

Final Exam

February 6th, 2013

First and last name: _____

ETH number: _____

Signature: _____

General Remarks

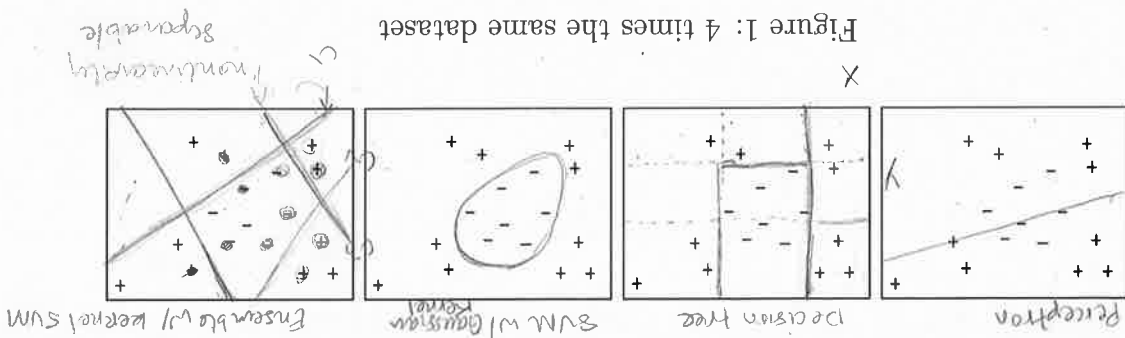
- You have 2 hours for the exam. There are five sections, each of which is worth 20 points. Scoring 100 points guarantees you a grade of six. In two sections you will find bonus questions, worth together 10 points. The bonus questions are a bit more difficult, we suggest you leave them to the end.
- 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:

Question 1: Assorted Questions (20 pts.)

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

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 \hat{f} , a prediction function which is **guaranteed** to have a smaller training error than f .

$$\hat{f}: \text{prediction fn} = \text{train error } 0.65$$

Regularized form of f and perform CV.

2 pts.

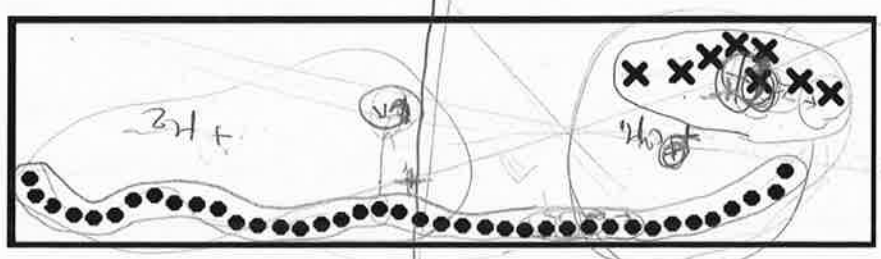
$$\hat{f} \rightarrow f_{\text{pred}}$$

\hat{f} can be evaluated by selecting a cost fn and performing a cross validation on the number of predicted fns to choose the best parameters.
 OR in ensemble methods update the weight of each classifier such as boosting.

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.



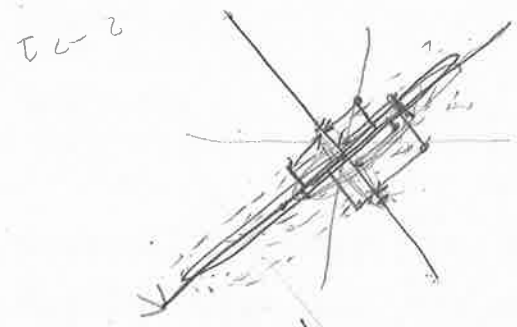
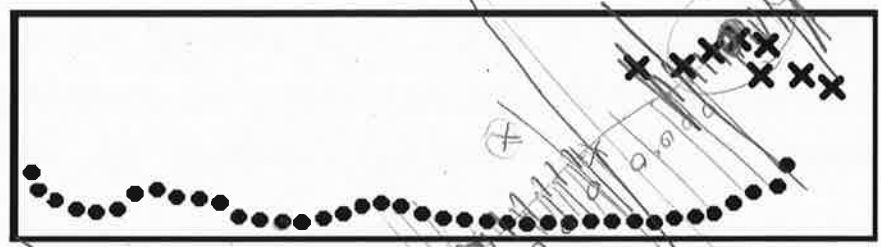
3 pts.

(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.

