

Final Project

Bryan Clark & Vamsi Ponnappalli

12/10/2020

Contents

Overview	1
Python Libraries	1
1. Fetch Data	2
2. Preprocess Data	3
2.1 Train/Test Split	3
2.2 Cross-Validation	3
3. Neural Networks	4
3.1 Baseline Model Performance	5
3.2 Deep Learning	6
4. Evaluation	8

Overview

In this term project, you will deploy Deep Learning models to build a classification model using RapidMiner to predict the sentiment of consumers towards US airlines based on their reviews expressed in the form of tweets. If you strongly prefer to use some other DL-based software/frameworks instead of RapidMiner, such as TensorFlow or PyTorch, let me know before starting the work. This is a group project, and you should work on it in the groups that you have formed already.

Python Libraries

```
# fetch data
import requests
from io import StringIO

# core
import numpy as np
import pandas as pd

# preprocessing
from sklearn.model_selection import train_test_split
from sklearn.model_selection import KFold

# baseline algorithms
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
```

```
# deep learning
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, BatchNormalization, Dropout
from tensorflow.keras.losses import mean_absolute_error as tf_mae
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.callbacks import EarlyStopping

# evaluation
from sklearn.model_selection import cross_validate
from sklearn.metrics import mean_absolute_error as skl_mae
```

1. Fetch Data

The data is provided to you in two versions:

1. The original version of the tweets (and their sentiments) is located at <https://drive.google.com/file/d/1atyRH5Yz7TU-2ziyZknfd7ib2LLwYeuv/view?usp=sharing>.
2. The preprocessed version of the tweets is located at <https://drive.google.com/file/d/1c96crlNZr7XiF3-9lmZ1nEJaY3MHTTz5/view?usp=sharing>, where text preprocessing and pre-training of the text embeddings of the tweets using autoencoders have already been done to make your life simpler. This preprocessed version contains the sentiments about the tweets in column 1 of the spreadsheet (either positive (1) or negative (0)) and the 8-dimensional pre-trained embeddings of the tweets (in columns 2 – 9 of the spreadsheet).

I recommend that you use the preprocessed version of the tweets since it will save you a lot of preprocessing work to build these embeddings that is non-trivial. However, if you like challenges, you can do preprocessing and building the embeddings using autoencoders yourself and, therefore, work directly with the “raw” tweets. As a “reward” for this extra work, you will be awarded 10 extra points (the max score of this project is 100) if you preprocess tweets yourself.

```
# url of dataset
google_drive_file_url = 'https://drivesds.google.com/file/d/' \
    + '1c96crlNZr7XiF3-9lmZ1nEJaY3MHTTz5/view?usp=sharing'

def fetch_google_drive_csv(google_drive_file_url):

    file_id = google_drive_file_url.split('/')[2]
    download_url = 'https://drive.google.com/uc?export=download&id=' + file_id
    url = requests.get(download_url)
    csv_raw = StringIO(url.text)
    return pd.read_csv(csv_raw)

data = fetch_google_drive_csv(google_drive_file_url)

data.head()
```

	sentiment	dimension1	dimension2	...	dimension6	dimension7	dimension8
## 0	1	-0.400418	0.293417	...	-0.034476	0.042176	-0.429317
## 1	1	-0.454608	-0.194998	...	0.064868	0.072154	0.629457
## 2	0	-0.515892	-0.120781	...	-0.155029	-0.306803	0.694974
## 3	1	0.047770	-0.230509	...	-0.229259	-0.835945	0.294148
## 4	0	-0.574353	-0.132517	...	-0.338480	0.202040	-0.100443
##							
##	[5 rows x 9 columns]						

```
data.info()

## <class 'pandas.core.frame.DataFrame'>
## RangeIndex: 55524 entries, 0 to 55523
## Data columns (total 9 columns):
##  #   Column      Non-Null Count  Dtype
## ---  -
##  0   sentiment    55524 non-null  int64
##  1   dimension1    55524 non-null  float64
##  2   dimension2    55524 non-null  float64
##  3   dimension3    55524 non-null  float64
##  4   dimension4    55524 non-null  float64
##  5   dimension5    55524 non-null  float64
##  6   dimension6    55524 non-null  float64
##  7   dimension7    55524 non-null  float64
##  8   dimension8    55524 non-null  float64
## dtypes: float64(8), int64(1)
## memory usage: 3.8 MB
```

2. Preprocess Data

Your task is to predict the score of the sentiment (positive or negative) between 0 and 1 based on the embeddings of the tweets specified in columns 2 – 9 of the pre-processed spreadsheet (or the original tweets if you decided to work with the raw tweeting data). To evaluate the performance of your model, please split the dataset into the train set and the test set in the 0.8:0.2 ratio and use cross-validation to calculate the prediction performance.

