

گزارش پروژه فاز اول

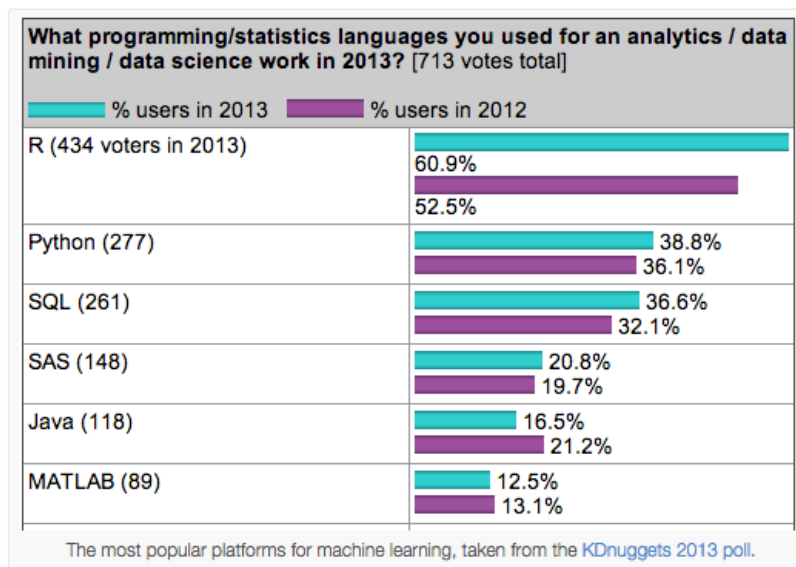
یادگیری

دکتر سلیمانی

تینا خواجه - ۹۳۲۱۰۷۶۱

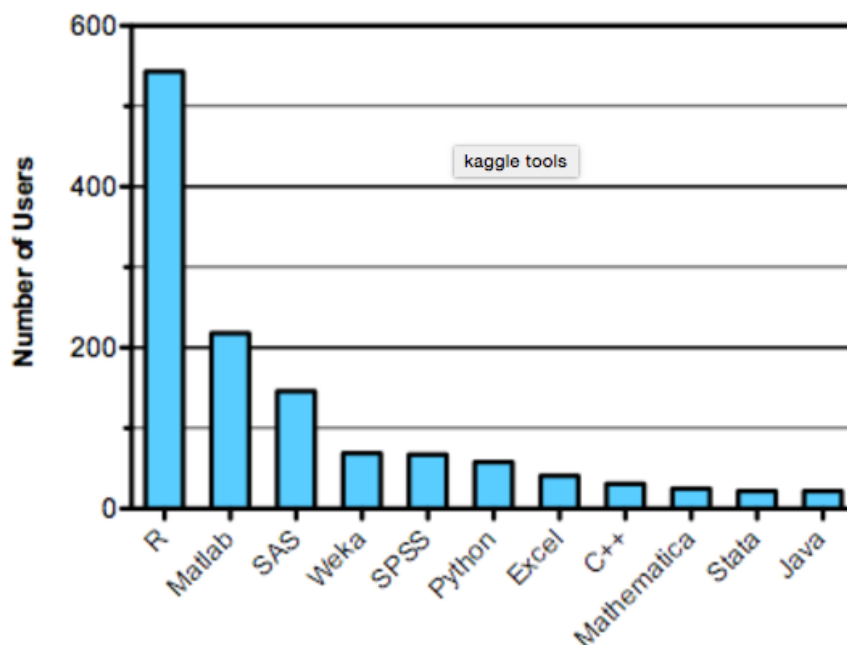
انتخاب زبان برنامه نویسی و کتابخانه مورد استفاده :

قبل از آنکه به انتخاب زبان برنامه نویسی بپردازیم به شناخت زبان های مختلفی که در این حیطه مورد استفاده است می پردازیم تا از میان آنها با توجه به نیاز خود یکی را برگزینیم. در یک نظر سنجی که با عنوان زبان برنامه نویسی مورد استفاده برای کار با داده ها که در سال ۲۰۱۳ برگزار گردید نتایج زیر بدست آمد.



همانطور که در شکل مشخص است زبان های R,Python,Java و Matlab از جمله زبان های برنامه نویسی پرطرف دار برای کار با داده ها هستند.

علاوه بر ای نظر سنجی Kaggle نتایجی را در سال ۲۰۱۱ با عنوان Kaggle's Favorite tools منتشر کرد که در آن مشخص می کند در مسابقات یادگیری ماشین برگزار شده شرکت کنندگان برا ی شرکت در مسابقات از چه زبان هایی استفاده کرده و کدامیک پر طرف دار تر هستند. این نتایج را در زیر مشاهده می نمایید:



همانطور که در این نتایج نیز مشخص است زبان برنامه نویسی R پر طرف دار ترین زبان می باشد. اما قبل از انتخاب زبان مورد نظر باید دقت داشته باشیم که برای پروژه پیش رو چه ابزار هایی لازم است و کدامیک از این زبان ها ما را به هدفمان نزدیک تر می کنند.

