

T.C.  
SAKARYA ÜNİVERSİTESİ  
FEN BİLİMLERİ ENSTİTÜSÜ  
BİLGİSAYAR VE BİLİŞİM MÜHENDİSLİĞİ ANABİLİM DALI  
BİLİŞİM TEKNOLOJİLERİ PR.(YL) (UZAKTAN EĞİTİM)

VERİ MADENCİLİĞİ UYGULAMALARI 2.ÖDEV  
BÜYÜK VERİ İLE VERİ MADENCİLİĞİ

Hazırlayan  
SEDAT ÖZTÜRK  
E235013168

Öğretim Üyesi  
Prof. Dr. NİLÜFER YURTAY

NİSAN 2024

## 1. Giriş

Günümüzde meydana gelen gelişmelerin temelini veri, enformasyon ve bilgi oluşturmaktadır. Veri üzerinde işlem yapılmamış küçük bir bilgi parçasıdır. Veriler tek başına bir anlam ifade etmezler. Enformasyon, karar vermek için değer olan ve organize edilmiş verilerin özetlenmesiyle elde edilir. Bilgi de enformasyon verilerinin düzenlenmiş ve analiz edilmesi sonucu değer kazanmış halidir. Bu kavramlar hayatımızın her aşamasında karşımıza çıkar ve oldukça hızlı bir şekilde artarlar.

## 2. Büyük Veri nedir?

Büyük veri, çeşitli kaynaklardan edinilen devasa büyüklükteki verilerin toplanması anlamına gelir. Toplanan veriler, analiz edilerek farklı teknik yöntemler ile anlamlı hale getirilmesi için veriler sınıflandırılır ve işlenir. Büyük veriler akışı hiç kesilmeyen yasa kaynaklarından elde edilen çeşitli veri kümelerinden oluşur.

Büyük verinin temel bileşenleri 5V olarak da bilinenen; çeşitlilik (Variety), hacim (Volume), hız (Velocity), doğruluk (Veracity), değer (Value) özelliklerini bulunmaktadır.

• Çeşitlilik (Variety): Verilerin birbirinden farklı formatlara sahip olmasıdır. Yapılandırılmamış veri, metin, e-posta, fotoğraf, video ve ses gibi birçok farklı biçimde gelir. Yapılandırılmış veri RDBS de (ilişkisel veritabanı) bulunan tablo formatındaki veriler ifade etmektedir.

• Hacim (Volume): Verilerin miktarı için kullanılan terimdir. Yıllar geçtikçe verilerin hacmi hızla artmaktadır.

• Hız (Velocity): Gelişen teknolojiyle birlikte daha hızlı şekilde hareket eder. Hız oranının yükseltilmesi değeri artırır.

• Doğruluk (Veracity): Veri güvencesi en önemli konulardan biri olup sağlıklı ve verimli sonuç alabilmek için hatalı verilerin temizlenmesi gerekir.

• Değer (Value): İşlenen veriler anlamlı ve hizmet ettiği kuruma değer yaratmasıdır.

## 3. Büyük Veri Teknolojileri

Büyük veri teknolojileri verinin depolanması, veri madenciliği, veri analizi ve veri görselleştirme olmak üzere 4 farklı gruba ayrılır.

### a. Veri Depolama:

Veri Depolama, verilerin saklanması ve yönetildiği teknolojileri kapsar. Zaman içerisinde bu verilerin saklamak büyük bir zorluk oluşturmış ve saklanması için belli başlı veri tabanları ve dosya sistemleri geliştirilmiştir. Büyük veri depolama ve analiz etmek için Hadoop, MapReduce ve Cassandra sistemleri örnek gösterilebilir. Bu sistemler yüksek işlem gücü ve çok sayıda eş zamanlı görevlere yönetme yeteneğine sahiptir.

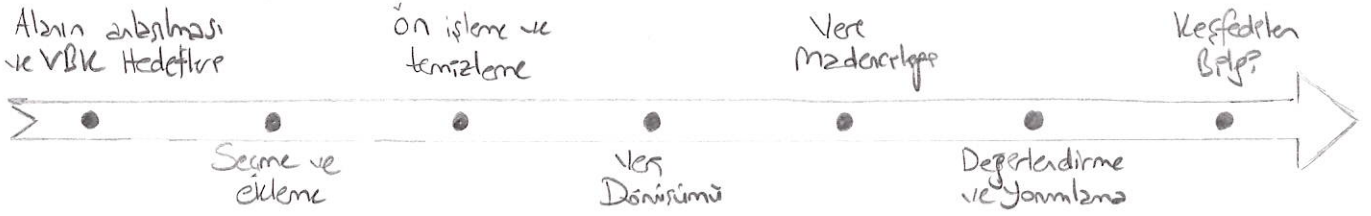
### b. Veri Madenciliği:

Veri Madenciliği, büyük ölçekli veri kümelerinde ihtiyaç yönelik veriler bulur, analiz eder ve işleyerek anlamlandırır. Araştırmacılar büyük verilerle arazi kullanmaksızın analiz edilemeyecekleri geçmişler ve bilgi keşfi kavramını ortaya çıkarmıştır. Bu kavram büyük miktarda veri içindeki değerli bilgilerle bulma, analiz etme ve kullanma sürecini ifade eder.



