

Machine Learning 1 — Final Exam

1 Preliminaries

- Please write your immatriculation number **but not your name** on *every* page you hand in.
- The exam is closed book. You may, however, take one A4 sheet of handwritten notes.
- The exam is limited to 2×60 minutes.
- If a question says “Describe in 2–3 sentences” or “Show your work” or something similar, these mean the same: give a succinct description or explanation.
- This exam consists of 8 pages, 15 problems. You can earn up to 44 points.

Problem 1 [3 points] Fill in your immatriculation number on every sheet you hand in. Make sure it is easily readable. Make sure you do **not** write your name on *any* sheet you hand in.

2 Linear Algebra and Probability Theory

Problem 2 [2 points] Let X and Y be two random variables. Show that

$$\text{var}[X + Y] = \text{var}[X] + \text{var}[Y] + 2\text{cov}[X, Y]$$

where $\text{cov}[X, Y]$ is the covariance between X and Y . You can use that

$$\begin{aligned}\text{var}[X] &= \text{E}[X^2] - \text{E}^2[X] \\ \text{cov}[X, Y] &= \text{E}[XY] - \text{E}[X]\text{E}[Y]\end{aligned}$$

We know [1 point]:

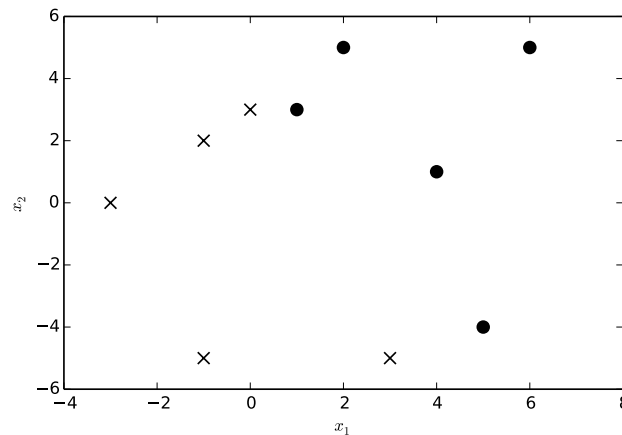
$$\begin{aligned}\text{var}[X] &= \text{E}[X^2] - \text{E}^2[X] \\ \text{cov}[X, Y] &= \text{E}[XY] - \text{E}[X]\text{E}[Y]\end{aligned}$$

Hence [1 point]:

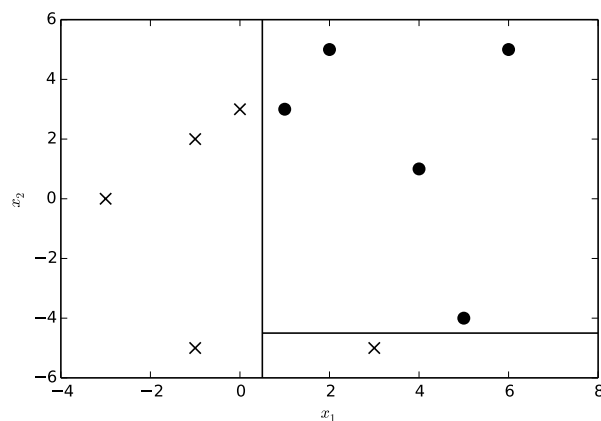
$$\begin{aligned}\text{var}[X + Y] &= \text{E}[(X + Y)^2] - \text{E}^2[X + Y] \\ &= \text{E}[X^2 + Y^2 + 2XY] - (\text{E}[X] + \text{E}[Y])^2 \\ &= \text{E}[X^2] + \text{E}[Y^2] + 2\text{E}[XY] - \text{E}^2[X] - \text{E}^2[Y] - 2\text{E}[X]\text{E}[Y] \\ &= \underbrace{\text{E}[X^2] - \text{E}^2[X]}_{=\text{var}[X]} + \underbrace{\text{E}[Y^2] - \text{E}^2[Y]}_{=\text{var}[Y]} + 2\underbrace{(\text{E}[XY] - \text{E}[X]\text{E}[Y])}_{=\text{cov}[X, Y]} \quad \square\end{aligned}$$

3 Decision Trees & kNN

You are given two-dimensional input data with corresponding targets in below plot.



