

Final Exam

January 29th, 2011

First and Last name: _____

ETH number: _____

Signature: _____

Topic	Max. Points	Points Achieved	Visum
1	Bayesian Inference	25	
2	Regression	25	
3	Bagging and Boosting	25	
4	SVM	25	
5	Mixture Models	25	
Total		125	

Grade:

This page has been intentionally left blank.

Question 1: Bayesian Inference (25 pts.)

Consider the task of estimating mean μ and variance σ^2 of a Gaussian density from n i.i.d. observations $\mathcal{X} = \{x_1, \dots, x_n\}$, $x \in \mathbb{R}$.

a) Write the maximum likelihood estimator for the mean $\hat{\mu}_{ML}(\mathcal{X})$ as an explicit function of the i.i.d. observations $\mathcal{X} = \{x_1, \dots, x_n\}$:
(please write the direct closed-form solution)

$$1. \hat{\mu}_{ML}(\mathcal{X}) = \arg \max_{\mu} p(\mathcal{X}|\mu) = \frac{1}{n} \sum_{i=1}^n x_i$$

1 pts.

2. Is this a biased estimator? YES \ NO

1 pts.

b) Write the maximum likelihood estimator for the variance $\hat{\sigma}_{ML}^2(\mathcal{X})$ as an explicit function of the i.i.d. observations $\mathcal{X} = \{x_1, \dots, x_n\}$:
(please write the direct closed-form solution)

$$1. \hat{\sigma}_{ML}^2(\mathcal{X}) = \arg \max_{\sigma^2} p(\mathcal{X}|\sigma^2) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

1 pts.

unbiased estimate of population variance: $\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2$

2. Is this a biased estimator? YES \ NO

1 pts.

Now, assume that the variance σ^2 is known.

c) Given a Gaussian prior over the mean (prior with zero mean and variance one), write the posterior density $p(\mu|\mathcal{X})$ as an explicit function of the single observation $\mathcal{X} = \{x_1\}$:
(please write the direct closed-form solution)

$$1. p(\mu|\mathcal{X}) = \frac{p(\mathcal{X}|\mu)p(\mu)}{p(\mathcal{X})} = \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\mu^2}{2\sigma^2}\right) \exp\left(-\frac{(x_1 - \mu)^2}{2\sigma^2}\right) = \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\mu^2 + (x_1 - \mu)^2}{2\sigma^2}\right)$$

2 pts.

2. Is the posterior a Gaussian density? YES \ NO

1 pts.

3. Is the posterior a Gaussian density for all priors? YES \ NO

1 pts.

you have to consider all existing priors, not only Gaussians

d) Given a Gaussian prior over the mean (prior with zero mean and variance one), write the maximum a posteriori estimator $\hat{\mu}_{MAP}(\mathcal{X})$ as an explicit function of the i.i.d. observations $\mathcal{X} = \{x_1, \dots, x_n\}$ (recall that

$$\int e^{-(a^2)x^2} dx = \sqrt{\pi}/a$$

(please write the direct closed-form solution)

2 pts.

$$\hat{\mu}_{MAP}(\mathcal{X}) = \arg \max_{\mu} p(\mu | \mathcal{X}) = \frac{n}{n+1} \frac{\mu + \sigma^2}{n+1} = \frac{1}{n+1} \sum_{i=1}^n x_i$$

2. Does the maximum a posteriori estimate equal the mean of the posterior?

YES/NO

1 pts.

Now consider a binary classification task from observations $\mathcal{X} = \{x_1, \dots, x_n\}$, with $x \in \mathbb{R}^D$. Assume that the likelihood of both classes is Gaussian (assume class prior π_i , mean μ_i , and covariance matrix Σ_i for class y_i , with $i = 1, 2$).

e) Write the discriminant $g_{y_1}(x) = p(y_1 | x)$ as an explicit function of class prior, mean and covariance:

(please write the direct closed-form solution)

3 pts.

$$g_{y_1}(x) = p(x | y_1, \mu_1, \Sigma_1) \cdot p(y_1 | \mu_1, \Sigma_1)$$

$$= \frac{1}{(2\pi)^{D/2} |\Sigma_1|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)\right) \cdot \pi_1$$

f) Assume that $\Sigma_1 = \Sigma_2 = \sigma^2 \mathbb{I}$, where \mathbb{I} denotes the identity matrix. Write the equation satisfied by the separating decision surface. The equation must be an explicit function of x_1 (the single observation), of class prior, means, and covariance:

(please write the direct closed-form solution)

2 pts.

$$1. 0 = g_{y_1}(x) - g_{y_2}(x) = \left[\frac{1}{(2\pi)^{D/2} |\Sigma_1|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)\right) \right] \cdot \pi_1$$

$$-\left[\frac{1}{(2\pi)^{D/2} |\Sigma_2|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2)\right) \right] \cdot \pi_2$$

2. In this case, is the decision surface constant, linear, quadratic, cubic, or something else? constant/linear/quadratic/cubic/other

1 pts.

YES/NO

3. Would it be possible to obtain a cylindrical decision surface?

Be linear.

