

Machine Learning

Summer Semester 2018, Homework 1

Prof. Dr. J. Peters, D. Tanneberg, B. Belousov



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Total points: 85 + 15 bonus

Due date: Wednesday, 02 May 2018 (before the lecture)

Name, Surname, ID Number

Zhiyuan Hu, 2863135, Yunlei Tang

Problem 1.1 Linear Algebra Refresher [20 Points]

a) Matrix Properties [5 Points]

A colleague of yours suggests matrix addition and multiplication are similar to scalars, thus commutative, distributive and associative properties can be applied. Is the statement correct? Prove it analytically or give counterexamples (for both operations) considering three matrices A, B, C of size $n \times n$.

For the matrix addition:

there is only commutative and commutative, no distributive. Both commutative and associative can be applied. To prove it, firstly, we define

$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & a_{ij} & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix} \quad B = \begin{pmatrix} b_{11} & \cdots & b_{1n} \\ \vdots & b_{ij} & \vdots \\ b_{n1} & \cdots & b_{nn} \end{pmatrix} \quad C = \begin{pmatrix} c_{11} & \cdots & c_{1n} \\ \vdots & c_{ij} & \vdots \\ c_{n1} & \cdots & c_{nn} \end{pmatrix}$$

commutative: $A + B = B + A$,

$$(A + B)_{ij} = a_{ij} + b_{ij}, \quad (B + A)_{ij} = b_{ij} + a_{ij} = a_{ij} + b_{ij} = (A + B)_{ij}$$

so $A + B = B + A$

associative: $(A + B) + C = A + (B + C)$,

$$((A + B) + C)_{ij} = (a_{ij} + b_{ij}) + c_{ij} = a_{ij} + (b_{ij} + c_{ij}) = (A + (B + C))_{ij}$$

so $(A + B) + C = A + (B + C)$

For the matrix multiplication:

the distributive and associative properties can be applied, but commutative not, to prove:
commutative:

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 4 & 3 & 2 \end{pmatrix} \quad B = \begin{pmatrix} 3 & 2 & 1 \\ 6 & 5 & 4 \\ 2 & 3 & 4 \end{pmatrix}$$

$$A * B = \begin{pmatrix} 21 & 21 & 21 \\ 54 & 51 & 48 \\ 34 & 29 & 24 \end{pmatrix} \quad B * A = \begin{pmatrix} 15 & 19 & 23 \\ 42 & 49 & 56 \\ 30 & 31 & 32 \end{pmatrix}$$

distributive:

$$(A + B) * C = A * C + B * C$$

$$((A + B) * C)_{ij} = \sum_{k=1}^n (a_{ik} + b_{ik}) * c_{kj} = \sum_{k=1}^n a_{ik} * c_{kj} + \sum_{k=1}^n b_{ik} * c_{kj} = (A * C + B * C)_{ij}$$

associative:

$$(A * B) * C = A * (B * C)$$

$$(A * B) * C = \sum_{j=1}^n (A * B)_{ij} * c_{jk} = \sum_{j=1}^n \sum_{l=1}^n a_{il} b_{lj} c_{jk} = \sum_{l=1}^n a_{il} \sum_{j=1}^n (b_{lj} * c_{jk}) = A * (B * C)$$

b) Matrix Inversion [6 Points]

Given the following matrix

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 1 & 2 & 4 \\ 1 & 4 & 5 \end{pmatrix}$$

analytically compute its inverse A^{-1} and illustrate the steps.

If we change the matrix in

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 1 & 2 & 4 \\ 1 & 2 & 5 \end{pmatrix}$$

is it still invertible? Why?

??

The inverse matrix can be calculate by the equation:

$$A^{-1} = \frac{A^*}{|A|}, \quad (A^*)_{ij}^T = (-1)^{i+j} * M_{ij}$$

A^* is the adjugate matrix, M_{ij} is minor of A , i.e. the determinant of the $(n-1) * (n-1)$ matrix that results from deleting row i and column j of A .

