

NLP

소설 작가 분류

natural

learning

processing

NLP

text

language

public

processed

understanding

automatic

linguistics

download
process

computer

retrieval

tag
typo

design

discourse
analysis

job

communicate
simulation
keywords

output

operating
typography
information
systems
human

telecommunications

interaction

coreference

programming technology automated
summarization evaluation
networks statistical

machine

connect

media

science

data
evolution

cloud

testing

intelligence

팀명 : Old Pine
영동대학교학과 노승찬
영동대학교학과 송정영

1. EDA

2. MODELING

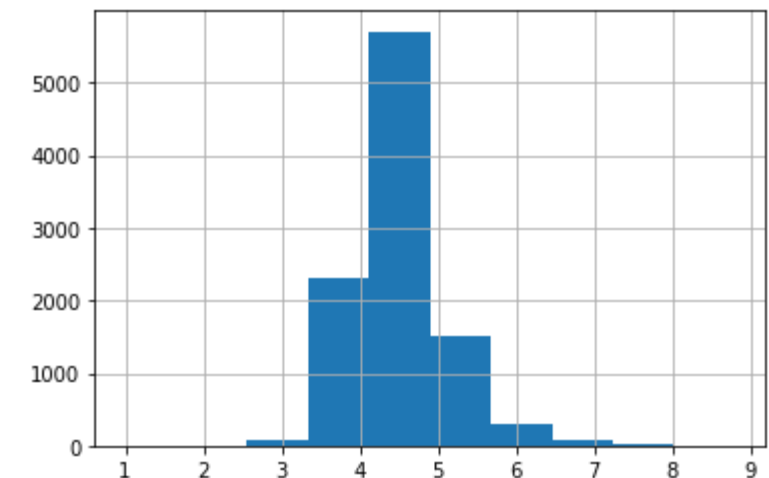
3. 결과 및 느낀 점



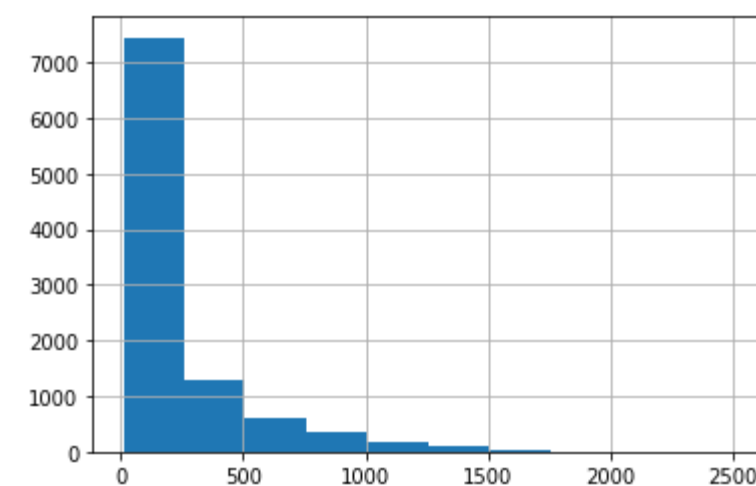
EDA

텍스트의 분포에 대한 EDA

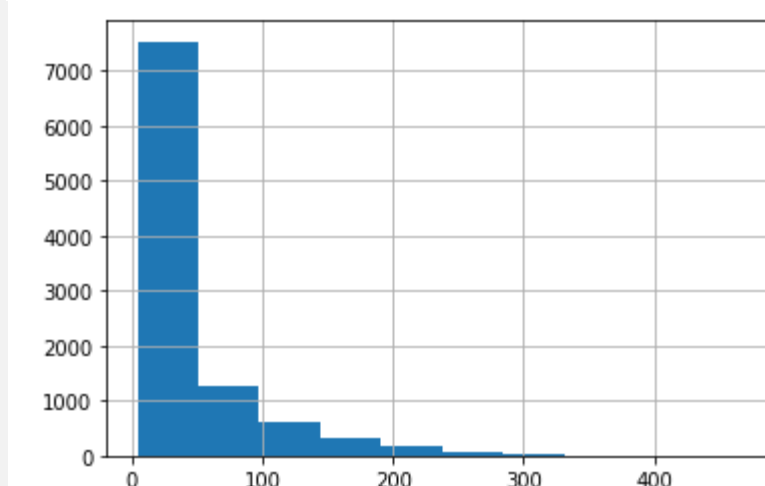
average word length (평균 단어 길이)



문자 길이 분포



word level length (단어 기준 길이 분포) ¶



이러한 수치들이 모델링에 관여하면 **성능이 향상**될까?