Under

$$-\frac{1}{2} x^T \Sigma^{-1} x + \dots$$

Not sure
if you
want to
integrate
instead of
differentiate

g) In the two-dimensional case subject to uniform class prior, write means (μ_1, μ_2) and covariances (Σ_1, Σ_2) such that the decision surface is a hyperplane. (any numerical instantiation which satisfies this constraint is acceptable)

$$\mu_1 = \frac{1}{n} \sum_{i=1}^n x_i, \quad \mu_2 = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\Sigma_1 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_1)(x_i - \mu_1)^T, \quad \Sigma_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_2)(x_i - \mu_2)^T$$

3 pts.

h) In the two dimensional case subject to uniform class prior, write means (μ_1, μ_2) and covariances (Σ_1, Σ_2) such that the decision surface is spherical. (any numerical instantiation which satisfies this constraint is acceptable)

$$\mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \mu_2 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Sigma_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \Sigma_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

3 pts.

Question 2: Regression (25 pts.)

Consider the linear regression model expressed in the homogeneous coordinates:

$$y = \beta^{(0)} + \sum_{i=1}^D \mathbf{x}^{(i)} \beta^{(i)} = \mathbf{x}^T \boldsymbol{\beta}$$

where $\mathbf{x} = (1, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(D)}) \in \mathbb{R}^{D+1}$ is the input variable and y is the corresponding target variable.

Assume that the input dataset is given by the matrix $\mathbf{X} \in \mathbb{R}^{N \times (D+1)}$ whose first column is 1. Then the linear regression model for all the observations is written as $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}$. Consider this linear regression model and answer the following questions:

a) For this problem, formally define the Residual Sum of Squares (RSS) cost function and write it down in matrix notation.

Answer:

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

b) Briefly motivate the connection between minimizing the Residual Sum of Squares and maximizing the likelihood of \mathbf{y} given \mathbf{X} .

Answer:

To minimize RSS is to reduce the error between the regressor and data points. Maximizing likelihood $\max_{\boldsymbol{\beta}} P(\mathbf{y}|\boldsymbol{\beta}) \Rightarrow$ finding the $\boldsymbol{\beta}$ that will "fit" the data.

$$\text{error of probability } P(\boldsymbol{\beta}) = \prod_{i=1}^N P(\epsilon_i|\boldsymbol{\beta}) = \prod_{i=1}^N P(y_i - \mathbf{x}_i^T \boldsymbol{\beta}|\boldsymbol{\beta})$$

2 pts.

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

$$\boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I})$$

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

$$P(\mathbf{y}_1, \dots, \mathbf{y}_N)$$

c) We minimize the RSS function and infer the model parameters as $\hat{\beta} = (X^T X)^{-1} X^T y$. Considering the matrix inverse operation, mathematically

describe why in practice we are interested in regularized models such as

ridge regression models rather than the given unregularized model.

Answer:

$$RSS(\beta) = (y - P^T X)(y - P^T X)^T$$

$$\|\beta\|_2^2 = \left(\sum_{i=1}^d \beta_i^2 \right)^2$$

In regularization, you assume β to be distributed by nature according to gaussian law.

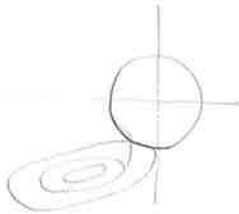
to choose β values that control complexity of model.
 p is distributed by nature, $p(\beta) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{\beta_i^2}{2\sigma_i^2}\right) \Rightarrow$ Take log of $p(\beta)$, $L_p = \left(\sum_{i=1}^d |x_i| p_i \right)^p$

$$\text{then arrive at } -\log p(\beta) = \sum_{i=1}^d \frac{\beta_i^2}{2\sigma_i^2} + \text{constant} \Rightarrow \arg \min_{\beta} \left\{ RSS(\beta) + \lambda \sum_{i=1}^d \beta_i^2 \right\}$$

3 pts.

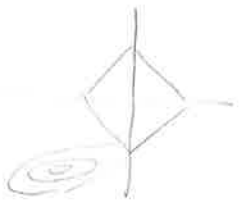
d) Formally define ridge and LASSO regression models. By depicting appropriate plots, demonstrate the difference between these two models in inferring the model parameters β .

Answer:



$$\text{ridge regression } \|\beta\|_2^2 = \left(\sum_{i=1}^d \beta_i^2 \right)^2$$

shrinks the regression coefficients by imposing a penalty on their size.



$$\text{LASSO} \Rightarrow L_1 \text{ penalty } \|\beta\|_1 = \left(\sum_{i=1}^d |\beta_i| \right)$$

favors β values to be exactly 0 when λ is large.
 sparseness penalty model with few coefficients non-vanishing.
 Because the least-square surface often hits the corners of the constraint surface.

$$\beta_1 \beta_2 + \beta_1 \beta_2$$

5 pts.

