# SAKARYA
## ÜNİVERSİTESİ

BİLGİSAYAR-BİLİŞİM BİLİMLERİ FAKÜLTESİ
BİLHİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Büyük Veriye Giriş dersi 2023-2024 Güz dönemi Proje raporu

Hazırladı: Fuad Garibli
Öğrenci Numarası: G201210558
Şube numarası: 2A

Bizden istenen Kafka ve Spark kullanılarak, gerçek zamanlı verilerin işlenmesine yönelik bir işlem hattını göstermektedir. Sistem, veri üretimi, kafka mesaj üretimi, Spark Yapılandırılmış Akış (structed streaming), Spark makine öğrenmesi (Spark ML) ve pipeline entegrasyonu için bileşenler içerir.

1) Öncelikle Veri setimizi tanıtalım:

Veri setimiz totalde 1510 satırdan ve 6 sütundan oluşan bir ev fiyatları listesidir. Kanadada Vancouver eyaletinden toplanmış bir verisetidir (housing.csv). .CSV dosya formatında olup içinde comma seperated values (noktayla ayrılmış veriler) vardır:



2) Bu veri setini, PyCharmda oluşturduğumuz "csvUpload.py" dosyasında olan kodlar sayesinde kafka producere aktarıyoruz.

{"SquareFeet": "1996", "Bedrooms": "5", "Bathrooms": "2", "Neighborhood": "Urban", "YearBuilt": "2002", "Price": "208832.5771638779"}
{"SquareFeet": "2434", "Bedrooms": "5", "Bathrooms": "1", "Neighborhood": "Suburb", "YearBuilt": "1979", "Price": "181513.9862070892"}
{"SquareFeet": "1950", "Bedrooms": "4", "Bathrooms": "3", "Neighborhood": "Suburb", "YearBuilt": "2013", "Price": "230751.60475472733"}
{"SquareFeet": "1288", "Bedrooms": "2", "Bathrooms": "1", "Neighborhood": "Urban", "YearBuilt": "1975", "Price": "206363.55286603371"}
{"SquareFeet": "2277", "Bedrooms": "5", "Bathrooms": "1", "Neighborhood": "Suburb", "YearBuilt": "1980", "Price": "161369.96935538173"}
{"SquareFeet": "1195", "Bedrooms": "2", "Bathrooms": "2", "Neighborhood": "Rural", "YearBuilt": "1971", "Price": "177762.60744982713"}
{"SquareFeet": "1597", "Bedrooms": "2", "Bathrooms": "3", "Neighborhood": "Suburb", "YearBuilt": "1993", "Price": "191429.0700920061"}
{"SquareFeet": "2857", "Bedrooms": "3", "Bathrooms": "2", "Neighborhood": "Suburb", "YearBuilt": "1975", "Price": "296301.481125016"}
{"SquareFeet": "2545", "Bedrooms": "2", "Bathrooms": "2", "Neighborhood": "Rural", "YearBuilt": "2014", "Price": "241398.23310669637"}
{"SquareFeet": "1516", "Bedrooms": "5", "Bathrooms": "3", "Neighborhood": "Urban", "YearBuilt": "1958", "Price": "169922.14559428632"}
{"SquareFeet": "2658", "Bedrooms": "2", "Bathrooms": "1", "Neighborhood": "Suburb", "YearBuilt": "1999", "Price": "267022.57863697683"}
{"SquareFeet": "2865", "Bedrooms": "3", "Bathrooms": "2", "Neighborhood": "Rural", "YearBuilt": "1998", "Price": "282967.21274589025"}
{"SquareFeet": "1864", "Bedrooms": "4", "Bathrooms": "3", "Neighborhood": "Rural", "YearBuilt": "1957", "Price": "166514.29821092534"}
{"SquareFeet": "2141", "Bedrooms": "4", "Bathrooms": "1", "Neighborhood": "Rural", "YearBuilt": "1995", "Price": "310230.40635129646"}
{"SquareFeet": "1805", "Bedrooms": "3", "Bathrooms": "1", "Neighborhood": "Urban", "YearBuilt": "2009", "Price": "236104.8196810374"}
{"SquareFeet": "2932", "Bedrooms": "5", "Bathrooms": "1", "Neighborhood": "Suburb", "YearBuilt": "1968", "Price": "276777.87852348736"}
{"SquareFeet": "2182", "Bedrooms": "5", "Bathrooms": "2", "Neighborhood": "Suburb", "YearBuilt": "1964", "Price": "248861.92116573846"}
{"SquareFeet": "1814", "Bedrooms": "5", "Bathrooms": "2", "Neighborhood": "Urban", "YearBuilt": "1976", "Price": "218993.10895829712"}

