

This exercise sheet contains problems from Parts II and III of the course.  
 Please submit the exercises highlighted in **red** before the deadline.

### PROBLEM 1

- (a) Let  $X$  be a random variable with  $E[X] = 0$ . Suppose that the moment-generating function of  $X^2$  is bounded at some point, that is,

$$E[e^{X^2}] \leq 2.$$

Prove that  $X$  satisfies the two-sided tail bound

$$P(|X| > t) \leq 2e^{(-t^2)} \text{ for all } t \geq 0.$$

- (b) Prove that if  $X$  is a non-negative random variable with expectation  $E[X]$ , then for all  $t > 0$ , we have  $P[X \geq t] \leq E[X]/t$ .
- (c) Recall Chernoff's inequality: Let  $X_i$  be independent Bernoulli random variables with success probability  $p_i$ . Consider their sum  $S_N = \sum_{i=1}^N X_i$  and denote its mean by  $\mu = E[S_N]$ . Then, for any  $t > \mu$ , we have

$$P(S_N \geq t) \leq e^{t-\mu} \left(\frac{\mu}{t}\right)^t.$$

Consider 200 independent coin flips. We wish to find an upper bound on the probability that the number of heads is greater or equal than 150. Use Chernoff's inequality.

- (d) Let  $X_i$ , for  $i = 1, \dots, n$ , be a random sample of a random variable  $X$ . Let  $X$  have mean  $\mu$  and variance  $\sigma^2$ . Find the size of the sample ( $n$ ), such that the probability that the difference between sample mean and true mean is smaller than  $\frac{\sigma}{10}$  is at least 0.95. Hint: Derive a version of the Chebyshev inequality for  $P(|X - \mu| \geq a)$  using Markov inequality.

### Solution.

- (a) It holds that

$$P(|X| > t) = P(e^{X^2} > e^{t^2}) \leq \frac{2}{e^{t^2}}.$$

(we just apply Markov in the last step)

- (b) It holds that

$$P(X \geq t) = P(t_{X \geq t} \geq t) = E[X_{X \geq t}] \leq E\left[\frac{X}{t}\right] = \frac{E[X]}{t}.$$

Alternative way:

$$\begin{aligned} E[X] &= \int_{-\infty}^{\infty} xf(x)dx = \int_0^{\infty} xf(x)dx \geq \int_t^{\infty} xf(x)dx \geq \int_t^{\infty} tf(x)dx = tP(X \geq t) \\ \implies E[X] &\geq tP(X \geq t) \implies \frac{E[X]}{t} \geq P(X \geq t) \end{aligned}$$

- (c) Chernoff gives  $e^{50} \left(\frac{2}{3}\right)^{150} = \left(\frac{8e}{27}\right)^{50}$ . It is not necessary to simplify this further.

- (d) Let  $\hat{X} = \frac{1}{n} \sum_{i=1}^n X_i$ . Then,  $E[\hat{X}] = \mu$  and  $\text{Var}[\hat{X}] = \frac{\sigma^2}{n}$ . Now, we need to determine  $n$  such that

$$P(|\hat{X} - \mu| \leq \frac{\sigma}{10}) \geq 0.95 \implies P(|\hat{X} - \mu| \geq \frac{\sigma}{10}) \leq 0.05$$

We can write the probability as:

$$\begin{aligned} P(\sqrt{(\hat{X} - \mu)^2} \geq \frac{\sigma}{10}) &= P((\hat{X} - \mu)^2 \geq \frac{\sigma^2}{100}) \leq \frac{\text{Var}[\hat{X}]}{\frac{\sigma^2}{100}} = \frac{\sigma^2}{n} \frac{100}{\sigma^2} = \frac{100}{n} \leq 0.05 \\ &\implies \frac{100}{0.05} \leq n \end{aligned}$$

Therefore, we need a sample size of  $n \geq 2000$ .

## PROBLEM 2

1. Estimation of diagonal covariances: Let  $(X_i)_{i=1,\dots,n}$  be an i.i.d. sequence of  $d$ -dimensional vectors, drawn from a zero-mean distribution with diagonal covariance matrix  $\Sigma = D$ . Consider the estimate  $\hat{D} = \text{diag}(\hat{\Sigma})$ , where  $\hat{\Sigma}$  is the usual sample covariance matrix. Suppose further that each component  $X_{ij}$  is sub-Gaussian with parameter at most  $\sigma = 1$ . Show the following:

- (a)  $X_{ij}^2$  is sub-exponential with parameters  $(2, 4)$ .
- (b)  $\sum_{i=1}^n X_{ij}^2$  is sub-exponential with parameters  $(2\sqrt{n}, 4)$
- (c) For each  $i = 1, \dots, d$ , we get

$$P\left(|\hat{D}_{ii} - D_{ii}| \geq t\right) \leq 2e^{-\frac{n}{8} \min\{t, t^2\}}.$$

2. Suppose that the random vector  $X \in \mathbb{R}^n$  has a  $N_n(\mu, \Sigma)$  distribution, where  $\Sigma$  is positive. Show the the random variable  $Y = (X - \mu)^T \Sigma^{-1} (X - \mu)$  is sub-exponential.

## Solution.

1. (a)  $X_{ij}^2$  is sub-exponential with parameters  $(2, 4)$ .

