

CS 446 / ECE 449 — Homework 1

ziyangx2

September 15, 2021

Instructions.

- Homework is due **Wednesday, September 15, at noon CST**; you have **3** late days in total for **all Homeworks**.
- Everyone must submit individually at gradescope under **hw1** and **hw1code**.
- The “written” submission at **hw1** **must be typed**, and submitted in any format gradescope accepts (to be safe, submit a PDF). You may use L^AT_EX, markdown, google docs, MS word, whatever you like; but it must be typed!
- When submitting at **hw1**, gradescope will ask you to **mark out boxes around each of your answers**; please do this precisely!
- Please make sure your NetID is clear and large on the first page of the homework.
- Your solution **must** be written in your own words. Please see the course webpage for full **academic integrity** information. You should cite any external reference you use.
- We reserve the right to reduce the auto-graded score for **hw1code** if we detect funny business (e.g., your solution lacks any algorithm and hard-codes answers you obtained from someone else, or simply via trial-and-error with the autograder).
- When submitting to **hw1code**, only upload **hw1.py** and **hw1_utils.py**. Additional files will be ignored.

1. Linear Regression.

- (a) Consider a linear regression problem with a dataset containing N data points $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, where $\mathbf{x}^{(i)} \in \mathbb{R}^d$. The accumulated loss function is given by:

$$L_{OLS}(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

where $\mathbf{y} \in \mathbb{R}^N$, $\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$ and $\mathbf{w} \in \mathbb{R}^{d+1}$.

- i. Find the Hessian matrix of $L_{OLS}(\mathbf{w})$.

Hint: You may want to use the fact that $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y})$

- ii. Recall that a twice-continuously differential function $f(\mathbf{x})$ is strictly convex i.f.f. its Hessian is positive definite for all \mathbf{x} . Prove that if N is less than the input dimension d , $L_{OLS}(\mathbf{w})$ can not be strictly convex.
- iii. No matter what \mathbf{X} is, prove that for $\forall \mathbf{w}_1, \mathbf{w}_2 \in \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} L_{OLS}(\mathbf{w})$, we have $\mathbf{X}\mathbf{w}_1 = \mathbf{X}\mathbf{w}_2$. Note that $\mathbf{w}_1, \mathbf{w}_2$ can be different.

Hint: Use the convexity of the loss function and the convex combinations of \mathbf{w}_1 and \mathbf{w}_2 .

- (b) Consider the same dataset with an L2-norm regularization added to the OLS loss function. Linear regression with L2 regularization is also called Ridge regression. Recall the composite loss function of ridge regression:

$$L_{ridge}(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- i. One advantage of ridge regression is that for a positive regularization constant ($\lambda > 0$), the matrix $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})$ is always invertible. Prove that the matrix $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})$ is invertible by showing that it's positive definite.
- ii. Knowing that $L_{ridge}(\mathbf{w})$ is a convex function, show that the estimator for ridge regression is:

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

Solution.

- (a) i. Let us assume:

$$L_{OLS} = \frac{1}{2} \sum_{i=1}^N L_i^2$$

where $L_i = \sum_{j=1}^{d+1} x_{i,j} w_j - y_i$

Take one case:

$$\frac{\partial L}{\partial w_1} = \sum_{i=1}^N L_i \cdot x_{i,1} \text{ and } \frac{\partial^2 L}{\partial w_1 \partial w_2} = \sum_{i=1}^N x_{i,1} \cdot x_{i,2}$$

