



باسمه تعالی



دانشگاه تهران - دانشکده مهندسی برق و کامپیوتر
درس: پردازش زبان طبیعی

گزارش تمرین کامپیوتری دوم

استاد: جناب آقای دکتر هشام فیلی

دستیار استاد: جناب باقرشاهی

دانشجو: مصطفی امیری

شماره دانشجویی: ۸۱۰۱۰۱۰۰۲

فهرست مطالب

۳	۱ سوال اول
۳	۱.۱ آماده سازی داده‌ها
۵	۲.۱ tf-idf
۷	۱.۲.۱ نتایج
۷	۳.۱ PPMI
۹	۱.۳.۱ نتایج
۱۱	۲ سوال دوم
۱۱	۱.۲ پیش پردازش
۱۴	۲.۲ آموزش مدل
۱۴	۱.۲.۲ تابع هزینه
۱۶	۲.۲.۲ ایجاد مدل و محاسبه گرادیان
۱۷	۳.۲ نتایج
۱۹	۳ سوال سوم
۱۹	۱.۳ آماده سازی بردار کلمات
۱۹	۲.۳ مدل
۲۱	۳.۳ داده‌ها
۲۲	۴.۳ آموزش و تست
۲۴	۵.۳ مقایسه با مدل بیز

فهرست تصاویر

۴	توزیع برچسب داده‌های آموزش و تست	۱
۷	نتایج تحلیل احساسی بر اساس بردار tf-idf	۲
۹	نتایج تحلیل احساسی بر اساس PPMI	۳
۱۲	نمودار فراوانی ۲۰۰ توکن پرتکرار	۴
۱۸	نمایش بردارهای تفاضل کلمات	۵
۲۳	نمودار تغییر تابع هزینه و صحت عملکرد مدل در حین فرایند آموزش	۶
۲۳	نتایج مدل logistic regression	۷
۲۴	نتایج مدل naive bayes	۸

۱ سوال اول

۱.۱ آماده سازی داده‌ها

ابتدا فایل داده‌ها را دریافت و در یک دیتافریم قرار دادم:

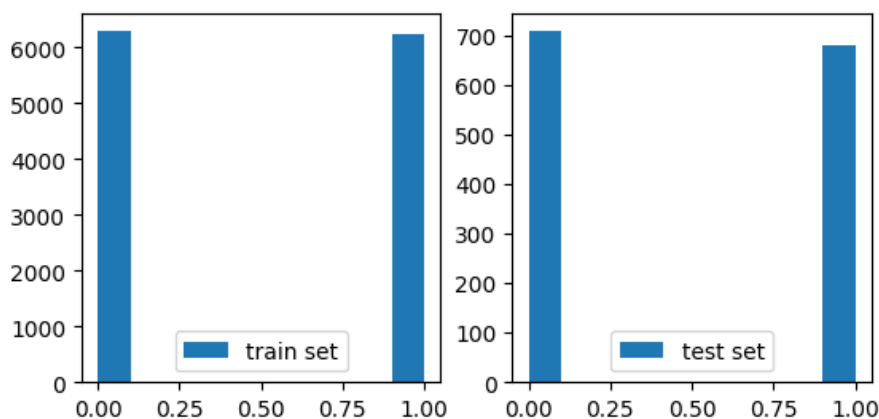
```
۱ snappfood = pd.read_csv(  
۲     "datasets\\Snapfood\\Snappfood - Sentiment Analysis.csv",  
۳     encoding="utf-8",  
۴     on_bad_lines="skip",  
۵     sep="\t")  
۶ snappfood = snappfood.drop('Unnamed: 0', axis=1)  
۷ snappfood.head()
```

سپس ۲۰ درصد از داده‌ها را انتخاب کردم و ۱۰ درصد از این مقدار را برای تست و بقیه را برای آموزش قرار دادم:

```
۱ train_ds, _ = train_test_split(snappfood, test_size=0.8)  
۲ train_ds, test_ds = train_test_split(train_ds, test_size=0.1)
```

سپس با استفاده از قطعه کد زیر نمودار هیستوگرام توزیع داده‌های آموزش و تست را برای اطمینان از مشابهت توزیع داده‌ها رسم کردم که نمودار آن به صورت شکل ۱ شد.

```
۱ fig, axes = plt.subplots(ncols=2)  
۲ fig.set_figheight(3)  
۳ axes[0].hist(train_ds["label_id"], label="train set")  
۴ axes[1].hist(test_ds['label_id'], label="test set")  
۵ axes[0].legend()  
۶ axes[1].legend()
```



شکل ۱: توزیع برچسب داده‌های آموزش و تست

برای پیش پردازش داده‌ها با استفاده از عبارات منظم حروفی که مانند «ی» در کلمه خییییییییییییییلی چندبار تکرار شده‌اند را حذف و با استفاده از کتابخانه هضم آنها را نرمال کردم و سپس روی داده‌ها این تابع را اعمال کردم. در ادامه نیز نمونه‌هایی از جملات پیش پردازش شده را در خروجی چاپ کردم.

```

1 def preproces(text):
2     regex = r'(\w)\1{2,}'
3     text = re.sub(regex, r"\1", text, 0, re.MULTILINE)
4     text = Normalizer().normalize(text)
5     return text
6
7 # applying preprocess function to data
8 train_ds["processed"] = train_ds['comment'].apply(preproces)
9 train_ds["tokens"] = train_ds['processed'].apply(WordTokenizer().tokenize)
10
11 test_ds["processed"] = test_ds['comment'].apply(preproces)
12 test_ds["tokens"] = test_ds['processed'].apply(WordTokenizer().tokenize)
13
14
15 # printing some example of processed data in console
16 for id in np.random.choice(train_ds.index, 20):
17     print(train_ds.loc[id]['comment'])
18     print(train_ds.loc[id]['processed'])

