

گزارش پروژه اول درس شناسایی الگو

برای انجام این پروژه، تمامی مراحل لازم و کارهای انجام شده را در سلولهای بالای کدها توضیح داده شده است. توجه شود که برای کلاسیفایرها و متریکها از ماژول `sklearn` استفاده شده است. همچنین برای بدست آوردن بردارهای `BOW` و `TF-IDF` نیز از ماژول `sklearn` استفاده کردیم. در دو روش `BOW` و `TF-IDF` 10000 تا ویژگی در نظر گرفته ایم. (زیرا برای مقادیر بیشتر مثلا 20000 تا با مشکل کمبود رم در `colab` مواجه شدیم). در روش `bert` هر `review` را به جمله هایش `tokenize` کردیم و هر جمله را به مدل `bert` دادیم. سپس میانگین بردارهای جمله های هر `review` را به عنوان بردار آن `review` در نظر گرفتیم. با توجه به طولانی بودن اجرای روش `svm` برای هیچ یک از حالتها موفق به اجرای این کلاسیفایر نشدیم.

نتیجی به صورت زیر است: (max هر ستون پررنگ شده است).

BOW	AUC	ROC	Precision	Recall	F1_score
Naïve Bayes	0.7707 891684 332626	(array([0. , 0.07959682, 1.], array([0. , 0.62117515, 1.]), array([2. , 1. , 0.])))	0.886415 52511415 52	0.62117 5152993 8803	0.730462 59495308 2
SVM					
Decision tree	0.7188 912443 50226	(array([0. , 0.274949, 1.]), array([0. , 0.71273149, 1.]), array([2. , 1. , 0.])))	0.721621 51216944 07	0.71273 1490740 3703	0.717148 95158369 22
Random Forest	0.7632 894684 21263	(array([0. , 0.17531299, 1.]), array([0. , 0.70189192, 1.]), array([2. , 1. , 0.])))	0.800145 91217910 72	0.70189 1924323 027	0.747805 33537884 61

bert	AUC	ROC	Precision	Recall	F1_score
Naïve Bayes	0.8359 065637 374505	(array([0. , 0.139 87441, 1.]), array([0. , 0.81168753 , 1.]), array([2. , 1., 0.])))	0.853005 46448087 43	0.81168 7532498 7001	0.831833 73982906 68
SVM					
Decision tree	0.7783 088676 452942	(array([0. , 0.218 11128, 1.]), array([0. , 0.77472901 , 1.]), array([2. , 1., 0.])))	0.780315 84884376 76	0.77472 9010839 5665	0.777512 39387431 53
Random Forest	0.8324 467021 319149	(array([0. , 0.144 75421, 1.]), array([0. , 0.80964761 , 1.]), array([2. , 1., 0.])))	0.848329 91073299 52	0.80964 7614095 4362	0.828537 51381441 61

TF_IDF(Removing stop words)	AUC	ROC	Precision	Recall	F1_score
Naïve Bayes	0.716 25134 99460 021	(array([0. , 0.18175273, 1. , , 0.61425543, 1. , , 1., 0.]), array([2., 1., 0.])))	0.771669 76533842 52	0.6142 554297 828087	0.684022 98338604 07
SVM					
Decision tree	0.705 77176 91292 349	(array([0. , 0.28838846, 1. , , 0.699932, 1. , , 1., 0.]), array([2., 1., 0.])))	0.708203 48860739	0.6999 320027 198912	0.704043 45202172 6
Random Forest	0.780 50877 96488 14	(array([0. , 0.14387425, 1. , , 0.7048918, 1. , , 1., 0.]), array([2., 1., 0.])))	0.830490 10367577 76	0.7048 918043 278268	0.762553 81752882 89

TF_IDF(without removing stop words)	AUC	ROC	Precision	Recall	F1_score
Naïve Bayes	0.7167313307467702	(array([0.18463261, 1.], array([0.61809528, 1.], array([2., 1., 0.])))	0.7699935223478998	0.6180952761889524	0.6857333037497226
SVM					
Decision tree	0.7026118955241789	(array([0.29338826, 1.], array([0.69861206, 1.], array([2., 1., 0.])))	0.704245796540462	0.6986120555177793	0.7014176137504518
Random Forest	0.7489900403983841	(array([0.16683333, 1.], array([0.66481341, 1.], array([2., 1., 0.])))	0.7993939976914197	0.6648134074637014	0.7259188085515255

از لحاظ معیار AUC در روش BOW کلاسیفایر naïve bayes از همه بهتر بوده است.

از لحاظ معیار AUC در روش bert کلاسیفایر naïve bayes باز هم از همه بهتر بوده است.

از لحاظ معیار AUC در روش tf idf بدون حذف و با حذف کلمات توقفی کلاسیفایر Random forest از همه بهتر بوده است.

در کل بهترین AUC را به ترتیب روش های bert و tf idf با حذف کلمات توقفی و BOW و tf idf بدون حذف کلمات توقفی است.

با مقیسه

از لحاظ معیار f1 score در روش BOW کلاسیفایر naïve bayes از همه بهتر بوده است.

از لحاظ معیار f1 score در روش bert کلاسیفایر naïve bayes باز هم از همه بهتر بوده است.

از لحاظ معیار f1 score در روش tf idf بدون حذف و با حذف کلمات توقفی کلاسیفایر Random forest از همه بهتر بوده است.

در کل بهترین f1 score را به ترتیب روش های bert و tf idf با حذف کلمات توقفی و BOW و tf idf بدون حذف کلمات توقفی است.

در نهایت بهترین روش استفاده از bert و کلاسیفایر naïve bayes است.