

traditional_ml

April 15, 2021

1 Emotion Classification in Texts using Scikit-learn

Classifying short messages into five emotion categories. We will prepare our dataset (nltk and regular expressions) and vectorize words using TF-IDF (term frequency-inverse document frequency) metric. Later we will use classifiers provided by scikit-learn and classify sentences into five emotion categories: joy, sadness, anger, fear, and neutral.

1.0.1 Workflow

- Importing Dataset
- Text Preprocessing
- Text Representation
- Classifiers: Naive Bayes, Linear Regression, Random Rorrrest, SVM
- Evaluation: F1 scores and Confussion Matrix
- Saving the Model

```
[1]: import pandas as pd
import numpy as np

# text preprocessing
from nltk import word_tokenize
from nltk.stem import PorterStemmer
from nltk.corpus import stopwords
import re

# plots and metrics
import matplotlib.pyplot as plt
from sklearn.metrics import accuracy_score, f1_score, confusion_matrix

# feature extraction / vectorization
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer

# classifiers
from sklearn.naive_bayes import MultinomialNB
from sklearn.linear_model import LogisticRegression, SGDClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import LinearSVC
from sklearn.pipeline import Pipeline
```

```
# save and load a file
import pickle
```

1.1 1. Import Dataset

Text-Emotion Dataset was split into training 70% and testing 30%

```
[2]: df_train = pd.read_csv('data/data_train.csv')
      df_test = pd.read_csv('data/data_test.csv')

      X_train = df_train.Text
      X_test = df_test.Text

      y_train = df_train.Emotion
      y_test = df_test.Emotion

      class_names = ['joy', 'sadness', 'anger', 'neutral', 'fear']
      data = pd.concat([df_train, df_test])

      print('size of training set: %s' % (len(df_train['Text'])))
      print('size of validation set: %s' % (len(df_test['Text'])))
      print(data.Emotion.value_counts())

      data.head()
```

```
size of training set: 7934
size of validation set: 3393
joy          2326
sadness      2317
anger        2259
neutral      2254
fear         2171
Name: Emotion, dtype: int64
```

```
[2]:      Emotion          Text
0  neutral  There are tons of other paintings that I thin...
1  sadness  Yet the dog had grown old and less capable , a...
2    fear   When I get into the tube or the train without ...
3    fear   This last may be a source of considerable disq...
4   anger   She disliked the intimacy he showed towards so...
```

1.1.1 *Plotting confusion matrix for later evaluation

```
[3]: def plot_confusion_matrix(y_true, y_pred, classes,
                                normalize=False,
                                title=None,
                                cmap=plt.cm.Blues):
    """
    This function prints and plots the confusion matrix.
    Normalization can be applied by setting `normalize=True`.
    """
    if not title:
        if normalize:
            title = 'Normalized confusion matrix'
        else:
            title = 'Confusion matrix, without normalization'

    # Compute confusion matrix
    cm = confusion_matrix(y_true, y_pred)

    if normalize:
        cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]

    fig, ax = plt.subplots()

    # Set size
    fig.set_size_inches(12.5, 7.5)
    im = ax.imshow(cm, interpolation='nearest', cmap=cmap)
    ax.figure.colorbar(im, ax=ax)
    ax.grid(False)

    # We want to show all ticks...
    ax.set(xticks=np.arange(cm.shape[1]),
           yticks=np.arange(cm.shape[0]),
           # ... and label them with the respective list entries
           xticklabels=classes, yticklabels=classes,
           title=title,
           ylabel='True label',
           xlabel='Predicted label')

    # Rotate the tick labels and set their alignment.
    plt.setp(ax.get_xticklabels(), rotation=45, ha="right",
              rotation_mode="anchor")

    # Loop over data dimensions and create text annotations.
    fmt = '.2f' if normalize else 'd'
    thresh = cm.max() / 2.
    for i in range(cm.shape[0]):
```

```

    for j in range(cm.shape[1]):
        ax.text(j, i, format(cm[i, j], fmt),
                ha="center", va="center",
                color="white" if cm[i, j] > thresh else "black")
fig.tight_layout()
return ax

