

Machine Learning

Linear Regression with
multiple variables

Multiple features

Multiple features (variables).

Size (feet ²)	Price (\$1000)
x	y
2104	460
1416	232
1534	315
852	178
...	...

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

Multiple features (variables).

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

$m = 47$

Notation:

- n = number of features $n = 4$
- $x^{(i)}$ = input (features) of i^{th} training example.
- $x_j^{(i)}$ = value of feature j in i^{th} training example.

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

$x_3^{(2)} = 2$

Hypothesis:

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

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

e.g. $\underline{h_{\theta}(x)} = \underline{80} + \underline{0.1}x_1 + \underline{0.01}x_2 + 3x_3 - 2x_4$

↑ ↑ ↑
age

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

For convenience of notation, define $x_0 = 1$. ($x_0^{(i)} = 1$)

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

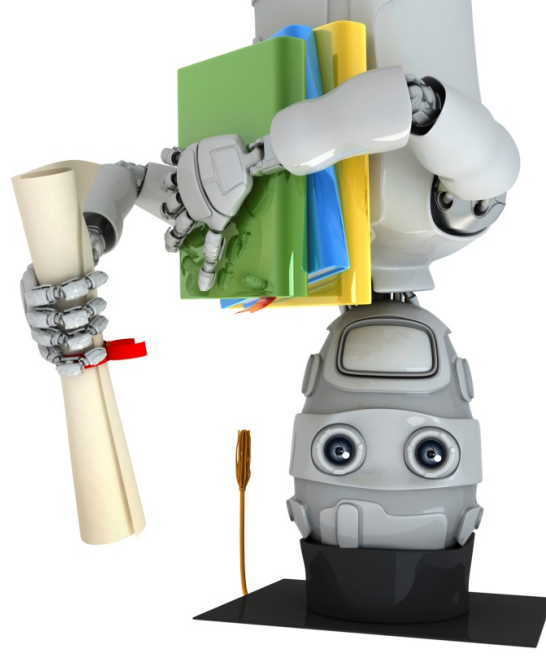
$$= \theta^T x$$

θ^T is a $(n+1) \times 1$ matrix.
 x is a $(n+1) \times 1$ vector.

Multivariate linear regression. \leftarrow

Linear Regression with
multiple variables

Gradient descent for
multiple variables



Machine Learning

Hypothesis: $h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ $\nearrow x_0 = 1$

Parameters: $\theta_0, \theta_1, \dots, \theta_n$ θ $n+1$ -dimensional vector

Cost function:

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

Gradient descent:

Repeat {

$\rightarrow \theta_j := \theta_j - \alpha \left[\frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n) \right]_{J(\theta)}$

(simultaneously update for every $j = 0, \dots, n$)

Gradient Descent

Previously (n=1):

Repeat {

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

$$\frac{\partial}{\partial \theta_0} J(\theta)$$

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

(simultaneously update θ_0, θ_1)

}

New algorithm ($n \geq 1$):

Repeat {

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

(simultaneously update θ_j for $j = 0, \dots, n$)

}

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})x_0^{(i)}$$

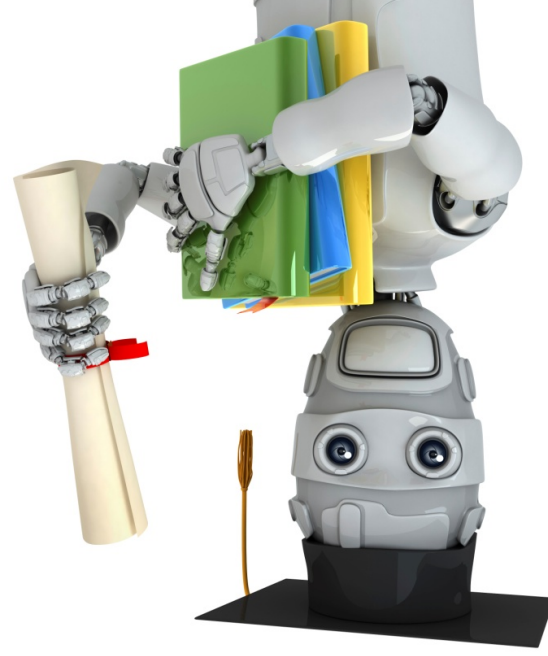
$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})x_1^{(i)}$$

$$\rightarrow \theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})x_2^{(i)}$$

...

