Local minima in deep networks: Till today it remains an open question Whether local minima with a high error rate relative to the global minima are common in practical deep networks. However, many recent studies seen to However, many recent studies seen to Indicate that most local minima have indicate and generalization characteristics error rates and generalization characteristics that are very similar to global minima.

Critical points: Given an arbitrary function, a point at which the gradient is the zero vector is called a critical point.



a randomly selected critical point being a minima is 1/3. Hence if we have a total of K critical points then on average we will have a total of <u>K</u> minimas-Now let's ansider a cost function in Q-dimensional space. In general, in a Q-dimensional space we an slice through a critical point ion 2-different axes. A critical Point can only be a local minima if it appears as a local minima in every single one of the

I one-dimensional subspaces. Hence using the result from earlier we have that the probability that a have that the probability that a randomly selected critical point in a Q-dimensional space is a local minima is $\frac{1}{3d}$. As d increases [DCa] minima become exponentially more rare.

Momentum:

In lecture, we observed that 1255 Surfaces with high curvature results In OSCILLATIONS of the SGLD updates. The oscillatory nature leads to a The oscillatory nature leads to a Slower convergence of the SGLD.

since SGD will more in the right direction but with sscillations so in order to damper out the scillationswe will use the average gradient to make the updates. $V \in dV - Eg$ $\Theta \in \Theta + V$

It is the velocity and more weight is applied to more recent gradients creatiz an exponentially decoying average of gradients.

Adaptive gradients: A major challenge for training deep networks is appropriately selecting the learning rate. The basic concept behind learning rate alaptation is that the learning rate alaptation is that the optimal learning rate is appropriately modified over the span of learning to achieve good convergence properties.

Adagrad:

- Adapt global learning rate over
 time using an accumulation of
 historical gradients.
- · Keep track of a learning rate for each parameter.

$$a \leftarrow a \leftarrow g \odot g$$

vector

· Parameters with largest gradients experience a rapil decrease in their learning rates while parameters with smaller gradients only observe a small derrease in their learning rates.

Rms Prop: · Use a exponentially weighted miving overage of gradients a G Ba+ (1-B) g Og where p is the decay factor Smaller the decay factor the shorter the effective winclow

• Keep track of an exponentially
weighted moving average of
the gradient
$$V \in B_1 V + (1 - B_1)g$$

h' = WX $= \begin{bmatrix} a & o \end{bmatrix} X$ h' = Wh' $= Wh' = \begin{bmatrix} a^{2} & o \end{bmatrix} X$ $= WWX = \begin{bmatrix} a^{2} & 0 \end{bmatrix} X$

$$h^{n} = W^{n} x$$

$$= \begin{bmatrix} a^{n} & 0 \end{bmatrix} x$$

$$b^{n} = a^{n} x_{1}$$

$$b^{n}_{2} = b^{n} x_{2}$$

Now,
$$J = h_1^n + h_2^n$$

 $\Rightarrow J = \alpha^n x_1 + \beta^n x_2$

$$\frac{\partial J}{\partial a} = na^{n} \dot{x}_{i}$$

$$\frac{\partial J}{\partial b} = nb^{n-1} \dot{x}_{2}$$

$$\frac{\partial J}{\partial b} = \left[na^{n-1} \dot{x}_{i} \right]$$

$$S_{2} \quad \nabla J = \left[na^{n-1} \dot{x}_{i} \right]$$

Now,
Assuming
$$x_1 = x_2 = 1$$
 and $a = 1, b = 2$, we

have

have

$$y = 1 + 2^{n}$$

 $\nabla y = \begin{bmatrix} n \\ n2^{n-1} \end{bmatrix}$
As $n \gg \infty$ then ∇y exploses

Assuming
$$X_1 = X_2 = 1$$
 and $\alpha = 0.5$, $b = 0.9$
we have

 $y = (0.5)^{n} + (0.9)^{n}$ $\nabla y = \left[n (0.5)^{n-1} \right]$ $n (0.9)^{n-1}$ $n (0.9)^{n-1}$

| As | $n \rightarrow \infty$ |
|----|------------------------|
| | |

From problem statement, we have the recorsive relation $V_t = (1 - \beta_1) \stackrel{t}{\underset{i=1}{\overset{t}{\sum}} \beta_i t^{-\hat{v}} g_i t$ $V_{t} = \delta_i t^{-\hat{v}} g_i t$ Let's do a change of variable in the above summation: $j = t - \hat{v}$

Then,

$$\hat{i}=1 \rightarrow \hat{j}=t-1$$

 $\hat{i}=t \rightarrow \hat{j}=0$
 $\hat{i}=t \rightarrow \hat{j}=0$
 $\hat{\beta}_1^{t-\hat{i}} \rightarrow \hat{\beta}_1^{\hat{j}}$
 $\hat{\beta}_1^{t-\hat{i}} \rightarrow \hat{\beta}_1^{\hat{j}}$
 $\hat{\beta}_1^{t-\hat{i}} \rightarrow \hat{\beta}_1^{\hat{j}}$

