

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ  
УНИВЕРСИТЕТ ИМЕНИ Н. Э. БАУМАНА  
Факультет информатики и систем управления  
Кафедра теоретической информатики и компьютерных технологий

Лабораторная работа №12  
по курсу «Автоматическая обработка текстов»  
«Методы автоматического разрешения  
лексической многозначности»

Выполнил:  
студент группы ИУ9-11М  
Беляев А. В.

Проверила:  
Лукашевич Н. В.

Москва 2018

# 1 Цель работы

Изучить способы разрешения лексической многозначности – выбора значения слова из набора описанных значений. Предлагается использовать подход, основанный на машинном обучении – метод «Наивного байесовского классификатора».

Необходимо выбрать многозначное слово, собрать контекст из 10 предложений, содержащих это слово в определенном значении и использовать его в качестве обучающей выборки.

В качестве признаков следует использовать сами слова (лексические признаки) и части речи близлежащих к целевому слову слов («частеречные» признаки).

# 2 Ход работы

На Рис. 1 и 2 приведены конкретные алгоритмы обучения и применения байесовского классификатора. Необходимо согласовать терминологию с предметной областью задачи в приведенных алгоритмах:

- Vocabulary – словарь (уникальных) признаков
- Doc – предложение
- Class – значение слова
- Term – признак

Под признаками в данном алгоритме понимаются предикаты.

```
TRAINBERNOULLNB(C, D)
1  V ← EXTRACTVOCABULARY(D)
2  N ← COUNTDOCS(D)
3  for each c ∈ C
4  do Nc ← COUNTDOCSINCLASS(D, c)
5     prior[c] ← Nc/N
6     for each t ∈ V
7     do Nct ← COUNTDOCSINCLASSCONTAININGTERM(D, c, t)
8        condprob[t][c] ← (Nct + 1)/(Nc + 2)
9  return V, prior, condprob
```

Рис. 1:

# 3 Текст программы

---

```
1 from collections import defaultdict
2 import re
3 from math import log
4 from typing import Optional
5 import pymorphy2
6
7 morph = pymorphy2.MorphAnalyzer()
8
```

```

APPLYBERNOULLINB( $\mathbb{C}$ ,  $V$ ,  $prior$ ,  $condprob$ ,  $d$ )
1   $V_d \leftarrow \text{EXTRACTTERMSFROMDOC}(V, d)$ 
2  for each  $c \in \mathbb{C}$ 
3  do  $score[c] \leftarrow \log prior[c]$ 
4    for each  $t \in V$ 
5    do if  $t \in V_d$ 
6      then  $score[c] += \log condprob[t][c]$ 
7      else  $score[c] += \log(1 - condprob[t][c])$ 
8  return  $\arg \max_{c \in \mathbb{C}} score[c]$ 

```

Рис. 2:

```

9 AMBIGUOUS_WORD = morph.parse('rak')[0].normal_form # make cyrillic
10
11
12 class Sentence:
13     def __init__(self, sentence: str, vocabulary: list, klass: str = None):
14         self.sent = sentence
15         self.words = string_to_words(sentence)
16         self.norm = normalize_words(self.words)
17         self.klass = klass
18         self.vector = {}
19         self.apply_pos_features() # use part-of-speech fts
20         self.apply_lexical_features(vocabulary) # use lexical features
21
22 # apply part-of-speech features
23 def apply_pos_features(self):
24     amb_word_index = self.norm.index(AMBIGUOUS_WORD)
25     word_left_2 = item_at_index_or_none(self.norm, amb_word_index - 2)
26     word_left_1 = item_at_index_or_none(self.norm, amb_word_index - 1)
27     word_right_1 = item_at_index_or_none(self.norm, amb_word_index + 1)
28     word_right_2 = item_at_index_or_none(self.norm, amb_word_index + 2)
29     self.vector = {
30         # word to the left is present
31         'l_exists': ft_word_exists(word_left_1),
32         # word to the right is present
33         'r_exists': ft_word_exists(word_right_1),
34         # pos features
35         # left -2
36         'l_2_noun': ft_noun(word_left_2),
37         'l_2_adjf': ft_adjf(word_left_2),
38         'l_2_verb': ft_verb(word_left_2),
39         'l_2_comp': ft_comp(word_left_2),
40         'l_2_numr': ft_numr(word_left_2),
41         'l_2_advb': ft_advb(word_left_2),
42         'l_2_prtf': ft_prtf(word_left_2),
43         # left -1
44         'l_1_noun': ft_noun(word_left_1),
45         'l_1_adjf': ft_adjf(word_left_1),
46         'l_1_verb': ft_verb(word_left_1),
47         'l_1_comp': ft_comp(word_left_1),
48         'l_1_numr': ft_numr(word_left_1),
49         'l_1_advb': ft_advb(word_left_1),