```

برای نمونه جمله:

● داغ و تازه عالی

به جمله زیر تبدیل می‌شود:

• داغ و تازه عالی

۲.۱ tf-idf

سپس برای هر کلمه یک بردار به اندازه تعداد نظرات ایجاد کردم و در درایه متناظر با آن نظر تعداد تکرار آن کلمه در آن نظر را قرار دادم. به این ترتیب بردار $term - frequency(tf)$ برای هر کلمه ساخته شد.

```

۱ vocabs_tf={}
۲ for id, doc in tqdm(enumerate(train_ds['tokens'])):
۳     for t in doc:
۴         if t not in stopwords_list():
۵             if t not in vocabs_tf:
۶                 vocabs_tf[t] = np.array([0]*len(train_ds['processed']))
۷                 vocabs_tf[t][id] += 1

```

با استفاده از قطعه کد زیر تعداد تکرار هر کلمه در نظرات مختلف و عبارت دیگر $document frequency$ را محاسبه کردم:

```

۱ vocabs_df = {}
۲ for v in vocabs_tf:
۳     vocabs_df[v] = 0
۴ for v in tqdm(vocabs_tf):
۵     vocabs_df[v] = sum([(int)(x>0) for x in vocabs_tf[v]])

```

مطابق روابط مربوط به محاسبه $tf - idf$ که در **فصل ششم** کتاب ژورافسکی آمده است بردار مربوط به هر کلمه را ساختم.

$$tf_{t,d} = \log_{10}(\text{count}(t, d) + 1)$$

$$idf_t = \log\left(\frac{N}{df_t}\right)$$

$$\Rightarrow w_{t,d} = tf_{t,d} \times idf_t$$

```

۱ tfidf = dict()
۲ N = len(train_ds)

```

```

۳ for t in vocabs_tf:
۴     tfidf[t] = np.log(vocabs_tf[t]+1)/np.log(10) * np.log(N/vocabs_df[t])/np.log
(10)

```

سپس با استفاده از تابع زیر یک متن را طبق رابطه $d = \frac{w_1+w_2+\dots+w_k}{k}$ به یک بردار تبدیل می‌کنم.

```

۱ def tfidf_vectorizer(tokens):
۲     count = 0
۳     vec = []
۴     for t in tokens:
۵         if t in tfidf:
۶             count+=1
۷             if len(vec)>0:
۸                 vec += tfidf[t]
۹             else:
۱0                vec = np.copy(tfidf[t])
۱۱     vec = np.array(vec)
۱۲     if count != 0:
۱۳         vec = vec/count
۱۴     else:
۱۵         for t in tfidf:
۱۶             vec = np.array([0] * len(tfidf[t]))
۱۷             break
۱۸     return vec

```

سپس تابع فوق را روی داده‌ها اعمال کرده و بعد از تبدیل داده‌ها به بردار یک طبقه‌بند بیز ساده را روی داده‌های آموزش، آموزش داده و روی داده‌های تست نتایج را بدست آوردم.

```

۱ test_ds["vec"] = test_ds["tokens"].apply(tfidf_vectorizer)
۲ train_ds["vec"] = train_ds["tokens"].apply(tfidf_vectorizer)
۳
۴ train_dss = train_ds.dropna()
۵ test_dss = test_ds.dropna()
۶
۷ clf = GaussianNB()
۸ clf.fit(np.array([np.array(x) for x in train_dss["vec"]]), train_dss["label_id"]
    ])
۹ prd = clf.predict(np.array([np.array(x) for x in test_dss["vec"]]))
۱۰

```

```
۱۱ print(classification_report(test_dss["label_id"], prd))
```

۱.۲.۱ نتایج

نتایج این قسمت روی داده‌های تست به صورت شکل ۲ است.

	precision	recall	f1-score	support
0.0	0.86	0.61	0.71	665
1.0	0.72	0.91	0.80	723
accuracy			0.76	1388
macro avg	0.79	0.76	0.76	1388
weighted avg	0.78	0.76	0.76	1388

شکل ۲: نتایج تحلیل احساسی بر اساس بردار tf-idf

۳.۱ PPMI

به ازاء هر دو کلمه تعداد حضور همزمان آن‌ها در یک نظر بوسیله قطعه کد زیر شمرده شده است:

```
۱ coaccurance = dict()
۲ for i in vocabs_tf:
۳     coaccurance[i] = dict()
۴     for j in vocabs_tf:
۵         coaccurance[i][j] = 0
۶ for doc in tqdm(train_ds["tokens"]):
۷     for comb in combinations(doc, 2):
۸         if comb[0] in vocabs_tf and comb[1] in vocabs_tf:
۹             coaccurance[comb[0]][comb[1]] += 1
۱۰            coaccurance[comb[1]][comb[0]] += 1
```

سپس تعداد مشاهده هر کلمه به تنهایی نیز محاسبه شده است.

```
۱ for i in tqdm(vocabs_tf):
۲     count = 0
۳     for j in vocabs_tf:
۴         count += coaccurance[i][j]
۵     coaccurance[i]["total"] = count
۶
۷ total = 0
```



```

۸ coaccurance["total"] = dict()
۹ for j in tqdm(vocabs_tf):
۱۰     count = 0
۱۱     for i in vocabs_tf:
۱۲         count += coaccurance[i][j]
۱۳     total += count
۱۴     coaccurance["total"][j] = count
۱۵ coaccurance["total"]["total"] = total

```

مقدار $PPMI$ برای هر دو کلمه نسبت به هم محاسبه و سپس با استفاده از این مقادیر بردار هر کلمه را ساخته‌ام:

```

۱ ppmi = dict()
۲ for i in tqdm(vocabs_tf):
۳     ppmi[i] = dict()
۴     for j in vocabs_tf:
۵         p_i_j = coaccurance[i][j] / coaccurance["total"]["total"]
۶         p_i = coaccurance[i]["total"] / coaccurance["total"]["total"]
۷         p_j = coaccurance["total"][j] / coaccurance["total"]["total"]
۸         ppmi[i][j] = 0
۹         if p_i != 0 and p_j != 0 and p_i_j != 0:
۱۰             ppmi[i][j] = max(np.log(p_i_j)-np.log(p_i)-np.log(p_j), 0)
۱۱
۱۲
۱۳ ppmi_vec = dict()
۱۴ for i in tqdm(vocabs_tf):
۱۵     ppmi_vec[i] = []
۱۶     for j in vocabs_tf:
۱۷         ppmi_vec[i].append(ppmi[i][j])

```

تابع زیر یک جمله را مشابه آنچه در بخش قبلی گفته شد به یک بردار تبدیل می‌کند:

```

۱ def ppmi_vectorizer(tokens):
۲     count = 0
۳     vec = []
۴     for t in tokens:
۵         if t in ppmi_vec:
۶             count+=1

```

```

۷         if len(vec)>0:
۸             vec += ppmi_vec[t]
۹         else:
۱۰             vec = np.copy(ppmi_vec[t])
۱۱     vec = np.array(vec)
۱۲     if count != 0:
۱۳         vec = vec/count
۱۴     else:
۱۵         for t in ppmi_vec:
۱۶             vec = np.array([0] * len(ppmi_vec[t]))
۱۷         break
۱۸     return vec

```

سپس تابع فوق را روی داده‌ها اعمال کرده و بعد از تبدیل داده‌ها به بردار یک طبقه‌بند بیز ساده را روی داده‌های آموزش، آموزش داده و روی داده‌های تست نتایج را بدست آوردم.

```

۱ test_dss["ppmi_vec"] = test_dss["tokens"].apply(ppmi_vectorizer)
۲ train_dss["ppmi_vec"] = train_dss["tokens"].apply(ppmi_vectorizer)
۳
۴ clf = GaussianNB()
۵ clf.fit(np.array([np.array(x) for x in train_dss["ppmi_vec"]]), train_dss["
    label_id"])
۶ prd = clf.predict(np.array([np.array(x) for x in test_dss["ppmi_vec"]]))
۷
۸ print(classification_report(test_dss["label_id"], prd))

```

۱.۳.۱ نتایج

نتایج حاصل روی داده‌های تست به صورت شکل ۳ است.

	precision	recall	f1-score	support
0.0	0.87	0.59	0.70	665
1.0	0.71	0.92	0.80	723
accuracy			0.76	1388
macro avg	0.79	0.75	0.75	1388
weighted avg	0.79	0.76	0.75	1388

شکل ۳: نتایج تحلیل احساسی بر اساس PPMI

نتایج دو روش تفاوت زیادی با هم ندارند و به میزان بسیار کمی روش *tfidf* نتایج بهتری داشته است.

۲ سوال دوم

۱.۲ پیش‌پردازش

ابتدا با قطعه کد زیر فایل داده شده را خوانده و در متغیر *doc* ذخیره کردم. ضمناً توضیحات اول فایل هم که از نوشته‌های شکسپیر نیست را حذف کردم.

```
1 with open("t8.shakespeare.txt", "r") as f:
2
3     doc = ""
4
5     for line in f.readlines():
6
7         doc += line
8
9
10 doc = doc[10520:]
```

در مرحله بعدی توضیحاتی که در بین ستاره یا << قرار داشتند را با استفاده از *regex* حذف کردم.

```
1 regex = r"\s+<[\s\d\-\(\)\A-Z_\.,:@;!\\"'+>+\s+"
2 result = re.sub(regex, r"", doc, 0, re.MULTILINE)
3 regex = r"\s+\*+[\s\d\-\(\)\A-Za-z_\.,:@!;\\"'+]+\*+\s+"
4 result = re.sub(regex, r"", result, 0, re.MULTILINE)
```

سپس فاصله‌های اضافه و اعداد و تکرارهای بیش از ۳ بار یک حرف را حذف کردم.

```
1 regex = r"(\s)\s+"
2 result = re.sub(regex, r"\1", result, 0, re.MULTILINE)
3 regex = r"\d+"
4 result = re.sub(regex, "", result, 0, re.MULTILINE)
5 regex = r'(\w)\1{2,}'
6 result = re.sub(regex, "", result, 0, re.MULTILINE)
```

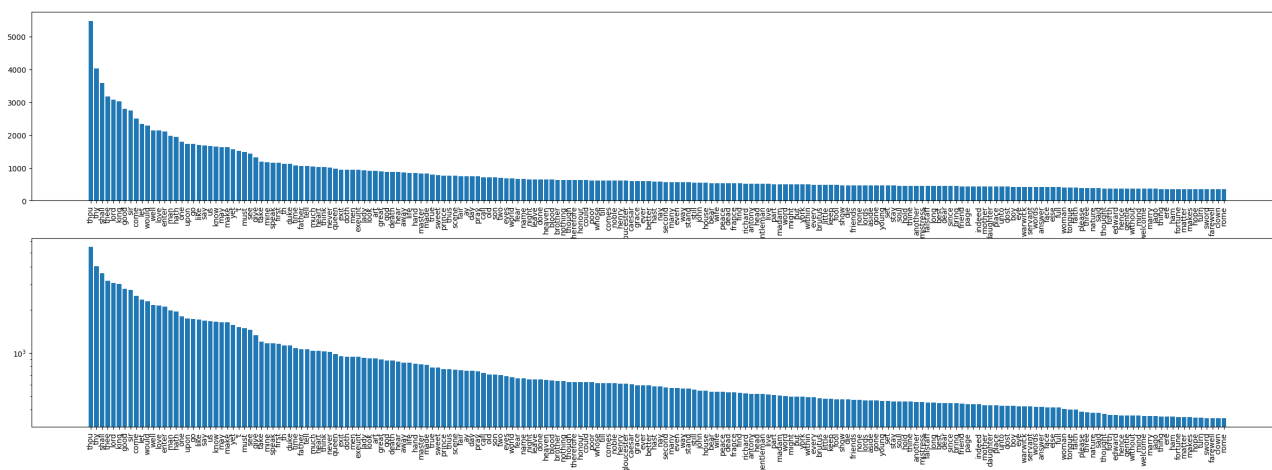
در مرحله بعد جملات و توکن‌های داخل هر جمله را از هم جدا کردم و برخی توکن‌های بدون استفاده را حذف کردم.

```
1 sentences = sent_tokenize(result)
2 sentences = [word_tokenize(s.lower()) for s in sentences]
3 sentences = [[(re.sub(r"[!,:;\?\\.\-\\[\]\(\)]", "", t)) for t in s if t not in ["'d",
4     "ll", "s", "!", ":", "?", ";", ",", ".", "[", "]", "(", "(", "(", "'"] +
5     stopwords_list] for s in sentences]
```

