

Overview of Machine Learning Methods and Applications

Yapay Öğrenme Metotları ve Uygulamaları

A. Taylan Cemgil, Bilgisayar Müh.

09.2017, İstanbul



<http://www.cmpe.boun.edu.tr/pilab>

Konu Başlıkları

- Büyük Veri - Yapay Öğrenme
- Kullanım örnekleri (Use Cases)
- Güdümlü Öğrenme (Supervised Learning)
 - Sınıflandırma (Classification)
- Güdümzsüz Öğrenme (Unsupervised Learning)
 - Öbekleme (Clustering)
 - Boyut indirgeme (Dimensionality Reduction)
- Büyük verilere uygulama (Scaling)
 - Mimariler/Büyük veri araçları
- Referanslar
- Sonuç

Big Data? (Büyük Veri)



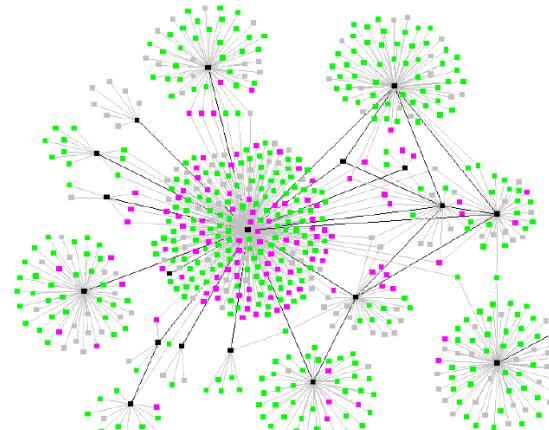
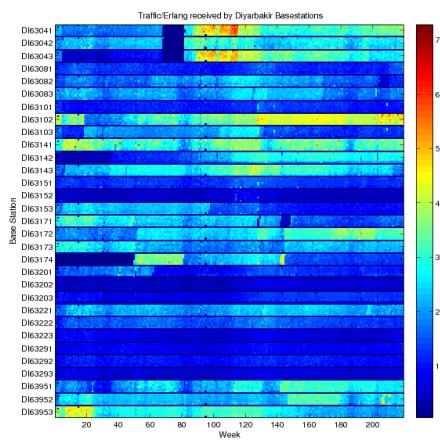
- David Hilbert (1900) – 20. yy Matematik için Sorular
- David Donoho (2000) – Analysis of Large Datasets
- Funding, Publicity → ~ 2010 Big Data

Büyük Veri

- Yazılım/Donanım Altyapısı (SW/HW Infrastructure)
- Görselleştirme/Etkileşim (Visualization/Interaction)
- Data Analytics ← Yapay Öğrenme

Yapay Öğrenme

- Hesaplama tabanlı metodlar topluluğu
 - Gizli örüntü (pattern) yakalama
 - Öngörüler üretmek
 - Belirsizlik altında karar desteği
 - Ham veriyi faydalı bilgiye dönüştürmek



Yapay Öğrenme, Veri Madenciliği, İstatistik



Yapay öğrenme ve Büyük veriler. Gerçekten yeni mi?

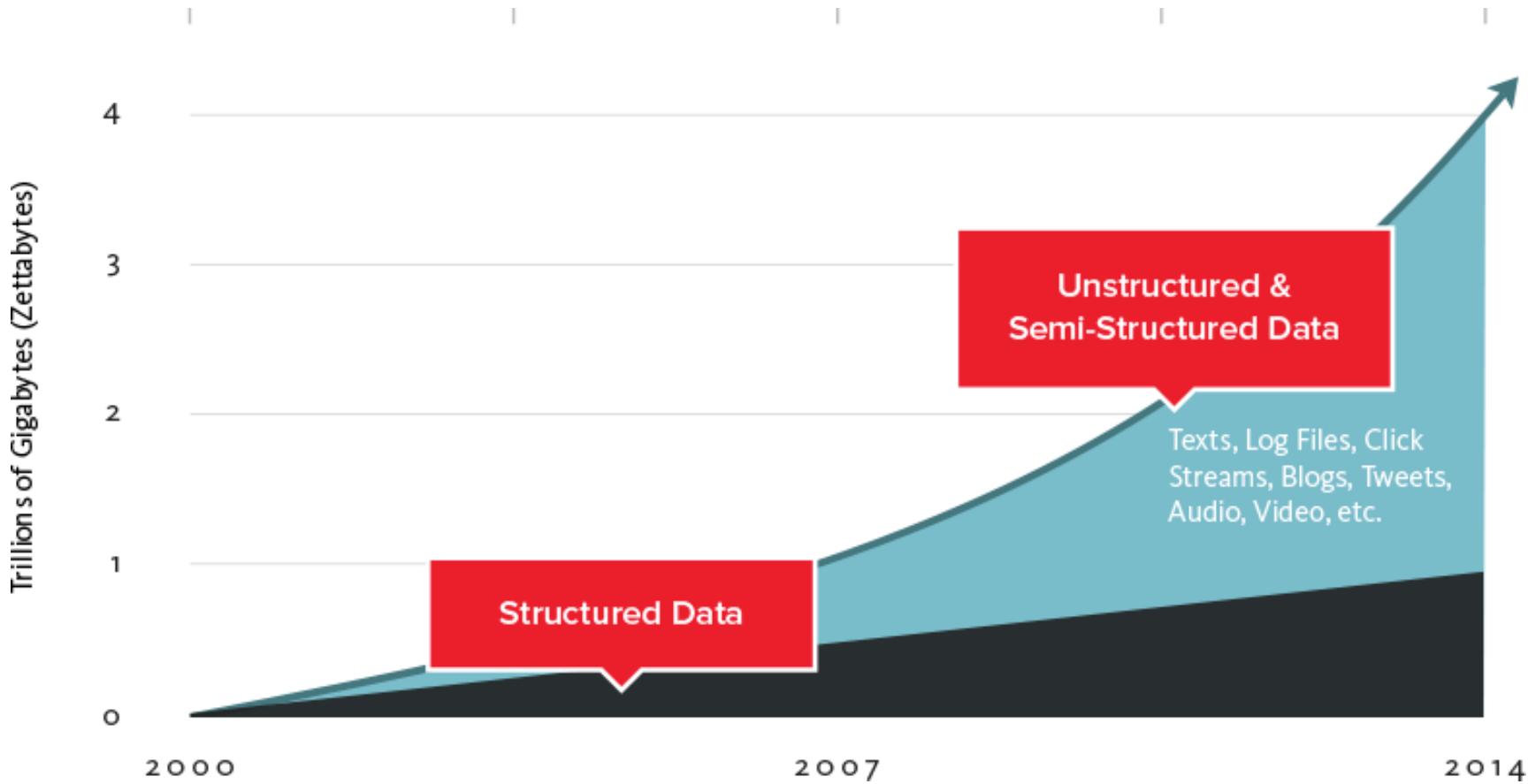
- Eski metotlara yeni bir bakış
- ... ve yenilerinin geliştirilmesi
- Büyük boyutun Laneti/Nimeti ‘Curse/Blessing of Dimensionality’
- Ucuzlayan Altyapı
 - Bulut Bilişim
 - Sensör Ağları, Nesnelerin İnterneti (“yeni veri”)
 - Hız (“gerçek zaman”)

Büyük Veri?

- Pragmatik Bakış Açısı
 - Küçük Veri: Naif algoritmalar uygulanabilir
 - Ortaboy Veri: Bir bilgisayar sisteminde işlenebilen miktarda
(Feasibly processed on one machine)
 - **Büyük Veri: Bir makinaya sığmayan miktarda**
- Karmaşık ilişkisel veri
 - ikili, üçlü veya daha üst seviyeden etkileşimler
- Hız, akan veri
- Yapılandırılmamış veri (blog metinleri/video/fotoğraf)
- 3V: Volume-Variety-Velocity

Büyük Veri?

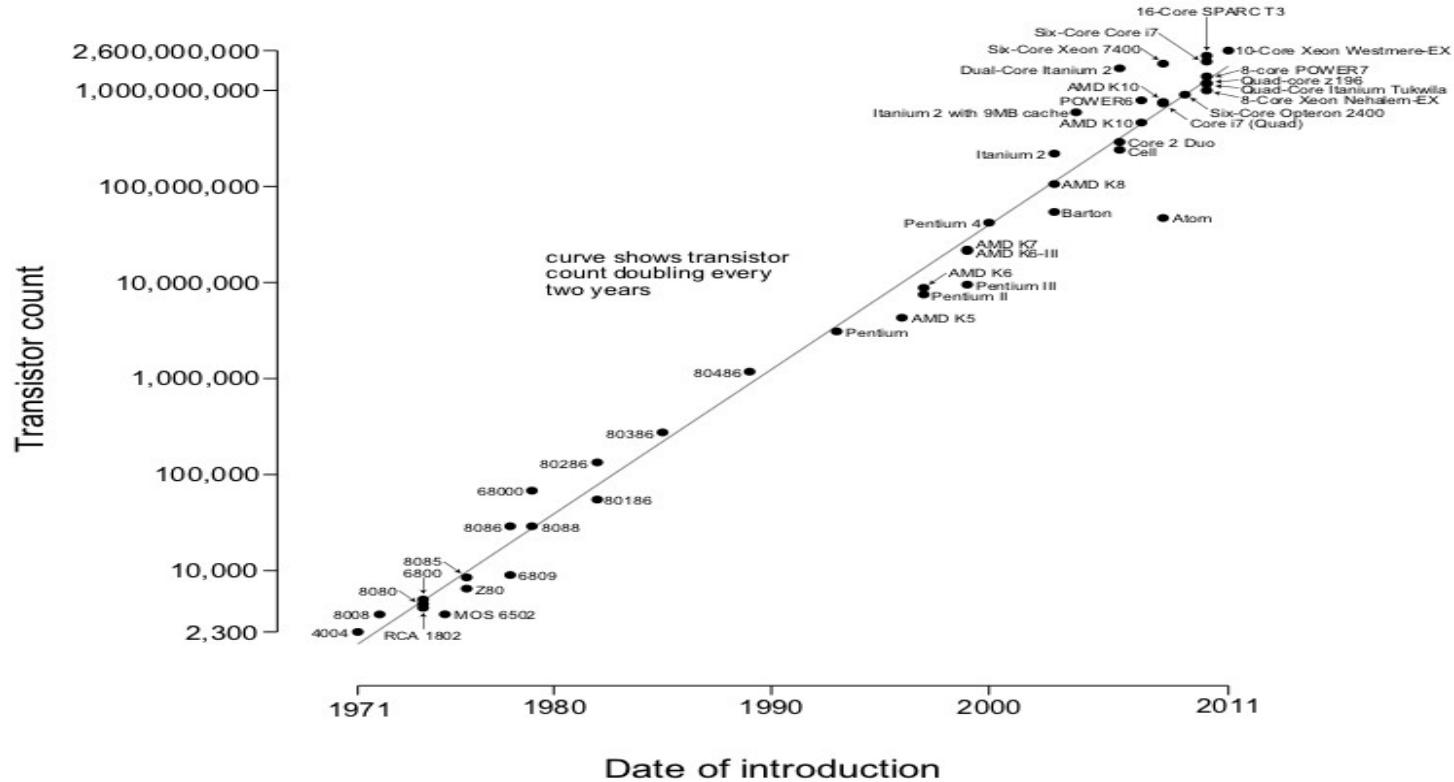
- <http://www.couchbase.com/nosql-resources/what-is-no-sql>



Moore Kanunu günü kurtarır mı?

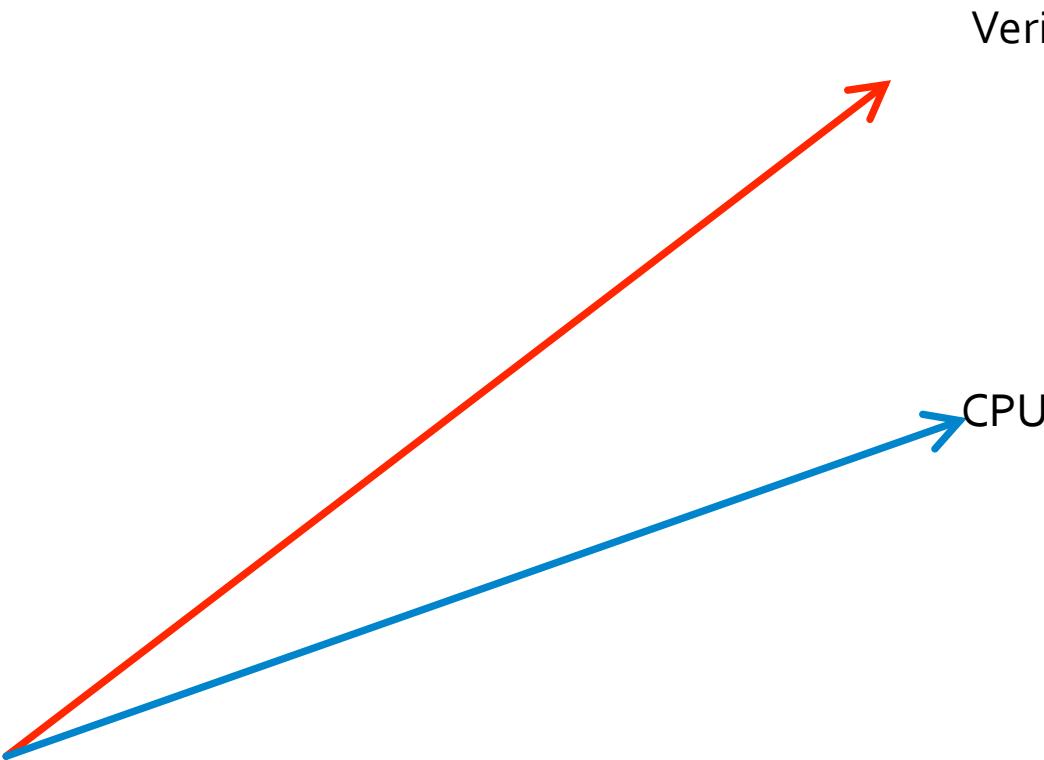
- "Transistor Count and Moore's Law - 2011" by Wgsimon - Own work. Licensed under CC BY-SA 3.0 via Wikimedia Commons - [https://commons.wikimedia.org/wiki/File:Transistor_Count_and_Moore%27s_Law_-_2011.svg](https://commons.wikimedia.org/wiki/File:Transistor_Count_and_Moore%27s_Law_-_2011.svg#/media/File:Transistor_Count_and_Moore%27s_Law_-_2011.svg)

Microprocessor Transistor Counts 1971-2011 & Moore's Law



Moore Kanunu günü kurtarır mı?

- Veri patlaması Moore Kanunundan hızlı
- Fiziksel Enerji bariyeri



Parkinson'un 1. Kanunu

- 'Data Expands to fill the space available for storage'
- Veri, kaplar her yeri



Veri = Enformasyon ≠ Bilgi

- Data = Information ≠ Knowledge

Enformasyon içinde boğulurken bilgiye açlık çekiyoruz

We are drowning in data and starving for knowledge

– J. Naisbitt

(Machine Learning, a probabilistic perspective, KP Murphy)

Kullanım Senaryoları: Tavsiye Sistemleri

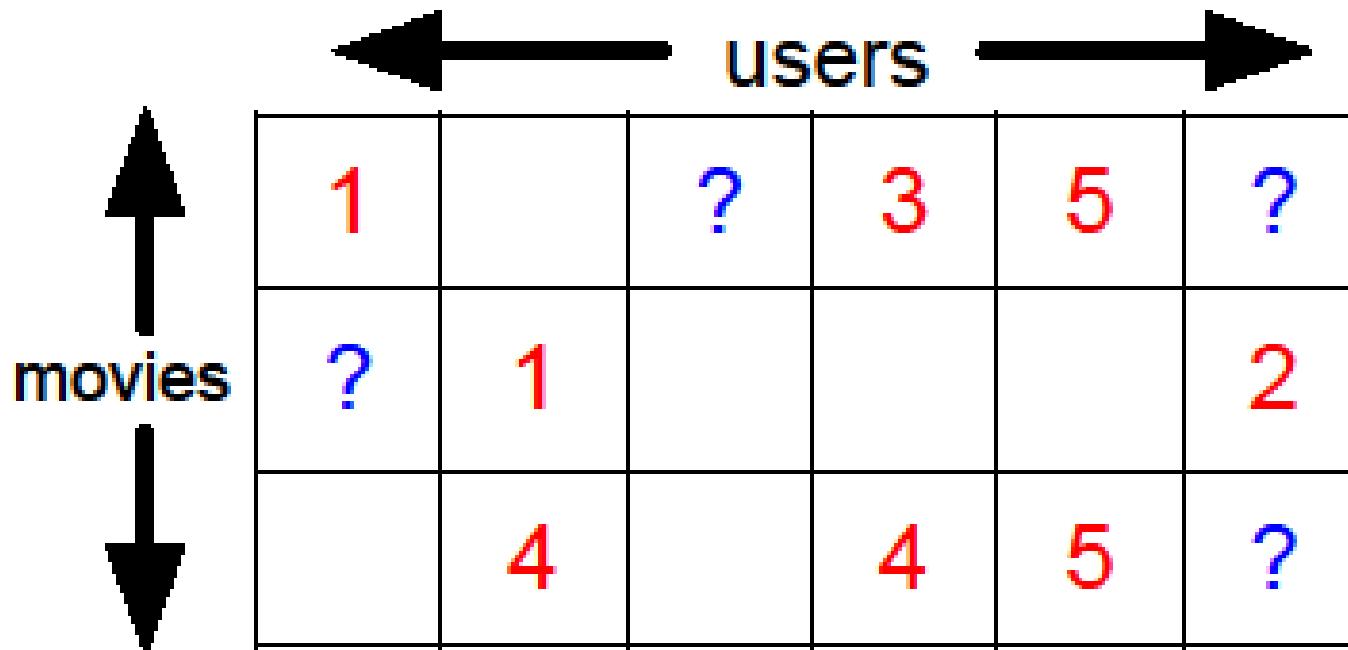


Kullanım Senaryoları: Perakende/ Tüketim

- Ürün Tavsiye Sistemleri
- Sepet Analizi (Market Basket Analysis)
- Olay/Aktivite/Davranış Analizi (Event/Activity/Behavior Analysis)
- Kampanya yönetimi ve eniyilemesi
- Tedarik zinciri yönetimi
- Pazar ve Tüketici ayrıştırması

