

Tutorial 3
Xiaojian Xu
02/15/2019

ESE415 OPTIMIZATION

Guideline



HW2 Recap



HW3 Recitation



Partial Derivative, Derivative,
Gradient, Directional
Derivative and Hessian

1. HW2 Recap

Lipschitz Constant vs. Largest Eigenvalue of Hessian Matrix

Problem 5

We consider an optimization problem that often arises in machine learning. We would like to assign a label $y \in \{-1, +1\}$ to a given input vector $\mathbf{x} \in \mathbb{R}^n$. For example, \mathbf{x} might denote a vectorized image with n pixels and y a class it corresponds to (e.g., $+1$ for *bird* and -1 for *not bird*). We seek to build a linear model

$$\varphi(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \cdots + w_n x_n = \mathbf{x}^\top \mathbf{w}_1 + w_0 = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}^\top \mathbf{w}, \quad (1)$$

where $\mathbf{w} = (w_1, w_0) \in \mathbb{R}^{n+1}$, with $\mathbf{w}_1 = (w_1, \dots, w_n) \in \mathbb{R}^n$, represent the weights. We train the model on a dataset of m data points $\{\mathbf{x}_i, y_i\}_{i=1, \dots, m}$, stored as a label vector $\mathbf{y} \in \{-1, +1\}^m$ and a data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \in \mathbb{R}^{n \times m}$. The columns of \mathbf{X} correspond to the input data-vectors \mathbf{x}_i . The training is performed by minimizing the following cost function

$$f(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^m (y_i - \varphi(\mathbf{x}_i; \mathbf{w}))^2 = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{w}\|^2, \quad (2)$$

where $\mathbf{A} = [\mathbf{X}^\top \mid \mathbf{1}] \in \mathbb{R}^{m \times (n+1)}$ is transposed and augmented version of the data-matrix, where every row is made of an input data-vector followed by a constant 1 as in eq. (1). Given a minimizer $\mathbf{w}^* \in \mathbb{R}^{n+1}$ of f , we can predict the label for any vector \mathbf{x} as

$$\hat{y} = \text{sgn}(\mathbf{x}^\top \mathbf{w}_1 + w_0) \quad \text{where} \quad \text{sgn}(x) \triangleq \begin{cases} -1 & \text{if } x < 0 \\ 0 & \text{if } x = 0 \\ +1 & \text{if } x > 0 \end{cases}$$

In this exercise, we will work with the dataset stored in `dataset.mat`. This file along with the initial script are stored in the archive `assignment2.zip`.

- Download, run, and study `ClassificationGM.m` for minimizing eq. (2) via the *gradient method (GM)*. Run the code with parameters $\gamma = \text{stepSize} = 10^{-4}$, $\text{maxIter} = 100$, and $\text{tol} = 10^{-6}$, and the random initialization $\mathbf{w}^0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ (i.e., the initialization consists of independent and identically distributed Normal random variables). What is the final value of f ? Include the generated plot with your solutions.
- Provide the expression for the gradient and the Hessian of f . Comment on the convexity of f .
- Analytically find the minimizer of f and comment on the uniqueness of the solution.
- Modify the script to implement *GM with backtracking line-search* (as described on Page 4 of Lecture 7). Run the algorithm using the step reduction parameter $\beta = 0.5$, initial $\gamma = \text{stepSize} = 1$, and fix all the other parameters as in (a). Comment on the convergence of the algorithm. What is the final value of f ? What is the final value of γ ? Include the generated plot with your solutions. Additionally, print and submit your code with your solutions.
- Compute the Lipschitz constant L of $\nabla f(\mathbf{w})$. What is the range of possible step-sizes for this problem?
- (Bonus) Modify the script in order to implement the *accelerated gradient method (AGM)*, with the following update rule for the acceleration parameter

$$\theta_0 = 1 \quad \text{and} \quad \theta_t = \frac{1 + \sqrt{1 + 4\theta_{t-1}^2}}{2}.$$

Run AGM using a fixed step size $\gamma = \text{stepSize} = 10^{-3}$ and fix all the other parameters as in (a). Comment on the convergence of the method. What is the final value of f ? Include the generated plot with your solutions. Additionally, print and submit your code with your solutions.