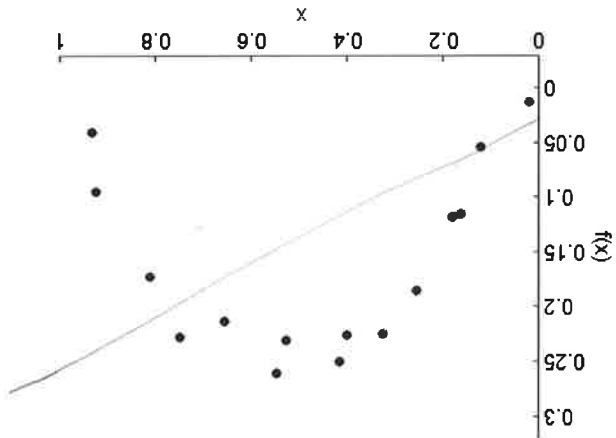
Good for stabilizing matrix positive-definite. Full rank. we can directly manipulate size of β . $\beta = (X^T X + \lambda I)^{-1} X^T y$

Now, we consider a general form of the problem where the set of observations $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, $\mathbf{x}_i \in \mathbb{R}^D$, $y_i \in \mathbb{R}$ are drawn i.i.d. from the joint distribution $P(\mathbf{X}, Y)$. The goal is to find the regression function $f \in \mathcal{F}$ such that the mean squared error $\mathbb{E}_{XY}[(Y - f(\mathbf{X}))^2]$ is minimal. Where the hypothesis class \mathcal{F} contains the set of all polynomial functions.

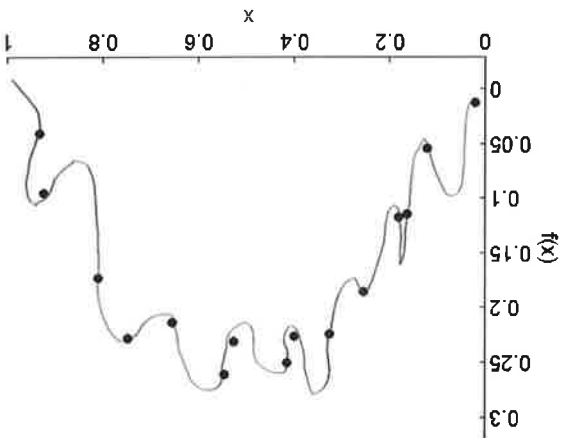
e) Choosing an inappropriate function f might lead to either underfitting or overfitting. For the depicted given data, draw relevant plots to show each of these situations. For the overfitting case, briefly explain what happens if we have more observations.

Answer:

underfitting



overfitting



4 pts.

3 pts.

$$E_{x,y} [f(x)^2] =$$

$$E_{x,y} [(Y - f(x))^2] = E_{x,y} [f(x)^2 - 2Yf(x) + Y^2]$$

Answer:

of variance + squared bias.

g) Expand the mean squared error $E_{XY}[(Y - f(X))^2]$ and write it in the form

2 pts.

$$\begin{aligned} \text{Bias} &= E[f(x)] - \hat{f}(x) \\ \text{variance} &= E[\hat{f}(x) - E[\hat{f}(x)]]^2 \\ \text{MSE} &= E[\hat{f}(x) - f(x)]^2 + E[f(x) - E[\hat{f}(x)]]^2 \end{aligned}$$

Answer: $\text{MSE} = \text{Bias}^2 + \text{variance}$.

f) Write down the mathematical definition of bias and variance.

We can split the mean squared error $E_{XY}[(Y - f(X))^2]$ into bias and variance, and find the best tradeoff between them.

1 pts.

i) According to the results, briefly explain why unbiased estimators remain unbiased after averaging.

Answer: Because taking the average of something gives you an unbiased. And in this case, if the estimator is already unbiased and taking the average of the different unbiased estimators again is unbiased.

3 pts.

h) Calculate the bias of the average estimator $\hat{f}(x)$ in terms of the bias of the estimators $\hat{f}_1(x), \dots, \hat{f}_B(x)$.

Hint: Start with the mathematical definition of the bias for $\hat{f}(x)$.

Answer:

$$\text{bias}(\hat{f}(x)) = \mathbb{E}[\hat{f}(x)] - \mathbb{E}[y|x] = \mathbb{E}\left[\frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)\right] - \mathbb{E}[y|x]$$

$$= \frac{1}{B} \sum_{b=1}^B \mathbb{E}[\hat{f}_b(x)] - \mathbb{E}[y|x] = \left[\frac{1}{B} \sum_{b=1}^B \text{bias}(\hat{f}_b(x)) \right]$$

In this section we study how the averaging changes the bias of a set of unbiased estimators. Assume that we are given a set of B estimators $\hat{f}_1(x), \hat{f}_2(x), \dots, \hat{f}_B(x)$. We take the average estimator by $\hat{f}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)$.

h) Calculate the bias of the average estimator $\hat{f}(x)$ in terms of the bias of the

estimators $\hat{f}_1(x), \dots, \hat{f}_B(x)$.

Hint: Start with the mathematical definition of the bias for $\hat{f}(x)$.

Answer:

$$\text{bias}(\hat{f}(x)) = \mathbb{E}[\hat{f}(x)] - \mathbb{E}[y|x] = \mathbb{E}\left[\frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)\right] - \mathbb{E}[y|x]$$

$$= \frac{1}{B} \sum_{b=1}^B \mathbb{E}[\hat{f}_b(x)] - \mathbb{E}[y|x] = \left[\frac{1}{B} \sum_{b=1}^B \text{bias}(\hat{f}_b(x)) \right]$$

Question 3: Bagging and Boosting (25 pts.)