در مرحله بعد تعداد تکرار هر توکن را شمرده و بر اساس فراوانی آنها را مرتب کردم:

```
۱ from collections import Counter
۲ freqs = Counter()
۳ for s in sentences:
۴     freqs.update(s)
۵ freqs_sorted = sorted(freqs.items(),
۶                        key=lambda p: (p[1], p[0]),
۷                        reverse=True)
```

در شکل ۴ نمودار فراوانی ۲۰۰ کلمه پرتکرار به صورت عادی و با مقیاس لگاریتمی رسم کردم



شکل ۴: نمودار فراوانی ۲۰۰ توکن پرتکرار

با در نظر گرفتن ۲۳۰۰۰ توکن فراوانی مابقی توکن‌ها کمتر از ۱ درصد کل خواهد بود.

```
۱ vocab_size = 23000
۲ sum_freq_pruned = sum(f for _, f in freqs_sorted[vocab_size:])
۳ print(sum_freq_pruned/sum(f for _, f in freqs_sorted)*100)
۴ # 0.9895637457690861
```

برای سایر توکن‌ها یک توکن *UNK* در نظر گرفتم و در ابتدای این دیکشنری اضافه کردم و فراوانی آن را نیز برابر فراوانی توکن‌های حذف شده در نظر گرفتم.

```
۱ freqs_sorted = [("UNK", sum_freq_pruned)] + freqs_sorted[:vocab_size-1]
```

با توجه به توضیحات کتاب ژورافسکی برای انتخاب نمونه‌های منفی به نسبت فراوانی توکن‌ها با توان ۰/۷۵ استفاده شده مطابق رابطه زیر که با قطعه کد زیر این نسب محاسبه و سپس برای تبدیل به اعداد صحیح تمام آنها

را مقیاس کرده و در یک جدول دیگر به تعداد مقدار مقیاس شده از آن توکن درج کردم.

$$P_{\alpha}(w) = \frac{count(w)^{\alpha}}{\sum_{w'} count(w')^{\alpha}}$$

```

۱ ns_table = {}
۲ sum_freq = 0
۳ for w, freq in freqs_sorted:
۴     ns_freq = freq ** 0.75
۵     ns_table[w] = ns_freq
۶     sum_freq += ns_freq
۷
۸ scaler = (1<<24) / sum_freq
۹ ns_table = [(w, freq, int(round(ns_table[w]*scaler))) for w, freq in
    freqs_sorted]
۱۰
۱۱ ns_table_expanded = []
۱۲ for i, (_, _, count) in enumerate(ns_table):
۱۳     ns_table_expanded.extend([i] * count)
۱۴ neg_ix = np.random.randint(len(ns_table_expanded), size=1000)
۱۵ _ = plt.hist([ns_table_expanded[k] for k in neg_ix], bins=100)

```

سپس یک دیکشنری برای تبدیل توکن‌ها به آیدی عددی ایجاد کردم.

```

۱ voc = { w:i for i, (w, _, _) in enumerate(ns_table) }

```

تابع زیر نیز وظیفه تولید داده‌های آموزشی را به عهده دارد. ساختار داده‌های تولیدی به این صورت است که یک پنجره را روی جمله حرکت می‌دهیم و به ازاء هر کلمه، هریک از کلمات مجاور آن و داخل پنجره را به عنوان بافتار و به تعداد ۴ نمونه منفی طبق توزیعی که در جدول تولید شده بدست آمده است انتخاب می‌کنیم و سپس آیدی آن کلمه جزء بافتار و آیدی این ۴ کلمه منفی را در یک آرایه قرار می‌دهیم. در نتیجه این تابع بسته‌ای از داده‌ها شامل یک داده هدف و آرایه‌ای که عضو اول آن کلمه مثبت و بقیه کلمات منفی هستند برمی‌گرداند. ضمناً توجه داریم که برچسب اولین داده ۱ و بقیه صفر است و بنابراین در خروجی این تابع برچسب را برمی‌گردانیم.

```

۱ def get_data(batch_size, sentences, window_size, ns_size=5):
۲     target = []
۳     context = []
۴

```

```

۵ for s in sentences:
۶     ps = np.random.random(size=len(s))
۷     encoded = [voc.get(t, 0) for t in s]
۸     for i in range(len(encoded)):
۹         start = max(0, i-window_size)
۱۰        end = min(i+window_size+1, len(encoded))
۱۱        for j in range(start, end):
۱۲            if j != i:
۱۳                target.append(encoded[i])
۱۴                context.append(np.append( np.array([encoded[j]]), ns_table_expanded[
np.random.randint(len(ns_table_expanded), size=ns_size)] ))
۱۵            if len(target) == batch_size:
۱۶                yield target, context
۱۷                target = []
۱۸                context = []
۱۹ yield target, context

