

---

Problem Set 2 — *Due Friday, October 26, before class starts*  
For the Exercise Sessions on Oct 12 and 19

---

Last name	First name	SCIPER Nr	Points

**Problem 1: Information Measures for Continuous Random Variables**

*Recommended Reading: Chapter 8 of the book by T. M. Cover and J. A. Thomas. “Elements of Information Theory,” Second Edition, Wiley, 2006. It is available for free download from the EPFL library.*

Find the differential entropy  $h(X)$  for the case where  $X$  is an exponentially distributed random variable of mean  $1/\lambda$ .

**Solution**

Let  $X$  be an exponentially distributed random variable with mean  $\frac{1}{\lambda}$ . Then the pdf of  $X$  is

$$f(x) = \lambda e^{-\lambda x} \text{ for } x \geq 0 \quad (1)$$

The differential entropy of  $X$  is

$$h(X) = - \int_0^\infty f(x) \log f(x) dx \quad (2)$$

$$= - \int_0^\infty \lambda e^{-\lambda x} \log(\lambda e^{-\lambda x}) dx \quad (3)$$

$$= - \left( \int_0^\infty \lambda e^{-\lambda x} \log \lambda dx + \int_0^\infty \lambda e^{-\lambda x} (-\lambda x) dx \right) \quad (4)$$

$$= - \log \lambda \int_0^\infty \lambda e^{-\lambda x} dx + \int_0^\infty x \lambda e^{-\lambda x} dx \quad (5)$$

$$= - \log \lambda + 1 \quad (6)$$

where the last step is because  $\int_0^\infty \lambda e^{-\lambda x} dx = \int_0^\infty f(x) dx = 1$  and  $\int_0^\infty x \lambda e^{-\lambda x} dx = \mathbb{E}[X] = \frac{1}{\lambda}$ .

**Problem 2: MMSE Estimation**

Consider the scenario where  $p(x|d) = de^{-dx}$ , for  $x \geq 0$  (and zero otherwise), that is, the observed data  $x$  is distributed according to an exponential with mean  $1/d$ . Moreover, the desired variable  $d$  itself is also exponentially distributed, with mean  $1/\mu$ .

(a) Find the MMSE estimator of  $d$  given  $x$ , and calculate the corresponding mean-squared error incurred by this estimator.

(b) Find the MAP estimator of  $d$  given  $x$ .

**Solution**

(a) Since  $d$  is exponentially distributed random variable with mean  $1/\mu$ , we have

$$p(d) = \mu e^{-\mu d} \quad (7)$$

Then the probability  $p(x, d) = p(d)p(x|d) = \mu e^{-\mu d} d e^{-dx} = \mu d e^{-(\mu+x)d}$ . Thus, we have

$$p_X(x) = \int_d p(x, d) = \frac{\mu}{(\mu + x)^2} \quad (8)$$

Given  $x$ , the probability of  $p(d|x) = (\mu + x)^2 d e^{-(\mu+x)d}$ , which is Gamma distribution  $\Gamma(2, \mu + x)$ .

The MMSE estimator of  $d$  given  $x$ ,  $\hat{d}_{MMSE}(x)$ , satisfies

$$\hat{d}_{MMSE}(x) = E[d|X = x] = \frac{2}{\mu + x} \quad (9)$$

and the mean-squared error is

$$\mathcal{E} = E_D[(d - \hat{d}_{MMSE})^2] \quad (10)$$

$$= E_X[E_D[(d - \hat{d}_{MMSE}(x))^2|X = x]] \quad (11)$$

$$= \int E_D[(d - \hat{d}_{MMSE}(x))^2|X = x] p_X(x) dx \quad (12)$$

$$= \int E_D[d^2 - 2d\hat{d}_{MMSE}(x) + \hat{d}_{MMSE}^2(x)|X = x] p_X(x) dx \quad (13)$$

$$= \int (E_D[d^2|X = x] - 2\hat{d}_{MMSE}(x) E_D[d|X = x] + \hat{d}_{MMSE}^2(x)) p_X(x) dx \quad (14)$$

$$= \int \left( \frac{6}{(\mu + x)^2} - \frac{4}{(\mu + x)^2} \right) p_X(x) dx \quad (15)$$

$$= \int \frac{2\mu}{(\mu + x)^4} dx \quad (16)$$

$$= \frac{2}{3\mu^2} \quad (17)$$

where  $E_D[d\hat{d}_{MMSE}(x)|X = x] = \hat{d}_{MMSE}^2(x)$  and  $E_D[d^2|X = x] = \text{Var}(D|X) + E[D|X] = \frac{6}{(\mu+x)^2}$  is because  $p(d|x)$  is Gamma distribution.