Bagging and boosting are two possible approaches to combine multiple models (classifiers or regressors) to achieve a composite model with improved performance. Both methods can be used in both a classification and a regression setting.

a) State two essential differences between bagging and boosting.

Answer:

1) Bagging varies the training sets using resampling

Boosting

Trains weak learners with the same training set on every iteration

2) gives the same importance to every prediction

weights the prediction of every weak classifier according to its accuracy.

2 pts.

b) Explain in terms of the bias-variance trade-off why the idea of combining models works.

Answer:

without combining models, either you are trying to battle with having a high bias and low variance or a low bias and high variance.

when you combine models, you are gaining a diverse group of classifiers that partition the sample space in a complex manner that tends to reduce both the bias and variance.

Bagging →
reduce variance
Boosting →
reduce bias.
Give more weight to
misclassified.

2 pts.

c) We now look at the problem of regression and how the combination of individual regression models can give better results. Left-column figures show the individual regression models, right-column figures show the true target function (solid) and the output of averaging the individual models to obtain a composite model (dashed).

1. The individual regression models (in Figure 1 and Figure 3) have been regularized using a regularization parameter λ , i.e. the cost function had the following form

$$RSS_{Ridge}(\beta) = \sum_{i=1}^n (t_i - \phi(\mathbf{x}_i)^\top \beta)^2 + \lambda \|\beta\|_2^2$$

The parameters used were $\lambda = 0.09$ and $\lambda = 13.5$.

Associate the regularization parameters to the figures in the left column.

Answer:

Figure 1: $\lambda = \dots 13.5 \rightarrow$ zero leading to large bias.
 large λ pulls the weight parameters toward

Figure 3: $\lambda = \dots 0.09 \rightarrow$ allows model become finely tuned to noise \rightarrow leading to large variance

2. Please interpret the figures in terms of the bias-variance trade-off

Answer:

Figure 1 and 2: Bias is high, variance is low
 rigid model.

Figure 3 and 4: Bias is low, variance is high

Figure 4 gives a low bias and variance.
 flexible models

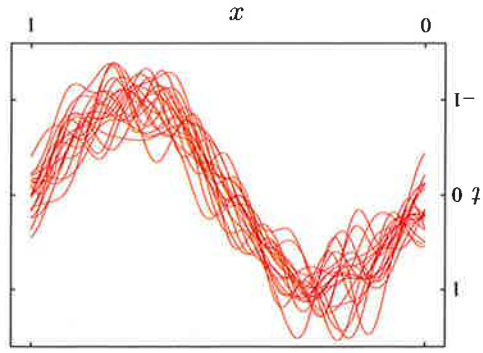


Figure 3

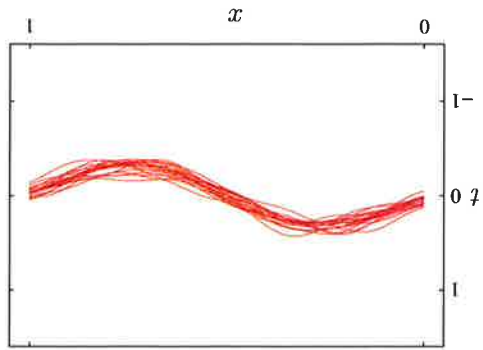


Figure 1

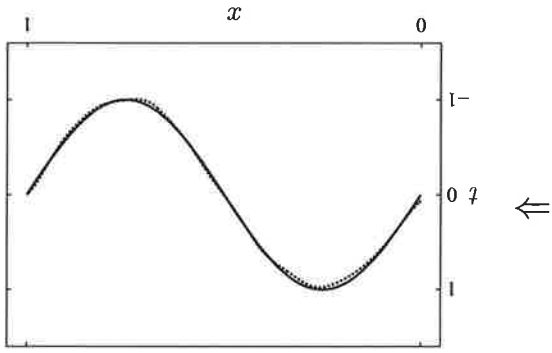


Figure 4

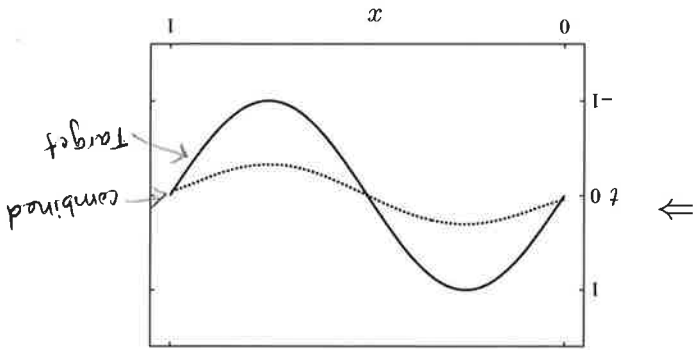


Figure 2

4 pts.

