

# 70015 Mathematics for Machine Learning: Exercises

Mark van der Wilk, Yingzhen Li\*  
`{m.vdwilk,yingzhen.li}@imperial.ac.uk`

October 21, 2022

## Contents

<b>1</b>	<b>Notation</b>	<b>1</b>
1.1	Sets	1
1.2	Probabilities	2
<b>2</b>	<b>Formula Sheet</b>	<b>3</b>
<b>3</b>	<b>Warm-up Exercises</b>	<b>3</b>
3.1	Probability Theory	3
3.2	Linear Algebra	5
<b>4</b>	<b>Lecture 1: Probability, Vectors, Differentiation</b>	<b>6</b>
<b>5</b>	<b>Lecture 2: Vector Differentiation</b>	<b>7</b>
<b>6</b>	<b>Lecture 3: Automatic Differentiation</b>	<b>8</b>
<b>7</b>	<b>Lecture 4: Probabilistic Modelling Principles</b>	<b>8</b>
<b>8</b>	<b>Lecture 5: Gradient Descent Convergence</b>	<b>9</b>
<b>9</b>	<b>Lecture <math>N</math>: Multivariate Probability</b>	<b>9</b>
<b>10</b>	<b>Warm-up Exercises Answers</b>	<b>10</b>
10.1	Warm-up Exercises	10
10.2	Linear Algebra	11
<b>11</b>	<b>Answers Lecture 2: Vector Differentiation</b>	<b>12</b>
<b>12</b>	<b>Answers Lecture 3: Automatic Differentiation</b>	<b>13</b>
<b>13</b>	<b>Answers Lecture 4: Probabilistic Modelling Principles</b>	<b>15</b>
<b>14</b>	<b>Answers Lecture 5: Gradient Descent Convergence</b>	<b>15</b>

## 1 Notation

### 1.1 Sets

Throughout this course, we will be using some standard mathematical notation which may be unfamiliar to some. It's ultimately not that special or even crucial to the overall argument, but it is compact (which is practical), and it helps somewhat with practising with expressing things mathematically. Wikipedia has good definitions on these things too.

---

\*Many thanks to teaching assistants Carles Balsells Rodas, and Alex Spies for their solutions and improvements to the document.

- Notation referring to sets of numbers, e.g. the natural numbers  $\mathbb{N} = \{0, 1, 2, \dots\}$ , integers  $\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$ , or real numbers  $\mathbb{R}$ .
- Vectors are sets containing  $n$  of some type of object, like reals. We denote the set of all such sets using a superscript notation. For example, all  $n$ -dimensional vectors becomes  $\mathbb{R}^n$ .
- With  $x \in \mathcal{S}$  we denote that  $x$  is an element of the set  $\mathcal{S}$ . This allows us to specify that a variable comes from a particular set (or, has a particular type), e.g.  $x \in \mathbb{R}^D$ .
- We sometimes use “set builder” notation. We did this informally above when defining  $\mathbb{N}$ ! Usually this works by specifying elements with some property, e.g.  $\mathbf{S} = \{2n | n \in \mathbb{N}\}$ , which means “all the elements  $2n$  such that  $n$  is a natural number”. This creates the set of all even positive whole numbers.
- We denote the union of two sets (the set with all elements that are in either set or both) as  $A \cup B$ . With set-builder notation this is  $A \cup B = \{x | x \in A \vee x \in B\}$ , where  $\vee$  means “or”.
- We denote the intersection of two sets (the set of all elements that are in both sets) as  $A \cap B = \{x | x \in A \wedge x \in B\}$ .
- For intervals of real numbers, we use brackets,  $[\cdot, \cdot]$ , to denote the elements in the set which are “greater than or equal to” and “less than or equal to” an element, respectively. We use parentheses,  $(\cdot, \cdot)$  to denote a strict lower bound or upper bound on the set, respectively. E.g.  $[1, 5]$  is equivalent to  $1 \leq x \leq 5, x \in \mathbb{R}$ .
- We use the symbol  $\neg$  to denote the complement of a set. Given a set containing all elements under consideration  $\Omega$ ,  $\neg A$  contains all elements of  $\Omega$  that are not in  $A$ , i.e.  $\neg A = \{x \in \Omega | x \notin A\}$ . We can also denote this as  $\neg A = \Omega \setminus A$ .

## 1.2 Probabilities

In this course we will use the notation for probabilities that is common in machine learning. The main advantage is that this notation is shorter, although it does leave certain things implicit. We include this to reduce confusion.

Consider a probability space  $(\Omega, \mathcal{E}, \mathbb{P})$  with sample space  $\Omega$  (all possible outcomes of a random procedure), event space  $\mathcal{E}$  (the set of all sets of outcomes that we assign a probability to), and probability function  $\mathbb{P} : \mathcal{E} \rightarrow [0, 1]$  (a function that assigns a probability to an event), with a random variable  $X : \Omega \rightarrow \mathbb{R}^D$ .

- With  $\mathbb{P}(E)$  we denote the probability of an event  $E \in \mathcal{E}$ , where  $E$  is a set of outcomes.
- Following the usual convention, we use the same notation when considering random variables, e.g.  $\mathbb{P}(X < 2)$  is short for  $\mathbb{P}(\{s \in \Omega : X(s) < 2\})$  (see §6.1 in 50008 *Probability & Statistics*).
- We usually work directly with random variables, and specify all properties using a probability mass function (pmf) or probability density function (pdf). For a specific outcome of the random variable  $\alpha$ , we write:

$$\mathbb{P}(X = \alpha) = p_X(\alpha) \quad \text{for a pmf } p_X(\cdot), \quad (1)$$

$$\mathbb{P}(X \in [a, b]) = \int_a^b p_X(\alpha) d\alpha \quad \text{for a pdf } p_X(\cdot) \text{ with } \alpha \in \mathbb{R}, \quad (2)$$

$$\mathbb{P}(X \in A) = \int_A p_X(\alpha) d\alpha \quad \text{for a pdf } p_X(\cdot) \text{ with } \alpha \in \mathbb{R}^D. \quad (3)$$

- Sometimes we may write vectors in boldface, i.e.  $\mathbf{x} \in \mathbb{R}^D$ . We won’t always though, so keep track of how we define variables!
- We generally denote outcomes of random variables without referring explicitly to the random variable itself. For example, when we refer to an outcome  $\mathbf{x}$ , we implicitly know there is a random variable that can take this value. We usually denote this as the capital, for example here  $X$ .
- Sometimes we abuse notation, and drop the random variable when denoting distributions when the argument of the function identifies it, e.g.  $p(\mathbf{x}) = p_X(\mathbf{x})$ .

- If we want to be explicit about the random variable that we are evaluating the density/mass of, I will write e.g.  $p_{X,Y}(\mathbf{x}, \mathbf{y}) = p_{X|Y}(\mathbf{x}|\mathbf{y})p_Y(\mathbf{y})$ .
- Expectations can be denoted in two ways:

$$\mathbb{E}_X[f(X)] \quad \text{to emphasise that } X \text{ is random, if it is clear what its distribution is,} \quad (4)$$

$$\mathbb{E}_{p(\mathbf{x})}[f(\mathbf{x})] \quad \text{to emphasise that we will be integrating over the distribution } p(\mathbf{x}). \quad (5)$$

In both cases this corresponds to the integral  $\int p(\mathbf{x})f(\mathbf{x})d\mathbf{x}$ .

- Often, densities and pmfs can be discussed in exactly the same way, if we think of the density of a discrete RV as a sum of delta functions. I.e.  $p(\mathbf{x}) = \sum_o \delta(\mathbf{x} - \mathbf{x}_o)p_o$ , where  $\{\mathbf{x}_o\}$  is the set of discrete possible outcomes that  $X$  can take, and  $p_o$  are their corresponding probabilities. This allows us to write an expectation as an integral, regardless of whether the RV is continuous or discrete, because for discrete RVs we get:

$$\mathbb{E}_{p(\mathbf{x})}[f(\mathbf{x})] = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x} = \int \sum_o \delta(\mathbf{x} - \mathbf{x}_o)p_o f(\mathbf{x})d\mathbf{x} = \sum_o f(\mathbf{x}_o)p_o. \quad (6)$$

(A delta function has the property that  $\int_A \delta(\mathbf{x})d\mathbf{x}$  is 1 if  $0 \in A$ , and 0 otherwise. Linearity of integrals still holds. It can often be seen as the limit of a Gaussian distribution with zero variance.)

