

Language models

CS 685, Spring 2021

Advanced Topics in Natural Language Processing

<http://brenocon.com/cs685>

https://people.cs.umass.edu/~brenocon/cs685_s21/

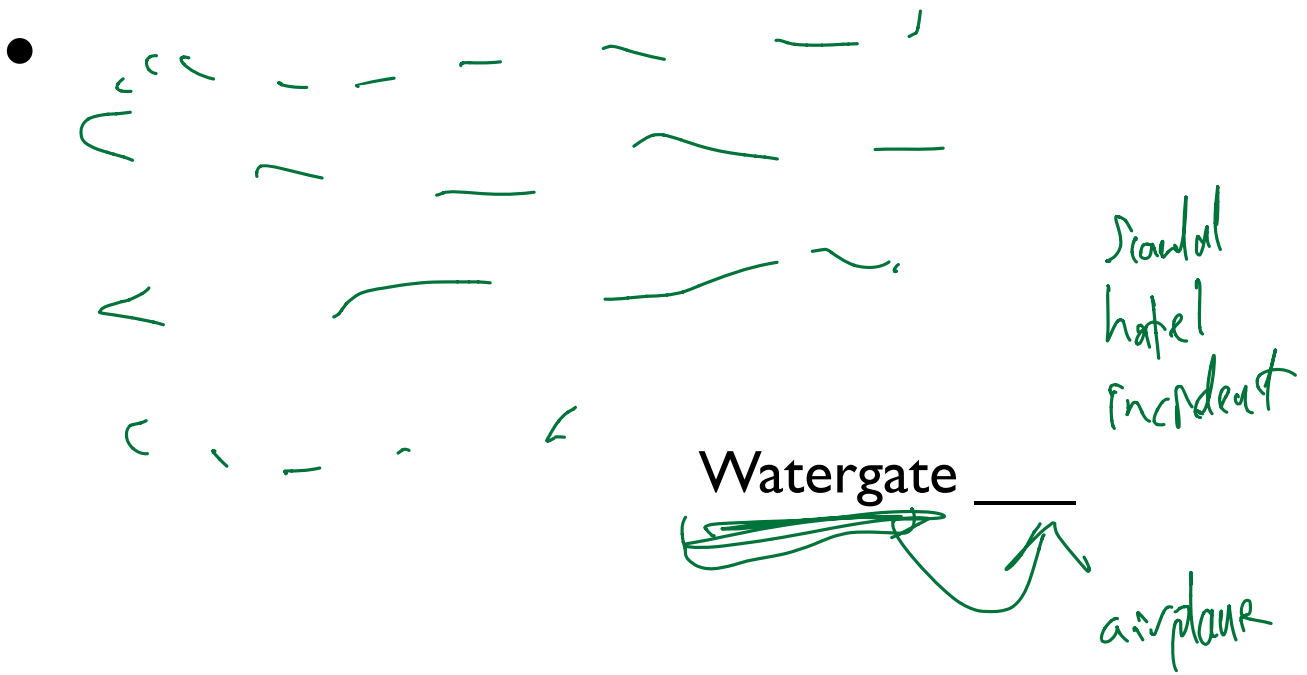
Brendan O'Connor

College of Information and Computer Sciences

University of Massachusetts Amherst

Announcements

- Hope HW1 is going well!
- Reading review #2 due this Monday
- I'll post some sample lit review topics





prosecutor

Watergate

↑ investigation
↓ conduct
Hotel
President

↓ airplane

- Ms. Yates's order was a remarkable rebuke by a government official to a sitting president, and it recalled the so-called Saturday Night Massacre in 1973, when President Richard M. Nixon fired his attorney general and deputy attorney general for refusing to dismiss the special prosecutor in the Watergate case

Language models

Can prob conv with quality??
frankly!!

$$P(\text{text}) = P(w_1, w_2 \dots w_N) = \prod_i P(w_i | w_1 \dots w_{i-1})$$

[history-based a.k.a. "causal"]

• **Generation:**

- Actually generate from the language distribution
- Evaluate quality of proposed translations/outputs

• **Unsupervised transfer learning:** Induce useful **word or token embeddings** for *other* NLP tasks

- Pretrained word embeddings (last lecture)
- Pretrained token (contextual) embeddings: ELMO, BERT..
- Usual assumption: train on very large corpus (>10M, >100M tokens)

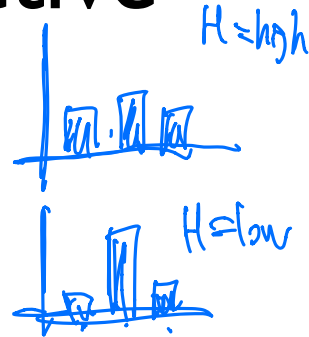
• **Model types**

- Short- or long-distance?
- Explicit features or neural representations?
- Model architecture to transmit sequential/structural information?

