# 1   Comparing Different Loss Functions [30 Points]

*Relevant materials: lecture 3 & 4*

We've discussed three loss functions for linear classification models so far:

- Squared loss: $L_{\text{squared}} = (1 - y\mathbf{w}^T\mathbf{x})^2$

- Hinge loss: $L_{\text{hinge}} = \max(0, 1 - y\mathbf{w}^T\mathbf{x})$

- Log loss: $L_{\text{log}} = \ln(1 + e^{-y\mathbf{w}^T\mathbf{x}})$

where $\mathbf{w} \in \mathbb{R}^n$ is a vector of the model parameters, $y \in \{-1, 1\}$ is the class label for datapoint $\mathbf{x} \in \mathbb{R}^n$, and we're including a bias term in $\mathbf{x}$ and $\mathbf{w}$. The model classifies points according to $\text{sign}(\mathbf{w}^T\mathbf{x})$.

Performing gradient descent on any of these loss functions will train a model to classify more points correctly, but the choice of loss function has a significant impact on the model that is learned.

**Problem A [3 points]:**  Squared loss is often a terrible choice of loss function to train on for classification problems. Why?

> **Solution A:** *Squared loss can be very inefficient during training for classification problems because datapoints that are far from the classification boundary but still classified correctly will result in a large error.*
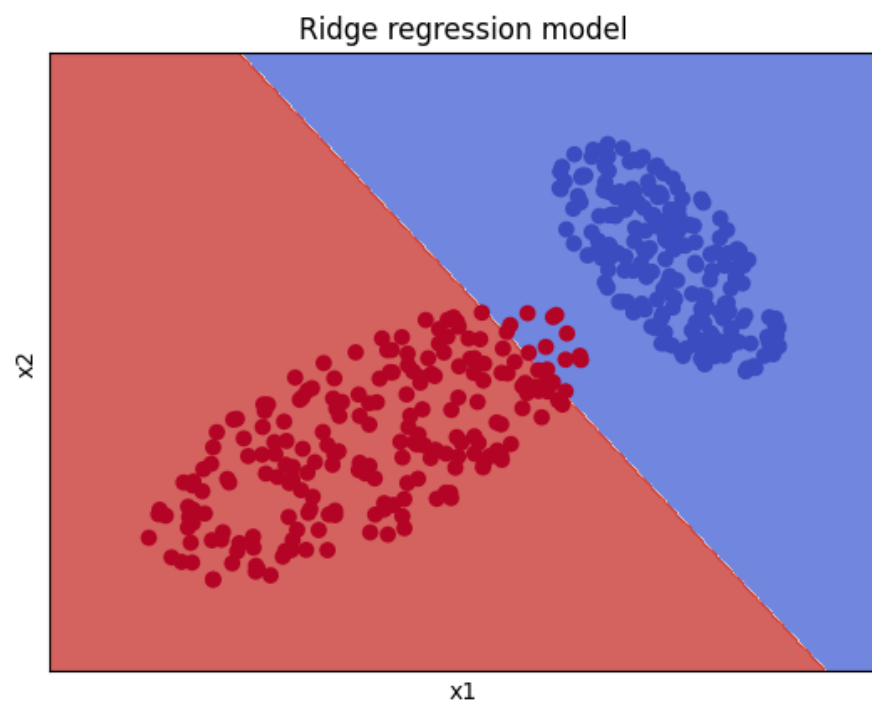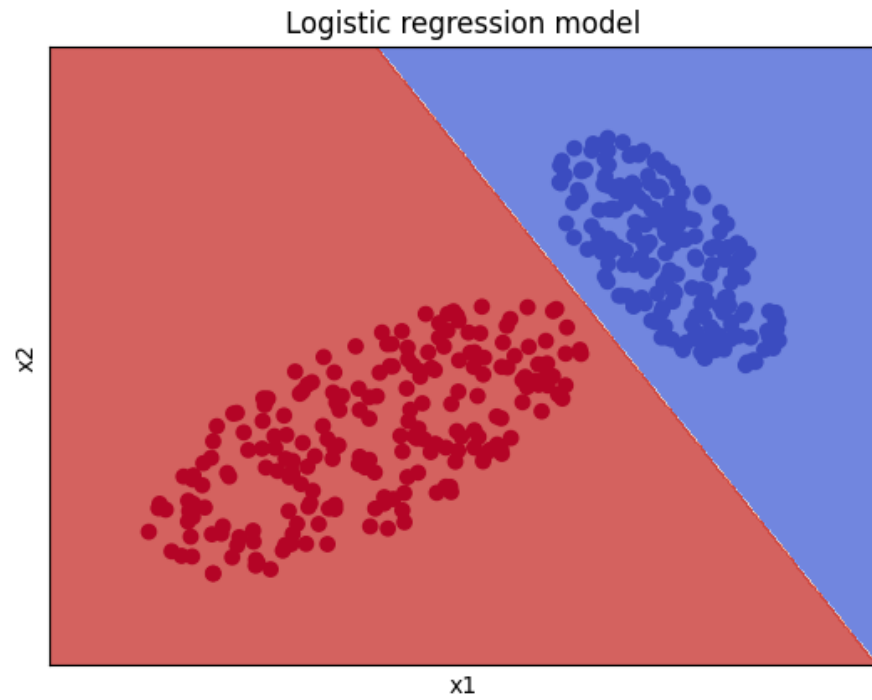
**Problem B [9 points]:**  A dataset is included with your problem set: `problem1data1.txt`. The first two columns represent $x_1, x_2$, and the last column represents the label, $y \in \{-1, +1\}$.

