

Machine Learning Exercise Sheet 1

Math Refresher

Group_369

Fan XUE – fan98.xue@tum.de

Xing ZHOU – xing.zhou@tum.de

Jianzhe LIU – jianzhe.liu@tum.de

October 27, 2021

Problem 1

According to the matrices multiplication rule and the dimension of function f 's return value, which is \mathbb{R} :

$$\mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{B} \in \mathbb{R}^{1 \times M}, \mathbf{c} \in \mathbb{R}^{N \times P}, \mathbf{D} \in \mathbb{R}^{Q \times 1}, \mathbf{E} \in \mathbb{R}^{N \times N}, \mathbf{F} \in \mathbb{R}^{1 \times 1}$$

Problem 2

$$f(\mathbf{x}) = \sum_{i=1}^N \sum_{j=1}^N x_i x_j M_{ij} = \sum_{i=1}^N x_i \left(\sum_{j=1}^N M_{ij} x_j \right) = \sum_{i=1}^N x_i (\mathbf{M}\mathbf{x})_i = \mathbf{x}^T \mathbf{M}\mathbf{x}$$

Problem 3

- a) If $M < N$ or $\text{rank}(\mathbf{A}) < N$, then the solution \mathbf{x} is not always unique. If $M > N$ or $\text{rank}(\mathbf{A}) < M$, there would be no solution \mathbf{x} for every $\mathbf{b} \in \mathbb{R}$.
If and only if $M = N$ and \mathbf{A} is full rank, i.e. \mathbf{A} is invertible, there is a unique solution \mathbf{x} for any choice of \mathbf{b} .

- b) No, because \mathbf{A} has an eigenvalue equal to 0, that means \mathbf{A} does not have full rank.

Problem 4

Since $\mathbf{B}\mathbf{A} = \mathbf{A}\mathbf{B} = \mathbf{I}$, \mathbf{A} is invertible. The determinant of \mathbf{A} is

$$\det \mathbf{A} = \prod_{i=1}^N \lambda_i \neq 0$$

, which indicates none of the eigenvalues of \mathbf{A} is 0.

Problem 5

1. $\lambda_i \geq 0 \Rightarrow \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$

\mathbf{A} is symmetric, there exists a orthonormal matrix \mathbf{Q} such that $\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T$. $\mathbf{\Lambda}$ is a diagonal matrix with the eigenvalues of \mathbf{A} in its diagonal.

For any $\mathbf{x} \in \mathbb{R}^N$, we have

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T \mathbf{x} = (\mathbf{Q}^T \mathbf{x})^T \mathbf{\Lambda} \mathbf{Q}^T \mathbf{x}$$

Notice that $\mathbf{Q}^T \mathbf{x} \in \mathbb{R}^N$, let $\mathbf{y} = \mathbf{Q}^T \mathbf{x}$, we have

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{y}^T \mathbf{\Lambda} \mathbf{y} = \begin{bmatrix} y_1 & \cdots & y_n \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n \lambda_i y_i^2$$

Because all y_i^2 and $\lambda_i \geq 0$, we have $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$.

2. $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0 \Rightarrow \lambda_i \geq 0$

Assume that \mathbf{x}_i and λ_i are one of the eigenvectors and the corresponding eigenvalue of \mathbf{A} . We have

$$\mathbf{A} \mathbf{x}_i = \lambda_i \mathbf{x}_i$$

then

$$\mathbf{x}_i^T \mathbf{A} \mathbf{x}_i = \mathbf{x}_i^T \lambda_i \mathbf{x}_i = \lambda_i \|\mathbf{x}_i\|^2 \geq 0$$

Because $\|\mathbf{x}_i\|^2 \geq 0$, we have $\lambda_i \geq 0$.

Problem 6

For any $\mathbf{A} \in \mathbb{R}^{M \times N}$, we have

$$\mathbf{x}^T \mathbf{B} \mathbf{x} = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} = (\mathbf{A} \mathbf{x})^T (\mathbf{A} \mathbf{x}) = \|\mathbf{A} \mathbf{x}\|^2 \geq 0$$

which means \mathbf{B} is positive semi-definite.

