# Linear Regression

Jia-Bin Huang

Virginia Tech

Spring 2019

# Recap: Machine learning algorithms

|  | **Supervised Learning** | **Unsupervised Learning** |
|---|---|---|
| **Discrete** | Classification | Clustering |
| **Continuous** | Regression | Dimensionality reduction |

# Today's plan: Linear Regression

- Model representation

- Cost function

- Gradient descent

- Features and polynomial regression

- Normal equation

# Linear Regression

- **Model representation**

- Cost function

- Gradient descent

- Features and polynomial regression

- Normal equation

Regression

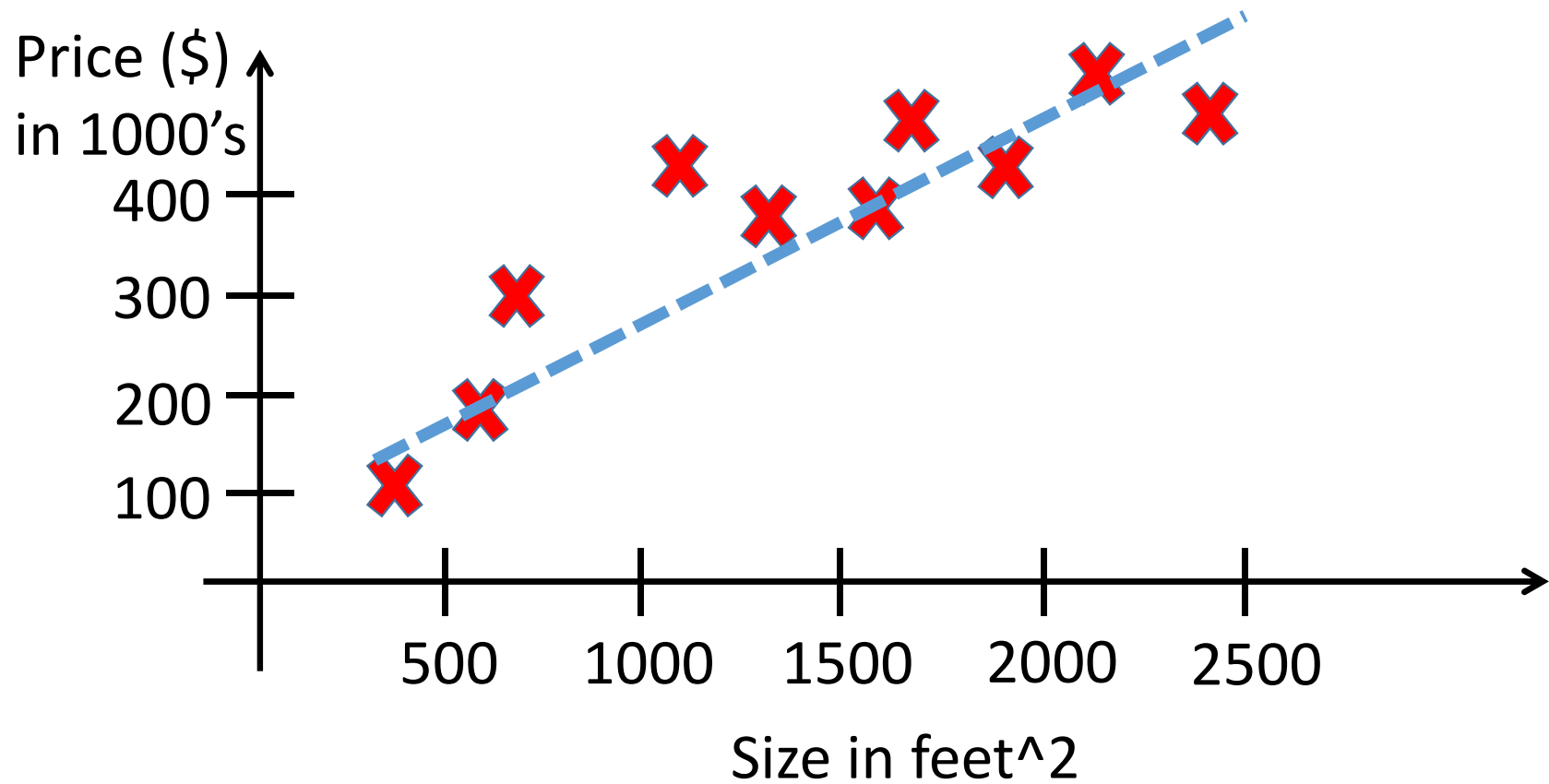real-valued output

Training set

Learning Algorithm

$x$ → $h$ → $y$
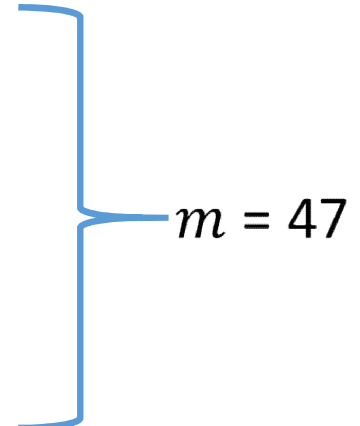
Size of house     Hypothesis     Estimated price

House pricing prediction

# Training set

| Size in feet^2 (x) | Price ($) in 1000's (y) |
| --- | --- |
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| … | … |

