



# BÜYÜK VERİ UYGULAMALARI – DERS 2

Doç. Dr. Yuriy Mishchenko

# PLAN

- Büyük veri nedir, kaynaklar nedir, kullanım alanları
- Örnekler
- Felsefesi ve temel yaklaşım
- Diagramlar
- Yaklaşım prensibi
- Ayrıntılı örnekler

# BÜYÜK VERİ KAVRAMI



# BÜYÜK VERİ NEDİR?

Büyük veri  
Excel'i krah  
edecek  
herşey

Küçük veri RAM'a  
sokulabilir, Büyük  
Veri bilgisayarı  
kapattır çünkü  
RAM yetmiyor



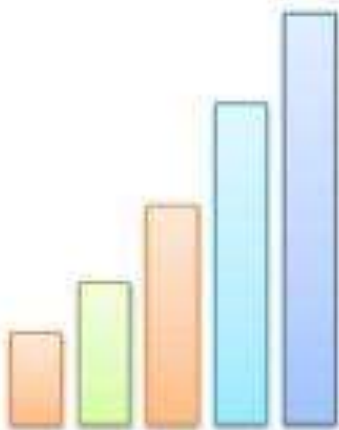
Diğer deęişle, Büyük Veri eski,  
alıştığımız metotlarla işletebilmek  
için fazla miktarda olan veriler  
demektir

[https://en.wikipedia.org/wiki/Pointy-haired\\_Boss](https://en.wikipedia.org/wiki/Pointy-haired_Boss)

# BÜYÜK VERİ NEDİR?

- Zamanımızda veri çok yüksek hızla üretilip toplanmaktadır
  - İnternette ziyaretçilerin tıklamaları
  - Alışveriş ödemeleri
  - Sensör kayıtları
  - Güvenlik kameraların kayıtları
  - GPS verileri
  - Sosyal media etkileşimleri
  - ...
- Bu tip verileri kaydedip işletmek ve analiz etmek gerçekten zor olmaya başlamıştır

# WWW'DEN VVV'E



**VOLUME**



**VELOCITY**



**VARIETY**

# WWW'DEN VVV'E

## ○ Volume (Hacim)

- Bugünkü iş/kurum/şirketlerin verileri inanılmaz miktardadır

## ○ Variety (Çeşitliği)

- İş/kurum/şirketlerin verilerinin karmaşıklığı artmaktadır
- Her gün yeni tür veriler toplanmaya başlamaktadır

## ○ Velocity (Hız)

- Verilerin toplama hızı artmaktadır
- Baze verilerin doğasına göre anlık işlenmesi ve tepki zorunludur – bu tür verilere “veri akışları” (data streams) denir



# BÜYÜK VERİNİN 4ÜNCÜ V

## ○ Veracity (Kalitesi)

- Veri toplamanın hızı artmakla beraber, verilerin kalitesi ve güvenilirlik düşmektedir
- Veriler sağlayan kaynaklardan yanlış, eksik, bozuk biçimde veriler gelebilmektedir



# BEKLENTİ

- Bugünkü sistemler, şirketler ve kurumlar Terabyte'ten Petabyte'e kadar rutin olarak bilgi üretmekte dir
- Bilgi, şirket/kurumun başarısı için büyük önem taşımakta dır
- İyi kararlar verebilmek için anlamlı verilerin var olması ve dikkate alınması şarttır

## Volume

- Terabyte
- Veri kayıtları
- İşlem kayıtları
- Tablolar
- Dosyalar
- Loglar

## Büyük Verinin Üç V

## Variety

- Yapılandırılmış
- Yapılandırılmamış
- Yarı-yapılandırılmış
- Karışık

## Velocity

- Batch veri analizi
- Neredeyse gerçek zamanlı veri
- Gerçek zamanlı veri
- Veri akışları

## ÇEŞİTLİĞİ

- XML dosyaları → yarı yapılandırılmış veri
- Word, PDF, TXT dosyaları → yapılandırılmamış veri
- Email metinleri → yapılandırılmamış
- Veri tabanları ve Excel tabloları → yapılandırılmış
- Sosyal media mesajları → yapılandırılmamış
- İşlem logları → yarı yapılandırılmış



# BÜYÜK VERİLERDEN - BÜYÜK BEKLENTİLER !



# KULLANIM ALANLARI

## ○ İnternet ve e-komerse

- Tavsiye (recommender) motorları
- Reklam hedeflenmesi
- Arama kalitesi artırma
- Yasal olmayan faaliyetlerin keşfetme

## ○ Telekom

- Müşterilerin memnuniyeti takip etme
- Telekom ağın performansı takip etme
- Telekom ağları optimizasyon
- Müşterilerin davranışları analizi
- Ağın çalışması analizi ve arıza durumlarının tahmin etme

# KULLANIM ALANLARI

## ○ Devlet

- Siber güvenlik
- Emniyet ve yasal araştırmaları
- Nüfus bilgi toplama ve analizi
- Ekonomik bilgi toplama ve analizi
- Diğer

## ○ Tıp ve medikal

- Tıp kayıtların işlenmesi
- Genetik araştırmaları
- Hizmet iyileştirilmesi
- İlaç/tedavi güvenliği araştırmaları

## KULLANIM ALANLARI

- Banka ve finans
  - Risk modellenmesi
  - Tehdit analizi
  - Dolandırıcılık keşfetme
  - Kredi skorlar
- Perakende
  - POS kayıtların analizi
  - Müşteri memnuniyet takibi
  - İmaj araştırmaları ve analizi



# SPESİFİK ÖRNEKLERİ: SPOR



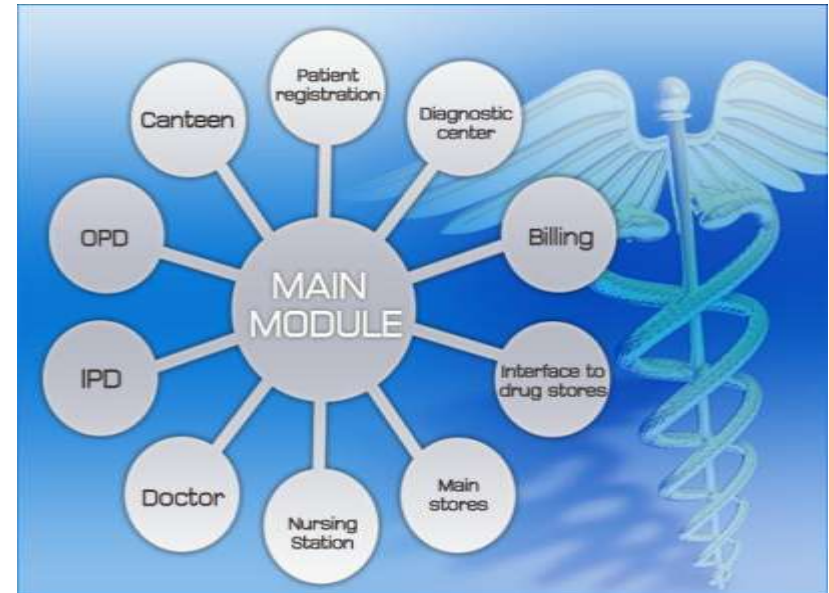
# SPESİFİK ÖRNEKLERİ: SPOR

- Büyük sporda büyük veri yaklaşımı bilet satışı, pazarlama ve reklam stratejilerinin geliştirilmesi için kullanılır,
- Bunun için sosyal media kullanılarak reklam kampanyaların verimliliği ve alması gereken yönler seçiliyor
- Spor takımları “büyük veri” modelleri oyun strateji, oyuncular seçimi vb konular için de kullanılır



## SPESİFİK ÖRNEKLERİ: Tıp

- Tıpta hastaların kayıtları analiz edilerek tedavi sonuçları ve daha iyi müdahale yapmak için yöntemler araştırılmakta
- Genel test sonuçları kullanılan teşhis tespiti yöntemleri de geliştirilmekte



