# Regression:
## Predicting House Prices

Emily Fox & Carlos Guestrin

Machine Learning Specialization

University of Washington

1

# Predicting house prices

# How much is my house worth?
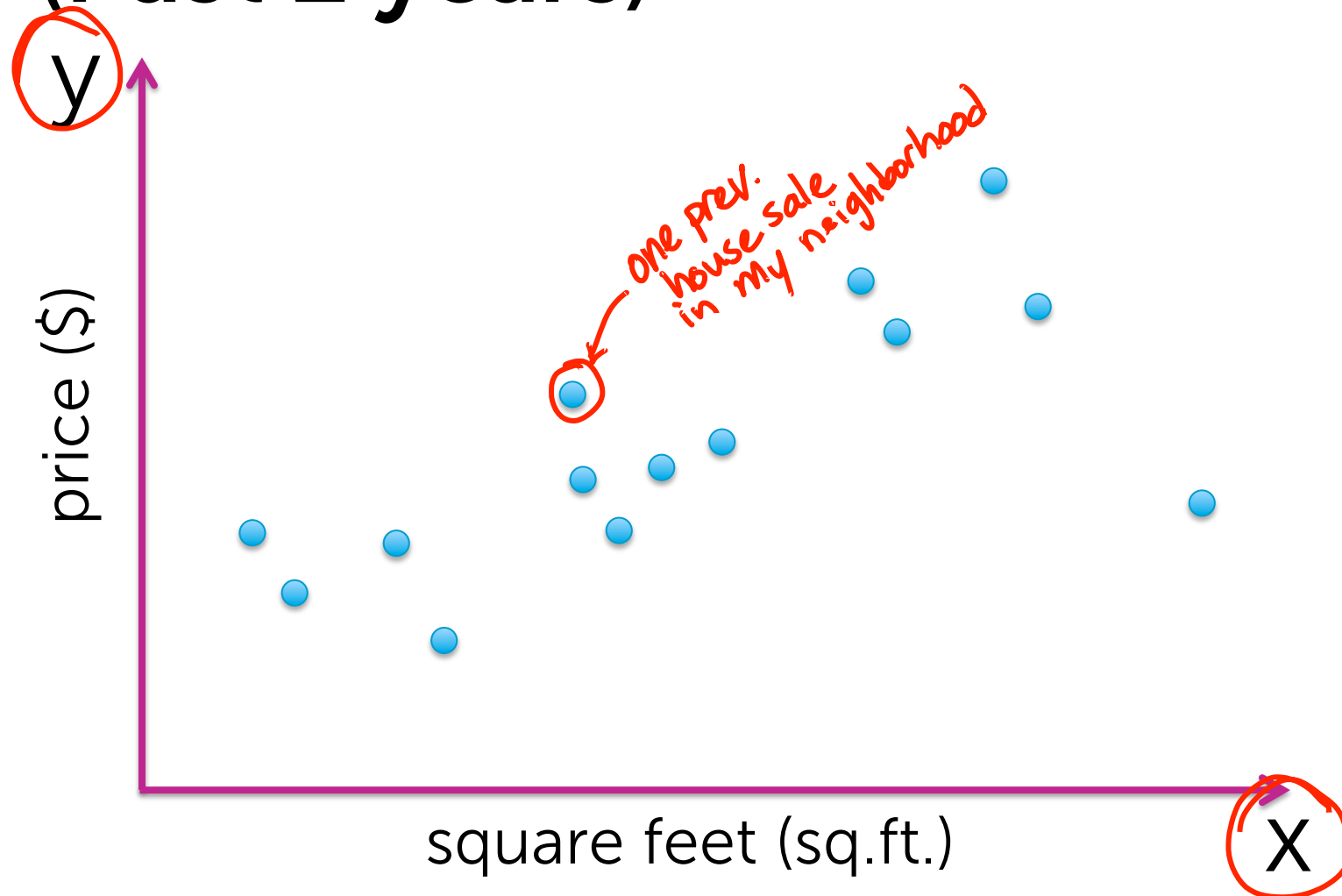


I want to list my house for sale

Machine Learning Specialization

# How much is my house worth?

$$ ????

Machine Learning Specialization

# Look at recent sales in my neighborhood

- How much did they sell for?

Machine Learning Specialization

# Plot recent house sales (Past 2 years)



**Terminology:**

x – feature, covariate, or predictor

y – observation or response

Machine Learning Specialization

# Predict your house by similar houses

y

price ($)

square feet (sq.ft.)

x

No house sold recently had *exactly* the same sq.ft.