HW2 Recap: problem 5 (e)

- In this case, Lipschitz constant L can be set as the largest eigenvalue of the hessian matrix $Hf(w)$. Here we show the proofs. The proofs contain three parts using following hints .

Hints(a)

Theorem 4 (Spectral decomposition). Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix. Then

- All of the eigenvalues of A are real.
- There are n mutually orthogonal, unit-norm eigenvectors u_1, \dots, u_n corresponding to n eigenvalues $\lambda_1, \dots, \lambda_n$ (with repetitions included). Therefore, the matrix A can be decomposed in the form $A = U\Lambda U^\top$, where $U = [u_1 \cdots u_n]$ is an orthogonal matrix and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$.
- We can represent A as

$$A = \sum_{i=1}^n \lambda_i u_i u_i^\top.$$

Hints(b)

Definition 3. A function f has a *Lipschitz continuous gradient* with constant $L > 0$, when

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|,$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. We will denote this class of functions by \mathcal{C}_L^1 .

a). First we show that for a symmetric matrix A , we
 $\|Ax\| = \lambda_{\max}\|x\|$ where λ_{\max} is maximum eigenvalue

*Note that a real symmetric matrix is diagonalizable by an orthogonal matrix.

\Rightarrow So we could have $A = U\Lambda U^T$ (as shown in Theorem 4 (2) above)

where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$.

∴ For $\forall x$,

$$\begin{aligned} x^T A x &= x^T U \Lambda U^T x \\ &= (U^T x)^T \cdot \Lambda \cdot (U^T x), \end{aligned}$$

Putting $y = U^T x \Rightarrow x^T A x = y^T \Lambda y$

$$\begin{aligned} \|Ax\|^2 &= (Ax)^T Ax = x^T A^T A x = x^T A^2 x = x^T (\Lambda U^T)(U \Lambda U^T) x \\ &= x^T U \Lambda^2 U^T x = (U^T x)^T \Lambda^2 \cdot (U^T x) \end{aligned}$$

Putting $y = U^T x$

$$\begin{aligned} \|Ax\|^2 &= y^T \Lambda^2 y = \sum_i \lambda_i^2 y_i^2 \\ &\leq \sum_i \lambda_{\max}^2 y_i^2 \\ &= \lambda_{\max}^2 \|y\|^2 \end{aligned}$$

$$\therefore \|y\|^2 = y^T y = (U^T x)^T \cdot (U^T x) = x^T U \cdot U^T x = x^T x = \|x\|^2$$

$$\Rightarrow \|Ax\|^2 \leq \lambda_{\max}^2 \|y\|^2 = \lambda_{\max}^2 \|x\|^2$$

$$\therefore \|Ax\| \leq \lambda_{\max} \|x\|$$

b). now we show that for function

$$f(w) = \frac{1}{2} \|y - Aw\|^2,$$

we have $\|A^T A(w_1 - w_2)\| \leq L \|w_1 - w_2\|, w_1, w_2 \in \mathbb{R}^n$

Proof:

$$f(w) = \frac{1}{2} \|y - Aw\|^2$$

$$\nabla f(w) = A^T(Aw - y) \quad (\text{gradient})$$

$$Hf(w) = A^T A \quad (\text{Hessian})$$

As shown in lecture notes, $f(w)$ has a Lipschitz continuous gradient with constant $L > 0$, when

$$\|\nabla f(w) - \nabla f(w')\| \leq L \|w_1 - w_2\|, \quad \forall w_1, w_2 \in \mathbb{R}^n$$

Plug in $\nabla f(w) = A^T(Aw - y)$, \Rightarrow

$$\|A^T(Aw_1 - y) - A^T(Aw_2 - y)\| \leq L \|w_1 - w_2\|$$

$$\|A^T A(w_1 - w_2)\| \leq L \|w_1 - w_2\|$$

c). Now we know that for function

$$f(w) = \frac{1}{2} \|y - Aw\|^2$$

setting $w_1 - w_2 \rightarrow x$, $A^T A \rightarrow A$.

$$\text{from (a)} \Rightarrow \|A^T A(w_1 - w_2)\| \leq \lambda_{\max} \|w_1 - w_2\| \quad \cdots \#1$$

λ_{\max} is the largest eigenvalue of $A^T A$.

$$\text{from (b)} \Rightarrow \|A^T A(w_1 - w_2)\| \leq L \|w_1 - w_2\| \quad \cdots \#2$$

L is the Lipschitz constant of $f(w)$.

