# Language Models and Recurrent Neural Networks

COSC 7336: Advanced Natural Language Processing
Fall 2017

# Announcements

★ Next Week no class (Mid-Semester Bash)
★ Assignment 2 is out
★ Mid term exam Oct. 27th
★ Paper presentation sign up coming up
★ Reminder to submit report, code and notebook for Assignment 1 today!

# Today's lecture

★ Short intro to language models

★ Recurrent Neural Networks

★ Demo: RNN for text

# Language Models

★ They assign a probability to a sequence of words:
- Machine Translation:
  - P(**high** winds tonite) > P(**large** winds tonite)
- Spell Correction:
  - The office is about fifteen **minuets** from my house
  - P(about fifteen **minutes** from) > P(about fifteen **minuets** from)
- Speech Recognition
  - P(I saw a van) >> P(eyes awe of an)
- Summarization, question-answering, OCR correction and many more!

# More Formally

★ Given a sequence of words predict the next one:
  ○ $P(w_5|w_1,w_2,w_3,w_4)$
★ Predict the likelihood of a sequence of words:
  ○ $P(W) = P(w_1,w_2,w_3,w_4,w_5 \ldots w_n)$
★ How do we compute these?

# Chain Rule

★ Recall the definition of conditional probabilities:

    ○ **p(B|A) = P(A,B)/P(A)**   Rewriting:   **P(A,B) = P(A)P(B|A)**

★ More variables:

    ○ P(A,B,C,D) = P(A)P(B|A)P(C|A,B)P(D|A,B,C)

★ The Chain Rule in General

    ○ $P(x_1,x_2,x_3,\ldots,x_n) = P(x_1)P(x_2|x_1)P(x_3|x_1,x_2)\ldots P(x_n|x_1,\ldots,x_{n-1})$

# Back to sequence of words

P("I am the fire that burns against the cold") = P(I) x P(am|I) x P(the|I am) x P(fire|I am the) x P(that| I am the fire) x P(burns| I am the fire that) x P(against| I am the fire that burns) x P(the| I am the fire that burns against) x P(cold| I am the fire that burns against the)

★ How do we estimate these probabilities?

*count*(I am the fire that burns against the cold)

*count*(I am the fire that burns against the)

# We shorten the context (history)

**Markov Assumption:**

P(cold | I am the fire that burns against the) ≈ P (cold | burns against the)

This is: $P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-N+1}^{n-1})$

When N = 1, this is a unigram language model:

$$P(w_1 w_2 \ldots w_n) \approx \prod_i P(w_i)$$

When k =2, this is a bigram language model:

$$P(w_1^n) \approx \prod_{k=1}^{n} P(w_k | w_{k-1})$$

# Are all our problems solved?

# Zeros

Training set:

… denied the allegations
… denied the reports
… denied the claims
… denied the request

P("offer" | denied the) = 0

Test set:

… denied the offer
… denied the loan

# Laplace Smoothing

★ Add one to all counts

★ Unigram counts:

$$P(w_i) = \frac{c_i}{N}$$

★ Laplace counts:

$$P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V}$$

★ Disadvantages:
  ○ Drastic change in probability mass

★ But for some tasks Laplace smoothing is a reasonable choice

# Leveraging Hierarchy of N-grams: Backoff

★ Adding one to all counts is too drastic

★ What about relying on shorter contexts?

  ○ Let's say we're tying to compute P(against | that burns), but count(that burns against) = 0

  ○ We backoff to a shorter context: P(against | that burns) ≈ P(against | burns)

★ We backoff to a lower n-gram

★ Katz Backoff (discounted backoff)

$$P_{\text{BO}}(w_n|w_{n-N+1}^{n-1}) = \begin{cases} P^*(w_n|w_{n-N+1}^{n-1}), & \text{if } C(w_{n-N+1}^n) > 0 \\ \alpha(w_{n-N+1}^{n-1})P_{\text{BO}}(w_n|w_{n-N+2}^{n-1}), & \text{otherwise.} \end{cases}$$

# Leveraging Hierarchy of N-grams: Interpolation

★ Better idea: why not always rely on lower order N-grams?

$$\hat{P}(w_n|w_{n-2}w_{n-1}) = \lambda_1 P(w_n|w_{n-2}w_{n-1})$$
$$+\lambda_2 P(w_n|w_{n-1})$$
$$+\lambda_3 P(w_n)$$

Subject to: $\sum_i \lambda_i = 1$

★ Even better, condition on context:

$$\hat{P}(w_n|w_{n-2}w_{n-1}) = \lambda_1(w_{n-2}^{n-1})P(w_n|w_{n-2}w_{n-1})$$
$$+\lambda_2(w_{n-2}^{n-1})P(w_n|w_{n-1})$$
$$+\lambda_3(w_{n-2}^{n-1})P(w_n)$$

# Absolute Discounting

| Bigram count in training | Bigram count in held out set |
|---|---|
| 0 | .0000270 |
| 1 | 0.448 |
| 2 | 1.25 |
| 3 | 2.24 |
| 4 | 3.23 |
| 5 | 4.21 |
| 6 | 5.23 |
| 7 | 6.21 |
| 8 | 7.21 |
| 9 | 8.26 |

For each bigram count in training data, what's the count in held out set?