$m = 47$

- Notation:
  - $m$ = Number of training examples
  - $x$ = Input variable / features
  - $y$ = Output variable / target variable
  - $(x, y)$ = One training example
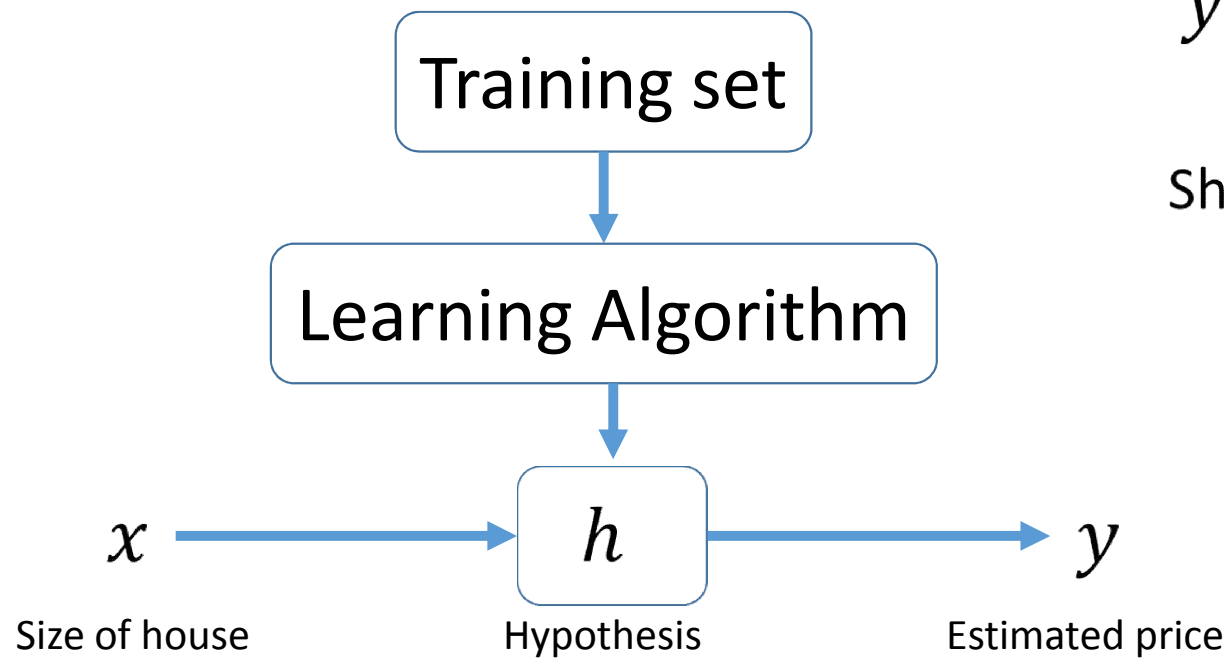  - $(x^{(i)}, y^{(i)})$ = $i^{th}$ training example
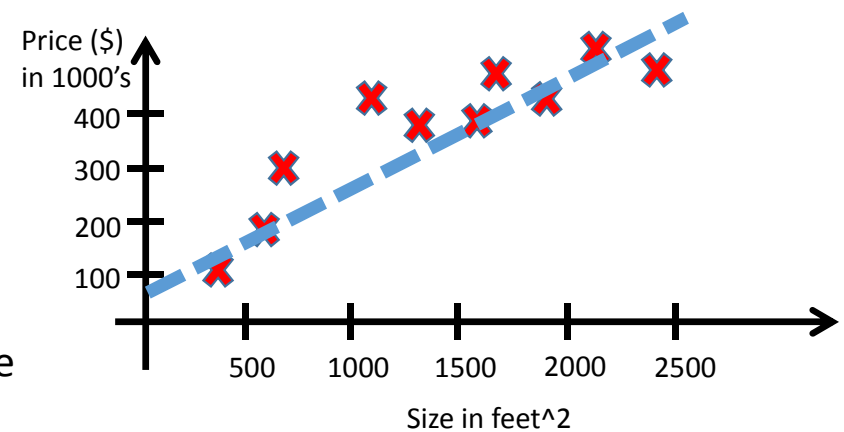
Examples:
$x^{(1)} = 2104$
$x^{(2)} = 1416$
$y^{(1)} = 460$

# Model representation

Training set

Learning Algorithm

$x$ → $h$ → $y$

Size of house        Hypothesis        Estimated price

$$y = h_\theta(x) = \theta_0 + \theta_1 x$$

Shorthand $h(x)$

Price ($) in 1000's

| | |
|---|---|
| 400 | |
| 300 | |
| 200 | |
| 100 | |

500   1000   1500   2000   2500

Size in feet^2

Univariate linear regression
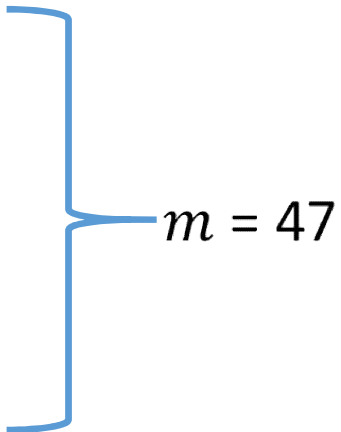
# Linear Regression

