

#### Final Exam

### Introduction to Machine Learning for Bioengineers BIO-322

1

# Firstname Lastname 123456

## February 2, 2023

- The exam lasts 180 minutes.
- Write all your answers in English in a legible way on the exam (no extra sheets).
- Use a dark (e.g. black or blue) pen or pencil.
- Use the last 4 pages of this exam as scratch space.
- 1 page handwritten notes is allowed (A4 one side).
- The handwritten notes must be yours; copies from other students are not allowed.
- No calculator or other electronic device is allowed.
- Put your bag including computer and cell phone to the indicated places.
- Have your student card displayed before you on your desk.
- Do not write your name or sciper number on any page.
- Check that your exam has 14 pages.

Respectez les consignes suivantes   Observe this guidelines   Beachten Sie bitte die unten stehenden Richtlinien								
choisir une réponse   select an answer Antwort auswählen	ne PAS choisir une réponse   NOT select an answer NICHT Antwort auswählen	Corriger une réponse   Correct an answer Antwort korrigieren						
ce qu'il ne faut <u>PAS</u> faire   what should <u>NOT</u> be done   was man <u>NICHT</u> tun sollte								



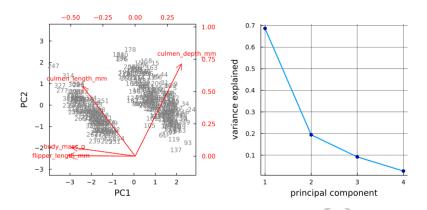
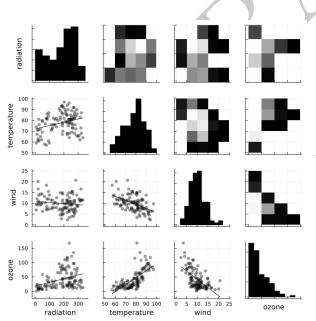


Figure 1: Principal Component Analysis of the Penguin Dataset.



The data contains measurements of solar radiation, temperature, wind speed, and cube root ozone concentration on 111 days at sites in the New York metropolitan region. Objective: predict ozone level.

- temperature in degree Fahrenheit.
- $\bullet$  wind speed in miles per hour
- ullet solar radiation in langleys
- cube root ozone level in ppb $\frac{1}{3}$ .

Figure 2: The ozone data set.



For each of the following questions there is one correct answer and you have two possibilities: tick one of the boxes or none. Every correct answer gives 2 positive point, every wrong answer 1 negative point, and no answer no point.

features: I numerical 1 categorical Question 1 (We want to predict the height of 10 weeks old corn seedlings, based on the average temperature (in  $^{\circ}$ C), the soil condition (dry or humid) and the fertilizer condition (none, biological, chemical). After one-hot coding relative to a standard level, the number of parameters of a linear regression (with intercept) is

> + 1 3 4 5 5

Suppose you did logistic regression and you get a test set  $((x_1 = 1, y_1 = A), (x_2 =$ Question 2  $-3, y_2 = A$ ,  $(x_3 = 4, y_3 = B), (x_4 = -1, y_4 = B)$ ). The predictions of your fitted machine on this test set are  $P(y_1 = A|x_1) = 0.7, P(y_2 = A|x_2) = 0.7, P(y_3 = A|x_3) = 0.6, P(y_4 = A|x_4) = 0.6.$ The AUC on this test set is

is test set is

for thresholds between

0 0.5 0.6 1 0.6 and 0.7 the
true positive rake is 1 and
the false positive
then neural networks with strong L2 regularization tend to have a lower bias but a rale
than neural networks without positive in the positive rate.

Question 3 higher variance than neural networks without regularization. reduces Variance

Wrong. Correct. and increases boas.

Cross validation can be used to select the number of hidden layers in a neural Question 4 network.

☐ Wrong. ☐ Correct.

Question 5 Consider the data generating process y = A if 2x - 1 > 1 and y = B otherwise. The irreducible error of this generator is larger than 1. This data generator x Wrong.  $\Box$  Correct. determinable  $\Rightarrow$  irreducible

With the Lasso with a large penalty  $\lambda$ , many parameters are exactly 0. Question 7

Wrong. Correct.

