



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**



Markov Chain

Definition: a discrete stochastic process with the Markov property :

$$P(X_t | X_{t-1}, \dots, X_1) = P(X_t | X_{t-1})$$

fully determined by a probability transition matrix P which defines the transition probabilities:

$$(P_{ij} = P(X_t = j | X_{t-1} = i))$$

and an initial probability distribution specified by the vector x where:

$$x_i = P(X_0 = i)$$

Should hopefully be clear why a bigram model is this

In general, we're often concerned about the **stationary distribution** of X_t over time

- Not so much in NLP

https://stephens999.github.io/fiveMinuteStats/markov_chains_discrete_intro.html



Stationary distribution

Corpus:

“i am a person .”

“a person am i .”

Stationary distribution: from the CPT alone, what is the probability that the word W_t occurring at time t will be a particular word w ?

- If you think of words as states, can also be: what percentage of our total time will we spend in each state?
- Important for some tasks, not really so much for NLP
- Can be solved by doing eigendecomposition on the CPT

CPT:

	i	am	a	person	.	[END]
[START]	0.5	0	0.5	0	0	0
i	0	0.5	0	0	0.5	0
am	0.5	0	0.5	0	0	0
a	0	0	0	1	0	0
person	0	0.5	0	0	0.5	0
.	0	0	0	0	0	1.0



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)
```

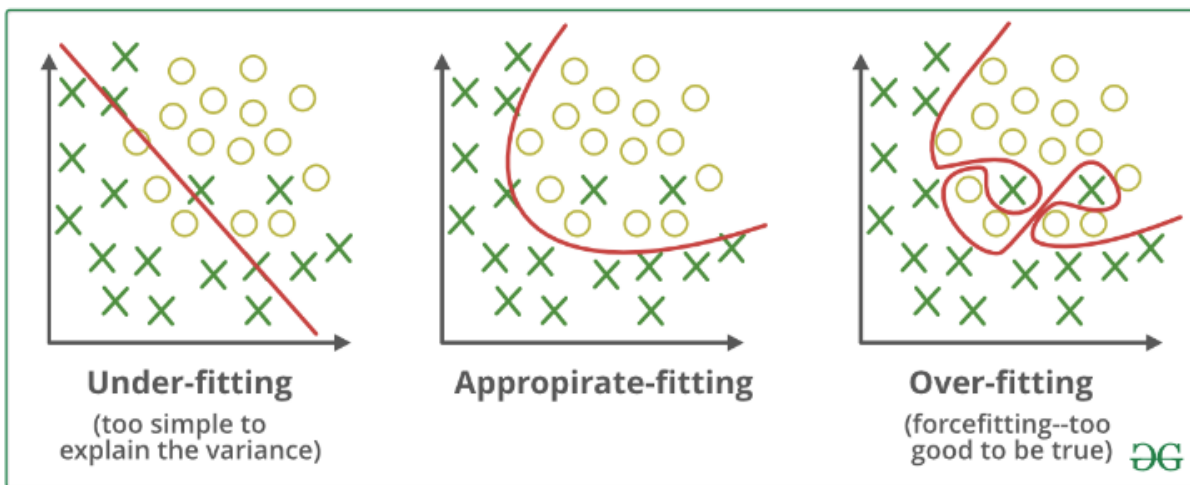
Accuracy: 0.742
Precision: 0.707
Recall: 0.842
F1: 0.769

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

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

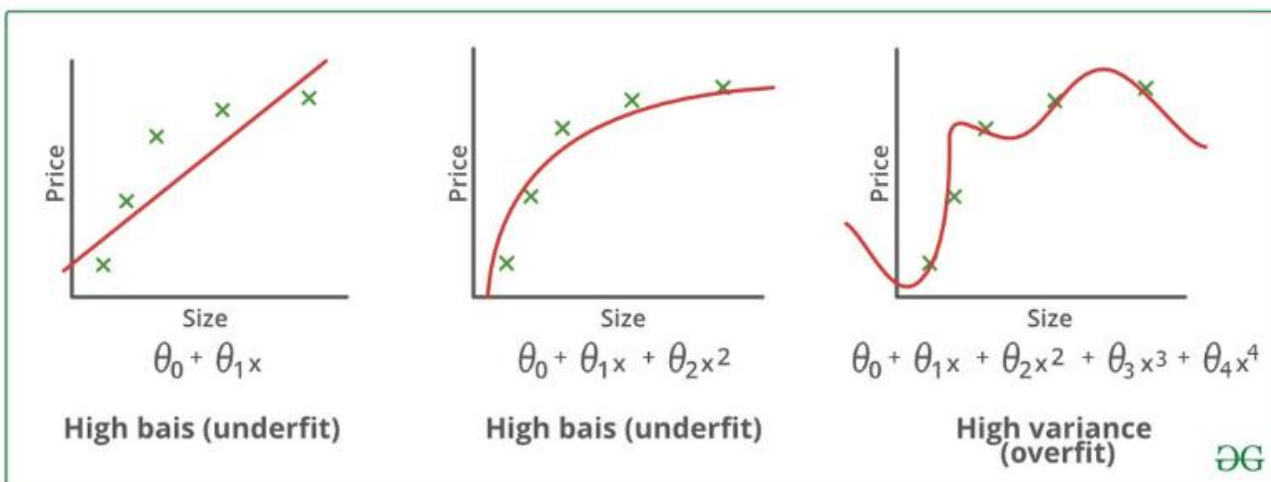


Overfitting and underfitting



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?