Information theory perspective

Entropy: uncertainty in distribution P
(obeys reasonable axioms)

$$H(P) = \sum_{\mathbf{x}} P(\mathbf{x}) \log \frac{1}{P(\mathbf{x})}$$



Cross-entropy: model P_M , test distribution P_T
(equiv. to average neg. log-likelihood)

P_T : Test $H'(P_T; P_M) = - \sum_{\mathbf{x}} P_T(\mathbf{x}) \cdot \log P_M(\mathbf{x})$
 P_M : Model

- Coding interpretation: average number of bits/nats
- Entropy of uniform V -sided die?

$V \text{ vocab PPL} = \exp(-\frac{1}{N} \sum \log \frac{1}{V}) = \exp(\log V) = V$

Information theory perspective

$$Ppl = \exp\left(-\frac{1}{N} \sum_{i=1}^N \log P(w_i | w_1 \dots w_{i-1})\right)$$

Perplexity = $S = \left[\frac{1}{V}, \frac{1}{V}, \dots, \frac{1}{V}\right]$

Cross-entropy: model P_M , test distribution P_T
 (equiv. to average neg. log-likelihood)

$$\exp\left(H'(P_T; P_M) = -\sum_{\mathbf{x}} P_T(\mathbf{x}) \cdot \log P_M(\mathbf{x})\right)$$

uniform dist with V tokens
~~Mean~~ PPL = $20,000 = V$

WSJ Penn Treebank
 $V = 20,000$
 1.5 M test tokens

- Unigram ($n = 1$): 962
- Bigram ($n = 2$): 170
- Trigram ($n = 3$): 109

$$P(w) = \frac{\#(w)}{\#tokens}$$

N-gram models

$x_1 x_2 x_3 \dots x_n$

$$\prod_i P(w_i | w_1 \dots w_{i-1}) = \prod_i P(w_i | w_{i-5} \dots w_{i-1})$$

- Markov assumption: only use *short distance* information, within a fixed window, say $k=5$ (Markov window)
- “N-gram LMs”: Markov models with count-based parameter fitting

joint stats
= N-gram counts

$$P(w_3 = \text{case} | w_1, w_2 = (\text{the, Watergate}))$$

$$\rightarrow \frac{\#(\text{the, Watergate, case})}{\#(\text{the, Watergate})}$$

Rel. freq.
Count estimator

Smoothing

$$\alpha = \frac{\alpha}{V}$$

$$\alpha = \epsilon \cdot 1$$

- Pseudocount Smoothing (Dirichlet prior) for parameter estimation:

$$P_{\text{smooth}}(w_m \mid w_{m-1}) = \frac{\text{count}(w_{m-1}, w_m) + \alpha}{\sum_{w' \in \mathcal{V}} \text{count}(w_{m-1}, w') + V\alpha}$$

- Note smoothing usually redistributes mass from seen words to unseen words

- Absolute Discounting: when $\text{count} > 0$, subtract d ($0 < d < 1$)

- Smoothing is important for many other word statistics-based preprocessing methods, like identifying multiword expressions a.k.a. collocations ("social security")

Interpolation

- Idea: higher Markov orders are more sparse.
So *combine* multiple order models
- Interpolation: weighted averaging ($\lambda \geq 0, \sum \lambda_n = 1$):

$$P_{\text{Interpolation}}(w_m \mid w_{m-1}, w_{m-2}) = \lambda_3 p_3^*(w_m \mid w_{m-1}, w_{m-2}) \\ + \lambda_2 p_2^*(w_m \mid w_{m-1}) \\ + \lambda_1 p_1^*(w_m).$$

- It's a generative model:

for Each token $w_m, m = 1, 2, \dots, M$ **do**:
draw the n -gram size $z_m \sim \text{Categorical}(\lambda)$;
draw $w_m \sim p_{z_m}^*(w_m \mid w_{m-1}, \dots, w_{m-z_m})$.

- If only we knew z , learning would be easy.
- So... use **EM!**

EM for the interpolation model

for Each token $w_m, m = 1, 2, \dots, M$ **do**:
draw the n -gram size $z_m \sim \text{Categorical}(\lambda)$;
draw $w_m \sim \mathbf{p}_{z_m}^*(w_m \mid w_{m-1}, \dots, w_{m-z_m})$.

