

Sentimental Analysis of IMDB movie reviews

```
In [2]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
from matplotlib import style
style.use('ggplot')
import re
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer
from nltk.corpus import stopwords
stop_words = set(stopwords.words('english'))
from nltk.cloud import wordcloud
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split

In [3]: df = pd.read_csv('IMDB Dataset.csv')
df.head()
```

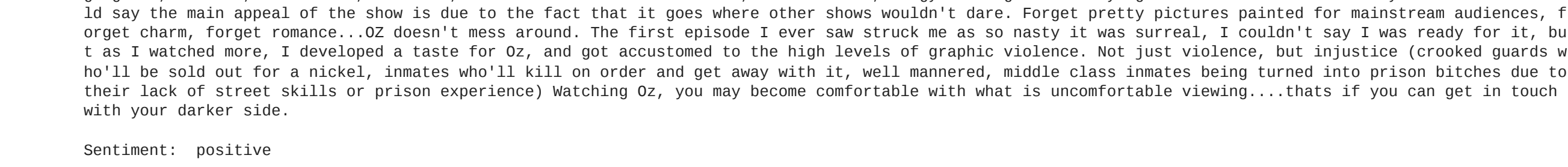
```
Out[3]:
```

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production. The...	positive
2	I thought this was a wonderful way to spend...	positive
3	Basically there's a family where a little boy ...	negative
4	Peter Mattei's "Love in the Time of Money" a...	positive

```
In [4]: df.shape
Out[4]: (50009, 2)
```

```
In [5]: df.info()
Out[5]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 50009 entries, 0 to 49999
Data columns (total 2 columns):
#   column            non-null count  dtype
#   ----            -
0    review           50009 non-null object
1    sentiment        50009 non-null object
dtypes: object(2)
memory usage: 781.4+ KB
```

```
In [6]: sns.countplot(x='sentiment', data=df)
plt.title('Sentiment distribution')
Text(0.5, 1.0, 'Sentiment distribution')
```



```
In [7]: for i in range(5):
print('review:', i, [1])
print(df['review'].iloc[i], "\n\n")
print('Sentiment:', df['sentiment'].iloc[i], "\n\n\n")

Review: [0]
One of the other reviewers has mentioned that we encounter, <br /><br />this being a variation on the Art her
/>The first thing that struck me about Oz was its brutality and unflinching scenes of violence, which set in right from the word GO. Trust me, this is n
of a show for the faint hearted or timid. This show pulls no punches with regards to drugs, sex or violence. In the classic use of the word <
/>I watched <br />This is called Oz as that is the nickname given to the Goudal Maximum Security State Penitentiary. It focuses mainly on Emerald City, an experimental
section of the prison where all the cells have glass fronts and face inwards, so privacy is not high on the agenda. En City is home to many, Aryans, Muslims,
is a terrifically written and performed piece. A masterful production about one of the great master's of comedy and his life. <br /><br />The realism really
is say the main appeal of the show is due to the fact that it goes where other shows wouldn't dare. Forget pretty pictures painted for mainstream audiences, f
orget charm, forget romance, Oz doesn't mess around. The first episode I ever saw struck me as so raucy it was surreal, I couldn't say I was ready for it, bu
t I watched <br />I developed a taste for Oz, and got accustomed to the high levels of graphic violence. Not just violence, but injustice (crooked guards w
ho'll be sold out for a nickel, inmates who'll kill on order and get away with it, well mannered, middle class inmates being turned into prison bitches due to
their lack of street skills or prison experience) watching Oz, you may become comfortable with what is uncomfortable viewing...that's if you can get in touch
with your darker side.

Sentiment: positive

Review: [1]
A wonderful little production. <br /><br />The filming technique is very unassuming- very old time-BBC fashion and gives a comforting, and sometimes discomfort
ing, sense of realism to the entire piece. <br /><br />The actors are extremely well chosen- Michael Sheen not only "has got all the polaris" but he has all t
