

# INTRODUCTION TO DATA SCIENCE

This lecture is  
based on course by E. Fox and C. Guestrin, Univ of Washington

26/10, 2/11,  
9/11 2022

WFAiS UJ, Informatyka Stosowana  
I stopień studiów

# Regression for predictions

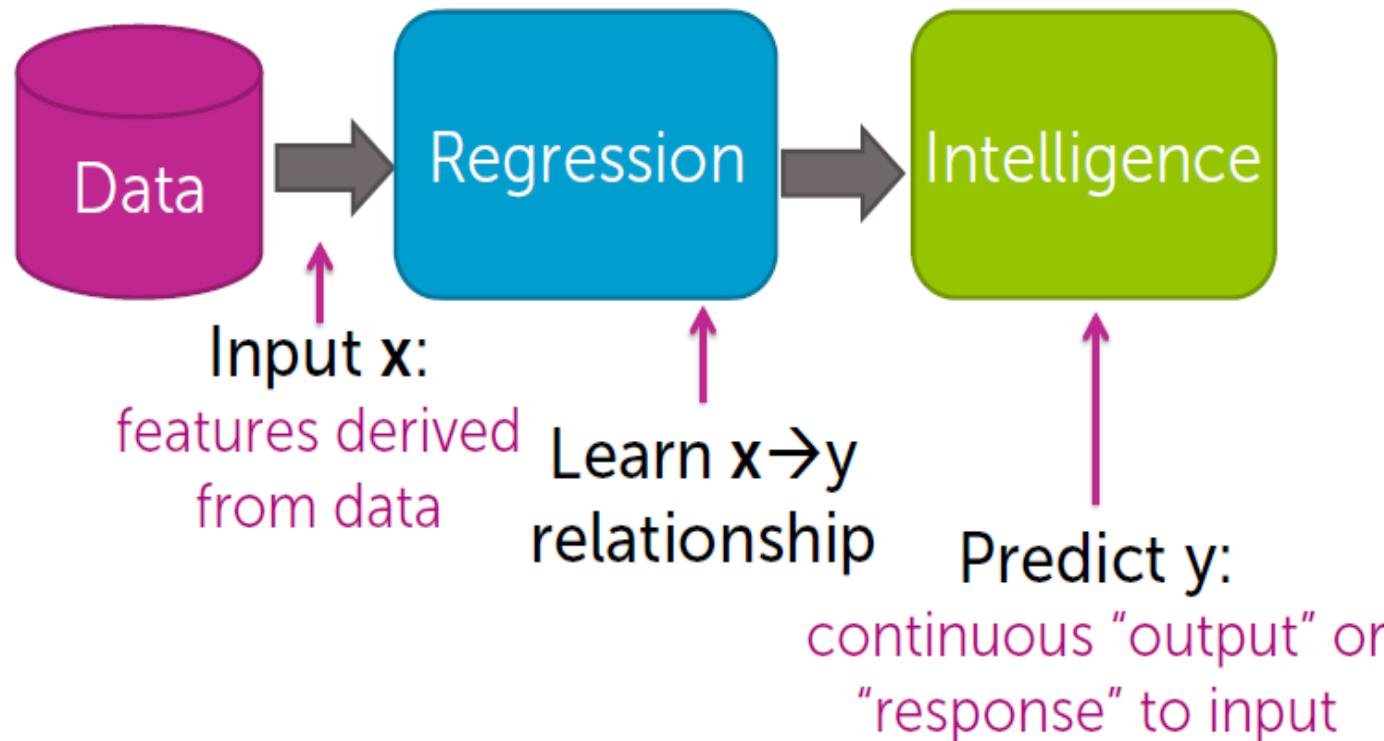
2

- **Simple regression**
- **Multiple regression**
- **Accessing performance**
- **Ridge regression**
- **Feature selection and lasso regression**
- **Nearest neighbor and kernel regression**

# What is regression?

3

From features to predictions



# Case study

4

## Predicting house prices



# Data

5



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



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



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



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



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

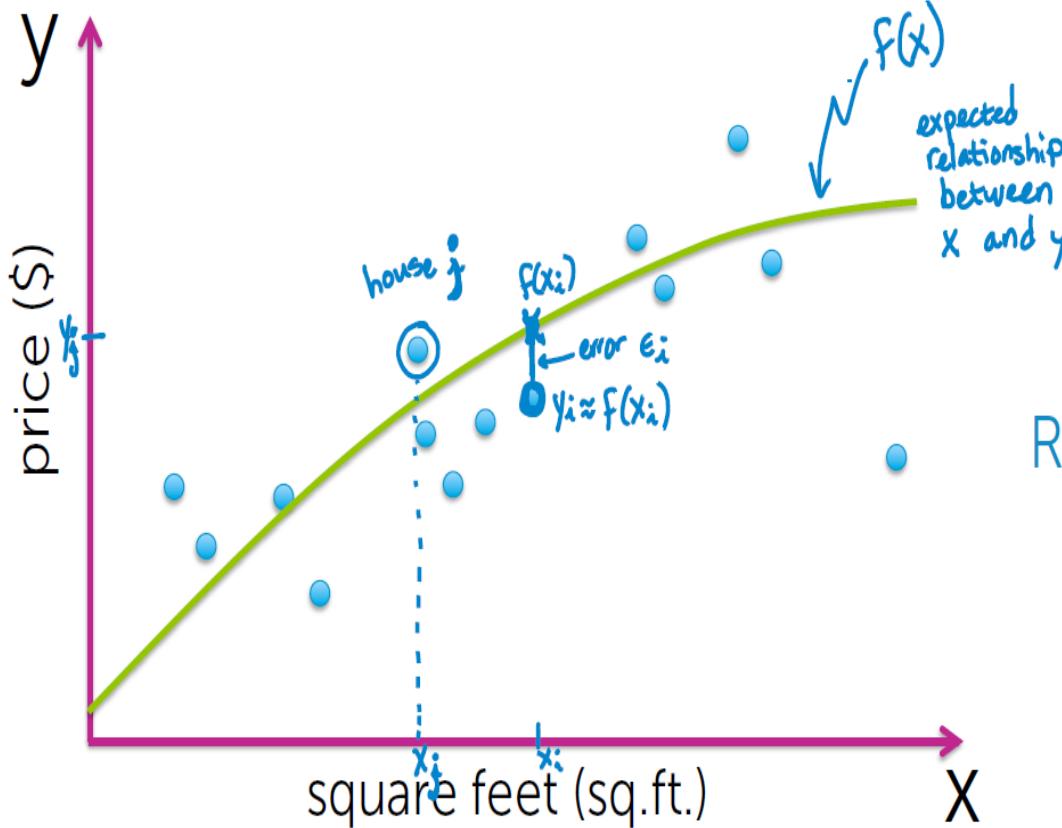
:

## Input vs output

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

# Model: assume functional relationship

6



„Essentially, all models are wrong but some are useful.”  
George Box, 1987.

Regression model:

$$y_i = f(x_i) + \epsilon_i$$

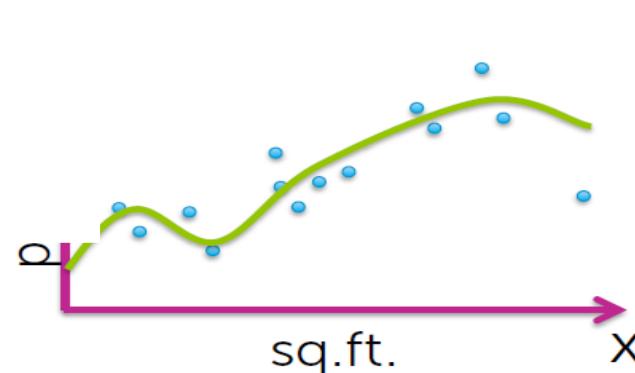
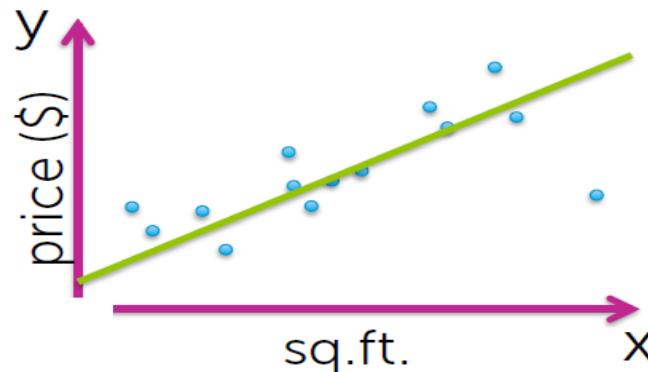
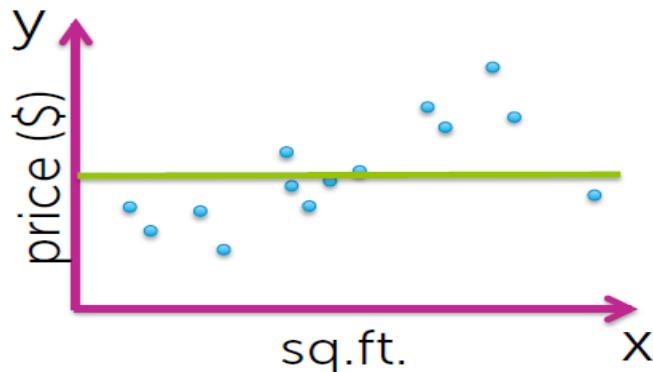
$E[\epsilon_i] = 0$  ← equally likely  
that error  
is + or -  
↑ expected value

↓  
 $y_i$  is equally  
likely to be above  
or below  $f(x_i)$

# Task 1:

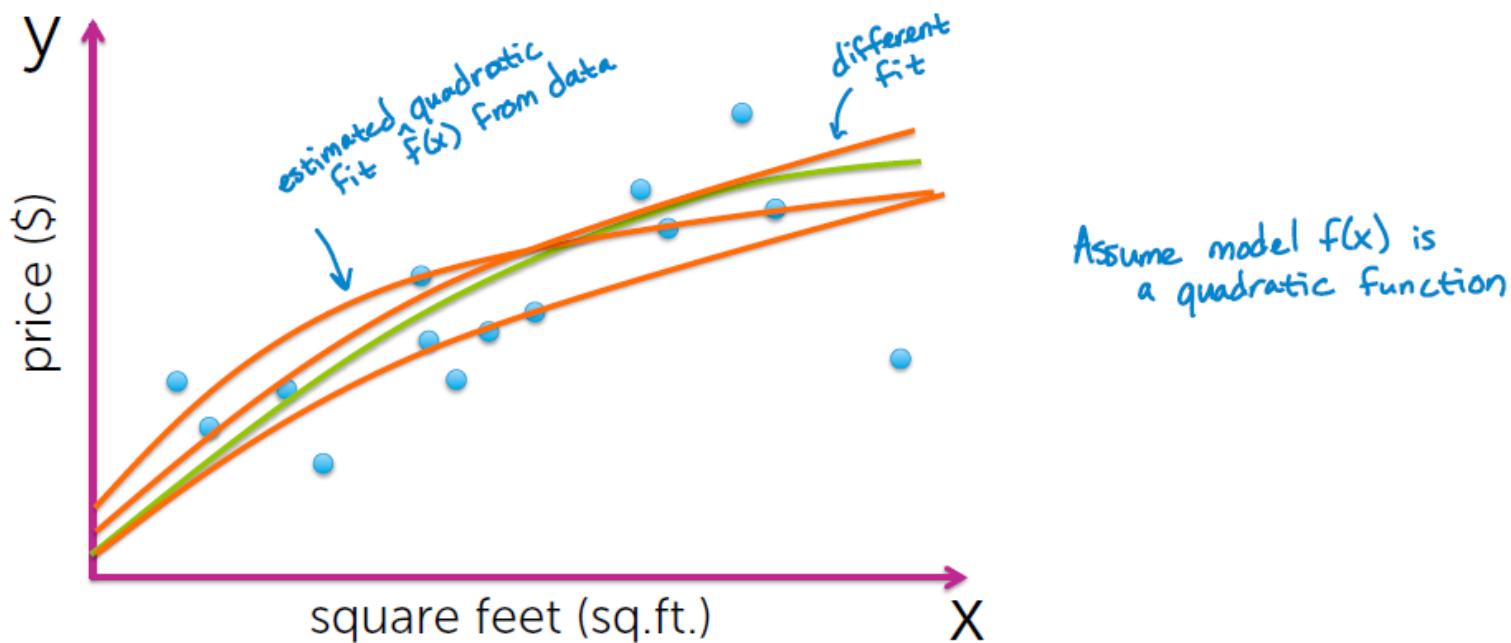
7

## Which model to fit?



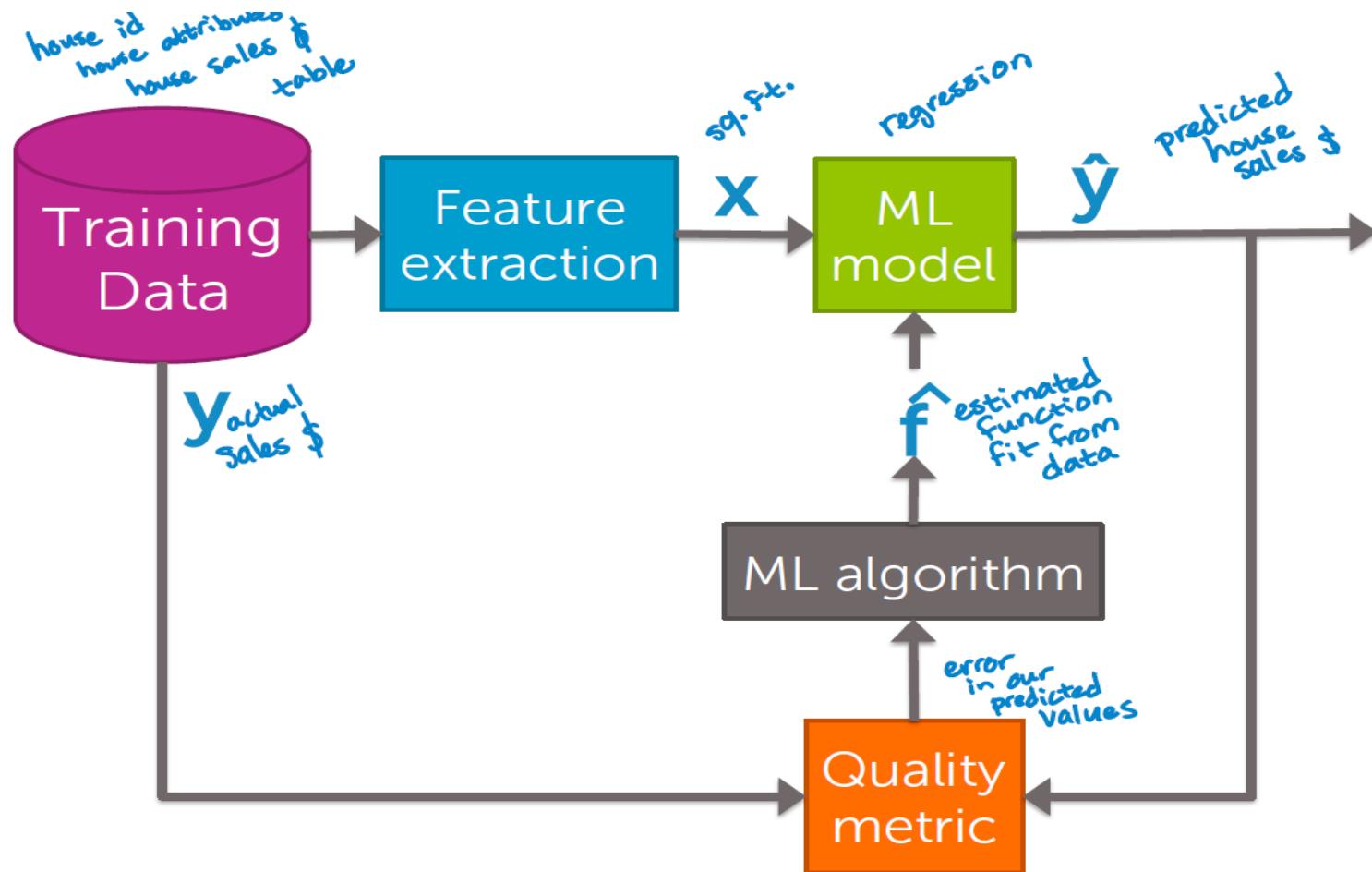
## Task 2:

**For a given model  $f(x)$  estimate function  $\hat{f}(x)$  from data**



# How it works: baseline flow chart

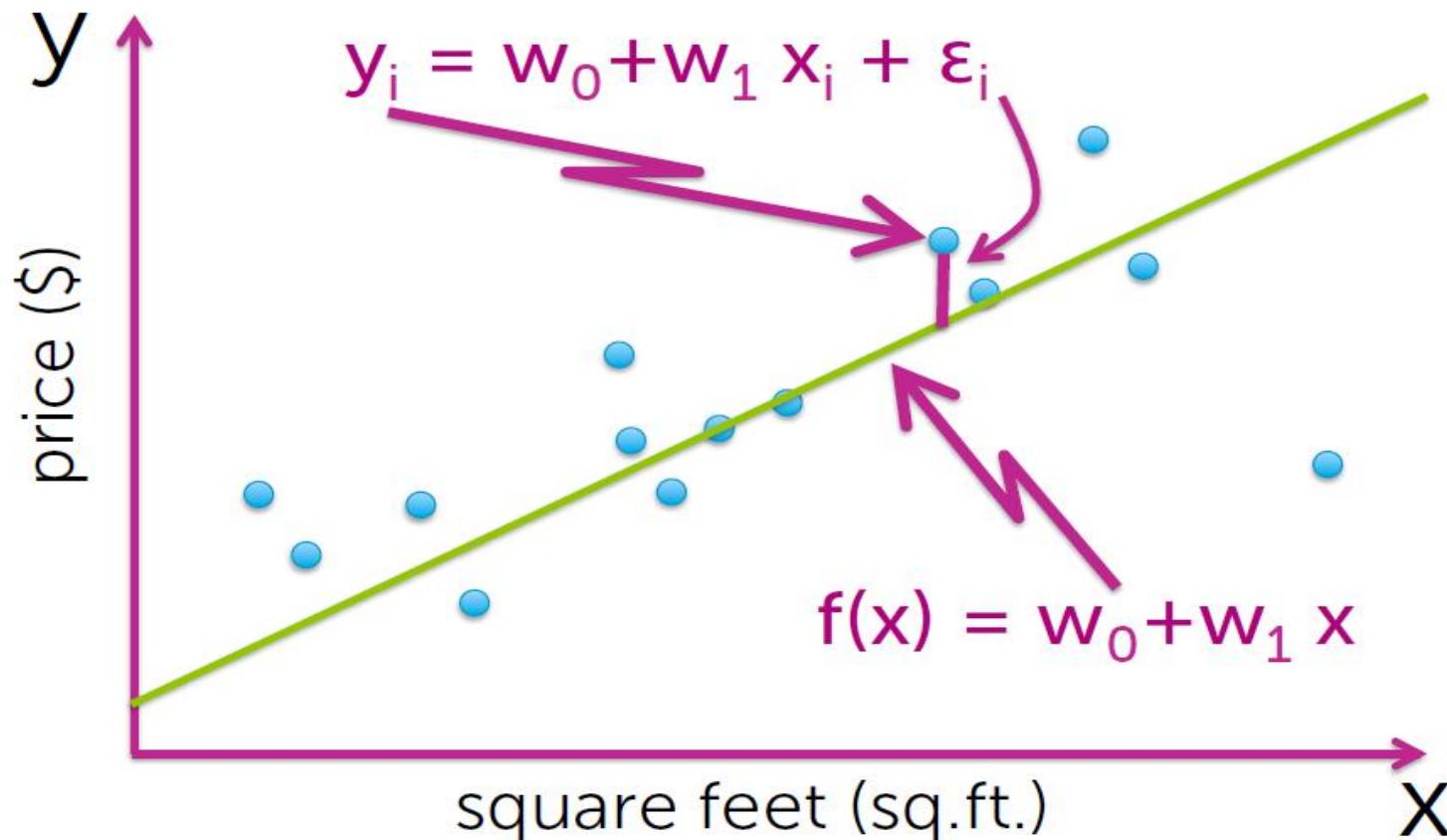
9



# SIMPLE LINEAR REGRESSION

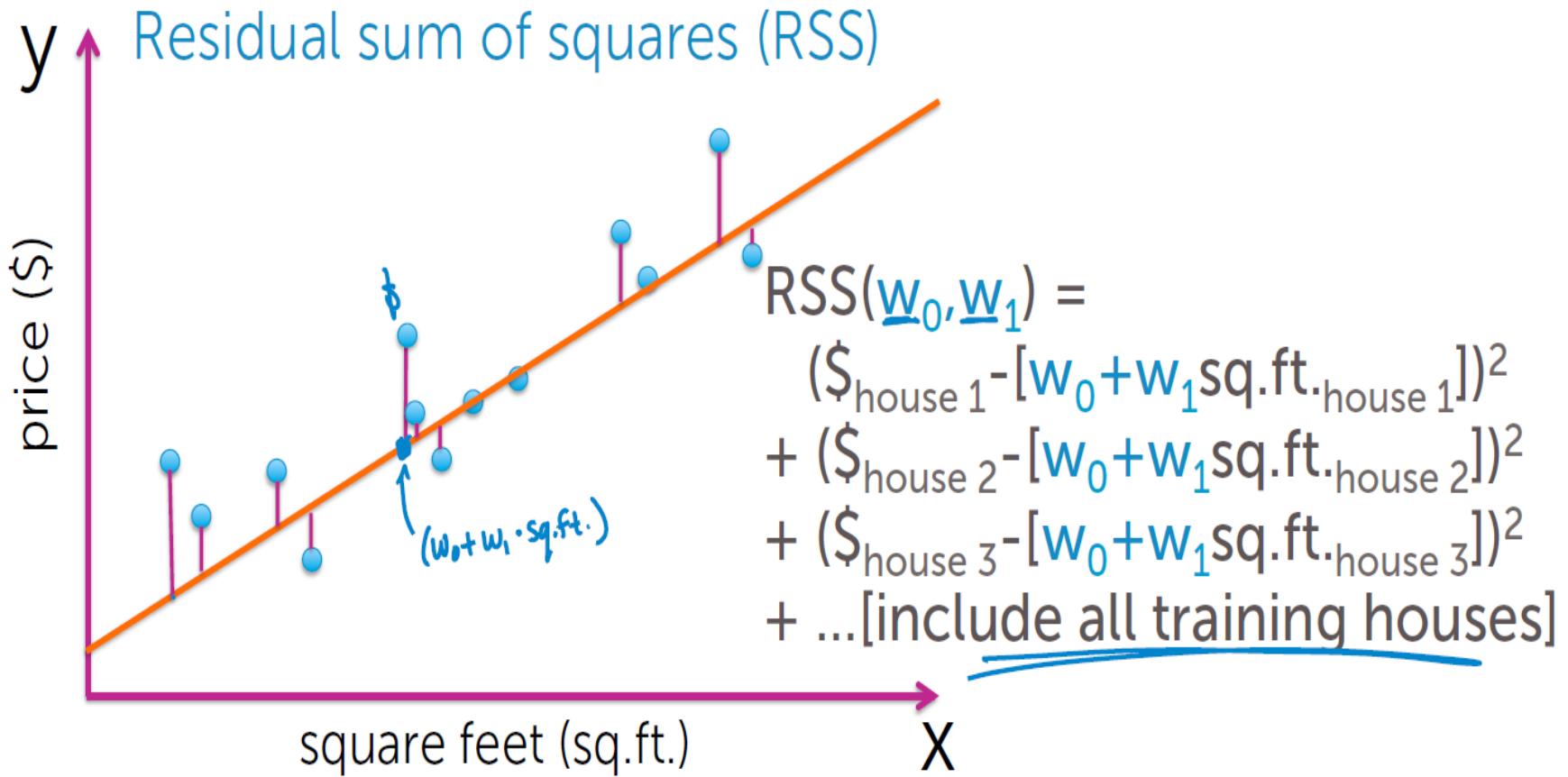
# Simple linear regression model

11



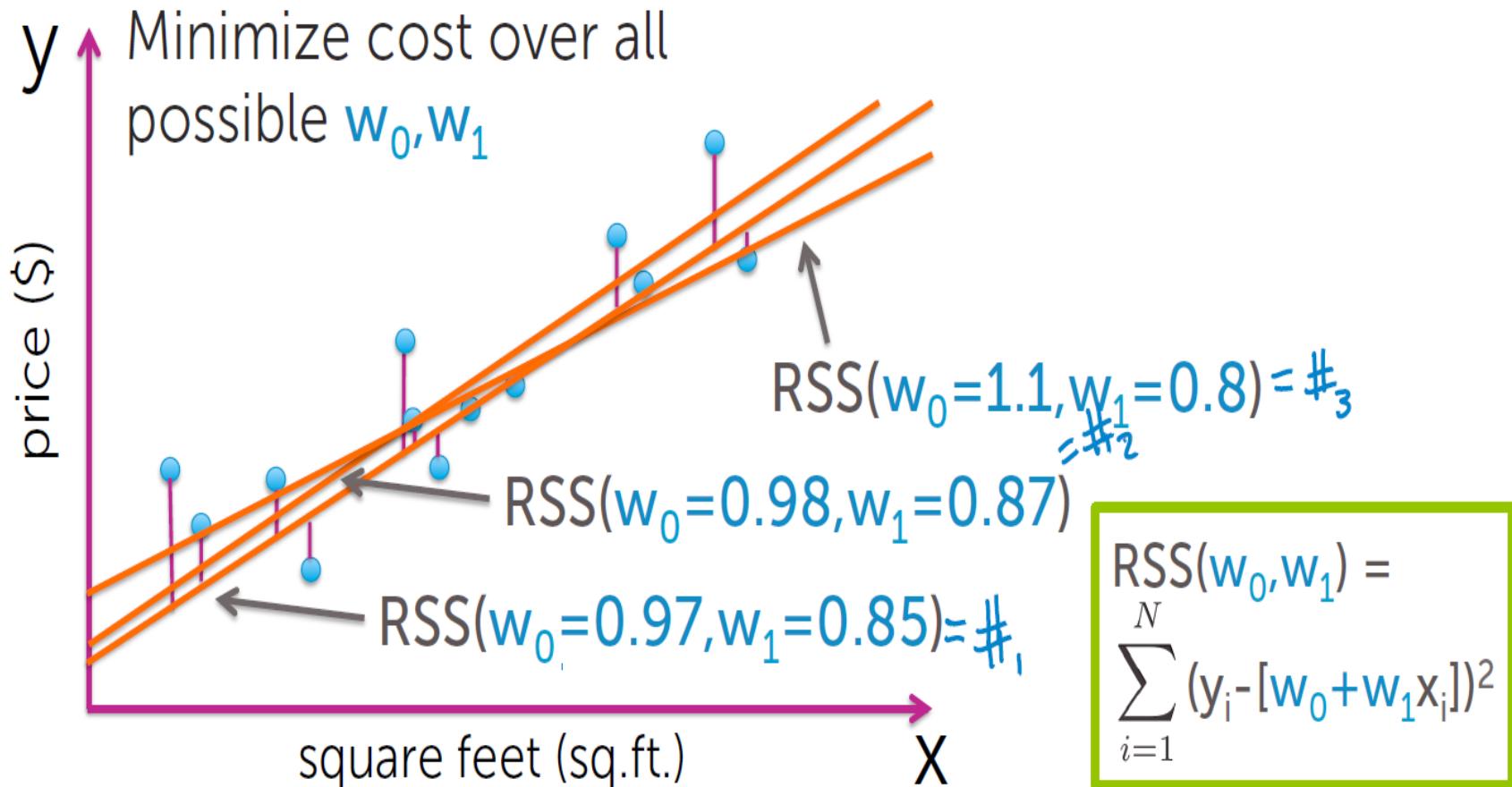
# The cost of using a given line

12



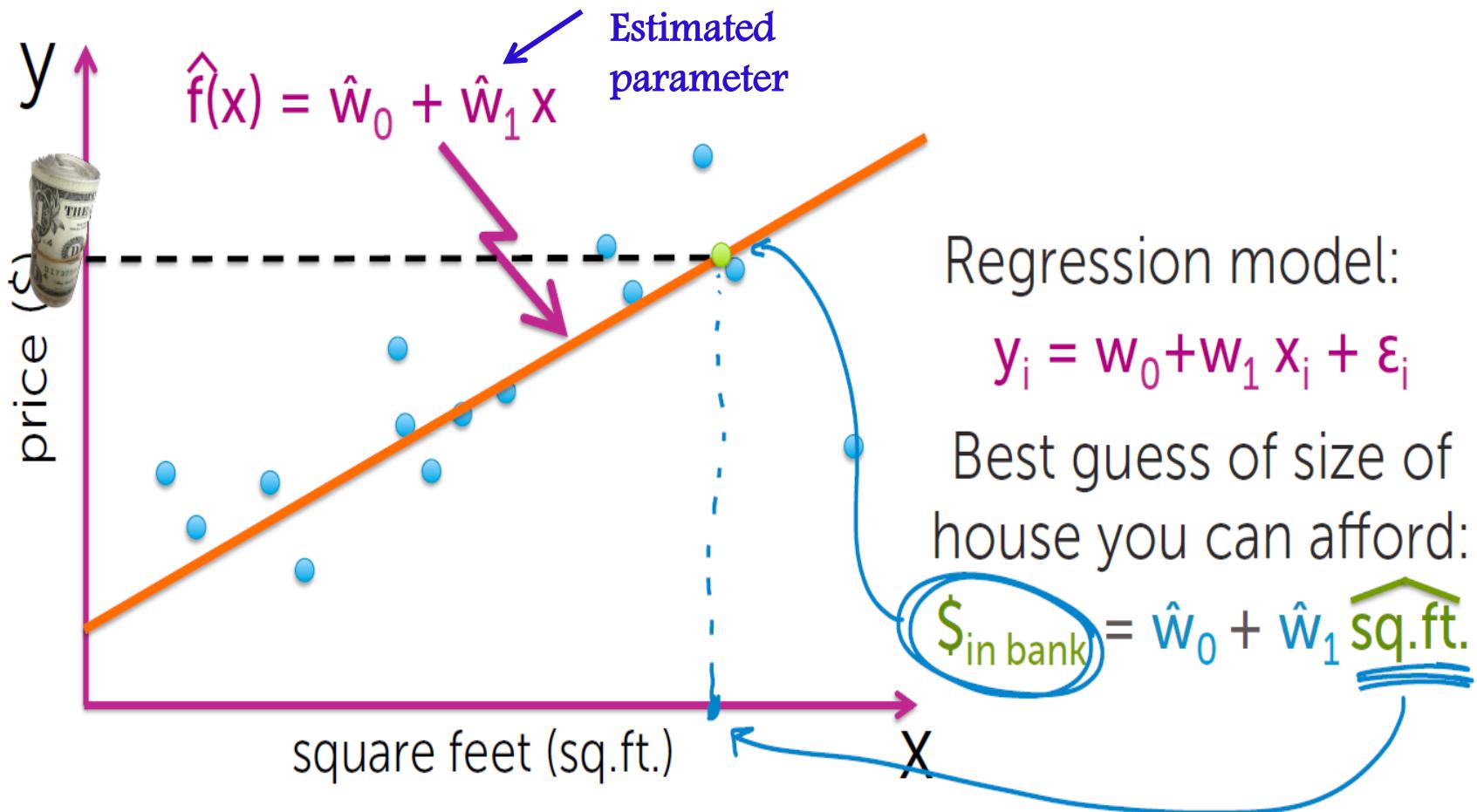
# Find „best” line

13



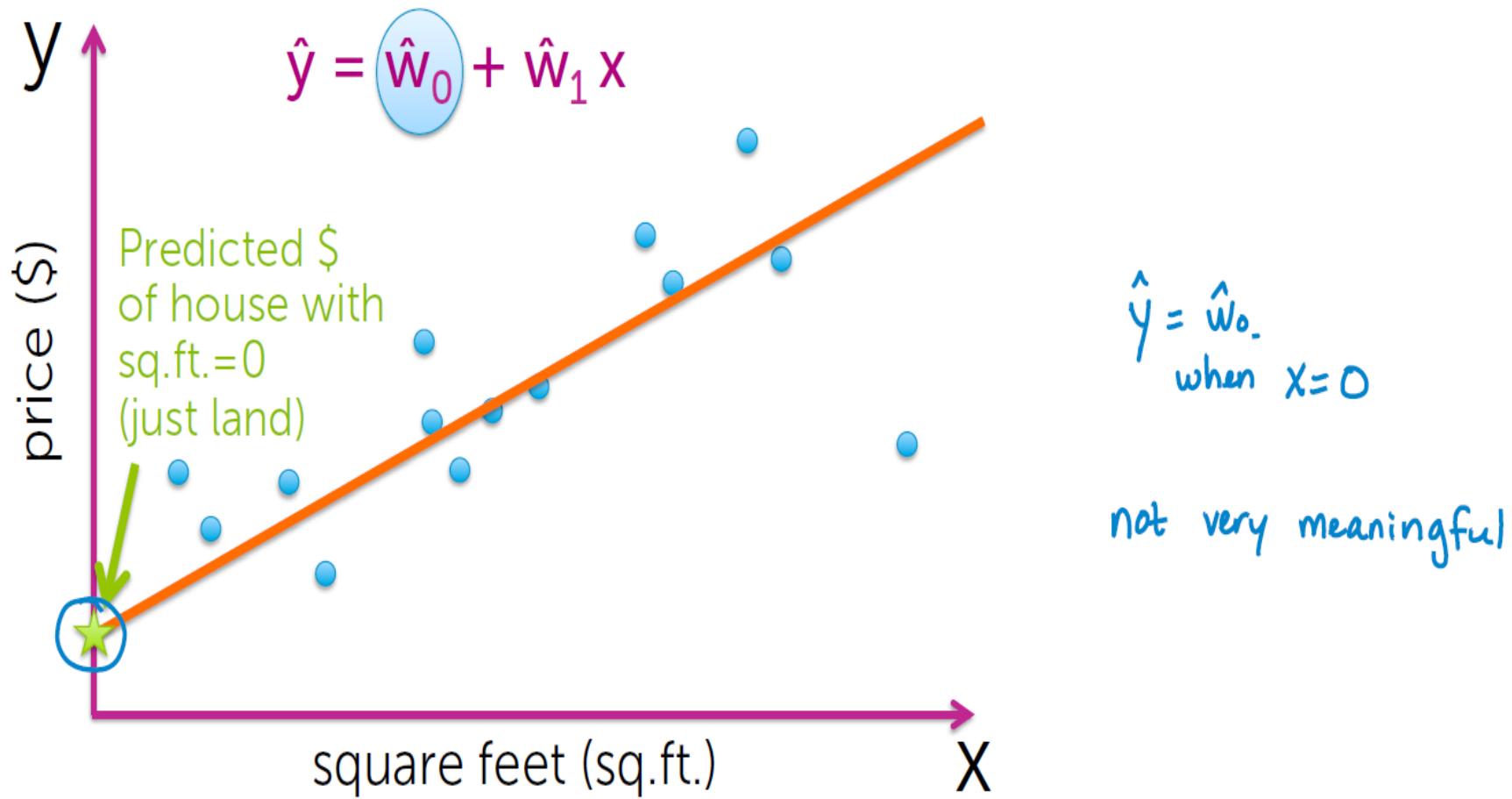
# Predicting size of house you can afford

14



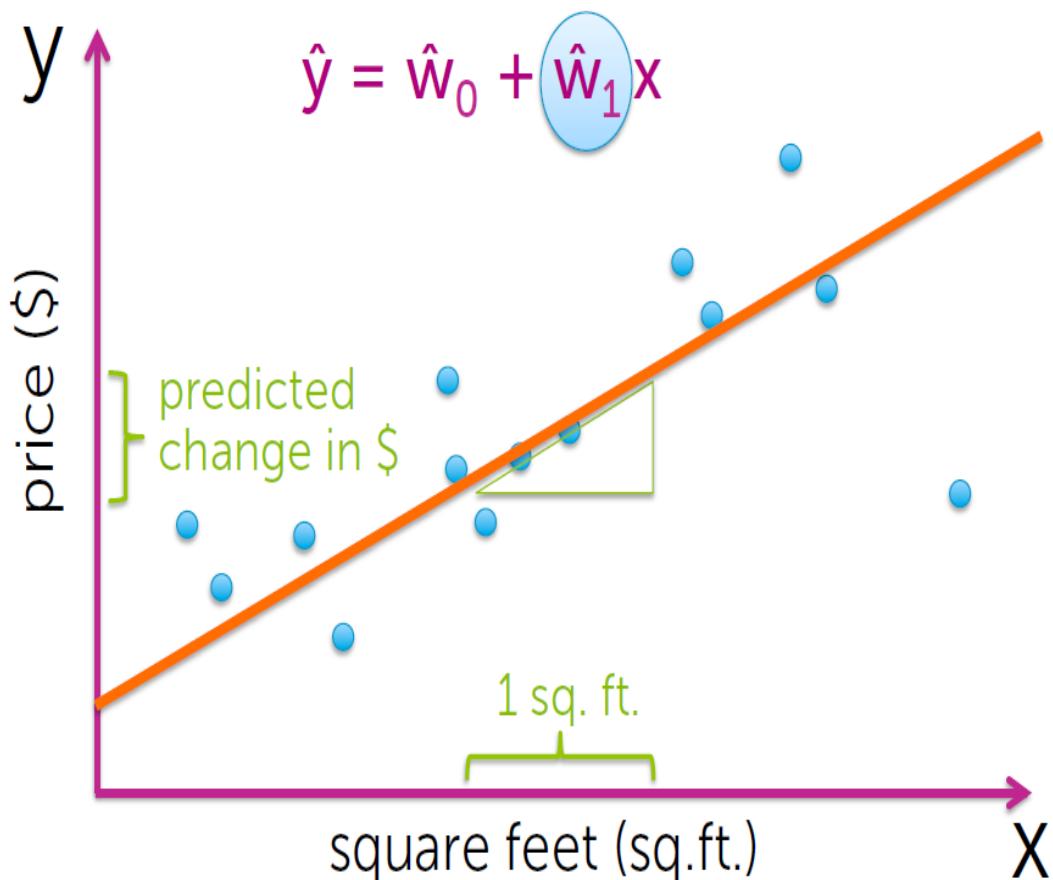
# Interpreting the coefficients

15



# Interpreting the coefficients

16



Magnitude of fit parameters depend on the units of both features and observations

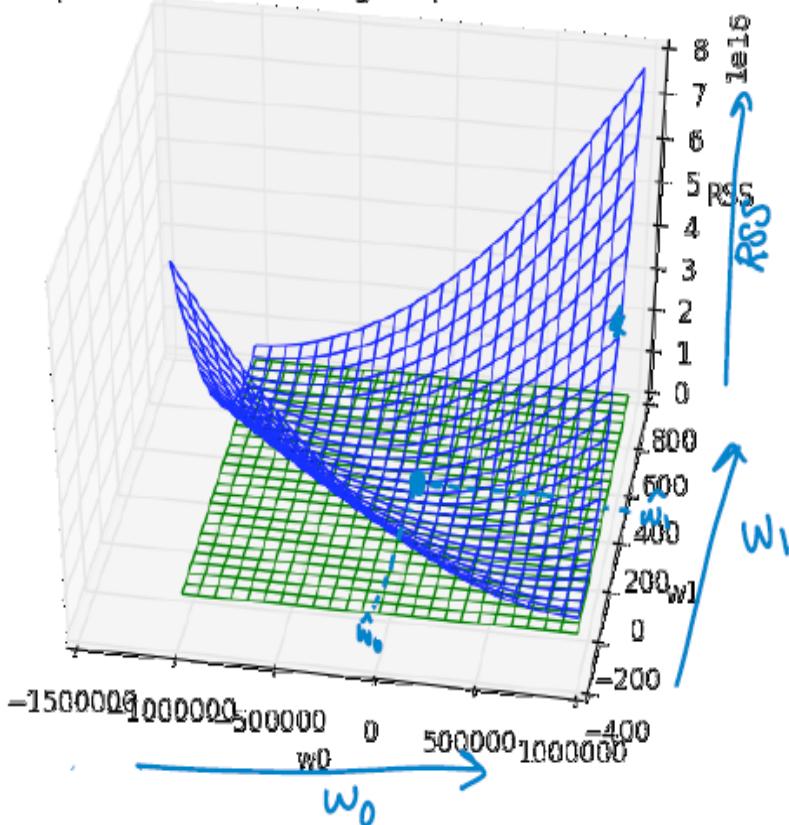
$$\begin{aligned}\hat{w}_1 &= \frac{\hat{y}_{1001 \text{ sq.ft.}} - \hat{y}_{1000 \text{ sq.ft.}}}{1001 \text{ sq.ft.} - 1000 \text{ sq.ft.}} \\ &= \hat{w}_0 + \hat{w}_1 \cdot 1001 \text{ sq.ft.} \\ &\quad - (\hat{w}_0 + \hat{w}_1 \cdot 1000 \text{ sq.ft.}) \\ &= \hat{w}_1\end{aligned}$$

predicted change in the output per unit change in input

# ML algorithm: minimising the cost

17

3D plot of RSS with tangent plane at minimum



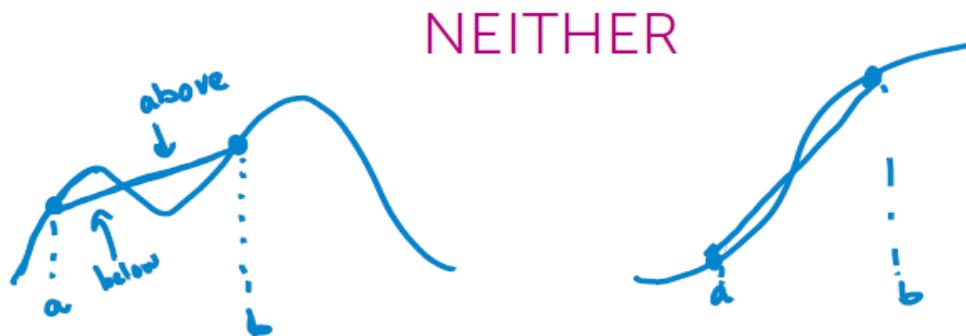
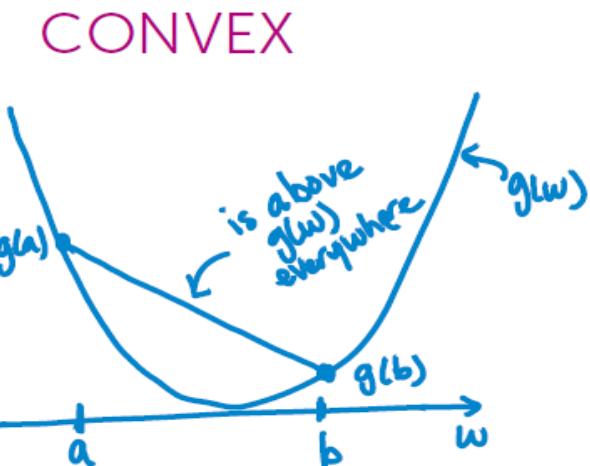
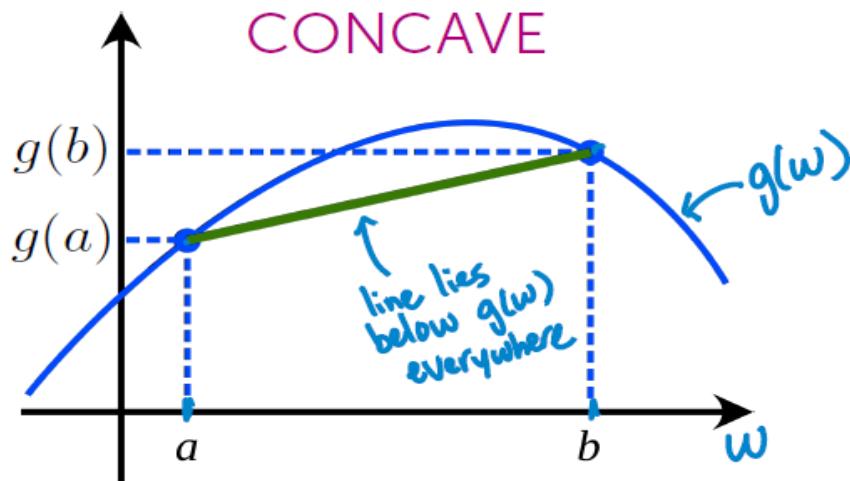
Minimize function  
over all possible  $w_0, w_1$

$$\min_{W_0, W_1} \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

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

# Convex/concave function

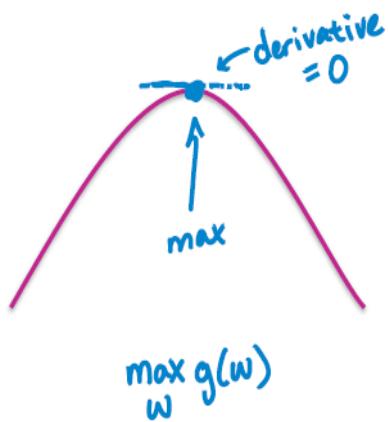
18



# Finding max/min analytically

19

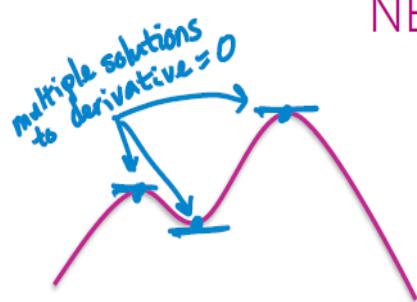
CONCAVE



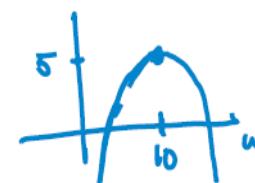
CONVEX



NEITHER



no solution to derivative = 0



Example:

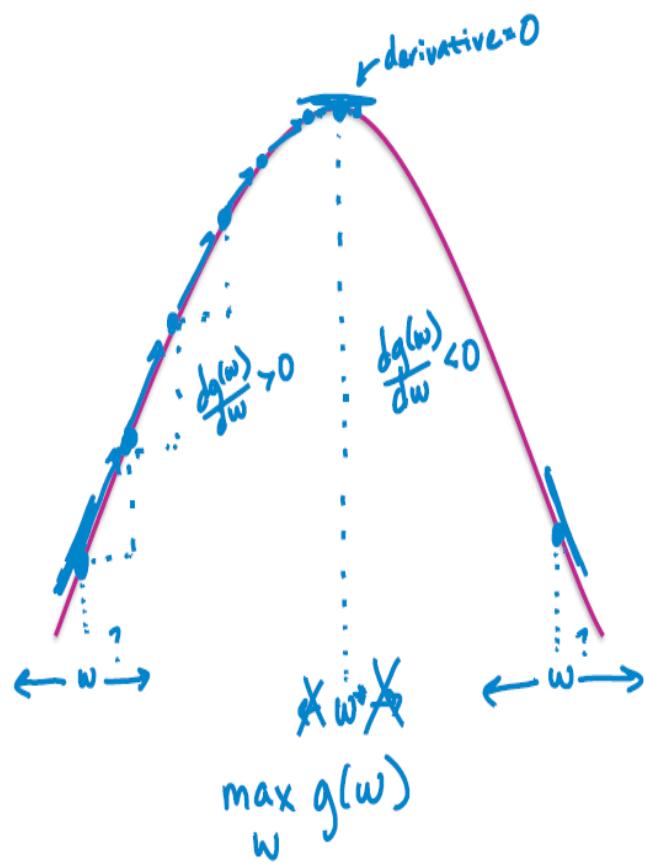
$$g(w) = 5 - (w - 10)^2$$

$$\begin{aligned} \frac{dg(w)}{dw} &= 0 - 2(w - 10)^1 \cdot 1 \\ &= -2w + 20 \end{aligned}$$

$$\begin{aligned} \text{Set derivative } &= 0: \\ -2w + 20 &= 0 \\ w &= 10 \end{aligned}$$

# Finding the max via hill climbing

20



Sign of the derivative is saying me what I want to do :move left or right or stay where I am

How do we know whether to move w to right or left?  
(inc. or dec. the value of w?)