\Rightarrow Set L_0 as the smallest Lipschitz constant.

then we could set $L_0 = \lambda_{\max}$, $L \geq \lambda_{\max}$

Proofs

2. HW2 Recitation (P1)

Problem

The goal of this problem is to give some intuition about different rates of convergence of $\{f(\mathbf{x}^t)\}_{t \in \mathbb{N}}$ to $f^* = f(\mathbf{x}^*)$. As a reminder, the rate of convergence is characterized via two values $p \in \mathbb{N}$ and $\beta \in [0, \infty)$. Specifically, the goal is to find the largest p such that

$$\beta = \limsup_{t \rightarrow \infty} \frac{|f(\mathbf{x}^t) - f^*|}{|f(\mathbf{x}^{t-1}) - f^*|^p} < \infty.$$

Let $\delta \in (0, 1)$ and consider $r_t \triangleq |f(\mathbf{x}^t) - f^*|$. Categorize the convergence of the following sequences as *sublinear*, *linear*, *superlinear*, and *quadratic*. Justify and explain your answer.

a) $r_t = \frac{\delta}{t}$

b) $r_t = \delta^t$

c) $r_t = \delta^{t^2}$

d) $r_t = \delta^{2^t}$

Hints

Definition 2. Let $\{f_k\}$ be a convergent sequence with $f_k \rightarrow f$. The *order of convergence* of $\{f_k\}$ is

$$p^* = \sup \left\{ p \in \mathbb{N} : \limsup_{k \rightarrow \infty} \frac{|f_k - f|}{|f_{k-1} - f|^p} < \infty \right\}.$$

The larger p^* implies faster convergence.

Definition 3. Let p denote the order of convergence of $\{f_k\}$ to f and let

$$\beta = \limsup_{k \rightarrow \infty} \frac{|f_k - f|}{|f_{k-1} - f|^p}.$$

We say that the convergence is

- *linear* or *geometric* if $p = 1$ and $\beta \in (0, 1)$.
- *superlinear* if $p = 1$ and $\beta = 0$.
- *sublinear* if $p = 1$ and $\beta = 1$.
- *quadratic* if $p = 2$ and $\beta < \infty$.

2. HW3 Recitation (P1)

Solution

- a) The convergence is *sublinear* as for $p = 1$, we have that

$$\frac{r_t}{r_{t-1}} = \frac{\frac{\delta}{t}}{\frac{\delta}{t-1}} = \frac{t-1}{t} = 1 - \frac{1}{t} \xrightarrow[t \rightarrow \infty]{} 1 = \beta.$$

Note that for any other $p > 1$, the value given by \limsup will be unbounded.