Algorithm 10 Expectation-maximization for interpolated language modeling

```
1: procedure ESTIMATE INTERPOLATED  $n$ -GRAM ( $\mathbf{w}_{1:M}, \{\mathbf{p}_n^*\}_{n \in 1:n_{\max}}$ )  
2:   for  $z \in \{1, 2, \dots, n_{\max}\}$  do ▷ Initialization  
3:      $\lambda_z \leftarrow \frac{1}{n_{\max}}$   
4:   repeat  
5:     for  $m \in \{1, 2, \dots, M\}$  do ▷ E-step  
6:       for  $z \in \{1, 2, \dots, n_{\max}\}$  do  
7:          $q_m(z) \leftarrow \mathbf{p}_z^*(w_m \mid \mathbf{w}_{1:m-}) \times \lambda_z$   
8:          $\mathbf{q}_m \leftarrow \text{Normalize}(\mathbf{q}_m)$   
9:       for  $z \in \{1, 2, \dots, n_{\max}\}$  do ▷ M-step  
10:       $\lambda_z \leftarrow \frac{1}{M} \sum_{m=1}^M q_m(z)$   
11:   until tired  
12:   return  $\lambda$ 
```

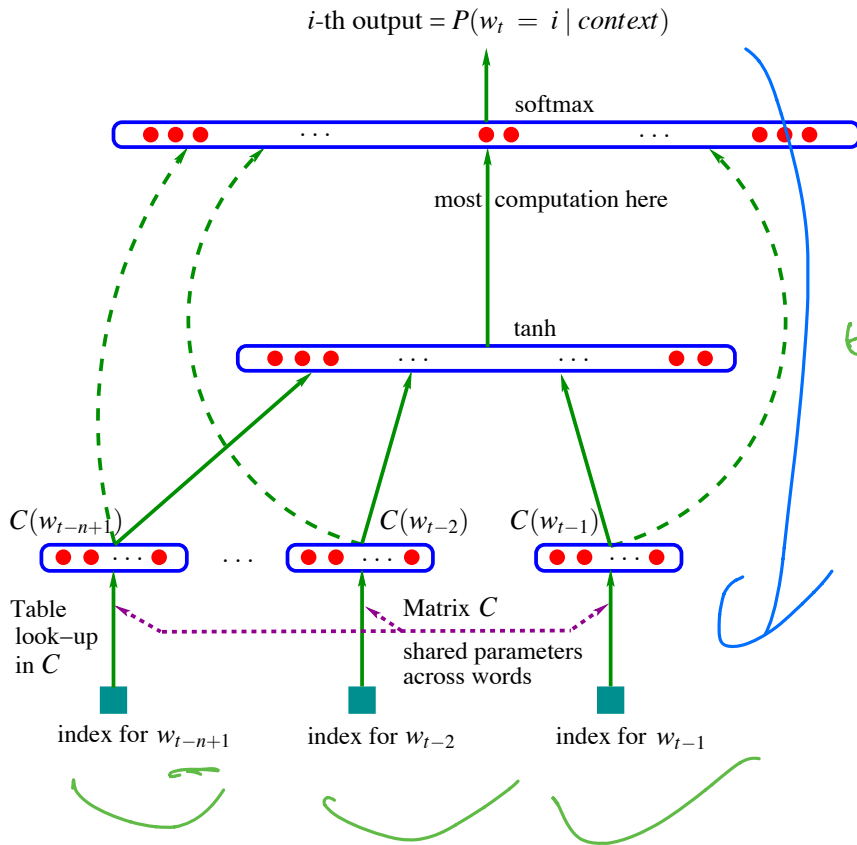
Logistic word prediction

- No more counts. Model next-word as softmax over the vocabulary.
- We can use anything to help predictions: features (Rosenfeld 1996) or neural nets (Bengio et al. 2003) to compose \mathbf{v}_u :

$$p(w | u) = \frac{\exp(\beta_w \cdot \mathbf{v}_u)}{\sum_{w' \in \mathcal{V}} \exp(\beta_{w'} \cdot \mathbf{v}_u)} \quad \begin{array}{l} \beta_w \in \mathbb{R}^K \\ \mathbf{v}_u \in \mathbb{R}^K \end{array}$$

- Can use character-level models (helps with out-of-vocabulary words), long-distance topical information, or any type of other information from the left context!

Bengio et al. 2003: Markov multilayer perceptron



Learn: C, W, U, H, d (chain rule)

$$\hat{P}(w_t | w_{t-1}, \dots, w_{t-n+1}) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}}$$

Output layer (softmax / log-linear)

shortcut linear layer

another hidden layer, size h

$$y = b + Wx + U \tanh(d + Hx)$$

Vocab output: log-probs size V

$$x = (C(w_{t-1}), C(w_{t-2}), \dots, C(w_{t-n+1}))$$

Lookup layer with concatenation: (kinda) hidden layer size $(n-1)m$

$C(i) \in \mathbb{R}^m$ Word embedding parameters

Pros/cons

V^k

Neural Networks (vs. Ngrams)

- ⇒ Fine tuning
- ⇒ Add more info
it out of vocab smoothness
- ⇒ Generalization!

Ngram Pros

- ⇒ Fast to train
- ⇒ Simple
- ⇒ Interpretable

Long-distance: Human LMs?

