

Ecorrection

Place student sticker here

Note:

- During the attendance check a sticker containing a unique code will be put on this exam.
- This code contains a unique number that associates this exam with your registration number.
- This number is printed both next to the code and to the signature field in the attendance check list.

Machine Learning

Exam: IN2064 / Endterm

Date: Thursday 13th February, 2020

Examiner: Prof. Dr. Stephan Günnemann

Time: 17:00 – 19:00

	P 1	P 2	P 3	P 4	P 5	P 6	P 7	P 8	P 9	P 10
I										

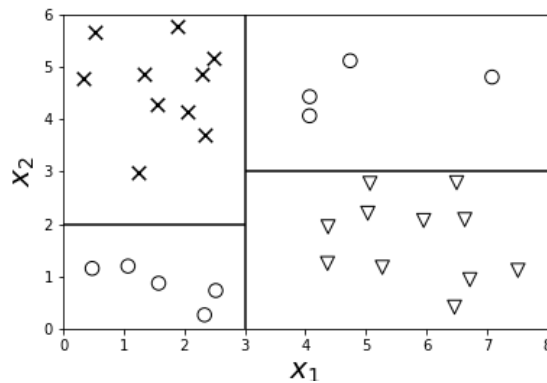
Working instructions

- This exam consists of **16 pages** with a total of **10 problems**.
Please make sure now that you received a complete copy of the exam.
- The total amount of achievable credits in this exam is 92 credits.
- Detaching pages from the exam is prohibited.
- Allowed resources:
 - A4 sheet of handwritten notes (two sides)
 - **no other materials (e.g. books, cell phones, calculators) are allowed!**
- Only write on the sheets given to you by supervisors. If you need more paper, ask the supervisors.
- Last two pages can be used as scratch paper.
- All sheets (including scratch paper) have to be returned at the end.
- **Only use a black or a blue pen (no pencils, red or green pens)!**
- Write your answers only in the provided solution boxes or the scratch paper.
- **For problems that say "Justify your answer" you only get points if you provide a valid explanation.**
- **For problems that say "Prove" you only get points if you provide a valid mathematical proof.**
- If a problem does not say "Justify your answer" or "Prove" it's sufficient to only provide the correct answer.
- Exam duration - 120 minutes.

Left room from _____ to _____ / Early submission at _____

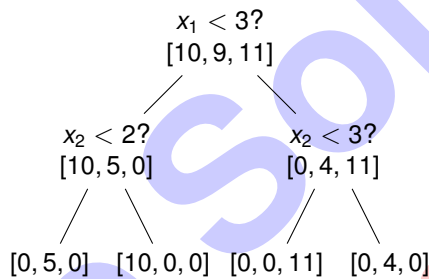
Problem 1 Decision Trees (12 credits)

You are given a dataset with points from three different classes and want to classify them based on a decision tree. The plot below illustrates the data points (class labels are indicated by the symbols \times , \circ , ∇) and the decision boundaries of a decision tree.



a) Draw the corresponding decision tree. Make sure that you include the feature (x_1 or x_2) and threshold of the split as well as the number of samples of each class that pass the corresponding inner node or leaf node.

0	
1	
2	
3	
4	



Each inner node of the tree must include: feature and threshold, number of samples of each class (value = [...]).

Each leaf node must include: number of samples of each class (value = [...]).

Point: root ✓ , inner nodes ✓ ✓ , leave node ✓

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3
<input type="checkbox"/>	4

b) Compute the Gini index of each node of your decision tree.

Note: Your answer may contain improper fractions (e.g. $\frac{33}{117}$)

$$\text{Gini index: } i_G(t) = \sum_{i \in C} \pi_i(1 - \pi_i) = 1 - \sum_{i \in C} \pi_i^2$$

$$\text{Root node: } i_G(t) = 1 - \left(\frac{10}{30}\right)^2 - \left(\frac{9}{30}\right)^2 - \left(\frac{11}{30}\right)^2 = \frac{598}{900} \approx 0.664 \quad \checkmark$$

$$\text{Left child of root: } i_G(t) = 1 - \left(\frac{10}{15}\right)^2 - \left(\frac{5}{15}\right)^2 - \left(\frac{0}{15}\right)^2 = \frac{100}{225} \approx 0.444 \quad \checkmark$$

