

Homework 3 solutions

DUE: SATURDAY, FEBRUARY 8, 11:59PM

For all the problems below, assume $(\Omega, \mathcal{F}, \mathbb{P})$ is a probability space.

Problem 3.1 (Warmup: eigenvalues and eigenvectors of $A + \lambda I$). Suppose $A = X\Lambda X^{-1} \in \mathcal{M}_{n \times n}(\mathbb{C})$ where $X \in \mathcal{M}_{n \times n}(\mathbb{C})$ is an invertible matrix having the eigenvectors of A as its columns, and $\Lambda \in \mathcal{M}_{n \times n}(\mathbb{C})$ is a diagonal matrix having the eigenvalues of A on its main diagonal.

a) Show that we can write the identity matrix $I_{n \times n}$ as

$$I_{n \times n} = X I_{n \times n} X^{-1}. \quad (3.1)$$

Conclude that for any $\lambda \in \mathbb{C}$,

$$\lambda I_{n \times n} = X(\lambda I_{n \times n})X^{-1}. \quad (3.2)$$

b) For any $\lambda \in \mathbb{C}$, show that

$$A + \lambda I_{n \times n} = X(\Lambda + \lambda I_{n \times n})X^{-1} \quad (3.3)$$

c) Use part b) to show that

$$(A + \lambda I_{n \times n})X = X(\Lambda + \lambda I_{n \times n}). \quad (3.4)$$

Use this identity to identify the eigenvectors and eigenvalues of $A + \lambda I_{n \times n}$.

d) Suppose A is not an invertible matrix but its eigenvalues all have non-negative real parts (in particular we're assuming that A does not have any negative eigenvalues). Use part c) to explain why $A + \lambda I_{n \times n}$ will always be invertible for any $\mathbb{R} \ni \lambda > 0$.

Solution.

a) We note that

$$X I_{n \times n} X^{-1} = X X^{-1} = I_{n \times n}, \quad (3.5)$$

therefore

$$\lambda I_{n \times n} = \lambda(X I_{n \times n} X^{-1}) = X(\lambda I_{n \times n})X^{-1}. \quad (3.6)$$

b) We note that by assumption, $A = X\Lambda X^{-1}$, therefore

$$A + \lambda I_{n \times n} = X\Lambda X^{-1} + X(\lambda I_{n \times n})X^{-1} = X(\Lambda + \lambda I_{n \times n})X^{-1}. \quad (3.7)$$

c) Since X is invertible, this shows that

d)

$$(A + \lambda I_{n \times n})X = X(\Lambda + \lambda I_{n \times n})X^{-1}X = X(\Lambda + \lambda I_{n \times n}). \quad (3.8)$$

By examining the columns of the LHS and RHS we see that the columns of X are eigenvectors of $A + \lambda I_{n \times n}$ and the diagonal elements of $\Lambda + \lambda I_{n \times n}$ are eigenvalues of A . Since X is invertible, we also know that $A + \lambda I_{n \times n}$ and $\Lambda + \lambda I_{n \times n}$ have the same characteristic polynomial. This follows from $\det X \det X^{-1} = 1$, which implies

$$\begin{aligned} p_{A+\lambda I_{n \times n}}(\alpha) &= \det(A + \lambda I_{n \times n} - \alpha I_{n \times n}) = \det(X(\Lambda + \lambda I_{n \times n} - \alpha I_{n \times n})X^{-1}) \\ &= \det(\Lambda + \lambda I_{n \times n} - \alpha I_{n \times n}) = p_{\Lambda+\lambda I_{n \times n}}(\alpha), \quad \alpha \in \mathbb{C}. \end{aligned} \quad (3.9)$$

This shows that $A + \lambda I_{n \times n}$ and $\Lambda + \lambda I_{n \times n}$ have the same eigenvalues. Since the eigenvalues of $\Lambda + \lambda I_{n \times n}$ are the entries on its main diagonal, as it is a diagonal matrix, we may conclude that the eigenvalues of $A + \lambda I_{n \times n}$ are the eigenvalues of A shifted uniformly by a factor of λ .

e) We note that a square matrix is invertible if and only if it does not have zero as an eigenvalue. By assumption, since all of the eigenvalues of A have non-negative real parts and $\lambda > 0$, by the result from the previous part we may conclude that all the eigenvalues of $A + \lambda I_{n \times n}$ have strictly positive real part, therefore $A + \lambda I_{n \times n}$ cannot have zero as an eigenvalue and is thus invertible.

