# Homework 1
# (Linear Regression, Optimization, Regularization)

For questions, please refer to Moodle.
Released on **8 March 2023**

---

GENERAL INSTRUCTIONS
- Submission of solutions is not mandatory but solving the exercises is highly recommended.
- The master solution will be released next week.

---

## Exercise 1: Convex functions

In the following exercise we consider real-valued functions $f : \operatorname{dom}(f) \to \mathbb{R}$, with $\operatorname{dom}(f) \subseteq \mathbb{R}^d$.

A function $f : \operatorname{dom}(f) \to \mathbb{R}$ is *convex* if (i) $\operatorname{dom}(f)$ is a convex set and (ii) for all $x, y \in \operatorname{dom}(f)$ and all $\lambda \in [0,1]$, we have

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \tag{1}$$

Geometrically, the condition means that the line segment connecting the points $(x, f(x)), (y, f(y)) \in \mathbb{R}^{d+1}$ lies point-wise above the graph of $f$; see Figure 1 below.
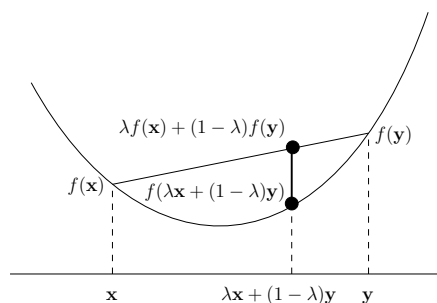


Figure 1: A convex function

### Questions

(a) A norm is a function $f : \mathbb{R}^d \to \mathbb{R}$ which satisfies the following three properties:

- $f(x) > 0$ for all $x \neq 0$
- $f(\theta x) = |\theta| f(x)$ for all $\theta \in \mathbb{R}$ and $x \in \mathbb{R}^d$
- $f(x + y) \leq f(x) + f(y)$ for all $x, y \in \mathbb{R}^d$

Show that any valid norm $f$ is a convex function using the definition of convexity in Equation 1.

---

**Solution:** Let $f$ be a norm satisfying the three properties above. We have

$$f(\lambda x + (1 - \lambda)y) \leq f(\lambda x) + f((1 - \lambda)y) = \lambda f(x) + (1 - \lambda)f(y) \tag{2}$$

where for the first step we used property (iii) and in the second step we used property (ii).

---

Consider now the function $f(x, y) = x^2 + y^2$. The graph of $f$ is the unit paraboloid in $\mathbb{R}^3$ which looks convex. However, verifying the condition in (1) directly is somewhat cumbersome.

To address this problem, we develop better ways to check convexity if the function under consideration is differentiable. In particular, suppose that $\mathrm{dom}(f)$ is open and that $f$ is differentiable, i.e. the gradient

$$\nabla f(x) := \left( \frac{\partial f}{\partial x_1}(x), \ldots, \frac{\partial f}{\partial x_p}(x) \right)$$

exists at every point $x \in \mathrm{dom}(f)$. Then, we will prove that an easier condition to verify convexity is the following:

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) \tag{3}$$

for all $x, y \in \mathrm{dom}(f)$. Geometrically, this means that for all $x \in \mathrm{dom}(f)$, the graph of $f$ lies above its tangent hyperplane at the point $(x, f(x))$ ; see Figure 2 below.
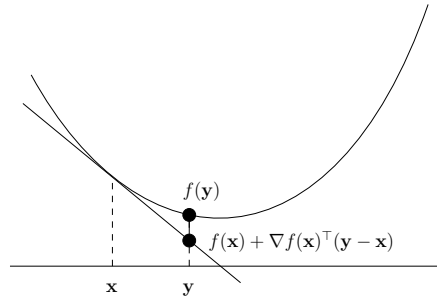


Figure 2: First-order characterization of convexity

## Questions

(b) Assume that $f : \mathbb{R}^d \to \mathbb{R}$ is differentiable. Prove that $f$ is convex if and only if $\mathrm{dom}(f)$ is convex and

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) \tag{4}$$

holds for all $x, y \in \mathrm{dom}(f)$.

---

**Solution:** Suppose that $f$ is convex, meaning that for $t \in [0, 1]$ we have

$$f(x + t(y - x)) = f((1 - t)x + ty) \leq (1 - t)f(x) + tf(y) = f(x) + t(f(y) - f(x)) \tag{5}$$

Dividing by $t$ and using differentiability at $x$, we get

$$
\begin{aligned}
f(y) &\geq f(x) + \frac{f(x + t(y - x)) - f(x)}{t} \\
&= f(x) + \frac{\nabla f(x)^\top t(y - x) + r(t(y - x))}{t} \\
&= f(x) + \nabla f(x)^\top (y - x) + \frac{r(t(y - x))}{t}
\end{aligned}
$$

where the error term $r(t(y - x))/t$ goes to 0 as $t \to 0$. The inequality

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x)$$

follows.

---

Now suppose this inequality holds for all $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$, let $\lambda \in [0,1]$, and define $\mathbf{z} := \lambda \mathbf{x} + (1 - \lambda)\mathbf{y} \in \text{dom}(f)$ (by convexity of $\text{dom}(f)$). Then we have

$$f(\mathbf{x}) \geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{z}),$$
$$f(\mathbf{y}) \geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\mathbf{y} - \mathbf{z}).$$

After multiplying the first inequality by $\lambda$ and the second one by $(1 - \lambda)$, the gradient terms cancel in the sum of the two inequalities, and we get

$$\lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \geq f(\mathbf{z}) = f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y})$$

(c) Suppose that $f$ is differentiable and convex. Prove that $\boldsymbol{x}^\star$ is a global minimum of $f$ if and only if $\nabla f(\boldsymbol{x}^\star) = 0$.

**Solution:** Suppose that $\nabla f(\boldsymbol{x}^\star) = 0$. According to Exercise 1b), we have

$$f(\boldsymbol{y}) \geq f(\boldsymbol{x}^\star) + \nabla f(\boldsymbol{x}^\star)^\top (\boldsymbol{y} - \boldsymbol{x}^\star) = f(\boldsymbol{x}^\star) \tag{6}$$

for all $\boldsymbol{y} \in \text{dom}(f)$, so $\boldsymbol{x}^\star$ is a global minimum.

The converse is also true and does not even require convexity. Suppose that $\nabla f(\boldsymbol{x}^\star)_i \neq 0$ for some $i$. For $t \in \mathbb{R}$, we define $\boldsymbol{x}(t) = \boldsymbol{x}^\star + t\boldsymbol{e}_i$, where $\boldsymbol{e}_i$ is the $i$-th unit vector. For $|t|$ sufficiently small, we have $\boldsymbol{x}(t) \in \text{dom}(f)$ since $\text{dom}(f)$ is open. Let $z(t) = f(\boldsymbol{x}(t))$. By the chain rule, $z'(0) = \nabla f(\boldsymbol{x}^\star)^\top \boldsymbol{e}_i = \nabla f(\boldsymbol{x}^\star)_i \neq 0$. Hence, $z$ decreases in one direction as we move away from $t = 0$, and this yields $f(\boldsymbol{x}(t)) < f(\boldsymbol{x}^\star)$ for some $t$, so $\boldsymbol{x}^\star$ is not a global minimum.

(d) Show that $f(x,y) = x^2 + y^2$ is a convex function and that the point $(0,0)$ is a global minimum.

**Solution:** For $f(x_1, x_2) = x_1^2 + x_2^2$ we have that the gradient is equal to

$$\nabla f(x_1, x_2) = (2x_1, 2x_2) \tag{7}$$

hence substituting into Equation 4 we have

$$y_1^2 + y_2^2 \geq x_1^2 + x_2^2 + 2x_1(y_1 - x_1) + 2x_2(y_2 - x_2) \tag{8}$$

which after some rearranging of terms is equivalent to

$$(y_1 - x_1)^2 + (y_2 - x_2)^2 \geq 0 \tag{9}$$

hence according to Exercise 1b), $f$ is convex.

