

Gradient Descent

- Simple optimization algorithm that can be used with any machine learning algorithm
- Gradient descent is simply used to find the values of a function's parameters (coefficients) that minimize a cost function as much as possible

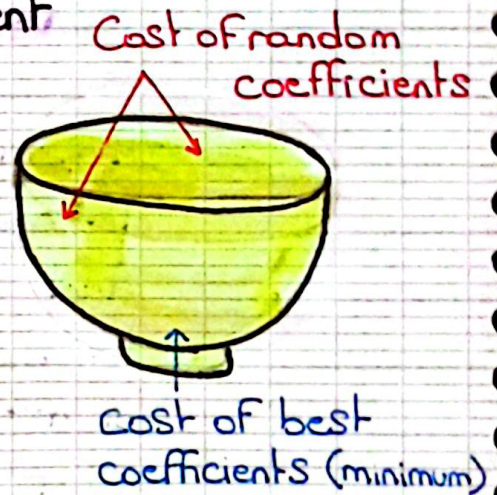
1) Intuition of Gradient Descent

• It's like a large bowl where we store some fruit.

• This bowl represents the plot of the cost function (f).

• A random position on the surface of the bowl is the cost of the current values of the coefficients.

• The bottom of the bowl is the cost of the best set of coefficients, that give the minimum of the cost function



• This is an iterative method, so we should move from a random point, to reach the minimum point.

• As we move we're varying the coefficients of the line until we reach the best line (or when we reach the minimum error)

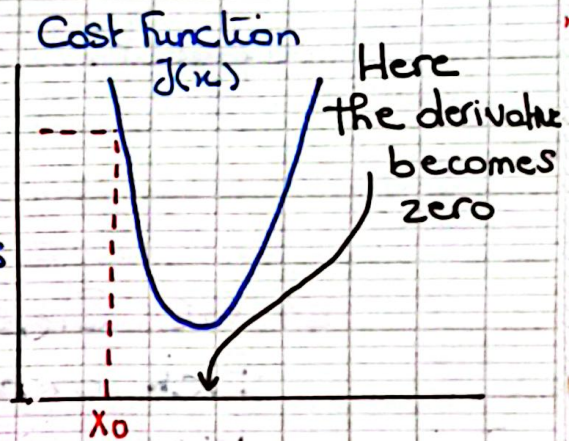
2) How does Gradient Descent Work?

Iterative process, where we start at a coefficient's initial point (x_0) and we move step by step until we reach a minimum.

The update rule of our position is given by the formula:

$$x_{t+1} = x_t - \alpha \frac{dJ}{dx} \Big|_{x_t}$$

Annotations:
- α : Learning rate
- x_t : parameter at previous step t
- x_{t+1} : parameter at step $t+1$
- $\frac{dJ}{dx} \Big|_{x_t}$: derivative of the cost function with $x = x_t$

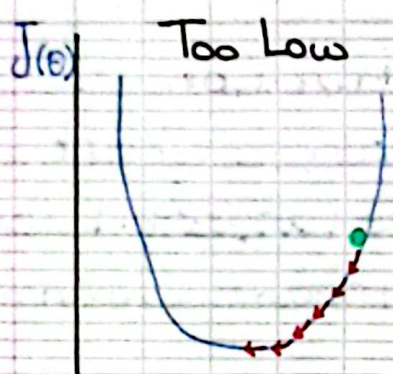


We stop when the values of x don't change a lot: $x_{t+1} \approx x_t$

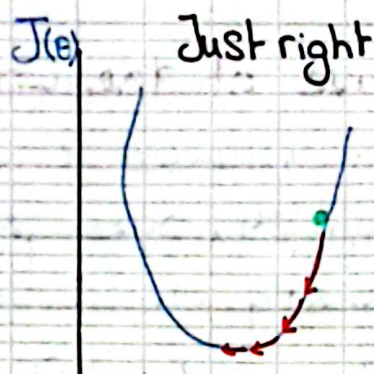
Since it's an iterative process, we can stop when we reach the maximum nb of iterations

3) The importance of the Learning Rate (α)

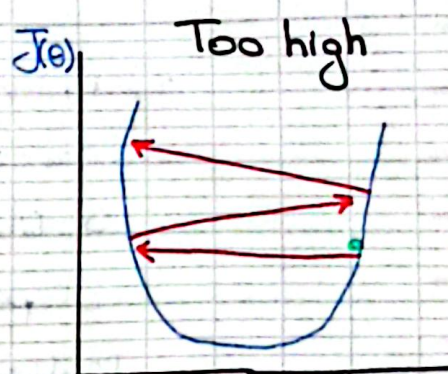
How big of steps a gradient descent algorithm takes into the direction of the local minimum is determined by the learning rate which determines the rate (how fast or slow) we will move towards the optimal coefficients/parameters



A small learning rate requires many updates before reaching the minimum point (very cost)



The optimal learning rate swiftly reaches the minimum point.



