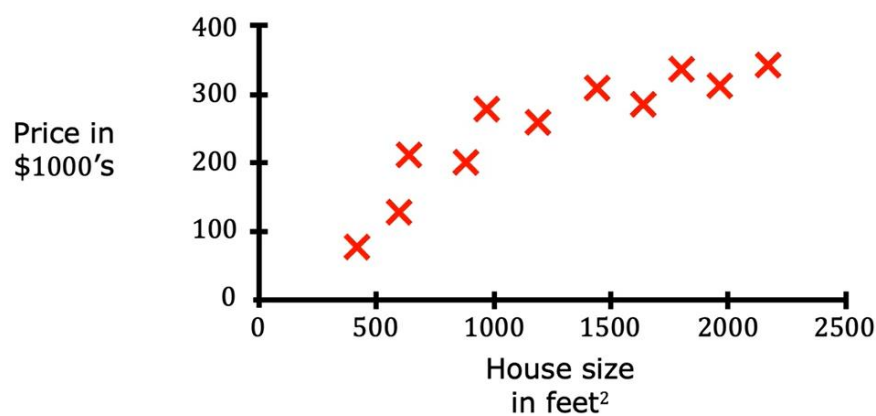


Input (X)	Output (Y)	Application
email	spam? (0/1)	spam filtering
audio	text transcripts	speech recognition
English	Spanish	machine translation
ad, user info	click? (0/1)	online advertising
image, radar info	position of other cars	self-driving car
image of phone	defect? (0/1)	visual inspection

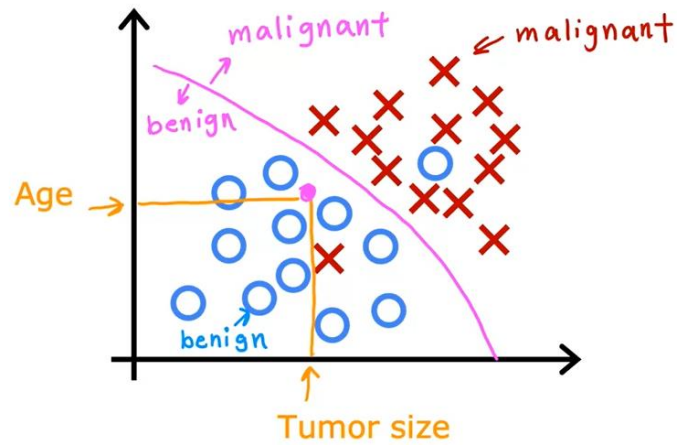
appropriate
corresponding output y.

Regression: Housing price prediction



is the size of the
house in square feet.

Two or more inputs



uniformity of the cell shape and so on.

Supervised learning

Learns from being given "right answers"

Regression

Predict a number

infinitely many possible outputs

Classification

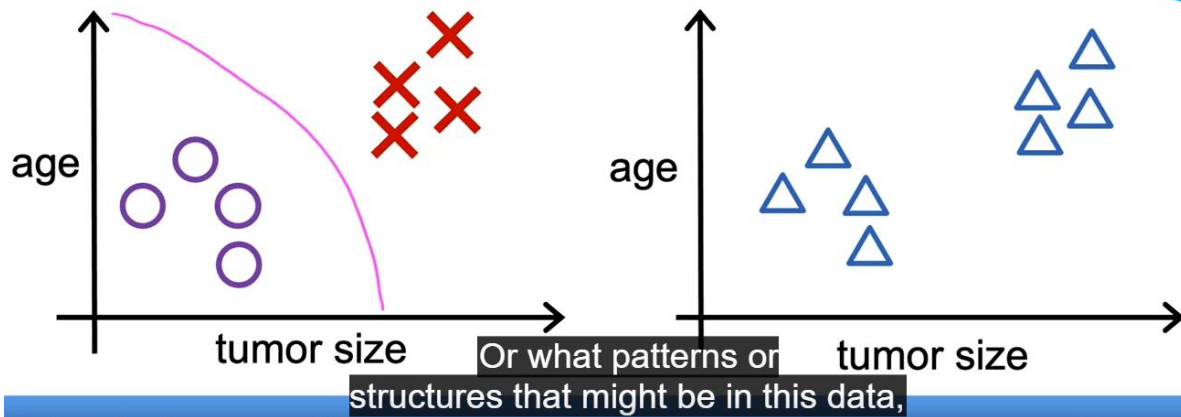
predict categories

small number of possible outputs

a category,
all of a small set of possible outputs.

