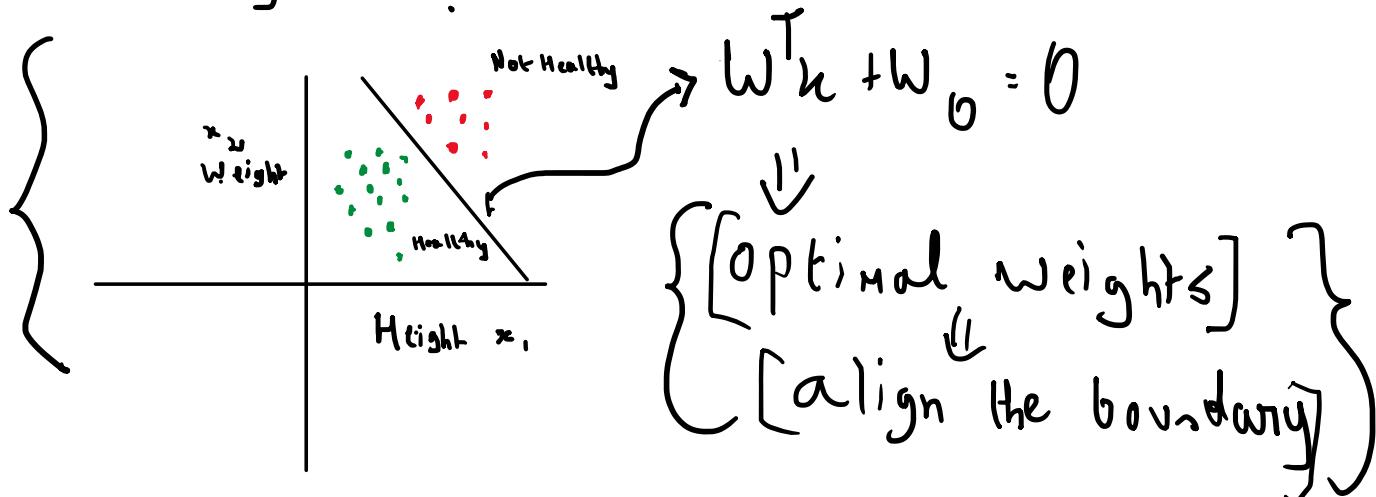
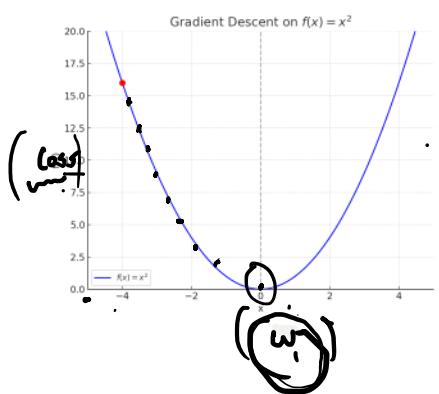


What will minimizing it do?



{ How do we minimize loss fn? } \Rightarrow Gradient Descent



for simplicity,
taking
 $L = w_i^2$

$$\left\{ \hat{w}^{t+1} = \hat{w}^t - \lambda \frac{dL}{dw} \right\}$$

learning
rate

initialise
a
random
weight

Gradient Descent Run

09 September 2025 19:22

$$w_1 \cdot w_2 \cdot w_0 \Rightarrow i) \text{ Initialize weights}$$

$$0.1x_1 - 0.05x_2 + 8.585$$

Student	Blood Sugar (x1 mg/dL)	BMI	y	Predicted \hat{y}
1	120	28	1	1
2	130	30	1	1
3	115	25	1	-1
4	85	20	1	-1
5	75	18	1	1
6	90	22	-1	1
7	95	24	-1	1
8	100	25	-1	-1
9	80	20	-1	1
10	85	21	-1	-1

$$\delta = -\sum \left(\frac{w^T x + w_0}{\|w\|} \right) y_i$$

$$\text{Loss} = -\sum \left(\frac{w^T x + w_0}{\|w\|} \right) y_i$$

Student	Blood Sugar (x1 mg/dL)	BMI	y	Predicted \hat{y}	Loss
1	120	28	1	1	0
2	130	30	1	1	0
3	115	25	1	-1	99.0
4	85	20	1	-1	74.2
5	75	18	1	1	0
6	90	22	-1	1	79.1
7	95	24	-1	1	84.1
8	100	25	-1	-1	0
9	80	20	-1	1	70.7
10	85	21	-1	-1	0

