

CS388: Natural Language Processing

Lecture 9: Language Modeling + Pretraining

Greg Durrett



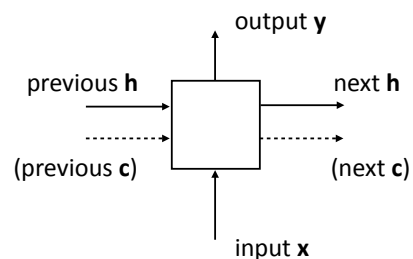
Administrivia

- ▶ Mini 2 out tonight (due Tuesday, October 8)
- ▶ Final project spec out next week

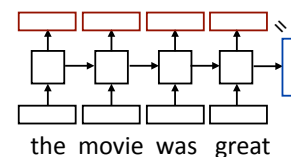


Recall: RNNs

- ▶ Cell that takes some input \mathbf{x} , has some hidden state \mathbf{h} , and updates that hidden state and produces output \mathbf{y} (all vector-valued)



Recall: RNN Abstraction



- ▶ **Encoding of the sentence** — can pass this a decoder or make a classification decision about the sentence
- ▶ **Encoding of each word** — can pass this to another layer to make a prediction (can also pool these to get a different sentence encoding)
- ▶ RNN can be viewed as a transformation of a sequence of vectors into a sequence of context-dependent vectors



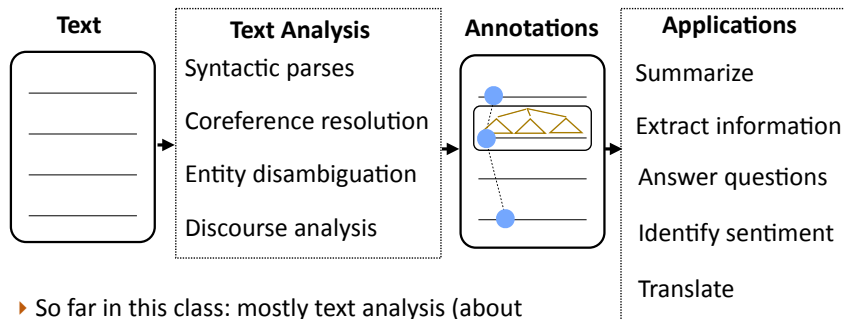
This Lecture

- ▶ Language modeling
 - ▶ N-gram models
 - ▶ Neural LMs
- ▶ LM-based pretraining: ELMo

Language Modeling



NLP Analysis Pipeline



- ▶ So far in this class: mostly text analysis (about tagging, classifying, etc. the structure of text)
- ▶ Haven't talked about text *generation* tasks



Challenges in Text Generation

- ▶ Dialogue, machine translation, summarization, etc.
 - ▶ What to say (content selection + content planning) and how to say it
- ▶ Template-based generation systems always generate fluent output
- ▶ For learned systems, how do we make sure language is plausible?
- ▶ Language models: place a distribution $P(\mathbf{w})$ over strings \mathbf{w} in a language
 - ▶ Next week: $P(\mathbf{T}, \mathbf{w})$ modeled by probabilistic context-free grammars
 - ▶ Today: autoregressive models $P(\mathbf{w}) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \dots$



N-gram Language Models

$$P(\mathbf{w}) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \dots$$

- n-gram models: distribution of next word is a multinomial conditioned on previous n-1 words $P(w_i|w_1, \dots, w_{i-1}) = P(w_i|w_{i-n+1}, \dots, w_{i-1})$

I visited San _____ put a distribution over the next word

$$P(w|\text{visited San}) = \frac{\text{count}(\text{visited San}, w)}{\text{count}(\text{visited San})}$$

Maximum likelihood estimate of this 3-gram probability from a corpus

- Just relies on counts, even in 2008 could scale up to 1.3M word types, 4B n-grams (all 5-grams occurring >40 times on the Web)



Smoothing N-gram Language Models

- What happens when we scale to longer contexts?

