S 378 Lecture 16

Announcements

Assignments: A4 only coding (but simplified), A5 no coding, FP in progress

Extra slip days

Zoom protocol: type Qs in chat, polls on Canvas, breakouts

Recap

$$N-gram model: P(w) = TTP(w; |w_{i-nfl/m}, w_{i-l})$$

Estimate counts from large corpora, smooth

Goals Recovert neural networks (RNN) 1. key RNN abstraction 2. Key properties of LSTM (long short-term memory) This lecture: LM, RNN definition Next lecture: training + implementation This sets the stage for machine translation, summarization, dialogue Neural Language Models sourie Discrim model for P(wilw, -. wi-1) all past words What we want neural net to look at w_- wi-1, place distribution over Wi

RNN: neural model for encoding sees of arbitrary length prev hidden yi output state 7 X, y, h vectors $h_{i-1} \rightarrow 7$ hi-1, Xi inputs ->Ci Yi, hi Outputs Ci-15 RNN tyles X- Current Elman network (UT ling PhD 1977) Vi = tanh (Uhi)

Virolled the RNN over this sent: p(wi) O his his his his hull in the dug matrix embed I saw the dug mul. F softmax forward() in PyTorch Mas a for loop, but otherwise this is just a special feed forward network At each step, hi captures the model's "picture" of the sentence so far $P(w; | w, ..., w_{i-1}) = \text{softmax}(Mh_{i-1})$ $h_{i-1} = RNN(w_{i-1})$ W, U, V, M, are our IVI-Ihl only params!

Training RNNs "Backpropagation through time" = backpropagation More next time Long Short-form memory network Vanishing gradient problem EO Containty backprop Vanishing gradient problem backprop Vanishing gradient problem Elman networks (not ESTMS) : ervor tem "dies out" Model Can't learn to feed info long distance

Key idea: gates $Elman: h_i = tanh (WX_i + Vh_{i-i})$ Gated: h; = h; -1 Of + func (x;) Oi 0 = element-wise mul f: vector E[0,1]d hi-1 ETRd f= 1 preserves hi: hi = 2 func(xi)oi s=0 Otherwise zerves out some parts

tloias - elementwise to F= sigmoid (W'X; +Whi-i) $\sigma = \frac{e^{2}}{1+e^{2}}$

 $\bar{I} = sigmoid \left(W_{X_i}^3 + W_{h_i-i}^4 \right)$ LSTM: hi \bar{C} : short-ferm 'cell' memory forget $\overline{FOE_i}$ $\overline{C_i} = \overline{C_{i-1}} \circ \overline{F} + func(x_i)$ h; = C; O O Canother outpute furg A Α

Source: Chris Olah blog post (link on course website)

Key properties: RINN seg of words (vectors) as input encodes into State h; LSTM: 2 hilden States (h, c) h vs. C distinction is not Critic Critica (better at venembering info for long Sequences by Using gates rather than tanh & matrix MUI.

Resources: Chris Olah blog PyTorch example 1stm-leeturepy Floydhub tutorial