

Московский государственный технический университет им. Н.Э. Баумана
Факультет «Информатика и системы управления»
Кафедра «Системы обработки информации и управления»



Лабораторная работа №5
по курсу «Методы машинного обучения»
«Предобработка и классификация текста»

ИСПОЛНИТЕЛЬ:

Ерохин И.А.
Группа ИУ5-24М

"__" _____ 2022 г.

Находим близость между словами и строим аналогию

```
In [110.]: print(model.similarity('Сперер_8', 'Сопро_8'))
0.512777
```

```
In [111.]: print(model.most_similar(positive=['праздник_8', 'паска_8'], negative=['ханука_8']))

[('брожествю_8', 0.5296328067779541), ('новоселье_8', 0.4783863425254822), ('воскресен
ие_8', 0.44232097268104553), ('сочельник_8', 0.439351886510849), ('именины_8', 0.43770
81902420044), ('счастье_8', 0.42595634857954995), ('христосовавшая_8', 0.410161018371
58203), ('благословение_8', 0.4094565510749817), ('хануи_8', 0.3997984230518341), ('пас
хальный_8', 0.3966675967803955)]
```

Обучим word2vec на наборе данных "fetch_20newsgroups"

```
In [112.]: import re
import pandas as pd
import numpy as np
from typing import Dict, Tuple
from sklearn.metrics import accuracy_score, balanced_accuracy_score
from sklearn.feature_extraction.text import CountVecorizer, TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import Pipeline
from nltk import WordPunctTokenizer
from nltk.corpus import stopwords
import nltk
nltk.download('stopwords')

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
```

```
Out[112.]: True
```

```
In [113.]: categories = ["rec.motorcycles", "rec.sport.baseball", "sci.electronics", "sci.med"]
newsgroups = fetch_20newsgroups(subset="train", categories=categories)
data = newsgroups['data']
```

```
In [114.]: # Подготавлим корпус
stop_words = stopwords.words('english')
tok = WordPunctTokenizer()
for line in newsgroups['data']:
    line1 = line.strip().lower()
    line1 = re.sub("[a-zA-Z]", "", line1)
    text_tok = tok.tokenize(line1)
    text_tok[1] = [w for w in text_tok if not w in stop_words]
    corpus.append(text_tok)
```

```
In [115.]: corpus[:5]
```

```
Out[115.]: [['nmendel',
'unix',
'amberst',
'edu',
'nathaniel',
'mendell',
'batting',
'bike',
'advice',
'organization',
'amberst',
'college',
'x',
'reader',
'tin',
'version',
'pi',
'lines',
'umms',
'quarius',
'ccom',
'edwards',
'subject',
'photography',
'reply',
'grance',
'equarius',
'rosemount',
'ccom',
'grant',
'edwards',
'organization',
'rosemount',
'inc',
'lines',
'trump',
'potting',
'host',
'aquarius',
'suprao',
'st',
'unocal',
'com',
'richard',
'ottolini',
'wlices',
'living',
'things',
'maintain',
'email',
'electric',
'fields',
'enhance',
'certain',
'chemical',
'reactions',
'promote',
'communication',
'states',
'cell',
'communicate',
'specialized',
'system',
'perhaps',
'uses',
'true',
'electric',
'fields',
'change',
'location',
'time',
'large',
'organism',
'also',
'true',
'special',
'photographic',
'applying',
'external',
'fields',
'interact',
'fields',
'resistances',
'caused',
'interesting',
'pictures',
'really',
'kirlian',
'photography',
'taking',
'pictures',
'corona',
'discharge',
'objects',
'animate',
'animate',
'fields',
'applied',
'objects',
'millions',
'time',
'larger',
'biologically',
'created',
'fields',
'want',
'created',
'biologically',
'created',
'electric',
'fields',
'got',
'use',
'low',
'noise',
'high',
'perhaps',
'sensors',
'typical',
'types',
'eggs',
'kirlian',
'photography',
'phun',
'physics',
'type',
'stuff',
'right',
'soaking',
'seal',
'extra',
'fine',
'steel',
'wool',
'liquid',
'oxygen',
'bitting',
'hammer',
'like',
'like',
'setup',
'fun',
'possibly',
'possibly',
'pictures',
'diagnostic',
'disease',
'better',
'understood',
'perhaps',
'probably',
'grant',
'edwards',
