

**Due: Friday, March 10 at 11:59 PM PST**

- Homework 4 consists of coding assignments and math problems.
- We prefer that you typeset your answers using  $\text{\LaTeX}$  or other word processing software. If you haven't yet learned  $\text{\LaTeX}$ , one of the crown jewels of computer science, now is a good time! Neatly handwritten and scanned solutions will also be accepted.
- In all of the questions, **show your work**, not just the final answer.
- **We will not provide points back with respect to homework submission errors.** This includes, but is not limited to: 1) not assigning pages to problems; 2) not including code in the write-up appendix; 3) not including code in the designated code Gradescope assignment; 4) not including Kaggle scores; 5) submitting code that only partially works; 6). submitting late regrade requests. **Please carefully read and follow the HW submission guidelines/reminders for Pages 1, 2, and 10 of HW 4.**
- **Start early; you can submit models to Kaggle only twice a day!**

**Deliverables:**

1. Submit your predictions for the test sets to Kaggle as early as possible. Include your Kaggle scores in your write-up. The Kaggle competition for this assignment can be found at
  - WINE: <https://www.kaggle.com/competitions/spring23-cs189-hw4-wine/>
2. Write-up: Submit your solution in **PDF** format to “Homework 4 Write-Up” in Gradescope.
  - On the first page of your write-up, please list students with whom you collaborated
  - Start each question on a new page. If there are graphs, include those graphs on the same pages as the question write-up. **DO NOT** put them in an appendix. We need each solution to be self-contained on pages of its own.
  - **Only PDF uploads to Gradescope will be accepted.** You are encouraged use  $\text{\LaTeX}$  or Word to typeset your solution. You may also scan a neatly handwritten solution to produce the PDF.
  - **Replicate all your code in an appendix.** Begin code for each coding question in a fresh page. Do not put code from multiple questions in the same page. When you upload this PDF on Gradescope, *make sure* that you assign the relevant pages of your code from appendix to correct questions.
  - While collaboration is encouraged, *everything* in your solution must be your (and only your) creation. Copying the answers or code of another student is strictly forbidden.

Furthermore, all external material (i.e., *anything* outside lectures and assigned readings, including figures and pictures) should be cited properly. We wish to remind you that consequences of academic misconduct are *particularly severe*!

3. Code: Submit your code as a .zip file to “Homework 4 Code”.

- **Set a seed for all pseudo-random numbers generated in your code.** This ensures your results are replicated when readers run your code. For example, you can seed numpy with `np.random.seed(189)`.
- Include a README with your name, student ID, the values of random seed (above) you used, and instructions for running (and compiling, if appropriate) your code.
- Do NOT provide any data files. Supply instructions on how to add data to your code.
- Code requiring exorbitant memory or execution time might not be considered.
- Code submitted here must match that in the PDF Write-up. The Kaggle score will not be accepted if the code provided a) does not compile or b) compiles but does not produce the file submitted to Kaggle.

**Notation:** In this assignment we use the following conventions.

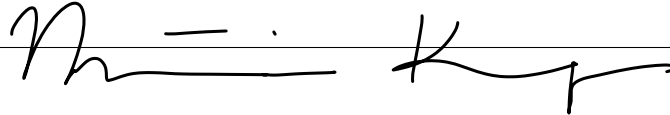
- Symbol “defined equal to” ( $\triangleq$ ) *defines* the quantity to its left to be the expression to its right.
- Scalars are lowercase non-bold:  $x, u_1, \alpha_i$ . Matrices are uppercase alphabets:  $A, B_1, C_i$ . Vectors (column vectors) are in bold:  $\mathbf{x}, \alpha_1, \mathbf{X}, \mathbf{Y}_j$ .
- $\|\mathbf{v}\|$  denotes the Euclidean norm (length) of vector  $\mathbf{v}$ :  $\|\mathbf{v}\| \triangleq \sqrt{\mathbf{v} \cdot \mathbf{v}}$ .  $\|A\|$  denotes the (operator) norm of matrix  $A$ , the magnitude of its largest singular value:  $\|A\| = \max_{\|\mathbf{v}\|=1} \|A\mathbf{v}\|$ .
- $[n] \triangleq \{1, 2, 3, \dots, n\}$ .  $\mathbf{1}$  and  $\mathbf{0}$  denote the vectors with all-ones and all-zeros, respectively.

# 1 Honor Code

Declare and sign the following statement (Mac Preview, PDF Expert, and FoxIt PDF Reader, among others, have tools to let you sign a PDF file):

*"I certify that all solutions are entirely my own and that I have not looked at anyone else's solution. I have given credit to all external sources I consulted."*

Signature: \_\_\_\_\_

A handwritten signature in black ink, appearing to be 'N. K.', written over a horizontal line.

I certify that all solutions are entirely my own and that I have not looked at anyone else's solution. I have given credit to all external sources I consult.

## 2 Logistic Regression with Newton's Method

Given examples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$  and associated labels  $y_1, y_2, \dots, y_n \in \{0, 1\}$ , the cost function for *unregularized* logistic regression is

$$J(\mathbf{w}) \triangleq - \sum_{i=1}^n \left( y_i \ln s_i + (1 - y_i) \ln(1 - s_i) \right)$$

where  $s_i \triangleq s(\mathbf{x}_i \cdot \mathbf{w})$ ,  $\mathbf{w} \in \mathbb{R}^d$  is a weight vector, and  $s(\gamma) \triangleq 1/(1 + e^{-\gamma})$  is the logistic function.

Define the  $n \times d$  design matrix  $X$  (whose  $i^{\text{th}}$  row is  $\mathbf{x}_i^\top$ ), the label  $n$ -vector  $\mathbf{y} \triangleq [y_1 \dots y_n]^\top$ , and  $\mathbf{s} \triangleq [s_1 \dots s_n]^\top$ . For an  $n$ -vector  $\mathbf{a}$ , let  $\ln \mathbf{a} \triangleq [\ln a_1 \dots \ln a_n]^\top$ . The cost function can be rewritten in vector form as

$$J(\mathbf{w}) = -\mathbf{y} \cdot \ln \mathbf{s} - (\mathbf{1} - \mathbf{y}) \cdot \ln(\mathbf{1} - \mathbf{s}).$$

Further, recall that for a real symmetric matrix  $A \in \mathbb{R}^{d \times d}$ , there exist  $U$  and  $\Lambda$  such that  $A = U\Lambda U^\top$  is the eigendecomposition of  $A$ . Here  $\Lambda$  is a diagonal matrix with entries  $\{\lambda_1, \dots, \lambda_d\}$ . An alternative notation is  $\Lambda = \text{diag}(\lambda_i)$ , where  $\text{diag}()$  takes as input the list of diagonal entries, and constructs the corresponding diagonal matrix. This notation is widely used in libraries like `numpy`, and is useful for simplifying some of the expressions when written in matrix-vector form. For example, we can write  $\mathbf{s} = \text{diag}(s_i) \mathbf{1}$ .

