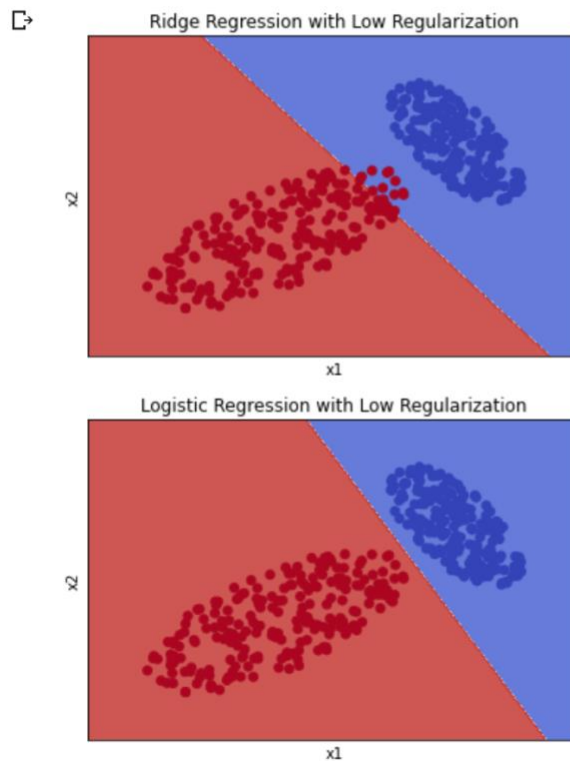


1. Comparing Different Loss Functions

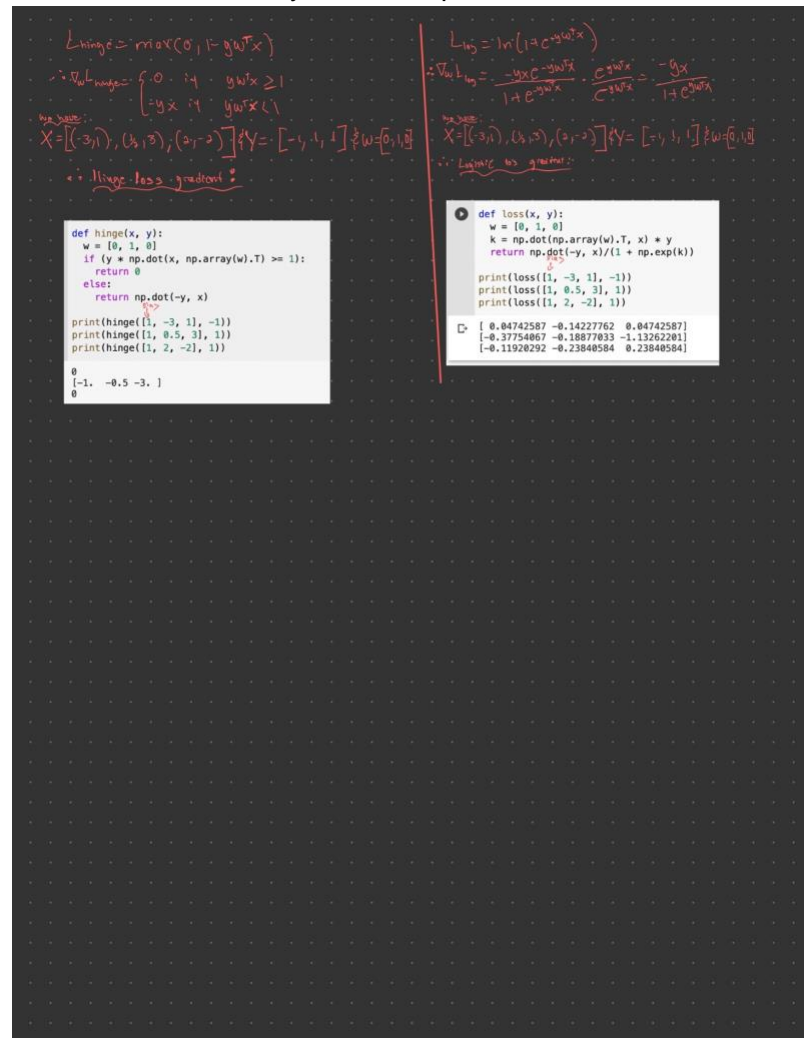
- a. Squared loss, commonly used for regression, is a poor choice for classification because it does not consider the discrete output nature and can attribute high loss to correctly classified points far from the line/hyperplane while attributing low loss to misclassified points close to it, leading to overfitting and suboptimal performance. Cross-entropy (Log) loss, on the other hand, is more appropriate for classification as it minimizes the difference between predicted class probabilities and the true class label and is less prone to overfitting.

```
✓ 0s ▶ ridge = Ridge(alpha = 1e-6)
ridge.fit(X, Y)
make_plot(X, Y, ridge, "Ridge Regression with Low Regularization", "Ridge.png")
logistic = LogisticRegression(C = np.Infinity)
logistic.fit(X, Y)
make_plot(X, Y, logistic, "Logistic Regression with Low Regularization", "Logistic.png")
```