In general:

$$(H_L)_{a,b} = \frac{\partial^2 L}{\partial w_a \partial w_b} = \sum_{i=1}^N x_{i,a} \cdot x_{i,b}$$

$$H_L = \begin{bmatrix} H_{1,1} & H_{1,2} & & H_{1,d+1} \\ H_{2,1} & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ H_{d+1,1} & \cdots & \cdots & H_{d+1,d+1} \end{bmatrix}$$

$$H_L = \begin{bmatrix} \sum_{i=1}^N x_{i1}x_{i1} & \sum_{i=1}^N x_{i1}x_{i2} & \sum_{i=1}^N x_{i1}x_{i,d+1} \\ \sum_{i=1}^N x_{i2}x_{i1} & \ddots & \vdots \\ \vdots & & \vdots \\ \sum_{i=1}^N x_{i,d+1}x_{i1} & \cdots & \sum_{i=1}^N x_{i,d+1}x_{i,d+1} \end{bmatrix}$$

Clearly therefore:

$$H_L = \mathbf{X}^\top \mathbf{X}$$

Or more mathematically,

$$H_L = \nabla_{\mathbf{w}}^2 L_{OLS}(\mathbf{w}) = \frac{\partial^2 L}{\partial \mathbf{w}^2} = \frac{\partial}{\partial \mathbf{w}} \mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{X}^\top \mathbf{X}$$

ii. Proof:

A symmetric matrix \mathbf{M} is positive definite if given a column vector \mathbf{v} , it satisfies $\mathbf{v}^\top \mathbf{M} \mathbf{v} > 0$

In this problem,

$$\mathbf{w}^\top H_L \mathbf{w} = \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} = (\mathbf{X}\mathbf{w})^\top \mathbf{X}\mathbf{w} = \|\mathbf{X}\mathbf{w}\|_2^2 \geq 0$$

Since N is less than d , there must be a vector \mathbf{w}' in the null space of \mathbf{X}

s.t. $\mathbf{X}\mathbf{w}' = 0$, i.e. $\|\mathbf{X}\mathbf{w}'\|_2 = 0$

Therefore, if N is less than the input dimension d , $L_{OLS}(\mathbf{w})$ can not be strictly convex.

iii. Proof:

$$\|\mathbf{X}\mathbf{w}' - \mathbf{y}\|_2^2$$

$$= \|\mathbf{X}\mathbf{w}' - \mathbf{X}\mathbf{w} + \mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

$$= \|\mathbf{X}\mathbf{w}' - \mathbf{X}\mathbf{w}\|^2 + 2(\mathbf{X}\mathbf{w}' - \mathbf{X}\mathbf{w})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) + \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

$$\text{Since } (\mathbf{X}\mathbf{w}' - \mathbf{X}\mathbf{w})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = (\mathbf{w}' - \mathbf{w})^\top (\mathbf{X}^\top \mathbf{X}\mathbf{w} - \mathbf{X}^\top \mathbf{y}) = 0$$

$$\|\mathbf{X}\mathbf{w}' - \mathbf{y}\|_2^2 = \|\mathbf{X}\mathbf{w}' - \mathbf{X}\mathbf{w}\|^2 + \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

Therefore,

$$\|\mathbf{X}\mathbf{w}' - \mathbf{X}\mathbf{w}\|^2 = 0$$

indicating that $\mathbf{X}\mathbf{w}' = \mathbf{X}\mathbf{w}$

. Also, according to the property of convexity: given two different \mathbf{w}_1 and \mathbf{w}_2 , if both of them belongs to $\arg \min_{\mathbf{w}} L_{OLS}(\mathbf{w})$, we have $\mathbf{X}\mathbf{w}_1 = \mathbf{X}\mathbf{w}_2$.

(b) i. Proof:

$$\mathbf{w}^\top (\mathbf{X}^\top \mathbf{x} + \lambda \mathbf{I}) \mathbf{w} = (\mathbf{X}\mathbf{w})^\top \mathbf{X}\mathbf{w} + \mathbf{w}^\top \lambda \mathbf{I} \mathbf{w} = \|\mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 > 0$$

By definition,

the matrix $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})$ is always positive definite for $\lambda > 0$.