```

۲.۲ آموزش مدل

ابتدا توابع کمکی زیر را برای محاسبه تابع هزینه و گرادیان تعریف کردم:

```

۱ def sigmoid(x):
۲     return 1/(1+np.exp(-x))
۳
۴ def log_sigmoid(x):
۵     return np.log(1/(1+np.exp(-x)))

```

۱.۲.۲ تابع هزینه

رابطه مربوط به تابع هزینه به صورت زیر است:

$$L_{CE} = - \left[\log \sigma(c_{pos}.w) + \sum_{i=1}^k \log \sigma(-c_{neg_i}.w) \right]$$

برای پیاده سازی این تابع هزینه با توجه به ساختار داده‌ها و برجسبها ابتدا مقدار $(-1)^{(1-y)}$ را محاسبه کردم. به این ترتیب حاصل برای داده اول ۱ و برای بقیه منفی ۱ خواهد بود. سپس این بردار را تغییر شکل داده و به اندازه ابعاد امبدینگ مدنظر مقادیر را برای ضرب کردن در بردار مربوط به

توکن‌ها تکرار کردم. بنابراین برای یک بسته از برجسب‌های داده مراحل انجام کار به صورت روابط زیر خواهد بود:

$$y = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ \dots \\ 1 & 0 & 0 & 0 & 0 \end{pmatrix} \Rightarrow$$

$$(-1)^{(1-y)} = \begin{pmatrix} 1 & -1 & -1 & -1 & -1 \\ 1 & -1 & -1 & -1 & -1 \\ \dots \\ 1 & -1 & -1 & -1 & -1 \end{pmatrix}$$

$$\begin{pmatrix} \begin{pmatrix} (1) & (-1) & (-1) & (-1) \end{pmatrix} \\ \begin{pmatrix} (1) & (-1) & (-1) & (-1) \end{pmatrix} \\ \dots \\ \begin{pmatrix} (1) & (-1) & (-1) & (-1) \end{pmatrix} \end{pmatrix}$$

$$\begin{pmatrix} \begin{pmatrix} (1 \ 1 \ \dots \ 1) & (-1 \ -1 \ \dots \ -1) & (-1 \ -1 \ \dots \ -1) & (-1 \ -1 \ \dots \ -1) \end{pmatrix} \\ \begin{pmatrix} (1 \ 1 \ \dots \ 1) & (-1 \ -1 \ \dots \ -1) & (-1 \ -1 \ \dots \ -1) & (-1 \ -1 \ \dots \ -1) \end{pmatrix} \\ \dots \\ \begin{pmatrix} (1 \ 1 \ \dots \ 1) & (-1 \ -1 \ \dots \ -1) & (-1 \ -1 \ \dots \ -1) & (-1 \ -1 \ \dots \ -1) \end{pmatrix} \end{pmatrix}$$

این کار با استفاده از کد زیر انجام می‌شود:

```
np.power(-1, 1-y).reshape(batch_size,1+neg_sample, 1).repeat(embed_size,axis=-1)
```

در مرحله بعد با ضرب این ماتریس در امبدینگ مربوط به $[c_{pos}, c_{neg_1}, \dots, c_{neg_k}]$ و ضرب اینها در بردار مربوط به امبدینگ کلمه هدف مقادیر مورد نظر در تابع هزینه تولید و با اعمال تابع $\log(\text{sigmoid})$ مقدار هزینه بدست می‌آید:

```
def loss(t, c, y, batch_size):
    n_c = np.power(-1, 1-y).reshape(batch_size,1+neg_sample, 1).repeat(embed_size,
        axis=-1)*c.reshape(batch_size,(1+neg_sample),embed_size)
```



```

۳ return -np.sum(log_sigmoid(np.sum(t.repeat(1+neg_sample, axis=0).reshape(
    batch_size,(1+neg_sample),embed_size)*n_c,axis=-1)),axis=-1)

```

۲.۲.۲ ایجاد مدل و محاسبه گرادیان

با استفاده از کد زیر مقادیر اولیه و ماتریس‌ها مدل مقدار دهی اولیه می‌شود.

```

۱ epochs = 10
۲ batch_size = 128
۳ embed_size = 100
۴ neg_sample = 4
۵ lr = 0.001
۶ T = np.random.random(size=(len(ns_table), embed_size))
۷ C = np.random.random(size=(len(ns_table), embed_size))

```

روابط مربوط به محاسبه گرادیان به صورت زیر است:

$$\begin{aligned}
 \frac{\partial L_{CE}}{\partial c_{pos}} &= (\sigma(c_{pos}.w) - 1)w \\
 \frac{\partial L_{CE}}{\partial c_{neg}} &= \sigma(c_{neg}.w)w \\
 \Rightarrow \frac{\partial L_{CE}}{\partial c} &= (\sigma(c.w) - y)w \\
 \frac{\partial L_{CE}}{\partial w} &= [\sigma(c_{pos}.w) - 1]c_{pos} + \sum_{i=1}^k [\sigma(c_{neg_i}.w)]c_{neg_i} \\
 \Rightarrow \frac{\partial L_{CE}}{\partial w} &= \sum [\sigma(c.w) - y]c
 \end{aligned}$$

این مقادیر طبق روابط زیر محاسبه و سپس وزن‌ها بر اساس این گرادیان‌ها به روز می‌شوند:

```

۱ c_diff = (sigmoid(np.sum(T[t].repeat(1+neg_sample, axis=0).reshape(
    epoch_batch_size,(1+neg_sample),embed_size)*r_c,axis=-1)) - y).reshape(
    epoch_batch_size,1+neg_sample,1).repeat(embed_size,axis=-1)*T[t].repeat(1+
    neg_sample,axis=0).reshape(epoch_batch_size,1+neg_sample,embed_size)
۲
۳ t_diff = np.sum((sigmoid(np.sum(T[t].repeat(1+neg_sample, axis=0).reshape(
    epoch_batch_size,(1+neg_sample),embed_size)*r_c,axis=-1)) - y).reshape(
    epoch_batch_size,1+neg_sample,1).repeat(embed_size,axis=-1)*r_c, axis=1)

