



UMSL

Department of  
Mathematics &  
Computer Science

## 22. Natural Language Processing

नमस्ते

Some text in  
'unknown' language



What language is this text?

## 22.1 Formal languages vs natural languages

- Formal languages (Java, Python, etc.) can be characterized by a set of rules
  - called a **grammar**
- Natural languages (English, Spanish, etc.) cannot be characterized as definitive set of rules

"Thanks for dinner. I've never seen potatoes cooked like that before."

- Natural languages are many times ambiguous
  - Natural languages are constantly changing
- Hence, **instead of deriving rules, we build approximate language models.**

# Can we build models to represent a natural language?

- This can then be used to classify text, generate text, etc.

## 22.1.1 N-gram character models

- **n-gram** is a sequence of written symbols of length  $n$ 
  - 1-gram is “unigram”
  - 2-gram is “bigram”
  - 3-gram is “trigram”
- A simple **language model** for **written text** is a probability distribution of  $n$ -grams
  - $P(c_{1:N})$  is the probability of a sequence of  $N$  characters,  $c_1$  through  $c_N$ .
  - e.g.  $P(\text{“the”}) = 0.027$  and  $P(\text{“zgq”}) = 0.00000002$
- An  $n$ -gram character model is a model of probability distribution of  $n$ -letter sequences

## 22.1.1 N-gram character models

- An n-gram model is defined as a **Markov chain of order n-1**
  - For example, 3-gram (or trigram) model is defined a Markov chain of order 2
  - **This means** - in a trigram model, the probability of a character  $c_i$  depends on two previous characters
  - i.e.  $P(c_i | c_{1:i-1}) = P(c_i | c_{i-2:i-1})$

"Thanks for dinner. I've never seen potatoes cooked like that before."

- We can define the probability of a sequence of characters  $P(c_{1:N})$  under the trigram model by first factoring with chain rule, and then using the Markov assumption

$$P(c_{1:N}) = \prod_{i=1}^N P(c_i | c_{1:i-1}) = \prod_{i=1}^N P(c_i | c_{i-2:i-1})$$

- For a trigram character model in a language with 100 characters  $P(c_i | c_{i-2:i-1})$  has a million entries (100x100x100)

How can we estimate this probability distribution? How much data do we need?

# 22.1.1 N-gram models for language identification

**Problem:** Given a text, determine what natural language it is written in.

- First build a trigram character model of each candidate language L
  - $P(c_i \mid c_{i-2:i-1}, L)$
  - For each L, the model is built by counting the trigrams in a **corpus of that language**
    - This gives  $P(\text{Text} \mid \text{Language})$ , e.g.  $P(\text{"Hello World"} \mid \text{English})$
- We know  $P(\text{Text} \mid \text{Language})$ , how can we get the most probable language?

$$\begin{aligned}\ell^* &= \underset{\ell}{\operatorname{argmax}} P(\ell \mid c_{1:N}) \\ &= \underset{\ell}{\operatorname{argmax}} P(\ell)P(c_{1:N} \mid \ell) \quad \text{Bayes' Rule: } P(A \mid B) = P(A) * P(B \mid A) / P(B) \\ &= \underset{\ell}{\operatorname{argmax}} P(\ell) \prod_{i=1}^N P(c_i \mid c_{i-2:i-1}, \ell) \quad \text{Markov chain of order 2}\end{aligned}$$

What happens to  $\ell^*$  if any of  $P(c_i \mid c_{i-2:i-1}, L)$  is 0?

## 22.1.1 N-gram models for language identification

The trigram models for each language can be learned from language corpus, but how to calculate  $P(L)$ ?

Option 1 - We may already have some estimates of  $P(L)$

- For example, if we are selecting a random Web page, we know English is the most likely language and the probability of Macedonian will be less than 1%

Option 2 - The exact number we select for these priors is not critical

- because the trigram model usually selects on language that is several order magnitude more probable than any other
  - we could do  $P(L) = 1$  for all languages when context is unknown
- 
- N-grams can be used to identify languages with > 99% accuracy!