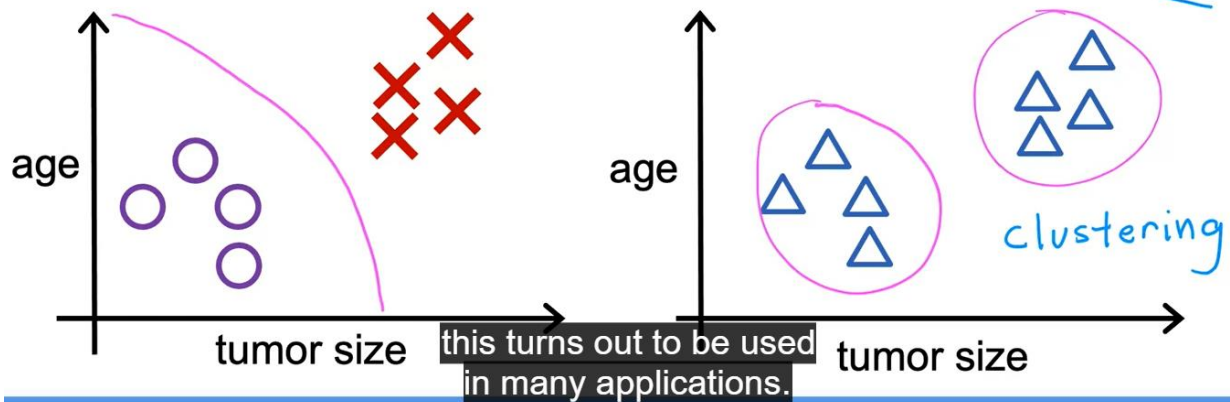
Supervised learning
Learn from data **labeled**
with the "right answers"

Unsupervised learning
Find something interesting
in **unlabeled** data.



Supervised learning
Learn from data **labeled**
with the "right answers"

Unsupervised learning
Find something interesting
in **unlabeled** data.



Unsupervised learning

Data only comes with inputs x , but not output labels y .
Algorithm has to find **structure** in the data.

Clustering

Group similar data points together.

Dimensionality reduction

Compress data using fewer numbers.

Anomaly detection

Find unusual data points.

In case anomaly detection and

Stanford ONLINE DeepLearning.AI Andrew Ng 1:27 / 3:39

training set

learning algorithm

x → f → \hat{y}

feature model prediction (estimated y) target

size → f → price (estimated)

features targets

"y-hat"

you'll also see a variation of regression where you'll

How to represent f ?

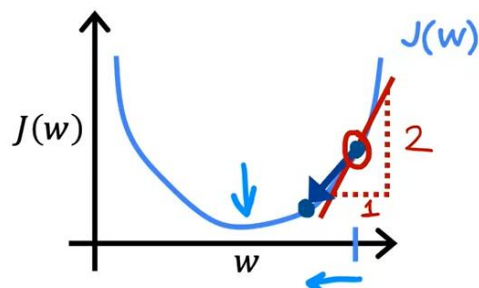
$$f_{w,b}(x) = wx + b$$
$$f(x)$$

single feature x size

Linear regression with one variable.

Univariate linear regression.

Stanford ONLINE DeepLearning.AI Andrew Ng 5:35 / 6:43

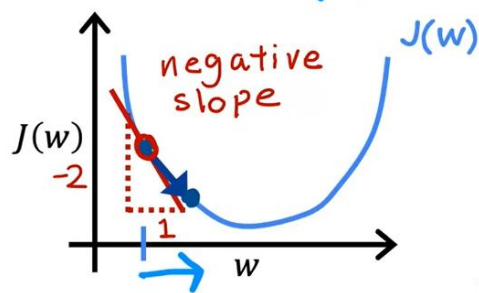


$$w = w - \alpha \frac{d}{dw} J(w)$$

> 0

$$w = w - \alpha \cdot (\text{positive number})$$

$$\frac{d}{dw} J(w) < 0$$



$$w = w - \alpha \cdot (\text{negative number})$$

Again, it looks like

Stanford ONLINE

DeepLearning.AI

Andrew Ng

$$w = w - \alpha \frac{d}{dw} J(w)$$

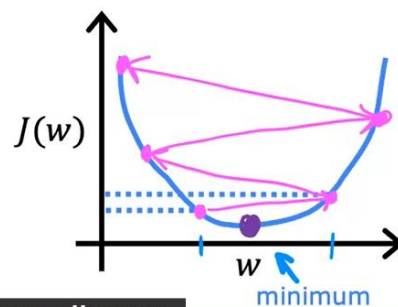
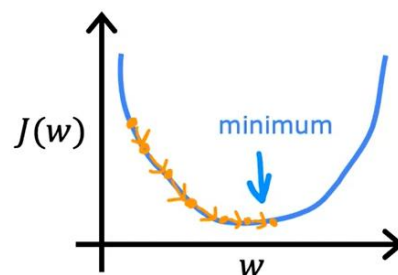
If α is too small...