We now consider *bagging* in a regression setting. In order to model a *target function* $h(x)$, we combine M *individual models* $y_m(x)$, $m = 1..M$ to obtain a *committee model* $y^{COM}(x)$. We model the *error* of each individual model using $\epsilon_m(x)$, $m = 1..M$. This question's final goal is to show that the error of the committee model is M times smaller than the average error of the individual models.

1. Write down the output of the committee model, which averages the individual model's outputs. Answer:

$$y_{COM}(x) = \dots = \frac{1}{M} \sum_{m=1}^M y_m(x)$$

2. Write down the error $\epsilon_m^m(x)$ of an individual model in terms of the target function $h(x)$ and the output of the individual model $y_m(x)$. Answer:

$$\epsilon_m^m(x) = h(x) - y_m(x)$$

3. Write down the expected squared error of an individual model. \mathbb{E}_x denotes the expectation value w.r.t. the distribution of x . Answer:

4. Write down the average of the expected squared errors made by the individual models

individual models

Answer: $E_{AV} = \dots \frac{1}{M} \sum_{m=1}^M E^x \left[\varepsilon_m(x)^2 \right]$

5. Write down the expected squared error made by the committee model, in terms of the output of the committee model and the target function.

Answer:

$$E_{COM} = \mathbb{E}_x \left[\left(h(x) - y_{COM}(x) \right)^2 \right]$$

$$= \mathbb{E}_x \left[\left(\frac{1}{M} \sum_{m=1}^M y_m(x) - h(x) \right)^2 \right]$$

Answer: $E_{COM} = \mathbb{E}_x \left[\left(h(x) - y_{COM}(x) \right)^2 \right]$

6. Show that $E_{COM} = \frac{1}{M} E_{AV}$. Hint: Use the assumption that the errors of the individual models are uncorrelated, i.e. $\mathbb{E}_x[\epsilon_m(x)\epsilon_l(x)] = 0, m \neq l$, and that the mean error of the individual models is zero, i.e. $\mathbb{E}_x[\epsilon_m(x)] = 0$.

of the individual models is zero, i.e. $\mathbb{E}[\epsilon^u(x)] = 0$.

Answer:

$$E_{\text{cor}} = E_x - \left[\frac{1}{N} \sum_{m=1}^M y_m(x) - h(x) \right]^2$$

$$\left(\left[y(x) - y_m(x) \right]^T E_x \sum_{m=1}^M \frac{W_m}{T} \right) = w_{om} \in$$

$$= \left(\frac{1}{T} \sum_{m=1}^M E_x [e_m(x)] \right)$$

$$E_{\text{com}} = \frac{1}{M} \sum_{i=1}^M E_{AV}^{(i)}$$

$$[y]_{\text{corr}} = \frac{1}{n} \sum_{i=1}^n E_i [e_i(x_i)^2]$$

10 pts.

$$E\left[\left(\frac{1}{n} \sum_{i=1}^n y_i^0(x) - h(x)\right) \left(\frac{1}{n} \sum_{i=1}^n y_i^1(x) - h(x)\right)\right] = y_1^2$$

We now turn our attention to boosting in a classification setting.

e) Given i.i.d. training data $\{(\mathbf{x}_1, c_1), (\mathbf{x}_2, c_2), \dots, (\mathbf{x}_n, c_n)\}$, $\mathbf{x}_i \in \mathbb{R}^D$, $c_i \in \{-1, 1\}$, and individual classifiers $y_m(\mathbf{x})$, $m = 1 \dots M$, we can define the following additive model

$$f_k(\mathbf{x}) = \sum_{j=1}^k \alpha_j y_j(\mathbf{x})$$

where α_j are weights, $k \leq M$.

Think of $f_k(\mathbf{x})$ as the *additive model* which uses the first $1..k$ individual classifiers.

1. Define the sum-of-squares error of $f_k(\mathbf{x})$ on the training dataset.

Answer:

$$Error_k = \sum_{i=1}^n \exp(-c_i f_k(\mathbf{x}_i))$$

error $f_k(\mathbf{x})$

2. Assume that the first $1..(k-1)$ classifiers $y_j(\mathbf{x})$ and weights α_j , $j = 1..(k-1)$ have already been determined, i.e. fixed.

Show that finding the optimal k -th classifier $y_k(\mathbf{x})$ and its weight α_k involves fitting the k -th classifier to the residual errors $f_{k-1}(\mathbf{x}_i) - c_i$ made by the $(k-1)$ -th additive model.

(Hint: Start from the expression $Error_k$ and recognize the terms $f_{k-1}(\mathbf{x}_i) - c_i$)

$$\sum_{i=1}^n \exp(-c_i f_k(\mathbf{x}_i))$$

$$\alpha_k = (f_{k-1}(\mathbf{x}_i) - c_i)$$

$$y_k(\mathbf{x})$$

$$\sum_{i=1}^n \exp(-c_i f_{k-1}(\mathbf{x}_i) - \alpha_k y_k(\mathbf{x}_i))$$

$$\sum_{i=1}^n \exp(-c_i f_{k-1}(\mathbf{x}_i) - \alpha_k y_k(\mathbf{x}_i))$$

Question 4: SVM (25 pts.)

In the following questions choose one answer only.



a) Why does a Support Vector Machine generalize better than a Perceptron?

