

گزارش کار:

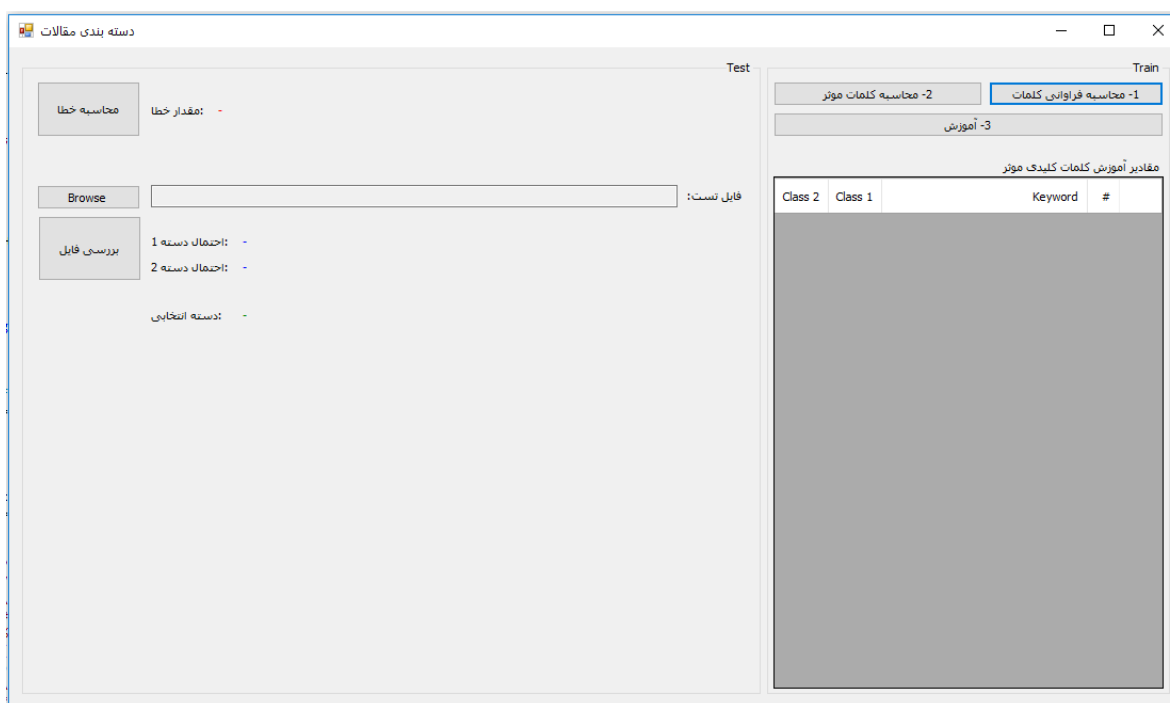
پروژه دسته بندی مقالات با استفاده از کلمات کلیدی موثر و طبقه بندی

فرایند آموزش

در فرایند آموزش، برنامه می بایست با استفاده از مقالاتی که به عنوان مقالات آموزشی در نظر گرفته شده اند استفاده نموده تا به سطحی از دانش برسد که بتواند مقالات را دسته بندی نماید.

در این پروژه دو دسته مقالات و برای هر دسته 100 مقاله آموزشی موجود می باشد که در مجموع 200 مقاله آموزشی در دسترس می باشد.

فرایند آموزش شامل 3 مرحله می باشد که در ادامه مشاهده می کنید:



1- **محاسبه فراوانی کلمات:** در این مرحله برای تمامی مقالات آموزشی در هر دسته، فراوانی کلمات را بدست می آوریم و در فایل به نام words_frequency.txt ذخیره می کنیم.

2- **محاسبه کلمات موثر:** در این مرحله می بایست با استفاده از مقدار فراوانی که برای کلمات در مرحله قبل بدست آمده است آنهایی که بیشترین تاثیر را برای انتخاب یک دسته دارند تعیین نمود که از فرمول زیر مقدار تاثیر هر کلمه در یک دسته بندی مشخص می شود:

$$\text{effectiveCount}(\text{class}) = \text{Sum}(\text{class}) - [\text{Sum}(\text{Sum}(\text{OtherClass}))]$$

effectiveCount(class): مقدار موثر کلمه کلیدی در کلاس

Sum(class): فروانی کلمه کلیدی در کلاس

Sum(Sum(OtherClass)): مجموع فراوانی کلمه کلیدی در دیگر کلاس ها

پس از محاسبه مقدار موثر کلمات هر دسته از مقالات آنها را در فایل به نام words_effective.txt ذخیره می کنیم.

همچنین تعداد 25 کلمه کلیدی که در هر کلاس دارای بیشترین مقدار موثر هستند را به عنوان خصوصیت های کلاس یا کلمات کلیدی نماینده آن کلاس در نظر می گیریم و در فایل class_attributes.txt ذخیره می کنیم.

