

**Preface:** This discussion is long because it is a compilation of relevant midterm and final questions from previous semesters. Work with your GSI in section to focus on the questions that will help you the most.

## 1 MLE, MAP, Linear Regression

### (a) (Sp16 Final)

(5) [3 pts] Lasso can be interpreted as least-squares linear regression where

- ☐ weights are regularized with the  $\ell_1$  norm
- ☐ the weights have a Gaussian prior
- ☐ weights are regularized with the  $\ell_2$  norm
- ☐ the solution algorithm is simpler

### (b) (Sp20 Midterm A)

(d) [4 pts] Suppose we perform least-squares linear regression, but we don't assume that all weight vectors are equally reasonable; instead, we use the maximum *a posteriori* method to impose a normally-distributed prior probability on the weights. Then we are doing

- ☐ A:  $L_2$  regularization
- ☐ C: logistic regression
- ☐ B: Lasso regression
- ☐ D: ridge regression

### (c) (Sp19 Midterm)

(d) [3 pts] Assuming we can find algorithms to minimize them, which of the following cost functions will encourage **sparse solutions** (i.e., solutions where many components of  $w$  are zero)?

- ☐  $\|Xw - y\|_2^2 + \lambda \|w\|_1$
- ☐  $\|Xw - y\|_2^2 + \lambda \cdot (\# \text{ of nonzero components of } w)$
- ☐  $\|Xw - y\|_2^2 + \lambda \|w\|_1^2$
- ☐  $\|Xw - y\|_2^2 + \lambda \|w\|_2^2$

### Q3. [10 pts] Maximum Likelihood Estimation

There are 5 balls in a bag. Each ball is either red or blue. Let  $\theta$  (an integer) be the number of blue balls. We want to estimate  $\theta$ , so we draw 4 balls **with replacement** out of the bag, replacing each one before drawing the next. We get “blue,” “red,” “blue,” and “blue” (in that order).

- (a) [5 pts] Assuming  $\theta$  is fixed, what is the likelihood of getting exactly that sequence of colors (expressed as a function of  $\theta$ )?
- (b) [3 pts] Draw a table showing (as a fraction) the likelihood of getting exactly that sequence of colors, for every value of  $\theta$  from zero to 5 inclusive.

$\theta$	$\mathcal{L}(\theta; \langle \text{blue, red, blue, blue} \rangle)$
0	?
1	?
2	?
3	?
4	?
5	?

- (c) [2 pts] What is the maximum likelihood estimate for  $\theta$ ? (Chosen among all integers; not among all real numbers.)

## 2 MVG, Gaussian Classification

### (a) (Sp20 Midterm B)

(b) [4 pts] Consider a random variable  $X \sim \mathcal{N}(\mu, \Sigma) \in \mathbb{R}^d$ , where the multivariate Gaussian probability density function (PDF) is axis-aligned,  $\Sigma$  is positive definite, and the standard deviation along coordinate axis  $i$  is  $\sigma_i$ . Select all that apply.

☐ A: The  $d$  features of  $X$  are uncorrelated but not necessarily independent

☐ C:  $\Sigma$  has a symmetric square root  $\Sigma^{1/2}$  with eigenvalues  $\sigma_1, \sigma_2, \dots, \sigma_d$

☐ B:  $\Sigma = \text{diag}(\sqrt{\sigma_1}, \sqrt{\sigma_2}, \dots, \sqrt{\sigma_d})$

☐ D:  $(X - \mu)^\top \Sigma^{-1} (X - \mu) \geq 0$

### (b) (Sp19 Midterm)

(i) [3 pts] Which of the following apply to **linear discriminant analysis**?

☐ You calculate the sample mean for each class

☐ It approximates the Bayes decision rule

☐ You calculate the sample covariance matrix using the mean of all the data points

☐ The model produced by LDA is never the same as the model produced by QDA

### (c) (Sp19 Midterm)

(s) [3 pts] Suppose you have a sample in which each point has  $d$  features and comes from class C or class D. The class conditional distributions are  $(X_i|y_i = C) \sim N(\mu_C, \sigma_C^2)$  and  $(X_i|y_i = D) \sim N(\mu_D, \sigma_D^2)$  for unknown values  $\mu_C, \mu_D \in \mathbb{R}^d$  and  $\sigma_C^2, \sigma_D^2 \in \mathbb{R}$ . The class priors are  $\pi_C$  and  $\pi_D$ . We use 0-1 loss.

☐ If  $\pi_C = \pi_D$  and  $\sigma_C = \sigma_D$ , then the Bayes decision rule assigns a test point  $z$  to the class whose mean is closest to  $z$ .

