



**ICT4SS**  
Information and Communications Technologies for Smart Societies

# Data science process

*Prof. Marco Mellia*



POLITECNICO  
DI TORINO

SmartData@Polito





POLITECNICO  
DI TORINO

# Big data hype?

**ICT4SS**  
Information and Communications Technologies for Smart Societies

## Big Data & Data Science

“Extracting meaning from very large quantities of data”

**ICT4SS**  
Information and Communications Technologies for Smart Societies

## Emergency management

EARTH OBSERVATIONS

UNMANNED AERIAL

HISTORICAL DATA

SEASONAL WEATHER FORECAST

SOCIAL MEDIA DATA STREAMS

FLOODS

3D Model

Improving Resilience to Emergencies Through Advanced Cyber Technologies

**iREACT**



The slide features the ICT4SS logo at the top right, which includes a stylized city skyline icon and the text "ICT4SS Information and Communications Technologies for Smart Societies". The main title "Emergency management" is centered above a diagram. The diagram shows a double-headed arrow connecting "FIRST RESPONDERS AND DECISION MAKERS" (represented by two people icons) and a server room image. Another double-headed arrow connects this central pair to "CITIZENS" (represented by two people icons holding phones) and a smartphone displaying a map application. Below the diagram are two screenshots of software interfaces: one showing a map with red lines and another showing a street scene with workers. The text "Improving Resilience to Emergencies Through Advanced Cyber Technologies" is at the bottom right, next to the iREACT logo.



The slide features the ICT4SS logo at the top right. The main title "User engagement" is centered above a comparison of two crowd scenes. On the left, a red box labeled "2005" contains a photograph of a large crowd of people at night, many of whom are looking towards a bright stage or screen. On the right, a red box labeled "2013" contains a photograph of a similar large crowd, but almost every person is holding up a smartphone or tablet to take a picture or video. The ICT4SS logo is located in the top right corner of the slide area.

**Who generates big data?**

• User Generated Content (Web & Mobile)

- E.g., Facebook, Instagram, Yelp, TripAdvisor, Twitter, YouTube

• Health and scientific computing

**ICT4SS**  
Information and Communications Technologies for Smart Societies

**Who generates big data?**

• Log files

- Web server log files, machine syslog files

• Internet Of Things

- Sensor networks, RFID, smart meters

**ICT4SS**  
Information and Communications Technologies for Smart Societies

## What is big data?

- Many different definitions

*"Data whose scale, diversity and complexity require new architectures, techniques, algorithms and analytics to manage it and extract value and hidden knowledge from it"*



## What is big data?

- Many different definitions

*"Data whose **scale, diversity** and **complexity** require new architectures, techniques, algorithms and analytics to manage it and extract value and hidden knowledge from it"*



## What is big data?

- Many different definitions

*"Data whose scale, diversity and complexity require new **architectures, techniques, algorithms** and **analytics** to manage it and extract value and hidden knowledge from it"*

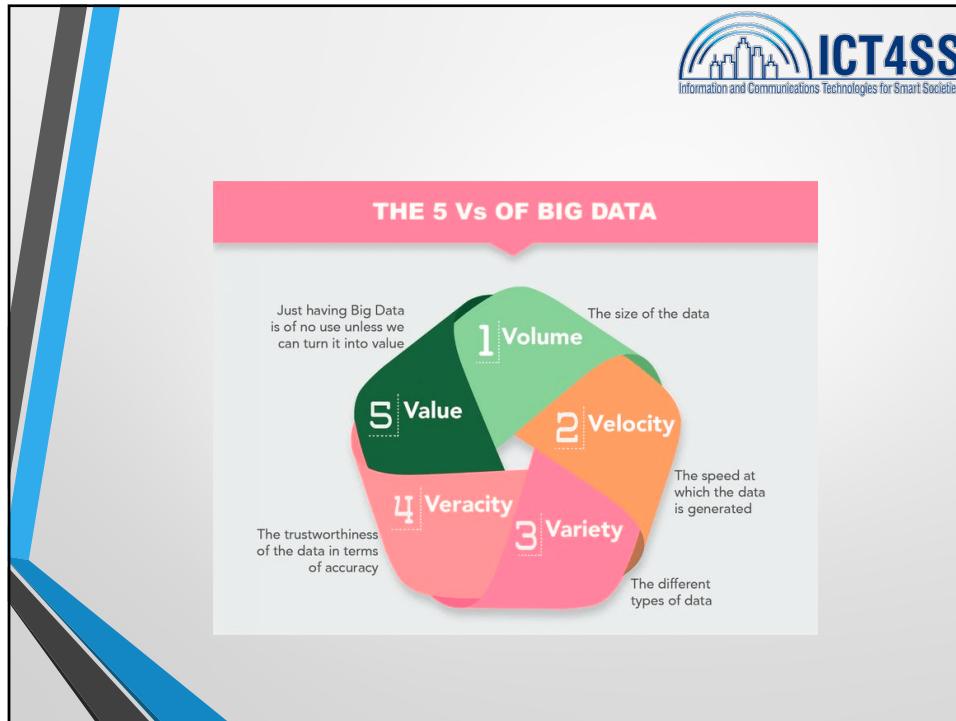


## What is big data?

- Many different definitions

*"Data whose scale, diversity and complexity require new architectures, techniques, algorithms and analytics to manage it and extract **value** and hidden **knowledge** from it"*





