



▼ Classification and Analysis of text data

Copyright 2018 The BUPT Zhengyuan Zhu.

Licensed under the Apache License, Version 2.0 (the "License").



[View all sources on GitHub](#)

Affiliation: BUPT

Author1:824zzy(计算机学院-2018140455-朱正源)

Author2:Regulusyy(计算机学院-2018140506-杨莹)

References

- [CNN wiki](#)
- [RNN wiki](#)
- [Support vector machine wiki](#)
- [python3:csv文件的读写](#)
- [pyhanlp 分词与词性标注](#)
- [python结巴分词、jieba加载停用词表](#)
- [中文常用停用词表](#)
- [python读取和存储dict\(\)与json格式文件](#)
- [Python爬虫之爬取动态页面数据](#)
- [824zzy（朱正源）的微博爬虫](#)
- [6 Easy Steps to Learn Naive Bayes Algorithm \(with codes in Python and R\)](#)
- <https://medium.com/jatana/report-on-text-classification-using-cnn-rnn-han-f0e887214d5f>
- [Naive Bayes Tutorial: Naive Bayes Classifier in Python](#)
- [Let's implement a Gaussian Naive Bayes classifier in Python](#)
- [Support Vector Machines with Scikit-learn](#)
- [python中sklearn实现交叉验证](#)
- [Practical Text Classification With Python and Keras](#)
- [Text Preprocessing - Keras](#)

▼ Crawler demo:

In this part, we will show you how to build a crawler from scratch. It is a little tricky but not difficult enough.

▼ Setup packages to Colab Virtual Machine.

- scrapy: especially for using XPATH to parse the html tree
- tqdm: a common tool for displaying the processing of ForLoop
- retrying: package for preventing connect lose
- requests: speed up for efficiency of crawler
- requests: basic package to get html

```
1 !pip install scrapy
2 !pip install tqdm
3 !pip install retrying
4 !pip install grequests
5 !pip install requests
```



Collecting scrapy

Downloading <https://files.pythonhosted.org/packages/5d/12/a6197eaf97385e961>

100% 256kB 11.7MB/s

▼ Import packages

Collecting Mediated -12 1 0 (from survey)

```
1 import requests
2 import grequests
3 import scrapy
4 from tqdm import tqdm
5 from retrying import retry
6 import time
7 import random
8 import json
9 import csv
10 import pickle
```

```
↳ /usr/local/lib/python3.6/dist-packages/grequests.py:21: MonkeyPatchWarning: Monkey-patching all is discouraged; if you're having issues using it, you may want to use curious_george.patch_all(thread=False, select=False)
```

```
Requirement already satisfied: lxml in /usr/local/lib/python3.6/dist-packages
```

▼ Auxiliary functions for crawler

[illegible]

▼ Generate Headers Randomly

Downloading <https://files.pythonhosted.org/packages/10/17/1d100a6aaa0aa450/>

```

1 userAgent_file = [
2 "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/535.1 (KHTML, like Gecko) Chro
3 "Mozilla/5.0 (Windows NT 6.1; WOW64; rv:6.0) Gecko/20100101 Firefox/6.0",
4 "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/534.50 (KHTML, like Gecko) Ver
5 "Opera/9.80 (Windows NT 6.1; U; zh-cn) Presto/2.9.168 Version/11.50",
6 "Mozilla/5.0 (compatible; MSIE 9.0; Windows NT 6.1; Win64; x64; Trident/5.0; .N
7 "Mozilla/4.0 (compatible; MSIE 8.0; Windows NT 6.1; WOW64; Trident/4.0; SLCC2;
8 "Mozilla/4.0 (compatible; MSIE 8.0; Windows NT 5.1; Trident/4.0; GTB7.0)",
9 "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 5.1)",
10 "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1)",
11 "Mozilla/5.0 (Windows; U; Windows NT 6.1; ) AppleWebKit/534.12 (KHTML, like Gec
12 "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 6.1; WOW64; Trident/5.0; SLCC2;
13 "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 6.1; WOW64; Trident/5.0; SLCC2;
14 "Mozilla/5.0 (Windows; U; Windows NT 6.1; en-US) AppleWebKit/534.3 (KHTML, like
15 "Mozilla/5.0 (compatible; MSIE 9.0; Windows NT 6.1; WOW64; Trident/5.0; SLCC2;
16 "Mozilla/5.0 (Windows NT 6.1) AppleWebKit/535.1 (KHTML, like Gecko) Chrome/13.0
17 "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 6.1; WOW64; Trident/5.0; SLCC2;
18 "Mozilla/5.0 (compatible; MSIE 9.0; Windows NT 6.1; WOW64; Trident/5.0)",
19 ]
20 class Headers:
21     @staticmethod
22     def getHeaders():
23         userAgentList = []
24         for line in userAgent_file:
25             userAgentList.append({
26                 'User-Agent': line.strip(),
27                 #'User-Agent': 'Mozilla/5.0 (compatible; Googlebot/2.1; +http:/
28                 'Referer': 'http://cn.bing.com/',
29                 'X-Forwarded-For': '%s.%s.%s.%s' % (
30                     random.randint(50, 250), random.randint(50, 250), random.ran
31                     'CLIENT-IP': '%s.%s.%s.%s' % (
32                         random.randint(50, 250), random.randint(50, 250), random.ran
33                 })
34             userAgent = random.sample(userAgentList, 1)
35             return userAgent[0]
36
37 # test case
38 print(Headers.getHeaders())

```

```
↳ {'User-Agent': 'Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 6.1; WOW64; Tri
```

▼ Base class for Crawlers

```
1 import json
2
3
4
5 class base_class(object):
6     def __init__(self, website, website_url, header=None):
7         """
8
9         """
10        self.cookies = {}
11        self.website = website
12        self.website_url = website_url
13        self.header = header
14        self.entrance_html = requests.get(self.website_url).content
15        self.sections = self.section_list()
16
17    def section_list(self):
18        pass
19
20    def section_urls(self):
21        pass
22
23    def parse_body(self):
24        pass
25
26    def ajax_url(self):
27        pass
28
29    def article_urls_dict(self):
30        pass
31
32    def display(self):
33        pass
```

▼ Jiandan Crawler

```
1 class Jiandan(base_class):
2     def __init__(self, **kwargs):
3         super(Jiandan, self).__init__(**kwargs)
4
5     def sec_subsec_dict(self, sections, sub_sections):
6         """build a dictionary according to sections and subsections
7
8         # Arguments:
9         sections: a list contains sections' string.
10        sub_sections: a list contains sub_sections' string.
11
12        # Returns:
13        sec_subsec_dict: a dictionary whose section as key and sub_sections as va
14
15        """
16        sec_subsec_dict = {}
17        for index, section in enumerate(sections):
18            sec_subsec_dict[section] = [sub_sections[index + 7 * i] for i in range(6)]
19        return sec_subsec_dict
20
21
22    def section_list(self):
23        """get section from each website. Sections could be seem as the label of
24        classification in the case of saving manpower.
25
26        # Arguments:
```

```

27     None:
28
29     # Returns:
30     section: a list contains section name that is a str
31     """
32     section = scrapy.Selector(text=self.entrance_html).xpath('//*[@id="header"]')
33     return section
34
35 def sub_section_list(self):
36     sub_sections = scrapy.Selector(text=self.entrance_html).xpath('//*[@id="hea
37     return sub_sections
38
39 def section_urls(self, sub_sections):
40     """get urls from each section in Jiandan website
41
42     # Arguments:
43     sub_sections: a list contains all the urls to be concatenated.
44
45     # Returns:
46     section_urls: a list contains all the urls after concatenate.
47     """
48     section_urls = ["http://jandan.net/tag/"+item for item in sub_sections]
49     return section_urls
50
51 @retry(stop_max_attempt_number=10)
52 def retry_dict_request(self, section_url):
53     reps = (grequests.get(section_url + "/page/" + str(page)) for page in range
54     atk_urls = []
55     for rep in grequests.map(reps):
56         try:
57             article_urls = [scrapy.Selector(text=rep.text).xpath('//*[@id="content"
58             atk_urls.append(article_urls)
59         except:
60             return atk_urls
61
62 def sub_sec_atk_dict(self, section_dict):
63     sec_atk_dict = {}
64     for section_name, section_url in section_dict.items():
65         start_time = time.time()
66         atk_urls = self.retry_dict_request(section_url)
67         sec_atk_dict[section_name] = atk_urls
68         print("Spended ", str(int(time.time() - start_time)), " seconds on subsec
69             section_name, " to get dictionary that section as key and article u
70     return sec_atk_dict
71
72 @retry(stop_max_attempt_number=10)
73 def retry_page_urls(self, page_urls, section, sub_section):
74     page_parsed_data = []
75     reps = (grequests.get(url) for url in page_urls)
76     for rep in grequests.map(reps):
77         title = scrapy.Selector(text=rep.text).xpath('//*[@id="content"]/div[2]/h
78         main_body = scrapy.Selector(text=rep.text).xpath('//*[@id="content"]/div[
79         body_str = ''
80         for sen in main_body:
81             body_str += sen
82         meta_tuple = (section[0], sub_section, title[0], body_str)
83         page_parsed_data.append(meta_tuple)
84     return page_parsed_data
85
86 def parsing_dict(self, sec_atk_dict, section_subsection_dict):
87     """ deal with ajax dynamic loading problem.
88
89     # Arguments:
90     sections: a list contains all the sections in website
91     website: a str represents target website
92     headers: a str we generate in previous stage
93
94     # Returns:
95     article_urls_dict: a dict whose keys are sections and values
96     are articles' urls for this section
97     """
98     parsed_data = []
99     for sub_section, article_urls in sec_atk_dict.items():
100         section = [k for k, v in section_subsection_dict.items() if sub_section i

