# Challenges on Neural Network Optimization:-

COL 865 . Deep Learning

Aug 28, 2017

## Key Problem:-

$\Rightarrow$ proxy for generalization error

$$\bar{J}(\theta) = \frac{1}{m} \sum_{i=1}^{m} \lambda \left[ f(x^{(i)}; \theta); y^{(i)} \right]$$

Expected loss over the training set

$\Downarrow$ over training examples

$$E\left[ \tilde{J}(\theta) = J(\theta) + \frac{1}{2} \theta^T \theta \right]$$

a. $\Lambda(\theta)$ $\Rightarrow$ regularizer

May not always

Not exactly an optimization problem

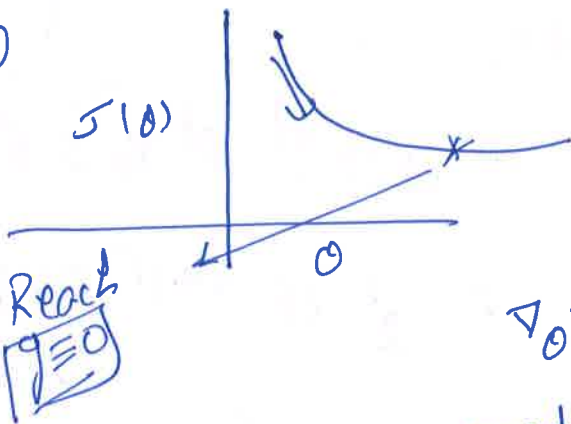$$J^*(\theta) = E_{(x,y) \sim p^*_{data}} \left[ L \left( f(x, \theta); y \right) \right]$$

$\Rightarrow$ gradient update rule

$$\theta^{(t+1)} \leftarrow \theta^t - \eta . \nabla_\theta J(\theta)$$

## Understand the ISSUES :-

How do we minimize a function ?

① Gradient Descent



①

$J(\theta)$

Reach $\boxed{g \equiv 0}$

Find the gradient

$J'(\theta) \equiv g$ &

$\nabla_\theta J(\theta) \equiv g$

& move opposite to direction of g.

Direction of steepest ascent

$\boxed{\begin{array}{l} ① \text{ Fixed /Static} \\ ② \text{ Dynamic} \end{array}}$

③ Should we always move along g?

Key Challenge:-

How much to move along −g?

Learning rate: η.

Alternate methods $\Rightarrow$ ① Adding momentum ② conjugate gradient

Second order Methods:-

$\Rightarrow J'(\theta)$

If $J(\theta)$ is quadratic, optimum in one step.

We want $\theta^*$ s.t $J'(\theta) = 0$

$\Rightarrow$ Newton's method to find zeros of a function-

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

$\longrightarrow$ computationally expensive

$$\boxed{H \equiv \text{Matrix of second order derivatives}}$$

Understanding $H \equiv$ curvature Becomes critical to any kind of optimization:- even gradient based.

---

* Illustration-

$\theta$ Using quadratic approximation:-

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

$$J(\theta) = J(\theta^{(0)}) + (\theta - \theta^{(0)})^T \underbrace{\nabla_\theta J(\theta^{(0)})}_{g}$$

$$+ \frac{1}{2}(\theta - \theta^{(0)})^T H (\theta - \theta^{(0)})$$

Let Now, $\boxed{\theta = \theta^{(0)} - \varepsilon \eta g}$ $\rightarrow$ using gradient descent

$$J(\theta^{(0)} - \eta g) = J(\theta^{(0)}) - \eta g^T \cdot g$$
$$+ \frac{1}{2} \eta g^T H g$$

$$J(\theta^{(\omega)} - \eta g) = J(\theta^{(\omega)}) - \eta g^T g + \frac{1}{2}\eta^2 g^T H g$$

At optimum values of

$$J(\theta^{(\omega)}) - \eta g) - J(\theta^{(\omega)}) \quad \text{is minimum as a}$$

function of $\eta \Rightarrow$

$$g^T g = \eta g^T H g$$

Note:- if $\theta \in \mathbb{R}$ then:-

$$\eta = \frac{1}{J''(\theta)} \Rightarrow \boxed{\theta^{(t+1)} = \theta^{(t)} - \frac{J'(\theta)}{J''(\theta)}}$$

