

STAT547U Diagnostic Problem Set

YOUR NAME

The purpose of this problem set is to ensure that you are fluent with linear algebra and basic probability theory, which will be necessary for the course. If these problems feels extremely difficult, then you will likely find this course to be technically overwhelming.

There are 3 questions (broken into subquestions) in this problem set for a total of 60 pts.

A quick note on notation. Variable names should use the following convention:

- Scalars (deterministic and random) will be represented by lowercase/non-bold letters (e.g. a , θ , etc.).
 - Vectors (deterministic and random) will be represented by lowercase/bold letters (e.g. \mathbf{a} , $\boldsymbol{\theta}$, etc.).
 - Matrices (deterministic and random) will be represented by uppercase/bold letters (e.g. \mathbf{A} , $\boldsymbol{\Theta}$, etc.).
-

1) Introduction to Functional Analysis (20 pts)

Let $\mu(\mathbf{x}_i)$ be a probability distribution over \mathbb{R}^d with mean 0 and covariance $\boldsymbol{\Sigma}$; i.e.

$$\mathbb{E}[\mathbf{x}_i] = 0, \quad \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top] = \boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}.$$

Consider the space of linear $\mathbb{R}^d \rightarrow \mathbb{R}$ functions $\mathcal{H} := \{f(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x} : \boldsymbol{\theta} \in \mathbb{R}^d\}$.

1. (4 points. L2 inner products) For any two functions $f(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x}$ and $f'(\mathbf{x}) = \boldsymbol{\theta}'^\top \mathbf{x}$, define the $L2$ inner product $\langle f, f' \rangle_{\mathcal{H}}$ as

$$\langle f, f' \rangle_{L2} := \mathbb{E}[f(\mathbf{x})f'(\mathbf{x})].$$

Show that this inner product is equal to $\boldsymbol{\theta}^\top \boldsymbol{\Sigma} \boldsymbol{\theta}'$.

2. (4 points. Orthonormal functions) An orthonormal basis of \mathcal{H} is a set of functions $\{\phi_1, \dots, \phi_d\}$ such that

- any function $f(\mathbf{x}) \in \mathcal{H}$ can be written as the linear combination of ϕ_1, \dots, ϕ_d and
- $\mathbb{E}[\phi_i(\mathbf{x})\phi_j(\mathbf{x})] = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}$.

Construct an orthonormal basis of \mathcal{H} — $\phi_1(\mathbf{x}), \dots, \phi_d(\mathbf{x})$ —in terms of $(\lambda_1, \mathbf{v}_1), \dots, (\lambda_d, \mathbf{v}_d)$ —the eigenvalues and eigenvectors of $\boldsymbol{\Sigma}$.

(Hint: This answer should take you ≈ 5 lines.)

3. (4 points. \mathcal{H} inner products) For any two functions $f(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x}$ and $f'(\mathbf{x}) = \boldsymbol{\theta}'^\top \mathbf{x}$, define the \mathcal{H} inner product $\langle f, f' \rangle_{\mathcal{H}}$ as

$$\langle f, f' \rangle_{L2} := \sum_{i=1}^d \frac{1}{\sigma_i} \langle f, \phi_i \rangle_{L2} \langle f', \phi_i \rangle_{L2},$$

where $\{\phi_1, \dots, \phi_d\}$ is any orthonormal basis of \mathcal{H} , and $\sigma_1, \dots, \sigma_d$ are scalars such that, for any $\mathbf{z} \in \mathbb{R}^d$ that is fixed (i.e. not random),

$$\int \phi_i(\mathbf{x})(\mathbf{x}^\top \mathbf{z}) d\mu(\mathbf{x}) = \mathbb{E}[\phi_i(\mathbf{x})(\mathbf{x}^\top \mathbf{z})] = \sigma_i \mathbf{z}$$

Show that this inner product is equal to $\boldsymbol{\theta}^\top \boldsymbol{\theta}'$.

(Hint: This proof should take you ≈ 6 lines.)

4. (4 points. $\mathcal{H} \subset L2$.) Now assume that $d \rightarrow \infty$. Assuming that the eigenvalues of $\boldsymbol{\Sigma}$ are summable (i.e. $\lim_{d \rightarrow \infty} \sum_{i=1}^d \lambda_i < \infty$), prove that—for any $f \in \mathcal{H}$ — $\|f\|_{L2} < \infty$ if $\|f\|_{\mathcal{H}} < \infty$, where

$$\|f\|_{L2} := \sqrt{\langle f, f \rangle_{L2}}, \quad \|f\|_{\mathcal{H}} := \sqrt{\langle f, f \rangle_{\mathcal{H}}}.$$

(Note: if you're smart with your linear algebra, this proof should take you ≈ 4 lines.)

5. (4 points. $L2 \not\subset \mathcal{H}$.) Show that the converse is not true. I.e., as $d \rightarrow \infty$, construct a function where $\|f\|_{L2} < \infty$ but $\|f\|_{\mathcal{H}} = \infty$. (You should continue to assume that $\lim_{d \rightarrow \infty} \sum_{i=1}^d \lambda_i < \infty$.)