Kullanım Senaryoları: Tavsiye Sistemleri

- Netflix: 18K film \times 500K kullanıcı %99 seyrek

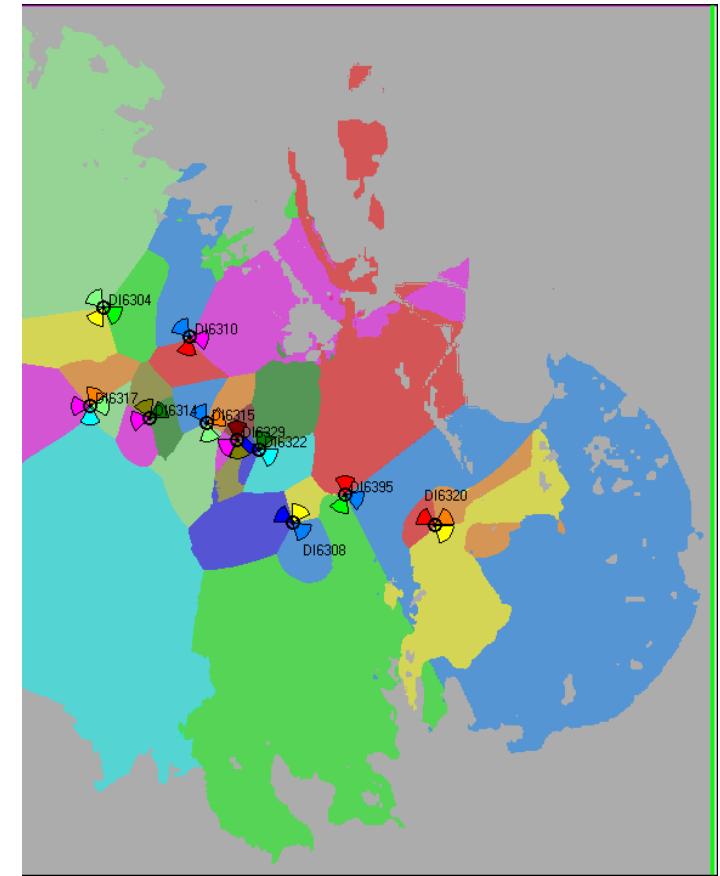
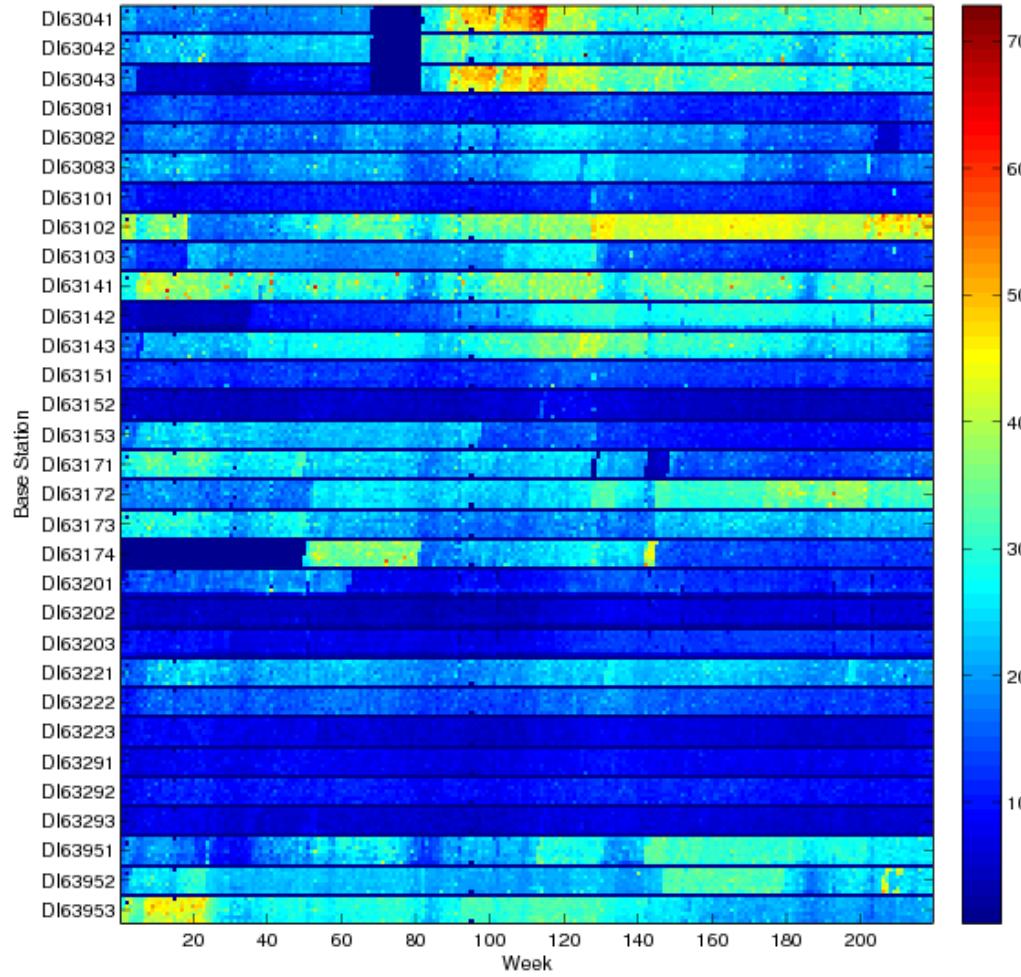


Kullanım Senaryoları: Haberleşme

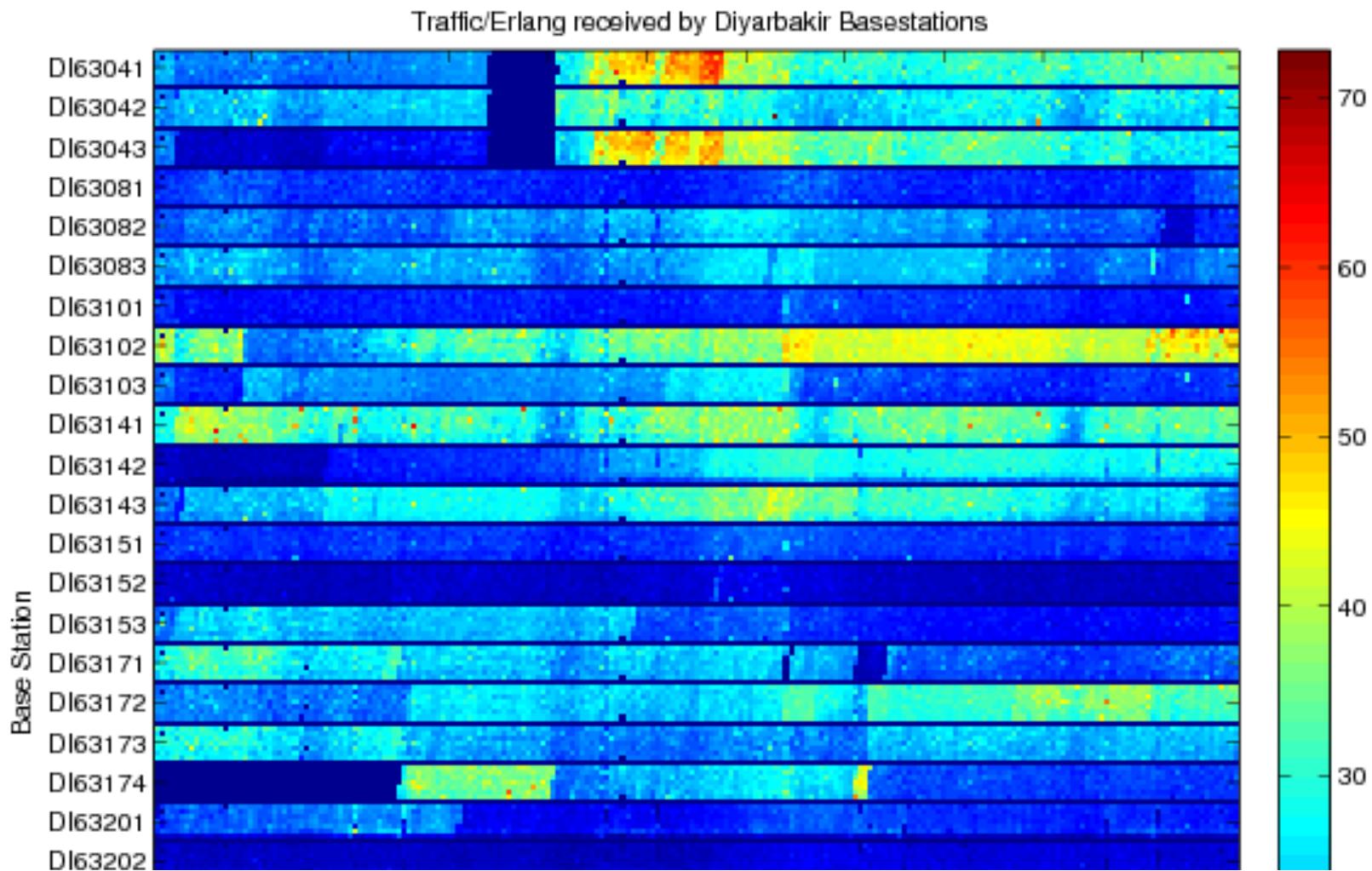
- Ağ izleme ve performans eniyileme
- Fiyatlandırma
- Müşteri ayrılma (Churn) tahmini
- Call Detail Record (CDR) Analizi
- (Mobile) Kullanıcı Davranış Analizi
- Siber güvenlik, DDOS saldırılarının tespiti ve önlenmesi
- Altyapı Planlaması

Kullanım Senaryoları, Haberleşme

Traffic/Erlang received by Diyarbakir Basestations

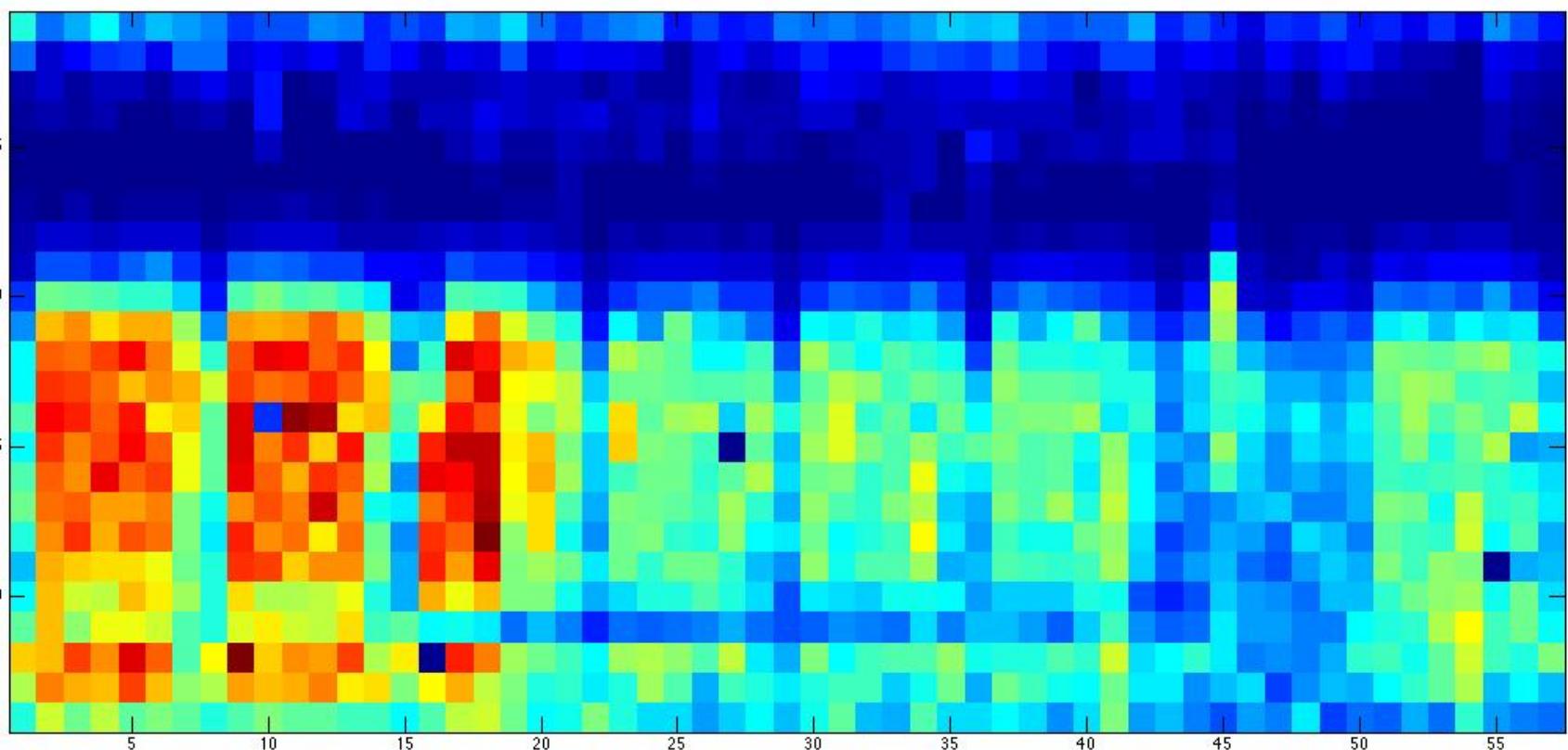


Kullanım Senaryoları

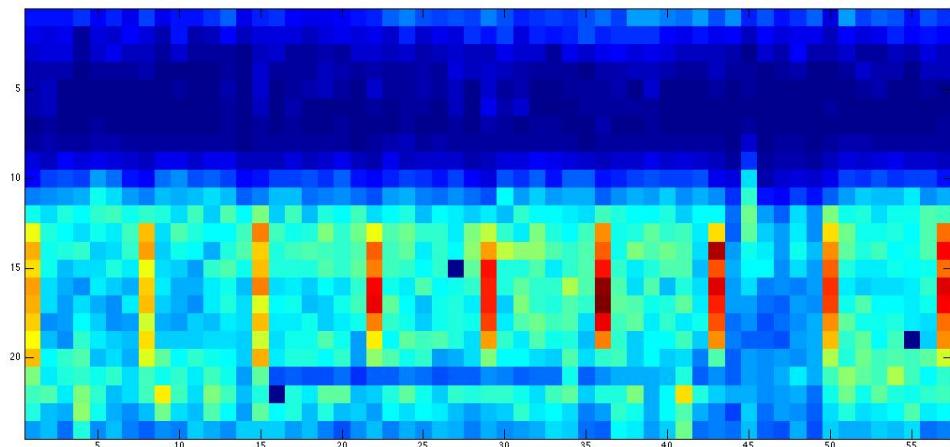
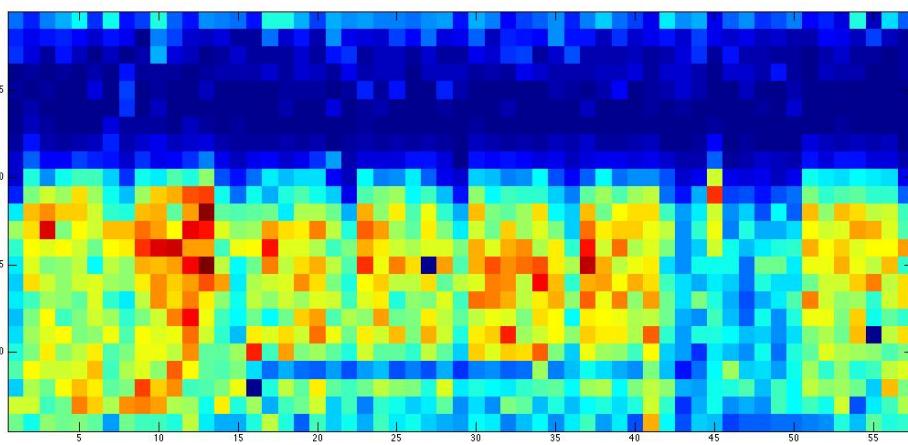
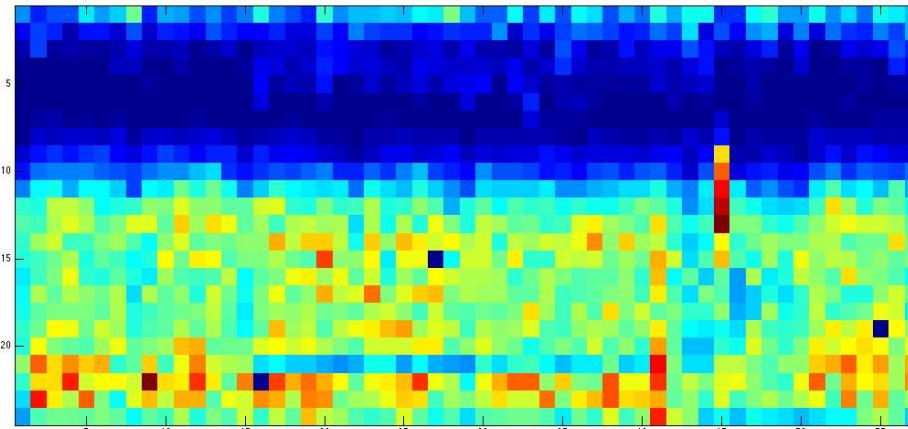