1. The selected support vectors tend to be very typical samples.
2. By supporting vectors it is not restricted to scalar input.
3. The requirement of maximal margins reduces the arbitrariness.
4. The support vector machine will have access to more training data.

2 pts.

b) What is the advantage of using kernels in support vector machines?

1. They tend to maximize the margins.
2. They enable non-linear separations.
3. They will increase the number of support vectors.
4. They reduce the risk of getting stuck in local minima.

2 pts.

Explain your answer to the following questions in 1-2 sentences.

c) Suppose that you have a binary SVM classifier with a linear kernel. Consider a vector x_i for which $y_i(\langle w, x_i \rangle + w_0) > 1$ ($x_i \in \mathbb{R}^D, y_i \in \{-1, +1\}$)

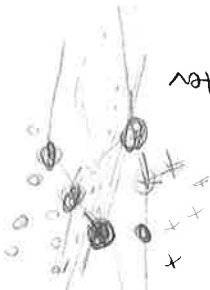
1. Is x_i correctly classified? YES/NO

2. If we remove x_i from the training set and re-train the classifier, will the decision boundary change or stay the same?

Answer:

The decision boundary might change depending if the training sample was a support vector

in the previous case.



3 pts.

3 pts.

$$\langle \phi(x), \phi(\bar{x}) \rangle = \phi(x_1, x_2) \cdot \phi(\bar{x}_1, \bar{x}_2)$$

Answer:

$$K(x, \bar{x}) = \langle \phi(x), \phi(\bar{x}) \rangle$$

$$x = (x_1, x_2) \quad \bar{x} = (\bar{x}_1, \bar{x}_2)$$

such that

Hint: One possible way to prove it is to find a feature mapping $\phi(x)$,Prove that $K(x, \bar{x})$ is a legitimate kernel function.

$$K(x, \bar{x}) = \begin{cases} 1 & \text{if } x = \bar{x} \\ 0 & \text{else} \end{cases}$$

e) Let \mathcal{X} be a finite set and consider the following function. For $x, \bar{x} \in \mathcal{X}$:

3 pts.

Answer:

What kind of separation is imposed by this kernel? (which points in the space will become close/far using this kernel).

$$K(h, \tilde{h}) = \sum_m \min(h_j, \tilde{h}_j) \quad h, \tilde{h} \in H$$

Consider using the following kernel function:

$$1. h = \{h_1, \dots, h_m\}, \quad \forall j: h_j \geq 0$$

$$2. \sum_{j=1}^L h_j = L \quad \text{for some } L \in \mathbb{N}.$$

class of histograms

d) We would like to build a binary SVM classifier for histograms. Let H denote the class of histograms, for each $h \in H$ the following properties hold:

L2-SVM use the square sum of the slack variables ξ_i in the objective function instead of the linear sum of the slack variables (squaring the hinge loss). Let $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ be a training set of examples and binary labels $y_i \in \{-1, +1\}$. The primal formulation of the L2-SVM is as follows

$$\min_{w, w_0, \xi} \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2$$

$$w = \sum_{i=1}^n \alpha_i$$

s.t.

$$y_i(\langle w, x_i \rangle + w_0) \geq 1 - \xi_i \quad i = 1, \dots, n$$

$$\xi_i \geq 0 \quad i = 1, \dots, n$$

$$y_i(\langle w, x_i \rangle + w_0) - 1 + \xi_i$$

f) Show that removing the last set of constraints: $\forall i : \xi_i \geq 0$ does not change the optimal solution to the primal problem.

Answer:

$$L(w, w_0, \xi, \lambda) = \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 + \sum_{i=1}^n \lambda_i (y_i(\langle w, x_i \rangle + w_0) - 1 + \xi_i) = 0$$

$$\frac{d}{d\xi} = 0$$

$$\frac{d}{d\xi} = 0 \Rightarrow \lambda_i = 0$$

$$\frac{d}{d\lambda} = 0$$

$$= \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 - \sum_{i=1}^n \lambda_i (y_i(\langle w, x_i \rangle + w_0) - 1 + \xi_i) = 0$$

4 pts.

$$w/out \cdot \frac{1}{2} \geq 0 \text{ constraint}$$

Now we would like to derive the dual form of the L_2 -SVM.

g) Write down the Lagrangian of the L_2 -SVM.

Answer:

$$L() =$$

2 pts.
h) Compute the derivatives of the Lagrangian with respect to the appropriate variables.
Answer:

3 pts.

We decided not to bother you with the remaining computation, instead we finished the derivation ourselves.

i) Below are 4 optimization problems, only 1 of which is the dual L2-SVM. Circle the optimization problem corresponding to the dual L2-SVM. Explain your choice by either shortly falsifying the other options or by showing the full derivation.

Hint: The full derivation is more time consuming.

Derive dual form.