$P(w|\text{to})$ *to* occurs 10M times in corpus

$P(w|\text{go to})$ *go to* occurs 100,000 times in corpus

$P(w|\text{to go to})$ *go to* occurs 10,000 times in corpus

$P(w|\text{want to go to})$ *want to go to*: only 100 occurrences

- Probability counts get very sparse, and we often want information from 5+ words away



Smoothing N-gram Language Models

I visited San _____ put a distribution over the next word

- Smoothing is very important, particularly when using 4+ gram models

$$P(w|\text{visited San}) = (1 - \lambda) \frac{\text{count}(\text{visited San}, w)}{\text{count}(\text{visited San})} + \lambda \frac{\text{count}(\text{San}, w)}{\text{count}(\text{San})} \quad \leftarrow \begin{array}{l} \text{smooth} \\ \text{this} \\ \text{too!} \end{array}$$

- One technique is “absolute discounting:” subtract off constant k from numerator, set lambda to make this normalize ($k=1$ is like leave-one-out)

$$P(w|\text{visited San}) = \frac{\text{count}(\text{visited San}, w) - k}{\text{count}(\text{visited San})} + \lambda \frac{\text{count}(\text{San}, w)}{\text{count}(\text{San})}$$



Kneser-Ney Smoothing

- Suppose context = *John visited Madagascar* and *Madagascar* has count 0

- Absolute discounting:

$$P(w|\text{John visited Madagascar}) = \lambda_1 \frac{\text{count}(\text{John visited Madagascar}, w)}{\text{count}(\text{John visited Madagascar})} + \dots + \lambda_4 \frac{\text{count}(w)}{\text{count}(\cdot)}$$

- First terms all end up being 0/0, so we back off to unigram probability

- What will the highest probability word be here?

- Instead: use probability based on number of unique contexts w occurs in (type counts):

$$P(w) = \frac{|\{p : \text{count}(p, w) > 0\}|}{\sum_w |\{p : \text{count}(p, w) > 0\}|}$$

- Kneser-Ney (1994): absolute discounting w/ type counts for lower-order probs



Engineering N-gram Models

- For 5+-gram models, need to store between 100M and 10B context-word-count triples

| (a) Context-Encoding | | | (b) Context Deltas | | | (c) Bits Required | | |
|----------------------|----------|-------|--------------------|------------|-------|-------------------|--------------|---------|
| w | c | val | Δw | Δc | val | $ \Delta w $ | $ \Delta c $ | $ val $ |
| 1933 | 15176585 | 3 | 1933 | 15176585 | 3 | 24 | 40 | 3 |
| 1933 | 15176587 | 2 | +0 | +2 | 1 | 2 | 3 | 3 |
| 1933 | 15176593 | 1 | +0 | +5 | 1 | 2 | 3 | 3 |
| 1933 | 15176613 | 8 | +0 | +40 | 8 | 2 | 9 | 6 |
| 1933 | 15179801 | 1 | +0 | +188 | 1 | 2 | 12 | 3 |
| 1935 | 15176585 | 298 | +2 | 15176585 | 298 | 4 | 36 | 15 |
| 1935 | 15176589 | 1 | +0 | +4 | 1 | 2 | 6 | 3 |

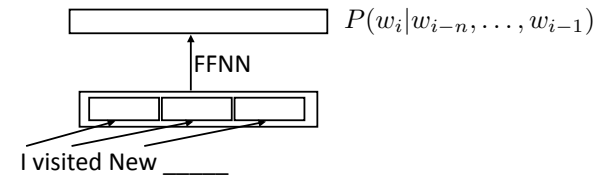
- Make it fit in memory by *delta encoding* scheme: store deltas instead of values and use variable-length encoding

Pauls and Klein (2011), Heafield (2011)



Neural Language Models

- Early work: feedforward neural networks looking at context

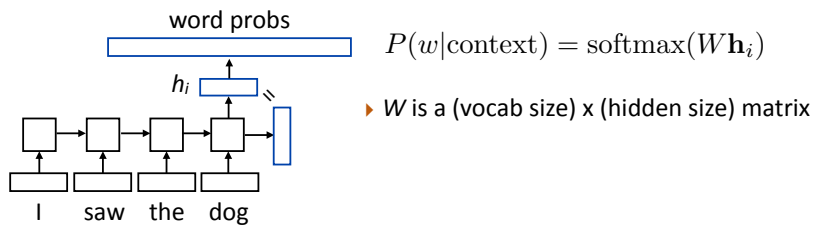


- Slow to train over lots of data!
- Still only look at a fixed window of information...can we use more?

