

Gradient Descent

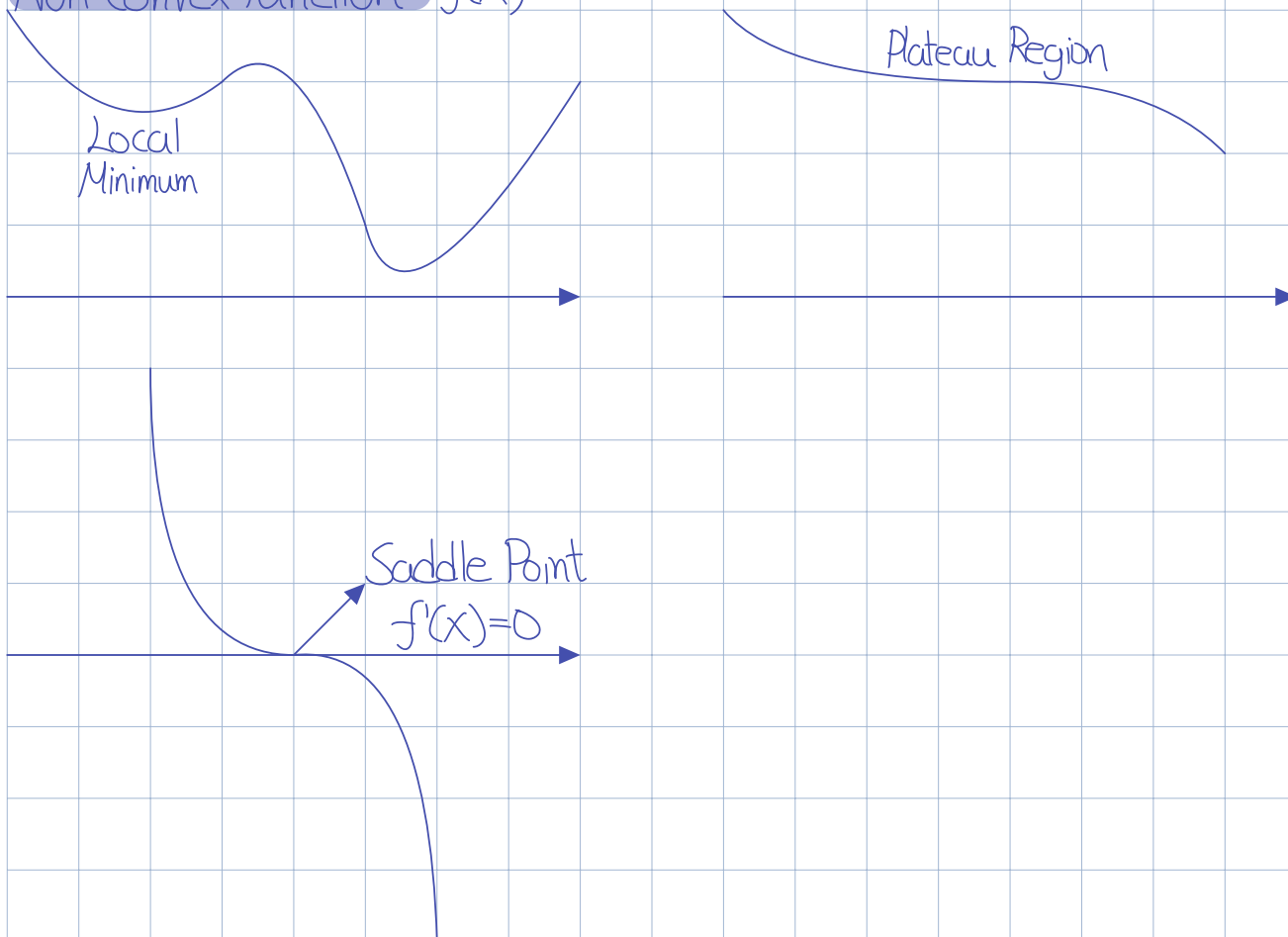
$$\min_{x \in \mathbb{R}^n} f(x) \quad x_{t+1} = x_t - \epsilon_t \nabla f(x_t)$$

