

## به نام خدا

**مسئله :** پیاده‌سازی سیستم تشخیص نفوذ براساس مقاله با عنوان A Study on NSL-KDD Dataset for

Intrusion Detection System Based on Classification Algorithms منتشر شده در سال ۲۰۱۵ است.

### مشخصات پروژه:

\_\_ زبان پایتون، نسخه 3.8.2

\_\_ کتابخانه numpy

\_\_ کتابخانه matplotlib

\_\_ کتابخانه scikit-learn

### توضیحات پروژه:

پیاده‌سازی پروژه شامل شش مرحله به شرح زیر است:

- مرحله یک – خواندن داده

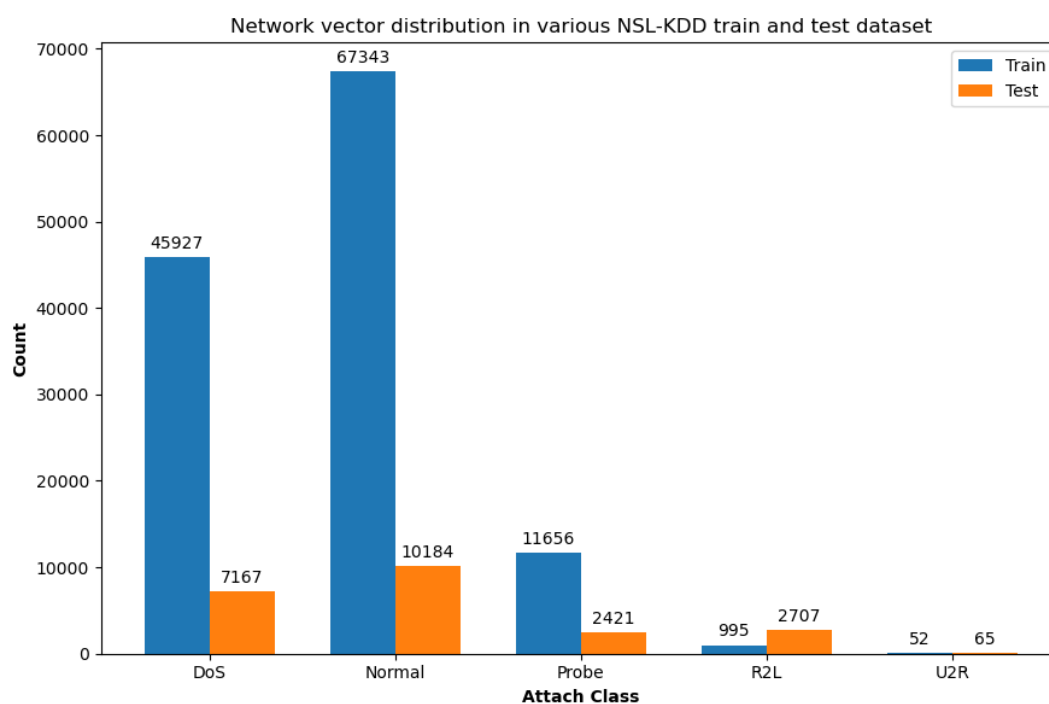
دو فایل KDDTrain+.txt ، KDDTest+.txt شامل مجموعه داده آموزشی و تست در پوشه data در مسیر پروژه قرار داده شده است. برای خواندن داده یک تابع با نام load\_NSL\_KDD در فایل datasets.py تعریف کردیم که تابع به عنوان ورودی مسیر فایل را می‌گیرد و خروجی نمونه های داده ، برچسب هر نمونه و نام ستون های نمونه را برمی گرداند. داخل تابع مسیر فایل چک می‌شود اگر فایلی با این مسیر موجود نبود خطا نمایش داده می‌شود وگرنه خط به خط از فایل خوانده و به وسیله کاما از هم مقادیر آن جدا می‌شود و داخل یک لیست ریخته می‌شود. در لیست یکی مانده به آخرین خانه مقدار برچسب است و بقیه خانه های لیست مقدار نمونه های ما را تشکیل می‌دهد. با توجه به مقاله برچسب ها بر پنج نوع (normal, Dos, Probe, R2L, U2R) است که 4 مورد Dos ، Probe ، R2L و U2R نوع حملات را مشخص می‌کند و هر کدام از این نوع حملات خود نیز شامل زیر بخش های دیگری هستند. در دیتا، برچسب ها شامل Normal

و زیربخش‌های حملات است چون مسئله ما از نوع طبقه‌بندی با پنج کلاس است نیاز به این داریم که برای هر برجسب از نوع زیربخش، نوع اصلی آن زیر بخش یعنی یکی از چهار مورد حملات را جایگزین کنیم و این جایگزینی در پیاده‌سازی در نظر گرفته شده است.

