

So many possible models - unigram, bigram, trigram, ...
 different values of λ - 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 corpus
- Step 3: Evaluate the model on the validation corpus, and 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

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

$$P(H) \rightarrow P(e|H) \rightarrow P(l|e) \rightarrow P(t|l) \dots$$

We know $P(\text{Text} | \text{Language})$, how can we get the most probable language?

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

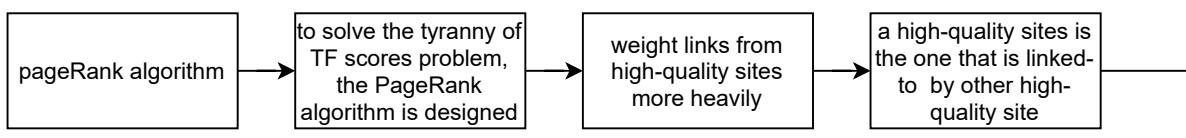
What happens to ℓ^* if any of $P(c_i | c_{i-2:i-1}, L)$ is 0?

the trigram models for each language can be learned from language corpus, but how to calculate $p(l)$

we may already have some estimates of $p(l)$
 the exact number we select for these priors is not critical

- Probability estimate is given by

$$\hat{P}(c_i | c_{i-2:i-1}) = \lambda_3 P(c_i | c_{i-2:i-1}) + \lambda_1 P(c_i | c_{i-1}) + \lambda_2 P(c_i | \text{start})$$
- The parameters λ_i can be fixed or they can be trained

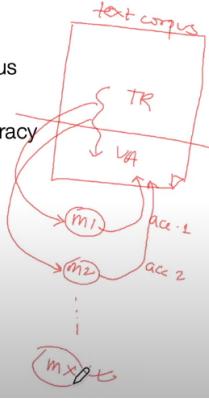


pageRank can be computed by an iterative procedure:
 start with all pages having $PR(p)=1$, and iterate the algorithm, updating ranks until they converge

trigram, interpolated smoothing with choose?

on:
us and a validation corpus
the training data
and obtain average accuracy
abilities are.

accuracy of the model



language modeling approach

Language Modeling approach (approach 1)

- Define one n-gram model for $P(\text{Message} | \text{Spam})$ by training on spam folder and one model for $P(\text{Message} | \text{Ham})$ by training on the inbox (basically, train two models)
- Classify a new message using Bayes' rule:

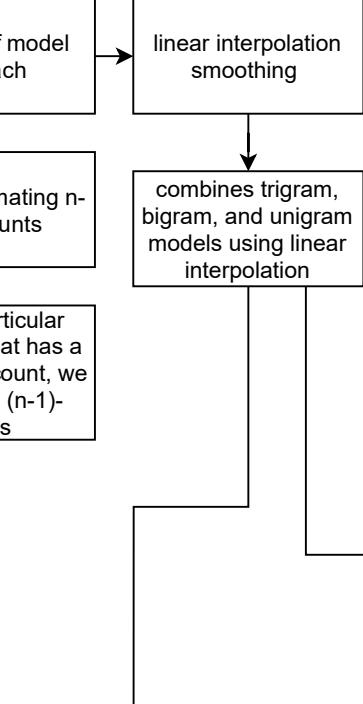
$P(c) \text{ is estimated by simply counting total number of spam and ham messages}$

$$\underset{c \in \{\text{spam, ham}\}}{\operatorname{argmax}} P(c | \text{message}) = \underset{c \in \{\text{spam, ham}\}}{\operatorname{argmax}} P(\text{message} | c) P(c)$$

Measuring acc

- Say that the
- The sum of a
- For a trigram small number
- If we have a underflow ca

we can describe probability of sequence using perplexity



$$\lambda_2 P(c_i | c_{i-1}) + \lambda_1 P(c_i)$$

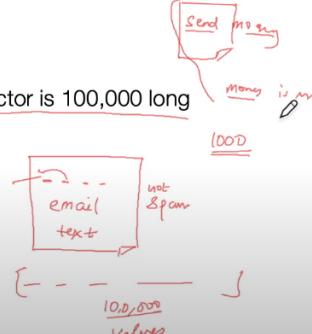
$$\lambda_3 + \lambda_2 + \lambda_1 = 1$$

ained.