- b. (See Code) The logistic regression boundary line is a better separator of the red and blue data points than the ridge regression boundary line which does not even completely separate the two as we still have misclassified red points. This is a result of squared loss, used by ridge regression, penalizing correctly classified points that are further from the boundary line hence shifting the boundary line towards those points to minimize loss and at the red data points have more points further from the boundary line hence the line is shifted to towards the red data points to minimize loss, however, it misclassified points at

the same time. However, log loss penalizes misclassified values and produces a lower loss for correctly classified points further from the boundary.

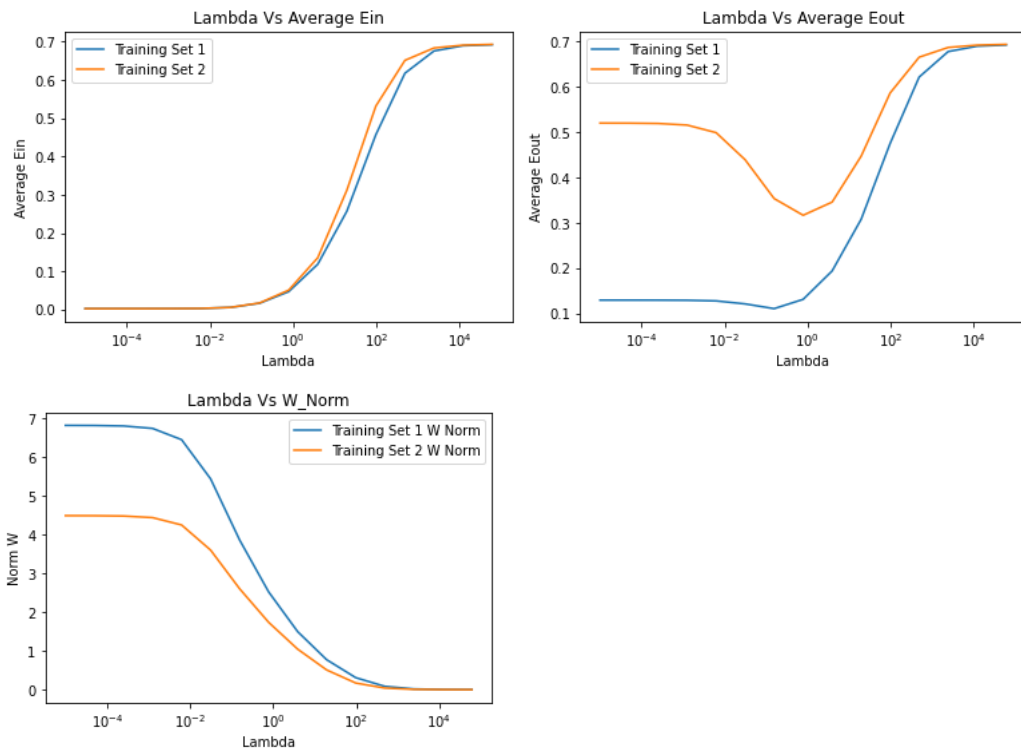


- c. The loss of for each point corresponds to the printed value for that input x in each function call starting after $x[0]$ which is the bias.
- d. The hinge loss gradient converges to zero when all points are correctly classified meaning $y(\text{transpose}(w))x \geq 1$ which results in 0. The log loss gradient converges to zero when $y(\text{transpose}(w))x$ goes to infinity which makes the equation from part c zero. For linearly separable data, training error can be minimized by scaling the weights for correctly classified points hence allowing both losses to converge to zero. However, this will not change the decision boundary, but it will minimize the loss from both functions as it will increase $y(\text{transpose}(w))x$ hence making it greater than or equal to 1 for hinge and closer to infinity for log loss hence reducing both losses.
- e. Minimizing L_{hinge} loss alone does not change the decision boundary and only focuses on minimizing the loss. However, by adding a term that constrains the magnitude of the weight vector, the optimization problem becomes one of minimizing the loss by maximizing the weight margin while still correctly classifying the data points. This

