

**Final Exam**

January 20th, 2015

First and Last name: \_\_\_\_\_

ETH number: \_\_\_\_\_

Signature: \_\_\_\_\_

**General Remarks**

- You have 2 hours for the exam. There are five sections.
- Write your answers directly on the exam sheets. At the end of the exam you will find supplementary sheets, feel free to separate them from the exam. If you submit the supplementary sheets, put your name and ETH number on top of each.
- Answer the questions in English. Do not use a pencil or red color pen.
- You may provide at most one valid answer per question. Invalid solutions must be canceled out clearly.

	Topic	Max. Points	Points	Signature
1	Kernels & SVMs	24		
2	Unsupervised learning	24		
3	Regression	24		
4	Bayesian inference	24		
5	Ensemble methods	24		
Total		120		

Grade: .....

*This page has been intentionally left blank.*

### Question 1: Kernels and SVMs (24 pts.)

Let us consider a dataset with the following positively labeled points:

$\{(2, 2); (2, -2); (-2, 2); (-2, -2)\}$ ; and negatively labeled points

$\{(1, 1); (1, -1); (-1, 1); (-1, -1)\}$ . We shall define a non-linear mapping from the input space to a new feature space according to the following transformation:

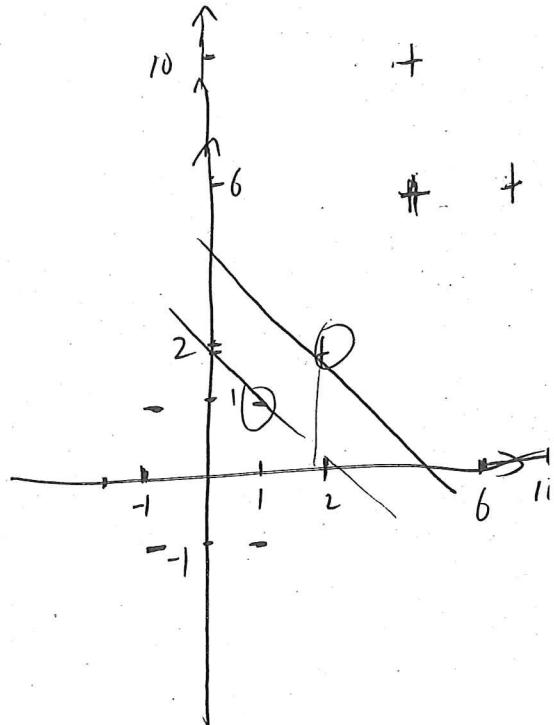
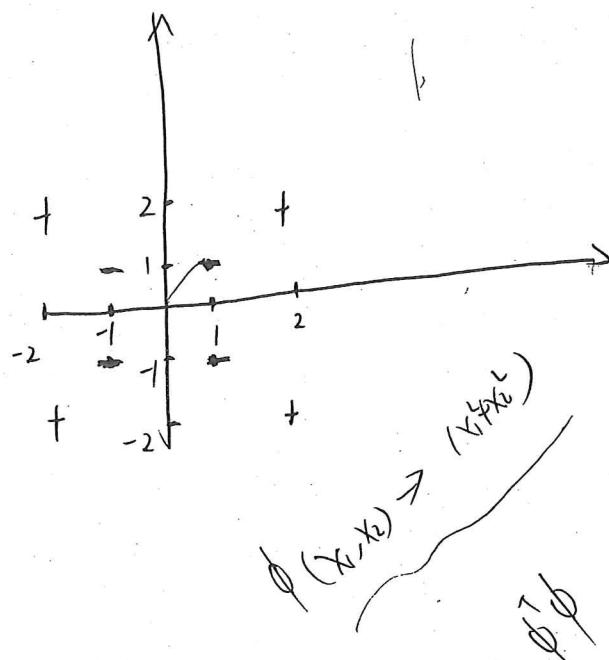
$$\Phi(x_1, x_2) = \begin{cases} (4 - x_2 + |x_1 - x_2|, 4 - x_1 + |x_1 - x_2|) & \text{if } \sqrt{x_1^2 + x_2^2} > 2 \\ (x_1, x_2) & \text{otherwise} \end{cases}$$

$x_1, x_2$  refers to the two feature dimensions.

1. Calculate the transformed sample points for each class. In two separate diagrams plot the original data points and the transformed data points. Use separate symbols for different classes (e.g.,  $+$ ,  $-$ ).

$$+ \quad \left\{ (2, 2); (6, 6); (6, 10); (6, 6) \right\}$$

$$- \quad \left\{ (1, 1); (1, -1); (-1, 1); (-1, -1) \right\}$$



3 pts.

$$s_1 = (2, 2) \quad s_2 = (1, 1)$$

2. In the above diagram identify the support vectors and denote them as  $s_i$ . Define an augmented feature space as  $\tilde{s} = \{s, 1\}$ . Our objective is to find the optimal Lagrangian dual variables  $\alpha_i$  for which the following constraints are defined

$$\alpha_1 \tilde{s}_1 \cdot \tilde{s}_1 + \alpha_2 \tilde{s}_1 \cdot \tilde{s}_2 = -1$$

$$\alpha_1 \tilde{s}_1 \cdot \tilde{s}_2 + \alpha_2 \tilde{s}_2 \cdot \tilde{s}_2 = +1.$$