```

بنابراین فرآیند آموزش به صورت زیر خواهد بود:

```

۱ for i in range(epochs):
۲     print("epoch {:}.format(i+1), end=" ")
۳     epoch_loss = 0
۴     bT = np.copy(T)
۵     bC = np.copy(C)
۶     for t, c in tqdm(get_data(batch_size, sentences, 5, neg_sample)):
۷         epoch_batch_size = len(t)
۸         y= np.array([[1]+[0]*neg_sample]*epoch_batch_size)
۹         #if epoch_batch_size == batch_size:
۱۰        epoch_loss = (epoch_loss * i + np.mean(loss(T[t], C[c], y, epoch_batch_size)
11        ))/(i+1)
12        r_c = C[c].reshape(epoch_batch_size,(1+neg_sample),embed_size)
13        c_diff = (sigmoid(np.sum(T[t].repeat(1+neg_sample, axis=0).reshape(
epoch_batch_size,(1+neg_sample),embed_size)*r_c,axis=-1)) - y).reshape(
epoch_batch_size,1+neg_sample,1).repeat(embed_size,axis=-1)*T[t].repeat(1+
neg_sample,axis=0).reshape(epoch_batch_size,1+neg_sample,embed_size)
14        t_diff = np.sum((sigmoid(np.sum(T[t].repeat(1+neg_sample, axis=0).reshape(
epoch_batch_size,(1+neg_sample),embed_size)*r_c,axis=-1)) - y).reshape(
epoch_batch_size,1+neg_sample,1).repeat(embed_size,axis=-1)*r_c, axis=1)
15        C[c] -= lr*c_diff.reshape(epoch_batch_size,(1+neg_sample),embed_size)
16        T[t] -= lr*t_diff
17    print("loss: {:.3f}".format(epoch_loss))
18    if bLoss < epoch_loss:
19        print("early stop")
20        T = bT
21        C = bC
22        break
23    bLoss = epoch_loss

```

۳.۲ نتایج

در مرحله بعد دو ماتریس را با هم جمع کرده و بردارهای گفته شده در سوال را رسم کردم.

```

۱ V = T+C
۲ from sklearn.decomposition import PCA
۳ pca = PCA(n_components=2, random_state=12)
۴ n_V = pca.fit_transform(V)

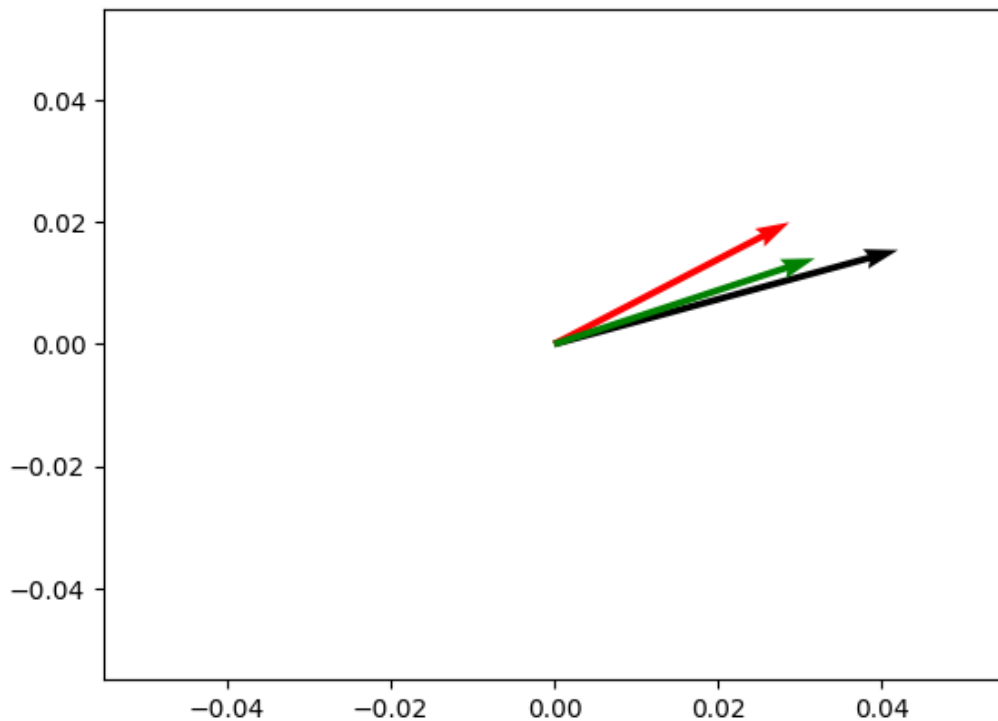
```

```

۵ km = n_V[[voc.get("king", 0)]] - n_V[[voc.get("man", 0)]]
۶ qw = n_V[[voc.get("queen", 0)]] - n_V[[voc.get("woman", 0)]]
۷ bs = n_V[[voc.get("brother", 0)]] - n_V[[voc.get("sister", 0)]]
۸ data = np.array([km[0], qw[0], bs[0]])
۹ origin = np.array([[0, 0, 0], [0, 0, 0]])
۱۰ plt.quiver(*origin, data[:, 0], data[:, 1], color=['black', 'red', 'green'],
    scale=1.5)

```

بردار سیاه تفاضل $king - man$ و بردار قرمز بردار $queen - woman$ و بردار سبز بردار $brother - sister$ را نشان می‌دهد.



شکل ۵: نمایش بردارهای تفاضل کلمات

علاوه بر ضمیمه کردن کدها، تمام موارد در **کولب** هم قابل دسترس است.