در این پروژه هدف ما یادگیری الگوریتم ها و جزییات دقیق پیاده سازی نیست و بیشتر به دنبال شناسایی رفتار هر یک از این روش ها بر روی داده ها هستیم.

ویژگی های هر یک از زبان هایی که کاندیدای انتخاب به عنوان استفاده در این پروژه بوده اند را در زیر می آوریم تا با بررسی هر یک در نهایت انتخاب درستی داشته باشیم:

MATLAB

یکی از خوبی های این زبان برنامه نویسی این است برای محاسبات برداری و ماتریسی طراحی شده و مشکلی با داده های برداری ندارد.

برای کاربرد هایی همچون پردازش سیگنال بسیار ارزشمند است و دارای کتابخانه های وسیعی می باشد. متلب ساده ترین و دقیق ترین زبان برای کار با ماتریس است و برای هر مسئله ای که به صورت بردار ویژگی عددی ظاهر شود بسیار مناسب کار می کند.

خوبی دیگری که این زبان دارد این است که مانند زبانی همچون C درگیر کار های low level در هنگام کد زدن نمی شویم.

این زبان برای مواردی که قصد یادگیری الگوریتم و پیاده سازی آن با جزییات توسط خود فرد می باشد مناسب است اما ابزار قوت مندی به عنوان مرجعی از توابع و روش های آماده و مورد استفاده کاربران محسوب نمی شود.

R

مشابه با متلب این زبان نیز محاسبات برداری و ماتریسی را با سرعت بالایی انجام می دهد. R زبانی است که با هدف آنالیز های آماری و تمرکز بر روی این اصل بوجود آمد است. این زبان برای کار با داده ها و استخراج ویژگی های آماری موجود در آنها بسیار مناسب می باشد و تمام ابزار لازم برای این کار از جمله توابع، نمودار ها و... را تحت عنوان پکیج های مختلفی در اختیار کاربران خود قرار می دهد. هدف اصلی زبان کار آنالیز داده ها می باشد و بعد از آن این زبان در جهتی گسترش یافت که برای کاربرد های یادگیری ماشین نیز مورد استفاده قرار بگیرد. اما باید دقت داشته باشیم هرچند که این زبان دارای پکیج های مختلفی با هدف فراهم آوردن الگوریتم ها و ابزار های مورد استفاده در زمینه یادگیری ماشین است، اما این پکیج ها نقاط ضعفی همچون نداشتن داکيومنت ها با توضیحات دقیق، نداشتن انسجام بین پکیج ها و... را دارند که منجر می شود گاهی کار بر در هنگام کار کردن با آنها دچار سر در گمی شود (همانطور که در ادامه و طبق تجربه اینجانب توضیح داده می شود) کار کردن با پکیج های مختلف R و رفتن از یک پکیج به پکیج دیگر معمولاً کاربر را با interface ای متفاوت از پکیج قبلی مواجه می کند و باعث سردر گمی می شود.

برای کارهایی همچون feature extraction مناسب نمی باشد و بهتر است در کنار R از ابزاری دیگر همچون vowpal wabbit برای این هدف استفاده کنیم.

Python

به تازگی پایتون جایگاه خوبی در بین برندگان رقابت های یادگیری ماشین پیدا کرده است که با استفاده از کتابخانه scikit-learn در رقابت ها شرکت کرده بودند. کتابخانه scikit-learn یکی از بهترین کتابخانه های موجود برای روش ها و الگوریتم های یادگیری ماشین محسوب می شود که راهنمای استفاده بسیار قوی همراه با توضیحات و مثال هایی را برای کاربران خود فراهم می آورد.

یکی از دغدغه هایی که هنگام انتخاب زبانی مثل python وجود دارد این است که سرعت اجرا و محاسبات بسیار کند باشد. این کند بودن سرعت برای اجرای الگوریتم هایی همچون یادگیری ماشین که با داده های برداری کار می کنند در زبان پایتون کاملاً مشاهده می شود. این کندی در سرعت اجرا را می توان از طریق تبدیل به زبان C حل کرد.

یکی دیگر از مزیت های کتابخانه scikit-learn در کنار جامع بود، کامل بود، توجه به جزئیات مربوط به یادگیری ماشین و همچنین داشتن داکيومنتی کامل و دقیق به همراه مثال های کافی و نیز محبوبیت در بین متخصصان مختلف، توجه به مشکل این زبان برنامه نویسی برای کند بودن اجرا است. این مشکل با تبدیل C حل شده و کاربران بدون داشتن دغدغه سرعت اجرای بسیار متفاوت با دیگر زبان ها می توانند از این کتابخانه استفاده نمایند.

در کنار توضیحات بالا که در باره مزیت های کتابخانه scikit-learn داده شد، این مزیت بزرگ را نیز باید در نظر بگیریم که زبان پایتون زبانی سطح بالا است و کاربران خود را از کار کردن با خیلی جزئیات راحت می کند و تجربه خوبی از برنامه نویسی را در اختیار کاربران خود قرار می دهد.