**Approach 1** For this, we consider that a sub-gaussian variable of parameter at most  $\sigma$  will be bounded for above by a gaussian variable. Let  $X \sim N(0, \sigma^2)$ , and further assume  $\sigma = 1$ . Now, consider that  $X^2$  follows a chi-squared distribution, and its moment generating function is defined as.

$$E\left[e^{\lambda(X^2-1)}\right] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{\lambda x^2} e^{-\frac{x^2}{2}} dx = \frac{e^{-\lambda}}{\sqrt{1-2\lambda}}, \text{ for } \lambda < \frac{1}{2}$$

The moment generating function is also obtain by using the gaussian distribution and considering  $E[X^2] = 1$ . Following the definition of sub-exponential we have:

$$E\left[e^{\lambda(X^2-1)}\right] \leq e^{\frac{\nu^2 \lambda^2}{2}} \text{ for all } \lambda^2 < \frac{1}{\alpha^2}$$

Now, considering  $\nu = 2$ , and  $\alpha = 4$ , we have that  $\lambda^2 < \frac{1}{16} \implies \lambda \in (-\frac{1}{4}, \frac{1}{4})$ . Therefore, the moment generation function previously calculated is bounded for these values of  $\lambda$ . With this it should hold that

$$\frac{e^{-\lambda}}{\sqrt{1-2\lambda}} \leq e^{\frac{2^2 \lambda^2}{2}} = e^{2\lambda^2} \text{ for all } \lambda^2 < \frac{1}{16}$$

Let's focus for  $\lambda \in (-\frac{1}{4}, \frac{1}{4})$ . Given that all terms are positive, we can square and reorder the inequality:

$$e^{-4\lambda^2-2\lambda} \leq 1 - 2\lambda$$

$$\implies -4\lambda^2 - 2\lambda \leq \log 1 - 2\lambda \implies 0 \leq \ln(1 - 2\lambda) + 4\lambda^2 + 2\lambda = f(\lambda)$$

It is easy to show that  $f(x)$  is a convex function in the domain of  $\lambda$ . Therefore we can calculate its minimum with first and second order condition. If  $\min f(\lambda) \geq 0$  for  $|\lambda| < \frac{1}{4}$ , the inequality holds and the variable is sub-exponential(2,4).

$$FOC : \frac{df(\lambda)}{d\lambda} = -\frac{2}{1-2\lambda} + 8\lambda + 2 = -2 + 8\lambda - 16\lambda^2 + 2 - 4\lambda = 0$$

$$4\lambda(1 - 4\lambda) = 0 \implies \lambda = 0 \vee \lambda = 1/4, \text{ we can see that } 1/4 \text{ is not a minimizer.}$$

The second derivative evaluated in  $\lambda = 0$  has a value of 4, therefore  $\lambda = 0$  is a proper minimizer of  $f(\lambda)$ , and  $f(0) = 0$ . Thus the inequality holds and the variable  $X^2$  is sub-exponential of parameter (2,4).

**Possible Alternative** Another way to approach the problem will be: Let  $Z = X^2 - E[X^2]$ . Then, we calculated its moment generation function using the Taylor expansion of the exponential,

$$E[e^{\lambda Z}] \leq E\left[1 + \sum_{k=1}^{\infty} \frac{\lambda^k Z^k}{k!}\right]$$

Following this, we can keep bounding using Jensen's inequality and the bounds available given that  $X$  is sub-gaussian the parameter at most 1.

- (b) Let  $Z_{ij} = X_{ij}^2$ , and therefore be sub-exponential with parameters (2,4). Now we compute the moment generating function:

$$E\left[e^{\lambda \sum_{i=1}^n (Z_{ij} - E[Z_{ij}])}\right] = \prod_{i=1}^n E\left[e^{\lambda (Z_{ij} - E[Z_{ij}])}\right]$$

Now, following the bounds obtained beforehand:

$$\begin{aligned} \prod_{i=1}^n E\left[e^{\lambda (Z_{ij} - E[Z_{ij}])}\right] &\leq \prod_{i=1}^n e^{v^2 \frac{\lambda^2}{2}} = e^{\sum_{i=1}^n (v^2 \frac{\lambda^2}{2})} \quad \forall |\lambda| \leq \frac{1}{4} \\ \implies E\left[e^{\lambda \sum_{i=1}^n (Z_{ij} - E[Z_{ij}])}\right] &\leq e^{(\sqrt{nv})^2 \frac{\lambda^2}{2}} \quad \forall |\lambda| \leq \frac{1}{4} \end{aligned}$$

Finally, we can conclude that  $\sum_{i=1}^n X_{ij}^2$  is sub-exponential with parameters  $(2\sqrt{n}, 4)$ .

- (c)  $\hat{D}_{ii}$  is the usual sample covariance matrix, and it's defined as  $\hat{D}_{ii} = \frac{1}{n} \sum_{i=1}^n x_{ij}^2$ . This estimator is unbiased, i.e.,  $E[\hat{D}_{ii}] = D$ . Also, following the previous exercises we have that  $\hat{D}_{ii}$  is sub-exponential de parameters  $(\frac{2}{\sqrt{n}}, \frac{4}{n})$ . Now, given sub-exponential concentration

$$P\left(|\hat{D}_{ii} - D_{ii}| \geq t\right) \leq 2 \exp^{-\frac{1}{2} \min\left\{\frac{t}{\alpha'}, \frac{t^2}{v^2}\right\}}.$$

Replacing  $nu$  and  $\alpha$  for their respective values, we get:

$$P\left(|\hat{D}_{ii} - D_{ii}| \geq t\right) \leq 2e^{-\frac{n}{8} \min\{t, t^2\}}.$$

2.  $Y = (X - \mu)^T \Sigma^{-1} (X - \mu)$ . Let's consider the spectral decomposition of  $\Sigma = Q\Lambda Q^T$ , where  $Q^T Q = I$ . Now,  $\Sigma^{-\frac{1}{2}}$  is then defined as  $Q\Lambda^{\frac{1}{2}}Q^T$ , and therefore  $\Sigma^{-1} = \Sigma^{-\frac{1}{2}}\Sigma^{-\frac{1}{2}}$ . Let  $Z = \Sigma^{-\frac{1}{2}}(X - \mu)$ , and corresponds to random variable that follows a normal standard distribution. Then,

$$Y = Z^T Z = \sum_i Z_i^2$$

Therefore, as presented in the lecture,  $Y \sim \chi^2(n)$  is a sub-exponential variable.