3. Partial derivative, derivative, gradient, directional derivative and Hessian

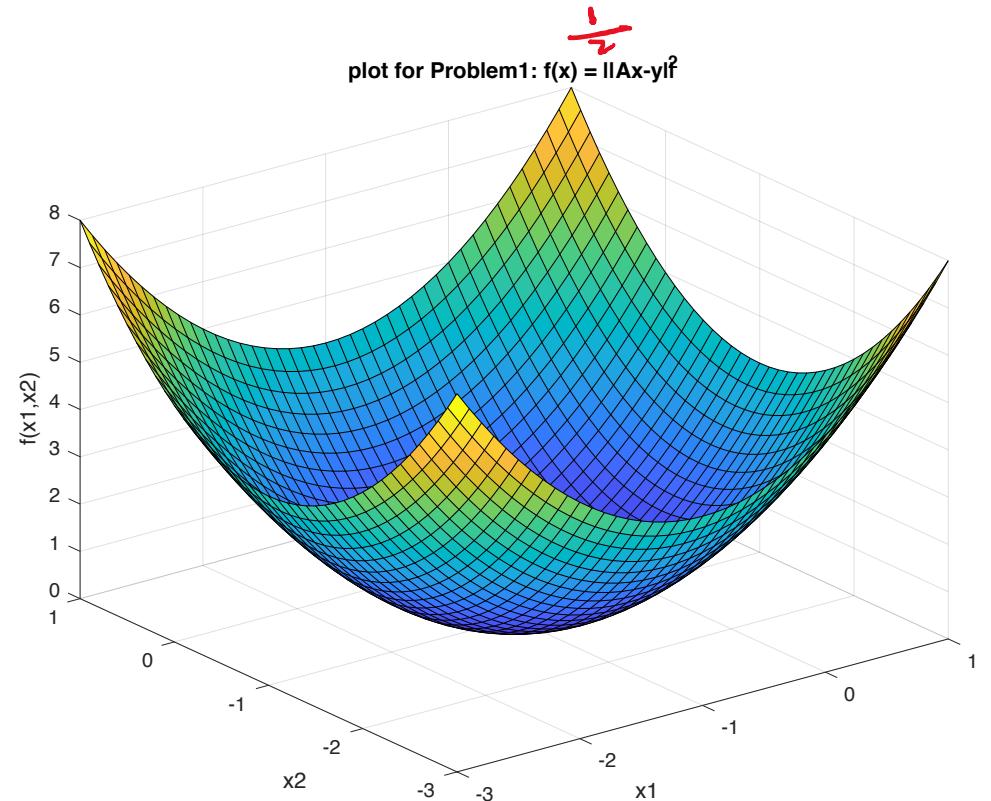
1. We will consider functions of form

$$f(\mathbf{x}) \quad \text{where} \quad \mathbf{x} \in \mathbb{R}^n$$

2. Vectors will be generally assumed to be laid out in the column form

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n \quad \text{and} \quad \mathbf{x}^T = [x_1 \dots x_n],$$

where the second notation indicates the transpose.



$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|^2 \quad \text{where} \quad \mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} -2 \\ 0 \end{bmatrix} \quad \text{with} \quad \mathbf{x} \in \mathbb{R}^2.$$

3.1 Partial Derivative

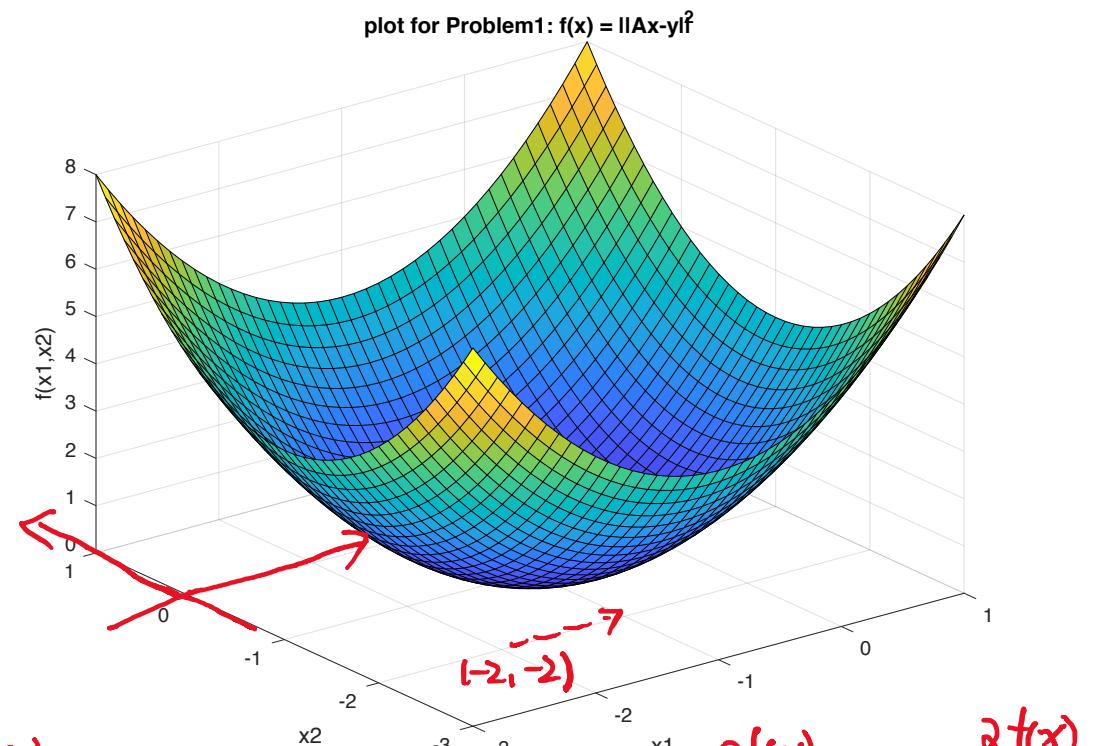
3. We denote and define the partial derivative of f with respect to coordinate x_i as

