

Final Exam

February 10th, 2010

First and Last name: _____

ETH number: _____

Signature: _____

General Remarks

- Please check that you have all 21 pages of this exam.
- Remove all material from your desk which is not permitted by the examination regulations.
- Fill in your first and last name and your ETH number and sign the exam. Place your student ID in front of you.
- You have 2 hours for the exam. There are five questions, where you can earn a total of 150 points. Scoring 120 points (equivalent to solving four questions) guarantees you a grade of six.
- Write your answers directly on the exam sheets. If you need more space, put your name and ETH number on top of each supplementary sheet.
- 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 Achieved	Visum
1	Density Estimation	30	
2	PCA	30	
3	Classifiers	30	
4	Regression	30	
5	Boosting	30	
Total		150	

Grade: _____

Question 1: Parametric Density Estimation (30 pts.)

When estimating densities from data, we assume that the data distribution is well approximated by a parametric model $P(x; \theta)$, where θ denotes the parameters of the model. Estimation of the density then reduces to estimation of the parameter θ .

a) In maximum-likelihood estimation (MLE), we usually start by assuming that the data is given as an i.i.d. sample $x_1, x_2, \dots, x_N \sim P(x; \theta)$. Please explain how this assumption simplifies the calculations.

Answer: ① i.i.d. allows us to estimate a probability for each event independent from other giving a product

$$P(x|\theta) = \prod_{i=1}^N P(x_i|\theta) = P(x_1|\theta) \cdot P(x_2|\theta) \cdot \dots \cdot P(x_N|\theta)$$

∴ This assumption would simplify to determine a mean and variance for the underlying distribution

b) We model a coin flip as a random variable X , where $X = 1$ denotes the outcome 'heads', and $X = 0$ denotes 'tails'. We model this event using the Bernoulli distribution, parametrized by μ that denotes $P(X = 1)$.

$$\text{Bern}(x|\mu) = \mu^x (1 - \mu)^{1-x}$$

Given a data set $\mathcal{X}_N = \{x_1, \dots, x_N\}$ of N i.i.d. observations of X , show us step by step how to compute the MLE for μ . Comment your calculations where necessary.

Answer:

Likelihood: $P(x|\theta) = \prod_{i=1}^N \text{Bern}(x_i|\mu) = \prod_{i=1}^N \mu^{x_i} (1 - \mu)^{1-x_i} = \mu^{\sum x_i} (1 - \mu)^{N - \sum x_i}$

Take log of likelihood fn: $\log P(x|\theta) = \log \left(\prod_{i=1}^N \text{Bern}(x_i|\mu) \right) = \sum_{i=1}^N \log \left(\mu^{x_i} (1 - \mu)^{1-x_i} \right)$

$$P(x|\mu) = \prod_{i=1}^N \left[\log[\mu^{x_i}] + \log[(1 - \mu)^{1-x_i}] \right]$$

Take derivative w.r.t to μ : $\nabla_{\mu} \log P(x|\mu) = \nabla_{\mu} \left[\sum_{i=1}^N x_i \log \mu + (1 - x_i) \log (1 - \mu) \right]$

$$= \sum_{i=1}^N x_i \cdot \frac{1}{\mu} + (1 - x_i) \cdot \frac{-1}{1 - \mu} = \sum_{i=1}^N \frac{x_i}{\mu} - \sum_{i=1}^N \frac{1 - x_i}{1 - \mu} = 0$$

8 pts.

$$\sum_{i=1}^N \frac{x_i}{\mu} = \sum_{i=1}^N \frac{1 - x_i}{1 - \mu} \Rightarrow \frac{1}{\mu} \sum_{i=1}^N x_i = \frac{1}{1 - \mu} \sum_{i=1}^N (1 - x_i)$$

$$\Rightarrow \frac{1}{\mu} \sum_{i=1}^N x_i = \frac{1}{1 - \mu} \left(\sum_{i=1}^N 1 - \sum_{i=1}^N x_i \right) \Rightarrow \frac{1}{\mu} \sum_{i=1}^N x_i = \frac{1}{1 - \mu} (N - \sum_{i=1}^N x_i)$$

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

$$\Rightarrow \frac{1}{1 - \mu} - 1 = \frac{\sum_{i=1}^N x_i}{N - \sum_{i=1}^N x_i} - 1 = \frac{\sum_{i=1}^N x_i - (N - \sum_{i=1}^N x_i)}{N - \sum_{i=1}^N x_i} = \frac{2 \sum_{i=1}^N x_i - N}{N - \sum_{i=1}^N x_i}$$

$$\Rightarrow \frac{2 \sum_{i=1}^N x_i - N}{N - \sum_{i=1}^N x_i} = 0 \Rightarrow 2 \sum_{i=1}^N x_i - N = 0 \Rightarrow \sum_{i=1}^N x_i = \frac{N}{2}$$

c) We flip our coin three times and obtain the dataset $\mathcal{X}_3 = \{1, 1, 1\}$.