- Model representation

- **Cost function**

- Gradient descent

- Features and polynomial regression

- Normal equation

# Training set

| Size in feet^2 (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

$m = 47$

- Hypothesis $h_\theta(x) = \theta_0 + \theta_1 x$
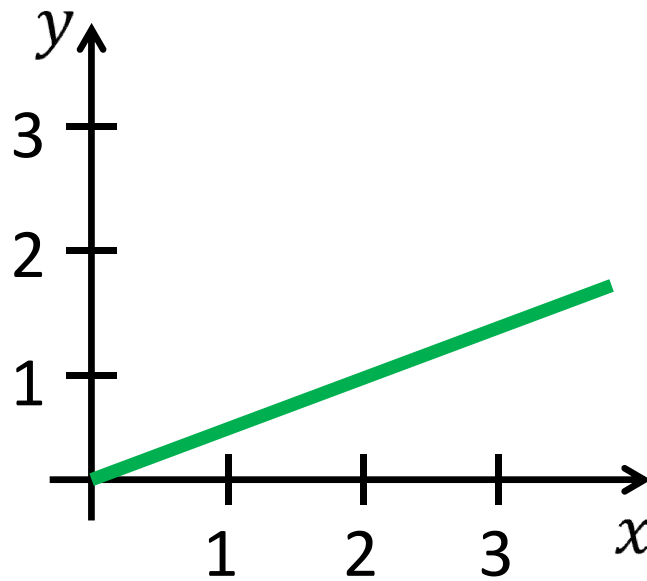
$\theta_0, \theta_1$: parameters/weights

How to choose $\theta_i$'s?

Slide credit: Andrew Ng

# Cost function

- Idea:

Choose $\theta_0, \theta_1$ so that $h_\theta(x)$ is close to $y$ for our training example $(x, y)$



$$\underset{\theta_0, \theta_1}{\text{minimize}} \ \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)^2$$

$$h_\theta\left(x^{(i)}\right) = \theta_0 + \theta_1 x^{(i)}$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)^2$$

$$\underset{\theta_0, \theta_1}{\text{minimize}} \ \boxed{J(\theta_0, \theta_1)} \ \textbf{Cost function}$$

Slide credit: Andrew Ng

# Simplified

- **Hypothesis:**

$$h_\theta(x) = \theta_0 + \theta_1 x \longrightarrow$$

- **Parameters:**

$$\theta_0, \theta_1 \longrightarrow$$

- **Cost function:**

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)^2 \longrightarrow$$

- **Goal:**

$$\underset{\theta_0, \theta_1}{\text{minimize}} \ J(\theta_0, \theta_1) \longrightarrow$$

- **Hypothesis:**

$$h_\theta(x) = \theta_1 x \qquad \theta_0 = 0$$

- **Parameters:**

$$\theta_1$$

- **Cost function:**

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)^2$$
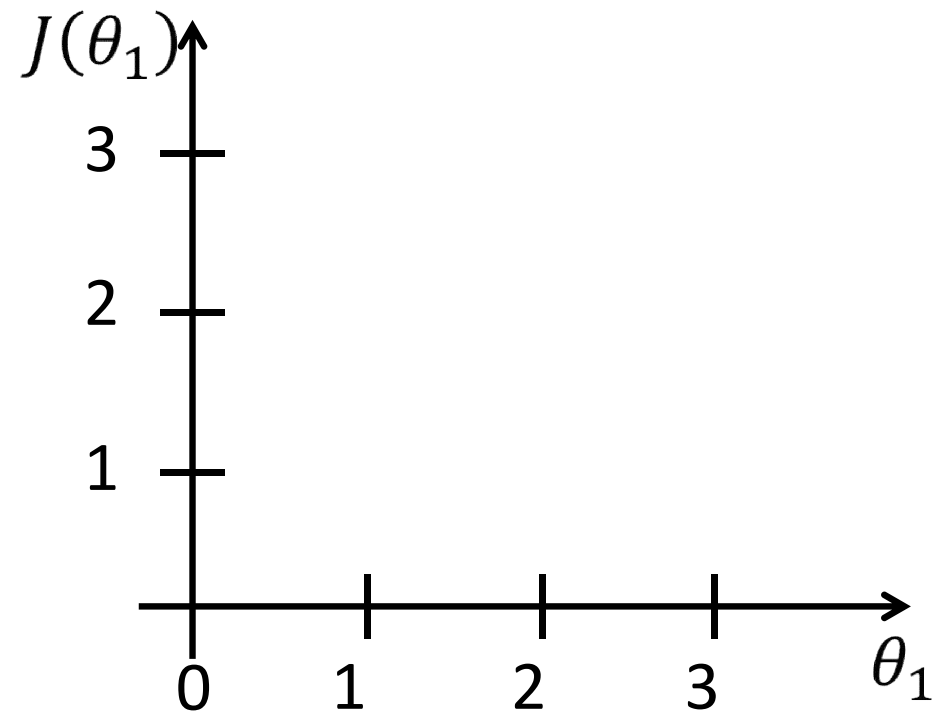
- **Goal:**

