

# Classification (14 points):

Each function is worth 2 points, except for 8 which is worth 4 points.

For each question that asks you to implement a function, implement it in the top cell where it is defined and then execute the function in the code cell provided below the question.

You should base your answers on the output.

You are allowed to implement and use additional functions. These would be defined and implemented in the cell directly below the questions they were implemented for.

All the textual answers should be based on and justified with output from the data in the code cell above.

For example, if the question asks about the correlation value, the code calculating it should appear above the answer, and the value should be in the output. The answers should be concise and written in your own words.

**Do Not Modify the Structure of this Notebook, don't add/remove/move cells or change their type (Code/Markdown)**

1. Read the feather file 'TrainQuestionsDFfeather.zstd' into a pandas dataframe.  
Use the function `train_test_split` to split the data into two sets, 75% for training and 25% for validation.  
Generate stratified samples with the `random_state=RANDOM_SEED`.
2. Implement the function `fit_tree_classifier(X, y, **decisiontree_kwargs)`, then use it to fit a decision tree classifier on the train dataset and generate data predictions. Display the predictions as a heatmap using the class.  
Use only the numerical columns as features (you can use the function from DataExploration).
3. Implement the function `evaluate_classification(y_true, y_predicted)`, then use it to evaluate the classification made by the decision tree classifier on the validation dataset.
4. Implement the function `fit_knn_classifier(X, y)` and use it to fit the model on the train data and then generate prediction for the validation data. Using the previous evaluation.
5. Now we turn to a different features set, we will utilize the text in the Title field of each sample to generate a features vector for the sample.  
You should apply the `TfidfVectorizer` to generate `tf-idf` (term frequency-inverse document frequency) features from the text in the Title field of each sample.  
Make sure to use the same vocabulary for both the training set and the validation. The vocabulary size determines the vector size, each entry in the vector represents the tf-idf value for the corresponding term.  
Implement the function `series_to_tfidf(sr, **tfidfvectorizer_kwargs)`, then generate `tf-idf` vectors for the training and validation sets and train two classifiers (decision tree and KNN) using the generated vectors.
6. Using the documentation for the two classification and the text feature extraction models and their different parameters. Find a combination of parameters that improves the accuracy for each model, and for at least one of the models the improvement should be by at least 5% on the validation dataset.  
Write your code below and describe the changes you made and the intuition behind them.  
For applying a systematic search, i.e. not just manually checking for different parameters.

```
In [1]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import 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 importMinMaxScaler
from sklearn.utils import shuffle
from sklearn.model_selection import GridSearchCV
from sklearn.feature_selection import SequentialFeatureSelector
from sklearn.metrics import roc_auc_score
```

```
In [3]: # TODO: Set the random seed as your student id (only numbers)
RANDOM_SEED = 395750
np.random.seed(RANDOM_SEED)
```

```
In [4]: # This cell is for functions given to you to use
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
    :param feather_file_name: a string or path like object
    :return: pandas dataframe
    """
    return pd.read_feather(feather_file_name)
```

```
In [5]: # This cell is for all the functions you are expected to implement.
# You should implement them here and only call them below when they are mentioned in a question.
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 id (case-insensitive) in their name
    :param df: pandas dataframe
    :return: pandas dataframe
    """
```

```
    df = df.select_dtypes(include=[int, float])
    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 pandas series as y.
    It creates a DecisionTreeClassifier classifier with random_state=RANDOM_SEED, fits it on X and y and return
    :param X: ndarray, pandas dataframe or a sparse matrix; data features
    :param y: array-like; target class labels
    :param decisiontree_kwargs: keyword arguments that will be passed to DecisionTreeClassifier class
    :return: a fitted DecisionTreeClassifier object
```

```
    # TODO: write your code here
    clf = DecisionTreeClassifier(random_state=RANDOM_SEED)
    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 series as y.
    It creates a KNeighborsClassifier classifier with random_state=RANDOM_SEED, fits it on X and y and return
    :param X: ndarray, pandas dataframe or a sparse matrix; data features
    :param y: array-like; target class labels
    :param knn_kwargs: keyword arguments that will be passed to KNeighborsClassifier class
    :return: a fitted KNeighborsClassifier object
```

```
    # TODO: write your code here
    neigh = KNeighborsClassifier(n_neighbors=5, weights='distance')
    neigh = fit(X, y)
    return neigh

def evaluate_classification(y_true, y_predicted):
    """
    The function receives two arrays of pandas series with the same length, the actual labels and the predicted labels.
    It prints the classification report and plots a confusion matrix as a heatmap using the class
    The plot should be readable (e.g. not overlapping labels or too small text)
    :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))
    plt.figure(figsize=(10, 10))
    cm = confusion_matrix(y_true, y_predicted, labels=sorted(set(y_true).union(y_predicted)))
    return cm

def series_to_tfidf(sr, **tfidfvectorizer_kwargs):
    """
    The function receives an array or a pandas series that contains text strings (a.k.a. documents).
    It then converts it to a matrix of TF-IDF features
    The function should return two objects:
    TfidfVectorizer object after it learned (fitted) the vocabulary and idf from the training set,
    :param sr: pandas series contains text strings
    :param tfidfvectorizer_kwargs: keyword 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(max_features=100,
                                min_df=5,
                                min_idf=5,
                                stop_words='english')
    vectors = vectorizer.fit_transform(sr)
    return vectorizer, vectors
```

