

Final Exam
February 8th, 2012

First and Last name: _____

ETH number: _____

Signature: _____

Topic	Max. Points	Points	Signature
1	Bayesian Inference and MLE	20	
2	Linear Classifiers and Kernels	20	
3	Ensemble Methods	20	
4	Regression, Bias and Variance	20	
5	Unsupervised Learning	20	
Total		100	

Grade:

Question 1: Bayesian Inference and Maximum Likelihood (20 pts.)

A telecommunication company needs to estimate the rate of the telephone calls in a small town in order to adjust its channel capacity. We model a length of a time-interval between telephone calls as a random variable x . Knowing that a sequence of calls is the realization of the Poisson process, we model the time it takes before the next call using the exponential distribution.

Consider the task of estimating the rate parameter λ of the exponential distribution from n i.i.d. observations $\mathcal{X} = \{x_i\}_{i=1}^n$, $x \in \mathbb{R}$

$$\text{Exp}(x|\lambda) = \lambda \exp(-\lambda x).$$

a) Write the maximum likelihood estimator for the rate $\hat{\lambda}_{ML}(\mathcal{X})$ as an explicit function of the i.i.d. observations $\mathcal{X} = \{x_i\}_{i=1}^n$:

(please write the direct closed-form solution and the derivation)

1. $\hat{\lambda}_{ML}(\mathcal{X}) = \arg \max_{\lambda} p(\mathcal{X}|\lambda) = \prod_{i=1}^n \lambda \exp(-\lambda x_i)$ **2 pts.**

$$\begin{aligned} \Rightarrow \log \left(\prod_{i=1}^n \lambda \exp(-\lambda x_i) \right) &= n \log \lambda - \lambda \sum_{i=1}^n x_i \\ \frac{d}{d\lambda} (n \log \lambda - \lambda \sum_{i=1}^n x_i) &= 0 \quad n \frac{1}{\lambda} - \sum_{i=1}^n x_i = 0 \quad \frac{\lambda}{n} = \sum_{i=1}^n x_i \end{aligned}$$

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

2. Where did you use the fact that the observations are i.i.d?

1 pt.

$\frac{1}{n}$ in the likelihood equation

To find the posterior density $p(\lambda|\mathcal{X})$ we need a prior on λ . We claim that a conjugate prior for the exponential distribution is the gamma distribution

old distribution $p(\lambda) = p(\lambda|x)$ distribution that x is drawn from

$$\text{Gamma}(\lambda|\alpha, \beta) = \frac{\Gamma(\alpha)}{\beta^\alpha} \lambda^{\alpha-1} \exp(-\lambda\beta)$$

where $\Gamma(\alpha) = \int_0^\infty \exp^{-t} t^{\alpha-1} dt$ is the gamma function.

b) 1. What does conjugate prior mean? 1 pt.

the prior distribution $p(\lambda)$ is conjugate to the likelihood distribution $p(x|\lambda)$ if multiplying these two distributions together and normalizing results in another distribution of the same form as the prior $p(\lambda)$.

2. Show that the gamma distribution is the conjugate prior of the exponential distribution. 2 pts.

$$= (\lambda \exp(-\lambda x)) \cdot \left(\frac{\Gamma(\alpha)}{\beta^\alpha} \lambda^{\alpha-1} \exp(-\lambda\beta) \right) = \lambda^\alpha \exp(-\lambda(x+\beta))$$

$$= \lambda^{1+\alpha-1} \exp(-\lambda(x+\beta))$$

c) * Given a Gamma prior over the rate λ (prior with parameters α and β), write the maximum a posteriori estimator $\hat{\lambda}_{\text{MAP}}(\mathcal{X})$ as an explicit function of the i.i.d. observations $\mathcal{X} = \{x_i\}_{i=1}^n$ (please write the direct closed-form solution)

$$\hat{\lambda}_{\text{MAP}}(\mathcal{X}) = \arg \max_{\lambda} p(\lambda|\mathcal{X}) = \prod_{i=1}^n \left[\exp(-\lambda x_i) \right] \left(\frac{\Gamma(\alpha)}{\beta^\alpha} \lambda^{\alpha-1} \exp(-\lambda\beta) \right)$$

$$\frac{d}{d\lambda} \left[\exp(-\lambda(x+\beta)) \right] = -\exp(-\lambda(x+\beta)) = 0$$

