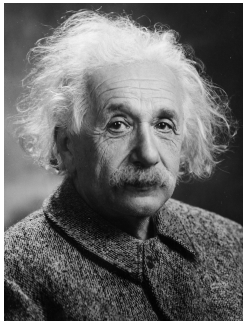# CS-568 Deep Learning

## Nazar Khan

PUCIT

Recurrent Neural Networks

*Everything should be made as simple as possible,
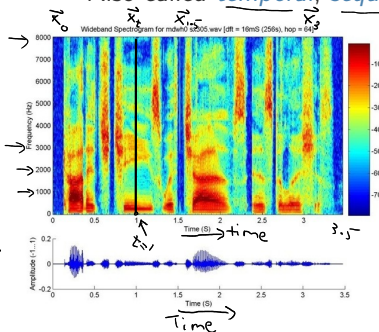but no simpler.*
Albert Einstein

Understanding Recurrent Neural Networks requires some effort and a correct perspective. Do not expect them to be as simple as linear regression.

# Static vs. Dynamic Inputs $\downarrow \boxed{I} \rightarrow label$

- *Static* signals, such as an image, do not change over time.
    - Ordered with respect to space.
    - Output depends on current input.
- *Dynamic* signals, such as text, audio, video or stock price change over time.   hate speech
    - Ordered with respect to time.
    - Output depends on current input as well as past (or even future) inputs.
    - Also called *temporal*, *sequential* or *time-series* data.


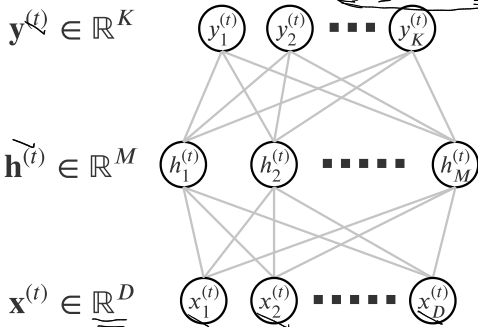
*Deep Learning*

# Context in Text

*The Taj* $\underset{\text{Mahal}}{\_\_\_\_}$ *was commissioned by Shah Jahan in* 1631, *to be built in the memory of* $\_his\_$ *wife Mumtaz Mahal, who died on 17 June* that *year, giving birth to* their *14th child, Gauhara Begum.* Construction *started in 1632, and the mausoleum was completed* $\_in\_$ *1643.*
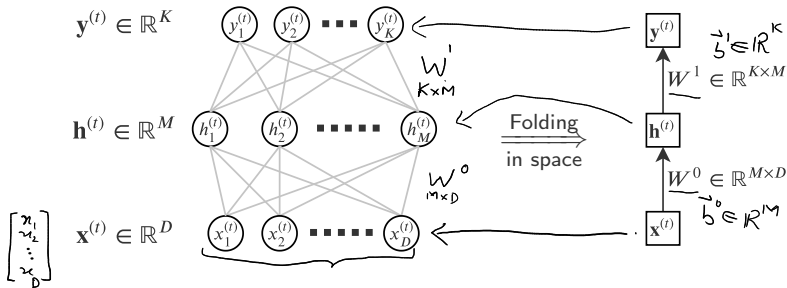
# Time-series Data

▶ A *single* input will be a *series of vectors* $\mathbf{x}^1, \underline{\mathbf{x}}^2, \ldots, \mathbf{x}^{\textcircled{T}}$.

time series
1/input



$$\mathbf{y}^{(t)} \in \mathbb{R}^K$$

$y_1^{(t)}$  $y_2^{(t)}$  ▪ ▪ ▪  $y_K^{(t)}$

$$\mathbf{h}^{(t)} \in \mathbb{R}^M$$

$h_1^{(t)}$  $h_2^{(t)}$  ▪ ▪ ▪ ▪ ▪  $h_M^{(t)}$

$$\mathbf{x}^{(t)} \in \underline{\underline{\mathbb{R}}}^D$$

$x_1^{(t)}$  $x_2^{(t)}$  ▪ ▪ ▪ ▪ ▪  $x_D^{(t)}$

Input component at time $t$ forward propagated through a network.

# Representational Shortcut 1 – Space Folding



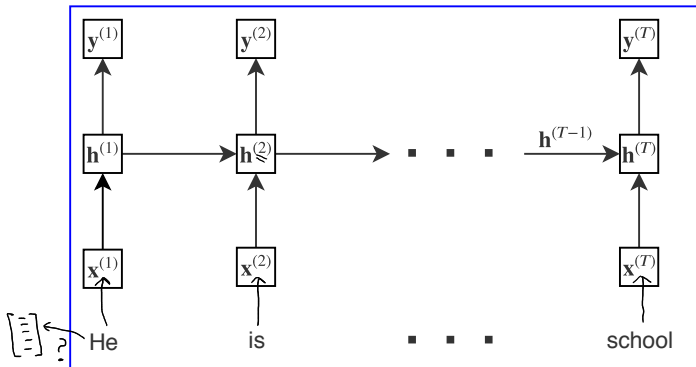Each box represents a layer of neurons.

$$\underset{M \times 1}{\vec{h}} = \tanh\left(\underset{M \times D}{W} \underset{D \times 1}{\vec{x}} + \underset{M \times 1}{\vec{b}}\right)$$

