

LANGUAGE MODELING

David Kauchak
CS159 – Fall 2020

some slides adapted from
Jason Eisner

Admin

How did assignment 1 finish up?

Assignment 2 out soon (two part assignment)

Class participation

Videos!

Independence

Two variables are independent if they do not affect each other

For two independent variables, knowing the value of one does not change the probability distribution of the other variable

- ▣ the result of the toss of a coin is independent of a roll of a dice
- ▣ price of tea in England is independent of the whether or not you get an A in NLP

Independent or Dependent?

You catching a cold and a butterfly flapping its wings in Africa

Miles per gallon and driving habits

Height and longevity of life

Independent variables

How does independence affect our probability equations/properties?



If A and B are independent, written $A \perp B$

- $P(A|B) = P(A)$
- $P(B|A) = P(B)$

What does that mean about $P(A,B)$?

Independent variables

How does independence affect our probability equations/properties?



If A and B are independent, written $A \perp B$

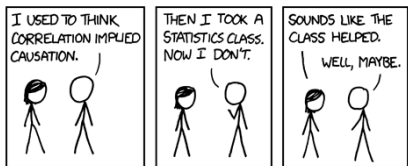
- $P(A|B) = P(A)$
- $P(B|A) = P(B)$
- $P(A,B) = P(A|B) P(B) = P(A) P(B)$
- $P(A,B) = P(B|A) P(A) = P(A) P(B)$

Conditional Independence

Dependent events can become independent given certain other events

Examples,

- height and length of life
- "correlation" studies
 - size of your lawn and length of life



<http://xkcd.com/552/>

Conditional Independence

Dependent events can become independent given certain other events

Examples,

- height and length of life
- "correlation" studies
 - size of your lawn and length of life

If A, B are conditionally independent given C $A \perp B | C$

- $P(A,B|C) = P(A|C) P(B|C)$
- $P(A|B,C) = P(A|C)$
- $P(B|A,C) = P(B|C)$
- but $P(A,B) \neq P(A)P(B)$

Assume independence

Sometimes we will assume two variables are independent (or conditionally independent) even though they're not

Why?

- Creates a simpler model
 - $p(X,Y)$ many more variables than just $P(X)$ and $P(Y)$
- May not be able to estimate the more complicated model

Language modeling

What does natural language look like?

More specifically in NLP, probabilistic model

$p(\text{ sentence })$

- $p(\text{"I like to eat pizza"})$
- $p(\text{"pizza like I eat"})$

Often is posed as: $p(\text{ word } | \text{ previous words })$

- $p(\text{"pizza"} | \text{"I like to eat"})$
- $p(\text{"garbage"} | \text{"I like to eat"})$
- $p(\text{"run"} | \text{"I like to eat"})$

Language modeling

How might these models be useful?

- Language generation tasks
 - machine translation
 - summarization
 - simplification
 - speech recognition
 - ...
- Text correction
 - spelling correction
 - grammar correction

Ideas?

$p(\text{"I like to eat pizza"})$

$p(\text{"pizza like I eat"})$

$p(\text{"pizza"} | \text{"I like to eat"})$

$p(\text{"garbage"} | \text{"I like to eat"})$

$p(\text{"run"} | \text{"I like to eat"})$

Look at a corpus

Three Google search results are shown, each with the Google logo, a search bar, and a search button. The first search is for the phrase "I like to eat pizza", showing approximately 189,000 results in 0.34 seconds. The second search is for "pizza like I eat", showing 5 results in 0.31 seconds. The third search is for "I like to eat", showing about 2,400,000 results in 0.33 seconds. Each search result includes a link to "Advanced search" and status indicators for "Instant is off" and "SafeSearch off".

Language modeling

I think today is a good day to be me

A Google search interface is shown with the search bar containing the sentence "I think today is a good day to be me". Below the search bar, there is a "Web" section with a "Show options..." link. A warning icon and text state: "No results found for 'I think today is a good day to be me'".

Language modeling is about dealing with data sparsity!