$$\hat{\lambda}_{\text{MAP}} = \frac{\alpha}{\sum_{i=1}^n x_i + \beta}$$

$$\hat{\lambda}_{\text{MAP}} = \frac{\alpha}{\sum_{i=1}^n x_i + \beta}$$

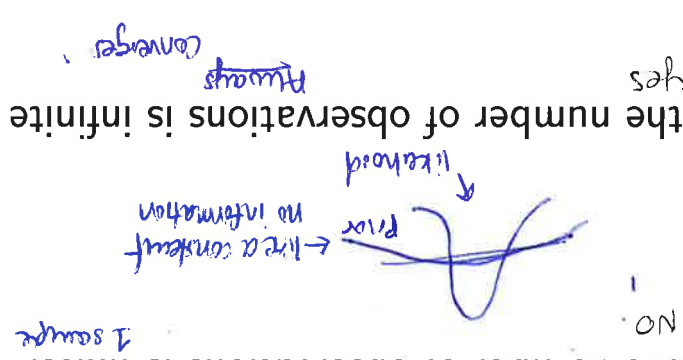
multivariate normal

$$\frac{1}{(2\pi)^{n/2} |\Sigma|} \exp\left(-\frac{1}{2}(\vec{x}-\vec{\mu})^T \Sigma^{-1}(\vec{x}-\vec{\mu})\right)$$

\Downarrow
 not Σ

d) When is the maximum likelihood estimator (MLE) equal to the maximum a posteriori (MAP) estimator given a set of i.i.d. observations $\mathcal{X} = \{x_i\}_{i=1}^n$?

1. If the number of observations is finite. 2 pts.
2. If the number of observations is infinite ($n \rightarrow \infty$). 2 pts.



e) Assume that you have a set of the i.i.d. observations $\mathcal{X} = \{x_i\}_{i=1}^n$ and that you can not decide which distribution to use for data description: the Gaussian distribution or the Beta distribution. If you use the Bayesian framework what you can look at? 2 pts.

* Now consider a binary classification task from a set of the i.i.d. observations $\mathcal{X} = \{\mathbf{x}_i, y_i\}_{i=1}^n$, with $\mathbf{x} \in \mathbb{R}^D$. Assume that the likelihood of both classes is Gaussian (assume class prior π_i , mean μ_i , and co-variance matrix Σ_i for class y_i , with $i = 1, 2$).

f) Recall that a discriminant function for class y_i is defined as:

$$g_{y_i}(\mathbf{x}) = p(y_i | \mathbf{x})$$

class decision: $g_{y_1}(\mathbf{x}) > g_{y_2}(\mathbf{x})$

How can you find a decision surface in terms of likelihood, prior and evidence? 1 pt.

$$p(y_i | \mathbf{x}) = p(\mathbf{x} | y_i) p(y_i)$$

Decision surface $g_1 = g_2$

$$p(y=1 | \mathbf{x}) \pi_1 > p(y=2 | \mathbf{x}) \pi_2$$

Decision discriminant fn = posterior, defines boundary if one is greater $\rightarrow p(\mathbf{x} | y_1) p(y_1) > p(\mathbf{x} | y_2) p(y_2)$

g) Assume that

$$\mu_1 = \mu_2 = \mu$$

$$\Sigma_1 = \frac{1}{2\lambda_1} \mathbb{I}, \quad \Sigma_2 = \frac{1}{2\lambda_2} \mathbb{I}$$

$$\lambda_1 > 0, \quad \lambda_2 > 0, \quad \lambda_1 \neq \lambda_2$$

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 the class prior, means and covariance:

(please write the solution in the polynomial form)

$$g_1(x) - g_2(x) = 0$$

3 pts.

$$\ln(P(x_1 | \omega_1) + \ln(P(x_1 | \omega_2)) - \ln(P(x_1 | \omega_2)) - \ln(P(x_1 | \omega_1))$$

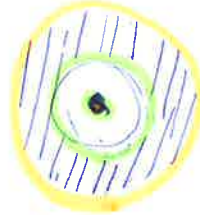
$$- \frac{1}{2} x_1^T \Sigma_1^{-1} x_1 + \frac{1}{2} x_1^T \Sigma_2^{-1} x_1 + w_1$$

$$w_1 = \frac{1}{2} \mu_1^T \Sigma_1^{-1} \mu_1 - \frac{1}{2} \mu_2^T \Sigma_2^{-1} \mu_2 - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln|\Sigma_1|$$

2. In the case described above, is the decision surface linear, parabolic, spherical, cylindrical, or something else?

linear | parabolic | spherical | cylindrical | other

Dependent on covariance matrix



Gaussian within a gaussian

ellipsoidal

2 pts.