'vote',
'rosemount',
'inc',
'well',
'taped',
'half',
'cooked',
'ill',
'conceived',
'grance',
'aquarius',
'rosemount',
'com',
'tax',
'deferred',
['lin',
'sun',
'scri',
'fau',
'edu',
'nemo',
'subject',
'bates',
'method',
'myopia',
'ray',
'lin',
'fau',
'fau',
'distribution',
'na',
'organization',
'scri',
'florida',
'state',
'university',
'lines',
'method',
'work',
'fired',
'heard',
'newsgroup',
'seveal',
'years',
'ago',
'got',
'hold',
'book',
'improve',
'sight',
'simple',
'daily',
'drills',
'relaxation',
'margaret',
'cothett',
'authorized',
'instructor',
'bates',
'method',
'published',
'talks',
'vision',
'improvement',
'relaxation',
'exercise',
'study',
'whether',
'method',
'actually',
'works',
'works',
'actually',
'shortening',
'previously',
'elongated',
'eyeball',
'increasing',
'lens',
'ability',
'flatten',
'lords',
'compensate',
'long',
'eyeball',
'since',
'myopia',
'result',
'eyeball',
'elongation',
'seems',
'logical',
'approach',
'correction',
'find',
'way',
'reverse',
'process',
'e',
'shorten',
'somehow',
'preferably',
'non',
'surgically',
'seem',
'studies',
'find',
'row',
'rick',
'works',
'changing',
'curvature',
'cornea',
'compensate',
'shape',
'eyeball',
'way',
'train',
'muscles',
'shorten',
'eyeball',
'back',
'correct',
'length',
'would',
'even',
'better',
'bates',
'idea',
'right',
'thanks',
'information',
['mcovingt',
'aisun',
'sai',
'uga',
'edu',
'michael',
'conviction',
'subject',
'sun',
'parts',
'time',
'mntp',
'potting',
'host',
'aisun',
'ai',
'uga',
'edu',
'organization',
'ai',
'program',
'university',
'georgia',
'athens',
'lines',
'parts',
'reminds',
'something',
'chemist',
'said',
'gram',
'dye',
'costs',
'dollar',
'comes',
'liter',
'jar',
'costs',
'dollar',
'e',
'charge',
'almost',
'exclusively',
'packaging',
'discovering',
'chemical',
'articular',
'case',
'byproduct',
'coat',
'almost',
'nothing',
'intrinsically',
'michael',
'conviction',
'associate',
'research',
'scientist',
'artificial',
'intelligence',
'programs',
'ai',
'uga',
'edu',
'university',
'georgia',
'phone',
'athens',
'georgia',
'e',
'amateur',
'radio',
'et',
'tm']],
['tamy',
'vandenboom',
'launchpad',
'unc',
'edu',
'tammy',
'vandenboom',
'subject',
'scott',
'spot',
'teactiles',
'mtp',
'posting',
'host',
'lambada',
'oit',
'unc',
'edu',
'organization',
'north',
'carolina',
'extended',
'bulletin',
'bates',
'service',
'distribution',
'na',
'lines',
'husband',
'wake',
'three',
'days',
'ago',
'smal',
'sore',
'spot',
'spot',
'size',
'nickel',
'bottom',
'side',
'nodes',
'lumps',
'little',
'sore',
'spot',
'says',
'reminds',
'bruise',
'feels',
'recollection',
'bitting',
'anything',
'like',
'would',
'cause',
'bruise',
'assures',
'remember',
'something',
'like',
'closes',
'might',
'somewhat',
'hypochondriac',
'sp',
'sure',
'gonna',
'die',
'thanks',
'opinions',
'expressed',
'necessarily',
'university',
'north',
'carolina',
'chapel',
'hill',
'campus',
'office',
'information',
'technology',
'experimental',
'bulletin',
'board',
'service',
'internet',
'launchpad',
'unc',
'edu']]
```

```
In [116.]: #time model_imdb = word2vec.Word2Vec(corpus, workers=4, min_count=10, window=10, samp=
CPU times: user 4.73 s, sys: 23.8 ms, total: 4.75 s
Wall time: 2.95 s
```

```
In [117.]: # Проверим, что модель обучилась
print(model_imdb.wv.most_similar(positive=['find'], topn=5))

[('using', 0.9903143644332886), ('circuits', 0.9873676300048828), ('used', 0.986776113
5101318), ('work', 0.9865034818649292), ('etc', 0.9864989591919263)]
```

```
In [118.]: def sentiment_2(v, c):
    model = Pipeline([
        ("vectorizer", v),
        ("classifier", c)])
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    print(accuracy_score_for_classes(y_test, y_pred))
```