## SPESİFİK ÖRNEKLERİ: E-KOMERS

- Online satıcılar inanılmaz miktarda kullanıcı ve ürün hakkında bilgilere sahiptir
- Bu bilgi, kullanıcıların davranışları analiz etmek, hedeflenen reklam üretmek ve alışverişteyken daha faydalı ürün önerileri yapmak için kullanılmakta





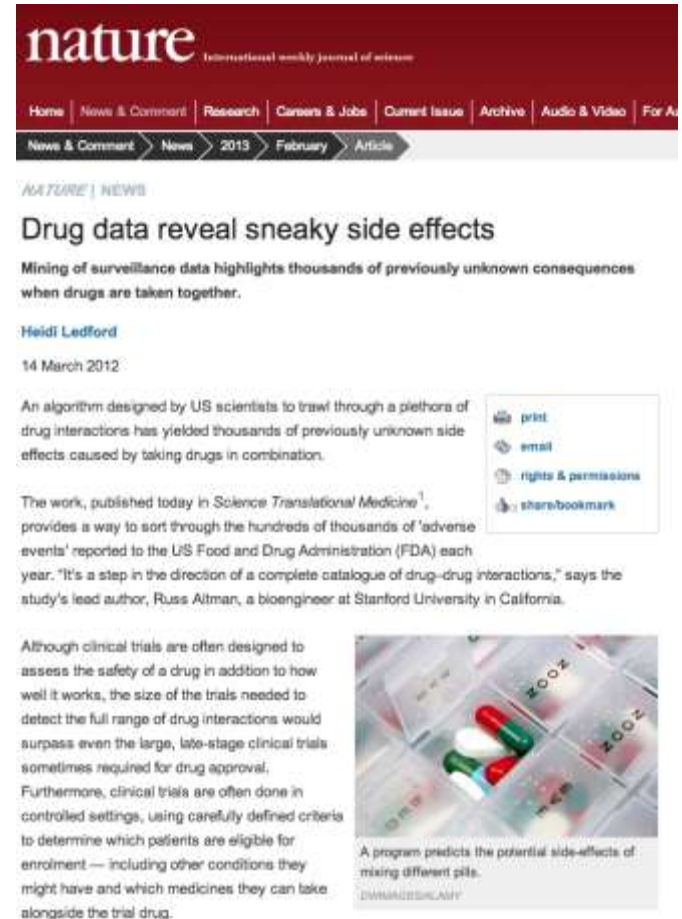
## SPESİFİK ÖRNEKLERİ: ONLINE

- Netflix online film seyretme hizmeti dir
- Kullanıcının daha önce seyrettiği film tarihçesi kullanılarak, kullanıcının ilgisi özel olarak çekebilecek film önerileri makine öğrenme yöntemleri yardımıyla Netflix'in seçtiği ünlü dür
- Netflix'in film deposu da 1Petabyte geçtiği bilinmekte



# SPESİFİK ÖRNEKLERİ: İLAÇ

- Tıpta ilaç ve tedavi güvenliği ve verimliği araştırmaları yeni değildir
- Bu araştırmalar, büyük veri kavramı temel olarak oluşturan veri modelleme ve veri analizi yaklaşımları yardımıyla gerçekleştirilmekte



# SPESİFİK ÖRNEKLERİ: ALIŞVERİŞ

- Target, Amerikadaki bir alışveriş zinciri, müşterilerinin alış tarihçesine bakarak hedeflenen reklam üretmeye çalışmakta
- 2012'de bu yöntemler kullanılarak aylesi bilmeyen bir genç kızın hamile olduğunu ve yaklaşık doğum tarihi tespit ettiği ile haberlere çıkmıştır





# SPESİFİK ÖRNEKLERİ: TARIM

- Tarım ve hayvan bilimi, tarımda kullanılan bitki ve hayvanların genetik bilgileri ve üreme tarihçeleri gibi veri analiz ederek, performansı artırmak için yapay seleksiyon planları modellenir ve kullanılmakta



# DAHA ÇOK ÖRNEK

## ○ Spor

- Basketbol oyunların planlanmasında veri analiz oldukça kullanılır
- Futbolda benzer eğilim görünmektedir

## ○ Eğlence sektöründe

- Bilgisayar oyunları müşteri veri analizine büyük önem vermektedir
- Yapılacak filmlerin seçiminde benzer eğilim vardır

## ○ Finans

- Viza otomatik ödeme bilgilerinin analizi ile dolandırıcılık keşfetme sistemleri geliştirmiştir

## ○ Google ve Facebook

- Kullanıcıların bilgilerinde veri madenciliği yaparak reklam ve benzer faaliyetleri hedefliyor

# DAHA DAHA ÇOK

- Tarım başkanlığı
  - Tarım şirket ve çiftliklerin verileri
  - Doğum, ölüm, taşınma, tedaviler, örnekler gibi verileri
- Enerji üretim
  - Elektrik enerji üretimi ve tüketim, en uygun dağıtım şekli, dinamik elektrik fiyatları, arza ihtimali, müşterin sayaç kurcalama
- Petrol ve madencilik
  - Jeolojik veriler, işlem veriler, lojistik, mühendislik
- Perakendeciler
  - Müşteri modellenmesi, önceki Target örneğine bakın
  - Satışlar ve hava, sezon vb durumlarla ilişkileri, lojistik ve stok yönetimi

## ORTAK DESEN ...

### Genel biçimde olan çeşitli veri serileri

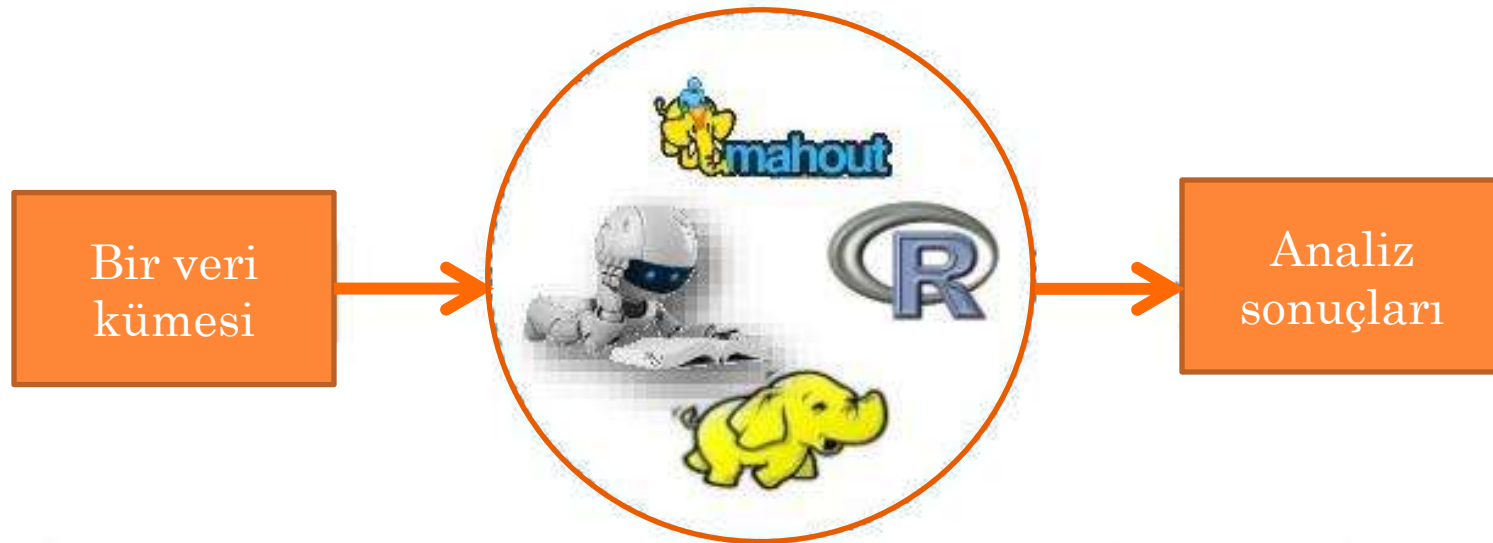
- A alışveriş kayıtları
- Üreme tarihçesi
- Süt/et üretimleri
- Çeşitli tıp test sonuçları
- Sosyal medya postları
- Doktora ziyaret kayıtları
- Haber parçaları
- ...



