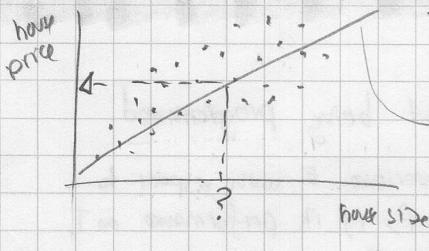


Linear regression of One Variable

(if discrete, it would be a classification problem) ②

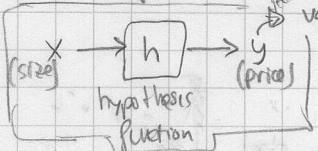


Supervised learning → we have a data-set of price/size
Regression → predict continuous output

Model

Training Set +

Learning Algorithm



(x_i, y_i) - one training example

$(x^{(2)}, y^{(2)})$ - 2nd training example

Size	Price
2104	460
1416	232
1534	315
:	:

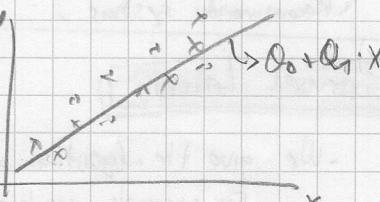
m: number of training examples

x: input variable (features)
y: output variable (target)

hypothesis function (one variable)

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

(for (x)) model parameters



Linear Regression with one variable
or
Univariate linear regression

How do we find θ_0 & θ_1 parameter values?

Informally, we want θ_0, θ_1 so that $h(x)$ is close to y for the training examples (x_i, y_i)

Formally:

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

\Downarrow

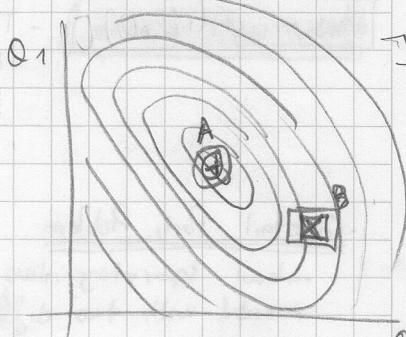
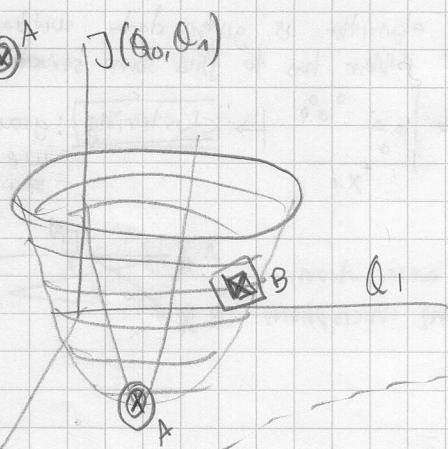
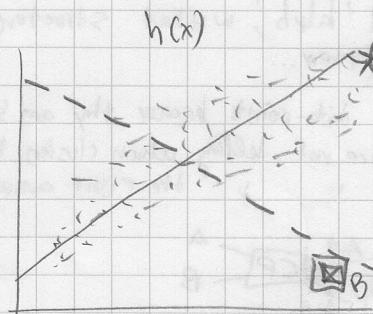
$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

} Objective function to minimize

} Cost function to minimize

Square Error function → most common for regression problems

another notation



$J(\theta_0, \theta_1)$
contour plot

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

Cost($h_{\theta}(x)$, y) → sum, over m, of cost

The cost is what we want the algorithm to penalize; in this case we penalize with the square error

Gradient Descent

Algorithm to minimize, here used to minimize $J(\theta_0, \dots, \theta_n)$ and get
(in general)

the right $\theta_0, \dots, \theta_n$ parameter values that better
match the regression line and the data-set.

Idea: ① initialize θ_0, θ_1 (to 0, for example)

② keep changing θ_0, θ_1 so J is reduced $\sim \sim \sim$ minimum of J will be found

Formally:

$$\text{assign max operator } (\text{if } \theta' = 1) \rightarrow \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)}) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

repeat until convergence {

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

}

learning
rate

when the derivative

is 0, we stop.