1. Read the feather file 'TrainQuestionsDFfeather.zstd' into a pandas dataframe.  
Use the function `train_test_split` to split the data into two sets, 75% for training and 25% for validation.  
Generate stratified samples with the `random_state=RANDOM_SEED`.

```
In [6]: # TODO: write your function calls and code here
df = read_feather_to_df('TrainQuestionsDFfeather.zstd')
X = df.drop(columns=['label'])
y = 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)
```

2. Implement the function `fit_tree_classifier(X, y, **decisiontree_kwargs)`, then use it to fit a decision tree classifier on the train dataset and generate prediction for the validation data. Using the previous evaluation.  
Use only the numerical columns as features and print the labels of the first 5 predictions (you can use the function from DataExploration).

```
In [7]: # TODO: write your function calls and code here
clf = fit_tree_classifier(select_numeric_non_id_columns(X_train), y_train, random_state=RANDOM_SEED, max_depth=10)
y_pred = clf.predict(X_test)
print(classification_report(y_test, y_pred))
```

```
precision    recall    f1-score   support

bayesian - 0.15  0.16  0.16    751
distributions - 0.15  0.16  0.16    750
hypothesis-testing - 0.17  0.17  0.17    751
logistic - 0.17  0.17  0.17    750
probability - 0.15  0.15  0.15    751
self-study - 0.18  0.16  0.17    750
time-series - 0.16  0.13  0.14    750
```

```
accuracy - 0.16  0.16  0.16    5253
macro avg - 0.16  0.16  0.16    5253
weighted avg - 0.16  0.16  0.16    5253
```

```
bayesian - 130 117 104 92 108 95 105
distributions - 129 128 111 100 115 82 96
hypothesis-testing - 125 115 129 109 103 84 86
logistic - 133 104 117 124 86 90 96
probability - 115 129 115 93 113 98 88
self-study - 143 120 137 101 108 120 75
time-series - 124 103 111 128 93 90 101
```

3. Implement the function `evaluate_classification(y_true, y_predicted)`, then use it to evaluate the classification made by the decision tree classifier on the validation dataset.

3.1 For which label did the model achieve the best result, and how many samples were classified correctly for that label?  
3.2 How many samples with the label 'bayesian' were classified as 'time-series'?  
3.3 Was the model successful?

```
In [8]: # TODO: write your function calls and code here
evaluate_classification(y_test, y_pred)
```

```
precision    recall    f1-score   support

bayesian - 0.15  0.16  0.16    751
distributions - 0.15  0.16  0.16    750
hypothesis-testing - 0.17  0.17  0.17    751
logistic - 0.17  0.17  0.17    750
probability - 0.15  0.15  0.15    751
self-study - 0.18  0.16  0.17    750
time-series - 0.16  0.13  0.14    750
```

```
accuracy - 0.16  0.16  0.16    5253
macro avg - 0.16  0.16  0.16    5253
weighted avg - 0.16  0.16  0.16    5253
```

```
bayesian - 130 117 104 92 108 95 105
distributions - 129 128 111 100 115 82 96
hypothesis-testing - 125 115 129 109 103 84 86
logistic - 133 104 117 124 86 90 96
probability - 115 129 115 93 113 98 88
self-study - 143 120 137 101 108 120 75
time-series - 124 103 111 128 93 90 101
```

The model achieved the best results for Bayesian label, correctly predicting 896 samples. The model incorrectly classified 14 samples as 'time-series' which were actually bayesian. With an overall model accuracy of 0.84, this model was reasonably successful at predicting labels.

