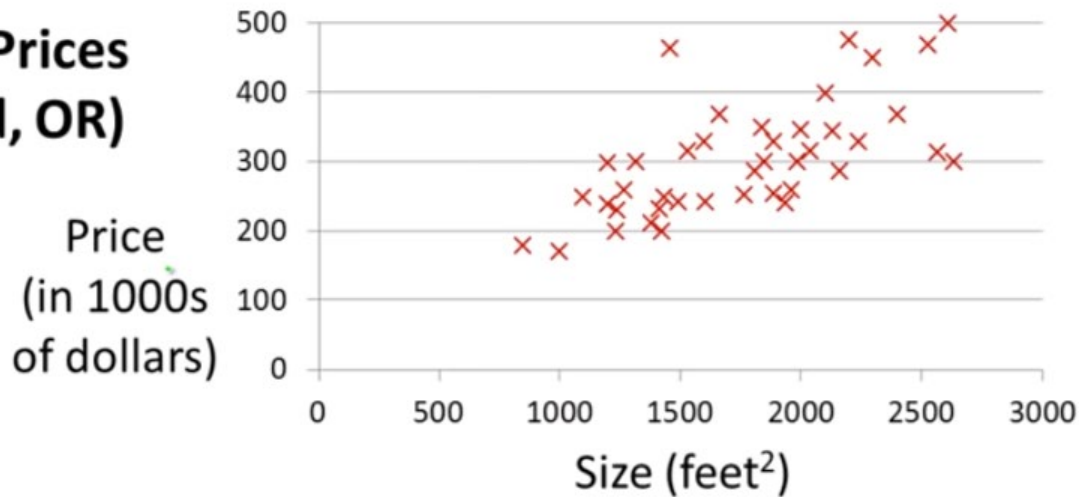


# Multivariate Linear Regression

## Housing Prices (Portland, OR)



(Single feature)

Size in feet <sup>2</sup> (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...

(Single variate linear regression)

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Predicting house price  $y$   
based on single feature of  
house size  $x$

## Multiple features (variables).

Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$x_1$	$x_2$	$x_3$	$x_4$	$y$
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...	...	...	...	...

### Notataion

$n$ : number of features

$x^i$ : features of  $i$ -th training sample

$x_j^i$ : value of feature  $j$  in  $i$ -th training sample

Hypothesis:

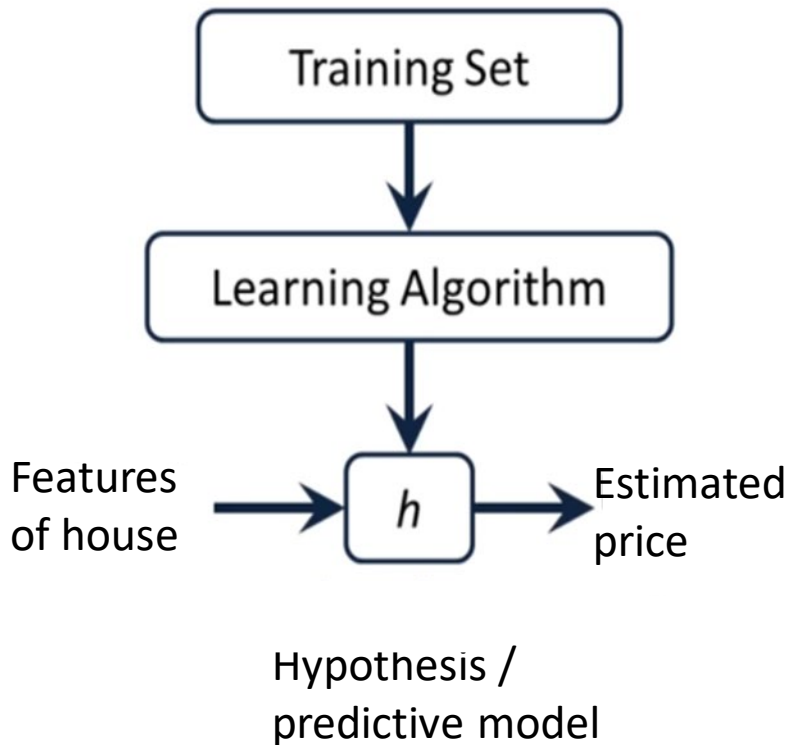
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \cdot \\ \cdot \\ \cdot \\ x_n \end{bmatrix} \in R^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \cdot \\ \cdot \\ \cdot \\ \theta_n \end{bmatrix} \in R^{n+1} \quad \theta_0: \text{bias, or} \\ \text{consider } x_0=1$$

