

صبغه بندی متن با پایگاه Keras و زبان برنامه نویسی پایتون

۱. مرحله اول آن همواره تعریف و فراخوانی پایگاه ها با پایتون است که برای ما پایتون اصلی ترین و معروف ترین پایگاه است. Pandas

در هر مرحله می توان پایگاه مورد نیاز را تعریف کرد اما بصورت متداول در ابتدا پایگاه ها تعریف می شود.

برای مثال حافظه علاوه بر پایگاه Pandas سایر پایگاه ها زیر نیز مورد نیاز است:

- ✓ sklearn library
- ✓ Keras library
- ✓ matplotlib library

۲. مرحله دوم خواندن متن و ذخیره آن در مقیادهای برنامه است که برای ما به تنهایی به تنهایی نام ستون اطلاعات این فایل متن (text, CSV, TSV, ...) در داخل یک DataFrame قرار می گیرد.

Python

```
import pandas as pd
```

آدرس فایل
مسئله را مشخص می کند

```
filepath_dict = {'yelp': 'data/sentiment_analysis/yelp_labelled.txt',
                  'amazon': 'data/sentiment_analysis/amazon_cells_labelled.txt',
                  'imdb': 'data/sentiment_analysis/imdb_labelled.txt'}
```

حالی ای؟
لیست را تبدیل
Dataframe

```
df_list = []
for source, filepath in filepath_dict.items():
    df = pd.read_csv(filepath, names=['sentence', 'label'], sep='\t')
    df['source'] = source # Add another column filled with the source name
    df_list.append(df)
df = pd.concat(df_list)
print(df.iloc[0])
```

به ازای هر آدرس و Source که وجود دارد، فایل مربوطه را می خواند و Source را هم به آن اضافه می کند و در لیست درج می کند.

برای مثال مطرح شده ۳ داده میانی وجود دارد که هر کدام از تعدادی جمله و label مربوطه با هم در این مثال احساس می شود و مشخص است که به ترتیب با ۱ و ۰ نشان داده شده، مشخص شده است.

در مرحله دوم خواندن متن، ایجاد یک DataFrame مشکل از ۳ متن است که باید Source متن به ابتدای DataFrame اضافه می شود مشخص می شود متن مربوطه به کدام Source است.

* نکته مهم این است که اگر دیتاست معادلی برای Train و Test صبغه بندی نشده

داریم باید هر دو را در Dataframe قرار دهیم تا به از Train و Test انجام شود.

۲- البته لزوماً نیاز به این کار در شبکه‌های محقق نیست و به pre-train شده، این توان که دایرکت و در مراحل بعدی از آن استفاده کرد.

۳- مرحله سوم ایجاد مدل پایه است. مطابق با نوع شبکه و داده‌ها باید اصطلاحات تعریف شود.

در برخی موارد آماده داده‌ها برای سیستم در مرحله‌ای می‌شود. Sklearn این اسطون را دارد که باید خطای داده‌های سیستم را یادگیری را برآورد.

۴- ایجاد مدل اولیه یا همان baseline Model بسیار در رویکردهای مبتنی بر یادگیری ماشین مهم است.

مبتنی از ماچانه نه عمل split سیستم و یادگیری را انجام می‌دهد Python

```
>>> from sklearn.model_selection import train_test_split
```

```
>>> df_yelp = df[df['source'] == 'yelp']
```

→ برای مثال با ستون فایل yelp عملیات انجام شده است

```
>>> sentences = df_yelp['sentence'].values
```

```
>>> y = df_yelp['label'].values
```

```
>>> sentences_train, sentences_test, y_train, y_test = train_test_split(
...     sentences, y, test_size=0.25, random_state=1000)
```

مقدار ستون Sentence داخل object رنج می‌شود

که نیل‌های مرتبط به ترتیب object به نام آرگومان می‌شود

داده‌های Train و Test و Label‌های مرتبط

از هم جدا می‌شوند آماده داده سیستم

برای Random Sampling استفاده می‌شود.

ولی optional است

۴- مرحله چهارم ایجاد شبکه عصبی یا هر شبکه دیگری است که به Train شود.

در این مثال هدف به شبکه عصبی محقق است.

با توجه به اینکه مثال ما ایجاد به شبکه عصبی (محقق) است. در این شبکه ما به طبقه

شبکه‌های عصبی (نورون‌های ورودی و خروجی) با بردار ویژگی‌ها (Feature Vector) تغذیه می‌شود و این مقادیر به لایه‌های پنهان (hidden layer) می‌رود.

به نکته مهم این است که شبکه به وزن (weight) و بایاس (bias)

نیاز دارد.

