

Quiz 2

Date: Monday, November 21, 2016

Name: _____

Andrew ID: _____

Department: _____

Guidelines:

1. **PLEASE DO NOT TURN THIS PAGE UNTIL INSTRUCTED.**
2. Write your name, Andrew ID, and department in the spaces provided above.
3. You have **60 minutes** for this exam.
4. This exam has **9 pages** on 9 (nine) sheets of paper, including this one.
5. This exam has a total of **50 possible points**, split between **5 “short” questions** and **2 “long” questions**. The points allocated to each question are indicated next to that question.
6. This exam is open notes. You may use any materials such as cheat sheets, class notes, etc. No electronic devices are permitted.
7. The questions vary in difficulty. The points allocated to a question do not entirely reflect its difficulty. Do not spend too much time on one question.
8. Questions only appear on one side of each sheet of paper. You may use any blank space for your answer or scratch work, but please clearly indicate your answers.

1 Short Questions

1. We wish to encode a dictionary of 4 symbols $\{a, b, c, d\}$ using a ternary alphabet $\{0, 1, 2\}$.

(a) [4 points] Identify the following 4 codes as Singular (S), Non-Singular but not uniquely decodable (NS), Uniquely Decodable but not instantaneous (UD), or Instantaneous (I).

i. $\{0, 1, 11, 21\}$ **Solution:** (NS)

ii. $\{01, 10, 11, 02\}$ **Solution:** (I)

iii. $\{0, 1, 2, 1\}$ **Solution:** (S)

iv. $\{0, 112, 11, 22\}$ **Solution:** (UD)

(b) [4 points] According to the IID source distribution

$$p(a) = 1/3 \quad p(b) = 1/9 \quad p(c) = 2/9 \quad p(d) = 1/3$$

(encoding based on that order) give a ternary arithmetic code for the sequence bcd . You may assume the decoder knows when to stop. (*Note: multiple valid answers exist; give any correctly decodable answer.*)

Solution: One valid solution (the Shannon-Fano-Elias code) is 101212.

2. Suppose we want to transmit an input $X = (X_1, X_2)$ across two parallel Gaussian channels with joint correlation matrix Σ , under a total power constraint $\mathbb{E}[\|X\|_2^2] \leq 3$. What distribution of the input X maximizes the rate if

(a) [**2 points**] $\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$.

Solution: Since By the water filling method, we have $p_1 + p_2 = 3$ and $p_1 + 1 = p_2 + 2$. Hence, $p_1 = 2$ and $p_2 = 1$. Thus, $X \sim \mathcal{N}\left(0, \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}\right)$.

(b) [**3 points**] $\Sigma = \begin{bmatrix} 1 & \sqrt{2} \\ \sqrt{2} & 2 \end{bmatrix}$. (*Hint: Σ has unit eigenvectors $v_1 = \left(\frac{1}{\sqrt{3}}, \sqrt{\frac{2}{3}}\right)$, $v_2 = \left(\sqrt{\frac{2}{3}}, -\frac{1}{\sqrt{3}}\right)$ and corresponding eigenvalues $\lambda_1 = 3$ and $\lambda_2 = 0$.)*)

Solution: Applying the water filling method to λ_1 and λ_2 , we should allocate all power along the vector v_2 . Thus,

$$X \sim \mathcal{N}(0, 3v_2v_2^T) = \mathcal{N}\left(0, \begin{bmatrix} 2 & -\sqrt{2} \\ -\sqrt{2} & 1 \end{bmatrix}\right).$$

3. [7 points] No justification is necessary for this problem.

(a) For each of the following tasks, would you use the Rate-Distortion (RD) method or the Information Bottleneck (IB) method?

i. Compress a dataset X of predictors while still being able to predict a response Y .

Solution: IB (we can penalize reducing $I(X; Y)$).

ii. Compress a video without significantly sacrificing video quality.

Solution: RD (we can encode degradation in video quality as a distortion $d(x, \hat{x})$).

(b) True or False: The objective functions defining the rate-distortion function, the information bottleneck method, and the capacity of a channel can all be optimized using the Blahut-Arimoto algorithm.

Solution: True; all these problems involve minimizing or maximizing mutual information subject to constraints. This can be done via the Blahut-Arimoto algorithm.