Kullanım Senaryoları, Haberleşme

- X: gün, Y: saat, renk: Kullanım Miktarı



Haberleşme

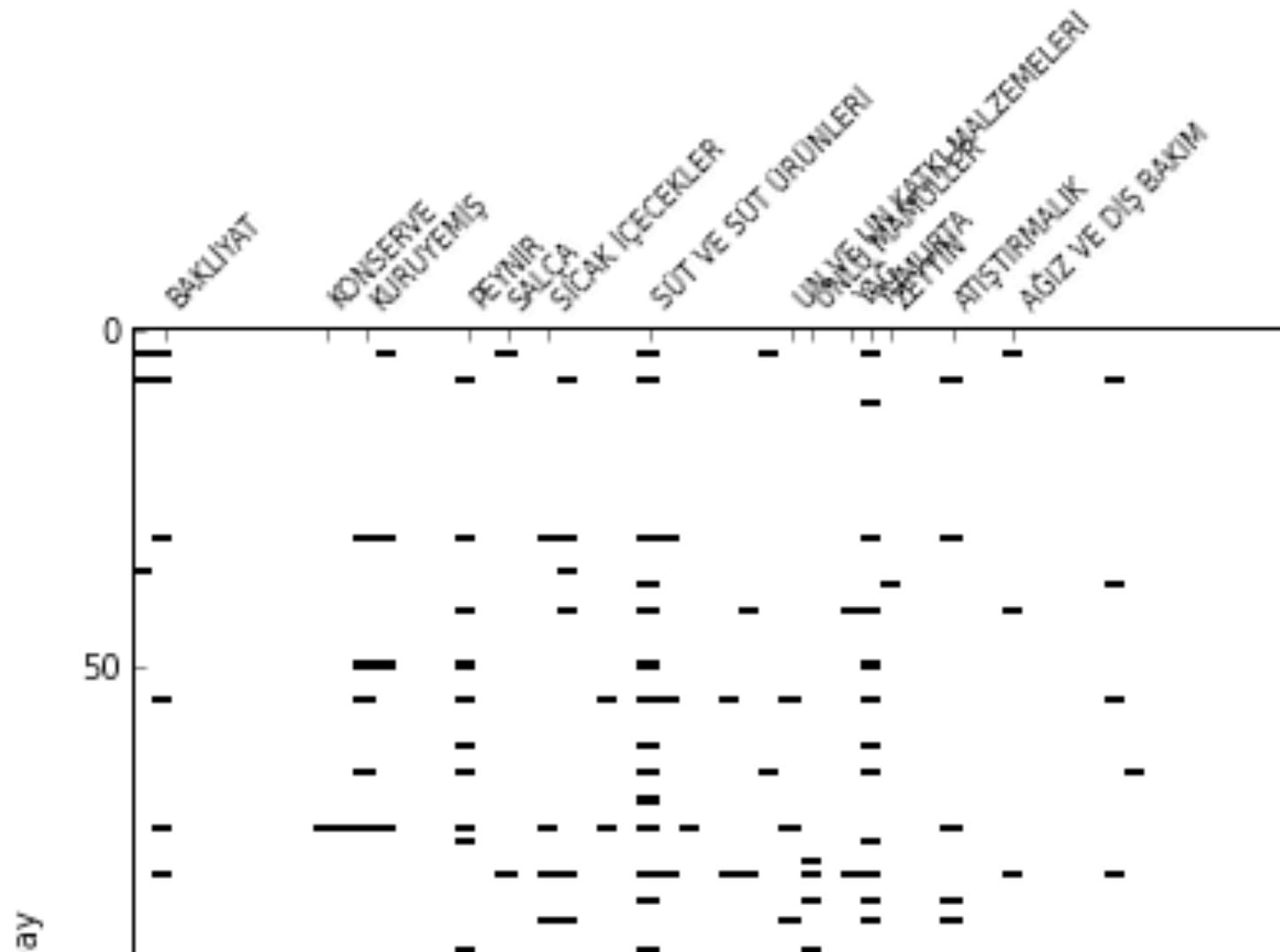


Kullanım Senaryoları: Finans/Ticaret/ Bankacılık

- Yolsuzluk (Fraud) Tespiti/Risk Kestirimi
- Insider Trading
- Yüksek hızda Trading
- Anomalite/Değişim noktası tanıma
- Sentiment

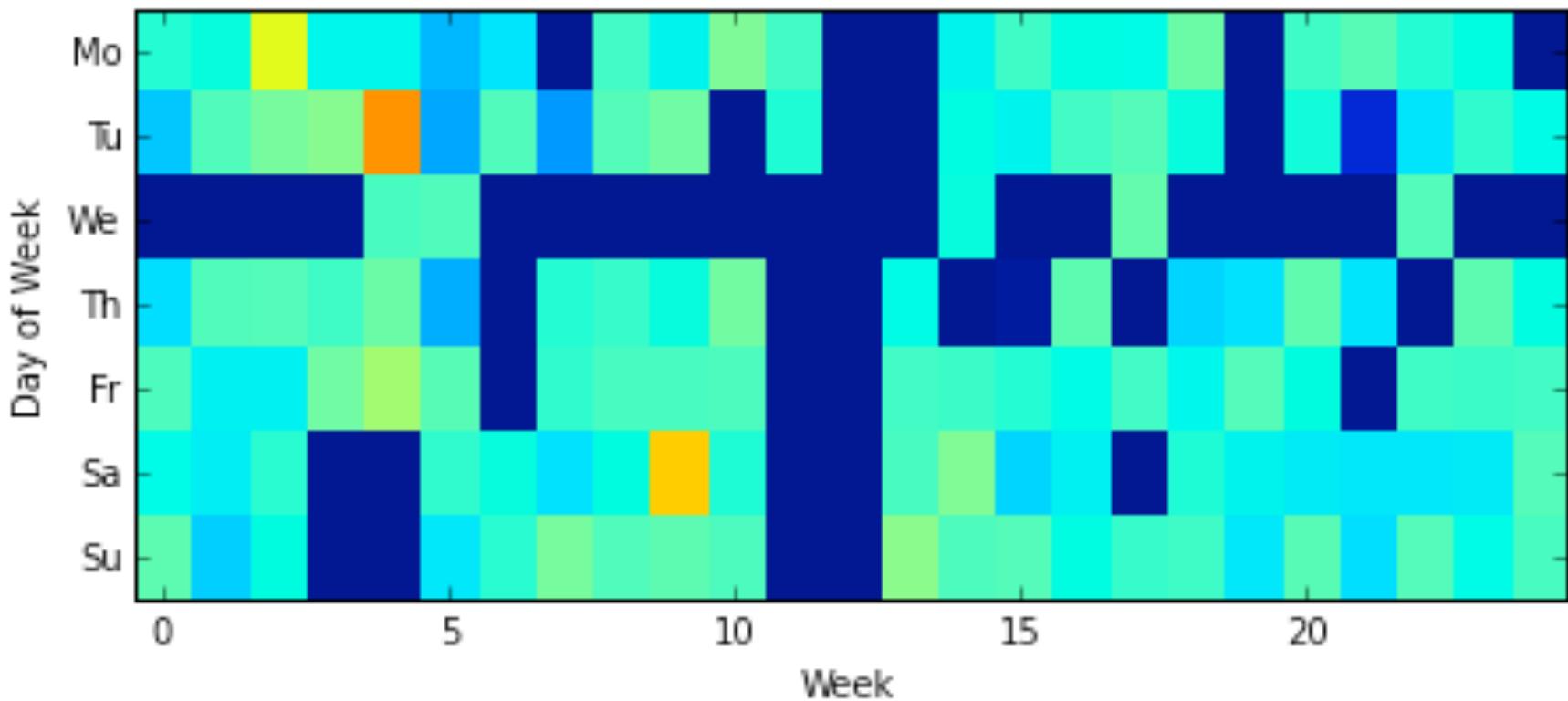


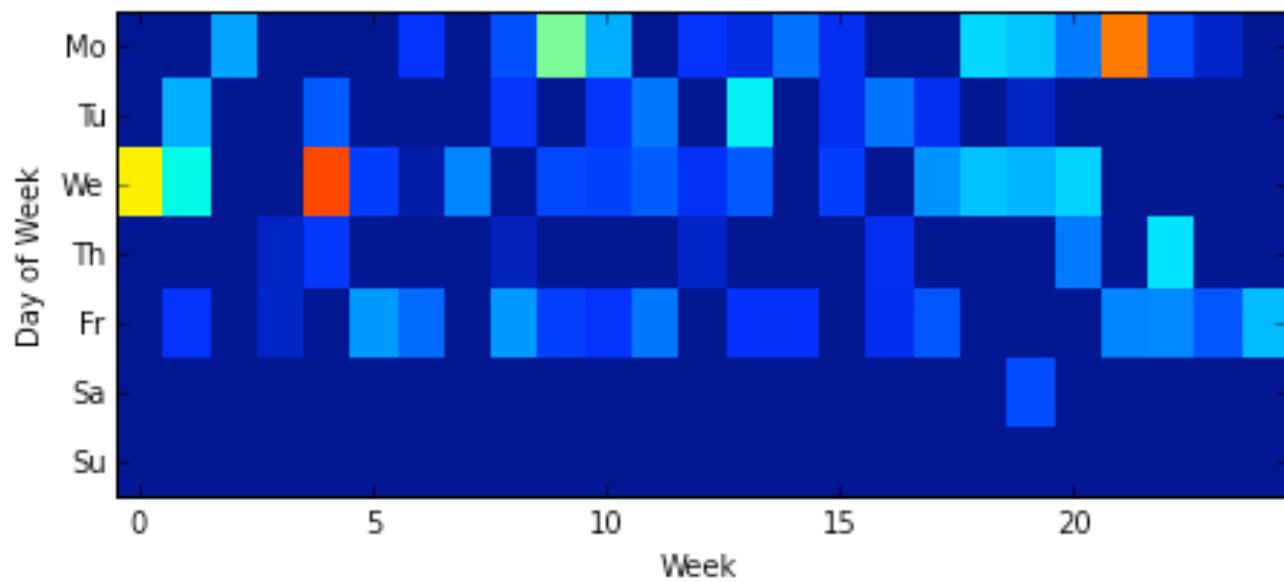
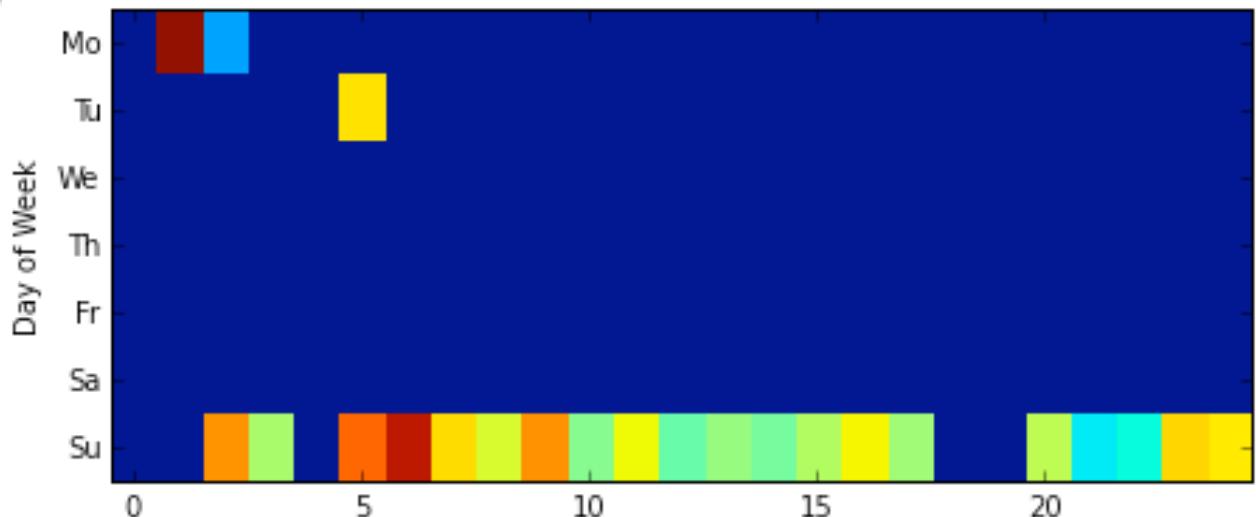
Customer Behavior Segmentation

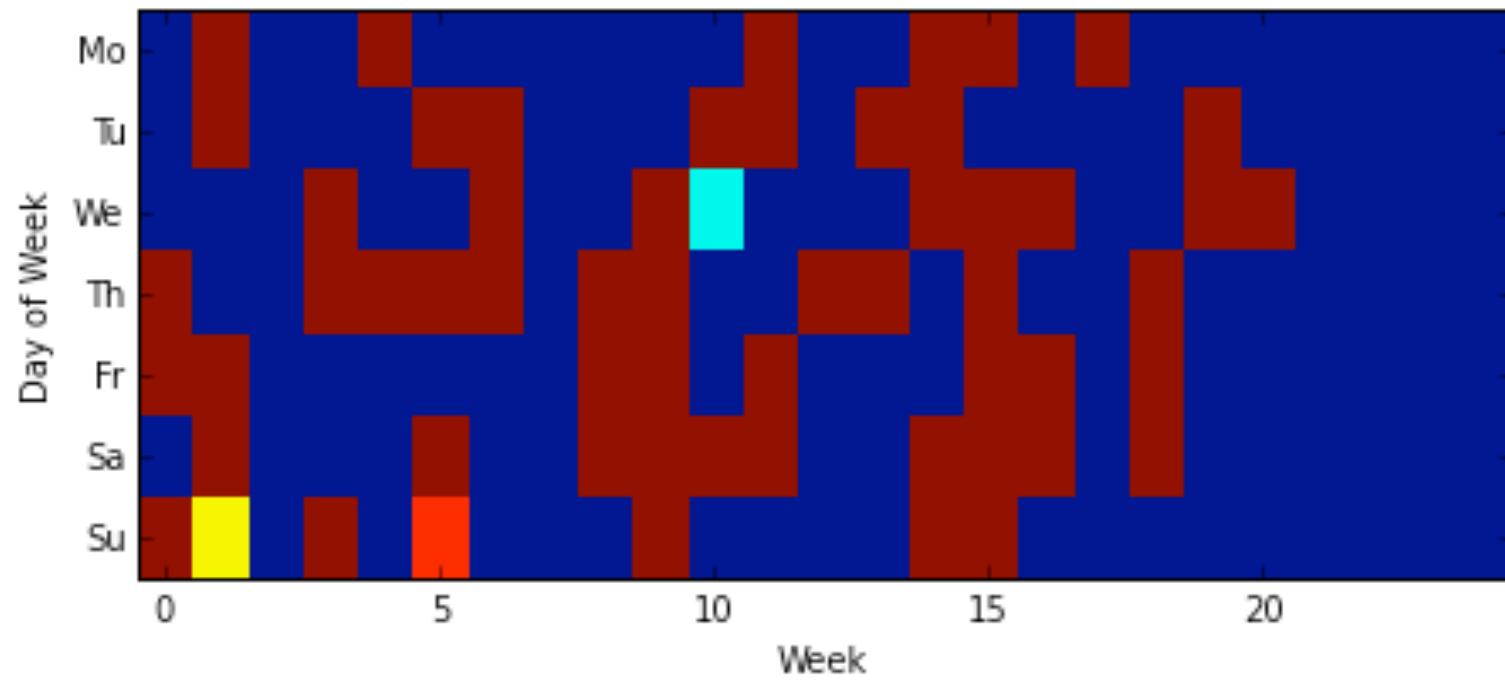


Customer Behavior Segmentation

Süpermarket: Haftanın Günlerine Göre Alışveriş miktarı

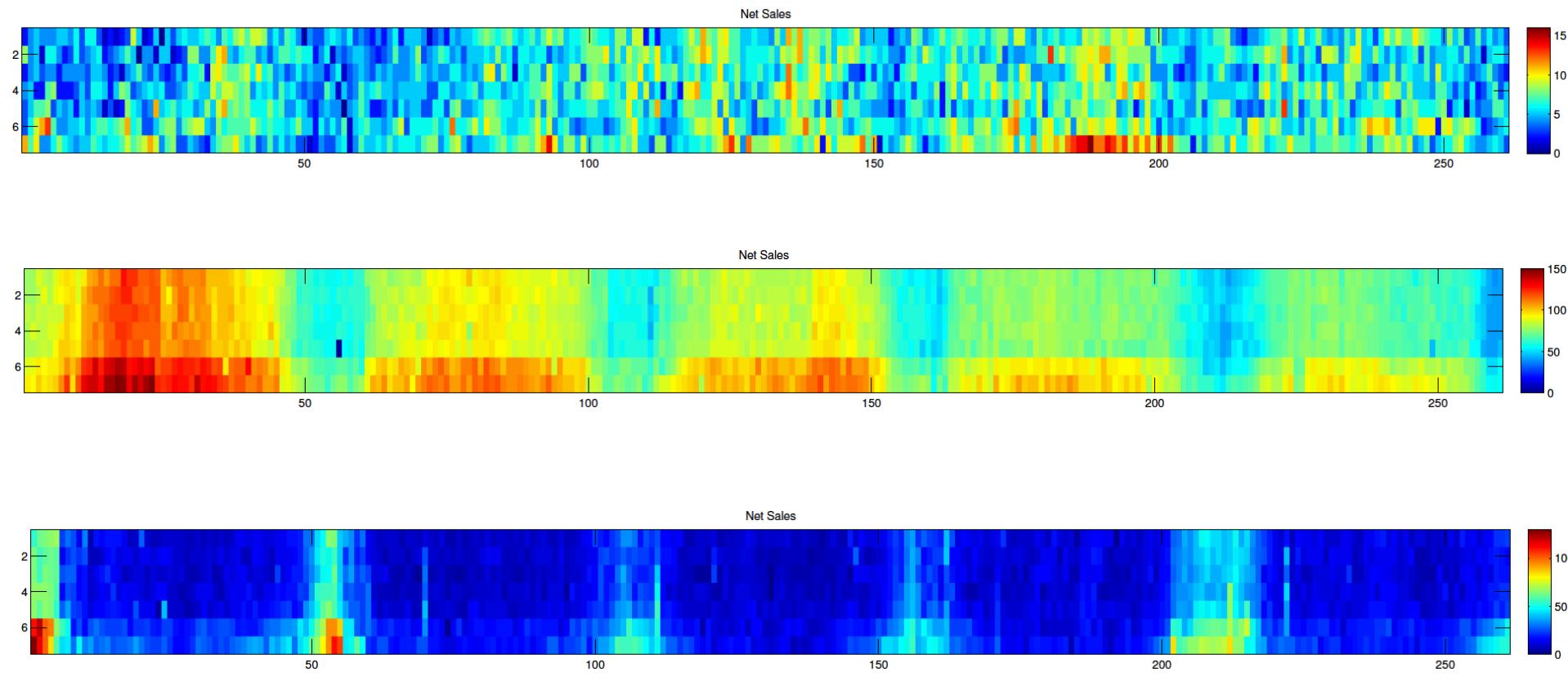






Müşteri Toplam Ciro: 198.538 TL, Toplam Maliyet: 208.516 TL

Demand Prediction: Gazete Satışları



Finansal Veri Analizi (Algosis)

gosis
ALGO Trader Demo

Ekle:

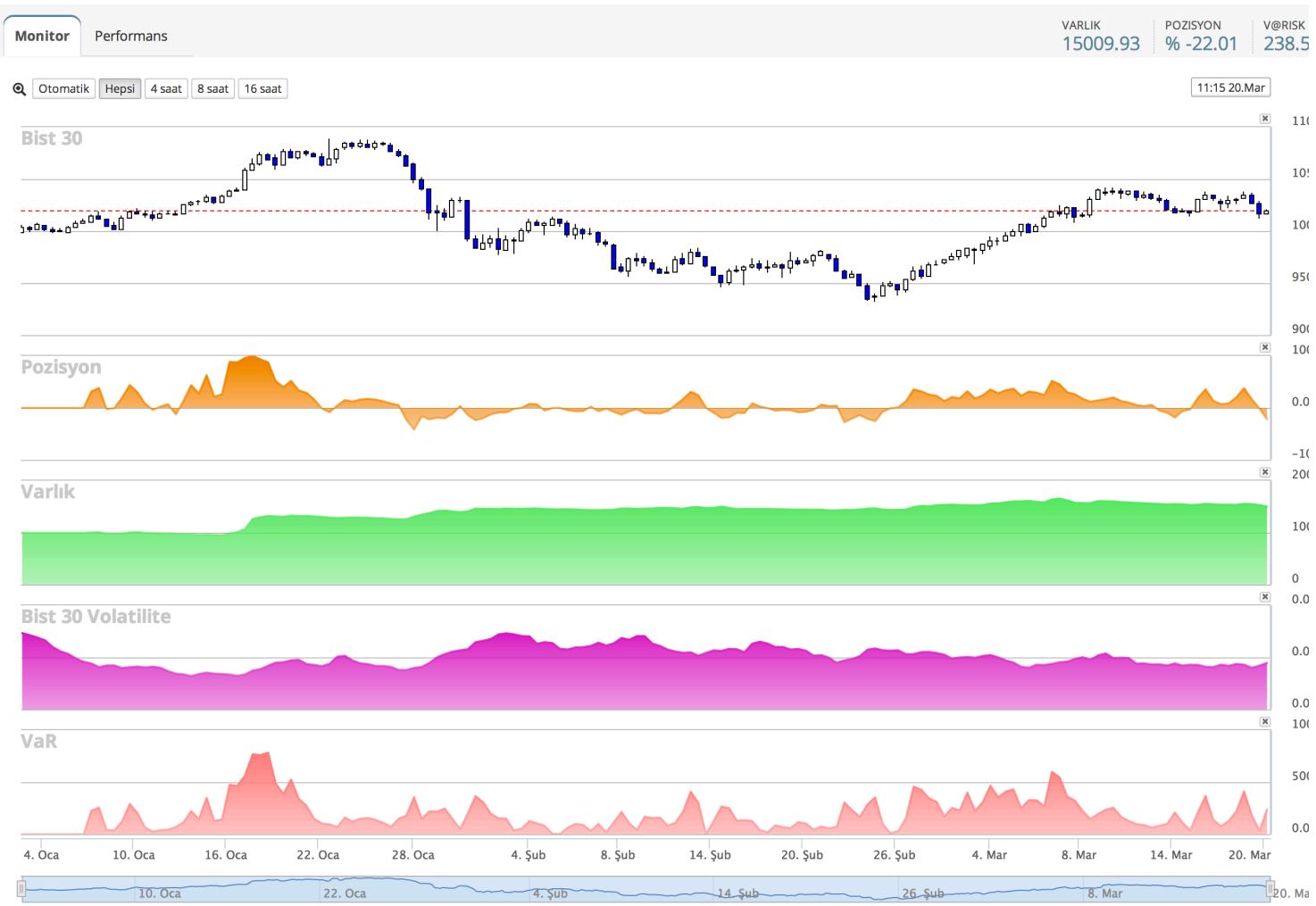
embol	Alış	Satış	Poz.	V@Risk	Volatil.
BIST 30	102025	102025	% -22.01	7.22	0.00

AlgоЕ AlgoRE AlgoARE

Algoritma Kontrol Paneli ON

	Durum (11:15)	Yeni Durum
Pozisyon	% 22 kısa	-
Risk	238.54	-

nir, BIST30, Alış 12 lot, G. F. 101650	19:40 19.Mar
nir, BIST30, Alış 11 lot, G. F. 101625	19:15 19.Mar
nir, BIST30, Satış 14 lot, G. F. 101450	19:05 19.Mar
nir, BIST30, Satış 7 lot, G. F. 101900	18:55 19.Mar
nir, BIST30, Alış 12 lot, G. F. 102000	18:50 19.Mar

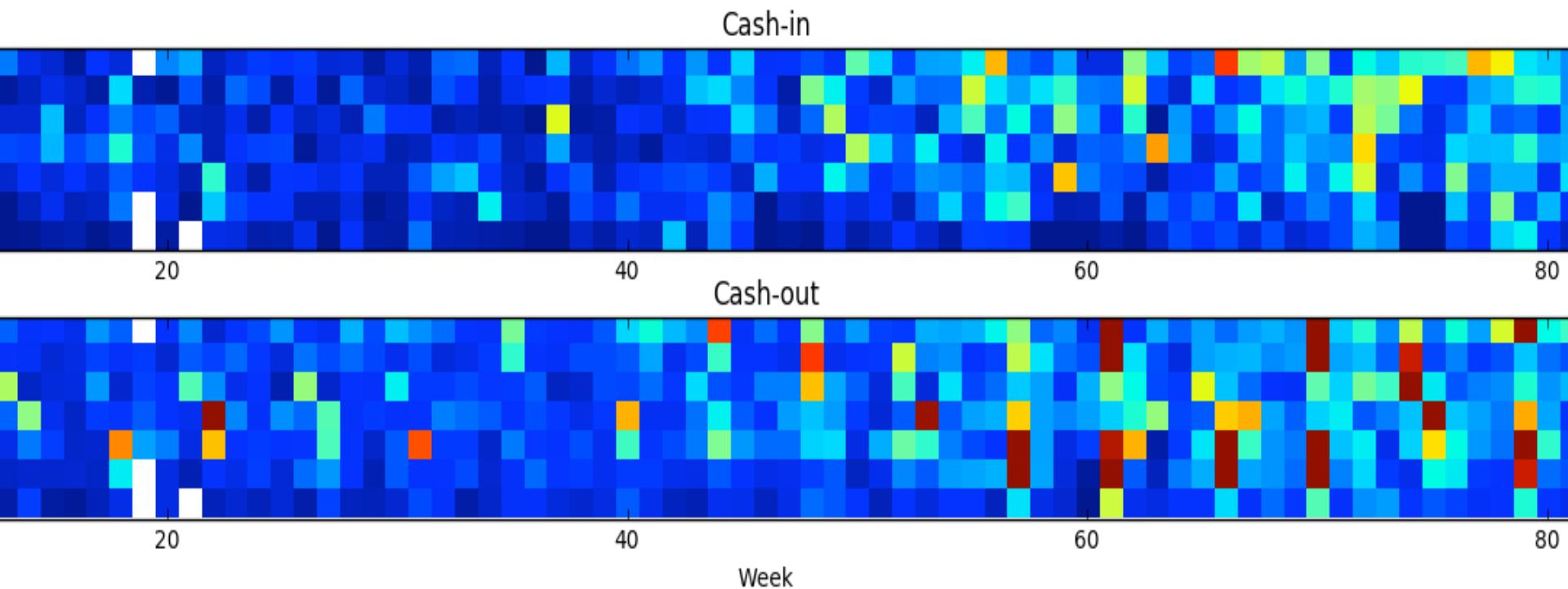


Bankacılık Operasyonları

- Veri kalitesi yönetimi, Eksik/anormal veri belirleme
- Müşteri davranışları kestirimi
- Şube, Çağrı merkezi, yoğunluk tahmini
- Sahtekarlık belirleme
- Kredi skorlama, Risk kestirimi
- Nakit transfer planlama
- Finansal Ürün tavsiyesi
- Kuyruk yönetimi

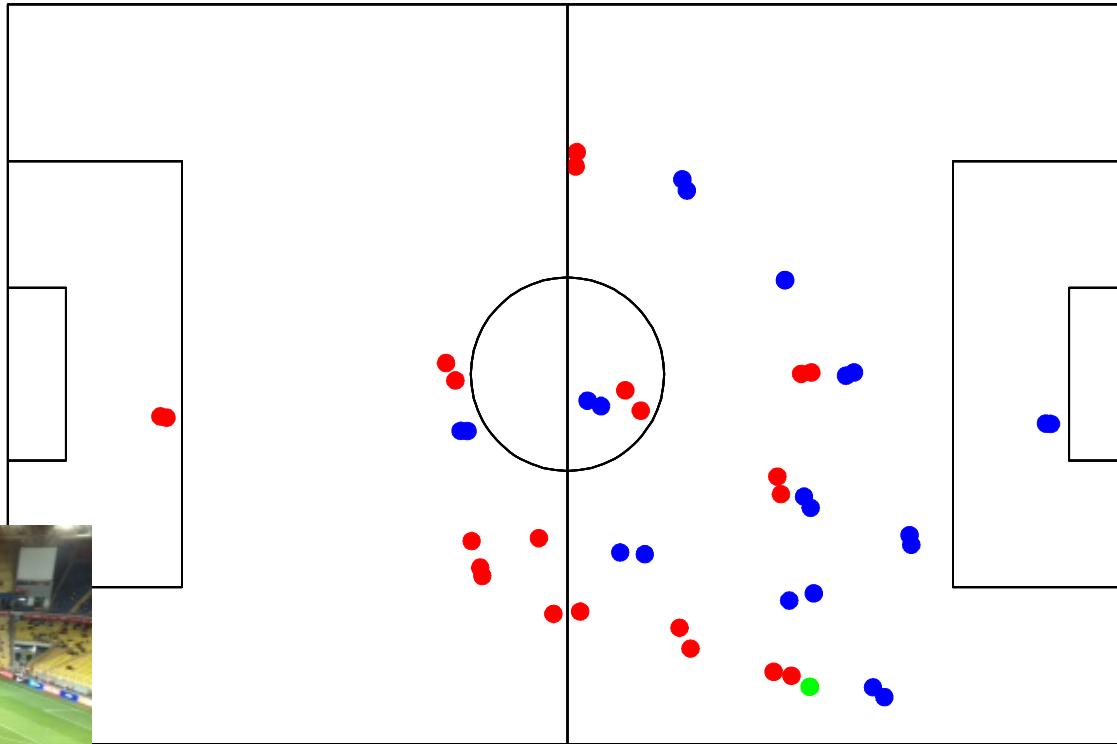
ATM Planlama

■ Para Çekme ve Yatırma davranışları



Spor Analitiği (Exatech)

00:43:55



■ Oyuncu Takibi

Kullanım Senaryosu örnekleri

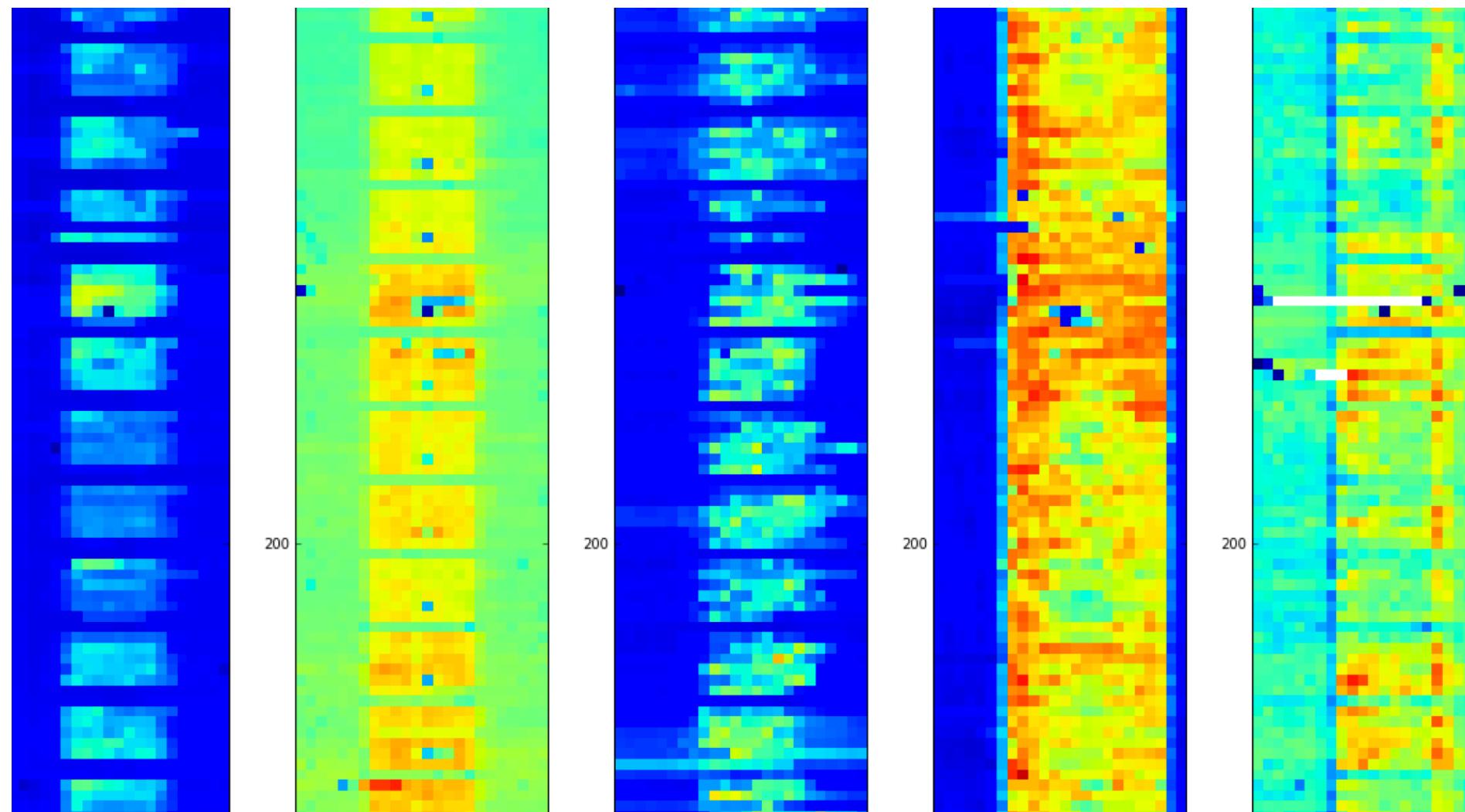
- Reklam Kişileştirme
 - Google ve Yahoo'nun temel gelir kaynağı

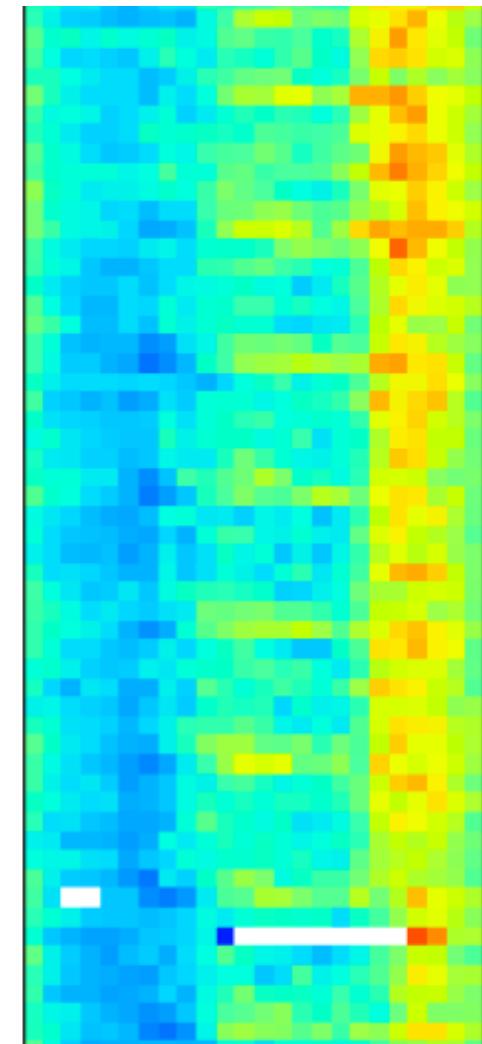
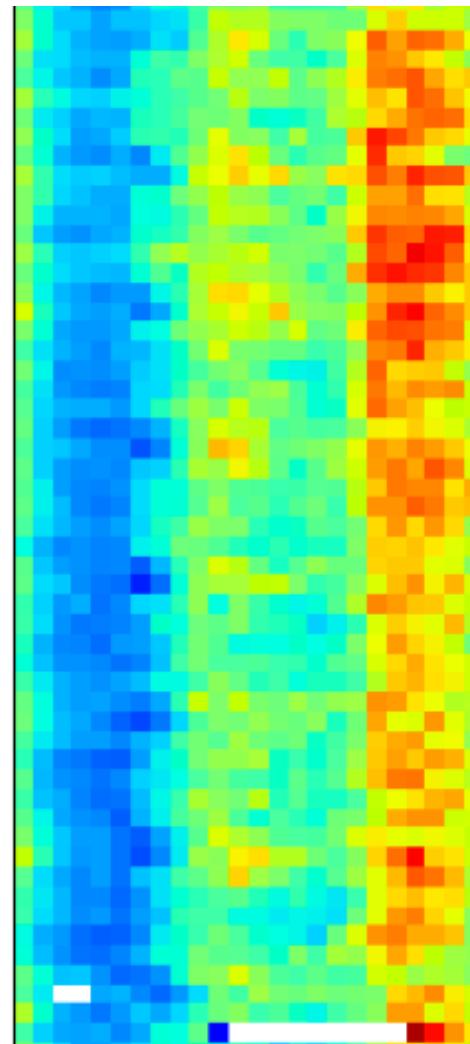
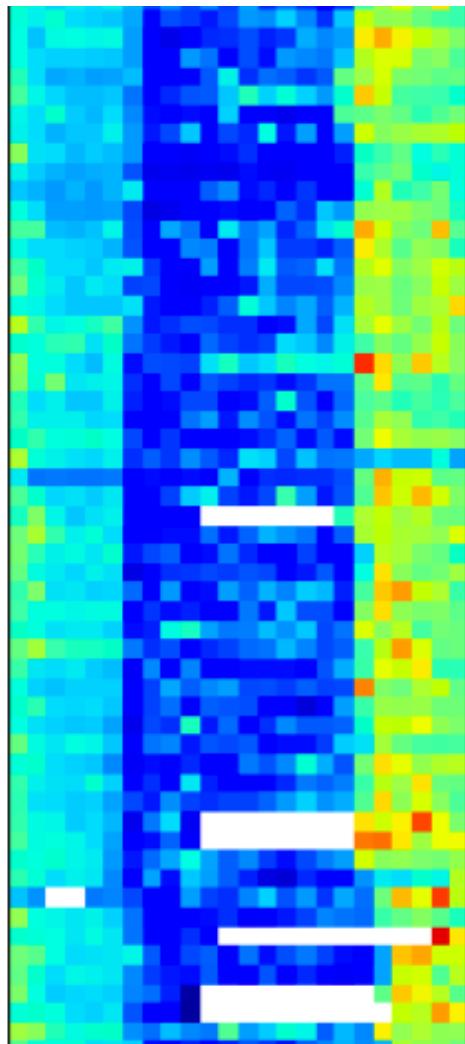
The screenshot shows a Gmail inbox with 19,360 messages. A specific email from 'Abu Dhabi's Mubadala' is highlighted with a red box. This message is a forwarded email from 'Yaser S. Abu-Mostafa @ ITU' dated Dec 25, 2012, at 10:30am. The forwarded message is from 'Zehra Cataltepe <cataltepe@itu.edu.tr>' on Sun, Dec 23, 2012 at 7:55 PM, with the subject 'Yaser S. Abu-Mostafa @ ITU, Dec 25, 2012, 10:30am'. The message body starts with '----- Forwarded message -----'. To the right of the main message, there is a sidebar for 'Ethem Alpaydin' with options like 'Colleagues', 'Show details', and 'Ads – Why these ads?'. The left sidebar includes links for 'Compose', 'Inbox (19,360)', 'Starred', 'Important', 'Sent Mail', and contact lists for 'arda.yurdakul', 'Baris Gokce', and 'Bruce Sharpe'. The top navigation bar shows the user 'Taylan Cemgil' with 6 unread messages.