he voices down pat. You can truly see the seamless editing guided by the references to Williams' diary entries, not only is it well worth the watching but
it is a terrifically written and performed piece. A masterful production about one of the great master's of comedy and his life. <br /><br />The realism really
comes home with the little things; the way of the guards which, rather than the traditional "dream" techniques remains solid then disappears. It plays
out knowledge about time money, power and success do to people in the different situations we encounter, <br /><br />this being a variation on the Art her
decorating every surface) are terribly well done.

Sentiment: positive
```

```
Review: [2]
This was a wonderful way to spend time on a too hot summer weekend, sitting in the air conditioned theater and watching a light-hearted comedy. The
plot is simplistic, but the dialogue is witty and the characters are likable (even the well bread suspected serial killer). While some may be disappointed wh
o in realize this is not Match Point 2: Risk Addiction, I thought it was proof that Woody Allen is still fully in control of the style many of us have grown
used to. I watched <br />This was the most I've laughed at one of Woody's comedies in years (here I say a decade?), while I've never been impressed with Scarlett Joh
anson, in this she managed to tone down her "sexy" image and jumped right into an average, but spirited young woman- <br /><br />this may not be the crown jewel
of his career, but it was wittier than "Devil wears Prada" and more interesting than "Superman" a great comedy to go see with friends.

Sentiment: positive
```

```
Review: [3]
Basically there's a family where a little boy (Jake) thinks there's a zombie in his closet & his parents are fighting all the time. <br /><br />this movie is a
load of crap, and a soap opera, and suddenly, Jake decides to become Rambo and kill the zombie. <br /><br />P.S. first of all when you're going to make a film you ha
ve to decide if it's a thriller or a drama. As a drama the movie is watchable. Parents are divorcing & arguing like in real life. And then we have Jake with his c
loset which totally ruins all the film. I expected to see a BODEGUYMAN similar movie, and instead I saw a meaningless thriller spots- <br /><br />
out of 18 stars for the well playing parents & decent dialogs. As for the shots with Jake: just ignore them.

Sentiment: negative
```

```
Review: [4]
Petter Mattei's "Love in the Time of Money" is a visually stunning film to watch. Mr. Mattei offers us a vivid portrait about human relations. This is a movie
that is a terrifically written and performed piece. A masterful production about one of the great master's of comedy and his life. <br /><br />The realism really
comes home with the little things; the different stages of loneliness each one inhabits. A big city is not exactly the best place in which human relations find sinc
erity. As one discerns is the case with most of the people we encounter <br /><br />the acting is good under Mr. Mattei's direction. Steve Buscemi,
Rosario Dawson, Carol Kane, Michael Imperioli, Adrian Grenier, and the rest of the talented cast, make these characters come alive. <br /><br />wish Mr. Mat
tei good luck and await anxiously for his next work.

Sentiment: positive
```

```
In [8]: def no_of_words(text):
words = text.split()
word_count = len(words)
return word_count

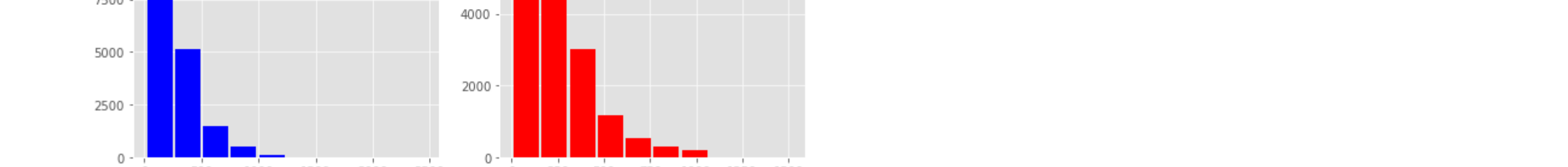
In [9]: df['word_count'] = df['review'].apply(no_of_words)

In [10]: df.head()
```

```
Out[10]:
```