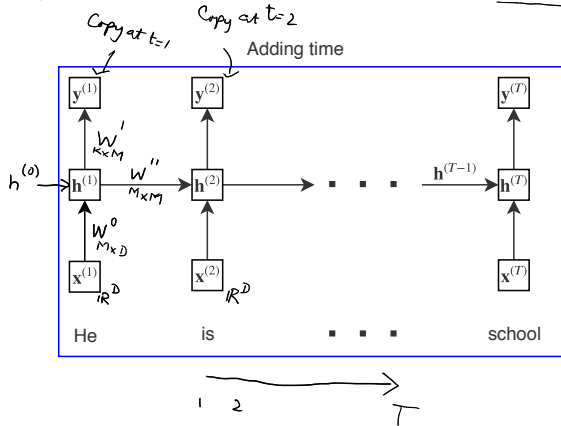# Recurrent Neural Networks

$y^{(2)}$ depends on $x^{(2)}, x^{(1)}$

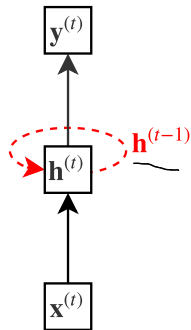$y^{(t)}$ depends on $\vec{x}^{(t)}, \vec{x}^{(t-1)}, ..., \vec{x}^{(1)}$

Adding time



- A recurrent neural network (RNN) makes hidden state at time $t$ directly dependent on the hidden state at time $t-1$ and therefore indirectly *on all previous times*.
- Output $\mathbf{y}_t$ depends on all that the network has already seen so far.

# Representational Shortcut 2 – Time Folding

# Recurrent Neural Networks

$$y = f(a)$$
$$y = f(z)$$

3 sets of weights



$\sigma$, tanh, $\underline{I}$, softmax          Pre-activation of layer 1 at time $t$.

$$\mathbf{y}^{(t)} = f(\overbrace{W^1 \mathbf{h}^{(t)} + \mathbf{b}_1}^{\mathbf{z}^{1(t)}})$$

$$\mathbf{h}^{(t)} = \tanh(\underbrace{W^0 \mathbf{x}^{(t)} + W^{11} \mathbf{h}^{(t-1)} + \mathbf{b}_0}_{\mathbf{z}^{0(t)}})$$

Pre-activation of layer 0 at time $t$.

scalar  $x_1^{(t)} \quad x_2^{(t)} \cdots x_D^{(t)}$
$\vec{x}^{(t)}$

scalar  $h_1^{(t)} \quad h_2^{(t)} \cdots h_M^{(t)}$
$\vec{h}^{(t)}$

# Sequence Mappings



One-to-many          Many-to-one

Many-to-many     Many-to-many delayed

# Loss Functions for Sequences

$$\begin{cases} y^{(1)}, t^{(1)} \\ y^{(2)}, t^{(2)} \\ \vdots \\ y^{(T)}, t^{(T)} \end{cases}$$

for 1 sample

$\vec{y}^{(t)}, \vec{t}^{(t)}$

▶ For recurrent nets, loss is between *series* of output and target vectors. That is $\mathcal{L}(\{\mathbf{y}_1, \ldots, \mathbf{y}_T\}, \{\mathbf{t}_1, \ldots, \mathbf{t}_T\})$.



Forward propagation in an RNN unfolded in time.

▶ Notice that loss $\mathcal{L}$ can be computed only after $\mathbf{y}_T$ has been computed.

## Loss Functions for Sequences

$$\mathcal{L} = \mathcal{L}(y^{(1)}, t^{(1)}) + \mathcal{L}(y^{(2)}, t^{(2)}) + \cdots + \mathcal{L}(y^{(T)}, t^{(T)})$$

- Loss is *not necessarily* decomposable.
- In the following, we will assume decomposable loss $\mathcal{L} = \sum_{t=1}^{T} \mathcal{L}(\mathbf{y}_t, \mathbf{t}_t)$.
- In both cases, as long as $\frac{\partial L}{\partial \mathbf{y}_t}$ has been computed, backpropagation can proceed.

# Backpropagation Through Time (BPTT)



Forward propagation in an RNN unfolded in time. Recurrence between hidden states through pre-activation $\mathbf{z}^{(t)}$ is shown in red.

# BPTT



$$\text{time} = T$$

$$\mathcal{L} \qquad y^{(T)}$$

$$\underline{\mathbf{y}}^{(T)} = f(\underline{W^o}\underline{\mathbf{h}}^{(T)} + \underline{\mathbf{b}}^o)$$

$$\underline{\mathbf{h}}^{(T)} = \tanh(\underline{\mathbf{z}}^{(T)})$$

$$\underline{\mathbf{z}}^{(T)} = W^i\underline{\mathbf{x}}^{(T)} + W^{11}\underline{\mathbf{h}}^{(T-1)} + \mathbf{b}^i$$

Forward propagation in an RNN unfolded in time. Recurrence between hidden states through pre-activation $\mathbf{z}^{(t)}$ is shown in red.

$x=h^0 \longrightarrow z^1 \rightarrow h^1 \xrightarrow{} z^2 \rightarrow h^2 \rightarrow z^3 \rightarrow h^3$

Take-home quiz 4

Find Derivative Formulas

# BPTT – Notational Clarity

▶ For notational clarity, at layer $l$, we will denote the pre-activation by $z^l$ and activation by $h^l$.

▶ So output layer $\mathbf{y}$ will be denoted by $\mathbf{h}^L$ in an $L$-layer network.

▶ Input will be denoted by $\mathbf{h}^0$.

▶ So forward propagation entails

$\longrightarrow \underline{\mathbf{h}^0} \rightarrow \underbrace{\underline{\mathbf{z}^1} \rightarrow \mathbf{h}^1}_{\text{layer } 1} \cdots \rightarrow \underbrace{\underline{\mathbf{z}^{L-1}} \rightarrow \mathbf{h}^{L-1}}_{\text{layer } L-1} \rightarrow \underbrace{\underline{\mathbf{z}^L} \rightarrow \mathbf{h}^L}_{\text{layer } L}.$

▶ For 2 layer network

$\rightarrow \underline{\mathbf{h}^2}^{,T} = \underline{f}(W^1\underline{\mathbf{h}^{1,T}} + \underline{\mathbf{b}^1})$

$\rightarrow \mathbf{h}^{1,T} = \underline{\tanh}(\underline{\mathbf{z}^{1,T}})$

$\rightarrow \underline{\mathbf{z}^{1,T}} = \underline{W^0}\underline{\mathbf{h}^{0,T}} + \underline{W^{11}\mathbf{h}^{1,T-1}} + \underline{\mathbf{b}^0}$



1 hidden layer.

# BPTT – Notational Clarity

# Multivariate Chain Rule

▶ The chain rule of differentiation states

$$\frac{df(u(x))}{dx} = \frac{df}{du}\frac{du}{dx}$$

when $f$ depends on $x$ through $u()$.

▶ The *multivariate* chain rule of differentiation states

$$\frac{df(u(x), v(x))}{dx} = \frac{\partial f}{\partial u}\frac{du}{dx} + \frac{\partial f}{\partial v}\frac{dv}{dx}$$

when $f$ depends on $x$ through $u()$ *and* through $v()$.

▶ Backpropagation is just an application of the multivariate chain rule.



$$\frac{\partial \mathcal{L}}{\partial w} = \sum$$

$$\frac{\partial L}{\partial \vec{y}} = \underbrace{\left[ \frac{\partial L}{\partial y_1} \quad \frac{\partial L}{\partial y_2} \quad \cdots \quad \frac{\partial L}{\partial y_K} \right]}_{1 \times K}$$

$$\frac{s}{M} \leftarrow M^T \qquad \frac{s}{v} \leftarrow v^T \qquad \frac{v}{s} \leftarrow v \qquad \frac{d\vec{v}}{dM}$$

$$\underbrace{\frac{dv_1}{dM}}_{M \times K}$$

$$\underbrace{\frac{\partial L}{\partial W}}_{M \times K} \overset{\text{scalar}}{=} \underbrace{\begin{bmatrix} \frac{\partial L}{\partial w_1} & \frac{\partial L}{\partial w_{21}} & \cdots & \frac{\partial L}{\partial w_{K1}} \\ & & \ddots & \\ & & & \frac{\partial L}{\partial w_{KM}} \end{bmatrix}}_{\text{matrix}}$$

$$\vec{y} = W \vec{x}$$
$$\underline{y_1} = w^1 \vec{x}$$

$$\implies \boxed{\; \underset{\downarrow}{W} \;}$$

**Detour**

Derivative of a loss function $\underset{\text{scalar}}{L}$ for vector output $\underset{K \times 1}{\vec{y}} = \underset{K \times M}{W} \underset{M \times 1}{\vec{x}}$ w.r.t matrix $\underset{K \times M}{W}$.

$$\underbrace{\nabla_W L}_{M \times K} = \underbrace{\nabla_{\vec{y}} L}_{1 \times K} \; \underbrace{\nabla_W \vec{y}}_{K \times (M \times K)}$$

$$\nabla_W \vec{y} = \quad \begin{array}{c} \nabla_W y_1 \\ \nabla_W y_K \end{array}$$
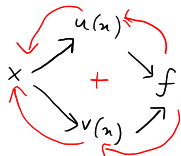
Dimension $\to$ M

Dimension $\to$ K     K $\leftarrow$ Dimension 1

$$\underset{M \times K}{\left\{ \begin{bmatrix} v_1 x_1 & v_2 x_1 & \cdots & v_K x_1 \\ & \ddots & & \\ v_1 x_M & \cdots & & v_K x_M \end{bmatrix} \right.} = \underset{M \times 1}{\vec{x}} \; \underset{1 \times K}{\vec{v}^T} \qquad K \times (M \times K) \qquad \underset{1 \times K}{\nabla_{\vec{y}} L}$$

