

Prof. J.M. Buhmann

Final Exam

February 6th, 2013

First and Last name: _____

ETH number: _____

Signature: _____

General Remarks

- You have 2 hours for the exam. There are five sections, each of which is worth 20 points. Scoring 100 points guarantees you a grade of six. In two sections you will find bonus questions, worth together 10 points. The bonus questions are a bit more difficult, we suggest you leave them to the end.
- 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	Assorted Questions	20		
2	Bayesian Inf., MAP and ML	20 + 5		
3	Supervised Learning	20		
4	Kernelized Ridge Regression	20		
5	Unsupervised Learning	20 + 5		
Total		100 + 10		

Grade:

This page has been intentionally left blank.

Question 1: Assorted Questions (20 pts.)

1. Figure 1 shows 4 times the same binary classification dataset.

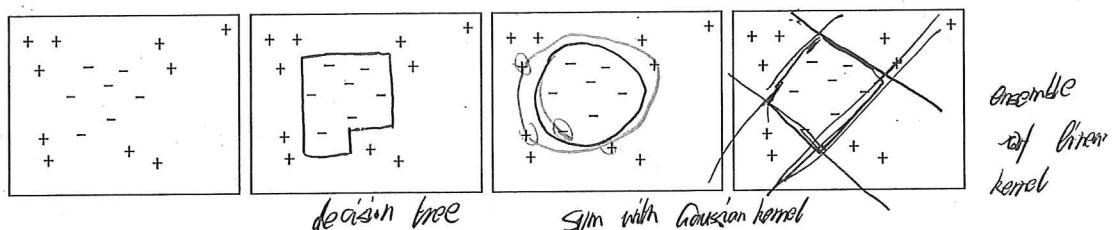


Figure 1: 4 times the same dataset

- (a) Cross all of the following algorithms/classifiers, which can achieve zero training error on this dataset.

- Perceptron
- Decision tree
- SVM with Gaussian kernel
- Ensemble of linear kernel SVMs

- (b) For each of the methods that can achieve zero training error, qualitatively depict a possible decision boundary (having zero error) in one of the plots of the dataset in Figure 1. Indicate which method belongs to which plot.

4 pts.

2. Let \mathcal{F} be an hypothesis class for a binary classification task and f be a randomly chosen prediction function, having a training error of 0.65, on some dataset S . Explain how to use f to obtain \tilde{f} , a prediction function which is **guaranteed** to have a smaller training error than f .

use the ensemble methods - adaboost,

for $b=1:B$:

train $C_b(x)$ on $w^{(b-1)}$ weighted dataset S

$$\text{compute } C_b = \frac{\sum_{i=1}^n w_i^{(b-1)} f_i}{\sum_{i=1}^n w_i^{(b-1)}}$$

$$\text{compute } \alpha_b = \frac{1}{2} \log \frac{1-\epsilon_b}{\epsilon_b}$$

reweight

$$w_i^{(b)} = w_i^{(b-1)} e^{\alpha_b (C_b(x_i) \neq y_i)}$$

$$C(x) = \frac{B}{\sum_{b=1}^B \alpha_b C_b(x)}$$

2 pts.

3. We consider applying the Viterbi algorithm to estimate a trajectory of an HMM with $|S|$ states over T time points. Assume that the number of states $|S|$ grows as $O(\sqrt{T})$. What is worst-case asymptotic computational complexity of the algorithm (as a function of T)?

3 pts.

4. For each of the following statements, circle the correct answer below.

- (a) The number of nodes in a decision tree is bounded by the number of features.

True/False

- (b) Boosting classifiers can in principle be done in a parallel manner.

True/False

- (c) Can the Baum-Welch algorithm be considered a type of Expectation Maximization procedure?

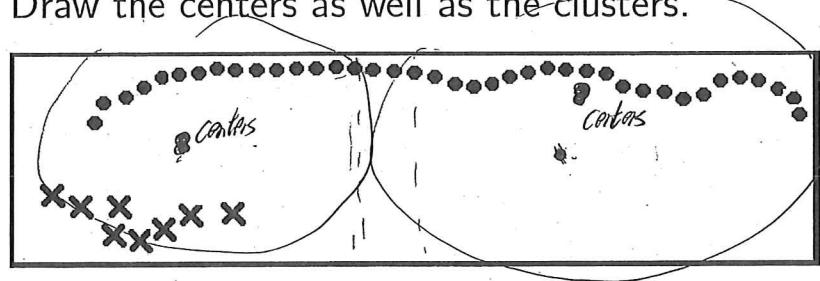
Always/Never/Only Sometimes

是 EM 算法在 HMM 的具体表现

4 pts.

