
Machine Learning A

2025-2026

Home Assignment 1, TA Version: Tasks only

Christian Igel and Yevgeny Seldin

Department of Computer Science
University of Copenhagen

The deadline for this assignment is **1 September 2025, 22:00**. You must submit your *individual* solution electronically via the Absalon home page.

A solution consists of:

- A PDF file with detailed answers to the questions, which may include graphs and tables if needed. Do *not* include your full source code in the PDF file, only selected lines if you are asked to do so.
- A .zip file with all your solution source code with comments about the major steps involved in each question (see below). Source code must be submitted in the original file format, not as PDF. The programming language of the course is Python.
- **IMPORTANT: Do NOT zip the PDF file**, since zipped files cannot be opened in *SpeedGrader*. Zipped PDF submissions will not be graded.
- Your PDF report should be self-sufficient. I.e., it should be possible to grade it without opening the .zip file. We do not guarantee opening the .zip file when grading.
- Your code should be structured such that there is one main file (or one main file per question) that we can run to reproduce all the results presented in your report. This main file can, if you like, call other files with functions, classes, etc.
- Handwritten solutions will not be accepted. Please use the provided LaTeX template to write your report.

1 Make Your Own (10 points)

Imagine that you would like to write a learning algorithm that would predict the final grade of a student in the Machine Learning course based on their profile, for example, their grades in prior courses, their study program, etc. Such an algorithm would have been extremely useful: we could save significant time on grading and predict the final grade when the student just signs up for the course. We expect that the students would also appreciate such service and avoid all the worries about their grades. Anyhow, if you were to make such an algorithm,

1. What profile information would you collect and what would be the sample space \mathcal{X} ?
2. What would be the label space \mathcal{Y} ?
3. How would you define the loss function $\ell(y', y)$?
4. Assuming that you want to apply K -Nearest-Neighbors, how would you define the distance measure $d(x, x')$?
5. How would you evaluate the performance of your algorithm? (In terms of the loss function you have defined earlier.)
6. Assuming that you have achieved excellent performance and decided to deploy the algorithm, would you expect any issues coming up? How could you alleviate them?

There is no single right answer to the question. The main purpose is to help you digest the definitions we are working with. Your answer should be short, no more than 2-3 sentences for each bullet point. For example, it is sufficient to mention 2-3 items for the profile information, you should not make a page-long list.

2 Digits Classification with K Nearest Neighbors (40 points)

Detailed Instructions We pursue several goals in this question:

- Get your hands on implementation of K -NN and practice vector operations in Python.
- Explore fluctuations of the validation error as a function of the size of a validation set. (**Task#1**)
- Explore the impact of data corruption on the optimal value of K . (**Task#2**, optional)

IMPORTANT: Please, remember to include axis labels, legends and appropriate titles in your plots!

Task #1 In order to explore fluctuations of the validation error as a function of the size of the validation set, we use the following construction:

- Implement a Python function `knn(training_points, training_labels, test_points, test_labels)` that takes as input a $d \times m$ matrix of training points `training_points`, where m is the number of training points and d is the dimension of each point ($d = 784$ in the case of digits), a vector `training_labels` of the corresponding m training labels, a $d \times n$ matrix `test_points` of n test points, and their labels `test_labels` (you will need to convert the labels from $\{5, 6\}$ to $\{-1, 1\}$). The function should return a vector of length m , where each element represents the average error of K -NN on the test points for the corresponding value of K for $K \in \{1, \dots, m\}$. **Include a printout of your implementation of the function in the report.** (Only this function, not all of your code, the complete code should be included in the `.zip` file.) Ideally, the function should have no for-loops, check the practical advice at the end of the question.
- Use the first m digits for training the K -NN model. Take $m = 50$.
- Consider five validation sets, where for $i \in \{1, \dots, 5\}$ the set i consists of digits $m + (i - 1) \times n + 1, \dots, m + i \times n$, and where n is the size of each of the five validation sets (we will specify n in a moment). The data split is visualized below.

