

PERAKENDE SATIŞ VERİLERİNDE ANOMALİ TESPİTİ

SELİN AVCI (211307026) NİSA AKSOY (211307089)
KOCAELİ ÜNİVERSİTESİ BİLİŞİM SİSTEMLERİ MÜHENDİSLİĞİ

Özet

Perakende Satış Verilerinde Anomali Tespiti isimli proje için önce bir veri seti seçilmiştir. Seçilen veri seti üzerinde veri temizleme, düzenleme, ön işleme, görselleştirme adımları yapılmıştır. Bu işlemler sonrasında Apache Kafka ve Apache Spark ile verilerin entegrasyonu sağlanmıştır. Makine öğrenmesi ve derin öğrenme teknikleri kullanılarak “Perakende Satış Verilerinde Anomali Tespiti” projesinin test sonuçları elde edilmiştir.

Veri Seti – Online Retail

Online Retail Veri Seti, Birleşik Krallık merkezli bir çevrim içi perakende şirketinin **1 Aralık 2010** ile **9 Aralık 2011** tarihleri arasındaki sekiz aylık alışveriş verilerini içermektedir. Veri seti, e-ticaret davranışlarını, müşteri trendlerini ve satış performansını analiz etmek için zengin bir kaynaktır. Aşağıda veri setine dair detaylar yer almaktadır:

Sütunlar

1. InvoiceNo (Fatura Numarası)

- Her işlem için atanmış benzersiz fatura numaraları.
- Örnek: 536365
- Bireysel alışverişlerin takibini sağlar.

2. StockCode (Ürün Kodu)

- Satılan ürünlere atanan kodlar.
- Örnek: 85123A

3. Description (Açıklama)

- Alınan ürünlerin isimleri veya açıklamaları.
- Örnek: WHITE HANGING HEART T-LIGHT HOLDER

4. Quantity (Miktar)

- Her üründen satın alınan birim sayısı.
- Örnek: 6

5. InvoiceDate (Fatura Tarihi)

- İşlemin gerçekleştiği tarih ve saat.
- Örnek: 12/1/2010 8:26

6. UnitPrice (Birim Fiyat)

- Ürünün bir birim fiyatı (GBP cinsinden).
- Örnek: 2.55

7. CustomerID (Müşteri Kimliği)

- Her müşteriye atanmış benzersiz kimlik numaraları.

- Örnek: 17850
- Müşteri davranışlarını ve sadakatini analiz etmek için kullanılabilir.

8. Country (Ülke)

- Siparişin verildiği ülke.
- Örnek: United Kingdom(Birleşik Krallık)

Veri Setinin Önemli Özellikleri

- **İşlem Hacmi:**
Veri seti toplamda 541.909 işlem kaydı içermektedir. Bu, trend analizi ve modelleme için büyük bir veri kaynağıdır.
- **Ürün Çeşitliliği:**
4.070 farklı ürün kodu yer almakta olup, geniş bir ürün yelpazesi sunmaktadır.
- **Müşteri Bilgisi:**
37.000+ benzersiz müşteri bilgisi, müşteri segmentasyonu ve kişiselleştirilmiş pazarlama analizi yapılmasını sağlar.
- **Negatif Miktarlar ve İadeler:**
Veri setinde negatif miktarlar (iade işlemleri) da yer almaktadır. Bu, hem iade süreçlerini hem de müşteri memnuniyetini analiz etmeye olanak tanır.

Veri Ön İşleme ve Görselleştirme

Veri ön işleme ve görselleştirme adımlarının uygulanabilmesi için önce çeşitli çıkarımlar yapılmıştır.

