# DATA SCIENCE WITH MACHINE LEARNING: REGRESSION

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

22/12 2020

WFAiS UJ, Informatyka Stosowana
I stopień studiów

# What is Data Science?

**Is mainly about extracting knowledge from data (terms "data mining" or "Knowledge Discovery in Databases" are highly related). It can be about analyzing trends, building predictive models, … etc.**

**Is an agglomerate of <span style="color:red">data collection, data modeling and analysis</span>, a decision making, and everything you need to know to accomplish your goals. Eventually, it boils down to the following fields/skills:**

- **Computer science:**

**Algorithms, programming (patterns, languages etc.), understanding hardware & operating systems, high-performance computing'**

- **Mathematical aspects:**

**Linear algebra, differential equations for optimization problems, statistics**

- **Few others:**

**<span style="color:red">Machine learning</span>, domain knowledge, and data visualization & communication skills**

# Data Science and Machine Learning?

**Machine learning** **algorithms are algorithms that learn (often predictive) models from data. I.e., instead of formulating "rules" manually, a machine learning algorithm will learn the model for you.**
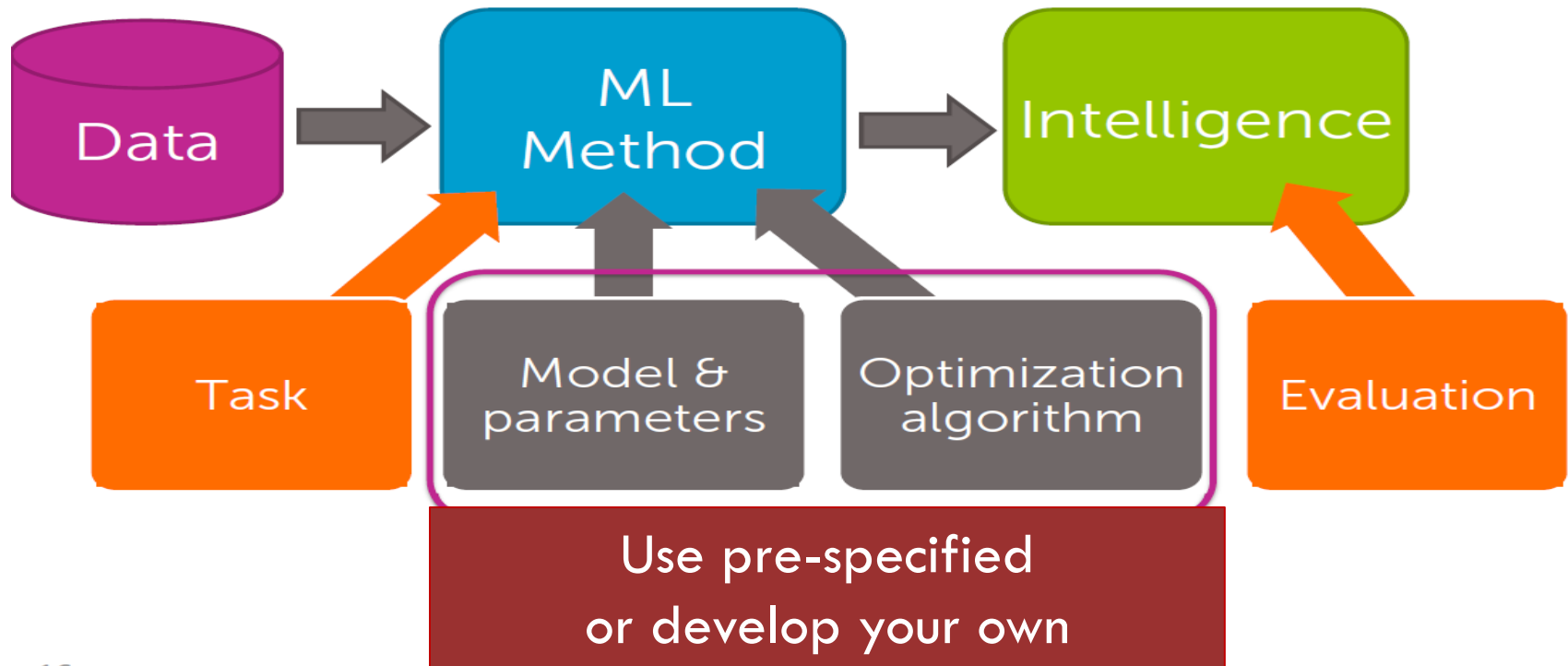
**Machine learning** **– at its core – is about the use and development of these learning algorithms.** **Data science** **is more about the extraction of knowledge from data to answer particular question or solve particular problems.**

**Machine learning is often a big part of a "data science" project****, e.g., it is often heavily used for exploratory analysis and discovery (clustering algorithms) and building predictive models (supervised learning algorithms). However, in** **data science****, you often also worry about the collection, wrangling, and cleaning of your data (i.e., data engineering), and eventually, you want to draw conclusions from your data that helps you solve a particular problem.**
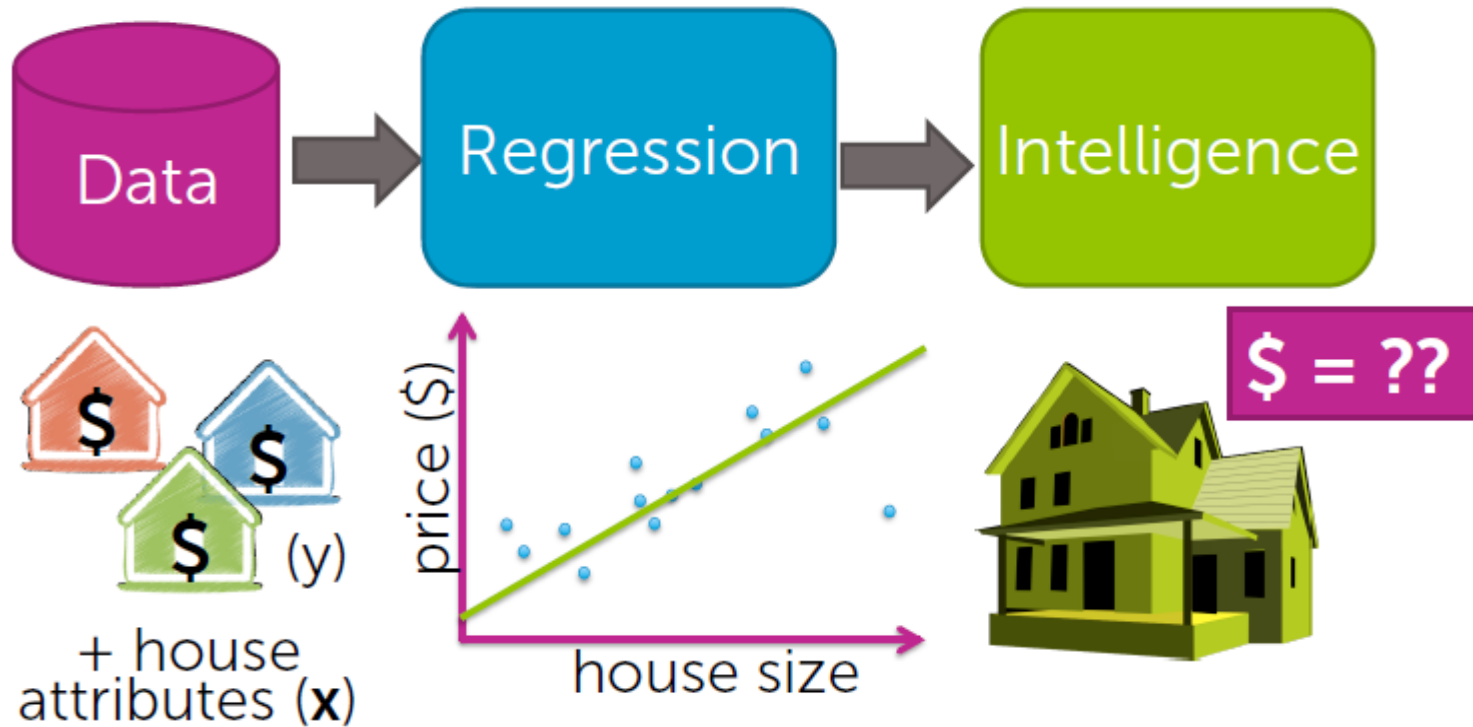
22/12 2020

# Deploing inteligence module

**Case studied are about building, evaluating, deploying inteligence in data analysis.**



22/12 2020

# Case study

Predicting house prices

# Prediction: Predicting house prices

**Models**
- Linear regression
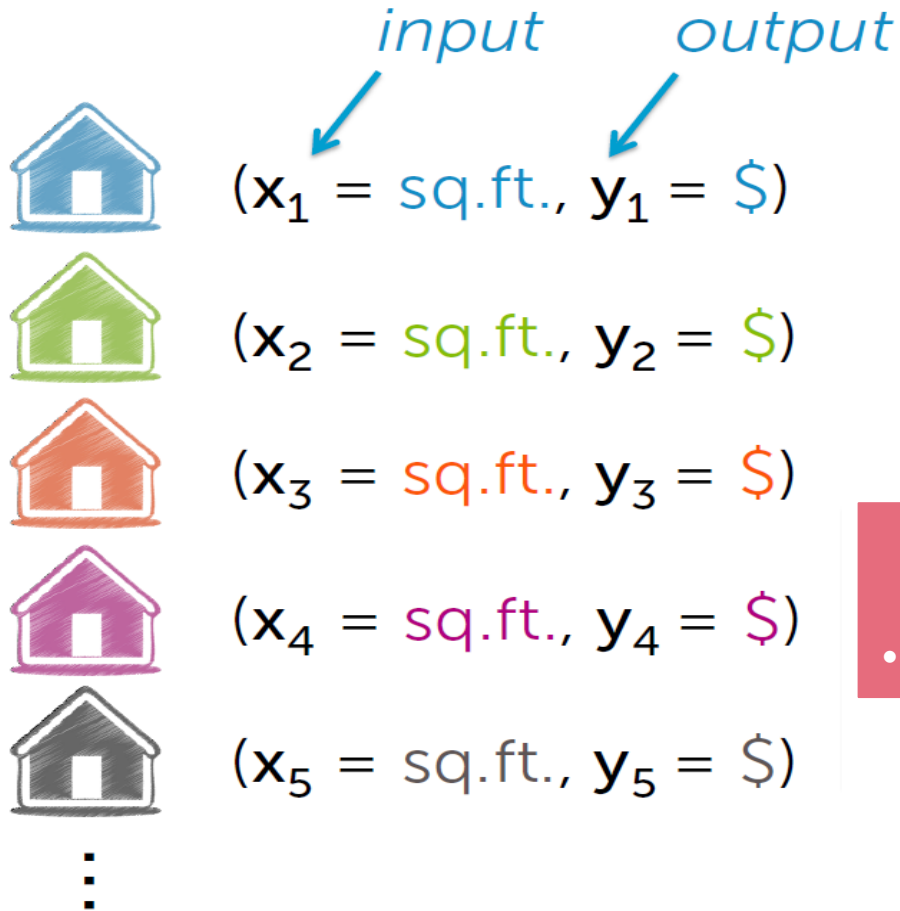- Regularization:
  Ridge (L2), Lasso (L1)

**Algorithms**
- Gradient descent
- Coordinate descent

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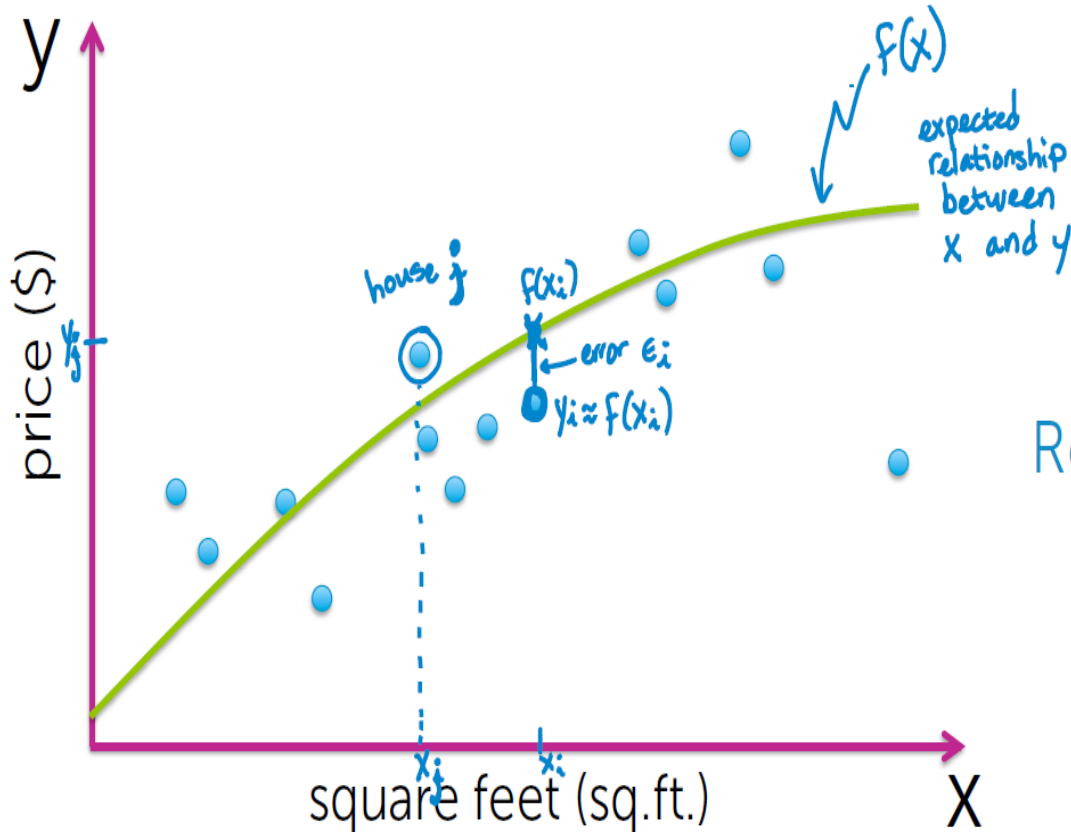- Loss functions, bias-variance tradeoff, cross-validation, sparsity, overfitting, model selection

13/10/2020

# 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 quantity of interest**
- **assume y can be predicted from x**

22/12 2020

# Model: assume functional relationship

"Essentially, all models are wrong but some are usefull." George Box, 1987.

y

price ($)

f(x)

expected relationship between x and y

house j

f($x_i$)

error $\epsilon_i$

$y_i \approx f(x_i)$

$x_j$    $x_i$

square feet (sq.ft.)

X

Regression model:

$$y_i = f(x_i) + \epsilon_i$$
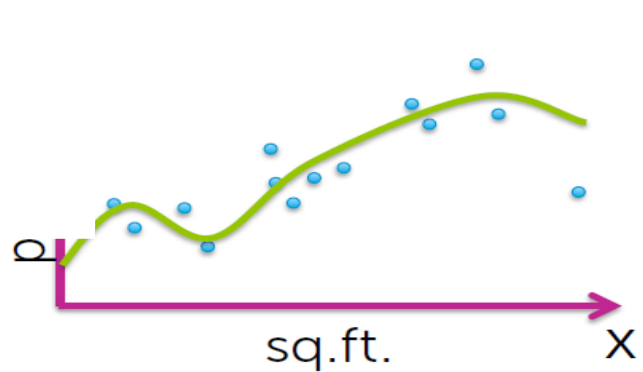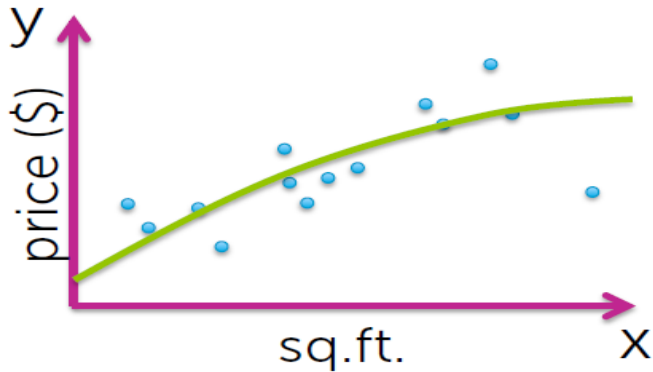
$$E[\epsilon_i] = 0$$ ← equally likely that error is + or −

expected value

⇓

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

22/12 2020

# Task 1:

**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

22/12 2020

# SIMPLE LINEAR REGRESSION

22/12 2020

# Simple linear regression model

$$y_i = w_0 + w_1 x_i + \varepsilon_i$$

$$f(x) = w_0 + w_1 x$$

price ($)

square feet (sq.ft.)

# The cost of using a given line

Residual sum of squares (RSS)



$$RSS(\underline{w}_0, \underline{w}_1) =$$
$$(\$_{house\ 1} - [w_0 + w_1 sq.ft._{house\ 1}])^2$$
$$+ (\$_{house\ 2} - [w_0 + w_1 sq.ft._{house\ 2}])^2$$
$$+ (\$_{house\ 3} - [w_0 + w_1 sq.ft._{house\ 3}])^2$$
$$+ ...[include\ all\ training\ houses]$$

22/12 2020

# Find „best" line

Minimize cost over all possible $w_0, w_1$

$RSS(w_0=1.1, w_1=0.8) = \#_3$
$= \#_2$

$RSS(w_0=0.98, w_1=0.87)$

$RSS(w_0=0.97, w_1=0.85) = \#_1$

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

y

price ($)

square feet (sq.ft.)     x

22/12 2020

# Interpreting the coefficients

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x$$

y

price ($)

Predicted $
of house with
sq.ft.=0
(just land)

square feet (sq.ft.)          X

$\hat{y} = \hat{w}_0$
   when $x = 0$

not very meaningful

22/12 2020

# Interpreting the coefficients

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x$$

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



price ($) vs square feet (sq.ft.)

predicted change in $

1 sq. ft.

$$\hat{\$}_{1001 \text{ sq.ft.}} - \hat{\$}_{1000 \text{ sq.ft.}}$$

$$= \hat{\cancel{w_0}} + \hat{w}_1 \cdot 1001 \text{ sq.ft.}$$

$$- (\hat{\cancel{w_0}} + \hat{w}_1 1000 \text{ sq.ft.})$$

$$= \hat{w}_1$$

predicted change in the output per unit change in input

22/12 2020

# ML algorithm: minimasing 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 $= q(w_0, w_1)$

22/12 2020

# 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

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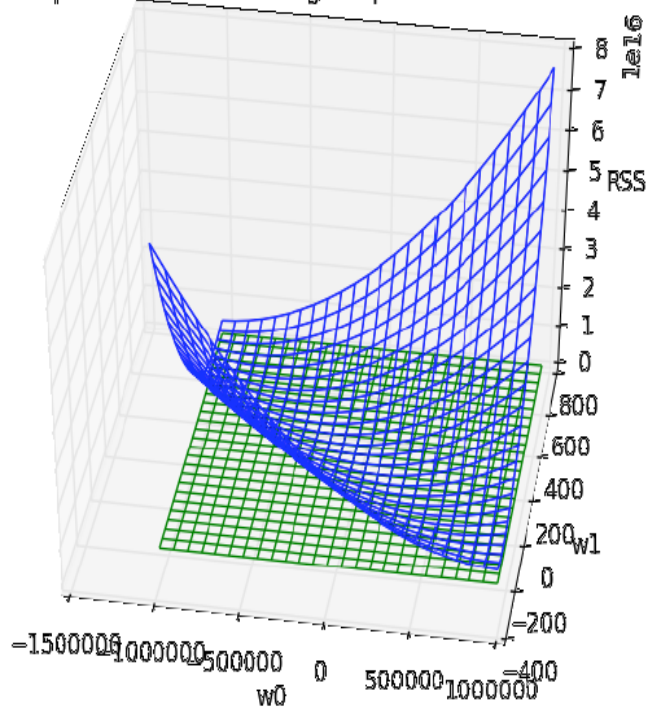value of $\epsilon$ depends on the data we are looking at

Algorithm:

**while** not <u>converged</u>

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

# Moving to multiple dimensions

3D plot of RSS with tangent plane at minimum

$$\nabla g(w) = \begin{bmatrix} \frac{\partial g}{\partial w_0} \\ \frac{\partial g}{\partial w_1} \\ \vdots \\ \frac{\partial g}{\partial w_p} \end{bmatrix}$$
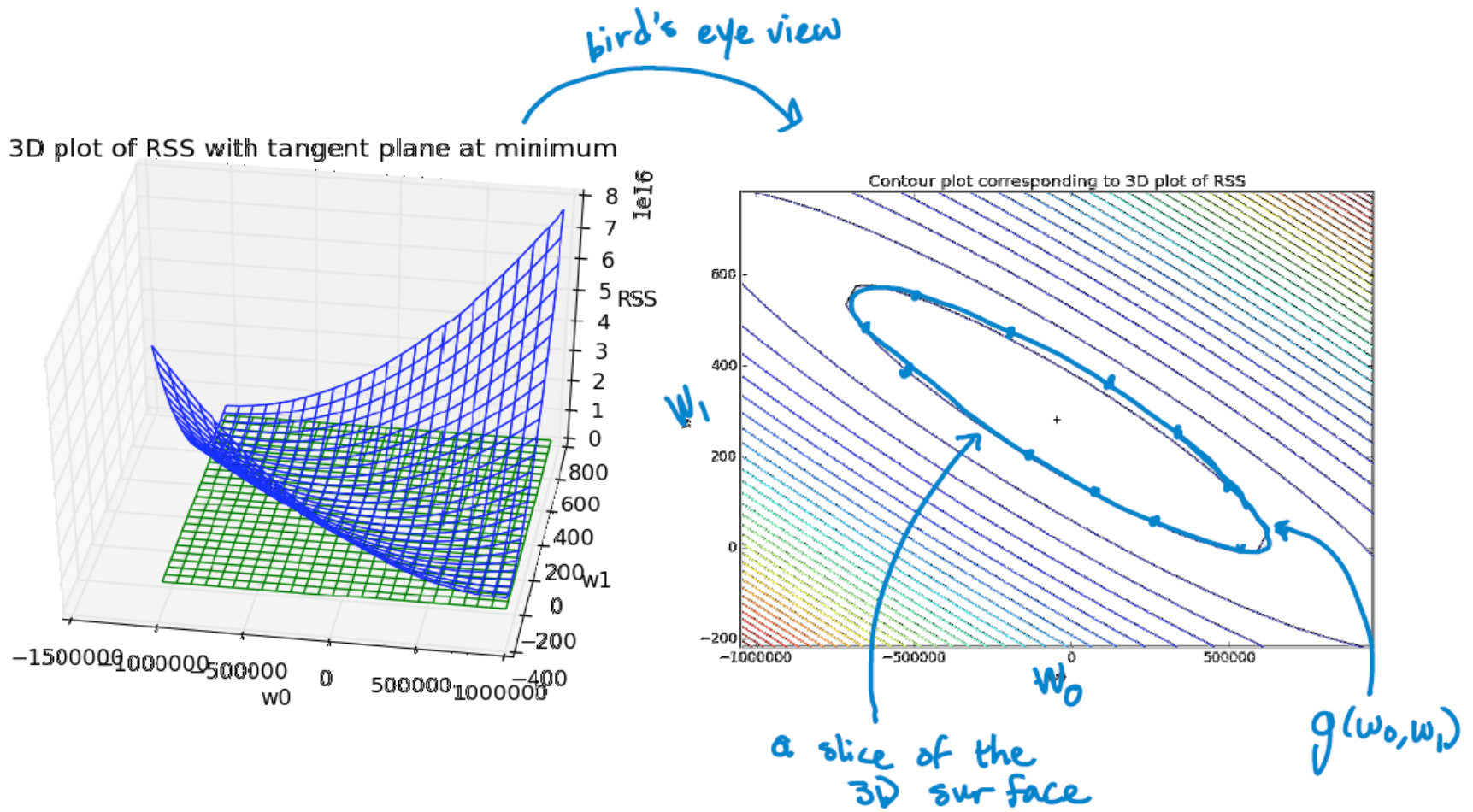
$(p+1)$-dimensional vector

gradient    $[w_0, w_1, ..., w_p]$

partial derivative is like a derivate with respect to $w_1$ treating all other variables as constant

22/12 2020

# Contour plots

bird's eye view

3D plot of RSS with tangent plane at minimum

$w_1$

a slice of the 3D surface

$g(w_0, w_1)$

$w_0$

22/12 2020

# Gradient descent

Contour plot corresponding to 3D plot of RSS
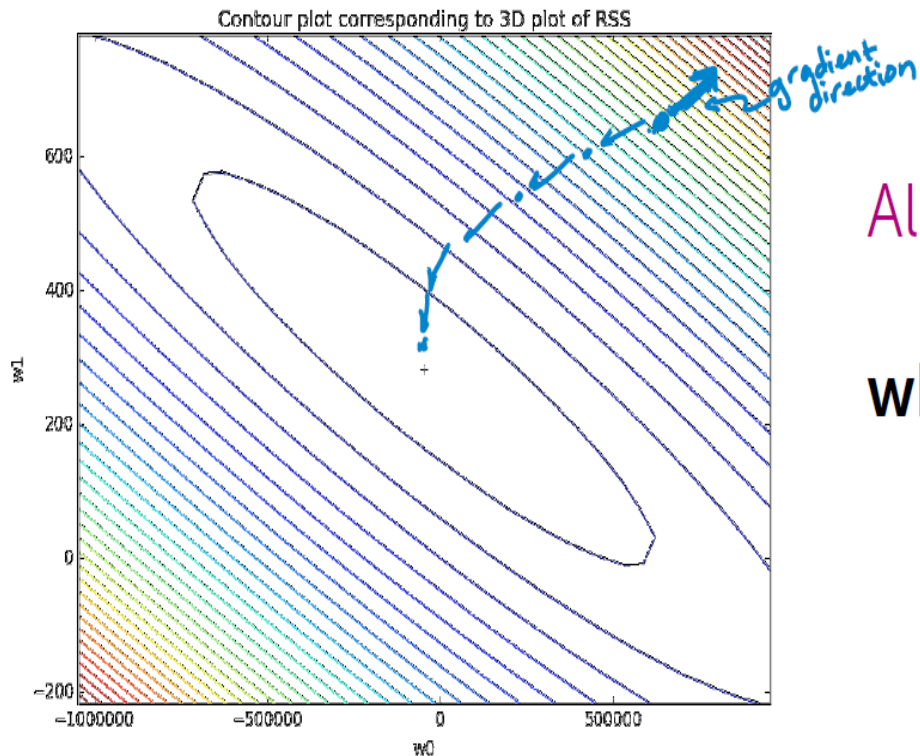
gradient direction

Algorithm:

**while** not converged

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

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

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

22/12 2020

# Compute the gradient

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

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

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

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

Putting it together:

$$\nabla RSS(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^{N} [y_i - (w_0 + w_1 x_i)] \\ -2 \sum_{i=1}^{N} [y_i - (w_0 + w_1 x_i)] x_i \end{bmatrix}$$
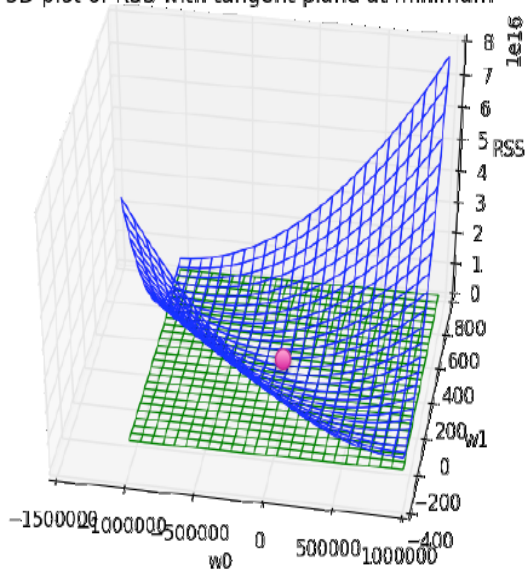
Taking the derivative w.r.t. $w_1$

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

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

22/12 2020

# Approach 1: set gradient to 0

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

*This method is called „Closed form solution"*

3D plot of RSS with tangent plane at minimum



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 sales price

estimate of the slope

average sq.ft.

bottom term:

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

plug in

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

Note:

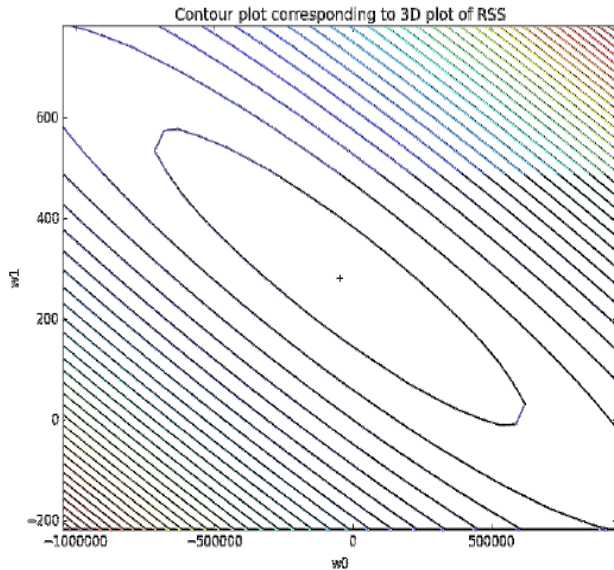$$\sum_{i=1}^{N} y_i$$

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

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

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

22/12 2020

# Approach 2: gradient descent

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



Contour plot corresponding to 3D plot of RSS

while not converged     $(-2) \cdot (-\eta)$

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

If overall, under predicting $\hat{y}_i$, then $\sum [y_i - \hat{y}_i]$ is positive
$\rightarrow W_0$ is going to increase
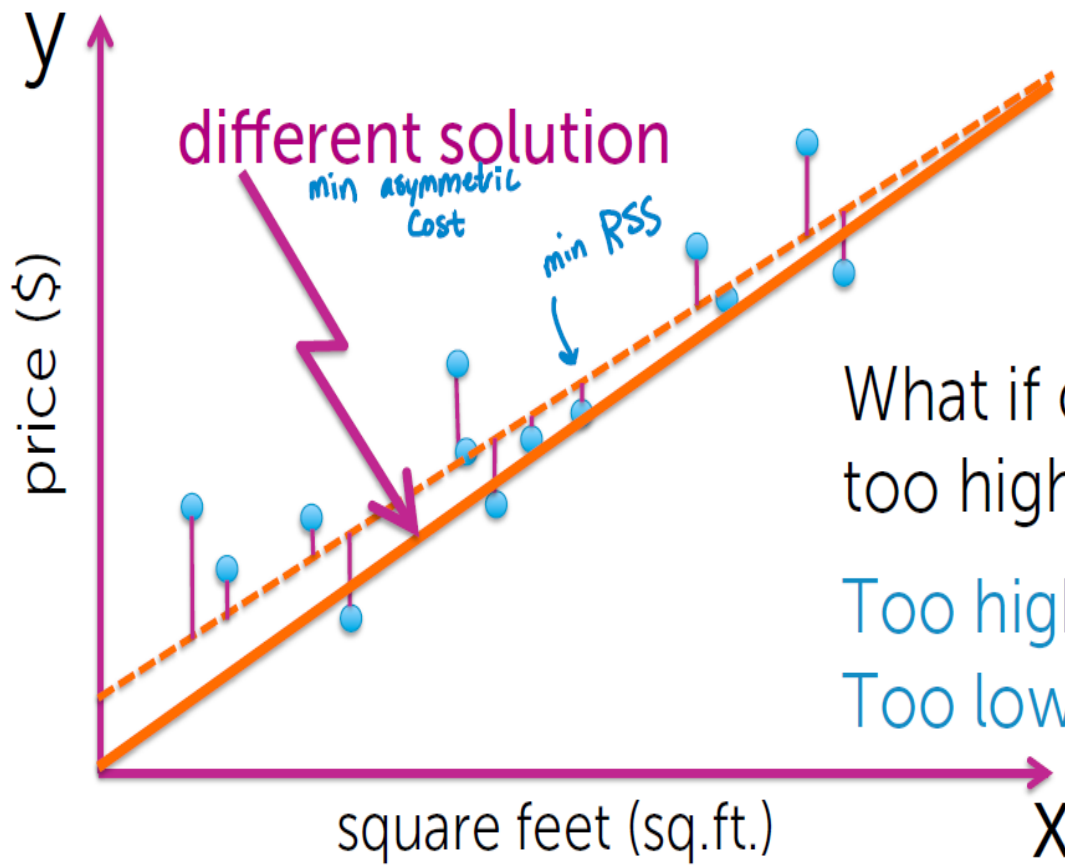similar intuition for $w_1$, but multiply by $x_i$

22/12 2020

# Comparing the approaches

- For most ML problems, cannot solve gradient = 0

- Even if solving gradient = 0 is feasible, gradient descent can be more efficient

- Gradient descent relies on choosing stepsize and convergence criteria

22/12 2020

# Asymmetric cost functions

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

y

price ($)

different solution

min asymmetric Cost

min RSS

square feet (sq.ft.)    x

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

Too high → no offers ($=0)
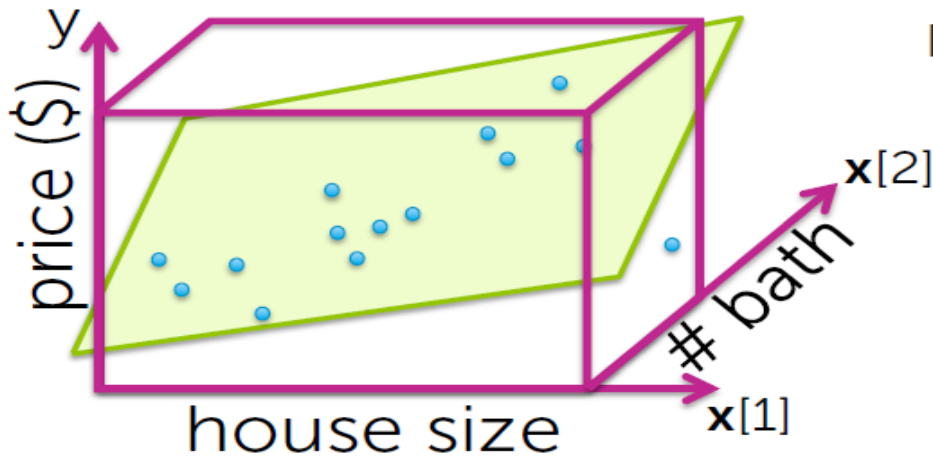Too low → offers for lower $

22/12 2020

# MULTIPLE REGRESSION

22/12 2020

# Multiple regression

Fit **more complex relationships** than just a line

Incorporate more inputs

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

22/12 2020

# Polynomial regression

Model:

$$y_i = w_0 + w_1 x_i + w_2 x_i^2 + \ldots + w_p x_i^p + \varepsilon_i$$

treat as different **features**

feature 1 = 1 (constant)

feature 2 = x

feature 3 = $x^2$

...

feature p+1 = $x^p$

parameter 1 = $w_0$

parameter 2 = $w_1$

parameter 3 = $w_2$

...

parameter p+1 = $w_p$

22/12 2020

# Other functional forms of one input

□ **Trends in time series**



$y_i$ = \$ of $i^{th}$ house sale
$t_i$ = month of $i^{th}$ house sale

House sales recorded monthly

On average, house prices tend to increase with time
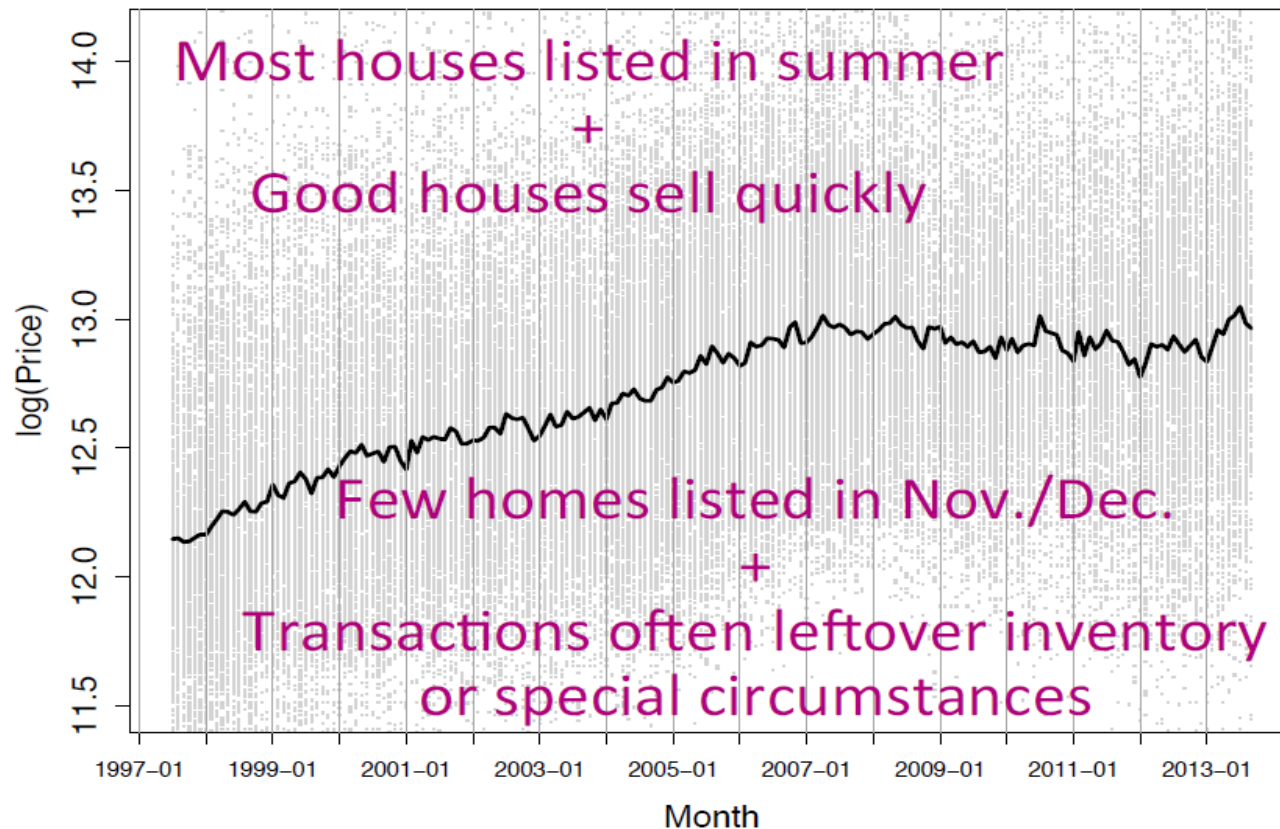
*This trend can be modeled with polynomial function.*

# Other functional forms of one input

□ **Seasonality**
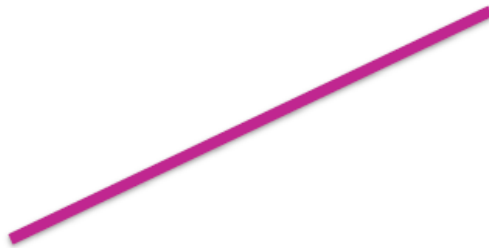
# Example of detrending

Model:

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

Linear trend

Unknown phase/shift

Seasonal component =
Sinusoid with period 12
(resets annually)

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

→ $\sin(2\pi t_i / 12 - \Phi) = \sin(2\pi t_i / 12)\cos(\Phi) - \cos(2\pi t_i / 12)\sin(\Phi)$
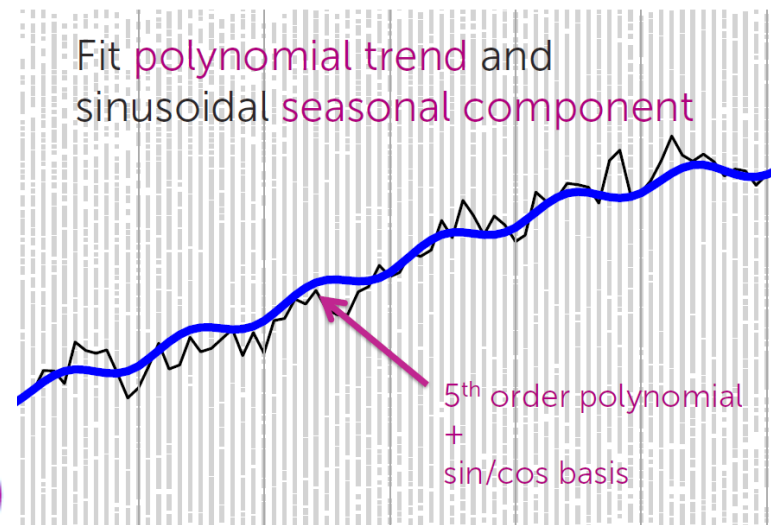
22/12 2020

# Example of detrending

Equivalently,

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

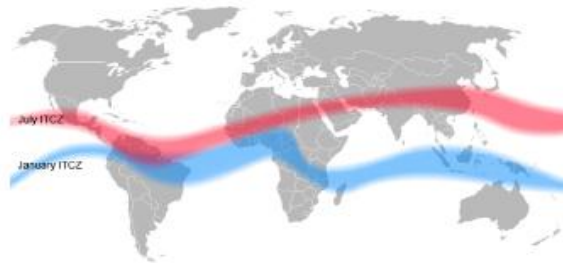*feature 1* = 1 (constant)

*feature 2* = t

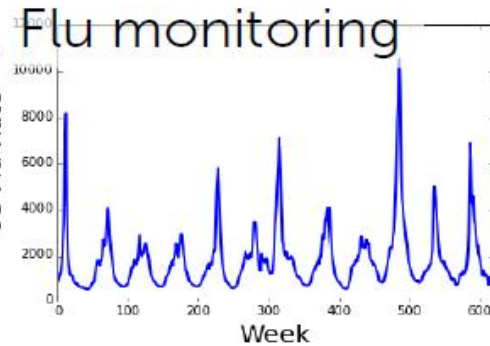*feature 3* = $\sin(2\pi t/12)$

*feature 4* = $\cos(2\pi t/12)$

Fit polynomial trend and sinusoidal seasonal component

5th order polynomial + sin/cos basis

22/12 2020

# Other examples of seasonality

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

Flu monitoring

Demand forecasting
(e.g., jacket purchases)

Motion capture data

22/12 2020

# Generic basic expansion

Model:

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

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

*feature 1* = $h_0(x)$...often 1 (constant)
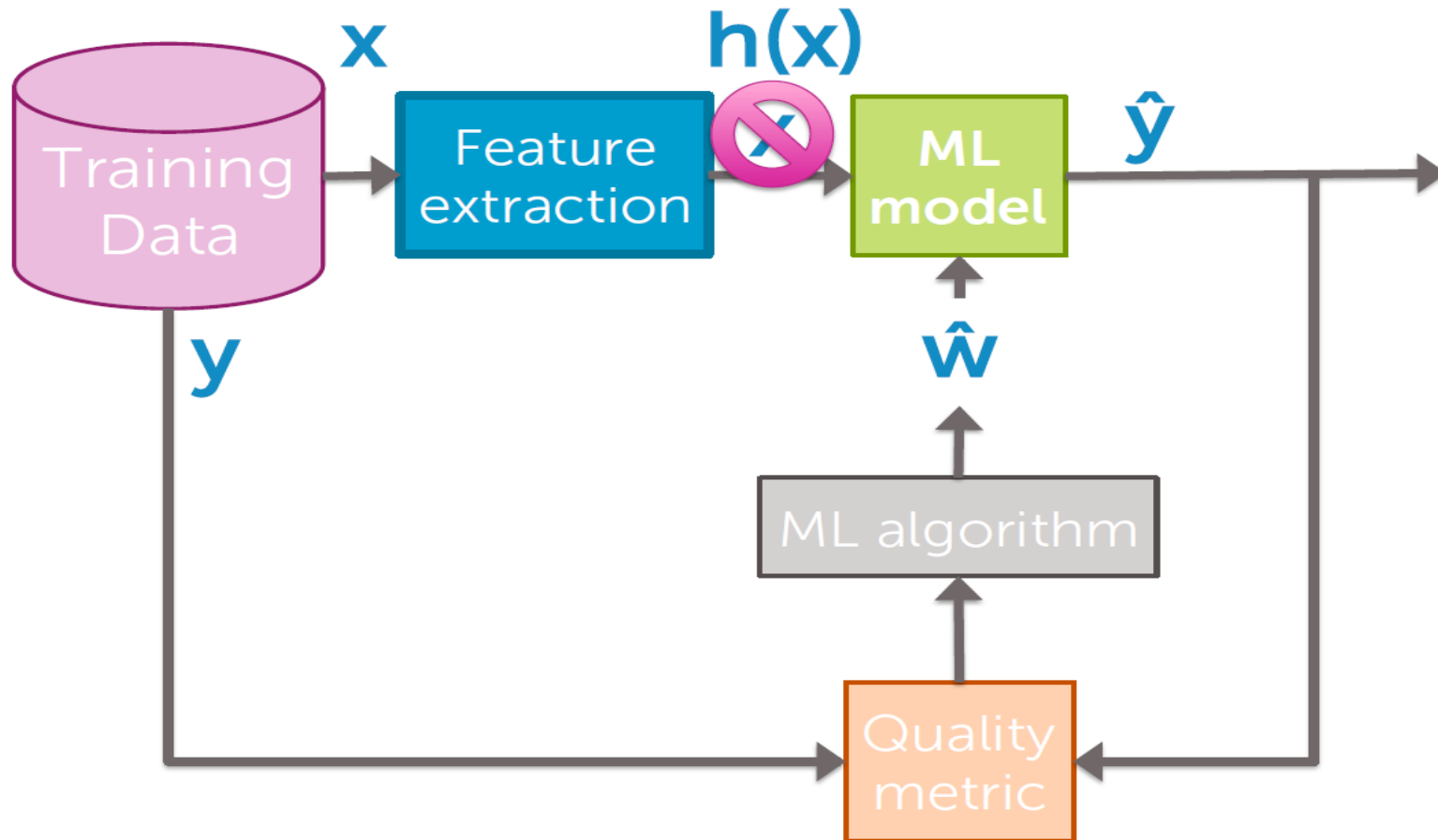
*feature 2* = $h_1(x)$... e.g., x

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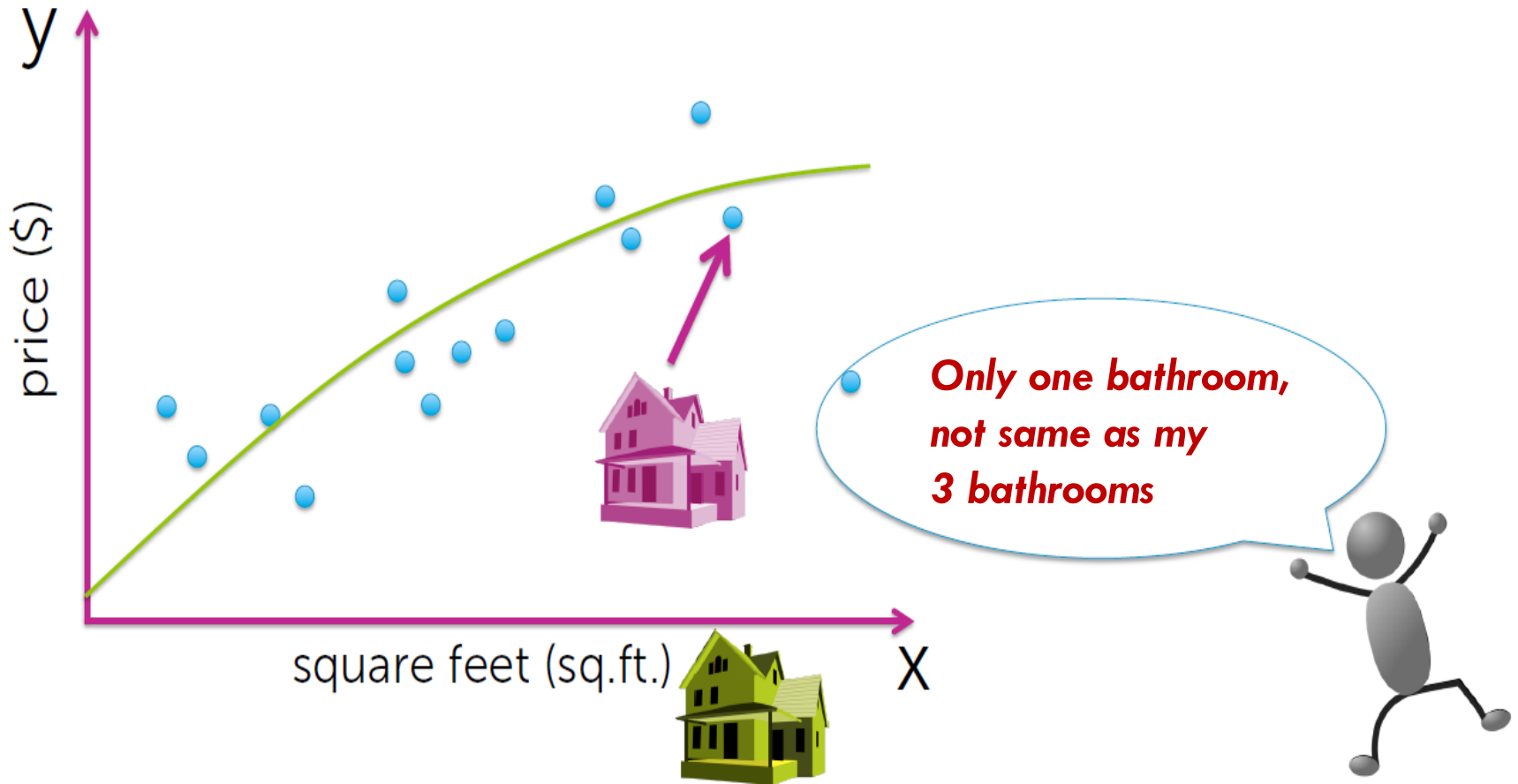
...

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

22/12 2020

# More realistic flow chart

22/12 2020

# Incorporating multiple inputs

22/12 2020

# Incorporating multiple inputs

$$f(\mathbf{x}) = w_0 + w_1 \text{ sq.ft.} + w_2 \text{ \#bath}$$

y

price ($)

# bathroms

$\mathbf{x}[2]$

square feet (sq.ft.)

$\mathbf{x}[1]$

*Many possible inputs*

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

22/12 2020

# General notation

Output: y ↙ scalar

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

↖ d-dim vector

Notational conventions:

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

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

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

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

22/12 2020

# Simple hyperplane

Model:

$$y_i = w_0 + w_1 \mathbf{x}_i[1] + \ldots + w_d \mathbf{x}_i[d] + \varepsilon_i$$

Noise term

feature 1 = 1

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

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

...

feature d+1 = $\mathbf{x}[d]$ ... e.g., lot size

22/12 2020

# More generally: D-dimensional curve

Model:

$$y_i = w_0 \, h_0(\mathbf{x}_i) + w_1 \, h_1(\mathbf{x}_i) + \ldots + w_D \, h_D(\mathbf{x}_i) + \varepsilon_i$$

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

## 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.
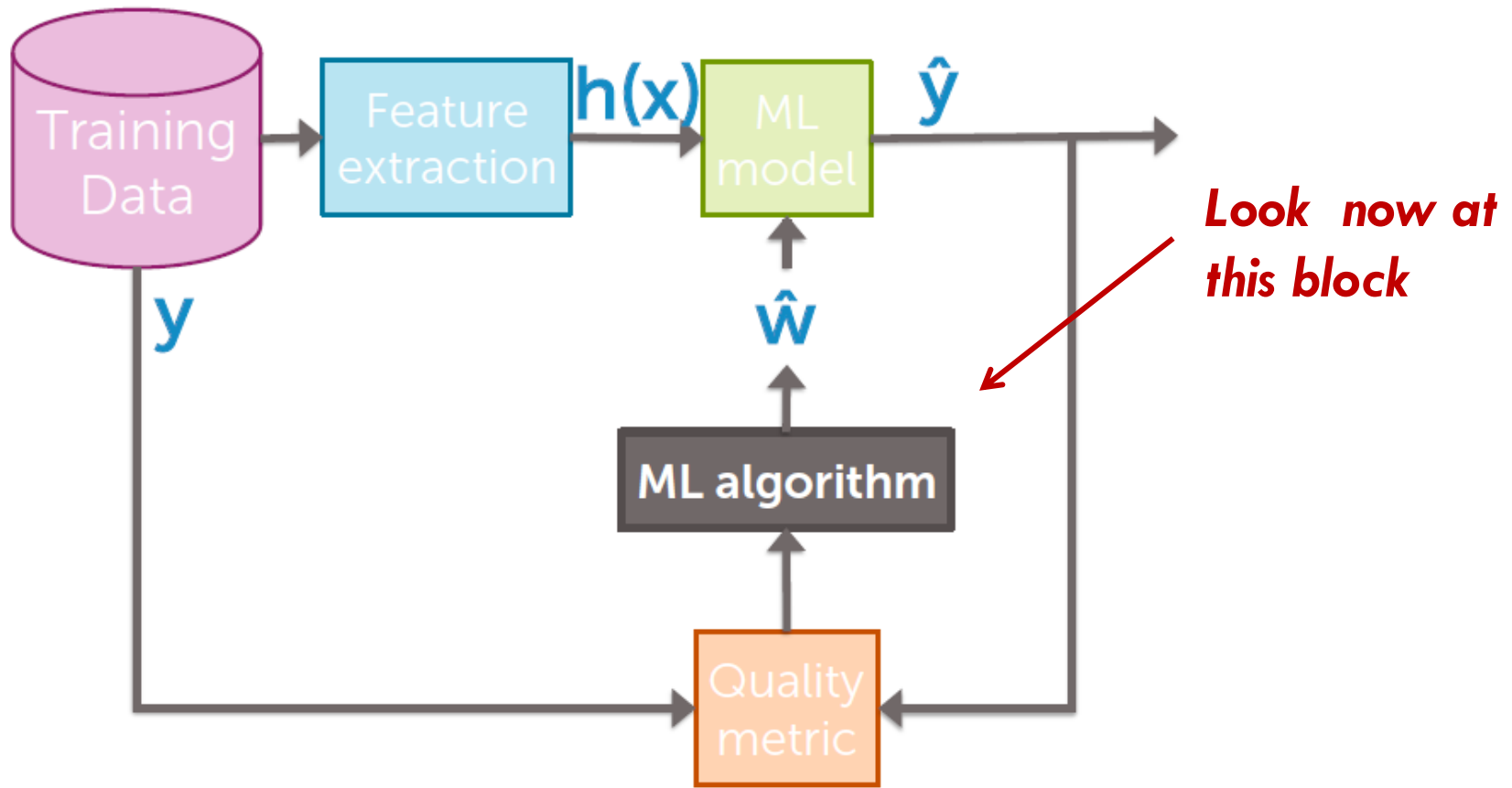feature 3 = $h_2(\mathbf{x})$ ... e.g., $\mathbf{x}[2]$ = \#bath
                          or, $\log(\mathbf{x}[7]) \, \mathbf{x}[2]$ = log(\#bed) x \#bath

...

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

22/12 2020

# Fitting in D-dimmensions

Look now at this block

# Rewriting in vector notation

For observation i

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

$$y_i = \boxed{\phantom{xxxxxxx}}\boxed{\phantom{x}}_{\substack{w_0\ w_1\ w_2\ \dots\ w_D}}^{\mathbf{w}^T} \;\boxed{\phantom{x}}_{\substack{h_0(x_i)\\ h_1(x_i)\\ h_2(x_i)\\ \vdots\\ h_D(x_i)}}^{h(x_i)} + \boxed{\varepsilon_i}$$

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

scalar

$$= \mathbf{w}^T h(x_i) + \varepsilon_i$$

$$= \boxed{\phantom{xxxxxxx}}_{\substack{h_0(x_i)\ h_1(x_i)\ \dots\ h_D(x_i)}}^{h^T(x_i)} \;\boxed{\phantom{x}}_{\substack{w_0\\ w_1\\ \vdots\\ w_D}}^{\mathbf{w}} + \boxed{\varepsilon_i}$$

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

22/12 2020

# Rewriting in matrix notation

For all observations together



$$\Rightarrow \boxed{y = Hw + \epsilon}$$

*Here is our ML algorithm*

22/12 2020

# Fitting in D-dimmensions

*Look now at this block*

©2015 Emily Fox & Carlos Guestrin

Machine Learning Specialization

22/12 2020

# Cost function in D-dimmension

**RSS in vector notation**



$$RSS(\mathbf{w}) = \sum_{i=1}^{N} (y_i - \underbrace{h^T(x_i)\,w}_{\hat{y}_i(w)})^2$$

$\hat{y}_i = $ [ green row vector $h^T(x_i)$: $h_0(x_i)\ h_1(x_i)\ \cdots\ h_D(x_i)$ ] $\times$ $\mathbf{w}$ = $\begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$
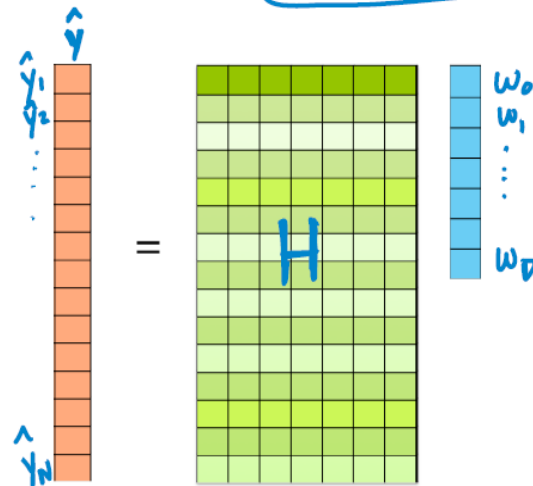
22/12 2020

# Cost function in D-dimmension

## RSS in matrix notation

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

$$= (y - H\omega)^\top (y - H\omega)$$

Why? (part 1)



$$\hat{y} = H\omega$$

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

$$\text{residual}_i = y_i - \hat{y}_i$$

22/12 2020

# Regression model for D-dimmension

**Gradient of RSS**

$$\nabla RSS(\mathbf{w}) = \nabla[(\mathbf{y}-\mathbf{Hw})^\top(\mathbf{y}-\mathbf{Hw})]$$

$$= -2\mathbf{H}^\top(\mathbf{y}-\mathbf{Hw})$$

Why? By analogy to 1D case:

$$\frac{d}{dw}(y-hw)(y-hw) = \frac{d}{dw}(y-hw)^2 = 2\cdot(y-hw)^1(-h)$$
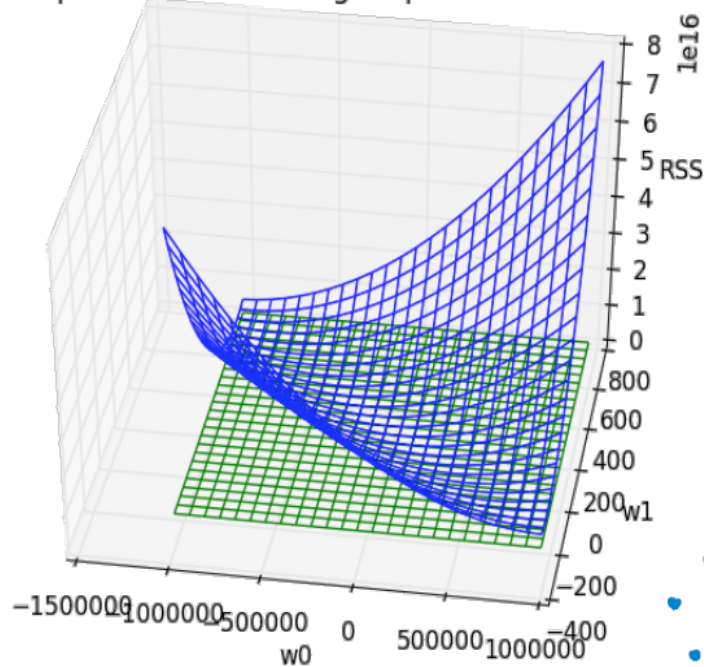
$$= -2h(y-hw)$$

scalars

22/12 2020

# Regression model for D-dimmension

## Approach 1: set gradient to zero



3D plot of RSS with tangent plane at minimum

*Closed form solution*

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

Solve for $\mathbf{w}$:

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

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

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

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

- $A^{-1} A = I$
- $I v = v$
- $I V = V$

22/12 2020

# Closed-form solution

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

*This matrix might not be invertible.*

# features = D

# features

D x D

# features

# obs = N

N

D

Invertible if:

In most cases is $N > D$

really,
# of linearly
ind. observations

Complexity of inverse:

$O(D^3)$

*This might not be CPU feasible.*

22/12 2020

# Regression model for D-dimmension

## Approach 2: gradient descent



Contour plot corresponding to 3D plot of RSS

*We initialise our solution somewhere and then …*

**while** not converged

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

$$-2\mathbf{H}^\top(\mathbf{y} - \mathbf{Hw})$$

$$\leftarrow w^{(t)} + 2\eta H^\top (\underbrace{y - Hw^{(t)}}_{\hat{y}(w^{(t)})})$$

22/12 2020

# Gradient descent

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

$$= \sum_{i=1}^{N} \left( y_i - w_0 h_0(x_i) - w_1 h_1(x_i) - \ldots - w_D h_D(x_i) \right)^2$$

Partial with respect to $w_j$

$$\sum_{i=1}^{N} 2 \left( y_i - w_0 h_0(x_i) - w_1 h_1(x_i) \ldots - w_D h_D(x_i) \right)^1 \cdot (-h_j(x_i))$$
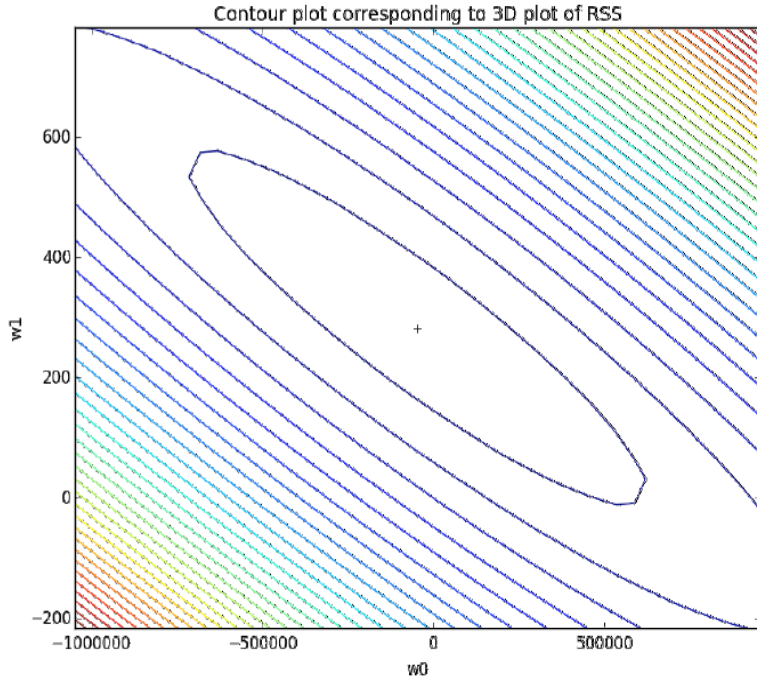
$$= -2 \sum_{i=1}^{N} h_j(x_i)(y_i - h(x_i)^T w)$$

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

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

22/12 2020

# Summary of gradient descent

*Extremely useful algorithm in several applications*



Contour plot corresponding to 3D plot of RSS

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

while $||\nabla RSS(\mathbf{w}^{(t)})|| > \varepsilon$ ← tolerance

$\sqrt{partial[0]^2 + \ldots + partial[D]^2}$

for $j=0,\ldots,D$

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

$w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta \, partial[j]$

$t \leftarrow t + 1$

22/12 2020

# ACCESSING PERFORMANCE

22/12 2020

# Measuring loss

"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:

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

Cost of using ŵ at x when y is true

actual value

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

**Symmetric loss functions**

Examples:
(assuming loss for underpredicting = overpredicting)