- مرحله دو – آنالیز اکتشافی داده

برای این مرحله در پیاده‌سازی، فایل `exploratory_analysis.py` تعریف کردیم که نمودار ستونی توزیع نوع حملات برای داده آموزشی و تست، نمودار ستونی توزیع پروتکل‌ها برحسب نوع حملات برای داده آموزشی و تست به صورت عکس در پوشه `result` ذخیره می‌کند.

نمودار ستونی توزیع نوع حملات

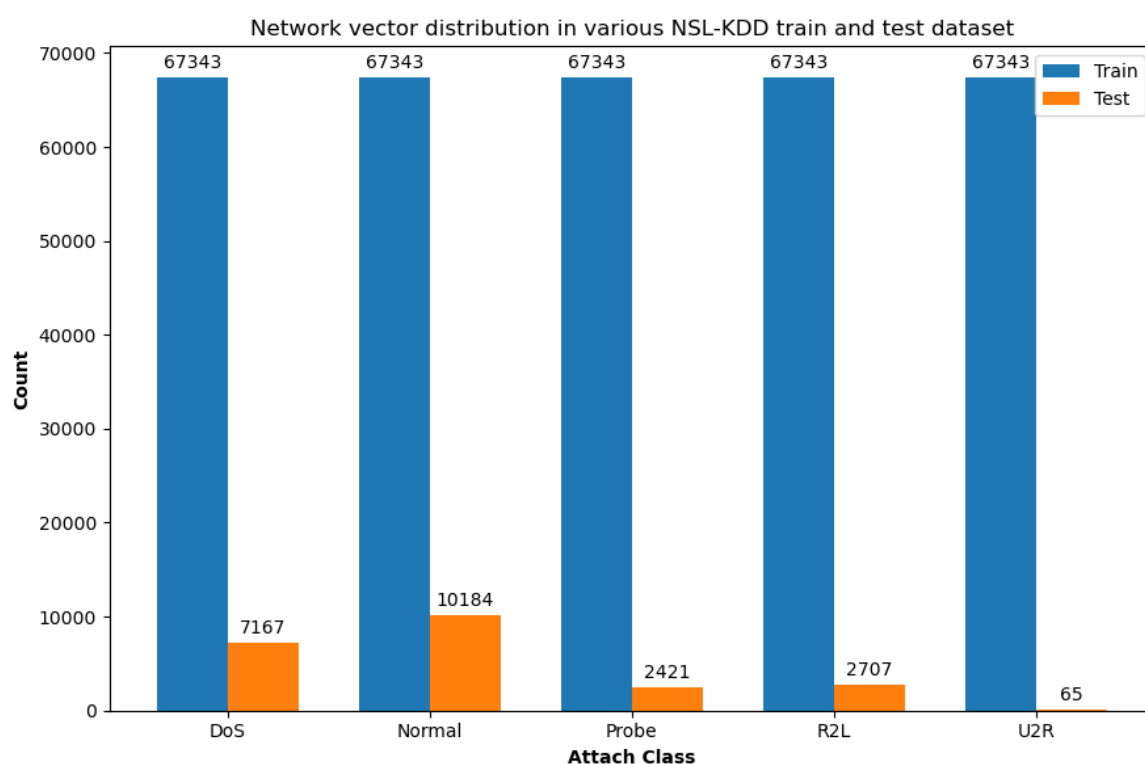


از نمودار بالا برداشت می‌شود که توزیع داده غیرمتعادل (imbalanced) است و این بالانس نبودن داده در میزان دقت سیستم تاثیر خواهد گذاشت که در قسمت خروجی نتایج را مشاهده خواهید کرد.

