

Homework 6

Due November 5 at 11 pm

1. (Spam detector) In order to build a spam classifier, we gather the following data. Each row is an email. The first column indicates whether it is spam or not. The remaining columns indicate whether it contains (✓) or not (✗) the word on top.

	Miracle	Alternative	Medicine	Basketball
Spam	✓	✓	✓	✓
Not spam	✗	✗	✓	✓
Spam	✓	✗	✗	✗
Not spam	✗	✓	✗	✗
Spam	✓	✗	✗	✗
Not spam	✗	✓	✗	✓
Spam	✓	✗	✗	✗
Spam	✗	✓	✓	✗
Not spam	✓	✓	✗	✓
Not spam	✗	✗	✓	✓

We use a Bernoulli random variable \tilde{y} to model whether the email is spam ($\tilde{y} = 1$) or not ($\tilde{y} = 0$), and a four-dimensional random vector \tilde{x} to indicate whether the i th word is ($\tilde{x}[i] = 1$) or not ($\tilde{x}[i] = 0$) in the email for $i \in \{1, 2, 3, 4\}$.

- (a) Your friend recommends that you estimate the conditional pmf of \tilde{y} given \tilde{x} and then maximize it to produce your estimate. What is the problem with this approach?

The problem with this approach is that if we are to compute conditional pmf of \tilde{y} given \tilde{x} , we will have to know the joint distribution of $\tilde{x}[1], \tilde{x}[2], \tilde{x}[3], \tilde{x}[4]$. If we are using the empirical pmf to estimate the distribution, we will have $2^4 = 16$ parameters to estimate, which seems ok in this problem but in reality the number of suspicious words in spam may be way lot more, where the parameter to be estimated will explode. Thus the curse of dimensionality will occur.

- (b) Apply naive Bayes to classify an email that reads *I hurt my foot playing basketball. Can you get me some medicine?*

We want to find the following:

$$\begin{aligned}
 MAP(\vec{x}) &= \arg \max_{y \in \{Spam, not \ spam\}} p_{\tilde{y}|\tilde{x}}(y|\vec{x}) \\
 &= \arg \max_{y \in \{Spam, not \ spam\}} \frac{\prod_{i=1}^4 p_{\tilde{x}[i]|\tilde{y}}(x[i]|y) p_{\tilde{y}}(y)}{\sum_{y \in \{Spam, not \ spam\}} \prod_{i=1}^4 p_{\tilde{x}[i]|\tilde{y}}(x[i]|y) p_{\tilde{y}}(y)} \\
 &= \arg \max_{y \in \{Spam, not \ spam\}} \frac{\prod_{i=1}^2 p_{\tilde{x}[i]|\tilde{y}}(0|y) p_{\tilde{y}}(y) \prod_{i=3}^4 p_{\tilde{x}[i]|\tilde{y}}(1|y) p_{\tilde{y}}(y)}{\sum_{y \in \{Spam, not \ spam\}} \prod_{i=1}^2 p_{\tilde{x}[i]|\tilde{y}}(0|y) p_{\tilde{y}}(y) \prod_{i=3}^4 p_{\tilde{x}[i]|\tilde{y}}(1|y) p_{\tilde{y}}(y)} \\
 &= \arg \max_{y \in \{Spam, not \ spam\}} \left\{ \frac{\frac{1}{2} \times \frac{1}{5} \times \frac{3}{5} \times \frac{2}{5} \times \frac{4}{5}}{\frac{1}{2} \times \frac{1}{5} \times \frac{3}{5} \times \frac{2}{5} \times \frac{4}{5} + \frac{1}{2} \times \frac{4}{5} \times \frac{2}{5} \times \frac{2}{5} \times \frac{4}{5}} \right. \\
 &\quad \left. , \frac{\frac{1}{2} \times \frac{4}{5} \times \frac{2}{5} \times \frac{2}{5} \times \frac{4}{5}}{\frac{1}{2} \times \frac{1}{5} \times \frac{3}{5} \times \frac{2}{5} \times \frac{4}{5} + \frac{1}{2} \times \frac{4}{5} \times \frac{2}{5} \times \frac{2}{5} \times \frac{4}{5}} \right\} \\
 &= not \ spam
 \end{aligned}$$

- (c) Apply naive Bayes to classify an email that reads *This alternative medicine is amazing, send us all your money!* Explain what shortcoming of the naive Bayes classifier is illustrated by this example.

