

Classification on Unseen Data (total 8 points)

In this final task, you should read the feather file "TestQuestionsDF.feather.zstd" into a pandas dataframe. Hereafter this will be referred to as the test_set.

You can assume that the test_set is a random sample from the same dataset as "TrainQuestionsDF.feather.zstd" (hereafter train_set). Your goal is to classify the data in the test_set and achieve the best **average f1-score** using the train_set. You are allowed to utilize any technique and model available in the scikit-learn library or the standard python libraries to do so. Pay particular attention to the lessons learned from your experiments in the Classification notebook -- any of these approaches can be used to construct the model you use for prediction. You can additionally choose to generate and/or construct any features from the available data. Remember that the test_set should be constructed with the same feature space as the train_set.

For example, features based on text should be represented with the same vocabulary on the test_set as the train_set.

To achieve a high f1 score on unseen data, remember to utilize all the techniques you've learned in the lectures, lectorials and practicals.

For this task, you are expected to submit the following:

1. This notebook with your code, the code should be well documented and must run without errors. There is no time limit, but it is a good practice to save the parameters of the best model and add an option to generate a model with those parameters.
2. The submitted solutions are reproducible, i.e. the submitted code can generate the submitted prediction files (2 points).
3. The highest (out of the 4 prediction files) achieved average f1-score is in the following range:

The file names should be SXXXXXX-A2-predictions-<n>.csv - where n is a running integer (1,2,3,4).

Your mark in this task will depend on the following:

1. The code is well documented, and the entire notebook runs without errors (1 points).
 2. The submitted solutions are reproducible, i.e. the submitted code can generate the submitted prediction files (2 points).
 3. The highest (out of the 4 prediction files) achieved average f1-score is in the following range:
- (0.8, 1] (5 points)
 - (0.7, 0.8] (4 points)
 - (0.65, 0.7] (3 points)
 - (0, 0.65] (1 point)

To support the reproducibility of your solution, use the random seed anywhere where the solution involves a random process.

```
In [1]: import numpy as np
import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import ConfusionMatrixDisplay, classification_report
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
```

```
In [2]: # TODO: Any additional (if needed) import statements should be in this cell
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.preprocessing import MinMaxScaler
from sklearn.utils import shuffle
from sklearn.model_selection import GridSearchCV
from sklearn.feature_selection import SequentialFeatureSelector
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import ShuffleSplit
from sklearn.ensemble import GradientBoostingClassifier
```

```
In [3]: # TODO: Set the random seed as your student id (only numbers)
RANDOM_SEED = 2922758
np.random.seed(RANDOM_SEED)
```

```
In [4]: def read_feather_to_df(feather_file_name):
    """
    The function expects to receive a path to feather file,
    it will read the file from the disk into a pandas dataframe
    """
    return pd.read_feather(feather_file_name)
```

```
In [5]: train_df = read_feather_to_df('TrainQuestionsDF.feather.zstd')
test_df = read_feather_to_df('TestQuestionsDF.feather.zstd')
test_df.head()
```

