



دانشگاه صنعتی امیرکبیر  
(پلیتکنیک تهران)  
دانشکده مهندسی کامپیوتر

## درس بیوانفورماتیک پروژه پایانی

امیرمهدی زرین نژاد

۹۷۳۱۰۸۷

گزارش کامل همراه با کد در نوت بوک نیز نوشته شده و ضمیمه شده است.

## - ابتدا بخش پیش پردازش را داریم:

در قدم اول توابع مربوط به **kmer** را بررسی می کنیم که جلوتر از آن ها استفاده می شود:

```
[58] def kmer_for_one_sequence(seq, k):  
    number_of_windows = len(seq) - k + 1  
    kmers = {}  
  
    alphabet = ['A', 'C', 'G', 'T']  
    products = [''.join(p) for p in itertools.product(alphabet, repeat=k)]  
    for product in products:  
        record = ''.join(product)  
        kmers[record] = 0  
    # print(kmers)  
  
    for i in range(number_of_windows):  
        current_window = seq[i:i + k]  
        # print('window:', current_window)  
        kmers[current_window] += 1  
  
    if number_of_windows:  
        for record in kmers:  
            kmers[record] = kmers[record] / number_of_windows  
    return kmers
```

این تابع برای اعمال به یک رشته پیاده سازی شده است.

حلقه **for** اول:

به این صورت کار می کند که ابتدا **k** یا همان سایز پنجره را همراه با رشته مورد نظر دریافت می کند؛ با توجه با الفبایی که برای رشته ها داریم یک دیکشنری از تمام حالت های **k** تایی الفبا ایجاد می کند و این دیکشنری را با صفر مقداردهی اولیه می کند.

حلقه **for** دوم:

سپس در رشته پیمایش می کند و در هر مرحله پنجره **k** تایی موجود و پیش رو را در نظر می گیرد. و برای این پنجره در دیکشنری، ۱ واحد اضافه می کند(در واقع با این کار تعداد هر کدام از پنجره های ممکن را در رشته ی داده شده می شمارد)

حلقه **for** سوم:

نهایتا در حلقه سوم مقادیر دیکشنری، تقسیم بر تعداد کل پنجره های ممکن می شوند و به این صورت نرمال سازی می شوند و اثر طول خنثی می شود.

```

✓ [59] def kmer(k, data):
0s     for index, row in data.iterrows():
         if len(row) > 1:
             updated_data = kmer_for_one_sequence(row[1], k)
             # row[1] = np.array(list(updated_data.values()))
             row[1] = list(updated_data.values())
         else:
             updated_data = kmer_for_one_sequence(row[0], k)
             row[0] = list(updated_data.values())

```

این تابع `kmer` را بر روی مجموعه‌ای از رشته‌ها اعمال می‌کند. (درواقع از تابع قبلی استفاده می‌کند و آن را به مجموعه‌ای از رشته‌ها اعمال می‌کند)

همانطور که در کد مشخص است بر روی مجموعه داده ورودی پیمایش می‌کند و `index` و `row` را برای هر سر استخراج می‌کند و رشته را از داخل `row` بدست می‌آورد.

و تابع `kmer_for_one_sequence` را بر روی هر رشته اعمال می‌کند و آن رشته را با برداری که از تابع قبلی بدست می‌آید جایگزین می‌کند

( چک می‌کند که:

اگر طول `row` بیش‌تر از ۱ بود: داده‌ی آموزشی است که هم `Type` دارد و هم `Sequence`

در غیر این صورت: داده‌ی تست که فقط `Sequence` دارد.

به این صورت رشته را پیدا می‌کند)

نهایتاً رشته‌های ما با یکسری بردار عددی نرمال‌شده، جایگزین می‌شوند.

حال که این توابع را تعریف کردیم، از آن‌ها همراه با یکسری تکنیک‌های دیگر استفاده می‌کنیم و بر مجموعه داده‌هایمان اعمال می‌کنیم تا به فرمت مناسب برای ارائه به مدل و پردازش در بیایند:

داده‌ها را می‌خوانیم:

```

▶ import pandas as pd
  # import seaborn as sns
  # import matplotlib.pyplot as plt

train_set = pd.read_csv("training_set.csv")
train_set

```

برچسب‌ها را با مقادیر عددیشان جایگزین می‌کنیم تا بتوان در مدل از آن‌ها استفاده کرد:

```
train_set.iloc[:, 0:1] = train_set.iloc[:, 0:1].replace('Class1', 1)
train_set.iloc[:, 0:1] = train_set.iloc[:, 0:1].replace('Class2', 2)
train_set.iloc[:, 0:1] = train_set.iloc[:, 0:1].replace('Class3', 3)
train_set.iloc[:, 0:1] = train_set.iloc[:, 0:1].replace('Class4', 4)
train_set.iloc[:, 0:1] = train_set.iloc[:, 0:1].replace('Class5', 5)
train_set.iloc[:, 0:1] = train_set.iloc[:, 0:1].replace('Class6', 6)

train_set
```

Type	Sequence
0	5 TACCACCTACGCTGACAATGGATGTTATTGTACCGATTGGAATTA...