Moreover, the gradient of $f$ at the point $(0,0)$ is $\nabla f(0,0) = \mathbf{0}$ and $f$ is differentiable on $\text{dom}(f) = \mathbb{R}$, hence according to Exercise 1c), the point is a global minimum.

A third way to verify convexity is through:

$$\nabla^2 f(x) \succeq 0 \tag{10}$$

Geometrically, this means that the graph of $f$ has non-negative curvature everywhere and hence looks like a bowl. We will come back to this when discussing ridge regression in Exercise 3.

**Questions**

(e) A continuously differentiable function $f : \mathbb{R}^d \to \mathbb{R}$ is called $\alpha$-strongly convex for some $\alpha > 0$, if for any points $x, y \in \mathbb{R}^d$ one has

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\alpha}{2} \|y - x\|^2.$$

If $f$ is twice differentiable, an equivalent condition is that for any point $x \in \mathbb{R}^d$, one has

$$\nabla^2 f(x) \succeq \alpha I$$

which means $\nabla^2 f(x) - \alpha I$ is positive semi-definite for all $x \in \mathbb{R}^d$. Prove that a strongly convex function admits a unique minimizer in $\mathbb{R}^d$.
*Hint: This is not an easy exercise. First prove that $f(x) \to \infty$ as $\|x\| \to \infty$ to show that there is some minimizer.*

---

**Solution:** First, let us prove that the function $f$ is coercive, i.e., $\lim_{\|x\| \to \infty} f(x) = \infty$. In the definition of strong convexity, by setting $x = 0$ we get

$$f(y) \geq f(0) + \nabla f(0)^\top y + \frac{\alpha}{2} \|y\|^2 \geq f(0) - \|\nabla f(0)\| \cdot \|y\| + \frac{\alpha}{2} \|y\|^2,$$

where we used the Cauchy-Schwartz inequality: $\nabla f(0)^\top y \geq -\|\nabla f(0)\| \cdot \|y\|$. The right-hand side of the equation above is a quadratic function of $\|y\|$ with a positive coefficient for the second degree term. Thus, it goes to infinity as $\|y\| \to \infty$. Hence, $f$ also goes to infinity.

Next, we prove that $f$ has a global minimum. Denote by $s = \inf_{x \in \mathbb{R}^d} f(x) < \infty$. Then, by the definition of the infimum, there exists a sequence $x_1, x_2, \ldots$ such that $f(x_n) \to s$. We claim that this sequence is bounded: otherwise, there would be a subsequence such that $\|x_{n_i}\| \to \infty$. But as $f$ is coercive, $f(x_{n_i}) \to \infty$, contradicting $f(x_{n_i}) \to s < \infty$. Hence, the sequence $x_1, x_2, \ldots$ is inside some bounded set. By compactness, we obtain that there exists a convergent subsequence. As $f$ is continuous, the $f$ value of this subsequence converges as well, meaning that the infimum is attained. That is, $\exists x_\infty : f(x_\infty) = s = \inf f(x)$.

Finally, we prove uniqueness. If $x$ and $y$ were two distinct global minima for $f$, then, by strong convexity, we have

$$f\left(\frac{x + y}{2}\right) < \frac{1}{2}(f(x) + f(y)) = \min f,$$

which results in a contradiction.

---

## Exercise 2: Linear regression

In the lecture we have learned how to fit an affine function to data by performing linear regression. In the tutorial on Friday we will discuss how to fit more general nonlinear functions to data. The goal of this exercise is to solidify our understanding of some of the concepts that have been touched upon.

Consider a dataset $\{(x_1, y_1), \ldots, (x_n, y_n)\} \subset \mathbb{R} \times \mathbb{R}$ and the hypothesis space of affine functions $H = \{w_0 + w_1 x : w_0, w_1 \in \mathbb{R}\}$. The error of a solution $f \in H$, i.e., $f(x) = w_0 + w_1 x$ for some $w_0, w_1 \in \mathbb{R}$ is given by

$$L(w_0, w_1) = \frac{1}{2n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i)^2. \tag{11}$$

## Questions

(a) Let us for a moment consider the simpler case where we fix $w_0 = 0$. Compute the optimal linear fit to the data by computing $w_1^* = \arg\min_{w_1 \in \mathbb{R}} L(0, w_1)$.

**Solution:** We start by computing the first derivative

$$\frac{\partial L(0, w_1)}{\partial w_1} = \frac{1}{n} \sum_{i=1}^{n} (w_1 x_i - y_i) x_i$$

and the second derivative

$$\frac{\partial^2 L(0, w_1)}{\partial w_1^2} = \frac{1}{n} \sum_{i=1}^{n} x_i^2.$$

We observe that the second derivative is greater or equal to zero for all $w_1 \in \mathbb{R}$. Thus, $L(0, w_1)$ is convex and we can find its global minimum (if it exists) by setting its first derivative to zero

$$\frac{1}{n} \sum_{i=1}^{n} (w_1 x_i - y_i) x_i = 0$$

$$\Leftrightarrow w_1 = \frac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2}.$$

If not all $x_i$ are equal to zero, $w_1^*$ is given by the equation above.

(b) Prove that, for $n \geq 2$ and $x_i \neq x_j$ for $i \neq j$, (11) is a strictly convex function with respect to $w = (w_0, w_1)$.

**Solution:** We are going to show that (11) is strictly convex by showing that its Hessian $\nabla^2 L(w_0, w_1)$ is positive definite. First we compute the gradient $\nabla L(w_0, w_1)$

$$\nabla L(w_0, w_1) = \frac{1}{n} \begin{pmatrix} \sum_{i=1}^{n} w_1 x_i + w_0 - y_i \\ \sum_{i=1}^{1} (w_1 x_i + w_0 - y_i) x_i \end{pmatrix}$$

Then, we can derive the gradient to obtain the Hessian

$$\nabla^2 L(w_0, w_1) = \frac{1}{n} \begin{pmatrix} n & \sum_{i=1}^{n} x_i \\ \sum_{i=1}^{n} x_i & \sum_{i=1}^{n} x_i^2 \end{pmatrix} = \frac{1}{n} \begin{pmatrix} n & \sum_{i=1}^{n} x_i \\ \sum_{i=1}^{n} x_i & \sum_{i=1}^{n} x_i^2 \end{pmatrix}$$

Now, a $2 \times 2$ symmetric matrix is positive definite if and only if both its trace and determinant are strictly greater than zero. First, we show that the trace is strictly positive:

$$\mathrm{Tr}(\nabla^2 L(w_0, w_1)) = \frac{1}{n}(n + \sum_{i=1}^{n} x_i^2) > 0 \quad \forall x_i \in \mathbb{R} \tag{12}$$

Next, we show that the determinant is strictly positive:

$$n^2 \det(\nabla^2 L(w_0, w_1)) = n \sum_{i=1}^{n} x_i^2 - \left( \sum_{i=1}^{n} x_i \right)^2 = \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} (x_i - x_j)^2 > 0 \tag{13}$$

since $x_i - x_j \neq 0$ for $i \neq j$.

(c) The unique global minimum of a strictly convex function can be computed by setting its gradient to zero. Compute the gradient

$$\nabla L(w_0, w_1) = \begin{pmatrix} \frac{\partial L(w_0, w_1)}{\partial w_0} \\ \frac{\partial L(w_0, w_1)}{\partial w_1} \end{pmatrix}. \tag{14}$$

**Solution:** We get

$$\nabla L(w_0, w_1) = \begin{pmatrix} \frac{\partial L(w_0, w_1)}{\partial w_0} \\ \frac{\partial L(w_0, w_1)}{\partial w_1} \end{pmatrix} = \begin{pmatrix} \frac{1}{n} \sum_{i=1}^{n} w_0 + w_1 x_i - y_i \\ \frac{1}{n} \sum_{i=1}^{n} (w_0 + w_1 x_i - y_i) x_i \end{pmatrix}.$$

(d) Compute the optimal parameters $(w_0^*, w_1^*) = \arg\min_{w_0, w_1 \in \mathbb{R}} L(w_0, w_1)$ by solving the linear system of equations obtained by setting (14) to zero, i.e., $\nabla L(w_0, w_1) = 0$.

**Solution:** From setting $\frac{\partial L(w_0, w_1)}{\partial w_0}$ to zero we get

$$w_0^* = \frac{\sum_{i=1}^{n} y_i - w_1 \sum_{i=1}^{n} x_i}{n} = \bar{y} - w_1 \bar{x}$$

where $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$ and $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$. Now, plugging $w_0^*$ into $\frac{\partial L(w_0, w_1)}{\partial w_1}$ and setting the resulting expression to zero we get

$$w_1^* = \frac{\sum_{i=1}^{n} x_i y_i - x_i \bar{y}}{\sum_{i=1}^{n} x_i^2 - x_i \bar{x}}$$

which can be rewritten (by adding zeros to both enumerator and denominator) to

$$w_1^* = \frac{\sum_{i=1}^{n} (x_i y_i - x_i \bar{y}) + \overbrace{\sum_{i=1}^{n} (\overline{xy} - \bar{x} y_i)}^{=0}}{\sum_{i=1}^{n} (x_i^2 - x_i \bar{x}) + \underbrace{\sum_{i=1}^{n} (\bar{x}^2 - \bar{x} x_i)}_{=0}}$$

$$= \frac{\sum_{i=1}^{n} (\bar{x} - x_i)(\bar{y} - y_i)}{\sum_{i=1}^{n} (\bar{x} - x_i)^2}$$

$$= \frac{\text{Cov}(x, y)}{\text{Var}(x)}.$$

Another way of writing (11) is in matrix notation

$$L(\boldsymbol{w}) = \frac{1}{2n} \|\boldsymbol{y} - \Phi \boldsymbol{w}\|^2,$$ (15)

where $\boldsymbol{y} = (y_1, \ldots, y_n)^T$ is the vector of target values, $\boldsymbol{w} = (w_0, w_1)^T$ is our weight vector and

$$\Phi = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$ (16)

is the data matrix.

For $n \geq 2$ different observations (15) is a strictly convex function and can be minimized by setting its gradient

$$\nabla L(\boldsymbol{w}) = \frac{1}{n} (\Phi^T \Phi \boldsymbol{w} - \Phi^T \boldsymbol{y})$$ (17)

to zero.

The benefit of (15) is that it straightforwardly generalizes to multiple inputs $x_i \in \mathbb{R}^d$ using

$$\Phi = \begin{pmatrix} 1 & x_{11} & \ldots & x_{1d} \\ 1 & x_{21} & \ldots & x_{2d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \ldots & x_{nd} \end{pmatrix},$$ (18)

where we have one row per data point and one column per input.

## Questions

(e) Provide necessary conditions for $\Phi$ such that $\Phi^T \Phi$ is invertible.

> **Solution:** A necessary and sufficient condition for invertibility of $\Phi^\top \Phi$ is:
>
> $$\Phi^\top \Phi \text{ is invertible} \iff \text{rank}(\Phi) = d + 1$$ (19)
>
> We first prove that $\text{rank}(\Phi) = d + 1 \implies \Phi^\top \Phi$ is invertible by defining a function $f$ as
>
> $$f(x) := x^\top \Phi^\top \Phi x = \|\Phi x\|^2$$
>
> which is positive semidefinite. Moreover, $f$ vanishes when $\Phi x = 0$. If $\Phi$ has full column rank, then $\Phi x = 0$ implies that $x = 0$, i.e., $f$ is positive definite. Hence, $\Phi^T \Phi$ is positive definite and thus invertible.
>
> We next show that $\Phi^\top \Phi$ is invertible $\implies \text{rank}(\Phi) = d + 1$ which concludes the proof. If $\Phi$ has not full column rank then there exists $v \in \mathbb{R}^d \backslash \{0\}$ with $\Phi v = 0$. It follows that $\Phi^\top \Phi v = 0$ and hence $\Phi^\top \Phi$ is not invertible.

(f) Show that if $\Phi^T \Phi$ is invertible, then $\boldsymbol{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \boldsymbol{y}$ is the unique minimum of $L(\boldsymbol{w})$ in (15).

> **Solution:** First, we compute the gradient:
>
> $$\nabla L(\boldsymbol{w}) = \frac{1}{n} (\Phi^T \Phi \boldsymbol{w} - \Phi^T \boldsymbol{y})$$ (20)

and by setting it to zero we get the following solution:

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T \boldsymbol{y} \tag{21}$$

Now, all is left to show is that the Hessian of $L(w)$ is positive definite, this implies strict convexity which in turn implies unique minimum.

First, we derive the Hessian:

$$\nabla_w^2 L(w) = \frac{1}{n} \Phi^\top \Phi \tag{22}$$

Now if $\Phi$ is full rank then $\Phi^\top \Phi$ is positive definite (see solution of exercise 2e) and hence the minimum is unique.

(g) Show that for $n < d + 1$ the regression problem

$$\min_{\boldsymbol{w} \in \boldsymbol{R}^{d+1}} \| \boldsymbol{y} - \Phi \boldsymbol{w} \|^2 \tag{23}$$

does not admit a unique solution.

**Solution:** If $n < d + 1$, we have rank$(\Phi) = \min(n, d + 1) = n$. Further, we know rank$\left(\Phi^T \Phi\right) = $ rank$(\Phi)$. By the rank-nullity theorem we have

$$\dim \left( \ker \left( \Phi^T \Phi \right) \right) = d + 1 - \text{rank} \left( \Phi^T \Phi \right) = d + 1 - n > 0$$

We have seen that our objective $\| \boldsymbol{y} - \Phi \boldsymbol{w} \|^2$ is convex. Thus, each minimal weight vector $w^* \in$ arg $\min_w \| \boldsymbol{y} - \Phi \boldsymbol{w} \|^2$ satisfies $\nabla \| \boldsymbol{y} - \Phi \boldsymbol{w}^* \|^2 = 0 \iff \Phi^T \Phi \boldsymbol{w}^* = \Phi^T \boldsymbol{y}$. Now, given a minimal weight vector $w^*$, we have $\Phi^T \Phi \left( \boldsymbol{w}^* + u \right) = \Phi^T \boldsymbol{y}$ for all $u \in \ker \left( \Phi^T \Phi \right)$. As $\dim \left( \ker \left( \Phi^T \Phi \right) \right) \geq 1$, there are infinitely many such $u \neq 0$.

Next, recall the gradient descent update as discussed in lecture: $\boldsymbol{w}^{t+1} = \boldsymbol{w}^t - \eta \nabla_w L(\boldsymbol{w}^t)$. We would like to compare the computational complexity of the closed-form solution $\boldsymbol{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \boldsymbol{y}$ against the one of the gradient descent algorithm. For the next exercises you may use the contraction inequality as discussed during lecture,

$$\| \boldsymbol{w}^{t+1} - \boldsymbol{w}^* \|_2 \leq \overbrace{\| I - \eta \Phi^\top \Phi \|_{op}}^{:=\rho} \| \boldsymbol{w}^t - \boldsymbol{w}^* \|_2$$

## Questions

(h) Prove that if $\Phi^\top \Phi$ is full rank and the learning rate satisfies $\eta < \frac{2}{\lambda_{max}(\Phi^\top \Phi)}$, where $\lambda_{max}$ is the largest eigenvalue, then $\rho < 1$. Conclude that gradient descent converges to the optimal solution $w^*$ as $t \to \infty$.

