

Gradient Descent

Data Science Immersive

Motivation

- OLS regression (one variable) $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$
- How do we actually find these beta values?
 - Linear algebra:

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{bmatrix} \quad A = \begin{bmatrix} b \\ m \end{bmatrix} \quad E = \begin{bmatrix} e_1 \\ e_2 \\ \dots \\ e_n \end{bmatrix}$$

This gives us the matrix equation: $Y = XA + E$.

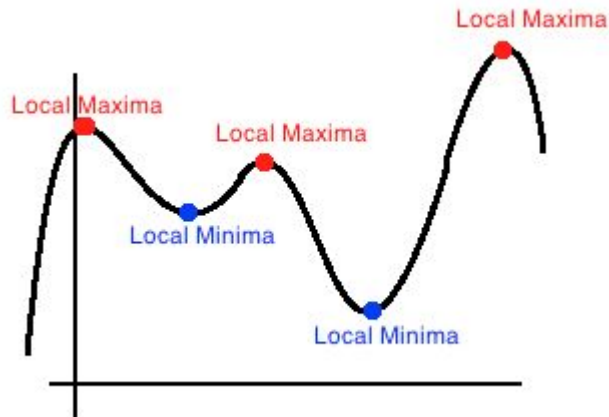
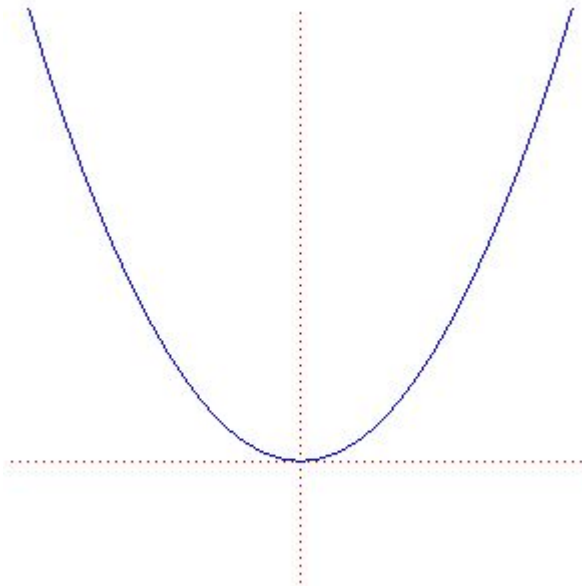
Actually, no.

- As it turns out, solving in this way although very convenient, a direct solution becomes more and more computationally difficult as data sets get bigger.
- “Minimize cost function”

