

1 Simple Bias-Variance Tradeoff

Consider a random variable X , which has unknown mean μ and unknown variance σ^2 . Given n i.i.d. realizations of training points $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ from the random variable, we wish to estimate the mean of X . We will call our estimate of X the random variable \hat{X} , which has mean $\hat{\mu}$. There are a few ways we can estimate μ given the realizations of the n samples:

1. Average the n sample points: $\frac{x_1 + x_2 + \dots + x_n}{n}$.
2. Average the n sample points and one sample point of 0: $\frac{x_1 + x_2 + \dots + x_n}{n+1}$.
3. Average the n sample points and n_0 sample points of 0: $\frac{x_1 + x_2 + \dots + x_n}{n+n_0}$.
4. Ignore the sample points: just return 0.

In the parts of this question, we will measure the *bias* and *variance* of each of our estimators. The *bias* is defined to be

$$E[\hat{X} - \mu]$$

and the *variance* is defined to be

$$\text{Var}[\hat{X}].$$

- (a) What is the bias of each of the four estimators above?
- (b) What is the variance of each of the four estimators above?
- (c) Suppose we have constructed an estimator \hat{X} from some samples of X . We now want to know how well \hat{X} estimates a fresh (new) sample of X . Denote this fresh sample by X' . Note that X' is an i.i.d. copy of the random variable X .

Derive a general expression for the expected squared error $E[(\hat{X} - X')^2]$ in terms of σ^2 and the bias and variance of the estimator \hat{X} . Similarly, derive an expression for the expected squared error $E[(\hat{X} - \mu)^2]$. Compare the two expressions and comment on the differences between them, if any.

- (d) For the following parts, we will refer to expected total error as $E[(\hat{X} - \mu)^2]$. It is a common mistake to assume that an unbiased estimator is always “best.” Let’s explore this a bit further. Compute the expected squared error for each of the estimators above.

- (e) Demonstrate that the four estimators are each just special cases of the third estimator, but with different instantiations of the hyperparameter n_0 .
- (f) What happens to bias as n_0 increases? What happens to variance as n_0 increases?
- (g) Say that $n_0 = \alpha n$. Find the setting for α that would minimize the expected total error, assuming you secretly knew μ and σ . Your answer will depend on σ , μ , and n .
- (h) For this part, let's assume that we had some reason to believe that μ *should be small* (close to 0) and σ *should be large*. In this case, what happens to the expression in the previous part?
- (i) In the previous part, we assumed there was reason to believe that μ *should be small*. Now let's assume that we have reason to believe that μ is not necessarily small, but *should be close to some fixed value μ_0* .

In terms of X and μ_0 , how can we define a new random variable X' such that X' is expected to have a small mean? Compute the mean and variance of this new random variable.

- (j) Draw a connection between α in this problem and the regularization parameter λ in the ridge-regression version of least-squares.

What does this problem suggest about choosing a regularization coefficient and handling our data-sets so that regularization is most effective? This is an open-ended question, so do not get too hung up on it.

(a) What is the bias of each of the four estimators above?

1. Average the n sample points: $\frac{x_1 + x_2 + \dots + x_n}{n}$.
2. Average the n sample points and one sample point of 0: $\frac{x_1 + x_2 + \dots + x_n}{n+1}$.
3. Average the n sample points and n_0 sample points of 0: $\frac{x_1 + x_2 + \dots + x_n}{n+n_0}$.
4. Ignore the sample points: just return 0.

$$\text{linearity: } E(aX+bY) = aE(X) + bE(Y)$$

$$\text{bias} = E(\hat{X}) - \mu$$

$$\begin{aligned} 1. \quad E(\hat{X}) &= E\left(\frac{1}{n}(X_1 + \dots + X_n)\right) = \frac{1}{n} \sum_i E(X_i) \\ &= \frac{1}{n} \sum_i \mu = \frac{1}{n} (n\mu) = \mu \\ \text{bias}(\hat{X}) &= 0 \end{aligned}$$

$$\begin{aligned} 2. \quad E(\hat{X}) &= E\left(\frac{1}{n+1}(X_1 + \dots + X_n)\right) = \frac{1}{n+1} \sum_i E(X_i) \\ &= \frac{1}{n+1} \sum_i \mu = \frac{1}{n+1} (n\mu) = \frac{n}{n+1} \mu \\ \text{bias}(\hat{X}) &= -\frac{1}{n+1} \mu \end{aligned}$$

$$\begin{aligned} 3. \quad E(\hat{X}) &= E\left(\frac{1}{n+n_0}(X_1 + \dots + X_n)\right) = \frac{1}{n+n_0} \sum_i E(X_i) \\ &= \frac{1}{n+n_0} \sum_i \mu = \frac{1}{n+n_0} (n\mu) = \frac{n}{n+n_0} \mu \\ \text{bias}(\hat{X}) &= -\frac{n_0}{n+n_0} \mu \end{aligned}$$

$$\begin{aligned} 4. \quad E(\hat{X}) &= E(0) = 0 \\ \text{bias}(\hat{X}) &= -\mu \end{aligned}$$

(b) What is the variance of each of the four estimators above?