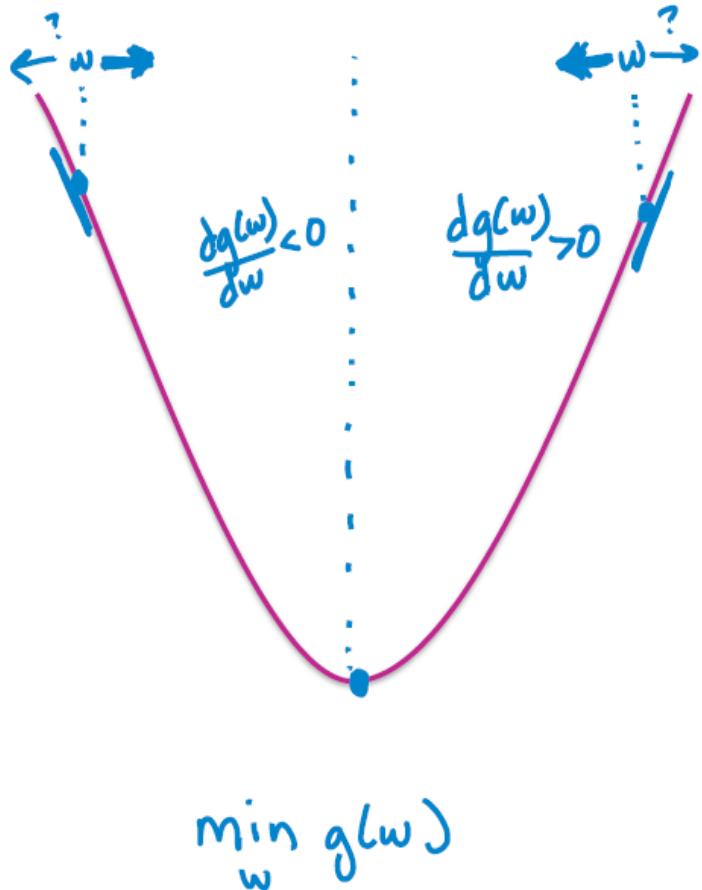
while not converged

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

iteration t  
stepsize

# Finding the min via hill descent

21



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

Algorithm:

**while** not converged

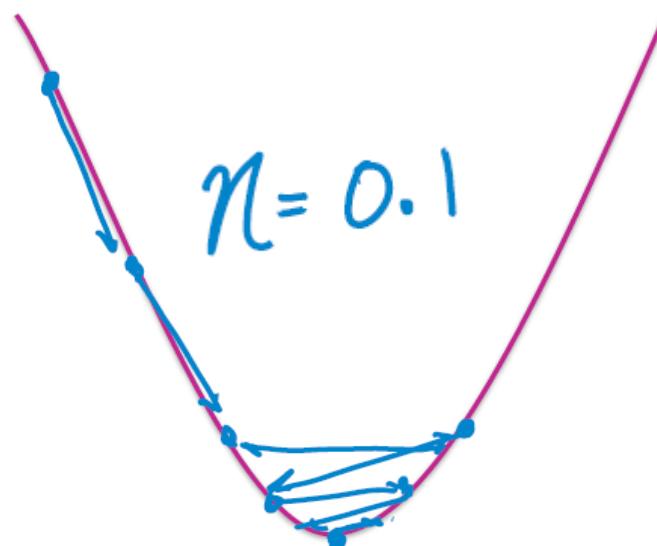
$$w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{dg}{dw} \Big|_{w^{(t)}}$$

# Choosing the step size (stepsize schedule)

22

## Fixed

Works well for strongly convex functions

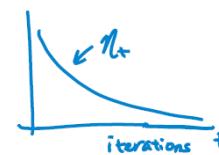
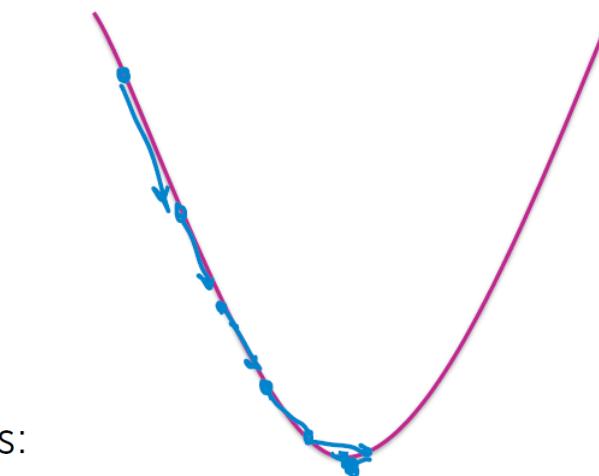


## Varying

Common choices:

$$\eta_t = \frac{\alpha}{t}$$

$$\eta_t = \frac{\alpha}{\sqrt{t}}$$



Try not to decrease  $\eta$  too fast

# Convergence criteria

23

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*

That will be „good enough”  
value of  $\epsilon$  depends on the data we are looking at

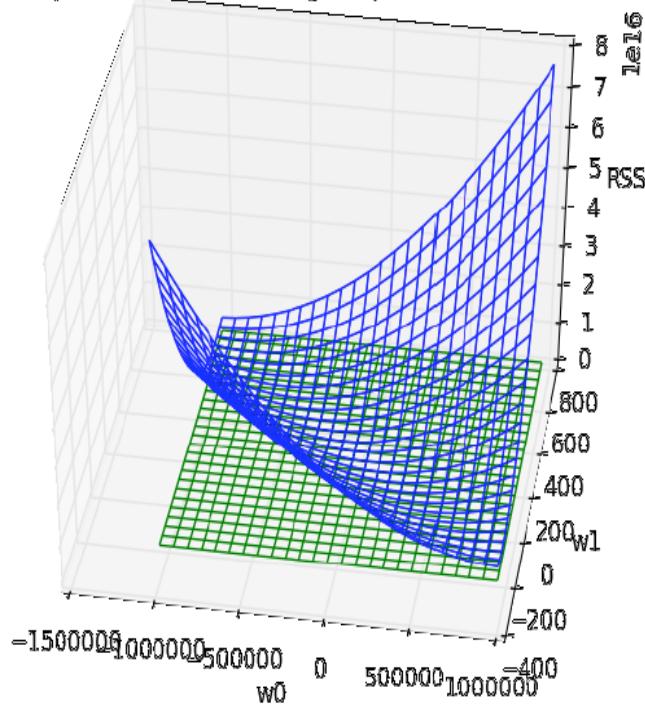
Algorithm:

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

# Moving to multiple dimensions

24

3D plot of RSS with tangent plane at minimum



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

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

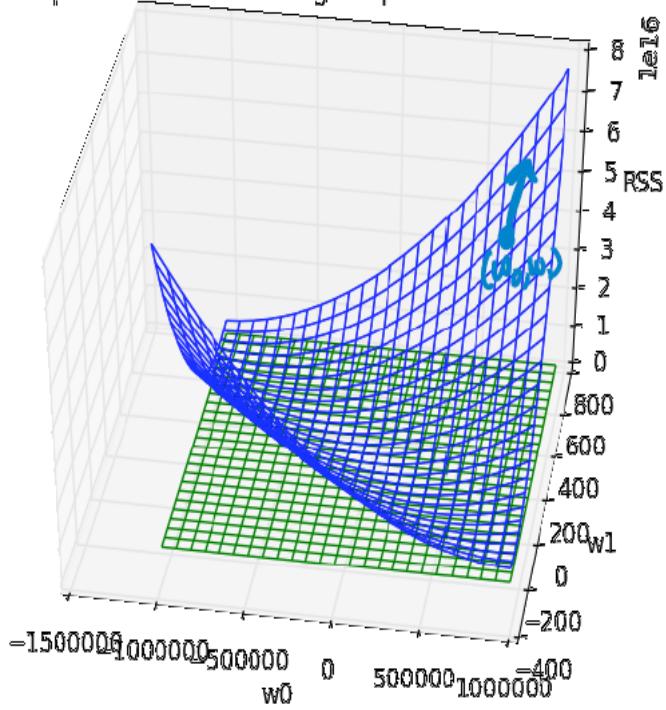
( $p+1$ ) -dimensional vector

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

# Gradient example

25

3D plot of RSS with tangent plane at minimum



$$g(\mathbf{w}) = 5w_0 + 10w_0w_1 + 2w_1^2$$

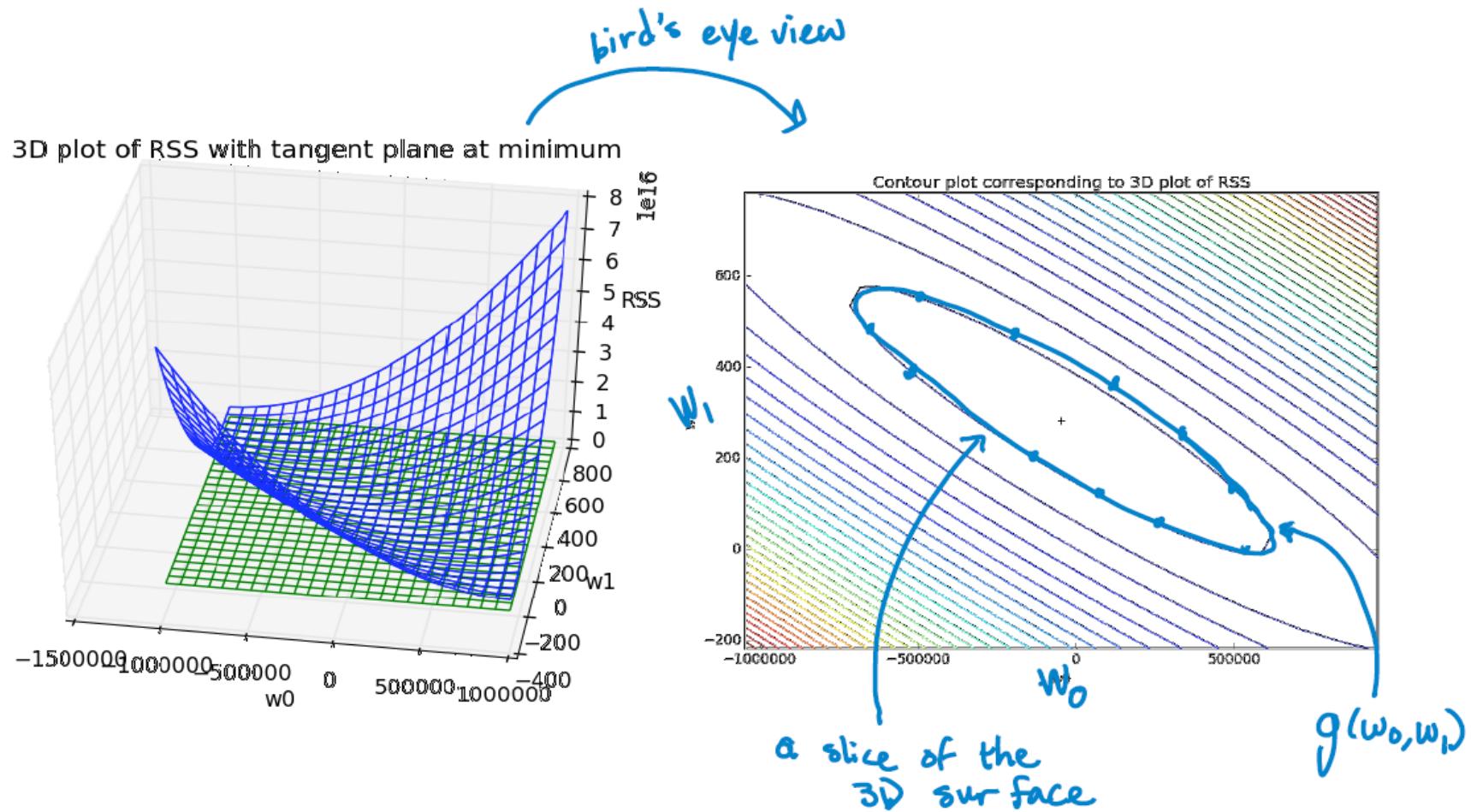
$$\frac{\partial g}{\partial w_0} = 5 + 10w_1$$

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

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

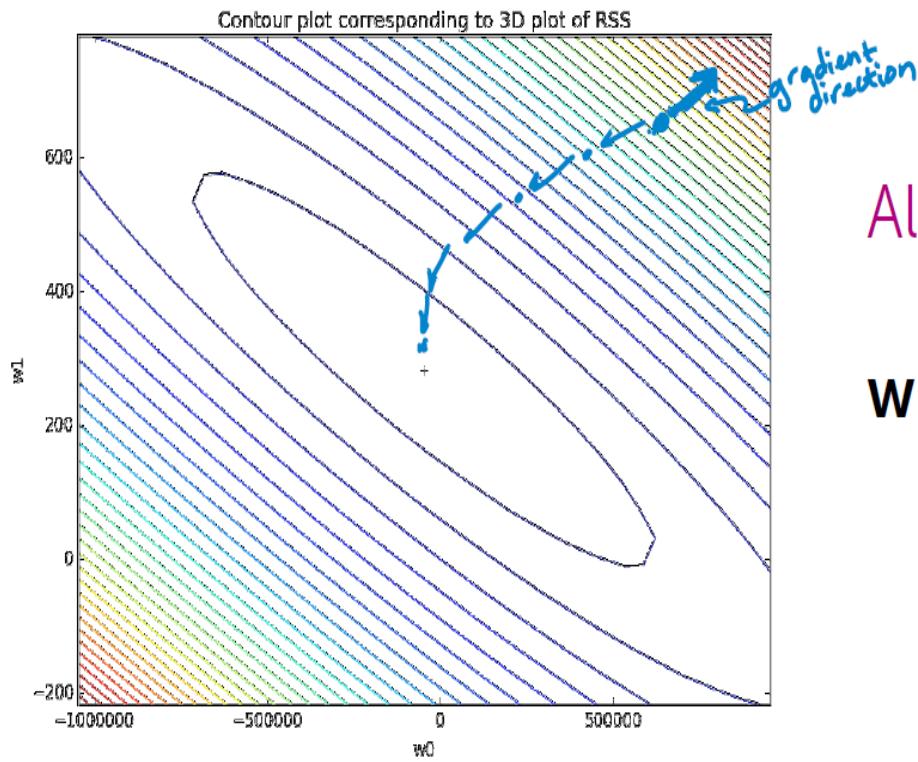
# Contour plots

26



# Gradient descent

27



Algorithm:

**while** not converged  
 $w^{(t+1)} \leftarrow w^{(t)} - \eta \nabla g(w^{(t)})$

$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \leftarrow \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} - \eta \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

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

# Compute the gradient

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

Taking the derivative w.r.t.  $w_0$

$$\begin{aligned} & \sum_{i=1}^N 2(y_i - [w_0 + w_1 x_i])^1 \cdot (-1) \\ &= -2 \sum_{i=1}^N (y_i - [w_0 + w_1 x_i]) \end{aligned}$$

Putting it together:

$$\nabla \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] x_i \end{bmatrix}$$

Taking the derivative w.r.t.  $w_1$

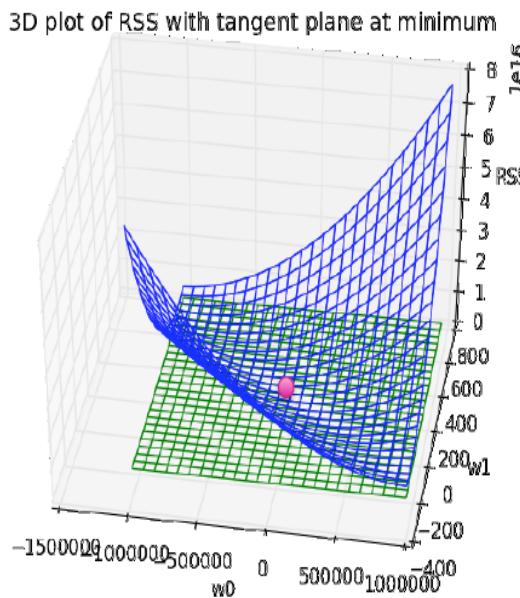
$$\begin{aligned} & \sum_{i=1}^N 2(y_i - [w_0 + w_1 x_i])^1 \cdot (-x_i) \\ &= -2 \sum_{i=1}^N (y_i - [w_0 + w_1 x_i]) \underline{x_i} \end{aligned}$$

# Approach 1: set gradient to 0

29

$$\nabla \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] x_i \end{bmatrix}$$

**This method is called „Closed form solution”**



top term:  $\hat{w}_0 = \frac{\sum_{i=1}^N y_i}{N} - \hat{w}_1 \frac{\sum_{i=1}^N x_i}{N}$

average house price  
average sq.-ft.

bottom term:

$$\sum y_i x_i - \hat{w}_0 \sum x_i - \hat{w}_1 \sum x_i^2 = 0$$
$$\hat{w}_1 = \frac{\sum y_i x_i - \frac{\sum y_i \sum x_i}{N}}{\sum x_i^2 - \frac{\sum x_i \sum x_i}{N}}$$

Note:

$$\sum_{i=1}^N y_i$$

$$\sum_{i=1}^N x_i$$

$$\sum_{i=1}^N y_i x_i$$

$$\sum_{i=1}^N x_i^2$$

# Approach 2: gradient descent

30

Interpreting the gradient:

$$\nabla_{\mathbf{w}_0, \mathbf{w}_1} \text{RSS}(\mathbf{w}_0, \mathbf{w}_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (\underline{\mathbf{w}_0 + \mathbf{w}_1 x_i})] \\ -2 \sum_{i=1}^N [y_i - (\mathbf{w}_0 + \mathbf{w}_1 x_i)] x_i \end{bmatrix} = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - \hat{y}_i(\mathbf{w}_0, \mathbf{w}_1)] \\ -2 \sum_{i=1}^N [y_i - \hat{y}_i(\mathbf{w}_0, \mathbf{w}_1)] x_i \end{bmatrix}$$

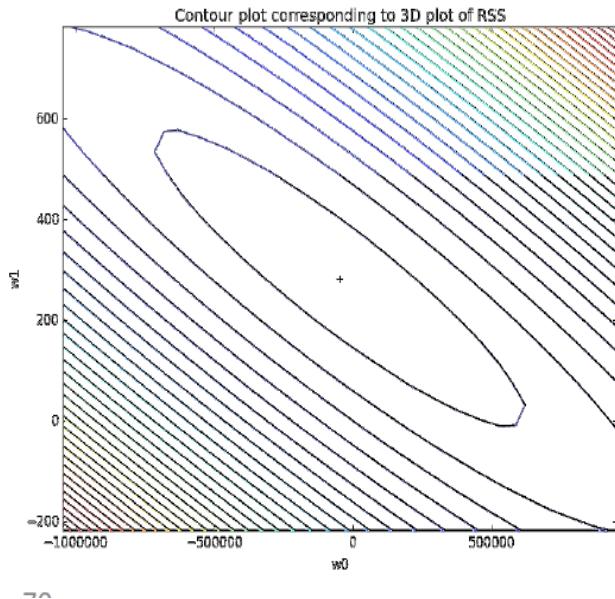
Annotations:

- actual house sales observation:  $y_i$
- predicted value:  $\hat{y}_i(\mathbf{w}_0, \mathbf{w}_1)$

# Approach 2: gradient descent

31

$$\nabla \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - \hat{y}_i(w_0, w_1)] \\ -2 \sum_{i=1}^N [y_i - \hat{y}_i(w_0, w_1)] x_i \end{bmatrix}$$



while not converged ( $\leftarrow$ )

$$\begin{bmatrix} w_0^{(t+1)} \\ w_1^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} w_0^{(t)} \\ w_1^{(t)} \end{bmatrix} + 2\eta \begin{bmatrix} \sum_{i=1}^N [y_i - \hat{y}_i(w_0^{(t)}, w_1^{(t)})] \\ \sum_{i=1}^N [y_i - \hat{y}_i(w_0^{(t)}, w_1^{(t)})] x_i \end{bmatrix}$$

If overall, underpredicting  $\hat{y}_i$ , then  $\sum [y_i - \hat{y}_i]$  is positive  
→  $w_0$  is going to increase  
similar intuition for  $w_1$ , but multiply by  $x_i$

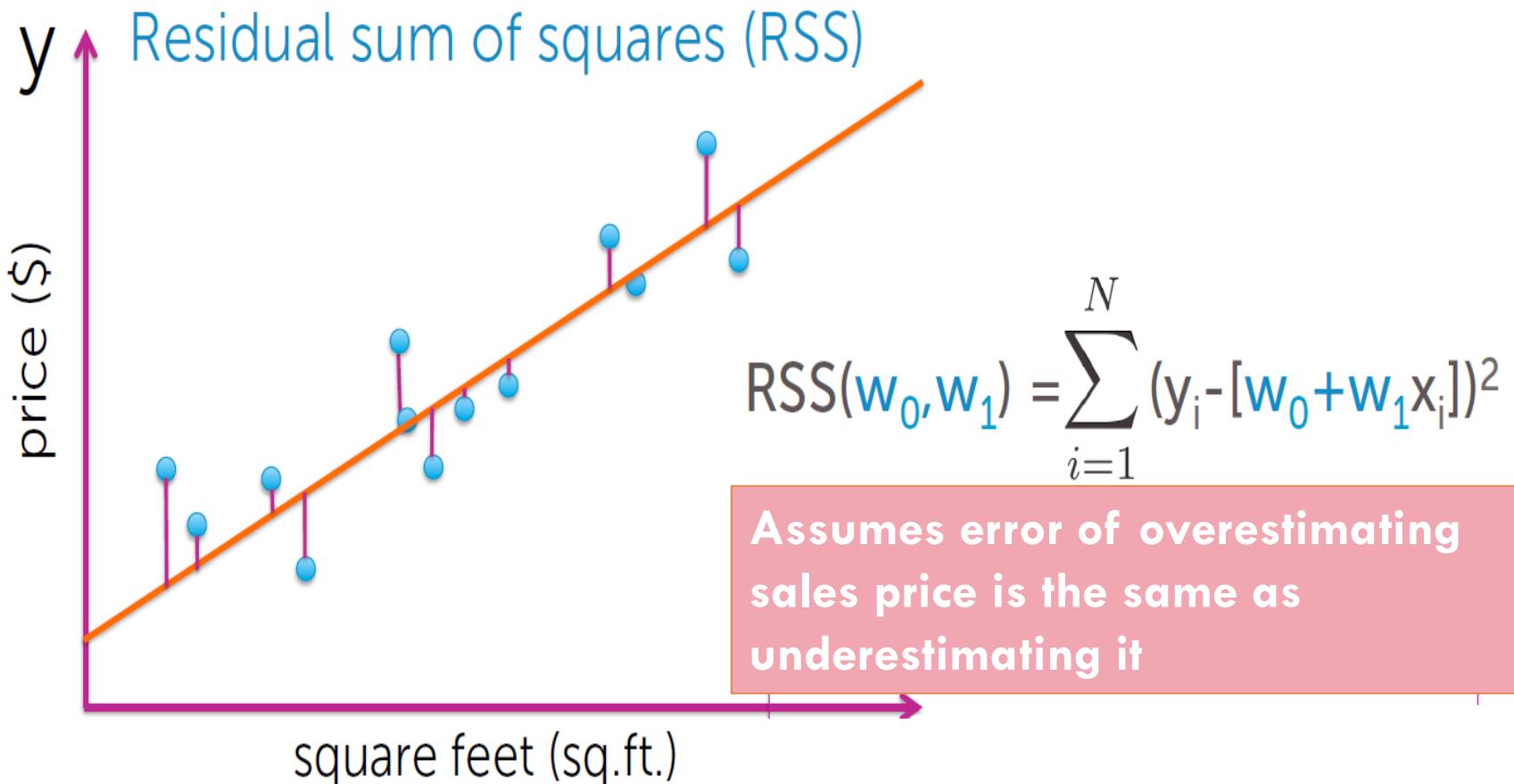
# Comparing the approaches

32

- For most ML problems, cannot solve  $\text{gradient} = 0$
- Even if solving  $\text{gradient} = 0$  is feasible, gradient descent can be more efficient
- Gradient descent relies on choosing stepsize and convergence criteria

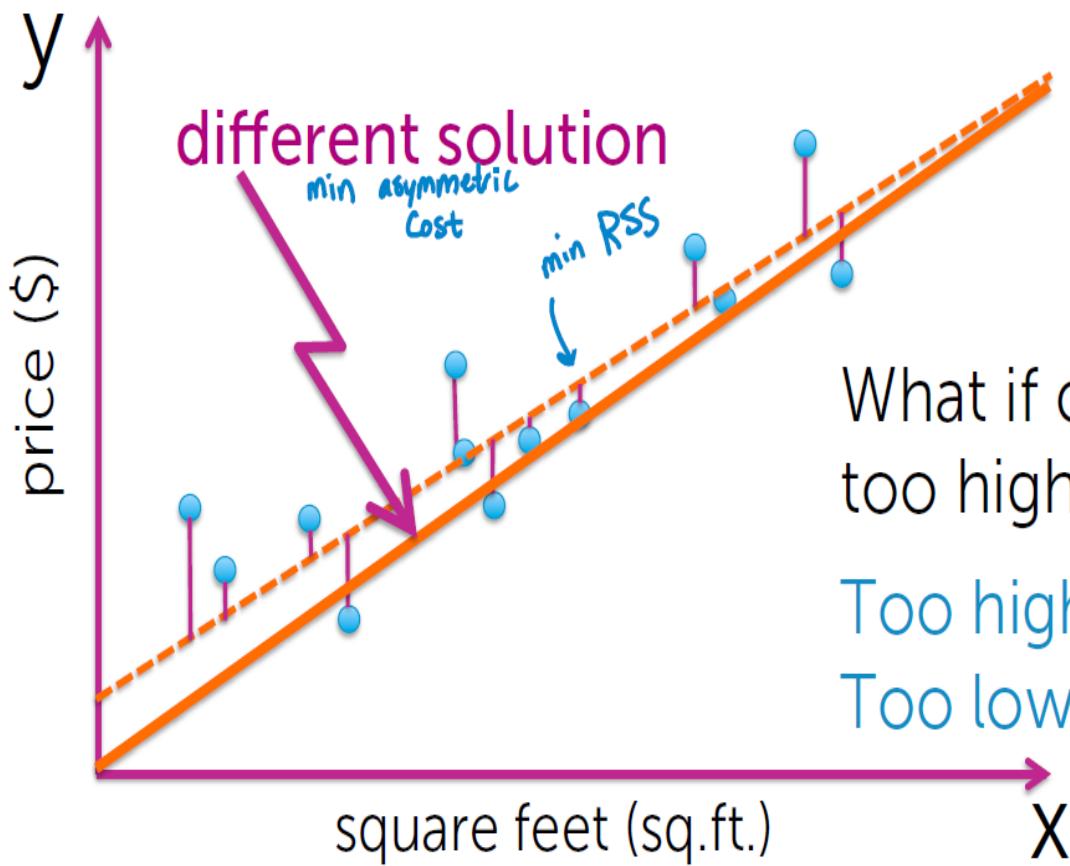
# Symmetric cost function

33



# Asymmetric cost functions

34



We can weight differently positive and negative errors in RSS calculations.

What if cost of listing house too high has bigger cost?

Too high  $\rightarrow$  no offers ( $\$=0$ )

Too low  $\rightarrow$  offers for lower  $\$$

# What you can do now

35

- Describe the input (features) and output (real-valued predictions) of a regression model
- Calculate a goodness-of-fit metric (e.g., RSS)
- Estimate model parameters to minimize RSS using gradient descent
- Interpret estimated model parameters
- Exploit the estimated model to form predictions
- Discuss the possible influence of high leverage points
- Describe intuitively how fitted line might change when assuming different goodness-of-fit metrics

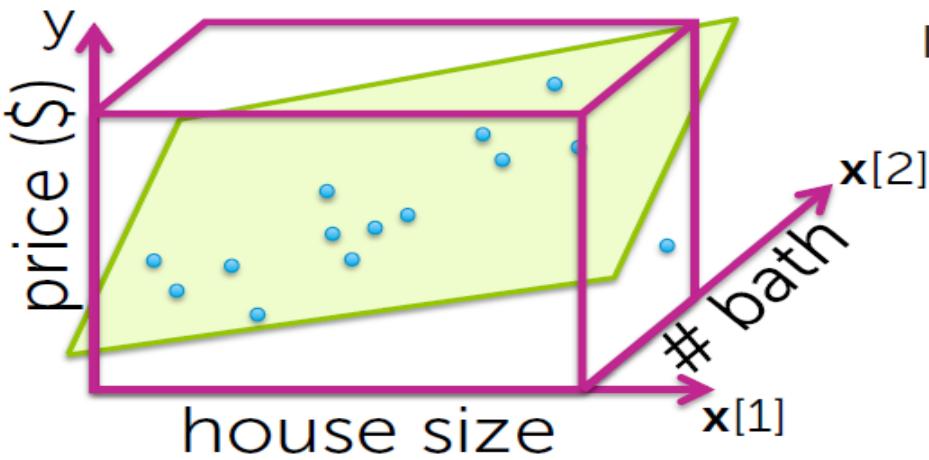
# MULTIPLE REGRESSION

# Multiple regression

37



Fit more complex relationships than just a line



Incorporate more inputs

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

# Polynomial regression

38

Model:

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

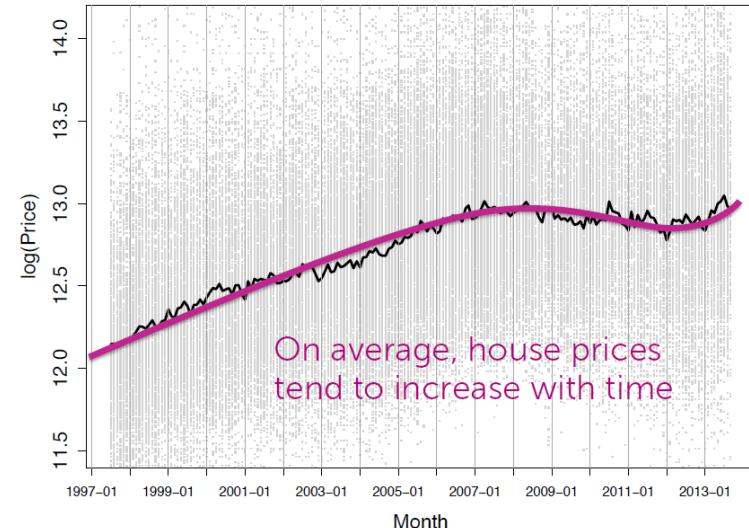
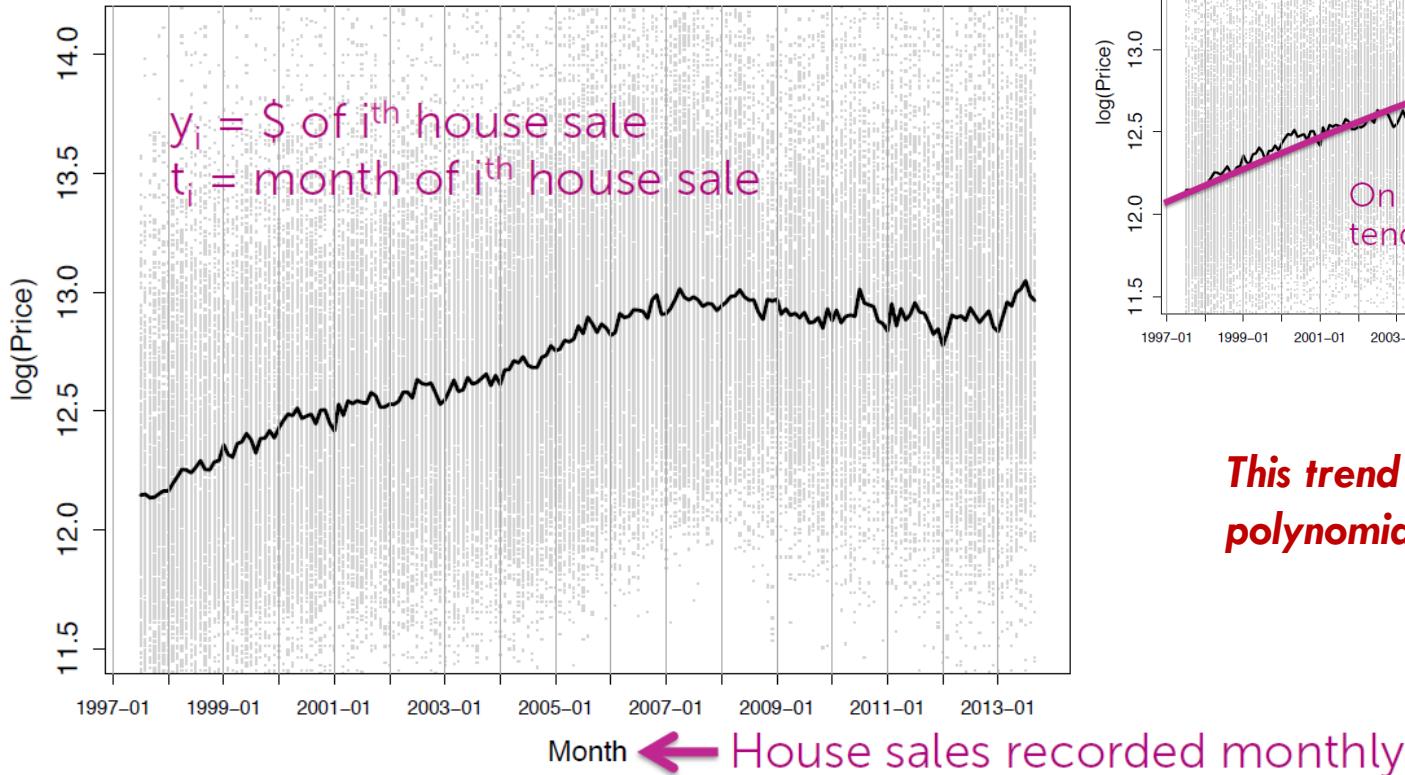


feature 1 = 1 (constant)	parameter 1 = $w_0$
feature 2 = $x$	parameter 2 = $w_1$
feature 3 = $x^2$	parameter 3 = $w_2$
...	...
feature $p+1 = x^p$	parameter $p+1 = w_p$

# Other functional forms of one input

39

## □ Trends in time series

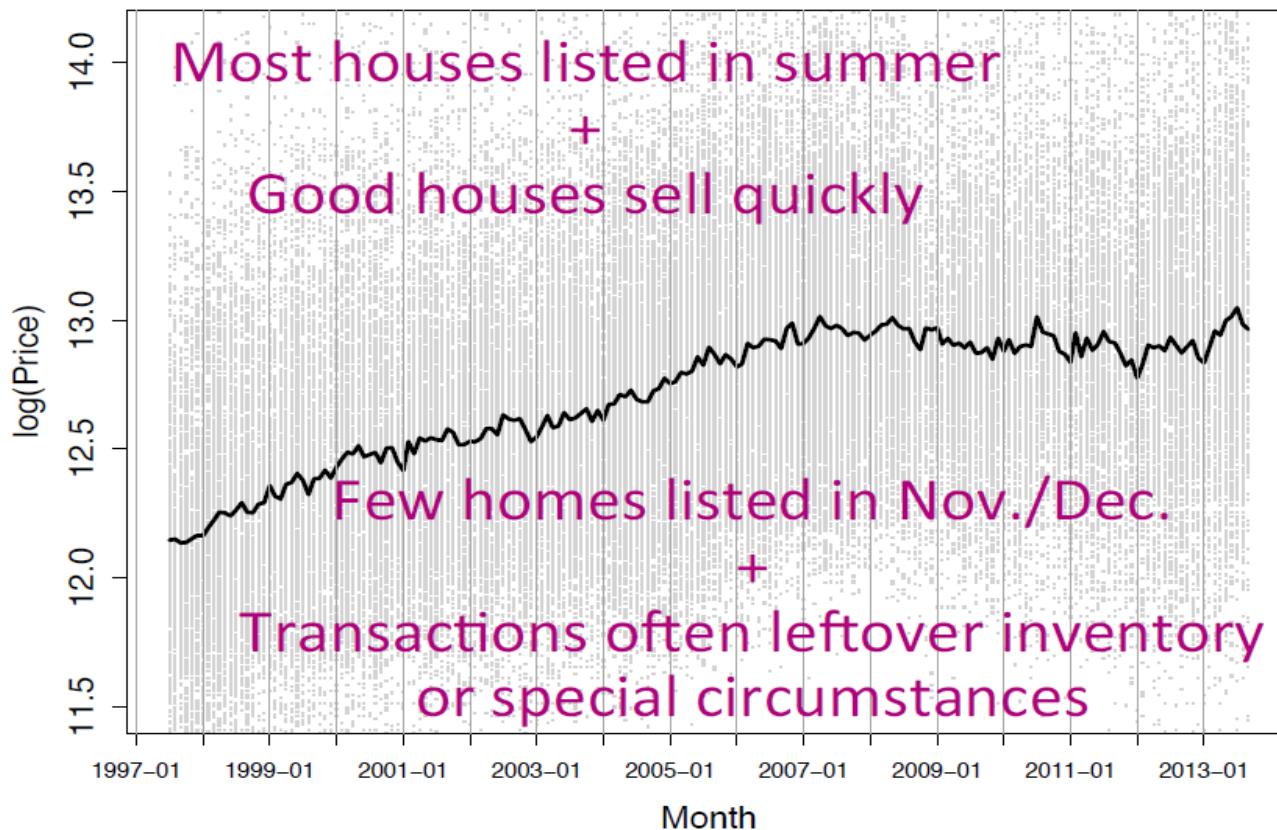


**This trend can be modeled with polynomial function.**

# Other functional forms of one input

40

## □ Seasonality



# Example of detrending

41

Model:

$$y_i = w_0 + w_1 t_i + w_2 \sin(2\pi t_i / 12 - \Phi) + \varepsilon_i$$

Linear trend

Seasonal component =  
Sinusoid with period 12  
(resets annually)

Unknown phase/shift

Trigonometric identity:  $\sin(a-b) = \sin(a)\cos(b) - \cos(a)\sin(b)$

$$\rightarrow \sin(2\pi t_i / 12 - \Phi) = \sin(2\pi t_i / 12)\cos(\Phi) - \cos(2\pi t_i / 12)\sin(\Phi)$$

# Example of detrending

42

Equivalently,

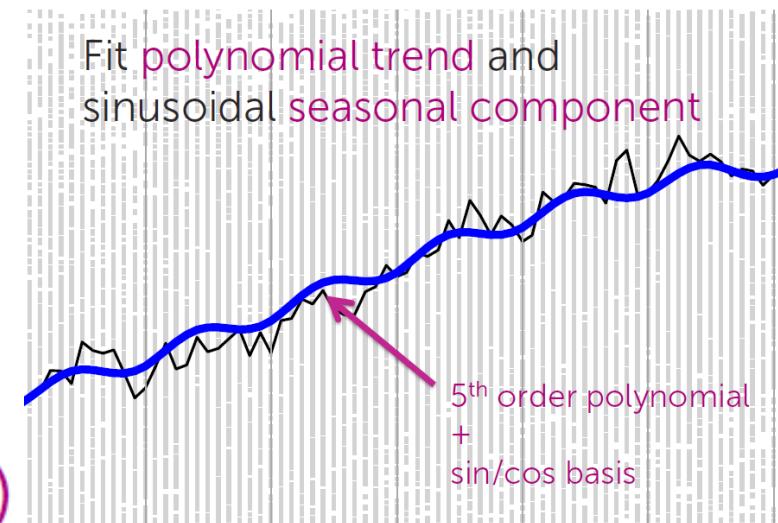
$$y_i = w_0 + w_1 t_i + w_2 \sin(2\pi t_i / 12) + w_3 \cos(2\pi t_i / 12) + \epsilon_i$$

*feature 1 = 1 (constant)*

*feature 2 = t*

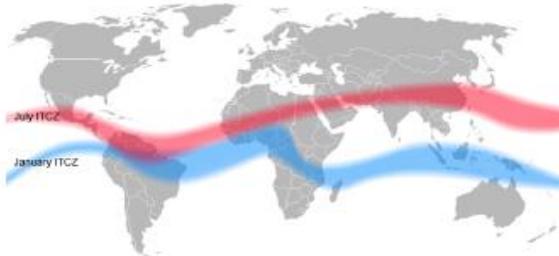
*feature 3 = sin(2πt/12)*

*feature 4 = cos(2πt/12)*

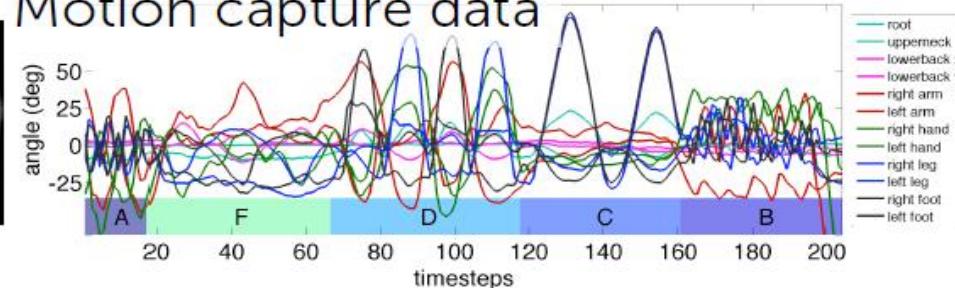
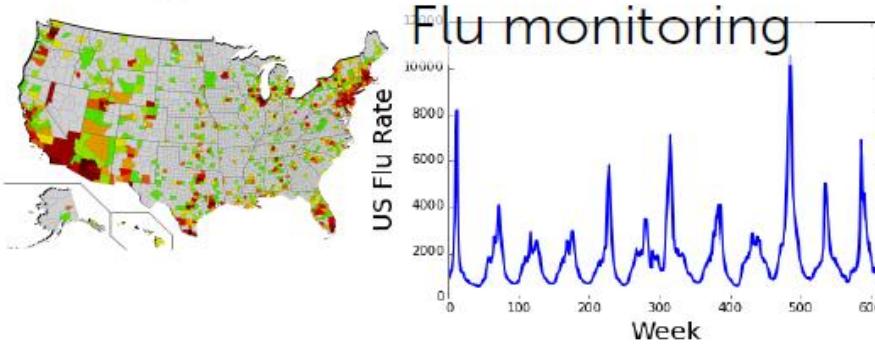


# Other examples of seasonality

43



Weather modeling  
(e.g., temperature, rainfall)



3

# Generic basic expansion

44

Model:

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

*feature 1* =  $h_0(x)$ ... often 1 (constant)

*feature 2* =  $h_1(x)$ ... e.g.,  $x$

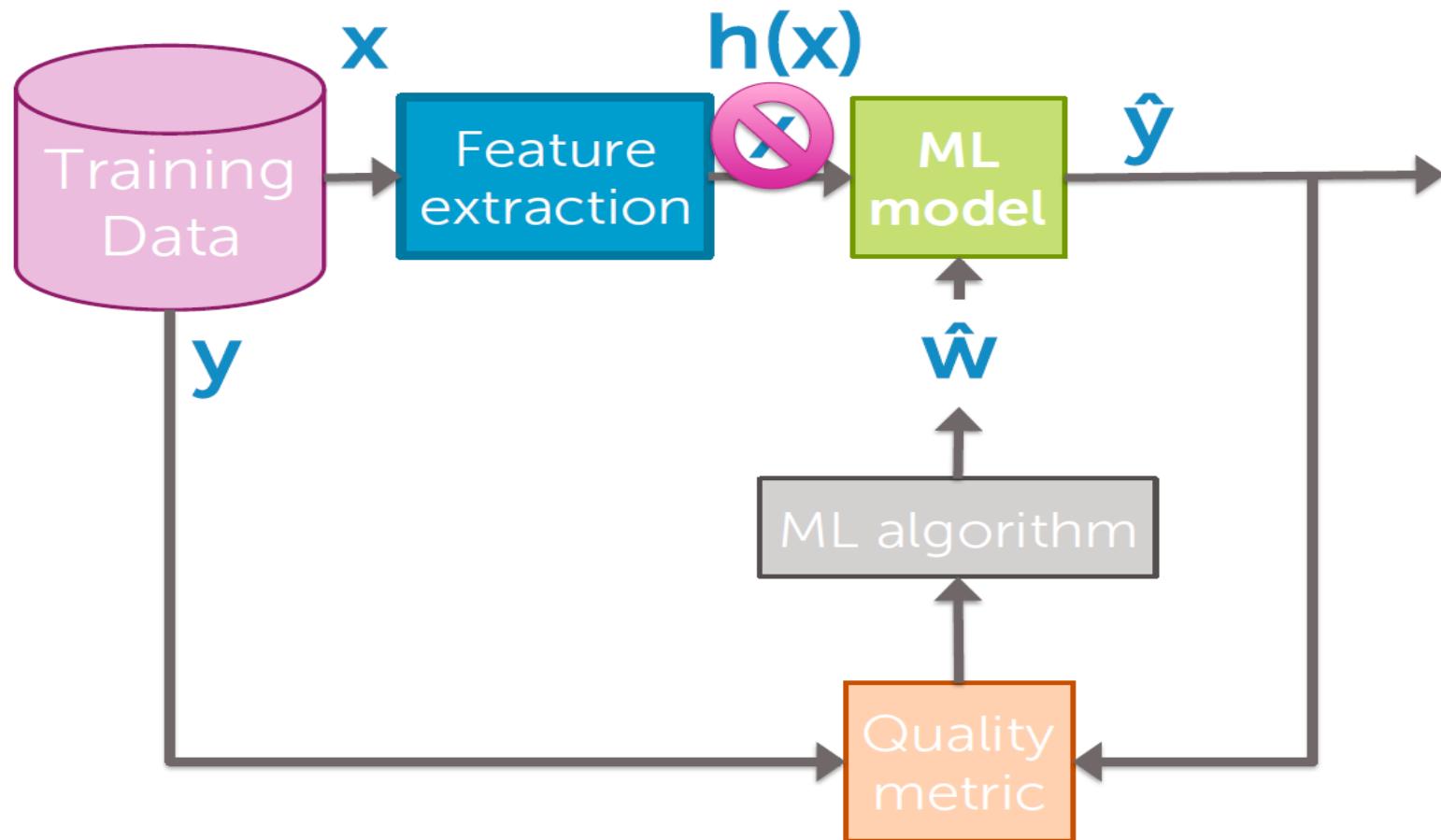
*feature 3* =  $h_2(x)$ ... e.g.,  $x^2$  or  $\sin(2\pi x/12)$

...

*feature D+1* =  $h_D(x)$ ... e.g.,  $x^p$

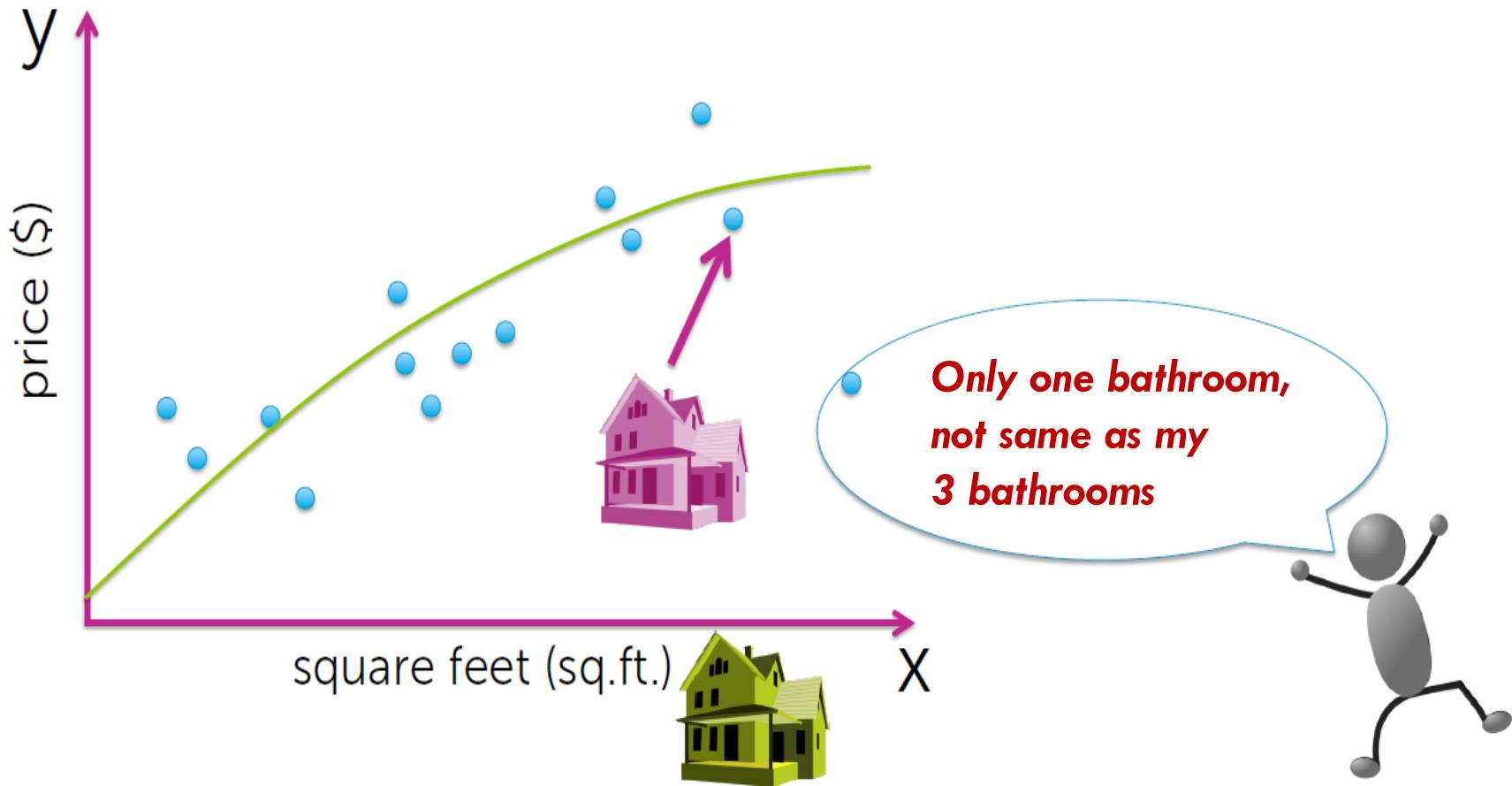
# More realistic flow chart

45



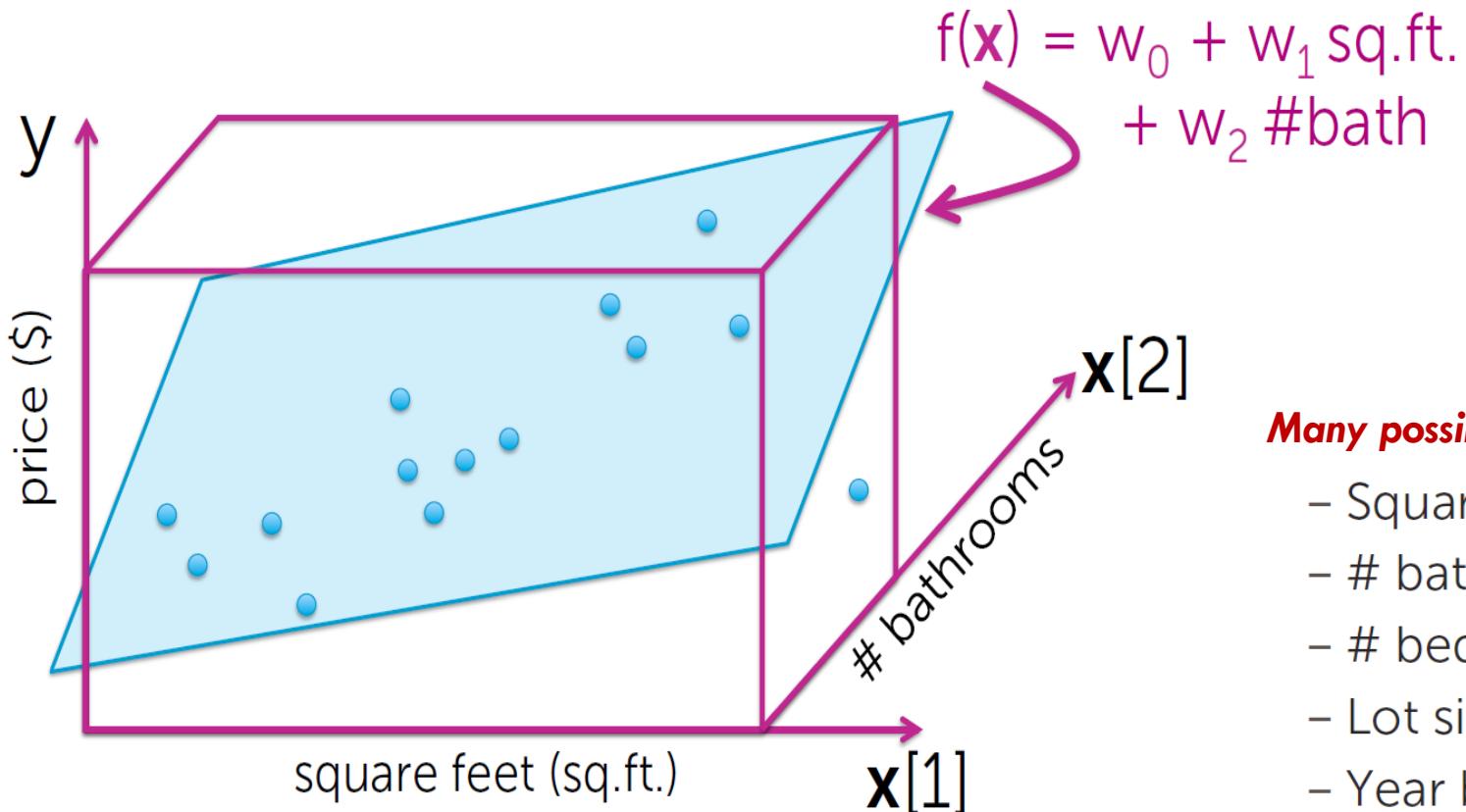
# Incorporating multiple inputs

46



# Incorporating multiple inputs

47



## Many possible inputs

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

# General notation

48

Output:  $y \leftarrow$  scalar

Inputs:  $\mathbf{x} = (\mathbf{x}[1], \mathbf{x}[2], \dots, \mathbf{x}[d])$   


Notational conventions:

$\mathbf{x}[j]$  =  $j^{\text{th}}$  input (scalar)

$h_j(\mathbf{x})$  =  $j^{\text{th}}$  feature (scalar)

$\mathbf{x}_i$  = input of  $i^{\text{th}}$  data point (vector)

$\mathbf{x}_i[j]$  =  $j^{\text{th}}$  input of  $i^{\text{th}}$  data point (scalar)

# Simple hyperplane

49

Model:

$$y_i = w_0 + w_1 x_i[1] + \dots + w_d x_i[d] + \epsilon_i$$

Noise term

*feature 1 = 1*

*feature 2 =  $x[1]$  ... e.g., sq. ft.*

*feature 3 =  $x[2]$  ... e.g., #bath*

...

