

# Regression recap

CS 229: Machine Learning

Emily Fox

Stanford University

January 10, 2024

©2024 Emily Fox

1

## How much is my house worth?



©2024 Emily Fox

CS 229: Machine Learning

2

1

# How much is my house worth?



©2024 Emily Fox

CS 229: Machine Learning

3

## Data



*input*      *output*  
 $(x_1 = \text{sq.ft.}, y_1 = \$)$



$(x_2 = \text{sq.ft.}, y_2 = \$)$



$(x_3 = \text{sq.ft.}, y_3 = \$)$



$(x_4 = \text{sq.ft.}, y_4 = \$)$



$(x_5 = \text{sq.ft.}, y_5 = \$)$

⋮

### Input vs. Output:

- $y$  is the quantity of interest
- assume  $y$  can be predicted from  $x$

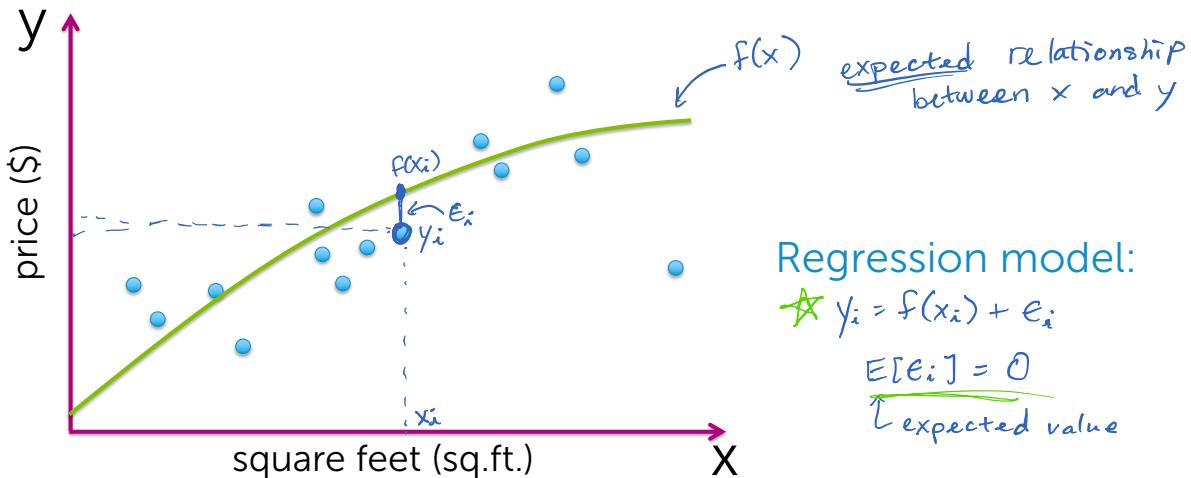
Example of supervised learning

©2024 Emily Fox

CS 229: Machine Learning

4

## Model – How we assume the world works



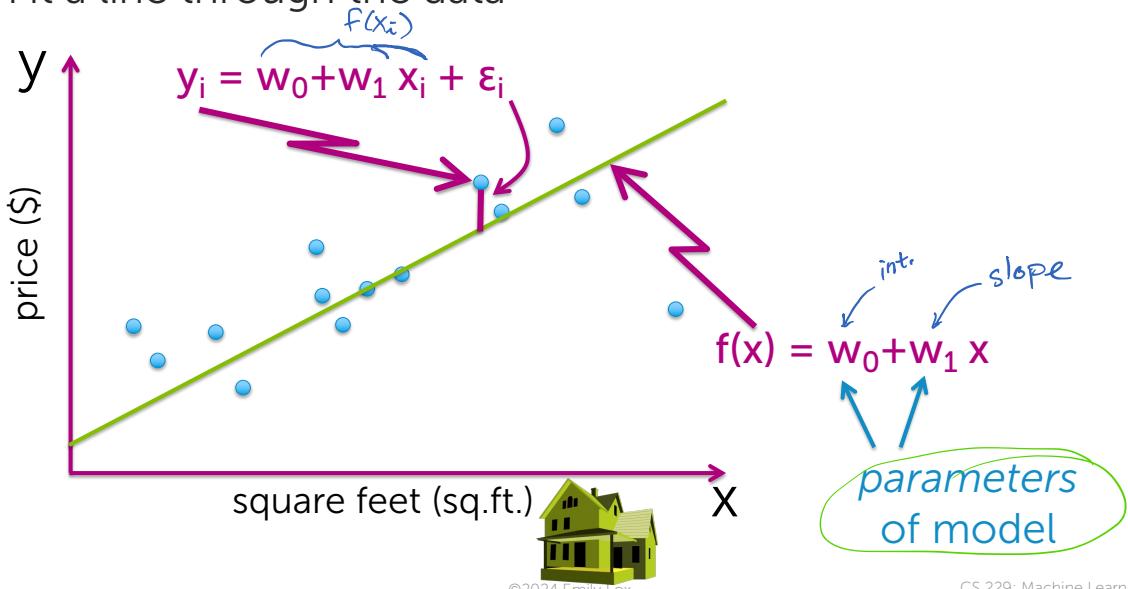
©2024 Emily Fox

CS 229: Machine Learning

5

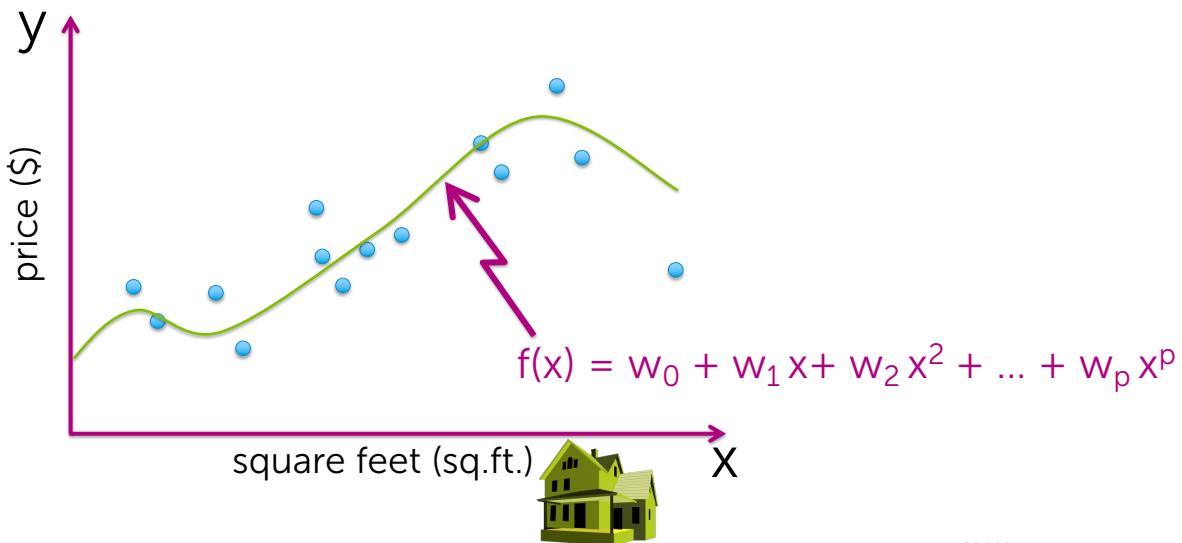
## Use a simple **linear** regression model

Fit a line through the data



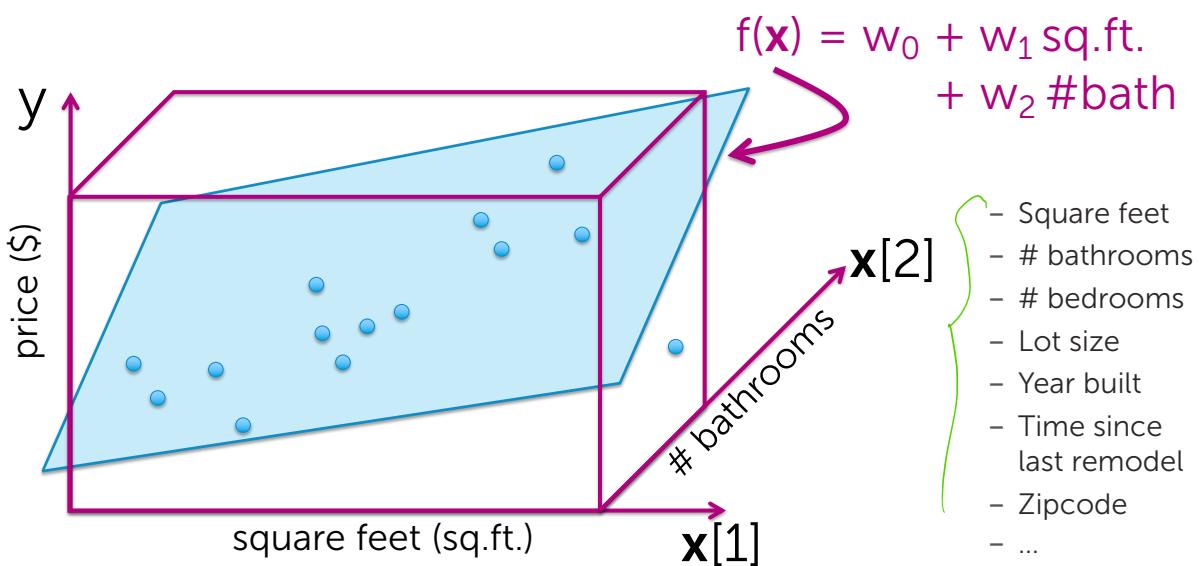
6

## Polynomial regression (still “linear” regression)



7

## Considering higher-dim input spaces



8

## Generic linear regression model

Model:

$$y_i = \mathbf{w}_0 h_0(\mathbf{x}_i) + \mathbf{w}_1 h_1(\mathbf{x}_i) + \dots + \mathbf{w}_D h_D(\mathbf{x}_i) + \varepsilon_i$$

$$= \sum_{j=0}^D \mathbf{w}_j h_j(\mathbf{x}_i) + \varepsilon_i = \mathbf{w}^\top \mathbf{h}(\mathbf{x}_i) + \varepsilon_i$$

feature 1 =  $h_0(\mathbf{x})$  ... e.g., 1

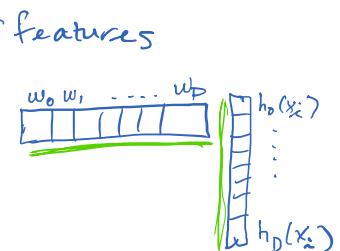
feature 2 =  $h_1(\mathbf{x})$  ... e.g.,  $\mathbf{x}[1]$  = sq. ft.

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

or,  $\log(\mathbf{x}[7]) \mathbf{x}[2] = \log(\#bed) \times \#bath$

...

feature  $D+1 = h_D(\mathbf{x})$  ... some other function of  $\mathbf{x}[1], \dots, \mathbf{x}[d]$



©2024 Emily Fox

CS 229: Machine Learning

9

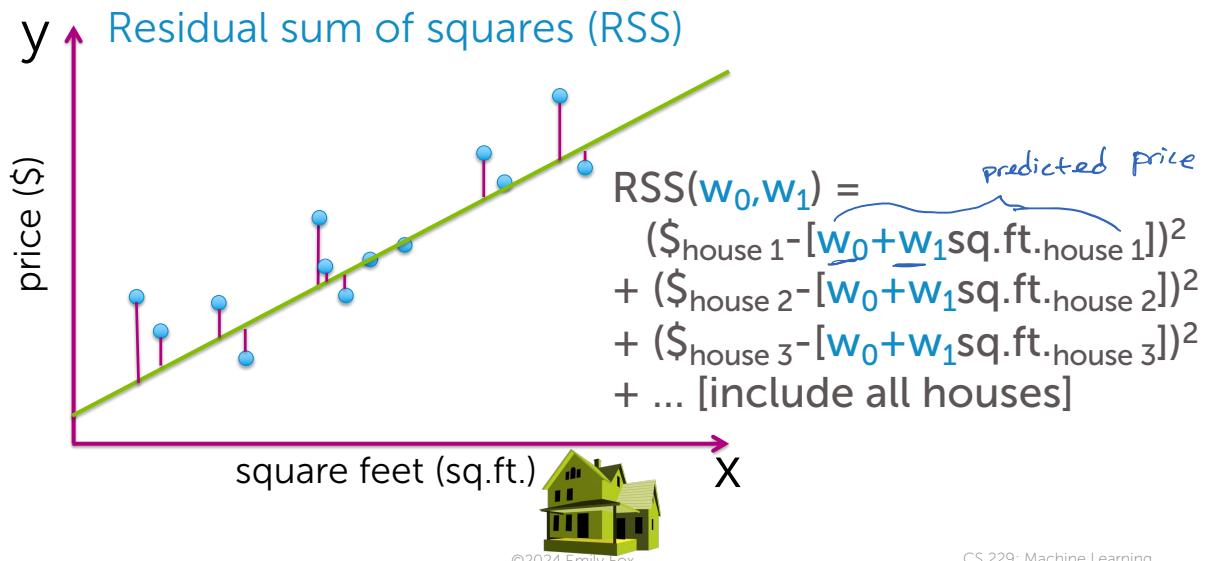
## Fitting the linear regression model

©2024 Emily Fox

CS 229: Machine Learning

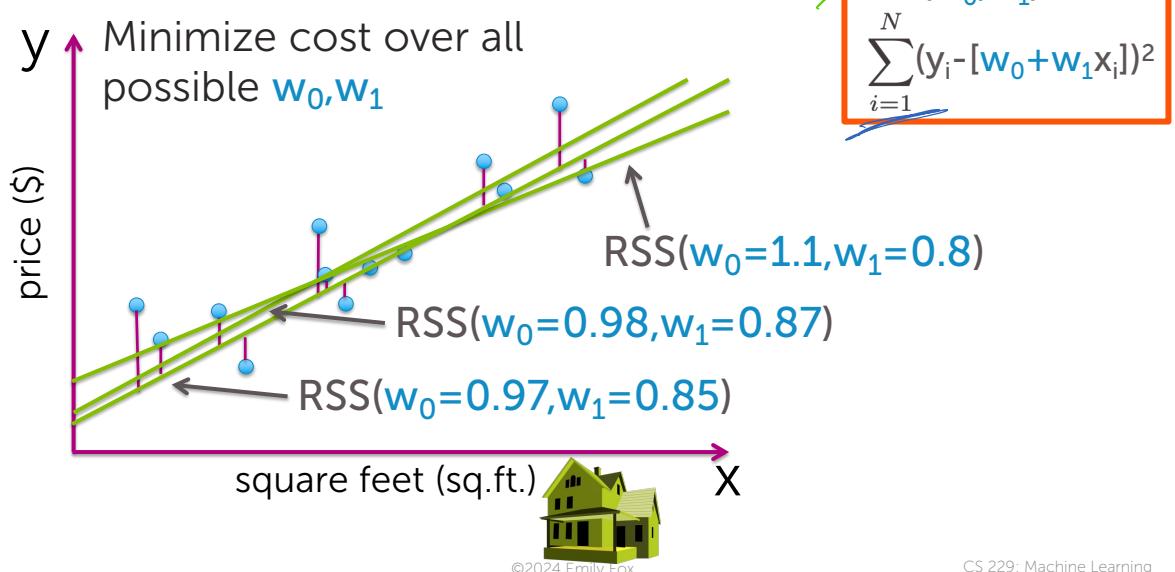
10

## "Cost" of using a given line (Loss function)



11

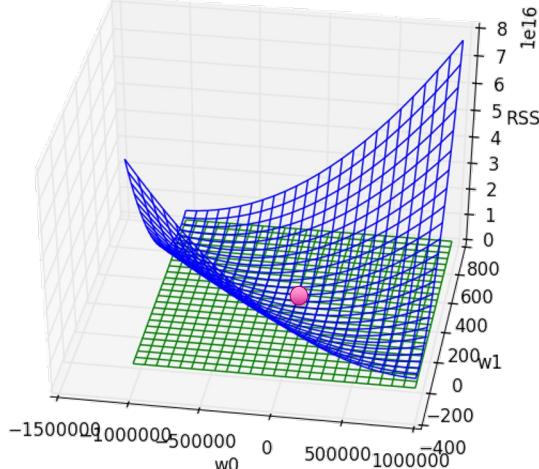
## Find "best" line



12

## Minimizing the cost

3D plot of RSS with tangent plane at minimum



Minimize function  
over all possible  $w_0, w_1$

$$\min_{w_0, w_1} \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

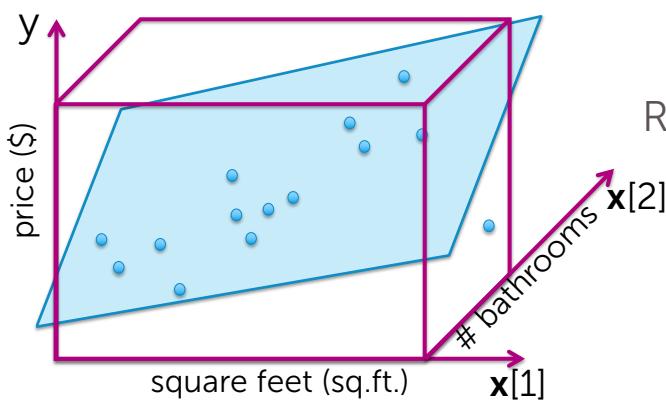
RSS( $w_0, w_1$ ) is a function  
of 2 variables

