

★ Now let's implement Lagrange's multiplier technique to minimize our Loss Function of G.D.

$$f(x) = L = - \sum \frac{\bar{w}^T \cdot \bar{x} + w_0}{\|\bar{w}\|} \cdot y_i \quad \text{s.t.} \quad \|\bar{w}\| = 1$$

An imp. note: Our constraint must be in the form:

$g(x): \underline{\hspace{2cm}} = 0$ \therefore Our constraint here will be:

$$g(x): \|\bar{w}\| - 1 = 0$$

★ \therefore Unconstrained form will be:

$$\arg \min_{\bar{w}, \lambda} f(x) + \lambda g(x) = \arg \min_{\bar{w}, \lambda} - \sum (\bar{w}^T \cdot \bar{x} + w_0) \cdot y_i + \lambda (\|\bar{w}\| - 1)$$

\therefore gradient $\nabla h = \begin{bmatrix} \frac{\partial}{\partial \bar{\omega}} h \\ \frac{\partial}{\partial \lambda} h \end{bmatrix}$

$\xrightarrow{\text{First Component}}$

$\xrightarrow{\text{Second Component}}$

Before calculating these components, let's see some interesting results:

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

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

$$\|\bar{\omega}\| = \sqrt{\bar{\omega}^t \cdot \bar{\omega}} \quad \text{--- (I)}$$

$$\frac{d}{dx} \sqrt{x} = \frac{1}{2\sqrt{x}}$$

--- (II)

$$\frac{\partial}{\partial \bar{x}} \bar{x}^t \cdot \bar{x} = 2\bar{x}$$

--- (III)

$$\underline{\text{First Component: }} \frac{\partial}{\partial \bar{\omega}} h = \frac{\partial}{\partial \bar{\omega}} \left(-\sum (\bar{\omega}^T \cdot \bar{x} + w_0) y_i + \lambda (\|\bar{\omega}\| - 1) \right)$$

$$= \underset{\bar{\omega}, \lambda}{\text{argmin}} \quad -\sum y_i (\bar{x} + 0) + \frac{\partial}{\partial \bar{\omega}} \lambda \|\bar{\omega}\| - 0$$

$$= \underset{\bar{\omega}, \lambda}{\text{argmin}} \quad -\sum y_i \cdot \bar{x} + \lambda \frac{\partial}{\partial \bar{\omega}} \left(\sqrt{\bar{\omega}^T \cdot \bar{\omega}} \right) \quad (\text{From (I)})$$

$$= \underset{\bar{\omega}, \lambda}{\text{argmin}} \quad -\sum y_i \bar{x} + \frac{\lambda}{2\sqrt{\bar{\omega}^T \cdot \bar{\omega}}} \cdot \frac{\partial}{\partial \bar{\omega}} (\bar{\omega}^T \cdot \bar{\omega}) \quad (\text{From (I)})$$

$$= \underset{\bar{\omega}, \lambda}{\text{argmin}} \quad -\sum y_i \bar{x} + \frac{\lambda \bar{\omega}}{\sqrt{\bar{\omega}^T \cdot \bar{\omega}}}$$

$$\frac{\partial}{\partial \bar{\omega}} h = \text{argmin}_{\bar{\omega}, \lambda} - \sum \gamma_i \cdot \bar{x} + \lambda \frac{\bar{\omega}}{\|\bar{\omega}\|}$$

* Second Component = $\frac{\partial}{\partial \lambda} h = \frac{\partial}{\partial \lambda} \lambda \cdot (\|\bar{\omega}\| - 1)$

$$= \frac{\partial}{\partial \lambda} \lambda \|\bar{\omega}\| - \frac{\partial}{\partial \lambda} \lambda = \|\bar{\omega}\| - 1$$

★ 3-variants of G.D. -

① Vanilla G.D. / Batch G.D.: $\omega_i^{t+1} = \omega_i^t - \eta (\nabla f)$

This will take us to the next guess from the initial guess. Suppose we need take 1000 iterations of taking next guess to reach the minima then let's call $n=1000$.

In each iteration to compute ∇f , we need to differentiate f 'd' no. of times where d = dimensions.

In each differentiation, we need to go through each x_i & y_i

So let $m = 10,000$ be our no. of datapoints. Hence to take each new guess Vanilla G.D. will do $m \cdot d$ computations ($10,000 \times 50 = 50,000$; assuming $d = 50$)

As we see it is computationally very expensive.

(2) Stochastic G.D. - This variant of GD doesn't consider all ' m ' datapoints to find the next guess but only 1 random datapoint.

As a result, this will need more iterations to converge.

③ Mini batch G.D.- This considers k datapoints (k can be any number) to compute the next guess.

As per CLT, the mean of k datapoints is close to actual mean hence, this variant of GD will have less randomness than Stochastic & hence will need less iterations than it.