Probabilistic Language modeling

A probabilistic explanation of how the sentence was generated

Key idea:

- ▣ break this generation process into smaller steps
- ▣ estimate the probabilities of these smaller steps
- ▣ the overall probability is the combined product of the steps

Language modeling

Many approaches:

- ▣ n-gram language modeling
 - ▣ Start at the beginning of the sentence
 - ▣ Generate one word at a time based on the previous words
- ▣ syntax-based language modeling
 - ▣ Construct the syntactic tree from the top down
 - ▣ e.g. context free grammar
 - ▣ eventually at the leaves, generate the words

Pros/cons?

n-gram language modeling

I think today is a good day to be me

Google "I think" Search

Web Show options... Results 1 - 10 of about 564,000,000 for "I think". (0.28 seconds)

Google "today is a good day" Search

Web Show options... Results 1 - 10 of about 10,100,000 for "today is a good day".

Google "to be me" Search

Web Show options... Results 1 - 10 of about 70,200,000 for "to be me".

Our friend the chain rule

Step 1: decompose the probability

$$P(\text{I think today is a good day to be me}) =$$

$$P(\text{I} \mid \langle \text{start} \rangle) \times$$

$$P(\text{think} \mid \text{I}) \times$$

$$P(\text{today} \mid \text{I think}) \times$$

$$P(\text{is} \mid \text{I think today}) \times$$

$$P(\text{a} \mid \text{I think today is}) \times$$

$$P(\text{good} \mid \text{I think today is a}) \times$$

...

How can we simplify these?

The n-gram approximation

Assume each word depends only on the previous $n-1$ words
(e.g. trigram: three words total)

$$P(\text{is} \mid \text{I think today}) \approx P(\text{is} \mid \text{think today})$$

$$P(\text{a} \mid \text{I think today is}) \approx P(\text{a} \mid \text{today is})$$

$$P(\text{good} \mid \text{I think today is a}) \approx P(\text{good} \mid \text{is a})$$

Estimating probabilities

How do we find probabilities?

$P(\text{is} \mid \text{think today})$

Get real text, and start counting (MLE)!

$$P(\text{is} \mid \text{think today}) = \frac{\text{count}(\text{think today is})}{\text{count}(\text{think today})}$$

Estimating from a corpus

Corpus of sentences
(e.g. gigaword corpus)

A vertical list of horizontal lines representing sentences, with a red question mark below it. A blue arrow points to a yellow box labeled "n-gram language model".

Estimating from a corpus

I am a happy Pomona College student .

↓ count all of the trigrams

```

<start> <start> I
<start> I am
I am a
am a happy
a happy Pomona
happy Pomona College
Pomona College student
College student .
student . <end>
. <end> <end>
    
```

why do we need <start> and <end>?

Estimating from a corpus

I am a happy Pomona College student .

↓ count all of the trigrams

```

<start> <start> I
<start> I am
I am a
am a happy
a happy Pomona
happy Pomona College
Pomona College student
College student .
student . <end>
. <end> <end>
    
```

Do we need to count anything else?

Estimating from a corpus

I am a happy Pomona College student .

↓ count all of the bigrams

```

<start> <start>
<start> I
I am
am a
a happy
happy Pomona
Pomona College
College student
student .
. <end>
    
```

$$p(c | a b) = \frac{\text{count}(a b c)}{\text{count}(a b)}$$

Estimating from a corpus

1. Go through all sentences and count trigrams and bigrams

- usually you store these in some kind of data structure

2. Now, go through all of the trigrams and use the count and the bigram count to calculate MLE probabilities

- do we need to worry about divide by zero?

Applying a model

Given a new sentence, we can apply the model

$$p(\text{Pomona College students are the best .}) = ?$$



$$p(\text{Pomona} | \langle \text{start} \rangle \langle \text{start} \rangle)$$

$$p(\text{College} | \langle \text{start} \rangle \text{Pomona})$$

$$p(\text{students} | \text{Pomona College})$$

⋮

$$p(\langle \text{end} \rangle | \langle \text{end} \rangle)$$

Generating examples

We can also use a trained model to generate a random sentence

Ideas?

