

Supervised Learning

- + Definitions
- + Linear Regression
- + Batch & stochastic gradient descent
- + Normal Equations

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Supervised learning

Prediction

$$h: X \rightarrow Y$$

Image

Contains car

TEXT

IS HAVE SPEECH?

Have Data

PRICE

Given: Training Set

$$\{(x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})\} \quad x^{(i)} \in X, y^{(i)} \in Y$$

Do: find "good" $h: X \rightarrow Y$ hypothesis

this job of training algorithm

WE USE h ON NEW DATA (x)

Call this Prediction, WE ARE VERY INTERESTED IN $x \notin$ TRAINING SET

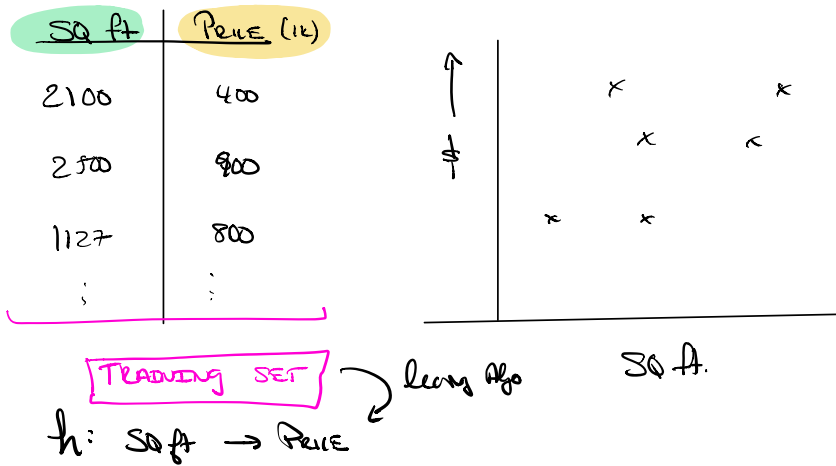
if y IS DISCRETE \Rightarrow classification
 y IS CONTINUOUS \Rightarrow Regression

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Example Data (House Prices)



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

How do we represent h_0 ?

$$h(x) = \theta_0 + \theta_1 x_1 \quad (\text{affine fn.})$$

	x_1 Size	x_2 Bedroom	x_3 Lot size	Price
$x^{(1)}$	2104	4 $x_2^{(1)}$	452	400
$x^{(2)}$	2500	3	302	900

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$$

$$= \sum_{j=0}^3 \theta_j x_j \quad \text{NB } x_0 \text{ identically 1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} \quad x^{(1)} = \begin{bmatrix} x_0^{(1)} \\ x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{bmatrix}$$

Annotations: $x_0^{(1)} = 1$, $x_1^{(1)}$ is Size, $x_2^{(1)}$ is Bedroom, $x_3^{(1)}$ is Lot size, $y^{(1)}$ is Price

PARAMETERS

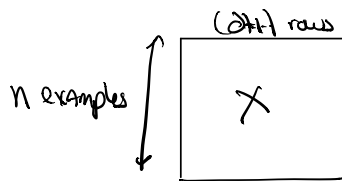
INPUTS / FEATURES

Output / Target

(x, y) is a training example

$(x^{(i)}, y^{(i)})$ is the i th example *i runs 1...n*

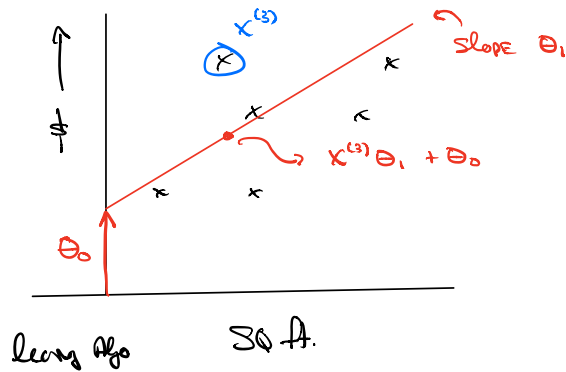
n examples and d features $\Rightarrow x^{(i)} \theta$ are $d+1$ dimensional