Explain how these constraints are derived. From the above constraint equations calculate the values of  $\alpha_1$  and  $\alpha_2$

Since for the dual function

$$\text{we have } w^* = \sum_i \alpha_i \tilde{s}_i$$

$$\text{for } s_1 \text{ we have } \tilde{s}_1^T w^* = 1$$

thus

$$\sum_i \alpha_i \tilde{s}_1^T \tilde{s}_i = 1$$

$$\text{for } s_2 \underbrace{-1 \times (\alpha_1 \tilde{s}_1 \cdot \tilde{s}_1 + \alpha_2 \tilde{s}_2 \cdot \tilde{s}_1)}_{= 1}$$

$$\text{for } s_2 \underbrace{(-\alpha_1 \tilde{s}_1 \cdot \tilde{s}_2 + \alpha_2 \tilde{s}_2 \cdot \tilde{s}_2)}_{= 1} = 1$$

4 pts.

3. Calculate the separating hyperplane  $w$  and plot the decision surface in your diagram from Question 1.1 showing the transformed sample points.

$$w = \sum \alpha_i (1, 1, 1) + \alpha_2 (2, 2, 1)$$

$$= -(1, 1, 1) + (2, 2, 1) = (1, 1, -3) \quad 5 \text{ pts.}$$

4. Given a test point  $x$  in the original feature space, write down the classification function after mapping it to the transformed feature space.

$$\begin{cases} \text{if } \operatorname{sgn} \left[ (1, 1, 1)^T (x_1, x_2, 1)^T \right] \leq 2 \\ \operatorname{sgn} \left[ (1, 1, 1)^T (4-x_2-x_1, 4-x_1+x_2, 1)^T \right] > 2 \end{cases} \quad \sqrt{x_1^2 + x_2^2} > 2 \quad 3 \text{ pts.}$$

(Q2, 1)

5. To which class is  $x = (4, 5)$  assigned? Discuss whether the classification seems intuitively correct. If not, what are the possible reasons for this discrepancy and what changes can improve the generalization of the classifier?

could not decide

wrong model is not suitable,  
the  
change to  
use the  
gaussian kernel

add more data

increase the data

feature space,  
kernel function  
appropriate.

5 pts.

6. Identify, with explanation, whether the following functions are kernels or not

$$(a) k(x, y) = \begin{cases} 10 & \text{if } x \neq y \\ 1 & \text{if } x = y \end{cases}$$

not considering (1,1) (2,2)

$$K = \begin{bmatrix} 1 & 10 \\ 10 & 1 \end{bmatrix} \quad \text{delta}(k) \geq 0$$

not semi-positive

2 pts.

$$(b) k(x, y) = \|x - y\|^2 + \|x + y\|^2$$

$$\begin{bmatrix} 4 & 1 & -1 \end{bmatrix}$$

not

considering to sample (1, 12)

$$K = \begin{bmatrix} 2 & 5 \\ 5 & 8 \end{bmatrix}$$

2 pts.

$$\begin{matrix} 5 \\ \text{delta } |k| < 0 \end{matrix}$$

**Question 2: Unsupervised learning (24 pts.)**1. *Histograms and Parzen estimators.*

- (a) Assume you have  $N$  data points  $x_i \in [0, 1], i \in 1, \dots, N$  (all the points are in the interval from 0 to 1). Give the definition of the histogram with  $K$  bins (use equidistant binning strategy). Explain how the resulting histogram can be used for density estimation. Will the "curse of dimensionality" be a problem in this case and why?

to divide the space  $[0, 1]$  to  $K$  <sup>intervals</sup>  $\frac{1}{K}$

which the corresponding histogram of the sample is

$$H_j := \# \{x \in S | x \in I_j\}$$

after normalization, we could get  $A := \frac{1}{N} (h_1, \dots, h_K)$

no, since it is one dimension,

**6 pts.**

(b) Let  $k$  be the unit cube function:

$$k(\mathbf{u}) = \begin{cases} 1 & \text{if } |u_i| \leq 1/2 \text{ for all } i \in 1, \dots, d \\ 0 & \text{otherwise} \end{cases}$$

A unit cube function (scaled to size  $h$ ) around a training data sample  $\mathbf{x}_j$  is obtained by  $k\left(\frac{\mathbf{x}-\mathbf{x}_j}{h}\right)$ . Using such unit cubes, we can estimate the probability density of a distribution that generated training data  $\{\mathbf{x}_j\}_{j=1}^n$  using the formula:

$$p(\mathbf{x}) = \frac{K(\mathbf{x})}{nV}$$

where  $K(\mathbf{x})$  is the number of samples falling within a window around  $\mathbf{x}$ ,  $n$  is the total number of training samples and  $V$  is the area of the window around  $\mathbf{x}$ . What are the values of  $K(\mathbf{x})$  and  $V$  for a Parzen window estimate in this setting?

$$p(\mathbf{x}) = \frac{1}{h} \frac{1}{V_n} \sum_{j=1}^n k\left(\frac{\mathbf{x}-\mathbf{x}_j}{h}\right)$$

$$V_n = h^d = h$$

$$K(\mathbf{x}) = \sum_{j=1}^n k\left(\frac{\mathbf{x}-\mathbf{x}_j}{h}\right)$$

5 pts.