(b) MAP estimator is

$$\hat{d}_{MAP}(x) = \arg \max_d p(d|x) \quad (18)$$

$$= \arg \max_d (\mu + x)^2 d e^{-(\mu+x)d} \quad (19)$$

$$= \arg \max_d d e^{-(\mu+x)d} \quad (20)$$

$$= \frac{1}{\mu + x} \quad (21)$$

as  $\frac{\partial}{\partial d} d e^{-(\mu+x)d} = 0$ , when  $d = \frac{1}{\mu+x}$ .

### Problem 3: Tweedie's Formula

For the special case where  $X = D + N$ , where  $N$  is Gaussian noise of mean zero and variance  $\sigma^2$ , *Tweedie's formula* says that the conditional mean (that is, the MMSE estimator) can be expressed as

$$\mathbb{E}[D|X = x] = x + \sigma^2 \ell'(x), \quad (22)$$

where

$$\ell'(x) = \frac{d}{dx} \log f_X(x), \quad (23)$$

where  $f_X(x)$  denotes the marginal PDF of  $X$ . In this exercise, we derive this formula.

(a) Assume that  $f_{X|D}(x|d) = e^{\alpha dx - \psi(d)} f_0(x)$  for some functions  $\psi(d)$  and  $f_0(x)$  and some constant  $\alpha$  (such that  $f_{X|D}(x|d)$  is a valid PDF for every value of  $d$ ). Define

$$\lambda(x) = \log \frac{f_X(x)}{f_0(x)}, \quad (24)$$

where  $f_X(x)$  is the marginal PDF of  $X$ , i.e.,  $f_X(x) = \int f_{X|D}(x|\delta) f_D(\delta) d\delta$ . With this, establish that

$$\mathbb{E}[D|X = x] = \frac{1}{\alpha} \frac{d}{dx} \lambda(x). \quad (25)$$

(b) Show that the case where  $X = D + N$ , where  $N$  is Gaussian noise of mean zero and variance  $\sigma^2$ , is indeed of the form required in Part (a) by finding the corresponding  $\psi(d)$ ,  $f_0(x)$ , and  $\alpha$ . Show that in this case, we have

$$\frac{f_0'(x)}{f_0(x)} = -\frac{x}{\sigma^2}, \quad (26)$$

and use this fact in combination with Part (a) to establish Tweedie's formula.

### Solution

This formula is due to M. C. K. Tweedie, "Function sof a statistical variate with given means, with special reference to Laplacian distributions," *Proc. Camb. Phil. Soc.*, Vol. 43 (1947), pp.41-49.

(a) Simply plugging in, we find

$$\frac{d}{dx} \lambda(x) = \frac{d}{dx} \log \left( \frac{\int f_{X|D}(x|\delta) f_D(\delta) d\delta}{f_0(x)} \right) \quad (27)$$

$$= \frac{d}{dx} \log \left( \frac{\int e^{\alpha \delta x - \psi(\delta)} f_0(x) f_D(\delta) d\delta}{f_0(x)} \right) \quad (28)$$

$$= \frac{d}{dx} \log \int e^{\alpha \delta x - \psi(\delta)} f_D(\delta) d\delta \quad (29)$$

$$= \frac{1}{\int e^{\alpha \delta x - \psi(\delta)} f_D(\delta) d\delta} \int \alpha \delta e^{\alpha \delta x - \psi(\delta)} f_D(\delta) d\delta \quad (30)$$

But since we know that

$$\int e^{\alpha \delta x - \psi(\delta)} f_0(x) f_D(\delta) d\delta = f_X(x), \quad (31)$$

we can rewrite

$$\frac{d}{dx}\lambda(x) = \frac{f_0(x)}{f_X(x)} \int \alpha \delta e^{\alpha \delta x - \psi(\delta)} f_D(\delta) d\delta \quad (32)$$

$$= \alpha \int \delta \underbrace{\frac{e^{\alpha \delta x - \psi(\delta)} f_0(x) f_D(\delta)}{f_X(x)}}_{f_{D|X}(d|x)} d\delta \quad (33)$$

$$= \alpha \mathbb{E}[D | X = x] \quad (34)$$

as claimed.