Gradient descent may be slow.

If α is too large...

Gradient descent may:

- Overshoot, never reach minimum
- Fail to converge, diverge



may fail to converge and may even diverge.

Stanford ONLINE

DeepLearning.AI

Andrew Ng

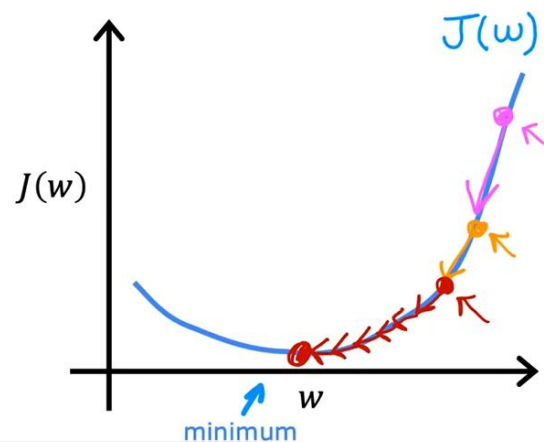
Can reach local minimum with fixed learning rate α

$$w = w - \underbrace{\alpha}_{\text{smaller}} \underbrace{\frac{d}{dw} J(w)}_{\text{not as large}} \underbrace{J(w)}_{\text{large}}$$

Near a local minimum,

- Derivative becomes smaller
- Update steps become smaller

Can reach minimum without decreasing learning rate α



So that's the gradient descent algorithm,

Stanford ONLINE

DeepLearning.AI

Andrew Ng

$$f_{\vec{w},b}(\vec{x}) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

$\vec{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n]$ parameters of the model
 b is a number

vector $\vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$

$$f_{\vec{w},b}(\vec{x}) = \underbrace{\vec{w} \cdot \vec{x}}_{\text{dot product}} + b = w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b$$

multiple linear regression

which has just one feature.

Stanford ONLINE

DeepLearning.AI

Andrew Ng

9:12 / 9:51

⏮ ⏪ ⏩ ⏭ ⚙ ⚡

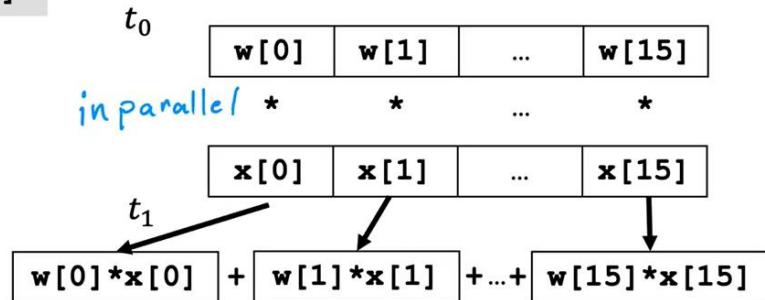
Without vectorization

```
for j in range(0,16):
    f = f + w[j] * x[j]
```

t_0 $f + w[0] * x[0]$
 t_1 $f + w[1] * x[1]$
 ...
 t_{15} $f + w[15] * x[15]$

Vectorization

```
np.dot(w,x)
```



efficient → scale to large datasets

now have to operate on.

Stanford ONLINE

DeepLearning.AI

Andrew Ng

Gradient descent

$\vec{w} = (w_1 \ w_2 \ \dots \ w_{16})$ ~~b~~ parameters

derivatives $\vec{d} = (d_1 \ d_2 \ \dots \ d_{16})$

```
w = np.array([0.5, 1.3, ... 3.4])
```

```
d = np.array([0.3, 0.2, ... 0.4])
```

compute $w_j = w_j - \underbrace{0.1}_{\text{learning rate } \alpha} d_j$ for $j = 1 \dots 16$

Without vectorization

$$w_1 = w_1 - 0.1d_1$$

$$w_2 = w_2 - 0.1d_2$$

\vdots

$$w_{16} = w_{16} - 0.1d_{16}$$

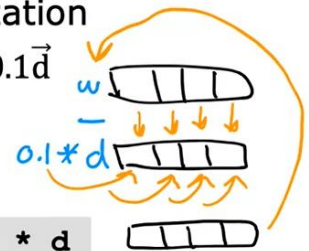
```
for j in range(0,16):
```

```
    w[j] = w[j] - 0.1 * d[j]
```