Newton's update

$$\twoheadrightarrow \eta = \frac{g^T H g}{g^T a}$$

$$\text{&} \quad \eta g^T g = \eta^2 g^T H g$$

$$\boxed{g = \alpha_i v_i} \quad \text{If } g \text{ aligns with an eigenvalue} \rightarrow v_i, \; g^T H g = \lambda_i \alpha_i^2 \; \text{#} \quad \boxed{\eta = \frac{g^T g}{g^T H g}} \longrightarrow \text{scales } g \text{ in the space of eigenvalues of}$$

Recall:- $\boxed{H = \theta^T \wedge \theta}$

$$\boxed{\eta = \frac{1}{\lambda_i}} \rightarrow \text{problem if } \frac{|\lambda_{max}|}{|\lambda_{max}|}$$

because

Eigenvalue decomposition

$$\boxed{g \equiv \alpha_1 v_1 + \alpha_2 v_2 + \cdots \alpha_n v_n}$$

$$\Rightarrow \left[\sum_i \alpha_i v_i\right] H \left[\sum_i \alpha_i v_i\right]$$

$$\left.\begin{array}{l}\left|\dfrac{\lambda_{max}}{\lambda_{min}}\right| \\ \text{then} \end{array}\right.$$ is very large, then $g$ may not be the optimal direction

$$\boxed{H v_i = \lambda_i v_i}$$

Conditioning number

Direction & second derivative?

$$g^T H g = \left[\sum_i \alpha_i v_i\right]\left[\sum_i \lambda_i \alpha_i v_i\right]$$

weighted sum of eigenvalues

$$d^T H d = d^T \theta \wedge \theta^T d$$

$$= (\theta d)^T \wedge (\theta^T d)$$

illustration:-



very steep



gradient may waste a lot of time in reaching the minima

Solution:-

Second order methods?

if ~~curvatures don~~"

if G & H does not change as much :- Problem if Curvature also varies.

→ large gradients → greater step

Further:-

$$\eta = \frac{g^T g}{g^T H g}$$

→ high curvature * shorter step

:- Quadratic approximation

other issues:-

if $g^T H g$ ~~decreas~~ decreases very fast ⇒ under estimate the step size.
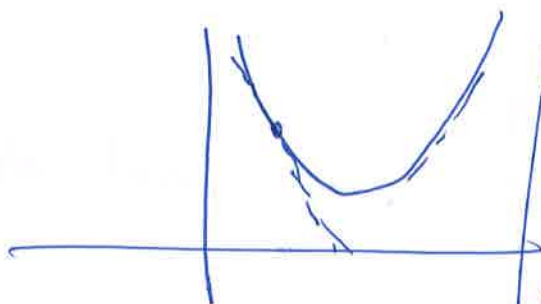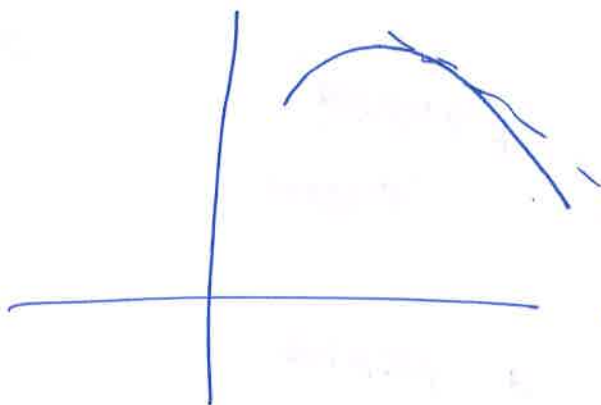
Learning very slow

If $g^T H g$ increases very fast:- over estimate the step size:- Overshoot the minima.