On this dataset, train both a logistic regression model and a ridge regression model to classify the points. (In other words, on each dataset, train one linear classifier using $L_{\text{log}}$ as the loss, and another linear classifier using $L_{\text{squared}}$ as the loss.) For this problem, you should use the logistic regression and ridge regression implementations provided within scikit-learn (logistic regression documentation) (Ridge regression documentation) instead of your own implementations. Use the default parameters for these classifiers except for setting the regularization parameters so that very little regularization is applied.

For each loss function/model, plot the data points as a scatter plot and overlay them with the decision boundary defined by the weights of the trained linear classifier. Include both plots in your submission. The template notebook for this problem contains a helper function for producing plots given a trained classifier.

What differences do you see in the decision boundaries learned using the different loss functions? Provide a qualitative explanation for this behavior.
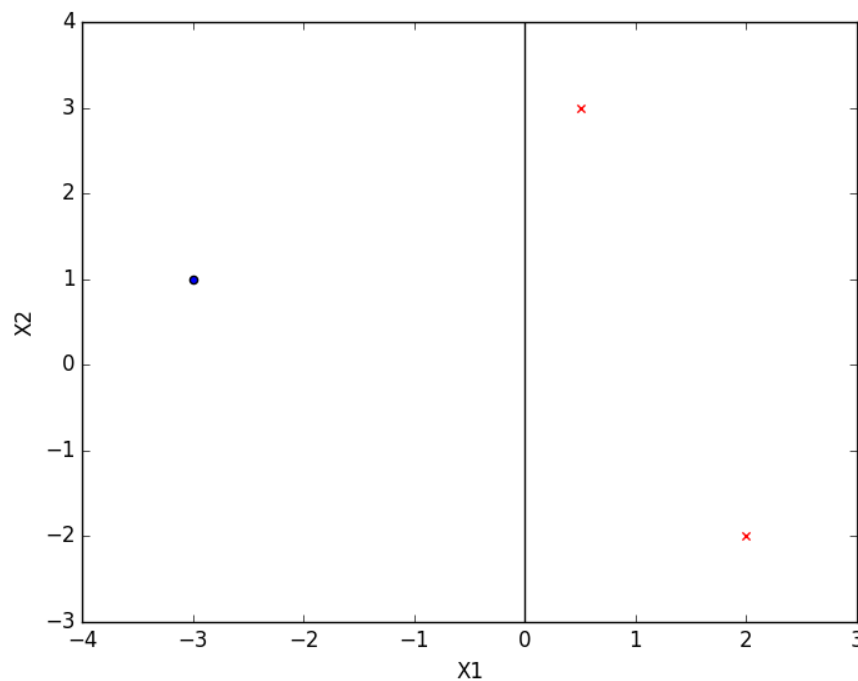
> **Solution B:** *Code: mantripragada_ishaan_set2_prob1.ipynb*

Logistic regression model



Ridge regression model

*The logisitic regression seems to have classified all the data points correctly while the ridge regression model misclassifies several datapoints. This may be due to the fact that ridge loss utilizes square loss. As described in part a, squared loss can be inefficient when classifying points very close or very far from the boundary. In this case, even though certain data points were classified incorrectly, they had a very small error associated with them.*

**Problem C [9 points]:** Leaving squared loss behind, let's focus on log loss and hinge loss. Consider the set of points $S = \{(\frac{1}{2}, 3), (2, -2), (-3, 1)\}$ in 2D space, shown below, with labels $(1, 1, -1)$ respectively.

Given a linear model with weights $w_0 = 0, w_1 = 1, w_2 = 0$ (where $w_0$ corresponds to the bias term), compute the gradients $\nabla_w L_{\text{hinge}}$ and $\nabla_w L_{\text{log}}$ of the hinge loss and log loss, and calculate their values for each point in S.



The example dataset and decision boundary described above. Positive instances are represented by red x's, while negative instances appear as blue dots.

**Solution C:**

1.c  $S = \{(\frac{1}{2}, 3), (2, -2), (-3, 1)\}$

$L_{hinge} = \max \{0, 1 - y_i \cdot f(x_i | w, b)\}$    $L_{log} = \log(1 + \exp(-y_i f(x_i)))$

$\nabla_{w_0} L_{hinge} = 0$ if $y \cdot f(x) \geq 1$     $\nabla_{w_0} L_{log} = \dfrac{y}{1 + \exp(y \cdot f(x))}$

$\nabla_{w_0} L_{hinge} = -y$ if $y \cdot f(x) < 1$

$\nabla_{w_1} L_{hinge} = -y x_1$ if $y \cdot f(x) < 1$     $\nabla_{w_1} L_{log} = \dfrac{y x_1}{1 + \exp(y \cdot f(x))}$

$\nabla_{w_2} L_{hinge} = -y x_2$ if $y \cdot f(x) < 1$     $\nabla_{w_2} L_{log} = \dfrac{y x_2}{1 + \exp(y \cdot f(x))}$

$(\frac{1}{2}, 3)$:   $\nabla_w L_{hinge} = (-1, -0.5, -3)$

$\nabla_w L_{log} = (-0.378, -0.189, -1.133)$

$(2, -2)$:   $\nabla_w L_{hinge} = (0, 0, 0)$

$\nabla_w L_{log} = (-0.119, -0.238, 0.238)$

$(-3, 1)$:   $\nabla_w L_{hinge} = (0, 0, 0)$

$\nabla_w L_{log} = (0.047, -0.142, 0.047)$

**Problem D [4 points]:**  Compare the gradients resulting from log loss to those resulting from hinge loss. When (if ever) will these gradients converge to 0? For a linearly separable dataset, is there any way to reduce or altogether eliminate training error without changing the decision boundary?

**Solution D:** *The hinge loss is simply the zero vector when $y \cdot f(x) > 1$. Therefore the gradients will eventually converge when $y_i \cdot (w^t x_i) > 1$. On the other hand, the log loss gets smaller when $|x|$ is large and is correctly classified. To reduce or eliminate training error altogether without changing the boundary, there must also be some regularization to the weights. For hinge loss, as long as all $y_i \cdot (w^t x_i)$ terms are greater than 1, there will be no training error. For log loss, the error cannot be fully eliminated, but it can be reduced by also adjusting the weights.*

**Problem E [5 points]:**  Based on your answer to the previous question, explain why for an SVM to be a "maximum margin" classifier, its learning objective must not be to minimize just $L_{hinge}$, but to minimize $L_{hinge} + \lambda \|w\|^2$ for some $\lambda > 0$.

---

(You don't need to prove that minimizing $L_{\text{hinge}} + \lambda\|w\|^2$ results in a maximum margin classifier; just show that the additional penalty term addresses the issues of minimizing just $L_{\text{hinge}}$.)

> **Solution E:** *An SVM's learning objective must be to minimize $L_{hinge} + \lambda\|w\|^2$ instead of just $L_{hinge}$ because, as stated in the previous part, this can be done simply by adjusting the weights such that all $y_i \cdot (w^t x_i)$ terms are greater than 1. This second term that the SVM must account for helps restrict and optimize the scaling of the weights with the margin.*

## 2  Effects of Regularization [40 Points]

*Relevant materials: Lecture 3 & 4*

For this problem, you are required to implement everything yourself and submit code (i.e. don't use scikit-learn but numpy is fine).