در کنار کتابخانه scikit-learn کتابخانه numpy کمک می کند که زبان پایتون تبدیل به زبان های مناسب برای یادگیری ماشین تبدیل شود.

Java

یکی از خوبی های این زبان برنامه نویسی سرعت خوب اجرای آن است. همچنین برای این زبان کتابخانه هایی

همچون Mahout , Weka

جمع بندی و نتیجه گیری :

در بالا زبان های مختلف را از جنبه های مختلف کتابخانه های موجود ، سرعت اجرا ، راحتی در برنامه نویسی و... بررسی کردیم. با وجود توضیحات بالا انتخابی باقی نمی ماند و واضح است که به سراغ کدام زبان باید برویم : پایتون!

با وجود مشخص بودن این موضوع به دلیل مشکل در نصب کتابخانه بر روی سیستم اینجانب در مرحله بعد به سراغ زبان برنامه نویسی R رفته و سعی در انجام پروژه به زبان R داشتم که بعد از یک روز به نامناسب بودن پکیج های یادگیری ماشین و نکات ذکر شده در جهت سردر گمی کاربر پی برده و دوباره به سراغ زبان برنامه نویسی پایتون آمده و هر طور بود کتابخانه را نصب نمودم. گزارشات پیش رو حاصل پیاده سازی به زبان پایتون می باشد.

پیاده سازی و نتایج :

دوکلاسه

feature selection :

در ابتدا ویژگی های مناسب را انتخاب می کنیم :

روش های مختلفی برای این کار وجود دارند که از این جمله می توان به روش های Removing features: with low variance, Univariate feature selection, Recursive feature elimination, Feature selection using SelectFromModel اشاره کرد. روش هایی که ویژگی را تنها با در نظر گرفتن خود داده انتخاب می کنند زیاد مورد علاقه نیستند و دوست داریم ویژگی ها را همه را با هم در نظر گرفته و تعیین کنیم که کدام ها خوب هستند و کدامها بد. برای انتخاب ویژگی های مناسب می توانیم از ماژول sklearn.feature_selection در کتابخانه scikit

learn استفاده نماییم. یکی از بهترین انواع انتخاب ویژگی Lasso است که جوابی sparse دارد و در نهایت آن ویژگی هایی که ضرایب غیر صفر برای آنها بدست آمده را انتخاب می نماید.

پارامتر C موجود در این تابع مشخص کننده میزان sparsity می باشد. فرم به صورت زیر قابل مشاهده است :

```
class sklearn.svm.LinearSVC(penalty='l2', loss='squared_hinge', dual=True, tol=0.0001, C=1.0, multi_class='ovr',
fit_intercept=True, intercept_scaling=1, class_weight=None, verbose=0, random_state=None,
max_iter=1000) \[source\]
```

پارامترهای استفاده شده در این حالت عبارتند از :

C=0.2, penalty="l1"

در این حالت تعداد ویژگی ها از ۴۱۴ به ۴۳ ویژگی کاهش یافت.

۱) دسته بندی دو کلاسه : svm

ابتدا برای در نظر گرفتن اهمیت برابر برای تمام داده ها، با توجه به اینکه داده ها را نمی شناسیم که آیا یکی بر دیگری اهمیت دارد یا نه، داده های را scale می کنیم :

برای این کار از پکیج sklearn.preprocessing استفاده می کنیم و با استفاده از تابع (MaxAbsScaler) داده ها را به بازه [-1,1] منتقل می نماییم. بدین ترتیب و از طریق یکسان کردن واریانس و رنج داده ها اثر همه در بدست آوردن مدل نهایی را یکسان می کنیم.

کد مربوط به این قسمت در برنامه نوشته شده عبارت است از:

```
import numpy as np
max_abs_scaler = preprocessing.MaxAbsScaler()
xtrainScaled = max_abs_scaler.fit_transform(xtrain)
```

همچنین علاوه بر داده های train داده های validation و همچنین test را نیز به همین صورت scale می نماییم.

برای اجرای الگوریتم svm بر روی داده ها به سراغ پکیج svm می ریم. برای استفاده از این الگوریتم از انواع مختلفی از تنظیمات استفاده می نماییم که در این حالت پارامترهای زیر را برای آن تغییر می دهیم :

kernel, C, gamma

در این حالت باید مدل های مختلف را بر روی داده train آموزش ببینیم و بهترین مدل را بر اساس کوچک ترین خطا بر روی داده های validation محاسبه نماییم. در صورتی که داده های train و validation جدا نبوند و می توانستیم از cross validation استفاده نماییم می توانستیم از تابع `sklearn.grid_search.GridSearchCV` برای پیدا کردن بهترین مقادیر پارامتر ها با یک دستور استفاده نماییم. اما به دلیل اینکه اینجا داده های train و validation جدا هستند باید ابتدا مدل را بر روی داده های train آموزش دهیم و بعد بر روی validation امتحان کنیم و مدل با پارامتر هایی را به عنوان بهترین مدل انتخاب کنیم که کمترین خطا بر اساس معیار مورد بررسی را دارد.

