

# COMP2610/6261 Assignment 1 2018 Solutions

## 1 Entropy and Mutual Information

1.
  - a [2 points]  $P(L = A) = \frac{4}{16}$  giving  $h(L = A) = 2$
  - b [2 points]  $P(V = 3) = \frac{4}{16}$  giving  $h(V = 3) = 2$
  - c [2 points]  $P(L = A|V = 1) = \frac{4}{8}$  giving  $h(L = A|V = 1) = 1$
  - d [2 points]  $P(V = 1|L = A) = 1$  giving  $h(V = 1|L = A) = 0$ , Knowing the letter determines the value hence there is no surprise in discovering the value.
  - e [2 points]  $H(L) = \frac{11}{4}$ ,  $H(V) = \frac{3}{2}$  and  $H(D) \approx 0.811$ .

2. [5 points]  $I(L; V) = H(V) - H(V|L) = H(V)$  as knowing the letter determines the value of the code meaning  $H(V|L) = 0$ . Hence  $I(L; V) = \frac{3}{2}$

$I(D; V) = H(V) - H(V|D)$  and we know that  $H(V) = \frac{3}{2}$ . If  $D = \text{decaf}$  then we have the following probability distribution for  $V$

$$P(V = 1|D = \text{decaf}) = \frac{8}{12} \quad P(V = 3|D = \text{decaf}) = \frac{3}{12} \quad P(V = 5|D = \text{decaf}) = \frac{1}{12}$$

Giving  $H(V|D = \text{decaf}) \approx 1.189$ . Similarly when  $D = \text{not decaf}$  then we obtain the following distribution for  $V$

$$P(V = 1|D = \text{not decaf}) = 0 \quad P(V = 3|D = \text{not decaf}) = \frac{1}{4} \quad P(V = 5|D = \text{not decaf}) = \frac{3}{4}$$

Giving  $H(V|D = \text{not decaf}) \approx 0.811$ . As  $P(D = \text{decaf}) = \frac{3}{4}$ ,  $H(V|D) \approx 1.094$  meaning  $I(V; D) \approx 0.406$ .

For the interpretation, note that  $P(V, D, L) = P(V)P(L|V)P(D|L, V) = P(V)P(L|V)P(D|L)$  as a tiles decaf status is determined by its letter. Hence we have a markov chain  $V \rightarrow L \rightarrow D$  so by the information processing theorem  $I(V; D) \leq I(V : L)$  as the above calculations have confirmed.

3.
  - a [2 points]  $P(D = \text{decaf}|C = \text{cheat}) = \frac{1}{2}$  and  $P(D = \text{decaf}|C = \text{not cheat}) = \frac{3}{4}$  giving  $H(D|C) \approx \frac{1}{2}(1 + 0.811) = 0.906$
  - b [3 points]  $H(C, D) = H(C) + H(D|C)$ .  $H(C) = 1$  giving  $H(C, D) \approx 1.905$
  - 3 [5 points]  $P(3 \text{ decaf tiles}|\text{cheat}) = (\frac{1}{2})^3$  and  $P(3 \text{ decaf tiles}|\text{not cheat}) = (\frac{3}{4})^3$ . Hence by Bayes rule

$$\begin{aligned} P(\text{cheat}|3 \text{ decaf tiles}) &= \frac{P(3 \text{ decaf tiles}|\text{cheat})P(\text{cheat})}{P(3 \text{ decaf tiles})} \\ &= \frac{\frac{1}{8}}{\frac{1}{8} + \frac{27}{64}} \\ &= \frac{8}{35} \end{aligned}$$

## 2 Mutual Information

1. We have

$$\begin{aligned} p(X = x) &= p(X = x|Y = 0)p(Y = 0) + p(X = x|Y = 1)p(Y = 1) \\ &= \frac{p(X = x|Y = 0) + p(X = x|Y = 1)}{2}. \end{aligned}$$

