CS 378 Lecture 16

Today - RNNS - LSTMS (the type of RNN you will be using) - Implementation - Implementation Announcements - AU out later today - Midtern back Monday Recop RNNs + Language modeling T SAW THE dog P(w | I saw the day) = softmax (Zh)

Training "Backpropagation through time" = backprop h ZOparams Cell (Q) -> W,V-> W,V-> P(W|--) (Q) -> W,V-> P(W|--) (X) -> V,V-> P(W|--) (X) -> V,V-> P(W|--) (S) = log p(V. back prop to get gradient WV-X: Multiple upon a Multiple updates for W, V, sum them together Backprop: 1055 -> gradient for each param

LSTMS

What we've seen so for is a recipe for RNNs Many types of RNNs

Long short-tern memory (LSTM)
1998

RNN state is a "short-term memory" LSTMs are better at venembering into for more fimesteps

Problem with Elman networks: vanishing gradient problem (exploding) Elman: $h_3 = \operatorname{Tanh}(WX_3 + V)$ tanh F

LSTM definition $Elman: \overline{h_i} = tanh(W\overline{X_i} + V\overline{h_{i-1}})$ Gated : h; = h: 0 f + function (x;, h: ...) oi prev state element vise multiplication f: forget gate ininput gate vector of values in [0,1] $\overline{h}_{i-1} \begin{bmatrix} 1 & f \\ 0 & f \\ 0 & 0 & -5 \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ $f = sigmoid (W\overline{x}_{i} + W^{2}\overline{h}_{i-1})$ $i = sigmoid \left(W^{(3)} \times W^{(n)} \overline{h} \right)$ $O: \underbrace{e^{\times}}_{l+e^{\times}} \longrightarrow$



LSTM: next page Real

forget tanh σ Xt gate (\mathbf{X}_{t-1}) Chris Olah blog LSTM: & weight matrices h, C in Pytorch, these are dealt with as اام h idden state state a tople hi -

frget bias Po[1: $f = sigmoid \left(W^{(1)} \overline{\chi}_{i} + W^{(2)} \overline{h}_{i-1} + b_{forget} \right)$ $i = sigmoid \left(W^{(3)} \overline{\chi}_{i} + W^{(4)} \overline{h}_{i-1} + b_{input} \right)$

[0,1,2] Closer to [1,1,2) than [0,1,3]

Most recent stiff is more velevant

Forget gete high: Charges forther back matter more