```

```

50         'l_1_prtf': ft_prtf(word_left_1),
51         # right +1
52         'r_1_noun': ft_noun(word_right_1),
53         'r_1_adjf': ft_adjf(word_right_1),
54         'r_1_verb': ft_verb(word_right_1),
55         'r_1_comp': ft_comp(word_right_1),
56         'r_1_numr': ft_numr(word_right_1),
57         'r_1_advb': ft_advb(word_right_1),
58         'r_1_prtf': ft_prtf(word_right_1),
59         # right +2
60         'r_2_noun': ft_noun(word_right_2),
61         'r_2_adjf': ft_adjf(word_right_2),
62         'r_2_verb': ft_verb(word_right_2),
63         'r_2_comp': ft_comp(word_right_2),
64         'r_2_numr': ft_numr(word_right_2),
65         'r_2_advb': ft_advb(word_right_2),
66         'r_2_prtf': ft_prtf(word_right_2)
67     }
68
69     # takes dictionary and turns every word 'x' into feature: 'sentence
    contains "x"'
70     def apply_lexical_features(self, vocabulary: list):
71         for word in set(vocabulary):
72             self.vector[word] = word in self.norm
73
74     def satisfies_feature(self, feature_name: str) -> bool:
75         return True == self.vector[feature_name]
76
77     # extracts features that evaluate to true(1) for given sentence
78     def extract_active_features(self) -> list:
79         fts = dict(filter(lambda ft: True == ft[1], self.vector.items()))
80         return list(fts.keys())
81
82     def __repr__(self):
83         return self.__str__()
84
85     def __str__(self):
86         return self.sent
87
88
89     # ===== Feature predicates =====
90     ft_noun = lambda w: 'NOUN' == tag(w)
91     ft_adjf = lambda w: 'ADJF' == tag(w)
92     ft_verb = lambda w: 'VERB' == tag(w)
93     ft_comp = lambda w: 'COMP' == tag(w)
94     ft_numr = lambda w: 'NUMR' == tag(w)
95     ft_advb = lambda w: 'ADVB' == tag(w)
96     ft_prtf = lambda w: 'PRTF' == tag(w)
97     ft_word_exists = lambda w: w is not None
98     # =====
99
100
101     def tag(w: str) -> Optional[str]:
102         return morph.parse(w)[0].tag.POS if w is not None else None
103

```

```

104 def item_at_index_or_none(lst: list, index: int) -> Optional[str]:
105     try:
106         return lst[index]
107     except IndexError:
108         return None
109
110 def read_sentences(filename: str) -> list:
111     contents = open(filename, 'r').read()
112     # remove empty strings
113     return list(filter(lambda s: '' != s, contents.split('\n')))
114
115 def string_to_words(s: str) -> list:
116     clean = re.sub(r'^[a-z]', ' ', s.casefold()) # make cyrillic
117     return re.compile(r'\s+').split(clean)
118
119 # normalizes words and removes functional words
120 def normalize_words(words: list) -> list:
121     tags_to_remove = ['NPRO', 'PRED', 'PREP', 'CONJ', 'PRCL', 'INTJ']
122     normalized = []
123     for w in words:
124         parsed = morph.parse(w)[0]
125         if parsed.tag.POS not in tags_to_remove:
126             normalized.append(parsed.normal_form)
127     return normalized
128
129 def extract_unique_features(sentences: list) -> list:
130     return sentences[0].vector.keys()
131
132 def count_sents_in_class(sentences: list, klass: str) -> int:
133     return len(list(filter(lambda s: klass == s.klass, sentences)))
134
135
136 def count_sents_in_class_satisfying_feature(sentences: list, klass: str,
137     feature: str) -> int:
138     sent_of_klass = list(filter(lambda s: klass == s.klass, sentences))
139     sent_satisfying_ft = list(filter(lambda s: s.satisfies_feature(feature)
140         , sent_of_klass))
141     return len(sent_satisfying_ft)
142
143
144 def train_bernoulli(klasses: list, sents: list) -> (dict, dict, dict):
145     V = extract_unique_features(sents)
146     N = len(sents)
147     cond_prod = defaultdict(dict)
148     prior = {}
149
150     for klass in klasses:
151         Nc = count_sents_in_class(sents, klass)
152         prior[klass] = Nc / N
153
154         for feature in V:
155             Nct = count_sents_in_class_satisfying_feature(sents, klass,
156                 feature)
157             cond_prod[feature][klass] = (Nct + 1) / (Nc + 2)
158     return V, prior, cond_prod