5. The following figures show a dataset of 48 objects from two different sources, represented by different symbols.

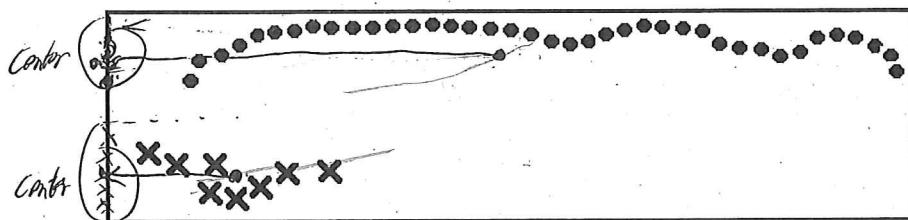
- (a) Sketch the optimal K -means solution on this dataset, for $K = 2$. Draw the centers as well as the clusters.



3 pts.

- (b) Consider reducing the dimensionality of the data to 1 before finding a 2-means solution. Propose an appropriate dimension reduction by drawing a projection line through the estimated center of mass of the data.

Now sketch the optimal 2-means solution on the dimension reduced data.



4 pts.

Question 2: Bayesian Inference, MAP and ML estimation (20 pts.)

1. Let $x_1, \dots, x_n \in \mathbb{R}$ be a dataset consisting of n samples which are assumed to be drawn iid from a normal distribution $\mathcal{N}(x|\mu, \sigma^2)$ in which the variance σ^2 and the mean μ are unknown. Demonstrate that the maximum likelihood estimation of μ can be performed without knowing the maximum likelihood estimate of the variance.

$$\begin{aligned} L &= \prod_{i=1}^n \log P(x_i|\mu, \sigma^2) \\ &= \sum_{i=1}^n \left(-\frac{(x_i-\mu)^2}{2\sigma^2} - \frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 \right) \end{aligned}$$

$$\frac{\partial L}{\partial \mu} = \sum_{i=1}^n \frac{2(x_i-\mu)}{2\sigma^2} = 0$$

$$\Rightarrow \sum_{i=1}^n (x_i - \mu) = 0$$

$$\Rightarrow \mu = \frac{1}{n} \sum_{i=1}^n x_i$$

8 pts.

2. Consider the following Maximum a Posteriori estimation task. The likelihood function is the normal distribution with unknown mean μ and variance $\sigma^2 = 1$. Let μ_0 and σ_0^2 respectively denote the mean and variance of the prior, and recall that the posterior has mean and variance respectively given by

$$\begin{aligned} \mu_n &= \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0 + \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \left(\frac{1}{n} \sum_{i=1}^n x_i \right), \\ \sigma_n^2 &= \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1}. \end{aligned} \tag{1}$$

(a) Show how to derive the above posterior formula for μ , from the prior and the likelihood function.

$$\begin{aligned}
 P(\mu, \sigma^2 | X) &\propto L(X | \mu, \sigma^2) P(\mu, \sigma^2 | \mu_0, \sigma_0^2) \\
 &\propto e^{-\frac{n}{2} \left[\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 + \frac{1}{\sigma_0^2} (\mu - \mu_0)^2 \right]} \\
 &e^{-\frac{1}{2} \left[\frac{n\mu^2 - 2\mu\sum x_i + \sum x_i^2}{\sigma^2} + \frac{\mu^2 - 2\mu\mu_0 + \mu_0^2}{\sigma_0^2} \right]} \\
 &\frac{\sigma_0^2 n \mu^2 - \sigma_0^2 \sum_{i=1}^n 2\mu x_i + \sigma_0^2 \sum_{i=1}^n x_i^2 + \sigma^2 \mu^2 - 2\mu \sigma^2 \mu_0 + \sigma^2 \mu_0^2}{\sigma^2 \sigma_0^2} \\
 &\mu = \frac{\sigma_0^2 \sum_{i=1}^n x_i + \sigma^2 \mu_0}{n\sigma_0^2 + \sigma^2} \\
 &\sigma^2 = \frac{\sigma_0^2 \sigma_0^2}{n\sigma_0^2 + \sigma^2} \\
 &= \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1}
 \end{aligned}$$

8 pts.