## 2 Formula Sheet

- Gaussian probability density function (pdf) with input  $\mathbf{x} \in \mathbb{R}^D$ , denoted as  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$  is

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (7)$$

- For a joint Gaussian density

$$p\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}; \begin{bmatrix} \mathbf{m}_\mathbf{x} \\ \mathbf{m}_\mathbf{y} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{\mathbf{xx}} & \boldsymbol{\Sigma}_{\mathbf{xy}} \\ \boldsymbol{\Sigma}_{\mathbf{yx}} & \boldsymbol{\Sigma}_{\mathbf{yy}} \end{bmatrix}\right), \quad (8)$$

we have the conditional density

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}\left(\mathbf{x}; \mathbf{m}_\mathbf{x} + \boldsymbol{\Sigma}_{\mathbf{xy}}\boldsymbol{\Sigma}_{\mathbf{yy}}^{-1}(\mathbf{y} - \mathbf{m}_\mathbf{y}), \boldsymbol{\Sigma}_{\mathbf{xx}} - \boldsymbol{\Sigma}_{\mathbf{xy}}\boldsymbol{\Sigma}_{\mathbf{yy}}^{-1}\boldsymbol{\Sigma}_{\mathbf{yx}}\right). \quad (9)$$

## 3 Warm-up Exercises

To start, here are some exercises which test knowledge which is assumed in the course.

### 3.1 Probability Theory

We assume that you are familiar with probability theory up to the Computing 2nd year 50008 *Probability & Statistics* course. Here are some questions to serve as a refresher. Students who are not familiar with this background should refer to the notes of 50008 *Probability & Statistics* or relevant chapters of [mml]. **We recommend you look at these questions when/before the course starts.** If you need a refresher, or if you do not know the notation, refer to the 50008 *Probability & Statistics* notes, or discuss with a TA.

**Question 1** (Set Theory and Probability). Using the three axioms of probability show that

- Write down the sample space of a dice. In your notation, use the set  $A$  to denote the event of a 3 or 4 occurring. What is the complement of  $A$ , denoted  $\neg A$ ?
- For a problem about lengths, we have a sample space  $\Omega = [0, 1]$ . For  $A = (0.3, 0.4]$ , what is  $\neg A$ ?
- $\mathbb{P}(\neg A) = 1 - \mathbb{P}(A)$
- $\mathbb{P}(\emptyset) = 0$ , where  $\emptyset$  is the empty set

e.  $0 \leq \mathbb{P}(A) \leq 1$

f.  $A \subseteq B \implies \mathbb{P}(A) \leq \mathbb{P}(B)$

*Hint:* Consider the following definition.  $B \setminus A = \{x \in B : x \notin A\}$

g.  $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$

h. (\*) if  $\{A_i\}_{i=1}^\infty \subseteq \Omega$  and  $A_i \subseteq A_{i+1} \forall i$  then:

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) = \lim_{i \rightarrow \infty} \mathbb{P}(A_i)$$

*Hint:* Use axiom 3. \*: The emphasis of this course isn't on these kinds of details, even though this should be doable with 1st-year calculus.

i. For two mutually exclusive events  $A, B$ , what is  $\mathbb{P}(A \cup B)$ ?

See **mml** for a general overview, and §4, §§5.1-5.4 of 50008 *Probability & Statistics* for more details.

**Question 2** (Independent events). **Independent events don't come up as much as independent random variables, so it's ok to just follow this answer, rather than spending lots of time on it.** When tossing two coins (where we care about the order), we have a sample space  $\Omega = \{HH, HT, TH, TT\}$ .

- What outcomes are contained in the event that corresponds to the the first coin being heads? We denote the event  $E_{1H}$ , and others similarly.
- If you assume that all outcomes have equal probability, show that  $E_{1H}$  and  $E_{2T}$  are independent.
- If you assume that  $E_{1H}$  and  $E_{2H}$  are independent and 0.5 each, show that all outcomes must have equal probability.

See §5.3.3 in 50008 *Probability & Statistics*.

**Question 3** (Random Variables). Consider throwing two fair dice.

- What is the sample space for all outcomes that you can get from throwing two dice? We specify the probability of each outcome to be the same.
- Define two random variables  $A, B$  which map the outcome to the face value on each die respectively. Find the probability mass function for  $A$  from the probability on outcomes. The answer will work from the definition of a random variable, but you will probably intuitively get the right answer as well.
- Show that  $A$  and  $B$  are independent.
- Define the random variable  $C = A + B$ . Derive the probability mass function of  $C$ .

See §6 of 50008 *Probability & Statistics*.

**Question 4** (Continuous Random Variables). Consider the random variable  $X$  with a probability density  $p(x) = C \cdot x$  when  $x \in [0, 1]$  and 0 elsewhere.

- Calculate  $C$ .
- Calculate  $\mathbb{P}(0.3 \leq X \leq 0.75)$ .
- Calculate  $\mathbb{P}(X \in [0.3, 0.75] \cup [0.8, 0.9])$ .
- Calculate  $\mathbb{E}_X[X]$ ,  $\mathbb{E}_X[X^2]$ ,  $\mathbb{V}_X[X]$ .

Check your answers by performing numerical integration, e.g. in Python.

See §6.3, §7 of 50008 *Probability & Statistics* or **mml**.

**Question 5** (Joint Discrete Random Variables). Consider two random variables  $A, C$ , where  $A$  is the outcome of one die, and  $C$  gives the sum of  $A$  and the sum of another die  $B$ .

- From intuition, write a table of  $\mathbb{P}(C = c|A = a)$ , which we use to denote the probability of  $C$  taking the value  $c$ , if we know that  $A$  has taken the value  $a$ .
- Write a table of  $\mathbb{P}(C = c, A = a)$ . To help you think it through, consider a tree of outcomes that can occur. This helps illustrate independence between outcomes, which helps you figure out when you can multiply probabilities.
- From the values in the table  $\mathbb{P}(C = c, A = a)$  find  $\mathbb{P}(2 \leq C \leq 4)$  and  $\mathbb{P}(2 \leq C \leq 4, 2 \leq A \leq 4)$ .

We will cover conditional probability more later, but for now just think it through.

**Question 6** (Multivariate Integration). Consider two continuous random variables  $X, Y$  with joint density  $p(x, y) = C \cdot (x^2 + xy)$  when  $x \in [0, 1]$  and  $y \in [0, 1]$ , and 0 elsewhere.

- Find  $C$ .
- Find  $\mathbb{P}(0.3 \leq X \leq 0.5)$ .
- Find  $\mathbb{P}(X < Y)$ . Perform the integration twice in both orders, once integrating over  $x$  first, once by integrating over  $y$  first.
- Bonus:** Convince yourself that you know how to do this for  $p(x, y, z) = C \cdot (x^2 + xyz)$  as well.

Check your answers by performing numerical integration, e.g. in Python.

**Question 7** (Statistics Terminology). Recall the following statistical terminology.

- What is a statistic?
- What is an estimator?
- What is a consistent estimator?
- What is a sample?

### 3.2 Linear Algebra

**Question 8** (Dot product). Compute  $\mathbf{x}^\top \mathbf{y}$  where  $\mathbf{x} = (1, -2, 5, -1)^\top$  and  $\mathbf{y} = (0, 4, -3, 7)^\top$ .

**Question 9** (Matrix product). Compute  $\mathbf{y} = A\mathbf{x}$  as well as the  $\ell_2$  norm of  $\mathbf{x}$  and  $\mathbf{y}$ , where

$$A = \begin{pmatrix} -1 & 4 & 7 & 2 \\ 3 & -2 & -1 & 0 \\ 5 & 3 & 0 & -1 \end{pmatrix}, \quad \mathbf{x} = (-3, 2, 1, 3)^\top.$$

**Question 10** (Basis). Which of the following set of vectors are basis for  $\mathbb{R}^2$ ?

- $\{(1, 1), (1, 0)\}$
- $\{(2, 4), (3, -1)\}$
- $\{(1, -1), (0, 2), (2, 1)\}$
- $\{(2, -1), (-2, 1)\}$
- $\{(0, 3)\}$

**Question 11** (Span of vectors). Which of the following points are within the span of  $\{(-1, 0, 2), (3, 1, 0)\}$ ?

- $(0, 1, 1)$
- $(1, 1, 4)$
- $(2, 1, 1)$
- $(-3, 4, 2)$
- $(0, 0, 0)$

**Question 12** (Rotation matrix in  $\mathbb{R}^2$ ). What is the  $2 \times 2$  matrix that rotates all the non-zero vectors in  $\mathbb{R}^2$  by  $45^\circ$  counter-clockwise?