```

1.2 2. Text Preprocessing

Here are some preprocessing steps to consider: * Removing noise: html markups, urls, non-ascii symbols, trailing whitespace etc. * Removing punctuation * Normalizing emoticons * Negation handling * Tokenization: split text into word tokens * Stopword removal * Stemming or lemmatization

However, most of these steps did not improve our classification results. Since our data was mostly taken from written dialogs it was almost ready to use.

```

[4]: def preprocess_and_tokenize(data):

    #remove html markup
    data = re.sub("<.*?>", "", data)

    #remove urls
    data = re.sub(r'http\S+', '', data)

    #remove hashtags and @names
    data= re.sub(r"#[\d\w\.]+", '', data)
    data= re.sub(r"@[\d\w\.]+", '', data)

    #remove punctuation and non-ascii digits
    data = re.sub("(\W|\d)", " ", data)

    #remove whitespace
    data = data.strip()

    # tokenization with nltk
    data = word_tokenize(data)

    # stemming with nltk
    porter = PorterStemmer()
    stem_data = [porter.stem(word) for word in data]

    return stem_data

```

1.3 3. Text Representation

Vectorizing text using Term Frequency technique (Term Frequency(TF) — Inverse Dense Frequency(IDF)) * Tekenize with our preprocess_and_tokenize * Find it's TF = (Number of repetitions