1. What is the MLE for $\hat{\mu}$ for dataset \mathcal{X}_3 ? $n=3$

Answer:

$$\hat{\mu} = \frac{\sum_{i=1}^n x_i}{n} = \frac{1+1+1}{3} = 1$$

2. State one problem with this estimate. What is a possible remedy (solution)?

Answer:

Problem: This estimate is biased. $\mu = 1$ versus $\mu = \frac{1}{2}$.
The MLE gave you a bias of $\mu = 1$ unfair coin toss.
High variance associated w/ a small population of values.

Remedy: you more samples to get a better, less unfair estimate unless you take many more trials

3 pts.

d) Assume that there are a total of m heads in \mathcal{X}_N , i.e., $\sum_{i=1}^N x_i = m$.

1. In terms of μ , N and m only, what is the probability of obtaining this dataset $P(\mathcal{X}_N | \mu)$?

Explain the terms of your solution. coming from the same dataset \mathcal{X}_N

Answer:

$$P(\mathcal{X}_N | \mu) = \prod_{i=1}^N \mu^{x_i} (1-\mu)^{1-x_i} = \mu^m (1-\mu)^{N-m}$$

$$\log P(\mathcal{X}_N | \mu) = \sum_{i=1}^N \log \mu^{x_i} (1-\mu)^{1-x_i} = \sum_{i=1}^N x_i \log \mu + \sum_{i=1}^N (1-x_i) \log (1-\mu)$$

2. What is the expected value of m ? (Hint: For independent random variables, the expectation of a sum is the sum of expectations)

Answer:

$$E[m] = E\left[\sum_{i=1}^N x_i\right] = \sum_{i=1}^N E[x_i] = \sum_{i=1}^N \mu = N\mu$$

probability of getting exactly $N-m$ tails

$$\mu^m (1-\mu)^{N-m}$$

$$\sum_{i=1}^N \log \mu^{x_i} (1-\mu)^{1-x_i}$$

Exponential property

$$\sum_{i=1}^N x_i = m$$

$$m = N\mu$$

$$\mu = \frac{m}{N}$$

$$\sum_{i=1}^N \log \mu^{x_i} (1-\mu)^{1-x_i} = \sum_{i=1}^N x_i \log \mu + \sum_{i=1}^N (1-x_i) \log (1-\mu)$$

$$\sum_{i=1}^N \log \mu^{x_i} (1-\mu)^{1-x_i} = \sum_{i=1}^N x_i \log \mu + \sum_{i=1}^N (1-x_i) \log (1-\mu)$$

5 pts.

Here, we use the (conjugate) prior given by the Beta distribution:

$$\text{Beta}(n|\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} n^{\alpha-1} (1-n)^{\beta-1},$$

where $\Gamma(x)$ denotes the gamma function (not necessary). Some useful properties:

$$\frac{g+x}{x} = (n) \mathbb{E}(u)$$

$$\frac{(1+\beta+x)(\beta+x)}{\beta^2} = \text{2. Variance } W(n)$$

e) With the mean and variance of μ , we can understand how the parameters influence the prior. Select the best (α_0, β_0) pair for the prior that emphasizes the following beliefs.

1. My coin is biased towards tails.

1. My coin is biased towards tails.

2. My coin is most probably biased (but I don't know heads or tails)

high-variance, low bias
not much info about
the system)

(f) Show that if the prior is Beta distributed

$$P(\mu) = \text{Beta}(\alpha_0, \beta_0),$$

then the posterior takes again the same functional dependency on μ as the prior, i.e.,

$$P(\mu|\chi_N)P(\mu) = P(\mu|\chi_N) \cdot 1_{I_{\beta_0+N-m-1}} \propto \mu^{\alpha_0+m-1} 1_{I_{\beta_0+N-m-1}} P(\mu|\chi_N)$$

$\beta_{\text{eta}} \Rightarrow \text{prior likelihood} \rightarrow \text{Bern}$

Answer: $\beta_{\alpha_0}(\mu | \alpha_0, \beta_0) = \frac{\Gamma(\alpha_0 + \beta_0)}{\Gamma(\alpha_0) \Gamma(\beta_0)} \mu^{\alpha_0 - 1} (1 - \mu)^{\beta_0 - 1}$

$$P(\mu | \chi_n) = \mu_m^{(1-\mu)N-m} \mu^{(1-\mu)N-m-1} \mu^{(1-\mu)N-m-1}$$

$$P(\mu | \chi_n) = \mu^{m + \alpha_0 - 1} (1 - \mu)^{p_0 + N - m - 1}$$

செ.பி.பி.

$$\text{mean} = \frac{0.5}{0.5 + 0.5} = \frac{1}{2}$$

$$\text{mean} = \frac{2}{2+2} = \frac{2}{4}$$

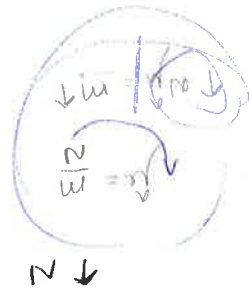
212

5 pts.

$$\text{var.} \cdot \frac{(1)^2 (2)}{0.25} = \frac{2}{0.25} = \frac{1}{\frac{1}{4}} = \frac{1}{8} \quad \left(\begin{matrix} 1 \\ 8 \end{matrix} \right)$$

As posterior $N \rightarrow \infty$, then becomes MLE equation.

$$E[W] \leftarrow (n/2) \phi = (n-1) \phi$$


$$\infty \leftarrow N$$

Jim

$$\left[\frac{1}{(n-1)!} \left(\frac{d}{dx} \right)^{n-1} \left(\frac{1}{x} \right) \right]_{x=0} = \frac{1}{(n-1)!} \left(\frac{d}{dx} \right)^{n-1} \left(\frac{1}{x} \right) \Big|_{x=0}$$

Answer:

the MLE?
 life is truly biased
 become
 if you need a
 Bayesian approach

2. What happens to the posterior as $N \rightarrow \infty$? How does this relate to the MLE?

with the prior information of a fair coin, this tells us

$$v = \frac{r}{v} \leftarrow \text{all W}$$

Fair coin: $P(X=1|X) = \frac{3+\alpha_0}{3+\alpha_0+\beta_0}$
 $\beta_0 = \beta_1 = 1$
 $\frac{1}{2} = \frac{\alpha}{\alpha + \beta}$
 $\alpha + \beta = 2$
 $\beta = 2 - \alpha \Rightarrow \beta = \alpha$
 $1 = 1$
 MLE \downarrow Posterior \downarrow Prior
 $1 - 0.8 - 0.5$

Answer:

1. For $\mathcal{X}_3 = \{1, 1, 1\}$, use a choice of (α_0, β_0) which assumes a fair coin, and obtain your estimate for $P(X = 1 | \mathcal{X})$. How does this estimate differ from the MLE?

(g) In fact, the posterior is again a Beta distribution, and we obtain $P(X=1|\mathcal{X}) = \frac{m+\alpha_0}{n+\alpha_0+\beta_0}$

$$\frac{0g' + 0x + N}{0x + u}$$

$$\tau_T = 10 \text{ ns}$$

Question 2: Principal Component Analysis (PCA) (30 pts.)

Principal component analysis is a widely applied method in data analysis and dimensionality reduction. Given a dataset $X \in \mathbb{R}^{D \times N}$ (observations as columns), where D is the number of dimensions and N is the number of observations, a linear transformation using an orthonormal matrix $P \in \mathbb{R}^{M \times D}$ is applied to make a *change of basis*, to obtain a (usually) lower-dimensional dataset $Y \in \mathbb{R}^{M \times N}$.

a) We begin by reviewing the steps of applying PCA to a dataset X . Please complete each step below by providing the appropriate formula to compute the desired quantity.

1. Define the zero-mean dataset \tilde{X} in terms of the original dataset X . For this purpose, introduce the data mean as a column-vector $\mu \in \mathbb{R}^{D \times 1}$.

Answer:

2. Define the covariance matrix C_X in terms of the zero-mean dataset \tilde{X} .

Answer:

3. Write down the eigen-decomposition of the covariance matrix C_X , in terms of the eigenvector matrix E (eigenvectors as columns) and the diagonal matrix of eigenvalues D .

Answer:

4. Define the PCA transformation matrix P in terms of the eigenvector matrix E . (Note: assume we don't want to reduce the dimensionality of the transformed dataset. Further assume that the eigenvectors in E have already been sorted according to the corresponding eigenvalues, in decreasing order.)

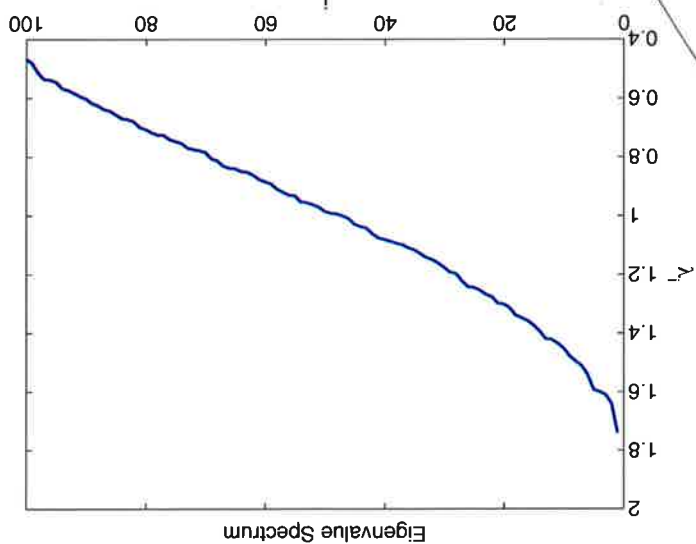
Answer:

5. Define the transformed dataset Y in terms of the original dataset X and the transformation matrix P .

Answer:

6 pts.

b) Assume we have applied PCA to some dataset ($D = 100$). We observe the following eigenvalue spectrum of the covariance matrix of the data. (λ_i : eigenvalues)



1. Is the intrinsic dimensionality of this dataset low or high? Why?

Answer:

2. Can this dataset be expressed in few dimensions with low approximation error? Why?

Answer:

3. If yes, which dimensionality (approximately) should be chosen for the transformed dataset and why?

Answer:

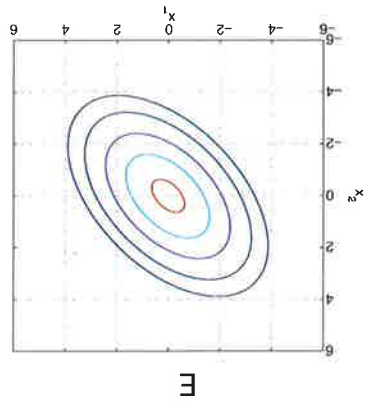
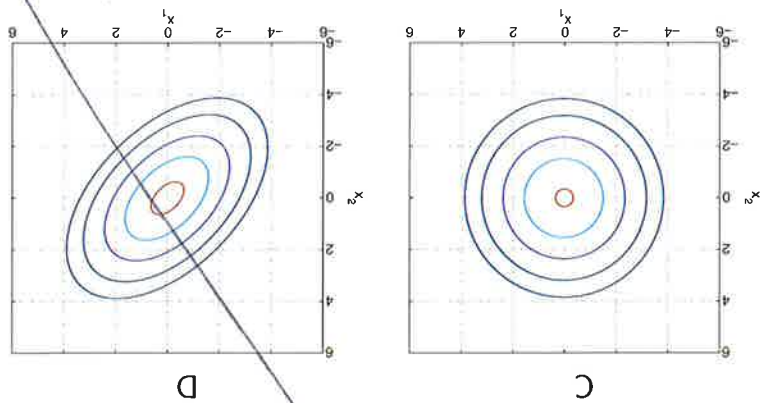
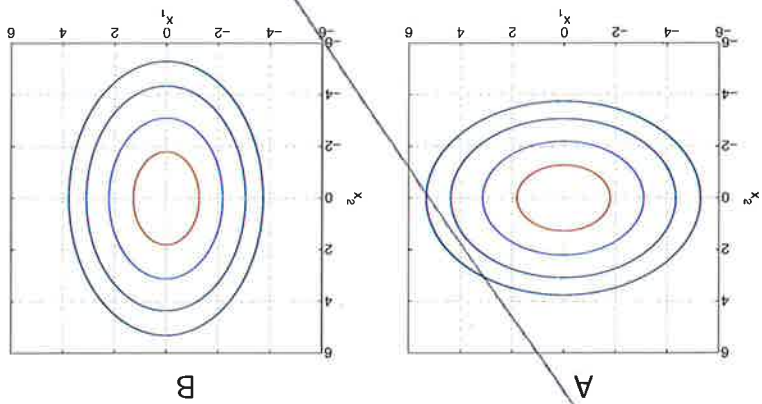
5 pts.

c) Assume you have observed 2D data $\mathbf{X} \in \mathbb{R}^{2 \times N}$ (observations as columns). The first row of \mathbf{X} corresponds to the first dimension x_1 , the second row corresponds to x_2 . For each of the three covariance matrices \mathbf{C}_X below, please choose the iso-line plot (A-E) corresponding to the covariance matrix. (Note the axis labels on the figures)

1. $\begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$ Answer: ()

2. $\begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$ Answer: ()

3. $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ Answer: ()



- d) PCA can be derived as follows: Given a dataset $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \in \mathbb{R}^{D \times N}$, find a new set of basis vectors such that the variance of the projected dataset is maximized. We begin by finding the first projection direction $\mathbf{e}_1 \in \mathbb{R}^D$, $\|\mathbf{e}_1\|_2 = 1$, which satisfies our goal.
1. Please describe in words why we are interested to find a projection direction such that the projected data has large variance.

Answer:

2. Write down formally the objective function for maximizing the variance of the projected dataset. For this purpose, define the variance of the projected dataset in terms of the original data \mathbf{x}_n and the first projection direction \mathbf{e}_1 . (Hint: introduce the mean μ of the data \mathbf{X} .)

Answer:

5 pts.

e) PCA transforms a dataset X into a dataset Y by defining a new basis using the eigenvectors of the covariance matrix C_X of the dataset X . With this particular choice of a new basis, the covariance matrix C_Y of the transformed dataset Y is *diagonalized*. (Note: For this exercise, assume a zero-mean dataset.)

1. Please explain in words, why we desire the covariance matrix of the transformed dataset to be diagonal.

Answer:

2. Show that $C_Y = PC_X P^T$, i.e., that the covariance matrix C_Y of the transformed dataset can be written in terms of the covariance matrix C_X of the original dataset.

Answer:

3. Show that the choice $P = E^T$ for the PCA transformation matrix, where E is the matrix of eigenvectors of the covariance matrix C_X , actually diagonalizes the covariance matrix C_Y of the transformed dataset Y . Use the eigen decomposition of C_X and the fact that $P^{-1} = P^T$.

Answer:

Question 3: Linear Classifiers (30 pts.)

The simplest representation of a linear discriminant function is obtained by taking a linear function of the input vector x so that

$$y(x) = w^T x + b$$

a) Define the term decision boundary. Explain what is the relation between the vector w and the decision boundary.