It will not be invertible for the new matrix, because the determinante $|A| = 0$, i.e. it does not have the full rank. And as to the pointed equation $A^* / |A|$ to calculate the inverse matrix, it's obvious not invertible.

$$A^{-1} = \frac{1}{|A|} \begin{pmatrix} 1 & 2 & 3 \\ 1 & 2 & 4 \\ 1 & 2 & 5 \end{pmatrix}$$

c) **Matrix Pseudoinverse [3 Points]**

Write the definition of the right and left Moore-Penrose pseudoinverse of a generic matrix $A \in \mathbb{R}^{n \times m}$.

Given $A \in \mathbb{R}^{2 \times 3}$, which one does exist? Write down the equation for computing it, specifying the dimensionality of the matrices in the intermediate steps.

d) **Eigenvectors & Eigenvalues [6 Points]**

What are eigenvectors and eigenvalues of a matrix A ? Briefly explain why they are important in Machine Learning.

Problem 1.2 Statistics Refresher [25 Points]

a) Expectation and Variance [8 Points]

Let Ω be a finite set and $P : \Omega \rightarrow \mathbb{R}$ a probability measure that (by definition) satisfies $P(\omega) \geq 0$ for all $\omega \in \Omega$ and $\sum_{\omega \in \Omega} P(\omega) = 1$. Let $f : \Omega \rightarrow \mathbb{R}$ be an arbitrary function on Ω .

1) Write the definition of expectation and variance of f and discuss if they are linear operators.

2) You are given a set of three dices $\{A, B, C\}$. The following table describes the outcome of six rollouts for these dices, where each column shows the outcome of the respective dice. (Note: assume the dices are standard six-sided dices with values between 1-6)

A	4	4	2	4	1	1
B	3	6	3	3	4	3
C	5	5	2	1	1	1

Estimate the expectation and the variance for each dice using unbiased estimators. (Show your computations).

3) According to the data, which of them is the “most rigged”? Why?

b) It is a Cold World [7 Points]

Consider the following three statements:

- a) A person with a cold has backpain 25% of the time.
- b) 4% of the world population has a cold.
- c) 15% of those who do not have a cold, still have backpain.

- 1) Identify random variables from the statements above and define a unique symbol for each of them.
- 2) Define the domain of each random variable.
- 3) Represent the three statements above with your random variables.
- 4) If you suffer from backpain, what are the chances that you suffer from a cold? (Show all the intermediate steps.)

c) Journey to THX1138 [10 Points]

After the success of the **Rosetta mission**, ESA decided to send a spaceship to rendezvous with the comet THX1138. This spacecraft consists of four independent subsystems A, B, C, D . Each subsystem has a probability of failing during the journey equal to $1/3$.

- 1) What is the probability of the spacecraft S to be in working condition (i.e., all subsystems are operational at the same time) at the rendezvous?
- 2) Given that the spacecraft S is not operating properly, compute analytically the probability that **only** subsystem A has failed.
- 3) Instead of computing the probability analytically, do a simple simulation experiment and compare the result to the previous solution. Include a snippet of your code.
- 4) An improved spacecraft version has been designed. The new spacecraft fails if the critical subsystem A fails, or any two subsystems of the remaining B, C, D fail. What is the probability that **only** subsystem A has failed, given that the spacecraft S is failing?

Problem 1.3 Optimization and Information Theory [40 Points + 15 Bonus]

a) Entropy [5 Points]

You work for a telecommunication company that uses a system to transmit four different symbols S_1, S_2, S_3, S_4 through time. In the current system, each symbol has a probability to occur according to the following table

	S_1	S_2	S_3	S_4
p_i	0.03	0.62	0.26	0.09

Compute the entropy of the system and write the minimum number of bits requires for transmission.

b) Constrained Optimization [25 Points]