Problem 7

a) We observe the first and second derivative of $f(x)$, we have

$$\begin{aligned} f'(x) &= ax + b \\ f''(x) &= a \end{aligned}$$

(i) a unique solution means there exists a $x \in \mathbb{R}$, such that $f'(x) = 0$ and $f''(x) > 0$, which is

$$a > 0$$

- (ii) infinitely many solutions means there exist infinitely many $x \in \mathbb{R}$, such that $f'(x) = 0$, which is

$$a = b = 0$$

- (iii) no solution means $f'(x) = 0$ have no solution or $f'(x) = 0$ have solution but at that point $f''(x) < 0$, which is

$$a = 0, b \neq 0 \quad \text{or} \quad a < 0$$

- b) let $f'(x) = ax + b = 0$, we have

$$x^* = -\frac{b}{a}$$

Problem 8

- a) Consider the term including x_k and x_l factors to take the partial derivative, under the constrain of $\mathbf{A} \in \mathbb{S}^N$, i.e. $A_{ij} = A_{ji}$:

$$\begin{aligned} \frac{\partial g(\mathbf{x})}{\partial x_k \partial x_l} &= \frac{\partial}{\partial x_k \partial x_l} \left(\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N A_{ij} x_i x_j + \sum_{i=1}^N b_i x_i + c \right) \\ &= \frac{\partial}{\partial x_k} \left(\frac{1}{2} \sum_{i=1}^N A_{il} x_i + \frac{1}{2} \sum_{j=1}^N A_{lj} x_j + b_l \right) \\ &= \frac{\partial}{\partial x_k} \left(\sum_{i=1}^N A_{il} x_i + b_l \right) = A_{kl} \end{aligned}$$

This result indicates the Hessian of $g(\mathbf{x})$ is:

$$\nabla^2 g(\mathbf{x}) = \mathbf{A}$$

To have a unique solution of this optimization problem, \mathbf{A} should be positive definite.

- b) Since $\mathbf{A} \in \mathbb{S}^N$, \mathbf{A} can be represented as $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$, in which \mathbf{U} is the matrix of orthonormal eigenvectors of \mathbf{A} , $\mathbf{\Lambda}$ is the diagonal matrix of eigenvalues of \mathbf{A} .

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \mathbf{x} = \mathbf{y}^T \mathbf{\Lambda} \mathbf{y} = \sum_{i=1}^N \lambda_i y_i^2$$

where $\mathbf{y} = \mathbf{U}^T \mathbf{x}$ (and since \mathbf{U} is full rank, any vector $\mathbf{y} \in \mathbb{R}^N$ can be represented in this form).

If there exists a negative eigenvalue of \mathbf{A} , which we assume $\lambda_k < 0$, then as we take $y_k \rightarrow \infty$, $g(\mathbf{x})$ would be infinite negative, which means the function has no minimum. Therefore the matrix \mathbf{A} is PSD should be well-defined.

- c) Since the matrix \mathbf{A} is PD, we could minimize the objective function by setting the gradient to zero. To get the gradient of $g(\mathbf{x})$ we first consider the term of the partial derivative of x_k :

$$\begin{aligned}\frac{\partial g(\mathbf{x})}{\partial x_k} &= \frac{\partial}{\partial x_k} \left(\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N A_{ij} x_i x_j + \sum_{i=1}^N b_i x_i + c \right) \\ &= \frac{1}{2} \sum_{i=1}^N A_{ik} x_i + \frac{1}{2} \sum_{j=1}^N A_{kj} x_j + b_k \\ &= \sum_{i=1}^N A_{ik} x_i + b_k\end{aligned}$$

therefore the gradient of $g(\mathbf{x})$ is:

$$\nabla g(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$$

By setting $\nabla g(\mathbf{x}) = 0$, we can get

$$\mathbf{x}^* = -\mathbf{A}^{-1}\mathbf{b}.$$

Problem 9

The equation can be **disproved** by following example: let's say we have a coin, and after one toss we define 3 events:

Event A: the coin is back-side. $p(a) = 0.5$

Event B: the coin is front-side. $p(b) = 0.5$

Event C: the coin is front-side. $p(c) = 0.5$

(Here we use an extreme example, where event B and event C are equivalent)

In this case, Event A and Event B (or Event C) are mutually exclusive, which means:

$$p(a|b, c) = p(a|c) = 0$$

But obviously,

$$p(a|b) = 0 \neq p(a) = 0.5$$

Problem 10

The equation can be **disproved** by following example: let's say we have 2 dices, and after rolling we define 3 events:

Event A: the first dice shows 1 or 2. $p(a) = \frac{1}{3}$

Event B: the second dice shows 1 or 3. $p(b) = \frac{1}{3}$

Event C: the result of both dices are 2 and 3. $p(c) = \frac{1}{18}$

In this case, Event A and Event B are independent, which means:

$$p(a|b) = p(a) = \frac{1}{3}$$

But when we see the other side of the equation,
when Event B and Event C are determined, $p(a) = 1$.

when only Event C is determined, $p(a) = \frac{1}{2}$.