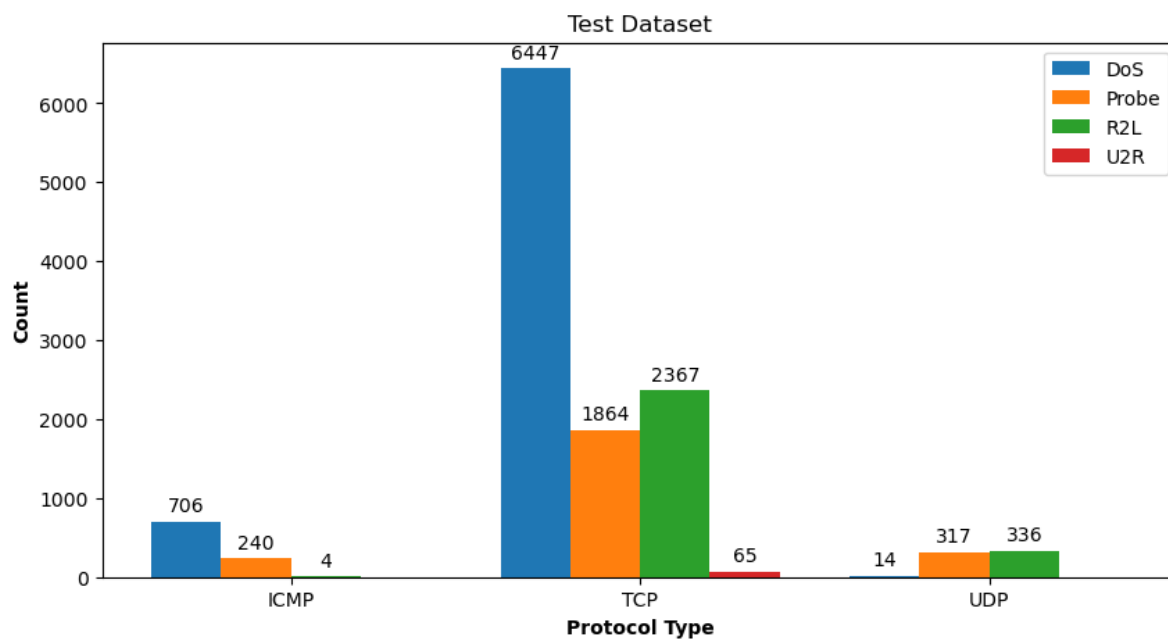
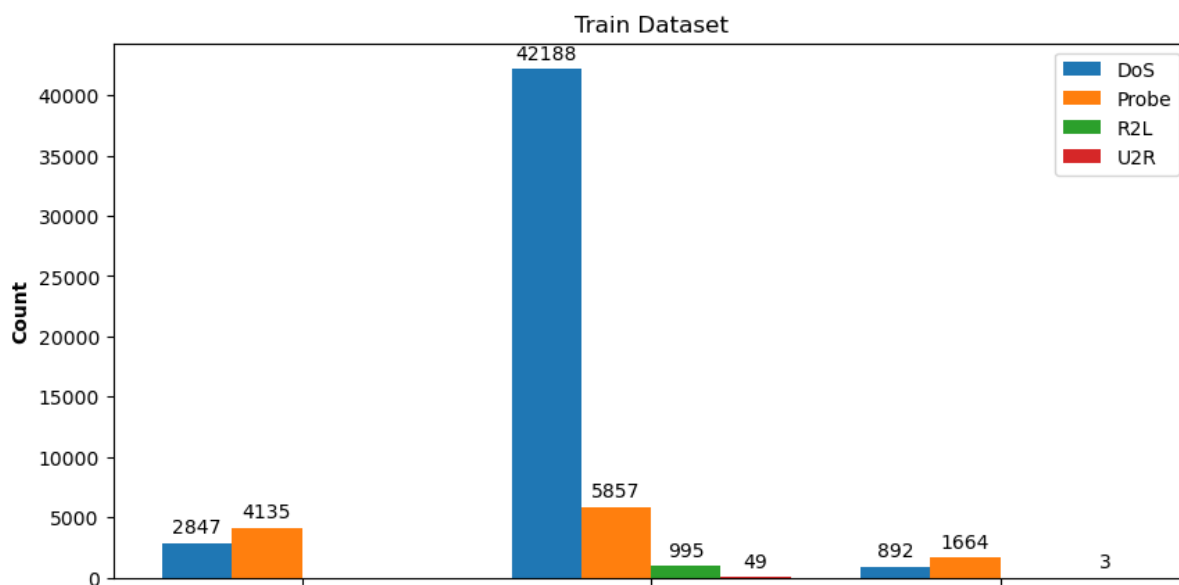
برای حل این مسئله با استفاده از روش SMOT حجم داده را افزایش دادیم تا توزیع داده بالانس شود. در پیاده‌سازی یک تابع با نام `resample_data` در فایل `util` تعریف کردم که ورودی داده و برجسب را می‌گیرد

و خروجی داده بالانس شده به همراه برچسب را برمی گردانند. داخل تابع فرکانس برچسبها را به دست می آوریم برچسبی به بیشترین فرکانس را دارد جدا می کنیم و برای بقیه برچسبها به تعداد بیشترین فرکانس از خود نمونه های متعلق به برچسب به صورت رندم نمونه تکراری اضافه می کنیم تا فرکانس آن برابر بیشترین فرکانس شود و داده بالانس شده را برمی گردانیم. این عملیات فقط برای داده آموزشی به دلیل آموزش مدل انجام می شود. نمودار ستونی داده بالانس شده را به صورت زیر است.

