

- Please do not open the exam before you are instructed to do so. Fill out the blanks below now.
- **Electronic devices are forbidden on your person**, including phones, laptops, tablet computers, headphones, and calculators. Turn your cell phone off and **leave all electronics at the front of the room**, or **risk getting a zero** on the exam. Exceptions are made for car keys and devices needed because of disabilities.
- When you start, the **first thing you should do is check that you have all 7 pages and all 4 questions**. The second thing is to please **write your initials at the top right of every page after this one** (e.g., write “JS” if you are Jonathan Shewchuk).
- The exam is closed book, closed notes except your one cheat sheet.
- You have **80 minutes**. (If you are in the Disabled Students’ Program and have an allowance of 150% or 200% time, that comes to 120 minutes or 160 minutes, respectively.)
- Mark your answers on the exam itself in the space provided. Do **not** attach any extra sheets. If you run out of space for an answer, write a note that your answer is continued on the back of the page.
- The total number of points is 100. There are 12 multiple choice questions worth 4 points each, and 3 written questions worth a total of 52 points.
- For multiple answer questions, fill in the bubbles for **ALL correct choices**: there may be more than one correct choice, but there is always at least one correct choice. **NO partial credit** on multiple answer questions: the set of all correct answers must be checked.

First name	
Last name	
SID	
Name and SID of student to your left	
Name and SID of student to your right	

Q1. [48 pts] Multiple Answer

Fill in the bubbles for **ALL correct choices**: there may be more than one correct choice, but there is always at least one correct choice. **NO partial credit**: the set of all correct answers must be checked.

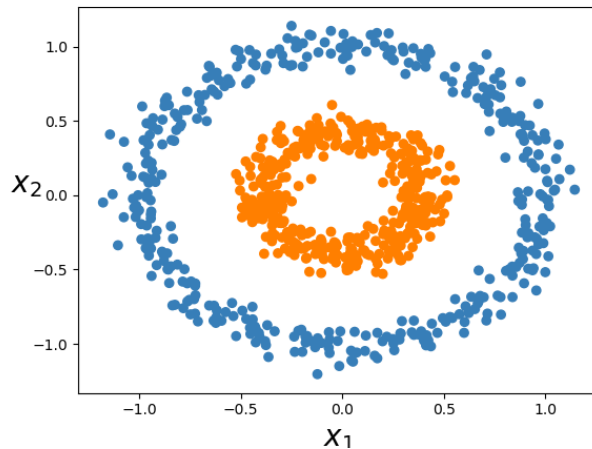
(a) [4 pts] We seek to find $w \in \mathbb{R}^d$ that minimizes a real-valued cost function $J(w)$. We know that J is continuous and smooth, and it has one and only one global minimum. (There are no other constraints on J .) Select the true statements about gradient descent on J .

- A: A step of gradient descent is $w \leftarrow w + \epsilon \nabla J(w)$, where $\epsilon > 0$ is the step size.
- B: The gradient descent algorithm will always converge to the global minimum of J if the step size ϵ is sufficiently small.
- C: If the global minimum of J is at the vector w^* , steps of gradient descent on J starting from $w = w^*$ will never change w .
- D: A step of gradient descent never causes $J(w)$ to increase.

(b) [4 pts] Which statements are true for every symmetric, real matrix $S \in \mathbb{R}^{n \times n}$?

- A: All the eigenvalues of S are real.
- B: S can be written as $S = A^2$, where A is symmetric and belongs to $\mathbb{R}^{n \times n}$.
- C: If S is positive semidefinite, then S is invertible.
- D: If all the eigenvalues of S are strictly less than zero, then S is invertible.

(c) [4 pts] You are given a two-class classification problem with the training points below. For each feature below, select it if adding it as a third feature (alongside x_1 and x_2) would make the two classes linearly separable.



- A: x_2^2
- B: $\|x\|_2$
- C: $\|x\|_2^3$
- D: $x_1 + x_2$

(d) [4 pts] Which statements are true of **Gaussian discriminant analysis for two-class classification**, specifically quadratic discriminant analysis (QDA) and linear discriminant analysis (LDA)? (Assume that there are no added features.)

