

# Bias-variance tradeoff (formally)

CS 229: Machine Learning

Emily Fox

Stanford University

January 17, 2024

©2024 Emily Fox

1

Recap: Assessing performance

©2024 Emily Fox

CS 229: Machine Learning

2

1

## Measuring loss

Loss function:

$$L(y, f_{\hat{w}}(\mathbf{x}))$$

$\hat{f}(\mathbf{x}) = \text{predicted value } \hat{y}$

*actual value*       *$f_{\hat{w}}(\mathbf{x})$*

Cost of using  $\hat{w}$  at  $x$   
when  $y$  is true

Examples: (assuming loss for underpredicting = overpredicting)

Absolute error:  $L(y, f_{\hat{w}}(\mathbf{x})) = |y - f_{\hat{w}}(\mathbf{x})|$

Squared error:  $L(y, f_{\hat{w}}(\mathbf{x})) = (y - f_{\hat{w}}(\mathbf{x}))^2$

3

©2024 Emily Fox

CS 229: Machine Learning

3

## Compute training error

Training error

= avg. loss on houses in **training set**

=  $\frac{1}{N} \sum_{i=1}^N L(y_i, f_{\hat{w}}(\mathbf{x}_i))$

$N$  ← # obs in training set

fit using training data

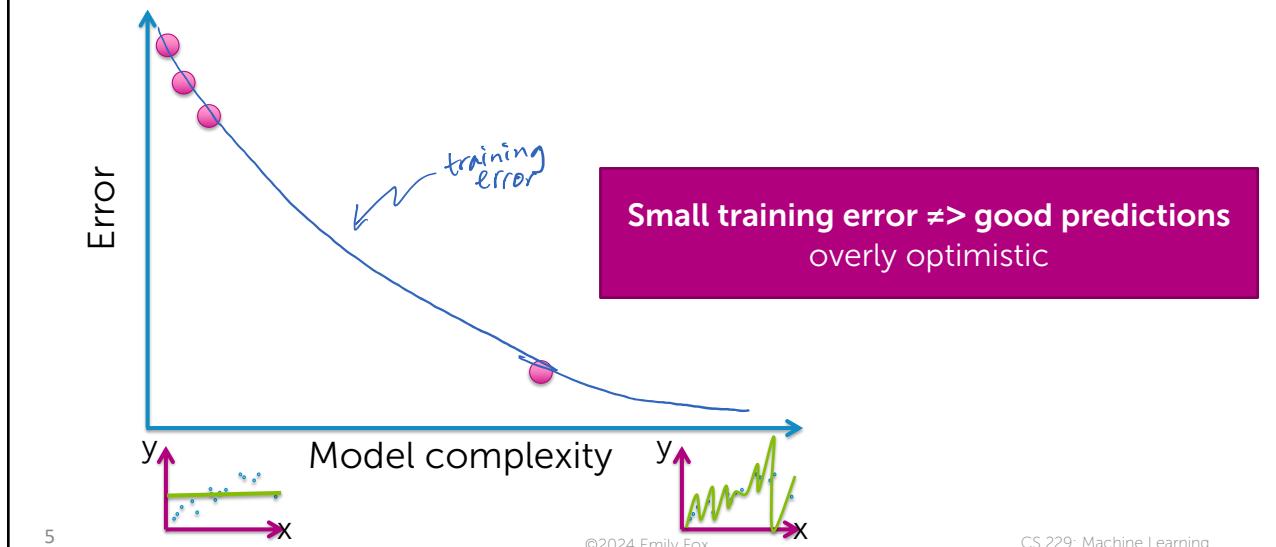
4

©2024 Emily Fox

CS 229: Machine Learning

4

## Training error vs. model complexity



5

## Generalization error definition

Really want estimate of loss over all possible (,\$) pairs

Formally:

average over all possible  
( $x,y$ ) pairs weighted by  
how likely each is

$$\text{generalization error} = \mathbb{E}_{x,y} [L(y, f_w(x))]$$

$$= \int L(y, f_w(x)) \underbrace{p(x, y)}_{p(y|x)p(x)} dx dy$$

fit using training data

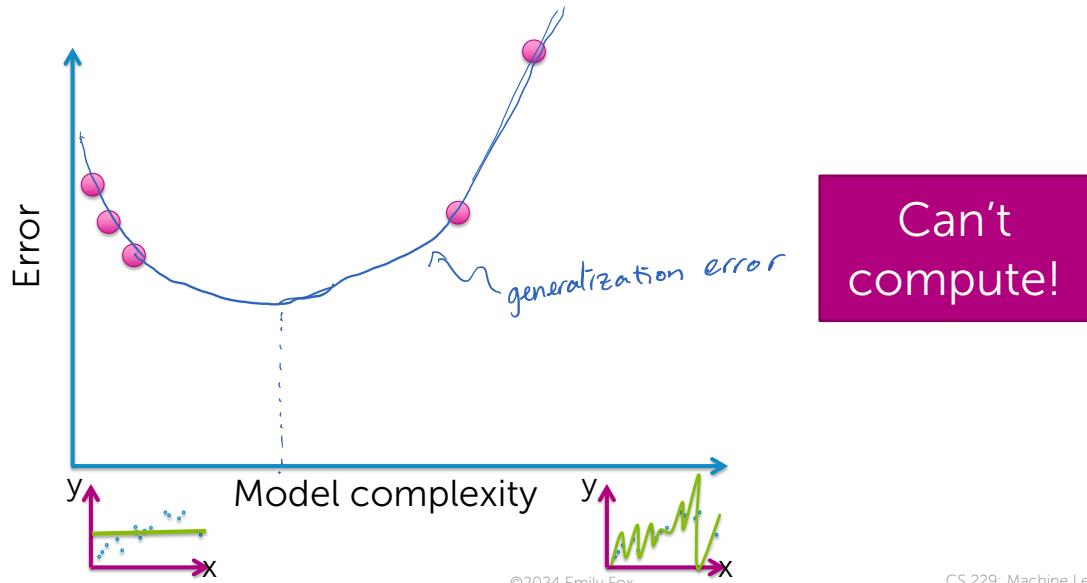
6

©2024 Emily Fox

CS 229: Machine Learning

6

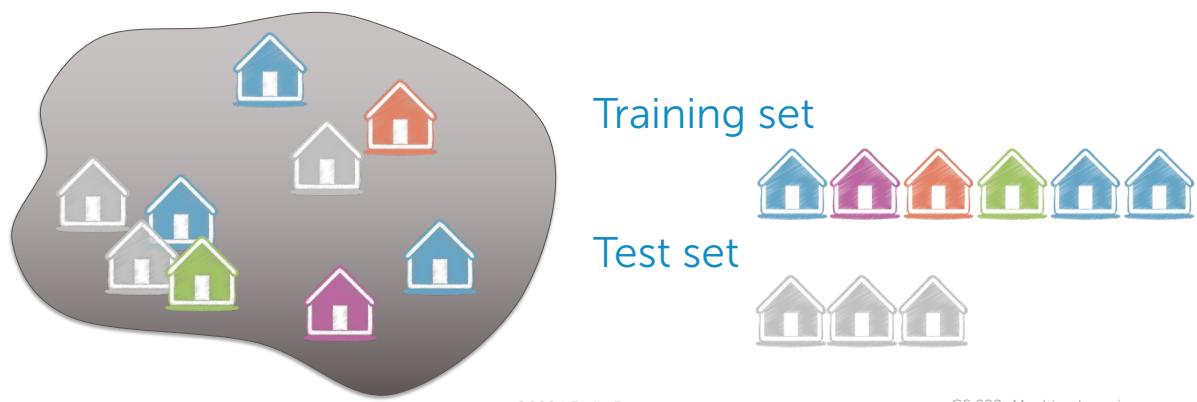
## Generalization error vs. model complexity