- **İlk 5 Satırın Gösterimi:** Veri setinin ilk 5 satırı görüntülenmiştir.
- **Veri Seti Bilgisi:** `pandas.DataFrame.info()` metodu kullanılarak veri setindeki sütun adları, veri türleri ve eksik değer sayıları görüntülenmiştir.
- **Eksik Değerlerin Tespiti:** Hangi sütunlarda eksik değer olduğu kontrol edilmiş ve bu eksik değerlerin toplamı listelenmiştir.
- **Betimsel İstatistiklerin Çıkarılması:** `DataFrame.describe()` metodu ile sayısal sütunların temel istatistikleri (ortalama, standart sapma, minimum, maksimum vb.) analiz edilmiştir. Ayrıca kategorik veriler için `unique`, `topve` `fraqqibi` bilgiler elde edilmiştir.

Bu bilgilere göre aşağıdaki adımlar uygulanmıştır.

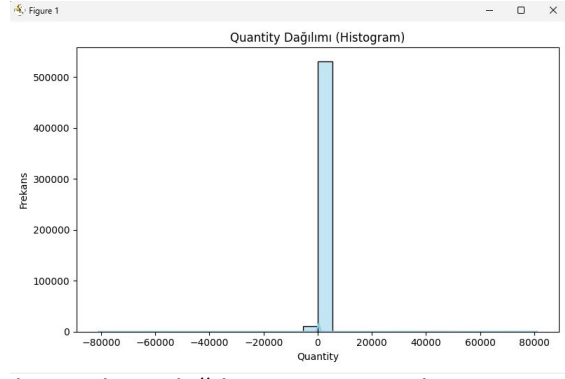
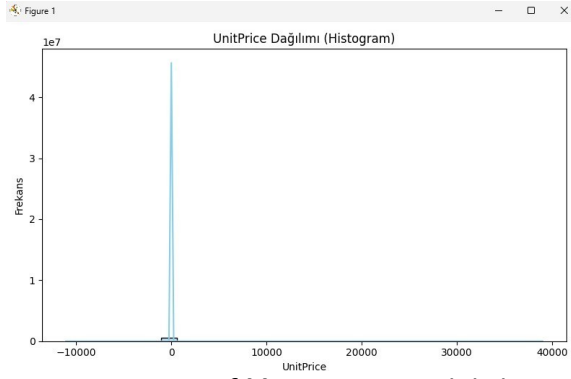
1. Eksik Değerlerin Temizlenmesi

- Eksik değerlere sahip satır veya sütunların temizlendiği belirtilmiştir. Bu işlem genellikle `dropna()` veya eksik değerlerin doldurulması için `fillna()` kullanılarak yapılır.
- Eksik değer temizliğinden sonra veri seti tekrar görüntülenmiştir.

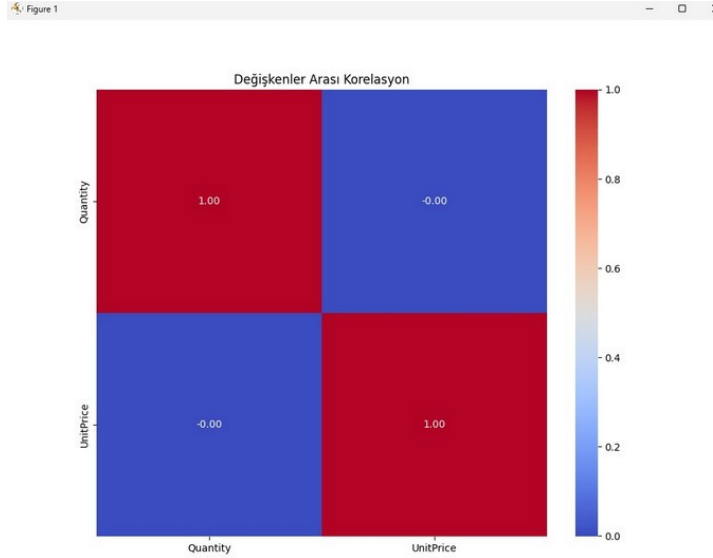
2. Veri Türlerinin Düzenlenmesi

- Bazı sütunların veri türleri (örneğin, `InvoiceDate`) uygun formatlara dönüştürülmüş olabilir (örneğin, tarih-saat formatı).
- Bu, genelde `pd.to_datetime()` gibi yöntemlerle yapılır.