In matrix form,

$$h_{\theta}(x) = x^T \theta$$

→ Multivariate linear regression



To represent  $h$  – need a structure/model & model parameters

Multivariate linear regression:

$$h_{\theta}(x) = x^T \theta$$

$x$ : features of house

$\theta$ : parameter of model

Hypothesis  $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$

Parameters  $\theta_0, \theta_1, \dots, \theta_n$

Cost function  $J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Update parameters  $\theta_0 \theta_1 \dots \theta_n$  using gradient descent until convergence

Simultaneously update  $\theta_j, j = 0, 1, \dots, n$ , according to

$$\theta_j := \theta_j - \eta \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \dots, \theta_n)$$

$\eta$  is a positive number, which is called the learning rate  
– too small, slow learning; too large, divergence

# Prepare Data

To make sure features are on a similar scale,  
e.g., within  $(-1, 1)$  range

# Feature scaling

Original features:

$x_1 = \text{size}(0 - 2000) \text{ sqf}$

$x_2 = \text{number of bedrooms } (1 - 5)$

Normalize by max feature:

$$x_1 = \frac{\text{size (sqf)}}{2000}$$

$$x_2 = \frac{\text{number of bedrooms}}{5}$$



Mean normalization:

$$x_1 = \frac{\text{size} - 1000}{2000}$$

$$x_2 = \frac{\# \text{ bedrooms} - 3}{5}$$

$$-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$$

Note:

Make the features approximately zero mean

Do not apply to  $x_0 = 1$

Use z score:

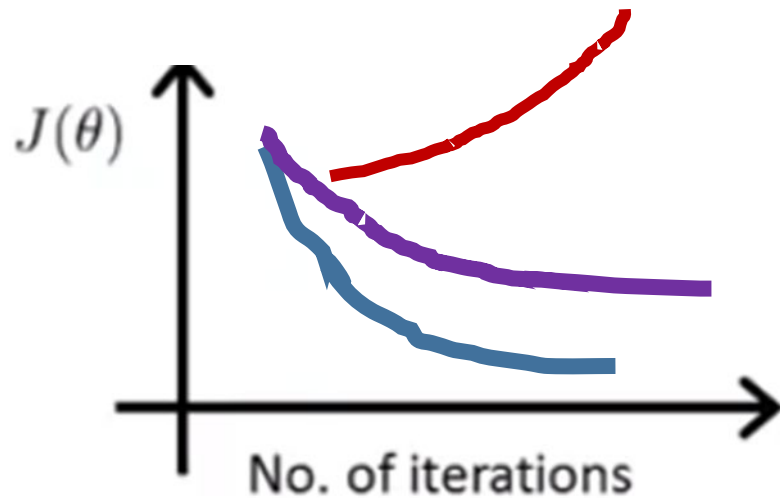
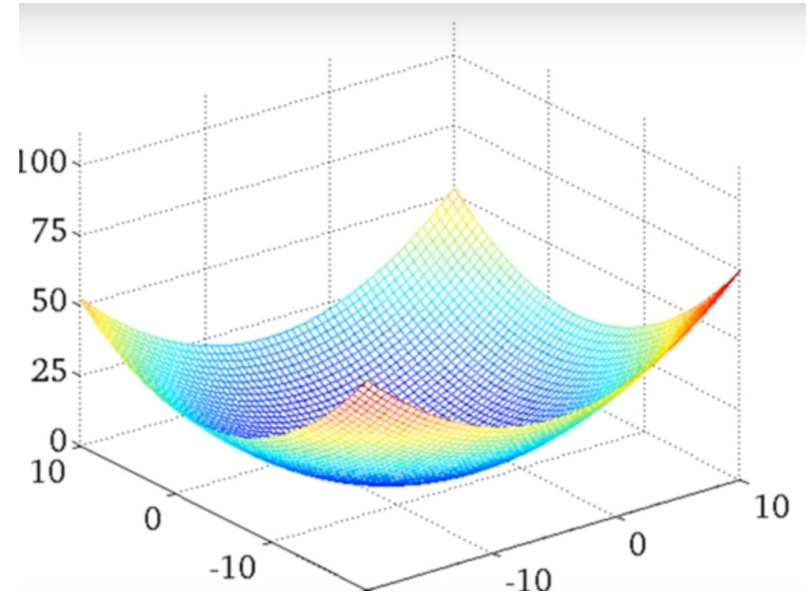
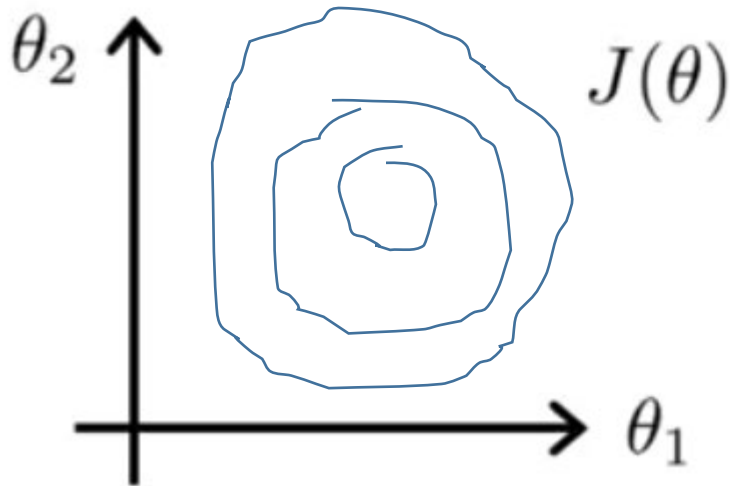
$$z = \frac{x - \mu}{\sigma}$$