## 22.1.2 Smoothing n-gram models

- The training corpus provides only an estimate of the true probability distribution.
- For common character sequences such as “the” or “\_th” English corpus will give a good estimate: about 1.5% of all trigrams
- However, “\_ht” or “xgz” are very uncommon
  - It is likely that the sequence will have a zero count when the training corpus is standard English
- Should we have  $P(\text{"_ht"}) = 0$ ?
  - If we do so, the  $P(\text{English} \mid \text{"The program issues an http request"})$  will be 0, which is wrong.
  - This is generalization error - We want our language models to generalize well to texts they haven't seen yet
  - Setting  $P(\text{"_ht"}) = 0$  means that our model thinks that “\_ht” is impossible in English
- Our language model needs to be adjusted!

What adjustments can we make?



## 22.1.2 Linear interpolation smoothing

- Smoothing is the process of adjusting the probability of low-frequency counts
- The Backoff model approach
  - We start by estimating n-gram counts
  - For any particular sequence that has a low (or zero) count, we back off to (n-1)-grams

$$P(\text{"zgx"})$$

- Linear interpolation smoothing is a backoff model
  - It combines trigram, bigram, and unigram models using linear interpolation
  - Probability estimate is given by

$$\hat{P}(c_i | c_{i-2:i-1}) = \lambda_3 P(c_i | c_{i-2:i-1}) + \lambda_2 P(c_i | c_{i-1}) + \lambda_1 P(c_i)$$

$$\lambda_3 + \lambda_2 + \lambda_1 = 1$$

- The parameters  $\lambda_i$  can be fixed or they can be trained.
- Example: Probability estimate of  $P(\text{"zgx"}) = 0.2 * P(\text{"zgx"}) + 0.3 * P(\text{"gx"}) + 0.5 * P(\text{"x"})$

## 22.1.3 Model evaluation

- So many possible models - unigram, bigram, trigram, interpolated smoothing with different values of  $\lambda$  - which model should we choose?
- We can evaluate the model with **cross-validation**:
  - Step 1: Randomly split the corpus into a training corpus and a validation corpus
  - Step 2: Determine the parameters of the model from the training data
  - Step 3: Evaluate the model on the validation corpus, and obtain average accuracy
    - i.e. Calculate how correct all the predicted probabilities are.
  - Step 4: Repeat step 1 to 3 a few times, say 10 times.
  - Step 5: Average the 10 accuracies - this will be the accuracy of the model

## 22.1.3 Model evaluation

- Measuring accuracy can be inconvenient
  - Say that the model predicts the probability of  $P(\text{"abc"})$ ,  $P(\text{"abd"})$ ,  $P(\text{"abe"})$ , etc.
  - The sum of all these probabilities must be 1 (it is a multi-class classification problem)
  - For a trigram model, if there are 1 million trigrams, the probability of each trigram (on average) is a small number 1/one million
  - If we have a 4-gram model, these probability values are very small numbers, and floating-point underflow can become an issue
- Alternately, we can describe the probability of a sequence using **perplexity**
  - $\text{Perplexity}(c_{1:N}) = P(c_{1:N})^{-1/N}$
  - Perplexity can be thought of as the reciprocal of probability, normalized by sequence length

## 22.1.3 Perplexity (example 1)

- Suppose there are 100 characters in a language L, and our unigram model says they are all equally likely, i.e.  $P(\text{"A"}) = 1/100$ ,  $P(\text{"B"}) = 1/100$ , etc.
- The perplexity of the model will be 100 for a sequence of any length
  - For a random sequence of length 1
    - Perplexity =  $P(\text{"A"})^{-1/1} = 0.01^{-1} = 100$
  - For a random sequence of length 2
    - Perplexity =  $P(\text{"AB"})^{-1/2} = (0.01 * 0.01)^{-1/2} = (0.01 * 0.01)^{-0.5} = 0.0001^{-0.5} = 100$
- If some characters were more likely than others, then the model will reflect it, with a perplexity less than 100.

$$\text{Perplexity}(c_{1:N}) = P(c_{1:N})^{-1/N}$$

# Perplexity (example 2)

Suppose there are 3 characters in a language L, and we have built a unigram model. The probabilities for the 3 characters are given by the models are  $P(\text{"A"}) = 0.25$ ,  $P(\text{"B"}) = 0.50$ , and  $P(\text{"C"}) = 0.25$ . What will be the perplexity of for the sequences "AAA" and "ABC"? How will the perplexity change for the two sequences if the probabilities were equal for the 3 characters?