Kullanım Senaryosu: Kamu Yönetimi

- Trafik Yönetimi (Urban Traffic Management)
- Enerji Dağıtım şebekesi yönetimi/eniyilemesi (Energy Grid Management/Optimization)
- Power Generation Management
- Çevre gözleme (Environment Monitoring)

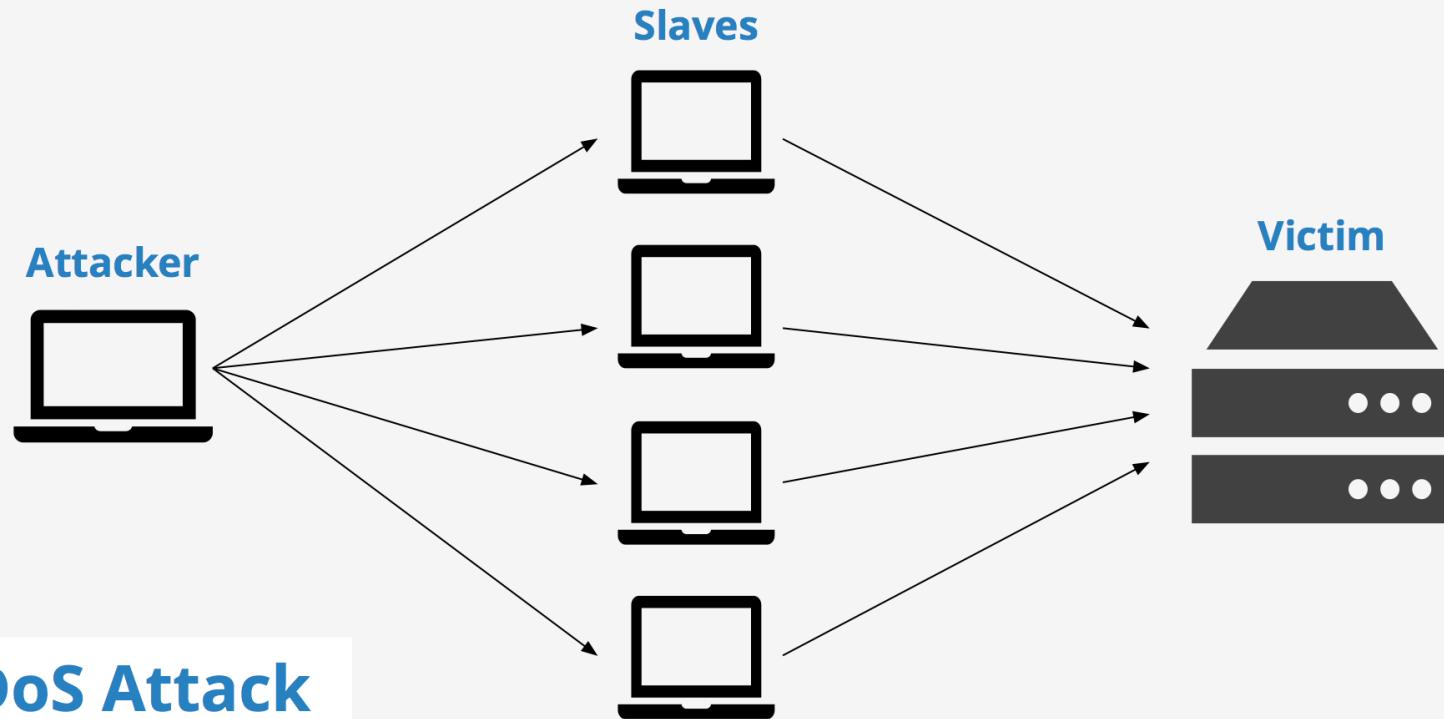
Elektrik Enerjisi Kullanımı





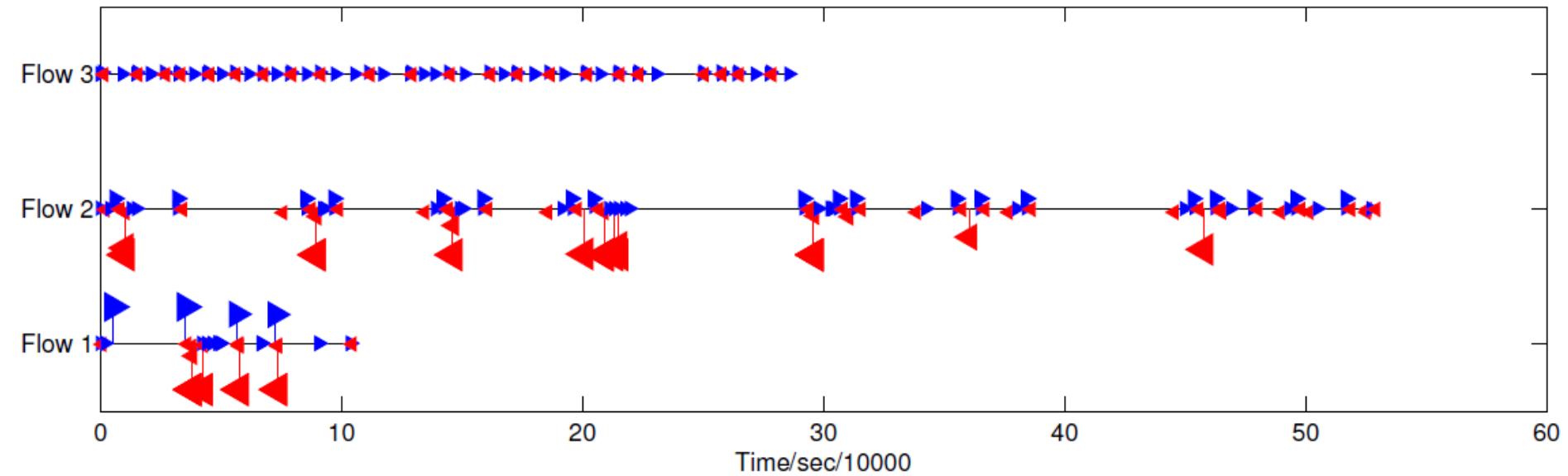
Cyber Security

- Distributed Denial of Service (DDoS)

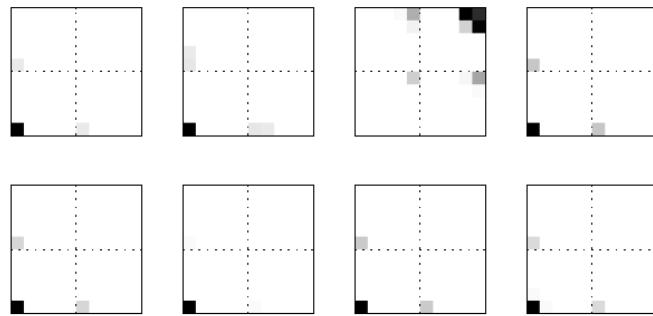
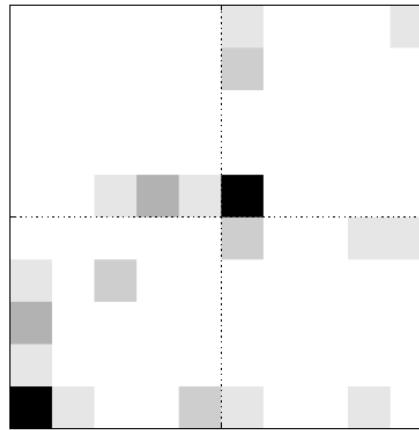
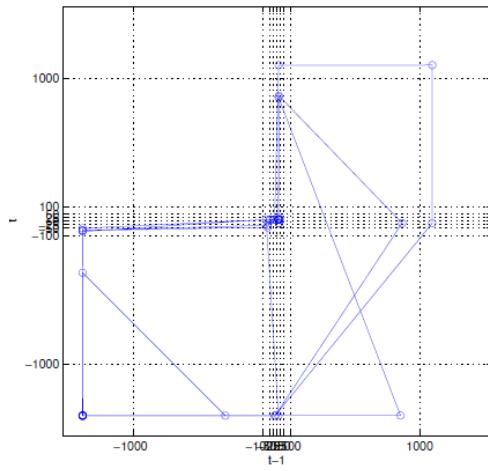


Source: <https://www.keycdn.com/support/ddos-attack/>

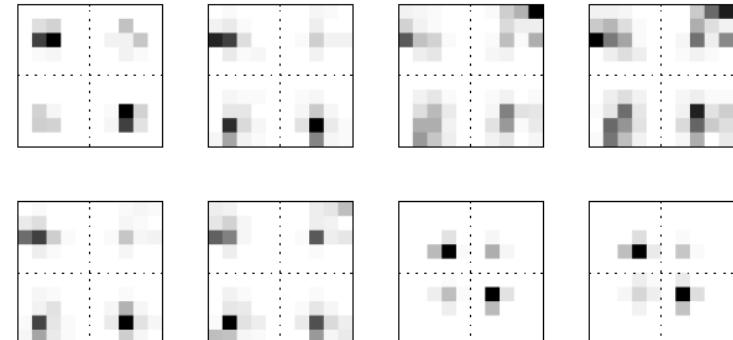
Paket Akışı



Öznitelikler



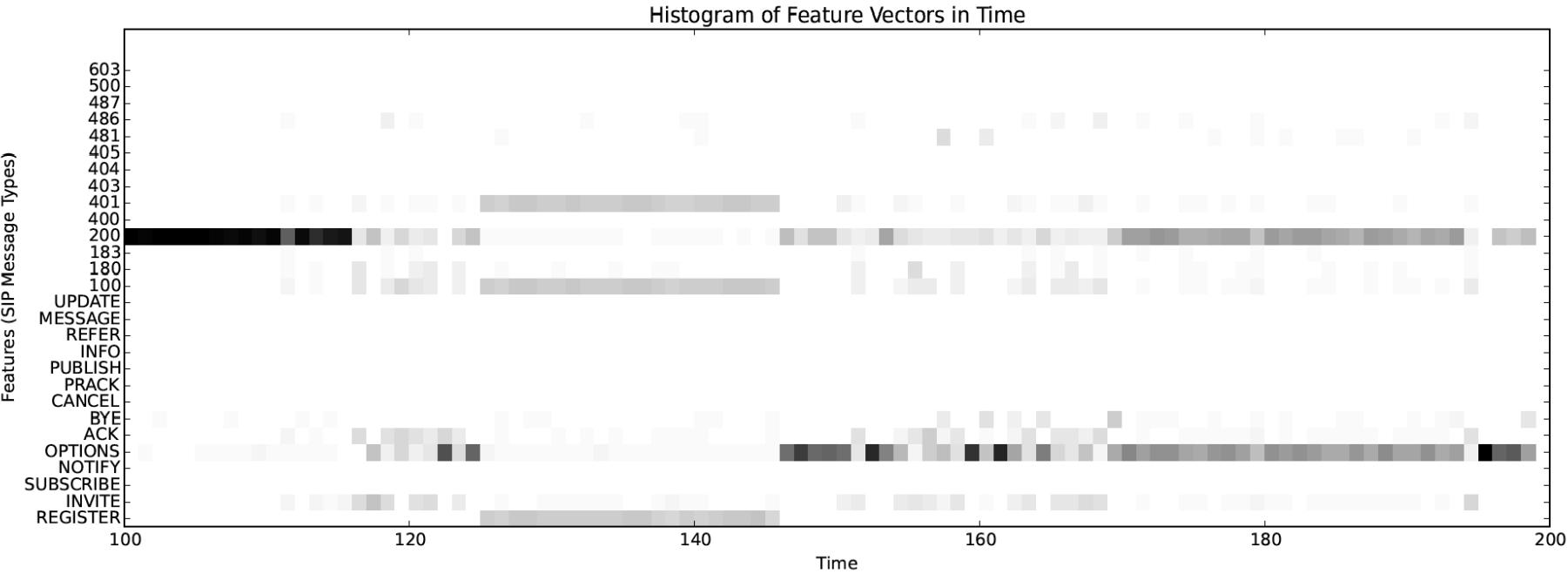
VIDEO



VIDEO₂

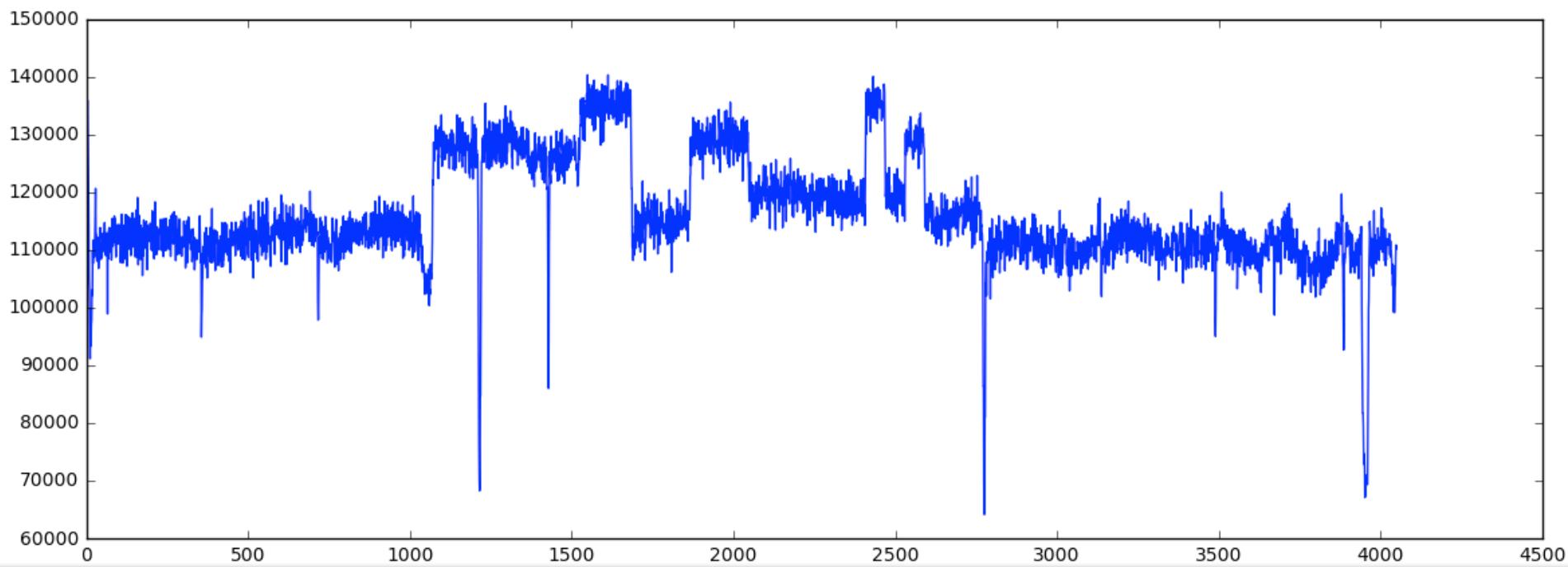
Cyber Security

- Example: Voice-over-IP, Monitor SIP traffic



Değişim Noktası Bulma

- 4050 nuclear magnetic resonance measurements taken from drill while drilling a well (Schlumberger)



Sağlık/Yaşam Bilimleri ve Biyoloji

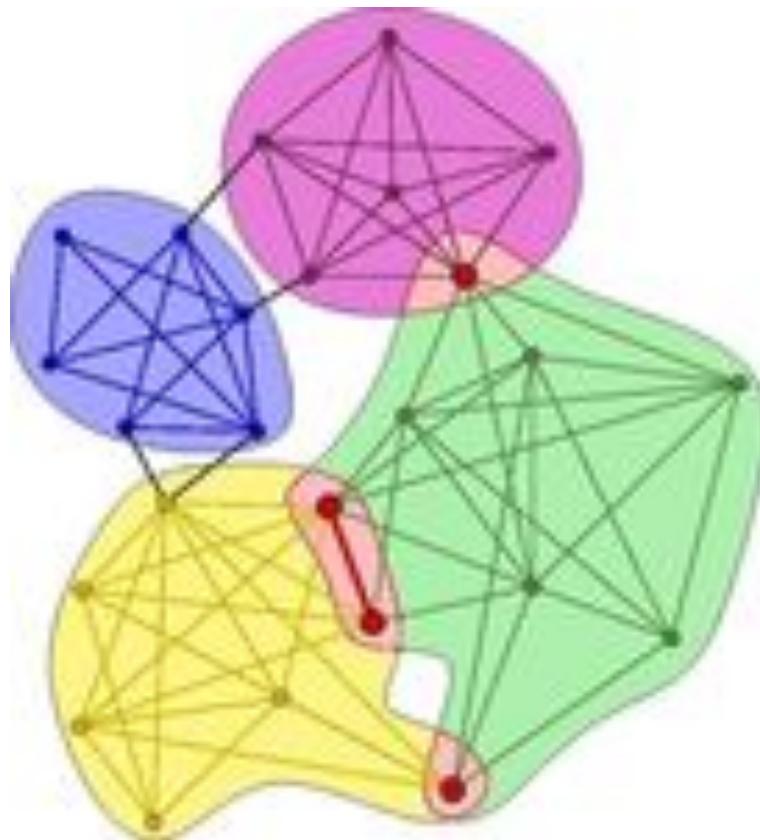
- Diagnosis and Medical Expert systems
- Health Insurance fraud detection
- Patient care quality and program analysis
- Drug discovery
- Remote Monitoring

Kullanım Senaryoları: Web

- Klik akısı öbekleme ve analizi
(Clickstream Segmentation and Analysis)
- İlan hedefleme/seçim/tahmin/eniyileme
(Ad Targeting/Selection)
- Klik Yolsuzluğu/Engelleme
(Click Fraud Detection/Prevention)
- Sosyal Ağ Analizi
- Müşteri Böülütlemesi
- Newsgroup/Blog/Sosyal Medya gündem takibi

Çizge/Ağ Analizi

- Sosyal Ağlarda gruplanmalar (source: matlab exchange)



+1 Arkadaşı Ekle



Duygu, İlgi, Eğilim Tahmini

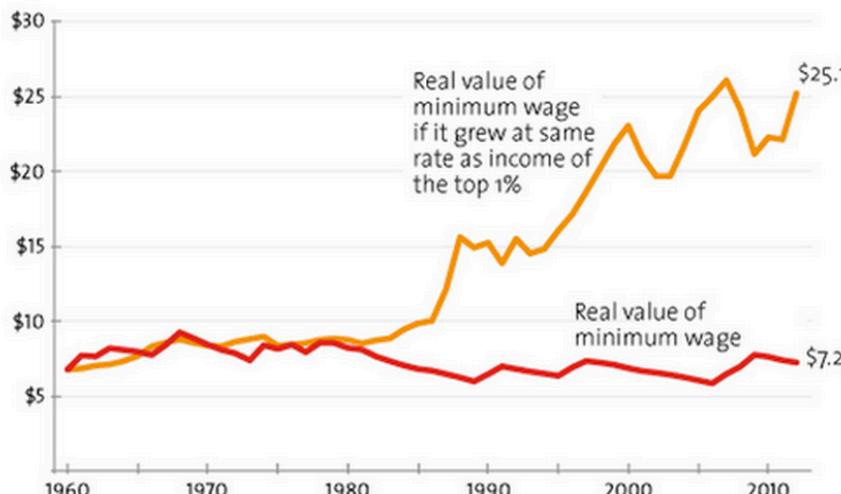


Conrad Hackett @conradhackett · 5h

Minimum wage would be \$25/hr if it grew like income of the top 1%

motherjones.com/politics/2013/...

