + Z_d) \\
&= \text{Var}(Z_1) + \dots + \text{Var}(Z_d) \\
&= \text{Var}(|X_1 - Y_1|^2) + \text{Var}(|X_2 - Y_2|^2) + \dots + \text{Var}(|X_d - Y_d|^2)
\end{aligned}$$

Since all random variables  $X_1, \dots, X_d$  and  $Y_1, \dots, Y_d$  are independently and uniformly form  $[0, 1]$ , using the answer from part (a), we get

$$\begin{aligned}
E(R) &= E(|X_1 - Y_1|^2) + E(|X_2 - Y_2|^2) + \dots + E(|X_d - Y_d|^2) \\
&= dE(Z) \\
&= \frac{d}{6}
\end{aligned}$$

$$\begin{aligned}
\text{Var}(R) &= \text{Var}(|X_1 - Y_1|^2) + \text{Var}(|X_2 - Y_2|^2) + \dots + \text{Var}(|X_d - Y_d|^2) \\
&= d\text{Var}(Z) \\
&= \frac{7d}{180}
\end{aligned}$$

(c)

Let MED be the maximum possible squared Euclidean distance between two points within the  $d$ -dimensional unit cube (i.e. the squared Euclidean distance between opposite corners of the cube  $(0, 0, \dots, 0)$  and  $(1, 1, \dots, 1)$ ).

$$\text{MED} = (1 - 0)^2 + (1 - 0)^2 + (1 - 0)^2 + \dots + (1 - 0)^2 = d$$

Compare  $E(R)$ ,  $\text{SD}(R)$  to MED, as dimension goes to infinity:

$$\lim_{d \rightarrow \infty} E(R) = \lim_{d \rightarrow \infty} \frac{d}{6} \rightarrow \infty$$

And the standard deviation relative to MED when dimension goes to infinity is

$$\lim_{d \rightarrow \infty} \frac{\sqrt{\text{Var}(R)}}{d} = \lim_{d \rightarrow \infty} \sqrt{\frac{7}{180} \frac{\sqrt{d}}{d}} \rightarrow 0$$

Therefore, in high dimensions, most points are far away, and approximately the same distance.

Q2.

(a)

Since  $p(x)$  is a probability mass function,  $0 \leq p(x) \leq 1$ .

Then  $\frac{1}{p(x)} \geq 1$  and by the property of log function,  $\log_2 \left( \frac{1}{p(x)} \right) \geq 0$ .

So  $p(x) \log_2 \left( \frac{1}{p(x)} \right) \geq 0$  for all possible  $x$ .

Since  $H(X) = \sum_x p(x) \log_2 \left( \frac{1}{p(x)} \right)$  where  $x$  is all possible values,  $H(X) \geq 0$ .

Hence  $H(X)$  is non-negative.

(b)

Since  $X, Y$  are independent random variables,  $p(x, y) = p(x)p(y)$

Then

$$H(X, Y) = \sum_x \sum_y p(x, y) \log_2 \left( \frac{1}{p(x, y)} \right)$$

$$\begin{aligned}
&= \sum_x \sum_y p(x)p(y) \log_2 \left( \frac{1}{p(x)p(y)} \right) \\
&= \sum_x \sum_y p(x)p(y) \left( \log_2 \left( \frac{1}{p(x)} \right) + \log_2 \left( \frac{1}{p(y)} \right) \right) \\
&= \sum_x \sum_y p(x)p(y) \log_2 \left( \frac{1}{p(x)} \right) + \sum_x \sum_y p(x)p(y) \log_2 \left( \frac{1}{p(y)} \right) \\
&= \sum_y p(y) \sum_x p(x) \log_2 \left( \frac{1}{p(x)} \right) + \sum_x p(x) \sum_y p(y) \log_2 \left( \frac{1}{p(y)} \right)
\end{aligned}$$

Since  $\sum_x p(x) = p(\chi) = 1 = \sum_y p(y)$

$$\begin{aligned}
&= \sum_x p(x) \log_2 \left( \frac{1}{p(x)} \right) + \sum_y p(y) \log_2 \left( \frac{1}{p(y)} \right) \\
&= H(X) + H(Y)
\end{aligned}$$

(c)

By the definition of conditional entropy,

$$\begin{aligned}
H(Y|X) &= - \sum_x \sum_y p(x,y) \log_2 p(y|x) = - \sum_x \sum_y p(x,y) \log_2 \frac{p(x,y)}{p(x)} \\
&= \sum_x \sum_y p(x,y) \log_2 \left( \frac{p(x,y)}{p(x)} \right)^{-1} = \sum_x \sum_y p(x,y) \log_2 \left( \frac{p(x)}{p(x,y)} \right) \\
&= \sum_x \sum_y p(x,y) \log_2(p(x)) - \sum_x \sum_y p(x,y) \log_2(p(x,y)) \\
&= \sum_x p(x) \log_2(p(x)) + \sum_x \sum_y p(x,y) \log_2 \left( \frac{1}{p(x,y)} \right) \text{ since } \sum_y p(x,y) = p(x) \\
&= \sum_x \sum_y p(x,y) \log_2 \left( \frac{1}{p(x,y)} \right) - \sum_x p(x) \log_2 \left( \frac{1}{p(x)} \right) \\
&= H(X,Y) - H(X)
\end{aligned}$$

Hence,  $H(X,Y) = H(X) + H(Y|X)$

(d)

From the question, we know  $KL(p||q) = \sum_x p(x) \log_2 \frac{p(x)}{q(x)} = - \sum_x p(x) \log_2 \frac{q(x)}{p(x)}$

To prove  $KL(p||q) \geq 0$ ,

we want to prove  $\sum_x p(x) \log_2 \frac{q(x)}{p(x)} \leq 0$  so that  $- \sum_x p(x) \log_2 \frac{q(x)}{p(x)} \geq 0$ .