Newton's method:- exact solution if function is quadratic.
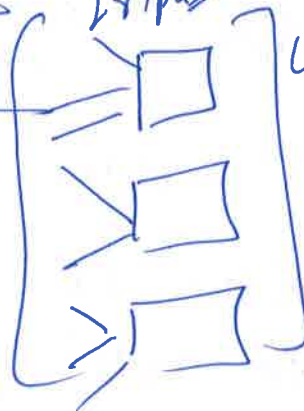
---

Other issues:-

Curvature Related:-

:- Multiple non linearities:- sharp transition in curvature



Model identifiability problem:- Unique combination of weights which result in global minima
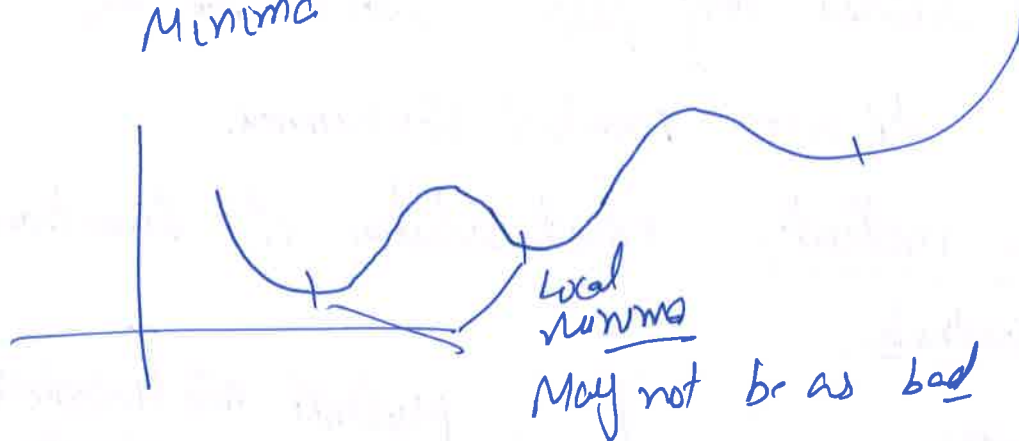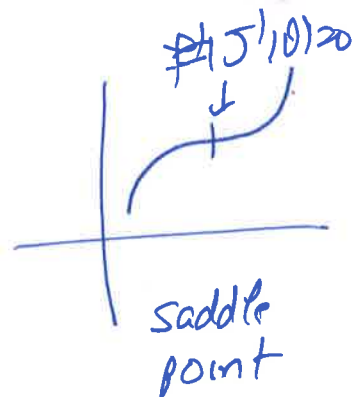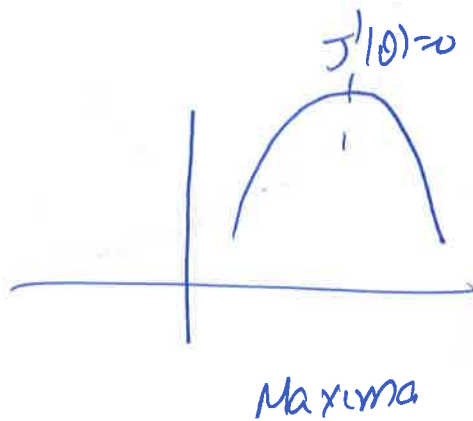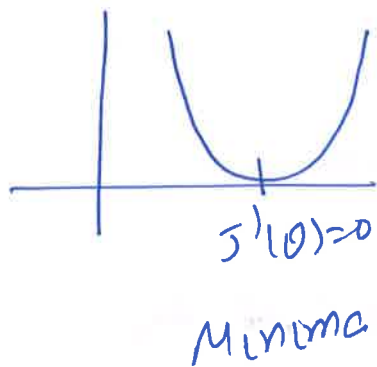
Inputs — unit's exchange unity



Similarly:-

$g(z) = |z|$ (Absolute value rectifier)

Multiply incoming weights by 2 & outgoing weights by 1/2.

Randomly initialize around zero. Break the symmetry.