where  $\mu$  is the mean of feature  $x$ , and  $\sigma$  is the standard deviation of  $x$ .

If  $x$  is a normal distribution,  $z$  is a standard normal distribution (0 mean and unit variance)

Making sure gradient descent is working properly (visualization tools can help)

# Proper use of learning curve



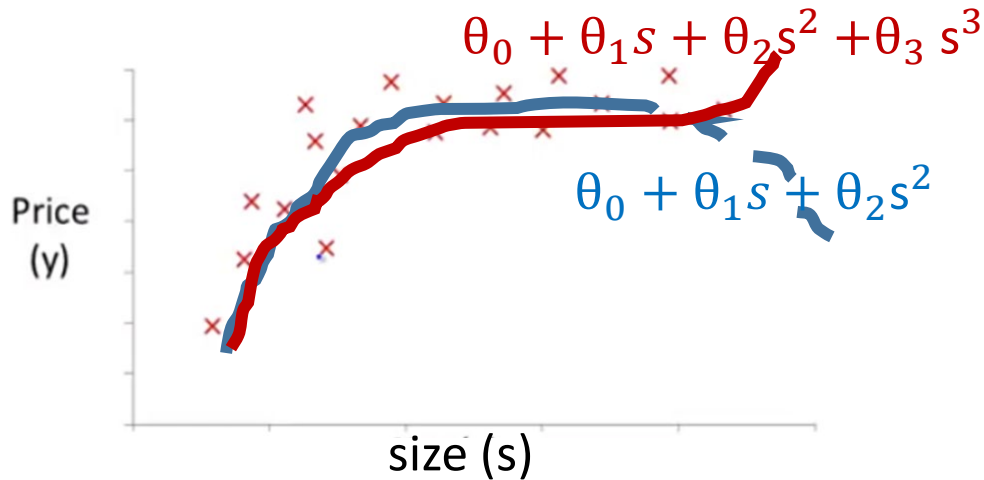
Select a range of learning rate

$$\eta = 0.001, 0.003, 0.01, 0.03, \dots$$

# Features and Polynomial Regression

(It is still linear regression but with polynomial features)

# Polynomial Regression



Still consider a single physical feature: size (s).