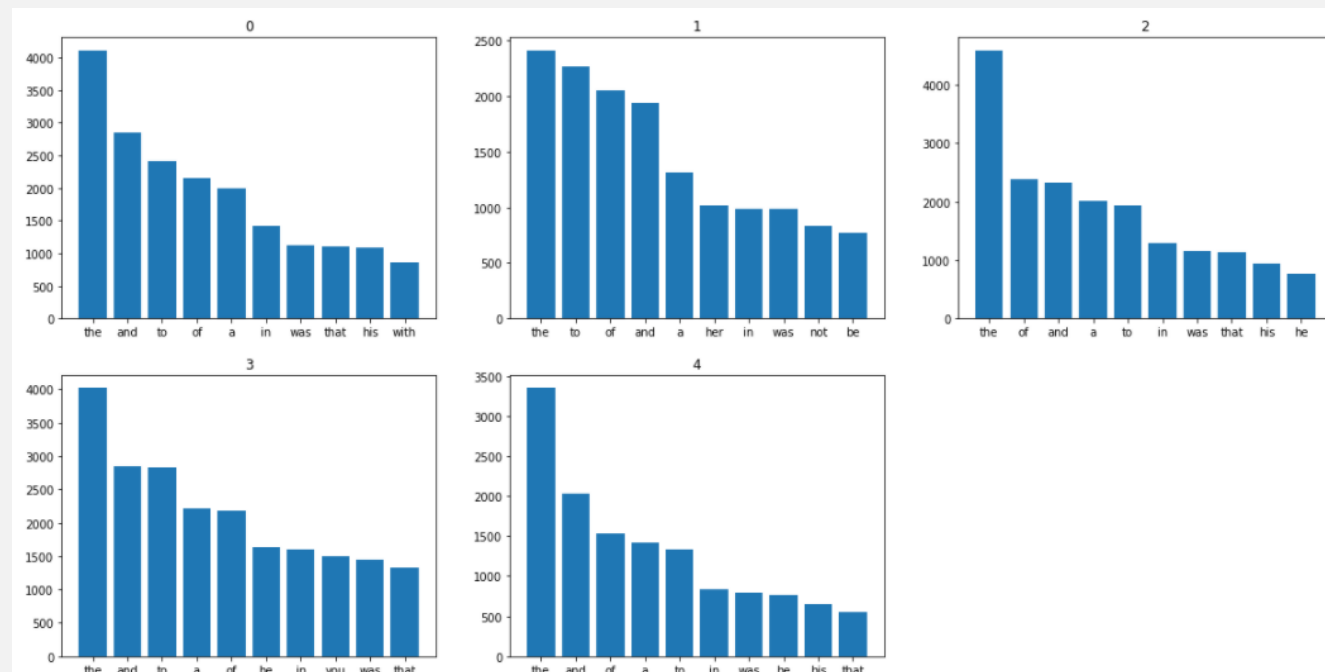
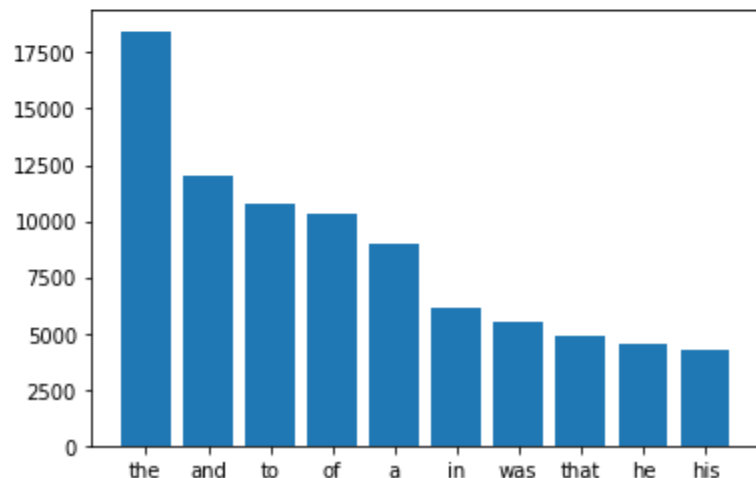
stopwords 분포

```
def plot_top_stopwords_barchart(text):
    stop=set(stopwords.words('english'))

    new= text.str.split()
    new=new.values.tolist()
    corpus=[word for i in new for word in i]
    from collections import defaultdict
    dic=defaultdict(int)
    for word in corpus:
        if word in stop:
            dic[word]+=1

    top=sorted(dic.items(), key=lambda x:x[1],reverse=True)[:10]
    x,y=zip(*top)
    plt.bar(x,y)
```

```
plot_top_stopwords_barchart(train['text'])
```



