

به نام خداوند بخشنده مهربان

پروژه پایانی

یادگیری تحت نظارت

اعضای گروه :

امیر محمد آقاجانی - ۴۰۰۵۲۱۰۶۳

سید محمد علی میرمحمدی - ۴۰۰۵۲۲۲۶۵

بخش اول : توضیح روند کلی

در ابتدا به توضیح کلی الگوریتم در حالت سه کلاس می پردازیم. در حالت سه کلاس هر کلمه در هر توئیت با توجه به برچسب توئیت شمرده می شود یعنی برای مثال داریم که کلمه not ۱۰ بار در کلاس مثبت، ۱۰۰ بار در کلاس منفی و ۵ بار در کلاس خنثی تکرار شده است. حال با استفاده از naïve bayes (بدون محاسبه آلفا که در داک آورده شده است) طبق فرمول احتمال پیشین (احتمال پیشین منظور این است که به طور برای کلاس خنثی بدون در نظر گرفتن هیچ گونه ویژگی صرفاً با توجه به تعداد توئیت های خنثی و تعداد کل توئیت ها چقدر احتمال دارد که این توئیت خنثی باشد) ضرب در احتمال رخ دادن این کلمه به شرط احتمال پیشین به ازای همه کلمات موجود در توئیت میتوانیم احتمال این که این توئیت از جنس احتمال پیشین باشد را نشان دهد. در واقع به طور مثال احتمال منفی بودن توئیت با توجه به تعداد تکرار کلمات آن در کلاس منفی ها و تعداد تکرار آن ها در کل را محاسبه کرده و برای دو کلاس دیگر نیز این احتمال را محاسبه میکنیم. حال هر احتمالی که بیشتر باشد را به عنوان پیش بینی برگردانیم.

بخش دوم : توضیح کد

ابتدا از فایل template شروع میکنیم. پس از کانستراکتور تابع train قرار دارد که به ازای هر توئیت تعداد توئیت ها و به ازای کلمات آنها تعداد کلمات متناسب با label هر توئیت را update میکند و این کار را به ازای تمام توئیت ها انجام میدهد. تابع calculate_prior به ازای هر کلاس صرفاً بر اساس تعداد توئیت های آن کلاس احتمال رخداد آن کلاس را بررسی میکند. تابع calculate_likelihood بر اساس تعداد کلمات در هر کلاس احتمال رخداد آن کلمه به شرط اینکه توئیت مورد نظر از کلاس خاصی باشد را محاسبه میکند.

تابع classify نیز در مواجهه با هر توئیت به ازای تمام کلاس‌ها احتمال اینکه این توئیت متعلق به آن کلاس‌ها باشد را محاسبه کرده و احتمال ماکزیمم را به عنوان خروجی برمی‌گرداند.

```
def train(self, data):
    # training process:
    # inputs: data(list) --> each item of list is a tuple
    # the first index of the tuple is a list of words and the second index is the label(positive, negative, or neutral)
    self.tweet_counts = len(data)
    for features, label in data:
        self.class_counts[label] += 1
        for feature in features:
            self.class_word_counts[label] += 1
            if (feature, label) in self.vocab.keys():
                self.vocab[(feature, label)] += 1
            else:
                self.vocab[(feature, label)] = 1

def calculate_prior(self, label):
    # calculate log prior
    # you can add some attributes to this method

    return log10(self.class_counts[label] / self.tweet_counts)

def calculate_likelihood(self, word, label):
    # calculate likelihood: P(word | label)
    # return the corresponding value

    if (word, label) in self.vocab.keys():
        return log10((self.vocab[(word, label)] + 1) / (self.class_word_counts[label] + 3))
    return log10(1 / (self.class_word_counts[label] + 3))

def classify(self, features):
    # predict the class
    # inputs: features(list) --> words of a tweet

    what_class = {'negative': 0, 'neutral': 0, 'positive': 0}

    for label in self.classes:
        for feature in features:
            what_class[label] += self.calculate_likelihood(feature, label)
            what_class[label] += self.calculate_prior(label)

    best_class = max(what_class, key=what_class.get)

    return best_class
```

حال به بررسی فایل run می‌پردازیم. در ابتدا تابع preprocess یک رشته (متن توئیت) را دریافت کرده آن را تمیز کرده و کاراکترهای اضافی را از آن حذف می‌کنیم. سپس کلماتی که معنی مثبت یا منفی ندارند مثل حروف ربط را حذف کرده و خروجی را به صورت لیستی از کلمات برمی‌گردانیم.

تابع load_data داده را از فایل csv خوانده و در نهایت به صورت لیستی از tuple های دوتایی (کلمات توئیت و label) برمی‌گرداند.

تابع load_test_data همین کار را برای داده های تست انجام می‌دهد. سپس در بخش بعد train انجام شده و زمان آن ثبت می‌شود و در بخش دوم validation انجام شده و دقت چاپ می‌شود و در بخش آخر نیز label گذاری صورت می‌گیرد.

```

def preprocess(tweet_string):
    # clean the data and tokenize it
    str = re.sub(r"[&!\"#$%&()*+,-./:;<=>?@[\\]^_{}~\n -' 0123456789\\]", "", tweet_string)
    str = unicode.decode(str)
    strs = str.split(' ')
    features = []
    stemmer = PorterStemmer()
    features = list(stemmer.stem(word) for word in strs)
    features = list(filter(None, features))
    li = ["to", "of", "as", "a", "an", "with", "for", "and", "or", "my", "your", "his", "her", "ours", "their", "me", "him", "us", "them", ""]
    final_features = []
    for word in features:
        if not (word in li or word.isnumeric()) and len(word) > 1:
            final_features.append(word)
    return final_features

def load_data(data_path):
    # load the train csv file and return the data
    grass_data = pd.read_csv(data_path)
    data = []
    for tweet in grass_data.iterrows():
        data.append((preprocess(str(tweet[1].text)), str(tweet[1].label_text)))
    return data

def load_test_data(data_path):
    grass_data = pd.read_csv(data_path)
    data = []
    for tweet in grass_data.iterrows():
        data.append(preprocess(str(tweet[1].text)))
    return data

```

```

# train your model and report the duration time
start = time()
train_data_path = "train_data.csv"
eval_data_path = "eval_data.csv"
test_data_path = "test_data_nolabel.csv"
classes = ['negative', 'neutral', 'positive']
nb_classifier = NaiveBayesClassifier(classes)
nb_classifier.train(load_data(train_data_path))
end = time()
print(f"training duration time:", round(end - start, 2), "second")

# check on evaluation data
eval_data = load_data(eval_data_path)
all_tweets = 0
correct_guess = 0
for eval_tweet, label in eval_data:
    all_tweets += 1
    guess_label = nb_classifier.classify(eval_tweet)
    if guess_label == label:
        correct_guess += 1
print(f"correctness percent on evaluation data: {round(correct_guess / all_tweets * 100, 2)}%")

# test on test data:
test_data = load_test_data(test_data_path)
number = 0
f = open("labels.txt", "a")
for test_tweet in test_data:
    guess_label = nb_classifier.classify(test_tweet)
    output = str(number) + ": " + guess_label + "\n"
    f.write(output)
    number += 1
f.close()

```

پایان