of word in a document) / (# of words in a document) * IDF = log(# of documents / # of documents containing the word)

```
[5]: # TFIDF, unigrams and bigrams
vect = TfidfVectorizer(tokenizer=preprocess_and_tokenize, sublinear_tf=True,
    →norm='l2', ngram_range=(1, 2))

# fit on our complete corpus
vect.fit_transform(data.Text)

# transform testing and training datasets to vectors
X_train_vect = vect.transform(X_train)
X_test_vect = vect.transform(X_test)
```

1.4 4. Classifiers

1.4.1 Naive Bayes

```
[6]: nb = MultinomialNB()

nb.fit(X_train_vect, y_train)

ynb_pred = nb.predict(X_test_vect)

from sklearn.metrics import accuracy_score, f1_score, confusion_matrix

print("Accuracy: {:.2f}%".format(accuracy_score(y_test, ynb_pred) * 100))
print("\nF1 Score: {:.2f}".format(f1_score(y_test, ynb_pred, average='micro') *
    →100))
print("\nCOnfusion Matrix:\n", confusion_matrix(y_test, ynb_pred))

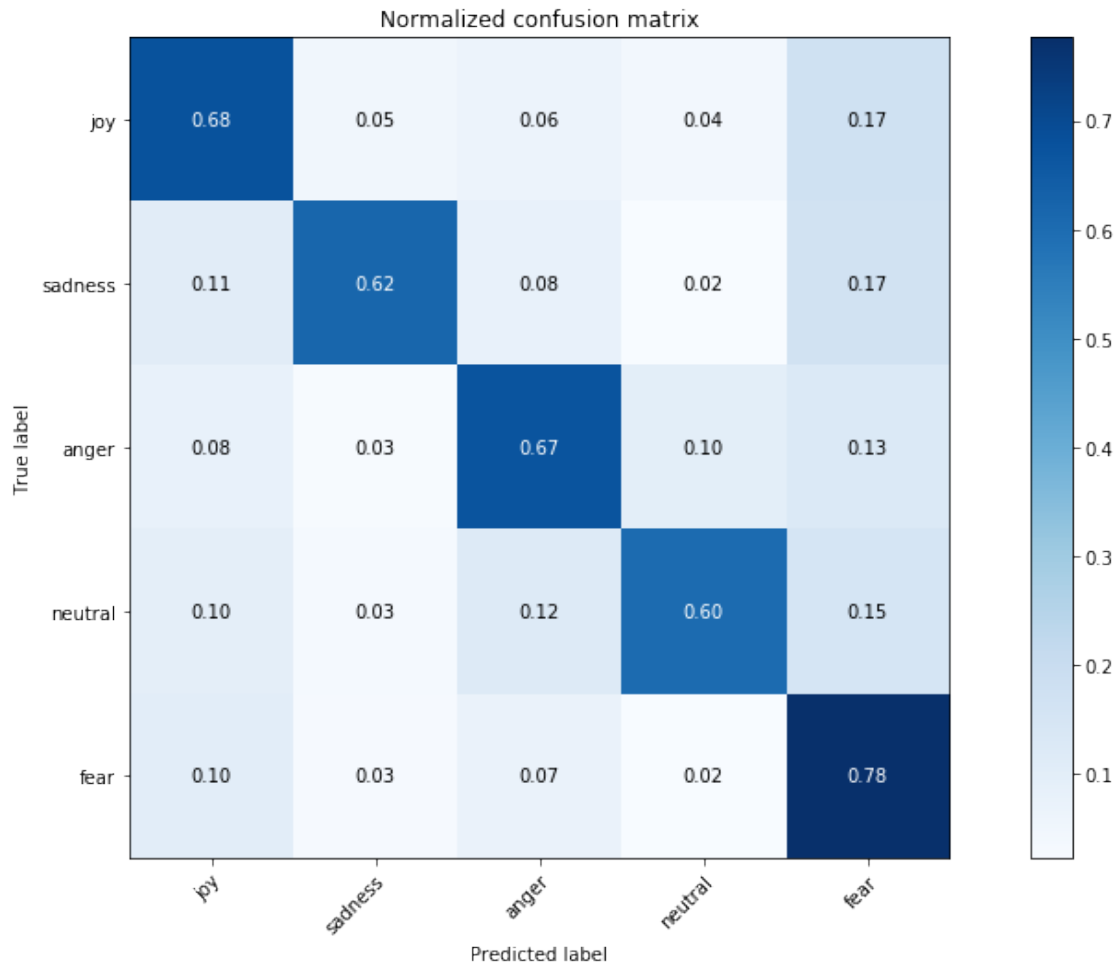
# Plot normalized confusion matrix
plot_confusion_matrix(y_test, ynb_pred, classes=class_names, normalize=True,
    →title='Normalized confusion matrix')
plt.show()
```