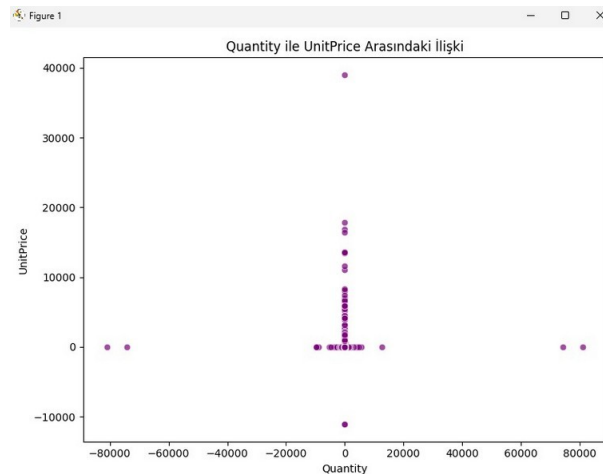
3. Veri Görselleştirme



- **Histogram Grafikleri:** Veri setindeki bazı sayısal sütunların dağılımını göstermek için histogramlar oluşturulmuş.
- **Korelasyon Haritası (Heatmap):** Sütunlar arasındaki ilişkiyi göstermek için korelasyon matrisi çıkarılmış ve görselleştirilmiş.



- **Scatter Plot:** Quantity ve UnitPrice sütunları arasındaki ilişkiyi göstermek için bir scatter plot (dağılım grafiği) oluşturulmuş.



4. Aykırı Değerlerin Tespiti ve İşlenmesi

- **Quantityve UnitPriceSütunları:** Aykırı değerler tespit edilmiş ve bu değerler listelenmiş.
- Bu işlem genellikle IQR (Interquartile Range), Z-score veya diğer aykırı değer tespit yöntemleriyle yapılır.
- Aykırı değerler ya temizlenmiş ya da yeniden işlenmiş olabilir.

5. Verinin Normalize Edilmesi veya Standartlaştırılması

- **Normalizasyon veya Standartlaştırma:** Sayısal sütunlar (örneğin, Quantityve UnitPrice) belirli bir aralıkta (örneğin, 0 ile 1 arasında) normalize edilmiş.
- Bu işlem genellikle MinMaxScalerveya StandardScalergibi yöntemlerle yapılır.

6. Kategorik Değişkenlerin Sayısallaştırılması (Encoding)

- Countrygibi kategorik değişkenler sayısal değerlere dönüştürülmüş. Örneğin, "United Kingdom" → 36gibi bir dönüşüm yapılmış.
- Bu işlem LabelEncoderveya OneHotEncoderkullanılarak yapılmış olabilir.

7. Temizlenmiş ve İşlenmiş Veri Setinin Hazırlanması

- Bütün işlemlerden sonra elde edilen veri seti ekrana yazdırılmıştır.

Apache Kafka – Apache Spark

Apache Kafka, yüksek hacimli verilerin hızlı ve güvenilir bir şekilde işlenmesini sağlayan bir mesajlaşma altyapısıdır. Genellikle genellikle gerçek zamanlı veri akışlarının merkezi bir platform olarak kullanılan Kafka'da:

- **Producer (Üretici):** Verileri Kafka'ya gönderir.
- **Consumer (Tüketici):** Kafka'dan verileri alır ve işler.
- **Topic:** Verilerin kategorilere ayrılarak saklandığı birimdir.

Apache Spark, büyük veri işlemleri için hızlı ve güçlü bir çerçevedir. Spark Streaming modülü, gerçek zamanlı veri akışlarını işlemek için kullanılır ve Kafka ile entegre olarak çalışabilir.

Perakende satış verilerinde anomali tespiti için Spark, anormal desenleri tespit etmek üzere verileri analiz eder ve sonuçları kaydeder.