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

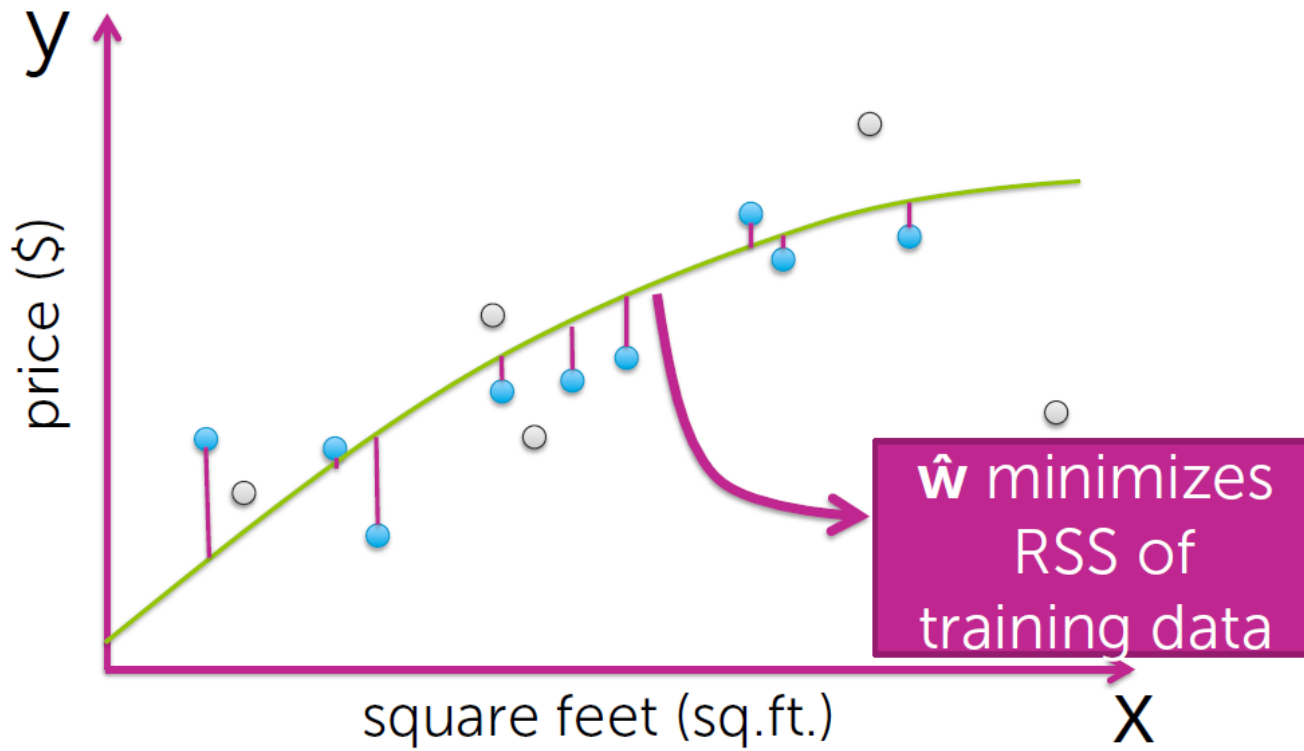Squared error: $L(y, f_{\hat{w}}(\mathbf{x})) = (y - f_{\hat{w}}(\mathbf{x}))^2$

22/12 2020

# Accessing the loss

**Use training data**



$\hat{w}$ minimizes RSS of training data

22/12 2020

# Compute training error

1. Define a loss function $L(y, f_{\hat{w}}(\mathbf{x}))$
   - E.g., squared error, absolute error,...
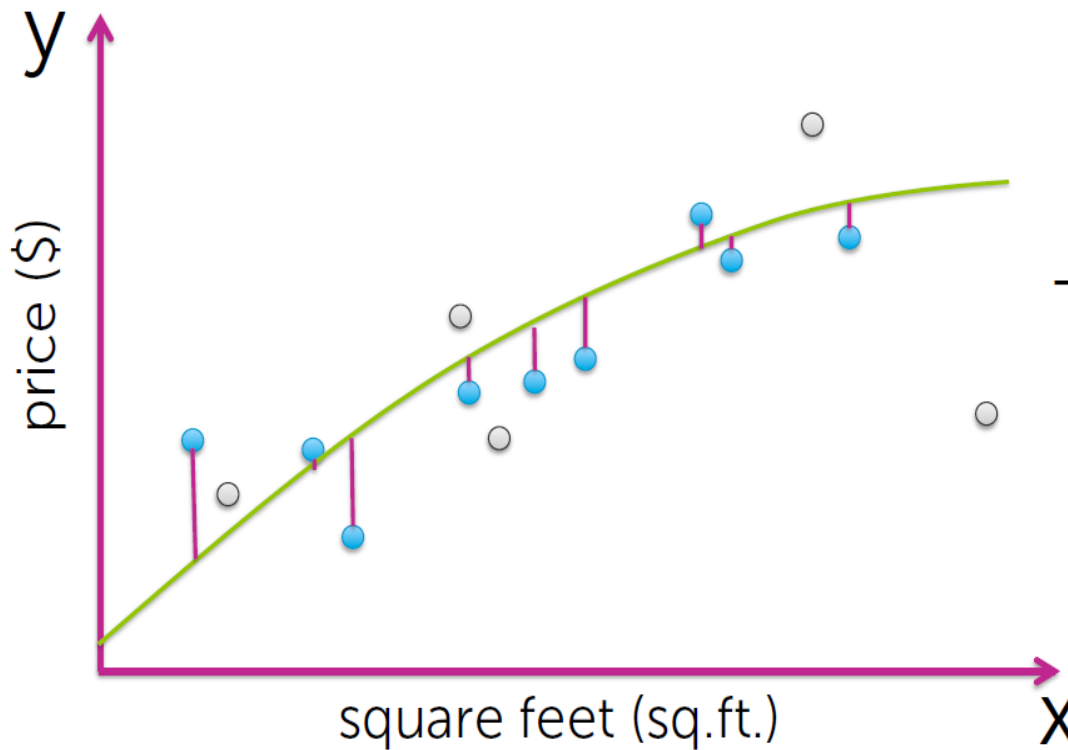
2. Training error
   = avg. loss on houses in training set
   $$= \frac{1}{N} \sum_{i=1}^{N} L(y_i, f_{\hat{w}}(\mathbf{x}_i))$$

fit using training data

22/12 2020

# Training error

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



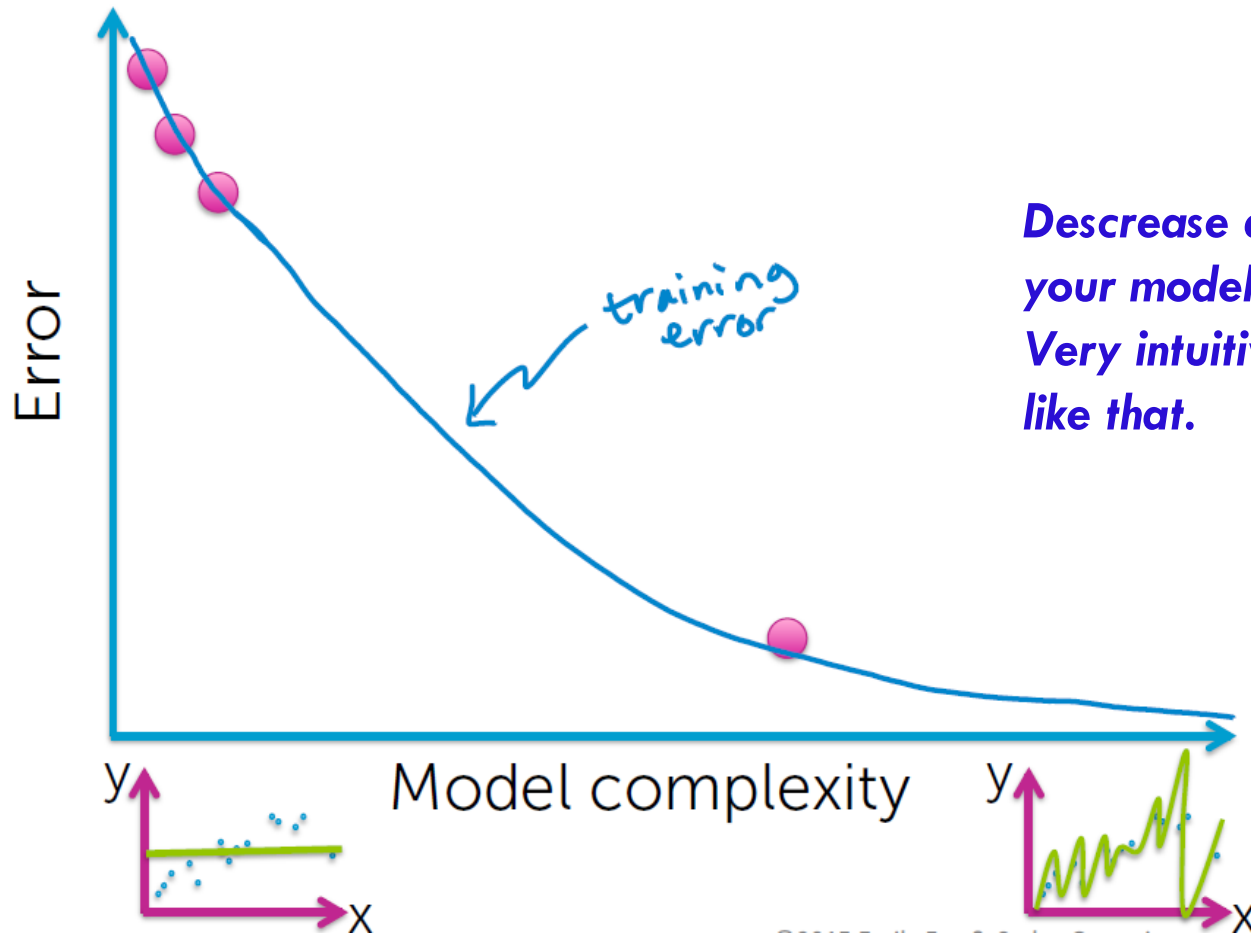*Convention is to take average here*

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

22/12 2020

# Training error vs. model complexity

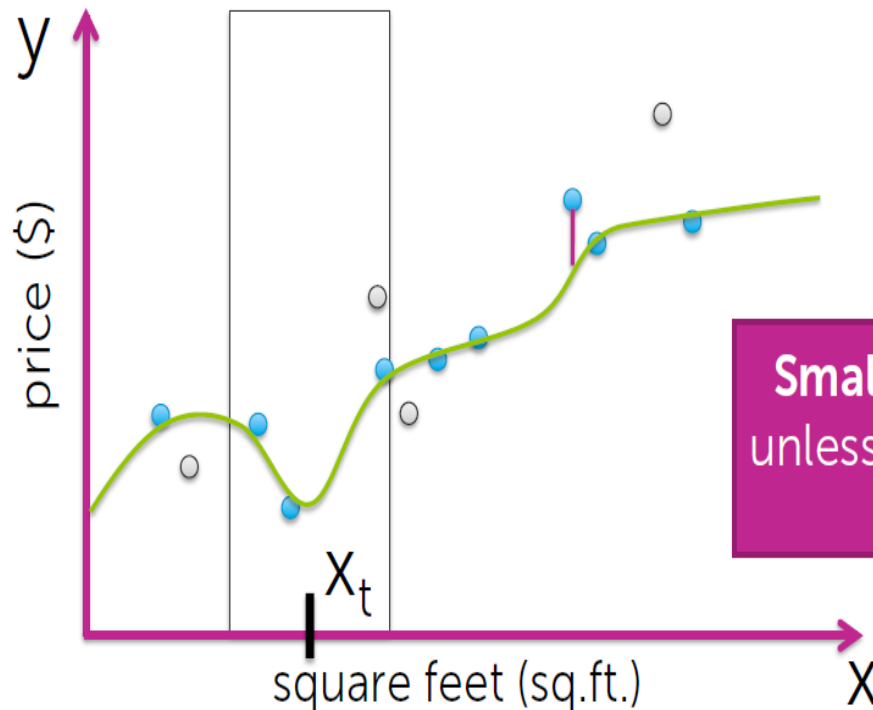*Descrease as you increase your model complexity. Very intuitive why it is like that.*

22/12 2020

# Is training error a good measure?

Issue: Training error is overly optimistic

because **ŵ** was fit to training data

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

**Small training error ≠> good predictions**
unless training data includes everything you
might ever see

22/12 2020

# Generalisation (true) error

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



Lots of houses
in neighborhood,
but not in dataset

22/12 2020

# Generalisation error vs model complexity

*However … in contrast to the training error, in practice we cannot really compute true generalisation error. We don't have data on all possible houses in the area.*

Can't compute!

22/12 2020

# Forming a test set

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

22/12 2020

# Compute test error

Test error

= avg. loss on houses in test set

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

↑
# test points

fit using training data

**has never seen test data!**

22/12 2020

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

*Test error: noisy version due to limited statistics.*

Overfitting if:

If there exists a model with estimated params $w'$ such that

① training error $(\hat{w})$
   $<$ training error $(w')$

② true error $(\hat{w})$
   $>$ true error $(w')$

Machine Learning Specialization

22/12 2020

# Training/test splits

| Training set | Test set |
|---|---|

Too few → **ŵ** poorly estimated

| Training set | Test set |
|---|---|

Too few → test error bad approximation of generalization error

| Training set | Test set |
|---|---|

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...)
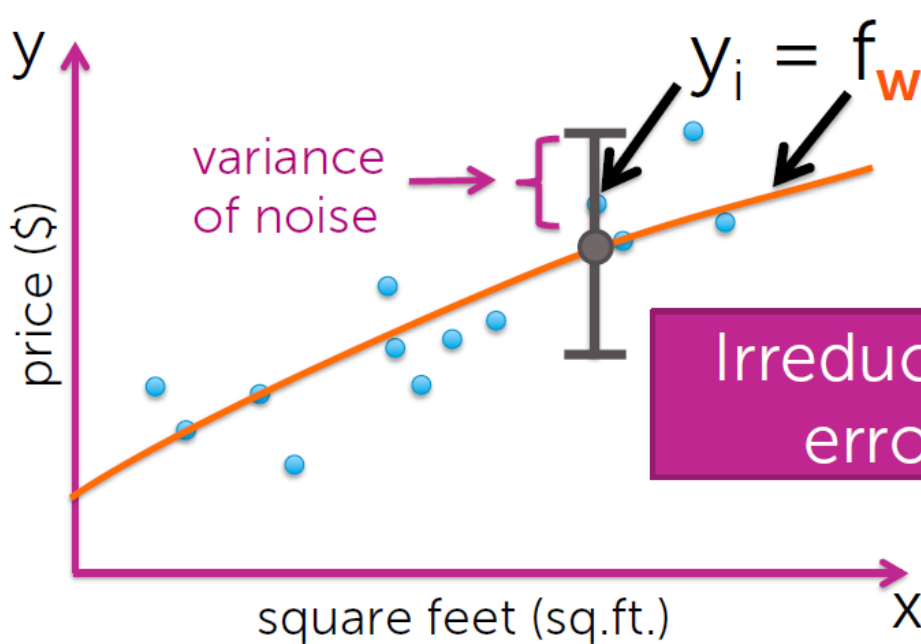
22/12 2020

# Three sources of errors

In forming predictions, there are 3 sources of error:

1. Noise

2. Bias

3. Variance

22/12 2020

# Data are inherently noisy

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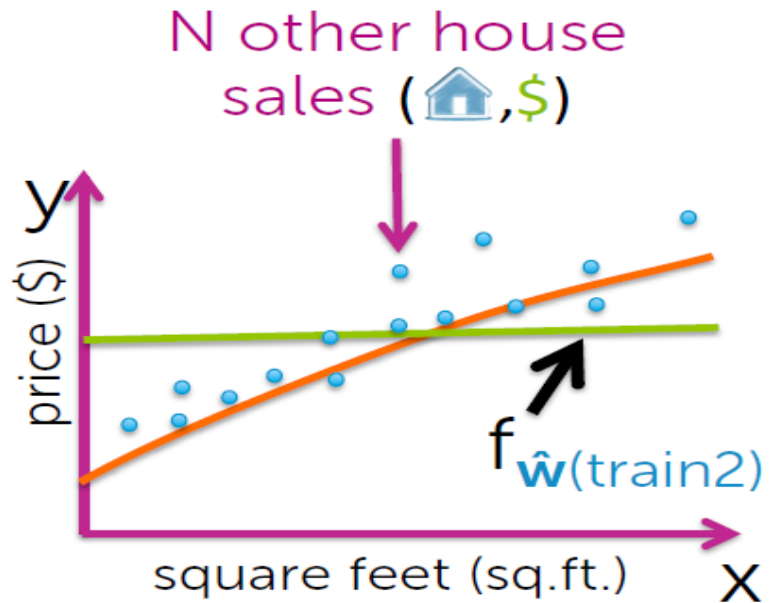


$$y_i = f_{\mathbf{w}(true)}(\mathbf{x}_i) + \varepsilon_i$$

variance of noise

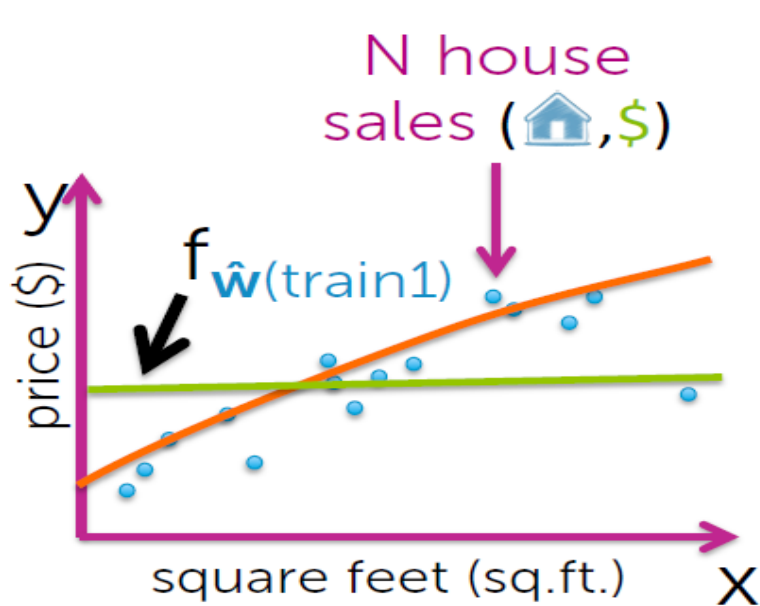price ($)

Irreducible error

square feet (sq.ft.)

*We cannot reduce it by chosing better model or procedure, It is beyond our control.*

22/12 2020

# Bias contribution

*This contribution we can control.*


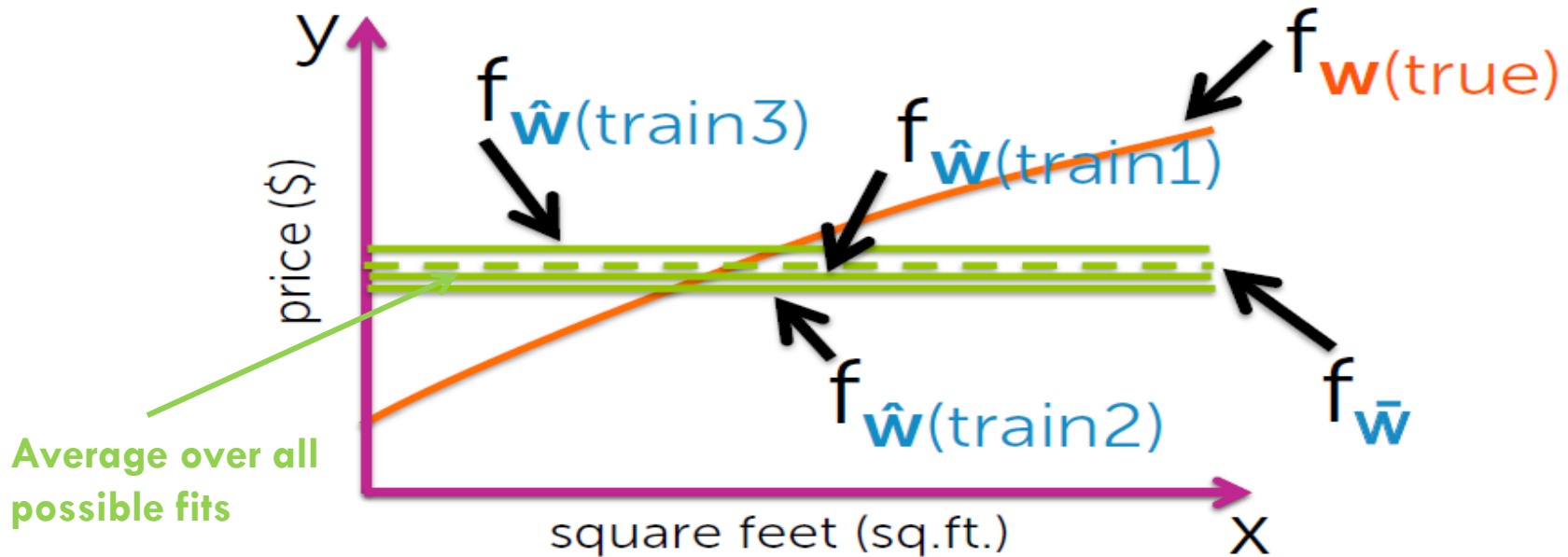
Assume we fit a constant function

# Bias contribution

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

# Bias contribution

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

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



low complexity
→
high bias

$f_{\mathbf{w}}(\text{true})$

$f_{\bar{\mathbf{w}}}$

y

price ($)

square feet (sq.ft.)

x

22/12 2020

# Variance contribution

How much do specific fits
vary from the expected fit?



