

In Supervised ML/DL, we have a training set of form:

$$D = \{z_n = (x_n, y_n)\}_{n=1}^N : \text{Data}$$

where  $z_n = (x_n, y_n) \sim p_{X \times Y}$  and  $p_{X \times Y}$  is an unknown underlying distribution.

We fix a Hypothesis class,  $H$ , to choose our model from, e.g., linear models, neural networks, etc.

We then define the Empirical Risk as:

$$\hat{R}(h) = R_D(h) = \frac{1}{N} \sum_{n=1}^N \ell(h(z_n)) : \text{Loss/Objective}$$

Empirical Risk Minimization (ERM) can be written as:

$$h_D^* = \underset{h \in H}{\operatorname{argmin}} R_D(h)$$

### Stochastic Gradient Descent (SGD)

Our goal is to minimize the empirical risk:

$$L(w) = \frac{1}{N} \sum_{n=1}^N \ell(w, z_n)$$

We know that in GD, our updates are:

$$w^{(t+1)} = w^{(t)} - \epsilon \nabla_w L(w^{(t)})$$

In Stochastic GD, we use:

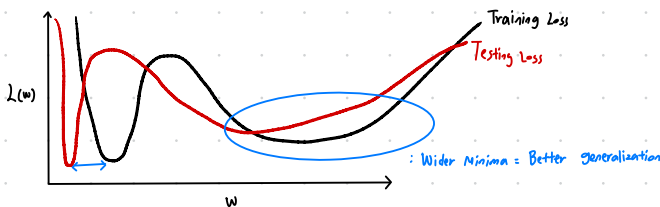
$$w^{(t+1)} = w^{(t)} - \epsilon \nabla_w \hat{L}_k(w^{(t)}), \text{ where } \hat{L}_k(w) = \frac{1}{b} \sum_{n \in \Omega_k} \ell(w, z_n)$$

Batch  
Batch size

What are the benefits of SGD?

1. Faster optimization
2. Better generalization

↳ Why? : 'Wide Minima' Phenomenon. "The wider the minimum, the better the performance on the test set"



## Finding Wider Minima

We can think of each gradient calculated from SGD as a noisy version of the actual GD:

$$U_k(w) = [\nabla L(w) - \nabla \hat{L}_k(w)] = \frac{1}{b} \sum_{n \in \mathcal{B}_k} [\nabla L(w) - \nabla L(w, z_n)]$$

Stochastic Gradient + Noise  
: Difference between  $\nabla L(w)$  and  $\nabla \hat{L}_k(w)$

Zero-mean and i.i.d. RV

CLT states that the sum (or average) of large number of i.i.d. random variables will be approximately normally distributed.

If we assume that  $U_k(w)$  has finite variance then, from

Central Limit Theorem (CLT) we have that:

$$U_k(w) \sim N(0, \sigma^2 I)$$

Hence, with the Gaussian assumption we have:

$$W(t+1) = W(t) - \epsilon \nabla L(W(t)) + \sqrt{\epsilon} \sqrt{\sigma^2} Z(t)$$

Standard Gaussian random variable

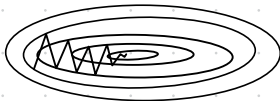
For small steps, we can write the discrete update above as a Stochastic Differential Equation (SDE):

$$dW_t = -\epsilon \nabla L(W_t) + \sqrt{\epsilon} \sigma dB_t$$

Brownian motion

What else can we do to get a flatter minima? : Momentum

## Momentum Update



Plain gradient update



With momentum

The momentum method maintains a running average of all gradients until the current step.

