#### CS 189 Spring 2013 Introduction to Machine Learning

# Midterm

- You have 1 hour 20 minutes for the exam.
- The exam is closed book, closed notes except your one-page crib sheet.
- Please use non-programmable calculators only.
- Mark your answers ON THE EXAM ITSELF. If you are not sure of your answer you may wish to provide a *brief* explanation. All short answer sections can be successfully answered in a few sentences AT MOST.
- For true/false questions, fill in the *True/False* bubble.
- For multiple-choice questions, fill in the bubbles for ALL CORRECT CHOICES (in some cases, there may be more than one). For a question with p points and k choices, every false positive wil incur a penalty of p/(k-1) points.

First name	
Last name	
SID	

For staff use only:			
Q1.	True/False	/14	
Q2.	Multiple Choice Questions	/21	
Q3.	Short Answers	/15	
	Total	/50	

## Q1. [14 pts] True/False

- (a) [1 pt] In Support Vector Machines, we maximize  $\frac{\|w\|^2}{2}$  subject to the margin constraints.  $\bigcirc$  True  $\bigcirc$  False
- (b) [1 pt] In kernelized SVMs, the kernel matrix K has to be positive definite.

   O True O False
- (c) [1 pt] If two random variables are independent, then they have to be uncorrelated. O True O False
- (d) [1 pt] Isocontours of Gaussian distributions have axes whose lengths are proportional to the eigenvalues of the covariance matrix.
  - $\bigcirc$  True  $\bigcirc$  False
- (e) [1 pt] The RBF kernel  $(K(x_i, x_j) = exp(-\gamma ||x_i x_j||^2))$  corresponds to an infinite dimensional mapping of the feature vectors.
  - $\bigcirc$  True  $\bigcirc$  False
- (f) [1 pt] If (X, Y) are jointly Gaussian, then X and Y are also Gaussian distributed.  $\bigcirc$  True  $\bigcirc$  False
- (g) [1 pt] A function f(x, y, z) is convex if the Hessian of f is positive semi-definite.  $\bigcirc$  True  $\bigcirc$  False
- (h) [1 pt] In a least-squares linear regression problem, adding an  $L_2$  regularization penalty cannot decrease the  $L_2$  error of the solution w on the training data.
  - $\bigcirc$  True  $\bigcirc$  False
- (i) [1 pt] In linear SVMs, the optimal weight vector w is a linear combination of training data points.  $\bigcirc$  True  $\bigcirc$  False
- (j) [1 pt] In stochastic gradient descent, we take steps in the exact direction of the gradient vector. O True O False
- (k) [1 pt] In a two class problem when the class conditionals P(x|y=0) and P(x|y=1) are modelled as Gaussians with different covariance matrices, the posterior probabilities turn out to be logistic functions.
  - ⊖ True ⊖ False
- (1) [1 pt] The perceptron training procedure is guaranteed to converge if the two classes are linearly separable.
   O True O False
- (m) [1 pt] The maximum likelihood estimate for the variance of a univariate Gaussian is unbiased.
   O True O False
- (n) [1 pt] In linear regression, using an  $L_1$  regularization penalty term results in sparser solutions than using an  $L_2$  regularization penalty term.
  - $\bigcirc$  True  $\bigcirc$  False

#### Q2. [21 pts] Multiple Choice Questions

(a) [2 pts] If  $X \sim \mathcal{N}(\mu, \sigma^2)$  and Y = aX + b, then the variance of Y is:

$$\bigcirc a\sigma^2 + b \qquad \bigcirc a^2\sigma^2 + b \qquad \bigcirc a\sigma^2 \qquad \bigcirc a^2\sigma^2$$

(b) [2 pts] In soft margin SVMs, the slack variables  $\xi_i$  defined in the constraints  $y_i(w^T x_i + b) \ge 1 - \xi_i$  have to be

$$\bigcirc <0 \qquad \bigcirc \le 0 \qquad \bigcirc >0 \qquad \bigcirc \ge 0$$

(c) [4 pts] Which of the following transformations when applied on  $X \sim \mathcal{N}(\mu, \Sigma)$  transforms it into an axis aligned Gaussian? ( $\Sigma = UDU^T$  is the spectral decomposition of  $\Sigma$ )

$$\bigcirc U^{-1}(X-\mu) \qquad \bigcirc UD(X-\mu)$$

$$\bigcirc (UD^{1/2})^{-1}(X-\mu) \qquad \bigcirc U(X-\mu) \qquad \bigcirc \Sigma^{-1}(X-\mu)$$