22/12 2020

# Variance contribution

How much do specific fits
vary from the expected fit?

Can specific fits
vary widely?
If so, erratic
predictions

$f_{\hat{w}}(train3)$

$f_{\hat{w}}(train1)$

$f_{\hat{w}}(train2)$

$f_{\bar{w}}$

y

price ($)

square feet (sq.ft.)

x

low complexity
→
low variance

22/12 2020

# Variance of high complexity models

## Assume we fit a high-order polynomial

*For each train remove few random houses*



high complexity → high variance

22/12 2020

# Bias of high complexity models
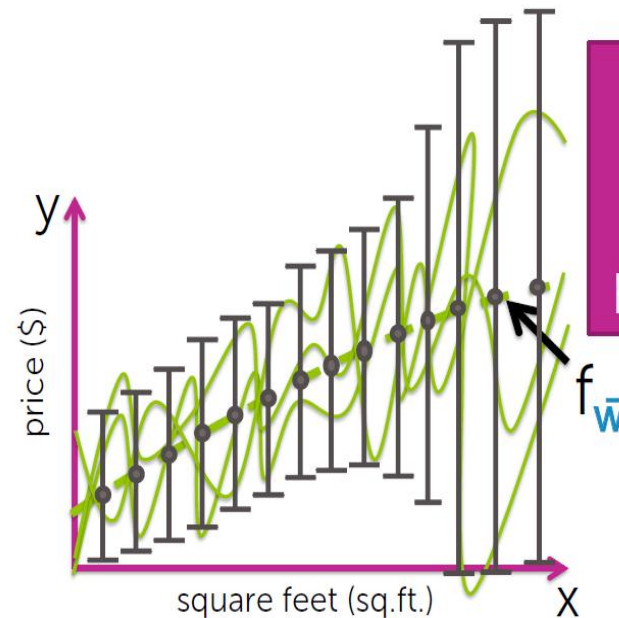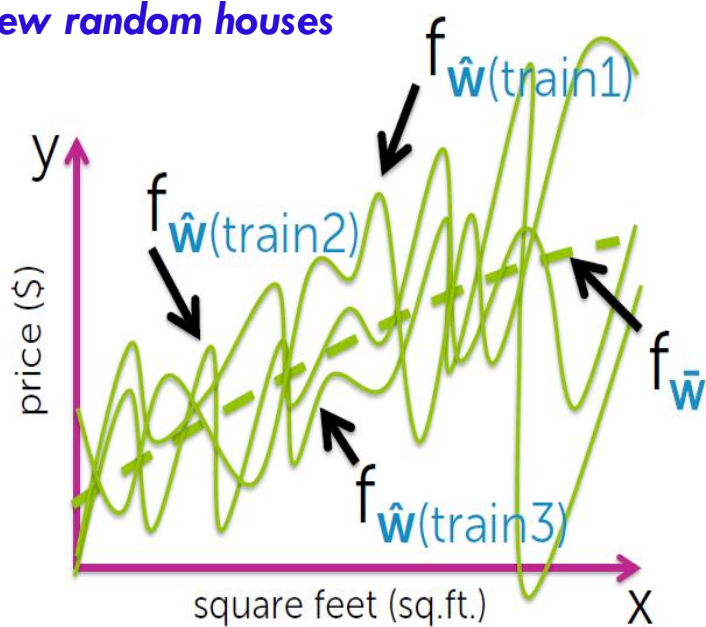
## Assume we fit a high-order polynomial

*For each train remove few random houses*



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

22/12 2020

# Bias –variance tradeoff

MSE = bias² +variance

sweet spot

variance

bias

Model complexity

**MSE = mean square error**

**Machine Learing
is all about this tradeoff**

*But….*

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

# Errors vs amount of data

for a fixed model complexity

ŵ not approx. well from few points

true error

In the limit true error = training error

with few data points, fixed complexity model can fit these points reasonably well

training error

In the limit, will flatten out to how well model can fit true relationship $f_{true}$

bias + noise

Error

# data points in training set

22/12 2020

# The regression/ML workflow

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

22/12 2020

# Hypothetical implementation

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

**Overly optimistic!**

2. **Model assessment**

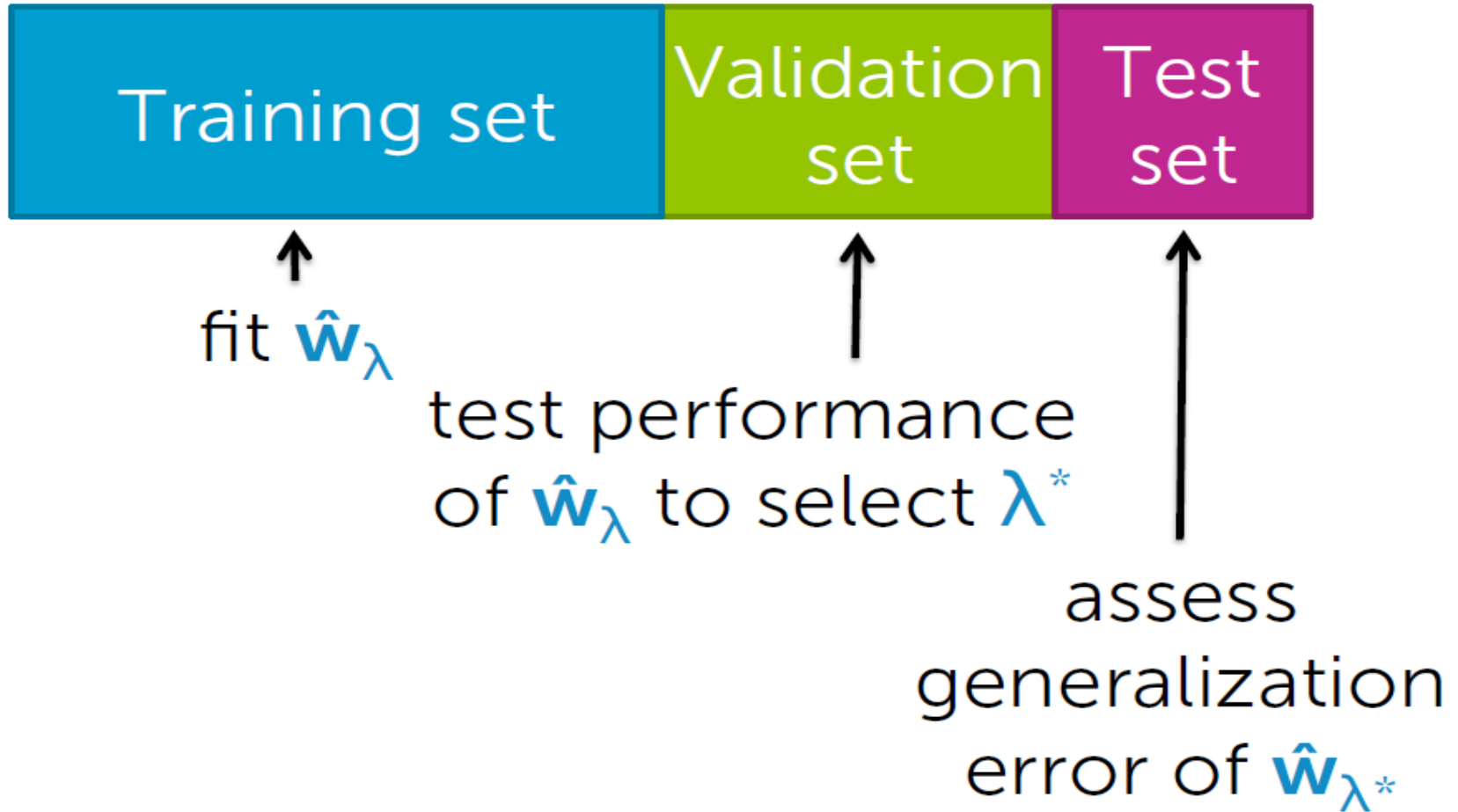Compute test error of $\hat{\mathbf{w}}_{\lambda^*}$ (fitted model for selected complexity $\lambda^*$) to approx. generalization error

22/12 2020

# Practical implementation

| Training set | Validation set | Test set |

fit $\hat{\mathbf{w}}_\lambda$

test performance of $\hat{\mathbf{w}}_\lambda$ to select $\lambda^*$

assess generalization error of $\hat{\mathbf{w}}_{\lambda^*}$

22/12 2020

# Typical splits

| Training set | Validation set | Test set |
|:---:|:---:|:---:|
| 80% | 10% | 10% |
| 50% | 25% | 25% |

# K-fold cross validation

## K-fold cross validation



$\hat{\mathbf{w}}_\lambda^{(5)}$

$error_5(\lambda)$

For k=1,...,K

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

Compute average error: $CV(\lambda) = \dfrac{1}{K} \displaystyle\sum_{k=1}^{K} error_k(\lambda)$

22/12 2020

# What value of K

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 $\lambda$
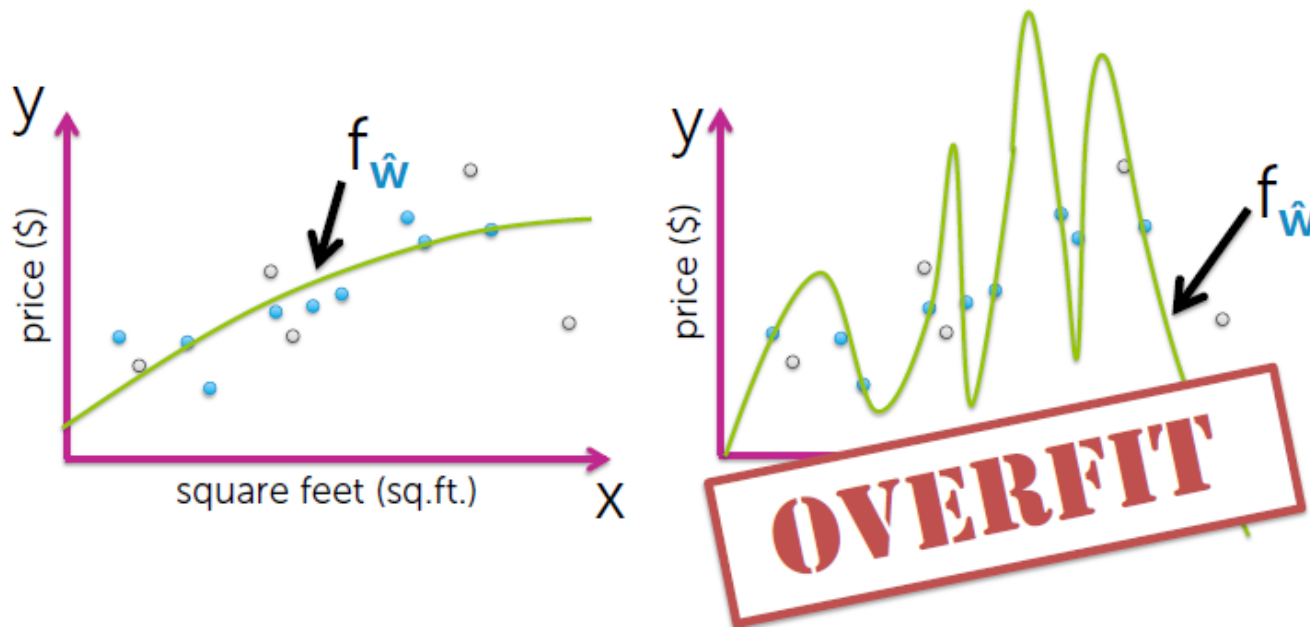
Typically, K=5 or 10

5-fold CV

10-fold CV

22/12 2020

# RIDGE REGRESSION

22/12 2020

# Flexibility of high-order polynomials

$$y_i = w_0 + w_1 x_i + w_2 x_i^2 + ... + w_p x_i^p + \varepsilon_i$$



***Symptoms for overfitting: often associated with very large value of estimated parameters $\hat{w}$***
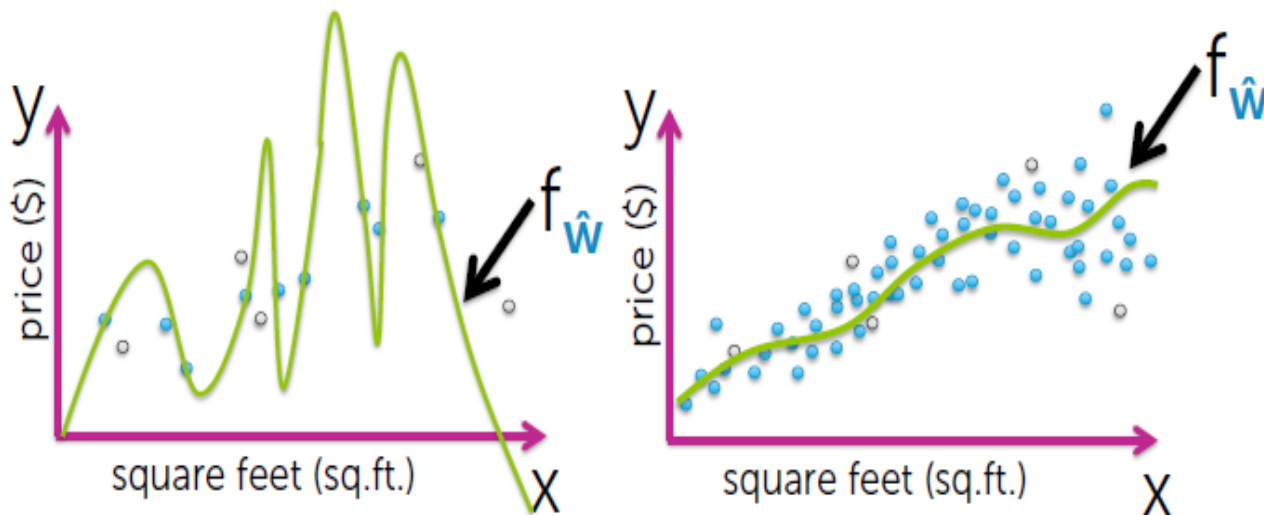
# How does # of observations influence overfitting?

**Few observations** (N small)
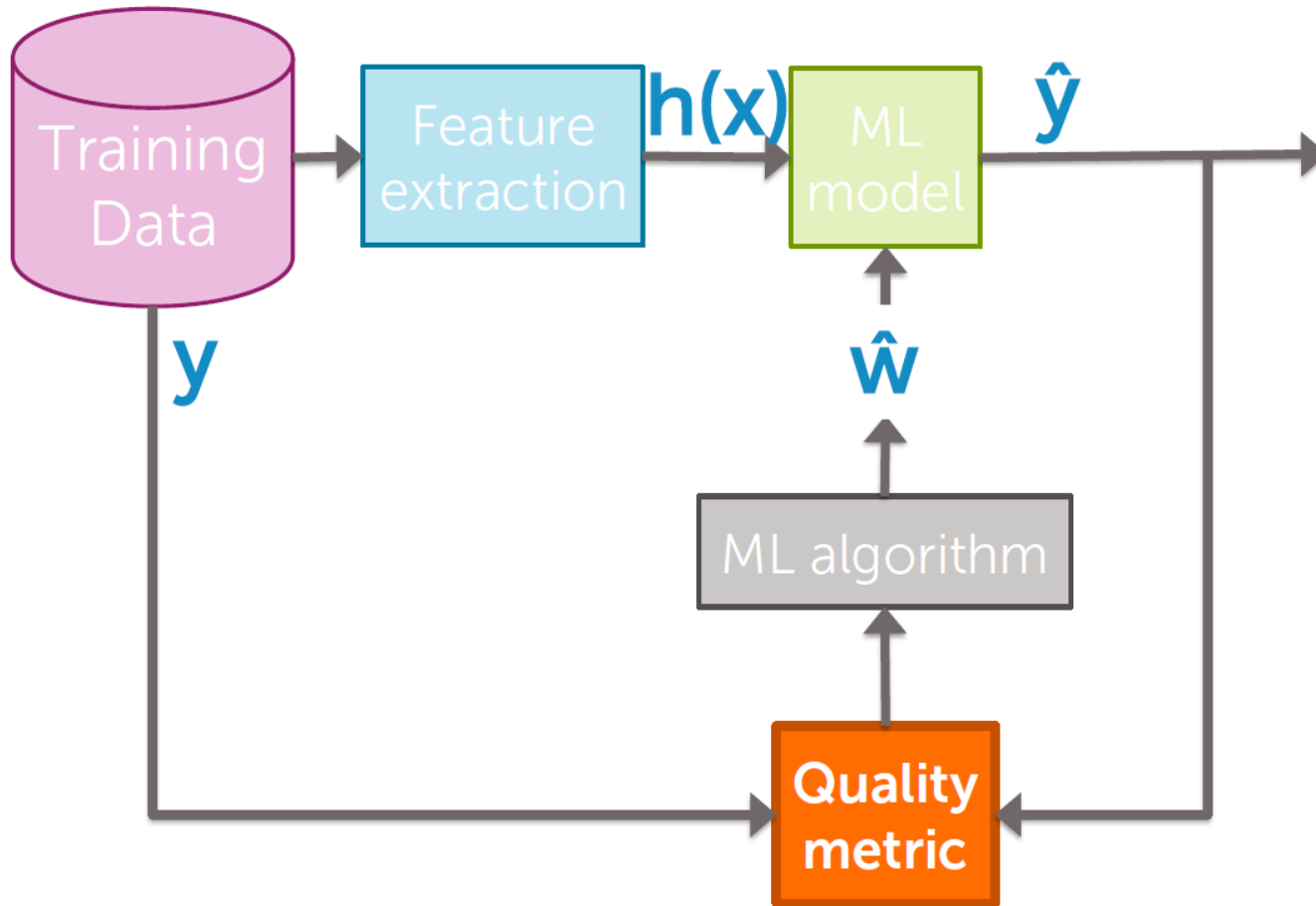→ rapidly overfit as model complexity increases

**Many observations** (N very large)
→ harder to overfit



22/12 2020

# Lets improve quality metric blok

# Desire total cost format

Want to balance:

i.   How well function fits data

ii.   Magnitude of coefficients

**want to balance**

Total cost =

    measure of fit + measure of magnitude of coefficients

small # = good fit to training data

small # = not overfit

22/12 2020

# Measure of magnitude of regression coefficients

What summary # is indicative of size of regression coefficients?

– Sum?  $W_0 = 1,527,301$  $W_1 = -1,605,253$

$W_0 + W_1 = \text{small } \#$

*But … the coefficients are very large*

– Sum of absolute value?

