

# Machine Learning and Data Mining I (DS 4400)

## Homework 1 Solution

Instructor: Ehsan Elhamifar

Due Date: January 28, 2020, 1:35pm

1) Let  $\mathbf{a} \in \mathbb{R}^n$  be an  $n$ -dimensional vector and let  $\mathbf{U} \in \mathbb{R}^{n \times n}$  be an orthonormal matrix, i.e.,  $\mathbf{U}^\top \mathbf{U} = \mathbf{U} \mathbf{U}^\top = \mathbf{I}_n$ . Show the following:

[ 8 points . 4 pts for each part ]

1.  $\text{trace}(\mathbf{a}\mathbf{a}^\top) = \|\mathbf{a}\|_2^2$ .

2.  $\|\mathbf{U}\mathbf{a}\|_2^2 = \|\mathbf{a}\|_2^2$ .  $\rightarrow$  This means the magnitude of a vector does not change after transformation by orthonormal matrix.

Solution: 1.  $\text{trace}(\mathbf{a}\mathbf{a}^\top) = \text{trace}(\mathbf{a}^\top \mathbf{a})$  [ we know  $\text{trace}(\mathbf{AB}) = \text{trace}(\mathbf{BA})$  ]  
 $= \text{trace}(\|\mathbf{a}\|_2^2) = \|\mathbf{a}\|_2^2$ .  $\swarrow$  [ Since  $\mathbf{a}^\top \mathbf{a} = \|\mathbf{a}\|_2^2 \in \mathbb{R}$  ]

2.  $\|\mathbf{U}\mathbf{a}\|_2^2 = (\mathbf{U}\mathbf{a})^\top (\mathbf{U}\mathbf{a}) = \mathbf{a}^\top \underbrace{\mathbf{U}^\top \mathbf{U}}_{\mathbf{I}_n} \mathbf{a} = \mathbf{a}^\top \mathbf{I}_n \mathbf{a}$  [ Since  $\mathbf{U}$  is orthonormal:  $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_n$  ]  
 $= \mathbf{a}^\top \mathbf{a} = \|\mathbf{a}\|_2^2$ .  $\swarrow$

2) Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{B} \in \mathbb{R}^{n \times n}$  be arbitrary but invertible matrices and let  $\alpha$  be a scalar. Show the following:

1.  $(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$ . [ 12 points , 4 pts each part ]

2.  $(\mathbf{A}^\top)^{-1} = (\mathbf{A}^{-1})^\top$ .

3.  $\text{trace}(\alpha \mathbf{A}) = \alpha \text{trace}(\mathbf{A})$ .

Solution: 1. The definition of inverse of a matrix  $\mathbf{Z}$  is  $\mathbf{Z}^{-1}$  so that  $\mathbf{Z}^{-1} \mathbf{Z} = \mathbf{Z} \mathbf{Z}^{-1} = \mathbf{I}_n$ .

For  $(\mathbf{B}^{-1} \mathbf{A}^{-1})$  to be inverse of  $\mathbf{AB}$ , we investigate  $(\mathbf{B}^{-1} \mathbf{A}^{-1})(\mathbf{AB}) \stackrel{?}{=} \mathbf{I}_n$

$$(\mathbf{B}^{-1} \mathbf{A}^{-1})(\mathbf{AB}) = \mathbf{B}^{-1} \underbrace{(\mathbf{A}^{-1} \mathbf{A})}_{\mathbf{I}_n} \mathbf{B} = \mathbf{B}^{-1} \mathbf{B} = \mathbf{I}_n \quad \swarrow$$

2. We know  $\mathbf{A} \mathbf{A}^{-1} = \mathbf{I}_n \xrightarrow[\text{of both sides}]{\text{take transpose}} (\mathbf{A} \mathbf{A}^{-1})^\top = \mathbf{I}_n^\top \rightarrow (\mathbf{A}^{-1})^\top \mathbf{A}^\top = \mathbf{I}_n$

As a result  $(\mathbf{A}^{-1})^\top$  must be the inverse of  $\mathbf{A}^\top$ :  $(\mathbf{A}^\top)^{-1} = (\mathbf{A}^{-1})^\top$ .  $\swarrow$

3.  $\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \rightarrow \text{trace}(\mathbf{A}) = \sum_{i=1}^n a_{ii}$

$$\alpha \mathbf{A} = \begin{bmatrix} \alpha a_{11} & \alpha a_{12} & \dots & \alpha a_{1n} \\ \alpha a_{21} & \alpha a_{22} & \dots & \alpha a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha a_{n1} & \alpha a_{n2} & \dots & \alpha a_{nn} \end{bmatrix} \rightarrow \text{trace}(\alpha \mathbf{A}) = \sum_{i=1}^n \alpha a_{ii} = \alpha \text{trace}(\mathbf{A}) \quad \swarrow$$

3) For vectors  $x \in \mathbb{R}^n$ ,  $a \in \mathbb{R}^n$  and matrices  $X \in \mathbb{R}^{n \times n}$ ,  $A \in \mathbb{R}^{n \times n}$ , show the following.