ensures that the decision boundary is not only optimal in terms of minimizing loss but also in terms of correctly classifying the data points.

2. Effects Of Regularization

- Adding the regularization term will not decrease the in-sample (training) error and it might actually increase the in-sample error (or keep it the same) because it is used to avoid overfitting by adding a constraint to the model to achieve simpler models that do not fit noise hence it cannot decrease the in-sample error and might actually increase it as a result of potential underfitting. Adding the regularization term might decrease the out of sample (test) error if there was overfitting prior to adding regularization penalty as the addition will prevent this overfitting. However, if the regularization term is too large, then it might actually increase the out of sample error as our model will be underfitting the data hence making it a bad generalizer.
- L0 regularization is rarely used because it is not continuous and not convex hence not differentiable therefore making it extremely difficult to optimize using our gradient methods.



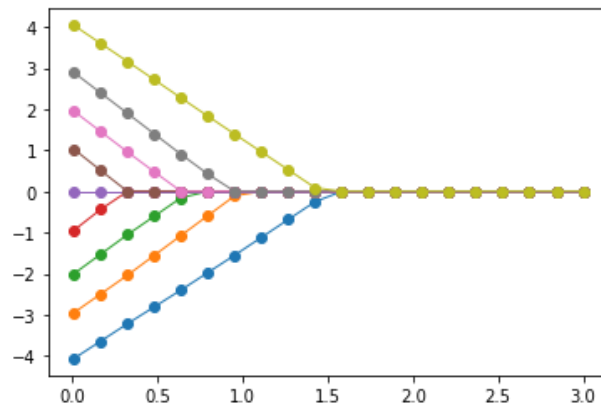
c.

(See Code)

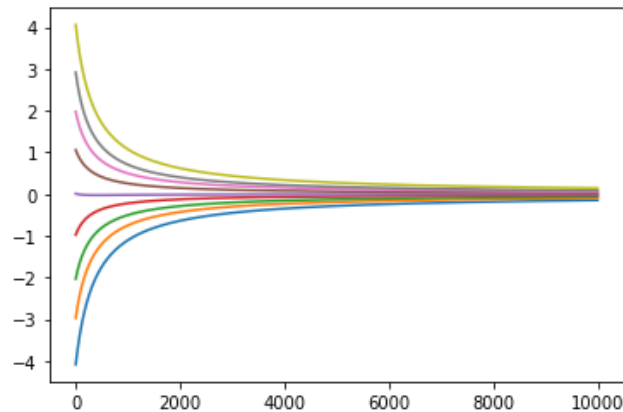
- Looking at the testing error we see that the wine training set with 100 points had a lower average Eout (out of sample error) for low lambda before increasing and starting to converge with the Eout of the dataset with 40 points for higher lambda towards the end. However, they had similar (minute differences in) Ein (Training Error) despite training set 1 having more training points. For training error, we see that for higher lambda values the training set with less points increases at a faster rate because it is smaller hence as the model gets more regularized, the training error for this set is more likely to change due to limits in model complexity.

- e. The training error (E_{in}) for wine_training1.txt increases with higher lambda values hence causing underfitting as lambda increases. The test error (E_{out}) is lower with no regularization and lower values for lambda implying no overfitting with the training set until around lambda equals 1 where it began to increase hence showing underfitting for higher lambda values.
- f. The L2 norm of w for wine_training1.txt decreases as lambda increases as the model is less constrained by the regularization term due to the equivalence between the norm constraint and regularization with the lagrangian.
- g. If the model were trained with wine_training2.txt, I would choose $\lambda = 0.00001 * (5^7)$ which is the point where E_{out} is the lowest hence allowing the model to best generalize.

3. Lasso (l1) vs. Ridge (l2) Regularization



a.