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \frac{1}{2C} \sum_{i=1}^n \alpha_i^2 \\ \text{s.t.:} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \forall i: \alpha_i \geq 0 \end{aligned}$$

$$\begin{aligned} \max_{\alpha, \xi} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + C \sum_{i=1}^n \xi_i \\ \text{s.t.:} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \forall i: 0 \leq \alpha_i \leq \frac{C}{\xi_i} \end{aligned}$$

$$\begin{aligned} \min_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \frac{1}{2C} \sum_{i=1}^n \alpha_i^2 \\ \text{s.t.:} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \forall i: \alpha_i \geq 0 \end{aligned}$$

$$\begin{aligned} \max_{\alpha, \xi} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + C \sum_{i=1}^n \xi_i \\ \text{s.t.:} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \forall i: \alpha_i \geq 0 \end{aligned}$$

doesn't include constraint for slack

Answer:
In the primal form you solved for w 's and plugged it into dual form $w = \sum \alpha_i y_i x_i$.

3 pts.

This page has been intentionally left blank.

Question 5: Mixture of Bernoulli models (25 pts.)

Let $\mathbf{x} = (x_1, \dots, x_D)^T \in \{0, 1\}^D$ be a D -dimensional random binary vector. We assume that every \tilde{x}_i is Bernoulli distributed with parameter μ_i , i.e.

$$p(x_i = 1; \mu_i) = \mu_i$$

$$p(x_i = 0; \mu_i) = 1 - \mu_i,$$

which can also be written as

$$p(x_i; \mu_i) = \mu_i^{x_i} (1 - \mu_i)^{1-x_i}.$$

Under the assumption of the x_i 's being independent, the distribution of the random vector $\mathbf{x} = (x_1, \dots, x_D)^T$ is

$$p(\mathbf{x}; \boldsymbol{\mu}) = \prod_{i=1}^D \mu_i^{x_i} (1 - \mu_i)^{1-x_i},$$

$$\text{Cov}(\mathbf{x}) = E[\mathbf{x}\mathbf{x}^T] - E[\mathbf{x}]E[\mathbf{x}]^T$$

with $\boldsymbol{\mu} = (\mu_1, \dots, \mu_D)^T$.

a) What is $E(\mathbf{x})$ and $\text{Cov}(\mathbf{x})$ in this setting?

Answer:

$$E(\mathbf{x}) = E\left[\prod_{i=1}^D \mu_i^{x_i} (1 - \mu_i)^{1-x_i}\right] = \log \mu_i^{x_i} (1 - \mu_i)^{1-x_i}$$

$$E[x] = \mu = \log(\mu_i^{x_i}) + \log((1 - \mu_i)^{1-x_i}) = x_i \log(\mu_i) + (1 - x_i) \log(1 - \mu_i)$$

$$= E[x_i \log(\mu_i)] + E[(1 - x_i) \log(1 - \mu_i)] = \log(\mu_i) E[x_i]$$

$$\text{Cov}(\mathbf{x}) = \text{Diag}[\mu_i(1 - \mu_i)] = I$$

b) Consider a Bernoulli mixture model for binary random vectors with K mixture components, i.e.

$$p(\mathbf{x}; \boldsymbol{\mu}, \mathbf{w}) = \sum_{k=1}^K w_k p(\mathbf{x}; \boldsymbol{\mu}_k)$$

prior weights likelihood

with weights $\mathbf{w} = (w_1, \dots, w_K)^T$ and Bernoulli parameters $(\boldsymbol{\mu}_k)_{k=1}^K = (\mu_{k1}, \dots, \mu_{kD})$.

(Side remark, which is not relevant to solve the question: The number of mixture components K cannot be larger than the dimension D of the random vector to ensure identifiability.)

What are natural constraints on w to obtain a proper probability distribution $p(\mathbf{x}; (\boldsymbol{\mu}_k)_k, \mathbf{w})$?

Answer:

c) Show that

$$\mathbb{E}(\mathbf{x}) = \sum_k w_k \bar{\boldsymbol{\mu}}_k \quad \text{and}$$

$$(\text{Cov}(\mathbf{x}))_{i,j} = \begin{cases} \sum_k w_k \mu_{ki}^2 - \mathbb{E}(\mathbf{x}_i)^2 & \text{if } i = j \\ \sum_k w_k \mu_{ki} \mu_{kj} - \mathbb{E}(\mathbf{x}_i) \mathbb{E}(\mathbf{x}_j) & \text{if } i \neq j \end{cases}$$

(the elements of the covariance matrix)

where μ_{ki} is the i -th component of $\boldsymbol{\mu}_k$.

Hints:

1. You do not need to include the last terms, $-\mathbb{E}(\mathbf{x}_i)^2$ and $-\mathbb{E}(\mathbf{x}_i) \mathbb{E}(\mathbf{x}_j)$, in your derivation, since these follow from the relation

$$\text{Cov}(\mathbf{x}) = \mathbb{E}(\mathbf{x} \mathbf{x}^T) - \mathbb{E}(\mathbf{x}) \mathbb{E}(\mathbf{x})^T.$$

Thus, you need to focus only on $\mathbb{E}(\mathbf{x} \mathbf{x}^T)$.

2. Consider $\mathbb{E}(\mathbf{x}_i \mathbf{x}_j)$ separately for the cases $i = j$ and $i \neq j$. Use the fact that $\mathbb{E}(\mathbf{x}_i \mathbf{x}_j) = P(\mathbf{x}_i = 1 \wedge \mathbf{x}_j = 1)$, since \mathbf{x} is a binary vector.