### İş kararı



## ORTAK DESEN ...



# NASIL

- İstatistik
- Olasılık
- Diferensiyal
- Linear cebir
- Algoritmalar
- Programlama
- ...

Ben bunu  
bildiysem ben de  
yapardım ...



# BÜYÜK VERİ FELSEFESİ

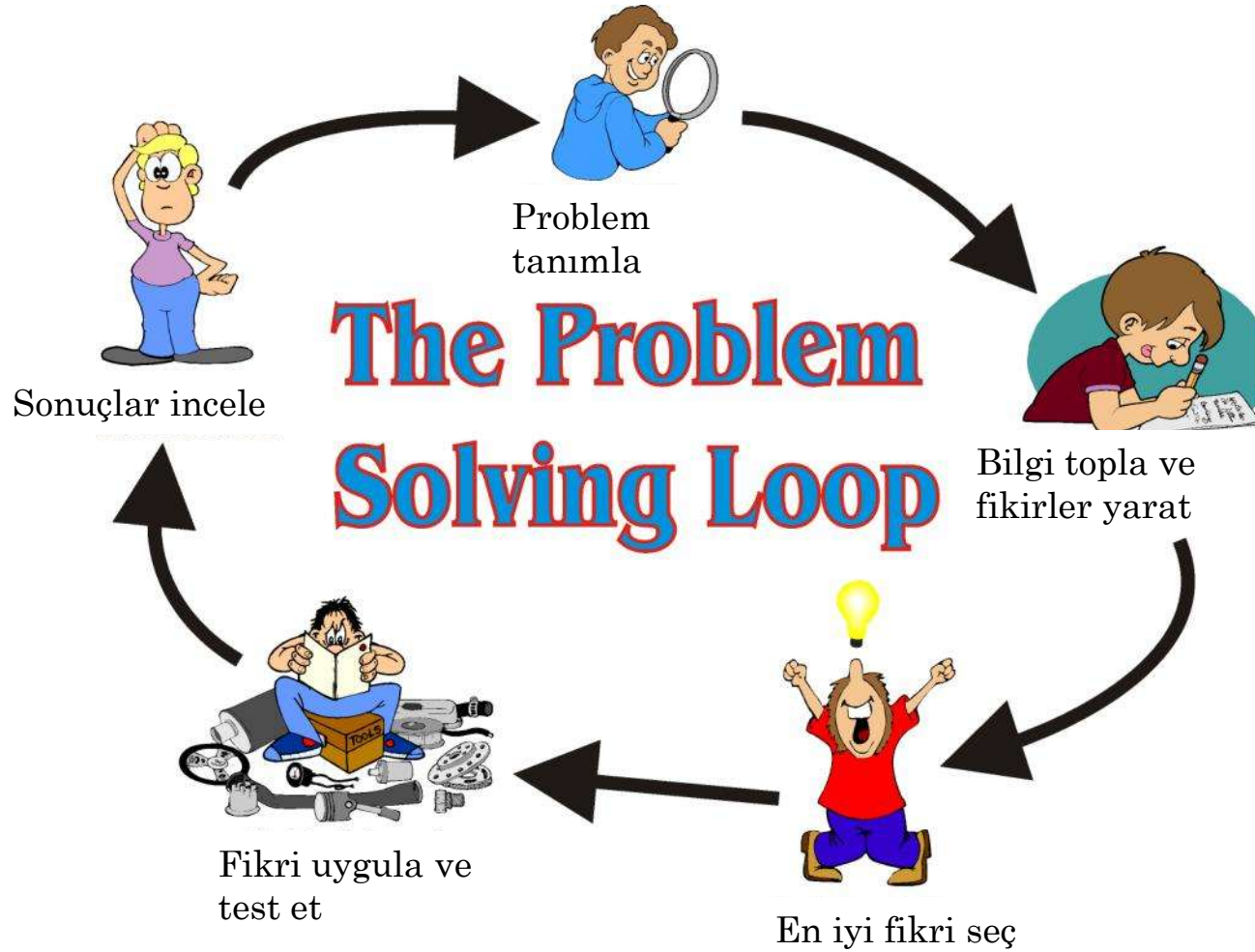




# İKİ PROBLEM ÇÖZME YAKLAŞIMI

- Alıştığımız problem çözme
  - Problemin mantığını anlamak
  - Mevcut kavramlar arasındaki ilişkiler kesinleştirmek
  - Mevcut olan faktörlerin muhtemel etkileri belirtmek
  - Belirli müdahale olduğunda problemde değişiklikleri tahmin etmek

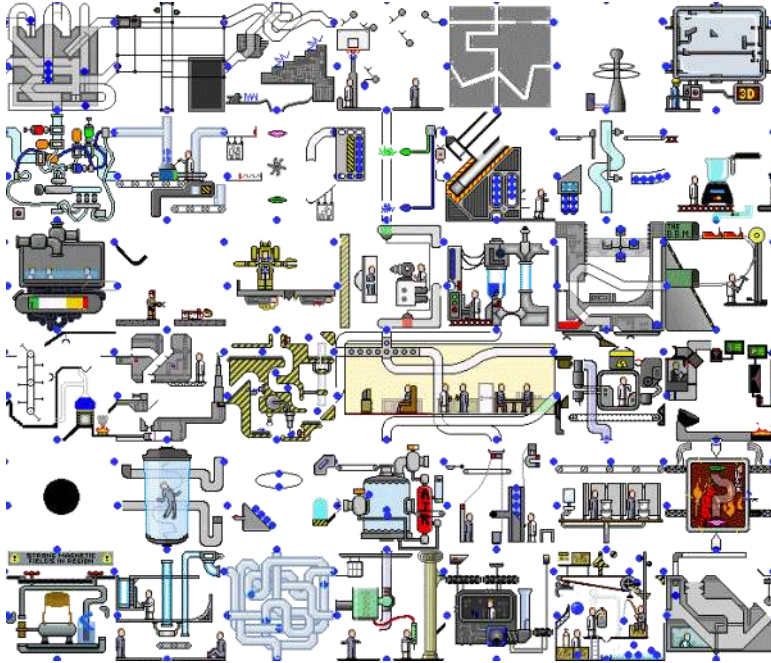
# ESKİ PROBLEM ÇÖZME YAKLAŞIMI (NORMAL)



## Bugünkü iş yönetiminde bu yaklaşımı uygulamak her gün daha zor oluyor

- Veri miktarı artmakta
- Dikkate alınabilir ilişki sayısı artmakta
- İlişkilere mevcut olan eleman sayısı artmakta
- Elemanların doğası ve biçimi karmaşık ve açık değil
- Model ilişkisi çok, doğası belirsiz ve karmaşık

# PROBLEM ÇÖZME



Bunu nasıl  
yaptığımız  
hoşuma  
gitmedi



Sen ne yapardın? Bunu başka bir  
şekilde  
yapabilirdik

Bu şeyi  
tamamen  
unutsaydık,  
nasıl olur?