Assume the data set  $((x_1 = 1, y_1 = 10), (x_2 = 3, y_2 = 0), (x_3 = 4, y_3 = 0),$ Question 8  $(x_4 = 8, y_4 = 0))$  should be fitted with a regression tree. In the first step of recursive binary splitting, the split is introduced at

X = 2

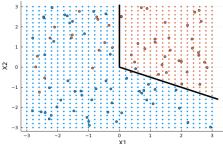
The result of K-Means clustering applied to un-standardized data is the same as Question 9 the result of K-Means clustering applied to standardized data.

Wrong. Correct.

Assume all points of a dataset with three predictors lie on a line. In this case, Question 10 the proportion of variance explained (PVE) is zero for the second and third principal component.

Wrong. Correct.

Question 11



The black decision boundary on the left could have been obtained by fitting the following method on the training data (shown as large blue and red dots).

a neural network classifier with one hidden layer

2 nearest neighbours classification

logistic regression none of these methods

The largest variance in the penguin dataset in figure 1 on page 2 is parallel to Question 12 the flipper length mm axis.

 $\oint_{\Lambda} \approx \left(-0.55, -0.5, -0.4, 0.8\right)$ Wrong.  $\Box$ Correct.

The task described in figure 2 on page 2 is a problem of type: Question 13

For a given data generating process and some training data, we find for a quadratic fit  $f_1(x) = 1 + 2x - 0.5x^2$  a training error of 0.5 and for a linear fit  $f_2(x) = 0.5 + 1.8x$  a training error of 1.5. This implies that the test loss of  $f_1$  for the given data generating process is equal or lower than the test loss of  $f_2$ .

for could be over fitting the training date.

| Wrong. | Correct.

Question 15 Early stopping in gradient descent can be used to find the value of the regularization constant  $\lambda$  for the Lasso.

Wrong.

Correct.

The data points in the penguin dataset in figure 1 on page 2 lie all very close (in Question 16 the sense of more than 95% of the dataset's variability lies within this subspace) to a

straight line 2D plane 3D linear subspace none of the other options

In a classification task with 3 classes we know that the data generating process has probabilities  $P(y=A|x=1)=P(y=B|x=1)=\frac{1}{4}$  and  $P(y=C|x=1)=\frac{1}{2}$ . How large is the test log-likelihood at  $x_0 = 1$  for a linear model with fitted parameters  $\hat{\theta}_{10} = \hat{\theta}_{11} = \hat{\theta}_{20} = \hat{\theta}_{21} = \hat{\theta}_{30} = \hat{\theta}_{31} = 0$ ?  $\frac{1}{4} \log \frac{1}{3} + \frac{1}{4} \log \frac{1}{3} = \log \frac{1}{3} = \log 3$   $\frac{1}{3} \log(\frac{1}{3})$   $\log(3)$  0

The task described in figure 2 on page 2 is a problem of type: Question 18

reinforcement learning

supervised learning unsupervised learning



2 | 3 | 4 (for grading)

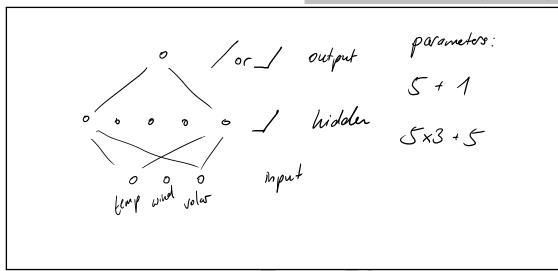
## Open Questions (56 Points)

Please write your answers to the following questions in the designated boxes. Do not tick the checkboxes for grading.

Question 19 Draw a sketch of a neural network with one hidden layer of 5 relu neurons to fit the data in figure 2 on page 2. Label the layers as "input", "hidden" and "output" layer and

0 1

indicate the activation function of each layer.



