

تمرین

اینترنت اشیا

استاد

محمد زارع

به کوشش

مهدی جوانشیر

تاریخ

آذر 1403



عنوان تمرین نرم افزاری: تمرین شناسایی حمله با RNN
برای شناسایی حملات سایبری و تحلیل رفتارهای نامتعارف از طریق شبکه‌های عصبی
بازگشتی (RNN) به تمرین زیر پردازید



هدف

شناسایی رفتارهای غیرطبیعی در ترافیک شبکه (شامل حملات مثل DDOS ، اسکن پورت و غیره) با استفاده از RNN.

مراحل انجام تمرین

1. جمع آوری داده ها

از یک دیتاست مناسب برای ترافیک شبکه استفاده کنید. دیتاست‌هایی مانند Cup KDD 99 یا CICIDS2017 مناسب هستند.

➤ این دیتاست‌ها باید شامل ویژگی‌هایی مانند زمان، منبع، مقصد، نوع پروتکل و نشانگر حمله یا عدم حمله باشند.

2. پیش‌پردازش داده‌ها

- داده‌ها را پاک‌سازی کنید و ویژگی‌های لازم را انتخاب کنید.
- داده‌ها را به شکل زمان‌بندی آماده کنید. به عنوان مثال، هر 100 نمونه از ترافیک را به عنوان یک ورودی برای شبکه RNN در نظر بگیرید.
- داده‌ها را به دو بخش آموزش و تست تقسیم کنید.

3. طراحی مدل RNN

- از یک معماری ساده RNN یا LSTM (Long Short- Term Memory) استفاده کنید.
- ورودی به مدل باید شامل ویژگی‌های ترافیک به همراه برچسب حمله یا عدم حمله باشد.
- می‌توانید از کتابخانه‌های TensorFlow یا PyTorch برای پیاده‌سازی استفاده کنید.

4. آموزش مدل

- مدل را با داده‌های آموزشی آموزش دهید.
- از تکنیک‌های تنظیم هایپرپارامتر و بهینه‌سازی استفاده کنید تا دقت مدل را افزایش دهید.

5. ارزیابی مدل

- مدل را با داده‌های تست ارزیابی کنید.
- از معیارهای مختلفی مانند دقت، یادآوری، F1-score و ROC-AUC برای اندازه‌گیری عملکرد مدل استفاده کنید.

6. تحلیل نتایج

- نتایج را تحلیل کنید و بررسی کنید که کجا مدل به درستی یا نادرستی پیش‌بینی کرده است.
- می‌توانید مهمترین ویژگی‌ها را شناسایی کنید که به شناسایی حمله کمک کرده‌اند.

نکات اضافی

- برای بهبود دقت، می‌توانید از تکنیک‌های متنوعی مانند Dropout یا Regularization استفاده کنید.
- همچنین می‌توانید با استفاده از شبکه‌های عصبی پیچشی (CNN) به عنوان ورودی به RNN، ویژگی‌های فضایی را نیز استخراج کنید.

- 1 وارد کردن کتابخانه‌های مورد نیاز #
- 2 `import pandas as pd` # برای پردازش داده‌ها به صورت ساختاریافته
- 3 `import numpy as np` # برای عملیات عددی و کار با آرایه‌ها
- 4 `import tensorflow as tf` # برای طراحی و اجرای مدل‌های یادگیری عمیق
- 5 `from sklearn.preprocessing import StandardScaler` # برای نرمال‌سازی داده‌ها
- 6 `from sklearn.model_selection import train_test_split` # برای تقسیم داده‌ها به مجموعه‌های آموزش و آزمایش
- 7 `from sklearn.metrics import classification_report, accuracy_score` # برای ارزیابی مدل
- 8 `import matplotlib.pyplot as plt` # برای رسم نمودارها
- 9 `import dask.dataframe as dd` # برای پردازش موازی داده‌های بزرگ

کتاب خانه ها و ماژول ها

(pd) pandas

- برای مدیریت داده‌ها در قالب DataFrame که دسترسی و عملیات روی داده‌ها را ساده‌تر می‌کند.
- در این کد از pandas برای انجام تغییرات روی ستون‌های داده (مانند تبدیل برچسب‌ها و ویژگی‌ها به مقادیر عددی) استفاده شده است.