7

## Forming a test set

Hold out some (, ) that are *not* used for fitting the model



8

## Compute test error

$$\text{generalization error} = E_{x,y} [L(y, f_{\hat{w}}(x))] \\ = \int L(y, f_{\hat{w}}(x)) p(y, x) dx dy \\ p(y|x) p(x)$$

$$x, y \stackrel{iid}{\sim} p(x, y)$$

### Test error

= avg. loss on houses in **test set**

$$= \frac{1}{N_{test}} \sum_{i \text{ in test set}} L(y_i, f_{\hat{w}}(x_i))$$

# test points      fit using **training data**

**has never seen  
test data!**

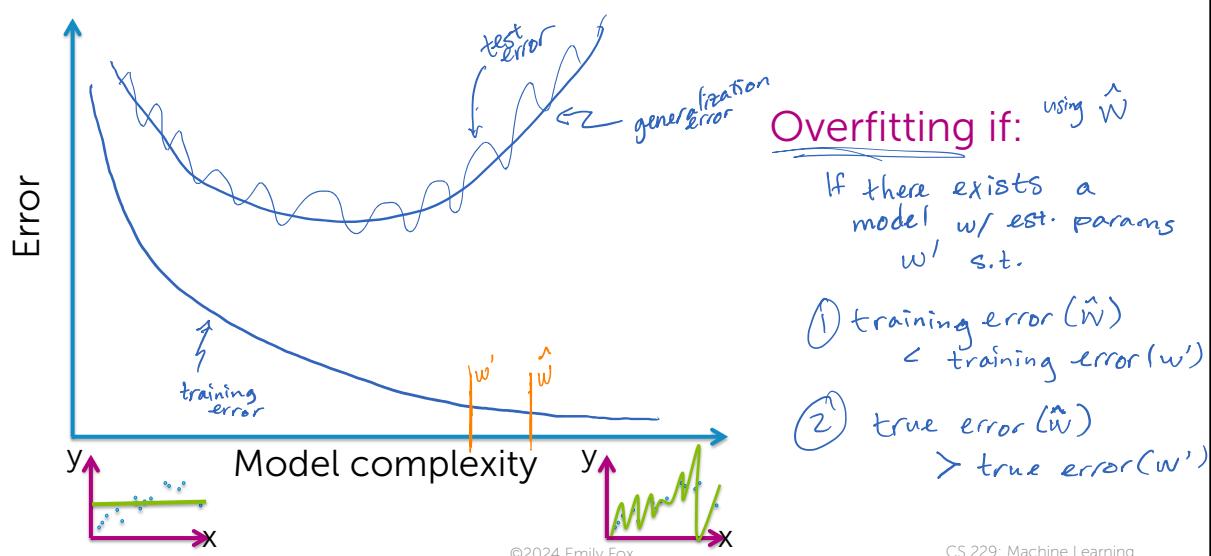
9

©2024 Emily Fox

CS 229: Machine Learning

9

## Training, true, & test error vs. model complexity



10

©2024 Emily Fox

CS 229: Machine Learning

10

## 3 sources of error + the bias-variance tradeoff

©2024 Emily Fox

CS 229: Machine Learning

11

## 3 sources of error

In forming predictions, there are 3 sources of error:

1. Noise
2. Bias
3. Variance

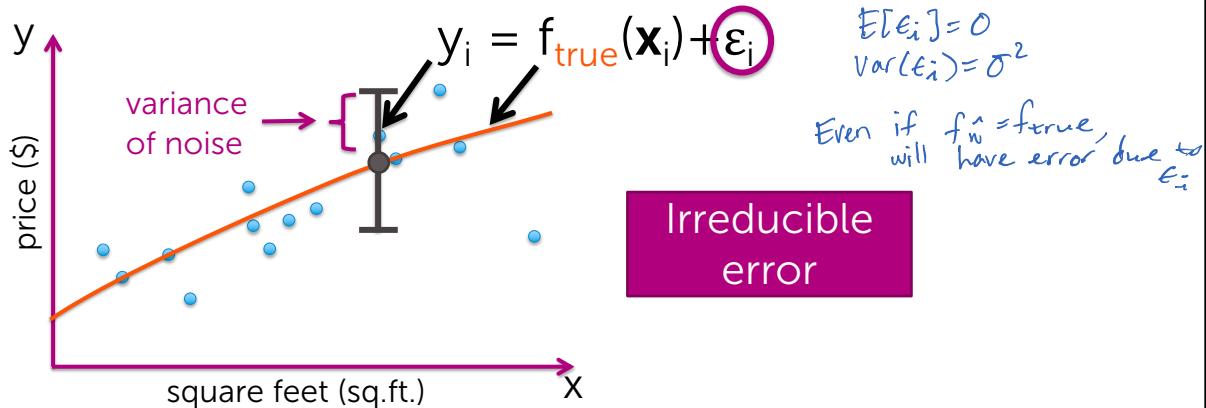
12

©2024 Emily Fox

CS 229: Machine Learning

12

## Data inherently noisy



13

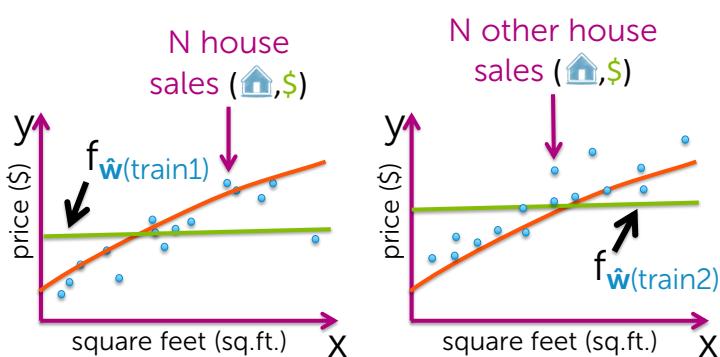
©2024 Emily Fox

CS 229: Machine Learning

13

## Bias contribution

Assume we fit a constant function



14

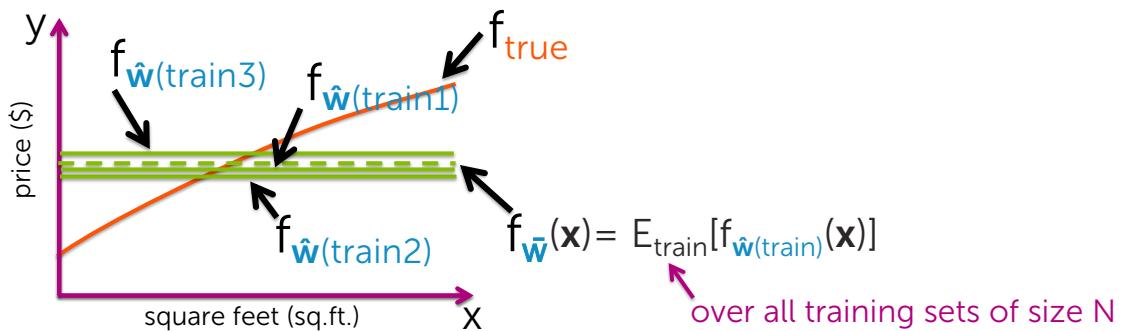
©2024 Emily Fox

CS 229: Machine Learning

14

## Bias contribution

Over all possible size N training sets,  
what do I expect my fit to be?



15

©2024 Emily Fox

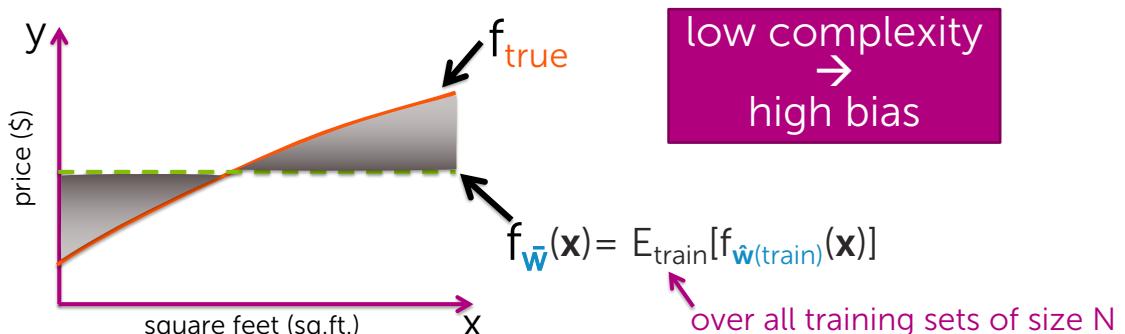
CS 229: Machine Learning

15

## Bias contribution

$$\text{Bias}(f_{\bar{w}}(x)) = f_{\text{true}}(x) - f_{\bar{w}}(x)$$

Is our approach flexible  
enough to capture  $f_{\text{true}}$ ?  
If not, error in predictions.