*feature  $d+1 = x[d]$  ... e.g., lot size*

# More generally: D-dimensional curve

50

Model:

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

More on notation

# observations  $(\mathbf{x}_i, y_i)$  :  $N$

# inputs  $\mathbf{x}[j]$  :  $d$

# features  $h_j(\mathbf{x})$  :  $D$

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

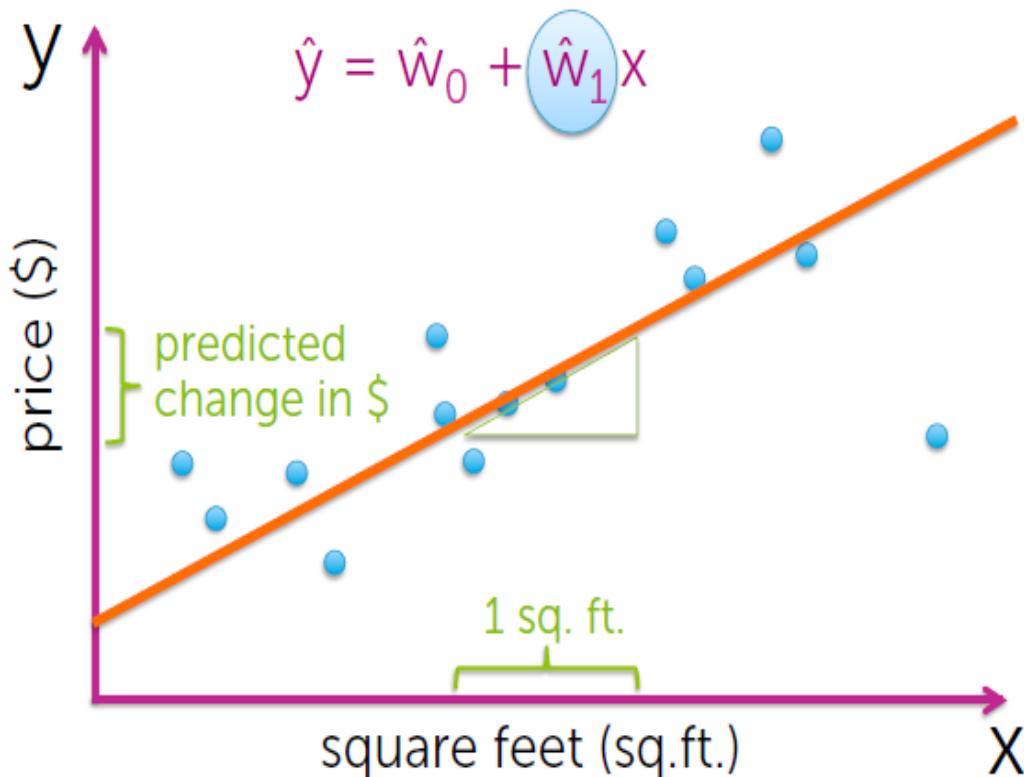
...

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

# Interpreting coefficients

51

## Simple linear regression



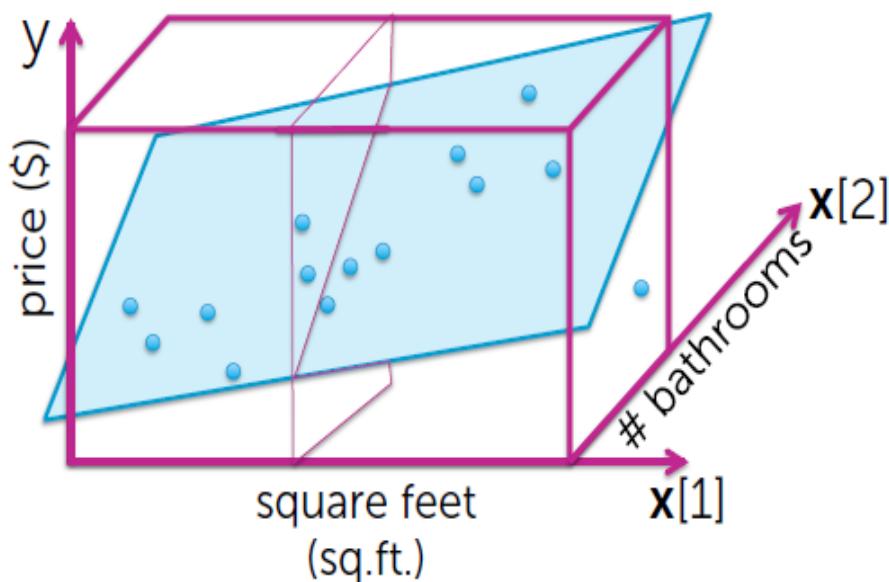
# Interpreting coefficients

52

Two linear features

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x[1] + \hat{w}_2 x[2]$$

fix



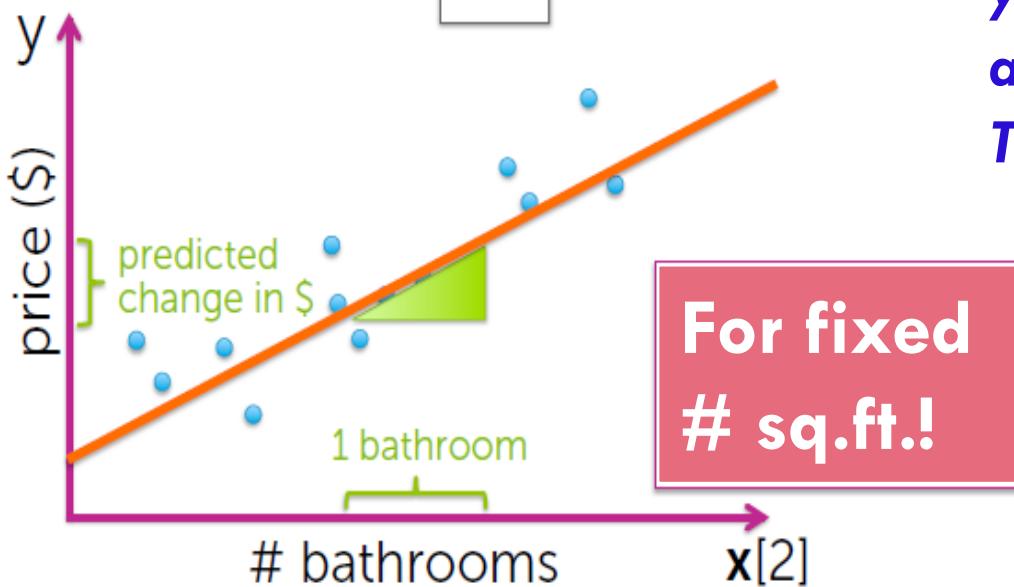
# Interpreting coefficients

53

Two linear features

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x[1] + \hat{w}_2 x[2]$$

fix



**But...**

*increasing #bathrooms  
for fixed #sq.ft will make  
your bedrooms smaller  
and smaller.*

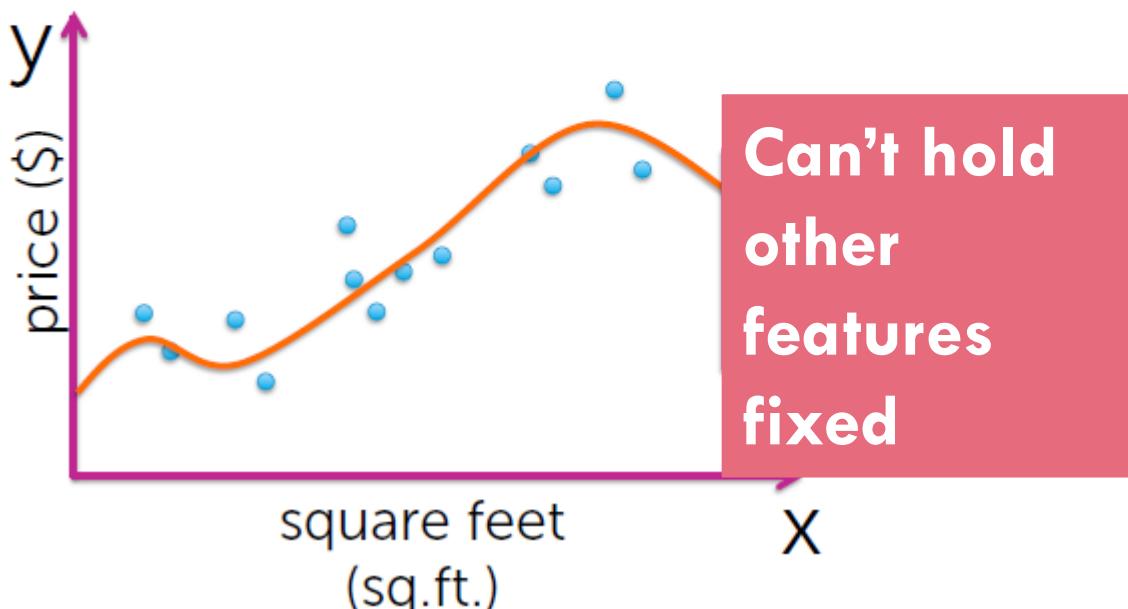
*Think about interpretation.*

# Interpreting coefficients

54

## Polynomial regression

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x + \dots + \hat{w}_j x^j + \dots + \hat{w}_p x^p$$



*Then ...  
can't interpret  
coefficients*

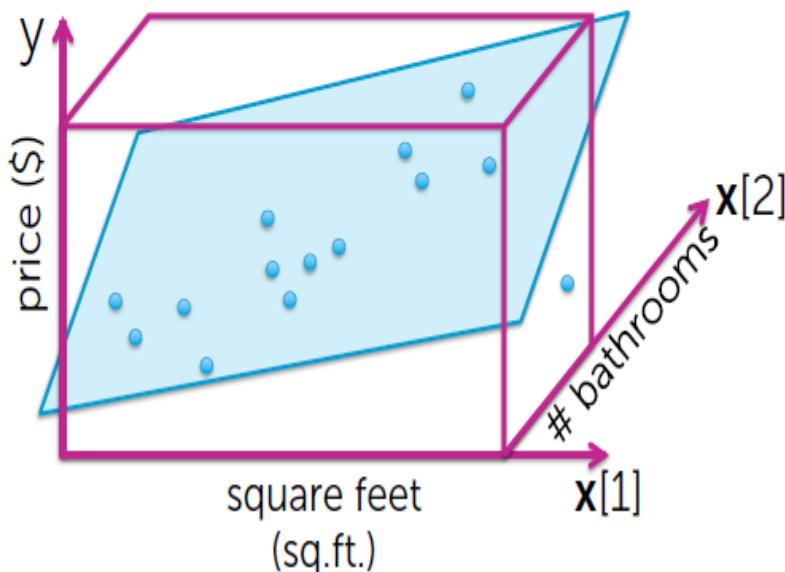
# Interpreting coefficients

55

## Multiple linear features

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x[1] + \dots + \hat{w}_j x[j] + \dots + \hat{w}_d x[d]$$

fix      fix      fix      fix



**But...**

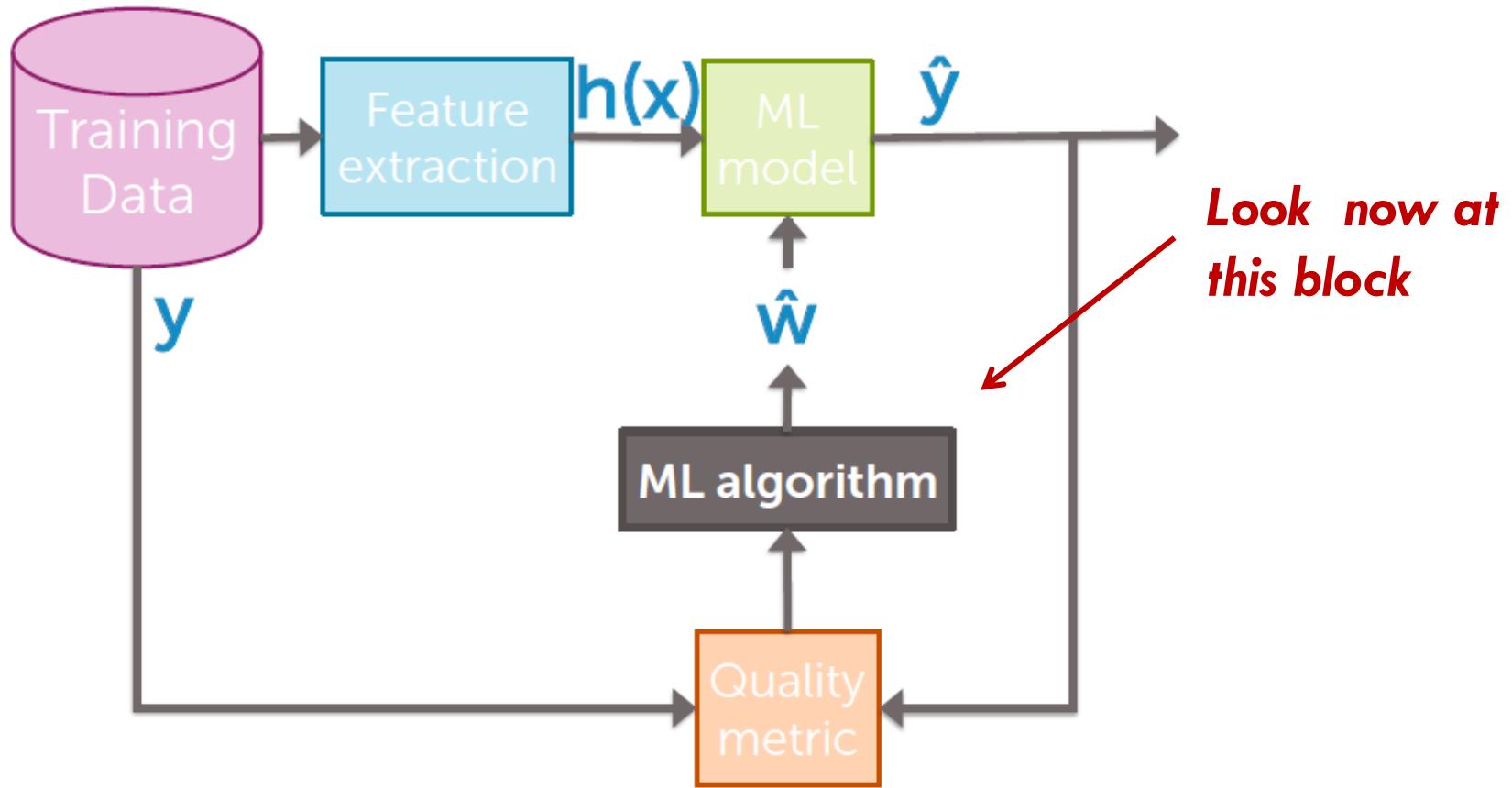
*increasing #bedrooms  
for fixed #sq.ft will make  
your bedrooms smaller  
and smaller.*

*You can end with negative  
coefficient. Might not be so  
if you removed #sq.ft from  
the model.*

***Think about interpretation  
in context of the model  
you put in.***

# Fitting in D-dimensions

56



# Rewriting in vector notation

57

For observation i

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

$$\begin{aligned} y_i &= \underbrace{\begin{bmatrix} w_0 & w_1 & w_2 & \dots & w_D \end{bmatrix}}_{w^T} \underbrace{h(x_i)}_{\text{vector}} + \varepsilon_i \\ &= \underbrace{w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i)}_{\text{scalar}} + \varepsilon_i \\ &= w^T h(x_i) + \varepsilon_i \end{aligned}$$
$$\begin{aligned} h(x_i) &= \begin{bmatrix} h_0(x_i) \\ h_1(x_i) \\ h_2(x_i) \\ \vdots \\ h_D(x_i) \end{bmatrix} \\ &= \begin{bmatrix} h^T(x_i) \end{bmatrix} \end{aligned}$$
$$\begin{aligned} h^T(x_i) &= \begin{bmatrix} h_0(x_i) & h_1(x_i) & \dots & h_D(x_i) \end{bmatrix} \\ &= h_0(x_i) w_0 + h_1(x_i) w_1 + \dots + h_D(x_i) w_D + \varepsilon_i \\ &= \begin{bmatrix} w_0 & w_1 & \dots & w_D \end{bmatrix} \begin{bmatrix} h^T(x_i) \end{bmatrix} + \varepsilon_i \end{aligned}$$

# Rewriting in matrix notation

58

For all observations together

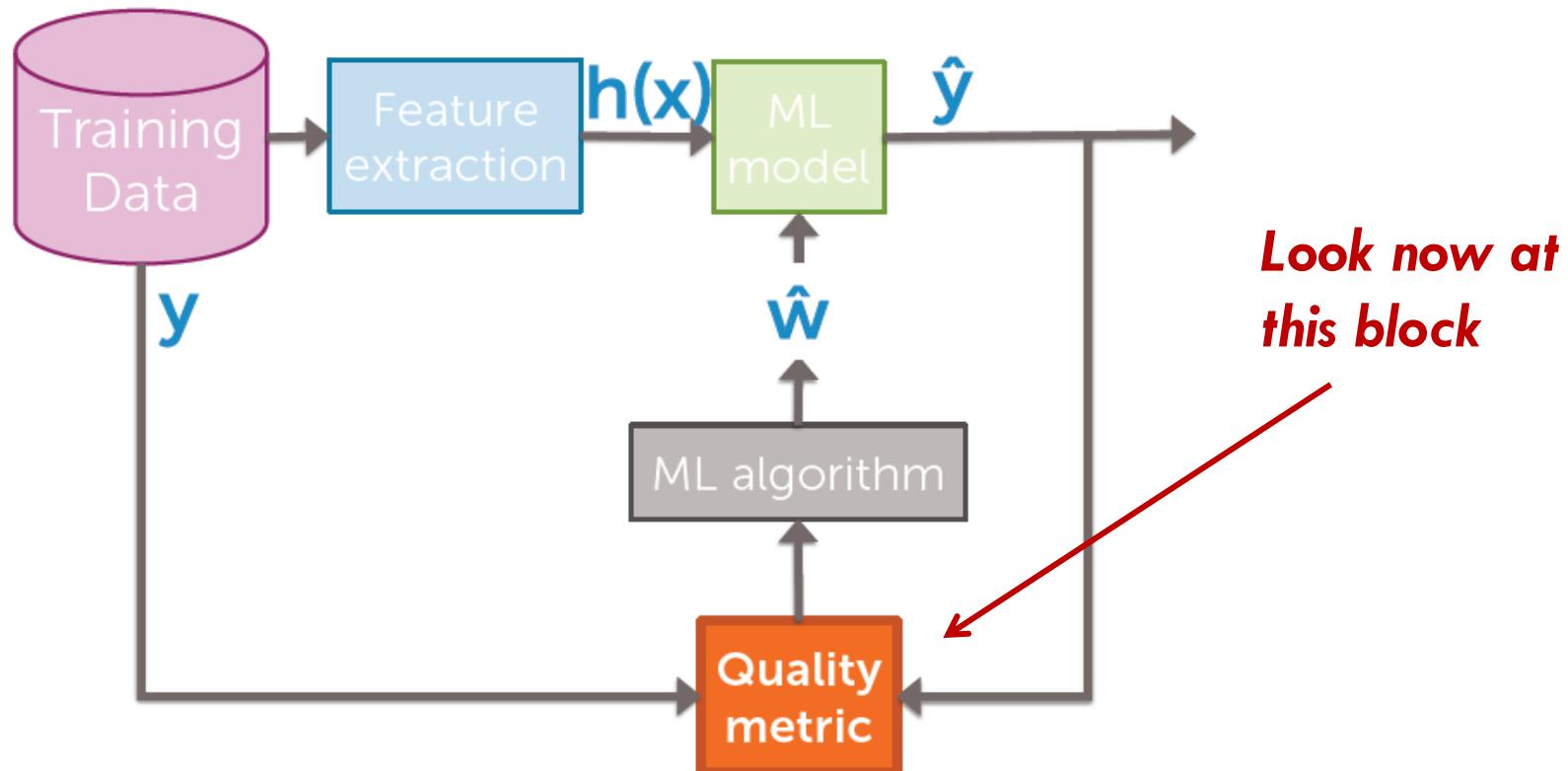
$$\begin{matrix} \mathbf{y} \\ y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{matrix} = \mathbf{H} \begin{matrix} \mathbf{w} \\ w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{matrix} + \begin{matrix} \boldsymbol{\epsilon} \\ \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_N \end{matrix}$$

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

Here is our  
ML algorithm

# Fitting in D-dimensions

59



46

©2015 Emily Fox & Carlos Guestrin

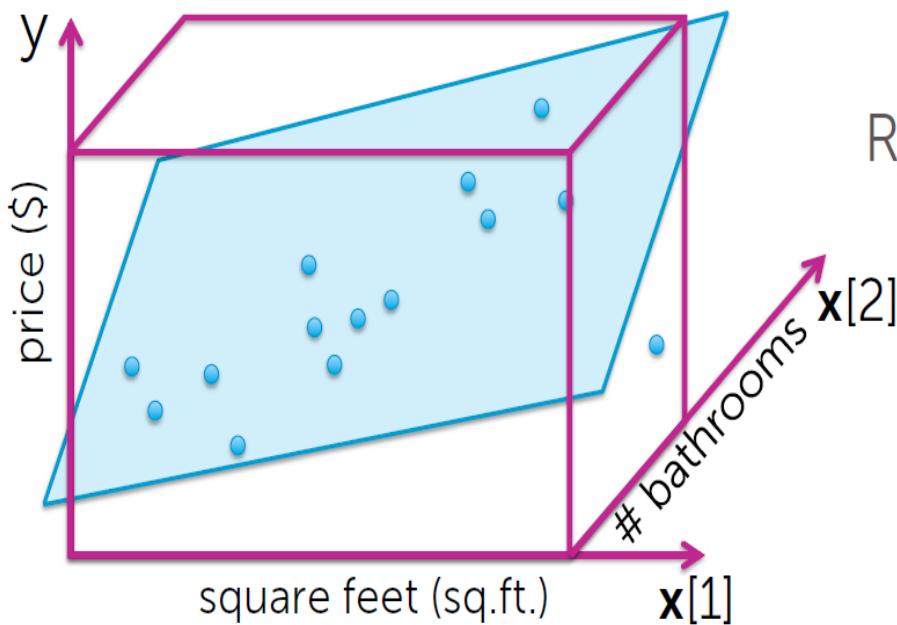
Machine Learning Specialization

26/10, 2/11, 9/11 2022

# Cost function in D-dimension

60

## RSS in vector notation



$$\text{RSS}(\underline{w}) = \sum_{i=1}^N (y_i - \underbrace{h^T(x_i) \underline{w}}_{\hat{y}_i(\underline{w})})^2$$
$$\hat{y}_i = \begin{bmatrix} h_0(x_i) & h_1(x_i) & \dots & h_p(x_i) \end{bmatrix}$$
$$\underline{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$$

# Cost function in D-dimension

61

## RSS in matrix notation

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

Why? (part 1)

$$\begin{matrix} \hat{\mathbf{y}} \\ \vdots \\ \hat{\mathbf{y}}_n \end{matrix} = \mathbf{H} \begin{matrix} \mathbf{w}_0 \\ \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_D \end{matrix}$$

$$\hat{\mathbf{y}} = \mathbf{H}\mathbf{w}$$
$$(\mathbf{y} - \tilde{\mathbf{H}}\mathbf{w}) = (\mathbf{y} - \hat{\mathbf{y}}) = \begin{bmatrix} \text{residual}_1 \\ \text{residual}_2 \\ \vdots \\ \text{residual}_N \end{bmatrix}$$

$\hat{\mathbf{y}}_i = y_i - \hat{y}_i$

# Regression model for D-dimension

62

## RSS in matrix notation

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

Why? (part 2)

residual <sub>1</sub>	residual <sub>2</sub>	residual <sub>3</sub>	...	residual <sub>N</sub>	residual <sub>1</sub>	residual <sub>2</sub>	residual <sub>3</sub>	...	residual <sub>N</sub>

$$\begin{aligned} & (\text{residual}_1^2 + \text{residual}_2^2 + \dots + \text{residual}_N^2) \\ &= \sum_{i=1}^N \text{residual}_i^2 \\ &\triangleq \text{RSS}(\mathbf{w}) \end{aligned}$$

# Regression model for D-dimension

63

## Gradient of RSS

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

Why? By analogy to 1D case:

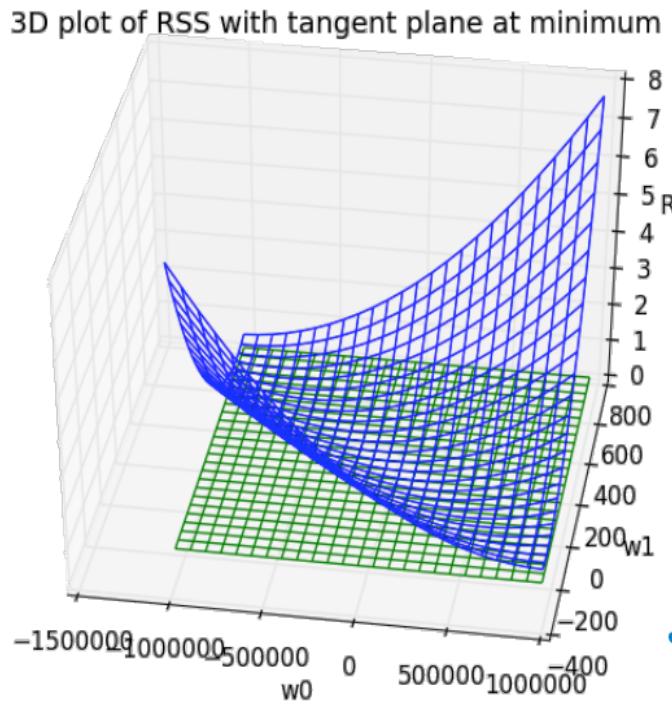
$$\begin{aligned}\frac{d}{dw} (y-hw)(y-hw) &= \frac{d}{dw} (y-hw)^2 = 2 \cdot (y-hw)' (-h) \\ &= -2h(y-hw)\end{aligned}$$

*↑  
scalars*

# Regression model for D-dimension

64

## Approach 1: set gradient to zero



**Closed form solution**

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

Solve for  $\mathbf{w}$ :

$$-\cancel{2}\mathbf{H}^T\mathbf{y} + \cancel{2}\mathbf{H}^T\mathbf{H}\hat{\mathbf{w}} = 0$$

$$\mathbf{H}^T\mathbf{H}\hat{\mathbf{w}} = \mathbf{H}^T\mathbf{y}$$

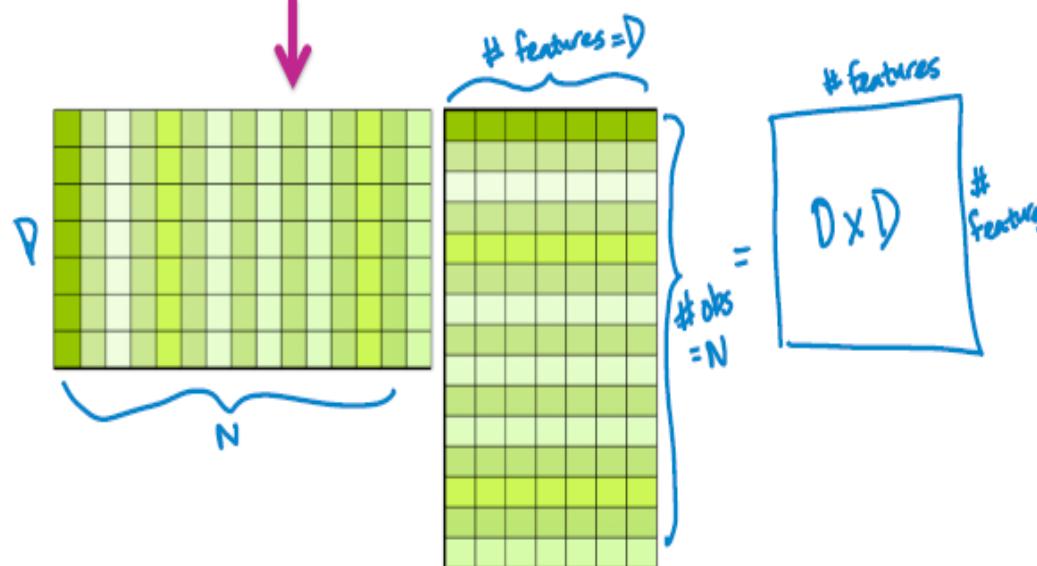
$$\underbrace{(\mathbf{H}^T\mathbf{H})^{-1}}_{I} \mathbf{H}^T\mathbf{H}\hat{\mathbf{w}} = (\mathbf{H}^T\mathbf{H})^{-1}\mathbf{H}^T\mathbf{y}$$

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

# Closed-form solution

65

$$\hat{w} = (\underbrace{H^T H}_{\text{D}})^{-1} H^T y$$



***This matrix might not be invertible.***

Invertible if:  
In most cases is  $N > D$

Complexity of inverse:

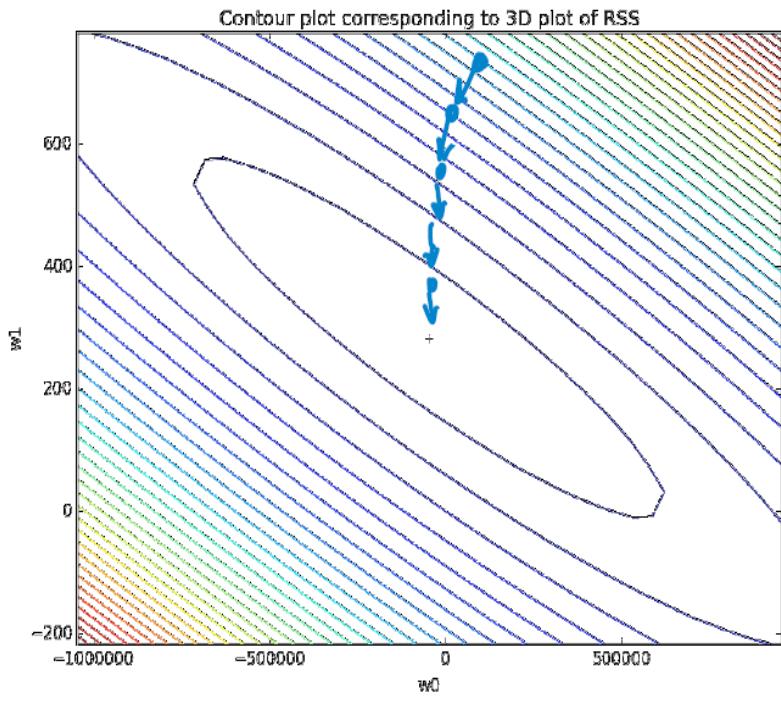
$$O(D^3)$$

***This might not be CPU feasible.***

# Regression model for D-dimension

66

## Approach 2: gradient descent



We initialise our solution somewhere and then ...

**while** not converged

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla \text{RSS}(\mathbf{w}^{(t)})$$
$$= \mathbf{w}^{(t)} + 2\eta \mathbf{H}^T (\mathbf{y} - \hat{\mathbf{y}}(\mathbf{w}^{(t)}))$$

# Gradient descent

67

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N (y_i - h(\mathbf{x}_i)^T \mathbf{w})^2$$
$$= \sum_{i=1}^N (y_i - w_0 h_0(x_i) - w_1 h_1(x_i) - \dots - w_D h_D(x_i))^2$$

Partial with respect to  $w_j$ .

$$\begin{aligned} & \sum_{i=1}^N 2(y_i - w_0 h_0(x_i) - w_1 h_1(x_i) - \dots - w_D h_D(x_i)) \\ & \quad \cdot (-\underline{h_j(x_i)}) \\ &= -2 \sum_{i=1}^N h_j(x_i) (y_i - h(\mathbf{x}_i)^T \mathbf{w}) \end{aligned}$$

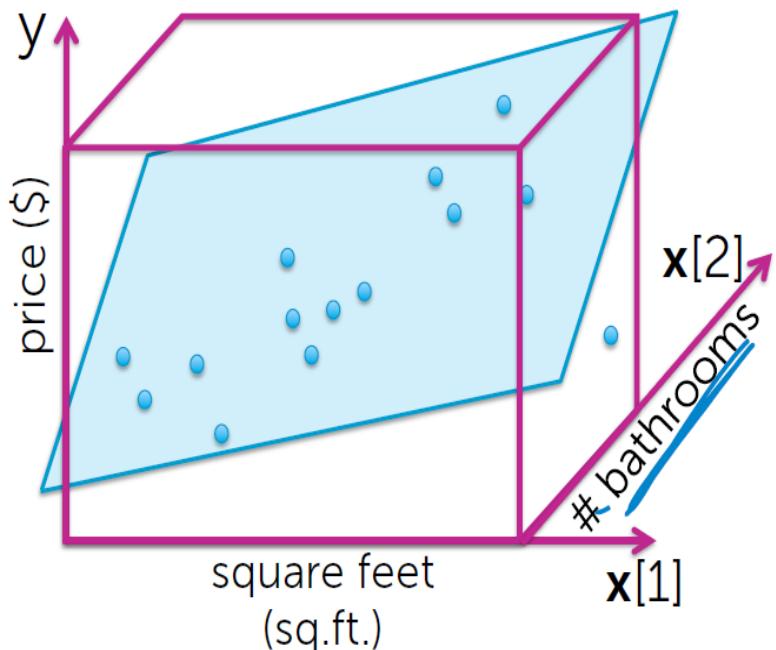
Update to  $j^{\text{th}}$  feature weight:

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta \left( -2 \sum_{i=1}^N h_j(x_i) (y_i - \underbrace{h^T(\mathbf{x}_i) \mathbf{w}^{(t)}}_{\hat{y}_i(\mathbf{w}^{(t)})}) \right)$$

# Regression model for D-dimension

68

## Interpreting elementwise



Update to  $j^{\text{th}}$  feature weight:

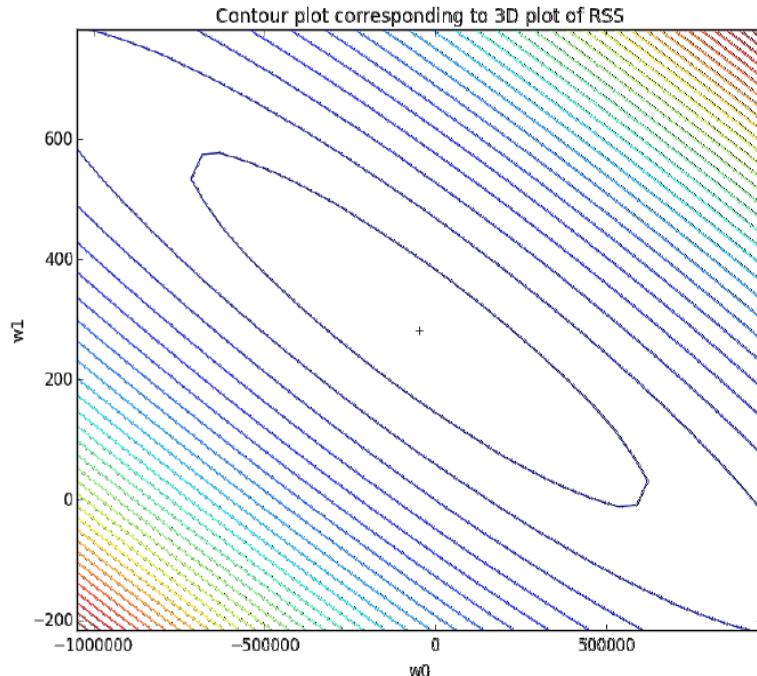
$$w_j^{(t+1)} \leftarrow w_j^{(t)} + 2\eta \sum_{i=1}^N h_i(\mathbf{x}_i)(y_i - \hat{y}_i(w^{(t)}))$$

If underestimating impact of # bath ( $\hat{w}_j^{(t)}$  is too small)  
then  $(y_i - \hat{y}_i(w^{(t)}))$  on average  
weighted by # bath will be positive  
 $\Rightarrow w_j^{(t+1)} > w_j^{(t)}$  (increase)

# Summary of gradient descent

69

**Extremely useful algorithm in several applications**



init  $\mathbf{w}^{(1)} = 0$  (or randomly, or smartly),  $t = 1$   
while  $\|\nabla \text{RSS}(\mathbf{w}^{(t)})\| > \epsilon$   $\epsilon$  tolerance  
 $\sqrt{\partial_1^2 + \dots + \partial_D^2}$   
for  $j = 0, \dots, D$   
 $\partial_j = -2 \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)}))$   
 $\mathbf{w}_j^{(t+1)} \leftarrow \mathbf{w}_j^{(t)} - \eta \partial_j$   
 $t \leftarrow t + 1$

# What you can do now

70

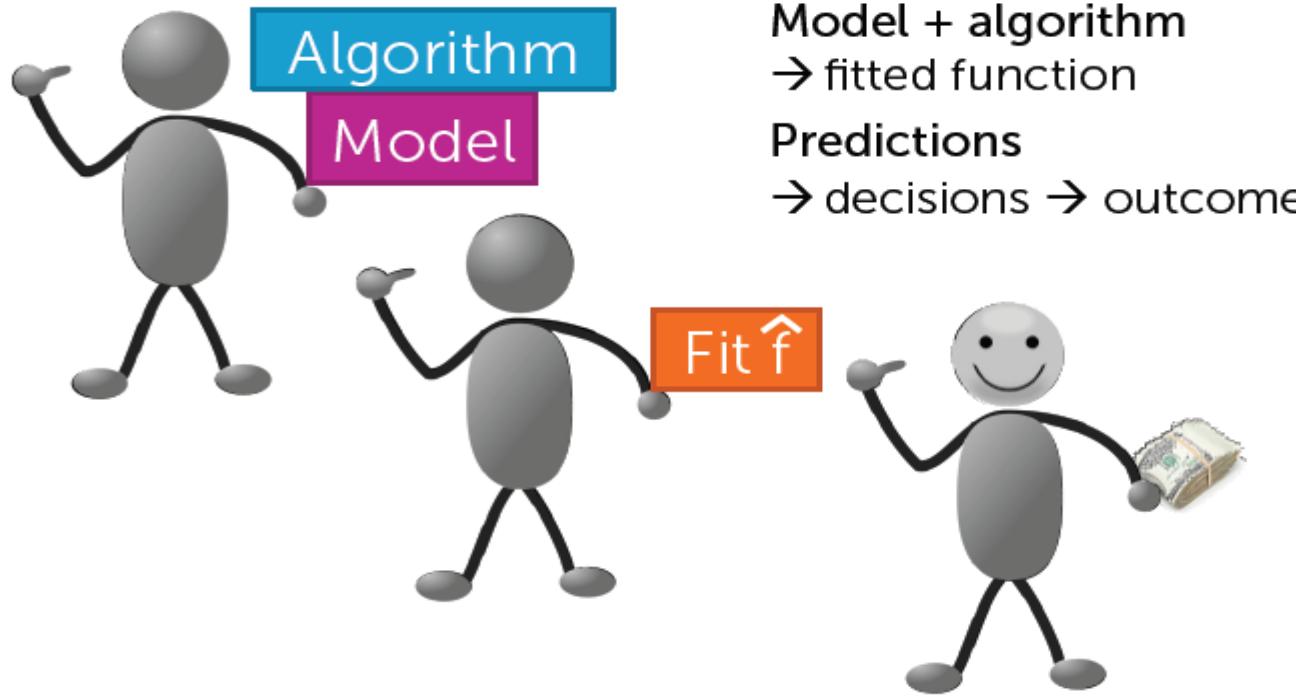
- Describe polynomial regression
- Detrend a time series using trend and seasonal components
- Write a regression model using multiple inputs or features thereof
- Cast both polynomial regression and regression with multiple inputs as regression with multiple features
- Calculate a goodness-of-fit metric (e.g., RSS)
- Estimate model parameters of a general multiple regression model to minimize RSS:
  - In closed form
  - Using an iterative gradient descent algorithm
- Interpret the coefficients of a non-featurized multiple regression fit
- Exploit the estimated model to form predictions
- Explain applications of multiple regression beyond house price modeling

# ACCESSING PERFORMANCE

# Assessing performance

72

Make predictions, get \$, right??



