

Constrained Optimization
[Optimisation]

Recap

$$\min_{\vec{w}, w_0} - \sum_{i=1}^n y_i \frac{\vec{w}^\top \vec{x}_i + w_0}{\|\vec{w}\|}$$

Since this is a multivariate function with $d+1$ variables (for ' d ' features), we cannot directly optimise it using calculus, hence

we will use gradient descent

Gradient Descent:

$$\text{iteration } t+1 \leftarrow \theta_i^{t+1} = \theta_i^t - n \frac{\text{learning rate}}{\partial \theta_i} \frac{\partial f}{\partial \theta_i} \rightarrow \begin{array}{l} \text{learning rate} \\ \text{loss function} \\ \partial \theta_i \rightarrow \text{parameter} \end{array}$$

For our loss function, we need

$$\frac{\partial f}{\partial w_1, w_2 \dots}, \quad \frac{\partial f}{\partial \underline{w}_0}$$

\equiv

where,

$$f(w_0, w_1, \dots, w_n) = - \sum_{i=1}^m y_i \frac{w^T x_i + w_0}{\|w\|}$$

Let's try to find derivatives !!.

It's principal way \rightarrow always works !!.

$$\frac{\partial f}{\partial w_1} = \frac{f(w_1 + \Delta, w_2 \dots, w_n) - f(w_1, w_2 \dots, w_n)}{\Delta}$$

$$\Delta \approx 0$$

However, this requires 2 computations of $f()$, which we know is expensive.

So for our problem, let's simplify the derivative before coding it.

2nd: manual differentiation

$$\begin{aligned}\frac{\partial f}{\partial w_i} &\rightarrow - \sum g_i \left(\frac{w_1 \partial l_1 + w_2 x_2 \dots + w_n}{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}} \right) \\ &= - \sum g_i \frac{f(w_i)}{g(w_i)}\end{aligned}$$

We can use quotient rule here, but that will be complicated & compute expensive.

Idea:

Consider equation of a hyperplane

$$\omega_1 x_1 + \omega_2 x_2 + \omega_0 = 0$$

if we divide it by some number

$$d = \sqrt{\omega_1^2 + \omega_2^2} = \|\vec{\omega}\|, \text{ we get}$$

$$\omega'_1 x_1 + \omega'_2 x_2 + \omega'_0 = 0$$

where

$$\omega'_1 = \omega_1 / \|\vec{\omega}\|, \quad \omega'_2 = \omega_2 / \|\vec{\omega}\|$$

Q: is this the same hyperplane?

→ Desmos

Hence our equation becomes

$$\omega^T x + w_0 = 0 \rightarrow \hat{\omega}^T x + \underbrace{\frac{w_0}{\|\omega\|}}_{w_0'} = 0$$

Going back,

$$\min_{\vec{\omega}, w_0} - \sum y_i (\omega^T x + w_0)$$

such that $\|\omega\| = \frac{1}{\sqrt{w_0'}}$
condition

This is called constrained optimisation

Now,

$$\frac{\partial f}{\partial \omega_i} = - \sum y_i \left[\frac{\omega_i x_i + w_0}{\sqrt{w_0'}} \right]$$

$$= - \sum y_i x_{1i}$$

$$\therefore \underset{j \neq 0}{\omega_j^{(t+1)}} = \omega_j^{(t)} - n \left(- \hat{\sum}_{i=1}^n y_i x_{ji} \right)$$

$$\frac{\partial f}{\partial \omega_0} = - \sum y_i (\cancel{\omega_1 x_1} + \cancel{\omega_2 x_2} - \frac{\omega_0}{J})$$

$$\therefore \frac{\partial f}{\partial \omega_0} = - \sum y_i$$

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

Hence we can see that f' is very easy to compute w.r.t all parameters.

General Form of Optimisation Problems

$\min \{ \max \}$

Operation

$\theta \leftarrow$ Parameters

Optimisation Function $\rightarrow f(\theta)$

such that / subject to:

constraint functions $\rightarrow g_1(\theta), g_2(\theta), \dots$

Eg:

$$\min_{\theta_1, \theta_2} f(\theta_1, \theta_2)$$

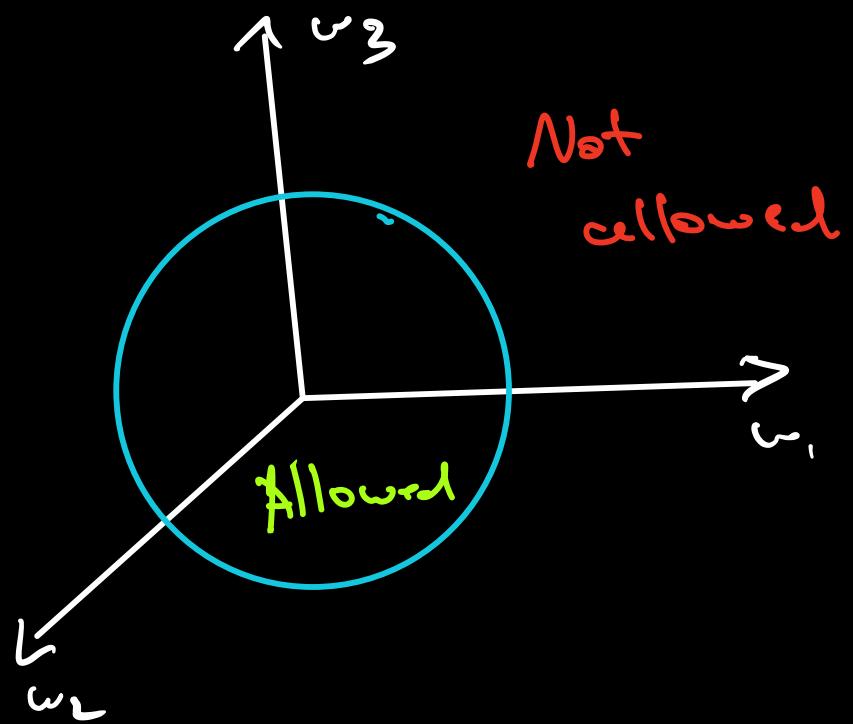
θ_1, θ_2

subject to $g_1(\theta), g_2(\theta_2), g_3(\theta_1, \theta_2), \dots$

We have to solve

$$\min_{\vec{w}, w_0} - \sum y_i (\vec{w}^\top \vec{x}_i + w_0) ; \text{ such that } \|\vec{w}\| - 1 = 0$$

Essentially we are restricting our search space to a hyper-sphere of radius = 1



Q: Can this be solved with simple (vacinal) gradient descent?

→ No, because when we update, we don't take care of the constraint,

$$w_{\text{new}} = w_{\text{old}} - \eta \nabla_w f$$

There is no guarantee that $\|w_{\text{new}}\| = \underline{1}$

We need to convert this to an unconstrained problem again. To do this we use a new technique called

Lagrange Multipliers

Lagrange Multipliers

$$\min_{\boldsymbol{\theta}} \quad f(\boldsymbol{\theta}) ; \text{ such that } \begin{aligned} g_1(\boldsymbol{\theta}) &= 0 \\ g_2(\boldsymbol{\theta}) &= 0 \\ &\vdots \\ g_n(\boldsymbol{\theta}) &= 0 \end{aligned}$$

$$\min_{\boldsymbol{\theta}} \quad f(\boldsymbol{\theta}) + \lambda_1 \downarrow g_1(\boldsymbol{\theta}) + \lambda_2 \downarrow g_2(\boldsymbol{\theta}) \dots$$

Lagrange mult. Lagrange mult.

minimising this unconstrained problem gives the same result as above.

Example

$$\min : \quad \rightarrow \quad x^2 + y^2$$

x, y

$$\text{such that, } \rightarrow x + 2y - 1 = 0$$

→ plot online

$$\text{Lagrangian} \rightarrow L(x, y, \lambda)$$

$$= x^2 + y^2 + \lambda(x + 2y - 1) = 0$$

$$\min_{x, y, \lambda} \rightarrow x^2 + y^2 + \lambda(x + 2y - 1) = 0$$

$$\frac{\partial L}{\partial x} = 2x + 1, \quad \frac{\partial L}{\partial y} = 2y + 2\lambda$$

