

- Recorded on Jan 17th Monday
- Replaces lecture on Wednesday 19th

Lecture 7: Regularization

- How to overcome overfitting.

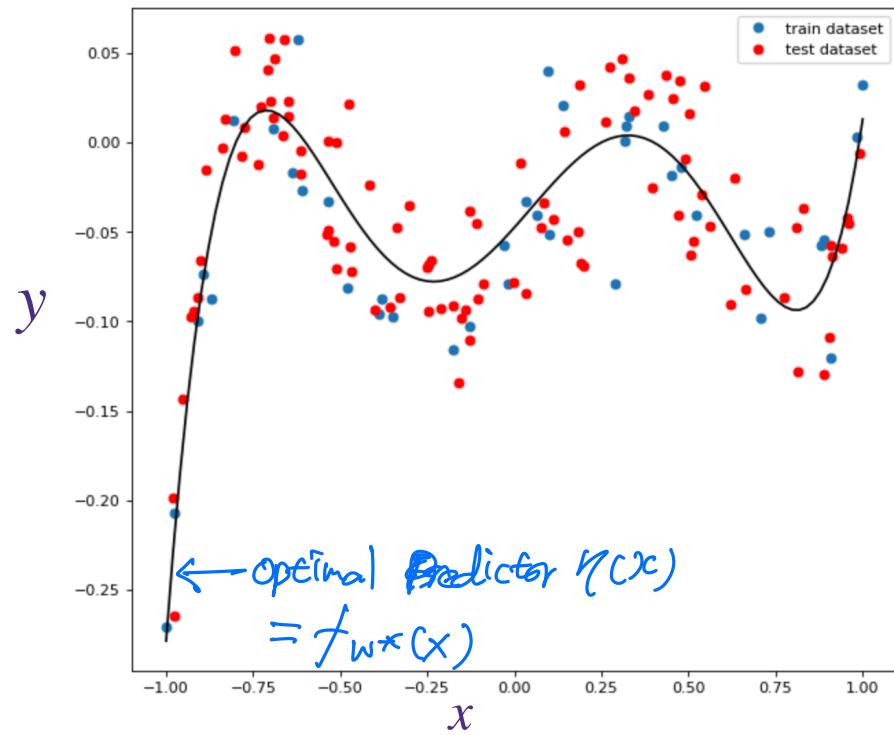
W

Recap: bias-variance tradeoff

- Consider 40 training examples and 100 test examples
i.i.d. drawn from degree-5 polynomial features

$$x_i \sim \text{Uniform}[-1, 1], y_i \sim f_{w^*}(x_i) + \epsilon_i, \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$f_{w^*}(x_i) = b^* + w_1^* x_i + w_2^* (x_i)^2 + w_3^* (x_i)^3 + w_4^* (x_i)^4 + w_5^* (x_i)^5$$



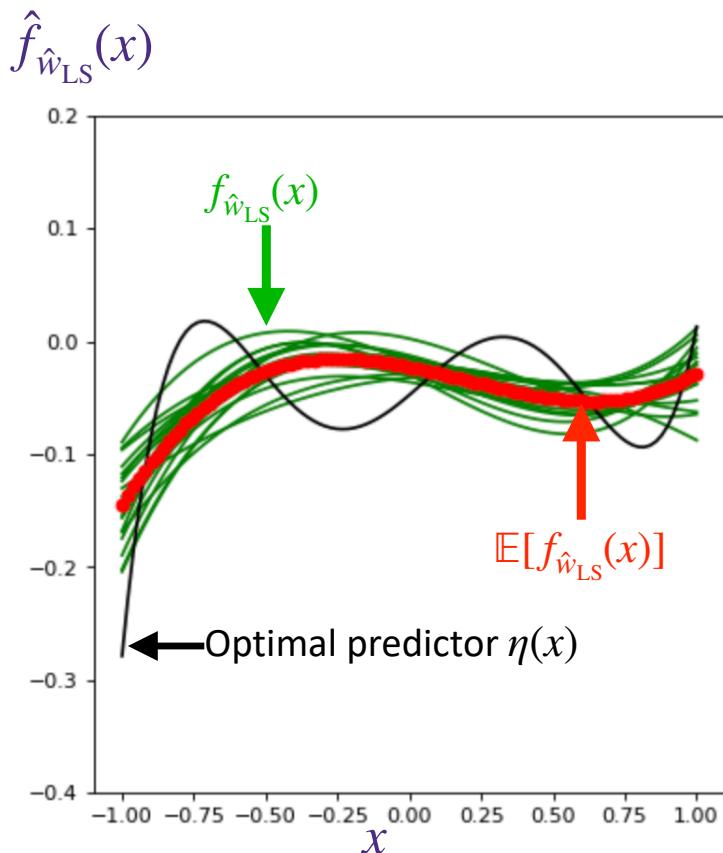
This is a linear model with features
 $h(x_i) = (x_i, (x_i)^2, (x_i)^3, (x_i)^4, (x_i)^5)$

$$f_w(x_i) = h(x_i)^T w + b$$

$$\widehat{w}_{\text{LS}} = \arg \min_{b \in \mathbb{R}, w \in \mathbb{R}^5} \sum_{i=1}^N (y_i - (h_w(x_i) + b))^2$$