**Question 13** (Linear equations). Given the following system of linear equations:

$$\begin{aligned}x + 2y &= 2 \\3x + 2y + 4z &= 5 \\-2x + y - 2z &= -1\end{aligned}$$

Answer the following questions:

- Writing this system in a matrix form  $A\mathbf{x} = \mathbf{b}$  with  $\mathbf{x} = (x, y, z)^\top$ . What are  $A$  and  $\mathbf{b}$ ?
- Solve this system, or show that the solution does not exist.
- What is the rank of  $A$ ?

**Question 14** (Eigen decomposition). Consider a matrix  $A \in \mathbb{R}^{d \times d}$  and assume it has an eigen decomposition of  $A = Q\Lambda Q^{-1}$  where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$ . When  $A$  is symmetric we also have  $Q^{-1} = Q^\top$ . Answer the following questions:

- If  $A$  is symmetric, show that  $\mathbf{x}^\top A \mathbf{x} \geq 0$  for any  $\mathbf{x} \in \mathbb{R}^{d \times 1}$  if and only if  $\lambda_i \geq 0$  for all  $i = 1, \dots, d$ .
- Show that  $\text{Tr}(A) = \sum_{i=1}^d \lambda_i$  where  $\text{Tr}(A)$  is the trace of  $A$ .
- Show that  $\det(A) = \prod_{i=1}^d \lambda_i$  where  $\det(A)$  is the determinant of  $A$ .
- Why an entry  $\lambda_i$  in the diagonal matrix  $\Lambda$  is one of the solutions for the equation  $A\mathbf{q} = \lambda\mathbf{q}$ ,  $\mathbf{q} \neq \mathbf{0}$ ?

## 4 Lecture 1: Probability, Vectors, Differentiation

**Question 15** (Vector notation). We define the probability density on the vector  $\mathbf{x} \in \mathbb{R}^3$  with all elements  $0 \leq x_k \leq 1$  as

$$p(\mathbf{x}) = \frac{1}{C}(x_1^2 + x_1x_2 + x_2^2 + 2x_2x_3). \quad (10)$$

Put this into notation that only uses  $\mathbf{x}$  as a single whole vector.

**Question 16** (Noise conditional independence). Consider the probability of the data in linear regression, for a fixed setting of the parameters  $\boldsymbol{\theta}$  and given inputs  $\mathbf{X} \in \mathbb{R}^{N \times D}$  where  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ :

$$p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{X}) = \mathcal{N}(\mathbf{y}; \boldsymbol{\theta}^\top \mathbf{x}, \sigma^2 \mathbf{I}) \quad (11)$$

Show that all  $y_n$ s are independent, for a fixed setting of the parameters  $\boldsymbol{\theta}$  and given inputs  $\mathbf{X}$ .

**Question 17** (Maximum likelihood revision). For a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ .

- Derive the probability distribution for  $N$  iid draws.
- Derive the maximum likelihood estimator for the mean  $\mu$ .

**Question 18** (Maximum likelihood and minimum loss). Show that the solution to the Maximum Likelihood estimator for linear regression is the same as the minimum squared loss estimator.

**Question 19** (MML 5.1-5.3). This is revision. Compute the derivatives for w.r.t.  $x$  for

- $f(x) = \log(x^4) \sin(x^3)$
- $f(x) = (1 + \exp(-x))^{-1}$
- $f(x) = \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

## 5 Lecture 2: Vector Differentiation

**Question 20** (Circle). Consider a vector function  $\mathbf{x}(t) = [\cos t \quad \sin t]^\top$ .

- Draw the set of points that this function passes through.
- To build intuition, draw the velocity vector at a few points by considering the direction that the point moves in.
- Find the derivative  $d\mathbf{x}/dt$ . Draw this vector for some point  $t$ .

**Question 21** (Index notation). Turn the following matrix-vector expressions into index notation:

- |                    |  |
|--------------------|--|
| a. $ABC\mathbf{x}$ | c. $\text{Tr}(AB)$                     |
| b. $\text{Tr}(A)$  | d. $\mathbf{y}^\top A^\top \mathbf{x}$ |

Turn the following index expressions back to matrix-vector notation:

- |                                      |                             |
|--------------------------------------|-----------------------------|
| a. $\sum_{ijk} A_{ij} B_{jk} C_{ki}$ | c. $x_i x_j$                |
| b. $b_i + \sum_j A_{ij} b_j$         | d. $\sum_j \delta_{ij} a_j$ |

**Question 22** (Index notation proofs). Using index notation, show that

- $\mathbf{x}^\top A \mathbf{y} = \mathbf{y}^\top A \mathbf{x}$  if  $A$  is symmetric, i.e.  $A = A^\top$ .
- $\mathbf{x}^\top \mathbf{y} = \text{Tr}(\mathbf{x}^\top \mathbf{y}) = \text{Tr}(\mathbf{y}^\top \mathbf{x})$ , for  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$ .
- $\text{Tr}(ABC) = \text{Tr}(CAB)$ .

**Question 23** (MML 5.5-5.6). First find the dimensions, then the Jacobian. It's probably easiest here to use index notation.

- $f(\mathbf{x}) = \sin(x_1) \cos(x_2)$ , find  $df/d\mathbf{x}$ .
- $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{y}$ , find  $df/d\mathbf{x}$ .
- $f(\mathbf{x}) = \mathbf{x} \mathbf{x}^\top$ , find  $df/d\mathbf{x}$ .
- $f(\mathbf{t}) = \sin(\log(\mathbf{t}^\top \mathbf{t}))$ , find  $df/d\mathbf{t}$ .
- $f(\mathbf{X}) = \text{Tr}(\mathbf{A} \mathbf{X} \mathbf{B})$  for  $\mathbf{A} \in \mathbb{R}^{D \times E}$ ,  $\mathbf{X} \in \mathbb{R}^{E \times F}$ ,  $\mathbf{B} \in \mathbb{R}^{F \times D}$ , find  $df/d\mathbf{X}$ .

**Question 24** (MML 5.7-5.8: Chain rule). Compute the derivatives  $df/d\mathbf{x}$  of the following functions.

- First, write out the chain rule for the given decomposition.
  - Give the shapes of intermediate results, and make clear which dimension(s) will be summed over.
  - Provide expressions for the derivatives, and describe your steps in detail. Providing an expression means specifying everything up to the point where you could implement it.
  - Give the results in vector notation if you can.
- $f(z) = \log(1 + z)$ ,  $z = \mathbf{x}^\top \mathbf{x}$ ,  $\mathbf{x} \in \mathbb{R}^D$ .
  - $f(\mathbf{z}) = \sin(\mathbf{z})$ ,  $\mathbf{z} = \mathbf{A} \mathbf{x} + \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{E \times D}$ . What sizes are  $\mathbf{x}$  and  $\mathbf{b}$ ?
  - $f(z) = \exp(-\frac{1}{2}z)$ ,  $z = \mathbf{y}^\top \mathbf{S}^{-1} \mathbf{y}$ ,  $\mathbf{y} = \mathbf{x} - \boldsymbol{\mu}$ .
  - $f(\mathbf{A}) = \text{Tr}(\mathbf{A})$ ,  $\mathbf{A} = \mathbf{x} \mathbf{x}^\top + \sigma^2 \mathbf{I}$ .
  - $f(\mathbf{z}) = \tanh(\mathbf{z})$ ,  $\mathbf{z} = \mathbf{A} \mathbf{x} + \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{M \times N}$ .
  - $f(\mathbf{A}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$ ,  $\mathbf{A} = \mathbf{x} \mathbf{x}^\top$ .

Remember: Generally, scalar functions are applied elementwise to vectors/matrices.