Model + algorithm  
→ fitted function  
**Predictions**  
→ decisions → outcome

# Assessing performance

73

## 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 = ???

# Measuring loss

74

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

Loss function:

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

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

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

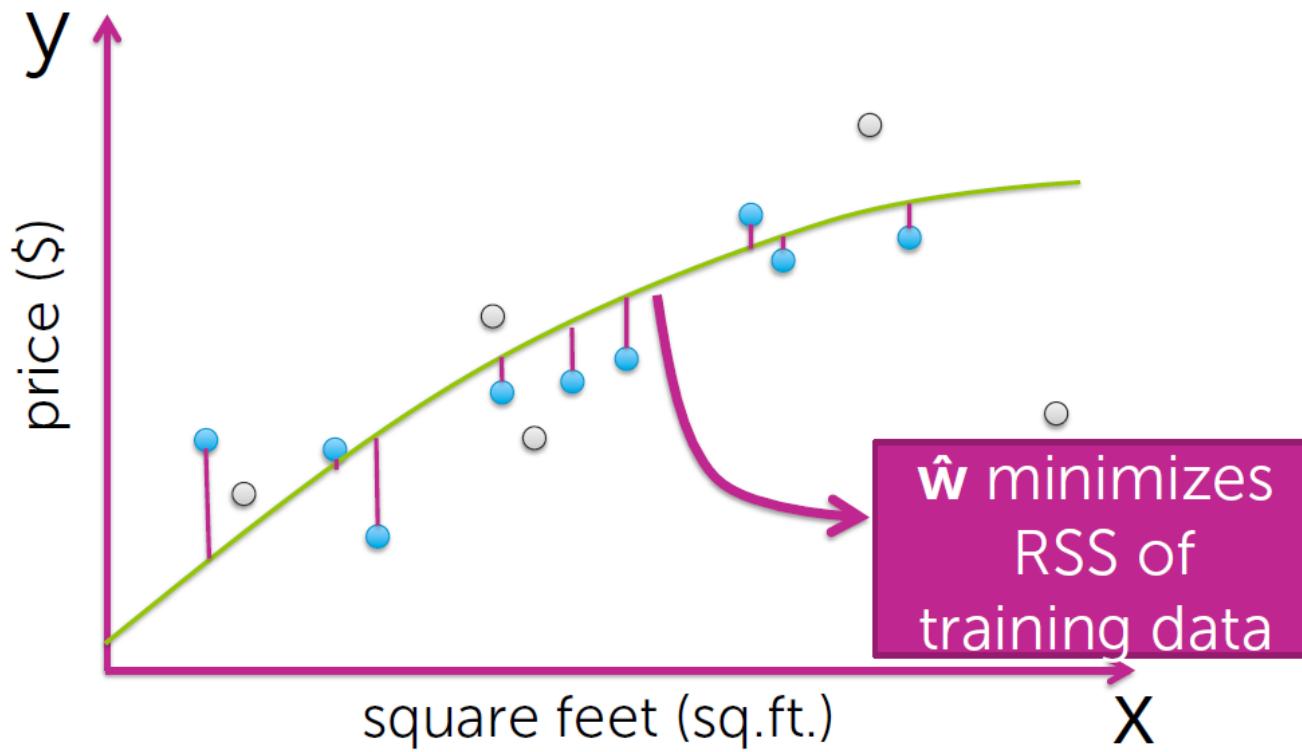
**Symmetric loss  
functions**



# Accessing the loss

75

## Use training data



# Compute training error

76

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

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

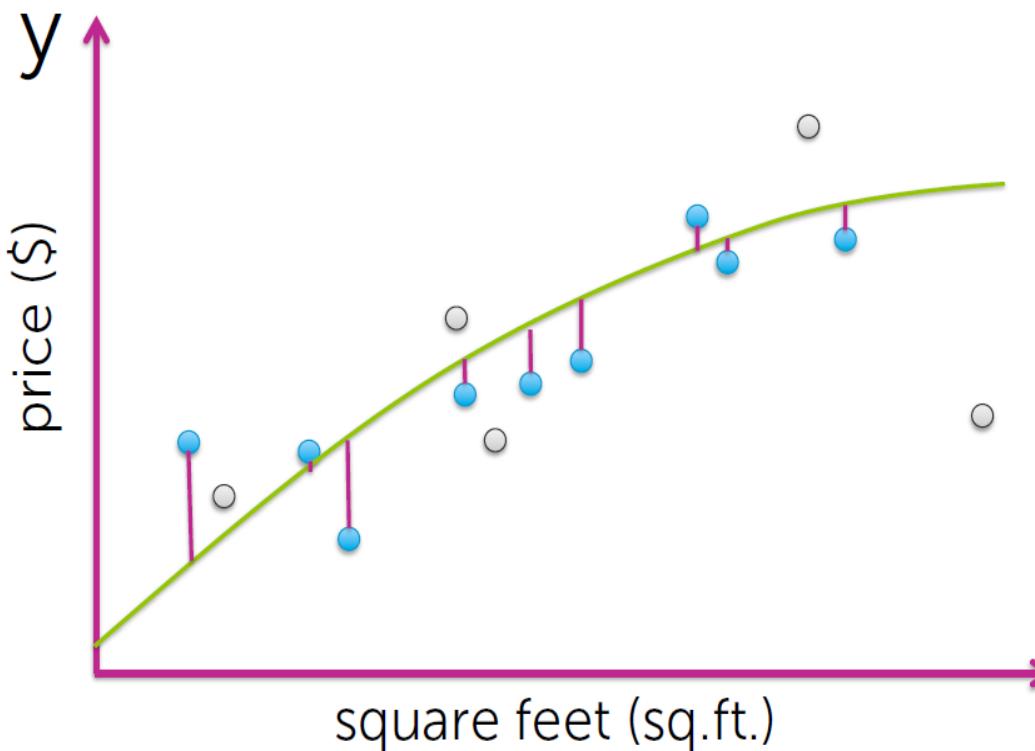


fit using training data

# Training error

77

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



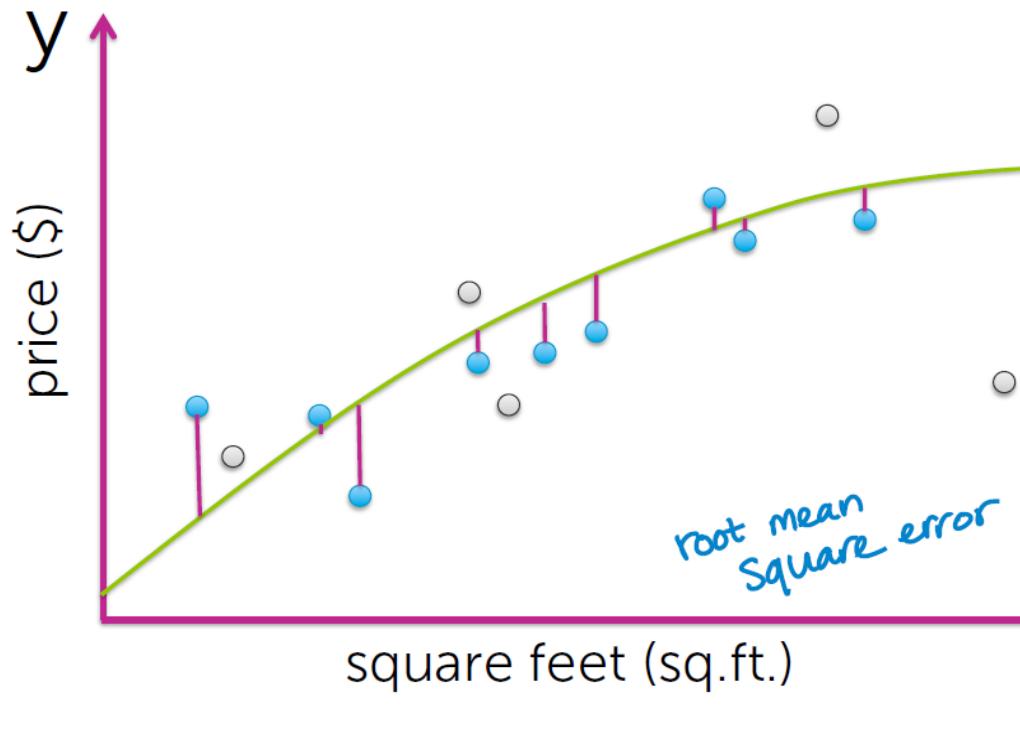
Convention is to take  
average here

Training error ( $\hat{w}$ ) =  $1/N * [ (\$_{train\ 1} - f_{\hat{w}}(\text{sq.ft.}_{train\ 1}))^2 + (\$_{train\ 2} - f_{\hat{w}}(\text{sq.ft.}_{train\ 2}))^2 + (\$_{train\ 3} - f_{\hat{w}}(\text{sq.ft.}_{train\ 3}))^2 + \dots \text{ include all training houses}]$

# Training error

78

**More intuitive is to take RMSE, same units as  $y$**

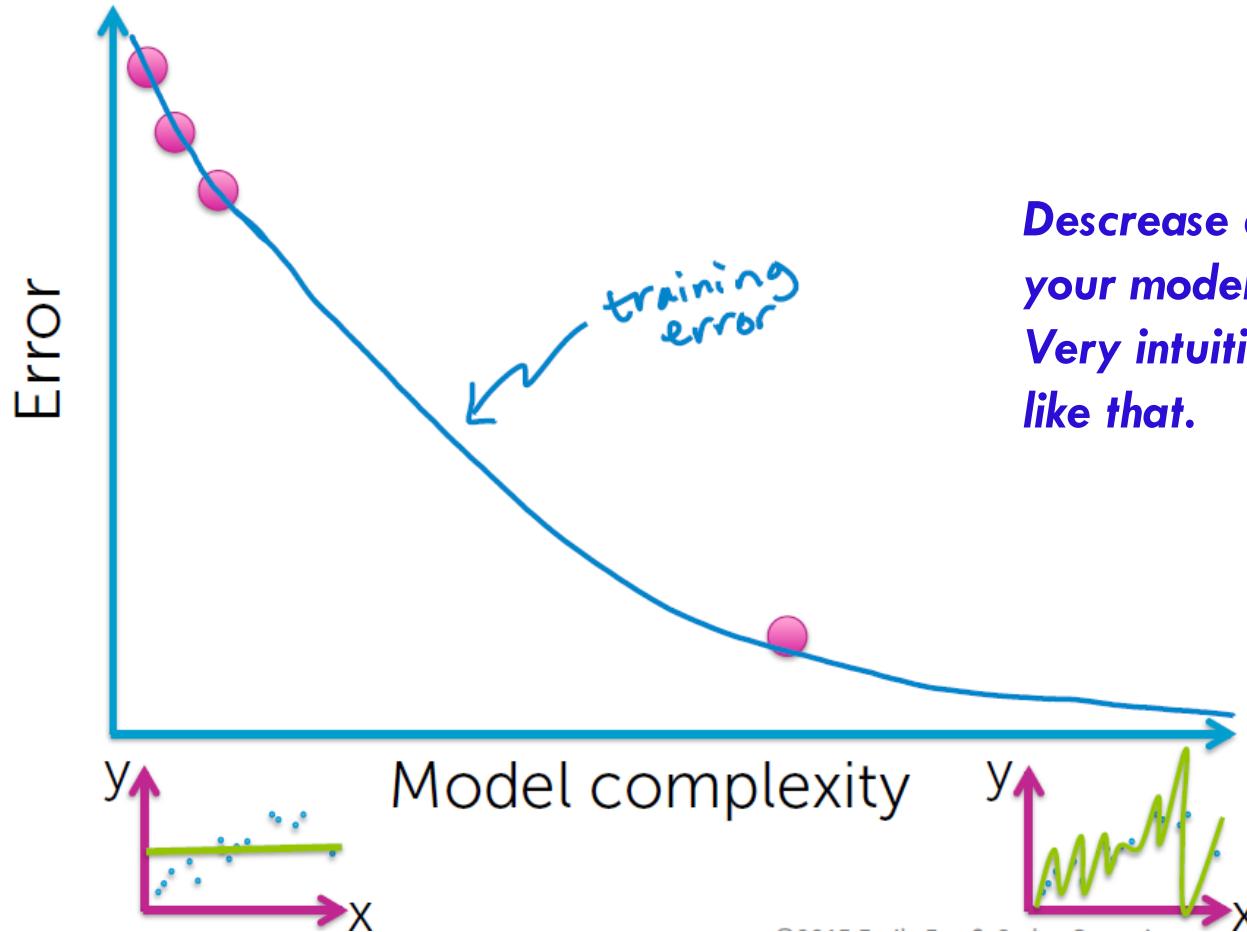


$$\text{Training error } (\hat{\mathbf{w}}) = \frac{1}{N} \sum_{i=1}^N (y_i - f_{\hat{\mathbf{w}}}(\mathbf{x}_i))^2$$

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

# Training error vs. model complexity

79

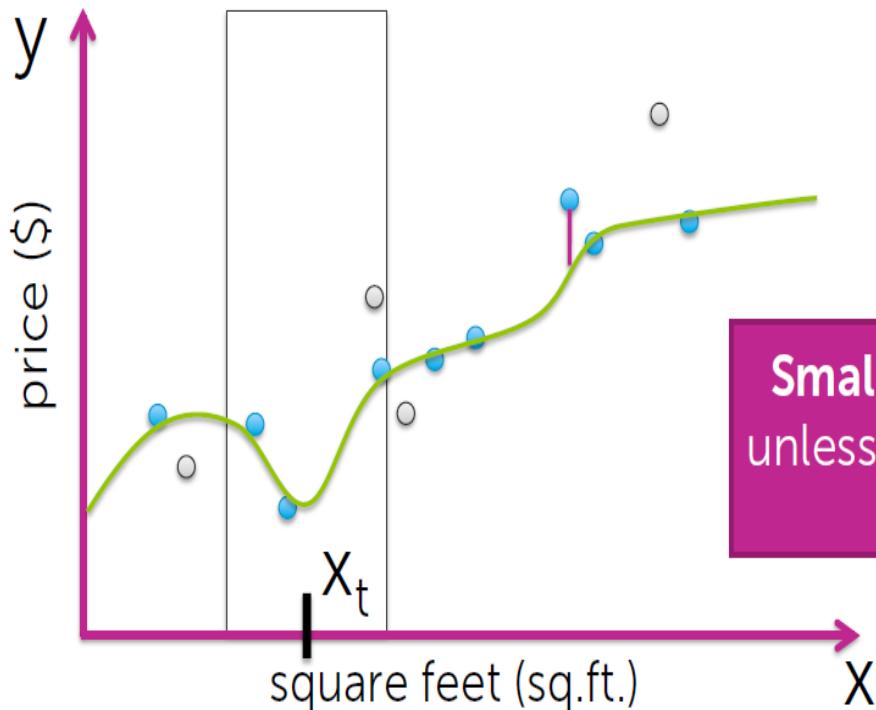


*Decrease as you increase  
your model complexity.  
Very intuitive why it is  
like that.*

# Is training error a good measure?

80

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



*Is there something particularly wrong  
about having  $x_t$ , square feet ???*

**Small training error  $\nRightarrow$  good predictions**  
unless training data includes everything you  
might ever see

# Generalisation (true) error

81

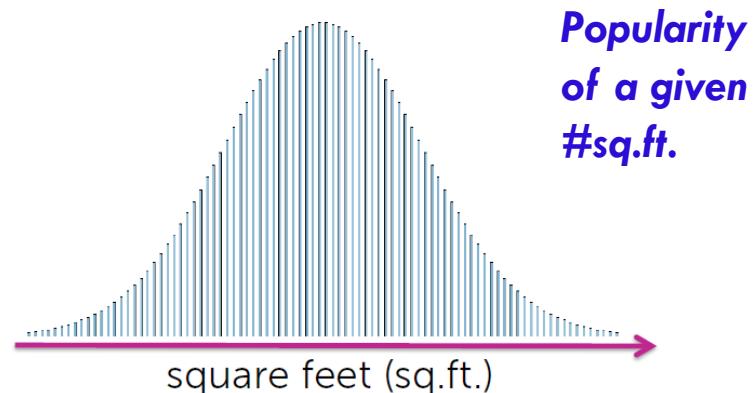
Really want estimate of loss  
over all possible (House, \$) pairs



# Distribution over house

82

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



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



# Generalisation error definition

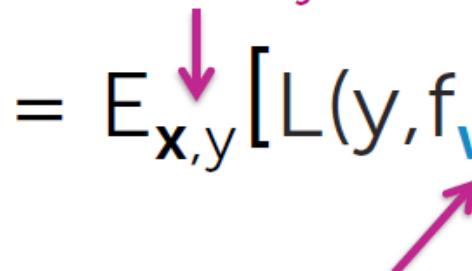
83

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

Formally:

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

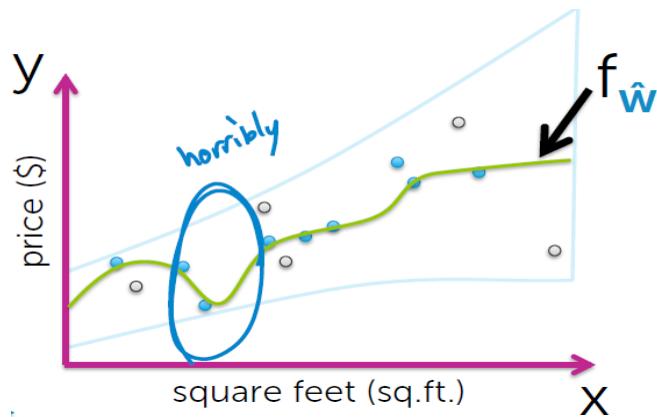
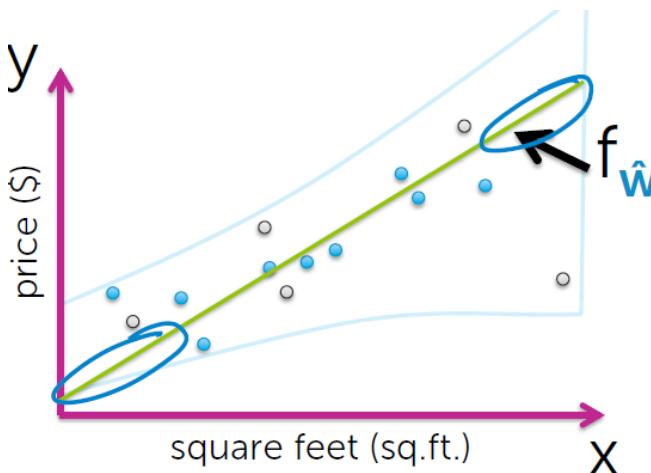
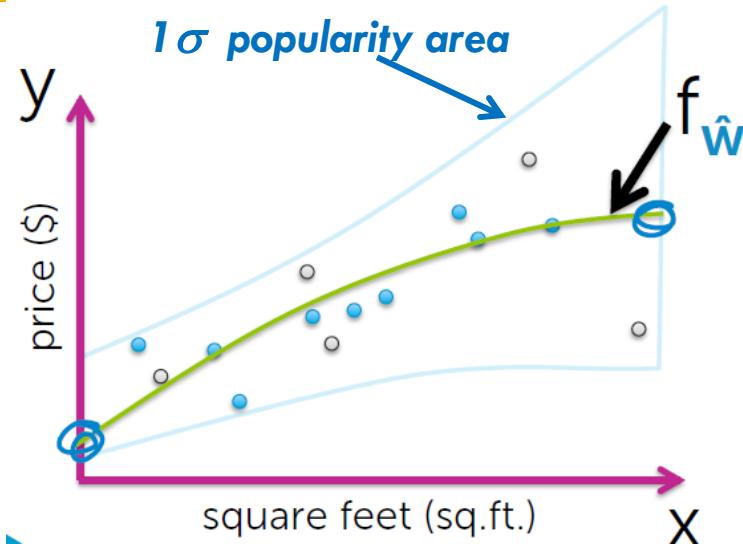
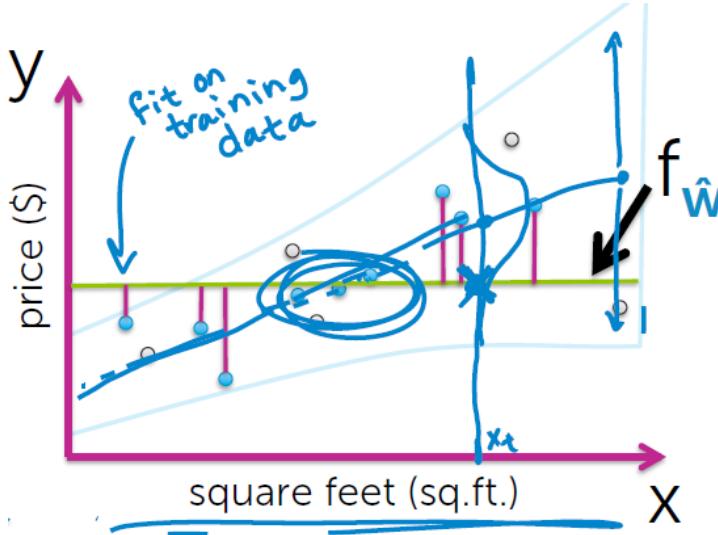
$$\text{generalization error} = E_{\mathbf{x},y} \left[ L(y, f_{\hat{\mathbf{w}}}(\mathbf{x})) \right]$$



fit using training data

# Generalisation error (weighted with popularity) vs model complexity

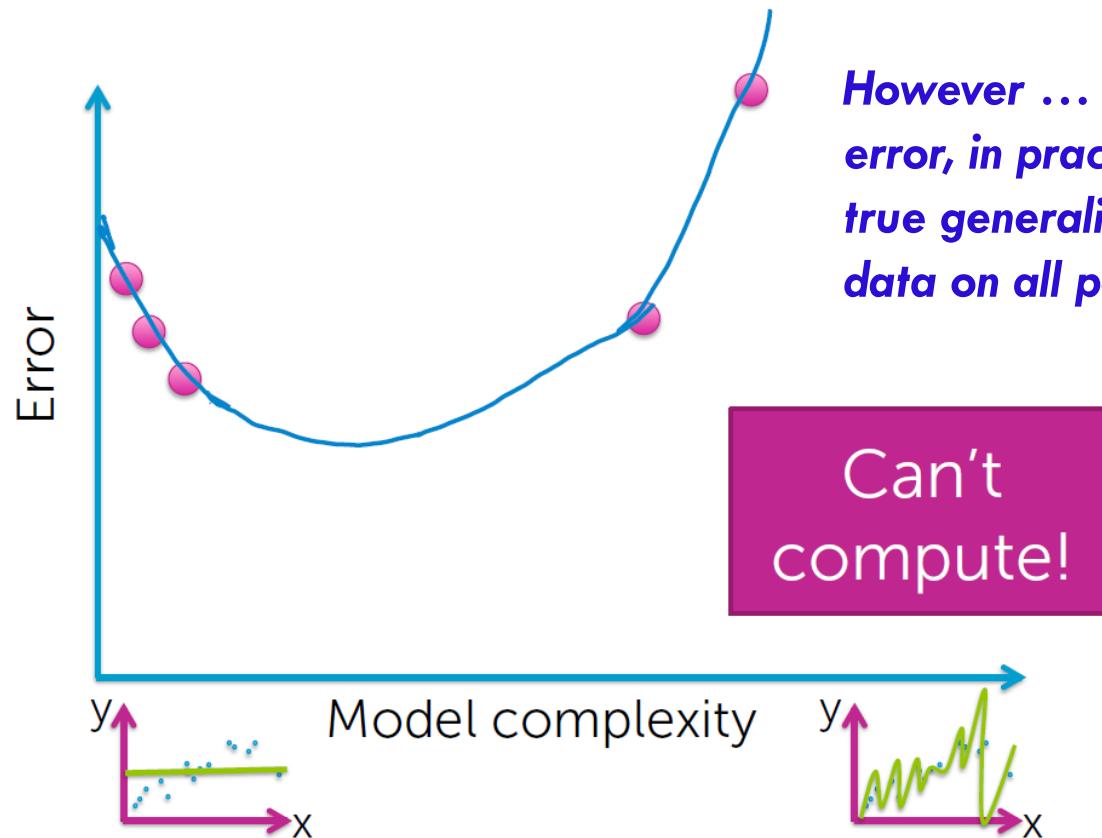
84



26/10, 2/11, 9/11 2022

# Generalisation error vs model complexity

85



# Forming a test set

86

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



We want to approximate generalisation error.

*Test set: proxy for „everything you might see“*

Training set



Test set



# Compute test error

87

## Test error

= avg. loss on houses in test set

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

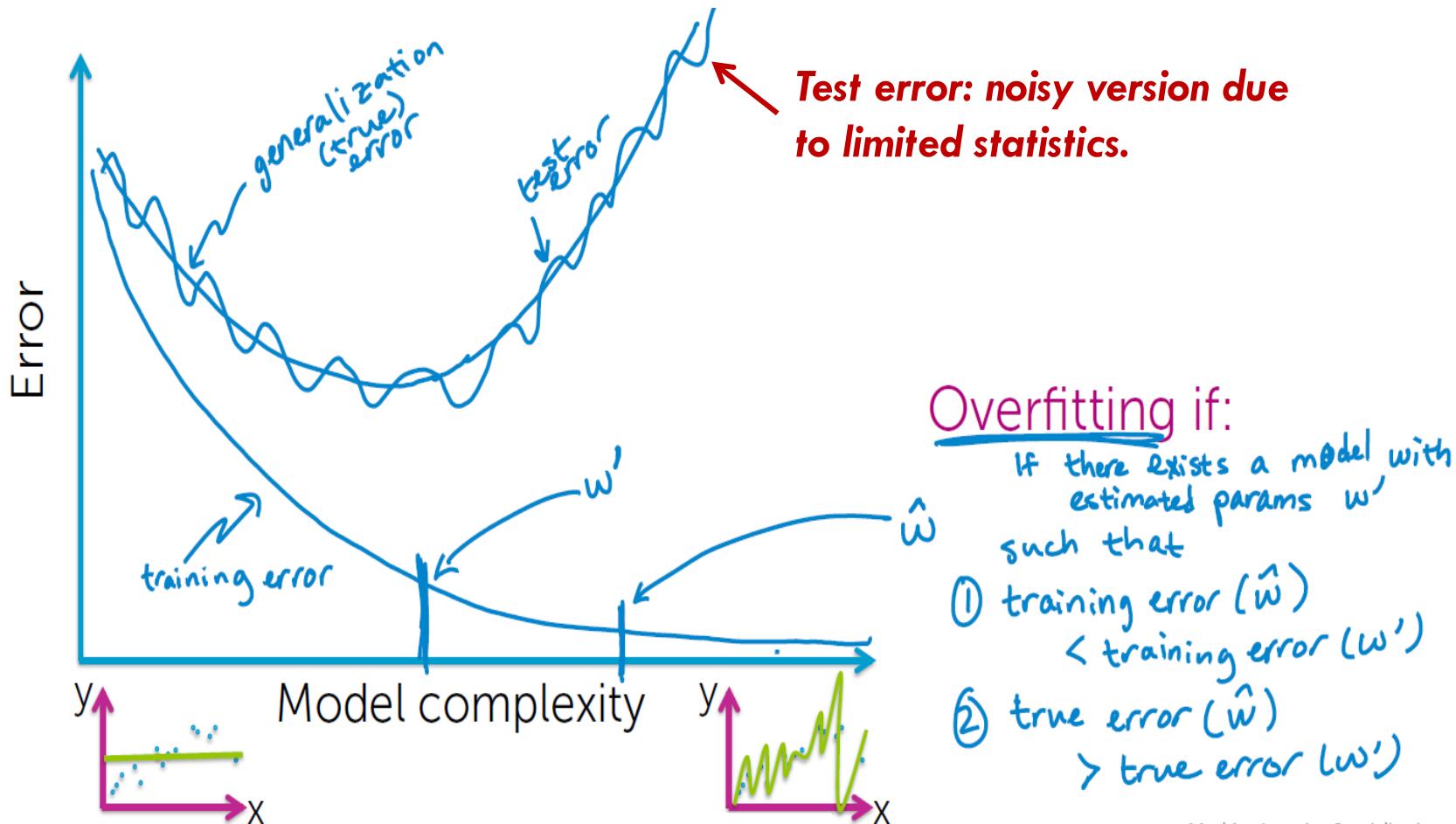
↑  
# test points

fit using training data

**has never seen  
test data!**

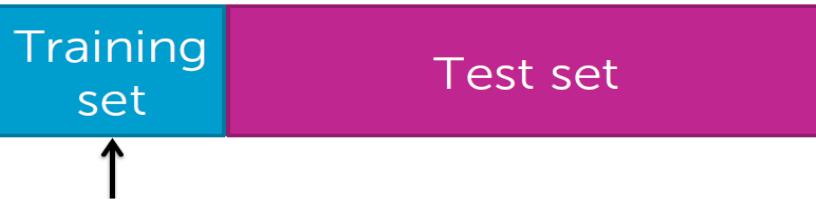
# Training, true and test error vs. model complexity. Notion of overfitting.

88

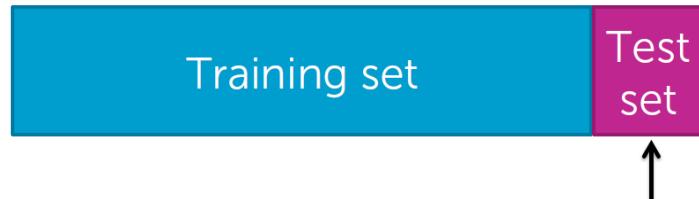


# Training/test splits

89



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



Too few → test error bad approximation of generalization error



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** (will see later...)

# Three sources of errors

90

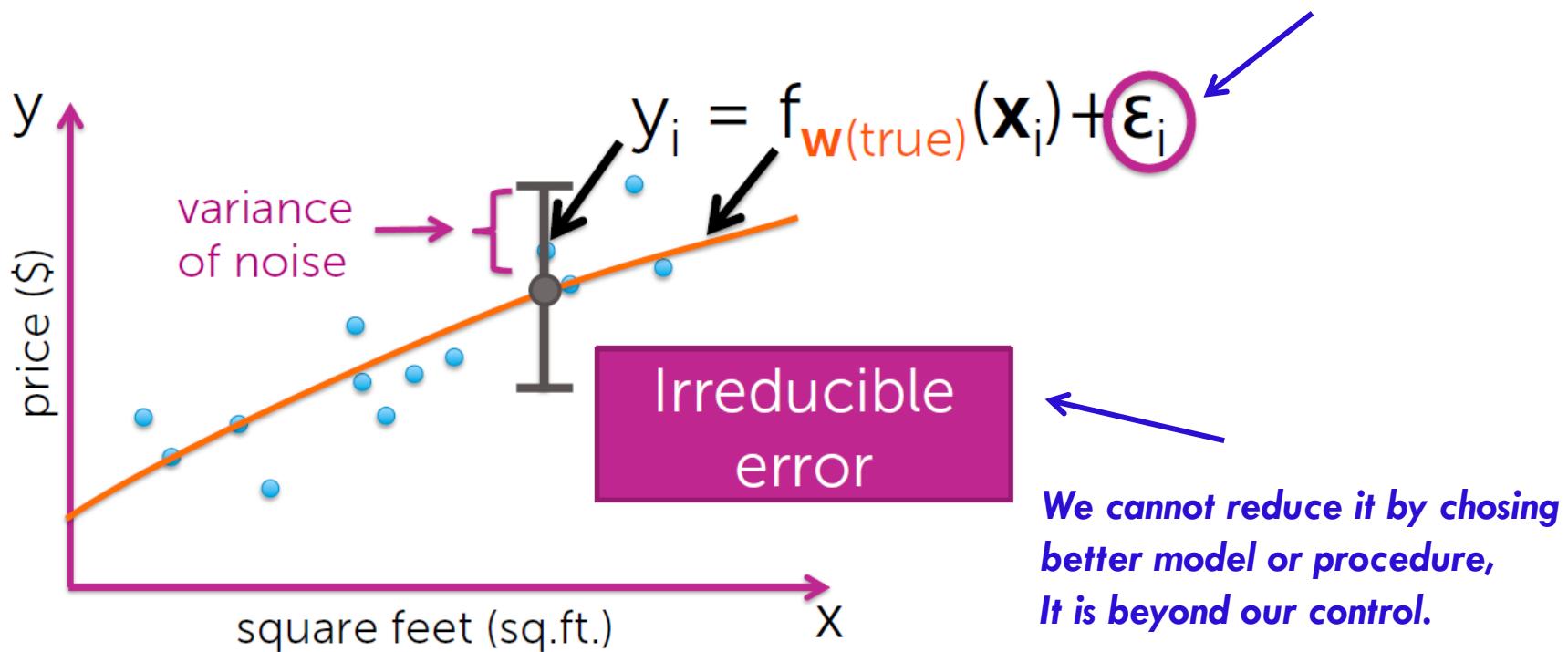
In forming predictions, there are 3 sources of error:

1. Noise
2. Bias
3. Variance

# Data are inherently noisy

91

*There is some true relationship between sq.ft and value of the house, specific to the given house.*

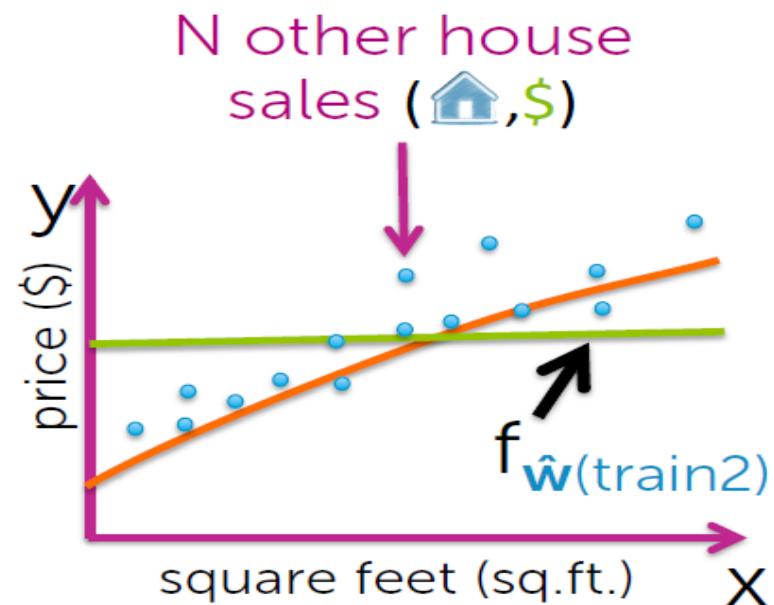
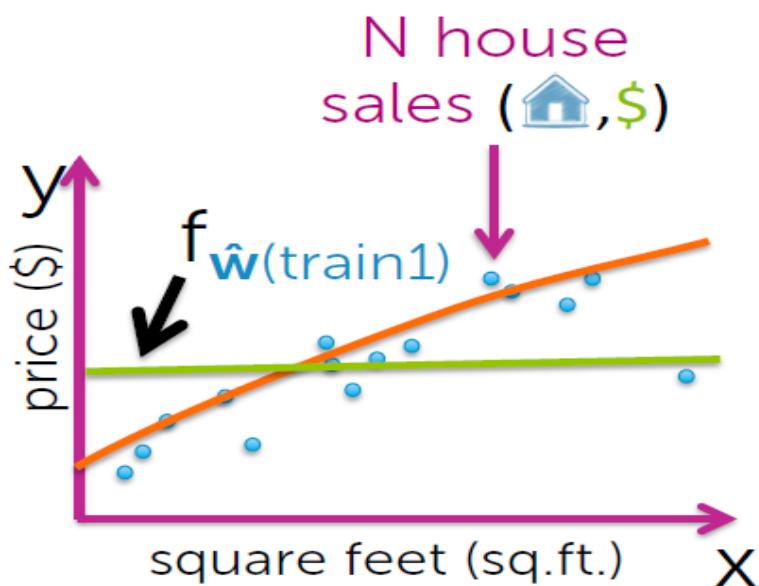


# Bias contribution

92

**This contribution we can control.**

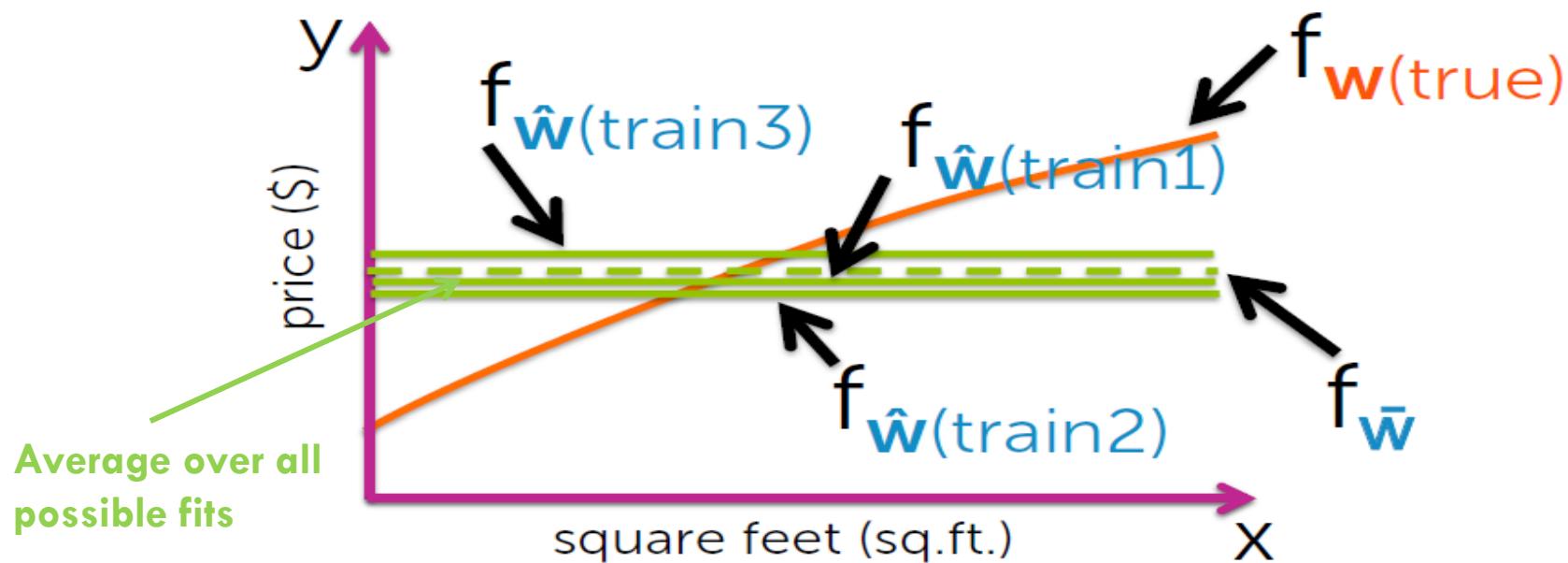
Assume we fit a constant function



# Bias contribution

93

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

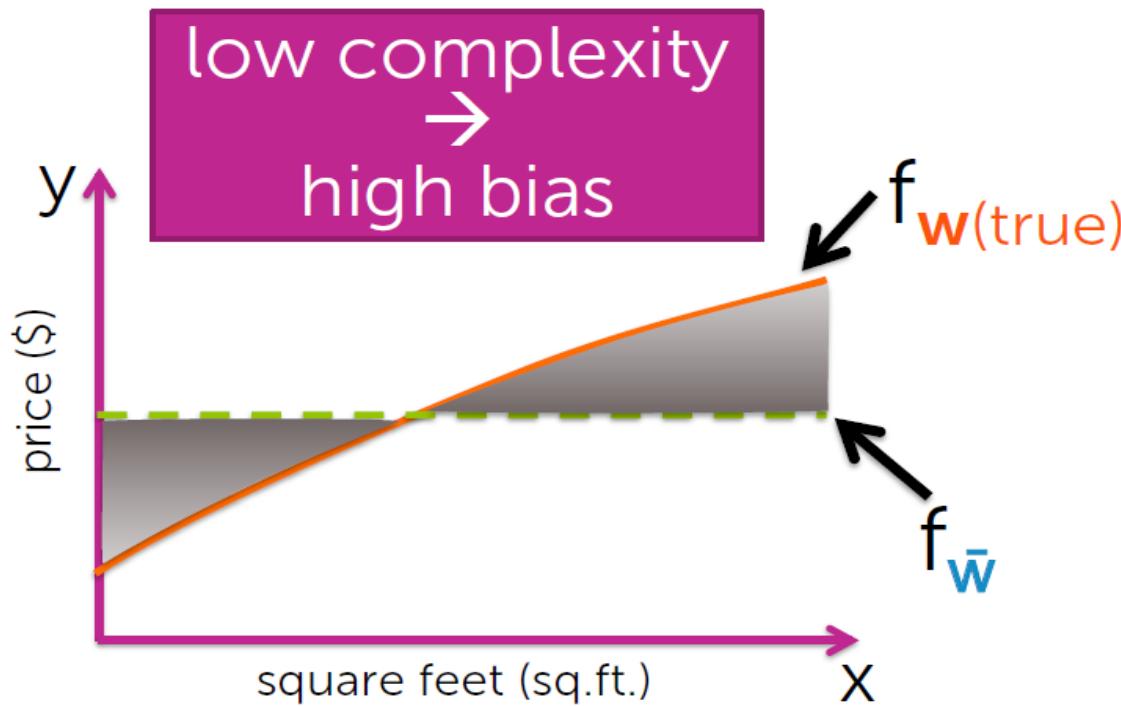


# Bias contribution

94

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

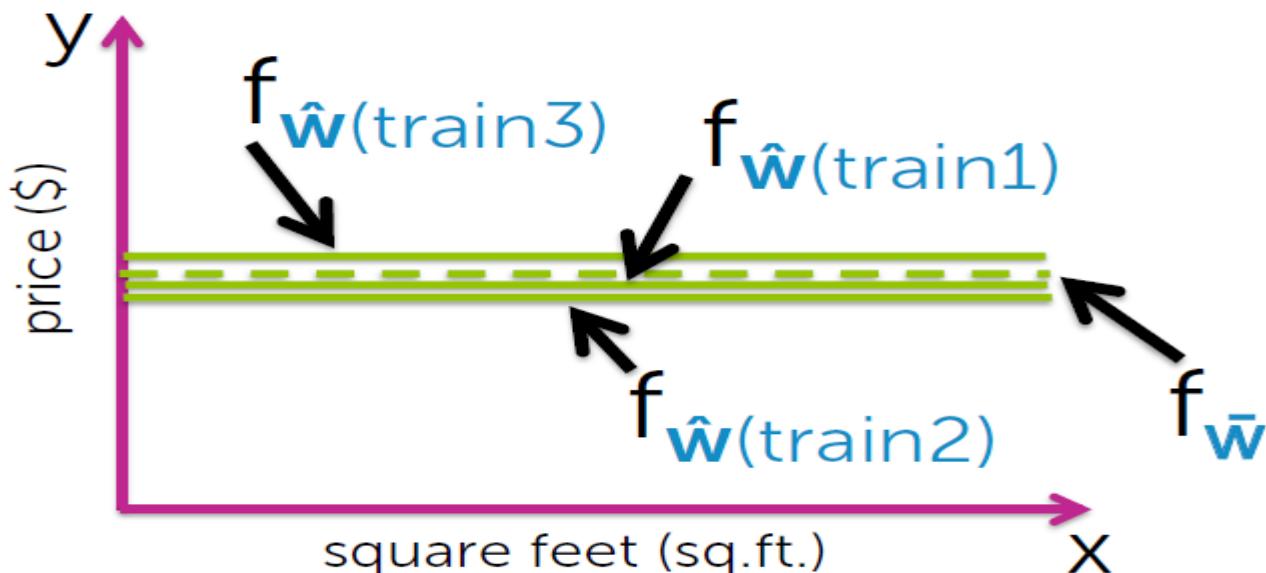
Is our approach flexible enough to capture  $f_{\mathbf{w}(\text{true})}$ ?  
If not, error in predictions.



# Variance contribution

95

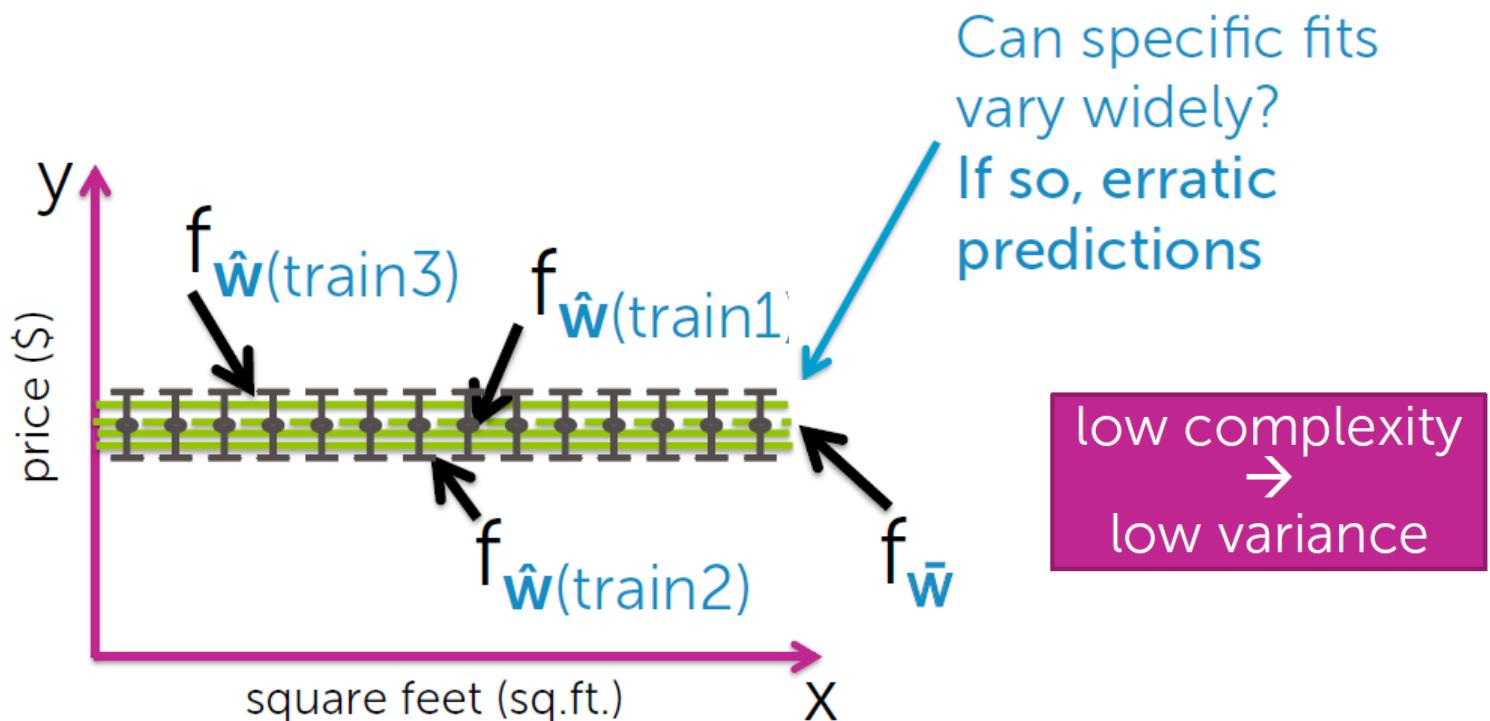
How much do specific fits vary from the expected fit?



# Variance contribution

96

How much do specific fits vary from the expected fit?

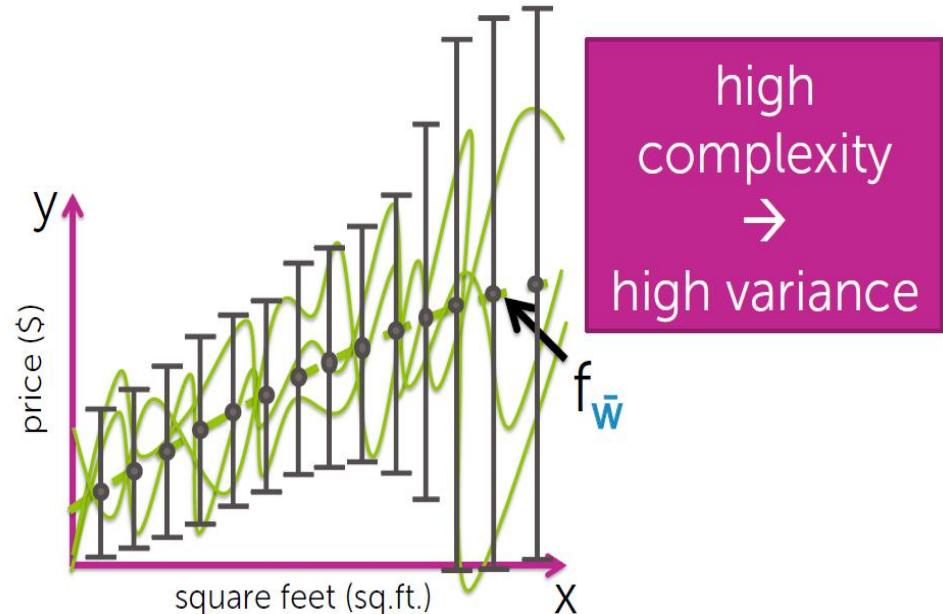
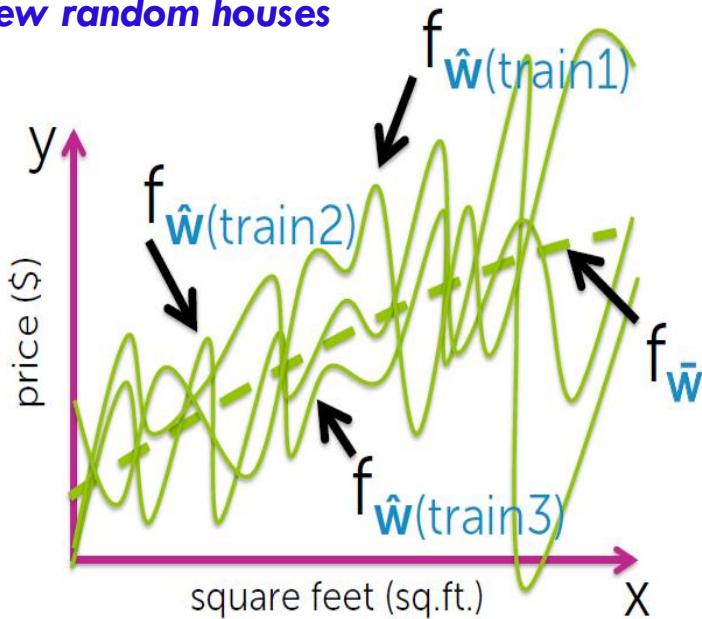


# Variance of high complexity models

97

Assume we fit a high-order polynomial

*For each train remove  
few random houses*

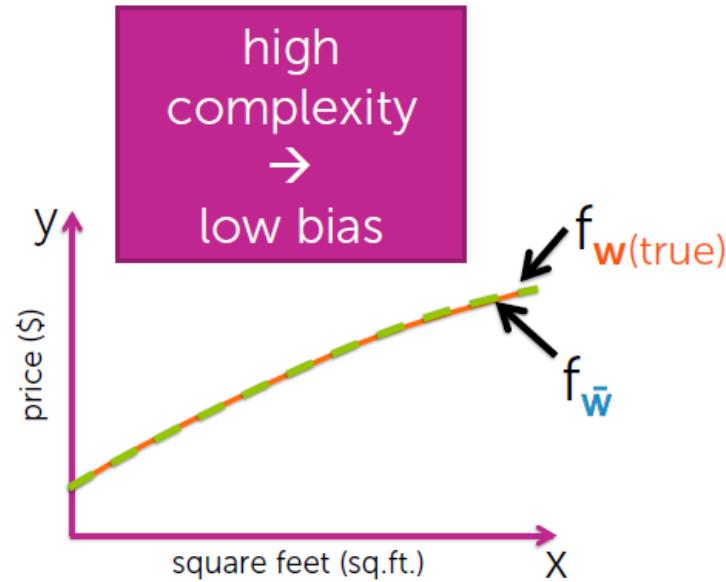
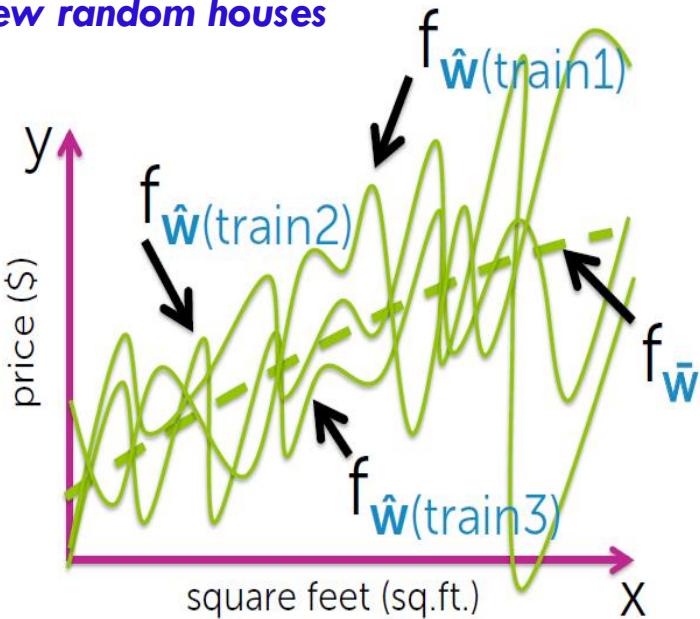


# Bias of high complexity models

98

Assume we fit a high-order polynomial

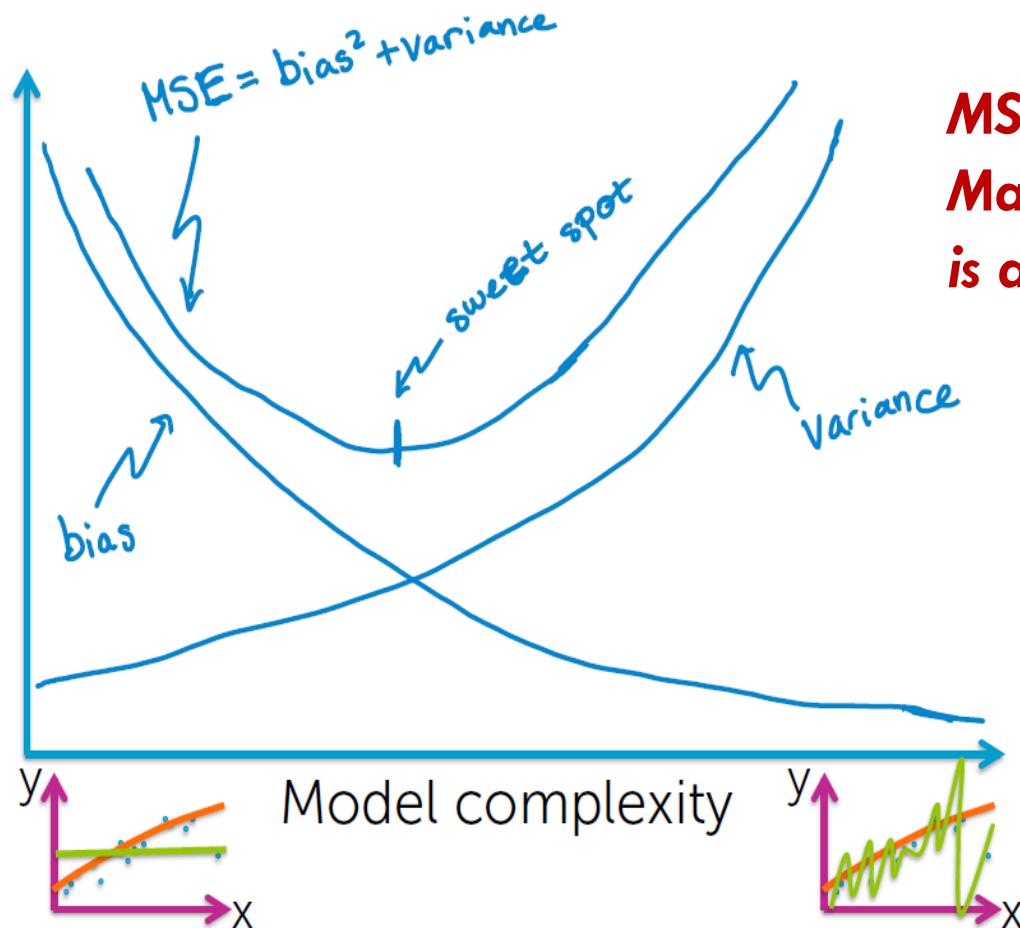
**For each train remove  
few random houses**



**High complexity models are very flexible,  
pick better average trends.**

# Bias –variance tradeoff

99



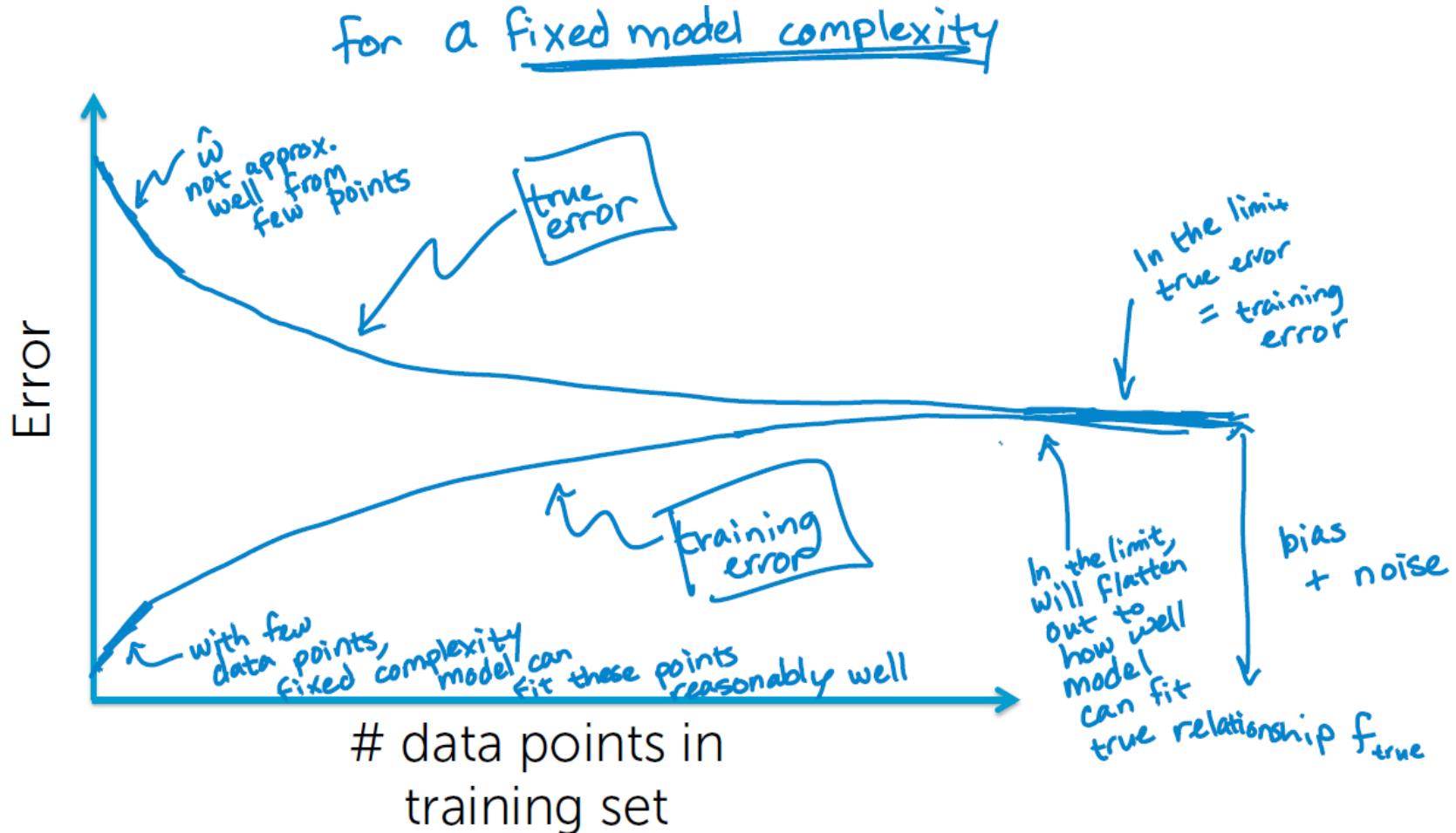
**MSE = mean square error**  
**Machine Learning**  
**is all about this tradeoff**

**But....**

Just like with  
generalization error,  
we cannot compute  
bias and variance

# Errors vs amount of data

100



# The regression/ML workflow

101

## 1. Model selection

Often, need to choose tuning parameters  $\lambda$  controlling model complexity (e.g. degree of polynomial)

## 2. Model assessment

Having selected a model, assess the generalization error

# Hypothetical implementation

102



## 1. Model selection

For each considered model complexity  $\lambda$  :

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

## 2. Model assessment

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

# Hypothetical implementation

103

Training set

Test set

## 1. Model selection

For each considered model complexity  $\lambda$  :

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

## 2. Model assessment

Overly optimistic!