**Problem A [4 points]:** In order to prevent over-fitting in the least-squares linear regression problem, we add a regularization penalty term. Can adding the penalty term decrease the training (in-sample) error? Will adding a penalty term always decrease the out-of-sample errors? Please justify your answers. Think about the case when there is over-fitting while training the model.

> **Solution A:** *Adding the penality term will not further decrease the in sample training error because this is already done through minimizing the loss. In the case when there is over-fitting during training, the regularization penalty term can decrease the out of sample errors sicne it makes the model less specific to the given dataset and more generalized. This is not always the case because too much generalization can result in underfitting.*

**Problem B [4 points]:** $\ell_1$ regularization is sometimes favored over $\ell_2$ regularization due to its ability to generate a sparse $w$ (more zero weights). In fact, $\ell_0$ regularization (using $\ell_0$ norm instead of $\ell_1$ or $\ell_2$ norm) can generate an even sparser $w$, which seems favorable in high-dimensional problems. However, it is rarely used. Why?

> **Solution B:** *$\ell_0$ is rarely used since it is not a continuous function. Since the $\ell_0$ norm is discontinuous, it does not have a well-defined gradient everywhere, making it incompatible with standard gradient-based optimization techniques.*

### Implementation of $\ell_2$ regularization:

We are going to experiment with regression for the Red Wine Quality Rating data set. The data set is uploaded on the course website, and you can read more about it here: https://archive.ics.uci.edu/ml/datasets/Wine. The data relates 13 different factors (last 13 columns) to wine type (the first column). Each column of data represents a different factor, and they are all continuous features. Note that the original data set has three classes, but one was removed to make this a binary classification problem.

Download the data for training and validation from the assignments data folder. There are two training sets, wine_training1.txt (100 data points) and wine_training2.txt (a proper subset of wine_training1.txt containing only 40 data points), and one test set, wine_validation.txt (30 data points). You will use the wine_validation.txt dataset to evaluate your models.

We will train a $\ell_2$-*regularized logistic regression* model on this data. Recall that the unregularized logistic error (a.k.a. log loss) is

$$E = -\sum_{i=1}^{N} \log(p(y_i|\mathbf{x}_i))$$

where $p(y_i = -1|\mathbf{x}_i)$ is

$$\frac{1}{1 + e^{\mathbf{w}^T\mathbf{x}_i}}$$

and $p(y_i = 1|\mathbf{x}_i)$ is

$$\frac{1}{1 + e^{-\mathbf{w}^T\mathbf{x}_i}},$$

where as usual we assume that all $\mathbf{x}_i$ contain a bias term. The $\ell_2$-regularized logistic error is

$$E = -\sum_{i=1}^{N} \log(p(y_i|\mathbf{x}_i)) + \lambda\mathbf{w}^T\mathbf{w}$$

$$= -\sum_{i=1}^{N} \log\left(\frac{1}{1 + e^{-y_i\mathbf{w}^T\mathbf{x}_i}}\right) + \lambda\mathbf{w}^T\mathbf{w}$$

$$= -\sum_{i=1}^{N} \left(\log\left(\frac{1}{1 + e^{-y_i\mathbf{w}^T\mathbf{x}_i}}\right) - \frac{\lambda}{N}\mathbf{w}^T\mathbf{w}\right).$$

Implement SGD to train a model that minimizes the $\ell_2$-regularized logistic error, i.e. train an $\ell_2$-regularized logistic regression model. Train the model with 15 different values of $\lambda$ starting with $\lambda_0 = 0.00001$ and increasing by a factor of 5, i.e.

$$\lambda_0 = 0.00001, \lambda_1 = 0.00005, \lambda_2 = 0.00025, ..., \lambda_{14} = 61,035.15625.$$

Some important notes: Terminate the SGD process after 20,000 epochs, where each epoch performs one SGD iteration for each point in the training dataset. You should shuffle the order of the points before each epoch such that you go through the points in a random order (hint: use `numpy.random.permutation`). Use a learning rate of $5 \times 10^{-4}$, and initialize your weights to small random numbers.