(np) numpy

- برای انجام عملیات عددی کارآمد، مانند تغییر شکل داده‌ها و بررسی وجود مقادیر NaN.
- آرایه‌های numpy به دلیل کارایی بالا در پردازش داده‌های عددی برای آماده‌سازی ورودی مدل‌ها استفاده می‌شوند.

(tf) tensorflow

- ساخت مدل RNN با لایه LSTM برای شناسایی الگوهای زمانی در داده‌ها.
- ابزارهای قدرتمندی برای طراحی لایه‌ها، کامپایل مدل، و آموزش آن فراهم می‌کند.

- `sklearn.preprocessing.StandardScaler`: برای نرمال‌سازی و استانداردسازی داده‌ها به کار می‌رود.
- `sklearn.model_selection.train_test_split`: تقسیم داده‌ها به مجموعه‌های آموزش و آزمایش.
- `sklearn.metrics.classification_report` و `accuracy_score`: برای ارائه معیارهای دقت، یادآوری (recall)، و F1-score. برای محاسبه درصد پیش‌بینی‌های صحیح.

(plt) matplotlib.pyplot

- رسم نمودار برای تحلیل و تجسم نتایج.
- برای نمایش گرافیکی تغییرات دقت مدل در طول `epoch`ها.

dask.dataframe

- داده‌های حجیم می‌توانند پردازش را کند کنند؛ `dask` امکان بارگذاری داده‌ها به صورت موازی و مدیریت منابع بهینه‌تر را فراهم می‌کند.

```
11 # مرحله 1: جمع‌آوری داده‌ها
12 print("مرحله 1: جمع‌آوری داده‌ها")
13
14 # مسیر فایل داده‌ها (فایل ورودی شامل داده‌های شبکه برای شناسایی نفوذ)
15 file_path = 'C:/Users/Mahdi/Desktop/KDDCup99.txt'
16
17 # برای پردازش موازی (پردازش سریع‌تر داده‌های بزرگ) Dask بارگذاری داده‌ها با استفاده از
18 data = dd.read_csv(file_path, delimiter=',', na_values='?', assume_missing=True)
19
20 # استخراج 10000 ردیف داده‌ها برای کاهش حجم پردازش
21 data = data.head(10000)
22
23 # نمایش پیام موفقیت در جمع‌آوری داده‌ها
24 print("داده‌ها با موفقیت جمع‌آوری شدند")
25
26 # نمایش چند ردیف اول داده‌ها برای بررسی
27 print(data.head())
```

جمع آوری داده ها

در این کد، داده ها از یک فایل متنی با نام `KDDCup99.txt` جمع آوری شده اند. فرآیند جمع آوری داده ها به این صورت انجام شده است:

بارگذاری داده ها با استفاده از Dask

- `dask` داده ها را به صورت موازی بارگذاری و پردازش می کند، بنابراین برای فایل های بزرگ کارایی بیشتری نسبت به `pandas` دارد.
- این قابلیت کمک می کند که با فایل های بسیار بزرگ بدون بارگذاری کامل در حافظه، کار کنیم.

عملکرد

- فایل ورودی با استفاده از `dask.dataframe.read_csv` بارگذاری می شود.
- `'=' delimiter`: مشخص می کند که داده ها با کاما جدا شده اند (فرمت CSV).
- `'?'=na_values`: مشخص می کند که مقادیر علامت `?` به عنوان مقادیر ناشناخته (`NaN`) در نظر گرفته شوند.
- `assume_missing=True`: برای اطمینان از اینکه ستون ها با داده های ناقص به درستی پردازش شوند.