What if minimum wage grew at the same rate as top incomes?



Based on income in 2012 dollars not including capital gains
Sources: Department of Labor, World Top Incomes Database

Mother Jones

RETWEETS
160

FAVORITES
82



Pivot @jeffposter1211 · 5h

@conradhackett somebody ought to paint an elaborate picture of what a beautiful world we would have if we applied that math



Emmanuel Acain @CitizenCainll · 4h

@conradhackett a good way to share that surplus value.



NikFromNYC @NikFromNYC · 4h

@conradhackett But would the top income, redistributed, at all cover the \$25 minimum wage? No! So your implication is moot.



N647 @night647 · 4h

@conradhackett is interesting you can see the trend picking up around 1985 wonder the conditions that produced this situation.



Doğal Dil İşleme

Its a btf nite, lukin for smth fun to do,
I think I wanna be w ma frnds.



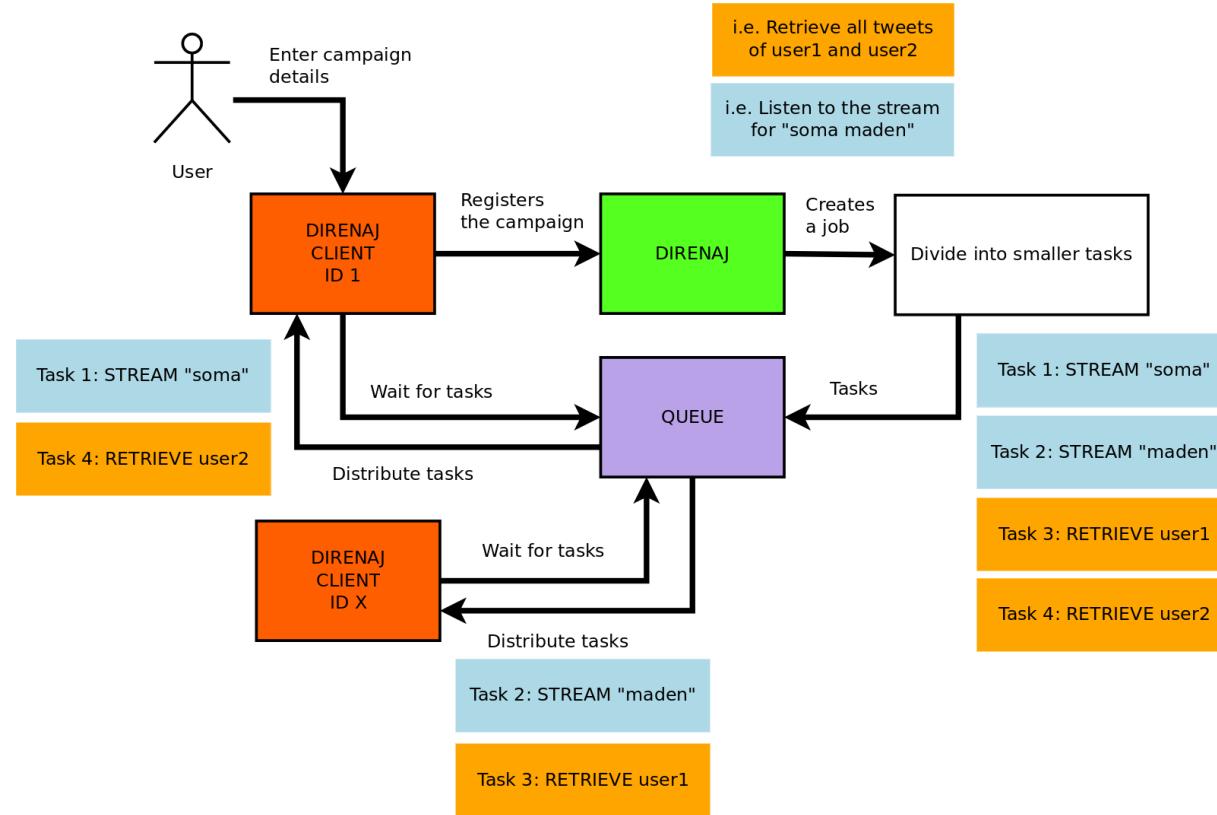
Its a beautiful night, looking for something fun to do,
I think I want to be with my friends.

Dun aksam okulun sitesi gene cokmustü.



Social Text Normalization, Çağıl Uluşahin ve Arzucan Özgür

Veri Toplama, Hesaplama ve Yazılım altyapıları



Catch Tweet streams
Crawl and Track Friend-Follower Graphs

Görselleştirme

Compose new Tweet...

Tweets

manovich @manovich · 2h
Tomnod lets people search satellite images for lost Malaysian plane
cnet.co/1ey0l0
View summary

Retweeted by HASTAC
Michael Widner @mwidner · 10h
Job: Program Coordinator for @HASTAC: hastac.org/opportunities/
Expand

Retweeted by Day of DH
Asher J. Klassen @AsherJKlassen · 4h
Got myself all registered 'n stuff on #DayofDH. Prepping the profile for #DayofDH2014, never been so pumped for April
dayofdh2014.matrix.msu.edu/members/ajklassen...
Expand

Retweeted by Dan Cohen
Dolly Jorgensen @DollyJorgensen · 9h
The Joy Ride for kids at the Los Angeles alligator farm at beginning of 20th c.
#envhist ASEH2014.pic.twitter.com/l5EtpfVac
Expand

designed by ALMILAKDAG. this is a mockup, please do not distribute.

ML for Big Data, Cemgil, 24.12.2012

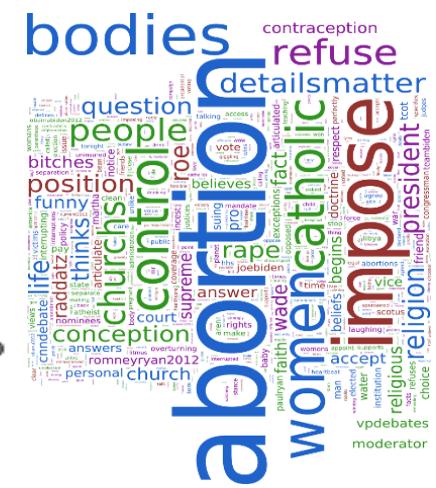
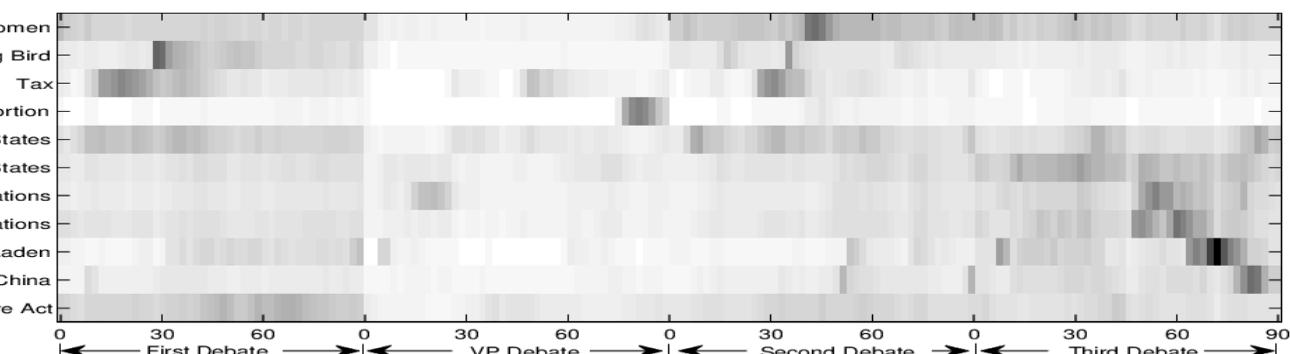
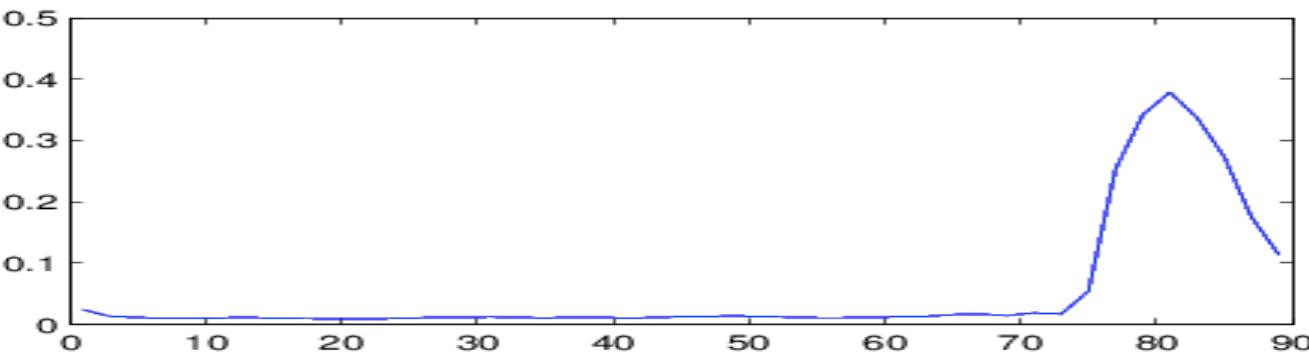
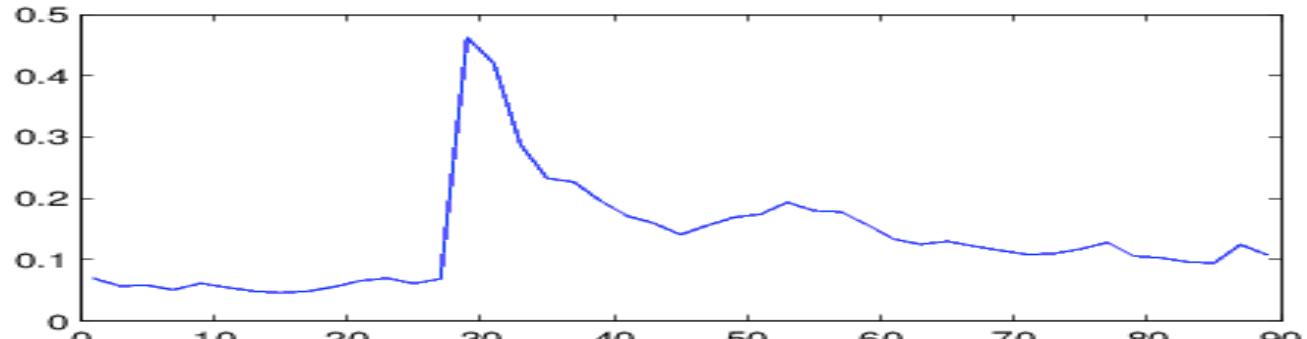
Konu Modellemesi (Topic M)



WIKIPEDIA
The Free Encyclopedia

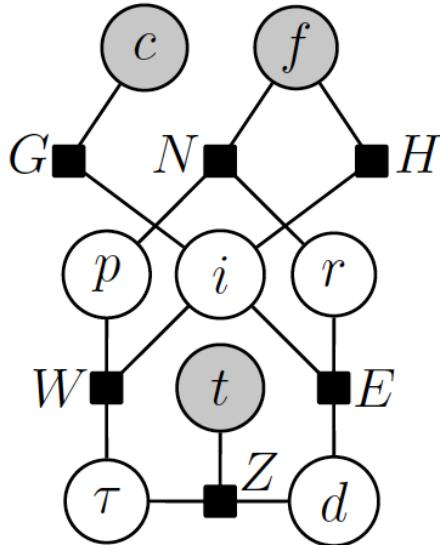
Dinamik Gündem Takibi

Yıldırım ve Üsküdarlı



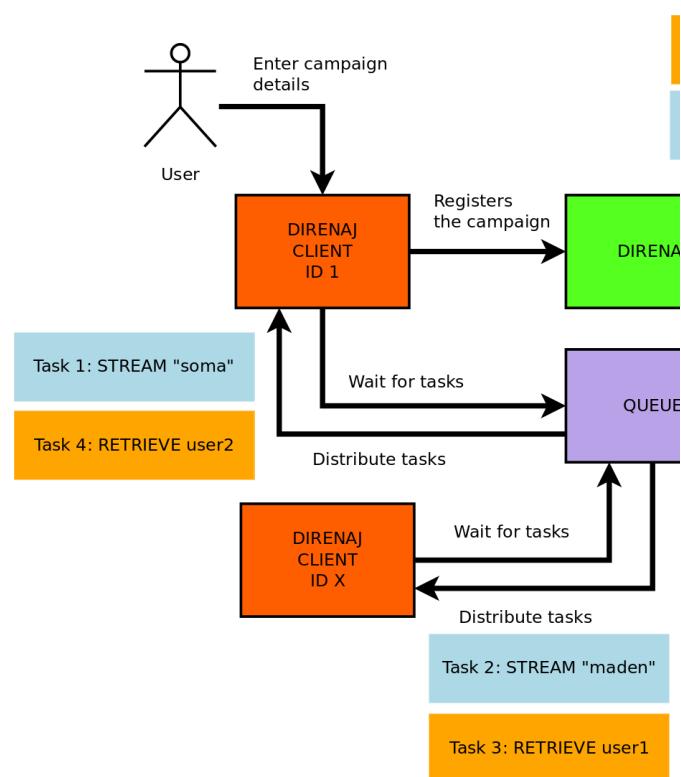
Modern Yapay Öğrenme

■ Modeller – Algoritmalar -- Sistemler



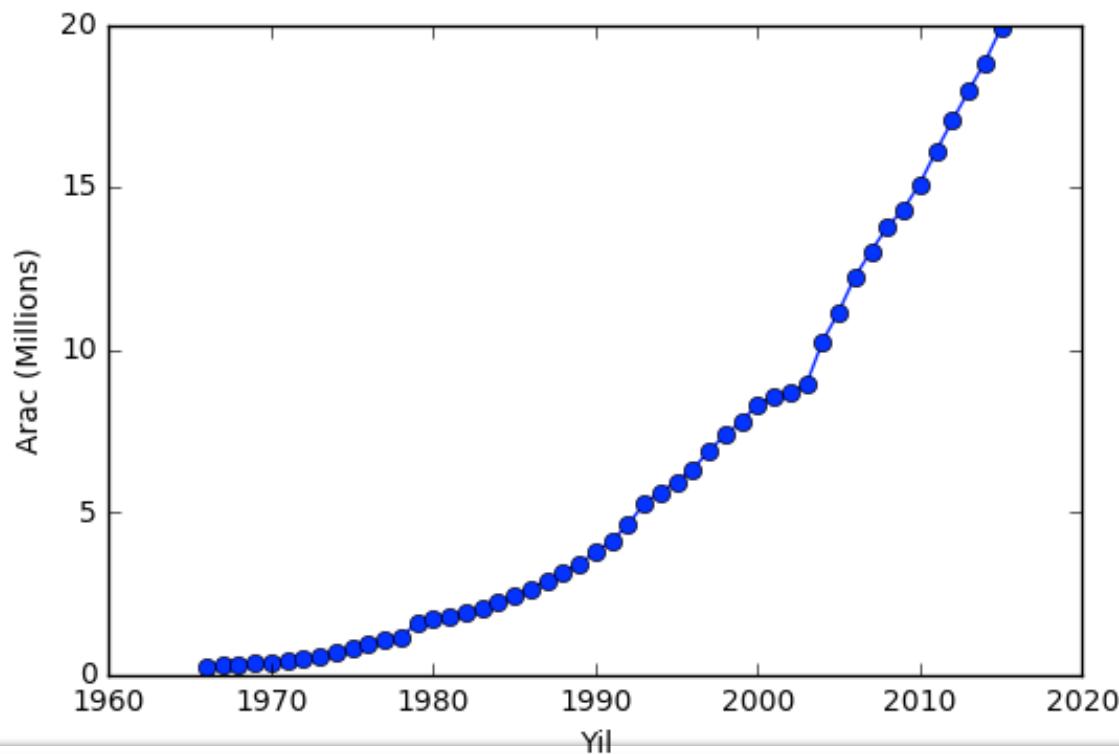
Algorithm 1: quad

```
input:  $z_0, \beta_1$ 
1 for  $t = 0, 1, 2, \dots$  do
2    $x_1 = z_t$ 
3   Compute  $H_t$ 
4   for  $k = 1, 2, \dots, c$ 
5     Choose a subs
6     Compute  $\nabla_{S_k}$ 
7      $x_{k+1} = \arg \min$ 
8   end
9    $z_{t+1} = x_{c+1}$ 
10  Set  $\beta_{t+1} \leq \beta_t$ 
11 end
```