Linear Regression with multiple variables

Gradient descent in
practice I: Feature Scaling



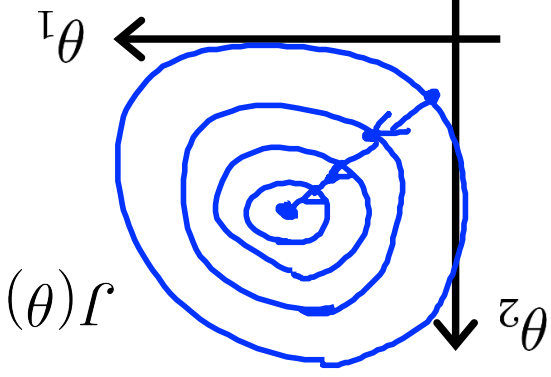
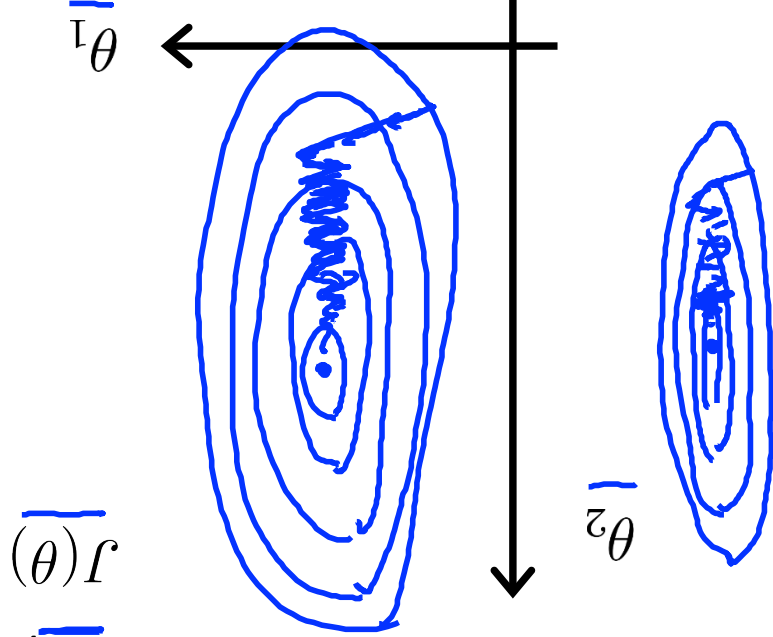
Machine Learning

Feature Scaling

Idea: Make sure features are on a similar scale.

E.g. x_1 = size (0-2000 feet²)

x_2 = number of bedrooms (1-5)



$$x_1 = \frac{\text{size (feet}^2\text{)}}{2000}$$

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

$$0 \leq x_1 \leq 1 \quad 0 \leq x_2 \leq 1$$

Feature Scaling

Get every feature into approximately a $[-1, 1]$ range.

$$x_0 = 1$$

$$\begin{aligned} &\checkmark \quad 0 \leq x_1 \leq 3 \\ &\checkmark \quad -2 \leq x_2 \leq 503 \end{aligned}$$

$$x_3 \leq \overline{100} \quad \checkmark$$

$$-0.0001 \leq x_4 \leq \boxed{1000.0} \quad \checkmark$$

$$\begin{aligned} &\checkmark \quad -3 \leq x_5 \leq 0.5 \\ &\checkmark \quad -\frac{1}{2} \leq x_6 \leq \frac{1}{2} \end{aligned}$$

