

İnönü Üniversitesi
Bilgisayar Mühendisliği Bölümü
Bitirme Projesi (2022-2023 GÜZ)

Proje başlığı	Trend Aramaların Sınıflandırılması
Öğrenci(ler)	Civan BOZAN - Hafta 8
Proje özeti	<p>Yöntem: <u>Apache Spark</u> veri işleme çerçevesi ile arama motorlarında sıklıkla aratılan ifadelerin, <u>Spark Streaming</u> özelliği ile gerçek zamanlı olarak işlenmesi ve <u>Spark MLlib</u> makine öğrenmesi kütüphanesi kullanılarak bu verilerin sınıflandırılması (kategorilere ayrılması) sağlanıp uygun formatta (json, xml..) gerçek zamanlı olarak sunulması ve geçmişe yönelik analizlerin gerçekleştirilebilmesi açısından bu verilerin <u>Apache HBase</u> veritabanı ile tutulması ve ilgili analizlerin bir <u>web sitesi</u> aracılığıyla servis edilmesi.</p> <p>Amaç: Kullanıcı eğilimlerinin tespit edilerek, ilgili kategoriler hakkında içerik üreten yazarlara ilham kaynağı olmak.</p>
Tamamlanması gereken özellikler	<ol style="list-style-type: none">1. Spark ve ilgili kütüphanelerin incelenmesi.2. HBase incelenmesi.3. HBase ve Spark entegrasyonu4. Arama motorlarından sıklıkla aratılan ifadelerin elde edilmesi.5. Elde edilen ifadelerin spark ile gerçek zamanlı olarak işlenmesi.6. Verilerin kategorilere ayrılması ve uygun formata dönüştürülmesi. CRISP-DM süreci<ol style="list-style-type: none">a. İşin Anlaşılmasıb. Verinin Anlaşılmasıc. Verinin Hazırlanmasıd. Modellemee. Değerlendirmef. Konuşlandırma7. Spark ile işlenmiş tüm verilerin HBase'de tutulması.8. Geçmiş verilerin web sitesi üzerinden sunulması.
Sonradan eklenmesi gereken özellikler	<ol style="list-style-type: none">1. Verilerin bir API aracılığıyla sunulması.2. Abonelik özelliği ile sadece ilgili kategorilere ait aramaların sunulduğu bir web sitesi.
Bu hafta gerçekleştirilenler	<ol style="list-style-type: none">1. Spark kümesi oluşturulması için gerekli hazırlıklar üzerinde çalışıldı. (Sanal makine, hadoop bileşenleri..)2. CRISP-DM sürecinin ilk 3 aşaması gerçekleştirildi.<ul style="list-style-type: none">• İşin Anlaşılması• Verinin Anlaşılması• Verinin Hazırlanması
Kaynaklar	<ol style="list-style-type: none">1. Veri seti: https://www.kaggle.com/datasets/savasy/multiclass-classification-data-for-turkish-tc32