	review	sentiment	word count
0	One of the other reviewers has mentioned that ...	positive	307
1	A wonderful little production. The...	positive	362
2	I thought this was a wonderful way to spend ti...	positive	166
3	Basically there's a family where a little boy ...	negative	138
4	Peter Mattei's "Love in the Time of Money" a...	positive	230

```
In [11]: fig, ax = plt.subplots(1,2, figsize=(16,8))
ax[0].hist(df[df['sentiment'] == 'positive']['review'].str.len(), label='Positive', color='blue', rwidth=0.9);
ax[0].legend(loc='upper right');
ax[1].hist(df[df['sentiment'] == 'negative']['review'].str.len(), label='Negative', color='red', rwidth=0.9);
ax[1].legend(loc='upper right');
fig.suptitle('Number of words in review')
plt.show()
```



```
In [12]: fig, ax = plt.subplots(1,2, figsize=(16,8))
ax[0].hist(df[df['sentiment'] == 'positive']['review'].str.len(), label='Positive', color='blue', rwidth=0.9);
ax[0].legend(loc='upper right');
ax[1].hist(df[df['sentiment'] == 'negative']['review'].str.len(), label='Negative', color='red', rwidth=0.9);
ax[1].legend(loc='upper right');
fig.suptitle('Number of words in review')
plt.show()
```



```
In [13]: df['review'] = df['review'].replace('positive', 1, inplace=True)
df['review'] = df['review'].replace('negative', 0, inplace=True)

In [14]: df.head()
```

```
Out[14]:
```

	review	sentiment	word count
0	One of the other reviewers has mentioned that ...	1	307
1	A wonderful little production. The...	1	362
2	I thought this was a wonderful way to spend ti...	1	166
3	Basically there's a family where a little boy ...	0	138
4	Peter Mattei's "Love in the Time of Money" a...	1	230

```
In [15]: def data_processing(text):
text = text.lower()
text = re.sub('<br />', '', text)
text = re.sub('http[s]?://[a-zA-Z]*[http[s]?*', '', text, flags=re.MULTILINE)
text = re.sub('<\/>', '', text)
text = re.sub('<\/>', '', text)
text_tokens = word_tokenize(text)
return text = [w for w in text_tokens if not w in stop_words]
return text = ".join(filtered_text)
```

```
In [16]: df['review'] = df['review'].apply(data_processing)
```

```
In [17]: duplicated_count = df.duplicated().sum()
print('Number of duplicate entries:', duplicated_count)
Number of duplicate entries: 421

In [18]: df = df.drop_duplicates('review')
```

```
In [19]: stemmer = PorterStemmer()
def stemming(data):
text = stemmer.stem(word) for word in data
return data

In [20]: df['review'] = df['review'].apply(lambda x: stemming(x))

In [21]: df['word_count'] = df['review'].apply(no_of_words)
df.head()
```

```
Out[21]:
```

	review	sentiment	word count
0	one reviewers mentioned watching 1 oz episode ...	1	168
1	wonderful little production filming technique ...	1	84
2	thought wonderful way spend time hot summer we...	1	86
3	basically there's family little boy jake thinks...	0	67
4	peter mattei love time money visually stunn...	1	125

```
In [22]: pos_reviews = df[df['sentiment'] == 1]
pos_reviews.head()
```

```
Out[22]:
```

	review	sentiment	word count
0	one reviewers mentioned watching 1 oz episode ...	1	168
1	wonderful little production filming technique ...	1	84
2	thought wonderful way spend time hot summer we...	1	86
3	peter mattei love time money visually stunn...	1	125
5	probably alltime favorite movie story selfless...	1	58

```
In [23]: text = ".join([word for word in pos_reviews['review']]
plt.figure(figsize=(20,15), facecolor='None')
wordcloud = WordCloud(max_words=500, width=1600, height=800).generate(text)
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.title('Most frequent words in positive reviews', fontsize=19)
plt.show()
```



```
In [24]: from collections import Counter
count = Counter()
for text in pos_reviews['review'].values:
for word in text.split():
count[word] +=1
count.most_common(15)
```