- (d) [2 pts] Consider the sigmoid function  $f(x) = 1/(1 + e^{-x})$ . The derivative f'(x) is
  - $\bigcirc f(x)\ln f(x) + (1 f(x))\ln(1 f(x)) \qquad \bigcirc f(x)(1 f(x))$  $\bigcirc f(x)\ln(1 f(x)) \qquad \bigcirc f(x)(1 + f(x))$
- (e) [2 pts] In regression, using an  $L_2$  regularizer is equivalent to using a \_\_\_\_\_ prior.
  - $\bigcirc \text{ Laplace, } 2\beta \exp(-|x|/\beta) \qquad \bigcirc \text{ Exponential, } \beta \exp(-x/\beta), \text{ for } x > 0$  $\bigcirc \text{ Gaussian with } \Sigma = cI, c \in R \qquad \bigcirc \Sigma \neq cI, c \in R)$
- (f) [2 pts] Consider a two class classification problem with the loss matrix given as  $\begin{pmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \end{pmatrix}$ . Note that  $\lambda_{ij}$  is the loss for classifying an instance from class j as class i. At the decision boundary, the ratio  $\frac{P(\omega_2|x)}{P(\omega_1|x)}$  is equal to:
  - $\bigcirc \ \frac{\lambda_{11} \lambda_{22}}{\lambda_{21} \lambda_{12}} \qquad \bigcirc \ \frac{\lambda_{11} \lambda_{21}}{\lambda_{22} \lambda_{12}} \qquad \bigcirc \ \frac{\lambda_{11} + \lambda_{22}}{\lambda_{21} + \lambda_{12}} \qquad \bigcirc \ \frac{\lambda_{11} \lambda_{12}}{\lambda_{22} \lambda_{21}}$
- (g) [2 pts] Consider the  $L_2$  regularized loss function for linear regression  $L(w) = \frac{1}{2} ||Y Xw||^2 + \lambda ||w||^2$ , where  $\lambda$  is the regularization parameter. The Hessian matrix  $\nabla_w^2 L(w)$  is
  - $\bigcirc X^T X \qquad \bigcirc 2\lambda X^T X \qquad \bigcirc X^T X + 2\lambda I \qquad \bigcirc (X^T X)^{-1}$
- (h) [2 pts] The geometric margin in a hard margin Support Vector Machine is
  - $\bigcirc \quad \frac{\|w\|^2}{2} \qquad \qquad \bigcirc \quad \frac{1}{\|w\|^2} \qquad \qquad \bigcirc \quad \frac{2}{\|w\|} \qquad \qquad \bigcirc \quad \frac{2}{\|w\|^2}$
- (i) [3 pts] Which of the following functions are convex?
  - $\bigcirc \sin(x) \qquad \bigcirc |x| \qquad \bigcirc \min(f_1(x), f_2(x)), \qquad \bigcirc \max(f_1(x), f_2(x)), \\ \text{where } f_1 \text{ and } f_2 \text{ are } \\ \text{convex} \qquad \text{convex}$

### Q3. [15 pts] Short Answers

(a) [4 pts] For a hard margin SVM, give an expression to calculate b given the solutions for w and the Lagrange multipliers  $\{\alpha_i\}_{i=1}^N$ .

- (b) Consider a Bernoulli random variable X with parameter p (P(X = 1) = p). We observe the following samples of X: (1, 1, 0, 1).
  - (i) [2 pts] Give an expression for the likelihood as a function of p.
  - (ii) [2 pts] Give an expression for the derivative of the negative log likelihood.
  - (iii) [1 pt] What is the maximum likelihood estimate of p?

•

(c) [6 pts] Consider the weighted least squares problem in which you are given a dataset  $\{\tilde{x}_i, y_i, w_i\}_{i=1}^N$ , where  $w_i$  is an importance weight attached to the  $i^{th}$  data point. The loss is defined as  $L(\beta) = \sum_{i=1}^N w_i (y_i - \beta^T x_i)^2$ . Give an expression to calculate the coefficients  $\tilde{\beta}$  in closed form. Hint: You might need to use a matrix W such that  $diag(W) = [w_1 w_2 \dots w_N]^T$  SCRATCH PAPER

SCRATCH PAPER