۳ سوال سوم

۱.۳ آماده سازی بردار کلمات

ابتدا فایل وزنهای *glove* را دریافت و آنها را برای هر کلمه با کد زیر به یک بردار تبدیل کردم. ضمناً در انتهای تمام بردارها مقدار ۱ را نیز برای محاسبه بایاس اضافه کردم.

```

۱ embed_size = 300
۲ glove = defaultdict()
۳ with open("glove.6B."+str(embed_size)+"d.txt", "r") as f:
۴     for line in f.readlines():
۵         l = line.split()
۶         glove[l[0]] = np.array([(float)(x) for x in l[1:]]+[1]) # 1 is for bias

```

۲.۳ مدل

ابتدا توابع کمکی زیر را برای محاسبه تابع هزینه و صحت عملکرد مدل تعریف کردم.

```

۱ def softmax(x):
۲     return np.exp(x)/(np.sum(np.exp(x), axis=-1).repeat(x.shape[-1], axis=-1).
        reshape(x.shape))
۳
۴ def loss(y_true, y_pred):
۵     return np.mean(np.sum(-y_true*np.log(y_pred),axis=-1))
۶
۷ def accuracy(y_true, y_pred):
۸     return sum([(int)(p) for p in np.argmax(y_pred, axis=-1) == np.argmax(y_true,
        axis=-1)])/len(y_true)

```

تابع هزینه مورد استفاده تابع هزینه *cross entropy* است. رابطه مربوط به این تابع هزینه به صورت زیر است:

$$L_{CE} = - \sum y_i \log(\hat{y}_i)$$

با توجه به اینکه یک بعد به بردار ویژگیها برای بایاس اضافه کردیم، روابط ریاضی مربوط به مدل به صورت زیر خواهد بود:

$$\mathbf{x} = (x_1, x_2, \dots, x_{emd}, 1)$$

$$\mathbf{z} = \mathbf{x} \times \mathbf{w}^T$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z})$$

برای محاسبه گرادیان تابع هزینه نسبت به وزنهای \mathbf{w} به صورت زیر عمل می‌کنیم:

$$\frac{\partial L_{CE}}{\partial \mathbf{w}} = \frac{\partial L_{CE}}{\partial \hat{\mathbf{y}}} \times \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}} \times \frac{\partial \mathbf{z}}{\partial \mathbf{w}}$$

با محاسبه این روابط و ساده کردن آنها به رابطه زیر می‌رسیم:

$$\frac{\partial L_{CE}}{\partial \mathbf{w}} = (\hat{\mathbf{y}} - \mathbf{y}) \times \mathbf{x}$$

با توجه به این روابط مدل به صورت زیر تعریف شده است:

```

1 class logisticmodel():
2     def __init__(self, n_features, n_categories):
3         self.W = np.random.random((n_categories, n_features+1))
4
5     def forward(self, X):
6         self.X = X
7         self.pred = softmax((X@self.W.T))
8         return self.pred
9
10    def backpropagate(self, y_true, lr):
11        self.W = self.W - lr * (self.pred - y_true).T@self.X
12
13    def train(self, X, y, X_val, y_val, batch_size, epochs, lr):
14        history = {"train_loss": [], "val_loss": [], "train_acc": [], "val_acc": []}
15        for i in range(epochs):
16            for b in range((int)(len(X)/batch_size)):
17                pred = self.forward(X[b*batch_size:min((b+1)*batch_size, len(X))])
18                self.backpropagate(y[b*batch_size:min((b+1)*batch_size, len(X))], lr)
19            if X_val:
20                pred = self.forward(X_val)
21                history["val_loss"].append(loss(y_val, pred))

```

```

۲۰     history["val_acc"].append(accuracy(y_val, pred))
۲۱     pred = self.forward(X)
۲۲     history["train_loss"].append(loss(y, pred))
۲۳     history["train_acc"].append(accuracy(y, pred))
۲۴
۲۵     if (i+1)%100 == 0 or i == 0:
۲۶         if X_val:
۲۷             print("epoch {:4} loss is: {:.3f} - accuracy is: {:.3f} - val_loss is:
{::.3f} - val_acc is: {:.3f}".format(i+1, history["train_loss"][-1], history[
"train_acc"][-1]*100, history["val_loss"][-1], history["val_acc"][-1]*100))
۲۸         else:
۲۹             print("epoch {:4} loss is: {:.3f} - accuracy is: {:.3f}".format(i+1,
history["train_loss"][-1], history["train_acc"][-1]*100))
۳۰     return history
۳۱ def test(self, X, y):
۳۲     pred = self.forward(X)
۳۳     print("=====")
۳۴     print("                        Test Report                        ")
۳۵     print("=====")
۳۶     print(classification_report(np.argmax(y_test, axis=-1), np.argmax(pred, axis
=-1)))

```

۳.۳ داده‌ها

با استفاده از قطعه کد زیر داده‌ها دریافت به سه قسمت آموزش، اعتبارسنجی و تست تقسیم شدند. ضمناً برچسب‌ها نیز در این مرحله به *one hot* تبدیل شدند.

```

۱ data= pd.read_csv("all-data.csv", on_bad_lines="skip", encoding="ISO-8859-1",
header=None)
۲
۳
۴ X_train, X_test, y_train, y_test = train_test_split(data[1], data[0], test_size
=0.1)
۵ X_train, X_val, y_train, y_val = train_test_split(X_train, y_train, test_size
=0.1)
۶
۷ onehot = OneHotEncoder()

```

```

۸ onehot.fit(np.array(y_train).reshape(-1,1))
۹ y_train = onehot.transform(np.array(y_train).reshape(-1,1)).toarray()
۱۰ y_val = onehot.transform(np.array(y_val).reshape(-1,1)).toarray()
۱۱ y_test = onehot.transform(np.array(y_test).reshape(-1,1)).toarray()

