

Bias-Variance Tradeoff & Model Selection

Announcements

HW5 and P5 are coming out

Recap on Bias-Variance Tradeoff

Denote $h_{\mathcal{D}}$ as the ERM solution on dataset \mathcal{D} w/ squared loss $\ell(h, x, y) = (h(x) - y)^2$

Recap on Bias-Variance Tradeoff

Denote $h_{\mathcal{D}}$ as the ERM solution on dataset \mathcal{D} w/ squared loss $\ell(h, x, y) = (h(x) - y)^2$

What we have shown is the Bias-Variance decomposition:

$$\mathbb{E}_{\mathcal{D}, x, y} (h_{\mathcal{D}}(x) - y)^2 = \mathbb{E}_{\mathcal{D}, x} (h_{\mathcal{D}}(x) - \bar{h}(x))^2 + \mathbb{E}_x (\bar{h}(x) - \bar{y}(x))^2 + \mathbb{E}_{x, y} (\bar{y}(x) - y)^2$$

$$\bar{h}(x) = \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(x)$$

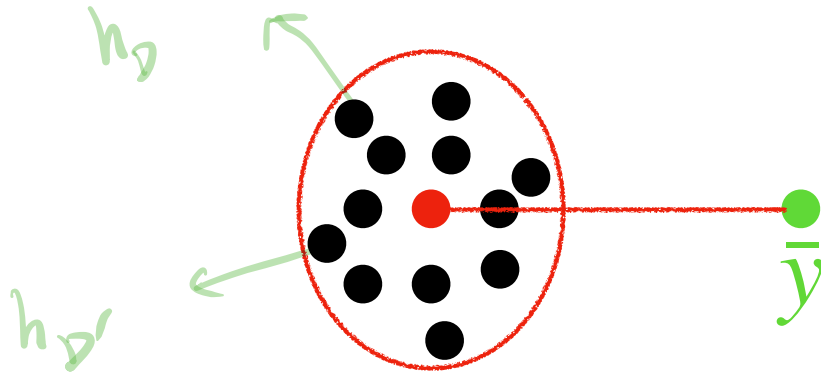
$$\bar{y}(x) = \mathbb{E}_{y|x} [y]$$

Recap on Bias-Variance Tradeoff

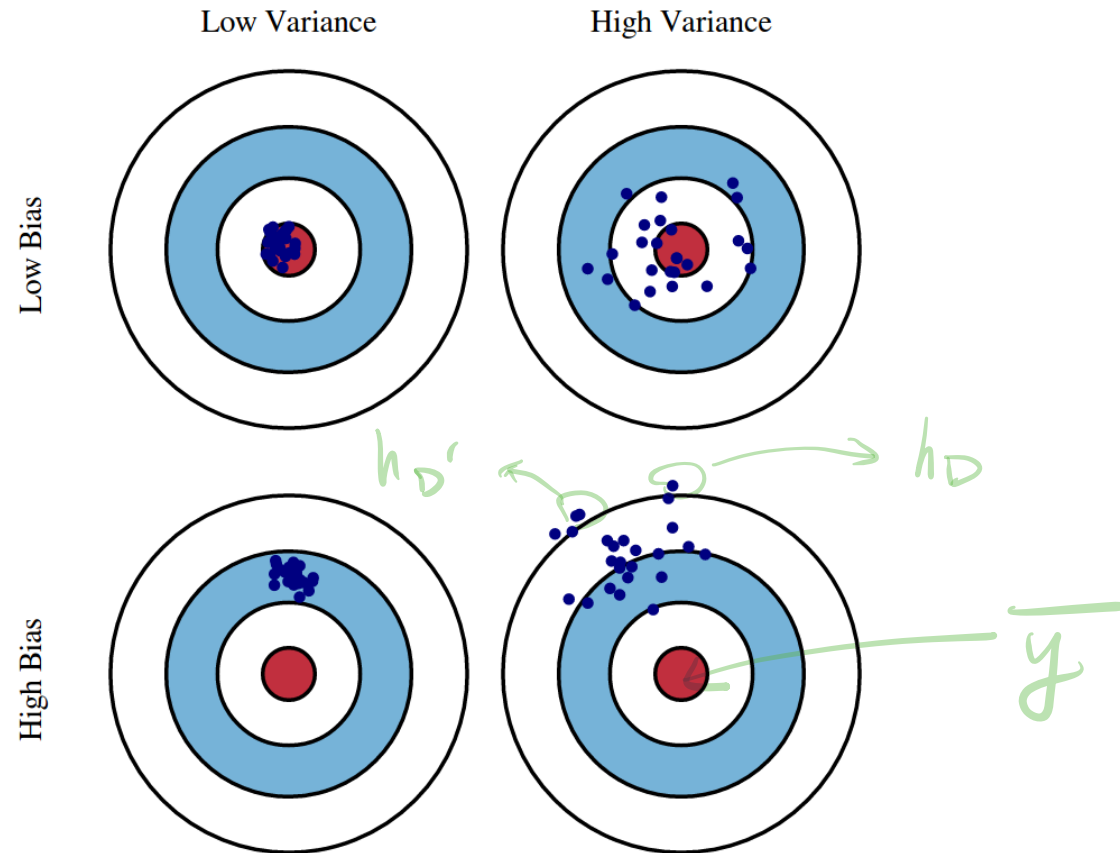
Denote $h_{\mathcal{D}}$ as the ERM solution on dataset \mathcal{D} w/ squared loss $\ell(h, x, y) = (h(x) - y)^2$

What we have shown is the Bias-Variance decomposition:

$$\mathbb{E}_{\mathcal{D}, x, y} (h_{\mathcal{D}}(x) - y)^2 = \mathbb{E}_{\mathcal{D}, x} (h_{\mathcal{D}}(x) - \bar{h}(x))^2 + \mathbb{E}_x (\bar{h}(x) - \bar{y}(x))^2 + \mathbb{E}_{x, y} (\bar{y}(x) - y)^2$$



Recap on Bias-Variance Tradeoff



Outline of Today

1. Bias & Variance tradeoff demo on Ridge Linear Regression

2. Derivation of Bias / Variance for Ridge LR

3. Model selection in practice (re-visit Cross Validation)

Ridge Linear regression w/ fixed features and Gaussian noises

Let us consider the case where features are fixed, i.e., x_1, \dots, x_n fixed (no randomness)

Ridge Linear regression w/ fixed features and Gaussian noises

Let us consider the case where features are fixed, i.e., x_1, \dots, x_n fixed (no randomness)

$$\text{But } y_i \sim (w^\star)^\top x_i + \epsilon_i, \epsilon_i \sim \mathcal{N}(0,1)$$

Ridge Linear regression w/ fixed features and Gaussian noises

Let us consider the case where features are fixed, i.e., x_1, \dots, x_n fixed (no randomness)

$$\text{But } y_i \sim (w^\star)^\top x_i + \epsilon_i, \epsilon_i \sim \mathcal{N}(0,1)$$

(This is called LR w/ fixed design)

Ridge Linear regression w/ fixed features and Gaussian noises

Let us consider the case where features are fixed, i.e., x_1, \dots, x_n fixed (no randomness)

$$\text{But } y_i \sim (w^\star)^\top x_i + \epsilon_i, \epsilon_i \sim \mathcal{N}(0,1)$$

(This is called LR w/ fixed design)

(So the only randomness of our dataset $\mathcal{D} = \{x_i, y_i\}$ is coming from the noises ϵ_i)

Ridge Linear regression

Ridge Linear Regression formulation

$$\hat{w} = \arg \min_w \sum_{i=1}^n (w^\top x_i - y_i)^2 + \lambda \|w\|_2^2$$

Ridge Linear regression

Ridge Linear Regression formulation

$$\hat{w} = \arg \min_w \sum_{i=1}^n (w^\top x_i - y_i)^2 + \lambda \|w\|_2^2$$

What we will show now:

Larger λ (model becomes “simpler”) \Rightarrow larger bias, but smaller variance

Ridge Linear regression

Ridge Linear Regression formulation

$$\hat{w} = \arg \min_w \sum_{i=1}^n (w^\top x_i - y_i)^2 + \lambda \|w\|_2^2$$

What we will show now:

Larger λ (model becomes “simpler”) \Rightarrow larger bias, but smaller variance

(Q: think about the case where $\lambda \rightarrow \infty$, what happens to \hat{w} ?)