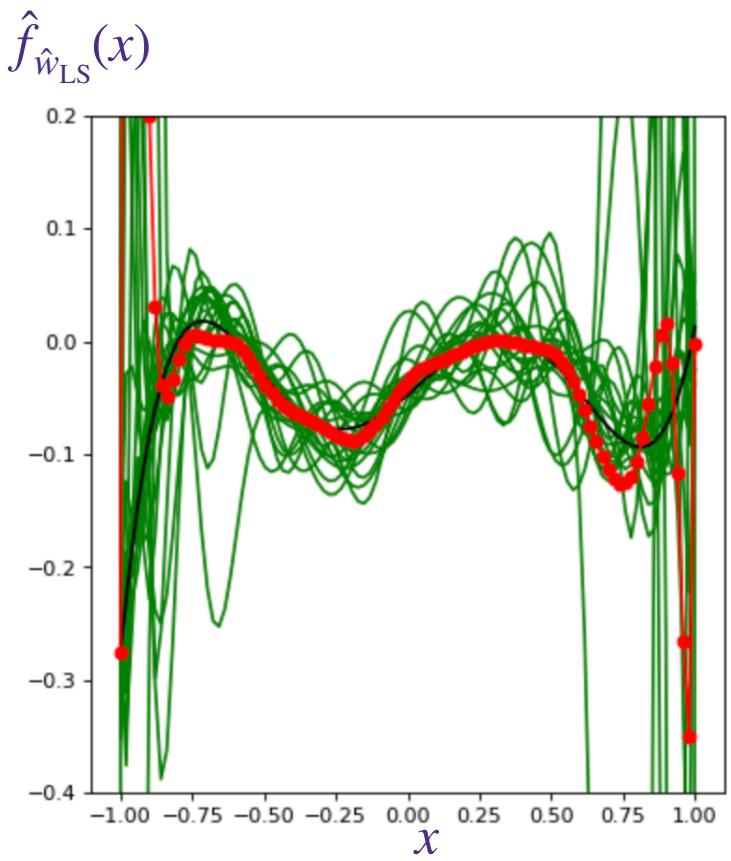
Recap: bias-variance tradeoff

With degree-3 polynomials, we underfit



```
current train error = 0.0036791644380554187  
current test error = 0.0037962529988410953
```

With degree-20 polynomials, we overfit



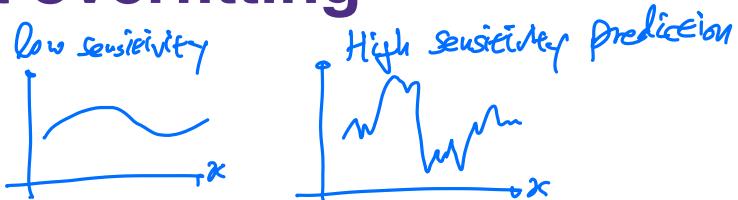
```
0.0005421686349568773  
0.14210029429557927
```

Sensitivity: how to detect overfitting

- For a linear model,

$$y \simeq b + w_1 x_1 + w_2 x_2 + \cdots + w_d x_d$$

if $|w_j|$ is large then the prediction is sensitive to small changes in x_j



- Large sensitivity leads to overfitting and poor generalization, and equivalently models that overfit tend to have large weights
- Note that b is a constant and hence there is no sensitivity for the offset b
- In **Ridge Regression**, we use a regularizer $\|w\|_2^2$ to measure and control the sensitivity of the predictor
$$\|w\|_2^2 = w_1^2 + w_2^2 + \cdots + w_d^2$$
- And optimize for small loss and small sensitivity, by adding a **regularizer** in the objective (assume no offset for now) with **regularization coefficient** $\lambda > 0$

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

- The regularization encourages solution w with smaller norm $\|w\|_2^2$, hence encouraging less overfitting
- Larger λ means more regularization
- The first term encourages fitting the training data

Minimizing the Ridge Regression Objective

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

$$\lambda(w) = \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

$$\nabla_w \lambda(w) = -2X^T(y - Xw) + 2\lambda \frac{\mathbb{I} \cdot w}{\|w\|_2}$$

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

$$= -2(X^T y - (X^T X + \lambda \mathbb{I}) \cdot w) \Big|_{w=\hat{w}_{ridge}} = 0$$

$$\hat{w}_{ridge} = (X^T X + \lambda \mathbb{I})^{-1} \cdot X^T y.$$

Shrinkage Properties

$$\widehat{w}_{\text{ridge}}^{(\lambda)} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda \|w\|_2^2$$
$$= (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T \mathbf{y} \quad , \quad \widehat{w}_{\text{LS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

For example, if $\mathbf{X}^T \mathbf{X} = n \mathbf{I}$, then

$$\widehat{w}_{\text{ridge}}^{(\lambda)} = \frac{1}{n+\lambda} \cdot \mathbf{X}^T \mathbf{y}$$
$$\widehat{w}_{\text{LS}} = \frac{1}{n} \mathbf{X}^T \mathbf{y}$$

$\widehat{w}_{\text{ridge}}^{(\lambda)} = \frac{n}{n+\lambda} \cdot \widehat{w}_{\text{LS}}$