$$\underset{\theta_0, \theta_1}{\text{minimize}} \ J(\theta_1)$$

# $h_\theta(x)$, function of $x$



# $J(\theta_1)$, function of $\theta_1$

# $h_\theta(x)$, function of $x$

# $J(\theta_1)$, function of $\theta_1$

# $h_\theta(x)$, function of $x$



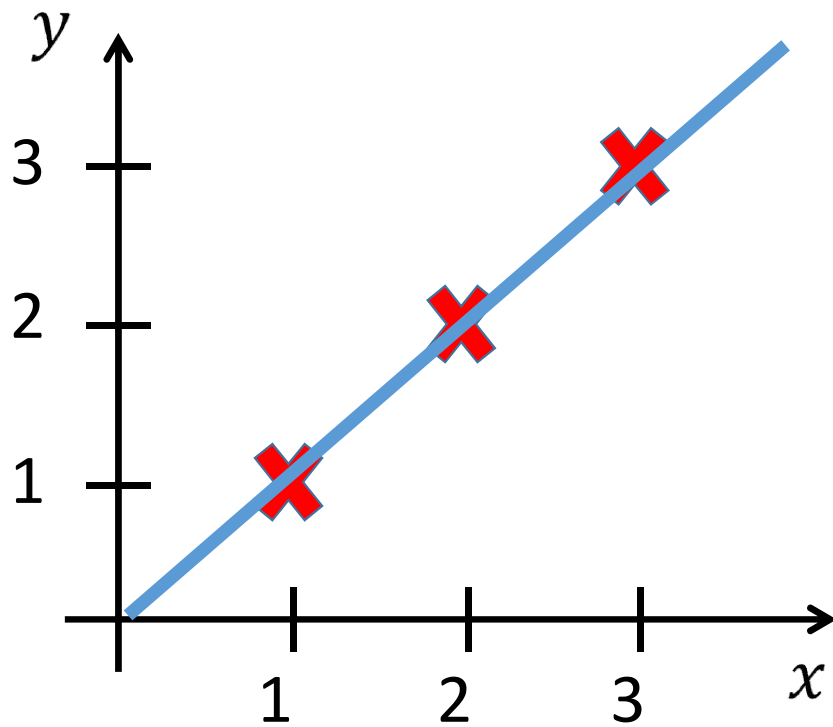# $J(\theta_1)$, function of $\theta_1$
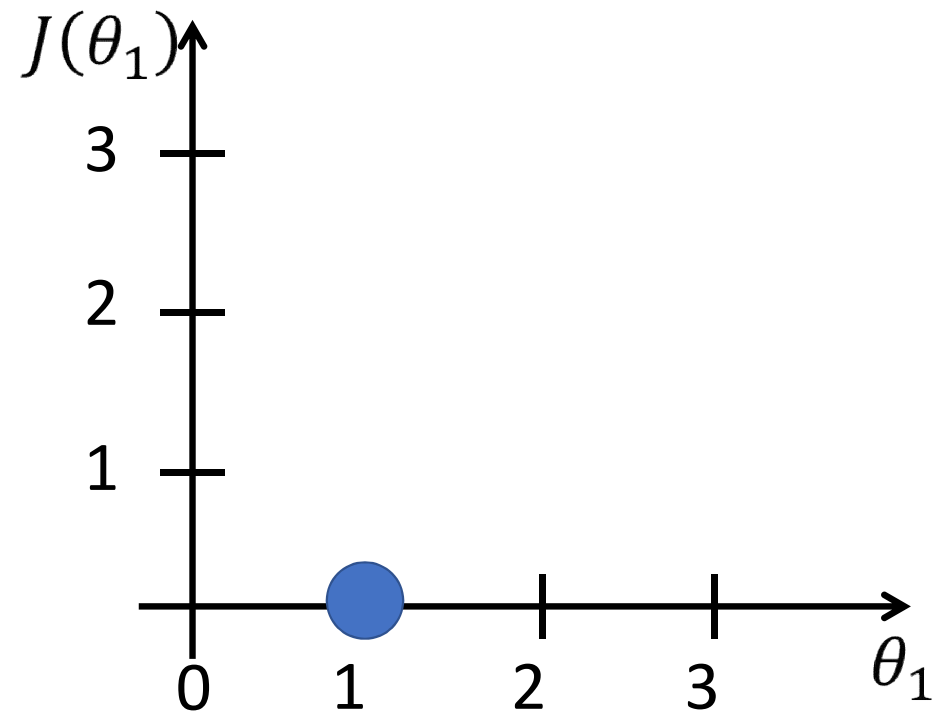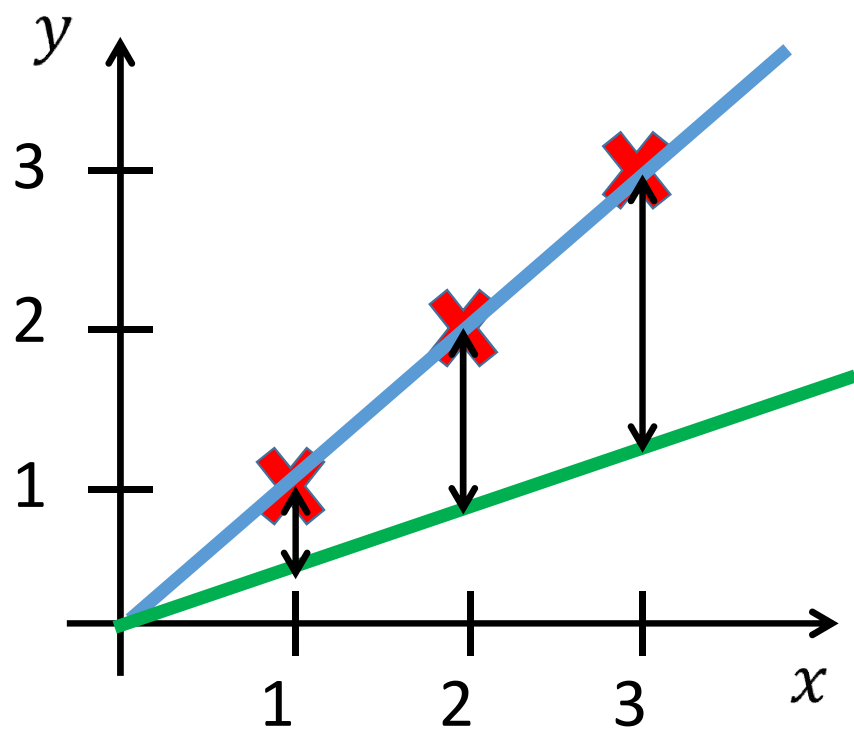
# $h_\theta(x)$, function of $x$

# $J(\theta_1)$, function of $\theta_1$

$h_\theta(x)$, function of $x$

$J(\theta_1)$, function of $\theta_1$

- **Hypothesis:** $h_\theta(x) = \theta_0 + \theta_1 x$

- **Parameters:** $\theta_0, \theta_1$

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

- **Goal:** $\underset{\theta_0, \theta_1}{\text{minimize}} \; J(\theta_0, \theta_1)$

# Cost function

# What are Contour Plots?

- Contour plots are a way to show a three-dimensional surface on a [two-dimensional plane](). It graphs two predictor variables X, Y on the y-axis and a [response variable]() Z as contours. These contours are sometimes called *z-slices* or *iso-response values*.