$\langle \text{start} \rangle \langle \text{start} \rangle$ _____

We have a distribution over all possible starting words

- $p(\text{A} | \langle \text{start} \rangle \langle \text{start} \rangle)$
- $p(\text{Apples} | \langle \text{start} \rangle \langle \text{start} \rangle)$
- $p(\text{I} | \langle \text{start} \rangle \langle \text{start} \rangle)$
- $p(\text{The} | \langle \text{start} \rangle \langle \text{start} \rangle)$
- ⋮
- $p(\text{Zebras} | \langle \text{start} \rangle \langle \text{start} \rangle)$

Draw one from this distribution

Generating examples

$\langle \text{start} \rangle \langle \text{start} \rangle$ Zebras _____

repeat!

- $p(\text{are} | \langle \text{start} \rangle \text{Zebras})$
- $p(\text{eat} | \langle \text{start} \rangle \text{Zebras})$
- $p(\text{think} | \langle \text{start} \rangle \text{Zebras})$
- $p(\text{and} | \langle \text{start} \rangle \text{Zebras})$
- ⋮
- $p(\text{mostly} | \langle \text{start} \rangle \text{Zebras})$

Generation examples

Unigram

are were that ères mammal naturally built describes jazz territory heteromyids film tenor prime live founding must on was feet negro legal gate in on beside . provincial san ; stephenson simply spaces stretched performance double-entry grove replacing station across to burma . repairing ères capital about double reached omnibus el time believed what hotels parameter jurisprudence words syndrome to ères profanity is administrators ères offices hilarius institutionalized remains writer royalty dennis , ères tyson , and objective , instructions seem timekeeper has ères valley ères " magnitudes for love on ères from allakaket , , ana central enlightened . to , ères is belongs fame they the corrected , . on in pressure %NUMBER% her flavored ères derogatory is won metacard indirectly of crop duty learn northbound ères ères dancing similarity ères named ères berkeley . . off-scale overtime . each mansfield stripes dānu traffic ossetic and at alpha popularity town

Generation examples

Bigrams

the wikipedia county , mexico .

maurice ravel . it is require that is sparta , where functions . most widely admired .

halogens chamiali cast jason against test site .

Generation examples

Trigrams

is widespread in north africa in june %NUMBER% %NUMBER% units were built by with .

jewish video spiritual are considered ircd , this season was an extratropical cyclone .

the british railways ' s strong and a spot .

Evaluation

We can train a language model on some data

How can we tell how well we're doing?

- for example
 - bigrams vs. trigrams
 - 100K sentence corpus vs. 100M
 - ...

Evaluation

A very good option: **extrinsic** evaluation

If you're going to be using it for machine translation

- ▣ build a system with each language model
- ▣ compare the two based on their approach for machine translation

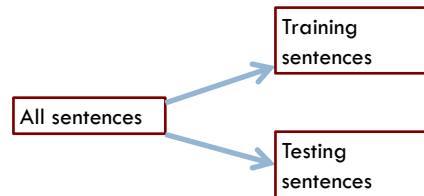
Sometimes we don't know the application

Can be time consuming

Granularity of results

Evaluation

Common NLP/machine learning/AI approach



Evaluation

n-gram
language
model

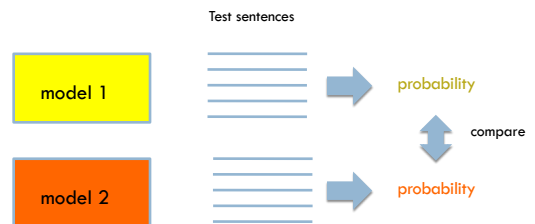
Test sentences



Ideas?

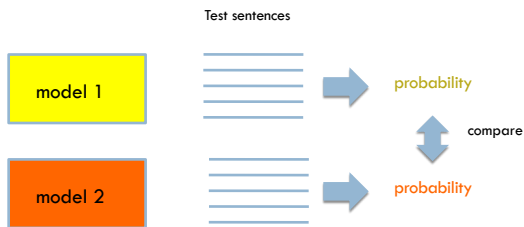
Evaluation

A good model should do a good job of predicting actual sentences



Evaluation

Pros: Fine for comparing two models
 Cons: Doesn't give us a sense of how well any model is doing



The problem

Which of these sentences will have a higher probability based on a language model?