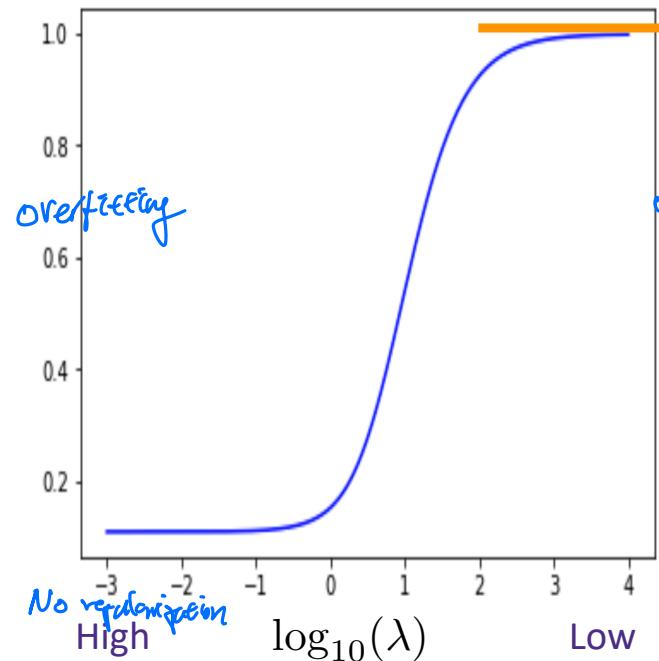
↳ shrinks \widehat{w}_{LS} .
Larger λ shrinks more.

Similar shrinking effect for general $\mathbf{X}^T \mathbf{X}$, which we do not go into details in class (come to my OH if interested).

- When $\lambda = 0$, this gives the least squares model
- This defines a family of models hyper-parametrized by λ
- Large λ means more regularization and simpler model
- Small λ means less regularization and more complex model

Ridge regression: minimize $\sum_{i=1}^n (w^T x_i + b - y_i)^2 + \lambda \|w\|_2^2$

training MSE $\frac{1}{n} \sum_{i=1}^n (y_i - (x_i^T \hat{w}_{\text{ridge}}^{(\lambda)} + \hat{b}_{\text{ridge}}))^2$



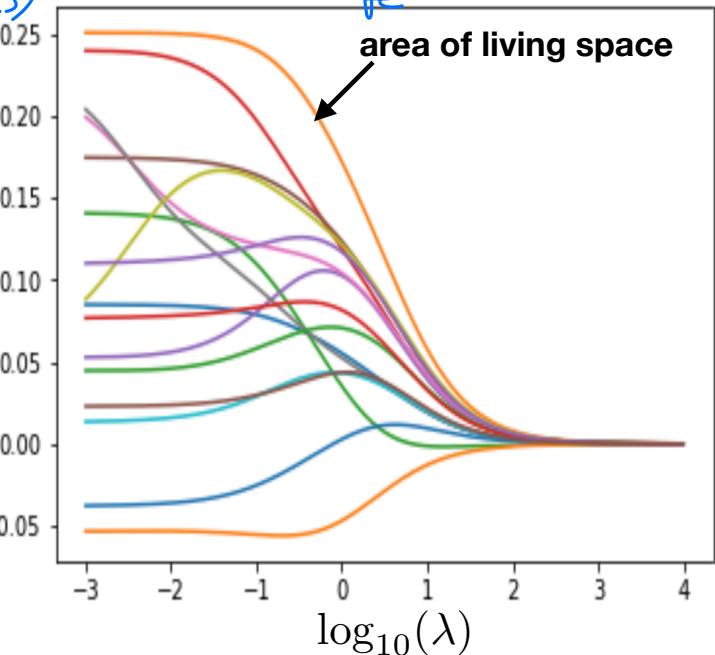
variance (Σx_i^2)

