

BLG 561 E FALL 2021

Deep Learning

26.10.2021

Görde ÜNAL

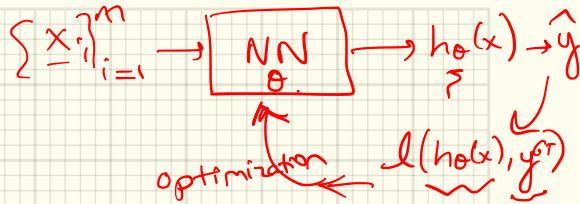
All ML algorithms : 3 components

1) Hypothesis Class: ✓

2) Loss function : ✓

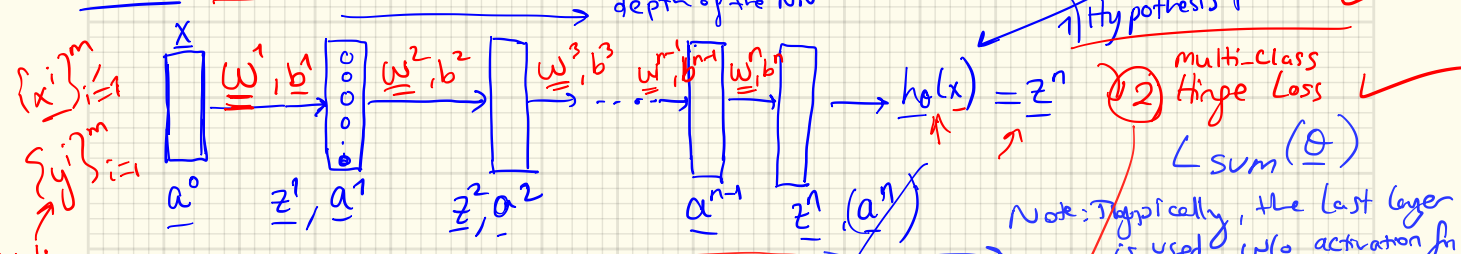
3) Optimization ✓

NN ←



Neural Networks (Fully-Connected, ANNs, Multi-Layered Perception (MLP))

Ex: Multi-Class Classification : n -layer ANN. → ① ✓



② multi-class Hinge Loss $L_{\text{sum}}(\theta)$
 Note: Typically, the last layer is used w/o activation fn.

$\theta = \{ \underline{w}^1, b^1, \underline{w}^2, b^2, \dots, \underline{w}^n, b^n \}$

$$L_{\text{sum}} = \sum_{i=1}^m l(h_{\theta}(x^i), y^{(i)}) = \sum_{i=1}^m \sum_j \max(0, \underbrace{h_{\theta}(x^i)_j}_{\text{ifn argument of vector fn } h} - h_{\theta}(x^i)_j y_i + \frac{1}{\uparrow})$$

③ Optimization / Training the NN: estimating parameters θ !

Today

3) How do we optimize the loss? $\theta^* = \arg \min_{\theta} \left(\sum_{i=1}^m \ell(h_{\theta}(x^i), y^i) \right)$

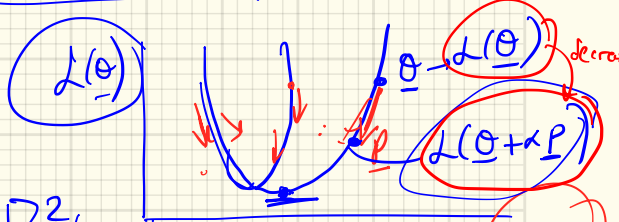
Q: Is this loss fn. convex? Not convex

$$h_{\theta}(x) = \underline{w}^n \left(\dots \underline{w}^3 f_2 \left(\underline{w}^2 f_1 \left(\underline{w}^1 x + \underline{b}^1 \right) + \underline{b}^2 \right) + \underline{b}^3 + \dots \right) + \underline{b}^n$$

Before $h_{\theta}(x) = \theta^T x$ linear hypothesis \rightarrow - - -

Deep learning models: we work w/ non-convex optimization problem

$L(\theta)$: overall loss function - P : search direction



T.S. expansion:

$$L(\theta + \alpha P) \approx L(\theta) + \underbrace{P \cdot \nabla_{\theta} L(\theta)} + \frac{1}{2} P^T \nabla_{\theta}^2 L(\theta) P$$