1.  $\frac{\partial a^T A x}{\partial x} = A^T a.$

[ 10 points , 6 pts each part ]

2.  $\frac{\partial \text{trace}(A^T X)}{\partial X} = A.$

3.  $\frac{\partial \|Ax\|^2}{\partial x} = 2A^T Ax.$

Solution: 1.  $\underbrace{A^T A}_{\triangleq b} x = (\underbrace{A^T A}_{\triangleq b})^T x = b^T x$  (by defining  $b \triangleq A^T A$ )

$$\frac{\partial b^T x}{\partial x} = \frac{\partial \sum_{i=1}^n b_i x_i}{\partial x} = \begin{pmatrix} \partial(b_1 x_1 + \dots + b_n x_n) / \partial x_1 \\ \vdots \\ \partial(b_1 x_1 + \dots + b_n x_n) / \partial x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} = b$$

$$\rightarrow \frac{\partial (A^T A x)}{\partial x} = \frac{\partial ((A^T A)^T x)}{\partial x} = \frac{\partial (b^T x)}{\partial x} = b = A^T A. \checkmark$$

2.  $f(x) \triangleq \text{trace}(A^T x) = \text{tr} \left[ \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 & \dots & x_{1n} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nn} \end{pmatrix} \right] = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_{ij}$

$$\frac{\partial f}{\partial x} = \begin{pmatrix} \frac{\partial f}{\partial x_{11}} & \dots & \frac{\partial f}{\partial x_{1n}} \\ \vdots & & \vdots \\ \frac{\partial f}{\partial x_{n1}} & \dots & \frac{\partial f}{\partial x_{nn}} \end{pmatrix} = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} = A. \checkmark$$

3.  $f(x) \triangleq \|Ax\|_2^2 = (Ax)^T (Ax) = x^T A^T A x = x^T B x$  (here, for simplicity, we defined  $B \triangleq A^T A$ ).

$$= [x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \vdots & \vdots & & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = [x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} b_{11}x_1 + b_{12}x_2 + \dots + b_{1n}x_n \\ b_{21}x_1 + b_{22}x_2 + \dots + b_{2n}x_n \\ \vdots \\ b_{n1}x_1 + b_{n2}x_2 + \dots + b_{nn}x_n \end{bmatrix}$$

$$= x_1(b_{11}x_1 + \dots + b_{1n}x_n) + x_2(b_{21}x_1 + \dots + b_{2n}x_n) + \dots + x_n(b_{n1}x_1 + \dots + b_{nn}x_n)$$

$$= b_{11}x_1^2 + \dots + b_{1n}x_1x_n + b_{21}x_1x_2 + \dots + b_{2n}x_2x_n + \dots + b_{n1}x_nx_1 + \dots + b_{nn}x_n^2 = \sum_{i=1}^n \sum_{j=1}^n b_{ij}x_ix_j$$

$$\frac{\partial f}{\partial x} = \begin{pmatrix} \partial f / \partial x_1 \\ \vdots \\ \partial f / \partial x_n \end{pmatrix} = \begin{pmatrix} 2b_{11}x_1 + (b_{12}+b_{21})x_2 + \dots + (b_{1n}+b_{n1})x_n \\ \vdots \\ (b_{n1}+b_{1n})x_1 + (b_{n2}+b_{2n})x_2 + \dots + 2b_{nn}x_n \end{pmatrix} = (B+B^T)x \stackrel{\text{using } B \triangleq A^T A}{=} [A^T A + (A^T A)^T]x$$

$$= [A^T A + A^T A]x = 2A^T A x. \checkmark$$

4) Determine whether each of the following functions is convex or not.

1.  $f(x) = (x - a)^2$ , for any real number  $a$ .

[ 15 points ; part 1 : 4 pts , part 2 : 4 pts , part 3 : 7 pts ]

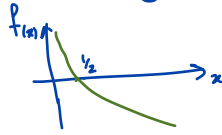
2.  $f(x) = -\log(2x)$ , with the domain  $x \in (0, +\infty)$ .

