

Lecture 2 [Linear Regression and Gradient Descent]

The process of supervised learning is that you have a [Training Set] that you feed to the learning algorithm whose job is to output a function to make predictions (a hypothesis)

When designing a learning algorithm, the first question you must ask is "how do you represent the hypothesis, H,

Linear Regression:

$$H(x) = O_0 + O_1x_1 + O_2x_2 + \dots$$

The hypothesis is going to have input of size x and outputs a number as a linear function of the size x

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

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 $\{n+1 \text{ dimensions}\}$

m = # of training examples
(# of rows in table above)

x = "inputs" / features $\{n+1 \text{ dimensions}\}$

y = "output" / target variable

(x, y) = training example

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

n = # of features

How do you choose the parameters?

The learning algorithm's job is to choose values for a parameter's data (θ) so that it can output a hypothesis

In the Linear Regression algorithm, also called ordinary least squares, you'll want to minimize the squared difference between what the hypothesis outputs

Choose θ s.t. $h(x) \approx y$ for training examples

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

Linear Regression: minimize θ
cost function

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



Batch / Stochastic Gradient Descent:

You can implement an algorithm that find the value of θ that minimizes the cost function $J(\theta)$

You're going to use an algorithm called Gradient Descent

Gradient Descent

Start with some θ (say $\theta = \overset{\text{vector of 0's}}{\begin{pmatrix} 0 \\ 0 \end{pmatrix}}$)
Keep changing θ to reduce $J(\theta)$



Gradient descent will take the steps that lead to the lowest possible value of $J(\theta)$ (local optima)

When you run Gradient Descent on Linear Regression, it turns out that there will not be local optimum

Each step in Gradient Descent is implemented as follows:

Gradient Descent

Start with some θ (say $\theta = \overset{\text{vector of 0's}}{\begin{pmatrix} 0 \\ 0 \end{pmatrix}}$)
Keep changing θ to reduce $J(\theta)$



$:=$ denotes assignment

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

d is the learning rate

One step of Gradient Descent:

(Repeat until convergence)

$$\theta_j := \theta_j - d \underbrace{\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}}_{\frac{\partial}{\partial \theta_j} J(\theta)}$$

(for $j=0, 1, \dots, n$)

Derivative

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2$$

$$= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y)$$

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

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

Partial derivative

Where n is the # of features

When you plot the cost function of $J(\theta)$ for a Linear Regression model, $J(\theta)$ will be a quadratic function (3D graphically looks like a bowl)

If you think about the contours of the function, it turns out that the direction of steepest descent is always at 90 degrees, it's always orthogonal to the contour direction

This algorithm on the cost function of $J(\theta)$ has only one global minimum and will always converge to the global minimum

Is an algorithm for fitting linear regression models

The Gradient Descent algorithm above calculates the derivative by summing over your entire training set, m , and sometimes this version of Gradient Descent has another name which is **Batch Gradient Descent**

The disadvantage of Batch Gradient Descent is that if you have a giant dataset, in order to make one update to your parameters, in order to take even a single step of Gradient Descent, you need to scan through your entire dataset and calculate the sum

There's an alternative to Batch Gradient Descent:

Algorithm: Stochastic Gradient Descent

Repeat {

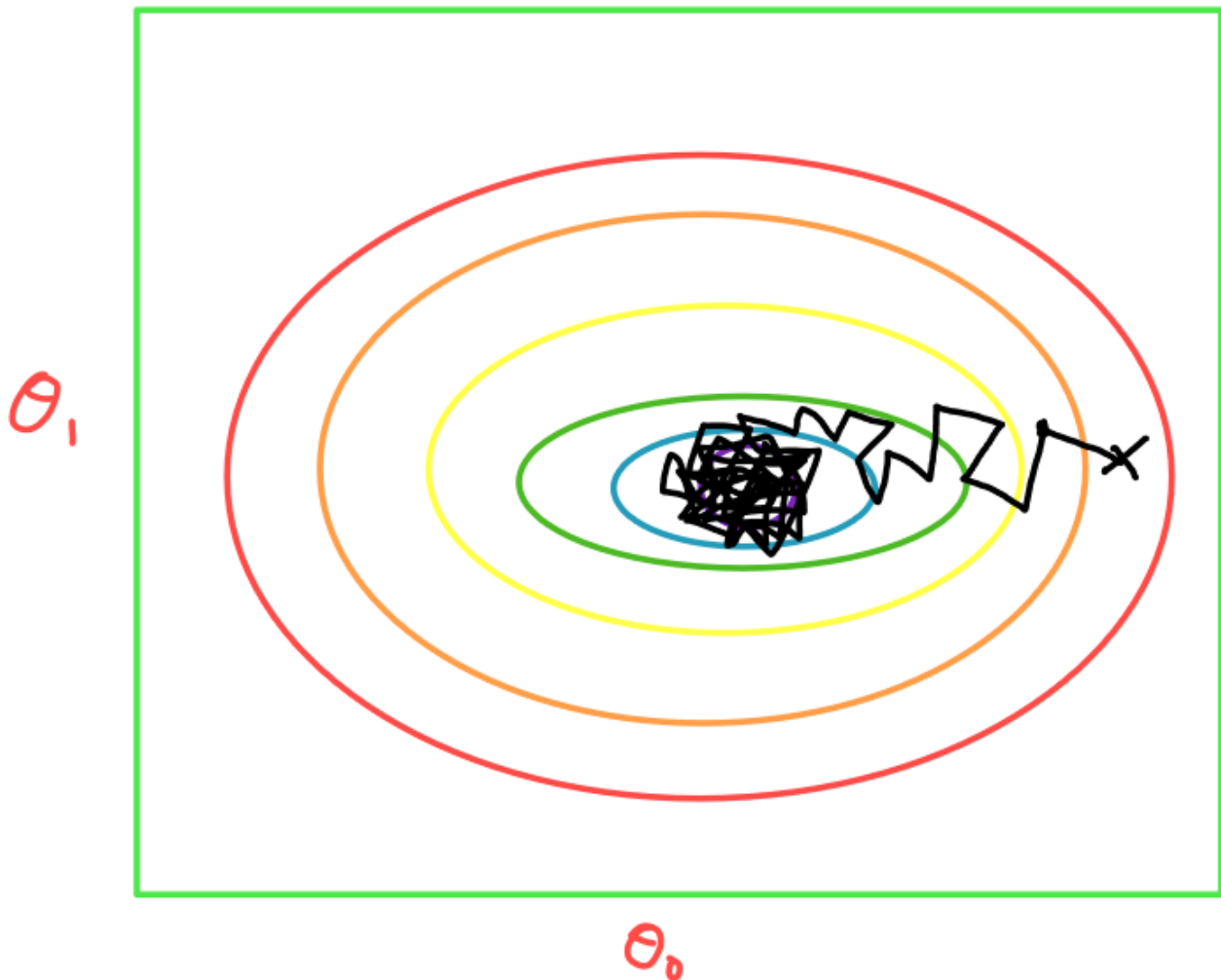
for $j=1$ to m {

$$\theta_j := \theta_j - d(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

}

}

Instead of scanning through all examples before you update the parameters θ even a little bit, in **Stochastic Gradient Descent**, instead in the inner loop of the algorithm, you loop through $j = 1$ to m of taking a gradient descent step using the derivative of just one



As you run Stochastic Gradient Descent, it takes a slightly noisy, slightly random path, but on average is headed toward the global minimum. It will never quite converge

When you have a very large dataset, it allows your implementation/algorithm to make much faster progress

When you have very large datasets, stochastic gradient descent is used much more in practice than batch gradient descent

Gradient Descent is an iterative algorithm

With Linear Regression specifically, there's a way to solve for the optimal value of the parameter's theta to just jump in one step to the global optimum without needing to use an iterative algorithm

Normal Equation:

It works only for Linear Regression

Matrix Derivation definition:

$\nabla_{\theta} J(\theta)$
 Derivative of $J(\theta)$
 w.r.t. θ

$\theta \in \mathbb{R}^{n+1}$
 $n+1$ dimensional vector

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{bmatrix}$$

3 component vector

$$A \in \mathbb{R}^{2 \times 2}, A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

$$f(A) = A_{11} + A_{12}^2, f: \mathbb{R}^{2 \times 2} \mapsto \mathbb{R}$$

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

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

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

Definition of a derivative of a matrix

• If A was a 2×2 matrix, then the derivative of $\nabla_A f(A)$ is also a 2×2 matrix

Derivation of the Normal Equation

Cost function: $J(\theta)$

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

minimize function

• then solve for θ

If A is square matrix ($A \in \mathbb{R}^{n \times n}$)

$$\text{tr } A = \text{sum of diagonal entries} = \sum_i A_{ii}$$

"tr(A)" "trace of A"

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

useful properties

• If, $f(A) = \text{tr } AB$, B is a fixed matrix

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

$$\text{tr } AB = \text{tr } BA$$

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

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

$$\approx \frac{d}{da} a^2 c = 2ac$$

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

$$\begin{aligned} \text{Design Matrix } X \theta &= \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(m)})^T \end{bmatrix} \theta = \begin{bmatrix} x^{(1)T} \theta \\ x^{(2)T} \theta \\ \vdots \\ x^{(m)T} \theta \end{bmatrix} \rightarrow y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \\ &= \begin{bmatrix} h_{\theta}(x^{(1)}) \\ \vdots \\ h_{\theta}(x^{(m)}) \end{bmatrix} \leftarrow x\theta \end{aligned}$$

column vector

vector of all predictions of the algorithm

$$x\theta - y = \begin{bmatrix} h_{\theta}(x^{(1)}) - y^{(1)} \\ \vdots \\ h_{\theta}(x^{(m)}) - y^{(m)} \end{bmatrix}$$

contains all the errors (difference between predictions and actual labels) that your learning algorithm makes under m examples

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

$$= \frac{1}{2} \nabla_{\theta} [\theta^T X^T x\theta - \theta^T X^T y - y^T x\theta + y^T y] \quad \cdot ((ax-b)(ax-b) = a^2 x^2 - axb - bax + b^2)$$

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

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

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

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

What if X is non-invertable?

That usually means that you have redundant features, that your features are linearly dependent, but if you use something called a pseudo-inverse you can kind of get the right answer. If you have linearly dependent features, it probably means you have the same feature repeated twice