$$\lambda^{(\infty)} = \frac{n}{\sum y_i^2}$$

$$\lambda = 0$$

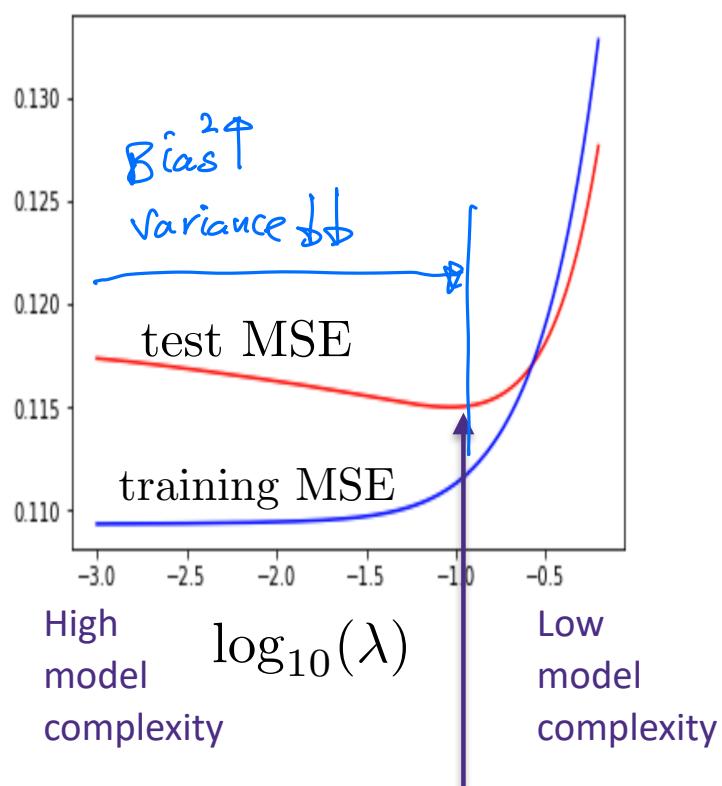
$$\hat{w}_{\text{ridge}}^{(0)} = 0$$

$$\|w\|_2^2 = 0$$

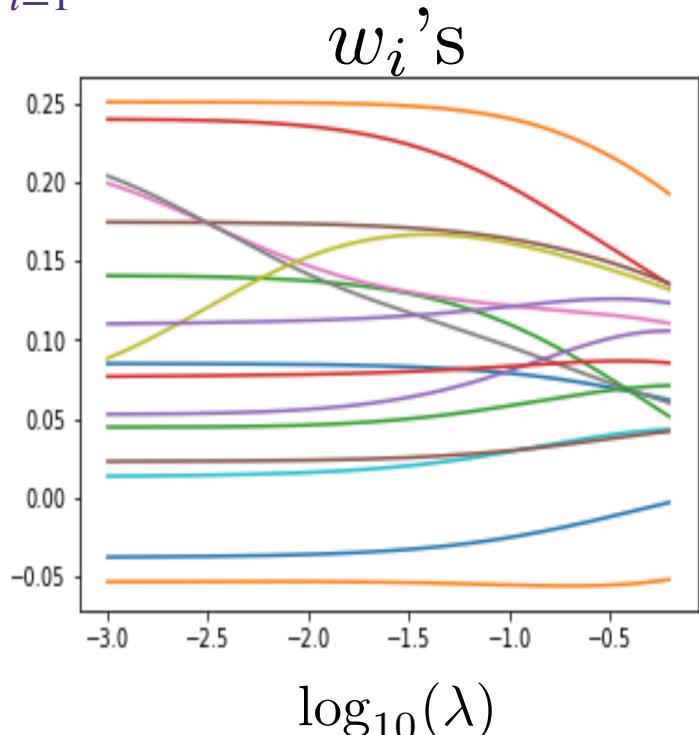


- Left plot: leftmost training error is with no regularization: 0.1093
- Left plot: rightmost training error is variance of the training data: 0.9991
- Right plot: called **regularization path**

Ridge regression: minimize $\sum_{i=1}^n (w^T x_i + b - y_i)^2 + \lambda \|w\|_2^2$



- this gain in test MSE comes from shrinking w's to get a less sensitive predictor
(which in turn reduces the variance)



Bias-variance tradeoff for ridge regression model

If $Y_i = X_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\mathbf{y} = \mathbf{X}w^* + \boldsymbol{\epsilon}$$

$$\underline{\eta(x)} = \mathbb{E}_{Y|X}[Y | X = x] = x^T w^*$$

$$\hat{w}_{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} = \underbrace{(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1}}_{=} \underbrace{(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T w^* + \boldsymbol{\epsilon})}_{= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} ((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) w^* - \lambda \mathbf{I} w^* + \mathbf{X}^T \boldsymbol{\epsilon})}$$

For example, if $\mathbf{X}^T \mathbf{X} = n \mathbf{I}$, then

$$= \frac{w^* - (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \cdot \lambda \cdot w^*}{\lambda \uparrow} + \frac{(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}}{\rightarrow \text{Variance}}$$

$$\hat{w}_{\text{ridge}} = w^* - \frac{\lambda}{n+\lambda} w^* + \frac{1}{n+\lambda} \mathbf{X}^T \boldsymbol{\epsilon}$$

H

Bias-variance tradeoff for ridge regression model

If $Y_i = \mathbf{X}_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\mathbf{y} = \mathbf{X}w^* + \boldsymbol{\epsilon}$$

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X=x] = x^T w^*$$

$$\begin{aligned}\hat{w}_{\text{ridge}} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T (\mathbf{X}w^* + \boldsymbol{\epsilon}) \\ &= w^* - (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \lambda w^* + (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}\end{aligned}$$

For example, if $\mathbf{X}^T \mathbf{X} = n \mathbf{I}$, then $\underbrace{(1 - \frac{\lambda}{n+\lambda}) w^*}_{\text{estimate is shrunk by regularizer}}$

$$\hat{w}_{\text{ridge}} = \underbrace{w^* - \frac{\lambda}{n+\lambda} w^*}_{\text{estimate is shrunk by regularizer}} + \underbrace{\frac{1}{n+\lambda} \mathbf{X}^T \boldsymbol{\epsilon}}_{\text{error due to noise}}$$

→ larger λ increases bias → larger λ decreases variance

Bias-variance tradeoff for ridge regression model

If $Y_i = X_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\mathbf{y} = \mathbf{X}w^* + \boldsymbol{\epsilon}$$