نمودار ستونی توزیع نوع حملات(بالانس شده)



## نمودار ستونی توزیع پروتکل‌ها برحسب نوع حملات



- مرحله سه - پیش پردازش:

برای این مرحله دو تابع `one_hot_encoding` ، `label_encoder` و `min_max_scaler` را در فایل `preprocessing.py` تعریف کردیم.

تابع `one_hot_encoding` به عنوان ورودی نمونه‌های داده آموزشی و تست را دریافت می‌کند و ویژگی‌هایی از داده که به صورت گروه‌بندی (categorical) هستند را به برداری از صفر و یک تبدیل می‌کند تا برای مرحله بعد برای انجام عملیات ریاضی مشکلی نداشته باشیم. در دیتا سه ویژگی `categorical` وجود دارد، ویژگی پروتکل (`protocol_type`) ، خدمات (`service`) و `Flag`. دیتا ورودی تابع 42 ویژگی است با اعمال `OneHotEncoder` به 123 ویژگی می‌رسد.

ویژگی پروتکل: 3 مقدار ، ویژگی خدمات : 79 مقدار ، ویژگی `flag`: 11 مقدار

ویژگی‌های بدون `categorical` در دیتا :  $42-3=39$

مجموع ویژگی‌ها بعد از اعمال `OneHotEncoder`:  $39+3+79+11=123$

تابع `label_encoder` به عنوان ورودی برچسب‌های داده آموزشی و تست را دریافت می‌کند و خروجی برچسب‌ها به صورت عدد را برمی‌گرداند. با استفاده از روش `labelEncoder` به هر برچسب یک مقدار عددی انتساب داده می‌شود.

['Normal' , 'Dos','Probe','R2L','U2R']

[0,1,2,3,4]

تابع `min_max_scaler` به عنوان ورودی نمونه را دریافت می‌کند و مقادیر هر ویژگی را از بازه

اصلی به بازه [0,1] تغییر می‌دهد تا میزان تاثیر هر ویژگی در مدل یکسان باشد.

- مرحله چهار - استخراج ویژگی:

برای این مرحله تابع `principal_component_analysis` را در فایل `feature_extraction.py` تعریف کردیم. تابع ورودی داده آموزشی و تست را دریافت می‌کند و خروجی با استفاده از تابع `PCA` از کتابخانه `Sklearn` از داده به تعداد تعیین شده ویژگی استخراج می‌کند و داده با تعداد ویژگی جدید برمی‌گرداند.

- مرحله پنج – آموزش مدل:

در این مرحله تابع `train_model` در فایل `run.py` تعریف کردیم. این تابع ورودی داده آموزشی و برچسب‌های آن را می‌گیرد و خروجی مدل آموزش دیده را برمی‌گرداند. داخل تابع با توجه به مقاله سه مدل `J48` یا همان `Naïve Bayes` ، `SVM` ، `Decision Tree C4.5` برای آموزش در نظر گرفته شده است. قطعا هر مدل نیاز به تنظیم کردن یک سری از پارامترها را دارد برای مقداردهی پارامترها دو حالت (دستی – خودکار) در نظر گرفتیم. اگر مایل به تنظیم دستی پارامترها دارید در فایل `config.py` مدل مورد نظر را مقدار `True` می‌دهید و پارامترهای تعریف شده را در همان فایل مقداردهی کنید و اگر مایل به تنظیم خودکار پارامترها هستید گزینه `gridsearch` را در فایل `config` مقدار `True` می‌دهید و برای مدل تابع `best_params_gridsearchcv` فراخوانی می‌شود و بعد از محاسبه بهترین پارامتر برگشت داده می‌شود و مدل با پارامتر مورد نظر آموزش دیده می‌شود و به عنوان خروجی تابع برگشت داده می‌شود.