(b) In this case, we have

$$f_{X|D}(x|d) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}} e^{-\frac{d^2}{2\sigma^2}} e^{\frac{1}{\sigma^2}xd}. \quad (35)$$

Pattern matching with the desired form

$$f_{X|D}(x|d) = e^{\alpha dx - \psi(d)} f_0(x), \quad (36)$$

it is quickly verified that  $\alpha = 1/\sigma^2$ , and

$$f_0(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}}, \quad (37)$$

and thus,

$$f'_0(x) = -\frac{1}{\sqrt{2\pi}\sigma} \frac{2x}{2\sigma^2} e^{-\frac{x^2}{2\sigma^2}}, \quad (38)$$

giving the claimed result.

Putting things together, we have

$$\mathbb{E}[D | X = x] = \frac{1}{\alpha} \frac{d}{dx} \lambda(x) = \sigma^2 \left( \frac{d}{dx} \log f_X(x) - \frac{d}{dx} \log f_0(x) \right) \quad (39)$$

$$= \sigma^2 \left( \frac{d}{dx} \log f_X(x) - \frac{f'_0(x)}{f_0(x)} \right) \quad (40)$$

$$= \sigma^2 \left( \frac{d}{dx} \log f_X(x) + \frac{x}{\sigma^2} \right) \quad (41)$$

$$= x + \sigma^2 \frac{d}{dx} \log f_X(x), \quad (42)$$

which is the claimed formula.

#### Problem 4: FIR Wiener Filter

Consider a (discrete-time) signal that satisfies the difference equation  $d[n] = 0.5d[n-1] + v[n]$ , where  $v[n]$  is a sequence of uncorrelated zero-mean unit-variance random variables. We observe  $x[n] = d[n] + w[n]$ , where  $w[n]$  is a sequence of uncorrelated zero-mean random variables with variance 0.5.

(a) (you may skip this at first and do it later — it is conceptually straightforward) Show that for this signal model, the autocorrelation function of the signal  $d[n]$  is

$$\mathbb{E}[d[n]d[n+k]] = \frac{4}{3} \left( \frac{1}{2} \right)^{|k|}, \quad (43)$$

and thus the autocorrelation function of the signal  $x[n]$  is

$$\mathbb{E}[x[n]x[n+k]] = \begin{cases} \frac{11}{6}, & \text{for } k = 0, \\ \frac{4}{3} \left(\frac{1}{2}\right)^{|k|}, & \text{otherwise.} \end{cases} \quad (44)$$

(b) We would like to find an (approximate) linear predictor  $\hat{d}[n+3]$  using only the observations  $x[n], x[n-1], x[n-2], \dots, x[n-p]$ . Using the Wiener Filter framework, determine the optimal coefficients for the linear predictor. Find the corresponding mean-squared error for your predictor.

(c) We would like to find a linear denoiser  $\hat{d}[n]$  using *all* of the samples  $\{x[k]\}_{k=-\infty}^{\infty}$ . Find the filter coefficients and give a formula for the incurred mean-squared error.

### Solution

(a) Let us start by writing out the recursion (where we use the general  $\alpha$ , and we can plug in  $\alpha = 1/2$  later on):

$$d[n] = \alpha d[n-1] + v[n] \quad (45)$$

$$= \alpha(\alpha d[n-2] + v[n-1]) + v[n] \quad (46)$$

$$= \dots$$

$$= \alpha^k d[n-k] + \sum_{i=0}^{k-1} \alpha^i v[n-i] \quad (47)$$

In particular, letting  $k$  tend to infinity, we thus observe (for  $|\alpha| < 1$ )

$$d[n] = \sum_{i=0}^{\infty} \alpha^i v[n-i]. \quad (48)$$

Hence, we find, for positive values of  $k$ ,

$$\mathbb{E}[d[n]d[n-k]] = \mathbb{E} \left[ \left( \alpha^k d[n-k] + \sum_{i=0}^{k-1} \alpha^i v[n-i] \right) d[n-k] \right] \quad (49)$$

$$= \alpha^k \mathbb{E}[d^2[n-k]] + \sum_{i=0}^{k-1} \alpha^i \mathbb{E}[v[n-i]d[n-k]] \quad (50)$$