2) Ridge(less) Linear Regression (20 pts)

Assume that we are performing linear regression to predict a real-valued response variable $y \in \mathbb{R}$ from a d -dimensional input variable $\mathbf{x} \in \mathbb{R}^d$ using the parameter vector $\boldsymbol{\theta}^*$:

$$y = \mathbf{x}^\top \boldsymbol{\theta}^*$$

Given training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, the ridge regression estimator $\hat{\boldsymbol{\theta}}_\lambda$ is given by:

$$\hat{\boldsymbol{\theta}}_\lambda = \left(\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{y}. \quad (1)$$

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ is a matrix where each *row* represents a training input (i.e. $\mathbf{X}^\top = [\mathbf{x}_1 \ \dots \ \mathbf{x}_n]$),
- $\mathbf{y} \in \mathbb{R}^n$ is the concatenation of the training responses, and
- $\lambda > 0$ is the ridge constant.

For all problems, assume that \mathbf{X} is full rank..

1. (5 points.) Using Woodbury's matrix inversion lemma, prove that $\hat{\boldsymbol{\theta}}_\lambda$ is equal to

$$\hat{\boldsymbol{\theta}}_\lambda = \mathbf{X}^\top \left(\lambda \mathbf{I} + \mathbf{X} \mathbf{X}^\top \right)^{-1} \mathbf{y}. \quad (2)$$

(Hint: after applying Woodbury, consider rewriting \mathbf{I} as $(\lambda \mathbf{I} + \mathbf{X} \mathbf{X}^\top)^{-1} (\lambda \mathbf{I} + \mathbf{X} \mathbf{X}^\top)$.)

2. (3 points.) We now have two formula for $\hat{\boldsymbol{\theta}}_\lambda$: Eq. (1) and Eq. (2). From a computational perspective, which formula is preferable when $n > d$? When $d < n$? Justify your answer in ≈ 2 sentences. (Hint: think about the matrices you have to invert.)
3. (3 points.) Consider the scenario where $d > n$; i.e. we have more features than training data. This scenario is often referred to as the *overparameterized regime*. Assuming that \mathbf{X} is full rank, show that the ridge estimator *interpolates* the training data as $\lambda \rightarrow 0$; i.e.

$$\lim_{\lambda \rightarrow 0} \mathbf{X} \hat{\boldsymbol{\theta}}_\lambda = \mathbf{y}.$$

This proof should take you ≈ 2 lines.

4. (3 points.) Now consider the scenario where $n < d$; i.e. we have more training data than features. This scenario is often referred to as the *underparameterized regime*. Assuming that \mathbf{X} is full rank, show that

$$\mathbf{X} \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top = \mathbf{U} \mathbf{U}^\top,$$

where \mathbf{U} is some $n \times d$ orthonormal matrix.

5. (3 points.) In ≈ 3 sentences, derive the eigenvalues and the corresponding eigenvectors of $\mathbf{U} \mathbf{U}^\top$.
6. (3 points.) Putting the last two results together, argue why—in general— $\hat{\boldsymbol{\theta}}_\lambda$ does not interpolate the training data in the *underparameterized ridgeless regime* ($\lambda = 0, n > d$).
-

3) Bias-Variance Tradeoff (20 pts)

Now imagine that $\mathbf{x}_1, \dots, \mathbf{x}_n$ from the previous problem are i.i.d. samples from a distribution $\mu(\mathbf{x})$ and $y_i \sim \mathcal{N}(\mathbf{x}_i^\top \boldsymbol{\theta}^*, \sigma^2)$ for each $i \in [1, n]$. Assume that the y_i are conditionally independent given \mathbf{x}_i . Let $\hat{\boldsymbol{\theta}}_\lambda$ be the ridge estimator given in Eq. (1) (or, equivalently, in Eq. (2)).

Consider a new (independent) test point (\mathbf{x}, y) , where $\mathbf{x} \sim \mu(\mathbf{x})$ and $y \sim \mathcal{N}(\mathbf{x}^\top \boldsymbol{\theta}^*, \sigma^2)$, the *risk* is defined as

$$\mathcal{R} := \mathbb{E} \left[\left(\mathbf{x}^\top \hat{\boldsymbol{\theta}}_\lambda - \mathbf{x}^\top \boldsymbol{\theta}^* \right)^2 \right].$$

- (8 points.) Decompose \mathcal{R} into two components that represent *squared bias* and *variance* of $\mathbf{x}^\top \hat{\boldsymbol{\theta}}_\lambda$. (Hint: consider adding and subtracting $\mathbf{x}^\top \mathbb{E}[\hat{\boldsymbol{\theta}}_\lambda]$ inside the parentheses.)
- (6 points.) Show that the squared bias of $\mathbf{x}^\top \hat{\boldsymbol{\theta}}_\lambda$ is 0 when $\lambda = 0$ and $n > d$. (This proof should take 3 lines when you write $\mathbb{E}[\hat{\boldsymbol{\theta}}_{\lambda=0}]$ as $\mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}]$ and rewrite \mathbf{y} in terms of \mathbf{X} .)
- (8 points.) Show that bias term of is nonzero when $d > n$ and $\lambda \rightarrow 0$. For this problem you can make the following assumptions:
 - Limits and expectations can be interchanged— i.e. you can assume that $\lim_{\lambda \rightarrow 0} \mathbb{E}[\hat{\boldsymbol{\theta}}_\lambda] = \mathbb{E}[\mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y}]$. No need to dive into the measure theory!
 - $\mathbb{E}[\mathbf{x} \mathbf{x}^\top] := \boldsymbol{\Sigma}$ is positive definite.

(There are many ways to prove this statement. Use your linear algebra and probability skills. Be creative!)