3.  $f(x) = e^{g(x)}$ , where  $g(x)$  is convex.

Solution: 1.  $f'(x) = 2(x-a) \rightarrow f''(x) = 2 \geq 0$ , Since  $f''(x) \geq 0 \forall x$ ,  $f(x)$  is convex.



2.  $f(x) = -\log 2x \rightarrow f'(x) = -\frac{2}{2x} = -\frac{1}{x} \rightarrow f''(x) = \frac{1}{x^2} \geq 0$ . Since  $f''(x) \geq 0 \forall x$ ,  $f(x)$  is convex.



3.  $f(x) = e^{g(x)} \rightarrow f'(x) = g'(x) e^{g(x)}$

$$\rightarrow f''(x) = g''(x) e^{g(x)} + (g'(x))^2 e^{g(x)} \quad \left( \begin{array}{c} \geq 0 \quad \geq 0 \quad \geq 0 \quad \geq 0 \end{array} \right) \quad \left( \begin{array}{c} \text{From convexity of } g, \text{ we know} \\ g''(x) \geq 0 \quad \forall x \end{array} \right)$$

$\rightarrow f''(x) \geq 0, \forall x \rightarrow f(x)$  is convex.

5) Consider the function  $f(x_1, x_2) = (x_1 - 2)^2 + (x_2 - 3)^2$ . Write down the expression of the gradient of  $f$ .

[31 points, part 2: 15 pts, rest each 4 pts]

1. Plot the function.
2. Write a python code that gets as input i) an initial point  $(x_1^0, x_2^0)$ , ii) the maximum number of iterations of the gradient descent, iii) learning rate  $\rho$ ; and runs the gradient descent algorithm starting from the initial point until convergence or until the maximum number of iterations is achieved. The output of the code must be the sequence of points (starting from initial point and ending with the last point) obtained by the gradient descent.
3. Use the learning rate  $\rho = 0.01$  and the initial point  $(x_1, x_2) = (1, 2)$ . Plot the sequence of obtained points. After how many iterations does GD converge?
4. Use the learning rate  $\rho = 0.5$  and the initial point  $(x_1, x_2) = (1, 2)$ . Plot the sequence of obtained points. After how many iterations does GD converge?
5. Use the learning rate  $\rho = 10$  and the initial point  $(x_1, x_2) = (1, 2)$ . Plot the sequence of obtained points. After how many iterations does GD converge?

$$f(x_1, x_2) = \underbrace{(x_1 - 2)^2}_{\geq 0} + \underbrace{(x_2 - 3)^2}_{\geq 0} \geq 0, \text{ when } x_1 = 2, x_2 = 3, f(2, 3) = 0 \rightarrow (x_1 = 2, x_2 = 3) \text{ is global minimizer.}$$

$$\frac{\partial f}{\partial x} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{pmatrix} = \begin{pmatrix} 2(x_1 - 2) \\ 2(x_2 - 3) \end{pmatrix}$$

$$\text{Running GD at iteration } t: \begin{pmatrix} x_1^{(t)} \\ x_2^{(t)} \end{pmatrix} = \begin{pmatrix} x_1^{(t-1)} \\ x_2^{(t-1)} \end{pmatrix} - \rho \begin{pmatrix} 2(x_1^{(t-1)} - 2) \\ 2(x_2^{(t-1)} - 3) \end{pmatrix}$$

6) Consider the training dataset

[ 16 points , part 1: 12 pts, part 2: 4 pts ] .

$\{(0.10, 0.15), (0.50, 0.40), (0.90, 0.85), (1.50, 1.62), (-0.20, -0.17), (-0.5, -0.42)\}$ ,

where in  $(\cdot, \cdot)$ , the first entry is the input variable,  $x$ , and the second entry is the output variable (response),  $y$ . Consider the regression model  $y = \theta x$ .

1. Write a python code that inputs the above data and outputs the optimal regression value  $\theta^*$ , using the closed-form solution.
2. Plot the data in 2D and plot the estimated line  $y = \theta^* x$ .

**Homework Submission Instructions:** Please submit the analytical part of the homework in class (put a paper copy of the HW on the instructor's desk before the class starts). Handwritten or printed documents generated by Latex/Word are both acceptable. Please submit all your plots and your Python code (.py file) via email, by the DEADLINE. To submit, please send an email to the instructor and cc the TAs.

- The title of your email must be “DS4400: HW01:Your-Last-Name”.
- Please attach a single zip file to your email that contains all python codes and plots and a readme file on how to run your files.
- Please name your zip file as “HW01:Your-Last-Name”.