**Question 25** (Hessian of Linear Regression). For the stationary point of linear regression, find the Hessian, and prove that it is positive definite, perhaps by making some assumptions. Discuss your assumptions.

## 6 Lecture 3: Automatic Differentiation

**Question 26** (Product rule). Consider the function  $f(a, b) = a \cdot b$ , where  $a = a(x), b = b(x)$ , i.e. unspecified functions of  $x$ .

- Show that by following forward mode autodiff, you effectively calculate the product rule.
- Show that if  $a(x) = x, b(x) = x$ , which means that the overall function  $f(x) = x^2$ , the gradient that is computed will be  $2x$ .

(Note from MvdW (autumn 2022): I somewhat messily described this on the board. The question is included here to provide a clearer explanation.)

**Question 27** (Multivariate Autodiff). This is a rather big question that should test your understanding of all material in the first three lectures. Consider the overall function  $f(\ell, X)$  consisting of the parts:

$$f = \mathbf{y}^\top (\mathbf{K}_1 + \mathbf{K}_2)^{-1} \mathbf{y}, \quad (12)$$

$$\mathbf{K}_a = \exp(\Lambda_a), \quad (13)$$

$$\Lambda_a = -\frac{\mathbf{D}_a}{2\ell_a^2}, \quad (14)$$

$$\mathbf{D}_a = (\mathbf{X}[:, \text{None}, a] - \mathbf{X}[\text{None}, :, a])^2, \quad (15)$$

where we use `numpy` broadcasting notation in the final equation.

- Given  $\ell \in \mathbb{R}^2$  and  $\mathbf{X} \in \mathbb{R}^{N \times 2}$ , find the shape of all intermediate computations.
- Draw the computational graph for  $f(\ell, X)$ .
- For forward and reverse mode differentiation, state which intermediate derivatives are computed at each step, and their computational and memory costs.

## 7 Lecture 4: Probabilistic Modelling Principles

**Question 28** (Training translation models). Imagine you want to train a neural network  $T_\theta(\cdot)$  to translate French words to English words. Assume you are given a dataset  $\mathcal{D} = \{(f_n, e_n)\}_{n=1}^N$  where  $f_n$  is a French word and  $e_n$  is an English word. Suppose the vocabulary of French and English is  $\mathcal{F}$  and  $\mathcal{E}$ , respectively.

- Assuming a probabilistic model  $p(e|T_\theta(f))$ , which distribution would you choose for this model?
- Continuing a), what is the corresponding MLE objective?

**Question 29** (Clustering). We consider a clustering task where given a dataset  $\mathcal{D} = \{x_1, \dots, x_N\}$ , we would like to group them into  $K$  clusters. The model we will use here is a Gaussian mixture model:

$$\text{GMM: } p(x|\theta) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \sigma^2), \quad \theta = \{\pi_k, \mu_k, \sigma^2\}_{k=1}^K.$$

- What is the MLE objective for this clustering task?
- Derive the gradient of the MLE objective w.r.t.  $\mu_k$ . What is the fixed-point equation for finding the optimal  $\{\mu_k\}$  parameters?

**Question 30** (Geometric interpretation of linear regression). Consider the following linear regression model:

$$y = \theta^\top \phi(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2).$$

For a given dataset  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ , Writing  $\Phi = (\phi(\mathbf{x}_1), \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_N))^\top$  and  $\mathbf{y} = (y_1, \dots, y_N)^\top$ , we have the optimal solution satisfies  $\theta^* = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$ . Show that by using the optimal parameter  $\theta^*$ , the prediction  $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_N)$ ,  $\hat{y}_n = (\theta^*)^\top \phi(\mathbf{x}_n)$  is the projection of  $\mathbf{y}$  onto the sub-space spanned by the columns of  $\Phi$ .

(Hint: consider singular value decomposition.)



## 8 Lecture 5: Gradient Descent Convergence

**Question 31** (Rayleigh quotient). The *Rayleigh quotient* is defined for a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  and a non-zero vector  $\mathbf{x} \in \mathbb{R}^{d \times 1}$ :

$$R(\mathbf{A}, \mathbf{x}) = \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\|\mathbf{x}\|_2^2}, \quad \|\mathbf{x}\|_2^2 = \mathbf{x}^\top \mathbf{x}.$$

Show that  $R(\mathbf{A}, \mathbf{x}) \in [\lambda_{\min}(\mathbf{A}), \lambda_{\max}(\mathbf{A})]$ .

This result immediately indicates that  $\lambda_{\min}(\mathbf{A})\|\mathbf{x}\|_2^2 \leq \mathbf{x}^\top \mathbf{A} \mathbf{x} \leq \lambda_{\max}(\mathbf{A})\|\mathbf{x}\|_2^2$ , which is used to prove gradient descent convergence.

**Question 32** (Gradient descent with pre-conditioning). Consider the following update rule named *pre-conditioned gradient descent*:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \gamma_t \mathbf{P}_t^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t).$$

Here  $\mathbf{P}_t$  is called *pre-conditioner* at time step  $t$ . We consider linear regression as an example, and assume constant learning rate and pre-conditioner, i.e.,  $\gamma_t = \gamma$  and  $\mathbf{P}_t = \mathbf{P}$  for all  $t$ . Show that with an appropriate choice of the pre-conditioner  $\mathbf{P}$ , we can achieve a robust selection of the learning rate  $\gamma$ , i.e., if the selected  $\gamma$  works for an initialisation  $\boldsymbol{\theta}_0$ , it will also work for all other initialisations.

Hints: you can follow the below steps to solve the question:

1. Work out the pre-conditioned gradient descent update in linear regression, and derive  $\boldsymbol{\theta}_t$  as a function of  $\boldsymbol{\theta}_0$ ,  $\gamma$ ,  $\mathbf{P}$  and the dataset  $(\mathbf{X}, \mathbf{y})$ ;
2. For a given  $\mathbf{P}$ , work out the learning rates  $\gamma_{\min}$  and  $\gamma_{\max}$  such that pre-conditioned gradient descent converges when  $\gamma < \gamma_{\min}$ , or diverges when  $\gamma \geq \gamma_{\max}$ ;
3. Select  $\mathbf{P}$  such that  $\gamma_{\min} = \gamma_{\max}$ , therefore there exist no interval (like  $[\gamma_{\min}, \gamma_{\max})$ ) such that convergence depends on initialisation when  $\gamma$  falls into such interval.

**Question 33** (Momentum gradient descent). Consider the following update rule named *momentum gradient descent*, with constant learning rate  $\gamma$  and momentum step-size  $\alpha$ :

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \gamma \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t) + \alpha \Delta \boldsymbol{\theta}_t, \\ \Delta \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t, \quad \Delta \boldsymbol{\theta}_0 = \mathbf{0}. \end{aligned}$$

Show that solving linear regression using momentum gradient descent, if converges, converges to  $\boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ .

Hint: follow the below steps and practice your linear algebra skills :)

1. Write down the update equations for the parameters  $\boldsymbol{\theta}_t$  and the momentum  $\Delta \boldsymbol{\theta}_t$ ;
2. Collect both terms as a long vector  $(\boldsymbol{\theta}_t^\top, \Delta \boldsymbol{\theta}_t^\top)^\top$ , and merge the two linear update equations in step 1 into one “joint” linear equation using block matrices;
3. Apply the analysis techniques for gradient descent convergence for linear regression to show the converged solution (if converges).

## 9 Lecture N: Multivariate Probability

**Question 34** (Vector independence). **While you can probably figure this one out already, we will discuss this in more detail later.** Consider the density on  $\mathbf{x} \in \mathbb{R}^4$  with all elements  $0 \leq x_k \leq 1$  as

$$p(\mathbf{x}) = \mathbf{x}^\top \begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{bmatrix} \mathbf{x}. \quad (16)$$

- a. Rewrite the density in terms of  $\tilde{\mathbf{x}} = [x_1, x_3, x_2, x_4]^\top$ . Note that you can do this by a substitution  $\mathbf{x} = \mathbf{P} \tilde{\mathbf{x}}$ , where  $\mathbf{P}$  is a permutation matrix. You will see that you just need to swap the relevant rows and columns of the matrix. However, make sure that you understand the mathematical steps that really show this.
- b. Divide up  $\mathbf{x}$  into two sub vectors  $\mathbf{y} = [x_2, x_4]^\top$  and  $\mathbf{z} = [x_1, x_3]^\top$ . Show that  $\mathbf{y} \perp \mathbf{z}$ , i.e. that they are independent.