4. Implement the function `fit_knn_classifier(X, y)` and use it to fit the model on the train data and then generate prediction for the validation data. Using the previous evaluation answer the following questions:

4.1 Which model achieved higher accuracy on the validation set?  
4.2 Can you identify a bias towards a certain class in the result?

```
In [9]: # TODO: write your function calls and code here
clf = fit_knn_classifier(select_numeric_non_id_columns(X_train), y_train)
y_pred_knn = neigh.predict(X_test)
evaluate_classification(y_test, y_pred_knn)
```

```
precision    recall    f1-score   support

bayesian - 0.14  0.18  0.16    751
distributions - 0.17  0.18  0.17    750
hypothesis-testing - 0.16  0.16  0.16    751
logistic - 0.15  0.15  0.15    750
probability - 0.15  0.15  0.15    751
self-study - 0.17  0.15  0.16    750
time-series - 0.17  0.16  0.16    750
```

```
accuracy - 0.16  0.16  0.16    5253
macro avg - 0.16  0.16  0.16    5253
weighted avg - 0.16  0.16  0.16    5253
```

```
bayesian - 136 111 93 108 104 96 103
distributions - 132 135 112 106 92 84 89
hypothesis-testing - 134 121 120 102 87 85 102
logistic - 142 105 108 110 86 100 99
probability - 129 127 110 104 110 81 90
self-study - 136 104 112 103 104 109 76
time-series - 152 109 112 91 97 94 117
```

The fit\_tree\_classifier model achieved a higher accuracy than the KNN on the validation set. There is a bias towards the 'bayesian' class in both models, with that class being predicted the most.

5. Now we turn to a different features set, we will utilize the text in the Title field of each sample to generate a features vector for the sample.  
You should apply the `TfidfVectorizer` to generate `tf-idf` (term frequency-inverse document frequency) features from the text in the Title field of each sample.  
Make sure to use the same vocabulary for both the training set and the validation. The vocabulary size determines the vector size, each entry in the vector represents the tf-idf value for the corresponding term.  
Implement the function `series_to_tfidf(sr, **tfidfvectorizer_kwargs)`, then generate `tf-idf` vectors for the training and validation sets and train two classifiers (decision tree and KNN) using the generated vectors. Answer the following questions:

5.1 Which model achieves higher accuracy on the validation set?  
5.2 For each model specify the label it gets most correct and most incorrect prediction for.

```
In [10]: # TODO: write your function calls and code here
tfidf = TfidfVectorizer(stop_words='english')
tfidf_train = tfidf.fit_transform(X_train['title'])
tfidf_test = tfidf.transform(X_test['title'])
clf = fit_tree_classifier(tfidf_train, y_train)
y_pred_tfidf_tree = clf.predict(tfidf_test)
evaluate_classification(y_test, y_pred_tfidf_tree)
clf = fit_knn_classifier(tfidf_train, y_train)
y_pred_tfidf_knn = clf.predict(tfidf_test)
evaluate_classification(y_test, y_pred_tfidf_knn)
```

```
precision    recall    f1-score   support

bayesian - 0.73  0.64  0.68    751
distributions - 0.71  0.56  0.67    750
hypothesis-testing - 0.65  0.65  0.65    751
logistic - 0.76  0.70  0.73    750
probability - 0.54  0.47  0.50    751
self-study - 0.38  0.42  0.35    750
time-series - 0.60  0.60  0.62    750
```

```
accuracy - 0.58  0.58  0.58    5253
macro avg - 0.60  0.58  0.59    5253
weighted avg - 0.60  0.58  0.59    5253
```

```
precision    recall    f1-score   support

bayesian - 0.63  0.58  0.61    751
distributions - 0.71  0.54  0.62    750
hypothesis-testing - 0.64  0.64  0.64    751
logistic - 0.72  0.68  0.70    750
probability - 0.53  0.46  0.49    751
self-study - 0.31  0.42  0.36    750
time-series - 0.71  0.57  0.63    750
```

```
accuracy - 0.58  0.56  0.56    5253
macro avg - 0.58  0.56  0.56    5253
weighted avg - 0.58  0.56  0.56    5253
```

```
bayesian - 478 33 22 16 49 115 38
distributions - 34 422 45 25 76 114 34
hypothesis-testing - 13 48 491 35 22 92 50
logistic - 22 13 38 528 18 91 40
probability - 37 84 30 18 356 193 33
self-study - 62 100 75 40 111 317 55
time-series - 17 28 53 33 29 139 451
```