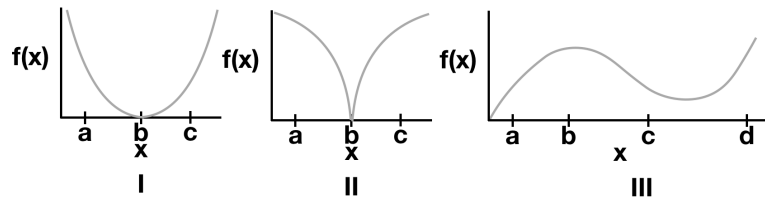
### PROBLEM 3

Convexity and norms: A norm  $\|\cdot\|$  over  $\mathbb{R}^n$  is defined by the properties:

- (i) non-negativity:  $\|x\| \geq 0$  for all  $x \in \mathbb{R}^n$  with equality if and only if  $x = 0$ ,
  - (ii) absolute scalability:  $\|ax\| = |a| \|x\|$  for all  $a \in \mathbb{R}$  and  $x \in \mathbb{R}^n$ ,
  - (iii) triangle inequality:  $\|x + y\| \leq \|x\| + \|y\|$  for all  $x, y \in \mathbb{R}^n$ .
- (a) Show that  $\|x\|_1 = \sum_{i=1}^n |x_i|$  is a norm. (Hint: for (iii), begin by showing that  $|a + b| \leq |a| + |b|$  for all  $a, b \in \mathbb{R}$ .) (Correspond to the penalty for LASSO.)
- (b) Show that  $f(x) = \left(\sum_{i=1}^n |x_i|^{1/2}\right)^2$  is not a norm. (Hint: it suffices to find two points in  $n = 2$  dimensions such that the triangle inequality does not hold.)
- (c) Show that  $\|x\|_2 := \left(\sum_{i=1}^n |x_i|^2\right)^{1/2}$  is a norm. (Correspond to the penalty for ridge regression.)

We say a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex on a set  $A$  if  $f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$  for all  $x, y \in A$  and  $\lambda \in [0, 1]$ .

- (d) For each of the functions below (I-III), state whether each one is convex on the given interval or state why not with a counterexample using any of the points a, b, c, d in your answer.



- 1 Function in panel I on  $[a, c]$
  - 2 Function in panel II on  $[a, c]$
  - 3 Function in panel III on  $[a, d]$
  - 4 Function in panel III on  $[c, d]$
- (e) For  $i = 1, \dots, n$  let  $\ell_i(w)$  be convex functions over  $w \in \mathbb{R}^d$  (e.g.,  $\ell_i(w) = (y_i - w^\top x_i)^2$ ),  $\|\cdot\|$  is any norm, and  $\lambda > 0$ . Show that

$$\sum_{i=1}^n \ell_i(w) + \lambda \|w\|$$

is convex over  $w \in \mathbb{R}^d$  (Hint: Show that if  $f, g$  are convex functions, then  $f(x) + g(x)$  is also convex.)

### Solution.

- (a) We show that (i)-(iii) hold.

- (i)  $\sum_{i=1}^n |x_i| \geq 0$ ;  $x = 0 \Rightarrow \sum_{i=1}^n |x_i| = 0$ ,  $\sum_{i=1}^n |x_i| = 0 \Rightarrow x_i = 0$  for all  $i \in \{1, \dots, n\}$ .
- (ii)  $\|ax\|_1 = \sum_{i=1}^n |ax_i| = |a| \sum_{i=1}^n |x_i| = |a| \|x\|_1$ .
- (iii)  $\|x + y\|_1 = \sum_{i=1}^n |x_i + y_i| \leq \sum_{i=1}^n |x_i| + \sum_{i=1}^n |y_i| = \|x\|_1 + \|y\|_1$ .

(b) We choose the points  $x = (0, 4)^T$  and  $y = (4, 0)^T$ . For these points, we obtain

$$f(x+y) = \left( \sum_{i=1}^n |x_i + y_i|^{12} \right)^2 = (4^{12} + 4^{12})^2 = 16$$

$$f(x) + f(y) = (4^{12})^2 + (4^{12})^2 = 8.$$

Hence, the triangle inequality does not hold and  $f(\cdot)$  cannot be a norm.

(c) We show that (i)-(iii) hold.

- (i)  $x_2 \geq 0; x = 0 \Rightarrow 0_2 = \left( \sum_{i=1}^n |0|^2 \right)^{12} = 0, 0_2 = 0 \Rightarrow \left( \sum_{i=1}^n |x_i|^2 \right)^{12} = 0 \Rightarrow \sum_{i=1}^n |x_i|^2 = 0 \Rightarrow x_i = 0$  for all  $i \in \{1, \dots, n\}$ .
- (ii)  $\|ax\|_2 = \left( \sum_{i=1}^n |ax_i|^2 \right)^{12} = |a| \left( \sum_{i=1}^n |x_i|^2 \right)^{12} = |a| \|a\|_2$ .
- (iii)  $\|x+y\|_2 = \left( \sum_{i=1}^n |x_i + y_i|^2 \right)^{12} = \left( \sum_{i=1}^n x_i^2 + 2 \sum_{i=1}^n x_i y_i + \sum_{i=1}^n y_i^2 \right)^{12} = \left( \sum_{i=1}^n x_i^2 + 2x^T y + \sum_{i=1}^n y_i^2 \right)^{12} = (x_2^2 + 2x^T y + y_2^2)^{12} \leq (x_2^2 + 2x_2 y_2 + y_2^2)^{12} = ((x_2 + y_2)^2)^{12} = x_2 + y_2$ , where we used the Cauchy-Schwarz inequality.