Hence, the matrix $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})$ is invertible.

ii. Proof:

$$L_{ridge}(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

$$\frac{\partial L_{ridge}(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w} = 0$$

$$\mathbf{X}^\top \mathbf{X} \hat{\mathbf{w}} - \mathbf{X}^\top \mathbf{y} + \lambda \hat{\mathbf{w}} = 0$$

$$(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \hat{\mathbf{w}} = \mathbf{X}^\top \mathbf{y}$$

Therefore, the estimator for ridge regression is:

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

2. Programming - Linear Regression.

Recall that the empirical risk in the linear regression method is defined as $\hat{\mathcal{R}}(\mathbf{w}) := \frac{1}{2N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}^{(i)} - y^{(i)})^2$, where $\mathbf{x}^{(i)} \in \mathbb{R}^d$ is a data point and $y^{(i)}$ is an associated label.

- (a) **Implement linear regression using gradient descent in the `linear_gd(X, Y, lrate, num_iter)` function of `hw1.py`.**

The arguments for this function are: `X` as the training features, a tensor with shape $N \times d$; `Y` as the training labels, an $N \times 1$ tensor; `lrate` as learning rate (step size), a float number; and `num_iter` as the number of iterations for gradient descent to run. The objective of this function is to find parameters \mathbf{w} that minimize the empirical risk $\hat{\mathcal{R}}(\mathbf{w})$ using gradient descent (only gradient descent). To keep consistent with the standard program and get correctly scored, **prepend** a column of ones to `X` in order to accommodate the bias term in \mathbf{w} , thus \mathbf{w} should have $d + 1$ entries. Then use $\mathbf{w} = 0$ as the initial parameters, and return

Hint. If you are new to machine learning or programming with pytorch, we offer some kind suggestions. First, try using the vector/matrix operations provided in pytorch and avoid using for-loops. This will improve both the efficiency and style of your program. Second, create your own test cases for debugging before submission. With very few samples in your own test case, it is convenient to compare the program output with your manual calculation. Third, to avoid matrix computation error, remember to check the shapes of tensors regularly.

Library routines: `torch.matmul` (`@`), `torch.tensor.shape`, `torch.tensor.t`, `torch.cat`, `torch.ones`, `torch.zeros`, `torch.reshape`.

- (b) **Implement linear regression by using the pseudo inverse to solve for \mathbf{w} in the `linear_normal(X, Y)` function of `hw1.py`.**

The arguments for this function are: `X` as the training features, a tensor with shape $N \times d$ tensor; `Y` as the training labels, an $N \times 1$ tensor. To keep consistent with the standard program and get correctly scored, **prepend** a column of ones to `X` in order to accommodate the bias term in \mathbf{w} , thus \mathbf{w} should have $d + 1$ entries.

Library routines: `torch.matmul` (`@`), `torch.cat`, `torch.ones`, `torch.pinv`.

- (c) **Implement the `plot_linear()` function in `hw1.py`.** Follow the steps below.

Use the provided function `hw1_utils.load_reg_data()` to generate a training set `X` and training labels `Y`. Then use `linear_normal()` to calculate the regression results \mathbf{w} . Eventually plot the points of dataset and regressed curve. Return the plot as output. Note that `plot_linear()` should return the figure object and you should **include the visualization in your written submission**.

Hint. If you are new to plotting machine learning visualizations, we offer some kind suggestions. `matplotlib.pyplot` is an “extremely” useful tool in machine learning, and we commonly refer to it as `plt`. Please first get to know the most basic usages by examples from its official website (such as scatter plots, line plots, etc.). As for our programming question specifically, you may divide and conquer it by first plotting the points in the dataset, then plotting the linear regression curve.

Library routines: `torch.matmul` (`@`), `torch.cat`, `torch.ones`, `plt.plot`, `plt.scatter`, `plt.show`, `plt.gcf` where `plt` refers to the `matplotlib.pyplot` library.

Solution.