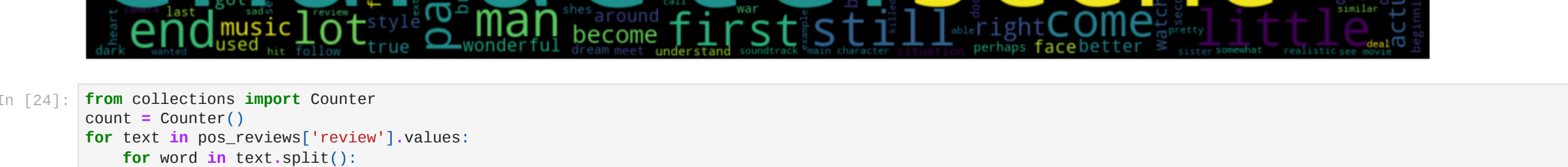
```
Out[24]:
[('movie', 47081),
('film', 34651),
('one', 24393),
('like', 16998),
('good', 14281),
('even', 13568),
('story', 12338),
('great', 11841),
('time', 11724),
('well', 10938),
('really', 10838),
('also', 10516),
('don't', 10320),
('even', 9818),
('much', 8971)]
```

```
In [25]: pos_words = pd.DataFrame(count.most_common(15))
pos_words.columns = ['word', 'count']
pos_words.head()
```

```
Out[25]:
```

	word	count
0	film	35285
1	movie	35280
2	one	25621
3	like	16998
4	good	14281

```
In [26]: px.bar(pos_words, x='count', y='word', title='Common words in positive reviews', color='word')
```

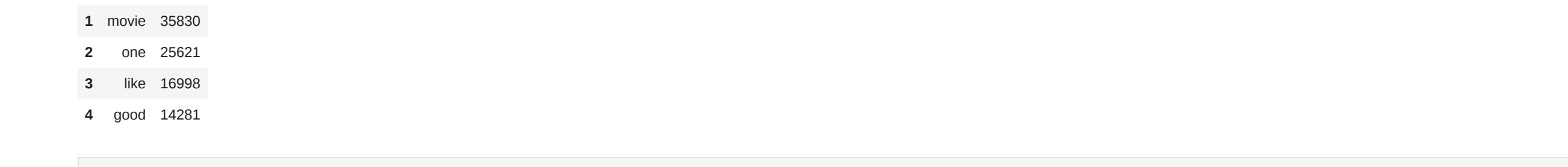


```
In [27]: neg_reviews = df[df['sentiment'] == 0]
neg_reviews.head()
```

```
Out[27]:
```

	review	sentiment	word count
3	basically there's family little boy jake thinks...	0	67
7	show amazing both movie and the film...	0	83
8	encouraged positive comments like looking fore...	0	64
9	phil allen one quirky film humor based mean...	0	93
11	saw movie 12 came recall scared scene big bl...	0	84

```
In [28]: text = ".join([word for word in neg_reviews['review']]
plt.figure(figsize=(20,15), facecolor='None')
wordcloud = WordCloud(max_words=500, width=1600, height=800).generate(text)
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.title('Most frequent words in negative reviews', fontsize=19)
plt.show()
```



```
In [29]: count = Counter()
for text in neg_reviews['review'].values:
for word in text.split():
count[word] +=1
count.most_common(15)
```

```
Out[29]:
[('movie', 47081),
('film', 34651),
('one', 24393),
('like', 16998),
('good', 14281),
('even', 13568),
('story', 12338),
('great', 11841),
('time', 11724),
('well', 10938),
('really', 10838),
('also', 10516),
('don't', 10320),
('even', 9818),
('much', 8971)]
```

```
In [30]: neg_words = pd.DataFrame(count.most_common(15))
neg_words.columns = ['word', 'count']
neg_words.head()
```

```
Out[30]:
```

