



More LMs and Naïve Bayes

CS 780/880 Natural Language Processing Lecture 8

Samuel Carton, University of New Hampshire



Last lecture

Key idea: Probabilistic language modeling

Concepts

- Conditional probability
- Chain rule
- N-gram models
- Uses of language models
 - Generation
 - Evaluation
- Perplexity

Unigram model

Basic idea: model the text as the individual words occurring independently

- Parametrized by corpus token frequencies

$$P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)})$$

What's the problem with this?

Bigram model

Basic idea: model text as words being dependent on **only** the prior word

- Parameterized by token co-occurrence frequencies

$$P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)} | w^{(i-1)})$$

A bigram model is a type of **Markov Chain**



Example: identifying the hater

Enron dataset: ~500,000 emails from Enron company, released due to scandal

- My corpus: 2000 emails each sampled from the top 5 emailers

Question: Who sent the email “You are a huge jerk and I hate you”?



NLTK Language model functionality

```
1 from nltk.util import bigrams
2
3 # NLTK has functionality for generating all N-grams from a sequence of tokens
4 list(bigrams(enron_df['tokenized'][0]))[0:10]
```

```
[('what', 'are'),
 ('are', 'you'),
 ('you', 'talking'),
 ('talking', 'about'),
 ('about', '?'),
 ('?', 'sandra'),
 ('sandra', 'dial'),
 ('dial', '04/24/2000'),
 ('04/24/2000', '03:36'),
 ('03:36', 'pm')]
```



NLTK Language model functionality

```
1 # We'll take a look at the first email in our dataset
2 print(enron_df['message_text'][0])
```

What are you talking about?

Sandra Dial
04/24/2000 03:36 PM
To: Chris Germany/HOU/ECT@ECT
cc: Scott Hendrickson/HOU/ECT@ECT
Subject: Re: Feb 00 (RE: Voice Mail)

Scott-- Thanks.

Dude (CG),

I knew you had to be involved with this somehow.... :-)

Can you help me with this, please. Thanks. Let me know if you need more info
(I really have no more info other than this e-mail thread). Thanks man.

S
x5-7213

----- Forwarded by Sandra Dial/HOU/ECT on 04/24/2000 03:30
PM -----

Scott Hendrickson
04/24/2000 03:31 PM
To: Sandra Dial/HOU/ECT@ECT
cc:
Subject: Re: Feb 00 (RE: Voice Mail)

Well, just in that amount of time, I found out more about the deal...

I was inherited from CES. CES had originally made the sale and when we
bought their book of business, we ended up with that sale. You may want to
speak with Chris Germany about it, he was very involved with the CES
assimilation.

Scott



NLTK Language model functionality

```
1 # It also has functionality for adding "beginning-of-sequence" and "end-of-sequence" tokens to a
2 # token sequence
3
4 from nltk.util import pad_sequence
5 list(pad_sequence(enron_df['tokenized'][0],
6                 pad_left=True,
7                 left_pad_symbol("<s>"),
8                 pad_right=True,
9                 right_pad_symbol("</s>"),
10                n=2))[0:10]
```

```
['<s>',
 'what',
 'are',
 'you',
 'talking',
 'about',
 '?',
 'sandra',
 'dial',
 '04/24/2000']
```




NLTK Language model functionality

```
1 # Both of these utilities are packaged into a convenient pipeline
2 # which will take in a collection of (tokenized) texts, pad each one
3 # with start/end tokens, and then count all n-grams up to the number you
4 # give it (in our case, 2)
5 from nltk.lm.preprocessing import padded_everygram_pipeline
6
7 # I am training this model on just the first text
8 train, vocab = padded_everygram_pipeline(2, enron_df['tokenized'][0:1])
```

```
1 # These things are generators, so they don't do anything until called on
2 print(train)
3 print(vocab)
```

```
<generator object padded_everygram_pipeline.<locals>.<genexpr> at 0x7cf18a443060>
<itertools.chain object at 0x7cf16ce91a20>
```