- it was the minimum!

- θ_j does not move

partial
derivative

the larger the derivative,
the larger the slope, so
the larger the jump in θ_j .

Simultaneously update θ_0, θ_1

$$\text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J$$

$$\text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J$$

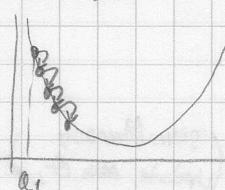
$$\theta_0 := \text{temp}_0$$

$$\theta_1 := \text{temp}_1$$

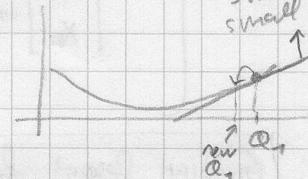
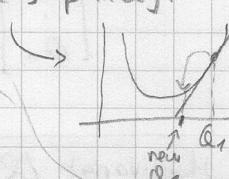
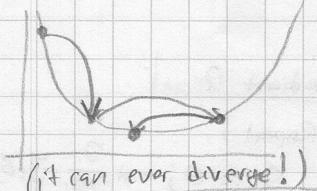
large slope
large derivative

small slope
small derivative

α small

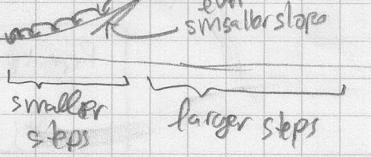


α large



we don't have to change the value of α as we advance, the derivative part already results in smaller steps the closer to the minimum we get

One problem of GD is that it might find different local minima...



repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \quad (\text{update } \theta_0, \theta_1)$$

$$\theta_1 := \theta_1 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \quad (\text{update } \theta_0, \theta_1)$$

}

simultaneously

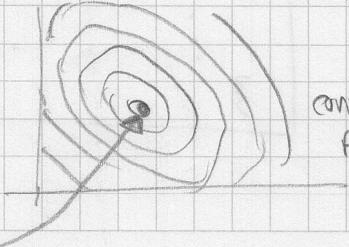
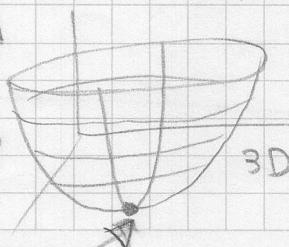
called "Batch" gradient
descent because at each
step we use the
entire set of training
samples ($\frac{m}{2}$)

Applying Gradient Descent to Univariate Linear Regression

o does not have the local minima problem because the $J(\theta)$ function

is Bowl-shaped

↓
single local minimum
= global minimum



Linear Regression with Multiple Variables

(or Multivariate LR)

Multiple variables = Multiple features

→ house size, num of floors, num of dorms, ...

	size	dorms	floors	years	price
2104	5	1	45	460	
1616	3	2	40	232	
1534	3	2	30	315	
:	:	:	:	:	
x_1	x_2	x_3	x_4	y	

Vector $X^{(2)}$

$n=4$ features

output variable
to predict

$x_3^{(2)}$

size dorms floors years

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

$$\theta_0 \cdot x_0 \\ \frac{\partial}{\partial} x_0 = 1$$

$$\left\{ \begin{array}{l} x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\ \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \end{array} \right\} \Rightarrow h_{\theta}(x) = \theta^T \cdot x = \begin{bmatrix} \theta_0 & \theta_1 & \theta_2 & \theta_3 & \theta_4 \end{bmatrix} \cdot \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Gradient Descent for Multivariate LR

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

parameters $\rightarrow \theta$

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

Gradient Descent:

$$\text{Repeat } \left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \right\} \quad \begin{array}{l} \text{(simultaneously)} \\ \text{(update each } \theta_j \text{)} \end{array}$$

Repeat

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \quad \begin{array}{l} \text{(simultaneously)} \\ \text{update for all } j \end{array}$$