$$= \alpha^k \mathbb{E}[d^2[n-k]], \quad (51)$$

because by assumption, the samples  $v[i]$  are uncorrelated, hence  $\mathbb{E}[v[m]d[n]] = 0$  whenever  $m > n$ . Moreover,

$$\mathbb{E}[d^2[n-k]] = \mathbb{E} \left[ \left( \sum_{i=0}^{\infty} \alpha^i v[n-k-i] \right)^2 \right] \quad (52)$$

$$= \sum_{i=0}^{\infty} \alpha^{2i} \mathbb{E}[(v[n-k-i])^2] + \sum_{i=0}^{\infty} \sum_{j=0, j \neq i}^{\infty} \alpha^{i+j} \mathbb{E}[v[n-k-i]v[n-k-j]] \quad (53)$$

$$= \sum_{i=0}^{\infty} \alpha^{2i} \mathbb{E}[(v[n-k-i])^2], \quad (54)$$

since, by assumption,  $\mathbb{E}[v[n-i]v[n-j]] = 0$  for  $i \neq j$ . Moreover, also by assumption,  $\mathbb{E}[(v[n])^2] = 1$ , for all  $n$ , hence, using the formula for the geometric series,

$$\mathbb{E}[d^2[n-k]] = \frac{1}{1-\alpha^2}. \quad (55)$$

Combining, we thus find, for positive values of  $k$ ,

$$R_d[k] = \frac{1}{1-\alpha^2} \alpha^k. \quad (56)$$

Finally, since  $R_d[-k] = R_d[k]$ , we conclude

$$R_d[k] = \frac{1}{1-\alpha^2} \alpha^{|k|}. \quad (57)$$

Plugging in  $\alpha = 1/2$ , we thus find

$$R_d[k] = \frac{4}{3} \left(\frac{1}{2}\right)^{|k|}. \quad (58)$$

Similarly, for  $x[n]$  satisfying difference equation  $x[n] = d[n] + w[n]$ :

$$\mathbb{E}[x[n]x[n-k]] = \mathbb{E}[(d[n] + w[n])(d[n-k] + w[n-k])] \quad (59)$$

$$= \underbrace{\mathbb{E}[d[n]d[n-k]]}_{=R_d[k]} + \underbrace{\mathbb{E}[d[n]w[n-k]]}_{=0} + \underbrace{\mathbb{E}[d[n-k]w[n]]}_{=0} + \underbrace{\mathbb{E}[w[n]w[n-k]]}_{=0.5\delta[k]}, \quad (60)$$

where  $\mathbb{E}[d[n]w[n-k]] = 0$  and  $\mathbb{E}[d[n-k]w[n]] = 0$  since  $d[n]$  and  $w[n]$  are uncorrelated by assumption, and since  $\mathbb{E}[w[n]] = 0$ . We have  $\mathbb{E}[w[n]w[n-k]] = 0.5\delta[k]$  because the samples of  $w[n]$  are uncorrelated and of variance 0.5. Hence,

$$R_x[k] = \begin{cases} \frac{4}{3} + \frac{1}{2} = \frac{11}{6}, & k = 0, \\ \frac{4}{3} \left(\frac{1}{2}\right)^{|k|}, & \text{otherwise.} \end{cases} \quad (61)$$

(b) From Part (a), we already have the covariance matrix of the observations. What we still need is the correlation between the observations and the desired signal, for  $\ell \geq 0$ ,

$$\begin{aligned} \mathbb{E}[d[n+3]x[n-\ell]] &= \mathbb{E}[d[n+3](d[n-\ell] + w[n-\ell])] \\ &= \mathbb{E}[d[n+3]d[n-\ell]] = \mathbb{E}[d[n]d[n-\ell-3]] = \frac{4}{3} \left(\frac{1}{2}\right)^{\ell+3} = \frac{1}{3} \left(\frac{1}{2}\right)^{\ell+1}. \end{aligned} \quad (62)$$

As in class, the key ingredients are the covariance matrix of the observations, which is found to be

$$R_x = \begin{bmatrix} \frac{11}{6} & \frac{2}{3} & \cdots & \frac{4}{3} \left(\frac{1}{2}\right)^p \\ \frac{2}{3} & \frac{11}{6} & \cdots & \frac{4}{3} \left(\frac{1}{2}\right)^{p-1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{4}{3} \left(\frac{1}{2}\right)^p & \frac{4}{3} \left(\frac{1}{2}\right)^{p-1} & \cdots & \frac{11}{6} \end{bmatrix} \quad (63)$$

and the correlation between the desired signal and the observations, which is found to be

$$\mathbf{r}_{dx} = \left[ \frac{1}{6} \quad \frac{1}{12} \quad \frac{1}{24} \quad \cdots \quad \frac{1}{3} \left(\frac{1}{2}\right)^{p+1} \right]^\top \quad (64)$$