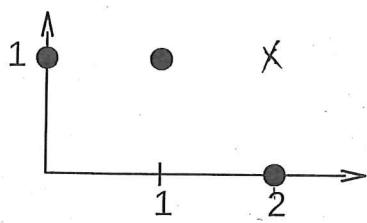


Figure 1: A figure for question (c).

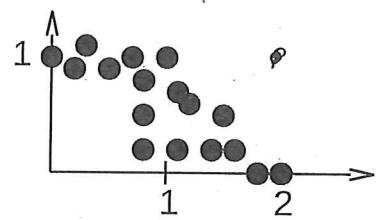


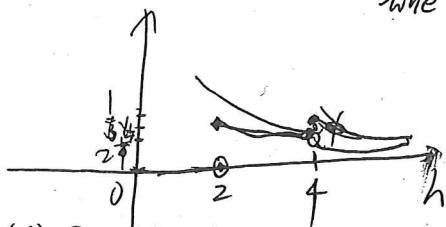
Figure 2: A figure for question (d).

- (c) Consider training data with  $d = 2$  dimensions and three samples  $(0, 1)$ ,  $(1, 1)$  and  $(2, 0)$  (see Figure 1):  
Using the Parzen window estimate (from the previous point), plot the value of  $p(w)$  for  $w = (2, 1)$  with respect to the parameter  $h$ .

$$P(w) = \frac{1}{3} \left( \frac{1}{h^2} \sum_{i=1}^3 k\left(\frac{\|w - x_i\|}{h}\right) \right) \quad \text{when } 0 < h \leq 2, P(w) = 0$$

$$\text{when } 2 \leq h < 4, P(w) = \frac{2}{3h^2}$$

$$\text{when } h \geq 4, P(w) = \frac{1}{h^2}$$



5 pts.

- (d) Consider the training data as shown in Figure 2:

Would you use the same Parzen window estimate for this training data and for the training data from the previous point? Motivate your decision.

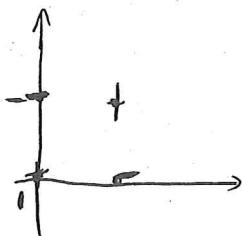
$h_0$ , since the train data became intensive

thus the  $h$  value could be smaller,  
~~to increase~~ thus not to lose detail

4 pts.

2. *K-Nearest Neighbours.*

Let the dataset for binary classification consist of only four points in  $\mathbb{R}^2$ :  $\mathbf{x}_1 = (0, 0), \mathbf{x}_2 = (0, 1), \mathbf{x}_3 = (1, 1), \mathbf{x}_4 = (1, 0)$ . The corresponding labels are the following:  $y_1 = 1, y_2 = 0, y_3 = 1, y_4 = 0$ . Compute leave-one-out cross-validation error of the 3-NN classifier.



$$\begin{aligned} & \frac{1}{4} \sum_{i=1}^4 \mathbb{I}_{y_i \neq g(\mathbf{x}_i)} \\ &= \frac{1}{4} (1 + 1 + 1 + 1) \\ &= 1 \end{aligned}$$

4 pts.

**Question 3: Regression (24 pts.)**

Consider the linear regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where  $\mathbf{y} \in \mathbb{R}^n, \mathbf{X} \in \mathbb{R}^{n \times p}, \boldsymbol{\beta} \in \mathbb{R}^p, \boldsymbol{\epsilon} \in \mathbb{R}^n$ .

1. (a) Write down the least squares estimator  $\hat{\boldsymbol{\beta}}$  as a function of  $\mathbf{y}$  and  $\mathbf{X}$ .

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

1 pts.

- (b) Under what statistical assumption is  $\hat{\boldsymbol{\beta}}$  unbiased?

$$E(\hat{\boldsymbol{\beta}}) = E((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y})$$

$$= \mathbf{B} \cdot \boldsymbol{\beta} + E(\boldsymbol{\epsilon})$$

$$\text{where } E(\boldsymbol{\epsilon}) = 0$$

2 pts.

- (c) Formally show that  $\hat{\boldsymbol{\beta}}$  is unbiased given the above assumption.

9 pts.

2. Name at least one reason why a gradient descent method would be preferred over the closed form solution for obtaining  $\hat{\beta}$ .

*it could be hard to compute the inverse of  $(X^T X)$   
the covariance matrix  
furthermore, when  
the data sample in this the singular matrix of the  $X$  function  
problem*

2 pts.

3. Explain qualitatively how the bias and variance of the estimator are related with under- and over-fitting.

*when over-fitting the bias is small while  
the variance become big*

*under fitting bias is big while variance small*

3 pts.

4. Consider using cross validation for estimating the generalization error.

- (a) Write down the value of the loss function for a new point  $x$  in the test set.

$$\text{loss} = \frac{1}{2} (\cancel{x}(X^T X)^{-1} \cancel{x}^T y' - y)^2$$

$$(\cancel{x}(X^T X)^{-1} \cancel{x}^T y' - y)^2$$

2 pts.

