# DataStar Machine Learning

**ADAX**

# Modules

# Session 2: Regression

15 Sept 2017

**ADAX**

# What is Regression?

# What is Regression?



From features to predictions

Data → Regression → Intelligence

Input **x**:
features derived from data

Learn **x→y** relationship

Predict **y**:
continuous "output" or "response" to input

# Stock prediction

- Predict the price of a stock (y)
- Depends on **x** =
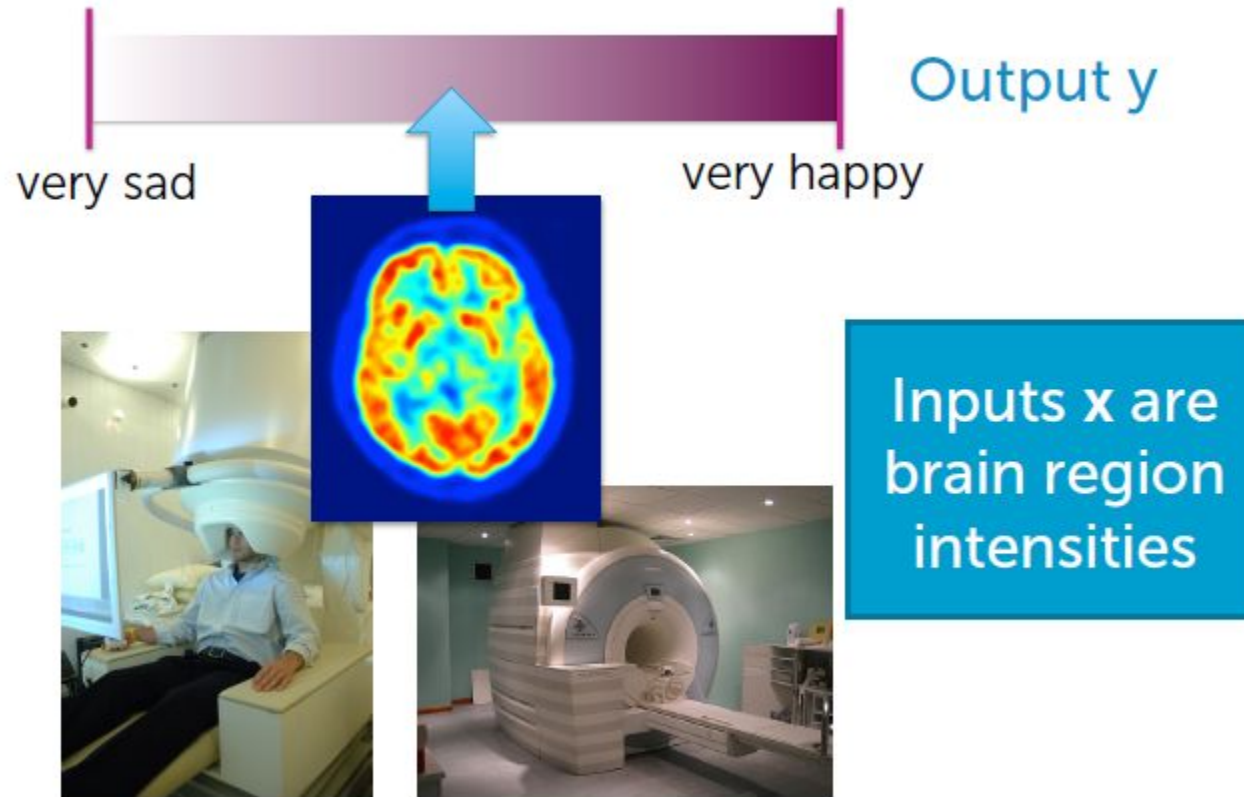  - Recent history of stock price
  - News events
  - Related commodities

# Tweet Popularity

- How many people will retweet your tweet? (y)
- Depends on $x$ = # followers,
  # of followers of followers,
  features of text tweeted,
  popularity of hashtag,
  # of past retweets,...

# Reading Your Mind



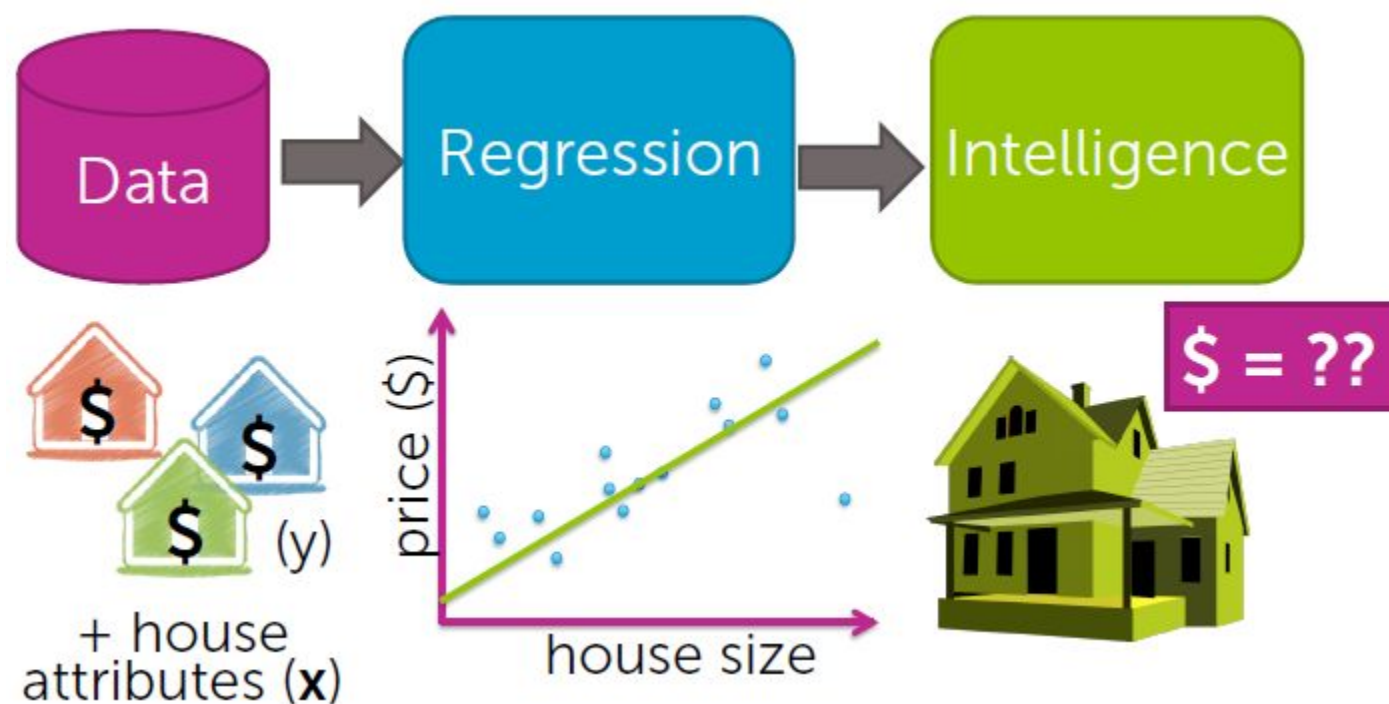Output y

very sad — very happy

Inputs **x** are brain region intensities

# Case Study: Predicting House Prices

# Simple Regression

Linear regression with one input

# How much is my house worth?

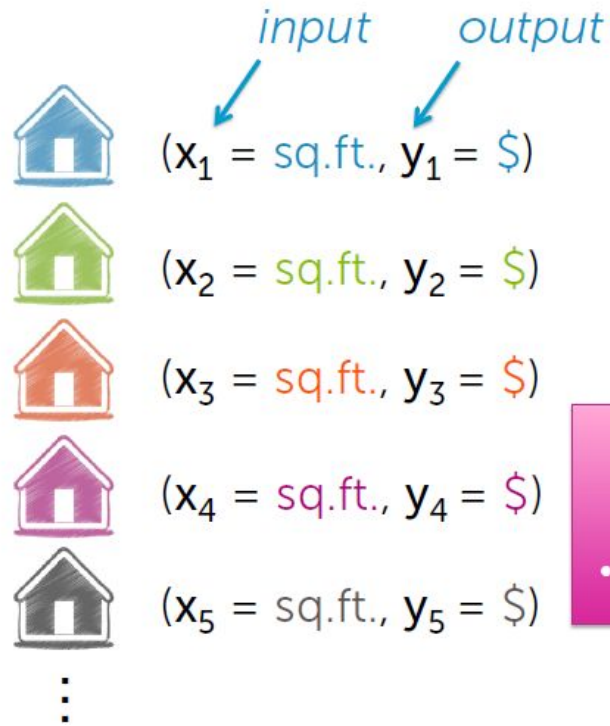# Look at recent sales in my neighborhood

- How much did they sell for?