Let

$$\begin{aligned} \mathbf{m} &= \left( \frac{p(X = \mathbf{a}|Y = 0) + p(X = \mathbf{a}|Y = 1)}{2}, \frac{p(X = \mathbf{a}|Y = 0) + p(X = \mathbf{a}|Y = 1)}{2}, \right. \\ &\quad \left. \frac{p(X = \mathbf{a}|Y = 0) + p(X = \mathbf{a}|Y = 1)}{2} \right) \\ &= \frac{\mathbf{p} + \mathbf{q}}{2} \\ &= \left( \frac{3}{8}, \frac{3}{8}, \frac{1}{4} \right). \end{aligned}$$

Now,

$$\begin{aligned} I(X; Y) &= D_{\text{KL}}(p(X, Y) \| p(X)p(Y)) \\ &= \sum_{x,y} -p(x, y) \cdot \log \frac{p(x, y)}{p(x)p(y)} \\ &= \sum_{x,y} -p(x, y) \cdot \log \frac{p(x|y)}{p(x)} \\ &= \sum_{x,y} -p(x|y) \cdot p(y) \cdot \log \frac{p(x|y)}{p(x)} \\ &= \sum_x -p(x|y = 0) \cdot \log \frac{p(x|y = 0)}{p(x)} + \sum_x -p(x|y = 1) \cdot \log \frac{p(x|y = 1)}{p(x)} \\ &= \frac{D_{\text{KL}}(\mathbf{p} \| \mathbf{m}) + D_{\text{KL}}(\mathbf{q} \| \mathbf{m})}{2}. \end{aligned}$$

It can be checked that

$$\begin{aligned}
D_{\text{KL}}(\mathbf{p} \parallel \mathbf{m}) &= - \sum_x p(x) \log \frac{m(x)}{p(x)} \\
&= \sum_x p(x) \log \frac{p(x)}{m(x)} \\
&= \frac{1}{2} \log \frac{4}{3} + \frac{1}{4} \log \frac{2}{3} + \frac{1}{4} \log 1 \\
&= 1 + \frac{1}{4} - \frac{3}{4} \log 3 \\
&\approx 0.0613.
\end{aligned}$$

$$\begin{aligned}
D_{\text{KL}}(\mathbf{q} \parallel \mathbf{m}) &= - \sum_x q(x) \log \frac{m(x)}{q(x)} \\
&= \sum_x q(x) \log \frac{q(x)}{m(x)} \\
&= \frac{1}{4} \log \frac{2}{3} + \frac{1}{2} \log \frac{4}{3} + \frac{1}{4} \log 1 \\
&\approx 0.0613.
\end{aligned}$$

Thus,

$$I(X; Y) \approx 0.0613.$$

2. We have established this above.

3. Let

$$\begin{aligned}
\mathbf{p}' &= (p(Z = \mathbf{a} | Y = 1), p(Z = \mathbf{b} | Y = 1), p(Z = \mathbf{c} | Y = 1)). \\
\mathbf{q}' &= (p(Z = \mathbf{a} | Y = 0), p(Z = \mathbf{b} | Y = 0), p(Z = \mathbf{c} | Y = 0)). \\
\mathbf{m}' &= \frac{\mathbf{p}' + \mathbf{q}'}{2}.
\end{aligned}$$

Then, from (ii),

$$I(Z; Y) = \frac{D_{\text{KL}}(\mathbf{p}' \parallel \mathbf{m}') + D_{\text{KL}}(\mathbf{q}' \parallel \mathbf{m}')}{2}.$$

From the problem statement,

$$\mathbf{p}' = \mathbf{q}$$

$$\mathbf{q}' = \mathbf{p}$$

and so

$$\mathbf{m}' = \mathbf{m}.$$

It thus follows that

$$I(Z; Y) = \frac{D_{\text{KL}}(\mathbf{q} \parallel \mathbf{m}) + D_{\text{KL}}(\mathbf{p} \parallel \mathbf{m})}{2} = I(X; Y).$$

This is intuitive, because mutual information measures the average reduction in uncertainty in  $Z$  when  $Y$  is known.  $Z$  can be seen as merely “relabelling” the outcomes of  $X$ , while sharing the probabilities. Therefore, the uncertainty in both random variables is the same, as is the mutual information with respect to  $Y$ .