Answer: A decision boundary is always \perp to the \vec{w} .
 Decision boundary: A hyperplane separating different classes (linearly separable). (points)

b) Consider the following set $S \subseteq \mathbb{R}^2 \times \{+1, -1\}$:

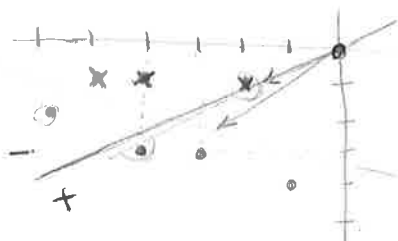
$$S = \{((3, 3), +1), ((1, 4), +1), ((2, 1), -1), ((5, 1), -1), ((6, 2), +1), ((4, 3), +1), ((4, 1), -1)\}$$

Prove that there exist no homogeneous half-space that separates S with no errors (recall that a half-space is homogeneous if its separating hyperplane passes through the origin). Assume that the classification of a point is given by $y(x) = \text{sign}(w^T x + b)$.

Answer:

3 points

$$\begin{aligned} x_1 & (2, 1), -1 \\ x_2 & (4, 3), +1 \\ x_3 & (6, 2), +1 \end{aligned}$$



$$y(x_1) = (3, 2) \cdot \begin{pmatrix} 2 \\ 1 \end{pmatrix} + 0 = \text{sign}(6 + 2) = +$$

$$y(x_2) = (3, 2) \cdot \begin{pmatrix} 4 \\ 3 \end{pmatrix} + 0 = \text{sign}(12 + 6) = +$$

5 pts.

The simple form of the online perceptron algorithm is defined as follows:

Initialize by setting $w^{(0)} = 0$.

$$\text{Update } w^{(t+1)} = \begin{cases} w^{(t)} & y_k(w^{(t)T} x_k) > 0 \\ w^{(t)} + y_k x_k & \text{otherwise.} \end{cases}$$

c) We would like to use the perceptron algorithm to find a separating half-space that makes no errors on the set S . Recall that you have just proven that this is impossible using a homogeneous half-space. Make the necessary adjustments and perform 7 iterations of the perceptron algorithm. For each iteration write down the training example and the current vector $w^{(t)}$. Assume that the points are given in the same order as they appear in S .

Answer:

Necessary adjustment is homogeneous coordinates

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}	x_{12}
1	1	4	1	1	3	3	1	1	0	0	0	0
2	2	1	1	1	3	3	1	1	0	0	0	0
3	5	1	1	1	3	3	1	1	0	0	0	0
4	6	2	1	1	3	3	1	1	0	0	0	0
5	4	3	1	1	3	3	1	1	0	0	0	0
6	2	2	1	1	3	3	1	1	0	0	0	0
7	1	4	1	1	3	3	1	1	0	0	0	0
8	2	1	1	1	3	3	1	1	0	0	0	0
9	5	1	1	1	3	3	1	1	0	0	0	0
10	6	2	1	1	3	3	1	1	0	0	0	0
11	4	3	1	1	3	3	1	1	0	0	0	0
12	2	2	1	1	3	3	1	1	0	0	0	0

5 pts.

$$\max \frac{1}{\|w\|} \rightarrow \min \|w\|$$

Consider the primal form of the soft margin SVM

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

s.t. $y_i(w^T x_i + b) \geq 1 - \xi_i, i = 1, \dots, n$

Distance of two lines = 1 (margin)

regularization = 1 (loss)

maximize margin

minimize mistakes over training examples

regularization term

allow error in classification

mistakes over training examples

d) Indicate which of the following statements hold as we increase the value of the parameter C (relative to any starting value, $C > 0$). Use

- D - if the validity of the statement depends on the situation,
- T - if the statement is necessary true
- N - if the statement is never true

$C = 1000 \rightarrow$ wide margin SVM
 $C = 10 \rightarrow$ wide margin
 $C = 0.1 \rightarrow$ widest margin
 $C = 0 \rightarrow$ widest margin
 $C = 0.1 \rightarrow$ widest margin
 $C = 1000 \rightarrow$ wide margin SVM
 $C = 10 \rightarrow$ wide margin
 $C = 0.1 \rightarrow$ widest margin
 $C = 0 \rightarrow$ widest margin

not in the constraint which is fixed.
 $\|w\|$ increases
 $\|w\|$ will not decrease
 $\|w\|$ will not decrease

- N [more points will be misclassified]
- T [The geometric margin will not increase]

10 pts.

12

4 pts.

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}$$

$$= (\alpha_1 y_1 (x_1 x_1) + b) + (\alpha_2 y_2 (x_2 x_2) + b) + (\alpha_3 y_3 (x_3 x_3) + b)$$

Answer: $x \cdot w + b = \sum_{i=1}^n \alpha_i y_i (x \cdot x_i) + b$ $K(x, z) = 1 + (x^T z)^2$

g) Use these values to compute the kernelized discriminant function $y(x)$.

$$\alpha_1 = 0.1, \alpha_2 = 0, \alpha_3 = 0.1, b = 1.5$$

Below are the α 's and b values

$$x_1 = (1, 1), y_1 = +1, x_2 = (3, 0), y_2 = -1, x_3 = (-2, 1), y_3 = -1.$$

We solved the dual SVM problem with the above kernel. We used the following training set

2 pts.

$$w^* = \sum_{i=1}^n \alpha_i \phi(x_i)$$

Answer:

f) Using the dual variables $\bar{\alpha}$'s and $\bar{\phi}(x)$, give an explicit form for the primal variable w .