## 10 Warm-up Exercises Answers

### 10.1 Warm-up Exercises

#### Question 1 – Set Theory and Probability

- a. We can choose any representation denoting the events, e.g. using abstract symbols  $\Omega = \{\square, \square\square, \square\square\square, \square\square\square\square, \square\square\square\square\square, \square\square\square\square\square\square\}$ .  
Alternatively, we can represent each of the outcomes as a number  $\Omega = \{1, 2, 3, 4, 5, 6\}$ .

Following the latter notation,  $A = \{3, 4\}$ , and  $\neg A = \{1, 2, 5, 6\}$ .

- b. Length problem with sample space  $\Omega = [0, 1]$ .

$$\neg A = [0, 0.3] \cup (0.4, 1]$$

- c.  $P(\neg A) = 1 - P(A)$

Since  $\neg A$  and  $A$  are mutually exclusive:  $A \cup \neg A = \Omega$  and  $A \cap \neg A = \emptyset$ .

By combining axiom 2 and 3:  $P(A) + P(\neg A) = P(A \cup \neg A) = P(\Omega) = 1$

Thus:  $P(\neg A) = 1 - P(A)$

- d.  $P(\emptyset) = 0$ , where  $\emptyset$  is the empty set

Given the sample space,  $\Omega$ , its complementary is the empty set  $\emptyset$ .

We use property (c) and axiom 2:  $P(\emptyset) = 1 - P(\Omega) = 1 - 1 = 0$ .

- e.  $0 \leq P(A) \leq 1$

We use property (c) and axiom 1.

Consider an event  $A$ , where  $P(A) \geq 0$  and  $P(\neg A) \geq 0$  by axiom 1.

Then,  $P(\neg A) = 1 - P(A) \geq 0 \implies 1 \geq P(A)$ .

By joining both inequalities,  $0 \leq P(A) \leq 1$ .

- f.  $A \subseteq B \implies P(A) \leq P(B)$

*Hint:* Consider the following definition.  $B \setminus A = \{x \in B : x \notin A\}$

Assume  $A \subseteq B$  and construct  $B$  as the union of two disjoint sets:  $B = B \setminus A \cup A$ .

Then,  $B \setminus A \cap A = \emptyset$  by definition of  $B \setminus A$ . By axiom 1, we have  $P(B \setminus A) \geq 0$ .

Use axiom 3:  $P(B) = P(B \setminus A) + P(A) \geq P(A) \implies P(A) \leq P(B)$ .

- g.  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

Define the union  $(A \cup B)$  in terms of two disjoint sets.  $(A \cup B) = A \cup B \setminus A$ , where  $A \cap B \setminus A = \emptyset$ .

Use axiom 3:  $P(A \cup B) = P(A) + P(B \setminus A)$ .

To compute  $P(B \setminus A)$ , we define  $B$  in terms of  $A$ , and the union of two disjoint sets:  $B = (B \cap A) \cup (B \setminus A)$ , where  $(B \cap A) \cap (B \setminus A) = \emptyset$  by definition.

Use axiom 3 again:  $P(B) = P(B \cap A) + P(B \setminus A) \implies P(B \setminus A) = P(B) - P(B \cap A)$ .

Finally:  $P(A \cup B) = P(A) + P(B \setminus A) = P(A) + P(B) - P(B \cap A)$ .

- h. (\*) if  $\{A_i\}_{i=1}^{\infty} \subseteq \Omega$  and  $A_{i-1} \subseteq A_i \quad \forall i > 0$  then:

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \lim_{i \rightarrow \infty} P(A_i)$$

*Hint:* Use axiom 3.

Let us define the following:  $A := \bigcup_{i=1}^{\infty} A_i$ . We would like to write  $A$  in terms of disjoint sets to use axiom 3.

$$A_{i-1} \subseteq A_i \quad \forall i > 0 \implies A = \bigcup_{i=1}^{\infty} A_i \setminus A_{i-1} \quad (17)$$

where the expression holds if we have  $A_0 = \emptyset$ . We regard 17 as starting with  $A_1$  and adding the new information from  $A_2, A_3, \dots$  (e.g  $A_2 \setminus A_1, A_3 \setminus A_2, \dots$ ).

$$P(A) = P\left(\bigcup_{i=1}^{\infty} A_i \setminus A_{i-1}\right) = \sum_{i=1}^{\infty} P(A_i \setminus A_{i-1}) \quad (\text{by axiom 3}) \quad (18)$$

$$P(A) = \sum_{i=1}^{\infty} P(A_i \setminus A_{i-1}) = \lim_{n \rightarrow \infty} \sum_{i=1}^n P(A_i \setminus A_{i-1}) \quad (\text{the infinite summation is a limit}) \quad (19)$$

From (f), we have  $P(A_i) = P(A_i \setminus A_{i-1}) + P(A_{i-1}) \implies P(A_i \setminus A_{i-1}) = P(A_i) - P(A_{i-1})$ . Then,

$$P(A) = \lim_{n \rightarrow \infty} \sum_{i=1}^n P(A_i) - P(A_{i-1}) = \lim_{n \rightarrow \infty} \left( \sum_{i=1}^n P(A_i) - \sum_{i=1}^{n-1} P(A_i) \right) = \lim_{n \rightarrow \infty} P(A_n) \quad (20)$$

where we used  $P(A_0) = P(\emptyset) = 0$  from (d).

In summary:

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = P(A) = \lim_{i \rightarrow \infty} P(A_i) \quad (21)$$

## Question 7 – Statistics Terminology

- A statistic is a function that is computed from data. For example, take a data set  $X = \{x_1, x_2, x_3, \dots\}$  where we compute the empirical mean  $\bar{X} = \frac{1}{|X|} \sum_n x_n$ .
- An estimator is a function of data that tries to estimate an unknown quantity. Estimators are statistics. Some statistics are also estimators. For example, if we have some data set from that is sampled from some unknown density  $p(x)$ , then its mean is unknown, and  $\bar{X}$  is an estimator of it.
- A consistent estimator finds the correct value of the unknown quantity if the dataset grows to infinity. We will prove that  $\bar{X}$  is a consistent estimate of  $\int p(x)xdx$  later on in the course.
- A sample from a random variable is an outcome of the random experiment it represents. For example, you can have a random variable representing the outcome of a coin toss. A sample from it would be heads or tails. We sampled a random variable independently many times, then the outcomes would occur with the frequency specified by the probability distribution of the random variable. Thinking about sampling outcomes from a random variable is often a helpful conceptual technique to think about randomness.

## 10.2 Linear Algebra

**Question 8**  $\mathbf{x}^\top \mathbf{y} = 1 \times 0 + (-2) \times 4 + 5 \times (-3) + (-1) \times 7 = 0 + (-8) + (-15) + (-7) = -30$ .

**Question 9**  $\mathbf{y} = (24, -14, -12)^\top$ ,  $\|\mathbf{x}\|_2 = \sqrt{23}$ ,  $\|\mathbf{y}\|_2 = \sqrt{916}$ .

Note that by definition the  $\ell_2$  norm of a vector is  $\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^\top \mathbf{x}}$ .

**Question 10** 1, 2.

A set of vectors  $\{\mathbf{b}_1, \dots, \mathbf{b}_K\}$  with  $\mathbf{b}_k \in \mathbb{R}^d$  can form a basis of  $\mathbb{R}^d$  iff  $K = d$  the vectors are linearly independent to each other.

**Question 11** 2, 5.

A point  $\mathbf{x} \in \mathbb{R}^d$  is in  $\text{span}(\{\mathbf{b}_1, \dots, \mathbf{b}_K\})$  with  $\mathbf{b}_k \in \mathbb{R}^d$  iff we can find  $a_1, \dots, a_K \in \mathbb{R}$  such that  $\mathbf{x} = \sum_{k=1}^K a_k \mathbf{b}_k$ .

