# Homework2 (late for 5 hours)
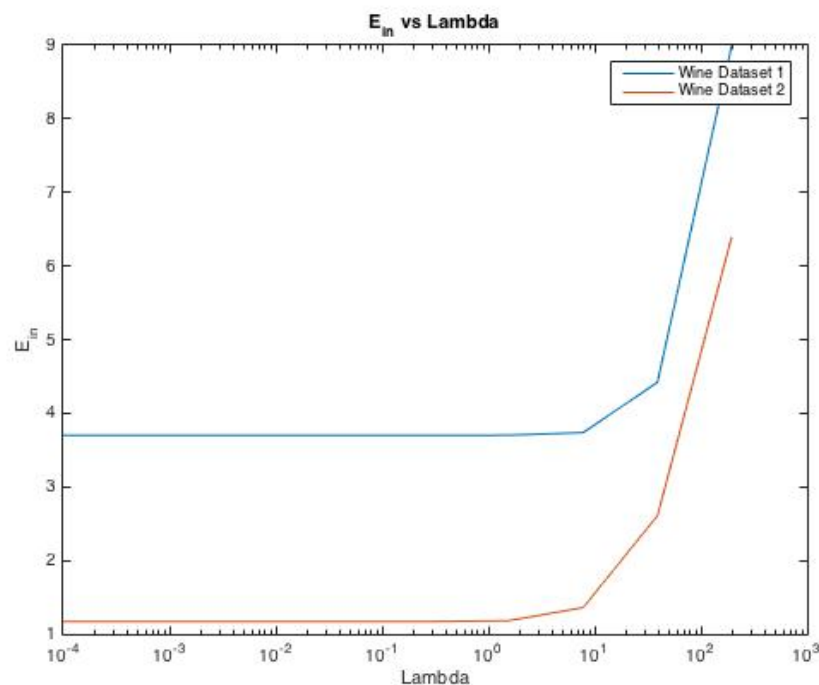
Leiya Ma
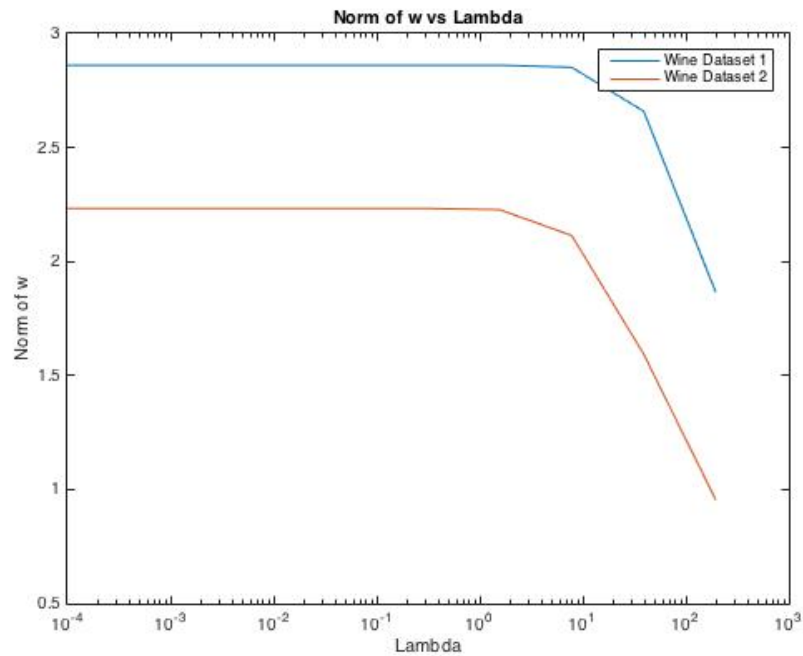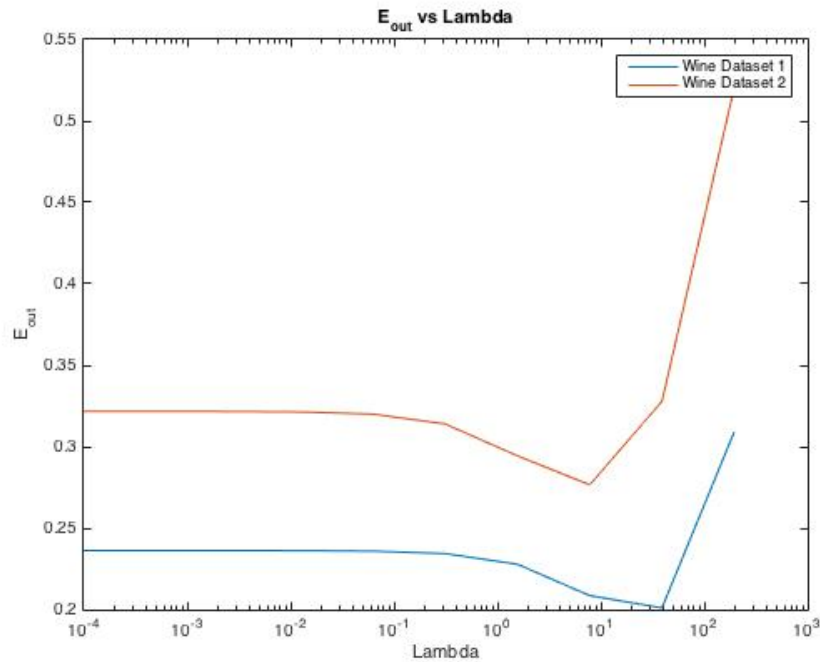
UID: 2018042

