
Exercise 01 of Machine Learning [IN 2064]

Name: Yiman Li
Matr-Nr: 03724352
cooperate with Kejia Chen(03729686)

Problem 1

In order for the matrix product to exist, the number of columns in the former matrix must equal the number of rows in the latter matrix, here the function f can be described as:

$$f(\mathbf{x}, \mathbf{y}, \mathbf{Z}) = (\mathbf{x}^T)^{1 \times M} \mathbf{A} \mathbf{y}^{N \times 1} + \mathbf{B} \mathbf{x}^{M \times 1} - (\mathbf{y}^T)^{1 \times N} \mathbf{C} \mathbf{Z}^{P \times Q} \mathbf{D} - (\mathbf{y}^T)^{1 \times N} \mathbf{E}^T \mathbf{y}^{N \times 1} + \mathbf{F}$$

Then using the aforementioned rules, we can get

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

Problem 2

Now that the function $f(\mathbf{x}) = \sum_{i=1}^N \sum_{j=1}^N x_i x_j M_{ij} \in \mathbb{R}^{N \times N}$, so we can simply rewrite the function as $f(\mathbf{x}) = \mathbf{x}^T \mathbf{M} \mathbf{x}$ just on the condition that the number of columns in the former matrix must equal the number of rows in the latter matrix in matrix product.

Problem 3

- a) When $R(\mathbf{A}) = R(\mathbf{A}, \mathbf{b}) = M$, then we can safely say that the solution is unique for each choice of \mathbf{b} . Here R represents "Rank of the Matrix".
- b) Notice that the matrix has a zero eigenvalue, which means the $R(\mathbf{A} < 5)$, so we can not safely draw a conclusion that this matrix is diagonalizable, so the Equation (1) can not always find a unique solution \mathbf{x} for any choice of \mathbf{b} .

Problem 4

Here we can use the property that the determinant of \mathbf{A} is equal to the product of its eigenvalues $|\mathbf{A}| = \prod_{i=1}^n \lambda_i$. Now that the matrix \mathbf{A} is invertible, which means that $|\mathbf{A}| \neq 0$, so we can say that one of the eigenvalues of matrix \mathbf{A} is zero since \mathbf{A} is invertible.

Problem 5

According to reference [4], since a symmetric matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ has orthogonal eigenvectors and is thus orthogonal, we can therefore represent \mathbf{A} as $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$, then we can show that

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

where $\mathbf{y} = \mathbf{U}^T \mathbf{x}$. Because y_i^2 is always positive, the sign of this expression depends entirely on the λ_i 's. So if all eigenvalues $\lambda_i \geq 0$, then it is positive semidefinite.

Conversely, if a matrix is positive semidefinite, then we assume that this matrix \mathbf{A} has a negative eigenvalue $\lambda < 0$ and its corresponding eigenvector \mathbf{x} , so we can get

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$$

then we multiply the two sides of the equation by the transpose of the eigenvector \mathbf{x}^T , which means

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \lambda \mathbf{x} = \lambda \mathbf{x}^T \mathbf{x} = \lambda \mathbf{x}^2 < 0$$

this result is contrary to the condition that the matrix is positive semidefinite, so the assumption that the matrix has a negative eigenvalue doesn't make sense. So in the end we can draw a conclusion that a positive semidefinite matrix has no negative eigenvalues.

Problem 6

Choose an arbitrary vector $\mathbf{x} \in \mathbb{R}^N$ for the matrix \mathbf{B} , the scalar value can be written as

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

So the matrix \mathbf{B} is positive semi-definite for any choice of \mathbf{A} .

Problem 7

a)

1. Using the second partial derivatives of function f , that is

$$d_x^2 f(x) = d_x(ax + b) = a$$

So when $a > 0$, then $f(x)$ becomes a strictly convex quadratic function, which has at most one global minimum.

