

Naive Bayes

Lecture 13

Centre for Data Science
Institute of Technical Education and Research
Siksha 'O' Anusandhan (Deemed to be University)
Bhubaneswar, Odisha, 751030



Overview

- 1 Dumb Spam filter
- 2 Sophisticated Spam filter
- 3 Implementation
- 4 Testing our model
- 5 Using our model
- 6 How could we get better performance
- 7 References

Dumb Spam filter

- Imagine a universe that consists of receiving a message chosen randomly from all possible messages.
- Let S be the event *the message is spam* and B be the event *the message contains the word bitcoin*.
- Bayes's theorem tells us that the probability that the message is spam conditional on containing the word bitcoin is:

$$P(S|B) = [P(B|S)P(S)]/[P(B|S)P(S) + P(B|\neg S)P(\neg S)]$$

Dumb Spam filter

- From below given equation,

$$P(S|B) = [P(B|S)P(S)]/[P(B|S)P(S) + P(B|\neg S)P(\neg S)]$$

- The numerator is the probability that a message is spam and contains bitcoin.
- The denominator is just the probability that a message contains bitcoin.
- We can think of this calculation as simply representing the proportion of bitcoin messages that are spam.

Dumb Spam filter

- If we have a large collection of messages we know are spam, and a large collection of messages we know are not spam, then we can easily estimate $P(B|S)$ and $P(B|\neg S)$. If we further assume that any message is equally likely to be spam or not spam (so that $P(S) = P(\neg S) = 0.5$), then:

$$P(S|B) = [P(B|S)]/[P(B|S) + P(B|\neg S)]$$

- Suppose, if 50% of spam messages have the word bitcoin, but only 1% of nonspam messages do, then the probability that any given bitcoin-containing email is spam is:

$$0.5/(0.5 + 0.01) = 98\%$$

Sophisticated Spam filter

- Imagine now that we have a vocabulary of many words, w_1, \dots, w_n .
- To move this into the realm of probability theory, we'll write X_i for the event *a message contains the word w_i* .
- Imagine that we've come up with an estimate $P(X_i|S)$ for the probability that a spam message contains the i th word
- A similar estimate $P(X_i|\neg S)$ for the probability that a nonspam message contains the i th word.
- The key to Naive Bayes is making the assumption that the presences (or absences) of each word are independent of one another, conditional on a message being spam or not.

Sophisticated Spam filter

- This assumption means that knowing whether a certain spam message contains the word *bitcoin* gives you no information about whether that same message contains the word *rolex*.
- In maths terms, this is equivalent to

$$P(X_1 = x_1, \dots, X_n = x_n | S) = P(X_1 = x_1 | S) \times \dots \times P(X_n = x_n | S)$$

Sophisticated Spam filter

- Imagine that our vocabulary consists only of the words *bitcoin* and *rolex*, and that half of all spam messages are for "earn bitcoin" and that the other half are for "authentic rolex".
- In this case, the Naive Bayes estimate that a spam message contains both *bitcoin* and *rolex* is:

$$P(x_1 = 1, X_2 = 1|S) = P(x_1 = 1|S)P(x_2 = 1|S) = 0.5 \times 0.5 = 0.25$$

- Since we've assumed away the knowledge that *bitcoin* and *rolex* actually never occur together.
- Despite the unrealisticness of this assumption, this model often performs well and has historically been used in actual spam filters.

Sophisticated Spam filter

- The Naive Bayes assumption allows us to compute each of the probabilities on the right simply by multiplying together the individual probability estimates for each vocabulary word.
- In practice, you usually want to avoid multiplying lots of probabilities together, to prevent a problem called *underflow*, in which computers don't deal well with floating-point numbers that are too close to 0.

Sophisticated Spam filter

- Imagine that in our training set the vocabulary word *data* only occurs in nonspam messages.
- Then we'd estimate $P(\text{data} \mid S) = 0$. The result is that our Naive Bayes classifier would always assign spam probability 0 to any message containing the word *data*, even a message like "*data on free bitcoin and authentic rolex watches*". To avoid this problem, we usually use some kind of smoothing.
- In particular, we'll choose a *pseudocount-k*- and estimate the probability of seeing the i th word in a spam message as:
$$P(X_i \mid S) = P(k + \text{number of spams containing } w_i) / (2k + \text{number of spams})$$

Implementation

- We have all the pieces we need to build our classifier.
- First, let's create a simple function to tokenize messages into distinct words.
- We'll first convert each message to lowercase, then use *re.findall* to extract "words" consisting of letters, numbers, and apostrophes.
- we'll use set to get just the distinct words.