Ridge Linear regression

Demonstration for 2d ridge linear regression

$$\mathcal{D} = \{x, y\}$$
$$y = w^*T x + \epsilon$$

1. We create 5000 datasets: $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_{5000}$,

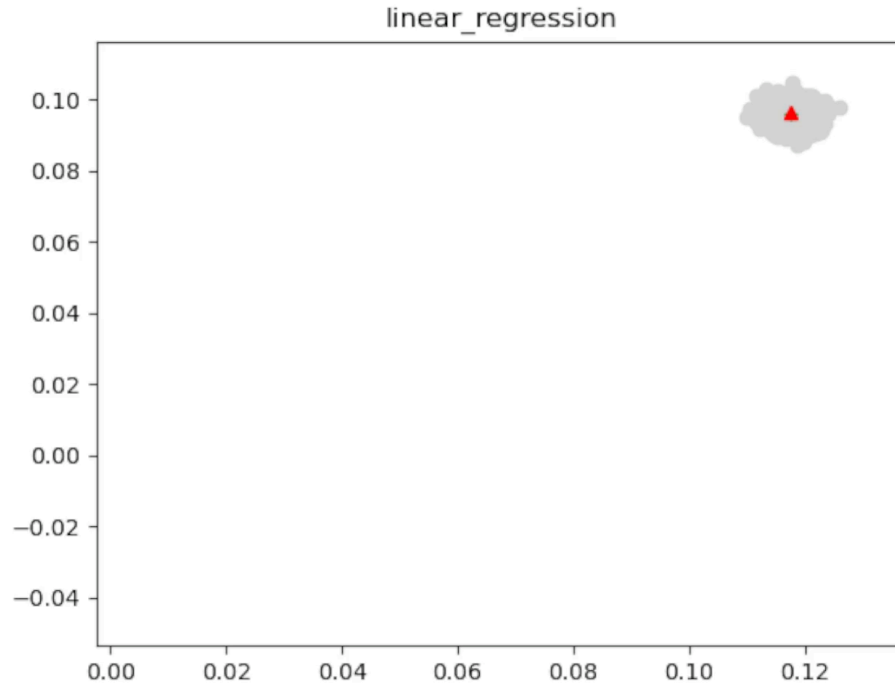
2. For a given λ , solve Ridge LR for each dataset, get $\hat{w}_1, \dots, \hat{w}_{5000}$

3. Estimate the mean $\bar{w} = \sum_i \hat{w}_i / 5000$ ✓

4. Plot $\hat{w}_1, \dots, \hat{w}_{5000}$, and mean \bar{w} , and the optimal w^*

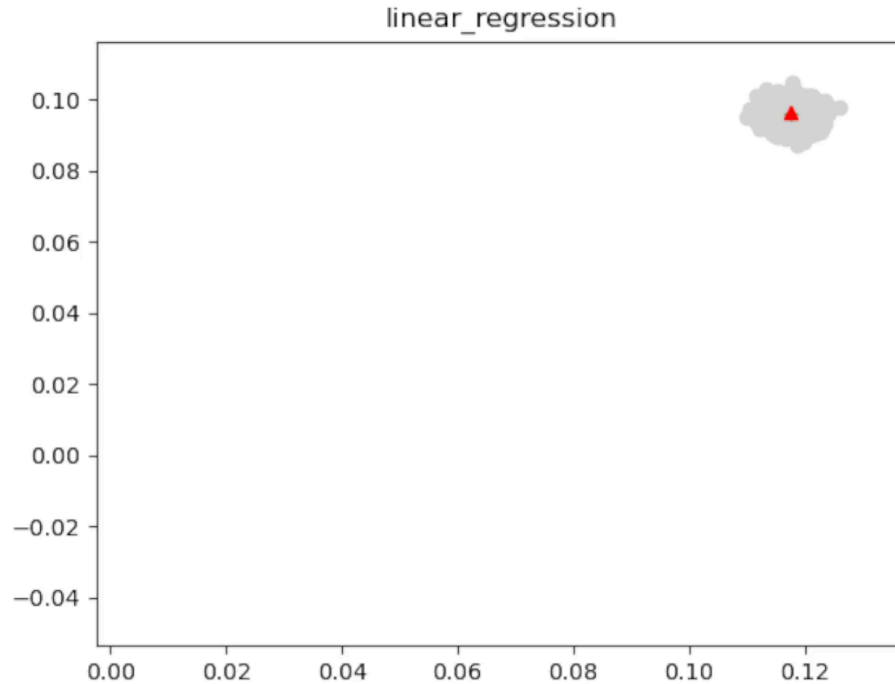
Ridge Linear regression

We start with $\lambda = 0$, and gradually increase λ to $+\infty$:



Ridge Linear regression

We start with $\lambda = 0$, and gradually increase λ to $+\infty$:



Outline of Today

1. Bias & Variance tradeoff demo on Ridge Linear Regression

2. Derivation of Bias / Variance for Ridge LR

2. Model selection in practice (re-visit Cross Validation)

Derivation of Bias and Variance for Ridge Linear regression

Denote $X = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$, $Y = [y_1, \dots, y_n]^T \in \mathbb{R}^n$, $\epsilon = [\epsilon_1, \dots, \epsilon_n]^T \in \mathbb{R}^n$

Ridge LR in matrix / vector form:

$$X = \begin{bmatrix} | & | & \dots & | \\ x_1 & x_2 & \dots & x_n \\ | & | & \dots & | \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Derivation of Bias and Variance for Ridge Linear regression

Denote $X = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$, $Y = [y_1, \dots, y_n]^T \in \mathbb{R}^n$, $\epsilon = [\epsilon_1, \dots, \epsilon_n]^T \in \mathbb{R}^n$

Ridge LR in matrix / vector form:

$$\hat{w} = \arg \min_w \|X^T w - Y\|_2^2 + \lambda \|w\|_2^2$$

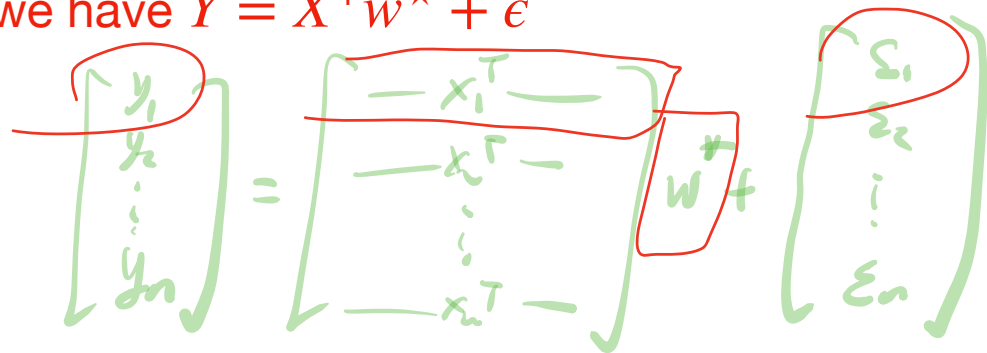
Derivation of Bias and Variance for Ridge Linear regression