```

```

101     print("we are parsing section: ", sub_section, " now\n")
102     start_time = time.time()
103     for page_urls in tqdm(article_urls):
104         page_parsed_data = self.retry_page_urls(page_urls, section, sub_section)
105         parsed_data.append(page_parsed_data)
106     print("Spended ", str(int(time.time()-start_time)), " seconds to parse fi
107         sub_section)
108
109     return parsed_data
110
111
112 def display(self):
113     """display result of each function in class.
114     """
115
116
117     sections = self.section_list()
118     print("Sections(class label) are:", sections)
119     print('\n' + '-'*50 + '\n')
120
121     sub_sections = self.sub_section_list()
122     print("Subsections are:", sub_sections)
123     print('\n' + '-'*50 + '\n')
124
125     section_subsection_dict = self.sec_subsec_dict(sections, sub_sections)
126     print("The Dictionary for section as key and subsection as value is ", secti
127     print('\n' + '-'*50 + '\n')
128
129     sub_section_urls = self.section_urls(sub_sections)
130     print("Subsections urls are: ", sub_section_urls)
131     print('\n' + '-'*50 + '\n')
132
133     section_dict = dict(zip(sub_sections, sub_section_urls))
134     section_article_dict = self.sub_sec_atk_dict(section_dict)
135
136     print(section_article_dict)
137
138     print('\n' + '-'*50 + '\n')
139     parsed_data = self.parsing_dict(section_article_dict, section_subsection_di
140
141     print("Some examples of tuple data")
142     for index in range(10):
143         print(parsed_data[index])
144
145     return parsed_data
146
147

```

```

1 header = Headers.getHeaders()
2 spider = Jiandan(website="煎蛋网", website_url="http://jandan.net/", header=head
3 print('\n'*2 + '='*50 + '\n'*2)
4 parsed_data = spider.display()
5 print('\n'*2 + '='*50 + '\n'*2)

```



```
=====
```

Sections(class label) are: ['科学', '技术', '极客', '脑洞', '故事', '人类', '折腾'

Subsections are: ['走进科学', 'TECH', 'GEEK', 'DIY', '冷新闻', '女性', '减肥', '无厘头研究', '人工智能', 'MEME', '艺术', '爷有钱', '熊孩子', '整形', '天文', '无人机', 'QUORA', '设计', '致富信息', '大丈夫', '旅游', 'NASA']

The Dictionary for section as key and subsection as value is {'科学': ['走进科学', 'TECH', 'GEEK', 'DIY', '冷新闻', '女性', '减肥', '无厘头研究', '人工智能', 'MEME', '艺术', '爷有钱', '熊孩子', '整形', '天文', '无人机', 'QUORA', '设计', '致富信息', '大丈夫', '旅游', 'NASA']

Subsections urls are: ['<http://jandan.net/tag/走进科学>', '<http://jandan.net/tag/TECH>', '<http://jandan.net/tag/GEEK>', '<http://jandan.net/tag/DIY>', '<http://jandan.net/tag/冷新闻>', '<http://jandan.net/tag/女性>', '<http://jandan.net/tag/减肥>', '<http://jandan.net/tag/无厘头研究>', '<http://jandan.net/tag/人工智能>', '<http://jandan.net/tag/MEME>', '<http://jandan.net/tag/艺术>', '<http://jandan.net/tag/爷有钱>', '<http://jandan.net/tag/熊孩子>', '<http://jandan.net/tag/整形>', '<http://jandan.net/tag/天文>', '<http://jandan.net/tag/无人机>', '<http://jandan.net/tag/QUORA>', '<http://jandan.net/tag/设计>', '<http://jandan.net/tag/致富信息>', '<http://jandan.net/tag/大丈夫>', '<http://jandan.net/tag/旅游>', '<http://jandan.net/tag/NASA>']

Spended	56	seconds on subsection	走进科学	to get dictionary that section as
Spended	10	seconds on subsection	TECH	to get dictionary that section as
Spended	64	seconds on subsection	GEEK	to get dictionary that section as
Spended	10	seconds on subsection	DIY	to get dictionary that section as
Spended	52	seconds on subsection	冷新闻	to get dictionary that section as
Spended	52	seconds on subsection	女性	to get dictionary that section as
Spended	17	seconds on subsection	减肥	to get dictionary that section as
Spended	35	seconds on subsection	无厘头研究	to get dictionary that section
Spended	10	seconds on subsection	人工智能	to get dictionary that section a
Spended	6	seconds on subsection	MEME	to get dictionary that section as
Spended	23	seconds on subsection	艺术	to get dictionary that section as
Spended	68	seconds on subsection	爷有钱	to get dictionary that section as
Spended	42	seconds on subsection	熊孩子	to get dictionary that section as
Spended	6	seconds on subsection	整形	to get dictionary that section as
Spended	10	seconds on subsection	天文	to get dictionary that section as
Spended	6	seconds on subsection	无人机	to get dictionary that section as
Spended	10	seconds on subsection	QUORA	to get dictionary that section as
Spended	26	seconds on subsection	设计	to get dictionary that section as
Spended	21	seconds on subsection	致富信息	to get dictionary that section a
Spended	9	seconds on subsection	大丈夫	to get dictionary that section as
Spended	8	seconds on subsection	旅游	to get dictionary that section as
Spended	7	seconds on subsection	NASA	to get dictionary that section as

▼ Write parsed data into CSV format

Spended 8 seconds on subsection 安全警示 to get dictionary that section as

```
1 import csv
2 import codecs
3
4 class File_IO(object):
5     def __init__(self, read_file=None, write_file=None, headers=None, data=None):
6         self.read_file = read_file
7         self.write_file = write_file
8         self.headers = headers
9         self.parsed_data = data
10
11     def read_csv(self):
12         """ reading CSV Files with Pandas
13
14         # Arguments:
15         name: file name of website without suffix, str
16
17         # Returns:
```

```

18         df: data frame contains
19         """
20
21         with codecs.open(self.read_file + ".csv", 'r', encoding='utf-8') as csv_fil
22             csv_reader = csv.reader(csv_file)
23             csv_data = [row for row in csv_reader]
24             csv_file.close()
25         return csv_data
26
27     def write_csv(self, seg=False):
28         """ writing CSV Files with Pandas
29
30         Arguments:
31             name: file name of website text data without suffix, str
32
33         Returns:
34             f_csv:
35
36             """
37         with codecs.open(self.write_file + '.csv', 'w', encoding='utf-8') as f:
38             f_csv = csv.writer(f)
39             if self.headers:
40                 f_csv.writerow(self.headers)
41             for row in self.parsed_data:
42                 if not seg:
43                     for item in row:
44                         f_csv.writerow(item)
45                 else:
46                     f_csv.writerow(row)
47
48             f.close()
49         return f_csv
50
51     def read_json(self):
52         """
53         Arguments:
54
55         Return:
56
57         """
58         with codecs.open(self.read_file + '.json', 'r') as f:
59             data = list()
60             for line in f:
61                 # load values of key
62                 data.append(json.load(line))
63             f.close()
64         return data
65
66
67
68     def write_json(self):
69         with codecs.open(self.write_file + '.json', 'w') as f:
70             json.dump(self.parsed_data, f, ensure_ascii=False)
71             f.close()
72         return f
73

```

```

1 header = ('类别', '子类别', '标题', '正文')
2 writer = File_IO(write_file="Jiandan_final", headers=header, data=parsed_data)
3 Jiandan = writer.write_csv()

```

100% |██████████| 1309/1309 [00:00<00:00, 1613.96it/s]

▼ Data Preprocessing

In this section we are going to use Jieba to deleting stop word and segementation.