4. We can construct a Markov chain  $Y \rightarrow X \rightarrow Z$ , which would imply that  $I(X; Y) \geq I(Z; Y)$ . To enforce strict equality, we must ensure that  $Z$  somehow discards information from  $X$ . For example, let  $Z$  be uniformly from  $\{\mathbf{a}, \mathbf{b}\}$  if  $X = \mathbf{a}$ , uniformly from  $\{\mathbf{a}, \mathbf{c}\}$  if  $X = \mathbf{b}$ , and uniformly from  $\{\mathbf{b}, \mathbf{c}\}$  if  $X = \mathbf{c}$ . Then, we have that

$$p(Z = \mathbf{a}|Y = y) = \frac{p(X = \mathbf{a}|Y = y) + p(Z = \mathbf{b}|Y = y)}{2}$$

$$p(Z = \mathbf{b}|Y = y) = \frac{p(X = \mathbf{a}|Y = y) + p(Z = \mathbf{c}|Y = y)}{2}$$

$$p(Z = \mathbf{c}|Y = y) = \frac{p(X = \mathbf{b}|Y = y) + p(Z = \mathbf{c}|Y = y)}{2},$$

so that

$$\mathbf{p}' = (3/8, 3/8, 1/4)$$

$$\mathbf{q}' = (3/8, 1/4, 3/8).$$

It can now be checked, using the formula from part (ii), that

$$I(Z; Y) = 0.0182 < I(X; Y).$$

5. We can construct a Markov chain  $Y \rightarrow Z \rightarrow X$ , where now  $X$  is a noisy version of  $Z$ . For example, let  $Z$  have conditional probability vectors

$$\mathbf{p}' = (1/2, 1/2, 0)$$

$$\mathbf{q}' = (1/2, 0, 1/2)$$

and say that  $X$  is uniform from  $\{\mathbf{a}, \mathbf{b}\}$  if  $Z = \mathbf{a}$ , uniform from  $\{\mathbf{a}, \mathbf{c}\}$  if  $Z = \mathbf{b}$ , and uniform from  $\{\mathbf{b}, \mathbf{c}\}$  if  $Z = \mathbf{c}$ . It may be verified that this results in the same  $\mathbf{p}, \mathbf{q}$  for  $X$  as in part (i); it may further be seen that  $I(Z; Y) = +\infty > I(X; Y)$ , with the infinite mutual information owing to the KL divergence blowing up due to the presence of zero entries in  $\mathbf{p}', \mathbf{q}'$ .

### 3 Uniform Length, Lossy Coding

Denote the albums by their first letters and let  $\mathcal{A} = \{A, B, C, D\}$ .

1. Bits required to uniformly code:

(a) Two bits ( $= \log_2 4$ ). Example uniform code:  $C = \{00, 01, 10, 11\}$ .

(b) The album Alina has 5 tracks so  $\lceil \log_2 5 \rceil = 3$  bits are required. Example code:  $C = \{000, 001, 010, 011, 100\}$ .

- (c) There are 44 tracks in total so  $\lceil \log_2 44 \rceil = 6$  bits are required.
2. The raw bit content is  $\log_2 44 \approx 5.46$ .
3. Let  $A$  denote ensemble of album chosen when tracks across all albums are picked uniformly at random.
- (a)  $\mathcal{A}_A = \{A, B, C, D\}$  and  $p = \{5/44, 12/44, 15/44, 12/44\}$ .
- (b) Raw bit content  $= \log_2 |\mathcal{A}_{A^4}| = \log_2 4^4 = 8$ .
- (c) Since  $A^4$  has 256 elements we need a  $\delta$  equal to the smallest probability for element of  $A^4$ . As album  $A$  has smallest probability, so too will  $AAAA \in \mathcal{A}_A$ . Thus choose  $\delta = P(AAAA) = (5/44)^4 \approx 0.0002$ .
- (d)  $H_\delta(A^4)$  will be zero when  $S_\delta$  contains only one element so if  $\delta$  is set so that there are two elements in  $S_\delta$  we are done. Since Coieda has the highest probability the sequence  $CCCC \in \mathcal{A}_{A^4}$  will have highest probability  $P(CCCC) = (15/44)^4$ . The next largest probability is any sequence with three  $C$ s and either  $B$  or  $D$  in it.  $P(BCCC) = (12/44)(15/44)^3$  so set  $\delta = (12/44)(15/44)^3 + (15/44)^4 \approx 0.024$ .
4. (a)  $H(A) = \frac{1}{44} (5 \log_2(5/44) + 2 \times 12 \log_2(12/44) + 15 \log_2(15/44)) \approx 1.91$
- (b) Items in typical set have  $P(a)$
- (c) The probability of sequences of length  $N = 100$  in  $T_{N\beta}$  for  $\beta = 0.1$  is no more than  $2^{-100(H(A)-0.1)} \approx 3.26 \times 10^{-55}$  and no less than  $2^{-100(H(A)+0.1)} \approx 3.11 \times 10^{-61}$ . This means  $\frac{1}{3.26 \times 10^{-55}} \approx 3.1 \times 10^{54} \leq |T_{N\beta}| \leq \frac{1}{3.11 \times 10^{-61}} \approx 3.2 \times 10^{60}$  and so there are approximately  $10^{57}$  elements.
- (d) No, because that rate is below the entropy of 1.91 so, by the source coding theorem, for large blocks a rate of 1.5 bits per title is not possible.