Email: lma2@caltech.edu

## 1. Effects of Regularization

- Question A:  Adding the penalty term will not decrease the training(in-sample) error. This is because adding the penalty term is equivalent to constraining the model complexity which will not bring down the in-sample error. Instead, constraining the model complexity must either increase the in-sample error or keep it unchanged.

- Question B: $\ell_0$ regulation is rarely used because it involves $\ell_0$ norm of w, which is flat or discontinuous everywhere. Therefore, it is very difficult to optimize using gradient descent algorithm.

- Question C:

E_out vs Lambda



Norm of w vs Lambda

- Question D: For small $\lambda$, $\mathrm{E_{in}}$ in wine_training2 is slightly lower than $\mathrm{E_{in}}$ in wine_training1; But for big $\lambda$, wine_training1 seems to give much lower $\mathrm{E_{in}}$. This means that more complex models are better fitting wine_training2 and simpler models are worse fitting wine_training2.

  The reason for the former statement is that the data amount in wine_training2 is much smaller than dataset1. Therefore, it is easier to more closely fit all the points with a more complex model in wine_training2.

The reason for the later statement is that when more simple and generous model are used, the dataset with fewer samples may be affected more severely and thus does not fit well.

- Question E: From the beginning to the midway of the graph, we can see overfitting. We can see this because $E_{out}$ is increasing from its local minimum when $\lambda$ is decreasing. This is overfitting because we are increasing the model complexity at the cost of accuracy when applied this model to a more generalized data.

  Starting from the midway to the end of the graph, we can see underfitting. This is illustrated by the fact as we increase the $\lambda$, $E_{out}$ and $E_{in}$ increases dramatically. When increasing $\lambda$, which means we apply simpler models to the data that it cannot capture the data accurately enough.

- Question F: We can see that the norm of w decreases as $\lambda$ increases while training with the data in wine_training1.txt. This it because when $\lambda$ increases, that means the propotion of the norm of w in the error function increases as well, and when we try to solve the optimization problem to minimize the error function, the norm of w gets more and more harmful. In summary, the larger $\lambda$ is, the stricter constraint is displayed on the norm of w.

- Question G:  if the model were trained with wine_training2.txt, I would choose $\lambda = 39.0625$. Because it has the smallest $E_{out}$.