[How many computations per update] ~ [individual losses] where n is no. of datapoints

$$0.1x_1 - 0.05x_2 + 8.585 \Rightarrow w_1 x_1 + w_2 x_2 + w_0$$

$$\text{Total Loss} = 407.1$$

$$w_1 = 0.1$$

$$w_2 = -0.05$$

Adding up losses of how many points - 10 points

$$w_0 = 8.585$$

{What happens next in order to minimize loss?}

$$\left[(\omega_1, \omega_2, \omega_0)^{\text{new}} = (\omega_1, \omega_2, \omega_0)^{\text{old}} - \eta \{ \nabla l \} \right]$$

After $\omega_1, \omega_2, \omega_0$ are updated

$$\left\{ \begin{array}{l} -0.35x_1 - 0.162x_2 + 8.580 \\ \omega_1 = -0.35 \\ \omega_2 = -0.162 \\ \omega_0 = 8.58 \end{array} \right\}$$

Student	x1	x2	y	f_new	\hat{y}_{new}	Loss
1	120	28	1	-0.35*120 - 0.162*28 + 8.58 = -37.956	-1	37.956
2	130	30	1	-45.06	-1	45.06
3	115	25	1	-35.955	-1	35.955
4	85	20	1	-26.64	-1	26.64
5	75	18	1	-23.346	-1	23.346
6	90	22	-1	-0.35*90 - 0.162*22 + 8.58 = -26.484	-1	0
7	95	24	-1	-28.338	-1	0
8	100	25	-1	-30.03	-1	0
9	80	20	-1	-24.86	-1	0
10	85	21	-1	-26.062	-1	0

$$\left[\text{loss} = 37.956 + 45.06 + 35.955 + 26.64 + 23.346 = 168.957 \right]$$

How many times will I repeat this?

How many times will I repeat this?

{ Until I feel I have
reached a minima }

{ [Vanilla / Batch Gradient Descent] }

[Stochastic Gradient Descent]

29 September 2025 21:26

Randomly initialize weights (1000 times)

Student	Blood Sugar (x1 mg/dL)	BMI	y	Predicted \hat{y}	Loss
1	128	28	1	1	0
2	130	30	1	1	0
3	115	25	1	-1	99.0
4	85	20	1	-1	74.2
5	75	18	1	1	0
6	90	22	-1	1	79.1
7	95	24	-1	1	84.1
8	100	25	-1	-1	0
9	80	20	-1	1	70.7
10	85	21	-1	-1	0

{ Dont take overall loss }

1 computation

$$\left\{ \begin{array}{l} w_1 = 0.01 \\ w_2 = 0.5 \\ w_0 = -9 \end{array} \right\}$$

$$\left\{ w_{\text{new}} = w_{\text{old}} + \eta \nabla l \rightarrow \left(- \frac{w_1^T u_i + w_0 \times y_i}{\|w\|} \right) \right\}$$

{ loss of only one point }

{ 1 computation per weight update }

Stochastic Gradient Descent

↓
Randomness/simulation

[Will we still get similar result?]