Let $P = -\alpha \nabla_{\theta} L(\theta)$ HOT

$$L(\theta + \alpha P) \approx L(\theta) - \alpha \underbrace{\|\nabla_{\theta} L(\theta)\|^2}_{\geq 0}$$

θ in millions billions

Recall: Newton method. Search dir: $\alpha \nabla_{\theta}^2 L$: Hessian

GRADIENT DESCENT Algorithm: We can use it whether $L(\theta)$ is convex or not!
(GD)

GD
Algorithm:

θ_0 : "initial" parameters \leftarrow
Iterate until "convergence" \leftarrow

$$\underline{\theta}^{k+1} \leftarrow \underline{\theta}^k - \alpha \cdot \nabla_{\theta} L(\underline{\theta})$$

Learning Rate (LR) : step size. } one of the important hyperparameters

eg. Newton method;
uses other search dir.

$$\underline{\theta} \leftarrow \underline{\theta} - \underbrace{(\nabla_{\theta}^2 L(\theta))^{-1}}_P \cdot \nabla_{\theta} L(\theta)$$

P : search direction.

Note: Deep NNs (Millions of parameters) $\theta_{100m \times 1}$ $\xrightarrow{\text{imagine computing}}$ $\nabla_{\theta}^2 L = \text{Hessian}!!$
 \hookrightarrow + Inverting !!

\exists Quasi-Newton methods:

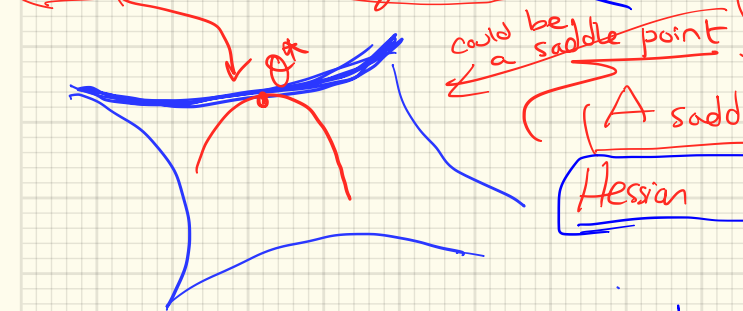
We approximate the Hessian thru gradients \rightarrow eg. BFGS.

Q. You have your $\mathcal{L}(\theta)$; is it enough to set $\nabla_{\theta} \mathcal{L}(\theta) = 0$ to obtain your minimizer of the loss function? _____

Thm: Necessary 1st order condition: If θ^* is a local minimizer & \mathcal{L} is continuously differentiable then $\nabla_{\theta} \mathcal{L}(\theta^*) = 0$

Any local minimizer $\theta^* \implies \nabla_{\theta} \mathcal{L}(\theta^*) = 0$ (Necessary Cond.)

θ^* is a stationary point $\longleftarrow \nabla_{\theta} \mathcal{L}(\theta^*) = 0$ (Sufficient Cond.)



(A saddle point is a stationary point.)

Hessian $\nabla_{\theta}^2 \mathcal{L}$ is not tr definite
 (its eigen values both > 0 and < 0)

Thm 2nd order Necessary Condition: Add $\nabla_{\theta}^2 \mathcal{L}$ exists, then $\nabla \mathcal{L}(\theta^*) = 0$
 Given if θ^* is a local minimizer $\implies \nabla^2 \mathcal{L}(\theta^*) \geq 0$

Thm: 2nd order Sufficient Conditions;

If $\nabla_{\theta}^2 L$ is continuous & positive definite & $\nabla_{\theta} L(\theta^*) = 0$
then θ^* is a local minimizer of L .

\therefore If $\nabla_{\theta} L = 0$ & $\nabla_{\theta}^2 L$ p.d. $\rightarrow \theta^*$ is a minimizer.