We want to find the following:

$$\begin{aligned}
MAP(\vec{x}) &= \arg \max_{y \in \{Spam, not \ spam\}} p_{\tilde{y}|\tilde{x}}(y|x) \\
&= \arg \max_{y \in \{Spam, not \ spam\}} \frac{\prod_{i=1}^4 p_{\tilde{x}[i]|\tilde{y}}(x[i]|y) p_{\tilde{y}}(y)}{\sum_{y \in \{Spam, not \ spam\}} \prod_{i=1}^4 p_{\tilde{x}[i]|\tilde{y}}(x[i]|y) p_{\tilde{y}}(y)} \\
&= \arg \max_{y \in \{Spam, not \ spam\}} \frac{\prod_{i=1,4} p_{\tilde{x}[i]|\tilde{y}}(0|y) p_{\tilde{y}}(y) \prod_{i=2,3} p_{\tilde{x}[i]|\tilde{y}}(1|y) p_{\tilde{y}}(y)}{\sum_{y \in \{Spam, not \ spam\}} \prod_{i=1,4} p_{\tilde{x}[i]|\tilde{y}}(0|y) p_{\tilde{y}}(y) \prod_{i=2,3} p_{\tilde{x}[i]|\tilde{y}}(1|y) p_{\tilde{y}}(y)} \\
&= \arg \max_{y \in \{Spam, not \ spam\}} \left\{ \frac{\frac{1}{2} \times \frac{1}{5} \times \frac{4}{5} \times \frac{2}{5} \times \frac{2}{5}}{\frac{1}{2} \times \frac{1}{5} \times \frac{4}{5} \times \frac{2}{5} \times \frac{2}{5} + \frac{1}{2} \times \frac{4}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{3}{5}}, \frac{\frac{1}{2} \times \frac{4}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{3}{5}}{\frac{1}{2} \times \frac{1}{5} \times \frac{4}{5} \times \frac{2}{5} \times \frac{2}{5} + \frac{1}{2} \times \frac{4}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{3}{5}} \right\} \\
&= not \ spam
\end{aligned}$$

We notice that this sentence is clearly spam since it contains provocative and monetary keyword but our naive bayes classifies it as "not spam". The shortcoming turns out to arise from the assumptions that we made in the model, which is conditional independence given the hypothesis. Since sometimes the words in a sentence can only make practical sense when they come in pairs or trigrams, the assumption basically interpret each word separately. Moreover, the naive bayes relies heavily on the quality of dataset, where since we only have four indicators, which cannot fully encompass the suspicious keywords that have higher probability in spam (for example, "send", "money").

2. (The Markov property) Let $\tilde{a}_1, \tilde{a}_2, \dots, \tilde{a}_n$ be a Markov chain, where for any $1 < i < n$, \tilde{a}_{i+1} is conditionally independent of $\tilde{a}_1, \dots, \tilde{a}_{i-1}$ given \tilde{a}_i , i.e.

$$p_{\tilde{a}_{i+1} | \tilde{a}_1, \dots, \tilde{a}_i} (a_{i+1} | a_1, a_2, \dots, a_i) = p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i), \quad (1)$$

for any values of a_1, a_2, \dots, a_n . Show that this implies that the future is conditionally independent from the past given the present:

$$p_{\tilde{a}_{i+1}, \tilde{a}_{i-1} | \tilde{a}_i} (a_{i+1}, a_{i-1} | a_i) = p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i) p_{\tilde{a}_{i-1} | \tilde{a}_i} (a_{i-1} | a_i), \quad (2)$$

for any $2 \leq i \leq n-1$ and any values of a_{i-1}, a_i, a_{i+1} . (Hint: First show that $p_{\tilde{a}_{i+1} | \tilde{a}_{i-1}, \tilde{a}_i} (a_{i+1} | a_{i-1}, a_i) = p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i)$ for a_{i-1}, a_i, a_{i+1} .)