©2024 Emily Fox

CS 229: Machine Learning

13

## RSS for multiple regression



$$\begin{array}{c} y \\ h \\ x \end{array} = \begin{array}{c} H \\ \vdots \\ h(x) \end{array} + \begin{array}{c} \omega \\ \vdots \\ \omega_p \end{array} + \begin{array}{c} \epsilon \\ \vdots \\ \epsilon_n \end{array}$$

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

*N actual obs*

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

*$\mathbf{y}$*

show this at home

©2024 Emily Fox

CS 229: Machine Learning

14

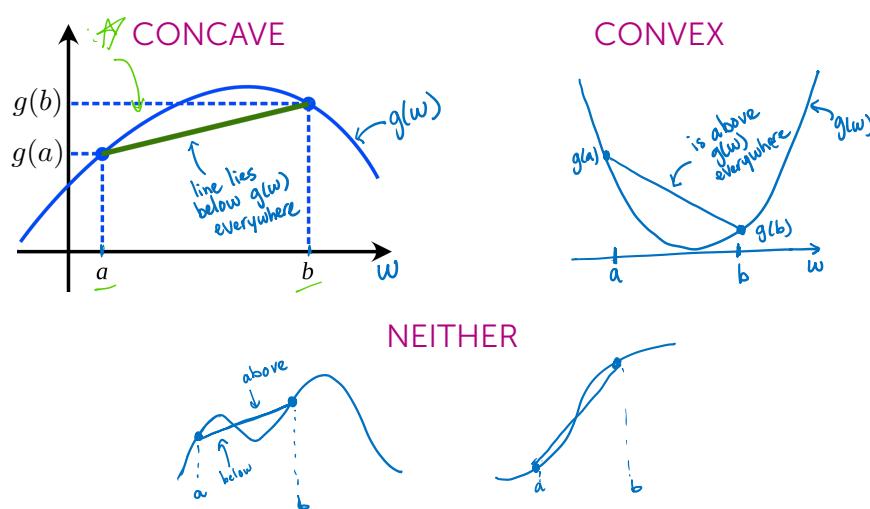
## Background on optimization

©2024 Emily Fox

CS 229: Machine Learning

15

## Convex/concave functions

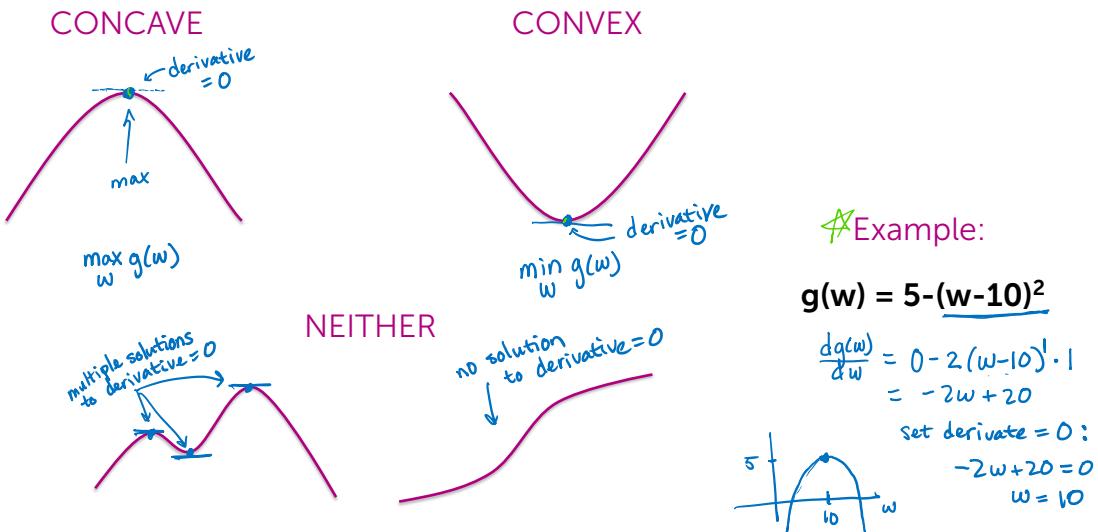


©2024 Emily Fox

CS 229: Machine Learning

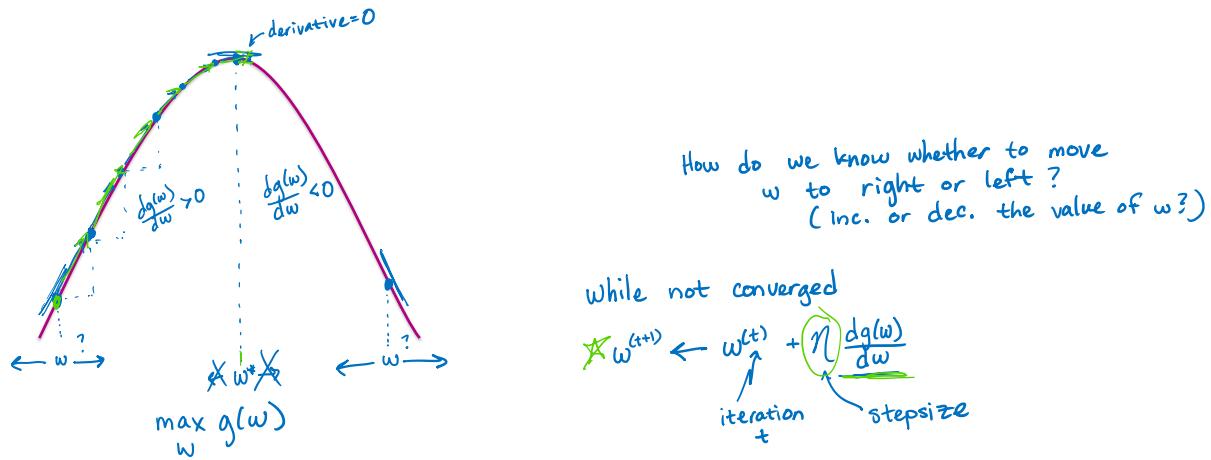
16

## Finding the max or min analytically



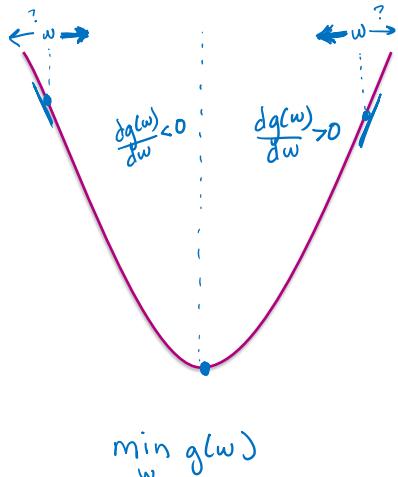
17

## Finding the max via gradient ascent ("hill climbing")



18

## Finding the min via gradient descent ("hill descent")



when derivative is positive, we want to decrease  $w$   
and when derivative is negative, we want to increase  $w$

Algorithm:

**while** not converged

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{dg}{dw} \Big|_{w^{(t)}}$$

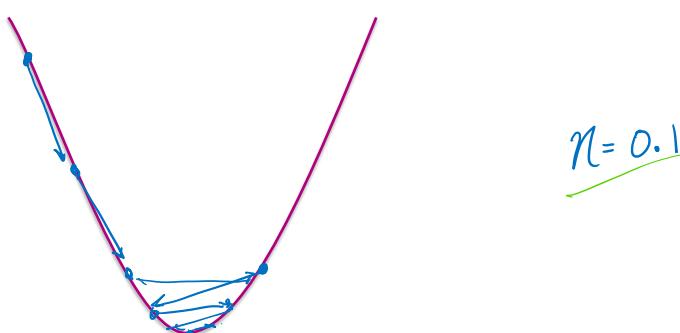
©2024 Emily Fox

CS 229: Machine Learning

19

## Choosing the stepsize— Fixed stepsize

$\eta$

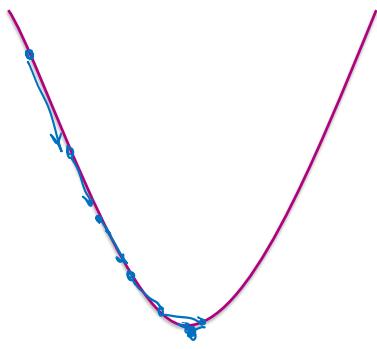


©2024 Emily Fox

CS 229: Machine Learning

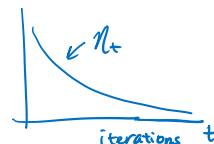
20

## Choosing the stepsize— Decreasing stepsize or stepsize schedule



Common choices:

$$\begin{cases} \eta_t = \frac{\alpha}{t} \\ \eta_t = \frac{\alpha}{\sqrt{t}} \end{cases}$$



©2024 Emily Fox

CS 229: Machine Learning

21

## Convergence criteria

For convex functions,  
optimum occurs when

$$\frac{dg(w)}{dw} = 0$$

In practice, stop when

$$\left| \frac{dg(w)}{dw} \right| < \epsilon$$

↑ threshold to be set

Algorithm:

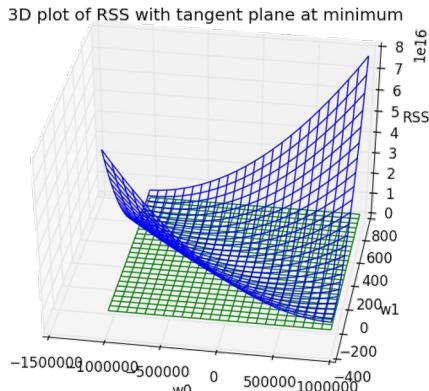
**while** not converged  
 $w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{dg}{dw} \Big|_{w^{(t)}}$

©2024 Emily Fox

CS 229: Machine Learning

22

# Moving to multiple dimensions: Gradients



$$\nabla g(\mathbf{w}) = \begin{bmatrix} \frac{\partial g}{\partial w_0} \\ \frac{\partial g}{\partial w_1} \\ \vdots \\ \frac{\partial g}{\partial w_p} \end{bmatrix}$$

gradient

$[w_0, w_1, \dots, w_p]$

$(p+1)$ -dimensional vector

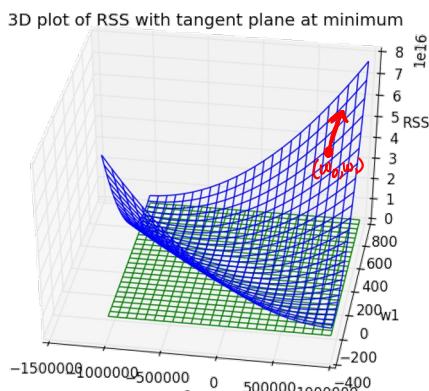
partial derivative is like a derivate with respect to  $w_i$ , treating all other variables as constants

©2024 Emily Fox

CS 229: Machine Learning

23

## Gradient example



$$g(\mathbf{w}) = 5w_0 + 10\underline{w_0 w_1} + 2w_1^2$$

$$\frac{\partial g}{\partial w_0} = 5 + 10w_1$$

$$\frac{\partial g}{\partial w_1} = 10w_0 + 4w_1$$

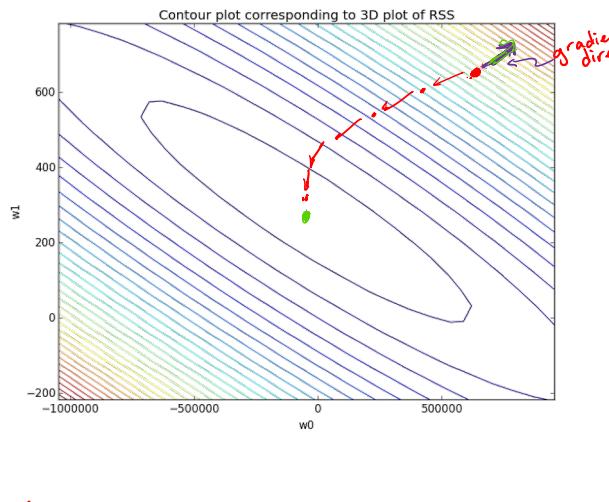
$$\nabla g(\mathbf{w}) = \begin{bmatrix} 5 + 10w_1 \\ 10w_0 + 4w_1 \end{bmatrix}$$

©2024 Emily Fox

CS 229: Machine Learning

24

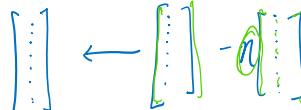
# Gradient descent



Algorithm:

**while** not converged

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



Convergence:  
 $\|\nabla g(\mathbf{w})\| < \epsilon$

©2024 Emily Fox

CS 229: Machine Learning

25

**Back to regression:**  
 Take the gradient of RSS with respect to  $\mathbf{w}$

©2024 Emily Fox

CS 229: Machine Learning

26

## Gradient of RSS

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_{11}} \\ \frac{\partial f}{\partial x_{12}} \\ \vdots \\ \frac{\partial f}{\partial x_{1n}} \end{bmatrix}$$

vector of partials

$$\begin{aligned}\nabla \text{RSS}(\mathbf{w}) &= \nabla[(\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w})] \\ &= -2\mathbf{H}^\top(\mathbf{y} - \mathbf{H}\mathbf{w}) \quad \leftarrow \text{vector of dim } D+1 \quad (\text{length } w)\end{aligned}$$

Why? By analogy to 1D case:

$$\frac{d}{dw} (\mathbf{y} - \mathbf{h}\mathbf{w})(\mathbf{y} - \mathbf{h}\mathbf{w}) = \frac{d}{dw} (\mathbf{y} - \mathbf{h}\mathbf{w})^2 = 2(\mathbf{y} - \mathbf{h}\mathbf{w})(-\mathbf{h}) = -2\mathbf{h}(\mathbf{y} - \mathbf{h}\mathbf{w})$$

↑↑  
scalars

©2024 Emily Fox

CS 229: Machine Learning

27

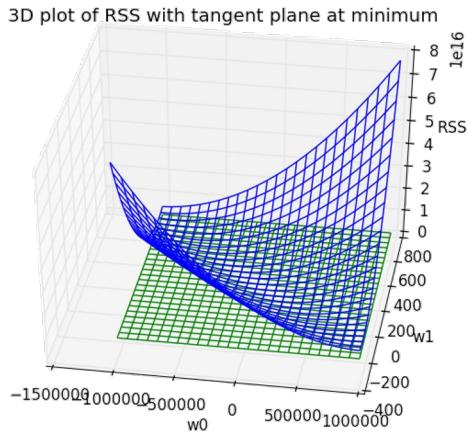
**Approach 1:**  
Set the gradient = 0

©2024 Emily Fox

CS 229: Machine Learning

28

## Closed-form solution



$$\nabla \text{RSS}(\mathbf{w}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) = 0$$

Solve for  $\mathbf{w}$ :

$$\begin{aligned} -2\mathbf{H}^T\mathbf{y} + 2\mathbf{H}^T\mathbf{H}\hat{\mathbf{w}} &= 0 \\ \mathbf{H}^T\mathbf{H}\hat{\mathbf{w}} &= \mathbf{H}^T\mathbf{y} \\ \hat{\mathbf{w}} &= (\mathbf{H}^T\mathbf{H})^{-1}\mathbf{H}^T\mathbf{y} \end{aligned}$$

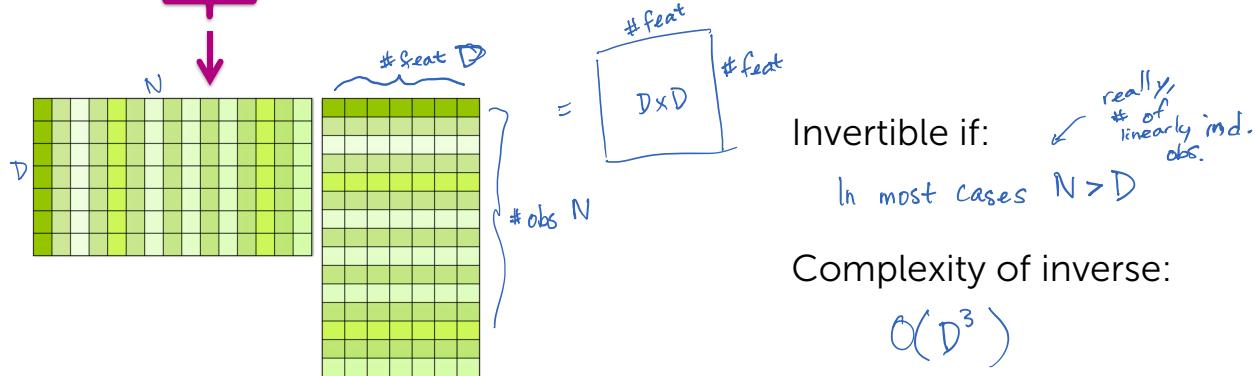
©2024 Emily Fox

CS 229: Machine Learning

29

## Closed-form solution

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



©2024 Emily Fox

CS 229: Machine Learning

30

## Approach 2: Gradient descent

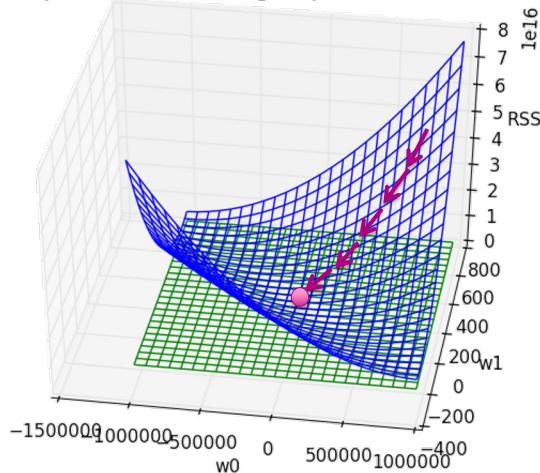
©2024 Emily Fox

CS 229: Machine Learning

31

## Gradient descent

3D plot of RSS with tangent plane at minimum



### Algorithm:

**while** not converged

$$\begin{aligned}
 w^{(t+1)} &\leftarrow w^{(t)} - \eta \nabla_{w^{(t)}} \text{RSS}(w^{(t)}) \\
 &\quad - 2H^T(y - Hw^{(t)}) \\
 &\leftarrow w^{(t)} + 2\eta H^T(y - Hw^{(t)}) \\
 &\quad \downarrow \hat{y}(w^{(t)}) \\
 &\quad \text{best guess} \\
 &\quad \text{of obs.} \\
 &\quad \text{using } w^{(t)}
 \end{aligned}$$

©2024 Emily Fox

CS 229: Machine Learning

32

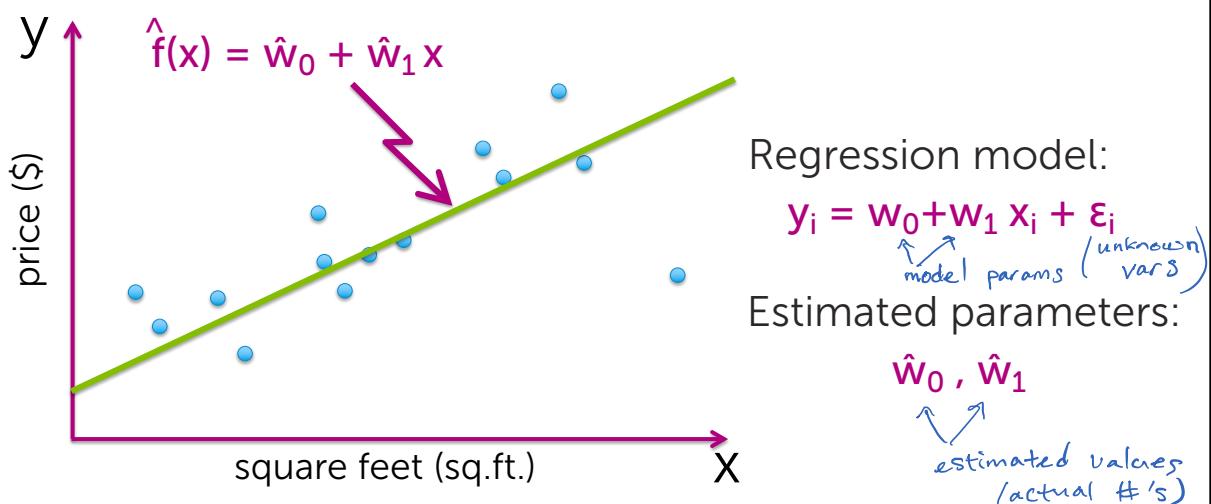
## Using the fitted line

©2024 Emily Fox

CS 229: Machine Learning

33

## Model vs. fitted line

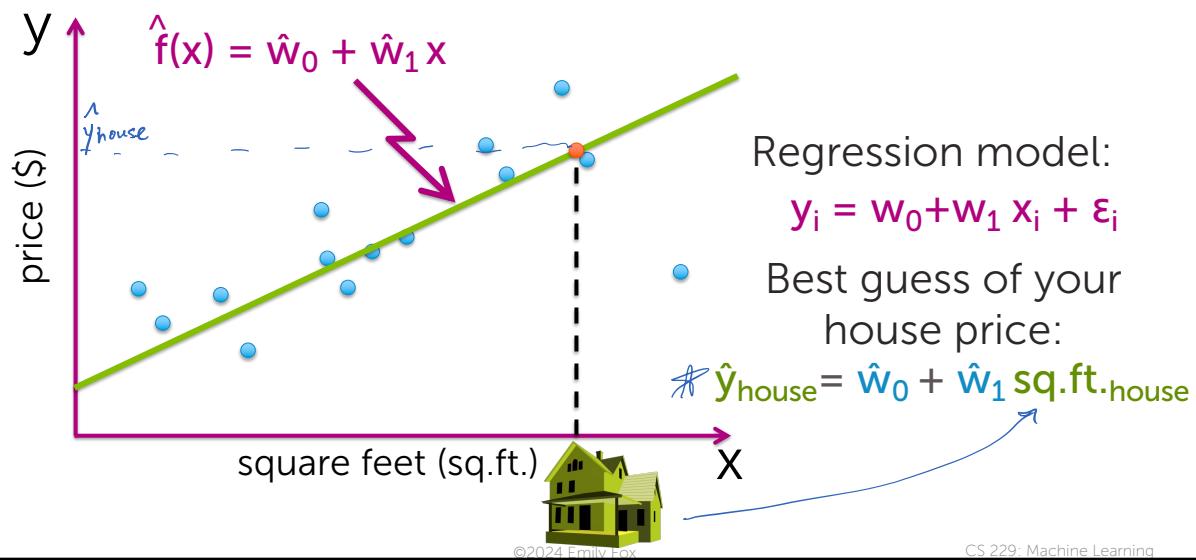


©2024 Emily Fox

CS 229: Machine Learning

34

## Seller: Predicting your house price

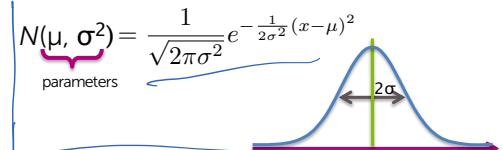
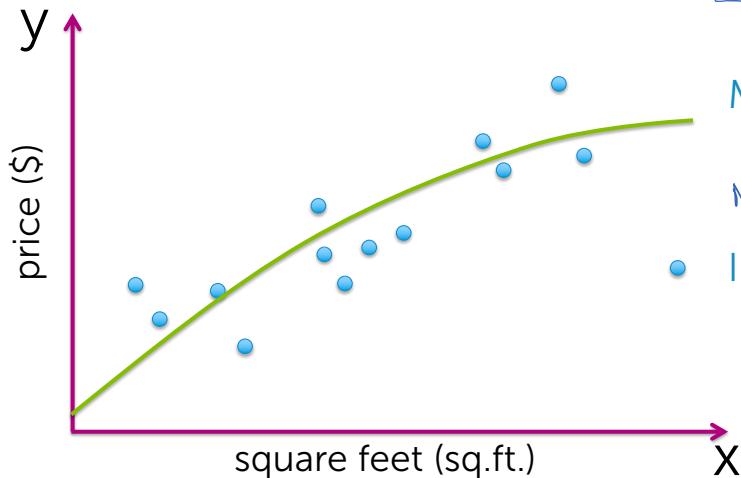


35

Why min RSS?

36

## Assuming Gaussian noise



Model for  $\epsilon_i$ :

$$\text{So far : } E[\epsilon_i] = 0$$

Now, stronger assumption:  
 $\epsilon_i \sim N(0, \sigma^2)$

- Implied distribution on  $y_i$ :

$$y_i = h^T(x_i) w + \epsilon_i$$

$\underbrace{h^T(x_i)}$  "given"  
 $\underbrace{w}$  const

$$y_i | x_i; w \sim N(h^T(x_i) w, \sigma^2)$$

©2024 Emily Fox

CS 229: Machine Learning

37

## Maximum likelihood estimate of params

$$\arg \max_w p(y | X; w) \stackrel{\epsilon_i \text{ iid}}{=} \arg \max_w \prod_{i=1}^N p(y_i | x_i; w) = \arg \max_w \ln \left\{ \prod_{i=1}^N p(y_i | x_i; w) \right\}$$

↑ strictly inc. func

Maximize log-likelihood wrt  $w$

$$\begin{aligned} \arg \max_w \ln p(y | X; w) &= \ln \left\{ \left( \frac{1}{\sigma \sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(y_i - h^T(x_i) w)^2}{2\sigma^2}} \right\} \\ &= \arg \max_w N \ln \left( \frac{1}{\sigma \sqrt{2\pi}} \right) - \sum_{i=1}^N \frac{(y_i - h^T(x_i) w)^2}{2\sigma^2} \\ &= \arg \min_w \sum_{i=1}^N (y_i - h^T(x_i) w)^2 \quad \leftarrow \text{MLE for Gauss model} = \min \text{ RSS! (LS line)} \end{aligned}$$

©2024 Emily Fox

CS 229: Machine Learning

38

## Summary of linear regression

©2024 Emily Fox

CS 229: Machine Learning

39

## What you can do now...

- Describe the input (features) and output (real-valued predictions) of a regression model
- Write a regression model using multiple inputs or features thereof
- Cast both polynomial regression and regression with multiple inputs as regression with multiple features
- Calculate a cost metric (e.g., RSS)
- Estimate model parameters of a general multiple regression model to minimize RSS:
  - In closed form
  - Using an iterative gradient descent algorithm
- Exploit the estimated model to form predictions
- Relate minimizing RSS to maximizing likelihood of a Gaussian model

©2024 Emily Fox

CS 229: Machine Learning

40

# Assessing performance, Bias-variance tradeoff

CS 229: Machine Learning

Emily Fox

Stanford University

January 10, 2024

©2024 Emily Fox

41

## Assessing performance

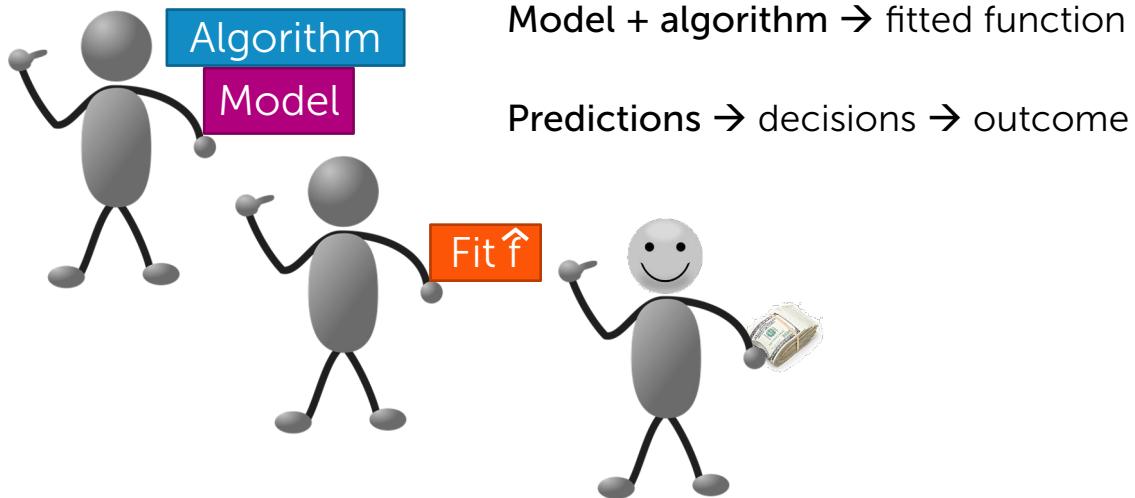
©2024 Emily Fox

CS 229: Machine Learning

42

21

## Make predictions, get \$, right??



43

©2024 Emily Fox

CS 229: Machine Learning

43

## Or, how much am I losing?

**Example:** Lost \$ due to inaccurate listing price

- Too low → low offers
- Too high → few lookers + no/low offers

How much am I **losing** compared to perfection?