Question 2: Linear Classifiers and Kernels (20 pts.)

a) Below is a list of algorithms which given a training set output a prediction function. Cross **all** of the algorithms that necessarily output a linear (in the original space) prediction function.

- ☒ Neural network
- ☒ Perceptron with learning rate $\eta = 1$
- ☐ SVM with radial basis kernel
- ☐ K-nearest neighbor classifier
- ☒ SVM with polynomial kernel with degree 1
- ☒ Ridge regression $w^T \phi(x)$

also works w/ kernels.

3 pts.

b) Recall the SVM problem. As a constrained optimization problem a solution can be obtained through both the primal and the dual form.

1. Given a primal solution for the SVM, write down the resulting classifier.

$$w = \sum \alpha_i y_i x_i$$

$$y = \text{sgn}(w^T x + w_0)$$

2. Given a dual solution for the SVM, write down the resulting classifier.

$$w = \sum \alpha_i y_i x_i \rightarrow w^T x + w_0 = y$$

$$y = \text{sgn}(w^T x + w_0)$$

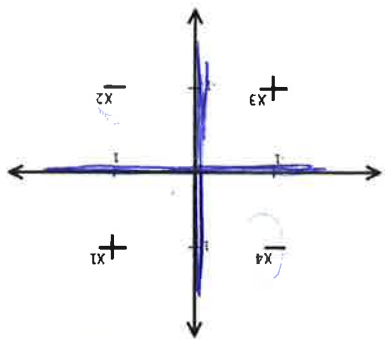
3. In practice, often the dual form of the SVM is solved to obtain a classifier. Provide one advantage of solving the dual SVM instead of the primal.

Even though the primal form it has been transformed into

$$k(x_i, y) = \langle \phi(x_i), \phi(y) \rangle$$

$$w = \sum \alpha_i y_i x_i$$

c) Let $S = \{x_i, y_i\}_{i=1}^4$ be the following training set



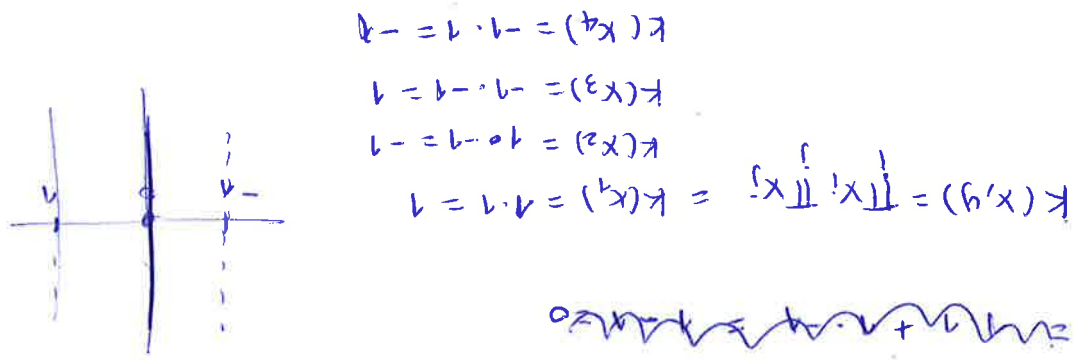
$$\begin{aligned} x_1 &= (1, 1) & y_1 &= 1 \\ x_2 &= (1, -1) & y_2 &= -1 \\ x_3 &= (-1, -1) & y_3 &= 1 \\ x_4 &= (-1, 1) & y_4 &= -1 \end{aligned}$$

$$\min \frac{1}{2} \|w\|^2 + \alpha(\dots)$$

Suppose that we trained an SVM on S and the resulting classifier $f(x)$ achieved zero training error. We ask you to provide an explicit description of $f(x)$ (a formula with numeric values).
Hint: Think of a suitable kernel function or alternatively a feature map.

$$\phi(x) = \langle \phi(x_1), \phi(x_2) \rangle$$

$$K(x_i, x_j) = \prod_{k=1}^4 x_{ik} x_{jk}$$



$$\begin{aligned} K(x_1, y) &= \prod_{k=1}^4 x_{1k} x_{jk} \\ K(x_1, x_1) &= 1 \cdot 1 = 1 \\ K(x_1, x_2) &= 1 \cdot (-1) = -1 \\ K(x_1, x_3) &= (-1) \cdot (-1) = 1 \\ K(x_1, x_4) &= (-1) \cdot 1 = -1 \end{aligned}$$