Out[5]:	Id	PostTypeId	AcceptedAnswerId	CreationDate	Score	ViewCount	Body	OwnerUserId	LastActivityDate	Title	Answer
0	2	1	59	2010-07-19 19:12:57.167	31	30036	<p>In many different statistical methods there...	24	2017-11-22 12:15:07.030	What is normality?	
1	30	1	55	2010-07-19 19:28:34.220	13	1620	<p>Which methods are used for testing random variables...	69	2011-05-12 18:38:27.547	Testing random variate generation algorithms	
2	298	1	-1	2010-07-20 13:11:50.297	161	312967	<p>Am I looking for a better behaved distribut...	125	2016-09-21 15:41:29.603	In linear regression, when is it appropriate to...	
3	870	1	956	2010-07-28 03:54:56.447	22	26750	<p>Given a list of p-values generated from ind...	520	2013-03-07 01:26:38.717	Multiple hypothesis testing correction with Be...	
4	881	1	1189	2010-07-28 08:19:51.733	5	919	<p>Here's something I've wondered about for a...	34	2011-03-28 08:58:47.087	Series expansion of a density function	

```
In [6]: # TODO: Write your code below.
# You can split it into as many cells and functions as you see fit to make it readable and well documented.
# df = train_df.set_index('Id')
# X_train, X_test, y_train, y_test = train_test_split(df.ix[:, ~df.columns.isin(["Label"])], df.Response)
# print(X_train)
```

```
In [7]: X = train_df.drop(columns=['Label'])
y = train_df.Label
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=RANDOM_SEED, stratify=y)
```

```
In [8]: def select_numeric_non_id_columns(df):
    """
    Return a subset of a DataFrame's columns based on the column dtypes,
    including only numerical columns and excluding columns with the string idf (case-insensitive) in their name
    :return: pd.DataFrame
    """
    # TODO: write your code here
    df = df.select_dtypes(include=['Int64'])
    df = df.loc[~, ~df.columns.str.contains('Id')]
    return df

def fit_tree_classifier(X, y, **decisiontree_kwargs):
    """
    The function receives a multidimensional array or a dataframe of features as X and one dimensional array or
    It creates a DecisionTreeClassifier classifier with random_state=RANDOM_SEED, fits it on X and y and returns
    :param X: ndarray, pd.DataFrame or a sparse matrix; data features
    :param y: ndarray-like; data class labels
    :param decisiontree_kwargs: key-word arguments that will be passed to DecisionTreeClassifier class
    :return: a fitted DecisionTreeClassifier object
    """
    # TODO: write your code here
    clf = DecisionTreeClassifier(**decisiontree_kwargs, random_state=RANDOM_SEED)
    clf = clf.fit(X, y)
    return clf

def fit_knn_classifier(X, y, **knn_kwargs):
    """
    The function receives a multidimensional array or a dataframe of features as X, on dimensional array or pandas
    It then creates a DecisionTreeClassifier classifier with random_state=RANDOM_SEED, fits it on X and y and returns
    :param X: ndarray, pd.DataFrame or a sparse matrix; data features
    :param y: array-like; data class labels
    :param knn_kwargs: key-word arguments that will be passed to KNeighborsClassifier class
    :return: a fitted KNeighborsClassifier object
    """
    # TODO: write your code here
    neigh = KNeighborsClassifier(**knn_kwargs)
    neigh = neigh.fit(X, y)
    return neigh

def evaluate_classification(y_true, y_predicted):
    """
    The function receives two arrays or pandas Series with the same length, the actual labels and the predicted
    It then prints the sklearn classification_report and plots a confusion matrix as a heatmap using the class
    :param y_true: array-like; ground truth data class labels
    :param y_predicted: array-like; predicted data class labels
    """
    # TODO: write your code here
    print(classification_report(y_true, y_predicted))
    mat=ConfusionMatrixDisplay.from_predictions(y_true, y_predicted, xticks_rotation=40)
    return
```

```
def fit_series_to_tfidf(sr, **tfidfvectorizer_kwargs):
    """
    The function receives a multidimensional array or a pandas Series that contains text strings (a.k.a documents).
    It then converts the documents into a matrix of TF-IDF features
    The function should return two objects: (fitted) the vocabulary and idf from the training set,
    TfidfVectorizer object after it learned (fitted) the vocabulary and idf from the training set,
    and a document-term matrix (the original documents array transformed into a TF-IDF features matrix).
    :param sr: pd.Series, contains text strings
    :param tfidfvectorizer_kwargs: key-word arguments that will be passed to TfidfVectorizer class
    :return: two objects, the fitted TfidfVectorizer object and the tf-idf document-term sparse matrix
    """
    # TODO: write your code here
    vectorizer = TfidfVectorizer(**tfidfvectorizer_kwargs)#max_features=100,
                                #min_df=5,
                                #stop_words='english'
    vectors = vectorizer.fit_transform(sr)
    return vectorizer, vectors

In [15]: ##text feature extraction with classifiers
#Random Forest with tfidf, looking for best parameters
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', RandomForestClassifier())])

parameters = {
    'tfidf__max_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': ((1, 1), (1, 2), (1, 3)),
    'tfidf__min_df': np.arange(5, 10, 1),
    'clf__max_depth': np.arange(1, 10, 1),
    'tfidf__stop_words': ['english'],
    'clf__random_state': [RANDOM_SEED],
    'clf__n_estimators': [50, 100]
    #'clf__min_samples_split': np.arange(1, 4, 1),
    #'tfidf__use_idf': (True, False) #means we score only with TF not IDF to-interesting this improved perform
}

grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train, y_train)
print("Best Tree TF-IDF score and parameters:")
print(grid_search_tune.best_score_)
print(grid_search_tune.best_params_)
```


[illegible]