*Hint: Recall matrix calculus identities.* The elements in **bold** indicate vectors.

$$\begin{aligned} \nabla_{\mathbf{x}} \alpha \mathbf{y} &= (\nabla_{\mathbf{x}} \alpha) \mathbf{y}^\top + \alpha \nabla_{\mathbf{x}} \mathbf{y} & \nabla_{\mathbf{x}} (\mathbf{y} \cdot \mathbf{z}) &= (\nabla_{\mathbf{x}} \mathbf{y}) \mathbf{z} + (\nabla_{\mathbf{x}} \mathbf{z}) \mathbf{y}; \\ \nabla_{\mathbf{x}} \mathbf{f}(\mathbf{y}) &= (\nabla_{\mathbf{x}} \mathbf{y}) (\nabla_{\mathbf{y}} \mathbf{f}(\mathbf{y})); & \nabla_{\mathbf{x}} g(\mathbf{y}) &= (\nabla_{\mathbf{x}} \mathbf{y}) (\nabla_{\mathbf{y}} g(\mathbf{y})); \end{aligned}$$

and  $\nabla_{\mathbf{x}} C \mathbf{y}(\mathbf{x}) = (\nabla_{\mathbf{x}} \mathbf{y}(\mathbf{x})) C^\top$ , where  $C$  is a constant matrix.

- 1 Derive the gradient  $\nabla_{\mathbf{w}} J(\mathbf{w})$  of cost  $J(\mathbf{w})$  as a matrix-vector expression. Also derive *all intermediate derivatives* in matrix-vector form. Do NOT specify them (**including the intermediates**) in terms of their individual components (e.g.  $\mathbf{w}_i$  for vector  $\mathbf{w}$ ).
- 2 Derive the Hessian  $\nabla_{\mathbf{w}}^2 J(\mathbf{w})$  for the cost function  $J(\mathbf{w})$  as a matrix-vector expression.
- 3 Write the matrix-vector update law for one iteration of Newton's method, substituting the gradient and Hessian of  $J(\mathbf{w})$ .
- 4 You are given four examples  $\mathbf{x}_1 = [0.2 \ 3.1]^\top$ ,  $\mathbf{x}_2 = [1.0 \ 3.0]^\top$ ,  $\mathbf{x}_3 = [-0.2 \ 1.2]^\top$ ,  $\mathbf{x}_4 = [1.0 \ 1.1]^\top$  with labels  $y_1 = 1, y_2 = 1, y_3 = 0, y_4 = 0$ . These points cannot be separated by a line passing through origin. Hence, as described in lecture, append a 1 to each  $\mathbf{x}_{i \in [4]}$  and use a weight vector  $\mathbf{w} \in \mathbb{R}^3$  whose last component is the bias term (called  $\alpha$  in lecture). Begin with initial weight  $\mathbf{w}^{(0)} = [-1 \ 1 \ 0]^\top$ . For the following, state only the final answer with four digits after the decimal point. You may use a calculator or write a program to solve for these, but do NOT submit any code for this part.

- (a) State the value of  $\mathbf{s}^{(0)}$  (the initial value of  $\mathbf{s}$ ).
- (b) State the value of  $\mathbf{w}^{(1)}$  (the value of  $\mathbf{w}$  after 1 iteration).
- (c) State the value of  $\mathbf{s}^{(1)}$  (the value of  $\mathbf{s}$  after 1 iteration).
- (d) State the value of  $\mathbf{w}^{(2)}$  (the value of  $\mathbf{w}$  after 2 iterations).

# Q1: Part 1

$$\begin{aligned}\nabla_w J(w) &= \frac{\partial}{\partial w} -y \cdot \ln s - (1-y) \cdot \ln(1-s) \\ &= -y \frac{\partial}{\partial w} [\ln s(xw)] - (1-y) \left[ \frac{\partial}{\partial w} \ln(1-s) \right]\end{aligned}$$

$$s_i = s(x_i \cdot w) = \frac{1}{1 + e^{x_i w}}$$

$$s'(\gamma) = \frac{d}{d\gamma} \frac{1}{1 + e^\gamma} = \frac{e^{-\gamma}}{(1 + e^{-\gamma})^2} = s(\gamma)(1 - s(\gamma)) \text{ by hint ④}$$

$$= -y \left[ \frac{x^T (\cancel{s(xw)})(1 - \cancel{s(xw)})}{\cancel{s(xw)}} \right] - (1-y) \left[ \frac{x^T (s(xw)(\cancel{1-s(xw)}))}{\cancel{1-s(xw)}} \right]$$

$$= -y [x^T (1 - s(xw))] - (1-y) [x^T (s(xw))]$$

$$= -y [x^T - x^T s(xw)] + (1-y) [x^T s(xw)]$$

$$= -y x^T + \cancel{y x^T s(xw)} + x^T s(xw) - \cancel{y x^T s(xw)}$$

$$= -y x^T + x^T s(xw)$$

$$= -x^T (y - s(xw))$$

$$\Rightarrow -x^T (y - s)$$

# Q1: Part 2

$$\nabla_{\omega}^2 J(\omega) = \nabla_{\omega} (\nabla_{\omega} J(\omega))$$

$$= \nabla_{\omega} (-X^T (y - s(X\omega)))$$

$$= \nabla_{\omega} (-X^T y + X^T s(X\omega))$$

$$= \nabla_{\omega} X^T s(X\omega)$$

$$= \nabla_{\omega} \sum X_i^T s_i$$

$$= \nabla_{\omega} \sum X_i^T s(X_i\omega)$$

$$= \nabla_{\omega} \sum X_i^T \left( \frac{1}{1 + e^{-X_i\omega}} \right)$$

$$= \sum X_i^T \frac{X_i e^{-X_i\omega}}{(1 + e^{-X_i\omega})^2}$$

$$= \sum X_i^T X_i \frac{e^{-X_i\omega}}{(1 + e^{-X_i\omega})^2}$$

$$= \sum X_i^T X_i \underbrace{\frac{1}{(1 + e^{-X_i\omega})}}_{s_i} \underbrace{\frac{e^{-X_i\omega}}{(1 + e^{-X_i\omega})}}_{1 - s_i}$$

$$= \sum X_i^T X_i (s_i)(1 - s_i)$$

$$= X^T \Omega X \quad \text{where} \quad \Omega = \begin{bmatrix} s_1(1-s_1) & 0 & \dots & 0 \\ 0 & s_2(1-s_2) & & \\ \vdots & & \ddots & \\ 0 & & & s_n(1-s_n) \end{bmatrix}$$

- 3 Write the matrix-vector update law for one iteration of Newton's method, substituting the gradient and Hessian of  $J(\mathbf{w})$ .