- $\text{Perplexity}(\text{"AAA"}) = P(\text{"AAA"})^{-1/3} = (0.25 * 0.25 * 0.25)^{-1/3} = 3.94$
- $\text{Perplexity}(\text{"ABC"}) = P(\text{"ABC"})^{-1/3} = (0.25 * 0.50 * 0.25)^{-1/3} = 3.13$  (less perplex = high probability)

If the probabilities were equal:

- $\text{Perplexity}(\text{"AAA"}) = P(\text{"AAA"})^{-1/3} = (0.33 * 0.33 * 0.33)^{-1/3} = 2.99$
- $\text{Perplexity}(\text{"ABA"}) = P(\text{"ABC"})^{-1/3} = (0.33 * 0.33 * 0.33)^{-1/3} = 2.99$

$$\text{Perplexity}(c_{1:N}) = P(c_{1:N})^{-1/N}$$

## 22.1.4 N-gram word models

- The vocabulary is much larger
  - There are only about 100 characters (for n-gram character models) but there can be tens of thousands of symbols
  - In English a sequence of letters is surrounded by spaces is a word, but in some languages, like Chinese, words are not separated by spaces
- How to deal with out of vocabulary words?
  - Add a new word to the vocabulary: <UNK> (unknown word)
  - Estimate the n-gram counts for <UNK> using the following process:
    - Scan the training corpus and the first time any individual word appears that is previously unknown, replace it with <UNK>
    - All subsequent appearances of the word remain unchanged
    - Then compute n-gram counts for the corpus as usual, treating <UNK> just like any other word
    - When an unknown word appears in a test set, look up its probability under <UNK>
  - This means, words such as “xfhdy” have only single time influence (it is counted only once) even if they appear multiple times

## 22.1.4 Unigram, bigram, and trigram word models

- Authors build a unigram, bigram, and trigram models over the words in the book and then randomly sampled sequences of words from the models
  - Unigram: logical are as are confusion a may right tries agent goal the was . . .
  - Bigram: systems are very similar computational approach would be represented . . .
  - Trigram: planning and scheduling are integrated the success of naive bayes model is . . .
  - Which sampling appears more close to AI language?
- Perplexity
  - Unigram model = 891
  - Bigram model = 142
  - Trigram model = 91
- If so, why should we not build 4-gram or 5-gram word models?

How to use n-gram word models to detect spam emails?



## 22.2 Spam detection (classify email as spam or 'ham')

Spam: Wholesale FashionWatches -57% today. Designer watches for cheap ...

Spam: You can buy ViagraFr\$1.85 All Medications at unbeatable prices! ...

Spam: WE CAN TREAT ANYTHING YOU SUFFER FROM JUST TRUST US ...

Spam: Sta.rt earn\*ing the salary yo,u d-eserve by o'btaining the prope,r crede'ntials!

Ham: The practical significance of hypertree width in identifying more ...

Ham: Abstract: We will motivate the problem of social identity clustering: ...

Ham: Good to see you my friend. Hey Peter, It was good to hear from you. ...

Ham: PDS implies convexity of the resulting optimization problem (Kernel Ridge ...

### Language Modeling approach (approach 1)

- Define **one n-gram model for  $P(\text{Message} \mid \text{Spam})$**  by training on spam folder and **one model for  $P(\text{Message} \mid \text{Ham})$**  by training on the inbox (basically, train two models)
- Classify a new message using Bayes' rule:
  - $P(c)$  is estimated by simply counting total number of spam and ham messages

$$\operatorname{argmax}_{c \in \{spam, ham\}} P(c \mid message) = \operatorname{argmax}_{c \in \{spam, ham\}} P(message \mid c) P(c)$$

# 22.2 Spam detection (classify email as spam or 'ham')

## Machine Learning Approach (approach 2)

- Represent the message as a set of feature/value pairs and apply classification algorithm 'h' to the feature vector X
- N-grams can be input as features
  - For example, for a unigram model, each word is a feature
  - The feature value is the count of the word in the email
  - If there are 100,000 words in the language model, then the feature vector is 100,000 long
  - For a short email message almost all features will have a count zero

	<b>Word 1</b>	<b>...</b>	<b>Word N</b>	<b>Spam</b>
Email 1	2	...	8	0
Email 2	5	...	5	1

- With unigram models, the notion of the order of words is lost

## 22.2 Spam detection (classify email as spam or 'ham')

- With bigrams and trigrams the number of features is squared or cubed
- It can be expensive to run algorithms on a very large feature vector, so often a process of **feature selection** is used
- With feature selection, often large number of features can be reduced to as little as a few hundred features
- Currently, we can detect spams with 98% to 99% accuracy!
  - Is this true? Why?