☐ If  $\sigma_C = \sigma_D$ , then the Bayes decision boundary is always linear.

☐ If  $\pi_C = \pi_D$ , then the Bayes decision rule is  $r^*(z) = \underset{A \in \{C, D\}}{\text{argmin}} \left( |z - \mu_A|^2 / (2\sigma_A^2) + d \ln \sigma_A \right)$

☐ If  $\sigma_C = \sigma_D$ , then QDA will always produce a linear decision boundary when you fit it to your sample.

## 4 Watermelons

1. (6 points) Finn lives on a watermelon farm and wants to classify whether a watermelon is sweet (labelled as  $y = 1$ ) or not sweet (labelled as  $y = 0$ ). Finn observes a  $d$ -dimensional feature vector  $x \in \mathbb{R}^d$  associated with the appearance and the smell of each watermelon. Before observing a watermelon, Finn's general prior is that a watermelon is sweet with probability  $p(y = 1) = \pi_1$  and not sweet with probability  $p(y = 0) = 1 - \pi_1$ .

Finn is a watermelon expert and knows that the conditional probability distribution of the watermelon features  $p(x|y = k)$  for class  $k$  (where  $k = 0, 1$ ) is a  $d$ -dimensional Gaussian distribution  $N(\mu_k, \Sigma)$  with mean  $\mu_k \in \mathbb{R}^d$  and covariance  $\Sigma \in \mathbb{R}^{d \times d}$ . Note that the same covariance matrix  $\Sigma$  is shared between the two classes.

$$f(x|y = k) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left\{-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1}(x - \mu_k)\right\} \quad (8)$$

Can you help Finn find out how likely a given watermelon is sweet? Write down a simplified expression of  $p(y = 1|x)$  as a function of  $x, \mu_0, \mu_1, \Sigma$ , and  $\pi_1$ . The expression should be in the form of  $s(w^T x + b)$  where  $s(z) = \frac{1}{1+e^{-z}}$  is the sigmoid function. Write down what  $w$  and  $b$  should be.

### 3 Logistic Regression, Optimization

(a) (Sp19 Midterm)

(e) [3 pts] Which of the following statements about **logistic regression** are correct?

- |   |   |
|---|---|
| <input type="radio"/> The cost function of logistic regression is convex              | <input type="radio"/> The cost function of logistic regression is concave   |
| <input type="radio"/> Logistic regression uses the squared error as the loss function | <input type="radio"/> Logistic regression assumes that each class's points are generated from a Gaussian distribution |

(b) (Sp20 Midterm B)

(d) [4 pts] For classification problems with two features ( $d = 2$ , test point  $z \in \mathbb{R}^2$ ), which of the following methods have posterior probability distributions of the form  $P(Y|X = z) = s(Az_1^2 + Bz_2^2 + Cz_1z_2 + Dz_1 + Ez_2 + F)$  where  $s$  is the logistic function  $s(\gamma) = \frac{1}{1+e^{-\gamma}}$  and  $A, B, C, D, E, F \in \mathbb{R}$  can all be nonzero?

- |   |   |
|---|---|
| <input type="radio"/> A: Logistic regression with linear features                   | <input type="radio"/> C: Logistic regression with quadratic features                |
| <input type="radio"/> B: Linear discriminant analysis (LDA) with quadratic features | <input type="radio"/> D: Quadratic discriminant analysis (QDA) with linear features |

(c) (Sp21 Midterm)

(c) [4 pts] Recall the logistic function  $s(\gamma)$  and its derivative  $s'(\gamma) = \frac{d}{d\gamma} s(\gamma)$ . Let  $\gamma^*$  be the value of  $\gamma$  that maximizes  $s'(\gamma)$ .

- |  |  |
|--|--|
| <input type="radio"/> A: $\gamma^* = 0.25$   | <input type="radio"/> C: $s'(\gamma^*) = 0.5$  |
| <input type="radio"/> B: $s(\gamma^*) = 0.5$ | <input type="radio"/> D: $s'(\gamma^*) = 0.25$ |

(d) (Sp19 Midterm)

(r) [3 pts] Let  $L_i(w)$  be the loss corresponding to a sample point  $X_i$  with label  $y_i$ . The update rule for **stochastic gradient descent** with step size  $\epsilon$  is

