

**Problem 0** [1 bonus mark]. Write your name and student ID (number) on every submitted answer sheet.

**Problem 1** [10 marks]. In a typical case, more training samples will result in training mean squared error (MSE) loss going up and validation error going down. Fill in the blanks with "up" and/or "down".

**Problem 2** [30 marks]. — discrete

The variance of a random variable  $X$  is defined as

$$\text{Var}_{X \sim P(X)}[X] = \mathbb{E}_{X \sim P(X)}[(X - \mathbb{E}_{X \sim P(X)}[X])^2].$$

a) Prove  $\text{Var}_{X \sim P(X)}[X] = \mathbb{E}_{X \sim P(X)}[X^2] - (\mathbb{E}_{X \sim P(X)}[X])^2$ .

b) Use the definition of expectation to prove its linearity:

$$\mathbb{E}_{X \sim P(X)}[af(X) + bg(X)] = a\mathbb{E}_{X \sim P(X)}[f(X)] + b\mathbb{E}_{X \sim P(X)}[g(X)]$$

c) Does linearity hold for variance? Prove or disprove the following statement:

$$\begin{aligned} & \text{Var}_{X \sim P(X)}[af(X) + bg(X)] \\ &= a \text{Var}_{X \sim P(X)}[f(X)] + b \text{Var}_{X \sim P(X)}[g(X)] \end{aligned}$$

$$\begin{aligned} \text{a) } \text{Var}_{X \sim P(X)}[X] &= \mathbb{E}_{X \sim P(X)}[X^2 + (\mathbb{E}_{X \sim P(X)}[X])^2 - 2X \mathbb{E}_{X \sim P(X)}[X]] \\ &= \mathbb{E}_{X \sim P(X)}[X^2] + (\mathbb{E}_{X \sim P(X)}[X])^2 - 2(\mathbb{E}_{X \sim P(X)}[X])^2 \\ &= \mathbb{E}_{X \sim P(X)}[X^2] - 2(\mathbb{E}_{X \sim P(X)}[X])^2 \end{aligned}$$

Dropping  $X \sim P(X)$  is okay.

$$\begin{aligned} \text{b) } \mathbb{E}_{X \sim P(X)}[af(X) + bg(X)] &= \sum_X [p(X) \cdot (af(X) + bg(X))] \\ &= a \sum_X p(X) f(X) + b \sum_X p(X) g(X) \\ &= a \mathbb{E}_{X \sim P(X)}[f(X)] + b \mathbb{E}_{X \sim P(X)}[g(X)] \end{aligned}$$

c). No. Counter-example:

$$\begin{aligned} \text{Let } X &\sim N(0, 1) \\ f(X) &= X \quad a = b = 1 \\ g(X) &= -X \end{aligned}$$

$$\text{Var}_{X \sim P(X)}[X - X] = \text{Var}_{X \sim P(X)}[0] = 0$$

$$\begin{aligned} \text{But } a \text{Var}_{X \sim P(X)}[f(X)] + b \text{Var}_{X \sim P(X)}[g(X)] \\ = a \cdot \sigma^2 + b \cdot \sigma^2 = 2 \end{aligned}$$

**Problem 3** [60 marks]. Consider (potentially problematic) maximum likelihood estimation (MLE) for feature selection. We assume

$t = \sum_{i=0}^d \alpha_i w_i x_i + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  and  $x_i$  is an input feature.

Both  $\alpha_i \in \{0, 1\}$  and  $w_i \in \mathbb{R}$  are model parameters. We denote column vectors by bold letters, for example,  $\alpha = (\alpha_0, \dots, \alpha_d)^\top$ .

a)  $\alpha_i = 1$  means the  $i$ th feature is selected, and  $\alpha_i = 0$  means the  $i$ th feature is discarded. Fill in the blanks with "selected" and/or "discarded".

b) Consider a training set  $\mathcal{D} = \{(x^{(m)}, t^{(m)})\}_{m=1}^M$ . Present the MLE principle and show its equivalence to mean squared error, denoted by  $J(\alpha, w)$ . Provide derivation steps.

c) Is the loss function convex or not? Briefly explain why or why not. No proof is needed.

No. dom  $J$  is not convex as  $\alpha_i$  is discrete

d) Closed-form solutions may still be possible for certain non-convex functions. Prove that  $\min_{\alpha, w: \alpha_i=1} J(\alpha, w) \leq \min_{\alpha, w: \alpha_i=0} J(\alpha, w)$ , where  $J(\alpha, w)$  is defined above.

e) Give a closed-form solution to the problem  $\min_{\alpha, w} J(\alpha, w)$ . Provide derivation steps, including the derivation for the optimal  $w$ .

f) Explain why such MLE cannot help feature selection.

because MLE will select all features

#### Hints:

- The optimization variables include  $\alpha_i$  and  $w_i$  for  $i = 0, \dots, d$ .
- The Gaussian assumption suggests  $p(\epsilon) = c_1 \exp\{-c_2 \epsilon^2\}$  for some positive constants  $c_1$  and  $c_2$ .
- For explanation questions, a few words would suffice. Then can be fit between question lines.

END OF THE EXAM

b) Maximum Likelihood estimation:   
 $\text{maximize}_{\alpha, w} \sum_{m=1}^M \log p(t^{(m)} | x^{(m)}; w)$    
 $\Leftrightarrow \text{maximize}_{\alpha, w} \sum_{m=1}^M \log c_1 \exp \left\{ -c_2 \left( t^{(m)} - \sum_{i=0}^d \alpha_i w_i x_i \right)^2 \right\}$    
 $\Leftrightarrow \text{minimize}_{\alpha, w} \frac{1}{2M} \sum_{m=1}^M \left( t^{(m)} - \sum_{i=0}^d \alpha_i w_i x_i \right)^2$    
 (Handwritten notes: "argmax... = argmax..." and "max... = max..." are okay. "not ideal but earns marks" is written next to the first equation. "constant may vary" is written below the second equation. "not okay" is written next to the third equation.)

d) For any optimum  $\alpha^*, w^*$  ( $\alpha_i=0$ ) that minimizes  $J(\alpha, w)$  we may have optimum  $\tilde{\alpha}^*, \tilde{w}^*$  such that

$$\tilde{\alpha}_i^* = 1, \tilde{w}_i^* = 0, \text{ and } \tilde{\alpha}_j^* = \alpha_j^*, \tilde{w}_j^* = w_j^* \text{ for } j \neq i$$

that achieves the same value of  $J(\alpha, w)$

$$\text{Therefore } \min_{\alpha, w: \alpha_i=1} J(\alpha, w) \leq \min_{\alpha, w: \alpha_i=0} J(\alpha, w)$$

e) a closed-form solution may have  $\alpha^* = 1$    
 Then the problem reduces to classic linear regression

$$J(w) = \frac{1}{2M} \|Xw - t\|^2$$

$$= \frac{1}{2M} (w^T X^T X w - 2 t^T X w + t^T t)$$

$$\frac{\partial J}{\partial w} = \frac{1}{2M} (2 X^T X w - 2 X^T t) \stackrel{!}{=} 0$$

$$w^* = (X^T X)^{-1} X^T t$$