Problem 3 [2 points] Sketch the decision boundaries of a maximally-trained decision tree classifier using misclassification rate.



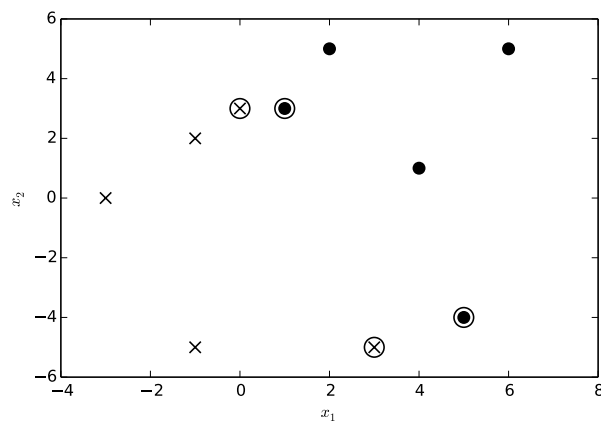
Problem 4 [2 points] Describe how this model can overfit the data. Describe how that problem can be solved or prevented.

Overfitting!

- prune the tree
- restrict the depth of the tree
- train an ensemble of slightly different trees \rightarrow random forests

Problem 5 [2 points] Perform 1-NN with leave-one-out cross validation on the data in the plot. Circle all points that are misclassified and write down the accuracy.

imat:



accuracy = 6/10

4 Linear Classification

Problem 6 [4 points] The decision boundary for some linear classifier on two-dimensional data crosses axis x_1 at 2 and x_2 at 5. First, write down the general form of a linear classifier model (how many parameters do you need, given the dimensions?). Calculate the coefficients (parameters).

$$w_0 + w_1x_1 + w_2x_2 = 0$$

$$w_0 + 2w_1 = 0$$

$$w_1 = -\frac{w_0}{2}$$

$$w_0 + 5w_2 = 0$$

$$w_2 = -\frac{w_0}{5}$$

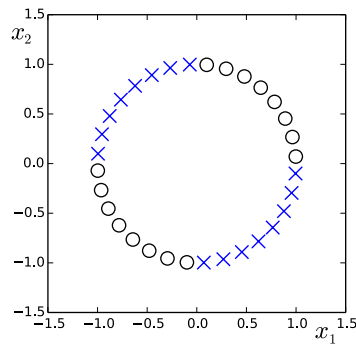
$$\text{set, e.g., } w_0 = -2$$

$$w_1 = 1$$

$$w_2 = \frac{2}{5}$$

or anything proportional

Problem 7 [2 points] Which basis function $\phi(x_1, x_2)$ makes the data in the example below linearly separable (crosses in one class, circles in the other)?



$$\phi(x_1, x_2) = x_1 x_2$$

5 Gaussian Processes

The posterior conditional for an MVN with distribution

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} \sim \mathcal{N} \left(\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

is given by

$$\begin{aligned} p(\mathbf{y}_1 | \mathbf{y}_2) &= \mathcal{N}(\mathbf{y}_1 | \mu_{1|2}, \Sigma_{1|2}) \\ \mu_{1|2} &= \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (\mathbf{y}_2 - \mu_2) \\ \Sigma_{1|2} &= \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \end{aligned}$$

Assume a noise-free GP with mean function

$$m(x) = 0$$

and covariance function

$$K(x, x') = 1 + (x - 2)(x' - 2).$$

You are given two data points $(0, 4) =: \mathbf{x}$.

Problem 8 [2 points] Compute the kernel matrix (aka covariance matrix) for \mathbf{x} .

$$\mathbf{K} = \begin{pmatrix} 5 & -3 \\ -3 & 5 \end{pmatrix}$$

Problem 9 [6 points] Given corresponding outputs $\mathbf{y} = (2, 6)$, compute the posterior function values for data points $\mathbf{x}_* = (0, 2, 4)$.