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

– Sum of squares ($L_2$ norm)

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

22/12 2020

# Consider specific total cost

Total cost =

measure of fit + measure of magnitude of coefficients

$RSS(\mathbf{w})$

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

22/12 2020

# Consider resulting objects

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

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

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

tuning parameter = balance of fit and magnitude

If $\lambda=0$:
reduces to minimizing $RSS(w)$, as before (old solution) $\longrightarrow \hat{w}^{LS} \leftarrow$ least squares

If $\lambda=\infty$:
For solutions where $\hat{w} \neq 0$, then total cost is $\infty$
If $\hat{w}=0$, then total cost $= RSS(0) \longrightarrow$ solution is $\hat{w}=0$

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

22/12 2020

# Ridge regression: bias-variance tradeoff

Large $\lambda$:

high bias, low variance

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

In essence, $\lambda$ controls model complexity
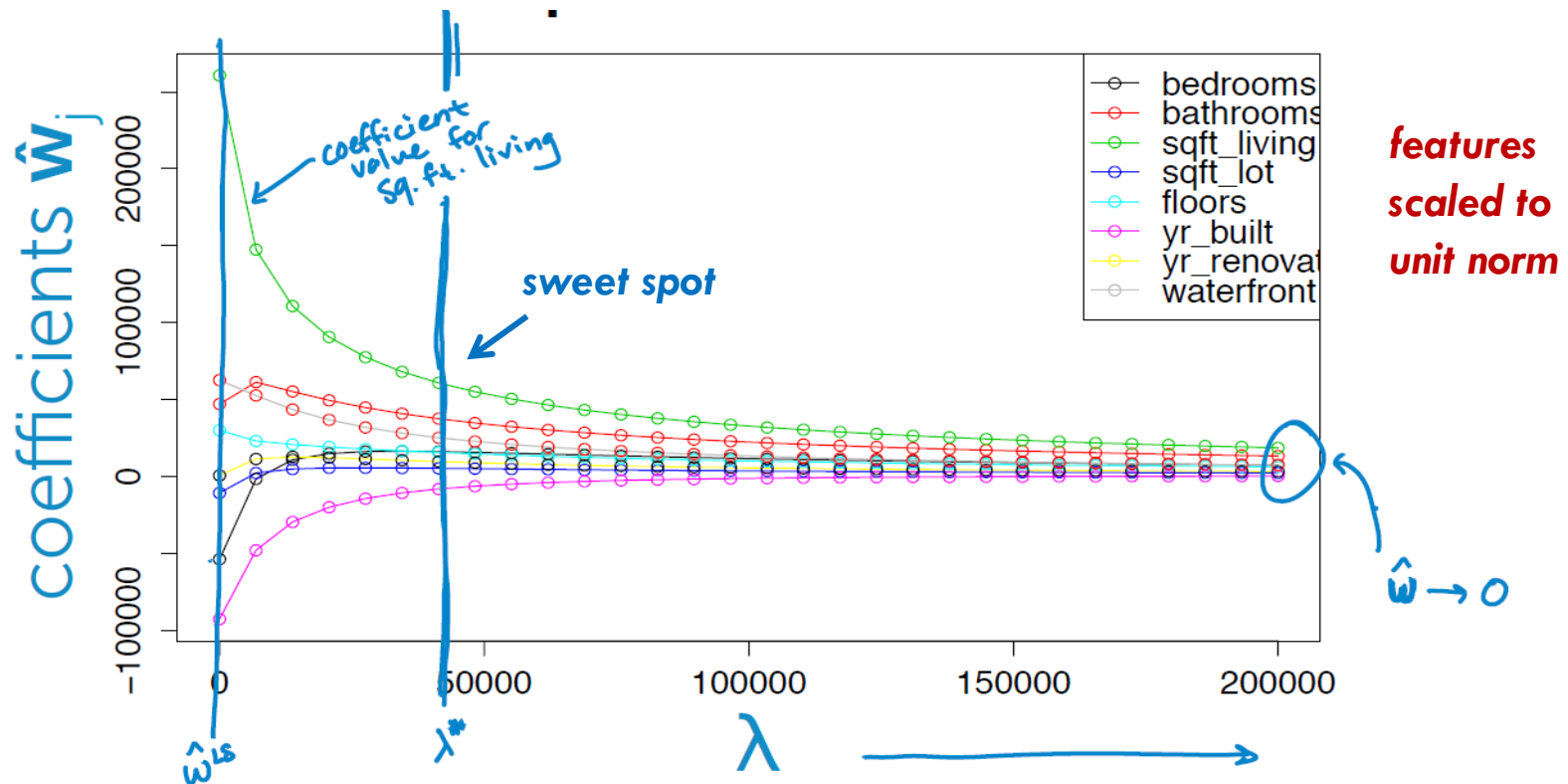
Small $\lambda$:

low bias, high variance

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

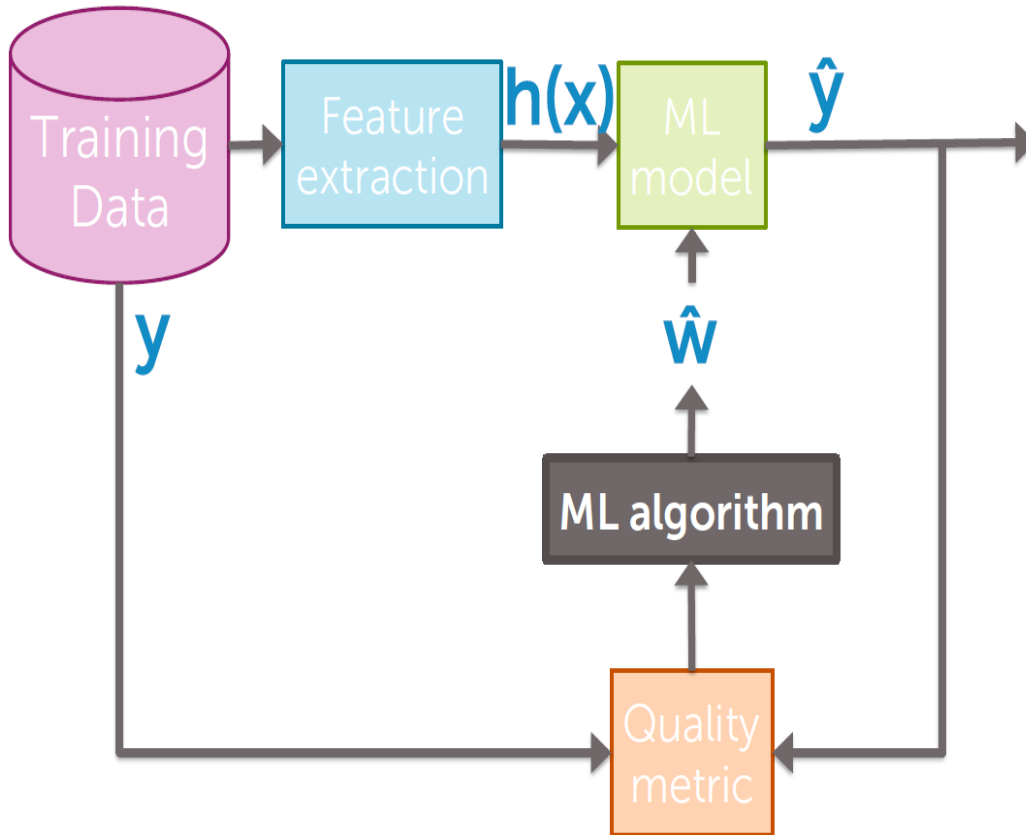22/12 2020

# Ridge regression: coefficients path

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



*features scaled to unit norm*

22/12 2020

# Flow chart

Model for all N observations together

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

22/12 2020

# Ridge regression: cost in matrix notation

In matrix form, ridge regression cost is:

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

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

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

$$= \begin{bmatrix} w_0 & w_1 & w_2 & \cdots & w_D \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$$

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

22/12 2020

# Gradient of ridge regresion cost

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

$$= \underbrace{[\nabla(\mathbf{y}-\mathbf{Hw})^\top(\mathbf{y}-\mathbf{Hw})]}_{-2\mathbf{H}^\top(\mathbf{y}-\mathbf{Hw})} + \lambda \underbrace{[\nabla\mathbf{w}^\top\mathbf{w}]}_{2\mathbf{w}}$$
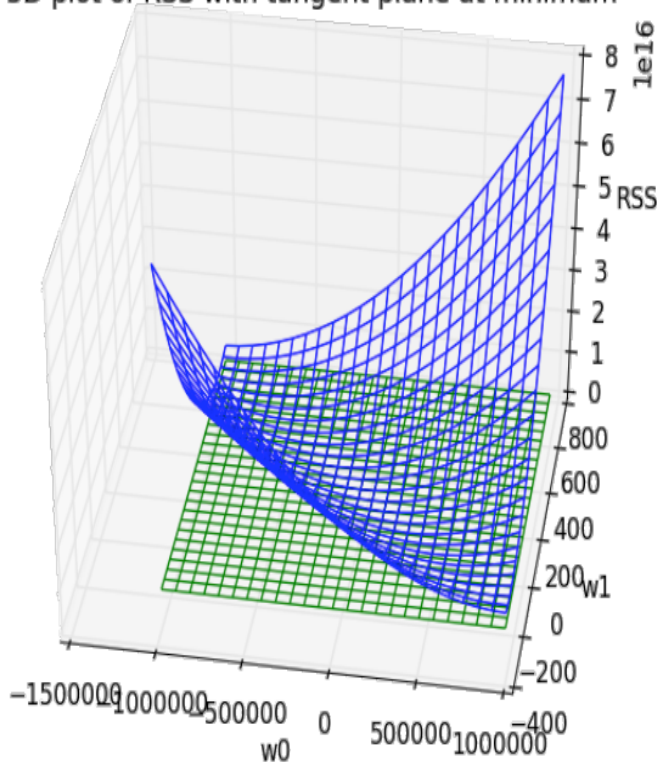
Why?   By analogy to 1d case…

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

22/12 2020

# Ridge regression:  closed-form solution

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 $\mathbf{w}$.

$$-H^T y + H^T H \hat{w} + \lambda I \hat{w} = 0$$

$$H^T H \hat{w} + \lambda I \hat{w} = H^T y$$

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

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

22/12 2020

# Ridge regression: gradient descent

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



Contour plot corresponding to 3D plot of RSS

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

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta *$$

Same as before (from RSS term)

$$\left[ -2 \sum_{i=1}^{N} h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)})) + 2\lambda w_j^{(t)} \right]$$

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

22/12 2020

# Summary of ridge regression algorithm

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

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

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

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

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

    $t \leftarrow t + 1$

22/12 2020

# How to handle the intercept

**Recall multiple regression model**

Model:

$$y_i = w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \ldots + w_D h_D(\mathbf{x}_i) + \varepsilon_i$$

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

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?

Standard ridge regression cost:

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

strength of penalty

Encourages intercept $w_0$ to also be small

Do we want a small intercept?
Conceptually, not indicative of overfitting...

22/12 2020

# Do we penalize intercept?

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

  Modified ridge regression cost:

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

- **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)
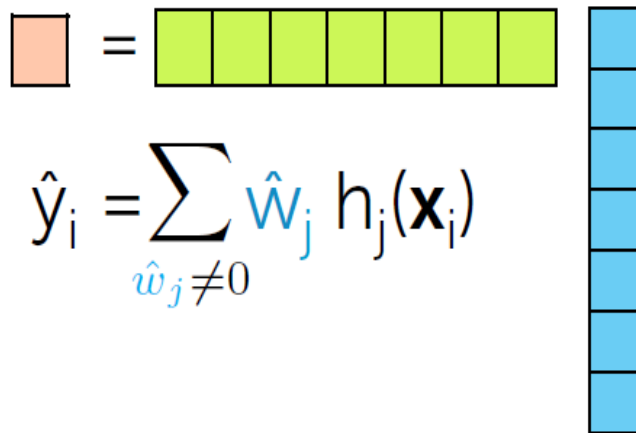
22/12 2020

# FEATURES SELECTION

# &

# LASSO REGRESSION

22/12 2020

# Why features selection?

## Efficiency:

- If 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:
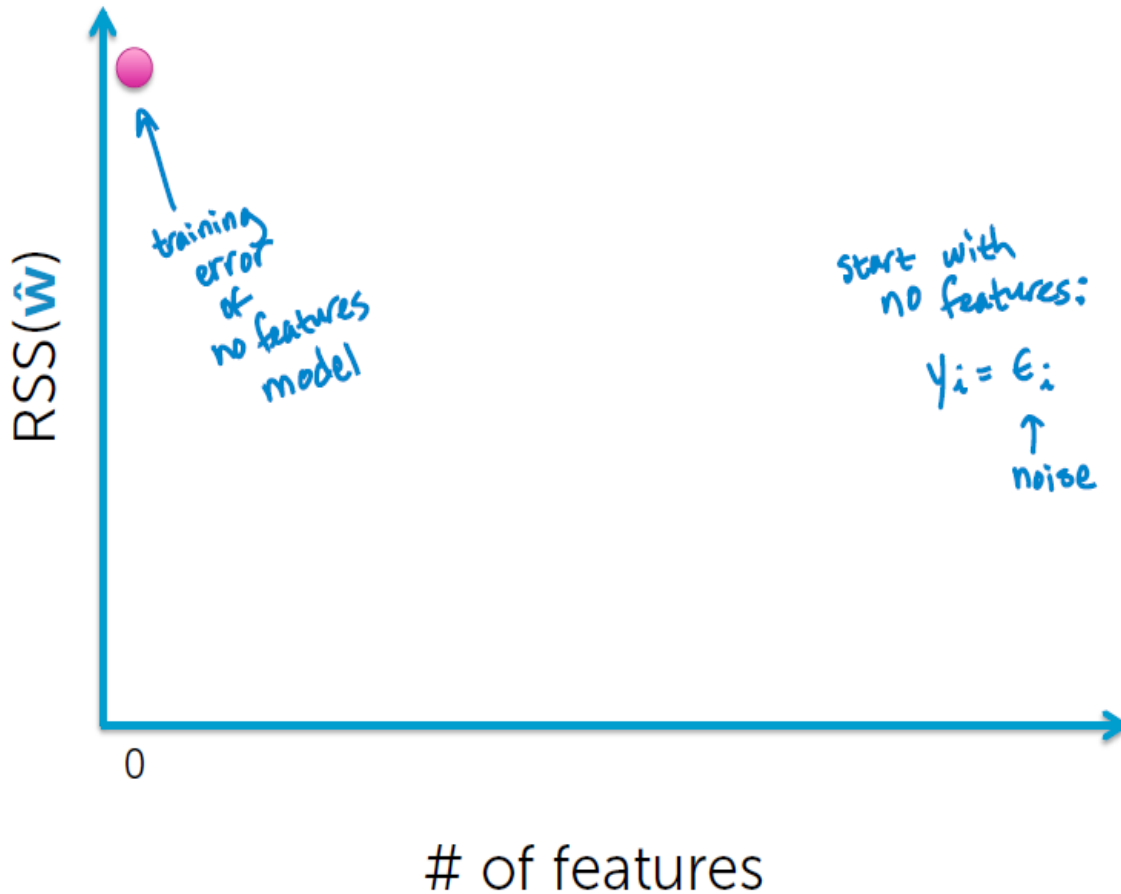
- Which features are relevant for prediction?

# Sparcity

## Housing application



| | |
|---|---|
| Lot size | Dishwasher |
| Single Family | Garbage disposal |
| Year built | Microwave |
| Last sold price | Range / Oven |
| Last sale price/sqft | Refrigerator |
| Finished sqft | Washer |
| Unfinished sqft | Dryer |
| Finished basement sqft | Laundry location |
| # floors | Heating type |
| Flooring types | Jetted Tub |
| Parking type | Deck |
| Parking amount | Fenced Yard |
| Cooling | Lawn |
| Heating | Garden |
| Exterior materials | Sprinkler System |
| Roof type | ⋮ |
| Structure style | |

# Find best model of size: 0

RSS($\hat{w}$)

training error of no features model

start with no features:

$y_i = \epsilon_i$

↑ noise

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

0

# of features

22/12 2020
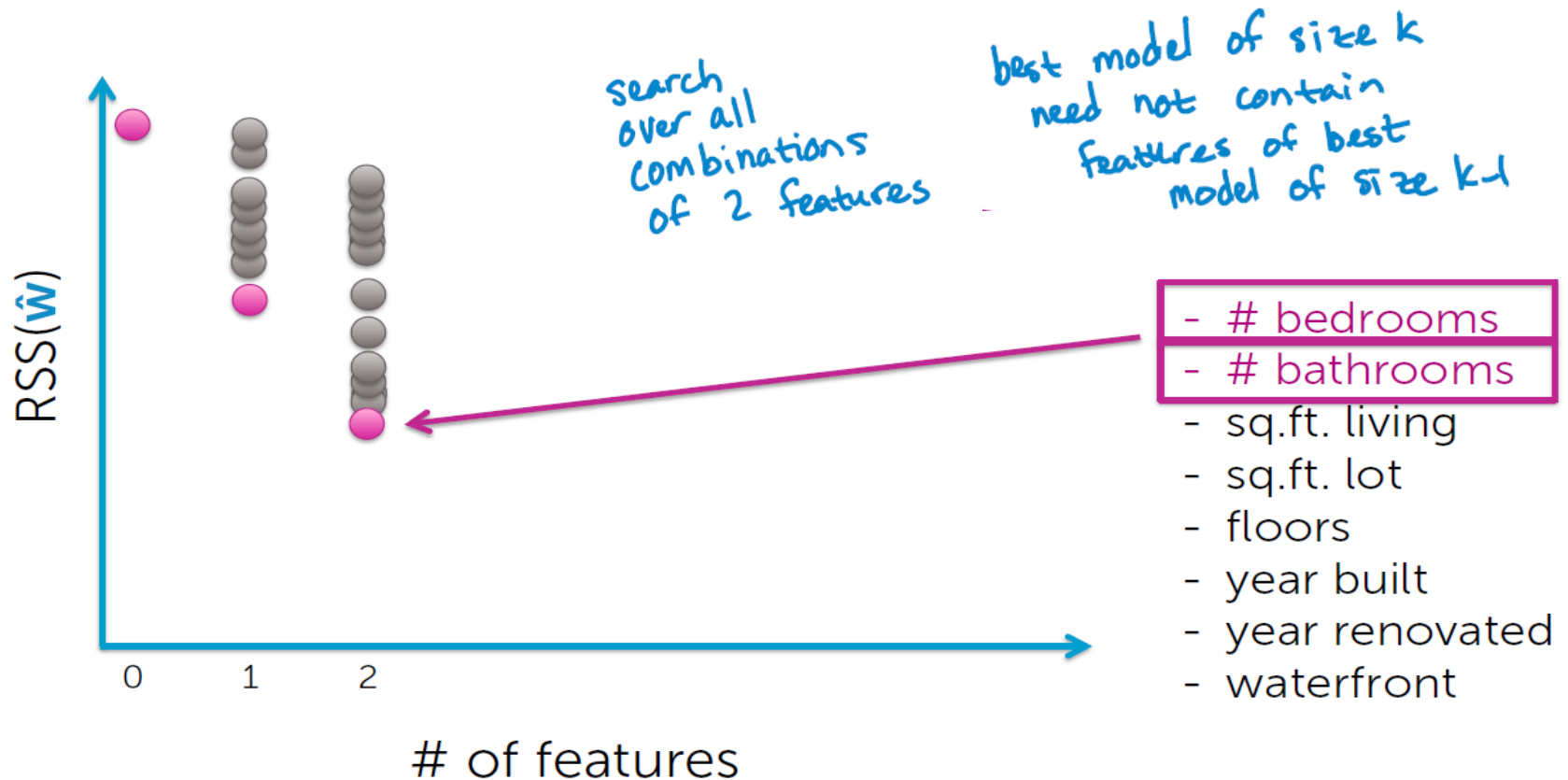
# Find best model of size: 1
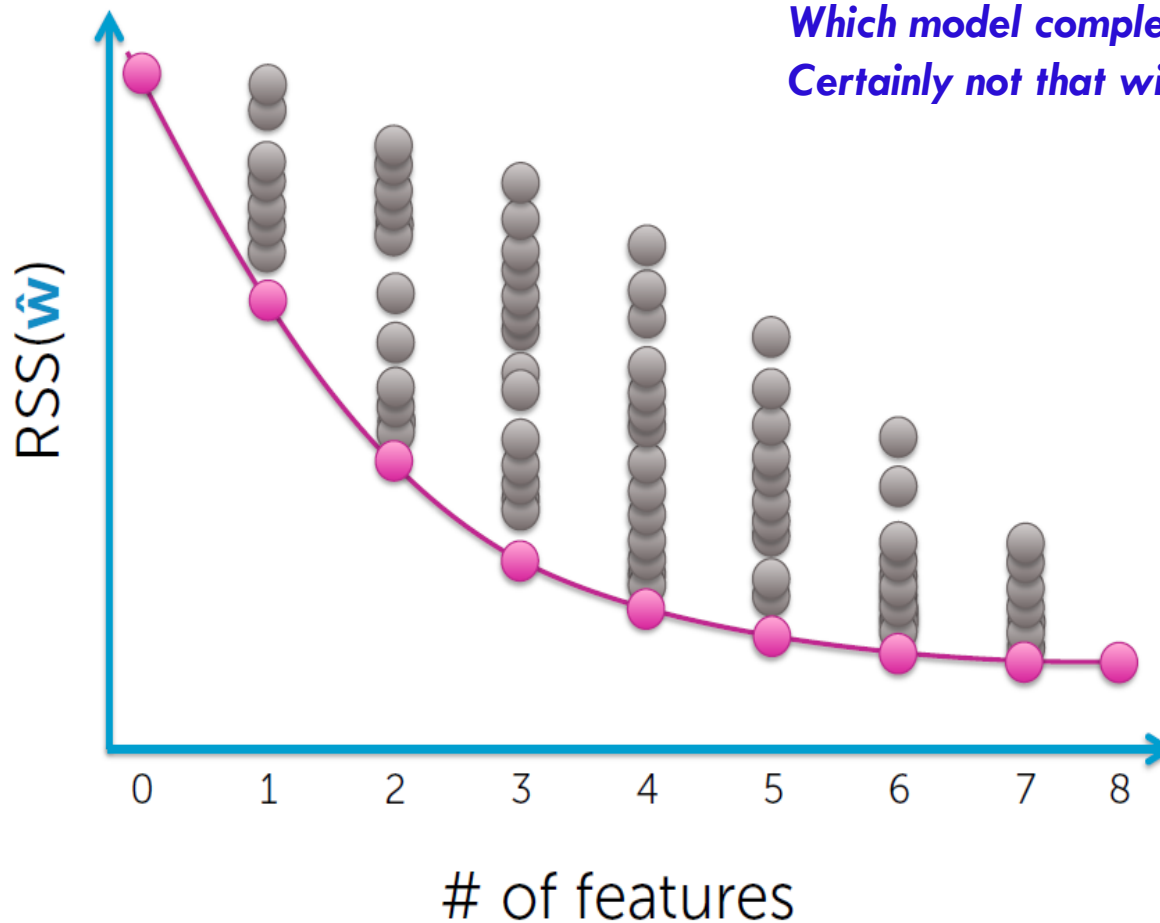
training of
error of
model fit
just with
#bed as feature

best fitting model
with only one feature

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

RSS($\hat{w}$)

0   1

# of features

22/12 2020

# Find best model of size: 2

**Note: not necessarily nested!**



search
over all
combinations
of 2 features

best model of size k
need not contain
features of best
model of size k-1

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

RSS($\hat{w}$)

