COL865
Deep Learning

Optimization :-

Algorithms with adaptive learning Rate :-

(Popular Algorithms) :-

① Adagrad :-

$$g^{(t+1)} = \frac{1}{m} \nabla_\theta J(\theta)$$

$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} L(f(x^{(i)}; \theta), y^{(i)})$

$$r^{(t+1)} = r^{(t)} + g \odot g^{(t+1)} \Rightarrow \quad \Leftarrow \text{ accumulation of second moment (squared gradient)}$$

$$\Delta\theta = - \frac{\eta}{\epsilon [s + \sqrt{r}]} \odot g^{(t+1)} \qquad \Rightarrow \text{ pointwise}$$

scale $\eta$ by $\frac{1}{[s + \sqrt{r}]}$   small parameter for stabilization  $\rightarrow (10^{-7})$ ??   scales $g_j$ by as :-

$$\boxed{\theta^{(t+1)} = \theta^{(t)} + \Delta\theta}$$

$$\boxed{\Delta\theta_j' = - \frac{\eta \; g_j^{(t)}}{s + \sqrt{r_j^{(t)}}}}$$

$$r_j^{(t)} = \sum_{t=1}^{t} [g_j^{(t)}]^2$$

{ parameters ~~Directions~~ with large partial derivative (accumulated) have a large decrease in learning rate.

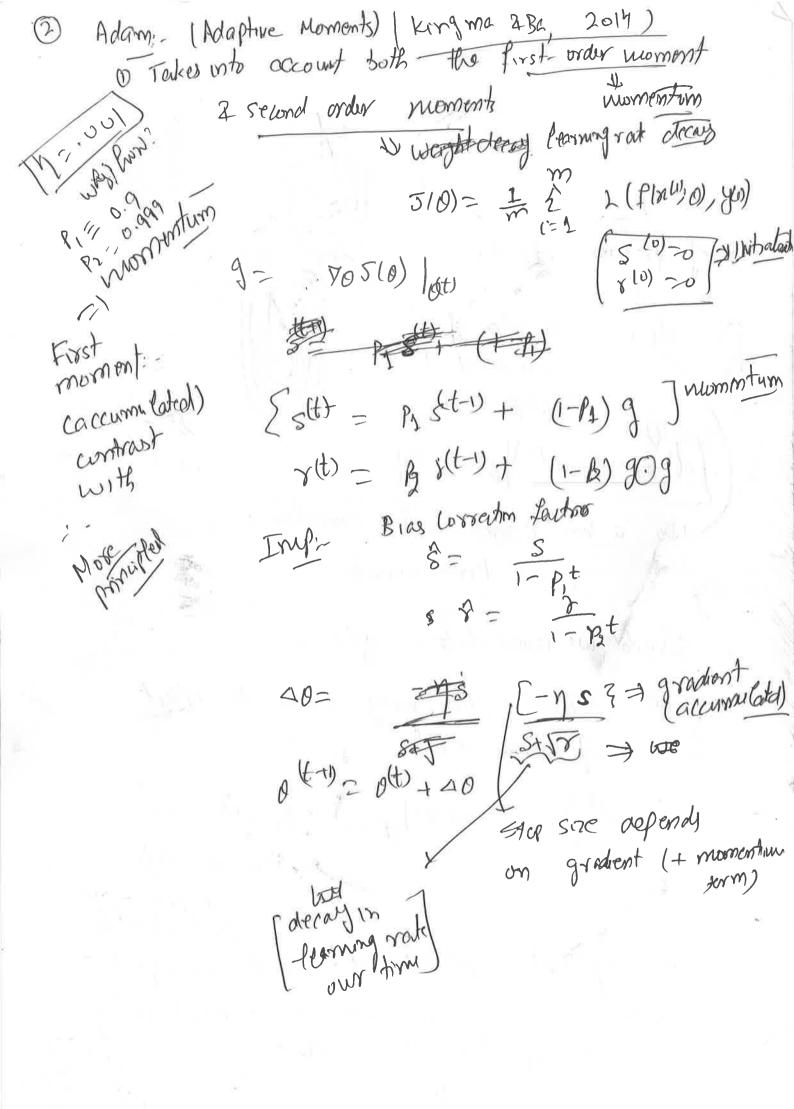parameters with small partial derivative (accumulated) have a slower decrease in learning rate.

↳ θ Amortized (happens over time).



Issue :- Premature & excessive decrease in the learning rate.

Stuck in local minima.

2. RMSProp:

$r=0$ in the beginning

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} L(f(x^{(i)};\theta), y^{(i)})$$

$$g^{(t)} = \nabla_\theta J(\theta)\big|_{\theta^{(t-1)}}$$

$$r^{(t)} \leftarrow \rho\, r^{(t-1)} + (1-\rho)\, g^{(t)} \odot g^{(t)} \quad\Big|\quad \text{as opposed to: } r^{(t)} = r^{(t-1)} + g^{(t)} \odot g^{(t)}$$

↙ Past effect decays exponentially.

↓ slower decrease in learning rate

$$\Delta\theta = \;\; \boxed{\phantom{x}} \;\; -\frac{\eta}{\delta + \sqrt{r}} \cdot g$$

$$\theta^{(t+1)} = \theta^{(t)} + \Delta\theta$$

Adagrad vs RMSProp — Better chances of escaping local minima

↖ can combine with (Nesterov) momentum

↓ Additional momentum

$v^{(t)} \quad \dot{v}^{(t)}$

$$\boxed{v^{(t+1)} = \alpha\, v^{(t)} - \frac{\eta}{\delta + \sqrt{r^{(t)}}} \odot g^{(t)}}$$

$$\boxed{\theta^{(t+1)} = \theta^{(t)} + v^{(t)}}$$

3.