## The Vs of big data: Volume

- Data volume increases exponentially over time
- 44x increase from 2009 to 2020
  - Digital data 35 ZB in 2020

terabytes | petabytes | exabytes | zettabytes

the amount of data stored by the average company today

**The Digital Universe 2009-2020**

Growing By A Factor Of 44

2009: 0.8 Zb

2020: 35.2 Zettabytes

**A TIDAL WAVE OF DATA**

Twitter: 8 TB per day

Facebook: 100 TB per day

Global Businesses: 1.8 Zb In 2011

US Library of Congress: 235 TB

Boeing: 640 TB per flight

McKinsey: 2.5 PB stored

**TWEETS PER DAY**

Number of tweets per day

Date

2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020

Exponential

**ICT4SS**  
Information and Communications Technologies for Smart Societies

**The Vs of big data: Velocity**

ICT4SS  
Information and Communications Technologies for Smart Societies

- Fast data generation rate
  - Streaming data
- Very fast data processing to ensure timeliness

The New York Stock Exchange captures **1 TB OF TRADE INFORMATION** during each trading session.

Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure.

By 2016, it is projected there will be **18.9 BILLION NETWORK CONNECTIONS** – almost 2.5 connections per person on earth.

**Velocity ANALYSIS OF STREAMING DATA**

**ONE SECOND ON THE INTERNET**

- 100,000 tweets
- 571 websites created
- 47,000 checkins posted online
- App
- g 2450 searches
- YouTube 3,000 photos shared
- \$ 1,000 transactions
- Facebook 684,478 shared items

<http://www.internetlivestats.com/>

**Real time processing**

Sensing

Wireless Sensor Networks

Computing

Crowdsourcing

Map data

Real time traffic info

**The Vs of big data: Variety**

ICT4SS  
Information and Communications Technologies for Smart Societies

- Various formats, types and structures
  - Numerical data, image data, audio, video, text, time series

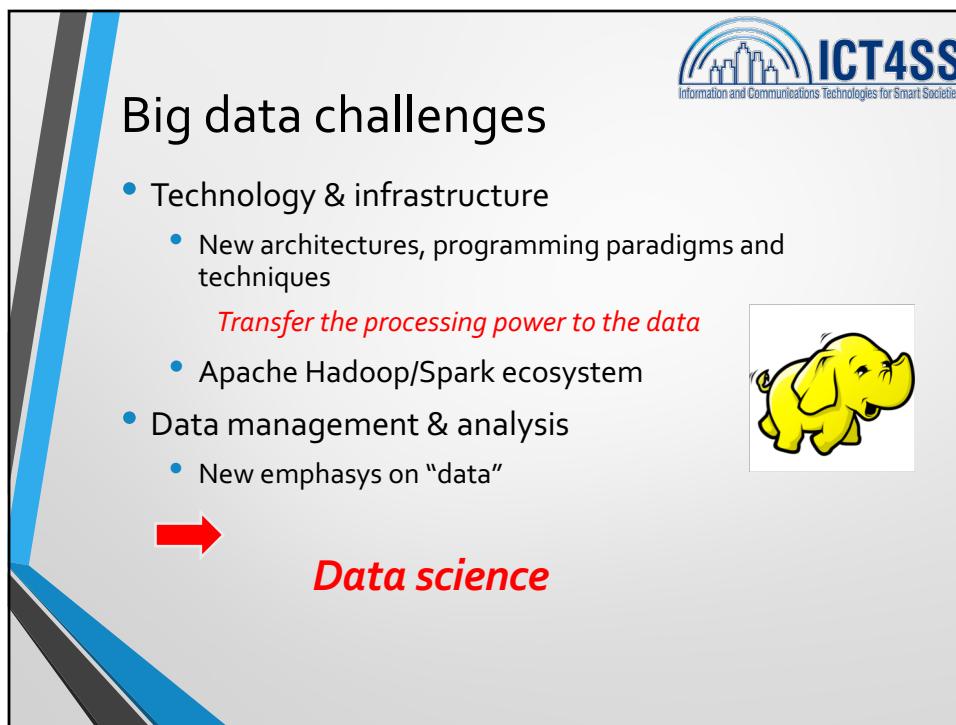
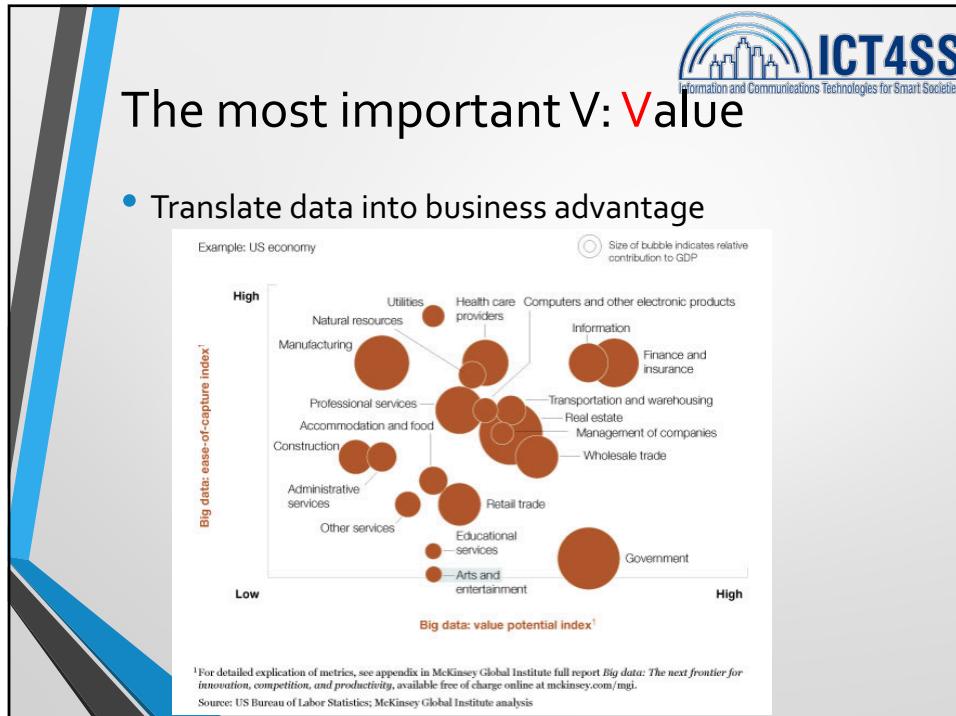
- A single application may generate many different formats