$$\text{Right child of root: } i_G(t) = 1 - \left(\frac{0}{15}\right)^2 - \left(\frac{4}{15}\right)^2 - \left(\frac{11}{15}\right)^2 = \frac{88}{225} \approx 0.391 \quad \checkmark$$

$$\text{Left leaf node: } i_G(t) = 1 - \left(\frac{0}{5}\right)^2 - \left(\frac{5}{5}\right)^2 - \left(\frac{0}{5}\right)^2 = 0$$

$$\text{Left-middle leaf node: } i_G(t) = 1 - \left(\frac{0}{5}\right)^2 - \left(\frac{5}{5}\right)^2 - \left(\frac{0}{5}\right)^2 = 0$$

$$\text{Right-middle leaf node: } i_G(t) = 1 - \left(\frac{10}{10}\right)^2 - \left(\frac{0}{10}\right)^2 - \left(\frac{0}{10}\right)^2 = 0$$

$$\text{Right leaf node: } i_G(t) = 1 - \left(\frac{0}{4}\right)^2 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2 = 0 \text{ (all leaf nodes)} \quad \checkmark$$

c) Assume you have a dataset with two-dimensional points from two different classes C_1 and C_2 . The points from class C_1 are given by $A = \{(i, i^2) \mid i \in \{1 \dots 100\}\} \subseteq \mathbb{R}^2$, while the points from class C_2 are $B = \{(i, \frac{125}{i}) \mid i \in \{1 \dots 100\}\} \subseteq \mathbb{R}^2$.

Construct a decision tree of minimal depth that assigns as many data points as possible to the correct class. Provide for each split the feature and corresponding thresholds. How many and which datapoints are misclassified?

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3
<input type="checkbox"/>	4

- Split root node based on feature $1 \leq 5 \quad \checkmark$
- Split both child nodes based on feature $2 \leq 25 \quad \checkmark$

One point (5, 25) is in both classes and misclassified. \checkmark

Decision tree / Calculation of thresholds / Picture. \checkmark

Problem 2 Probabilistic inference (4 credits)

0
1
2
3
4

We are interested in estimating a discrete parameter z that can take values in $\{1, 2, 3, 4\}$.

- We place a categorical prior on z , that is $p(z \mid \pi) = \text{Categorical}(z \mid \pi)$ with $\pi = (0.1, 0.05, 0.85, 0.0)$.
- We choose the following likelihood function: $p(x \mid z) = \text{Exponential}(x \mid 2^z) = 2^z \exp(-x2^z)$.
- We have observed one sample $x = 32$.

What is the posterior probability that z is equal to 4, i.e. what is $p(z = 4 \mid x, \pi)$? Justify your answer.

Using the Bayes formula ✓ (for correctly writing the Bayes formula)

$$\begin{aligned} p(z = 4 \mid x, \pi) &\propto p(x \mid z = 4)p(z = 4 \mid \pi) \\ &\propto 2^4 \exp(-32 \cdot 2^4) \cdot 0 \\ &\propto 0 \end{aligned}$$

Since the prior probability $p(z = 4 \mid \pi)$ equals to zero, the posterior probability $p(z = 4 \mid x, \pi)$ equals to zero as well ✓✓✓.

minus ✓ for each mistake in the formulas, even if the final answer is correct

Problem 3 Probabilistic inference (8 credits)

0
1
2
3
4
5
6
7
8

We are interested in estimating the parameter $\theta \in \mathbb{R}$ of the following probabilistic model:

$$p(x \mid \theta) = \exp(\theta - x - \exp(\theta - x)).$$

We have observed a single sample $x \in \mathbb{R}$ drawn from the above model. Derive the maximum likelihood estimate (MLE) of the parameter θ . Justify your answer.

The maximum likelihood estimate of θ is defined as

$$\begin{aligned} \theta_{MLE} &= \arg \max_{\theta} p(x \mid \theta) \checkmark \checkmark \\ &= \arg \max_{\theta} \log p(x \mid \theta) \checkmark \\ &= \arg \min_{\theta} -\log p(x \mid \theta) \\ &= \arg \min_{\theta} \left(-\theta + x + \frac{\exp(\theta)}{\exp(x)} \right). \end{aligned}$$

Clearly, this is a convex function of θ . To minimize, compute the derivative ✓, set it to zero ✓ and solve for θ .

$$\begin{aligned} \frac{\partial}{\partial \theta} \left(-\theta + x + \frac{\exp(\theta)}{\exp(x)} \right) &= -1 + \frac{\exp(\theta)}{\exp(x)} \stackrel{!}{=} 0 \\ \Leftrightarrow \frac{\exp(\theta)}{\exp(x)} &\stackrel{!}{=} 1 \checkmark \checkmark \text{ for correctly computing and simplifying} \\ \Leftrightarrow \theta_{MLE} &= x \checkmark \text{ for the correct final answer} \end{aligned}$$

Therefore, $\theta_{MLE} = x$.