- مرحله شش – ارزیابی مدل:

برای این مرحله تابع `evaluate` در فایل `run.py` تعریف کردیم. تابع ورودی داده آموزشی تست، برچسب آن، مدل آموزش دیده ، عنوان ماتریس درهم ریختگی عکس ، نام برچسب‌ها و مسیر ذخیره عکس ماتریس درهم ریختگی را دریافت می‌کند. داخل تابع به وسیله مدل برچسب داده تست را پیش‌بینی می‌کند. میزان دقت مدل را برای برچسب داده تست و برچسب پیش‌بینی شده محاسبه و نمایش می‌دهد و ماتریس درهم ریختگی را ایجاد و در مسیر مشخص شده ذخیره می‌کند.

برای راحتی کار فایل `config.py` تعریف کردیم که نیازی به تغییر فایل‌های دیگر نباشد.

متغیر `path_train_NSL_KDD` مسیر داده آموزشی

متغیر `path_test_NSL_KDD` مسیر داده تست

متغیر `balance_data` برای اینکه داده آموزشی بالانس شود یا نه اگر مایل به بالانس کردن داده هستید مقدار

`True` را به آن انتساب دهید

متغیر `normalize` برای تغییر بازه اعداد داده به بازه  $[0,1]$  است.

متغیر `max_feature` تعداد ویژگی‌هایی که باید به وسیله `pca` استخراج شود (طبق مقاله 6 در نظر گرفته شده

است)

متغیر `model_J48` ، `model_SVM` ، `model_NaiveBayes` برای تعیین مدل برای آموزش در نظر گرفته شده است

برای اجرا فقط به متغیر مورد نظر مقدار `True` انتساب دهید.

متغیر `gridsearch` برای تنظیم خودکار پارامترهای مدل در نظر گرفته شده است.

متغیر `max_depth` تنظیم دستی پارامتر مدل `J48` است.

متغیرهای `gamma` ، `C` ، `kernel` برای تنظیم دستی پارامتر `SVM` است.

## خروجی:

مدل J48 – داده غیرمتعادل (imbalance)

```

C:\> Command Prompt

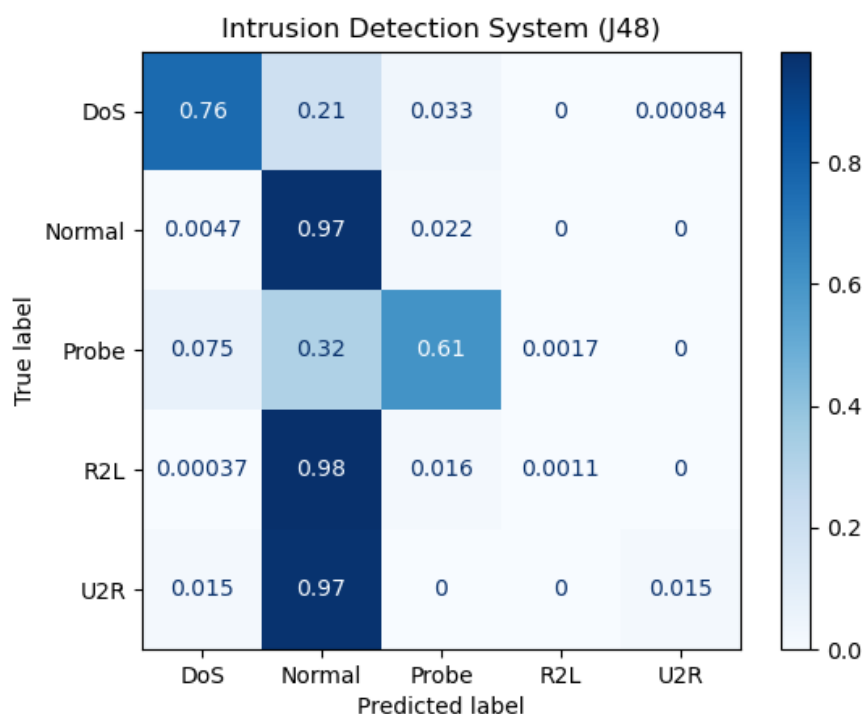
C:\Users\z.poorghorban\Atom\Intrusion Detection System>python run.py
Reading data...
Preprocessing data...
Extract features by PCA method...
Training J48 (Decision Tree C4.5)...
Evaluating the results...
accuracy:0.75

      precision    recall  f1-score   support

   DoS           0.96       0.76       0.85        7167
  Normal          0.67       0.97       0.79       10184
   Probe          0.74       0.61       0.67        2421
    R2L           0.43       0.00       0.00       2707
    U2R           0.14       0.02       0.03         65

 accuracy          0.75          0.75          0.75       22544
 macro avg          0.59          0.47          0.47       22544
weighted avg          0.74          0.75          0.70       22544

C:\Users\z.poorghorban\Atom\Intrusion Detection System>_
  