1. Average the n sample points: $\frac{x_1 + x_2 + \dots + x_n}{n}$.
2. Average the n sample points and one sample point of 0: $\frac{x_1 + x_2 + \dots + x_n}{n+1}$.
3. Average the n sample points and n_0 sample points of 0: $\frac{x_1 + x_2 + \dots + x_n}{n+n_0}$.
4. Ignore the sample points: just return 0.

properties: $\text{Var}(\alpha X) = \alpha^2 \text{Var}(X)$
 $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$

$$\begin{aligned} 1. \text{Var}(\hat{X}) &= \text{Var}\left(\frac{1}{n}(X_1 + \dots + X_n)\right) = \frac{1}{n^2} \sum_i \text{Var}(X_i) \\ &= \frac{1}{n^2} \sum_i \sigma^2 = \frac{1}{n^2} (n\sigma^2) = \frac{1}{n} \sigma^2 \end{aligned}$$

$$\begin{aligned} 2. \text{Var}(\hat{X}) &= \text{Var}\left(\frac{1}{n+1}(X_1 + \dots + X_n)\right) = \frac{1}{(n+1)^2} \sum_i \text{Var}(X_i) \\ &= \frac{1}{(n+1)^2} \sum_i \sigma^2 = \frac{1}{(n+1)^2} (n\sigma^2) = \frac{n}{(n+1)^2} \sigma^2 \end{aligned}$$

$$\begin{aligned} 3. \text{Var}(\hat{X}) &= \text{Var}\left(\frac{1}{n+n_0}(X_1 + \dots + X_n)\right) = \frac{1}{(n+n_0)^2} \sum_i \text{Var}(X_i) \\ &= \frac{1}{(n+n_0)^2} \sum_i \sigma^2 = \frac{1}{(n+n_0)^2} (n\sigma^2) = \frac{n}{(n+n_0)^2} \sigma^2 \end{aligned}$$

$$4. \text{Var}(\hat{X}) = \text{Var}(0) = 0$$

- (c) Suppose we have constructed an estimator \hat{X} from some samples of X . We now want to know how well \hat{X} estimates a fresh (new) sample of X . Denote this fresh sample by X' . Note that X' is an i.i.d. copy of the random variable X .

Derive a general expression for the expected squared error $E[(\hat{X} - X')^2]$ in terms of σ^2 and the bias and variance of the estimator \hat{X} . Similarly, derive an expression for the expected squared error $E[(\hat{X} - \mu)^2]$. Compare the two expressions and comment on the differences between them, if any.

$$\begin{aligned}
 E[(\hat{X} - X')^2] &= E[\hat{X}^2 - 2\hat{X}X' + X'^2] \\
 &= E[\hat{X}^2] + E[X'^2] - 2E[\hat{X}X'] \\
 &= (\text{Var}(\hat{X}) + E[\hat{X}]^2) + (\text{Var}(X') + E[X'])^2 \\
 &\quad - 2E[\hat{X}X'] \\
 &= (E[\hat{X}]^2 - 2E[\hat{X}X'] + E[X']^2) \\
 &\quad + \text{Var}(\hat{X}) + \text{Var}(X') \\
 &= (E[\hat{X}] - E[X'])^2 + \text{Var}(\hat{X}) + \text{Var}(X') \\
 &= (E[\hat{X}] - \mu)^2 + \text{Var}(\hat{X}) + \sigma^2 \\
 &= \text{bias}(\hat{X})^2 + \text{Var}(\hat{X}) + \sigma^2
 \end{aligned}$$

$$\begin{aligned}
 E[(\hat{X} - \mu)^2] &= E[\hat{X}^2 - 2\mu\hat{X} + \mu^2] \\
 &= E[\hat{X}^2] - 2\mu E[\hat{X}] + \mu^2 \\
 &= (\text{Var}(\hat{X}) + E[\hat{X}]^2) - 2\mu E[\hat{X}] + \mu^2 \\
 &= (E[\hat{X}]^2 - 2\mu E[\hat{X}] + \mu^2) + \text{Var}(\hat{X}) \\
 &= (E[\hat{X}] - \mu)^2 + \text{Var}(\hat{X}) \\
 &= \text{bias}(\hat{X})^2 + \text{Var}(\hat{X})
 \end{aligned}$$

The expected squared error $E[(\hat{X} - \mu)^2]$ is the same as $E[(\bar{X} - \mu)^2]$, except it doesn't contain the term σ^2 . We refer to this term as the irreducible error b/c it comes from the natural variation in X , which cannot be changed by changing our model.

- (d) For the following parts, we will refer to expected total error as $E[(\hat{X} - \mu)^2]$. It is a common mistake to assume that an unbiased estimator is always "best." Let's explore this a bit further.