Onu böyle  
yapalım

# ALTERNATIF (MODERN) YAKLAŞIM

- Yeni problem çözmeye yaklaşım
  - Belirli problemi için genel bir parametrelili model oluştur
  - Bu modelin parametreleri mevcut verilerden tahmin et
  - Müdahalenin sonucunu bu modelden tahmin et

# ALTERNATIF PROBLEM ÇÖZME YAKLAŞIMI (YENİ)

Gerçek problem ve  
dahil olacak  
değişkenler tanımla

$a, c, X, Y, \xi, \mu, \dots$

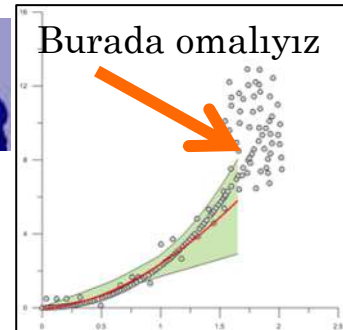
Bir genel parametrelili  
model seç

Parametreleri tahmin  
et ve model kullan

$$Y = a_1 X_1 + a_2 X_2 + a_3 X_3 + \dots$$



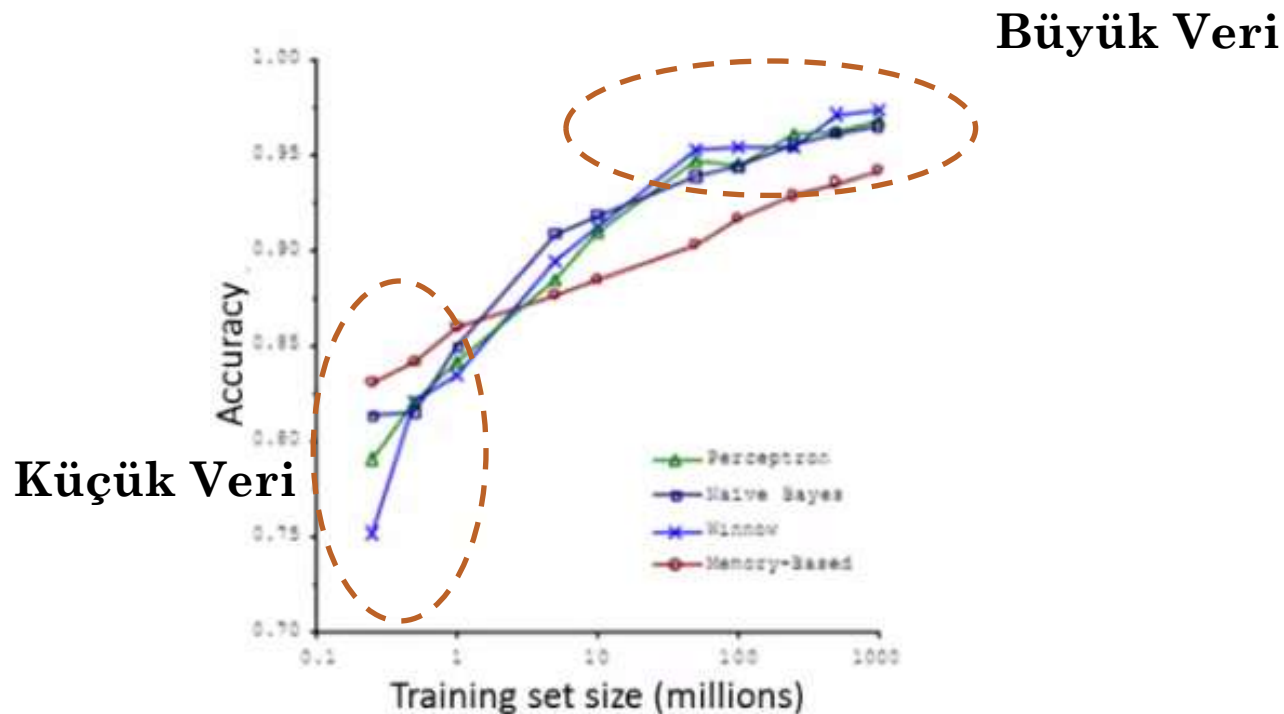
Burada omalıyız





# BÜYÜK VERİ “TEOREMI”

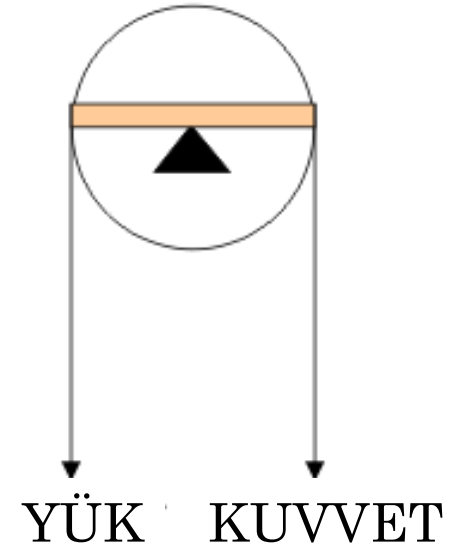
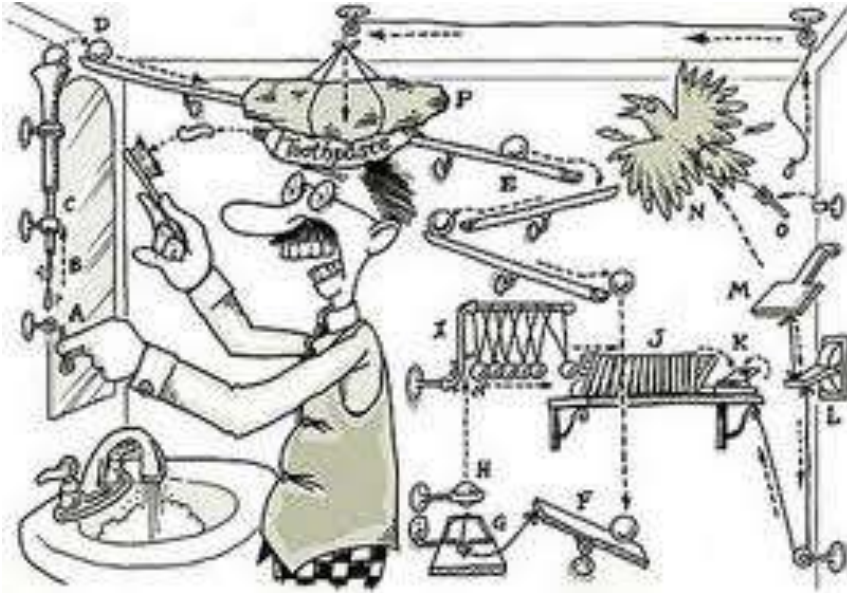
- Veri miktarı büyük olduğunda, basit modellerin performansı ve çok karmaşık model performansı arasında büyük fark yok





# BÜYÜK VERİ “TEOREMI”

Büyük veri rejiminde genel basit modeller, probleme özel tasarlanmış ve çok karmaşık olan modeller'den genelde daha başarılıdır