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

# Hypothetical implementation

104



**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

# Practical implementation

105

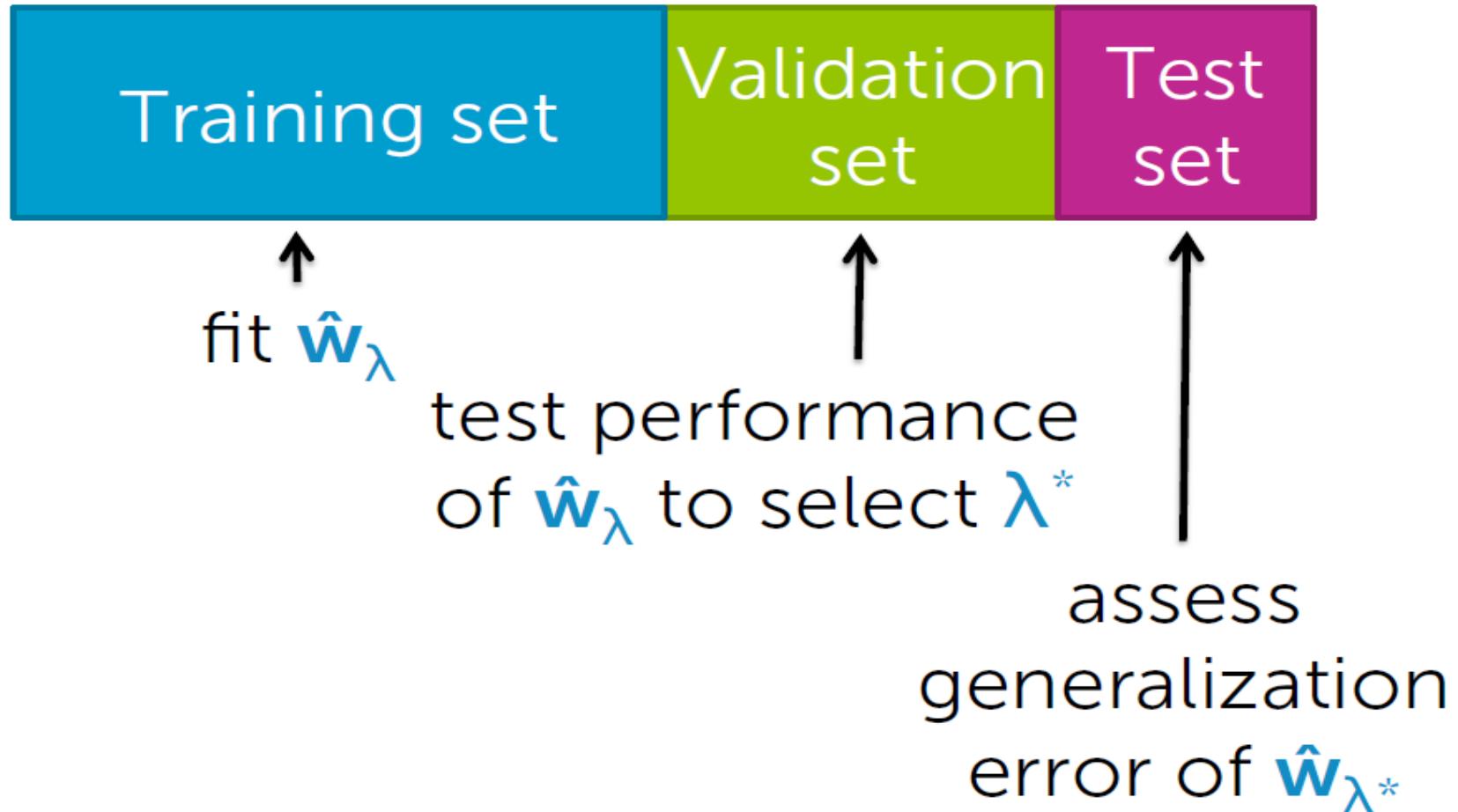


**Solution:** Create two “test” sets!

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

# Practical implementation

106



# Typical splits

107



# What you can do now

108

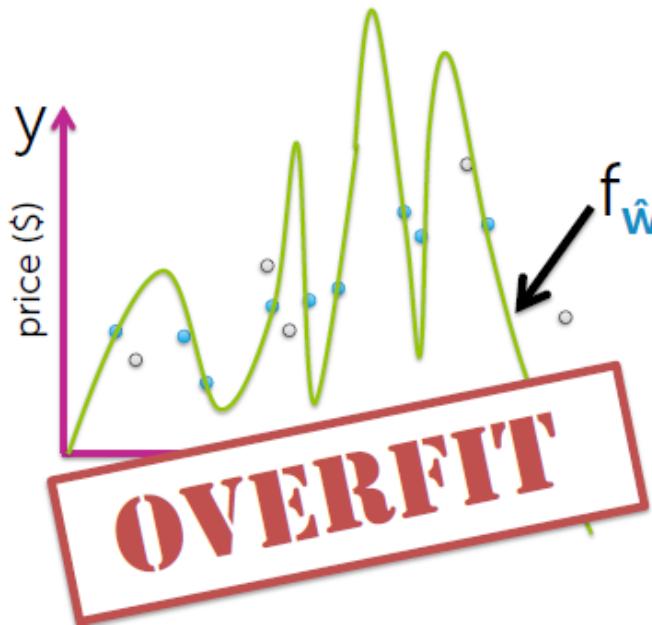
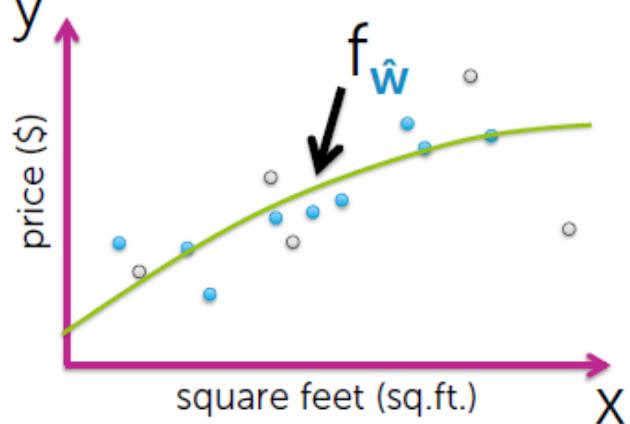
- Describe what a loss function is and give examples
- Contrast training, generalization, and test error
- Compute training and test error given a loss function
- Discuss issue of assessing performance on training set
- Describe tradeoffs in forming training/test splits
- List and interpret the 3 sources of avg. prediction error
  - Irreducible error, bias, and variance
- Discuss issue of selecting model complexity on test data and then using test error to assess generalization error
- Motivate use of a validation set for selecting tuning parameters (e.g., model complexity)
- Describe overall regression workflow

# RIDGE REGRESSION

# Flexibility of high-order polynomials

110

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



**Symptoms for overfitting: often associated with very large value of estimated parameters  $\hat{w}$**

# Overfitting with many features

111

Not unique to polynomial regression,  
but also if **lots of inputs (d large)**

Or, generically,

**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
- ...

# How does # of observations influence overfitting?

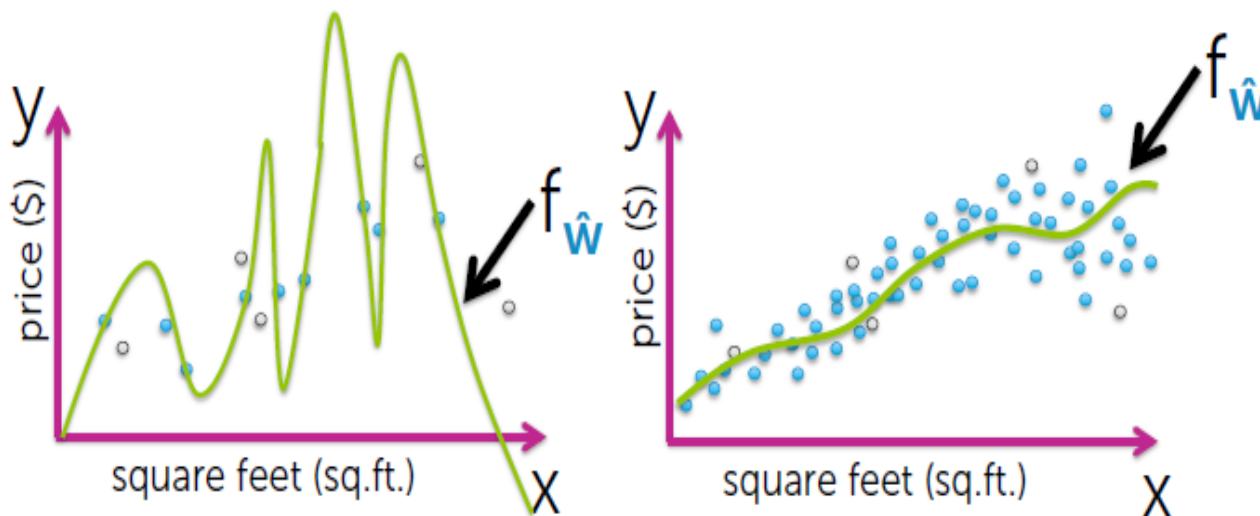
112

Few observations ( $N$  small)

→ rapidly overfit as model complexity increases

Many observations ( $N$  very large)

→ harder to overfit



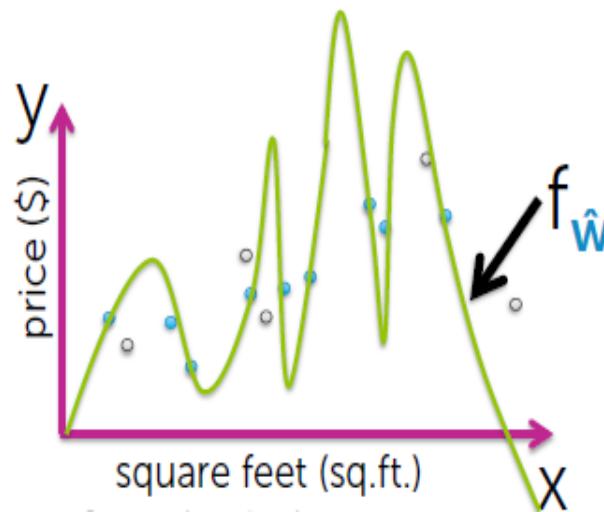
# How does # of inputs influence overfitting?

113

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

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

HARD



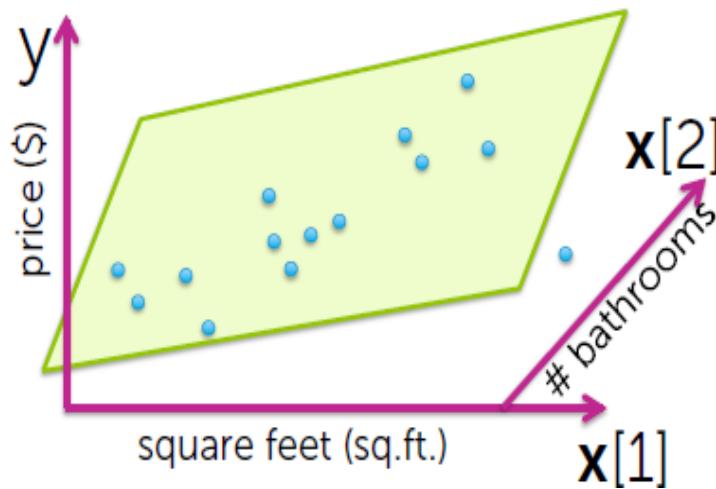
# How does # of inputs influence overfitting?

114

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

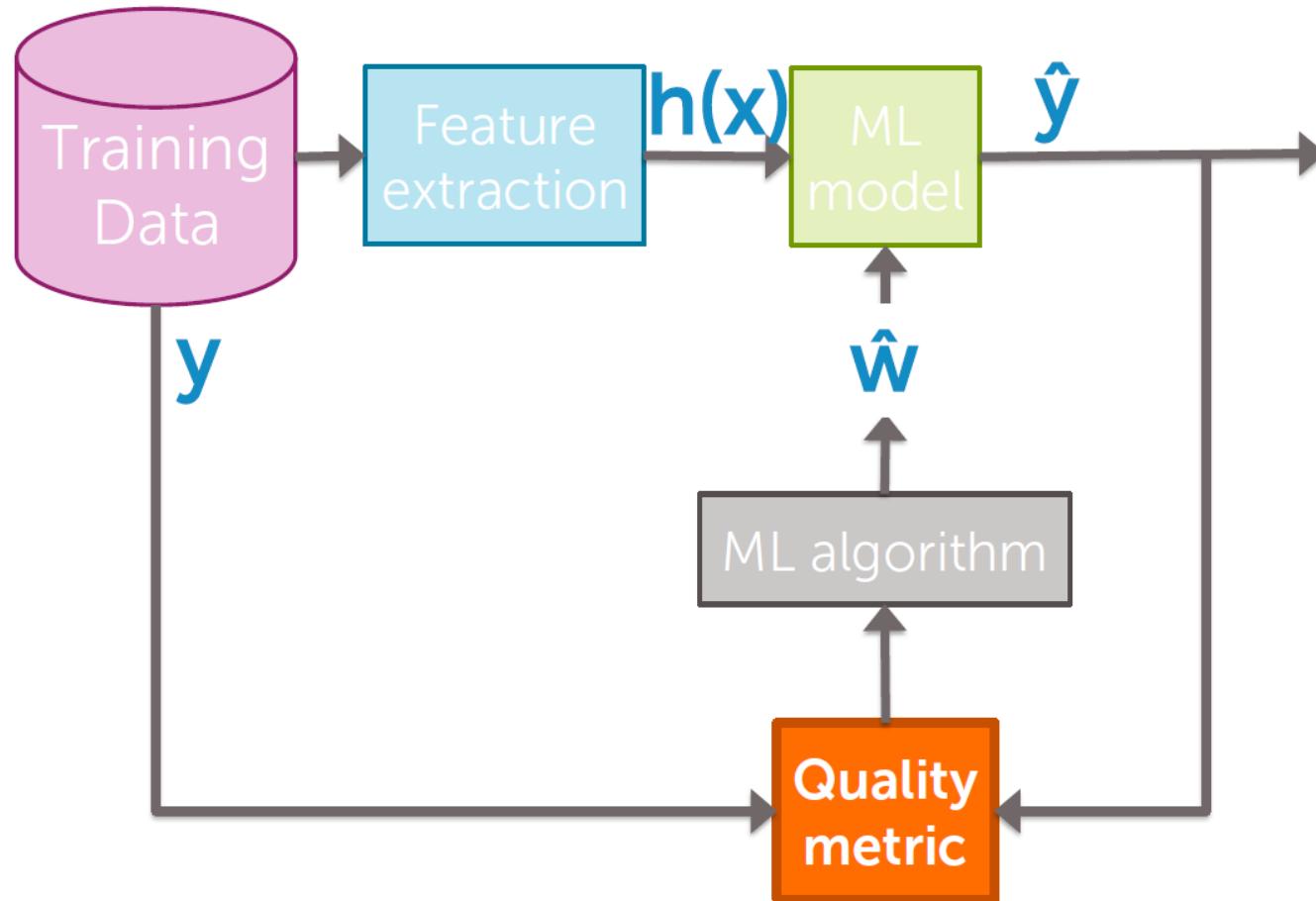
Data must include examples of all possible (sq.ft., #bath, #bed, lot size, year,..., \$) combos to avoid overfitting

MUCH!!!  
HARDER



# Lets improve quality metric blok

115



# Desire total cost format

116

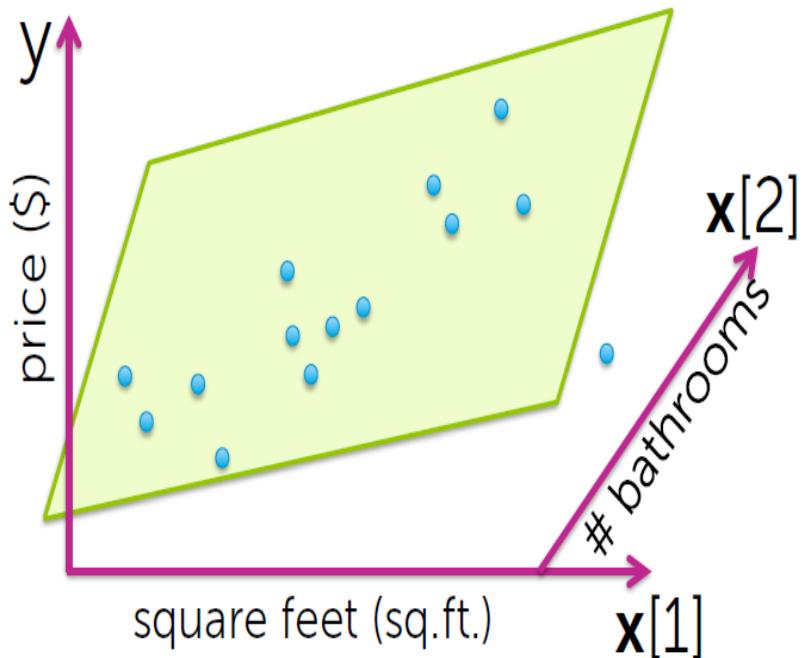
Want to balance:

- i. How well function fits data
- ii. Magnitude of coefficients



# Measure of fit to training data

117



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

↑ pred. value using  $\mathbf{w}$

small RSS → model fitting training data well

# Measure of magnitude of regression coefficients

118

What summary # is indicative of size of regression coefficients?

- Sum?  $w_0 = 1,527,301 \quad w_1 = -1,605,253$

$$w_0 + w_1 = \text{small } \#$$

*But ... the coefficients are very large*

- Sum of absolute value?

$$|w_0| + |w_1| + \dots + |w_D| = \sum_{j=0}^D |w_j| \triangleq \|w\|_1, \quad L_1 \text{ norm} \quad \dots \text{discuss more in next module}$$

- Sum of squares ( $L_2$  norm)

$$w_0^2 + w_1^2 + \dots + w_D^2 = \sum_{j=0}^D w_j^2 \triangleq \|w\|_2^2 \quad L_2 \text{ norm} \quad \dots \boxed{\text{focus of this module}}$$

# Consider specific total cost

119

Total cost =

measure of fit + measure of magnitude

of coefficients

$\text{RSS}(\mathbf{w})$

$\|\mathbf{w}\|_2^2$

# Consider resulting objectives

120

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

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

↑ tuning parameter = balance of fit and magnitude

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

If  $\lambda=0$ :

reduces to minimizing  $\text{RSS}(w)$ , as before (old solution)  $\rightarrow \hat{w}^{\text{LS}}$  ← least squares

If  $\lambda=\infty$ :

for solutions where  $\hat{w} \neq 0$ , then total cost is  $\infty$

If  $\hat{w}=0$ , then total cost =  $\text{RSS}(0)$   $\rightarrow$  solution is  $\hat{w}=0$

If  $\lambda$  in between: Then  $0 \leq \|\hat{w}\|_2^2 \leq \|\hat{w}^{\text{LS}}\|_2^2$

# Ridge regression: bias-variance tradeoff

121

Large  $\lambda$ :

high bias, low variance

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

In essence,  $\lambda$   
controls model  
complexity

Small  $\lambda$ :

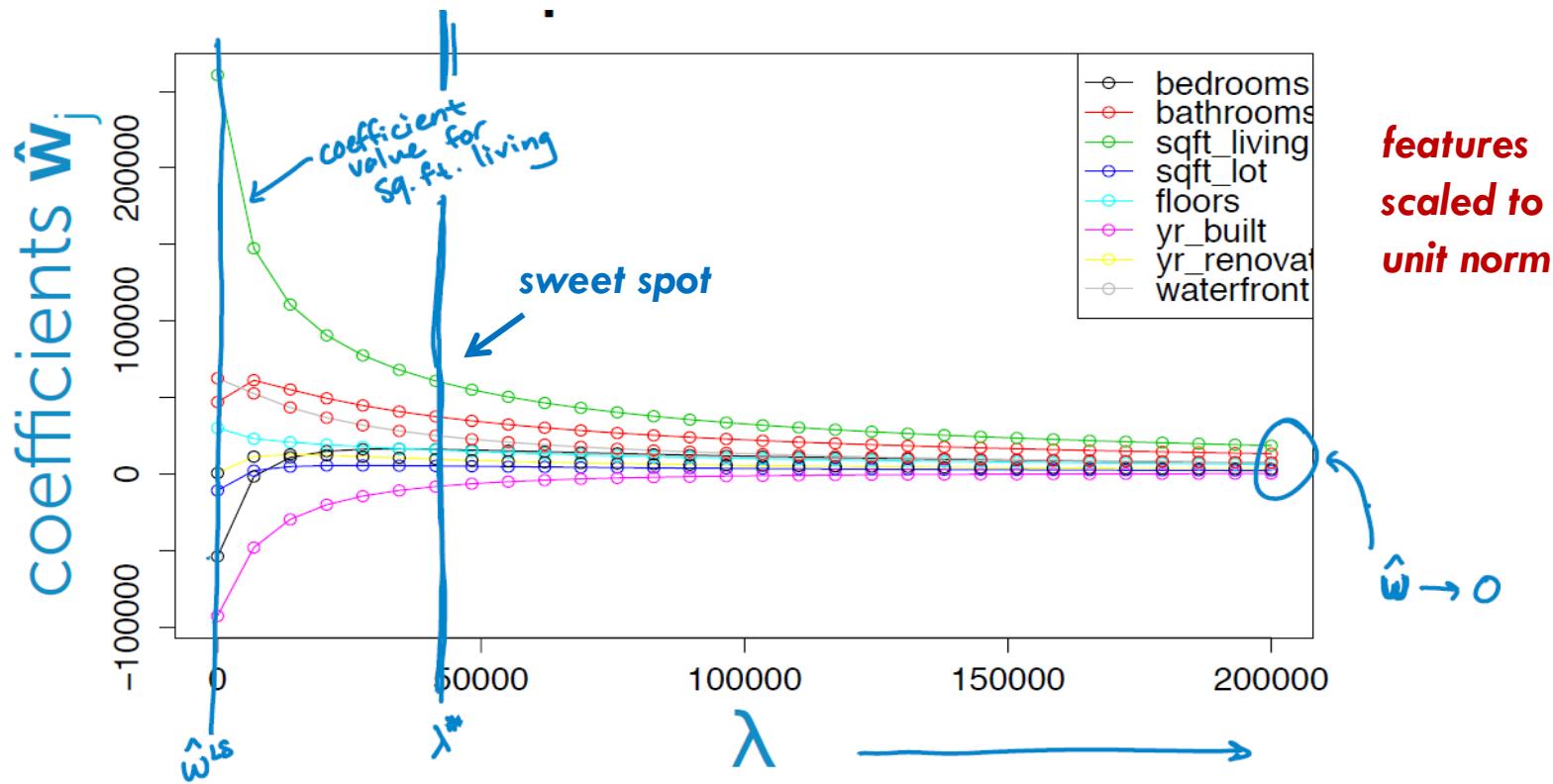
low bias, high variance

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

# Ridge regression: coefficients path

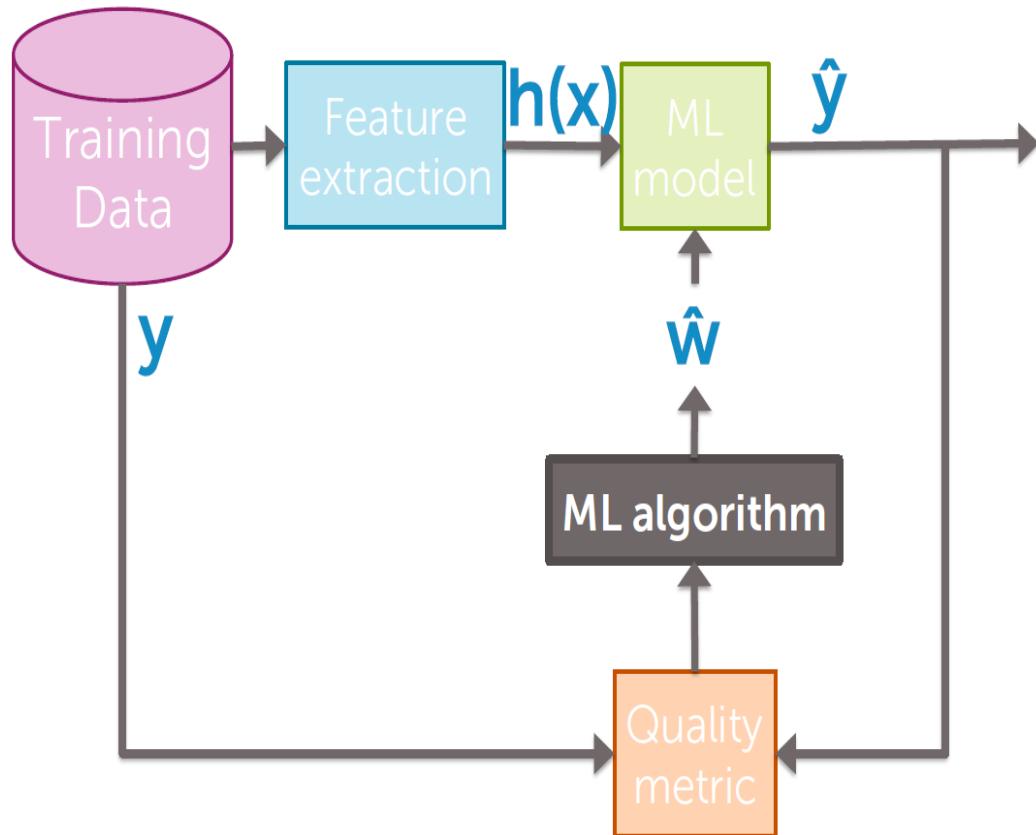
122

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



# Flow chart

123



Model for all N observations together

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

The equation shows the model structure:  $\mathbf{y}$  is a vertical vector of observed values,  $\mathbf{H}$  is a matrix of features,  $\mathbf{w}$  is a vector of weights, and  $\boldsymbol{\varepsilon}$  represents the error term.

# Ridge regression: cost in matrix notation

124

In matrix form, ridge regression cost is:

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \\ = (\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w}$$

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

$$= \begin{matrix} \text{---} \\ w_0 & w_1 & w_2 & \dots & w_D \end{matrix} \begin{matrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{matrix}$$

$$= \mathbf{w}^\top \mathbf{w}$$

# Gradient of ridge regression cost

125

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

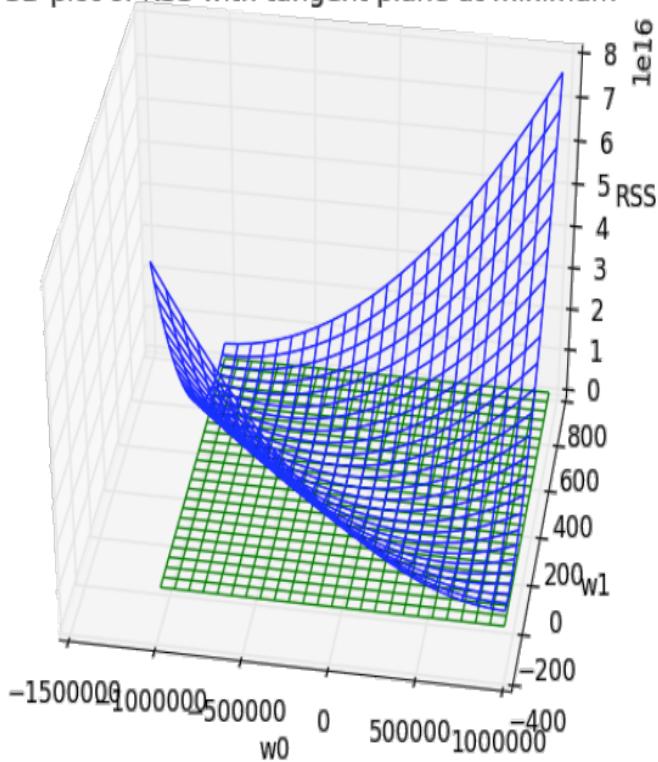
Why? By analogy to 1d case...

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

# Ridge regression: closed-form solution

126

3D plot of RSS with tangent plane at minimum



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

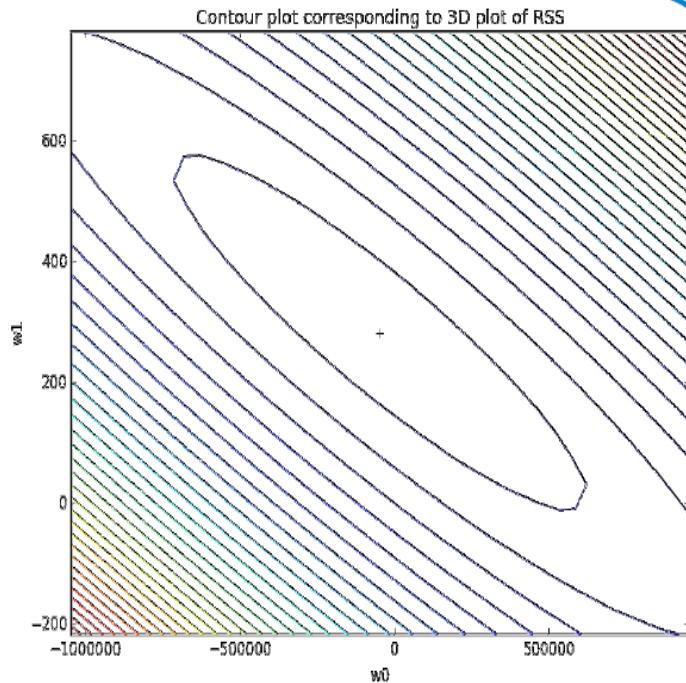
Solve for  $\hat{\mathbf{w}}$ :

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

# Ridge regression: gradient descent

127

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



Update to  $j^{\text{th}}$  feature weight:

$$w_j^{(t+1)} \leftarrow \underline{w_j^{(t)}} - \eta * \underbrace{-2 \sum_{i=1}^N (\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)}))}_{\substack{\text{Same as before} \\ (\text{from RSS term})}} + 2\lambda w_j^{(t)}$$

*new term, comes from the  $j^{\text{th}}$  component of  $2\lambda\mathbf{w}$*

# Summary of ridge regression algorithm

128

init  $\mathbf{w}^{(1)} = 0$  (or randomly, or smartly),  $t=1$

**while**  $\|\nabla \text{RSS}(\mathbf{w}^{(t)})\| > \epsilon$

**for**  $j=0, \dots, D$

$$\text{partial}[j] = -2 \sum_{i=1}^N (\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)}))$$

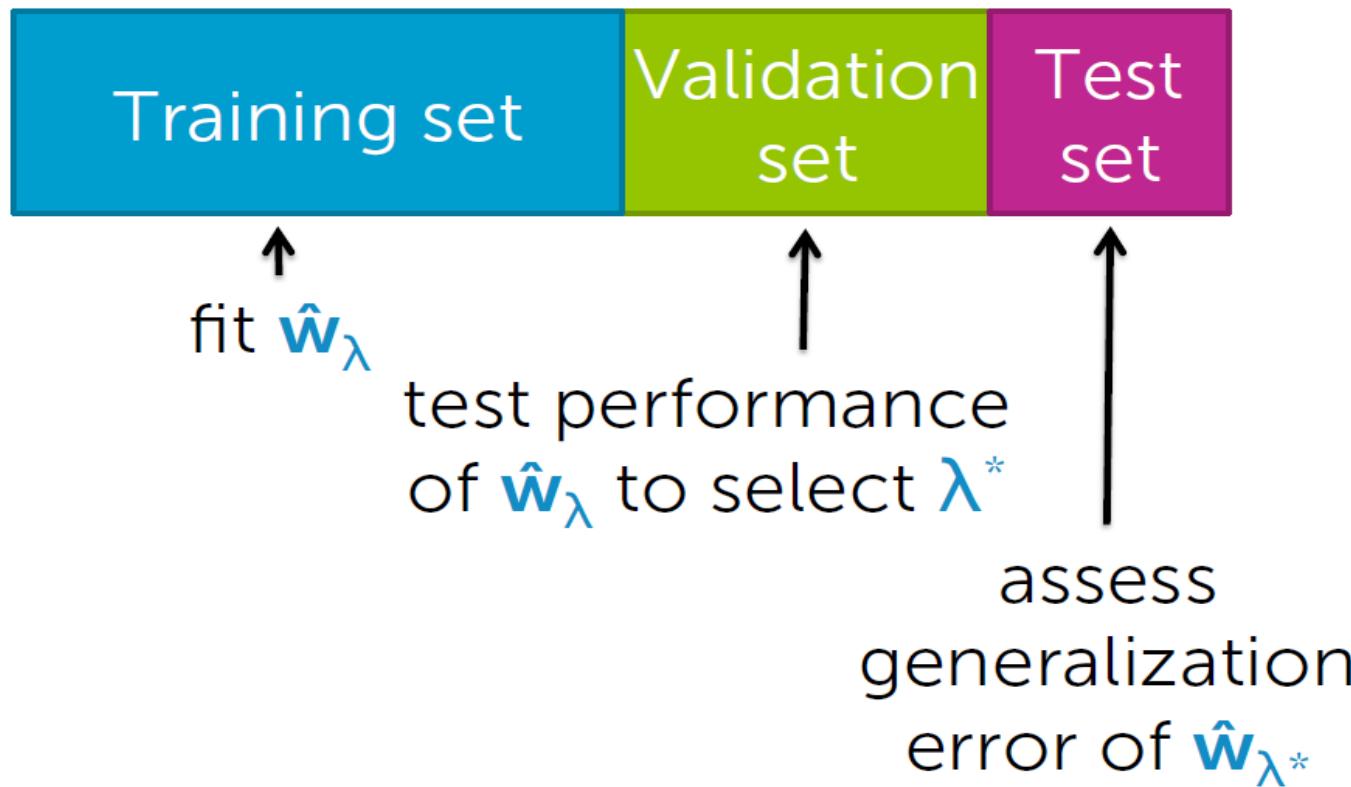
$$\mathbf{w}_j^{(t+1)} \leftarrow (1 - 2\eta\lambda)\mathbf{w}_j^{(t)} - \eta \text{partial}[j]$$

$$t \leftarrow t + 1$$

# How to choose $\lambda$

129

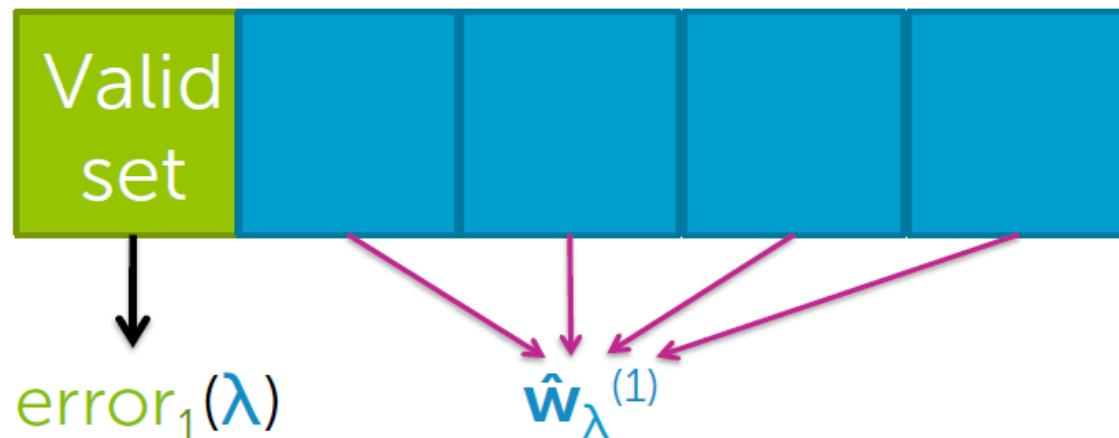
If sufficient amount of data...



# How to choose $\lambda$

130

## K-fold cross validation



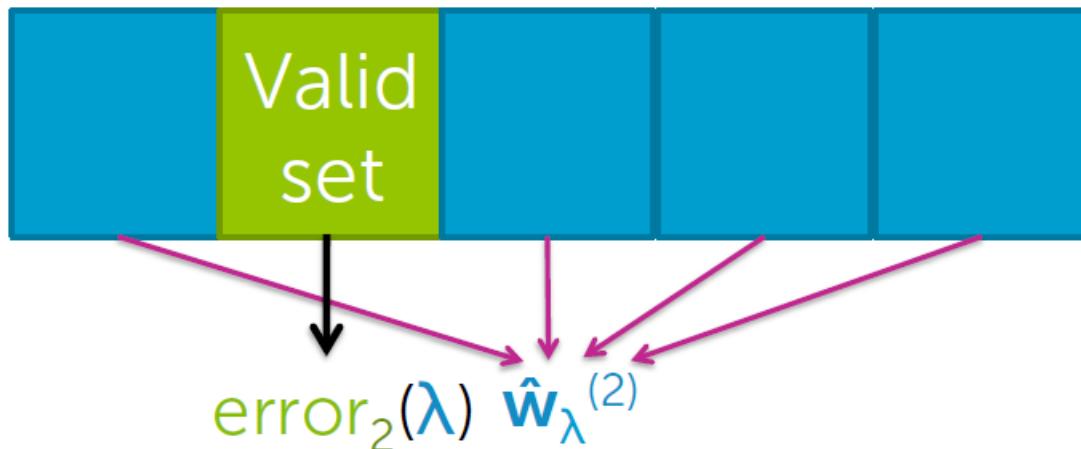
For  $k=1, \dots, K$

1. Estimate  $\hat{w}_\lambda^{(k)}$  on the training blocks
2. Compute error on validation block:  $\text{error}_k(\lambda)$

# How to choose $\lambda$

131

## K-fold cross validation



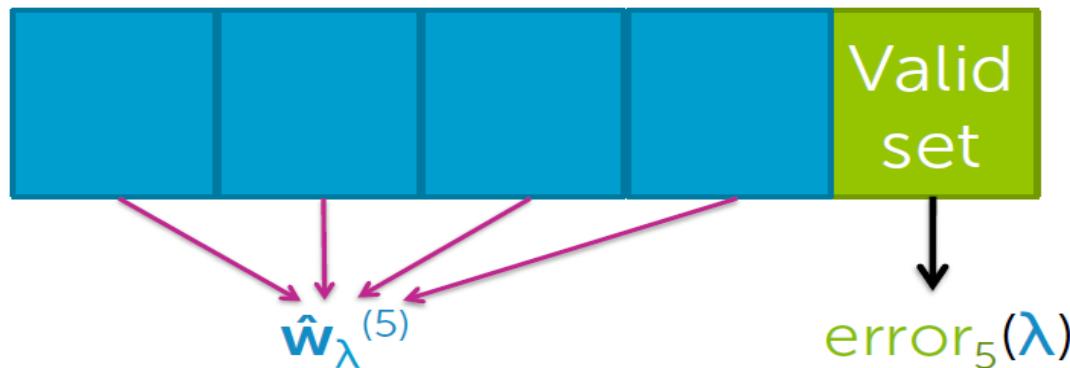
For  $k=1, \dots, K$

1. Estimate  $\hat{w}_\lambda^{(k)}$  on the training blocks
2. Compute error on validation block:  $\text{error}_k(\lambda)$

# How to choose $\lambda$

132

## K-fold cross validation



For  $k=1, \dots, K$

1. Estimate  $\hat{w}_\lambda^{(k)}$  on the training blocks
2. Compute error on validation block:  $\text{error}_k(\lambda)$

Compute average error:  $\text{CV}(\lambda) = \frac{1}{K} \sum_{k=1}^K \text{error}_k(\lambda)$

# How to choose $\lambda$

133

## K-fold cross validation



Repeat procedure for each choice of  $\lambda$

Choose  $\lambda^*$  to minimize  $CV(\lambda)$

# What value of K

134

Formally, the best approximation occurs for validation sets of size 1 ( $K=N$ )

leave-one-out  
cross validation

Computationally intensive

- requires computing  $N$  fits of model per  $\lambda$

Typically,  $K=5$  or 10

5-fold CV

10-fold CV

# How to handle the intercept

135

## 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.,  $\mathbf{x}[d]$*

# Do we penalize intercept?

136

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...

# Do we penalize intercept?

137

- **Option 1: don't penalize intercept**

Modified ridge regression cost:

$$\text{RSS}(\mathbf{w}_0, \mathbf{w}_{\text{rest}}) + \lambda \|\mathbf{w}_{\text{rest}}\|_2^2$$

- **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)

# What you can do now

138

- 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
- Implement K-fold cross validation to select the ridge regression tuning parameter  $\lambda$

# FEATURES SELECTION & LASSO REGRESSION

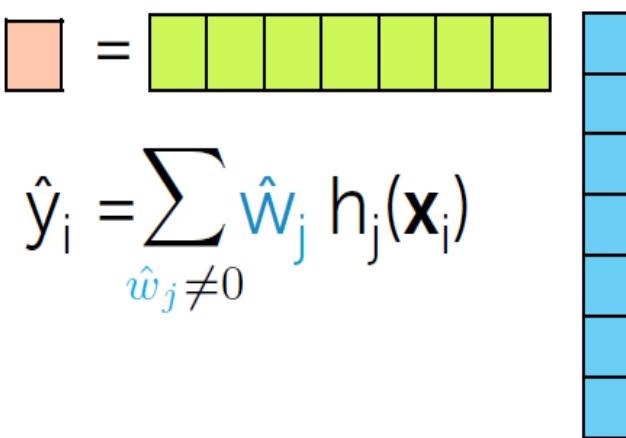
# Why features selection?

140

## Efficiency:

- If  $\text{size}(\mathbf{w}) = 100B$ , each prediction is expensive
- If  $\hat{\mathbf{w}}$  sparse, computation only depends on # of non-zeros

many zeros

$$\hat{y}_i = \sum_{\hat{w}_j \neq 0} \hat{w}_j h_j(\mathbf{x}_i)$$


## Interpretability:

- Which features are relevant for prediction?