[Yes we will get a similar result]

{ G.D \Rightarrow Access to all points } My entire dataset

1. Batch Gradient Descent

- * Computes the gradient of the loss function across the entire dataset:

$$\nabla_{\theta} L(\theta) = \sum_{i=1}^N \nabla_{\theta} \ell(x_i, y_i; \theta)$$

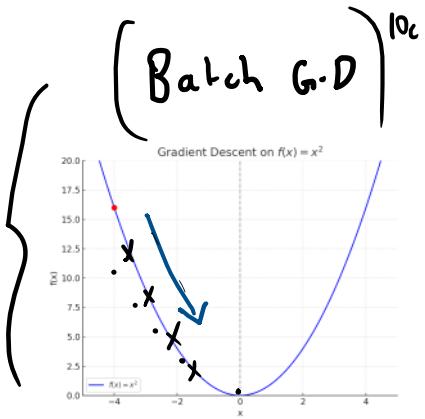
- This is precise because you're using all data points, but for **large datasets**, computing this sum at every step is expensive.
- Each update happens **after seeing the whole dataset**, so the model updates less frequently.

2. Stochastic Gradient Descent

- * Uses only one data point (or a small batch) to estimate the gradient:

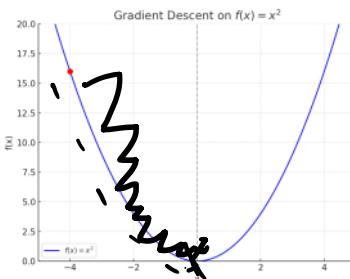
$$\nabla_{\theta} \ell(x_i, y_i; \theta)$$

- [The gradient is **noisy**] because it's based on one sample rather than the whole dataset.
- But it's much faster per update because you avoid summing over N points.



(single points)

