

Lecture 14. Advanced optimization Methods for DL.

1. Optimization methods in DL. (chapter 3, The science of DL)

① First-order methods

GD/SGD:

$$w^i = w^{i-1} + \underbrace{(-\epsilon)}_{\text{learning rate}} \cdot \underbrace{\nabla_w L}_0$$

First-order derivative.

② Second-order methods.

Newton's method:

$$\Delta W_{\text{nst}} = - \underbrace{\left[\nabla_w^2 L \right]^{-1}}_0 \nabla_w L$$

newton step:

Quasi-Newton

$$\nabla^2 L = \begin{bmatrix} \nabla^2 L \\ \nabla^2 L \end{bmatrix}_{n \times n} \quad L: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Hessian matrix,

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \quad \nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}_{n \times 1}$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m \quad \nabla f = \begin{pmatrix} \frac{\partial f^1}{\partial x_1} & \dots & \frac{\partial f^1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f^m}{\partial x_1} & \dots & \frac{\partial f^m}{\partial x_n} \end{pmatrix}$$

$\begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{pmatrix}$

③ Evolutionary Learning / Genetic algorithms.

fitness function.

2. Adaptive GD.

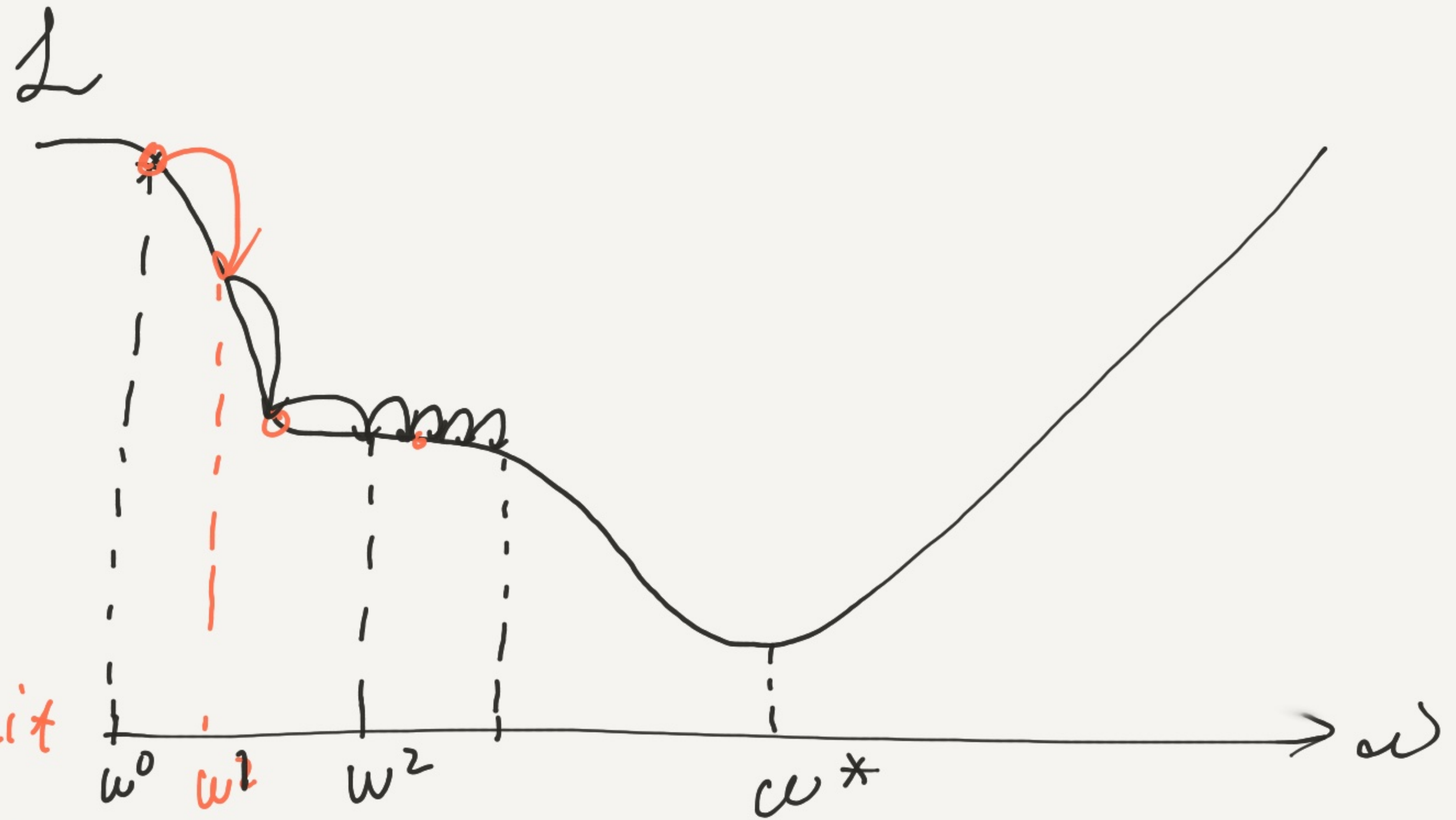
GD/SGD:

$$\omega^0, \omega^1, \omega^2, \dots, \omega^K = \omega^*$$

$$\omega^i = \omega^{i-1} + (\underbrace{\epsilon}_{\text{predefined}} \cdot \underbrace{(-\nabla_{\omega} L)}_{\text{not unit vector}})$$

predefined

not unit vector.



Step size is determined by both the ϵ and $\nabla_{\omega} L$, if $\nabla_{\omega} L$ is very small (flat region), our model converges very slow.

2.1. SGD with momentum.

△ goal is to accelerate learning by accumulating a moving average of past gradients.

$$\text{SGD: } w^{i+1} = w^i + \underbrace{\varepsilon \cdot (-\nabla_w L)}_{\Delta W_{bp} \rightarrow \text{combined step size}} = w^i + \Delta W_{bp}$$

$$\left. \begin{array}{l} w^0 \rightarrow \Delta W_{bp}^0 = \varepsilon \cdot (-\nabla L) \\ w^1 \rightarrow \Delta W_{bp}^1 \\ w^2 \rightarrow \Delta W_{bp}^2 \\ \vdots \\ w^{k-1} \rightarrow \Delta W_{bp}^{k-1} \end{array} \right\}$$

moving average: m

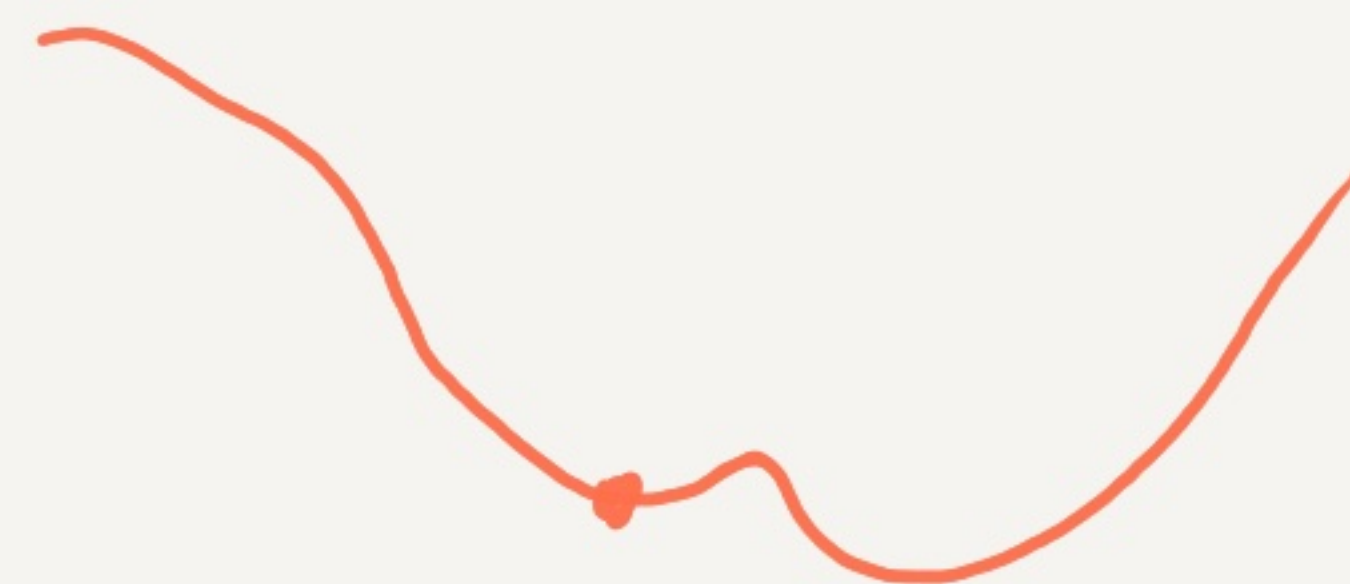
$$\left\{ \begin{array}{l} m_0 = \Delta W_{bp}^0 \\ m_1 = (\Delta W_{bp}^0 + \Delta W_{bp}^1) / 2 = \frac{m_0}{2} + \frac{\Delta W_{bp}^1}{2} \\ m_2 = (\Delta W_{bp}^0 + \Delta W_{bp}^1 + \Delta W_{bp}^2) / 3 \\ \quad = \frac{2m_1}{3} + \frac{\Delta W_{bp}^2}{3} \\ \vdots \\ m^{k-1} = \frac{(k-2) \cdot m_{k-2}}{k-1} + \frac{\Delta W_{bp}^{k-1}}{k-1} \end{array} \right.$$