# Sparcity

141

## Housing application



Lot size  
Single Family

Year built

Last sold price

Last sale price/sqft

Finished sqft

Unfinished sqft

Finished basement sqft

# floors

Flooring types

Parking type

Parking amount

Cooling

Heating

Exterior materials

Roof type

Structure style

Dishwasher  
Garbage disposal

Microwave

Range / Oven

Refrigerator

Washer

Dryer

Laundry location

Heating type

Jetted Tub

Deck

Fenced Yard

Lawn

Garden

Sprinkler System

:

# Sparcity

142

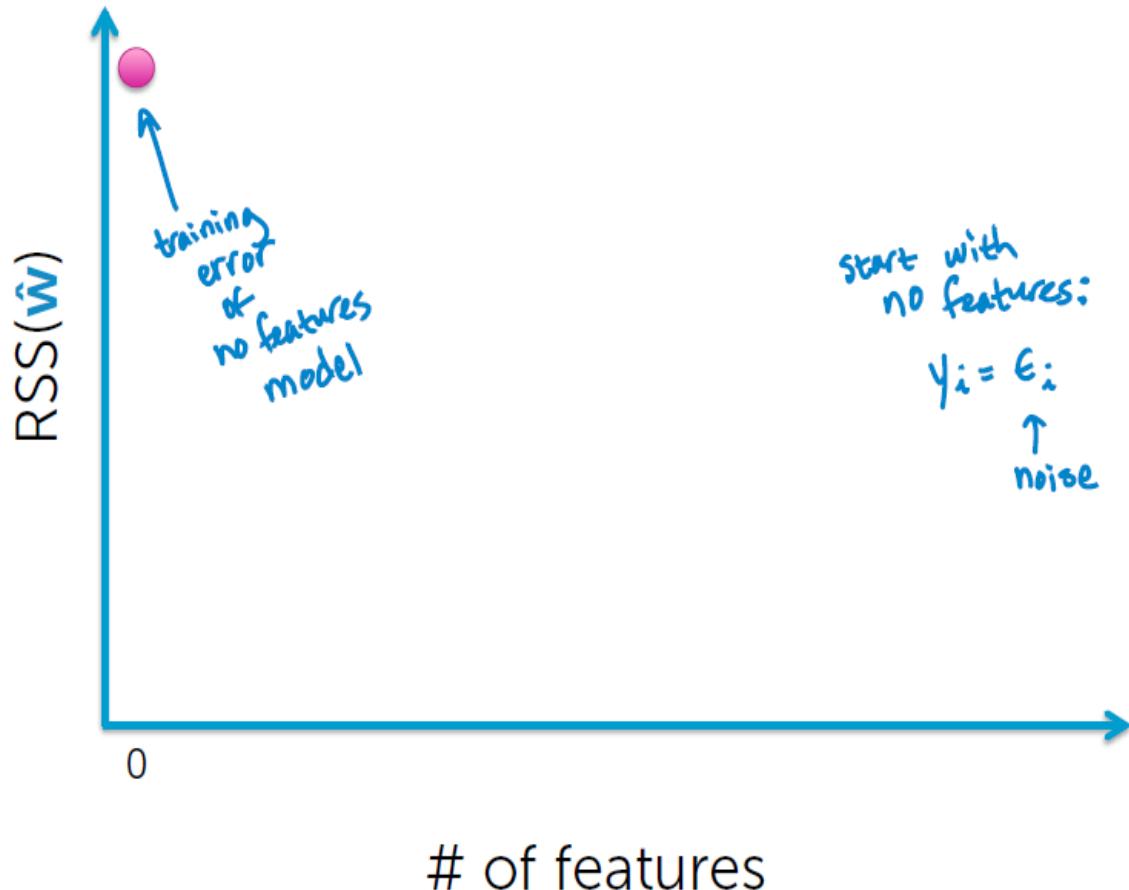
## Reading your mind



Activity in which  
brain regions  
can predict  
happiness?

# Find best model of size: 0

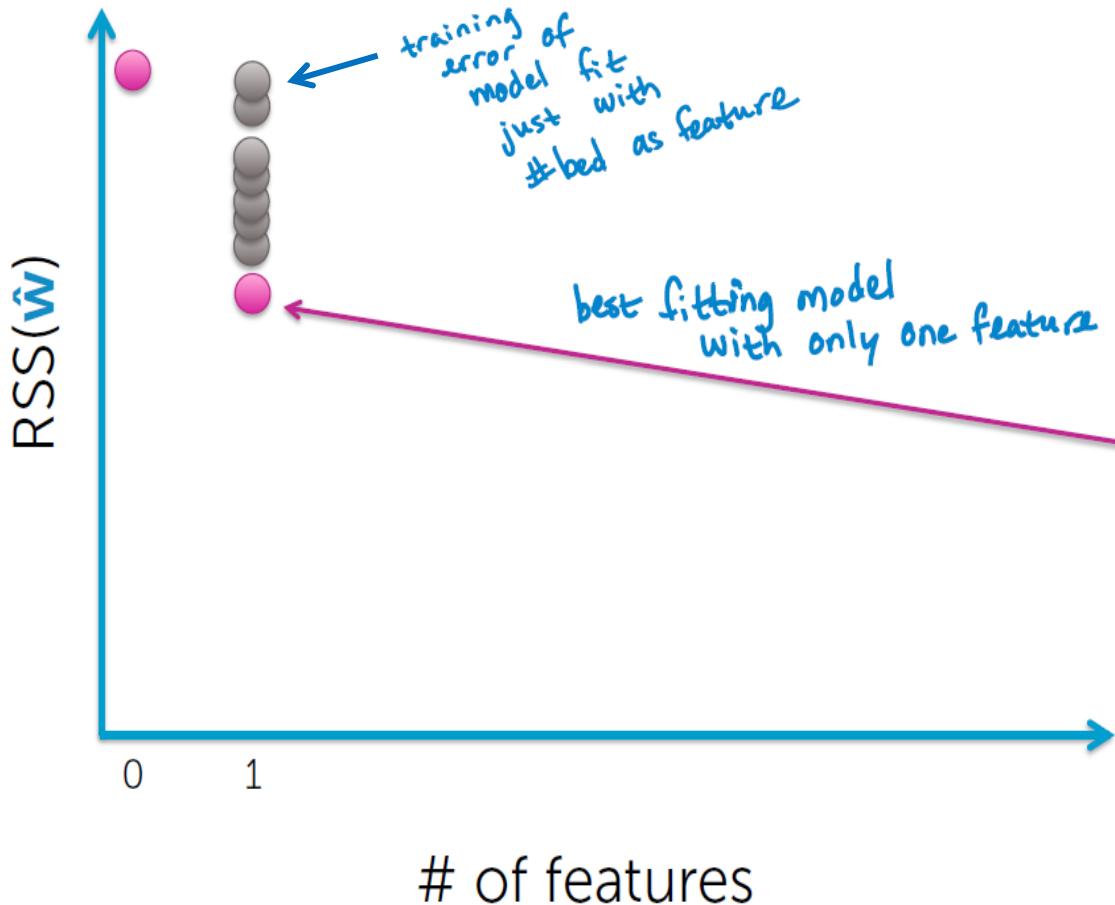
143



- ~~- # bedrooms  
- # bathrooms  
- sq.ft. living  
- sq.ft. lot  
- floors  
- year built  
- year renovated  
- waterfront~~

# Find best model of size: 1

144

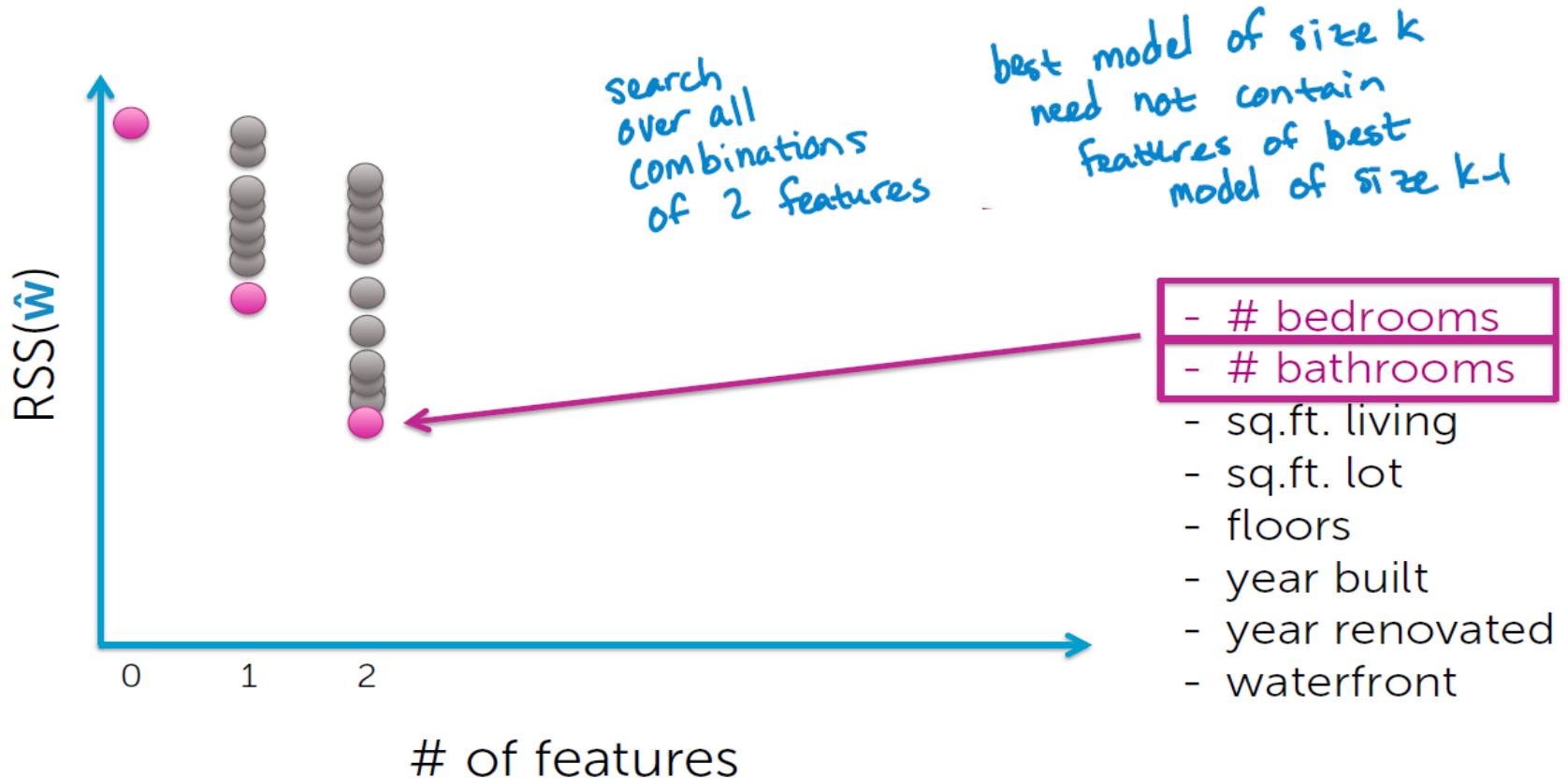


- # bedrooms
- # bathrooms
- **sq.ft. living**
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

# Find best model of size: 2

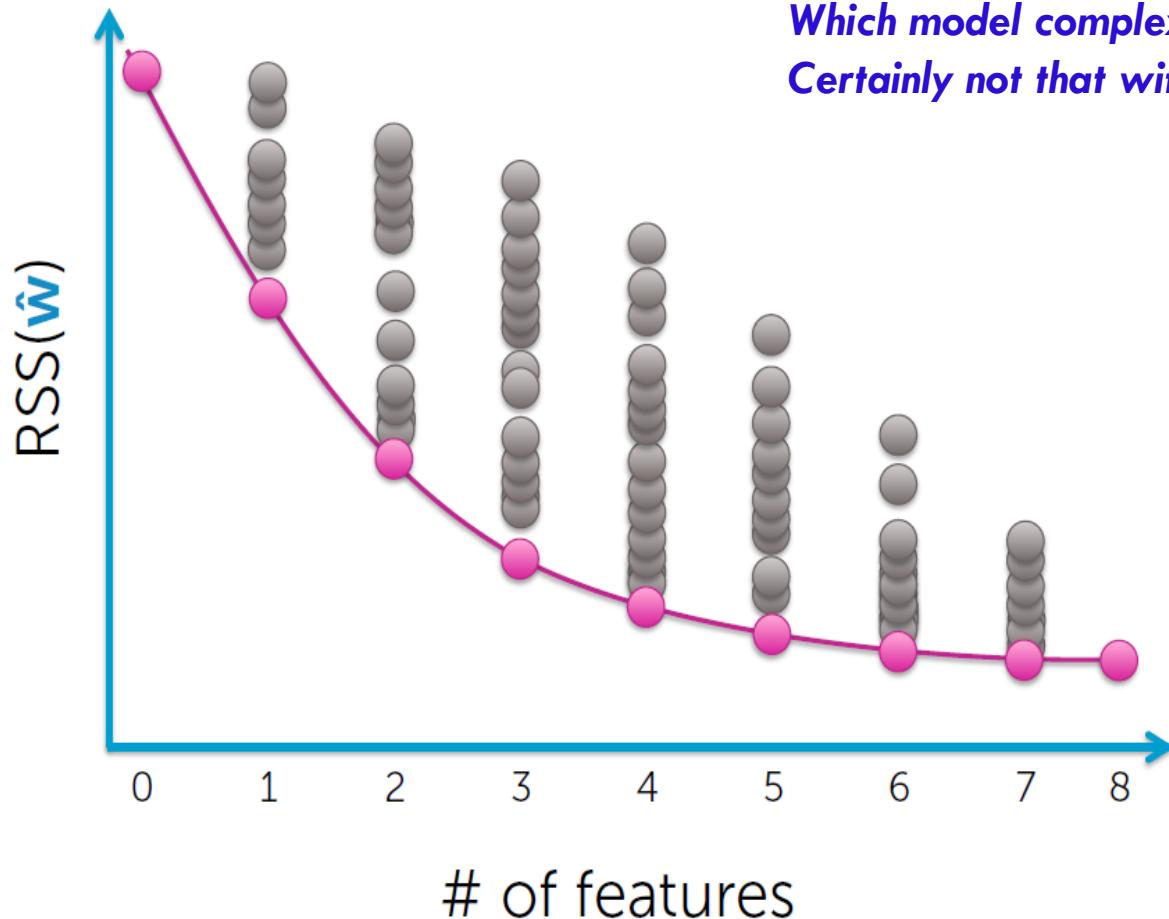
145

**Note: not necessarily nested!**



# Find best model of size: N

146



*Which model complexity to choose?  
Certainly not that with the smallest training error!*

- # bedrooms
- # bathrooms
- sq.ft. living
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

# Choosing model complexity

147

Option 1: Assess on validation set

Option 2: Cross validation

Option 3+: Other metrics for penalizing  
model complexity like BIC...

# Complexity of „all subsets”

148

How many models were evaluated?

- each indexed by features included

$$y_i = \varepsilon_i$$

$$y_i = w_0 h_0(x_i) + \varepsilon_i$$

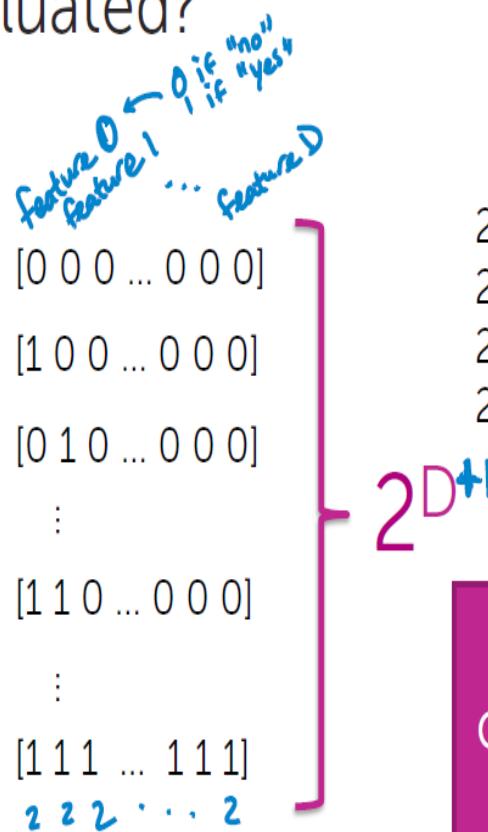
$$y_i = w_1 h_1(x_i) + \varepsilon_i$$

:

$$y_i = w_0 h_0(x_i) + w_1 h_1(x_i) + \varepsilon_i$$

:

$$y_i = w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_D h_D(x_i) + \varepsilon_i$$



$$2^8 = 256$$

$$2^{30} = 1,073,741,824$$

$$2^{1000} = 1.071509 \times 10^{301}$$

$2^{100B}$  = HUGE!!!!!!

Typically,  
computationally  
infeasible

# Greedy algorithm

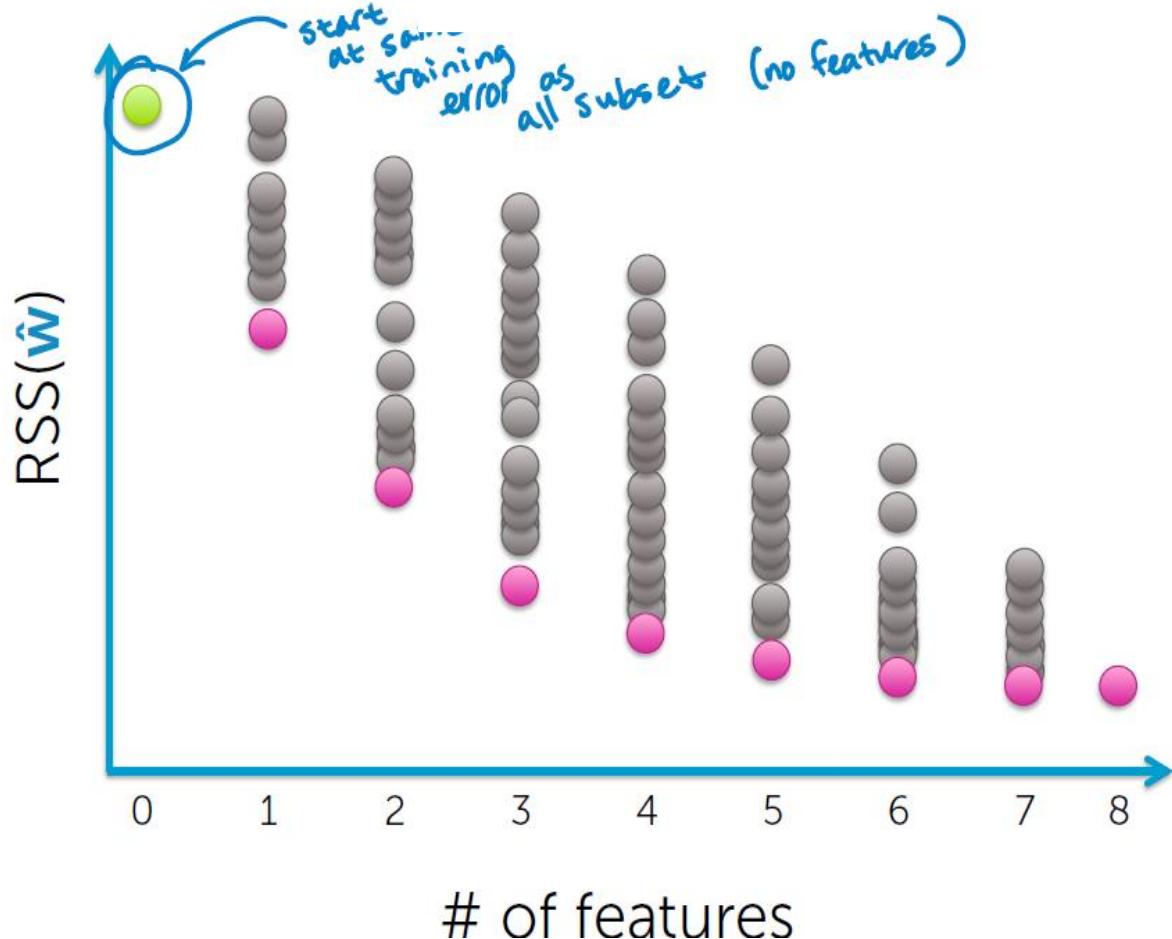
149

## Forward stepwise algorithm

1. Pick a dictionary of features  $\{h_0(\mathbf{x}), \dots, h_D(\mathbf{x})\}$ 
  - e.g., polynomials for linear regression
2. Greedy heuristic:
  - i. Start with empty set of features  $F_0 = \emptyset$   
(or simple set, like just  $h_0(\mathbf{x})=1 \rightarrow y_i = w_0 + \epsilon_i$ )
  - ii. Fit model using current feature set  $F_t$  to get  $\hat{\mathbf{w}}^{(t)}$
  - iii. Select next best feature  $h_{j^*}(\mathbf{x})$ 
    - e.g.,  $h_j(\mathbf{x})$  resulting in lowest training error  
when learning with  $F_t + \{h_j(\mathbf{x})\}$
  - iv. Set  $F_{t+1} \leftarrow F_t + \{h_{j^*}(\mathbf{x})\}$
  - v. Recurse

# Visualizing greedy algorithm

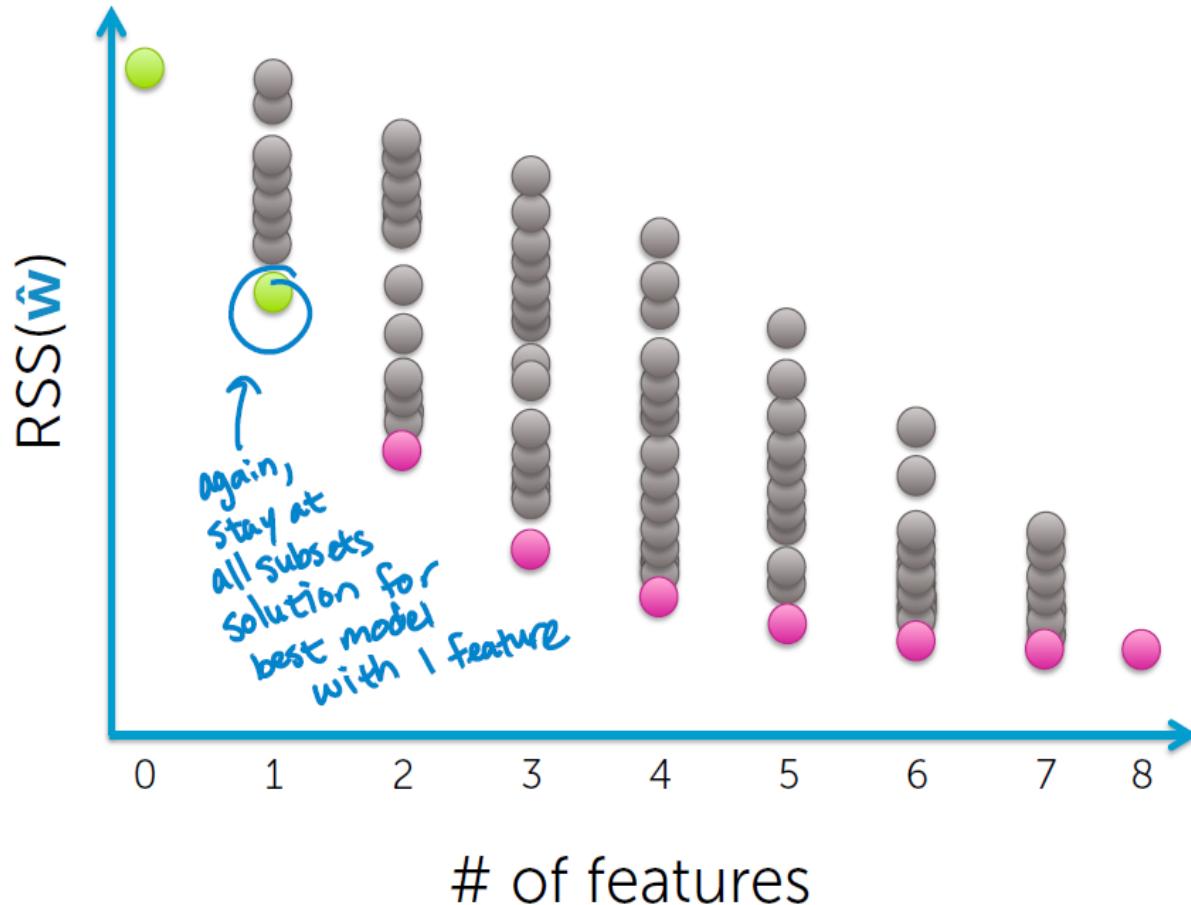
150



- # bedrooms
- # bathrooms
- sq.ft. living
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

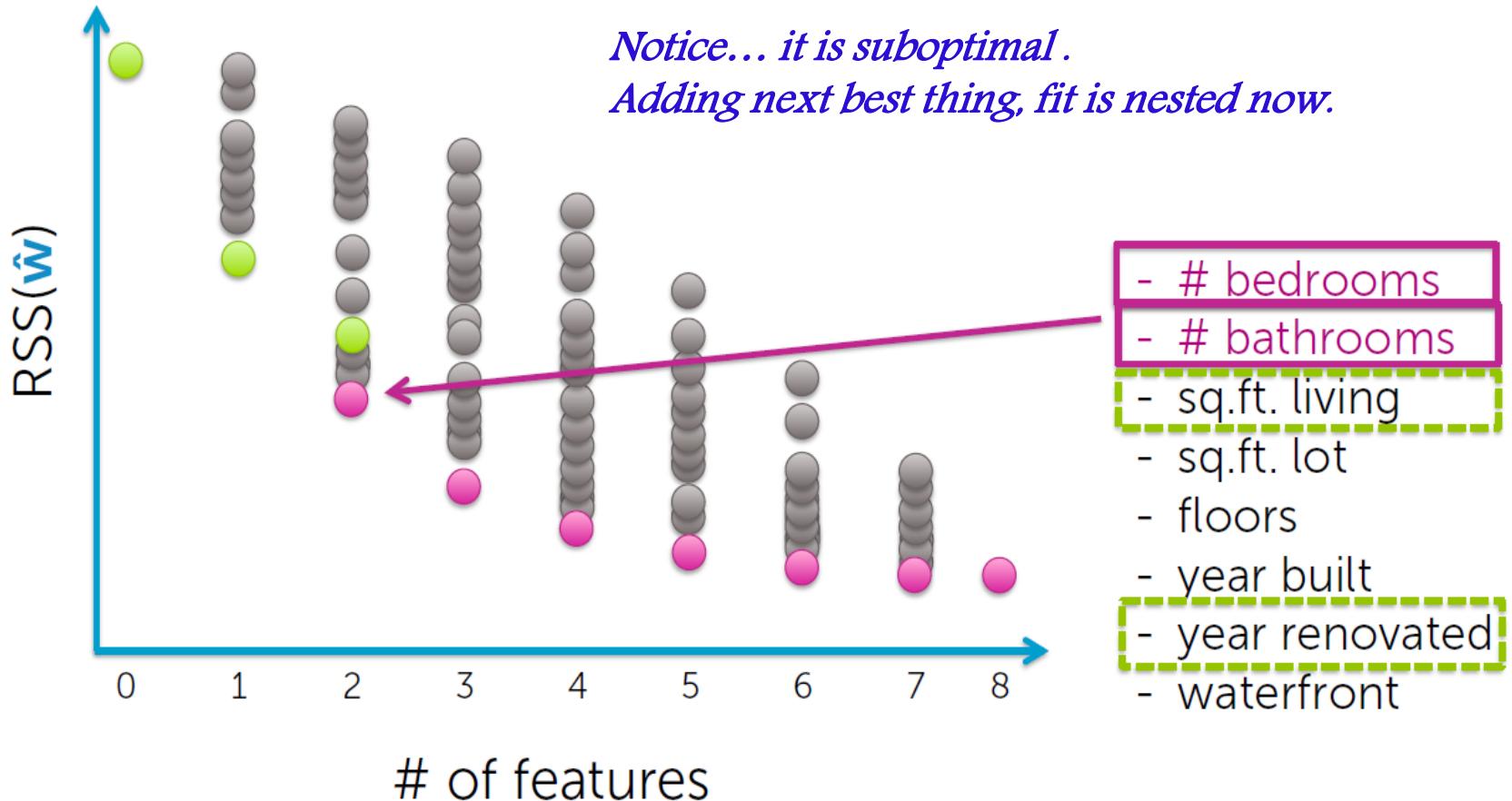
# Visualizing greedy algorithm

151



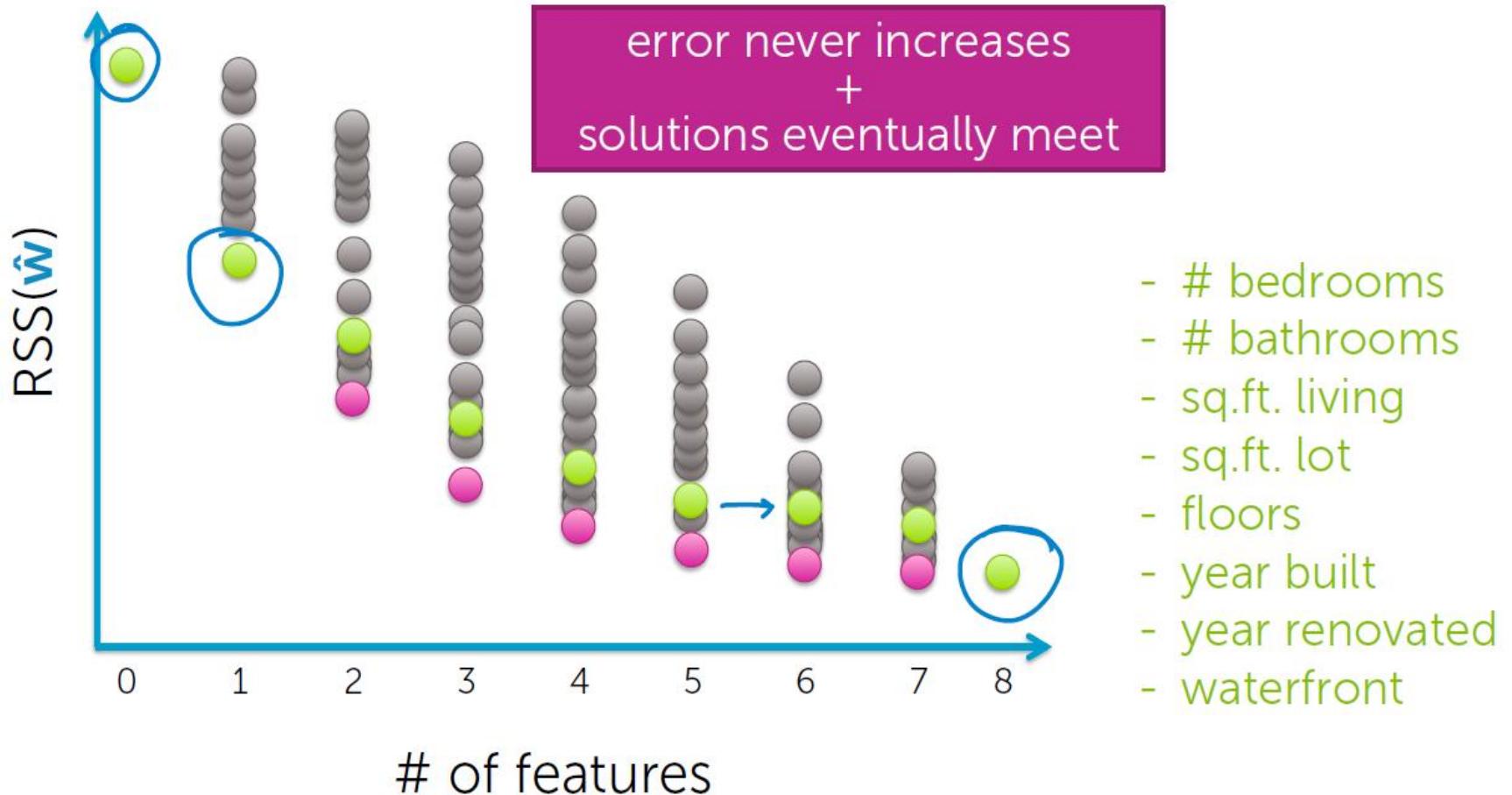
# Visualizing greedy algorithm

152



# Visualizing greedy algorithm

153



# When do we stop?

154

When training error is low enough?

No!

When test error is low enough?

No!

***Use validation set or cross validation!***

# Complexity of forward stepwise

155

How many models were evaluated?

- 1<sup>st</sup> step, D models
- 2<sup>nd</sup> step, D-1 models (add 1 feature out of D-1 possible)
- 3<sup>rd</sup> step, D-2 models (add 1 feature out of D-2 possible)
- ...

How many steps?

- Depends
- At most D steps (to full model)

$$O(D^2) \ll 2^D$$

for large D

# Other greedy algorithms

156

Instead of starting from simple model  
and always growing...

**Backward stepwise:**

Start with full model and iteratively remove  
features least useful to fit

**Combining forward and backward steps:**

In forward algorithm, insert steps to remove  
features no longer as important

*Lots of other variants, too.*

# Using regularisation for features selection

157

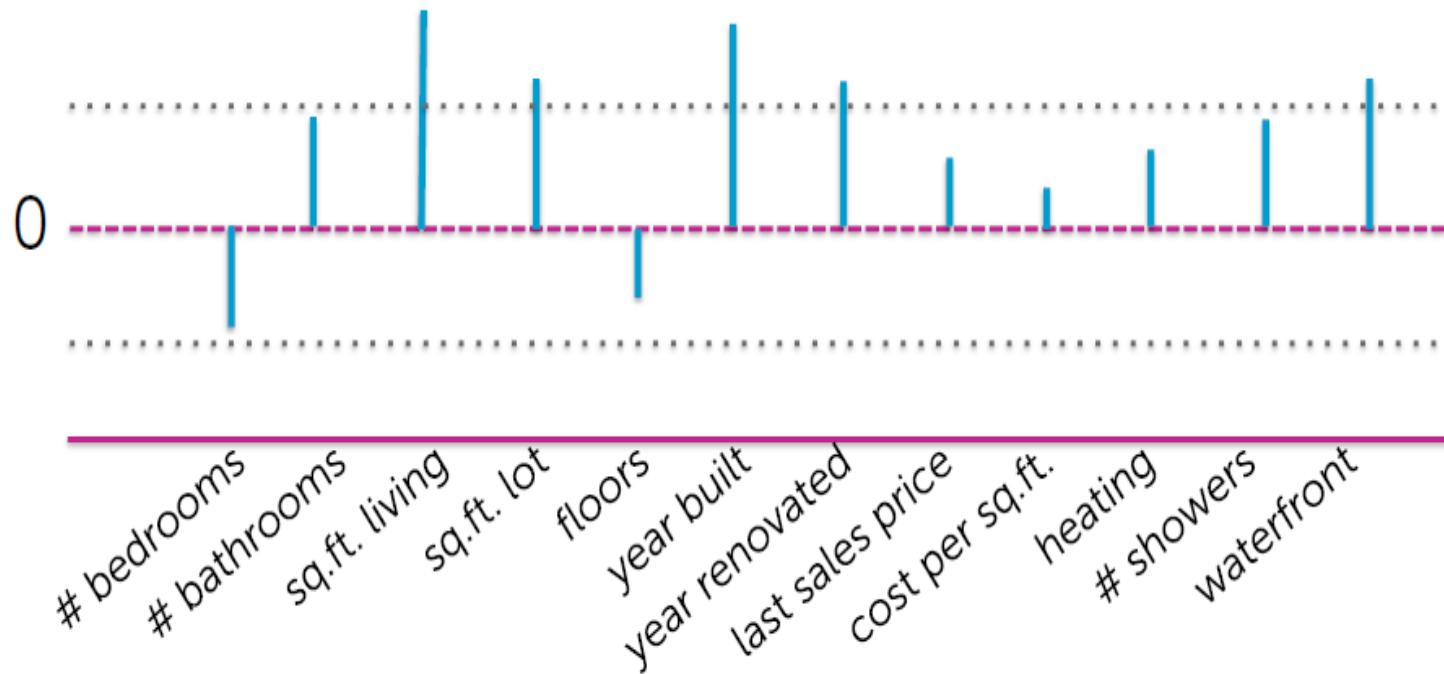
Instead of searching over a **discrete** set of solutions, can we use **regularization**?

- Start with full model (all possible features)
- “Shrink” some coefficients **exactly to 0**
  - i.e., knock out certain features
- Non-zero coefficients indicate “selected” features

# Thresholding ridge coefficients?

158

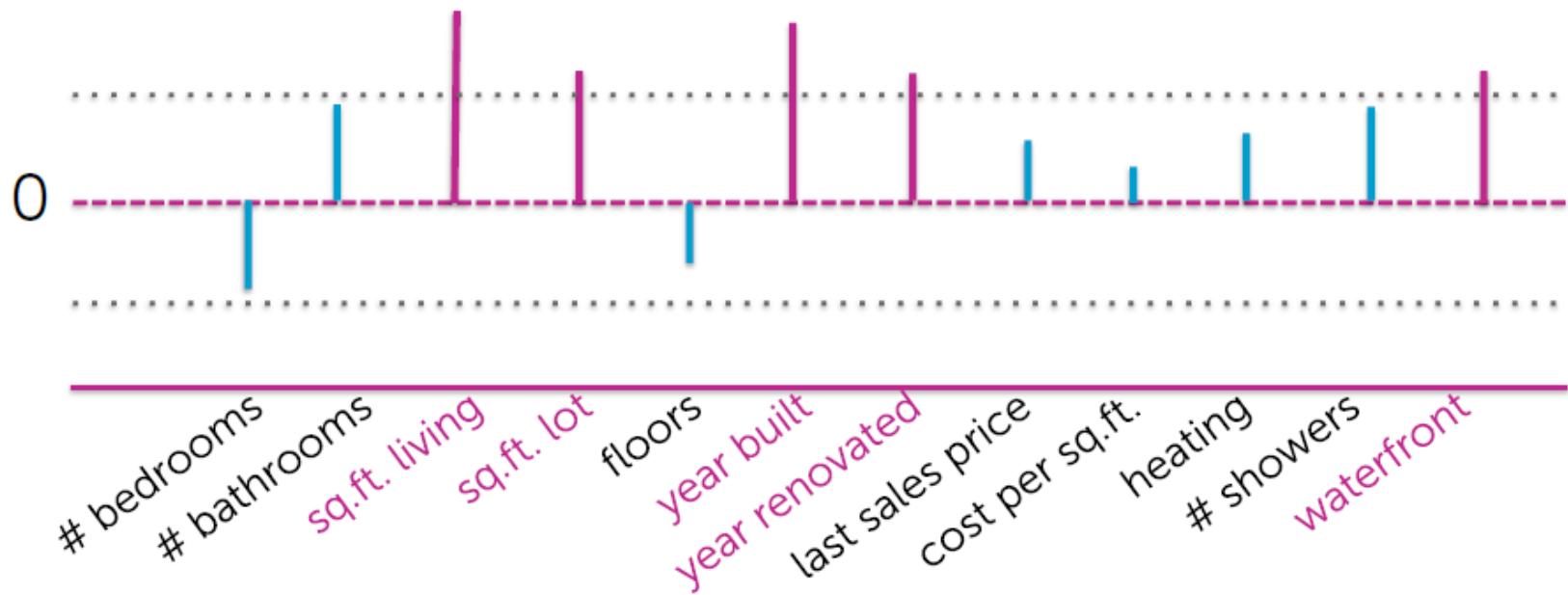
Why don't we just set small ridge coefficients to 0?



# Thresholding ridge coefficients?

159

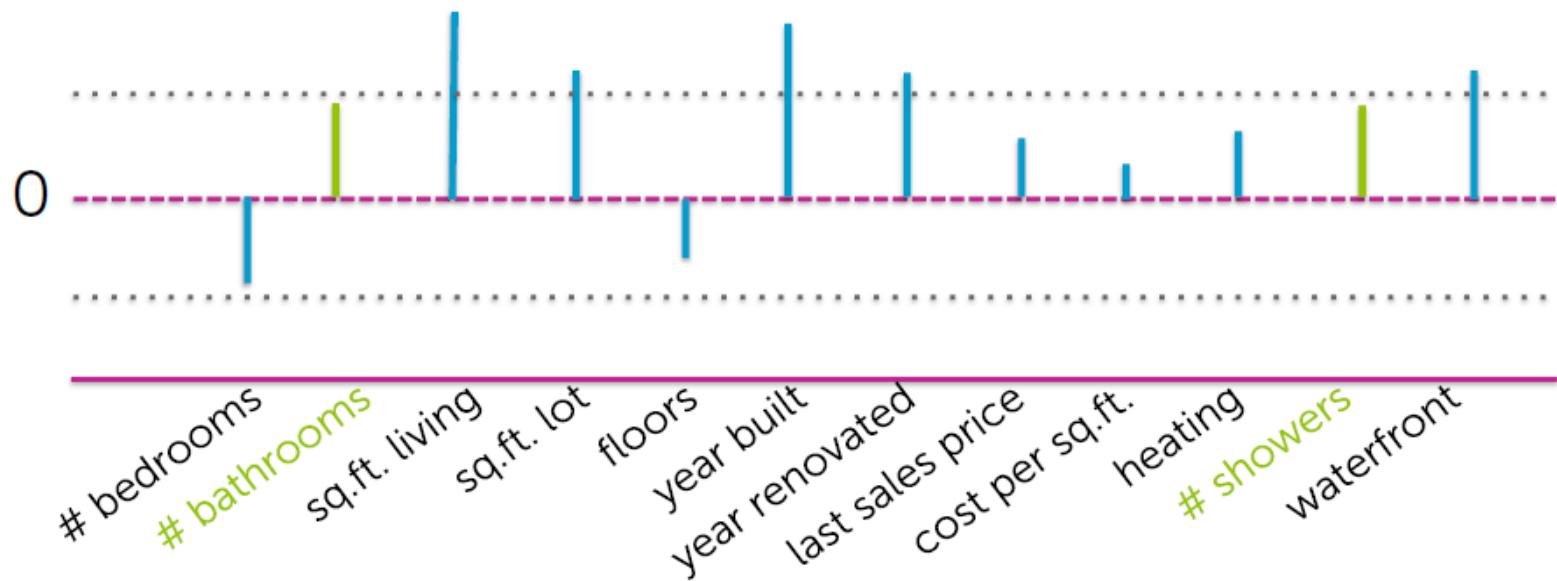
Selected features for a given threshold value



# Thresholding ridge coefficients?

160

Let's look at two related features...

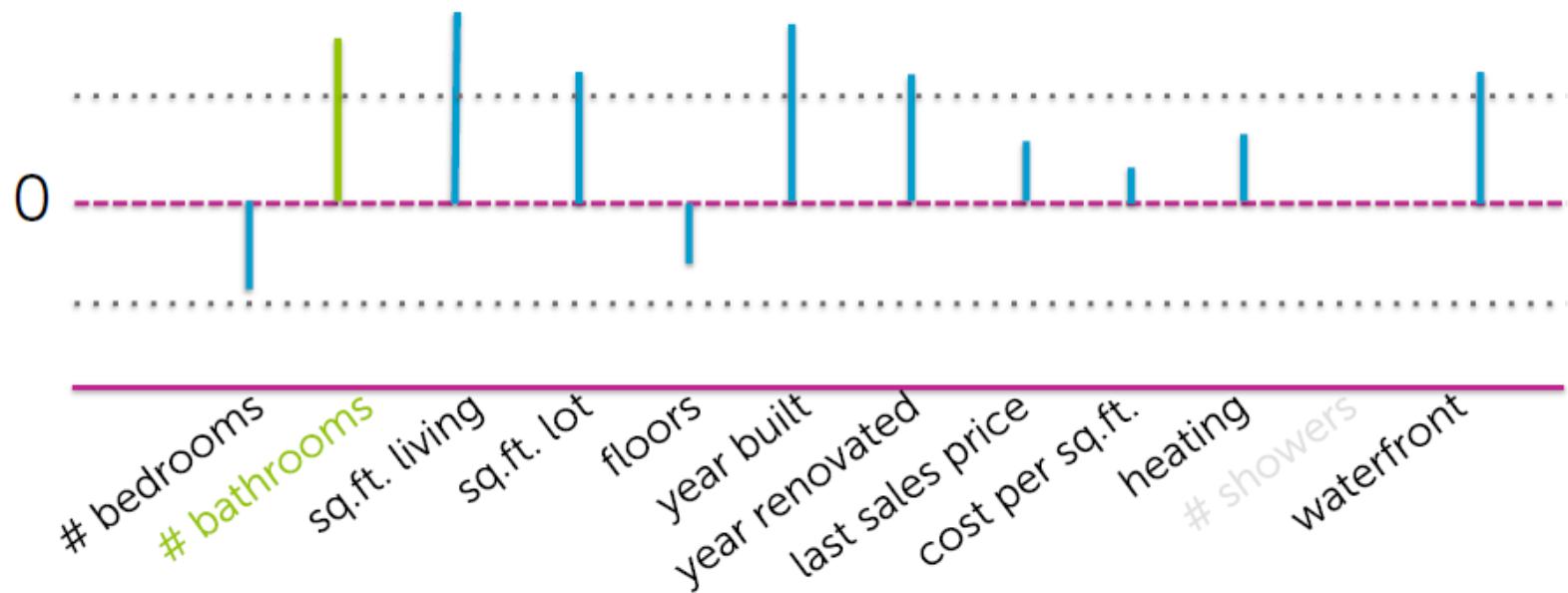


Nothing measuring bathrooms was included!

# Thresholding ridge coefficients?

161

If only one of the features had been included...



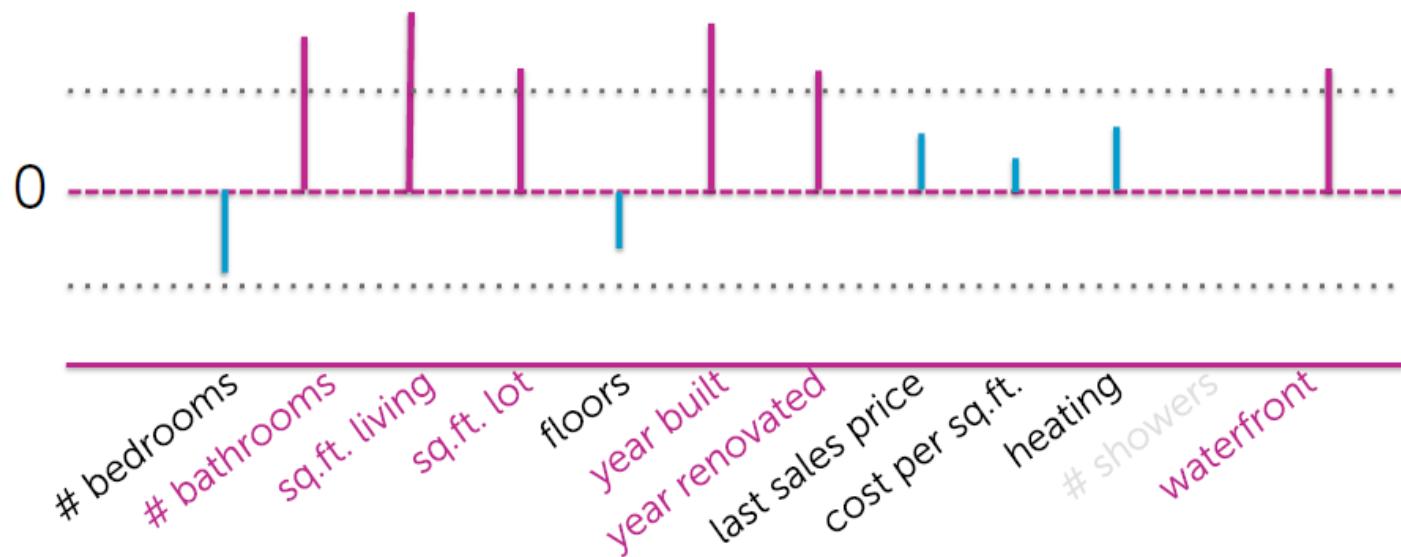
*Remember:*

*this is linear model. If we assume that #showers = #bathrooms and remove one of them from the model, coefficients will sum up.*

# Thresholding ridge coefficients?

162

Would have included bathrooms in selected model



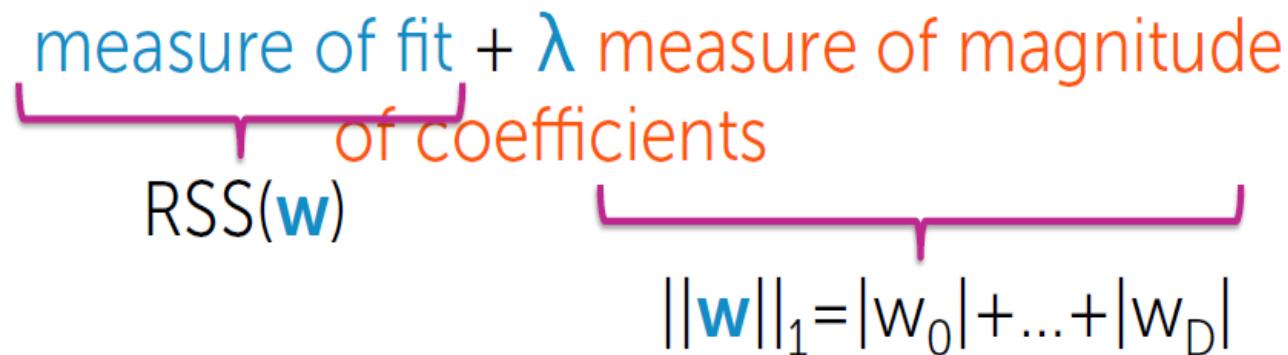
Can regularization lead directly to sparsity?

# Try this cost instead of ridge ...

163

Total cost =

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

  
of coefficients

$$\text{RSS}(\mathbf{w})$$
$$\|\mathbf{w}\|_1 = |w_0| + \dots + |w_D|$$

Lasso regression  
(a.k.a.  $L_1$  regularized regression)

Leads to  
sparse  
solutions!