□

Problem 3.2 (Ridge regression part II). Recall that if an $m \times n$ matrix A has independent columns, then $A^T A$ is an invertible $n \times n$ square matrix. Thus, the least squares solution $\hat{\mathbf{x}}$ satisfying $A^T A \hat{\mathbf{x}} = A^T \mathbf{b}$ can be written as

$$\hat{\mathbf{x}} = (A^T A)^{-1} A^T \mathbf{b}. \quad (3.10)$$

However, if A represents a data matrix where its rows represent samples and its columns represents features, it is very common for the columns of A to be dependent (or very close to being dependent) if there is high correlation between the feature variables. So in practice, the matrix $A^T A$ is usually very close to being non-invertible. In statistics, this is referred to as the phenomenon of *multicollinearity*.

One way to combat multicollinearity is by replacing the standard least squares estimator with the ridge estimator

$$\hat{\mathbf{x}}_{\text{ridge}} = (A^T A + \lambda I_{n \times n})^{-1} A^T \mathbf{b}, \quad (3.11)$$

where $\mathbb{R} \ni \lambda > 0$ is a positive real number. Our goal in this problem is to show that for any $\lambda > 0$, the ridge estimator is well-defined, by showing that the matrix $A^T A + \lambda I_{n \times n}$ is invertible.

- a) Use the spectral theorem to show that there exists an orthogonal matrix Q and a diagonal matrix Λ such that

$$A^T A = Q \Lambda Q^T. \quad (3.12)$$

- b) Show that one can write

$$A^T A + \lambda I_{n \times n} = Q(\Lambda + \lambda I_{n \times n})Q^T. \quad (3.13)$$

- c) Using the previous problem, how are the eigenvalues of $A^T A + \lambda I_{n \times n}$ related to eigenvalues of $A^T A$?
d) Conclude that for any $\lambda > 0$, the $n \times n$ matrix $A^T A + \lambda I_{n \times n}$ is invertible. This shows that the ridge coefficient $\hat{\mathbf{x}}_{\text{ridge}}$ is well-defined for any $\lambda > 0$.

Solution.

- a) We note that $A^T A$ is an $n \times n$ real symmetric matrix since A is real and

$$(A^T A)^T = A^T (A^T)^T = A^T A. \quad (3.14)$$

So by the spectral theorem, there exists an orthogonal matrix Q and a diagonal matrix Λ for which

$$A^T A = Q \Lambda Q^T. \quad (3.15)$$

- b) Note that since Q is orthogonal, we have $Q^{-1} = Q^T$. Therefore $Q Q^T = I$. So $\lambda I = Q(\lambda I)Q^T$, and

$$A^T A + \lambda I_{n \times n} = Q \Lambda Q^T + Q(\lambda I_{n \times n})Q^T = Q(\Lambda + \lambda I_{n \times n})Q^T. \quad (3.16)$$

- c) We learned from last week's homework that the eigenvalues of $A^T A + \lambda I_{n \times n}$ are the eigenvalues of $A^T A$ shifted by a factor of λ .
d) We've shown in lecture that the eigenvalues of $A^T A$ are non-negative. Since $\lambda > 0$, we have $\lambda_i + \lambda > 0$ for each eigenvalue λ_i of A , so the eigenvalues of $A^T A + \lambda I_{n \times n}$ are strictly positive. Since $A^T A + \lambda I_{n \times n}$ does not have 0 as an eigenvalue, this implies that $A^T A + \lambda I_{n \times n}$ is invertible.

□

Problem 3.3 (Another problem on correlation). Let $\mathbb{Z} \ni n \geq 2$. Suppose $X_1, \dots, X_n : \Omega \rightarrow \mathbb{R}$ are random variables for which the pairwise correlation coefficients are all equal to $\rho \in \mathbb{R}$. In other words, for all $i, j \in \{1, \dots, n\}$, we have $\text{Corr}(X_i, X_j) = \rho$ for $i \neq j$. In this problem our goal is to find the range of possible values of ρ .