- |   |   |
|---|---|
| <input type="radio"/> $w_{\text{new}} \leftarrow w - \epsilon \nabla_{X_i} L_i(w)$              | <input type="radio"/> $w_{\text{new}} \leftarrow w - \epsilon \nabla_w L_i(w)$              |
| <input type="radio"/> $w_{\text{new}} \leftarrow w - \epsilon \sum_{i=1}^n \nabla_{X_i} L_i(w)$ | <input type="radio"/> $w_{\text{new}} \leftarrow w - \epsilon \sum_{i=1}^n \nabla_w L_i(w)$ |

## Q5. [10 pts] Logistic Regression with One Feature

We are given another sample in which each point has only one feature. Consider a binary classification problem in which sample values  $x \in \mathbb{R}$  are drawn randomly from two different class distributions. The first class, with label  $y = 0$ , has its mean to the left of the mean of the second class, with label  $y = 1$ . We will use a modified version of logistic regression to classify these data points. We model the posterior probability at a test point  $z \in \mathbb{R}$  as

$$P(y = 1|z) = s(z - \alpha),$$

where  $\alpha \in \mathbb{R}$  is the sole parameter we are trying to learn and  $s(\gamma) = 1/(1 + e^{-\gamma})$  is the logistic function. The decision boundary is  $z = \alpha$  (because  $s(z) = \frac{1}{2}$  there).

We will learn the parameter  $\alpha$  by performing gradient descent on the logistic loss function (a.k.a. cross-entropy). That is, for a data point  $x$  with label  $y \in \{0, 1\}$ , we find the  $\alpha$  that minimizes

$$J(\alpha) = -y \ln s(x - \alpha) - (1 - y) \ln(1 - s(x - \alpha)).$$

- (a) [5 pts] Derive the stochastic gradient descent update for  $J$  with step size  $\epsilon > 0$ , given a sample value  $x$  and a label  $y$ . Hint: feel free to use  $s$  as an abbreviation for  $s(x - \alpha)$ .

- (b) [3 pts] Is  $J(\alpha)$  convex over  $\alpha \in \mathbb{R}$ ? Justify your answer.

- (c) [2 pts] Now we consider multiple sample points. As  $d = 1$ , we are given an  $n \times 1$  design matrix  $X$  and a vector  $y \in \mathbb{R}^n$  of labels. Consider batch gradient descent on the cost function  $\sum_{i=1}^n J(\alpha; X_i, y_i)$ . There are circumstances in which this cost function does not have a minimum over  $\alpha \in \mathbb{R}$  at all. What is an example of such a circumstance?

## 4 SVM

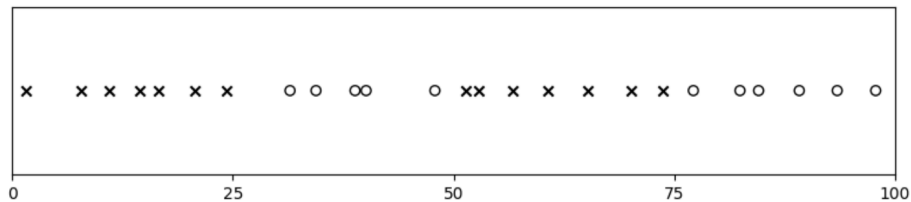
(a) (Sp16 Final)

(14) [3 pts] Which of the following can help to reduce overfitting in an SVM classifier?

- ☐ Use of slack variables
- ☐ High-degree polynomial features
- ☐ Normalizing the data
- ☐ Setting a very low learning rate

(b) (Sp19 Midterm)

(k) [3 pts] Suppose you are given the one-dimensional data  $\{x_1, x_2, \dots, x_{25}\}$  illustrated below and you have only a **hard-margin support vector machine** (with a fictitious dimension) at your disposal. Which of the following modifications can give you 100% training accuracy?



- ☐ Centering the data
- ☐ Add a feature  $x_i^2$
- ☐ Add a feature that is 1 if  $x \leq 50$ , or  $-1$  if  $x > 50$
- ☐ Add two features,  $x_i^2$  and  $x_i^3$

(c) (Sp21 Midterm)

(b) [4 pts] Which of the following changes would commonly cause an SVM's margin  $1/\|w\|$  to shrink?