- A: QDA for isotropic Gaussians (i.e., with the same variance in all directions) becomes the centroid method when the prior probabilities of the two classes are equal.
- B: QDA is more likely to overfit than LDA when the number of training points is small.
- C: LDA for isotropic Gaussians (i.e., with the same variance in all directions) becomes the centroid method when the prior probabilities of the two classes are equal.
- D: LDA for anisotropic Gaussians can produce nonlinear decision boundaries.

(i) [4 pts] Which statements are true about **ridge regression and Lasso** (with $\lambda > 0$ for both).

A: Ridge regression has a unique solution if and only if the design matrix has full rank.

C: Ridge regression can be formulated as a linear programming problem.

B: There are points in feature space where the gradient of Lasso's cost function is not defined.

D: One of Lasso's virtues is its tendency to set some weights to zero.

(j) [4 pts] Two classes of observations are drawn from two **univariate normal** distributions: $D_1 = \mathcal{N}(\mu_1, \sigma_1^2)$ for class 1 and $D_2 = \mathcal{N}(\mu_2, \sigma_2^2)$ for class 2. We know the parameters and prior probabilities of each class (the priors may or may not be equal), and we construct their Bayes classifier (with a 0-1 loss function). Which statements are true of the **Bayes optimal decision boundary**?

A: It might be \emptyset (no points).

C: It might have exactly three points.

B: It might have exactly two points.

D: It might have exactly ten points.

(k) [4 pts] Which of the following classifiers are guaranteed to assign the same classes to the test data if we apply to all points (training and test points) an invertible linear transformation that **whitens** the training points? (By "the same classes," we mean the same predictions as if we didn't whiten the data.)

A: Soft-margin support vector machine

C: Quadratic discriminant analysis

B: k -nearest neighbor classifier

D: Linear discriminant analysis

(l) [4 pts] Consider a **continuous uniform distribution** $\mathcal{U}[0, b]$, from which we draw a random real number between 0 and b . The probability density function (PDF) of $\mathcal{U}[0, b]$ is

$$f(x) = \begin{cases} 1/b & 0 \leq x \leq b, \\ 0 & \text{otherwise.} \end{cases}$$

We want to use **maximum likelihood estimation** to estimate the parameter b . We draw three points at random from $\mathcal{U}[0, b]$ and obtain $x_1 = 44.4$, $x_2 = 8$, and $x_3 = 41.2$. What is the maximum likelihood estimate \hat{b} of b ?

A: $\hat{b} = 44.4 + 8 + 41.2$.

C: $\hat{b} = (44.4 + 8 + 41.2)/3$.

B: $\hat{b} = 44.4$.

D: $\hat{b} = \frac{4}{3} \cdot 44.4$.

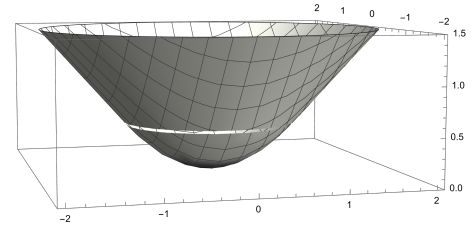
Extra space: if you need extra space for your answer to a written problem on pages 5–7, you may write here. **Be sure to write "see page 4" under the unfinished answer!**

Q2. [18 pts] Optimizing Huber Loss

Given a vector prediction $z \in \mathbb{R}^k$, a vector true label $y \in \mathbb{R}^k$, a fixed constant $\delta > 0$, and the ℓ_2 -norm $\|v\| = \sqrt{v^\top v}$, the Huber ℓ_2 -loss function is

$$L_\delta(z, y) = \begin{cases} \frac{1}{2}\|z - y\|^2, & \|z - y\| \leq \delta, \\ \delta \cdot (\|z - y\| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases}$$

At right is a plot of $L_1(z, 0)$ for $k = 2$. This loss can be used for regression where the regression function returns a k -dimensional vector.



The Huber ℓ_2 -loss is designed to be similar to the loss function $\|z - y\|$ (i.e., the Euclidean distance, which in 1D we call the absolute loss); but unlike the Euclidean distance, it is smooth at the minimum, $z = y$. For a fixed y , the Huber ℓ_2 -loss is quadratic in z in a small region near y , but it is shaped like a cone farther away from y . The Huber ℓ_2 -loss is continuous and convex.