کاهش حجم داده ها

- از کل مجموعه داده، فقط 10,000 ردیف انتخاب می شود
- کاهش حجم داده ها به منظور ساده تر و سریع تر کردن پردازش اولیه (به خصوص در مراحل تست و دیباگ) ..


```

28 مرحله 2: پیش‌پردازش داده‌ها #
29 print("\nمرحله 2: پیش‌پردازش داده‌ها")
30
31 به مقادیر عددی (labels) تعریف نقشه‌برداری برای تبدیل برچسب‌های متنی #
32 label_mapping = {
33     'normal': 0, # رفتار عادی
34     'buffer_overflow': 1, 'loadmodule': 1, 'perl': 1, 'neptune': 1, 'smurf': 1,
35     'guess_passwd': 1, 'pod': 1, 'teardrop': 1, 'portsweep': 1, 'ipsweep': 1,
36     'land': 1, 'ftp_write': 1, 'back': 1, 'imap': 1, 'satan': 1, 'phf': 1,
37     'nmap': 1, 'multihop': 1, 'warezmaster': 1, 'warezclient': 1, 'spy': 1,
38     'rootkit': 1 # رفتارهای غیرعادی (حملات)
39 }
40
41 اعمال نقشه‌برداری به ستون برچسب‌ها #
42 data['label'] = data['label'].map(label_mapping)
43
44 دارد Null حذف ردیف‌هایی که برچسب آنها مقدار #
45 data = data.dropna(subset=['label'])
46
47 تبدیل ویژگی‌های متنی به کدهای عددی (برای استفاده در مدل یادگیری ماشین) #
48 data['protocol_type'] = data['protocol_type'].astype('category').cat.codes
49 data['service'] = data['service'].astype('category').cat.codes
50 data['flag'] = data['flag'].astype('category').cat.codes
51
52 (y) و برچسب‌ها (X) جدا کردن ویژگی‌ها #
53 X = data.drop(columns=['label']) # حذف ستون برچسب‌ها برای استفاده به‌عنوان ویژگی
54 y = data['label'] # برچسب‌ها (هدف)
55
56 نرمال‌سازی داده‌های ویژگی‌ها برای بهبود عملکرد مدل #
57 scaler = StandardScaler()
58 X = scaler.fit_transform(X)

```

پیش پردازش داده ها

داده ها را برای استفاده در مدل یادگیری عمیق آماده می کند. هدف این مرحله اطمینان از این است که داده ها به فرم مناسب برای مدل های یادگیری ماشین و یادگیری عمیق تبدیل شوند.

تعریف نقشه برداری برای برجسب ها

- این نقشه برداری برجسب های متنی موجود در داده ها را به مقادیر عددی تبدیل می کند.
- مقدار 0 برای رفتارهای عادی (Normal)، مقدار 1 برای رفتارهای غیرعادی (حملات).
- با استفاده از label_mapping، مقادیر ستون label به مقادیر عددی (0 یا 1) تبدیل می شود.

تبدیل ویژگی های متنی به مقادیر عددی

- ستون های متنی (مانند protocol_type، service و flag) به کدهای عددی تبدیل می شوند.
- مدل ها فقط می توانند مقادیر عددی را پردازش کنند، بنابراین مقادیر متنی باید به کدهای عددی تبدیل شوند.