Denote $X = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$, $Y = [y_1, \dots, y_n]^T \in \mathbb{R}^n$, $\epsilon = [\epsilon_1, \dots, \epsilon_n]^T \in \mathbb{R}^n$

Ridge LR in matrix / vector form:

$$\hat{w} = \arg \min_w \|X^T w - Y\|_2^2 + \lambda \|w\|_2^2$$

Since $y_i = (w^*)^T x_i + \epsilon_i$ we have $Y = X^T w^* + \epsilon$



A handwritten diagram illustrating the equation $Y = X^T w^* + \epsilon$. On the left, a vertical column vector Y is shown with elements y_1, y_2, \dots, y_n . This is equal to the product of a matrix X^T and a vector w^* . The matrix X^T is shown with rows $x_1^T, x_2^T, \dots, x_n^T$. The vector w^* is shown to the right of the matrix. Finally, the sum of the matrix product and a vector ϵ is shown on the right, with elements $\epsilon_1, \epsilon_2, \dots, \epsilon_n$. Red circles and boxes highlight the y_i terms, the x_i^T rows, and the ϵ_i terms.

$$\min_{\omega} \|X^T \omega - Y\|_2^2 + \lambda \|\omega\|_2^2$$

The Expectation of the Ridge LR solution

Recall we have closed form solution for Ridge LR

$$\hat{w} = (XX^T + \lambda I)^{-1}XY = (XX^T + \lambda I)^{-1}X(X^T w^* + \epsilon)$$

$$Y = X^T w^* + \epsilon$$

The Expectation of the Ridge LR solution

Recall we have closed form solution for Ridge LR

$$\hat{w} = (XX^T + \lambda I)^{-1}XY = (XX^T + \lambda I)^{-1}X(X^T w^* + \epsilon)$$

Source of the randomness of \hat{w}

The Expectation of the Ridge LR solution

Recall we have closed form solution for Ridge LR

$$\hat{w} = (XX^T + \lambda I)^{-1}XY = (XX^T + \lambda I)^{-1}X(X^T w^* + \epsilon)$$

Source of the randomness of \hat{w}

Let us compute the average $\bar{w} := \mathbb{E}_\epsilon[\hat{w}]$:

The Expectation of the Ridge LR solution

Recall we have closed form solution for Ridge LR

$$\hat{w} = (XX^T + \lambda I)^{-1}XY = (XX^T + \lambda I)^{-1}X(X^T w^* + \epsilon)$$

Source of the randomness of \hat{w}

Let us compute the average $\bar{w} := \mathbb{E}_\epsilon[\hat{w}]$:

$$\mathbb{E}_\epsilon[\hat{w}] = (XX^T + \lambda I)^{-1}X[\underbrace{X^T w^*}_{=0} + \mathbb{E}_\epsilon[\epsilon]]$$

The Expectation of the Ridge LR solution

Recall we have closed form solution for Ridge LR

$$\hat{w} = (XX^T + \lambda I)^{-1}XY = (XX^T + \lambda I)^{-1}X(X^T w^* + \epsilon)$$

Source of the randomness of \hat{w}

Let us compute the average $\bar{w} := \mathbb{E}_\epsilon[\hat{w}]$:

$$\mathbb{E}_\epsilon[\hat{w}] = (XX^T + \lambda I)^{-1}X[X^T w^* + \mathbb{E}_\epsilon[\epsilon]]$$

$$= (XX^T + \lambda I)^{-1}XX^T w^*$$

$$\stackrel{=0}{=} (XX^T)^{-1}XX^T = I$$

The Expectation of the Ridge LR solution

Recall we have closed form solution for Ridge LR

$$\hat{w} = (XX^T + \lambda I)^{-1}XY = (XX^T + \lambda I)^{-1}X(X^T w^* + \epsilon)$$

Source of the randomness of \hat{w}

Let us compute the average $\bar{w} := \mathbb{E}_\epsilon[\hat{w}]$:

$$\mathbb{E}_\epsilon[\hat{w}] = (XX^T + \lambda I)^{-1}X[X^T w^* + \mathbb{E}_\epsilon[\epsilon]]$$

$$= (XX^T + \lambda I)^{-1}XX^T w^*$$

$$= \underbrace{(XX^T + \lambda I)^{-1}}_{= I} \underbrace{(XX^T + \lambda I - \lambda I)}_{= I} w^*$$

The Expectation of the Ridge LR solution

Recall we have closed form solution for Ridge LR

$$\hat{w} = (XX^T + \lambda I)^{-1}XY = (XX^T + \lambda I)^{-1}X(X^T w^* + \epsilon)$$

Source of the randomness of \hat{w}

Let us compute the average $\bar{w} := \mathbb{E}_\epsilon[\hat{w}]$:

$$\mathbb{E}_\epsilon[\hat{w}] = (XX^T + \lambda I)^{-1}X[X^T w^* + \mathbb{E}_\epsilon[\epsilon]]$$

$$= (XX^T + \lambda I)^{-1}XX^T w^*$$

$$= (XX^T + \lambda I)^{-1}(XX^T + \lambda I - \lambda I)w^* = w^* - \lambda(XX^T + \lambda I)^{-1}w^*$$

The Bias of Ridge Linear regression

$$\bar{w} = \mathbb{E}[\hat{w}] = w^* - \lambda (XX^T + \lambda I)^{-1} \lambda w^*$$

Bias term: $\sum_{i=1}^n \left(\underbrace{(\bar{w} - w^*)^T}_{\text{Bias}} x_i \right)^2$

$$\bar{w} - w^* = -\lambda (XX^T + \lambda I)^{-1} \lambda w^*$$

The Bias of Ridge Linear regression

$$\bar{w} = \mathbb{E}[\hat{w}] = w^* - \lambda (XX^T + \lambda)^{-1} \lambda w^*$$

$$\text{Bias term: } \sum_{i=1}^n ((\bar{w} - w^*)^T x_i)^2$$

$$= \sum_{i=1}^n \underbrace{((\lambda(XX^T + \lambda)^{-1} w^*)^T x_i)^2}_{\bar{w} - w^*}$$

$$= \sum_{i=1}^n w^{*T} \lambda (XX^T + \lambda)^{-1} x_i x_i^T \lambda (XX^T + \lambda)^{-1} w^*$$
$$= w^{*T} \lambda (XX^T + \lambda)^{-1} \left(\sum_{i=1}^n x_i x_i^T \right) \lambda (XX^T + \lambda)^{-1} w^*$$

$\underbrace{\sum_{i=1}^n x_i x_i^T}_{= XX^T}$

The Bias of Ridge Linear regression

$$\bar{w} = \mathbb{E}[\hat{w}] = w^\star - \lambda(XX^\top + \lambda)^{-1}\lambda w^\star$$

$$\text{Bias term: } \sum_{i=1}^n \left((\bar{w} - w^\star)^\top x_i \right)^2$$

$$= \sum_{i=1}^n \left((\lambda(XX^\top + \lambda)^{-1}w^\star)^\top x_i \right)^2$$

$$= \lambda^2(w^\star)^\top (XX^\top + \lambda I)^{-1} XX^\top (XX^\top + \lambda I)^{-1} w^\star$$