- a) Let $A \in \mathcal{M}_{n \times n}(\mathbb{R})$ be the correlation matrix of the random vector $\mathbf{X} : \Omega \rightarrow \mathbb{R}^n$ defined via $\mathbf{X} = (X_1 \dots X_n)^T$. Show that

$$A = (1 - \rho)I_{n \times n} + \begin{pmatrix} \rho & \rho & \cdots & \rho \\ \rho & \rho & \cdots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \cdots & \rho \end{pmatrix} = (1 - \rho)I_{n \times n} + \rho B, \quad (3.17)$$

where $B \in \mathcal{M}_{n \times n}(\mathbb{R})$ is a matrix with all entries equal to 1.

- b) Show that if $\lambda \in \mathbb{C}$ is an eigenvalue of B , then either $\lambda = 0$ or $\lambda = n$.
c) Use (3.17) and part b) to show that if $\lambda \in \mathbb{C}$ is an eigenvalue of A , then either $\lambda = 1 - \rho$ or $\lambda = 1 + (n - 1)\rho$.
d) Use part c) and the fact that A is positive semi-definite to show that

$$-\frac{1}{n-1} \leq \rho \leq 1. \quad (3.18)$$

- e) Let $n = 3$. Construct explicit examples of random variables $X_1, X_2, X_3 : \Omega \rightarrow \mathbb{R}$ for which the minimum and maximum values of ρ in (3.18) are achieved. (Hint: use Homework 2 Problem 5)

Solution. For part a), we note that

$$A = \begin{pmatrix} 1 & \rho & \cdots & \rho \\ \rho & 1 & \cdots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \cdots & 1 \end{pmatrix} = \begin{pmatrix} (1 - \rho) + \rho & \rho & \cdots & \rho \\ \rho & (1 - \rho) + \rho & \cdots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \cdots & (1 - \rho) + \rho \end{pmatrix} = (1 - \rho)I_{n \times n} + \rho B. \quad (3.19)$$

For part b), suppose

$$\mathbf{v} = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} \in \mathbb{C}^n \setminus \{\mathbf{0}\} \quad (3.20)$$

is an eigenvalue of B with eigenvalue $\lambda \in \mathbb{C}$. Then we have

$$v_1 + v_2 + \cdots + v_n = \lambda v_1 \quad (3.21)$$

$$v_1 + v_2 + \cdots + v_n = \lambda v_2, \quad (3.22)$$

$$\vdots \quad (3.23)$$

$$v_1 + v_2 + \cdots + v_n = \lambda v_n. \quad (3.24)$$

This implies

$$\lambda v_1 = \lambda v_2 = \cdots = \lambda v_n, \quad (3.25)$$

so either $\lambda = 0$ or $v_1 = v_2 = \cdots = v_n$. In the latter case, we have $v_1 + v_2 + \cdots + v_n = nv_i = \lambda v_i$ for all $1 \leq i \leq n$. Since at least one of the v_i 's is non-zero, we must have $\lambda = n$.

For part c), we note that if $\mathbf{v} \in \mathbb{C}^n \setminus \{\mathbf{0}\}$ is an eigenvector of B with eigenvalue λ , then we have

$$A\mathbf{v} = (1 - \rho)I_{n \times n}\mathbf{v} + \rho B\mathbf{v} = (1 - \rho)\mathbf{v} + \rho\lambda\mathbf{v} = (1 - \rho + \lambda\rho)\mathbf{v}. \quad (3.26)$$

This shows that \mathbf{v} is an eigenvector of A with eigenvalue $1 - \rho + \lambda\rho$. Conversely, if \mathbf{v} is an eigenvector of A with eigenvalue λ , then we have

$$\rho B\mathbf{v} = A\mathbf{v} - (1 - \rho)\mathbf{v} = \lambda\mathbf{v} - (1 - \rho)\mathbf{v} = (\lambda - 1 + \rho)\mathbf{v}. \quad (3.27)$$