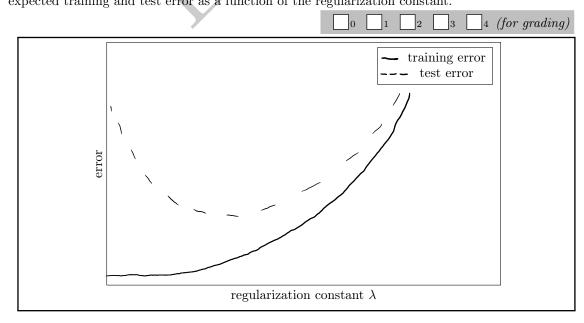
Question 20 How many parameters (including biases) does the neural network of the previous

question have?

0	1	2	3	4	(for	grading

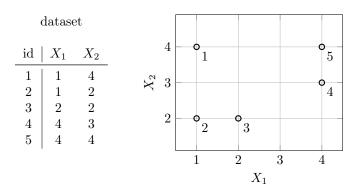
$$5 \times 3 + 5 + 5 + 1 = 26$$

Question 21 Suppose you use L2 regularization in linear regression. Draw a sketch of the expected training and test error as a function of the regularization constant.



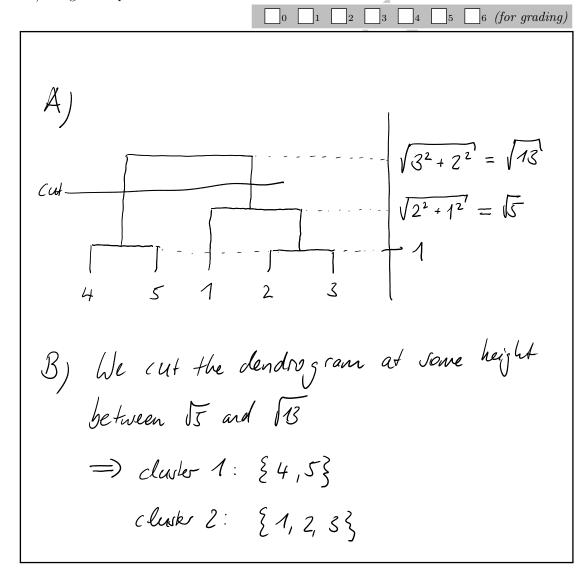


#### Question 22

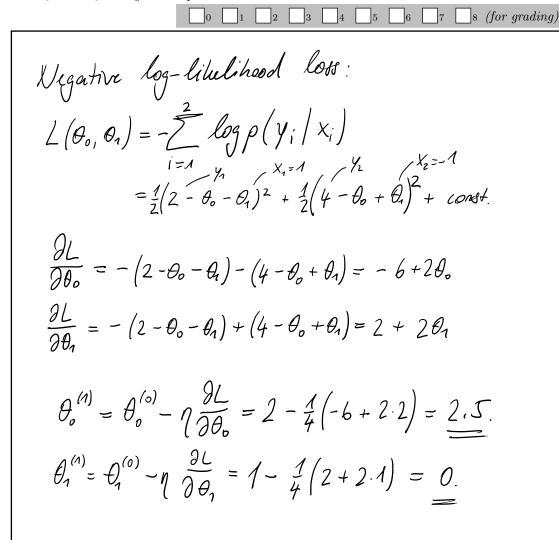


With hierarchical clustering we want to assign each data point in the above dataset to one of two clusters. In the figure the observation identity is given by the number next to the data point. We use the Euclidean distance and complete linkage.

- A) Determine the dendrogram and compute the heights at which branches merge.
- B) Assign each point to one of two clusters.



Question 23 You are given the data set  $((x_1 = 1, y_1 = 2), (x_2 = -1, y_2 = 4))$ . Assume the data comes from a generator with normally distributed (Gaussian) noise and you would like to fit this data with the family of probability densities  $p(y|x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y-\theta_0-\theta_1x)^2}{2}\right)$ . Compute one update-step of gradient descent on the negative log-likelihood loss with learning rate  $\eta = \frac{1}{4}$ , assuming initial guess  $\theta_0^{(0)} = 2, \theta_1^{(0)} = 1$ .