Mark correct answers with a cross



To undo a cross, completely fill out the answer option

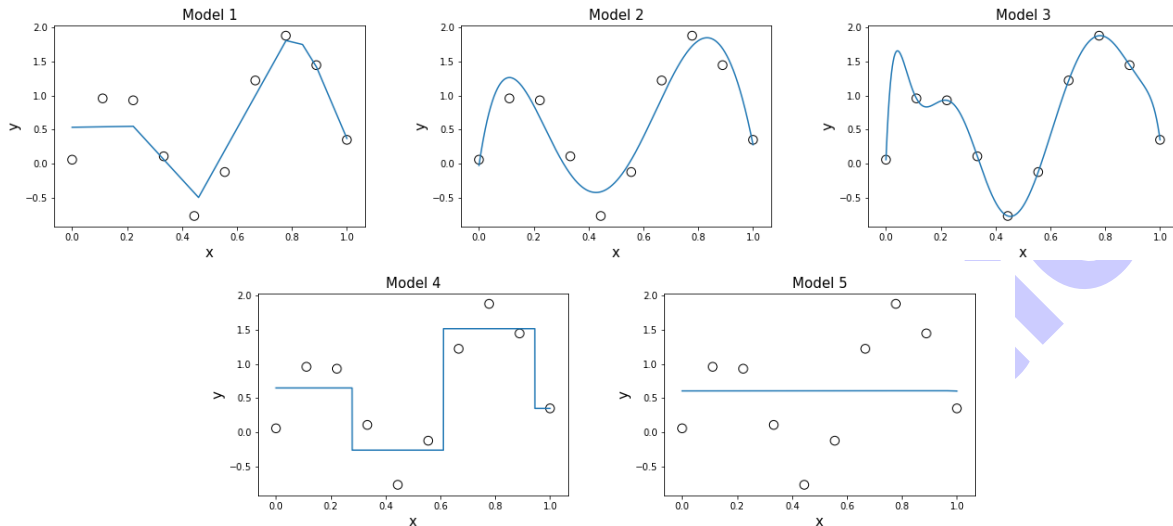


To re-mark an option, use a human-readable marking



Problem 4 Regression (10 credits)

The following five plots show five different regression models fitted to the same dataset. Your task is to assign each of the plots to the corresponding model.



For each subtask ✓✓ if only the correct option is chosen, otherwise 0.

a) Polynomial regression (degree = 5), no regularization

☐ Model 1

☒ Model 2

☐ Model 3

☐ Model 4

☐ Model 5

b) Polynomial regression (degree = 10), no regularization

☐ Model 1

☐ Model 2

☒ Model 3

☐ Model 4

☐ Model 5

c) Polynomial regression (degree = 50), L_2 regularization with $\lambda = 10^3$

☐ Model 1

☐ Model 2

☐ Model 3

☐ Model 4

☒ Model 5

d) Feed-forward neural network with ReLU activation functions, no regularization

☒ Model 1

☐ Model 2

☐ Model 3

☐ Model 4

☐ Model 5

e) Decision tree of depth 2

☐ Model 1

☐ Model 2

☐ Model 3

☒ Model 4

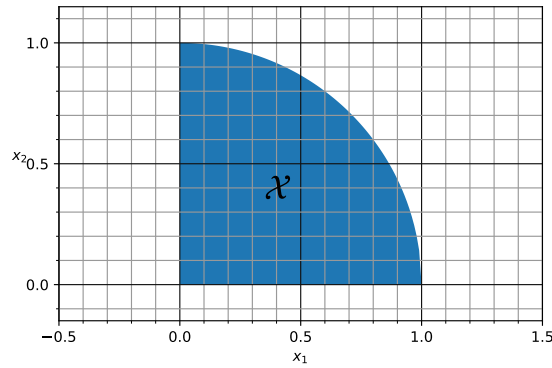
☐ Model 5

Problem 5 Convex optimization (18 credits)

Consider the set $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^D : \|\mathbf{x}\|_2 \leq 1 \text{ and } x_i \geq 0 \text{ for all } i = 1, \dots, D\}$ with $1 < D \in \mathbb{N}$.