# Lasso regression

164

Just like ridge regression, solution is governed by a continuous parameter  $\lambda$

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

↑ tuning parameter = balance of fit and sparsity

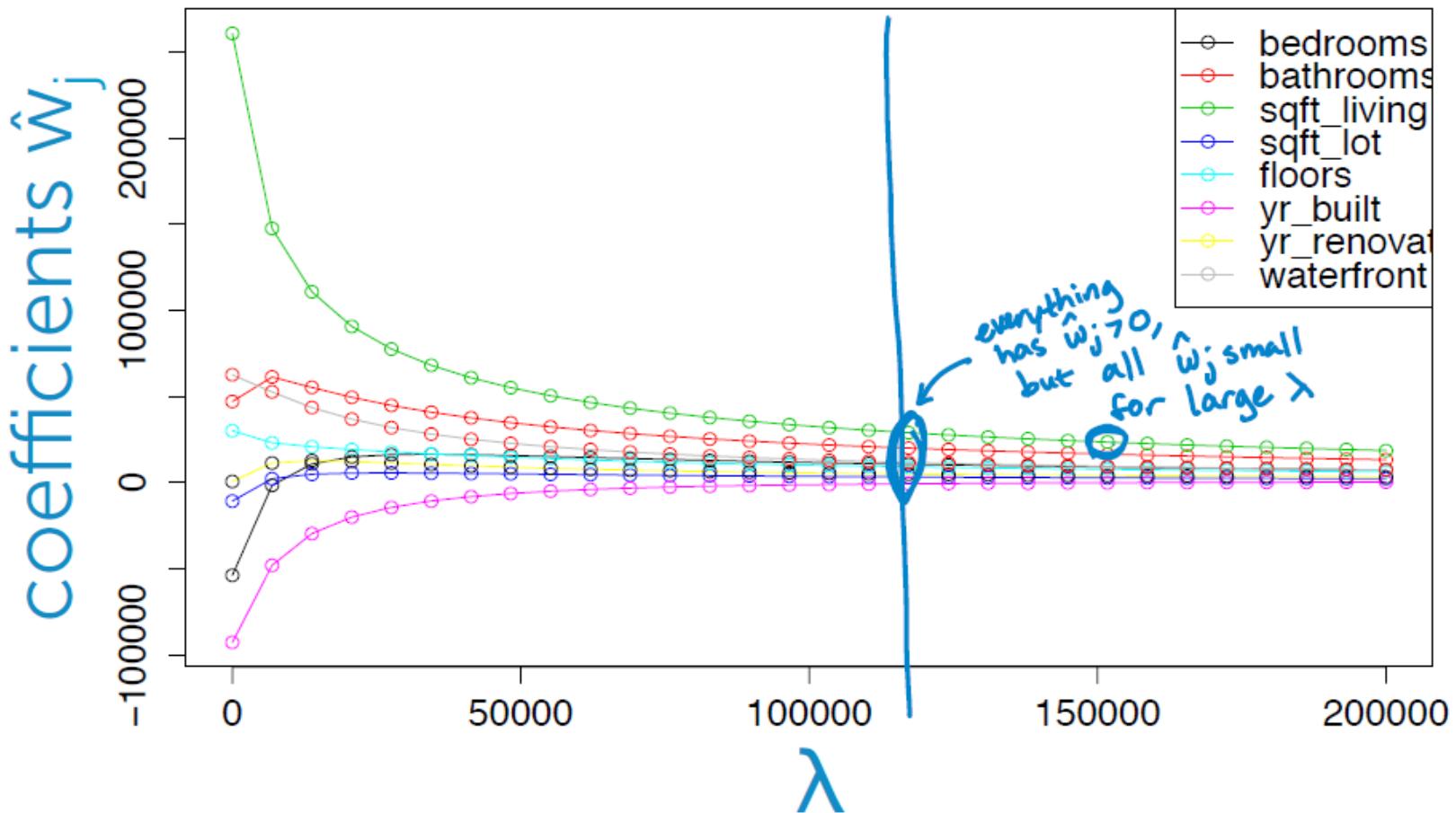
If  $\lambda=0$ :  $\hat{\mathbf{w}}^{\text{lasso}} = \hat{\mathbf{w}}^{\text{LS}}$  (unregularized solution)

If  $\lambda=\infty$ :  $\hat{\mathbf{w}}^{\text{lasso}} = \mathbf{0}$

If  $\lambda$  in between:  $\mathbf{0} \leq \|\hat{\mathbf{w}}^{\text{lasso}}\|_1 \leq \|\hat{\mathbf{w}}^{\text{LS}}\|_1$

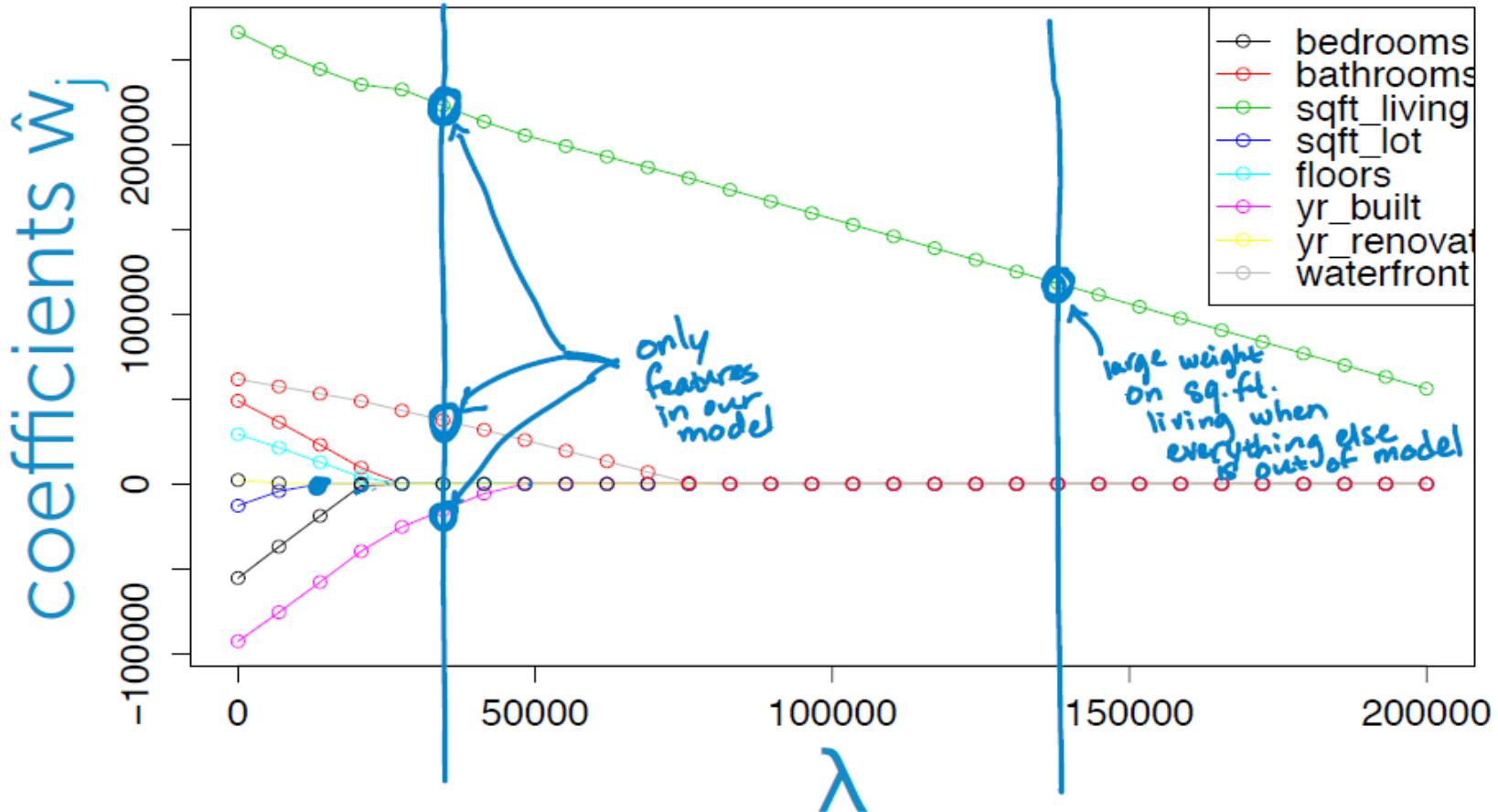
# Coefficient path: ridge

165



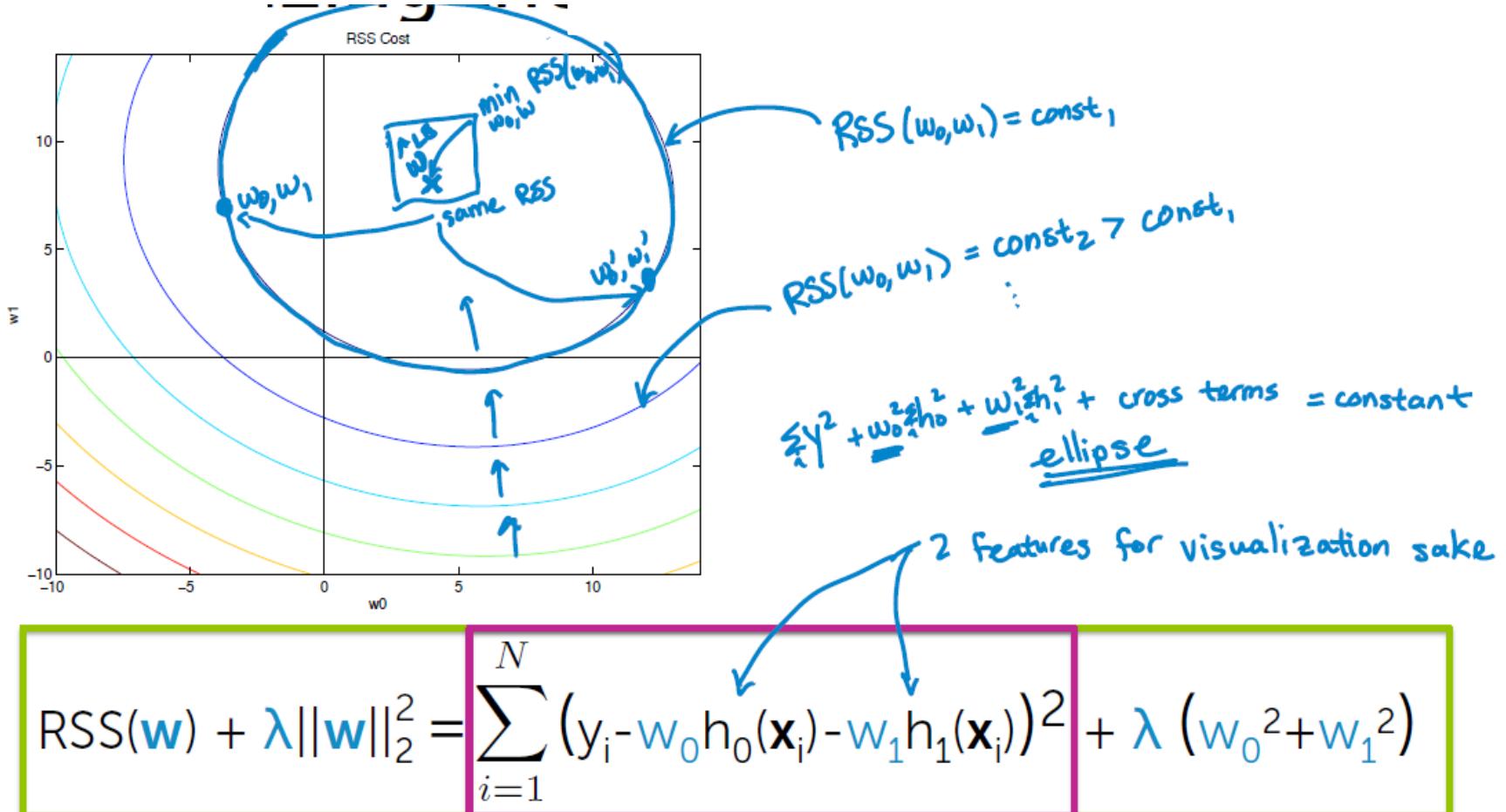
# Coefficient path: lasso

166



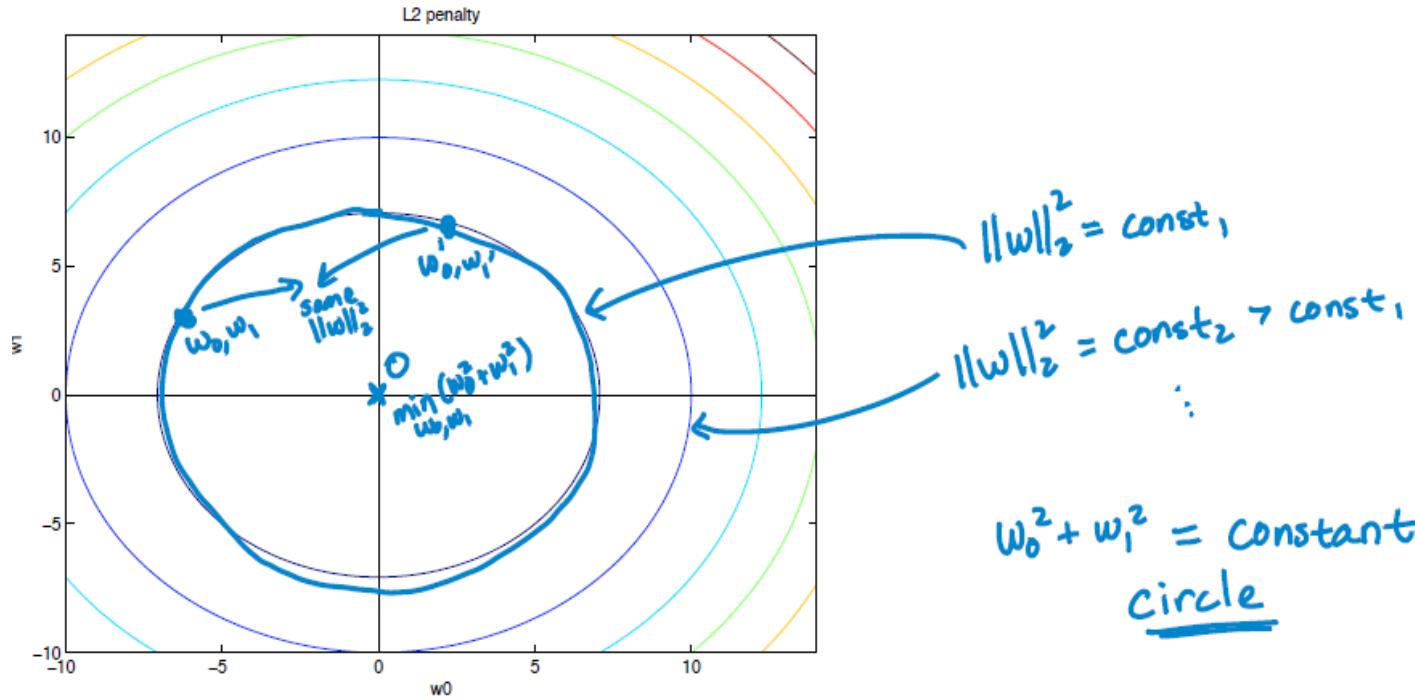
# Visualising ridge cost in 2D

167



# Visualising ridge cost in 2D

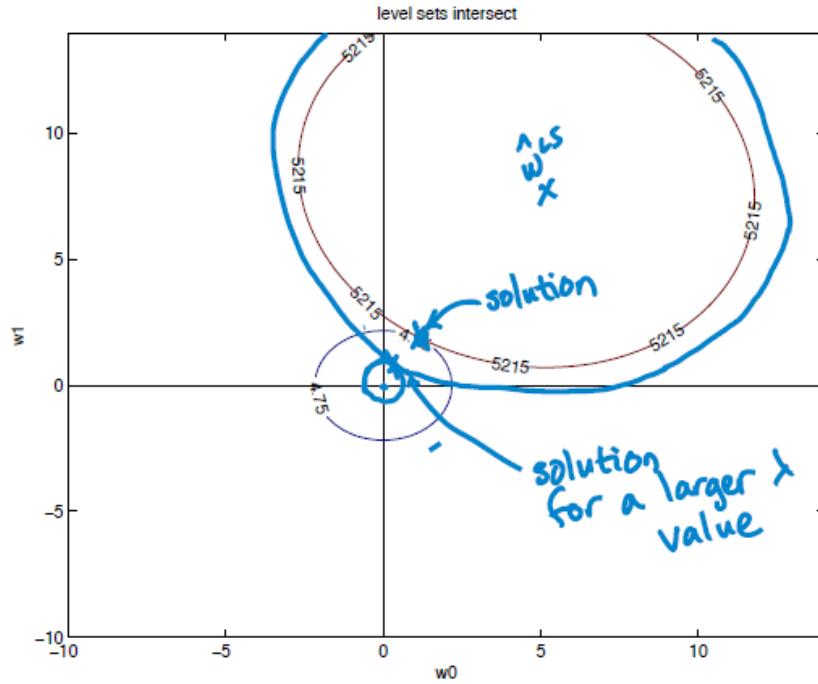
168



$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \underline{\lambda} (w_0^2 + w_1^2)$$

# Visualising ridge cost in 2D

169

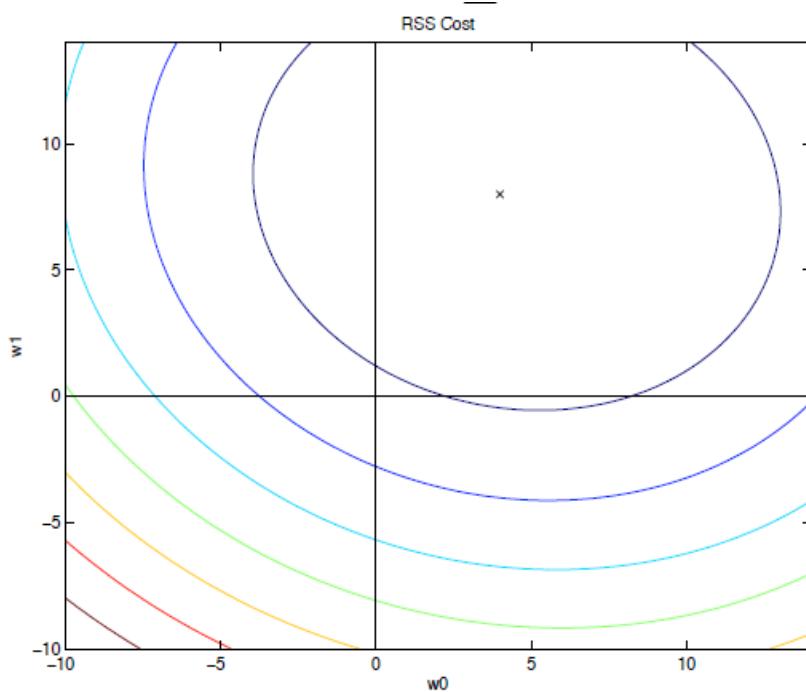


For a specific  $\lambda$  value,  
some balance between  
RSS and  $\|w\|_2^2$

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \lambda (w_0^2 + w_1^2)$$

# Visualising lasso cost in 2D

170

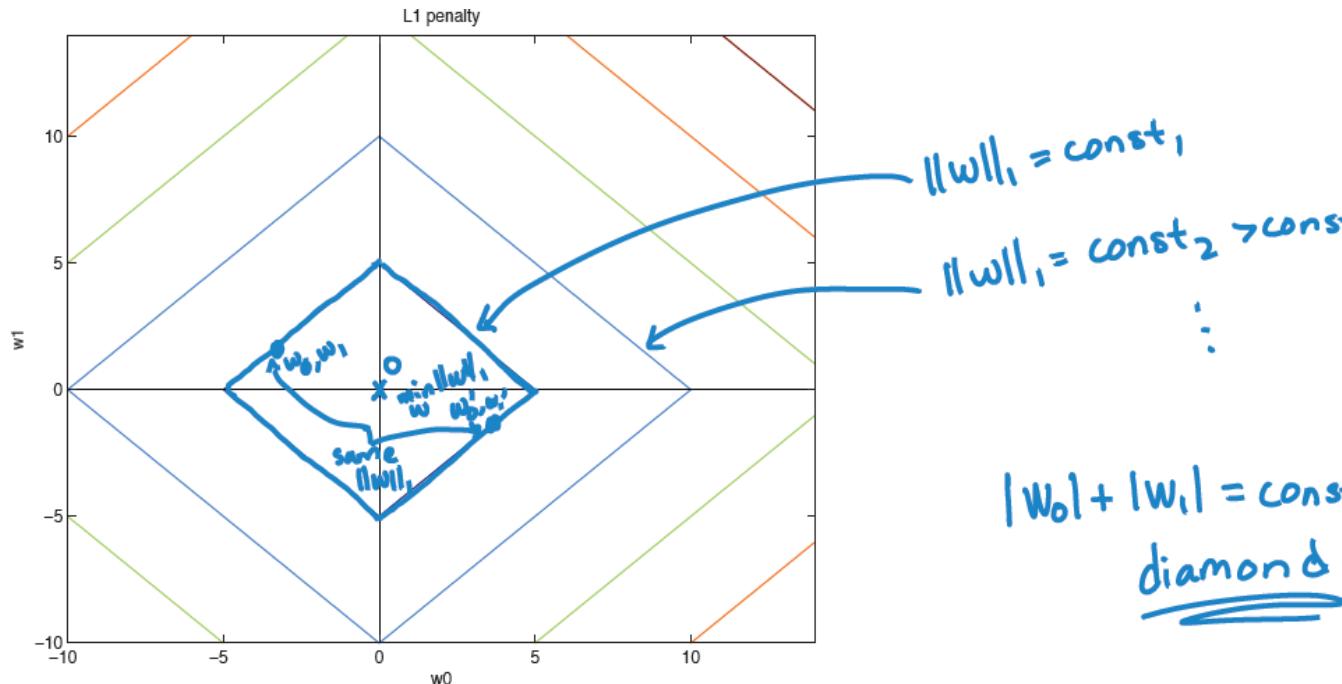


RSS contours for  
lasso are exactly  
the same as  
those for ridge!

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \lambda (|w_0| + |w_1|)$$

# Visualising lasso cost in 2D

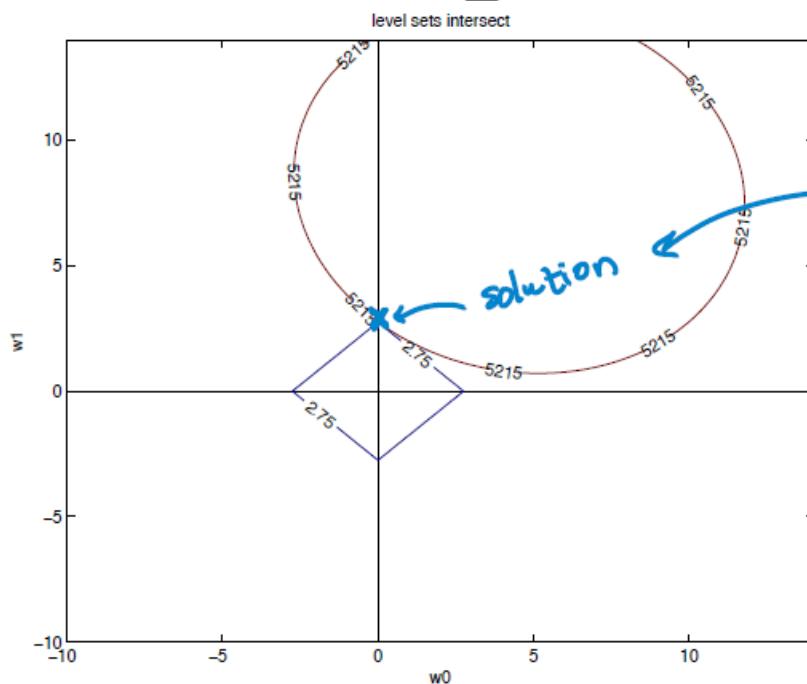
171



$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \lambda (|w_0| + |w_1|)$$

# Visualising lasso cost in 2D

172



For a specific value of  $\lambda$ ,

We are getting sparse solution,  
the  $w_0 = 0$

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \lambda (|w_0| + |w_1|)$$

# How we optimise for objective

173

To solve for  $\hat{\mathbf{w}}$ , previously took gradient of total cost objective and either:

- 1) Derived closed-form solution
- 2) Used in gradient descent algorithm

# Optimise for lasso objective

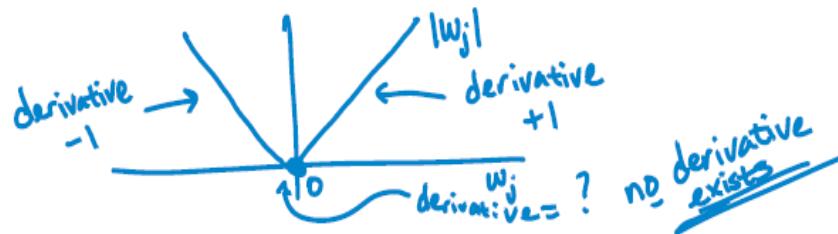
174

Lasso total cost:  $\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$

$\|\mathbf{w}\|_1 = \sum_{j=0}^p |w_j|$

Issues:

- 1) What's the derivative of  $|w_j|$ ?



gradients → subgradients

- 2) Even if we could compute derivative, no closed-form solution

can use subgradient descent

# Coordinate descent

175

Goal: Minimize some function  $g$

$$\min_w g(w)$$

$$g(w) = g(w_0, w_1, \dots, w_D)$$

when keeping others fixed

Often, hard to find minimum for all coordinates, but easy for each coordinate

Coordinate descent:

Initialize  $\hat{w} = 0$  (or smartly...)

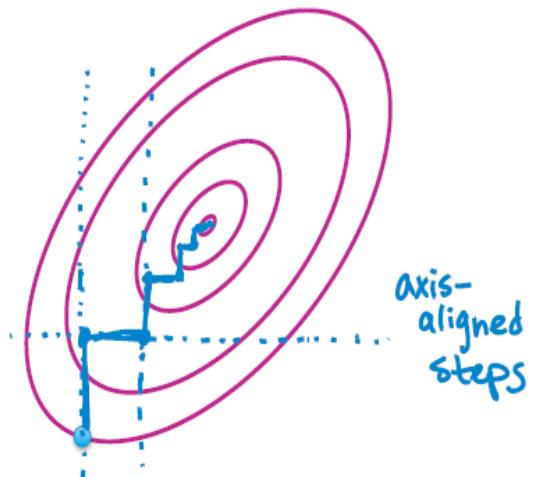
while not converged

pick a coordinate  $j$

$$\hat{w}_j \leftarrow$$

$$\min_w g(\hat{w}_0, \dots, \hat{w}_{j-1}, w_j, \hat{w}_{j+1}, \dots, \hat{w}_D)$$

values from previous iterations  
just min over  $j$ th coordinate



# Comments on coordinate descent

176

How do we pick next coordinate?

- At random ("random" or "stochastic" coordinate descent), round robin, ...

No stepsize to choose!

Super useful approach for *many* problems

- Converges to optimum in some cases (e.g., "strongly convex")
- Converges for lasso objective

# Normalizing features

177

## Normalizing features

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

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

Normalizer:  $Z_j$

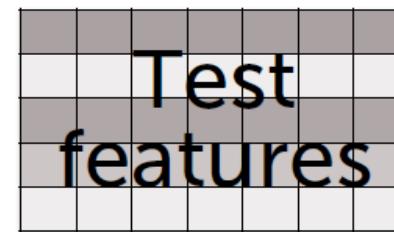
Apply same training scale factors  
to test data:

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

Normalizer:  $Z_j$

apply to test point

summing over training points



# Optimising least squares objective

178

## One coordinate at a time

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N \left( y_i - \sum_{j=0}^D w_j h_j(\mathbf{x}_i) \right)^2$$

normalized  
features

Fix all coordinates  $w_{-j}$  and take partial w.r.t.  $w_j$

$\frac{\partial}{\partial w_j} \text{RSS}(\mathbf{w}) = -2 \sum_{i=1}^N h_j(\mathbf{x}_i) \left( y_i - \sum_{k=0}^D w_k h_k(\mathbf{x}_i) \right)$

$= -2 \sum_{i=1}^N h_j(\mathbf{x}_i) \left( y_i - \underbrace{\sum_{k \neq j} w_k h_k(\mathbf{x}_i)}_{\text{all } w_k \text{ for } k \neq j} - w_j h_j(\mathbf{x}_i) \right)$

$= -2 \sum_{i=1}^N h_j(\mathbf{x}_i) \left( y_i - \underbrace{\sum_{k \neq j} w_k h_k(\mathbf{x}_i)}_{\text{all } w_k \text{ for } k \neq j} \right) + 2 w_j \boxed{\sum_{i=1}^N h_j(\mathbf{x}_i)^2}$

$= -2 p_j + 2 w_j$

↑  
 $p_j \triangleq \sum_{i=1}^N h_j(\mathbf{x}_i)^2$

↑  
 $\text{by definition of normalized features, } = 1$

↑  
 $\text{id optimization coordinate by coordinate}$

# Optimising least squares objective

179

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N \left( y_i - \sum_{j=0}^D w_j h_j(\mathbf{x}_i) \right)^2$$

Set partial = 0 and solve

$$\frac{\partial}{\partial w_j} \text{RSS}(\mathbf{w}) = -2\rho_j + 2w_j = 0$$
$$\hat{w}_j = \rho_j$$

# Coordinate descent for least squares regression

180

Initialize  $\hat{\mathbf{w}} = 0$  (or smartly...)

while not converged

for  $j=0,1,\dots,D$

compute:  $\rho_j = \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\hat{\mathbf{w}}_{-j}))$

set:  $\hat{w}_j = \rho_j$

residual  
*without feature j*

$\uparrow$

$\hat{y}_i(\hat{\mathbf{w}}_{-j})$

$\uparrow$

*prediction without feature j*

Measure of the correlation between  $w_j$  and the residual without this feature.

# How to access convergence

181

Initialize  $\hat{\mathbf{w}} = 0$  (or smartly...)

while not converged

for  $j=0,1,\dots,D$

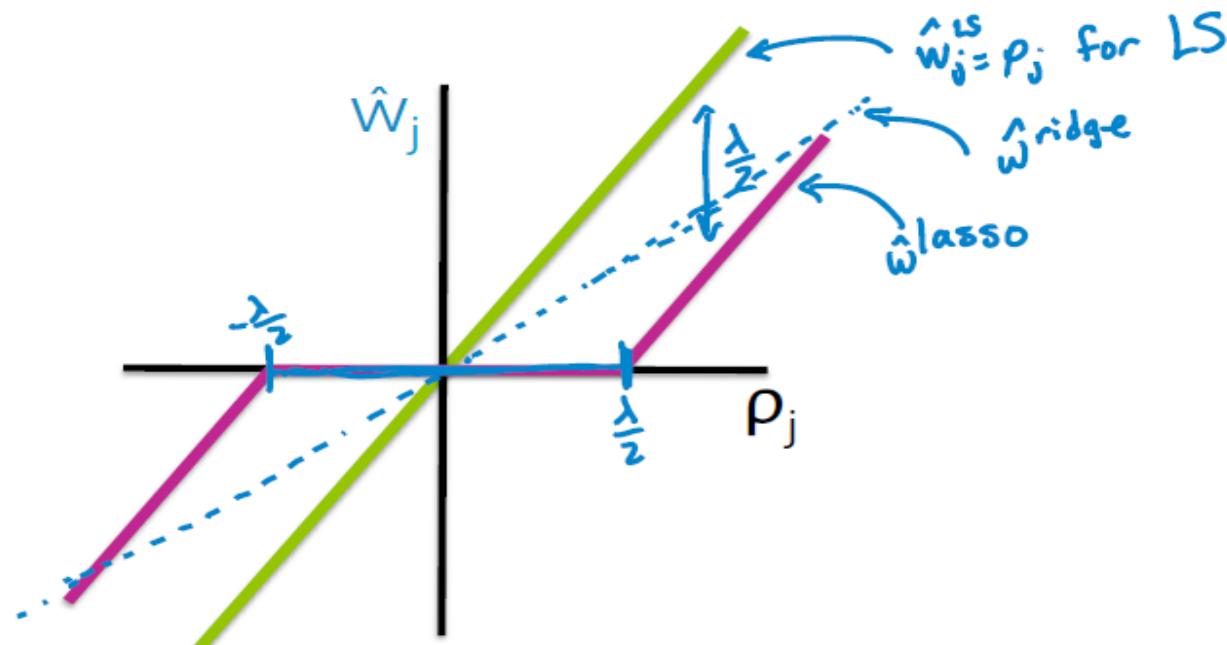
compute:  $\rho_j = \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\hat{\mathbf{w}}_{-j}))$

set:  $\hat{\mathbf{w}}_j = \begin{cases} \rho_j + \lambda/2 & \text{if } \rho_j < -\lambda/2 \\ 0 & \text{if } \rho_j \text{ in } [-\lambda/2, \lambda/2] \\ \rho_j - \lambda/2 & \text{if } \rho_j > \lambda/2 \end{cases}$

# Soft thresholding

182

$$\hat{w}_j = \begin{cases} \rho_j + \lambda/2 & \text{if } \rho_j < -\lambda/2 \\ 0 & \text{if } \rho_j \text{ in } [-\lambda/2, \lambda/2] \\ \rho_j - \lambda/2 & \text{if } \rho_j > \lambda/2 \end{cases}$$



# Convergence criteria

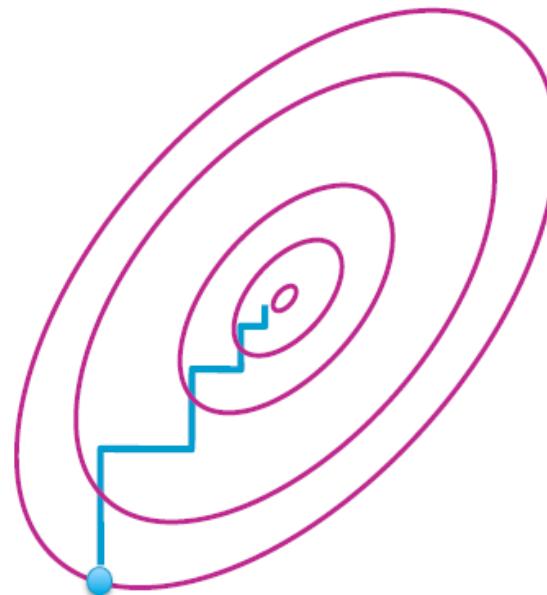
183

When to stop?

For convex problems, will start to take smaller and smaller steps

Measure size of steps taken in a full loop over all features

- stop when  $\max \text{ step} < \epsilon$



# Other lasso solvers

184

Classically: Least angle regression ([LARS](#)) [Efron et al. '04]

Then: [Coordinate descent](#) algorithm

[Fu '98, Friedman, Hastie, & Tibshirani '08]

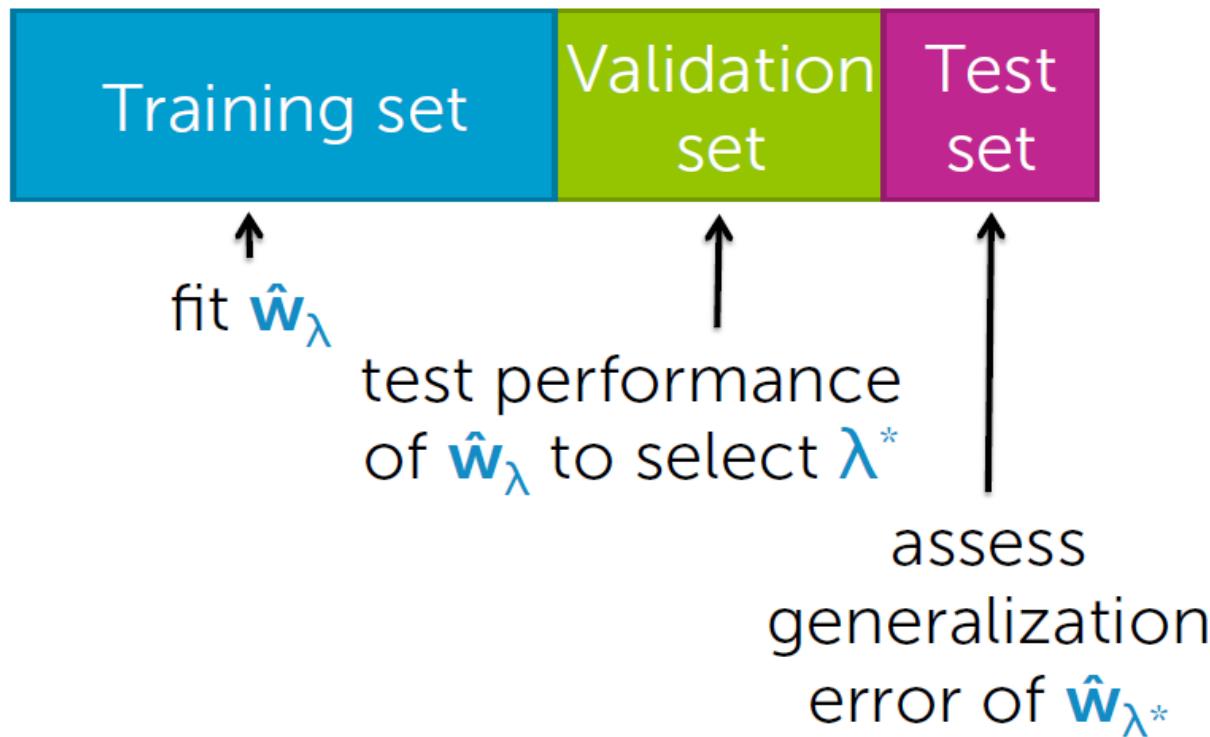
Now:

- Parallel CD (e.g., Shotgun, [\[Bradley et al. '11\]](#))
- Other parallel learning approaches for linear models
  - Parallel stochastic gradient descent ([SGD](#)) (e.g., Hogwild! [\[Niu et al. '11\]](#))
  - Parallel independent solutions then [averaging](#) [\[Zhang et al. '12\]](#)
- Alternating directions method of multipliers ([ADMM](#)) [\[Boyd et al. '11\]](#)

# How do we chose $\lambda$

185

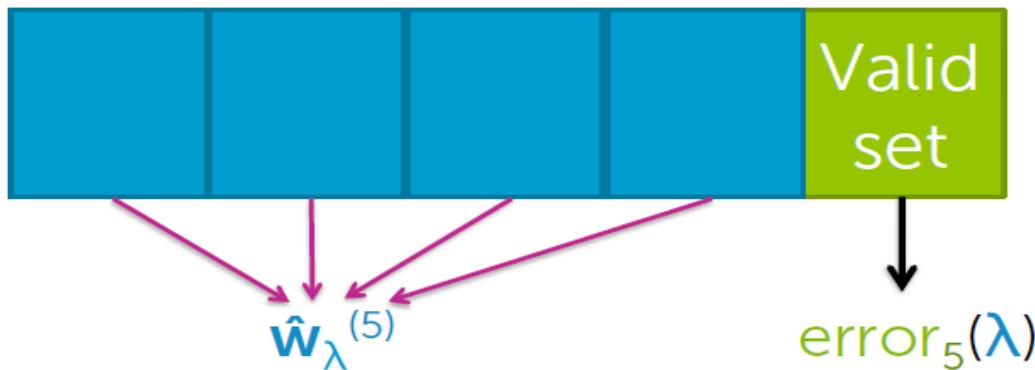
If sufficient amount of data...



# How do we chose $\lambda$

186

## K-fold cross validation



For  $k=1, \dots, K$

1. Estimate  $\hat{w}_\lambda^{(k)}$  on the training blocks
2. Compute error on validation block:  $\text{error}_k(\lambda)$

Compute average error:  $\text{CV}(\lambda) = \frac{1}{K} \sum_{k=1}^K \text{error}_k(\lambda)$

# How do we chose $\lambda$

187

## Choosing $\lambda$ via cross validation

Cross validation is choosing the  $\lambda$  that provides best predictive accuracy

Tends to favor less sparse solutions, and thus smaller  $\lambda$ , than optimal choice for feature selection

c.f., "Machine Learning: A Probabilistic Perspective", Murphy, 2012 for further discussion

# Impact of feature selection and lasso

188

Lasso has changed machine learning,  
statistics, & electrical engineering

But, for feature selection in general, be **careful about interpreting selected features**

- selection only considers features included
- sensitive to correlations between features
- result depends on algorithm used
- there are theoretical guarantees for lasso under certain conditions

# What you can do now

189

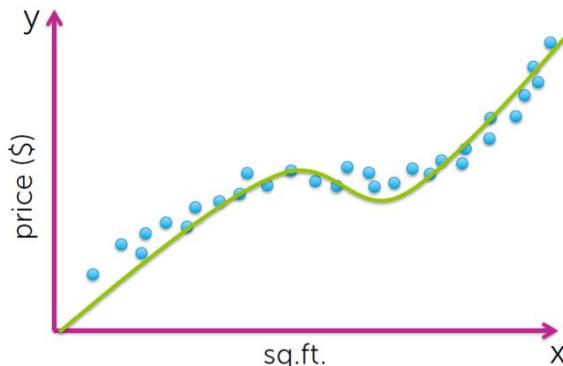
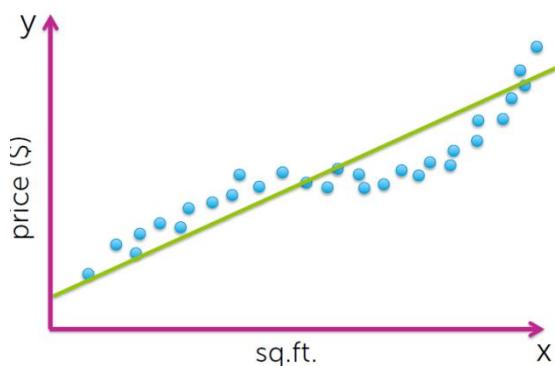
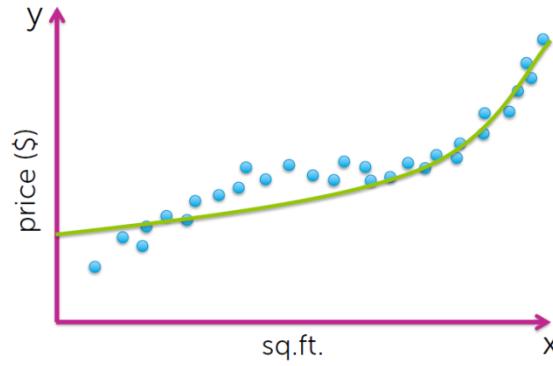
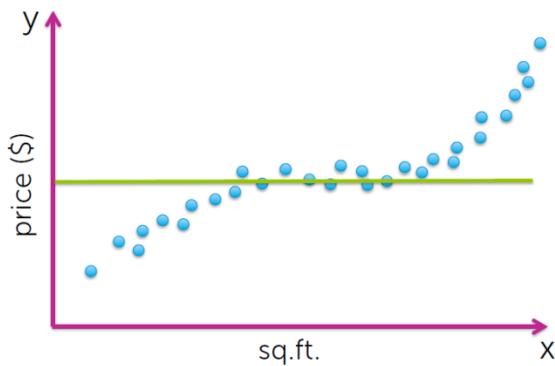
- Perform feature selection using “all subsets” and “forward stepwise” algorithms
- Analyze computational costs of these algorithms
- Contrast greedy and optimal algorithms
- Formulate lasso objective
- Describe what happens to estimated lasso coefficients as tuning parameter  $\lambda$  is varied
- Interpret lasso coefficient path plot
- Contrast ridge and lasso regression
- Describe geometrically why L1 penalty leads to sparsity
- Estimate lasso regression parameters using an iterative coordinate descent algorithm
- Implement K-fold cross validation to select lasso tuning parameter  $\lambda$

# NONPARAMETRIC REGRESSION

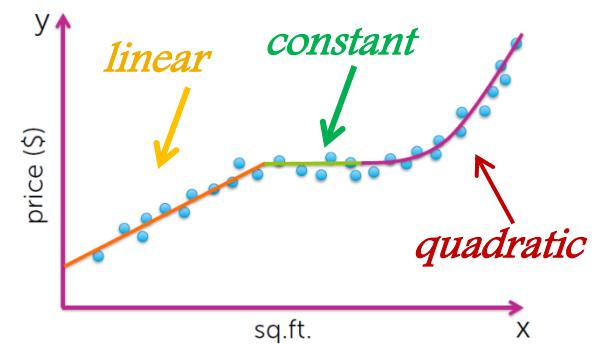
# Fit globally vs fit locally

191

## Parametric models



*Below ...  
 $f(x)$  is not really  
a polynomial function*



# What alternative do we have?

192

If we:

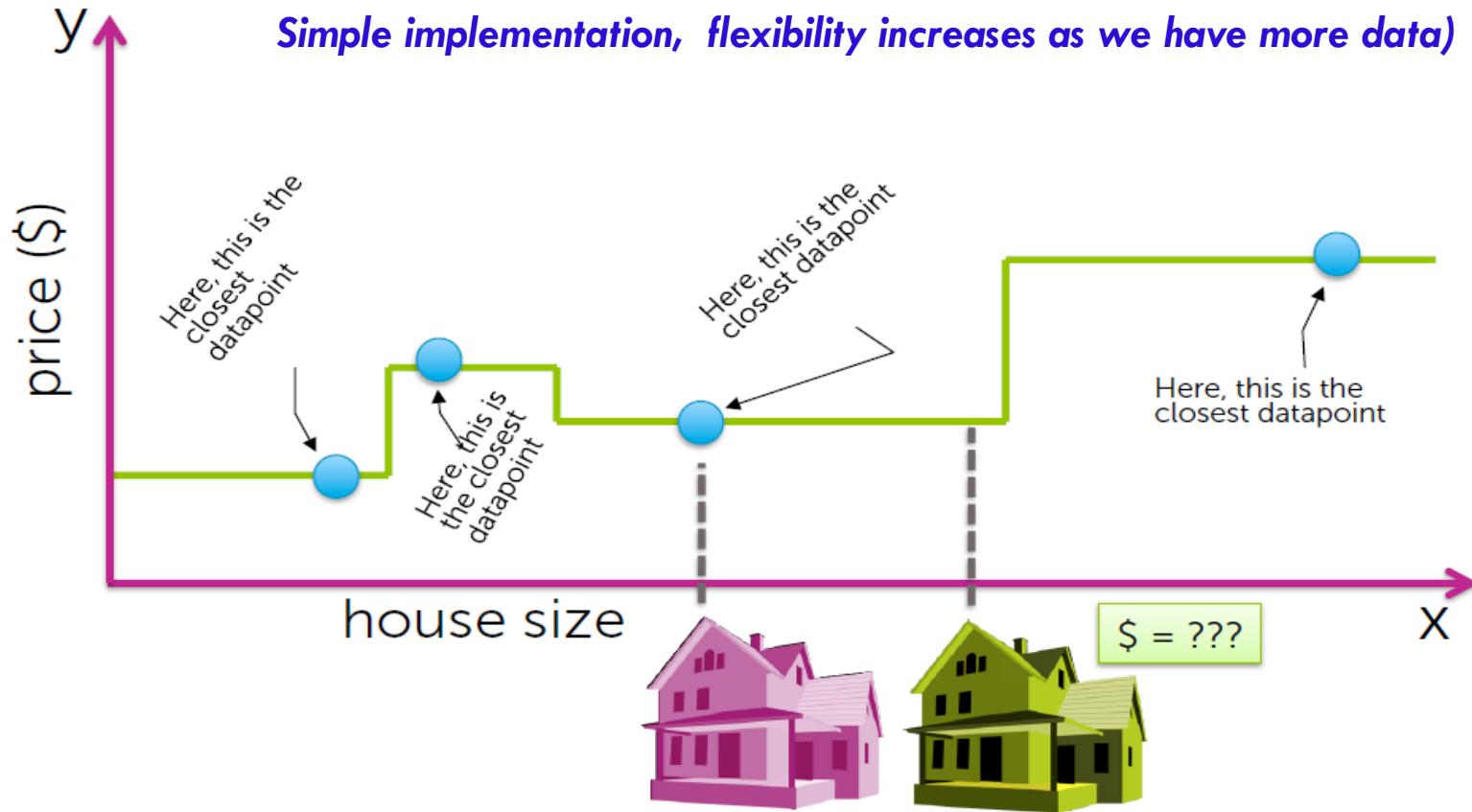
- Want to allow flexibility in  $f(\mathbf{x})$  having local structure
- Don't want to infer "structural breaks"

What's a simple option we have?

- Assuming we have plenty of data...

