

Machine Learning Exercise Sheet 1

Math Refresher

The machine learning lecture relies heavily on your knowledge of undergraduate mathematics, especially linear algebra and probability theory. You should think of this exercise sheet as a test to see if you meet the prerequisites for taking this course. If you struggle with a large fraction of the exercises you should reconsider taking this lecture at this point and instead first prepare by taking a course that reinforces your mathematical foundations (e.g. "Basic Mathematical Tools for Imaging and Visualization" (IN2124)).

Homework

Reading

We strongly recommend that you review the following documents to refresh your knowledge. You should already be familiar with most of their content from your previous studies.

- Linear algebra <http://cs229.stanford.edu/section/cs229-linalg.pdf> (except sections 4.4, 4.5, 4.6), and http://ee263.stanford.edu/notes/matrix_crimes.pdf (common linear algebra mistakes)
- Probability theory <http://cs229.stanford.edu/summer2020/cs229-prob.pdf>

Linear Algebra

Notation. We use the following notation in this lecture:

- Scalars are denoted with lowercase letters, e.g. a , x , μ .
- Vectors are denoted with bold lowercase letters, e.g. \mathbf{a} , \mathbf{x} , $\boldsymbol{\mu}$.
- Matrices are denoted with bold uppercase letters, e.g. \mathbf{A} , \mathbf{X} , $\boldsymbol{\Sigma}$.
- \mathbb{R}^N denotes N -dimensional Euclidean space, i.e. the set of N -dimensional vectors with real-valued entries. For example, $\mathbf{x} = (2, \sqrt{2}, 6.5, -7)^T$ is an element of \mathbb{R}^4 , which we denote as $\mathbf{x} \in \mathbb{R}^4$.
- $\mathbb{R}^{M \times N}$ is the set of matrices with M rows and N columns. For example, the matrix $\mathbf{A} = \begin{pmatrix} 2 & 3 & 1 \\ 1 & 4 & 5 \end{pmatrix}$ is an element of $\mathbb{R}^{2 \times 3}$, which we denote as $\mathbf{A} \in \mathbb{R}^{2 \times 3}$.
- A function $f : \mathcal{X} \rightarrow \mathcal{Y}$ maps elements of the set \mathcal{X} into the set \mathcal{Y} . An example would be a function $f : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ defined as $f(x, y) = 2x^2 + xy - 4$.

Problem 1: Let $\mathbf{x} \in \mathbb{R}^M$, $\mathbf{y} \in \mathbb{R}^N$ and $\mathbf{Z} \in \mathbb{R}^{P \times Q}$. The function $f : \mathbb{R}^M \times \mathbb{R}^N \times \mathbb{R}^{P \times Q} \rightarrow \mathbb{R}$ is defined as

$$f(\mathbf{x}, \mathbf{y}, \mathbf{Z}) = \mathbf{x}^T \mathbf{A} \mathbf{y} + \mathbf{B} \mathbf{x} - \mathbf{y}^T \mathbf{C} \mathbf{Z} \mathbf{D} - \mathbf{y}^T \mathbf{E}^T \mathbf{y} + \mathbf{F}.$$

What should be the dimensions (shapes) of the matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{E}, \mathbf{F}$ for the expression above to be a valid mathematical expression?

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

Problem 2: Let $\mathbf{x} \in \mathbb{R}^N, \mathbf{M} \in \mathbb{R}^{N \times N}$. Express the function $f(\mathbf{x}) = \sum_{i=1}^N \sum_{j=1}^N x_i x_j M_{ij}$ using **only** matrix-vector multiplications.

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{M} \mathbf{x}$$

Problem 3: Let $\mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{x} \in \mathbb{R}^N$ and $\mathbf{b} \in \mathbb{R}^M$. We are interested in solving the following system of linear equations for \mathbf{x}

$$\mathbf{A} \mathbf{x} = \mathbf{b} \quad (1)$$

- a) Under what conditions does the system of linear equations have a **unique** solution \mathbf{x} for **any** choice of \mathbf{b} ?

