

Gradient Descent

x	y
1	1
2	2
3	3

$$\hat{y} = \theta_0 + \theta_1 x$$

$$e_i = y_i - \hat{y}_i = y_i - (\theta_0 + \theta_1 x_i)$$

$$e_1 = 1 - \theta_0 - \theta_1$$

$$e_2 = 2 - \theta_0 - 2\theta_1$$

$$e_3 = 3 - \theta_0 - 3\theta_1$$

$$e_1^2 = (1 - \theta_0 - \theta_1)^2$$

⋮

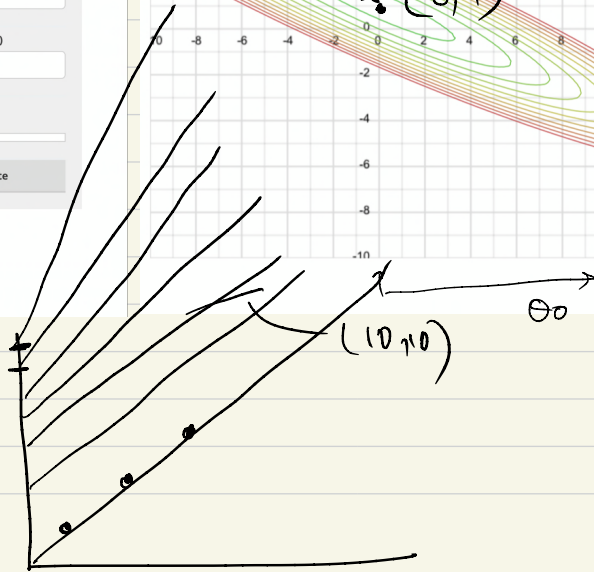
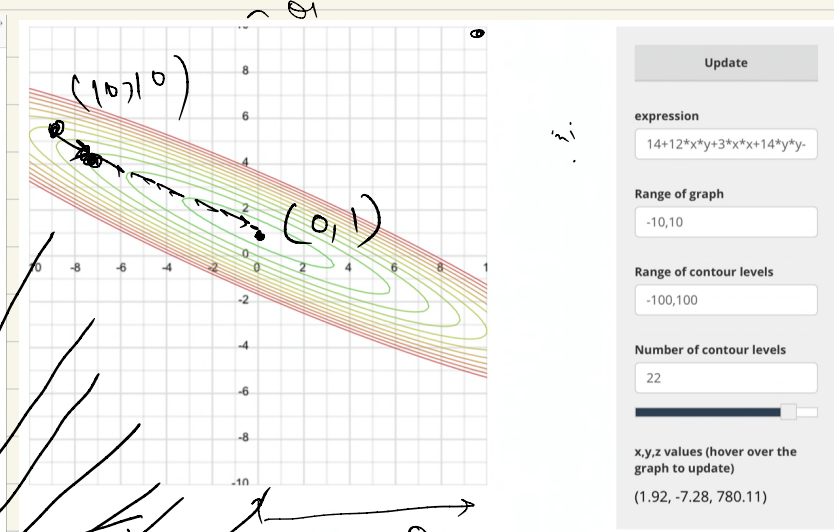
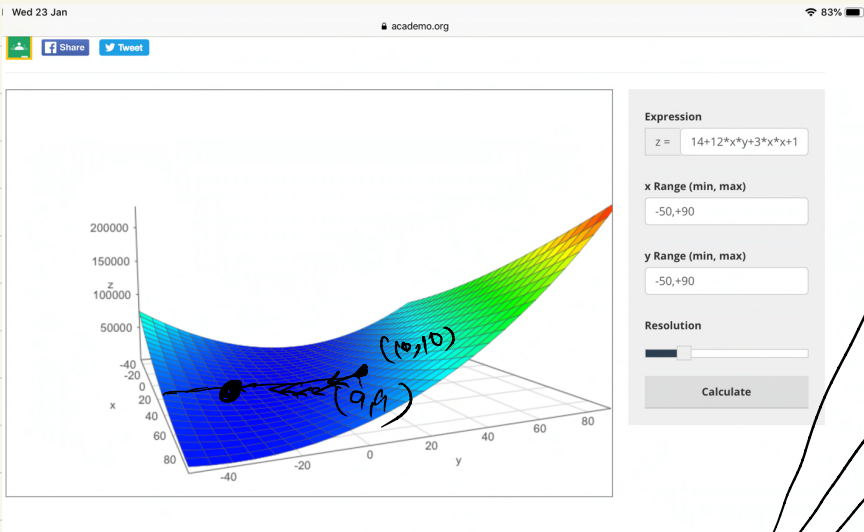
$$\sum e_i^2 = f(\theta_0, \theta_1) = (1 + \theta_0^2 + \theta_1^2 - 2\theta_0 - 2\theta_1 + 2\theta_0\theta_1) + (4 + \theta_0^2 + 4\theta_1^2 - 4\theta_0 + 4\theta_0\theta_1 - 8\theta_1) + (9 + \theta_0^2 + 9\theta_1^2 - 6\theta_0 - 18\theta_1 + 6\theta_0\theta_1)$$