작가별로 자주 사용하는 STOPWORDS가 있지 않을까?



작가별 STOPWORDS의 빈도로 작가를 파악하는 것에 도움이 되지 않을까?



MODELING

전처리

모델링

앙상블

A

불용어 제거 x
+
Embedding

- >> Model 1 : RNN(LSTM + GRU)
- >> Model 2 : CNN

B

추가 피쳐
+
TF-IDF
+
Embedding

- >> Model 3 : Light GBM
- >> Model 4 : XGBoost

Stacking(메타 러너 : Light GBM)

A. RNN & CNN

전처리

```

1 #알파벳과 숫자를 제외한 나머지 제거
2 def alpha_num(text):
3     return re.sub(r'[^A-Za-z0-9,! ]', '', text)
4
5 #불용어 제거
6 def remove_stopwords(text):
7     final_text = []
8     for i in text.split():
9         if i.strip().lower() not in stopwords:
10             final_text.append(i.strip())
11     return " ".join(final_text)
12
13 stopwords = [
14     # "the", "a", "about", "above", "after", "again", "against", "all", "am", "an", "and", "any", "are", "as",
15     # "at", "be", "because", "been", "before", "being", "below", "between", "both", "but", "by", "could",
16     # "did", "do", "does", "doing", "down", "during", "each", "few", "for", "from", "further", "had", "has",
17     # "have", "having", "he", "he'd", "he'll", "he's", "her", "here", "here's", "hers", "herself", "him", "himself",
18     # "his", "how", "how's", "i", "i'd", "i'll", "i'm", "i've", "if", "in", "into", "is", "it", "it's", "its", "itself",
19     # "let's", "me", "more", "most", "my", "myself", "nor", "of", "on", "once", "only", "or", "other", "ought", "our", "ours",
20     # "ourselves", "out", "over", "own", "same", "she", "she'd", "she'll", "she's", "should", "so", "some", "such", "than", "that",
21     # "that's", "the", "their", "theirs", "them", "themselves", "then", "there", "there's", "these", "they", "they'd", "they'll",
22     # "they're", "they've", "this", "those", "through", "to", "too", "under", "until", "up", "very", "was", "we", "we'd", "we'll",
23     # "we're", "we've", "were", "what", "what's", "when", "when's", "where", "where's", "which", "while", "who", "who's", "whom",
24     # "why", "why's", "with", "would", "you", "you'd", "you'll", "you're", "you've", "your", "yours", "yourself", "yourselves"
25 ]

```

```

[ ] 1 #작가 분류이기 때문에 불용어를 제거 하지 않는 것이 성능향상에 도움됨
2
3 # train['text'] = train['text'].str.lower()
4 # test['text'] = test['text'].str.lower()
5 # train['text'] = train['text'].apply(alpha_num)
6 # .apply(remove_stopwords)
7 # test['text'] = test['text'].apply(alpha_num)
8 # .apply(remove_stopwords)