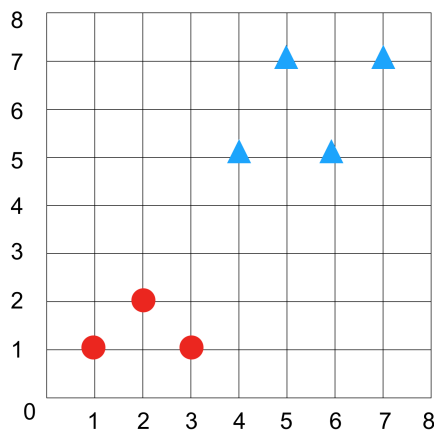
- ☐ A: Soft margin SVM: increasing the value of  $C$
- ☐ C: Soft margin SVM: decreasing the value of  $C$
- ☐ B: Hard margin SVM: adding a sample point that violates the margin
- ☐ D: Hard margin SVM: adding a new feature to each sample point

## Q2. [20 pts] Hard-Margin Support Vector Machines

Recall that a **maximum margin classifier**, also known as a hard-margin support vector machine (SVM), takes  $n$  training points  $X_1, X_2, \dots, X_n \in \mathbb{R}^d$  with labels  $y_1, y_2, \dots, y_n \in \{+1, -1\}$ , and finds parameters  $w \in \mathbb{R}^d$  and  $\alpha \in \mathbb{R}$  that satisfy a certain objective function subject to the constraints

$$y_i(X_i \cdot w + \alpha) \geq 1, \quad \forall i \in \{1, \dots, n\}.$$

For parts (a) and (b), consider the following training points. Circles are classified as positive examples with label  $+1$  and triangles are classified as negative examples with label  $-1$ .



- (a) [3 pts] Which points are the support vectors? Write it as  $\begin{bmatrix} \text{horizontal} \\ \text{vertical} \end{bmatrix}$ . E.g., the bottom right circle is  $\begin{bmatrix} 3 \\ 1 \end{bmatrix}$ .
- (b) [4 pts] If we add the sample point  $x = \begin{bmatrix} 5 \\ 1 \end{bmatrix}$  with label  $-1$  (triangle) to the training set, which points are the support vectors?

For parts (c)–(f), forget about the figure above, but assume that there is at least one sample point in each class and that the sample points are linearly separable.

- (c) [2 pts] Describe the geometric relationship between  $w$  and the decision boundary.
- (d) [2 pts] Describe the relationship between  $w$  and the margin. (For the purposes of this question, the margin is just a number.)
- (e) [4 pts] Knowing what you know about the hard-margin SVM objective function, explain why for the optimal  $(w, \alpha)$ , there must be at least one sample point for which  $X_i \cdot w + \alpha = 1$  and one sample point for which  $X_i \cdot w + \alpha = -1$ .
- (f) [5 pts] If we add new features to the sample points (while retaining all the original features), can the optimal  $\|w_{\text{new}}\|$  in the enlarged SVM be greater than the optimal  $\|w_{\text{old}}\|$  in the original SVM? Can it be smaller? Can it be the same? Explain why! (Most of the points will be for your explanation.)



## 5 Featurization, Kernels

### (a) (Sp19 Midterm)

(q) [3 pts] You are given four sample points  $X_1 = [-1, -1]^\top$ ,  $X_2 = [-1, 1]^\top$ ,  $X_3 = [1, -1]^\top$ , and  $X_4 = [1, 1]^\top$ . Each of them is in class C or class D. For what feature representations are the lifted points  $\Phi(X_i)$  *guaranteed* to be **linearly separable** (with no point lying exactly on the decision boundary) for every possible class labeling?

☐  $\Phi(x) = [x_1, x_2, 1]$

☐  $\Phi(x) = [x_1^2, x_2^2, x_1, x_2, 1]$

☐  $\Phi(x) = [x_1, x_2, x_1^2 + x_2^2, 1]$

☐  $\Phi(x) = [x_1^2, x_2^2, x_1 x_2, x_1, x_2, 1]$

### (b) (Sp20 Final)

(10) [4 pts] Which of the following statement(s) about **kernels** are true?

☐ A: The dimension of the lifted feature vectors  $\Phi(\cdot)$ , whose inner products the kernel function computes, can be infinite.

☐ B: For any desired lifting  $\Phi(x)$ , we can design a kernel function  $k(x, z)$  that will evaluate  $\Phi(x)^\top \Phi(z)$  more quickly than explicitly computing  $\Phi(x)$  and  $\Phi(z)$ .

☐ C: The kernel trick, when it is applicable, speeds up a learning algorithm if the number of sample points is substantially less than the dimension of the (lifted) feature space.

☐ D: If the raw feature vectors  $x, y$  are of dimension 2, then  $k(x, y) = x_1^2 y_1^2 + x_2^2 y_2^2$  is a valid kernel.

### (c) (Sp16 Final)

(7) [3 pts] The kernel trick

☐ can be applied to every classification algorithm

☐ is commonly used for dimensionality reduction

☐ changes ridge regression so we solve a  $d \times d$  linear system instead of an  $n \times n$  system, given  $n$  sample points with  $d$  features

☐ exploits the fact that in many learning algorithms, the weights can be written as a linear combination of input points

## Q4. [10 pts] Kernels

(1) [2 pts] What is the primary motivation for using the kernel trick in machine learning algorithms?

