

→ Learning Rate (α)

too small: gradient descent is too slow

too big: "big steps", may lead to overshooting and failure to converge

⇒ Gradient Descent for Linear Regression

model $\begin{cases} h_\theta(x) = \theta_0 + \theta_1 x \\ J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \end{cases}$ → apply GD to get lowest cost param

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \rightarrow j=0 \rightarrow \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$
$$\rightarrow j=1 \rightarrow \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

$$\begin{cases} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)} \end{cases}$$

* "Batch" Gradient Descent:
each iteration goes through
the entire training set

Week 2

⇒ Multivariate Linear Regression

X_n feature contribute to output value Y

Notation: $x^{(i)}$ is vector containing all features in i -th example

Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

$$\hookrightarrow \theta_0 = 1 \rightarrow h_\theta(x) = \sum_{i=0}^n x_i \theta_i = \theta^T x$$

vector with all θ_i vector with all x_i

Gradient Descent:

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

→ Feature Scaling

Features being on a similar scale can make gradient descent work faster

get every feature in range $[-1, 1]$ → divide by maximum value

Mean Normalization: $x_i = \frac{x_i - \text{average value}}{\text{maximum range}}$

→ Convergence method

establish a threshold ϵ and if $J(\theta)$ decreases by less than ϵ in one iteration it has converged * hard to determine ϵ

plot no. iterations $\times J(\theta)$ → helpful when choosing α

* It's possible to create new features that enable the fitting of polynomial functions to the data

Normal Equations

- Solves for optimum value directly
 - Take derivative and equal to zero
 - Defining X as a matrix containing all features and y as the solution vector
- $$\theta = (X^T X)^{-1} X^T y$$
- No need to choose α or to iterate
 - Does not scale well to a large number of features