# Another Issue:—



$J'(\theta)=0$

Minima

$J'(\theta)=0$

Maxima

$\text{but } J'(\theta)=0$

saddle point



Local minima

May not be as bad

## Saddle Points:—  In multiple dimensions

$$\nabla_\theta J(\theta)=0$$

$$\exists i \quad \frac{\partial^2 J(\theta)}{\partial \theta_i^2} > 0 \qquad \cancel{\text{Minima}} \quad \cancel{\text{Sorry}}$$

$$\exists j \quad \frac{\partial^2 J(\theta)}{\partial \theta_j^2} \le 0 \qquad \Rightarrow \text{Maxima}$$

$$\boxed{d^T H d \ge 0 \qquad \forall d} \qquad \therefore \quad \text{Local Minima}$$

$$\lambda_i \ge 0 \quad \forall i$$

$$\begin{cases} d^T H d \le 0 \quad \forall d:- \\ \qquad \lambda_i \le 0 \; \forall i \qquad \text{local maxima} \end{cases}$$

otherwise saddle point $\Rightarrow$ Much more likely to occur than minima/maxima

## Algorithms:-

### Stochastic gradient Descent:-

- Batch Gradient Descent
- Mini-batch gradient Descent -

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

$$= E[$$

$$J_{st}^{S\#}(\theta)$$
$$= \frac{1}{r} \sum_{i=1}^{r} L(f(x^{(i)}, \theta), y^{(i)})$$

⇓

mini batch size

$\overline{\text{minibatch} > 1}$

$O(r)$ time

### Advantages:-

$\nabla_\theta J(\theta)$   $O(m)$ time

$\nabla_\theta J_{st}(\theta)$

⇓ Approximation

Approximate loss using a smaller number of examples.

only movement on the rough direction is important



Each update: $O(m)$ time

Guaranteed to converge to local optima with sufficiently small $\gamma$.

Each update: $O(r)$ time

Considerations in choosing $r$?

↳ $r$ :- small gradient approximate
↳ $r$ :- large Increased computation time
⇒ can be parallelized
→ $r$ :- small :- approximates generalization
error (until you do
a single pass)

Section 8.3 :-
⌐→ variations to # gradient based methods
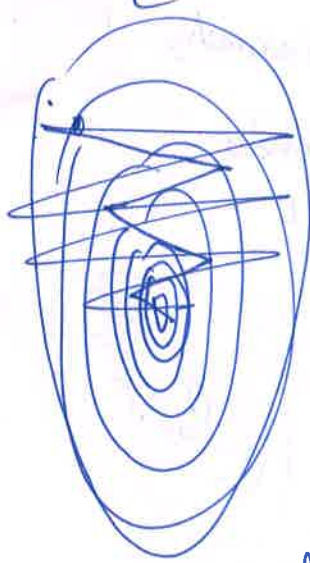└→ second order methods.

Momentum :-

$\nabla_\theta J(\theta)$ :- gradient
Update Rule :-

decay (exponentially decays wrt # iterations).

⇒ Accumulation of previous gradient

$v^{(t+1)}$

$$
\begin{bmatrix}
v = \alpha v - \eta \nabla_\theta J(\theta) \\
{}^{(t+1)} \quad {}^{(t)} \\
\theta = \theta + \alpha v
\end{bmatrix}
$$

$\begin{bmatrix} \text{USE Along with} \\ \text{SGD} \end{bmatrix}$

$v^{(t+1)} = \alpha v^{(t)} - \eta g$
$= \alpha(\alpha v^{(t-1)} - \eta g)$
$\qquad - \eta g$

If $\nabla_\theta J(\theta)$ are aligned
$\equiv g$
then
$- \eta g :-$
$- \eta g\alpha - \eta g$
$- \eta g\alpha^2 - \eta g\alpha - \eta g$

$- \eta \left( \dfrac{\|g\|}{1 - \alpha} \right)$

Nesterov Momentum :-

$$
\begin{bmatrix}
v^{(t+1)} = \alpha v^{(t)} - \eta \nabla_\theta J(\theta^{(t)} + \alpha v^{(t)}) \\
\theta^{(t+1)} = \theta^{(t)} + \alpha v^{(t+1)}
\end{bmatrix}
$$