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

$$\mu_{MAP} = \arg \max_{\mu} P(x_1, \dots, x_n | \mu, \sigma^2) P(\mu | \mu_0, \sigma_0^2)$$

prior

Posterior of μ_{MAP}

$p(x_1, \dots, x_n | \mu_0, \sigma_0^2) \Rightarrow$ can throw away because doesn't depend on μ .

$$= \arg \max_{\mu} p(x_1, \dots, x_n | \mu, \sigma^2) \cdot p(\mu | \mu_0, \sigma_0^2)$$

μ_{MAP} is a linear combination of μ_0 according to ML

$$= \arg \max_{\mu} \left[\frac{1}{n} \sum_{i=1}^n \exp \left(-\frac{(x_i - \mu)^2}{2\sigma^2} \right) \right] \cdot \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp \left(-\frac{(\mu - \mu_0)^2}{2\sigma_0^2} \right)$$

likelihood

multiplying 2 parts together.

scaling doesn't affect max μ value

$$= \arg \max_{\mu} \exp \left(-\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2} - \frac{(\mu - \mu_0)^2}{2\sigma_0^2} \right) = \arg \min_{\mu} \left[\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2} + \frac{(\mu - \mu_0)^2}{\sigma_0^2} \right]$$

such that $\frac{\partial \mu_{MAP}}{\partial \mu} = 0$

$$= 2 \left[\sum_{i=1}^n -\frac{x_i}{\sigma^2} - \frac{\mu_0}{\sigma_0^2} + \frac{\mu_{MAP}}{\sigma^2} + \frac{\mu_{MAP}}{\sigma_0^2} \right] = 0 \Rightarrow \mu_{MAP} = \frac{n\sigma^2}{n\sigma^2 + \sigma_0^2} \left(\frac{1}{n} \sum_{i=1}^n x_i \right) + \frac{\sigma_0^2}{n\sigma^2 + \sigma_0^2} \mu_0$$

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 ?

$$\mu_{MAP} = \frac{5\pi + \mu_0}{5\pi + 1}$$

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

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

cost *center* *radius* *trained classifier*

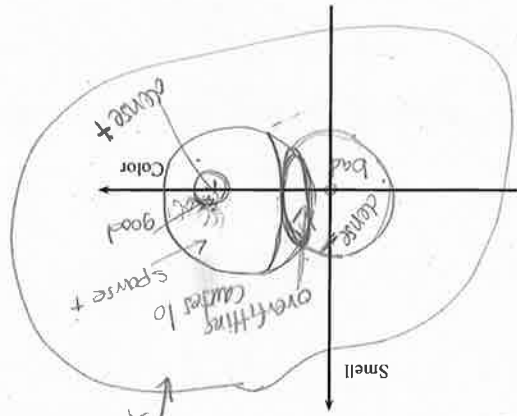
We then sold h^* to Migros as part of a watermelon test kit.

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
 know that the distribution of the watermelons is not 50/50 and you have to consider the additional assumption that getting 100/30 and give you a dense solution for good. Added noise to true distribution



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

$$\Omega(w_1, w_2, r) = R^2 \text{ in } L_2\text{-norm}$$

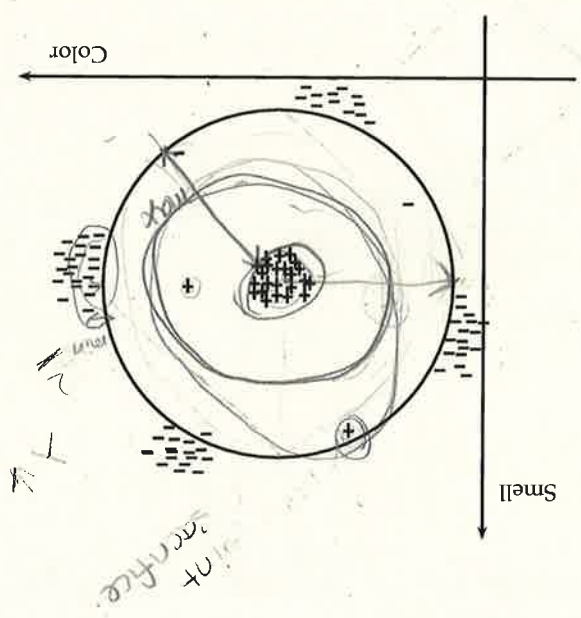


Figure 2: Watermelons dataset and h^*

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 size the circle will decrease.

Variance ↓ 150

4 pts.

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

$$\min_{w, b, \xi} \sum \xi_i^2 + \frac{\lambda}{2} \|w\|^2 \quad \text{primal}$$

$$\text{s.t.} \quad \xi_i = \dots \quad (6)$$

2 pts.

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

objective of Lagrangian is to minimize the problem.

$$\begin{aligned} \max_{f(x)} \quad & L(\tilde{x}, \lambda) = f(\tilde{x}) + \lambda g(\tilde{x}) \\ \text{s.t.} \quad & g(\tilde{x}) = 0 \\ & \tilde{x} \in \mathbb{R}^D \\ & \nabla g(\tilde{x}) \perp g(\tilde{x} + \epsilon) \approx g(\tilde{x}) + \epsilon^T \nabla g(\tilde{x}) \implies \|\epsilon\| \rightarrow 0 \end{aligned}$$

solve the Lagrangian

$$\frac{\partial L}{\partial x} = 0 \quad \frac{\partial L}{\partial \lambda} = g(x) = 0$$

(c) Derive the dual optimization problem.