**Question 12** The rotation matrix is

$$\begin{pmatrix} \cos \frac{\pi}{4} & -\sin \frac{\pi}{4} \\ \sin \frac{\pi}{4} & \cos \frac{\pi}{4} \end{pmatrix}.$$

**Question 13** a) The matrix  $A$  and vector  $\mathbf{b}$  are

$$A = \begin{pmatrix} 1 & 2 & 0 \\ 3 & 2 & 4 \\ -2 & 1 & -2 \end{pmatrix}, \quad \mathbf{b} = (2, 5, 1)^\top.$$

b) The inverse of  $A$  is

$$A^{-1} = \begin{pmatrix} 2/3 & -1/3 & -2/3 \\ 1/6 & 1/6 & 1/3 \\ 7/12 & 5/12 & 1/3 \end{pmatrix}.$$

Therefore we have  $\mathbf{x} = A^{-1}\mathbf{b} = (-1, 3/2, 43/12)^\top$ .

c)  $\text{rank}(A) = 3$ : as  $A$  is invertible, it must have full rank.

**Question 14** a) When  $A$  is symmetric, then  $A = Q\Lambda Q^\top$ , and  $\mathbf{x}^\top A \mathbf{x} = \mathbf{x}^\top Q\Lambda Q^\top \mathbf{x} = (Q^\top \mathbf{x})^\top \Lambda (Q^\top \mathbf{x})$ . As  $Q$  is an orthonormal matrix, we have  $\mathbf{x} \rightarrow Q^\top \mathbf{x}$  a one-to-one mapping. Therefore we have

$$\mathbf{x}^\top A \mathbf{x} = \mathbf{z}^\top \Lambda \mathbf{z} = \sum_{i=1}^d \lambda_i z_i^2, \quad \mathbf{z} = (z_1, \dots, z_d)^\top = Q^\top \mathbf{x}.$$

Therefore  $\mathbf{x}^\top A \mathbf{x} \geq 0 \Leftrightarrow \sum_{i=1}^d \lambda_i z_i^2 \geq 0$ . This is true for any  $\mathbf{x} \in \mathbb{R}^{d \times 1}$  if and only if  $\lambda_i \geq 0$  for all  $i = 1, \dots, d$ .

b) We use the permutation invariance property of matrix trace to show the result:

$$\text{Tr}(A) = \text{Tr}(Q\Lambda Q^{-1}) = \text{Tr}(Q^{-1}Q\Lambda) = \text{Tr}(\Lambda) = \sum_{i=1}^d \lambda_i.$$

c) We use the product rule of matrix determinant to show the result:

$$\det(A) = \det(Q\Lambda Q^{-1}) = \det(Q)\det(\Lambda)\det(Q^{-1}) = \det(Q)\det(\Lambda)\det(Q)^{-1} = \det(\Lambda) = \prod_{i=1}^d \lambda_i.$$

d) Let us assume the statement is false, i.e., there exists a solution  $\lambda^* \neq \lambda_i, \forall i = 1, \dots, d$  for the equation  $A\mathbf{q} = \lambda\mathbf{q}, \mathbf{q} \neq 0$ . Then we can rewrite the equation as

$$A\mathbf{q} = \lambda^*\mathbf{q} \Rightarrow (A - \lambda^*I)\mathbf{q} = \mathbf{0} \Rightarrow Q(\Lambda - \lambda^*I)Q^{-1}\mathbf{q} = \mathbf{0}.$$

By definition, the column vectors of  $Q$  forms a basis of  $\mathbb{R}^d$ . Notice that the diagonal entries of  $\Lambda - \lambda^*I$  are non-zero as we assume  $\lambda^* \neq \lambda_i$ . This indicates a contradiction to the assumption of  $\mathbf{q} \neq 0$ :

$$Q(\Lambda - \lambda^*I)Q^{-1}\mathbf{q} = \mathbf{0} \Rightarrow Q^{-1}\mathbf{q} = \mathbf{0} \Rightarrow \mathbf{q} = \mathbf{0}.$$

## 11 Answers Lecture 2: Vector Differentiation

**Question 20 – Circle** Answer discussed in lectures.

**Question 25 – Hessian of Linear Regression** The objective function and gradient w.r.t.  $\boldsymbol{\theta}$  (see lectures) for Linear Regression is

$$L(\boldsymbol{\theta}) = \|\mathbf{y} - \Phi(X)\boldsymbol{\theta}\|^2, \quad \frac{dL}{d\boldsymbol{\theta}} = 2(\Phi(X)\boldsymbol{\theta} - \mathbf{y})^\top \Phi(X). \quad (22)$$

We begin by finding the Hessian, i.e. the matrix containing all second partial derivatives. We need to do this in index notation, as the vector conventions of our vector chain rule break down. So we first write

the derivative in index notation, and then we take the derivative again, after which we return to vector notation:

$$\frac{\partial}{\partial \theta_j} \left( \frac{\partial L}{\partial \theta_i} \right) = \frac{\partial}{\partial \theta_j} \left( 2 \sum_k \left( \sum_m \Phi_{km} \theta_m - y_k \right) \Phi_{ki} \right) = \frac{\partial}{\partial \theta_j} \left( 2 \sum_k \left( \sum_m \Phi_{km} \theta_m - y_k \right) \Phi_{ki} \right) \quad (23)$$

$$= 2 \sum_{km} \Phi_{km} \delta_{mj} \Phi_{ki} = 2 \sum_k \Phi_{kj} \Phi_{ki}, \quad (24)$$

$$\implies \mathbf{H}_\theta(L) = 2\Phi(X)^\top \Phi(X). \quad (25)$$

The Hessian doesn't depend on the parameter  $\theta$ , so if we prove that the matrix is positive definite, then the local where  $\frac{dL}{d\theta} = 0$  (see lecture slides) will be a minimum. For a matrix to be PD, we need  $\mathbf{v}^\top \mathbf{H} \mathbf{v} > 0$  for all  $\mathbf{v}$ . We substitute our Hessian into  $\mathbf{H}$  to prove this

$$\mathbf{v}^\top \mathbf{H} \mathbf{v} = 2\mathbf{v}^\top \Phi(X)^\top \Phi(X) \mathbf{v} \quad (26)$$

$$= \mathbf{w}^\top \mathbf{w} = \sum_i w_i^2, \quad \text{with } \mathbf{v} = \Phi(X) \mathbf{v}. \quad (27)$$

This already shows that  $\mathbf{v}^\top \mathbf{H} \mathbf{v} \geq 0$ , with equality if there exists a  $\mathbf{v}$  such that  $\Phi(X) \mathbf{v} = 0$ . So now we need to prove that *there cannot be* a  $\mathbf{v}$  for which  $\Phi(X) \mathbf{v} = 0$ . If  $\text{rank } \Phi(X) \geq M$ , then this will not happen, by the rank-nullity theorem [mml].

At this point, we need to assume this is the case. For full marks though, you should state the implications on the problem at hand, rather than in abstract maths. One *necessary* implication of this is that  $N \geq M$ . This is only a necessary condition, rather than a sufficient one, since even if  $N \geq M$ ,  $\Phi(X)$  can still have many linearly dependent rows. This will at least happen if you observe repeated input points. However, to prove more than this, you need more information about  $\Phi(X)$ .<sup>1</sup>

So to summarise, we could prove that **if  $\text{rank } \Phi(X) \geq M$ , which at least needs  $N \geq M$ , then Linear Regression has a single minimum solution.**

If we are coding up a linear regression problem, and we want to check numerically for a *specific* regression problem whether there is a unique solution, we can compute the eigenvalues of  $\Phi(X)^\top \Phi(X)$ , and see if they are all positive. This implies a PD Hessian because

$$\mathbf{v}^\top \mathbf{H} \mathbf{v} = \mathbf{v}^\top \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1} \mathbf{v} \quad (\text{eigenvalue decomposition}) \quad (28)$$

$$= \mathbf{v}^\top \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top \mathbf{v} \quad (\mathbf{H} = \mathbf{H}^\top, \text{ so } \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1} = (\mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1})^\top, \text{ so } \mathbf{Q}^{-1} = \mathbf{Q}^\top) \quad (29)$$

$$= \mathbf{z}^\top \mathbf{\Lambda} \mathbf{z}, \quad (30)$$

which is only  $> 0$  if all the elements in the diagonal matrix  $\mathbf{\Lambda}$  are positive.

