

# پروژه آنالیز بلادرنگ تراکنش‌های بانکی و کشف تقلب

آرمین افضلی  
محمد صالح پژند

## فهرست مطالب

|     |   |   |
|-----|---|---|
| ۱   | مقدمه و شرح پروژه                             | ۲ |
| ۱.۱ | هدف پروژه                                     | ۲ |
| ۱.۲ | دیتاست  | ۲ |
| ۲   | معماری سیستم                                  | ۲ |
| ۲.۱ | جريان داده                                    | ۲ |
| ۲.۲ | اجزای سیستم                                   | ۲ |
| ۳   | گام اول: تولید جریان داده                     | ۳ |
| ۳.۱ | پیاده‌سازی <i>Kafka Producer</i>              | ۳ |
| ۳.۲ | منطق تخصیص کاربر                              | ۳ |
| ۳.۳ | کد پیاده‌سازی                                 | ۳ |
| ۳.۴ | خروجی تولیدکننده (نتایج واقعی)                | ۳ |
| ۴   | گام دوم: پردازش و تحلیل                       | ۴ |
| ۴.۱ | دسته‌بندی تراکنش‌ها                           | ۴ |
| ۴.۲ | تشخیص کاربران لیست سیاه ( <i>JOIN</i> )       | ۴ |
| ۴.۳ | تشخیص بات با پنجره زمانی ( <i>Windowing</i> ) | ۴ |
| ۵   | گام سوم: نمایش نتایج                          | ۵ |
| ۵.۱ | خروجی تشخیص بات (نتایج واقعی)                 | ۵ |
| ۵.۲ | آمار کلی (نتایج واقعی)                        | ۵ |
| ۵.۳ | آمار نهایی                                    | ۵ |
| ۶   | پیاده‌سازی <i>Docker</i>                      | ۶ |
| ۶.۱ | سروریس‌های کانتینری                           | ۶ |
| ۶.۲ | دستورات اجرا                                  | ۶ |

# ۱. مقدمه و شرح پروژه

## ۱.۱. هدف پروژه

هدف این پروژه پیاده‌سازی یک پایپلاین کلان‌داده (*Big Data Pipeline*) است که داده‌های واقعی تراکنش‌های بانکی را به صورت بلاذرنگ پردازش کرده و ناهنجاری‌های زیر را شناسایی می‌کند:

- تراکنش‌های مشکوک از کاربران شناخته شده متقلب (*User\_Hacker*)
- حملات بات (*Bot Attack*) با استفاده از تحلیل پنجره زمانی (*Macro/Micro/Normal*)
- دسته‌بندی تراکنش‌ها بر اساس مبلغ (*BLOCKED*)
- کاربران لیست سیاه با وضعیت

## ۲. دیتاست

این پروژه از دیتاست *Credit Card Fraud Detection* از سایت *Kaggle* استفاده می‌کند:

## ۲. معماهی سیستم

### ۲.۱. جریان داده

داده‌ها در این پایپلاین به صورت زیر جریان می‌یابند:

*creditcard.csv* → *Kafka Producer* → *Kafka Topic* → *Spark Streaming* → *Output*

### ۲.۲. اجزای سیستم

| توضیحات  | کامپوننت               |
|--|------------------------|
| سرویس هماهنگ‌کننده کافکا   | <i>Zookeeper</i>       |
| صف پیام توزیع شده با تاپیک <i>bank_transactions</i>                      | <i>Apache Kafka</i>    |
| اسکریپت پایتون برای خواندن <i>CSV</i> و ارسال با شبیه‌سازی تقلب/بات      | <i>Kafka Producer</i>  |
| هماهنگ‌کننده کلاستر اسپارک   | <i>Spark Master</i>    |
| اجراکننده تسک‌های پردازشی  | <i>Spark Worker</i>    |
| اپلیکیشن <i>Structured Streaming</i> با دسته‌بندی، لیست سیاه و تشخیص بات | <i>Spark Processor</i> |

جدول ۱: اجزای معماهی سیستم

## ۳. گام اول: تولید جریان داده

### ۳.۱. پیاده‌سازی *Kafka Producer*

اسکریپت تولیدکننده فایل *creditcard.csv* را خط به خط خوانده و تراکنش‌ها را به کافکا ارسال می‌کند.

## ۲.۳. منطق تخصیص کاربر

| عمل                | User_ID        | شرط       | سناریو   |
|--------------------|----------------|-----------|----------|
| ارسال یکبار        | User_Hacker    | Class = 1 | تقلب     |
| ارسال ۵ بار متوالی | User_Bot       | احتمال ۲٪ | حمله بات |
| ارسال یکبار        | User_0001-1000 | پیشفرض    | عادی     |

جدول ۲: منطق تخصیص کاربر

## ۳.۳. کد پیاده‌سازی

```
if transaction_class == 1:
    user_id = 'User_Hacker'
elif random.random() < 0.02:
    user_id = 'User_Bot'
    for i in range(5):
        producer.send(topic, record)
else:
    user_id = random.choice(USER_POOL)
```