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

| | Word 1 | ... | Word N | Spam |
|---------|--------|-----|--------|------|
| Email 1 | 2 | 0 | 8 | 0 |
| Email 2 | 5 | 0 | 5 | 1 |



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

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

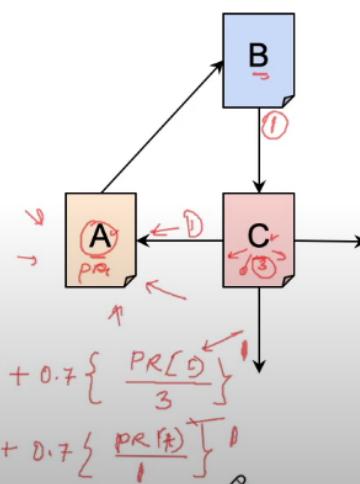
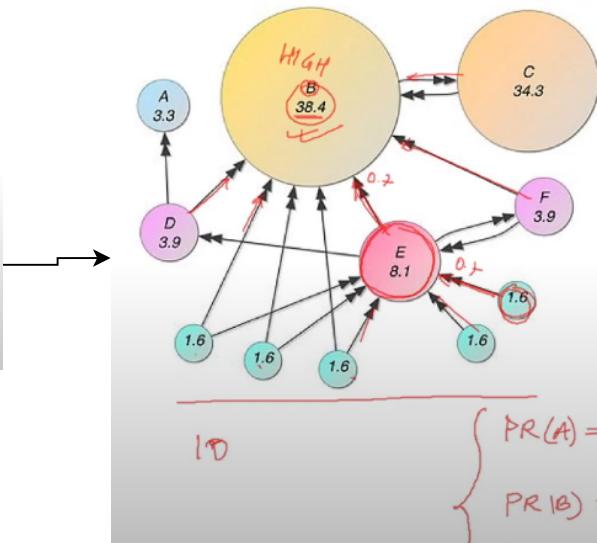
Three facts

- (1) Term
- (2) Inve
- (3) Len

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

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

constant



$$\left\{ \begin{aligned} PR(A) &= \frac{1 - 0.7}{3} + 0.7 \left\{ \frac{PR(D)}{3} \right\} \\ PR(B) &= \dots + 0.7 \left\{ \frac{PR(F)}{1} \right\} \end{aligned} \right.$$

accuracy can be inconvenient

model predicts the probability of $P("abc")$, $P("abd")$, $P("abe")$, etc.

all these probabilities must be 1 (it is a multi-class classification problem)

in model, if there are 1 million trigrams, the probability of each trigram (on average) is a

er 1/million

4-gram model, these probability values are very small numbers, and floating-point

can become an issue

be the
of a
using
y

Authors build a unigram, bigram, and trigram models over the words in the book
and then randomly sample 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?

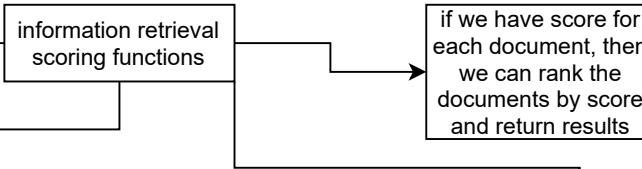


Perplexity

- Unigram model = 891 ✓
- Bigram model = 142 ✓
- Trigram model = 91 ↗

If so, why should we not build 4-gram or 5-gram word models?

IR is the task of
finding documents
that are relevant to a
user's need for
information



ors affect the weight of each query term:

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

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"

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

scoring function takes
a document and a
query and returns a
numeric score

for $d = 0.7$. Calculate the PageRank of the

graph below.

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

0.2

↓

3

↓

4

↓

5

↓

6

↓

7

↓

a page that
links to p

- $PR(p)$ 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 outgoing links on page in_i
- d is a damping factor constant

$$PR(\alpha) = 1 + 0.7 \cdot \frac{PR(\alpha')}{\gamma}$$

