

Sentiment analysis part 2

Module 2



Plan for today

- Text representation
 - Sub-words
 - Byte-Pair Encoding (BPE)
- Categorisation methods:
 - Review of Bayesian classification
 - Multinomial Naïve Bayes
 - Smoothing
 - Laplace for classification
 - Kneser-Ney for n-grams
 - Other categorization methods: multilabel and stream categorization
- Sentiment analysis



Text representation



Questions

- Difference between One Hot and BOW encoding
- What additional information is captured by n-gram, but not BOW?
- Which is faster: stemming or lemmatisation? Which is more accurate?
- Should stop words be removed?
 - Which are stop words in this sentence?
 "The well is very deep"



Dealing with rare words in BOW

- Consider the following scenario:
 - Word **w1** (called out-of-vocabulary word OOV) occurs in test instance, but not in training set. How can we use similar words to possibly improve the classification?
- 1. Find word vector for w1.
- 2. Find the most similar word w1' to w1 in training set vocabulary.
- 3. Add w1' to the test instance.
- 4. Classify the instance
- May not always work, e.g. WV king, queen.
- Use Byte Pair Encoding for subwords



Sub-Words

Word representation cannot handle unseen word or rare word well. Character-level representation is one of the solutions to overcome out-of-vocabulary (OOV). However, it may be too fine-grained and miss some important information.

Subword is in between word and character. It is not too finegrained while able to handle unseen words and rare words.

- Byte Pair Encoding (BPE: Jurafsky et al. ch. 2.4.3)
- WordPiece: similar to BPE (Jurafsky, ch. 13.2.1)
- We will discuss BPE



Tokeniser: Byte Pair Encoding: BPE

The dataset is divided into training and test

Token learner

- The algorithm begins with the set of symbols equal to the set of characters. Each word
 is represented as a sequence of characters plus a special end-of-word symbol.
- At each step of the algorithm, we count the number of symbol pairs, find the most frequent pair ('A', 'B'), and replace it with the new merged symbol ('AB').
- Add the new merged token to the vocabulary
- We continue to count and merge, creating new longer and longer character strings, until we've done k merges; k is a parameter of the algorithm.

Token segmenter

We apply character strings found in the training set to the test set



Byte Pair Encoding example

Corpus and Vacabulary:

Corpus

Vocabulary

```
_, d, e, i, l, n, o, r, s, t, w
```



Byte Pair Encoding example cont.

Step 1 Corpus

Vocabulary

Step 2 Corpus

Vocabulary

```
5 low__ _, d, e, i, l, n, o, r, s, t, w, er, er_
2 lowest_
6 newer_
3 wider_
2 new
```



Byte Pair Encoding example cont.

Step 3 Corpus

```
2 lowest
6 ne w er
3 w i d er
2 ne w
```

Vocabulary

```
5 l o w , d, e, i, l, n, o, r, s, t, w, er, er_, ne
```

8 steps, k #merges k=8

Next 5 steps merges

```
(ne, w)
(1, 0)
(10, W)
(new, er)
(low, )
```

Current Vocabulary

```
_, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new
_, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo
_, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo, low
_, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo, low, newer_
_, d, e, i, l, n, o, r, s, t, w, er, er_, ne, new, lo, low, newer_, low_
```



BPE algorithm

This is token learner part

```
V \leftarrow all unique characters in C # initial set of tokens is characters for i = 1 to k do # merge tokens k times t_L, t_R \leftarrow Most frequent pair of adjacent tokens in C # make new token by concatenating V \leftarrow V + t_{NEW} # update the vocabulary
```

Replace each occurrence of t_L , t_R in C with t_{NEW}

function BYTE-PAIR ENCODING(strings C, number of merges k) **returns** vocab V

Figure 2.13 The token learner part of the BPE algorithm for taking a corpus broken up into individual characters or bytes, and learning a vocabulary by iteratively merging tokens. Figure adapted from Bostrom and Durrett (2020).

and update the corpus

This is token segmenter part

Apply learned tokens to a new text n same order as learned tokens

return V

- Frequencies of words don't matter in segmenter
- Replace shorter tokens with longer from the learned tokens
- Rare words will get segmented and common subwords identified.