$$RSS = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

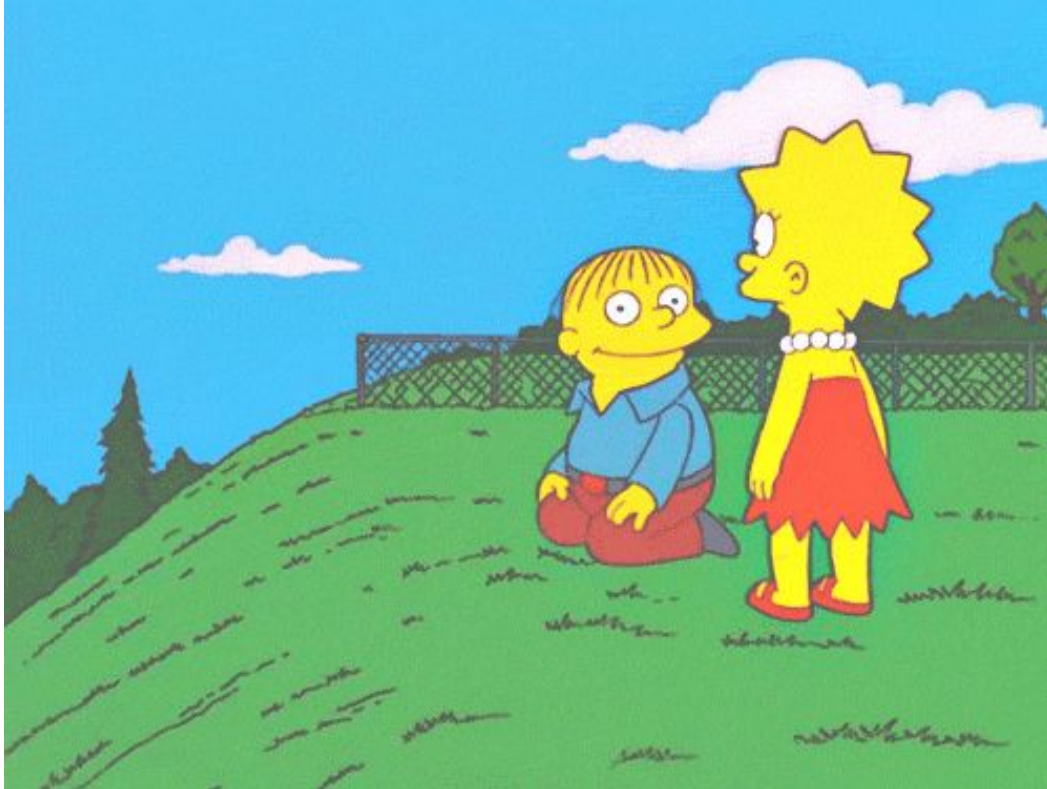
Minimizing a Function

Parabola

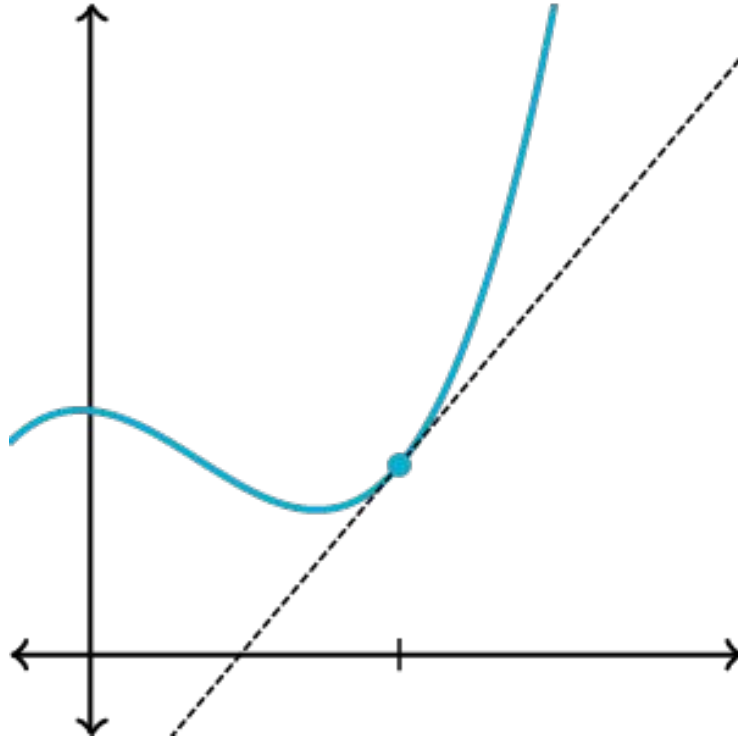


- To find the maxima and minima of a function, find where the derivative equals zero. Very easy for a convex function like the one on the left.
- But, what if you can't use this method?

Minimizing a Function



Derivatives



Derivatives

$$\lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}$$

$$\lim_{h \rightarrow 0} \frac{f(a + h) - f(a)}{h}$$

Derivatives

- Chain Rule

$$\frac{d}{dx}f(g(x)) = f'(g(x))g'(x)$$

- Power rule

$$\frac{d}{dx}x^n = nx^{n-1}$$

- Example:

$$\frac{d}{dx}x^2$$

Gradient Descent formula

Gradient descent algorithm

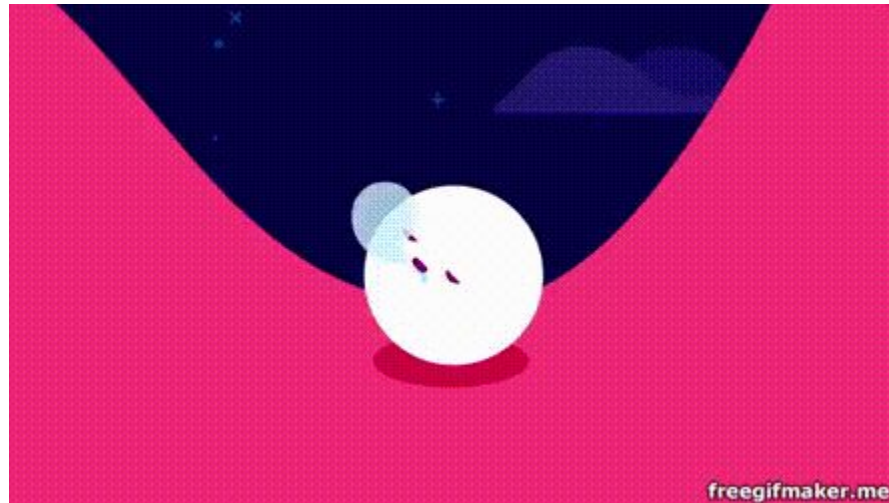
repeat until convergence {
 $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ (for $j = 0$ and $j = 1$)
}

- The weird equal sign
- Thetas
- Partial derivative
- Alpha
- J

Learning rate, α

Potential problems:

- Overshooting the minimum so GD diverges instead of converging
- Getting stuck at a local minimum
- Not approaching minimum fast enough



Gradient Descent for Linear Regression

Gradient descent algorithm

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

Cost

$$J(\theta) = 1/2m \sum_{i=1}^m (h(\theta)^{(i)} - y^{(i)})^2$$

Simple Linear Regression with GD

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

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) * x_i$$

Code Example