Convex Function

$$\|\nabla f(x_t)\|_2 = 0$$

Can be the stopping criteria

Non Convex Function $f(x)$



Rule of Thumb

1. Large number of iterations
2. $f(x)$ must be small
3. $\|\nabla f(x_t)\|_2$ must be close to 0

Momentum Based Model

Recap Basic SGD Rule

$$E_n(w) = \frac{1}{N} \sum_{n=1}^N E_n(w)$$

$$g_t = \nabla E_n(w_t)$$

$$w_{t+1} = w_t + \epsilon_t g_t$$

SGD + Momentum

Velocity Vector: v_t $v_0 = 0$

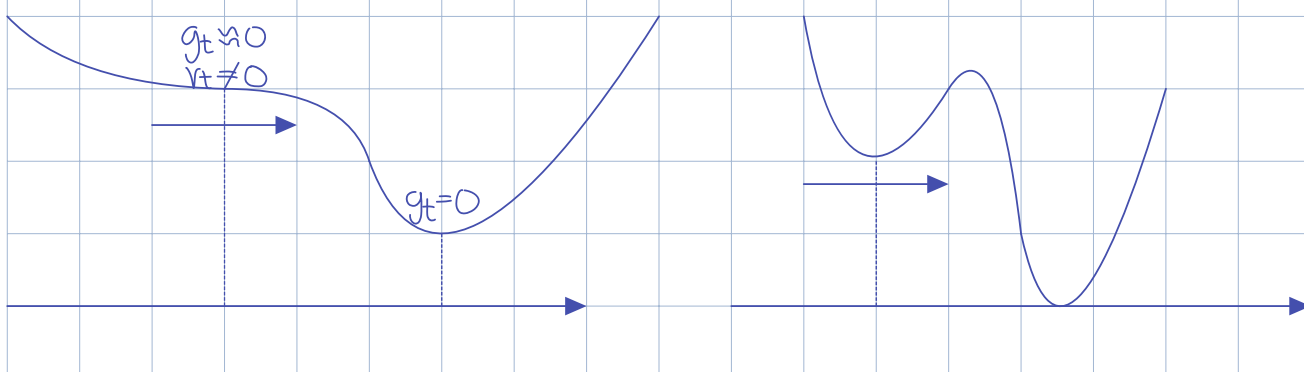
$$g_t = \nabla E_n(w_t)$$

$\mu \equiv$ momentum coefficient

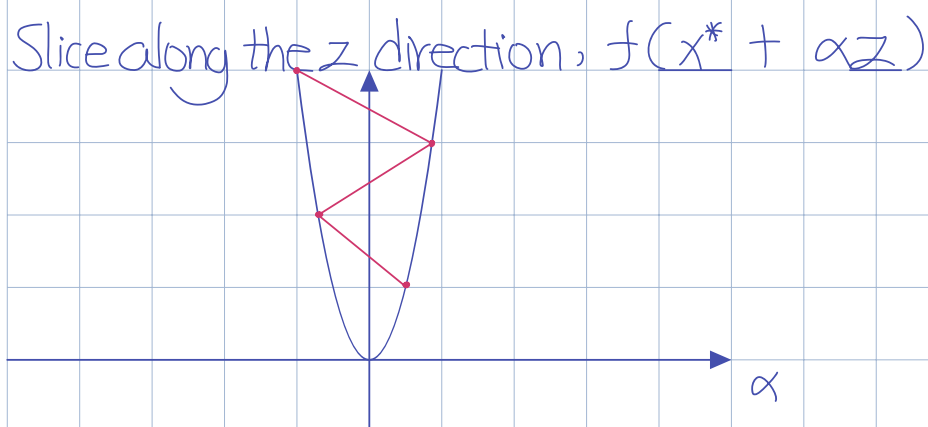
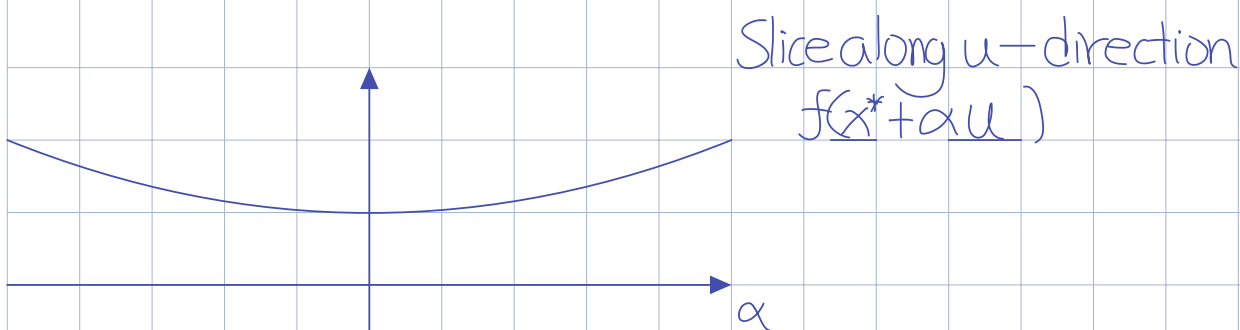
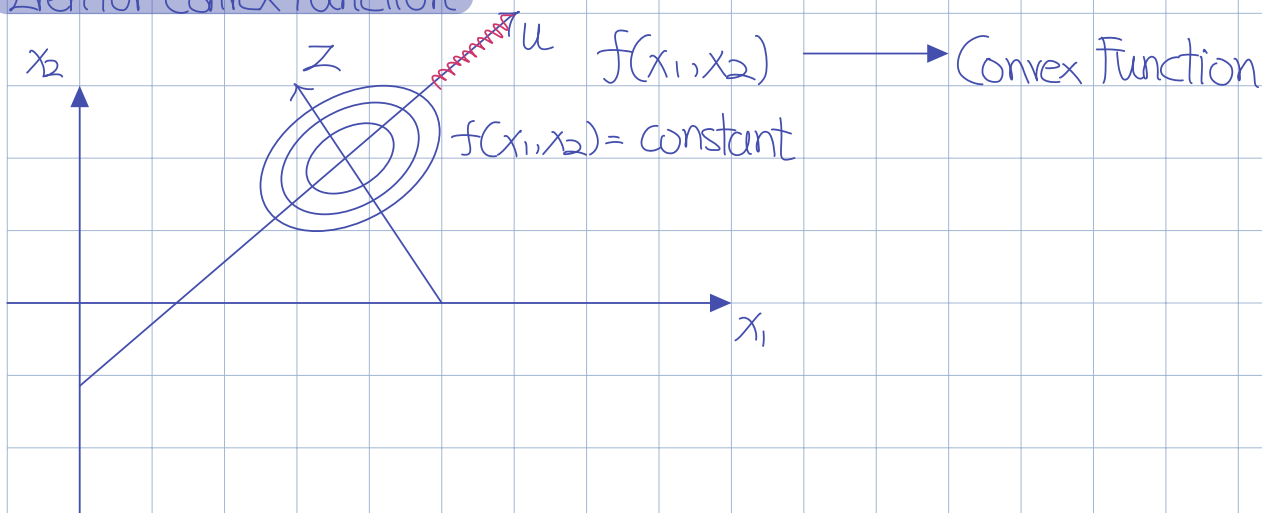
$$v_t = -\epsilon_t g_t + \mu v_{t-1}$$