Proje kapsamında sırasıyla aşağıdaki işlemler yapılır:

- Kafka'dan retail_saleskonusundan veriler alınır (Consumer).

```
Created topic retail_sales.
WARNING: Due to limitations in metric names, topics with a period ('.') or underscore ('_') could collide
but not both.
Created topic normal_data.
Created topic anomalies.
selinavci@Mac kafka_2.12-3.6.0 %
```

- Negatif değer içeren anomaliler anomalies konusuna, normal veriler normal_data konusuna gönderilir (Consumer).

```
Message delivered to normal_data [0]
Message delivered to normal_data [0]
Message delivered to normal_data [0]
Message delivered to anomalies [0]
Message delivered to normal_data [0]
Message delivered to normal_data [0]
Message delivered to normal_data [0]
selinavci@Mac kafka_2.12-3.6.0 % bin/kafka-topics.sh --list --bootstrap-server localhost:9092

__consumer_offsets
anomalies
normal_data
retail_sales
```

- CSV dosyasındaki perakende satış verileri okunur.
- Her veri Kafka'daki retail_sales konusuna gönderilir (Producer)

SparkSession ve SparkContext Oluşturma:

- Bu adımda, Apache Spark uygulamasını başlatmak için `SparkSession` ve `SparkContext` oluşturulur. Spark, büyük veri analizini paralel bir şekilde gerçekleştirir ve veri işleme işlevlerini sağlar. Burada, Spark'ı "Retail Anomaly Detection" adıyla başlatıyoruz ve veri işlemeyi hızlandırmak için yapılandırmalar ekliyoruz.

Veri Hazırlama

- Verilerin işlenmesi ve model eğitimi için uygun hale getirilmesi gerekmektedir. Bu işlevde:
- - Veriler CSV dosyasından okunur.
- - `Quantity` ve `UnitPrice` sıfırdan büyük olan satırlar seçilir.
- - Eksik değerler (`na.drop()`) temizlenir.
- - `TotalPrice` (toplam fiyat) sütunu, `Quantity` ve `UnitPrice` çarpılarak hesaplanır.
- - Sonrasında, bu veriler bir `VectorAssembler` kullanılarak sayısal özellikler (`Quantity`, `UnitPrice`, `TotalPrice`) tek bir vektöre dönüştürülür.
- - `label` adlı yeni bir sütun eklenir; burada, eğer `Quantity > 10` ve `UnitPrice > 5` ise etiket 1 (normal), aksi takdirde 0 (anormal) olur.

SparkSession ve Random Forest Modeli Eğitimi

- Bu fonksiyon, eğitim verilerini kullanarak bir Random Forest sınıflandırıcı modeli oluşturur. Random Forest, birçok karar ağacından oluşan ve sınıflandırma yapabilen bir makine öğrenmesi modelidir. Burada modelin özellikleri:
- - `featuresCol` olarak özelliklerin bulunduğu sütun belirlenir (`features`).
- - `labelCol` olarak etiketlerin bulunduğu sütun belirlenir (`label`).
- - 100 ağaçla model eğitilir ve model geri döndürülür.
- 4. Autoencoder Modeli Eğitimi
- Autoencoder, verinin sıkıştırılması ve yeniden yapılandırılması temeline dayanan bir yapay sinir ağı modelidir. Burada:

- Autoencoder, giriş verisini düşük boyutlu bir uzaya sıkıştırıp tekrar eski haline getirmeye çalışır.
- Modelin girdi boyutu, veri setinin özellik sayısına (özellik vektörü uzunluğu) göre belirlenir.
- Modelin yapısı, 64, 32 ve 64 nörondan oluşan 3 katmandan meydana gelir ve çıktı, orijinal veriyi yeniden oluşturmayı hedefler.
- Model, `mse` (mean squared error) kayıp fonksiyonu ile eğitilir ve eğitim tamamlandığında `autoencoder_model.h5` dosyası olarak kaydedilir.