Implementation

```
1 from typing import Set
2 import re
3 def tokenize(text: str) -> Set[str]:
4     text = text.lower()
5     all_words = re.findall("[a-z0-9]+", text)
6     return set(all_words)
7 assert tokenize("Data Science is science") == {"data", "science", "is"}
```

- We'll also define a type for our training data:

```
1 from typing import NamedTuple
2 class Message(NamedTuple):
3     text: str
4     is_spam: bool
```

Implementation

- As our classifier needs to keep track of tokens, counts, and labels from the training data, we'll make it a class.
- we refer to nonspam emails as *ham* emails.

```
1 from typing import List, Tuple, Dict, Iterable
2 import math
3 from collections import defaultdict
4
5 class NaiveBayesClassifier:
6     def __init__(self, k: float = 0.5) -> None:
7         self.k = k # smoothing factor
8
9         self.tokens: Set[str] = set()
10        self.token_spam_counts: Dict[str, int] =
defaultdict(int)
11        self.token_ham_counts: Dict[str, int] =
defaultdict(int)
12        self.spam_messages = self.ham_messages = 0
```

Implementation

- we'll give it a method to train it on a bunch of messages. First, we increment the *spam_messages* and *ham_messages* counts.
- Then we tokenize each message text, and for each token we increment the *token_spam_counts* or *token_ham_counts* based on the message type.

Implementation

```
1 #This is method of class NaiveBayesClassifier
2 def train(self, messages: Iterable[Message]) -> None:
3     for message in messages:
4         # Increment message counts
5         if message.is_spam:
6             self.spam_messages += 1
7         else:
8             self.ham_messages += 1
9
10        # Increment word counts
11        for token in tokenize(message.text):
12            self.tokens.add(token)
13            if message.is_spam:
14                self.token_spam_counts[token] += 1
15            else:
16                self.token_ham_counts[token] += 1
```

Implementation

- We'll want to predict $P(\text{spam} \text{ — token})$.
- As we saw earlier, to apply Bayes's theorem we need to know $P(\text{token} \text{ — spam})$ and $P(\text{token} \text{ — ham})$ for each token in the vocabulary.
- So we'll create a "private" helper function to compute it.

Implementation

```
1 #This is also a method of class NaiveBayesClassifier
2 def _probabilities(self, token: str) -> Tuple[float,
3 float]:
4     """returns P(token | spam) and P(token | not spam)"""
5     spam = self.token_spam_counts[token]
6     ham = self.token_ham_counts[token]
7
8     p_token_spam = (spam + self.k) / (self.spam_messages
9 + 2 * self.k)
10    p_token_ham = (ham + self.k) / (self.ham_messages + 2
    * self.k)
11
12    return p_token_spam, p_token_ham
```

Implementation

- we're ready to write our predict method.

```
1 #This is also a method of class NaiveBayesClassifier
2 def predict(self, text: str) -> float:
3     text_tokens = tokenize(text)
4     log_prob_if_spam = log_prob_if_ham = 0.0
5     for token in self.tokens:
6         prob_if_spam, prob_if_ham = self.
          _probabilities(token)
7         if token in text_tokens:
8             log_prob_if_spam += math.log(prob_if_spam)
9             log_prob_if_ham += math.log(prob_if_ham)
10        else:
11            log_prob_if_spam += math.log(1.0 -
          prob_if_spam)
12            log_prob_if_ham += math.log(1.0 -
          prob_if_ham)
13        prob_if_spam = math.exp(log_prob_if_spam)
14        prob_if_ham = math.exp(log_prob_if_ham)
15    return prob_if_spam / (prob_if_spam + prob_if_ham)
```

Testing our model

- Now we will check our model is working or not properly.

```
1 messages = [Message("spam rules", is_spam=True),
2               Message("ham rules", is_spam=False),
3               Message("hello ham", is_spam=False)]
4
5 model = NaiveBayesClassifier(k=0.5)
6 model.train(messages)
7
8 assert model.tokens == {"spam", "ham", "rules", "hello"}
9 assert model.spam_messages == 1
10 assert model.ham_messages == 2
11 assert model.token_spam_counts == {"spam": 1, "rules": 1}
12 assert model.token_ham_counts == {"ham": 2, "rules": 1, "
    hello": 1}
```