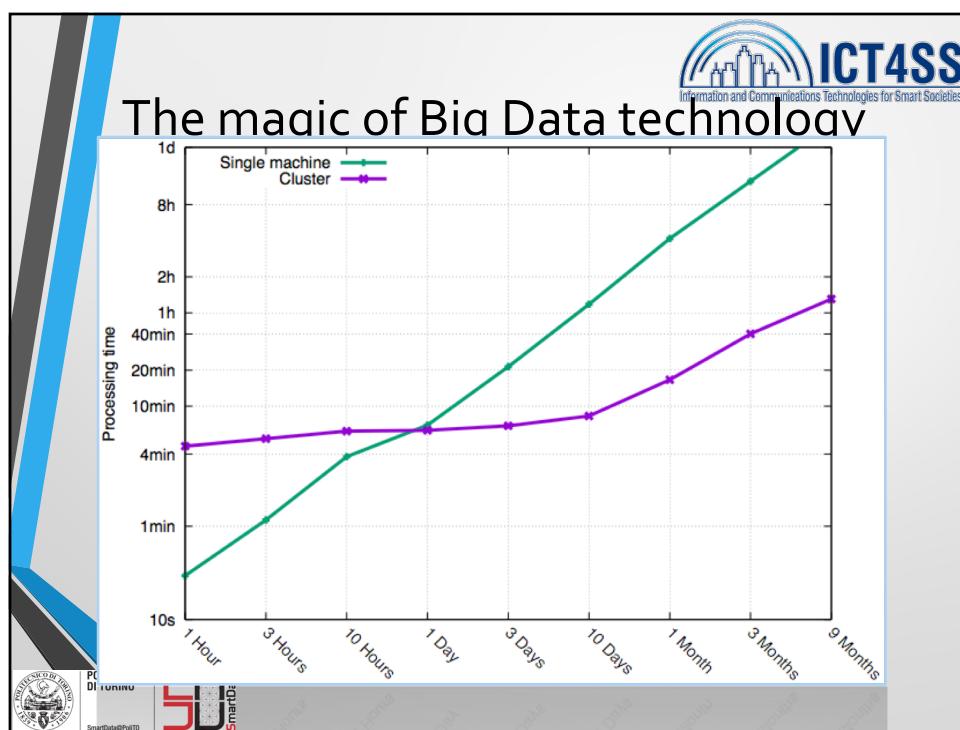
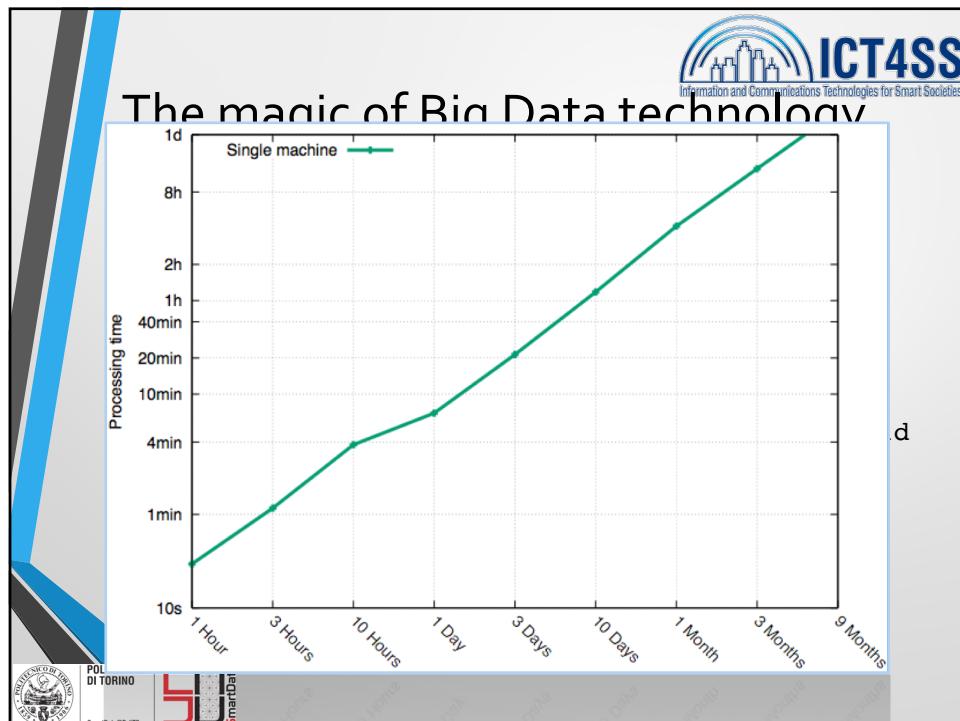
**The Vs of big data: Veracity**