کرنل های مختلفی که می توانیم استفاده کنیم اما به دلیل اینکه هیچ اطلاعی از داده ها نداریم که کدام کرنل برای آن مناسب تر است و همچنین اینکه اجرای تمام کرنل ها با توجه به زمان اجرا و اجرای برنامه بر روی لپ تاپ غیر ممکن می باشد به همین دلیل چون کرنل rbf حالت کلی تری دارد و برای وقتی که اطلاعی از داده ها نداریم معمولا جواب مناسبی می دهد به سراغ این کرنل رفته و به دنبال پیدا کردن بهترین نوع پارامتر ها برای آن می گردیم.

البته این قابلیت را با استفاده از این کتابخانه خواهیم داشت که کرنل توسط کاربر مشخص شود ولی از آن استفاده نمی کنیم.

پارامتر C نیز بیان می کند که چقدر missclassification را جریمه می کنید و این میزان بر روی margin اثر گذار است هرچه C بزرگ تر باشد margin کوچک تر می شود. برای مشخص کردن بهترین C مقادیر مختلف 10^3 تا 10^{-6} را چک می کنیم و C ای که بهترین معیار را بر روی داده ها داشته باشد انتخاب می کنیم و سپس به سراغ پیدا کردن بهترین gamma با توجه به این c خاص می گردیم. برای پیدا کردن بهترین c نیز پارامتر گاما را برای تمام حالات ۰/۸ قرار می دهیم. نتایج دو معیار در جدول زیر آمده است :

C	10^{-6}	10^{-5}	10^{-4}	10^{-3}	10^{-2}	0.1	1	10	100	1000
Accuracy	0.6189 999999 999999 9	0.6189 999999 999999 9	0.6189 999999 999999 9	0.6189 999999 999999 9	0.7072 000000 000000 5	0.8042 000000 000000 3	0.8299 999999 999999 6	0.8461 999999 999999 5	0.8329 999999 999999 6	0.8145 999999 999999 9
F-Measure	0.7646 695491 043854 5	0.7646 695491 043854 5	0.7646 695491 043854 5	0.7646 695491 043854 5	0.8078 235757 416644 1	0.8534 650501 421943 1	0.8695 518723 143032	0.8463	0.8673 550436 854645 7	0.8522 709163 346614 5

همانطور که در جدول بالا نیز مشخص می باشد بهترین c عبارت است از ۱۰. حال با در نظر گرفتن این مقدار به عنوان c به دنبال مقدار gamma می رویم که بهترین دقت را به ما بدهد نتایج را در جدول زیر می توان مشاهده بنمایید:

gamma	10 ⁻⁶	10 ⁻⁵	10 ⁻⁴	10 ⁻³	10 ⁻²	0.1	1	10
Accuracy	0.618999 99999999 999	0.618999 99999999 999	0.753199 99999999 998	0.792399 99999999 999	0.807200 00000000 003	0.846199 99999999 995	0.849999 99999999 998	0.618999 99999999 999
F- Measure	0.764669 54910438 545	0.764669 54910438 545	0.828611 11111111 108	0.841961 02314250 909	0.851326 34176434 308	0.879674 54232514 486	0.886122 07713331 304	0.764669 54910438 545

همانطور که در بالا قابل مشاهده است بهترین مقدار برای gamma عدد یک است که منجر به بدست آوردن بهترین معیار شده است. در نهایت بهترین دقتی که برای داده های validation به دست آمده است عبارت است از دو خانه رنگی شده.

همچنین cache size را که میزان RAM اختصاص داده شده به برنامه می باشد را 500 قرار می دهیم به این امید که سرعت اجرا بالا رود. که البته همچنان سرعت پایین است و برای یک بار اجرا و بدست آوردن نتایج بالا حدودا به یک ساعت زمان احتیاج است.

برای محاسبه معیار دقت از تابع `sklearn.metrics.accuracy_score` استفاده می کنیم که prototype تابع به صورت زیر می باشد :

```
sklearn.metrics.accuracy_score(y_true, y_pred, normalize=True, sample_weight=None)
```

همچنین برای حالت دو کلاسه نیز از معیار f1 استفاده می کنیم که معادل تابع زیر می باشد :

```
sklearn.metrics.f1_score(y_true, y_pred, labels=None, pos_label=1, average='binary', sample_weight=None)
```

حال زمان آن رسیده که با توجه به بهترین مدلی که بدست آورده ایم برچسب های داده های test را پیش بینی نماییم و نتایج را در فایل بنویسیم. پیش بینی مقادیر نهایی در فایل `bin_test_svm` برای این قسمت قرار داده شده

و ارسال گردیده است همانطور که در فایل کد این قسمت با نام bin_svm.py نیز آمده است برای خواندن اطلاعات از دستور زیر استفاده نمایید.

```
import pickle
f = open('bin_test_svm', 'r')
yTestPredic = pickle.load(f)
```