```

설명

- RNN의 경우 어떠한 전처리로 하지 않음
- CNN은 특수문자만 제거

RNN

```
[ ] 1 def get_model():
2     # with tf.device('/device:GPU:0'):
3     model = Sequential([
4         Embedding(vocab_size, embedding_dim, input_length=max_length),
5         # Embedding(47256, 300, weights=[embedding_matrix], input_length=max_length, trainable=False),
6         Bidirectional(GRU(16, activation='tanh', return_sequences=True)),
7         # Dropout(0.1),
8         Bidirectional(LSTM(16)),
9         Dense(n_class, activation='softmax')
10    ])
11    model.compile(loss='categorical_crossentropy', optimizer=Adam(learning_rate=0.005))
12    return model

[ ] 1 p_val = np.zeros((trn.shape[0], n_class))
2   p_tst = np.zeros((tst.shape[0], n_class))
3   for i, (i_trn, i_val) in enumerate(cv.split(trn, y), 1):
4       print(f'training model for CV #{i}')
5       clf = get_model()
6
7       es = EarlyStopping(monitor='val_loss', min_delta=0.001, patience=2,
8                           verbose=1, mode='min', baseline=None, restore_best_weights=True)
9
10      # rlr = ReduceLROnPlateau(monitor='val_loss', factor=0.5,
11                               # patience=2, min_lr=1e-6, mode='min', verbose=1)
12
13      clf.fit(trn[i_trn],
14              to_categorical(y[i_trn]),
15              validation_data=(trn[i_val], to_categorical(y[i_val])),
16              epochs=30,
17              callbacks=[es
18                          # ,rlr
19                      ])
20
21      p_val[i_val, :] = clf.predict(trn[i_val])
22      p_tst += clf.predict(tst) / n_fold
```

설명

- 임베딩층과 GRU층 LSTM층으로 구성
- Grid search로 노드의 개수 탐색
- Early Stopping 사용
- 마지막 output layer의 활성화 함수는 'soft max'

>> Model 1 : RNN(LSTM + GRU)

CNN

```

1 def get_model():
2     # with tf.device('/device:GPU:0'):
3     model = Sequential([
4         Embedding(vocab_size, embedding_dim, input_length=max_length),
5         SpatialDropout1D(0.7),
6         Conv1D(32, 8, padding="valid", activation="relu", strides=3),
7         # Dropout(.3),
8         Conv1D(32, 8, padding="valid", activation="relu", strides=3),
9         GlobalMaxPooling1D(),
10        # Dense(64, activation='relu'),
11        Dense(n_class, activation='softmax')
12    ])
13    model.compile(loss='categorical_crossentropy', optimizer=Adam(learning_rate=.01))
14    return model

[ ] 1 p_val = np.zeros((trn.shape[0], n_class))
2     p_tst = np.zeros((tst.shape[0], n_class))
3     for i, (i_trn, i_val) in enumerate(cv.split(trn, y), 1):
4         print(f'training model for CV #{i}')
5         clf = get_model()
6         es = EarlyStopping(monitor='val_loss', min_delta=0.001, patience=5,
7                             verbose=1, mode='min', baseline=None, restore_best_weights=True)
8         rlr = ReduceLROnPlateau(monitor='val_loss', factor=0.5,
9                                 patience=3, min_lr=1e-6, mode='min', verbose=1)
10        clf.fit(trn[i_trn],
11                to_categorical(y[i_trn]),
12                validation_data=(trn[i_val], to_categorical(y[i_val])),
13                epochs=100,
14                batch_size=512,
15                callbacks=[es
16                            # ,rlr
17                        ])
18    )
19    p_val[i_val, :] = clf.predict(trn[i_val])
20    p_tst += clf.predict(tst) / n_fold

```

» Model 2 :CNN

설명

- 베이스 라인의 CNN 코드에서 SpatialDropout1D를 추가하여 과적합 방지 및 성능 향상도모.
(노드를 드랍해버리는 것이 아닌 몇몇 컬럼 자체를 드랍)