If $\rho \neq 0$, then the above calculations show that A, B share the same eigenvectors and the eigenvalues of A are $1 - \rho$ and $1 + (n - 1)\rho$. If $\rho = 0$, then $A = I_{n \times n}$, which implies that its eigenvalues are all equal to $1 = 1 - \rho = 1 + (n - 1)\rho$. Since A is positive semi-definite, we must have $1 - \rho \geq 0$ and $1 + (n - 1)\rho \geq 0$, which implies that $-\frac{1}{n-1} \leq \rho \leq 1$. For the last part, we note that if $\rho = 1$, we can simply consider random variables X_1, X_2, X_3 for which $X_1 = X_2 = X_3$. If $\rho = -\frac{1}{3-1} = -\frac{1}{2}$, we learned from Homework 2 Problem 5 that every positive semi-definite matrix is a covariance

matrix, in the sense that we can construct a random vector $(X_1, X_2, X_3) : \Omega \rightarrow \mathbb{R}^3$ for which the covariance matrix is A . Following the construction from Homework 2 Problem 5, we first orthogonally diagonalize the matrix

$$A = \begin{pmatrix} 1 & -\frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & -\frac{1}{2} & 1 \end{pmatrix}. \quad (3.28)$$

From part c) we learned that the eigenvalues of A are either $\frac{3}{2}$ or 0. To find eigenvectors corresponding to $\lambda = 0$, we note that

$$A \sim \begin{pmatrix} 2 & -1 & -1 \\ -1 & 2 & -1 \\ -1 & -1 & 2 \end{pmatrix} \sim \begin{pmatrix} 1 & -2 & 1 \\ 0 & -3 & 3 \\ 0 & 3 & -3 \end{pmatrix} \sim \begin{pmatrix} 1 & -2 & 1 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{pmatrix} \sim \begin{pmatrix} 1 & 0 & -1 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{pmatrix} \quad (3.29)$$

Therefore a basis for $\text{Null}(A)$ is

$$\mathcal{B}_{\text{Null}(A)} = \left\{ \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \right\}. \quad (3.30)$$

Note that since the row sums of A are all 0, the vector $\mathbf{v}_1 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$ is guaranteed to be a corresponding eigenvector; row reducing A shows that this is the only eigenvector corresponding to $\lambda = 0$.

Next we note that

$$A - \frac{3}{2}I_{3 \times 3} = \begin{pmatrix} -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \end{pmatrix} \sim \begin{pmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}. \quad (3.31)$$

This implies that a basis for $\text{Null}(A - \frac{3}{2}I_{3 \times 3})$ is given by

$$\mathcal{B}_{\text{Null}(A - \frac{3}{2}I_{3 \times 3})} = \left\{ \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \right\}. \quad (3.32)$$

To construct an orthogonal basis for this subspace, we perform Gram-Schmidt:

$$\mathbf{v}_2 = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix}, \quad \mathbf{q}_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \quad (3.33)$$

$$\mathbf{v}_3 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} - \frac{1}{2} \left(\begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \right) \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \quad (3.34)$$

$$= \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} - \frac{1}{2} \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 1 \\ 1 \\ -2 \end{pmatrix}, \quad \mathbf{q}_3 = \frac{1}{\sqrt{6}} \begin{pmatrix} 1 \\ 1 \\ -2 \end{pmatrix}. \quad (3.35)$$

This shows that $A = Q\Lambda Q^T$ for

$$Q = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} \\ 0 & -\frac{2}{\sqrt{6}} & \frac{1}{\sqrt{3}} \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \frac{3}{2} & 0 & 0 \\ 0 & \frac{3}{2} & 0 \\ 0 & 0 & 0 \end{pmatrix}. \quad (3.36)$$

Let $\mathbf{Z} : \Omega \rightarrow \mathbb{R}^3$ be a random vector with mean $\mathbf{0}$ and covariance matrix I . Then from Homework 2 Problem 5, it follows that the random vector $\mathbf{X} = Q\Lambda^{1/2}\mathbf{Z}$ has mean $\mathbf{0}$ and covariance matrix A . \square

Problem 3.4 (Row stochastic matrices). Later we will encounter a class of matrices referred to as *stochastic matrices* or sometimes *Markov matrices* when we discuss finite-state Markov chains. Stochastic matrices in this context are used to model the transition probabilities of a discrete dynamical system.

A matrix $M \in \mathcal{M}_{n \times n}(\mathbb{R})$ is said to be a row stochastic matrix if all of its entries are non-negative and the sum of the entries in each row is equal to 1.