The Bias of Ridge Linear regression

$$\text{Bias} = \lambda^2 (w^*)^\top (XX^\top + \lambda I)^{-1} XX^\top (XX^\top + \lambda I)^{-1} w^*$$

The Bias of Ridge Linear regression

$$\text{Bias} = \lambda^2 (w^*)^T (XX^T + \lambda I)^{-1} XX^T (XX^T + \lambda I)^{-1} w^*$$

Eigendecomposition on $XX^T = U\Sigma U^T$

$$\lambda^2 (w^*)^T U (\Sigma + \lambda I)^{-1} \underbrace{U^T U}_{XX^T} \underbrace{\Sigma U^T U}_{(XX^T + \lambda I)^{-1}} (\Sigma + \lambda I)^{-1} U^T w^* = \begin{bmatrix} | & | & | \\ u_1 & u_2 & u_d \\ | & | & | \end{bmatrix} \begin{bmatrix} \sigma_1 \\ \sigma_2 \\ \dots \\ \sigma_d \end{bmatrix} U^T$$

$$= \lambda^2 (w^*)^T U \underbrace{(\Sigma + \lambda I)^{-1} \Sigma (\Sigma + \lambda I)^{-1}}_{XX^T + \lambda I = U(\Sigma + \lambda I)U^T} U^T w^* = \begin{bmatrix} \sigma_1 + \lambda & & & \\ & \sigma_2 + \lambda & & \\ & & \dots & \\ & & & \sigma_d + \lambda \end{bmatrix}$$

The Bias of Ridge Linear regression

$$\text{Bias} = \lambda^2 (w^*)^\top (XX^\top + \lambda I)^{-1} XX^\top (XX^\top + \lambda I)^{-1} w^*$$

$$\text{Eigendecomposition on } XX^\top = U \Sigma U^\top$$

$$= (w^*)^\top U \begin{bmatrix} \frac{\sigma_1}{(\sigma_1/\lambda + 1)^2} & 0 & 0 \dots \\ 0 & \frac{\sigma_2}{(\sigma_2/\lambda + 1)^2} & 0 \dots \\ \dots & \dots & \dots \\ 0, & \dots & \frac{\sigma_d}{(\sigma_d/\lambda + 1)^2} \end{bmatrix} U^\top w^*$$

$$(\Sigma + \lambda I)^{-1} \Sigma (\Sigma + \lambda I)^{-1}$$

The Bias of Ridge Linear regression

$$\text{Bias} = \lambda^2 (w^*)^T (XX^T + \lambda I)^{-1} XX^T (XX^T + \lambda I)^{-1} w^*$$

Eigendecomposition on $XX^T = U \Sigma U^T$

when $\lambda \rightarrow +\infty \Rightarrow \frac{\sigma_i}{(\sigma_i/\lambda + 1)^2} \approx \sigma_i$

$$= (w^*)^T U \begin{bmatrix} \frac{\sigma_1}{(\sigma_1/\lambda + 1)^2} & 0 & 0 \dots \\ 0 & \frac{\sigma_2}{(\sigma_2/\lambda + 1)^2} & 0 \dots \\ \dots & \dots & \dots \\ 0, & \dots & \frac{\sigma_d}{(\sigma_d/\lambda + 1)^2} \end{bmatrix} U^T w^*$$

Q: how does bias behave when $\lambda \rightarrow +\infty$

$$\begin{aligned} \text{Bias} &= w^{*T} U \Sigma U^T w^* \\ &= w^{*T} X X^T w^* \\ &= \sum_{i=1}^d \left(w^{*T} x_i \right)^2 \end{aligned}$$

The Bias of Ridge Linear regression

$$\text{Bias} = \lambda^2 (w^*)^T (XX^T + \lambda I)^{-1} XX^T (XX^T + \lambda I)^{-1} w^*$$

Eigendecomposition on $XX^T = U\Sigma U^T$

$\rightarrow 0$ when $\lambda \rightarrow \infty \Rightarrow 0^T$

$$= (w^*)^T U \begin{bmatrix} \frac{\sigma_1}{(\sigma_1/\lambda + 1)^2} & 0 & 0 \dots \\ 0 & \frac{\sigma_2}{(\sigma_2/\lambda + 1)^2} & 0 \dots \\ \dots & \dots & \dots \\ 0, & \dots & \frac{\sigma_d}{(\sigma_d/\lambda + 1)^2} \end{bmatrix} U^T w^*$$

$\rightarrow 0$, when $\lambda \rightarrow 0$

Q: how does bias behave when $\lambda \rightarrow +\infty$

Q: how does bias behave when $\lambda \rightarrow 0$

Bias $\rightarrow 0$

The Variance of Ridge Linear regression

$$\bar{w} = \mathbb{E}[\hat{w}] = (XX^T + \lambda I)^{-1} XX^T w^*$$

$$\hat{w} = (XX^T + \lambda I)^{-1} X (X^T w^* + \varepsilon) = Y = X^T w^* + \varepsilon$$

The Variance of Ridge Linear regression

$$\bar{w} = \mathbb{E}[\hat{w}] = (XX^T + \lambda I)^{-1}XX^T w^*$$

Variance term: $\sum_{i=1}^n \mathbb{E}(\hat{w}^T x_i - \bar{w}^T x_i)^2$

↑
Random solutions

← Avg solution

$$= \sum_{i=1}^n \mathbb{E} \left((\hat{w} - \bar{w})^T x_i x_i^T (\hat{w} - \bar{w}) \right)$$

$$= \mathbb{E} \left[(\hat{w} - \bar{w})^T \sum_{i=1}^n x_i x_i^T (\hat{w} - \bar{w}) \right]$$

$$= \mathbb{E} \left[(\hat{w} - \bar{w})^T XX^T (\hat{w} - \bar{w}) \right]$$

$$= \mathbb{E} \left[(XX^T + \lambda I)^{-1} XX^T w^* + \varepsilon - (XX^T + \lambda I)^{-1} XX^T w^* \right]^T XX^T \left[(XX^T + \lambda I)^{-1} XX^T w^* + \varepsilon - (XX^T + \lambda I)^{-1} XX^T w^* \right]$$

$$= \mathbb{E} \left[\varepsilon^T XX^T \varepsilon \right] = \text{Tr} \left(XX^T \Sigma \right)$$

The Variance of Ridge Linear regression

$$\bar{w} = \mathbb{E}[\hat{w}] = (XX^T + \lambda I)^{-1}XX^T w^*$$

Variance term: $\sum_{i=1}^n \mathbb{E}(\hat{w}^T x_i - \bar{w}^T x_i)^2$

$$= \sum_{i=1}^d \sigma_i^2 / (\sigma_i + \lambda)^2$$

The Variance of Ridge Linear regression

$$\bar{w} = \mathbb{E}[\hat{w}] = (XX^T + \lambda I)^{-1}XX^T w^*$$

Variance term: $\sum_{i=1}^n \mathbb{E}(\hat{w}^T x_i - \bar{w}^T x_i)^2$

$$= \sum_{i=1}^d \sigma_i^2 / (\sigma_i + \lambda)^2$$

(Optional — tedious but basic
computation, see note)

The Variance of Ridge Linear regression

$$\bar{w} = \mathbb{E}[\hat{w}] = (XX^T + \lambda I)^{-1}XX^T w^*$$

Variance term: $\sum_{i=1}^n \mathbb{E}(\hat{w}^T x_i - \bar{w}^T x_i)^2$

$$= \sum_{i=1}^d \sigma_i^2 / (\sigma_i + \lambda)^2$$

(Optional — tedious but basic computation, see note)