(d) (a) convex

(b) not convex on  $[b, c]$ : Every function value in  $(b, c)$  is greater than the corresponding convex combination of the function values  $f(b)$  and  $f(c)$ .

(c) not convex on  $[a, c]$ : Every function value in  $(a, c)$  is greater than the corresponding convex combination of the function values  $f(a)$  and  $f(c)$ .

(d) convex

(e) Let  $f$  and  $g$  be two convex functions. Then, for every  $\lambda \in [0, 1]$ ,

$$\begin{aligned} (f+g)(\lambda x + (1-\lambda)y) &= f(\lambda x + (1-\lambda)y) + g(\lambda x + (1-\lambda)y) \\ &\leq \lambda f(x) + (1-\lambda)f(y) + \lambda g(x) + (1-\lambda)g(y) \\ &= \lambda(f(x) + g(x)) + (1-\lambda)(f(y) + g(y)) \\ &= \lambda(f+g)(x) + (1-\lambda)(f+g)(y), \end{aligned}$$

which shows that  $f+g$  is convex. With this, it is now sufficient to show that every norm is a convex function:

$$\lambda x + (1-\lambda)y \stackrel{(iii)}{\leq} \lambda x + (1-\lambda)y \stackrel{(ii)}{=} \lambda x + (1-\lambda)y.$$

#### PROBLEM 4

The following questions should be answered without referring to external materials. Briefly justify your answers with a few words.

- (a) How does lasso regression differ from ridge regression?
- (b) Why do least squares fail in high dimensions?
- (c) In a LASSO regression, if the regularization parameter  $\lambda$  is very high, what happens to the estimated regression coefficients?
- (d) True or False: The LASSO is a convex optimization problem.

#### Solution.

- (a) LASSO regression can achieve feature selection (by setting feature weights to exactly zero) while ridge regression cannot.
  - (b)  $X'X$  has no longer full rank when  $p > n$ . The OLS results in infinitely many solutions, leading to over-fitting in HD.
  - (c) The model can shrink the coefficients of uninformative features to exactly zero.
  - (d) True (see Problem 3).
- 

## PROBLEM 5

Ridge Regression: Consider the linear regression model

$$y = X\beta_0 + \varepsilon$$

with  $y \in \mathbb{R}^n$  and  $X \in \mathbb{R}^{n \times d}$  and  $\varepsilon \in \mathbb{R}^n$  some random noise vector. The ridge regression estimator is employed when  $\text{rk}(X'X) < d$ . It is defined for a given parameter  $\lambda > 0$  by

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^d}{\text{argmin}} \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2. \quad (1)$$

- (a) Show that for any  $\lambda > 0$  the solution to the minimization problem (1) is

$$\hat{\beta} = (X'X + \lambda I_{d \times d})^{-1} X'y.$$

You may use that  $\frac{\partial}{\partial \beta} \beta' X' X \beta = 2X' X \beta$ . Note also that you need to argue why  $X'X + \lambda I_{d \times d}$  is invertible.

- (b) Compute the bias  $E(\hat{\beta}) - \beta_0$ .

## Solution.

- (a) We take the derivative with respect to  $\beta$  and set it to zero.

$$\begin{aligned} & \frac{\partial}{\partial \beta} (y - X\beta)^2 + \lambda \beta^2 \\ &= \frac{\partial}{\partial \beta} ((y - X\beta)'(y - X\beta) + \lambda \beta' \beta) \\ &= \frac{\partial}{\partial \beta} (y'y - \beta' X'y - y' X \beta + \beta' X' X \beta + \lambda \beta' \beta) \\ &= \frac{\partial}{\partial \beta} (-2\beta' X'y + \beta' X' X \beta + \lambda \beta' \beta) \\ &= -2X'y + 2X' X \beta + 2\lambda I \beta \stackrel{!}{=} 0 \\ &\Rightarrow (X'X + \lambda I) \beta = X'y \\ &\Rightarrow \hat{\beta} = (X'X + \lambda I)^{-1} X'y. \end{aligned}$$

(b)

$$\begin{aligned}
 E(\hat{\beta}) - \beta_0 &= (X'X + \lambda I)^{-1} X' E(y) - \beta_0 \\
 &= (X'X + \lambda I)^{-1} X' E(X\beta_0 + \varepsilon) - \beta_0 \\
 &\stackrel{E(\varepsilon)=0}{=} (X'X + \lambda I)^{-1} X' X \beta_0 - \beta_0 \\
 &= (X'X + \lambda I)^{-1} (X'X - (X'X + \lambda I)) \beta_0 \\
 &= (X'X + \lambda I)^{-1} (-\lambda I) \beta_0 \\
 &= -\lambda (X'X + \lambda I)^{-1} \beta_0.
 \end{aligned}$$

## PROBLEM 6

Consider the sub-Gaussian sequence model

$$Y = \theta + \sigma \varepsilon,$$

where  $\varepsilon \in \mathbb{R}^n$  consists of independent mean-zero 1-sub-Gaussian components,  $\theta \in \mathbb{R}^n$  and  $\sigma > 0$ . The soft-thresholding operator, defined for  $v \in \mathbb{R}$  by

$$S_\lambda(v) = \begin{cases} v - \lambda, & v > \lambda, \\ 0, & |v| \leq \lambda, \\ v + \lambda, & v < -\lambda \end{cases}$$

gives the soft-thresholding estimator (when applied elementwise)  $\hat{\theta} := S_\lambda(Y)$ . We suppose further that  $\theta$  is  $s$ -sparse, meaning that  $\|\theta\|_0 = \sum_{j=1}^n 1_{\{\theta_j \neq 0\}} = s$ .