Kafka Entegrasyonu

- Bu adımda, bir veri akışının işlenmesi için Apache Kafka kullanılır. Kafka, gerçek zamanlı veri akışlarını işlemek için kullanılan bir platformdur.
- `process_stream` fonksiyonu, verileri Kafka'dan alır, gerekli veriyi çıkarır ve bu veriler üzerinde işlem yapar.
- Veriler önce JSON formatında alınır, sonra özellikler `VectorAssembler` ile birleştirilir.
- Random Forest modelinden tahminler yapılır.
- Autoencoder ile verinin yeniden yapılandırılması gerçekleştirilir ve MSE (mean squared error) hesaplanarak anomali tespiti yapılır. MSE değeri belirli bir eşik değeri aşarsa, veri anormal olarak sınıflandırılır.

```

1 from pyspark.sql import SparkSession

/dev/fd/13:18: command not found: compdef
selinavci@Mac retail-anomaly-detection % python3 spark_model.py
Setting default log level to 'WARN'.
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
24/12/19 18:05:58 WARN NativeCodeLoader: Unable to load native-heapoop library for your platform... using builtin-java classes where applicable
Training Random Forest...
Training Autoencoder...
/Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages/keras/src/layers/core/dense.py:87: UserWarning: Do not pass an 'input_shape' argument to a layer. When using Sequential models, prefer using an 'Input(shape)' object as the first layer in the model instead.
  super().__init__(activity_regularizer=activity_regularizer, **kwargs)
Epoch 1/10
9958/9958 4s 391us/step - loss: 39916.2773
Epoch 2/10
9958/9958 4s 378us/step - loss: 6204.8315
Epoch 3/10
9958/9958 4s 383us/step - loss: 35009.9844
Epoch 4/10
9958/9958 4s 378us/step - loss: 128580.4453
Epoch 5/10
9958/9958 4s 394us/step - loss: 45902.8398
Epoch 6/10
9958/9958 4s 375us/step - loss: 52640.1211
Epoch 7/10
9958/9958 4s 380us/step - loss: 45013.2656
Epoch 8/10
9958/9958 4s 381us/step - loss: 19758.6445
Epoch 9/10
9958/9958 4s 375us/step - loss: 33851.4883
Epoch 10/10
9958/9958 4s 385us/step - loss: 13651.9158
WARNING:absl:You are saving your model as an H5 file via 'model.save()' or 'keras.save_model(model)'. This file format is considered legacy. We recommend using instead the native Keras format, e.g. 'model.save('my_model.keras')' or 'keras.save_model(model, 'my_model.keras')'.
Evaluating Random Forest...
Random Forest Accuracy: 1.00
Starting Kafka Stream Processing...
/Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages/pyspark/streaming/context.py:72: FutureWarning: DStream is deprecated as of Spark 3.4.0. Migrate to Structured Streaming.
  warnings.warn(
24/12/19 18:16:53 WARN ReceiverSupervisorImpl: Restarting receiver with delay 2000 ms: Socket data stream had no more data
24/12/19 18:16:53 ERROR ReceiverTracker: Deregistered receiver for stream 0: Restarting receiver with delay 2000ms: Socket data stream had no more data
[Stage 22:] (0 + 1) / 1]

```

Main İşlevi

- Veriler yüklenir ve eğitim ile test setlerine ayrılır.
- Random Forest modeli eğitilir ve sonuçları değerlendirir.
- Autoencoder modeli eğitilir.
- Son olarak, Random Forest modelinin doğruluğu değerlendirilir.

- Kafka akış işleme fonksiyonu başlatılır ve gerçek zamanlı veri işlenmeye başlanır.

Spark Jobs (?)