2. When $a = b = 0$, then $f(x)$ is a constant function, so each point of x is the solution of the minimum value of $f(x)$.
3. When $a < 0$, then $f(x)$ becomes a strictly concave function; or when $a = 0, b \neq 0$, whose minimum value lies on infinity. In both situation, there is no solution for the optimization.

b) Now we try to equal the first partial derivatives of function f to zero, that is

$$d_x f(x) = ax + b = 0$$

So the point $x = -\frac{b}{a}$ minimizes the objective function.

Problem 8

a) Now that the matrix is symmetric and positive semidefinite, so the Hessian $\nabla_x^2 g(x)$ of the objective function is shown as below:

$$\nabla_x^2 g(x) = \nabla_x^2 \left[\left(\frac{1}{2} \right) \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b} \mathbf{x}^T + c \right] = \mathbf{A}$$

Similar to **Problem 7**, when $\mathbf{A} > 0$, which means that the matrix is positive definite, then the optimization problem becomes a strictly convex function (See definition in reference [5])

b) If \mathbf{A} is positive definite, then we can find a unique solution for the optimization at the point where the gradient of the function is zero; if \mathbf{A} is positive semidefinite, then we might find the solution at infinity. However, when \mathbf{A} has a negative eigenvalue, which means that the matrix is indefinite, so we may not find the solution for the optimization.

c) Now that \mathbf{A} is positive definite (PD), so we just try to equal the gradient of the function to zero to figure out the result, that is:

$$\nabla_x g(x) = \nabla_x \left[\left(\frac{1}{2} \right) \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b} \mathbf{x}^T + c \right] = \mathbf{A} \mathbf{x} + \mathbf{b} = 0$$

So we can get to the point that when the matrix \mathbf{A} is positive definite (PD), then this optimization problem $g(x)$ will have a unique solution at the point $\mathbf{x} = -\mathbf{A}^{-1} \mathbf{b}$.

Problem 9

According to

$$p(A|B, C) = \frac{p(A, B, C)}{p(B, C)} = \frac{\frac{p(A, B, C)}{p(C)}}{\frac{p(B, C)}{p(C)}} = \frac{p(A, B|C)}{p(B|C)} = p(A|C)$$

that is

$$p(A, B|C) = p(A|C)p(B|C)$$

so we can say that two events A and B are conditionally independent given an event C with $P(C) > 0$. However, according to reference [3], conditional independence cannot lead to the concept of independence. Here is an example. Assuming that a box contains two coins: a regular coin and one irregular coin with $(P(H) = 1)$. Now choose a coin at random and toss it twice. Define the following events.

- A = First coin toss results in an H.
- B = Second coin toss results in an H.
- C = The regular coin has been chosen.

From this example we can get that $p(A|B, C) = p(A|C) = \frac{1}{2}$, then we can calculate that

$$p(A, B) = p(A, B|C) + p(A, B|\bar{C}) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} + 1 \cdot 1 \cdot \frac{1}{2} = \frac{5}{8}$$

whereas

$$p(A) = p(A|C) + p(A|\bar{C}) = \frac{1}{2} \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = \frac{3}{4}$$

$$p(B) = p(B|C) + p(B|\bar{C}) = \frac{1}{2} \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = \frac{3}{4}$$

On this condition $p(A, B) \neq p(A)p(B)$, so this statement is wrong.

Problem 10

When $p(A|B, C) = p(A|C)$, we can get the information that A and B are conditionally independent when given C. But generally speaking, conditional independence neither implies (nor is it implied by) independence, which is the information in $p(A|B) = p(A)$. Consider rolling a die and let

$$A = \{1, 2\}, B = \{2, 4, 5\}, C = \{1, 4\}$$

so we can get

$$p(A) = \frac{1}{3}, p(B) = \frac{1}{2}, p(A, B) = \frac{1}{6} = p(A)p(B)$$

which means A and B are independent. But we can also figure that

$$p(A|B, C) = \frac{p(\{1, 2\})}{p(\{4\})} = 0, p(A|C) = \frac{p(\{1, 2\})}{p(\{1, 4\})} = \frac{1}{2}$$

so the statement is false.