Compute the expected squared error for each of the estimators above.

$$E[(\hat{X} - \mu)^2] = \text{bias}(\hat{X})^2 + \text{var}(\hat{X})$$

$$1. E[(\hat{X} - \mu)^2] = \sigma^2 + \frac{1}{n} \sigma^2 = \frac{1}{n} \sigma^2$$

$$2. E[(\hat{X} - \mu)^2] = \left(-\frac{1}{n+1} \mu\right)^2 + \frac{n}{(n+1)^2} \sigma^2 = \frac{1}{(n+1)^2} (\mu^2 + n\sigma^2)$$

$$3. E[(\hat{X} - \mu)^2] = \left(-\frac{n_0}{n+n_0} \mu\right)^2 + \frac{n}{(n+n_0)^2} \sigma^2 = \frac{1}{(n+n_0)^2} (n_0^2 \mu^2 + n\sigma^2)$$

$$4. E[(\hat{X} - \mu)^2] = (-\mu)^2 + 0 = \mu^2$$

- (e) Demonstrate that the four estimators are each just special cases of the third estimator, but with different instantiations of the hyperparameter n_0 .

$$1. n_0 = 0 \rightarrow \frac{1}{n+n_0} (x_1 + \dots + x_n) = \frac{1}{n} (x_1 + \dots + x_n)$$

$$2. n_0 = 1 \rightarrow \frac{1}{n+n_0} (x_1 + \dots + x_n) = \frac{1}{n+1} (x_1 + \dots + x_n)$$

$$4. n_0 \rightarrow \infty \rightarrow \lim_{n_0 \rightarrow \infty} \frac{1}{n+n_0} (x_1 + \dots + x_n) = 0$$

- (f) What happens to bias as n_0 increases? What happens to variance as n_0 increases?

$$\text{bias}(\hat{x}) = -\frac{n_0}{n+n_0} \mu$$

$$\text{Var}(\hat{x}) = \frac{n}{(n+n_0)^2} \sigma^2$$

As n_0 increases, the variance of \hat{x} decreases.
 The bias of \hat{x} may inc. or dec. depending
 on the value of μ .

Therefore, there is a trade-off b/t bias &
 variance. An unbiased estimator is not
 always optimal, nor is an estimator w/
 low variance always optimal.

- (g) Say that $n_0 = \alpha n$. Find the setting for α that would minimize the expected total error, assuming you secretly knew μ and σ . Your answer will depend on σ , μ , and n .

$$\begin{aligned}
 E[(\hat{x} - \mu)^2] &= \frac{1}{(n + n_0)^2} (n_0^2 \mu^2 + n \sigma^2) \\
 &= \frac{1}{(n + \alpha n)^2} (\alpha^2 n^2 \mu^2 + n \sigma^2) \\
 &= \frac{1}{(1 + \alpha)^2} \left(\frac{1}{n^2}\right) (\alpha^2 n^2 \mu^2 + n \sigma^2) \\
 &= \frac{1}{(1 + \alpha)^2} (\alpha^2 \mu^2 + \frac{1}{n} \sigma^2) \\
 &= \frac{\alpha^2}{(1 + \alpha)^2} \mu^2 + \frac{1}{(1 + \alpha)^2} \frac{\sigma^2}{n}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial \alpha} E[(\hat{x} - \mu)^2] &= \frac{2\alpha(1 + \alpha)^2 - 2\alpha^2(1 + \alpha)}{(1 + \alpha)^4} \mu^2 - \frac{2}{(1 + \alpha)^3} \frac{\sigma^2}{n} \\
 &= \frac{2\alpha(1 + \alpha) - 2\alpha^2}{(1 + \alpha)^3} \mu^2 - \frac{2}{(1 + \alpha)^3} \frac{\sigma^2}{n} \\
 &= \frac{2\alpha}{(1 + \alpha)^3} \mu^2 - \frac{2}{(1 + \alpha)^3} \frac{\sigma^2}{n}
 \end{aligned}$$

$$\begin{aligned}
 \frac{2\hat{\alpha}}{(1 + \hat{\alpha})^3} \mu^2 - \frac{2}{(1 + \hat{\alpha})^3} \frac{\sigma^2}{n} &= 0 \\
 2\mu^2 \hat{\alpha} - \frac{2}{n} \sigma^2 &= 0 \\
 \hat{\alpha} &= \frac{\sigma^2}{n \mu^2}
 \end{aligned}$$

- (h) For this part, let's assume that we had some reason to believe that μ should be small (close to 0) and σ should be large. In this case, what happens to the expression in the previous part?

$$\hat{\mu} = \frac{\sigma^2}{n\mu^2} \quad \text{If } \mu \text{ is small + } \sigma \text{ is large,}\\ \text{then } \hat{\mu} \text{ would be large.}$$

- (i) In the previous part, we assumed there was reason to believe that μ should be small. Now let's assume that we have reason to believe that μ is not necessarily small, but should be close to some fixed value μ_0 .

In terms of X and μ_0 , how can we define a new random variable X' such that X' is expected to have a small mean? Compute the mean and variance of this new random variable.

$$X' := X - \mu_0$$

$$E(X') = E(X - \mu_0) = E(X) - \mu_0 = \mu - \mu_0 \approx 0 \quad \checkmark \mu \text{ close to } \mu_0$$

$$\text{Var}(X') = \text{Var}(X - \mu_0) = \text{Var}(X) - \text{Var}(\mu_0) = \sigma^2 - 0 = \sigma^2$$

- (j) Draw a connection between α in this problem and the regularization parameter λ in the ridge-regression version of least-squares.