Thus, the optimal coefficient should be

$$\mathbf{w} = [w_0 \quad w_1 \quad \cdots \quad w_p]^\top = R_x^{-1} \mathbf{r}_{dx} \quad (65)$$

As we have seen in class, the estimation error is simply given by

$$\begin{aligned} \mathcal{E} &= \mathbb{E}[d^2[n+3]] - \mathbf{r}_{dx}^H R_x \mathbf{r}_{dx} \\ &= \frac{4}{3} - \mathbf{r}_{dx}^H R_x^{-1} \mathbf{r}_{dx} \end{aligned} \quad (66)$$

$$(67)$$

(c) The correlation between the observations and the desired signal, for all  $\ell$ ,

$$\begin{aligned}\mathbb{E}[d[n]x[n-\ell]] &= \mathbb{E}[d[n](d[n-\ell] + w[n-\ell])] \\ &= \mathbb{E}[d[n]d[n-\ell]]\end{aligned}\tag{68}$$

$$= \frac{4}{3} \left(\frac{1}{2}\right)^{|\ell|}.\tag{69}$$

The covariance matrix of  $x[n]$  becomes matrix with infinite dimensions with entry

$$R_x = \begin{bmatrix} \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \frac{11}{9} & \frac{2}{3} & \frac{1}{3} & \frac{1}{6} & \dots \\ \dots & \frac{2}{3} & \frac{11}{9} & \frac{2}{3} & \frac{1}{3} & \dots \\ \dots & \frac{1}{3} & \frac{2}{3} & \frac{11}{9} & \frac{2}{3} & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix}\tag{70}$$

and the correlation between the desired signal and the observations, which is found to be

$$\mathbf{r}_{dx} = \left[ \dots \quad \frac{1}{3} \quad \frac{2}{3} \quad \frac{4}{3} \quad \frac{2}{3} \quad \frac{1}{3} \quad \dots \right]^T\tag{71}$$

Thus, the optimal coefficient should be

$$\mathbf{w} = R_x^{-1} \mathbf{r}_{dx}\tag{72}$$

and the estimation error is simply given by

$$\begin{aligned}\mathcal{E} &= \mathbb{E}[d^2[n]] - \mathbf{r}_{dx}^H R_x \mathbf{r}_{dx} \\ &= \frac{4}{3} - \mathbf{r}_{dx}^H R_x^{-1} \mathbf{r}_{dx}\end{aligned}\tag{73}$$

$$\tag{74}$$

### Problem 5: Bounding The Exploration Bias

(a) Let  $X_1, X_2, \dots, X_n \sim \text{i.i.d. } \mathcal{N}(0, 1)$ . Let  $Y = \arg\max_i X_i$  and  $T \in \{1, 2, \dots, n\}$  is such that

$$P_{T|Y}(t|y) = \begin{cases} p, & t = y \\ \frac{1-p}{n-1}, & t \neq y \end{cases} \text{ for some } p \in [0, 1].$$

1. Compute  $I(X; T)$  where  $X = (X_1, X_2, \dots, X_n)$ . (Hint: write  $I(X; T) = H(T) - H(T|X)$ . What is the marginal distribution of  $T$ ?)

(b) Let  $X_1, \dots, X_4 \sim \text{i.i.d. } \mathcal{N}(0, 1)$  and  $X_5 \sim \mathcal{N}(0, 4)$ . Let  $Y$  and  $T$  be as in part (a) with  $p = 0.3$ .

1. Show that  $\mathbf{Pr}(Y = 5) = \int_{-\infty}^{\infty} \frac{1}{\sqrt{8\pi}} (1 - Q(x))^4 e^{-x^2/8} dx$  (where  $Q(x) = \int_x^{\infty} \frac{1}{\sqrt{2\pi}} e^{-u^2/2} du$ ), and find a corresponding numerical approximation (using Mathematica, for example).
2. Using the previous numerical approximation, find the marginal distributions  $P_Y$  and  $P_T$ .