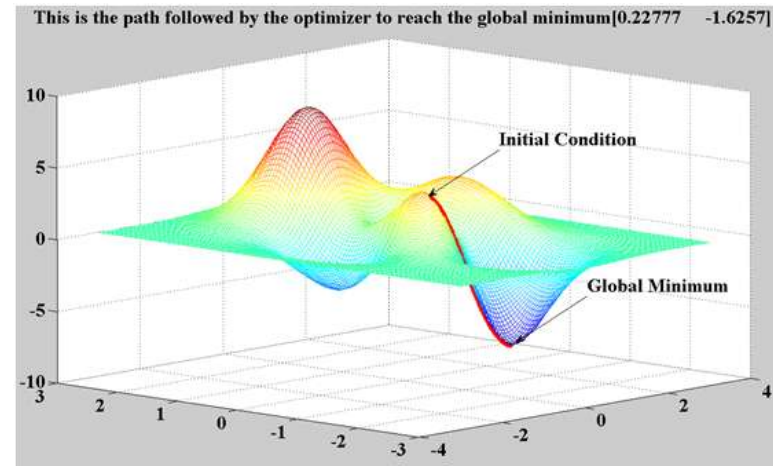


## BÜYÜK VERİ TEOREMINDEN UYGULAMALARINA

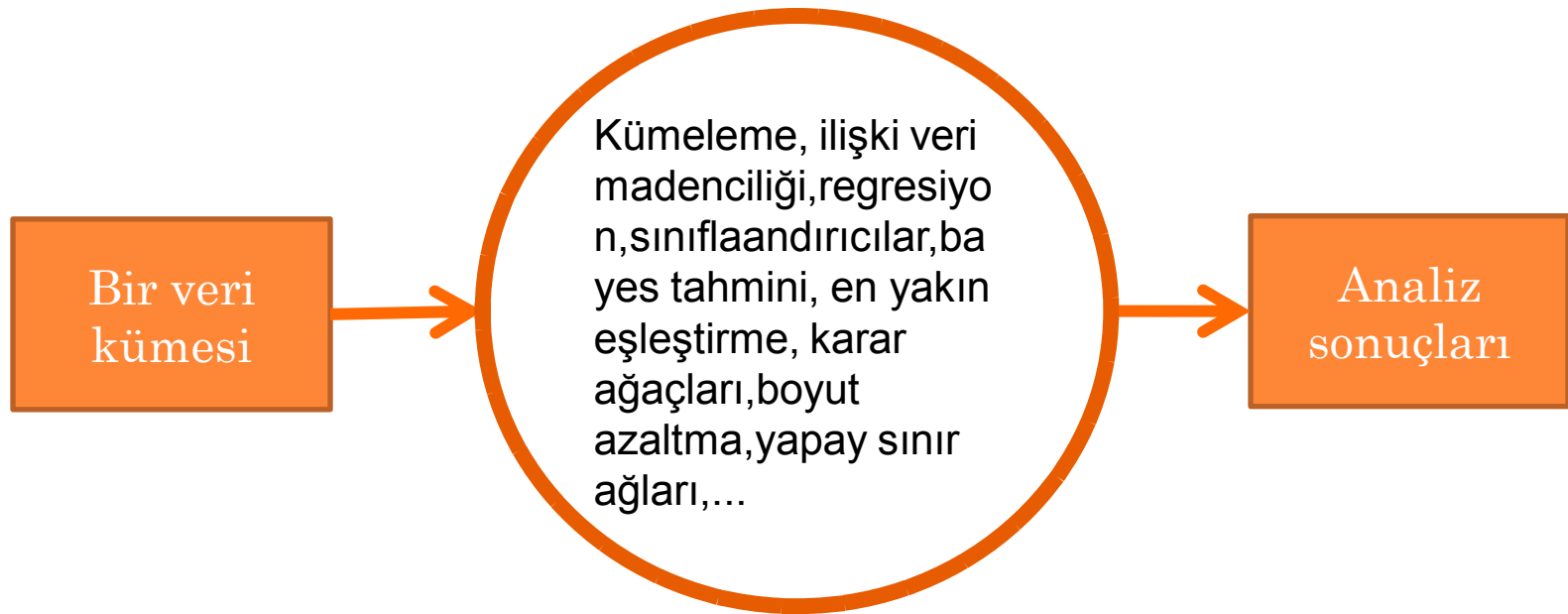
- Teknolojinin gelişimiyle mevcut olmaya başlayan büyük miktarda olan iş yönetimi ile alakalı verilerden faydalanmak için ...
- İş yönetimi ile ilgili soruların modellenmesi için genel makine öğrenme aletleri kullanılır ve ...
- Bu modellerin sonuçları iş yönetim kararları vermek için kullanılır

# BÜYÜK VERİ İÇİN MAKİNE ÖĞRENME ALETLERİ

- Kümeleme
- Regresiyon
- İlişkisel kural madenciliği
- Sınıflandırma
- Bayes tahmini
- En yakın eşleştirme
- Karar ağaçları
- Boyut azaltma
- Yapay sınır ağları