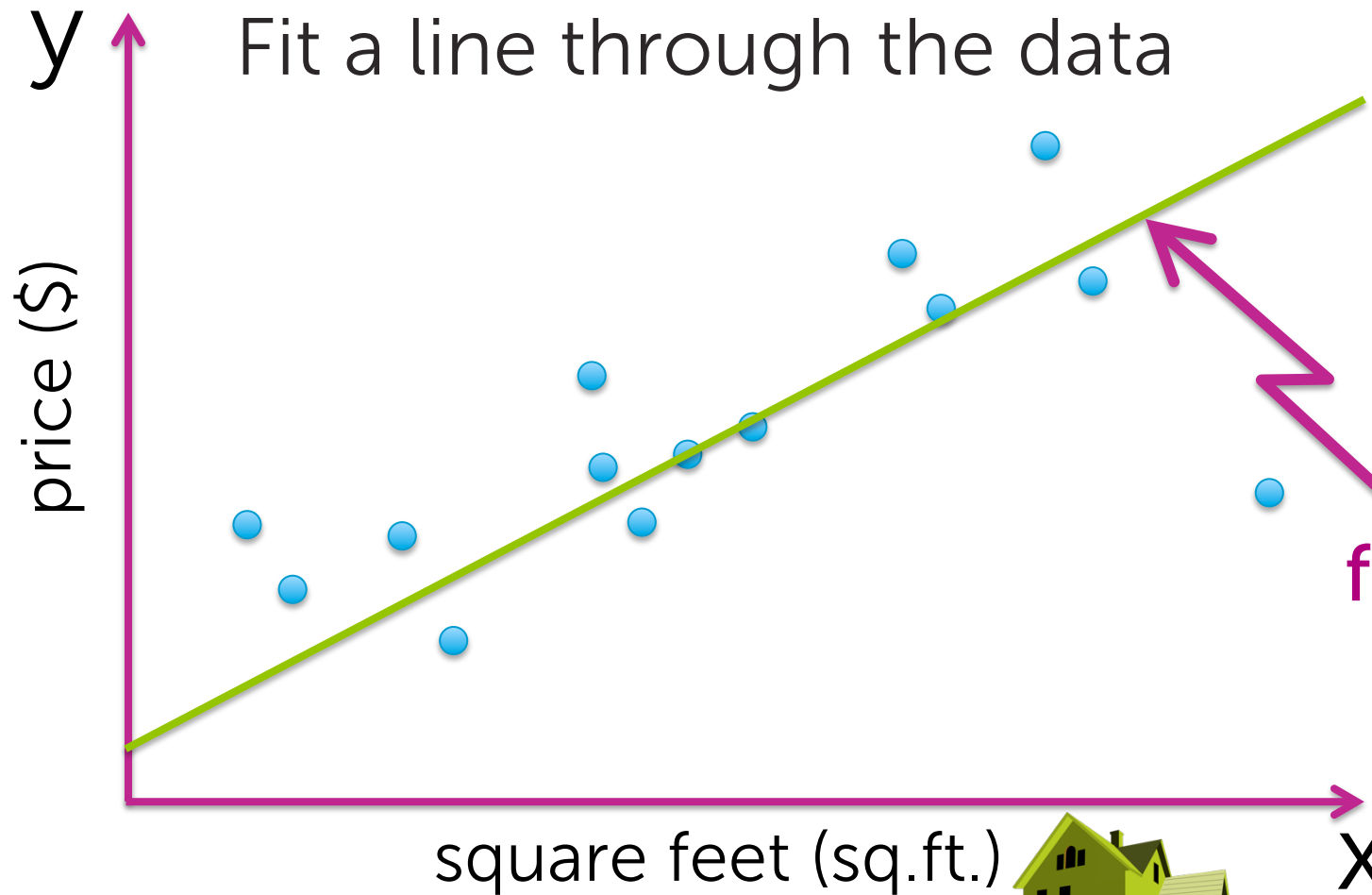Machine Learning Specialization

# Predict your house by similar houses



- Look at average price in range
- **Still only 2 houses!**
- Throwing out info from all other sales
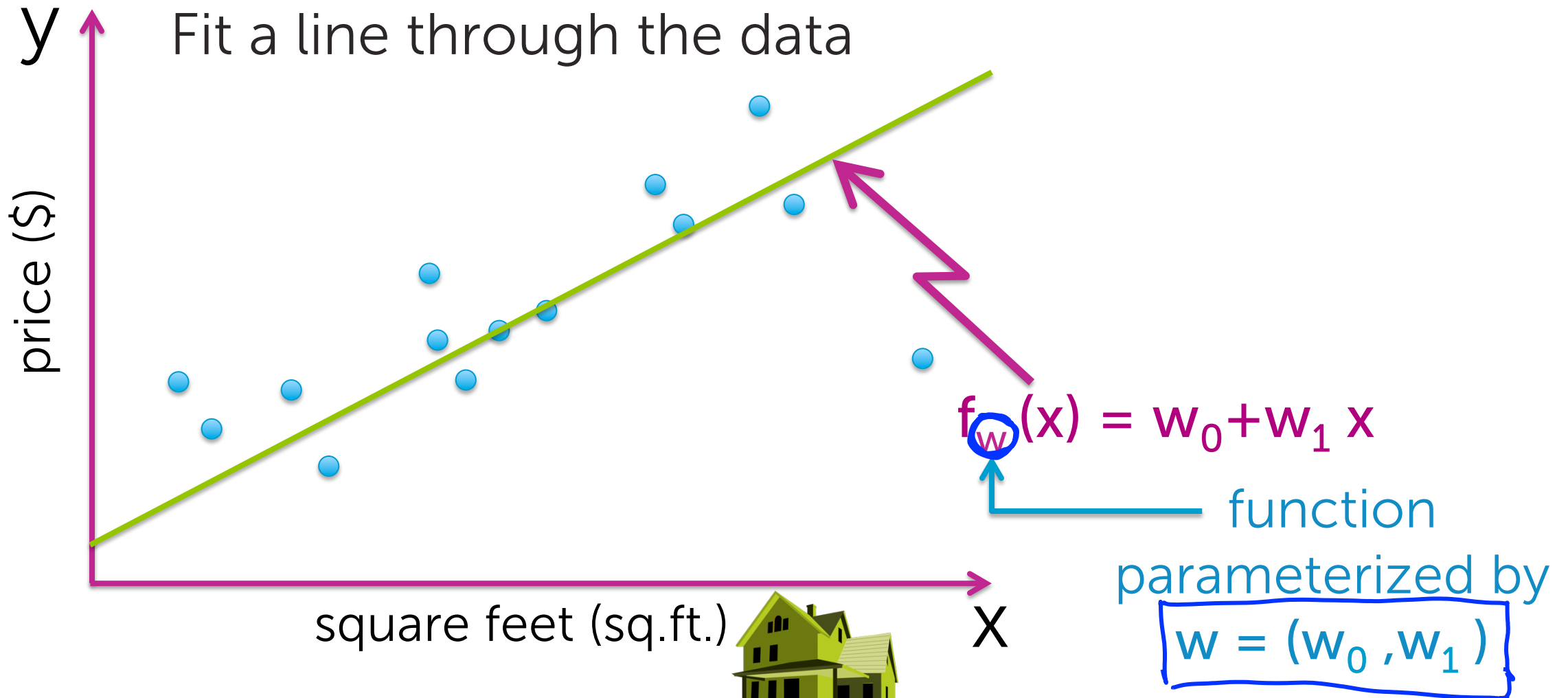
Machine Learning Specialization

# Linear regression

# Use a **linear** regression model

y

Fit a line through the data

price ($)

square feet (sq.ft.)

x

$f(x) = w_0 + w_1 x$

intercept  slope

*parameters of model*

Machine Learning Specialization

# Use a **linear** regression model

Fit a line through the data



y — price ($)

x — square feet (sq.ft.)

$f_w(x) = w_0 + w_1 x$

function parameterized by

$w = (w_0, w_1)$

# Which line?



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

different parameters w

square feet (sq.ft.)