- (b) Formally define the generalization error and suggest a way how to estimate it.

*for the test set with  $n$  sample*

$$\text{generalization error} = \frac{1}{n} \sum_{i=1}^n (x_i^T \hat{\beta} - y_i)^2$$

3 pts.

5. What are the advantages of the LASSO over ridge regression?

it LASSO estimates are known to be sparse  
with few coefficients non-vanishing, thus  
make it more efficient and  
be better at interpreting

2 pts.

feature selection

**Question 4: Bayesian inference and maximum likelihood (24 pts.)**

1. Let  $X_1, \dots, X_n$  be  $n$  mutually independent and normally distributed random variables,  $X_i \sim \mathcal{N}(\mu, \sigma^2)$ , with a common unknown mean and common known variance. Let  $Y$  be a sum of the random variables  $Y = \sum_{i=1}^n X_i$ .

- (a) How is the likelihood  $p(y|\mu, \sigma^2)$  distributed? Derive the density function.

*Hint:* The moment generating function of a univariate normal distribution is  $M_X(t) = \exp(t\mu + \frac{1}{2}\sigma^2 t^2)$ .

$$p(y|\mu, \sigma^2) = \exp\left(\mu n + \frac{1}{2}\sigma^2 \sum_{i=1}^n y_i^2\right)$$

$$= \exp\left(\mu n + \frac{1}{2}\sigma^2 Y^2\right)$$

Since  $X_i$  are iid

$Y$  is also a  $\text{d}$  normally distributed  
variables

$$E(Y) = \sum_{i=1}^n E(X_i) = n\mu$$

$$\text{Var}(Y) = \sum_{i=1}^n \text{Var}(X_i) = n\sigma^2$$

$$p(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi n}} e^{-\frac{1}{2n\sigma^2} (Y-n\mu)^2}$$

4 pts.

- (b) Determine the maximum a posteriori estimate of  $\mu$  given  $y$ . Assume that the prior distribution is a normal distribution  $N(\mu|\mu_0, \sigma_0^2)$ .

$$\max_{\mu, \sigma^2} p(\mu | Y) \propto p(Y | \mu, \sigma^2) \cdot p(\mu | \mu_0, \sigma_0^2)$$

$$e^{-\frac{(Y-\mu)^2}{2\sigma^2}} \cdot e^{-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}}$$

$$\max N(\mu | \mu_n, \sigma_n^2)$$

$$=\mu_n$$

$$\mu_n = \frac{\frac{Y}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}}{\frac{1}{\sigma^2} + \frac{1}{\sigma_0^2}} = \frac{Y\sigma_0^2 + \mu_0\sigma^2}{n\sigma^2 + \sigma_0^2}$$

**5 pts.**

(c) Determine the probability density function of the posterior  $p(\mu|y, \mu_0, \sigma_0, \sigma)$ .

$$p(\mu|y, \mu_0, \sigma_0, \sigma)$$

$$\cong N(\mu|\mu_h, \sigma_h^2)$$

$$\cong e^{-\left(\frac{h}{2\sigma^2} + \frac{1}{2\sigma_0^2}\right)} (\mu - \mu_h)^2$$

$$\frac{h}{2\sigma^2} + \frac{1}{2\sigma_0^2} = \frac{1}{2\sigma_h^2}$$

$$\Rightarrow \sigma_h^2 = \frac{1}{\frac{h}{\sigma^2} + \frac{1}{\sigma_0^2}}$$

$$= \frac{\sigma^2 \sigma_0^2}{h \sigma_0^2 + \sigma^2}$$

$$\Rightarrow p = \frac{\sqrt{h \sigma_0^2 + \sigma^2}}{\sqrt{2\pi} \sigma_h^2} \cdot e^{-\frac{(h \sigma_0^2 + \sigma^2)(\mu - \frac{y_0 \sigma_0^2 + \sigma^2}{h \sigma_0^2 + \sigma^2})^2}{2 \sigma_h^2}}$$

5 pts.

2. (a) Let  $\mathcal{X} = \{x_1, \dots, x_n\}$  be a dataset consisting of  $n$  samples which are drawn i.i.d. from a normal distribution  $\mathcal{N}(x|\mu, \sigma^2)$  with unknown mean and unknown variance. Recall that the maximum likelihood estimate of the mean is given by:

$$\hat{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$$

Determine the maximum likelihood estimate of the variance  $\hat{\sigma}_{ML}^2$ . Which estimator has a larger bias, the mean or the variance? Why?

$$\log L(x) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu_{ML})^2}{\sigma^2} - \frac{n}{2} \log n$$

$$\frac{\partial \log L(x)}{\partial \sigma^2} = -\frac{n}{2} \frac{2\sigma^{-1}}{2\sigma^2} + \frac{\sum (x_i - \mu_{ML})^2}{2\sigma^4} = 0$$

$$\Rightarrow \hat{\sigma}_{ML}^2 = \frac{\sum (x_i - \mu_{ML})^2}{n} \\ = \frac{1}{n} \sum_i^n \left( x_i - \frac{1}{n} \sum_j^n x_j \right)^2 = \frac{1}{n} \overline{(n-1)\sigma^2}$$

$$E(\mu_{ML}) = \frac{1}{n} \cdot n\mu = \mu$$

$$\text{bias}_{\mu_{ML}} = 0$$

while  $\hat{\sigma}_{ML}^2$  is biased, since  $\mu_{ML}$  is unknown

3 pts.