NLTK Language model functionality

```
1 # Both of these utilities are packaged into a convenient pipeline
2 # which will take in a collection of (tokenized) texts, pad each one
3 # with start/end tokens, and then count all n-grams up to the number you
4 # give it (in our case, 2)
5 from nltk.lm.preprocessing import padded_everygram_pipeline
6
7 # I am training this model on just the first text
8 train, vocab = padded_everygram_pipeline(2, enron_df['tokenized'][:1])
```

```
1 # These things are generators, so they don't do anything until called on
2 print(train)
3 print(vocab)
```

```
<generator object padded_everygram_pipeline.<locals>.<genexpr> at 0x7cf18a443060>
<itertools.chain object at 0x7cf16ce91a20>
```



NLTK Language model functionality

```
1 # Once we have training corpus and a vocabulary, we can create a MLE model to
2 # fit across it
3 from nltk.lm import MLE
4 lm = MLE(2)
5
6 #Then we fit the model
7 lm.fit(train, vocab)
8 print(lm.vocab)
9 print(len(lm.vocab))
```

```
<Vocabulary with cutoff=1 unk_label='<UNK>' and 110 items>
110
```



NLTK Language model functionality

```
1 # We can look up a piece of text in the vocabulary
2 lm.vocab.lookup(enron_df['tokenized'][0])[0:10]
```

```
('what',
 'are',
 'you',
 'talking',
 'about',
 '?',
 'sandra',
 'dial',
 '04/24/2000',
 '03:36')
```

```
1 # And any token it doesn't recognize it will replace with <UNK>
2 lm.vocab.lookup(enron_df['tokenized'][1])[0:10]
```

```
('i', '<UNK>', 'we', '<UNK>', 'about', 'this', 'the', 'other', '<UNK>', '.')
```



NLTK Language model functionality

```
1 # We can look up a piece of text in the vocabulary
2 lm.vocab.lookup(enron_df['tokenized'][0])[0:10]
```

```
('what',
 'are',
 'you',
 'talking',
 'about',
 '?',
 'sandra',
 'dial',
 '04/24/2000',
 '03:36')
```

```
1 # And any token it doesn't recognize it will replace with <UNK>
2 lm.vocab.lookup(enron_df['tokenized'][1])[0:10]
```

```
('i', '<UNK>', 'we', '<UNK>', 'about', 'this', 'the', 'other', '<UNK>', '.')
```



NLTK Language model functionality

```
1 # We can see how many ngrams it found
2 print(lm.counts)
```

<NgramCounter with 2 ngram orders and 443 ngrams>

```
1 # We can count how many instances of the 'a' unigram it found
2 lm.counts['the']
```

3

```
1 # And how many counts of the "in that" bigram
2 lm.counts[['in']]['that']
```

1

NLTK Language model functionality



```
1 # We can look up probabilities of particular unigrams under the model!  
2 lm.score('the')
```

```
0.013513513513513514
```

```
1 lm.score('a')
```

```
0.0
```

```
1 # And we can do the same for bigrams  
2 lm.score('that', ['in'])
```

```
1.0
```



NLTK Language model functionality

```
1 # It's often useful to report log-probabilities rather than raw probabilities
2 # because of underflow
3 lm.logscore('the')
```

-6.209453365628949

```
1 # Although this can result in weird arithmetic when dealing with unknown tokens
2 lm.logscore('a')
```

-inf



NLTK Language model functionality

```
1 # Finally, we can generate new sequences according to the model
2 lm.generate(1, random_seed=3)
```

'00'

```
1 lm.generate(10, random_seed=7)
```

['</s>', '---', '---', '---', '---', '---', '---', '---', '-', ')']



Modeling our Enron data

```
1 from nltk.lm.preprocessing import padded_everygram_pipeline
2
3 # First I train a global model on all the emails
4 global_train, global_vocab = padded_everygram_pipeline(2, enron_train_df['tokenized'])
5
6 # Use a so-called "Lidstone" model to do the add-one smoothing I talked about last lecture
7 from nltk.lm import Lidstone
8 global_lm = Lidstone(order=2, gamma=1)
9
10 global_lm.fit(global_train, global_vocab)
```

```
1 # Here's an email generated by our global LM
2 display(' '.join(global_lm.generate(100, text_seed='<s>', random_seed=8)))
```

', vladimir gorny , a consumer welfare (hour : dicarlo , ccampbell @ ect , 2000.= if the talks between air credits will attempt to discover the escrow as of the talking to : 1.0 content-type : 7bit x-from : pverde_5_devers interchg_id : can the borland database ! an gsa would remain responsible for ban kruptcy filing in abilene , john.anderson @ enron.com , '' < /o=enron/ou=na/cn=recipients/cn=pdavis1 > , sarah novosel/corp/enron @ enron.com , monika.c ausholli @ enron employees at \$ 50 books . i am on behalf of students and appointments . = and=20 electricity rates could review what .'