What does this problem suggest about choosing a regularization coefficient and handling our data-sets so that regularization is most effective? This is an open-ended question, so do not get too hung up on it.

The α parameter acts like a regularization parameter by decreasing variance at the risk of potentially increasing the bias.

2 The Ridge Regression Estimator

Recall the ridge regression estimator for $\lambda > 0$,

$$\widehat{\theta}_\lambda := \arg \min_{\theta} \|X\theta - y\|_2^2 + \lambda \|\theta\|_2^2.$$

Let

$$X = UDV^\top = \sum_i d_i u_i v_i^\top$$

be the singular value decomposition of X . Here U and V are orthogonal matrices, meaning that $U^\top U = I$ and $V^\top V = I$. D is a diagonal matrix.

- (a) Show that the optimal weight vector $\widehat{\theta}_\lambda$ can be expressed in the form

$$\widehat{\theta}_\lambda = V\Sigma U^\top y$$

where Σ is a diagonal matrix with $\Sigma_{ii} = \frac{d_i}{d_i^2 + \lambda}$. Equivalently, we can write $\widehat{\theta}_\lambda$ as

$$\widehat{\theta}_\lambda = \sum_{i=1}^d \frac{d_i}{d_i^2 + \lambda} v_i \langle u_i, y \rangle.$$

- (b) Show that

$$\|\widehat{\theta}_\lambda\|_2^2 = \sum_{i:d_i>0} \left(\frac{d_i}{d_i^2 + \lambda} \right)^2 (u_i^\top y)^2.$$

- (c) Recall the least-norm least-squares solution is $\widehat{\theta}_{LN,LS}$ from Discussion Section 6. Show that if $\widehat{\theta}_{LN,LS} = 0$, then $\widehat{\theta}_\lambda = 0$ for all $\lambda > 0$. *Hint:* Recall that in Discussion 6 we showed that $\widehat{\theta}_{LN,LS} = \sum_{i:d_i>0} d_i^{-1} \langle u_i, y \rangle v_i$. This shows that in the case where the least-norm least square solution is zero, the ridge regression solution is also zero.
- (d) Show that if $\widehat{\theta}_{LN,LS} \neq 0$, then the function $f(\lambda) = \|\widehat{\theta}_\lambda\|_2^2$ is strictly decreasing and strictly positive on $(0, +\infty)$.
- (e) Show that

$$\lim_{\lambda \rightarrow 0^+} \widehat{\theta}_\lambda \rightarrow \widehat{\theta}_{LN,LS}.$$

Note that just because the limit of the ridge-regression objective as $\lambda \rightarrow 0^+$ is the least-squares objective, this does not immediately guarantee that the limit of the ridge solution is the least-squares solution.

- (f) In light of the above, why do you think that people describe the ridge regression as “controlling the complexity” of the solution $\widehat{\theta}_\lambda$?

Recall the ridge regression estimator for $\lambda > 0$,

$$\widehat{\theta}_\lambda := \arg \min_{\theta} \underbrace{\|X\theta - y\|_2^2 + \lambda \|\theta\|_2^2}_{f(\theta)}.$$

Let

$$X = UDV^\top = \sum_i d_i u_i v_i^\top$$

be the singular value decomposition of X . Here U and V are orthogonal matrices, meaning that $U^\top U = I$ and $V^\top V = I$. D is a diagonal matrix.

- (a) Show that the optimal weight vector $\widehat{\theta}_\lambda$ can be expressed in the form

$$\widehat{\theta}_\lambda = V\Sigma U^\top y$$

where Σ is a diagonal matrix with $\Sigma_{ii} = \frac{d_i}{d_i^2 + \lambda}$. Equivalently, we can write $\widehat{\theta}_\lambda$ as

$$\widehat{\theta}_\lambda = \sum_{i=1}^d \frac{d_i}{d_i^2 + \lambda} v_i \langle u_i, y \rangle.$$

$$\begin{aligned}\nabla_{\theta} f(\theta) &= \nabla_{\theta} \left(\theta^T X^T X \theta - 2\theta^T X^T y + y^T y + \lambda \theta^T \theta \right) \\&= 2X^T X \theta - 2X^T y + 2\lambda \theta \\&= 2((X^T X + \lambda I)\theta - X^T y)\end{aligned}$$

$$2 \left((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \hat{\boldsymbol{\theta}}_\lambda - \mathbf{X}^T \mathbf{y} \right) = 0$$

$$\hat{\Theta}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$= (V D V^T V D V^T + \lambda I)^{-1} V D V^T y$$

$$= \left(V(D^2 + \lambda I) V^T \right)^{-1} V D U^T y$$

$$= \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{V}^\top \mathbf{V} \mathbf{D} \mathbf{U} \mathbf{U}^\top \mathbf{y}$$

$$= V \underbrace{(D^2 + \lambda I)^{-1} D U^\top y}_{\Sigma}$$

$U + V$ are
orthogonal

11

(b) Show that

$$\|\widehat{\theta}_\lambda\|_2^2 = \sum_{i:d_i>0} \left(\frac{d_i}{d_i^2 + \lambda} \right)^2 (u_i^\top y)^2.$$

$$\widehat{\theta}_\lambda = V \Sigma U^\top y = \sum_{i=1}^d \left(\frac{d_i}{d_i^2 + \lambda} \right) v_i u_i^\top y$$

$$\begin{aligned} \|\widehat{\theta}_\lambda\|_2^2 &= \widehat{\theta}_\lambda^\top \widehat{\theta}_\lambda = y^\top U \Sigma V^\top V \Sigma U^\top y = y^\top U \Sigma^2 U^\top y \\ &= (y^\top U) \Sigma^2 (y^\top U)^\top \\ &= \sum_{i=1}^d \left(\frac{d_i}{d_i^2 + \lambda} \right)^2 (y^\top U)_{ii} (y^\top U)_{ii} \\ &= \sum_{i=1}^d \left(\frac{d_i}{d_i^2 + \lambda} \right)^2 (y^\top u_i)^2 \end{aligned}$$

\mathcal{D} is a diagonal matrix containing the singular values of X , so $d_i \geq 0$. Therefore,

$$\|\widehat{\theta}_\lambda\|_2^2 = \sum_{i:d_i>0} \left(\frac{d_i}{d_i^2 + \lambda} \right)^2 (y^\top u_i)^2$$

- (c) Recall the least-norm least-squares solution is $\widehat{\theta}_{LN,LS}$ from Discussion Section 6. Show that if $\widehat{\theta}_{LN,LS} = 0$, then $\widehat{\theta}_\lambda = 0$ for all $\lambda > 0$. Hint: Recall that in Discussion 6 we showed that $\widehat{\theta}_{LN,LS} = \sum_{i:d_i>0} d_i^{-1} \langle u_i, y \rangle v_i$. This shows that in the case where the least-norm least square solution is zero, the ridge regression solution is also zero.

$$\widehat{\theta}_{LN,LS} = \sum_{i:d_i>0} \frac{1}{d_i} v_i u_i^T y$$

$$\widehat{\theta}_\lambda = \sum_{i:d_i>0} \left(\frac{d_i}{d_i^2 + \lambda} \right) v_i u_i^T y$$

$$\begin{aligned} \widehat{\theta}_{LN,LS} = 0 &\Rightarrow u_i^T y = 0 \text{ for all } i \text{ s.t. } d_i > 0 \\ &\Rightarrow \widehat{\theta}_\lambda = 0 \end{aligned}$$

- (d) Show that if $\widehat{\theta}_{LN,LS} \neq 0$, then the function $f(\lambda) = \|\widehat{\theta}_\lambda\|_2^2$ is strictly decreasing and strictly positive on $(0, +\infty)$.

$$f(\lambda) := \|\widehat{\theta}_\lambda\|_2^2 = \sum_{i:d_i>0} \left(\frac{d_i}{d_i^2 + \lambda} \right)^2 (u_i^T y)^2$$

$$\widehat{\theta}_{LN,LS} \neq 0 \Rightarrow u_i^T y \neq 0 \text{ for some } i \text{ s.t. } d_i > 0$$

For $\lambda > 0$, as λ inc., $f(\lambda)$ dec. $f(\lambda)$ is a sum of a strictly positive & non-negative terms, so $f(\lambda)$ is strictly positive.

(e) Show that

$$\lim_{\lambda \rightarrow 0^+} \hat{\theta}_\lambda \rightarrow \hat{\theta}_{LN,LS}.$$

Note that just because the limit of the ridge-regression objective as $\lambda \rightarrow 0^+$ is the least-squares objective, this does not immediately guarantee that the limit of the ridge solution is the least-squares solution.

$$\begin{aligned}\hat{\theta}_\lambda &= \sum_{i:d_i > 0} \left(\frac{d_i}{d_i^2 + \lambda} \right) v_i u_i^T y \\ \lim_{\lambda \rightarrow 0^+} \hat{\theta}_\lambda &= \sum_{i:d_i > 0} \left(\lim_{\lambda \rightarrow 0^+} \frac{d_i}{d_i^2 + \lambda} \right) v_i u_i^T y \\ &= \underbrace{\sum_{i:d_i > 0} \frac{1}{d_i} v_i u_i^T y}_{\hat{\theta}_{LN,LS}}\end{aligned}$$

(f) In light of the above, why do you think that people describe the ridge regression as “controlling the complexity” of the solution $\hat{\theta}_\lambda$?

(c) $\|\hat{\theta}_\lambda\|_2$ is strictly pos. & strictly dec.

(e) $\lim_{\lambda \rightarrow 0^+} \hat{\theta}_\lambda = \hat{\theta}_{LN,LS}$

The norm of $\hat{\theta}_\lambda$ can be viewed as its complexity. Increasing λ decreases the complexity of the ridge regression sol'n compared to the LN-LS sol'n.