✓
Stochastic

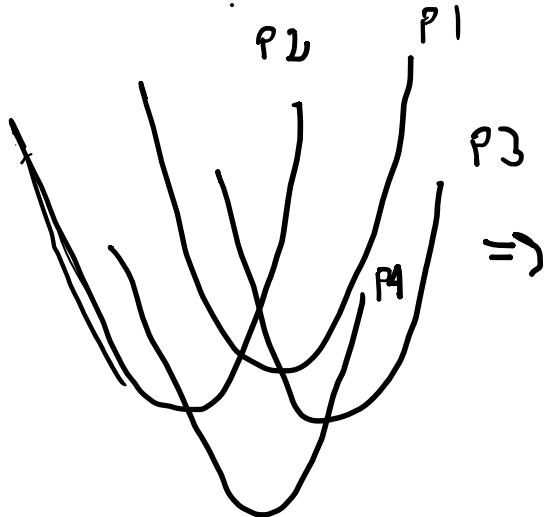


in individual points
have
individual loss
shapes

If I need to update weights
100 times for regular Gr-D,

Variety into
gradient descent

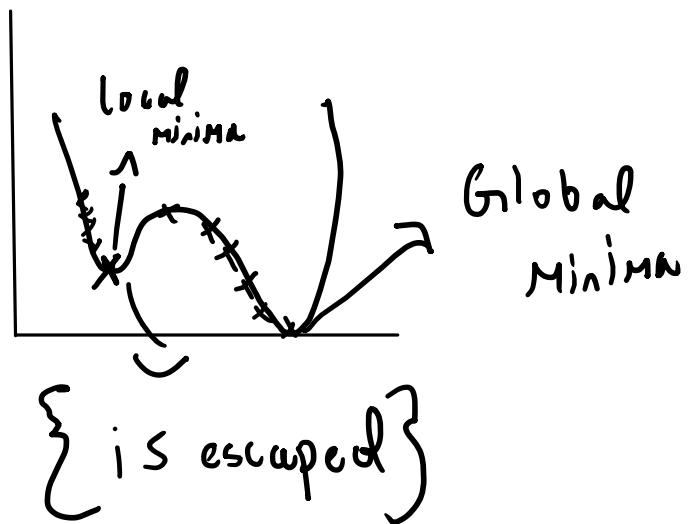
For stochastic, I will need to do



Variety is precious

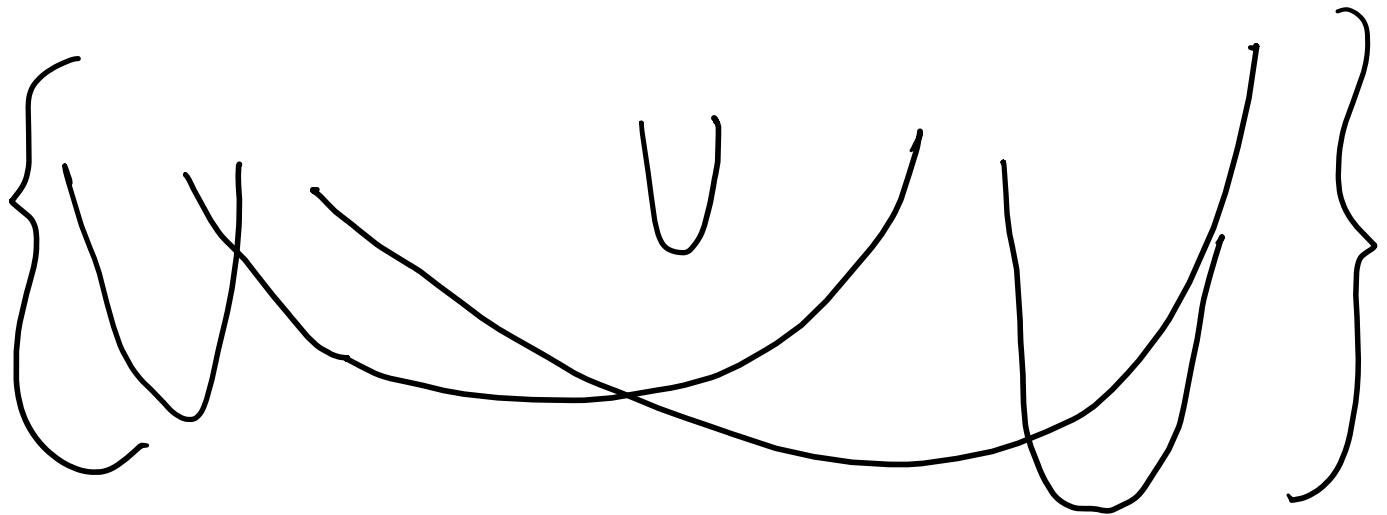
Learn better

It {might} help us solve the [non-convex problem]



{ Disadvantage ? } Double edged sword

When each datapoint is too different
from each other



{Variety can backfire}

Chance of missing global Minima

(Vanilla / Batch GD)

If 100 points,

$\rightarrow n_{100}$ points per weight update

(Stochastic GD)

\rightarrow 1 point per weight update

Subset of the points

10 points at a time } \Rightarrow randomly

$$\therefore \left\{ \begin{array}{l} \text{10} \\ \text{---} \\ i=1 \end{array} \right\} \left(\frac{\mathbf{w}^T \mathbf{x}_i + w_0}{\|\mathbf{w}\|} \right) y_i$$

{ Mini-Batch Gradient Descent }

{Best of both worlds}

Interview :-

- { i) Loss ✓ }
- ii) Speed ✓
- iii) Noise ✓ }

Epoch vs iterations

10 September 2025 06:52

Student	Blood Sugar (x1 mg/dL)	BMI	y	Predicted \hat{y}
1	120	28	1	1
2	130	30	1	1
3	115	25	1	-1
4	85	20	1	-1
5	75	18	1	1
6	90	22	-1	1
7	95	24	-1	1
8	100	25	-1	-1
9	80	20	-1	1
10	85	21	-1	-1

100

(Interview)

(Epoch)

One round of visiting

all points in my dataset

100 points ~

Stochastic Gradient Descent

How many updates per epoch \Rightarrow (100)

(1 epoch)