▼ Mount to Google Drive

```
1 # Install the PyDrive wrapper & import libraries.
2 # This only needs to be done once in a notebook.
3 !pip install -U -q PyDrive
4 from pydrive.auth import GoogleAuth
5 from pydrive.drive import GoogleDrive
6 from google.colab import auth
7 from oauth2client.client import GoogleCredentials
8
9 # Authenticate and create the PyDrive client.
10 # This only needs to be done once in a notebook.
11 auth.authenticate_user()
12 gauth = GoogleAuth()
13 gauth.credentials = GoogleCredentials.get_application_default()
14 up_drive = GoogleDrive(gauth)
15
16 from google.colab import drive
17 down_drive = drive.mount('/gdrive')
```

➞ Go to this URL in a browser: <https://accounts.google.com/o/oauth2/auth?client>

Enter your authorization code:

.....

Mounted at /gdrive

```
1 # check out mount result
2 print("Origin dir is:")
3 !ls
4 import os
5 os.chdir("../gdrive/My Drive")
6 print("Now the dir has changed to:")
7 !ls
```

➞ Origin dir is:

adc.json sample_data

Now the dir has changed to:

action_data_car_racing_0.npy	MLP_model.h5
action_data_car_racing_1.npy	NB_model.m
action_data_car_racing_2.npy	papers
auth.ipynb	pdf
BiLSTM.h5	SVM_model.m
'CAIC Shared Code'	Tensorflow_eager_mode.ipynb
'Code Chips'	test_case2.ipynb
'Colab Notebooks'	testcase3.ipynb
Conv_model.h5	Untitled0.ipynb
Dconv_model.h5	vae_weights.h5
homeworks	中文停用词表.txt
Jiandan_final.csv	哈工大停用词表.txt
Jiandan_re_segment.csv	四川大学机器智能实验室停用词库.txt
Jiandan_segment.csv	百度停用词表.txt

▼ Segmentation and Delete Stop Words

▼ Download stop word dictionary from Internet

```
1 !pip install jieba
2 # download a stop word dictionary from network
3 !wget https://raw.githubusercontent.com/goto456/stopwords/master/%E7%99%BE%E5%B
```

```
4 !wget https://raw.githubusercontent.com/goto456/stopwords/master/%E5%9B%9B%E5%B
5 !wget https://raw.githubusercontent.com/goto456/stopwords/master/%E5%93%88%E5%B
6 !wget https://raw.githubusercontent.com/goto456/stopwords/master/%E4%B8%AD%E6%9
7 !ls
```



Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 151.101.0.
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|151.101.0.

▼ Apply dictionary to corpus

```
1  ## basic demo for segmentation
2  import jieba
3  from tqdm import tqdm
4  import re
5
6  f_o = File_IO(read_file='Jiandan_final')
7  corpus = f_o.read_csv()
8  print("Example of corpus: ", corpus[1])
9
10 print("-----")
11 docs = [doc[1]+doc[2]+doc[3] for doc in corpus[1:]]
12 labels = [doc[0] for doc in corpus[1:]]
13
14 print("Example of Segment result: ", [word for word in jieba.cut(docs[3])])
15 print("-----")
16 segmented_corpus = [jieba.cut(doc) for doc in docs]
17
18 stop_word_dict = ["中文停用词表.txt", "哈工大停用词表.txt",
19                  "四川大学机器智能实验室停用词库.txt", "百度停用词表.txt"]
20
21 stop_words_list = []
22 # delete stopping word
23 for each_dict in stop_word_dict:
24     with open(each_dict, 'r') as word_file:
25         stop_words = word_file.readlines()
26         each_stop_words = [item[:-1] for item in stop_words]
27         stop_words_list += each_stop_words
28
29
30 stop_words_list = list(set(stop_words_list))
31 print("\n停用词表: ", stop_words_list)
32 print("-----")
33
34 final_corpus = []
35 for index, segment_list in enumerate(segmented_corpus):
36     final = ''
37     for seg in segment_list:
38         if seg not in stop_words_list:
39             final = final + seg + ' '
40     final_corpus.append([final[:-1], labels[index]])
41
42
43 print("Example of final corpus: ", final_corpus[1])
44 print("-----")
45 print(len(final_corpus))
46 f_i = File_IO(write_file='Jiandan_segment', data=f, headers=['context', 'label'])
47 seg_corpus = f_i.write_csv(seg=True)
48
```

☞ Example of corpus: ['科学', '走进科学', '为什么袋熊的便便是方块状的', '有屎以来, 只有屎', '科学', '走进科学', '为什么袋熊的便便是方块状的', '有屎以来, 只有屎']

Example of Segment result: ['走进', '科学', '超级', '沙漠', '里', '的', '生命', '走进', '科学', '超级', '沙漠', '里', '的', '生命']

停用词表: ['替', '上去', '长话短说', '若非', 'x', '越是', '从古至今', '据悉', 'don']

Example of final corpus: ['走进 科学 走近 科学 为啥 加州 山火 烧足 一个月 11 月 号', '走进 科学 走近 科学 为啥 加州 山火 烧足 一个月 11 月 号']

942480

▼ Extract processed data as DataFrame

```

1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.multiclass import OneVsRestClassifier
5
6 print("-----")
7
8 df = pd.read_csv("Jiandan_segment.csv")

```

➞ -----

▼ Naive Bayes implementation

In this section, we implement Naive Bayes as Baseline.

▼ Import Libraries

```

1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.svm import LinearSVC
3 from sklearn import svm
4 from sklearn.pipeline import Pipeline
5 from sklearn.metrics import accuracy_score
6 from sklearn import metrics
7 from sklearn.model_selection import StratifiedKFold
8 from sklearn.decomposition import KernelPCA
9 from sklearn.externals import joblib

```

▼ Implementation of Naive Bayes

```

1 from collections import Counter
2 from collections import defaultdict
3 from math import log
4
5
6 class MultinomialNB:
7     """Hybrid implementation of Naive Bayes.
8     Supports discrete and continuous features.
9     """
10
11     def __init__(self, extract_features, use_smoothing=True):
12         """Create a naive bayes classifier.
13         :param extract_features: Callback to map a feature vector to discrete a
14         :param use_smoothing: Whether to use smoothing when calculating probabi
15         """
16         self.priors = defaultdict(dict)
17
18         self.label_counts = Counter()
19         self.discrete_feature_vectors = DiscreteFeatureVectors(use_smoothing)
20         self.continuous_feature_vectors = ContinuousFeatureVectors()
21         self.extract_features = extract_features
22         self._is_fitted = False
23
24     def fit(self, design_matrix, target_values):
25         """Fit model according to design matrix and target values.
26         :param design_matrix: Training examples with dimension m x n,
27                             where m is the number of examples,
28                             and n is the number of features.
29         :param target_values: Target values with dimension m,
30                             where m is the number of examples.
31         :return: self
32         """
33         for i, training_example in enumerate(design_matrix):

```

```

34         label = target_values[i]
35         self.label_counts[label] += 1
36         features = self._extract_features(training_example)
37         for j, feature in enumerate(features):
38             if feature.is_continuous():
39                 self.continuous_feature_vectors.add(label, j, feature)
40             else:
41                 self.discrete_feature_vectors.add(label, j, feature)
42
43         total_num_records = len(target_values)
44         for label in set(target_values):
45             self.priors[label] = self.label_counts[label] / total_num_records
46             self.continuous_feature_vectors.set_mean_variance(label)
47
48         self._is_fitted = True
49         return self
50
51     def predict(self, test_set):
52         """Predict target values for test set.
53         :param test_set: Test set with dimension m x n,
54                         where m is the number of examples,
55                         and n is the number of features.
56         :return: Predicted target values for test set with dimension m,
57                 where m is the number of examples.
58         """
59         self._check_is_fitted()
60
61         predictions = []
62         for i in range(len(test_set)):
63             result = self.predict_record(test_set[i])
64             predictions.append(result)
65         return predictions
66
67     def predict_record(self, test_record):
68         """Predict the label for the test record.
69         Maximizes the log likelihood to prevent underflow.
70         :param test_record: Test record to predict a label for.
71         :return: The predicted label.
72         """
73         self._check_is_fitted()
74
75         log_likelihood = {k: log(v) for k, v in self.priors.items()}
76         for label in self.label_counts:
77             features = self._extract_features(test_record)
78             for i, feature in enumerate(features):
79                 probability = self._get_probability(i, feature, label)
80                 try:
81                     log_likelihood[label] += log(probability)
82                 except ValueError as e:
83                     pass
84         return max(log_likelihood, key=log_likelihood.get)
85
86     def _check_is_fitted(self):
87         if not self._is_fitted:
88             raise NotFittedError(self.__class__.__name__)
89
90     def _get_probability(self, feature_index, feature, label):
91         if feature.is_continuous():
92             probability = self.continuous_feature_vectors.probability(label,
93                                                                    feature_i
94                                                                    feature_index,
95                                                                    feature,
96                                                                    self.label_
97                                                                    self.label_
98                                                                    self.label_
99         return probability