0    1    2

# of features

22/12 2020

# Find best model of size: N

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

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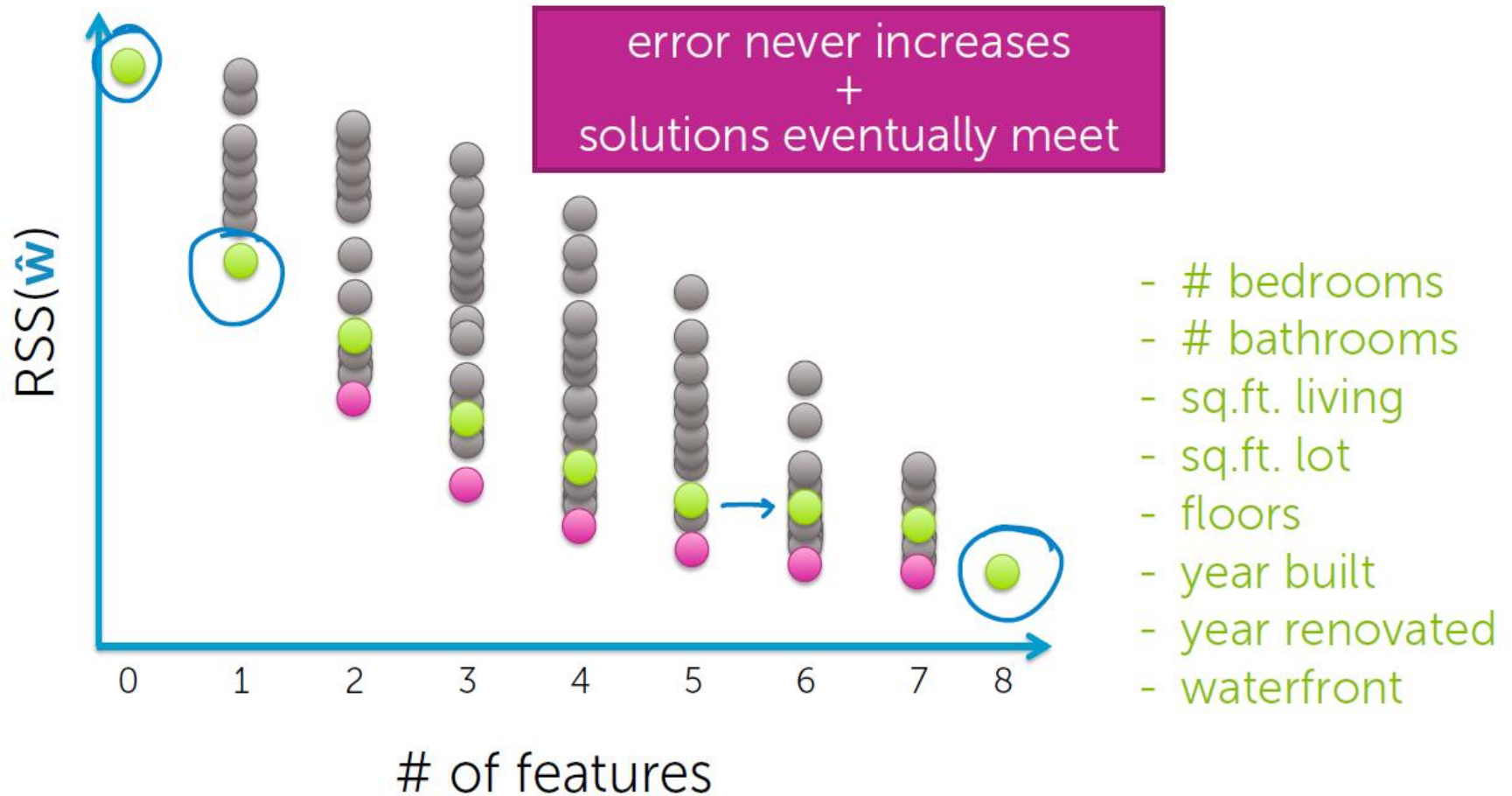- floors
- year built
- year renovated
- waterfront

22/12 2020

# Choosing model complexity

Option 1: Assess on validation set

Option 2: Cross validation

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

22/12 2020

# Complexity of „all subsets"

## How many models were evaluated?

- each indexed by features included

$y_i = \varepsilon_i$         [0 0 0 ... 0 0 0]

$y_i = w_0 h_0(x_i) + \varepsilon_i$      [1 0 0 ... 0 0 0]

$y_i = w_1 h_1(x_i) + \varepsilon_i$      [0 1 0 ... 0 0 0]

⋮        ⋮

$y_i = w_0 h_0(x_i) + w_1 h_1(x_i) + \varepsilon_i$     [1 1 0 ... 0 0 0]

⋮        ⋮

$y_i = w_0 h_0(x_i) + w_1 h_1(x_i) + ... + w_D h_D(x_i) + \varepsilon_i$     [1 1 1 ... 1 1 1]

*(handwritten annotations: feature 0 ← 0 if "no", 1 if "yes"; feature 1 ... feature D; 2 2 2 ... 2)*

$2^{D+1}$

$2^8 = 256$
$2^{30} = 1{,}073{,}741{,}824$
$2^{1000} = 1.071509 \times 10^{301}$
$2^{100B} = $ HUGE!!!!!!

**Typically, computationally infeasible**

22/12 2020

# Greedy algorithm

## Forward stepwise algorithm

1. Pick a dictionary of features $\{h_0(\mathbf{x}),...,h_D(\mathbf{x})\}$
   - e.g., polynomials for linear regression
2. Greedy heuristic:
   i. Start with empty set of features $F_0 = \varnothing$
      (or simple set, like just $h_0(\mathbf{x})=1$ → $y_i = w_0 + \varepsilon_i$)
   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

22/12 2020

# Visualizing greedy algorithm

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

22/12 2020

# Visualizing greedy algorithm

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

22/12 2020

# Visualizing greedy algorithm

*Notice… it is suboptimal .*
*Adding next best thing, fit is nested now.*

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

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

0  1  2  3  4  5  6  7  8

# of features

22/12 2020

# Visualizing greedy algorithm

error never increases
+
solutions eventually meet

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

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

# of features

0  1  2  3  4  5  6  7  8

22/12 2020

# Complexity of forward stepwise

How many models were evaluated?

– 1st step, D models

– 2nd step, D-1 models (add 1 feature out of D-1 possible)

– 3rd step, D-2 models (add 1 feature out of D-2 possible)

– ...

How many steps?

- Depends

- At most D steps (to full model)

$$O(D^2) << 2^D$$
for large D

22/12 2020

# Other greedy algorithms

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.*

22/12 2020

# Using regularisation for features selection

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

22/12 2020

# Thresholding ridge coefficients?

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

# Thresholding ridge coefficients?

Selected features for a given threshold value

# Thresholding ridge coefficients?

Let's look at two related features...

Nothing measuring bathrooms was included!

22/12 2020

# Thresholding ridge coefficients?

## 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.

22/12 2020

# Thresholding ridge coefficients?

Would have included bathrooms in selected model

Can regularization lead directly to sparsity?

22/12 2020

# Try this cost instead of ridge …

Total cost =

measure of fit + $\lambda$ measure of magnitude of coefficients

RSS(**w**)

$$\|\mathbf{w}\|_1 = |w_0| + \ldots + |w_D|$$

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

Leads to
**sparse**
solutions!

22/12 2020

# Lasso regression

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

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

tuning parameter = balance of fit and **sparsity**

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

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

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

22/12 2020

# Coefficient path: ridge

22/12 2020

# Coefficient path: lasso

# NONPARAMETRIC REGRESSION

22/12 2020

# Fit globaly vs fit locally

**Parametric models**



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

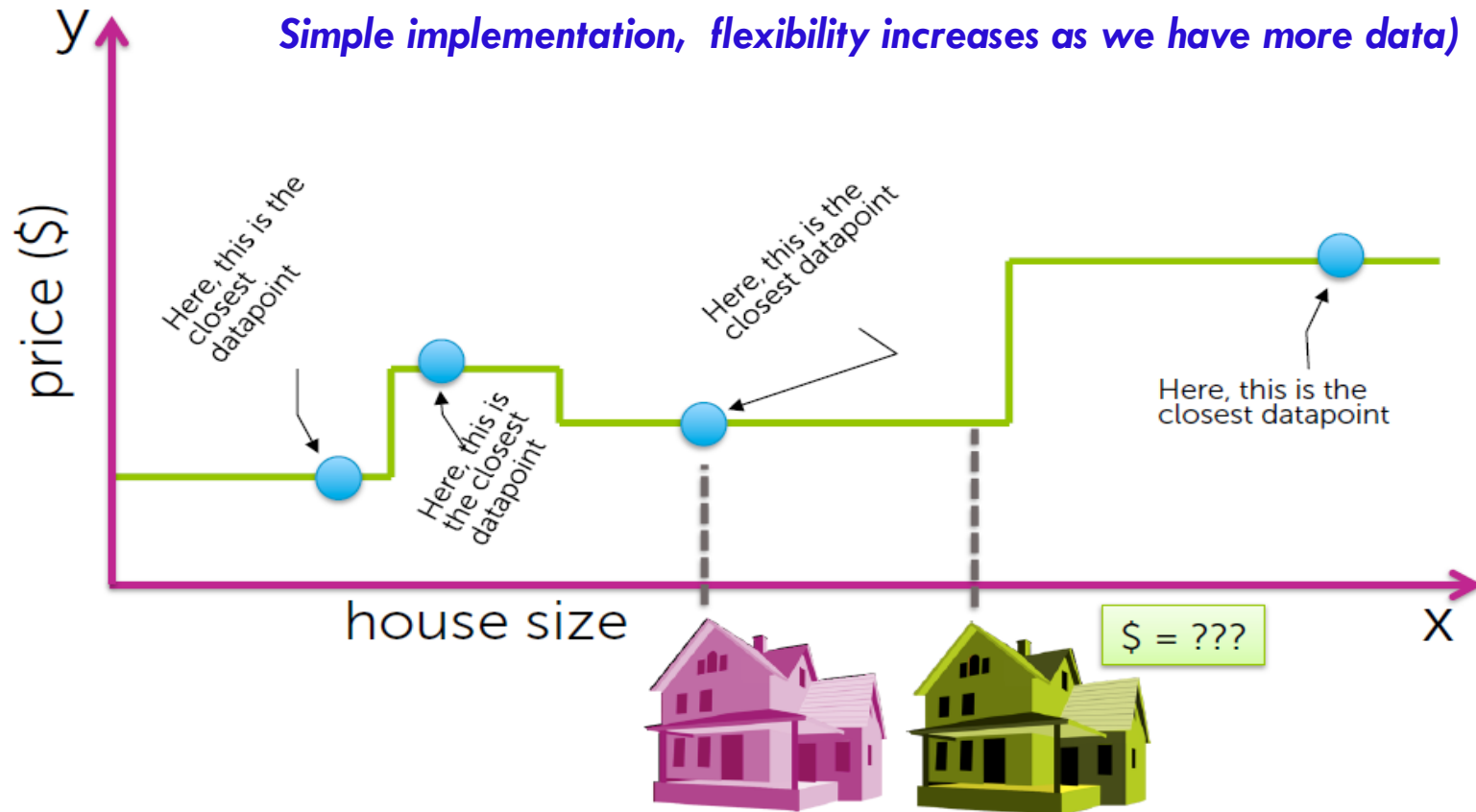linear     constant

quadratic

# What alternative do we have?

If we:

– Want to allow flexibility in f(**x**) having local structure

– Don't want to infer "structural breaks"

What's a simple option we have?

– Assuming we have plenty of data...

22/12 2020

# Nearest Neighbor & Kernel Regression (nonparametric approach)



Simple implementation, flexibility increases as we have more data)

22/12 2020

# Fit locally to each data point

Predicted value = "closest" $y_i$

1 nearest neighbor (1-NN) regression

Here, this is the closest datapoint

Here, this is the closest datapoint

Here, this is the closest datapoint

Here, this is the closest datapoint

y

price ($)

sq.ft.

x

22/12 2020

# What people do naturally…

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

$ = ???

$ = 850k

22/12 2020

# 1-NN regression more formally

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

Query point: $\mathbf{x}_q$ ← 🏠 $ ?

*big lime green house*

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

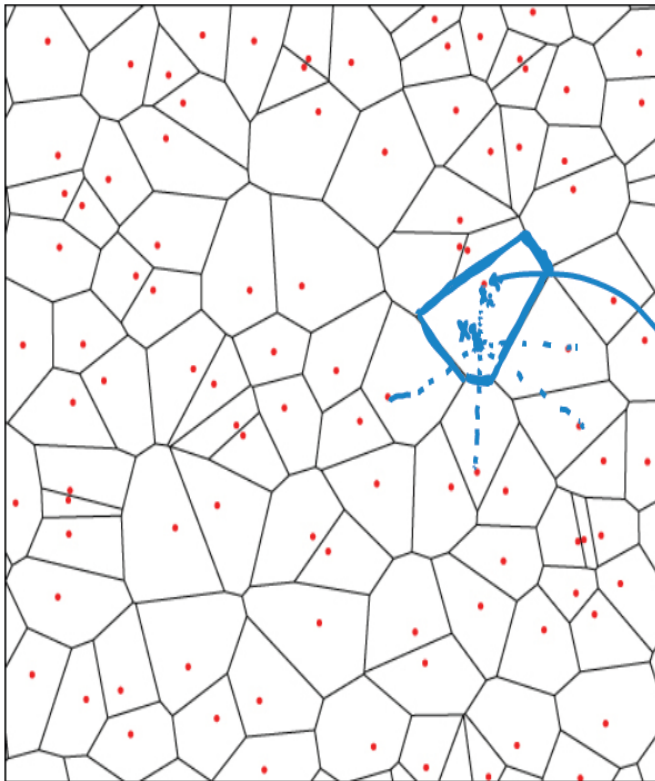$$X_{NN} \leftarrow \min_i \ distance(X_i, X_q)$$

*big pink house*

2. Predict

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

*sales price of big pink house*

**Transition point**

y

price ($)

Here, this is the closest datapoint

Here, this is the closest datapoint

$y_{NN}$

Here, this is the closest datapoint

$\hat{y}_q$

Here, this is the closest datapoint

$X_{NN}$

sq.ft.

*pink house*

$x_q$ *green house*

X

# Visualizing 1-NN in multiple dimensions

## Voronoi tesselation (or diagram):

- Divide space into N regions, each containing 1 datapoint

- Defined such that any **x** in region is "closest" to region's datapoint

observation

$X_q$ closer to $X_i$ than any other $X_j$ for $j \neq i$.

**Don't explicitly form!**

22/12 2020

# Distance metrics: Notion of „closest"

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

22/12 2020

# Weighting housing inputs

Some inputs are more relevant than others

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

22/12 2020

# Scaled Euclidan distance

Formally, this is achieved via

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

weight on each input
(defining relative importance)

Other example distance metrics:
– Mahalanobis, rank-based, correlation-based, cosine similarity, Manhattan, Hamming, …

22/12 2020

# Different distance metrics

Euclidean distance

Manhattan distance

22/12 2020

# Performing 1-NN search

- Query house:

- Dataset:

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

22/12 2020

# 1-NN algorithm

closest house

Initialize **Dist2NN** $= \infty$, 🏠 $= \varnothing$

For i=1,2,...,N

query house

  Compute: $\delta$ = **distance**(🏠$_i$ , 🏠$_q$)

   If $\delta$ < **Dist2NN**

   set    🏠   🏠$_i$

   set **Dist2NN** $= \delta$

Return most similar house 🏠

closest house
to query house

22/12 2020

# 1-NN in practice

*1–NN fit*

*function*

Fit looks good for data dense in x and low noise

Not great at interpolating over large regions...
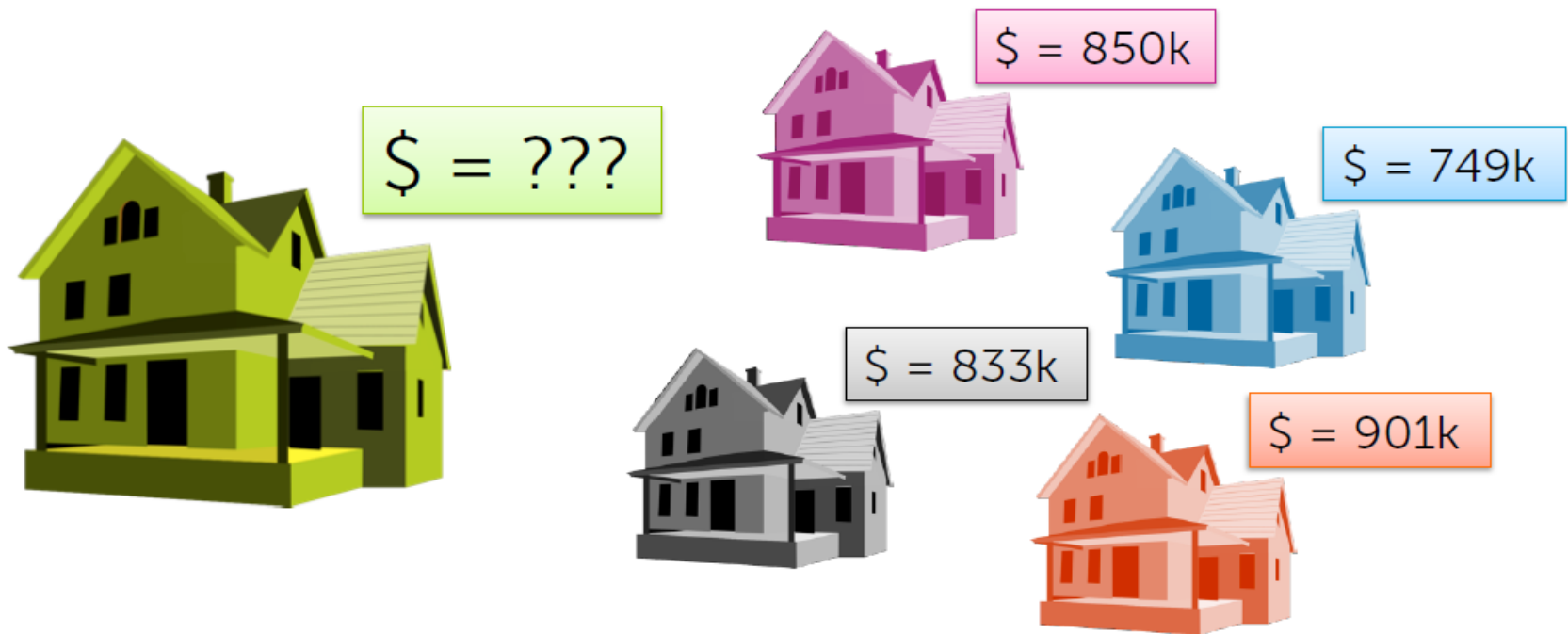
Fits can look quite wild... Overfitting?

*1-NN  sensitive to noise in the data*

22/12 2020

# Get more „comps"

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

$ = ???

$ = 850k

$ = 749k

$ = 833k

$ = 901k

22/12 2020

# K-NN regression more formally

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

Query point: $\mathbf{x}_q$

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

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

2. Predict

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

22/12 2020

# K-NN more formally

- Query house: 

- Dataset: 

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



22/12 2020

# K-NN algorithm

Initialize **Dist2kNN** = sort($\delta_1, ..., \delta_k$) ← list of sorted distances

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

= sort( 🏠 , ..., 🏠$_1$ 🏠$_k$ ) ← list of sorted houses

For i=k+1,...,N

Compute: $\delta$ = **distance**( 🏠$_i$ , 🏠$_q$ )

*query house*

If $\delta$ < **Dist2kNN**[k]

find j such that $\delta$ > **Dist2kNN**[j-1] but $\delta$ < **Dist2kNN**[j]

remove furthest house and shift queue:

[j+🏠] = [j:k🏠]

*insert new NN*

**Dist2kNN**[j+1:k] = **Dist2kNN**[j:k-1]

set **Dist2kNN**[j] = $\delta$ and 🏠 = 🏠$_i$

Return k most similar houses 🏠 ← *closest houses to query house*

22/12 2020

# K-NN in practice

Nearest Neighbors Kernel (K = 30)

Much more reasonable fit in the presence of noise

*All k-NN for a specific red point*

Boundary & sparse region issues

# K-NN in practice

Nearest Neighbors Kernel (K = 30)

**Discontinuities!**
Neighbor either in or out

# Issues with discontinuities

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.