- (b) Let  $\mathcal{X} = \{x_1, \dots, x_n\}$ ,  $x_{i:n} \geq 0$  be a dataset consisting of  $n$  samples which are drawn i.i.d. from a Weibull distribution:

$$\text{Weibull}(x_i|c, \lambda) = \left(\frac{c}{\lambda}\right) \left(\frac{x_i}{\lambda}\right)^{c-1} \exp\left(-\left(\frac{x_i}{\lambda}\right)^c\right).$$

Determine the derivative of the log-likelihood  $\frac{\partial}{\partial c} (\ln p(\mathcal{X}|c, \lambda))$ . Can you find an analytic expression for the maximum likelihood estimate  $\hat{c}_{\text{ML}}$ ? If not, suggest a method to compute  $\hat{c}_{\text{ML}}$ .

$$\begin{aligned} \ln p(\mathcal{X}|c, \lambda) &= \sum_{i=1}^n \ln(p(x_i|c, \lambda)) \\ &= \sum_{i=1}^n \left[ \ln c - \ln \lambda + (c-1) \ln x_i - (c-1) \ln \lambda - \left(\frac{x_i}{\lambda}\right)^c \right] \\ &= \cancel{n \ln c} - n \ln \lambda + n(c-1) \ln x_i - n \cancel{\left(\frac{x_i}{\lambda}\right)^c} \\ f(c) &= \frac{\frac{\partial \ln p(\mathcal{X}|c, \lambda)}{\partial c}}{f(c)} = \frac{n}{c} - n \ln \lambda + n \ln x_i - n \cancel{\left(\frac{x_i}{\lambda}\right)^c} \cancel{\ln \frac{x_i}{\lambda}} = 0 \\ \text{no analytic expression for } c_{\text{ML}} \end{aligned}$$

can use numerical ~~new~~ method to

approximate it. gradient ~~descent~~  
~~gradient~~

$$f(c_1) > 0 \quad f(c_2) < 0$$

the find  $f\left(\frac{c_1+c_2}{2}\right) > 0$  or  $< 0$  4 pts.  
where

if  $> 0$

$$\text{let } c_1 \leftarrow \frac{c_1+c_2}{2}$$

$$\text{and } c_2 \leftarrow \frac{c_1+c_2}{2}$$

17

see the fib between  
value of

if it  $> 0$  ... in tel  $(f(c))$

3. Consider a binomial distribution as likelihood function of  $x$ :

$$p(x|\theta, n) = \frac{n!}{x!(n-x)!} \theta^x (1-\theta)^{n-x},$$

Determine the posterior distribution  $p(\theta|x, n)$  (written in direct closed form) using the following beta distribution as a prior:

$$p(\theta|\alpha, \beta) = B(\alpha, \beta)^{-1} \theta^{\alpha-1} (1-\theta)^{\beta-1}.$$

The beta function  $B(\alpha, \beta)$  is a normalization constant.

$$\begin{aligned} p(\theta|x, n) &= \cancel{p(x|\theta, n)} \cdot p(\theta|\alpha, \beta) \\ &\propto \cancel{\frac{n!}{x!(n-x)!}} \theta^x (1-\theta)^{n-x} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &\propto \theta^{\cancel{\alpha}+x-1} (1-\theta)^{\cancel{\beta}+n-x-1} \\ p(\theta|x, n) &= B(\alpha', \beta')^{-1} \theta^{\alpha'-1} (1-\theta)^{\beta'-1} \\ \alpha' &= \alpha + x \quad \beta' = \beta + n - x \end{aligned}$$

3 pts.

### Question 5: Ensemble methods (24 pts.)

#### 1. Bagging

One of the possible theoretical arguments for bagging (bootstrap aggregation) is the fact that the expected error of an averaged model (i.e. committee) is one order less than an average expected error of individual models under certain assumptions.

However, this error reduction fails often in reality. Please briefly explain why.

bias to large

$$\begin{aligned}
 & \text{Var} = \left( \frac{1}{B} \sum_{i=1}^B f_i^*(x) - \frac{1}{B} E \sum_{i=1}^B f_i^*(x) \right)^2 \quad \text{it is assume that each} \\
 & = \left( \frac{1}{B} \sum_{i=1}^B f_i^*(x) - E f_i^*(x) \right)^2 \quad \text{that covariances} \\
 & = \left( \frac{1}{B} \sum_{i=1}^B f_i^*(x) - E f_i^*(x) \right)^2 \quad \text{are small} \\
 & = \frac{1}{B^2} \sum_{i=1}^B \text{Var}(f_i^*(x)) + \frac{1}{B^2} \sum_{i \neq j} \text{Cov}(f_i^*(x), f_j^*(x)) \\
 & \approx \frac{\sigma^2}{B} \quad \text{only under the assumption} \quad 2 \text{ pts.} \\
 & \text{from sh subset}
 \end{aligned}$$