- (a) code
(b) code
(c) The plot is shown below:

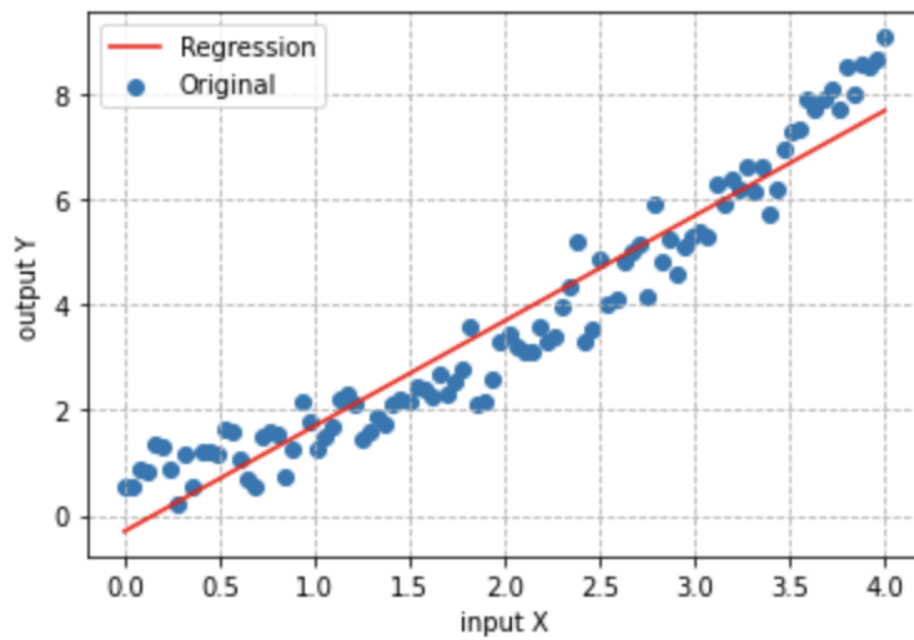


Figure 1: Solution to Problem 2(c).

3. Programming - Logistic Regression.

Recall the empirical risk $\hat{\mathcal{R}}$ for logistic regression (as presented in lecture 3):

$$\hat{\mathcal{R}}_{\log}(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)})).$$

Here you will minimize this risk using gradient descent.

- (a) In your **written submission**, derive the gradient descent update rule for this empirical risk by taking the gradient. Write your answer in terms of the learning rate (step size) η , previous parameters \mathbf{w} , new parameters \mathbf{w}' , number of examples N , and training examples $\mathbf{x}^{(i)}$. Show all of your steps.
- (b) Implement the `logistic()` function in `hw1.py`. You are given as input a training set `X`, training labels `Y`, a learning rate (step size) `lr`, and number of gradient updates `num_iter`. Implement gradient descent to find parameters \mathbf{w} that minimize the empirical risk $\hat{\mathcal{R}}_{\log}(\mathbf{w})$. Perform gradient descent for `num_iter` updates with a learning rate (step size) of `lr`. Same as previous questions, initialize $\mathbf{w} = 0$, return \mathbf{w} as output, and prepend `X` with a column of ones.

Library routines: `torch.matmul` (`@`), `torch.tensor.t`, `torch.exp`.

- (c) Implement the `logistic_vs_ols()` function in `hw1.py`. Use `hw1_utils.load_logistic_data()` to generate a training set `X` and training labels `Y`. Run `logistic(X,Y)` from part (b) taking `X` and `Y` as input to obtain parameters \mathbf{w} (use the defaults for `num_iter` and `lr`). Also run `linear_gd(X,Y)` from Problem 2 to obtain parameters \mathbf{w} . Plot the decision boundaries for your logistic regression and least squares models along with the data `X`. [Note: As we learned in the class that the decision rule of Least Squares and Logistic Regression for predicting the class label is $\text{sign}(\hat{\mathbf{w}}^\top \mathbf{x})$, the decision boundary can be obtained from $\hat{\mathbf{w}}^\top \mathbf{x} = 0$, i.e., for $d = 2$, we have $x_2 = -(\hat{w}_0 + \hat{w}_1 \times x_1) / \hat{w}_2$.] Include the visualizations in your **written submission**. Which model appears to classify the data better? Explain in the **written submission** that why you believe it is better for this problem.