0
1

a) Draw \mathcal{X} on the provided axes below for $D = 2$.



0
1
2
3
4

b) Write down the function $\pi_{\mathcal{X}}$ projecting an arbitrary point $\mathbf{p} \in \mathbb{R}^2$ on \mathcal{X} for the case $D = 2$.

Note: if you decide to split \mathbb{R}^2 into regions and consider them separately then you have to describe the regions analytically (just a reference to your plot from a) will not be sufficient).

We consider the following disjoint decomposition of $\mathbb{R}^2 = \mathcal{X} \cup \mathcal{X}_0 \cup \mathcal{X}_1$ with

- $\mathcal{X}_0 = \{\mathbf{x} \in \mathbb{R}^2 : x_1 < 0 \text{ or } x_2 < 0\}$, here $\pi_{\mathcal{X}}$ is the same as the projection on $[0, 1]^2$, that is

$$\pi_{\mathcal{X}}(\mathbf{x}) = \begin{pmatrix} \min(1, \max(0, x_1)) \\ \min(1, \max(0, x_2)) \end{pmatrix} \text{ for } \mathbf{x} \in \mathcal{X}_0,$$

- $\mathcal{X}_1 = \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\| > 1, x_1 \geq 0 \text{ and } x_2 \geq 0\}$, here $\pi_{\mathcal{X}}$ is the the projection on $\mathcal{B}_1(\mathbf{0})$, therefore

$$\pi_{\mathcal{X}}(\mathbf{x}) = \frac{\mathbf{x}}{\|\mathbf{x}\|_2} \text{ for } \mathbf{x} \in \mathcal{X}_1,$$

✓ ✓ ✓ if the projection $\pi_{\mathcal{X}}$ is correct for the points $\mathbf{x} \notin \mathcal{X}$.

- and finally $\pi_{\mathcal{X}}(\mathbf{x}) = \mathbf{x}$ for $\mathbf{x} \in \mathcal{X}$.

✓ for the case $\mathbf{x} \in \mathcal{X}$.

From now on we consider the setting with an arbitrary $1 < D \in \mathbb{N}$.

c) Prove that \mathcal{X} is convex.

0
1
2
3

$$\begin{aligned}\mathcal{X} &= \{\mathbf{x} \in \mathbb{R}^D : \|\mathbf{x}\|_2 \leq 1, x_i \geq 0 \text{ for all } i = 1, \dots, D\} \\ &= \underbrace{\{\mathbf{x} \in \mathbb{R}^D : \|\mathbf{x}\|_2 \leq 1\}}_{\text{unit ball } \mathcal{B}_1(\mathbf{0}) \checkmark} \cap \underbrace{\{\mathbf{x} \in \mathbb{R}^D : 0 \leq x_i \leq 1 \text{ for all } i = 1, \dots, D\}}_{\text{cube } [0,1]^D \checkmark}\end{aligned}$$

Therefore, \mathcal{X} is an intersection of two convex sets and hence convex. \checkmark

If proven by definition of convexity:

$\checkmark \checkmark$ for proving that an intermediate point \mathbf{x}_λ satisfies $\|\mathbf{x}_\lambda\|_2 \leq 1$,

\checkmark for proving that $(\mathbf{x}_\lambda)_i \geq 0$ for all i .

d) Fill in the space in the box below using mathematical notation with a description of the vertices of \mathcal{X} . Note that just writing down the definition of $\text{vert}(\mathcal{X})$ will not be sufficient.

0
1
2

$\text{vert}(\mathcal{X}) =$

$$\{\mathbf{x} \in \mathbb{R}^D : \|\mathbf{x}\|_2 = 1 \text{ and } x_i \geq 0 \text{ for all } i = 1, \dots, D\} \checkmark \cup \{\mathbf{0}\} \checkmark$$



e) Find the maximum of the following constrained optimization problem. Justify your answer, all properties of the objective function and \mathcal{X} that you use should be clearly stated and derived from the previous tasks or results considered in the course.

Hint: results from c) and d) might help you.