در حالت غیر قابل مرئی است - (output layer) که می تواند به یک مقدار خود داده باشد

اگر صفتی نبی داشته باشد با نوری است ← لایه هورنی که نور دارد

اند صبیحه بیدار شد چندان دلم دارد ← لایه درونی به تعداد طاس ها نمود خواهه طاس .

تعداد لایه‌های پهنای توده متغیر و متفاوت باشد به سبب سستی دارد

در سبب این activation function مهم است - انواع مختلفی دارد ← Relu
معمولاً برای لایه های پنهان

Backpropagation: همین روش را برعکس می‌کنیم و در جهت برعکس می‌رویم. برای سیمای سگ.

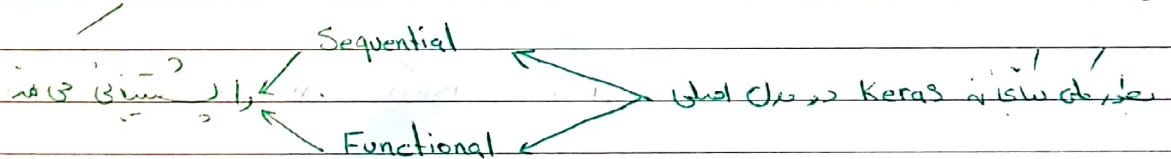
Softmax ↙
برای محاسبه چند کلاس

loss function ↘
انواع مختلفی دارد ↘
cross entropy ↙

لایه Keras دستورات ساده‌ای برای تعریف شبکه دارد. این لایه‌ها متعلق به Tensorflow و Theano است و دستورات ساده‌تری دارد.

در مرحله چهارم در واقع مسئله تحلیل می شود تا نیازهای استخراج می شود و در مرحله پنجم این مدل با کمک نتایج آن طراحی می شود.

متون حول ساختاری تدریجی (Sequential) دارند و ساختار Keras آن، الیبتی



در این مثال مدل Sequential است، مدل Sequential داده‌ها را به صورت توالی می‌گیرد.

ابعاد مدل رفوع مدل مهم است - و برای سبب تعریف مدل

* ابعاد ورودی بستگی به ویژگی‌ها و بردار ویژگی‌ها دارد .
input_dim \rightarrow number of features

عبارت میں اعداد لایہ دی سہلہ نہ تہہ تہہ حرف می شود

input_dim => Number of features

model = Sequential()

model.add(layers.Dense(10, input_dim=input_dim, activation='relu'))

* دستور model.summary() می‌تواند خلاصه اطلاعاتی از شبکه را نمایش دهد

برای تبدیل داده‌ها به داده‌های قابل فهم توسط ماشین، به روش تبدیل آن‌ها به بردار به هر طریقی که عددی دهه است که اطلاعات اضافی ندارد و خیلی مناسب نیست

* به دلیل نیاز داشتن اطلاعات اضافی Vectorize کردن زیاد مفید نیست

تعداد مراحل Train و داده Test باید مشخص شود.

Python

```
>>> from sklearn.feature_extraction.text import CountVectorizer
```

```
>>> vectorizer = CountVectorizer()
>>> vectorizer.fit(sentences_train)
```

```
>>> x_train = vectorizer.transform(sentences_train)
```

```
>>> x_test = vectorizer.transform(sentences_test)
```

```
>>> x_train
```

```
<750x1714 sparse matrix of type '<class 'numpy.int64''
with 7368 stored elements in Compressed Sparse Row format>
```

خودش عبارت را
Tokenize می‌کند
و به طایفه‌ها (کلمات) تبدیل می‌کند.
نصفه و کلمات را
امکان را حذف می‌کند

تعداد sort کردن
به هر طریقی که
عدد مناسبی می‌دهد

در صفحه بالا عبارت داخل object می‌باشد. Sentences_train و Sentences_test قرار دارند

* این کار به نوعی سیستمی که دارش است و باید چانه‌ها می‌باشد nltk می‌تواند این کار را انجام دهد

حالا داده‌ها آماده است و باید برای درسیته است

Python

```
>>> from keras.models import Sequential
```

```
>>> from keras import layers
```

```
>>> input_dim = x_train.shape[1] # Number of features
```

```
>>> model = Sequential()
```

```
>>> model.add(layers.Dense(10, input_dim=input_dim, activation='relu'))
```

```
>>> model.add(layers.Dense(1, activation='sigmoid'))
```

```
Using TensorFlow backend.
```

Vectorize ایجاد شده است
batch size درسی

سینه
میان در
لایه دارد

دستوری `shape[]` ابعاد دوری به راستی می‌دهد برای سبب ما باید مشکل
کند

`input_dim = x_train.shape[1]`

در نتیجه مهم این ابعاد دوری را مورد نیاز لایه اول سبب عصبی است. لایه‌های دیگر خودکار
ایجاد می‌شود.