2 pts.

$$\phi(x) = (1, x_1, x_2, \sqrt{2}x_1x_2, x_2^2)$$

$$\phi(x)^T \phi(z) = (1, x_1, x_2, \sqrt{2}x_1x_2, x_2^2)^T \cdot (1, z_1, z_2, \sqrt{2}z_1z_2, z_2^2) = 1 + ((x_1z_1)^2 + 2x_1z_1x_2z_2 + (x_2z_2)^2)$$

Answer: $K(x, z) = 1 + (x_1z_1 + x_2z_2)^2$

e) What is the feature representation $\phi(x)$ corresponding to this kernel for $x \in \mathbb{R}^2$?

We would like to solve the SVM problem using the kernel $K(x_i, x_j) = 1 + (x_i^T x_j)^2$.

$$\max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j) \quad \text{s.t.} \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n$$

The dual soft margin SVM is given below

Question 4: Regression (30 pts.)

In univariate regression analysis, we study the statistical dependency of the output variable $y \in \mathbb{R}$ on the input variable $x \in \mathbb{R}$.
 a) Complete the equation

$$y = \sum_{i=1}^n w_i \phi_i(x) + w_0$$

nonlinear regression

where the right hand side consists of two terms. Define each term in the equation.
 Note: Be general, don't restrict yourself to linear least-squares regression.

Answer:

$\phi_i(x)$ are basis fns
 w_i = weights/coefficient

w_0 = offset

$$E[w^T x + w_0 + \epsilon] = 0$$

$\epsilon \sim (0, \sigma^2)$

3 pts.

b) Regression analysis is only meaningful if there actually is a true dependency between x and y . Prove that if x and y are statistically independent, there is no linear dependency:

1. Assume that x and y are statistically independent.

2. Begin with the definition of covariance,

$$\text{cov}(x, y) = E[(x - \mu_x)(y - \mu_y)], \quad f: x \rightarrow y$$

where E is the expectation operator and μ_x is the mean of x .

3. Show that x and y are uncorrelated.

Note: Comment your calculations where necessary.
 mean of x .

Answer:

$$\begin{aligned} E[(x - \mu_x)(y - \mu_y)] &= E[xy - x\mu_y - y\mu_x + \mu_x\mu_y] \\ &= E[xy] - E[x\mu_y] - E[y\mu_x] + E[\mu_x\mu_y] \\ &= E[xy] - \mu_y E[x] - \mu_x E[y] + \mu_x\mu_y \\ &= E[xy] - \mu_y\mu_x - \mu_x\mu_y + \mu_x\mu_y \\ &= E[xy] - \mu_x\mu_y = 0 \end{aligned}$$

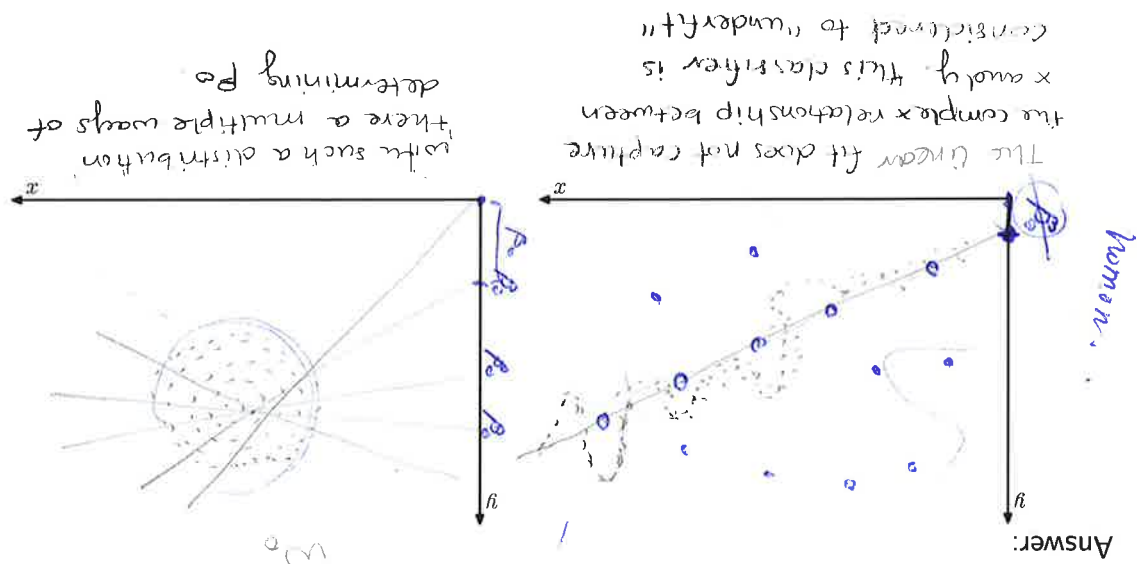
independent

Covariance is 0, which means uncorrelated.