**Perfect predictions:** Loss = 0

**My predictions:** Loss = ???

44

©2024 Emily Fox

CS 229: Machine Learning

44

## Measuring loss

Loss function:

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

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

actual value

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

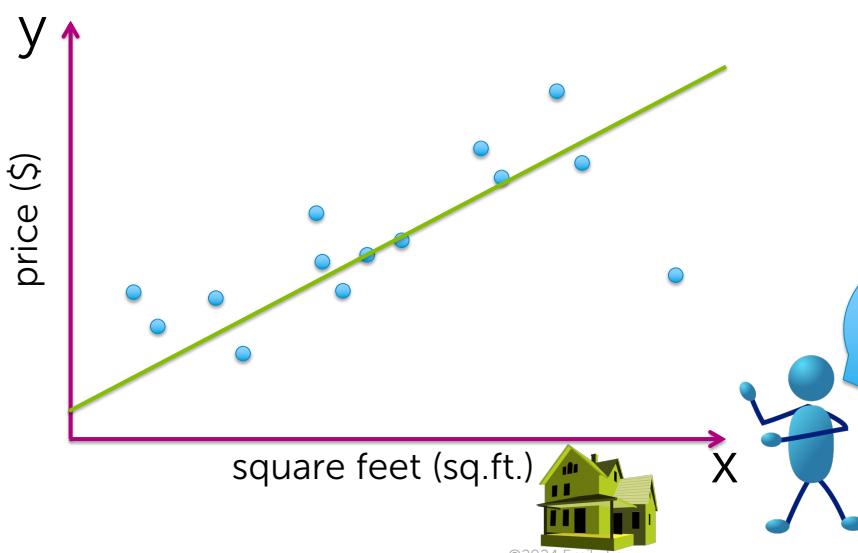
45

@2024 Emily Fox

CS 229: Machine Learning

45

## Fit data with a line or ... ?



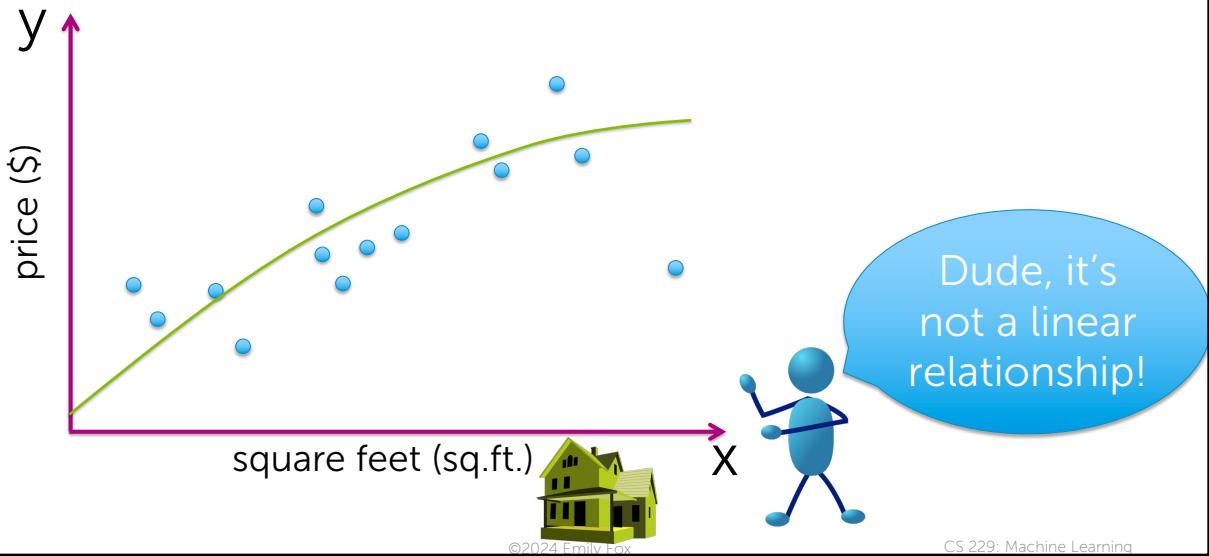
46

@2024 Emily Fox

CS 229: Machine Learning

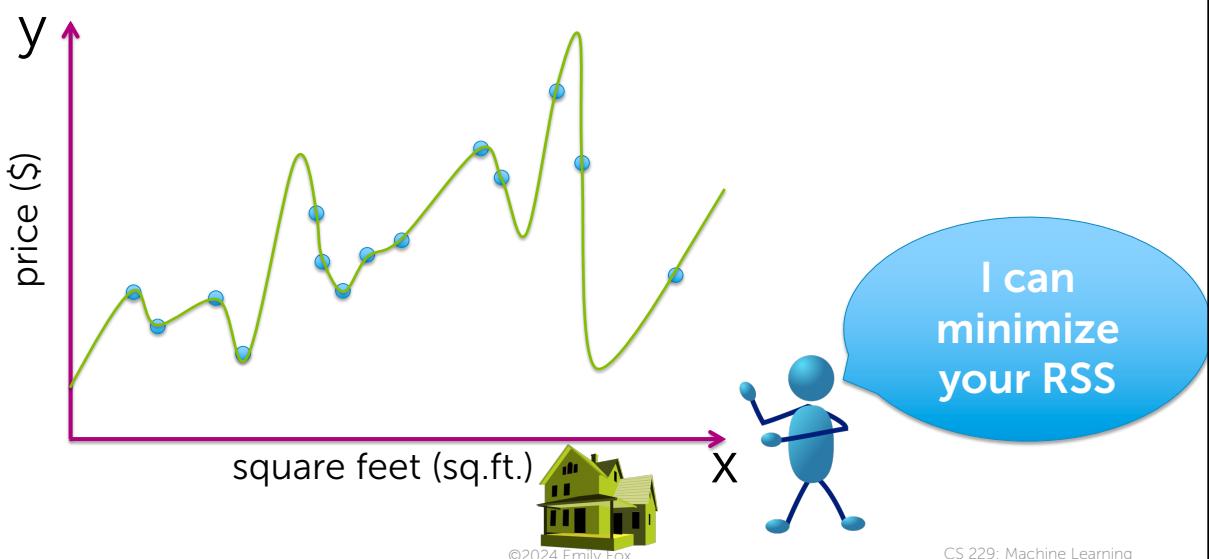
46

## What about a quadratic function?



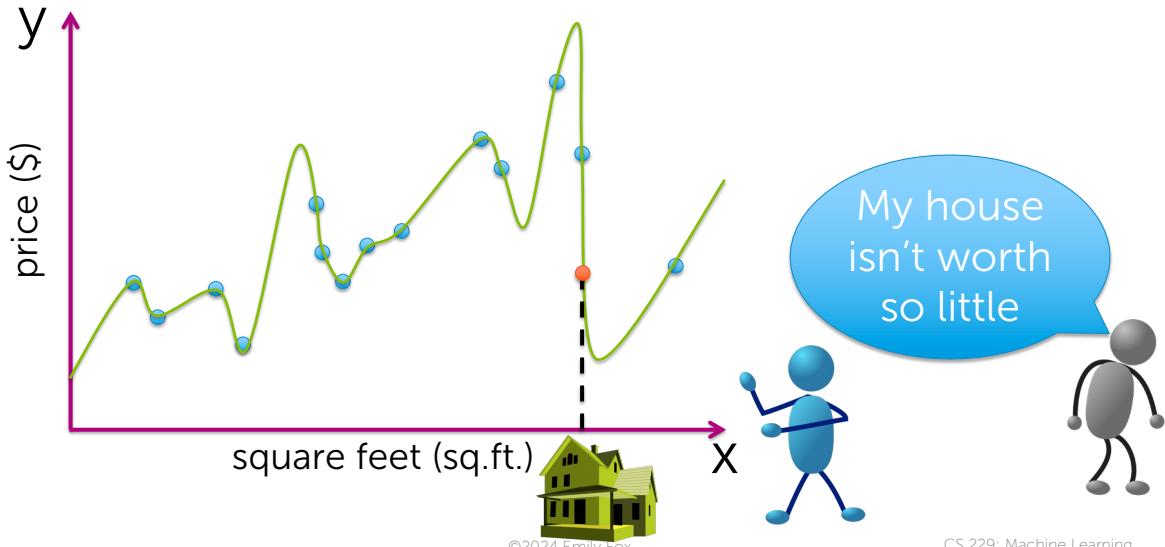
47

## Even higher order polynomial



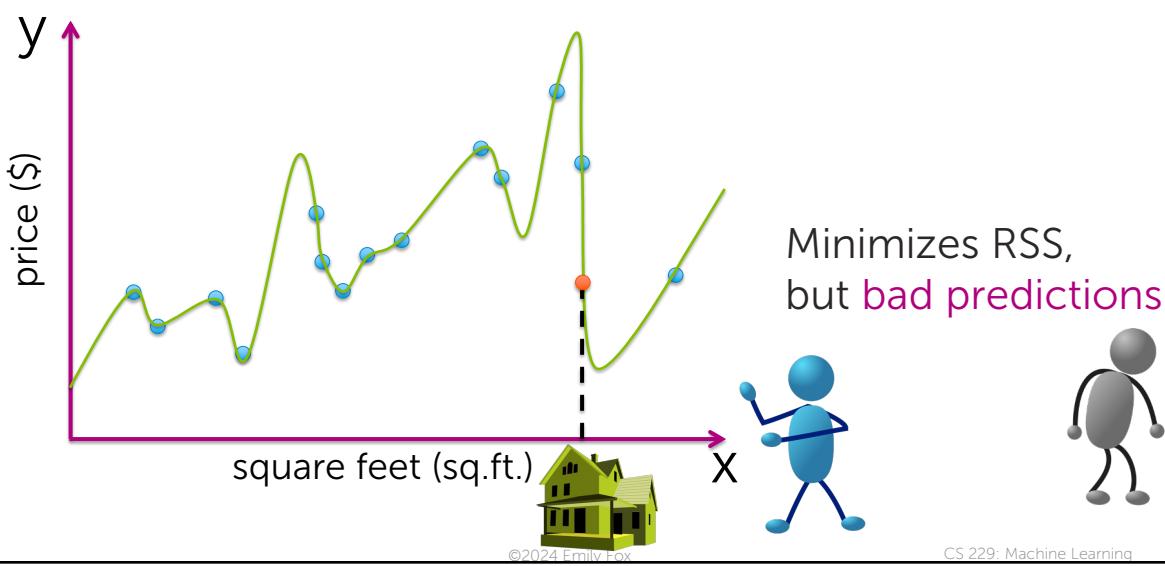
48

## Do you believe this fit?



49

## Do you believe this fit?



50

“Remember that all models are wrong; the practical question is how wrong do they have to be to not be useful.” George Box, 1987.

©2024 Emily Fox

CS 229: Machine Learning

51

## Assessing the loss

©2024 Emily Fox

CS 229: Machine Learning

52

26

## Assessing the loss

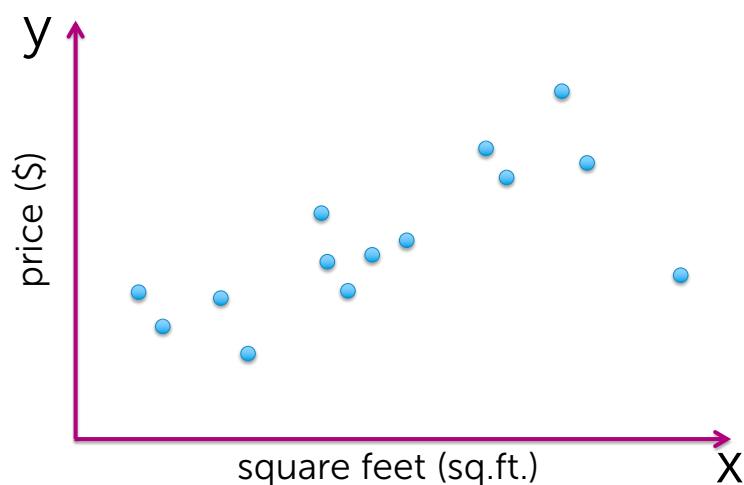
### Part 1: Training error

©2024 Emily Fox

CS 229: Machine Learning

53

## Define training data



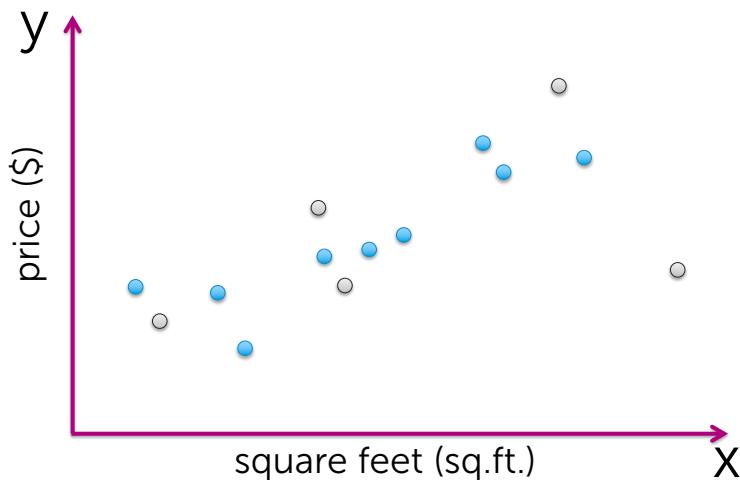
54

©2024 Emily Fox

CS 229: Machine Learning

54

## Define training data



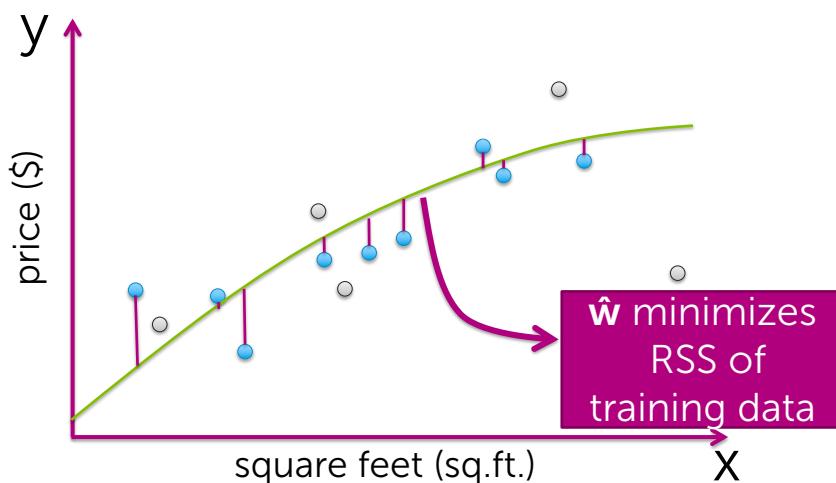
55

©2024 Emily Fox

CS 229: Machine Learning

55

**Example:**  
**Fit quadratic to minimize RSS**



56

©2024 Emily Fox

CS 229: Machine Learning

56

## Compute training error

1. Define a loss function  $L(y, f_{\hat{w}}(\mathbf{x}))$ 
  - E.g., squared error, absolute error,...

### 2. Training error

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

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

fit using training data

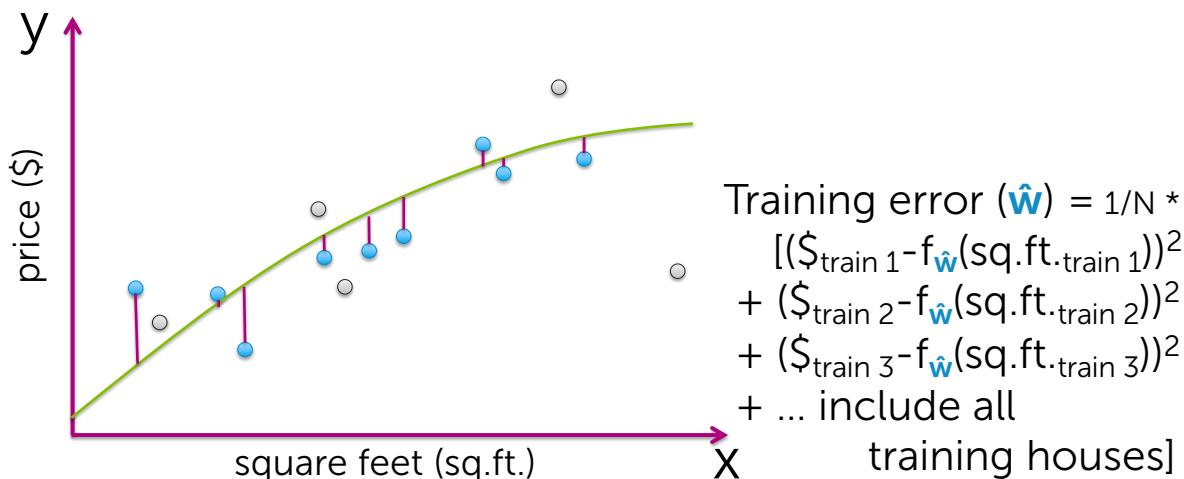
57

@2024 Emily Fox

CS 229: Machine Learning

57

**Example:**  
Use **squared error** loss  $(y - f_{\hat{w}}(\mathbf{x}))^2$



58

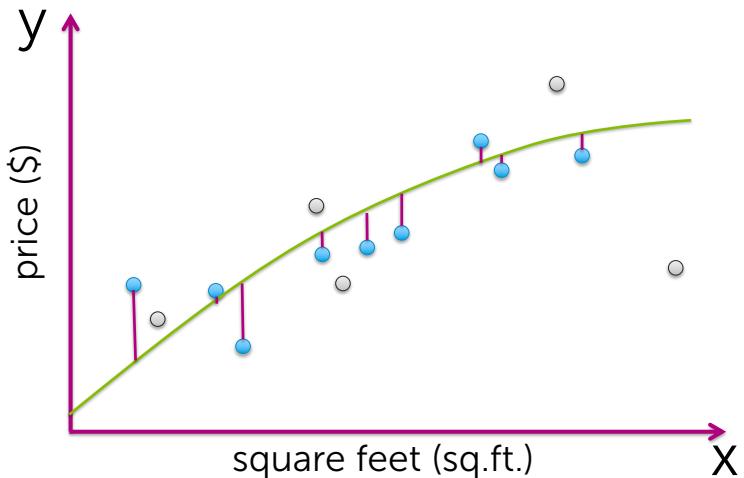
@2024 Emily Fox

CS 229: Machine Learning

58

**Example:**

Use squared error loss  $(y - f_{\hat{w}}(x))^2$



$$\text{Training error } (\hat{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - f_{\hat{w}}(\mathbf{x}_i))^2$$

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - f_{\hat{w}}(\mathbf{x}_i))^2}$$

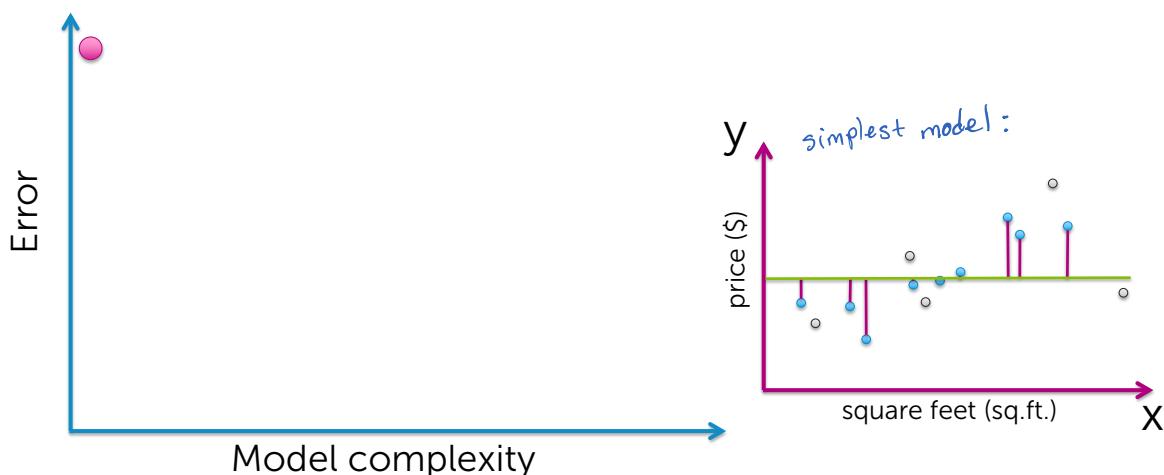
59

@2024 Emily Fox

CS 229: Machine Learning

59

## Training error vs. model complexity



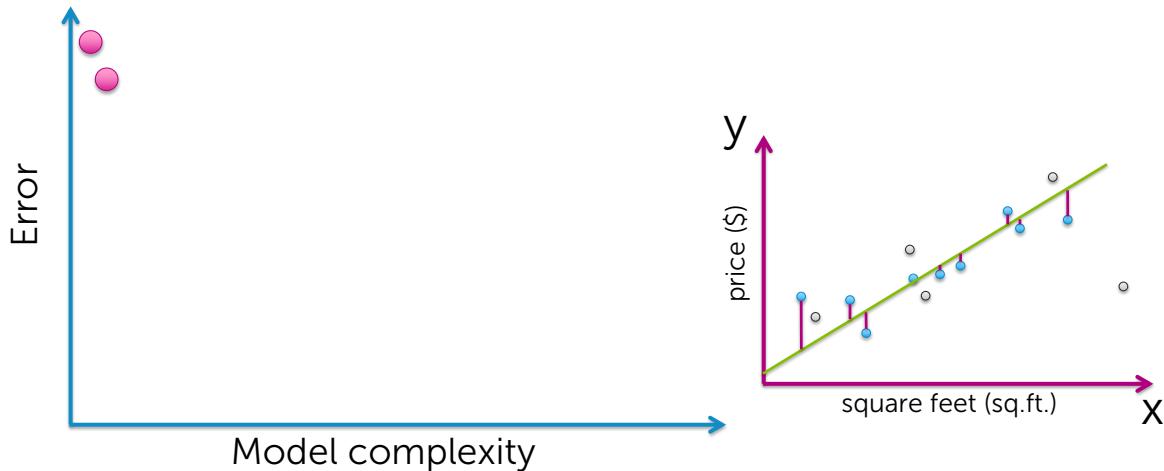
60

@2024 Emily Fox

CS 229: Machine Learning

60

## Training error vs. model complexity



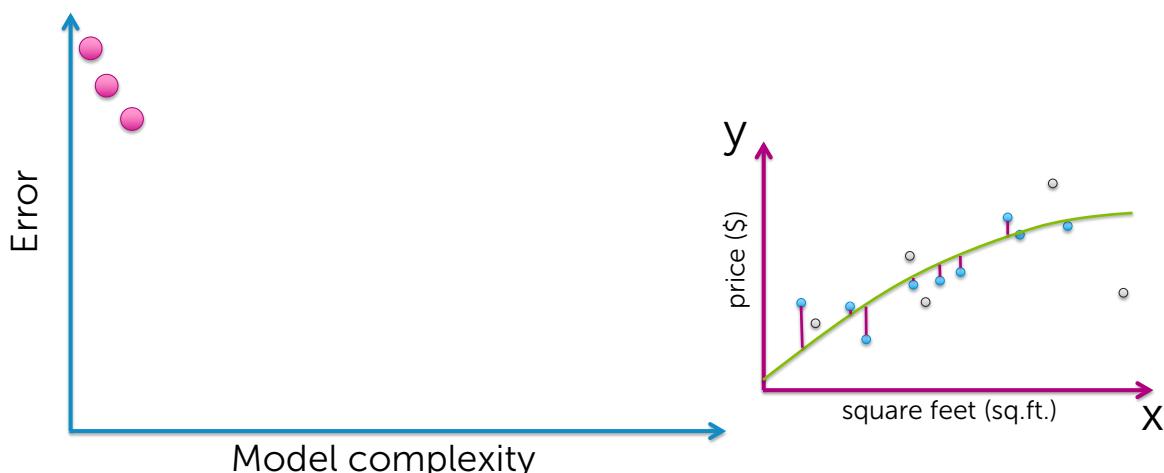
61

©2024 Emily Fox

CS 229: Machine Learning

61

## Training error vs. model complexity



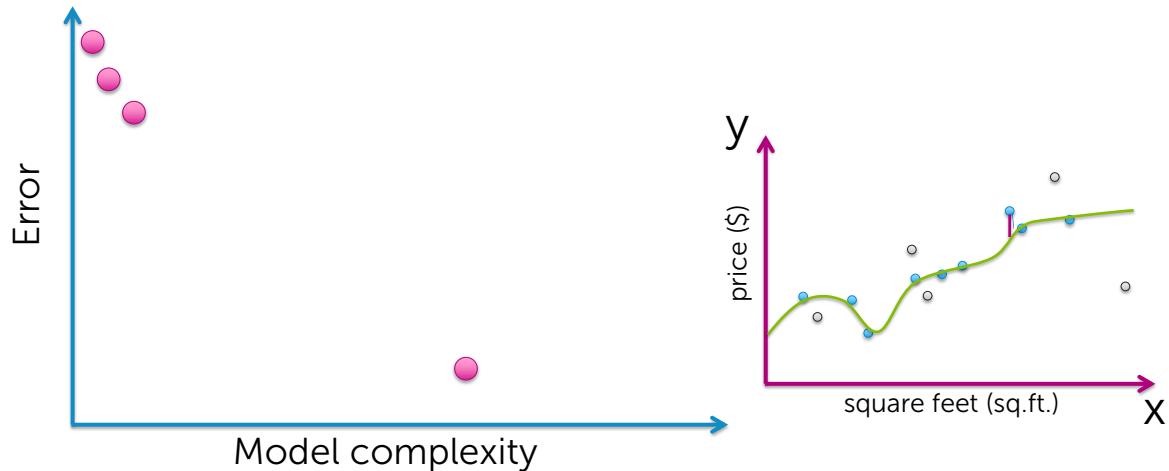
62

©2024 Emily Fox

CS 229: Machine Learning

62

## Training error vs. model complexity



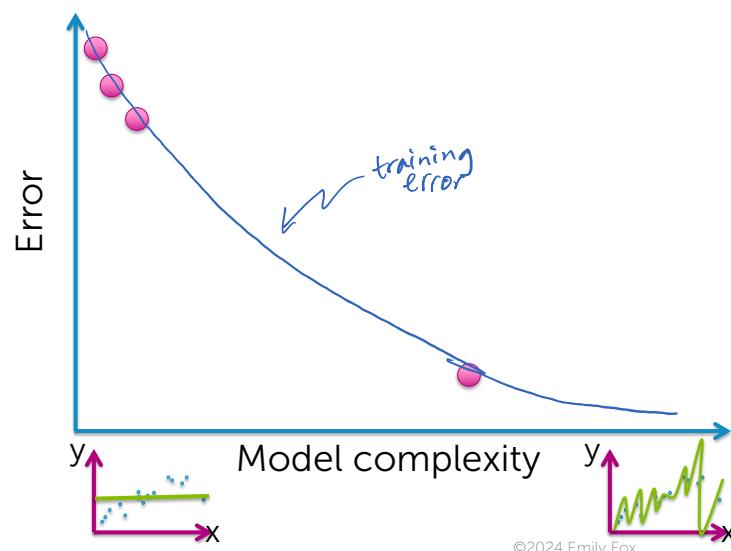
63

©2024 Emily Fox

CS 229: Machine Learning

63

## Training error vs. model complexity



64

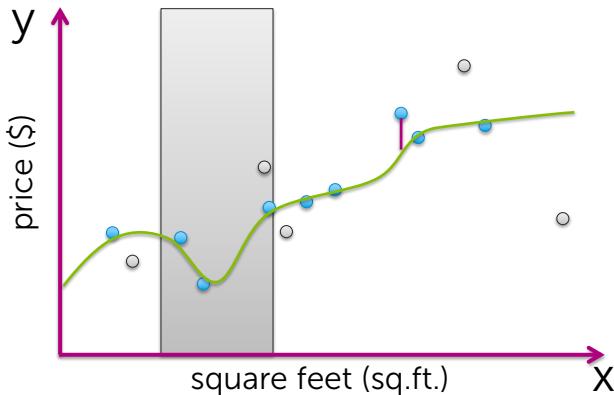
©2024 Emily Fox

CS 229: Machine Learning

64

## Is training error a good measure of predictive performance?

How do we expect to perform on a new house?



65

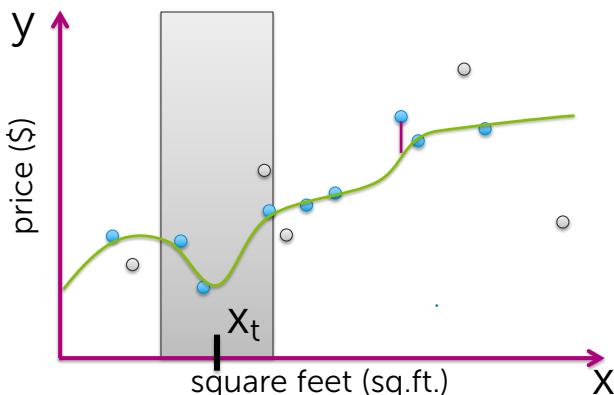
@2024 Emily Fox

CS 229: Machine Learning

65

## Is training error a good measure of predictive performance?

Is there something particularly bad about having  $x_t$  sq.ft.??



66

@2024 Emily Fox

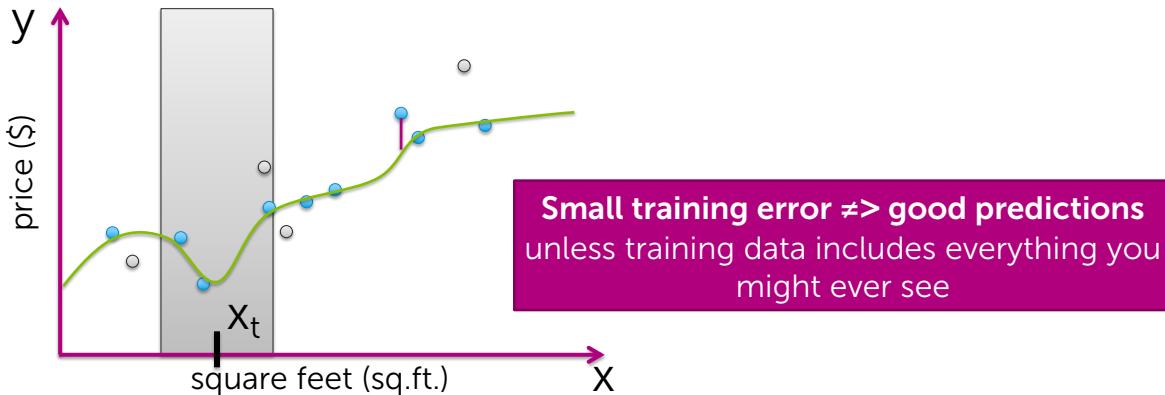
CS 229: Machine Learning

66

# Is training error a good measure of predictive performance?

Issue:

Training error is overly optimistic...  $\hat{w}$  was fit to training data



67

@2024 Emily Fox

CS 229: Machine Learning

67

Assessing the loss  
Part 2: Generalization (true) error

@2024 Emily Fox

CS 229: Machine Learning

68

## Generalization error

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



69

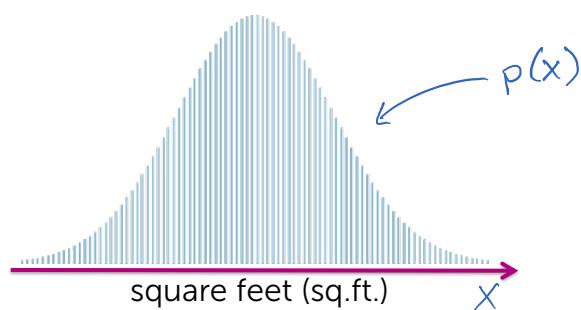
@2024 Emily Fox

CS 229: Machine Learning

69

## Distribution over houses

In our neighborhood, houses of what # sq.ft. () are we likely to see?