## ۴.۳. خروجی تولیدکننده (نتایج واقعی)

```
Connected to Kafka at kafka:29092
FRAUD DETECTED! Transaction TXN_1770283968053 - User: User_Hacker, Amount: $179.66
BOT ATTACK SIMULATION - Sending 5 rapid transactions
Bot attack completed - User: User_Bot, Amount: $18.13 x 5
Processed 10000 transactions (Normal: 9054, Fraud: 36, Bot Attacks: 182)
```

## ۴. گام دوم: پردازش و تحلیل

### ۱.۴. دسته‌بندی تراکنش‌ها

یک ستون جدید به نام *Category* بر اساس مبلغ تراکنش ایجاد می‌شود:

| توضیحات           | شرط                   | دسته   |
|-------------------|-----------------------|--------|
| تراکنش‌های بزرگ   | Amount > \$500        | Macro  |
| تراکنش‌های کوچک   | Amount < \$20         | Micro  |
| تراکنش‌های معمولی | \$20 < Amount < \$500 | Normal |

جدول ۳: دسته‌بندی تراکنش‌ها بر اساس مبلغ

```
categorized = stream.withColumn("Category",
    when(col("Amount") > 500, "Macro")
    .when(col("Amount") < 20, "Micro")
    .otherwise("Normal"))
```

## ۲.۴. تشخیص کاربران لیست سیاه (JOIN)

یک DataFrame دستی شامل کاربران لیست سیاه ایجاد شده و با استریم زنده JOIN می‌شود:

```
blacklist_data = [( 'User_Hacker' , 'Known Fraud Actor' , '2024-01-01' )]  
blacklist_df = spark.createDataFrame(blacklist_data , schema)  
  
joined = transactions.join(blacklist_df , on='User_ID' , how='left')  
  
result = joined.withColumn('User_Status',  
    when(col('Block_Reason').isNotNull() , 'BLOCKED')  
    .otherwise('ACTIVE'))
```

## ۳.۴. تشخیص بات با پنجره زمانی (Windowing)

از قابلیت پنجره‌بندی اسپارک برای تشخیص الگوهای فرکانس بالا استفاده می‌شود:

- مدت پنجره: ۱۰ ثانیه (لغزان با فاصله ۵ ثانیه)
- آستانه: بیش از ۴ تراکنش در پنجره
- هشدار: BOT\_DETECTED

```
windowed = stream.withWatermark('event_time' , '20 seconds')  
.groupBy(  
    window(col('event_time') , '10 seconds' , '5 seconds'),  
    col('User_ID'))  
)  
.agg(count('*').alias('transaction_count'))  
.filter(col('transaction_count') > 4)  
.withColumn('Alert_Status' , lit('BOT_DETECTED'))
```

## ۵. گام سوم: نمایش نتایج

### ۱.۵. خروجی تشخیص بات (نتایج واقعی)

| Batch: 119 | window_start        | window_end          | User_ID  | transaction_count | Alert_Status |
|------------|---------------------|---------------------|----------|-------------------|--------------|
|            | 2026-02-05 09:33:15 | 2026-02-05 09:33:25 | User_Bot | 5                 | BOT_DETECTED |
|            | 2026-02-05 09:33:10 | 2026-02-05 09:33:20 | User_Bot | 5                 | BOT_DETECTED |

## ۲.۵. آمار کلی (نتایج واقعی)

| total_transactions | macro_count | micro_count | normal_count | blocked_count | fraud_count | avg_amount |
|--------------------|-------------|-------------|--------------|---------------|-------------|------------|
| 175                | 4           | 80          | 91           | 1             | 1           | 68.33      |

## ۳.۵. آمار نهایی

| معیار                                    | مقدار   |
|--|---------|
| کل تراکشن‌های پردازش شده                 | +10,300 |
| تراکشن‌های عادی                          | 9,327   |
| تراکشن‌های تقلیبی ( <i>User_Hacker</i> ) | 38      |
| رویدادهای حمله بات                       | 187     |
| تراکشن‌های بات ( $187 \times 5$ )        | 935     |

جدول ۴: آمار نهایی پردازش

## ۶. پیاده‌سازی Docker

### ۱.۶. سرویس‌های کانتینری

| سرویس           | ایمیج              | کاربرد              |
|-----------------|--------------------|---------------------|
| zookeeper       | cp-zookeeper:7.5.0 | هماهنگی کافکا       |
| kafka           | cp-kafka:7.5.0     | بروکر پیام          |
| spark-master    | apache/spark:3.5.0 | هماهنگ‌کننده کلاستر |
| spark-worker    | apache/spark:3.5.0 | اجرای کننده تسک     |
| kafka-producer  | Custom Python      | استریم داده         |
| spark-processor | apache/spark:3.5.0 | پردازش استریم       |

جدول ۵: سرویس‌های کانتینری

## ۲.۶. دستورات اجرا

```
podman compose up -d --build  
podman compose logs -f kafka-producer  
podman compose logs -f spark-processor  
podman compose down -v
```