The fit\_knn\_classifier model achieved a higher accuracy than the KNN on the validation set. There is a bias towards the 'bayesian' class in both models, with that class being predicted the most.

6. Using the documentation for the two classification and the text feature extraction models and their different parameters. Find a combination of parameters that improves the accuracy for each model, and for at least one of the models the improvement should be by at least 5% on the validation dataset.  
Write your code below and describe the changes you made and the intuition behind them.  
For applying a systematic search, i.e. not just manually checking different values for different parameters.

```
In [11]: # TODO: write your function calls and code here
# Decision tree tuning
params = {
    'max_depth': [2, 3, 5, 7, 9, 10, 12],
    'min_samples_leaf': [15, 10, 20, 50, 100, 150, 160, 175, 185, 209, 210, 212],
    'criterion': ['gini', 'entropy']
}
```

```
grid_search = GridSearchCV(estimator=fit_tree_classifier(select_numeric_non_id_columns(X_train), y_train),
                           param_grid=params, cv=4, n_jobs=-1, verbose=1, scoring='accuracy')
grid_search.fit(select_numeric_non_id_columns(X_train), y_train)
print('Best Tree score and parameters:')
print(grid_search.best_params_)
```

```
# KNN tuning
grid_params = {
    'n_neighbors': [33, 39, 40, 41, 42, 43, 49, 50, 55],
    'weights': ['uniform', 'distance'],
    'metric': ['minkowski', 'euclidean', 'manhattan']
}
```

```
gs = GridSearchCV(estimator=fit_knn_classifier(select_numeric_non_id_columns(X_train), y_train),
                  param_grid=grid_params, cv=4, n_jobs=-1, verbose=1)
grid_params = gs.best_params_
print('Best KNN score and parameters:')
print(grid_search.best_params_)
```

```
# Text feature extraction with classifiers
from sklearn.pipeline import Pipeline
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', DecisionTreeClassifier())
])
```

```
parameters = {
    'tfidf__min_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
    'tfidf__max_df': np.arange(1, 10, 1),
    'clf__min_df': np.arange(1, 10, 1),
    'tfidf__stop_words': ['english', 'none'],
    'tfidf__use_idf': (True, False)
}
```

```
grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train['title'], y_train)
print('Best KNN TF-IDF score and parameters:')
print(grid_search_tune.best_params_)
```

```
# KNN with tfidf prev accuracy = 0.56
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', KNeighborsClassifier())
])
```

```
parameters = {
    'tfidf__min_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
    'tfidf__max_df': np.arange(1, 42, 2),
    'tfidf__stop_words': ['english', 'none'],
    'tfidf__use_idf': (True, False)
}
```

```
grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train['title'], y_train)
print('Best KNN TF-IDF score and parameters:')
print(grid_search_tune.best_params_)
```

```
# Decision tree tuning
params = {
    'max_depth': [2, 3, 5, 7, 9, 10, 12],
    'min_samples_leaf': [15, 10, 20, 50, 100, 150, 160, 175, 185, 209, 210, 212],
    'criterion': ['gini', 'entropy']
}
```

```
grid_search = GridSearchCV(estimator=fit_tree_classifier(select_numeric_non_id_columns(X_train), y_train),
                           param_grid=params, cv=4, n_jobs=-1, verbose=1, scoring='accuracy')
grid_search.fit(select_numeric_non_id_columns(X_train), y_train)
print('Best Tree score and parameters:')
print(grid_search.best_params_)
```

```
# KNN tuning
grid_params = {
    'n_neighbors': [33, 39, 40, 41, 42, 43, 49, 50, 55],
    'weights': ['uniform', 'distance'],
    'metric': ['minkowski', 'euclidean', 'manhattan']
}
```

```
gs = GridSearchCV(estimator=fit_knn_classifier(select_numeric_non_id_columns(X_train), y_train),
                  param_grid=grid_params, cv=4, n_jobs=-1, verbose=1)
grid_params = gs.best_params_