# BÜYÜK VERİ FELSEFESİ



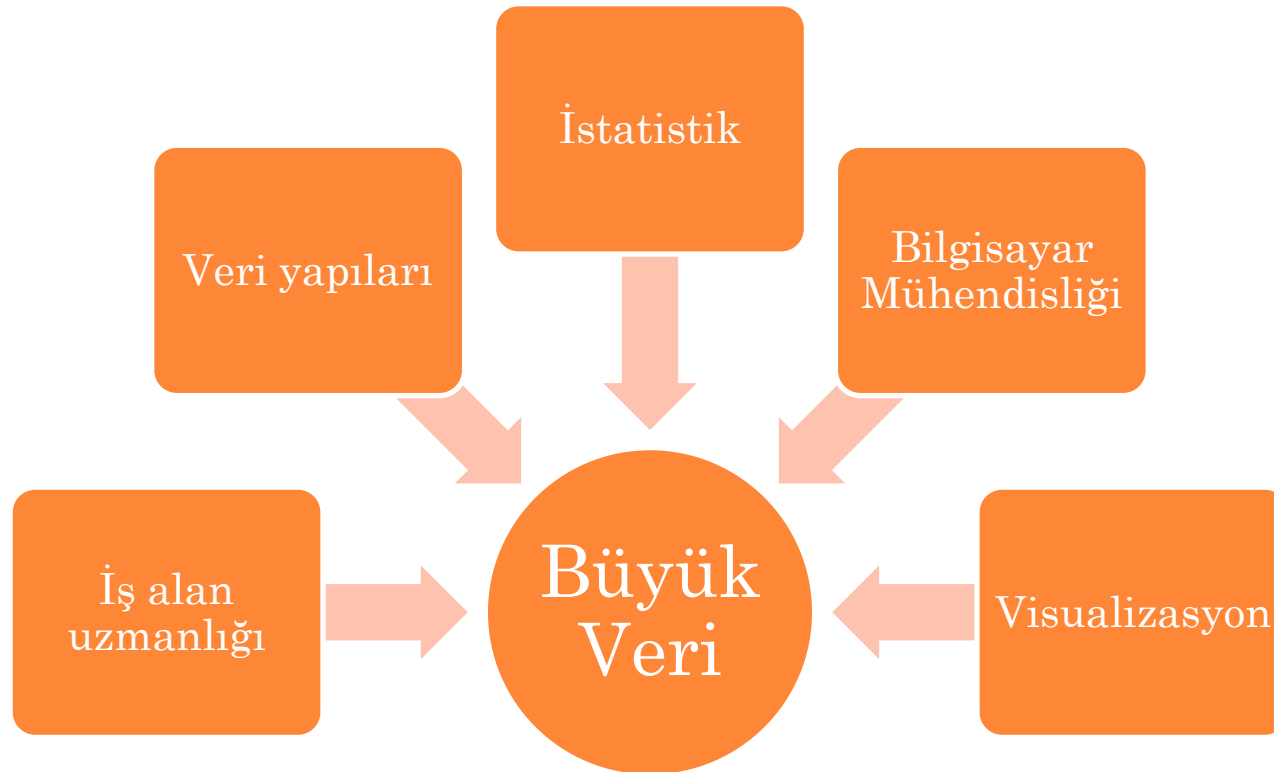
# BÜYÜK VERİ İÇERİSİ



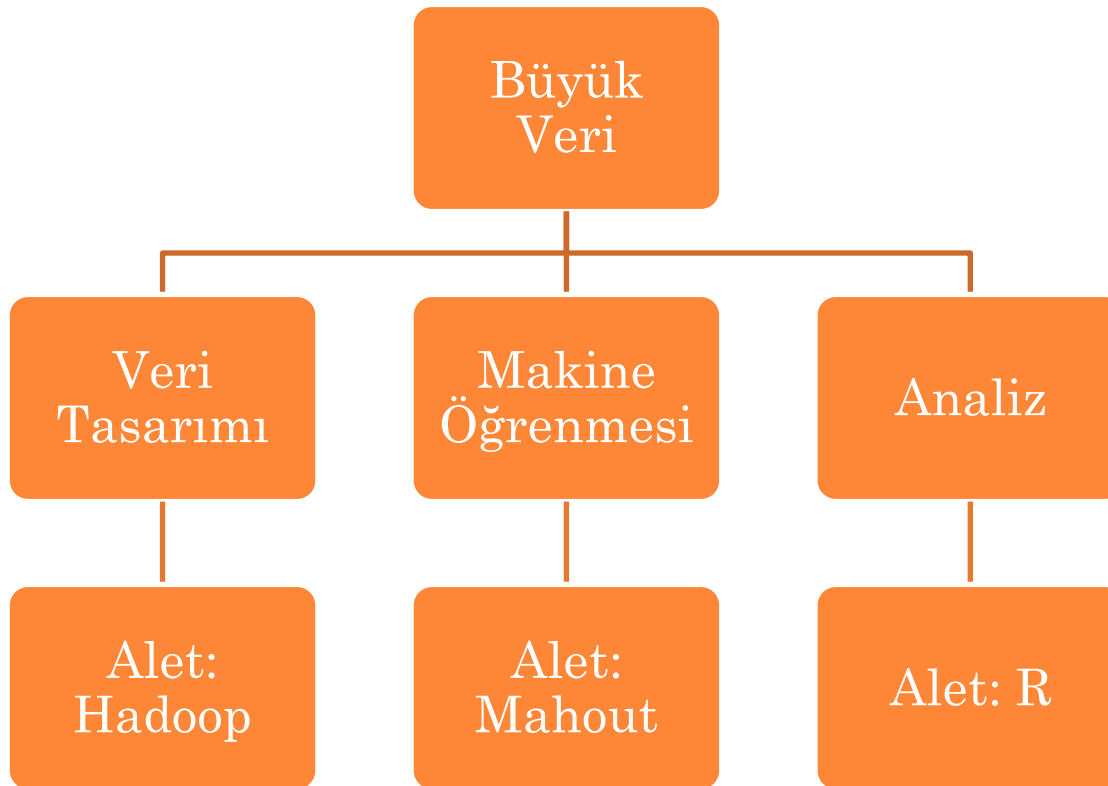
# BÜYÜK VERİ'NİN SORUNLARI

- Veri toplama
  - hangi veriler kullanılacak?
- Veri depolama
  - TerraByte/PettaByte veri nerede kaydedilecek?
- Veri transferleri
  - TB/PB veri nasıl transfer edilecek?
- Veri sorgulanması
  - TB/PB veri setleri nasıl sorgulanacak?
- Veri analizi ve sonuç çıkartma
  - Karmaşık ilişkiler nasıl ortaya çıkartılacak?
- Sonuç bilgilendirme
  - Karmaşık sonuçlar nasıl incelenip bilgilendirecek?