- A contour plot is appropriate if you want to see how some value Z changes as a [function]() of two inputs, X and Y: z = f(x, y).

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

How do we find good $\theta_0, \theta_1$ that minimize $J(\theta_0, \theta_1)$?

Slide credit: Andrew Ng

# Linear Regression

- Model representation

- Cost function

- **Gradient descent**

- Features and polynomial regression

- Normal equation

# Gradient descent

Have some function $J(\theta_0, \theta_1)$

Want $\underset{\theta_0, \theta_1}{\text{argmin}}\ J(\theta_0, \theta_1)$

Outline:

- Start with some $\theta_0, \theta_1$
- Keep changing $\theta_0, \theta_1$ to reduce $J(\theta_0, \theta_1)$
  until we hopefully end up at minimum

$J(\theta_0, \theta_1)$

$\theta_0$

$\theta_1$

Slide credit: Andrew Ng

# Gradient descent

Repeat until convergence{

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \text{(for } j = 0 \text{ and } j = 1\text{)}$$

}

$\alpha$: Learning rate (step size)

$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$: derivative (rate of change)

# Gradient descent

**Correct: simultaneous update**

$\text{temp0} := \theta_0 - \alpha \dfrac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

$\text{temp1} := \theta_1 - \alpha \dfrac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$

$\theta_0 := \text{temp0}$

$\theta_1 := \text{temp1}$

**Incorrect:**

$\text{temp0} := \theta_0 - \alpha \dfrac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

$\theta_0 := \text{temp0}$

$\text{temp1} := \theta_1 - \alpha \dfrac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$

$\theta_1 := \text{temp1}$

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

# Learning rate



Big learning rate

Small learning rate

# Gradient descent for linear regression

Repeat until convergence{

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \text{(for } j = 0 \text{ and } j = 1)$$

}

- Linear regression model

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

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)^2$$

# Computing partial derivative

- $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)^2$

$$= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^{m} \left( \theta_0 + \theta_1 x^{(i)} - y^{(i)} \right)^2$$

- $j = 0:\ \ \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)$

- $j = 1:\ \ \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right) x^{(i)}$

# Gradient descent for linear regression

Repeat until convergence{

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right) x^{(i)}$$

}

Update $\theta_0$ and $\theta_1$ simultaneously



Fit at iteration 0

# Batch gradient descent

- "Batch": Each step of gradient descent uses all the training examples

Repeat until convergence{

$m$: Number of training examples

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right) x^{(i)}$$

}

# Mini Batch gradient descent

Say $b = 10, m = 1000$.

Repeat {

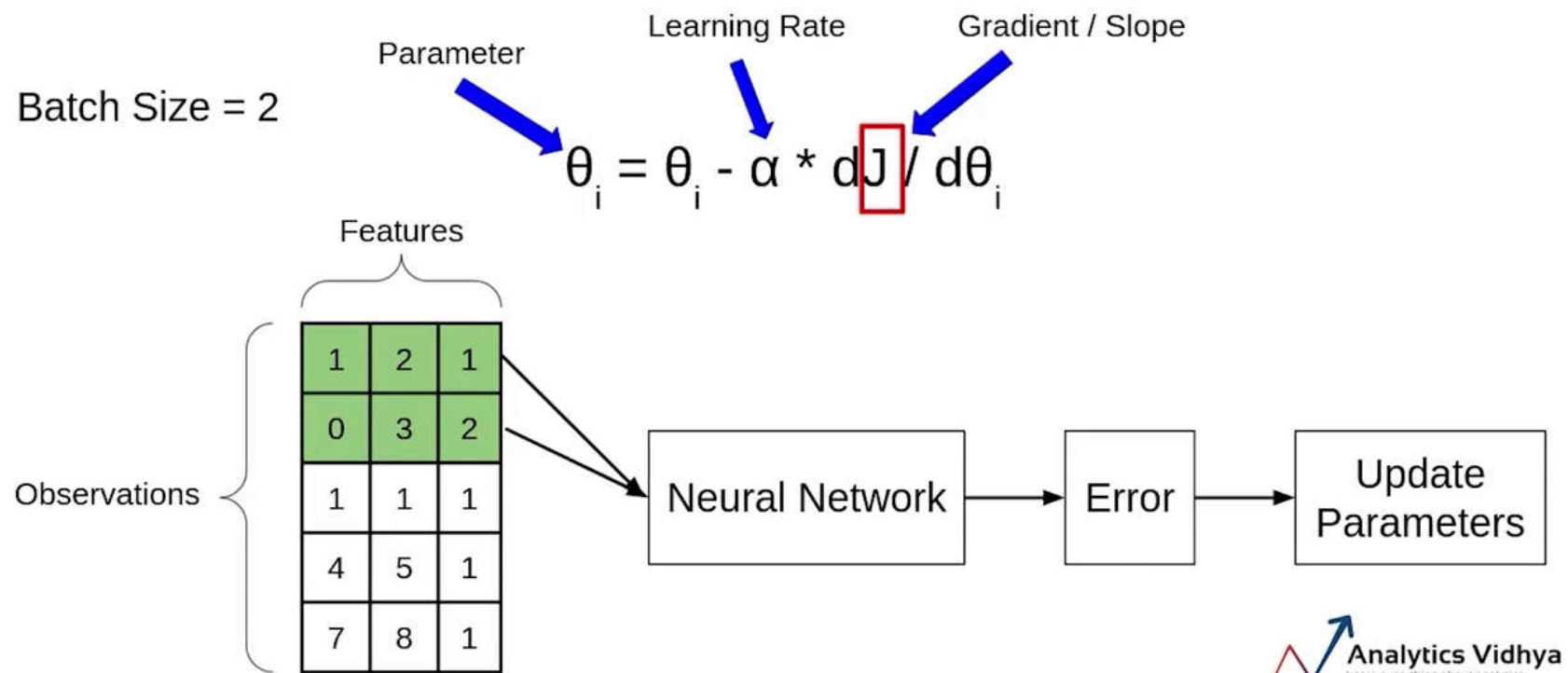    for $i = 1, 11, 21, 31, \ldots, 991$ {