5 pts.

c) You have shown that eq. (1) consists of two terms. Illustrate two conceptually different cases where a linear least-squares fit is inappropriate for the true dependency between x and y , one related to the first term in eq. (1), and one related to the second term in eq. (1):

1. Draw a scatter plot of a data sample.
2. Fit the linear least-squares solution into the sample (qualitatively).
3. Explain why the solution does not capture the true dependency between x and y .



4 pts. So far, we considered a single input variable x . In multivariate linear regression, we assume a weighted linear functional dependency

$$y = w_0 + w_1 x_1 + \dots + w_D x_D$$

(2) between D input variables $x_1, \dots, x_D \in \mathbb{R}$ and output variable y .
 Define the input data matrix \mathbf{X} , weight vector \mathbf{w} and output data vector \mathbf{y} such that eq. (2) for all N data points of the sample $(x_1^n, x_2^n, \dots, x_D^n, y^n), n = 1, \dots, N$ can be written as

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

Answer:

$$\mathbf{X} := \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1D} \\ x_{21} & x_{22} & \dots & x_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{ND} \end{bmatrix} \quad \begin{matrix} D+1 \\ N \times (D+1) \end{matrix}$$

$$\mathbf{w} := \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix} \quad \begin{matrix} D+1 \\ D \times 1 \end{matrix}$$

$$\mathbf{y} := \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad \begin{matrix} N \times 1 \end{matrix}$$

3 pts.

The linear least-squares prediction \hat{y} is the orthogonal projection of y onto the space spanned by the columns of X , and can be written as

$$\hat{y} = X(X^T X)^{-1} X^T y \Rightarrow \hat{y} = W^T X \Rightarrow W = (X^T X)^{-1} X^T y$$

e) Analytically derive P , where \hat{y} is the least-squares prediction vector. To achieve this:

1. Derive the optimal weight vector \hat{w} that minimizes the least-squares objective function

$$\|Xw - y\|_2^2 \quad (3)$$

Note: we have $\frac{\partial}{\partial w} w^T S w = 2S w$ for symmetric matrix S , and $\frac{\partial}{\partial w} w^T a = a$.

2. Using your result for \hat{w} , give the definition for P .

$$\text{Answer: } \arg \min_w \|Xw - y\|_2^2 = (Xw - y)^T (Xw - y)$$

$$= w^T X^T X w - w^T X^T y - y^T X w - y^T y \quad \frac{d}{dw} w^T X^T X w - w^T X^T y - y^T X w = 2X^T X w - 2X^T y = 0$$

$$0 = 2X^T X w - 2X^T y$$

$$X^T X w = X^T y$$

$$w = (X^T X)^{-1} X^T y$$

$$y = X^T P y$$

$$X = n \times d$$

$$w = d \times 1$$

$$y = n \times 1$$



f) Prove that P is an orthogonal projector, by showing that $P = PP$ and $P = P^T$.
 Answer: $(X(X^T X)^{-1} X^T)^T = (X^T)^T (X(X^T X)^{-1})^T = X(X(X^T X)^{-1})^T = X(X^T X)^{-1} X^T = P$

$$P^T = (X(X^T X)^{-1} X^T)^T = X^T (X(X^T X)^{-1})^T = X(X(X^T X)^{-1})^T = X(X^T X)^{-1} X^T = P$$

2 pts.

g) Computing \hat{w} involves the inversion of the symmetric matrix $S := X^T X$, which we assume to be positive definite. Give a mathematical condition (in terms of the eigenvalues of S), when this inversion is numerically unstable.

Answer: Defn of a positive definite matrix \rightarrow eigendecomposition AA^T diagonal elements, eigenvalues $\lambda_i > 0$ positive and real, if not

then inversion is unstable.

All its eigenvalues are

positive

$$\begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \end{bmatrix}$$

2 pts.

h) Ridge regression finds the regularized weight vector \hat{w}_{ridge} that minimizes eq. (3) with an additional norm penalty,

$$\|Xw - y\|_2^2 + \lambda \|w\|_2^2.$$

Explain why choosing $\lambda > 0$ improves stability of the inversion:

choosing $\lambda > 0$ will

shrink w 's to not grow so big.

1. Derive the regularized optimal weight vector \hat{w}_{ridge} .

Answer:

$$\arg \min_w \|Xw - y\|_2^2 + \lambda \|w\|_2^2 \Rightarrow (w^T X^T X w - w^T X^T y - y^T X w + \lambda w^T w) \frac{d}{dw}$$

$$0 = 2X^T X w - 2X^T y + \lambda w$$

$$2X^T y = 2X^T X w + \lambda w$$

$$2X^T y = w(2X^T X + \lambda)$$

$$w_{\text{ridge}} = (X^T X + \frac{\lambda}{2})^{-1} X^T y$$

2 pts.

2. What is the effect of increasing λ on

(a) the norm $\|\hat{w}_{\text{ridge}}\|_2$ and

(b) the numerical stability of computing \hat{w}_{ridge} ?

Provide arguments for your answers.

Answer:

a) increasing λ will shrink w coefficients to not grow so big and prevent overfitting

b)

4 pts.

Question 5: Bagging and AdaBoost (30 pts.)

Bagging and AdaBoost are two ensemble methods used frequently in classification.

a) State two essential differences between bagging and AdaBoost.