Library routines: `torch.linspace`, `plt.scatter`, `plt.plot`, `plt.show`, `plt.gcf`.

Solution.

(a)

$$\begin{aligned} \nabla_{\mathbf{w}} f(\mathbf{w}) &= \frac{\partial \hat{\mathcal{R}}_{\log}(\mathbf{w})}{\partial \mathbf{w}} \\ \nabla_{\mathbf{w}} f(\mathbf{w}) &= \frac{1}{N} \sum_{i=1}^N \frac{\exp(-y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)})}{1 + \exp(-y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)})} (-y^{(i)} \mathbf{x}^{(i)}) \\ \mathbf{w}' &= \mathbf{w} - \eta \cdot \nabla_{\mathbf{w}} f(\mathbf{w}_t) \end{aligned}$$

Gradient descent update rule:

$$\mathbf{w}' = \mathbf{w} - \eta \cdot \frac{1}{N} \sum_{i=1}^N \frac{\exp(-y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)})}{1 + \exp(-y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)})} (-y^{(i)} \mathbf{x}^{(i)})$$

(b) code

(c) The plot is shown below:

From the plot, we can find that for this dataset, the logistics model performs better than linear regression model according to the maximum margin principle. Personally speaking, logistics model fits more for those binary classification problems.

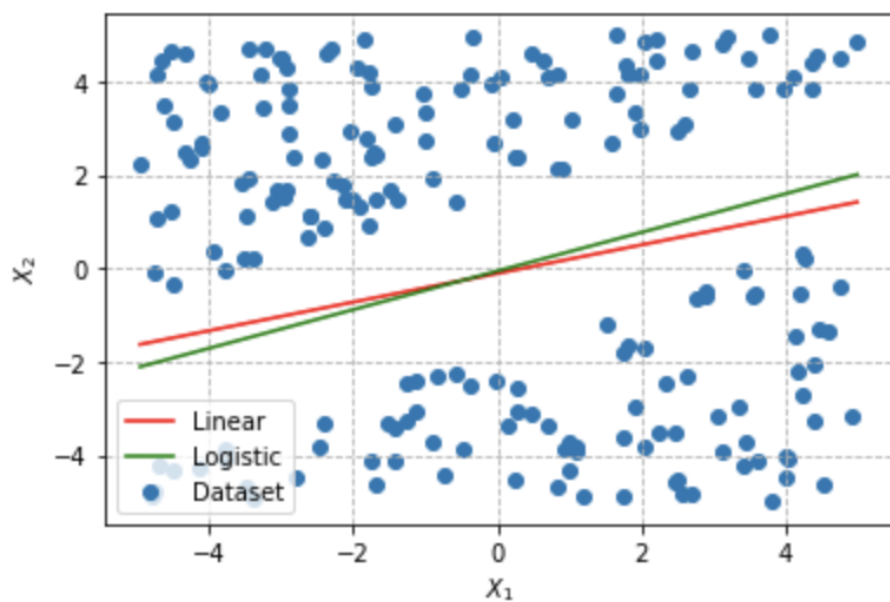


Figure 2: Solution to Problem 3(c).

4. Convexity, Lipschitz Continuity, and Smoothness

(a) Convexity

- i. Show that if a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, then for any matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ and vector $\mathbf{b} \in \mathbb{R}^n$, the function $g : \mathbb{R}^m \rightarrow \mathbb{R}$ defined by $g(\mathbf{x}) = f(\mathbf{A}\mathbf{x} + \mathbf{b})$ is convex, where $\mathbf{x} \in \mathbb{R}^m$.
- ii. Prove that if the differentiable function f is λ -strongly convex and the differentiable function g is convex then $f + g$ is λ -strongly convex.
- iii. Given m convex functions $\{f_i : \mathbb{R}^n \rightarrow \mathbb{R}\}_{i=1}^m$, denote

$$f(\mathbf{x}) = \max_{i \in [m]} f_i(\mathbf{x}),$$

where $[m] = \{1, 2, \dots, m\}$. Prove that f is convex.