Too large of a learning rate causes drastic updates which lead to divergent behaviors. (risk overshooting min)

4) GD Steps

- Pick an initial value for learning rate α
- Pick an initial value for the parameter x_0
- Evaluate the derivative over x_t
- Update the parameter x_{t+1} using the GD update formula
- Repeat until we reach convergence (i.e. minimum cost)

$$x_{t+1} = x_t - \alpha \frac{dJ}{dx} \Big|_{x_t}$$

Derivative of the Cost Function

The cost function in a linear regression problem is mainly the mean squared error (MSE). In the MSE, we compute the derivatives as follows in general form.

$$\text{Cost} = \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{pred}})^2 = \frac{1}{n} \sum_{i=1}^n (y_i - (a_0 + a_1 x_i))^2$$

Our task is to find a_0 and a_1 that minimize the cost, but how we do this?

- 1) We should take the partial derivatives of the cost with respect to a_0 and a_1
- 2) Set the partial derivatives equal to 0
- 3) Solve for a_0 and a_1

1) Partial Derivatives

Let's start taking the partial derivative with respect to a_0

$$\begin{aligned} \frac{\partial \text{cost}}{\partial a_0} &= \frac{\partial}{\partial a_0} \sum_{i=1}^n (y_i - (a_0 + a_1 x_i))^2 \\ &= \sum_{i=1}^n \frac{\partial}{\partial a_0} (y_i - (a_0 + a_1 x_i))^2 \end{aligned}$$

derivative of sum is the sum of derivatives
 $u^n = n u^{n-1} \cdot u'$

as we take the derivative with respect to a_0 we treat y_i, a_1, x_i as constants

$$= \sum_{i=1}^n 2 (y_i - (a_0 + a_1 x_i)) \cdot \frac{\partial}{\partial a_0} (y_i - (a_0 + a_1 x_i))$$

$\frac{\partial}{\partial a_0} (y_i - (a_0 + a_1 x_i)) = -1$

$$= \sum_{i=1}^n 2 (y_i - (a_0 + a_1 x_i)) \cdot (-1)$$

$$= - \frac{2}{n} \sum_{i=1}^n (y_i - (a_0 + a_1 x_i))$$

Now with respect to a_1

$$\begin{aligned} \frac{d \text{cost}}{d a_1} &= \frac{d}{d a_1} \frac{1}{n} \sum_{i=1}^n (y - (a_0 + a_1 x))^2 \\ &= \frac{1}{n} \sum_{i=1}^n \frac{d}{d a_1} (y - (a_0 + a_1 x))^2 \\ &= \frac{1}{n} \sum_{i=1}^n 2 (y - (a_0 + a_1 x)) \cdot \frac{d}{d a_1} (y - (a_0 + a_1 x)) \\ &= \frac{1}{n} \sum_{i=1}^n 2 (y - (a_0 + a_1 x)) \cdot (-x) \\ &= -\frac{2}{n} \sum_{i=1}^n x (y - (a_0 + a_1 x)) \end{aligned}$$

The derivatives of sum is the sum of the derivatives
 $u^n = n u^{n-1} \cdot u'$
 $\frac{d}{d a_1} (y - (a_0 + a_1 x)) = -x$

as we take the derivative in respect to a_1 , we treat y, a_0 as constants.

2) Set partial derivatives equal to 0

a. $-\frac{2}{n} \sum_{i=1}^n (y - (a_0 + a_1 x)) = 0$ solve this for a_0

b. $-\frac{2}{n} \sum_{i=1}^n x (y - (a_0 + a_1 x)) = 0$
 \uparrow Substitute that here for a_1

a) $\sum (y - (a_0 + a_1 x)) = 0$
 $[\sum y - \sum a_0 - \sum a_1 x] = 0$ with respect the summation a_0 and a_1 are constants so they can come outside
 $[\sum y - n a_0 - a_1 \sum x] = 0$
 $n a_0 = \sum y - a_1 \sum x$
 $a_0 = \bar{y} - a_1 \bar{x}$

b) $\sum x (y - (a_0 + a_1 x)) = 0$
 $\sum x (y - (\bar{y} - a_1 \bar{x} + a_1 x)) = 0$
 $\sum x (y - \bar{y} - a_1 (x - \bar{x})) = 0$
 $\sum x (y - \bar{y}) - \sum a_1 x (x - \bar{x}) = 0$
 $\sum x (y - \bar{y}) = a_1 \sum x (x - \bar{x}) \Rightarrow a_1 = \frac{\sum x (y - \bar{y})}{\sum x (x - \bar{x})}$