KNN:

برای پیاده سازی این قسمت به سراغ تابع

```
class sklearn.neighbors.KNeighborsRegressor(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30,
p=2, metric='minkowski', metric_params=None, n_jobs=1, **kwargs) [source]
```

می رویم که یک پیاده سازی از KNN است و در آن پیدا کردن تعدا همسایه ها مشکل به نظر می رسد. در ابتدا از روش استفاده شده در مرحله قبل برای feature extraction استفاده می کنیم. پارامتر هایی که برای این مسئله تنظیم می کنیم عبارتند از :

`n_neighbor = sqrt(number of features)`

`weights = distance`

distance در واقع میزان اثر گذاری هر یک از نمونه ها را مشخص می نماید که متناسب است با عکس

فاصله.

همانطور که گفته شده پیدا کردن بهترین K کار ساده ای نیست و به دلیل اینکه هیچ اطلاعاتی در باره داده ها نداریم نمی توانیم یک بازه مشخص کرده و د ر آن بازه به دنبال K بگردیم. یک قانون سر انگشتی وجود دارد که جذر تعداد نمونه ها در نظر بگیریم. مشکل این classifire هم همین است که انتخاب K مهم است و به شدت می تواند نتیجه را تحت تاثیر قرار دهد و قانون خاصی برای آن نداریم. به منظور انتخاب درست مقدار K مقدار دو معیار خواسته شده را برای k های مختلف محاسبه نمودیم :

K	5	10	25	50	100	250
Accuracy	0.839799999 99999999	0.849400000 00000004	0.848999999 99999998	0.847400000 00000004	0.837600000 00000001	0.821799999 99999997
F-Measure	0.873679230 40529882	0.882067345 34064202	0.884006759 8709479	0.883919062 83280086	0.877637130 8016877	0.868253733 55019975

K	500	1000	2500	5000	10000
Accuracy	0.804000000 00000005	0.7792	0.728400000 00000005	0.660399999 99999999	0.618999999 99999999
F-Measure	0.858135495 07817033	0.844507042 25352113	0.817374932 75954798	0.783639143 73088672	0.764669549 10438545

بنابراین همانطور که از اطلاعات درج شده در جدول مشخص است بهترین مقدار k که البته دارای تفاوت اندکی در دو معیار محاسبه شده و مشخص شده در بالا است، می تواند مقداری برابر با ۲۰ در نظر گرفته شود. پس مقادیر y را با این k تخمین می زنیم و در فایل `bin_test_knn` ثبت می نماییم. همانطور که در قسمت قبل نیز اشاره شد برای بدست آوردن اطلاعات درون فایل از دستورات ذکر شده در قسمت قبل استفاده نمایید.

Naive bayes :

این روش احتیاج به scaling ندارد و این کار بر روی نتیجه اثری نخواهد داشت بنابراین این کار را بر روی داده های این قسمت انجام نمی دهیم. همچنین این روش با توجه به ماهیت خود که ویژگی ها را مستقل از هم دیگر در نظر می گیرد و در نتیجه دارای سرعت اجرای بالایی است احتیاج به انتخاب ویژگی ها ندارد. برای اجرای الگوریتم از کلاس GaussianNB استفاده می کنیم که در آن الگوریتم GaussianNaive bayse را برای اجرای classification استفاده می کند. likelihood ویژگی ها فرض می شود که دارای فرم زیر می باشد.

$$P(x_i | y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right)$$

برای اینکه اطمینان حاصل شود که این روش بریا بهبود عملکرد خود احتیاج به کاهش ابعاد ندارد و روش کاهش ابعاد را پیاده سازی کرده و معیار های بدست آمده را با داده های بدون کاهش مقایسه نمودیم. که در هر دو روش خود داده ها معیار های بهتری را گزارش می نمودند. که البته این موضوع از با کمی دقت بر روی روش naive bayse قابل فهم است.

معیار های دقت بدست آمده در این حالت نیز برابر است با :

Accuracy = 0.79500000000000004

F-Measure = 0.84472049689440998

Logistic Regression :

برای این قسمت از کلاس و توابع مربوط به آن با نام `sklearn.linear_model.LogisticRegression` استفاده می نماییم. prototype تابع به صورت زیر می باشد :

```
class sklearn.linear_model.LogisticRegression(penalty='l2', dual=False, tol=0.0001, C=1.0, fit_intercept=True, intercept_scaling=1, class_weight=None, random_state=None, solver='liblinear', max_iter=100, multi_class='ovr', verbose=0, warm_start=False, n_jobs=1) ¶
```

[\[source\]](#)

ماهیت روش logistic regression به گونه ای است که احتیاج به نرمال سازی داده ها ندارد و این کار اثری بر روی الگوریتم ندارد پس پیش پردازشی در این حالت نداریم.