② Adam:- (Adaptive Moments) ( Kingma & Ba, 2014 )

① Takes into account both ~~the~~ first order moment
↓
2 second order moments  momentum
↓
~~weight decay~~ learning rate decay

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} L(f(x^{(i)}; \theta), y_{(i)})$$

$\eta = .001$
why how?

$P_1 = 0.9$
$P_2 = 0.999$
momentum
⇒)

$\begin{pmatrix} s^{(0)} = 0 \\ r^{(0)} = 0 \end{pmatrix}$ ⇒ Initialized

$$g = \nabla_\theta J(\theta) \big|_{\theta^{(t)}}$$

First moment:-
(accumulated)
contrast with

$$\cancel{s^{(t)} = P_1 \frac{s^{(t-1)}}{} + (1-P_1)}$$

$$\left\{ s^{(t)} = P_1 \, s^{(t-1)} + (1-P_1) \, g \quad \right] \text{momentum}$$

$$r^{(t)} = P_2 \, s^{(t-1)} + (1-P_2) \, g \odot g$$

More principled

Imp:- Bias correction factors
$$\hat{s} = \frac{s}{1 - P_1^t}$$

$$s \quad \hat{r} = \frac{r}{1 - P_2^t}$$

$$\Delta \theta = \frac{\cancel{-\eta \hat{s}}}{s + \sqrt{r}} \quad \left[ -\eta s \right\} \Rightarrow \left\{ \begin{array}{c} \text{gradient} \\ \text{accumulated} \end{array} \right.$$

$$\frac{}{s + \sqrt{r}} \Rightarrow \text{we}$$

$$\theta^{(t+1)} = \theta^{(t)} + \Delta\theta$$

Step size depends on gradient (+ momentum term)

$\left[ \begin{array}{c} \text{WD} \\ \text{decay in} \\ \text{learning rate} \\ \text{over time} \end{array} \right]$

# Conjugate Gradient Method:-

Each successive iteration
might undo the effect

of previous iteration.

$d_{t-1}^{(t-1)}$:- Direction of
movement at
step $(t-1)$

$$d_{(t-1)}^{(t-1)} \cdot \nabla_\theta J(\theta)\big|_{\theta^{(t)}} = 0$$

→ Do a line search in this direction
Find minima.

Should we move $d_{dt}^{(t)} = g^{(t)}$ ?

∴ But $d^{(t)}$ may undo the effect
of previous iteration.

$$d^{(t)} = \nabla_\theta J(\theta)\big|_{\theta^{(t)}} + \beta_t \, d^{(t-1)}$$

⇓ Some left
over from previous
search direction

we say that $d^{(t-1)}$ & $d^{(t)}$ are conjugate

given G    so ≠ the semi-definite matrix

H if $\boxed{d^{(t-1)^T} \, H \, d^{(t)} = 0}$    → useful for
solving a set
of linear equations

○ Progressively move towards the solution

For us :- if $H$ is constant (function is quadratic)

$d^{(t)} \; d^{(1)} \; d^{(2)} \; -- \; d^{(t)}$ be obtained such that

$$d^{(t-1)T} H \; d^{(t)} = 0$$

Does not undo the effect of previous iterations

then :- gradient along $d^{(t-1)}$

if gradient along $d^{(t-1)} = 0$ at $d^{(t-1)}$

then gradient along $d^{(t)}$ is also

$d^{(t)}$ can be found analytically

$d^{(t-1)}$ is also 0 at $d^{(t)}$.

But ∤H∤

Further :- Conjugacy is preserved across

Use information from the Hessian

$$\left[ \; d^{(t')} H \; d^{(t)} = 0 \qquad \forall \; t' \leq < t \right]$$

if $d^{(t-1)} H d^{(t)} = 0 \qquad \forall t$

$d^{(1)}, d^{(2)} \; -- \; d^{(n)}$ will span the entire space.

Further :-

$$d^{(t)} = \nabla_\theta J(\theta^{(t)}) \frac{\nabla_\theta J(\theta^{(t)})}{\nabla_\theta J(\theta^{(t-1)})^T} \nabla_\theta J(\theta^{(t-1)})$$

⟹ After $n$ steps, we will reach the minimum following conjugate directions

[ Scaled conjugate gradient method ]

Non-linear conjugate gradient method :- 

Reset $\beta^{(t)} = 0$ occasionally

BFGS:- Newton's method:-

$$\theta^{(t+1)} = \theta^{(t)} - H^{-1}g$$

↓

Approximate

by another matrix $M(t)$ which successively becomes __better__.

$\begin{bmatrix} Boyden \\ Fletcher \\ Goldfarb \\ Shanno \end{bmatrix}$

$L\overline{BFGS}$

Do not store the matrix (approximation) explicitly in memory.

__Batch normalization:__ - Do it later

Newton's method:- escaping saddle points.



~~Do~~ Scale down the grad movement in directions in which gradient moves upwards.

$$\theta^{(t+1)} = \theta^{(t)} - \left[ \left[ H(f(\theta^{(t)})) + \alpha I \right]^{-1} \right] \nabla_\theta f(\theta^{(t)})$$

Attempt to make eigenvalues +ve.

④ Parameter Initialization Strategies :-

     ↳ ① Initializing $w's$ → close to zero! — Randomly (ranges)?

        ② Initializing biases $b's$

        ☞

$$x^{(t)} = W^T \quad \rightarrow \text{Activation function}$$

$$\boxed{x^{(t)} = h\left[W^{(t)} x^{(t-1)} + b^{(t)}\right]}$$

$$\boxed{\text{Read section } 8.4}$$