- Translate the definition above to the following: show that $M \in \mathcal{M}_{n \times n}(\mathbb{R})$ is row stochastic iff. all the entries of M are non-negative and $M\mathbf{1}_{n \times 1} = \mathbf{1}_{n \times 1}$, where $\mathbf{1}_{n \times 1} \in \mathbb{R}^n$ is the column vector with all entries equal to 1.
- Show that if $M_1, \dots, M_k \in \mathcal{M}_{n \times n}(\mathbb{R})$ are row stochastic matrices, then the product $\prod_{i=1}^k M_i = M_1 \cdots M_k \in \mathcal{M}_{n \times n}(\mathbb{R})$ is also a row stochastic matrix.

Solution. We note that for any matrix $M \in \mathcal{M}_{n \times n}(\mathbb{R})$, if we write M in terms of its rows as

$$M = \begin{pmatrix} \mathbf{m}_1^T \\ \vdots \\ \mathbf{m}_n^T \end{pmatrix}, \quad (3.37)$$

then

$$M\mathbf{1}_{n \times 1} = \begin{pmatrix} \mathbf{m}_1^T \mathbf{1}_{n \times 1} \\ \vdots \\ \mathbf{m}_n^T \mathbf{1}_{n \times 1} \end{pmatrix}, \quad (3.38)$$

where $\mathbf{m}_i^T \mathbf{1}_{n \times 1}$ is the sum of the entries in the i th row of M . Therefore it follows that M is row stochastic iff. all the entries of M are non-negative and $M\mathbf{1}_{n \times 1} = \mathbf{1}_{n \times 1}$.

For part b), we note that it suffices to prove this for $k = 2$, the general case follows from a simple induction argument. Suppose $M_1, M_2 \in \mathcal{M}_{n \times n}(\mathbb{R})$ are row stochastic matrices. Then we have

$$(M_1 M_2)\mathbf{1}_{n \times 1} = M_1(M_2\mathbf{1}_{n \times 1}) = M_1\mathbf{1}_{n \times 1} = \mathbf{1}_{n \times 1}. \quad (3.39)$$

This shows that $M_1 M_2$ is row stochastic. \square

Problem 3.5 (Linear regression and cloning datasets).

Suppose one is working on a dataset with m samples, p features, and 1 target and sets up a linear regression model with a design matrix $X \in \mathcal{M}_{m \times (p+1)}(\mathbb{R})$, a target variable $\mathbf{y} \in \mathbb{R}^m$ and tries to solve for the least squares regressor $\hat{\beta} \in \mathbb{R}^{p+1}$. After solving for the least squares regressor $\hat{\beta}$, they then decided to “clone the data” and run the regression again to see if anything changes. For example, if the original dataset had 3 samples with one target and one predictor, then the cloned dataset would have 6 samples:

x	y
0	2
1	2
2	8

x	y
0	2
1	2
2	8
0	2
1	2
2	8

FIGURE 1. Original dataset on the left vs doubled dataset on the right

In general with $X \in \mathcal{M}_{m \times (p+1)}(\mathbb{R})$, $\mathbf{y} \in \mathbb{R}^m$, this means that instead of looking for the least squares estimator $\hat{\mathbf{x}} \in \mathbb{R}^{p+1}$ to $X\beta = \mathbf{y}$, they instead try to look for the least squares estimator $\hat{\beta}_2 \in \mathbb{R}^{p+1}$ for $X_2\beta = \mathbf{y}_2$, where

$$X = \begin{pmatrix} X \end{pmatrix} \in \mathcal{M}_{m \times (p+1)}(\mathbb{R}), \quad X_2 = \begin{pmatrix} X \\ X \end{pmatrix} \in \mathcal{M}_{2m \times (p+1)}(\mathbb{R}), \quad \mathbf{y} = \begin{pmatrix} \mathbf{y} \end{pmatrix} \in \mathbb{R}^m, \quad \mathbf{y}_2 = \begin{pmatrix} \mathbf{y} \\ \mathbf{y} \end{pmatrix} \in \mathbb{R}^{2m}. \quad (3.40)$$

For example, for the dataset in Figure 1, under the standard simple linear regression model $\mathbf{y} = X\beta + \varepsilon$ where we assume the random vector $\varepsilon \sim N(0, \sigma^2 I)$, we would set up

$$X = \begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{pmatrix}, \quad X_2 = \begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \\ 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 2 \\ 2 \\ 8 \end{pmatrix}, \quad \mathbf{y}_2 = \begin{pmatrix} 2 \\ 2 \\ 8 \\ 2 \\ 2 \\ 8 \end{pmatrix}. \quad (3.41)$$

Intuitively, nothing should really change because no new information has been added to the original dataset. Is this true? Let's investigate.