print('Best KNN score and parameters:')
print(grid_search.best_params_)
```

```
# Text feature extraction with classifiers
from sklearn.pipeline import Pipeline
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', DecisionTreeClassifier())
])
```

```
parameters = {
    'tfidf__min_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
    'tfidf__max_df': np.arange(1, 10, 1),
    'clf__min_df': np.arange(1, 10, 1),
    'tfidf__stop_words': ['english', 'none'],
    'tfidf__use_idf': (True, False)
}
```

```
grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train['title'], y_train)
print('Best KNN TF-IDF score and parameters:')
print(grid_search_tune.best_params_)
```

```
# KNN with tfidf prev accuracy = 0.56
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', KNeighborsClassifier())
])
```

```
parameters = {
    'tfidf__min_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
    'tfidf__max_df': np.arange(1, 42, 2),
    'tfidf__stop_words': ['english', 'none'],
    'tfidf__use_idf': (True, False)
}
```

```
grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train['title'], y_train)
print('Best KNN TF-IDF score and parameters:')
print(grid_search_tune.best_params_)
```

```
# Decision tree tuning
params = {
    'max_depth': [2, 3, 5, 7, 9, 10, 12],
    'min_samples_leaf': [15, 10, 20, 50, 100, 150, 160, 175, 185, 209, 210, 212],
    'criterion': ['gini', 'entropy']
}
```

```
grid_search = GridSearchCV(estimator=fit_tree_classifier(select_numeric_non_id_columns(X_train), y_train),
                           param_grid=params, cv=4, n_jobs=-1, verbose=1, scoring='accuracy')
grid_search.fit(select_numeric_non_id_columns(X_train), y_train)
print('Best Tree score and parameters:')
print(grid_search.best_params_)
```

```
# KNN tuning
grid_params = {
    'n_neighbors': [33, 39, 40, 41, 42, 43, 49, 50, 55],
    'weights': ['uniform', 'distance'],
    'metric': ['minkowski', 'euclidean', 'manhattan']
}
```

```
gs = GridSearchCV(estimator=fit_knn_classifier(select_numeric_non_id_columns(X_train), y_train),
                  param_grid=grid_params, cv=4, n_jobs=-1, verbose=1)
grid_params = gs.best_params_
print('Best KNN score and parameters:')
print(grid_search.best_params_)
```

```
# Text feature extraction with classifiers
from sklearn.pipeline import Pipeline
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', DecisionTreeClassifier())
])
```

```
parameters = {
    'tfidf__min_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
    'tfidf__max_df': np.arange(1, 10, 1),
    'clf__min_df': np.arange(1, 10, 1),
    'tfidf__stop_words': ['english', 'none'],
    'tfidf__use_idf': (True, False)
}
```

```
grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train['title'], y_train)
print('Best KNN TF-IDF score and parameters:')
print(grid_search_tune.best_params_)
```

```
# KNN with tfidf prev accuracy = 0.56
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', KNeighborsClassifier())
])
```

```
parameters = {
    'tfidf__min_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
    'tfidf__max_df': np.arange(1, 42, 2),
    'tfidf__stop_words': ['english', 'none'],
    'tfidf__use_idf': (True, False)
}
```

```
grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train['title'], y_train)
print('Best KNN TF-IDF score and parameters:')
print(grid_search_tune.best_params_)
```

```
# Decision tree tuning
params = {
    'max_depth': [2, 3, 5, 7, 9, 10, 12],
    'min_samples_leaf': [15, 10, 20, 50, 100, 150, 160, 175, 185, 209, 210, 212],
    'criterion': ['gini', 'entropy']
}
```

```
grid_search = GridSearchCV(estimator=fit_tree_classifier(select_numeric_non_id_columns(X_train), y_train),
                           param_grid=params, cv=4, n_jobs=-1, verbose=1, scoring='accuracy')
grid_search.fit(select_numeric_non_id_columns(X_train), y_train)
print('Best Tree score and parameters:')
