

Lec 12 GD / SGD non-convex func

1. For convex functions:

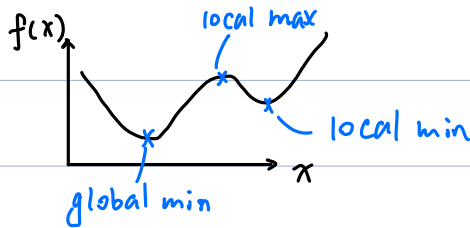
e.g. logistic regression. linear regression

the $E_{in}(\underline{w})$ is convex w.r.t the parameter w .

\Rightarrow we will get a global minimum when $\nabla E_{in}(w) = 0$

2. Non-convex functions:

\hookrightarrow e.g.



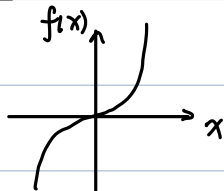
\hookrightarrow those points are called saddle points:

$\nabla f(\underline{x}) = \underline{0}$, but \underline{x} is neither a local min or local max

\uparrow
include local min & global min

\hookrightarrow example in 1-D space:

$$f(x) = x^3$$



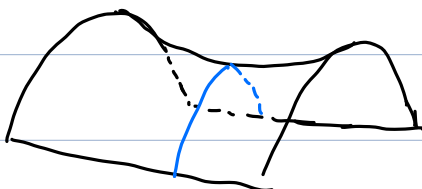
\Rightarrow very slow progress at saddle point

\Rightarrow Need to pre-set No. of iteration runs out

\hookrightarrow In n-D case:

① 1-D saddle point in any direction

② local min in some direction, and local max in some other direction



saddle shape

③ both ① and ②

⇒ highly likely to have saddle points in high-dimensional space

↳ stopping conditions:

• aim: to reduce the chance of stopping at a saddle point

• is saddle point if

- ① $\nabla E_{in}(\underline{w}) = 0$
- ② $\nabla E_{in}(\underline{w})$ is small
- ③ No. of iteration is large

⇒ solution: SGD with momentum

3. SGD with Momentum: (By Polyak, 1964)

↳ recall basic SGD:

$$\underline{g}_k = \nabla E_{in}(\underline{w}_k)$$

$$\underline{w}_{k+1} = \underline{w}_k - \epsilon_k \underline{g}_k$$

↳ SGD with momen

$$\underline{g}_k = \nabla E_{in}(\underline{w}_k)$$

$$\underline{v}_k = -\epsilon_k \underline{g}_k + \mu \underline{v}_{k-1}$$

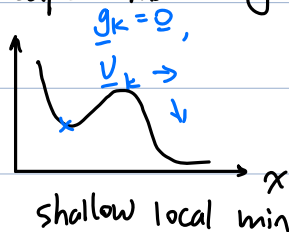
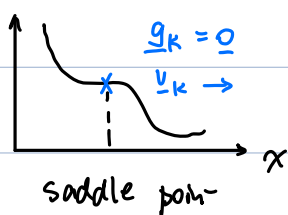
$$\underline{w}_{k+1} = \underline{w}_k + \underline{v}_k$$

// commonly $\mu = 0.9$

→ accumulation of prev momentums

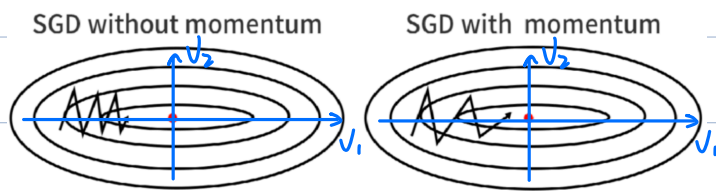
↳ usage:

① Momentum helps SGD escape that regions



"heavy ball momentum"

② Momentum can lead to faster convergence,
even for convex function



$f(x)$, $d=2$.

GD strongly favors
the v_2 direction
 \Rightarrow Not ideal

↓
Moving along v_1 direction

4. Nesterov momentum (1983)

$$\hookrightarrow \underline{v}_k = -\epsilon_k \nabla_{\text{en}}(\underline{w}_k + \mu \underline{v}_{k-1}) + \mu \underline{v}_{k-1}$$

$$\underline{w}_{k+1} = \underline{w}_k + \underline{v}_k$$

\hookrightarrow compare:

original version is : $\underline{v}_k = -\epsilon_k \nabla_{\text{en}}(\underline{w}_k) + \mu \underline{v}_{k-1}$
here is different

\hookrightarrow Can prove better convergence.

For convex func,

\rightarrow Full GD: distance between \underline{w}_k and \underline{w}^* is $O(\frac{1}{k})$

\rightarrow with Nesterov momentum : $O(\frac{1}{k^2})$