Güdümlü Öğrenme: Regresyon

i	Araç Sayısı (y)	Yıl (x)
1	231977	1966
...		
49	19882069	2015



Güdümlü Öğrenme: Regresyon

i	Araç Sayısı (y)	Yıl (x)
1	231977	1966
...		
49	19882069	2015

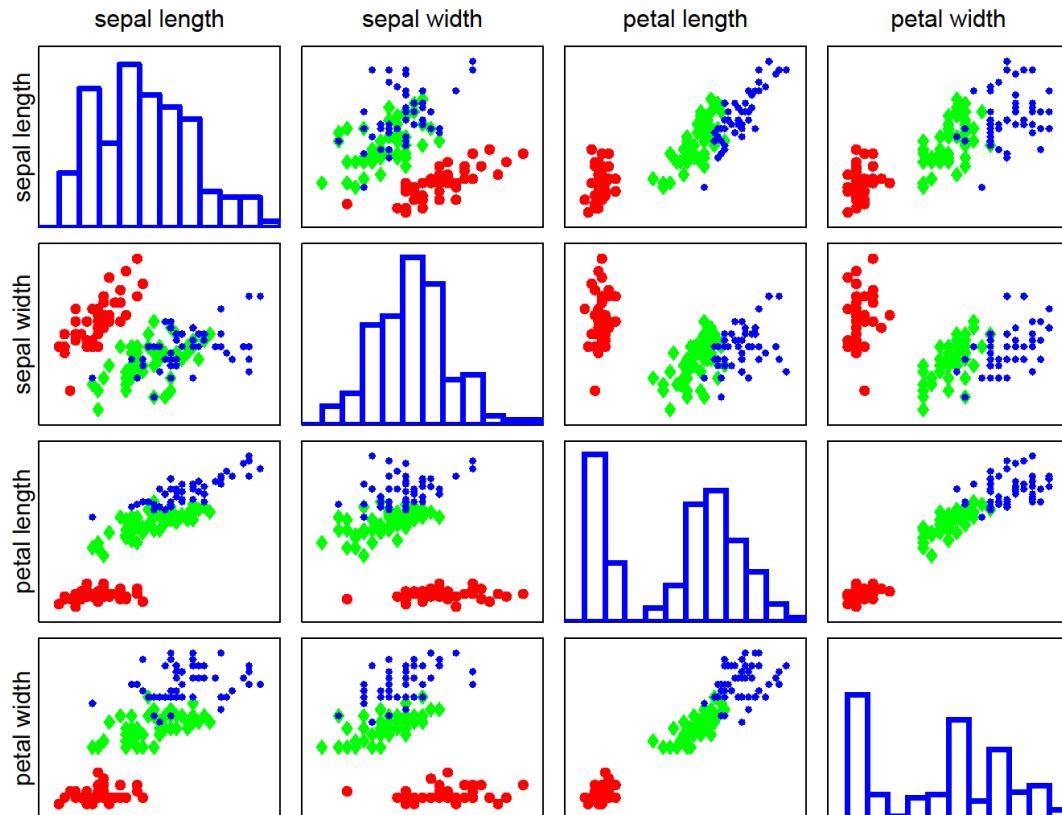
$$y \approx w_1 x + w_0$$

$$y \approx w_2 x^2 + w_1 x + w_0$$

$$y \approx f(x; w)$$

Güdümlü Öğrenme (Supervised Learning)

Sınıflandırma

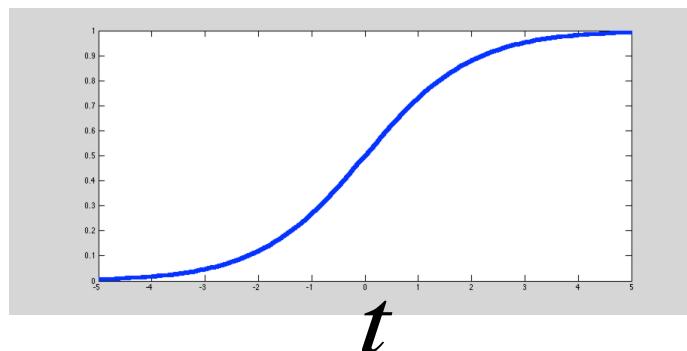


Sınıflandırma: Lojistik Regresyon

Öznitelik 1	Öznitelik 2	Öznitelik 3	Öznitelik 4	Sınıf
5.1	4.3	2.1	0.3	0
5.7	3.5	3.2	0.8	0
3.4	5.2	0.4	0.6	1
X ₁	X ₂	X ₃	X ₄	y

$$y \approx \sigma(x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4)$$

$\sigma(t)$



Öznitelik Mühendisliği (Feature Engineering)

- Belirli bir problem için uygun sayısal bir gösterim bulmak
- Örn: e-mail Spam/NoSpam filtering
 - X_1 = [email metni 'Rolex' içeriyor mu?]
 - X_2 = [email metni 'http://' içeriyor mu?]
 - X_3 = email metindeki büyük/küçük harf sayıları oranı
 - ...
 - X_{100000} = [Gönderen kişi adres defterinde var mı?]

Risk Kestirimimi(Probability of Default PD)

Öznitelikler

Demografik değişkenler

Gelir düzeyi

Müşteri olma süresi ...

Macroekonomik değişkenler

İşsizlik katsayısı ...

Davranışsal değişkenler (Müşteriye özel)

İşlem sayısı

Geri ödeme gecikmesi ...

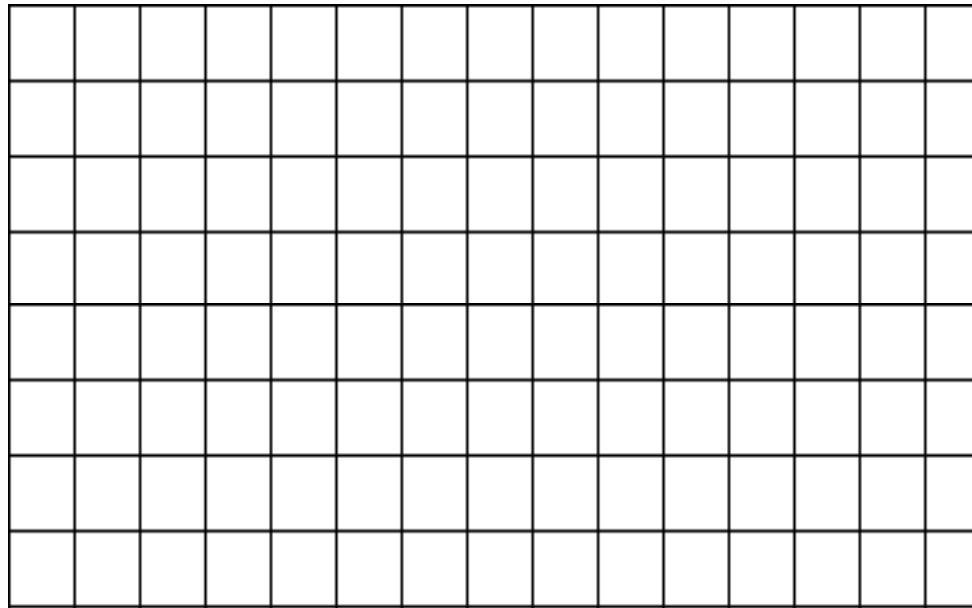
Time customer with bank (years)	-0.00250**
Time with bank unknown +	-0.342**
Income (log)	-0.146**
Income unknown +	-1.46**
Number of cards	-0.0610**
Time at current address	-0.00129
Employment + :	
Self-employed	0.303**
Homemaker	0.072
Retired	0.111
Student	-0.035
Unemployed	0.231
Part time	-0.365**
Other	-0.037
<i>Excluded category:</i> Employed	
Age + : 18 to 24	0.074
25 to 29	-0.058
30 to 33	0.010

Büyük Ölçekte sınıflandırma

- Reklam Tahmini, YouTube video sıralaması
- Bir kullanıcının bir reklamı/videoyu klikleme ihtimali nedir?
- A Reliable Effective Terascale Linear Learning System, Agarwal et.al. 2012

Öznitelikler = 16 M

Örnek sayısı 17 Milyar



3TB Veri
kümesi
1000 Makina

Dağıtık Algoritmalar

1. Her düğümde sıralı öğrenme kullanarak bir parametre bul
2. AllReduce kullanarak ortalama hesapla
3. Her düğümde, gradyan toplamlarını hesapla
4. AllReduce kullanarak gradyanları topla.
5. L-BFGS kullanarak parametreleri güncelle ve 3. adıma dön

Yapay Sinir Ağları, Derin Öğrenme

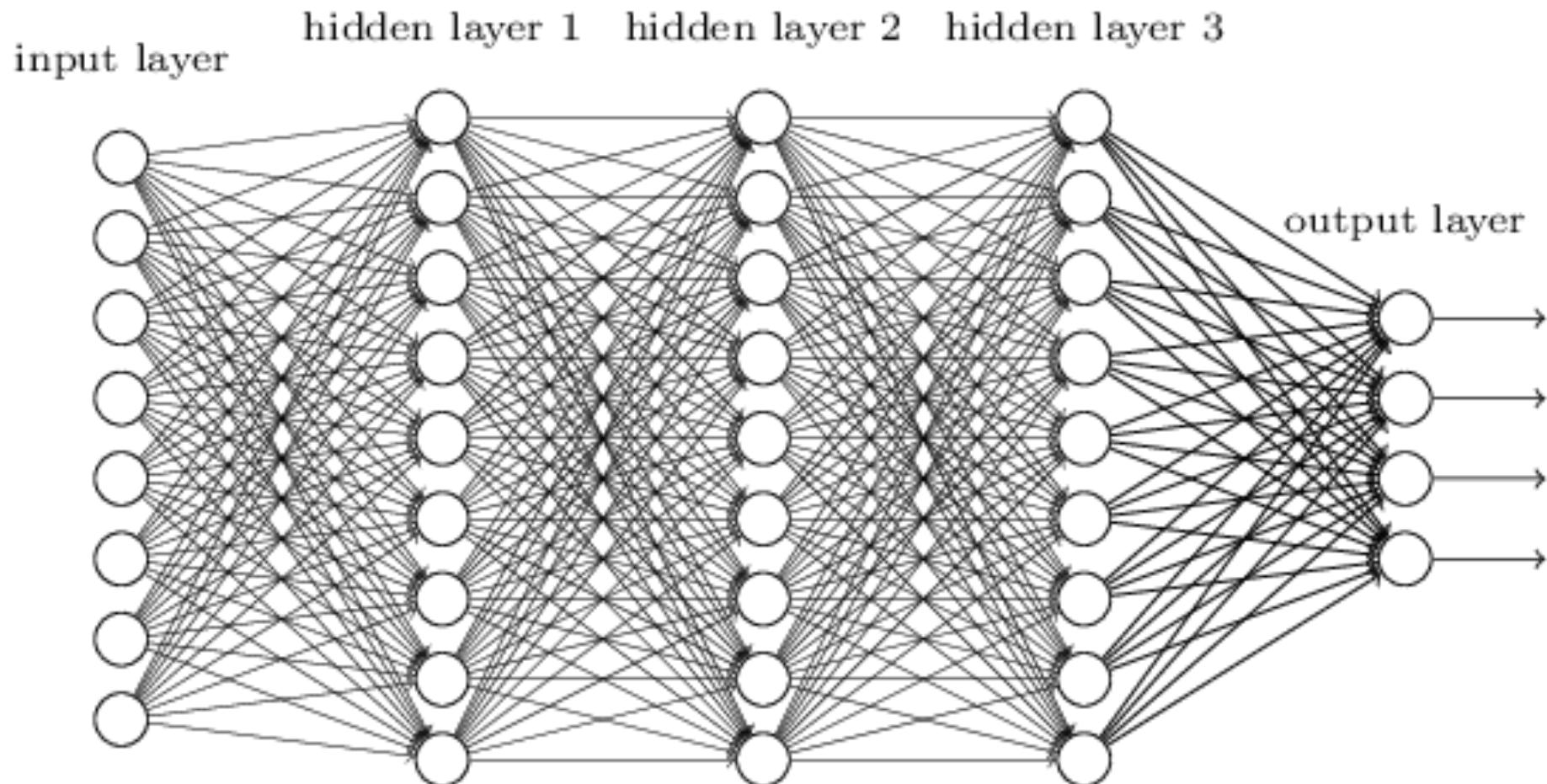


Figure:<http://neuralnetworksanddeeplearning.com/chap5.html>

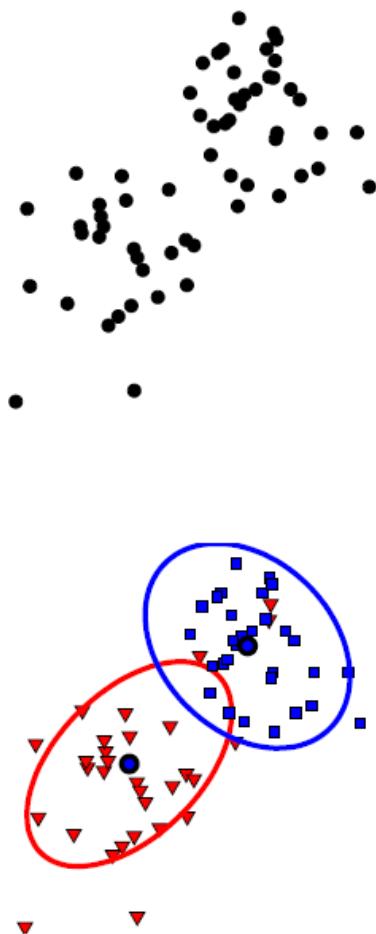
Derin Öğrenme araçları

- **Theano** CPU/GPU symbolic expression compiler in python
 - <http://deeplearning.net/software/theano/index.html>
- **Caffe** deep learning framework
- **Torch** Matlab-like environment for state-of-the-art machine learning algorithms in lua
- **Tensorflow** open source software library for numerical computation using data flow graphs
- ... and many others

Güdümsüz Öğrenme

- Öbekleme (Clustering)
- Boyut Düşürme (Dimensionality Reduction)
- Göreselleştirme (Visualization)

Öbekleme



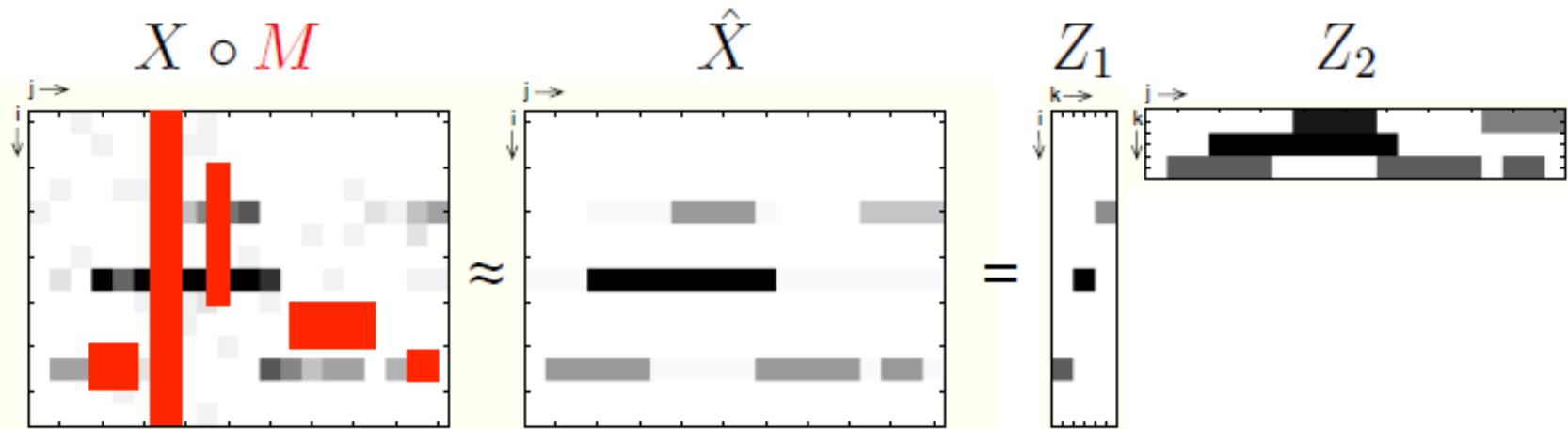
Boyut Düşürme

■ Kelime-Kitap

	j.caesar	hamlet	othello	macbeth	rom&jul	sonnets
caesar	270	2	1	1	0	0
brutus	379	1	0	0	0	0
malcolm	0	0	0	60	0	0
muse	0	0	1	1	0	16
:						
love	34	68	80	19	150	195
friend	23	14	18	5	13	16
the	610	1148	759	733	682	446
traitor	1	0	0	5	1	0
traitors	9	0	1	3	0	0
:						
napkin	0	1	3	0	0	0
sword	15	16	10	14	8	1
laptop	0	0	0	0	0	0