$$\mu \approx 0.9$$

$$w_{t+1} = w_t + v_t$$



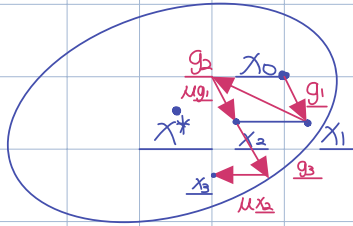
Momentum Based Method Can Led to Faster Convergence Even for Convex Function



Some direction, small gradient
Some direction, large gradient

Application of Momentum Method ($\epsilon_t = 1$)

Step 1, $v_0 = 0$



$$g_1 = -\nabla f(x_0)$$

$$v_1 = g_1$$

$$x_1 = x_0 + v_1$$

Step 2,

$$g_2 = -\nabla f(x_1)$$

$$v_2 = g_2 + \mu g_1$$

$$x_2 = x_1 + v_2$$

Step 3,

$$g_3 = -\nabla f(x_2)$$

$$v_3 = g_3 + \mu g_2$$

$$x_3 = x_2 + v_3$$

Step 4,

$$g_4 = -\nabla f(x_3)$$

$$v_4 = g_4 + \mu g_3$$

$$x_4 = x_3 + v_4$$

Nesterov Momentum

$$Y_t = \mu Y_{t-1} - \epsilon_t \nabla E_n(m_t + \mu Y_{t-1})$$

Intermediate Point Between
 m_t and m_{t+1}

$$m_{t+1} = m_t + Y_t$$

$$= \underbrace{m_t + \mu Y_{t-1}} - \epsilon_t \nabla E_n(m_t + \mu Y_{t-1})$$

For Convex Functions, "Convergence Rate" (1983)

1. GD, $O(1/k)$

2. Nesterov Momentum, $O(1/k^2)$

AdaGrad (adaptive step size)

$$g_t = \nabla E_n(m_t) \in \mathbb{R}^d$$

$$r_t = r_{t-1} + g_t \odot g_t \quad \text{Componentwise Vector Multiplication}$$

$$r_t(j) = \sum_{q=1}^t g_q^2(j), j=1, 2, \dots, d$$

$$m_{t+1}(j) = m_t(j) - \frac{\epsilon_t}{\sqrt{r_t(j)}} g_t(j)$$

RMS-PROP

$$\underline{r_t} = (1 - \beta) \underline{r_{t-1}} + \beta \underline{g(t) \odot g(t)}$$

β = decay factor

$$\underline{w_{t+1}} = \underline{w_t} - \frac{\epsilon_t}{\sqrt{\underline{r(t)}}} \odot \underline{g(t)}$$

ADAM (RMS-PROP + Momentum)

$$\underline{g_t} = \nabla \text{en}(\underline{w_t})$$

$$\underline{r_t} = \beta \underline{r_{t-1}} + (1 - \beta) \underline{g_t \odot g_t}$$

$$\underline{v_t} = \mu \underline{v_{t-1}} - \frac{\epsilon_t}{\sqrt{\underline{r(t)}}} \odot \underline{g(t)}$$

$$\underline{w_{t+1}} = \underline{w_t} + \underline{v_t}$$