Approx. a 0.75 difference!

UNIVERSITY of
HOUSTON

UNIVERSIDAD
NACIONAL
DE COLOMBIA

# Absolute Discounting

$$P_{\text{AbsoluteDiscounting}}(w_i|w_{i-1}) = \frac{C(w_{i-1}w_i) - d}{\sum_v C(w_{i-1}v)} + \lambda(w_{i-1})P(w_i)$$

# How much do we want to trust unigrams?

★ Instead of P(w): "How likely is w"

★ $P_{continuation}(w)$: "How likely is w to appear as a novel continuation?

★ For each word, count the number of bigram types it completes

$$P_{\text{CONTINUATION}}(w) = \frac{|\{v : C(vw) > 0\}|}{|\{(u', w') : C(u'w') > 0\}|}$$

# Interpolated Kneser-Ney

★ **Intuition:** Use estimate of $P_{\text{CONTINUATION}}(w_i)$

$$P_{\text{KN}}(w_i|w_{i-n+1}^{i-1}) = \frac{\max(c_{KN}(w_{i-n+1}^i) - d, 0)}{\sum_v c_{KN}(w_{i-n+1}^{i-1}v)} + \lambda(w_{i-n+1}^{i-1})P_{KN}(w_i|w_{i-n+2}^{i-1})$$

# Evaluating Language Models

★ Ideal: Evaluate on end task (extrinsic)

★ Intrinsic evaluation: use perplexity:

$$\mathrm{PP}(W) \;=\; \sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i|w_1 \dots w_{i-1})}}$$

★ Perplexity is inversely proportional to the probability of W

Recurrent Neural Networks

# Recurrent NNs

★ Neural networks with memory
★ Feed-forward NN: output exclusively depends on the current input
★ Recurrent NN: output depends on current and previous states
★ This is accomplished through lateral/backward connections which carry information while processing a sequence of inputs



(source: http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

# Character-level language model



(source: http://karpathy.github.io/2015/05/21/rnn-effectiveness/)

# Sequence learning alternatives

(source: http://karpathy.github.io/2015/05/21/rnn-effectiveness/)

# Network unrolling

# Backpropagation through time (BPTT)

# BPTT is hard

★ The vanishing and the exploding gradient problem

★ Gradients could vanish (or explode) when propagated several steps back

★ This makes difficult to learn long-term dependencies.



Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. On the difficulty of training Recurrent Neural Networks. Proc. of ICML, abs/1211.5063.

# Long term dependencies

# Long short-term memory (LSTM)

★   LSTM networks solve the problem of long-term dependency problem.
★   They use gates that allow to keep memory through long sequences and be updated only when required.
★   Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." Neural computation 9, no. 8 (1997): 1735-1780.

# Conventional RNN vs LSTM

# Forget Gate

★ Controls the flow of the previous internal state Ct-1

★ $f_t=1 \Rightarrow$ keep previous state

★ $f_t=0 \Rightarrow$ forget previous state

(source: http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

UNIVERSIDAD
NACIONAL
DE COLOMBIA

# Input Gate

★ Controls the flow of the input state xt

★ $i_t = 1 \Rightarrow$ take input into account

★ $i_t = 0 \Rightarrow$ ignore input

(source: http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

# Current state calculation



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

(source: http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

HOUSTON

NACIONAL
DE COLOMBIA

# Output Gate

★ Controls the flow of information from the internal state $x_t$ to the outside $h_t$

★ $o_t=1 \Rightarrow$ allows internal state out

★ $o_t=0 \Rightarrow$ doesn't allow internal state out



(source: http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

# Gated Recurrent Unit (GRU)



$$z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$$

$$r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$$

$$\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using rnn encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078.

# The Unreasonable Effectiveness of Recurrent Neural Networks

★ Famous [blog](#) entry from Andrej Karpathy (UofS)

★ Character-level language models based on multi-layer LSTMs.

★ Data:
  ○ Shakspare plays
  ○ Wikipedia
  ○ LaTeX
  ○ Linux source code

# Algebraic geometry book in LaTeX





UNIVERSITY of HOUSTON

UNIVERSIDAD NACIONAL DE COLOMBIA

# Linux source code

```c
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
  int error;
  if (fd == MARN_EPT) {
    /*
     * The kernel blank will coeld it to userspace.
     */
    if (ss->segment < mem_total)
      unblock_graph_and_set_blocked();
    else
      ret = 1;
    goto bail;
  }
  segaddr = in_SB(in.addr);
  selector = seg / 16;
  setup_works = true;
  for (i = 0; i < blocks; i++) {
    seq = buf[i++];
    bpf = bd->bd.next + i * search;
    if (fd) {
      current = blocked;
```

# LSTM language model demo