Mean normalization

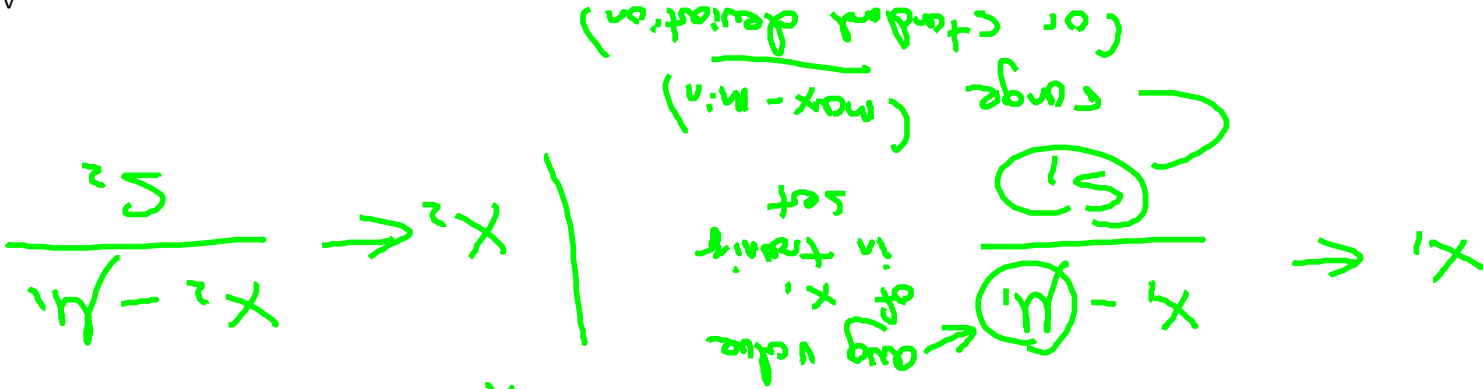
Replace x_i with $x_i - \mu_i$ to make features have approximately zero mean
(Do not apply to $x_0 = 1$).

E.g. $x_1 = \frac{size - 1000}{2000}$

$x_2 = \frac{\#bedrooms - 2}{4}$

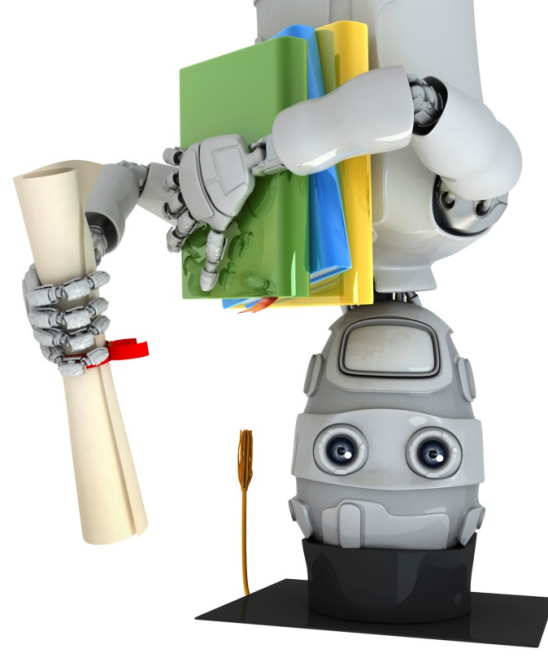
$-0.5 \leq x_1 \leq 0.5$ $-0.5 \leq x_2 \leq 0.5$