Проверка качества работы модели word2vec

```
In [119.]: class EmbeddingVectorizer(object):
    """
    Для текста усредним вектора входящих в него слов
    """
    def __init__(self, model):
        self.model = model
    def fit(self, X, y):
        return self
    def transform(self, X):
        return np.array([np.mean([
            self.model[w] for w in words if w in self.model])
            for words in X])
```

```
In [120.]: def accuracy_score_for_classes(y_true: np.ndarray,
y_pred: np.ndarray) -> Dict[int, float]:
    """
    Вычисление метрики accuracy для каждого класса
    y_true - истинные значения классов
    y_pred - предсказанные значения классов
    Возвращает словарь: ключ - метка класса,
    значение - accuracy для данного класса
    """
    # Для удобства функции сформируем Pandas DataFrame
    d = {'t': y_true, 'p': y_pred}
    df = pd.DataFrame(data=d)
    # Ищем классы
    classes = np.unique(y_true)
    # Результирующий словарь
    res = dict()
    # Проходим по всем классам
    for c in classes:
        # отфильтруем данные, которые соответствуют
        # текущей метке класса в истинных значениях
        temp_data_fit = df[df['t']==c]
        # расчет accuracy для заданной метки класса
        temp_acc = accuracy_score(
            temp_data_fit['t'].values,
            temp_data_fit['p'].values)
        # сохранение результата в словарь
        res[c] = temp_acc
    return res
```

```
def print_accuracy_score_for_classes(y_true: np.ndarray,
y_pred: np.ndarray):
    """
    Вывод метрики accuracy для каждого класса
    """
    accs = accuracy_score_for_classes(y_true, y_pred)
    if len(accs)>0:
        print("Метрика Accuracy")
        for i in accs:
            print("{} \t {}".format(i, accs[i]))
```

```
In [121.]: # Обучающая и тестовая выборки
boundary = 1500
X_train = corpus[boundary:]
X_test = corpus[boundary:]
y_train = newsgroups['target'][boundary:]
y_test = newsgroups['target'][boundary:]
```

```
In [122.]: sentiment_2(EmbeddingVectorizer(model_imdb.wv), LogisticRegression(C=5.0))

/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:818: Converge
Warning: lbfgs failed to converge (status=1):
STOP: FINAL No. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solvers options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
ExtraWarning:MSG=LOGISTIC_SOLVER_CONVERGENCE_MSG,
Метка Accuracy
0 0.8333333333333334
1 0.9223300970873787
2 0.7247706422018348
3 0.7192982456140351
```

Как видно из результатов проверки качества моделей, лучшее качество показал CountVecorizer