Lecture 6: Representing Words

Kai-Wei Chang CS @ UCLA <u>kw@kwchang.net</u>

Couse webpage: https://uclanlp.github.io/CS269-17/



ML in NLP

Bag-of-Words with N-grams

N-grams: a contiguous sequence of n tokens from a given piece of text



http://recognize-speech.com/language-model/n-gram-model/comparison



Language model

Probability distributions over sentences (i.e., word sequences) $P(W) = P(w_1 w_2 w_3 w_4 \dots w_k)$ Can use them to generate strings $P(w_k \mid w_2 w_3 w_4 \dots w_{k-1})$ Rank possible sentences P("Today is Tuesday") > P("Tuesday Today is") P("Today is Tuesday") > P("Today is Los Angeles")



N-Gram Models

- Unigram model: $P(w_1)P(w_2)P(w_3) \dots P(w_n)$ Bigram model: $P(w_1)P(w_2|w_1)P(w_3|w_2) \dots P(w_n|w_{n-1})$ Trigram model: $P(w_1)P(w_2|w_1)P(w_3|w_2,w_1) \dots P(w_n|w_{n-1}w_{n-2})$ N-gram model:
- $P(w_1)P(w_2|w_1)...P(w_n|w_{n-1}w_{n-2}...w_{n-N})$



Random language via n-gram

http://www.cs.jhu.edu/~jason/465/PowerPo int/lect01,3tr-ngram-gen.pdf

Collection of n-gram

https://research.googleblog.com/2006/08/al I-our-n-gram-are-belong-to-you.html



N-Gram Viewer



https://books.google.com/ngrams



ML in NLP

How to represent words?

N-gram -- cannot capture word similarity

Word clusters

- Brown Clustering
- Part-of-speech tagging

Continuous space representation Word embedding



Brown Clustering

- Similar to language model But, basic unit is "word clusters"
- Intuition: similar words appear in similar context
- Recap: Bigram Language Models

$$◊ P(w_0, w_1, w_2, ..., w_n) = P(w_1 | w_0)P(w_2|w_1) ... P(w_n | w_{n-1}) = Π_{i=1}^n P(w_i | w_{i-1})$$

 w_0 is a dummy word representing "begin of a sentence"



a dog is chasing a cat"
 P(w₀, "a", "dog", ..., "cat")
 P("a" | w₀)P("dog" | "a") ... P("cat" | "a")

Assume Every word belongs to a cluster





Assume every word belongs to a cluster

"a dog is chasing a cat"





6501 Natural Language Processing

Assume every word belongs to a cluster

"a dog is chasing a cat"





Assume every word belongs to a cluster
 "the boy is following a rabbit"





Assume every word belongs to a cluster

"a fox was chasing a bird"





Brown Clustering

Let C(w) denote the cluster that w belongs to
"a dog is chasing a cat"





Brown clustering model

P("a dog is chasing a cat")

 $= P(C("a")|C_0) P(C("dog")|C("a")) P(C("dog")|C("a"))...$ P("a"|C("a"))P("dog"|C("dog"))...



Brown clustering model

P("a dog is chasing a cat")

 $= P(C("a")|C_0) P(C("dog")|C("a")) P(C("dog")|C("a"))...$ P("a"|C("a"))P("dog"|C("dog"))...

In general

 $P(w_{0}, w_{1}, w_{2}, ..., w_{n})$ $= P(C(w_{1}) | C(w_{0}))P(C(w_{2})|C(w_{1})) ... P(C(w_{n}) | C(w_{n-1}))$ $P(w_{1}|C(w_{1})P(w_{2}|C(w_{2})) ... P(w_{n}|C(w_{n}))$ $= \Pi_{i=1}^{n} P(C(w_{i}) | C(w_{i-1}))P(w_{i} | C(w_{i}))$



Model parameters

 $P(w_0, w_1, w_2, ..., w_n) = \prod_{i=1}^{n} P(C(w_i) | C(w_{i-1})) P(w_i | C(w_i))$



Model parameters

$$P(w_0, w_1, w_2, \dots, w_n)$$

- $= \prod_{i=1}^{n} P(C(w_i) \mid C(w_{i-1})) P(w_i \mid C(w_i))$
- ✤ A vocabulary set W
- ♦ A function $C: W \rightarrow \{1, 2, 3, \dots k\}$
 - A partition of vocabulary into k classes
- ♦ Conditional probability P(c' | c) for $c, c' \in \{1, ..., k\}$
- ♦ Conditional probability P(w | c) for $c, c' \in \{1, ..., k\}, w \in c$

 $\boldsymbol{\theta}$ represents the set of conditional probability parameters C represents the clustering



Log likelihood

- $LL(\theta, C) = \log P(w_0, w_1, w_2, ..., w_n | \theta, C)$ = $\log \prod_{i=1}^{n} P(C(w_i) | C(w_{i-1})) P(w_i | C(w_i))$ = $\sum_{i=1}^{n} [\log P(C(w_i) | C(w_{i-1})) + \log P(w_i | C(w_i))]$
- Maximizing $LL(\theta, C)$ can be done by alternatively update θ and C
 - 1. $\max_{\theta \in \Theta} LL(\theta, C)$
 - 2. $\max_{C} LL(\theta, C)$