avg size = 1000
avg #bedrooms = 1.5



Linear Regression with multiple variables

Gradient descent in
practice II: Learning rate



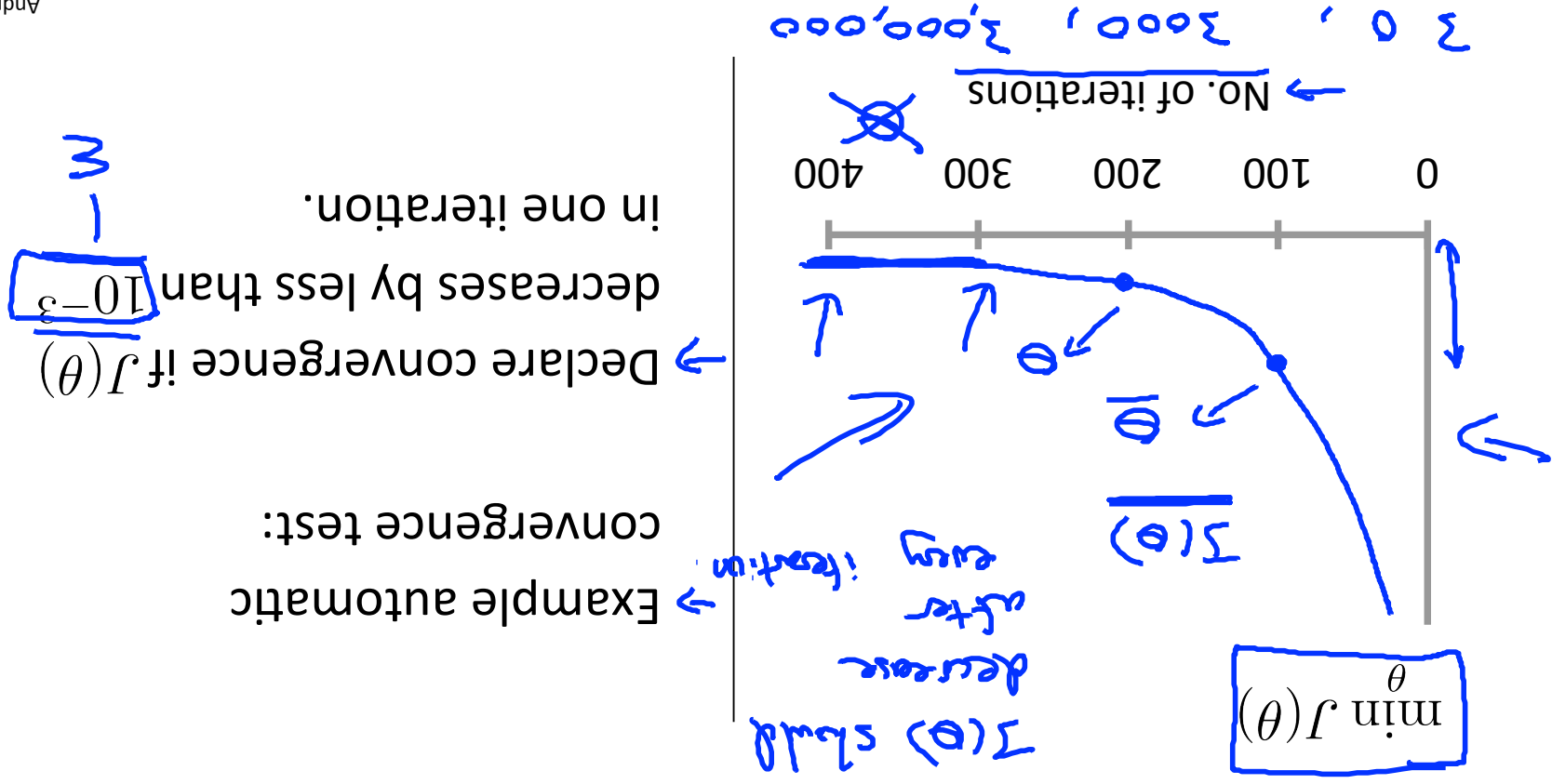
Machine Learning

Gradient descent

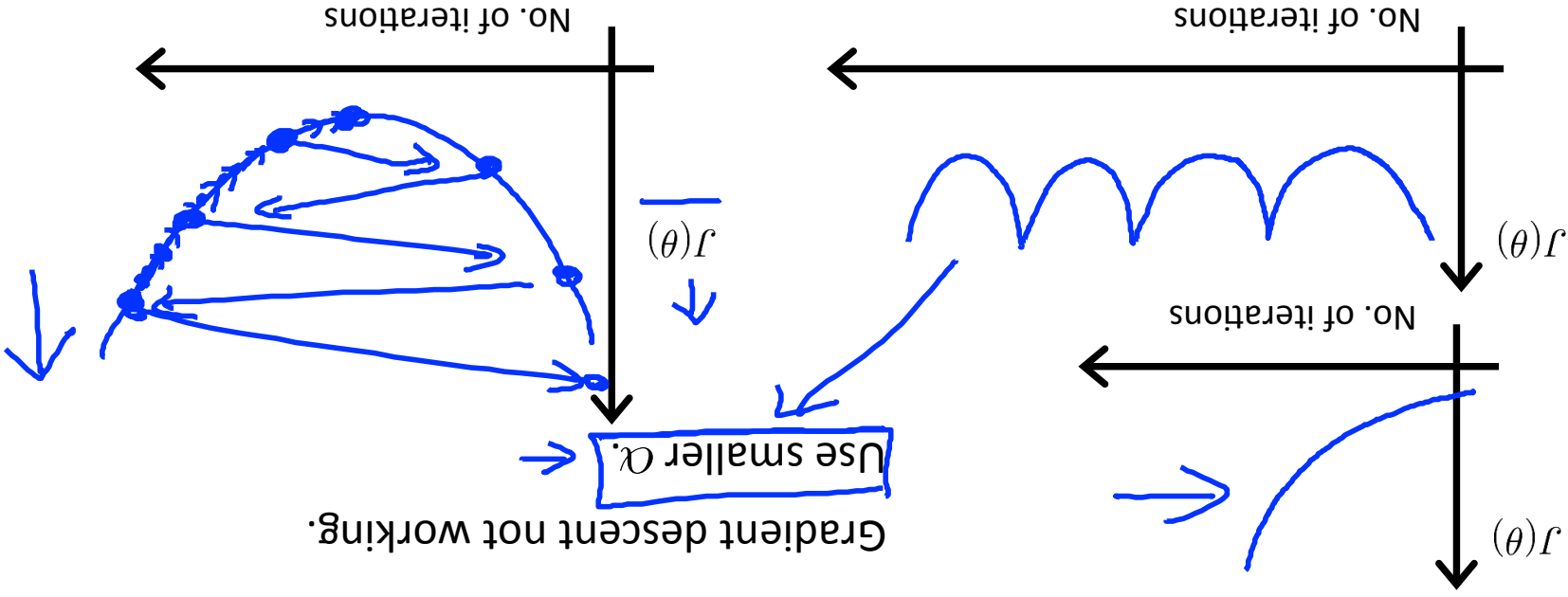
$$\leftarrow \theta_j := \theta_j - \alpha \frac{\partial \theta_j}{\partial} f(\theta)$$

- “Debugging”: How to make sure gradient descent is working correctly.
- How to choose learning rate α .

Making sure gradient descent is working correctly.



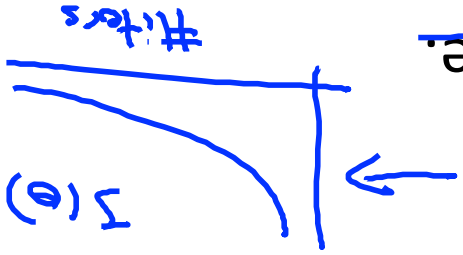
Making sure gradient descent is working correctly.



- For sufficiently small α , $J(\theta)$ should decrease on every iteration.