لایه‌های پس از دیگری تعریف می‌شوند در نتیجه در سبب مشکل است.

Python

`>>> model.compile(loss='binary_crossentropy',
... optimizer='adam',
... metrics=['accuracy'])` → متریک ارزیابی

`>>> model.summary()`
خلاصه‌های مدل را نشان می‌دهد

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 10)	17150
dense_2 (Dense)	(None, 1)	11
Total params: 17,161		
Trainable params: 17,161		
Non-trainable params: 0		

برای شروع فرآیند `Training` به از تعریف مدل و لایه‌ها، فرآیند یادگیری نیز باید تعریف
شود.

حال دیگر شروع عملیات `Training` امکان پذیر است.

Python دستور شروع عملیات
یادگیری

`>>> history = model.fit(x_train, y_train,
... epochs=100, → تعداد iteration
... verbose=False,
... validation_data=(x_test, y_test)
... batch_size=10)`
← تعداد Sample های که در هر epoch استفاده می‌شود

فرآیند `iterative` برای یادگیری داریم.

با یک مدل مبتنی بر Train انجام خواهد شد. حالا به مدل ارزیابی شود.

Python

ارزیابی مدل یا مکرمه شد

```
>>> loss, accuracy = model.evaluate(X_train, y_train, verbose=False)
>>> print("Training Accuracy: {:.4f}".format(accuracy))
>>> loss, accuracy = model.evaluate(X_test, y_test, verbose=False)
>>> print("Testing Accuracy: {:.4f}".format(accuracy))
Training Accuracy: 1.0000
Testing Accuracy: 0.7960
```

سپس برای تعریف سبب، Train و در انتها تست سبب مراحل زیر انجام شد:

- 1- Select Model
- 2- Model add layers (1..N)
- 3- compile Model
- 4- Fit Model
- 5- Evaluate Model

در مثال فوق Text منطبق بر Vectorize شد و اطلاعات اولیه دلیلی نباشد.
مطابق آن Text فراداده‌های ترتیبی است. تمام دوری مرحله قبل در مثال مبتنی بر
عنوان Vector به سبب داده شد. (به مدل مثال مبتنی Bag of word (Bow)
(گفته می‌شود).

word as a
vector
Character
as a vector

N-gram

امروزه برای Vectorize کردن داده‌های دوری و بردار دارد.
مثالی برای word as a vector Hot encode - one است.
که برداری به طول حالات داریم و هر کلمه یک بردار دارد.

با یک one-hot encoding داده‌های غیر عددی به عدد تبدیل می‌شوند تا برای ماشین مبتنی بر
باشد.

به این روش بیت میتواند اطلاعات را به بردار تبدیل می‌کند و اطلاعات غیر دلیلی
از ترسها حالات ندارد.

امروزه معمول است.

این روش توسط ماتریکس SKlearn تبدیل می‌شود.

بهترین روش برای نمایش داده‌های غیر عددی به خاصیت ترسب word embedding دارند.
است.

مثال دیگر: همان داده‌ها / آشنایی با لغت word embedding پیاده سازی می شود.

در مورد word embedding داده‌ها با ایجاد لغت ایجاد می شود اما اطلاعات بیشتر در این داده‌ها وجود دارد.
بردار طاقتهای دیگر می شود که به همین خاطر به آن dense word vectors هم گفته می شود.

* نکته Keras دستورانی برای تعریف word embedding دارد اما با توجه به نوع مسئله word embedding می تواند متفاوت باشد و خود برنامه نویس آن ها را تعریف کند.
یادگیری آن به یادگیری استن نسبت اما ساختار آمار می باشد اما یادگیری آن به هر دو نسبت به یک ساختار هندسی است.

Map semantic meaning into a Geometric space.

embedding Space

دوربرد برای پیاده سازی word embedding وجود دارد.
نکته Keras از ویژگی بردارهای ترتیبی و متن ساده پشتیبانی می کند.
در مرحله قبل از یادگیری (pre-train) به صورت سیستم درج شده است.
در مرحله سبک عصبی ایجاد شود.

دستورات زیر برای پیش بردارها استفاده می شود که روش متداول ترین نسبت به استفاده از آن در بدنه سبک عصبی است.
احتمال بازیابی طاقتهای وکتور آن ها نیز وجود دارد.