$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

        (for every $j = 0, \ldots, n$)

    }

}

# Mini Batch gradient descent

# Stochastic Batch gradient descent

**Algorithm 2:** Pseudo-code for SGD

**Function** SGD:

    Set epsilon as the limit of convergence

    **for** $i = 1, \ldots, m$ **do**

        **for** $j = 0, \ldots, n$ **do**

            **while** $|\omega_{j+1} - \omega_j| < epsilon$ **do**

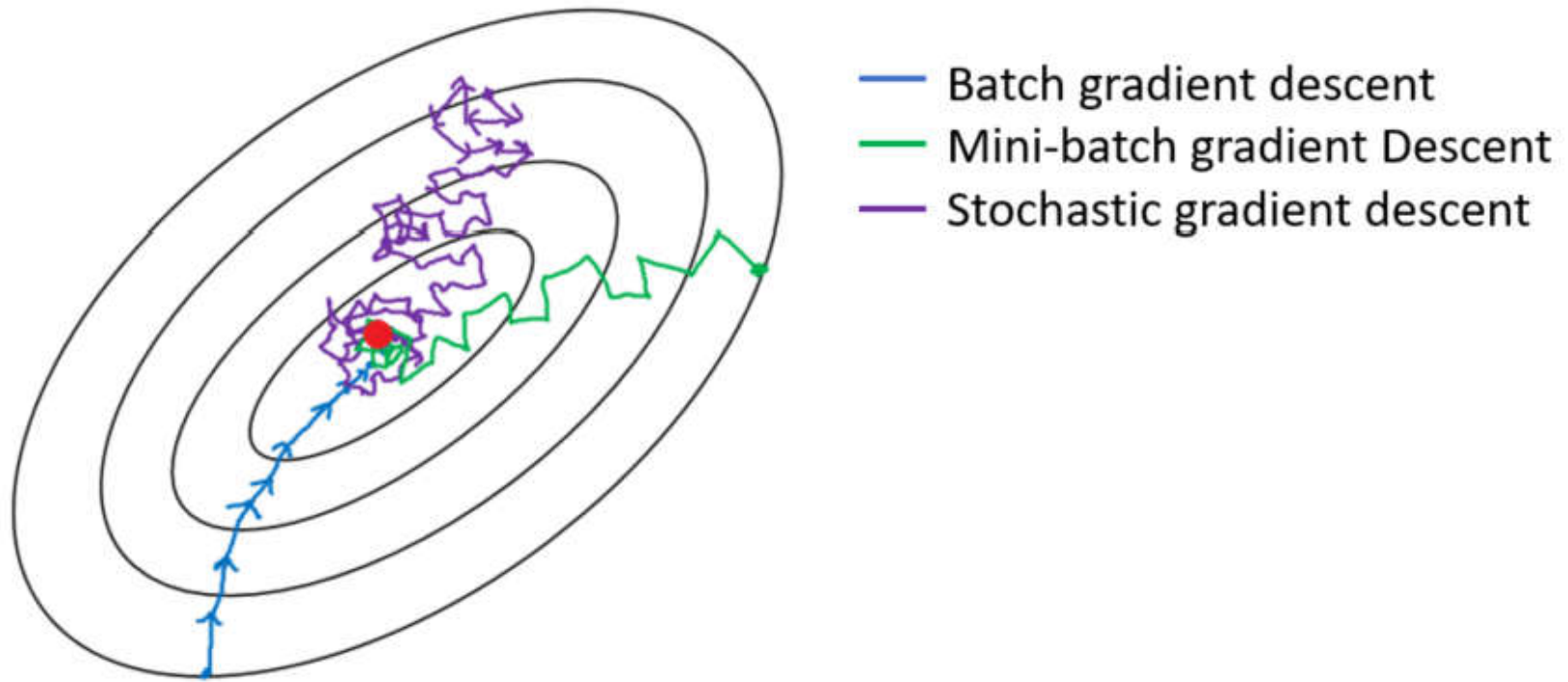                $\omega_{j+1} := \omega_j - \alpha \cdot (h_\omega(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)};$

            **end**

        **end**

    **end**

# GD VS MGD VS SGD

# Linear Regression

- Model representation

- Cost function

- Gradient descent

- **Features and polynomial regression**

- Normal equation

# Training dataset

| Size in feet^2 (x) | Price ($) in 1000's (y) |
|:---:|:---:|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

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

Slide credit: Andrew Ng

# Multiple features (input variables)

| Size in feet² (x₁) | Number of bedrooms (x₂) | Number of floors (x₃) | Age of home in years (x₄) | Price ($) in 1000's (y) |
|---|---|---|---|---|
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| ... | | | | ... |

