

ابتدا کتابخانه ها و ماژول های مورد نیاز را import میکنیم. سپس دیتای هر فولدر را خواندیم و در قالب یک دیتافریم همراه با ستون lable ذخیره کردیم. در مرحله پیش پردازش ابتدا داده ها را نرمالیز کردم که خروجی آن بعد از نرمالیز کردن چاپ شده. سپس بوسیله ی تابع split هر ایمیل را در قالب یک لیست توکنایز کردم و بعد از حذف stop word ها، به کمک کتابخانه hazm عملیات حذف پیشوند، پسوند و یافتن ریشه کلمات شبیه بهم را انجام دادم. در قسمت وکتورایز کردن به روش tfidf ابتدا کلمات یکتا در داده ی آموزشی و تست را پیدا کردم و به کمک تابع tfidf_vectorize عملیات وکتورایز را انجام میدهم. به وسیله ی تابع chi2 از کتابخانه sklearn هم 500 فیچر که تاثیرگذاری بیشتری دارند را پیدا میکنیم. دو تابع knn پیاده سازی شد که شباهت را با cos similarty و tfidf score محاسبه میکنند. دو تابع knn را یک بار روی وکتورهای اصلی و یک بار روی وکتورهایی که توسط chi2 بدست آوردیم از $k=2$ تا 50 امتحان میکنیم تا بهترین k و بیشترین دقت را بدست آوریم. دقت تابع tfidf_knn روی وکتوری که با chi2 بدست آوردیم 10 درصد بیشتر از حالت عادی شد. ولی دقت تابع cos_sim_knn روی هر دوتا وکتور تقریباً 95 درصد بود و در مجموع از tfidf_knn بهتر عمل کرد. پس مدل نهایی ما روی وکتور منتخب chi2 و cos_sim_knn با $k=3$ آموزش می بیند.