## Bilgi Keşif Kavramı:

Bilgi Keşif kavramı 7 aşamadan oluşmaktadır.



- Uygulama alanının anlaşılması ve geliştirilmesi: Bu aşamada uygulama alanı anlamaya çalışılarak çeşitli hazırlıklar yapılır ve bilgi keşfinde ki amaç geliştirilir.
- Seçme ve ekleme: Veri tabanından alınan analiz ile ilgili verilerden, probleme ilişkin olan verilere seçme sürecidir.
- Ön işleme ve veri temizliği: Eksik, çelişkili ve tutarsız verilerin temizleme sürecidir.
- Veri Dönüşümü: Bu aşamada verilen uygun formlara dönüştürülüp veri madenciliğinde kullanılabilir hale getirme sürecidir.
- Veri Madenciliği: Hazırlanan veriler üzerinden amacına göre Veri Madenciliği Algoritmalarının uygulanma sürecidir.
- Değerlendirme ve Yorumlama: Bu aşamada oluşan sonuçlar değerlendirilerek birinci aşamada belirlenen amaç ile alakalı kurulmalıdır ve elde edilen bilgi amaçla yönelik kullanılmalıdır.
- Keşfedilen Bilgi: Veri madenciliği elde edilmiş bilginin kullanıcıya sunulmasıdır.

Presto, Apache spark, Apache Flink, Apache Hadoop veri madenciliği için kullanılan en popüler programlardır. Özellikle presto verileri ayrı bir analiz sistemine taşımaya gerek kalmadan depolandığı yerde sorgulanabilir.

## Veri Madenciliği Yöntemleri:

Veri madenciliği yöntemlerine 3 grupta toplayabiliriz. Sınıflama ve regresyon, kümeleme ve birleştirebilir kuralları modelleridir.

Sınıflama ve Regresyon: Sınıflama, verilen önceden belirlenen çıktıları uygun olarak ayrıştırılmasını sağlayan bir tekniktir. Çıktılar, önceden bilindiği için sınıflama, veri kümesinin denetimli olarak öğrenir.

Modelde kullanılan başlıca teknikler: Karar Ağaçları, Yapay Sinir Ağları, Genetik Algoritmalar, K-En Yakın Komşu, Regresyon Analizi, Naive-Bayes, Kaba Kümeler

Kümeleme: Nesnelere benzerliklerle gruplama sürecine kümeleme denir. Kümeleme analiz, temel amacı nesnelere sahip oldukları karakteristیک özelliklere baz alınarak gruplamak olan çok disiplinli teknikler grubudur.

Genel olarak başlıca kümeleme yöntemleri: Bölme yöntemleri, Hiyerarşik yöntemler, Yapanıya tabanlı yöntemler, İzgara tabanlı yöntemler, Model tabanlı yöntemler

Birleştirebilir kuralları: Büyük veri kümeleri arasında birleştirebilir ilişkiler bulur. Bu yöntem büyük miktardaki ilişkili kayıtlardan ilgili birleştirebilir ilişkileri keşfetmek ve şirketlerin karar alma süreçlerini daha verimli hale getirmektedir.

### C. Veri Analizi:

Veri analizi, yararlı bilgilere keşfetmek, sonuç çıkarmak ve karar vermek desteklemek amacıyla verileri incelemek, temizlemek, dönüştürmek ve modellemek için kullanılan bir süreçtir. Veri analizi teknolojisinin sunan isimler arasında Kafka, Splunk, Spark v.b. yazılım platformları yer alıyor.

### d. Veri Görselleştirme:

Veri görselleştirme, verinin grafikler, tablolar veya haritalar sayesinde görsel bir dile dönüştürülmesi anlamına gelir. Görselleştirmede amaç, istatistiksel ve derin bilgilerin klasik formatta sunulan kompleks verileri, kolay alınılabilecek grafik arayüzler ile rahat anlaşılır hale getirmektir. Bu sayede verilerdeki bilgiler, ortaya çıkan trendler ve düşüşler kolayca görülüyor ve yorumlanabiliyor. Bu teknolojiye, verilerin görselleştirilmesini sağlayan Tableau, Microsoft Power BI, Qlik Sense, Chart Blocks, Plotly, Looker v.b. yazılım programları örnek olarak gösterilebilir.