I like to eat banana peels .

I like to eat banana peels with peanut butter.

The problem

Which of these sentences will have a higher probability based on a language model?

I like to eat banana peels .

I like to eat banana peels with peanut butter.

Since probabilities are multiplicative (and between 0 and 1), they get smaller for longer sentences.

The solution: perplexity

$$prob(w_{1..n}) = \prod_{i=1}^n p(w_i | w_{1..i-1})$$

average the probabilities



geometric mean

$$PP(w_{1..n}) = \sqrt[n]{\frac{1}{\prod_{i=1}^n p(w_i | w_{1..i-1})}}$$

Calculating perplexity in practice

$$\begin{aligned} \log \left(\sqrt[n]{\frac{1}{\prod_{i=1}^n p(w_i | w_{1:i-1})}} \right) &= \log \left(\left(\frac{1}{\prod_{i=1}^n p(w_i | w_{1:i-1})} \right)^{1/n} \right) \\ &= \frac{\log \left(\frac{1}{\prod_{i=1}^n p(w_i | w_{1:i-1})} \right)}{n} \\ &= \frac{-\log \left(\prod_{i=1}^n p(w_i | w_{1:i-1}) \right)}{n} \\ &= -\frac{\sum_{i=1}^n \log p(w_i | w_{1:i-1})}{n} \end{aligned}$$

What is this?

Calculating perplexity in practice

$$\begin{aligned} \log \left(\sqrt[n]{\frac{1}{\prod_{i=1}^n p(w_i | w_{1:i-1})}} \right) &= \log \left(\left(\frac{1}{\prod_{i=1}^n p(w_i | w_{1:i-1})} \right)^{1/n} \right) \\ &= \frac{\log \left(\frac{1}{\prod_{i=1}^n p(w_i | w_{1:i-1})} \right)}{n} \\ &= \frac{-\log \left(\prod_{i=1}^n p(w_i | w_{1:i-1}) \right)}{n} \\ &= -\frac{\sum_{i=1}^n \log p(w_i | w_{1:i-1})}{n} \end{aligned}$$

Average logprob per word!

Calculating perplexity

$$\begin{aligned} PP(w_{1:n}) &= \sqrt[n]{\frac{1}{\prod_{i=1}^n p(w_i | w_{1:i-1})}} \\ &= 10^{-\frac{\sum_{i=1}^n \log_{10} p(w_i | w_{1:i-1})}{n}} \end{aligned}$$

- This is often how it's calculated (and how we'll calculate it)
- Avoid underflow from multiplying too many small probabilities together

Another view of perplexity

Weighted average branching factor

- number of possible next words that can follow a word or phrase
- measure of the complexity/uncertainty of text (as viewed from the language models perspective)

Smoothing

What if our test set contains the following sentence, but one of the trigrams never occurred in our training data?

$P(\text{I think today is a good day to be me}) =$

$P(\text{I} \mid \langle \text{start} \rangle \langle \text{start} \rangle x)$

$P(\text{think} \mid \langle \text{start} \rangle \text{I}) x$

$P(\text{today} \mid \text{I think}) x$

$P(\text{is} \mid \text{think today}) x$

$P(\text{a} \mid \text{today is}) x$

$P(\text{good} \mid \text{is a}) x$

...

If any of these has never been seen before, prob = 0!

A better approach

$p(z \mid x y) = ?$

Suppose our training data includes

... x y a ...

... x y d ...

... x y d ...

but never: xyz

We would conclude

$p(a \mid x y) = 1/3?$

$p(d \mid x y) = 2/3?$

$p(z \mid x y) = 0/3?$

Is this ok?

Intuitively, how should we fix these?