(b) Let $\sigma_0^2 = \pi$ and $x_i = 1$ for $i = 1, \dots, 5$. What is the numerical value of the maximum a posteriori estimate of μ ?

$$\begin{aligned}
 \mu_h &= \frac{1}{5\pi + 1} \mu_0 + \frac{5\pi}{5\pi + 1} \left(\frac{1}{5} \sum_{i=1}^5 x_i \right) \\
 &= \frac{\mu_0}{5\pi + 1} + \frac{5\pi}{5\pi + 1} \\
 &= \frac{5\pi + \mu_0}{5\pi + 1}
 \end{aligned}$$

4 pts.

3. **Bonus question:** Let μ_{ML} and μ_{MAP} respectively denote the maximum likelihood estimator and the maximum a posteriori estimator for μ . Calculate the following:

$$\lim_{\sigma_0^2 \rightarrow \infty} \mathbb{E}[\mu_{MAP}] - \mathbb{E}[\mu_{ML}] = ?$$

$$\mathbb{E}\left[\frac{\frac{n}{\sigma^2} \bar{x}_i}{\frac{n}{\sigma^2} + 1}\right]$$

$$\underset{\sigma^2 \rightarrow \infty}{\mathbb{E}} \left[\frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0 + \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \left(\frac{1}{n} \sum_{i=1}^n x_i \right) - \frac{\bar{x}_i}{n} \right]$$

$$= 0$$

5 pts.

Question 3: Supervised Learning

This question is concerned with classification of watermelons into 'good' watermelons (+1) and 'bad' ones (-1). Watermelons can be distinguished based only on their color and smell. Let \mathcal{H} be the class of all *circles* in \mathbb{R}^2 . We associate a classification rule with each $h \in \mathcal{H}$: the interior of the circle is classified as 'good' and outside of the circle is 'bad'.

Given $\{(x_i, y_i)_{i=1}^n | x_i \in \mathbb{R}^2, y_i \in \{1, -1\}\}$, a labeled sample of watermelons, we used the following criterion for the parameters of $h^* \in \mathcal{H}$:

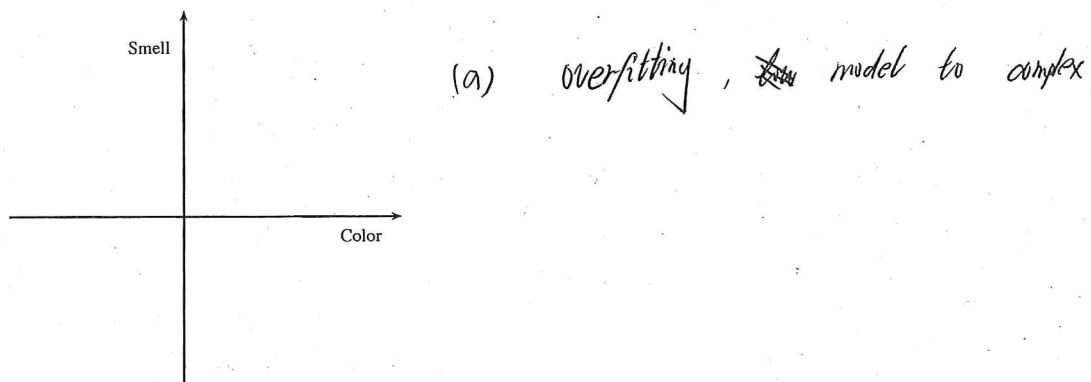
$$w_1^*, w_2^*, r^* = \underset{w_1, w_2, r}{\operatorname{argmin}} \sum_{i=1}^n \exp(-y_i[r^2 - ((x_{i1} - w_1)^2 + (x_{i2} - w_2)^2)]) \quad (2)$$

We then sold h^* to Migros as part of a watermelon test kit.

