

Instructor: Preethi Jyothi

#### (R)NN-based Language Models

#### Lecture 12



#### Word representations in Ngram models

- discrete space involving the vocabulary
- unseen Ngrams
- space?

In standard Ngram models, words are represented in the

Limits the possibility of truly interpolating probabilities of

Can we build a representation for words in the continuous

#### Word representations

- 1-hot representation: •
  - being 1
- similarity
- E.g. dog  $\rightarrow$  {-0.02, -0.37, 0.26, 0.25, -0.11, 0.34}

• Each word is given an index in  $\{1, \ldots, V\}$ . The 1-hot vector  $f_i \in R^v$  contains zeros everywhere except for the i<sup>th</sup> dimension

1-hot form, however, doesn't encode information about word

Distributed (or continuous) representation: Each word is associated with a dense vector. Based on the "distributional hypothesis".

#### Word embeddings

- also referred to as "word embeddings"
  - Low dimensional
  - Similar words will have similar vectors
- properties (glad is similar to gladly, etc.)

These distributed representations in a continuous space are

• Word embeddings capture semantic properties (such as man is to woman as boy is to girl, etc.) and morphological

#### Word embeddings

FRANCE	JESUS	XBOX	REDDISH	SCRATCHED	MEGABITS
AUSTRIA	GOD	AMIGA	GREENISH	NAILED	OCTETS
BELGIUM	SATI	PLAYSTATION	BLUISH	SMASHED	MB/S
GERMANY	CHRIST	MSX	PINKISH	PUNCHED	BIT/S
ITALY	SATAN	IPOD	PURPLISH	POPPED	BAUD
GREECE	KALI	SEGA	BROWNISH	CRIMPED	CARATS
SWEDEN	INDRA	PSNUMBER	GREYISH	SCRAPED	$_{\rm KBIT/S}$
NORWAY	VISHNU	HD	GRAYISH	SCREWED	MEGAHERTZ
EUROPE	ANANDA	DREAMCAST	WHITISH	SECTIONED	MEGAPIXELS
HUNGARY	PARVATI	GEFORCE	SILVERY	SLASHED	$_{\rm GBIT}/{\rm S}$
SWITZERLAND	GRACE	CAPCOM	YELLOWISH	RIPPED	AMPERES

#### **Relationships learned from embeddings**

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

#### **Bilingual embeddings**



#### Word embeddings

- These distributed representations in a continuous space are also referred to as "word embeddings"
  - Low dimensional
  - Similar words will have similar vectors
- Word embeddings capture semantic properties (such as man is to woman as boy is to girl, etc.) and morphological properties (glad is similar to gladly, etc.)
- The word embeddings could be learned via the first layer of a neural network [B03].

### Word embeddings





- Embedding layer •
- One or more middle/hidden layers ullet
- Softmax output layer



## **Continuous space language models**



#### **NN language model**

- Project all the words of the context  $h_j = w_{j-n+1}, ..., w_{j-1}$  to their dense forms
- Then, calculate the language model probability  $Pr(w_j = il h_j)$  for the given context h<sub>i</sub>



discrete continuous representation: representation: indices in wordlist P dimensional vectors LM probabilities for all words

#### **NN language model**

- Dense vectors of all the words in context are concatenated forming the first hidden layer of the neural network
- Second hidden layer:  $\bullet$

$$d_j = tanh(\Sigma m_{jl}c_l + b_j) \forall j = 1, ..., H w$$

Output layer:

$$o_i = \sum v_{ij}d_j + b'_i \quad \forall i = 1, \dots, N$$

→ softmax output from the ith pi neuron  $\rightarrow Pr(w_j = i | h_j)$ 



#### NN language model

$$L = \sum_{i=1}^{N} t_i \log p_i + \epsilon \left( \sum_{kl} m_{kl}^2 + \sum_{ik} v_{ik}^2 \right)$$

- the training instance, 0 elsewhere)
- the distribution estimated by the NN
- Second part: Regularization term

Model is trained to minimise the following loss function:

Here, t<sub>i</sub> is the target output 1-hot vector (1 for next word in

First part: Cross-entropy between the target distribution and

## **Decoding with NN LMs**

- large vocabulary ASR systems:
  - 1. Lattice rescoring
  - 2. Shortlists

Two main techniques used to make the NN LM tractable for

#### Use NN language model via lattice rescoring



- Lattice Graph of possible word s
   Ngram backoff LM
- Each lattice arc has both acoustic/language model scores.
- LM scores on the arcs are replaced by scores from the NN LM

• Lattice — Graph of possible word sequences from the ASR system using an

anguage model scores. I by scores from the NN LM

## **Decoding with NN LMs**

# large vocabulary ASR systems:

1. Lattice rescoring

2. Shortlists

Two main techniques used to make the NN LM tractable for

- Softmax normalization of the output layer is an expensive operation esp. for large vocabularies
- Solution: Limit the output to the s most frequent words. •
  - LM probabilities of words in the short-list are calculated by the NN
  - LM probabilities of the remaining words are from Ngram • backoff models

