

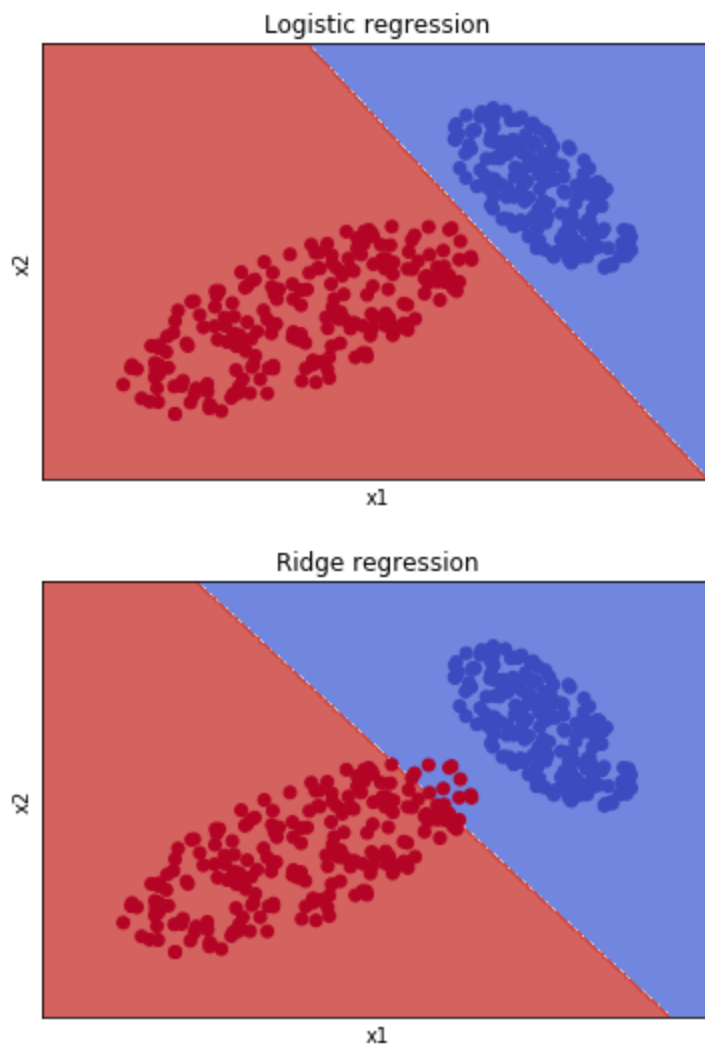
CS 155 PS 2

1 Comparing Different Loss Functions

Problem A

Squared loss is a terrible choice of loss function to train on for classification problems because it penalizes large values even if classified correctly. Since it only cares about how far the value of $\mathbf{w}^T \mathbf{x}$ is from the target (-1 or +1), points that are well classified but far from the target (i.e. $y = +1$, $\mathbf{w}^T \mathbf{x} = +5$) will have high loss simply because of the distance.

Problem B



The decision boundary for logistic regression perfectly classifies the points, but the decision boundary for ridge regression is shifted to the bottom left, misclassifying some of the red points. This is because the squared loss function penalizes correctly classified points that are far from the boundary, causing the boundary to shift towards such points to decrease the loss values for them. On the other hand, log loss

function only penalizes misclassified values, according to their magnitude, and furthermore produces lower loss for correctly classified points farther from the boundary.

Problem C

$$\nabla_{\mathbf{w}} L_{\text{hinge}} = \begin{cases} \vec{0} & \text{if } y \vec{w}^T \vec{x} > 1 \\ -y \vec{x} & \text{if } y \vec{w}^T \vec{x} \leq 1 \end{cases}$$

Given $w_0 = 0, w_1 = 1, w_2 = 0,$

$$\nabla_{\mathbf{w}} L_{\text{hinge}} = \begin{cases} \vec{0} & \text{if } y x_1 > 1 \\ -y \vec{x} & \text{if } y x_1 \leq 1 \end{cases}$$

$$\nabla_{\mathbf{w}} L_{\text{hinge}} |_{x=(\frac{1}{2}, 3), y=1} = (-1, -\frac{1}{2}, -3)$$

$$\nabla_{\mathbf{w}} L_{\text{hinge}} |_{x=(2, -2), y=1} = (0, 0, 0)$$

$$\nabla_{\mathbf{w}} L_{\text{hinge}} |_{x=(-3, 1), y=-1} = (0, 0, 0)$$

$$\nabla_{\mathbf{w}} L_{\text{log}} = \frac{-y \vec{x}}{1 + e^{y \vec{w}^T \vec{x}}}$$