- But if α is too small, gradient descent can be slow to converge.

Summary:

- If α is too small: slow convergence.
- If α is too large: $J(\theta)$ may not decrease on every iteration; may not converge. (Slow converge also possible)



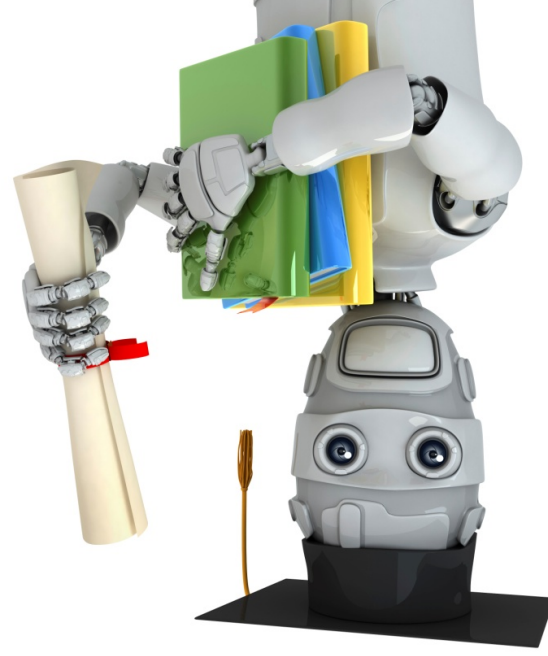
To choose α , try

$\dots, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, \dots$

Handwritten annotations above the sequence include arrows pointing to each value and labels: $1 \times$ above 0.001, $2 \times$ above 0.003, $2 \times$ above 0.01, $3 \times$ above 0.03, $3 \times$ above 0.1, and $3 \times$ above 0.3.