3)Veri üretiminin süresini kontrol etmek için UploadByLimit.py scripti, belirli bir süre boyunca Kafka'ya veri gönderir. Bu, Kafka'ya aktarılan veri miktarını test etmek ve sınırlamak için kullanılacaktır. Bu süre zarfında Streaming yapabilecek duruma geliyoruz:
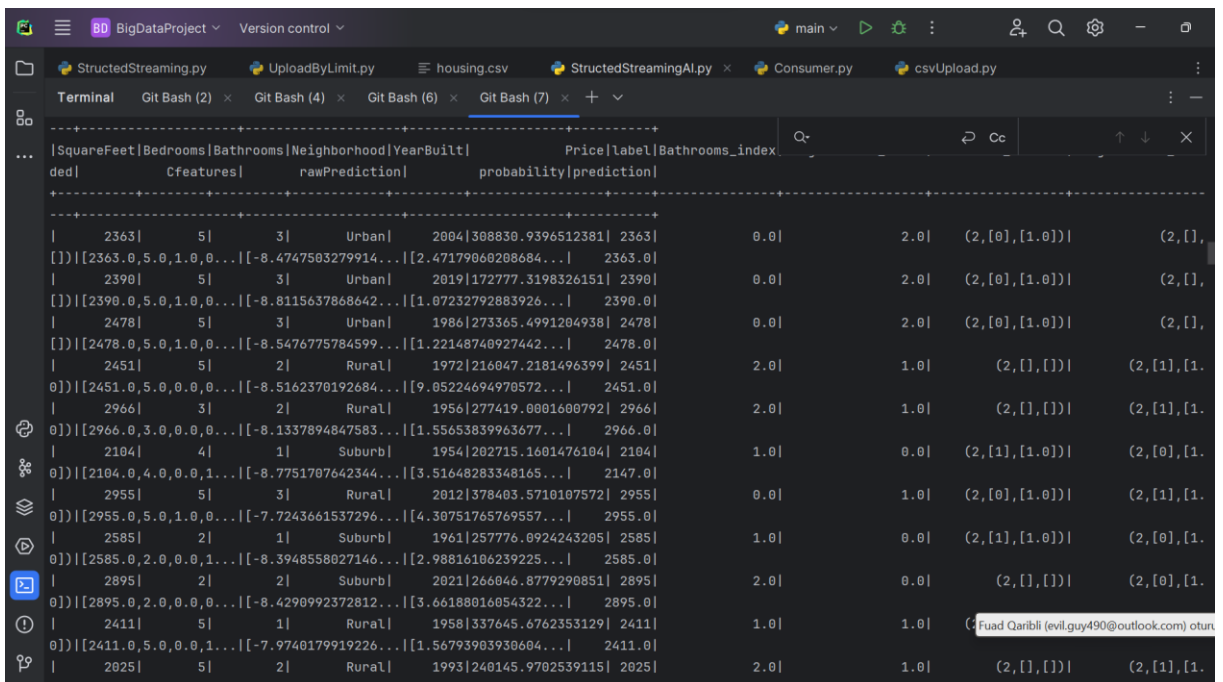
3) StructedStreaming.py çalıştığı zaman artık veriler kafka topic'e yüklenerek bizim Streamingimiz gerçekleştiriyor. Buna uygun çıktı aşağıda gösterilmiştir:



4) StructedStreamingAI.py dosyası yukarda tanımlanan işlemleri yaparak, üstüne Machine Learning tekniklerini kullanarak bize herhangi bir senede inşa edilmiş evin alanına göre fiyatını tahmin ediyor. Çıktısı aşağıda verilmiştir:

## 5) Bazı Spark LocalHost çıktıları aşağıda verilmiştir:

**Streaming Query Statistics**

Running batches for **9 minutes 55 seconds** since **2023/12/26 15:31:56** (40 completed batches)

**Name:** <no name>
**Id:** 94041982-64bf-47a2-b4ce-1e5c78e9b454
**RunId:** 65d917ec-6026-4c68-b165-fc7b77e3a947



**Stages for All Jobs**

**Completed Stages:** 40

▾ **Completed Stages (40)**

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

| Stage Id ▾ | Description | Submitted | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|---|---|---|---|---|---|---|---|---|
| 39 | id = 94041982-64bf-47a2-b4ce-1e5c78e9b454 runId = 65d917ec-6026-4c68-b165-fc7b77e3a... start at NativeMethodAccessorImpl.java:0 +details | 2023/12/26 15:38:12 | 0.5 s | 1/1 | | | | |
| 38 | id = 94041982-64bf-47a2-b4ce-1e5c78e9b454 runId = 65d917ec-6026-4c68-b165-fc7b77e3a... start at NativeMethodAccessorImpl.java:0 +details | 2023/12/26 15:38:10 | 0.2 s | 1/1 | | | | |
| 37 | id = 94041982-64bf-47a2-b4ce-1e5c78e9b454 runId = 65d917ec-6026-4c68-b165-fc7b77e3a... start at NativeMethodAccessorImpl.java:0 +details | 2023/12/26 15:38:09 | 0.2 s | 1/1 | | | | |
| 36 | id = 94041982-64bf-47a2-b4ce-1e5c78e9b454 runId = 65d917ec-6026-4c68-b165-fc7b77e3a... start at NativeMethodAccessorImpl.java:0 +details | 2023/12/26 15:38:08 | 0.2 s | 1/1 | | | | |
| 35 | id = 94041982-64bf-47a2-b4ce-1e5c78e9b454 runId = 65d917ec-6026-4c68-b165-fc7b77e3a... start at NativeMethodAccessorImpl.java:0 +details | 2023/12/26 15:38:06 | 0.3 s | 1/1 | | | | |
| 34 | id = 94041982-64bf-47a2-b4ce-1e5c78e9b454 runId = 65d917ec-6026-4c68-b165-fc7b77e3a... start at NativeMethodAccessorImpl.java:0 +details | 2023/12/26 15:38:05 | 0.2 s | 1/1 | | | | |
| 33 | id = 94041982-64bf-47a2-b4ce-1e5c78e9b454 runId = 65d917ec-6026-4c68-b165-fc7b77e3a... start at NativeMethodAccessorImpl.java:0 +details | 2023/12/26 15:38:04 | 0.2 s | 1/1 | | | | |
| 32 | id = 94041982-64bf-47a2-b4ce-1e5c78e9b454 runId = 65d917ec-6026-4c68-b165-fc7b77e3a... start at NativeMethodAccessorImpl.java:0 +details | 2023/12/26 15:38:03 | 0.2 s | 1/1 | | | | |

## 6) Akış diagramı: