

EE 248 HW#1 - Convexity and Gradient Methods

Due date: Sunday October 17, 2021 midnight over iLearn.

- You only need to submit solutions for 2 out of the 3 problems.
- If you submit all 3 solutions you get 10% bonus. If your effort is not sufficient, you may not get it (at TA's discretion).
- Only 1 problem of TA's choice will be graded (out of the ones you submitted).
- You are free to discuss with others but you are not allowed to look at other's HW solutions. It has to be 100% your writing, your work.
- Please L^AT_EX the solutions.

Problem 1

Part a. Let $(y_i, \mathbf{x}_i)_{i=1}^n$ be a dataset with $y \in \mathbb{R}, \mathbf{x} \in \mathbb{R}^d$. Fix $\lambda \geq 0$ and consider the ridge regression problem

$$\mathcal{L}(\boldsymbol{\beta}) = \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + \frac{\lambda}{2} \|\boldsymbol{\beta}\|_2^2.$$

- Derive $\nabla \mathcal{L}(\boldsymbol{\beta})$ and $\nabla^2 \mathcal{L}(\boldsymbol{\beta})$. Make sure to use matrix/vector notation where $\mathbf{X} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ and $\mathbf{y} = [y_1 \ \dots \ y_n] \in \mathbb{R}^n$.
- Determine the precise conditions under which $\mathcal{L}(\boldsymbol{\beta})$ is a strongly-convex function.

Part b. Now consider the loss function

$$\bar{\mathcal{L}}(\boldsymbol{\beta}) = \frac{1}{2} \sum_{i=1}^n (y_i - |\mathbf{x}_i^\top \boldsymbol{\beta}|)^2.$$

- Given a choice of $(\mathbf{x}_i)_{i=1}^n$, construct $(y_i)_{i=1}^n$ such that the loss function $\bar{\mathcal{L}}(\boldsymbol{\beta})$ is convex. What is the most general set (to which $(y_i)_{i=1}^n$ belongs) that guarantees convexity?
- Given a choice of $(\mathbf{x}_i)_{i=1}^n$, construct $(y_i)_{i=1}^n$ such that $\bar{\mathcal{L}}(\boldsymbol{\beta})$ is nonconvex.

Solution a. $\nabla \mathcal{L}(\boldsymbol{\beta}) = \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) + \lambda \boldsymbol{\beta}$. $\nabla^2 \mathcal{L}(\boldsymbol{\beta}) = \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}$.

The Problem is strongly convex iff $\nabla^2 \mathcal{L}(\boldsymbol{\beta}) > 0$ i.e. Hessian is positive definite. For this to happen, we either need

- $\lambda > 0$: positive ridge regression or
- $\mathbf{X}^\top \mathbf{X} > 0$: \mathbf{X} is rank d . This is same as $n \geq d$ and \mathbf{X} is full rank.

Solution b: Q1. Expanding the function we find $(y_i - |\mathbf{x}_i^\top \boldsymbol{\beta}|)^2 = y_i^2 - 2y_i |\mathbf{x}_i^\top \boldsymbol{\beta}| + (\mathbf{x}_i^\top \boldsymbol{\beta})^2$. Observe that y_i^2 is a constant, $(\mathbf{x}_i^\top \boldsymbol{\beta})^2$ is a convex function and $|\mathbf{x}_i^\top \boldsymbol{\beta}|$ is a convex function as well. However, $-2y_i |\mathbf{x}_i^\top \boldsymbol{\beta}|$ can be convex or concave depending on y_i .

Main observation is that if we have $y_i \leq 0$ for all $1 \leq i \leq n$, then the problem is convex because each $-2y_i |\mathbf{x}_i^\top \boldsymbol{\beta}|$ term becomes convex.

Solution b: Q2. If all of $(\mathbf{x}_i)_{i=1}^n$ are zero problem is convex. Suppose one of the examples $(\mathbf{x}_i)_{i=1}^n$ are nonzero. Then, we can choose some $\boldsymbol{\beta}$ so that $\mathbf{X}\boldsymbol{\beta} \neq 0$ and set $\mathbf{y} = |\mathbf{X}\boldsymbol{\beta}|$. This \mathbf{y} choice leads to nonconvexity. The reason is that $\bar{\mathcal{L}}(\boldsymbol{\beta}) = \bar{\mathcal{L}}(-\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{y} - |\mathbf{X}\boldsymbol{\beta}|\|_{\ell_2}^2 = 0$ so $\pm\boldsymbol{\beta}$ are global minima. However $\bar{\mathcal{L}}(0) = \frac{1}{2} \|\mathbf{y}\|_{\ell_2}^2 > 0$ which contradicts with the convexity since 0 is linear combination of $\pm\boldsymbol{\beta}$.

Solution b: Q1: Most general set. Here, the cleanest answer is all choices of y_i that makes $f(\boldsymbol{\beta}) = -\sum_{i=1}^n y_i |\mathbf{x}_i^\top \boldsymbol{\beta}|$ convex.

- This automatically covers the case y_i are nonpositive.
- Conversely, if y_i 's are all non-negative and there is at least one pair (y, \mathbf{x}) obeying $y > 0$ and $\mathbf{x} \neq 0$, then $f(\boldsymbol{\beta})$ is nonconvex because

$$f(\mathbf{x}) = f(-\mathbf{x}) \leq y|\mathbf{x}^\top \mathbf{x}| = y\|\mathbf{x}\|_{\ell_2}^2 < f(0) = 0.$$

I say the convexity of $f(\boldsymbol{\beta})$ is the cleanest answer because $\bar{\mathcal{L}}(\boldsymbol{\beta}) = 0.5\|\mathbf{y}\|_{\ell_2}^2 + f(\boldsymbol{\beta}) + f'(\boldsymbol{\beta})$ where $f'(\boldsymbol{\beta}) = 0.5\sum_{i=1}^n (\mathbf{x}_i^\top \boldsymbol{\beta})^2$. $0.5\|\mathbf{y}\|_{\ell_2}^2$ term is constant and $f'(\boldsymbol{\beta})$ has no importance. To see the latter observe that $\bar{\mathcal{L}}$ is convex if and only if f is convex.

- If $f(\boldsymbol{\beta})$ is convex, $f'(\boldsymbol{\beta})$ is also convex so $\bar{\mathcal{L}}(\boldsymbol{\beta})$ is convex.
- If $f(\boldsymbol{\beta})$ is not convex, there is $\boldsymbol{\beta}, \boldsymbol{\beta}'$ such that $\alpha f(\boldsymbol{\beta}) + (1-\alpha)f(\boldsymbol{\beta}') < f(\boldsymbol{\beta}_\alpha)$ where $\boldsymbol{\beta}_\alpha = \alpha\boldsymbol{\beta} + (1-\alpha)\boldsymbol{\beta}'$. Let ε be a small scalar of our choice. Observe that

$$\bar{\mathcal{L}}(\varepsilon\boldsymbol{\beta}) = 0.5\|\mathbf{y}\|_{\ell_2}^2 + \varepsilon f(\boldsymbol{\beta}) + \varepsilon^2 f'(\boldsymbol{\beta}).$$

Thus, we find

$$(1/\varepsilon)[\alpha\bar{\mathcal{L}}(\varepsilon\boldsymbol{\beta}) + (1-\alpha)\bar{\mathcal{L}}(\varepsilon\boldsymbol{\beta}') - \bar{\mathcal{L}}(\varepsilon\boldsymbol{\beta}_\alpha)] = \underbrace{[\alpha f(\boldsymbol{\beta}) + (1-\alpha)f(\boldsymbol{\beta}') - f(\boldsymbol{\beta}_\alpha)]}_{<0} + \underbrace{\varepsilon[\alpha f'(\boldsymbol{\beta}) + (1-\alpha)f'(\boldsymbol{\beta}') - f'(\boldsymbol{\beta}_\alpha)]}_{\rightarrow 0 \text{ as } \varepsilon \rightarrow 0}$$

Thus, letting $\varepsilon \rightarrow 0$, we find that $\bar{\mathcal{L}}(\boldsymbol{\beta})$ is not convex.

Problem 2

Part a. Recall that a convex function f may not be differentiable but it admits subgradients. The set of all subgradients of a function f is called the subdifferential and is denoted by $\partial f(\mathbf{a})$. Consider the ℓ_1 norm of a vector $\mathbf{a} \in \mathbb{R}^d$ which is defined as $\|\mathbf{a}\|_1 = \sum_{i=1}^d |a_i|$. Determine the set $\partial \|\mathbf{a}\|_1$.

Part b. Suppose you wish to pick a subgradient $\mathbf{g} \in \partial \|\mathbf{a}\|_1$ such that ℓ_1 norm reduces as much as possible when you take an ε small step in the $-\mathbf{g}$ direction. Determine this optimal \mathbf{g} at a given \mathbf{a} .

Part c. Consider the lasso problem which is ℓ_1 -penalized least-squares regression

$$\mathcal{L}(\boldsymbol{\beta}) = \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + \lambda \|\boldsymbol{\beta}\|_1.$$

Suppose you wish to solve $\arg \min_{\boldsymbol{\beta}} \mathcal{L}(\boldsymbol{\beta})$ using an iterative algorithm with constant learning rate η . Based on Parts a. & b., derive how to take a step from the current iterate $\boldsymbol{\beta}$ to obtain the next iterate $\hat{\boldsymbol{\beta}}$.

Hint: Part of the loss is not differentiable so you need to choose the best sub-gradient.

Part d. People use ℓ_1 norm for a reason right? It turns out, it is because ℓ_1 norm is the smallest p for which ℓ_p quasi-norm is convex. Let us define the ℓ_p quasi-norm of a vector $\|\mathbf{a}\|_p = (\sum_{i=1}^d |a_i|^p)^{1/p}$ for all $p > 0$. Prove that ℓ_p quasi-norm is a nonconvex function for all $0 < p < 1$.

Solution a. Here a key observation is that ℓ_1 norm is decomposable into individual coordinates since $\|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i|$. Thus subgradient is the concatenation of the subgradients of the individual absolute value functions because $\partial\|\mathbf{x}\|_1/\partial x_i = \partial|x_i|/\partial x_i$. Let $\text{sgn}(\mathbf{a})$ be the sign vector which is 1 if $a_i > 0$, -1 if $a_i < 0$ and 0 else. Subdifferential is given by

$$\partial\|\mathbf{a}\|_1 = \{\mathbf{x} \mid \|\mathbf{x}\|_\infty \leq 1, x_i = \text{sgn}(a_i), \text{ whenever } a_i \neq 0\}.$$

In words, subgradient is equal to $\text{sgn}(\mathbf{a})$ over the nonzero entries and the remaining entries are arbitrary with absolute value bounded by 1.

Solution b. $\mathbf{g} = \text{sgn}(\mathbf{a})$. If \mathbf{g} had a nonzero entry in a coordinate $a_i = 0$, then we would have $\|\mathbf{a} - \varepsilon \mathbf{g}\|_1 > \|\mathbf{a} - \varepsilon \text{sgn}(\mathbf{a})\|_1$ because $|\mathbf{g}_i + a_i| = |\mathbf{g}_i| > 0$ would increase the ℓ_1 norm. Thus, optimal \mathbf{g} has to be $\text{sgn}(\mathbf{a})$.

Solution c. Any subgradient can be written as $g_{\mathcal{L}}(\beta) = 0.5 \nabla \|\mathbf{y} - \mathbf{X}\beta\|_{\ell_2}^2 + \lambda \mathbf{g}_{\|\cdot\|_1}(\beta)$. Following part (b), we choose $g_{\mathcal{L}}(\beta) = \mathbf{X}^\top (\mathbf{X}\beta - \mathbf{y}) + \lambda \text{sgn}(\beta)$.

Solution d. If $d = 1$, all norms are equal to absolute value and convex. So assume $d \geq 2$. Choose $\mathbf{x} = [1 \ 0]^\top$ and $\mathbf{y} = [0 \ 1]^\top$ and $\mathbf{z} = [0.5 \ 0.5]^\top$. For $p < 1$, convexity is violated via

$$\frac{\|\mathbf{x}\|_p + \|\mathbf{y}\|_p}{2} = 1 < 2^{1/p-1} = 2^{1/p} 0.5 = \|\mathbf{z}\|_p.$$

Problem 3

Part a. Let $\mathcal{C} \in \mathbb{R}^d$ be a closed convex set. Let $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ two points. Define the projection operator to be the unique point $\Pi_{\mathcal{C}}(\mathbf{a}) = \arg \min_{\mathbf{v} \in \mathcal{C}} \|\mathbf{a} - \mathbf{v}\|_2$. Prove that projection is contractive, that is,

$$\|\Pi_{\mathcal{C}}(\mathbf{a}) - \Pi_{\mathcal{C}}(\mathbf{b})\|_2 \leq \|\mathbf{a} - \mathbf{b}\|_2. \quad (\text{CONTRACTION})$$

Hint: Use the fact that projection satisfies $\langle \mathbf{v} - \Pi_{\mathcal{C}}(\mathbf{a}), \mathbf{a} - \Pi_{\mathcal{C}}(\mathbf{a}) \rangle \leq 0$ for any $\mathbf{v} \in \mathcal{C}$.

Significance: In the special case of $\mathbf{b} \in \mathcal{C}$, this result implies: “Projecting \mathbf{a} onto \mathcal{C} is guaranteed to return a closer point to \mathbf{b} ” and forms the basis of iterative algos that use gradient-descent followed by projection.

Part b. Set \mathcal{C} to be the set of s -sparse signals, that is, $\mathbf{v} \in \mathcal{C}$ iff \mathbf{v} has at most s -nonzero entries for some $s \leq d$. Construct an example where (CONTRACTION) is violated for \mathcal{C} .

Part c. Consider the ℓ_∞ norm-ball defined as $\mathcal{C} = \{\mathbf{v} \in \mathbb{R}^d \mid |v_i| \leq 1 \text{ for all } 1 \leq i \leq d\}$. Derive (the analytical expression for) the projection onto this ball i.e. $\Pi_{\mathcal{C}}(\mathbf{x})$.

Part d. In Problem 2a. you calculated the subdifferential of the ℓ_1 norm. Now, let us set \mathcal{C} to be this subdifferential i.e. $\mathcal{C} = \partial\|\mathbf{a}\|_1$. Derive the projection onto this subdifferential i.e. $\Pi_{\mathcal{C}}(\mathbf{x})$. Comment on any similarities to the norm-ball projection in Problem 3c. (**Hint:** What happens when $\mathbf{a} = 0$?)

Solution a (easiest to understand is by drawing on a paper). Let $\mathbf{a}' = \Pi_{\mathcal{C}}(\mathbf{a})$, $\mathbf{b}' = \Pi_{\mathcal{C}}(\mathbf{b})$. Let P denote all points between \mathbf{a}', \mathbf{b}' . Note that P is in the convex set \mathcal{C} and $\mathbf{b}' = \Pi_P(\mathbf{b})$ is the projection of \mathbf{b} onto $P \subset \mathcal{C}$. Consider the two dimensional plane \mathcal{S} induced by the vectors $\mathbf{a} - \mathbf{a}'$ and $\mathbf{a} - \mathbf{b}$. Obtain $\mathbf{b}_2 = \Pi_{\mathcal{S}}(\mathbf{b})$ which projects \mathbf{b} onto the plane (now you can draw $\mathbf{a}, \mathbf{b}_2, \mathbf{a}', \mathbf{b}'$ on a paper). Since \mathbf{b}_2 is a subspace projection and $P \subset \mathcal{S}$, we still have that \mathbf{b}' is the closest point i.e. $\mathbf{b}' = \Pi_P(\mathbf{b}_2) = \Pi_P(\mathbf{b})$.

Let \mathcal{S}' be the 1 dimensional subspace P lies on. Using the hint observe that

$$\|\mathbf{a} - \mathbf{b}\|_{\ell_2} \geq \|\mathbf{a} - \mathbf{b}_2\|_{\ell_2} \geq \|\Pi_{\mathcal{S}'}(\mathbf{a}) - \Pi_{\mathcal{S}'}(\mathbf{b}_2)\|_{\ell_2} \geq \|\mathbf{a}' - \mathbf{b}'\|_{\ell_2}$$

Solution b. Let us set $d = 2$ and $s = 1$. Consider $\mathbf{a} = [1 \ 0.9]$ and $\mathbf{b} = [0.9 \ 1]$. We have that $\|\mathbf{a} - \mathbf{b}\|_{\ell_2}^2 = 0.02$. On the other hand $\mathbf{a}' = \Pi_{\mathcal{C}}(\mathbf{a}) = [1 \ 0]$ and $\mathbf{b}' = \Pi_{\mathcal{C}}(\mathbf{b}) = [0 \ 1]$ thus $\|\mathbf{a}' - \mathbf{b}'\|_{\ell_2}^2 = 2 > \|\mathbf{a} - \mathbf{b}\|_{\ell_2}^2$. Violation for general $s < d$ can be done with similar construction (except for $s = d$ where sparse set \mathcal{C} is the whole \mathbb{R}^d and convex).

Solution c. The key observation again is that this projection is decomposable i.e. it can be done entrywise. The ℓ_∞ projection is given by

$$\Pi_\infty(\mathbf{x}_i) = \min(|\mathbf{x}_i|, 1) \cdot \text{sgn}(\mathbf{x}_i) \quad \text{for } 1 \leq i \leq d.$$

Solution d. At the special case of $\mathbf{a} = 0$, ℓ_1 subdifferential is exactly the ℓ_∞ norm ball. This connection arises from the fact that ℓ_∞ and ℓ_1 are so-called *dual norms* of each other. Let \mathcal{S} be the set of nonzero coordinates of \mathbf{a} . Noticing entrywise decomposability again, the projection is same as Π_∞ over \mathcal{S}^c and it is simply $\text{sgn}(\mathbf{a}_i)$ over \mathcal{S} as follows

$$\Pi_{\mathcal{C}}(\mathbf{x}_i) = \begin{cases} \min(|\mathbf{x}_i|, 1) \cdot \text{sgn}(\mathbf{x}_i) & i \in \mathcal{S}^c \\ \text{sgn}(\mathbf{a}_i) & i \in \mathcal{S} \end{cases}$$

Vector notation is $\Pi_{\mathcal{C}}(\mathbf{x}) = \text{sgn}(\mathbf{a}) + \Pi_\infty(\mathbf{x}_{\mathcal{S}^c})$ where $\mathbf{x}_{\mathcal{S}^c}$ returns the entries over \mathcal{S}^c and zeros over \mathcal{S} .