(See Code) As the regularization parameter increases, the number of model weights that are exactly zero increases with Lasso regression. For Ridge regression asymptotically approaches zero as the regularization parameter increases without ever reaching zero.

b. (B and C)

$$(B.1) \quad \argmin \|y - Xw\|^2 + \lambda \|w\|_1$$

$$\nabla_{\tilde{w}} = \partial(y - X\tilde{w}) : (X^T) + \lambda = 0$$

$$(y - X\tilde{w}) = \frac{\lambda}{\partial X}$$

$$y - \frac{\lambda}{\partial X} = X\tilde{w}$$

$$\frac{\partial^2 y - \lambda}{\partial^2 X} = w \quad \therefore \quad w = (\partial^2 y - \lambda) (\partial^2 X^T X)^{-1}$$

(B.2) Yes and that value is $\lambda = \|\partial^2 y\|$

$$(C.1) \quad \argmin \|y - Xw\|^2 + \lambda \|w\|_2^2$$

$$\nabla_w = \partial X^T (y - Xw) + \partial \lambda(w) = 0$$

$$y - Xw = \frac{\lambda w}{X^T X}$$

$$\therefore w = \frac{X^T y}{X^T X + \lambda I} \quad \therefore \quad w = (\tilde{X}^T \tilde{y}) (\tilde{X}^T \tilde{X} + \lambda I)^{-1}$$

(C.2) There does not exist such a value for λ where $w=0$.

4. Convexity and Lipschitz Continuity

a. (a and b).

(4a) we will show $y = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \in C$ for $n \geq 1$.
 using induction. by the problem statement, definition of convexity we have this holds true for $n=2$. It obviously holds for $n=1$ as that is equivalent to having each point $y \in C$ in C . Now assume it holds for $n-1$.

$$z = \frac{\theta_1 x_1}{\theta_1 + \theta_2 + \dots + \theta_{n-1}} + \dots + \frac{\theta_{n-1} x_{n-1}}{\theta_1 + \theta_2 + \dots + \theta_{n-1}} \in C$$

$$\text{and } y = (\theta_1 + \theta_2 + \dots + \theta_{n-1})z + \theta_n x_n \in C.$$

which gives us the line from z to x_n . Putting it in the convex set of C hence proven by induction that $\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \in C$.

Alternatively: By definition a convex set is a set where for any two points s in the set, the line segment connecting the two is contained within the set. As C is a convex set with $x_1, \dots, x_n \in C$, then the line segment between any two points in C is also contained within C hence as $\theta_i \geq 0$ with $\theta_1 + \dots + \theta_n = 1$.

$\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ is a convex combination of x_1, x_2, \dots, x_n hence meaning it is part of their connecting line segment. As the line segment between any two points in C is also contained within C (by definition) and $\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ is a convex combination of x_1, x_2, \dots, x_n it is definitely in C hence shown.

(4b) To show this, we need to

- show that \exists a real valued constant $L \geq 0 \forall x_1, x_2 \in D_1$ we have:

$$\star = \|f(t_1, \dots, t_n; 0, \dots, 0)(x_1) - f(t_1, \dots, t_n; 0, \dots, 0)(x_2)\| \leq L \|x_1 - x_2\|$$

where $f: D_f \rightarrow D_f$ is L -Lipschitz continuous.

for all $x_1, x_2 \in D_1$ $\|f_1(x_1) - f_1(x_2)\| \leq L \|x_1 - x_2\|$.

hence by this definition we can transform

$$\text{our function in } \star \text{ to get } \|f_1(t_1, \dots, t_n; 0, \dots, 0)(x_1) - f_1(t_1, \dots, t_n; 0, \dots, 0)(x_2)\| \leq L \|f_1(t_1, \dots, t_n; 0, \dots, 0)(x_1) - f_1(t_1, \dots, t_n; 0, \dots, 0)(x_2)\| \quad \triangleright$$

$$\leq L^2 \|f_1(t_1, \dots, t_n; 0, \dots, 0)(x_1) - f_1(t_1, \dots, t_n; 0, \dots, 0)(x_2)\|, \text{ and so on } \leq L^n \|x_1 - x_2\|, \text{ making}$$

f_1, f_2, \dots, f_n as L_i -Lipschitz continuous w $L_i = L$.