Modeling our Enron data

```
1 # Then I try training an individual LM for each person
2 personal_lms = {}
3
4 from nltk.lm import Lidstone
5
6 #I use the pandas groupby function to split up the training set by name
7 for group_name, group_df in enron_train_df.groupby('name'):
8     group_train, group_vocab = padded_everygram_pipeline(2, group_df['tokenized'])
9     group_lm = Lidstone(order=2, gamma=12)
10    group_lm.fit(group_train, group_vocab) # Use the same (global) vocabulary for everyone
11    personal_lms[group_name] = group_lm
```

```
1 personal_lms
```

```
{'chris': <nltk.lm.models.Lidstone at 0x7cf168b35f60>,
'jeff': <nltk.lm.models.Lidstone at 0x7cf1683fd120>,
'kay': <nltk.lm.models.Lidstone at 0x7cf167ba5ae0>,
'pete': <nltk.lm.models.Lidstone at 0x7cf16722a2f0>,
'vince': <nltk.lm.models.Lidstone at 0x7cf16729d000>}
```



Modeling our Enron data

```
1 new_email = "you are a huge jerk and I hate you"
```

```
1 import numpy as np
2 def score_new_email(email_text:str, lm:nltk.lm.models.MLE):
3     tokens = preprocess_enron(email_text) #tokenize the new email
4     tokens = list(pad_sequence(tokens, # add start/end tokens
5                               pad_left=True,
6                               left_pad_symbol("<s>",
7                               pad_right=True,
8                               right_pad_symbol("</s>",
9                               n=2))
10    print(tokens)
11    tokens = lm.vocab.lookup(tokens) # Look up the tokens in the vocab
12    print(tokens)
13    scores = []
14    for i in range(1, len(tokens)): # Look up score of each token given previous token
15        scores.append(lm.logscore(tokens[i], [tokens[i-1]]))
16
17    return np.sum(scores), scores
```



Modeling our Enron data

```
1 for name, personal_lm in personal_lms.items():
2     score, scores = score_new_email(new_email, personal_lm)
3     print(name, np.round(score,2))
4     print('\tScores',np.round(scores,2))
5     print('---')
```

```
['<s>', 'you', 'are', 'a', 'huge', 'jerk', 'and', 'i', 'hate', 'you', '</s>']
('<s>', 'you', 'are', 'a', 'huge', '<UNK>', 'and', 'i', 'hate', 'you', '</s>')
chris -127.87
      Scores [-12.04 -10.31 -12.23 -13.79 -13.99 -13.99 -9.69 -13.8 -13.99 -14.01]
---
['<s>', 'you', 'are', 'a', 'huge', 'jerk', 'and', 'i', 'hate', 'you', '</s>']
('<s>', 'you', 'are', 'a', 'huge', 'jerk', 'and', 'i', 'hate', 'you', '</s>')
jeff -135.48
      Scores [-12.85 -11.56 -12.66 -13.88 -14.84 -14.84 -10.54 -14.74 -14.73 -14.85]
---
['<s>', 'you', 'are', 'a', 'huge', 'jerk', 'and', 'i', 'hate', 'you', '</s>']
('<s>', 'you', 'are', 'a', 'huge', '<UNK>', 'and', 'i', 'hate', 'you', '</s>')
kay -126.97
      Scores [-12.65 -8.78 -12.95 -13.86 -14.05 -14.05 -9.66 -13.5 -14.05 -13.41]
---
['<s>', 'you', 'are', 'a', 'huge', 'jerk', 'and', 'i', 'hate', 'you', '</s>']
('<s>', 'you', 'are', 'a', '<UNK>', '<UNK>', 'and', 'i', '<UNK>', 'you', '</s>')
pete -113.7
      Scores [-11.44 -11.36 -11.36 -11.42 -11.35 -11.35 -11.36 -11.35 -11.35 -11.36]
---
['<s>', 'you', 'are', 'a', 'huge', 'jerk', 'and', 'i', 'hate', 'you', '</s>']
('<s>', 'you', 'are', 'a', 'huge', '<UNK>', 'and', 'i', 'hate', 'you', '</s>')
vince -134.97
      Scores [-14.71 -9.92 -12.97 -14.14 -14.7 -14.7 -9.92 -14.72 -14.7 -14.5 ]
---
```

