

## Automatic Speech Recognition (CS753)

Lecture 17: RNN language models + Introduction to Kaldi

Instructor: Preethi Jyothi Mar 20, 2017

## Recap

• Language models using feedforward neural networks with fixed length context  $i-\text{th output} = P(w_t = i \mid context)$ 



Image from: Bengio et al., "A neural probabilistic language model", JMLR, 03

## Recap

- Language models using feedforward neural networks with fixed length context: Provide significant improvements in ASR performance over Ngram based language models
- But fixed length context needs to be specified before training
- Recurrent neural networks (RNNs) do not use limited size of context
- Build RNN-based language models

## Simple RNN language model



Current word, x<sub>t</sub>
Hidden state, s<sub>t</sub>
Output, y<sub>t</sub>

$$s_t = f(Ux_t + Ws_{t-1})$$
$$o_t = \operatorname{softmax}(Vs_t)$$

• RNN is trained using the cross-entropy criterion

Image from: Mikolov et al., "Recurrent neural network based language model", Interspeech 10

#### RNN-LMs

- Optimizations used for NNLMs are relevant to RNN-LMs as well (rescoring Nbest lists or lattices, using a shortlist, etc.)
- Perplexity reductions over Kneser-Ney models:

Model	# words	PPL	WER
KN5 LM	200K	336	16.4
KN5 LM + RNN 90/2	200K	271	15.4
KN5 LM	1M	287	15.1
KN5 LM + RNN 90/2	1M	225	14.0
KN5 LM	6.4M	221	13.5
KN5 LM + RNN 250/5	6.4M	156	11.7

Image from: Mikolov et al., "Recurrent neural network based language model", Interspeech 10

## Training RNN-LMs

- RNN-LMs are trained using backpropagation through time (BPTT): Unfold the RNN in time + train the unfolded RNN using backpropagation + mini-batch gradient descent
- Main issues with BPTT: Exploding and vanishing gradients
  - Exploding gradients: Gradients can increase exponentially over time during backpropagation. Clip values of gradients to handle this.
  - Vanishing gradients: Magnitude of gradients approach very tiny values as we propagate gradients back in time. Architectures like Long Short Term Memory (LSTMs) networks can handle this.

## LSTM-LMs

p(w|h)



 Vanilla RNN-LMs unlikely to show full potential of recurrent models due to issues like vanishing gradients

 LSTM-LMs: Similar to RNN-LMs except use LSTM units in the 2nd hidden (recurrent) layer

Image from: Sundermeyer et al., "LSTM NNs for Language Modeling", 10

output layer

2nd hidden layer

projection layer

input layer

#### Comparing RNN-LMs with LSTM-LMs



Image from: Sundermeyer et al., "LSTM NNs for Language Modeling", 10

#### Character-based RNN-LMs

"e" "" target chars: "]" "o" 0.5 0.2 0.1 1.0 2.2 0.3 -1.5 0.5 output layer -3.0 1.9 -1.0 -0.1 -1.1 1.2 2.2 4.1 W\_hy 0.3 -0.3 1.0 0.1 W hh hidden layer -0.1 0.3 -0.5 0.9 0.9 0.1 -0.3 0.7 W\_xh 0 0 0 0 1 0 0 input layer 0 0 1 0 0 0 0 input chars: "h" " 33 "" "e"

Image from: http://karpathy.github.io/2015/05/21/rnn-effectiveness/

# Generate text using a trained character-based LSTM-LM

VIOLA:

WHY, SALISBURY MUST FIND HIS FLESH AND THOUGHT THAT WHICH I AM NOT APS, NOT A MAN AND IN FIRE, TO SHOW THE REINING OF THE RAVEN AND THE WARS TO GRACE MY HAND REPROACH WITHIN, AND NOT A FAIR ARE HAND, THAT CAESAR AND MY GOODLY FATHER'S WORLD; WHEN I WAS HEAVEN OF PRESENCE AND OUR FLEETS, WE SPARE WITH HOURS, BUT CUT THY COUNCIL I AM GREAT, MURDERED AND BY THY MASTER'S READY THERE MY POWER TO GIVE THEE BUT SO MUCH AS HELL: SOME SERVICE IN THE NOBLE BONDMAN HERE, WOULD SHOW HIM TO HER WINE.

## The Big Picture

## Putting it all together: How do we recognise an utterance?

- *A:* speech utterance
- *O<sub>A</sub>*: acoustic features corresponding to the utterance *A*

$$W^* = \underset{W}{\operatorname{arg\,max}} \Pr(O_A|W) \Pr(W)$$

- Return the word sequence that jointly assigns the highest probability to  $O_A$
- How do we estimate  $Pr(O_A|W)$  and Pr(W)?
- How do we decode?





 All the free parameters (means, covariances, mixture weights, transition probabilities) are learned using EM (Baum-Welch) algorithm

#### Language Model

$$W^* = \underset{W}{\operatorname{arg\,max}} \operatorname{Pr}(O_A|W) \operatorname{Pr}(W)$$

$$\Pr(W) = \Pr(w_1, w_2, \dots, w_N)$$
$$= \Pr(w_1) \dots \Pr(w_N | w_{N-m+1}^{N-1})$$

m-gram language model

 Further optimized using smoothing and interpolation with lower-order Ngram models

#### Decoding: Search

$$W^* = \underset{W}{\operatorname{arg\,max}} \operatorname{Pr}(O_A|W) \operatorname{Pr}(W)$$

$$W^* = \underset{w_1^N, N}{\operatorname{arg\,max}} \left\{ \left[ \prod_{n=1}^N \Pr(w_n | w_{n-m+1}^{n-1}) \right] \left[ \sum_{q_1^T, w_1^N} \prod_{t=1}^T \Pr(O_t | q_t, w_1^N) \Pr(q_t | q_{t-1}, w_1^N) \right] \right\}$$

$$\begin{array}{l} \text{Viterbi} \\ \approx \underset{w_1^N, N}{\operatorname{arg\,max}} \left\{ \left[ \prod_{n=1}^N \Pr(w_n | w_{n-m+1}^{n-1}) \right] \left[ \underset{q_1^T, w_1^N}{\max} \prod_{t=1}^T \Pr(O_t | q_t, w_1^N) \Pr(q_t | q_{t-1}, w_1^N) \right] \right\} \end{array}$$

- Viterbi approximation divides the above optimisation problem into subproblems that allows the efficient application of dynamic programming
- Search space still very huge for LVCSR tasks! Use approximate decoding techniques (A\* decoding, beam-width decoding, etc.) to visit only promising parts of the search space

#### ASR Search



## Kaldi Toolkit & Assignment 2