	<pre>import nltk.data from nltk import word_tokenize, sent_tokenize sdg_names = pd.read_excel("./SDGtrainingset.xlsx") sdg_names = sdg_names.drop([e], 2], axis=0) sdg_names = sdg_names.drop([e], 2], axis=0) sdg_names = sdg_names.sdg_names.sdg_name.tolist() enbed_url = "https://tfnbb.dew/google/universal-sentence-encoder/4" embed = hub.load(embed_url) def get_text_df(ffile_name): text_df = pd.read_csv(file_name, sep = "\t", quotechar='") col_names = text_df.columns.values[0].split('\t') text_df rool names] = text_df.columns.values[0].axis = 1, inplace=true) text_df - text_df.columns.values[0].axis = 1, inplace=true) text_df = text_df.columns.values[0].axis = 1, inplace=true) text_df = text_df.query("agreement > 0.5 and (labels_positive - labels_negative) : int, 'labels_negative': int, 'lapecement': float}, copy=True) text_df = text_df.query("agreement > 0.5 and (labels_positive - labels_negative) > 2") text_df = get_text_df(file_name) text_df = get_text_df(f</pre>
	<pre>Bigram Only def run classifier_count(text_df, classifier_algorithm, ngram_range): docs = text_df.text categories = text_df.sdg X_train, X_ttest, y_train, y_test = train_test_split(docs, categories, test_size=0.33, random_state=42) X_train.count_vectorizer = Countvectorizer(ngram_range-ngram_range, stop_words = "english", min_df = 7) X_train.count_vectorizer.fit(X_train) X_train.count_vector = X_train_count_vectorizer.transform(X_train) X_test_count_vector = X_train_count_vectorizer.transform(X_test) if classifier_algorithm == RidgeClassifier: clf = classifier_algorithm== RidgeClassifier: clf = classifier_algorithm== RidgeClassifier: clf = classifier_algorithm== M_FClassifier: clf = classifier_algorithm= M_FClassifier: clf = classifier_algorithm= M_FClassifier: clf = classifier_algorithm= M_FClassifier: clf = classifier_algorithm() clf = clf.fit(X_train_count_vector, y_train) y_pred = clf.predict(X_test_count_vector) accuracy = accuracy_score(y_test, y_pred) precision = precision_score(y_test, y_pred, average='macro') fl = fl_score(y_test, y_pred, average='macro') fl = fl_score(y_test, y_pred, average='macro') fl = fl_score(y_test, y_pred, average='macro') print(metrics_classification_report(y_test, y_pred, digits = 4)) run classifier_count(text_df, MultinomialMB, (2, 2))</pre>
	2 0.7898 0.567 0.8981 359 3 0.8976 0.7931 0.8903 667 4 0.8196 0.7931 0.8216 0.8216 876 6 0.7932 0.8378 0.8149 444 7 0.7519 0.8226 0.8287 0.8149 444 7 0.7519 0.8296 0.8296 0.7888 716 8 0.5781 0.4990 0.4799 335 9 0.5788 0.4832 0.5247 327 10 0.6429 0.4532 0.5316 278 11 0.6658 0.6493 0.6532 0.5316 278 11 0.6658 0.6493 0.6539 442 12 0.7515 0.4980 0.5976 250 13 0.6946 0.7963 0.6539 442 14 0.8239 0.5431 0.6546 287 15 0.7293 0.6144 0.6631 386 16 0.7997 0.9489 0.8679 1977 accuracy weighted avg 0.7254 0.8339 0.7593 8141 run_classifier_count(text_df, RidgeClassifier_ (2,2)) precision recall fi-score support 1 0.6275 0.6377 0.6325 494 2 0.6161 0.5799 0.9799 0.7595 811 5 0.7866 0.7861 0.7919 0.7255 871 6 0.7861 0.7861 0.7919 0.7255 871 6 0.7861 0.7861 0.7919 0.7255 871 6 0.7861 0.7861 0.7919 0.7255 871 6 0.7861 0.7861 0.7919 0.7255 871 6 0.7861 0.7861 0.7919 0.7255 871 6 0.7861 0.7861 0.7919 0.7255 871 6 0.7861 0.7861 0.7863 0.7818 444