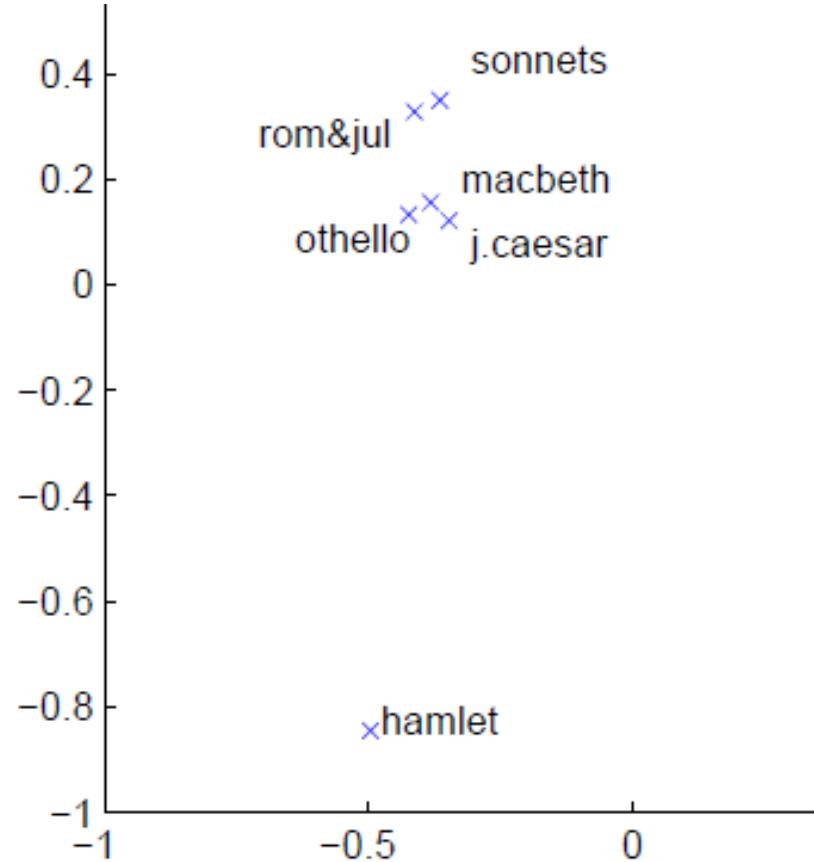
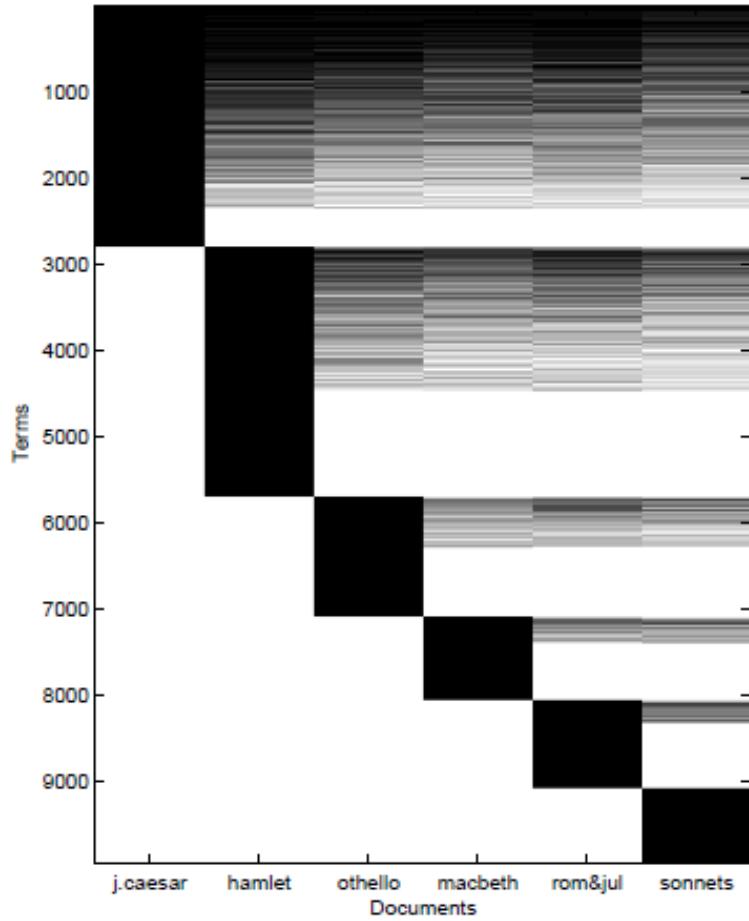
Matrix Factorizations

$$X(i, j) \approx \sum_k Z_1(i, k)Z_2(k, j)$$



$$(Z_1, Z_2)^* = \arg \min_{Z_1, Z_2} D(X || Z_1 Z_2) + \lambda R(Z_1, Z_2)$$

Kelime Dokuman Matrisi



Yapay Öğrenme ve Olasılık Teorisi

- Olasılık Teorisi
 - Olasılık teorisi, sağduyunun hesaplamaya dökülmüş halinden başka bir şey değildir. – P. Laplace
 - Probability theory is nothing but common sense reduced to calculation – P. Laplace
- Grafik Modeller, Olasılıksal Uzman Sistemler
- Zaman Serileri
 - Örnek: 'Network flow' sınıflandırma

Tavsiye Sistemleri



Matris tamamlama

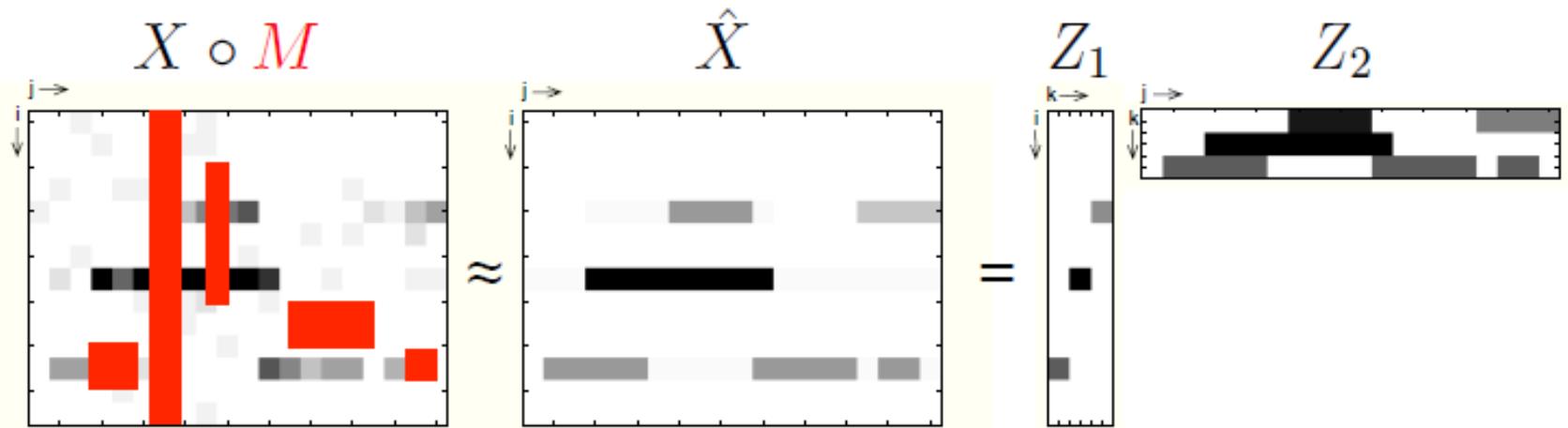
- Netflix: 18K film \times 500K kullanıcı %99 seyrek

The diagram illustrates a sparse matrix representation for movie recommendation. It consists of a grid of 6 rows and 6 columns. The columns are labeled "users" at the top, and the rows are labeled "movies" on the left. A vertical double-headed arrow on the left indicates the rows represent movies, and a horizontal double-headed arrow at the top indicates the columns represent users.

1		?	3	5	?
?	1				2
	4		4	5	?

Matris ve Tensor Ayrıştırma

$$X(i, j) \approx \sum_k Z_1(i, k)Z_2(k, j)$$



$$(Z_1, Z_2)^* = \arg \min_{Z_1, Z_2} D(X || Z_1 Z_2) + \lambda R(Z_1, Z_2)$$

Tavsiye Sistemleri

	1	?	3	4
	2	4	6	8
	1.5	3	?	6.1

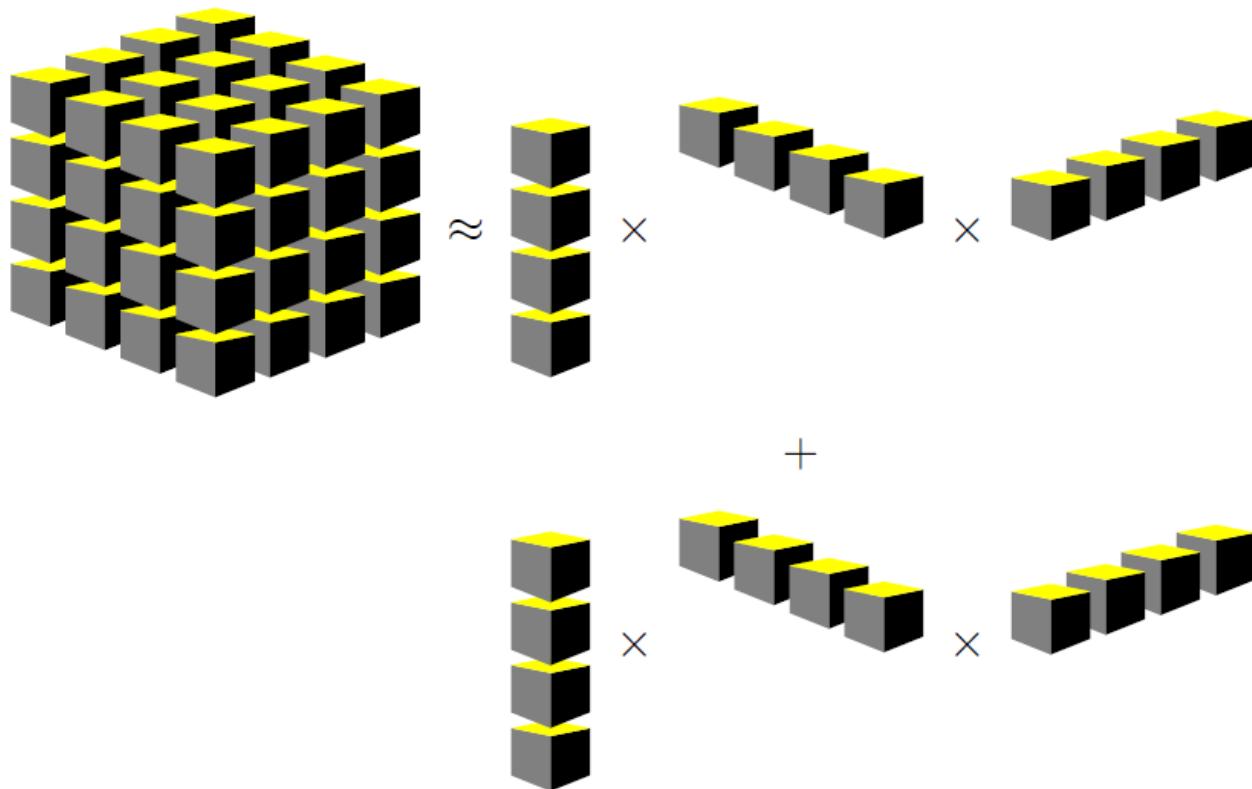
Tavsiye Sistemleri: Öğrenme

	1	2	3	4
1	1	?	3	4
2	2	4	6	8
1.5	1.5	3	?	6.1

Tavsiye Sistemleri

	1	2	3	4
1	1	2	3	4
2	2	4	6	8
1.5	1.5	3	4.5	6.1

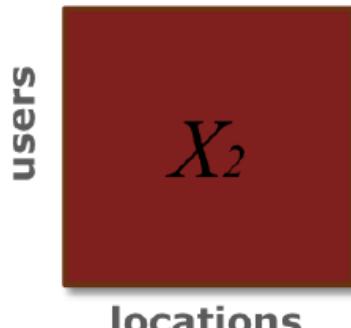
Tensor Ayrıştırma



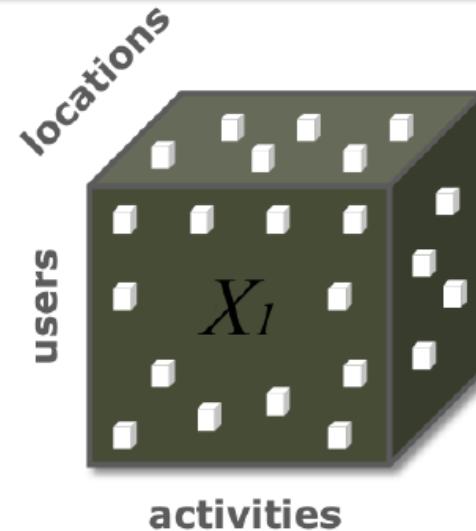
$$X(i, j, k) \approx \sum_r Z_1(i, r) Z_2(j, r) Z_3(k, r)$$

- Recommendation : User \times Location \times Activity
- EEG : Channel \times Frequency \times Time \times Subject
- Music : Channel \times Frequency \times Time \times Song
- Text : Documents \times Terms \times Language
- fMRI : Voxels \times Run \times Time \times Trial
- Movies : Movies \times Users \times Time
- Email : Users \times Messages \times Time

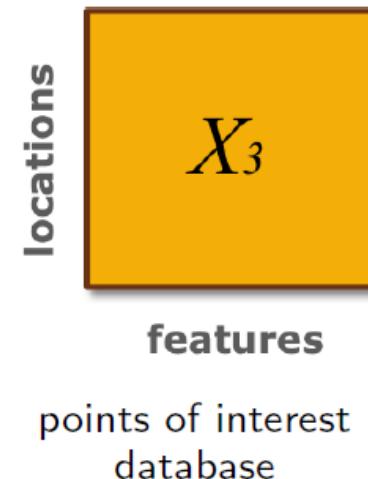
Veri Birleştirme (Sensor Fusion)



the user-location from GPS trajectory



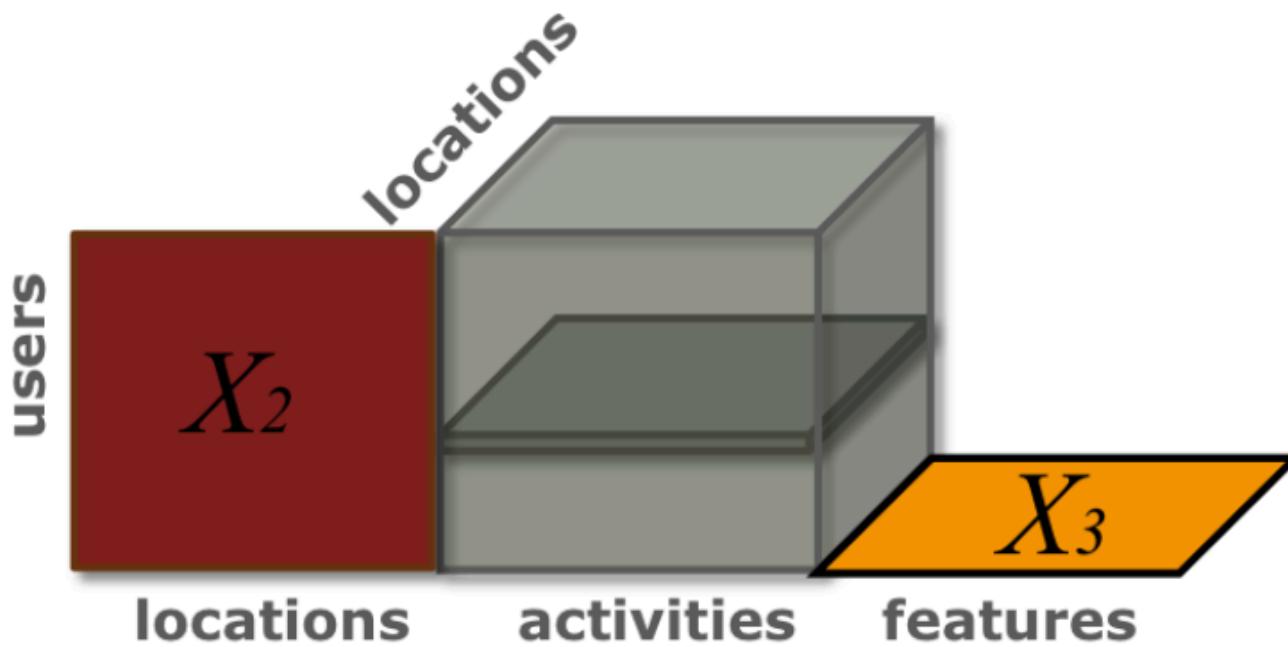
user-location-activity data
(activity: 'Food', 'Shopping',
'Movie', 'Sports',
'Amusement')



points of interest
database

- Mahremiyet gözeten yapay öğrenme
 - Dağıtık Sistemlerde öğrenme
 - Veri değil parametre paylaşımı

Cold Start



- Coupled models outperform low-rank approximation models

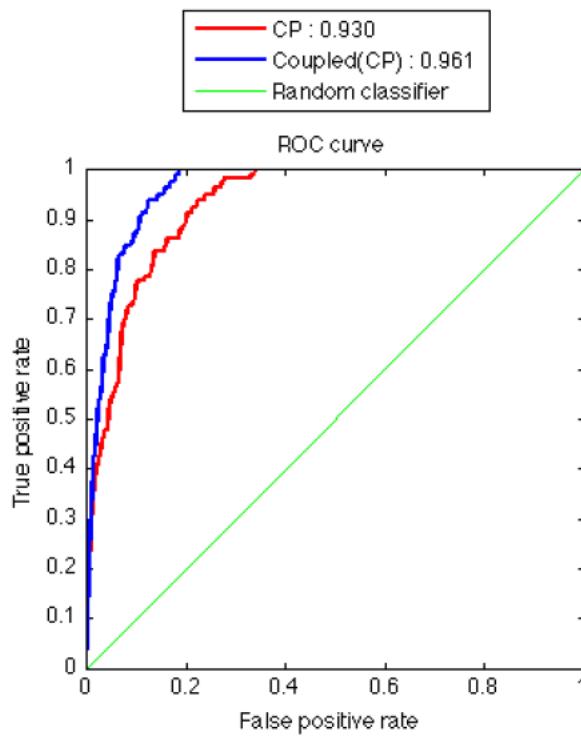


Figure 1: 40% missing

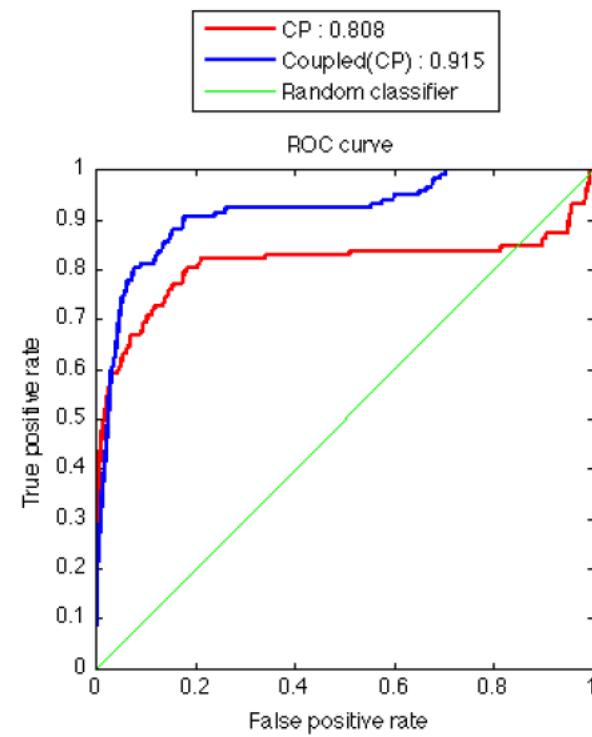


Figure 2: 80% missing

Paralel İşleme Platformları (BBL2011)

Platform	Communication Scheme	Data size
Peer-to-Peer	TCP/IP	Petabytes
Virtual Clusters	MapReduce / MPI	Terabytes
HPC Clusters	MPI / MapReduce	Terabytes
Multicore	Multithreading	Gigabytes
GPU	CUDA	Gigabytes
FPGA	HDL	Gigabytes

Slide from ICML 2011 tutorial Langford et. al.

Büyük Veri Platformları (Açık Kaynak)

- Hadoop/MapReduce



- Apache Spark



- ...



Özet

- Veri ≠ Bilgi
- Algoritma tasarıımında yeni bakış açıları
- Büyük veri: veri noktaları arasındaki ilişkiler ve etkileşimler
- Yapay Öğrenme : bir çok uygulamada olgun teknoloji
- Bilgisayar bilimleri eğitimi:
 - Daha çok Matematik, Fizik, İstatistik ve Sosyal Bilimler etkileşimi
- Büyük veri = Büyük Potansiyel

Teşekkürler, Sorular