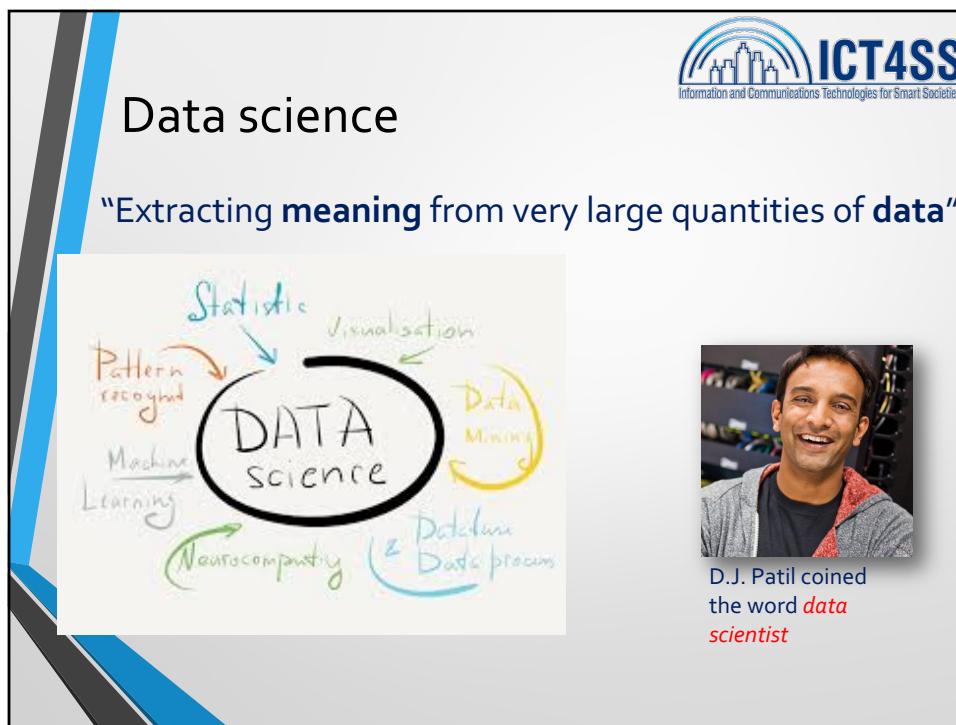
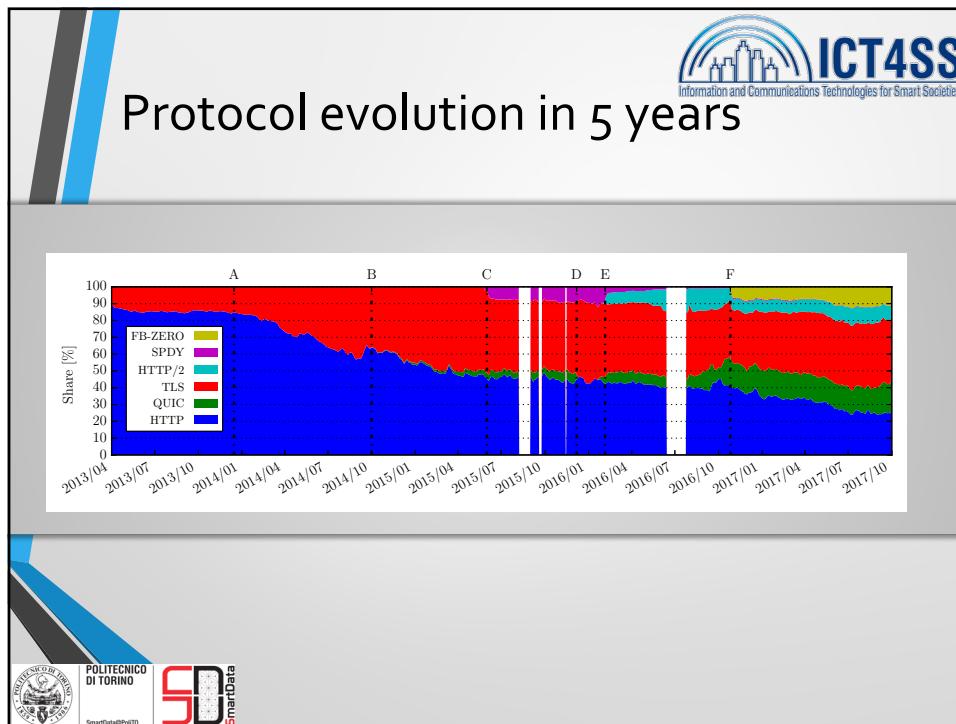
ICT4SS  
Information and Communications Technologies for Smart Societies