Bayes Rule





Conditional probability

When two variables may be dependent, then their joint probability is expressed as follows:

$$P(X, Y) = P(Y)P(X|Y) = P(X)P(Y|X)$$

If they happen to be independent, then $P(X|Y) = P(X)$ and $P(Y|X) = P(Y)$, so

$$P(X, Y) = P(Y)P(X) = P(X)P(Y)$$



Bayes Rule

It follows from

$$P(X, Y) = P(Y)P(X|Y) = P(X)P(Y|X)$$

that

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$



Examples

$$P(\text{Lung cancer}|\text{Cough}) = \frac{P(\text{Cough}|\text{Lung cancer})P(\text{Lung cancer})}{P(\text{Cough})}$$

$$P(\text{Conspiracy}|\text{Event}) = \frac{P(\text{Event}|\text{Conspiracy})P(\text{Conspiracy})}{P(\text{Event})}$$

$$= \frac{P(\text{Barista likes you}|\text{Smiles when they give you coffee})P(\text{Barista likes you})}{P(\text{Smiles when they give you coffee})}$$



Relative probabilities

Often we only care about the relative probability of two possible outcomes, rather than their true probability:

$$P(\text{Lung cancer}|\text{Cough}) \text{ vs. } P(\text{COVID}|\text{Cough})$$

$$\frac{P(\text{Cough}|\text{Lung cancer})P(\text{Lung cancer})}{P(\text{Cough})} \text{ vs. } \frac{P(\text{Cough}|\text{COVID})P(\text{COVID})}{P(\text{Cough})}$$

Because we only care about the relative value, we can ignore the denominator

$$P(\text{Cough}|\text{Lung cancer}) \approx P(\text{Cough}|\text{COVID}) \approx 1.0$$

~50 million COVID cases in 2022, ~300k new lung cancer cases in 2023

So $P(\text{COVID}) = .15$, and $P(\text{Lung cancer}) = 0.001$

So $P(\text{COVID}|\text{Cough})$ is **150 times** higher than $P(\text{Lung cancer}|\text{Cough})$

<https://www.cancer.org/cancer/lung-cancer/about/key-statistics.html>

https://covid.cdc.gov/covid-data-tracker/#trends_totalcases_select_00



Base rate fallacy

A lot of fallacious thinking comes from ignoring the **base rates** $P(X)$ and $P(Y)$ in $\frac{P(Y|X)P(X)}{P(Y)}$

$$P(\text{Hypothesis} | \text{Rare event}) = \frac{P(\text{Rare event} | \text{Hypothesis})P(\text{Hypothesis})}{P(\text{Rare event})}$$

$P(\text{Hypothesis})$ is often lower than you think

$P(\text{Rare event})$ is often higher than you think

https://en.wikipedia.org/wiki/Base_rate_fallacy

https://en.wikipedia.org/wiki/List_of_cognitive_biases

Naïve Bayes





Application to text

Classification:

$$P(\text{Class} \mid \text{Words})$$

$$P(\text{Class 0} \mid \text{Words}) \text{ vs. } P(\text{Class 1} \mid \text{Words})$$

$$\frac{P(\text{Words} \mid \text{Class 0})P(\text{Class 0})}{P(\text{Words})} \text{ vs. } \frac{P(\text{Words} \mid \text{Class 1})P(\text{Class 1})}{P(\text{Words})}$$

We can ignore $P(\text{Words})$, but how do we calculate:

- $P(\text{Words} \mid \text{Class 0})$
- $P(\text{Class 0})$
- $P(\text{Words} \mid \text{Class 1})$
- $P(\text{Class 1})$

Application to text

$$P(\text{Class } 0) = \frac{\# \text{Class } 0}{\# \text{Class } 0 + \# \text{Class } 1}$$

- And likewise for class 1

$P(\text{Words} \mid \text{Class } 0)$

- Build an n-gram model of all texts for which class is Class 0
- Use this model to estimate $P(\text{Words} \mid \text{Class } 0)$
- And likewise for Class 1