Hint: remember that an equivalent characterization of the operator norm in terms of the eigenvalues of $A$ is given by $\| A \|_{op} = \max \{ |\lambda| : \lambda \text{ eigenvalue of } A \}$

**Solution:**

First, we note that if $\lambda_1, \ldots, \lambda_d$ are the eigenvalues of the matrix $\Phi^\top \Phi$, then the eigenvalues of $I - \eta \Phi^\top \Phi$ are exactly $1 - \eta \lambda_1, \ldots, 1 - \eta \lambda_d$. Now, using the characterisation of operator norm given in the hint, we have

$$\rho = \max \left\{ \left| \lambda_{\max} \left( I - \eta \Phi^\top \Phi \right) \right|, \left| \lambda_{\min} \left( I - \eta \Phi^\top \Phi \right) \right| \right\}$$
$$= \max \left\{ \left| 1 - \eta \lambda_{\min} \left( \Phi^\top \Phi \right) \right|, \left| 1 - \eta \lambda_{\max} \left( \Phi^\top \Phi \right) \right| \right\}.$$

Under a smaller step size $\eta \leq \frac{1}{\lambda_{\max}(\Phi^\top \Phi)}$, both of the terms in the expression are non-negative and hence we can drop the absolute value function

$$\rho = \max \left\{ 1 - \eta \lambda_{\min} \left( \Phi^\top \Phi \right), 1 - \eta \lambda_{\max} \left( \Phi^\top \Phi \right) \right\}$$
$$= 1 - \eta \lambda_{\min} \left( \Phi^\top \Phi \right)$$

Therefore, $\rho < 1$, as we have assumed that $\Phi^\top \Phi$ is full rank and hence $\lambda_{\min} \left( \Phi^\top \Phi \right) > 0$. On the other hand, if $\eta > \frac{1}{\lambda_{\max}(\Phi^\top \Phi)}$, we have

$$\rho = \max \left\{ \left| 1 - \eta \lambda_{\min} \left( \Phi^\top \Phi \right) \right|, \left| 1 - \eta \lambda_{\max} \left( \Phi^\top \Phi \right) \right| \right\}$$
$$= \max \left\{ 1 - \eta \lambda_{\min} \left( \Phi^\top \Phi \right), \eta \lambda_{\min} \left( \Phi^\top \Phi \right) - 1, \eta \lambda_{\max} \left( \Phi^\top \Phi \right) - 1 \right\}$$
$$= \max \left\{ 1 - \eta \lambda_{\min} \left( \Phi^\top \Phi \right), \eta \lambda_{\max} \left( \Phi^\top \Phi \right) - 1 \right\}$$

If we want $\rho < 1$, two conditions shall be satisfied: (i) $1 - \eta \lambda_{\min} \left( \Phi^\top \Phi \right) < 1$, which is already fulfilled, as we assumed $\lambda_{\min} \left( \Phi^\top \Phi \right) > 0$, and (ii) $\eta \lambda_{\max} \left( \Phi^\top \Phi \right) - 1 < 1$, which holds if $\eta < \frac{2}{\lambda_{\max}(\Phi^\top \Phi)}$. Hence, it is clear that both conditions are satisfied under our assumptions.

Finally, observe that if $\rho < 1$, the upper bound on $\left\| w^{t+1} - w^* \right\|_2$ converges linearly to zero as $t \to \infty$.

(i) Assume that the stepsize $\eta$ is such that $\| I - \eta \Phi^T \Phi \|_{op} < 1$. Compute the number of gradient steps $\tau$ and the overall complexity required to obtain a solution $w^\tau$ that satisfies $\| w^\tau - w^* \| < \varepsilon$, where $w^\tau$ is the parameter vector computed by gradient descent after $\tau$ steps.

**Solution:** Let $\rho = \| I - \eta \Phi^T \Phi \|_{op}$. We want to find $\tau$ such that $\| w^\tau - w^* \| < \epsilon$. Applying the contraction inequality to $\| w^\tau - w^* \|$ $\tau$ times gives $\| w^\tau - w^* \| \leq \rho^\tau \| w^0 - w^* \|$. Because $\| w^\tau - w^* \| \leq \rho^\tau \| w^0 - w^* \|$ we have that $\| w^\tau - w^* \| < \epsilon$ when $\rho^\tau \| w^0 - w^* \| < \epsilon$. Bringing $\| w^0 - w^* \|$ to the other side and taking the logarithm on both sides yields

$$\rho^\tau \left\| w^0 - w^* \right\| < \epsilon$$
$$\Leftrightarrow \rho^\tau < \frac{\epsilon}{\| w^0 - w^* \|}$$
$$\Leftrightarrow \tau \log \rho < \log \epsilon - \log \left\| w^0 - w^* \right\|$$
$$\Leftrightarrow \tau > \frac{\log \epsilon - \log \left\| w^0 - w^* \right\|}{\log \rho},$$

where in the last line the inequality is flipped because $\log \rho$ is smaller than zero for $0 < \rho < 1$. Let's denote $p = d + 1$. Computing the gradient at a given step $t + 1$, i.e., computing $\nabla L \left( w^t \right) = \frac{1}{n} \left( \Phi^T \Phi w^t - \Phi^T y \right)$, $O \left( p^2 \right)$ operations to compute $\Phi^T \Phi w^t$ and once $O \left( np^2 + np \right) = O \left( np^2 \right)$ operations to compute $\Phi^T \Phi$ and $\Phi^T y$ at the beginning of the algorithm. Subtracting

$\eta \nabla L\left(\boldsymbol{w}^t\right)$ to $\boldsymbol{w}^t$ has negligible cost of $O(p)$. Therefore, in order to compute a solution $\boldsymbol{w}^\tau$ that satisfies $\left\|\boldsymbol{w}^\tau - \boldsymbol{w}^*\right\| < \epsilon$ gradient descent requires $\tau = \left\lceil \frac{\log \epsilon - \log\left\|\boldsymbol{w}^0 - \boldsymbol{w}^*\right\|}{\log \rho} \right\rceil$ steps, resulting in a computational complexity of $O\left(np^2 + \frac{\log \epsilon - \log\left\|\boldsymbol{w}^0 - \boldsymbol{w}^*\right\|}{\log \rho} p^2\right)$.

(j) For the linear regression loss function defined in Equation 15, prove that $L(\boldsymbol{w}^{t+1}) < L(\boldsymbol{w}^t)$ for small enough stepsize $\eta$.

**Solution:** Without loss of generality, let us assume that $\nabla L\left(\boldsymbol{w}^t\right)$ is nonzero (otherwise the gradient descent makes no move, and we are already at a stationary point). Using the multivariate Taylor expansion, we have

$$
\begin{aligned}
L\left(\boldsymbol{w}^{t+1}\right) &= L\left(\boldsymbol{w}^t - \eta \nabla L\left(\boldsymbol{w}^t\right)\right) \\
&= L\left(\boldsymbol{w}^t\right) - \eta \left\langle \nabla L\left(\boldsymbol{w}^t\right), \nabla L\left(\boldsymbol{w}^t\right)\right\rangle + R\left(-\eta \nabla L\left(\boldsymbol{w}^t\right)\right) \\
&= L\left(\boldsymbol{w}^t\right) - \eta \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\|^2 + R\left(-\eta \nabla L\left(\boldsymbol{w}^t\right)\right)
\end{aligned}
\tag{24}
$$

where $R : \mathbb{R}^d \to \mathbb{R}$ is the remainder term. This term is of $o\left(\eta \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\|\right)$, i.e.,

$$
\lim_{\eta \to 0} \frac{R\left(-\eta \nabla L\left(\boldsymbol{w}^t\right)\right)}{\eta \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\|} = 0.
$$

Then, by definition of the limit, for all $\epsilon > 0$ there exists $\eta > 0$ small enough such that

$$
\left| \frac{R\left(-\eta \nabla L\left(\boldsymbol{w}^t\right)\right)}{\eta \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\|} \right| \le \epsilon
$$

or

$$
\left|R\left(-\eta \nabla L\left(\boldsymbol{w}^t\right)\right)\right| \le \epsilon \eta \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\|.
$$

Applying this to (24) and picking $\epsilon < \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\|$ we get

$$
\begin{aligned}
L\left(\boldsymbol{w}^{t+1}\right) &= L\left(\boldsymbol{w}^t\right) - \eta \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\|^2 + R\left(-\eta \nabla L\left(\boldsymbol{w}^t\right)\right) \\
&\le L\left(\boldsymbol{w}^t\right) - \eta \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\|^2 + \epsilon \eta \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\| \\
&= L\left(\boldsymbol{w}^t\right) + \eta \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\| \left(\epsilon - \left\|\nabla L\left(\boldsymbol{w}^t\right)\right\|\right) \\
&< L\left(\boldsymbol{w}^t\right),
\end{aligned}
$$

(k) Now, say you are free to choose a constant stepsize. What is the minimum number of iterations $\tau$ required to obtain a solution $\boldsymbol{w}^\tau$ that satisfies $\left\|\boldsymbol{w}^\tau - \boldsymbol{w}^*\right\| < \varepsilon$? How does it depend on the maximum and minimum eigenvalues $\lambda_{max}$, $\lambda_{min}$ of the matrix $\Phi^T \Phi$?

