

Multiple Features Linear Regression

Given 4 features, x_1 : size

x_2 : number of bedrooms

x_3 : number of floors

x_4 : age of house

x_j = j-th feature

n = number of features

$X^{(i)}$ = features of i-th training examples

$x_j^{(i)}$ = value of j-th feature in the i-th training example

$$f_{w,b}(x) = wX + b = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b$$

$$= w_1x_1 + \dots + w_nx_n + b = \vec{w} \cdot \vec{x} + b$$

$$\vec{w} = [w_1, w_2, \dots, w_n] \quad 1 \times n$$

dot product

$$\vec{x} = [x_1, x_2, \dots, x_n] \quad \text{for } m \text{ samples, } n \times m$$

It's called Multiple Linear Regression (NOT multivariate regression).

Gradient Descent for Multiple Linear Regression with Vectorization

Given $\vec{w} = [w_1, \dots, w_n]$, $\vec{x} = [x_1, \dots, x_n]$ and b (still a number)

repeat until convergence { n features ($n \geq 2$) }

$$w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_1^{(i)} \rightarrow \frac{\partial J(\vec{w}, b)}{\partial w_1}$$

$$\vdots$$

$$w_n = w_n - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_n^{(i)}$$

Simultaneously update

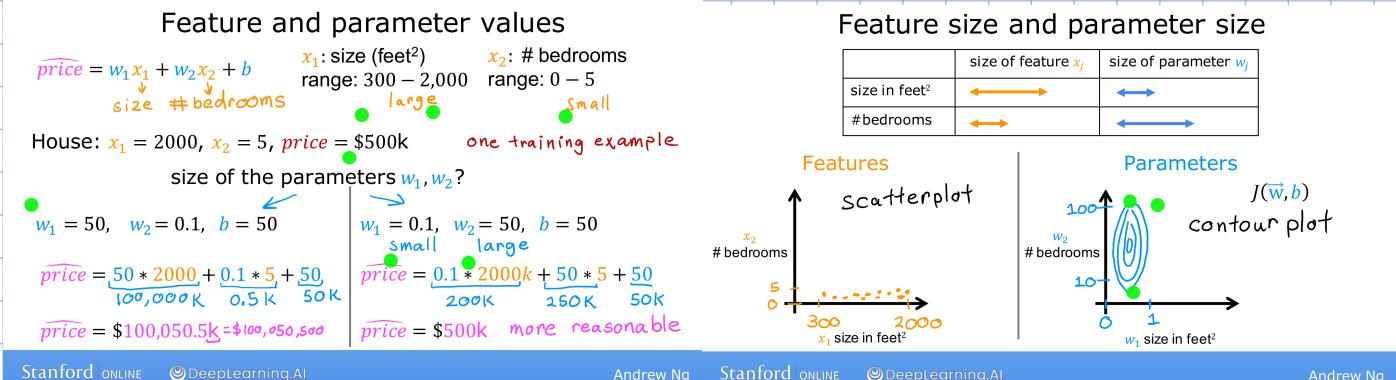
w_j (for $j=1, \dots, n$) and b

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

}

Gradient Descent in Practice

Feature Scaling:

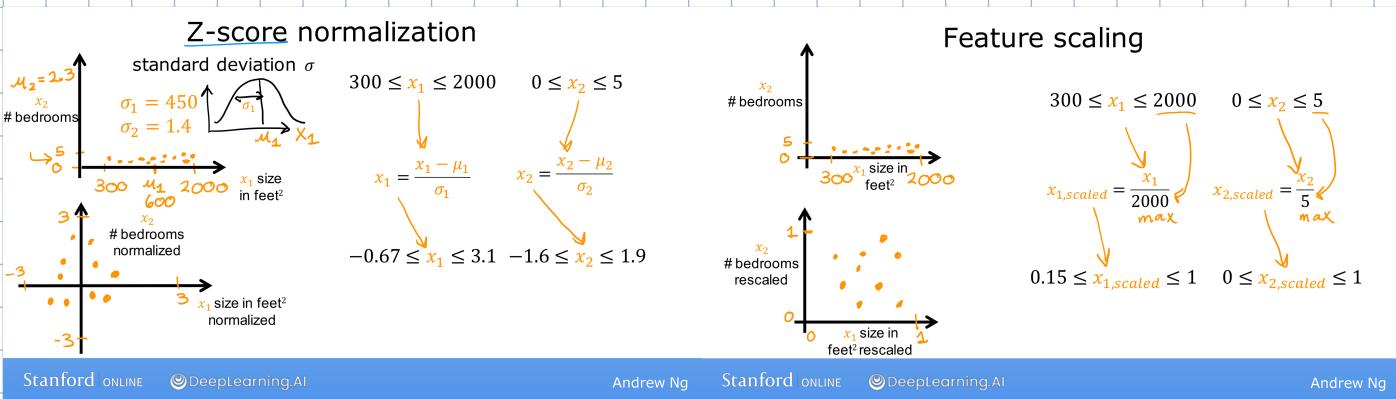
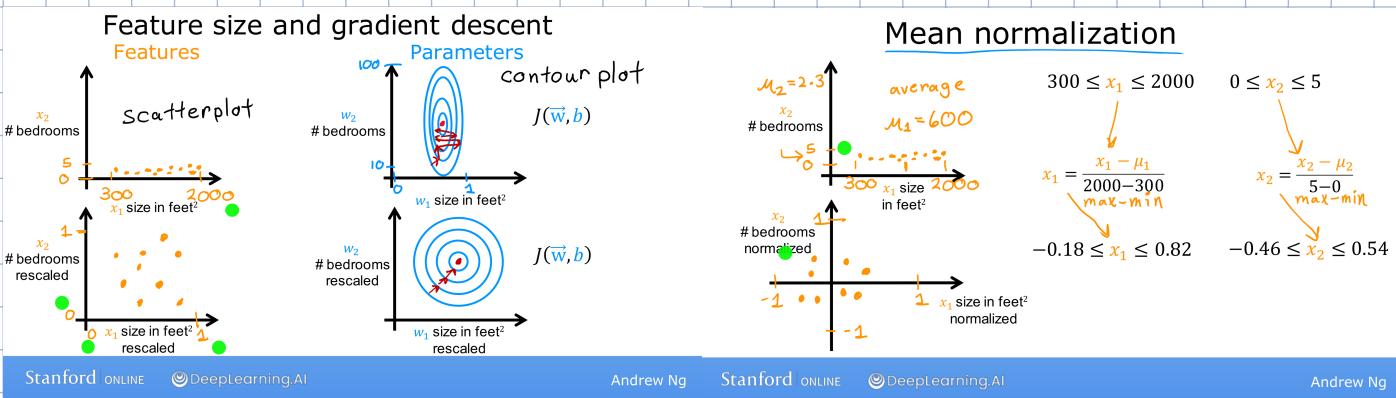


The larger the scale of the feature, the smaller the weight.

The smaller the scale of the feature, the bigger the weight.

This is BAD for gradient descent, because it may end up bouncing back and forth for long time.

We can rescale the feature. But what should the output range be? $0 \sim +1, -1 \sim +1, -3 \sim +3$, etc



Here, the output ranges are all different after rescaling. And they're all fine. The rule of thumb is aiming for $-1 \leq x_j \leq 1$ for each feature X_j . Here is a list of examples.