Testing our model

- We'll also (laboriously) go through our Naive Bayes logic by hand, and make sure that we get the same result.

```
1 text = "hello spam"
2 probs_if_spam = [
3     (1 + 0.5) / (1 + 2 * 0.5),      # "spam"    (present)
4     1 - (0 + 0.5) / (1 + 2 * 0.5),  # "ham"     (not present)
5     1 - (1 + 0.5) / (1 + 2 * 0.5),  # "rules"  (not present)
6     (0 + 0.5) / (1 + 2 * 0.5)]      # "hello"  (present)
7 probs_if_ham = [
8     (0 + 0.5) / (2 + 2 * 0.5),      # "spam"    (present)
9     1 - (2 + 0.5) / (2 + 2 * 0.5),  # "ham"     (not present)
10    1 - (1 + 0.5) / (2 + 2 * 0.5),  # "rules"  (not present)
11    (1 + 0.5) / (2 + 2 * 0.5)]      # "hello"  (present)
12 p_if_spam=math.exp(sum(math.log(p) for p in probs_if_spam))
13 p_if_ham=math.exp(sum(math.log(p) for p in probs_if_ham))
14 # Should be about 0.83
15 assert model.predict(text)==p_if_spam/(p_if_spam+p_if_ham)
```

Using our model

- A popular dataset is the **SpamAssassin public corpus**. We'll look at the files prefixed with *20021010*.
- After downloading the data you should have three folders: *spam*, *easy_ham*, and *hard_ham*.
- Each folder contains many emails, each contained in a single file.
- To keep things really simple, we'll just look at the subject lines of each email.

Using our model

- When we look through the files, they all seem to start with "Subject:". So we'll look for that

```
1 import glob, re
2 # modify the path to wherever you've put the files
3 path = 'spam_data/**/*.txt'
4 data: List[Message] = []
5 # glob.glob returns every filename that matches the
  wildcarded path
6 for filename in glob.glob(path):
7     is_spam = "ham" not in filename
8     # There are some garbage characters in the emails;
  the errors='ignore'
9     # skips them instead of raising an exception.
10    with open(filename, errors='ignore') as email_file:
11        for line in email_file:
12            if line.startswith("Subject:"):
13                subject = line.lstrip("Subject: ")
14                data.append(Message(subject, is_spam))
15                break # done with this file
```

Using our model

- Now we can split the data into training data and test data, and then we're ready to build a classifier.

```
1 import random
2 from scratch.machine_learning import split_data
3 random.seed(0) # just so you get the same answers as me
4 train_messages, test_messages = split_data(data, 0.75)
5 model = NaiveBayesClassifier()
6 model.train(train_messages)
```

Using our model

- Let's generate some predictions and check how our model does

```
1 from collections import Counter
2 predictions = [(message, model.predict(message.text)) for
3               message in test_messages]
4 # Assume that spam_probability > 0.5 corresponds to spam
5   prediction
6 # and count the combinations of (actual is_spam ,
7   predicted is_spam)
8 confusion_matrix = Counter((message.is_spam ,
9                             spam_probability > 0.5) for message, spam_probability
10                            in predictions)
11 print(confusion_matrix)
```


Using our model

- This gives 84 true positives (spam classified as "spam")
- 25 false positives (ham classified as "spam")
- 703 true negatives (ham classified as "ham")
- 44 false negatives (spam classified as "ham").
- This means our precision is $84 / (84 + 25) = 77\%$, and our recall is $84 / (84 + 44) = 65\%$, which are not bad numbers for such a simple model.

How could we get better performance

- Look at the message content, not just the subject line. You'll have to be careful how you deal with the message headers.
- Our classifier takes into account every word that appears in the training set, even words that appear only once. Modify the classifier to accept an optional *min_count* threshold and ignore tokens that don't appear at least that many times.
- The tokenizer has no notion of similar words (e.g., cheap and cheapest). Modify the classifier to take an optional stemmer function that converts words to equivalence classes of words.

[1] Joel Grus. Data Science from Scratch. First Principles with Python. Second Edition. O'REILLY, May 2019.

Thank You