$$(X^T \Omega X) e = X^T (y - s)$$

$$e = (X^T \Omega X)^{-1} X^T (y - s)$$

$$\omega_{n+1} \leftarrow \omega_n + e$$

$$\leftarrow \omega_n + (X^T \Omega X)^{-1} X^T (y - s)$$

$$\Rightarrow \omega_{n+1} \leftarrow \omega_n + (X^T \Omega X)^{-1} X^T (y - s)$$



4 You are given four examples  $\mathbf{x}_1 = [0.2 \ 3.1]^\top$ ,  $\mathbf{x}_2 = [1.0 \ 3.0]^\top$ ,  $\mathbf{x}_3 = [-0.2 \ 1.2]^\top$ ,  $\mathbf{x}_4 = [1.0 \ 1.1]^\top$  with labels  $y_1 = 1, y_2 = 1, y_3 = 0, y_4 = 0$ . These points cannot be separated by a line passing through origin. Hence, as described in lecture, append a 1 to each  $\mathbf{x}_{i \in [4]}$  and use a weight vector  $\mathbf{w} \in \mathbb{R}^3$  whose last component is the bias term (called  $\alpha$  in lecture). Begin with initial weight  $\mathbf{w}^{(0)} = [-1 \ 1 \ 0]^\top$ . For the following, state only the final answer with four digits after the decimal point. You may use a calculator or write a program to solve for these, but do NOT submit any code for this part.

- State the value of  $\mathbf{s}^{(0)}$  (the initial value of  $\mathbf{s}$ ).
- State the value of  $\mathbf{w}^{(1)}$  (the value of  $\mathbf{w}$  after 1 iteration).
- State the value of  $\mathbf{s}^{(1)}$  (the value of  $\mathbf{s}$  after 1 iteration).
- State the value of  $\mathbf{w}^{(2)}$  (the value of  $\mathbf{w}$  after 2 iterations).

a.)

$$\left. \begin{aligned} \mathbf{x}_1 &= [0.2 \ 3.1 \ 1]^\top \\ \mathbf{x}_2 &= [1.0 \ 3.0 \ 1]^\top \\ \mathbf{x}_3 &= [-0.2 \ 1.2 \ 1]^\top \\ \mathbf{x}_4 &= [1.0 \ 1.1 \ 1]^\top \end{aligned} \right\} \begin{bmatrix} 0.2 & 3.1 & 1 \\ 1.0 & 3.0 & 1 \\ -0.2 & 1.2 & 1 \\ 1.0 & 1.1 & 1 \end{bmatrix} = \mathbf{X}$$

$$\mathbf{X} \mathbf{w}^{(0)} = \mathbf{s}^{(0)}$$

$$\mathbf{X} \mathbf{w}^{(0)} = \begin{bmatrix} 0.2 & 3.1 & 1 \\ 1.0 & 3.0 & 1 \\ -0.2 & 1.2 & 1 \\ 1.0 & 1.1 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2.9 \\ 2 \\ 1.4 \\ 0.1 \end{bmatrix}$$

$$= \begin{bmatrix} s(2.9) \\ s(2.0) \\ s(1.4) \\ s(0.1) \end{bmatrix} = \begin{bmatrix} 0.9478 \\ 0.8808 \\ 0.8022 \\ 0.5250 \end{bmatrix} = \mathbf{s}^{(0)}$$

b.)

$$\omega_1 \leftarrow \omega_0 + (X^T \Omega X)^{-1} X^T (y - s)$$

$$\omega_1 = \text{array}([1.32465198, 3.04991697, -6.82910388])$$

$$\omega^{(1)} = \begin{bmatrix} 1.3247 \\ 3.0499 \\ -6.8291 \end{bmatrix}$$

c.)

$$X \omega^{(1)} = s^{(1)} = \text{array}([0.94737826, 0.97455097, 0.03124556, 0.10437391])$$

$$\Rightarrow s^{(1)} = \begin{bmatrix} 0.9474 \\ 0.9746 \\ 0.0312 \\ 0.1044 \end{bmatrix}$$

d.)

$$\omega^{(2)} = \text{array}([1.06902047, 3.52252931, -7.56418584])$$

$$\Rightarrow \omega^{(2)} = \begin{bmatrix} 1.0690 \\ 3.5225 \\ -7.5642 \end{bmatrix}$$

### 3 Wine Classification with Logistic Regression

The wine dataset `data.mat` consists of 6,497 sample points, each having 12 features. The description of these features is provided in `data.mat`. The dataset includes a training set of 6,000 sample points and a test set of 497 sample points. Your classifier needs to predict whether a wine is white (class label 0) or red (class label 1).

Begin by normalizing the data with each feature's mean and standard deviation. You should use training data statistics to normalize both training and validation/test data. Then add a fictitious dimension. Whenever required, it is recommended that you tune hyperparameter values with cross-validation.

Please set a random seed whenever needed and **report it**.

**Use of automatic logistic regression libraries/packages is prohibited for this question.** If you are coding in python, it is better to use `scipy.special.expit` for evaluating logistic functions as its code is numerically stable, and doesn't produce NaN or MathOverflow exceptions.

- 1 *Batch Gradient Descent Update.* State the batch gradient descent update law for logistic regression **with  $\ell_2$  regularization**. As this is a “batch” algorithm, each iteration should use *every training example*. You don't have to show your derivation. You may reuse results from your solution to question 4.1.
- 2 *Batch Gradient Descent Code.* Implement your batch gradient descent algorithm for logistic regression and include your code here. Choose reasonable values for the regularization parameter and step size (learning rate), specify your chosen values in the write-up, and train your model from question 3.1. Shuffle and split your data into training/validation sets and mention the random seed used in the write-up. Plot the value of the cost function versus the number of iterations spent in training.
- 3 *Stochastic Gradient Descent (SGD) Update.* State the SGD update law for logistic regression with  $\ell_2$  regularization. Since this is not a “batch” algorithm anymore, each iteration uses *just one* training example. You don't have to show your derivation.
- 4 *Stochastic Gradient Descent Code.* Implement your stochastic gradient descent algorithm for logistic regression and include your code here. Choose a suitable value for the step size (learning rate), specify your chosen value in the write-up, and run your SGD algorithm from question 3.3. Shuffle and split your data into training/validation sets and mention the random seed used in the write-up. Plot the value of the cost function versus the number of iterations spent in training.  
  
Compare your plot here with that of question 3.2. Which method converges more quickly? Briefly describe what you observe.
- 5 Instead of using a constant step size (learning rate) in SGD, you could use a step size that slowly shrinks from iteration to iteration. Run your SGD algorithm from question 3.3 with a step size  $\epsilon_t = \delta/t$  where  $t$  is the iteration number and  $\delta$  is a hyperparameter you select

empirically. Mention the value of  $\delta$  chosen. Plot the value of cost function versus the number of iterations spent in training.

How does this compare to the convergence of your previous SGD code?