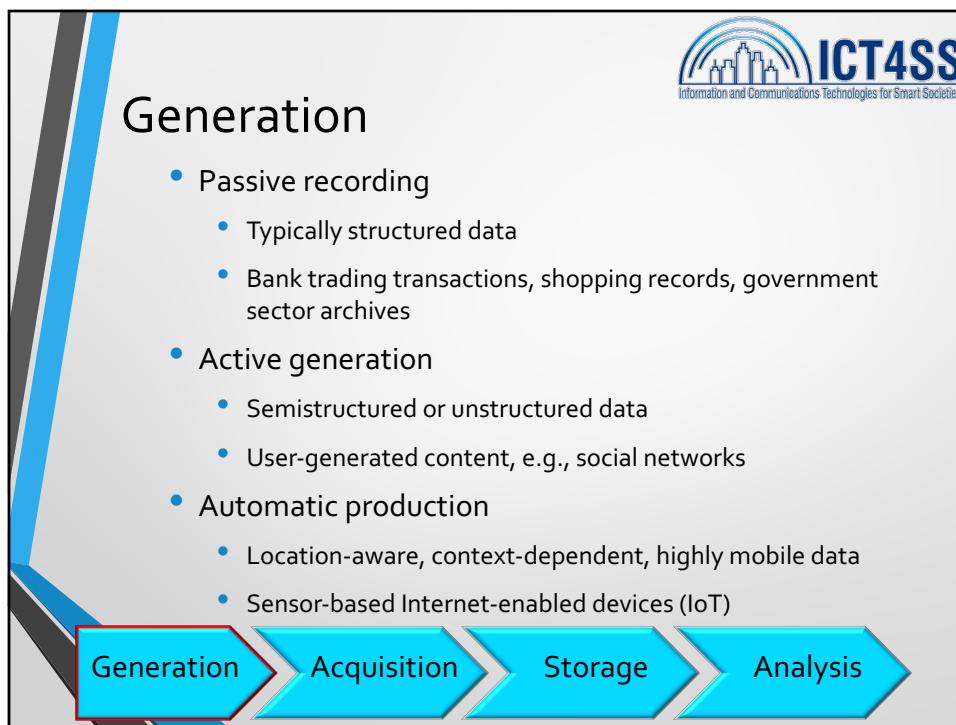
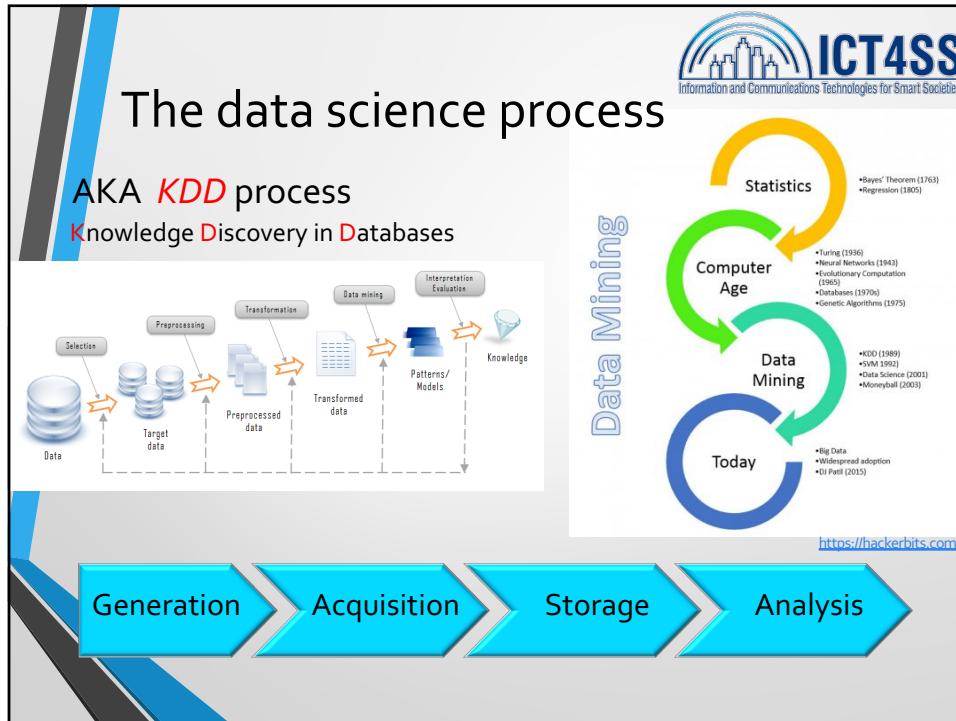
- Data quality

**Reliability**  
Format  
Sufficiency  
Conciseness  
**Timeliness**  
Flexibility  
Completeness  
Consistency  
Level-of-detail  
Informativeness  
**Accuracy**  
Currency  
Precision  
Efficiency  
Quantitableness  
Usefulness  
Usableness  
Clarity Content  
**Relevance**  
Comparability  
Scope  
Interpretability  
Importance











## Acquisition

- Collection
  - Pull-based, e.g., web crawler
  - Push-based, e.g., video surveillance, click stream
- Transmission
  - Transfer to data center over high capacity links
- Preprocessing
  - Integration, cleaning, redundancy elimination



```

graph LR
    G[Generation] --> A[Acquisition]
    A --> S[Storage]
    S --> A
    A --> A[Analysis]
  