## 2. Lasso($\ell_1$) vs. Ridge($\ell_2$) Regularization

- Question A:

  I. Using the given priors and Bayes rule, we have

  $$p(w|\mathcal{D}) = \frac{p(\mathcal{D}|w)p(w|\lambda)}{p(\mathcal{D})}$$

  we first have

  (here the $argmax$ and $argmin$ all means to find the w which maximum or minimum the function after them, I don't know how to write the w below the argmax and argmin in work equations so first state it here.)

  $$\widehat{w} = argmax\ p(w|\mathcal{D})$$

  then we can replace it using the Bayes rule, and ignore the denominator because it has nothing to do with w

  $$\widehat{w} = argmax\ p(\mathcal{D}|w)p(w|\lambda)$$

$$\widehat{w} = argmax \ p(\mathcal{D}|w) \prod_{j=1}^{\mathcal{D}} \frac{\lambda}{2} e^{-\lambda|w_j|}$$

we can replace the right side with log function because it is monotonically increasing,

$$\widehat{w} = argmax \ \log \left( p(\mathcal{D}|w) \prod_{j=1}^{\mathcal{D}} \frac{\lambda}{2} e^{-\lambda|w_j|} \right)$$

$$\widehat{w} = argmin - \log \left( p(\mathcal{D}|w) \prod_{j=1}^{\mathcal{D}} \frac{\lambda}{2} e^{-\lambda|w_j|} \right)$$

$$\widehat{w} = argmin - \left( \log p(\mathcal{D}|w) + \sum_{j=1}^{\mathcal{D}} \log \frac{\lambda}{2} e^{-\lambda|w_j|} \right)$$

$$\widehat{w} = argmin - \left( \log p(\mathcal{D}|w) + \sum_{j=1}^{\mathcal{D}} \log \frac{\lambda}{2} + \sum_{j=1}^{\mathcal{D}} \log e^{-\lambda|w_j|} \right)$$

we can ignore the second item as it has nothing to do with w,

$$\widehat{w} = argmin - \left( \log p(\mathcal{D}|w) + \sum_{j=1}^{\mathcal{D}} \log e^{-\lambda|w_j|} \right)$$

$$\widehat{w} = argmin - \left( \log p(\mathcal{D}|w) + \sum_{j=1}^{\mathcal{D}} -\lambda|w_j| \right)$$

$$\widehat{w} = argmin - \log p(\mathcal{D}|w) + \lambda \|w\|_1 )$$


II. Using the given priors and Bayes rule, we have

$$p(w|\mathcal{D}) = \frac{p(\mathcal{D}|w)p(w|\lambda)}{p(\mathcal{D})}$$

we first have

$$\widehat{w} = argmax \ p(w|\mathcal{D})$$

then we can replace it using the Bayes rule, and ignore the denominator because it has nothing to do with w

$$\widehat{w} = argmax \ p(\mathcal{D}|w)p(w|\lambda)$$

$$\widehat{w} = argmax \ p(\mathcal{D}|w) \prod_{j=1}^{\mathcal{D}} \sqrt{\frac{\lambda}{\pi}} e^{-\lambda w_j^2}$$

we can replace the right side with log function because it is monotonically increasing,

$$\hat{w} = argmax \; \log\left(p(\mathcal{D}|w)\prod_{j=1}^{\mathcal{D}}\sqrt{\frac{\lambda}{\pi}}e^{-\lambda w_j^2}\right)$$

$$\hat{w} = argmin \; -\log\left(p(\mathcal{D}|w)\prod_{j=1}^{\mathcal{D}}\sqrt{\frac{\lambda}{\pi}}e^{-\lambda w_j^2}\right)$$

$$\hat{w} = argmin \; -\left(\log p(\mathcal{D}|w) + \log\prod_{j=1}^{\mathcal{D}}\sqrt{\frac{\lambda}{\pi}}e^{-\lambda w_j^2}\right)$$