Linear Regression with multiple variables

Features and polynomial regression



Machine Learning

Housing prices prediction

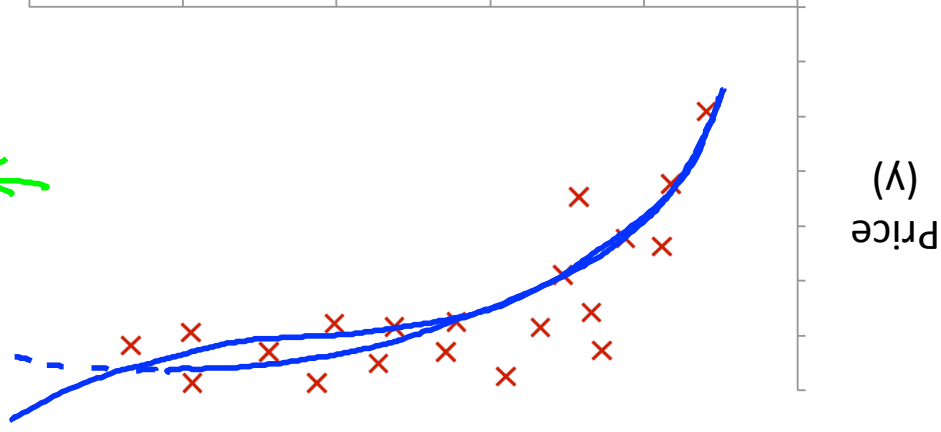
$$h_{\theta}(x) = \theta_0 + \theta_1 \times \overbrace{\text{frontage}}^{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$



$$\underline{\text{Area}} \quad x = \text{frontage} \times \text{depth}$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad \downarrow \text{land area}$$

Polynomial regression



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

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

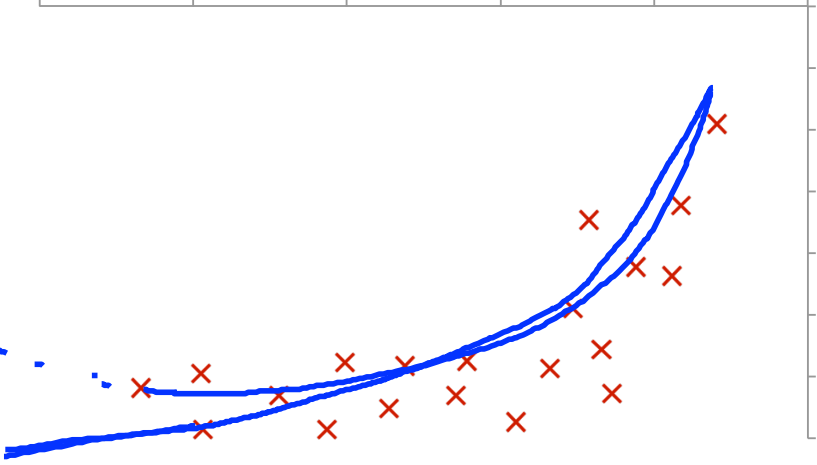
$$\theta_0 + \theta_1 \overline{size} + \theta_2 \overline{size}^2 + \theta_3 \overline{size}^3$$

$$\begin{cases} x_1 = size \\ x_2 = size^2 \\ x_3 = size^3 \end{cases}$$

Size: 1-1550
Size²: 1-1000,000
Size³: 1-10⁹

Choice of features

Price
(y)



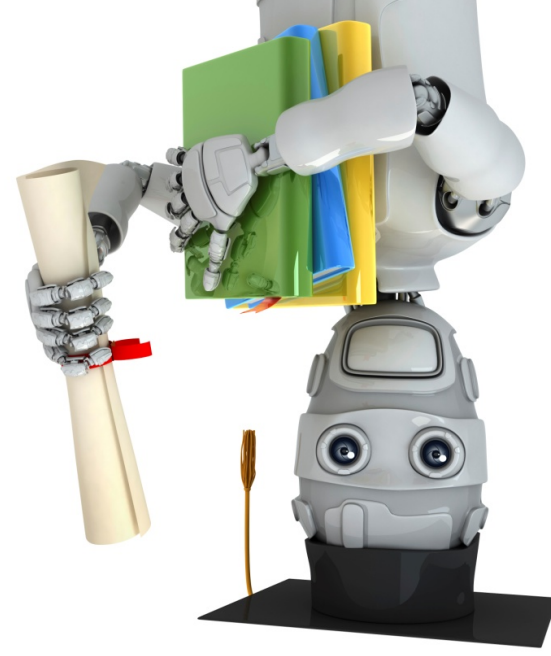
Size (x)