Python

>>>

```
>>> from keras.preprocessing.text import Tokenizer
```

```
>>> tokenizer = Tokenizer(num_words=5000)
```

```
>>> tokenizer.fit_on_texts(sentences_train)
```

```
>>> X_train = tokenizer.texts_to_sequences(sentences_train)
```

```
>>> X_test = tokenizer.texts_to_sequences(sentences_test)
```

```
>>> vocab_size = len(tokenizer.word_index) + 1 # Adding 1 because of reserved 0 index
```

```
>>> print(sentences_train[2])
```

```
>>> print(X_train[2])
```

Of all the dishes, the salmon was the best, but all were great.

```
[11, 43, 1, 171, 1, 283, 3, 1, 47, 26, 43, 24, 22]
```


Python

```
>>> for word in ['the', 'all', 'happy', 'sad']:
...     print('{}: {}'.format(word, tokenizer.word_index[word]))
the: 1
all: 43
happy: 320
sad: 450
```

برای بهینه‌سازی به طبات

در مسئله این است که طبات داخل هر جمله تعداد متفاوتی دارند و در واقع طول جمله‌ها متفاوت خواهد بود که از `pad_sequence` استفاده می‌کنیم که طول‌های خالی را به صفر پر می‌کند تا طول یکسان شود.

Python

```
>>> from keras.preprocessing.sequence import pad_sequences
```

```
>>> maxlen = 100 → حداکثر طول
```

{ اعمال padding با } `>>> X_train = pad_sequences(X_train, padding='post', maxlen=maxlen)`
`>>> X_test = pad_sequences(X_test, padding='post', maxlen=maxlen)`

مقدار، دستور روی داده

Train و Test

به صورت هم طول

شکل

```
>>> print(X_train[0, :])
[ 1 10  3 282 739 25  8 208 30 64 459 230 13  1 124  5 231  8
 58  5 67  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]
```

امکان استفاده از `embedding layer` در بین شبکه عصبی هم وجود دارد که با این دستور انجام می‌شود.

مقادیر زیر صاف باید برای آن مشخص شود.

input_dim → اندازه طبات
 output_dim → اندازه بردار خروجی
 input_length → اندازه ترتیب (Sequence)

```
→ model.add(layers.Embedding( ))
```

```
→ model.add(layers.Flatten(1))
```

}

word2vec → Google

Glove → Stanford NLP Group

Subject

Date

عده بر word embedding پاده سازن شده در حین هی تبلی به Vectorize
word embedding نایه Keras مای انجام است که بصورت یولته داخل که
انجام سه (jointly) امکان استفاده از word embedding های که از سبل حالیه شده
شده وجود دارد.

مهم ترین و معروف ترین متد موجود و پاده سازن شده برای آن وجود دارد و دیلرد word2vec
توسط یافته توسط گوگل است و دیلرد Glove نیز معروف است

word2vector → use neural network

Glove → use co-occurrence matrix

* دیلرد word2vec دقیق تر و سریع تر است اما دیلرد Glove رابط تر مای برریت
است

میتونه هم این است که اگر بخواهم دیلرد حالیه دیلردن را برای word embedding استفاده
کنم این امکان توسط نایه Gensim وجود دارد.

→ امکان استفاده از مدل های معروف شده توسط برنامه نویسی شده وجود دارد

ایمپادیرن به ایادی هم را بصورت ماتریس و در انبار و دیلرد را هم در نظر بگیریم و در لایه
عرف لایه های شبکه به عنوان embedding قرار دهیم

برلنریند ماتریس با اعداد Random به عنوان embedding به سبل داده شده است
و به جای آن برنامه نویسی می تواند ماتریس خود را از دیلردن و مای خود در نظر خود بخواند
و در این لایه وارد کند

Python

```
model = Sequential()
embedding_dim = 50
numpy_matrix = np.random.rand(1747, 50)
model.add(layers.Embedding(vocab_size, embedding_dim,
                           weights=[numpy_matrix],
                           input_length=maxlen,
                           trainable=True))
model.add(layers.GlobalMaxPool1D())
model.add(layers.Dense(10, activation='relu'))
model.add(layers.Dense(1, activation='sigmoid'))
model.compile(optimizer='adam',
              loss='binary_crossentropy',
              metrics=['accuracy'])
model.summary()
history = model.fit(X_train, y_train,
                   epochs=50,
                   verbose=False,
                   validation_data=(X_test, y_test),
                   batch_size=10)
loss, accuracy = model.evaluate(X_train, y_train, verbose=False)
print("Training Accuracy: {:.4f}".format(accuracy))
loss, accuracy = model.evaluate(X_test, y_test, verbose=False)
print("Testing Accuracy: {:.4f}".format(accuracy))
```

ماتریس Random می نه در
لایه embedding به سبل
داده می شود