Mini batch Gradient Descent

Batch size = 10

How many updates per epoch \Rightarrow 10

10 batches

$x_1 \rightarrow x_{10}$

$x_n \rightarrow x_{n+9}$

$x_{20} \rightarrow x_{30}$

Iterations \rightarrow Number of times the

Weight is updated

Interview :-

{ Difference b/w epoch and iteration }

2 epochs

If my batch size is 10

{ 20 updates / Iterations }

Coding Gradient Descent

11 September 2025 00:01

[https://colab.research.google.com/drive/14aBoDBQly6rRAoajnvnf1Ct-JgTO--uj?authuser=1
#scrollTo=c157b062-2759-4cc0-aed0-333eb871838e](https://colab.research.google.com/drive/14aBoDBQly6rRAoajnvnf1Ct-JgTO--uj?authuser=1#scrollTo=c157b062-2759-4cc0-aed0-333eb871838e)

How can I minimize a function $[3x^2 + 4]?$

$$\frac{d(3x^2+4)}{dx} = 0$$

$$6x = 0$$

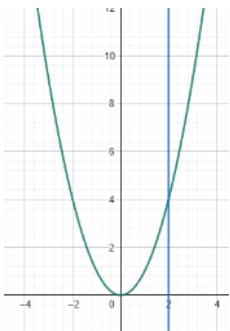
$$x = 0$$

$$f(x) = 3 \times 0 + 4 = 4$$

Constrained Optimization

04 September 2025 18:52

{ Minimize $f(x) = x^2$ such that $\underline{x=2}$ }



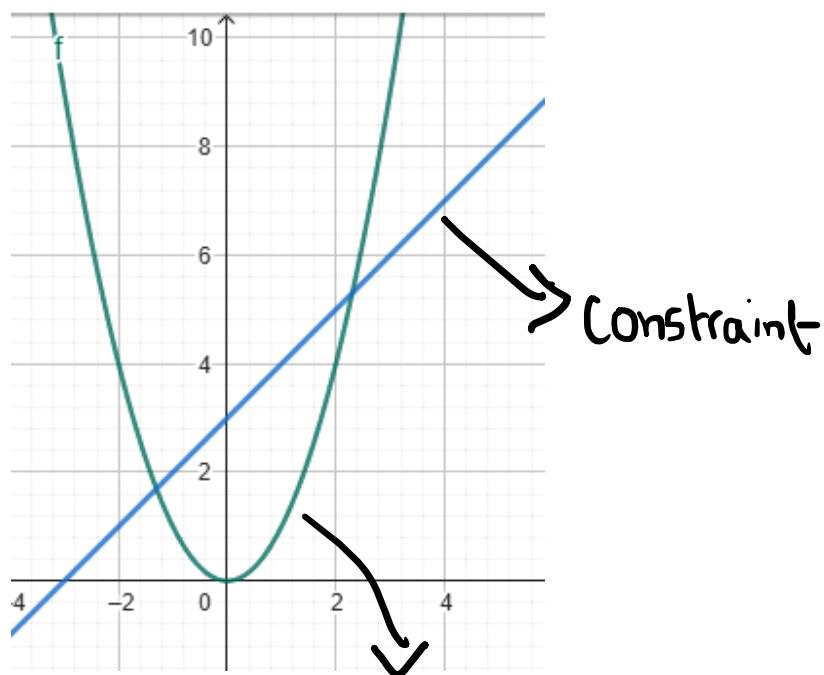
Minimize $f(x) = x^2$ such that

$$y = x + 3$$

$$x^2 = x + 3$$

$$x^2 - x - 3 = 0$$

$$\left\{ \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \right\}$$