16

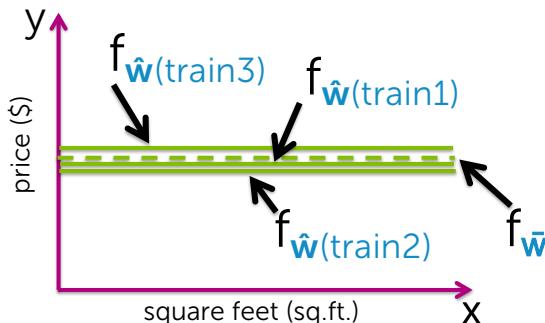
©2024 Emily Fox

CS 229: Machine Learning

16

## Variance contribution

How much do specific fits vary from the expected fit?



17

©2024 Emily Fox

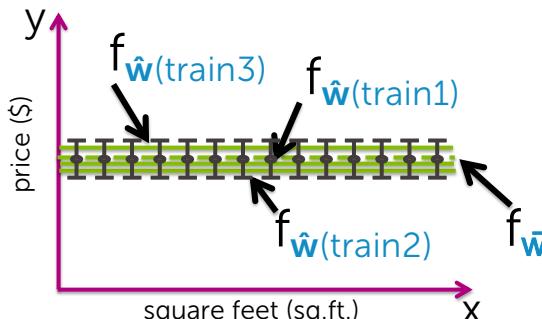
CS 229: Machine Learning

17

## Variance contribution

$$\text{var}(x) = E[(x - E[x])^2]$$

How much do specific fits vary from the expected fit?



$$\text{var}(f_{\hat{w}}(x)) = E_{\text{train}}[(f_{\hat{w}}(\text{train})(x) - f_{\bar{w}}(x))^2]$$

fit on a specific training dataset  
 over all training sets of size N  
 deviation of specific fit from expected fit at  $x$   
 what I expect to learn over all training sets

18

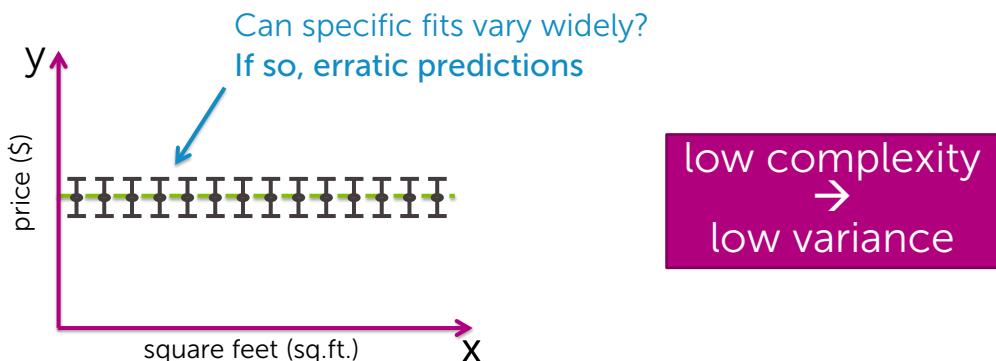
©2024 Emily Fox

CS 229: Machine Learning

18

## Variance contribution

How much do specific fits vary from the expected fit?