Q4 (c) (replace ~~2~~). part (c) by setting  $c = \frac{1}{2}$ .

We need to compute  $EL_X(X-c) = \int_{-\infty}^{\infty} L_X(x-c) f(x) dx$  where  $f(\cdot)$  is the density of  $X$ .

Observe that  $L_X(x-c) = \begin{cases} x(x-c) & x \geq c \\ (x-1)(x-c) & x \leq c. \end{cases}$

$$\begin{aligned} \text{Thus } EL_X(X-c) &= \int_{-\infty}^c (x-1)(x-c) f(x) dx + \int_c^{\infty} x(x-c) f(x) dx \\ &= (x-1) \int_{-\infty}^c x f(x) dx - (x-1)c \int_{-\infty}^c f(x) dx \\ &\quad + x \int_c^{\infty} x f(x) dx - xc \int_c^{\infty} f(x) dx \\ &= (x-1) E_{-\infty}^c - (x-1)c F(c) + x E_c^{\infty} - xc(1-F(c)) \quad (*) \end{aligned}$$

where  $F(c)$  is the cumulative distribution of  $X$ , and  $E_a^b := \int_a^b x f(x) dx$ .

Observe that by Leibnitz's rule, (or just fundamental theorem of calculus)

$$\frac{\partial}{\partial c} E_{-\infty}^c = \frac{\partial}{\partial c} \int_{-\infty}^c \phi(x) dx = \phi(c)$$

where  $\phi(x) = x f(x)$

$$\text{and } \frac{\partial}{\partial c} E_c^{\infty} = \frac{\partial}{\partial c} \int_c^{\infty} \phi(x) dx = -\phi(c).$$

Then from (\*) we have

$$\begin{aligned} \frac{\partial}{\partial c} EL_X(X-c) &= (x-1) \phi(c) - (x-1)c f(c) - (x-1) F(c) \\ &\quad - xc f(c) + x c f(c) - x(1-F(c)) \\ &= c f(c) [(x-1) - (x-1) - x + x] - (x-1) F(c) - x(1-F(c)) \end{aligned}$$

Set to zero:

$$-x F(c) + F(c) - x + x F(c) = 0$$

$$\Rightarrow F(c) = x \Rightarrow c = \frac{1}{2}$$

4 b).

$$\begin{aligned}
 |\mu - m| &= E(X - m) \\
 &\leq E|X - m| \\
 &\leq E|X - \mu| \\
 &\leq \sqrt{E(X - \mu)^2} \\
 &= \sigma
 \end{aligned}$$

(Tenser)

(from 4a)

(Tenser)

(Def<sup>n</sup> of  $\sigma$ )

4(d).  $h_x(x) = 0$  only if  $x = 0$ .

Thus  $E h_x(X - c) = 0$  only if  $X = c$ .