OK

$$-1 \leq x_j \leq 1$$

$$-3 \leq x_j \leq 3$$

$$-0.3 \leq x_j \leq 0.3$$

Okay

$$0 \leq x_1 \leq 3$$

okay

$$-2 \leq x_2 \leq 0.5$$

rescale

$$-100 \leq x_3 \leq 100$$

rescale

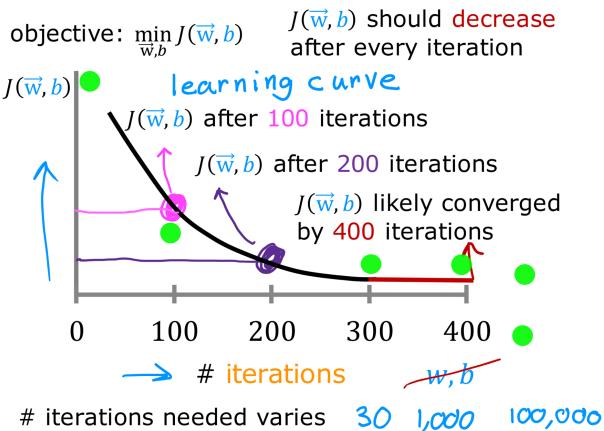
$$-0.01 \leq x_4 \leq 0.01$$

rescale

$$98 \leq x_5 \leq 105$$

Checking Gradient Descent for Convergence

Make sure gradient descent is working correctly



Automatic convergence test

Let ε "epsilon" be 10^{-3} .
0.001

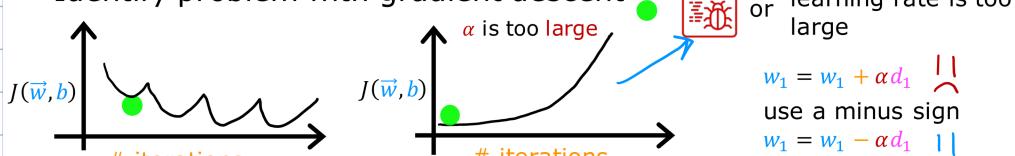
If $J(\vec{w}, b)$ decreases by $\leq \varepsilon$ in one iteration, declare convergence.
(found parameters \vec{w}, b to get close to global minimum)

Stanford ONLINE DeepLearning.AI

Andrew Ng

Choosing the Learning Rate α

Identify problem with gradient descent

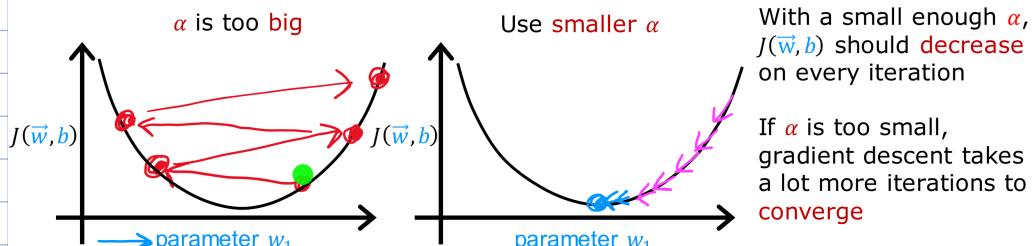


or learning rate is too large

$$w_1 = w_1 + \alpha d_1$$

$$w_1 = w_1 - \alpha d_1$$

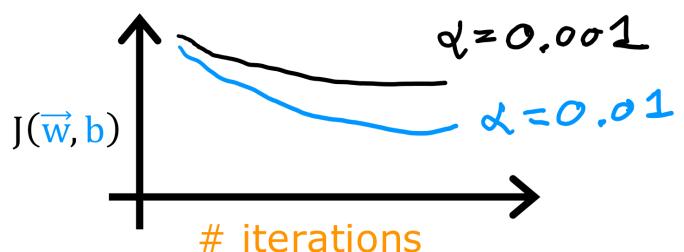
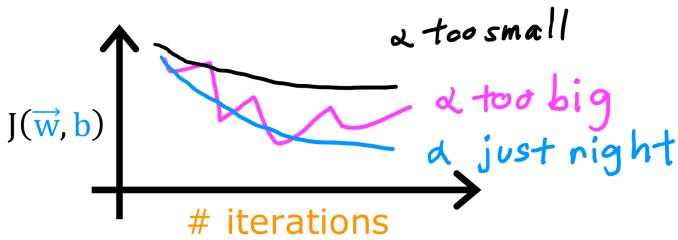
Adjust learning rate