```



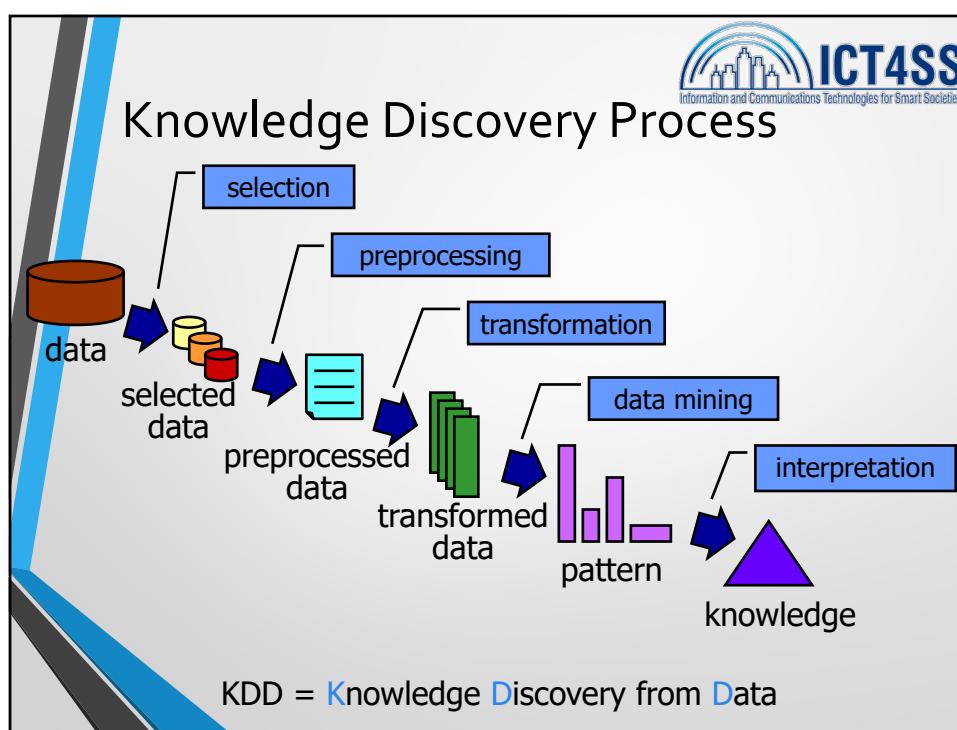
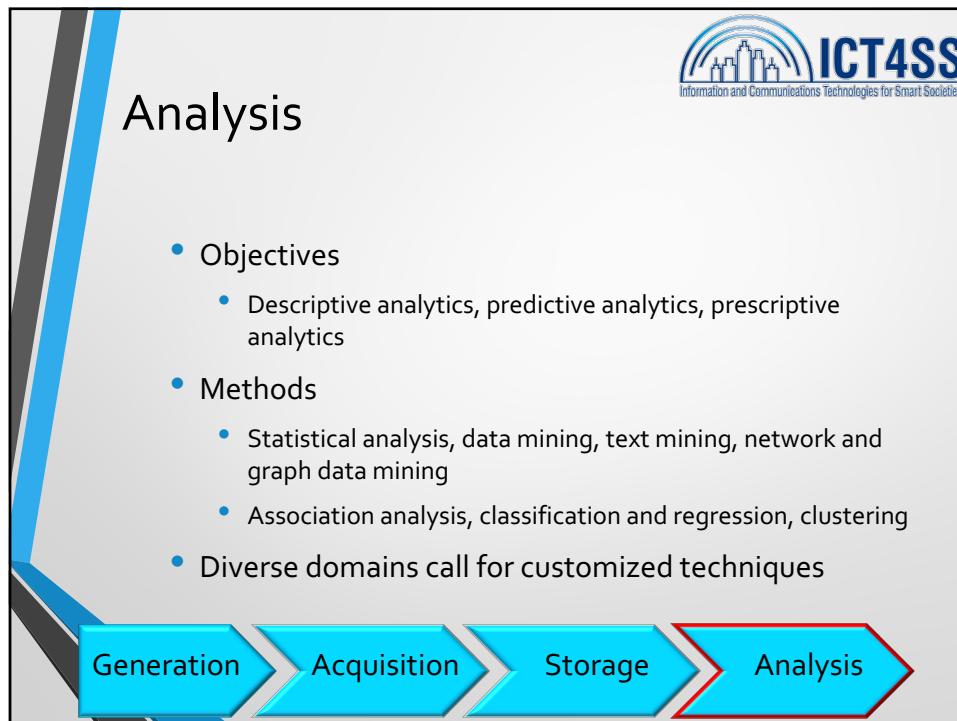
## Storage

