Hypothesis: $h_\theta(x) = \Theta^T x = \Theta_0 x_0 + \Theta_1 x_1 + \Theta_2 x_2 + \dots + \Theta_n x_n$

Parameters: $\Theta_0, \Theta_1, \dots, \Theta_n \cong \Theta$   $n+1$ dimensional vector

Cost function

$$J(\Theta_0, \Theta_1, \dots, \Theta_n) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J(\theta)$$

Gradient descent:

Repeat {

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

}

$n \geq 1$

Repeat {
$\frac{\partial}{\partial \Theta_j} J(\theta)$

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

}

Feature Scaling

Make sure features on a similar scale

$x_1 = \frac{\text{size (feet}^2)}{2000}$

$x_2 = \frac{\# \text{ bedrooms}}{5}$

Get every feature into approximately $-1 \leq x_i \leq 1$ range

$0 \leq x_1 \leq 3$ ✓

$-2 \leq x_2 \leq 0.5$ ✓          $-3$ to $3$ ✓

$-100 \leq x_3 \leq 100$ ✗        $-\frac{1}{3}$ to $\frac{1}{3}$ ✓

$-.0001 \leq x_4 \leq 0.0001$ ✗

mean normalization

replace $X_i$ with $X_i - \mu_i$ to make features have approximately zero mean

E.g. $X_1 = \frac{size - 1000}{2000}$

$X_1 \leftarrow \dfrac{X_1 - \boxed{\mu_1}}{\boxed{S_1}}$ ← avg value of $x_1$ in training set

range(max - min)
(or standard deviation)
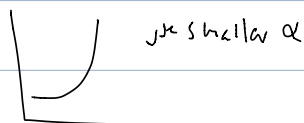
$\dfrac{X_1 - \mu_1}{S_1}$  $\mu_1 = 81$  $S_1 = 25$

$\Theta_j := \Theta_j - \alpha \frac{\partial}{\partial \Theta_j} J(\Theta)$

making sure GD working correctly

automatic convergence test

declare convergence if $J(\Theta)$ decreases by less than $\boxed{10^{-3}}$ in one iteration



use smaller $\alpha$

For sufficiently small $\alpha$, $J(\Theta)$ should decrease on every iteration

But if $\alpha$ is too small, gradient descent trust slowly

If $\alpha$ is too small: slow convergence
If $\alpha$ is too large: $J(\Theta)$ may not decrease on every iteration; may not converge

Normal equation: Method to solve for $\Theta$ analytically

Intuition: If $1D (\Theta \in \mathbb{R})$

$J(\Theta) = a\Theta^2 + b\Theta + c$

$\frac{\partial}{\partial \Theta} J(\Theta) = \dots \overset{set}{=} 0$

Solve for $\theta$

$\theta \in \mathbb{R}^{n+1}$   $J(\theta_0, \theta_1, ..., \theta_m) = \frac{1}{2n} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

$\frac{\partial}{\partial \theta_j} J(\theta) \overset{set}{=} 0$ (for every $j$)

Solve for $\theta_0, \theta_1, ..., \theta_n$

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

Octave: $pinv(X' * X) * X' * y$

m training examples, n features

| Gradient Descent | Normal Equation |
|---|---|
| · Need to choose $\alpha$ | · No need to choose $\alpha$ |
| · Needs many iterations | · Don't need to iterate |
| · Works well even | · Need to compute $(X^T X)^{-1}$ $n \times n$   $O(n^3)$ |
| when n is large | · Slow if n is very large |
| $n = 10^6$ | $n = 100 \checkmark$ |
| | $n = 1000 \checkmark$ |
| | $\longleftarrow n = 10000 \sim$ |