Given $w_0 = 0, w_1 = 1, w_2 = 0,$

$$\nabla_{\mathbf{w}} L_{\text{log}} = \frac{-y \vec{x}}{1 + e^{y x_1}}$$

$$\nabla_{\mathbf{w}} L_{\text{log}} |_{x=(\frac{1}{2}, 3), y=1} = \frac{-1}{1 + e^{\frac{1}{2}}} (1, \frac{1}{2}, 3) = (-.378, -.189, -1.133)$$

$$\nabla_{\mathbf{w}} L_{\text{log}} |_{x=(2, -2), y=1} = \frac{1}{1 + e^2} (-1, -2, 2) = (-.119, -.238, .238)$$

$$\nabla_{\mathbf{w}} L_{\text{log}} |_{x=(-3, 1), y=-1} = \frac{1}{1 + e^3} (1, -3, 1) = (.047, -.142, .047)$$

Problem D

For the points with absolute value of x_1 greater than or equal to 1, the gradients for the hinge loss are zero (since they're all correctly classified). For the same points, the gradients for the log loss are small but not zero. The gradients for the log loss decrease as they get farther from the decision boundary (as long as classified correctly). The gradients for the hinge loss become zero as soon as $y * \mathbf{w}^T \mathbf{x}$ becomes greater than or equal to 1 (outside of margin on the correct side), while the gradients for the log loss converge to zero as $y * \mathbf{w}^T \mathbf{x}$ goes to infinity (as far from the margin as possible on the correct side). Scaling up the weight vector (including the bias term) will reduce the training error if the points are correctly classified,

while preserving the decision boundary, since it will increase the value of $y * \mathbf{w}^T \mathbf{x}$. The hinge loss will be completely eliminated if $y * \mathbf{w}^T \mathbf{x}$ is greater than or equal to 1, and the log loss will be reduced as $y * \mathbf{w}^T \mathbf{x}$ gets scaled up higher.

Problem E

Minimizing just L_{hinge} could be done by scaling up \mathbf{w} without changing the decision boundary, which is what we want to control in order to classify the points correctly and maximize the margin. The penalty term addresses this issue by constraining the magnitude of \mathbf{w} .