```

با استفاده از قطعه کد زیر و مشابه کاری که در سوال اول انجام شد، هر کدام از داده‌ها به یک بردار تبدیل شد.

```

۱ def sentence2vec(sentence):
۲     vec = []
۳     stop = nltk.corpus.stopwords.words('english')
۴     for w in word_tokenize(sentence.lower()):
۵         v = glove.get(w, glove.get('unk'))
۶         vec.append(v)
۷     vec = np.array(vec)
۸     return np.mean(vec, axis=0)
۹
۱۰ X_train = [sentence2vec(x) for x in X_train]
۱۱ X_val = [sentence2vec(x) for x in X_val]
۱۲ X_test = [sentence2vec(x) for x in X_test]

```

۴.۳ آموزش و تست

با استفاده از قطعه کد زیر مدل آموزش داده شد که نمودار تابع هزینه و صحت مدل نیز در شکل ۶ نمایش داده شده است.

```

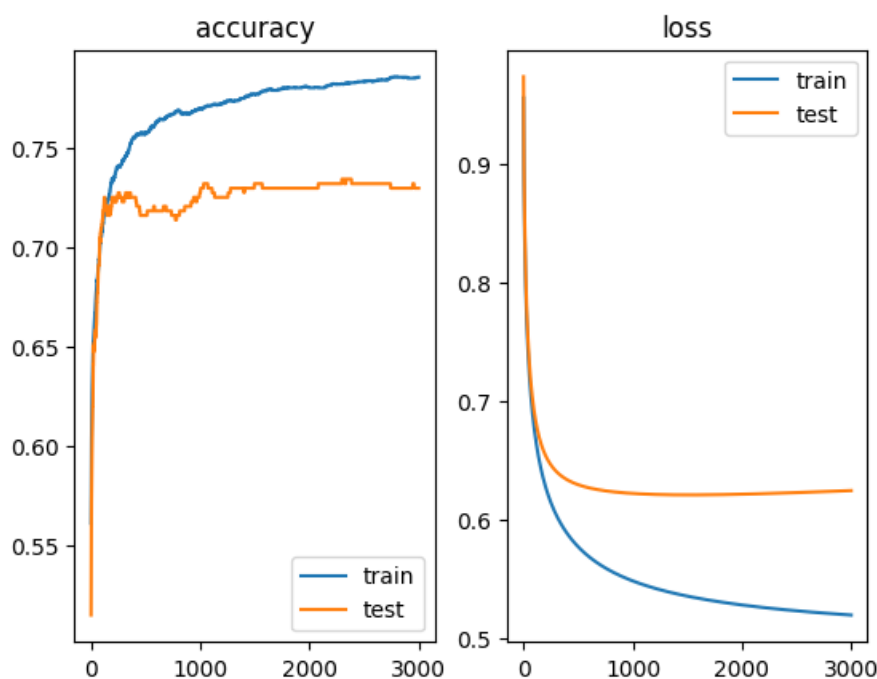
۱ batch_size = 64
۲ categories = 3
۳ epochs = 3000
۴ lr = 5e-4
۵
۶ model = logisticmodel(embed_size, categories)
۷ history = model.train(X_train, y_train, X_val, y_val, batch_size, epochs, lr)
۸ fig, ax = plt.subplots(ncols=2)
۹ ax[0].plot(history["train_acc"], label="train")
۱۰ ax[0].plot(history["val_acc"], label="test")
۱۱ ax[0].legend()
۱۲ ax[0].set_title("accuracy")
۱۳ ax[1].plot(history["train_loss"], label="train")

```

```

۱۴ ax[1].plot(history["val_loss"], label="test")
۱۵ ax[1].legend()
۱۶ ax[1].set_title("loss")
۱۷ plt.show()
۱۸ model.test(X_test, y_test)

```



شکل ۶: نمودار تغییر تابع هزینه و صحت عملکرد مدل در حین فرایند آموزش

نهایتاً نتایج روی داده‌های تست به صورت شکل ۷ شد.

Test Report				
	precision	recall	f1-score	support
0	0.64	0.52	0.57	62
1	0.74	0.82	0.78	283
2	0.58	0.51	0.54	140
accuracy			0.69	485
macro avg	0.65	0.61	0.63	485
weighted avg	0.68	0.69	0.68	485

شکل ۷: نتایج مدل logistic regression

همانطور که نتایج هم پیداست عملکرد روی داده‌هایی که نمونه بیشتر دارند بهتر و روی داده‌های با نمونه‌های کمتر ضعیف‌تر است.

۵.۳ مقایسه با مدل بیز

با استفاده از قطعه کد زیر همین کار با استفاده از مدل بیز هم انجام شد و نتایج در شکل ۸ نشان داده شده است.

```
۱ from sklearn.naive_bayes import GaussianNB
۲ clf = GaussianNB()
۳ clf.fit(X_train, [np.argmax(y) for y in y_train])
۴ print(classification_report(np.argmax(y_test, axis=-1), clf.predict(X_test)))
```

	precision	recall	f1-score	support
0	0.29	0.47	0.36	62
1	0.70	0.84	0.76	283
2	0.43	0.14	0.21	140
accuracy			0.59	485
macro avg	0.47	0.48	0.44	485
weighted avg	0.57	0.59	0.55	485

شکل ۸: نتایج مدل naive bayes

همانطور که مشخص است نتایج حاصل از مدل بیز ضعیفتر است که به نظر می‌رسد با توجه به غیر متوازن بودن داده‌ها و فرض استقلال ویژگی‌ها و ابعاد بالای بردار ویژگی این عملکرد ضعیفتر منطقی باشد.

علاوه بر ضمیمه کردن کدها، تمام موارد در **کولب** هم قابل دسترس است.

کلیه موارد مربوط به این تمرین در **گیت هاب** نیز در دسترس است.