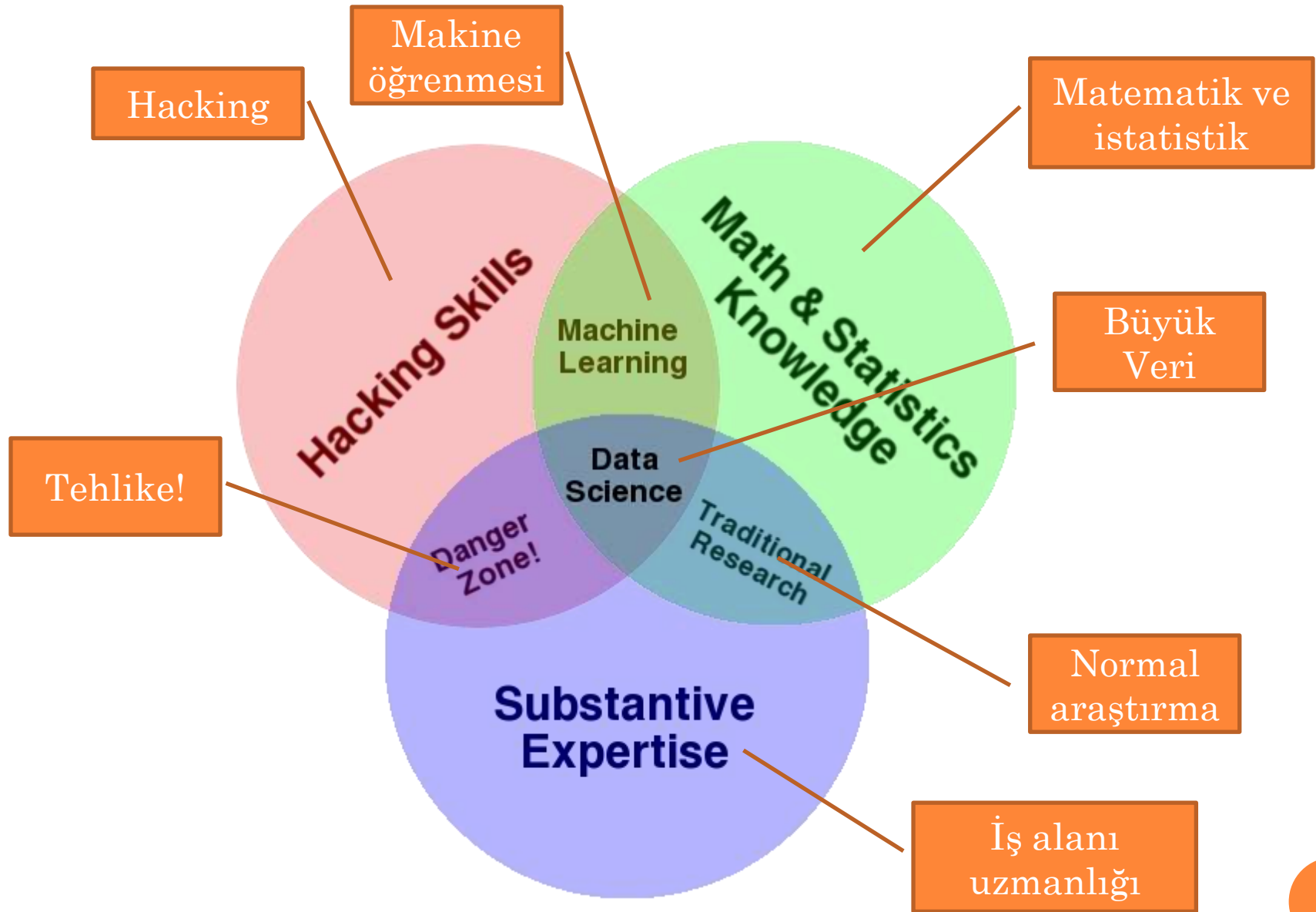
# BÜYÜK VERİ BİLEŞENLERİ



# BÜYÜK VERİ AŞAMALARI



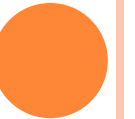




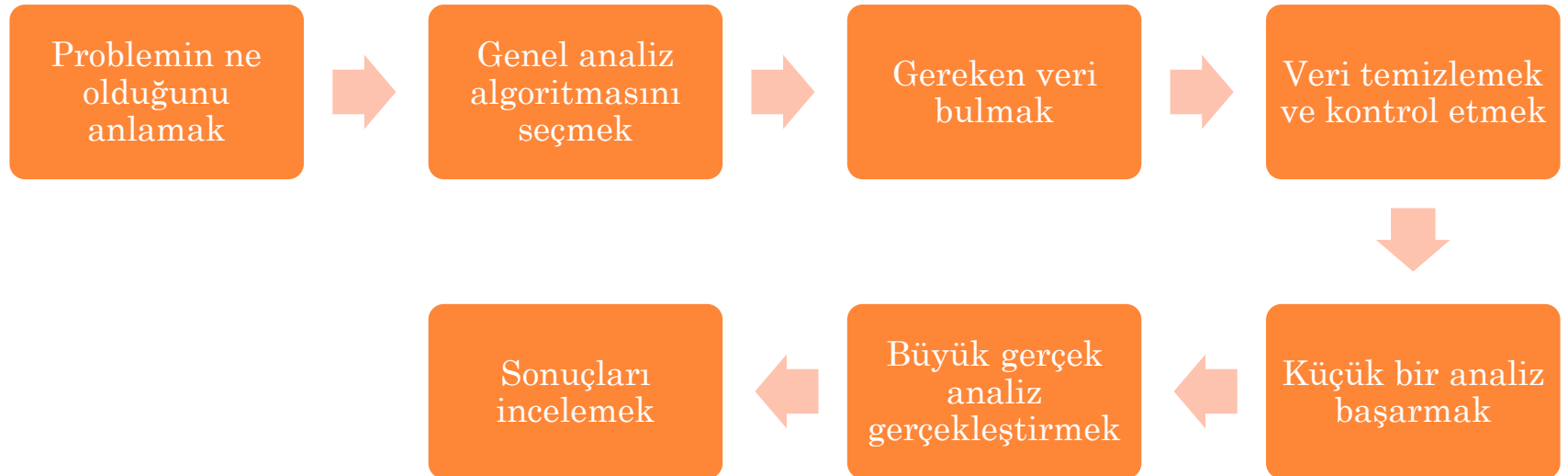
# BÜYÜK VERİ'NİN İKİ TARAFI

- Veri analizi/makine öğrenmesi/istatistik
  - Verilerden kavrama yaratabilmek
- Büyük ölçekli bilgi işlem/bilgisayar programlama
  - Verilerin büyük hacimleri işletebilmek
- Son yıllarda Google Prediction API, Microsoft Azura ML, Amazon ML, BigML gibi Big Veri bulut çözümleri bu ikisi işi oldukça kolaylaştırmıştı

# DETAYLI ÖRNEKLER



# BÜYÜK VERİ UYGULAMALARI

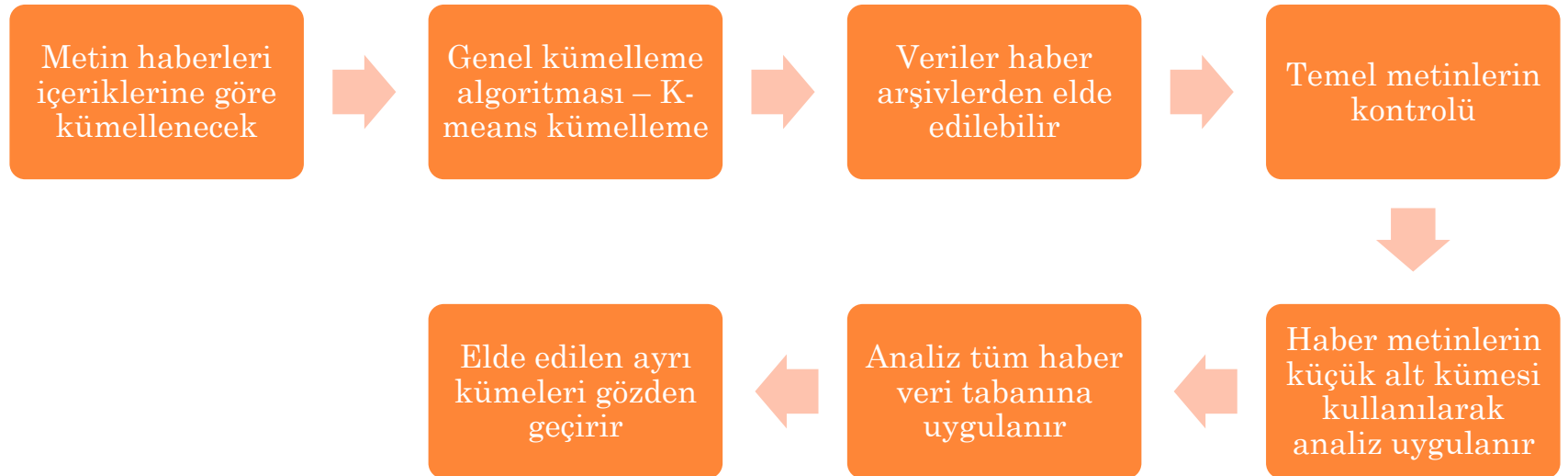


# MEDIA KULLANIM ÖRNEĞİ

## ○ Problem tanımı:

- Bir büyük media şirketi, 1980 den itibaren belirli konuda popüler haberleri incelemek istiyor.
- Analizi kolaylaştırmak için, tüm haberleri benzer kümelerine ayırdırmak istiyor

# MEDIA KULLANIM ÖRNEĞİ



# ALİŞVERİŞ KULLANIM ÖRNEĞİ

## ○ Problem tanımı

- Bir magaza müşterilerinin davranışı ve satın alınları anlamak istiyor. Bu bilgi müşterilerine daha iyi hizmet sağlamak için kullanılacak.
- Analiz magazaya, müşterilerin bir ürünle beraber başka ürünlerin satın aldığını ortaya çıkabilir; bu halde söz konusu ürüne uygulanacak kampanya ilişkili ürünlerin satışı artırıp magazanın geliri etkileyebilir



# ALİŞVERİŞ KULLANIM ÖRNEĞİ



# Tıp Kullanım Örneği

## ○ Problem tanımı

- Bir hastane hastaların demografik özellikleri ile beraber çeşitli test sonuçları ve hastalığın başlangıcı arasında muhtemelen ilişkiyi ortaya çıkartmak istiyor
- Bunları kullanılarak tahmini müdahale stratejileri tasarımlamak istiyor

# Tıp Kullanım Örneği



# SOSİYAL MEDIA KULLANIM ÖRNEĞİ

## ○ Problem tanımı

- Sosyal media araştırma şirketi Facebook'taki “sıcak” konuları analiz ederek onları aşağıdaki kategorilerine sınıflandırmak istiyor:
  - Giyisi (ayakabı, kıyafet, saat, takı, ...)
  - Sanat (Kitaplar, film, DVD, VCD, müzik)
  - Fotokameralar
  - Olaylar (seyahat, konser, film,...)
  - Sağlık (güzellik, spa,...)
  - Ev (mutfak, mobilya, bahçe,...)
  - Teknoloji (bilgisayar, laptop, tablet, smartphone,...)

# SOSİYAL MEDIA KULLANIM ÖRNEĞİ



# ALETLER



# GOOGLE PREDICTION API





# DEMO

# AZURA ML

## Azure Machine Learning Service

Data -> Predictive model -> Operational web API in minutes

Data



Blobs and Tables  
Hadoop (HDInsight)  
Relational DB (Azure SQL DB)



Integrated development environment  
for Machine Learning



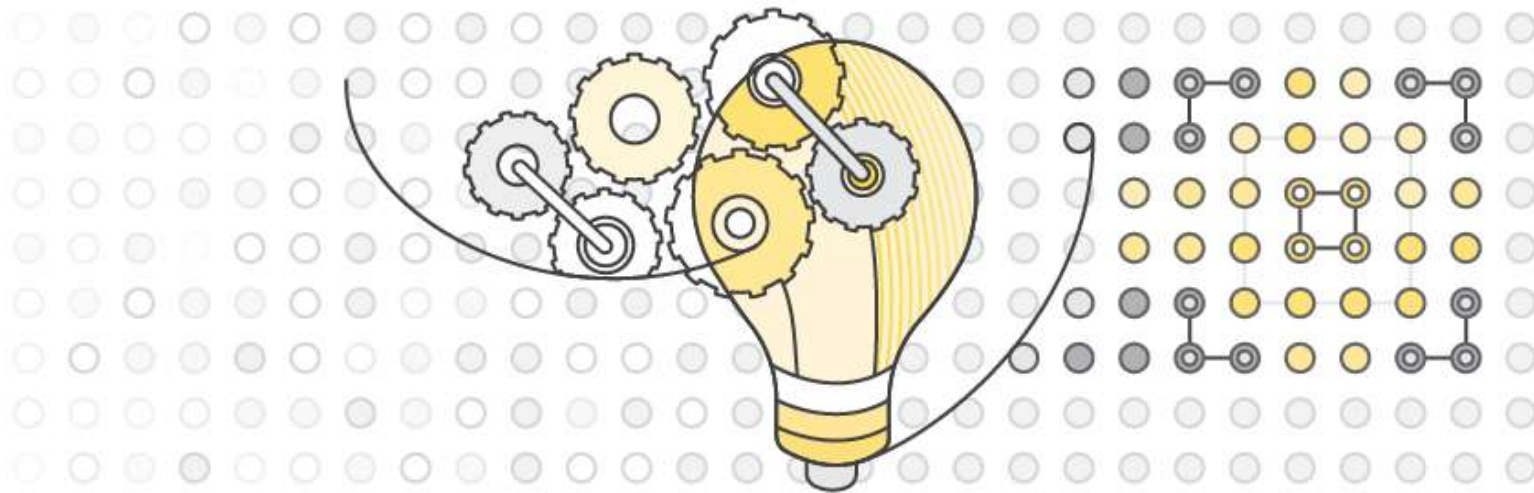
Model is now a web  
service that is callable