Suppose  $Q_x(X) = \min_{c \in \mathbb{R}} E h_x(X - c) = q \neq 0$

Consider  $Y = \beta X$ , for some  $\beta > 0$ .

$$Q_x(Y) = Q_x(\beta X) = \min_c E h_x(\beta X - c)$$

$$\text{Now } h_x(\beta x) = \beta h_x(x)$$

The minimizing  $c$  value will differ, but we see that

$$Q_x(\beta X) = \beta Q_x(X)$$

$$\text{Confer } \sigma(\beta X) = \beta \sigma(X).$$

$\sigma^2$  does not always exist. But every random variable has quantiles.

(There is a subtlety about the distribution of  $X$  has atomic components, but the key point is that  $Q_x$  always exists.)

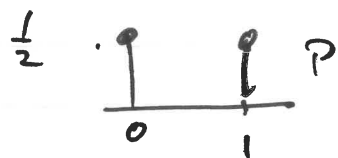
$Q_x(X)$  comes from to  $CVar_x(X - E(X))$

The conditional value at risk of  $X - E(X)$

- See [R. Tyrrell Rockafellar & Stan Uryasev, "The Fundamental Risk Quadrangle" in Risk Management, Optimization and Statistical Estimation, Surveys in Operations Research and Management Science 18(1-2), 33-53 (2013)]



4e. Consider two discrete distributions



$$H = \log_2 \frac{2}{1} = 1$$

equal mixture is



$$H = 4 \cdot \frac{1}{4} \log_2 4 = 2$$

$$\text{So } H(P) = H(Q) = 1$$

$$H\left(\frac{1}{2}P + \frac{1}{2}Q\right) = 2.$$

5). Assume w.l.o.g.  $\mu = E(X) = 0$ .  
 We assume  $X \in [a, b]$  (as assumed in part (d)).

Can write  $X = \alpha b + (1-\alpha)a$  where  
 $\alpha = \frac{(X-a)}{b-a}$ .

Since  $g(t) = e^{tx}$  is convex

$$e^{tX} \leq \alpha e^{tb} + (1-\alpha) e^{ta} \\
= \frac{X-a}{b-a} e^{tb} + \frac{b-X}{b-a} e^{ta}.$$

Take expectations of both sides to obtain:

$$E(e^{tX}) \leq -\frac{a}{b-a} e^{tb} + \frac{b}{b-a} e^{ta} = e^{g(t(b-a))}$$

where  $u = t(b-a) \in g(u) = -\delta u + \log(1-\delta + \delta e^u)$ ,  
 $\delta = \frac{-a}{b-a}$

6) Since  $E(X) = 0$ ,  $a \leq 0 \leq b$ .

$$\frac{\partial^2}{\partial u^2} g(u) = \frac{\delta e^u}{1-\delta + \delta e^u} - \frac{\delta^2 (e^u)^2}{(1-\delta + \delta e^u)^2} \\
= \frac{\delta e^u (1-\delta + \delta e^u) - \delta^2 (e^u)^2}{(1-\delta + \delta e^u)^2}$$

Denominator is always non-negative.

Numerator can be expanded as

$$\delta e^u - \delta^2 e^u + \delta^2 (e^u)^2 - \delta^2 (e^u)^2 \\
= \delta e^u (1-\delta). \text{ But } \delta = \frac{-a}{b-a} \leq \frac{-a}{a} = 1$$

Thus  $\frac{d^2}{du^2} f(u) \geq 0$  &  $f$  is convex.

Observe that  $\frac{d^2}{du^2} f(u) = \frac{\gamma e^u [(1-\gamma+\gamma e^u) - \gamma e^u]}{(1-\gamma+\gamma e^u)^2}$   
 $= \frac{\gamma e^u (1-\gamma)}{(1-\gamma+\gamma e^u)^2} \quad \gamma \in [0,1]$

One can check numerically that  $\frac{d^2}{du^2} f(u) \leq \frac{1}{4}$   
 $\forall u > 0$ .

By Taylor's theorem there is  $\xi \in (0, u)$  such that

$$g(u) = g(0) + u g'(0) + \frac{u^2}{2} g''(\xi) = \frac{u^2}{2} g''(\xi)$$

Since  $g(0) = g'(0) = 0$ . Hence

$$g(u) \leq \frac{u^2}{8} = \frac{t^2 (b-a)^2}{8}$$

Thus  $E e^{tx} \leq e^{g(u)} \leq e^{t^2 (b-a)^2 / 8}$ .

c)  $P(X > \varepsilon) = P(e^X > e^\varepsilon) = P(e^{tx} > e^{t\varepsilon})$   
for any  $t \geq 0$

But  $P(e^{tx} > e^{t\varepsilon}) \leq e^{-t\varepsilon} E(e^{tx})$ .