	8
	accuracy
	most_significant_feature_for_class(X_train_count_vectorizer, X_train_count_vector, Count_ridge_clf, [2, 5, 8, 16], n=10) SDG 2 : food security
	Spo 5 : Women rights 0.00222067 Spo 5 : gender mainstreaming 0.00160507 Spo 5 : gender mainstreaming 0.00160507 Spo 5 : lathith care -0.00200067 Spo 5 : climate change -0.00200067 Spo 5 : climate change -0.0017078 Spo 5 : climate change -0.0017078 Spo 5 : international law -0.00127078 Spo 5 : climate change -0.0017078 Spo 5 : porty reduction -0.0017078 Spo 5 : long term -0.00160801 Spo 5 : porty reduction -0.0016081 Spo 5 : sub saharan -0.000082774 Spo 5 : sub saharan -0.00082774 Spo 5 : developed countries -0.000831374 Spo 5 : developed countries -0.000831374 Spo 5 : developed countries -0.00083188 Spo 8 : labour market -0.0053858 Spo 8 : young people -0.0011303 Spo 8 : labour force -0.0011303 Spo 8 : labour force -0.0011303 Spo 8 : labour force -0.0011303 Spo 8 : job search -0.00128646 Spo 8 : job search -0.00108935 Spo 8 : job search -0.00108935 Spo 8 : gender equality -0.0019284 Spo 8 : pay gap -0.000866199 Spo 8 : energy efficiency -0.00093642 Spo 8 : energy efficiency -0.00093642 Spo 8 : soil economic -0.00093642 Spo 8 : soil economic -0.00093643
	Soc 8 Numan rights
]: _ ;	classifiers = { 'RidgeClassifier' RidgeClassifier(), 'MULTIONNIALINB': MULTIONNIALINB() } results_df = pd_DataFrame({ 'Classifier': ['RidgeClassifier', 'MLPClassifier', 'MLPClassifier', 'MultinomialNB'], 'Accuracy': [0.4882, 0.6873, 0.7583], 'Presision': [0.1482, 0.6832, 0.6873, 0.7583], 'Presision': [0.1482, 0.6293, 0.6835], 'Presults_df RidgeClassifier
,	* Unigram Only run classifier count(text df, MultinomialNB, (1,1)) precision recall fi-score support 1 0.7458 0.8877 0.7555 494 2 0.7547 0.8822 0.7934 359 3 0.9274 0.8816 0.9039 667 4 0.9334 0.9173 0.9253 871 5 0.8889 10.8896 0.8747 444 7 0.8723 0.8898 0.8896 0.8747 444 7 0.8723 0.8898 0.8996 0.8747 444 7 0.8723 0.8898 0.8996 0.8747 444 7 0.8723 0.8898 0.8996 0.8747 444 7 0.8723 0.8898 0.8996 0.8755 716 8 0.6398 0.6542 0.9953 335 9 0.7311 0.7982 0.7520 0.7580 0.895 342 10 0.6398 0.6519 0.6923 278 11 0.7970 0.8439 0.8988 442 12 0.8720 0.7360 0.7983 259 13 0.717 0.8124 0.8155 0.7622 432 14 0.9237 0.8614 0.8915 267 15 0.8839 0.8256 0.8258 396 16 0.9560 0.9675 0.9617 1977 accuracy macro avg 0.8213 0.8142 0.8156 3141 weighted avg 0.8213 0.8491 0.8478 8141 run_classifier_count(text_df, RidgeClassifier, (1,1))
	1
]: (5
]: _	results.df Classifier Acturacy Precision Recall Fl. Score Regoçtassifier 0.8733 0.8898 0.7958 0.8912 MultinomialNB 0.8421 0.8213 0.8142 0.8213 0.8142 MultinomialNB 0.8421 0.8213 0.8142 0.8216 - As we expect, metrics are much better when doing unigram only. Suprisingly, MLPClassifier performed the worst on bigram only, and the best for unigram only. Tfidf Vectorizer - Unigram and Bigram def run classifier tfidf(text.df, classifier.algorithm, ngram_range):