$$\frac{\partial L}{\partial \lambda} = x + 2y - 1$$

$$\therefore 2x + \lambda = 0 \rightarrow \lambda = -2x$$

$$2y + 2x = 0 \rightarrow \lambda = -y$$

$$x + 2y - 1 = 0 \quad \therefore y = 2x$$

$$\therefore x + 2(2x) - 1 = 0$$

$$\therefore 5x - 1 = 0$$

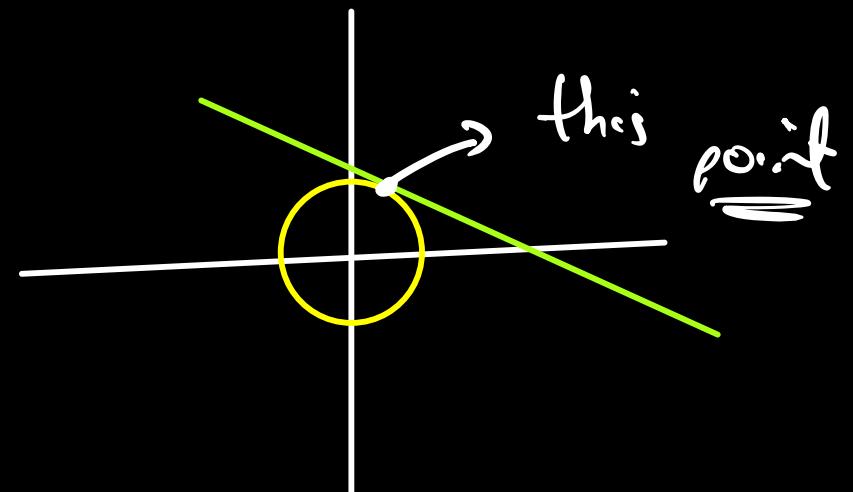
$$\therefore x = 1/5$$

$$\therefore y = 2/5$$

$$\therefore \lambda = -2/5$$

$$\frac{1}{2s} + \frac{2}{2s} = \frac{3}{s}$$

desired result



Intuition (Very high level)

→ Why does this work?

$$\min_{x, y, \lambda} L = f(x, y) + \lambda(g(x, y))$$

if λ is large, and $g(x, y)$ is > 0
 $\lambda(g(x, y))$ would be large,
hence L would not be min

The algorithm is incentivised to make
 $g(x, y) \approx 0$, hence the constrained
is forced to be satisfied.

Exact proofs of why this works out
are beyond our scope.

↳ Note that there are conditions
which need to be satisfied for this
to work, but we won't discuss them
here.

GD with Lagrangian

$$\underset{\vec{\omega}, \omega_0, \lambda}{\min} L(\omega, \omega_0, \lambda)$$

update

$$\omega_i^{+} = \omega - \gamma \frac{\partial L}{\partial \omega_i};$$

$$\lambda^{+} = \lambda + \gamma \frac{\partial L}{\partial \lambda} \quad \begin{array}{l} \text{max lambda to} \\ \text{force } g(\vec{\omega}, \omega_0) \rightarrow 0 \end{array}$$

Assessment Doubt

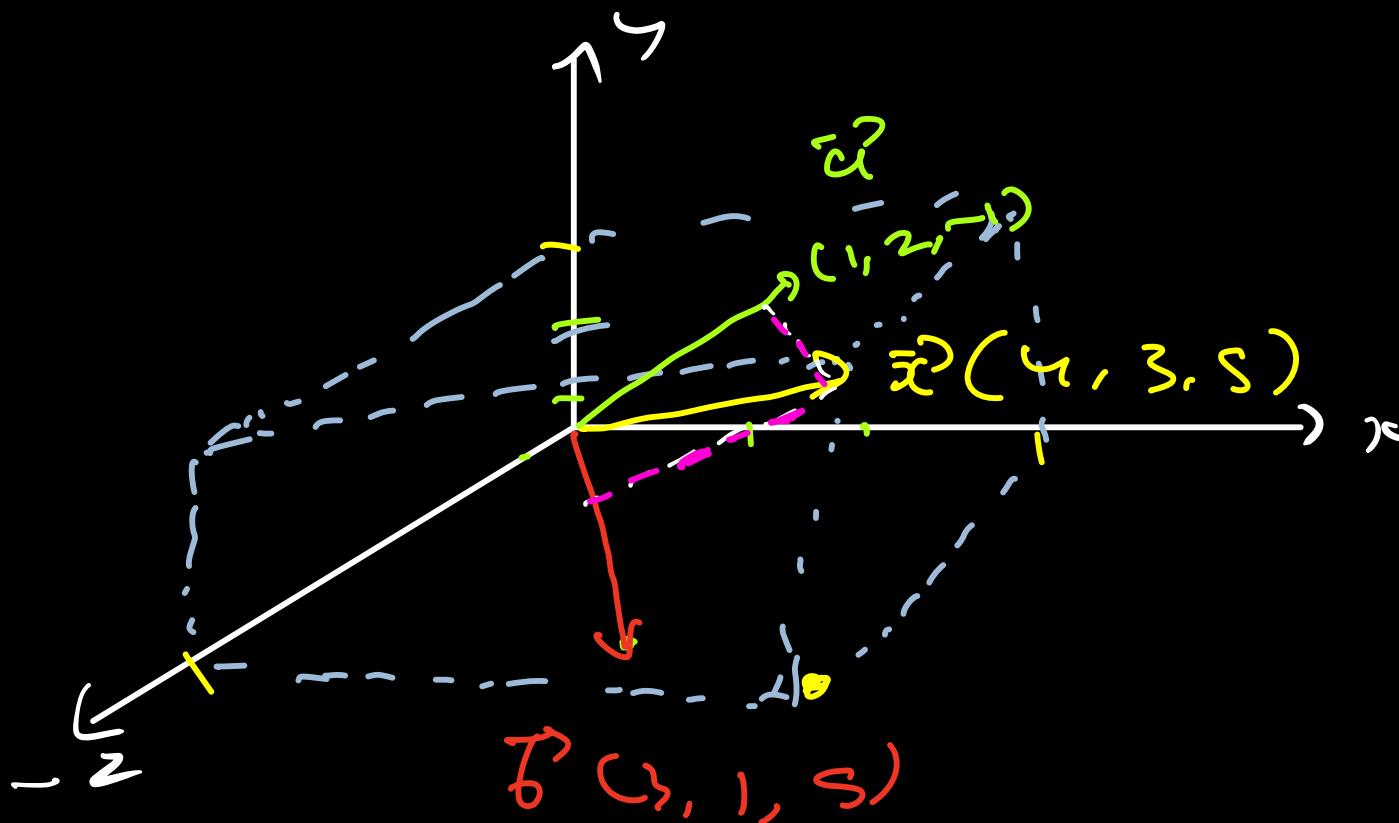
Q.1 $\vec{a} = (4, 3, 5)$
 $\vec{b} = (1, 2, -1)$