- 6 *Kaggle*. Train your *best* classifier on the entire training set and submit your prediction on the test sample points to Kaggle. As always for Kaggle competitions, you are welcome to add or remove features, tweak the algorithm, and do pretty much anything you want to improve your Kaggle leaderboard performance **except** that you may not replace or augment logistic regression with a wholly different learning algorithm. Your code should output the predicted labels in a CSV file.

Report your Kaggle username and your best score, and briefly describe what your best classifier does to achieve that score.

### Q3: Part 1

The update rule will only differ by adding the  $l_2$  term to the cost function

$$J(w) = -y \cdot \ln s - (1-y) \cdot \ln(1-s) + \underbrace{\lambda \|w\|^2}_{l_2 \text{ regularizer}}$$

$$\begin{aligned}\nabla_w J(w) &= \frac{\partial}{\partial w} [-y \cdot \ln(s) - (1-y) \ln(1-s)] + \frac{\partial}{\partial w} \lambda \|w\|^2 \\ &= \underbrace{-X^T(y-s)}_{\text{part 2.1}} + 2\lambda w\end{aligned}$$

Thus, we see for the update rule:

$$w_{i+1} := w_i - \alpha (\nabla_w J(w))$$

$$\rightarrow w_{i+1} := w_i - \alpha (-X^T(y-s) + 2\lambda w)$$

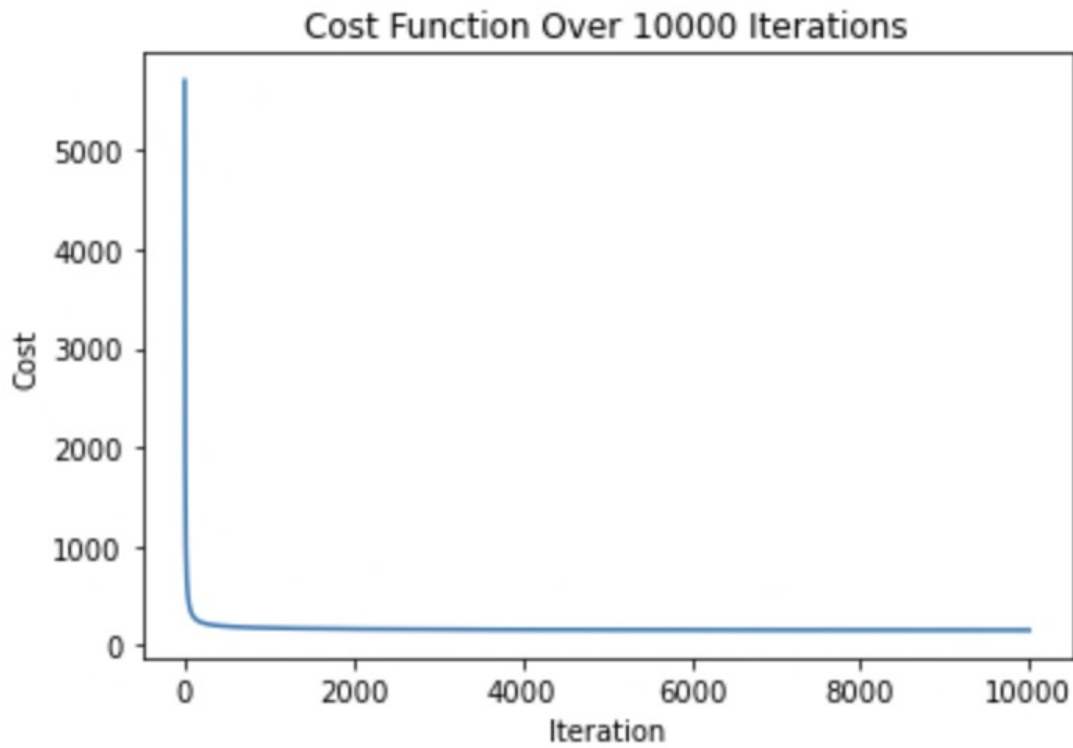
update  
rule

### Q3: Part 2

$\alpha$ : 0.0001

$\lambda$ : 0.1

Batch Gradient Descent



`Random.seed(10)`

### Q3 : Part 3

Thus, we see for the update rule for stochastic GD:

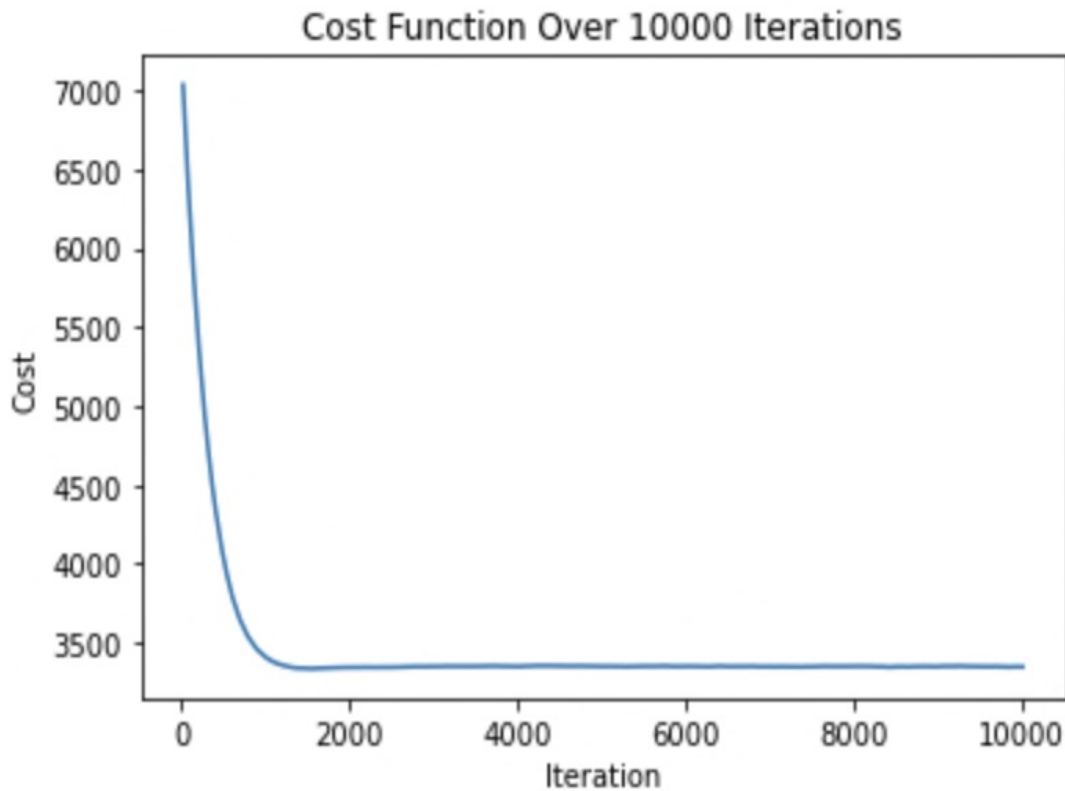
$$\begin{aligned}\omega_{n+1} &:= \omega_n - \alpha \nabla_w (J(w)) \\ &= \omega_n - \alpha (-X_j^T (y_j - s_j) + 2\lambda w) \\ &= \omega_n - \alpha (-X_j^T (y_j - s(x_j, w)) + 2\lambda w)\end{aligned}$$

$$\Rightarrow \omega_n = \alpha (-X_j^T (y_j - s(x_j, w)) + 2\lambda w)$$

### Q3: Part 4

$\alpha : 0.0001$   
 $\lambda : 10$

Stochastic Gradient  
Descent



Random.seed(10)

Batch Gradient Descent converges more quickly compared to Stochastic Descent. The cost is also higher overall for Stochastic than Batch.