نرمال سازی داده ها با StandardScaler

- نرمال سازی داده ها باعث می شود که میانگین هر ویژگی برابر صفر و انحراف معیار برابر یک باشد.
- این کار به مدل کمک می کند تا بهتر یاد بگیرد، به خصوص وقتی که مقادیر ویژگی ها در مقیاس های مختلف باشند..

```

59 # مرحله 3: طراحی مدل
60 print("\n مرحله 3: طراحی مدل")
61 # (20%) تقسیم داده‌ها به مجموعه‌های آموزش (80%) و آزمایش
62 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
63
64 # در برجسب‌های آموزشی و آزمایشی NaN بررسی عدم وجود مقادیر
65 assert not np.any(np.isnan(y_train)), "y_train contains NaN values"
66 assert not np.any(np.isnan(y_test)), "y_test contains NaN values"
67
68 # RNN تغییر شکل داده‌ها به 3 بعد برای ورودی مدل
69 X_train = X_train.reshape((X_train.shape[0], 1, X_train.shape[1])) # تغییر شکل داده‌های آموزشی
70 X_test = X_test.reshape((X_test.shape[0], 1, X_test.shape[1])) # تغییر شکل داده‌های آزمایشی
71
72 # LSTM با استفاده از RNN تعریف مدل
73 model = tf.keras.Sequential([
74     tf.keras.layers.Input(shape=(X_train.shape[1], X_train.shape[2])), # لایه ورودی با شکل داده
75     tf.keras.layers.LSTM(64, return_sequences=False), # با 64 واحد LSTM لایه
76     tf.keras.layers.Dense(32, activation='relu'), # 32 نرون و فعال‌سازی Dense لایه
77     tf.keras.layers.Dense(1, activation='sigmoid') # لایه خروجی با فعال‌سازی سیگموئید (برای
78     طبقه‌بندی دودویی)
79 ])
80
81 # کامپایل مدل با تابع هزینه و بهینه‌ساز
82 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
83
84 # (Overfitting) برای جلوگیری از آموزش بیش از حد Early Stopping تعریف
85 early_stopping = tf.keras.callbacks.EarlyStopping(monitor='val_loss', patience=3,
86 restore_best_weights=True)
87
88 # epoch برای کاهش تدریجی نرخ یادگیری در هر Learning Rate Scheduler تعریف
89 lr_scheduler = tf.keras.callbacks.LearningRateScheduler(
    lambda epoch: 1e-3 * 0.7 ** epoch, verbose=True)

```

طراحی مدل

این تکه کد به مرحله طراحی مدل مربوط است که در آن یک مدل یادگیری عمیق بر پایه شبکه عصبی بازگشتی (RNN) با استفاده از لایه‌های LSTM طراحی و آماده استفاده می‌شود.

بررسی و تغییر شکل داده‌ها برای استفاده در RNN

داده‌ها به شکل سه‌بعدی تبدیل می‌شوند که ساختار موردنیاز برای ورودی به مدل‌های RNN است:

- بعد اول: تعداد نمونه‌ها (Samples).
- بعد دوم: تعداد گام‌های زمانی (Time Steps). در اینجا برابر 1 است، زیرا هر نمونه یک بردار ویژگی است.
- بعد سوم: تعداد ویژگی‌ها (Features) که در هر گام زمانی ورودی به مدل هستند.

تعریف مدل شبکه عصبی بازگشتی (RNN)

- تعریف شکل ورودی مدل به صورت (گام‌های زمانی، ویژگی‌ها). اینجا گام‌های زمانی برابر 1 و تعداد ویژگی‌ها برابر تعداد ستون‌های X_{train} است.
- لایه LSTM: برای پردازش داده‌های ترتیبی یا سری‌های زمانی بسیار مناسب است. یک لایه Long Short-Term Memory (LSTM) با 64 واحد تعریف می‌شود.
- تابع فعال‌سازی (ReLU): غیرخطی بودن را به مدل اضافه می‌کند و به آن کمک می‌کند تا روابط پیچیده‌تری را یاد بگیرد.
- لایه خروجی با فعال‌سازی Sigmoid: مقادیر خروجی را به بازه $[0, 1]$ محدود می‌کند که برای طبقه‌بندی دودویی مناسب است.

کامپایل مدل

کامپایل مدل: مشخص می‌کند که مدل چگونه یاد می‌گیرد و ارزیابی می‌شود.

➤ بهینه‌ساز (Optimizer)

- لایه خروجی با فعال‌سازی Sigmoid: مقادیر خروجی را به بازه $[0, 1]$ محدود می‌کند که برای طبقه‌بندی دودویی مناسب است.

- تابع هزینه (Loss Function): `binary_crossentropy` برای مسائل طبقه‌بندی دودویی استفاده می‌شود.
- معیار ارزیابی (Metrics): `accuracy` برای محاسبه درصد پیش‌بینی‌های صحیح.