]: [cif = classifier_algorithm() clf = clf, fit(X_trai_ffidf_vector, y_train) y_pred = clf.predict(X_test_tfidf_vector) accuracy = accuracy_score(y_test, y_pred, average='macro') precision = precision_score(y_test, y_pred, average='macro') fi = fl_score(y_test, y_pred, average='macro') fi = fl_score(y_test, y_pred, average='macro') print(metrics.classification_report(y_test, y_pred, digits = 4)) run_classifier_tfidf(text_off, MultinowialNB, (1,2)) precision
]:	run_classifier_tfidf(text_df, RidgeClassifier, (1,2)) precision
]: (1 0.8153 0.8401 0.8275 494 2 0.8842 0.8719 0.8700 359 3 0.9259 0.9310 0.9250 667 4 0.9328 0.9564 0.9261 0.9262 667 4 0.9328 0.9564 0.9064 871 5 0.8896 0.9201 0.9066 875 6 0.9954 0.9954 0.9954 0.9954 0.9064 444 7 0.8945 0.9128 0.9128 0.9952 0.8032 716 8 0.699 0.6299 0.6492 335 9 0.8121 0.8196 0.8158 327 10 0.8896 0.6107 0.6315 273 11 0.819 0.8733 0.8521 442 12 0.8789 0.7840 0.8288 259 13 0.8645 0.8565 0.8665 432 14 0.9533 0.9176 0.9361 267 15 0.9934 0.9562 0.8792 306 16 0.9675 0.9684 0.9689 1077 accuracy
]: _	0 RidgeClassifier 0.8872 0.8729 0.8557 0.8628 1 MLPClassifier 0.8839 0.8636 0.8538 0.8538 0.8538 Bigram Only ru_classifier_tfidf(text_df, MultinomialNB, (2,2)) precision recall f1-score support 1 0.6920 0.7368 0.7137 494 2 0.8700 0.4847 0.6225 359 3 0.7818 0.8096 0.7911 667 4 0.7147 0.8944 0.7945 871 5 0.6126 0.8756 0.7209 876 6 0.8304 0.8378 0.8341 444 7 0.6269 0.8729 0.7297 716 8 0.8026 0.1821 0.2968 335 9 0.7379 0.2324 0.3555 327
	10
]: [
	docs = text_df.text categories = text_df.sdg X_train, X_test, y_train, y_test = train_test_split(docs, categories, test_size=0.33, random_state=42) X_train_tridf_vectorizer = Tridfvectorizer(ngram_range=(2,2), stop_words = "english", min_df=7) X_train_tridf_vectorizer_fit(X_train) X_train_tridf_vector = X_train_tridf_vectorizer_transform(X_test) X_train_tridf_vector = X_train_tridf_vectorizer_transform(X_test) **tridf_ridge_clf = Ridge_clsssifier(tol=0.2, solver="ymarse_cg") **tridf_ridge_clf = Ridge_clsssifier(tol=0.2, solver="ymarse_cg") **tridf_ridge_clf = tridf_ridge_clf.fit(X_train_tridf_vector, y_train) **most_significant_feature_for_class(X_train_tridf_vector, y_train) **most_significant_feature_for_class(X_train_tridf_vector, y_train) **most_significant_feature_for_class(X_train_tridf_vectorizer, X_train_tridf_vector, tridf_ridge_clf, [2, 3, 6, 8,], n=10) **most_significant_feature_for_class(X_train_tridf_vector, Y_train_tridf_vector, Y_train_tridf_vector, Y_train_tridf_vector, Y_train_tridf_vector, Y_train_tridf_vector, Y_train_tridf_
	SDG 2 : climate finance
	SD6 6 : water quality 0.00477423 SD6 6 : iver basin 0.00239355 SD6 6 : drinking water 0.00237797 SD6 6 : water sanitation 0.00234192 SD6 6 : water sector 0.00186882 SD6 6 : water sector 0.00151079 SD6 6 : human rights -0.00113892 SD6 6 : international law -0.000741593 SD6 6 : international law -0.000662528 SD6 6 : developing countries -0.00056258 SD6 6 : sustainable development -0.000560395 SD6 6 : labour market -0.000560395 SD6 8 : labour force -0.0004943608 SD6 8 : labour market -0.000560395 SD6 8 : social security -0.00095614 SD6 8 : social security -0.00095614 SD6 8 : social security -0.00095614