Answer:

Bagging

choice of data
varies the training sets using
resampling to train weak learners
gives the same importance to every
prediction

AdaBoost

Trains weak learners with the same training set
on every iteration
weights the prediction of every weak classifier
according to its accuracy

2 pts.

b) Explain how bagging reduces the variance of an estimator?

Answer:

bagging uses the method of boot strap to create different training sets
from the original distribution. in order to reduce the variance, this
simulates increasing sample size.

3 pts.

c) State two factors that may influence the performance of AdaBoost?

Answer:

1) if the number of outliers is large, then the emphasis placed
on these outliers can become large.

2 pts.

d) How can we use AdaBoost to detect outliers (misclassified or ambiguously labeled examples)?

Answer:

By finding the classifier with the highest
weight

2 pts.

2 pts.

performance.

AdaBoost assigns a weight to each classifier based on

but instead place weights on training samples.

Answer: The don't directly place coefficients with each weak learner,

weak learners?

f) What is the difference between AdaBoost and this algorithm in the way they combine the

Let $W_b = \sum_{n=1}^N w_{b,n}$ be the sum of the weights at the end of iteration b .

$$\hat{c}_B(x) = \text{sgn}\left(\sum_{b=1}^B c_b(x)\right).$$

The ensemble classifier decides by combining the weak learners as follows:

- The algorithm finds a new weak learner c_b on the weighted training set that is guaranteed to have an error of at most $\frac{1}{2} - \gamma$ on the weighted training set.
- We then update the weights of the training samples:
 - If x_n is misclassified by c_b , i.e., if $c_b(x_n) \neq y_n$: Set $w_{b,n} := \beta w_{b-1,n}$.
 - Otherwise the weight remains unchanged: $w_{b,n} := w_{b-1,n}$.

For each iteration $b = 1$ to B :

ated weight $w_{b,n}$. At the beginning, $w_{0,n} = 1$ for $n = 1, \dots, N$. Let $0 < \gamma < \frac{1}{2}$ and $\beta = \frac{1-2\gamma}{1+2\gamma}$.
 We consider a dataset $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ of N i.i.d. samples, where $y_n \in \{-1, +1\}$. After each iteration b of this algorithm, each training sample x_n has an associ-

Consider the following simplified ensemble learning algorithm:

4 pts.

For a binary classification problem (± 1) , assume we have c_1, \dots, c_B successive base classifiers obtained via AdaBoost, and let $\alpha_1, \dots, \alpha_B$ be the corresponding weights. The ensemble classifier is $\hat{c}_B(x) = \text{sgn}\left(\sum_{b=1}^B \alpha_b c_b(x)\right)$.
 Reusing the same base classifiers c_1, \dots, c_B , we now train a linear classifier, by finding new $\tilde{\alpha}_b$'s that minimize the average exponential loss on the training examples, and obtain $\tilde{c}_B(x) = \text{sgn}\left(\sum_{b=1}^B \tilde{\alpha}_b c_b(x)\right)$.
 e) Is \tilde{c}_B equivalent to \hat{c}_B ? Justify your answer.
 Answer:

$$\hat{c}_B(x) = \text{sgn}\left(\sum_{b=1}^B \alpha_b c_b(x)\right) \Rightarrow \tilde{c}_B(x) = \text{sgn}\left(\sum_{b=1}^B \tilde{\alpha}_b c_b(x)\right)$$

 = linear classifier.
 yes, it is equivalent because the sum of classifiers given, weights is a linear fn of

g) Derive a bound on W_B , the sum of the weights after the final iteration. Given that at each iteration b , the sum of the weights increases by at most a factor of $1 + 2\gamma$:

$$\frac{W_b}{W_{b-1}} \leq 1 + 2\gamma,$$

show that for the final sum of weights it holds that

$$W_B \leq N(1 + 2\gamma)^B.$$

Answer:

$$W_b \leq (1 + 2\gamma)^{b-1} W_{b-1}$$

4 pts.

h) Prove the following lower bound on the final weight of a training sample, if it is misclassified by the ensemble:

$$\text{If } \hat{c}_B(x_n) \neq y_n, \text{ then } w_{B,n} \geq \beta^{B/2}.$$

Answer:

4 pts.

i) Finally, we will show a bound on the empirical error rate on the training data. Let $M = |\{n : \hat{c}_B(x_n) \neq y_n\}|$ be the number of training samples misclassified by the ensemble classifier.

1. Plugging the two previous bounds together, we have

$$M\beta^{B/2} \leq \sum_{n=1}^N w_{n,B} = W_B \leq N(1 + 2\gamma)^B.$$

Starting from this inequality, show that the empirical error rate on the training data is bounded by the following:

$$\frac{M}{N} \leq e^{-2B\gamma^2}.$$

(Hint: note that $(1+a)^c \leq e^{ac}$ for $c \geq 0$, and recall that $\beta = \frac{1+2\gamma}{1-2\gamma}$)

Answer:

2. Explain what happens to this bound when we
- (a) increase B (the number of training iterations),
 - (b) increase γ (have a very good learner).

Answer:

7 pts.

This page has been intentionally left blank.

①

mean

$$\frac{x}{n} + p$$