B. LightGBM & XGBoost

변수 추가

```
1  ## Number of words in the text ##
2  train["num_words"] = train["text"].apply(lambda x: len(str(x).split()))
3  test["num_words"] = test["text"].apply(lambda x: len(str(x).split()))
4
5  ## Number of unique words in the text ##
6  train["num_unique_words"] = train["text"].apply(lambda x: len(set(str(x).split())))
7  test["num_unique_words"] = test["text"].apply(lambda x: len(set(str(x).split())))
8
9  ## Number of characters in the text ##
10 train["num_chars"] = train["text"].apply(lambda x: len(str(x)))
11 test["num_chars"] = test["text"].apply(lambda x: len(str(x)))
12
13 ## Number of stopwords in the text ##
14 train["num_stopwords"] = train["text"].apply(lambda x: len([w for w in str(x).lower().split() if w in stop_words]))
15 test["num_stopwords"] = test["text"].apply(lambda x: len([w for w in str(x).lower().split() if w in stop_words]))
16
17 ## Number of punctuations in the text ##
18 train["num_punctuations"] = train["text"].apply(lambda x: len([c for c in str(x) if c in string.punctuation]))
19 test["num_punctuations"] = test["text"].apply(lambda x: len([c for c in str(x) if c in string.punctuation]))
20
21 ## Number of title case words in the text ##
22 train["num_words_upper"] = train["text"].apply(lambda x: len([w for w in str(x).split() if w.isupper()]))
23 test["num_words_upper"] = test["text"].apply(lambda x: len([w for w in str(x).split() if w.isupper()]))
24
25 ## Number of title case words in the text ##
26 train["num_words_title"] = train["text"].apply(lambda x: len([w for w in str(x).split() if w.istitle()]))
27 test["num_words_title"] = test["text"].apply(lambda x: len([w for w in str(x).split() if w.istitle()]))
28
29 ## Average length of the words in the text ##
30 train["mean_word_len"] = train["text"].apply(lambda x: np.mean([len(w) for w in str(x).split()]))
31 test["mean_word_len"] = test["text"].apply(lambda x: np.mean([len(w) for w in str(x).split()]))
```

설명

- 아래와 같은 변수 추가
 - 단어 개수
 - 불용어 개수
 - 글자 개수
 - 특수문자 개수
 - 대문자 개수
 - 단어 길이의 평균
- Kaggle 내에 있는 프로젝트
참고

Feature Engineering

```

1 def runMNB(train_X, train_y, test_X, test_y, test_X2):
2     model = naive_bayes.MultinomialNB()
3     model.fit(train_X, train_y)
4     pred_test_y = model.predict_proba(test_X)
5     pred_test_y2 = model.predict_proba(test_X2)
6     return pred_test_y, pred_test_y2, model

1 ### Fit transform the tfidf vectorizer ###
2 tfidf_vec = TfidfVectorizer(ngram_range=(1,5), analyzer='char')
3 full_tfidf = tfidf_vec.fit_transform(train_df['text'].values.tolist())
4 train_tfidf = tfidf_vec.transform(train_df['text'].values.tolist())
5 test_tfidf = tfidf_vec.transform(test_df['text'].values.tolist())
6
7 cv_scores = []
8 pred_full_test = 0
9 pred_train = np.zeros([train_df.shape[0], 5])
10 kf = model_selection.KFold(n_splits=5, shuffle=True, random_state= 42)
11 for dev_index, val_index in kf.split(train_X):
12     dev_X, val_X = train_tfidf[dev_index], train_tfidf[val_index]
13     dev_y, val_y = train_y[dev_index], train_y[val_index]
14     pred_val_y, pred_test_y, model = runMNB(dev_X, dev_y, val_X, val_y, test_tfidf)
15     pred_full_test = pred_full_test + pred_test_y
16     pred_train[val_index,:] = pred_val_y
17     cv_scores.append(metrics.log_loss(val_y, pred_val_y))
18 print("Mean cv score : ", np.mean(cv_scores))
19 pred_full_test = pred_full_test / 5.
20
21 # add the predictions as new features #
22 train_df["nb_tfidf_char_eap"] = pred_train[:,0]
23 train_df["nb_tfidf_char_hpl"] = pred_train[:,1]
24 train_df["nb_tfidf_char_mws"] = pred_train[:,2]
25 test_df["nb_tfidf_char_eap"] = pred_full_test[:,0]
26 test_df["nb_tfidf_char_hpl"] = pred_full_test[:,1]
27 test_df["nb_tfidf_char_mws"] = pred_full_test[:,2]

