

Midterm Exam

- Total duration: 90 minutes.
- You **can** use one page as a cheat sheet.
- You **cannot** consult your notes, textbooks, Google, or any other form of external help.
- Maximum points: 60. Any score above 60 will be rounded to 60.
- Once you are finished, please scan, or take a picture of, your answers and upload on NYUClasses before 8pm ET. You will have to include your cheat sheet, if you used one. No late submissions will be accepted.
- Good luck and stay safe!

-
1. **(3 points)** Please write down the time at the *start* and *end* of your exam. The difference should not exceed 90 minutes. Please also write down your *name* and *signature* below; by doing so, you are affirming the NYU Tandon School of Engineering student code of conduct.

2. **(10 points)** This is a slight variant of a homework problem. Let $\{x_1, x_2, \dots, x_n\}$ be a set of points in d -dimensional space, and let $\{p_1, p_2, \dots, p_n\}$ denote a probability distribution over the integers $[1, 2, \dots, n]$. Suppose we wish to produce a single point estimate $\mu \in \mathbb{R}^d$ that minimizes the *weighted* squared-error:

$$L(\mu) = p_1 \|x_1 - \mu\|_2^2 + p_2 \|x_2 - \mu\|_2^2 + \dots + p_n \|x_n - \mu\|_2^2$$

Find a closed form expression for the optimal μ and prove that your answer is correct.

3. **(15 points)** This is a variant of a homework problem. Suppose x is a d -dimensional input, w is a d -dimensional variable, and λ is a regularization parameter.

a. Show that the minimizer of the squared-error loss with ℓ_1 regularizer:

$$L(w) = \frac{1}{2} \|x - w\|_2^2 + \lambda \|w\|_1$$

is given by:

$$w_i^* = \begin{cases} x_i - \lambda & \text{if } x_i > \lambda, \\ x_i + \lambda & \text{if } x_i < -\lambda, \\ 0 & \text{otherwise.} \end{cases}$$

b. Show that the minimizer of the squared-error loss with ℓ_2 regularizer:

$$L(w) = \frac{1}{2} \|x - w\|_2^2 + \lambda \|w\|_2^2$$

is given by:

$$w_i^* = \left(\frac{1}{1 + 2\lambda} \right) x_i.$$

- c. In class, we argued via contour plots that greater ℓ_2 regularization encourages “small” solutions, while greater ℓ_1 regularization encourages “sparse” solutions. Mathematically justify why that is the case by examining the structure of the optimal solutions derived above.

4. **(10 points)** The following represents python code for an algorithm that attempts to perform linear regression. (a) Identify the algorithm. (b) Explain why this algorithm may not converge as implemented below, and identify the line in the algorithm that makes this happen. (c) Suggest a way to fix this algorithm.

```
def optim_alg(init, steps, grad):
    xs = [init]
    for step in steps:
        xs.append(xs[-1] - step * grad(xs[-1]))
    return xs

def linear_reg_grad(X, y, w):
    return X.T.dot(X.dot(w) - y)

input_to_optim_alg = lambda w: linear_reg_grad(X, y, w)
learning_rates = np.arange(start=1, stop=300, step=3)
ws = optim_alg(w0, learning_rates, input_to_optim_alg)
```

5. **(10 points)** To combat the COVID-19 pandemic, an enterprising NYU Tandon graduate student decides to build a logistic regression model to predict the conditional likelihood of a person being one of two states – *infected* or *clear* – based on daily forehead temperature measurements over the last 30 days. Fortunately, a dataset of such measurements for a population of 100,000 persons is available.
- Identify the parameters of the problem (number of samples n , data dimension d , number of classes k .)
 - If X and y denote the arrays that encode the training data points and labels, what are the sizes of X and y ?
 - Starting from the definition of conditional likelihood, derive the loss function used to train the model. You can assume the probabilities can be modeled as a sigmoid function.

6. **(15 points)** Suppose we are given real-valued scalar data (i.e., $d = 1$) belonging to one of two classes. We are given a set of three data samples with negative labels, $X_- = \{0, 1, -1\}$, and a set of three data samples with positive labels, $X_+ = \{-3, 3, -2\}$. Our goal is to build a classifier for this dataset. We will show that kernel methods are particularly useful in this case.
- Argue that no perfect linear separator in the original space can exist.
 - Argue that if the data is mapped via the two-dimensional feature mapping $\phi(u) = (u, u^2)$, then a perfect linear separator exists.
 - Explicitly draw the maximum-margin linear separator in the new feature space, and mark the closest points nearest to this linear separator.
 - Calculate the equation of the maximum-margin separator.