Hint: for arbitrary $\mathbf{c} \in \mathbb{R}^D$ the maximum of the constrained problem $\max_{\mathbf{x}} \mathbf{c}^T \mathbf{x}$ subject to $\|\mathbf{x}\|_2 = 1$ is $\|\mathbf{c}\|_2$.

$$\begin{aligned} & \text{maximize}_{\mathbf{x}} \quad \sum_{i=1}^D x_i + e^{\|\mathbf{x}\|_2^2} \\ & \text{subject to} \quad \mathbf{x} \in \mathcal{X} \end{aligned}$$

We will utilize the following results:

- \mathcal{X} is convex (from c).
- Function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ with $f(\mathbf{x}) = \sum_{i=1}^D x_i + e^{\|\mathbf{x}\|_2^2}$ is convex since the function $h : \mathbb{R}^D \rightarrow [0, \infty)$ with $h(\mathbf{x}) = \|\mathbf{x}\|_2^2$ is convex and the function $g : [0, \infty) \rightarrow \mathbb{R}$ with $g(x) = e^x$ is convex and non-decreasing. Therefore, using the convexity preserving operations we see that $g(h(\mathbf{x})) = e^{\|\mathbf{x}\|_2^2}$ and linear $l(\mathbf{x}) = \sum_{i=1}^D x_i$ are both convex and hence $f(\mathbf{x}) = l(\mathbf{x}) + g(h(\mathbf{x}))$ is convex as well.
- Maximum of a convex function over a convex domain is attained at one of its vertices (from the lecture).

✓✓✓✓ for a correct argumentation that it is enough to consider the optimization problem on $\text{vert}(\mathcal{X})$ only.

Now it is sufficient to look for the maximum value of f over $\text{vert}(\mathcal{X})$ that consists of $\{\mathbf{0}\}$ and a set where for each point it holds that $\|\mathbf{x}\|_2 = 1$ and $x_i \geq 0$ for all i (from d).

Case 0: $\mathbf{x} = \mathbf{0}$, then $f(\mathbf{x}) = f(\mathbf{0}) = 0 + e^0 = 1$.

Case 1: $\mathbf{x} \in B := \{\mathbf{x} \in \mathbb{R}^D : \|\mathbf{x}\|_2 = 1 \text{ and } x_i \geq 0 \text{ for all } i = 1, \dots, D\}$ then the largest value that f attains at these points is

$$\max \{f(\mathbf{x}) : \|\mathbf{x}\|_2 = 1, x_i \geq 0\} = e + \max \{\mathbf{c}^T \mathbf{x} : \|\mathbf{x}\|_2 = 1\} \text{ where } \mathbf{c} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}.$$

Using the provided hint we get that the maximum of f over B is $e + \|\mathbf{c}\|_2 = e + \sqrt{D}$. Since $\sqrt{D} + e > 1$ we get that the maximum of f over \mathcal{X} is $\sqrt{D} + e$.

✓✓✓✓ if all vertices were considered and the correct results were obtained using valid argumentation.

Problem 6 Kernels (10 credits)

Let $\mathbf{A} \in \mathbb{R}^{D \times D}$ be a positive semi-definite matrix and consider for $p \in \mathbb{N}$ the following function

$$k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}, \quad k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^T \mathbf{A} \mathbf{x}_2 + 1)^p.$$

a) Prove that k is a valid kernel using kernel preserving operations known from the course.

Solution 1:

- ✓ $k_1(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{A} \mathbf{x}_2$ is a kernel since \mathbf{A} is positive semi-definite (kernel is known from the course).
- ✓ $k_2(\mathbf{x}_1, \mathbf{x}_2) = 1$ is a kernel since $k_2(\mathbf{x}_1, \mathbf{x}_2) = \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2)$ for $\phi(\mathbf{x}) = 1$.
- ✓ $k_3(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{A} \mathbf{x}_2 + 1$ is a kernel since it is a sum of two kernels (kernel preserving operation known from the lecture)
- ✓ Finally, subsequently applying the rule that a product of two kernels is a kernel for $k_3 * k_3$ we get that k is a kernel as well.