# "Cost" of using a given line

Residual sum of squares (RSS)



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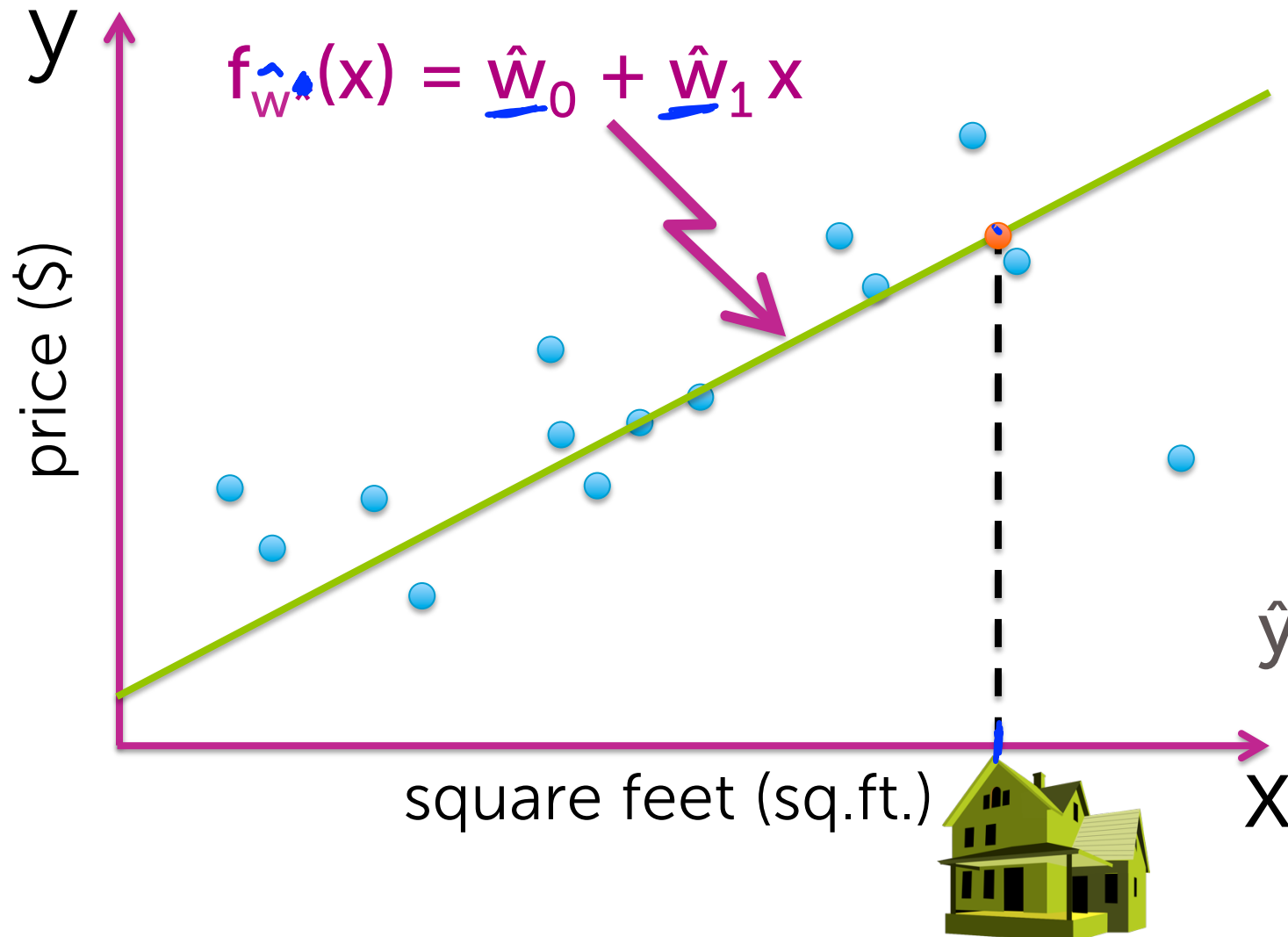$$+ ... [include\ all\ houses]$$

# Find "best" line

Minimize cost over all possible $w_0, w_1$

$y$ ↑

price ($)

