

Learning algorithms, linear vs. quadratic

Supervised learning = "right answers" given

regression = predict continuous valued output (price)

classification = discrete valued output (0 or 1), could be $1 \rightarrow N$ (multi-class)

infinite # of features, other algorithms

Unsupervised learning, not told what to do, given to extract find structure

clustering algorithm = find clusters, used to organize large computer clusters

Cocktail party problem = find structure and separate sounds (audio processing), 1 line of code

$[W, s, v] = \text{Svd}((\text{repmat}(\text{sum}(x.^* 1), \text{size}(x, 1), 1).^* x)^* x');$

Octave (open-source) or MatLab

Arthur Samuel (1959): Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed

Tom Mitchell (1998): Well-posed Learning Problem: A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E

Machine Learning algorithms:

- Supervised Learning

- Unsupervised Learning

Others: Reinforcement Learning, Evolutionary Systems

Model and Cost Function / Linear Regression w/ 1 Variable

Regression = predict real-valued output / classification = discrete-valued output

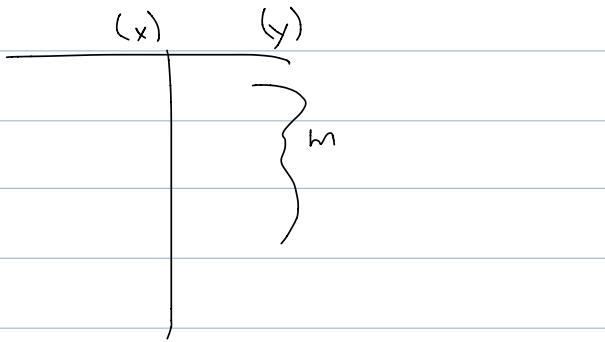
m = Number of training examples

$x^{(i)}$ s = "input" variable / features

$y^{(i)}$ s = "output" variable / "target" variable

(x, y) - One training example

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



training set

How do we represent h ?



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

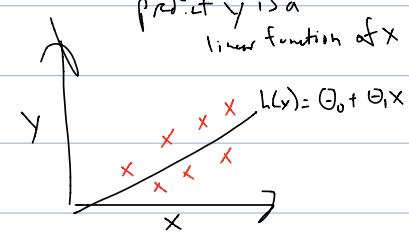
learning algorithm

Shorthand: $h(x)$



size of house \rightarrow h \rightarrow Estimated price
(x) \quad hypothesis ($\text{estimated value of } y$)

h maps from x 's to y 's



predict y is a linear function of x

Linear regression with one variable (x)

Univariate linear regression
one variable

$$\text{Hypothesis: } h_{\theta}(x) = \theta_0 + \theta_1 x$$

θ_i 's: Parameters

How to choose?

$$(x^{(i)}, y^{(i)})$$

Idea: Choose θ_0, θ_1 so that $h_{\theta}(x)$ is close to y for our training examples (x, y)

minimize $\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$ $m = \# \text{ of training examples}$

\uparrow $\arg \min_{\theta_0, \theta_1} h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$

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

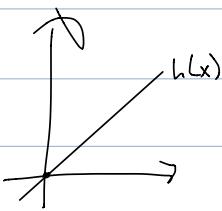
$$\underset{\theta_0, \theta_1}{\text{minimize}} \underset{\text{cost function}}{J(\theta_0, \theta_1)} = \text{goal}$$

squared error function

Simpl: fixed $h_{\theta}(x) = \theta_0 + \theta_1 x$

$$\theta_0 = 0$$

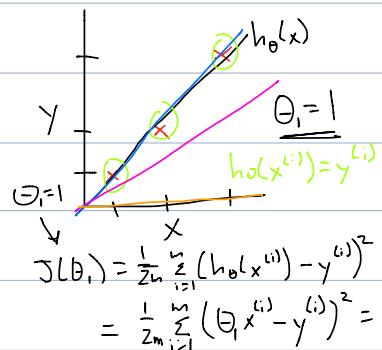
$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



minimize $J(\theta_1)$
 θ_1

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

(for fixed θ_0 , this is a function of x)