```





```

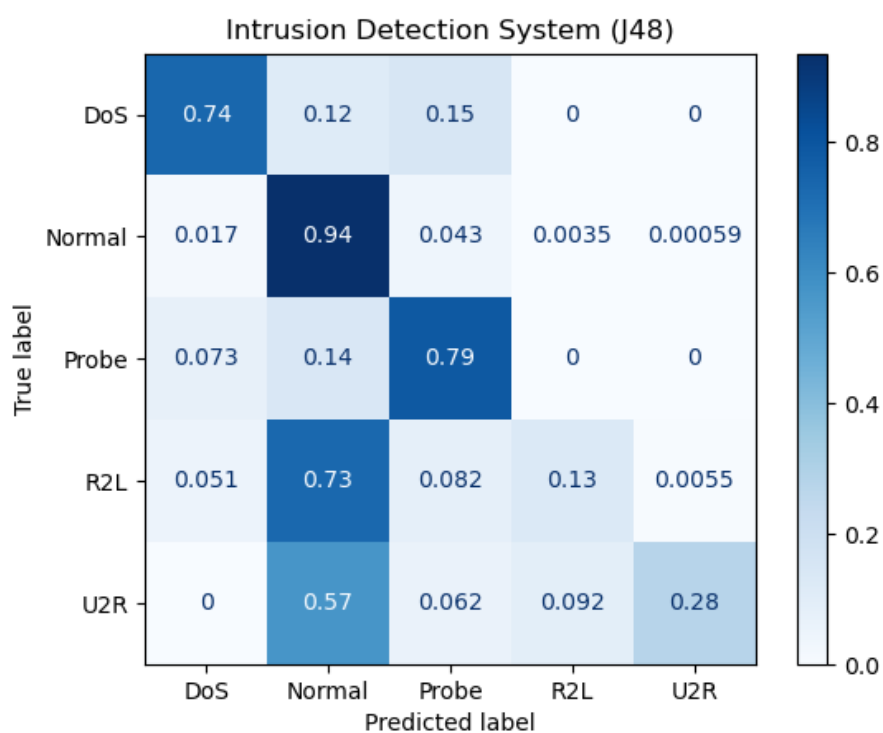
C:\>
(venv) C:\Users\z.poorghorban\Atom\Intrusion Detection System>python run.py
Reading data...
Balancing train dataset...
Preprocessing data...
Extract features by PCA method...
Training J48 (Decision Tree C4.5)...
Evaluating the results...
accuracy:0.76
precision    recall  f1-score   support

   DoS        0.91     0.74     0.82     7167
  Normal        0.75     0.94     0.83    10184
   Probe        0.52     0.79     0.63     2421
    R2L        0.89     0.13     0.23     2707
    U2R        0.46     0.28     0.35         65

 accuracy          0.76    22544
 macro avg         0.71    22544
weighted avg         0.79    22544

(venv) C:\Users\z.poorghorban\Atom\Intrusion Detection System>