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



$$h_{\theta}(x) = \sum_{j=0}^2 \theta_j x_j \quad \text{WANT TO CHOOSE } \theta \text{ st. } h_{\theta}(x) \approx y$$

IDEA: $J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$ Cost function (least squares)

Assignment Project Exam Help

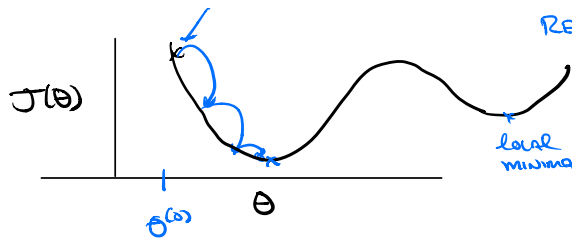
$$\min_{\theta} J(\theta)$$

<https://powcoder.com>

Add WeChat powcoder

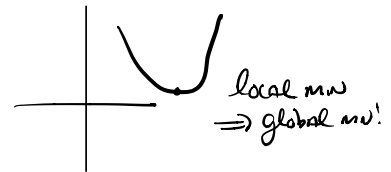
Gradient Descent

START $\theta^{(0)}$ AT RANDOM OR ZERO



REDUCE using Gradient

if J is nice (convex)



$$\theta^{(0)} := 0$$

$$\theta_j^{(t+1)} := \theta_j^{(t)} - \alpha \frac{\partial J(\theta^{(t)})}{\partial \theta_j}$$

$$j = 1, \dots, d$$

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$$h_{\theta}(x^{(i)}) = \theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_d x_d^{(i)}$$

$$\frac{\partial h_{\theta}(x^{(i)})}{\partial \theta_j} = x_j^{(i)}$$

$$\theta_j^{(t+1)} := \theta_j^{(t)} - \alpha \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

SOMETIMES WRITE AS $\theta^{(t+1)} := \theta^{(t)} - \alpha \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$

vector notation

BATCH VERSUS STOCHASTIC MINIBATCH

$$\Theta^{(t+1)} := \Theta^{(t)} - \alpha \sum_{i=1}^n (h_{\Theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

Minibatch: Randomly select $b < n$ points AND estimate gradient

1. Pick b points $\{i_1, \dots, i_b\} = B$
- 2.

$$\Theta^{(t+1)} := \Theta^{(t)} - \alpha \sum_{k \in B} (h_{\Theta}(x^{(k)}) - y^{(k)}) x^{(k)}$$

One detail Scale α AND α_b differently.

Tradeoff: Noisier BUT much faster

faster: Imagine if training set contains 100 copies of same point

\Rightarrow Not as ridiculous as it seems (near copies)

<https://powcoder.com>

How do you choose B ? Simply, whatever makes

Add WeChat powcoder

Normal Equation

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_0} J(\theta) \\ \frac{\partial}{\partial \theta_1} J(\theta) \\ \vdots \end{bmatrix}$$

$$A \in \mathbb{R}^{2 \times 2}$$

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$f: A \rightarrow \mathbb{R}$$

then

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial a_{11}} & \frac{\partial f}{\partial a_{12}} \\ \frac{\partial f}{\partial a_{21}} & \frac{\partial f}{\partial a_{22}} \end{bmatrix}$$

Now, we want to find minimum

$$\nabla_{\theta} J(\theta) = \vec{0} \quad (\nabla_{\theta} J(\theta) \in \mathbb{R}^{d+1})$$

<https://powcoder.com>

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

Add WeChat powcoder

$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times d} \quad \text{Design matrix}$$

$$X\theta = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(n)} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \end{bmatrix} = \begin{bmatrix} h_{\theta}(x^{(1)}) \\ \vdots \\ h_{\theta}(x^{(n)}) \end{bmatrix}$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$$\text{then } J(\theta) = \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

$$\nabla_{\theta} J(\theta) = X^T X \theta - X^T y = 0 \Rightarrow \theta = (X^T X)^{-1} X^T y$$

Optimal value.