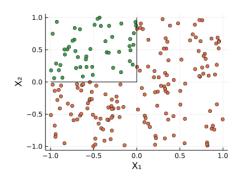


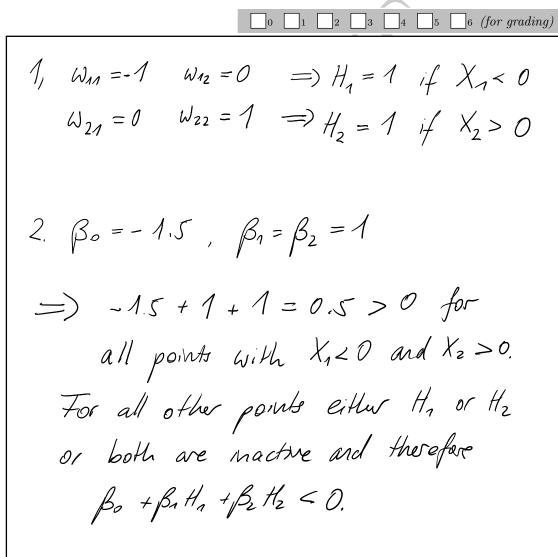


#### Question 24

On the right you see a decision problem with non-linear decision boundary (black line).

- 1. Determine a feature representation with 2 features  $H_1$  = heaviside( $w_{11}X_1 + w_{12}X_2$ ),  $H_2$  = heaviside( $w_{21}X_1 + w_{22}X_2$ ) such that the decision boundary is linear in the feature representation. Give your answer in the form  $w_{ij} = \dots$  for all i = 1, 2, j = 1, 2. The heaviside function is given by heaviside(x) = 1 if x > 0, heaviside(x) = 0 otherwise.
- 2. Show that the decision boundary is indeed linear, i.e. show that there are regression coefficients  $\beta_0, \beta_1, \beta_2$  such that  $\beta_0 + \beta_1 H_1 + \beta_2 H_2 > 0$  for all points with  $X_1 < 0$  and  $X_2 > 0$  and  $\beta_0 + \beta_1 H_1 + \beta_2 H_2 < 0$  for all other points.





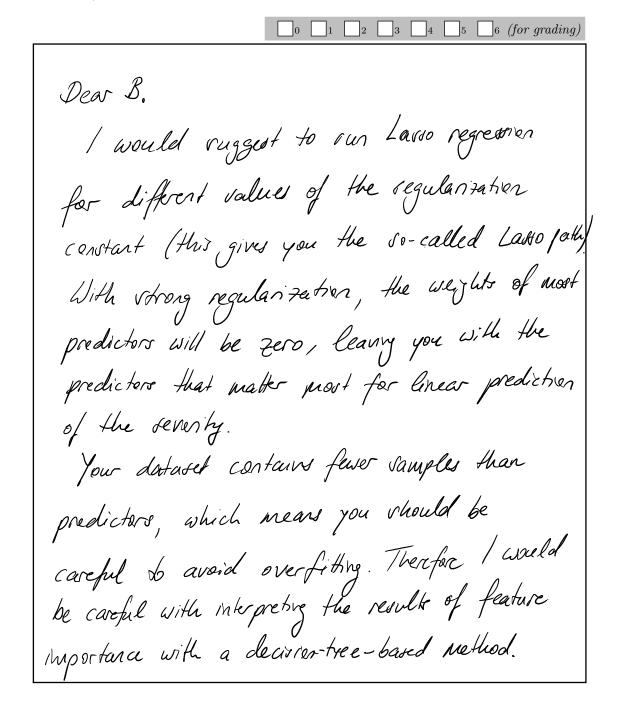


Question 25 Suppose you receive the following email from a colleague.