2 Effects of Regularization

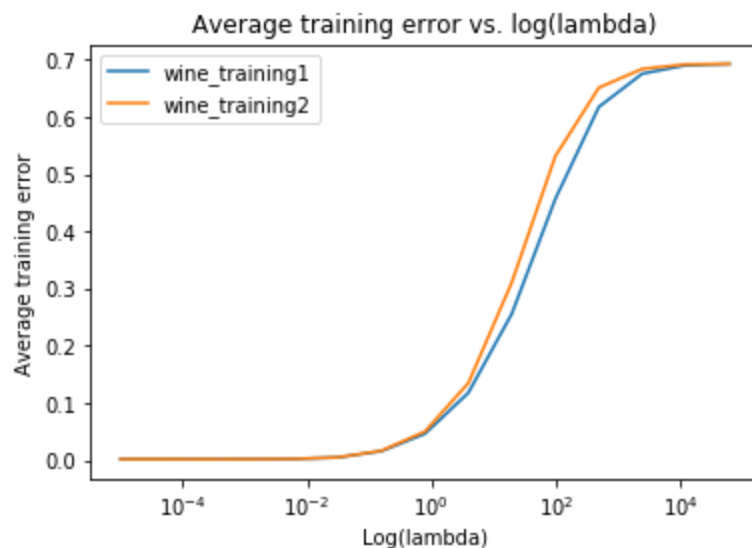
Problem A

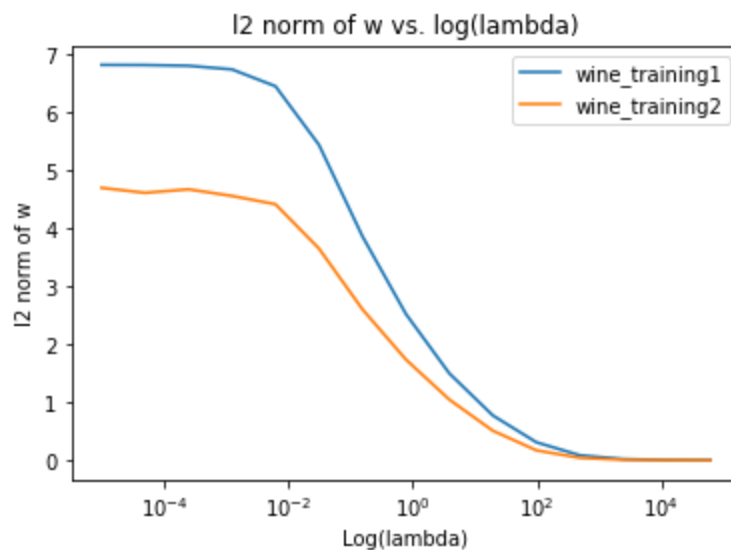
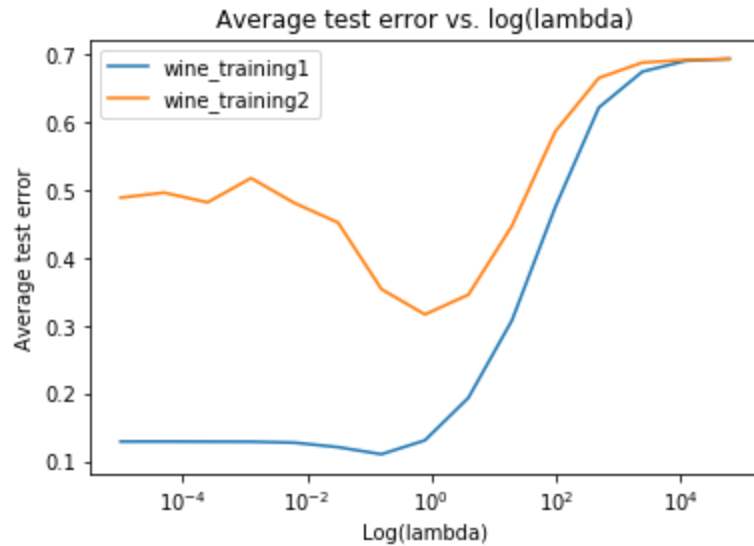
Adding the penalty term can't decrease the training error, since the term generally prevents the model from fitting perfectly to the training sample. The training error can only stay the same or increase due to a smaller hypothesis set. Adding a penalty term will decrease the out-of-sample errors only when there is overfitting (such as when the training set is small and model complexity is high). If the regularization parameter is too large, the out-of-sample errors can increase due to underfitting.

Problem B

The regularization term of the objective function for l_0 regularization is discontinuous and non-convex, making the optimization very difficult.

Problem C



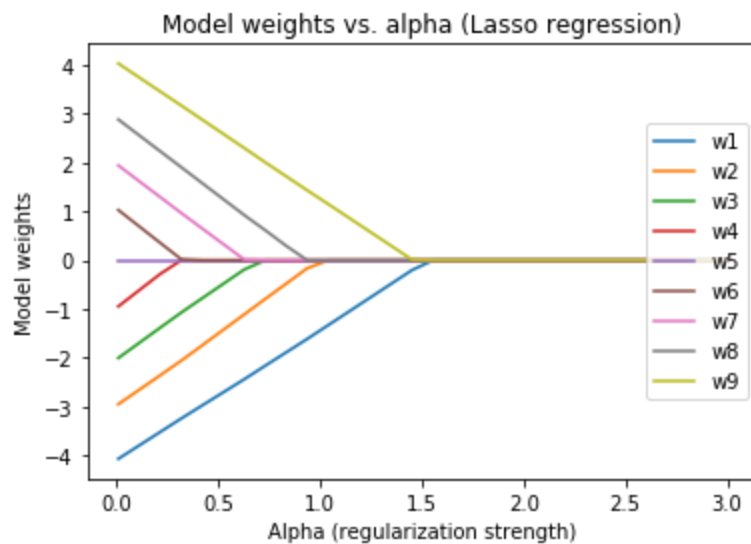


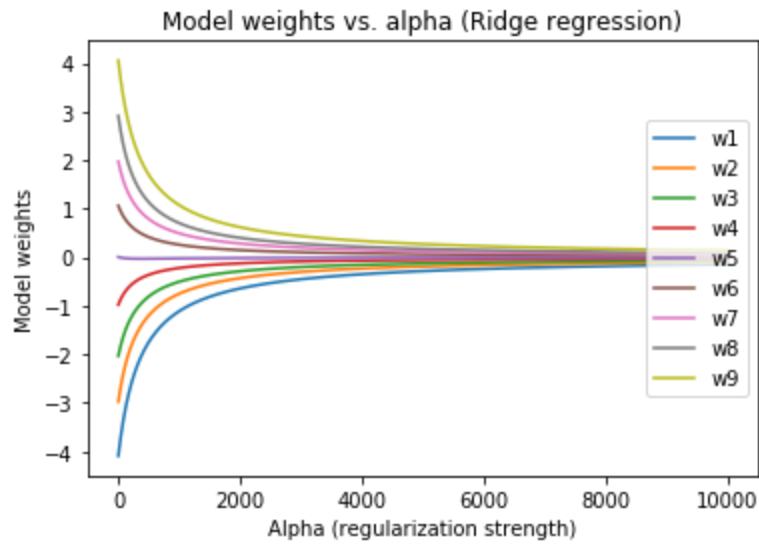
Problem D

The difference between the training errors for the two training sets is minimal. As λ increases, the training error for wine_training2 increases at a slightly faster rate than wine_training1. Since wine_training1 has fewer data points, its training error is more susceptible to the change (the limitation of the model) as the model gets more regularized. The test error for wine_training2 is much higher with low λ , because its small sample size causes it to overfit. With the optimal amount of regularization, the test error for wine_training2 goes down, yet it is always higher (until the errors converge) than the test error for wine_training1 which is a better approximation of the true distribution.