# Note

- “Word embeddings,” an alternative to one-hot word vectors, are the trend
  - Particularly in deep learning methods
- But, N-grams are a powerful, unavoidable feature engineering tools
  - Particularly, when using lightweight, shallow text processing models such as logistic regression and random forests.

## 22.3.1 Information retrieval (IR) scoring functions

- IR is the task of finding documents that are relevant to a user's need for information
  - For example, search engines on the World Wide Web
- A **scoring function** takes a document and a query and returns a numeric score
  - the most relevant documents have the highest scores
- If we have score for each documents (for an input query):
  - then we can rank the documents by the score and return the results



## 22.3.1 BM25 scoring function

- Say we want to search “farming in Kansas” in a corpus of documents
  - BM25 score is a linear weighted combination of scores for each of the words that make up the query
  - The total score is the sum of score for each of the three words
- Three factors affect the weight of each query term:
  - (1) Term Frequency (TF) - the frequency with which the query term appears in a document
    - For example, the documents that mention “farming” will appear more frequently
  - (2) Inverse Document Frequency Term (IDF)
    - The word “in” appears in almost every document, so it has high document frequency, and thus low inverse document frequency, and so it is not as important to the query as “farming” or “Kansas”
  - (3) Length of the document
    - A million word document will mention all query words, but may not actually be about the query; instead, a short document that mentions all the words is a much better candidate

## 22.3.1 BM25 scoring function

$$BM25(d_j, q_{1:N}) = \sum_{i=1}^N \left\{ IDF(q_i) \cdot \frac{TF(q_i, d_j) \cdot (k + 1)}{TF(q_i, d_j) + k \cdot (1 - b + b \cdot \frac{|d_j|}{L})} \right\}$$

i.e. calculate BM25 for each term/word in the query and sum them up  
[for example  $BM25(\text{"president lincoln"}) = BM25(\text{"president"}) + BM25(\text{"lincoln"})$ ]

$$L = \sum_i |d_i| / N$$

$$IDF(q_i) = \log_e \frac{N - DF(q_i) + 0.5}{DF(q_i) + 0.5}$$

- The equation gives the BM25 score for a query consisting of the words  $q_{1:N}$  given a document  $d_j$
- $TF(q_i, d_j)$  is the count of the number of times word  $q_i$  appears in document  $d_j$
- $|d_j|$  is the length of the document  $d_j$  in words
- $L$  is the average document length in the corpus (see above)
- $k$  and  $b$  can be tuned by cross validation, typical values are  $k = 2$  and  $b = 0.75$
- $IDF(q_i)$  is the inverse document frequency of word  $q_i$
- $DF(q_i)$  gives the number of documents that contain the word  $q_i$
- $N$  (in BM25 equation) is the number of terms in query
- $N$  (in IDF equation) is the total number of documents

## 22.3.1 BM25 scoring function

$$BM25(d_j, q_{1:N}) = \sum_{i=1}^N \left\{ IDF(q_i) \cdot \frac{TF(q_i, d_j) \cdot (k + 1)}{TF(q_i, d_j) + k \cdot (1 - b + b \cdot \frac{|d_j|}{L})} \right\}$$

i.e. calculate BM25 for each term/word in the query and sum them up  
[for example  $BM25(\text{"president lincoln"}) = BM25(\text{"president"}) + BM25(\text{"lincoln"})$ ]

$$L = \sum_i |d_i| / N$$

$$IDF(q_i) = \log_e \frac{N - DF(q_i) + 0.5}{DF(q_i) + 0.5}$$

- The equation gives the BM25 score for a query consisting of the words  $q_{1:N}$  given a document  $d_j$
- $TF(q_i, d_j)$  is the count of the number of times word  $q_i$  appears in document  $d_j$
- $|d_j|$  is the length of the document  $d_j$  in words
- $L$  is the average document length in the corpus (see above)
- $k$  and  $b$  can be tuned by cross validation, typical values are  $k = 2$  and  $b = 0.75$
- $IDF(q_i)$  is the inverse document frequency of word  $q_i$
- $DF(q_i)$  gives the number of documents that contain the word  $q_i$
- $N$  (in BM25 equation) is the number of terms in query
- $N$  (in IDF equation) is the total number of documents

### Problem:

A search query “president lincoln” is being scored against a document D. The terms “president” and “lincoln” both appear just once in D. The length of this document D is 90% of the average length of all documents in the corpus. There are 40,000 documents that contain the term “president” and there are 300 documents that contain the term “lincoln”. Assume that IDF for all queries is 1, and  $k = 1.2$  and  $b = 0.75$ .