square feet (sq.ft.)   $x$

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

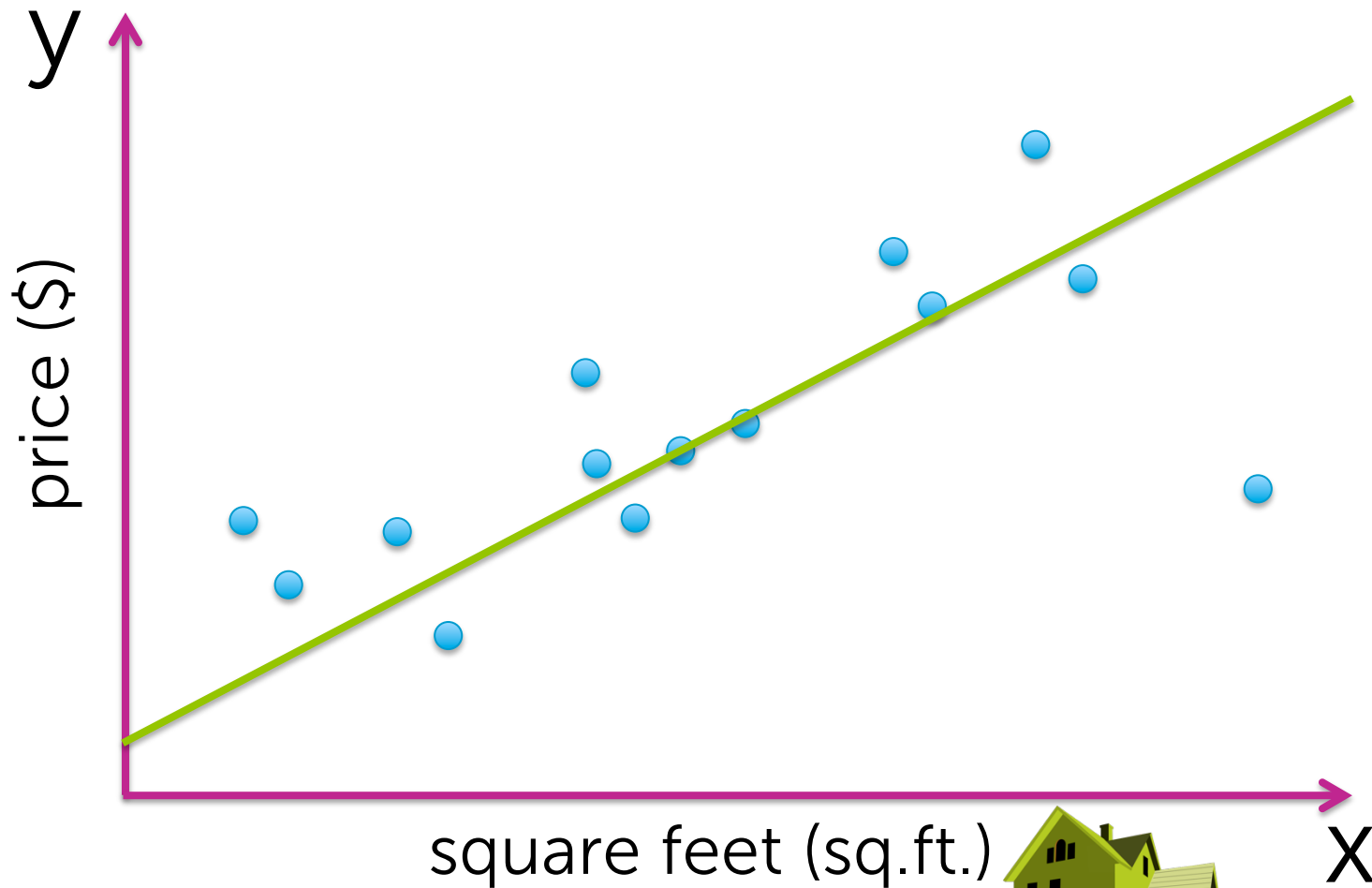$$\hat{W} = (\hat{w}_0, \hat{w}_1)$$

# Predicting your house price



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

Best guess of your house price:

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

y

price ($)