(a) Show that if  $\lambda \geq \sigma \|\varepsilon\|_{\max}$ , then

$$\|\hat{\theta} - \theta\|_2^2 \leq 4s\lambda^2$$

(b) Show that

$$\mathbb{P} \left( \|\hat{\theta} - \theta\|_2 > 2\sqrt{s}\lambda \right) \leq \frac{1}{2n}$$

$$\text{for } \lambda = 2\sqrt{\sigma^2 \log(2n)}.$$

## Solution.

(a)

$$\begin{aligned}
 \hat{\theta}_j - \theta_j &= S_\lambda(Y_j) - \theta_j = \begin{cases} Y_j - \lambda - \theta_j & Y_j > \lambda \\ -\theta_j & Y_j \in [-\lambda, \lambda] \\ Y_j + \lambda - \theta_j & Y_j < -\lambda \end{cases} \\
 &= \begin{cases} \sigma \varepsilon_j - \lambda & Y_j > \lambda \\ -\theta_j & Y_j \in [-\lambda, \lambda] \\ \sigma \varepsilon_j + \lambda & Y_j < -\lambda \end{cases}
 \end{aligned}$$

$$\begin{aligned}\hat{\theta} - \theta_2^2 &= \sum_{j=1}^n \left( (\sigma\varepsilon_j - \lambda) \mathbb{1}_{\{\sigma\varepsilon_j + \theta_j > \lambda\}} - \theta_j \mathbb{1}_{\{\sigma\varepsilon_j + \theta_j \in [-\lambda, \lambda]\}} + (\sigma\varepsilon_j + \lambda) \mathbb{1}_{\{\sigma\varepsilon_j + \theta_j < -\lambda\}} \right)^2 \\ &= \sum_{j=1}^n \left( (\sigma\varepsilon_j - \lambda)^2 \mathbb{1}_{\{\sigma\varepsilon_j + \theta_j > \lambda\}} + \theta_j^2 \mathbb{1}_{\{\sigma\varepsilon_j + \theta_j \in [-\lambda, \lambda]\}} + (\sigma\varepsilon_j + \lambda)^2 \mathbb{1}_{\{\sigma\varepsilon_j + \theta_j < -\lambda\}} \right),\end{aligned}$$

where for  $\theta_j = 0$ , the summand is equal to 0, since  $\sigma\varepsilon_j \leq \max |\sigma\varepsilon_j| \leq \lambda$  by assumption, so that, assuming that the  $\theta_j \neq 0$  for the first  $s$  indices, we obtain

$$= \sum_{j=1}^s 4\lambda^2 = s4\lambda^2,$$

since

$$\begin{aligned}(\sigma\varepsilon_j - \lambda)^2 &= \sigma^2\varepsilon_j^2 + \lambda^2 - 2\lambda\sigma\varepsilon_j \\ &\leq \sigma^2\varepsilon^2 + \lambda^2 + 2\lambda\sigma\varepsilon \\ &\leq \lambda^2 + \lambda^2 + 2\lambda^2 = 4\lambda^2.\end{aligned}$$

The computation works analogously for  $(\sigma\varepsilon_j + \lambda)^2$ .

(b)

$$\begin{aligned}P(\hat{\theta} - \theta_2 > 2\lambda\sqrt{s}) &= P(\hat{\theta} - \theta_2 > 4\lambda^2 s) \\ &\leq P(\{\lambda < \sigma|\varepsilon_1|\} \cup \dots \cup \{\lambda < \sigma|\varepsilon_n|\}) \\ &\leq \sum_{i=1}^n P(\lambda < \sigma|\varepsilon_i|) \\ &= \sum_{i=1}^n P(|\varepsilon_i| > \frac{1}{\sigma}) \\ &\leq \sum_{i=1}^n 2e^{-(\frac{1}{\sigma})^2 2} \\ &= n2e^{-(\frac{1}{\sigma})^2 2} = e^{\log(2n) - (\frac{1}{\sigma})^2 2},\end{aligned}$$

where the first inequality holds due to a):  $\lambda \geq \sigma\varepsilon_{\max} \Rightarrow \hat{\theta} - \theta_2^2 \leq 4s\lambda^2$  or  $\{\lambda \geq \sigma\varepsilon_{\max}\} \subset \{\hat{\theta} - \theta_2^2 \leq 4s\lambda^2\}$ . Note that with  $\{\lambda \geq \sigma\varepsilon_{\max}\} = \{\lambda \geq \sigma|\varepsilon_1|\} \cap \dots \cap \{\lambda \geq \sigma|\varepsilon_n|\}$ . it is also true that  $\{\hat{\theta} - \theta_2^2 \leq 4s\lambda^2\}^C \subset \{\lambda \geq \sigma\varepsilon_{\max}\}^C = (\{\lambda \geq \sigma|\varepsilon_1|\} \cap \dots \cap \{\lambda \geq \sigma|\varepsilon_n|\})^C = \{\lambda < \sigma|\varepsilon_1|\} \cup \dots \cup \{\lambda < \sigma|\varepsilon_n|\}$ .