**Solution:** In order to achieve the minimum number of required iterations $\tau$, we first need to determine the optimal learning rate $\eta$. The operator norm associated with the L2 norm is the spectral norm. Thus, we have

$$
\rho = \left\|I - \eta \Phi^T \Phi\right\| = \max\left(\left|1 - \eta \lambda_{\max}\right|, \left|1 - \eta \lambda_{\min}\right|\right)
$$

where $\lambda_{\min}$ and $\lambda_{\max}$ denote the smallest and the largest eigenvalue of $\Phi^T \Phi$. Now, in order to

determine the optimal choice of $\eta$ we first derive an expression for $\rho$ as a function of $\eta$. We have

$$
\begin{aligned}
\rho(\eta) &= \max\left(|1 - \eta\lambda_{\max}|, |1 - \eta\lambda_{\min}|\right) \\
&= \max\left(1 - \eta\lambda_{\max}, \eta\lambda_{\max} - 1, 1 - \eta\lambda_{\min}, \eta\lambda_{\min} - 1\right) \\
&= \max\left(\eta\lambda_{\max} - 1, 1 - \eta\lambda_{\min}\right)
\end{aligned}
$$

where the last equality holds because for $0 \leq \lambda_{\min} \leq \lambda_{\max}$ we have $\eta\lambda_{\max} - 1 > \eta\lambda_{\min} - 1$ and $1 - \eta\lambda_{\min} > 1 - \eta\lambda_{\max}$. Therefore, we have either $\rho = \eta\lambda_{\max} - 1$, i.e., $\eta\lambda_{\max} - 1 \geq 1 - \eta\lambda_{\min}$ which is equivalent to $\eta \geq \frac{2}{\lambda_{\min}+\lambda_{\max}}$, or $\rho = 1 - \eta\lambda_{\min}$, i.e., $1 - \eta\lambda_{\min} > \eta\lambda_{\max} - 1$ which is equivalent to $\eta < \frac{2}{\lambda_{\min}+\lambda_{\max}}$. As a result we get

$$
\rho(\eta) = \begin{cases}
1 - \eta\lambda_{\min} & \text{for } 0 < \eta < \frac{2}{\lambda_{\min}+\lambda_{\max}}, \\
\eta\lambda_{\max} - 1 & \text{for } \frac{2}{\lambda_{\min}+\lambda_{\max}} \leq \eta < \frac{1}{\lambda_{\max}}
\end{cases},
$$

which is a piecewise affine function in $\eta$ that attains its minimum at $\eta^* = \frac{2}{\lambda_{\min}+\lambda_{\max}}$. The corresponding value of $\rho$ is

$$
\begin{aligned}
\rho(\eta^*) &= 1 - \frac{2\lambda_{\min}}{\lambda_{\min} + \lambda_{\max}} \\
&= \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\min} + \lambda_{\max}}.
\end{aligned}
$$

Plugging the optimal $\rho(\eta^*)$ into the expression for $\tau$ from exercise (h) yields

$$
\tau = \left\lceil \frac{\log \epsilon - \log\|w^0 - w^*\|}{\log(\lambda_{\max} - \lambda_{\min}) - \log(\lambda_{\min} + \lambda_{\max})} \right\rceil
$$

(l) Compare the computational complexity of gradient descent to the one required to solve the linear system of equations $\Phi^T\Phi w^* = \Phi^T y$ in closed form.

**Solution:** The complexity of gradient descent is $O\left(np^2 + \tau p^2\right)$. In comparison, the complexity of computing $w^*$ in closed form, which is dominated by the cost of the matrix inversion $\left(\Phi^T\Phi\right)^{-1}$ is in $O\left(p^3\right)$. Furthermore, $w^*$ can be also computed by solving the linear system of equations $\Phi^T\Phi w = \Phi^T y$, which is dominated by the $O\left(np^2\right)$ cost of computing $\Phi^T\Phi$ as solving the resulting $p \times p$ linear system of equations is in $O\left(p^2\right)$ and the computation of $\Phi^T y$ in $O(np)$. As a consequence, in terms of computational complexity it is best to compute $w^*$ by solving the linear system of equations $\Phi^T\Phi w = \Phi^T y$ that is obtained by setting the gradient to zero. Comparing the cost of gradient descent to the cost of computing the closed form by multiplying with the inverse $\left(\Phi^T\Phi\right)^{-1}$ shows that gradient descent can be preferable if $p$ is significantly larger than $n + \tau$.

# Exercise 3: Bias variance tradeoff

In the lecture we have seen the bias variance decomposition of the expected squared error. In this exercise, we will first derive this decomposition in the case of linear regression and we will then use it to study ridge regression. Consider a dataset with fixed data matrix $\Phi \in \mathbb{R}^{n \times (d+1)}$ as shown in Equation 18 and noisy labels $y \in \mathbb{R}^n$ with $y_i = w^\top\phi(x_i) + \epsilon_i$. We define the noise $\epsilon_i$ such that $\mathbb{E}[\epsilon_i] = 0$ and $\text{Var}(\epsilon_i) = \sigma^2$. We assume that we have used this training data to train a linear regressor with estimator $\hat{w}$. We now consider the expected squared error of this estimate defined as $\mathbb{E}_\epsilon\left[\|w - \hat{w}\|^2\right]$. In this exercise, $\|\cdot\|$ always refers to the Euclidean norm. Remember the definitions of bias and variance as given in the

lecture:

$$\text{Bias}(\hat{w})^2 = \|w - \mathbb{E}_\epsilon[\hat{w}]\|^2 \tag{25}$$

$$\text{Var}(\hat{w}) = \mathbb{E}_\epsilon\left[\|\hat{w} - \mathbb{E}_\epsilon[\hat{w}]\|^2\right] \tag{26}$$

## Questions

(a) Prove that $\mathbb{E}_\epsilon\left[\|w - \hat{w}\|^2\right] = \text{Bias}(\hat{w})^2 + \text{Var}(\hat{w})$.

**Solution:**

$$\begin{aligned}
\mathbb{E}_\epsilon\left[\|w - \hat{w}\|^2\right] &= \mathbb{E}_\epsilon\left[\|w - \mathbb{E}_\epsilon[\hat{w}] + \mathbb{E}_\epsilon[\hat{w}] - \hat{w}\|^2\right] \\
&= \mathbb{E}_\epsilon\left[\|w - \mathbb{E}_\epsilon[\hat{w}]\|^2\right] + 2(w - \mathbb{E}_\epsilon[\hat{w}])^\top \underbrace{\mathbb{E}_\epsilon\left[(\mathbb{E}_\epsilon[\hat{w}] - \hat{w})\right]}_{=0} + \mathbb{E}_\epsilon\left[\|\mathbb{E}_\epsilon[\hat{w}] - \hat{w}\|^2\right] \\
&= \|w - \mathbb{E}_\epsilon[\hat{w}]\|^2 + \mathbb{E}_\epsilon\left[\|\hat{w} - \mathbb{E}_\epsilon[\hat{w}]\|^2\right] = \text{Bias}(\hat{w})^2 + \text{Var}(\hat{w})
\end{aligned}$$

We can see that even if our estimator $\hat{w}$ is unbiased, we can still incur a large error if it has high variance. Similarly, even if the estimator has small variance, we can incur a large error if it is strongly biased. This is commonly referred to as the bias variance tradeoff.