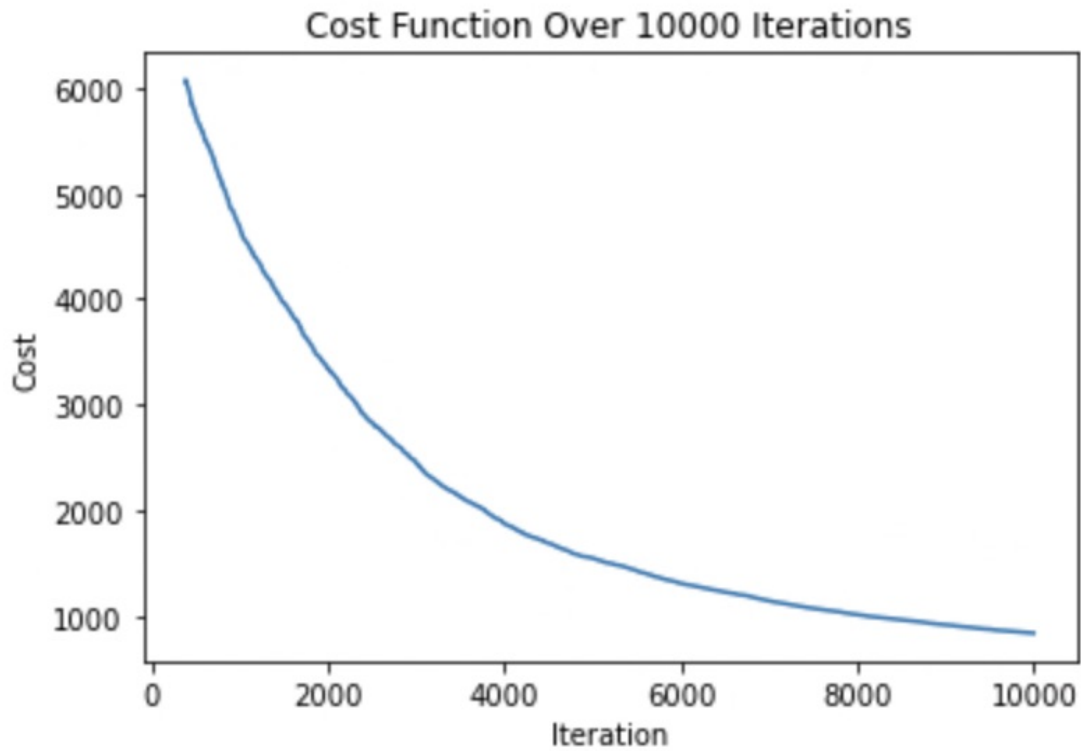


### Q3: Part 5

$\alpha$  : 0.001

$\lambda$  : 0.1

Decreasing Step  
Size



`Random.seed(10)`

## Q5: Part 6

Kaggle username: nadaleek

Wine Score: 0.95564

lambda: 0.1

alpha: 0.0001

random.seed(10)

I used batch gradient descent for my classifier. It produces the most accuract results compared to stochastic gradient descent. I calculated the cost function and iterated though 50000 times to train the data.

```
test = wineData["X_test"]
lam = 0.1
alpha = .0001
w = np.ones((test.shape[1], 1))
iterations = 50000

X = wineData["X"]
y = wineData["y"]

w_final, cost = gradient_descent(alpha, lam, X, y, iterations)

prediction = (sigmoid(test @ w_final) >= 0.5).astype(int)
prediction = prediction.reshape(-1,)
prediction
```

## 4 A Bayesian Interpretation of Lasso

Suppose you are aware that the labels  $y_{i \in [n]}$  corresponding to sample points  $\mathbf{x}_{i \in [n]} \in \mathbb{R}^d$  follow the density law

$$f(y_i | \mathbf{x}_i, \mathbf{w}) = \frac{1}{\sigma \sqrt{2\pi}} e^{-(y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 / (2\sigma^2)}$$

where  $\sigma > 0$  is a known constant and  $\mathbf{w} \in \mathbb{R}^d$  is a random parameter. Suppose further that experts have told you that

- each component of  $\mathbf{w}$  is independent of the others, and
- each component of  $\mathbf{w}$  has the Laplace distribution with location 0 and scale being a known constant  $b$ . That is, each component  $w_i$  obeys the density law  $f(w_i) = e^{-|w_i|/b} / (2b)$ .

Assume the outputs  $y_{i \in [n]}$  are independent from each other.

Your goal is to find the choice of parameter  $\mathbf{w}$  that is *most likely* given the input-output examples  $(\mathbf{x}_i, y_i)_{i \in [n]}$ . This method of estimating parameters is called *maximum a posteriori* (MAP); Latin for “*maximum [odds] from what follows*.”

1. Derive the *posterior* probability density law  $f(\mathbf{w} | (\mathbf{x}_i, y_i)_{i \in [n]})$  for  $\mathbf{w}$  up to a *proportionality constant* by applying Bayes’ Theorem and substituting for the densities  $f(y_i | \mathbf{x}_i, \mathbf{w})$  and  $f(\mathbf{w})$ . Don’t try to derive an exact expression for  $f(\mathbf{w} | (\mathbf{x}_i, y_i)_{i \in [n]})$ , as the denominator is very involved and irrelevant to maximum likelihood estimation.
2. Define the log-likelihood for MAP as  $\ell(\mathbf{w}) \triangleq \ln f(\mathbf{w} | \mathbf{x}_{i \in [n]}, y_{i \in [n]})$ . Show that maximizing the MAP log-likelihood over all choices of  $\mathbf{w}$  is the same as minimizing  $\sum_{i=1}^n (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1$  where  $\|\mathbf{w}\|_1 = \sum_{j=1}^d |w_j|$  and  $\lambda$  is a constant. Also give a formula for  $\lambda$  as a function of the distribution parameters.

1. Derive the *posterior* probability density law  $f(\mathbf{w} | (\mathbf{x}_i, y_i)_{i \in [n]})$  for  $\mathbf{w}$  up to a *proportionality constant* by applying Bayes' Theorem and substituting for the densities  $f(y_i | \mathbf{x}_i, \mathbf{w})$  and  $f(\mathbf{w})$ . Don't try to derive an exact expression for  $f(\mathbf{w} | (\mathbf{x}_i, y_i)_{i \in [n]})$ , as the denominator is very involved and irrelevant to maximum likelihood estimation.