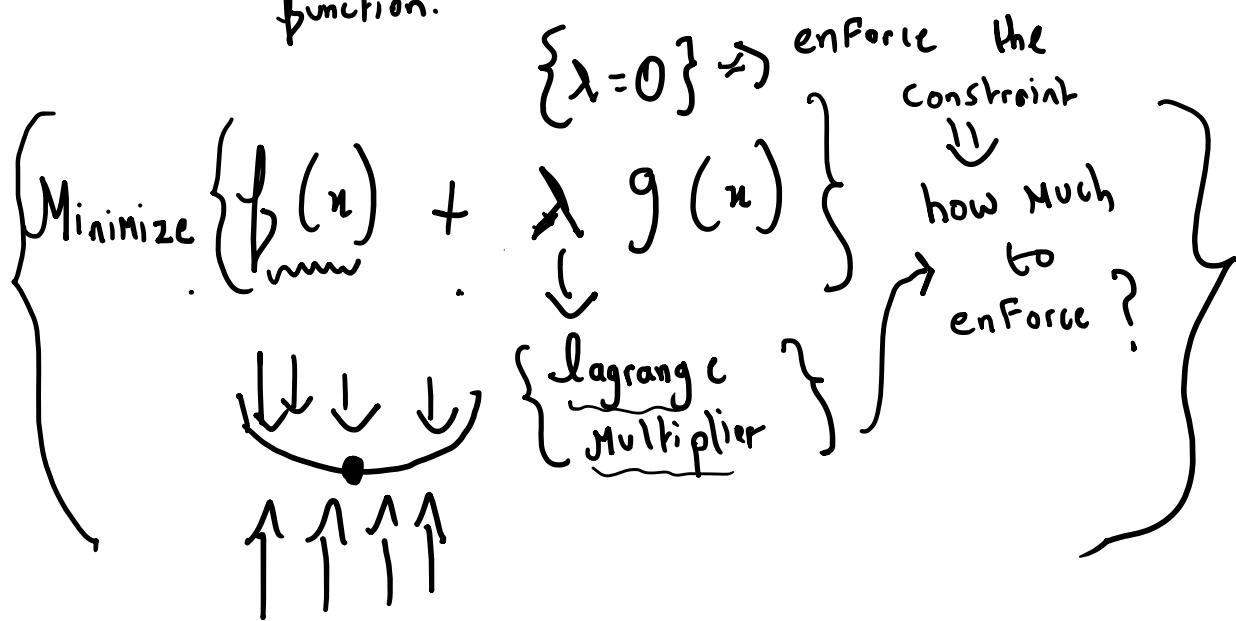


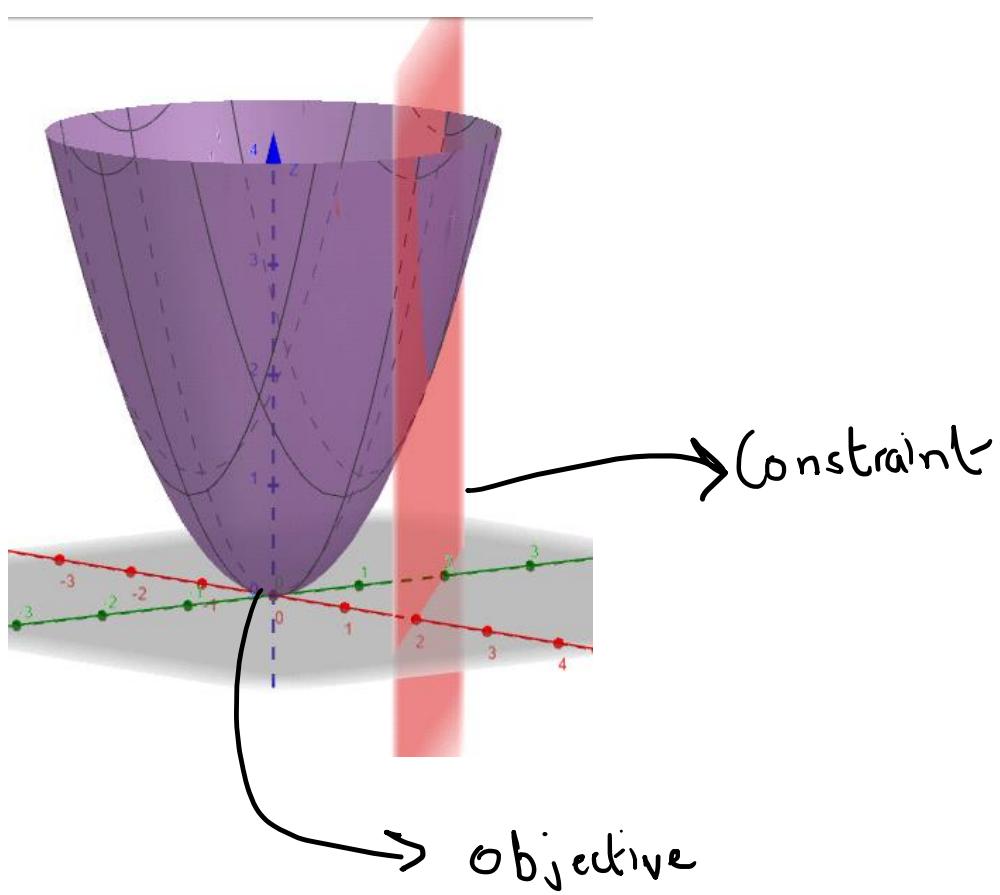
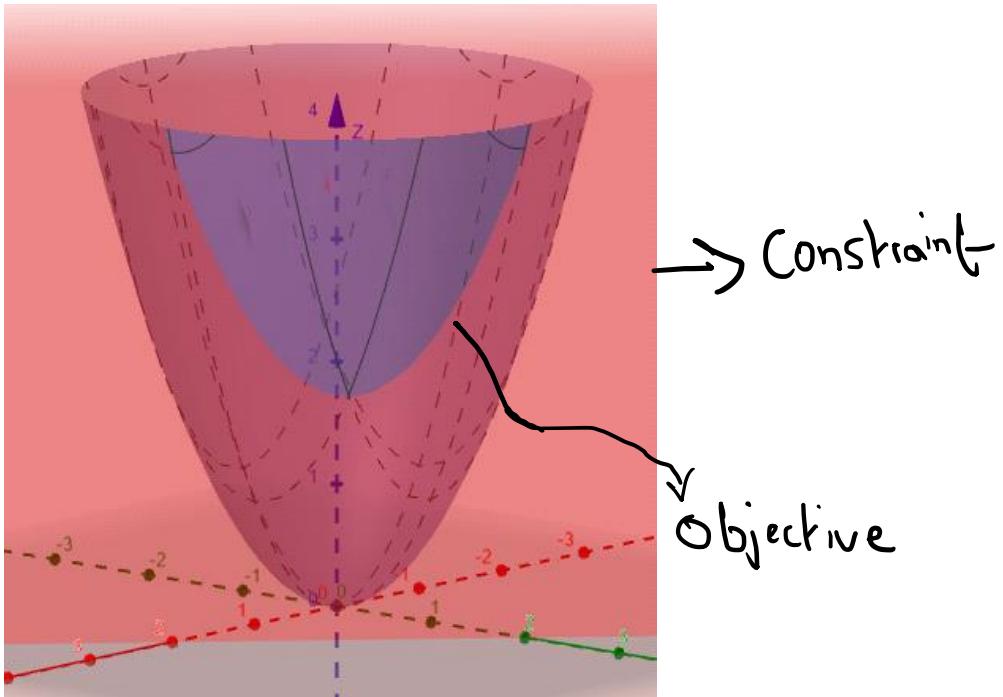
Objective

$$\begin{aligned} \text{Minimize } & (x^2 + y^2) \Rightarrow f(x) \\ (x+y-2=0) & \Rightarrow g(x) \end{aligned}$$

Instead of substituting.

Add the constraint to the objective function.





x, y, λ

$$\underline{x^2} + \underline{y^2} - \lambda (x + y - 2)$$

$$1^2 + 1^2 - \lambda (1+1-2)$$

$$(1+1-2 \times 0 = 2)$$

Minimum value is $Z = 2$

on
 $f(x,y) = 2$

where $x,y = (1,1)$

(Lagrangian multipliers?
Where?)