Notation:

$n$ = Number of features

$x^{(i)}$ = Input features of $i^{th}$ training example

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

$x_3^{(2)} =?$

$x_3^{(4)} =?$

# Hypothesis

Previously:  $$h_\theta(x) = \theta_0 + \theta_1 x$$

Now:  $$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

- For convenience of notation, define $x_0 = 1$
  ($x_0^{(i)} = 1$ for all examples)

- $x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in R^{n+1}$ $\qquad$ $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in R^{n+1}$

- $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$
  $\qquad = \theta^\top x$

# Gradient descent

- Previously ($n = 1$)

Repeat until convergence{

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right) x^{(i)}$$

}

- New algorithm ($n \geq 1$)

Repeat until convergence{

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right) x_j^{(i)}$$
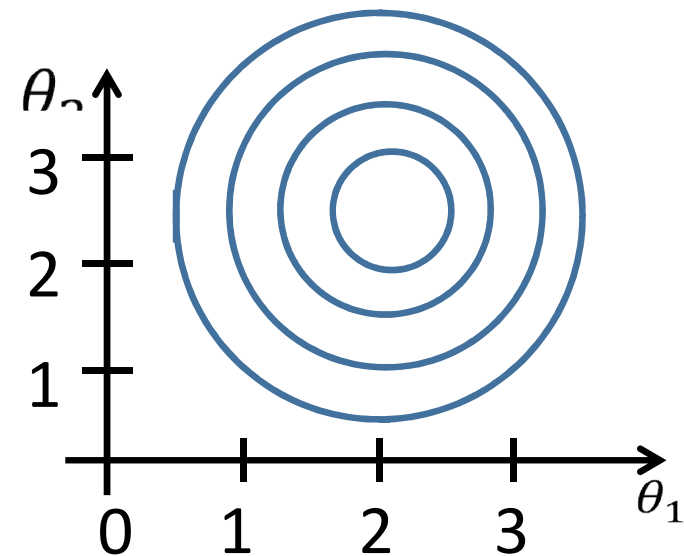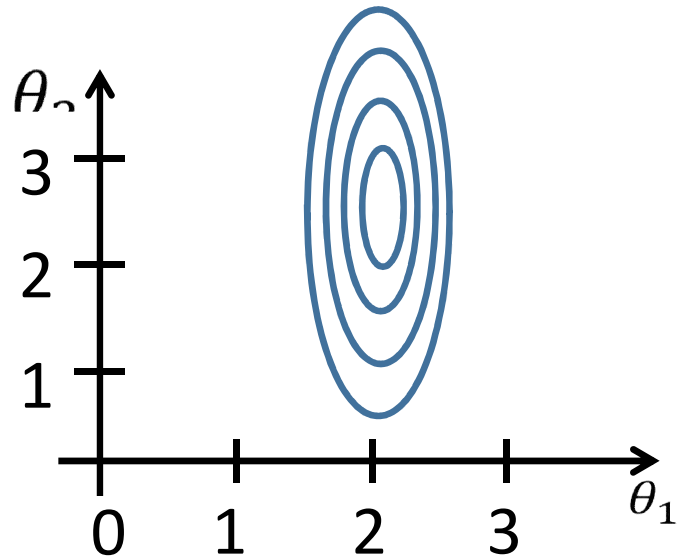
}

Simultaneously update
$\theta_j$, for $j = 0, 1, \cdots, n$

# Gradient descent in practice: Feature scaling

- Idea: Make sure features are on a similar scale (e.g,. $-1 \leq x_i \leq 1$)
- E.g. $x_1$ = size (0-2000 feat^2)

  $x_2$ = number of bedrooms (1-5)

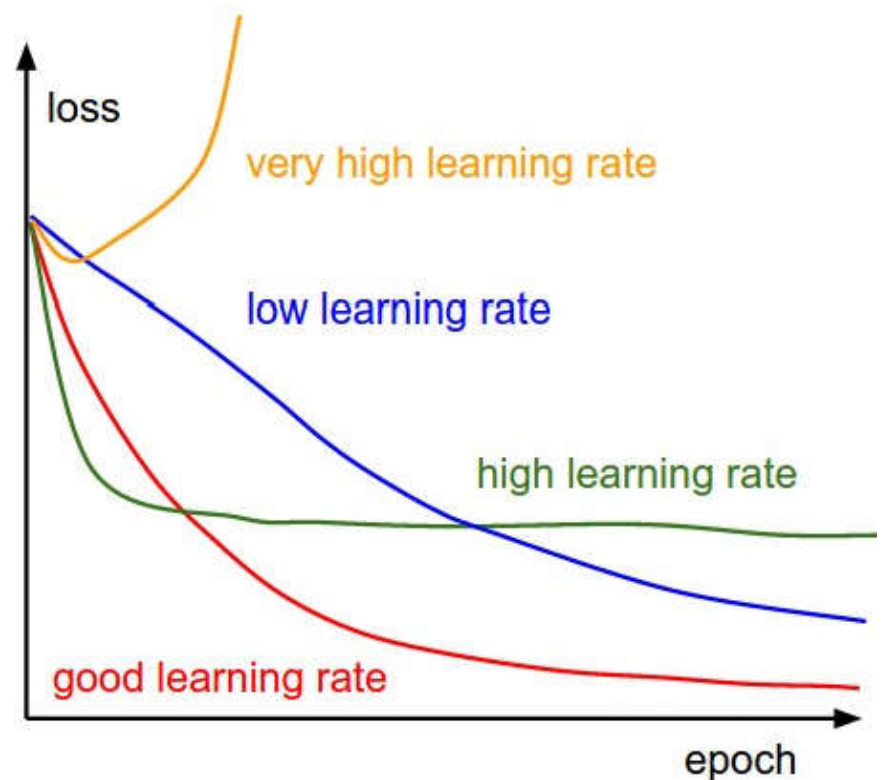# Gradient descent in practice: Learning rate