SGD with momentum:

$$w^{i+1} = w^i + \Delta w^i$$

$$\Delta w^i = \frac{2}{\sigma} (\Delta w^{i-1}) + \frac{(1-\sigma)}{\sigma} \Delta W_{bp}^i$$

$$0 < \sigma < 1$$



2.2 AdaGrad (2011)

It aims to adapt ϵ by scaling inversely proportional to the square root of accumulated gradients.

$$\begin{aligned} \text{Step size: } \Delta w &= \frac{1}{\sqrt{r} + \delta} (\epsilon \cdot -1 \cdot \nabla_w L) \\ &= \underbrace{\frac{\epsilon}{\sqrt{r} + \delta}}_{\text{actual learning rate}} \underbrace{(-\nabla_w L)}_{\text{initial learning rate}} \end{aligned}$$

r : accumulated gradients.

$$\begin{aligned} r^i &= r^{i-1} + \underbrace{(\nabla_w L)^2}_{\text{size of the gradient}} \xrightarrow{\text{L}^2 \text{ norm for vectors.}} \\ &= \sum_{m=1}^{i-1} (\nabla_w L)^2 \end{aligned}$$

δ : small positive constants, $(10^{-7}, 10^{-6})$
to avoid division by zero.

2-3. RMS prop (Hinton, 2012)

Root mean square propagation
RMSprop.

Basic idea: Discard history from extreme past.

$$\Delta W = \frac{\varepsilon}{\sqrt{\underbrace{r+2}_{\text{constant}}}} \cdot (-\nabla_w L)$$

$$r^{i+1} = \underbrace{\rho \cdot r^i}_{\substack{\text{historical} \\ \text{decay rate}}} + (1-\rho) \cdot (\nabla_{w^i} L)^2 \quad 0 < \rho < 1$$

$$r^2 = \rho r^1 + (1-\rho) (\nabla_{w^1} L)^2$$

$$r^3 = \rho r^2 + (1-\rho) (\nabla_{w^2} L)^2$$

$$= \rho^2 r^1 + \rho(1-\rho) (\nabla_{w^1} L)^2 + (1-\rho) (\nabla_{w^2} L)^2$$

$$r^n = \rho^n \cdot r^1 + \dots$$

Contribution of r^1 is determined by.

$$\begin{aligned} 0.9 &= \rho, \\ 0.81 &= \rho^2, \\ 0.729 &= \rho^3, \\ 0.6561 &= \rho^4, \\ &\vdots \\ &\rho^n \end{aligned} \quad \left. \begin{array}{l} \rho > \rho^2 > \rho^3 \\ \downarrow \end{array} \right\}$$