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \triangleq \lim_{\delta \rightarrow 0} \frac{f(\mathbf{x} + \delta e_i) - f(\mathbf{x})}{\delta},$$

where $e_i \in \mathbb{R}^n$ is a vector that is 0 in all coordinates except at i , where it is 1.

Highlight:

- Partial derivative implies the changing rate on one coordinate direction at a certain point .
- Its sign implies the increasement or reduction of the function and its magnitude implies the changing speed (rate).



$$\frac{\partial f(x)}{\partial x_1} = 2x_1 + 2$$

$$\frac{\partial f(x)}{\partial x_2} = 2x_2 + 2$$

at point $(-2, -2) \Rightarrow$

$$\frac{\partial f(x)}{\partial x_1} = -2 = \frac{\partial f(x)}{\partial x_2}$$

3.2 Derivative

4. We denote the derivative of a vector-valued function $\mathbf{y}(x) \in \mathbb{R}^m$ with respect to $x \in \mathbb{R}^n$

$$\frac{\partial \mathbf{y}}{\partial x} = D\mathbf{y}(x) \triangleq \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}.$$

Highlight:

- Derivative here is defined for the functions $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$;
- Set $y_i = g_i(x)$, each row can be regarded as the gradient of $\nabla g_i(x)^T$;
- When $m = 1$, it becomes gradient.

Example:

$$f(x) = Ax, \quad A \in \mathbb{R}^{m \times n}, \quad x \in \mathbb{R}^{n \times 1}$$

$$\Rightarrow f(x) \in \mathbb{R}^{m \times 1} \quad (\text{vector})$$

$$A = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$f(x) = [x_1 + 4x_2, 2x_1 + 5x_2, 3x_1 + 6x_2]^T$$

3.3 Gradient

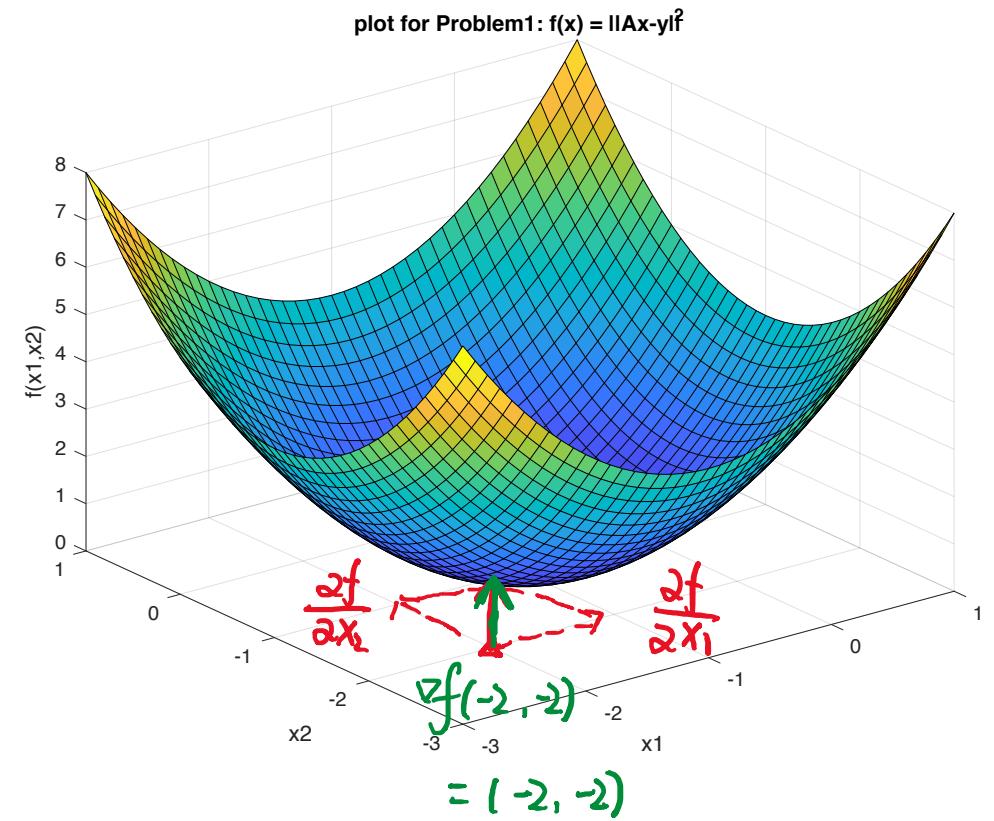
5. Gradient of f at $\mathbf{x} \in \mathbb{R}^n$