What is the BM25 score for the query against the document D?

### Solution:



## 22.3.4 The problem of tyranny of TF scores

Does Term Frequency (e.g. BM25 scoring) always work?

- If the query is “IBM”, how do we make sure that IBM’s home page, `ibm.com`, is the first result, even if another page mentions the term “IBM” more frequently?
- One possible solution: Check for in-links
  - Say that each in-link is a vote for the quality of the linked-to page
  - The **more in-links in a page, the more likely that we will score it higher**
  - `ibm.com` will have many in-links (links to the page), so it should be ranked higher (this assumes that our algorithm can rank pages considering the number of links from other pages)
  - But, if we counted in-links, then it would be possible for a Web spammer to **create a network of pages and have them all point to a page of her choosing**, increasing the score of that page

**ibm.com**

`ibm.com/ddd`  
`ibm.com/xxx`  
`ibm.com/ddd`  
`ibm.com/xxx`  
`ibm.com/ddd`  
`ibm.com/xxx`  
`ibm.com/ddd`

**spam.com**

`ibm.com/ddd`  
`ibm.com/xxx`  
`ibm.com/ddd`  
`ibm.com/xxx`  
`ibm.com/ddd`  
`ibm.com/xxx`  
`ibm.com/ddd`

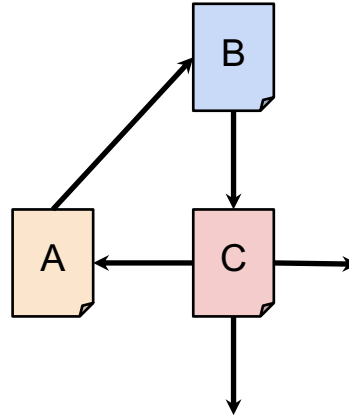
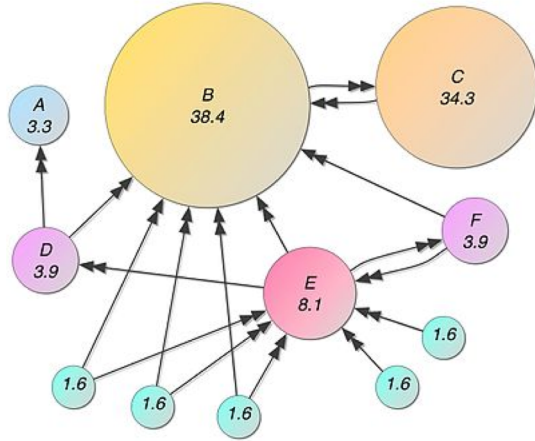
## 22.3.4 The PageRank algorithm

- To solve the tyranny of TF scores problem, the PageRank algorithm is designed
  - Idea: weight links from high-quality sites more heavily
- A high-quality site is the one that is linked-to by other high-quality site
  - yes, a recursive definition!
- PageRank can be computed by an iterative procedure: start with all pages having  $PR(p) = 1$ , and iterate the algorithm, updating ranks until they converge

$$PR(p) = \frac{1-d}{N} + d \sum_i \frac{PR(in_i)}{C(in_i)}$$

# PageRank algorithm (example)

Given the number of web pages  $N = 3$ , and the damping parameter  $d = 0.7$ . Calculate the PageRank of the pages A, B, and C. Links between the pages are shown in the graph below.



$$PR(p) = \frac{1 - d}{N} + d \sum_i \frac{PR(in_i)}{C(in_i)}$$

- $PR(p)$  is the PageRank of page  $p$
- $N$  is the total number of pages in the corpus
- $in_i$  are the pages that link in to  $p$
- $C(in_i)$  is the count of the total number of out-links on page  $in_i$
- $d$  is a damping factor constant

# Summary

- Probabilistic language models based on n-grams recover a surprising amount of information about a language
- Language models can have millions of features, so feature selection and preprocessing of data to reduce noise is important
- Text classification can be done with Bayes n-gram models or with any classification algorithms (machine learning)
- BM25 scoring function and PageRank use slightly different techniques to rank documents/pages