#### Shortlist

#### **Results**

Table 3 Perplexities on the 200	3 eva	luation data for the back-off	and the hyl	orid LM as a function of the size of the	CTS training data
CTS corpus (words)		7.2 N	[	12.3 M	27.3 M
<i>In-domain data only</i> Back-off LM Hybrid LM		62.4 57.0		55.9 50.6	50.1 <b>45.5</b>
Interpolated with all da Back-off LM Hybrid LM	ita	53.0 50.8		51.1 48.0	<b>47.5</b> 44.2
Eval03 word error rate	<ul> <li>28</li> <li>26</li> <li>24</li> <li>22</li> <li>20</li> <li>18</li> </ul>	System 1 25.27% 24.51% 24.51% 23.70% 23.04% 22.19% 7.2M	ystem 2 22.329 21.779	backoff LM, CTS data hybrid LM, CTS data backoff LM, CTS+BN data hybrid LM, CTS+BN data	
		in-dom	ain LM tr	aining corpus size	



[S07]: Schwenk et al., "Continuous space language models", CSL, 07

#### word2vec (to learn word embeddings)





#### **Bias in word embeddings**

#### Gender bias in word embeddings

When you envision a nurse, a woman most likely pops into your mind. If you imagine an accomplished executive, on the other hand, it's quite likely you're thinking about a man.

It's not just you, though. The machine learning algorithms that target ads at us, prune our search results, or sort resumes for recruiters are all plagued by gendered stereotypes.

Algorithms that model natural language transform words into vectors, and similar words should be near each other in this vector space. Unfortunately, our models have learned to capture the biases present in the real-life data on which we train them. In word embedding space, for example, the relationship between "he" and "she" mirrors that of "programmer" and "homemaker". When we train our machine learning models on embeddings like these, a recruiter searching for "programmers" will leave female resumes at the bottom of the pile.

The following visualization allows you to define a category and see how words in that category relate to gender, by projecting them onto the axis representing gender in word embedding space. It is intended to encourage you to think critically about the tools you use and to consider carefully before treating anything as a black box. Feel free to explore the categories, or create your own.

Choose your word category...

Personality

Or type your own words...

charming, accomplished, ambitious, impressive





#### Longer word context?

- compute an Ngram probability  $Pr(w_i = i|h_i)$  (where  $h_i$ ) encodes the Ngram history)
- We know Ngrams are limiting:
- networks (RNNs)

```
What have we seen so far: A feedforward NN used to
```

Alice who had attempted the assignment asked the lecturer

• How can we predict the next word based on the entire sequence of preceding words? Use recurrent neural

#### Simple RNN language model



• Current word, x<sub>t</sub> Hidden state, st Output, yt

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

 RNN is trained using the cross-entropy criterion

#### **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



#### LSTM-LMs

•

output layer

2nd hidden layer

projection layer

input layer

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

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



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

#### **Character-based RNN-LMs**



Image from: <u>http://karpathy.github.io/2015/05/21/rnn-effectiveness/</u> Good tutorial available at https://github.com/yunjey/pytorch-tutorial/blob/master/tutorials/02-intermediate/language\_model/main.py#L30-L50





# 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.

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

#### Generate text using an LM trained on Obama speeches

Good morning. One of the borders will be able to continue to be here today. We have to say that the partnership was a partnership with the American people and the street continually progress that is a process and distant lasting peace and support that they were supporting the work of concern in the world. They were in the streets and communities that could have to provide steps to the people of the United States and Afghanistan. In the streets — the final decade of the country that will include the people of the United States of America. Now, humanitarian crisis has already rightly achieved the first American future in the same financial crisis that they can find reason to invest in the world.

Thank you very much. God bless you. God bless you. Thank you.

#### **NN trained on Trump's speeches (now defunct)**



I'm a Neural Network trained on Trump's transcripts. Priming text in []s. Donate (gofundme.com/deepdrumpf) to interact! Created by @hayesbh.

Joined March 2016

....



 $\sim$ When I have to build a hotel, we're bombing the hell out of them. Lots of money. To those suffering, I say vote for Donald. #SyriaStrikes

Photos and videos m/profile\_images/705464735353991168/d4eBpkKr\_400x400.jpg

 $\bigcirc$ 

4

wing	Followers	Likes
7	25.9K	19

Usually that's a bad sign of things to come.

1 36 125

#### DeepDrumpf @DeepDrumpf · 7 Apr 2017

M 172 1 62 CA



### **Common RNNLM training tricks**

- SGD fares very well on this task (compared to other optimizers like) Adagrad, Adam, etc.).
- Use dropout regularization
- Truncated BPTT
- Use mini batches to aggregate gradients during training
  - In batched RNNLMs, process multiple sentences at the same time •
  - Handle variable length sequences using padding and masking
  - To be judicious about padding, sort the sentences in the corpus by length before creating batches

### **Spotlight: Regularizing and Optimizing LSTM Language** Models (Merity et al. 2018)

- No special model, just better regularisation + optimization
- Dropout on recurrent connections and embeddings •
- SGD w/ averaging triggered when model is close to • convergence
- Weight tying between embedding and softmax layers •
- Reduced embedding sizes
- https://github.com/salesforce/awd-lstm-lm

### **Spotlight: On the State of the art of Evaluation** in Neural Language Models (Melis et al., 2018)



