

# Stochastic Gradient Descent

Thursday, January 30, 2020

13:09

So far the gradient descent we've seen is

"Batch Gradient Descent"

—— where we use all of in-sample data

there are more types/configurations of gradient descent:

→ "Mini-Batch" GD:

use a subset of total in-sample for gradient descent

"Stochastic Gradient Descent (SGD)"

⊕ save computation per update

⊖ less accurate

⊕ faster to use many smaller "noisy" GDs than a few accurate GDs

⊕ "noisy" helps with generalization.

extreme case of mini batch: Single sample

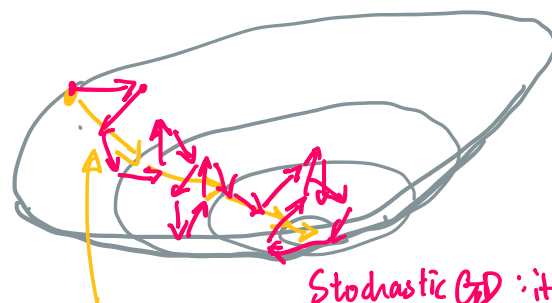
$$\underline{w}(t+1) = \underline{w}(t) - \eta_t \nabla_{e_{n(t)}}(\underline{w}(t))$$

↑  
just a single error

the average of individual gradient of error is the actual gradient.

$$\mathbb{E}[\nabla_{e_{n(t)}}(\underline{w}(t))] = \frac{1}{N} \sum_{n=1}^N \nabla_{e_n}(\underline{w}(t))$$

terminology: "Epoch" — one run of GD through ALL training sample



Batch GD

Stochastic GD: it "hovers" around the minimum (erratic behaviour)

To fix: Have variable learning rate  $\eta$ :

condition:  $\sum_{t=1}^{\infty} \eta_t = \infty$  and  $\sum_{t=1}^{\infty} \eta_t^2 < \infty$  guarantees convergence

method

①  $\eta_t = (1 - \frac{t}{T}) \eta_0 + \frac{t}{T} \eta_c$ ,  $\eta_c \approx 0.01 \eta_0$   
 $t = 1, 2, 3, \dots, T$

method

②  $\eta_t = \frac{\delta}{\sqrt{t}}$

method

③ additive learning rate:

learning rate depends on - properties of  $w$   
- features  
- directions  
etc.

To stop:

- ① Enough iterations
- ② Error difference threshold  
(but computing error will affect computation time)

Adaptive Learning Rate SGD:

① AdaGrad

→ downscale  $\eta_t$  by square root of sum of all historical squares of the gradient

gradient est. =  $\hat{g}_t$

element wise multiplication

accumulate:  $\underline{r}_t = \underline{r}_{t-1} + \hat{g}_t \otimes \hat{g}_t$

update  $\underline{w}(t+1) = \underline{w}(t) - \underbrace{\eta}_{\text{current } \frac{1}{\sigma}} \underbrace{\left( \frac{1}{\delta + \sqrt{\underline{r}_t}} \right)}_{\text{element wise}} \otimes \hat{g}_t$

prevent  $\frac{1}{0}$

↑  
element wise  
squareroot  
& division

## ② RMS Prop

gradient  $\underline{g}_t$

accumulate

$$\underline{r}_t = \rho \underline{r}_{t-1} + (1-\rho) \hat{\underline{g}}_t \otimes \hat{\underline{g}}_t$$

"forget" factor  $\Rightarrow$  slows down the shrinkage

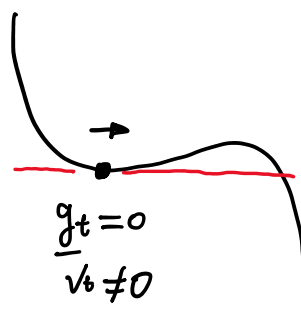
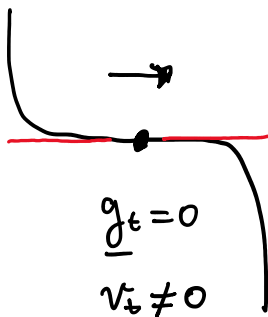
update: same as AdaGrad

Exponentially decaying moving average filter

act as momentum

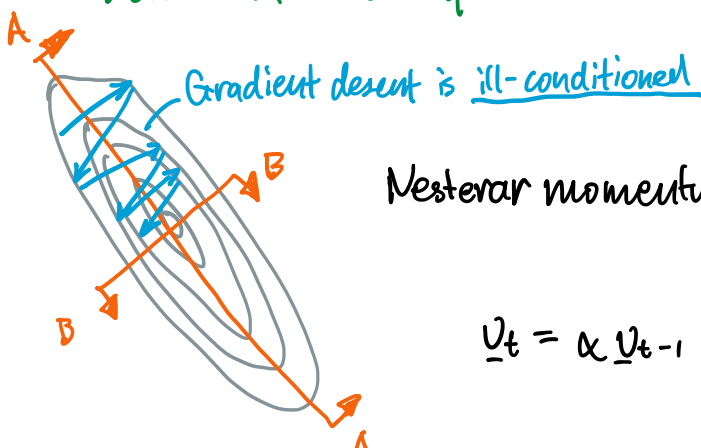
"Momentum" :  $\underline{v}_t = \alpha \underline{v}_{t-1} - \eta_t \underline{g}_t$

update:  $\underline{w}(t+1) = \underline{w}(t) + \underline{v}_t$



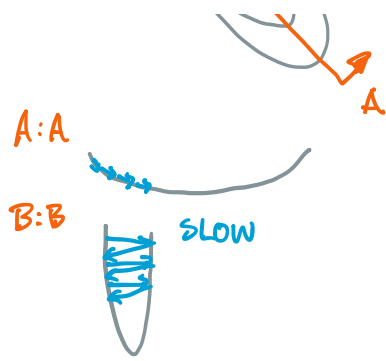
Momentum helps solving non-convex surfaces,

Momentum also helps:



Nesterov momentum:

$$\underline{v}_t = \alpha \underline{v}_{t-1} + \eta_t \cdot \nabla_{\underline{w}(t)} (\underline{w}(t) + \alpha \underline{v}_{t-1})$$



$\underline{w}(t+1) = \underline{w}(t) + \underline{v}(t)$   
 ↑  
 SGD with minibatch size 1  
 update  $\underline{w}(t+1) = \underline{w}(t) + \underline{v}(t)$