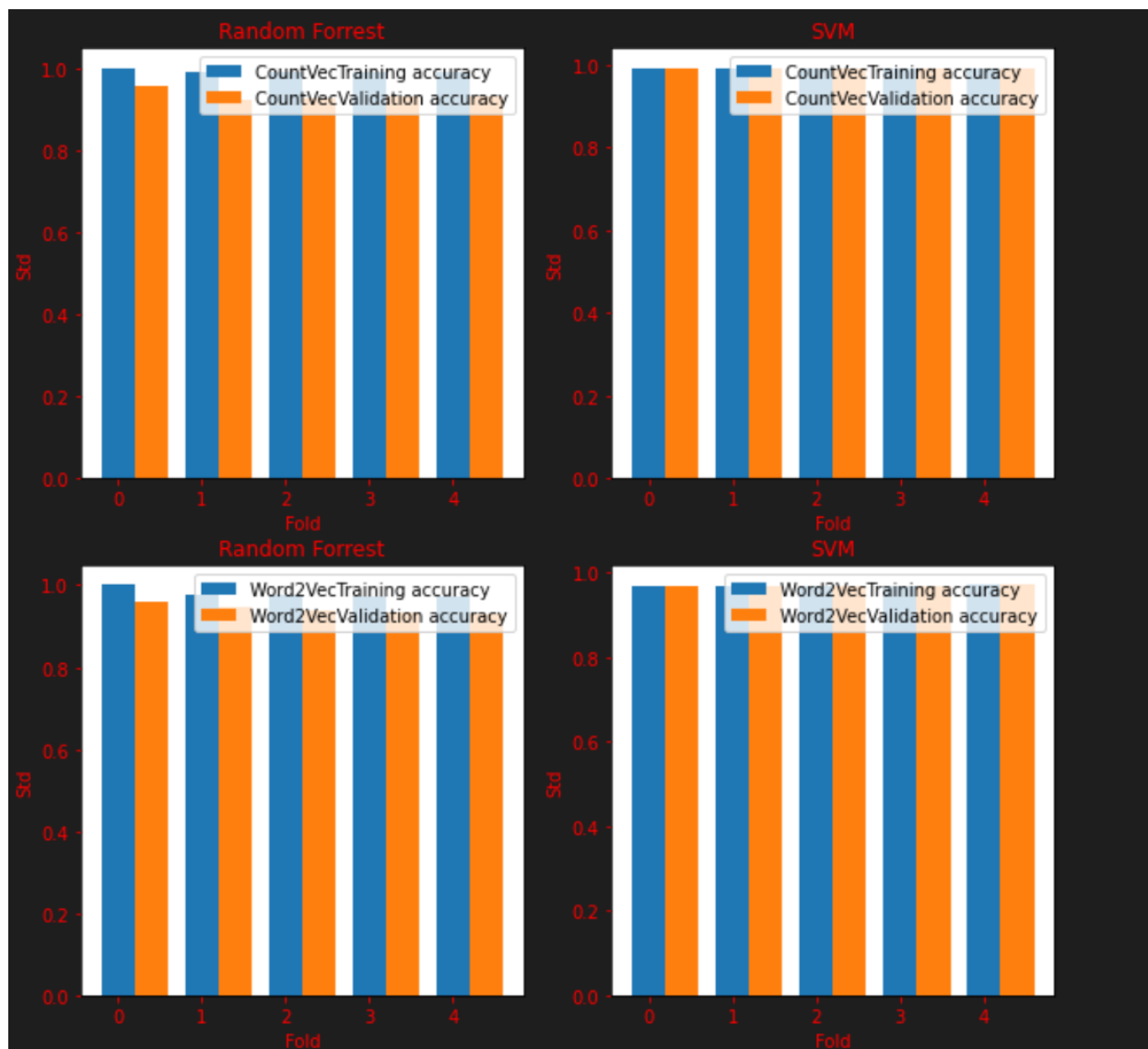


I generate features for using Word2Vectorizer taking in around 100 words and tokens which tokenize the text of each article and then predicts the transformed text values in test.csv.