$$= 14 + 3\theta_0^2 + 14\theta_1^2 - 12\theta_0 - 28\theta_1 + 12\theta_0\theta_1$$

$$3x^2 + 14y^2 - 2x - 21y + 12xy + 14$$

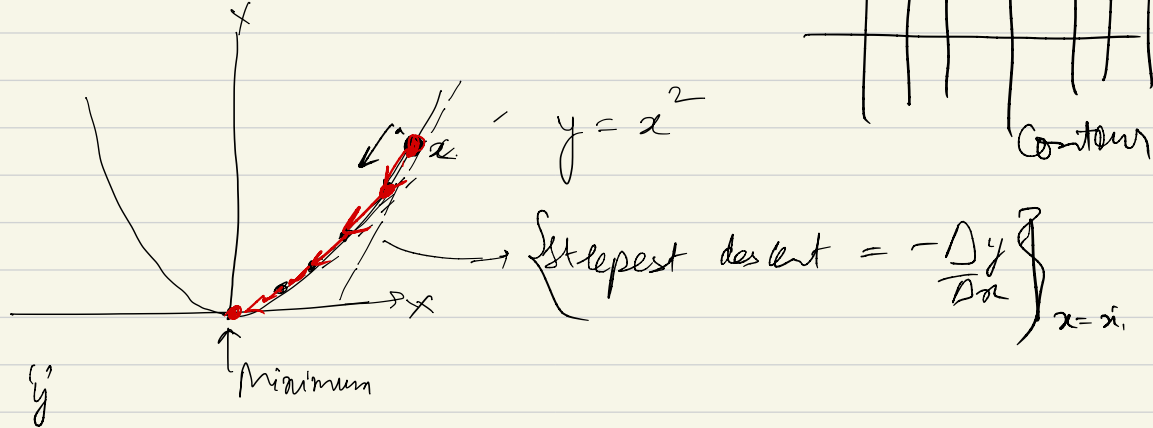
3D Plot

Contour Plot



GRADIENT DESCENT

* General optimisation technique



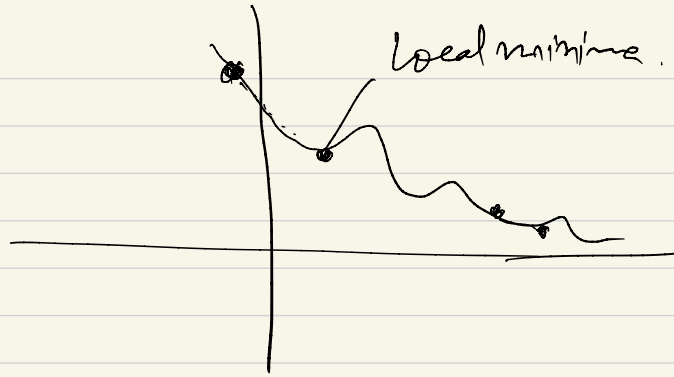
Question: Find minimum y

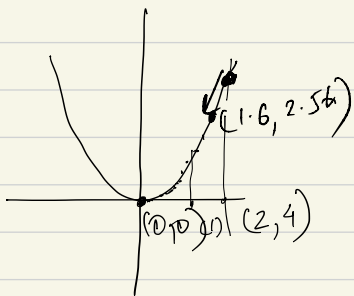
① start with some (x_0)

② Move from (x_0) to x_1 s.t. $y_1 < y_0$.

③ repeat ②

$$x = x - \alpha \frac{dy}{dx} |_{x_1}$$





$$\alpha = 0.1$$

$$x_0 = 2$$

$$y_0 = 4$$

$$\frac{\Delta y}{\Delta x} = 2x$$

$$\therefore y = x^2$$

Solving y. D.

$$x_1 = x_0 - 2 \frac{\Delta y}{\Delta x} \Big|_{x_0} = x_0 - 2x_0 = -8x_0$$

$$x_1 = 0.8 \times 2 = 1.6$$

$$x_2 = x_1 - 2x_1 = -8x_1 = 1.28$$

$$x_N = (0.8)^N x_0$$

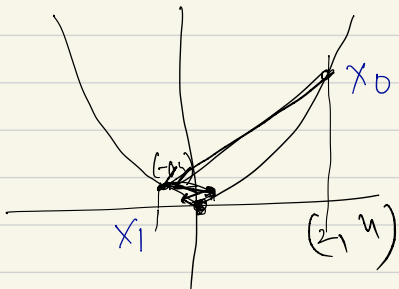
$$x_N \rightarrow 0 \text{ for high } N.$$