We first show the probability in the hint and try to generalize it:

Proof. We know from the bayes theorem that:

$$\begin{aligned} p_{\tilde{a}_{i+1} | \tilde{a}_{i-1}, \tilde{a}_i} (a_{i+1} | a_{i-1}, a_i) &= \frac{p_{\tilde{a}_{i+1}, \tilde{a}_{i-1}, \tilde{a}_i} (a_{i+1}, a_{i-1}, a_i)}{p_{\tilde{a}_{i-1}, \tilde{a}_i} (a_{i-1}, a_i)} && \text{(Bayes theorem)} \\ &= \frac{\sum_{a_1, a_2, \dots, a_{i-2}} p_{\tilde{a}_1} (a_1) \prod_{n=1}^{i-1} p_{\tilde{a}_{n+1} | \tilde{a}_n} (a_{n+1} | a_n) \cdot p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i)}{\sum_{a_1, a_2, \dots, a_{i-2}} p_{\tilde{a}_1} (a_1) \prod_{n=1}^{i-1} p_{\tilde{a}_{n+1} | \tilde{a}_n} (a_{n+1} | a_n)} \\ & && \text{(Chain Rule)} \\ &= p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i) \end{aligned}$$

Generally speaking, we could have $p_{\tilde{a}_{i+1} | \tilde{a}_{i-1}, \tilde{a}_i} (a_{i+1} | a_{i-k}, \dots, a_{i-1}, a_i) = p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i)$ for $a_{i-k}, \dots, a_{i-1}, a_i, a_{i+1}$ (where $i-k \geq 1$), which means that if we know the present, then the past information doesn't matter. \square

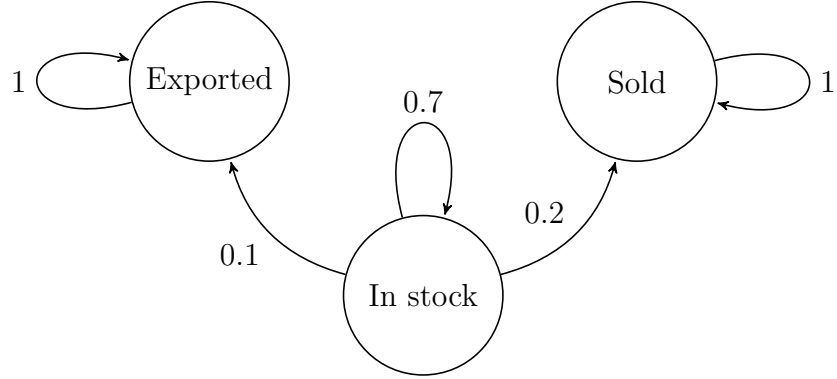
Once we have this, we can proceed with our proof as follows:

Proof.

$$\begin{aligned} p_{\tilde{a}_{i+1}, \tilde{a}_{i-1} | \tilde{a}_i} (a_{i+1}, a_{i-1} | a_i) &= \frac{p_{\tilde{a}_{i+1} | \tilde{a}_i, \tilde{a}_{i-1}} (a_{i+1} | a_i, a_{i-1}) p_{\tilde{a}_i, \tilde{a}_{i-1}} (a_{i-1} | a_i) p_{\tilde{a}_i} (a_i)}{p_{\tilde{a}_i} (a_i)} && \text{(Bayes Rule, Chain Rule)} \\ &= \frac{p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i) p_{\tilde{a}_i, \tilde{a}_{i-1}} (a_{i-1} | a_i) p_{\tilde{a}_i} (a_i)}{p_{\tilde{a}_i} (a_i)} && \text{(Using former lemma)} \\ &= p_{\tilde{a}_{i+1} | \tilde{a}_i} (a_{i+1} | a_i) p_{\tilde{a}_i, \tilde{a}_{i-1}} (a_{i-1} | a_i) && \text{(Arithmetic Operations)} \end{aligned}$$

\square

3. (Mobile phones) A company that makes mobile phones wants to model the sales of a new model they have just released. At the moment 90% of the phones are in stock, 10% have been sold locally and none have been exported. Based on past data, the company determines that each day a phone is sold with probability 0.2 and exported with probability 0.1. We define the following time-homogeneous Markov chain with three states to model this:



- (a) What is the limit of the state vector π_i as $i \rightarrow \infty$?