(  $P(Y > s) \leq \frac{1}{s} E(Y)$  ) <sup>by Markov's  $\neq$ .</sup>

This holds for any  $t > 0$ , & Thus

$$P(X > \varepsilon) \leq \inf_{t > 0} e^{-t\varepsilon} E(e^{tx})$$

d) Assume w.l.o.g. that  $\mu = E(X_i) = 0$ .

$$\begin{aligned} P(|\bar{X}_n| \geq \varepsilon) &= P(\bar{X}_n \geq \varepsilon) + P(\bar{X}_n \leq -\varepsilon) \\ &= P(\bar{X}_n \geq \varepsilon) + P(-\bar{X}_n \geq \varepsilon) \end{aligned}$$

Now using earlier results, we have

$$P(\bar{X}_n \geq \varepsilon) = P\left(\sum_{i=1}^n X_i \geq n\varepsilon\right) = P\left(e^{\sum_{i=1}^n t X_i} \geq e^{tn\varepsilon}\right)$$

$$= P\left(e^{t \sum_{i=1}^n X_i} \geq e^{tn\varepsilon}\right)$$

$$\leq e^{-tn\varepsilon} E\left(e^{t \sum_{i=1}^n X_i}\right)$$

$$= e^{-tn\varepsilon} E\left(\prod_{i=1}^n e^{t X_i}\right)$$

$$= e^{-tn\varepsilon} \prod_{i=1}^n E(e^{t X_i})$$

(by independence)

$$= e^{-tn\varepsilon} (E(e^{t X_i}))^n$$

(for any  $i \in [n]$ )

By part (b)  $E(e^{t X_i}) \leq e^{t^2(b-a)^2/8}$

Then  $P(\bar{X}_n \geq \varepsilon) \leq e^{-tn\varepsilon} e^{t^2 n(b-a)^2/8} \quad (*)$

This holds for all  $t$ , so we can optimize over  $t$

$$\frac{d}{dt} e^{-tn\varepsilon} \frac{t^2 n(b-a)^2}{8}$$

$$= -n\varepsilon e^{-tn\varepsilon} \frac{t^2 n(b-a)^2}{8} + e^{-tn\varepsilon} \frac{2t n(b-a)^2}{8}$$

Take  $\frac{d}{dt}$   $-tn\varepsilon + t^2 n(b-a)^2/8$   
(exp is constant).

$$\frac{d}{dt} (-tn\varepsilon + t^2 n(b-a)^2/8)$$

$$= -n\varepsilon + 2tn(b-a)^2/8$$

Set to zero & solve for  $t$

$$\frac{2tn(b-a)^2}{8} = n\varepsilon$$

$$\Rightarrow t = \frac{4\varepsilon}{(b-a)^2}$$

Substitute this value of  $t$ , so (\*) becomes

$$P(\bar{X}_n \geq \varepsilon) \leq e^{-\frac{4\varepsilon}{(b-a)^2} \cdot n\varepsilon} e^{+\frac{16\varepsilon^2}{(b-a)^4} \cdot \frac{(b-a)^2 n}{8}}$$

$$= e^{-\frac{4\varepsilon^2 n}{(b-a)^2}} \cdot e^{-\frac{2\varepsilon^2 n}{(b-a)^2}}$$

$$e^{-\frac{4n\varepsilon^2}{(b-a)^2}} \cdot e^{+\frac{2\varepsilon^2 n}{(b-a)^2}}$$

$$= e^{-2n\varepsilon^2/(b-a)^2}$$

$$= e$$

A similar argument holds for  $P(\bar{X}_n \leq -\varepsilon)$

So we obtain the result by  $P(A \text{ or } B) = P(A) + P(B)$   
given  $A, B$  mutually exclusive

5e). Bernoulli  $\Rightarrow (b-a) = 1$ .