OK, what if we start from $x_0 = -2$



$$x_1 = x_0 - 2 \frac{\Delta y}{\Delta x} \Big|_{x_0} = 0.8x_0 = -1.1$$

Q: what if ' α ' is large



$$\alpha = 0.6$$

$$x_0 = 2$$

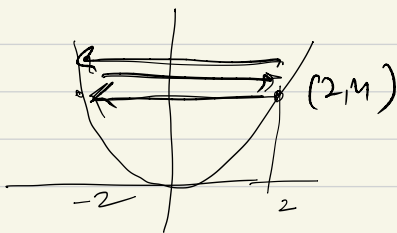
$$x_1 = x_0 - \alpha \frac{\Delta y}{\Delta x} \Big|_{x_0} = x_0 - (0.6)(2x) \Big|_{x_0} = -0.2x_0$$

$$x_2 = -0.2x_1 = (-0.4)(-0.2) = -0.08$$

what if $\alpha = 1$

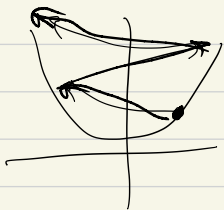
$$x_1 = x_0 - (1)(2x) \Big|_{x_0} = -x_0$$

$$x_2 = -x_1$$



$$\alpha = 2$$

$$x_1 = x_0 - (2)(2m) \Big|_{x_0} = -3x_0$$



DIVERGING

$$\alpha = 10^{-3} \quad (\alpha \text{ TOO SMALL?})$$

$$x_1 = x_0 - (10^{-3})(2m) \Big|_{x_0} \approx x_0$$

TOO MANY STEPS



α TOO SMALL \rightarrow LOW

α TOO LARGE \rightarrow DIVERGE

Many
iterations

$$\alpha = \frac{K}{L} \frac{1}{\sqrt{A}}$$

$$\frac{K}{\sqrt{AK_0}}$$

G.D. for lin reg.

$$\hat{y}_i = \theta_0 + \theta_1 x_i$$

$$e_i = y_i - \hat{y}_i = y_i - (\theta_0 + \theta_1 x_i)$$

$$\sum e_i^2 = \sum (y_i - (\theta_0 + \theta_1 x_i))^2$$

① start with random values of θ_0 & θ_1

till convergence:

$$\theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} (\sum e_i^2)$$

$$\theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} (\sum e_i^2)$$

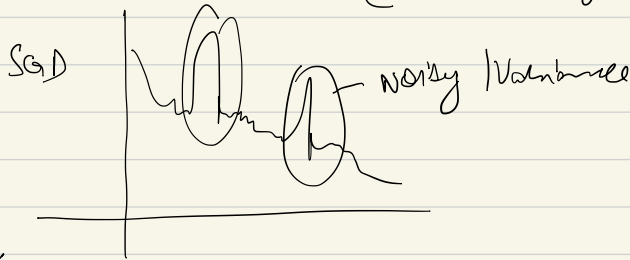
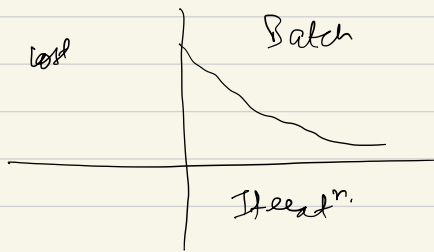
} simultaneously

$$\text{Now } \frac{\partial}{\partial \theta_0} (\sum e_i^2) = 2 \sum (y_i - (\theta_0 + \theta_1 x_i)) (-1) \quad \& \quad \frac{\partial}{\partial \theta_1} (\sum e_i^2) = 2 \sum (y_i - (\theta_0 + \theta_1 x_i)) (-x_i)$$

Gradient Descent Variants

① Vanilla (or Batch) : update params after going through all data (Slower update)

② S.G.D : " " " " for each data point (more noisy)



mini batch (Best of both worlds)

Gradient Descent

- ① Good for online setting (more data).
- ② Good for large data

Normal Equations

- ① Good for simple
- ② No need to worry about learning rate, etc.
- ③ Non trivial to solve (impossible...)

Projected Gradient Descent

Q) Learn a fit where $\theta_0, \dots, \theta_d$ are all ≥ 0

$$\theta_i = \max\left(\theta_i - \alpha \frac{\Delta \xi(\theta_1, \dots, \theta_i, \dots, \theta_d)}{\Delta \theta_i}, 0\right)$$