$$\vec{y} = W \vec{x}$$

Dim 1 $\downarrow$     Dim 2 $\to$

$$\nabla_W y_1 = \begin{bmatrix} \frac{\partial y_1}{\partial w_{11}} & \frac{\partial y_1}{\partial w_{21}} & \cdots \\ \frac{\partial y_1}{\partial w_{11}} & & \\ \vdots & & \\ \frac{\partial y_1}{\partial w_{1M}} & & \end{bmatrix} = 0$$
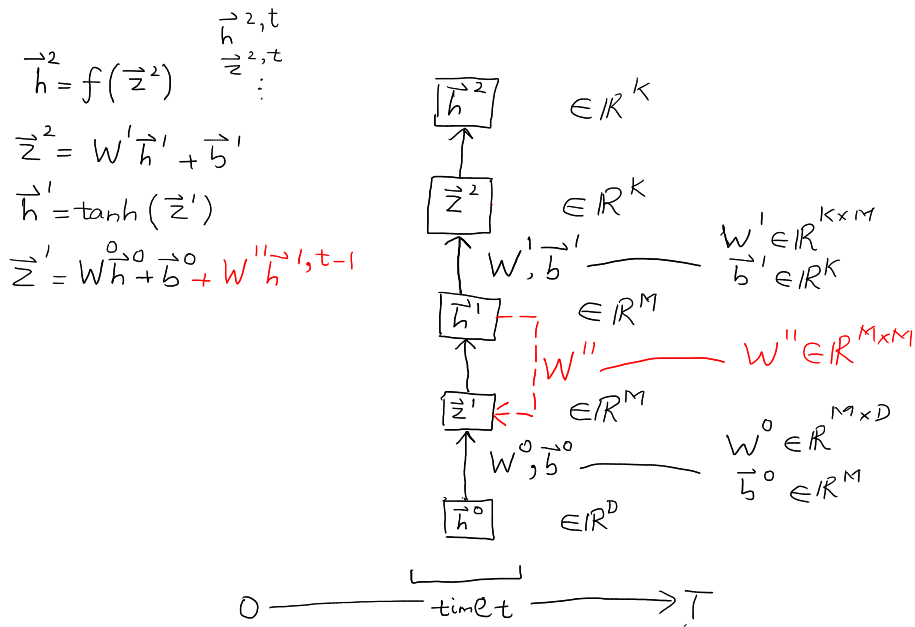
$$= \begin{bmatrix} \vec{x} & 0 \end{bmatrix}$$

$$\vec{v}^T$$

# Back propagation through time (BPTT)

$\vec{h}^2 = f(\vec{z}^2)$

$\vec{z}^2 = W' \vec{h}^1 + \vec{b}'$

$\vec{h}^1 = \tanh(\vec{z}^1)$

$\vec{z}^1 = W^0 \vec{h}^0 + \vec{b}^0 + \textcolor{red}{W'' \vec{h}^{1,t-1}}$

$\boxed{\vec{h}^{2,t}} \quad \in \mathbb{R}^K$    $\vec{h}^{2,t}$, $\vec{z}^{2,t}$ ...

$\boxed{\vec{z}^2} \quad \in \mathbb{R}^K$

$W', \vec{b}' \quad\longrightarrow\quad W' \in \mathbb{R}^{K \times M}$ , $\vec{b}' \in \mathbb{R}^K$

$\boxed{\vec{h}^1} \quad \in \mathbb{R}^M$

$\textcolor{red}{W'' \longrightarrow W'' \in \mathbb{R}^{M \times M}}$

$\boxed{\vec{z}^1} \quad \in \mathbb{R}^M$

$W^0, \vec{b}^0 \quad\longrightarrow\quad W^0 \in \mathbb{R}^{M \times D}$ , $\vec{b}^0 \in \mathbb{R}^M$

$\boxed{\vec{h}^0} \quad \in \mathbb{R}^D$

$0 \xrightarrow{\qquad \text{time } t \qquad} T$

We need **5** derivatives.

① $\nabla_{W'} L$    ③ $\nabla_{W''} L$    ④ $\nabla_{W^0} L$

② $\nabla_{\vec{b}'} L$              ⑤ $\nabla_{\vec{b}^0} L$

$L$ connects to $\vec{z}^{2,1}, \vec{z}^{2,2}, \dots \vec{z}^{2,T}$ and $W'$

$\sigma, \tanh, \text{ReLU}, \dots$

$h_1 = f(z_1) \quad h_2 = f(\vec{z}) \quad h_k = f(z_k)$

---

① $\nabla_{W'} L$    $W'$ affects $L$ through $\vec{z}^2$ at each time step $t$.

So, $\underbrace{\nabla_{W'} L}_{M \times K} = \sum_{t=T}^{1} \overbrace{\underbrace{\nabla_{\vec{z}^{2,t}} L}_{1 \times K}}^{\delta\text{-values}} \underbrace{\nabla_{W'} \vec{z}^{2,t}}_{K \times (M \times K)} = \sum_{t=T}^{1} \underbrace{\vec{h}^{1,t}}_{M \times 1} \underbrace{\nabla_{\vec{z}^{2,t}} L}_{1 \times K}$

where

$\underbrace{\nabla_{\vec{z}^{2,t}} L}_{1 \times K} = \underbrace{\nabla_{\vec{h}^{2,t}} L}_{1 \times K} \underbrace{\nabla_{\vec{z}^{2,t}} \vec{h}^{2,t}}_{K \times K} = \nabla_{\vec{h}^{2,t}} L$