In Summary: In deep learning: We want $\nabla_{\theta} L(\theta)$ to vanish.
Search θ^* from $\nabla_{\theta} L(\theta) = 0$.

GD: Repeat (Loss: uses all samples)
$$L = \frac{1}{m} \sum_{i=1}^m l(h_{\theta}(x^i), y^i)$$

For $i=1:m$

$$\text{grad}^{(i)} \leftarrow \nabla_{\theta} l(h_{\theta}(x^i), y^i)$$

Update the parameters

$$\theta \leftarrow \theta - \alpha \frac{1}{m} \sum_{i=1}^m \text{grad}^{(i)}$$

until "convergence"

SGD: Stochastic GD

For $i=1:m$ ↑ introduces randomness.

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(h_{\theta}(x^i), y^i)$$

grad⁽ⁱ⁾ = loss for 1 sample

Noisy gradient update.

In practice : we use MiniBatch GD; a compromise btw GD & SGD.

⇒ SGD

Let B : batchsize
 ep. $B = \begin{matrix} 8 \\ 16 \\ 32 \\ 64 \\ 128 \end{matrix}$
 another hyperparameter you have to fix.

$$L_{MB} = \frac{1}{B} \sum_{i=1}^B l(h_{\theta}(x^i), y^i)$$

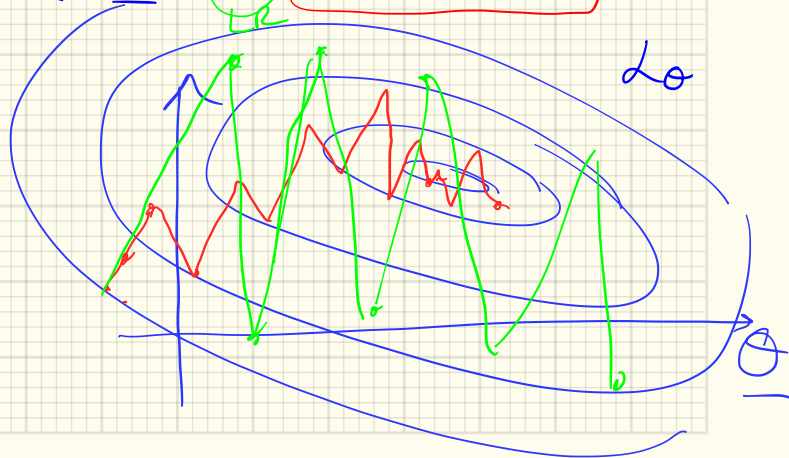
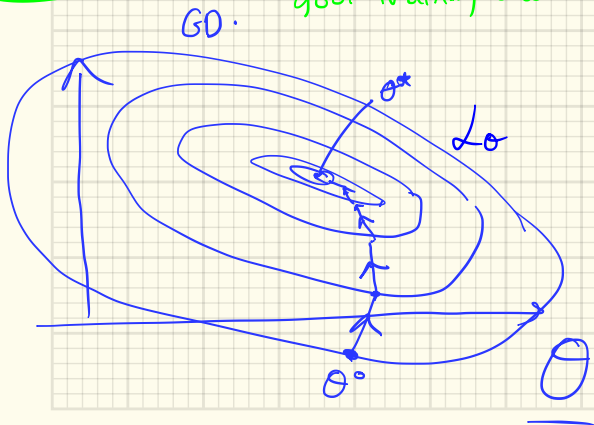
Note:
 Vanilla SGD uses $B=1$ not used in practice

SGD (minibatch): For $k=1 \dots \lfloor \frac{M}{B} \rfloor$ (# minibatches)