3 The Bias-Variance Tradeoff for Ridge Regression

Recall the statistical model for ridge regression from lecture. We have a set of sample points $\{x_i, y_i\}_{i=1}^n$ and Gaussian noise z_i . Our model follows, where the rows of X are x_i .

$$Y = Xw^* + z$$

Throughout this problem, you may assume $X^\top X$ is invertible. Recall both least-squares estimators we studied.

$$\begin{aligned} w_{\text{ols}} &= \min_{w \in \mathbb{R}^d} \|Xw - y\|_2^2 \\ w_{\text{ridge}} &= \min_{w \in \mathbb{R}^d} \|Xw - y\|_2^2 + \lambda \|w\|_2^2 \end{aligned}$$

- (a) Write the solution for $w_{\text{ols}}, w_{\text{ridge}}$. No need to derive it.
- (b) Let $\widehat{w} \in \mathbb{R}^d$ denote any estimator of w_* . In the context of this problem, an estimator $\widehat{w} = \widehat{w}(X, Y)$ is any function which takes the data X and a realization of Y , and computes a guess of w_* .

Define the MSE (mean squared error) of the estimator \widehat{w} as

$$\text{MSE}(\widehat{w}) := E\|\widehat{w} - w_*\|_2^2.$$

Above, the expectation is taken with respect to the randomness inherent in z . Define $\widehat{\mu} := E\widehat{w}$. Show that the MSE decomposes as

$$\text{MSE}(\widehat{w}) = \|\widehat{\mu} - w_*\|_2^2 + \text{Tr}(\text{Cov}(\widehat{w})).$$

Hint: Expectation and trace commute, so $E[\text{Tr}(A)] = \text{Tr}(E[A])$ for any square matrix A .

- (c) Show that

$$E[w_{\text{ols}}] = w_*, \quad E[w_{\text{ridge}}] = (X^\top X + \lambda I_d)^{-1} X^\top X w_*.$$

That is, w_{ols} is an *unbiased* estimator of w_* , whereas w_{ridge} is a *biased* estimator of w_* .

- (d) Let $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_d$ denote the d eigenvalues of the matrix $X^\top X$ arranged in non-increasing order. First, argue that the smallest eigenvalue, γ_d , is positive (i.e. $\gamma_d > 0$). Then, show that

$$\text{Tr}(\text{Cov}(w_{\text{ols}})) = \sigma^2 \sum_{i=1}^d \frac{1}{\gamma_i}, \quad \text{Tr}(\text{Cov}(w_{\text{ridge}})) = \sigma^2 \sum_{i=1}^d \frac{\gamma_i}{(\gamma_i + \lambda)^2}.$$

Finally, use these formulas to conclude that

$$\text{Tr}(\text{Cov}(w_{\text{ridge}})) < \text{Tr}(\text{Cov}(w_{\text{ols}})).$$

Hint: For the ridge variance, consider writing $X^\top X$ in terms of its eigendecomposition $U\Sigma U^\top$.

$$w_{\text{ols}} = \min_{w \in \mathbb{R}^d} \|Xw - y\|_2^2$$

$$w_{\text{ridge}} = \min_{w \in \mathbb{R}^d} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$

(a) Write the solution for $w_{\text{ols}}, w_{\text{ridge}}$. No need to derive it.

$$w_{\text{ols}} = (X^T X)^{-1} X^T y$$

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

(b) Let $\widehat{w} \in \mathbb{R}^d$ denote any estimator of w_* . In the context of this problem, an estimator $\widehat{w} = \widehat{w}(X, Y)$ is any function which takes the data X and a realization of Y , and computes a guess of w_* .

Define the MSE (mean squared error) of the estimator \widehat{w} as

$$\text{MSE}(\widehat{w}) := E\|\widehat{w} - w_*\|_2^2.$$

Above, the expectation is taken with respect to the randomness inherent in z . Define $\widehat{\mu} := E\widehat{w}$. Show that the MSE decomposes as

$$\text{MSE}(\widehat{w}) = \|\widehat{\mu} - w_*\|_2^2 + \text{Tr}(\text{Cov}(\widehat{w})).$$

Hint: Expectation and trace commute, so $E[\text{Tr}(A)] = \text{Tr}(E[A])$ for any square matrix A .

$$\begin{aligned} E(\|\widehat{w} - w^*\|_2^2) &= E(\|\widehat{w} - w^* + \widehat{\mu} - \widehat{\mu}\|_2^2) \\ &= E(\|(\widehat{w} - \widehat{\mu}) - (w^* - \widehat{\mu})\|_2^2) \\ &= E(\|\widehat{w} - \widehat{\mu}\|_2^2 - 2(\widehat{w} - \widehat{\mu})^T (w^* - \widehat{\mu}) + \|w^* - \widehat{\mu}\|_2^2) \\ &= E(\|\widehat{w} - \widehat{\mu}\|_2^2) + E(\|w^* - \widehat{\mu}\|_2^2) \\ &\quad - 2 E[(\widehat{w} - \widehat{\mu})^T (w^* - \widehat{\mu})] \end{aligned}$$

Note that b/c the expectation is taken w.r.t. the randomness inherent in z , \widehat{w} is the only r.v.