Accuracy: 67.02%

F1 Score: 67.02

COnfusion Matrix:

[[469	32	44	28	120]
[73	420	55	16	115]
[56	18	475	68	90]
[61	20	76	385	96]
[68	20	48	15	525]]



1.4.2 Random Forrest

```
[7]: rf = RandomForestClassifier(n_estimators=50)
      rf.fit(X_train_vect, y_train)

      yrf_pred = rf.predict(X_test_vect)

      print("Accuracy: {:.2f}%".format(accuracy_score(y_test, yrf_pred) * 100))
      print("\nF1 Score: {:.2f}".format(f1_score(y_test, yrf_pred, average='micro') * 100))
      print("\nConfusion Matrix:\n", confusion_matrix(y_test, yrf_pred))

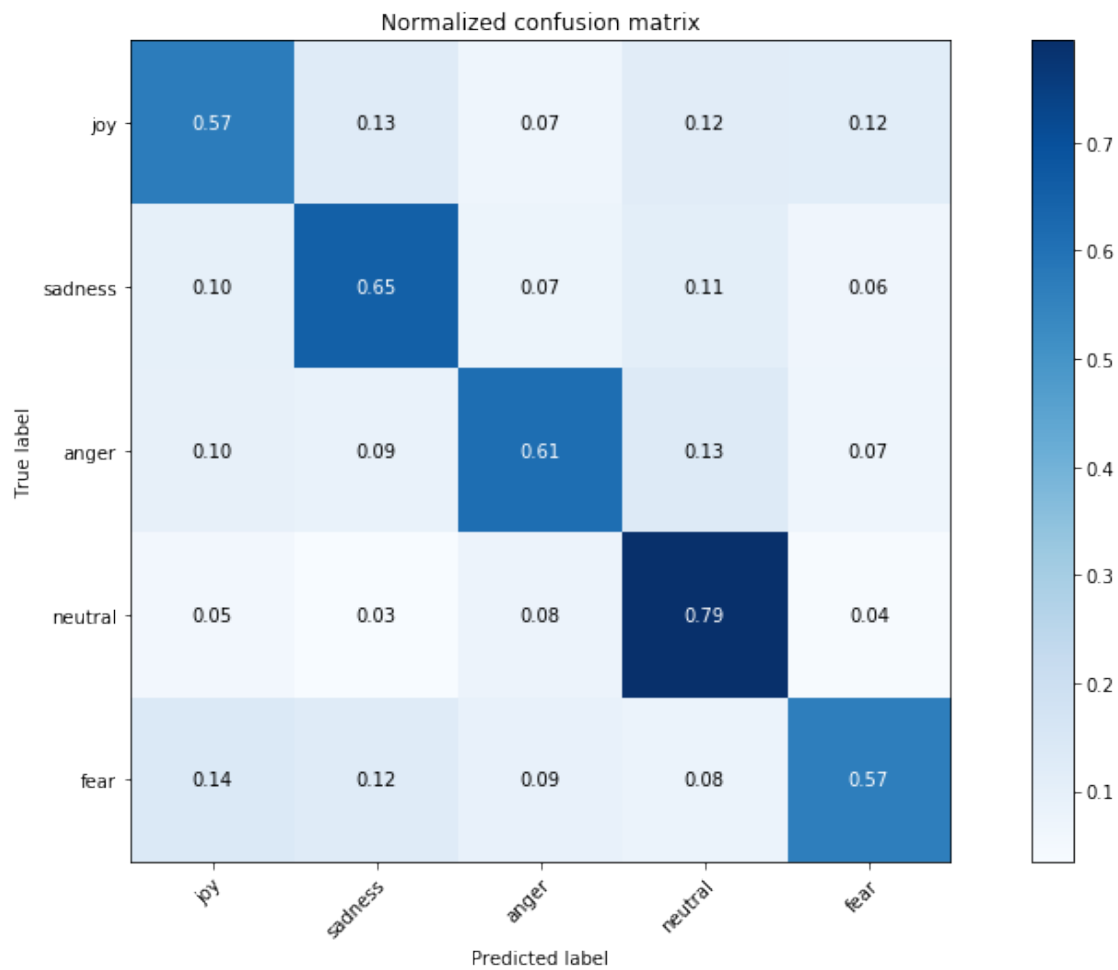
      plot_confusion_matrix(y_test, yrf_pred, classes=class_names, normalize=True,
                             title='Normalized confusion matrix')
      plt.show()
```

Accuracy: 63.72%

F1 Score: 63.72

Confusion Matrix:

```
[[396  87  47  83  80]
 [ 66 444  50  75  44]
 [ 68  62 433  95  49]
 [ 34  22  50 507  25]
 [ 94  84  62  54 382]]
```



1.4.3 Logistic Regression

```
[8]: log = LogisticRegression(solver='lbfgs', multi_class='auto', max_iter=200)
log.fit(X_train_vect, y_train)

ylog_pred = log.predict(X_test_vect)
```

```

print("Accuracy: {:.2f}%".format(accuracy_score(y_test, ylog_pred) * 100))
print("\nF1 Score: {:.2f}".format(f1_score(y_test, ylog_pred, average='micro') * 100))
print("\nConfusion Matrix:\n", confusion_matrix(y_test, ylog_pred))

plot_confusion_matrix(y_test, ylog_pred, classes=class_names, normalize=True,
    title='Normalized confusion matrix')
plt.show()