19

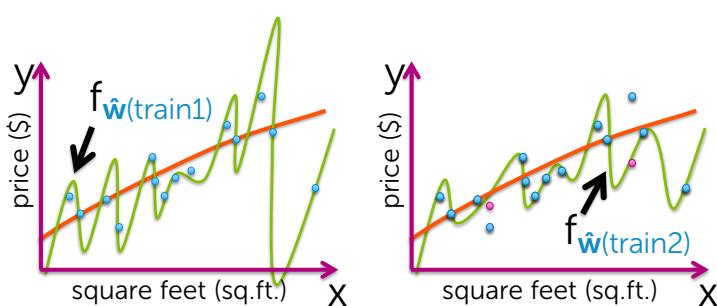
©2024 Emily Fox

CS 229: Machine Learning

19

## Variance of high-complexity models

Assume we fit a high-order polynomial



20

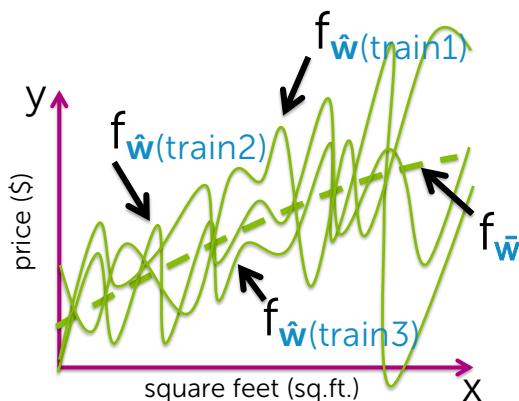
©2024 Emily Fox

CS 229: Machine Learning

20

## Variance of high-complexity models

Assume we fit a high-order polynomial



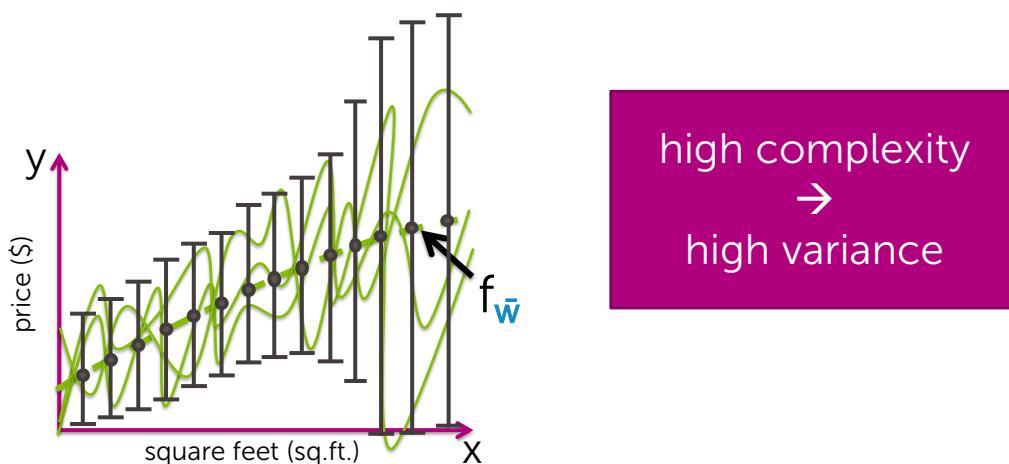
21

©2024 Emily Fox

CS 229: Machine Learning

21

## Variance of high-complexity models



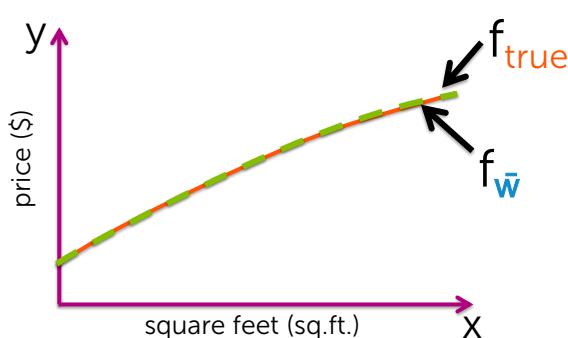
22

©2024 Emily Fox

CS 229: Machine Learning

22

## Bias of high-complexity models



high complexity  
→  
low bias

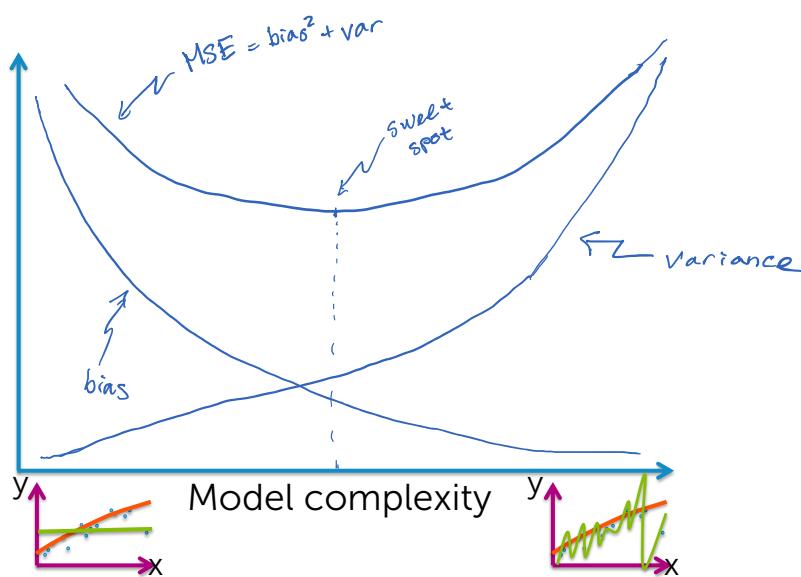
23

©2024 Emily Fox

CS 229: Machine Learning

23

## Bias-variance tradeoff



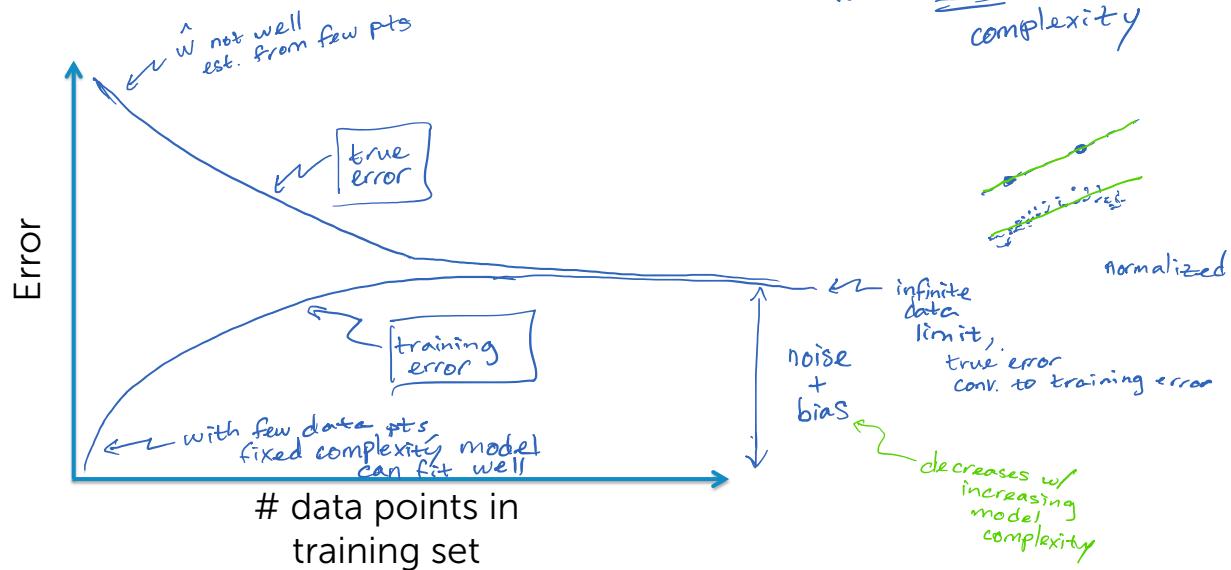
24

©2024 Emily Fox

CS 229: Machine Learning

24

## Error vs. amount of data



25

©2024 Emily Fox

CS 229: Machine Learning

25

Why 3 sources of error?  
A formal derivation

©2024 Emily Fox

CS 229: Machine Learning

26

# Deriving expected prediction error

Expected prediction error

$$= E_{\text{train}} [\text{generalization error of } \hat{\mathbf{w}}(\text{train})]$$

$$= E_{\text{train}} [E_{\mathbf{x}, y} [L(y, f_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x}))]]$$

1. Look at specific  $\mathbf{x}$
2. Consider  $L(y, f_{\hat{\mathbf{w}}}(\mathbf{x})) = (y - f_{\hat{\mathbf{w}}}(\mathbf{x}))^2$

Expected prediction error at  $\mathbf{x}$

$$= E_{\text{train}} [E_{y|\mathbf{x}} [(y - f_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x}))^2]]$$

27

©2024 Emily Fox

CS 229: Machine Learning

27

# Deriving expected prediction error

Expected prediction error at  $\mathbf{x}$

$$= E_{\text{train}} [E_{y|\mathbf{x}} [(y - f_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x}))^2]]$$

$$= E_{\text{train}} [E_{y|\mathbf{x}} [((y - f_{\text{true}}(\mathbf{x})) + (f_{\text{true}}(\mathbf{x}) - f_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x})))^2]]$$

$$= E_{\text{train}} [E_{y|\mathbf{x}} [(y - f)^2]] + 2 E_{\text{train}} [E_{y|\mathbf{x}} [(y - f)(f - \hat{f})]]$$

$$+ E_{\text{train}} [E_{y|\mathbf{x}} [(f - \hat{f})^2]]$$

$$\approx \underline{\text{MSE}(\hat{f})} \text{ mean square error}$$

$$= \sigma^2 + \underline{\text{MSE}(\hat{f})}$$

Shorthand:  
 $f_{\text{true}} \rightarrow f$   
 $f_{\hat{\mathbf{w}}(\text{train})} \rightarrow \hat{f}$

$$E[(a+b)^2] = E[a^2] + 2E[ab] + E[b^2]$$

$$E[ab] = E[a]E[b]$$

if  $a, b$  uncorr (or, indep.)

28

©2024 Emily Fox

CS 229: Machine Learning

# Equating MSE with bias and variance

$$\begin{aligned}
 \text{MSE}(\mathbf{x}) &= E_{\text{train}}[(f_{\text{true}}(\mathbf{x}) - \hat{f}_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x}))^2] \\
 &= E_{\text{train}}[((f_{\text{true}}(\mathbf{x}) - f_{\bar{\mathbf{w}}}(\mathbf{x})) + (f_{\bar{\mathbf{w}}}(\mathbf{x}) - \hat{f}_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x})))^2] \\
 &= E_{\text{train}}[(f_{\text{true}} - \bar{f})^2] + 2E_{\text{train}}[(f_{\text{true}} - \bar{f})(\bar{f} - \hat{f})] + E_{\text{train}}[(\bar{f} - \hat{f})^2] \\
 &\quad \underbrace{=} \text{bias}^2(\hat{f}) \text{ by defn} \quad \underbrace{=} \text{not "f" or "train"} \quad \underbrace{=} \text{var}(\hat{f}) \\
 &= \text{bias}^2(\hat{f}) + \text{var}(\hat{f})
 \end{aligned}$$

29

©2024 Emily Fox

CS 229: Machine Learning

# Putting it all together

## Expected prediction error at $\mathbf{x}$

$$= \sigma^2 + \text{MSE}(f_{\hat{w}}(\mathbf{x}))$$

$$\equiv \sigma^2 + [\text{bias}(f_{\hat{w}}(\mathbf{x}))]^2 + \text{var}(f_{\hat{w}}(\mathbf{x}))$$

## 3 sources of error

30

©2024 Emily Fox

CS 229: Machine Learning

30

## Summary of assessing performance

©2024 Emily Fox

CS 229: Machine Learning

31

## What you can do now...

- Describe what a loss function is and give examples
- Contrast training, generalization, and test error
- Compute training and test error given a loss function
- Discuss issue of assessing performance on training set
- Define overfitting in terms of training and generalization (or, in practice, test) error
- Describe tradeoffs in forming training/test splits
- List and interpret the 3 sources of avg. prediction error
  - irreducible error, bias, and variance
- Derive avg. prediction error in terms of irreducible error, bias, and variance
- Sketch:
  - training & generalization/test error vs. model complexity
  - training & generalization error vs. # of training data points