- Automatic convergence test
- $\alpha$ too small: slow convergence
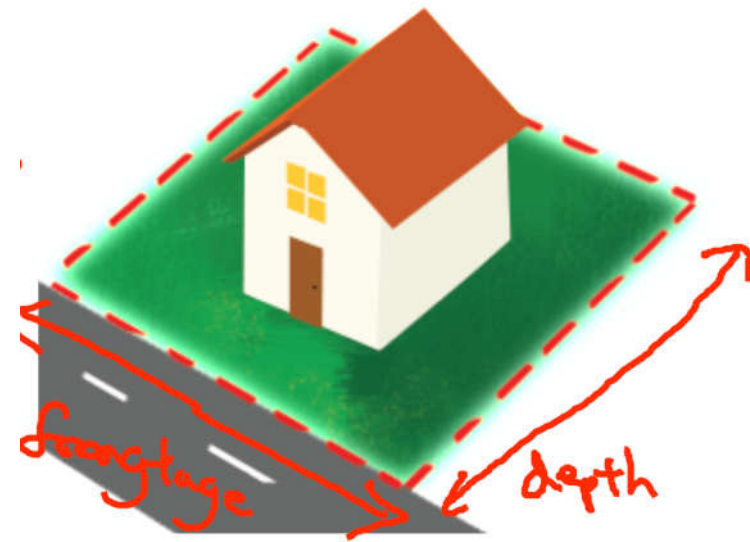- $\alpha$ too large: may not converge

- To choose $\alpha$, try

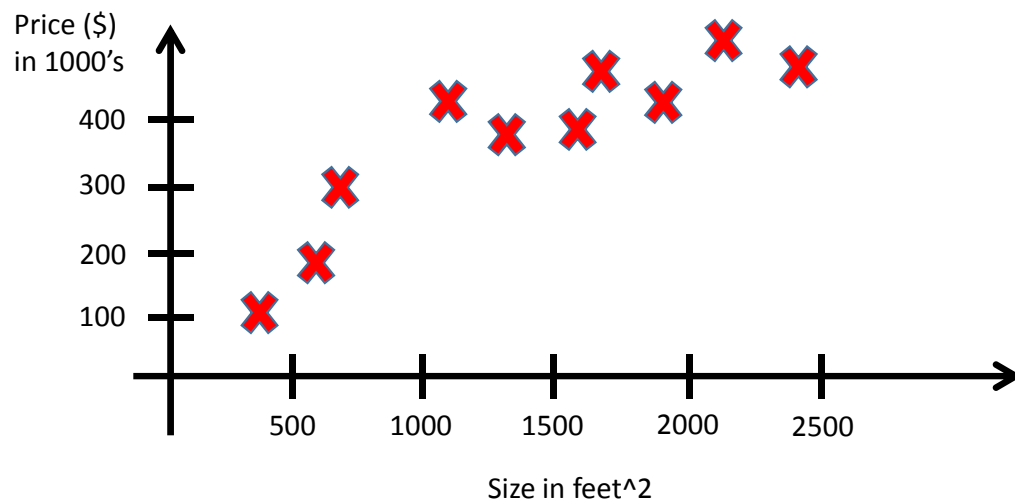0.001, … 0.01, …, 0.1, … , 1



Image credit: CS231n@Stanford

# House prices prediction

- $h_\theta(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth}$

- Area
  $x = \text{frontage} \times \text{depth}$

- $h_\theta(x) = \theta_0 + \theta_1 x$

# Polynomial regression



$x_1 = (size)$

$x_2 = (size)^2$

$x_3 = (size)^3$

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

  $= \theta_0 + \theta_1 (size) + \theta_2 (size)^2 + \theta_3 (size)^3$

# Linear Regression

- Model representation

- Cost function

- Gradient descent

- Features and polynomial regression

- **Normal equation**

| $(x_0)$ | Size ... | ... | Number ... $(x_3)$ | Age of home (years) $(x_4)$ | Price ($) in 1000's (y) |
|---|---|---|---|---|---|
| 1 | 2104 | 5 | 1 | 45 | 460 |
| 1 | 1416 | 3 | 2 | 40 | 232 |
| 1 | 1534 | 3 | 2 | 30 | 315 |
| 1 | 852 | 2 | 1 | 36 | 178 |
| | ... | | | | ... |

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

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

# $m$ training examples, $n$ features

## Gradient Descent

- Need to choose $\alpha$
- Need many iterations
- Works well even when $n$ is large

## Normal Equation

- No need to choose $\alpha$
- Don't need to iterate
- Need to compute
$$(X^\top X)^{-1}$$
- Slow if $n$ is very large

# Things to remember

- **Model representation**

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n = \theta^\top x$$

- **Cost function**  $J(\theta) = \frac{1}{2m}\sum_{i=1}^{m}\left(h_\theta\left(x^{(i)}\right) - y^{(i)}\right)^2$

- **Gradient descent for linear regression**

Repeat until convergence $\{\theta_j := \theta_j - \alpha \frac{1}{m}\sum_{i=1}^{m}\left(h_\theta\left(x^{(i)}\right) - y^{(i)}\right) x_j^{(i)}\}$

- **Features and polynomial regression**

Can combine features; can use different functions to generate features (e.g., polynomial)

- **Normal equation**  $\theta = (X^\top X)^{-1} X^\top y$