# Regression fundamentals

## Data, Model, Task

**Data**



input    output

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

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

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

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

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

**Input vs. Output:**
- y is the quantity of interest
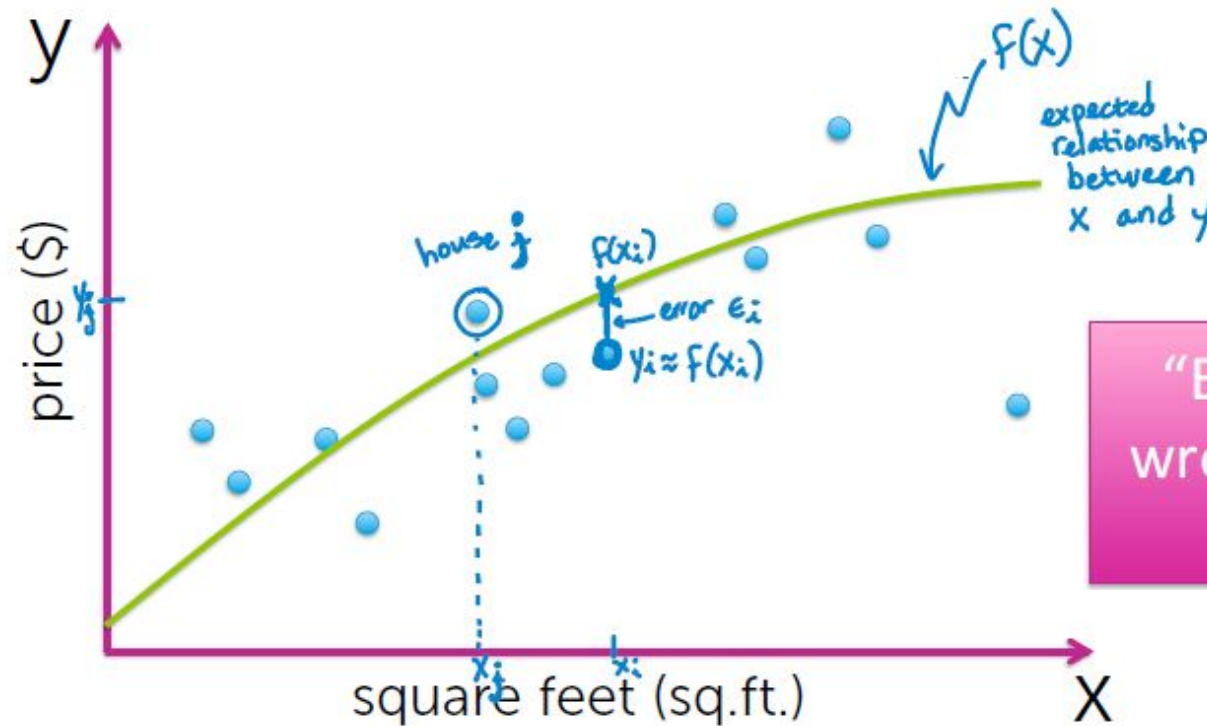- assume y can be predicted from x

# Regression fundamentals

## Data, Model, Task

**Model**

How we
assume the
world works



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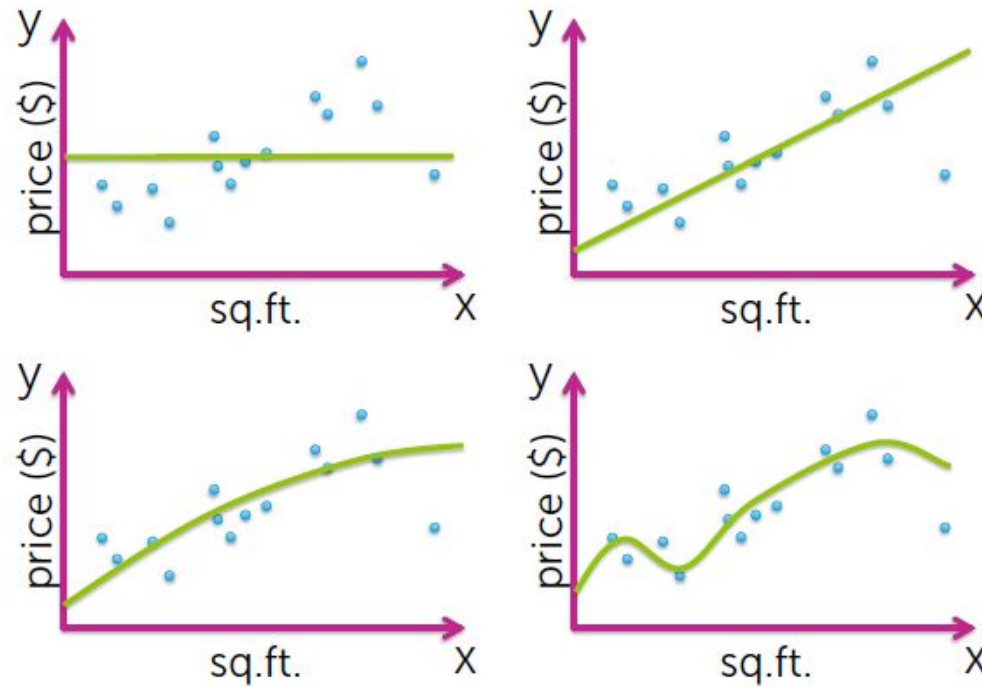George Box, 1987.

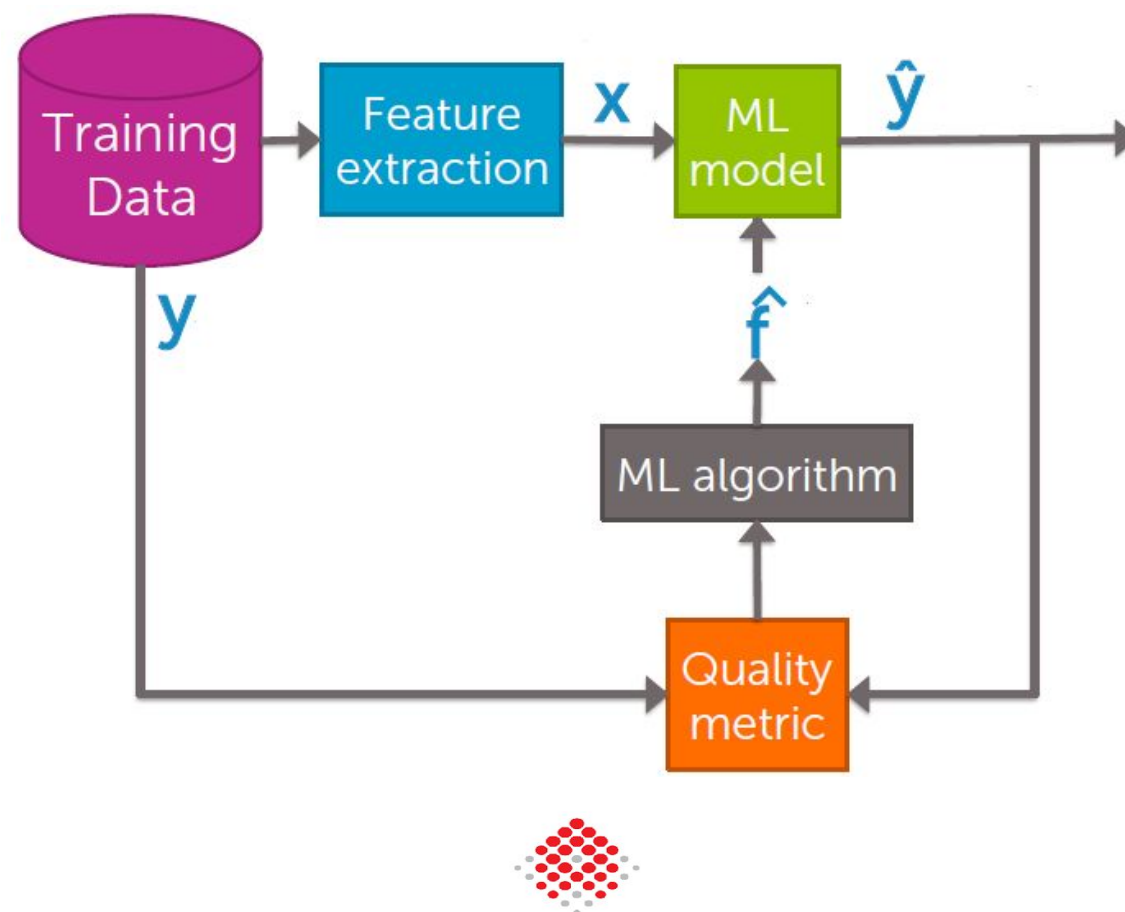# Regression fundamentals

**Data, Model, Task**

**Task**

Which model f(x) is the best fit?



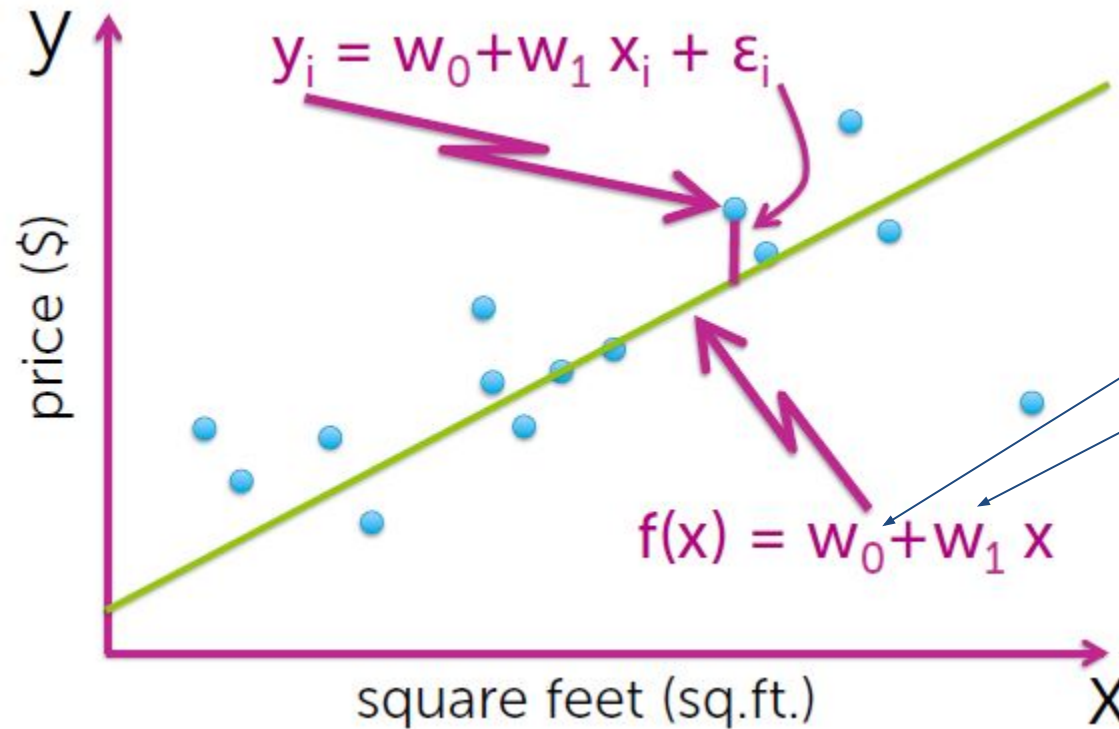From the model that we estimate, we can do a lot more other tasks

# How this works
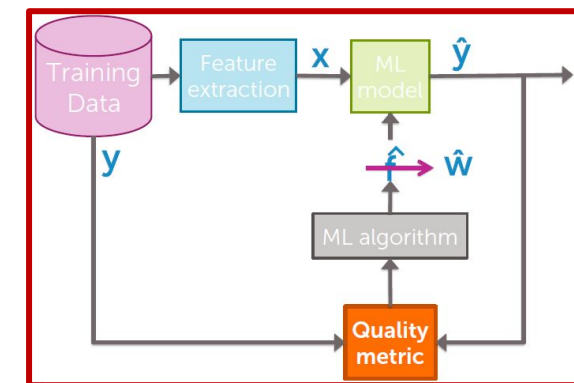
# Simple linear regression



What's the equation of a line?