2.1 Train/Test Split

```
target = 'sentiment'

X_train, X_test, y_train, y_test = train_test_split(
    data.drop(target, axis=1), # predictors
    data[target], # target
    test_size=0.2,
    random_state=90
)
```

```
X_train.shape
```

```
## (44419, 8)
```

```
X_test.shape
```

```
## (11105, 8)
```

2.2 Cross-Validation

2.2.1 Cross-Validation Splits

```
cv_splits = KFold(n_splits=5, random_state=90, shuffle=True)

fold_num = 1
for train, test in cv_splits.split(X_train, y_train):
```

```

print(f'Fold #{fold_num}: Train shape: {X_train.iloc[train].shape},',
      f'Test shape: {X_train.iloc[test].shape}')
fold_num += 1

```

```

## Fold #1: Train shape: (35535, 8), Test shape: (8884, 8)
## Fold #2: Train shape: (35535, 8), Test shape: (8884, 8)
## Fold #3: Train shape: (35535, 8), Test shape: (8884, 8)
## Fold #4: Train shape: (35535, 8), Test shape: (8884, 8)
## Fold #5: Train shape: (35536, 8), Test shape: (8883, 8)

```

2.2.2 Keras Cross-Validation Wrapper

```

def cross_val_score_keras(X_train, y_train, cv_splitter, model, batch_size, num_epochs):

    # model configuration
    loss_function = tf_mae
    optimizer = Adam()
    callback = EarlyStopping(monitor='loss', patience=10)

    # cross-validation scores
    cv_results = list()

    # K-fold Cross Validation model evaluation
    for train, test in cv_splitter.split(X_train, y_train):

        # compile the model
        model.compile(loss=loss_function,
                      optimizer=optimizer)

        # fit data to model
        history = model.fit(X_train.iloc[train], y_train.iloc[train],
                            batch_size=batch_size,
                            epochs=num_epochs,
                            callbacks=[callback],
                            verbose=0)

        # generate generalization metrics
        score = model.evaluate(X_train.iloc[test], y_train.iloc[test], verbose=0)
        cv_results.append(score)

    return np.array(cv_results)

```

3. Neural Networks

You can use any neural network model you like for this classification task. In particular, you may start with a simple single fully connected network as a “baseline” and then try to use more complex models, including CNN and RNN based models, to achieve better performance results than this simple baseline model. Your goal is to reach the mean absolute error of at least 0.48, which should not be too difficult. If you want to be more ambitious, you can try to reach the mean absolute error of 0.47 (medium difficulty), or even 0.46 (this is difficult). The higher accuracy you get, the more points you will be awarded.

In addition to the simple NN baseline mentioned above, you should also build another basic baseline, such as a logistic regression model (similar to the one we used in the RapidMiner Lab done in the class) and

compare the performance results of your DL-based model with that baseline. The expectation is that the more sophisticated DL-model should outperform simple baselines.

3.1 Baseline Model Performance

3.1.1 Logistic Regression

```
model_lr = LogisticRegression()
cv_results_lr = cross_validate(
    model_lr,
    X_train,
    y_train,
    scoring='neg_mean_absolute_error',
    cv=cv_splits
)
cv_lr_mu = -cv_results_lr['test_score'].mean()
cv_lr_sd = cv_results_lr['test_score'].std()

print(f'Logistic Regression MAE: {cv_lr_mu:.4f} +/- {cv_lr_sd:.4f}\n')
```

Logistic Regression MAE: 0.4366 +/- 0.0041

Time to execute this code chunk: 0 s

3.1.2 Random Forest

```
model_rf = RandomForestClassifier(n_estimators=500)
cv_results_rf = cross_validate(
    model_rf,
    X_train,
    y_train,
    scoring='neg_mean_absolute_error',
    cv=cv_splits
)
cv_rf_mu = -cv_results_rf['test_score'].mean()
cv_rf_sd = cv_results_rf['test_score'].std()

print(f'Random Forest MAE: {cv_rf_mu:.4f} +/- {cv_rf_sd:.4f}\n')
```

Random Forest MAE: 0.4202 +/- 0.0048

Time to execute this code chunk: 4 s

3.1.3 Simple Neural Network

```
# define the model architecture
model_nn_0 = Sequential([
    Dense(32, activation='relu', input_dim=8),
    Dense(1, activation='sigmoid')
])

# cross-validation
cv_results_nn_0 = cross_val_score_keras(
    X_train, y_train,
    cv_splits,
```

```

    model_nn_0,
    batch_size=50,
    num_epochs=100
)

# summary statistics
cv_nn_0_mu = cv_results_nn_0.mean()
cv_nn_0_sd = cv_results_nn_0.std()
print(f'Simlpe NN Baseline Model #0 MAE: {cv_nn_0_mu:.4f} +/- {cv_nn_0_sd:.4f}\n')

## Simlpe NN Baseline Model #0 MAE: 0.4703 +/- 0.0046

Time to execute this code chunk: 19 s

# model paramters
batch_size = 250
num_epochs = 100
loss_function = tf_mae
optimizer = Adam()
callback = EarlyStopping(monitor='loss', patience=10)

# architecture
model_nn_0 = Sequential([
    Dense(32, activation='relu', input_dim=8),
    Dense(1, activation='sigmoid')
])

# compile the model
model_nn_0.compile(loss=loss_function,
                   optimizer=optimizer)

# fit data to model
history_nn0 = model_nn_0.fit(X_train, y_train,
                             batch_size=batch_size,
                             epochs=num_epochs,
                             callbacks=[callback],
                             verbose=0)