$$\nabla f(\mathbf{x}) = [\mathbf{D}f(\mathbf{x})]^T = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix},$$

where we adopted a column-layout convention for convenience.

Highlight:

- Gradient here is defined for the functions $f: \mathbb{R}^n \rightarrow \mathbb{R}$;
- It is a vector consists of partial derivative.
- It implies the direction that function f changes fastest.



3.4 Directional Derivative

Definition 1. The directional derivative of $f(\mathbf{x})$ along $\mathbf{u} \in \mathbb{R}^n$, $\|\mathbf{u}\| = 1$, is

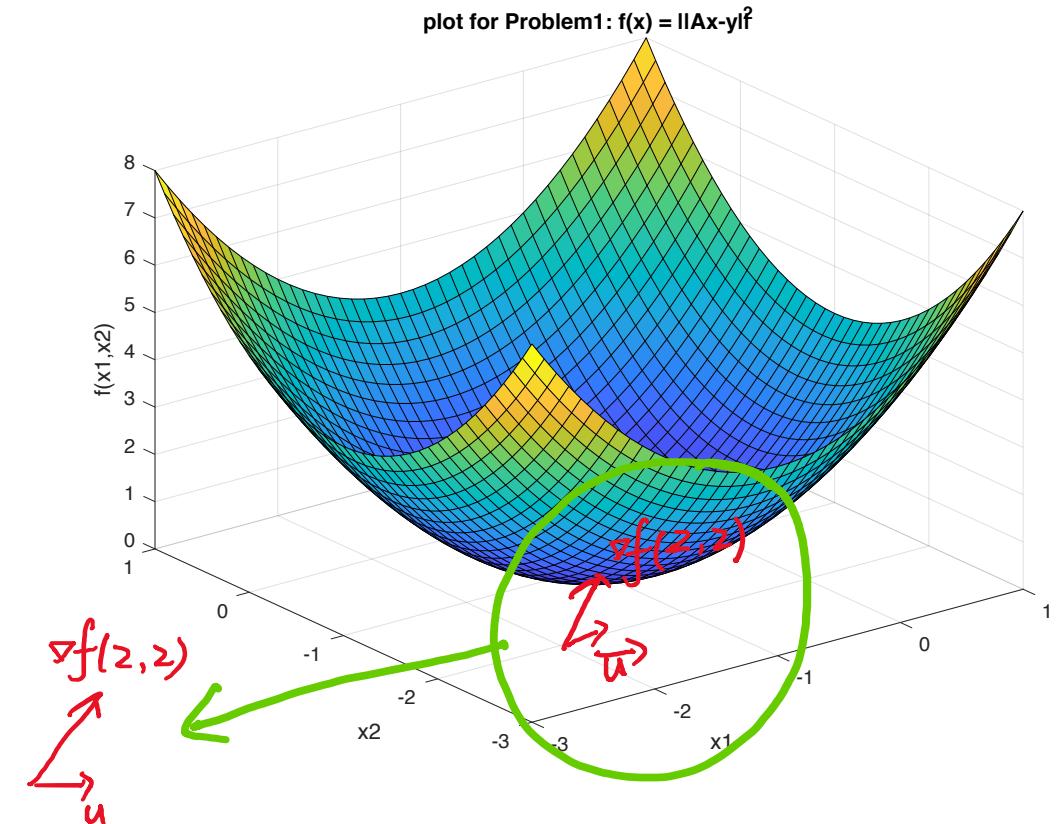
$$\nabla_{\mathbf{u}} f(\mathbf{x}) = \lim_{\gamma \rightarrow 0} \frac{f(\mathbf{x} + \gamma \mathbf{u}) - f(\mathbf{x})}{\gamma} = \nabla f(\mathbf{x})^T \mathbf{u}.$$