32

©2024 Emily Fox

CS 229: Machine Learning

32

16

# Ridge Regression:

## Regulating overfitting when using many features

CS 229: Machine Learning

Emily Fox

Stanford University

January 17, 2024

©2024 Emily Fox

33

## Overfitting of polynomial regression

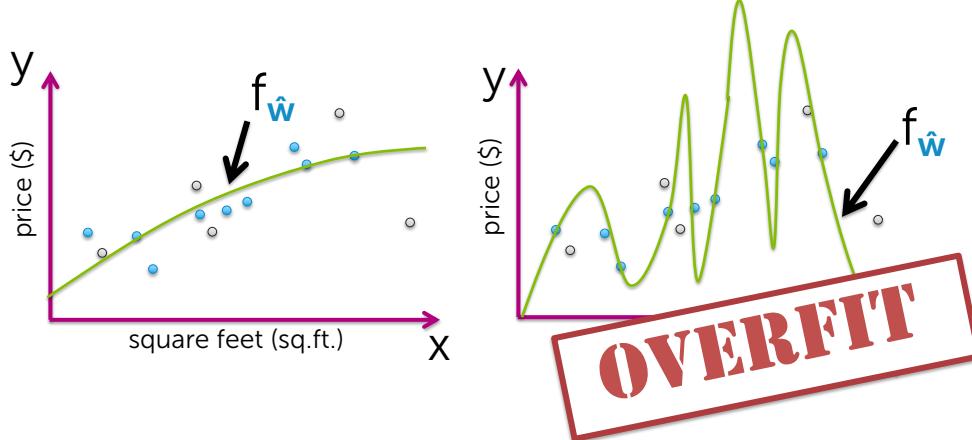
©2024 Emily Fox

CS 229: Machine Learning

34

## Flexibility of high-order polynomials

$$y_i = w_0 + w_1 x_i + w_2 x_i^2 + \dots + w_p x_i^p + \epsilon_i$$



35

©2024 Emily Fox

CS 229: Machine Learning

35

## Symptom of overfitting

Often, overfitting associated with very large estimated parameters  $\hat{w}$

36

©2024 Emily Fox

CS 229: Machine Learning

36

## Overfitting of linear regression models more generically

©2024 Emily Fox

CS 229: Machine Learning

37

## Overfitting with many features

Not unique to polynomial regression

Generically, can happen with  
**lots of features (D large)**

$$y_i = \sum_{j=0}^D w_j h_j(\mathbf{x}_i) + \varepsilon_i$$

- Square feet
- # bathrooms
- # bedrooms
- Lot size
- Year built
- ...
- .

38

©2024 Emily Fox

CS 229: Machine Learning

38

19

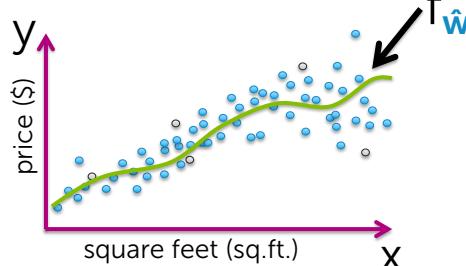
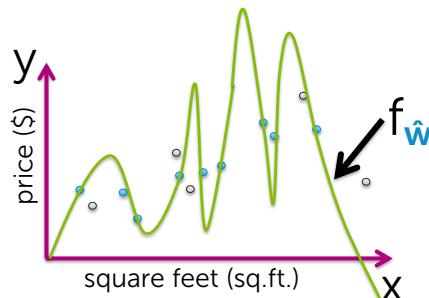
## How does # of observations influence overfitting?

Few observations (N small)

→ rapidly overfit as model complexity increases

Many observations (N very large)

→ harder to overfit



39

©2024 Emily Fox

CS 229: Machine Learning

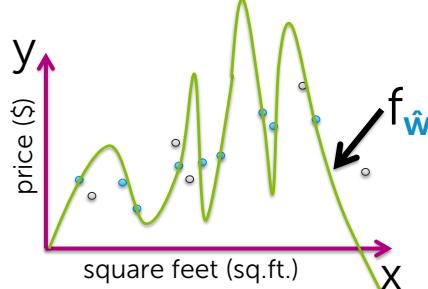
39

## How does # of inputs influence overfitting?

1 input (e.g., sq.ft.):

Data must include representative examples of all possible (sq.ft., \$) pairs to avoid overfitting

**HARD**



40

©2024 Emily Fox

CS 229: Machine Learning

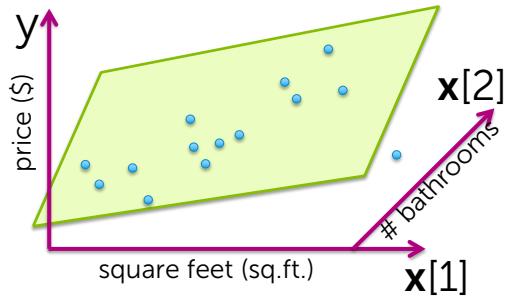
40

## How does # of inputs influence overfitting?

**d inputs** (e.g., sq.ft., #bath, #bed, lot size, year,...):

Data must include examples of all possible  
(sq.ft., #bath, #bed, lot size, year,..., \$) combos  
to avoid overfitting

**MUCH!!!  
HARDER**



41

CS 229: Machine Learning

41

Adding term to cost (loss) function  
to prefer small coefficients

©2024 Emily Fox

CS 229: Machine Learning

42

## Desired total cost format

Want to balance:

- How well function fits data
- Magnitude of coefficients

$$\text{Total cost} = \text{measure of fit} + \text{measure of magnitude of coefficients}$$

want to balance

↑ small # = good fit to training data

↑ small # = not overfit

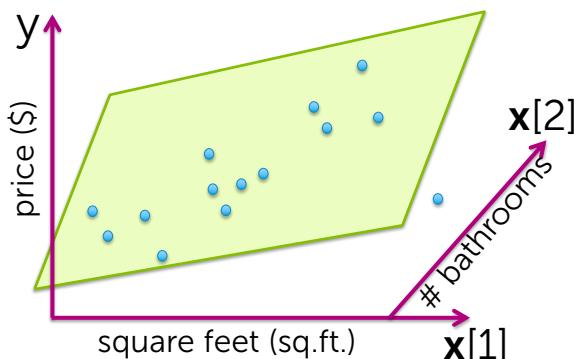
43

©2024 Emily Fox

CS 229: Machine Learning

43

## Measure of fit to training data



$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N (y_i - h(\mathbf{x}_i)^T \mathbf{w})^2$$

pred. value using  $\mathbf{w}$

small RSS → model fitting training data well

44

©2024 Emily Fox

CS 229: Machine Learning

44

## Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

- Sum?  $w_0 = 1,527,301 \quad w_1 = -1,605,253$   
 $w_0 + w_1 = \text{small } \# \quad X$
- Sum of absolute value?  
 $\sum_{j=0}^D |w_j| \triangleq \|w\|_1 \quad L_1 \text{ norm... discuss next lecture}$
- Sum of squares ( $L_2$  norm)  
 $\sum_{j=0}^D w_j^2 \triangleq \|w\|_2^2 \quad \boxed{\begin{array}{l} L_2 \text{ norm...} \\ \text{focus of this lecture} \end{array}}$

45

©2024 Emily Fox

CS 229: Machine Learning

45

## Consider specific total cost

Total cost =  
measure of fit + measure of magnitude of coefficients

46

©2024 Emily Fox

CS 229: Machine Learning

46

## Consider specific total cost

Total cost =

$$\text{measure of fit} + \text{measure of magnitude of coefficients}$$

$\text{RSS}(\mathbf{w})$ 
 $\|\mathbf{w}\|_2^2$

47

©2024 Emily Fox

CS 229: Machine Learning

47

## Consider resulting objective

What if  $\hat{\mathbf{w}}$  selected to minimize

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

↑ tuning parameter = balance of fit and magnitude

If  $\lambda=0$ :

reduces to  $\min \text{RSS}(\mathbf{w})$ , old soln,  $\rightarrow \hat{\mathbf{w}}^{\text{LS}}$  (least squares)

If  $\lambda=\infty$ :

For solns where  $\hat{\mathbf{w}} \neq 0$ , then total cost =  $\infty$   
 If  $\hat{\mathbf{w}} = 0$ , then total cost =  $\text{RSS}(0) \rightarrow \hat{\mathbf{w}} = 0$

If  $\lambda$  in between:

$$\text{Then } 0 \leq \|\hat{\mathbf{w}}\|_2^2 \leq \|\hat{\mathbf{w}}^{\text{LS}}\|_2^2$$

48

©2024 Emily Fox

CS 229: Machine Learning

## Consider resulting objective

What if  $\hat{\mathbf{w}}$  selected to minimize

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

↑ tuning parameter = balance of fit and magnitude

Ridge regression  
(a.k.a  $L_2$  regularization)

49

©2024 Emily Fox

CS 229: Machine Learning

49

## Bias-variance tradeoff

Large  $\lambda$ :

high bias, low variance

(e.g.,  $\hat{\mathbf{w}} = 0$  for  $\lambda = \infty$ )

Small  $\lambda$ :

low bias, high variance

(e.g., standard least squares (RSS) fit of high-order polynomial for  $\lambda = 0$ )

In essence,  $\lambda$   
controls model  
complexity

50

©2024 Emily Fox

CS 229: Machine Learning

50

## Revisit polynomial fit demo

What happens if we refit our high-order polynomial, but now using ridge regression?