We obtain

$$P(\hat{\theta} - \theta_2 > 2\sqrt{s}\sqrt{\log(2n)\sigma^2}) \leq e^{\log(2n) - 2\log(2n)} = e^{-\log(2n)} = \frac{1}{2n}.$$

## PROBLEM 7

For the orthogonal case, i.e. when  $X'X = I$ , derive the following explicit forms for estimators,



(a) For ridge:

$$\hat{\beta}^{Ridge} = \hat{\beta}^{OLS} / (1 + \lambda).$$

(b) For lasso:

$$\hat{\beta}_i^{Lasso} = \text{sign}(\hat{\beta}_i^{OLS}) (|\hat{\beta}_i^{OLS}| - \lambda)_+,$$

where  $\hat{\beta}^{OLS}$  is the regular OLS estimator and  $\hat{\beta}_i^{OLS}$  its  $i$ th component. Note that the results can differ depending on how one chooses the multiplicative constants. The solutions in Problem 7 are based on the following objective functions:

$$\begin{aligned}\hat{\beta}^{Ridge} &= \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}. \\ \hat{\beta}^{LASSO} &= \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}.\end{aligned}$$

### Solution.

(a) For ridge regression, we know

$$\begin{aligned}\hat{\beta}^{Ridge} &= (X^T X + \lambda I)^{-1} X^T y \\ &= \frac{1}{1 + \lambda} X^T y = \frac{\hat{\beta}^{OLS}}{(1 + \lambda)}.\end{aligned}$$

(b) For lasso let us write the objective with matrices:

$$\begin{aligned}\hat{\beta}^{LASSO} &= \left\{ \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\} \\ &= \left\{ \frac{1}{2} (y^T y - 2y^T X\beta + \beta^T X^T X\beta) + \lambda \|\beta\|_1 \right\} \equiv \left\{ -y^T X\beta + \frac{1}{2} \beta^T \beta + \lambda \|\beta\|_1 \right\} \\ &= \left\{ \lambda \|\beta\|_1 - \beta^T \hat{\beta}^{OLS} + \frac{1}{2} \beta^T \beta \right\} \\ &= \left\{ \sum_i \lambda |\beta_i| - \beta_i \hat{\beta}_i^{OLS} + \frac{1}{2} \beta_i^2 \right\}.\end{aligned}$$

We can see that the problem is separable, thus it can be solve for each individual  $i$  separately. We have two cases:

- When  $\hat{\beta}_i^{OLS} \geq 0$ , we have that the optimal solution follows  $\beta_i^* \geq 0$ . It can be show that if  $\beta^* < 0$ , the exist a new solution within an  $\varepsilon$ -neighborhood of  $\beta^*$  with better objective, contradicting the optimality of  $\beta^*$ . Thus, the problem to solve is reduced to

$$\min_{\beta \geq 0} \left\{ \beta_i (\lambda - \hat{\beta}_i^{OLS}) + \frac{1}{2} \beta_i^2 \right\}.$$

And it's optimal value is achieved when  $\beta^* = \hat{\beta}_i^{OLS} - \lambda$ . However, as  $\beta \geq 0$  we need to define the solution fo only when  $\beta^*$  is non-negative, i.e.,  $\beta^* = (\hat{\beta}_i^{OLS} - \lambda)_+$ .

- Analogously, when  $\hat{\beta}_i^{OLS} \leq 0$ , we have that the optimal solution follows  $\beta_i^* \leq 0$ . Then, we now solve

$$\min_{\beta \leq 0} \left\{ -\beta_i (\lambda + \hat{\beta}_i^{OLS}) + \frac{1}{2} \beta_i^2 \right\}.$$

which has solution  $\beta^* = (\hat{\beta}_i^{OLS} + \lambda)_-$ .

In both cases the solution can be written as:

$$\hat{\beta}_i^{Lasso} = \text{sign}(\hat{\beta}_i^{OLS})(|\hat{\beta}_i^{OLS}| - \lambda)_+,$$

## PROBLEM 8

Consider the linear regression problem

$$y = X\beta + \varepsilon,$$

with  $y \in \mathbb{R}^n$ ,  $X = (x_{ij})_{i=1,\dots,n;j=1,\dots,p} \in \mathbb{R}^{n \times p}$ ,  $\beta \in \mathbb{R}^p$ , and  $\varepsilon \in \mathbb{R}^n$ .

Suppose we have an orthogonal design matrix, i.e.  $X^T X = I_{p \times p}$ .

- Write down the classical ordinary least squares estimator under the assumption of an orthogonal design matrix. Denote each component of the vector as  $\hat{\beta}_i^{OLS}$ , with  $i = 1, \dots, p$ .
- Then, the Ridge regression problem can be written as

$$\sum_{i=1}^p (\beta_i - \hat{\beta}_i^{OLS})^2 + \lambda \sum_{i=1}^p \beta_i^2.$$

Derive the Ridge regression estimator for the  $i$ -th component in terms of  $\hat{\beta}_i^{OLS}$ .

## Solution.

- Since we know that the OLS estimator is  $\hat{\beta}_{OLS} = (X^T X)^{-1} X^T y$ , and since  $(X^T X)^{-1}$  is the unit matrix, this simplifies to  $\hat{\beta}_{OLS} = X^T y$ . Then,

$$\hat{\beta}_i^{OLS} = \sum_{j=1}^n x_{ji} y_j.$$

- By taking the first derivative with respect to  $\beta_i$ , we get

$$\hat{\beta}_i^{RIDGE} = \frac{1}{\lambda + 1} \hat{\beta}_i^{OLS}.$$

Note why we can write ridge problem in terms of OLS when we have orthogonal design:

$$\begin{aligned} (y - X\beta)'(y - X\beta) &= y'y + \beta'\beta - 2y'X\beta \\ &= y'y + \beta'\beta - 2\hat{\beta}_{OLS}'\beta \\ &= y'y + \beta'\beta - 2\hat{\beta}_{OLS}'\beta + \hat{\beta}_{OLS}'\hat{\beta}_{OLS} - \hat{\beta}_{OLS}'\hat{\beta}_{OLS} \\ &= y'y + \beta'\beta - 2\hat{\beta}_{OLS}'\beta + \hat{\beta}_{OLS}'\hat{\beta}_{OLS} - y'XX'y \\ &= (\hat{\beta}_{OLS} - \beta)'(\hat{\beta}_{OLS} - \beta) + y'(I - XX')y \end{aligned}$$

Since the last term is independent of  $\beta$ , we only need to consider the first term.