3- محاسبه مقادیر فراوانی کلمات مورد استفاده در فرمول بیز (آموزش): در این مرحله برای هر یک از کلمات کلیدی نماینده کلاس که مجموعاً 50 کلمه هستند (25 کلمه کلیدی به ازای هر دسته)، فراوانی در هر دسته یا کلاس را محاسبه می کنیم تا از این مقادیر در فرمول بیز برای مقالات تستی بتوان جهت دسته بندی استفاده نمود. برای محاسبه این مقدار در هر کلاس تعداد مقالاتی که دارای یک کلمه کلیدی هستند را می شماریم. در آخر برای تمامی کلمات کلیدی مقدار فراوانی آن در مقالات هر کلاس را خواهیم داشت.

Class 2	Class 1	Keyword	#
5	100	graphics	1
6	100	comp	2
26	55	that	3
33	40	this	4
19	74	crabapple	5
2	19	image	6
14	35	some	7
0	7	polygon	8
100	100	from	9
18	29	about	10
0	16	animation	11
7	20	line	12
19	54	europa	13
19	54	gtefsd	14
0	10	points	15
3	14	point	16
11	33	there	17
10	30	would	18
0	15	program	19
1	12	robert	20

فرایند تست

روش بررسی مقاله تستی: برای بررسی مقاله تستی و تعیین کلاس آن مراحل زیر انجام می شود:

- 1- تعیین کلمات کلیدی مقاله تستی
- 2- انتخاب کلمه کلیدی های معرف هر کلاس از میان کلمات کلیدی مقاله تستی
- 3- انتخاب نسبت فراوانی کلمه کلیدهای انتخاب شده به ازای هر کلاس که در مرحله آموزش برای آن بدست آمده است.
- 4- محاسبه احتمال وجود مقاله در یک کلاس با استفاده از احتمالات بدست آمده برای کلمات کلیدی آن در همان کلاس

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

$\hat{y} = C_k$ for some k as follows:

5- مقایسه احتمالات بدست آمده وجود مقاله به ازای هر کلاس و انتخاب بزرگترین احتمال به عنوان کلاس انتخاب شده برای مقاله تستی.

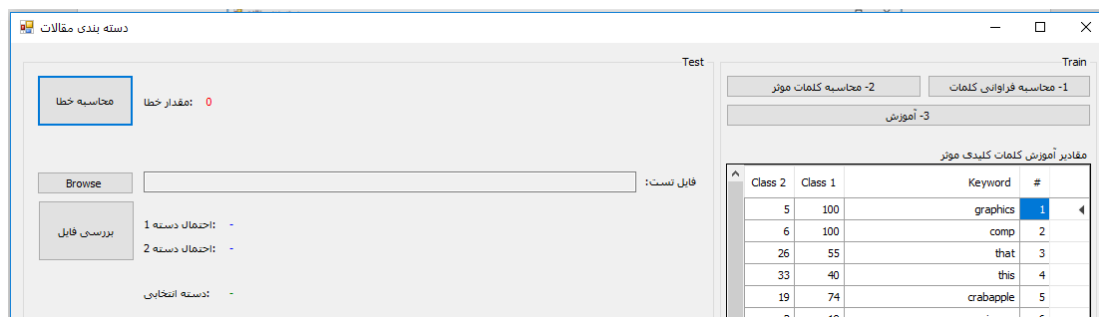
$$\hat{y} = \operatorname{argmax}_{k \in \{1, \dots, K\}} p(C_k) \prod_{i=1}^n p(x_i | C_k).$$

در فرایند تست دو عملیات پیاده سازی شده است که در زیر مشاهده می کنید:

1- **محاسبه خطا:** کلاس انتخاب شده برای مقاله های تستی که 10 مقاله برای دسته اول و 10 مقاله برای دسته دوم می باشد که در مجموع 20 مقاله تستی موجود است را با کلاس معرفی شده برای آنها مقایسه می کنیم و در نهایت از فرمول زیر برای نمایش نسبت خطا استفاده می کنیم. مقدار بدست آمده برای خطای محاسبه شده عددی بین 0 و 1 می باشد که هر چه کوچکتر باشد یعنی سیستم دارای خطای کمتری است. همانطور که مشاهده می کنید برای برنامه پیاده سازی شده عدد 0 بدست آمده که نشان می دهد برنامه تمامی مقالات تستی را در گروه درست طبقه بندی کرده است.

تعداد طبقه بندی نادرست

تعداد کل مقالات تستی



2- بررسی يك فایل انتخابي توسط کاربر: کاربر كه فایل مقاله را انتخاب مي كند و سیستم با روش طبقه بندی بیز احتمال وجود مقاله در هر دو دسته به روشي كه بیشتر توضیح داده شد محاسبه مي كند و دسته اي را كه احتمال بیشتری دارد را به عنوان کلاس انتخابي پیشنهاد مي دهد. كه در شکل زیر مشاهده مي كنید. لازم به ذکر است به علت مساوي بودن تعداد مقالات آموزشی و در نتیجه مساوي بودن احتمال وجود مقالات در هر کلاس كه برابر 0.5 مي باشد در محاسبه احتمال صرف نظر شده است.