```

▼ Build Pipeline of Naive Bayes with TF-IDF as Text feature extraction

In information retrieval or text mining, the term frequency – inverse document frequency (also called tf-idf), is a well know method to evaluate how important is a word in a document. tf-idf are is a very

interesting way to convert the textual representation of information into a Vector Space Model (VSM), or into sparse features.

```
1 # Define a pipeline combining a text feature extractor with multi lable classif
2 NB_pipeline = Pipeline([
3     ('tfidf', TfidfVectorizer()),
4     ('clf', OneVsRestClassifier(MultinomialNB(fit_prior=True, class
5     ]).
```

▼ Fit & Predict

▼ K-Fold Cross Validation

The general idea of K cross tests is to roughly divide the data into K sub-samples. One sample is taken as the verification data and the remaining k-1 samples are taken as the training data.

Measure the model with Stratified k-fold

Than KFold StratifiedKFold () this function is used, the advantage of the k data data set on a percentage basis, the percentage for each category in the training set and test set are the same, so that can not have a certain categories of data in the training set and test set is not this kind of situation, also not in training all in the test set, this will lead to worse results.

```
1 X = df.context
2 y = df.label
3 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=824)
4 for train_index, test_index in skf.split(X, y):
5
6     X_train, X_test = X[train_index], X[test_index]
7     y_train, y_test = y[train_index], y[test_index]
8     NB_pipeline.fit(X_train, y_train)
9     prediction = NB_pipeline.predict(X_test)
10    print(metrics.classification_report(y_test, prediction))
```



	precision	recall	f1-score	support
人类	1.00	0.01	0.01	370
技术	0.96	0.06	0.11	394
折腾	0.84	0.02	0.05	653
故事	0.57	0.73	0.64	1412
极客	1.00	0.10	0.19	543
科学	0.50	0.97	0.66	2136
脑洞	0.91	0.29	0.44	778
micro avg	0.54	0.54	0.54	6286
macro avg	0.83	0.31	0.30	6286
weighted avg	0.70	0.54	0.45	6286

	precision	recall	f1-score	support
人类	1.00	0.02	0.03	370
技术	0.92	0.03	0.05	394
折腾	0.93	0.02	0.04	653
故事	0.58	0.75	0.65	1411
极客	0.98	0.08	0.16	543
科学	0.50	0.97	0.67	2136
脑洞	0.92	0.30	0.45	778
micro avg	0.55	0.55	0.55	6285
macro avg	0.83	0.31	0.29	6285
weighted avg	0.71	0.55	0.45	6285

	precision	recall	f1-score	support
人类	1.00	0.02	0.03	370
技术	0.96	0.06	0.11	394
折腾	1.00	0.01	0.03	653

▼ One-Fold Validation for demonstrating

Save model to Google Drive

```
1 import time
2 X = df.context
3 y = df.label
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.5, random
5 NB_pipeline.fit(X_train, y_train.)
6 joblib.dump(NB_pipeline, "NB_model.m")
```

📁 ['NB_model.m']

训练集 1.000 0.01 0.03 653

Test for one fold

测试集 0.50 0.97 0.66 2136

```
1 X = df.context
2 y = df.label
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.5, random
4
5 NB_pipeline = joblib.load("NB_model.m")
6 prediction = NB_pipeline.predict(X_test)
7 print(metrics.classification_report(y_test, prediction))
```

📁

	precision	recall	f1-score	support
人类	1.00	0.03	0.06	27715
技术	0.96	0.06	0.11	29630
折腾	0.97	0.02	0.04	48908
故事	0.61	0.75	0.67	106037
极客	1.00	0.10	0.18	40656
科学	0.50	0.98	0.67	160056
脑洞	0.96	0.34	0.50	58238
micro avg	0.56	0.56	0.56	471240
macro avg	0.86	0.33	0.32	471240

▼ Support Vector Machine

In machine learning, support vector machines (SVMs, also support vector networks) are supervised learning models with associated learning algorithms that analyze data used for classification and regression analysis. Given a set of training examples, each marked as belonging to one or the other of two categories, an SVM training algorithm builds a model that assigns new examples to one category or the other, making it a non-probabilistic binary linear classifier (although methods such as Platt scaling exist to use SVM in a probabilistic classification setting). An SVM model is a representation of the examples as points in space, mapped so that the examples of the separate categories are divided by a clear gap that is as wide as possible. New examples are then mapped into that same space and predicted to belong to a category based on which side of the gap they fall.

In addition to performing linear classification, SVMs can efficiently perform a non-linear classification using what is called the kernel trick, implicitly mapping their inputs into high-dimensional feature spaces.

When data is unlabelled, supervised learning is not possible, and an unsupervised learning approach is required, which attempts to find natural clustering of the data to groups, and then map new data to these formed groups. The support vector clustering algorithm, created by Hava Siegelmann and Vladimir Vapnik, applies the statistics of support vectors, developed in the support vector machines algorithm, to categorize unlabeled data, and is one of the most widely used clustering algorithms in industrial applications

```
1 # Define a pipeline combining a text feature extractor with multi label classif
2 l_svm = Pipeline([
3     ('tfidf', TfidfVectorizer()),
4     ('clf', OneVsRestClassifier(LinearSVC()))
5 ])
```

▼ Test Performance of Basic Linear Kernel

```
1 for train_index, test_index in skf.split(X, y):
2     X_train, X_test = X[train_index], X[test_index]
3     y_train, y_test = y[train_index], y[test_index]
4     l_svm.fit(X_train, y_train)
5     prediction = l_svm.predict(X_test)
6     print(metrics.classification_report(y_test, prediction))
```



weighted avg	0.79	0.79	0.79	6286
	precision	recall	f1-score	support
人类	0.85	0.62	0.71	370
技术	0.71	0.60	0.65	394
折腾	0.68	0.63	0.65	653
故事	0.84	0.91	0.87	1411
极客	0.73	0.56	0.63	543
科学	0.80	0.87	0.83	2136
脑洞	0.82	0.83	0.83	778
micro avg	0.79	0.79	0.79	6285
macro avg	0.77	0.72	0.74	6285
weighted avg	0.79	0.79	0.79	6285

	precision	recall	f1-score	support
人类	0.89	0.69	0.78	370
技术	0.74	0.60	0.66	394
折腾	0.68	0.64	0.66	653
故事	0.85	0.91	0.88	1411
极客	0.70	0.55	0.62	542
科学	0.81	0.89	0.85	2136
脑洞	0.85	0.87	0.86	778
micro avg	0.80	0.80	0.80	6284
macro avg	0.79	0.73	0.76	6284
weighted avg	0.80	0.80	0.80	6284

	precision	recall	f1-score	support
人类	0.85	0.64	0.73	369
技术	0.71	0.64	0.67	393
折腾	0.67	0.63	0.65	653
故事	0.83	0.90	0.86	1411
极客	0.73	0.57	0.64	542
科学	0.81	0.88	0.84	2136
脑洞	0.84	0.86	0.85	777
micro avg	0.80	0.80	0.80	6281
macro avg	0.78	0.73	0.75	6281
weighted avg	0.79	0.80	0.79	6281

	precision	recall	f1-score	support
人类	0.88	0.69	0.78	369
技术	0.70	0.58	0.64	393
折腾	0.70	0.59	0.64	652
故事	0.83	0.91	0.87	1411
极客	0.68	0.54	0.60	542
科学	0.79	0.87	0.83	2136
脑洞	0.83	0.86	0.85	777

▼ Further More: Deep Learning Method!

weighted avg	0.79	0.79	0.79	6286
--------------	------	------	------	------

▼ Word Embedding as Feature Extraction

Text is considered a form of sequence data similar to time series data that you would have in weather data or financial data. In the previous BOW model, you have seen how to represent a whole sequence of words as a single feature vector. Now you will see how to represent each word as vectors. There are various ways to vectorize text, such as:

- Words represented by each word as a vector
- Characters represented by each character as a vector
- N-grams of words/characters represented as a vector (N-grams are overlapping groups of

▼ Preprocessing text with Keras Tokenizer

This class allows to vectorize a text corpus, by turning each text into either a sequence of integers (each integer being the index of a token in a dictionary) or into a vector where the coefficient for each token could be binary, based on word count, based on tf-idf...

```
1 import pandas as pd
2 from keras.preprocessing.text import Tokenizer
3 from keras.models import Sequential
4 from keras.layers import Dense
5 from keras.utils import to_categorical
6 from keras.models import load_model
7
8 df = pd.read_csv("Jiandan_segment.csv")
9 X = df.context
10 # convert Chinese to number
11 label_dict = {'人类': 0,
12               '技术': 1,
13               '折腾': 2,
14               '故事': 3,
15               '极客': 4,
16               '科学': 5,
17               '脑洞': 6,
18               }
19 y = df.label
20 for i, l in enumerate(y):
21     y[i] = label_dict[l]
22
23 print("check out numerize label:")
24 y_cate = to_categorical(y, num_classes=7)
25 print(y_cate[:10])
26 print("-----")
27
28
29 # Tokenize X_train data todo:change num_words
30 x_tokenizer = Tokenizer(num_words=5000)
31 x_tokenizer.fit_on_texts(X)
32 X_token = x_tokenizer.texts_to_sequences(X)
33 vocab_size = len(x_tokenizer.word_index) + 1
34
35 print("Examples of text after tokenized is: ", X_token[0])
```



check out previous labels

▼ Build basic MLP as Deep Learning Baseline

[0. 0. 0. 0. 0. 1. 0.1

```
1 !apt-get install python-pydot python-pydot-ng graphviz
2 !pip install pydot graphviz pydot-ng
```

```
↳ Reading package lists... Done
Building dependency tree
Reading state information... Done
graphviz is already the newest version (2.40.1-2).
Suggested packages:
  python-pyparsing-doc
The following NEW packages will be installed:
  python-pydot python-pydot-ng python-pyparsing
0 upgraded, 3 newly installed, 0 to remove and 8 not upgraded.
Need to get 91.5 kB of archives.
After this operation, 443 kB of additional disk space will be used.
Get:1 http://archive.ubuntu.com/ubuntu bionic/main amd64 python-pyparsing all
Get:2 http://archive.ubuntu.com/ubuntu bionic/universe amd64 python-pydot all
Get:3 http://archive.ubuntu.com/ubuntu bionic/universe amd64 python-pydot-ng
Fetched 91.5 kB in 1s (161 kB/s)
Selecting previously unselected package python-pyparsing.
(Reading database ... 110377 files and directories currently installed.)
Preparing to unpack .../python-pyparsing_2.2.0+dfsg1-2_all.deb ...
Unpacking python-pyparsing (2.2.0+dfsg1-2) ...
Selecting previously unselected package python-pydot.
Preparing to unpack .../python-pydot_1.2.3-1_all.deb ...
Unpacking python-pydot (1.2.3-1) ...
Selecting previously unselected package python-pydot-ng.
Preparing to unpack .../python-pydot-ng_1.0.0-3_all.deb ...
Unpacking python-pydot-ng (1.0.0-3) ...
Setting up python-pyparsing (2.2.0+dfsg1-2) ...
Setting up python-pydot-ng (1.0.0-3) ...
Setting up python-pydot (1.2.3-1) ...
Requirement already satisfied: pydot in /usr/local/lib/python3.6/dist-package
Requirement already satisfied: graphviz in /usr/local/lib/python3.6/dist-pack
Requirement already satisfied: pydot-ng in /usr/local/lib/python3.6/dist-pack
Requirement already satisfied: pyparsing>=2.1.4 in /usr/local/lib/python3.6/c
```

```
1 from keras.models import Sequential
2 from keras import layers
3 from keras.utils import plot_model
4 from IPython.display import SVG
5 from keras.utils.vis_utils import model_to_dot
6 from keras import regularizers
7 from keras.preprocessing.sequence import pad_sequences
8 import numpy as np
```

```
1 embedding_dim = 50
2
3 mlp = Sequential()
4 mlp.add(layers.Embedding(input_dim=vocab_size,
5                           output_dim=embedding_dim,
6                           input_length=400))
7 # mlp.add(layers.Flatten())
8 mlp.add(layers.GlobalMaxPool1D())
9 mlp.add(layers.Dense(128, activation='relu', kernel_regularizer = regularizers.
10 mlp.add(layers.Dropout(0.5))
11 mlp.add(layers.Dense(64, activation='relu', kernel_regularizer = regularizers.l
12 mlp.add(layers.Dropout(0.5))
13 mlp.add(layers.Dense(7, activation='softmax'))
14 mlp.compile(optimizer='adam',
```

```

15         loss='categorical_crossentropy',
16         metrics=['accuracy'],
17     )
18 mlp.summary()

```



Layer (type)	Output Shape	Param #
=====		
embedding_3 (Embedding)	(None, 400, 50)	4875050

global_max_pooling1d_3 (Glob	(None, 50)	0

dense_7 (Dense)	(None, 128)	6528

dropout_5 (Dropout)	(None, 128)	0

dense_8 (Dense)	(None, 64)	8256

dropout_6 (Dropout)	(None, 64)	0

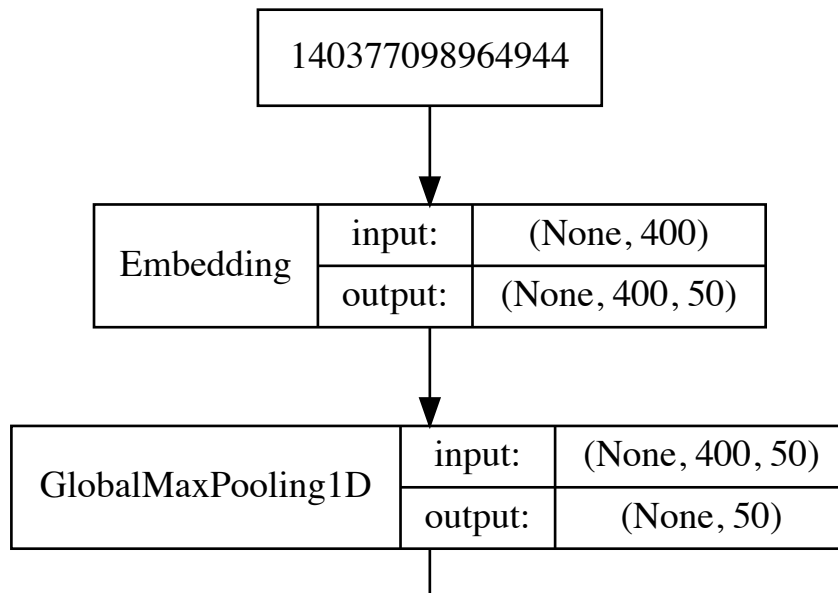
dense_9 (Dense)	(None, 7)	455
=====		
Total params: 4,890,289		
Trainable params: 4,890,289		
Non-trainable params: 0		

```

1 SVG(model_to_dot(mlp, show_shapes=True, show_layer_names=False).create(prog='do

```





▼ Measure performance of MLP

```

1 import matplotlib.pyplot as plt
2 plt.style.use('ggplot')
3
4 def plot_history(history):
5     acc = history.history['acc']
6     val_acc = history.history['val_acc']
7     loss = history.history['loss']
8     val_loss = history.history['val_loss']
9     x = range(1, len(acc) + 1)
10
11     plt.figure(figsize=(12, 5))
12     plt.subplot(1, 2, 1)
13     plt.plot(x, val_acc, 'r', label='Validation acc')
14     plt.title('Validation accuracy')
15     plt.legend()
16     plt.subplot(1, 2, 2)
17     plt.plot(x, val_loss, 'r', label='Validation loss')
18     plt.title('Validation loss')
19     plt.legend()
  