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

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

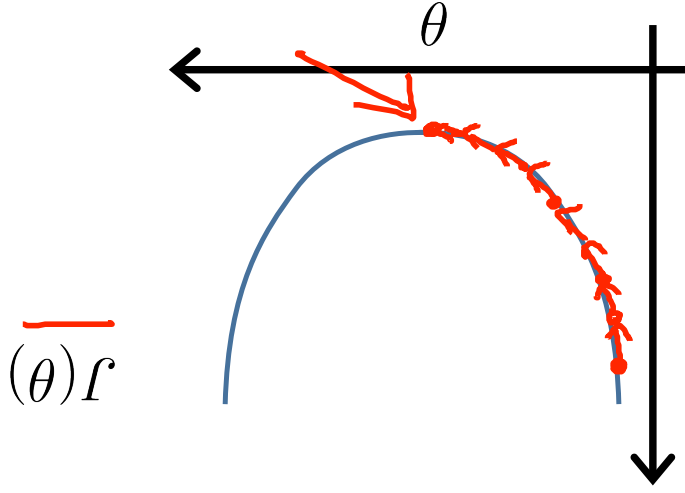
Linear Regression with multiple variables

Normal equation



Machine Learning

Gradient Descent



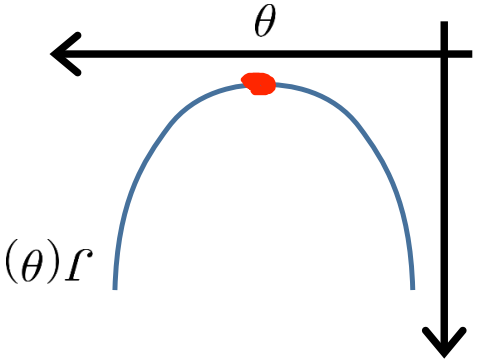
Normal equation: Method to solve for θ
analytically.

Intuition: If 1D ($\theta \in \mathbb{R}$)

$$\rightarrow J(\theta) = a\theta^2 + b\theta + c$$

$$\frac{d}{d\theta} J(\theta) = \dots \stackrel{\text{set}}{=} 0$$

Solve for θ



$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_m) = \frac{1}{m} \sum_{i=1}^m z^{(i)}_2 - y^{(i)}_2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \quad (\text{for every } j)$$

Solve for $\theta_0, \theta_1, \dots, \theta_n$

Examples: $m = 4$.

	x_0	x_1	x_2	x_3	x_4	y
Size (feet ²)	1	2104	5	1	45	460
Number of bedrooms	1	1416	3	2	40	232
Number of floors	1	1534	3	2	30	315
Age of home (years)	1	852	2	1	36	178
Price (\$1000)						

$\bar{X} =$ $\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 2104 & 1416 & 1534 & 852 & 45 \\ 5 & 3 & 3 & 2 & 1 \\ 1 & 2 & 2 & 1 & 36 \end{bmatrix}$ $m \times (n+1)$
 $\bar{y} =$ $\begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$ $m \times 1$ vector
 $\theta = (X^T X)^{-1} X^T y$

m examples $(x_{(1)}, y_{(1)}), \dots, (x_{(m)}, y_{(m)})$; n features.

$$\underline{x}_{(i)} = \begin{bmatrix} x_{(i)}^0 \\ x_{(i)}^1 \\ x_{(i)}^2 \\ \vdots \\ x_{(i)}^n \end{bmatrix} \in \mathbb{R}^{n+1}$$