You may run into numerical instability issues (overflow or underflow). One way to deal with these issues is by normalizing the input data $X$. Given the column for the $j$th feature, $X_{:,j}$, you can normalize it by setting $X_{ij} = \frac{X_{ij} - \overline{X_{:,j}}}{\sigma(X_{:,j})}$ where $\sigma(X_{:,j})$ is the standard deviation of the $j$th column's entries, and $\overline{X_{:,j}}$ is the mean of the $j$th column's entries. Normalization may change the optimal choice of $\lambda$; the $\lambda$ range given above corresponds to data that has been normalized in this manner. If you treat the input data differently, simply plot enough choices of $\lambda$ to see any trends.
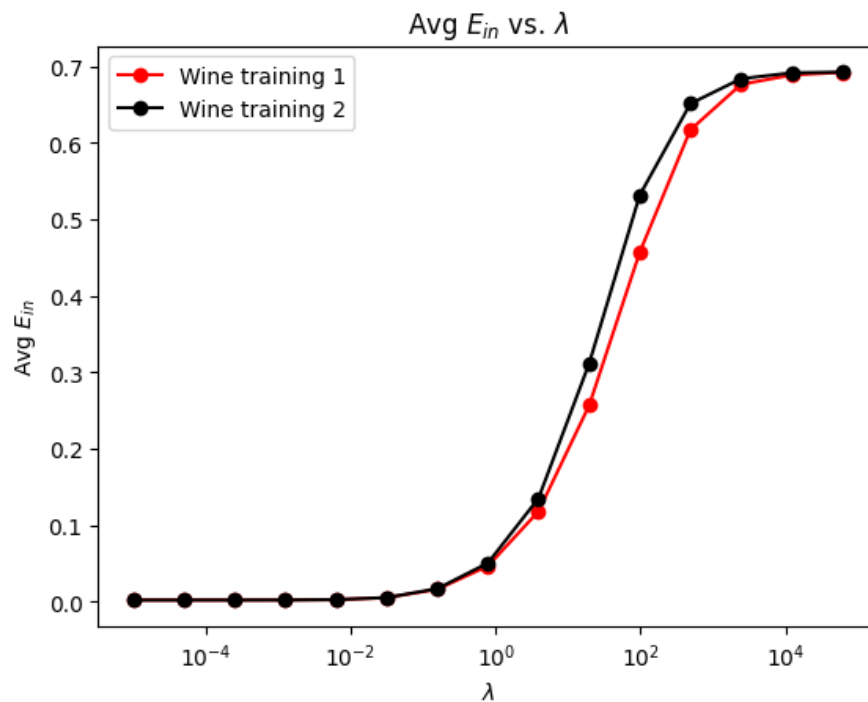
**Problem C [16 points]:**   Do the following for both training data sets (wine_training1.txt and wine_training2.txt) and attach your plots in the homework submission (use a log-scale on the horizontal axis):
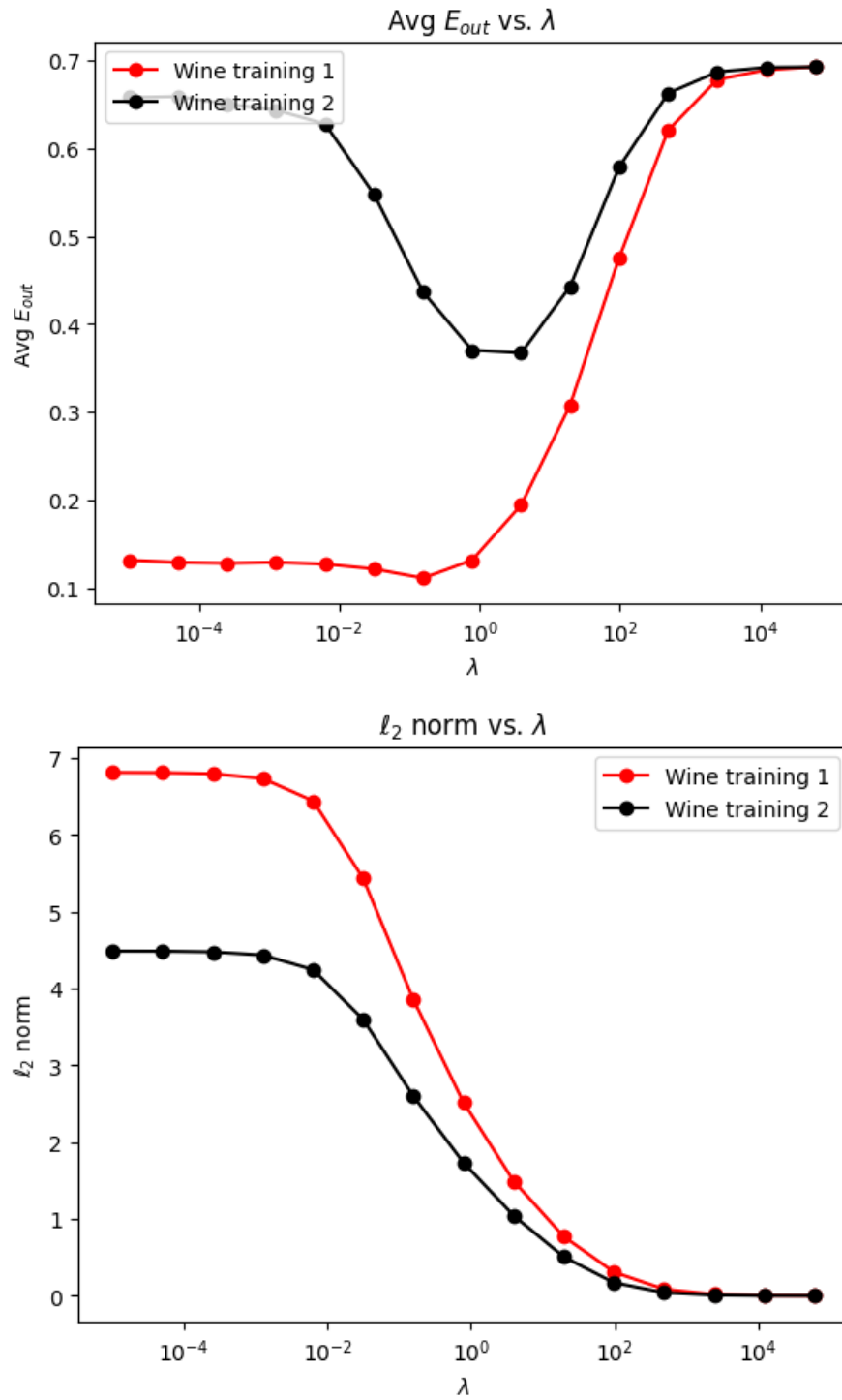
**i.** Plot the average training error ($E_{\text{in}}$) versus different $\lambda$s.

