

★ Recap with new notations

In gradient descent $x_2 = x_1 - \eta \frac{d}{dx} f(x)$

The above formula is only for a univariate function.

For a bivariate function it will be:

$$x_2 = x_1 - \eta \left(\frac{\partial f(x, y)}{\partial x} + \frac{\partial f(x, y)}{\partial y} \right)$$

New notation: $x_1 \rightarrow \theta^1$ & $x_2 \rightarrow \theta^2$

$$\theta^{t+1} = \theta^t - \eta \sum \frac{\partial f(x, y)}{\partial \theta}$$

$$\theta_2 = \theta_1 - \eta \cdot \frac{\partial f(x, y)}{\partial x}$$

$$\theta_3 = \theta_2 - \eta \frac{\partial f(x, y)}{\partial y}$$

$$\theta^{t+1} = \theta^t - \eta \frac{\partial}{\partial \theta} f(x_k)$$

where k is a random number
between 1 to n

A Constrained Optimization Problem

Loss fn:

$$\sum \frac{\vec{\omega}^T \cdot \vec{x} + \omega_0 \cdot y_i}{\|\vec{\omega}\|}$$

If $\|\vec{\omega}\| = 1$ (\leftarrow constraint) then it will make it very easy to calculate differentiation of loss function.

$$\underset{\vec{\omega}, \omega_0}{\text{argmin}} \sum \frac{\vec{\omega}^T \cdot \vec{x} + \omega_0 \cdot y_i}{\|\vec{\omega}\|} \quad \underline{\text{S.T.}} \quad \|\vec{\omega}\| = 1$$

Now $\|\vec{\omega}\| = 1$

$$\frac{\partial L}{\partial \omega_1} = - \frac{\partial}{\partial \omega_1}$$

$$= -x_1 y_1$$

$$\frac{\partial L}{\partial \omega_2} = -x_2 y_2$$

$$\frac{\partial L}{\partial \omega} = -\sum x_i y_i$$

$$\frac{(\omega_1 x_1 y_1 + \cancel{\omega_2 x_2 y_2} + \cancel{\omega_n x_n} + \cancel{\omega_0})}{1}$$

$$\omega_j^{t+1} = \omega_j^t - \eta (-\sum x_i y_i)$$

For $j \neq 0$

For $j=0$ (or for ω_0)

$$\omega_0^{t+1} = \omega_0^t - \eta (-\sum y_i)$$

★ Suppose our function that we want to minimize is $f(x) = x^2 - 3x - 3$ & we put a constraint that our minima must also satisfy $g(x) = x^2 - 2x - 3$

$\therefore \underset{x}{\text{argmin}} f(x) \text{ s.t. } g(x)$; s.t. = 'such that' or 'subject to' \textcircled{I}

Goal: To do something so that $f(x)$ is also minimized and constraint is also taken care of.

Solution: Lagrange's multipliers

Instead of computing \textcircled{I} , we will do this:

Optimize
 x, λ

$$f(x) + \lambda g(x)$$

Lagrange's Multiplier

Hey! Minimize
the $f(x)$!

Keep also $g(x)$
in mind...

Why did we convert 'such that ...' form into above form?

Because:

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



$$\frac{d}{dx} (f(x) \text{ s.t. } g(x)) = \text{Not possible}$$

∴ Our problem is converted to an
'unconstrained' form.

★ Generalized form of the previous formula for
 n constraints $(g_1(x), g_2(x), \dots, g_n(x))$:

$$\begin{array}{ll} \text{Optimize} & f(x) + \lambda_1 g_1(x) + \lambda_2 g_2(x) + \dots + \lambda_n g_n(x) \\ x, \lambda_1, \lambda_2, \dots, \lambda_n \end{array}$$

∴ Doing this for our example, let 'h' be:

$$\arg \min_{x, \lambda} f(x) + \lambda g(x) = \arg \min_{x, \lambda} x^2 - 3x - 3 + \lambda x^2 - 2\lambda x - 3\lambda$$