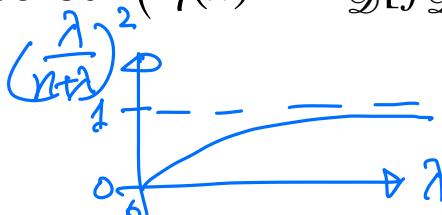
$$\eta(x) = \mathbb{E}_{Y|X}[Y|X=x] = x^T w^*$$

For example, if $\mathbf{X}^T \mathbf{X} = n\mathbf{I}$, then

$$\hat{f}_{\mathcal{D}}(x) = \boxed{x^T w^* - \frac{\lambda}{n+\lambda} x^T w^* + \frac{1}{n+\lambda} x^T \mathbf{X}^T \epsilon}$$

$\mathbb{E} \rightarrow x^T w^* - \frac{1}{n+\lambda} x^T w^*$
 $\mathbb{E}(\epsilon) = 0$

- Irreducible error: $\mathbb{E}_{X,Y}[(Y - \eta(x))^2 | X = x] = \sigma^2$
- Bias squared: $(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2 = (\cancel{x^T w^*} - \cancel{\frac{\lambda}{n+\lambda} x^T w^*})^2 = \frac{\lambda^2}{(n+\lambda)^2} \cdot (x^T w^*)^2$



$$= \boxed{\frac{\lambda^2}{(n+\lambda)^2}} (x^T w^*)^2$$

- ① Bias decreases with increasing sample size
- ② Bias increases with increasing λ

Bias-variance tradeoff for ridge regression model

If $Y_i = X_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\mathbf{y} = \mathbf{X}w^* + \boldsymbol{\epsilon}$$

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X=x] = x^T w^*$$

For example, if $\mathbf{X}^T \mathbf{X} = n\mathbf{I}$, then

$$\hat{f}_{\mathcal{D}}(x) = x^T w^* - \frac{\lambda}{n+\lambda} x^T w^* + \frac{1}{n+\lambda} x^T \mathbf{X}^T \boldsymbol{\epsilon}$$

• Variance:
$$\begin{aligned} \mathbb{E}_{\mathcal{D}}[(\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2] &= \mathbb{E}\left[\frac{1}{(n+\lambda)^2} x^T \mathbf{X}^T \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T \mathbf{X} x\right] \\ &= \frac{\sigma^2}{(n+\lambda)^2} \cdot x^T \mathbb{E}[\mathbf{X}^T \mathbf{X}] x \\ &= \frac{\sigma^2 n}{(n+\lambda)^2} \|x\|_2^2 \end{aligned}$$

- Variance decreases with increasing sample size
- Variance decrease with increasing λ

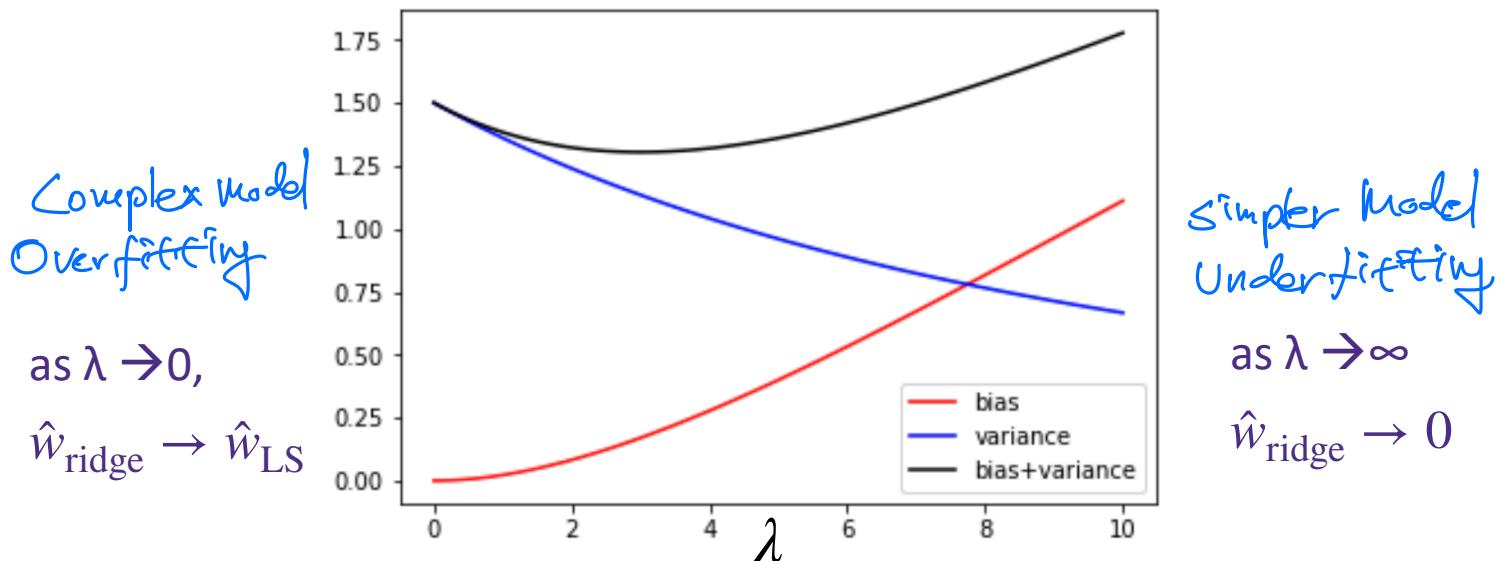
Bias-Variance Properties

- Ridge regressor: $\hat{w}_{ridge} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda \|w\|_2^2$
- True error