$M < N$ or $\text{rank}(\mathbf{A}) < N$? No, the solution \mathbf{x} would not always be unique.
 $M > N$ or $\text{rank}(\mathbf{A}) < M$? No, a solution \mathbf{x} would not exist for every $\mathbf{b} \in \mathbb{R}^M$.
Hence, $M = N$ and \mathbf{A} has full rank and is therefore invertible.

- b) Assume that $M = N = 5$ and that \mathbf{A} has the following eigenvalues: $\{-5, 0, 1, 1, 3\}$. Does Equation 1 have a unique solution \mathbf{x} for any choice of \mathbf{b} ? Justify your answer.

No, because \mathbf{A} has an eigenvalue 0 and therefore does not have full rank.

Problem 4: Let $\mathbf{A} \in \mathbb{R}^{N \times N}$. Assume that there exists a matrix $\mathbf{B} \in \mathbb{R}^{N \times N}$ such that $\mathbf{B} \mathbf{A} = \mathbf{A} \mathbf{B} = \mathbf{I}$. What can you say about the eigenvalues of \mathbf{A} ? Justify your answer.

By definition, \mathbf{B} is the inverse of \mathbf{A} . Therefore, \mathbf{A} is invertible, i.e. the determinant of \mathbf{A} is not equal to zero, i.e. none of the eigenvalues of \mathbf{A} are equal to zero.

Problem 5: A symmetric matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ is positive semi-definite (PSD) if and only if for any $\mathbf{x} \in \mathbb{R}^N$ it holds that $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$. Prove that a symmetric matrix \mathbf{A} is PSD **if and only if** it has no negative eigenvalues.

1. No negative eigenvalues \Rightarrow PSD:

Since \mathbf{A} is symmetric we can choose orthonormal eigenvectors and then express any vector in \mathbb{R}^N as $\mathbf{x} = \sum_{i=1}^N w_i \mathbf{v}_i$, with the eigenvectors \mathbf{v}_i and some coefficients w_i . Hence,

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_{i=1}^N \sum_{j=1}^N w_i w_j \mathbf{v}_i^T \mathbf{A} \mathbf{v}_j = \sum_{i=1}^N \sum_{j=1}^N w_i w_j \mathbf{v}_i^T \lambda_j \mathbf{v}_j = \sum_{i=1}^N \sum_{j=1}^N w_i w_j \lambda_j \delta_{ij} = \sum_{i=1}^N w_i^2 \lambda_i \geq 0,$$

since the eigenvalues $\lambda_i \geq 0$. $\delta_{ij} = \begin{cases} 0 & \text{if } i \neq j, \\ 1 & \text{if } i = j \end{cases}$ denotes the Kronecker delta.

2. PSD \Rightarrow no negative eigenvalues:

Let \mathbf{v} be an eigenvector of \mathbf{A} with eigenvalue λ . By the PSD property of \mathbf{A} we have

$$0 \leq \mathbf{v}^T \mathbf{A} \mathbf{v} = \mathbf{v}^T \lambda \mathbf{v} = \lambda \|\mathbf{v}\|_2^2.$$

Because the Euclidean norm $\|\mathbf{v}\|_2^2$ is non-negative we have $\lambda \geq 0$. □

Problem 6: Let $\mathbf{A} \in \mathbb{R}^{M \times N}$. Prove that the matrix $\mathbf{B} = \mathbf{A}^T \mathbf{A}$ is positive semi-definite for any choice of \mathbf{A} .

$$\mathbf{x}^T \mathbf{B} \mathbf{x} = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} = (\mathbf{A} \mathbf{x})^T (\mathbf{A} \mathbf{x}) = \|\mathbf{A} \mathbf{x}\|_2^2 \geq 0.$$

Since the norm is always non-negative.

Calculus

Problem 7: Consider the following function $f: \mathbb{R} \rightarrow \mathbb{R}$

$$f(x) = \frac{1}{2}ax^2 + bx + c$$

We are interested in solving the following optimization problem

$$\min_{x \in \mathbb{R}} f(x)$$

- a) Under what conditions does this optimization problem have (i) a unique solution, (ii) infinitely many solutions or (iii) no solution? Justify your answer.