We denote the state vector at the beginning to be $\pi_1 = \begin{bmatrix} 0.9 \\ 0.1 \\ 0 \end{bmatrix}$ and the transition

matrix to be $T = \begin{bmatrix} 0.7 & 0 & 0 \\ 0.2 & 1 & 0 \\ 0.1 & 0 & 1 \end{bmatrix}$. We do a spectrum decomposition and could get

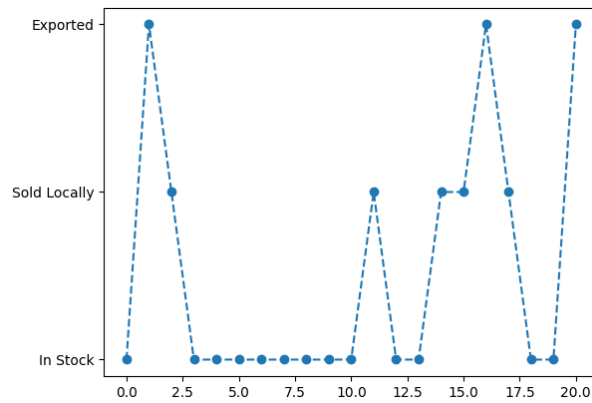
$$T = Q\Lambda Q^{-1} = \begin{bmatrix} 3/\sqrt{14} & 0 & 0 \\ -2/\sqrt{14} & 1 & 0 \\ -1/\sqrt{14} & 0 & 1 \end{bmatrix} \begin{bmatrix} 0.7 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \sqrt{14}/3 & 0 & 0 \\ 2/3 & 1 & 0 \\ 1/3 & 0 & 1 \end{bmatrix}.$$

We then apply the following:

$$\begin{aligned} \lim_{i \rightarrow \infty} \pi_i &= T^{i-1} \pi_1 = \begin{bmatrix} 3/\sqrt{14} & 0 & 0 \\ -2/\sqrt{14} & 1 & 0 \\ -1/\sqrt{14} & 0 & 1 \end{bmatrix} \lim_{i \rightarrow \infty} \begin{bmatrix} 0.7^{i-1} & 0 & 0 \\ 0 & 1^{i-1} & 0 \\ 0 & 0 & 1^{i-1} \end{bmatrix} \begin{bmatrix} \sqrt{14}/3 & 0 & 0 \\ 2/3 & 1 & 0 \\ 1/3 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0.9 \\ 0.1 \\ 0 \end{bmatrix} \\ &= \left(\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 2/3 & 1 & 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 1/3 & 0 & 1 \end{bmatrix} \right) \begin{bmatrix} 0.9 \\ 0.1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ 0.7 \\ 0.3 \end{bmatrix} \end{aligned}$$

- (b) Simulate the Markov chain and plot the evolution of the state vector.

Attached is the link to the code(click it): [Jupyter notebook pdf link](#)



```

1  import numpy as np
2  import matplotlib.pyplot as plt
3  %matplotlib inline
4
5  # Problem 3(a): Limit of State Vector
6  T = np.array([[0.7,0,0],[0.2,1,0],[0.1,0,1]])
7  eigenvalues, eigenvectors = np.linalg.eig(T)
8  Q = eigenvectors
9  Q.inv = np.linalg.inv(eigenvectors)
10
11 # Problem 3(b): Simulate the Markov Chain
12 x0 = np.array([0.9,0.1,0])
13 T = np.array([[0.7,0,0],[0.2,1,0],[0.1,0,1]])
14 state_name = ["In-Stock", "Sold-Locally", "Exported"]
15
16
17 def one_step(state, transition_matrix):
18     return np.random.choice(3, 1, p = transition_matrix[:, state].flatten())[0]
19
20
21 def simulate_discrete_markov_chain(start, transition_matrix, step):
22     trace = [start]
23     for i in range(step):
24         trace.append(one_step(start, transition_matrix))
25
26     return trace
27
28 def plot_trace(state_name, start_prob, transition_matrix, step):
29     start_state = np.random.choice(len(state_name), 1, p = start_prob)[0]
30     trace = simulate_discrete_markov_chain(start_state, transition_matrix, step)
31     print(trace)
32     plt.plot(np.arange(0, step + 1), trace, "o—")
33     plt.yticks([0,1,2], state_name)
34
35     plt.show()
36
37 plot_trace(state_name, x0, T, 20)

```

4. (Smoothing) In the application of naive Bayes, instead of computing the conditional pmf $p_{\tilde{x}[i]|\tilde{y}}(x[i]|y)$ empirically, we usually compute the conditional pmf as follows:

$$p_{\tilde{x}[i]|\tilde{y}}(\tilde{x}[i]|y) = \frac{\text{number of samples with label } y \text{ and having feature } x[i] + m}{\text{number of samples with label } y + md},$$

where m is a smoothing constant, and d is the number of features.