↳ Gradient Descent Example

Given a set of data points representing the height and the corresponding weights of students we need to find a linear relation through Gradient Descent

Height	Weight
160	85
175	88
180	92
190	97

First, we plot the data points, and pass a random line $y = a_0 + a_1 x$ to initialize the parameters (a_0 : bias, a_1 : slope)

Assume we choose

a_0 -initial = 0.2 and
 a_1 -initial = 0.2, we can compute the predicted weight as $0.2 + 0.2 \times \text{Height}$

Height	Weight	W-pred
160	85	32.2
175	88	35.2
180	92	36.2
190	97	38.2

(we can see that the predicted values are not close at all to our real weights: So we need to work on that)

Now we can calculate the squared error (MSE) for each example and then compute the function

Height	Weight	W-pred	error	Error Squared
160	85	32.2	52.8	...
175	88	35.2	52.8	...
180	92	36.2	55.8	...
190	97	38.2	58.8	...

$$\text{Cost} = \frac{1}{4} * \text{Sum}(\text{Error Squared})$$

$$\frac{\partial \text{cost}}{\partial a_0} = -\frac{2}{4} \sum (\text{error}) = -110.1$$

$$\frac{\partial \text{cost}}{\partial a_1} = -\frac{2}{4} \sum \text{Height} (\text{error}) = -1945.2$$

} derivatives of the cost function

• Now we perform the update equations using a learning rate of 0.01

$$x_{t+1} = x_t - \alpha \frac{dT}{dx} |_{x_t}$$

$$a_0 = a_{0 \text{ initial}} - 0.01 \times \frac{\partial \text{cost}}{\partial a_0} = 0.2 - 0.01 \times -110.1 = 1.301$$

$$a_1 = a_{1 \text{ initial}} - 0.01 \times \frac{\partial \text{cost}}{\partial a_1} = 0.2 + 0.01 \times 19452 = 194.72$$

This is only 1 iteration of GD

→ Gradient Descent Tips

1) Take into consideration the cost versus the time (by plotting cost vs time). Analyzing the plot, we expect that a well performing rate in descent is going to decrease during each run/iteration. If it doesn't decrease we should try reducing the learning rate.

2) We should take a small value for the learning rate (0.1, 0.01, 0.001 or 0.0001)

3) The algorithm will reach the minimum cost faster if the shape of the cost function is not skewed and distorted.

We can achieve this by rescaling all of the input variables (x) to the same range (0-1) which is by performing Normalization

→ Gradient Descent for logistic regression

The same concept of GD will be applied here, but the derivatives will change.

$$C(y_i, y'_i) = -\frac{1}{n} \sum_{i=1}^n y_i \log(\sigma(a_0 + a_1 x_i)) + (1 - y_i) \log(1 - \sigma(a_0 + a_1 x_i))$$

We need to calculate the derivatives of the cost function with respect to the different parameters involved in the prediction process

$$\frac{\partial C}{\partial a_0} = -\frac{1}{n} \sum_{i=1}^n [y_i - \sigma(a_0 + a_1 x_i)]$$

$$\frac{\partial C}{\partial a_1} = -\frac{1}{n} \sum_{i=1}^n [y_i - \sigma(a_0 + a_1 x_i)] x_i$$

• If the prediction involves several features, we need to calculate the derivatives with respect to the different parameters that we have:

• Here's an example of 2 features x_1 and x_2

$$\frac{\partial C}{\partial a_0} = -\frac{1}{n} \sum [\gamma_i - \sigma(a_0 + a_1 x_{1i} + a_2 x_{2i})]$$

$$\frac{\partial C}{\partial a_1} = -\frac{1}{n} \sum [\gamma_i - \sigma(a_0 + a_1 x_{1i} + a_2 x_{2i})] x_{1i}$$

$$\frac{\partial C}{\partial a_2} = -\frac{1}{n} \sum [\gamma_i - \sigma(a_0 + a_1 x_{1i} + a_2 x_{2i})] x_{2i}$$

Note that: x_1 and x_2 refer to the features,
and "i" refers to the index of the data example

• Afterwards, we just perform the updates as usual using the famous GD update formula

$$a_i := a_i - \text{learning-rate} \times \frac{\partial \text{Cost}}{\partial a_i}$$

• We keep performing the updates until we reach the stopping criteria that we set (convergence, maximum number of iterations...)