$$\mathbb{E}_{Y|X,\mathcal{D}}[(y - x^T \hat{w}_{ridge})^2 | x] = \sigma^2 + \frac{\lambda^2}{(n + \lambda)^2} (w^T x)^2 + \frac{\sigma^2 n}{(n + \lambda)^2} \|x\|_2^2$$

Bias-squared Variance

$$d=10, n=20, \sigma^2 = 3.0, \|w\|_2^2 = 10$$



What you need to know...

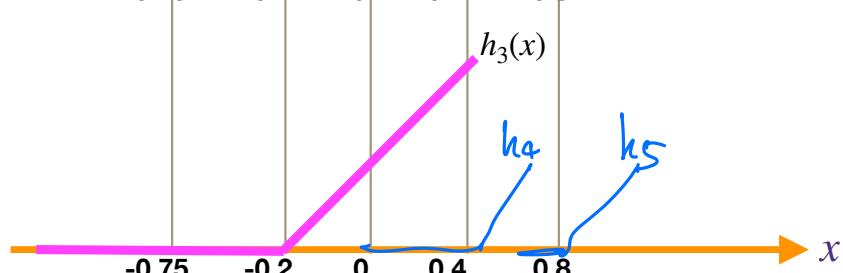
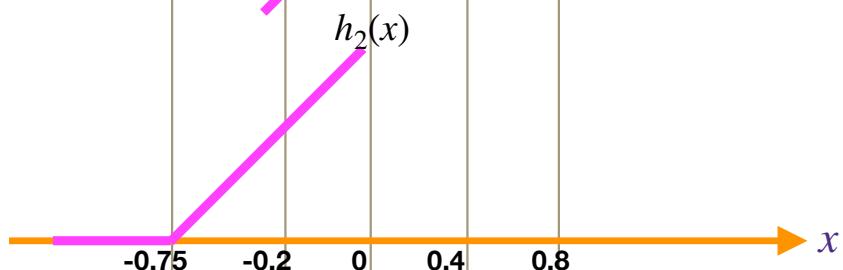
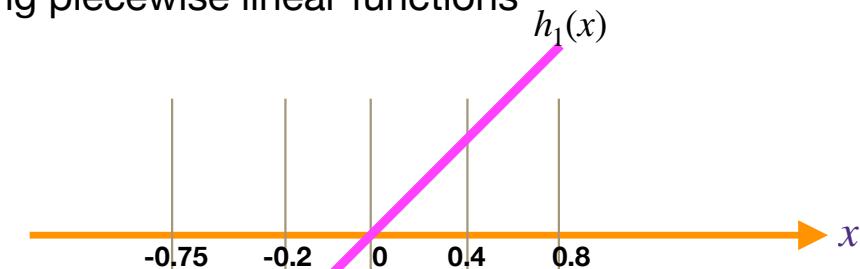
- > Regularization
 - Penalizes complex models towards simpler models
- > Ridge regression
 - L₂ penalized least-squares regression $\lambda \cdot \|\omega\|_2^2$
 - Regularization parameter trades off model complexity with training error
 - Never regularize the offset!

Example: piecewise linear fit

- we fit a linear model for $x \in [-1,1]$:
 $f(x) = b + w_1 h_1(x) + w_2 h_2(x) + w_3 h_3(x) + w_4 h_4(x) + w_5 h_5(x)$
- with a specific choice of features using piecewise linear functions

$$h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ h_3(x) \\ h_4(x) \\ h_5(x) \end{bmatrix} = \begin{bmatrix} x \\ [x + 0.75]^+ \\ [x + 0.2]^+ \\ [x - 0.4]^+ \\ [x - 0.8]^+ \end{bmatrix}$$