- (a) What is the problem with naive Bayes when there is no smoothing ($m = 0$) and how does smoothing alleviate this?

When there is no smoothing, if the number of samples with label y could be zero, which results in zero division error when writing code.

- (b) Complete the notebook *spam.ipynb*. Which m results in the highest accuracy on the test set?

It seems that when there is no smoothing, the prediction accuracy is maximized, and the bigger the smoothing constant is, the lower the accuracy is. The maximum accuracy occurs when $m = 0$, where zero division error happens. The parameter table is listed below for your reference:

m	accuracy
0	0.932
0.01	0.927
0.1	0.924
1	0.922
10	0.905
100	0.88
1000	0.845

Attached is the link to the code(click it): [Jupyter notebook pdf link](#)

- (c) In the code, we compare the joint pmf $p_{\tilde{y},\tilde{x}[1],\dots,\tilde{x}[d]}(y,\tilde{x}[1],\dots,\tilde{x}[d])$ for $y \in \{0,1\}$. Why is this equivalent to Equation (4.146) in the notes? Why do we need to apply the logarithm in practice?

Since ultimately we are comparing $p_{\tilde{y}|\tilde{x}[1],\dots,\tilde{x}[d]}(y,\tilde{x}[1],\dots,\tilde{x}[d])$ for $y \in \{0,1\}$ and that $p_{\tilde{y}|\tilde{x}[1],\dots,\tilde{x}[d]}(y,\tilde{x}[1],\dots,\tilde{x}[d]) = \frac{p_{\tilde{y},\tilde{x}[1],\dots,\tilde{x}[d]}(y,\tilde{x}[1],\dots,\tilde{x}[d])}{p_{\tilde{x}[1],\dots,\tilde{x}[d]}(\tilde{x}[1],\dots,\tilde{x}[d])}$ by bayes formula(conditional probability formula). Since the denominator is the same across $y = 0, 1$, we have that it's enough just to compare $p_{\tilde{y}|\tilde{x}[1],\dots,\tilde{x}[d]}(y,\tilde{x}[1],\dots,\tilde{x}[d])$.