- (a) [7 pts] Compute the gradient $\nabla_z L_\delta(z, y)$ of the ℓ_2 -Huber loss (for a fixed δ and y).
- (b) [4 pts] If we optimize z with gradient descent on L_δ (for a fixed y), **what learning rate (step size) ϵ** guarantees that we will eventually reach the exact minimum (rather than just inching closer and closer forever)? **Why?**
- (c) [2 pts] Suppose we use Newton's method to find a z that minimizes the ℓ_2 -Huber loss. (Technically, the Hessian of L_δ with respect to z is not defined where $\|z - y\| = \delta$, but we fix that by simply using the Hessian of $\frac{1}{2}\|z - y\|^2$ at those points.)
If we start at a point $z = z_0$ that satisfies $\|z_0 - y\| < \delta$, what will the value of z be after one step of Newton's method?
Why?
- (d) [3 pts] If we start at a point $z = z_0$ that satisfies $\|z_0 - y\| > \delta$, what will one step of Newton's method do? **Why?** (Note: for this question and the next one, we want a qualitative answer; you don't need to calculate a Hessian.)
- (e) [2 pts] Suppose we add an ℓ_2 regularization term $\lambda\|z\|^2$ to the ℓ_2 -Huber loss, with $\lambda > 0$, and perform one step of Newton's method on the ℓ_2 -regularized ℓ_2 -Huber loss. How do your answers to (c) and (d) change (qualitatively), and **why?**

Q3. [17 pts] Quadratic Discriminant Analysis

We want to predict whether a person prefers vanilla or chocolate ice cream based on a single feature: their age. We suspect that the ages of vanilla-lovers are normally distributed, and so are the ages of chocolate-lovers, so we build a classifier with **quadratic discriminant analysis (QDA)** and a **0-1 loss function**. Our survey of 13 random people turns up 8 vanilla lovers and 5 chocolate lovers of the following ages.

Vanilla: [21, 26, 27, 28, 30, 30, 31, 31]

Chocolate: [15, 18, 21, 22, 24]

- (a) [17 pts] Please do QDA. Determine the **distribution parameters and prior probabilities** of vanilla lovers and chocolate lovers (as exact, simplified integers or fractions). Then determine the **probability that a person of age x prefers vanilla over chocolate** (substituting the numbers so your answer is an exact, simplified function of x , which can include logistic functions $s(\cdot)$ or exponentials). Also, determine the **decision boundary** (as one or more numbers written as simplified expressions, possibly with logarithms and fractions). Show all your work! (Hint: as MNIST taught us, $28^2 = 784$.)

Q4. [17 pts] Estimating the Noise in Linear Regression

In Lecture 12 we suggested a model of reality in which we want to determine a linear natural law $g(z) = v^\top z$ ($g =$ “ground truth”) mapping each data point $z \in \mathbb{R}^d$ to a label in \mathbb{R} . (For simplicity, we don’t use a bias term α in this question; assume our natural law satisfies $g(0) = 0$.) But the measurements are noisy, so what we get is an $n \times d$ design matrix X and a vector $y \in \mathbb{R}^n$ of labels such that $y_i = v^\top X_i + \epsilon_i$, where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ is random noise (and each ϵ_i is independent of the others). In class we applied maximum likelihood estimation (MLE) to justify using the mean of the squared losses as the cost function for linear regression to compute a weight vector $w \in \mathbb{R}^d$ that is an estimate of v . Now we will use MLE to estimate σ^2 , the variance of the measurement noise.

- (a) [5 pts] Write the **likelihood function** $\mathcal{L}(\sigma; y, X, v)$ for obtaining the labels y_i given fixed values of X and v . (Note: for the purposes of this problem, X and v are **not** random. There should be no μ or other unlisted parameters in your answer.)
- (b) [6 pts] Write the **log likelihood function** $\ell(\sigma; y, X, v)$ and **find the value of σ^2 that maximizes ℓ** . Show your work. (Note: you do not need to prove it’s a maximum.)
- (c) [2 pts] What formula that you’re familiar with does your optimal value of σ^2 look like? (“The mean variance of the labels y_i ” doesn’t count. It’s something else too.)
- (d) [4 pts] Unfortunately, we don’t know the value of v . How should we estimate σ^2 , given that we cannot obtain v ? **Write an estimate of your estimate for σ^2 expressed solely in terms of X , y , and n , with no v** . (This is your maximum likelihood estimator $\hat{\sigma}^2$ for the true σ^2 .) For full points, write your final answer in matrix notation with no summation. You may assume that X has rank d .