Naïve Bayes vs Multinomial Naïve Bayes



Exercise on Bayesian rule

$$P(h|X) = \frac{P(X|h) \cdot P(h)}{P(X)}, \quad likelihood = P(X|h) \cdot P(h)$$
$$h_{MAP} = \arg\max_{h \in H} P(h|X) = \arg\max_{h \in H} \frac{P(X|h) \cdot P(h)}{P(X)} = \arg\max_{h \in H} P(X|h) \cdot P(h)$$

A patient takes a lab test and the result comes back positive. The test returns a correct positive result (+) in only 98% of the cases in which the disease is actually present, and a correct negative result (-) in only 97% of the cases in which the disease is not present. Furthermore, 0.008 of the entire population have this cancer.

P(cancer) = 0.008, P(not cancer) = 0.992, these are called prior probabilities

$$P(+ | cancer) = 0.98 P(- | cancer) = 0.02$$

P(+ | not cancer) = 0.03 P(- | not cancer) = 0.97

Does the patient have cancer or not? To answer this, calculate P1 = P(not cancer $| + \rangle$ and P2 = P(cancer $| + \rangle$, and then calculate likelihood ratio P1/P2. If ratio > 1, the patient does not have cancer (according to this calculation)

How to Classify Using Naïve Bayes

1. Using training dataset:

- a) For each class, calculate P(h) as $\frac{n_h}{n}$, where n is number of instances, n_h is number of instances in class y for hypothesis h
- b) For each feature in training set, calculate $P(x_j|h)$ as $\frac{n_j}{n_h}$, where n_h is number of instances in class h, n_j is number of instances in class h with feature x_j
- c) Store the above calculations for later

2. Given test instance with features : $X = \{x_1, ..., x_j, ...\}$:

- a) Retrieve all P(h) and all $P(x_i|h)$ for these features
- b) Calculate $P(X|h) \cdot P(h) = P(h) * P(x_1|h) * \cdots * P(x_j|h) * \cdots$ for all features x in the test instance, and all hypotheses h, corresponding to all m classes. So we will have a vector of m values, one value for each class
- c) Take arg max of that vector, as the predicted class



NB calculations: text example

Data (words are just letters):

Spam	Ham
b b b	a b a b a b a
b c b b c c	acacaca
a a a	a a a
b a b a b	-

We want to classify a new document into HAM or SPAM



BNB calculations: example

	ಡ	b	د	class
ei	0	l	Ó	+
وء	0	l	1	+
ez	l	0	0	+
٤4	l	ι	0	+
وح	ι	ι	0	_
eb	l	0	1	_
67	ι	0	Ō	_
68	0	0	0	_

- BNB training:
 - For each class: num docs with word/num all docs
 - Spam: P(a|+), P(b|+), P(c|+) = (2/4, 3/4, 1/4) = (0.5; 0.75; 0.25)
 - Ham: P(a|-), P(b|-), P(c|-) = (3/4, 1/4, 1/4) = (0.75; 0.25; 0.25)

Testing for instance x=(a b -) = (1,1,0):

Classify x=(a b -) using MAP:

- P(x|+) = 0.5 *0.75 *(1-0.25) = 0.28;
- P(x|-) = 0.75 *0.25*(1-0.25) = 0.14
- Result: SPAM

Classify using ML: same, since P(-) = P(+) = 0.5



BNB calculations: example with 0

	ત	b	د	class
ei	0	l	Ó	+
ez	0	ι	1	+
ez	0	0	0	+
64	0	l	0	+
وح	ι	0	0	_
eb	ι	0	1	_
67	ι	0	O	_
C8	0	0	0	_

- BNB training:
 - For each class: num docs with word/num all docs
 - Spam: P(a|+), P(b|+), P(c|+) = (0/4, 3/4, 1/4) = (0.0; 0.75; 0.25)
 - Ham: P(a|-), P(b|-), P(c|-) = (3/4, 0/4, 1/4) = (0.75; 0.0; 0.25)

Testing for instance x=(a b -) = (1,1,0):

Classify x=(a b -) using MAP:

- P(x|+) = 0.0 *0.75 *(1-0.25) = 0;
- P(x|-) = 0.75 *0.0*(1-0.25) = 0
- Result: ?



Smoothing



Why we need Smoothing?

- Some words do not appear in some training classes,
 e.g. if c column has all 0's in (+)
 - Solution: probabilistic smoothing (Laplace smoothing we will talk about soon)
- Some words appear in test but not in training set (out of vocab words)
 - Use broader vocab + smoothing
 - Make these words as <UNK> and treat them as any other word.
 - Remove these words from test
- Problem with n-gram: zero probability
 - Apply discounting



Laplace Smoothing for classification

If a test case has an attribute that does not occur in every class? Then the probability value for that attribute is 0, so all is 0!