تعریف Learning Rate Scheduler

عملکرد: نرخ یادگیری مدل را به صورت پویا در هر دوره کاهش می‌دهد.

➤ تابع کاهش نرخ یادگیری: کاهش تدریجی نرخ یادگیری به مدل کمک می‌کند تا به یک مینیمم پایدارتر برسد.

```

90 # مرحله 4: آموزش مدل
91 print("\nمرحله 4: آموزش مدل")
92
93 # آموزش مدل با داده‌های آموزشی و اعتبارسنجی بر اساس داده‌های آزمایشی
94 history = model.fit(X_train, y_train, epochs=20, batch_size=128, validation_data=(X_test,
95 y_test),
96
97 callbacks=[early_stopping, lr_scheduler])
98
99 # مرحله 5: ارزیابی مدل
100 print("\nمرحله 5: ارزیابی مدل")
101
102 # پیش‌بینی بر روی داده‌های آزمایشی
103 y_pred = (model.predict(X_test) > 0.5).astype("int32") # تبدیل احتمال به مقادیر 0 یا 1
104
105 # (F1-Score شامل دقت، یادآوری، و) تولید و نمایش گزارش طبقه‌بندی
106 print(classification_report(y_test, y_pred))
107
108 # محاسبه و نمایش دقت مدل
109 accuracy = accuracy_score(y_test, y_pred)
110 print(f"دقت مدل: {accuracy * 100:.2f}%")

```

آموزش و ارزیابی مدل

داده‌های آموزشی و اعتبارسنجی

- X_{train} و y_{train} : داده‌های آموزشی که مدل بر اساس آنها یادگیری انجام می‌دهد.
- $validation_data=(X_test, y_test)$: داده‌های آزمایشی که برای بررسی عملکرد مدل در طول آموزش استفاده می‌شوند. این داده‌ها در حین آموزش برای تنظیم مدل استفاده شده، اما در یادگیری مستقیم دخالت ندارند.

تعداد دوره‌ها: (Epochs)

➤ epochs=20: مدل در حداکثر 20 دوره (دوره‌های کامل بر داده‌های آموزشی) آموزش داده می‌شود.

اندازه دسته (Batch Size):

➤ batch_size=128: داده‌های آموزشی در گروه‌های 128 تایی پردازش می‌شوند. این به کاهش مصرف حافظه و تسریع آموزش کمک می‌کند.

پیش‌بینی بر روی داده‌های آزمایشی

- model.predict(X_test): مدل، احتمال پیش‌بینی شده برای هر نمونه در داده‌های آزمایشی را برمی‌گرداند (مقدار بین 0 و 1).
- ($0.5 <$): اگر احتمال پیش‌بینی شده بزرگتر از 0.5 باشد، مقدار 1 (حمله) و در غیر این صورت مقدار 0 (عادی) تعیین می‌شود.
- astype("int32"): مقادیر پیش‌بینی شده به نوع عدد صحیح تبدیل می‌شوند.

گزارش طبقه‌بندی

یک گزارش کامل از معیارهای عملکرد مدل تولید و نمایش داده می‌شود:

- Precision (دقت): نسبت پیش‌بینی‌های صحیح از میان تمام پیش‌بینی‌های مثبت.
- Recall (بازیابی): نسبت نمونه‌های مثبت که به درستی شناسایی شده‌اند.
- F1-Score: میانگین موزون دقت و بازیابی.
- Support: تعداد نمونه‌های هر کلاس.