### Solution

(a1) Since  $X_1, \dots, X_n$  are i.i.d and  $Y = \arg \max_i X_i$ , we have  $P_Y(Y = i) = \frac{1}{n}$ , for all  $i \in [n]$ . The

marginal distribution of  $T$  is

$$P_T(T = t) = \sum_{i=1}^n P_{T|Y}(T = t|Y = i)P_Y(Y = i) \quad (75)$$

$$= P_{T|Y}(T = t|Y = t)P_Y(Y = t) + \sum_{i=1, i \neq t}^n P_{T|Y}(T = t|Y = i)P_Y(Y = i) \quad (76)$$

$$= p \frac{1}{n} + (n-1) \frac{1-p}{n-1} \frac{1}{n} \quad (77)$$

$$= \frac{1}{n} \quad (78)$$

Hence,  $T$  is uniformly distributed over  $\{1, \dots, n\}$ ,

$$H(T) = - \sum_{t=1}^n P_T(T = t) \log P_T(T = t) = \log n \quad (79)$$

Additionally, given  $X_1, \dots, X_n$ ,  $Y$  is fixed, which means  $H(Y|X) = 0$ . And given  $Y$ ,  $X$  and  $T$  are independent.

$$H(T|X) = H(T, Y|X) = H(Y|X) + H(T|X, Y) = H(T|Y) = -p \log p - (n-1) \frac{1-p}{n-1} \log \frac{1-p}{n-1} \quad (80)$$

Therefore,

$$I(X; T) = H(T) - H(T|X) = \log n + p \log p + (1-p) \log \frac{1-p}{n-1} \quad (81)$$

(a2)

$$|\mathbb{E}[X_T]| \leq \sqrt{2I(X; T)} \quad (82)$$

(b1) Since  $Y = 5$  means  $X_5$  is the largest one and  $X_1, \dots, X_4$  are i.i.d.

$$\mathbf{Pr}(Y = 5) = \mathbf{Pr}(X_5 > X_1, X_5 > X_2, X_5 > X_3, X_5 > X_4) \quad (83)$$

$$= \mathbb{E}_{X_5}[\mathbf{Pr}(X_5 > X_1, X_5 > X_2, X_5 > X_3, X_5 > X_4)|X_5] \quad (84)$$

$$= \mathbb{E}_{X_5}[\mathbf{Pr}(x_5 \geq X_1|X_5 = x_5)\mathbf{Pr}(x_5 \geq X_2|X_5 = x_5)\mathbf{Pr}(x_5 \geq X_3|X_5 = x_5)\mathbf{Pr}(x_5 \geq X_4|X_5 = x_5)] \quad (85)$$

$$= \mathbb{E}_{X_5}[\mathbf{Pr}(x_5 \geq X_1|X_5 = x_5)^4] \quad [X_1, X_2, X_3, X_4 \text{ are i.i.d}] \quad (86)$$

$$= \int_{-\infty}^{\infty} \mathbf{Pr}(X_1 \leq x)^4 \mathbf{Pr}(X_5 = x) dx \quad (87)$$

$$= \int_{-\infty}^{\infty} (1 - Q(x))^4 \frac{1}{\sqrt{8\pi}} e^{-\frac{x^2}{8}} dx \quad (88)$$

$$\simeq 0.31 \quad (89)$$

Thus,  $\forall i \in \{1, 2, 3, 4\}$ , we have  $\mathbf{Pr}(Y = i) = \frac{1 - \mathbf{Pr}(Y=5)}{4} \simeq 0.1725$ .

$$\mathbf{Pr}(T = 5) = \mathbf{Pr}(T = 5, Y = 5) + \mathbf{Pr}(T = 5, Y \neq 5) \quad (90)$$

$$= \mathbf{Pr}(T = 5|Y = 5)\mathbf{Pr}(Y = 5) + \sum_{i \neq 5} \mathbf{Pr}(T = 5|Y = i)\mathbf{Pr}(Y = i) \quad (91)$$

$$\simeq 0.3 \times 0.31 + 4 \times \frac{1 - 0.3}{4} \times 0.1725 = 0.2137 \quad (92)$$



### Problem 6: Gibbs Algorithm

Let  $\mathcal{X}$  be the sample space,  $\mathcal{W}$  the hypothesis space, and let  $\ell : \mathcal{W} \times \mathcal{X} \rightarrow \mathbb{R}_+$  be a corresponding loss function. On a dataset  $D = (X_1, X_2, \dots, X_n)$ , the empirical risk for a hypothesis  $w$  is given by  $L_D(w) = \frac{1}{n} \sum_{i=1}^n \ell(w, X_i)$ . We saw in class that  $I(D; W)$  can be used to bound the generalization error. Hence, we can use it as a *regularizer* in empirical risk minimization.