w is assigned to w minus 0.1 times d. Behind the scenes,

With vectorization

$$\vec{w} = \vec{w} - 0.1\vec{d}$$

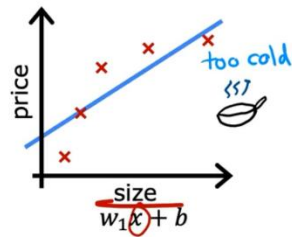


Stanford ONLINE

DeepLearning.AI

Andrew Ng

Regression example



underfit

- Does not fit the training set well

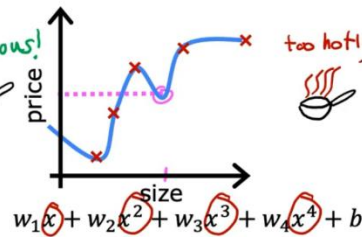
high bias



just right

- Fits training set pretty well

generalization



overfit

- Fits the training set extremely well

high variance

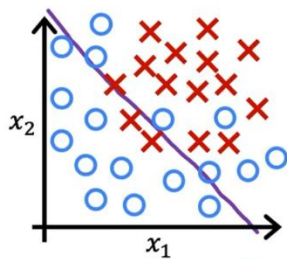
that seems to be just right.

Stanford ONLINE

DeepLearning.AI

Andrew Ng

Classification

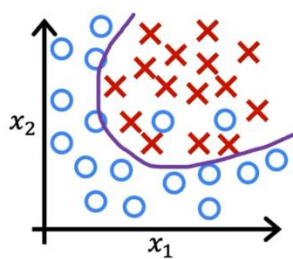


$$z = w_1x_1 + w_2x_2 + b$$

$$f_{\vec{w},b}(\vec{x}) = g(z)$$

g is the sigmoid function

underfit high bias

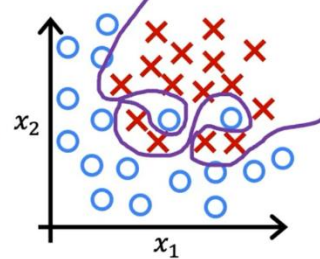


$$z = w_1x_1 + w_2x_2$$

$$+ w_3x_1^2 + w_4x_2^2$$

$$+ w_5x_1x_2 + b$$

just right



$$z = w_1x_1 + w_2x_2$$

$$+ w_3x_1^2x_2 + w_4x_1^2x_2^2$$

$$+ w_5x_1^2x_2^3 + w_6x_1^3x_2$$

$$+ \dots + b$$

overfit

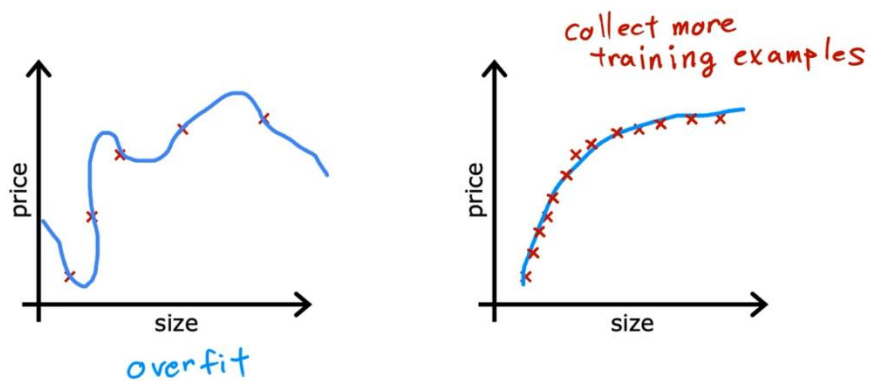
doesn't look like it'll generalize well to new examples.

Stanford ONLINE

DeepLearning.AI

Andrew Ng

Collect more training examples



been sold in this location,

Stanford ONLINE

DeepLearning.AI

Andrew Ng

Select features to include/exclude



the information by throwing away some of the features.

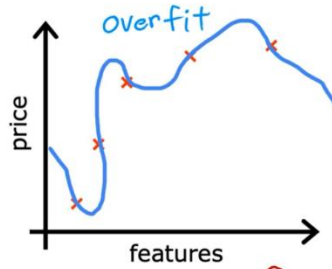
Stanford ONLINE

DeepLearning.AI

Andrew Ng

Regularization

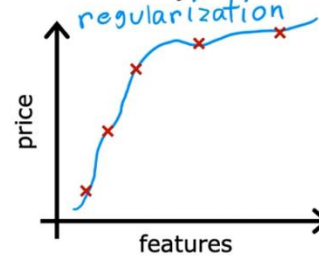
Reduce the size of parameters w_j



$$f(x) = 28x - 385x^2 + 39x^3 - 174x^4 + 100$$

large values for w_j

eliminate feature



$$f(x) = 13x - 0.23x^2 + 0.000014x^3 - 0.0001x^4 + 10$$

small values for w_j

we normally just reduce the size of the w_j parameters,

Addressing overfitting

Options

1. Collect more data
2. Select features
 - Feature selection in course 2
3. Reduce size of parameters
 - "Regularization" next videos!

for training learning algorithms,