Smoothing the estimates

Basic idea:

$p(a \mid x y) = 1/3?$ *reduce*

$p(d \mid x y) = 2/3?$ *reduce*

$p(z \mid x y) = 0/3?$ *increase*

Discount the positive counts somewhat

Reallocate that probability to the zeroes

Remember, it needs to stay a probability distribution

Other situations

$p(z \mid x y) = ?$

Suppose our training data includes

... x y a ... (100 times)

... x y d ... (100 times)

... x y d ... (100 times)

but never: x y z

Suppose our training data includes

... x y a ...

... x y d ...

... x y d ...

... x y ... (300 times)

but never: x y z

Is this the same situation as before?

Smoothing the estimates

Should we conclude

$$p(a | xy) = 1/3? \text{ reduce} \quad p(c | a b) = \frac{\text{count}(a b c)}{\text{count}(a b)}$$

$$p(d | xy) = 2/3? \text{ reduce}$$

$$p(z | xy) = 0/3? \text{ increase}$$

Readjusting the estimate is particularly important if:

- the denominator is small ...
 - 1/3 probably too high, 100/300 probably about right
- numerator is small ...
 - 1/300 is probably too high, 100/300 probably about right

Add-one (Laplacian) smoothing

xya	1	1/3	2	2/29
xyb	0	0/3	1	1/29
xyc	0	0/3	1	1/29
xyd	2	2/3	3	3/29
xye	0	0/3	1	1/29
...				
xyz	0	0/3	1	1/29
Total xy	3	3/3	29	29/29

Add-one (Laplacian) smoothing

300 observations instead of 3 – better data, less smoothing

xya	100	100/300	101	101/326
xyb	0	0/300	1	1/326
xyc	0	0/300	1	1/326
xyd	200	200/300	201	201/326
xye	0	0/300	1	1/326
...				
xyz	0	0/300	1	1/326
Total xy	300	300/300	326	326/326

Add-one (Laplacian) smoothing

What happens if we're now considering a vocabulary of 20,000 words?

xya	1	1/3	2	2/29
xyb	0	0/3	1	1/29
xyc	0	0/3	1	1/29
xyd	2	2/3	3	3/29
xye	0	0/3	1	1/29
...				
xyz	0	0/3	1	1/29
Total xy	3	3/3	29	29/29

Add-one (Laplacian) smoothing

20,000 words, not 26 letters

see the abacus	1	1/3	2	2/20003
see the abbot	0	0/3	1	1/20003
see the abduct	0	0/3	1	1/20003
see the above	2	2/3	3	3/20003
see the Abram	0	0/3	1	1/20003
...				
see the zygote	0	0/3	1	1/20003
Total	3	3/3	20003	20003/20003

Any problem with this?

Add-one (Laplacian) smoothing

An "unseen event" is a 0-count event

The probability of an unseen event is 19998/20003

- add one smoothing thinks it is very likely to see a novel event

The problem with add-one smoothing is it gives too much probability mass to unseen events

see the abacus	1	1/3	2	2/20003
see the abbot	0	0/3	1	1/20003
see the abduct	0	0/3	1	1/20003
see the above	2	2/3	3	3/20003
see the Abram	0	0/3	1	1/20003
...				
see the zygote	0	0/3	1	1/20003
Total	3	3/3	20003	20003/20003

The general smoothing problem

			modification	probability
see the abacus	1	1/3	?	?
see the abbot	0	0/3	?	?
see the abduct	0	0/3	?	?
see the above	2	2/3	?	?
see the Abram	0	0/3	?	?
...			?	?
see the zygote	0	0/3	?	?
Total	3	3/3	?	?

Add-lambda smoothing

A large dictionary makes novel events too probable.

Instead of adding 1 to all counts, add $\lambda = 0.01$?

- This gives much less probability to novel events

see the abacus	1	1/3	1.01	1.01/203
see the abbot	0	0/3	0.01	0.01/203
see the abduct	0	0/3	0.01	0.01/203
see the above	2	2/3	2.01	2.01/203
see the Abram	0	0/3	0.01	0.01/203
...			0.01	0.01/203
see the zygote	0	0/3	0.01	0.01/203
Total	3	3/3	203	