**ii.** Plot the average test error ($E_{\text{out}}$) versus different $\lambda$s using wine_validation.txt as the test set.

**iii.** Plot the $\ell_2$ norm of **w** versus different $\lambda$s.

You should end up with three plots, with two series (one for wine_training1.txt and one for wine_training2.txt) on each plot. Note that the $E_{\text{in}}$ and $E_{\text{out}}$ values you plot should not include the regularization penalty — the penalty is only included when performing gradient descent.

**Solution C:** *Code: mantripragada_ishaan_set2_prob2.ipynb*

Avg $E_{out}$ vs. $\lambda$

$\ell_2$ norm vs. $\lambda$

**Problem D [4 points]:** Given that the data in wine_training2.txt is a subset of the data in wine_training1.txt, compare errors (training and test) resulting from training with wine_training1.txt (100 data points) versus wine_training2.txt (40 data points). Briefly explain the differences.

> **Solution D:** *The training and test curves are similar, but the training error is slightly higher for wine_training2. The testing error curves have more differences between the datasets. For wine_training1, the testing error is much lower than wine_training2. The reason for this may be because it had more data to train on. We can also see indications of overfitting for wine_training2 for larger lambdas.*

**Problem E [4 points]:** Briefly explain the qualitative behavior (i.e. over-fitting and under-fitting) of the training and test errors with different $\lambda$s while training with data in wine_training1.txt.

> **Solution E:** *The high out of sample error and low in sample error at low $\lambda$ values indicates overfitting in this region. On the other hand at higher $\lambda$ values, we see indications of underfitting as both the in sample and out of sample errors become close to equal.*

**Problem F [4 points]:** Briefly explain the qualitative behavior of the $\ell_2$ norm of **w** with different $\lambda$s while training with the data in wine_training1.txt.

> **Solution F:** *For wine_training1, the $\ell_2$ norm decreases as $\lambda$ decreases. This makes sense because as the regularization increases, the regularization term becomes minimized, which results in a smaller value for the weight vector **w**.*

**Problem G [4 points]:** If the model were trained with wine_training2.txt, which $\lambda$ would you choose to train your final model? Why?

> **Solution G:** *I would choose $\lambda = 10^0 = 1$ because the out of sample error is at its lowest.*

# 3   Lasso ($\ell_1$) vs. Ridge ($\ell_2$) Regularization [30 Points]
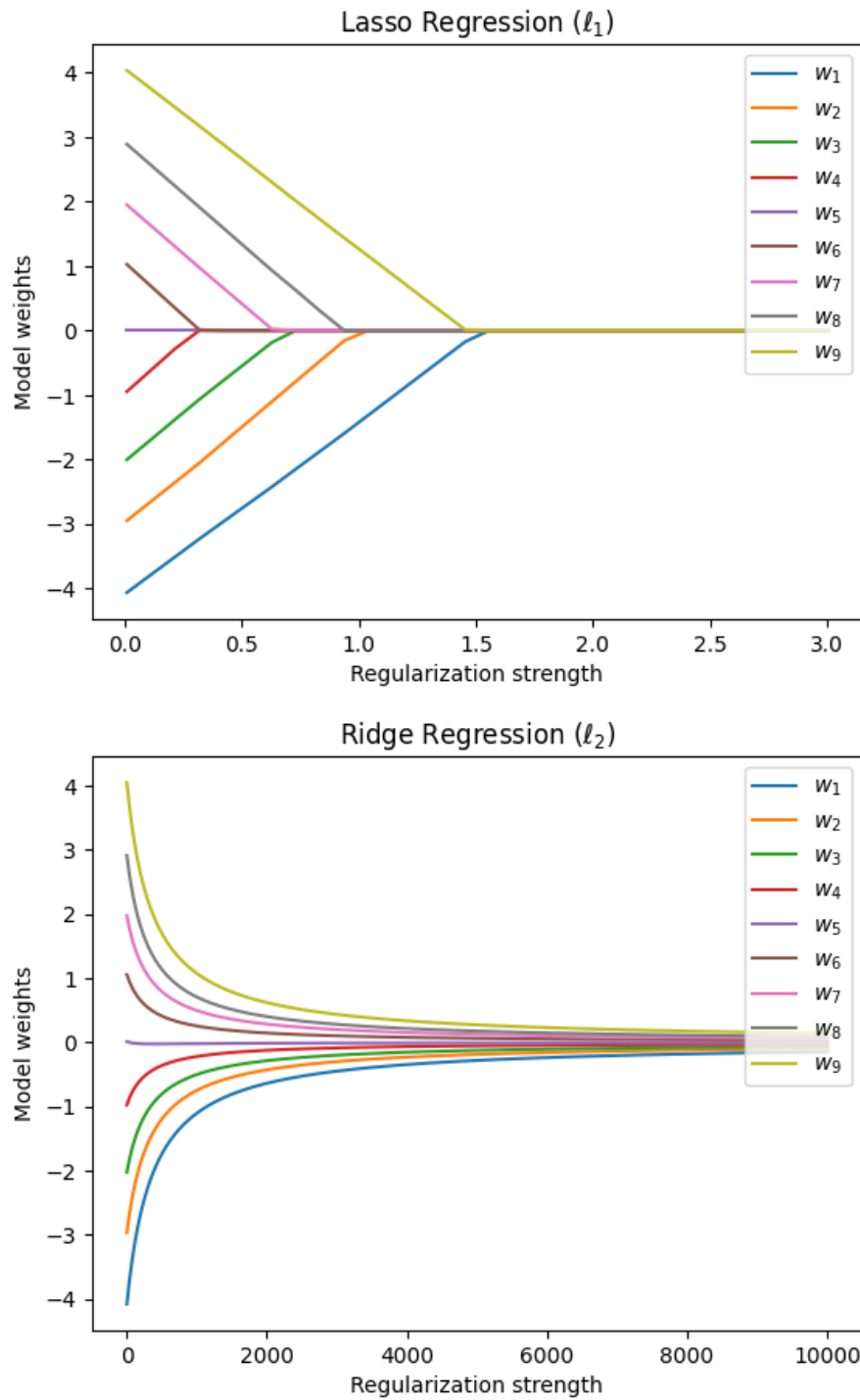
*Relevant materials: Lecture 3*

For this problem, you may use the scikit-learn (or other Python package) implementation of Lasso and Ridge regression — you don't have to code it yourself.