In this case use Laplace smoothing:

$$P(x|y) = \frac{n_{yx}}{n_y} \Longrightarrow \widehat{P}(x|y) = \frac{\gamma + n_{yx}}{\gamma \cdot d + n_y}, \ x \in (x_1, \dots, x_d) \text{ d-dimensional multinomial distribution}$$

 n_y no. of instances from training in class y no. of instances from class y with value x for attribute X no. of distinct features for that attribute (number of dimensions), γ is usually taken as 1.

Every feature must be represented in each class with non-zero probability.



BNB example with Laplace smoothing

	ત્ર	b	د	class	
ei	0	l	Ó	+	
وء	0	l	1	+	
ez	١	0	0	+	
e4	l	ι	0	+	
√,	١	l	1	+	
V2	0	0	Ō	+	
وح	١	ι	0	_	
وړ وړ	l	0	1	_	
e7	ι	0	Ō	_	
CP	0	0	0		
٧	١	l	1	_	
٧٧	Q	O	O	_	

Smoothing added 2 rows in each class, need to add row with 0's, because each feature must have prob>0 in each class.

And this is binomial, so 2 dimensions



BNB example with Laplace smoothing

	ત્ર	b	د	class	
ei	0	l	Ó	+	
وء	0	ι	1	+	
ez	l	0	0	+	
e4	l	ι	0	+	
√,	ı	l	1	+	
٧ ₂	0	0	D	+	
وح	l	ι	0	_	
وړ وړ	l	0	1	_	
e7	ι	0	O	_	
ሮ ያ	0	0	0	_	
٧	١	ι	1	_	
٧٧	Q	O	0	_	

- BNB training:
 - For each class: num docs with word/num all docs
 - Spam: P(a|+), P(b|+), P(c|+) = (3/6, 4/6, 2/6) = (0.5; 0.67; 0.33)
 - Ham: P(a|-), P(b|-), P(c|-) = (4/6, 2/6, 2/6) = (0.67; 0.33; 0.33)

Testing for instance (a b -) = (1,1,0):

Classify (a b -) using MAP:

- P(c|+) = ?
- P(c|-) = ?
- Result:



BNB example with Laplace smoothing

	ત	b c	class	
ei	0	1 0	+	
وء	0	ιį	+	
ez	l	0 0	+	
e4	l	10	+	
√,	١	1 1	+	
٧ ₂	0	0 0	+	
وح	١	10	_	
وړ وړ	l	0 1	_	
e7	ι	00	_	
C8	0	00		
٧	1	l I	_	
٧٧	O	00	_	

- BNB training:
 - For each class: num docs with word/num all docs
 - Spam: P(a|+), P(b|+), P(c|+) = (3/6, 4/6, 2/6) = (0.5; 0.67; 0.33)
 - Ham: P(a|-), P(b|-), P(c|-) = (4/6, 2/6, 2/6) = (0.67; 0.33; 0.33)

Testing for instance x=(a b -) = (1,1,0):

Classify x=(a b -) using MAP:

- P(x|+) = 0.5 *0.67 *(1-0.33) = 0.224;
- P(x|-) = 0.67 *0.33*(1-0.33) = 0.148
- Result: SPAM (same as before)

Classify using ML: same, since P(-) = P(+) = 0.5



Laplace Smoothing for n-grams

Probabilities of some n-grams are 0.

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})} \Longrightarrow \hat{P}(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)+1}{C(w_{n-1})+|V|}$$

C function is count of words, V is the vocabulary

Another version of Laplace: k-smoothing, k is optimised on devset

$$\widehat{P}(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n) + k}{C(w_{n-1}) + k \times |V|}$$

But Laplace smoothing does not work well for n-grams



More on Smoothing

- Problem with n-gram: zero probability
- Back-off smoothing: if n-gram has zero probability, back-off to n-1 gram, until reaching non-zero probability. But this method has a problem: mixing probabilities for lower and higher order n-grams may sum to greater than 1. We need to discount higher order probabilities to save for lower order.
- Undiscounted back-off smoothing can be used for very large language models, called stupid back-off (③)
- Interpolation smoothing: instead of simple back-off, combine n-1, n-2, ..., smoothing by weighted averaging probabilities.
- Kneser-Ney method addresses this problem and combines both.