(b) Lipschitzness and Smoothness

We say a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^d$ is ρ -Lipschitz if $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$, it holds that $\|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2 \leq \rho \|\mathbf{x}_1 - \mathbf{x}_2\|_2$.

- i. Prove that if $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g : \mathbb{R}^m \rightarrow \mathbb{R}^d$ are ρ -Lipschitz functions, then the composite $g \circ f : \mathbb{R}^n \rightarrow \mathbb{R}^d$ defined by $(g \circ f)(\mathbf{x}) = g(f(\mathbf{x}))$ is ρ^2 -Lipschitz.
- ii. Given a differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ whose gradient is β -Lipschitz, prove that for $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ we have

$$f(\mathbf{y}) - f(\mathbf{x}) \leq \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

Hint: You are not required to follow the hints, but please consider them if you have no idea for proof. (1) Define a tool function $\phi(t) = f((1-t)\mathbf{x} + t\mathbf{y})$, thus $f(\mathbf{y}) - f(\mathbf{x}) = \phi(1) - \phi(0) = \int_0^1 \phi'(t) dt$ (figure it out); (2) If you get stuck at the final steps, taking a look at the Cauchy-Schwarz inequality may be helpful.

Solution.

(a) i. Proof:

$$\begin{aligned} g(\alpha \mathbf{x} + (1-\alpha)\mathbf{x}') &= f(\mathbf{A}(\alpha \mathbf{x} + (1-\alpha)\mathbf{x}') + \mathbf{b}) \\ g(\alpha \mathbf{x} + (1-\alpha)\mathbf{x}') &= f(\mathbf{A}(\alpha \mathbf{x} + (1-\alpha)\mathbf{x}') + \alpha \mathbf{b} + (1-\alpha)\mathbf{b}) = f(\alpha(\mathbf{A}\mathbf{x} + \mathbf{b}) + (1-\alpha)(\mathbf{A}\mathbf{x}' + \mathbf{b})) \\ \text{Since function } f \text{ is convex,} \\ f(\alpha(\mathbf{A}\mathbf{x} + \mathbf{b}) + (1-\alpha)(\mathbf{A}\mathbf{x}' + \mathbf{b})) &\leq \alpha f(\mathbf{A}\mathbf{x} + \mathbf{b}) + (1-\alpha)f(\mathbf{A}\mathbf{x}' + \mathbf{b}) \\ &= \alpha g(\mathbf{x}) + (1-\alpha)g(\mathbf{x}') \end{aligned}$$

We can get:

$$g(\alpha \mathbf{x} + (1-\alpha)\mathbf{x}') \leq \alpha g(\mathbf{x}) + (1-\alpha)g(\mathbf{x}')$$

By definition, the function $g(\mathbf{x})$ is convex as well.

ii. Proof:

f is λ -strongly convex and g is convex, indicating that:

$$\begin{aligned} f(\mathbf{x}') &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{x}' - \mathbf{x}) + \frac{\lambda}{2} \|\mathbf{x}' - \mathbf{x}\|^2 \\ g(\mathbf{x}') &\geq g(\mathbf{x}) + \nabla g(\mathbf{x})^\top (\mathbf{x}' - \mathbf{x}) \end{aligned}$$

\therefore Combine:

$$f(\mathbf{x}') + g(\mathbf{x}') \geq (f(\mathbf{x}) + g(\mathbf{x})) + \nabla(f(\mathbf{x}) + g(\mathbf{x}))^\top (\mathbf{x}' - \mathbf{x}) + \frac{\lambda}{2} \|\mathbf{x}' - \mathbf{x}\|^2$$