$\underset{\substack{\text{Jacobian} \\ \text{Matrix} \\ \downarrow \\ \text{Derivatives of outputs} \\ \text{w.r.t inputs.} \\ y = f(x) \\ \frac{dy}{dx}}}{} \quad \begin{bmatrix} \frac{\partial h_1}{\partial z_1} & \frac{\partial h_1}{\partial z_2} & \cdots & \frac{\partial h_1}{\partial z_K} \\ \frac{\partial h_2}{\partial z_1} & \frac{\partial h_2}{\partial z_2} & \cdots & \frac{\partial h_2}{\partial z_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_K}{\partial z_1} & \frac{\partial h_K}{\partial z_2} & \cdots & \frac{\partial h_K}{\partial z_K} \end{bmatrix}^{2,t}$

$K \times K$

- For $\sigma, \tanh, \text{ReLU}$, Jacobian matrix is diagonal.
- For softmax, Jacobian matrix will be full.

---

② $\nabla_{\vec{b}'} L$    Follow the reasoning for $\nabla_{W'} L$ above but note that $\nabla_{\vec{b}'} \vec{z}^{2,t} = I_{K \times K}$

So, $\nabla_{\vec{b}'} L = \sum_{t=T}^{1} \nabla_{\vec{z}^{2,t}} L$

---

③ $\nabla_{W''} L$    $W''$ affects $L$ through $\vec{z}^1$ at each time $t$.

$\underbrace{\nabla_{W''} L}_{M \times M} = \sum_{t=T}^{1} \underbrace{\nabla_{\vec{z}^{1,t}} L}_{1 \times M} \underbrace{\nabla_{W''} \vec{z}^{1,t}}_{M \times (M \times M)} = \sum_{t=T}^{1} \underbrace{\textcolor{red}{\vec{h}^{1,t-1}}}_{M \times 1} \underbrace{\nabla_{\vec{z}^{1,t}} L}_{1 \times M}$

where $\nabla_{\vec{z}^{1,t}} L = \nabla_{\vec{h}^{1,t}} L \underbrace{\nabla_{\vec{z}^{1,t}} \vec{h}^{1,t}}_{\substack{M \times M \\ \text{Jacobian} \\ \text{matrix}}} \begin{bmatrix} \frac{\partial h_1}{\partial z_1} & \frac{\partial h_1}{\partial z_2} & \cdots & \frac{\partial h_1}{\partial z_M} \\ \frac{\partial h_2}{\partial z_1} & \frac{\partial h_2}{\partial z_2} & \cdots & \frac{\partial h_2}{\partial z_M} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_M}{\partial z_1} & \frac{\partial h_M}{\partial z_2} & \cdots & \frac{\partial h_M}{\partial z_M} \end{bmatrix}^{1,t}$

$M \times M$

$\vec{h}^{1,t}$ affects $L$ through $\vec{z}^{2,t}$ at each time $t$ and through $\vec{z}^{1,t+1}$ at each time $t+1$.



So,

$$\underbrace{\nabla L}_{\vec{h}^{1,t}} = \underbrace{\nabla L}_{\vec{z}^{2,t}} \underbrace{\nabla \vec{z}^{2,t}}_{\vec{h}^{1,t}} + \underbrace{\nabla L}_{\vec{z}^{1,t+1}} \underbrace{\nabla \vec{z}^{1,t+1}}_{\vec{h}^{1,t}}$$

$$\underbrace{1 \times M}_{} = \underbrace{\nabla L}_{\vec{z}^{2,t}} \underbrace{W^1}_{} + \underbrace{\nabla L}_{\vec{z}^{1,t+1}} \underbrace{W^{11}}_{} \longleftarrow \text{This term will not be present for } t = T.$$

$$\underbrace{\phantom{xx}}_{1 \times K} \underbrace{\phantom{xx}}_{K \times M} \qquad \underbrace{\phantom{xx}}_{1 \times M} \underbrace{\phantom{xx}}_{M \times M}$$

④ $\underline{\nabla L \atop W^0}$

$W^0$ affects $L$ through $\vec{z}^1$ at each time $t$.

So, $$\underbrace{\nabla L}_{W^0} = \sum_{t=T}^{1} \underbrace{\nabla L}_{\vec{z}^{1,t}} \underbrace{\nabla \vec{z}^{1,t}}_{W^0} = \sum_{t=T}^{1} \underbrace{\vec{h}^{0,t}}_{D \times 1} \underbrace{\nabla L}_{\vec{z}^{1,t}}$$

$$\underbrace{\phantom{xx}}_{D \times M} \qquad \underbrace{\phantom{xx}}_{1 \times M} \underbrace{\phantom{xx}}_{M \times (D \times M)} \qquad \underbrace{\phantom{xx}}_{1 \times M}$$

⑤ $\underline{\nabla L \atop \vec{b}^0}$

Follow the reasoning for $\nabla L \atop W^0$ above and note that $\nabla_{\vec{b}^0} \vec{z}^{1,t} = \underset{M \times M}{I}$

So, $$\underbrace{\nabla L}_{\vec{b}^0} = \sum_{t=T}^{1} \underbrace{\nabla L}_{\vec{z}^{1,t}}$$

$$\underbrace{\phantom{xx}}_{1 \times M} \qquad \underbrace{\phantom{xx}}_{1 \times M}$$

Now you have $\underline{\underline{ALL}}$ the derivatives required to train an RNN with 1 hidden layer.