# Nearest Neighbor & Kernel Regression (nonparametric approach)

193

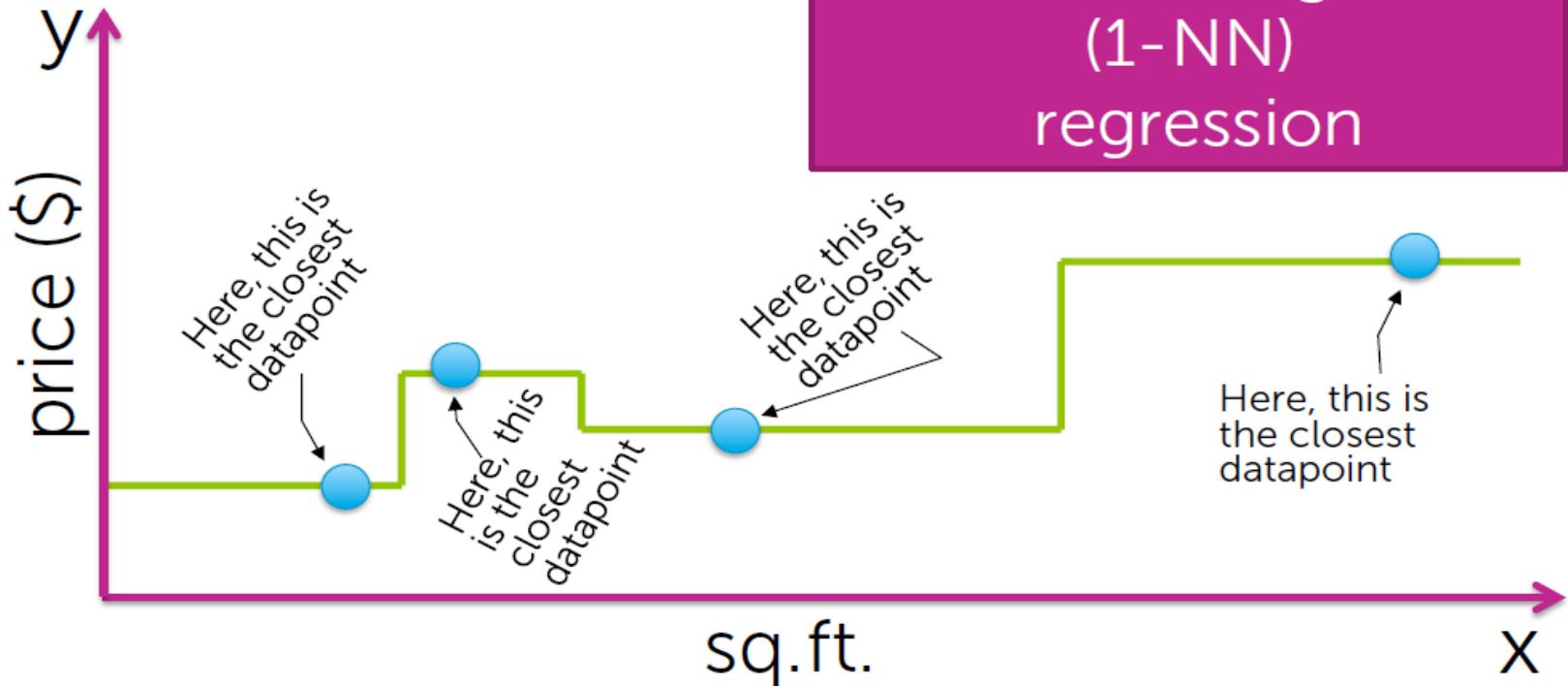


# Fit locally to each data point

194

Predicted value = “closest”  $y_i$

1 nearest neighbor  
(1-NN)  
regression



# What people do naturally...

195

Real estate agent assesses value by finding sale of most similar house



# 1-NN regression more formally

196

Dataset of (, \$) pairs:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

Query point:  $\mathbf{x}_q \leftarrow$  \$ ?  
big lime green house

1. Find "closest"  $\mathbf{x}_i$  in dataset

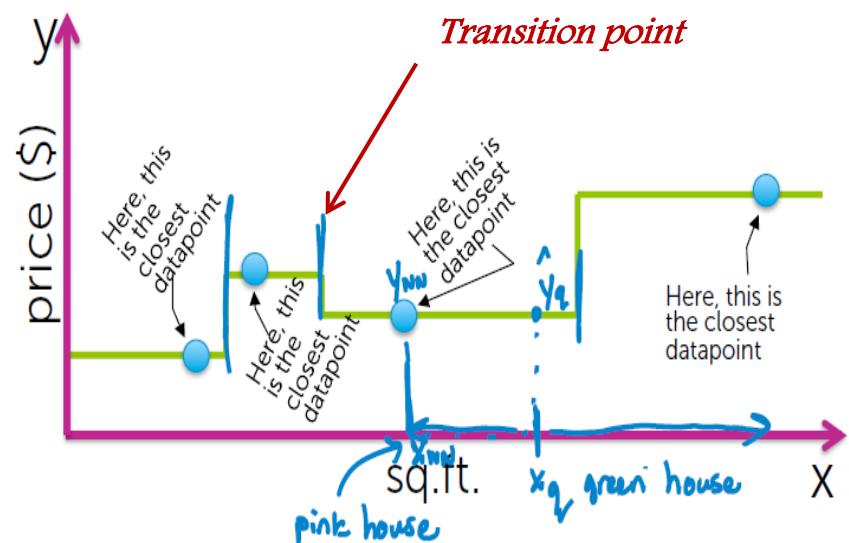
$$x_{NN} \leftarrow \min_i \underline{\text{distance}}(\mathbf{x}_i, \mathbf{x}_q)$$

big pink house

2. Predict

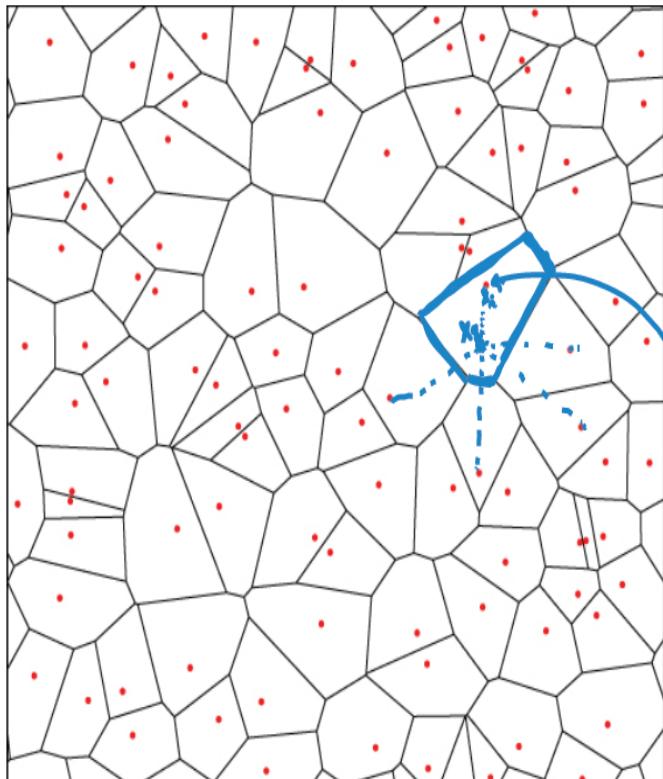
$$\hat{y}_q = y_{NN}$$

sales price of big pink house



# Visualizing 1-NN in multiple dimensions

197



Voronoi tessellation  
(or diagram):

- Divide space into  $N$  regions, each containing 1 datapoint
- Defined such that any  $x$  in region is “closest” to region’s datapoint

$x_i$  closer to  $x_i$   
than any other  
 $x_j$  for  $j \neq i$ .

Don't explicitly form!

# Distance metrics: Notion of „closest”

198

In 1D, just Euclidean distance:

$$\text{distance}(x_j, x_q) = |x_j - x_q|$$

In multiple dimensions:

- can define many interesting distance functions
- most straightforwardly, might want to weight different dimensions differently

# Weighting housing inputs

199

Some inputs are more relevant than others



**# bedrooms**  
**# bathrooms**  
**sq.ft. living**  
sq.ft. lot  
floors  
**year built**  
year renovated  
**waterfront**



# Scaled Euclidean distance

200

Formally, this is achieved via

$$\text{distance}(\mathbf{x}_j, \mathbf{x}_q) = \sqrt{a_1(\mathbf{x}_j[1]-\mathbf{x}_q[1])^2 + \dots + a_d(\mathbf{x}_j[d]-\mathbf{x}_q[d])^2}$$

weight on each input  
(defining relative importance)

Other example distance metrics:

- Mahalanobis, rank-based, correlation-based, cosine similarity, Manhattan, Hamming, ...

# Different distance metrics

201



# Performing 1-NN search

202

- Query house:



- Dataset:



- **Specify:** Distance metric
- **Output:** Most similar house



# 1-NN algorithm

203

Initialize **Dist2NN** =  $\infty$ ,  =  $\emptyset$

For  $i=1,2,\dots,N$

Compute:  $\delta = \text{distance}(\text{house}_i, \text{query house})$

If  $\delta < \text{Dist2NN}$

set



set **Dist2NN** =  $\delta$

Return most similar house

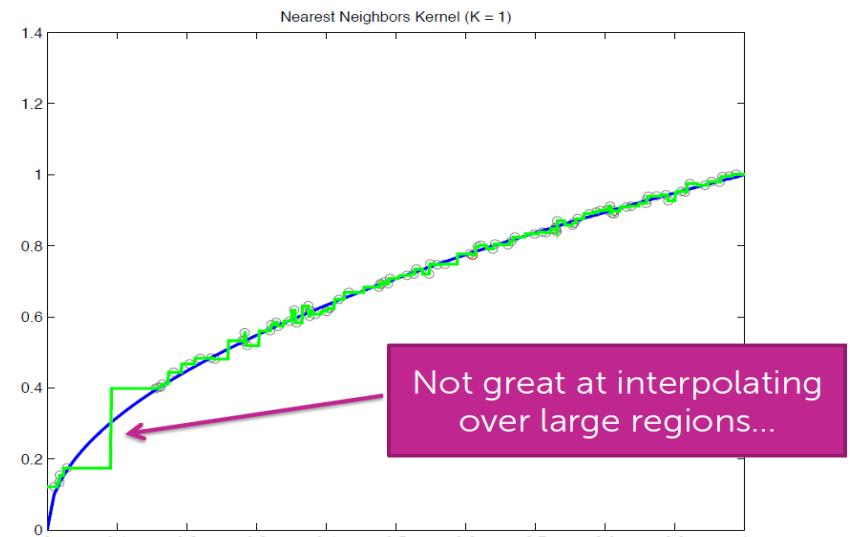
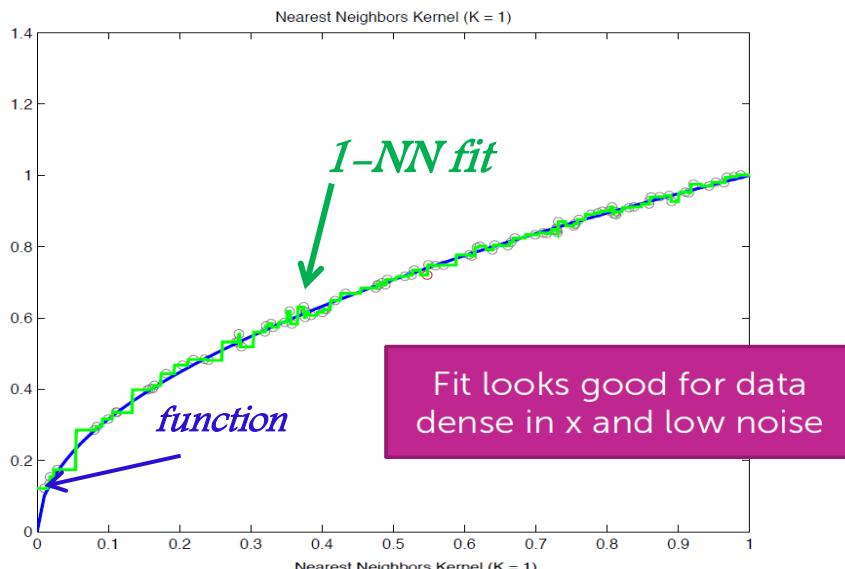


closest house  
to query house



# 1-NN in practice

204

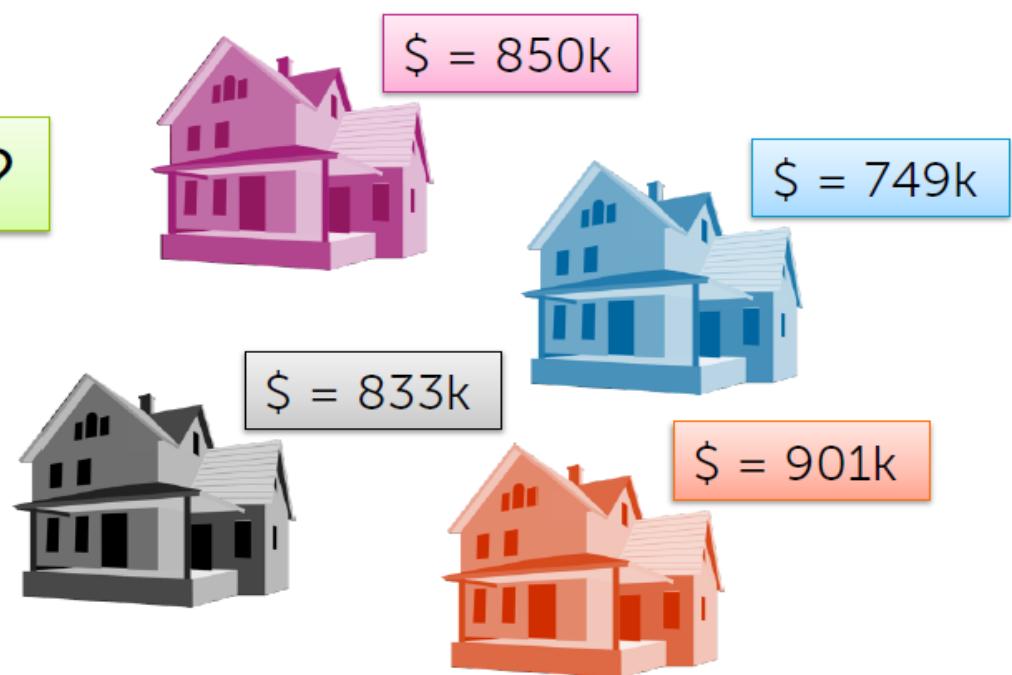


1-NN sensitive to noise in the data

# Get more „comps”

205

More reliable estimate if you base estimate off of a larger set of comparable homes



# K-NN regression more formally

206

Dataset of (, \$) pairs:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

Query point:  $\mathbf{x}_q$

1. Find  $k$  closest  $\mathbf{x}_i$  in dataset

$(x_{NN_1}, x_{NN_2}, \dots, x_{NN_k})$  such that for any  $x_i$  not in nearest neighbor set,  
 $distance(x_i, x_q) \geq distance(x_{NN_k}, x_q)$

2. Predict

$$\begin{aligned}\hat{y}_q &= \frac{1}{k} (y_{NN_1} + y_{NN_2} + \dots + y_{NN_k}) \\ &= \frac{1}{k} \sum_{j=1}^k y_{NN_j}\end{aligned}$$

# K-NN more formally

207

- Query house:



- Dataset:



- **Specify:** Distance metric
- **Output:** Most similar houses



# K-NN algorithm

208

sort first **k houses**  
by distance to query house

Initialize **Dist2kNN** = sort( $\delta_1, \dots, \delta_k$ )  $\leftarrow$  list of sorted distances  
= SO   $\dots$ ,   $_1$    $_k$   $\leftarrow$  list of sorted houses

For  $i = k+1, \dots, N$

Compute:  $\delta = \text{distance}(\text{house}_i, \text{house}_q)$  

If  $\delta < \text{Dist2kNN}[k]$

find  $j$  such that  $\delta > \text{Dist2kNN}[j-1]$  but  $\delta < \text{Dist2kNN}[j]$

remove furthest house and shift queue:

$[j+1: \text{house}] = [j:k: \text{house}]$

$\text{Dist2kNN}[j+1:k] = \text{Dist2kNN}[j:k-1]$

set  $\text{Dist2kNN}[j] = \delta$  and

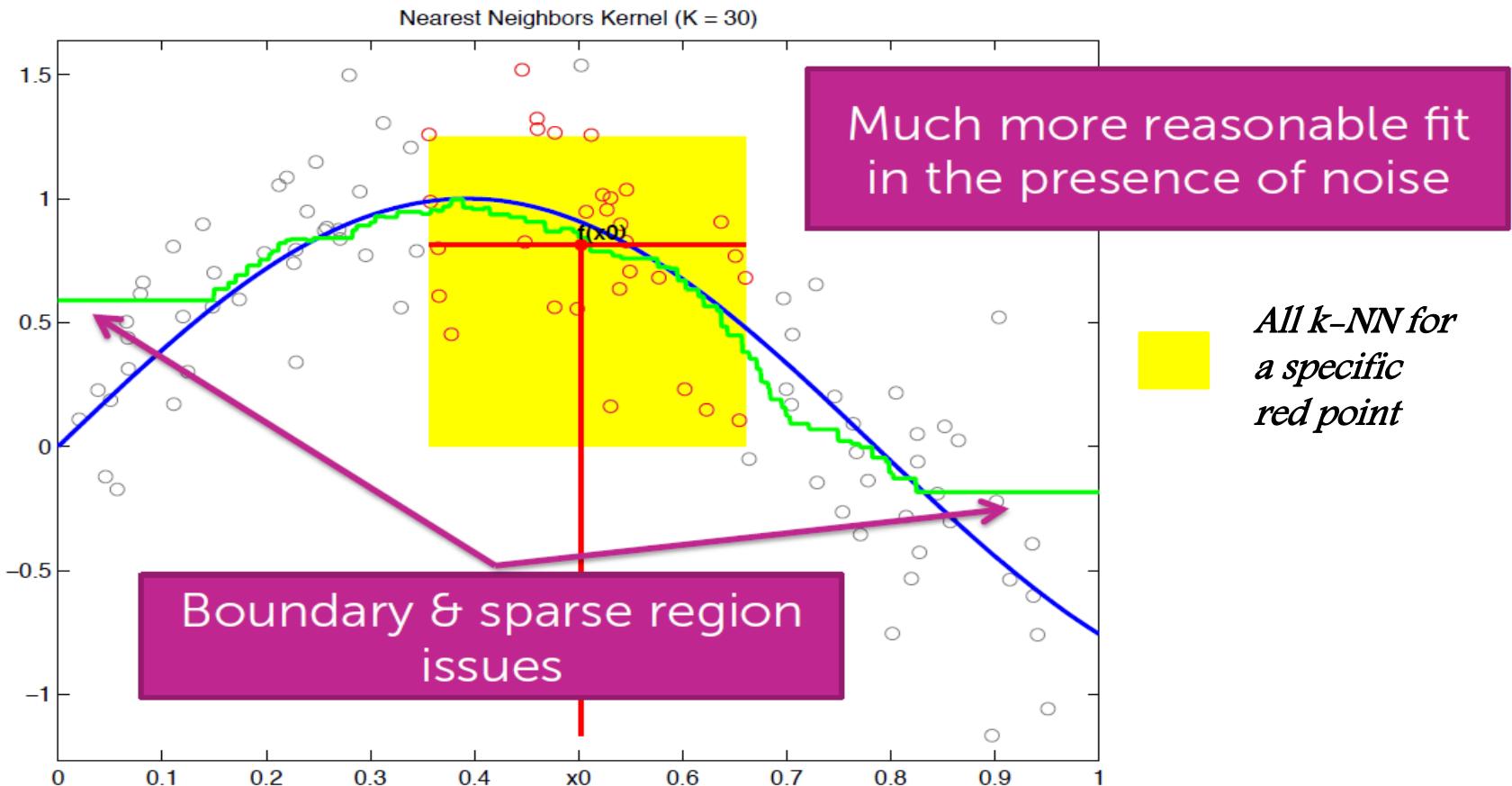


closest houses  
to query house

Return **k** most similar houses

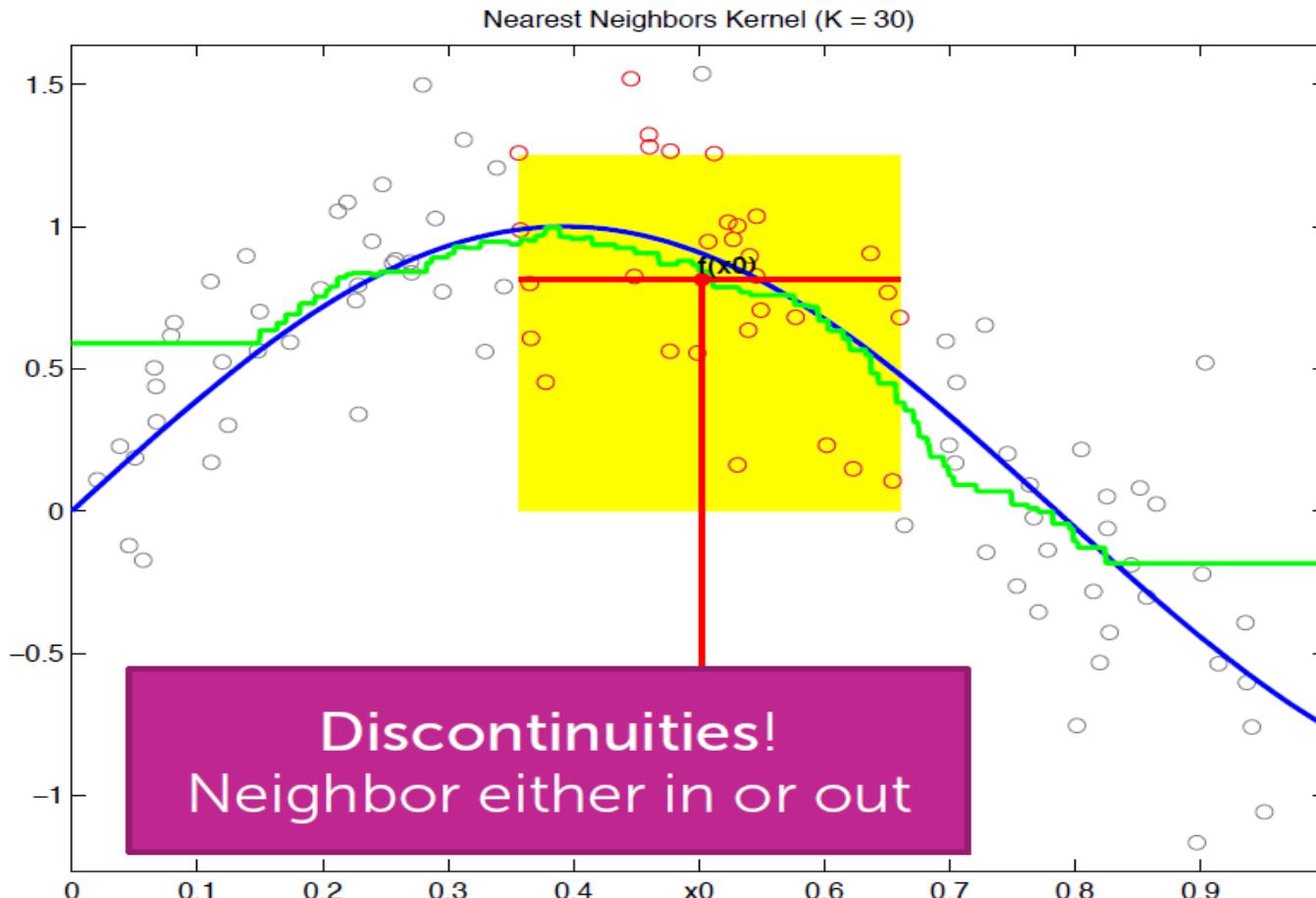
# K-NN in practice

209



# K-NN in practice

210



# Issues with discontinuities

211

Overall predictive accuracy might be okay, but...

For example, in housing application:

- If you are a buyer or seller, this matters
- Can be a jump in estimated value of house going just from 2640 sq.ft. to 2641 sq.ft.
- Don't really believe this type of fit

# Weighted k-NN

212

Weigh more similar houses more than those less similar in list of k-NN

Predict:

weights on NN

$$\hat{y}_q = \frac{c_{qNN1}y_{NN1} + c_{qNN2}y_{NN2} + c_{qNN3}y_{NN3} + \dots + c_{qNNk}y_{NNk}}{\sum_{j=1}^k c_{qNNj}}$$

# How to define weights

213

Want weight  $c_{qNNj}$  to be small when  
distance( $\mathbf{x}_{NNj}, \mathbf{x}_q$ ) large

and  $c_{qNNj}$  to be large when  
distance( $\mathbf{x}_{NNj}, \mathbf{x}_q$ ) small

Simple method :

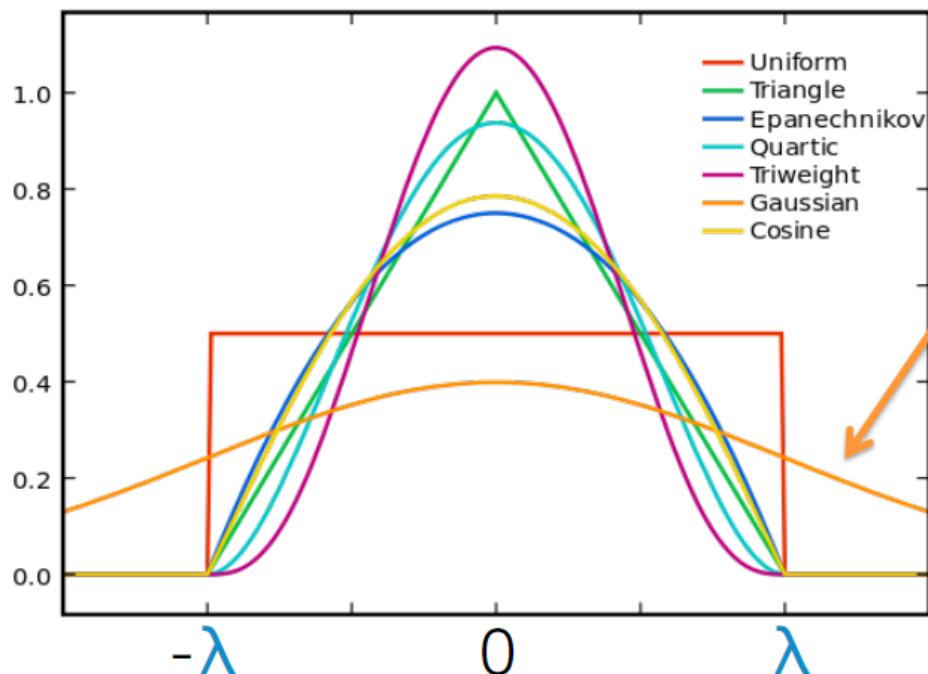
$$c_{q,NNj} = \frac{1}{\text{distance}(\mathbf{x}_j, \mathbf{x}_q)}$$

# Kernel weights for $d=1$

214

Define:  $c_{qNNj} = \text{Kernel}_\lambda(|x_{NNj} - x_q|)$

simple isotropic case



Gaussian kernel:  
 $\text{Kernel}_\lambda(|x_i - x_q|) = \exp(-(x_i - x_q)^2/\lambda)$

Note: never exactly 0!

Kernel drives how the weights will decay, if at all, as a function of the distance.

# Kernel regression

Nadaraya-Watson  
kernel weighted average

215

Instead of just weighting NN, weight *all* points

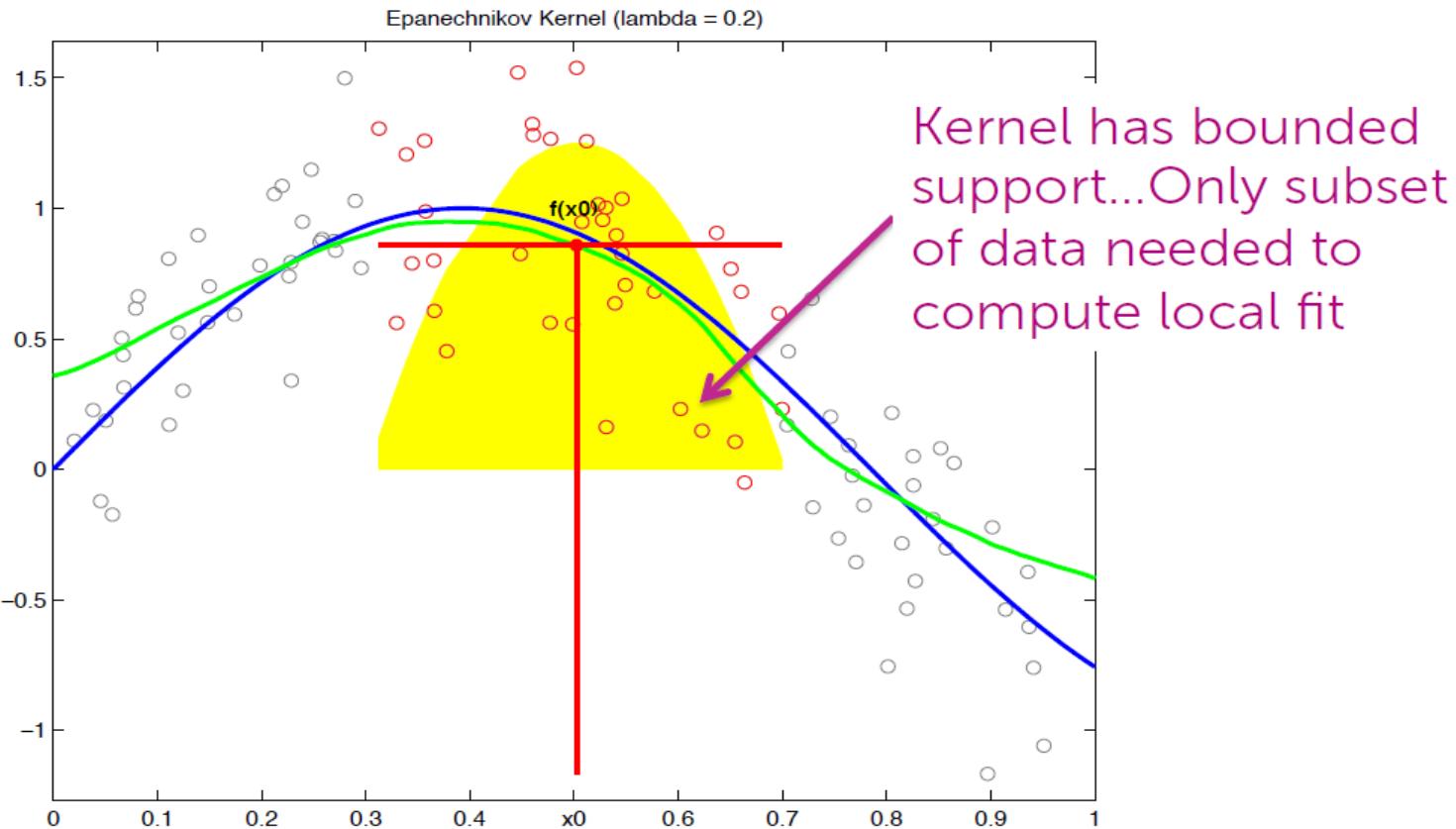
Predict:

weight on each datapoint

$$\hat{y}_q = \frac{\sum_{i=1}^N c_{qi} y_i}{\sum_{i=1}^N c_{qi}} = \frac{\sum_{i=1}^N \text{Kernel}_{\lambda}(\text{distance}(\mathbf{x}_i, \mathbf{x}_q)) * y_i}{\sum_{i=1}^N \text{Kernel}_{\lambda}(\text{distance}(\mathbf{x}_i, \mathbf{x}_q))}$$

# Kernel regression in practice

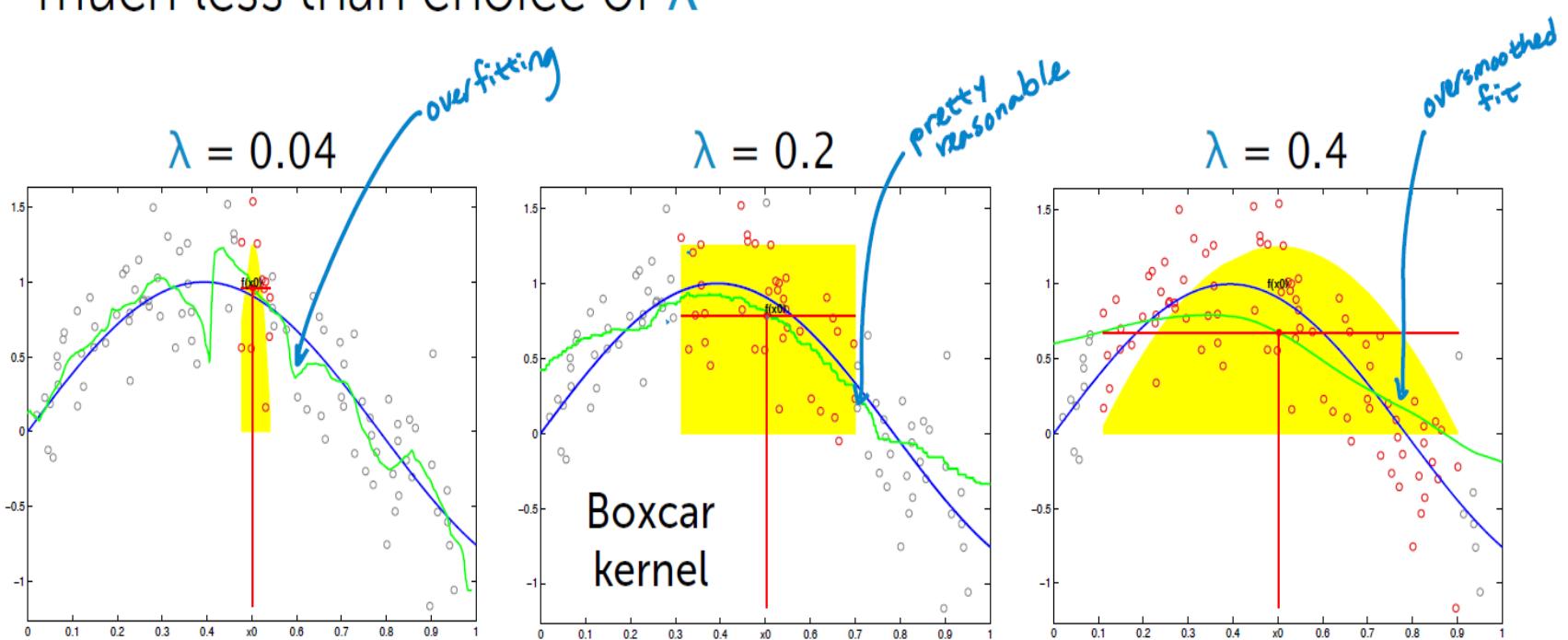
216



# Choice of bandwidth $\lambda$

217

Often, choice of kernel matters  
much less than choice of  $\lambda$



# Choosing $\lambda$ (or $k$ on k-NN)

218

How to choose? Same story as always...

Cross Validation

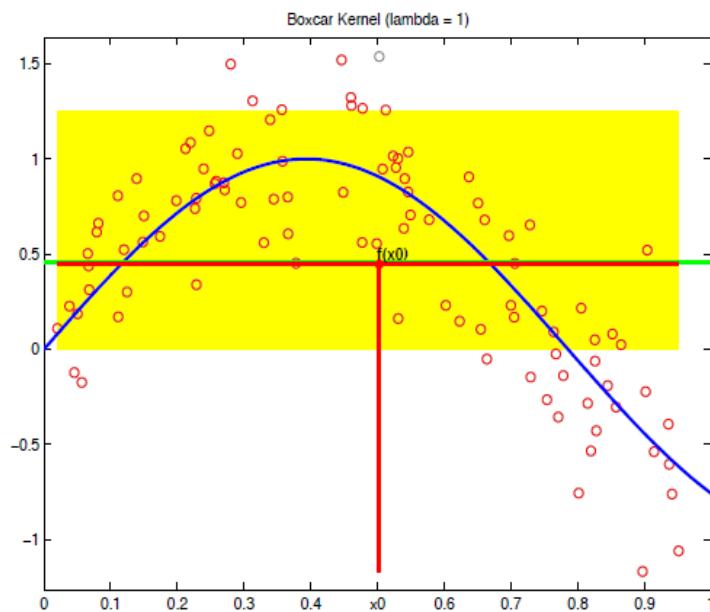
# Contrasting with global average

219

A globally constant fit weights all points equally

$$\hat{y}_q = \frac{1}{N} \sum_{i=1}^N y_i = \frac{\sum_{i=1}^N c y_i}{\sum_{i=1}^N c}$$

equal weight on each datapoint



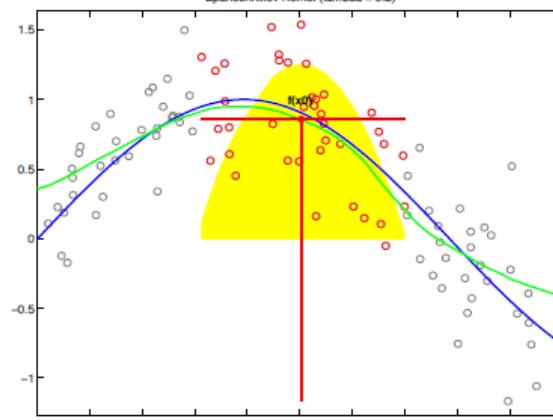
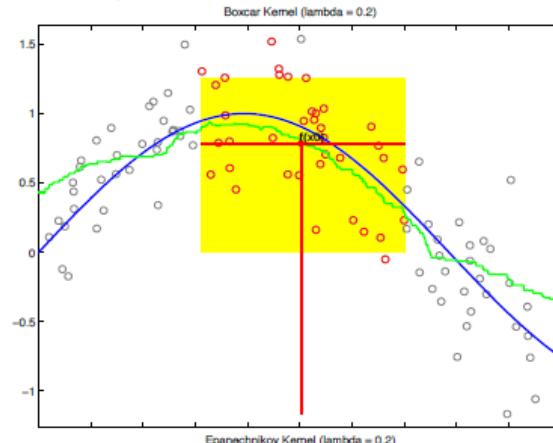
# Contrasting with global average

220

Kernel regression leads to **locally constant fit**

- slowly add in some points and  
and let others gradually die off

$$\hat{y}_q = \frac{\sum_{i=1}^N \text{Kernel}_\lambda(\text{distance}(x_i, x_q)) * y_i}{\sum_{i=1}^N \text{Kernel}_\lambda(\text{distance}(x_i, x_q))}$$



# Local linear regression

221

So far, discussed fitting constant function locally at each point

→ “locally weighted averages”

Can instead fit a line or polynomial locally at each point

→ “locally weighted linear regression”

# Local regression rules of thumb

222

- Local linear fit reduces bias at boundaries with minimum increase in variance
- Local quadratic fit doesn't help at boundaries and increases variance, but does help capture curvature in the interior
- With sufficient data, local polynomials of odd degree dominate those of even degree

Recommended default choice:  
**local linear regression**

# Nonparametric approaches

223

k-NN and kernel regression are examples of **nonparametric** regression

General goals of nonparametrics:

- Flexibility
- Make few assumptions about  $f(\mathbf{x})$
- Complexity can grow with the number of observations  $N$

Lots of other choices:

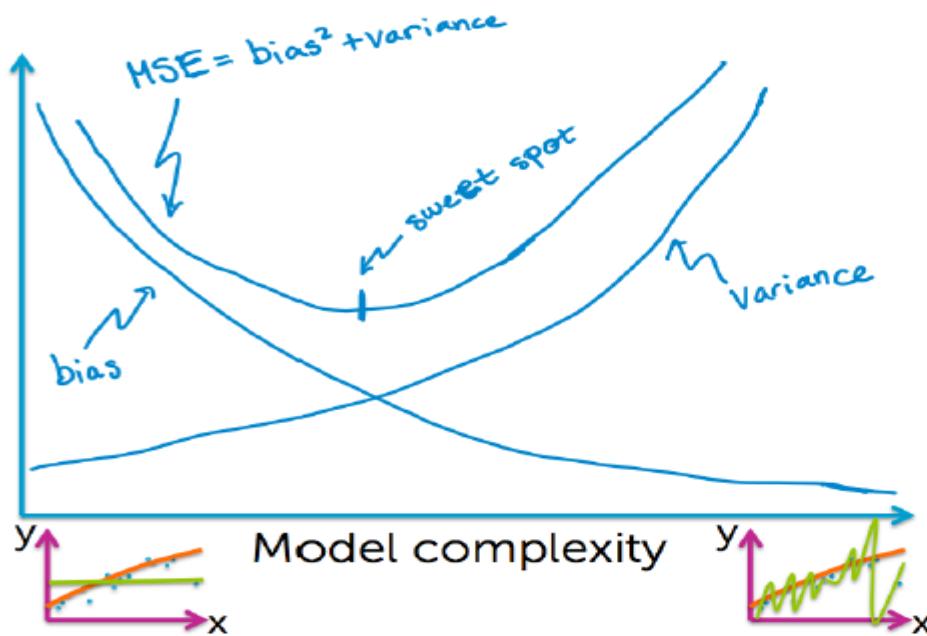
- Splines, trees, locally weighted structured regression models...

# Limiting behaviour of NN

224

## Noiseless setting ( $\varepsilon_i = 0$ )

In the limit of getting an infinite amount of noiseless data, the MSE of 1-NN fit goes to 0

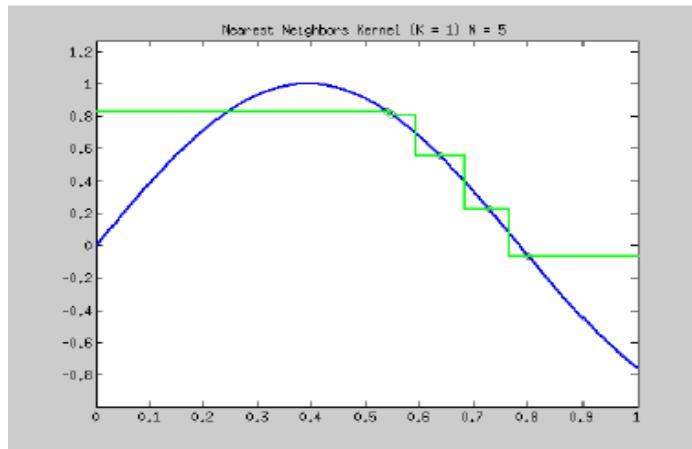


# Limiting behaviour of NN

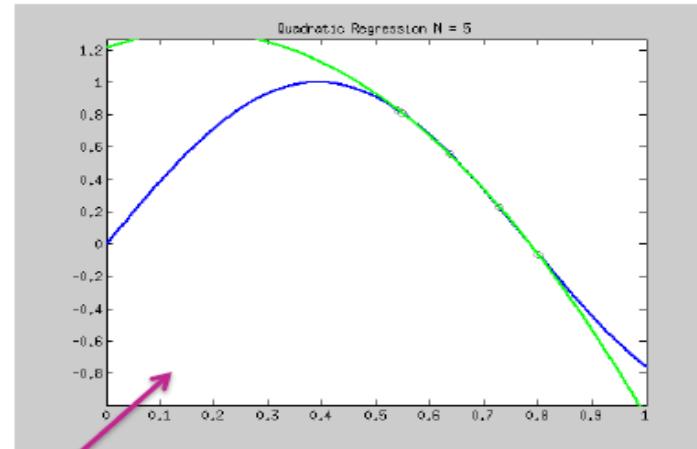
225

## Noiseless setting ( $\varepsilon_i = 0$ )

In the limit of getting an infinite amount of noiseless data, the MSE of 1-NN fit goes to 0



1-NN fit

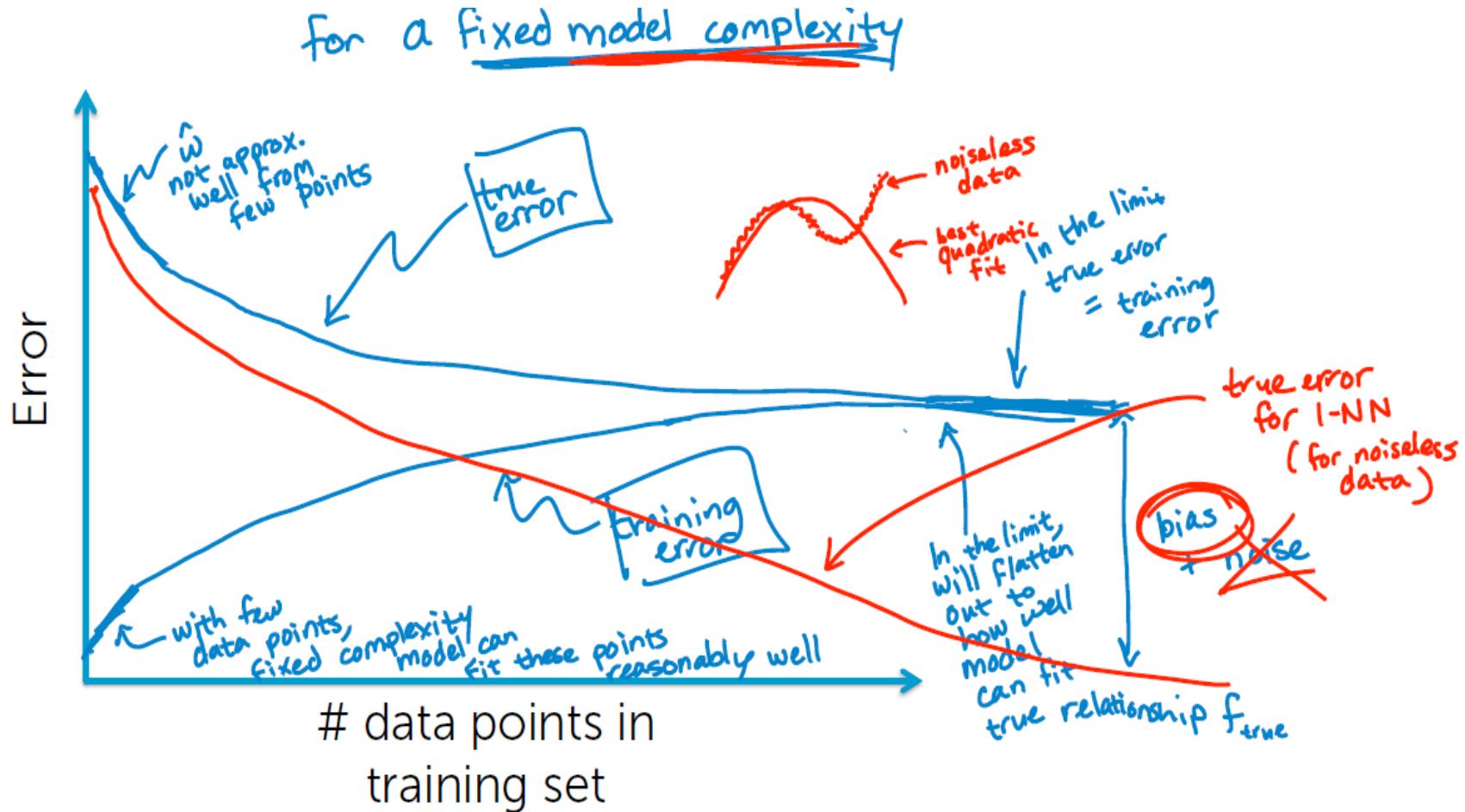


Quadratic fit

Not true for parametric models!

# Error vs amount of data

226

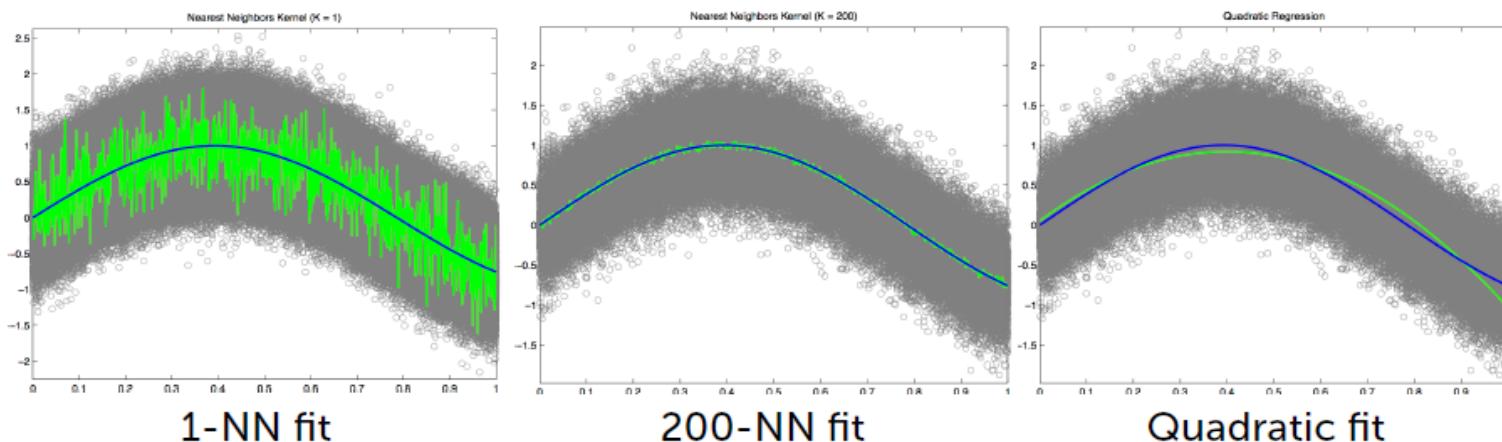


# Limiting behaviour of NN

227

## Noisy data setting

In the limit of getting an infinite amount of data,  
the MSE of NN fit goes to 0 if  $k$  grows, too



# Issues: NN and kernel methods

228

NN and kernel methods work well when the data cover the space, but...

- the more dimensions  $d$  you have, the more points  $N$  you need to cover the space
- need  $N = O(\exp(d))$  data points for good performance

This is where parametric models become useful...

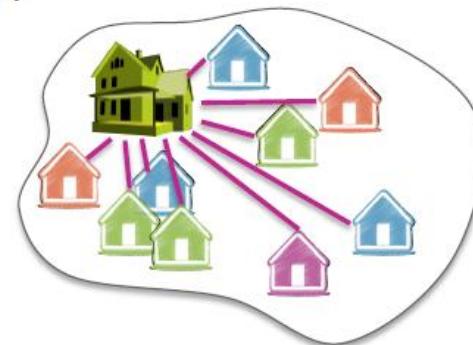
# Issues: Complexity of NN search

229

Naïve approach: Brute force search

- Given a query point  $\mathbf{x}_q$
- Scan through each point  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$
- $O(N)$  distance computations per 1-NN query!
- $O(N \log k)$  per k-NN query!

What if N is huge???  
(and many queries)



Will talk more about efficient methods in  
[Clustering & Retrieval](#) course

# We have discussed how to

230

- Motivate the use of nearest neighbor (NN) regression
- Define distance metrics in 1D and multiple dimensions
- Perform NN and k-NN regression
- Analyze computational costs of these algorithms
- Discuss sensitivity of NN to lack of data, dimensionality, and noise
- Perform weighted k-NN and define weights using a kernel
- Define and implement kernel regression
- Describe the effect of varying the kernel bandwidth  $\lambda$  or # of nearest neighbors  $k$
- Select  $\lambda$  or  $k$  using cross validation
- Compare and contrast kernel regression with a global average fit
- Define what makes an approach nonparametric and why NN and kernel regression are considered nonparametric methods
- Analyze the limiting behavior of NN regression

# Summarising

231

## Models

- Linear regression
- Regularization: Ridge (L2), Lasso (L1)
- Nearest neighbor and kernel regression

## Algorithms

- Gradient descent
- Coordinate descent

## Concepts

- Loss functions, bias-variance tradeoff, cross-validation, sparsity, overfitting, model selection, feature selection