 $V_{t} = (1 - B_{1}) \sum_{i=0}^{t-1} B_{i}^{i} \partial_{t} - \hat{J}$ Itence, Now, taking expectation of both siles we get $E[V_t] = E[(1-B_i) = B_i^{2} B_i^{2} + B_i^{2}]$ Since EFJ is a linear operator, so $E[V_{t}] = (I - B_{1}) \sum_{i=0}^{t-1} E[B_{i}^{3} \partial_{t} - j]$ Since Bi is a deterministic quantity $FEV_{t}] = (1 - \beta_{1}) \sum_{j=0}^{t-1} \beta_{j}^{3} FE g_{t-3}]$ 50 From problem Statement, $EEg_{t-j}J = M, \quad j = 0, \dots, t-1$

So,

$$EEV_{t}] = (I-\beta_{1}) \stackrel{t-1}{\underset{j=0}{\sum}} \beta_{1}^{j} M$$

 $= (I-\beta_{1}) M \stackrel{t-1}{\underset{j=0}{\sum}} \beta_{1}^{j}$
Recall,
 $I_{j} \beta_{1} \beta_{1}^{2} \beta_{1}^{3} \dots \beta_{1}^{3}$, β_{1}^{j}
 $I_{j} \beta_{1} \beta_{1}^{2} \beta_{1}^{3} \dots \beta_{1}^{j}$
 $is a geometric series with
is a geometric series with
 $a = I \quad and \quad r = \beta_{1} \dots \text{ Then Using}$
 $d = I \quad and \quad r = \beta_{1} \dots \text{ Then Using}$
 $the summation of geometric series$
 $\frac{t-1}{\sum_{j=0}^{j}} \beta_{1}^{j} = \frac{I-\beta_{1}}{I-\beta_{1}}$$

Then,

$$EEV_{t}] = (1-B_{1})M (1-B_{1}^{t})$$

$$1-B_{1}^{t}$$

$$= EEV_{t}] = M(1-B_{1}^{t})$$
Hence,

$$\frac{1}{1-B_{1}^{t}} EEV_{t}] = M = EE9t$$

Similarly)

$$q_t = (1-\beta_2) \sum_{i=1}^{t} \beta_2^{t-i} g_i^2$$

 $q_t = charge of Variables$
like before
 $q_t = (1-\beta_2) \sum_{j=0}^{t-1} \beta_2^j g_{t-j}^2$

Taking EEJ of both sides and simplifying $E[a+] = (1-B_2) \sum_{j=0}^{t-1} B_2 E[9+3]$

 $Since, \quad \gamma = S$ EL 9t-3 = SSo, $E[a_{t}] = (1-\beta_{2}) S \sum_{j=0}^{t-1} \beta_{2}$ Using the summation of geometric $E\left[\alpha t\right] = \left(1 - \beta z\right) S \cdot \frac{\left(1 - \beta z\right)}{1 - \beta z}$ series, $E[a_{t}] = (I - \beta_{2}^{t})S$



Recall that for AER^{mxn}, it's
infinity norm is Defined as
$$||A||_{\infty} = \max_{i=1}^{max} \sum_{j=1}^{2} |q_{ij}|$$

Observe that,

$$\sum_{j=1}^{2} |a_{ij}^{*}| \leftarrow Sum of the
absolute values
of elements in
ith row
is the maximum absolute row
is the maximum absolute row$$

som.



Then,

$$||A||_{\mathcal{B}} = \max\left(21+2+4\right), 23+12\beta, 22+1+2\beta$$

$$= 23$$
Now, coming back to our problem

$$W = \left(\begin{array}{c} w_{11}, w_{12}, \dots, w_{1n} \\ w_{21}, w_{22}, \dots, w_{2n} \\ \vdots & \vdots \\ \vdots & \vdots \\ w_{n1}, w_{n2}, \dots, w_{nn} \end{array}\right)$$

$$\frac{\partial ||w||_{\infty}}{\partial w_{11}} = \begin{bmatrix} \frac{\partial ||w||_{\infty}}{\partial w_{11}} & \frac{\partial ||w||_{\infty}}{\partial w_{1n}} \\ \vdots \\ \vdots \\ \frac{\partial ||w||_{\infty}}{\partial w_{n1}} & \frac{\partial ||w||_{\infty}}{\partial w_{nn}} \end{bmatrix}$$
Suppose, | tell you that the

$$\frac{\partial ||w||_{\infty}}{\partial w_{n1}} + \frac{\partial ||w||_{\infty}}{\partial w_{nn}} = \frac{\partial ||w||_{\infty}}{\partial w_{nn}}$$
Suppose, | tell you that the

$$\frac{\partial ||w||_{\infty}}{\partial w_{n1}} + \frac{\partial ||w||_{\infty}}{\partial w_{nn}} = \frac{\partial ||w||_{\infty}}{\partial w_{nn}}$$

$$= ||w||_{\infty} = \frac{\partial ||w||_{\infty}}{\partial z} + \frac{\partial ||w||_{\infty}}{\partial w_{nn}} = \frac{\partial ||w||_{\infty}}{\partial z}$$

$$= ||w||_{\infty} = \frac{\partial ||w||_{\infty}}{\partial z} + \frac{\partial ||w||_{\infty}}{\partial w_{nn}} = \frac{\partial ||w||_{\infty}}{\partial z}$$

Then,