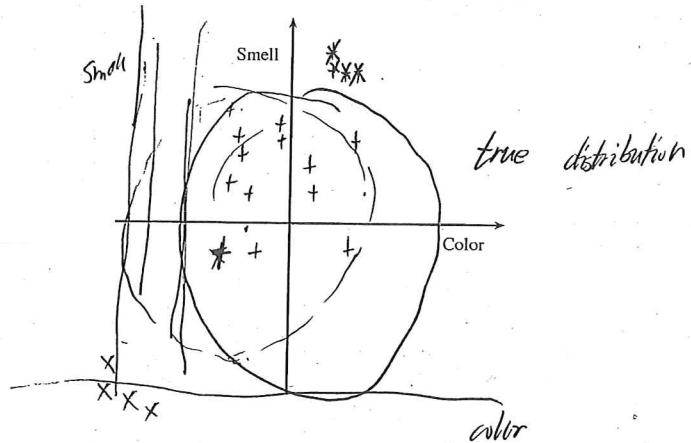
Unfortunately h^* did not meet the expectations, it misclassified a non-negligible proportion of the watermelons used at test time.

1. For each of the following additional assumptions:
 - (a) Give a possible explanation for h^* performing poorly
 - (b) Draw a training set, the prediction function h^* , and the true distribution (if needed) that demonstrate your explanation.

Additional assumption: h^* had a very low training error



Additional assumption: The sample size was large, and h^* had training error of ~ 0.4



8 pts.

2. Suggest a way to measure the empirical variance of the classifier h^* , given that we are out of budget for obtaining more watermelon samples.

use cross-validation K-fold

f $h = 1 \dots k$
from $C_k(x)$ based on the $n - \frac{k}{K}$ dataset

$$\frac{1}{n} \sum_{i=1}^n I(f_{h(i)}(x_i) \neq y_i)$$

the empirical variance

$$\frac{1}{n} \sum_{i=1}^n (y_i - f_{h(i)}(x_i))^2$$

$x_i \in \mathcal{X}_{h(i)}$ fold data set

$$\frac{1}{k} \sum_{h=1}^k \left[\left(w_{1h} - \frac{1}{n} \sum_{i=1}^n w_{ih} \right)^2 + \left(w_{2h} - \frac{1}{n} \sum_{i=1}^n w_{ih} \right)^2 \right] + \left(\frac{1}{K} \sum_{h=1}^K \left(\frac{f_{h(i)}}{K} - \frac{1}{n} \sum_{i=1}^n f_{h(i)} \right) \right)^2 + \left(\frac{1}{K} \sum_{h=1}^K \frac{h}{K} \right)^2$$

3 pts.

Figure 2 depicts the dataset we had, and h^* that we got using equation (2). To improve h^* we decided to add a regularizing term, the new criterion will be

$$\underset{w_1, w_2, r}{\operatorname{argmin}} \sum_{i=1}^n \exp(-y_i[r^2 - ((x_{i1} - w_1)^2 + (x_{i2} - w_2)^2)]) + \lambda \Omega(w_1, w_2, r) \quad (3)$$

Where $\Omega(w_1, w_2, r)$ is the regularizer. We ask you to suggest a suitable regularizer.

$$(w_1 - w_{10})^2 + (w_2 - w_{20})^2 + (r - r_0)^2$$

are the prior knowledge
we have about the shape
of the boundary

4. (a) Draw on Figure 2 the regularized solution you envision.
 (b) Write down the mathematical term of the regularizer. Explain your answer.

$$\Omega(w_1, w_2, r) = (w_1 - w_{10})^2 + (w_2 - w_{20})^2 + (r - r_0)^2$$

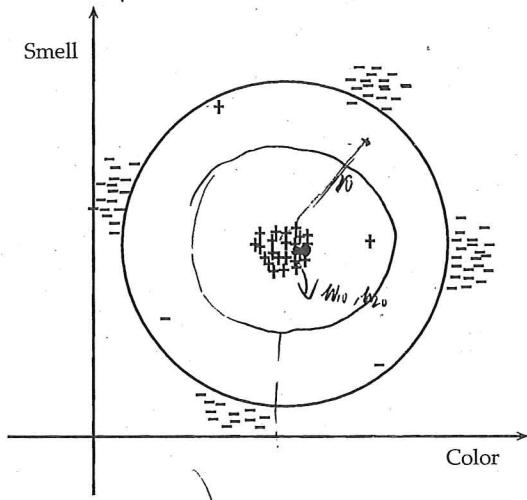


Figure 2: Watermelons dataset and h^*

we assume r_0 not that large
 thus penelize the way to
 make r in the decision boundary large.
 we also made prior assumption
 about the typical characteristic of a good
 watermelon. thus give the w_{10} and
 w_{20}