```

```

1 n_comp = 20
2 svd_obj = TruncatedSVD(n_components=n_comp, algorithm='arpack')
3 svd_obj.fit(train_tfidf)
4 train_svd = pd.DataFrame(svd_obj.transform(train_tfidf))
5 test_svd = pd.DataFrame(svd_obj.transform(test_tfidf))
6
7 train_svd.columns = ['svd_char_'+str(i) for i in range(n_comp)]
8 test_svd.columns = ['svd_char_'+str(i) for i in range(n_comp)]
9 train_df = pd.concat([train_df, train_svd], axis=1)
10 test_df = pd.concat([test_df, test_svd], axis=1)
11 del full_tfidf, train_tfidf, test_tfidf, train_svd, test_svd

```

- 나이브 베이즈 모델의 예측값을 피처로 추가
- 정보를 압축하고 컴팩트하게 표현하기 위해 SVD 사용하여 차원축소
- TF-IDF에 svd 형상을 만들어 특징 세트에 추가
- Kaggle에서 소스 필사. 의미에 대한 학습 필요

Light GBM

```
[ ] 1 def runLGBM(train_X, train_y, test_X, test_y=None, test_X2=None, seed_val = 42, child=1):
2     param = {}
3     param['objective'] = 'multiclass'
4     param['boosting_type'] = 'gbdt'
5     param['subsample_freq'] = 5
6     param['max_depth'] = 30
7     param['num_leaves'] = 100
8     param['num_class'] = 5
9     param['colsample_bytree']=0.7
10    param['subsample'] = 0.8
11    param['min_data_in_leaf'] = 64
12    param['metric'] = 'multi_logloss'
13    param['subsample_for_bin'] = 23000
14    param['min_child_weight'] = child
15    param['learning_rate'] = 0.01
16    param['seed'] = seed_val
17    n_estimators = 20000
18
19    # plst = param.items()
20    lgbmtrain = lgbm.Dataset(train_X, label=train_y)
21
22    if test_y is not None:
23        lgbmtest = lgbm.Dataset(test_X, label=test_y)
24        # watchlist = [ (lgbmtrain, 'train'), (lgbmtest, 'test') ]
25        model = lgbm.train(param, lgbmtrain, n_estimators, valid_sets= [lgbmtrain, lgbmtest], early_stopping_rounds=50, verbose_eval= 100)
26    else:
27        lgbmtest = lgbm.Dataset(test_X)
28        model = lgbm.train(plst, lgbmtrain, num_rounds)
29
30    pred_test_y = model.predict(test_X)
31    if test_X2 is not None:
32        pred_test_y2 = model.predict(test_X2)
33    return pred_test_y, pred_test_y2, model
```