- Explain briefly why $\text{Col}(X) \neq \text{Col}(X_2)$ yet the dimensions of $\text{Col}(X)$ and $\text{Col}(X_2)$ are the same.
- Show that $\text{Null}(X) = \text{Null}(X_2)$.
- Assuming that the original X matrix has independent columns, show that the unique least squares solution $\hat{\beta}$ solving $X^T X \hat{\beta} = X^T \mathbf{y}$ is the same as the unique least squares regressor solving $(X_2)^T X_2 \hat{\beta} = (X_2)^T \mathbf{y}_2$, however in the cloned system the $\text{RSS}(\hat{\beta})$ is larger by a factor of 2. What about the mean squared error (MSE) $\frac{\text{RSS}(\hat{\beta})}{m}$ in the original and cloned systems? Use this observation to explain why the MSE might be preferred over RSS as a more reliable metric in practice.
- What about the coefficient of determination R^2 , the unbiased estimator $\hat{\sigma}^2 = \frac{\text{RSS}(\hat{\beta})}{m-p-1}$ for σ^2 , the standard errors $\text{SE}(\hat{\beta}_i)$, and also the t -statistics $t_i = \frac{\hat{\beta}_i - \beta_i}{\text{SE}(\hat{\beta}_i)}$ for $1 \leq i \leq 2$? What would happen to the confidence intervals for β_i if we were to use the t -statistics from the cloned system?

Solution.

- a) We note that $\text{Col}(X) \subseteq \mathbb{R}^m$ whereas $\text{Col}(X_2) \subseteq \mathbb{R}^{2m}$, therefore $\text{Col}(X) \neq \text{Col}(X_2)$. However,

$$\text{rref}(X_2) = \begin{pmatrix} \boxed{\text{rref}(X)} \\ \boxed{0_{m \times n}} \end{pmatrix}, \quad (3.42)$$

so X and X_2 admit the same number of pivots, therefore $\dim \text{Col}(X) = \dim \text{Col}(X_2)$.

- b) We note that by (3.42), $\text{rref}(X)\mathbf{x} = \mathbf{0}_m$ if and only if $\text{rref}(X_2)\mathbf{x} = \mathbf{0}_{2m}$ since adding an equation of the form $0x_1 + \dots + 0x_n = 0$ does not change the underlying solution set. Therefore $\text{Null}(X) = \text{Null}(\text{rref}(X)) = \text{Null}(\text{rref}(X_2)) = \text{Null}(X_2)$.

- c) We note that by block multiplication,

$$(X_2)^T X_2 = \begin{pmatrix} \boxed{X^T} & \boxed{X^T} \end{pmatrix} \begin{pmatrix} \boxed{X} \\ \boxed{X} \end{pmatrix} = X^T X + X^T X = 2X^T X, \quad (3.43)$$

and

$$(X_2)^T \mathbf{y}_2 = \begin{pmatrix} \boxed{X^T} & \boxed{X^T} \end{pmatrix} \begin{pmatrix} \boxed{\mathbf{y}} \\ \boxed{\mathbf{y}} \end{pmatrix} = X^T \mathbf{y} + X^T \mathbf{y} = 2X^T \mathbf{y}. \quad (3.44)$$

Therefore $X^T X \hat{\mathbf{x}} = X^T \mathbf{y} \iff (X_2)^T X_2 \hat{\mathbf{x}} = (X_2)^T \mathbf{y}_2$. However,

$$\mathbf{y}_2 - X_2 \hat{\mathbf{x}} = \begin{pmatrix} \boxed{\mathbf{y}} \\ \boxed{\mathbf{y}} \end{pmatrix} - \begin{pmatrix} \boxed{X} \\ \boxed{X} \end{pmatrix} \hat{\mathbf{x}} = \begin{pmatrix} \boxed{\mathbf{y} - X \hat{\mathbf{x}}} \\ \boxed{\mathbf{y} - X \hat{\mathbf{x}}} \end{pmatrix} \implies \|\mathbf{y}_2 - X_2 \hat{\mathbf{x}}\|^2 = 2 \|\mathbf{y} - X \hat{\mathbf{x}}\|^2, \quad (3.45)$$

so the error increases by a factor of 2 in the cloned system. However, if we consider the mean squared error, since there are $2m$ samples in the cloned system and m samples in the original system, we have

$$\frac{1}{2m} \|\mathbf{y}_2 - X_2 \hat{\mathbf{x}}\|^2 = \frac{2}{2m} \|\mathbf{y} - X \hat{\mathbf{x}}\|^2 = \frac{1}{m} \|\mathbf{y} - X \hat{\mathbf{x}}\|^2, \quad (3.46)$$

therefore the mean squared error does not change. In practice, this means that the MSE is a more robust measure of error, as it is not affected by artificial duplication of data.