2.5 Long Distance (Triggers)

2.5.1 Evidence for Long Distance Information

Shannon Game at IBM [Mercer and Roukos 92]. A “Shannon game” program was implemented at IBM, where a person tries to predict the next word in a document while given access to the entire history of the document. The performance of humans was compared to that of a trigram language model. In particular, the cases where humans outsmarted the model were examined. It was found that in 40% of these cases, the predicted word, or a word related to it, occurred in the history of the document.

- Cognitive science & human behavioral evidence can inspire better NLP modeling
- Inspecting *differences* in two models’ performance (here, human-vs-machine; can also do machine-vs-machine)

Long-distance conditioning

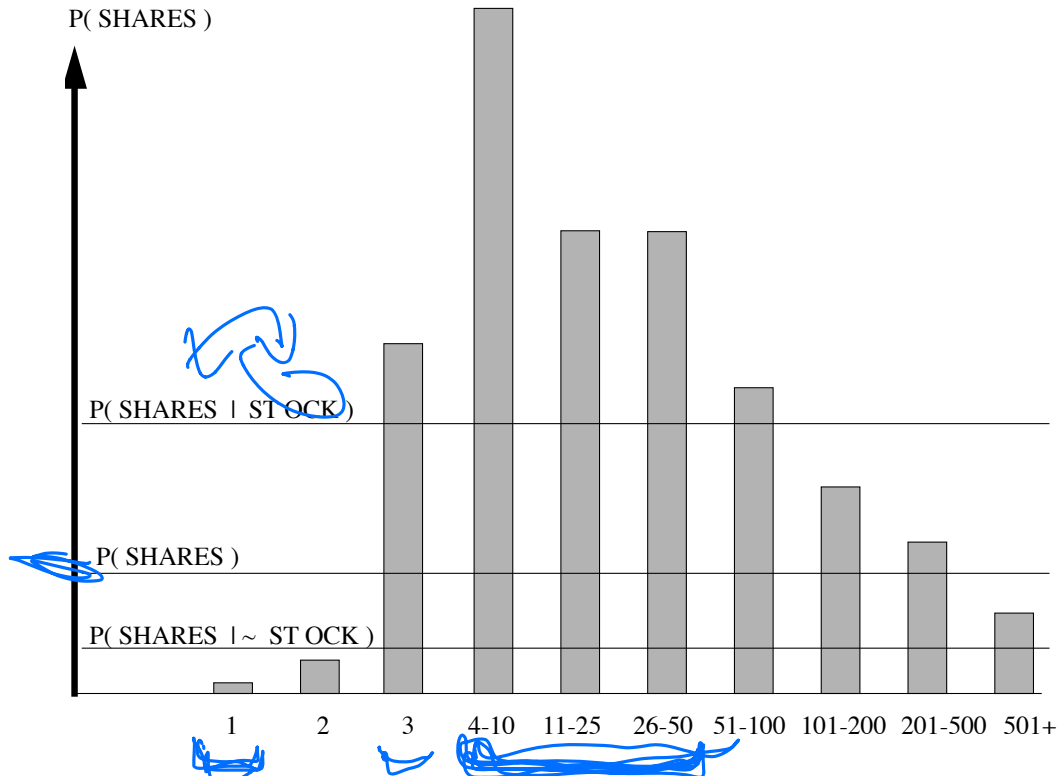


Figure 2: Probability of 'SHARES' as a function of the distance from the last occurrence of 'STOCK' in the same document. The middle horizontal line is the unconditional probability. The top (bottom) line is the probability of 'SHARES' given that 'STOCK' occurred (did not occur) before in the document.

faral eddit

[Rosenfeld 1996]

Recurrent neural networks

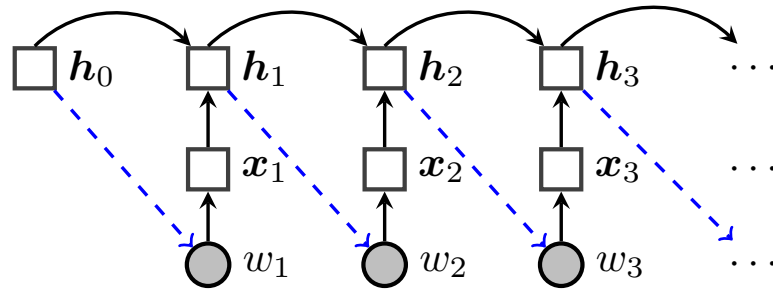


Figure 6.1: The recurrent neural network language model, viewed as an “unrolled” computation graph. Solid lines indicate direct computation, dotted blue lines indicate probabilistic dependencies, circles indicate random variables, and squares indicate computation nodes.

- Idea: extend a feedforward net to sequential data by iterating an NN at each position.
- Theoretically, an RNN can learn *any* update function. (Represent any Turing machine!)