- (a) First, show that given any joint distribution  $P_{XY}$  on  $\mathcal{X} \times \mathcal{Y}$  and marginal distribution  $Q$  on  $\mathcal{Y}$ ,  $D(P_{XY} || P_X P_Y) \leq D(P_{XY} || P_X Q)$ .

Since we cannot directly compute  $D(P_{DW} || P_D P_W)$ , we will use  $D(P_{DW} || P_D Q)$  as a proxy, where  $Q$  is a distribution on  $\mathcal{W}$ .

- (b) Let

$$P_{W|D}^* = \operatorname{argmin}_{P_{W|D}} \left( \mathbb{E}[L_D(W)] + \frac{1}{\beta} D(P_{DW} || P_D Q) \right).$$

1. Show that

$$\min_{P_{W|D}} \left( \mathbb{E}[L_D(W)] + \frac{1}{\beta} D(P_{DW} || P_D Q) \right) = \mathbb{E}_D \left[ \min_{P_{W|D=d}} \left( \mathbb{E}[L_d(W)] + \frac{1}{\beta} D(P_{W|D=d} || Q) \right) \right].$$

2. Show that the minimizer on the right-hand side  $P_{W|D=d}^*$  is given by

$$P_{W|D=d}^* = \frac{e^{-\beta L_d(w)} Q(w)}{\mathbb{E}_Q [e^{-\beta L_d(W)}]}.$$

This is known in the literature as the Gibbs algorithm. (Hint: Write  $\mathbb{E}[\beta L_d(W)] = \mathbb{E}[\log e^{\beta L_d(W)}]$ , combine with the KL divergence term and use non-negativity of KL divergence.)

3. Show that  $P_{W|D=d}^*$  is  $2\beta/n$ -differential private if  $\ell \in [0, 1]$ .

### Solution

- (a) For any marginal distribution  $Q$  on  $\mathcal{Y}$ ,

$$D(P_{XY} || P_X P_Y) - D(P_{XY} || P_X Q) = \sum_{x,y} P_{XY}(x, y) \left( \log \frac{P_{XY}(x, y)}{P_X(x) P_Y(y)} - \log \frac{P_{XY}(x, y)}{P_X(x) Q(y)} \right) \quad (93)$$

$$= \sum_{x,y} P_{XY}(x, y) \log \frac{Q(y)}{P_Y(y)} \quad (94)$$

$$= \sum_y P_Y(y) \log \frac{Q(y)}{P_Y(y)} \quad (95)$$

$$\stackrel{(*)}{\leq} \log \sum_y P_Y(y) \frac{Q(y)}{P_Y(y)} \quad (96)$$

$$= \log \sum_y Q(y) = 0 \quad (97)$$

where  $(*)$  is because  $\log(x)$  is a concave function of  $x$ .

(b1)

$$\min_{P_{W|D}} \left( \mathbb{E}[L_D(W)] + \frac{1}{\beta} D(P_{DW} || P_D Q) \right) \quad (98)$$

$$= \min_{P_{W|D}} \left( \mathbb{E}_D[\mathbb{E}[L_D(W)|D=d]] + \frac{1}{\beta} \sum_{w,d} P_{W|D}(w|d) P_D(d) \log \frac{P_{W|D}(w|d) P_D(d)}{P_D(d) Q} \right) \quad (99)$$

$$= \min_{P_{W|D}} \left( \mathbb{E}_D[\mathbb{E}[L_D(W)|D=d]] + \frac{1}{\beta} \sum_{w,d} P_{W|D}(w|d) P_D(d) \log \frac{P_{W|D}(w|d)}{Q} \right) \quad (100)$$

$$= \min_{P_{W|D}} \left( \mathbb{E}_D[\mathbb{E}[L_D(W)|D=d]] + \mathbb{E}_D\left[\frac{1}{\beta} D(P_{W|D} || Q) | D=d\right] \right) \quad (101)$$

$$= \mathbb{E}_D \left[ \min_{P_{W|D=d}} \left( \mathbb{E}[L_d(W)] + \frac{1}{\beta} D(P_{W|D=d} || Q) \right) \right] \quad (102)$$