(c) Let X be a random variable, and suppose $T(X)$ is a sufficient statistic for a parameter θ . Which of the following statements is always true?

i. $H(\theta|X) = H(\theta|T(X))$.

Solution: True; recall that T is sufficient if $I(\theta; T(X)) = I(\theta; X)$. Thus, $H(\theta|X) = H(\theta) - I(\theta; X) = H(\theta) - I(\theta; T(X)) = H(\theta|T(X))$.

ii. $H(X|\theta) \geq H(T(X)|\theta)$.

Solution: True; since $I(X; \theta) = I(\theta; T(X))$ and $H(X) \geq H(T(X))$, $H(X|\theta) = H(X) - I(X; \theta) \geq H(T(X)) - I(T(X); \theta) = H(T(X)|\theta)$.

iii. If T is minimal, then T is unique.

Solution: False; any bijection of $T(X)$ is also minimal and sufficient.

iv. $P(X|T(X))$ is independent of θ .

Solution: True; this is one of several equivalent definitions of sufficiency; it also follows, for example, from the Fisher-Neyman factorization theorem.

4. [5 points] Consider a discrete privacy mechanism \mathcal{M} over a space of data sets \mathcal{X} (e.g., $\mathcal{M} : \mathcal{X} \rightarrow \mathbb{N}$). Show that, if \mathcal{M} is ε -differentially private, then

$$\sup_{x, x' \in \mathcal{X}: x \sim x'} D_{KL}(\mathcal{M}(x) || \mathcal{M}(x')) \leq \varepsilon,$$

where $x \sim x'$ denotes that x and x' differ in a single entry and the KL divergence D_{KL} is over the randomness of \mathcal{M} .

Solution: By definition of ε -differential privacy, if \mathcal{M} is ε -differentially private, then

$$\sup_{S \subseteq \mathcal{X}, x \sim x'} \frac{\mathbb{P}[\mathcal{M}(x) \in S]}{\mathbb{P}[\mathcal{M}(x') \in S]} \leq e^\varepsilon.$$

Applying this for each singleton $S = \{M\} \subseteq \mathcal{X}$,

$$D_{KL}(\mathcal{M}(x) || \mathcal{M}(x')) = \mathbb{E}_{M \sim \mathcal{M}} \left[\log \left(\frac{\mathbb{P}[\mathcal{M}(x) = M]}{\mathbb{P}[\mathcal{M}(x') = M]} \right) \right] \leq \mathbb{E}_{M \sim \mathcal{M}} [\varepsilon] \leq \varepsilon.$$

5. [5 points] Consider universal prediction in the context of online linear regression. At each time point t , we
1. observe a predictor $x_t \in \mathbb{R}^D$.
 2. output a prediction $\hat{y}_t \in \mathbb{R}$.
 3. observe a true $y_t \in \mathbb{R}$.
 4. suffer squared error loss $\ell(y_t, \hat{y}_t) = (y_t - \hat{y}_t)^2$.

Suppose we consider all linear hypotheses, i.e., our hypothesis space can be represented by $\Theta = \mathbb{R}^D$, where $w \in \Theta$ predicts $\langle w, x \rangle$ for each input $x \in \mathbb{R}^D$. Explain how to predict each \hat{y}_t , using the exponential weights algorithm with a prior p_0 over Θ and learning rate η .

Solution: At each time point t , we sample a linear hypothesis w_t from the probability distribution

$$p_{t-1}(w) \propto p_0(w) e^{-\eta \sum_{s=1}^{t-1} (\langle w, x_s \rangle - y_s)^2}.$$

We then make the prediction $\hat{y}_t = \langle w_t, x_t \rangle$.

Alternative Solution: Because squared error loss is convex in its first argument, as an alternative to sampling (which works in more general settings), we can use the expectation of the posterior, i.e. set $w_t = \mathbb{E}_{w \sim p_{t-1}} [w]$, and then still predict $y_t = \langle w_t, x_t \rangle$. In particular, by Jensen's inequality, this always has lower expected loss than random sampling:

$$\left(y_t - \left\langle \mathbb{E}_{w \sim p_{t-1}} [w], x_t \right\rangle \right)^2 = \left(\mathbb{E}_{w \sim p_{t-1}} [y_t - \langle w, x_t \rangle] \right)^2 \leq \mathbb{E}_{w \sim p_{t-1}} [(y_t - \langle w, x_t \rangle)^2].$$

2 Long Questions

1. [10 points] Consider a channel C that takes a binary input X and returns a binary output Y , according to the following conditional distribution:

x	$\mathbb{P}[Y = 0 X = x]$	$\mathbb{P}[Y = 1 X = x]$
0	1	0
1	0.5	0.5

What is the capacity (in bits) of C ?