We obtain a solution by setting the derivative to zero, i.e. $f'(x) = ax + b = 0$, and checking if the second derivative is positive, i.e. $f''(x) = a > 0$. Hence, we obtain

- (i) a single solution if $a > 0$,
- (ii) infinitely many solutions if $a = b = 0$, and
- (iii) no solution if $a = 0$, $b \neq 0$ or $a < 0$.

- b) Assume that the optimization problem has a unique solution. Write down the closed-form expression for x^* that minimizes the objective function, i.e. find $x^* = \arg \min_{x \in \mathbb{R}} f(x)$.

This is an important question that we will encounter in multiple disguises throughout the lecture. Since we know that $a > 0$ we can solve it by simply setting the derivative to zero:

$$f'(x) = ax + b = 0 \quad \Leftrightarrow \quad x = -\frac{b}{a}$$

Problem 8: Consider the following function $g : \mathbb{R}^N \rightarrow \mathbb{R}$

$$g(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c$$

where $\mathbf{A} \in \mathbb{R}^{N \times N}$ is a symmetric, PSD matrix, $\mathbf{b} \in \mathbb{R}^N$ and $c \in \mathbb{R}$.

We are interested in solving the following optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^N} g(\mathbf{x})$$

- a) Compute the Hessian $\nabla^2 g(\mathbf{x})$ of the objective function. Under what conditions does this optimization problem have a unique solution?

The Hessian is defined as the matrix of partial derivatives (see Section 4.2 of the Stanford Linear Algebra Review and Reference).

$$\begin{aligned} \partial_{x_l} \partial_{x_k} g(\mathbf{x}) &= \partial_{x_l} \partial_{x_k} \left(\frac{1}{2} \sum_i \sum_j x_i A_{ij} x_j + \sum_i b_i x_i + c \right) = \partial_{x_l} \left(\frac{1}{2} \sum_i A_{ik} x_i + \frac{1}{2} \sum_j A_{kj} x_j + b_k \right) \\ &= \partial_{x_l} \left(\sum_j A_{kj} x_j + b_k \right) = A_{kl} \\ \nabla^2 g(\mathbf{x}) &= \mathbf{A} \end{aligned}$$

A differentiable function has a unique minimum if the Hessian $\nabla^2 g(\mathbf{x})$ is positive definite for all \mathbf{x} . For our case, this means that $g(\mathbf{x})$ has a unique minimum if \mathbf{A} is positive definite.

Note that if we extend the definition of PSD to non-symmetric matrices we would obtain $\nabla^2 g(\mathbf{x}) = \frac{1}{2}(\mathbf{A} + \mathbf{A}^T)$, which is the symmetric part of the matrix \mathbf{A} and PSD if \mathbf{A} is PSD. The rest of the exercise would work exactly the same way.

- b) Why is it necessary for a matrix \mathbf{A} to be PSD for the optimization problem to be well-defined?
Hint: What happens if \mathbf{A} has a negative eigenvalue?

The problem is ill-defined if it does not have a minimum. We assume that the matrix \mathbf{A} is not PSD and use the logical equivalence $a \Rightarrow b \Leftrightarrow \neg b \Rightarrow \neg a$.

Not being PSD means that the matrix \mathbf{A} has a negative eigenvalue. The eigenvector associated with that eigenvalue gives a direction in which the function always decreases, i.e.

$$a \mathbf{v}^T \mathbf{A} \mathbf{v} + \mathbf{b}^T \mathbf{v} + c = a^2 \mathbf{v}^T \lambda \mathbf{v} + \mathbf{b}^T \mathbf{v} + c = a^2 \lambda \|\mathbf{v}\|_2^2 + \mathbf{b}^T \mathbf{v} + c,$$

with the eigenvalue $\lambda < 0$, the corresponding eigenvector \mathbf{v} , and $a \in \mathbb{R}$. As we take $a \rightarrow \infty$ the first term will dominate at some point $a > a_0$ since $\mathbf{b}^T \mathbf{v} + c \in \mathcal{O}(a)$, with $\mathcal{O}(g(x)) = \{f \mid \forall c > 0, \exists x_0 > 0, \forall x > x_0 : |f(x)| \leq c|g(x)|\}$. Because $\lambda < 0$ this function will thus keep decreasing and therefore does not have a minimum. In summary, if \mathbf{A} and therefore the Hessian is not PSD the problem does not have a solution, i.e. g does not have a minimum. Hence the problem is not well-defined.

