# Machine Learning algorithms

## Machine learning:

"The field of study that gives computers the ability to learn without being explicitly programmed." by Arthur Samuel

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E"

## Linear regression

- **Linear regression**: predict a real-value output based on input value. For example, house price prediction

- **Multivariate linear regression**: predict a real-value output based on multiple input values.

- **Cost function**: mean square error (MSE)

**Gradient Descent**

Previously (n=1):

Repeat {

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

$$\boxed{\frac{\partial}{\partial \theta_0} J(\theta)}$$

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

(simultaneously update $\theta_0, \theta_1$)

}

New algorithm $(n \geq 1)$:

Repeat {

$$\sqrt{} \frac{\partial}{\partial \theta_j} J(\theta)$$

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

(simultaneously update $\theta_j$ for $j = 0, \ldots, n$)

$x_0^{(i)} = 1$

}

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_2^{(i)}$$

...

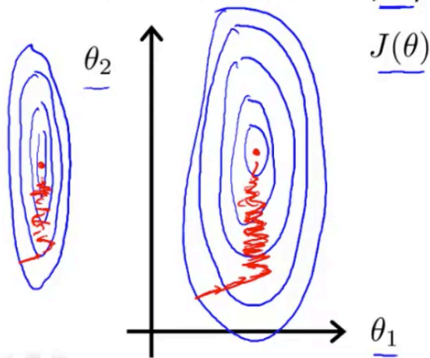- **Practical tricks for Gradient descent**:

  • **Feature scaling**: make sure features are on a similar scale.

    - Gradient descent has a hard time to find its way and take a much <span style="color:red">more convoluted</span> path to the minimum if features are not in the same scale.

    - On the contrary, if the features are in the same scale, Gradient descent takes a <span style="color:red">much simpler and more direct</span> path to the minimum

**Feature Scaling**
Idea: Make sure features are on a similar scale.
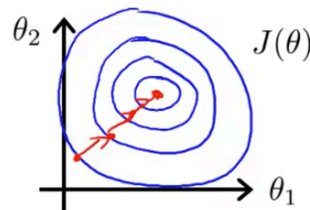
E.g. $x_1$ = size (0-2000 feet$^2$) ←
$x_2$ = number of bedrooms (1-5) ←

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

$x_2 = \dfrac{\text{number of bedrooms}}{5}$

$0 \leq x_1 \leq 1 \qquad 0 \leq x_2 \leq 1$

$\theta_2$ $J(\theta)$ $\theta_1$

$\theta_2$ $J(\theta)$ $\theta_1$

Andrew Ng

    - Get every feature into approximately a [-1 1] range, [-3 3], [-1/3 1/3] range is still fine, no worries if the features are not in the exact same range

    - Mean normalization: to make features with zero mean

**Mean normalization**

Replace $x_i$ with $x_i - \mu_i$ to make features have approximately zero mean (Do not apply to $x_0 = 1$).

E.g. $x_1 = \dfrac{size - 1000}{2000}$

Averyg size = 100

$x_2 = \dfrac{\#bedrooms - 2}{5}$

1-5 bedrooms

$-0.5 \leq x_1 \leq 0.5$ $-0.5 \leq x_2 \leq 0.5$

$x_1 \leftarrow \dfrac{x_1 - \mu_1}{s_1}$

avg value of $x_1$ in training set

range (max - min) (or standard deviation)

$x_2 \leftarrow \dfrac{x_2 - \mu_1}{s_2}$

Andrew Ng

2

- How to make sure **gradient descent is working correctly**?

  - Cost function J($\theta$) should decrease after every iteration (plot J($\theta$) vs no. of iterations)

    - If J($\theta$) is not decreasing, probably learning rate $\alpha$ is too large, should decrease the learning rate $\alpha$

    - If $\alpha$ is small enough, J($\theta$) should decrease on every iteration

    - But if $\alpha$ is too small, then it will take a long time to converge

    - Try a number of $\alpha$, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1…, plot J($\theta$) vs no. of iterations

  - Automatic convergence test: declare convergence if J($\theta$) decreases by less than a certain threshold after every iteration (the threshold is hard to choose most of the time)

- **Polynomial regression**: We can **change the behavior or curve** of our hypothesis function by making it a quadratic, cubic or square root function (or any other form).

  - **Feature engineering**:

    - create new features based on the problem

    - Feature scaling is important with respect to new features (new features might be in a much larger or smaller magnitude)

- **Normal equation method** is another alternative for gradient descent

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

| Gradient Descent | Normal Equation |
|---|---|
| Need to choose alpha | No need to choose alpha |
| Needs many iterations | No need to iterate |
| O $(kn^2)$ | O $(n^3)$, need to calculate inverse of $X^T X$ |
| Works well when n is large | Slow if n is very large |

- Noninvertible (singular)

  • Causes

    - redundant features (linearly dependent)

    - Too many features (No. Training samples less than No. of features )

  • Pseudo inverse will solve this problem

  • Inverse

- **Coding steps**

Simple linear regression

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

parameters/weights: $\theta_0, \theta_1$

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

Find the parameters that'll minimize cost function $J(\theta_0, \theta_1)$:

  ✷ Gradient descent algorithm
   △ starting with random $\theta_0, \theta_1$
   △ Updates the weights until reach minimum of $J(\theta_0, \theta_1)$

   $\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$     $j=0$, for $\theta_0$
   $j=1$, for $\theta_1$

   Note: weights should be updated simultaneously.
   e.g.  $\theta_{0, next} = \theta_{0, current} - \alpha \frac{\partial}{\partial \theta_0} J(\theta_{0, current}, \theta_{1, current})$

      $\theta_{1, next} = \theta_{1, current} - \alpha \frac{\partial}{\partial \theta_1} J(\theta_{0, current}, \theta_{1, current})$

  ✷ Normal equation method

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