$y_i$  constant

By Bayes Rule,

$$f(\mathbf{w} | (\mathbf{x}_i, y_i)_{i \in [n]}) = \frac{f(y_i | \mathbf{x}_i, \mathbf{w}) f(\mathbf{w} | \mathbf{x}_i)}{f(y_i | \mathbf{x}_i)}$$

$f(\mathbf{w} | \mathbf{x}_i) = f(\mathbf{w})$

$$= \frac{f(y_i | \mathbf{x}_i, \mathbf{w}) f(\mathbf{w})}{f(y_i | \mathbf{x}_i)}$$

$$f(w_i) = \frac{e^{-|w_i|/b}}{2b} \longrightarrow f(\mathbf{w}) = \frac{e^{-\sum_{k=1}^d |w_k|/b}}{(2b)^d}$$

$$f(y_i | \mathbf{x}_i, \mathbf{w}) = \frac{1}{\sigma \sqrt{2\pi}} e^{-(y_i - \mathbf{w} \mathbf{x}_i)^2 / 2\sigma^2}$$

$$\Rightarrow \frac{e^{-(y_i - \mathbf{w} \mathbf{x}_i)^2 / 2\sigma^2}}{\sigma \sqrt{2\pi}} \cdot \frac{e^{-\sum_{k=1}^d |w_k|/b}}{(2b)^d}$$

$f(y_i | \mathbf{x}_i)$  ← constant since all  $y_i$  are independent

$$\Rightarrow \frac{e^{-(y_i - \mathbf{w} \mathbf{x}_i)^2 / 2\sigma^2}}{\sigma \sqrt{2\pi}} \cdot \frac{e^{-\sum_{k=1}^d |w_k|/b}}{(2b)^d}$$

2. Define the log-likelihood for MAP as  $\ell(\mathbf{w}) \triangleq \ln f(\mathbf{w} | \mathbf{x}_{i \in [n]}, y_{i \in [n]})$ . Show that maximizing the MAP log-likelihood over all choices of  $\mathbf{w}$  is the same as minimizing  $\sum_{i=1}^n (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1$  where  $\|\mathbf{w}\|_1 = \sum_{j=1}^d |w_j|$  and  $\lambda$  is a constant. Also give a formula for  $\lambda$  as a function of the distribution parameters.

$$\begin{aligned}
 \max_{\mathbf{w}} \ell(\mathbf{w}) &= \max_{\mathbf{w}} \ln f(\mathbf{w} | \mathbf{x}_{i \in [n]}, y_{i \in [n]}) \\
 &= \max_{\mathbf{w}} \ln \left[ \frac{f(y_i | x_i, \mathbf{w}) f(\mathbf{w})}{f(y_i | x_i)} \right] \\
 &= \max_{\mathbf{w}} \ln(f(y_i | x_i, \mathbf{w})) + \ln(f(\mathbf{w})) - \ln(\cancel{f(y_i | x_i)}) \quad \text{not depend. on } \mathbf{w} \\
 &= \max_{\mathbf{w}} \ln \left( \frac{1}{\sigma \sqrt{2\pi}} e^{-(y_i - \mathbf{w} x_i)^2 / 2\sigma^2} \right) + \ln \left( \frac{e^{-\sum_{k=1}^d |w_k|/b}}{(2b)^d} \right) \quad \leftarrow = \|\mathbf{w}\|_1 \\
 &= \max_{\mathbf{w}} \ln(\cancel{\frac{1}{\sigma \sqrt{2\pi}}}) + \ln(e^{-(y_i - \mathbf{w} x_i)^2 / 2\sigma^2}) + \ln(\cancel{\frac{1}{(2b)^d}}) + \ln(e^{-\|\mathbf{w}\|_1 / b}) \\
 &= \max_{\mathbf{w}} -(y_i - \mathbf{w} x_i)^2 / 2\sigma^2 + (-\|\mathbf{w}\|_1 / b) \\
 &= \max_{\mathbf{w}} \frac{-(y_i - \mathbf{w} x_i)^2}{2\sigma^2} - \frac{\|\mathbf{w}\|_1}{b} \\
 &= \max_{\mathbf{w}} -(y_i - \mathbf{w} x_i)^2 - \frac{2\sigma^2 \|\mathbf{w}\|_1}{b} \\
 &= \max_{\mathbf{w}} -(y_i - \mathbf{w} x_i)^2 - \lambda \|\mathbf{w}\|_1 \\
 &= \min_{\mathbf{w}} (y_i - \mathbf{w} x_i)^2 + \lambda \|\mathbf{w}\|_1 \quad \text{where } \lambda = \frac{2\sigma^2}{b}
 \end{aligned}$$

## 5 $\ell_1$ -regularization, $\ell_2$ -regularization, and Sparsity

You are given a design matrix  $X$  (whose  $i^{\text{th}}$  row is sample point  $\mathbf{x}_i^\top$ ) and an  $n$ -vector of labels  $\mathbf{y} \triangleq [y_1 \dots y_n]^\top$ . For simplicity, assume  $X$  is whitened, so  $X^\top X = nI$ . Do not add a fictitious dimension/bias term; for input  $\mathbf{0}$ , the output is always 0. Let  $\mathbf{x}_{*i}$  denote the  $i^{\text{th}}$  column of  $X$ .

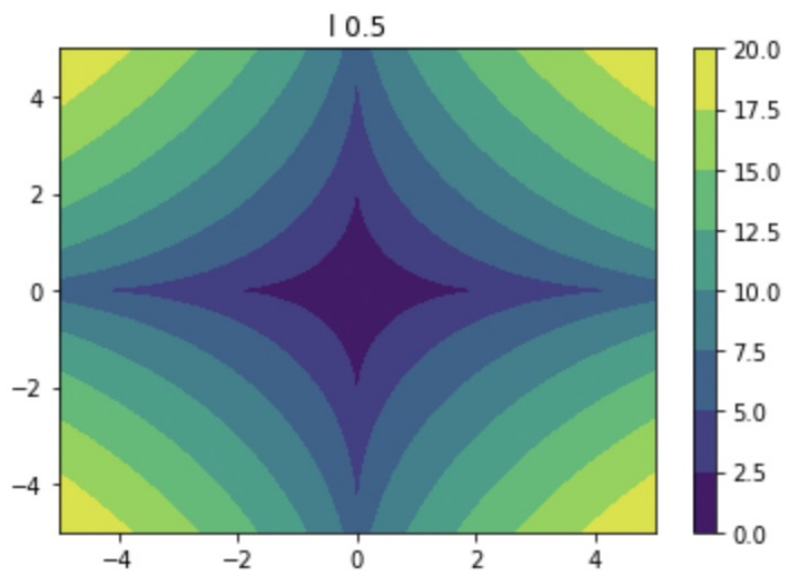
1. The  $\ell_p$ -norm for  $w \in \mathbb{R}^d$  is defined as  $\|w\|_p = (\sum_{i=1}^d |w_i|^p)^{1/p}$ , where  $p > 0$ . Plot the isocontours with  $w \in \mathbb{R}^2$ , for the following norms.  
(a)  $\ell_{0.5}$     (b)  $\ell_1$     (c)  $\ell_2$