و **kmer** را اعمال می‌کنیم:

**kmer** را بر روی رشته‌های مجموعه داده‌مان اعمال می‌کنیم و به این صورت بردارهای **kmer** متناظرشان را جایگزین می‌کنیم.

در بخش‌های مختلف کد از  $k=2$  برای پیاده‌سازی **kmer** استفاده شده. این مقدار از تست مقادیر مختلف بدست آمده و نتیجه شده. درواقع مدل شبکه عصبی که ایجاد کردیم با چند مقدار مختلف  $k$  نتیجه مناسب را می‌دهد اما  $k=2$  برای مدل **mlp** ما از بقیه بهتر بود. درواقع این برتری زمانی مشخص شد که یک شبکه عصبی ساده با تعداد پارامترها و نورون‌های کم داشتیم. در تست‌های مکرر حالت‌های مختلف و مقادیر مختلف  $k$  مشاهده شد که در یک شبکه عصبی سبک با افزایش  $k$  درصد خطا بر روی داده‌های **develop** افزایش می‌یابد. اما با  $k=2$  شبکه‌ی ساده‌ی ما می‌توانست درصد زیادی از داده‌های ارزیابی را درست تشخیص دهد و فقط تعداد کمی را با تفاوت تنها ۱ واحد متفاوت تشخیص می‌داد. اما با افزایش  $k$  مقدار خطا بیش‌تر و بزرگ‌تر می‌شد. (علت این امر هم این است که پنجره‌های کوچک‌تر زیربخش‌های بیش‌تری از رشته‌ها را بررسی می‌کنند و به عبارتی جزئی‌تر رشته را بررسی می‌کنند. این مطلب را می‌توان با درکنار هم قرار دادن یک  $k$  خیلی کوچک و یک  $k$  خیلی بزرگ به خوبی درک کرد)

```
kmer(k=2, data=train_set)
```

Type	Sequence
0	5 [0.07471585615800261, 0.04713992919694429, 0.0...
1	5 [0.0692064083457526, 0.043964232488822655, 0.0...
2	3 [0.03645266594124048, 0.06641050054406965, 0.0...
3	6 [0.0641025641025641, 0.05042735042735043, 0.03...
4	0 [0.0408976737398593, 0.0795795973955594, 0.04...

سپس مجموعه داده و برجسب های آموزشی را در پارامترهای `train_data` و `train_label` میریزیم. و نهایتا داده ها و برجسب را به صورت لیست درمی آوریم تا در مدل بتوانیم استفاده کنیم.

```
train_data = train_set.iloc[:, 1:2]
train_labels = train_set.iloc[:, 0:1]

train_data
```

	Sequence
0	[0.07471585615800261, 0.04713992919694429, 0.0...
1	[0.0692064083457526, 0.043964232488822655, 0.0...
2	[0.09245000504161618, 0.08814959351498835, 0.0...

```
[71] train_data = list(train_data['Sequence'])
train_labels = list(train_labels['Type'])

print(train_data)
print(train_labels)

[[0.07471585615800261, 0.04713992919694429, 0.04574250046580958, 0.09288242966275387, 0.07881498043599776, 0.050307434
5, 5, 3, 6, 2, 3, 3, 1, 4, 3, 1, 1, 2, 4, 1, 2, 3, 1, 5, 2, 2, 1, 3, 5, 4, 5, 4, 3, 4, 5, 5, 4, 5, 4, 5, 5, 1, 1, 6,
```

همین کارها و مراحل پیش پردازش را برای مجموعه داده های `test` و `develop` هم انجام می دهیم تا قابل استفاده، پردازش و تحلیل بشوند.

بعد از اعمال پیش پردازش بر مجموعه داده ها گام بعدی فرا می رسد:

## - ایجاد مدل، پردازش، یادگیری، ارزیابی و تست

ما از MLP استفاده کردیم و مدل را ایجاد کردیم:

ابتدا یک مدل دسته بندی MLP ایجاد می کنیم. این مدل ۳ لایه پنهان دارد که هر کدام ۶۴ نورون دارند. از تابع فعال ساز `relu` استفاده شده است که غشبه مارا غیر خطی می کند. رندم استیت هسته ی شروع و مقداردهی رندم اولیه پارامترها را تعیین می کند و مکس ایتر هم حداکثر تعداد پیمایش ها را تعیین می کند.

این مقادیر پارامترها باتوجه به توزیع داده ها خوب عمل می کنند و با تست و تحلیل محاسبه شده اند. برای برخی از پارامترها مقدار ا تیمم وجود دارد به این صورت که بیش تر یا کم تر از یک مقدار خاص، عملکرد مدل کاهشی می شود و ما با پیدا کردن آن نقطه خاص می توانیم عملکرد مناسب مدل را بدست بیاوریم.

در طراحی مدل باتوجه به کم بودن داده های آموزش از `earlystopping` می توانیم استفاده کنیم تا از بیش برازش جلوگیری شود.