If any of the linear algebra was unfamiliar, I recommend looking at chapter 2 in **mml**, particularly §2.3, §2.6, and §2.7, or the 1st year linear algebra course.

## 12 Answers Lecture 3: Automatic Differentiation

**Question 26 – Product rule** We begin by drawing the computational graph (fig. 1). We now find the primal trace and the forward tangent trace:

$$v_0 = x \quad \frac{\partial v_0}{\partial x} = 1 \quad (31)$$

$$v_1 = a(x) \quad \frac{\partial v_1}{\partial x} = \frac{\partial v_1}{\partial v_0} \frac{\partial v_0}{\partial x} = \frac{\partial a(x)}{\partial x} \quad (32)$$

$$v_2 = b(x) \quad \frac{\partial v_2}{\partial x} = \frac{\partial v_2}{\partial v_0} \frac{\partial v_0}{\partial x} = \frac{\partial b(x)}{\partial x} \quad (33)$$

$$v_3 = v_1 \cdot v_2 \quad \frac{\partial v_3}{\partial x} = \sum_{j \in \text{inputs}(3)} \frac{\partial v_3}{\partial v_j} \frac{\partial v_j}{\partial x} = v_2 \frac{\partial a(x)}{\partial x} + v_1 \frac{\partial b(x)}{\partial x}. \quad (34)$$

This means that for any  $x$ , forward mode autodiff calculates the derivative to be:

$$\frac{df}{dx} = b(x) \frac{da(x)}{dx} + a(x) \frac{db(x)}{dx}. \quad (35)$$

<sup>1</sup>A case that is harder to think about is if you observe points that make the feature vectors  $\phi(\mathbf{x}_n)$  linearly dependent. One example is if you have a 2D input with  $\phi(\mathbf{x}) = \mathbf{x}^\top$ , and all your input points lie on a line.

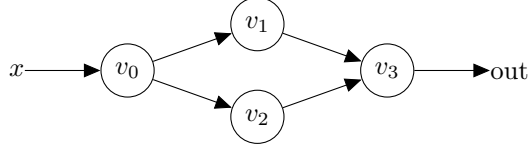


Figure 1: Computational graph for Question 26 – Product rule.

Which is the product rule.

If we substitute  $a(x) = x, b(x) = x$ , then we obtain  $df/dx = 2x$ , as expected.

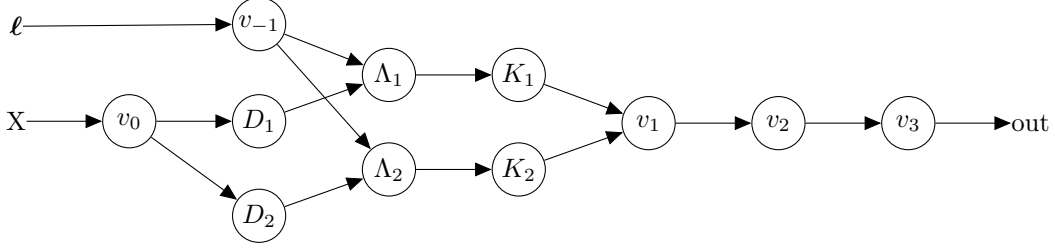


Figure 2: Computational graph for Question 27 – Multivariate Autodiff, where we define  $v_1 = K_1 + K_2$ ,  $v_2 = v_1^{-1}$ , and  $v_3 = \mathbf{y}^\top v_1 \mathbf{y}$ .

#### Question 27 – Multivariate Autodiff

a.

$$D_a : \mathbb{R}^{N \times 2} \rightarrow \mathbb{R}^{N \times N} \quad (36)$$

$$\Lambda_a : \mathbb{R}^{N \times N} \times \mathbb{R}^2 \rightarrow \mathbb{R}^{N \times N} \quad (37)$$

$$K_a : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^{N \times N} \quad (38)$$

$$f : \mathbb{R}^{N \times N} \times \mathbb{R}^{N \times N} \rightarrow \mathbb{R} \quad (39)$$

b. See fig. 2. We use names evident from the question for some nodes, but give new names to some additional intermediate notes.

c. Let's first consider **forward mode** for the derivatives w.r.t. X.

$$v_{-1} = \ell \quad \dot{v}_{-1,iab} = \frac{\partial [v_{-1}]_i}{\partial X_{ab}} = 0 \quad v_{-1} \in \mathbb{R}^{2 \times (N \times 2)}, O(N) \quad (40)$$

$$v_0 = X \quad \dot{v}_0 = \frac{\partial [v_0]_{ij}}{\partial X_{ab}} = \delta_{ia} \delta_{jb} \quad \dot{v}_0 \in \mathbb{R}^{(N \times 2) \times (N \times 2)}, O(N^2) \quad (41)$$

$$D_z = \dots \quad \frac{\partial [D_z]_{nm}}{\partial v_{0ij}} = \frac{\partial}{\partial v_{0ij}} (v_{0nz} - v_{0mz})^2 = 2(v_{0nz} - v_{0mz})(\delta_{ni} \delta_{zj} - \delta_{mi} \delta_{zj}) \quad \dot{D}_z \in \mathbb{R}^{(N \times N) \times (N \times 2)} \quad (42)$$

$$\dot{D}_{znab} = \left[ \frac{\partial D_z}{\partial v_0} \dot{v}_0 \right]_{nab} \quad O(N) \text{ for sum, so total } O(N^4). \quad (43)$$

$$= 2(v_{0nz} - v_{0mz})(\dot{v}_{0nzab} - \dot{v}_{0mzab}) \quad \text{Structure allows } O(N^3). \quad (44)$$

$$\Lambda_z = -\frac{D_z}{2v_{-1z}^2} \quad \frac{\partial \Lambda_{zij}}{\partial D_{znm}} = -\frac{D_z}{2v_{-1z}^2} \delta_{in} \delta_{jm} \quad \dot{\Lambda}_z \in \mathbb{R}^{(N \times N) \times (N \times 2)} \quad (45)$$

$$\frac{\partial \Lambda_{zij}}{\partial v_{-1,k}} = \frac{D_{zij}}{v_{-1z}^3} \delta_{zk} \quad (46)$$

$$\dot{\Lambda}_{zijab} = \left[ \frac{\partial \Lambda_z}{\partial D_z} \dot{D}_z + \frac{\partial \Lambda_z}{\partial v_{-1}} \dot{v}_{-1} \right]_{ijab} \quad O(N^2) \text{ for sum, so total } O(N^5). \quad (47)$$

d. To be continued...

## 13 Answers Lecture 4: Probabilistic Modelling Principles

### Question 28

- a. Choose a categorical distribution. Let  $T_{\theta}(\cdot) : \mathcal{F} \rightarrow \mathbb{R}^{|\mathcal{E}|}$  maps a French word  $f$  to a real-value vector of length  $|\mathcal{E}|$ . Then define

$$p(e|T_{\theta}(f)) = \text{Categorical}(\text{softmax}(T_{\theta}(f))).$$

- b. The MLE objective is (if we write  $e_n$  using one-hot encoding:  $e_n = (0, \dots, 0, 1, 0, \dots, 0)$ )

$$\theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p(e_n | T_{\theta}(f_n)) = \arg \max_{\theta} \frac{1}{N} \sum_{n=1}^N \log \frac{\exp[T_{\theta}(f_n)]^{\top} e_n}{\sum_{i=1}^{|\mathcal{E}|} \exp[T_{\theta}(f_n)_i]}$$

### Question 29

- a. With i.i.d. assumption, the MLE objective is:

$$L(\theta) = \frac{1}{N} \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \sigma^2).$$

- b. The gradient of the MLE objective w.r.t.  $\mu_k$  is

$$\nabla_{\mu_k} L(\theta) = \frac{1}{N} \sum_{n=1}^N \frac{\pi_k \mathcal{N}(x_n; \mu_k, \sigma^2)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \sigma^2)} \frac{x_n - \mu_k}{\sigma^2}.$$

Setting  $\nabla_k L(\theta) = 0$  for all  $k$ , we have the fixed-point equation as

$$\mu_k = \frac{1}{N} \sum_{n=1}^N \frac{\pi_k \mathcal{N}(x_n; \mu_k, \sigma^2)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \sigma^2)} x_n.$$

**Question 30** Using matrix-vector notations we have  $\hat{\mathbf{y}} = \Phi \theta^* = \Phi (\Phi^{\top} \Phi)^{-1} \Phi^{\top} \mathbf{y}$ . Writing the SVD of  $\Phi = \mathbf{U} \Sigma \mathbf{V}^{\top}$ , it is easy to show that  $\Phi (\Phi^{\top} \Phi)^{-1} \Phi^{\top} = \mathbf{U} \mathbf{U}^{\top}$ . Notice that  $\mathbf{U}$  contains basis vectors which span to the same subspace spanned by the column vectors of  $\Phi$ .

## 14 Answers Lecture 5: Gradient Descent Convergence

**Question 31** As  $\mathbf{A}$  is symmetric, we can write the eigen-decomposition formula as  $\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{\top}$  with  $\mathbf{Q}$  containing an orthonormal basis of  $\mathbb{R}^{d \times 1}$ . Then using the fact that  $\mathbf{Q} \mathbf{Q}^{\top} = \mathbf{I}$ , we can define  $\mathbf{z} = \mathbf{Q}^{\top} \mathbf{x}$  and rewrite the Rayleigh quotient as:

$$R(\mathbf{A}, \mathbf{x}) = \frac{\mathbf{x}^{\top} \mathbf{Q} \Lambda \mathbf{Q}^{\top} \mathbf{x}}{\mathbf{x}^{\top} \mathbf{Q} \mathbf{Q}^{\top} \mathbf{x}} = \frac{\mathbf{z}^{\top} \Lambda \mathbf{z}}{\mathbf{z}^{\top} \mathbf{z}}. \quad (48)$$

As  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$  is a diagonal matrix, we have (writing  $\mathbf{z} = (z_1, \dots, z_d)^{\top}$ )

$$\mathbf{z}^{\top} \Lambda \mathbf{z} = \sum_{i=1}^d \lambda_i z_i^2. \quad (49)$$

Therefore the Rayleigh quotient can be written as the following weighted average of the eigenvalues

$$R(\mathbf{A}, \mathbf{x}) = \sum_{i=1}^d \frac{z_i^2}{\|\mathbf{z}\|_2^2} \lambda_i, \quad \text{with} \quad \sum_{i=1}^d \frac{z_i^2}{\|\mathbf{z}\|_2^2} = 1. \quad (50)$$

In summary, these derivations indicate that the Rayleigh quotient is bounded as

$$\begin{aligned} \lambda_{\min}(\mathbf{A}) &\leq R(\mathbf{A}, \mathbf{x}) \leq \lambda_{\max}(\mathbf{A}) \\ \Rightarrow \lambda_{\min}(\mathbf{A}) \|\mathbf{x}\|_2^2 &\leq \mathbf{x}^{\top} \mathbf{A} \mathbf{x} \leq \lambda_{\max}(\mathbf{A}) \|\mathbf{x}\|_2^2, \end{aligned} \quad (51)$$

where  $\lambda_{\min}(\mathbf{A})$  and  $\lambda_{\max}(\mathbf{A})$  are the smallest and largest eigenvalues of  $\mathbf{A}$ , respectively.

**Question 32** For linear regression we have  $\nabla_{\theta} L(\theta) = \mathbf{X}^{\top}(\mathbf{X}\theta_t - \mathbf{y})$ . Therefore in this case the pre-conditioned gradient descent update rule is

$$\theta_{t+1} = \theta_t - \frac{\gamma}{\sigma^2} \mathbf{P}^{-1} \mathbf{X}^{\top} (\mathbf{X}\theta_t - \mathbf{y}).$$

Under the formula of Arithmetico-geometric sequence, we have

$$\theta_t = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{P}^{-1} \mathbf{X}^{\top} \mathbf{X})^t (\theta_0 - \theta^*) + \theta^*, \quad \theta^* = (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{y}.$$

So this means if pre-conditioned gradient descent converges, it will converge to the right answer  $\theta^*$ .

Now we need to show that pre-conditioned gradient descent converges with the right choices of  $\gamma$  and  $\mathbf{P}$ . This requires us to analyse the eigenvalues of the matrix  $(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{P}^{-1} \mathbf{X}^{\top} \mathbf{X})^2$ , and for a given  $\mathbf{P}$ :

$$\lambda \text{ is an eigenvalue of } \mathbf{P}^{-1} \mathbf{X}^{\top} \mathbf{X} \Rightarrow (1 - \frac{\gamma}{\sigma^2})^2 \text{ is an eigenvalue of } (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{P}^{-1} \mathbf{X}^{\top} \mathbf{X})^2.$$

Following the same idea as to prove convergence for gradient descent, we have the learning rates bounds

$$\gamma_{min} = 2\sigma^2 / \lambda_{max}(\mathbf{P}^{-1} \mathbf{X}^{\top} \mathbf{X}), \quad \gamma_{max} = 2\sigma^2 / \lambda_{min}(\mathbf{P}^{-1} \mathbf{X}^{\top} \mathbf{X}).$$

Now to set  $\gamma_{min} = \gamma_{max}$ , we should choose  $\mathbf{P}$  such that  $\lambda_{min}(\mathbf{P}^{-1} \mathbf{X}^{\top} \mathbf{X}) = \lambda_{max}(\mathbf{P}^{-1} \mathbf{X}^{\top} \mathbf{X})$ , and an easy way to do so is to choose  $\mathbf{P} \propto \mathbf{X}^{\top} \mathbf{X}$ .

One can show that for linear regression problems, the Hessian matrix of  $L(\theta)$  is  $\nabla_{\theta}^2 L(\theta) \propto \mathbf{X}^{\top} \mathbf{X}$ . In general for a given loss function  $L(\theta)$  we often set  $\mathbf{P}_t = \nabla_{\theta_t}^2 L(\theta_t)$  if it can be computed in a fast way.

**Question 33** First note that  $\nabla_{\theta} L(\theta_t) = \frac{1}{\sigma^2} \mathbf{X}^{\top} (\mathbf{X}\theta_t - \mathbf{y})$ . The update equations for both the parameter and the momentum are

$$\begin{aligned} \theta_{t+1} &= \theta_t - \gamma \nabla_{\theta} L(\theta_t) + \alpha \Delta \theta_t = \theta_t - \frac{\gamma}{\sigma^2} \mathbf{X}^{\top} (\mathbf{X}\theta_t - \mathbf{y}) + \alpha \Delta \theta_t \\ \Delta \theta_{t+1} &= \theta_{t+1} - \theta_t = \alpha \Delta \theta_t - \frac{\gamma}{\sigma^2} \mathbf{X}^{\top} (\mathbf{X}\theta_t - \mathbf{y}). \end{aligned} \tag{52}$$

Now collecting both equations together into a “joint” linear equation:

$$\begin{bmatrix} \theta_{t+1} \\ \Delta \theta_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^{\top} \mathbf{X} & \alpha \mathbf{I} \\ -\frac{\gamma}{\sigma^2} \mathbf{X}^{\top} \mathbf{X} & \alpha \mathbf{I} \end{bmatrix} \begin{bmatrix} \theta_t \\ \Delta \theta_t \end{bmatrix} + \begin{bmatrix} \frac{\gamma}{\sigma^2} \mathbf{X}^{\top} \mathbf{y} \\ \frac{\gamma}{\sigma^2} \mathbf{X}^{\top} \mathbf{y} \end{bmatrix}. \tag{53}$$

Then we can apply the derivation of arithmetico-geometric sequences again, and show that

$$\begin{bmatrix} \theta_t \\ \Delta \theta_t \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^{\top} \mathbf{X} & \alpha \mathbf{I} \\ -\frac{\gamma}{\sigma^2} \mathbf{X}^{\top} \mathbf{X} & \alpha \mathbf{I} \end{bmatrix}^t \begin{bmatrix} \theta_0 - \theta^* \\ \Delta \theta_0 \end{bmatrix} + \begin{bmatrix} \theta^* \\ \mathbf{0} \end{bmatrix}, \tag{54}$$

with  $\theta^* = (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{y}$ . This equation also says if momentum GD converges, the momentum  $\Delta \theta_t$  will vanish to  $\mathbf{0}$ , which is as expected as  $\Delta \theta_t = \theta_t - \theta_{t-1} \rightarrow \mathbf{0}$ .