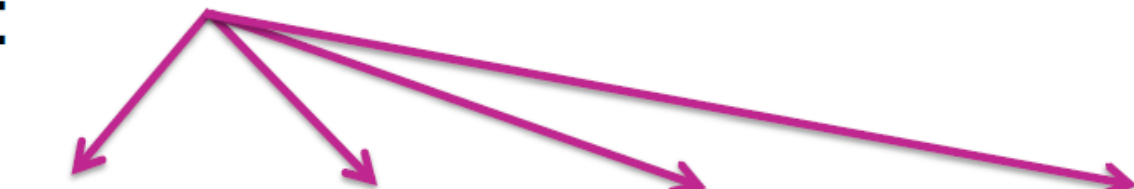– Don't really believe this type of fit

# Weighted k-NN

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

weights on NN

Predict:

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

22/12 2020

# How to define weights

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

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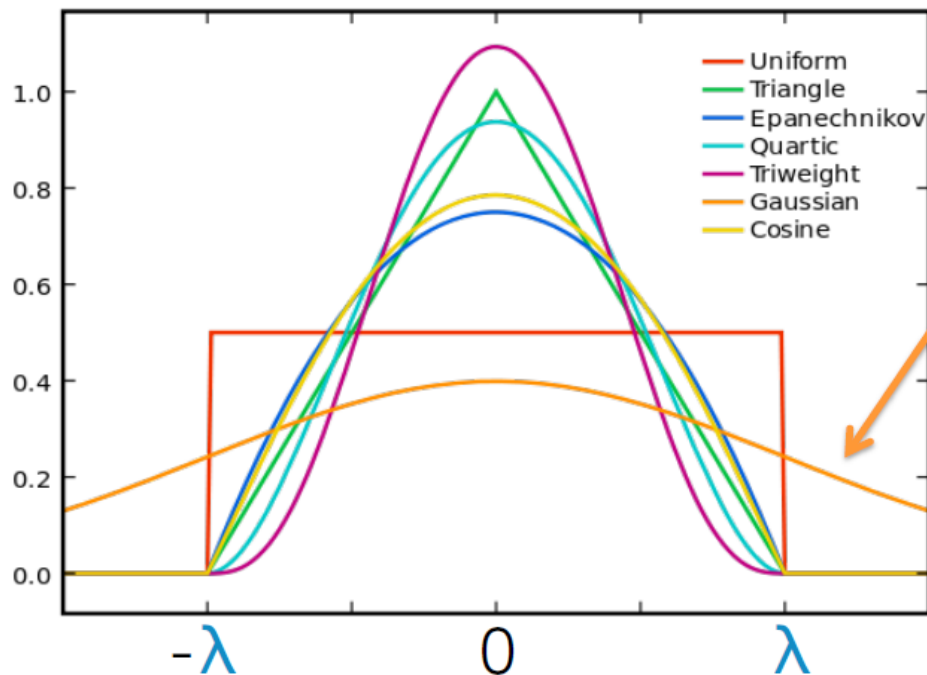distance($\mathbf{x}_{NNj}, \mathbf{x}_q$) small

Simple method:

$$c_{q\,NN_j} = \frac{1}{\text{distance}(x_j, x_q)}$$

22/12 2020

# Kernel weights for d=1

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

**simple isotropic case**



**Gaussian kernel:**

$$Kernel_\lambda(|x_i-x_q|) = \exp(-(x_i-x_q)^2/\lambda)$$

**Note:** never exactly 0!

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

22/12 2020

# Kernel regression

Nadaraya–Watson
kernel weighted average

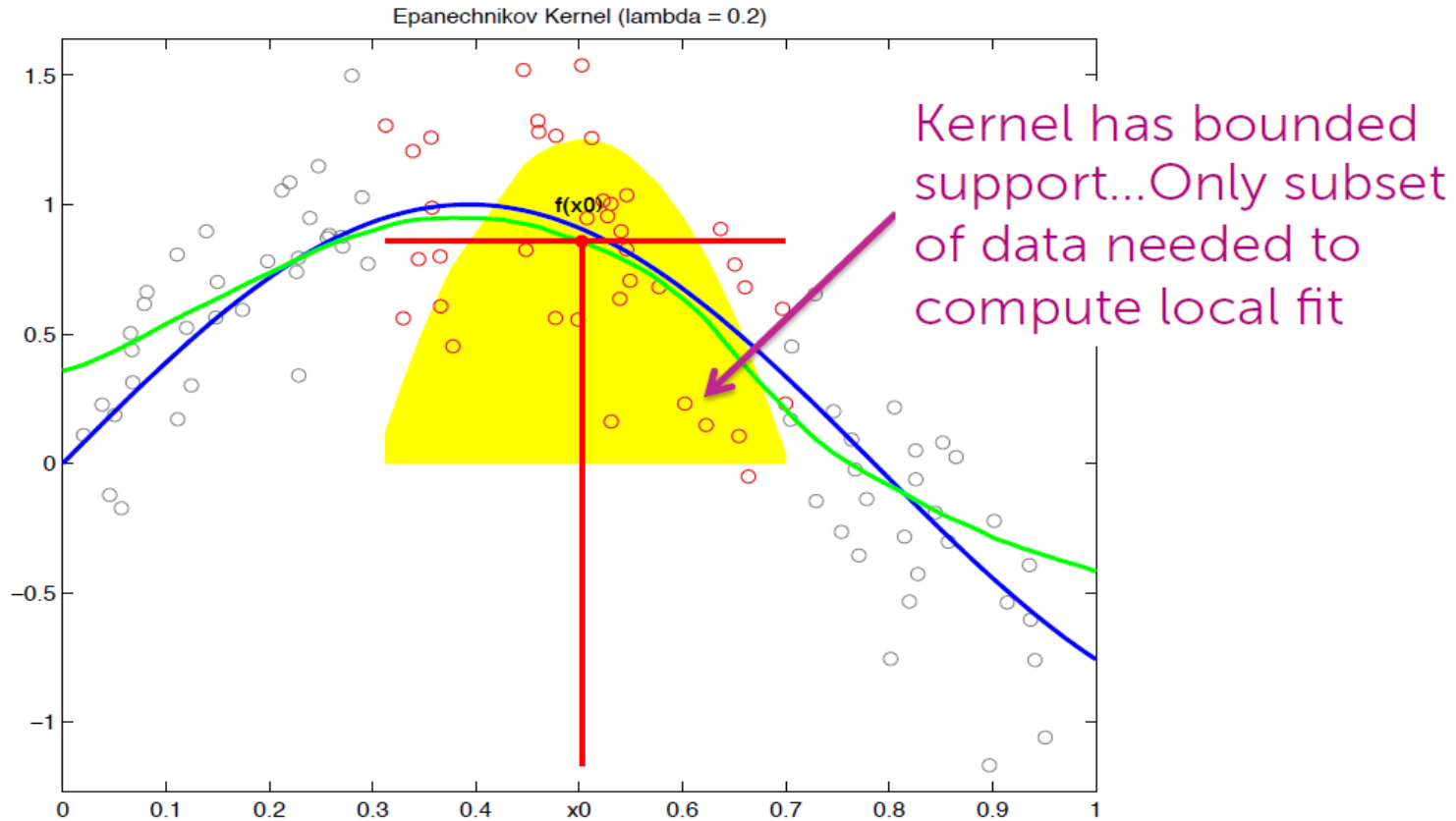Instead of just weighting NN, weight *all* points

Predict:

weight on each datapoint

$$\hat{y}_q = \frac{\sum_{i=1}^{N} c_{qi} y_i}{\sum_{i=1}^{N} c_{qi}} = \frac{\sum_{i=1}^{N} \text{Kernel}_\lambda(\text{distance}(\mathbf{x}_i, \mathbf{x}_q)) * y_i}{\sum_{i=1}^{N} \text{Kernel}_\lambda(\text{distance}(\mathbf{x}_i, \mathbf{x}_q))}$$
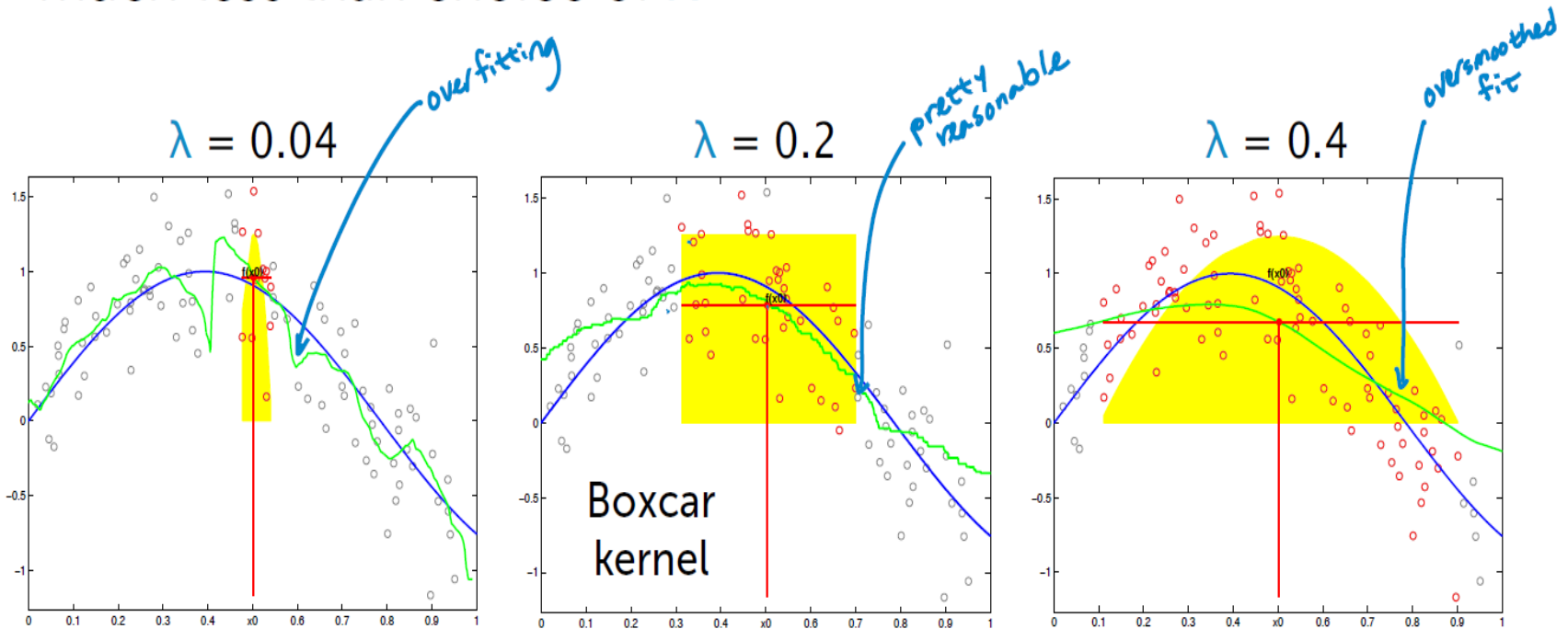
22/12 2020

# Kernel regression in practice

Epanechnikov Kernel (lambda = 0.2)

f(x0)

Kernel has bounded support...Only subset of data needed to compute local fit

# Choice of bandwith λ

Often, choice of kernel matters
much less than choice of λ



λ = 0.04    *overfitting*

λ = 0.2    *pretty reasonable*

λ = 0.4    *oversmoothed fit*

Boxcar kernel

22/12 2020

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

How to choose?  Same story as always...

Cross Validation

# Contrasting with global average

A **globally constant fit** weights all points equally

equal weight on each datapoint

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



Boxcar Kernel (lambda = 1)

22/12 2020

## Kernel regression leads to locally constant fit

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

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



45

# Local linear regression

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

- 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

22/12 2020

# Nonparametric approaches

k-NN and kernel regression are examples
of nonparametric regression

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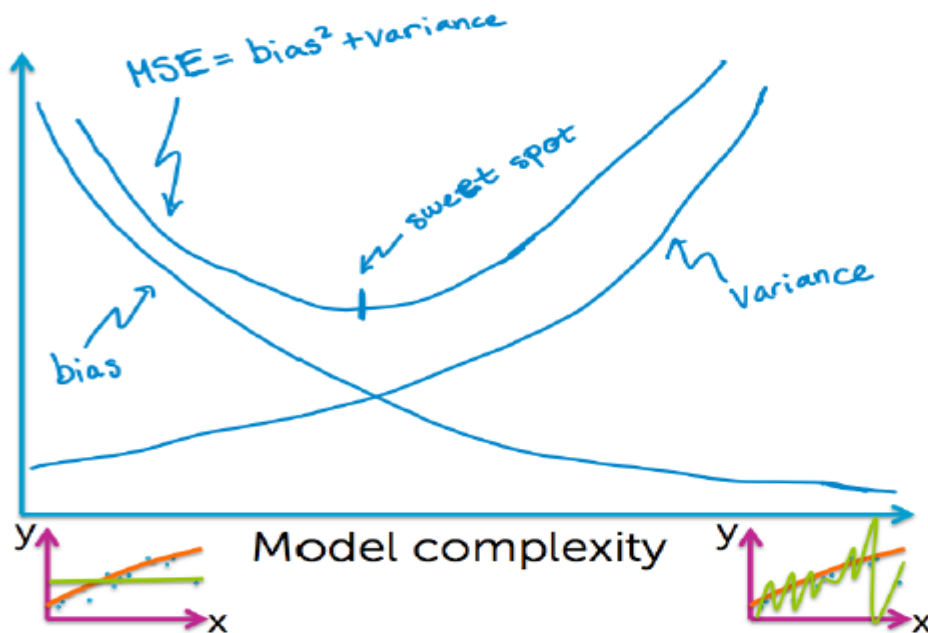- Complexity can grow with the number of observations N

Lots of other choices:
- Splines, trees, locally weighted structured regression models...

22/12 2020

## 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
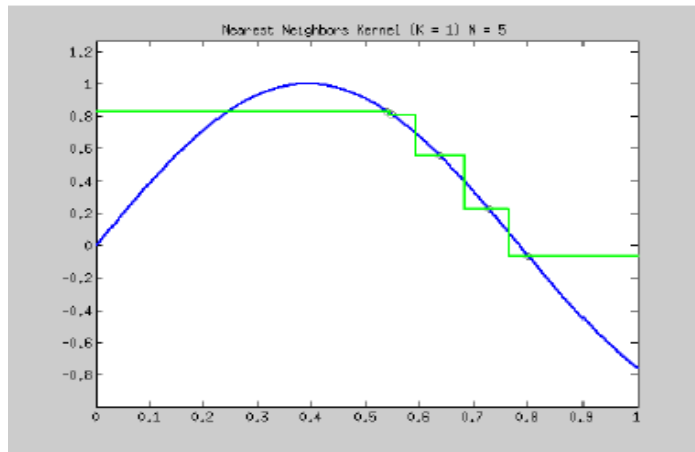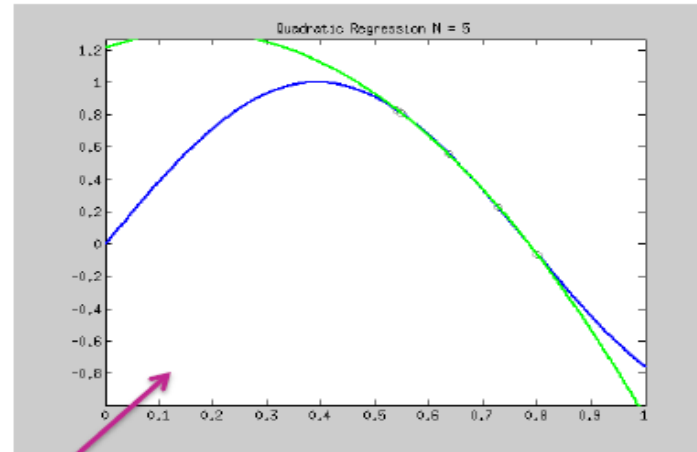


22/12 2020

# Limiting behaviour of NN

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

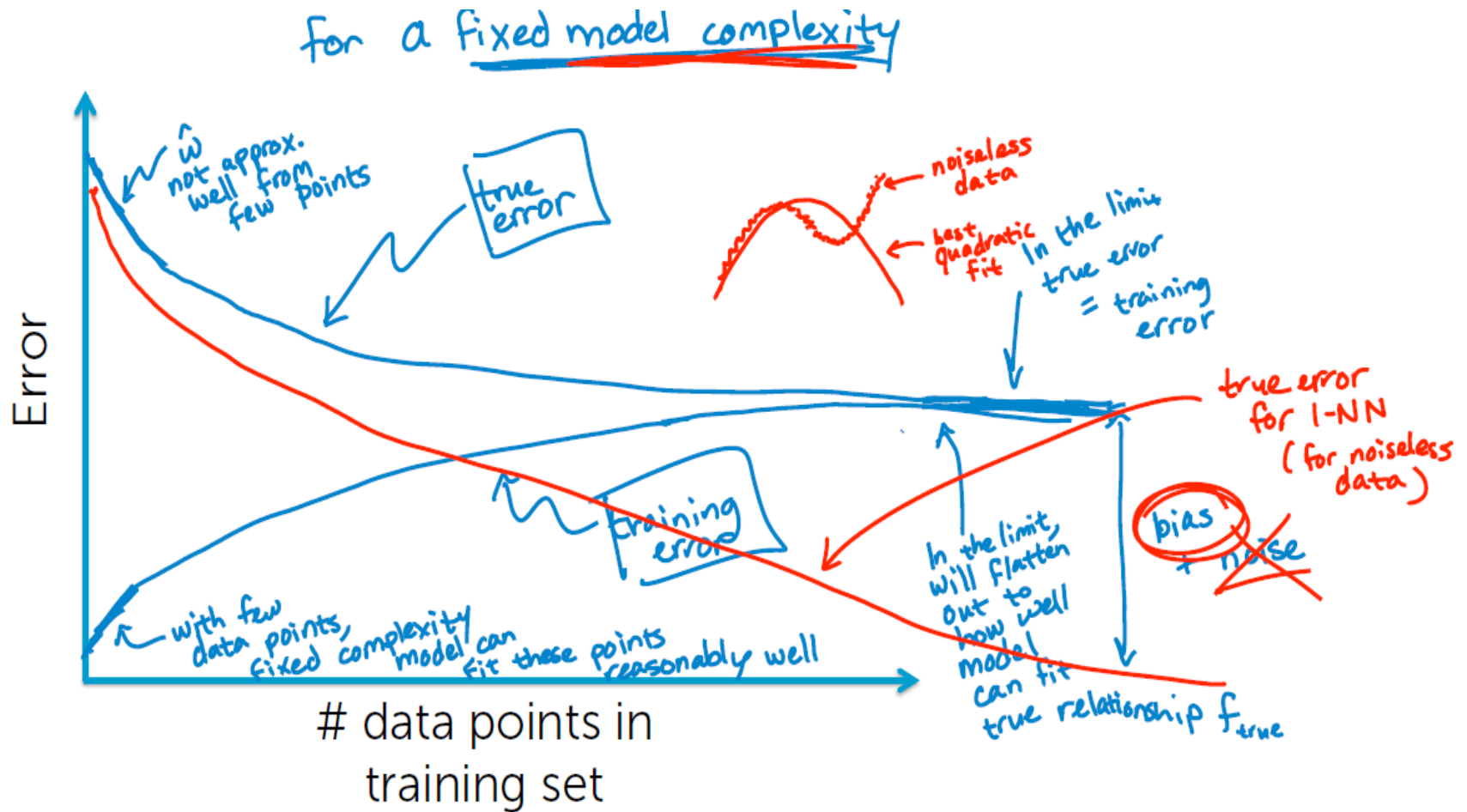In the limit of getting an infinite amount of noiseless data, the MSE of 1-NN fit goes to 0



1-NN fit



Quadratic fit

Not true for parametric models!

# Error vs amount of data
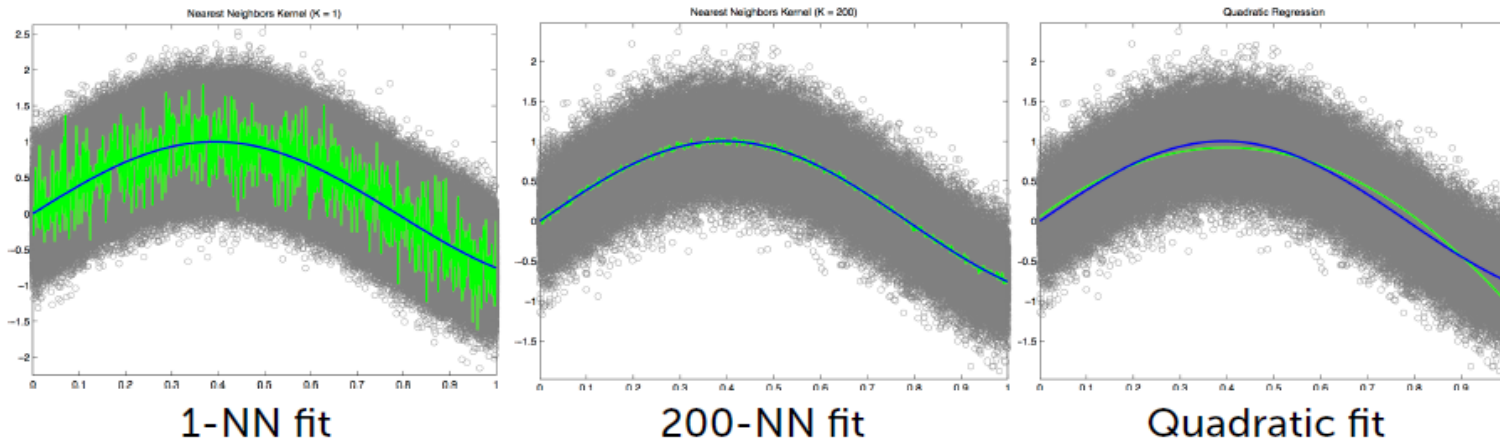
for a fixed model complexity

Error

# data points in training set

$\hat{w}$ not approx. well from few points

true error

noiseless data

best quadratic fit

In the limit true error = training error

true error for 1-NN (for noiseless data)

training error

with few data points, fixed complexity model can fit these points reasonably well

In the limit, will flatten out to how well model can fit true relationship $f_{true}$

bias + noise

22/12 2020

## 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



| 1-NN fit | 200-NN fit | Quadratic fit |

# Issues: NN and kernel methods

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

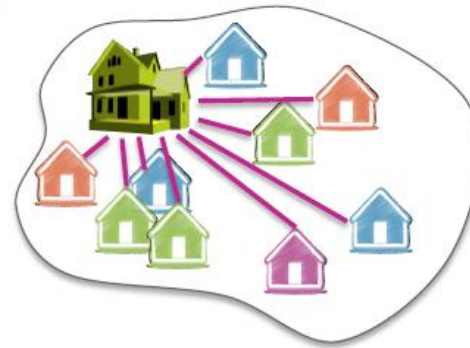This is where parametric models become useful...

# Issues: Complexity of NN search

Naïve approach: **Brute force search**
- Given a query point $x_q$
- Scan through each point $x_1, x_2, ..., x_N$
- O(N) distance computations per 1-NN query!
- O(Nlogk) per k-NN query!

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

Will talk more about efficient methods in
Clustering & Retrieval course

# Summarising

**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

22/12 2020