- 5 pts.
5. Assume that the true distribution of watermelons consists of high density regions visible in Figure 2, plus sparse outliers. Explain what happens to the variance of h^* as we increase λ compared to some starting value $\lambda_0 > 0$.

~~the variance become smaller, it give negative incentive for h^* to be a large circle. thus make r smaller and close to r_0 at the same time~~

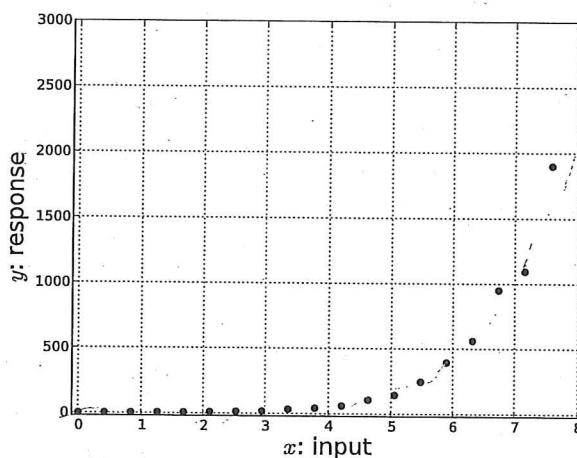
4 pts.

give constraints to the w_1 and w_2 ,
 make it ~~close~~ to a smaller
 distance to w_{10}, w_{20} , after

Question 4: Kernelized Ridge Regression (20 pts.)

Recall the regression setting: Given input vectors \mathbf{x}_i , and output (response) variables y_i , the goal is to find a functional relation between them, often expressed with a weight vector \mathbf{w} and bias b .

- Below is a dataset with one dimensional input variables x , and response variable y . Your task is to find a kernel function $K(x_i, x_j)$, such that you can use a linear regression method in the kernel space.



$$e^{x_1} e^{x_2}$$

$$K(x_i, x_j) = \dots \quad \cancel{(e^{x_i})} \cancel{(e^{x_j})} \quad e^{(x_i + x_j)} \quad 3 \text{ pts.}$$

- You will now derive a kernelized version of ridge regression by introducing a feature transform $\Phi(\mathbf{x}_i)$. This should allow a non-linear regression solution, for datasets such as the one depicted above.

Recall the formulation of ridge regression as an optimization problem:

$$\min_{\mathbf{w}, b} \quad \sum_i (y_i - \mathbf{w}^\top \mathbf{x}_i - b)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (4)$$

- We replace the inputs \mathbf{x}_i in Equation (4) with the vectors of the features in the kernel space and rewrite the problem as a

constrained optimization problem by introducing the new variables ξ_i . Write down the equality constraint in Equation (6).

$$\min_{\mathbf{w}, b, \xi} \sum_i \xi_i^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (5)$$

$$\text{s.t.} \quad \xi_i = \sqrt{y_i - \mathbf{w}^\top \phi(x_i) - b} \quad (6)$$

$$y_i - \mathbf{w}^\top \phi(x_i) - b$$

2 pts.

- (b) Write down the Lagrangian of this new optimization problem using α as the dual variable.

$$L(\mathbf{w}, b, \xi) = \sum_i \rho_i + \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_i \alpha_i (\xi_i - y_i + \mathbf{w}^\top \phi(x_i) + b)$$

- (c) Derive the dual optimization problem.