```

```

1 X_train, X_test, y_train, y_test = train_test_split(X_token, y_cate, test_size=
2
3 X_train, X_test = pad_sequences(X_train, padding='post', maxlen=400),\
4                               pad_sequences(X_test, padding='post', maxlen=400)
5
6 history = mlp.fit(X_train, y_train,
7                  epochs=12,
8                  verbose=0,
9                  validation_data=(X_test, y_test),
10                 batch_size=32)
11
12 mlp.save('MLP_model.h5')
13 loss, accuracy = mlp.evaluate(X_train, y_train, verbose=False)
14 print("Training Accuracy: {:.4f}".format(accuracy))
15 loss, accuracy = mlp.evaluate(X_test, y_test, verbose=False)
16 print("Testing Accuracy: {:.4f}".format(accuracy))
17 plot_history(history)
18
19
20 mlp_train = mlp.predict_classes(X_train)
21 mlp_pred = mlp.predict_classes(X_test)
22 print("Training results:")
23 print(metrics.classification_report(mlp_train, np.argmax(y_train, axis=1)))
24 print("Testing results:")
25 print(metrics.classification_report(mlp_pred, np.argmax(y_test, axis=1)))
26
  
```

☞ Training Accuracy: 0.9555

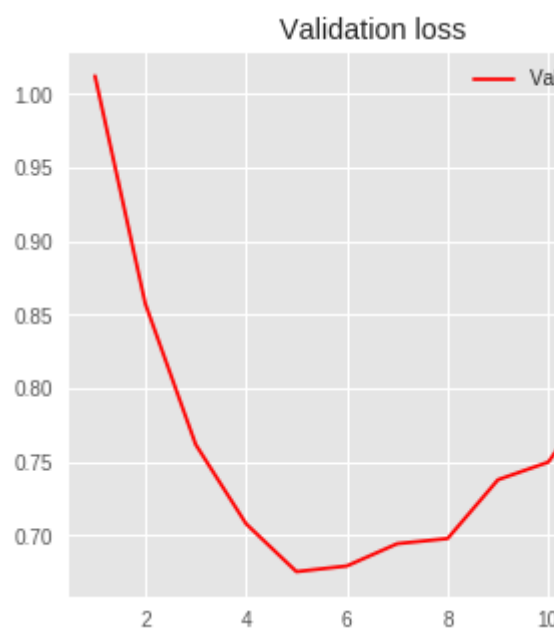
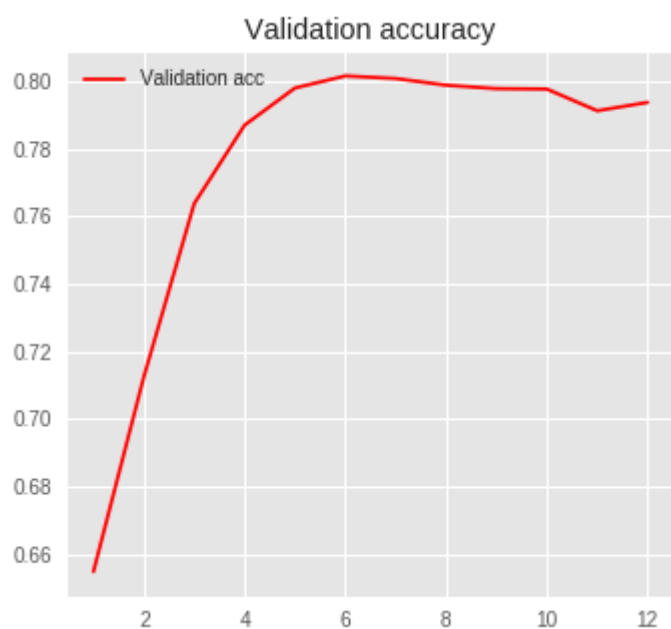
Testing Accuracy: 0.7937

Training results:

	precision	recall	f1-score	support
0	0.92	0.97	0.94	866
1	0.86	0.88	0.87	994
2	0.92	0.94	0.93	1594
3	0.99	0.98	0.99	3618
4	0.86	0.87	0.87	1342
5	0.99	0.97	0.98	5392
6	0.97	0.97	0.97	1902
micro avg	0.96	0.96	0.96	15708
macro avg	0.93	0.94	0.94	15708
weighted avg	0.96	0.96	0.96	15708

Testing results:

	precision	recall	f1-score	support
0	0.73	0.86	0.79	791
1	0.45	0.51	0.48	834
2	0.59	0.69	0.63	1397
3	0.95	0.91	0.93	3633
4	0.59	0.57	0.58	1406
5	0.87	0.81	0.84	5768
6	0.81	0.85	0.83	1879
micro avg	0.79	0.79	0.79	15708
macro avg	0.71	0.74	0.73	15708
weighted avg	0.80	0.79	0.80	15708



CNN: Convolutional Neural Network

In deep learning, a convolutional neural network (CNN, or ConvNet) is a class of deep neural networks, most commonly applied to analyzing visual imagery.

CNNs use a variation of multilayer perceptrons designed to require minimal preprocessing.[1] They are also known as shift invariant or space invariant artificial neural networks (SIANN), based on their

shared-weights architecture and translation invariance characteristics.[2][3]

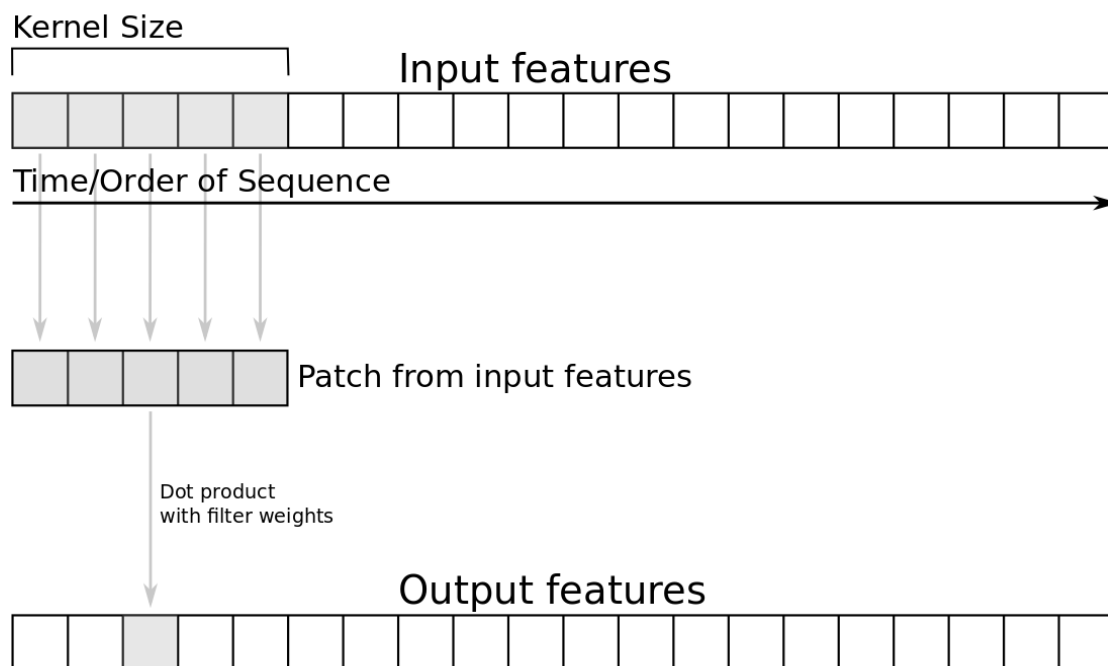
Convolutional networks were inspired by biological processes[4] in that the connectivity pattern between neurons resembles the organization of the animal visual cortex. Individual cortical neurons respond to stimuli only in a restricted region of the visual field known as the receptive field. The receptive fields of different neurons partially overlap such that they cover the entire visual field.

CNNs use relatively little pre-processing compared to other image classification algorithms. This means that the network learns the filters that in traditional algorithms were hand-engineered. This independence from prior knowledge and human effort in feature design is a major advantage.

They have applications in image and video recognition, recommender systems,[5] image classification, medical image analysis, and natural language processing.