Q: how does Var behave when $\lambda \rightarrow +\infty$

$$\frac{\sigma_i^2}{(\sigma_i + \lambda)^2} \rightarrow 0 \Rightarrow \text{Var} = 0$$

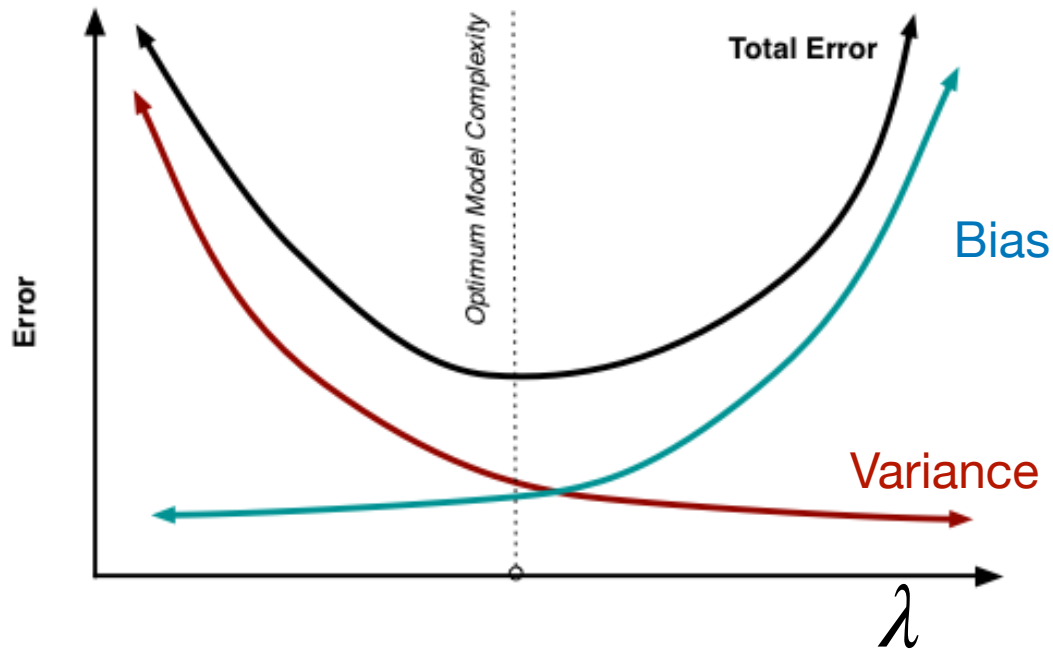
Q: how does Var behave when $\lambda \rightarrow 0$

$$\text{Var} = \sum_{i=1}^d \frac{\sigma_i^2}{\sigma_i^2} = d$$

Ridge Linear regression

Tuning λ allows us to control the generalization error of Ridge LR solution:

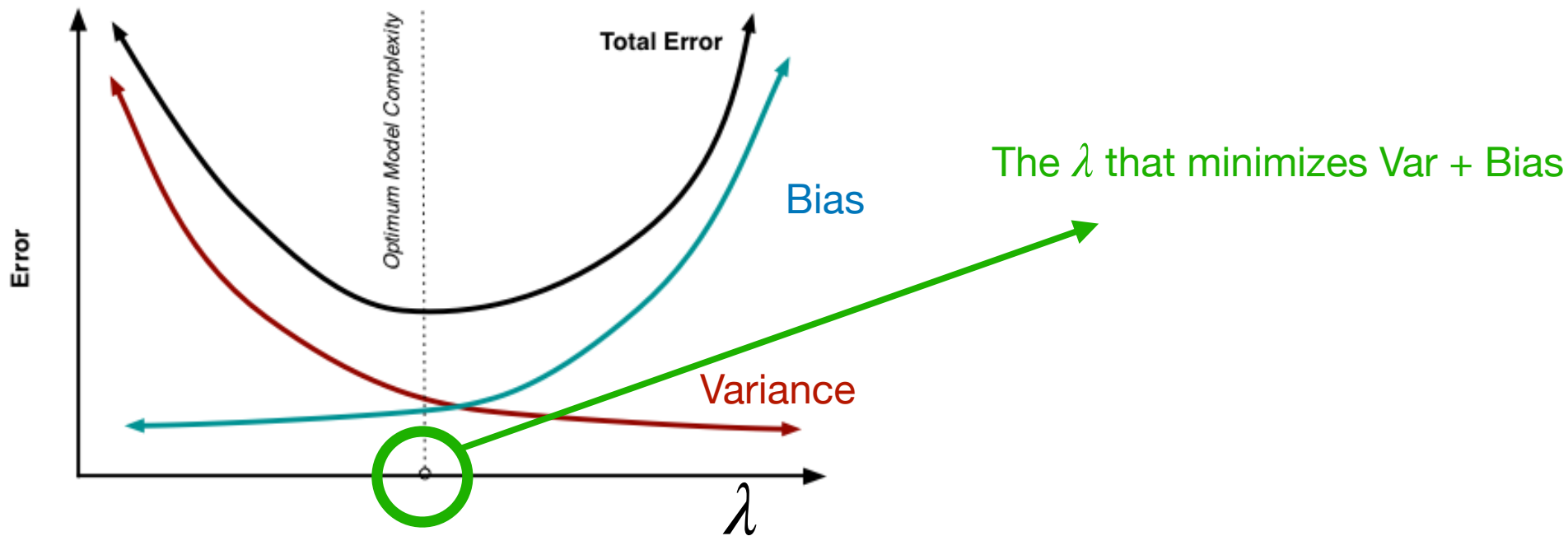
$$\mathbb{E}(\hat{w}^T x - y)^2 = \text{Variance} + \text{Bias} + \text{Inherent noise}$$



Ridge Linear regression

Tuning λ allows us to control the generalization error of Ridge LR solution:

$$\mathbb{E}(\hat{w}^T x - y)^2 = \text{Variance} + \text{Bias} + \text{Inherent noise}$$



Outline of Today

1. Bias & Variance tradeoff demo on Ridge Linear Regression

2. Derivation of Bias / Variance for Ridge LR

2. Model selection in practice (re-visit Cross Validation)

How to select the best model from data

Examples:

1. Select the right order of polynomials for regression

How to select the best model from data

Examples:

1. Select the right order of polynomials for regression
2. Select the right ridge regularization weight λ

How to select the best model from data

Examples:

1. Select the right order of polynomials for regression
2. Select the right ridge regularization weight λ
3. Select the right penalty for slack variables in soft SVM (i.e., the C parameter)

How to select the best model from data

Examples:

1. Select the right order of polynomials for regression

 2. Select the right ridge regularization weight λ

3. Select the right penalty for slack variables in soft SVM (i.e., the C parameter)

Select the right λ for Ridge LR

Cross Validation revisit:

Split the data into K folds

For $i = 1$ to K :

|

Select the right λ for Ridge LR

Cross Validation revisit:

Split the data into K folds

For $i = 1$ to K :

$$\hat{w}_{\tau} = \text{Ridge LR}(\mathcal{D}_{-i}, \lambda),$$

\mathcal{D}_{-i} as the data set that leaves out fold i

Select the right λ for Ridge LR

Cross Validation revisit:

Split the data into K folds

For $i = 1$ to K :

$$\hat{w}_k^i = \text{Ridge LR}(\mathcal{D}_{-i}, \lambda),$$

$$\epsilon_{\text{vad},k} = \sum_{x,y \in \mathcal{D}_i} (\hat{w}_i^T x - y)^2 / |\mathcal{D}_i|$$

evaluating on the hold out fold i

Select the right λ for Ridge LR

Cross Validation revisit:

Split the data into K folds

For $i = 1$ to K :

$$\hat{w}_k = \text{Ridge LR}(\mathcal{D}_{-i}, \lambda),$$

$$\epsilon_{\text{vad};k} = \sum_{x,y \in \mathcal{D}_i} (\hat{w}_i^\top x - y)^2 / |\mathcal{D}_i|$$

Output avg validation error $\bar{\epsilon}_\lambda = \sum_{i=1}^K \epsilon_{\text{vad};i} / K$

Select the right λ for Ridge LR

Cross Validation revisit:

Split the data into K folds

For $i = 1$ to K :

$$\hat{w}_i = \text{Ridge LR}(\mathcal{D}_{-i}, \lambda),$$

$$\epsilon_{\text{val};i} = \sum_{x,y \in \mathcal{D}_i} (\hat{w}_i^\top x - y)^2 / |\mathcal{D}_i|$$

$\approx \mathbb{E}_{x,y \sim P} (\hat{w}_i^\top x - y)^2$, i.e., test error of \hat{w}_i

Output avg validation error $\bar{\epsilon}_\lambda = \sum_{i=1}^K \epsilon_{\text{val};i} / K$

Select the right λ for Ridge LR

Cross Validation revisit:

Split the data into K folds

For $i = 1$ to K :

$$\hat{w}_k = \text{Ridge LR}(\mathcal{D}_{-i}, \lambda),$$

$$\epsilon_{\text{vad};k} = \sum_{x,y \in \mathcal{D}_i} (\hat{w}_i^\top x - y)^2 / |\mathcal{D}_i|$$

$$\approx \mathbb{E}_{x,y \sim P} (\hat{w}_i^\top x - y)^2, \text{ i.e., test error of } \hat{w}_i$$

$$\approx \mathbb{E}_{\mathcal{D}} \left[\mathbb{E}_{x,y \sim P} (\hat{w}_{\mathcal{D}}^\top x - y)^2 \right], \text{ i.e.,}$$

Generalization error of Ridge LR w/ λ

Output avg validation error $\bar{\epsilon}_\lambda = \sum_{i=1}^K \epsilon_{\text{vad};i} / K$

Select the right λ for Ridge LR

By numerating a set of possible $\lambda \in \mathbb{R}^+$, we select the one that has the smallest Cross-Valid error:

Select the right λ for Ridge LR

By numerating a set of possible $\lambda \in \mathbb{R}^+$, we select the one that has the smallest Cross-Valid error:

For λ in [1e-5, 1e-4, ... 1e4, 1e5]:

Select the right λ for Ridge LR

By numerating a set of possible $\lambda \in \mathbb{R}^+$, we select the one that has the smallest Cross-Valid error:

For λ in [1e-5, 1e-4, ... 1e4, 1e5]:

Split the data into K folds

For $i = 1$ to K :

$$\left| \begin{array}{l} \hat{w}_k = \text{Ridge LR}(\mathcal{D}_{-i}, \lambda), \\ \epsilon_{\text{vad};k} = \sum_{x,y \in \mathcal{D}_i} (\hat{w}_i^\top x - y)^2 / |\mathcal{D}_i| \end{array} \right.$$

Output avg validation error $\bar{\epsilon}_\lambda = \sum_{i=1}^K \epsilon_{\text{vad};i} / K$

Select the right λ for Ridge LR

By numerating a set of possible $\lambda \in \mathbb{R}^+$, we select the one that has the smallest Cross-Valid error:

For λ in [1e-5, 1e-4, ... 1e4, 1e5]:

Split the data into K folds

For $i = 1$ to K :

$$\left| \begin{array}{l} \hat{w}_k = \text{Ridge LR}(\mathcal{D}_{-i}, \lambda), \\ \epsilon_{\text{vad};k} = \sum_{x,y \in \mathcal{D}_i} (\hat{w}_i^\top x - y)^2 / |\mathcal{D}_i| \end{array} \right.$$

Output avg validation error $\bar{\epsilon}_\lambda = \sum_{i=1}^K \epsilon_{\text{vad};i} / K$

Select $\lambda^\star = \arg \min_{\lambda} \bar{\epsilon}_\lambda$

Select the right λ for Ridge LR

By numerating a set of possible $\lambda \in \mathbb{R}^+$, we select the one that has the smallest Cross-Valid error:

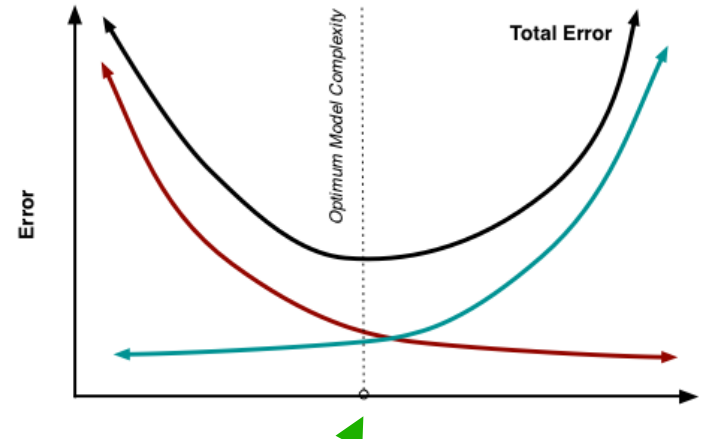
For λ in $[1e-5, 1e-4, \dots, 1e4, 1e5]$:

Split the data into K folds

For $i = 1$ to K:

$$\left| \begin{array}{l} \hat{w}_k = \text{Ridge LR}(\mathcal{D}_{-i}, \lambda), \\ \epsilon_{\text{vad};k} = \sum_{x,y \in \mathcal{D}_i} (\hat{w}_i^\top x - y)^2 / |\mathcal{D}_i| \end{array} \right.$$

Output avg validation error $\bar{\epsilon}_\lambda = \sum_{i=1}^K \epsilon_{\text{vad};i} / K$



Select $\lambda^* = \arg \min_{\lambda} \bar{\epsilon}_\lambda$

Select the right λ for Ridge LR

By numerating a set of possible $\lambda \in \mathbb{R}^+$, we select the one that has the smallest Cross-Valid error:

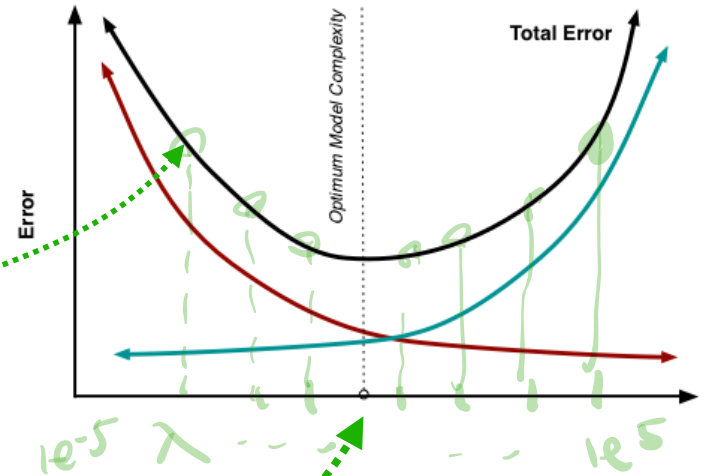
For λ in $[1e-5, 1e-4, \dots, 1e4, 1e5]$:

Split the data into K folds

For $i = 1$ to K :

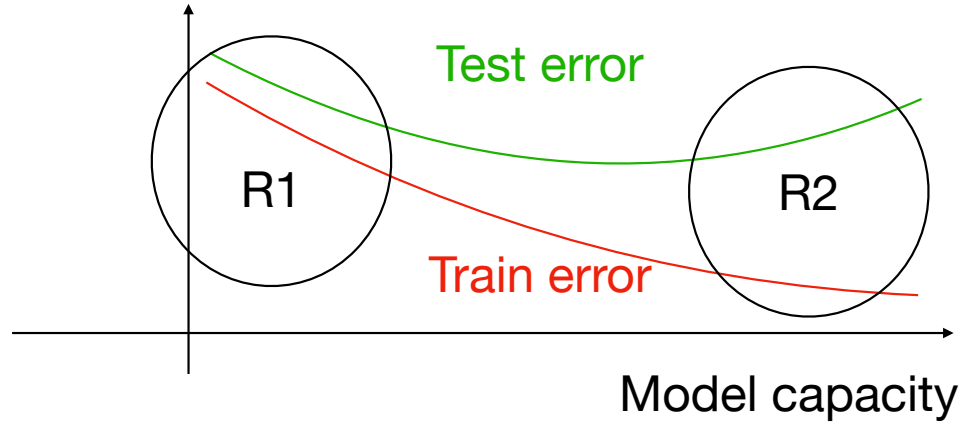
$$\begin{cases} \hat{w}_k = \text{Ridge LR}(\mathcal{D}_{-i}, \lambda), \\ \epsilon_{vad;k} = \sum_{x,y \in \mathcal{D}_i} (\hat{w}_i^\top x - y)^2 / |\mathcal{D}_i| \end{cases}$$

Output avg validation error $\bar{\epsilon}_\lambda = \sum_{i=1}^K \epsilon_{vad,i} / K$



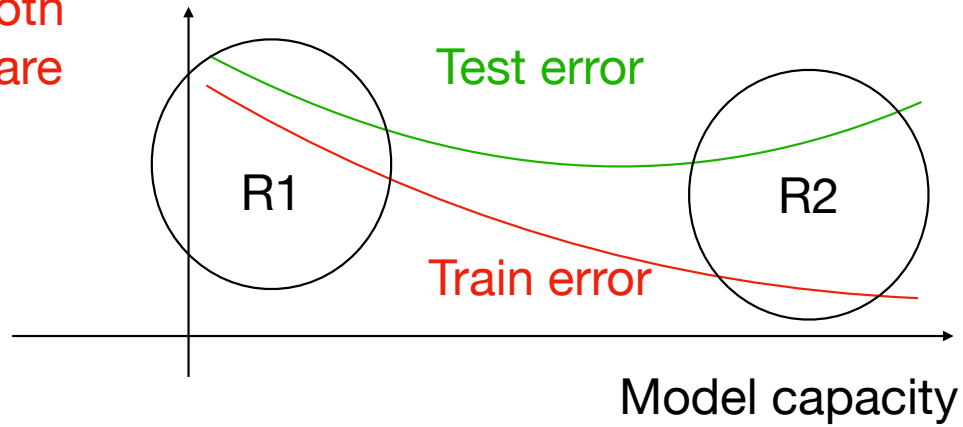
Select $\lambda^* = \arg \min_{\lambda} \bar{\epsilon}_\lambda$

Practical Suggestions for combating over/under fitting



Practical Suggestions for combating over/under fitting

R1: Underfitting (both train and test errs are large)

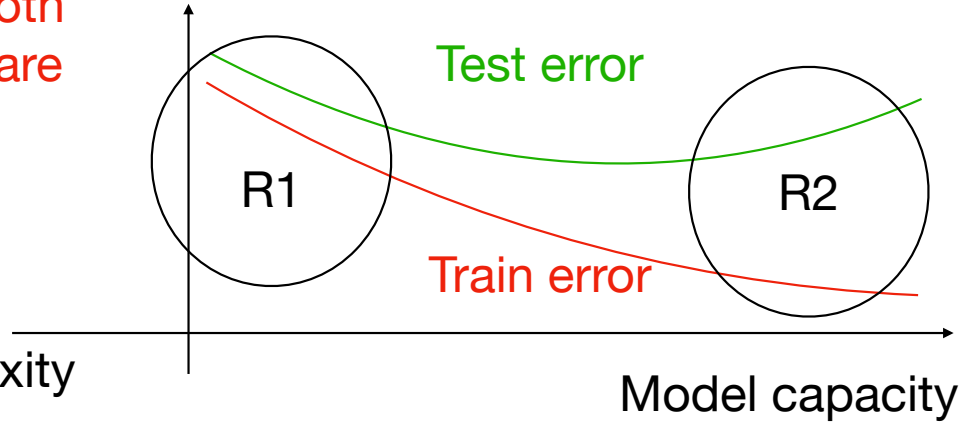


Practical Suggestions for combating over/under fitting

R1: Underfitting (both train and test errs are large)

Suggestions:

1. Increase complexity of models

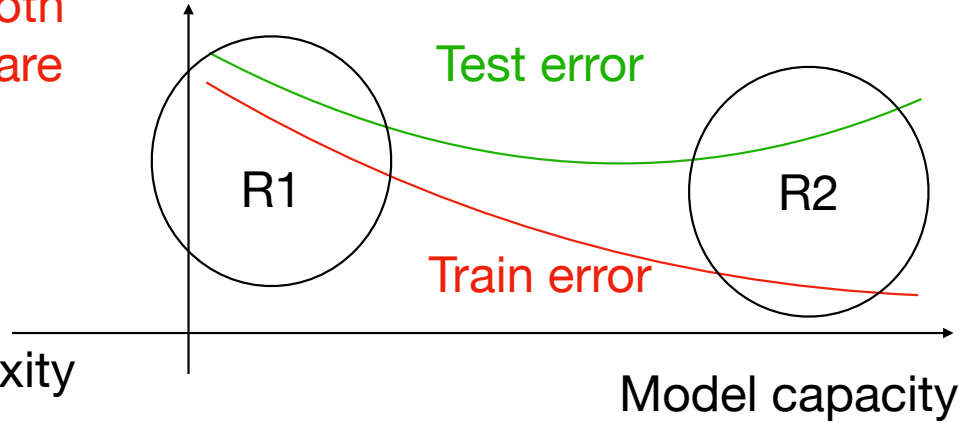


Practical Suggestions for combating over/under fitting

R1: Underfitting (both train and test errs are large)

Suggestions:

1. Increase complexity of models
2. More features

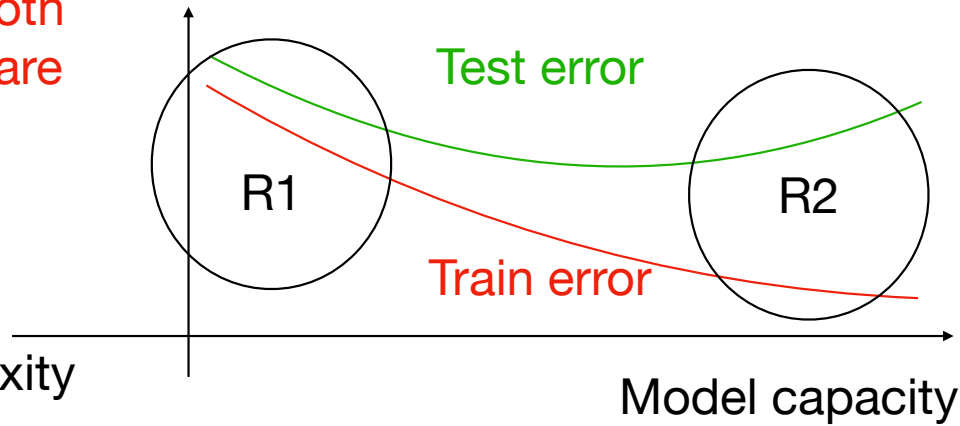


Practical Suggestions for combating over/under fitting

R1: Underfitting (both train and test errs are large)