#### 2. AdaBoost

- (a) Denote the  $n$  given training samples as  $(x^{(i)}, y^{(i)})$ ,  $i = 1, \dots, n$ , where  $x^{(i)} \in \mathbb{R}^d$  are points and  $y^{(i)} \in \{-1, +1\}$  are labels. Denote the weak learners as  $c_b(x) \in \{-1, +1\}$ ,  $b = 1, \dots, M$ . Sketch the classical AdaBoost algorithm in pseudocode and highlight the training error optimized at each step:

initialize  $w_i^{(0)} = \frac{1}{n}$

for  $b = 1 : B$  do

train  $c_b(x)$  based on training samples with weight  $w_i^{(b)}$

compute the training error

$$\text{err}_b = \frac{\sum w_i^{(b)} I(c_b(x_i) \neq y_i)}{\sum w_i^{(b)}}$$

compute  $\alpha_b = \log \frac{1 - \text{err}_b}{\text{err}_b}$

update  $w_i^{(b+1)} = \frac{w_i^{(b)} e^{\alpha_b I(c_b(x_i) \neq y_i)}}{\sum w_i^{(b)}}$

and for return the combined classifier is  $H(x) = \operatorname{sgn} \left( \sum_{i=1}^B \alpha_i c_i(x) \right)$  4 pts.

- (b) In some formulations of AdaBoost it is required that the training error  $\epsilon_b$  of each weak binary classifier is  $< 0.5$ . How could one guarantee fulfilling such a requirement (i.e. what should be done if  $\epsilon_b > 0.5$ )?

if  $\epsilon_b > 0.5$ , in this case  $\alpha_b = \log \frac{1-\epsilon_b}{\epsilon_b} < 0$

thus this base classifier will be assigned  
a negative weight in the combined classifier,

the reverse one would have a error rate less than its  
base classifier

2 pts.

- (c) Explain how AdaBoost could be used for feature selection. Which weak learners would you propose for that?

use decision stumps to train the model

the Variable importance is measured by

how much ~~err~~ it the  $\Delta_b$  for each variable

2 pts.

used to conduct split/branch. we can choose the  
feature that possess large  $\Delta_b$  for training

### 3. Decision Boundaries of AdaBoost

Assume now that the data are 2-dimensional and we use decision stumps as weak learners:

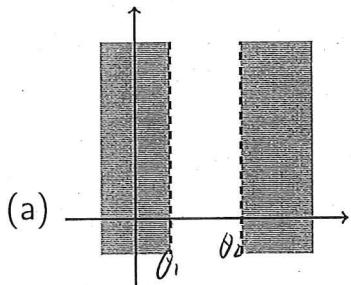
$$c_b(\mathbf{x}) = c_b(\mathbf{x}|k, \theta, \gamma) := \begin{cases} +\gamma & x_k > \theta, \\ -\gamma & \text{otherwise} \end{cases}$$

Here,  $k = 1, 2$  is the dimension, along which the decision is made, and  $\gamma \in \{-1, +1\}$  is the parameter, which provides labelling "direction". A decision stump is a simple binary classifier which slices classes according to a decision line orthogonal to the  $i$ -th axis, and orientation of classification is given by  $\gamma$ . The  $k, \theta$  and  $\gamma$  are *individual* for each step  $b$  (i.e. may vary from step to step).

The final classifier is a linear combination of the trained weak classifiers:

$$c_{\text{AdaBoost}}(\mathbf{x}) = \text{sgn}\left(\sum_{b=1}^M \alpha_b c_b(\mathbf{x})\right)$$

Write down the combination of decision stumps for the following decision boundaries, or explain why it's impossible (Figure convention: decision boundary is a dashed line, and positive class is a shaded region):



not possible

let's assume

the first  $\theta_1$ , second axis  $\theta_2$

$$\text{for } C_1(x) \quad k=1 \quad \gamma=-1 \\ \theta=\theta_1$$

$$\frac{2}{b=1} \alpha_b C_b(x)$$

$$= \text{sgn}(\alpha_1 \theta_1) >$$

$$\alpha_1 > 0$$

$$C_2(x) \quad k=1 \quad \gamma=1 \\ \theta=\theta_2$$

which is not possible.

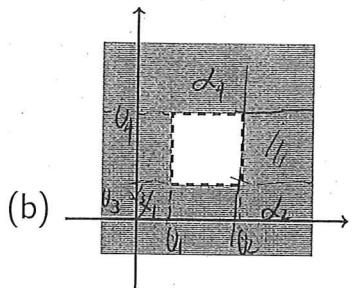
in this case

the part that  $x > \theta_2$

$$\text{we have } \frac{2}{b=1} \alpha_b C_b(x)$$

$$3\gamma(\alpha_2 - \alpha_1) > 0$$

$$\alpha_2 > \alpha_1$$

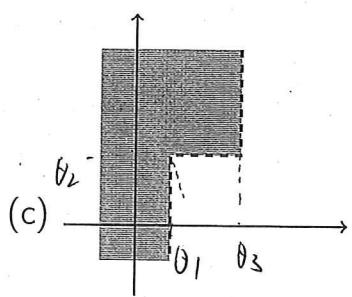


$$d_2 = d_1 \phi - d_4 - d_3 > 0$$

$$d_1 - d_2 - d_3 - d_4 > 0$$

not fully

not possible



$$C_1(x) \quad k=1$$

$$\theta = \theta_1$$

$$Y = -1$$

$$C_2(x) \quad k=2$$

$$\theta = \theta_2$$

$$Y = 1$$

$$C_3(x) \quad k=1$$

$$\theta = \theta_3$$

$$Y = -1$$

3 pts.