Kneser-Ney Smoothing elements

- Kneser-Ney is based on P continuation P_{cont}
 - Discount is based on the number of different contexts word w has appeared after other words
 - E.g., AA, AB, AC, AD, AC, AB, CB, DB, BA, CB: B appears in 3 different contexts out of 7 contexts, therefore, $P_{cont} = \frac{3}{7}$
 - Intuition words that have appeared in more contexts in the past are more likely to appear in some new context
 - General formula $P_{cont}(w_i) = \frac{|\{w': 0 < c(w', w_i)\}|}{|\{(w', w''): 0 < c(w', w'')\}|}$
- Discount d, $0 \le d \le 1$
- Constant λ that distributes probability between 1-gram and 2-gram (for 2-grams)

$$\lambda(w_{i-1}) = \frac{d}{\sum_{v} C(w_{i-1} v)} |\{w : C(w_{i-1} w) > 0\}|$$



Advanced: Kneser-Ney Smoothing

- The final probability function using KN smoothing formula is recursive, it calculates probability recursively for sequences until it gets to unigram.
- Formula for bigram is

$$P_{KN}(w_i|w_{i-1}) = \frac{\max(C(w_{i-1}w_i) - d, 0)}{C(w_{i-1})} + \lambda(w_{i-1})P_{KN}(w_i)$$

Where
$$\lambda(w_{i-1}) = \frac{d}{\sum_{v} C(w_{i-1} v)} |\{w: C(w_{i-1} w) > 0\}|$$

 $P_{KN}(w_i) = P_{cont}(w_i)$ for bi-gram

General formula for n-gram in Jurafsky ch 3.7.2



Advanced: KN Smoothing full formula

General formula for n-gram (<u>Kneser–Ney smoothing – Wikipedia</u>)

$$p_{KN}(w_i|w_{i-n+1}^{i-1}) = \frac{\max(c(w_{i-n+1}^{i-1},w_i) - \delta, 0)}{\sum_{w'} c(w_{i-n+1}^{i-1},w')} + \delta \frac{|\{w': 0 < c(w_{i-n+1}^{i-1},w')\}|}{\sum_{w'} c(w_{i-n+1}^{i-1},w')} p_{KN}(w_i|w_{i-n+2}^{i-1})$$

 δ in this formula is our d, which is a discount constant, usually 0.75.

 w_{i-n+1}^{i-1} is a sequence of words in n-gram from the one before last token in n-gram w_i back to the first token in n-gram

For example, for n = 4, we have $w_{i-3}, w_{i-2}, w_{i-1}$ before token w_i

- For bi-gram P_{KN} is reduced to just $P_{KN}(w_i) = P_{cont}(w_i) = \frac{|\{w': 0 < c(w', w_i)\}|}{|\{(w', w''): 0 < c(w', w'')\}|}$
- The full formula is recursive and goes from n-gram to one-gram.



Multinomial Naïve Bayes



Multinomial Naïve Bayes (MNB)

- What information was not used in NB method so far?
- So far we treated occurrence of a word as binary: it is there or it is not. This is called Bernoulli Naïve Bayes method (BNB), because it follows Bernoulli distribution

$$p(x) = P[X = x] = \begin{cases} q = 1 - p & x = 0 \\ p & x = 1 \end{cases}$$

$$f(k;p) = p^k (1-p)^{1-k} \quad ext{for } k \in \{0,1\}$$

k is number of outcomes

- But if we count how many times a word (token) appears in an instance, this
 is called Multinomial Naïve Bayes method (MNB)
- More complex calculations, but generally performs better

MNB example with Laplace smooting

	•	a	b	٢	class	
	e, e _z	0	3	0	+	
	در	0	3	3	+	
	e3	3	0	0	+	
	e4	2	3	0	+	
>	V	l	0	0	+	
	٧ ₂	0	I	0	+	
	∨ 3	0	0	1	+	
	e 5	4	3	0	_	
	eb	4	0	3	_	
	e 7	3	0	0	_	
	68	٥	0	0	_	
	٧,	١	0	ð	_	
	Y 2	0	I	0	~	
	٧3	0	0	l	_	

Smoothing added 3 features, 1 for each class, because this is multinomial 3-dim multinomial distribution



MNB example with Laplace smooting

	•	a	Ь	٢	class	
	ور ور وع	0	3	0	+	
	ور	0	3	3	+	
	e3	3	0	0	+	
	e ₄	2	3	0	+	
G) V	l	0	0	+	
	٧ُو	0	I	0	+	
	∨ 3	0	0	1	+	
	e 5	4	3	0	-	
	eb	4	0	3	_	
	e 7	3	0	0	_	
	68	0	0	O	_	
П	٧,	١	0	Õ	~	П
	Y 2	0	I	0	~	
	v ₃	0	0	1		