محاسبه دقت مدل

➤ accuracy_score(y_test, y_pred): درصد نمونه‌های آزمایشی که به درستی پیش‌بینی شده‌اند.

```

109 مرحله 6: تحلیل نتایج #
110 print("\nمرحله 6: تحلیل نتایج")
111
112 ها برای آموزش و اعتبارسنجی epoch رسم نمودار دقت در طول #
113 plt.figure(figsize=(10, 6)) # تنظیم ابعاد نمودار
114 plt.plot(history.history['accuracy'], label='Training Accuracy') # نمودار دقت آموزشی
115 plt.plot(history.history['val_accuracy'], label='Validation Accuracy') # نمودار دقت اعتبارسنجی
116 plt.title('Model Accuracy Over Epochs') # عنوان نمودار
117 plt.xlabel('Epochs') # برچسب محور افقی
118 plt.ylabel('Accuracy') # برچسب محور عمودی
119 plt.legend() # افزودن راهنما به نمودار
120 plt.show() # نمایش نمودار

```

نتیجه گیری

جمع آوری داده‌ها

- داده‌ها از فایل KDDCup99.txt با استفاده از Dask بارگذاری شدند.
- برای کاهش بار پردازشی، تنها 10,000 نمونه داده مورد استفاده قرار گرفت.
- این داده‌ها شامل ویژگی‌های مختلف ترافیک شبکه و نوع رفتار (عادی یا غیرعادی) بودند.

پیش‌پردازش داده‌ها

- تبدیل برچسب‌ها: انواع حملات به مقدار عددی 1 (غیرعادی) و رفتارهای عادی به 0 تبدیل شدند.
- پردازش ویژگی‌ها:
- داده‌های متنی (مثل پروتکل‌ها) به مقادیر عددی تبدیل شدند.
- داده‌ها با استفاده از StandardScaler نرمال‌سازی شدند.
- تقسیم داده‌ها: داده‌ها به مجموعه‌های آموزشی (80٪) و آزمایشی (20٪) تقسیم شدند.
- تغییر شکل: داده‌ها به قالب سه‌بعدی تبدیل شدند تا با معماری شبکه‌های عصبی بازگشتی (RNN) سازگار باشند.

طراحی مدل

- یک مدل RNN با لایه‌های زیر طراحی شد:
- لایه LSTM: برای یادگیری توالی‌های زمانی داده‌ها.
- لایه Dense: برای ترکیب ویژگی‌های یادگرفته شده.
- لایه خروجی: با فعال‌سازی سیگموئید برای پیش‌بینی طبقه‌بندی دودویی.
- از Early Stopping و Learning Rate Scheduler برای بهبود فرآیند آموزش و جلوگیری از Overfitting استفاده شد.

آموزش مدل

- مدل با 20 دوره (epoch) و حجم دسته‌ای (batch size) برابر 128 آموزش داده شد.
- داده‌های آزمایشی برای اعتبارسنجی استفاده شدند.
- یادگیری نرخ (Learning Rate) به تدریج کاهش یافت.

ارزیابی مدل

- گزارش طبقه‌بندی (Classification Report):
- مدل عملکرد بسیار خوبی داشته است:
- Precision (دقت): درصد پیش‌بینی‌های درست در مقایسه با تمام پیش‌بینی‌های مثبت.
- برای هر دو کلاس (عادی و غیرعادی) برابر 1.00 (100٪) است.
- Recall (بازخوانی): درصد نمونه‌های مثبت که به درستی شناسایی شده‌اند.
- برای کلاس عادی: 1.00.
- برای کلاس غیرعادی: 0.99.
- F1-Score: میانگین هارمونیک دقت و بازخوانی، که برای هر دو کلاس نزدیک به 1.00 است.
- Overall Accuracy (دقت کل): دقت نهایی مدل 99.85٪ است که نشان‌دهنده عملکرد عالی مدل است.

مدل طراحی شده به طور موفقیت آمیزی حملات شبکه را شناسایی می‌کند.

دقت 99.85٪ بیانگر این است که مدل در یادگیری رفتارهای ترافیک شبکه و شناسایی نفوذ بسیار قوی عمل کرده است. مراحل مختلف از جمع‌آوری داده‌ها تا پیش‌پردازش، طراحی، آموزش، و ارزیابی بدون خطا اجرا شده‌اند. این سیستم می‌تواند به عنوان بخشی از یک راه‌حل امنیت سایبری در شناسایی نفوذها و رفتارهای غیرعادی در شبکه استفاده شود.