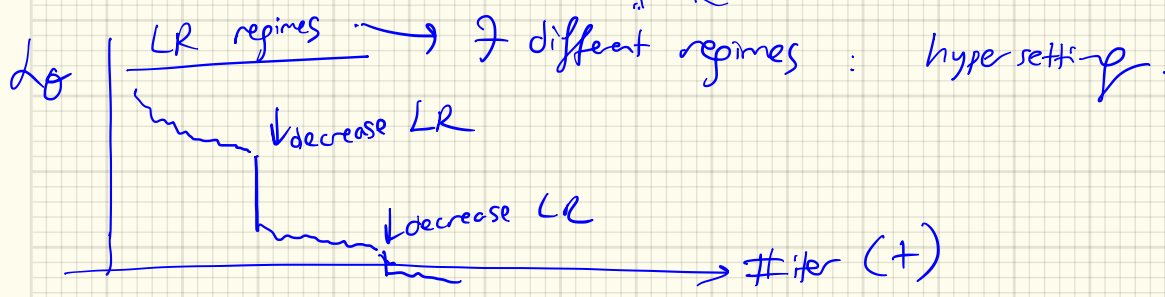
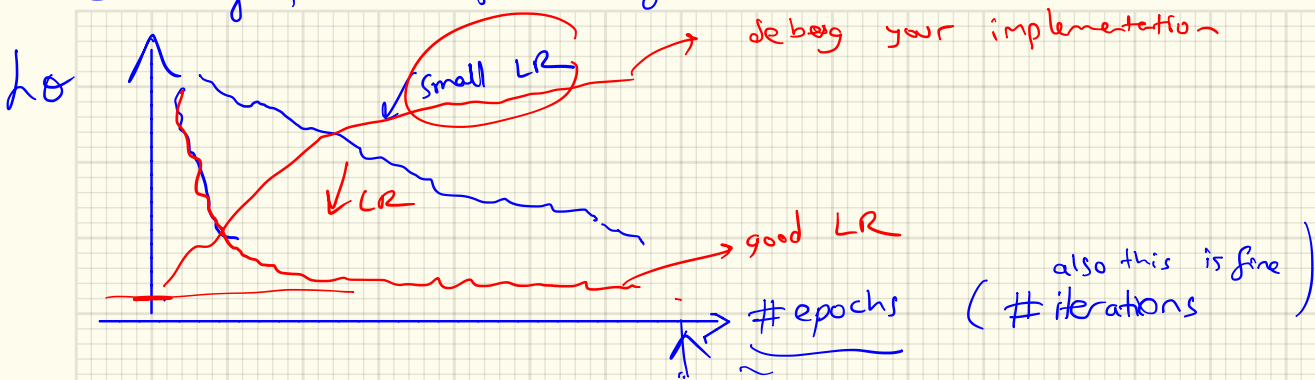
1 epoch: once you've seen all the samples in your training data

$$\underline{\theta} \leftarrow \underline{\theta} - \alpha \nabla_{\underline{\theta}} L_{MB}(\underline{\theta})$$

∇_{θ^i}



Checking for convergence in your optimization

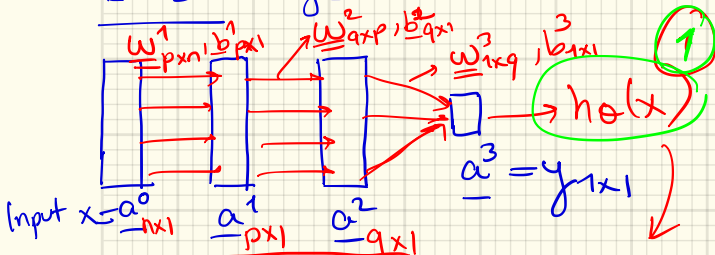


→ Always check your loss curves.

Q: How do we compute the gradient: $\nabla_{\theta} \ell(\text{hol}x, y)$? Backpropagation Optimization for NN

Q: set of parameters of the NN: $\underline{w}^1, \underline{b}^1, \underline{w}^2, \underline{b}^2, \dots, \underline{w}^k, \underline{b}^k$.

