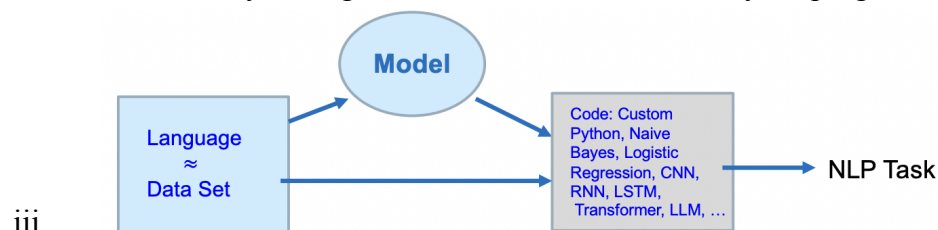


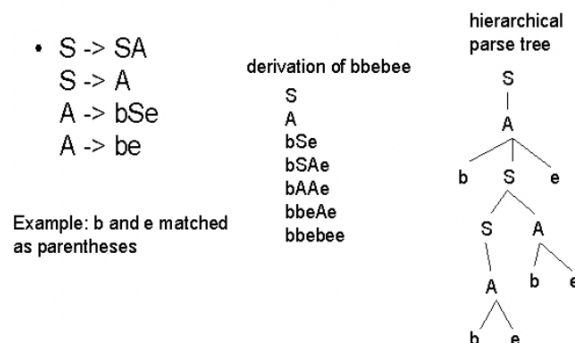
Language Models, Bag-of-Words, N-grams, Skip-Grams

1. Language Models

- a. A Language Model is a simplified representation of a language which facilitates an NLP task, where
 - i. A language (potentially infinite) is approximated by a (finite) data set; and
 - ii. The model is a set of (simplified) assumptions about the language, embodied by the algorithms and data structures of your program.



- b. NLP systems rely on models to capture knowledge of about a language, and as a representation for texts which facilitate an NLP task.
 - i. Example: Context-Free Grammars (Chomsky, Backus-Naur)



ii.

2. Language Modeling: Word Representations

- a. Before diving into the subject of Language Models, let's prepare a bit by talking about an essential component of language modeling...
- b. Last time we discussed how to represent characters, as integer ASCII codes in the range $[0 \dots 127]$. But how do we represent words?
- c. Bad idea: word = sequence of ASCII codes
- d. Why is this bad?
 - i. Variable length (every word length has to be the same, which leads to having lots of _ after the word) → ex) dog_____
 - ii. Information contains confusing correspondences:
 1. “to” very similar to “too” “dog” same chars as “god”
- e. (Neural networks will find these difficult to learn.)

- | Color | | Red | Green | Blue |
|-------|---|-----|-------|------|
| Red | | 1 | 0 | 0 |
| Green | → | 0 | 1 | 0 |
| Blue | | 0 | 0 | 1 |
| Green | | 0 | 1 | 0 |

- h. Word Embedding (we'll do these in a few weeks) – putting words that are similar close to each other



- i. We generically call sequences/lists/arrays by the standard term vector. For simplicity, we often write them as Python lists here

a. Basic Idea of One-Hot Encoding:

- Vocabulary list: ["John", "likes", "Mary", "movies", "to", "too", "watch"]

One-Hot Encodings:

“movies” \Rightarrow [0, 0, 0, 1, 0, 0, 0]

"likes" $\Rightarrow [0, 1, 0, 0, 0, 0, 0]$

vi. Plus:

1. Vectors have same length
2. Spelling is completely irrelevant

1. Very long vectors (typically 10,000 or more) → can eliminate many words that do not matter much for later

4. Language Models: Bag-of-Words

a. Examples of Models: Bag of Words (BOW)

- i. The BOW model represents a text (sentence, sequence of words, entire corpus) as a multiset (bag) of all words in the text, i.e, just the vocabulary, no information about order of words!

1. Useful in classification (Spam Vs Ham)

- ii. Text: "John likes to watch movies. Mary likes movies too."