Proof. We consider the Taylor's expansion $f(\mathbf{x} + \gamma \mathbf{g}) = f(\mathbf{x}) + \gamma \nabla f(\mathbf{x})^T \mathbf{g} + o(\gamma)$, where $o(\gamma)/\gamma \rightarrow 0$ when $\gamma \rightarrow 0$. Then, we have that

$$\frac{f(\mathbf{x} + \gamma \mathbf{g}) - f(\mathbf{x})}{\gamma} = \nabla f(\mathbf{x})^T \mathbf{g} + \frac{o(\gamma)}{\gamma}.$$

Highlight:

- Directional derivative is a scalar;
- It shows the projection of the gradient on a certain direction \mathbf{u} ;
- It measures the changing speed at the direction \mathbf{u} .
- $\mathbf{u} = -\nabla f(\mathbf{x}) / \| \nabla f(\mathbf{x}) \|$ points in the direction of the fastest local decrease of f (this is how gradient method works)



3.5 Hessian

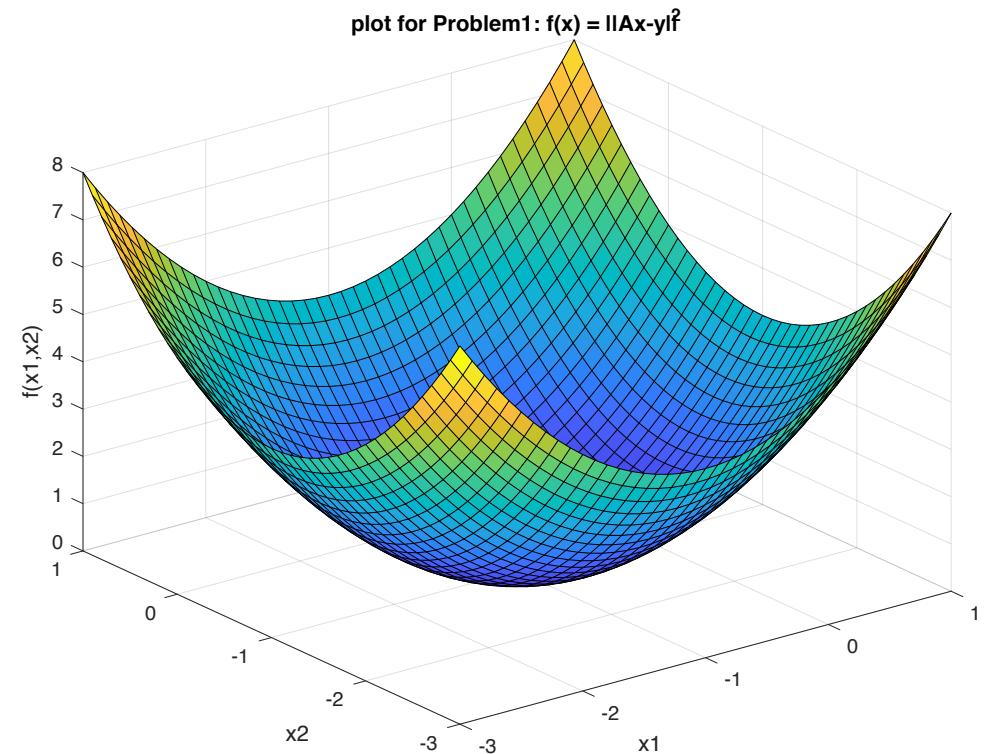
6. Hessian of f at $\mathbf{x} \in \mathbb{R}^n$

$$Hf(\mathbf{x}) = D^2f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix},$$

where we note that Hessian is a symmetric matrix, i.e., $Hf(\mathbf{x}) = [Hf(\mathbf{x})]^\top$.

Highlight:

- Hessian is defined here is for function $f: \mathbb{R}^n \rightarrow \mathbb{R}$;
- Hessian is a matrix with size $n \times n$;



Thanks