2 pts.

#### 4. Decision Trees

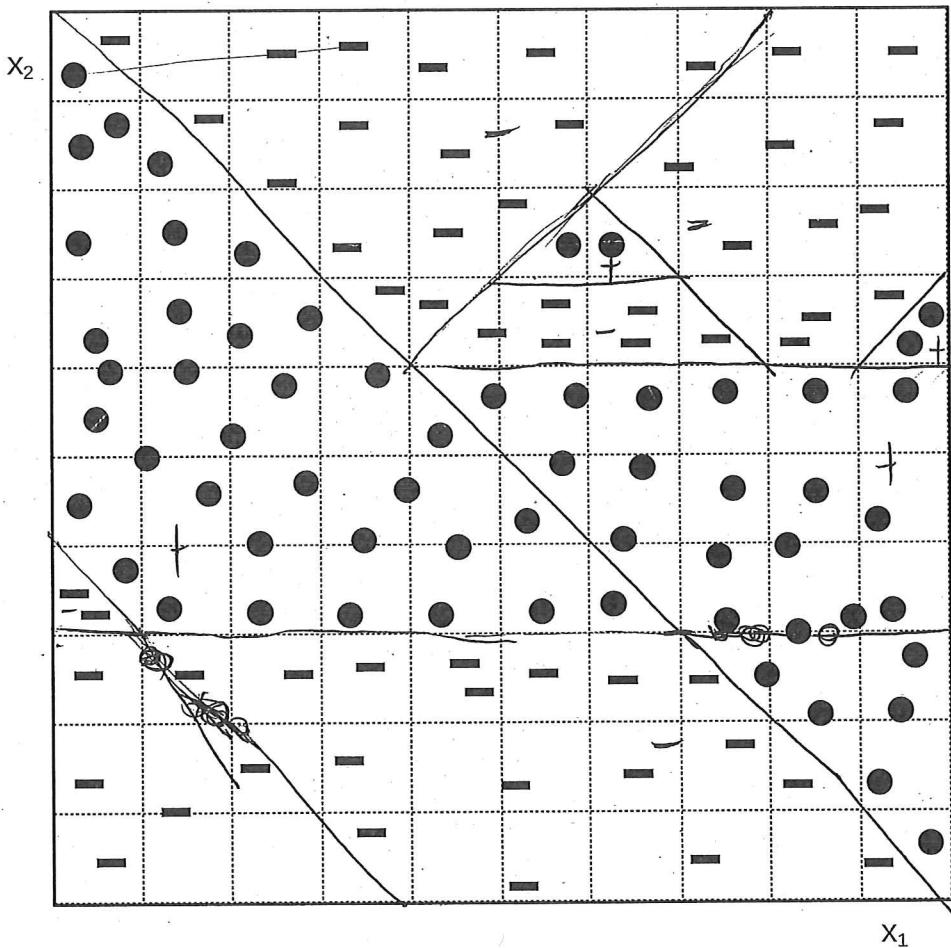
In this exercise, we again consider a binary classification task for a 2-dimensional dataspace. Below you can see a training set for this task.