Vocabulary list: ["John", "likes", "Mary", "movies", "to", "too", "watch"]
0 1 2 3 4 5 6

BOW model of text: [1, 2, 1, 2, 1, 1, 1]

iii.

- iv. The multiplicity (frequency) of words do not really matter

1. Ex) "The movie is bad bad bad" is similar to "The movie is bad"

- v. Do this for large number of data

- vi. Alternate BOW representations

1. We might only consider the presence (0/1) of a word, not its frequency (as if "Set of Words")
 2. Since most BOW vectors are sparse, we might want to store them as a dictionary

3. { "John" : 1, "likes" : 2, "to" : 1, "watch" : 1, "movies" : 2, "Mary" : 1, "too" : 1 }

- vii. What is the relationship between One-Hot Encodings and a BOW model?

1. The BOW model of a text is the array sum of the one-hot encodings

"John"	[1, 0, 0, 0, 0, 0, 0]
"likes"	[0, 1, 0, 0, 0, 0, 0]
"to"	[0, 0, 0, 0, 1, 0, 0]
"watch"	[0, 0, 0, 0, 0, 0, 1]
"movies"	[0, 0, 0, 1, 0, 0, 0]
"Mary"	[0, 0, 1, 0, 0, 0, 0]
"likes"	[0, 1, 0, 0, 0, 0, 0]
"movies"	[0, 0, 0, 1, 0, 0, 0]
"too"	[0, 0, 0, 0, 0, 1, 0]

+

[1, 2, 1, 2, 1, 1, 1]

2.

b. Reducing words

- i. Very common words such as "the" "an"
- ii. Very uncommon words

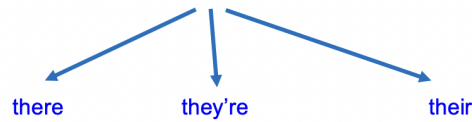
5. Language Models: Probabilistic Language Models

- a. Probabilistic Language Model: Assign a probability to text components (letters, words, sentences,)

- b. This is very useful to work with the ambiguous nature of human languages, and very amenable to computation:

- i. “Given K choices for some ambiguous input, choose the most probable one.”

“I went to **they're** house on Sunday.”



- c. Which is most likely?

- i. Which one appears the most in the corpus? If you feed wrong corpus, you will get wrong

6. Language Models: Vector Space Language Models

- a. Vector space models use a vector in M-dimensional space to represent
 - i. Words
 - ii. Sentences
 - iii. Texts



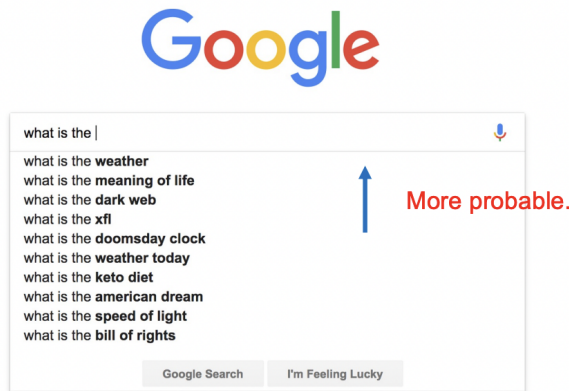
- b.
- c. This is the current SOTA (“Stage Of The Art”) for language modeling. → very useful computation since words that are similar are close to each other

7. Probabilistic Language Models

- a. Main Idea of PLMs: Assign a probability to a sequence of words. Why?
 - i. Some words are more probable than others
- b. Machine Translation:
 - i. EX) Translating from Japanese to English and there are two options
 - ii. $P(\text{high winds tonight}) > P(\text{large winds tonight})$
 - 1. Look for every occurrence of those words in the corpus
- c. Spelling Correction:
 - i. The office is about fifteen minuets from my house
 - ii. $P(\text{about fifteen minutes from}) > P(\text{about fifteen minuets from})$

