

# Gradient Ascent

---

## Maximum Likelihood Refresher

Our first algorithm for estimating parameters is called Maximum Likelihood Estimation (MLE). The central idea behind MLE is to select that parameters ( $\theta$ ) that make the observed data the most likely.

The data that we are going to use to estimate the parameters are going to be  $n$  independent and identically distributed (IID) samples:  $X_1, X_2, \dots, X_n$ .

## Likelihood

First we define the likelihood of our data give parameters  $\theta$ :

$$L(\theta) = \prod_{i=1}^n f(X_i|\theta)$$

This is the probability of all of our data. It evaluates to a product because all  $X_i$  are independent. Now we chose the value of  $\theta$  that maximizes the likelihood function. Formally  $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta)$ .

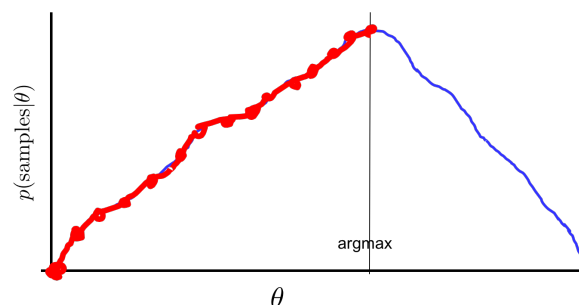
A cool property of  $\operatorname{argmax}$  is that since  $\log$  is a monotone function, the  $\operatorname{argmax}$  of a function is the same as the  $\operatorname{argmax}$  of the  $\log$  of the function! That's nice because logs make the math simpler. Instead of using likelihood, you should instead use log likelihood:  $LL(\theta)$ .

$$LL(\theta) = \log \prod_{i=1}^n f(X_i|\theta) = \sum_{i=1}^n \log f(X_i|\theta)$$

To use a maximum likelihood estimator, first write the log likelihood of the data given your parameters. Then chose the value of parameters that maximize the log likelihood function.  $\operatorname{Argmax}$  can be computed in many ways. Most require computing the first derivative of the function.

## Gradient Ascent Optimization

In many cases we can't solve for  $\operatorname{argmax}$  mathematically. Instead we use a computer. To do so we employ an algorithm called gradient ascent (a classic in optimization theory). The idea behind gradient ascent is that if you continuously take small steps in the direction of your gradient, you will eventually make it to a local maxima. In the case of logistic regression you can prove that the result will always be a global maxima.



Start with theta as any initial value (often 0). Then take many small steps towards a local maxima. The new theta after each small step can be calculated as:

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \frac{\partial LL(\theta^{\text{old}})}{\partial \theta_j^{\text{old}}}$$

Where  $\eta$  is the magnitude of the step size that we take. If you keep updating  $\theta$  using the equation above you will (often) converge on good values of  $\theta$ . Here is the gradient ascent algorithm in pseudo-code:

**Initialize:**  $\theta_j = 0$  for all  $0 \leq j \leq m$

**Repeat many times:**

**gradient[j] = 0** for all  $0 \leq j \leq m$

**For each training example (x, y):**

**Update gradient[j] for current training example**

$\theta_j \mathrel{+}= \eta * \text{gradient}[j]$  for all  $0 \leq j \leq m$

## Linear Transform Plus Noise

MLE is an algorithm that can be used for any probability model with a derivable likelihood function. As an example lets estimate the parameter  $\theta$  in a model where there is a random variable  $Y$  such that  $Y = \theta X + Z$ ,  $Z \sim N(0, \sigma^2)$  and  $X$  is an unknown distribution.

In the case where you are told the value of  $X$ ,  $\theta X$  is a number and  $\theta X + Z$  is the sum of a gaussian and a number. This implies that  $Y|X \sim N(\theta X, \sigma^2)$ . Our goal is to chose a value of  $\theta$  that maximizes the probability IID:  $(X_1, Y_1), (X_2, Y_2), \dots (X_n, Y_n)$ .

We approach this problem by first finding a function for the log likelihood of the data given  $\theta$ . Then we find the value of  $\theta$  that maximizes the log likelihood function. To start, use the PDF of a Normal to express the probability of  $Y|X, \theta$ :

$$f(Y_i|X_i, \theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(Y_i - \theta X_i)^2}{2\sigma^2}}$$

Now we are ready to write the likelihood function, then take its log to get the log likelihood function:

$$L(\theta) = \prod_{i=1}^n f(Y_i, X_i | \theta)$$

Let's break up this joint

$$= \prod_{i=1}^n f(Y_i|X_i, \theta) f(X_i)$$

$f(X_i)$  is independent of  $\theta$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(Y_i - \theta X_i)^2}{2\sigma^2}} f(X_i)$$

Substitute in the definition of  $f(Y_i|X_i)$

$$\begin{aligned}
LL(\theta) &= \log L(\theta) \\
&= \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(Y_i - \theta X_i)^2}{2\sigma^2}} f(X_i) && \text{Substitute in } L(\theta) \\
&= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(Y_i - \theta X_i)^2}{2\sigma^2}} + \sum_{i=1}^n \log f(X_i) && \text{Log of a product is the sum of logs} \\
&= n \log \frac{1}{\sqrt{2\pi}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \theta X_i)^2 + \sum_{i=1}^n \log f(X_i)
\end{aligned}$$

Remove constant multipliers and terms that don't include  $\theta$ . We are left with trying to find a value of  $\theta$  that maximizes:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} - \sum_{i=1}^n (Y_i - \theta X_i)^2$$

To run this argmax we first need to find the derivative of the function with respect to  $\theta$ .

$$\begin{aligned}
\frac{\partial}{\partial \theta} - \sum_{i=1}^n (Y_i - \theta X_i)^2 &= - \sum_{i=1}^n \frac{\partial}{\partial \theta} (Y_i - \theta X_i)^2 \\
&= - \sum_{i=1}^n 2(Y_i - \theta X_i)(-X_i)
\end{aligned}$$

This first derivative can be plugged into gradient ascent to give our final algorithm:

**Initialize:**  $\theta = 0$

**Repeat many times:**

**gradient** = 0

**For each training example (x, y):**

**gradient** += 2(y -  $\theta$  x) (-x)

$\theta$  +=  $\eta$  \* **gradient**