Training data (m points)	Validation set #1 (n points)	Validation set #2 (n points)	Validation set #3 (n points)	Validation set #4 (n points)	Validation set #5 (n points)
--------------------------------	------------------------------------	------------------------------------	------------------------------------	------------------------------------	------------------------------------

- Calculate the validation error for each of the sets as a function of K , for $K \in \{1, \dots, m\}$. Plot the validation error for each of the five validation sets as a function of K in the same figure (you will get five lines in the figure).
- Execute the experiment above with $n \in \{10, 20, 40, 80\}$. You will get four figures for the four values of n , with five lines in each figure. **Include these four figures in your report.**
- Create a figure where for each $n \in \{10, 20, 40, 80\}$ you plot the variance of the validation error over the five validation sets, as a function of K . You will get four lines in this figure, one for each n . **Include this figure in your report.** (Clarification, in case you got confused: fix n and K , then you have five numbers corresponding to validation errors on the five validation sets. You should compute the variance of these five values. Now keep n fixed, take $K \in \{1, \dots, m\}$, and compute the variance as a function of K , i.e., compute it for each K separately. This gives you one line. And then each $n \in \{10, 20, 40, 80\}$ gives you a line, so you get four lines.)
- What can you say about fluctuations of the validation error as a function of n ? **Answer in the report.**
- What can you say about the prediction accuracy of K -NN as a function of K ? **Answer in the report.**
- A high-level comment: a more common way of visualizing variation of outcomes of experiment repetitions is to plot the mean and error bars, but this form of visualization makes it too easy for humans to ignore the error bars and concentrate just on the mean, see the excellent book of Kahneman (2011). The visualization you are asked to provide in this question makes it hard to ignore the variation.

3 Regression (50 points)

Consider the data $S = \{(x_1, y_1), \dots, (x_N, y_N)\}$ in the file `PCB.dt`¹, which contains the concentrations of polychlorinated biphenyl (PCB) residues in lake trouts from the Cayuga Lake, NY, as reported by Bache et al. (1972). “The ages of the fish were accurately known, because the fish are annually stocked as yearlings and distinctly marked as to year class. Each whole fish was mechanically chopped, ground, and thoroughly mixed, and 5-gram samples taken. The samples were treated and PCB residues in parts per million (ppm) were estimated using column chromatography” (Bates and Watts, 1988).

Each line in `PCB.dt` is one training pattern. The first number is the input (x), the age of the fish in years, and the second is the corresponding output / target / label (y), the PCB concentration (in ppm).

Tasks

1. (5 points) Implement the linear regression algorithm as described in the lecture. For vector and matrix operations, such as computing the inverse of a matrix, you can use high-level (library) functions (e.g., from NumPy). Show the implementation of the linear regression algorithm in your report (you need not show boilerplate code, loading of the data, etc.).
2. (10 points) The task is to build a non-linear model $h : \mathbb{R} \rightarrow \mathbb{R}$ of the data in `PCB.dt`. The model should be of the form

$$h(x) = \exp(ax + b) \tag{1}$$

with parameters $a, b \in \mathbb{R}$.

The parameters can be learned using linear regression in the following way. Before applying linear regression, transform the output data (the y values) by applying the natural logarithm.

Then build an affine linear model using the transformed targets. By doing so, you effectively learn the (non-linear) model (1).

That is, you have to perform the following steps:

- Construct the data set $S' = \{(x_1, \ln y_1), \dots, (x_N, \ln y_N)\}$.
- Fit a model $h'(x) = ax + b$ to S' .
- Obtain the final model as $h(x) = \exp(h'(x))$

Build the model, report the two parameters of the model as well as the mean-squared-error of the model h computed over the training data set S .

¹Download from <https://github.com/christian-igel/ML/tree/main/data>

3. (10 points) In this particular example, linear regression without the non-linear transformation gives a lower error. Do not get confused by this: In the above procedure we train for a low error “on logarithmic scale”, and minimizing this error may not minimize the error on the original scale, which is the error you should consider in this exercise. The idea is that the data looks more linear on logarithmic scale, this is why a model of the form (1) (a standard *allometric* equation, Huxley and Teissier, 1936) has been considered in literature.

Provide an example showing that

$$\arg \min_{a,b} \sum_{i=1}^N (y_i - h(x_i))^2 \neq \arg \min_{a,b} \sum_{i=1}^N (\ln y_i - (ax_i + b))^2 . \quad (2)$$

Note two important things: First, the transformation is non-linear. Second, we are considering a transformation of the output space that changes the error/objective function – the learning goal. In task 6, we change the input representation and compare approaches using the same objective function.

Provide an instructive example with two data points. An easy way is to assume that the output is independent of the input and a data set $S = \{(x_1, y_1), (x_2, y_2)\}$ with $x_1 = x_2$ and $y_1 \neq y_2$.