Will consider a few settings of  $\lambda$  ...

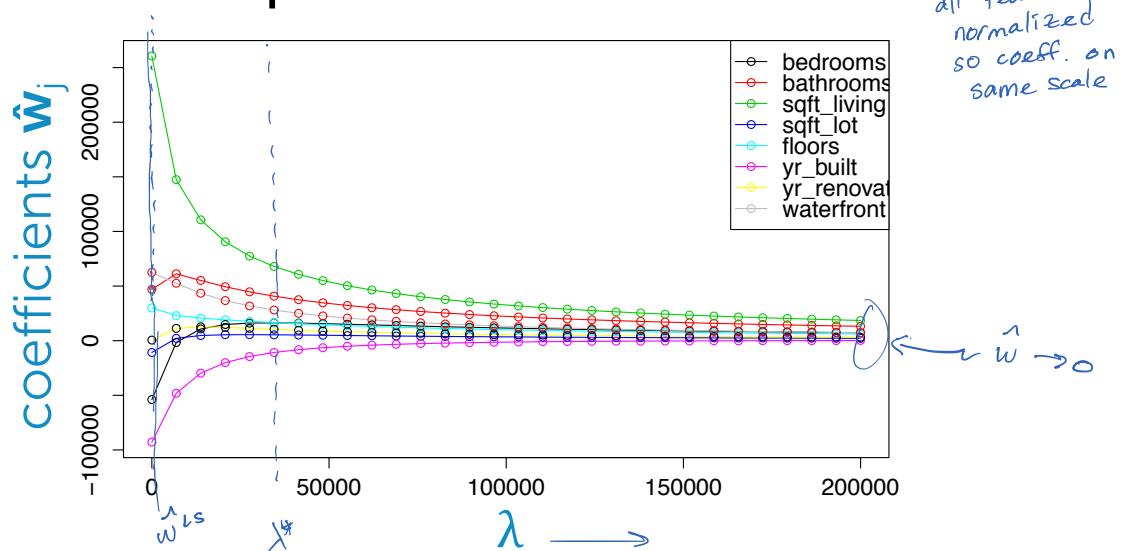
51

©2024 Emily Fox

CS 229: Machine Learning

51

## Coefficient path



52

©2024 Emily Fox

CS 229: Machine Learning

52

## Fitting the ridge regression model (for given $\lambda$ value)

©2024 Emily Fox

CS 229: Machine Learning

53

### Gradient of ridge regression cost

$$\begin{aligned} \nabla [\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2] &= \nabla [(\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w}] \\ &= \underbrace{\nabla [(\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w})]}_{-2\mathbf{H}^\top(\mathbf{y} - \mathbf{H}\mathbf{w})} + \underbrace{\lambda \nabla [\mathbf{w}^\top \mathbf{w}]}_{2\mathbf{w}} \end{aligned}$$

$$\begin{aligned} \|\mathbf{w}\|_2^2 &= w_0^2 + w_1^2 + \dots + w_D^2 \\ &= \mathbf{w}^\top \mathbf{w} \\ &= [w_0 \dots w_D] \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix} \end{aligned}$$

**Why?** By analogy to 1d case...

$\mathbf{w}^\top \mathbf{w}$  analogous to  $w^2$  and derivative of  $w^2 = 2w$

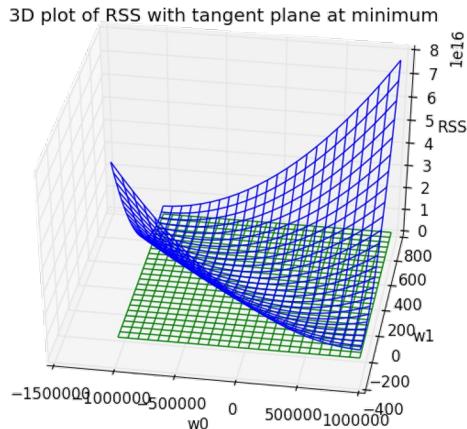
54

©2024 Emily Fox

CS 229: Machine Learning

54

## Ridge closed-form solution



$$\nabla \text{cost}(\mathbf{w}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) + 2\lambda \mathbf{I}\mathbf{w} = 0$$

Solve for  $\mathbf{w}$ :

$$\hat{\mathbf{w}}_{\text{ridge}} = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^T \mathbf{y}$$

{ "regularizer"

55

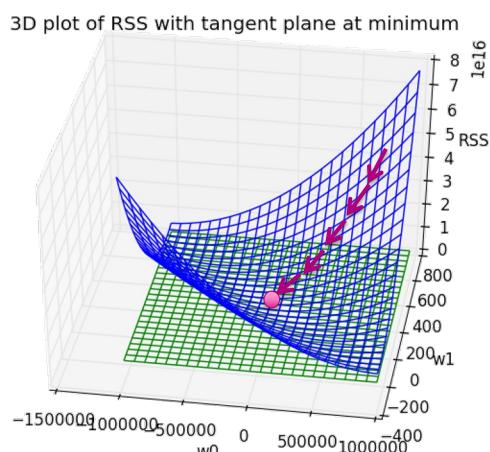
©2024 Emily Fox

CS 229: Machine Learning

55

## Gradient descent for ridge regression

$$\nabla \text{cost}(\mathbf{w}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) + 2\lambda \mathbf{w}$$



Algorithm:

**while** not converged

$$\begin{aligned} \mathbf{w}^{(t+1)} &\leftarrow \mathbf{w}^{(t)} - \eta \nabla \text{cost}(\mathbf{w}^{(t)}) \\ &\quad - 2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}^{(t)}) + 2\lambda \mathbf{w}^{(t)} \\ &\leftarrow \underbrace{(1 - 2\eta\lambda)}_{-\eta\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}^{(t)})} \mathbf{w}^{(t)} \end{aligned}$$

56

©2024 Emily Fox

CS 229: Machine Learning

56

## How to choose $\lambda$

©2024 Emily Fox

CS 229: Machine Learning

57

## The regression/ML workflow

### 1. Model selection

Need to **choose tuning parameters  $\lambda$**  controlling model complexity

### 2. Model assessment

Having selected a model, **assess generalization error**

58

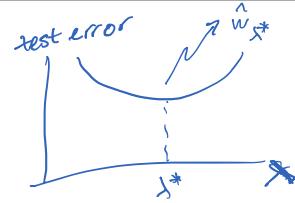
©2024 Emily Fox

CS 229: Machine Learning

58

29

## Hypothetical implementation



### 1. Model selection

For each considered  $\lambda$  :

- Estimate parameters  $\hat{w}_\lambda$  on training data
- Assess performance of  $\hat{w}_\lambda$  on test data
- Choose  $\lambda^*$  to be  $\lambda$  with lowest test error

Overly optimistic!