Since floating multiplication takes a huge amount of time, typically 5 clock cycles per elements whereas the floating addition only takes around 3 clock cycles per element. In other words, it is faster. Moreover, logarithm functions will preserve monotonicity of functions, which means $\max\{a,b\} = \max\{\log(a),\log(b)\}$, so applying it won't change the MAP result.

```
1 import numpy as np
2 from collections import Counter
3
4 def data_loader():
```

```

5     # load data
6     with open('train', 'r') as file:
7         train_data = file.read().split('\n')[:-1]
8     with open('test', 'r') as file:
9         test_data = file.read().split('\n')[:-1]
10    return train_data, test_data
11
12    def parser(datum):
13        # extract labels and words
14        email_addr, label, words = datum.split('-', 2)
15        words = words.split()
16        # transform words into dictionary
17        word_dict = dict(zip([words[i] for i in
18            range(0, len(words), 2)],
19            [int(words[i+1]) for i in range(0, len(words), 2)]))
20        # transform label into 0, 1
21        if label == 'ham':
22            label = 0
23        elif label == 'spam':
24            label = 1
25        else:
26            raise ValueError
27        return label, word_dict
28
29    def data_preprocessing(train_data, test_data):
30        y_train = np.zeros(len(train_data))
31        y_test = np.zeros(len(test_data))
32        x_train = []
33        x_test = []
34        for i, datum in enumerate(train_data):
35            label, word_dict = parser(datum)
36            y_train[i] = label
37            x_train.append(word_dict)
38        for i, datum in enumerate(test_data):
39            label, word_dict = parser(datum)
40            y_test[i] = label
41            x_test.append(word_dict)
42        return x_train, y_train, x_test, y_test
43
44    def compute_empirical_pmf_y(y_train):
45        # compute distribution P(y=1), P(y=0)
46        # TODO
47        return np.sum(y_train == 1) / len(y_train), np.sum(y_train == 0) / len(y_train)
48
49    def m_estimation_conditional_probability(x_train_frts, y_train, num_vocab, m):
50        # compute P(x_j|y=1) and P(x_j|y=0) for j = 1, ..., d
51        # TODO
52        p_on_spam = np.zeros((2, x_train_frts.shape[1]))
53        p_on_ham = np.zeros((2, x_train_frts.shape[1]))
54
55        temp_df = pd.DataFrame(np.hstack([x_train_frts, y_train.reshape(-1,1)]))
56        temp_df.columns = temp_df.columns = [*list(temp_df.columns[:-1], "label")]
57
58        for j in range(num_vocab):
59            column_data_0 = (temp_df[temp_df["label"]==0])[j]
60            column_data_1 = (temp_df[temp_df["label"]==1])[j]
61
62            # P(x_j = 1|y=0) by definition
63            p_on_ham[1, j] = (sum(column_data_0 > 0) + m) / (len(column_data_0) + m * num_vocab)
64            # P(x_j = 0|y=0) by the fact that conditional probability is a valid pmf
65            p_on_ham[0, j] = 1 - p_on_ham[1, j]
66
67
68
69            # P(x_j = 1|y=1)
70            p_on_spam[1, j] = (sum(column_data_1 > 0) + m) / (len(column_data_1) + m * num_vocab)
71            # P(x_j = 0|y=1)
72            p_on_spam[0, j] = 1 - p_on_spam[1, j]
73
74        return p_on_spam, p_on_ham
75
76
77    def log_estimated_probability(p_spam, p_ham, p_on_spam_m, p_on_ham_m, x_frts):
78        # compute log(P(y=1, x_1, x_2, ..., x_n)) and log(P(y=0, x_1, x_2, ..., x_n))
79        # hint: log(P(y, x_1, x_2, ..., x_n)) = log(P(y=1)) + log(P(x_1 | y)) + ... + log(P(x_d | y))
80        # TODO
81
82        log_p_spam, log_p_ham = np.zeros((len(x_frts),)), np.zeros((len(x_frts),))
83
84        for i in range(len(x_frts)):
85            log_p_spam[i] = np.sum([np.log(p_spam)
86                + [np.log(p_on_spam_m[1 if x_frts[i, j] > 0 else 0, j]) for j in range(x_frts.shape[1])]])
87            log_p_ham[i] = np.sum([np.log(p_ham)
88                + [np.log(p_on_ham_m[1 if x_frts[i, j] > 0 else 0, j]) for j in range(x_frts.shape[1])]])
89
90        return log_p_spam, log_p_ham
91
92    def accuracy(y_true, y_pred):
93        # calculate accuracy
94        # TODO
95
96        return np.sum(y_true == y_pred) / len(y_true)

```

```

97
98 from sklearn.feature_extraction import DictVectorizer
99 # load data
100 train_data, test_data = data_loader()
101
102 # extract labels to 0,1 and features to dictionary
103 x_train, y_train, x_test, y_test = data_preprocessing(train_data, test_data)
104
105 # transform word dicts to feature vectors
106 vectorizer = DictVectorizer(sparse=False)
107 x_train_frt = vectorizer.fit_transform(x_train)
108 x_test_frt = vectorizer.transform(x_test)
109 print(x_train_frt.shape, x_test_frt.shape)
110
111 def pipeline(x_train_frt, y_train, x_test_frt, y_test, m):
112     p_spam, p_ham = compute_empirical_pmf(y_train)
113     p_on_spam_m, p_on_ham_m = m_estimation_conditional_probability(x_train_frt, y_train, x_train_frt.shape[1],
114                                                                     p_spam, p_ham, x_test_frt)
115     log_p_spam, log_p_ham = log_estimated_probability(p_spam, p_ham, p_on_spam_m, p_on_ham_m, x_test_frt)
116     test_pred = (log_p_spam > log_p_ham)
117     print(str(m) + ":" + str(accuracy(y_test, test_pred)))
118
119 m_grid = [0, 0.01, 0.1, 1, 10, 100, 1000]
120 for m in m_grid:
121     pipeline(x_train_frt, y_train, x_test_frt, y_test, m)

```