Clients



# AMAZON ML



# REAL STATISTICS EXCEL

- “İleri istatistik MS Excel’inize...”



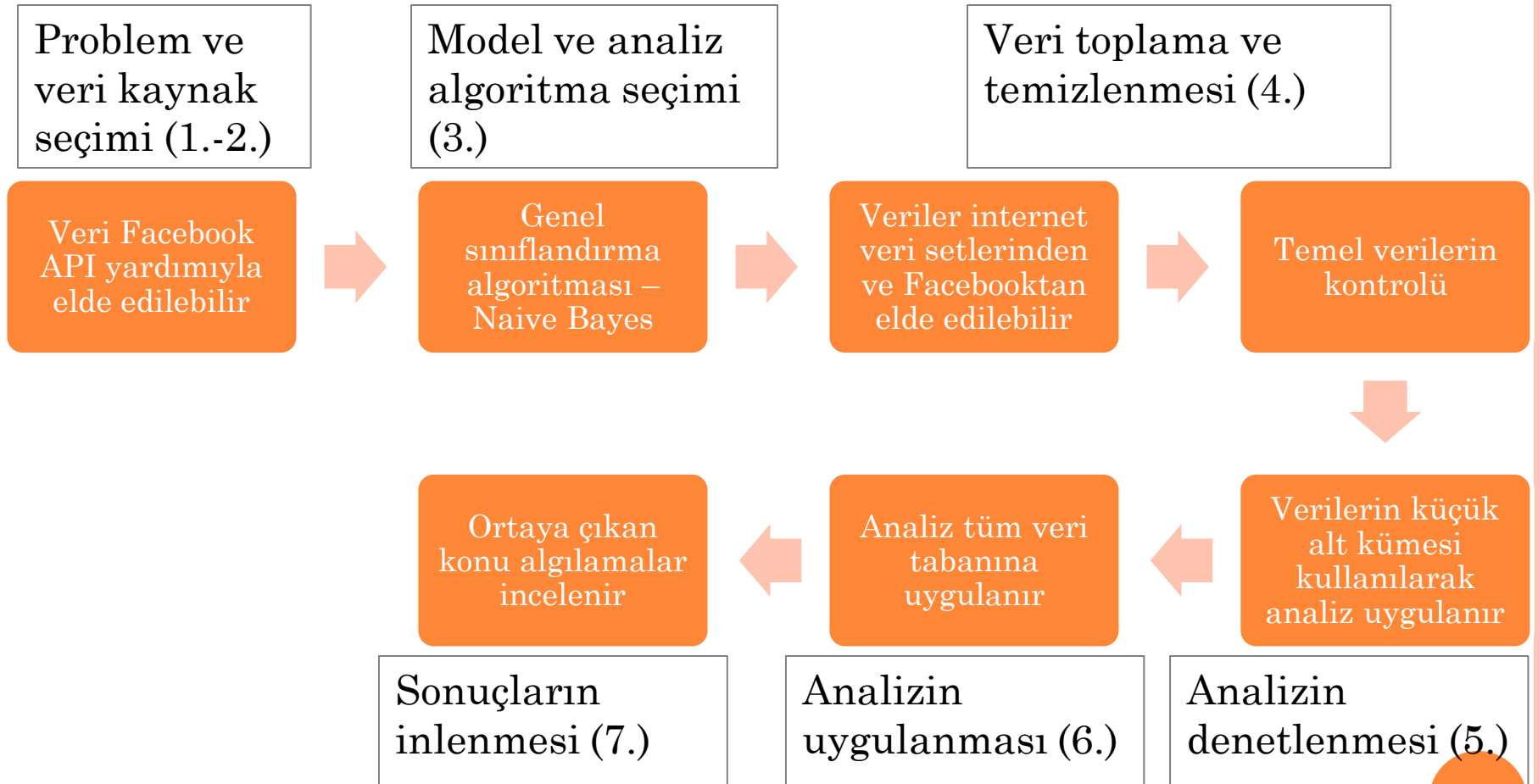
## DAHA ...

- Open Source kaynaklar: Hadoop, Apache Mahout, MapReduce, TensorFlow Google, H2O on Hadoop, Spark MLib, Weka-Java, SHOGUN, WSO2 Data Analytics, DARWIN, scikit, FAIR Facebook, vb.
- Daha ileri ve düşük seviyeli ARGE için

# BÜYÜK VERİ UYGULAMA TARIFI

1. Mevcut iş sorunları ve veri kaynakları inceleyip özetleyin
2. Uygun bir sorun modelleme problemi olarak tanımlayın ve mevcut veri kaynaklarından ilgili olabilecek kaynakları belirtin
3. Sorunu ve veri bağlayacak bir genel makine öğrenme modeli ve analiz algoritmasını seçin
4. Gereken veri toplayın ve kalite açısından inceleyin, gerekirse --- veri temizleyin
5. İstedığınız analiz işleminin adımları seçip, küçük bir denetleme veri kümesi üzerinde uygulayın
6. Analizinizi tüm verilerinize uygulayın
7. Analiz sonuçları uygun görsel şekilde gösterin ve orijinal soruna dair sonuçlar çıkartın

# SOSİYAL MEDIA ÖRNEĞİ





# BÜYÜK VERİ UYGULAMA TARIFI

1. Mevcut iş sorunları ve veri kaynakları inceleyip özetleyin
2. Uygun bir sorun modelleme problemi olarak tanımlayın ve mevcut veri kaynaklarından ilgili olabilecek kaynakları belirtin
3. Sorunu ve veri bağlayacak bir genel makine öğrenme modeli ve analiz algoritmasını seçin
4. Gereken veri toplayın ve kalite açısından inceleyin, gerekirse --- veri temizleyin
5. İsteddiğiniz analiz işleminin adımları seçip, küçük bir denetleme veri kümesi üzerinde uygulayın
6. Analizinizi tüm verilerinize uygulayın
7. Analiz sonuçları uygun görsel şekilde gösterin ve orijinal soruna dair sonuçlar çıkartın

# ÖDEV/DÖNEM PROJE ÖN-ÇALIŞMASI

- **Öbür derse kadar lütfen, iş alanınızdan bir problemi seçin**
  - Mesleki alanınızdan
  - Kolay ulaşılabilir veri (elektronik halde)
  - İlginç bir soru
- **10 dakika geçmeyen, derste sunulmak üzere bir tanıtım sunuşu hazırlayın**

# ÖDEV/DÖNEM PROJE ÖN-ÇALIŞMASI

## ○ **Dönem projesi**

- Problemin tanımı
- Modelin tanımı
- Veri elde edilmesi
- Bulut araçları yardımıyla analiz edilmesi
- Sonuçların incelenmesi