Remember that while bias drives underfitting, variance is related to overfitting. During training, we aim to find the model with the lowest expected generalization error which requires joint optimization over these two components. As seen in the lecture, regularization is frequently used to control model complexity which reduces variance at the expense of increased bias. Ridge regression is a regularized version of linear regression. The optimization problem with parameter $\lambda > 0$ is given by Equation 27 which is optimizing over the loss we have seen in Exercise 2 but with an added regularization term $\lambda\|w\|^2$.

$$w^*_{\text{ridge}} = \arg\min_w L_{\text{ridge}}(w) = \arg\min_w \left[\sum_{i=1}^n (y_i - w^\top \phi(x_i))^2 + \lambda\|w\|^2\right] \tag{27}$$

## Questions

(b) Prove that the solution to Equation 27 is unique for any matrix $\Phi$ by showing that $L_{\text{ridge}}$ is strictly convex.

**Solution:** In matrix notation, we can write $L_{\text{ridge}}(w) = \|\Phi w - y\|^2 + \lambda\|w\|^2$ (analogously to what we did for unregularized linear regression). We first compute the gradient $\nabla_w L_{\text{ridge}}(w)$ by using the chain rule as:

$$\nabla_w L_{\text{ridge}}(w) = 2(\Phi w - y)^\top \Phi + 2\lambda w^\top = 2w^\top \left(\Phi^\top \Phi + \lambda I_d\right) - 2y^\top \Phi$$

Moreover, we can obtain the Hessian $\nabla_w^2 L_{\text{ridge}}(w)$ by differentiating $\nabla_w L_{\text{ridge}}(w)$ again:

$$\nabla_w^2 L_{\text{ridge}}(w) = 2\left(\Phi^\top \Phi + \lambda I_d\right).$$

Now, to prove that $L_{\text{ridge}}$ is strictly convex we need that for some $\alpha > 0$:

$$\nabla_w^2 L_{\text{ridge}}(w) - \alpha I_d \succeq 0 \tag{28}$$

and for $\alpha = 2\lambda$ the above condition is equivalent to:

$$2\Phi^\top \Phi \succeq 0 \tag{29}$$

which holds for any $\Phi$ since for any non-zero vector $v \in \mathbb{R}^d$, we have $v^\top \Phi^\top \Phi v = \|\Phi v\|^2 \geq 0$.

Hence, $L_{\text{ridge}}$ is $2\lambda$-strongly convex, and from Exercise 1e) we can conclude that $L_{\text{ridge}}$ has a unique minimizer.

In the following, we want to develop a better understanding of the bias variance tradeoff by comparing ridge regression to unregularized linear regression.

## Questions

(c) Derive the closed-form solution $w^*_{\text{ridge}}$ by computing the unique minimizer defined in Equation 27.

> **Solution:** We know from task (b) that $\nabla_w L_{\text{ridge}}(w) = 2w^\top \left(\Phi^\top \Phi + \lambda I_d\right) - 2y^\top \Phi$ and by setting this to 0, we find:
>
> $$w^*_{\text{ridge}} = \left(\Phi^\top \Phi + \lambda I_d\right)^{-1} \Phi^\top y$$
>
> Note that the regularization term in the optimization problem admits a solution for any matrix $\Phi$ even in the case where $\Phi^\top \Phi$ is not be invertible.

(d) Show that its bias is given by $\text{Bias}(w^*_{\text{ridge}}) = \|\lambda \left(\Phi^\top \Phi + \lambda I_d\right)^{-1} w\|$. How does it compare to the bias of unregularized linear regression?

> **Solution:** Remember that the bias depends on the noisy labels $y$ with noise $\epsilon$ as defined in the beginning of Exercise 3. We find
>
> $$\begin{aligned}
> \text{Bias}(w^*_{\text{ridge}}) &= \|w - \mathbb{E}_\epsilon[w^*_{\text{ridge}}]\| \\
> &= \|w - \mathbb{E}_\epsilon\left[\left(\Phi^\top \Phi + \lambda I_d\right)^{-1} \Phi^\top y\right]\| \\
> &= \|w - \left(\Phi^\top \Phi + \lambda I_d\right)^{-1} \Phi^\top \mathbb{E}_\epsilon[y]\| \\
> &= \|w - \left(\Phi^\top \Phi + \lambda I_d\right)^{-1} \Phi^\top \Phi w\| \\
> &= \|\left(\Phi^\top \Phi + \lambda I_d\right)^{-1} \left(\Phi^\top \Phi + \lambda I_d\right) w - \left(\Phi^\top \Phi + \lambda I_d\right)^{-1} \Phi^\top \Phi w\| \\
> &= \|\lambda \left(\Phi^\top \Phi + \lambda I_d\right)^{-1} w\|
> \end{aligned}$$

(e) We next focus on the variance of the unregularized estimator as derived in Exercise 2. Remember that this estimator is defined as $w^* = (\Phi^\top \Phi)^{-1} \Phi^\top y$ assuming that $\Phi$ is full rank. Show that its variance is given by

$$\text{Var}(w^*) = \sum_{i=1}^{m} \frac{\sigma^2}{\sigma_i^2}$$

where the $\sigma_i$ are the singular values of $\Phi$ and $m = \min(n, d+1)$.

*Hint: Remember that $\text{Var}(\hat{w}) = \mathbb{E}_\epsilon\left[\|\hat{w} - \mathbb{E}_\epsilon[\hat{w}]\|^2\right]$ and use the fact that $A = Tr(A)$ if $A$ is a real number and that the trace is invariant under cyclic permutations $Tr(ABC) = Tr(BCA) = Tr(CAB)$.*

**Solution:** We begin by showing that $\mathrm{Var}(w^*) = \mathbb{E}_\epsilon\left[\|(\Phi^\top\Phi)^{-1}\Phi^\top\epsilon\|^2\right]$.

$$\mathrm{Var}(w^*) = \mathbb{E}_\epsilon\left[\|w^* - \mathbb{E}_\epsilon[w^*]\|^2\right]$$
$$= \mathbb{E}_\epsilon\left[\|(\Phi^\top\Phi)^{-1}\Phi^\top y - \mathbb{E}_\epsilon[(\Phi^\top\Phi)^{-1}\Phi^\top y]\|^2\right]$$
$$= \mathbb{E}_\epsilon\left[\|(\Phi^\top\Phi)^{-1}\Phi^\top y - (\Phi^\top\Phi)^{-1}\Phi^\top\mathbb{E}_\epsilon[y]\|^2\right]$$
$$= \mathbb{E}_\epsilon\left[\|(\Phi^\top\Phi)^{-1}\Phi^\top(\Phi w + \epsilon) - (\Phi^\top\Phi)^{-1}\Phi^\top\mathbb{E}_\epsilon[(\Phi w + \epsilon)]\|^2\right]$$
$$= \mathbb{E}_\epsilon\left[\|(\Phi^\top\Phi)^{-1}\Phi^\top(\Phi w + \epsilon) - (\Phi^\top\Phi)^{-1}\Phi^\top\Phi w\|^2\right]$$
$$= \mathbb{E}_\epsilon\left[\|(\Phi^\top\Phi)^{-1}\Phi^\top\epsilon\|^2\right]$$

If we expand this term and define the empirical covariance matrix $\hat{\Sigma} = \Phi^T\Phi$, we get

$$\mathbb{E}_\epsilon\left[\|(\Phi^\top\Phi)^{-1}\Phi^\top\epsilon\|^2\right] = \mathbb{E}_\epsilon\left[\epsilon^\top\Phi\left((\Phi^\top\Phi)^{-1}\right)^\top(\Phi^\top\Phi)^{-1}\Phi^\top\epsilon\right]$$
$$= \mathbb{E}_\epsilon\left[\epsilon^\top\Phi(\Phi^\top\Phi)^{-1}(\Phi^\top\Phi)^{-1}\Phi^\top\epsilon\right]$$
$$= \mathbb{E}_\epsilon\left[\epsilon^\top\Phi\hat{\Sigma}^{-1}\hat{\Sigma}^{-1}\Phi^\top\epsilon\right]$$