---

### Training an RNN

Notice that an RNN is a very very deep network.

$\vec{h}^2 = f(\vec{z}^2)$

$\vec{z}^2 = W^1 \vec{h}^1 + \vec{b}^1$

$\vec{h}^1 = \tanh(\vec{z}^1)$

$\vec{z}^1 = \overset{0}{W} \overset{0}{\vec{h}} + \overset{0}{\vec{b}} \; {\color{red}+ \; W^{11} \vec{h}^{1,t-1}}$

$\overset{\vec{h}^{2,t}}{\underset{\vec{z}^{2,t}}{\vdots}}$

$\delta_\cdot = \underbrace{h'(a_k)}_{\substack{\sigma \quad \tanh \\ \leq \frac{1}{4} \quad \leq 1}} \sum_{k=1}^{K} w_{kj} \underbrace{\delta_k}_{\leq 1}$

Vanishing gradients

Exploding gradients $\longleftarrow \boxed{W^{11}}$

evals. $\lambda$

---

### Benefit of Recurrent architecture over standard MLP

N-bit addition



2-bits

$\begin{array}{r} 11 \\ + 11 \\ \hline 110 \end{array}$

Draw backs

— Lots of training examples.

  Exponentially many in terms of $n$

— Generalization $\Longrightarrow$ Training errors.

— Will work for n-bit numbers only.



Bit 1    Bit 2

— 8 training examples

— Perfect answers

— Will work for numbers of any size.

| C | $B_1$ | $B_2$ | Y |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 |

Parity  (# 1s even or odd)

0/1

$\boxed{\text{MLP}}$

$\boxed{|\,|\,|\,|\,|\,|}$

Same drawbacks as above.

1



Next bit

$\underline{0010}$         $\underline{00101}$
1                $\underline{10}$

---

Bidirectional RNNs



① $F_{prop}\left(\vec{h}^{0,1}, \vec{h}^{0,2}, \ldots, \vec{h}^{0,T}, \vec{h}^{1,0}\right) \rightarrow \vec{h}_f^{1,1}, \vec{h}_f^{1,2}, \ldots, \vec{h}_f^{1,T}$

② $F_{prop}\left(\underbrace{\vec{h}^{0,T}, \vec{h}^{0,T-1}, \ldots, \vec{h}^{0,1}}_{\text{Flip input seq.}}, \vec{h}^{1,T+1}\right) \rightarrow \underbrace{\vec{h}_b^{1,1}, \vec{h}_b^{1,2}, \ldots, \vec{h}_b^{1,T}}_{\text{after flipping in time.}}$

③ Compute $\vec{h}^{2,t}$ using $\vec{h}_f^{1,t}$ and $\vec{h}_b^{1,t}$.

## BPTT

▶ After dropping $L$, $T$ for clarity and using the multivariate chain rule

$$\frac{\partial \mathcal{L}}{\partial z_i} = \frac{\partial \mathcal{L}}{\partial h_1}\frac{\partial h_1}{\partial z_i} + \frac{\partial \mathcal{L}}{\partial h_2}\frac{\partial h_2}{\partial z_i} + \cdots + \frac{\partial \mathcal{L}}{\partial h_K}\frac{\partial h_K}{\partial z_i}$$

$$= \underbrace{\left[\begin{array}{cccc} \frac{\partial \mathcal{L}}{\partial h_1} & \frac{\partial \mathcal{L}}{\partial h_2} & \cdots & \frac{\partial \mathcal{L}}{\partial h_K} \end{array}\right]}_{\nabla_{\mathbf{h}}\mathcal{L}} \underbrace{\left[\begin{array}{c} \frac{\partial h_1}{\partial z_i} \\ \frac{\partial h_2}{\partial z_i} \\ \vdots \\ \frac{\partial h_K}{\partial z_i} \end{array}\right]}_{\nabla_{z_i}\mathbf{h}} = \nabla_{\mathbf{h}}\mathcal{L}\nabla_{z_i}\mathbf{h}$$

## BPTT

▶ This allows us to write

$$
\underbrace{\nabla_{\mathbf{z}} \mathcal{L}}_{1 \times K} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial z_1} & \frac{\partial \mathcal{L}}{\partial z_2} & \cdots & \frac{\partial \mathcal{L}}{\partial z_K} \end{bmatrix}
$$

$$
= \begin{bmatrix} \nabla_{\mathbf{h}} \mathcal{L} \nabla_{z_1} \mathbf{h} & \nabla_{\mathbf{h}} \mathcal{L} \nabla_{z_2} \mathbf{h} & \ldots & \nabla_{\mathbf{h}} \mathcal{L} \nabla_{z_K} \mathbf{h} \end{bmatrix}
$$

$$
= \underbrace{\nabla_{\mathbf{h}} \mathcal{L}}_{1 \times K} \underbrace{\begin{bmatrix} \nabla_{z_1} \mathbf{h} & \nabla_{z_2} \mathbf{h} & \ldots & \nabla_{z_K} \mathbf{h} \end{bmatrix}}_{K \times K} = \nabla_{\mathbf{h}} \mathcal{L} \nabla_{\mathbf{z}} \mathbf{h}
$$