```

3.2 Deep Learning

3.2.1 DL Model #1

```

# define the model architecture
model_nn_1 = Sequential([
    Dense(8, activation='relu', input_dim=8),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(1, activation='sigmoid')
])

```

```

# cross-validation
cv_results_nn_1 = cross_val_score_keras(
    X_train, y_train,
    cv_splits,
    model_nn_1,
    batch_size=250,
    num_epochs=100
)

# summary statistics
cv_nn_1_mu = cv_results_nn_1.mean()
cv_nn_1_sd = cv_results_nn_1.std()
print(f'Deep Learning Model #1 MAE: {cv_nn_1_mu:.4f} +/- {cv_nn_1_sd:.4f}\n')

```

Deep Learning Model #1 MAE: 0.4703 +/- 0.0046

Time to execute this code chunk: 8 s

```

# model paramters
batch_size = 250
num_epochs = 100
loss_function = tf_mae
optimizer = Adam()
callback = EarlyStopping(monitor='loss', patience=10)

# architecture
model_nn_1 = Sequential([
    Dense(8, activation='relu', input_dim=8),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(8, activation='relu'),
    Dense(1, activation='sigmoid')
])

# compile the model
model_nn_1.compile(loss=loss_function,
                   optimizer=optimizer)

# fit data to model
history_nn1 = model_nn_1.fit(X_train, y_train,
                             batch_size=batch_size,
                             epochs=num_epochs,
                             callbacks=[callback],
                             verbose=0)

```

Time to execute this code chunk: 3 s

3.2.2 DL Model #2

```

# define the model architecture
model_nn_2 = Sequential([

```

```

Dense(64, activation='relu', input_dim=8),
BatchNormalization(),
Dense(32, activation='relu'),
BatchNormalization(),
Dense(16, activation='relu'),
Dense(1, activation='sigmoid')
])

# cross-validation
cv_results_nn_2 = cross_val_score_keras(
    X_train, y_train,
    cv_splits,
    model_nn_2,
    batch_size=500,
    num_epochs=100
)

# summary statistics
cv_nn_2_mu = cv_results_nn_2.mean()
cv_nn_2_sd = cv_results_nn_2.std()
print(f'Deep Learning Model #1 MAE: {cv_nn_2_mu:.4f} +/- {cv_nn_2_sd:.4f}\n')

```

Deep Learning Model #1 MAE: 0.4055 +/- 0.0105

Time to execute this code chunk: 26 s

```

# model paramters
batch_size = 250
num_epochs = 100
loss_function = tf_mae
optimizer = Adam()
callback = EarlyStopping(monitor='loss', patience=10)

# compile the model
model_nn_2.compile(loss=loss_function,
                   optimizer=optimizer)

# fit data to model
history_nn2 = model_nn_2.fit(X_train, y_train,
                             batch_size=batch_size,
                             epochs=num_epochs,
                             callbacks=[callback],
                             verbose=0)

```

Time to execute this code chunk: 3 s

4. Evaluation

After you build your neural network, apply the trained deep learning model to the test set and evaluate its performance using the accuracy measures.

```

# refit baseline models with full data and predict holdout data
_ = model_lr.fit(X_train, y_train)
preds_lr = model_lr.predict(X_test)

```



```

_ = model_rf.fit(X_train, y_train)
preds_rf = model_rf.predict(X_test)

# predict test test for NN models
preds_nn_0 = model_nn_0.predict(X_test).reshape(-1)
preds_nn_1 = model_nn_1.predict(X_test).reshape(-1)
preds_nn_2 = model_nn_2.predict(X_test).reshape(-1)

# evaluation of test predictions
mae_test_lr = skl_mae(y_test, preds_lr)
mae_test_rf = skl_mae(y_test, preds_rf)
mae_test_nn_0 = skl_mae(y_test, preds_nn_0)
mae_test_nn_1 = skl_mae(y_test, preds_nn_1)
mae_test_nn_2 = skl_mae(y_test, preds_nn_2)

test_results_df = pd.DataFrame({
    'Model': ['logistic regression', 'random forest', 'simple NN', 'DL #1', 'DL #2'],
    'MAE': [mae_test_lr, mae_test_rf, mae_test_nn_0, mae_test_nn_1, mae_test_nn_2]
}).sort_values('MAE').reset_index(drop=True)

test_results_df

```

```

##           Model      MAE
## 0           DL #2  0.411826
## 1    random forest  0.417019
## 2  logistic regression  0.429896
## 3           DL #1  0.467447
## 4        simple NN  0.467447

```