After an upgrade of the system, your boss asks you to change the probabilities of transmission in order to maximize the entropy. However, the new system has the following constraint

$$4 = \sum_{i=1}^4 2p_i i.$$

- 1) Formulate it as a constrained optimization problem. Do you need to include additional constraints beside the one above?
- 2) Write down the Lagrangian of the problem. Use one Lagrangian multiplier per constraint.
- 3) Compute the partial derivatives of the Lagrangian above for each multiplier and the objective variable. Is it easy to solve it analytically?
- 4) Formulate the dual function of this constrained optimization problem. Solve it analytically.
- 5) Name one technique for numerically solve these problems and briefly describe it.

c) Numerical Optimization [10 Points]

Rosenbrock's function (to be minimized) is defined as

$$f(\mathbf{x}) = \sum_{i=1}^{n-1} \left[100(x_{i+1} - x_i^2)^2 + (x_i - 1)^2 \right].$$

Write in Python a simple gradient descent algorithm and simulate it for 10,000 steps on Rosenbrock's function with $n = 20$. Attach a snippet of your algorithm, discuss the effects of the learning rate and attach a plot of your learning curve with your best learning rate.

d) **Natural Gradient [10 Bonus Points]**

Let $\theta \in \mathbb{R}^n$ be a parameter vector and $J: \mathbb{R}^n \rightarrow \mathbb{R}$ a cost function. The negative gradient $-\nabla J(\theta)$ is sometimes called the *steepest descent direction*. But is it really? To be able to claim that it is *the* steepest descent direction, we should compare it to other descent directions and pinpoint what is so unique about the negative gradient direction.

Covariant gradient. A fair way to compare descent directions is to make a small step of fixed length, say ε , in every direction $\Delta\theta$ and check which direction leads to the greatest decrease in $J(\theta)$. Since we assume that the step size is small, we can evaluate the decrease in $J(\theta)$ using its first-order Taylor approximation

$$J(\theta + \Delta\theta) - J(\theta) \approx \nabla J(\theta)^T \Delta\theta.$$

To make precise what we mean by ‘small’ step size, we need to introduce a norm (or a distance) in the space of parameters θ . A good choice, that among other advantages captures the intuition that some parameters may influence the objective function more than others, is the generic quadratic norm

$$\|\Delta\theta\|^2 = \frac{1}{2} \Delta\theta^T F(\theta) \Delta\theta$$

with a positive-definite matrix $F(\theta)$; note that in general F may depend on θ .

1) Find the direction $\Delta\theta$ that yields the largest decrease in the linear approximation of $J(\theta)$ for a fixed step size ε . Does this direction coincide with $-\nabla J(\theta)$? The direction that you found is known as the negative covariant gradient direction.

Natural gradient. In statistical models, parameter vector θ often contains parameters of a probability density function $p(x; \theta)$ (for example, mean and covariance of a Gaussian density); thus, the cost function J depends on θ indirectly through $p(x; \theta)$. This two-level structure gives a strong hint as to what matrix F to pick for measuring the distance in the parameter space in the most ‘natural’ way. Namely, one can carry over the notion of ‘distance’ between probability distributions $p(x; \theta + \Delta\theta)$ and $p(x; \theta)$ (which is known from information theory to be well captured by the Kullback-Leibler divergence) to the distance between the corresponding parameter vectors $\theta + \Delta\theta$ and θ .

2) Obtain the quadratic Taylor approximation of the KL divergence from $p(x; \theta)$ to $p(x; \theta + \Delta\theta)$ in the form

$$KL(p(x; \theta + \Delta\theta) || p(x; \theta)) \approx \frac{1}{2} \Delta\theta^T F(\theta) \Delta\theta.$$

Covariant gradient with the matrix $F(\theta)$ that you found is known as the natural gradient.

e) **Gradient Descent Variants [5 Bonus Points]**

Throughout this class we have seen that gradient descent is one of the most used optimization techniques in Machine Learning. This question asks you to deepen the topic by conducting some research by yourself.

1) There are several variants of gradient descent, namely *batch*, *stochastic* and *mini-batch*. Each variant differs in how much data we use to compute the gradient of the objective function. Discuss the differences among them, pointing out pros and cons of each one.

2) Many gradient descent optimization algorithms use the so-called *momentum* to improve convergence. What is it? Is it always useful?