Naïve Bayes

Basic idea: apply Bayes rule to find relative likelihoods of $P(\text{Class } 0 \mid \text{Words})$ vs. $P(\text{Class } 1 \mid \text{Words})$, using **unigram model** for $P(\text{Words} \mid \text{Class } C)$

So if we consider words = $\{w_0, w_1, \dots, w_N\}$:

$$P(\text{Class } 0 \mid \text{Words}) \propto P(\text{Class } 0) \prod_{i=1}^N P(w_i \mid \text{Class } 0)$$

$$P(\text{Class } 1 \mid \text{Words}) \propto P(\text{Class } 1) \prod_{i=1}^N P(w_i \mid \text{Class } 1)$$



Read the SST-2 dataset

```
1 display(dev_df)
```

	sentence	label
0	it 's a charming and often affecting journey .	1
1	unflinchingly bleak and desperate	0
2	allows us to hope that nolan is poised to emba...	1
3	the acting , costumes , music , cinematography...	1
4	it 's slow -- very , very slow .	0
...
867	has all the depth of a wading pool .	0
868	a movie with a real anarchic flair .	1
869	a subject like this should inspire reaction in...	0
870	... is an arthritic attempt at directing by ca...	0
871	looking aristocratic , luminous yet careworn i...	1

872 rows × 2 columns



Preprocess and vectorize the data

```
1 from nltk import PorterStemmer
```

```
1 # for this dataset, the tokenization has already been done for us
2 stemmer = PorterStemmer()
3 def preprocess(s):
4 | return ' '.join([stemmer.stem(token) for token in s.strip().split(' ')])
```

```
1 train_df['preprocessed'] = train_df['sentence'].apply(preprocess)
2 dev_df['preprocessed'] = dev_df['sentence'].apply(preprocess)
```

```
1 from sklearn.feature_extraction.text import CountVectorizer
```

```
1 # Why are we using a CountVectorizer here instead of TF-IDF?
2
3 vectorizer = CountVectorizer()
4 train_X = vectorizer.fit_transform(train_df['preprocessed'])
5 dev_X = vectorizer.transform(dev_df['preprocessed'])
```



Build and evaluate the model

```
1 from sklearn.naive_bayes import MultinomialNB
```

```
1 # See https://scikit-learn.org/stable/modules/generated/sklearn.naive\_bayes.MultinomialNB.html#sklearn.naive\_bayes.MultinomialNB  
2 # for hyperparameter options  
3  
4 model = MultinomialNB()
```

```
1 model.fit(train_X, train_df['label'])
```

```
MultinomialNB()
```

```
1 dev_py = model.predict(dev_X)
```

```
1 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
```

```
1 def evaluate_predictions(y, py):  
2     print(f'Accuracy: {accuracy_score(y, py):.3f}')  
3     print(f'Precision: {precision_score(y, py):.3f}')  
4     print(f'Recall: {recall_score(y, py):.3f}')  
5     print(f'F1: {f1_score(y, py):.3f}')
```



Build and evaluate the model

Naïve Bayes:

```
1 # Evaluating on the dev set
2 evaluate_predictions(dev_df['label'], dev_py)
```

Accuracy: 0.807
Precision: 0.794
Recall: 0.840
F1: 0.816

```
1 # Evaluating on a sample of the training set
2 train_py = model.predict(train_X[0:1000])
3 evaluate_predictions(train_df['label'].iloc[0:1000], train_py)
```

Accuracy: 0.891
Precision: 0.894
Recall: 0.906
F1: 0.900

K-nearest-neighbors:

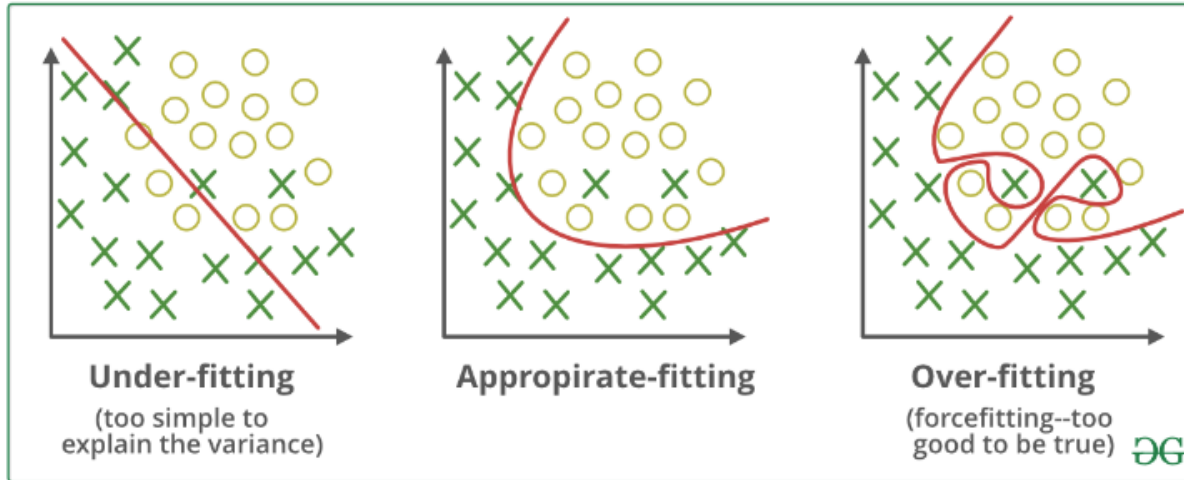
```
1 evaluate_model(dev_X, dev_y, classifier)
```

Accuracy: 0.742
Precision: 0.707
Recall: 0.842
F1: 0.769

```
1 evaluate_model(train_X[0:1000], train_y[0:1000], classifier)
```

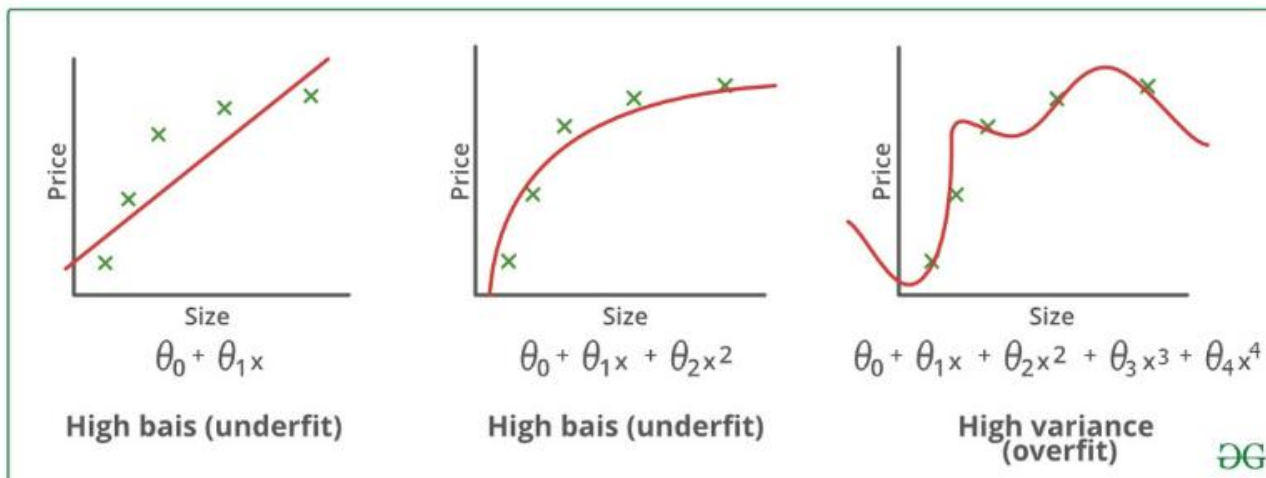
Accuracy: 0.948
Precision: 0.945
Recall: 0.959
F1: 0.952

Overfitting and underfitting



Overfitting: model overly tuned to quirks of the training data—doesn't generalize

Underfitted: model not tuned enough to training data—doesn't capture data structure



Related (but not identical) to **bias-variance trade-off**

- High bias \rightarrow underfitting
- High variance \rightarrow overfitting



Explaining the model

```
1 # We can get the (log) probability of each word for each class
2 print('Word log-probs:')
3 display(model.feature_log_prob_)
4
5 # It will have one row for each class and one column for each word in the vocabulary
6 print('Log-probs matrix shape:')
7 display(model.feature_log_prob_.shape)
8
```

```
Word log-probs:
array([[ -11.12495518,  -8.09240893, -10.20866444, ..., -12.51124954,
        -10.71949007, -10.90181162],
       [-11.03503481,  -9.50897851, -10.15956608, ..., -11.25817837,
        -12.64447273, -11.95132555]])
Log-probs matrix shape:
(2, 10106)
```