where

$$
\nabla_{\mathbf{z}} \mathbf{h} = \begin{bmatrix} \frac{\partial h_1}{\partial z_1} & \frac{\partial h_1}{\partial z_2} & \cdots & \frac{\partial h_1}{\partial z_K} \\ \frac{\partial h_2}{\partial z_1} & \frac{\partial h_2}{\partial z_2} & \cdots & \frac{\partial h_2}{\partial z_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_K}{\partial z_1} & \frac{\partial h_K}{\partial z_2} & \cdots & \frac{\partial h_K}{\partial z_K} \end{bmatrix}
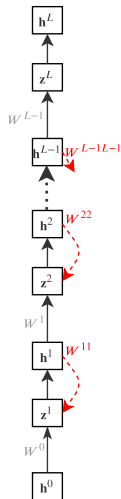$$

is the $K \times K$ *Jacobian matrix* of $\mathbf{h}$ with respect to $\mathbf{z}$.

▶ For simple, per-component activation functions (tanh, ReLU), Jacobian will be diagonal. For softmax, it will be dense.

# BPTT

Derivative of loss $\mathcal{L}$ for any time $t$ w.r.t weights of any layer $l$ can be computed as

$$\underbrace{\nabla_{W^{l-1}}\mathcal{L}\Big|_t}_{K \times M} = \begin{bmatrix} \dfrac{\partial \mathcal{L}}{\partial W_{11}^{l-1}} & \dfrac{\partial \mathcal{L}}{\partial W_{12}^{l-1}} & \cdots & \dfrac{\partial \mathcal{L}}{\partial W_{1M}^{l-1}} \\[2ex] \dfrac{\partial \mathcal{L}}{\partial W_{21}^{l-1}} & \dfrac{\partial \mathcal{L}}{\partial W_{22}^{l-1}} & \cdots & \dfrac{\partial \mathcal{L}}{\partial W_{2M}^{l-1}} \\[2ex] \vdots & \vdots & \ddots & \vdots \\[2ex] \dfrac{\partial \mathcal{L}}{\partial W_{K1}^{l-1}} & \dfrac{\partial \mathcal{L}}{\partial W_{K2}^{l-1}} & \cdots & \dfrac{\partial \mathcal{L}}{\partial W_{KM}^{l-1}} \end{bmatrix}$$
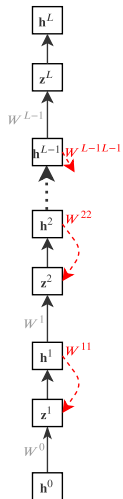
## BPTT

Loss $\mathcal{L}$ for time $t$ depends on $W^{l-1}$ through $\mathbf{z}^{l,t}$

$$\underbrace{\nabla_{W^{l-1}}\mathcal{L}}_{K \times M}\bigg|_t = \begin{bmatrix} \dfrac{\partial \mathcal{L}}{\partial z_1^{l,t}}\dfrac{\partial z_1^{l,t}}{\partial W_{11}^{l-1}} & \dfrac{\partial \mathcal{L}}{\partial z_1^{l,t}}\dfrac{\partial z_1^{l,t}}{\partial W_{12}^{l-1}} & \cdots & \dfrac{\partial \mathcal{L}}{\partial z_1^{l,t}}\dfrac{\partial z_1^{l,t}}{\partial W_{1M}^{l-1}} \\[2em] \dfrac{\partial \mathcal{L}}{\partial z_2^{l,t}}\dfrac{\partial z_2^{l,t}}{\partial W_{21}^{l-1}} & \dfrac{\partial \mathcal{L}}{\partial z_2^{l,t}}\dfrac{\partial z_2^{l,t}}{\partial W_{22}^{l-1}} & \cdots & \dfrac{\partial \mathcal{L}}{\partial z_2^{l,t}}\dfrac{\partial z_2^{l,t}}{\partial W_{2M}^{l-1}} \\[2em] \vdots & \vdots & \ddots & \vdots \\[2em] \dfrac{\partial \mathcal{L}}{\partial z_K^{l,t}}\dfrac{\partial z_K^{l,t}}{\partial W_{K1}^{l-1}} & \dfrac{\partial \mathcal{L}}{\partial z_K^{l,t}}\dfrac{\partial z_K^{l,t}}{\partial W_{K2}^{l-1}} & \cdots & \dfrac{\partial \mathcal{L}}{\partial z_K^{l,t}}\dfrac{\partial z_K^{l,t}}{\partial W_{KM}^{l-1}} \end{bmatrix}$$

## BPTT

$$
\underbrace{\nabla_{W^{l-1}}\mathcal{L}\Big|_t}_{K \times M} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial z_1^{l,t}} h_1^{l-1,t} & \frac{\partial \mathcal{L}}{\partial z_1^{l,t}} h_2^{l-1,t} & \cdots & \frac{\partial \mathcal{L}}{\partial z_1^{l,t}} h_M^{l-1,t} \\ \frac{\partial \mathcal{L}}{\partial z_2^{l,t}} h_1^{l-1,t} & \frac{\partial \mathcal{L}}{\partial z_2^{l,t}} h_2^{l-1,t} & \cdots & \frac{\partial \mathcal{L}}{\partial z_2^{l,t}} h_M^{l-1,t} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \mathcal{L}}{\partial z_K^{l,t}} h_1^{l-1,t} & \frac{\partial \mathcal{L}}{\partial z_K^{l,t}} h_2^{l-1,t} & \cdots & \frac{\partial \mathcal{L}}{\partial z_K^{l,t}} h_M^{l-1,t} \end{bmatrix}
$$