4. (5 points) Plot the data and the model output. The plot must have proper axis labels and a legend. In *all* plots in this assignment, plot the logarithm of the PCB concentration versus the (not transformed) age.
5. (10 points) Compute the coefficient of determination R^2 as

$$1 - \frac{\sum_{i=1}^N (y_i - h(x_i))^2}{\sum_{i=1}^N (y_i - \bar{y})^2} , \quad (3)$$

where \bar{y} denotes the mean of the training labels. Discuss this quantity. What does it mean if R^2 is 1 and especially if R^2 is 0? Can R^2 be negative?

6. (10 points) Now let us build a non-linear model

$$h(x) = \exp(a\sqrt{x} + b) \quad (4)$$

with $a, b \in \mathbb{R}$. This is the same as before plus additionally applying the non-linear input transformation $x \mapsto \sqrt{x}$ to the input data.

One can view this as the inputs x being mapped to a feature space \mathcal{Z} by the *feature map* $\phi(x) = \sqrt{x}$.

Build the model and report the mean-squared-error. Then plot the target (on logarithmic scale) and the model output *over the original inputs* (linear scale). That is, the unit of the x -axis should be years.

Compute R^2 for the new model and the transformed labels. Discuss the result in comparison to the previous model.

Deliverables: Source code (main parts also in the report); plot of data and model output; mean-squared error, parameters of regression model, discussion of R^2 ; counterexample showing the difference when fitting in non-linearly transformed output space; mean-squared error of model with transformed inputs, plot of data and the model given by the second model (note axis scaling in years), comparison of R^2 values of the two different models

References

- C. A. Bache, J. W. Serum, W. D. Youngs, and D. J. Lisk. Polychlorinated biphenyl residues: Accumulation in cayuga lake trout with age. *Science*, 177(4055):1191–1192, 1972.
- D. M. Bates and D. G. Watts. *Nonlinear Regression Analysis and Its Applications*, volume 2. Wiley New York, 1988.
- J. S. Huxley and G. Teissier. Terminology of relative growth. *Nature*, 137(3471):780–781, 1936.
- D. Kahneman. *Thinking, fast and slow*. Farrar, Straus and Giroux, New York, 2011.

Appendix: Python Tips

Before you begin, ensure you have the following installed:

- Python (recommended version 3.x)
- NumPy library
- Matplotlib library

Loading and Visualizing MNIST-5-6-Subset Data

```
import numpy as np
import matplotlib.pyplot as plt
```

```

# Load the data from MNIST-5-6-Subset.txt
# Change the path as needed
data_file_path = "MNIST-5-6-Subset/MNIST-5-6-Subset.txt"
data_matrix = np.loadtxt(data_file_path).reshape(1877, 784)

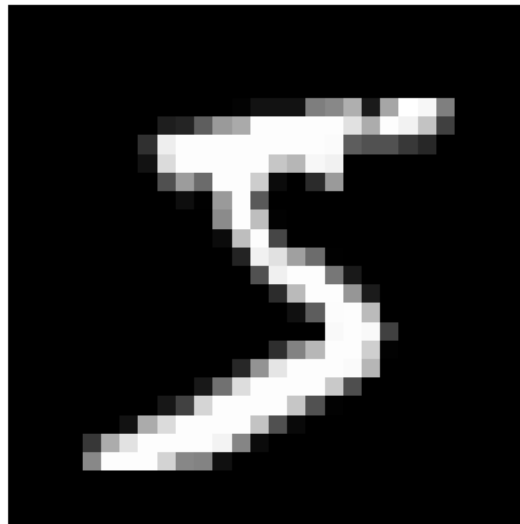
# Load the labels from MNIST-5-6-Labels.txt
# Change the path as needed
labels_file_path = "MNIST-5-6-Subset/MNIST-5-6-Subset-Labels.txt"
labels = np.loadtxt(labels_file_path)

# Assuming you want to visualize the first image
# Change the index as needed
image_index = 0
image_data = data_matrix[image_index]
selected_label = int(labels[image_index])

# Visualize the image using Matplotlib
# We transpose the image to make the number look upright.
plt.imshow(image_data.reshape(28,28).transpose(1,0), cmap='gray')
plt.title(f"Label: {selected_label}")
plt.axis('off') # Turn off axis
plt.show()