▼ Conv1D

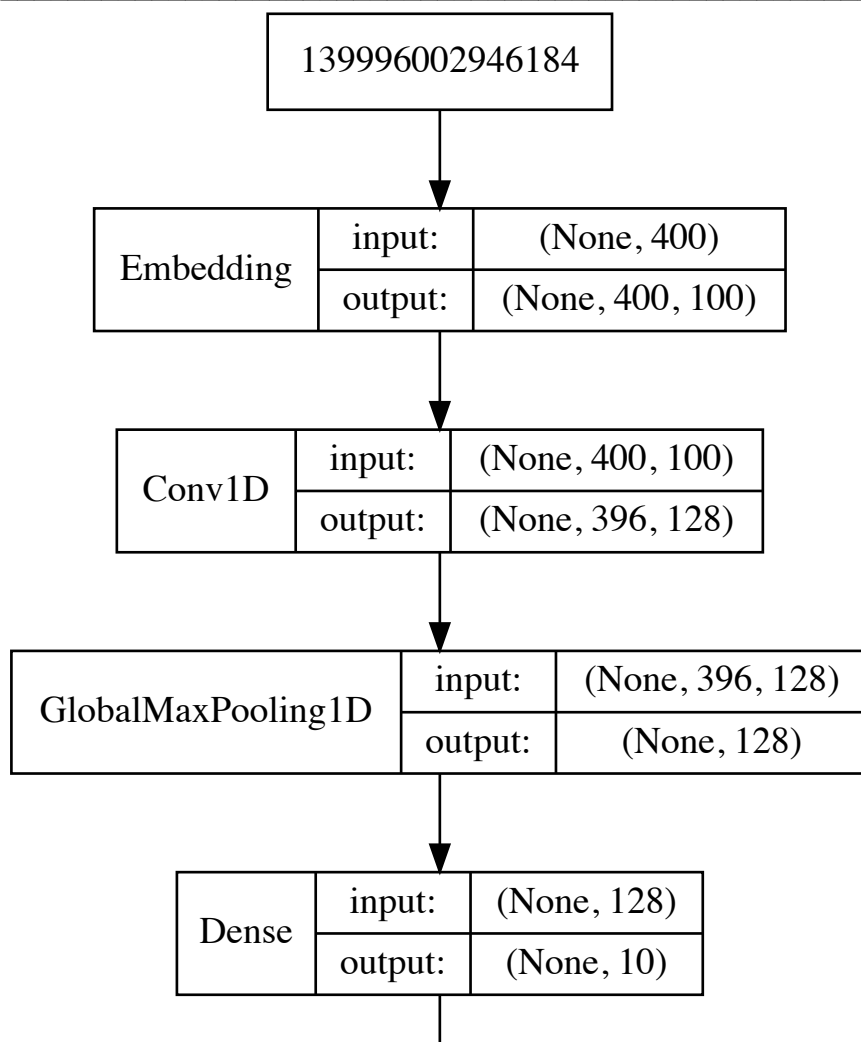
Keras offers again various Convolutional layers which you can use for this task. The layer you'll need is the Conv1D layer. This layer has again various parameters to choose from.



```
1 embedding_dim = 100
2
3 Conv = Sequential()
4 Conv.add(layers.Embedding(vocab_size, embedding_dim, input_length=400))
5 Conv.add(layers.Conv1D(128, 5, activation='relu'))
6 Conv.add(layers.GlobalMaxPooling1D())
7 Conv.add(layers.Dense(10, activation='relu'))
8 Conv.add(layers.Dense(7, activation='softmax'))
9 Conv.compile(optimizer='adam',
10              loss='categorical_crossentropy',
11              metrics=['accuracy'])
12 Conv.summary()
13 SVG(model_to_dot(Conv, show_shapes=True, show_layer_names=False).create(prog='d
```



Layer (type)	Output Shape	Param #
embedding_3 (Embedding)	(None, 400, 100)	9750100
conv1d_4 (Conv1D)	(None, 396, 128)	64128
global_max_pooling1d_1 (Glob	(None, 128)	0
dense_6 (Dense)	(None, 10)	1290
dense_7 (Dense)	(None, 7)	77
Total params: 9,815,595		
Trainable params: 9,815,595		
Non-trainable params: 0		



▼ Measure the performance of basic CNN

```

1 from keras.preprocessing.sequence import pad_sequences
2 import numpy as np
3
4 # X_train, X_test, y_train, y_test = train_test_split(X_token, y_cate, test_siz
5
6 # X_train, X_test = pad_sequences(X_train, padding='post', maxlen=400),\
7 #                               pad_sequences(X_test, padding='post', maxlen=400)
8
9 history = Conv.fit(X_train, y_train,
10                   epochs=15,
11                   verbose=0,

```



```
12         validation_data=(X_test, y_test),
13         batch_size=32)
14
15 Conv.save('MLP_model.h5')
16 loss, accuracy = Conv.evaluate(X_train, y_train, verbose=False)
17 print("Training Accuracy: {:.4f}".format(accuracy))
18 loss, accuracy = Conv.evaluate(X_test, y_test, verbose=False)
19 print("Testing Accuracy: {:.4f}".format(accuracy))
20 plot_history(history)
21
22
23 Conv_train = Conv.predict_classes(X_train)
24 Conv_pred = Conv.predict_classes(X_test)
25 print("Training results:")
26 print(metrics.classification_report(Conv_train, np.argmax(y_train, axis=1)))
27 print("Testing results:")
28 print(metrics.classification_report(Conv_pred, np.argmax(y_test, axis=1)))
29
```



Training Accuracy: 0.9999

Testing Accuracy: 0.8516

Training results:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	912

▼ Deeper CNN

Let us check out whether a deeper CNN can improve the performance of text classification.

	5	1.00	1.00	1.00	5292
--	---	------	------	------	------

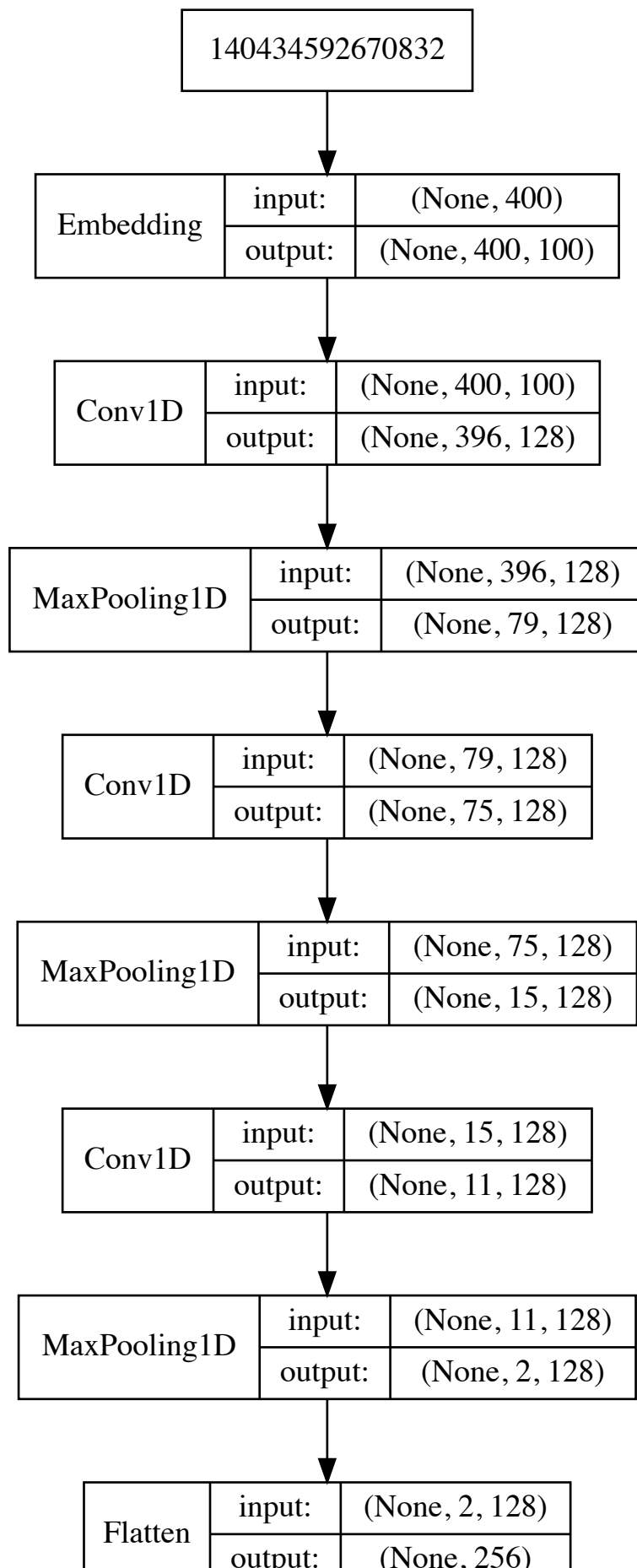
```
1 embedding_dim = 100
2
3 D_Conv = Sequential()
4 D_Conv.add(layers.Embedding(vocab_size, embedding_dim, input_length=400))
5 D_Conv.add(layers.Conv1D(128, 5, activation='relu'))
6 D_Conv.add(layers.MaxPooling1D(5))
7 D_Conv.add(layers.Conv1D(128, 5, activation='relu'))
8 D_Conv.add(layers.MaxPooling1D(5))
9 D_Conv.add(layers.Conv1D(128, 5, activation='relu'))
10 D_Conv.add(layers.MaxPooling1D(5))
11 D_Conv.add(layers.Flatten())
12 D_Conv.add(layers.Dense(32, activation='relu'))
13 D_Conv.add(layers.Dense(7, activation='softmax'))
14 D_Conv.compile(optimizer='adam',
15               loss='categorical_crossentropy',
16               metrics=['accuracy'])
17 D_Conv.summary()
```