	word	count
0	movie	47001
1	film	34601
2	one	24301
4	even	14759

```
In [31]: px.bar(neg_words, x='count', y='word', title='Common words in negative reviews', color='word')
```



```
In [32]: X = df['review']
Y = df['sentiment']

In [33]: vect = TfidfVectorizer()
X = vect.fit_transform(df['review'])

In [34]: x_train, x_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, random_state=42)
```

```
In [35]: print('Size of x_train:', (x_train.shape))
print('Size of y_train:', (y_train.shape))
print('Size of x_test:', (x_test.shape))
print('Size of y_test:', (y_test.shape))
Size of x_train: (34704, 221707)
Size of y_train: (34704,)
Size of x_test: (14874, 221707)
Size of y_test: (14874,)
```

```
In [36]: x_train = x_train[:20800]
y_train = y_train[:20800]
x_test = x_test[5000]
y_test = y_test[5000]
```

```
In [37]: print('Size of x_train:', (x_train.shape))
print('Size of y_train:', (y_train.shape))
print('Size of x_test:', (x_test.shape))
print('Size of y_test:', (y_test.shape))
Size of x_train: (20800, 221707)
Size of y_train: (20800,)
Size of x_test: (5000, 221707)
Size of y_test: (5000,)
```

```
In [38]: x_train = x_train.toarray()
x_test = x_test.toarray()
```

```
In [39]: from keras.models import Sequential
from keras.layers import Dense

In [40]: model = Sequential()
model.add(Dense(units=16, activation='relu', input_dim=x_train.shape[1]))
model.add(Dense(units=8, activation='relu'))
model.add(Dense(units=1, activation='sigmoid'))

In [41]: model.compile(optimizer='rmsprop', loss='binary_crossentropy', metrics=['accuracy'])

In [42]: history = model.fit(x_train, y_train, batch_size=16, epochs=15)
```

```
Epoch 1/15
208/208 [=====] - 6s 19ms/step - loss: 0.8778 - accuracy: 0.7190
Epoch 2/15
208/208 [=====] - 5s 25ms/step - loss: 0.5571 - accuracy: 0.9168
Epoch 3/15
208/208 [=====] - 5s 25ms/step - loss: 0.3855 - accuracy: 0.9535
Epoch 4/15
208/208 [=====] - 5s 25ms/step - loss: 0.1908 - accuracy: 0.9715
Epoch 5/15
208/208 [=====] - 5s 25ms/step - loss: 0.0882 - accuracy: 0.9795
Epoch 6/15
208/208 [=====] - 5s 25ms/step - loss: 0.0848 - accuracy: 0.9875
Epoch 7/15
208/208 [=====] - 5s 25ms/step - loss: 0.0811 - accuracy: 0.9915
Epoch 8/15
208/208 [=====] - 5s 25ms/step - loss: 0.0804 - accuracy: 0.9968
Epoch 9/15
208/208 [=====] - 5s 25ms/step - loss: 0.0816 - accuracy: 0.9986
Epoch 10/15
208/208 [=====] - 5s 26ms/step - loss: 0.0878 - accuracy: 0.9985
Epoch 11/15
208/208 [=====] - 5s 25ms/step - loss: 0.0855 - accuracy: 0.9990
Epoch 12/15
208/208 [=====] - 5s 25ms/step - loss: 0.0838 - accuracy: 0.9995
Epoch 13/15
208/208 [=====] - 5s 27ms/step - loss: 0.0824 - accuracy: 0.9995
Epoch 14/15
208/208 [=====] - 5s 25ms/step - loss: 0.0822 - accuracy: 0.9995
```