```



```

Command Prompt

C:\Users\z.poorghorban\Atom\Intrusion Detection System>python run.py
Reading data...
Preprocessing data...
Extract features by PCA method...
Training SVM...
Evaluating the results...
accuracy:0.71

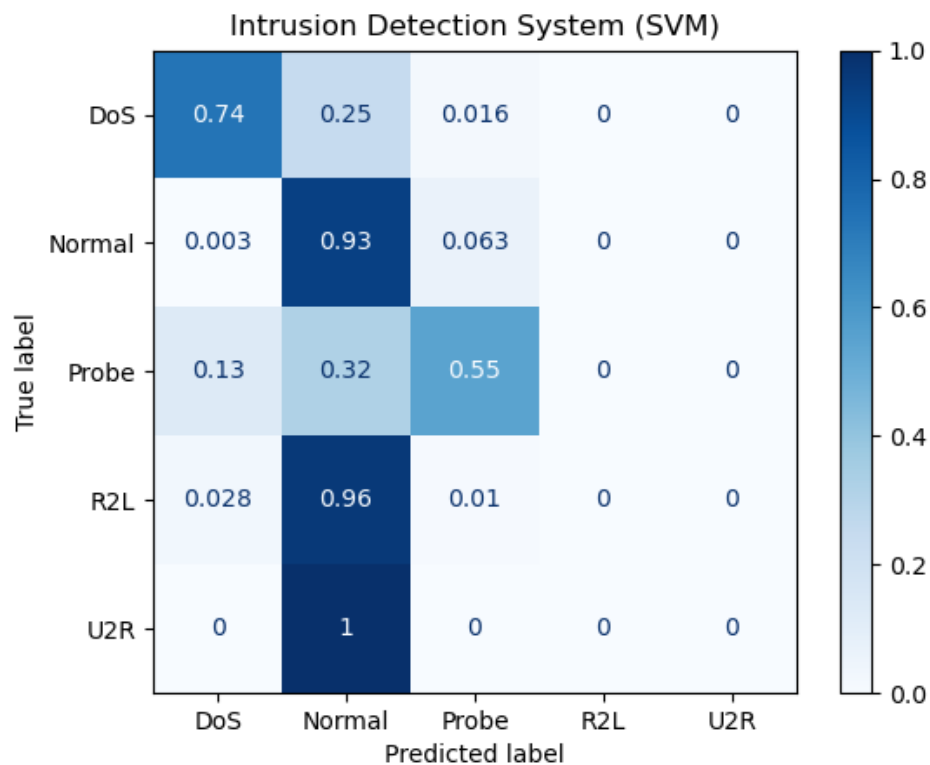
      precision    recall  f1-score   support

   DoS           0.93      0.74      0.82       7167
  Normal          0.65      0.93      0.76      10184
   Probe          0.63      0.55      0.58       2421
    R2L           0.00      0.00      0.00       2707
    U2R           0.00      0.00      0.00         65

 accuracy          0.71          0.71          0.67      22544
 macro avg          0.44      0.44      0.43      22544
weighted avg          0.65      0.71      0.67      22544

C:\Users\z.poorghorban\Atom\Intrusion Detection System>_

```



```

Command Prompt

(venv) C:\Users\z.poorghorban\Atom\Intrusion Detection System>python run.py
Reading data...
Balancing train dataset...
Preprocessing data...
Extract features by PCA method...
Training SVM...
Evaluating the results...
accuracy:0.77

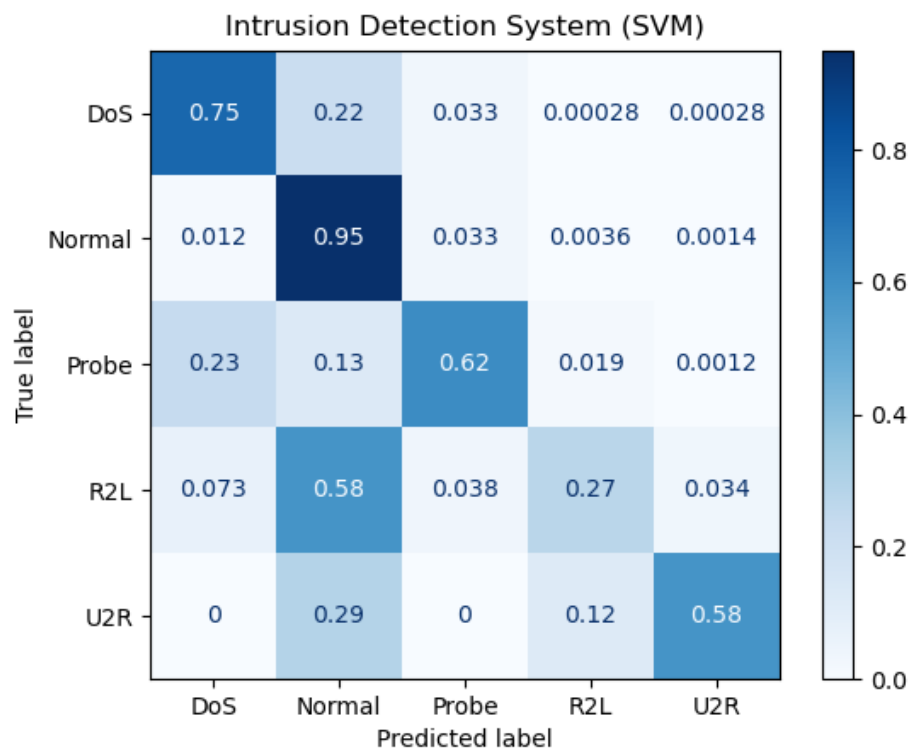
              precision    recall  f1-score   support

   DoS              0.86        0.75        0.80        7167
  Normal              0.74        0.95        0.83       10184
   Probe              0.69        0.62        0.65        2421
   R2L              0.89        0.27        0.42       2707
   U2R              0.26        0.58        0.36         65

 accuracy              0.77        0.77        0.75       22544
  macro avg              0.69        0.63        0.61       22544
 weighted avg              0.79        0.77        0.75       22544

(venv) C:\Users\z.poorghorban\Atom\Intrusion Detection System>_