Layer (type)	Output Shape	Param #
=====		
embedding_1 (Embedding)	(None, 400, 100)	9750100
<hr/>		
conv1d_1 (Conv1D)	(None, 396, 128)	64128
<hr/>		
max_pooling1d_1 (MaxPooling1D)	(None, 79, 128)	0
<hr/>		
conv1d_2 (Conv1D)	(None, 75, 128)	82048
<hr/>		
max_pooling1d_2 (MaxPooling1D)	(None, 15, 128)	0
<hr/>		
conv1d_3 (Conv1D)	(None, 11, 128)	82048
<hr/>		
max_pooling1d_3 (MaxPooling1D)	(None, 2, 128)	0
<hr/>		
flatten_1 (Flatten)	(None, 256)	0
<hr/>		
dense_1 (Dense)	(None, 32)	8224
<hr/>		
dense_2 (Dense)	(None, 7)	231
=====		
Total params: 9,986,779		
Trainable params: 9,986,779		
Non-trainable params: 0		
<hr/>		

```
1 SVG(model_to_dot(D_Conv, show_shapes=True, show_layer_names=False).create(prog=
```





```
1 from keras.preprocessing.sequence import pad_sequences
2 import numpy as np
3
4 history = D_Conv.fit(X_train, y_train,
5                     epochs=1,
6                     verbose=0,
7                     validation_data=(X_test, y_test),
```

```
8 | batch_size=32)
```

```
1 D_conv = load_model('Dconv_model.h5')
2 loss, accuracy = D_Conv.evaluate(X_train, y_train, verbose=False)
3 print("Training Accuracy: {:.4f}".format(accuracy))
4 loss, accuracy = D_Conv.evaluate(X_test, y_test, verbose=False)
5 print("Testing Accuracy: {:.4f}".format(accuracy))
6 # plot_history(history)
7
8
9 D_Conv_train = D_Conv.predict_classes(X_train)
10 D_Conv_pred = D_Conv.predict_classes(X_test)
11 print("Training results:")
12 print(metrics.classification_report(D_Conv_train, np.argmax(y_train, axis=1)))
13 print("Testing results:")
14 print(metrics.classification_report(D_Conv_pred, np.argmax(y_test, axis=1)))
```

☞ Training Accuracy: 0.9221

Testing Accuracy: 0.8240

Training results:

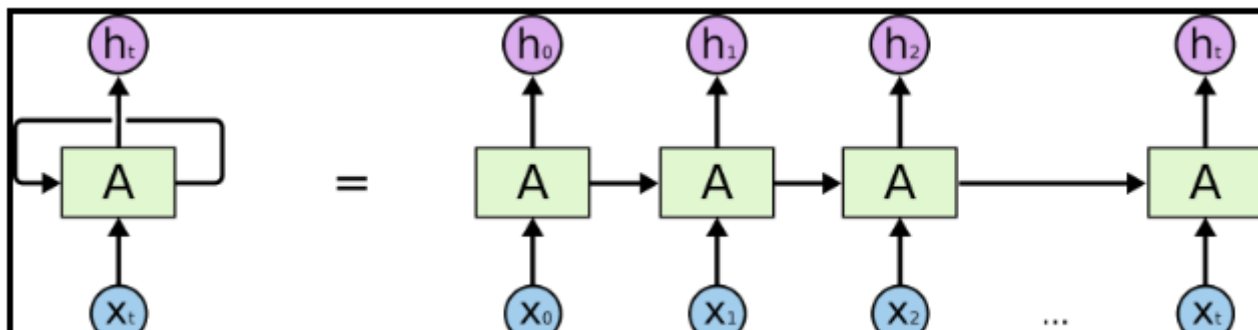
	precision	recall	f1-score	support
0	0.90	0.71	0.79	1154
1	0.77	0.72	0.74	1092
2	0.93	0.92	0.92	1640
3	0.98	1.00	0.99	3526
4	0.64	0.82	0.72	1056
5	0.98	0.96	0.97	5355
6	0.95	0.97	0.96	1885
micro avg	0.92	0.92	0.92	15708
macro avg	0.88	0.87	0.87	15708
weighted avg	0.93	0.92	0.92	15708

Testing results:

	precision	recall	f1-score	support
0	0.77	0.50	0.61	1444
1	0.56	0.49	0.52	1091
2	0.74	0.73	0.73	1665
3	0.98	1.00	0.99	3399
4	0.38	0.57	0.46	901
5	0.91	0.91	0.91	5424
6	0.84	0.93	0.88	1784
micro avg	0.82	0.82	0.82	15708
macro avg	0.74	0.73	0.73	15708
weighted avg	0.83	0.82	0.82	15708

▼ RNN

A recurrent neural network (RNN) is a class of artificial neural network where connections between nodes form a directed graph along a sequence. This allows it to exhibit temporal dynamic behavior for a time sequence. Unlike feedforward neural networks, RNNs can use their internal state (memory) to process sequences of inputs. This makes them applicable to tasks such as unsegmented, connected handwriting recognition or speech recognition.



```

1 LSTM = Sequential()
2 LSTM.add(layers.Embedding(vocab_size, 100, input_length=400))
3 LSTM.add(layers.Bidirectional(layers.CuDNNLSTM(100)))
4 LSTM.add(layers.Dense(32, activation='relu'))
5 LSTM.add(layers.Dense(7, activation='softmax'))
6
7 LSTM.compile(optimizer='adam',
8               loss='categorical_crossentropy',
9               metrics=['accuracy'])
10 LSTM.summary()

```



Layer (type)	Output Shape	Param #
=====		
embedding_2 (Embedding)	(None, 400, 100)	9750100

bidirectional_1 (Bidirection	(None, 200)	161600

dense_4 (Dense)	(None, 32)	6432

dense_5 (Dense)	(None, 7)	231
=====		
Total params: 9,918,363		
Trainable params: 9,918,363		
Non-trainable params: 0		

```

1 SVG(model_to_dot(LSTM, show_shapes=True, show_layer_names=False).create(prog='d

```



139996281477216

input: (None, 400)

```
1 from keras.preprocessing.sequence import pad_sequences
2 import numpy as np
3
4 X_train, X_test, y_train, y_test = train_test_split(X_token, y_cate, test_size=
5
6 X_train, X_test = pad_sequences(X_train, padding='post', maxlen=400),\
7                               pad_sequences(X_test, padding='post', maxlen=400)
8
9 history = LSTM.fit(X_train, y_train,
10                   epochs=10,
11                   verbose=0,
12                   validation_data=(X_test, y_test),
13                   batch_size=128)
14
15 loss, accuracy = LSTM.evaluate(X_train, y_train, verbose=False)
16 print("Training Accuracy: {:.4f}".format(accuracy))
17 loss, accuracy = LSTM.evaluate(X_test, y_test, verbose=False)
18 print("Testing Accuracy: {:.4f}".format(accuracy))
19 plot_history(history)
20
21
22 LSTM_train = LSTM.predict_classes(X_train)
23 LSTM_pred = LSTM.predict_classes(X_test)
24 print("Training results:")
25 print(metrics.classification_report(LSTM_train, np.argmax(y_train, axis=1)))
26 print("Testing results:")
27 print(metrics.classification_report(LSTM_pred, np.argmax(y_test, axis=1)))
28
```



Training Accuracy: 0.9999

Testing Accuracy: 0.9004

Training results:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	912
1	1.00	1.00	1.00	1014
2	1.00	1.00	1.00	1632
3	1.00	1.00	1.00	3585
4	1.00	1.00	1.00	1353
5	1.00	1.00	1.00	5292
6	1.00	1.00	1.00	1920
micro avg	1.00	1.00	1.00	15708
macro avg	1.00	1.00	1.00	15708
weighted avg	1.00	1.00	1.00	15708

Testing results:

	precision	recall	f1-score	support
0	0.86	0.84	0.85	968
1	0.60	0.68	0.64	851
2	0.91	0.89	0.90	1669
3	0.99	0.99	0.99	3499
4	0.66	0.72	0.69	1239
5	0.94	0.93	0.94	5456
6	0.94	0.91	0.92	2026
micro avg	0.90	0.90	0.90	15708
macro avg	0.84	0.85	0.85	15708
weighted avg	0.90	0.90	0.90	15708