$$\begin{aligned}
\mathbf{K}^{-1} &= \begin{pmatrix} \frac{5}{16} & \frac{3}{16} \\ \frac{3}{16} & \frac{5}{16} \end{pmatrix} \\
\mathbf{K}_*^T &= \begin{pmatrix} 5 & -3 \\ 1 & 1 \\ -3 & 5 \end{pmatrix} \\
\mu_{f_*|y} &= m(\mathbf{x}_*) + \mathbf{K}_*^T \mathbf{K}^{-1} (\mathbf{y} - m(\mathbf{x})) \\
&= 0 + \begin{pmatrix} 1 & 0 \\ 1/2 & 1/2 \\ 0 & 1 \end{pmatrix} \left(\begin{pmatrix} 2 \\ 6 \end{pmatrix} - 0 \right) \\
&= \begin{pmatrix} 2 \\ 4 \\ 6 \end{pmatrix}
\end{aligned}$$

Problem 10 [4 points] Which other algorithm does this resemble? Please describe the corresponding feature space.

K is a linear kernel, and we are basically solving linear regression.

6 Neural networks

Problem 11 [2 points] Geoffrey has a data set with input $X \in \mathbb{R}^2$ and output $Y \in \mathbb{R}^1$. He has neural network A with one hidden layer and 9 neurons in that layer, which can fit the data. However, he does not know how good the model is, so he also tests neural network B with two hidden layers and three neurons for each of these layers. Both models have biases for the hidden units only.

- How many free parameters do the two models have? Show your calculation, not just the result.
- What are the pros and cons of model A compared to model B? Mention at least one pro and one con.

Model A and model B have 36 and 24 free parameters, respectively. Pros: fewer parameters for B. Cons: more complicated structure may prone to local minimum.

Problem 12 [4 points] You know that the sum of squared errors is related to the Gaussian distribution—differently put, if you assume a normal distribution of the data around their expectation, the maximum likelihood estimate (MLE) is reached when the summed squared errors is minimised.

The same is true for a Laplace distribution and the sum of absolute errors. In particular, if the data observes a Laplacian distribution

$$p(\mathbf{z}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N p(z_n|x_n, \mathbf{w}, \beta) = \prod_{n=1}^N \frac{1}{2\beta} \exp\left(-\frac{|z_n - y(x_n, \mathbf{w})|}{\beta}\right)$$

then minimising the summed absolute errors

$$\sum_{n=1}^N |z_n - y(x_n, \mathbf{w})|$$

leads to MLE. In these equations, \mathbf{x} is the vector of all inputs, x_n is the input of sample n , while $y(x_n, \mathbf{w})$ is the neural network prediction on x_n . Then, z_n is the desired output for x_n .

Show that the MLE of the Laplace distribution minimises the sum of absolute errors.

Taking the negative logarithm, we obtain the error function

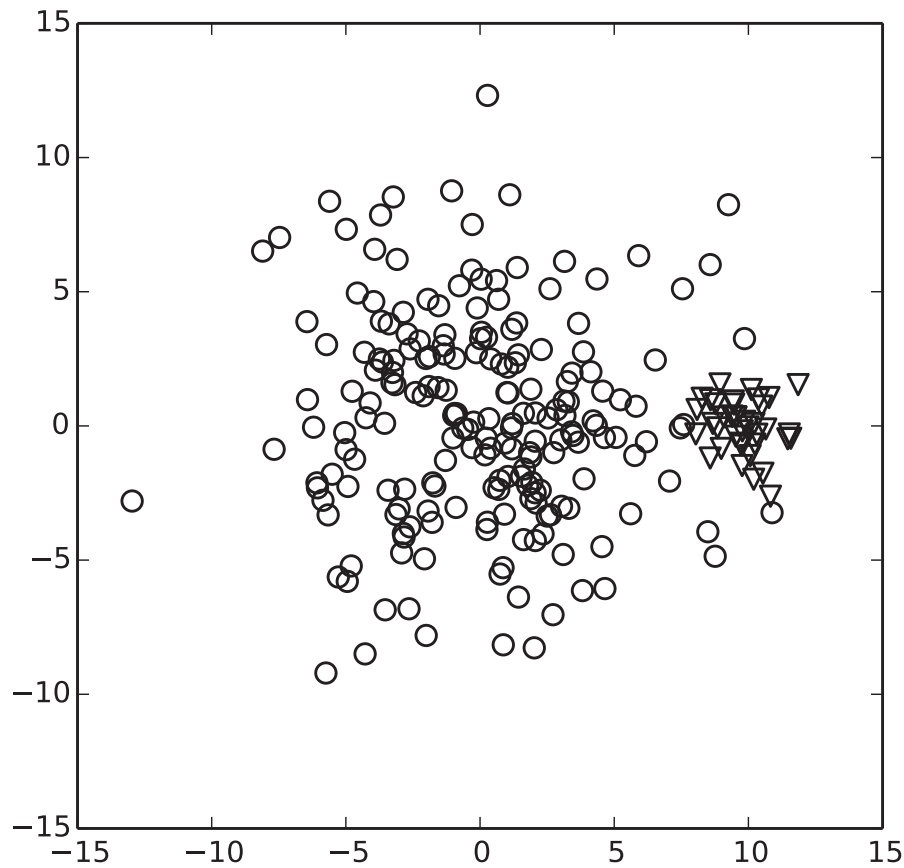
$$\frac{\beta}{2} \sum_{n=1}^N |z_n - y(x_n, \mathbf{w})| - \frac{N}{2} \ln \beta + \frac{N}{2} \ln(2\pi).$$