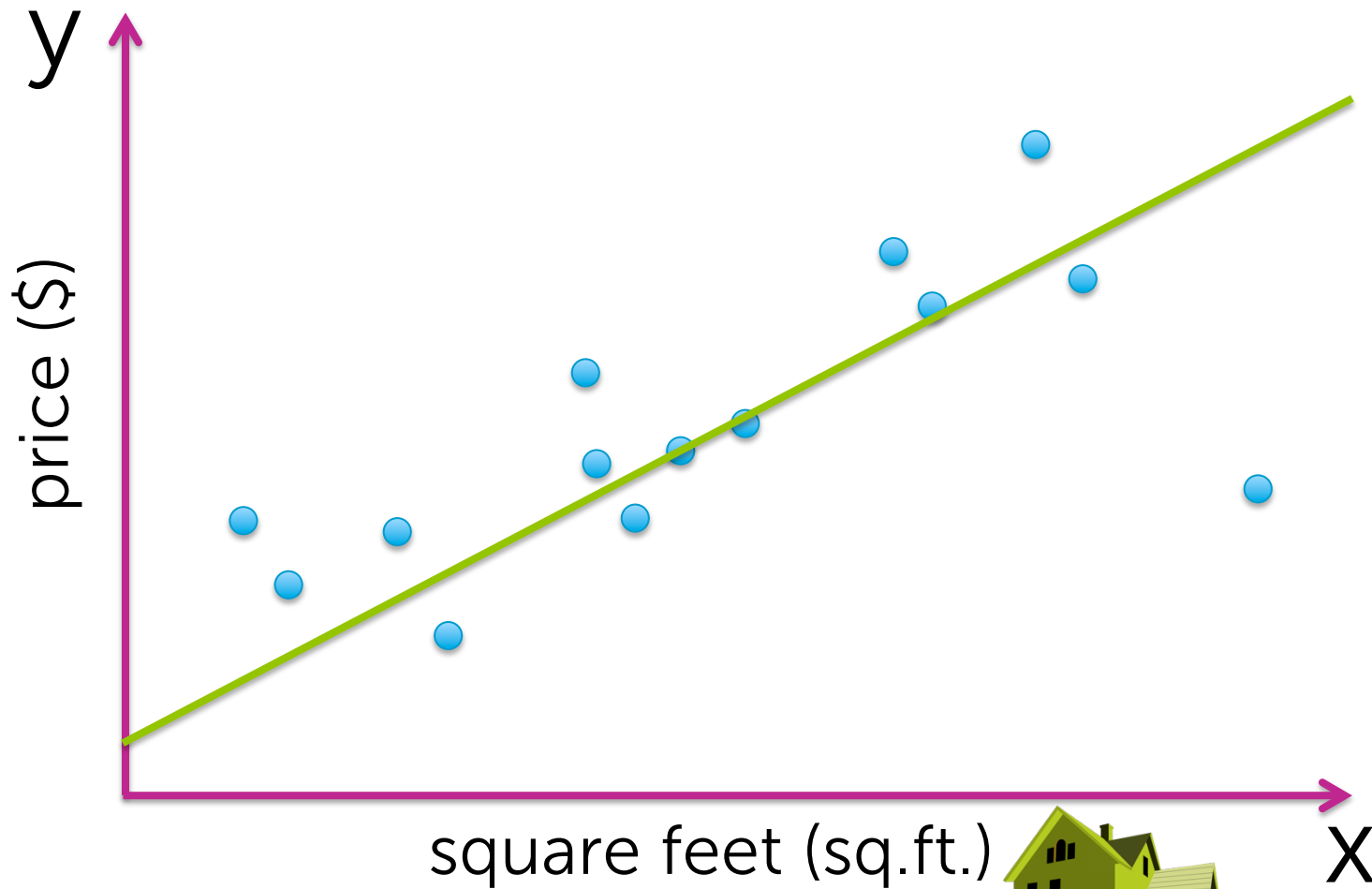square feet (sq.ft.)