What is the main (qualitative) difference between the covariance derived in a) and the one for the mixture model?

Answer:

5 pts.

d) In order to derive an EM-method to determine the mixture of Bernoulli parameters, we introduce a latent binary vector $\mathbf{z} = (z_1, \dots, z_K)^T \in \{0, 1\}^K$ linked with \mathbf{x} , such that $\sum_k z_k = 1$ and $z_k = 1$ if \mathbf{x} is generated by the k -th mixture component. Hence, we have

$$p(\mathbf{x} | \mathbf{z}) = \prod_{k=1}^K p(\mathbf{x}; \boldsymbol{\mu}_k)^{z_k} \quad \text{and} \quad p(\mathbf{z}) = \prod_{k=1}^K w_k^{z_k}.$$

Derive $p(\mathbf{x})$ by marginalizing over \mathbf{z} . Answer:

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{x} | \mathbf{z}) \cdot p(\mathbf{z})}{p(\mathbf{x})} \quad p(\mathbf{x}) = \int p(\mathbf{x} | \mathbf{z}) \cdot p(\mathbf{z}) d\mathbf{z}$$

e) M-step: let $\mathbf{X} := (\mathbf{x}_n)_n \in \mathbb{R}^{D \times N}$ be a sequence of N i.i.d. samples drawn from the mixture of Bernoulli distribution. We have a corresponding matrix of the latent variables $\mathbf{Z} := (\mathbf{z}_n)_n \in \mathbb{R}^{K \times N}$. Write down the joint log-likelihood function $\ln p(\mathbf{X}, \mathbf{Z}; (\boldsymbol{\mu}_k)_k, \mathbf{w})$ and derive the M-step by maximizing the joint log-likelihood over the unknown parameters $(\boldsymbol{\mu}_k)_k$ and \mathbf{w} . Hint: don't forget to introduce a Lagrange multiplier for the constraint on \mathbf{w} .

Answer:

6 pts.

f) E-step: derive the expectation of the joint log-likelihood function with respect to \mathbf{Z} . Hints:

1. Derive $p(\mathbf{Z} | \mathbf{X})$ first via $p(\mathbf{z})$ and $p(\mathbf{x} | \mathbf{z})$ given above.
2. Consider $\mathbb{E}(z_{nk}) = p(z_{nk} = 1)$, since $z_{nk} \in \{0, 1\}$.
3. Observe that z_{nk} appears only linearly in the joint log-likelihood function.

Answer:

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{x} | \mathbf{z}) \cdot p(\mathbf{z})}{p(\mathbf{x})} = \frac{\prod_k p(x_k | \mu_k)^{z_k} \cdot \prod_k w_k^k}{\prod_k p(x_k)}$$

log-likelihood fn:

$$\log p(\mathbf{x} | \mathbf{z}) = \sum_k \sum_i z_{ik} \log p(x_i | \mu_k) + \sum_k \log p(z_k)$$

samples mixtures

$$\sum_k (x_{nk} \ln \mu_k)$$

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

$$E[S^2]$$

Primal form.

$$L(w, w_0, \frac{1}{2}) = \frac{1}{2} w^2 + \frac{1}{2} w_0^2 - \sum_{i=1}^N \alpha_i (y_i w_i x_i + w_0 - 1 + \frac{1}{2}) - \sum_{i=1}^N \lambda_i \frac{1}{2} = 0$$

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^N \alpha_i y_i x_i = 0 \quad w = \sum_{i=1}^N \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial w_0} = w_0 - \sum_{i=1}^N \alpha_i = 0$$

$$\frac{\partial L}{\partial \frac{1}{2}} = c \sum_{i=1}^N \frac{1}{2} - \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \alpha_i = 0 \rightarrow c \sum_{i=1}^N \frac{1}{2} = \sum_{i=1}^N \lambda_i + \sum_{i=1}^N \alpha_i = 0 \Rightarrow \sum_{i=1}^N \frac{1}{2} = \frac{1}{c} \left(\sum_{i=1}^N \lambda_i + \sum_{i=1}^N \alpha_i \right)$$

Primal form:

$$\max_{\alpha} \frac{1}{2} \left(\sum_{i=1}^N \alpha_i y_i x_i \right)^2 + \frac{1}{2} \left(\sum_{i=1}^N \alpha_i y_i \right)^2 - \sum_{i=1}^N \alpha_i (y_i x_i + w_0 - 1 + \frac{1}{2}) + \sum_{i=1}^N \lambda_i \frac{1}{2}$$

$$-\sum_{i=1}^N \alpha_i \frac{1}{2} - \sum_{i=1}^N \lambda_i \frac{1}{2} = 0$$

$$\Rightarrow -\frac{1}{2} \sum_{i=1}^N \alpha_i y_i x_i + \frac{1}{2} \sum_{i=1}^N \alpha_i y_i^2 + \frac{1}{2} \sum_{i=1}^N \alpha_i^2 - \sum_{i=1}^N \alpha_i w_0 + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \lambda_i \frac{1}{2} = 0$$