Note that convexity is not related to the problem being well-defined. A non-convex function can still have a global minimum.

- c) Assume that the matrix \mathbf{A} is positive definite (PD). Write down the closed-form expression for \mathbf{x}^* that minimizes the objective function, i.e. find $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^N} g(\mathbf{x})$.

We solve this by setting the gradient to zero (as in the previous exercise). For this, we first calculate the gradient:

$$\begin{aligned}\partial_{x_k} g(\mathbf{x}) &= \partial_{x_k} \left(\frac{1}{2} \sum_i \sum_j x_i A_{ij} x_j + \sum_i b_i x_i + c \right) \\ &= \frac{1}{2} \sum_i A_{ik} x_i + \frac{1}{2} \sum_j A_{kj} x_j + b_k = \sum_j A_{kj} x_j + b_k, \\ \nabla g(\mathbf{x}) &= \mathbf{A}\mathbf{x} + \mathbf{b}\end{aligned}$$

Setting $\nabla g(\mathbf{x}) = 0$, we obtain

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

Since \mathbf{A} is symmetric and PD it has no zero eigenvalues (see Exercise 5) and is therefore invertible.

Probability Theory

Notation. We use the following notation in our lecture

- For conciseness and to avoid clutter, we use $p(x)$ to denote multiple things
 1. If X is a discrete random variable, $p(x)$ denotes the probability mass function (PMF) of X at point x (usually denoted as $p_X(x)$ or $p(X = x)$ in the statistics literature).
 2. If X is a continuous random variable, $p(x)$ denotes the probability density function (PDF) of X at point x (usually denoted as $f_X(x)$ in the statistics literature).
 3. If $A \in \Omega$ is an event, $p(A)$ denotes the probability of this event (usually denoted as $\Pr(\{A\})$ or $\mathbb{P}(\{A\})$ in the statistics literature)

You will mostly encounter (1) and (2) throughout the lecture. Usually, the meaning is clear from the context.

- Given the distribution $p(x)$, we may be interested in computing the expected value $\mathbb{E}_{p(x)}[f(x)]$ or, equivalently, $\mathbb{E}_X[f(x)]$. Usually, it is clear with respect to which distribution we are computing the expectation, so we omit the subscript and simply write $\mathbb{E}[f(x)]$.
- $x \sim p$ means that x is distributed (sampled) according to the distribution p . For example, $x \sim \mathcal{N}(\mu, \sigma^2)$ (or equivalently $p(x) = \mathcal{N}(x|\mu, \sigma^2)$) means that x is distributed according to the normal distribution with mean μ and variance σ^2 .

Problem 9: Prove or disprove the following statement

$$p(a|b, c) = p(a|c) \Rightarrow p(a|b) = p(a)$$

Disprove by counterexample: We have a coin and do not know if it is fair, i.e. $p(A = T) = 0.5$, or unfair, i.e. $p(A = T) = 1$. C denotes the event whether the coin is fair $C = F$ or unfair $C = U$. Both of these events have equal probability $p(C = F) = p(C = U) = 0.5$. We perform two coin tosses