MNB training:

- –For each class: word count/all words count
- -Spam: P(a|+), P(b|+), P(c|+) = (6/20, 10/20, 4/20) = (0.3, 0.5, 0.2)
- -Ham: P(a|-), P(b|-), P(c|-) = (12/20, 4/20, 4/20) = (0.6, 0.2, 0.2)



MNB calculations: example

Testing instance: $(a \ a \ a \ b) = (3a \ 1b \ 0c) = (3,1,0)$:

- Multiple occurrences of a word in doc
- We model this as:

There are k different unique words in a bag of words.

Probability of drawing k^{th} word is p_k .

We draw n of those words from the bag with replacement.

What is the probability that we get x_1 of k_1 word AND x_2 of k_2 AND ..., etc.? (sum(x_i) = n). It follows multinomial probability mass function (general formula):

$$f(x_1,\ldots,x_k;n,p_1,\ldots,p_k) = \Pr(X_1 = x_1 ext{ and } \ldots ext{ and } X_k = x_k) \ = egin{cases} rac{n!}{x_1!\cdots x_k!}p_1^{x_1} imes\cdots imes p_k^{x_k}, & ext{ when } \sum_{i=1}^k x_i = n \ 0 & ext{ otherwise,} \end{cases}$$



MNB calculations: example

```
Our instance: (#a=3,#b=1,#c=0):

n=4,

k = 3 \mid (a,b,c)

x1=3, x2=1, x3=0,

likelihood = argmax(P(X|h))
```

Therefore:

$$P(x|+) = ?$$

 $P(x|-) = ?$



MNB calculations: example

Our instance: (#a=3,#b=1,#c=0):

n=4, x1=3, x2=1, x3=0,

Therefore:

 $P(x|+) = 4!/(3!*1!*0!)*0.3^3*0.5^1*0.2^0 = 0.054 \text{ for SPAM}$

 $P(x|-) = 4!/(3!*1!*0!)*0.6^3*0.2^1*0.2^0 = 0.1728 = for HAM$

So by arg max likelihood this is ham

BNB: same instance (a a a b -)=(a b -) = (1,1,0): spam, different than MNB, because of 3 a's.



Sentiment analysis and classfication



Sentiment, aspect, stance

"@united currently on board so not now. Check in was terrible. Staff rude. Expensive luggage thrown around etc."

Identify aspects in this tweet



How to use sentiment lexicons

- Lexicons are kind of dictionaries of words labelled with emotion or sentiment
- E.g. NRC: words labelled with emotions, e.g.

```
abandon fear abandon joy
```

- Score words by summing positive and negative emotion labels
- VADER (Valence Aware Dictionary and sEntiment Reasoner)
 - Lexicon and Tool Set
 - Validated by multiple independent human judges,
 - Suitable for tweets and longer texts.
 - Hard to use on its own, better to use NLTK API or the tools.
 - Additional information and the full VADER dataset and tools are in MyUni Online Learning.



How to classify sentiment

- Use lexicons, e.g. VADER on its own (unsupervised) or combined with a classifier
- Once you have polarity of the word (=sentiment), add a feature to the text, e.g.
 "__NEG___" and use it in training.
- (Jurafsky&Martin, ch. 4.4) prepend the prefix NOT to every word after a token of logical negation (n't, not, no, never) until the next punctuation mark. Thus the phrase "didn't like this movie, but I" becomes "didn't NOT_like NOT_this NOT_movie, but I"
- Can parse sentence into phrases (will discuss in future modules)
- (better to test this method on devset before applying)





Exercise on Bayesian rule: solution

P(cancer) = 0.008, P(not cancer) = 0.992, these are called prior probabilities P(+ | cancer) = 0.98 P(- | cancer) = 0.02 P(+ | not cancer) = 0.03 P(- | not cancer) = 0.97 Does the patient have cancer or not? To answer this, calculate P1 = P(not cancer | +) and P2 = P(cancer | +), and then calculate likelihood ratio P1/P2.

P1 = P(nc | +) = P(+|nc)*P(nc) =
$$0.03*0.992$$

P2 = P(c|+) = P(+|c)*P(c) = $0.98*0.008$

Answer: P1/P2 = 3.8