Suggestions:

1. Increase complexity of models
2. More features
3. Using Boosting (we will see it later)

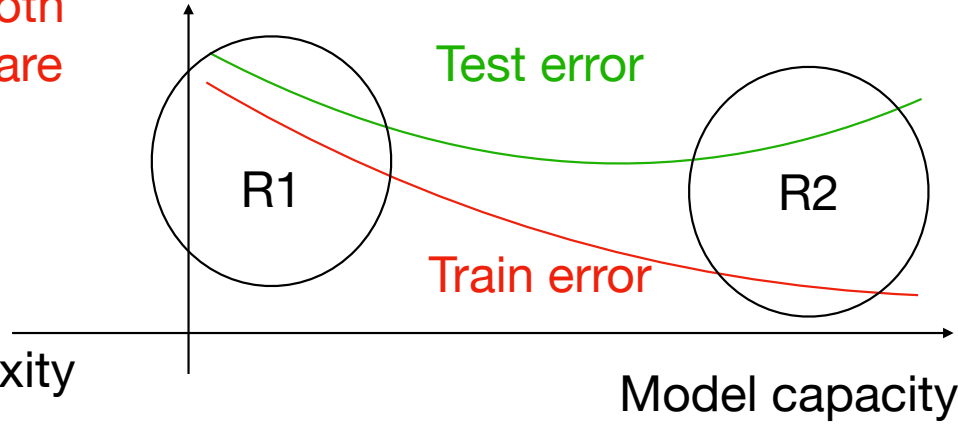


Practical Suggestions for combating over/under fitting

R1: Underfitting (both train and test errs are large)

Suggestions:

1. Increase complexity of models
2. More features
3. Using Boosting (we will see it later)



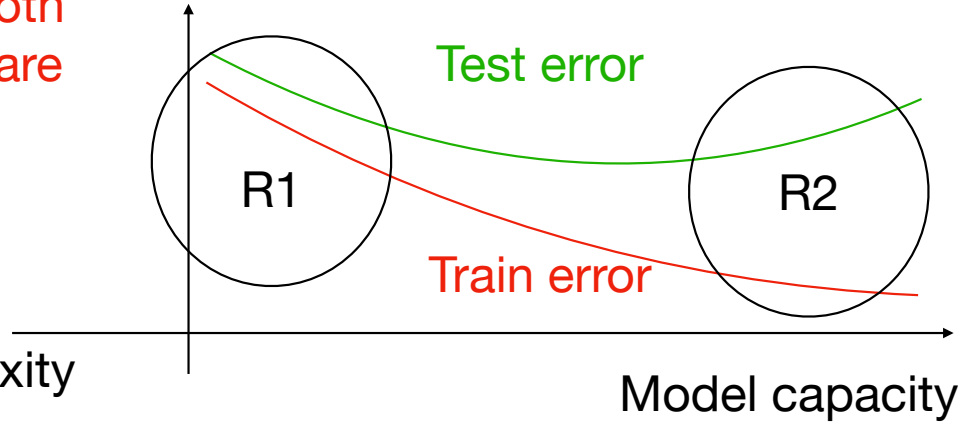
R2: overfitting (small train err but large test err)

Practical Suggestions for combating over/under fitting

R1: Underfitting (both train and test errs are large)

Suggestions:

1. Increase complexity of models
2. More features
3. Using Boosting (we will see it later)



R2: overfitting (small train err but large test err)

Suggestions:

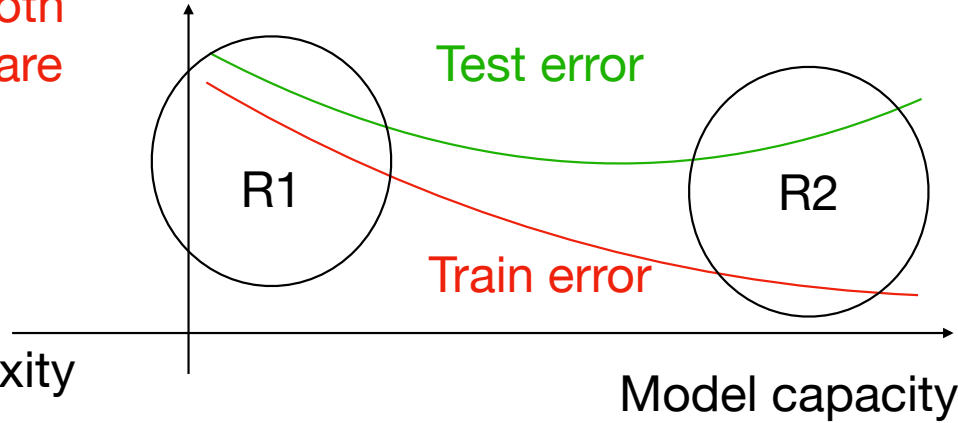
1. More train data

Practical Suggestions for combating over/under fitting

R1: Underfitting (both train and test errs are large)

Suggestions:

1. Increase complexity of models
2. More features
3. Using Boosting (we will see it later)



R2: overfitting (small train err but large test err)

Suggestions:

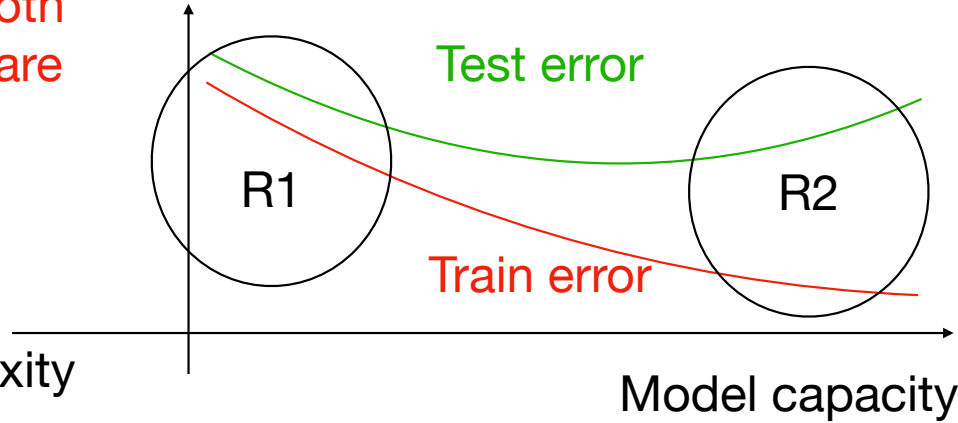
1. More train data
2. Reduce model capacity

Practical Suggestions for combating over/under fitting

R1: Underfitting (both train and test errs are large)

Suggestions:

1. Increase complexity of models
2. More features
3. Using Boosting (we will see it later)



R2: overfitting (small train err but large test err)

Suggestions:

1. More train data
2. Reduce model capacity
3. Using Bagging (we will see it later)