As this is a real number, we can replace it with its trace (see hint). Using the invariance of the trace under cyclic permutations, we find

$$\epsilon^\top\Phi\hat{\Sigma}^{-1}\hat{\Sigma}^{-1}\Phi^\top\epsilon = \mathrm{Tr}(\epsilon^\top\Phi\hat{\Sigma}^{-1}\hat{\Sigma}^{-1}\Phi^\top\epsilon) = \mathrm{Tr}(\epsilon\epsilon^\top\Phi\hat{\Sigma}^{-1}\hat{\Sigma}^{-1}\Phi^\top)$$

Taking the expectation and exploiting the invariance of the trace under cyclic permutations, we finally arrive at

$$\mathrm{Var}(w^*) = \mathbb{E}_\epsilon\left[\mathrm{Tr}(\epsilon\epsilon^\top\Phi\hat{\Sigma}^{-1}\hat{\Sigma}^{-1}\Phi^\top)\right]$$
$$= \mathrm{Tr}\left(\mathbb{E}_\epsilon\left[\epsilon\epsilon^\top\right]\Phi\hat{\Sigma}^{-1}\hat{\Sigma}^{-1}\Phi^\top\right)$$
$$= \sigma^2\mathrm{Tr}\left(\Phi\hat{\Sigma}^{-1}\hat{\Sigma}^{-1}\Phi^\top\right)$$
$$= \sigma^2\mathrm{Tr}\left(\Phi^\top\Phi\hat{\Sigma}^{-1}\hat{\Sigma}^{-1}\right)$$
$$= \sigma^2\mathrm{Tr}\left(\hat{\Sigma}\hat{\Sigma}^{-1}\hat{\Sigma}^{-1}\right)$$
$$= \sigma^2\mathrm{Tr}\left(\hat{\Sigma}^{-1}\right)$$
$$= \sum_{i=1}^{d+1}\frac{\sigma^2}{\sigma_i^2}$$

where we have used that $\mathrm{Var}(\epsilon_i) = \sigma^2$ and $\sigma_i$ are the singular values of $\Phi$. The last equality follows from the fact that $\mathrm{Tr}(\hat{\Sigma}^{-1})$ is equal to the sum of its eigenvalues and the eigenvalues of $\hat{\Sigma}^{-1}$ are equal to the inverse of the eigenvalues of $\hat{\Sigma}$, which in turn are equal to the singular values of $\Phi$ squared.

(f) Similarly, it can be shown that the ridge estimator has variance

$$\mathrm{Var}(w^*_{\mathrm{ridge}}) = \sigma^2\sum_{i=1}^{m}\frac{\sigma_i^2}{(\sigma_i^2 + \lambda)^2}$$

where the $\sigma_i$ are the singular values of $\Phi$ and $m = \min(n, d+1)$. Compare the variance of the unregularized estimator to that of the ridge estimator.

**Solution:** We can see that for $\lambda = 0$ we have $\text{Var}(w^*) = \text{Var}(w^*_{\text{ridge}})$. For $\lambda > 0$, we find that $\text{Var}(w^*_{\text{ridge}}) < \text{Var}(w^*)$ and in the extreme case where $\lambda \to \infty$, $\text{Var}(w^*_{\text{ridge}}) \to 0$.

(g) How does the choice of $\lambda$ affect bias and variance? How do bias and variance behave as $\lambda \to 0$ and $\lambda \to \infty$?

**Solution:** To summarize, we have the following:

$$\text{Bias}(w^*) = 0$$

$$\text{Var}(w^*) = \sum_{i=1}^{m} \frac{\sigma^2}{\sigma_i^2}$$

$$\text{Bias}(w^*_{\text{ridge}}) = \|\lambda \left(\Phi^\top \Phi + \lambda I_d\right)^{-1} w\|$$

$$\text{Var}(w^*_{\text{ridge}}) = \sigma^2 \sum_{i=1}^{m} \frac{\sigma_i^2}{(\sigma_i^2 + \lambda)^2}$$

As discussed in the lecture, regularization lowers the variance at the expense of introducing a bias. To see this, we will examing the limits $\lambda \to 0$ and $\lambda \to \infty$.

As $\lambda \to 0$, we find that $w^*_{\text{ridge}} \to w^*$ and hence $\text{Bias}(w^*_{\text{ridge}}) \to 0$ and $\text{Var}(w^*_{\text{ridge}}) \to \text{Var}(w^*)$.

For $\lambda \to \infty$, we can see by computing the limits that $\text{Bias}(w^*_{\text{ridge}}) \to \|w\|$ and $\text{Var}(w^*_{\text{ridge}}) \to 0$.

Hence, increasing $\lambda$ corresponds to reducing the variance and increasing the bias simultaneously while decreasing $\lambda$ reduces the bias but increases the variance. This tradeoff corresponds to the tradeoff between underfitting (low variance, high bias) and overfitting (low bias, high variance) and hence the choice of $\lambda$ has to be considered carefully when fine-tuning a model. We have to jointly minimize over the bias squared and the variance of our estimator to find the the solution with the lowest squared error.

## Exercise 4: Gradient Descent for Linear Regression (Bonus)

In this exercise, we are going to prove that under mild conditions the gradient descent algorithm for ordinary linear regression problems converges to the solution with minimum norm. This is a very good exercise to practice your linear algebra skills. To help you prove this argument, we have divided the complete proof into smaller chunks. As in the lecture, suppose $X \in \mathbb{R}^{n \times d}$ is the data matrix and $y \in \mathbb{R}^n$ is the response vector. The goal is to find a vector $w \in \mathbb{R}^d$ such that $L(w) := \frac{1}{2}\|Xw - y\|^2$ is minimized. For this, we use the gradient descent algorithm: starting from an initial vector $w^0$, the iterates of the gradient descent algorithm for a step size $\eta$ are

$$w^{k+1} = w^k - \eta \nabla L(w^k), \quad k = 0, 1, \dots$$

(a) By computing the gradient of $L$, confirm that

$$w^{k+1} = (I - \eta X^\top X)w^k + \eta X^\top y.$$

**Solution:** We have seen in the lecture that

$$\nabla L(w^k) = X^\top X w^k - X^\top y$$

and if we insert this into the equation for the gradient descent algorithm, we find that:

$$w^{k+1} = w^k - \eta(X^\top X w^k - X^\top y) = (I - \eta X^\top X)w^k + \eta X^\top y$$

(b) By using induction on $k$, prove that

$$w^k = \underbrace{(I - \eta X^\top X)^k w^0}_{(A)} + \underbrace{\eta \left( \sum_{j=0}^{k-1} (I - \eta X^\top X)^j \right) X^\top y}_{(B)}$$

(30)

**Solution:** Remember that an induction proof consists of two steps: we first show that a given equation holds for the base case (here $k = 0$) and we then show that if it holds for any given $k$, it must also hold for the next step at $k + 1$.

For $k = 0$, we find that part (B) in Equation (30) is 0 due to the empty set of summation indices $j$ and we hence find that:
$$w^0 = (I - \eta X^\top X)^0 w^0 = w^0$$

Assuming that Equation (30) holds for $k$, we can prove it for $k + 1$ by plugging it into the gradient descent equation and verifying that Equation (30) holds for $k + 1$ as well:

$$
\begin{aligned}
w^{k+1} &= (I - \eta X^\top X)w^k + \eta X^\top y \\
&= (I - \eta X^\top X)\left((I - \eta X^\top X)^k w^0 + \eta \left(\sum_{j=0}^{k-1}(I - \eta X^\top X)^j\right)X^\top y\right) + \eta X^\top y \\
&= (I - \eta X^\top X)^{k+1} w^0 + \eta \left(\sum_{j=0}^{k-1}(I - \eta X^\top X)^{j+1}\right)X^\top y + \eta X^\top y \\
&= (I - \eta X^\top X)^{k+1} w^0 + \eta \left(\sum_{j=1}^{k}(I - \eta X^\top X)^{j}\right)X^\top y + \eta X^\top y \\
&= (I - \eta X^\top X)^{k+1} w^0 + \eta \left(\sum_{j=0}^{k}(I - \eta X^\top X)^{j}\right)X^\top y
\end{aligned}
$$

From (b) it is clear that powers of the matrix $I - \eta X^\top X$ play an important role in understanding what happens to $w^k$ when $k$ is large. It is usual to look at the eigenvalues of a matrix when studying its powers. Hence, we start by the SVD of $X = U\Sigma V^\top$, where $U \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{d \times d}$ are orthogonal matrices and $\Sigma \in \mathbb{R}^{n \times d}$ is a rectangular diagonal matrix with non-negative real numbers $\sigma_1, \ldots, \sigma_n$ on its diagonal. *From this part onwards, we focus on the over-parameterized case where $n < d$.*

(c) Verify that the eigenvalue decomposition of $I - \eta X^\top X$ is $V(I - \eta\Lambda)V^\top$, where $\Lambda \in \mathbb{R}^{d \times d}$ is a diagonal matrix whose first $n$ diagonal entries are $\sigma_i^2$ and the rest are zero.