x

Machine Learning Specialization

# Adding higher order effects

# Fit data with a line or ... ?



You show
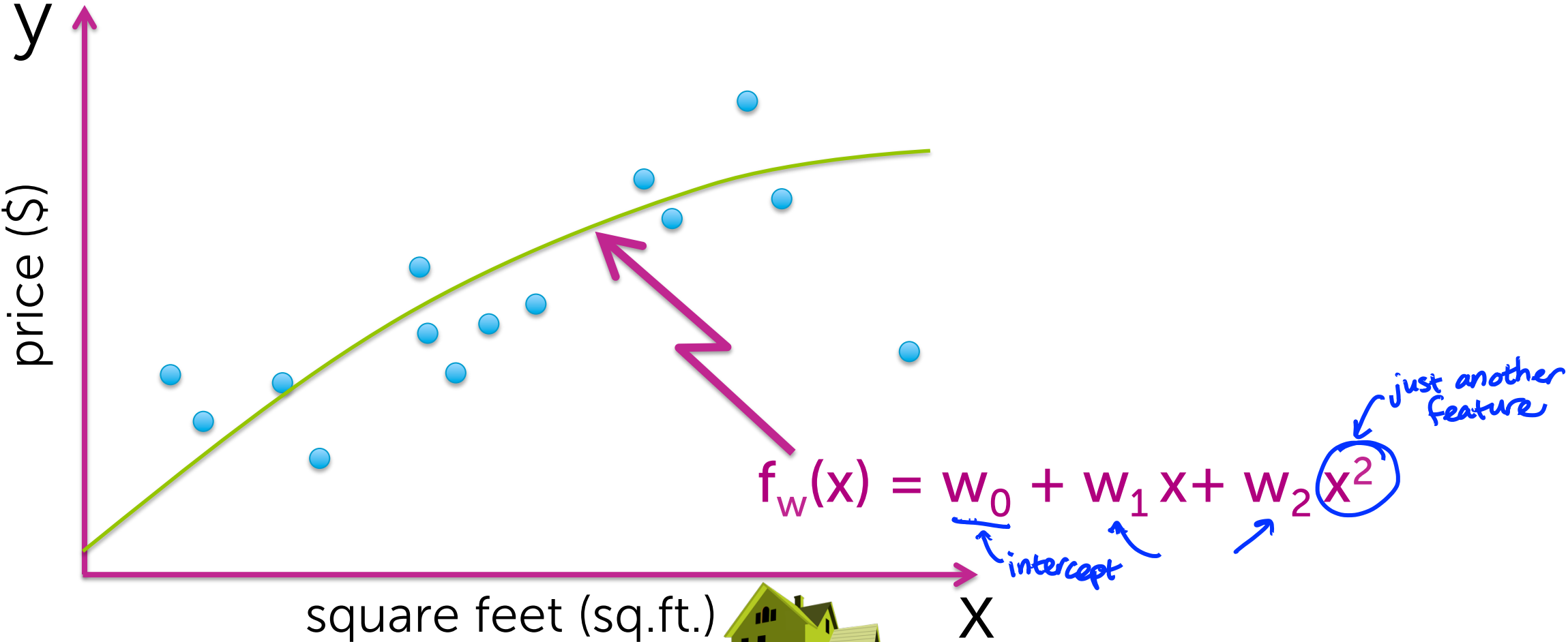your friend
your analysis

# Fit data with a line or ... ?

Machine Learning Specialization

# What about a quadratic function?

# What about a quadratic function?



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

intercept

just another feature

# Even higher order polynomial

©2015 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Do you believe this fit?

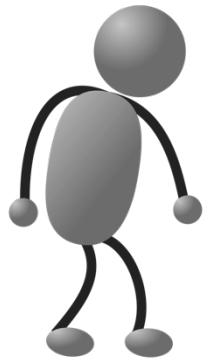Machine Learning Specialization

# Evaluating overfitting via training/test split

# Do you believe this fit?



y

price ($)

square feet (sq.ft.)

x

Minimizes RSS,
but bad predictions

# What about a quadratic function?



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

y — price ($)

x — square feet (sq.ft.)

# 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

Machine Learning Specialization

# Training error



y

price ($)

square feet (sq.ft.) x

Minimize to find $\hat{w}$

Training error ($w$) =
$(\$_{train\ 1} - f_w(sq.ft._{train\ 1}))^2$
$+ (\$_{train\ 2} - f_w(sq.ft._{train\ 2}))^2$
$+ (\$_{train\ 3} - f_w(sq.ft._{train\ 3}))^2$
$+ ...$ [include all training houses]

Machine Learning Specialization

# Test error

$y$

price ($)

square feet (sq.ft.) $x$

Assess predictions using $\hat{w}$

Test error $(\hat{w})$ =
$$(\$_{\text{test 1}}-f_{\hat{w}}(\text{sq.ft.}_{\text{test 1}}))^2$$
$$+ (\$_{\text{test 2}}-f_{\hat{w}}(\text{sq.ft.}_{\text{test 2}}))^2$$
$$+ (\$_{\text{test 3}}-f_{\hat{w}}(\text{sq.ft.}_{\text{test 3}}))^2$$
$$+ ... [\text{include all}$$
$$\text{test houses}]$$

# Training/Test Curves



Error

Model complexity

test error($\tilde{w}$)

training error($\hat{w}$)

linear    quadratic    13th order    . . .

©2015 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Adding other features

# Predictions just based on house size



y

price ($)

square feet (sq.ft.)

x

Only 1 bathroom!
Not same as my
3 bathrooms

Machine Learning Specialization

# Add more features

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



y

price ($)

square feet (sq.ft.)