$$\hat{\vec{a}} = \left(\frac{1}{\sqrt{6}}, \frac{2}{\sqrt{6}}, \frac{-1}{\sqrt{6}} \right) \quad \hat{\vec{b}} = \left(\frac{1}{\sqrt{35}}, \frac{2}{\sqrt{35}}, \frac{-1}{\sqrt{35}} \right)$$

$$\vec{x} = \frac{\vec{a} \cdot \vec{a}}{\|\vec{a}\|} \cdot \hat{\vec{a}} + \frac{\vec{a} \cdot \vec{b}}{\|\vec{b}\|} \cdot \hat{\vec{b}}$$

$$= \frac{5}{\sqrt{6}} \cdot \hat{\vec{a}} + \frac{40}{\sqrt{35}} \cdot \hat{\vec{b}}$$

also $\vec{a} \cdot \vec{b} = 0 \quad , \quad \therefore \vec{a} \perp \vec{b}$



Further \rightarrow Find eqn of the plane
 containing \vec{a} , \vec{b} and \vec{c} \rightarrow online calc
 $(x - 8y - 5z + 0 = 0)$
 \therefore put \vec{r} in it $\rightarrow ((4) - 8(3) - 5(5) + 0$
 $= 44 - 24 - 25$
 $= -5 \neq 0$ \therefore x is not on plane

This means that we cannot recreate the full point by projection on those 2 axis. Some information is lost, that's okay. This will form the foundation of our next topic $\rightarrow \underline{\text{PCA}}$

Q.2

Let $f(w, w_0) = y (w^T x + w_0)$. For a single datapoint $x = (1, 2, 1)$ and $y = -1$, which of the following are the correct update equations for w and w_0 if we are using gradient descent?

Options :

- a) $w^{(t+1)} = w^t + \eta x, \quad w_0^{(t+1)} = w_0^{(t)} + \eta \cdot 1$
- b) $w^{(t+1)} = w^t - \eta x, \quad w_0^{(t+1)} = w_0^{(t)} - \eta \cdot 1$
- c) $w^{(t+1)} = w^t + \eta x, \quad w_0^{(t+1)} = w_0^{(t)} - \eta \cdot 1$
- d) $w^{(t+1)} = w^t - \eta x, \quad w_0^{(t+1)} = w_0^{(t)} + \eta \cdot 1$

$$\omega^{(t+1)} = \omega^t - n \nabla f(\vec{\omega})$$

$$\omega_0^{(t+1)} = \omega^t - n \nabla f(\omega^*)$$

$$f = -y_i \left(\frac{\omega_1 x_1 + \omega_2 x_2 + \dots + \omega_n}{\|\omega\|} \right)$$

$$\therefore \frac{\partial f}{\partial \omega_i} = -y_i (x_1 + 0 + 0 + \dots + 0)$$

$$\frac{\partial f}{\partial \omega_i} = x_i y_i$$

$$\therefore \frac{\partial f}{\partial \omega_0} = -y_i (0 + 0 + \dots + 1)$$

$$\frac{\partial f}{\partial \omega_0} = -y_i$$

$$\therefore \omega^{(t+1)} = \omega^t - \eta(-(-1)^x) = \boxed{\omega^t - \eta x}$$

$$\omega_o^{(t+1)} = \omega_o^t - \eta(-(-1)) = \boxed{\omega_o^t - \eta}$$