- d) We note that the coefficient of determination R^2 is defined via

$$R^2 = 1 - \frac{2\text{RSS}_{\text{original}}(\hat{\beta})}{2\text{TSS}_{\text{original}}}, \quad (3.47)$$

Since the scaling factor is the same in both the RSS and TSS, the coefficient of determination does not change. The unbiased estimator for σ^2 in the cloned system is $\frac{2\text{RSS}_{\text{original}}}{2m-p-1}$, which means that it is roughly unchanged if $p-1$ is small.

For the standard errors, recall that $\hat{\beta} \sim N(\beta, \sigma^2(X^T X)^{-1}) \approx N(\beta, \hat{\sigma}^2(X^T X)^{-1})$. For $p = 1$, we have

$$\hat{\sigma}^2(X^T X)^{-1} = \frac{\hat{\sigma}^2}{(m-1)s_x^2} \begin{pmatrix} \frac{1}{m} \sum_{i=1}^m x_i^2 & -\frac{1}{m} \sum_{i=1}^m x_i \\ -\frac{1}{m} \sum_{i=1}^m x_i & 1 \end{pmatrix} \quad (3.48)$$

In the cloned system, s_x^2 and $\hat{\sigma}^2$ are roughly the same, however the scaling factor $\frac{1}{m-1}$ is replaced by $\frac{1}{2m-1}$, which means the standard errors in the cloned system get scaled down by roughly a factor of $\sqrt{2}$. As a result, the t -statistics in the cloned system scales up by roughly a factor of $\sqrt{2}$, and the confidence intervals for β_i in the cloned system shrinks by roughly a factor of $\sqrt{2}$.

For general p , the scaling factors for $\hat{\sigma}^2$, the t -statistics, and the standard errors will also be affected by the value of p , however tracking this explicitly is more difficult as it depends on $(X^T X)^{-1}$.

□

Problem 3.6 (Characteristic functions and affine transformations of Gaussian random variables). Let $\mathbf{X} : \Omega \rightarrow \mathbb{R}^n$ be a random vector. The *characteristic function* of \mathbf{X} is defined as the function $\phi_{\mathbf{X}} : \mathbb{R}^n \rightarrow \mathbb{C}$ defined via

$$\phi_{\mathbf{X}}(\mathbf{t}) = \mathbb{E}[e^{i\mathbf{t} \cdot \mathbf{X}}] = \int_{\mathbb{R}^n} e^{i\mathbf{t} \cdot \mathbf{x}} f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}. \quad (3.49)$$

For example, if $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$, then the characteristic function $\phi_{\mathbf{X}} : \mathbb{R}^n \rightarrow \mathbb{C}$ of \mathbf{X} is given by

$$\phi_{\mathbf{X}}(\mathbf{t}) = \exp\left(i\boldsymbol{\mu} \cdot \mathbf{t} - \frac{1}{2}\mathbf{t}^T \Sigma \mathbf{t}\right). \quad (3.50)$$

If $n = 1$, then this says that if $X \sim N(\mu, \sigma^2)$, then the characteristic function $\phi_X : \mathbb{R} \rightarrow \mathbb{C}$ of X is given by

$$\phi_X(t) = e^{i\mu t - \frac{1}{2}\sigma^2 t^2}. \quad (3.51)$$

The characteristic function is a powerful tool as it allows one to prove many results with ease. Below are two properties of the characteristic function that we will use in this problem.

- If $\mathbf{X}_1, \dots, \mathbf{X}_n : \Omega \rightarrow \mathbb{R}^n$ are independent random vectors, then the characteristic function of their sum is the product of their individual characteristic functions.
- The characteristic function of a random vector uniquely determines its distribution: if \mathbf{X} and \mathbf{Y} are random vectors such that $\phi_{\mathbf{X}} = \phi_{\mathbf{Y}}$, then $F_{\mathbf{X}} = F_{\mathbf{Y}}$, where $F_{\mathbf{X}}$ and $F_{\mathbf{Y}}$ are the cumulative distribution functions of \mathbf{X} and \mathbf{Y} , respectively.