Consider the set of weak classifiers having the following decision boundaries

$$x_2 = \beta x_1 + \gamma,$$

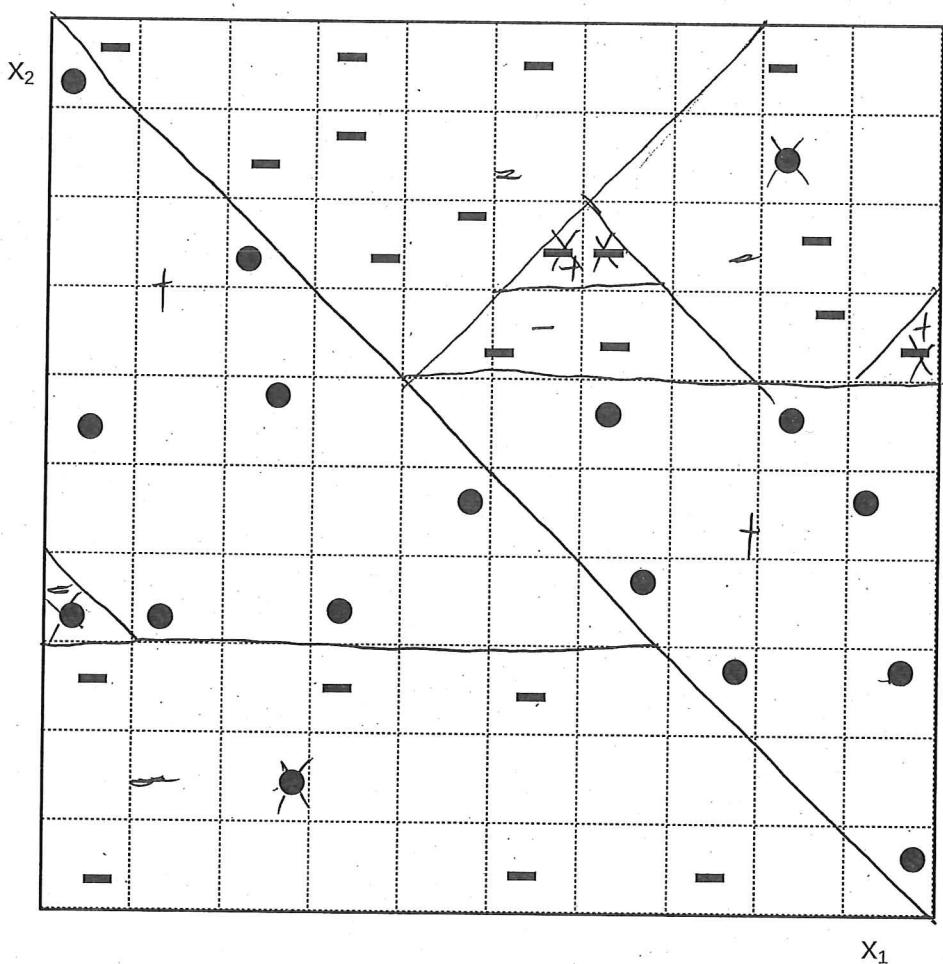
where each weak classifier has *its own* parameters  $\beta \in \{-1, 0, 1\}$  and  $\gamma \in \mathbb{R}$ .

- (a) Sketch graphically the flow of a decision tree algorithm with the above-mentioned weak classifiers (i.e. draw the respective decision bounds) on the given training set, until zero training error is achieved.



2 pts.

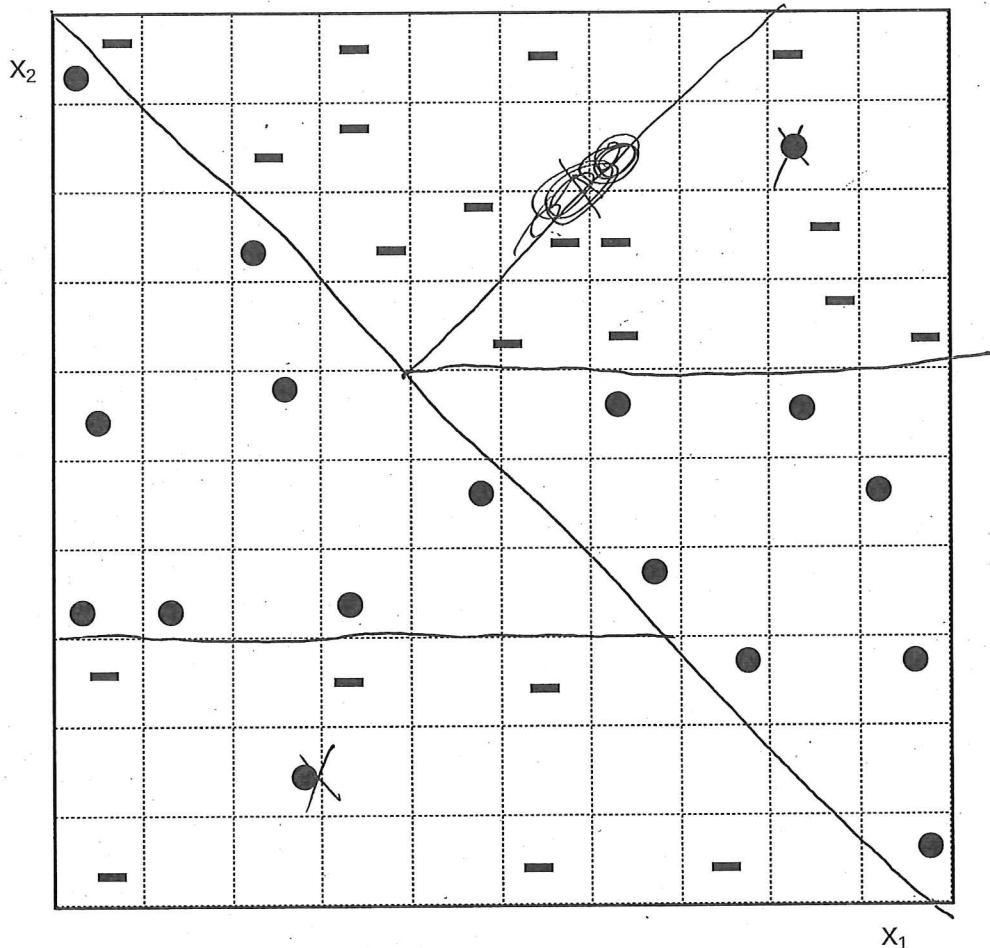
(b) Transfer the obtained boundary to the test set below and compute the test error.



1 pts.

38

- decrease the complexity
- (c) Explain what problems you see with the obtained classifier. How could you modify it in order to make it generalize better? Draw the new one on the test set (it is the same as above) and compute the error.



in this case , 3 pts.

only two misclassification

## Supplementary Sheet