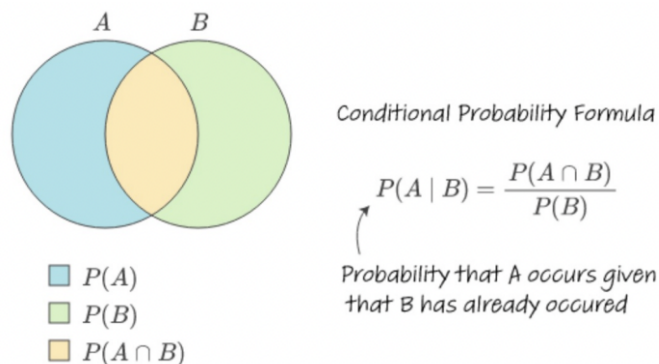
- d. Speech Recognition:
 - i. $P(\text{I saw a van}) \gg P(\text{eyes awe of an})$
 - ii. Could be any sequence and length of words (corpus might not have enough examples to capture all the data)
- e. Summarization, Q&A, etc.
- f. The main task: compute the probability of a sentence or sequence of words:
 - i. $P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$
 - ii. Words might not occur or occur \rightarrow handle later
- g. Subtask: compute the conditional probability of the next word
 - i. $P(w_5 | w_1, w_2, w_3, w_4)$ “The weather is ? “
- h. A model that computes either of these: $P(W)$ or $P(w_n | w_1, w_2 \dots w_{n-1})$ is called a probabilistic language model (often, just “language model”).

You have seen these before!



- i.
 - j. It is possible to apply this framework to any sequence (predicting the next sequence given a sequence)
 - i. Letters in a word
 - ii. Pitches in a melody
 - iii. Phonemes in a voice signal
 - iv. Sentences in a paragraph
 - v. Topics in a discourse.
 - k. The amount of data is very significant because we need to consider all the likelihoods
 - l. How to compute $P(W)$?
 - i. How to compute this joint probability
 - 1. $P(\text{I went to their house on Sunday})$
 - ii. Intuition: let's rely on the Chain Rule of Probability
8. Calculating Probabilities
- a. Recall the definition of conditional probabilities
 - i. $P(A|B) = P(B,A) / P(B)$
 - ii. Rewriting: $P(B,A) = P(B) * P(A|B)$

- iii. For this LM, we will think of B,A as a sequence: B happens, then A happens. Thus, $P(A|B)$ = “given that B has happened, what is the probability that A happens”?



b.

c. More variables:

i. $P(A,B,C,D) = P(A) \times P(B|A) \times P(C|A,B) \times P(D|A,B,C)$

d. General Chain Rule:

i. $P(x_1, x_2, x_3, \dots, x_n) = P(x_1) \times P(x_2|x_1) \times P(x_3|x_1, x_2) \times \dots \times P(x_n|x_1, \dots, x_{n-1})$

e. The Chain Rule applied to compute joint probability of words in sentence:

i. “I went to their house on Sunday.”

$$P(w_1 w_2 \dots w_n) = \prod_{1 \leq i \leq n} P(w_i | w_1 w_2 \dots w_{i-1})$$

ii.

iii. $P(\text{I went to their house on Sunday.}) =$

$$P(I) \times P(\text{went} | I) \times P(\text{to} | I \text{ went}) \times P(\text{their} | I \text{ went to}) \\ \times P(\text{house} | I \text{ went to their}) \times P(\text{on} | I \text{ went to their house}) \\ \times P(\text{Sunday} | I \text{ went to their house on})$$

f. How to estimate these probabilities?

g. Could we just count and divide?

$$P(\text{Sunday} | I \text{ went to their house on}) =$$

$$\frac{\text{Count}(I \text{ went to their house on Sunday})}{\text{Count}(I \text{ went to their house on})}$$

i.

ii. Not realistic

iii. In an infinite set of sentences, the probability of any distinct sequence of words is 0. So a data set is a small sample which hopefully represents the essential features of the language.

iv. But realistic data sets never have enough sample sequences, and sequences might be very long or simply not exist in your data.