X (design matrix)

$$\begin{bmatrix} (x_{(1)})^T & & \\ & \ddots & \\ (x_{(m)})^T & & \end{bmatrix} = \begin{bmatrix} x_{(1)}^0 & x_{(1)}^1 & \dots & x_{(1)}^n \\ x_{(2)}^0 & x_{(2)}^1 & \dots & x_{(2)}^n \\ \vdots & \vdots & \ddots & \vdots \\ x_{(m)}^0 & x_{(m)}^1 & \dots & x_{(m)}^n \end{bmatrix} \quad \begin{matrix} n \times (n+1) \\ m \times 2 \end{matrix}$$

$$= \begin{bmatrix} \underline{x}_{(1)} \\ \underline{x}_{(2)} \\ \vdots \\ \underline{x}_{(m)} \end{bmatrix}$$

E.g. If $\underline{x}_{(i)} = \begin{bmatrix} x_{(i)}^0 \\ x_{(i)}^1 \\ \vdots \\ x_{(i)}^n \end{bmatrix}$ then $\underline{x}_{(i)}^T \underline{x}_{(i)} = x_{(i)}^0^2 + x_{(i)}^1^2 + \dots + x_{(i)}^n^2$

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

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

$(X^T X)^{-1}$ is inverse of matrix $X^T X$.

Set $A = X^T X$

$$(X^T X)^{-1} = A^{-1}$$

Octave: `pinv(X'*X)*X'*y`

$$\text{pinv}(X^T X) * X^T * y$$

min $J(\theta)$

$$X' \quad \left| \quad \begin{array}{l} 0 \leq x_1 \leq 1 \\ 0 \leq x_2 \leq 1000 \\ 0 \leq x_3 \leq 10^{-5} \end{array} \right.$$

\bar{m} training examples, \bar{n} features.

Gradient Descent

- \rightarrow Need to choose α .
- \rightarrow Needs many iterations.

• Works well even when \bar{n} is large.

$$\bar{n} \approx 10^6$$

\rightarrow

Normal Equation

- \rightarrow No need to choose α .

- \rightarrow Don't need to iterate.

• Need to compute

$$\frac{1}{n \times n} (X^T X - I)$$

- Slow if \bar{n} is very large.

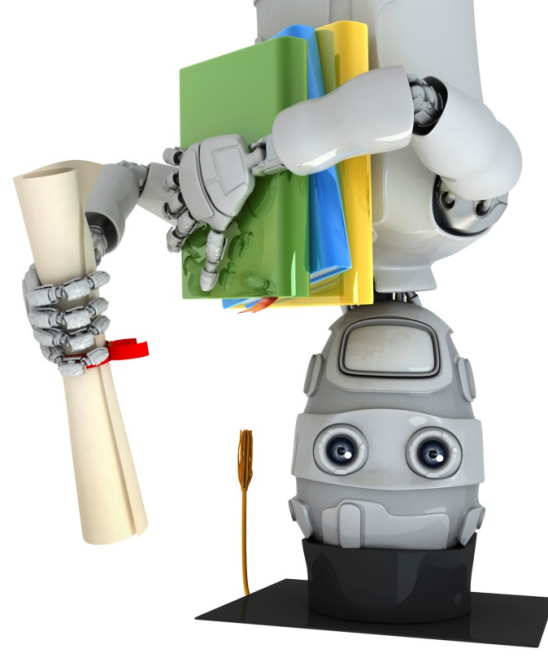
$$n = 100$$

$$n = 1000$$

$$n \approx 10000$$

Linear Regression with
multiple variables

Normal equation
and non-invertibility
(optional)



Machine Learning

Normal equation

$$(X^T X)^{-1} X^T y$$

$$X^T X$$

- What if $X^T X$ is non-invertible? (singular/degenerate)

- Octave: $\text{pinv}(X, X)$ and $\text{pinv}(X, Y)$



What if $X^T X$ is non-invertible?

- Redundant features (linearly dependent).

E.g. $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ = size in feet

~~x_2 = size in m²~~

$$x_1 = (3.28)^2 x_2$$

- Too many features (e.g. $m \leq n$).

$\rightarrow n = 10$
 $\rightarrow n = 100$
 $\odot \in \mathbb{R}^{101}$

- Delete some features, or use regularization.

↓
later