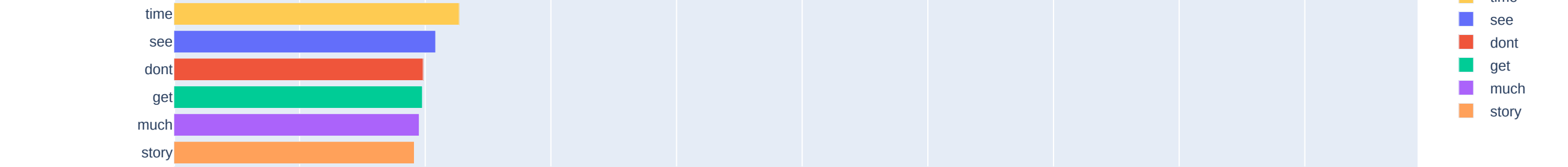
```
In [43]: model.summary()
```

```
Model: "sequential_1"
Layer (type) Output Shape Param #
-----
dense_3 (Dense) (None, 16) 3547328
dense_2 (Dense) (None, 8) 136
dense_1 (Dense) (None, 1) 9
```

```
Total params: 3,547,473
Trainable params: 3,547,473
Non-trainable params: 0
```

```
In [44]: test_loss, test_acc = model.evaluate(x_test, y_test)
print('Test loss:', test_loss)
print('Test accuracy:', test_acc)
16/16 [=====] - 0s 5ms/step - loss: 0.5546 - accuracy: 0.8320
Test loss: 0.55462270413269
Test accuracy: 0.832089171681377
```

```
In [45]: plt.figure(figsize=(16,8))
plt.subplot(1,2,1)
plt.plot(history.history['loss'], color='r', label='loss')
plt.xlabel('Training Loss')
plt.ylabel('Number of epochs')
plt.legend()
plt.subplot(1,2,2)
plt.plot(history.history['accuracy'], color='b', label='accuracy')
plt.xlabel('Training Accuracy')
plt.ylabel('Number of epochs')
plt.legend()
plt.show()
```



```
In [46]: model = Sequential()
model.add(Dense(units=16, activation='relu', input_dim=x_train.shape[1]))
model.add(Dense(units=8, activation='relu'))
model.add(Dense(units=1, activation='sigmoid'))

In [47]: model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

In [48]: history = model.fit(x_train, y_train, batch_size=16, epochs=15)
```

```
Epoch 1/15
208/208 [=====] - 2s 12ms/step - loss: 0.6563 - accuracy: 0.6105
Epoch 2/15
208/208 [=====] - 2s 12ms/step - loss: 0.2878 - accuracy: 0.9548
Epoch 3/15
208/208 [=====] - 2s 12ms/step - loss: 0.0427 - accuracy: 0.9978
Epoch 4/15
208/208 [=====] - 2s 12ms/step - loss: 0.0031 - accuracy: 0.9995
Epoch 5/15
208/208 [=====] - 3s 13ms/step - loss: 0.0044 - accuracy: 1.0000
Epoch 6/15
208/208 [=====] - 3s 14ms/step - loss: 0.0036 - accuracy: 1.0000
Epoch 7/15
208/208 [=====] - 3s 14ms/step - loss: 0.0032 - accuracy: 1.0000
Epoch 8/15
208/208 [=====] - 2s 12ms/step - loss: 7.771e-04 - accuracy: 1.0000
Epoch 9/15
208/208 [=====] - 2s 12ms/step - loss: 5.761e-04 - accuracy: 1.0000
Epoch 10/15
208/208 [=====] - 2s 12ms/step - loss: 3.7625e-04 - accuracy: 1.0000
Epoch 11/15
208/208 [=====] - 3s 15ms/step - loss: 1.6541e-04 - accuracy: 1.0000
Epoch 12/15
208/208 [=====] - 3s 13ms/step - loss: 1.3109e-04 - accuracy: 1.0000
Epoch 13/15
208/208 [=====] - 3s 14ms/step - loss: 1.0546e-04 - accuracy: 1.0000
```

```
In [49]: model.summary()
```

```
Model: "sequential_1"
Layer (type) Output Shape Param #
-----
dense_4 (Dense) (None, 8) 136
dense_3 (Dense) (None, 1) 9
```

```
Total params: 3,547,473
Trainable params: 3,547,473
Non-trainable params: 0
```

```
In [50]: test_loss, test_acc = model.evaluate(x_test, y_test)
print('Test loss:', test_loss)
print('Test accuracy:', test_acc)
16/16 [=====] - 0s 7ms/step - loss: 0.4168 - accuracy: 0.8460
Test loss: 0.416836608817505
Test accuracy: 0.846800157356262
```

```
In [51]: plt.figure(figsize=(16,8))
plt.subplot(1,2,1)
plt.plot(history.history['loss'], color='r', label='loss')
plt.xlabel('Training Loss')
plt.ylabel('Number of epochs')
plt.legend()
plt.subplot(1,2,2)
plt.plot(history.history['accuracy'], color='b', label='accuracy')
plt.xlabel('Training Accuracy')
plt.ylabel('Number of epochs')
plt.legend()
plt.show()
```