print(grid_search.best_params_)
```

```
# KNN tuning
grid_params = {
    'n_neighbors': [33, 39, 40, 41, 42, 43, 49, 50, 55],
    'weights': ['uniform', 'distance'],
    'metric': ['minkowski', 'euclidean', 'manhattan']
}
```

```
gs = GridSearchCV(estimator=fit_knn_classifier(select_numeric_non_id_columns(X_train), y_train),
                  param_grid=grid_params, cv=4, n_jobs=-1, verbose=1)
grid_params = gs.best_params_
print('Best KNN score and parameters:')
print(grid_search.best_params_)
```

```
# Text feature extraction with classifiers
from sklearn.pipeline import Pipeline
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', DecisionTreeClassifier())
])
```

```
parameters = {
    'tfidf__min_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
    'tfidf__max_df': np.arange(1, 10, 1),
    'clf__min_df': np.arange(1, 10, 1),
    'tfidf__stop_words': ['english', 'none'],
    'tfidf__use_idf': (True, False)
}
```

```
grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train['title'], y_train)
print('Best KNN TF-IDF score and parameters:')
print(grid_search_tune.best_params_)
```

```
# KNN with tfidf prev accuracy = 0.56
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', KNeighborsClassifier())
])
```

```
parameters = {
    'tfidf__min_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
    'tfidf__max_df': np.arange(1, 42, 2),
    'tfidf__stop_words': ['english', 'none'],
    'tfidf__use_idf': (True, False)
}
```

```
grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train['title'], y_train)
print('Best KNN TF-IDF score and parameters:')
print(grid_search_tune.best_params_)
```

```
# Decision tree tuning
params = {
    'max_depth': [2, 3, 5, 7, 9, 10, 12],
    'min_samples_leaf': [15, 10, 20, 50, 100, 150, 160, 175, 185, 209, 210, 212],
    'criterion': ['gini', 'entropy']
}
```

```
grid_search = GridSearchCV(estimator=fit_tree_classifier(select_numeric_non_id_columns(X_train), y_train),
                           param_grid=params, cv=4, n_jobs=-1, verbose=1, scoring='accuracy')
grid_search.fit(select_numeric_non_id_columns(X_train), y_train)
print('Best Tree score and parameters:')
print(grid_search.best_params_)
```

```
# KNN tuning
grid_params = {
    'n_neighbors': [33, 39, 40, 41, 42, 43, 49, 50, 55],
    'weights': ['uniform', 'distance'],
    'metric': ['minkowski', 'euclidean', 'manhattan']
}
```

```
gs = GridSearchCV(estimator=fit_knn_classifier(select_numeric_non_id_columns(X_train), y_train),
                  param_grid=grid_params, cv=4, n_jobs=-1, verbose=1)
grid_params = gs.best_params_
print('Best KNN score and parameters:')
print(grid_search.best_params_)
```

```
# Text feature extraction with classifiers
from sklearn.pipeline import Pipeline
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', DecisionTreeClassifier())
])
```

```
parameters = {
    'tfidf__min_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
    'tfidf__max_df': np.arange(1, 10, 1),
    'clf__min_df': np.arange(1, 10, 1),
    'tfidf__stop_words': ['english', 'none'],
    'tfidf__use_idf': (True, False)
}
```

```
grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train['title'], y_train)
print('Best KNN TF-IDF score and parameters:')
print(grid_search_tune.best_params_)
```

```
# KNN with tfidf prev accuracy = 0.56
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', KNeighborsClassifier())
])
```

```
parameters = {
    'tfidf__min_df': (0.25, 0.5, 0.75),
    'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
    'tfidf__max_df': np.arange(1, 42, 2),
    'tfidf__stop_words': ['english', 'none'],
    'tfidf__use_idf': (True, False)
}
```

```
grid_search_tune = GridSearchCV(pipeline, parameters, cv=2, n_jobs=-1, verbose=3)
grid_search_tune.fit(X_train['title'], y_train)
print('Best KNN TF-IDF score and parameters:')
print(grid_search_tune.best_params_)
```