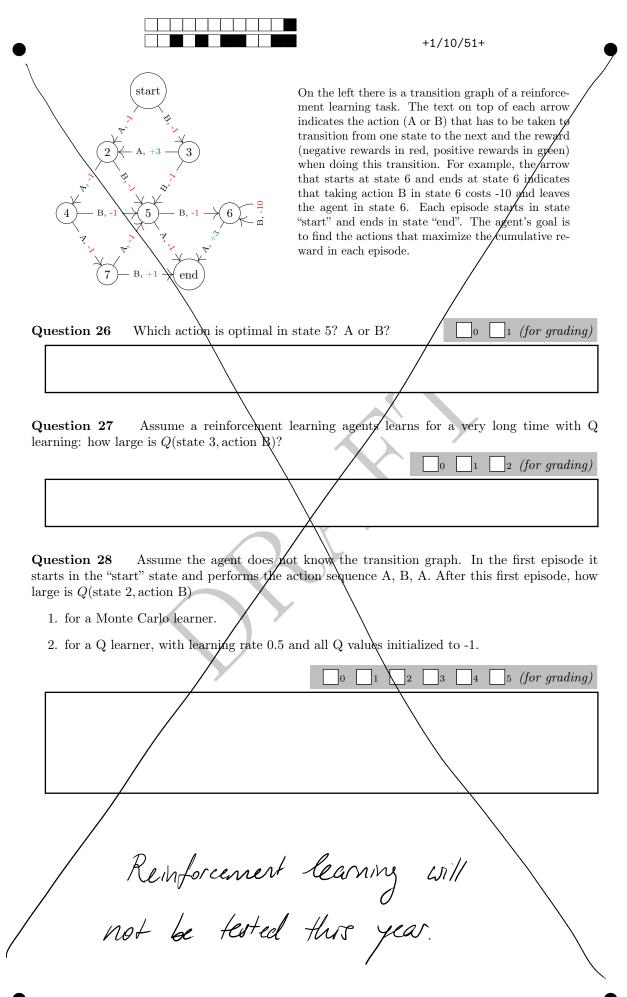
Dear Firstname,

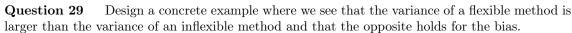
For my current project I work with a database of 71 patients with different severity of a certain disease and 112 measured variables from the analysis of the blood etc. I would like to use machine learning to determine which of the 112 measured variables are unlikely to be related with the severity of the disease. Can you suggest one machine learning method you would use for this task? Please explain, how and why you would use this method.

Cheers, B.



9



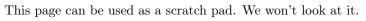


Here is what we mean by a "concrete example": we should see with numbers what you want to show. For example, if the task were to "Show with a concrete example that predictions under linear regression are invariant under scaling of the training input data but predictions with ridge regression are not" a valid answer would be:

Assume data set  $((x_1=0,y_1=1),(x_2=1,y_1=2))$ . For linear regression the solution would be  $\theta_0=1$  and  $\theta_1=(2-1)/(1-0)=1$ . Therefore, predictions for  $x_0=2$  would be  $\theta_1x_0+\theta_0=3$ . If we scale the input data by s=2, the solution would be  $\tilde{\theta}_0=1$  and  $\tilde{\theta}_1=(2-1)/(2-0)=\frac{1}{2}$ . Therefore, predictions for  $\tilde{x}_0=2\cdot 2=4$  would be  $\tilde{\theta}_1\tilde{x}_0+\tilde{\theta}_0=3$ , the same as before. On the other hand, for ridge regression with  $\lambda=1$  the solution for the unscaled case would be... etc.

8 (for grading) Assume a data generator produces datasets with  $\{(x_1 = 0, y_1 = \varepsilon_1), (x_2 = 1, y_2 = 2 + \varepsilon_2)\}$ where En and Ez have mean O and variance 1. Ul use k-nearest neighors regression with U=1 for a flexible method and k=2 for au inflexible method (in this context). For h= 1 the bias at Xo = 0  $(f(x) - E_f(\hat{f}(x_0)))^2 = (0 - 0)^2 = 0$  and  $Var f(x_0) = Var(\epsilon_1) = 1$ For k=2 the bias at Xo=0  $(0-\frac{f(0)+f(1)}{2})^2=1$  and  $Var \int (X_0) = Var \left(\frac{\mathcal{E}_1 + \mathcal{E}_2}{2}\right) = \frac{1}{4} Var \left(\mathcal{E}_1\right) + \frac{1}{4} Var \left(\mathcal{E}_2\right) = \frac{1}{2}$ 

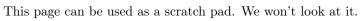




If you use this page for results that we should look at write this explicitly and clearly.

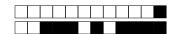






This page can be used as a scratch pad. We won't look at it. If you use this page for results that we should look at **write this explicitly and clearly**.





This page can be used as a scratch pad. We won't look at it. If you use this page for results that we should look at **write this explicitly and clearly**.