Explaining the model

```
1 # We can identify the words that were the biggest distinguishers by calculating
2 # the diff between the two rows
3 word_prob_diffs = model.feature_log_prob_[0] - model.feature_log_prob_[1]
4 word_prob_diffs
```

```
array([-0.08992036,  1.41656958, -0.04909837, ..., -1.25307117,
        1.92498266,  1.04951392])
```

```
1 # And then we can use numpy.argsort() and numpy.abs() to find the indices of the
2 # words with the biggest diff (positive or negative)
3 import numpy as np
4 sorted_diff_indices = np.argsort(np.abs(word_prob_diffs))
5 sorted_diff_indices
```

```
array([6103, 9942, 7833, ..., 6692, 6721, 9402])
```

```
1 # Numpy argsort always goes in ascending order, so to get the top K indices
2 # we have to grab the last K indices
3
4 # We can use -1 as the third part of our slice, to get these back in reverse order
5 k= 10
6 top_k_indices = sorted_diff_indices[:-k:-1]
```



Explaining the model

```
1 # Then we can find the words and values associated with those indices
2 vocab = vectorizer.get_feature_names_out()
3 top_words = vocab[top_k_indices]
4 top_diffs = word_prob_diffs[top_k_indices]
5
6 print(f'Top {k} distinguishing words in our Naive Bayes classifier')
7 for word, diff in zip(top_words, top_diffs):
8 | print(f'\tWord: "{word}" - Diff: {diff:.3f}')
```

```
Top 10 distinguishing words in our Naive Bayes classifier
Word: "unfunni" - Diff: 4.861
Word: "poorli" - Diff: 4.698
Word: "pointless" - Diff: 4.677
Word: "tiresom" - Diff: 4.588
Word: "eleg" - Diff: -4.410
Word: "unnecessari" - Diff: 4.410
Word: "badli" - Diff: 4.382
Word: "embrac" - Diff: -4.355
Word: "inept" - Diff: 4.338
```

Interpreting log-probability differences



If:

$$\log(P(w_i | \text{class 0})) - \log(P(w_i | \text{class 1})) = 4.8$$

Then:

$$\frac{P(w_i | \text{class 0})}{P(w_i | \text{class 1})} = e^{4.8} = 2.718^{4.8} = 121.51$$

Meaning that w_i (“unfunny” in this case) is **121.51** times more likely to occur in class 0 than in class 1



Explaining individual predictions

```
1 sentence = 'the movie was pretty awful : not good at all .'
2 preprocessed_sentence = preprocess(sentence)
3 sentence_x = vectorizer.transform([preprocessed_sentence])
4 py = model.predict(sentence_x)
5 print(f'Model prediction for "{preprocessed_sentence}": {py}')
```

Model prediction for "the movi wa pretti aw : not good at all .": [0]

```
1 # we can find the vocab indices for the tokens in the sentence
2 sentence_tokens = preprocessed_sentence.split(' ')
3 token_indices = [vectorizer.vocabulary_[token] for token in sentence_tokens \
4 | | | | | | | | | | if token in vectorizer.vocabulary_] #there's one or two stopwords to ignore
5 token_indices
```

[8892, 5827, 9702, 6828, 680, 6067, 3791, 615, 339]



Explaining individual sentences

Is this overfitting?

```
1 # Then we can do the same thing as we did with the top indices above
2 sentence_words = vocab[token_indices]
3 sentence_diffs = word_prob_diffs[token_indices]
4 print(f'Class probability differences for tokens in the sentence:')
5 for word, diff in zip(sentence_words, sentence_diffs):
6 | print(f'\tWord: "{word}" - Diff: {diff:.3f}')
```

```
Class probability differences for tokens in the sentence:
  Word: "the" - Diff: -0.032
  Word: "movi" - Diff: 0.196
  Word: "wa" - Diff: 0.836
  Word: "pretti" - Diff: -0.643
  Word: "aw" - Diff: 1.574
  Word: "not" - Diff: 0.679
  Word: "good" - Diff: -1.021
  Word: "at" - Diff: 0.080
  Word: "all" - Diff: 0.111
```



Concluding thoughts

Naïve bayes: application of Bayes Rule + unigram language modeling to classification

Huge deal in 1998

Limitations?