User: selinavci
Total Uptime: 11 min
Scheduling Mode: FIFO
Active Jobs: 1
Completed Jobs: 14

Event Timeline

Active Jobs (1)

Page: 11 Pages. Jump to 1. Show 100 items in a page. Go

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
14	Streaming job running receiver 0 start at NativeMethodAccessorImpl.java:0 (kill)	2024/12/19 18:06:52	9,6 min	0/1	0/1 (1 running)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Completed Jobs (14)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
13	collectAsMap at MulticlassMetrics.scala:61 collectAsMap at MulticlassMetrics.scala:61	2024/12/19 18:06:51	1 s	2/2	16/16
12	collect at /Users/selinavci/Desktop/retail-anomaly-detection/spark_model.py:102 collect at /Users/selinavci/Desktop/retail-anomaly-detection/spark_model.py:102	2024/12/19 18:06:06	5 s	1/1	8/8
11	collectAsMap at RandomForest.scala:663 collectAsMap at RandomForest.scala:663	2024/12/19 18:06:04	1 s	2/2	16/16
10	collectAsMap at RandomForest.scala:663	2024/12/19 18:06:03	1 s	2/2	16/16

- Bu işlem hattı, perakende verileri üzerinde hem sınıflandırma (Random Forest) hem de anomali tespiti (Autoencoder) yapar, ve sonuçları Kafka üzerinden işlemeye olanak tanır. Bu, özellikle alışveriş verileri gibi büyük veri setleriyle çalışırken yararlı olabilir.

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
13	collectAsMap at MulticlassMetrics.scala:61 collectAsMap at MulticlassMetrics.scala:61	2024/12/19 18:06:51	1 s	2/2	16/16
12	collect at /Users/selinavci/Desktop/retail-anomaly-detection/spark_model.py:102 collect at /Users/selinavci/Desktop/retail-anomaly-detection/spark_model.py:102	2024/12/19 18:06:06	5 s	1/1	8/8
11	collectAsMap at RandomForest.scala:663 collectAsMap at RandomForest.scala:663	2024/12/19 18:06:04	1 s	2/2	16/16
10	collectAsMap at RandomForest.scala:663 collectAsMap at RandomForest.scala:663	2024/12/19 18:06:03	1 s	2/2	16/16
9	collectAsMap at RandomForest.scala:663 collectAsMap at RandomForest.scala:663	2024/12/19 18:06:01	1 s	2/2	16/16
8	collectAsMap at RandomForest.scala:663 collectAsMap at RandomForest.scala:663	2024/12/19 18:06:00	0,9 s	2/2	16/16
7	collectAsMap at RandomForest.scala:663 collectAsMap at RandomForest.scala:663	2024/12/19 18:05:58	2 s	2/2	16/16
6	collectAsMap at RandomForest.scala:1054 collectAsMap at RandomForest.scala:1054	2024/12/19 18:05:58	0,6 s	2/2	11/11
5	aggregate at DecisionTreeMetadata.scala:125 aggregate at DecisionTreeMetadata.scala:125	2024/12/19 18:05:57	0,8 s	1/1	8/8
4	take at DecisionTreeMetadata.scala:119 take at DecisionTreeMetadata.scala:119	2024/12/19 18:05:57	0,3 s	1/1	1/1
3	take at DatasetUtils.scala:193 take at DatasetUtils.scala:193	2024/12/19 18:05:57	36 ms	1/1 (1 skipped)	1/1 (8 skipped)
2	take at DatasetUtils.scala:193 take at DatasetUtils.scala:193	2024/12/19 18:05:55	2 s	1/1	8/8
1	csv at NativeMethodAccessorImpl.java:0 csv at NativeMethodAccessorImpl.java:0	2024/12/19 18:05:53	1 s	1/1	8/8
0	csv at NativeMethodAccessorImpl.java:0 csv at NativeMethodAccessorImpl.java:0	2024/12/19 18:05:53	0,2 s	1/1	1/1

Perakende Satış Verilerinde Anomali Tespiti Github Link

- <https://github.com/SelinAVCI09/big-data-project.git>
- <https://www.kaggle.com/datasets/vijayuv/onlineretail>

-

Kaynakça:

- <https://spark.apache.org/downloads.html>
- <https://kafka.apache.org/downloads>
- <https://www.python.org/>