70

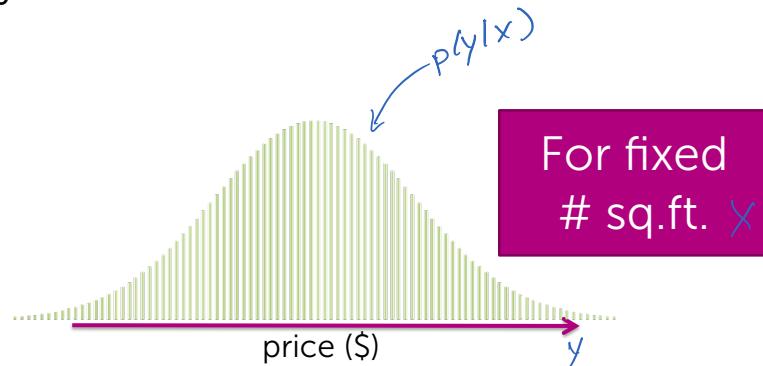
@2024 Emily Fox

CS 229: Machine Learning

70

## Distribution over sales prices

For houses with a given # sq.ft. (🏠), what house prices (\$) are we likely to see?



71

@2024 Emily Fox

CS 229: Machine Learning

71

## 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

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

$$= \int L(y, f_w(x)) p(x, y) dx dy$$

$p(x)$  fit using training data

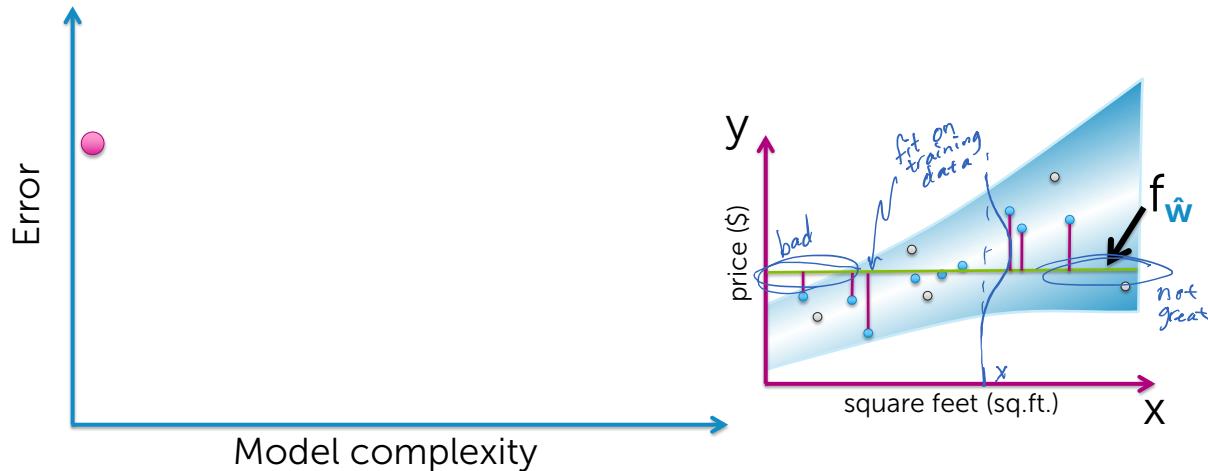
72

@2024 Emily Fox

CS 229: Machine Learning

72

## Generalization error vs. model complexity



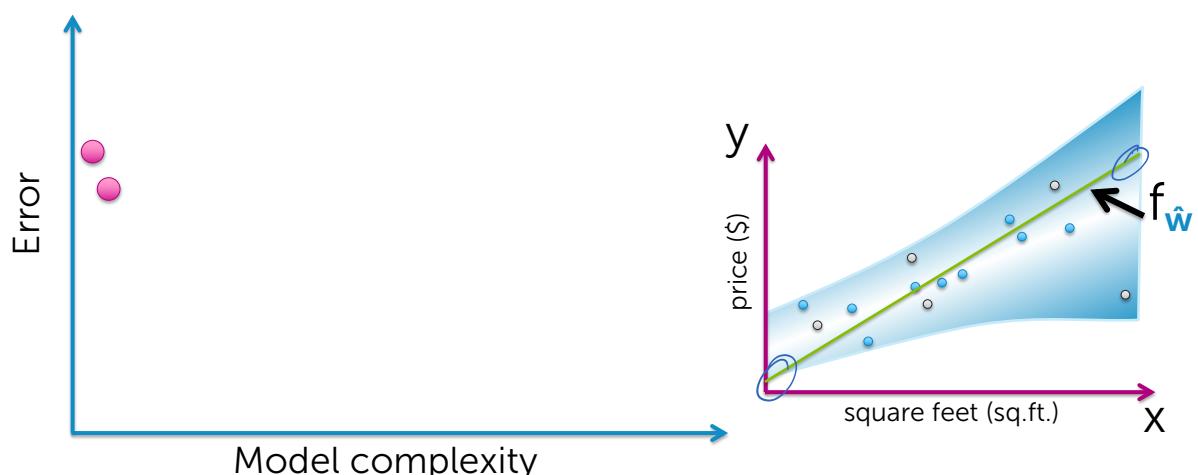
73

©2024 Emily Fox

CS 229: Machine Learning

73

## Generalization error vs. model complexity



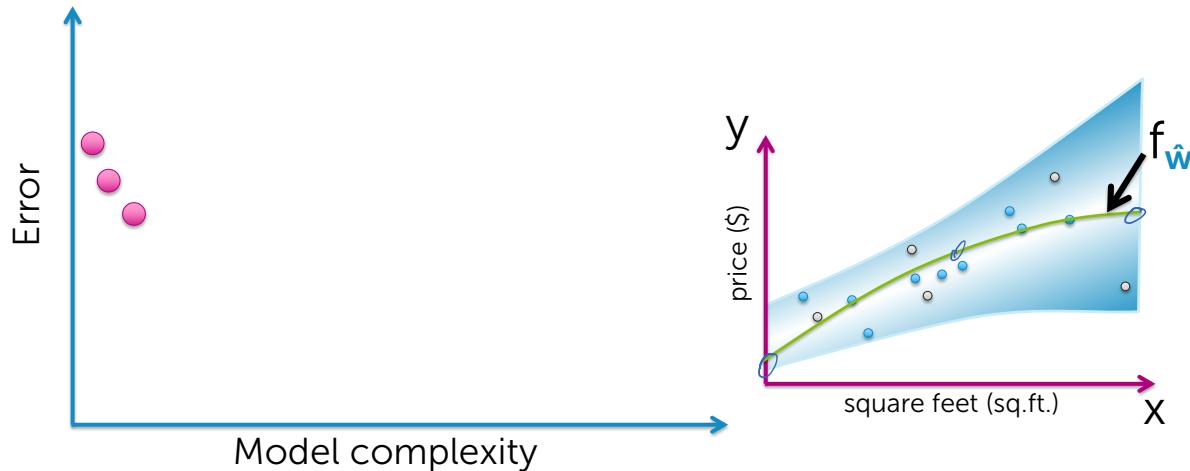
74

©2024 Emily Fox

CS 229: Machine Learning

74

## Generalization error vs. model complexity



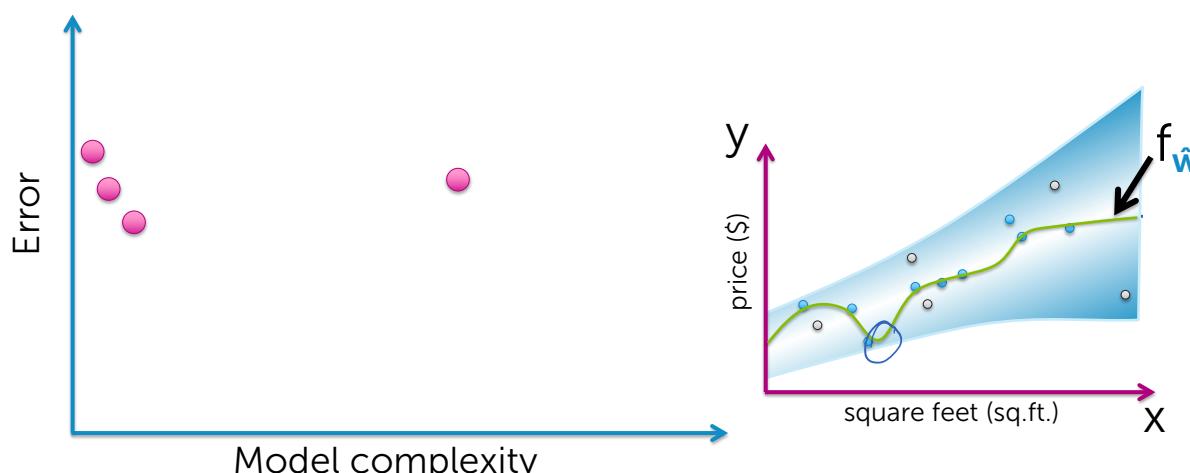
75

©2024 Emily Fox

CS 229: Machine Learning

75

## Generalization error vs. model complexity



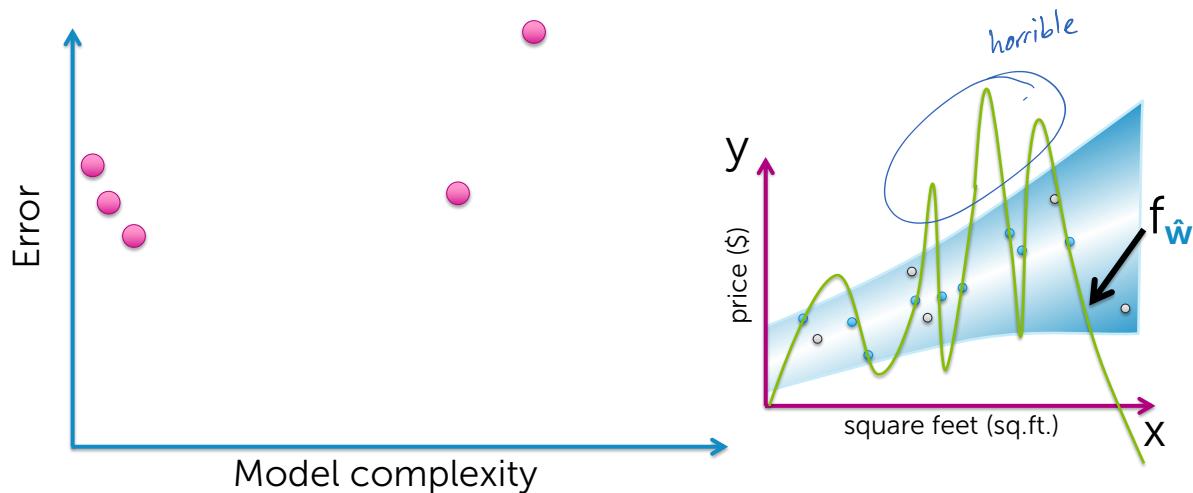
76

©2024 Emily Fox

CS 229: Machine Learning

76

## Generalization error vs. model complexity



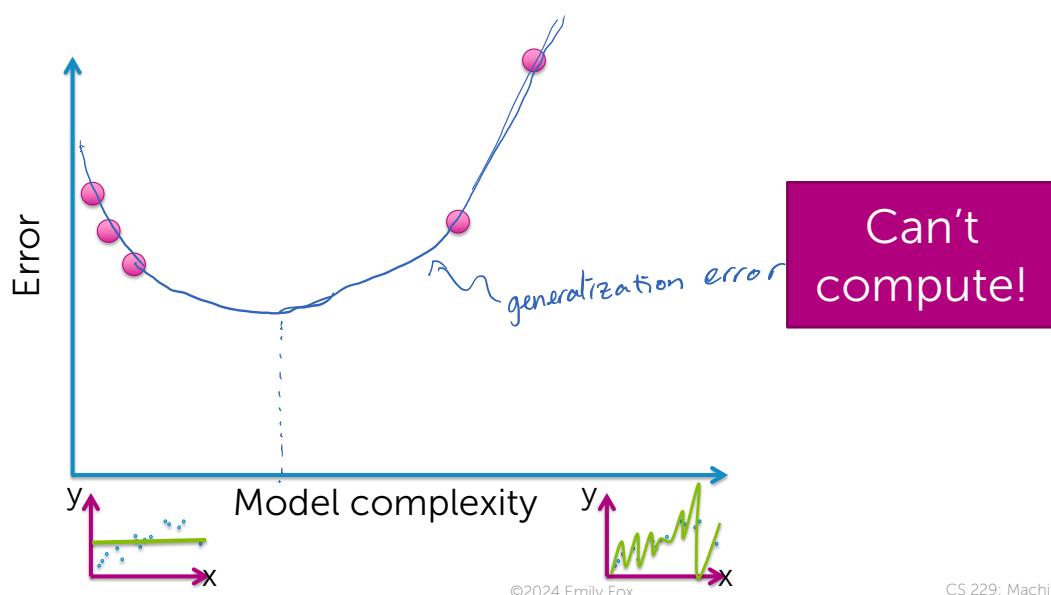
77

@2024 Emily Fox

CS 229: Machine Learning

77

## Generalization error vs. model complexity



78

@2024 Emily Fox

CS 229: Machine Learning

78

## Assessing the loss Part 3: Test error

©2024 Emily Fox

CS 229: Machine Learning

79

## Approximating generalization error

Wanted estimate of loss over all possible (, ) pairs



Approximate by looking at houses not in training set

80

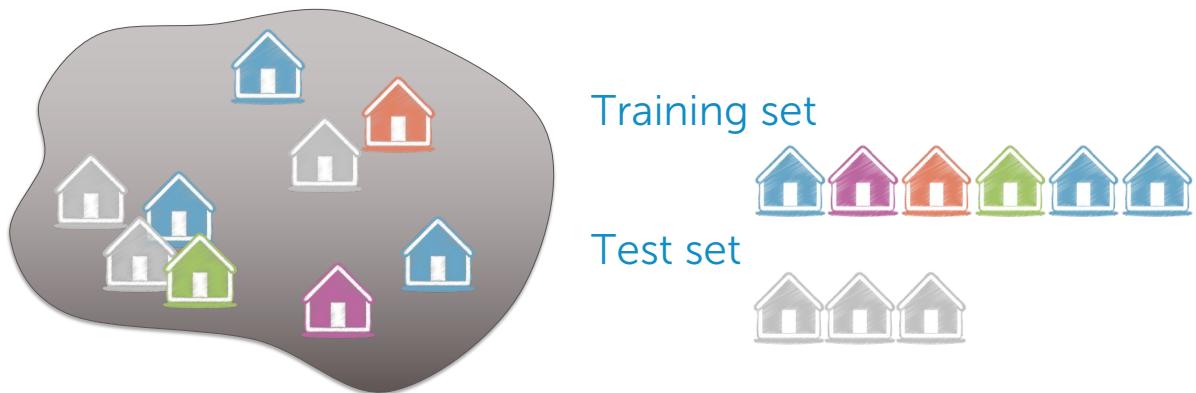
©2024 Emily Fox

CS 229: Machine Learning

80

## Forming a test set

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



81

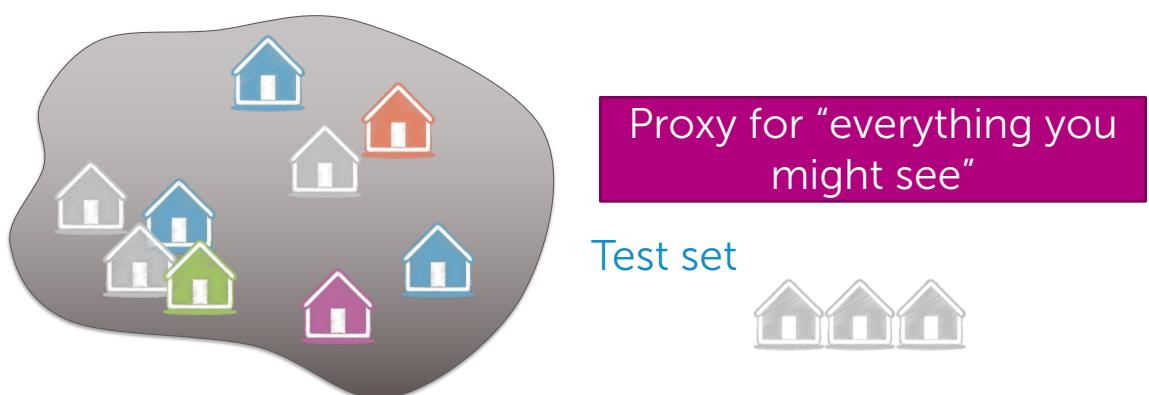
@2024 Emily Fox

CS 229: Machine Learning

81

## Forming a test set

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



82

@2024 Emily Fox

CS 229: Machine Learning

82

# Compute test error

$$\text{generalization error} = E_{x,y} [L(y, f_w(x))] \\ = \int L(y, f_w(x)) \underbrace{p(y, x)}_{p(y|x) p(x)} dx dy$$