(b2) Given  $D = d$ , we know that  $P_W(w) = \sum_{d'} P_{W|D}(w|d') P_D(d') = P_{W|D}(w|d)$ .

$$\arg \min_{P_{W|D=d}} \left( \mathbb{E}[L_d(W)] + \frac{1}{\beta} D(P_{W|D=d} || Q) \right) \quad (103)$$

$$= \arg \min_{P_{W|D=d}} \left( \mathbb{E}[\beta L_d(W)] + D(P_{W|D=d} || Q) \right) \quad (104)$$

$$= \arg \min_{P_{W|D=d}} \left( \mathbb{E}[\log e^{\beta L_d(W)}] + D(P_{W|D=d} || Q) \right) \quad (105)$$

$$= \arg \min_{P_{W|D=d}} \left( \sum_w \log e^{\beta L_d(w)} P_W(w) + \sum_w P_{W|D}(w|d) \log \frac{P_{W|D}(w|d)}{Q(w)} \right) \quad (106)$$

$$= \arg \min_{P_{W|D=d}} \left( \sum_w \log e^{\beta L_d(w)} P_{W|D}(w|d) + \sum_w P_{W|D}(w|d) \log \frac{P_{W|D}(w|d)}{Q(w)} \right) \quad (107)$$

$$= \arg \min_{P_{W|D=d}} \left( \sum_w P_{W|D}(w|d) (\log e^{\beta L_d(w)} + \log \frac{P_{W|D}(w|d)}{Q(w)}) \right) \quad (108)$$

$$= \arg \min_{P_{W|D=d}} \left( \sum_w P_{W|D}(w|d) \log \frac{P_{W|D}(w|d)}{Q(w) e^{-\beta L_d(w)}} \right) \quad (109)$$

$$= \arg \min_{P_{W|D=d}} \left( \sum_w P_{W|D}(w|d) \log \frac{P_{W|D}(w|d)}{Q(w) e^{-\beta L_d(w)}} \frac{\mathbb{E}_Q[e^{-\beta L_d(W)}]}{\mathbb{E}_Q[e^{-\beta L_d(W)}]} \right) \quad (110)$$

$$= \arg \min_{P_{W|D=d}} D \left( P_{W|D} || \frac{Q(w) e^{-\beta L_d(w)}}{\mathbb{E}_Q[e^{-\beta L_d(W)}]} \right) - \log \mathbb{E}_Q[e^{-\beta L_d(W)}] \quad (111)$$

$$= \arg \min_{P_{W|D=d}} D \left( P_{W|D} || \frac{Q(w) e^{-\beta L_d(w)}}{\mathbb{E}_Q[e^{-\beta L_d(W)}]} \right) \quad (112)$$

$$= \frac{Q(w) e^{-\beta L_d(w)}}{\mathbb{E}_Q[e^{-\beta L_d(W)}]} \quad (113)$$

The reason why we added  $\mathbb{E}_Q[e^{-\beta L_d(W)}]$  as a normalization term is  $P_{W|D}$  has to be a valid pmf, i.e.  $\sum_w P_{W|D}(w|d) = 1$ . However, the scaled version of  $Q$ ,  $Q(w) e^{-\beta L_d(w)}$ , may not be a valid pmf.

(b3) Suppose  $d$  and  $d'$  differ at  $j$ -th entry only. Hence,

$$e^{-\beta L_d(w)} e^{\beta L_{d'}(w)} = e^{-\frac{\beta}{n} (e(w_j, X_j) - e(w_j, X'_j))} \leq e^{\beta/n} \quad (114)$$

Similarly,

$$\frac{\mathbb{E}_Q[e^{-\beta L_d(w)}]}{\mathbb{E}_Q[e^{-\beta L_{d'}(w)}]} \leq \frac{\mathbb{E}_Q[e^{\frac{\beta}{n}} e^{-\beta L_{d'}(w)}]}{\mathbb{E}_Q[e^{-\beta L_{d'}(w)}]} \leq e^{\frac{\beta}{n}} \quad (115)$$

Thus, we have  $\frac{P_{W|D=b}^*}{P_{W|D=d'}^*} \leq e^{2\beta/n}$  and  $P_{W|D=d}^*$  is  $2\beta/n$ -differential private if  $l \in [0, 1]$ .