پارامتر هایی که برای داده هایمان ست خواهیم کرد `penalty` و `C` می باشد که با توجه به معیار دقت بدست آمده بر روی ترکیب آنها تصمیم می گیریم کدام یک را برگزینیم. جدول زیر مقادیر مختلف `C` و `penalty` را بررسی نموده است که نتایج در آن آورده شده اند :

Penalty = L1

C	0.0001	0.001	0.01	0.1	1	10
Accuracy	0.381000000000000001	0.618999999999999999	0.728400000000000005	0.788399999999999999	0.823400000000000002	0.823400000000000002
F-Measure	0.0	0.76466954910438545	0.80677290836653381	0.83792892156862742	0.86175043056207923	0.86162043566839053

Penalty = 'L2'

C	0.0001	0.001	0.01	0.1	1	10
Accuracy	0.618999999999999999	0.618999999999999999	0.726400000000000005	0.7984	0.820999999999999995	0.823999999999999995
F-Measure	0.76466954910438545	0.76466954910438545	0.81573275862068972	0.84805547181187824	0.86143365846106212	0.862500000000000004

همانطور که در بالا مشاهده می کنیم بیشینه مقدار به ازای $L2$ و نیز $c=10$ بدست می آید. در نهایت خروجی برای داده های test محاسبه می شود و در فایل bin_test_logistic قرار می گیرد. برای خواندن اطلاعات فایل همانطور که در قسمت اول نیز اشاره شد به صورت زیر عمل می کنیم

```
import pickle
f = open('bin_test_logistic', 'r')
yTestPredic = pickle.load(f)
```

(۲) چند کلاسه :

SVM :

برای این حالت نیز در ابتدا به سراغ انتخاب ویژگی می رویم. مشابه حالت دو کلاسه در اینجا نیز از تابع LinearSVC با پارامتر $C=0.02$ استفاده می کنیم. که حاصل آن استخراج همچنین برای این حالت نیز مشابه با حالت دو کلاسه ابتدا باید داده ها را scale نماییم تا نتیجه بهتری را بگیریم و وزن تمان ویژگی ها برابر باشد.

در ابتدا با استفاده از linearSVC پیاده سازی می کنیم که در واقع همان svm است با کرنل linear نتایج برای C های مختلف به صورت زیر می باشند :

Linear SVM(one vs rest)						
C	0.0001	0.001	0.01	0.1	1	10
Accuracy	0.21379999999999999	0.39760000000000001	0.49059999999999998	0.54679999999999995	0.56699999999999995	0.56940000000000002
Macro_F1	0.094932583822372629	0.25277775060149371	0.35196556779517379	0.43757001072783108	0.47390230277022316	0.47846796291208266
Micro_F1	0.21379999999999999	0.39760000000000006	0.49059999999999998	0.54679999999999995	0.56699999999999995	0.56940000000000002

اما همانطور که مشخص است در بهترین حالت دقت ۰/۵ است برای معیار های مختلف که دقت مناسبی به نظر نمی رسد برای همین به سراغ svm غیر خطی می رویم. مطابق با قبل به سراغ کرنل rbf می ریم. همچنان احتیاج به انتخاب مناسب دو پارامتر C و گاما داریم.

مثل قبل از میان مقادیر مختلف c که در جدول زیر قابل مشاهده می باشد و با قرار دادن متغیر گاما به صورت ثابت در این قسمت معیار های سنجش محاسبه می نماییم. و از بین آنها بهترین C را بر می گزینیم :

RBF SVM(choose C)

C	0.0001	0.001	0.01	0.1	1	10	100	1000
Accuracy	0.042200 00000000 0001	0.042200 00000000 0001	0.1158	0.543799 99999999 995	0.604999 99999999 998	0.617399 99999999 995	0.604400 00000000 005	0.604199 99999999 996
Macro_F1	0.001883 31481258 33991	0.001883 31481258 33991	0.050362 02171209 0518	0.433332 61366689 863	0.542105 80848942 946	0.6175	0.549657 69707987 389	0.549384 28715956 8
Micro_F1	0.042200 00000000 0001	0.042200 00000000 0001	0.1158	0.543799 99999999 995	0.604999 99999999 998	0.6176	0.604400 00000000 005	0.604199 99999999 996

همانطور که در بالا مشخص است، مقداری از c که منجر به مقدار بهتر معیار ها می شود عبارت است از ۱۰. بنابراین با این مقدار به سراغ تعیین پارامتر گاما می ریم.