$$\hat{w} = argmin \; -\left(\log p(\mathcal{D}|w) + \sum_{j=1}^{\mathcal{D}}\log\sqrt{\frac{\lambda}{\pi}} + \sum_{j=1}^{\mathcal{D}}\log e^{-\lambda w_j^2}\right)$$

we can ignore the second item as it has nothing to do with w,

$$\hat{w} = argmin \; -\left(\log p(\mathcal{D}|w) + \sum_{j=1}^{\mathcal{D}}\log e^{-\lambda w_j^2}\right)$$

$$\hat{w} = argmin \; -\log p(\mathcal{D}|w) + \lambda\|w\|_2^2$$

III. we have

$$y \sim \mathcal{N}(Xw, \sigma^2 I), \mathcal{D} \text{ contains X and y}$$

we start with

$$\hat{w} = argmax \; p(\mathcal{D}|w)$$

The multivariate normal distribution has a density function:

$$\frac{1}{\sqrt{(2\pi)^2|\sigma^2 I|}}e^{-\frac{1}{2}(y-Xw)^T(\sigma^2 I)^{-1}(y-Xw)}$$

This density function describes the relative likelihood for generating the dataset $\mathcal{D}$ given a certain weight vector w, which is obviously equivalent to $p(\mathcal{D}|w)$.

$$\hat{w} = argmax \; \frac{1}{\sqrt{(2\pi)^2|\sigma^2 I|}}e^{-\frac{1}{2}(y-Xw)^T(\sigma^2 I)^{-1}(y-Xw)}$$

Exclude the irrelevant item,

$$\hat{w} = argmax \; e^{-\frac{1}{2}(y-Xw)^T(\sigma^2 I)^{-1}(y-Xw)}$$

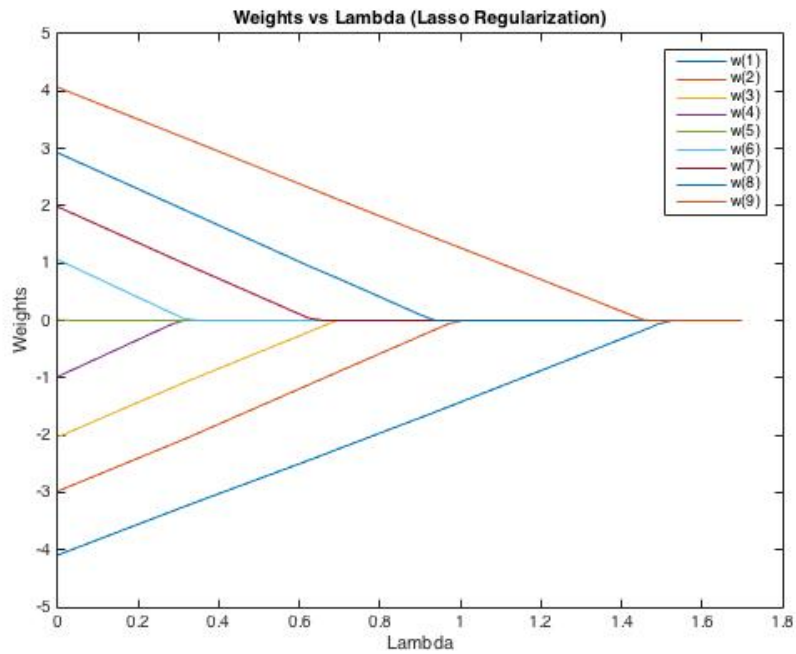$$\hat{w} = argmax \; \log\left(e^{-\frac{1}{2}(y-Xw)^T(\sigma^2 I)^{-1}(y-Xw)}\right)$$