$$\Delta W^{(k)} = \beta \Delta W^{(k-1)} - \epsilon \nabla W_{\text{Loss}}(W^{(k-1)})^T$$

maintains a running average

- Typical  $\beta$  value is 0.9

Ex)

$$\Delta W^{(1)} = \beta \Delta W^{(0)} - \epsilon \nabla W_L(W^{(0)})^T$$

$$\Delta W^{(2)} = \beta \Delta W^{(1)} - \epsilon \nabla W_L(W^{(1)})^T$$

$$= \beta (\beta \Delta W^{(0)} - \epsilon \nabla W_L(W^{(0)})^T) - \epsilon \nabla W_L(W^{(1)})^T$$

$$= \beta^2 \Delta W^{(0)} - \beta \epsilon \nabla W_L(W^{(0)})^T - \epsilon \nabla W_L(W^{(1)})^T$$

$$\Delta W^{(3)} = \beta \Delta W^{(2)} - \epsilon \nabla W_L(W^{(2)})^T$$

$$= \beta (\beta^2 \Delta W^{(0)} - \beta \epsilon \nabla W_L(W^{(0)})^T - \epsilon \nabla W_L(W^{(1)})^T) - \epsilon \nabla W_L(W^{(2)})^T$$

$$= \beta^3 \Delta W^{(0)} - \beta^2 \epsilon \nabla W_L(W^{(0)})^T - \beta \epsilon \nabla W_L(W^{(1)})^T - \epsilon \nabla W_L(W^{(2)})^T$$

$$= \beta^3 \Delta W^{(0)} - \beta^2 \epsilon \nabla W_L(W^{(0)})^T - \beta \epsilon \nabla W_L(W^{(1)})^T - \epsilon \nabla W_L(W^{(2)})^T$$

And usually  $\beta = 0.9 < 1$

$$\text{Vanilla GD} : W(t+1) = W(t) - \epsilon \nabla W_L(W(t))$$

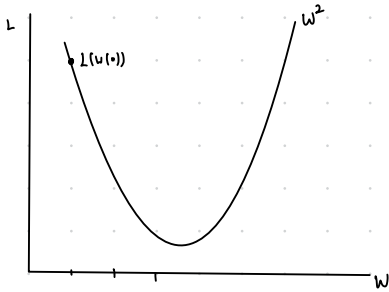
$$W(1) = W(0) - \epsilon \nabla W_L(W(0))$$

$$W(2) = W(1) - \epsilon \nabla W_L(W(1))$$

$$= W(0) - \epsilon \nabla W_L(W(0)) - \epsilon \nabla W_L(W(1))$$

$$W(3) = W(2) - \epsilon \nabla W_L(W(2))$$

$$= W(0) - \epsilon \nabla W_L(W(0)) - \epsilon \nabla W_L(W(1)) - \epsilon \nabla W_L(W(2))$$



$$L + L(w(1)) = L(-4) = 16, \beta = 0.9, \epsilon = 0.01$$

Momentum

$$\Delta W^{(k)} = \beta \Delta W^{(k-1)} - \epsilon \nabla_w \text{Loss}(w^{(k-1)})^T$$

↳ maintains a running average

$$\Delta W(0) = -4 \rightarrow$$

$$\Delta W(1) = \beta \Delta W(0) - \epsilon \nabla_w L(w(0)) = -4\beta + 8\epsilon = -3.52$$

$$\Delta W(2) = \beta \Delta W(1) - \epsilon \nabla_w L(w(1)) = -3.088$$

$$\Delta W(3) = -2.69$$

$$\Delta W(4) = -2.03$$

$$\Delta W(5) = -1.75$$

⋮

$$\Delta W(18) = -0.0005$$

Vanilla GD

$$w(0) = -4$$

$$w(1) = w(0) - \epsilon \nabla_w L(w(0)) = -4 + 8\epsilon = -3.992$$

$$w(2) = w(1) - \epsilon \nabla_w L(w(1)) = -3.984$$

⋮

$$w(500) = w(499) - \epsilon \nabla_w L(w(499)) = -0.019$$

GD w/ Momentum converges much faster than Vanilla GD.

However, since Momentum has running average aspect it sometimes overshoots.