Solution 2:

- ✓ $k_1(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{A} \mathbf{x}_2$ is a kernel since \mathbf{A} is positive semi-definite. Therefore there exists a map ϕ_A such that $k_1(\mathbf{x}_1, \mathbf{x}_2) = \phi_A(\mathbf{x}_1)^T \phi_A(\mathbf{x}_2)$.
- Polynomial kernel $k_p(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^T \mathbf{x}_2 + 1)^p$ is a kernel. ✓ if this fact is used later in the solution
- ✓ ✓ Using the composition rule (kernel preserving operation known from the lecture) we know that $k_p(\phi_A(\mathbf{x}_1), \phi_A(\mathbf{x}_2)) = (\mathbf{x}_1^T \mathbf{A} \mathbf{x}_2 + 1)^p = k(\mathbf{x}_1, \mathbf{x}_2)$ is again a kernel.

b) For the special case of $p = 2$ and $\mathbf{A} = \mathbf{I}$ (identity matrix) write down the corresponding feature map $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$ such that

$$k(\mathbf{x}_1, \mathbf{x}_2) = \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2).$$

What is the dimension M of the feature space in this case?

$$\phi(\mathbf{x}) = \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \vdots \\ \sqrt{2}x_D \\ x_1x_1 \\ x_1x_2 \\ \vdots \\ x_1x_D \\ \vdots \\ x_2x_1 \\ \vdots \\ x_Dx_D \end{pmatrix} \quad \text{and} \quad M = 1 + D + D^2.$$

Comment: in this case we get the quadratic kernel and the feature map includes constant, linear and all quadratic terms we can produce from the initial features.

✓ for the constant term, ✓ ✓ for the linear terms, ✓ ✓ ✓ for the quadratic terms, -✓ if M is wrong, -✓ if only the case of $D = 2$ is considered.

Problem 7 Deep learning (8 credits)

The code snippet below shows an implementation of two functions $f : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$.

Given two input vectors $\mathbf{x} \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^N$, we perform the following computations:

$$z = f(\mathbf{x}, \mathbf{y})$$
$$t = g(z)$$

The code below uses backpropagation to compute $\frac{\partial t}{\partial \mathbf{x}}$ and $\frac{\partial t}{\partial \mathbf{y}}$ (similarly to how we did it in Tutorial 9: Deep Learning I). However, some code fragments are missing. Your task is to complete the missing code fragments.

Note: It's also fine to write your answer using pseudocode (we won't deduct points for small Python syntax errors, etc.).

```
class F:
    def forward(self, x, y):
        self.cache = (x, y)
        #####
        # MISSING CODE FRAGMENT #1
        #####
        return out

    def backward(self, d_out):
        # x, y are np.arrays of shape [N]
        x, y = self.cache
        N = len(x)
        # np.ones(N) returns a vector of ones of shape [N]
        d_y = (x + np.ones(N)) * d_out
        d_x = (y - np.ones(N)) * d_out
        return d_x, d_y

def sigmoid(a):
    return 1 / (1 + np.exp(-a))

class G:
    def forward(self, z):
        self.cache = z
        return sigmoid(z)

    def backward(self, d_out):
        z = self.cache
        #####
        # MISSING CODE FRAGMENT #2
        #####
        return d_z

# Example usage
f = F()
g = G()
x = np.array([1., 2., 3])
y = np.array([-2., 3., -1.])

z = f.forward(x, y)
t = g.forward(z)
d_z = g.backward(1.0)
d_x, d_y = f.backward(d_z)
```

a) Complete the MISSING CODE FRAGMENT #1

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3
<input type="checkbox"/>	4

Option 1:

```
out = np.dot(x, y) + np.sum(y) - np.sum(x)
```

Option 2:

```
N = len(x)
ones = np.ones(N)
out = np.dot(x, y) - np.dot(x, ones) + np.dot(y, ones)
```

Option 3:

```
N = len(x)
ones = np.ones(N)
out = np.dot(x + ones, y - ones)
```

Two points for $x^T y$ ✓✓, one point for each of the sums ✓✓

b) Complete the MISSING CODE FRAGMENT #2

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3
<input type="checkbox"/>	4

Option 1:

```
d_z = sigmoid(z) * sigmoid(-z) * d_out
```

Option 2:

```
d_z = sigmoid(z) * (1 - sigmoid(z)) * d_out
```

Three points for the derivative of sigmoid ✓✓✓, one point for not forgetting d_out ✓

Problem 8 SVD and linear regression (8 credits)

0 ☐
1 ☐
2 ☐
3 ☐
4 ☐
5 ☐
6 ☐
7 ☐
8 ☐

You want to perform linear regression on a data set with features $\mathbf{X} \in \mathbb{R}^{N \times D}$ and targets $\mathbf{y} \in \mathbb{R}^N$. Assume that you have already computed the SVD of the feature matrix $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. Additionally, assume that \mathbf{X} has full rank.