$x_1$

$x_2$

\# bathrooms

# How many features to use?

- Possible choices:
  - Square feet
  - # bathrooms
  - # bedrooms
  - Lot size
  - Year built
  - ...
- **See Regression Course!**

Machine Learning Specialization

# Other regression examples

# Salary after ML specialization



hard work

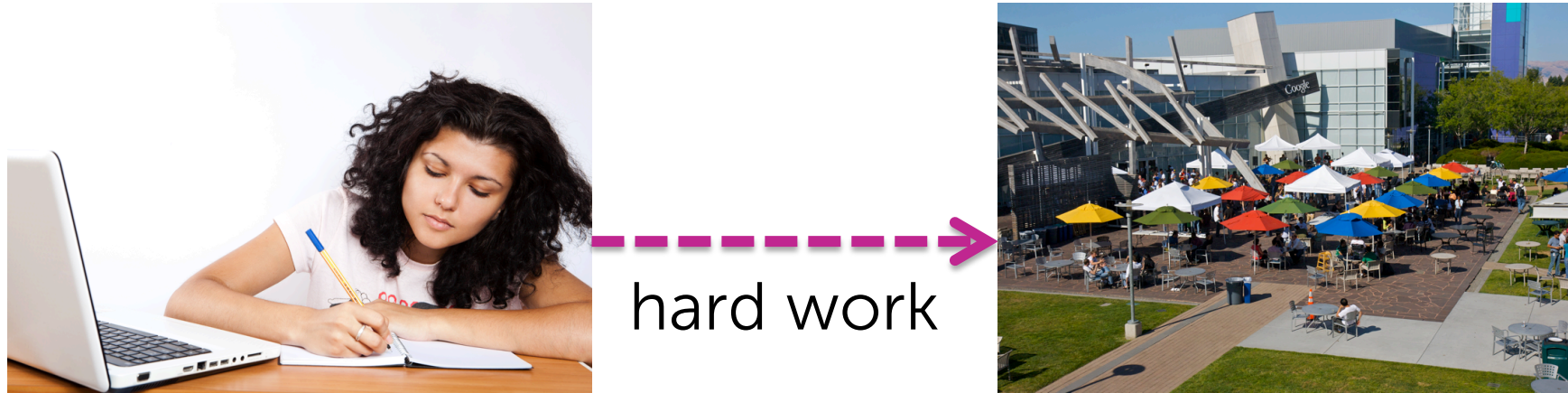- How much will your salary be? (**y** = $$)
- Depends on **x** = performance in courses, quality of capstone project, # of forum responses, …
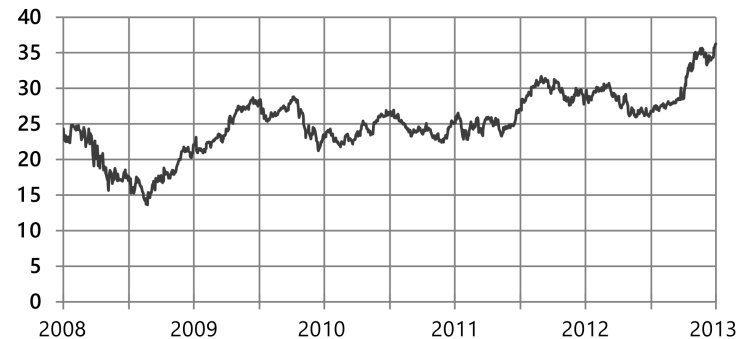
# Salary after ML specialization



hard work

$$\hat{y} = \hat{w}_0 + \hat{w}_1 \text{ performance } +$$
$$\hat{w}_2 \text{ capstone } + \hat{w}_3 \text{ forum}$$

informed by other students who
completed specialization

# Stock prediction

- Predict the price of a stock

- Depends on
  - Recent history of stock price
  - News events
  - Related commodities

Machine Learning Specialization

# Tweet popularity

- How many people will retweet your tweet?
- Depends on # followers,
   # of followers of followers,
      features of text tweeted,
         popularity of hashtag,
            # of past retweets,...

Machine Learning Specialization

# Smart houses

- Smart houses have many distributed sensors

- What's the temperature at your desk? (no sensor)
  - Learn spatial function to predict temp

- Also depends on
  - Thermostat setting
  - Blinds open/closed or window tint
  - Vents
  - Temperature outside
  - Time of day

# Summary for regression

# What you can do now...

- 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 by minimizing RSS (algorithms to come...)
- Exploit the estimated model to form predictions
- Perform a training/test split of the data
- Analyze performance of various regression models in terms of test error
- Use test error to avoid overfitting when selecting amongst candidate models
- Describe a regression model using multiple features
- Describe other applications where regression is useful

Machine Learning Specialization