$$\begin{aligned}
E[(\hat{w} - \mu)^T (w^* - \hat{\mu})] &= E[\hat{w}^T w^* - \hat{w}^T \hat{\mu} - \hat{\mu}^T w^* - \hat{\mu}^T \hat{\mu}] \\
&= w^{*T} E[\hat{w}] - \hat{\mu}^T E[\hat{w}] - \mu^T w^* - \hat{\mu}^T \hat{\mu} \\
&= \hat{\mu}^T w^* - \hat{\mu}^T \hat{\mu} - \hat{\mu}^T w^* - \hat{\mu}^T \hat{\mu} \\
&= 0
\end{aligned}$$

$$E[\|w^* - \hat{\mu}\|_2^2] = \|\hat{\mu} - w^*\|_2^2$$

$$\begin{aligned}
E[\|\hat{w} - \hat{\mu}\|_2^2] &= E[(\hat{w} - \hat{\mu})^T (\hat{w} - \hat{\mu})] \\
&= E[\text{Tr}((\hat{w} - \hat{\mu})^T (\hat{w} - \hat{\mu}))] \\
&= E[\text{Tr}((\hat{w} - \hat{\mu})(\hat{w} - \hat{\mu})^T)] \quad \downarrow \text{hint!} \\
&= \text{Tr}(E[(\hat{w} - \hat{\mu})(\hat{w} - \hat{\mu})^T]) \\
&= \text{Tr}(\text{Cov}(\hat{w}))
\end{aligned}$$

$$E[\|\hat{w} - w^*\|_2^2] = \text{Tr}(\text{Cov}(\hat{w})) + \|\hat{\mu} - w^*\|_2^2$$

(c) Show that

$$E[w_{\text{ols}}] = w_*, \quad E[w_{\text{ridge}}] = (X^\top X + \lambda I_d)^{-1} X^\top X w_*.$$

That is, w_{ols} is an *unbiased* estimator of w_* , whereas w_{ridge} is a *biased* estimator of w_* .

$$\begin{aligned} w_{\text{ols}} &= (X^\top X)^{-1} X^\top Y \\ &= (X^\top X)^{-1} X^\top (X w^* + Z) \\ &= (X^\top X)^{-1} (X^\top X) w^* + (X^\top X)^{-1} X^\top Z \\ &= w^* + (X^\top X)^{-1} X^\top Z \end{aligned}$$

Now Z is the only r.v., so

$$E[w_{\text{ols}}] = w^* + (X^\top X)^{-1} X^\top E[Z]$$

Assuming Z is a zero-mean Gaussian r.v.,

$$E[w_{\text{ols}}] = w^* + 0 = w^*$$

$$\begin{aligned} w_{\text{ridge}} &= (X^\top X + \lambda I)^{-1} X^\top Y \\ &= (X^\top X + \lambda I)^{-1} X^\top (X w^* + Z) \\ &= (X^\top X + \lambda I)^{-1} X^\top X w^* + (X^\top X + \lambda I)^{-1} X^\top Z \end{aligned}$$

Again, Z is the only r.v., so

$$E[w_{\text{ridge}}] = (X^\top X + \lambda I)^{-1} X^\top X w^* + (X^\top X + \lambda I)^{-1} X^\top E[Z]$$

Assuming Z is a zero-mean Gaussian r.v.,

$$E[w_{\text{ridge}}] = (X^\top X + \lambda I)^{-1} X^\top X w^*$$

- (d) Let $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_d$ denote the d eigenvalues of the matrix $X^T X$ arranged in non-increasing order. First, argue that the smallest eigenvalue, γ_d , is positive (i.e. $\gamma_d > 0$). Then, show that

$$\text{Tr}(\text{Cov}(w_{\text{ols}})) = \sigma^2 \sum_{i=1}^d \frac{1}{\gamma_i}, \quad \text{Tr}(\text{Cov}(w_{\text{ridge}})) = \sigma^2 \sum_{i=1}^d \frac{\gamma_i}{(\gamma_i + \lambda)^2}.$$

Finally, use these formulas to conclude that

$$\text{Tr}(\text{Cov}(w_{\text{ridge}})) < \text{Tr}(\text{Cov}(w_{\text{ols}})).$$

Hint: For the ridge variance, consider writing $X^T X$ in terms of its eigendecomposition $U \Sigma U^T$.

$X^T X$ is necessarily a symmetric PSD matrix. B/c we assume that $X^T X$ is invertible, it is really PD. This means all its evals are strictly positive, so $\gamma_d > 0$.

$$\text{Tr}(\text{Cov}(w_{\text{ols}})) = \text{Tr} \left(E((w_{\text{ols}} - E(w_{\text{ols}}))(w_{\text{ols}} - E(w_{\text{ols}}))^T) \right)$$

$$w_{\text{ols}} = w^* + (X^T X)^{-1} X^T z$$

$$E(w_{\text{ols}}) = w^*$$

$$w_{\text{ols}} - E(w_{\text{ols}}) = (X^T X)^{-1} X^T z$$

$$\begin{aligned} \text{Tr}(\text{Cov}(w_{\text{ols}})) &= \text{Tr} \left(E \left((X^T X)^{-1} X^T z z^T X (X^T X)^{-1} \right) \right) \\ &= \text{Tr} \left((X^T X)^{-1} X^T E(z z^T) X (X^T X)^{-1} \right) \end{aligned}$$