$k=3$ - layer NN (ANN = FCN = MLP)



eg. LS-loss $\ell(a^3, y) = \frac{1}{2} (y - a^3)^2$

Loss $\ell(\text{hol}x)$

$\theta = (\underline{w}^1, \underline{b}^1, \underline{w}^2, \underline{b}^2, \underline{w}^3, \underline{b}^3)$

Forward Pass

Compute $\underline{a}^1, \underline{a}^2, \underline{a}^3, \dots, \underline{a}^k$; $\ell(\underline{a}^k, y)$

$$\underline{a}^1 = f_1(\underline{w}^1 \underline{a}^0 + \underline{b}^1)$$

$$\underline{a}^2 = f_2(\underline{w}^2 \underline{a}^1 + \underline{b}^2)$$

$$\underline{a}^3 = f_3(\underline{w}^3 \underline{a}^2 + \underline{b}^3)$$

z^3

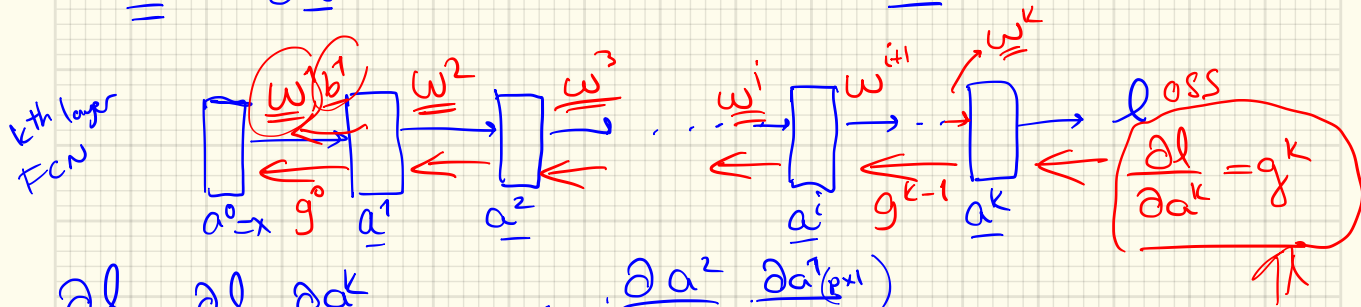
Also $z^i = \underline{w}^i \underline{a}^{i-1} + \underline{b}^i$

1 Hypothesis function for each layer in ANN (FCN)

3) Optimization (BACKPROP) $\left. \frac{\partial l}{\partial \underline{w}^i}, \frac{\partial l}{\partial \underline{b}^i} \right\} \equiv \frac{\partial l}{\partial \underline{\theta}} = \underline{\delta}$

Backward Pass: Based on the chain rule:

$$\frac{\partial l}{\partial \underline{w}^1} = \frac{\partial l}{\partial a^k} \frac{\partial a^k}{\partial a^{k-1}} \dots \frac{\partial a^3}{\partial a^2} \frac{\partial a^2}{\partial a^1} \frac{\partial a^1}{\partial \underline{w}^1} \leftarrow$$



$$\frac{\partial l}{\partial \underline{b}^1} = \frac{\partial l}{\partial a^k} \frac{\partial a^k}{\partial a^{k-1}} \dots \frac{\partial a^2}{\partial a^1} \frac{\partial a^1}{\partial b^1(\text{ex1})}$$

$$\rightarrow \frac{\partial a^i}{\partial b^i} = \frac{\partial f}{\partial z^i} \cdot \frac{\partial z^i}{\partial b^i} = \text{diag } f'(z^i) \cdot \underline{1} \rightarrow \text{vector of 1's}$$

diag matrix

$$\rightarrow \frac{\partial a^i}{\partial a^{i-1}} = \frac{\partial f}{\partial z^i} \cdot \frac{\partial z^i}{\partial a^{i-1}} = \text{diag } f'(z^i) \cdot \underline{w}^i \leftarrow$$

$$\frac{\partial a^i}{\partial \underline{w}^i} = \frac{\partial f}{\partial z^i} \cdot \frac{\partial z^i}{\partial \text{vec}(\underline{w}^i)} \rightarrow$$

$$\rightarrow \underline{z}^i = \underline{\underline{I}} \underline{\underline{W}}^i \underline{a}^{i-1} + \underline{b}^i$$

$$\underline{z}^i = (\underline{a}^{i-1})^T \otimes \underline{\underline{I}} \text{vec}(\underline{\underline{W}}^i) + \underline{b}^i$$

$$\frac{\partial \underline{z}^i}{\partial \text{vec}(\underline{\underline{W}}^i)} = (\underline{a}^{i-1})^T \otimes \underline{\underline{I}} \quad \leftarrow$$

vector

$$\Rightarrow \underline{g}^i \triangleq \frac{\partial l}{\partial \underline{a}^i}$$