```

```

156
157 def apply_bernoulli(klasses: list, V: list, prior: dict, cond_prob: dict,
158     new_sent: Sentence) -> dict:
159     Vd = new_sent.extract_active_features()
160     score = {}
161
162     for klass in klasses:
163         score[klass] = log(prior[klass])
164         for feature in V:
165             if feature in Vd:
166                 score[klass] += log(cond_prob[feature][klass])
167             else:
168                 score[klass] += log(1 - cond_prob[feature][klass])
169
170     return score
171
172 def main():
173     # make dictionary of normalized words from all train files
174     train_1 = read_sentences('c1.txt')
175     train_2 = read_sentences('c2.txt')
176     vocabulary = set()
177     for s in train_1 + train_2:
178         normalized = normalize_words(string_to_words(s))
179         vocabulary.update(normalized)
180     vocabulary = list(vocabulary)
181
182     # make sentences (Docs) for class 1
183     klass_1 = 'cancer'
184     sent_1 = list(map(lambda s: Sentence(s, vocabulary, klass_1), train_1))
185
186     # make sentences (Docs) for class 2
187     klass_2 = 'crayfish'
188     sent_2 = list(map(lambda s: Sentence(s, vocabulary, klass_2), train_2))
189
190     V, prior, cond_prod = train_bernoulli([klass_1, klass_2], sent_1+sent_2)
191
192     test_sentences = list(map(lambda s: Sentence(s, vocabulary),
193         read_sentences('test.txt')))
194
195     for sent in test_sentences:
196         scores = apply_bernoulli([klass_1, klass_2], V, prior, cond_prod,
197             sent)
198
199         res = klass_1 if max(scores.values())==scores[klass_1] else klass_2
200         print(f'{res}: {sent}: cancer:{scores[klass_1]:.3f}, crayfish:{
201             scores[klass_2]:.3f}')
202
203 if __name__ == '__main__':
204     main()

```

---

Листинг 1: Исходный код программы

## 4 Результаты тестирования

В качестве многозначного слова было выбрано слово *Рак*. Первое значение – бо-

Таблица 1: Результаты работы

Предложение	Cancer, без ча- стеречн	Crayfish, б/чр.	Cancer, часте- речн	Crayfish, часте- речн.
В скандинавских страна питаются вареными раками	-23.178	<b>-21.516</b>	-33.980	<b>-31.258</b>
Известно, что рак может вызываться некоторыми профессиональными вредностями	-18.383	<b>-18.019</b>	<b>-25.765</b>	-26.125
Длина тела широкопалого рака может достигать 25 см и более	-18.383	<b>-18.019</b>	-27.980	<b>-26.462</b>
Употребление экологически чистых продуктов уменьшает вероятность появления рака	-20.781	<b>-20.417</b>	<b>-26.882</b>	-28.213
Понятие рак распространено среди игроков многопользовательских игр, в частности Dota	<b>-18.383</b>	-18.808	<b>-23.211</b>	-24.206

лезнь. Обучающая выборка представлена в Таблице 2. Второе значение – животное. Обучающая выбора - в Таблице 3.

Тестовая выборка и результаты представлены в Таблице 1. Выборка составлена следующим образом:

- 2 предложения, в которых есть слова из соответствующих обучающих выборок.
- 1 предложение, в котором нет слов из обучающих выборок.
- 1 предложение, в котором есть слова из обеих обучающих выборок.
- 1 предложение из другой предметной области (впрочем, может быть отнесено сразу в обе области)

В первом случае (без учета частеречных признаков) алгоритм верно «угадал» одного рака и один рак.

Во втором случае (с учетом) алгоритм верно справился во всех случаях – как простых, так и сложных (без пересечения по словам).

Можно сделать вывод о том, что использование частеречных признаков улучшает результаты как в абсолютном так и относительном значении: значения без учета частеречных признаков ближе друг к другу и сильно подвержены колебаниям. Значения с учетом чр. признаков подвержены им в меньшей степени.

Слово вне обоих контекстов было определено одинаково.

## 5 Выводы

В ходе работы были изучены методы машинного обучения для разрешения лексической неоднозначности. Ожидаемые результаты совпали с реальными.

Таблица 2: Обучающая выборка класса 1

Основной упор в разрабатываемой BostonGene системе делается на интеграцию исследований терапии рака  
 Точные причины развития рака неизвестны  
 Общими симптомами рака являются резкое снижение веса, потеря аппетита, плохое самочувствие  
 Самым распространенным является рак легких  
 группа экспертов по изучению рака составляла классификацию веществ, входящих в продукты питания  
 В клетках рака с высокой частотой возникают разломы  
 Рак может поражать практически любую систему, орган или ткань организма  
 Проблема рака имеет огромное значение  
 Показатели смертности от рака в разных странах сильно варьируют  
 Курение является наиболее значимой причиной рака, определяя его возникновение

Таблица 3: Обучающая выборка класса 2

Цвет рака зависит от условия обитания и состава воды  
 В пищу раков стали употреблять очень много лет назад  
 Молодые раки питаются только растениями  
 Поймать раков в их привычной среде обитания довольно просто  
 Ловить раков башмаком придумали достаточно давно  
 Цвет вареных раков должен быть равномерный и насыщенно красный  
 Регулярное использование мяса раков улучшает самочувствие людей после затяжной простуды  
 Употребление раков может причинить вред организму только в случае, если злоупотреблять ими  
 В состав мяса раков входят различные микроэлементы, которые насыщают организм  
 Очень вкусными получаются раки вареные в пиве