Show how we can compute the optimal linear regression weights \mathbf{w}^* in $\mathcal{O}(ND^2)$ operations by using the result of the SVD.

Hint: Matrix operations have the following asymptotic complexity

- Matrix multiplication \mathbf{AB} for arbitrary $\mathbf{A} \in \mathbb{R}^{P \times Q}$ and $\mathbf{B} \in \mathbb{R}^{Q \times R}$ takes $\mathcal{O}(PQR)$
- Matrix multiplication \mathbf{AD} for an arbitrary $\mathbf{A} \in \mathbb{R}^{P \times Q}$ and a diagonal $\mathbf{D} \in \mathbb{R}^{Q \times Q}$ takes $\mathcal{O}(PQ)$
- Matrix inversion \mathbf{C}^{-1} for an arbitrary matrix $\mathbf{C} \in \mathbb{R}^{M \times M}$ takes $\mathcal{O}(M^3)$
- Matrix inversion \mathbf{D}^{-1} for a diagonal matrix $\mathbf{D} \in \mathbb{R}^{M \times M}$ takes $\mathcal{O}(M)$

$$\begin{aligned}\mathbf{w}^* &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \\ &= ((\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T))^{-1} (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T \mathbf{y} = \\ &= (\mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^{-1} \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} = \\ &= (\mathbf{V}\mathbf{\Sigma}^2 \mathbf{V}^T)^{-1} \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} = \\ &= (\mathbf{V}^T)^{-1} \mathbf{\Sigma}^{-2} \mathbf{V}^T \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} = \\ &= \mathbf{V}\mathbf{\Sigma}^{-2} \mathbf{V}^T \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} = \\ &= \mathbf{V}\mathbf{\Sigma}^{-1} \mathbf{U}^T \mathbf{y}\end{aligned}$$

Multiplication $\mathbf{a} = \mathbf{U}^T \mathbf{y}$ takes $\mathcal{O}(N \cdot D \cdot 1)$

Multiplication $\mathbf{b} = \mathbf{\Sigma}^{-1} \mathbf{a}$ takes $\mathcal{O}(D)$

Multiplication $\mathbf{w} = \mathbf{Vb}$ takes $\mathcal{O}(D^2)$

In total, $\mathcal{O}(ND + D + D^2) = \mathcal{O}(ND)$ if $N > D$.

✓ for $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$,

✓ for $\mathbf{U}^T \mathbf{U} = \mathbf{I}$

✓ for $(\mathbf{V}\mathbf{\Sigma}^2 \mathbf{V}^T)^{-1} = (\mathbf{V}^T)^{-1} \mathbf{\Sigma}^{-2} \mathbf{V}^{-1}$

✓ for $(\mathbf{V}^T)^{-1} \mathbf{\Sigma}^{-2} \mathbf{V}^{-1} = \mathbf{V}\mathbf{\Sigma}^{-2} \mathbf{V}^T$

✓ for $\mathbf{V}^T \mathbf{V} = \mathbf{I}$

✓ for the final answer $\mathbf{w}^* = \mathbf{V}\mathbf{\Sigma}^{-1} \mathbf{U}^T \mathbf{y}$

minus ✓ for any mistake

✓ ✓ for correct big-O analysis

There is a mistake in the task, saying " $\mathcal{O}(ND^2)$ " instead of " $\mathcal{O}(ND)$ ". The task was graded taking this into account.

Problem 9 K-Means (10 credits)

Let $\gamma_i \in \mathbb{R}^D$ for $i = 1, \dots, K$ be a set of K points more than 4 apart, i.e. $\|\gamma_i - \gamma_j\|_2 > 4$ for all $i \neq j$. Consider K non-empty datasets \mathcal{X}_i each contained within a unit ball around γ_i , i.e. $\|\mathbf{x} - \gamma_i\|_2 \leq 1$ for all $\mathbf{x} \in \mathcal{X}_i$. Let $\mathcal{X} = \bigcup_{i=1}^K \mathcal{X}_i$ be the combined dataset.