→ upstream gradients
i=1, ..., k

Start w/ this $\underline{g}^k, \underline{g}^{k-1}, \dots, \underline{g}^1$: want a recursive computation:

$$\underline{g}^{i-1} \triangleq \frac{\partial l}{\partial \underline{a}^{i-1}} = \frac{\partial l}{\partial a^k} \cdot \frac{\partial a^k}{\partial a^{k-1}} \cdots \frac{\partial a^{i+1}}{\partial a^i} \cdot \frac{\partial a^i}{\partial \underline{a}^{i-1}} = (\underline{\underline{W}}^i)^T \text{diag}(f'_{(z)} - \underline{g}^i$$

→ Target derivatives: $\underline{\underline{V}}_{\underline{\underline{W}}^i} l, \underline{V}_{\underline{b}} l, i=1, \dots, k$

We are using this identity

$$\text{vec}(\underline{\underline{A}} \underline{\underline{B}} \underline{\underline{C}}) = (\underline{\underline{C}}^T \otimes \underline{\underline{A}}) \text{vec} \underline{\underline{B}}$$

Recall \otimes Kronecker product

$$\underline{\underline{C}} = (\underline{\underline{A}}_{m \times n} \otimes \underline{\underline{B}}_{p \times q}) \in \mathbb{R}^{mp \times nq}$$

$$= \begin{bmatrix} a_{11} \underline{\underline{B}} & a_{12} \underline{\underline{B}} & \dots & a_{1n} \underline{\underline{B}} \\ a_{21} \underline{\underline{B}} & & & \\ \vdots & & & \\ a_{m1} \underline{\underline{B}} & & & a_{mn} \underline{\underline{B}} \end{bmatrix}$$

$$\nabla_{\underline{\underline{w}}^i} l = \underbrace{g^{i-1}}_{\text{upstream gradient}} \cdot \underbrace{\text{diag } f_i' (a^{i-1})^T \otimes \underline{\underline{I}}}_{\text{local gradient: } \frac{\partial a^i}{\partial \underline{\underline{w}}^i}} \rightarrow \text{vec}(\underline{\underline{w}}^i)$$

Summarize :

1) Forward Pass : $\underline{\underline{a}}^i = f_i \left(\underbrace{\underline{\underline{w}}^i \underline{\underline{a}}^{i-1} + \underline{\underline{b}}^i}_{\underline{\underline{z}}^i} \right)$



Backpropagation Algorithm: Parameters we search for for a k -layer NN $\theta = \{ \underline{w}^1, b^1, \dots, \underline{w}^k, b^k \}$

1) Forward Pass: Compute $\underline{a}^1, \dots, \underline{a}^k, \ell(\underline{a}^k, y)$ ←

Also, compute local gradients: $\frac{\partial a^i}{\partial a^{i-1}}, \frac{\partial a^i}{\partial b^i}, \frac{\partial a^i}{\partial \underline{w}^i}$ ✓

2) Backward Pass: Compute the upstream Gradients.

Compute recursively: g^k, \dots, g^1 ($g^k = \frac{\partial \ell}{\partial a^k}$) (g^{k-1}) $\rightarrow g^{k-2}$

$$\underline{g}^i = (\underline{w}^i)^T \text{diag } f'(z^i) \cdot g^{i-1}$$

3) Compute the gradients of the loss fn. w.r.t. θ :

$$\nabla_{\underline{w}^i} \ell \approx \frac{\partial \ell}{\partial (\text{vec } \underline{w}^i)} = \left[\underline{g}^i \text{diag } f'(z^i) (\underline{a}^{i-1})^T \otimes \underline{I} \right]$$

upstream grad
local grad.
(convert back to a matrix)

$$\nabla_{b^i} \ell = \frac{\partial \ell}{\partial b^i} = g^i \text{diag } f'(z^i) \underline{1}_{p \times p} \rightarrow \text{vector.}$$