RBF SVM(choose gamma, C = 10)

gamma	0.00001	0.0001	0.001	0.01	0.1	1	10
Accuracy	0.0422000 000000000 01	0.2560000 000000000 1	0.5590000 000000000 5	0.5887999 999999999 9	0.6173999 999999999 5	0.6191999 999999999 7	0.0422000 000000000 01
Macro_F1	0.0018833 148125833 991	0.1478235 460454012 4	0.4613336 004464463 9	0.5280101 834346472 3	0.5635132 877850447 3	0.5568637 65065711	0.0018833 148125833 991
Micro_F1	0.0422000 000000000 01	0.2560000 000000000 1	0.5590000 000000000 5	0.5887999 999999999 9	0.6173999 999999999 5	0.6191999 999999999 7	0.0422000 000000000 01

همانطور که در بالا قابل مشاهده است به ازای $\text{gamma} = 1$ معیار ها مقادیر مناسبی خواهند داشت. بنابراین هر چند دقت به اندازه دلخواه بالا نیست اما پارامتر ها به صورت زیر تنظیم می شوند :

$C=10, \text{gamma}=1$

همانطور که می بینیم کرنل rbf بهتر از linear عمل کرده است پس برای پیش بینی به سراغ کرنل rbf می

رویم.

مقادیر پیش بینی شده با این پارمتر ها در فایل با نام multiClass_test_svm ذخیره شده است. برای

خواندن اطلاعات فایل از دستور زیر می توان استفاده کرد:

```
f = open('multiClass_test_svm', 'r')
```

```
yTestPredic = pickle.load(f)
```

KNN

برای این قسمت نیز ابتدا ویژگی های مورد احتیاج را انتخاب می کنیم. و سپس با استفاده از تابع

NearestNeighbors استفاده شده در قسمت دو کلاسه دسته بندی را انجام می دهیم.

مثل قبل چالشی که با آن مواجه هستیم انتخاب درست k می باشد. از میان مقادیر مختلف k نتیجه معیار

های مختلف را مشخص کرده و در زیر آورده ایم.

KNN(choose K)

K	5	10	20	25	50	100
Accuracy	0.556599999 99999998	0.581400000 00000003	0.595199999 99999995	0.5948	0.581799999 99999998	0.571999999 99999995
Macro_F1	0.556599999 99999998	0.581400000 00000003	0.5953	0.5948	0.581799999 99999998	0.571999999 99999995
Micro_F1	0.496085878 25795325	0.510309151 16049758	0.5954	0.522547073 81123155	0.501578119 55670861	0.475334247 11303635

KNN(choose K)-1

K	250	500	1000	2500	5000	10000
Accuracy	0.5444	0.523000000 00000002	0.498800000 00000002	0.447799999 99999998	0.391000000 00000001	0.264799999 99999998
Macro_F1	0.5444	0.523000000 00000002	0.498800000 00000002	0.447799999 99999998	0.391000000 00000001	0.264799999 99999998
Micro_F1	0.436653093 22616946	0.405214551 78261467	0.369956766 35682828	0.309474211 4277011	0.251118491 84820673	0.134891313 01037066

همانطور که در بالا می بینیم بیشینه مقدار معیار های مختلف به ازای $K = 20$ بدست می آید. به همین دلیل classfire را بر روی مجموعه داده های train و validation (به منظور بالا رفتن دقت از ترکیب داده های train و validation استفاده می کنیم) آموزش می دهیم و با استفاده از آن برچسب ها را پیش بینی می کنیم.

دلیل اینکه کلا میزان دقت برای داده های چند کلاسه پایین تر از دو کلاسه آمده است این است که با افزایش تعداد کلاس ها به 40 کلاس و ثابت باقی ماندن تعداد نمونه ها پیدا کردن الگوی درست در داده برای classification سخت تر می شود.

مقادیر پیش بینی شده با این پارمتر ها در فایل با نام multiClass_test_knn ذخیره شده است. برای خواندن اطلاعات فایل از دستور زیر می توان استفاده کرد:

```
f = open('multiClass_test_knn', 'r')
yTestPredic = pickle.load(f)
```

Logistic Regression :

در حالت چند کلاسه تابع logistic regression این کتابخانه از روش one vs rest استفاده می کند. برای این کار مشابه با روش دو کلاسه عمل می کنیم و خطای validation را محاسبه می کنیم و با ترکیب داده های train و validation و به دست آوردن مجموعه بزرگ تری برای آموزش مدل را آموزش داده و برچسب ها را پیش بینی می نماییم.

معیار های دقت خواسته شده برای مقادیر مختلف C و penalti در زیر آورده شده است.

Penalty = 'L1'

C	0.001	0.01	0.1	1	10	100
Accuracy	0.042200000 000000001	0.124799999 99999999	0.474999999 99999998	0.633000000 00000001	0.656399999 99999998	0.615399999 99999995
Macro_F1	0.001883314 8125833991	0.032629437 332303356	0.349612520 9948191	0.558293256 56513251	0.603586704 44417484	0.563584421 19169105
Micro_F1	0.042200000 000000001	0.124799999 99999999	0.475000000 00000003	0.633000000 00000001	0.656399999 99999998	0.615399999 99999995

Penalty = 'L2'-1

C	0.001	0.01	0.1	1	10	100
Accuracy	0.186	0.353200000 00000001	0.537200000 00000001	0.635000000 00000001	0.6714	0.656000000 00000003
Macro_F1	0.079932101 744003753	0.216759512 17294889	0.405547645 94562093	0.554785914 64972493	0.6714	0.601527691 45099272
Micro_F1	0.186	0.353200000 00000001	0.537200000 00000001	0.635000000 00000001	0.6714	0.656000000 00000003

همانطور که در تصویر مشخص است به ازای $\text{Penalty} = \text{L2}$ و همچنین $C = 10$ بهترین جواب را داریم. پس مقادیر خروجی test را با این پارامتر ها آموزش داده و نتیجه را در فایل با نام multiClass_test_logistic ثبت می نماییم.