Now consider a centroid initialization procedure similar to k-means++, though it deterministically chooses the data point farthest away from all previous centroids. That means it initializes the cluster centers μ_i as

$$\mu_i = \begin{cases} \text{random sample from } \mathcal{X} & \text{if } i = 1 \\ \arg \max_{\mathbf{x} \in \mathcal{X}} \min_{j \in \{1, \dots, i-1\}} \|\mathbf{x} - \mu_j\|_2 & \text{if } i \in \{2, \dots, K\} \end{cases}$$

a) Explain why this deterministic k-means++ initialization of K clusters assigns each μ_i to a different ball, i.e. for $i \neq j$ such that $\mu_i \in \mathcal{X}_{i'}$ and $\mu_j \in \mathcal{X}_{j'}$ it holds that $i' \neq j'$.

0
1
2
3
4
5

The first centroid is placed into a random ball. Now assume that a subsequent centroid μ_j would be placed into the same ball as some previous μ_i . That means that $\|\mu_j - \mu_i\|_2 \leq 2$ because each $\mathcal{X}_{i'}$ is contained within a ball of radius 1. However, because we are placing K centroids into K balls and one ball has been chosen twice, there is at least one ball $\mathcal{X}_{i'}$ that does not have a centroid in it. $\mathcal{X}_{i'}$ is non-empty, so it has an element $\mathbf{x} \in \mathcal{X}_{i'}$ which is more than 2 away from any previously chosen centroid because the ball centers are more than 4 apart and the data points can deviate at most 1 from their closest ball center. But that contradicts the construction of μ_j as the data point that is the furthest away from any previously chosen centroid, in particular μ_i . Therefore, every centroid is placed in a different ball.

✓✓✓ if the student shows an understanding that the statement is true because the balls are far apart in some way

✓✓ if the student works out the crux of the argument clearly, i.e. ≤ 2 vs. > 2 or in words "at most 2" vs. "more than 2"

-✓ for mistakes, erroneous statements or arguments that do not make explicit use of ≤ 2 vs. > 2 and thus would work in the setting $\|\gamma_i - \gamma_j\|_2 > 3.9$ just the same

b) Assuming a), explain why k-means clustering of \mathcal{X} with K clusters and our deterministic k-means++ initialization recovers the underlying structure of the data, i.e. all data points $\mathbf{x} \in \mathcal{X}_i$ will be assigned to the same centroid for all i .

0
1
2
3
4
5

Without loss of generality, let the centroids be assigned such that $\mu_i \in \mathcal{X}_i$ for all i . Then each data point $\mathbf{x} \in \mathcal{X}_i$ will be assigned to μ_i because $\|\mu_i - \mathbf{x}\|_2 \leq 2$ and $\|\mu_j - \mathbf{x}\|_2 > 2$ for all i and $j \neq i$, see a). Because updating centroid μ_i is a convex combination of data points $\mathbf{x} \in \mathcal{X}_i$, μ_i stays within the bounding ball of \mathcal{X}_i for any i . So the assignments stay the same and k-means terminates in the next iteration.

✓✓✓✓ for arguing convincingly that each \mathbf{x} is assigned to the centroid in its ball because again ≤ 2 vs. > 2

✓ for stating that and why the algorithm converges/terminates

-✓ for mistakes, erroneous statements or arguments that do not make explicit use of ≤ 2 vs. > 2 and thus would work in the setting $\|\gamma_i - \gamma_j\|_2 > 3.9$ just the same

Problem 10 Differential Privacy & Fairness (4 credits)

0 ☐
1 ☐
2 ☐

a) What is the *robustness to post-processing* property of Differential Privacy?

You can apply any function on the output from a ϵ -DP mechanism and the new output remains ϵ -DP
✓ as long as you do not touch again the data. ✓
Only one point if "do not touch again the data" is missing.

0 ☐
1 ☐
2 ☐

b) Suppose that we require that the same *percentage* of applicants from two different groups must receive a loan. What fairness criterion are we implementing? What is the biggest downside/con of this criterion?

Demographic Parity, also called Independence and Statistical Parity. ✓

Any of the following downsides are acceptable:

- Rules out the perfect predictor when base rates are different across groups. ✓
- Laziness: We can trivially satisfy the criterion if we give loan to qualified people from one group and random people from the other. ✓
- Too strong, if one group is much more likely to repay compared to the other group. ✓

Additional space for solutions—clearly mark the (sub)problem your answers are related to and strike out invalid solutions.

Sample Solution

Correction Notes

Sample Solution

Correction Notes