$$\begin{aligned} & \text{maximize}_{\mathbf{w}, b, \xi} \quad \sum_i \alpha_i^2 - \frac{\sum_i \alpha_i \phi(x_i)}{2\lambda} \\ & \text{subject to} \quad \mathbf{w}^\top \phi(x_i) + b = y_i - \sum_i \alpha_i \end{aligned}$$

$$\frac{\partial L}{\partial \mathbf{w}} = \lambda \mathbf{w} + \sum_i \alpha_i \phi(x_i) = 0$$

$$\mathbf{w} = \frac{-\sum_i \alpha_i \phi(x_i)}{\lambda}$$

$$\frac{\partial L}{\partial b} = \sum_i \alpha_i = 0$$

$$\frac{\partial L}{\partial \xi_i} = 2\rho_i + \alpha_i = 0$$

$$\xi_i = -2\rho_i \quad \rho_i = -\frac{\xi_i}{2}$$

$$L = \sum_i \frac{\xi_i^2}{4} + \frac{\lambda}{2} \frac{\|\sum_i \alpha_i \phi(x_i)\|^2}{\lambda^2} + \sum_i \alpha_i \left(-\frac{\xi_i}{2} - y_i + \frac{-\sum_j \alpha_j \phi(x_j)^\top \phi(x_i)}{\lambda} \right)$$

8 pts.

$$= \sum_i \frac{\xi_i^2}{4} + \frac{\sum_j \alpha_j \phi(x_j)^\top \phi(x_i)}{2\lambda} + -\frac{13}{2} \sum_i \xi_i^2 - \sum_i \xi_i y_i - i \sum_j \alpha_j \phi(x_j)^\top \phi(x_i)$$

$$= -\frac{\sum_i \xi_i^2}{4} - \frac{\sum_j \alpha_j \phi(x_j)^\top \phi(x_i)}{2\lambda} - \sum_i \xi_i y_i$$

3. Express the dual problem in terms of the kernel function $K(\mathbf{x}_i, \mathbf{x}_j)$.

$$W(\alpha) = -\frac{1}{4} \sum_i \alpha_i^2 - \frac{\sum_i \alpha_i K(\mathbf{x}_i, \mathbf{x}_i)}{\sum_i \alpha_i} - \sum_i \alpha_i y_i$$

subject to $\sum_i \alpha_i = 0$ ~~$A \alpha = 0$~~

3 pts.

4. Given the optimal solution of the dual problem α^* and a new point \mathbf{x}_k , write down the equation to compute y_k .

$$y_k = \alpha^* \psi(\mathbf{x}_k) + b$$

~~$b = \frac{1}{n} \sum_i \alpha_i^* y_i$~~

$$b = \frac{1}{n} \left[\sum_i \alpha_i^* y_i - \sum_i \alpha_i^* \psi(\mathbf{x}_i) \right]$$

$$\alpha^* = -\frac{\sum_i \alpha_i^* \psi(\mathbf{x}_i)}{n}$$

4 pts.

$$y_k = \frac{-\sum_i \alpha_i^* K(\mathbf{x}_i, \mathbf{x}_k)}{n}$$

when conditional number ratio of the largest eigenvalue to the least eigenvalue

- f) Model inference (computing model parameters) for the RSS cost function requires the inversion of a matrix: $(\mathbf{X}^T \mathbf{X})^{-1}$.
- Specify a mathematical condition when this inversion is numerically unstable.
 - Describe in your own words under which circumstance this instability happens during a practical application of regression.
 - Are we then in risk of under- or over-fitting?
 - Also comment qualitatively on the bias and variance of the model parameter's estimation in this circumstance.

3.5 pts.

over-fitting

in that case,

variance will be ~~big~~

the bias will be small

no regularization,

when the input data sample are correlated

highly

in the extreme case,

when the data sample \mathbf{x} are the same

which means $x_1 \dots x_N$ are the same

x_1

x_2

x_3

$$(y_i - \sum_{j=1}^{D+1} x_{ij} \beta_j) = 0$$

we can find ∞ infinite number of β who meet this situation.

feature selection

g) We can stabilize the inversion by reducing the model complexity.

- Explain how Ridge Regression (RR) limits the model complexity.
 the $\lambda\beta^T$ factor somehow represents the prior knowledge of the distribution for β , to be $N(0, \lambda^{-1})$.