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

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

Regression coefficients

# "Cost" of using a given line



Residual sum of squares (RSS)

What's this residual?

$RSS(w_0,w_1) =$
$(\$_{house\ 1}-[w_0+w_1sq.ft._{house\ 1}])^2$
$+ (\$_{house\ 2}-[w_0+w_1sq.ft._{house\ 2}])^2$
$+ (\$_{house\ 3}-[w_0+w_1sq.ft._{house\ 3}])^2$
$+ ...$[include all training houses]

price ($)

$(w_0+w_1 \cdot sq.ft.)$

square feet (sq.ft.)

# Find the "best" line



Minimize cost over all possible $w_0, w_1$

$RSS(w_0=1.1, w_1=0.8)$

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

$RSS(w_0=0.97, w_1=0.85)$

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

y

price ($)

square feet (sq.ft.)

x

# Model vs. Fitted line

Let's say
$w_0 = -44850$
$w_1 = 280.76$



$\hat{f}(x) = \hat{w}_0 + \hat{w}_1 x$

Regression model:

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

price ($)

square feet (sq.ft.)

# Seller: Predicting house price



$$\hat{f}(x) = \hat{w}_0 + \hat{w}_1 x$$

y

price ($)

square feet (sq.ft.)

X

Regression model:

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

Best guess of your house price:

$$\hat{y}_{house} = \hat{w}_0 + \hat{w}_1 \, sq.ft._{house}$$

# Buyer: Predicting size of house

# A concrete example



$\hat{f}(x) = -44850 + 280.76\,x$

Predict $ of 2,640 sq.ft. house:

$-44850 + 280.76 * 2,640$

$= \$696,356$

Predict sq.ft. of $859.000 sale:

$(859000+44850)/ 280.76$

$= 3,219$ sq.ft.

# Making sense of the coefficients



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

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

price ($)

square feet (sq.ft.)

**Note:**

y = w0 when x=0
Not very meaningful!

# Making sense of the coefficients



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

y

price ($)

predicted change in $

1 sq. ft.

square feet (sq.ft.)

X

$w_1$

Predicted change in the output per unit change in input

# What-ifs



$$\hat{f}(x) = -\$44{,}850 + 280.76\ (\$/\text{sq.ft.})\ x$$

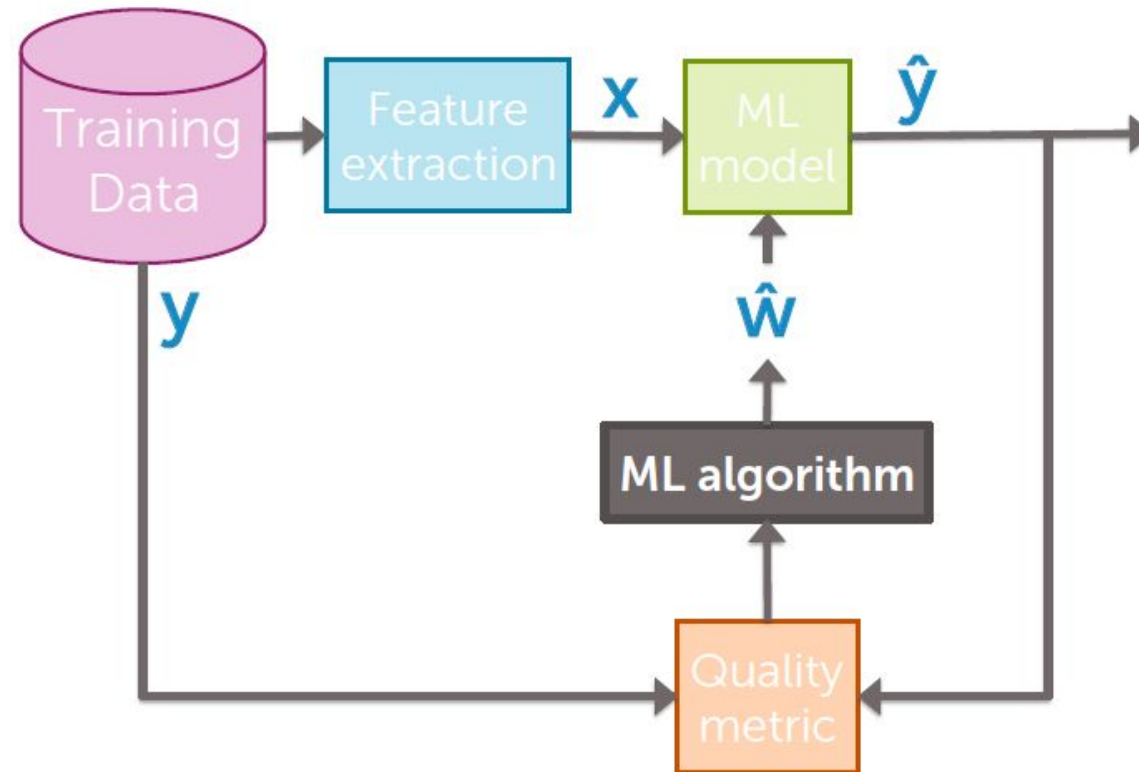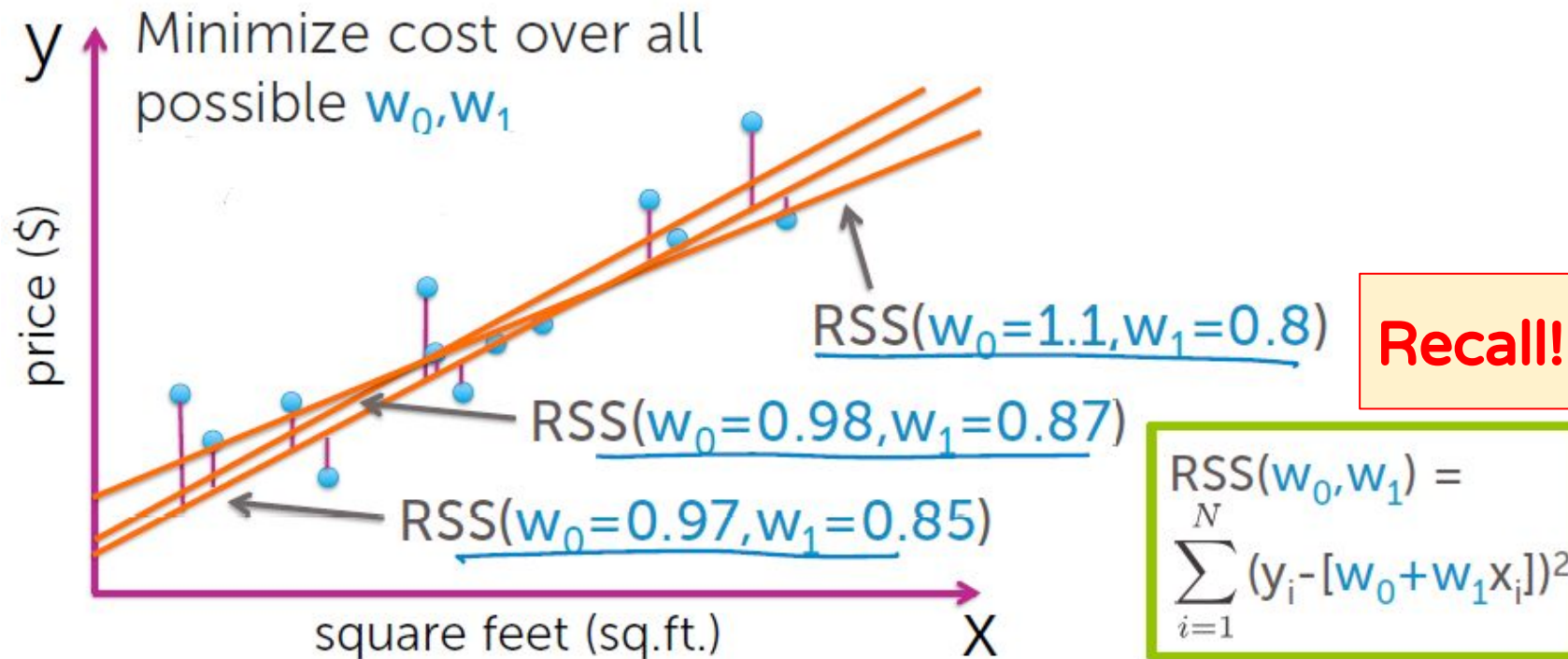**What if…**

1. What if house was measured in square meters?
2. Price was measured in RM? Euros?

# The Algorithm

# Find "best" line



Minimize cost over all possible $w_0, w_1$

y — price ($)
x — square feet (sq.ft.)

$RSS(w_0=1.1, w_1=0.8)$

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

$RSS(w_0=0.97, w_1=0.85)$

**Recall!**

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

# Find "best" line



Minimize cost over all possible $w_0, w_1$

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

**CONVEX**
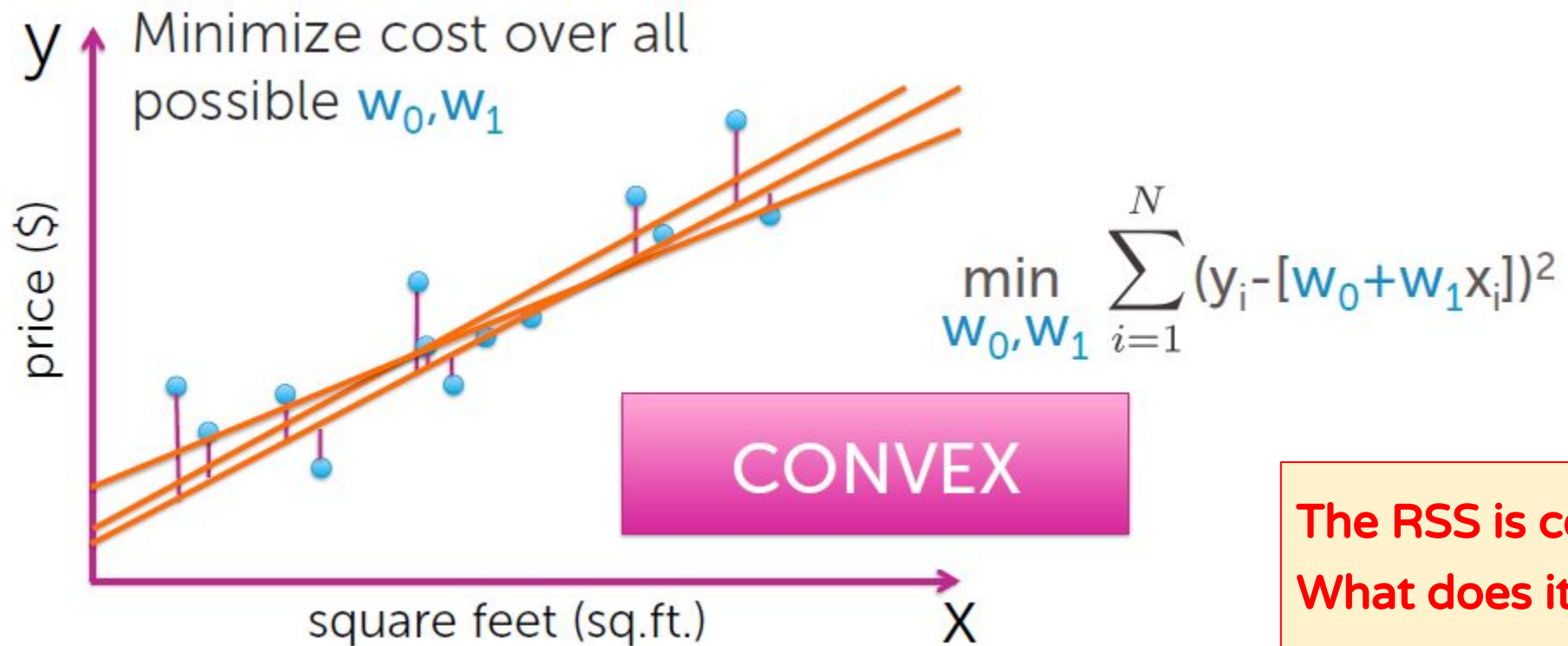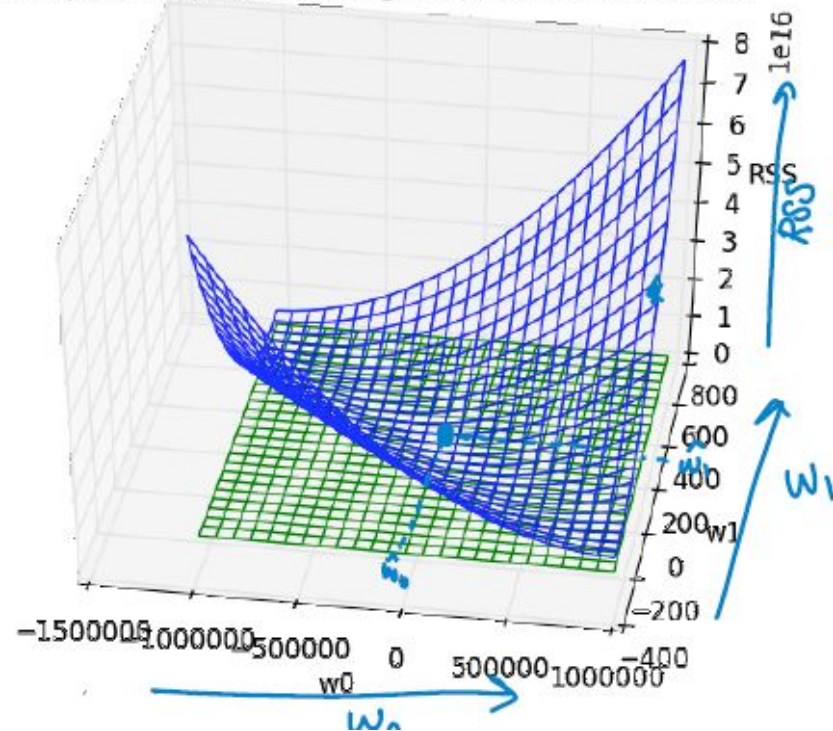
The RSS is convex.
What does it mean?

# Minimizing a cost function



3D plot of RSS with tangent plane at minimum
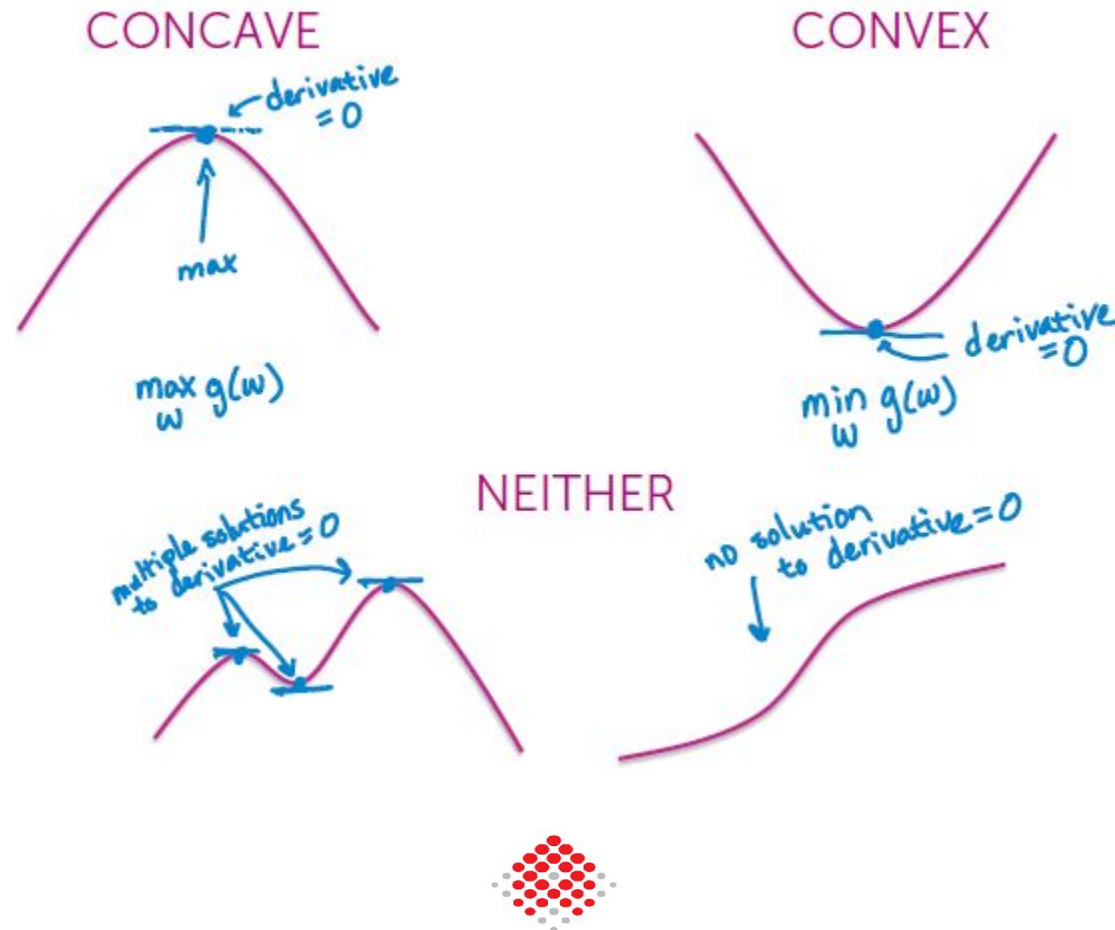
Minimize function over all possible $w_0, w_1$

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

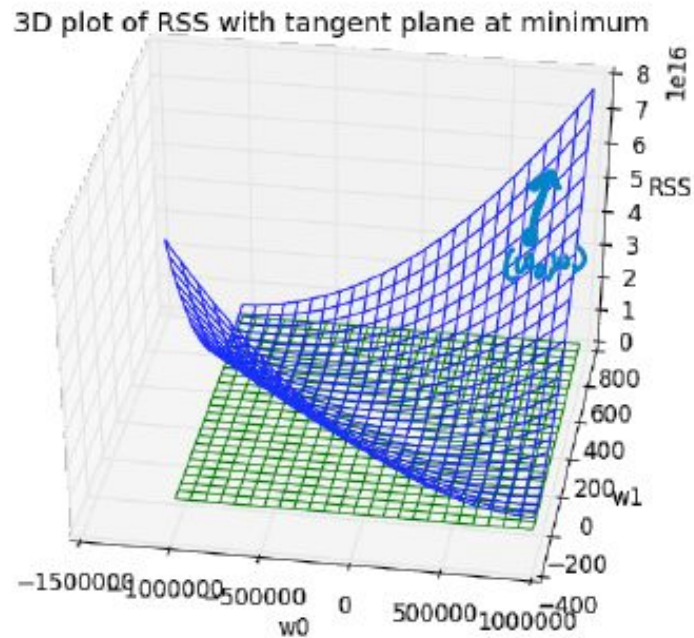RSS($w_0, w_1$) is a function of 2 variables