Note:
Forward Pass uses multiply by \underline{w}^i

Backward Pass uses multiply by $\underline{w}^{i \top}$

4) Update the Parameters:

$$\underline{w}^{i(t)} \leftarrow \underline{w}^{i(t-1)} - \alpha \nabla_{\underline{w}^i} \ell^{(t-1)}$$

$$\underline{b}^{i(t)} \leftarrow \underline{b}^{i(t-1)} - \alpha \nabla_{b^i} \ell^{(t-1)}$$

} for each layer i 's parameters
 $\forall i = 1, \dots, k$ layers.
 for a k -layer NN.

(Variants of GD) SGD w/ Momentum:

Idea: Smooth out your gradients to speed up GD.

At iteration t : Compute $\nabla_{\mathbf{w}} \mathcal{L} \rightarrow \underbrace{\nabla_{\mathbf{w}} \mathcal{L}}_{d\mathbf{w}}$, $\underbrace{\nabla_{\mathbf{b}} \mathcal{L}}_{d\mathbf{b}}$ on the current minibatch

Add Momentum: $dV_{\mathbf{w}} = 0$
 $dV_{\mathbf{b}} = 0$

GD w/ momentum algo.

$$\begin{aligned} dV_{\mathbf{w}}^t &= \beta_1 \cdot \underbrace{dV_{\mathbf{w}}^{t-1}}_{\text{previous gradients}} + (1 - \beta_1) d\mathbf{w}^t \\ dV_{\mathbf{b}}^t &= \beta_1 \cdot \underbrace{dV_{\mathbf{b}}^{t-1}}_{\text{previous gradients}} + (1 - \beta_1) d\mathbf{b}^t \end{aligned}$$

$d\mathbf{w}^t, d\mathbf{b}^t$

current gradients at t

(β_1 : a new hyperparameter)

eg. $\beta_1 = 0.9$

Now, update the parameters:

$$\mathbf{w}^t = \mathbf{w}^{t-1} - \alpha \cdot dV_{\mathbf{w}}^t$$

$$\mathbf{b}^t = \mathbf{b}^{t-1} - \alpha dV_{\mathbf{b}}^t$$

Note: A bias correction is added to affect initial iterations

divide by $\frac{dV^t}{1 - \beta_1^t}$ → bias term

$$t=1 \rightarrow \frac{1}{1 - (0.9)^1} \Rightarrow 1$$

:

$$t=10 \Rightarrow \frac{1}{1 - (0.9)^{10}} \approx 1$$

After a few iterations bias term has no effect

RMS_prop GD: At iteration t : Want to compute \underline{dw} , \underline{db} :

Compute:

$$\underline{sdw} = \beta_2 \underline{sdw} + (1 - \beta_2) (dw)^2 \leftarrow \text{current gradients}$$

$$\underline{sdb} = \beta_2 \underline{sdb} + (1 - \beta_2) (db)^2$$

Update Rules:

$$\underline{w} \leftarrow \underline{w} - \alpha \frac{\underline{dw}}{\sqrt{\underline{sdw} + \epsilon}} \leftarrow \text{like normalizing by std. dev. of the gradient.}$$

$$b \leftarrow b - \alpha \frac{db}{\sqrt{\underline{sdb} + \epsilon}}$$

Adam optimizer: Combination of Momentum & RMS prop:

Note: Use bias-corrected versions

Update

At iteration t :

$$\underline{w} \leftarrow \underline{w} - \alpha \frac{dVw^{corr}}{\sqrt{\underline{sdw}^{corr} + \epsilon}}$$

$$b \leftarrow b - \alpha \frac{dVb^{corr}}{\sqrt{\underline{sdb}^{corr} + \epsilon}}$$

$$dVw^{corr} = \frac{dVw}{(1 - \beta_1^t)}$$

Initialize: $dVw = 0$
 $dVb = 0$
 $ds_w = 0$
 $ds_b = 0$

At each iteration t : Compute $\underline{dw}, \underline{db}$ w/ current minibatch
current gradients
 dVw, dVb
 ds_w, ds_b