**Use of automatic libraries/packages for computing norms is prohibited for the question.**

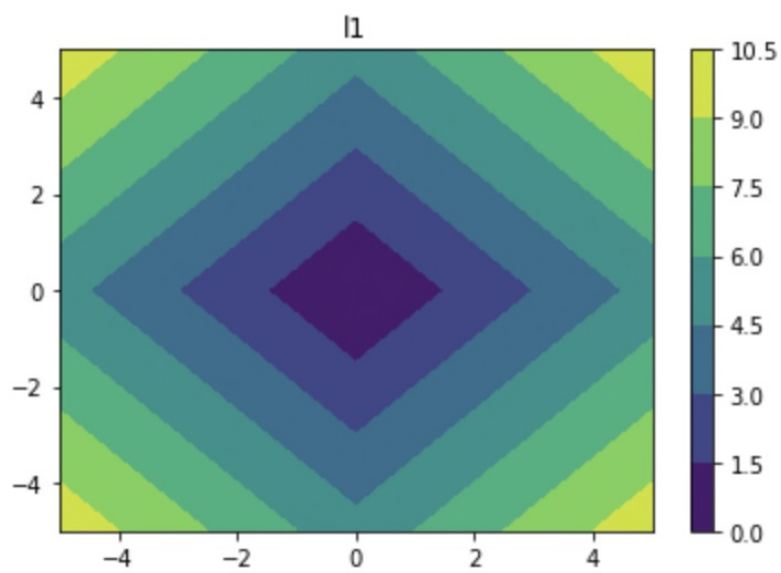
2. Show that the cost function for  $\ell_1$ -regularized least squares,  $J_1(\mathbf{w}) \triangleq \|X\mathbf{w} - \mathbf{y}\|^2 + \lambda\|\mathbf{w}\|_1$  (where  $\lambda > 0$ ), can be rewritten as  $J_1(\mathbf{w}) = \|\mathbf{y}\|^2 + \sum_{i=1}^d f(\mathbf{x}_{*i}, \mathbf{w}_i)$  where  $f(\cdot, \cdot)$  is a suitable function whose first argument is a vector and second argument is a scalar.
3. Using your solution to part [2](#) derive necessary and sufficient conditions for the  $i^{\text{th}}$  component of the optimizer  $\mathbf{w}^*$  of  $J_1(\cdot)$  to satisfy each of these three properties:  $w_i^* > 0$ ,  $w_i^* = 0$ , and  $w_i^* < 0$ .
4. For the optimizer  $\mathbf{w}^\#$  of the  $\ell_2$ -regularized least squares cost function  $J_2(\mathbf{w}) \triangleq \|X\mathbf{w} - \mathbf{y}\|^2 + \lambda\|\mathbf{w}\|^2$  (where  $\lambda > 0$ ), derive a necessary and sufficient condition for  $\mathbf{w}_i^\# = 0$ , where  $\mathbf{w}_i^\#$  is the  $i$ th component of  $\mathbf{w}^\#$ .
5. A vector is called *sparse* if most of its components are 0. From your solution to part [3](#) and [4](#) which of  $\mathbf{w}^*$  and  $\mathbf{w}^\#$  is more likely to be sparse? Why?

Q 5 Part 1

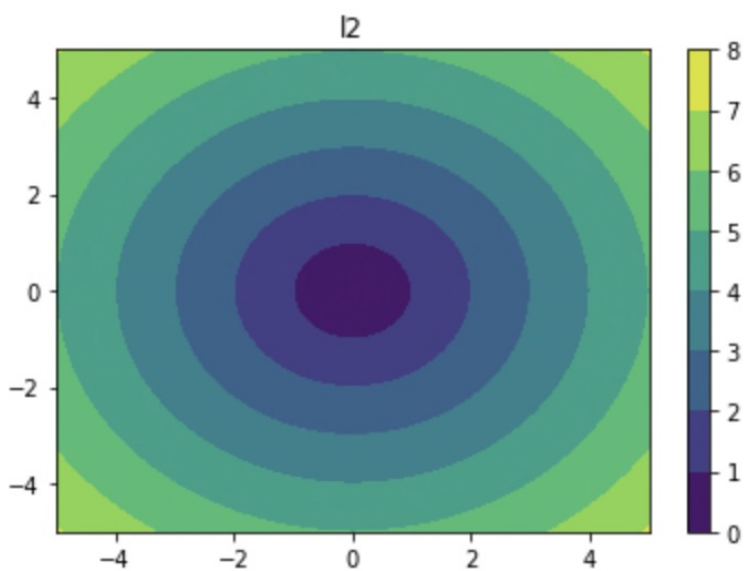
$\|w\|_{0.5} \longrightarrow$



$\|w\|_1 \longrightarrow$



$\|w\|_2 \longrightarrow$



Q5  
Part 2

2. Show that the cost function for  $\ell_1$ -regularized least squares,  $J_1(\mathbf{w}) \triangleq \|X\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1$  (where  $\lambda > 0$ ), can be rewritten as  $J_1(\mathbf{w}) = \|\mathbf{y}\|^2 + \sum_{i=1}^d f(\mathbf{x}_i, \mathbf{w}_i)$  where  $f(\cdot, \cdot)$  is a suitable function whose first argument is a vector and second argument is a scalar.

$$f(\text{vec}, \text{scal})$$

$$\begin{aligned} J_1(\mathbf{w}) &= \|X\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1 \\ &= (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) + \lambda \sum_{i=1}^d w_i \\ &= (\mathbf{w}^T X^T - \mathbf{y}^T) (X\mathbf{w} - \mathbf{y}) + \lambda \sum_{i=1}^d w_i \\ &= \underbrace{\mathbf{w}^T X^T X \mathbf{w}}_{n\mathbf{I}} - \mathbf{w}^T X^T \mathbf{y} - \mathbf{y}^T X \mathbf{w} + \mathbf{y}^T \mathbf{y} + \lambda \sum_{i=1}^d w_i \\ &= n \|\mathbf{w}\|^2 - 2 \mathbf{w}^T X^T \mathbf{y} + \|\mathbf{y}\|^2 + \lambda \sum_{i=1}^d w_i \\ &= \|\mathbf{y}\|^2 + n \sum_{i=1}^d w_i^2 + \lambda \sum_{i=1}^d w_i - 2 \sum_{i=1}^d w_i x_i y_i \\ &= \|\mathbf{y}\|^2 + \sum_{i=1}^d n w_i^2 + \lambda w_i - 2 w_i x_i y_i \end{aligned}$$

$$\text{Let } f(\mathbf{x}, w) = n w^2 + \lambda w - 2 w (\mathbf{y}^T \mathbf{x})$$

$$\Rightarrow \|\mathbf{y}\|^2 + f(\mathbf{x}, w)$$



3. Using your solution to part 2, derive necessary and sufficient conditions for the  $i^{\text{th}}$  component of the optimizer  $\mathbf{w}^*$  of  $J_1(\cdot)$  to satisfy each of these three properties:  $w_i^* > 0$ ,  $w_i^* = 0$ , and  $w_i^* < 0$ .