$$[a]^+ \triangleq \max\{a, 0\}$$

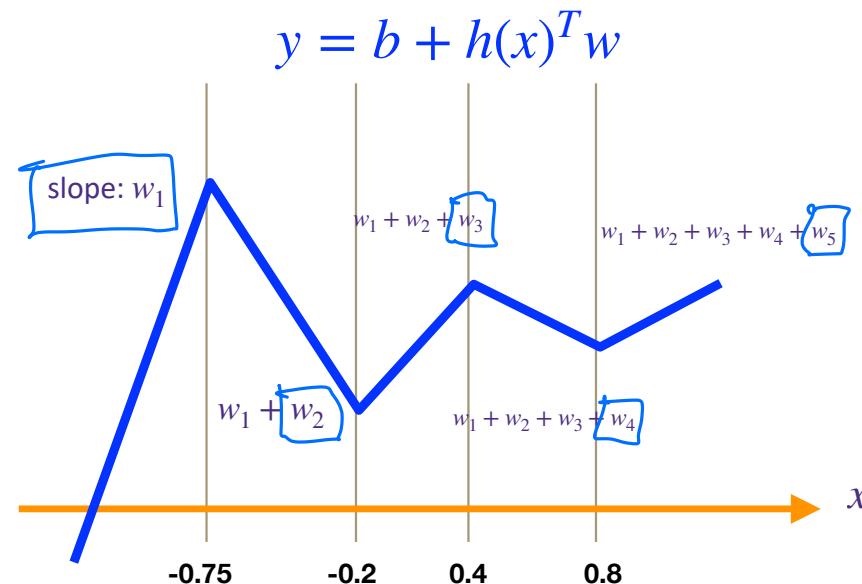


Example: piecewise linear fit

- we fit a linear model:
 $f(x) = b + w_1 h_1(x) + w_2 h_2(x) + w_3 h_3(x) + w_4 h_4(x) + w_5 h_5(x)$
- with a specific choice of features using piecewise linear functions

$$h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ h_3(x) \\ h_4(x) \\ h_5(x) \end{bmatrix} = \begin{bmatrix} x \\ [x + 0.75]^+ \\ [x + 0.2]^+ \\ [x - 0.4]^+ \\ [x - 0.8]^+ \end{bmatrix}$$

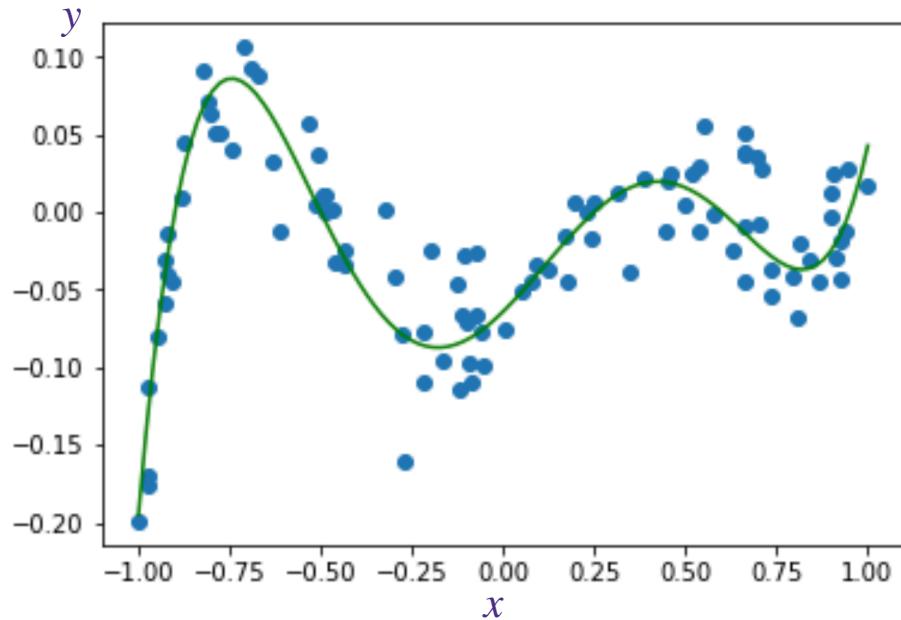
$$[a]^+ \triangleq \max\{a, 0\}$$



the weights capture the change in the slopes

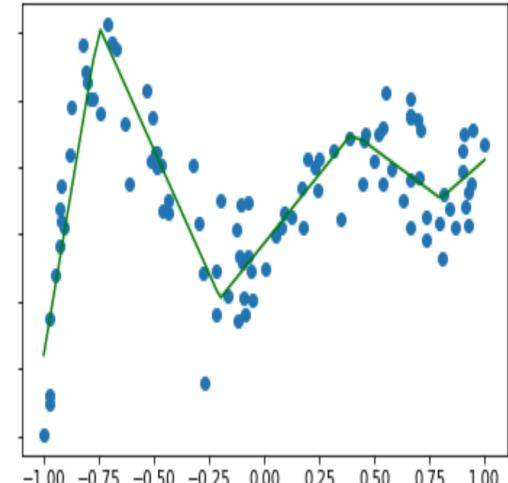
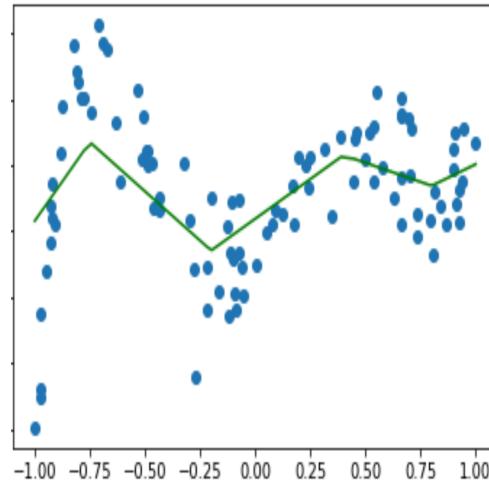
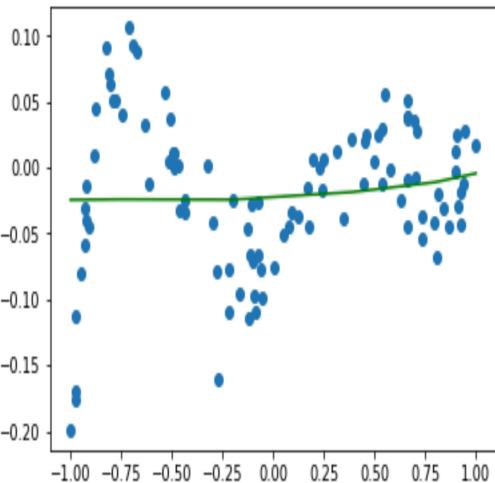
Example: piecewise linear fit