A and B . When we know which coin we have, coin tosses A and B are of course independent, i.e. $p(a|b, c) = p(a|c)$. However, if we do not observe C , then

$$\begin{aligned}
 p(A = T) &= p(A = T|C = F)p(C = F) + p(A = T|C = U)p(C = U) = \frac{3}{4}, \\
 p(A = T|B = T) &= \frac{p(A = T, B = T)}{p(B = T)} \\
 &= \frac{p(A = T, B = T|C = F)p(C = F) + p(A = T, B = T|C = U)p(C = U)}{p(B = T)} \\
 &= \frac{p(A = T|B = T, C = F)p(B = T|C = F)p(C = F) + p(A = T|B = T, C = U)p(B = T|C = U)p(C = U)}{p(B = T)} \\
 &= \frac{1/2 \cdot 1/2 \cdot 1/2 + 1 \cdot 1 \cdot 1/2}{3/4} = \frac{5/8}{3/4} = \frac{5}{6}.
 \end{aligned}$$

Therefore, $p(A = T) \neq p(A = T|B = T)$ and A and B are not independent. \square

Problem 10: Prove or disprove the following statement

$$p(a|b) = p(a) \Rightarrow p(a|b, c) = p(a|c)$$

Disprove by counterexample: Let the random variables A and B denote two independent dice rolls. $C = A + B$ denotes the sum of these dice rolls. Clearly, A and B are independent and therefore $p(a|b) = p(a)$. However, when we observe their sum the two become dependent. E.g. if we observe $C = 3$, then $p(A = 1) = 1/2$ and $p(A = 2) = 1/2$. However, if we observe $B = 2$, then $p(A = 1|B = 2) = 1$. \square

This fact is quite interesting. It means that two independent random variables can *become dependent* when observing a different variable, which is known as the *explaining away effect* (or selection bias, Berkson's paradox) in the literature. For example, if students are admitted to a university either by having good school grades or by showing excellent athletic performance, then these two attributes will be negatively correlated in the student population, even though they are independent in general.

Problem 11: You are given the joint PDF $p(a, b, c)$ of three continuous random variables. Show how the following expressions can be obtained using the rules of probability

1. $p(a)$
2. $p(c|a, b)$
3. $p(b|c)$

$$\begin{aligned}
 p(a) &= \int \int p(a, b, c) db dc \\
 p(c | a, b) &= \frac{p(a, b, c)}{p(a, b)} = \frac{p(a, b, c)}{\int p(a, b, c) dc}
 \end{aligned}$$

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

Problem 12: Researchers have developed a test which determines whether a person has a rare disease. The test is fairly reliable: if a person is sick, the test will be positive with 95% probability, if a person is healthy, the test will be negative with 95% probability. It is known that $\frac{1}{1000}$ of the population have this rare disease. A person (chosen uniformly at random from the population) takes the test and obtains a positive result. What is the probability that the person has the disease?

We can use Bayes' theorem to solve this. We denote the event of having the disease by D and a positive test by T .

$$p(D|T) = \frac{p(T|D)p(D)}{p(T)} = \frac{p(T|D)p(D)}{p(T|D)p(D) + p(T|\neg D)p(\neg D)} = \frac{0.95 \cdot 0.001}{0.95 \cdot 0.001 + 0.05 \cdot 0.999} = 0.0187$$

Problem 13: Let $X \sim \mathcal{N}(\mu, \sigma^2)$, and $f(x) = ax + bx^2 + c$. What is $\mathbb{E}[f(x)]$?

Using the linearity of expectation and $\sigma^2 = \mathbb{E}[x^2] - \mathbb{E}[x]^2$ we obtain

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

Problem 14: Let $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$, and $g(\mathbf{x}) = \mathbf{A}\mathbf{x}$ (where $\mathbf{A} \in \mathbb{R}^{N \times N}$). What are the values of the following expressions:

- $\mathbb{E}[g(\mathbf{x})]$,
- $\mathbb{E}[g(\mathbf{x})g(\mathbf{x})^T]$,
- $\mathbb{E}[g(\mathbf{x})^T g(\mathbf{x})]$,
- the covariance matrix $\text{Cov}[g(\mathbf{x})]$.

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

Note that we have used the fact that taking the trace of a scalar does not have any effect, the trace's cyclic property and that it is defined as a sum and can therefore be pulled out of the expectation when proving $\mathbb{E}[g(\mathbf{x})^T g(\mathbf{x})]$.