 $\max_{\theta \in \Theta} LL(\theta, C)$

$$LL(\theta, C) = \log P(w_0, w_1, w_2, ..., w_n | \theta, C)$$

= $\log \prod_{i=1}^{n} P(C(w_i) | C(w_{i-1})) P(w_i | C(w_i))$
= $\sum_{i=1}^{n} [\log P(C(w_i) | C(w_{i-1})) + \log P(w_i | C(w_i))]$

♦ P(c' | c) =
$$\frac{\#(c',c)}{\#c}$$

This part is the same as training a POS tagging model

♦ P(w | c) = $\frac{\#(w,c)}{\#c}$

See section 9.2:

http://ciml.info/dl/v0_99/ciml-v0_99-ch09.pdf



$$\max_{C} LL(\theta, C)$$

$$\max_{C} \sum_{i=1}^{n} \left[\log P(C(w_{i}) \mid C(w_{i-1})) + \log P(w_{i} \mid C(w_{i})) \right]$$
$$= n \sum_{c=1}^{k} \sum_{c'=1}^{k} p(c,c') \log \frac{p(c,c')}{p(c)p(c')} + G$$

where G is a constantSee classnote here:
http://web.cs.ucla.edu/~kwchang/teaching
/NLP16/slides/classnote.pdf

$$p(c,c') = \frac{\#(c,c')}{\sum_{c,c'} \#(c,c')} , \quad p(c) = \frac{\#(c)}{\sum_{c} \#(c)}$$

★ *c*: cluster of w_i , *c*':cluster of w_{i-1}

$$\frac{p(c,c')}{p(c)p(c')} = \frac{p(c|c')}{p(c)}$$
 (mutual information)



Algorithm 1

- Start with |V| clusters each word is in its own cluster
- The goal is to get k clusters
- We run |V|-k merge steps:
 - Pick 2 clusters and merge them
 - ***** Each step pick the merge maximizing $LL(\theta, C)$
- Cost? (can be improved to $O(|V|^3)$) $O(|V|-k) \quad O(|V|^2) \quad O(|V|^2) = O(|V|^5)$ #Iters #pairs compute LL



Algorithm 2

- m : a hyper-parameter, sort words by frequency
- Take the top m most frequent words, put each of them in its own cluster $c_1, c_2, c_3, \dots c_m$

♦ For
$$i = (m + 1) ... |V|$$

- ♦ Create a new cluster c_{m+1} (we have m+1 clusters)
- ♦ Choose two cluster from m+1 clusters based on LL(θ, C) and merge ⇒ back to m clusters
- ♦ Carry out (m-1) final merges \Rightarrow full hierarchy
- Running time O(|V|m² + n), n=#words in corpus

Example clusters (Brown+1992)

Friday Monday Thursday Wednesday Tuesday Saturday Sunday weekends Sundays Saturdays June March July April January December October November September August people guys folks fellows CEOs chaps doubters commies unfortunates blokes down backwards ashore sideways southward northward overboard aloft downwards adrift water gas coal liquid acid sand carbon steam shale iron great big vast sudden mere sheer gigantic lifelong scant colossal man woman boy girl lawyer doctor guy farmer teacher citizen American Indian European Japanese German African Catholic Israeli Italian Arab pressure temperature permeability density porosity stress velocity viscosity gravity tension mother wife father son husband brother daughter sister boss uncle machine device controller processor CPU printer spindle subsystem compiler plotter John George James Bob Robert Paul William Jim David Mike anyone someone anybody somebody feet miles pounds degrees inches barrels tons acres meters bytes director chief professor commissioner commander treasurer founder superintendent dean custodian liberal conservative parliamentary royal progressive Tory provisional separatist federalist PQ had hadn't hath would've could've should've must've might've asking telling wondering instructing informing kidding reminding bothering thanking deposing that tha theat

head body hands eyes voice arm seat eye hair mouth



6501 Natural Language Processing

Example Hierarchy(Miller+2004)

lawyer	1000001101000
newspaperman	100000110100100
stewardess	100000110100101
toxicologist	10000011010011
slang	1000001101010
babysitter	100000110101100
conspirator	1000001101011010
womanizer	1000001101011011
mailman	10000011010111
salesman	100000110110000
bookkeeper	1000001101100010
troubleshooter	10000011011000110
bouncer	10000011011000111
technician	1000001101100100
janitor	1000001101100101
saleswoman	1000001101100110
Nilse	101101110010010101011100
Mastag	101101110010010101011100
Conoroli	1011011100100101010111010
Gan	1011011100100101010111110
Harley Davidson	1011011100100101010111110
Enfield	10110111001001010101111110
Conne	10110111001001010101111111
Microsoft	10110111001001010101111111
Ventritex	101101110010010110010
Tractabal	1011011100100101100110
Suppose	1011011100100101100111
WordPerfect	1011011100100101101000
wordt cricer	1011011100100101101000
John	10111001000000000
Consuelo	10111001000000001
Jeffrey	10111001000000010
Kenneth	1011100100000001100
Phillip	10111001000000011010