Naive bayes

در این قسمت مشابه با قسمت دو کلاسه عمل می کنیم. دقت های خواسته شده عبارتند از :

Accuracy = 0.5966000000000002

F1Macro = 0.56125307629578935

F1Micro = 0.5966000000000002

مقادیر خروجی test را با این پارامترها آموزش داده و نتیجه را در فایلی با نام multiClass_test_naive

ثبت می نمایم.

: Regression

: Linear Regression

برای این قسمت پیاده سازی رگرسیون خطی را با regularization term از نوع ridge انجام میدهیم. در

این حالت خطی و ساده باید پارامتر alpha که بیان کننده میزان جریمه ضرایب است را مشخص نماییم.

برای بدست آوردن معیار RMSE نیز از تابع mean_squared_error موجود در sklearn.metrics استفاده

می کنیم و از آن جذر می گیریم.

MAE نیز مخفف mean absolute error است. که برا یابین معیار نیز از تابع mean_absolute_error

استفاده می کنیم.

در ابتدا بدون transformation از فضای ویژگی استفاده می کنیم و به ازای alpha های مختلف میزان خطا

را محاسبه نموده تا کمترین خطا را مشخص نماییم.

alpha	0.00001	0.0001	0.001	0.01	0.1	1	10	100	1000
RMSE	15.4219 8750412 214	15.4219 7920579 222	15.4218 9626587 9724	15.4210 7119233 6685	15.4132 4028802 4154	15.3670 1140550 973	15.6119 1486373 7745	17.7505 3275830 4594	21.9754 2223991 7446
MAE	10.7678 3465813 4829	10.7678 2364427 8733	10.7677 1352608 4265	10.7666 1910718 8298	10.7561 9234106 7393	10.6921 8737915 2946	10.8427 6764571 3504	12.5812 4779252 7186	15.3101 3049662 9936

میزان کمینه خطا به ازای $\alpha = 1$ است.

همچنین رگرسیون را با regularization term از نوع lasso نیز پیاده سازی کردیم که دقت بالاتری نسبت به قبل دریافت نمودیم.

پس سراغ polynomial transformation می رویم....!!!! اما به دلیل بالا بودن تعداد ویژگی ها و عدم توانایی لپ تاپ اینجانب و قفل کردن سیستم در حین اجرا امکان ثبت نتایج ممکن نبود. بنابراین با همین حالت ساده مقادیر خروجی را پیش بینی می کنیم.

KNN Regression

برای این مورد نیز از تابع نوشته شده در کتابخانه استفاده می کنیم. و میزان اثر گذاری هر همسایه را از طریق عکس میزان فاصله ای که با نقطه مورد نظر دارد دخیل می کنیم. در این حالت نیز باید تعداد همسایه ها را به گونه ای مشخص نماییم که کمترین خطای پیش بینی را داشته باشیم.

K	5	10	20	50	100
RMSE	13.8317490338 42394	13.43145538 8996383	13.28395205 2376959	13.57462933 0814451	13.97680091 1501243
MAE	7.32074279999 99997	7.484569400 0000007	7.690298099 9999997	8.207610240 0000011	8.713486460 0000004

همانطور که در جدول بالا مشاهده می کنیم مقدار کمینه دو خطا یکی نیست، برای RMSE تعداد ۲۰ همسایه مناسب است و برای MEA تعداد ۵ همسایه. حال باید کمی دقت کنیم که این تفاوت ناشی از چیست؟ این تفاوت ناشی از ماهیت این دو خطا است. RMSE وزن بیشتری به خطا های بزرگ می دهد در حالی که MAE برای تمام خطا ها وزن یکسانی را در نظر می گیرد. اینجا دوست نداریم که پیش بینی مان پرت باشد و خیلی تفاوت داشته باشیم بنابراین سعی می کنیم که خطا های بزرگ را کمینه کنیم، پس به سراغ RMSE می رویم و با تعداد ۲۰ همسایه پیش بینی را انجام می دهیم. نتایج در فایلی با نام reg_test_KNNReg قرار دارند.

توضیح درباره کد :

همانطور که در بالا مشخص است برای پیدا کردن پارامتر های مختلف بسته به نوع مسئله و گاهی احتیاج به چک کردن بیش از یک راه حل برای پیدا کردن جواب بهتر کد ها معمولا دارای بلاک های کامنت شده هستند که این قسمت های کامنت شده در واقع قسمت هایی از کد می باشند که برای بدست آوردن پارامتر های بهینه مورد استفاده بوده اند.