Hint: Check that $I(X; Y)$ is concave in p and find the maximizer.

Solution: Since X is binary, without loss of generality, $X \sim \text{Bernoulli}(p)$ for some $p \in [0, 1]$. Then, $Y \sim \text{Bernoulli}(p/2)$. Also, $H(Y|X = 0) = 0$ and $H(Y|X = 1) = 1$. Thus,

$$\begin{aligned}
 I(X; Y) &= H(Y) - H(Y|X) \\
 &= H(p/2) - ((1 - p)H(Y|X = 0) + pH(Y|X = 1)) \\
 &= H(p/2) - p \\
 &= -(1 - p/2) \log(1 - p/2) - (p/2) \log(p/2) - p
 \end{aligned} \tag{1}$$

From line (1), it is easy to see that $I(X; Y)$ is concave in p . Thus, at the optimum p^* ,

$$0 = \left. \frac{d}{dp} I(X; Y) \right|_{p=p^*} = \frac{1}{2} (\log(1 - p^*/2) - \log(p^*/2)) - 1 = \frac{1}{2} \log(2/p^* - 1) - 1.$$

Solving for p^* gives $p^* = 2/5$. Thus,

$$I(X; Y) = \frac{4}{5} \log(5/4) + \frac{1}{5} \log(5) - 2/5 = \log(5) - 2 \approx 0.32.$$

One should also note that, if $p \in \{0, 1\}$, then $I(X; Y) = 0$, so these are clearly not optima.

2. [10 points] Consider n IID samples $(X_1, Y_1), \dots, (X_n, Y_n) \sim \mathcal{N}(0, \Sigma)$, where

$$\Sigma = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$$

for some $\rho \in [-1, 1]$.

- (a) For a fixed $\rho_* \in (0, 1)$, use Le Cam's method to derive a minimax lower bound for testing the null hypothesis $\rho = 0$ against the alternative hypothesis $|\rho| > \rho_*$.
- (b) For a given constant $c \in (0, 1)$, what is largest value of ρ_* such that, according to your lower bound, every test has worst-case error probability at least c ?

Hint: The KL divergence between Gaussians $\mathcal{N}(0, \Sigma_0)$ and $\mathcal{N}(0, \Sigma_1)$ over \mathbb{R}^2 is

$$D_{KL}(\mathcal{N}(0, \Sigma_0), \mathcal{N}(0, \Sigma_1)) = \frac{1}{2} \left(\log \frac{|\Sigma_0|}{|\Sigma_1|} - 2 + \text{tr}(\Sigma_0^{-1} \Sigma_1) \right),$$

where $|\Sigma| = \Sigma_{1,1}\Sigma_{2,2} - \Sigma_{1,2}\Sigma_{2,1}$ and $\text{tr}(\Sigma) = \Sigma_{1,1} + \Sigma_{2,2}$ denote the determinant and trace of Σ .

Solution:

- (a) By the hint, the KL divergence between the distribution P_0 of (X_1, Y_1) when $\rho = 0$ and the distribution P_1 of (X_1, Y_1) when $\rho = \rho_*$ is $D_{KL}(P_0, P_1) = -\frac{1}{2} \log(1 - \rho_*^2)$.

Le Cam's method then gives a lower bound of

$$\frac{1}{2} - \frac{1}{2} \sqrt{\frac{n}{2} D_{KL}(P_0, P_1)} = \frac{1}{2} - \frac{1}{4} \sqrt{-n \log(1 - \rho_*^2)}. \quad (2)$$

- (b) One can check that, by setting $\rho_* = \sqrt{1 - \exp\left(-\frac{4(1-2c)^2}{n}\right)}$, the lower bound (2) is c .

This space is intentionally blank. You may use it for scratch work.

Please do not mark below this line.

Problem	Max	Points
S1	8	
S2	5	
S3	7	
S4	5	
S5	5	
L1	10	
L2	10	
Total	50	