نهایتا داده های آموزشی همراه با برجسب هایشان در مدل فیت شدند و مدل با آن ها آموزش می بیند.

```
from keras.backend import dropout

clf = MLPClassifier(hidden_layer_sizes=(64, 64, 64), activation="relu", random_state=1, early_stopping=True)

clf.fit(train_data, train_labels)

MLPClassifier(early_stopping=True, hidden_layer_sizes=(64, 64, 64),
              random_state=1)
```

حال داده‌های develop را به مدل می‌دهیم تا پیش‌بینی کند و نتایج را نشان می‌دهیم.

```
[81] development_predicts = clf.predict(development_data)

print(development_predicts)

[5 2 2 6 3 6 3 5 1 6 5 6 6 2 6 5 2 5 5 3 4 2 4 6 3 2 3 5 4 5 4 3 5 2 1 1 5
 2 2 5 6 1 6 1 1 3 5 5 5 4 1 6 1 3 1 6 6 1 3 4 2 4 2 2 4 6 1 1 3 1 3 1 3 1
 5 1 5 2 4 4 2 4 1 2 6 6 2 3 6 5 6 4 4 6 4 2 4 5 5 3 1 3 5 2 5 4 3 2 6 3 1
 3 6 6 1 2 3 3 4 2 5 4 2 3 4 2 1 5 4 6 4 4 4 6 1 4 3 2 6 2 1 2 3 2 3 4 3 6
 5 6 4 4 5 1 5 4 3 1 6 4 3 6 5 5 5 5 3 1 2 4 3 3 1 2 6 2 6 1 1 1]
```

سپس تفاوت نتایج پیش‌بینی شده‌ی مدل را با مقادیر واقعی نشان می‌دهیم که برای همه داده‌های develop برابر صفر است و یعنی مدل عملکرد ۱۰۰ درصدی از خود نشان داده‌است.

این امتیاز هم در همین بخش با `r2_score` محاسبه و نمایش داده‌شده است.

```
[82] print(development_predicts - development_labels)
print("\nThe Score with ", (r2_score(development_predicts, development_labels)))

[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0]

The Score with 1.0
```

همین کار را برای داده‌های تست انجام می‌دهیم:

```
from keras.backend import dropout
clf = MLPClassifier(hidden_layer_sizes=(64, 64, 64), activation="relu", random_state=1, early_stopping=True)

clf.fit(train_data, train_labels)

MLPClassifier(early_stopping=True, hidden_layer_sizes=(64, 64, 64),
              random_state=1)
```

```
[110] test_predicts = clf.predict(test_data)

print("answers =\n", (test_predicts))

answers =
[5 5 3 1 6 1 6 3 4 1 4 3 2 4 4 3 6 3 3 3 4 1 5 6 1 5 3 3 3 3 1 3 6 1 6 1 3
 5 5 6 5 4 5 5 3 6 4 5 1 2 4 3 4 3 3 6 1 6 5 3 1 3 3 3 3 6 4 5 6 6 6 3 2 6
 1 1 4 6 6 2 5 4 1 1 3 5 6 4 3 5 3 1 1 6 5 3 5 1 4 1 3 3 2 2 1 3 4 3 3 3 3 3
 1 4 1 4 4 3 3 3 4 6 1 6 4 5 3 6 2 1 3 3 6 5 2 2 3 2 1 3 3 5 1 6 1 2 2 1 3
 1 1 6 6 2 6 3 5 5 3 6 3 1 2 4 4 3 3 5 4 1 1 2 3 6 1 3 2 6 4 5 6 5 6 5 5 6
 3 3 3 6 4 6 3 3 6 6 6 1 5 1 1 3 3 1 1 1 5 3 2 2 1 1 1 1 3 6 3 3 4 5 6 4
 1 4 3 4 6 5 4 6 5 5 6 6 4 5 3 6 6 3 1 1 4 4 4 5 1 2 6 1 6 3 2 2 1 2 3 6 6
 5 1 3 6 1 3 3 1 1 5 4 1 3 1 2 1 1 3 3 3 6 1 6 4 1 3 4 3 3 4 3 1 3 6 6 2 1
 6 2 1 6 6 1 4 1 6 5 6 3 4 2 4 3 1 6 5 5 2 3 5 4 5 1 3 3 4 1 1 1 1 2 1 5 2
 3 3 4 2 5 1 5 4 1 1 3 3 2 5 4 4 3 2 5 1 2 1 2 1 2 2 3 5 6 1 6 3 1 6 4 1 1
 1 2 3 5 6 3 1 1 4 2 1 6 3 6 3 4 2 6 6 3 3 2 2 6 3 1 3 2 5 1]
```

با توجه به نتایجی که از `quera` بدست آمد می‌بینیم مدل ما عملکرد بسیار بالا و تقریباً ۱۰۰ درصدی دارد که نشان می‌دهد پیش پردازش، نرمال‌سازی، ساخت و استفاده از مدل به خوبی صورت گرفته است.

باتشكر از توجه شما.

اميرمهدى زرین نژاد