### 2. Model assessment

Compute test error of  $\hat{w}_{\lambda^*}$  (fitted model for selected  $\lambda^*$ ) to approx. generalization error

59

©2024 Emily Fox

CS 229: Machine Learning

59

## Hypothetical implementation



**Issue:** Just like fitting  $\hat{w}$  and assessing its performance both on training data

- $\lambda^*$  was selected to minimize test error (i.e.,  $\lambda^*$  was fit on test data)
- If test data is not representative of the whole world, then  $\hat{w}_{\lambda^*}$  will typically perform worse than test error indicates

60

©2024 Emily Fox

CS 229: Machine Learning

60

## Practical implementation



**Solution:** Create two “test” sets!

1. Select  $\lambda^*$  such that  $\hat{w}_{\lambda^*}$  minimizes error on validation set
2. Approximate generalization error of  $\hat{w}_{\lambda^*}$  using test set

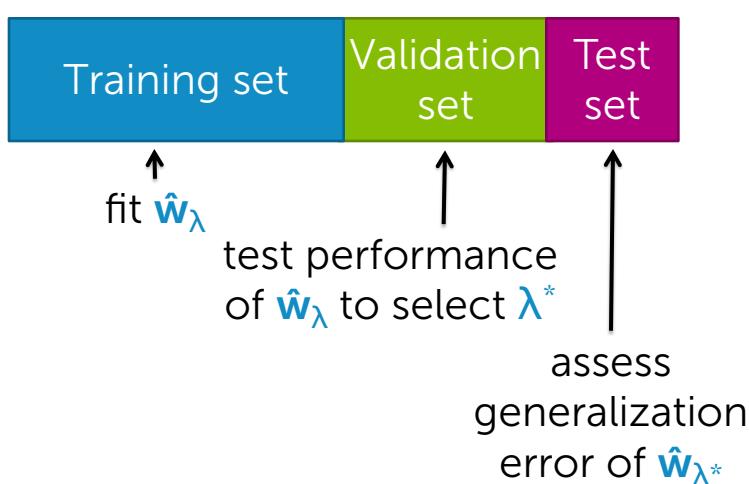
61

©2024 Emily Fox

CS 229: Machine Learning

61

## Practical implementation



62

©2024 Emily Fox

CS 229: Machine Learning

62

## Typical splits

Training set	Validation set	Test set
80%	10%	10%
50%	25%	25%

63

©2024 Emily Fox

CS 229: Machine Learning

63

## Summary for ridge regression

©2024 Emily Fox

CS 229: Machine Learning

64

## What you can do now...

- Describe what happens to magnitude of estimated coefficients when model is overfit
- Motivate form of ridge regression cost function
- Describe what happens to estimated coefficients of ridge regression as tuning parameter  $\lambda$  is varied
- Interpret coefficient path plot
- Estimate ridge regression parameters:
  - In closed form
  - Using an iterative gradient descent algorithm
- Use a validation set to select the ridge regression tuning parameter  $\lambda$
- Handle intercept and scale of features with care

65

©2024 Emily Fox

CS 229: Machine Learning

65

Fitting the ridge regression model  
(for given  $\lambda$  value) in more detail

**OPTIONAL**

©2024 Emily Fox

CS 229: Machine Learning

66

## Step 1:

Rewrite total cost in matrix notation

©2024 Emily Fox

CS 229: Machine Learning

67

## Recall matrix form of RSS

Model for all N observations together

$$\mathbf{y} = \mathbf{H} \mathbf{w} + \boldsymbol{\epsilon}$$

$$\begin{aligned} \text{RSS}(\mathbf{w}) &= \sum_{i=1}^N (y_i - h(\mathbf{x}_i)^T \mathbf{w})^2 \\ &= (\mathbf{y} - \mathbf{H}\mathbf{w})^T (\mathbf{y} - \mathbf{H}\mathbf{w}) \end{aligned}$$

68

©2024 Emily Fox

CS 229: Machine Learning

68

34

# Rewrite magnitude of coefficients in vector notation

$$\|w\|_2^2 = w_0^2 + w_1^2 + w_2^2 + \dots + w_D^2$$

$$= \begin{array}{c} \boxed{\phantom{0}} \quad \boxed{\phantom{0}} \quad \boxed{\phantom{0}} \quad \boxed{\phantom{0}} \quad \boxed{\phantom{0}} \\ w_0 \quad w_1 \quad - \quad - \quad w_D \end{array} \quad \begin{array}{c} \boxed{\phantom{0}} \\ \vdots \\ \boxed{\phantom{0}} \\ \vdots \\ \boxed{\phantom{0}} \\ w_0 \\ w_1 \\ \vdots \\ w_D \end{array}$$

$$= w^T w$$

69

©2024 Emily Fox

CS 229: Machine Learning

# Putting it all together

In matrix form, ridge regression cost is:

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

$$= (\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w}$$

70

©2024 Emily Fox

CS 229: Machine Learning

70

## Step 2:

### Compute the gradient

©2024 Emily Fox

CS 229: Machine Learning

71

## Gradient of ridge regression cost

$$\begin{aligned}\|\mathbf{w}\|_2 &= \sqrt{w_1^2 + \dots + w_p^2} \\ &= \sqrt{\mathbf{w}^\top \mathbf{w}}\end{aligned}$$

$$\begin{aligned}\nabla [\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2] &= \nabla [(\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w}] \\ &= \underbrace{\nabla [(\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w})]}_{-2\mathbf{H}^\top(\mathbf{y} - \mathbf{H}\mathbf{w})} + \lambda \underbrace{\nabla [\mathbf{w}^\top \mathbf{w}]}_{2\mathbf{w}}\end{aligned}$$

**Why?** By analogy to 1d case...

$\mathbf{w}^\top \mathbf{w}$  analogous to  $w^2$  and derivative of  $w^2 = 2w$

72

©2024 Emily Fox

CS 229: Machine Learning

72

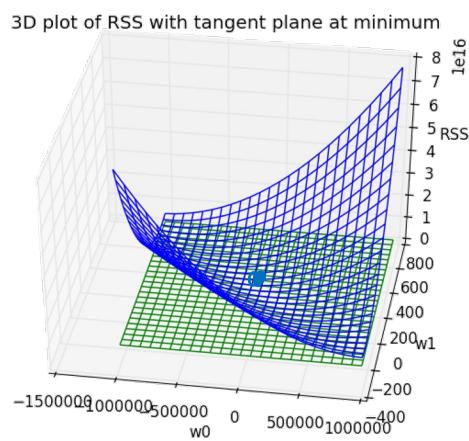
## Step 3, Approach 1: Set the gradient = 0

©2024 Emily Fox

CS 229: Machine Learning

73

## Ridge closed-form solution



$$\nabla \text{cost}(\mathbf{w}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) + 2\lambda \mathbf{I}\mathbf{w} = 0$$

add

Solve for  $\mathbf{w}$ :

$$\begin{aligned}
 -\mathbf{H}^T\mathbf{y} + \underline{\mathbf{H}^T\mathbf{H}}\hat{\mathbf{w}} + \underline{\lambda \mathbf{I}}\hat{\mathbf{w}} &= 0 \\
 (\mathbf{H}^T\mathbf{H} + \lambda \mathbf{I})\hat{\mathbf{w}} &= \mathbf{H}^T\mathbf{y} \\
 \hat{\mathbf{w}}^{\text{ridge}} &= (\mathbf{H}^T\mathbf{H} + \lambda \mathbf{I})^{-1}\mathbf{H}^T\mathbf{y}
 \end{aligned}$$

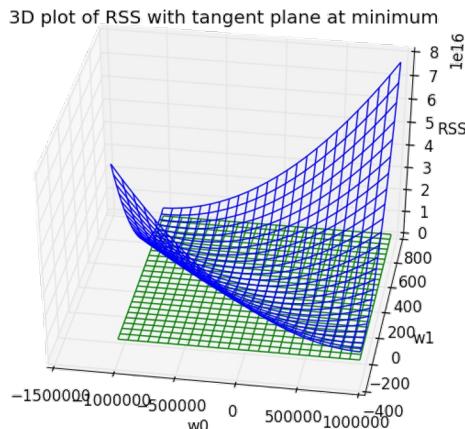
74

©2024 Emily Fox

CS 229: Machine Learning

74

## Interpreting ridge closed-form solution



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

If  $\lambda=0$ :  $\hat{\mathbf{w}}_{\text{ridge}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y} = \hat{\mathbf{w}}^{\text{LS}}$