Closure property:

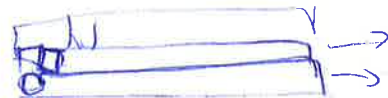
$$K(x, y) = f(x) \cdot f(y)$$

d) Consider a training set $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^D$ and $y_i \in \{-1, 1\}$

1. Briefly describe a leave one out (LOO) procedure for estimating the error of an SVM classifier on S .

2 pts

$$R_{cv} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\text{indicator fcn}}$$



number of errors
10

2. What is the LOO error?

1 pts

$$R_{cv} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{y_i \neq C^{(n)}(\mathbf{x}_i)}$$

3. Suppose that we trained an SVM classifier on the entire dataset S , denote by sv the set of support vectors, $sv = \{\mathbf{x}_j | \alpha_j > 0\}$.

For the same value of C , prove that the LOO error is bounded by $\frac{n}{|sv|}$ i.e.

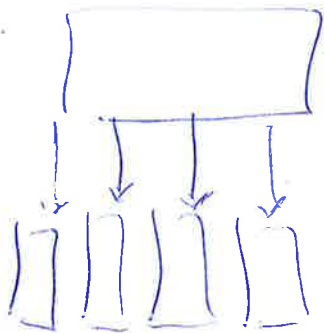
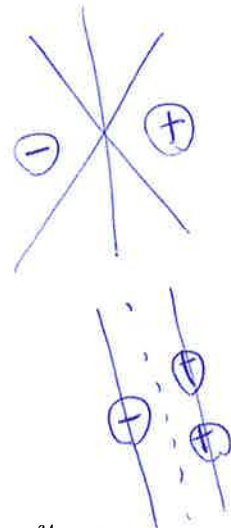
$$\text{LOO error} \leq \frac{n}{|sv|}$$

cardinality = number of sv

4 pts.

$$R_{cv} = \frac{1}{n} \sum_{i=1}^n R_{-i}$$

$$R_{-i} = \frac{1}{n} \sum_{j \in sv} R_{-i} \leq \frac{1}{n} |sv| \cdot C$$



Question 3: Bagging and Boosting (20 pts.)

a) Answer precisely the following questions.

1. Are bagging and Boosting Bayesian approaches? Why?

No Bayesian approaches. → pg. 410 has the

Bagging construct new parameters ← approximate parameters prior

1 pt.

2. How is it possible to detect outliers with AdaBoost?

Super large weights at the end of the training.

to one classifier will detect

1 pt.

3. From the frequentist perspective, bagging is motivated by the tradeoff between two terms. Which?

Bias, variance

1 pt.

4. AdaBoost has an alternative interpretation which is based on the minimization of a certain cost function. Which function?

exponential loss

1 pt.

5. AdaBoost aims at selecting the best approximation to which ratio?

log-odds ratio

2 pt.

6. Why is the standard form of AdaBoost limited to binary classification?

1 pt.

$$C_b = \arg \min \left(\sum_{b=1}^B \alpha_b C_b(x) \right)$$

7. How could one parallelize bagging?

subsets

1 pt.

8. Name a design property of the base classifiers of AdaBoost which impacts the overall predictive power.

weights to classifiers.

1 pt.

weakness of base learners
strong and weak learners

9. Under which conditions does AdaBoost yield good results even when the base classifiers exhibit an individual performance that is only slightly better than that purely due to chance?

when uncorrelated to different classifiers

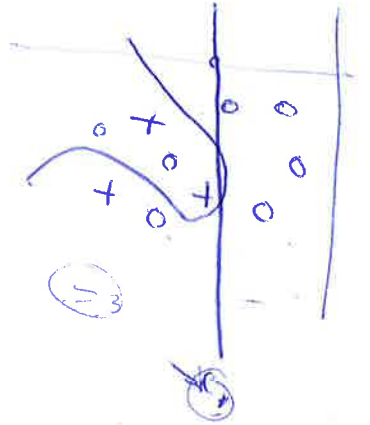
2 pt.

classifiers give different outputs

and answers
B is much bigger than A.

and much larger dataset

2 pt.



uncorrelated
different

