

# EEL709: Re-Major Test

July 8, 2013

Maximum Marks: 62

**Instructions: All working must be clearly shown, with no missing or assumed steps. Whenever words like ‘obtain’, ‘derive’, or ‘compute’ occur, you should make explicit your entire process for doing so. Your answers should be self-sufficient, not requiring reference to any other materials.**

1. (a) It is sometimes said that the Bayesian view of probability is subjective, whereas the frequentist view is objective. Assess this statement; to what extent do you agree with it? Can you give an example of a pair of probabilistic statements to illustrate the difference in question? [4]

(b) A coin is tossed 5 times, and 5 heads are observed. Would either a frequentist or a Bayesian infer that the coin has heads on both sides? Why or why not? [2]

2. Here we explore a regression model where the noise variance is a function of the input (variance increases as a function of input). Specifically

$$y = wx + \epsilon,$$

where the noise  $\epsilon$  is normally distributed with mean 0 and standard deviation  $\sigma x^2$ . The value of  $\sigma$  is assumed known and the input  $x$  is restricted to the interval  $[1, 4]$ . We can write the model more compactly as  $y \sim \mathcal{N}(wx, \sigma^2 x^4)$ .

- (a) How is the ratio  $y/x$  distributed for a fixed (constant)  $x$ ? [2]

(b) Suppose we now have  $N$  training points and targets  $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ , where each  $x_n$  is chosen at random from  $[1, 4]$  and the corresponding  $y_n$  is subsequently sampled from  $y_n \sim \mathcal{N}(wx_n, \sigma^2 x_n^4)$ . Obtain the maximum likelihood estimate for  $w$  as a function of the training data. [3]

(c) What are the bias (i.e., difference between expected and actual value) and variance of the estimator for  $w$  just obtained, as a function of  $N$  and  $\sigma^2$  for fixed inputs  $x_1, \dots, x_N$ ? Can you suggest a method for reducing the variance, even if it involves increasing the bias? [5]

(d) Now supposing I put a prior distribution on  $w$ :  $w \sim \mathcal{N}(0, \alpha^{-1})$ , for some fixed  $\alpha$ . Obtain the posterior distribution for  $w$ , given the same data set as above; also compute the maximum a posteriori estimate. [3]

(e) What are the bias and variance of this estimator? What do you infer from this about the role of  $\alpha$  in controlling the bias-variance tradeoff? [5]

(Some potentially useful relations: if  $z \sim \mathcal{N}(\mu, \sigma^2)$ , then  $az \sim \mathcal{N}(a\mu, \sigma^2 a^2)$  for fixed  $a$ . If  $z_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  and  $z_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$  and they are independent, then  $\text{Var}(z_1 + z_2) = \sigma_1^2 + \sigma_2^2$ .)

3. Here we will look at methods for selecting input features for a logistic regression model

$$P(y = 1 | \mathbf{x}, \mathbf{w}) = \sigma(w_0 + w_1 x_1 + w_2 x_2).$$

The available training examples are very simple, involving only binary valued inputs:

Number of copies	$x_1$	$x_2$	$y$
5	1	1	0
10	0	1	1
10	1	0	1
10	0	0	0

So, for example, there are 5 copies of  $\mathbf{x} = (1, 1)^T$  in the training set, all labeled  $y = 0$ . The correct label is actually a deterministic function of the two features:  $y = 0$  if  $x_1 = x_2$  and 1 otherwise. We define greedy selection in this context as follows: we start with no features (train only with  $w_0$ ) and successively try to add new features provided that each addition strictly improves the training log-likelihood. We use no other stopping criterion.

(a) Could greedy selection add either  $x_1$  or  $x_2$  in this case? [2]

(b) What is the classification error on the training examples that we could achieve by including both  $x_1$  and  $x_2$  in the logistic regression model? [2]

(c) Suppose we define another possible feature to include, a function of  $x_1$  and  $x_2$ . Which of the following features, if any, would permit us to correctly classify all the training examples when used in combination with  $x_1$  and  $x_2$  in the logistic regression model:  $x_1 - x_2$ ,  $x_1x_2$ ,  $x_2^2$ ? [3]

(d) Could the greedy selection method choose this feature as the first feature to add when the available features are  $x_1$ ,  $x_2$  and your choice of the new feature? [2]

4. We have seen in class that a soft-margin separating hyperplane of the form  $\mathbf{w}^T \mathbf{x} + b = 0$  can be obtained via optimising the following Lagrangian:

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i + \sum_{i=1}^N \mu_i [1 - \xi_i - y_i (\mathbf{w}^T \mathbf{x}_i + b)] - \sum_{i=1}^N \lambda_i \xi_i,$$

where  $\{(\mathbf{x}_1, y_1); (\mathbf{x}_2, y_2); \dots; (\mathbf{x}_N, y_N)\}$  denote the observed data points, with  $y_i \in \{-1, 1\}$ ;  $\boldsymbol{\xi} = \{\xi_1, \xi_2, \dots, \xi_N\}$  are the slack variables, and  $\boldsymbol{\mu} = \{\mu_1, \mu_2, \dots, \mu_N\}$  and  $\boldsymbol{\lambda} = \{\lambda_1, \lambda_2, \dots, \lambda_N\}$  are the Lagrange multipliers.