Since  $-\log_2$  is a concave function, using Jensen's inequality  $\mathbb{E}[\phi(X)] \leq \phi(\mathbb{E}[X])$ , we get

$$- \sum_x p(x) \log_2 \frac{q(x)}{p(x)} = \sum_x p(x) \left( -\log_2 \frac{q(x)}{p(x)} \right)$$

$$\begin{aligned}
&= E_{X \sim p} \left[ -\log_2 \frac{q(x)}{p(x)} \right] \geq -\log_2 E_{X \sim p} \left[ \frac{q(x)}{p(x)} \right] \\
&= -\log_2 \sum_x p(x) \frac{q(x)}{p(x)} = -\log_2 \sum_x q(x)
\end{aligned}$$

Since  $q$  is a probability distribution,  $\sum_x q(x) = 1 \rightarrow -\log_2 \sum_x q(x) = 0$

So  $-\sum_x p(x) \log_2 \frac{q(x)}{p(x)} \geq -\log_2 \sum_x q(x) = 0$

Therefore  $KL(p||q) = -\sum_x p(x) \log_2 \frac{q(x)}{p(x)} \geq 0$  i.e.  $KL(p||q)$  is non-negative.

(e)

By the definition of KL divergence,

$$KL(p(x,y)||p(x)p(y)) = \sum_x \sum_y p(x,y) \log_2 \frac{p(x,y)}{p(x)p(y)}$$

By the definition of information gain and entropy

$$\begin{aligned}
I(Y; X) &= H(Y) - H(Y|X) = \sum_y p(y) \log_2 \left( \frac{1}{p(y)} \right) + \sum_x \sum_y p(x,y) \log_2 \frac{p(x,y)}{p(x)} \\
&= \sum_x \sum_y p(x,y) \log_2 \left( \frac{1}{p(y)} \right) + \sum_x \sum_y p(x,y) \log_2 \frac{p(x,y)}{p(x)} \text{ since } \sum_x p(x,y) = p(y) \\
&= \sum_x \sum_y p(x,y) \log_2 \left( \frac{p(x,y)}{p(x)p(y)} \right)
\end{aligned}$$

Hence,

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

3.

(b)

Accuracy for depth 16 and criteria entropy : 0.7346938775510204

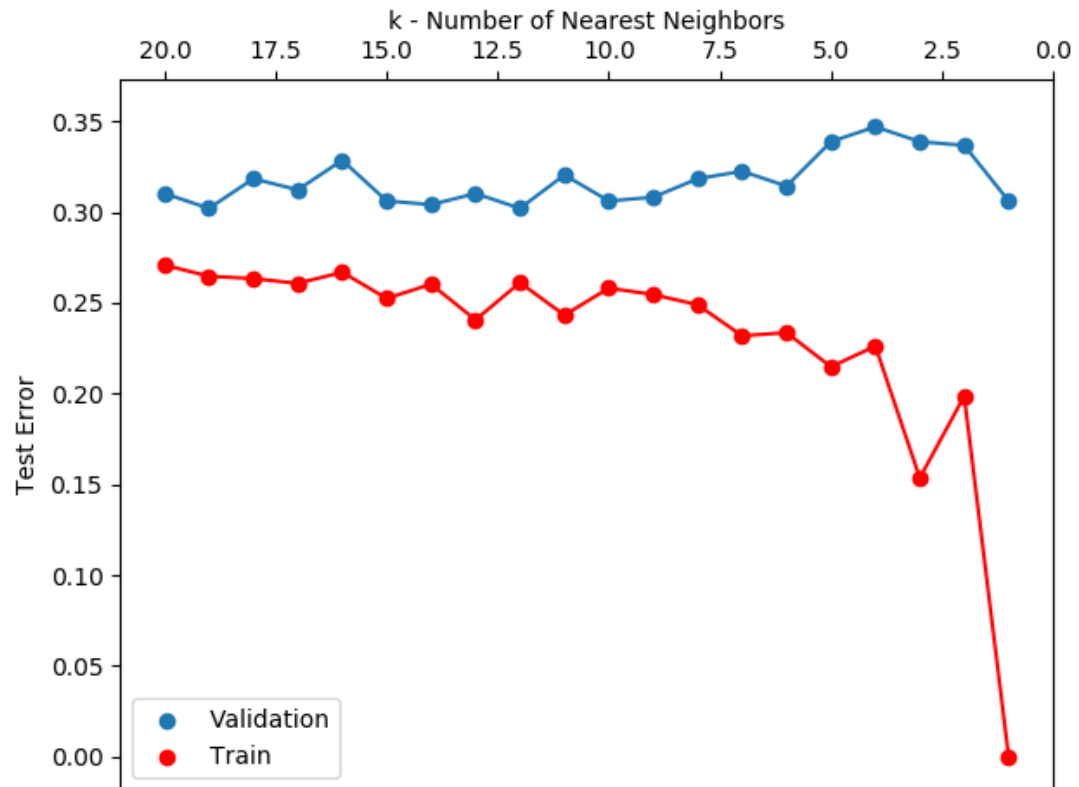
Accuracy for depth 32 and criteria entropy : 0.753061224489796



The information gain of 'good': 0.0004942945179704354

(e)

The relationship between number of nearest neighbours and test error



Accuracy of the best KNN model on the test dataset: 0.6551020408163265