Maximising the likelihood function is equivalent to minimising the error given by the sum of absolutes

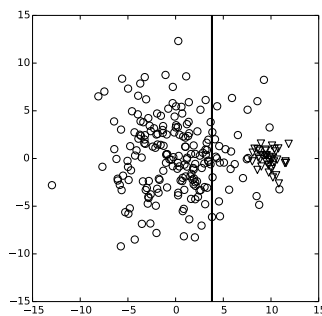
$$\sum_{n=1}^N |z_n - y(x_n, \mathbf{w})|.$$

7 Unsupervised learning

Problem 13 [2 points] Consider the plot below. The two classes are circles and triangles. What significance do class labels have for k-means? Draw the resulting decision boundary in the plot for k-means with two centroids.

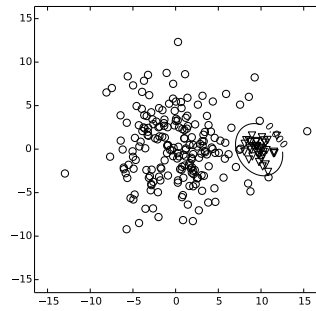


Exact position is not important, but should not separate the two classes perfectly.



Problem 14 [3 points] Would the separation be different using the EM algorithm with a Gaussian mixture model using two components and individual full covariance matrices? What would it look like? The left cluster has 200 points and the right cluster has 40 points. Draw qualitatively in the above figure.

EM for GMM would be able to distinguish better between the two classes. k-means only fits clusters of the same size. EM for GMM can model different sizes for each cluster. The new boundary is now a circle around the right cluster:



Problem 15 [4 points] The likelihood for ICA is

$$\ell(\mathbf{W}) = \sum_{i=1}^m \sum_{j=1}^n \log p_{\mathbf{S}_i}(\mathbf{w}_i^T \mathbf{x}) + \log |\mathbf{W}|$$

When calculating the gradient for this likelihood we need to compute the derivative of $\log p_{\mathbf{S}_i}(\mathbf{w}_i^T \mathbf{x})$. Here we are using $p_{\mathbf{S}_i}(s) \approx \sigma'(s)$ where the sigmoid function is given as

$$\sigma(a) = \frac{1}{1 + e^{-a}}.$$

with the derivative

$$\sigma'(a) \equiv \frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a)).$$

Show that $\frac{d}{ds} \log \sigma'(s) = 1 - 2\sigma(s)$.

$$\frac{d}{ds} \log \sigma'(s) = \frac{\sigma''(s)}{\sigma'(s)} = \frac{\frac{d}{ds} \frac{e^{-s}}{(1+e^{-s})^2}}{\frac{e^{-s}}{(1+e^{-s})^2}} = \frac{\frac{e^{-s}(e^{-s}-1)}{(1+e^{-s})^3}}{\frac{e^{-s}}{(1+e^{-s})^2}} = \frac{e^{-s}-1}{1+e^{-s}} = \frac{e^{-s}+1-2}{1+e^{-s}} = 1 - \frac{2}{1+e^{-s}} = 1 - 2\sigma(s)$$