(a) By deriving and substituting in the values of the hyperplane parameters and the slack variables that minimise  $L$ , obtain the dual of this Lagrangian, as a function of just the multipliers  $\boldsymbol{\mu}$ . [5]

(b) Also obtain the constraints on  $\boldsymbol{\mu}$  for the dual. How do they differ from the hard-margin case? [2]

(c) Let  $f : \mathbb{R}^D \mapsto \mathbb{R}$  be a real-valued function, where  $D$  is the dimensionality of the  $\mathbf{x}_i$ . Is the function  $K(\mathbf{x}_1, \mathbf{x}_2) = f(\mathbf{x}_1)^2 + f(\mathbf{x}_2)^2$  guaranteed to be a valid kernel? Prove either way. [2]

5. Consider a simple example, where a burglar alarm at my house ( $A$ ) can be set off by a burglary ( $B$ ), or an earthquake ( $E$ ), or a hurricane ( $H$ ). I have two neighbours, John ( $J$ ) and Mary ( $M$ ), either of whom could call me in case the alarm goes off.

(a) Draw a Bayesian network to represent the causal relationships between these six binary random variables. [2]

(b) Write down the factorisation of the full joint distribution represented by your network. Also specify at least three of the conditional independencies implied by this factorisation. [2]

(c) Give an instance in this network of the *explaining away* property, i.e., when a particular variable is observed then another pair of variables which were previously independent, become conditionally dependent. [1]

(d) Show that if our model is such that the alarm always (deterministically) goes off whenever there is an earthquake:

$$\begin{aligned} P(A = 1 | B = 1, E = 1, H = 1) &= P(A = 1 | B = 0, E = 1, H = 1) \\ &= P(A = 1 | B = 1, E = 1, H = 0) = P(A = 1 | B = 0, E = 1, H = 0) = 1, \end{aligned}$$

then  $P(B = 1|A = 1, E = 1) = P(B = 1)$  and  $P(H = 1|A = 1, E = 1) = P(H = 1)$ , i.e., observing an earthquake provides a full explanation for the alarm. [3]

6. Consider a setting where, over 3 successive days, when I get back home in the evening I observe the grass on my lawn to be either *wet* or *dry*. Because I work far away from home, I could not observe what the daytime weather was like on those 3 days, but I know that each day it was either *sunny* or *rainy*. Suppose also that I know the following: if the weather was rainy, the probability of the grass being wet in the evening is 0.9; if it was sunny, this probability is 0.2 (there is a sprinkler which the gardener switches on sometimes); if it is rainy one day, then the probability of rain the next day is 0.3; if it is sunny, then the probability of rain the next day is 0.1; and finally, the probability of it being rainy to start with is 0.1.

(a) Draw an appropriate Hidden Markov Model to represent this situation. Specify clearly your notation for random variables, and the corresponding initialisation, emission, and transition probabilities. [2]

(b) Suppose my actual observations over the 3 days are  $\{wet, dry, wet\}$ . Based on this and my specified model, I wish to estimate the probability that the weather on the 3rd day was sunny. Use the forward-backward algorithm to compute this. Referring to the notation used in class, which  $\alpha$  value(s) do you need to evaluate for this purpose? Show the steps of the recursion involved in doing so. [4]

7. Suppose we have a data set  $\{(\mathbf{x}_1, t_1); (\mathbf{x}_2, t_2); \dots; (\mathbf{x}_N, t_N)\}$ , where the  $\mathbf{x}_i \in \mathbb{R}^n$  are  $n$ -dimensional feature vectors, and the  $t_i \in \{\mathcal{C}_1, \mathcal{C}_2\}$  are categorical class labels. Further suppose that we adopt the following model for the class priors and the class-conditional likelihoods:

$$\begin{aligned} p(\mathcal{C}_1) &= \theta, \\ p(\mathcal{C}_2) &= 1 - \theta, \\ p(\mathbf{x}|\mathcal{C}_k) &= \frac{1}{(2\pi)^{n/2}} \frac{1}{|\boldsymbol{\Sigma}_k|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}. \quad (k \in \{1, 2\}) \end{aligned}$$

Obtain an expression for the posterior distribution of the class label for a given data point  $\mathbf{x}$ . What kind of separation boundary does this specify between the classes? Under what conditions does this boundary become linear? [6]