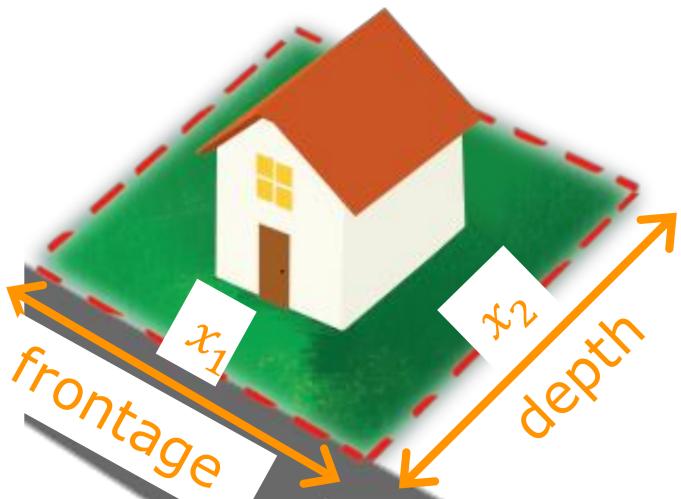
Stanford ONLINE DeepLearning.AI

Andrew Ng

Values of α to try: 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1



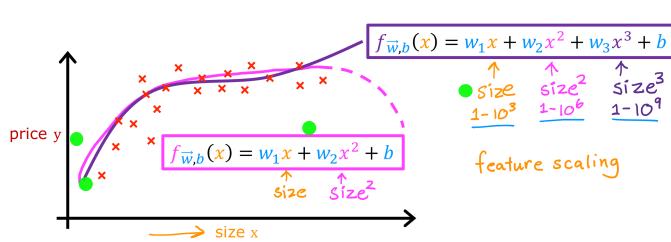
Feature Engineering



We must use our knowledge, creativity and intuition to **design** new features, by transforming or combining original features.

Polynomial Regression:

Polynomial regression



Stanford ONLINE DeepLearning.AI Andrew Ng

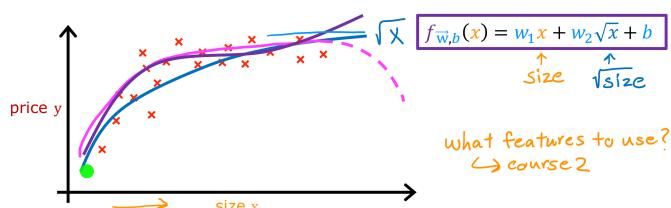
$$\text{Naive model: } f(X) = w_1 x_1 + w_2 x_2 + b$$

Better, define a new feature area X_3

$$X_3 = X_1 X_2$$

$$\text{Thus, new model } f(X) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

Choice of features



Stanford ONLINE DeepLearning.AI Andrew Ng

Extra attention must be put on the scale of designed features X^2 and X^3 . They need to be rescaled properly.

How do I know which feature to use?

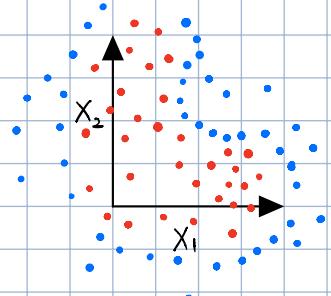
Should it be a straight line, parabolic or a random curve?

One issue with linear regression is that I may have to hand-engineer features.

For instance,

We can define $g(\theta_1 X_1 + \theta_2 X_2 + \theta_3 X_1 X_2 + \theta_4 X_1^2 X_2 + \dots)$

There are just too many possible features.



Another example is given an $50 \times 50 \times 3$ RGB image.

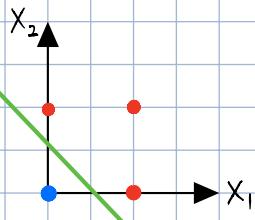
There are 7500 features. I.e. every pixel X_i is a feature.

If creating a new feature $X_i X_j$ from a given pixel X_i and another pixel X_j ,

then there will be $\frac{7500^2}{2} \approx 3$ million, which is just too many features.

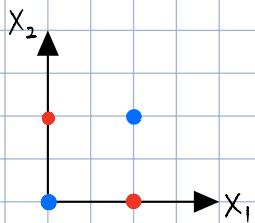
The solution is using Deep Neural Networks.

Here's another example where neural networks can create features automatically from input.



This case can be easily solved with linear regression.

Decision boundary



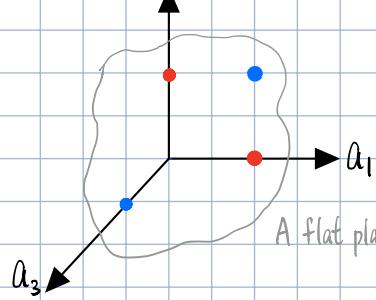
This case cannot be solved with linear regression.

The training data is not linearly separable. Note: non-linear neuron (sigmoid, tanh) can learn non-linear features, but linear neuron can't.

a_2

We can use examples above with two features X_1, X_2 as input to a neural network.

Then it may automatically create features a_1, a_2, a_3 to form a decision boundary.



A flat plane as decision boundary.