Problem E

The training error for wine_training1 only increases with increasing λ , showing the underfitting behavior (the model is limited in fitting the training data perfectly) as λ increases. The test error for wine_training1 is quite low even with no regularization, and reaches its minimum at a λ value





As the regularization parameter increases, the number of model weights that are exactly zero increases with Lasso regression. The model weights for Ridge regression asymptotically reach zero as the regularization parameter increases, but none of the weights become exactly zero.

Problem B

i. To minimize $\|\vec{y} - \vec{x}w\|^2 + \lambda\|w\|$, with respect to w , we want to differentiate this function w.r.t. w and set it equal to 0. Since $\lambda\|w\|$ isn't differentiable, we will use its subgradient definition:

$$\nabla_w \lambda\|w\| = \begin{cases} -\lambda & \text{if } w < 0 \\ +\lambda & \text{if } w > 0 \\ [-\lambda, +\lambda] & \text{if } w = 0 \end{cases}$$

$$\text{Thus } \frac{d}{dw} (\|\vec{y} - \vec{x}w\|^2 + \lambda\|w\|)$$

$$= -2\vec{x}^T(\vec{y} - \vec{x}w) + \lambda \text{ if } w > 0, \quad -2\vec{x}^T(\vec{y} - \vec{x}w) - \lambda \text{ if } w < 0$$

Solving for w using $\frac{d}{dw} (\|\vec{y} - \vec{x}w\|^2 + \lambda\|w\|) = 0$ gives

$$w = (2\vec{x}^T\vec{y} - \lambda)(2\vec{x}^T\vec{x})^{-1} \quad \text{if } w > 0$$

$$w = (2\vec{x}^T\vec{y} + \lambda)(2\vec{x}^T\vec{x})^{-1} \quad \text{if } w < 0$$

$w > 0$ only if $2\vec{x}^T\vec{y} > \lambda$, given $\lambda \geq 0$

$w < 0$ only if $2\vec{x}^T\vec{y} < -\lambda$, given $\lambda \geq 0$

$w = 0$ only if $-\lambda \leq 2\vec{x}^T\vec{y} \leq \lambda$

Thus

$$w = \begin{cases} (2\vec{x}^T\vec{y} - \lambda)(2\vec{x}^T\vec{x})^{-1} & \text{if } 2\vec{x}^T\vec{y} > \lambda \\ (2\vec{x}^T\vec{y} + \lambda)(2\vec{x}^T\vec{x})^{-1} & \text{if } 2\vec{x}^T\vec{y} < -\lambda \\ 0 & \text{if } -\lambda \leq 2\vec{x}^T\vec{y} \leq \lambda \end{cases}$$

ii. Yes; as shown above, $w = 0$ if $-\lambda \leq 2\vec{x}^T\vec{y} \leq \lambda$

Thus the smallest value for λ s.t. $w = 0$ is $|2\vec{x}^T\vec{y}|$.

If negative λ values are allowed, $-\infty$ would be the smallest such value.

iii. To minimize $\|\vec{y} - \vec{X}\vec{w}\|^2 + \lambda \|\vec{w}\|_2^2$ wrt. \vec{w} , we want to take the gradient of this function wrt. \vec{w} and set it equal to zero.

$$\nabla_{\vec{w}}(\|\vec{y} - \vec{X}\vec{w}\|^2 + \lambda \|\vec{w}\|_2^2) \\ = -2\vec{X}^T(\vec{y} - \vec{X}\vec{w}) + \lambda\vec{w} = 0$$

Solving for \vec{w} gives

$$\vec{w} = (\vec{X}^T\vec{X} + \lambda\mathbf{I})^{-1}\vec{X}^T\vec{y}$$

iv. Given $w \neq 0$ when $\lambda = 0$,

we know that $\vec{X} \neq 0$ and $\vec{y} \neq 0$, since

$$w = \vec{X}^T\vec{y}(\vec{X}^T\vec{X})^{-1} \text{ when } \lambda = 0$$

and either \vec{X} or \vec{y} being 0 causes w to be 0.

Therefore, regardless of λ value, w can't be 0.

Intuitively, λ is in the denominator and can't make the expression 0.

\therefore There does not exist a value for $\lambda > 0$ s.t. $w = 0$.