```

Accuracy: 69.35%

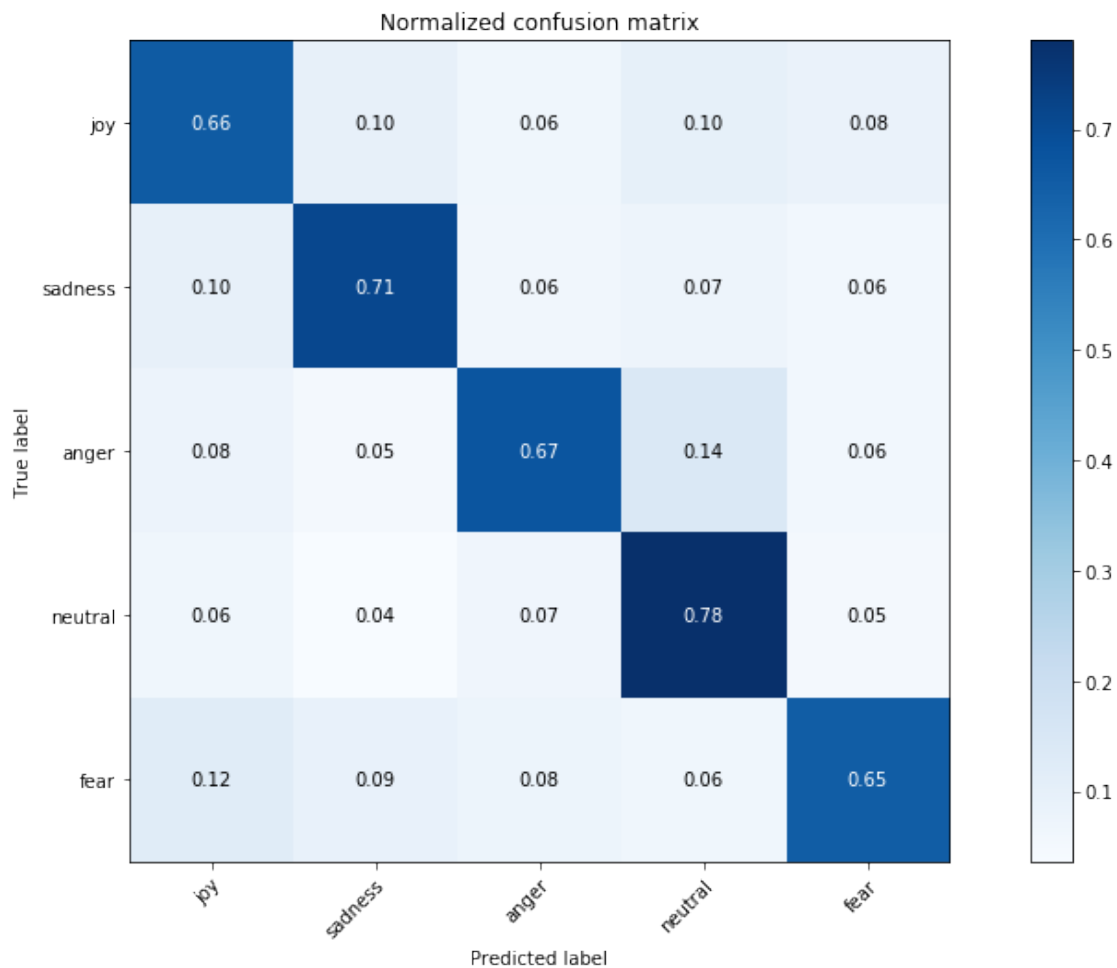
F1 Score: 69.35

Confusion Matrix:

```

[[456  67  44  68  58]
 [ 65 483  42  50  39]
 [ 56  34 476 101  40]
 [ 41  23  42 498  34]
 [ 82  60  51  43 440]]