(2) [4 pts] Prove that for every design matrix  $X \in \mathbb{R}^{n \times d}$ , the corresponding kernel matrix is positive semidefinite.

(3) [2 pts] Suppose that a regression algorithm contains the following line of code.

$$\mathbf{w} \leftarrow \mathbf{w} + X^\top M X X^\top \mathbf{u}$$

Here,  $X \in \mathbb{R}^{n \times d}$  is the design matrix,  $\mathbf{w} \in \mathbb{R}^d$  is the weight vector,  $M \in \mathbb{R}^{n \times n}$  is a matrix unrelated to  $X$ , and  $\mathbf{u} \in \mathbb{R}^n$  is a vector unrelated to  $X$ . We want to derive a dual version of the algorithm in which we express the weights  $\mathbf{w}$  as a linear combination of samples  $X_i$  (rows of  $X$ ) and a dual weight vector  $\mathbf{a}$  contains the coefficients of that linear combination. Rewrite the line of code in its dual form so that it updates  $\mathbf{a}$  correctly (and so that  $\mathbf{w}$  does not appear).

(4) [2 pts] Can this line of code for updating  $\mathbf{a}$  be kernelized? If so, show how. If not, explain why.

## 6 Bias-Variance Decomposition, Decision Theory

### (a) (Sp16 Final)

(20) [3 pts] How does the bias-variance decomposition of a ridge regression estimator compare with that of ordinary least squares regression? (Select one.)

- |  |   |
|--|---|
| <input type="checkbox"/> Ridge has larger bias, larger variance  | <input type="checkbox"/> Ridge has smaller bias, larger variance  |
| <input type="checkbox"/> Ridge has larger bias, smaller variance | <input type="checkbox"/> Ridge has smaller bias, smaller variance |

### (b) (Sp19 Midterm)

(j) [3 pts] Which of the following are reasons why you might adjust your model in ways that increase the bias?

- |   |  |
|---|--|
| <input type="radio"/> You observe high training error and high validation error | <input type="radio"/> You observe low training error and high validation error |
| <input type="radio"/> You have few data points                                  | <input type="radio"/> Your data are not linearly separable                     |

### (c) (Sp20 Midterm A)

(e) [4 pts] Which of the following statements regarding ROC curves are true?

- |   |  |
|---|--|
| <input type="radio"/> A: the ROC curve is monotonically increasing  | <input type="radio"/> C: the ROC curve is concave  |
| <input type="radio"/> B: for a logistic regression classifier, the ROC curve's horizontal axis is the posterior probability used as a threshold for the decision rule | <input type="radio"/> D: if the ROC curve passes through $(0, 1)$ , the classifier is always correct (on the test data used to make the ROC curve) |

## Q4. [20 pts] Finding Bias, Variance, and Risk

For  $z \in \mathbb{R}$ , you are trying to estimate a true function  $g(z) = 2z^2$  with **least-squares regression**, where the regression function is a line  $h(z) = wz$  that goes through the origin and  $w \in \mathbb{R}$ . Each sample point  $x \in \mathbb{R}$  is drawn from the **uniform distribution on**  $[-1, 1]$  and has a corresponding label  $y = g(x) \in \mathbb{R}$ . There is no noise in the labels. We train the model with **just one sample point**! Call it  $x$ , and assume  $x \neq 0$ . We want to apply the bias-variance decomposition to this model.

- (a) [3 pts] In one sentence, why do we expect the bias to be large?
- (b) [6 pts] What is the bias of your model  $h(z)$  as a function of a test point  $z \in \mathbb{R}$ ? (*Hint: start by working out the value of the least-squares weight  $w$ .*) Your final bias should not include an  $x$ ; work out the expectation.
- (c) [6 pts] What is the variance of your model  $h(z)$  as a function of a test point  $z \in \mathbb{R}$ ? Your final variance should not include an  $x$ ; work out the expectation.
- (d) [5 pts] Let  $R(h, z)$  be the risk (expected loss) for a fixed, arbitrary test point  $z \in \mathbb{R}$  with the noise-free label  $g(z)$  (where the expectation is taken over the distribution of values of  $(x, y)$ ). What is the mathematical relationship between the risk  $R(h, z)$ , the bias of  $h(z)$  at  $z$ , and the variance of  $h(z)$  at  $z$ ? What are the values (as numbers) of these three quantities for  $z = 1$ ?