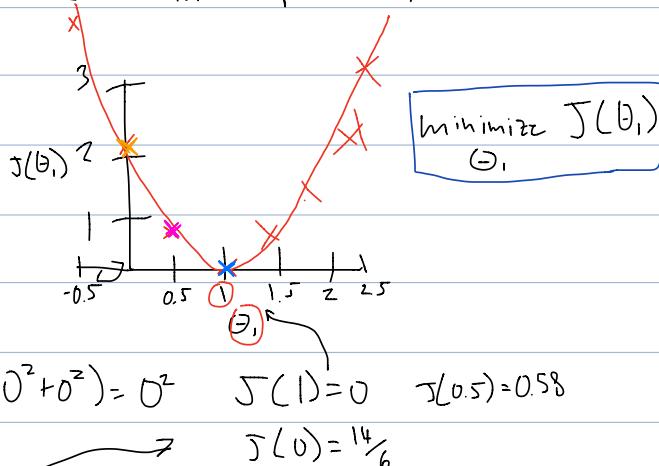


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} (\theta_1^2 + 0^2 + 0^2) = \theta_1^2$$

$$\frac{1}{2m} [(1)^2 + (1)^2 + (3)^2] = \frac{14}{236}$$

$$J(\theta)$$

(function of the parameter θ_1)



$$J(1) = 0 \quad J(0.5) = 0.58 \quad J(0) = \frac{14}{6}$$

1 parameter $\theta_1 = \text{bow}$

2 parameters $\theta_0, \theta_1 = \text{bow (contour)}$

Contour plots \approx Contour figures

Parameter Learning / Gradients

Gradient descent (linear regression, and other stuff)

Have some function $J(\theta_0, \theta_1, \dots, \theta_n)$

Want $\min J(\theta_0, \theta_1) \quad \min J(\theta_0, \dots, \theta_n)$

Θ_0, Θ_1 $\Theta_0, \dots, \Theta_n$ \dots

Outline

- Start with some Θ_0, Θ_1 ($\rightarrow \Theta_0 = 0, \Theta_1 = 0$)
- Keep changing Θ_0, Θ_1 to reduce $J(\Theta_0, \Theta_1)$ until hopefully end up at minimum (local)

Gradient descent algorithm

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 (step size) derivative term

Assignment

 $a := b$ $a := a + 1$

Truth Assertion

 $a = b \leftarrow \text{same}$

Correct: Simultaneous update natural

$$\text{temp}0 := \Theta_0 - \alpha \frac{\partial}{\partial \Theta_0} J(\Theta_0, \Theta_1)$$

$$\text{temp}1 := \Theta_1 - \alpha \frac{\partial}{\partial \Theta_1} J(\Theta_0, \Theta_1)$$

$$\Theta_0 := \text{temp}0$$

$$\Theta_1 := \text{temp}1$$

Incorrect (Not gradient descent, Flat algo)

$$\text{temp}0 := \Theta_0 - \alpha \frac{\partial}{\partial \Theta_0} J(\Theta_0, \Theta_1)$$

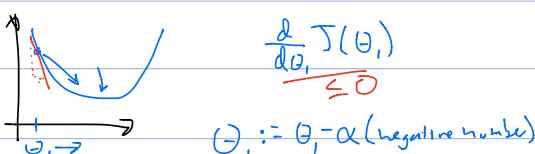
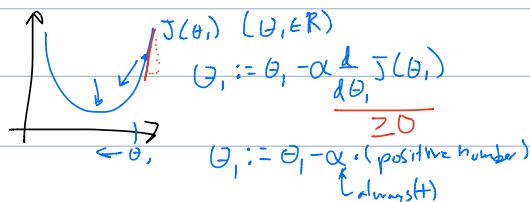
$$\Theta_0 := \text{temp}0$$

$$\text{temp}1 := \Theta_1 - \alpha \frac{\partial}{\partial \Theta_1} J(\Theta_0, \Theta_1)$$

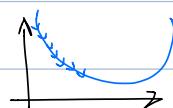
$$\Theta_1 := \text{temp}1$$

$$\min_{\Theta_1} J(\Theta_1)$$

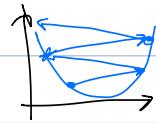
$$\Theta_1 \in \mathbb{R}$$



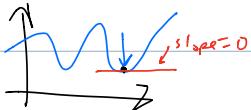
If α is too small, gradient descent can be slow



If α is too large, gradient can overshoot minimum. It may fail to converge, or even diverge.



If at local optimum, derivative is 0, so don't move



Gradient descent can converge to local minimum, even with the learning rate α fixed

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time



Apply gradient descent to Linear Regression Model

Need partial derivatives term

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

$$= \frac{\partial}{\partial \theta_0} \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

Repeat until convergence { $\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

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

update

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

θ_0 and θ_1
simultaneously

}

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

Cost function is always a "convex" function (bowl shape)

always get global optimum

usually initialize θ_0, θ_1

"Batch" Gradient Descent

"Batch": Each step of gradient descent uses all the training examples

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