Problem 11

This problem is based on the concept of probability density function, more details can be found in reference [1]. According to the concept, we can directly write down the corresponding formular as below:

(1)

$$p(a) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p(a, b, c) db dc$$

(2)

$$p(c|a, b) = \frac{p(a, b, c)}{p(a, b)} = \frac{p(a, b, c)}{\int_{-\infty}^{\infty} p(a, b, c) dc}$$

(3)

$$p(b|c) = \frac{p(b, c)}{p(c)} = \frac{\int_{-\infty}^{\infty} p(a, b, c) da}{\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p(a, b, c) da db}$$

Problem 12

The probability that a person has a positive test result is

$$\frac{1}{1000} \times 0.95 + \frac{999}{1000} \times 0.05 = \frac{5090}{100000}$$

So when man obtains a positive result, his/her probability of having the disease is

$$\frac{\frac{1}{1000} \times 0.95}{\frac{1}{1000} \times 0.95 + \frac{999}{1000} \times 0.05} = \frac{19}{1018} = 0.0187$$

Problem 13

Since the mean value of a Gaussian distribution is μ , using the property that $Var[X] = E[X^2] + E[X]^2$, we can easily figure out that $E[f(x)] = a\mu + b(\mu^2 + \sigma^2) + c$

Problem 14

Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$, according to reference 2, we have the following conclusion:

$$E[\mathbf{Ax}] = \mathbf{AE}[\mathbf{x}] \quad (1)$$

$$E(\mathbf{xx}^T) = \Sigma + \mathbf{mm}^T \quad (2)$$

$$Var[\mathbf{Ax}] = \mathbf{A}Var[\mathbf{x}]\mathbf{A}^T \quad (3)$$

$$E[\mathbf{x}^T \mathbf{Ax}] = Tr(\mathbf{A} \Sigma) + \mathbf{m}^T \mathbf{Am} \quad (4)$$

- Using the Equation 1, we can get that

$$E[g(\mathbf{x})] = E[\mathbf{Ax}] = \mathbf{AE}[\mathbf{x}] = \mathbf{A}\mu$$

- Since $g(\mathbf{x}) = \mathbf{Ax} \sim \mathcal{N}(\mathbf{A}\mu, \mathbf{A} \Sigma \mathbf{A}^T)$, so according to the Equation 2, we can get

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

- Since $g(\mathbf{x}) = \mathbf{Ax} \sim \mathcal{N}(\mathbf{A}\mu, \mathbf{A} \Sigma \mathbf{A}^T)$, the changing the Equation 4 by replacing \mathbf{A} by an Identity matrix \mathbf{I} , we can get the result as below:

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

- Here we use the Equation 3, we can simply get the result that:

$$Cov[g(\mathbf{x})] = Cov[\mathbf{Ax}, \mathbf{Ax}] = \mathbf{ACov}[\mathbf{x}, \mathbf{x}] \mathbf{A}^T = \mathbf{A}Var(\mathbf{x}) \mathbf{A}^T = \mathbf{A} \Sigma \mathbf{A}^T$$

References

- [1] Arian Maleki and Tom Do. Review of probability theory. <http://cs229.stanford.edu/summer2019/cs229-prob.pdf>.
- [2] K. B. Petersen and M. S. Pedersen. The matrix cookbook, nov 2012.
- [3] Hossein Pishro-Nik. Conditional independence. https://www.probabilitycourse.com/chapter1/1_4_4_conditional_independence.php.
- [4] Gilbert Strang. *Introduction to Linear Algebra*. Wellesley-Cambridge Press, Wellesley, MA, fourth edition, 2009.
- [5] Wikipedia. Convex function. https://en.wikipedia.org/wiki/Convex_function.