- we fit a linear model:
 $f(x) = b + w_1h_1(x) + w_2h_2(x) + w_3h_3(x) + w_4h_4(x) + w_5h_5(x)$
- with a specific choice of features using piecewise linear functions



Example: piecewise linear fit (ridge regression)

$$L(w) = \sum_{i=1}^n (y_i - w^\top h(x_i) - b)^2 + \lambda \cdot \|w\|_2^2$$



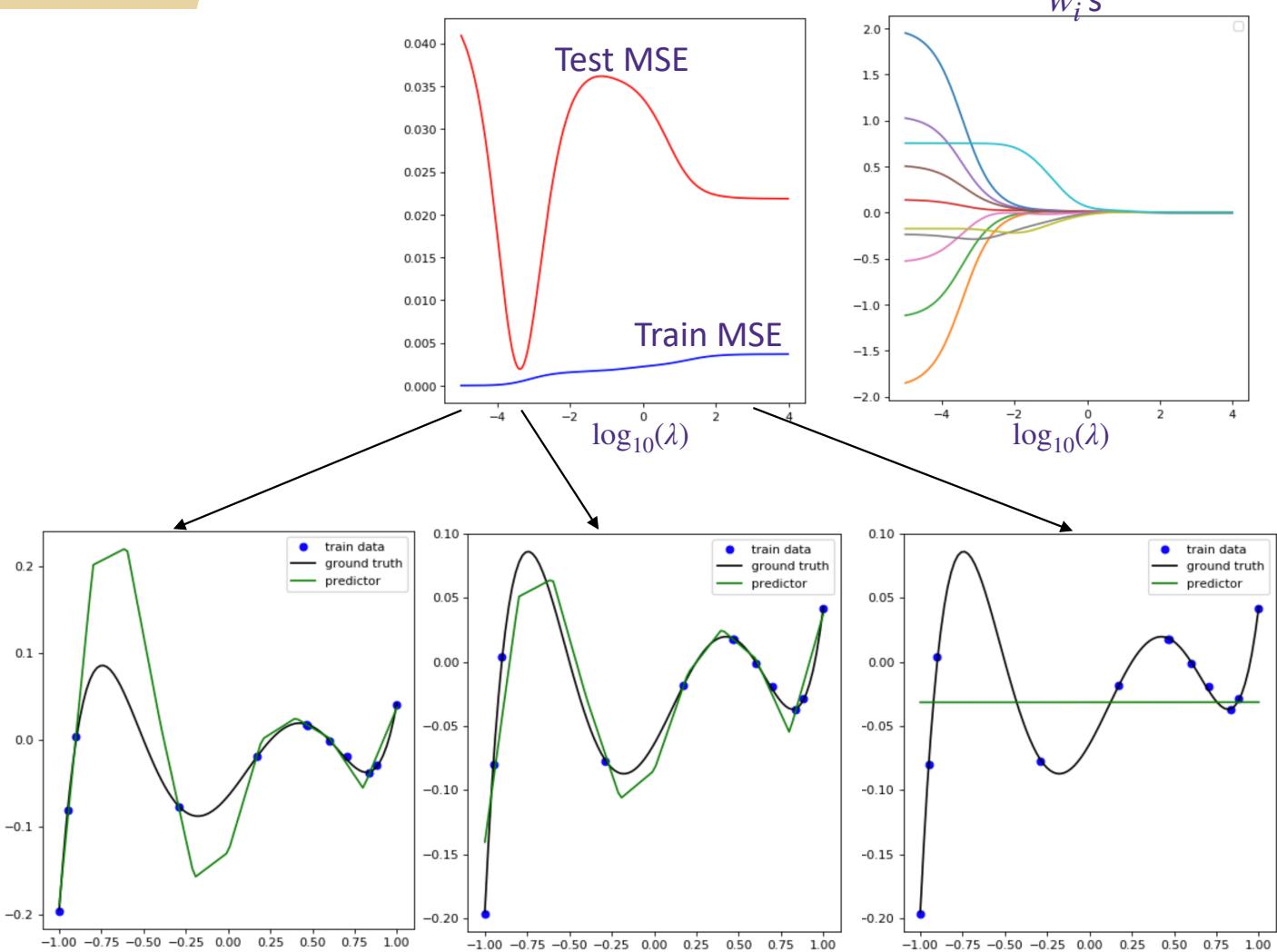
$$\lambda = 1$$

regularization ↑

prefer $w_1, w_2, \dots \approx 0$

We do not observe overfitting, as $d=5$ and $n=100$

Can avoid overfitting even $w \in \mathbb{R}^{10}$ and n=11 samples



Questions?