### PROBLEM 9

Recall the Ridge regression optimization problem

$$\hat{\beta} = \min_{\beta} \left\{ \sum_{i=1}^n \left( y_i - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}.$$

Considering that  $n = 2$ ,  $p = 2$ ,  $x_{11} = x_{12}$ ,  $x_{21} = x_{22}$ .

- Write down the Ridge regression optimization problem in the described setting, and simplify it as much as possible.
- Show that, in this setting, the Ridge coefficient estimators satisfy  $\hat{\beta}_1 = \hat{\beta}_2$ .
- Write out the LASSO optimization problem in this setting.

### Solution.

a)

$$\begin{aligned} & \sum_{i=1}^n \left( y_i - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \\ &= \sum_{i=1}^n \left( y_i - \beta_1 x_{i1} - \beta_2 x_{i2} \right)^2 + \lambda (\beta_1^2 + \beta_2^2) \\ &= (y_1 - \beta_1 x_{11} - \beta_2 x_{12})^2 + (y_2 - \beta_1 x_{21} - \beta_2 x_{22})^2 + \lambda (\beta_1^2 + \beta_2^2) \\ &= (y_1 - \beta_1 x_{11} - \beta_2 x_{11})^2 + (y_2 - \beta_1 x_{22} - \beta_2 x_{22})^2 + \lambda (\beta_1^2 + \beta_2^2). \\ \Rightarrow \hat{\beta} &= \min_{\beta_1, \beta_2} \left\{ (y_1 - \beta_1 x_{11} - \beta_2 x_{11})^2 + (y_2 - \beta_1 x_{22} - \beta_2 x_{22})^2 + \lambda (\beta_1^2 + \beta_2^2) \right\}. \end{aligned}$$

b)

$$\begin{aligned} & \frac{d}{d\beta_1} [(y_1 - \beta_1 x_{11} - \beta_2 x_{11})^2 + (y_2 - \beta_1 x_{22} - \beta_2 x_{22})^2 + \lambda (\beta_1^2 + \beta_2^2)] \\ &= 2(y_1 - \beta_1 x_{11} - \beta_2 x_{11})(-x_{11}) + 2(y_2 - \beta_1 x_{22} - \beta_2 x_{22})(-x_{22}) + 2\lambda \beta_1 \\ &= 2[-y_1 x_{11} - y_2 x_{22} + \beta_1(x_{11}^2 + x_{22}^2) + \beta_2(x_{11}^2 + x_{22}^2) + \lambda \beta_1]. \end{aligned}$$

$$\begin{aligned} & \frac{d}{d\beta_2} [(y_1 - \beta_1 x_{11} - \beta_2 x_{11})^2 + (y_2 - \beta_1 x_{22} - \beta_2 x_{22})^2 + \lambda (\beta_1^2 + \beta_2^2)] \\ &= 2(y_1 - \beta_1 x_{11} - \beta_2 x_{11})(-x_{11}) + 2(y_2 - \beta_1 x_{22} - \beta_2 x_{22})(-x_{22}) + 2\lambda \beta_2 \\ &= 2[-y_1 x_{11} - y_2 x_{22} + \beta_1(x_{11}^2 + x_{22}^2) + \beta_2(x_{11}^2 + x_{22}^2) + \lambda \beta_2]. \end{aligned}$$

$$\begin{aligned} 2[-y_1 x_{11} - y_2 x_{22} + \beta_1(x_{11}^2 + x_{22}^2) + \beta_2(x_{11}^2 + x_{22}^2) + \lambda \beta_2] &\stackrel{!}{=} 0 \\ 2[-y_1 x_{11} - y_2 x_{22} + \beta_1(x_{11}^2 + x_{22}^2) + \beta_2(x_{11}^2 + x_{22}^2) + \lambda \beta_1] &\stackrel{!}{=} 0 \\ \Rightarrow \beta_1 &= \beta_2. \end{aligned}$$

$$\text{c) } \hat{\beta} = \min_{\beta_1, \beta_2} \left\{ (y_1 - \beta_1 x_{11} - \beta_2 x_{11})^2 + (y_2 - \beta_1 x_{22} - \beta_2 x_{22})^2 + \lambda |\beta_1| + \lambda |\beta_2| \right\}.$$

### PROBLEM 10

The LASSO problem is not always strictly convex, and thus does not necessarily have a unique solution according to standard convexity theory. However, we can define a modified problem, known as the elastic-net optimization problem, that is always strictly convex:

$$\min_{\beta} \left( \|y - X\beta\|_2^2 + \alpha_1 \|\beta\|_2^2 + \alpha_2 \|\beta\|_1 \right) \quad (2)$$

where  $\alpha_1, \alpha_2$  are nonnegative tuning parameters. Besides ensuring uniqueness for all  $X$ , the elastic-net combines some of the desirable predictive characteristics of Ridge regression with the sparsity properties of LASSO.

Show how one can turn this into a LASSO problem, using an augmented version of  $X$  and  $y$ .

### Solution.

Consider  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ , and  $\beta \in \mathbb{R}^p$ .

- First, let  $\bar{X}$  be an augmented version of  $X$  defined as:

$$\bar{X} = \begin{pmatrix} X \\ \gamma I_{p \times p} \end{pmatrix}.$$

Following this, we augment  $y$  by a zero vector of dimension  $p$ , i.e.,

$$\bar{y} = \begin{pmatrix} y \\ 0_{p \times 1} \end{pmatrix}.$$

Then we have

$$\|\bar{y} - \bar{X}\beta\|_2^2 = \left\| \begin{pmatrix} y - X\beta \\ 0_p - \gamma I_{p \times p} \beta \end{pmatrix} \right\|_2^2 = \left\| \begin{pmatrix} y - X\beta \\ -\gamma\beta \end{pmatrix} \right\|_2^2 = \|y - X\beta\|_2^2 + \gamma^2 \|\beta\|_2^2 \quad (3)$$

- Consider now the LASSO problem for  $\bar{y}$  and  $\bar{X}$

$$\min_{\beta} \|\bar{y} - \bar{X}\beta\|_2^2 + \alpha_2 \|\beta\|_1$$

which, by making use of (3), is

$$\min_{\beta} \|y - X\beta\|_2^2 + \gamma^2 \|\beta\|_2^2 + \alpha_2 \|\beta\|_1$$

Therefore, by choosing  $\gamma = \sqrt{\alpha_1}$  we get the original problem.