```



```

C:\Users\z.poorghorban\Atom\Intrusion Detection System>python run.py
Reading data...
Preprocessing data...
Extract features by PCA method...
Training Naive Bayes...
Evaluating the results...
accuracy:0.67

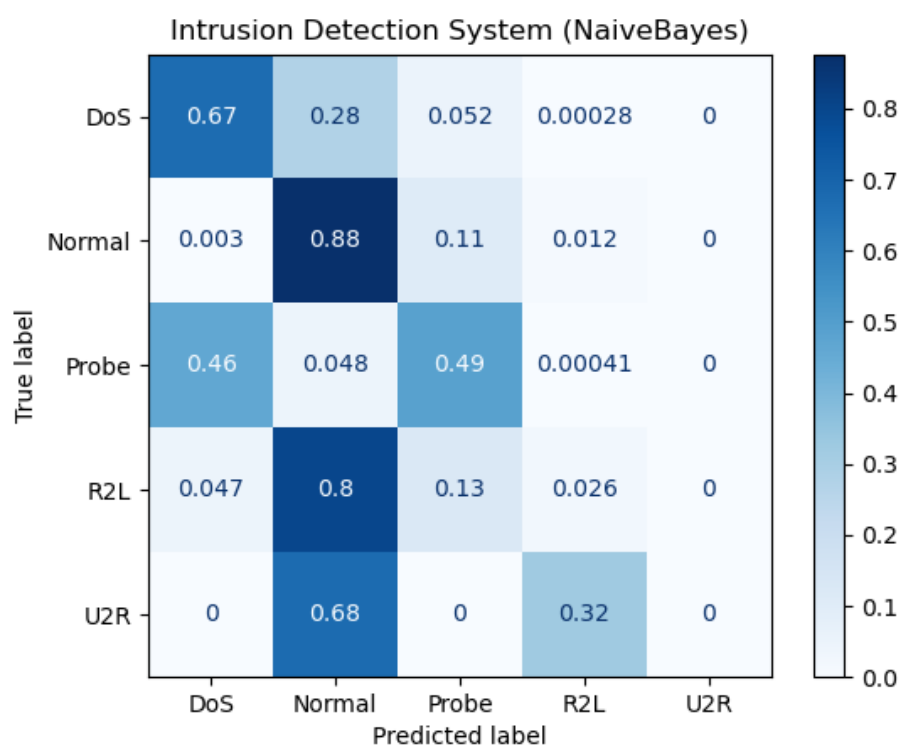
```

	precision	recall	f1-score	support
DoS	0.79	0.67	0.73	7167
Normal	0.68	0.88	0.76	10184
Probe	0.39	0.49	0.44	2421
R2L	0.33	0.03	0.05	2707
U2R	0.00	0.00	0.00	65
accuracy			0.67	22544
macro avg	0.44	0.41	0.39	22544
weighted avg	0.64	0.67	0.63	22544

```

C:\Users\z.poorghorban\Atom\Intrusion Detection System>

```



```

Command Prompt

<venv> C:\Users\z.poorghorban\Atom\Intrusion Detection System>python run.py
Reading data...
Balancing train dataset...
Preprocessing data...
Extract features by PCA method...
Training Naive Bayes...
Evaluating the results...
accuracy:0.72

      precision    recall  f1-score   support

   DoS           0.80      0.62      0.70       7167
  Normal          0.73      0.91      0.81      10184
   Probe          0.57      0.57      0.57       2421
    R2L           0.69      0.40      0.51      2707
    U2R           0.13      0.54      0.22         65

 accuracy          0.72      0.72      0.71      22544
  macro avg          0.58      0.61      0.56      22544
 weighted avg          0.73      0.72      0.71      22544

<venv> C:\Users\z.poorghorban\Atom\Intrusion Detection System>_

```