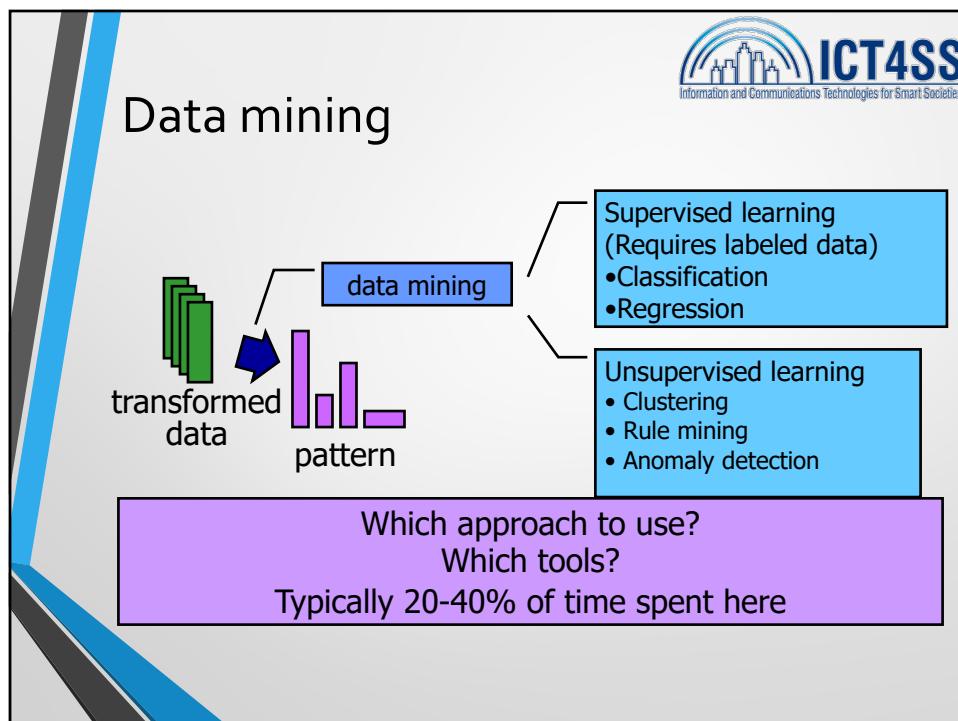
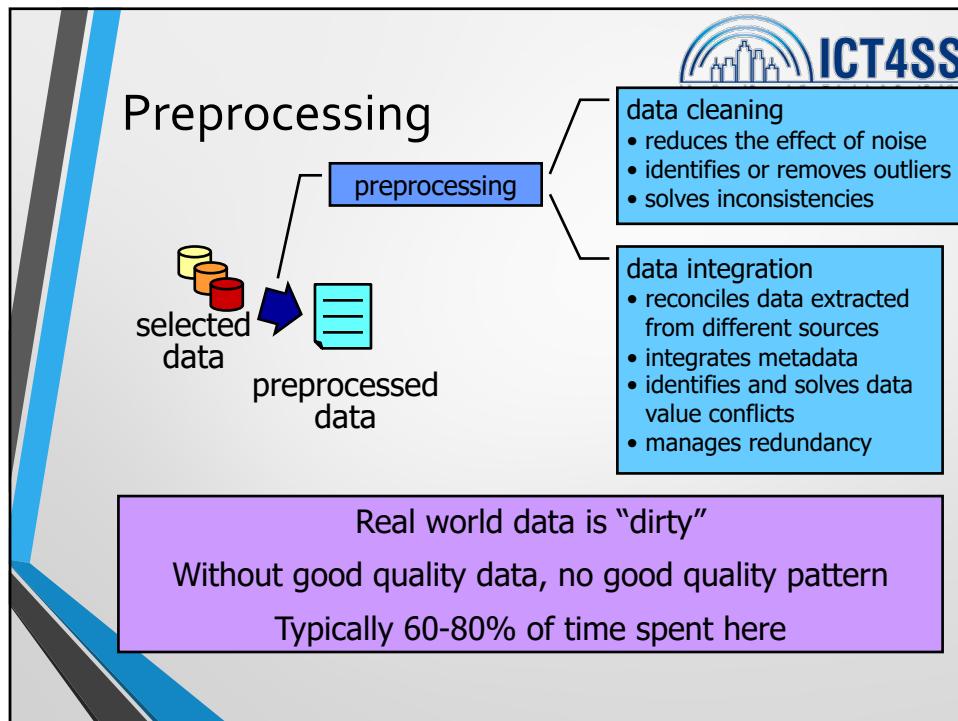
- Storage infrastructure
  - Storage technology, e.g., HDD, SSD
  - Networking architecture, e.g., DAS, NAS, SAN
- Data management
  - File systems (HDFS), key-value stores (Memcached), column-oriented databases (Cassandra), document databases (MongoDB)
- Programming models
  - Map reduce, stream processing, graph processing



```

graph LR
    G[Generation] --> A[Acquisition]
    A --> S[Storage]
    S --> A
    A --> A[Analysis]
  
```





**Classification**

ICT4SS  
Information and Communications Technologies for Smart Societies

- Objectives
  - prediction of a class label
  - definition of an interpretable model of a given phenomenon

The diagram shows a flow from 'training data' (represented as a grid of colored squares) through a 'model' (represented by a purple arrow) to 'classified data' (another grid of colored squares). Below the model, 'unclassified data' (also a grid of colored squares) is shown being processed by the model to produce three output categories represented by colored arrows pointing to separate grids.

**Classification techniques**

ICT4SS  
Information and Communications Technologies for Smart Societies

- Decision trees
- Classification rules
- Association rules
- Neural Networks
- Naïve Bayes and Bayesian Networks
- k-Nearest Neighbours (k-NN)
- Support Vector Machines (SVM)
- ...

The diagram includes three main parts: a decision tree for survival analysis, a neural network diagram showing layers and connections, and a CNN diagram showing input, convolution, pooling, and fully connected layers processing an image of a boat.

**Classification**



Depth Image

180 x 154 x 64

C1 Layer: 58 convolutions, 1x1 kernel, stride 2, max pooling

P1/LCN1 Layer: 72 convolutions, 3x3 kernel, stride 2, max pooling

C2 Layer: 24 convolutions, 3x3 kernel, stride 2, max pooling

P2/LCN2 Layer: 34 convolutions, 3x3 kernel, stride 2, max pooling

C3 Layer: 10 convolutions, 3x3 kernel, stride 2, max pooling

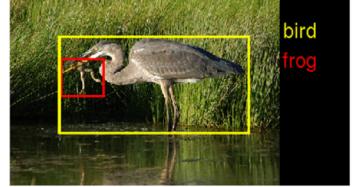
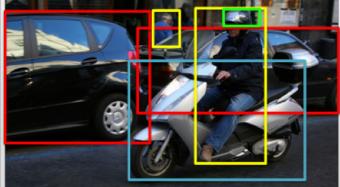
P3/LCN3 Layer: 64 convolutions, 3x3 kernel, stride 2, max pooling

F1: 16 convolutions, 3x3 kernel, stride 2, max pooling

F2: 3 convolutions, 3x3 kernel, stride 2, max pooling

Output Labels: 3

Feature Extraction

bird  
frog

Person  
Car  
Motorcycle  
Helmet

**Evaluation of classification techniques**

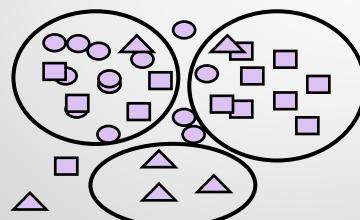


- Accuracy
  - quality of the prediction
- Interpretability
  - model interpretability
  - model compactness
- Robustness
  - noise, missing data