$$\boxed{w_i^* > 0}$$

By 5.2,

$$\longrightarrow \nabla_{\mathbf{w}} n\mathbf{w}^2 + \lambda \|\mathbf{w}\| - 2\mathbf{y} \mathbf{x}_i$$

$$= 2nw_i^* + \lambda - 2y x_i$$

Setting this equal to 0 and solving for  $w_i^*$ :

$$2nw_i^* + \lambda = 2y x_i$$

$$\longrightarrow 2nw_i^* = 2y x_i - \lambda$$

$$\longrightarrow w_i^* = \frac{2y x_i - \lambda}{2n}$$

$$\longrightarrow \frac{2y x_i - \lambda}{2n} > 0$$

$$\longrightarrow 2y x_i - \lambda > 0$$

$$\longrightarrow y x_i - \frac{\lambda}{2} > 0$$

$$\text{Thus, } w_i^* > 0 \text{ iff } y \cdot x_i - \frac{\lambda}{2} > 0$$

$$\boxed{w_i^* = 0}$$

By the first part,

$$w_i^* = \frac{2y x_i - \lambda}{2n}$$

$$\rightarrow \frac{2y x_i - \lambda}{2n} = 0$$

$$\rightarrow 2y x_i - \lambda = 0$$

$$\rightarrow 2y x_i = \lambda$$

$$\rightarrow y x_i = \frac{\lambda}{2}$$

$$\text{Thus, } w_i^* = 0 \text{ iff } y \cdot x_i = \frac{\lambda}{2}, \text{ or } y \cdot x_i - \frac{\lambda}{2} = 0$$

$$w_i^* < 0$$

$$\text{Again, } w_i^* = \frac{2y x_i - \lambda}{2n}$$

$$\rightarrow \frac{2y x_i - \lambda}{2n} < 0$$

$$\rightarrow 2y x_i - \lambda < 0$$

$$\rightarrow y x_i - \frac{\lambda}{2} < 0$$

$$\rightarrow y x_i < \frac{\lambda}{2}$$

$$\Rightarrow w_i^* < 0 \text{ iff } y^T x_i - \frac{\lambda}{2} < 0 \text{ or } y \cdot x_i < \frac{\lambda}{2}$$

4. For the optimizer  $\mathbf{w}^\#$  of the  $\ell_2$ -regularized least squares cost function  $J_2(\mathbf{w}) \triangleq \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2$  (where  $\lambda > 0$ ), derive a necessary and sufficient condition for  $\mathbf{w}_i^\# = 0$ , where  $\mathbf{w}_i^\#$  is the  $i$ th component of  $\mathbf{w}^\#$ .

Taking the gradient, setting it to 0, and solving for  $w_i^\#$ ,

$$\begin{aligned} J_2(\mathbf{w}) &= (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \|\mathbf{w}\|^2 = \mathbf{w}^\top \underbrace{\mathbf{X}^\top \mathbf{X}}_{n\mathbf{I}} \mathbf{w} - 2\mathbf{y}^\top \mathbf{X}\mathbf{w} + \mathbf{y}^\top \mathbf{y} + \lambda \|\mathbf{w}\|^2 \\ &= n \|\mathbf{w}\|^2 - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \|\mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2 \quad \leftarrow \text{part 5.2} \end{aligned}$$

Taking gradient & setting to 0 ↴

$$\nabla_{\mathbf{w}} J_2(\mathbf{w}) = 2n\mathbf{w}^\# - 2\mathbf{X}^\top \mathbf{y} + 2\lambda \mathbf{w}^\# = 0$$

$$\rightarrow 2n\mathbf{w}^\# + 2\lambda \mathbf{w}^\# = 2\mathbf{X}^\top \mathbf{y}$$

$$\rightarrow \mathbf{w}^\# (2n + 2\lambda) = 2\mathbf{X}^\top \mathbf{y}$$

$$\rightarrow \mathbf{w}^\# = \frac{2\mathbf{X}^\top \mathbf{y}}{2n + 2\lambda}$$

$$\rightarrow w_i^\# = \frac{\mathbf{y} \cdot \mathbf{x}_i}{n + \lambda}$$

Thus, for  $w_i^\# = 0$ , then  $\mathbf{y} \cdot \mathbf{x}_i = 0$

$$\Rightarrow w_i^\# = 0 \quad \text{iff} \quad \mathbf{y} \cdot \mathbf{x}_i = 0$$

5. A vector is called *sparse* if most of its components are 0. From your solution to part 3 and 4, which of  $\mathbf{w}^*$  and  $\mathbf{w}^\#$  is more likely to be sparse? Why?

$$\text{Part 3 \& 4} \quad \begin{cases} w_i^* = 0 & \text{iff } y \cdot x_i = \frac{\lambda}{2} \\ w_i^\# = 0 & \text{iff } y \cdot x_i = 0 \end{cases}$$

$w_i^*$  is more likely to be sparse. That is because it is more likely for  $y \cdot x_i$  to be equal to some value rather than it being 0. Thus, since it is less likely that  $y$  and  $x_i$  are orthogonal, then  $y \cdot x_i = \frac{\lambda}{2}$  is more likely.

## Submission Checklist

Please ensure you have completed the following before your final submission.

At the beginning of your writeup...

1. Have you copied and hand-signed the honor code specified in Question 1?
2. Have you listed all students (Names and ID numbers) that you collaborated with?

In your writeup for Question 3...

1. Have you included your **Kaggle Score** and **Kaggle Username**?

At the end of the writeup...

1. Have you provided a code appendix including all code you wrote in solving the homework?

## Executable Code Submission

1. Have you created an archive containing all “.py” files that you wrote or modified to generate your homework solutions?
2. Have you removed all data and extraneous files from the archive?
3. Have you included a README file in your archive containing any special instructions to reproduce your results?

## Submissions

1. Have you submitted your written solutions to the Gradescope assignment titled **HW4 Write-Up** and selected pages appropriately?
2. Have you submitted your executable code archive to the Gradescope assignment titled **HW4 Code**?
3. Have you submitted your test set predictions for **Wine** dataset to the appropriate Kaggle challenge?

Congratulations! You have completed Homework 4.