Count Vectorizer

Random Forest Training Accuracy	SVM Training Accuracy
1.0	0.9899328859060403
0.9932885906040269	0.9903523489932886
0.9922818791946308	0.9900727069351231
0.991490891658677	0.9890939597315437
0.9903057419835943	0.9893456375838927
Random Forest Validation Accuracy	SVM Validation Accuracy
0.9530201342281879	0.9899328859060403
0.9429530201342282	0.9903523489932886
0.9395973154362416	0.9900727069351231
0.9367209971236817	0.9890939597315437
0.9366144668158092	0.9893456375838927
Random Forest Training Std	SVM Training Std
0.0	0.010067114093959717
0.009491366190423444	0.009665871327015844
0.009467608040044198	0.009954825230048956
0.00984931966026094	0.011184808411051482
0.01108815925929372	0.010906362927596636
Random Forest Training Accuracy	SVM Training Accuracy
0.0	0.010067114093959717
0.008219764237527433	0.010171439306193063
0.00925103943093303	0.017220139058546576
0.018504562862529973	0.01860355441404083
0.01887992378996593	0.018938189341640883



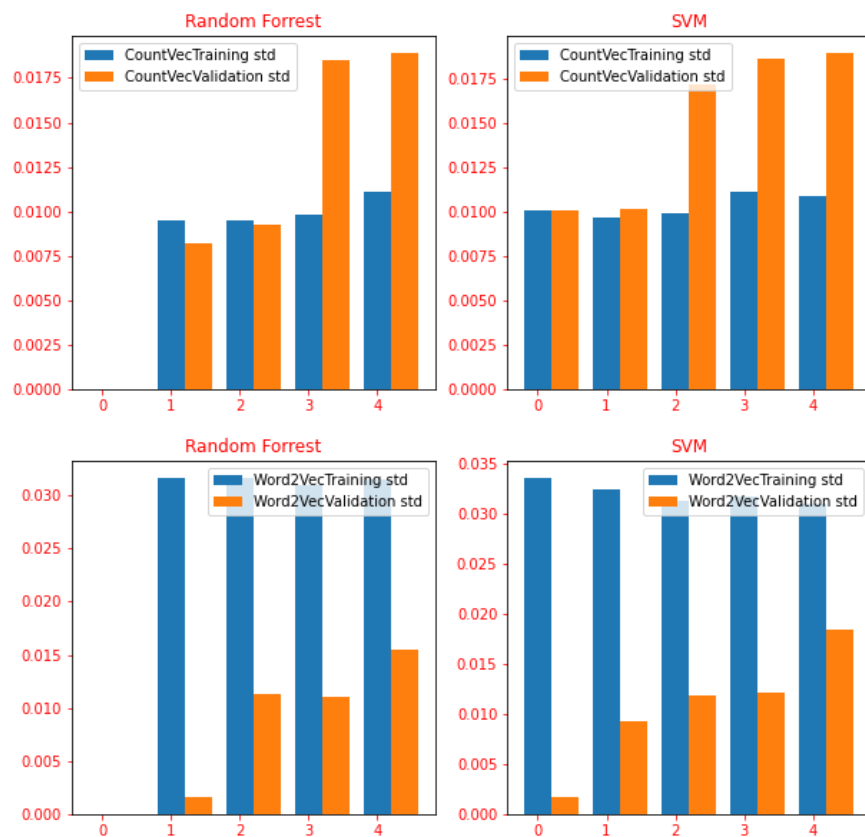
Word2Vectorizer

Random Forest Training Accuracy	SVM Training Accuracy
1.0	0.9664429530201342
0.9776286353467561	0.9677013422818792
0.9741610738255033	0.9688199105145414
0.9732742090124641	0.9684354026845639
0.9719425801640568	0.9691275167785236

Random Forest Validation Accuracy	SVM Validation Accuracy
0.9563758389261745	0.9664429530201342
0.9541387024608502	0.9677013422818792
0.9449664429530202	0.9688199105145414
0.941994247363375	0.9684354026845639
0.935868754660701	0.9691275167785236

Random Forest Training Std	SVM Training Std
0.0	0.03355704697986578
0.03163788730141153	0.03234764871610417
0.03168609834661176	0.0312939954506365
0.03095895315263141	0.03166303853898056
0.03145721121177398	0.03101486353800494

Random Forest Training Accuracy	SVM Training Accuracy
0.0	0.001677852348993314
0.0015818943650706	0.009304141364437428
0.01135002317133405	0.01185101795325338
0.01101546049192333	0.012099165353905995
0.015535172748694736	0.018367694205911403



After looking at all the data we find that SVM model with word2vect has the highest accuracy, so I decided to use SVM with Word2vect word embedding to predict my labels. The accuracy seems to be in 99 percent tile with average accuracy and the standard deviation seems higher with training accuracy, but not validation accuracy.

	Pred
ArticleId	
568	politics
178	sport
258	tech
1200	entertainment
1789	entertainment
...	...
827	entertainment
1936	politics
2039	politics
1051	entertainment
215	sport

735 rows × 1 columns