These essentially follows from properties of the *Fourier transform*, which we will explore later in the course.

- Show that if X_1, \dots, X_n are independent random variables and $X_i \sim N(\mu_i, \sigma_i^2)$, then their sum $X = X_1 + \dots + X_n$ is also a Gaussian random variable with mean $\mu = \mu_1 + \dots + \mu_n$ and variance $\sigma^2 = \sigma_1^2 + \dots + \sigma_n^2$. What is the exact distribution of their mean $\bar{X} = \frac{X}{n}$?
- Suppose $\mathbf{X} : \Omega \rightarrow \mathbb{R}^n$ is a random vector and $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$. Let $A \in \mathcal{M}_{n \times n}(\mathbb{R})$ and $\mathbf{b} \in \mathbb{R}^n$. Show that the random vector $A\mathbf{X} + \mathbf{b}$ is also Gaussian with mean $A\boldsymbol{\mu} + \mathbf{b}$ and covariance matrix $A\Sigma A^T$.

Solution. We note that the characteristic function of X is $\phi_X : \mathbb{R} \rightarrow \mathbb{C}$ given by

$$\begin{aligned} \phi_X(t) &= \mathbb{E}[e^{itX}] = \mathbb{E}[e^{it(X_1 + \dots + X_n)}] = \mathbb{E}[e^{itX_1} \dots e^{itX_n}] = \mathbb{E}[e^{itX_1}] \dots \mathbb{E}[e^{itX_n}] \\ &= \prod_{i=1}^n \phi_{X_i}(t) = \prod_{i=1}^n e^{i\mu_i t - \frac{1}{2}\sigma_i^2 t^2} = e^{i\mu t - \frac{1}{2}\sigma^2 t^2}, \end{aligned} \quad (3.52)$$

where

$$\mu = \mu_1 + \dots + \mu_n, \quad \sigma^2 = \sigma_1^2 + \dots + \sigma_n^2. \quad (3.53)$$

Since the characteristic function determines the distribution, we have that $X \sim N(\mu, \sigma^2)$. The distribution of the mean $\bar{X} = \frac{X}{n}$ is then $N\left(\frac{\mu}{n}, \frac{\sigma^2}{n^2}\right)$.

For part b), we note that the characteristic function of $A\mathbf{X} + \mathbf{b}$ is $\phi_{A\mathbf{X} + \mathbf{b}} : \mathbb{R}^n \rightarrow \mathbb{C}$ given by

$$\begin{aligned} \phi_{A\mathbf{X} + \mathbf{b}}(\mathbf{t}) &= \mathbb{E}[\exp(i\mathbf{t} \cdot (A\mathbf{X} + \mathbf{b}))] = \mathbb{E}[\exp(i\mathbf{t} \cdot A\mathbf{X}) \exp(i\mathbf{t} \cdot \mathbf{b})] \\ &= \exp(i\mathbf{t} \cdot \mathbf{b}) \mathbb{E}[\exp(i\mathbf{t} \cdot A\mathbf{X})] = \exp(i\mathbf{t} \cdot \mathbf{b}) \mathbb{E}[\exp(iA^T \mathbf{t} \cdot \mathbf{X})] = \exp(i\mathbf{t} \cdot \mathbf{b}) \phi_{\mathbf{X}}(A^T \mathbf{t}) \\ &= \exp(i\mathbf{t} \cdot \mathbf{b}) \exp\left(i\boldsymbol{\mu} \cdot A^T \mathbf{t} - \frac{1}{2}(A^T \mathbf{t})^T \Sigma (A^T \mathbf{t})\right) = \exp\left(i(A\boldsymbol{\mu} + \mathbf{b}) \cdot \mathbf{t} - \frac{1}{2}\mathbf{t}^T (A\Sigma A^T) \mathbf{t}\right). \end{aligned} \quad (3.54)$$

Since the characteristic function determines the distribution, it follows that $A\mathbf{X} + \mathbf{b} \sim N(A\boldsymbol{\mu} + \mathbf{b}, A\Sigma A^T)$. \square