To minimize this, we need to find its gradient

$$\nabla h = \begin{bmatrix} \frac{\partial}{\partial x} h \\ \frac{\partial}{\partial \lambda} h \end{bmatrix} = \begin{bmatrix} 2x - 3 + 2\lambda x - 2\lambda \\ x^2 - 2x - 3 \end{bmatrix}$$

At minimum gradient is 0.

$$\therefore \nabla h = 0$$

$$\therefore 2x + 2\lambda x - 2\lambda - 3 = 0 \text{ --- (A)}$$

$$x^2 - 2x - 3 = 0 \Rightarrow x^2 - 3x + x - 3 = 0$$

$$x(x-3) + 1(x-3) = 0$$

$$\Rightarrow \therefore x = 3 \text{ OR } x = -1$$

Putting these values in (A)

$$\text{FOA } x = 3$$

$$6 + 6\lambda - 2\lambda - 3 = 0$$

$$4\lambda + 3 = 0$$

$$\lambda = \frac{-3}{4}$$

$$\underline{2} \quad \text{FOA } x = -1$$

$$-2 - 2\lambda - 2\lambda - 3 = 0$$

$$-4\lambda - 5 = 0$$

$$\therefore \lambda = \frac{-5}{4}$$

The two points are: $\left(3, \frac{-3}{4}\right)$ ~~and~~ $\left(-1, \frac{-5}{4}\right)$

value of $f(x) + \lambda g(x)$: \downarrow -3 \downarrow 1

★ Implementing this technique for our loss function -
Then our $f(x) = - \sum_{i=1}^n \frac{\vec{\omega}^T \cdot \vec{x} + \omega_0}{\|\vec{\omega}\|} \cdot y_i$ is our constraint

$$\text{is } \|\vec{\omega}\| = 1 \Rightarrow g(x) : \|\vec{\omega}\| - 1 = 0$$

$$\underset{\vec{\omega}, \lambda}{\text{argmin}} f(x) + \lambda g(x) = - \sum_{i=1}^n \vec{\omega}^T \cdot \vec{x} + \omega_0 + \lambda (\|\vec{\omega}\| - 1)$$

Let's find gradient

$$\nabla h = \begin{bmatrix} \frac{\partial}{\partial \vec{\omega}} h \\ \frac{\partial}{\partial \lambda} h \end{bmatrix}$$

First component: $\partial h / \partial \vec{\omega}$

$$\|\vec{\omega}\| = \sqrt{\omega_1^2 + \omega_2^2 + \dots + \omega_n^2}$$

$$= \sqrt{[\omega_1, \omega_2, \dots, \omega_n] \begin{bmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_n \end{bmatrix}}$$

$$\|\vec{\omega}\| = \sqrt{\vec{\omega}^T \cdot \vec{\omega}}$$

$$\frac{\partial \sqrt{x}}{\partial x} = \frac{1}{2\sqrt{x}}$$

$$\frac{\partial}{\partial \vec{x}} \vec{x}^T \cdot \vec{x} = 2\vec{x}$$

$$\frac{\partial h}{\partial \omega} = - \sum \frac{\partial}{\partial \vec{\omega}} \vec{\omega}^T \cdot \vec{x} + \omega_0 + \lambda (\|\vec{\omega}\| - 1)$$

$$= - \sum \vec{x} + 0 + \frac{\partial}{\partial \vec{\omega}} \lambda (\|\vec{\omega}\| - 1)$$

$$\frac{\partial}{\partial \vec{\omega}} \lambda (\|\vec{\omega}\| - 1) = \lambda \frac{\partial}{\partial \vec{\omega}} \|\vec{\omega}\| + 0$$

$$= \lambda \frac{\partial}{\partial \vec{\omega}} \sqrt{\vec{\omega}^T \cdot \vec{\omega}}$$

$$\frac{\partial}{\partial \vec{w}} \lambda (\|\vec{w}\| - 1) = \lambda \frac{1}{2\sqrt{\vec{w}^T \cdot \vec{w}}} \cdot \frac{\partial}{\partial \vec{w}} \vec{w}^T \cdot \vec{w}$$

$$= \lambda \frac{2 \vec{w}}{2\sqrt{\vec{w}^T \cdot \vec{w}}}$$

$$= \lambda \frac{\vec{w}}{\|\vec{w}\|}$$

$$\frac{\partial}{\partial \vec{w}} h = -\sum \vec{x} + \lambda \frac{\vec{w}}{\|\vec{w}\|}$$

☆ second component of the gradient is: $\frac{\partial h}{\partial \lambda}$

$$= \frac{\partial}{\partial \lambda} \left(- \sum_{i=1}^n \vec{\omega}^T \cdot \vec{x} + \omega_0 + \lambda (\|\vec{\omega}\| - 1) \right)$$

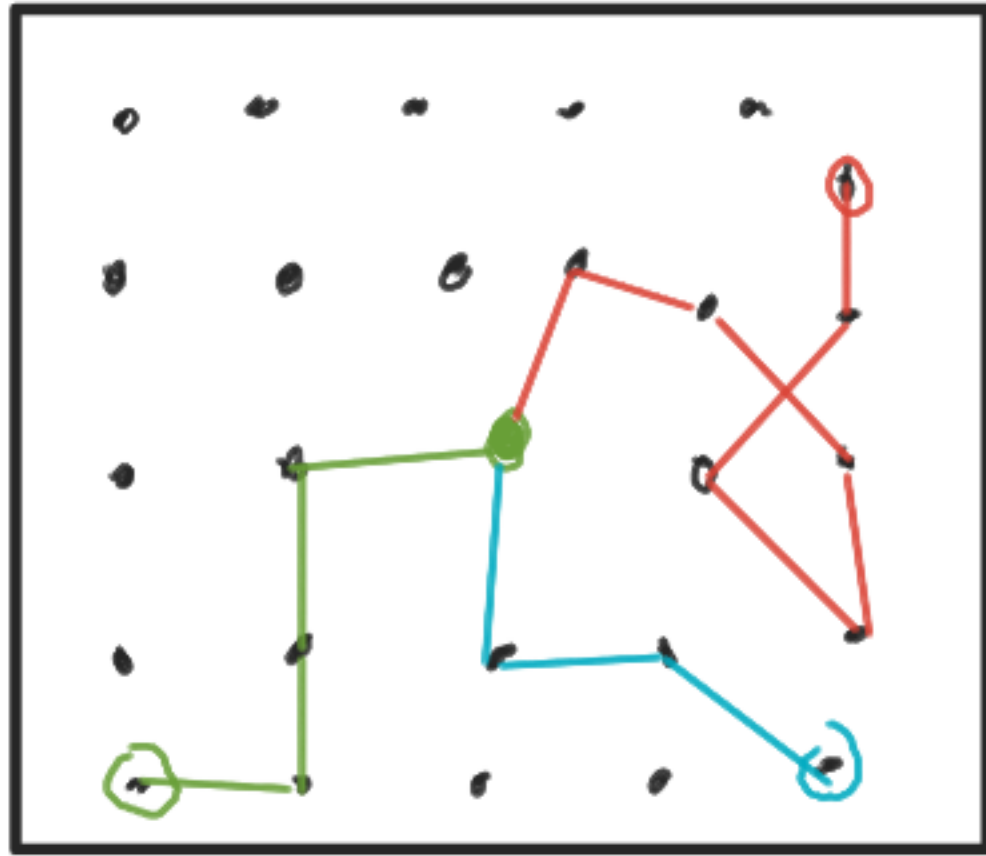
$$= - \sum_{i=1}^n 0 + 0 + \frac{\partial}{\partial \lambda} \lambda (\|\vec{\omega}\| - 1)$$

$$= - \sum_{i=1}^n (\|\vec{\omega}\| - 1)$$

★ Quick Recap:

Gradient Descent: (1) Vanilla G.D. / G.D. / Batch G.D.

For finding each next guess it needs N iterations



(2) Stochastic G.D. - For each next guess, this takes only one iteration, but it will need more guesses to reach to minima.

(3) Minibatch G.D. - For each next guess

minibatch G.D. takes ' k ' iterations where $1 < k < N$.