### PROBLEM 11

Consider a linear regression problem where  $p \gg n$ , and assume that the rank of  $X$  is  $n$ . Let the SVD of  $X = UDV^T = RV^T$ , where  $R$  is  $n \times n$  nonsingular, and  $V$  is  $p \times n$  with orthonormal columns.

- Show that there are infinitely many least-squares solutions all with zero residuals.
- Show that the Ridge-regression estimate for  $\beta$  can be written as

$$\hat{\beta}_{\lambda} = V(R^T R + \lambda I)^{-1} R^T y. \quad (4)$$

- (c) Show that when  $\lambda = 0$ , the solution  $\hat{\beta}_0 = VD^{-1}U^T y$  has zero residuals, and is unique in that it has the smallest Euclidean norm among all zero-residual solutions.

### Solution.

- (a) Since  $X \in \mathbb{R}^{n \times p}$  has rank  $n \leq p$ , then exists  $v \in \mathbb{R}^p \neq 0$  such that  $Xv = 0$ . Let  $\hat{\beta}$  be a zero residual solution for the least-squares problems, i.e.,  $\hat{\beta} = \min_{\beta} \{ \|y - X\beta\|_2^2 \}$ . Then, for every  $k \in \mathbb{R}$ , we have :

$$\|y - X(\hat{\beta} + kv)\|_2^2 = \|y - X\hat{\beta} - Xkv\|_2^2 \equiv \|y - X\hat{\beta}\|_2^2.$$

Therefore, there are infinitely many least-squares solutions all with zero residuals.

- (b) We know that the Ridge-regression estimator  $\beta$  is computed as:

$$\beta = (X^T X + \lambda I)^{-1} X^T y$$

Therefore, it solves the equation  $X^T (y - X\beta) = \lambda\beta$ . Then, by (4),  $X^T (y - X\hat{\beta}_\lambda) = \lambda\hat{\beta}_\lambda$ . Working on the left side of this last equation, we have:

$$\begin{aligned} X^T (y - X\hat{\beta}_\lambda) &= X^T (y - XV(R^T R + \lambda I)^{-1} R^T y) \\ &\stackrel{X=RV^T}{=} VR^T (y - RV^T V(R^T R + \lambda I)^{-1} R^T y) \\ &\stackrel{V^T V=I}{=} VR^T (y - R(R^T R + \lambda I)^{-1} R^T y) \\ &= V (R^T y - R^T R(R^T R + \lambda I)^{-1} R^T y) \\ &= V (I - R^T R(R^T R + \lambda I)^{-1}) R^T y \\ &= V \left( I - (R^T R + \lambda I - \lambda I) (R^T R + \lambda I)^{-1} \right) R^T y \\ &= V \left( I - (R^T R + \lambda I) (R^T R + \lambda I)^{-1} + (\lambda I) (R^T R + \lambda I)^{-1} \right) R^T y \\ &= V \left( I - I + \lambda I (R^T R + \lambda I)^{-1} \right) R^T y \\ &= V \left( \lambda I (R^T R + \lambda I)^{-1} \right) R^T y \\ &= \lambda V \left( (R^T R + \lambda I)^{-1} \right) R^T y \\ &= \lambda \hat{\beta}_\lambda \end{aligned}$$

Therefore, the Ridge-regression estimate for  $\beta$  can be written as

$$\hat{\beta}_\lambda = V(R^T R + \lambda I)^{-1} R^T y.$$

- (c) • Zero residual implies that  $y = X\beta$ .

$$\begin{aligned}
 X\hat{\beta}_0 &= XVD^{-1}U^T y \\
 &= UDV^T VD^{-1}U^T y \\
 &= UDD^{-1}U^T y \\
 &= UU^T y \\
 &= y.
 \end{aligned}$$

Then,  $\hat{\beta}_0$  has zero residual.

- If the solution is not unique, we can construct a zero residual solution as follows:

$$\beta = \hat{\beta}_0 + v,$$

with  $v \in \mathbb{R}^p \neq 0$ . For  $\beta$  to be zero residual it needs to satisfy  $X\beta = y$ , i.e.,  $X(\hat{\beta}_0 + v) = y$ . Given that  $\hat{\beta}_0$  already has zero residual, we have

$$Xv = RV^T v = 0.$$

Taking into consideration that  $R$  is  $n \times n$  nonsingular, we have then that  $V^T v = 0$ . Now, taking the Euclidean norm (squared) of  $\beta$ , we have

$$\|\beta\|_2^2 = (\hat{\beta}_0 + v)^T (\hat{\beta}_0 + v) \tag{5}$$

$$= \hat{\beta}_0^T \hat{\beta}_0 + v^T v + 2\hat{\beta}_0^T v \tag{6}$$

$$= \hat{\beta}_0^T \hat{\beta}_0 + v^T v + 2y^T UD^{-1}V^T v \tag{7}$$

$$= \hat{\beta}_0^T \hat{\beta}_0 + v^T v + 0 \tag{8}$$

$$= \|\hat{\beta}_0\|_2^2 + v^T v. \tag{9}$$

Finally, since  $v^T v > 0$  we have that  $\|\beta\|_2^2 > \|\hat{\beta}_0\|_2^2$ , i.e.,  $\hat{\beta}_0$  has the smallest Euclidean norm.