# The idea of gradients

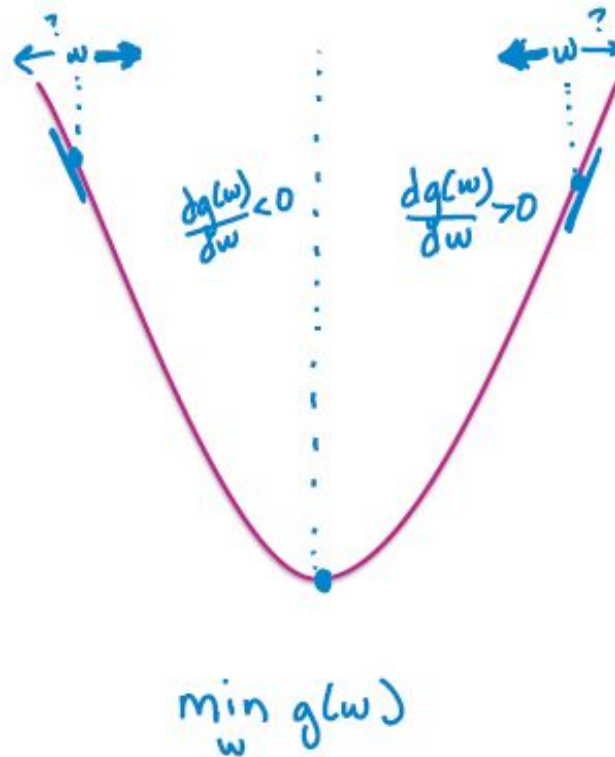# Gradient example



3D plot of RSS with tangent plane at minimum

$$g(w) = 5w_0 + 10 w_0 w_1 + 2w_1^2$$

$$\nabla g(w) =$$

# Descending the hill

What we want is to get w which minimizes the cost function g(w)

$\frac{dg(w)}{dw} < 0$  $\frac{dg(w)}{dw} > 0$

$\min_{w} g(w)$

Step size

While not converged

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \left.\frac{dg}{dw}\right|_{w^{(t)}}$$
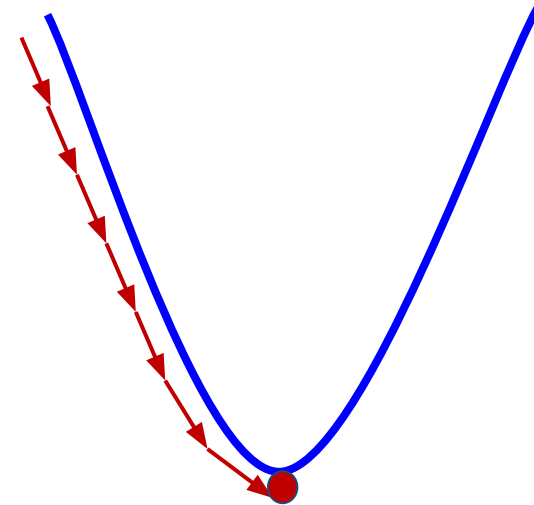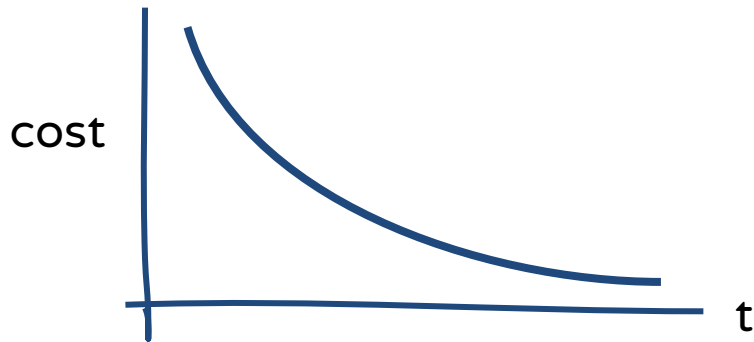
The value of w needs to be iteratively updated by "descending" the gradient ⇒ Gradient Descent

# Will it ever converge at 0?

For convex functions, optimum occurs when $\left| \dfrac{dg}{dw} \right| = 0$

In practice, stop when $\left| \dfrac{dg}{dw} \right| < \varepsilon$



Is a fixed step size desirable?

# Compute the gradient of RSS

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

**Taking the derivative w.r.t. $w_1$ and $w_2$**

Putting it together:

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

# Method 1: Set gradient = 0

$$\nabla RSS(w_0, w_1) = \begin{bmatrix} -2\sum_{i=1}^{N}[y_i - (w_0 + w_1 x_i)] \\ -2\sum_{i=1}^{N}[y_i - (w_0 + w_1 x_i)]x_i \end{bmatrix} = 0$$

**Can you solve these 2 equations to get a closed-form solution?**

# Method 2: Gradient Descent

Interpreting the gradient:

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

**So plug this back to the Gradient Descent formula**

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

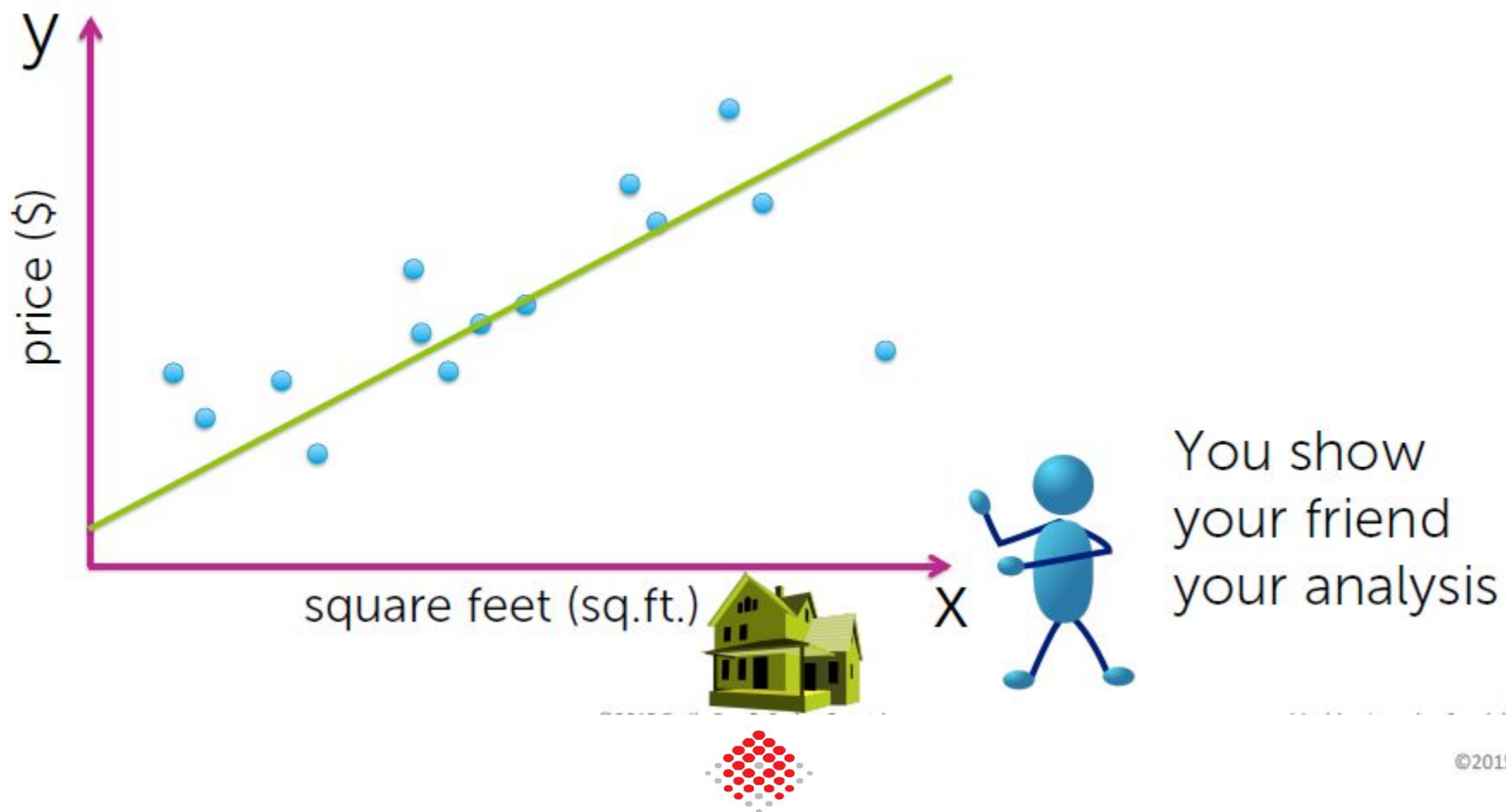**Gradient descent relies on choosing step size and convergence criteria**
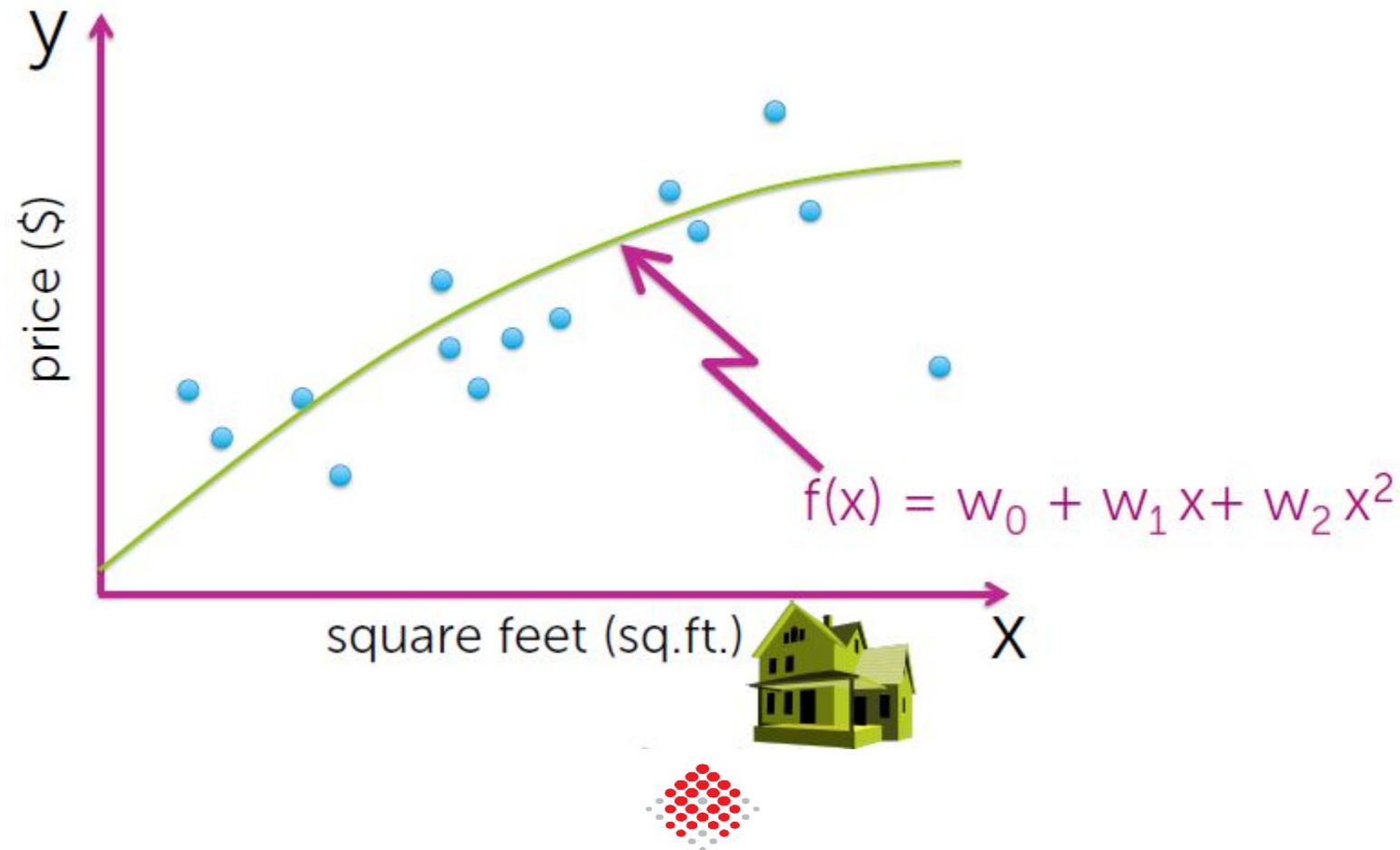
# Multiple Regression