**Solution:** As $X^\top X = V\Sigma^\top \Sigma V^\top$ and $I = VV^\top$ (since $V$ is orthogonal), we can write

$$I - \eta X^\top X = V(I - \eta\Sigma^\top\Sigma)V^\top = V(I - \eta\Lambda)V^\top,$$

where $\Lambda$ is the matrix as described in the problem. Note that since $V$ is orthogonal and $I - \eta\Lambda$ is diagonal, it follows that we have just derived the eigendecomposition of $I - \eta X^\top X$.

(d) Denote by $\sigma_{\max} := \max \sigma_i$. Observe that if $\eta \leq 1/\sigma_{\max}^2$, all eigenvalues of $I - \eta X^\top X$ will be non-negative.

> **Solution:** The entries on the diagonal (and hence, the eigenvalues) of $I - \eta \Lambda$ are $1 - \eta \sigma_i^2$. Hence, if $\eta \leq 1/\sigma_{\max}^2$, they are all non-negative.

(e) Compute $(I - \eta X^\top X)^k$ in closed form for any $k \geq 1$ based on $V$, $\eta$ and $\Lambda$.

> **Solution:** The blessing of having the eigendecomposition shows up here: we have
> $$(I - \eta X^\top X)^k = V(I - \eta \Lambda)^k V^\top.$$

We now compute parts (A) and (B) in Equation (30) separately. *From now on, we assume that $\eta \leq 1/\sigma_{\max}^2$.*

(f) If $v^i$ is an eigenvector of $X^\top X$ corresponding to the eigenvalue $\sigma_i^2$, compute $(I - \eta X^\top X)^k v^i$. Describe what happens when $k \to \infty$. (*Hint: you have to consider two cases: when $\sigma_i = 0$ and $\sigma_i > 0$*)

> **Solution:** Note that the eigenvectors of $X^\top X$ and $I - \eta X^\top X$ are the same. Hence, if $v^i$ is an eigenvector for $X^\top X$ corresponding to eigenvalue $\sigma_i^2$, it is also an eigenvector for $I - \eta X^\top X$ with eigenvalue $(1 - \eta \sigma_i^2)$. Hence,
> $$(I - \eta X^\top X)^k v^i = (1 - \eta \sigma_i^2)^k v^i.$$
> If $\sigma_i = 0$, then $(I - \eta X^\top X)^k v^i = v^i$ for any $k$.
> If $\sigma_i > 0$, then $1 - \eta \sigma_i^2 < 1$ and as $k \to \infty$, $(1 - \eta \sigma_i^2)^k \to 0$ and hence $(I - \eta X^\top X)^k v^i \to \mathbf{0}$.

(g) Based on the last step, compute part (A) when $k \to \infty$. (*Hint: decompose $\boldsymbol{w}^0 = \boldsymbol{v} + \boldsymbol{u}$, where $\boldsymbol{u} \in \ker(X)$ and $\boldsymbol{v} \in \ker(X)^\perp$*)

> **Solution:** It is a fact from linear algebra that the right singular vectors (columns of $V$) corresponding to zero singular values form a basis of the kernel of the matrix.
>
> In our case, this means that $\sigma_i = 0$ only if $v^i \in \ker(X)$. Thus, if we decompose $\boldsymbol{w}^0 = \boldsymbol{v} + \boldsymbol{u}$ as in the hint, we have that $\boldsymbol{u} = \sum_{i:\sigma_i=0} \alpha_i v^i$ (since $\boldsymbol{u} \in \ker(X)$ and $v^i$ form a basis for the kernel), and $\boldsymbol{v} = \sum_{j:\sigma_j>0} \beta_j v^j$. We can thus write
> $$(I - \eta X^\top X)^k \boldsymbol{w}^0 = \sum_{i:\sigma_i=0} \alpha_i (I - \eta X^\top X)^k v^i + \underbrace{\sum_{j:\sigma_j>0} \beta_j (I - \eta X^\top X)^k v^j}_{\to 0 \quad \text{as} \quad k\to\infty}$$

(h) Show that part (B) is equal to
$$V\left(\eta \sum_{j=0}^{k-1} (I - \eta \Lambda)^j\right) \Sigma^\top U^\top \boldsymbol{y}$$

> **Solution:**
> $$\eta \left(\sum_{j=0}^{k-1} (I - \eta X^\top X)^j X^\top \boldsymbol{y}\right) = \eta \left(\sum_{j=0}^{k-1} V(I - \eta \Lambda)^j V^\top\right) V \Sigma^\top U^\top \boldsymbol{y} = V\left(\eta \sum_{j=0}^{k-1} (I - \eta \Lambda)^j\right) \Sigma^\top U^\top \boldsymbol{y}$$

(i) Compute (B) when $k \to \infty$. (*Hint: treat zero and positive singular values separately.*)

**Solution:** For simplicity, define
$$M := \sum_{j=0}^{k-1} (I - \eta \Lambda)^j.$$

Note that $M$ is a diagonal $d \times d$ matrix. Its $i$th entry on the diagonal is

$$1 + (1 - \eta \sigma_i^2) + \cdots + (1 - \eta \sigma_i^2)^{k-1} = \begin{cases} k & \text{if} \quad \sigma_i = 0 \\ \frac{1-(1-\eta\sigma_i^2)^k}{\eta\sigma_i^2} & \text{if} \quad 0 < \sigma_i < 1/\sqrt{\eta} \end{cases}$$

Now note that $M\Sigma^\top$ is rectangular diagonal, and its $i$th element on the diagonal will be

$$M_{ii}\Sigma_{ii} = \begin{cases} 0 & \text{if} \quad \sigma_i = 0 \\ \frac{1-(1-\eta\sigma_i^2)^k}{\eta\sigma_i} & \text{if} \quad 0 < \sigma_i < 1/\sqrt{\eta} \end{cases}$$

Hence, when $k \to \infty$, the matrix $\eta M\Sigma^\top$ converges to the rectangular diagonal matrix $N$ with diagonal entries
$$N_{ii} = \begin{cases} 0 & \text{if} \quad \sigma_i = 0 \\ \frac{1}{\sigma_i} & \text{if} \quad \sigma_i > 0 \end{cases}.$$

(j) Prove that the limit computed above equals $X^+ y$, where $X^+$ is the Moore-Penrose pseudo-inverse.

**Solution:** It is known that $VNU^\top$ is indeed the Moore-Penrose pseudo-inverse. Hence, part (B) equals $X^+ y$ when $k \to \infty$. Further, it can be proved that the Moore-Penrose pseudo-inverse provides the solution $w = X^+ y$ with minimal Euclidean norm $\|w\|$ among all possible solutions.

(k) Notice that the above argument also works for the case where $n \geq d$. Make the necessary adjustments and prove that gradient descent initialized at **0** and with a small enough step size converges to the correct solution in under-parameterized setting.

**Solution:** The only thing that changes is the dimensions of $\Lambda$. The rest of the proof follows.