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

## Test error

= avg. loss on houses in test set

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

# test points      fit using training data

**has never seen  
test data!**

83

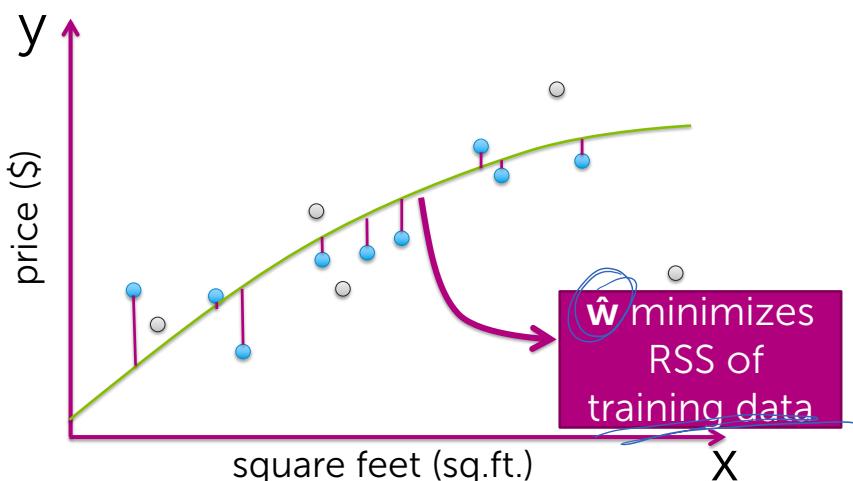
©2024 Emily Fox

CS 229: Machine Learning

83

## Example:

As before, fit quadratic to training data



84

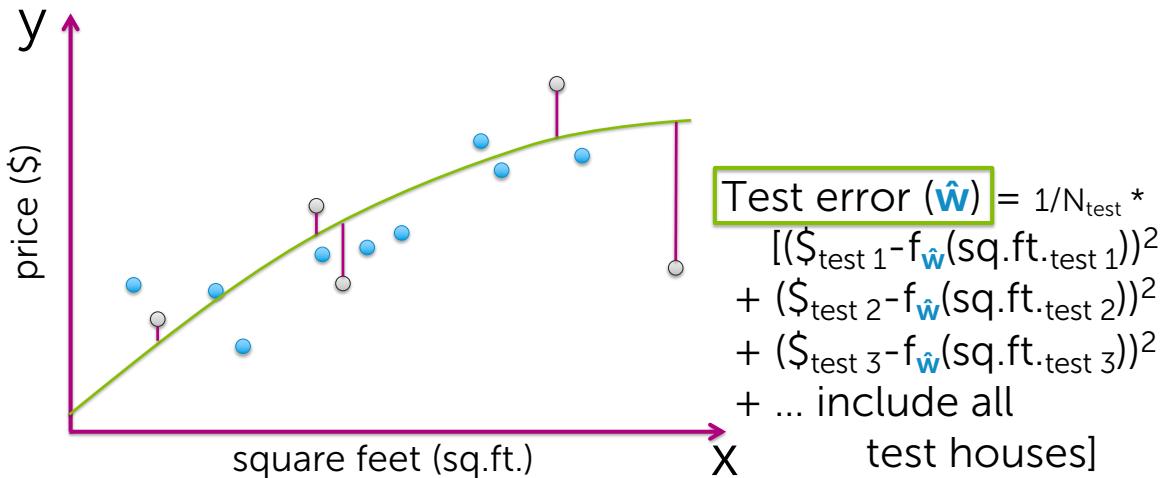
©2024 Emily Fox

CS 229: Machine Learning

84

**Example:**

As before, use **squared error loss**  $(y - f_{\hat{w}}(x))^2$



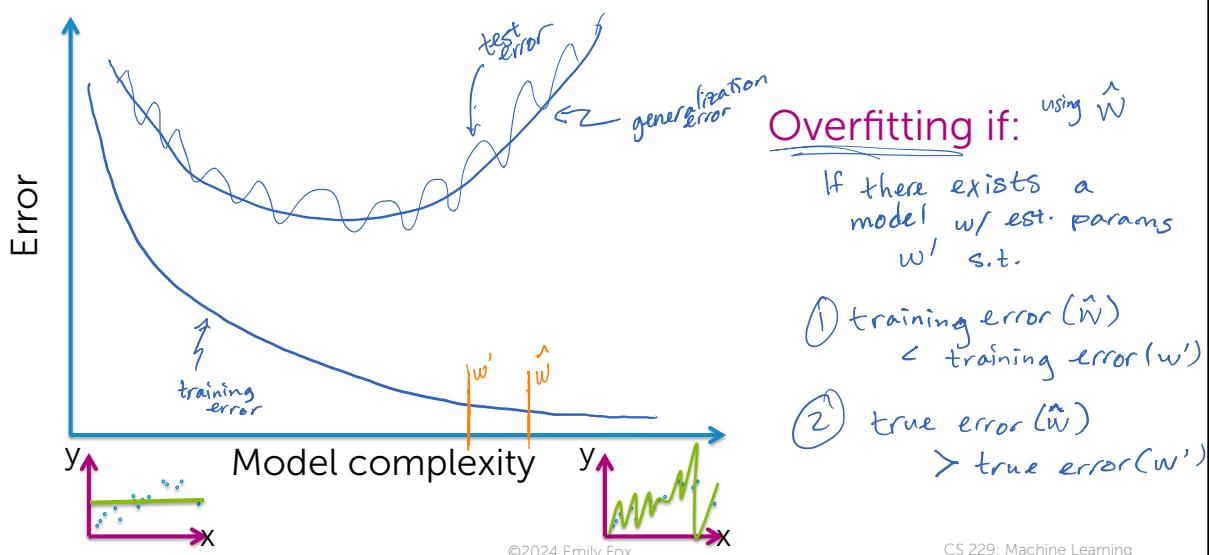
85

@2024 Emily Fox

CS 229: Machine Learning

85

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



86

@2024 Emily Fox

CS 229: Machine Learning

86

## Training/test split

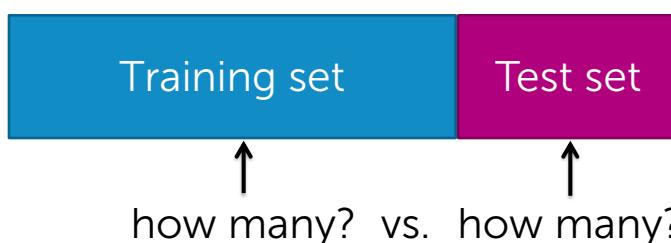
4

©2024 Emily Fox

CS 229: Machine Learning

87

## Training/test splits



©2024 Emily Fox

CS 229: Machine Learning

88

88

44

## Training/test splits



↑  
Too few →  $\hat{w}$  poorly estimated

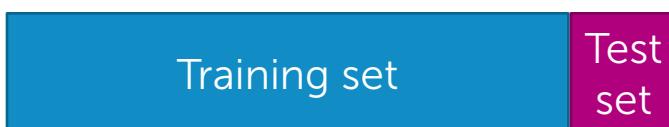
89

©2024 Emily Fox

CS 229: Machine Learning

89

## Training/test splits



↑  
Too few → test error bad approximation of generalization error

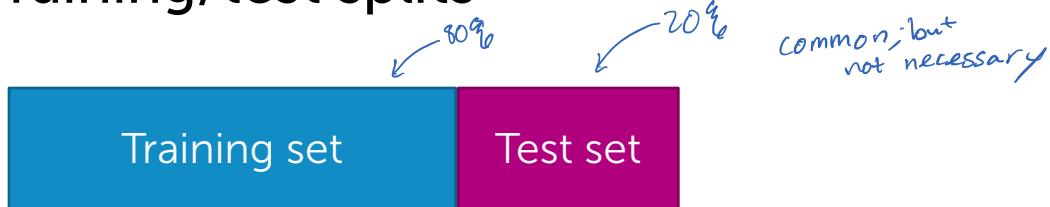
90

©2024 Emily Fox

CS 229: Machine Learning

90

## Training/test splits



Typically, just enough test points to form a reasonable estimate of generalization error

If this leaves too few for training, other methods like **cross validation**

91

@2024 Emily Fox

CS 229: Machine Learning

91

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

@2024 Emily Fox

CS 229: Machine Learning

92

## 3 sources of error

In forming predictions, there are 3 sources of error:

1. Noise
2. Bias
3. Variance

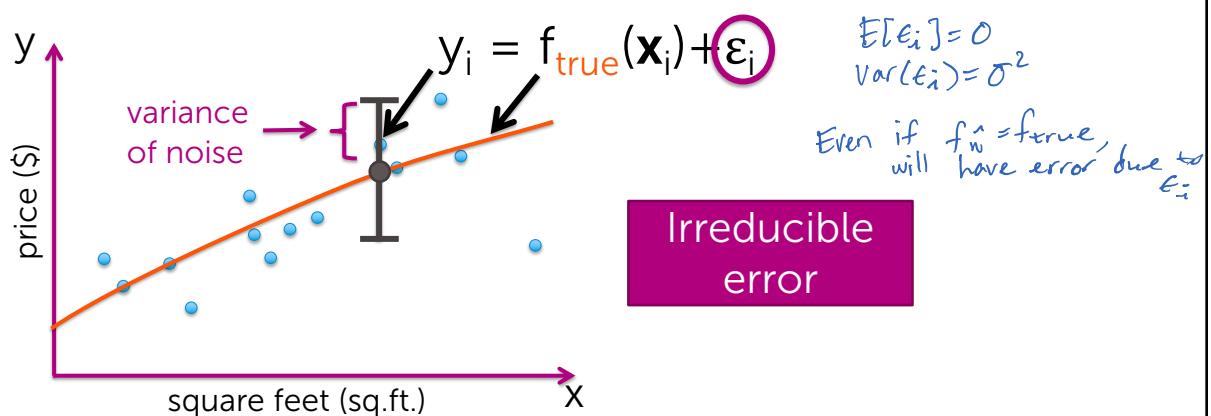
93

@2024 Emily Fox

CS 229: Machine Learning

93

## Data inherently noisy



94

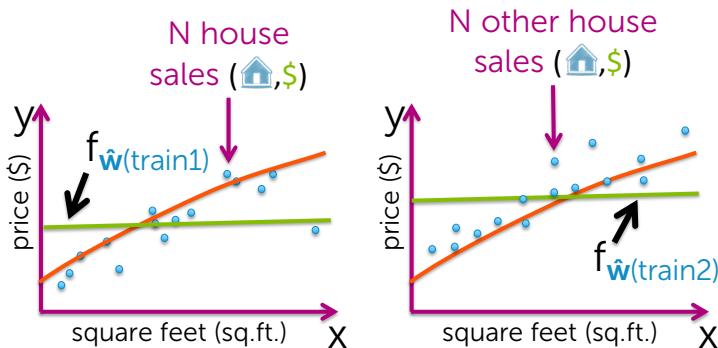
@2024 Emily Fox

CS 229: Machine Learning

94

## Bias contribution

Assume we fit a constant function



95

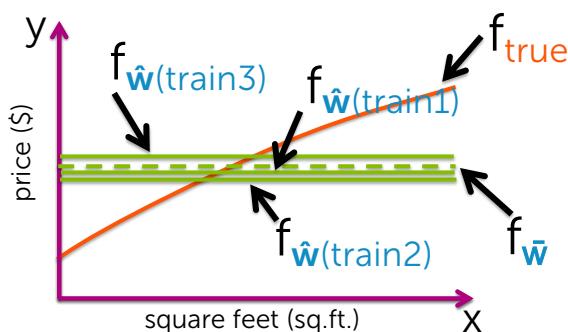
@2024 Emily Fox

CS 229: Machine Learning

95

## Bias contribution

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



96

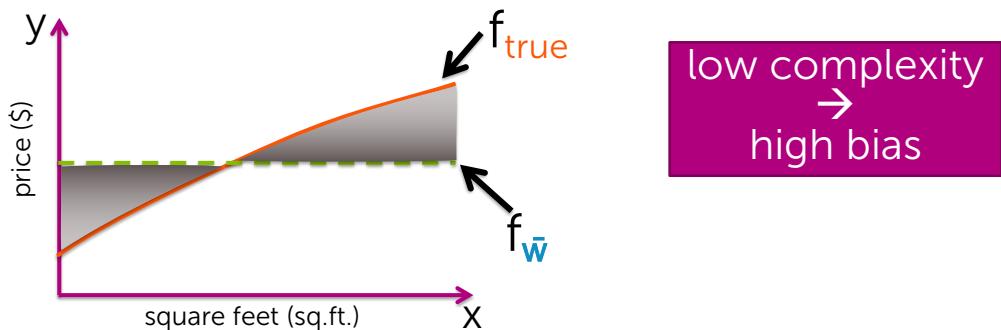
@2024 Emily Fox

CS 229: Machine Learning

96

## Bias contribution

$$\text{Bias}(\mathbf{x}) = f_{\text{true}}(\mathbf{x}) - f_{\bar{\mathbf{w}}}(\mathbf{x}) \quad \begin{matrix} \text{Is our approach flexible} \\ \leftarrow \text{enough to capture } f_{\text{true}}? \\ \text{If not, error in predictions.} \end{matrix}$$



97

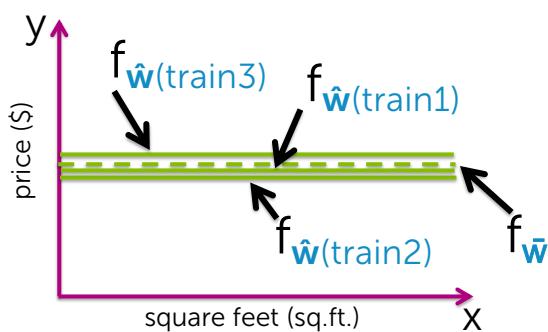
©2024 Emily Fox

CS 229: Machine Learning

97

## Variance contribution

How much do specific fits vary from the expected fit?