Let :

$$x_1 = (\text{size}) \text{ or } x_1 = s$$

$$x_2 = (\text{size})^2 \text{ or } x_2 = s^2$$

$$x_3 = (\text{size})^3 \text{ or } x_3 = s^3$$

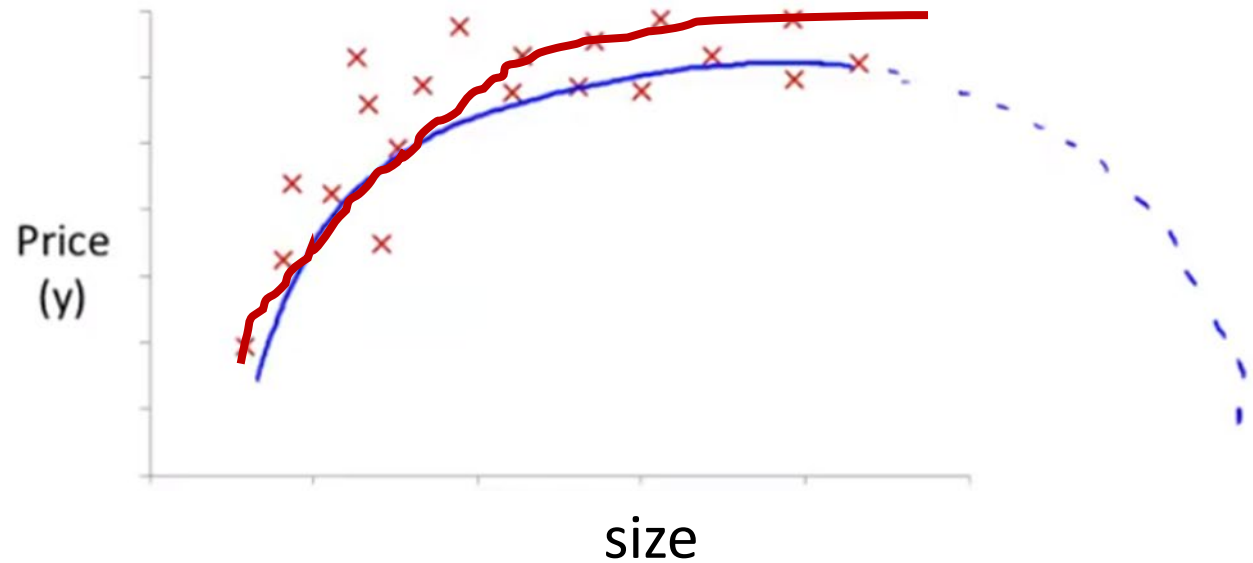
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

Single variate polynomial regression as multi-variate linear regression

$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n \\ &= \theta_0 + \theta_1 s + \theta_2 s^2 + \cdots + \theta_n s^n \end{aligned}$$

## Choice of features



$$h_{\theta}(\text{size}) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2$$

$$h_{\theta}(\text{size}) = \theta_0 + \theta_1(\text{size}) + \theta_2\sqrt{(\text{size})}$$

# Solve for $\theta$

- Gradient descent
- Normal equation (closed form solution)



How to optimize (minimize) a cost/loss function by tuning neural network weights for a predicted value to approach an actual value?

- Stochastic gradient descent (such as the LMS by Widrow and Hoff)
- Batch gradient descent as in linear regression, quadratic problem with least squares solution in closed form (may not be the best for large data set)
- Mini-batches (as in many deep networks)

$m$  **examples**  $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$  ;  $n$  **features**.

Examples:  $m = 4$ .

Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$x_1$	$x_2$	$x_3$	$x_4$	$y$
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_4 x_4$$

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \quad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

$$X\theta = y$$

$$\theta = (X^T X)^{-1} X^T y$$

$m$  **training examples**,  $n$  **features**.

Gradient Descent

- Need to choose  $\eta$
- Needs many iterations.

Normal Equation

- No need to choose  $\eta$
- Don't need to iterate.

Use when  $n$  is large