- h. Markov Assumption: Finite history!
- i. Only consider N-1 words of left context, for some fixed N.
- j. So, if N = 2 , “I went to their house on Sunday” becomes
 - i. I went

went to
 to their
 their house
 house on
 on Sunday

- k. If N = 3, “I went to their house on Sunday” becomes
 - i. I went to

went to their
 to their house
 their house on
 house on Sunday

- l. Terminology: An N-Gram is a sequence of N contiguous words from the data set.
 - i. unigram = 1-gram
 - ii. bigram = 2-gram
 - iii. trigram = 3-gram, etc
- m. Markov Assumption: Only consider N-1 words of left context.
- n. Thus, for a sequence of length M,

$$P(w_1 w_2 \dots w_M) \approx \prod_{N \leq i \leq M-N} P(w_i | w_{i-N+1} \dots w_{i-1})$$

Bigram Example (N = 2)

$P(<s> \text{ I went to their house on Sunday } </s>) =$

$$\begin{aligned}
 &P(\text{I} | <s>) \times P(\text{went} | \text{I}) \times P(\text{to} | \text{went}) \times P(\text{their} | \text{to}) \\
 &\times P(\text{house} | \text{their}) \times P(\text{on} | \text{house}) \\
 &\times P(\text{Sunday} | \text{on}) \times P(</s> | \text{Sunday})
 \end{aligned}$$

- o.
 - i. <s> is where the sentence begins, so you lose the information that the start of the sentence is “I”

$$P(\text{I} | <s>) \approx \frac{C(<s> \text{ I})}{C(<s>)}$$

$$P(\text{went} | \text{I}) \approx \frac{C(\text{I went})}{C(\text{I})}$$

- p.
- q. Note that this calculation involves finding the number of occurrences of an N-gram and of an (N-1)-gram (the prefix)!

Trigram Example (N = 3)

$P(<s> \text{ I went to their house on Sunday } </s>) =$

$$\begin{aligned} &P(\text{ I } | <s>) \times P(\text{ went } | \text{ I }) \times P(\text{ to } | \text{ went }) \times P(\text{ their } | \text{ to }) \\ &\times P(\text{ house } | \text{ their }) \times P(\text{ on } | \text{ house }) \\ &\times P(\text{ Sunday } | \text{ on }) \times P(</s> | \text{ Sunday }) \end{aligned}$$

$$P(\text{ went } | <s> \text{ I }) = \frac{C(<s> \text{ I went })}{C(<s> \text{ I })}$$

$$P(\text{ to } | \text{ I went }) = \frac{C(\text{ I went to })}{C(\text{ I went })}$$

r.

s. In trigram, the beginning of the sentence has no length

i. Start with bigram for once → go trigram

t. Similarly, when doing pentagram,

i. Unigram → bigram → trigram → etc.. → pentagram

u. Remarks: This is almost trivial to code after you have separated your text into words and sentences.

v. For small N, it will be reasonably efficient.

w. BUT, it does NOT capture the recursive/nesting structure inherent in complex sentences:

i. My friend Bill, who went to the same high school that I did –Pennsbury, which is in Fairless Hills in PA—lives in his car, and he called me yesterday (or the day before, I forget).

x. Example

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

<s> I am Sam </s>

<s> Sam I am </s>

<s> I do not like green eggs and ham </s>

$$P(\text{ I } | <s>) = \frac{2}{3} = .67$$

$$P(\text{ Sam } | <s>) = \frac{1}{3} = .33$$

$$P(\text{ am } | \text{ I }) = \frac{2}{3} = .67$$

$$P(</s> | \text{ Sam }) = \frac{1}{2} = 0.5$$

$$P(\text{ Sam } | \text{ am }) = \frac{1}{2} = .5$$

$$P(\text{ do } | \text{ I }) = \frac{1}{3} = .33$$

i.

9. Generative Language Models

a. A clever feature of this model is that it can easily generate sentences.

b. For bigrams, after calculating the probability of all bigrams appearing in the data.

i. Pick a probable <s> w1

ii. Pick a probable bigram w1 w2

iii. ... etc. ...

iv. End when you generate a bigram wk </s>

c. You may not want to always choose the most likely, or you will not be able to generate many sentences! → it will always generate the same sentence

d. So choose randomly from the probability distribution of next words.