$$\hat{w} = argmax \; \left(-\frac{1}{2}(y-Xw)^T(\sigma^2 I)^{-1}(y-Xw)\right)$$

$$\widehat{w} = argmin \ \frac{1}{2}(y - Xw)^T(y - Xw)$$
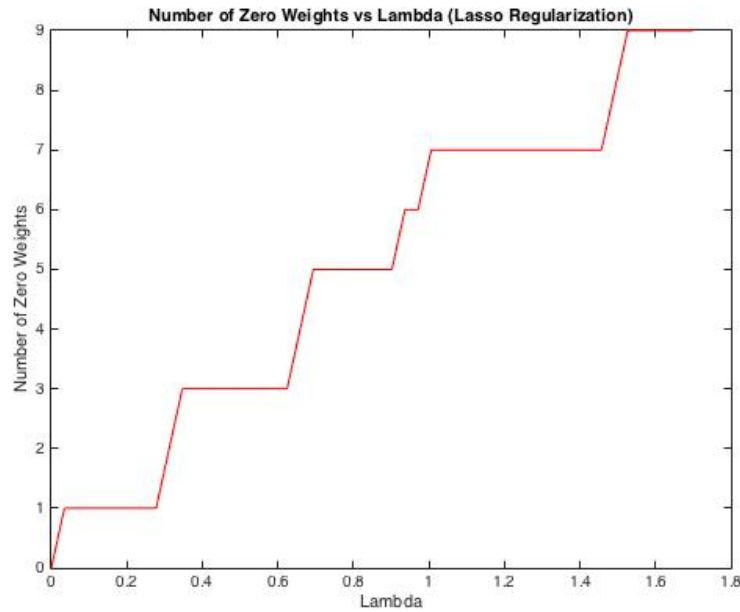
$$\widehat{w} = argmin \ \|(y - Xw)\|_2^2$$

- Question B: I.



II.

III. As regularization parameter with Lasso Regression increases, all the estimated weights go to zero (before $\lambda = 1.8$). As the regularization parameter with Ridge Regression increases, none of the estimated weights go to zero. The plot below shows this trend.



Number of Zero Weights vs Lambda (Lasso Regularization)

- Question C:

I. For $\text{argmin}\|y - w^T x\|^2 + \lambda \|w\|_1$, we want to find the w where its corresponding $\nabla f(x)$ equals 0.

$$f(x) = \|y - w^T x\|^2 + \lambda \|w\|_1$$

$$\nabla f(x) = -2X^T(y - Xw) \pm \lambda I$$

$$\nabla f(x) = -2X^T y + 2X^T Xw + \frac{\lambda w}{|w|} = 0$$

$$w = (2X^T X)^{-1}(2X^T y \mp \lambda I)$$

II. we want to explore whether there exist a value for $w_1 = 0$ and find out the smallest such value. We already got this from I.

$$w = (2X^T X)^{-1}(2X^T y \mp \lambda I)$$

therefore, we want

$$w = (2X^T y \mp \lambda I) = 0$$

now we have two situation that when w>0, $0 = 2X^T y - \lambda I \implies \lambda = 2X^T y$
when w<0, $0 = 2X^T y + \lambda I \implies \lambda = -2X^T y$; but notice that when w<0, $X^T y < 0$, therefore, our two situation can convert to a uniform one,

$$\lambda = \|2X^T y\|$$

which is the smallest value.

III. For $\text{argmin}\|y - w^T x\|^2 + \lambda\|w\|_2^2$, we want to find the w where its corresponding $\nabla f(x)$ equals 0.

$$f(x) = \|y - w^T x\|^2 + \lambda\|w\|_2^2$$
$$\nabla f(x) = 2X^T(Xw - y) + 2\lambda w = 0$$
$$2X^T Xw - 2X^T y + 2\lambda w = 0$$

therefore, we have

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

IV. when $\lambda = 0, w_i \neq 0$. If we want to find a value of $\lambda$ such that w=0,

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

we have

$$(X^T X + \lambda I)^{-1} X^T y = 0$$

$$(X^T X + \lambda I)^{-1} = 0$$

from this equation, we cannot find out the $\lambda$ such that w=0 because the inverse of some matrix cannot equal 0. Therefore, there is no such $\lambda$.