That is to say:

$$p(a|b, c) = 1 \neq p(a|c) = \frac{1}{2}$$

Problem 11

When the joint PDF $p(a, b, c)$ of three continuous random variables are given, we can solve the following problems according to the definition of the probability density function and conditional Probability:

1. $p(a)$ can be obtained by integrating b and c

$$p(a) = \iint p(a, b, c) db dc$$

2. Since a , b and c are independent to each other, so:

$$p(c|a, b) = \frac{p(a, b, c)}{p(a, b)} = \frac{p(a, b, c)}{\int p(a, b, c) dc}$$

3. According to the question 1 and question 2:

$$p(b|c) = \frac{p(b, c)}{p(c)} = \frac{\int p(a, b, c) da}{\iint p(a, b, c) da db}$$

Problem 12

The probability that the person has the disease can be calculated by following steps:
If he has the disease while being tested positive, then:

$$p(\text{disease}) * p(\text{positive with disease}) = 0.1\% * 95\%$$

If he has no disease while being tested positive, then:

$$p(\text{healthy}) * p(\text{positive when healthy}) = 99.9\% * 5\%$$

According to these 2 equations, the probability that he really has the disease is:

$$\frac{0.1\% * 95\%}{0.1\% * 95\% + 99.9\% * 5\%} = 1.87\%$$

(This probabilistic model is also known as a typical example in Bayes' theorem)

Problem 13

If $X \sim \mathcal{N}(\mu, \sigma^2)$, then we have following equations:

$$\mathbb{E}(x) = \mu$$

$$\text{var}(x) = \mathbb{E}(x^2) - \mathbb{E}(x)^2 = \sigma^2$$

Then for $f(x) = ax + bx^2 + c$, we have:

$$\mathbb{E}(f(x)) = \mathbb{E}(ax) + \mathbb{E}(bx^2) + \mathbb{E}(c) = a\mu + b(\mu^2 + \sigma^2) + c$$

Problem 14

a)

$$\mathbb{E}[g(\mathbf{x})] = \mathbb{E}[\mathbf{A}\mathbf{x}] = \mathbf{A}\mathbb{E}[\mathbf{x}] = \mathbf{A}\boldsymbol{\mu}$$

b)

$$\mathbb{E}[g(\mathbf{x})g(\mathbf{x})^T] = \mathbb{E}[\mathbf{A}\mathbf{x}(\mathbf{A}\mathbf{x})^T] = \mathbb{E}[\mathbf{A}\mathbf{x}\mathbf{x}^T\mathbf{A}^T] = \mathbf{A}\mathbb{E}[\mathbf{x}\mathbf{x}^T]\mathbf{A}^T$$

we have $(\mathbf{x}\mathbf{x}^T)_{ij} = x_i x_j$, so $(\mathbb{E}[\mathbf{x}\mathbf{x}^T])_{ij} = \mathbb{E}[x_i x_j]$.

if $i = j$, we have

$$\mathbb{E}[x_i^2] = \text{Var}[x_i] + (\mathbb{E}[x_i])^2 = \Sigma_{ii} + \mu_i^2$$

if $i \neq j$, x_i and x_j are independent, we have

$$\mathbb{E}[x_i x_j] = \mathbb{E}[x_i]\mathbb{E}[x_j] = \mu_i \mu_j$$

Combining the two cases, we have

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^T$$

So

$$\mathbb{E}[g(\mathbf{x})g(\mathbf{x})^T] = \mathbf{A}\mathbb{E}[\mathbf{x}\mathbf{x}^T]\mathbf{A}^T = \mathbf{A}(\boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^T)\mathbf{A}^T$$

c)

$$\mathbb{E}[g(\mathbf{x})^T g(\mathbf{x})] = \mathbb{E}[(\mathbf{A}\mathbf{x})^T \mathbf{A}\mathbf{x}] = \mathbb{E}[\mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x}]$$

Because $\mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x} \in \mathbb{R}$, we have

$$\begin{aligned} \mathbb{E}[\mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x}] &= \mathbb{E}[\text{Tr}(\mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x})] = \mathbb{E}[\text{Tr}(\mathbf{A}\mathbf{x}\mathbf{x}^T \mathbf{A}^T)] \\ &= \text{Tr}(\mathbb{E}[\mathbf{A}\mathbf{x}\mathbf{x}^T \mathbf{A}^T]) \\ &= \text{Tr}(\mathbf{A}(\boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^T)\mathbf{A}^T) \end{aligned}$$

d)

$$\begin{aligned}\text{Cov}[g(\boldsymbol{x})] &= \mathbb{E}[g(\boldsymbol{x})g(\boldsymbol{x})^T] - \mathbb{E}[g(\boldsymbol{x})](\mathbb{E}[g(\boldsymbol{x})])^T \\ &= \boldsymbol{A}(\boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^T)\boldsymbol{A}^T - \boldsymbol{A}\boldsymbol{\mu}(\boldsymbol{A}\boldsymbol{\mu})^T \\ &= \boldsymbol{A}\boldsymbol{\Sigma}\boldsymbol{A}^T\end{aligned}$$