- ✓ Provide another approach to limiting the model complexity. add constraint for β , just like LASSO.
- Demonstrate the stabilization mathematically by writing down the Ridge Regression solution and argue with the Eigenvalues of the matrix that needs to be inverted.

feature selection

in RR the β would be $(X^T X + \lambda I)^{-1} X^T Y$ 4 pts.

can be express

by $\sum_{j=1}^d \frac{u_j}{\sqrt{\lambda}} \frac{du_j^T}{d\beta} u_j^T Y$ in this case, we add the shrinkage factor which mean the small value will be balance by λ . thus adding stability and limit the

$$\begin{aligned}\hat{\beta} &= X(X^T X + \lambda I)^{-1} X^T Y \\ &= U D (D^2 + \lambda I)^{-1} D U^T Y \\ &= \sum_{j=1}^d u_j \frac{du_j^T}{d\beta} u_j^T Y\end{aligned}$$

h) Comment qualitatively on the bias and variance of Ridge Regression as compared to the RSS estimator.

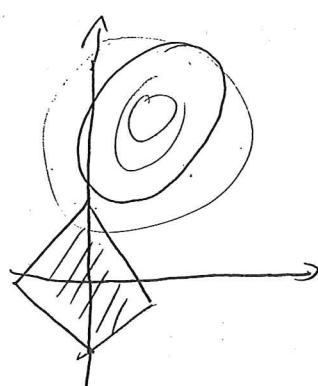
1.5 pts.

the bias becomes larger

while the variance decreases

i) By depicting an appropriate plot (including notation), provide a graphical argument why LASSO favors sparse solutions.

1 pt.



Question 5: Unsupervised Learning (20 pts.)

- a) For each of the following non-parametric approaches mention the most important parameter that influences the smoothness of the results: **3 pts.**

- histograms

of bins (the length of bins)

- Parzen window estimates

size of W_h h_n

- nearest neighbor estimates.

number of K

- b) Determine whether the following statements are true or false. Briefly explain your answer. **2 pts.**

- Non-parametric estimation methods are less sensitive than parametric approaches to model misspecification.

Yes, since non-parametric may does not specify an distribution in advance, putting more emphasis on data

- In histograms, by changing the dimensionality the number of required bins (to keep the resolution) increases linearly with dimension.

No, polynomially, d is the dimension, $(N)^d$

- c) Consider Hidden Markov Models (HMMs). For each of the following algorithms determine if it solves a supervised or an unsupervised problem. Explain your answer.

4 pts.

- Viterbi algorithm

*in supervised learning we only observe
the emitted sequence while the
hidden path is unknown
since we know the evidence variable*

- Baum-Welch algorithm.

unsupervised

- d) In this section we study the k -means clustering method.

1. Mention at least two main differences between k -means and Gaussian Mixture Model (GMM) clustering methods.

k -means doesn't assume the underlying distribution of the data. 2 pts.

*GMM assume the model of the distribution, which is gaussian -
not only give the classification, by also give the
probability of the specific point assigned to different clusters.*

Consider the k -means cost function defined as:

$$R^{km} = \sum_{n=1}^N \sum_{l=1}^k r_{nl} \|\mathbf{x}_n - \boldsymbol{\mu}_l\|^2. \quad (2)$$

Here $\boldsymbol{\mu}_l$ denotes the l -th centroid and $r_{nl} \in \{0, 1\}$ indicates the assignment of object \mathbf{x}_n to the l -th cluster.

2. Write down the assignment update step (E-step) for the k -means algorithm.

$$C^{(t)}(x)_k = \underset{c(x) \in \{1, \dots, k\}}{\operatorname{argmin}} \|x - u_k\|^2$$

2 pts.

3. Derive the centroids update step (M-step). We expect you to write down all intermediate steps of the centroid update derivation.

4 pts.

$$\begin{aligned} D &= \sum_i^{n_d} (x_i - u_k)^2 \\ \frac{\partial D}{\partial u_k} &= -2 \sum_i^{n_d} x_i + n_d u_k = 0 \end{aligned}$$

thus $u_k = \frac{1}{n_d} \sum_{x: \text{assigned}} x$ with $n_d = H(c(x), k)$

until change to $c(x), y$ vanish.

4. Show that the k -means algorithm always converges.

3 pts.

since there are only finite way of assigning sample to k groups

every iteration, the configuration changes
and distortion improves

Supplementary Sheet