The two most commonly-used regularized regression models are Lasso ($\ell_1$) regression and Ridge ($\ell_2$) regression. Although both enforce "simplicity" in the models they learn, only Lasso regression results in sparse weight vectors. This problem compares the effect of the two methods on the learned model parameters.

**Problem A [12 points]:**   The tab-delimited file problem3data.txt on the course website contains 1000 9-dimensional datapoints. The first 9 columns contain $x_1, \ldots, x_9$, and the last column contains the target value $y$.

**i.** Train a linear regression model on the problem3data.txt data with Lasso regularization for regularization strengths $\alpha$ in the vector given by `numpy.linspace(0.01, 3, 30)`. On a single plot, plot each of the model weights $w_1, ..., w_9$ (ignore the bias/intercept) as a function of $\alpha$.

**ii.** Repeat **i.** with Ridge regression, and this time using regularization strengths $\alpha \in \{1, 2, 3, \ldots, 1e4\}$.

**iii.**   As the regularization parameter increases, what happens to the number of model weights that are exactly zero with Lasso regression? What happens to the number of model weights that are exactly zero with Ridge regression?

---

**Solution A:** *Code: mantripragada_ishaan_set2_prob3.ipynb*

---

**Problem B [9 points]:**

**i.** In the case of 1-dimensional data, Lasso regression admits a closed-form solution. Given a dataset containing $N$ datapoints, each with $d = 1$ feature, solve for

$$\arg\min_{w}\|\mathbf{y} - \mathbf{x}w\|^2 + \lambda\|w\|_1,$$

where $\mathbf{x} \in \mathbb{R}^N$ is the vector of datapoints and $\mathbf{y} \in \mathbb{R}^N$ is the vector of all output values corresponding to these datapoints. Just consider the case where $d = 1$, $\lambda \geq 0$, and the weight $w$ is a scalar.

This is linear regression with Lasso regularization.

---

**Solution B.i:**

$$3.b.i \quad \arg\min_{w}\|y - xw\|^2 + \lambda\|w\|_1$$

$$\nabla_w\left(\|y - xw\|^2 + \lambda\|w\|_1\right) = 0$$
$$-2x^T(y - xw) + \lambda = 0$$
$$-2x^Ty + 2xx^Tw + \lambda = 0$$
$$2xx^Tw = 2x^Ty - \lambda$$
$$\boxed{w = \frac{2x^Ty - \lambda}{2xx^T}}$$

---

**ii.** In this question, we continue to consider Lasso regularization in 1-dimension. Now, suppose that $w \neq 0$ when $\lambda = 0$. Does there exist a value for $\lambda$ such that $w = 0$? If so, what is the smallest such value?

---

**Solution B.ii:**

$$3.b.ii \quad 0 = \frac{2x^Ty - \lambda}{2xx^T} \implies \boxed{\lambda = 2x^Ty}$$

---

**Problem C [9 points]:**

**i.** Given a dataset containing $N$ datapoints each with $d$ features, solve for

$$\arg\min_{\mathbf{w}}\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda\|\mathbf{w}\|_2^2$$

where $\mathbf{X} \in \mathbb{R}^{N \times d}$ is the matrix of datapoints and $\mathbf{y} \in \mathbb{R}^N$ is the vector of all output values for these datapoints. Do so for arbitrary $d$ and $\lambda \geq 0$.

This is linear regression with Ridge regularization.

**Solution C.i:**



**ii.** In this question, we consider Ridge regularization in 1-dimension. Suppose that $w \neq 0$ when $\lambda = 0$. Does there exist a value for $\lambda > 0$ such that $w = 0$? If so, what is the smallest such value?

**Solution C.ii:**