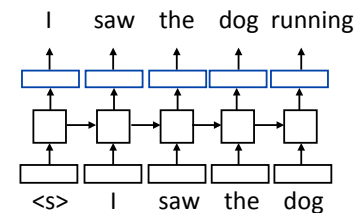
Mnih and Hinton (2003)



RNN Language Modeling



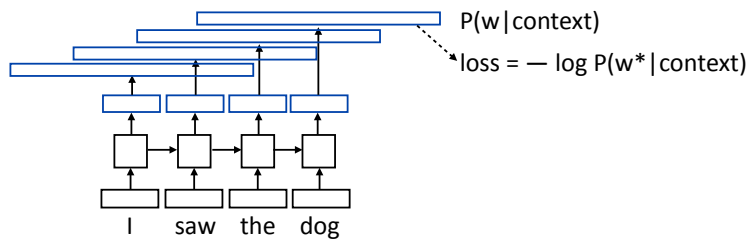
Training RNNLMs



- Input is a sequence of words, output is those words shifted by one,
- Allows us to efficiently batch up training across time (one run of the RNN)



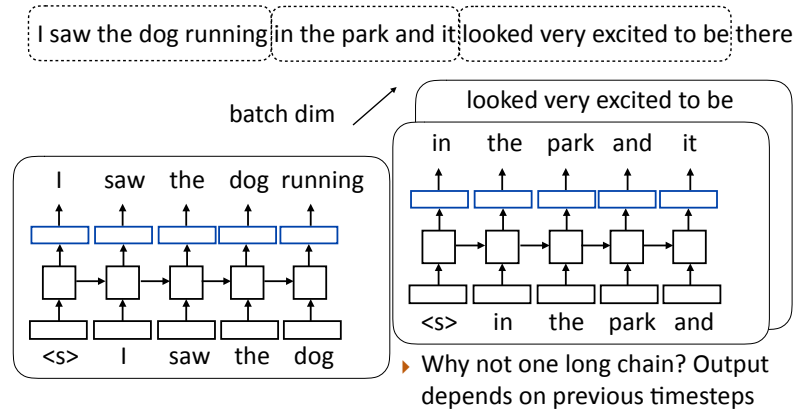
Training RNNLMs



- ▶ Total loss = sum of negative log likelihoods at each position
- ▶ Backpropagate through the network to simultaneously learn to predict next word given previous words at all positions



Batched LM Training



LM Evaluation

- ▶ Accuracy doesn't make sense — predicting the next word is generally impossible so accuracy values would be very low
- ▶ Evaluate LMs on the likelihood of held-out data (averaged to normalize for length)

$$\frac{1}{n} \sum_{i=1}^n \log P(w_i | w_1, \dots, w_{i-1})$$
- ▶ Perplexity: $\exp(\text{average negative log likelihood})$. Lower is better
 - ▶ Suppose we have probs 1/4, 1/3, 1/4, 1/3 for 4 predictions
 - ▶ Avg NLL (base e) = 1.242 Perplexity = 3.464



Results

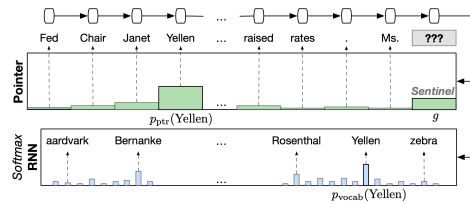
- ▶ Evaluate on Penn Treebank: small dataset (1M words) compared to what's used in MT, but common benchmark
- ▶ Kneser-Ney 5-gram model with cache: PPL = 125.7
- ▶ LSTM: PPL ~ 60-80 (depending on how much you optimize it)
- ▶ Melis et al.: many neural LM improvements from 2014-2017 are subsumed by just using the right regularization (right dropout settings). So LSTMs are pretty good
 - ▶ Main tricks: changing some tricky dropout settings + how params are structured (sizes, etc.)

Merity et al. (2017), Melis et al. (2017)



Limitations of LSTM LMs

- ▶ Need some kind of pointing mechanism to repeat recent words



Merity et al. (2016)

$$p(\text{Yellen}) = g(p_{\text{vocab}}(\text{Yellen})) + (1 - g)p_{\text{ptr}}(\text{Yellen})$$

- ▶ Transformers can do this (will discuss later in the course)
- ▶ SOTA: GPT-2 transformer trained on 40GB of English text.
PTB perplexity = 65.85 with 117M params => 35.76 W/ 1542M params



Applications of Language Modeling

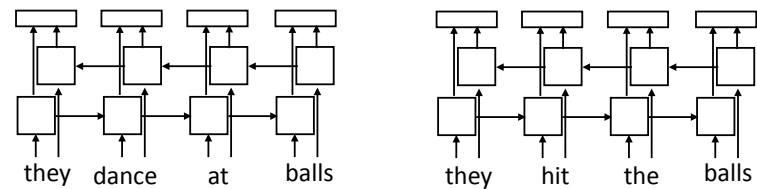
- ▶ All generation tasks: translation, dialogue, text simplification, paraphrasing, etc.
- ▶ Grammatical error correction
- ▶ Predictive text
- ▶ Pretraining!