Assuming z is composed of uncorrelated, zero-mean Gaussian r.v.s w/ variance σ^2 , $E(z z^T) = \sigma^2 I$.

$$\begin{aligned}\text{Tr}(\text{Cov}(w_{\text{ols}})) &= \text{Tr}\left(\sigma^2(X^T X)^{-1} X^T X (X^T X)^{-1}\right) \\ &= \text{Tr}\left(\sigma^2 (X^T X)^{-1}\right) \\ &= \sigma^2 \text{Tr}\left((X^T X)^{-1}\right)\end{aligned}$$

The trace of a matrix is the sum of its evals.

If $X^T X$ is a symmetric PD matrix w/ evals

$\gamma_1, \dots, \gamma_d$, then the evals of $(X^T X)^{-1}$ are $\frac{1}{\gamma_1}, \dots, \frac{1}{\gamma_d}$.

$$\text{Tr}(\text{Cov}(w_{\text{ols}})) = \sigma^2 \sum_{i=1}^d \frac{1}{\gamma_i}$$

$$\text{Tr}(\text{Cov}(w_{\text{ridge}})) = \text{Tr}\left(E((w_{\text{ridge}} - E(w_{\text{ridge}}))(w_{\text{ridge}} - E(w_{\text{ridge}}))^T)\right)$$

$$w_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T X w^+ + (X^T X + \lambda I)^{-1} X^T z$$

$$E(w_{\text{ridge}}) = (X^T X + \lambda I)^{-1} X^T X w^+$$

$$w_{\text{ridge}} - E(w_{\text{ridge}}) = (X^T X + \lambda I)^{-1} X^T z$$

$$\begin{aligned}\text{Tr}(\text{Cov}(w_{\text{ridge}})) &= \text{Tr}\left(E\left((X^T X + \lambda I)^{-1} X^T z z^T X (X^T X + \lambda I)^{-1}\right)\right) \\ &= \text{Tr}\left((X^T X + \lambda I)^{-1} X^T E(z z^T) X (X^T X + \lambda I)^{-1}\right)\end{aligned}$$

Assuming z is composed of uncorrelated, zero-mean Gaussian r.v.s w/ variance σ^2 , $E(z z^T) = \sigma^2 I$.

$$\begin{aligned}\text{Tr}(\text{Cov}(w_{\text{ridge}})) &= \text{Tr}(G^2(X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}) \\ &= G^2 \text{Tr}((X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1})\end{aligned}$$

B/c $X^T X$ is a symmetric matrix, it admits a spectral decomp. $X^T X = U \Sigma U^T$.

$$\begin{aligned}(X^T X + \lambda I)^{-1} &= (U \Sigma U^T + \lambda I)^{-1} \\ &= (U (\Sigma + \lambda I) U^T)^{-1} \\ &= U (\Sigma + \lambda I)^{-1} U^T\end{aligned}$$

$$\begin{aligned}(X^T X + \lambda I)^{-1} X^T X &= U (\Sigma + \lambda I)^{-1} U^T U \Sigma U^T \\ &= U (\Sigma + \lambda I)^{-1} \Sigma U^T\end{aligned}$$

$$\begin{aligned}(X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} &= U (\Sigma + \lambda I)^{-1} \Sigma U^T U (\Sigma + \lambda I)^{-1} U^T \\ &= U (\Sigma + \lambda I)^{-1} \Sigma (\Sigma + \lambda I)^{-1} U^T\end{aligned}$$

$$\begin{aligned}\text{Tr}(\text{Cov}(w_{\text{ridge}})) &= G^2 \text{Tr}(U (\Sigma + \lambda I)^{-1} \Sigma (\Sigma + \lambda I)^{-1} U^T) \\ &= G^2 \text{Tr}(U^T U (\Sigma + \lambda I)^{-1} \Sigma (\Sigma + \lambda I)^{-1}) \\ &= G^2 \text{Tr}((\Sigma + \lambda I)^{-1} \Sigma (\Sigma + \lambda I)^{-1}) \\ &= G^2 \text{Tr}(\Sigma (\Sigma + \lambda I)^{-2}) \\ &= G^2 \sum_{i=1}^d \frac{\gamma_i}{(\gamma_i + \lambda)^2}\end{aligned}$$

$$\text{Tr}(\text{Cov}(w_{\text{ols}})) = \sigma^2 \sum_{i=1}^d \frac{1}{Y_i}$$

$$\text{Tr}(\text{Cov}(w_{\text{ridge}})) = \sigma^2 \sum_{i=1}^d \frac{Y_i}{(Y_i + \lambda)^2}$$

If $\lambda = 0$, $\text{Tr}(\text{Cov}(w_{\text{ridge}})) = \text{Tr}(\text{Cov}(w_{\text{ols}}))$

If $\lambda > 0$, $\text{Tr}(\text{Cov}(w_{\text{ridge}})) < \text{Tr}(\text{Cov}(w_{\text{ols}}))$