

TaskB_SVM

November 27, 2024

```
[15]: import evaluate
import numpy as np
import pandas as pd
from nltk.tokenize import sent_tokenize, word_tokenize
from gensim.models import Word2Vec
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.model_selection import KFold
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import confusion_matrix
from sklearn import svm
from tqdm import tqdm
from keras.preprocessing import sequence
from sklearn.decomposition import PCA
from nltk.probability import FreqDist
```

```
[2]: """
Utility functions.
"""

def f1_score(tp, fp, fn):
    return (2 * tp) / (2 * tp + fp + fn)

def precision_score(tp, fp):
    return tp / (tp + fp)

def accuracy_score(tp, fp, tn, fn):
    return (tp + tn) / (tp + fp + tn + fn)

def recall_score(tp, fn):
    return tp / (tp + fn)

def flatten(matrix):
    flat_list = []
    for row in matrix:
        flat_list += row
    return flat_list
```

```
[3]: """
Download dataset SubtaskA.jsonl from
https://github.com/mbzuai-nlp/M4GT-Bench.
"""
DATA_PATH = "C:/Users/Admin/Desktop/cse847_proj/SubtaskB.jsonl"

# initialize dataset
df = pd.read_json(DATA_PATH, lines=True)
df = df[['text', 'label', 'model']]
print(df)
```

```

                                text  label  \
0      We consider a system of many polymers in solut...      2
1      We present a catalog of 66 YSOs in the Serpens...      2
2      Spectroscopic Observations of the Intermediate...      2
3      We present a new class of stochastic Lie group...      2
4      ALMA as the ideal probe of the solar chromosph...      2
...
122806  Title: The Unsung Heroes: Seagoing Cowboys and...      0
122807  Title: The Benefits of Autonomy: Student-led P...      0
122808  The Electoral College system, established by t...      0
122809  In the ever-evolving landscape of education, c...      0
122810  When faced with critical decisions, the wise o...      0

                                model
0      cohere
1      cohere
2      cohere
3      cohere
4      cohere
...
122806  gpt-3.5-turbo
122807  gpt-3.5-turbo
122808  gpt-3.5-turbo
122809  gpt-3.5-turbo
122810  gpt-3.5-turbo

[122811 rows x 3 columns]
```

```
[28]: """
Evaluate model using count/TFIDF vectorization.
"""
results = []
def run_cv(model, X, y, count_vectorizer, tfidf_transformer=None):
    k_fold = KFold(n_splits=K_FOLDS, shuffle=True, random_state=777)
    for train, test in tqdm(k_fold.split(X, y)):
        # split fold into training & testing sets
        # train = train[:10000]
```

```

# test = test[:1000]
X_train, y_train, X_test, y_test = X[train], y[train], X[test], y[test]

# fit & transform data sets
print("Count vectorizing...")
X_train = count_vectorizer.fit_transform(X_train)
X_test = count_vectorizer.transform(X_test)

if tfidf_transformer:
    print("TFIDF transforming...")
    X_train = tfidf_transformer.fit_transform(X_train)
    X_test = tfidf_transformer.transform(X_test)

# train the model
print("Fitting the model...")
model.fit(X_train, y_train)

# test the model
print("Predicting the model...")
y_hat = model.predict(X_test)

# evaluate the model
results.append(metric.compute(predictions=y_hat, references=y_test))

# analyze the run results
results_df = pd.DataFrame.from_records(results).mean()

return results_df

```

```

[29]: """
Train and evaluate SVM classifier model using count/TFIDF vectorization.
"""

# consts
MAX_FEATURES = 3000
K_FOLDS = 3
MIN_DF = 2
MAX_DF = 0.7
NGRAM_RANGE = (1, 1)
ANALYZER = 'word'

# init model
model = svm.SVC(
    verbose=True,
    max_iter=-1,
    kernel='linear',
)

```

```

# define metric
metric = evaluate.load("accuracy")

# load the data set
X = np.array(df.text)
y = np.array(df.label)

# init vectorizer and transformer
count_vectorizer = CountVectorizer(
    min_df=MIN_DF,
    max_df=MAX_DF,
    max_features=MAX_FEATURES,
    tokenizer=word_tokenize,
    token_pattern=None,
    ngram_range=NGRAM_RANGE,
    # strip_accents=STRIP_ACCENTS,
    # stop_words=STOP_WORDS,
)
tfidf_transformer = TfidfTransformer()

# run cross validation
results = run_cv(model, X, y, count_vectorizer, tfidf_transformer)
print(f"# model={model}, k_folds={K_FOLDS}, max_features={MAX_FEATURES}, \u
    \u2192min_df={MIN_DF}, max_df={MAX_DF}, "
        f"ngram_range={NGRAM_RANGE}\n{results}")

```

0it [00:00, ?it/s]

Count vectorizing...

TFIDF transforming...

Fitting the model...

[LibSVM]Predicting the model...

1it [1:27:02, 5222.69s/it]

Count vectorizing...

TFIDF transforming...

Fitting the model...

[LibSVM]Predicting the model...

2it [2:53:44, 5210.51s/it]

Count vectorizing...

TFIDF transforming...

Fitting the model...

[LibSVM]Predicting the model...

3it [4:21:07, 5222.40s/it]

```

# model=SVC(kernel='linear', verbose=True), k_folds=3, max_features=3000,
min_df=2, max_df=0.7, ngram_range=(1, 1)

```

```
accuracy    0.79992  
dtype: float64
```

```
[ ]:
```