

Lecture 2

Saturday, April 18, 2020 8:33 PM

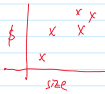
April 18th

Linear Regression and Gradient Descent

Example: housing dataset

given Size and Price

Size	Price (thous)
2104	460
1416	252
1634	315
807	128
...	...

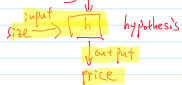


Goal: fit a straight

How?

Training Set

Learning Algorithm



How to represent h ?

$$h(x) = \theta_0 + \theta_1 x$$

You can have more features.

Size	# of rooms	Price
...

More Notations:

m = # of training examples
(# of rows from dataset)

X = "Inputs" / features

y = "Output" / target variables

(x, y) = training example

$(x^{(i)}, y^{(i)})$ = i -th training example

Example: using above datasets

$$x^{(1)} = 2104$$

$$x^{(1)} = 1416$$

n = # of features

more generally, x_i is size

x_i is # bedrooms

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

$$h(x) = \sum_{j=0}^n \theta_j x_j$$

where $x_0 = 1$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

θ = parameters

the goal of the learning algorithms is to choose parameter θ that allows you to make good predictions about price of houses.

How to choose θ ?

choose θ such that $h(x)$ is close to y for training example

notation for representing h depends on input x and parameter θ

$$h_{\theta}(x) = h(x)$$

We want to minimize the equation $(h_{\theta}(x) - y)^2$

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

$$J(\theta)$$

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

One implementation:

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

assignment notation

learning rate

partial derivative

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

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

$$\frac{\partial J(\theta)}{\partial \theta_j} = 2 \left[\frac{1}{2} (h_{\theta}(x) - y) \right] \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y)$$

$$= (h_{\theta}(x) - y) \cdot \left[\frac{\partial}{\partial \theta_j} [\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n - y] \right]$$

if

all the terms become zero except the j th term, which leaves us with $\theta_j x_j$ and the derivative gives us x_j

$$= (h_{\theta}(x) - y) \cdot x_j$$

Use gradient descent to minimize

start with some initial value

θ (say $\theta = \vec{0}$)

Keep changing θ to reduce $J(\theta)$

```

Loop {
  for i=1 to m, {
     $\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$  (for every j)
  }
}

```

Stochastic gradient descent

looking at random training example at each iteration
 and look for improvement in the next iteration
 would be helpful if m , # of training is large.
 don't have to scan entire dataset (compared to batch gradient descent)

The normal equations

Notations =

$\nabla_\theta J(\theta)$ derivative of $J(\theta)$
 with respect to θ
 $\theta \in \mathbb{R}^{n+1}$

$$= \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{bmatrix}$$

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

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

$f(A)$
function

$$f: \mathbb{R}^{2 \times 2} \rightarrow \mathbb{R}$$

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

$$\text{if } f(A) = A_{11} + A_{12}^2$$

$$\text{then } f\left(\begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}\right) = 5 + 6^2$$

$$\longrightarrow \nabla_A f(A) = \begin{bmatrix} 1 & 2A_{12} \\ 0 & 0 \end{bmatrix}$$

$$\nabla_\theta J(\theta) \stackrel{\text{set}}{=} 0$$

More Notation: if A is a $n \times n$ square matrix
 then the trace of A is defined
 to be its diagonal entries:

$$\text{tr } A = \sum_{i=1}^n A_{ii}$$

some properties of $\text{tr } A$:

$$\text{tr } A = \text{tr } A^T$$

$$\text{if } f(A) = \text{tr } AB$$

$$\text{tr } (A+B) = \text{tr } A + \text{tr } B$$

$$\text{then } \nabla_A f(A) = B^T$$

$$\text{tr } aA = a \text{tr } A$$

$$\text{tr } ABC = \text{tr } ACB$$

$$\nabla_A \text{tr } AA^T C = C + C^T A$$

Express

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \quad \text{in Matrix vector notation}$$

$$X\theta = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} = \begin{bmatrix} x^{(1)T} \theta \\ x^{(2)T} \theta \\ \vdots \\ x^{(m)T} \theta \end{bmatrix} = \begin{bmatrix} h_\theta(x^{(1)}) \\ \vdots \\ h_\theta(x^{(m)}) \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\dots \dots \dots [h_\theta(x^{(i)}) - y^{(i)}] \quad Z^T Z = \sum z^2$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$h_{\theta}(x^{(i)})$$

$$X\theta - y = \begin{bmatrix} h_{\theta}(x^{(1)}) - y^{(1)} \\ \vdots \\ h_{\theta}(x^{(m)}) - y^{(m)} \end{bmatrix} \quad Z^T Z = \sum_i z^2$$

$$\text{thus: } J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

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

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

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

$$= \frac{1}{2} \nabla_{\theta} [\theta^T X^T X\theta - \theta^T X^T y - y^T X\theta + y^T y]$$

$$= \frac{1}{2} [X^T X\theta + X^T X\theta - X^T y - X^T y]$$

$$= X^T X\theta - X^T y \stackrel{\text{set } 0}{=}$$

$$X^T X\theta = X^T y \quad \text{"Normal equation"}$$

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

$$\nabla_A^T f(A) = (\nabla_A f(A))^T$$

$$\nabla_{A^T}^T ABA^T C = B^T A^T C^T + BA^T C$$

$$\nabla_{A^T}^T ABA^T C = CAB + C^T AB^T$$

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

$$= \frac{1}{2} \nabla_{\theta} (X\theta^T X\theta - X\theta^T \vec{y} - \vec{y}^T X\theta + \vec{y}^T \vec{y}) \rightarrow \frac{1}{2} \nabla_{\theta} (\theta^T X^T X\theta - \theta^T X^T \vec{y} - \vec{y}^T X\theta + \vec{y}^T \vec{y})$$