$$
= \underbrace{\begin{bmatrix} \frac{\partial \mathcal{L}}{\partial z_1^{l,t}} \\ \frac{\partial \mathcal{L}}{\partial z_2^{l,t}} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial z_K^{l,t}} \end{bmatrix}}_{\nabla_{\mathbf{z}^{l,t}}^{Tr}\mathcal{L}} \underbrace{\begin{bmatrix} h_1^{l-1,t} & h_2^{l-1,t} & \cdots & h_M^{l-1,t} \end{bmatrix}}_{\nabla_{W^{l-1}}^{Tr}\mathbf{z}^{l,t}} = \left( \nabla_{\mathbf{z}^{l,t}}^{Tr}\mathcal{L} \right) \left( \nabla_{W^{l-1}}^{Tr}\mathbf{z}^{l,t} \right)
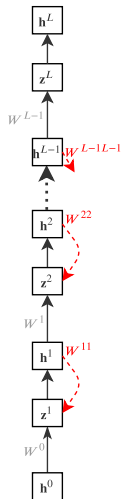$$

# BPTT

▶ Start from final output $\mathbf{h}^{L,T} = f\left(\mathbf{z}^{L,T}\right)$. Assuming $\nabla_{\mathbf{h}^{L,T}}\mathcal{L}$ has been computed,

$$\underbrace{\nabla_{\mathbf{z}^{L,T}}\mathcal{L}}_{1 \times K} = \underbrace{\nabla_{\mathbf{h}^{L,T}}\mathcal{L}}_{1 \times K} \underbrace{\nabla_{\mathbf{z}^{L,T}}\mathbf{h}^{L,T}}_{K \times K}$$

▶ Derivative of loss w.r.t weights of final layer *at time T* can now be computed as

$$\underbrace{\nabla_{W^{L-1}}\mathcal{L}}_{K \times M_{L-1}}\Big|_{T} = \left(\nabla_{\mathbf{z}^{L,T}}^{Tr}\mathcal{L}\right)\left(\mathbf{h}^{L-1,T}\right)^{Tr}$$

# BPTT



- Viewed as a function of $W^{L-1}$, loss can be written over time as

$$\mathcal{L}\left(\mathbf{z}^{L,1}(W^{L-1}), \mathbf{z}^{L,2}(W^{L-1}), \ldots, \mathbf{z}^{L,T}(W^{L-1})\right)$$

- Therefore, using multivariate chain rule

$$\nabla_{W^{L-1}}\mathcal{L} = \left(\left(\nabla_{\mathbf{z}^{L,T}}^{Tr}\mathcal{L}\right)\left(\mathbf{h}^{L-1,T}\right)^{Tr} + \left(\nabla_{\mathbf{z}^{L,T-1}}^{Tr}\mathcal{L}\right)\left(\mathbf{h}^{L-1,T-1}\right)^{Tr} + \ldots + \left(\nabla_{\mathbf{z}^{L,1}}^{Tr}\mathcal{L}\right)\left(\mathbf{h}^{L-1,1}\right)^{Tr}\right)$$

$$= \sum_{1}^{t=T}\left(\nabla_{\mathbf{z}^{L,t}}^{Tr}\mathcal{L}\right)\left(\mathbf{h}^{L-1,t}\right)^{Tr}$$
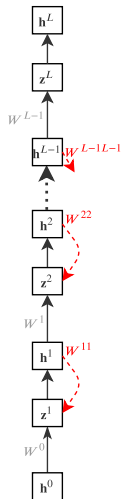
## BPTT

▶ Similarly, for each layer $l$ from 1 to $L$

$$\nabla_{W^{l-1}}\mathcal{L} = \sum_{1}^{t=T} \left(\nabla_{\mathbf{z}^{l,t}}^{Tr}\mathcal{L}\right) \left(\mathbf{h}^{l-1,t}\right)^{Tr} \tag{1}$$

# BPTT

$$\underbrace{\nabla_{\mathbf{h}^{L-1}} \mathcal{L} \Big|_{T}}_{1 \times M_{L-1}} = \underbrace{\nabla_{\mathbf{z}^{L}} \mathcal{L} \Big|_{T}}_{1 \times M_{L}} \underbrace{W^{L-1}}_{M_{L} \times M_{L-1}}$$

$$\underbrace{\nabla_{\mathbf{h}^{L-1}} \mathcal{L} \Big|_{t}}_{1 \times M_{L-1}} = \underbrace{\nabla_{\mathbf{z}^{L}} \mathcal{L} \Big|_{t}}_{1 \times M_{L}} \underbrace{W^{L-1}}_{M_{L} \times M_{L-1}} + \underbrace{\nabla_{\mathbf{z}^{L-1}} \mathcal{L} \Big|_{t+1}}_{1 \times M_{L-1}} \underbrace{W^{L-1,L-1}}_{M_{L-1} \times M_{L-1}}$$

# RNN Variations

THQ 4



1 hidden state     2 hidden states     Skip connections