The running average step:

- Get longer in directions where gradient retains the same sign
- Become shorter in directions where the sign keeps flipping

At any iteration, to compute the current step:

- First compute the gradient step at the current location
- Then add in the scaled previous step. (which is actually a running average)

**Nesterov's Accelerated Gradient**: Momentum: Gradient step  $\rightarrow$  extend previous step  
Nesterov: extend previous step  $\rightarrow$  then compute gradient

It changes the order of operations.

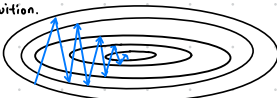
At any iteration, to compute the current step:

- First extend the previous step
- Then compute the gradient step at the resultant position
- Add the two to obtain the final step.

Nesterov's Accelerated Gradient converges much faster.

Smoothing the trajectory: Normalizing steps by second moment.

Intuition.



: Huge oscillations in the vertical direction even while it's trying to make progress in the horizontal direction

Here, we need optimization that goes vertically slow and horizontally fast.

Two popular methods that embody this principle: Adam and RMSProp

**RMSProp**: Updates are by parameters

Notation

- Derivative of loss w.r.t any individual parameter  $w$  is shown  $\partial_w D$
- The squared derivative is  $\partial_w^2 D = (\partial_w D)^2$  Not the second derivative
- The mean squared derivative is running estimate of the average squared derivative.  $\mathbb{E}[\partial_w^2 D]$

Modified Update rule:

- Scale down updates with large mean squared derivatives
- Scale up updates with small mean squared derivatives

RMSProp is a variant on the basic mini batch SGD algorithm.

Procedure:

- Maintain a running estimate of the mean squared value of derivatives for each parameter
- Scale update of the parameter by the inverse of the root mean squared derivative

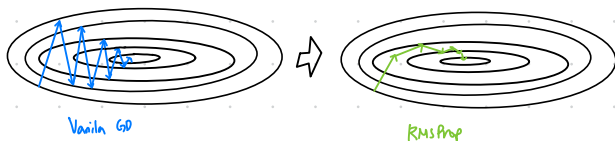
$\mathbb{E}[\partial_w^2 D]_k = \gamma \mathbb{E}[\partial_w^2 D]_{k-1} + (1-\gamma)(\partial_w^2 D)_k$  : running estimate of the mean squared value of derivatives for each parameter

$$w_{k+1} = w_k - \frac{\eta}{\sqrt{\mathbb{E}[\partial_w^2 D]_k + \epsilon}} \partial_w D$$

$\swarrow$  Preventing  $\sqrt{\mathbb{E}[\partial_w^2 D]_k}$  not goes to 0.

where  $\eta$ : learning rate

Hence, if  $\sqrt{\mathbb{E}[\partial_w^2 D]_k + \epsilon}$  : running average is small, scale  $\uparrow$  : fastening  
 large, scale  $\downarrow$  : slowing down



Adam : RMSProp + Momentum

: RMSProp only considers a second-moment normalized version of the current gradient

ADAM utilizes a smooth version of the momentum-augmented gradient

↳ Considers both first and second moment

Procedure:

- Maintain a running estimate of the mean derivative for each parameter
- Maintain a running estimate of the mean squared value of derivatives for each parameter
- Scale update of the parameter by the inverse of root mean squared derivative

$$m_k = \beta m_{k-1} + (1-\beta) (\partial_w D)_k$$

$$V_k = \gamma V_{k-1} + (1-\gamma) (\partial_w^2 D)_k$$

$$\hat{m}_k = \frac{m_k}{1-\beta^k}, \quad \hat{V}_k = \frac{V_k}{1-\gamma^k}$$

$$w_{k+1} = w_k - \frac{\eta}{\sqrt{\hat{V}_k + \epsilon}} \hat{m}_k$$

Typical Parameters:

- RMSProp:  $\beta=0.001$ ,  $\gamma=0.9$
- ADAM:  $\beta=0.001$ ,  $\beta=0.9$ ,  $\gamma=0.999$