Linear regression with multiple features
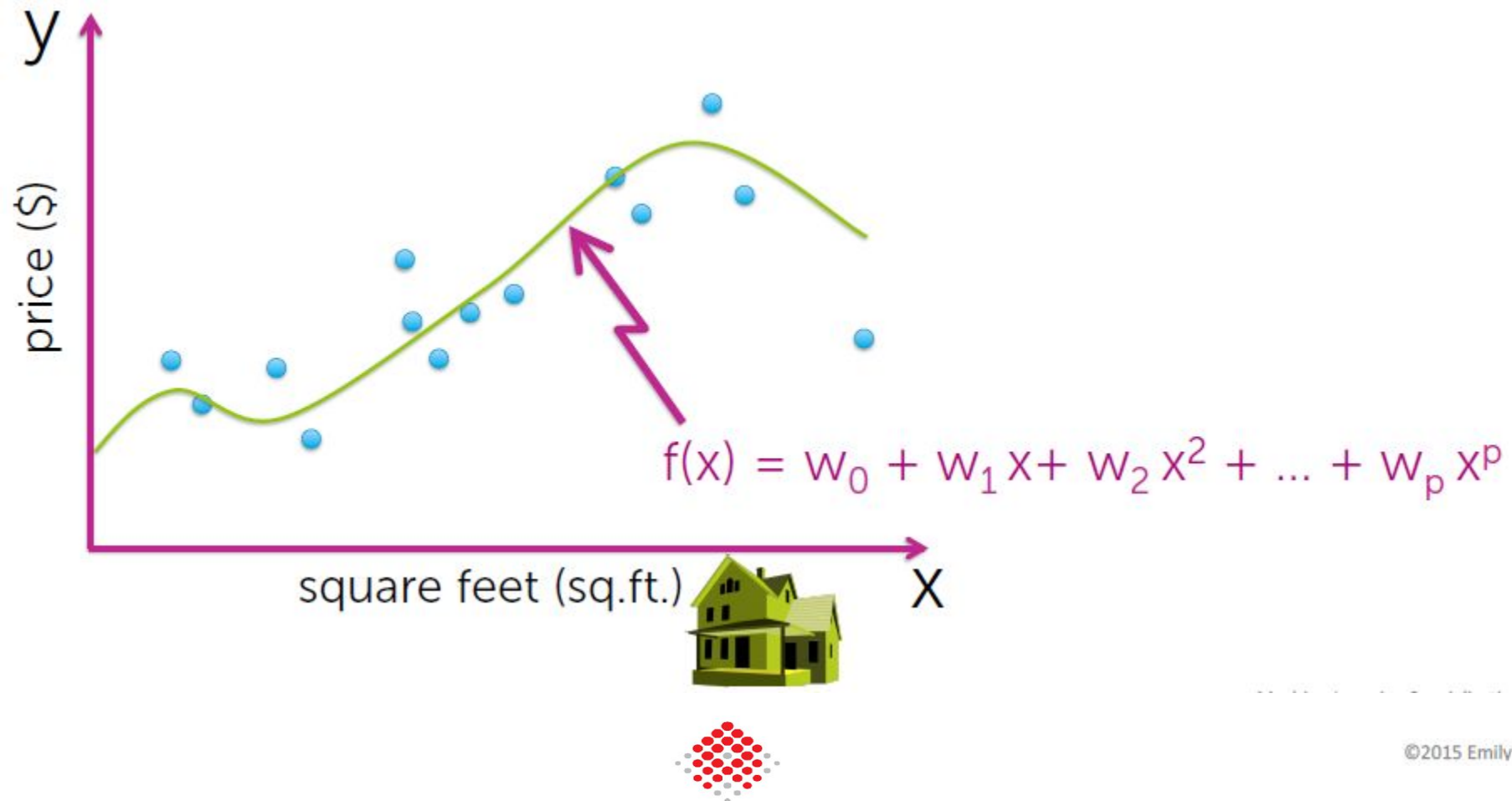
# Fit with a line or...?

# How about a quadratic function?



$$f(x) = w_0 + w_1 x + w_2 x^2$$

# Or even higher order polynomials?



$$f(x) = w_0 + w_1 x + w_2 x^2 + \ldots + w_p x^p$$

# 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)     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 = $w_p$

# Generic basis expansion

Model:

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

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

$j^{th}$ *feature*

$j^{th}$ *regression coefficient or weight*

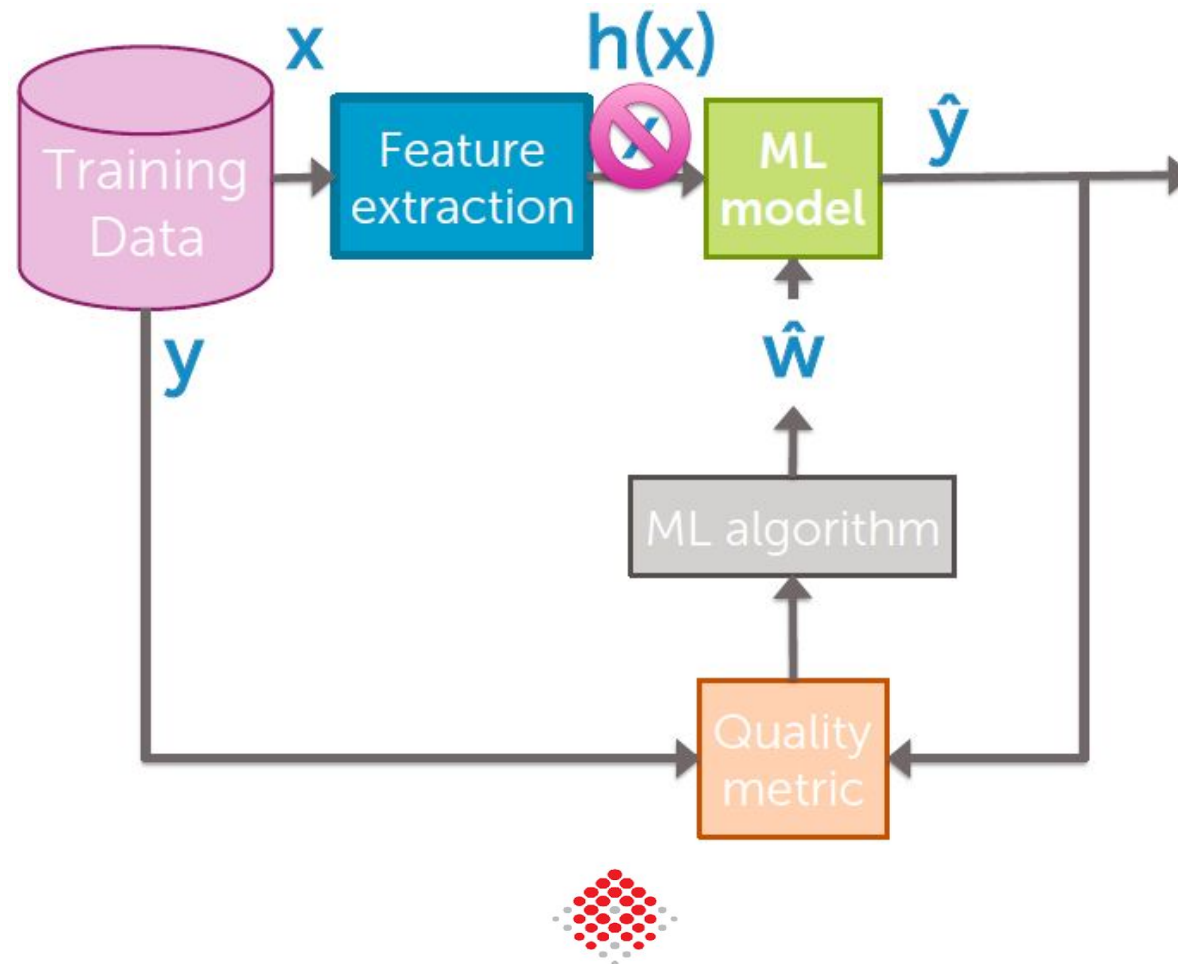feature 1 = $h_0(x)$...often 1 (constant)