If  $\lambda=\infty$ :  $\hat{\mathbf{w}}_{\text{ridge}} = \mathbf{0}$  ← because it's like dividing by  $\infty$

75

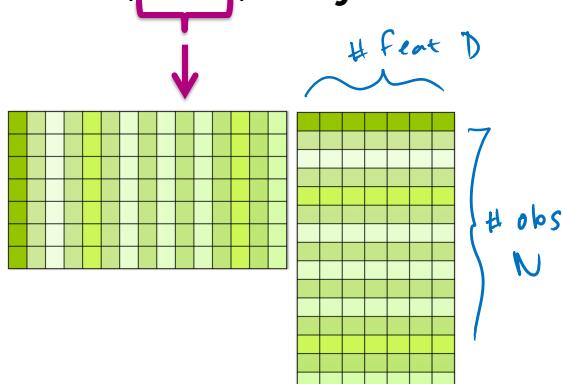
©2024 Emily Fox

CS 229: Machine Learning

75

## Recall discussion on previous closed-form solution

$$\hat{\mathbf{w}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y}$$



Invertible if:

In general,  
(# linearly independent obs)  
 $N > D$

Complexity of inverse:  
 $O(D^3)$

76

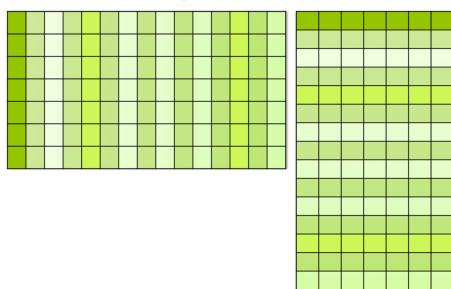
©2024 Emily Fox

CS 229: Machine Learning

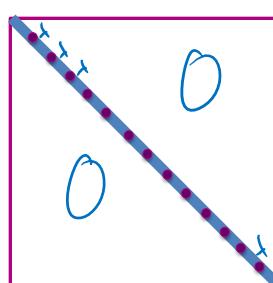
76

## Discussion of ridge closed-form solution

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



+



$\lambda \mathbf{I}$  is making  $\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I}$  more "regular"  
 $\rightarrow$  "regularization"

Invertible if:  
 Always if  $\lambda > 0$ ,  
 even if  $N < D$

Complexity of  
 inverse:  
 $O(D^3)$ ...  
 big for large  $D$ !

77

©2024 Emily Fox

CS 229: Machine Learning

77

## Step 3, Approach 2: Gradient descent

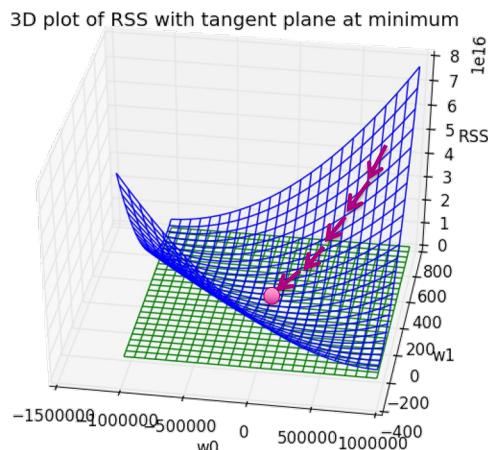
©2024 Emily Fox

CS 229: Machine Learning

78

## Gradient descent for ridge regression

$$\nabla \text{cost}(\mathbf{w}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) + 2\lambda\mathbf{w}$$



79

©2024 Emily Fox

CS 229: Machine Learning

Algorithm:

**while** not converged

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla \text{cost}(\mathbf{w}^{(t)})$$

$$-2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}^{(t)}) + 2\lambda\mathbf{w}^{(t)}$$

How to handle the intercept

PRACTICALITIES

80

©2024 Emily Fox

CS 229: Machine Learning

## Recall multiple regression model

Model:

$$\begin{aligned} y_i &= w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i) + \varepsilon_i \\ &= \sum_{j=0}^D w_j h_j(\mathbf{x}_i) + \varepsilon_i \end{aligned}$$

*feature 1 =  $h_0(\mathbf{x})$ ...often 1 (constant)*

*feature 2 =  $h_1(\mathbf{x})$ ... e.g.,  $\mathbf{x}[1]$*

*feature 3 =  $h_2(\mathbf{x})$ ... e.g.,  $\mathbf{x}[2]$*

...

*feature  $D+1 = h_D(\mathbf{x})$ ... e.g.,  $\mathbf{x}[d]$*

81

©2024 Emily Fox

CS 229: Machine Learning

81

## If constant feature (intercept)...

$$y_i = w_0 + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i) + \varepsilon_i$$

In matrix notation for  $N$  observations:

$$\mathbf{y} = \mathbf{H} \mathbf{w} + \boldsymbol{\varepsilon}$$

82

©2024 Emily Fox

CS 229: Machine Learning

82

## Do we penalize intercept?

Standard ridge regression cost:

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

↗ strength of penalty

Encourages intercept  $w_0$  to also be small

Do we want a small intercept?

Conceptually, not indicative of overfitting...

83

©2024 Emily Fox

CS 229: Machine Learning

83

## Option 1: Don't penalize intercept

Modified ridge regression cost:

$$\text{RSS}(w_0, \mathbf{w}_{\text{rest}}) + \lambda \|\mathbf{w}_{\text{rest}}\|_2^2$$

How to implement this in practice?

84

©2024 Emily Fox

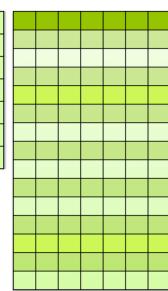
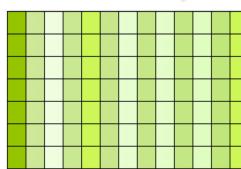
CS 229: Machine Learning

84

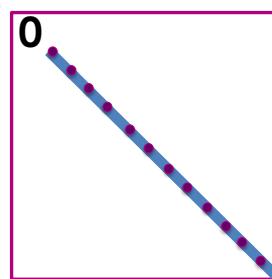
## Option 1: Don't penalize intercept

– Closed-form solution –

$$\hat{\mathbf{w}} = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I}^{\text{mod}})^{-1} \mathbf{H}^T \mathbf{y}$$



+



85

©2024 Emily Fox

CS 229: Machine Learning

85

## Option 1: Don't penalize intercept

– Gradient descent algorithm –

```

while || $\nabla$  RSS( $\mathbf{w}^{(t)}$ )|| >  $\epsilon$ 
  for  $j=0, \dots, D$ 
     $\text{partial}[j] = -2 \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)}))$ 
    if  $j == 0$ 
       $\mathbf{w}_0^{(t+1)} \leftarrow \mathbf{w}_0^{(t)} - \eta \text{partial}[j]$ 
    else
       $\mathbf{w}_j^{(t+1)} \leftarrow (1 - 2\eta\lambda)\mathbf{w}_j^{(t)} - \eta \text{partial}[j]$ 
     $t \leftarrow t + 1$ 
  
```

86

©2024 Emily Fox

CS 229: Machine Learning

86

## Option 2: Center data first

If data are first **centered about 0**, then favoring small intercept not so worrisome

**Step 1:** Transform  $y$  to have 0 mean

**Step 2:** Run ridge regression as normal  
(closed-form or gradient algorithms)

87

©2024 Emily Fox

CS 229: Machine Learning

87

## Feature normalization

**PRACTICALITIES**

©2024 Emily Fox

CS 229: Machine Learning

88

# Normalizing features

Scale training columns (**not rows!**) as:

$$\underline{h_j(x_k)} = \frac{h_j(\mathbf{x}_k)}{\sqrt{\sum_{i=1}^N h_j(\mathbf{x}_i)^2}}$$

Normalizer:  
 $Z_j$

Apply same training scale factors to test data:

$$\underline{h_j(x_k)} = \frac{h_j(\mathbf{x}_k)}{\sqrt{\sum_{i=1}^N h_j(\mathbf{x}_i)^2}}$$

Normalizer:  
 $Z_j$

*apply to  
test point*

*summing over training points*



89

©2024 Emily Fox

CS 229: Machine Learning