» Model 3

XGBoost

```
def runXGB(train_X, train_y, test_X, test_y=None, test_X2=None, seed_val = 42, child=1, colsample=0.3):
    param = {}
    param['objective'] = 'multi:softprob'
    param['eta'] = 0.1
    param['max_depth'] = 10
    param['silent'] = 1
    param['booster'] = 'dart'
    param['num_class'] = 5
    param['eval_metric'] = "mlogloss"
    param['min_child_weight'] = child
    param['subsample'] = 0.8
    param['colsample_bytree'] = colsample
    param['seed'] = seed_val
    param['tree_method'] = 'gpu_hist'
    num_rounds = 20000

    plst = list(param.items())
    xgtrain = xgb.DMatrix(train_X, label=train_y)

    if test_y is not None:
        xgtest = xgb.DMatrix(test_X, label=test_y)
        watchlist = [ (xgtrain, 'train'), (xgtest, 'test') ]
        model = xgb.train(plst, xgtrain, num_rounds, watchlist, early_stopping_rounds=50, verbose_eval= 100)
    else:
        xgtest = xgb.DMatrix(test_X)
        model = xgb.train(plst, xgtrain, num_rounds)

    pred_test_y = model.predict(xgtest, ntree_limit = model.best_ntree_limit)
    if test_X2 is not None:
        xgtest2 = xgb.DMatrix(test_X2)
        pred_test_y2 = model.predict(xgtest2, ntree_limit = model.best_ntree_limit)
    return pred_test_y, pred_test_y2, model
```

» Model 5

스택킹(메타 러너 : Light GBM)

```

1  # p_trn = np.zeros((stk_trn.shape[0], n_class))
2  p_val = np.zeros((stk_trn.shape[0], n_class))
3  p_tst = np.zeros((stk_tst.shape[0], n_class))
4  for i, (i_trn, i_val) in enumerate(cv.split(stk_trn, y), 1):
5      print(f'training model for CV #{i}')
6      clf = lgb.LGBMClassifier(objective='multiclass',
7                              n_estimators=10000,
8                              learning_rate=0.01,
9                              boosting_type='gbdt',
10                             max_depth=5,
11                             feature_fraction=0.4,
12                             min_child_weight=0.01,
13                             num_leaves=30,
14                             random_state=seed,
15                             n_jobs=-1,
16                             verbose=100)
17      clf.fit(stk_trn[i_trn], y[i_trn],
18             eval_set=[(stk_trn[i_val], y[i_val])],
19             eval_metric='multi_logloss', early_stopping_rounds=100,
20             verbose=100)
21      # p_trn[i_trn, :] = clf.predict_proba(stk_trn[i_trn])
22      p_val[i_val, :] = clf.predict_proba(stk_trn[i_val])
23      p_tst += clf.predict_proba(stk_tst) / n_fold
24  print()
25  print('models:', model_names)
26  print(clf)
27  # print(f'train cv accuracy : {accuracy_score(y, np.argmax(p_trn, axis=1)) :.6f}')
28  print(f'valid cv accuracy : {log_loss(y, p_val) :.6f}')

```

설명

- Grid search를 통해 parameter tuning
- XGBoost, DNN, RandomForest를 시도하였으나 LightGBM의 경우가 성능 가장 좋음.
- 최종 스택킹 CV Logloss : 0.41...

전처리

모델링

앙상블

A

불용어 제거 x
+
Embedding

- >> Model 1 : RNN(LSTM + GRU)
- >> Model 2 : CNN

B

추가 피쳐
+
TF-IDF
+
Embedding

- >> Model 3 : Light GBM
- >> Model 4 : XGBoost

Stacking(메타 러너 : Light GBM)

결과 및 느낀 점

최종 스코어(private 기준)
log loss : 0.23471(27위)

- Pre-trained embedding을 사용하였으나 성능 하락.
- 오히려 Lemmatization, 불용어 제거 등의 전처리하지 않은 데이터셋의 성능이 더 우수.(RNN 기준)
- 어텐션에 대한 시도를 못해본 것이 아쉬움.
- NLP관련 지식이 부족하여 이론적으로 접근하는 것이 아닌 다양한 방법을 모두 시도하여 시간적으로 많은 투자가 있었던 것이 아쉬움.
- NLP 분야에 입문하는 것에 큰 도움이 되었음.

참고 사이트

<https://www.kaggle.com/sudalairajkumar/simple-feature-engg-notebook-spooky-author>

<https://jiho-ml.com/>

참고 서적

핸즈온 머신러닝 : 사이킷런, 케라스, 텐서플로를 활용한 머신러닝, 딥러닝 완벽 실무

텐서플로와 머신러닝으로 시작하는 자연어 처리

한 학기간 좋은 수업과 강의 제공해주셔서
진심으로 감사드립니다.