$$L(\tilde{x}, w, b, \xi) = \sum \xi_i^2 + \frac{\lambda}{2} \|\tilde{w}\|^2 + \sum \alpha_i (y_i - \tilde{w}^T \phi_i - b - \frac{\epsilon}{2})$$

$$\tilde{w} = \begin{pmatrix} \frac{w}{b} \\ \frac{x}{1} \end{pmatrix}$$

$$\frac{\partial L}{\partial \tilde{w}} = -\sum \alpha_i \phi_i + \lambda \tilde{w} = 0 \implies \tilde{w} = \frac{\sum \alpha_i \phi_i}{\lambda}$$

$$\frac{\partial L}{\partial b} = \sum \alpha_i - \alpha_1 = 0 \implies \alpha_1 = \sum \alpha_i$$

$$\|\tilde{w}\|^2 = \frac{\lambda}{2} \cdot \frac{1}{\lambda^2} \left\| \sum \alpha_i \phi_i \right\|^2 = \frac{1}{2\lambda} \left\| \sum \alpha_i \phi_i \right\|^2$$

$$L(\tilde{x}) = -\frac{1}{2} \sum \alpha_i^2 + \sum \alpha_i y_i - \frac{1}{2\lambda} \left\| \sum \alpha_i \phi_i \right\|^2 \implies \min_{\alpha} L(\tilde{x})$$

you want to get the objective f(x) on α .

This is why you derive on α .

8 pts.

Question 5: Unsupervised Learning (20 pts.)

- In this section we study non-parametric density estimation of an arbitrary point x . We consider some small region \mathcal{R} containing x . In the class we have seen the following generic formula for density estimation:

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

where K denotes the number of data points falling inside the region \mathcal{R} and V shows the volume of the region. n is the number of data points in the sample set $\mathcal{S} = \{x_1, \dots, x_n\}$.

- Consider the following Gaussian distribution to be used as a Parzen window function:

$$\phi(x - x_j) = \frac{1}{(2\pi)^{1/2}} \exp\left(-\frac{(x - x_j)^2}{2}\right) \quad (7)$$

What are K and V for this window function?

$$K = \sum_{j=1}^n \phi(x - x_j)$$

$V = 1 \leftarrow$ because a probability distribution

- This particular choice of a window function leads to underfitting. Add a parameter to increase the model complexity.

Add variance parameter, $= h$, controls smoothness

$$\phi\left(\frac{h}{(x - x_j)^2}\right) = \frac{1}{(2\pi)^{1/2}} \exp\left(-\frac{(x - x_j)^2}{2h^2}\right)$$

2 pts.

2. We consider a mixture of K poisson distributions and perform the Expectation-Maximization (EM) algorithm to compute the unknown parameters. The log-likelihood function of n independent objects for mixture of K Poisson distribution is defined as:

$$P(x; \lambda) = \sum_{i=1}^n \log \sum_{c=1}^K \pi_c f(x_i; \lambda_c)$$

where π_c 's are the mixture weights and λ_c 's are the parameters of K Poisson distributions: $f(x_i; \lambda_c)$ is defined as:

$$f(x; \lambda_c) = \frac{\lambda_c^x e^{-\lambda_c}}{x!}$$

(a) Introduce the latent indicator variables necessary for maximizing the log-likelihood function.

$$\log P(x; \lambda, \pi) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k f(x_n; \lambda_k)$$

$z \in \mathbb{R}^{N \times K}$, each row is a vector for a point

$z_{nk} = \begin{cases} 1 & \text{component } k \text{ is responsible for generated } x_n \\ 0 & \text{otherwise} \end{cases}$

Conditional expectation

(b) Calculate the expectation of the latent variables. Provide a

Bayesian interpretation for your answer.

E-step: estimate posterior distribution of data

$$p(z|x, \theta^{old}) = E_z(z|...)$$

$$M\text{-step: } \theta^{new} = \arg \max_{\theta} Q(\theta, \theta^{old})$$

θ

$$Q(\theta, \theta^{old}) = \sum_z p(z|x, \theta^{old}) \cdot \log p(x, z|\theta)$$

joint distribution

$$E[z_{nk}] = p(z_{nk}=1|x_n) = \frac{p(z_{nk}=1) \cdot p(x_n|z_{nk}=1)}{\sum_{k=1}^K p(z_{nk}=1) \cdot p(x_n|z_{nk}=1)} = p(x_n|\lambda_k)$$

responsibility

$$p(z_n|\pi) = \prod_{k=1}^K \pi_k^{z_{nk}}$$

$$\Rightarrow p(x_n|z_n, \lambda_k) = \prod_{k=1}^K p(x_n|\lambda_k)^{z_{nk}}$$

$$\sum_{z_1} p(x_n, z_n) = \sum_{k=1}^K \pi_k \cdot p(x_n|\lambda_k)$$

5 pts.

Supplementary Sheet

NO "fixed" "right solution", as in real life. \rightarrow the accuracy is the only goal.

Google solution for insights and do it better \rightarrow do not reinvent the wheel.
CV scheme is really important

SVM \rightarrow assume vectorial space behind data.

Q/A session before exam.