98

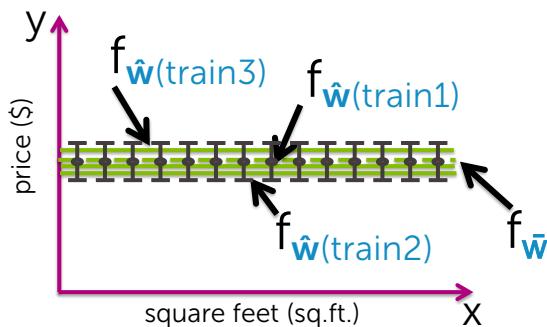
©2024 Emily Fox

CS 229: Machine Learning

98

## Variance contribution

How much do specific fits vary from the expected fit?



99

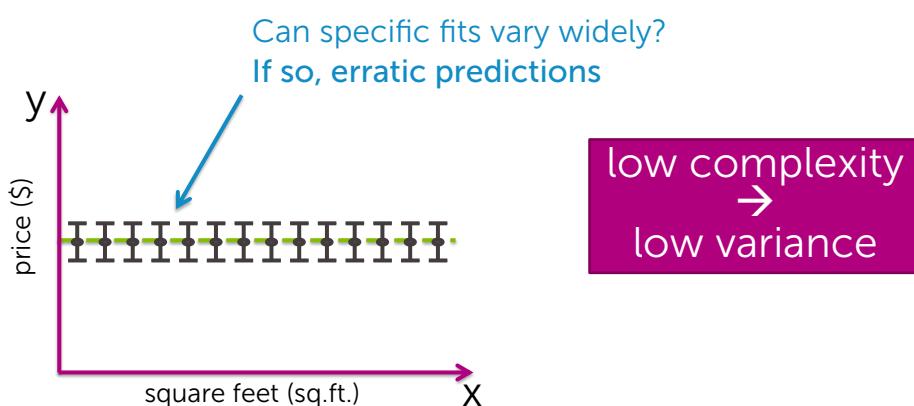
@2024 Emily Fox

CS 229: Machine Learning

99

## Variance contribution

How much do specific fits vary from the expected fit?



100

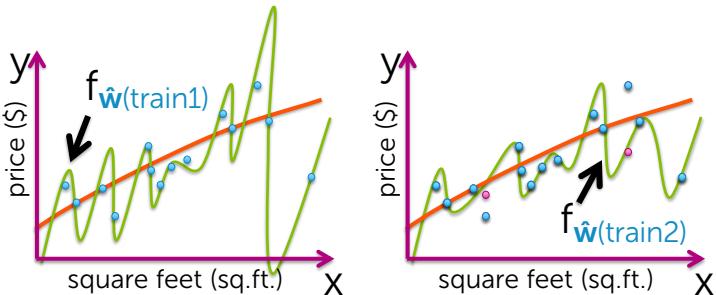
@2024 Emily Fox

CS 229: Machine Learning

100

## Variance of high-complexity models

Assume we fit a high-order polynomial



101

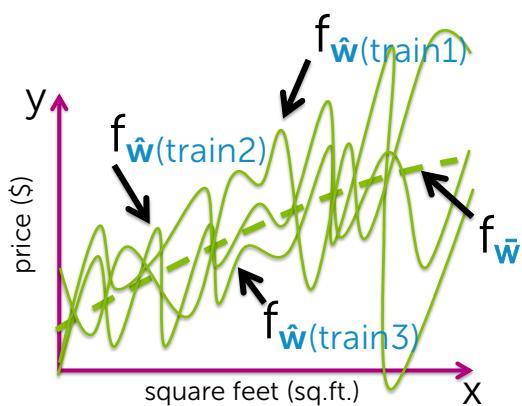
@2024 Emily Fox

CS 229: Machine Learning

101

## Variance of high-complexity models

Assume we fit a high-order polynomial



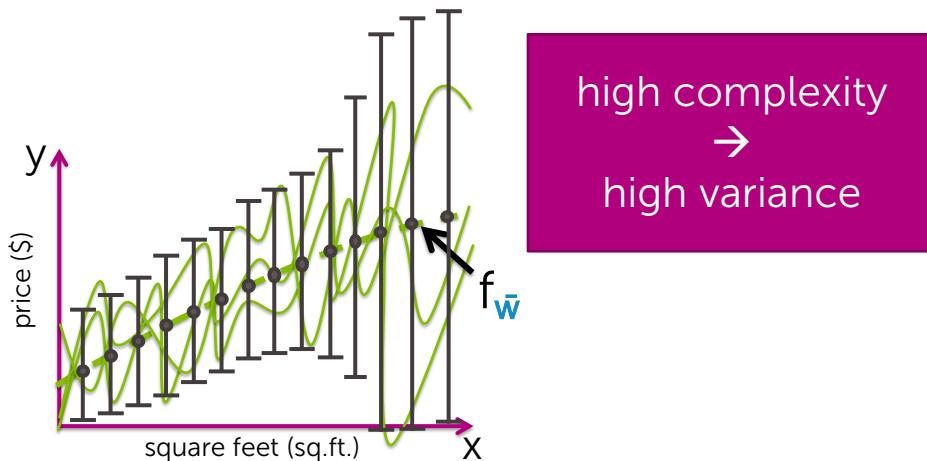
102

@2024 Emily Fox

CS 229: Machine Learning

102

## Variance of high-complexity models



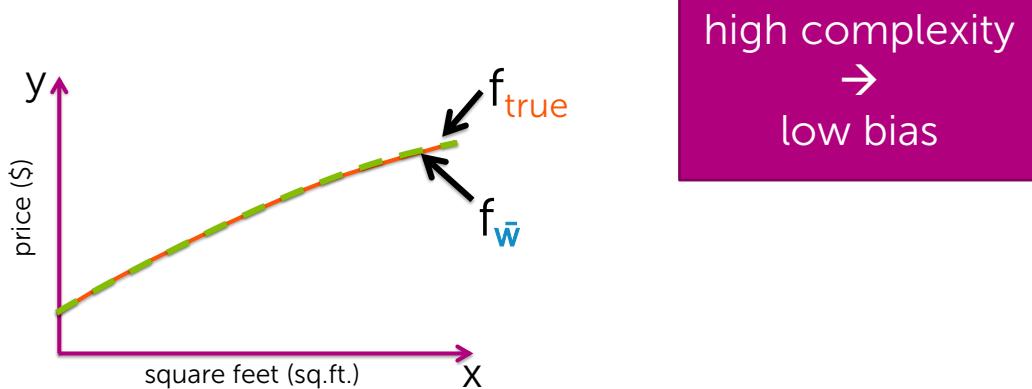
103

©2024 Emily Fox

CS 229: Machine Learning

103

## Bias of high-complexity models



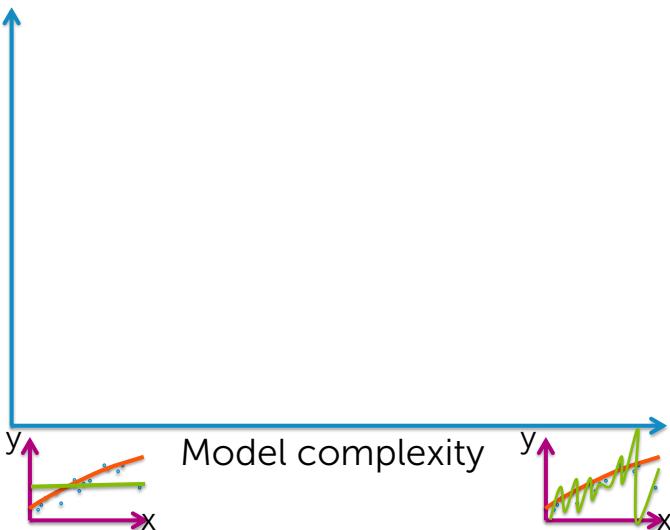
104

©2024 Emily Fox

CS 229: Machine Learning

104

## Bias-variance tradeoff



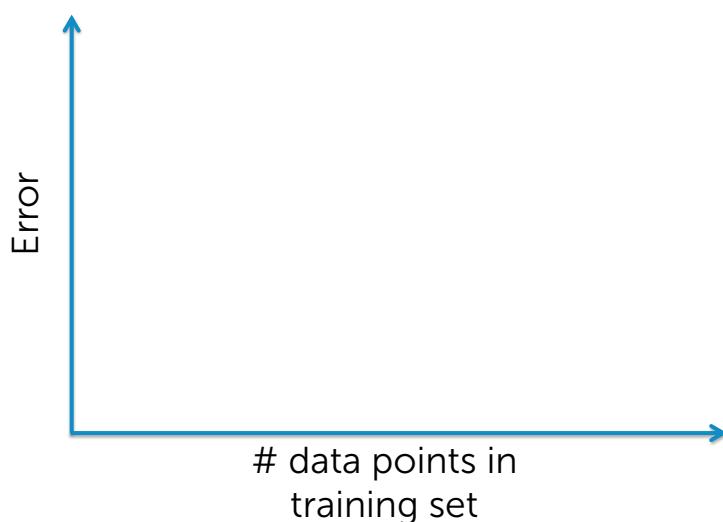
105

@2024 Emily Fox

CS 229: Machine Learning

105

## Error vs. amount of data



106

@2024 Emily Fox

CS 229: Machine Learning

106

## Why 3 sources of error? A formal derivation

©2024 Emily Fox

CS 229: Machine Learning

107

## Bias of function estimator, formally

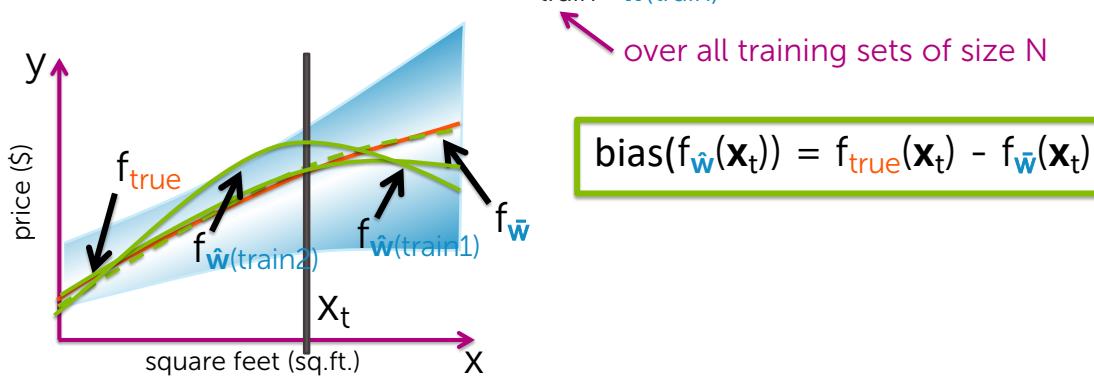
Average estimated function =  $f_{\bar{w}}(x)$

True function =  $f_{true}(x)$

$$E_{train}[f_{\hat{w}(\text{train})}(x)]$$

over all training sets of size N

$$\text{bias}(f_{\hat{w}}(x_t)) = f_{true}(x_t) - f_{\bar{w}}(x_t)$$



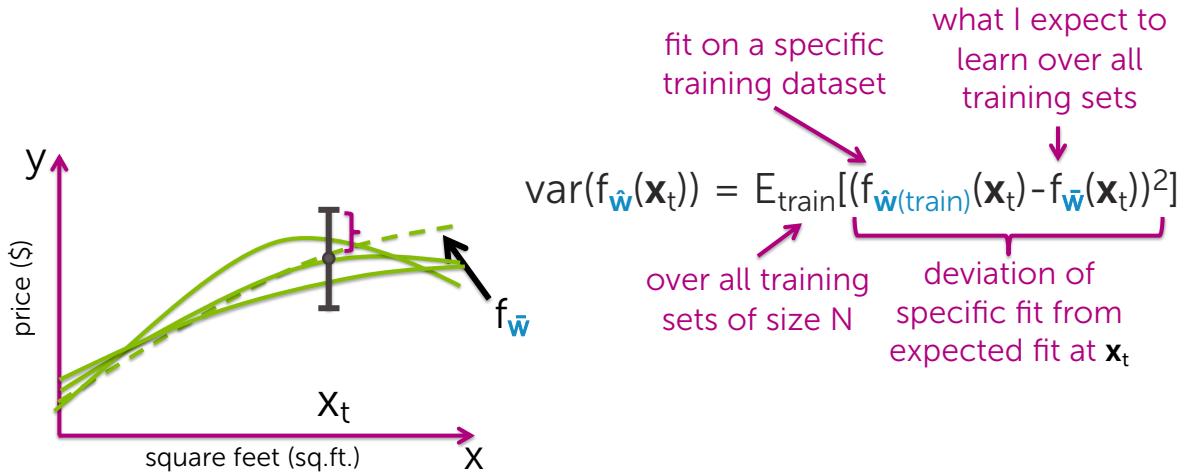
108

©2024 Emily Fox

CS 229: Machine Learning

108

## Variance of function estimator, formally



109

@2024 Emily Fox

CS 229: Machine Learning

109

## Deriving expected prediction error

Expected prediction error

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

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

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

Expected prediction error at  $x_t$

$= E_{\text{train}} [E_{y_t|x_t} [(y_t - f_{\hat{w}(\text{train})}(x_t))^2]]$

110

@2024 Emily Fox

CS 229: Machine Learning

110

## Deriving expected prediction error

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

$$\begin{aligned} &= E_{\text{train}}[E_{y_t|\mathbf{x}_t}[(y_t - f_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x}_t))^2]] \\ &= E_{\text{train}}[E_{y_t|\mathbf{x}_t}[(y_t - f_{\text{true}}(\mathbf{x}_t)) + (f_{\text{true}}(\mathbf{x}_t) - f_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x}_t)))^2]] \end{aligned}$$

111

@2024 Emily Fox

CS 229: Machine Learning

111

## Equating MSE with bias and variance

$$\begin{aligned} \text{MSE}[f_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x}_t)] &= E_{\text{train}}[(f_{\text{true}}(\mathbf{x}_t) - f_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x}_t))^2] \\ &= E_{\text{train}}[((f_{\text{true}}(\mathbf{x}_t) - f_{\bar{\mathbf{w}}}(\mathbf{x}_t)) + (f_{\bar{\mathbf{w}}}(\mathbf{x}_t) - f_{\hat{\mathbf{w}}(\text{train})}(\mathbf{x}_t)))^2] \end{aligned}$$

112

@2024 Emily Fox

CS 229: Machine Learning

112

## Putting it all together

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

$$\begin{aligned} &= \sigma^2 + \text{MSE}[f_{\hat{\mathbf{w}}}(\mathbf{x}_t)] \\ &= \sigma^2 + [\text{bias}(f_{\hat{\mathbf{w}}}(\mathbf{x}_t))]^2 + \text{var}(f_{\hat{\mathbf{w}}}(\mathbf{x}_t)) \end{aligned}$$

3 sources of error

113

©2024 Emily Fox

CS 229: Machine Learning

113

## Summary of assessing performance

114

©2024 Emily Fox

CS 229: Machine Learning

## 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

115

©2024 Emily Fox

CS 229: Machine Learning

115