Hence, $f + g$ is λ -strongly convex as well.

iii. Proof:

Since all the m functions are convex,

$$\max_i^m f_i(\alpha \mathbf{x} + (1-\alpha)\mathbf{x}') \leq \max_i^m \alpha f_i(\mathbf{x}) + \max_i^m (1-\alpha) f_i(\mathbf{x}') = \alpha \max_i^m f_i(\mathbf{x}) + (1-\alpha) \max_i^m f_i(\mathbf{x}')$$

Therefore,

$$\max_i^m f_i(\alpha \mathbf{x} + (1-\alpha)\mathbf{x}') \leq \alpha \max_i^m f_i(\mathbf{x}) + (1-\alpha) \max_i^m f_i(\mathbf{x}')$$

By definition, the function $f(\mathbf{x}) = \max_i^m f_i(\mathbf{x})$ is convex as well.

(b) i. Proof:

Since function f and g are both ρ -Lipschitz, then

$$\|g(f(\mathbf{x}_1)) - g(f(\mathbf{x}_2))\|_2 \leq \rho \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2$$

$$\|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2 \leq \rho \|\mathbf{x}_1 - \mathbf{x}_2\|_2$$

Therefore,

$$\|g(f(\mathbf{x}_1)) - g(f(\mathbf{x}_2))\|_2 \leq \rho^2 \|\mathbf{x}_1 - \mathbf{x}_2\|_2$$

indicating that function g is ρ^2 -Lipschitz with respect to \mathbf{x} .

ii. Proof:

Suppose $\phi(t) = f((1-t)\mathbf{x} + t\mathbf{x}')$,

and $\phi'(t) = \nabla f(t\mathbf{x}' + (1-t)\mathbf{x})^\top (\mathbf{x}' - \mathbf{x})$

then $\phi'(0) = \nabla f(\mathbf{x})^\top (\mathbf{x}' - \mathbf{x})$

then $f(\mathbf{x}') - f(\mathbf{x}) = \phi(1) - \phi(0) = \int_0^1 \nabla f(t\mathbf{x}' + (1-t)\mathbf{x})^\top (\mathbf{x}' - \mathbf{x}) dt$

$f(\mathbf{x}') - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{x}' - \mathbf{x})$

$= \phi(1) - \phi(0) - \phi'(0) \cdot 1$

$= \int_0^1 \nabla f(t\mathbf{x}' + (1-t)\mathbf{x})^\top (\mathbf{x}' - \mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{x}' - \mathbf{x}) dt$

$= \int_0^1 (\nabla f(t\mathbf{x}' + (1-t)\mathbf{x}) - \nabla f(\mathbf{x}))^\top (\mathbf{x}' - \mathbf{x}) dt$

$\leq \int_0^1 \|\nabla f(t\mathbf{x}' + (1-t)\mathbf{x}) - \nabla f(\mathbf{x})\|_2 \cdot \|\mathbf{x}' - \mathbf{x}\|_2 \cdot dt$, by Cauchy-Schwarz inequality

Gradient of f is β -Lipschitz, indicating that

$\|\nabla f(t\mathbf{x}' + (1-t)\mathbf{x}) - \nabla f(\mathbf{x})\|_2 \leq \beta \|t(\mathbf{x}' - \mathbf{x})\|_2$

therefore,

$f(\mathbf{x}') - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{x}' - \mathbf{x})$

$\leq \int_0^1 \beta \|t(\mathbf{x}' - \mathbf{x})\|_2 \cdot \|\mathbf{x}' - \mathbf{x}\|_2 dt$

$= \|\mathbf{x}' - \mathbf{x}\|_2^2 \cdot \int_0^1 \beta t dt$

$= \|\mathbf{x}' - \mathbf{x}\|_2^2 \cdot \frac{\beta}{2}$

Hence,

$f(\mathbf{x}') - f(\mathbf{x}) \leq \nabla f(\mathbf{x})^\top (\mathbf{x}' - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2$