Pretraining / ELMo



Recall: Context-dependent Embeddings

- ▶ How to handle different word senses? One vector for *balls*



- ▶ Train a neural language model to predict the next word given previous words in the sentence, use its internal representations as word vectors

Peters et al. (2018)



ELMo

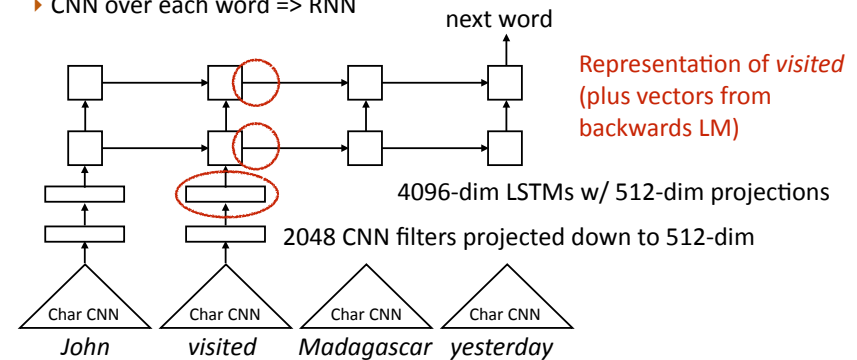
- ▶ Key idea: language models can allow us to form useful word representations in the same way word2vec did
- ▶ Take a powerful language model, train it on large amounts of data, then use those representations in downstream tasks
- ▶ What do we want our LM to look like?

Peters et al. (2018)



ELMo

- ▶ CNN over each word => RNN

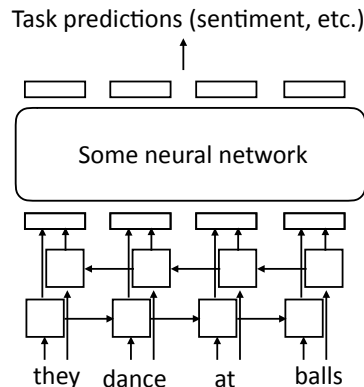


Peters et al. (2018)



How to apply ELMo?

- ▶ Take those embeddings and feed them into whatever architecture you want to use for your task
- ▶ *Frozen* embeddings: update the weights of your network but keep ELMo's parameters frozen
- ▶ *Fine-tuning*: backpropagate all the way into ELMo when training your model



Peters, Ruder, Smith (2019)



Results: Frozen ELMo

| TASK | PREVIOUS SOTA | | OUR BASELINE | ELMo + BASELINE | INCREASE (ABSOLUTE/RELATIVE) |
|-------|----------------------|--------------|--------------|-----------------|------------------------------|
| SQuAD | Liu et al. (2017) | 84.4 | 81.1 | 85.8 | 4.7 / 24.9% |
| SNLI | Chen et al. (2017) | 88.6 | 88.0 | 88.7 ± 0.17 | 0.7 / 5.8% |
| SRL | He et al. (2017) | 81.7 | 81.4 | 84.6 | 3.2 / 17.2% |
| Coref | Lee et al. (2017) | 67.2 | 67.2 | 70.4 | 3.2 / 9.8% |
| NER | Peters et al. (2017) | 91.93 ± 0.19 | 90.15 | 92.22 ± 0.10 | 2.06 / 21% |
| SST-5 | McCann et al. (2017) | 53.7 | 51.4 | 54.7 ± 0.5 | 3.3 / 6.8% |

- ▶ Massive improvements across 5 benchmark datasets: question answering, natural language inference, semantic role labeling (discussed later in the course), coreference resolution, named entity recognition, and sentiment analysis



How to apply ELMo?

| Pretraining | Adaptation | NER CoNLL 2003 | SA SST-2 | Nat. lang. MNLI | inference SICK-E | Semantic textual similarity SICK-R | MRPC | STS-B |
|---------------|------------|-------------------|-------------|--------------------|---------------------|---------------------------------------|------|-------|
| Skip-thoughts | ❄️ | - | 81.8 | 62.9 | - | 86.6 | 75.8 | 71.8 |
| ELMo | 🔥 | 91.7 | 91.8 | 79.6 | 86.3 | 86.1 | 76.0 | 75.9 |
| | Δ=🔥-❄️ | 91.9 | 91.2 | 76.4 | 83.3 | 83.3 | 74.7 | 75.5 |
| | | 0.2 | -0.6 | -3.2 | -3.3 | -2.8 | -1.3 | -0.4 |

► How does frozen (❄️) vs. fine-tuned (🔥) compare?

► Recommendations:

| Conditions | | | Guidelines |
|------------|--------|--------------|--|
| Pretrain | Adapt. | Task | |
| Any | ❄️ | Any | Add many task parameters |
| Any | 🔥 | Any | Add minimal task parameters ⚠️ Hyper-parameters |
| Any | Any | Seq. / clas. | ❄️ and 🔥 have similar performance |
| ELMo | Any | Sent. pair | use ❄️ |
| BERT | Any | Sent. pair | use 🔥 |

Peters, Ruder, Smith (2019)



Why is language modeling a good objective?

- “Impossible” problem but bigger models seem to do better and better at distributional modeling (no upper limit yet)
- Successfully predicting next words requires modeling lots of different effects in text

Context: My wife refused to allow me to come to Hong Kong when the plague was at its height and –” “Your wife, Johanne? You are married at last ?” Johanne grinned. “Well, when a man gets to my age, he starts to need a few home comforts.

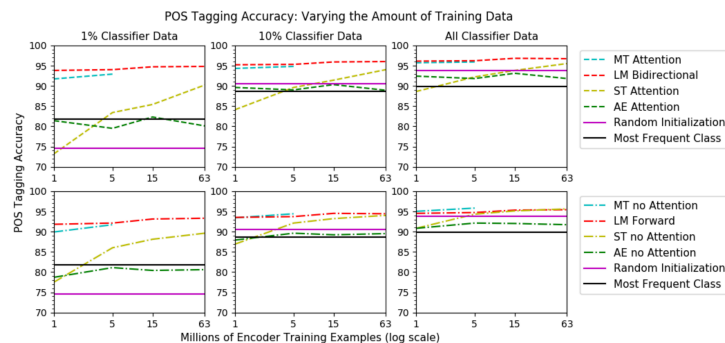
Target sentence: After my dear mother passed away ten years ago now, I became

Target word: lonely

- LAMBADA dataset (Papernot et al., 2016): explicitly targets world knowledge and very challenging LM examples
- Coreference, Winograd schema, and much more



Why is language modeling a good objective?



Zhang and Bowman (2018)



Why did this take time to catch on?

- Earlier version of ELMo by the same authors in 2017, but it was only evaluated on tagging tasks, gains were 1% or less
- Required: training on lots of data, having the right architecture, significant hyperparameter tuning



Probing ELMo

- From each layer of the ELMo model, attempt to predict something: POS tags, word senses, etc.
- Higher accuracy => ELMo is capturing that thing more nicely

| Model | F ₁ |
|----------------------------|----------------|
| WordNet 1st Sense Baseline | 65.9 |
| Raganato et al. (2017a) | 69.9 |
| Iacobacci et al. (2016) | 70.1 |
| CoVe, First Layer | 59.4 |
| CoVe, Second Layer | 64.7 |
| biLM, First layer | 67.4 |
| biLM, Second layer | 69.0 |

Table 5: All-words fine grained WSD F₁. For CoVe and the biLM, we report scores for both the first and second layer biLSTMs.

| Model | Acc. |
|-------------------------|-------------|
| Collobert et al. (2011) | 97.3 |
| Ma and Hovy (2016) | 97.6 |
| Ling et al. (2015) | 97.8 |
| CoVe, First Layer | 93.3 |
| CoVe, Second Layer | 92.8 |
| biLM, First Layer | 97.3 |
| biLM, Second Layer | 96.8 |

Table 6: Test set POS tagging accuracies for PTB. For CoVe and the biLM, we report scores for both the first and second layer biLSTMs.



Analysis

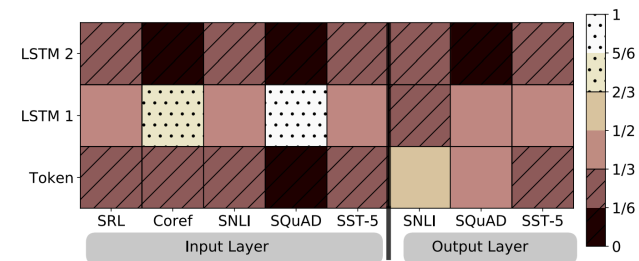


Figure 2: Visualization of softmax normalized biLM layer weights across tasks and ELMo locations. Normalized weights less than 1/3 are hatched with horizontal lines and those greater than 2/3 are speckled.



Takeaways

- Language modeling involves predicting the next word given context. Several techniques to do this, more later in the course
- Learning a neural network to do this induces useful representations for other tasks, similar to word2vec/GloVe
- Next class: interpreting neural network models