```

Label: 5



Setting up a figure with axis labels, legend and title

```
# Dummy data, x and y
x = np.arange(0, 20.1, 0.1)
y = np.sin(x) + np.random.normal(0, 0.2, len(x))
some_parameter = 54

# Initialise figure (fig) and axis (ax)
fig, ax = plt.subplots(figsize=(8,5))

# Plot in axis, add label to data
ax.plot(x, y, label='Dummy data') # (*)

# Set labels and title
ax.set_xlabel('X axis')
ax.set_ylabel('Y axis')
ax.set_title(f'Dummy data with some parameter = {some_parameter}')

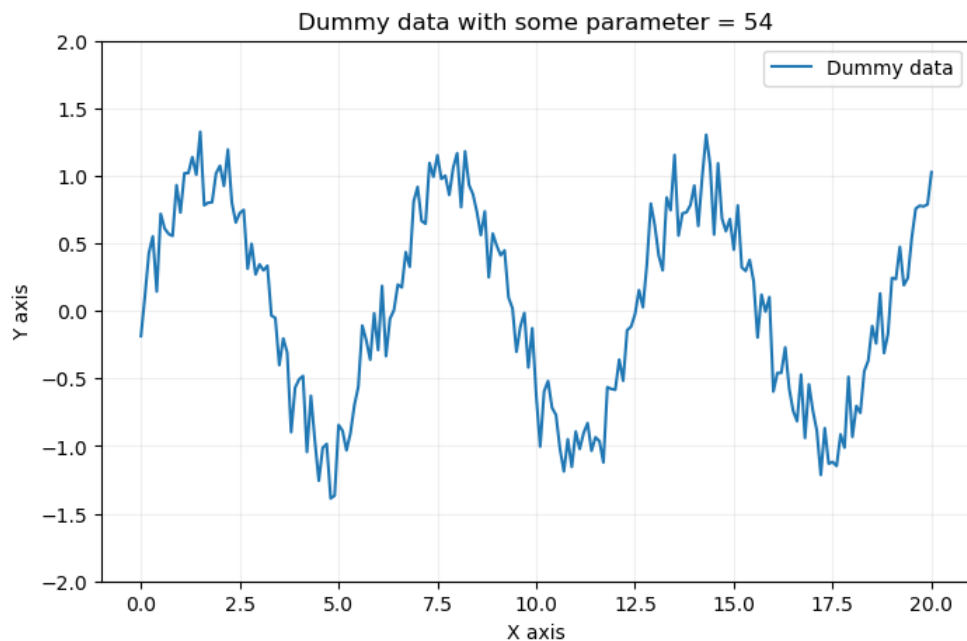
# Add grid
ax.grid(alpha=0.2)

# Set axes limits
ax.set_ylim(-2,2)

# Add legend (remember to label the data as shown above (*))
ax.legend()

# Show plot
plt.show()

# Save plot to some local path
fig.savefig('validation_err.png')
```



Other useful Numpy functions: cumsum, sort and argsort

```
# Creating an example array
data = np.array([5, 2, 8, 1, 6])

# 1)
# Calculating cumulative sum using cumsum
cumulative_sum = np.cumsum(data)
print("Original data:", data)
print("Cumulative sum:", cumulative_sum)
# Documentation for np.cumsum: https://numpy.org/doc/stable/reference/generated/numpy.cumsum.html

# 2)
# Sorting the array using sort
sorted_data = np.sort(data)
print("\nOriginal data:", data)
print("Sorted data:", sorted_data)
# Documentation for np.sort: https://numpy.org/doc/stable/reference/generated/numpy.sort.html
```

```

# 3)
# Getting indices that would sort the array using argsort
sorted_indices = np.argsort(data)
print("\nOriginal data:", data)
print("Sorted indices:", sorted_indices)
# Documentation for np.argsort: https://numpy.org/doc/stable/reference/generated/numpy.argsort.html

# 4)
# Accessing elements in sorted order using sorted indices
sorted_data_using_indices = data[sorted_indices]
print("\nOriginal data:", data)
print("Sorted data using indices:", sorted_data_using_indices)

```

1)
Original data: [5 2 8 1 6]
Cumulative sum: [5 7 15 16 22]

2)
Original data: [5 2 8 1 6]
Sorted data: [1 2 5 6 8]

3)
Original data: [5 2 8 1 6]
Sorted indices: [3 1 0 4 2]

4)
Original data: [5 2 8 1 6]
Sorted data using indices: [1 2 5 6 8]