Repeat

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

(simultaneously)
(update all j)

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

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

⋮

feature scaling so every feature fits approx. into $-1 \leq x_i \leq +1$

Practical
notes

$$x_i \leftarrow \frac{x_i}{\max}$$

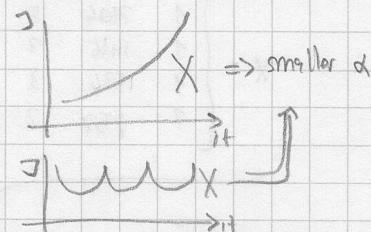
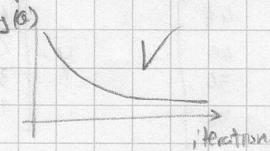
mean normalization → so every feature has approx. zero mean

$$x_i \leftarrow \frac{x_i - \bar{x}_i}{S_i}$$

S_i → range, or standard dev.

Gradient Descent
evaluates better

debugging → plots of the value of $J(\theta)$ versus iterations



value of alpha? → debug J

$\alpha = 0.001$ $\xrightarrow{x_3} 0.01$ $\xrightarrow{x_3} 0.1$ $\xrightarrow{x_3} 1$... } plot $J(\theta)$ vs iterations
then decrease $\alpha = 0.003$ $\xrightarrow{x_3} 0.03$ $\xrightarrow{x_3} 0.3$ $\xrightarrow{x_3} 3$... }

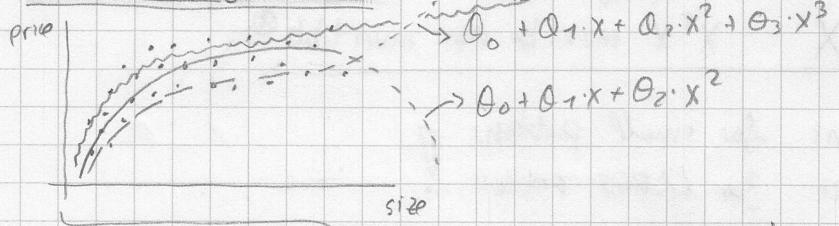
Choosing the right features (x_i)

• if we have width and height, better work with area

$$h_\theta(x) = \theta_0 + \theta_1 \cdot w + \theta_2 \cdot h \Rightarrow h_\theta(x) = \theta_0 + \theta_1 \cdot a$$

} combining
features

Polynomial Regression



We can obtain this regression with the Multivariate Linear Regression

$$h_\theta(x) = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \theta_3 \cdot x_3 = \\ - \theta_0 + \theta_1 \cdot (\text{size}) + \theta_2 \cdot (\text{size}^2) + \theta_3 \cdot (\text{size}^3)$$

$$x_1 = \text{size}$$

$$x_2 = \text{size}^2$$

$$x_3 = \text{size}^3$$

} be careful with scaling! apply
feature scaling!

Normal Equation

Method to solve for θ analytically (alternative to Gradient Descent Algorithm)
 ↳ in one 'mathematical' step

	size	bedrooms	floors	years	price
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_0 x_1 x_2 x_3 x_4 | y
 ↓ ↓ ↓ ↓ ↓

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

$$\boxed{\theta = (X^T X)^{-1} \cdot X^T y} \quad \rightsquigarrow \theta \text{ values that minimize } J(\theta)$$

Gradient Descent or Normal Equation ?

Need to choose α	X	V	No parameters to choose
Needs many iterations	X	V	No iterations
Feature scaling	V	V	No need to scale features
Works well for large n (features) ($n \geq 10000$)	V	X	Slow if n large \rightarrow Needs to compute $(X^T X)^{-1}$ $(n \leq 10000) \sim O(n^3)$
	- X		$X^T X$ might be non invertible

Normal Equations for small problems

Gradient Descent for LARGE problems !

④ Octave 'pinv' will manage it

- if it has redundant features
 - if $m \leq n$
 - ↳ reduce features (n)
 - ↑ regularization
- num
data
experiments