```



1.4.4 Linear Support Vector

```
[9]: svc = LinearSVC(tol=1e-05)
      svc.fit(X_train_vect, y_train)

      ysvm_pred = svc.predict(X_test_vect)

      print("Accuracy: {:.2f}%".format(accuracy_score(y_test, ysvm_pred) * 100))
      print("\nF1 Score: {:.2f}".format(f1_score(y_test, ysvm_pred, average='micro') *
      →100))
      print("\nCOnfusion Matrix:\n", confusion_matrix(y_test, ysvm_pred))

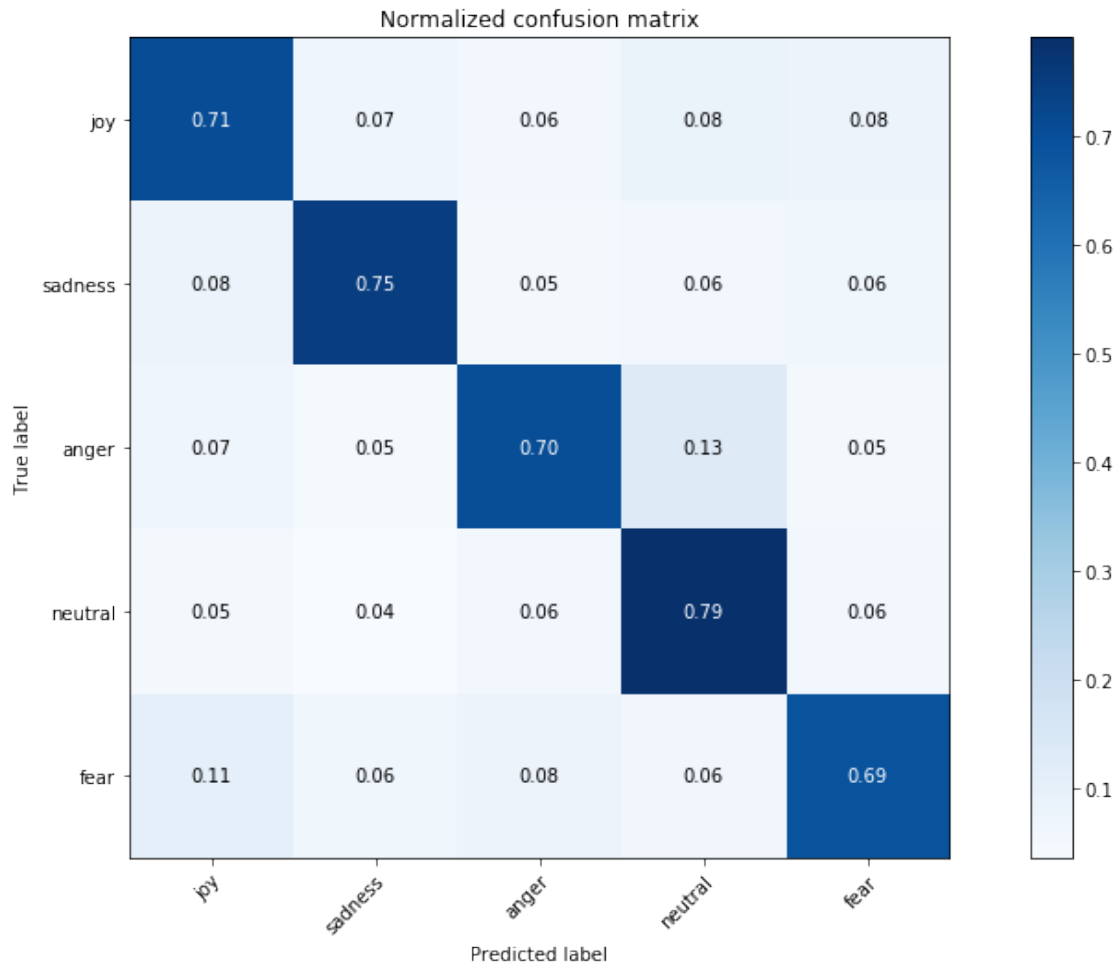
      plot_confusion_matrix(y_test, ysvm_pred, classes=class_names, normalize=True,
      →title='Normalized confusion matrix')
      plt.show()
```

Accuracy: 72.71%

F1 Score: 72.71

COnfusion Matrix:

```
[[490  49  41  58  55]
 [ 53 508  34  40  44]
 [ 50  33 498  91  35]
 [ 34  23  38 505  38]
 [ 72  43  53  42 466]]
```



1.5 4. Saving the tf-idf + SVM Model

```
[10]: #Create pipeline with our tf-idf vectorizer and LinearSVC model
svm_model = Pipeline([
    ('tfidf', vect),
    ('clf', svc),
])
```

```
[11]: # save the model
filename = 'models/tfidf_svm.sav'
pickle.dump(svm_model, open(filename, 'wb'))
```

```
[12]: model = pickle.load(open(filename, 'rb'))

message = 'delivery was hour late and my pizza is cold!'
model.predict([message])
```

```
[12]: array(['anger'], dtype=object)
```