feature 2 = $h_1(x)$... e.g., $x$
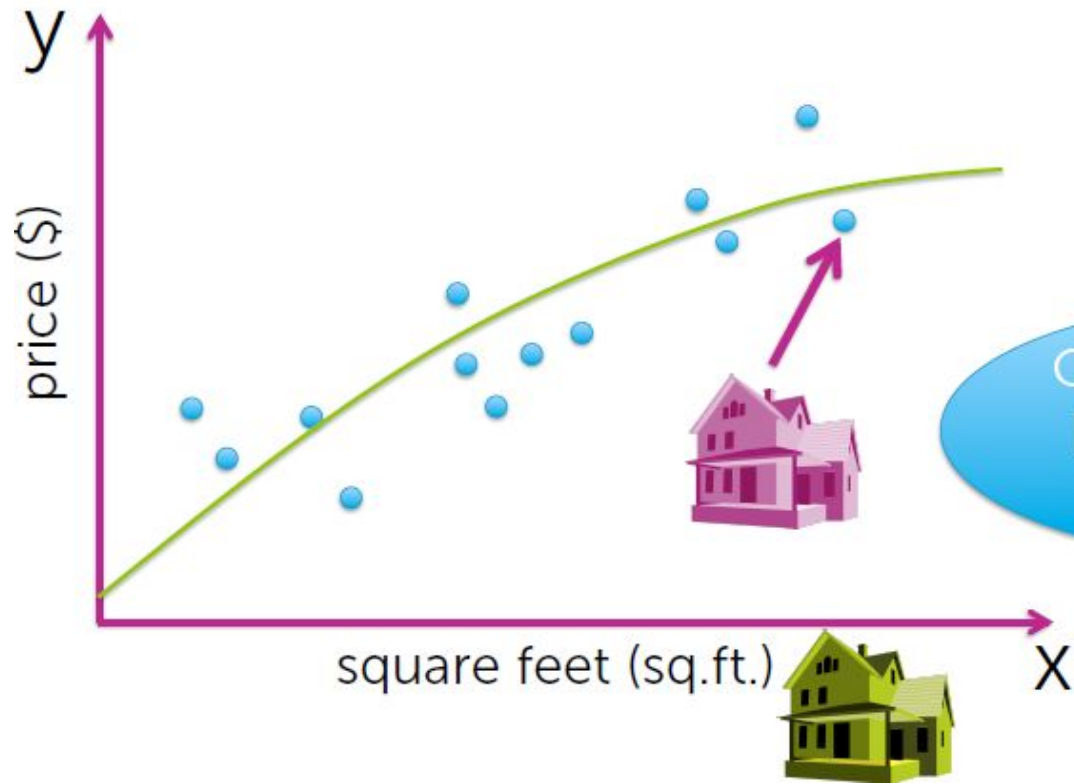
feature 3 = $h_2(x)$... e.g., $x^2$

...

feature D+1 = $h_D(x)$... e.g., $x^p$
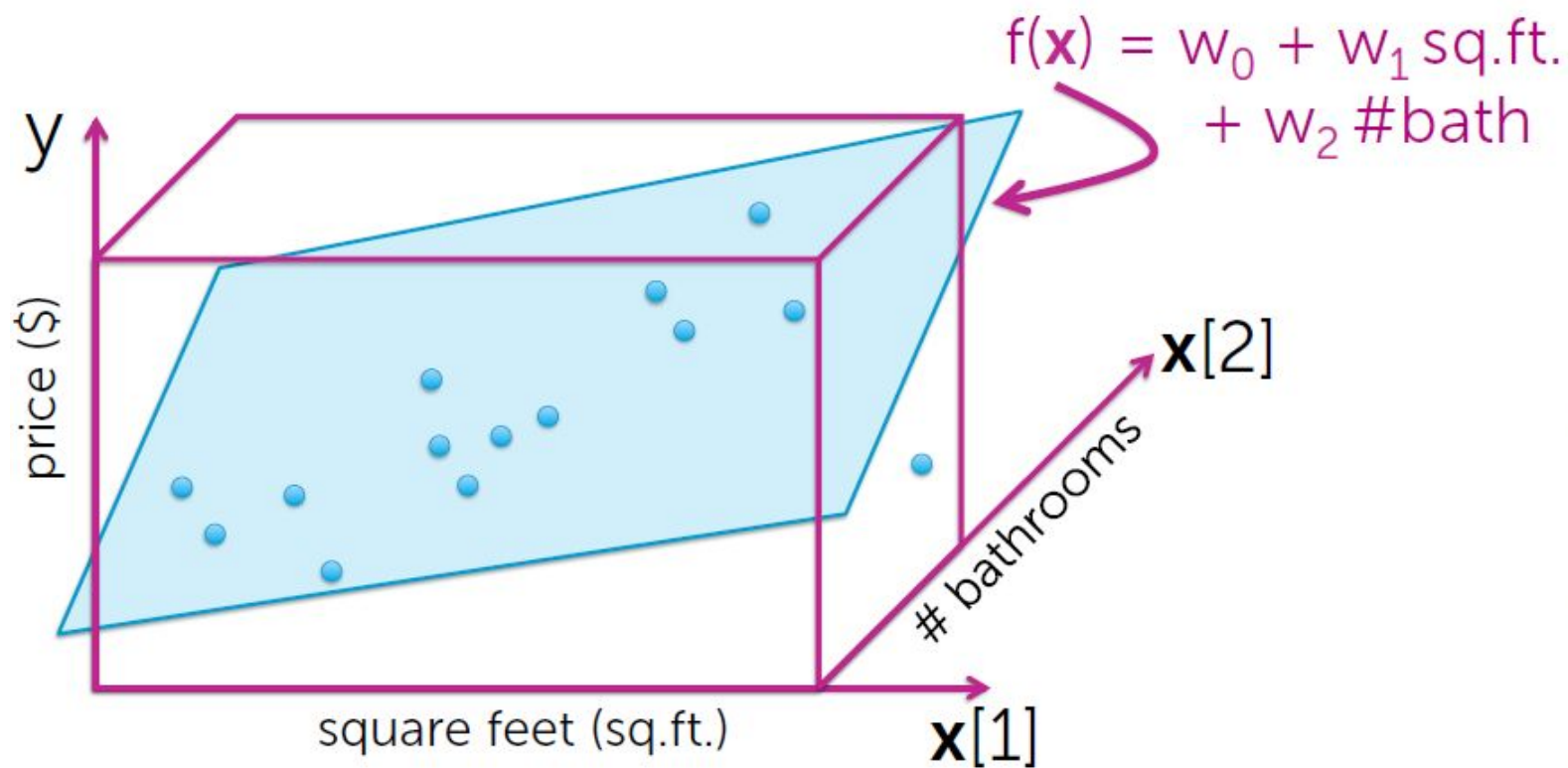
# No longer a single input

# Big house, but...



©2015 Emily Fox & Carlos Guestrin

# Add more inputs

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

**Any many more**
– Square feet
– # bathrooms
– # bedrooms
– Lot size
– Year built
– ...

price ($)

square feet (sq.ft.)

# bathrooms

$\mathbf{x}[2]$

$\mathbf{x}[1]$

y

# Planes and Hyperplanes



Model:

$$y_i = w_0 + w_1 x_i[1] + \ldots + w_d x_i[d] + \varepsilon_i$$

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

# Generically...a D-dimensional curve

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

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

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 $\mathbf{x}[1],\ldots, \mathbf{x}[d]$

# Some common notations

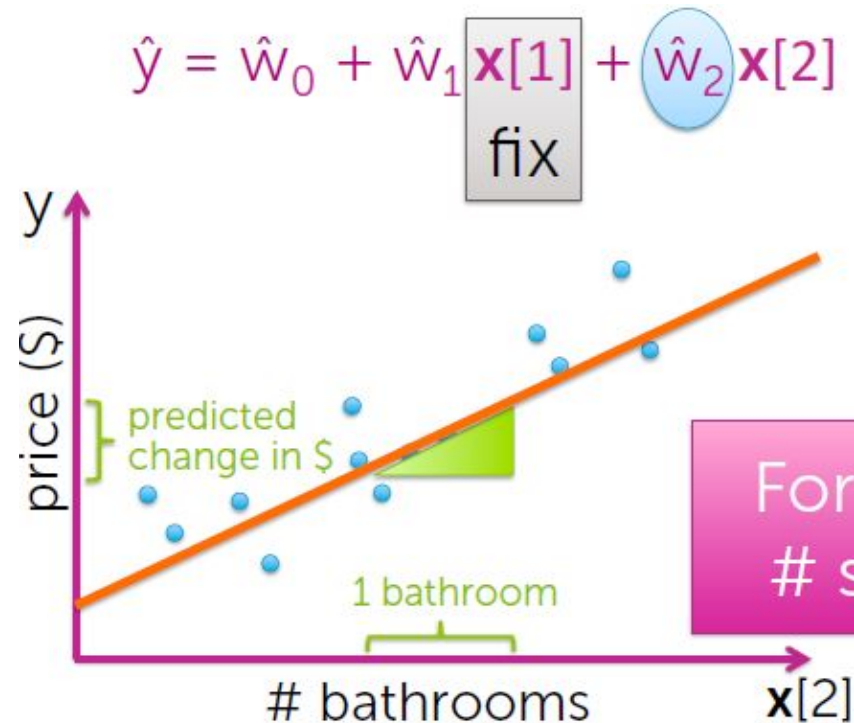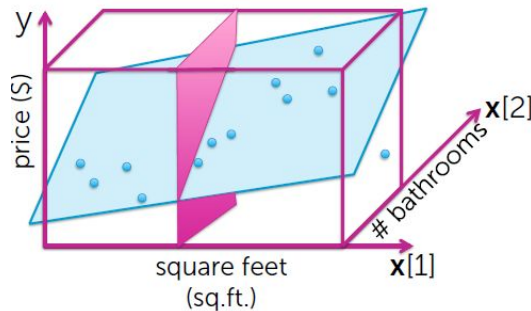# observations $(\mathbf{x}_i, y_i)$ : N

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

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

# Interpreting the coefficients

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

fix

For 2 linear features

We can fix one of the two features to see the other



©2015 Emily Fox & Carlos Guestrin

# Interpreting the coefficients

For 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

Can't hold other features fixed!

price ($)

square feet (sq.ft.)

y

x

# Fitting D-dimensional curves

Rewrite in matrix notation

For observation i

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

# Fitting D-dimensional curves

For all observations together



$$y = Hw + \varepsilon$$

# Recap: Cost of using a line



Residual sum of squares (RSS)

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

price ($) — y

square feet (sq.ft.) — x

$\hat{y}_i$

$x_i$

# RSS for multiple regression



$$RSS(\mathbf{w}) = \sum_{i=1}^{N} (y_i - \quad )^2$$

$$\hat{y}_i = \begin{array}{|c|c|c|c|c|c|c|}\hline & & & & & & \\ \hline\end{array} \begin{array}{|c|}\hline \\ \\ \\ \\ \\ \\ \\ \hline\end{array}$$

# RSS in matrix notation

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

$$= (\mathbf{y} - \mathbf{Hw})^T(\mathbf{y} - \mathbf{Hw})$$

Why?

| residual$_1$ | residual$_2$ | residual$_3$ | ... | residual$_N$ |
|---|---|---|---|---|

| residual$_1$ |
|---|
| residual$_2$ |
| residual$_3$ |
| ... |
| residual$_N$ |

# Gradient of RSS

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

$$= -2\mathbf{H}^{\top}(\mathbf{y}-\mathbf{Hw})$$
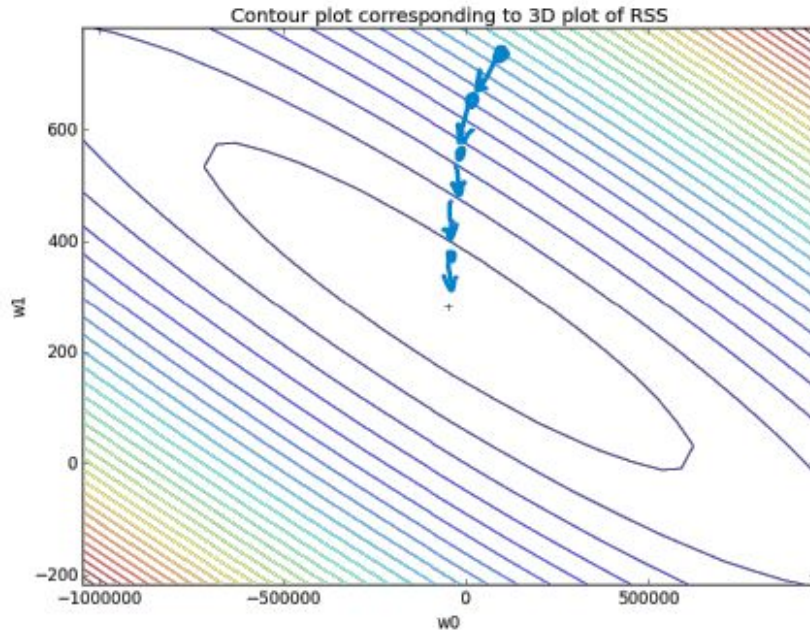
Why? By analogy to 1D case:

# Approaches to get w

Approach 1: Set gradient = 0 and solve it to get closed-form solution

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

# Gradient Descent

Contour plot corresponding to 3D plot of RSS

**while** not converged

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

$$-2H^T(y - Hw)$$

# Summary of GD for multiple regression



Contour plot corresponding to 3D plot of RSS

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

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

    for $j=0,...,D$

    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]

    $t \leftarrow t + 1$

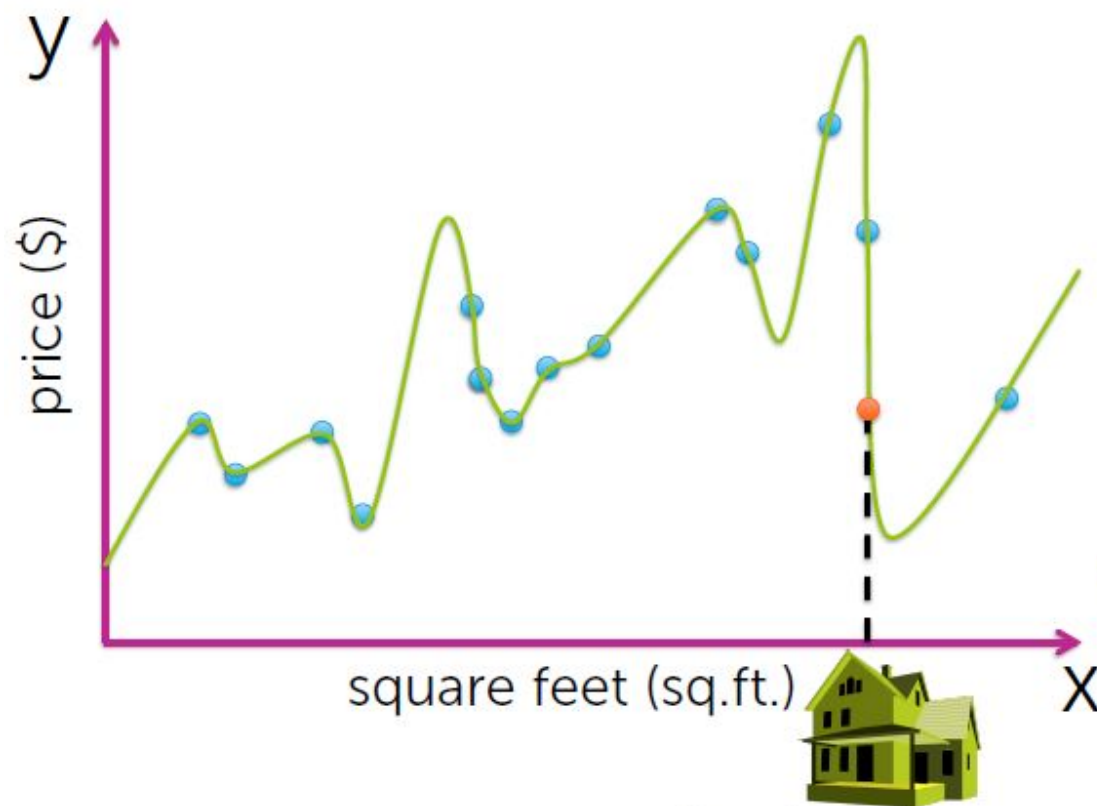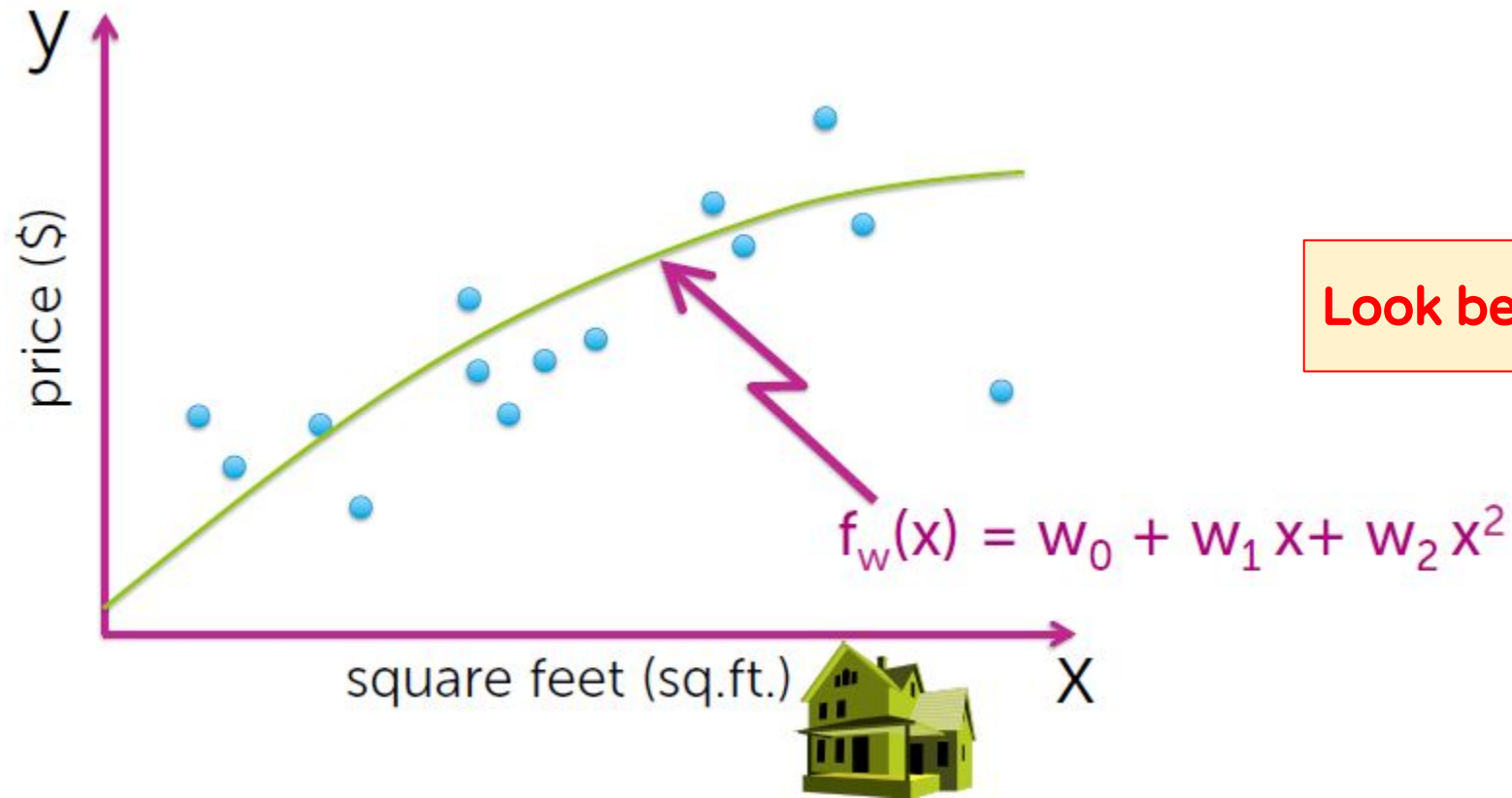# Evaluating overfitting via training/test split

# Do you believe this fit?

# What about a quadratic function?



$$f_w(x) = w_0 + w_1 x + w_2 x^2$$

Look better now?

# How to choose model order/complexity



- Want good predictions, but can't observe future
- **Simulate predictions**
  1. Remove some houses
  2. Fit model on remaining
  3. Predict heldout houses

# Training/test split



**Terminology:** – training set
– test set

# Training error

# Go To Exercises

# Anyone tried Kaggle?