- Efficiency
  - model building time
  - classification time
- Scalability
  - training set size
  - attribute number

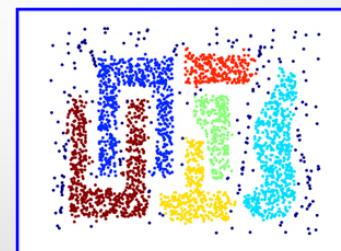
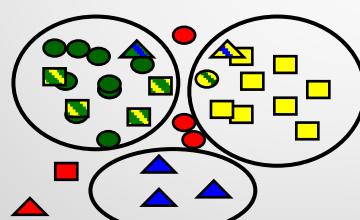
## Clustering

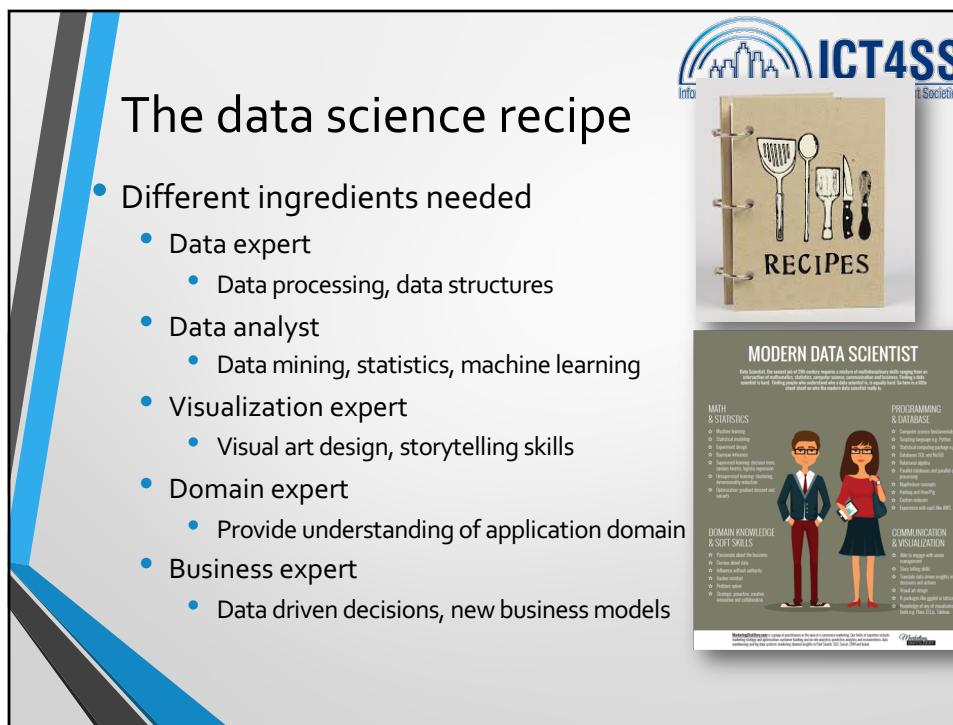
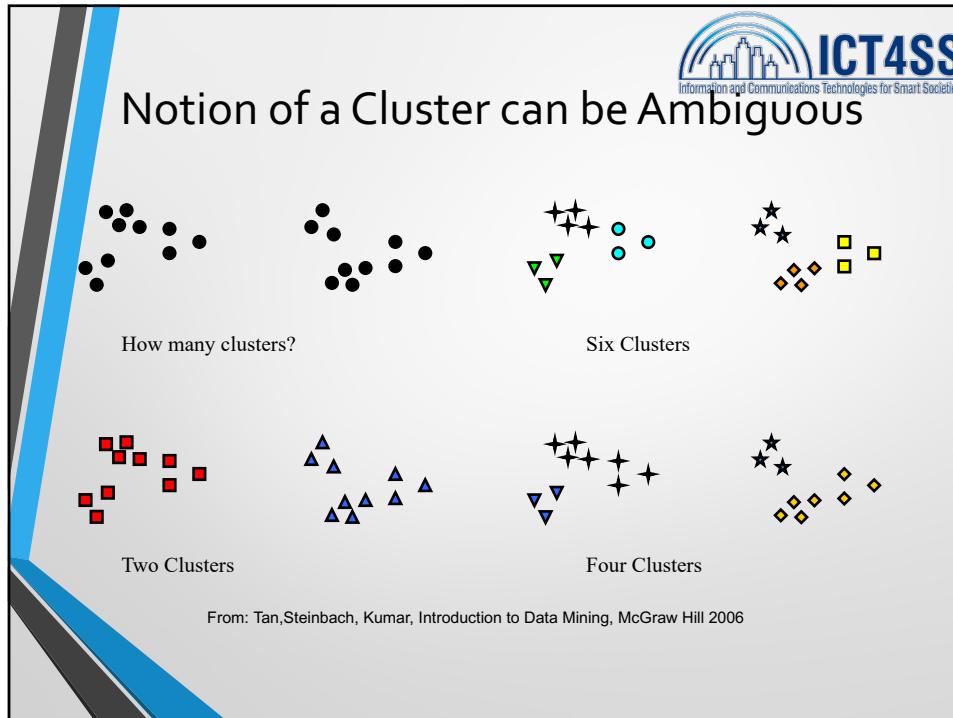
- Objectives
  - detecting groups of similar data objects
  - identifying exceptions and outliers

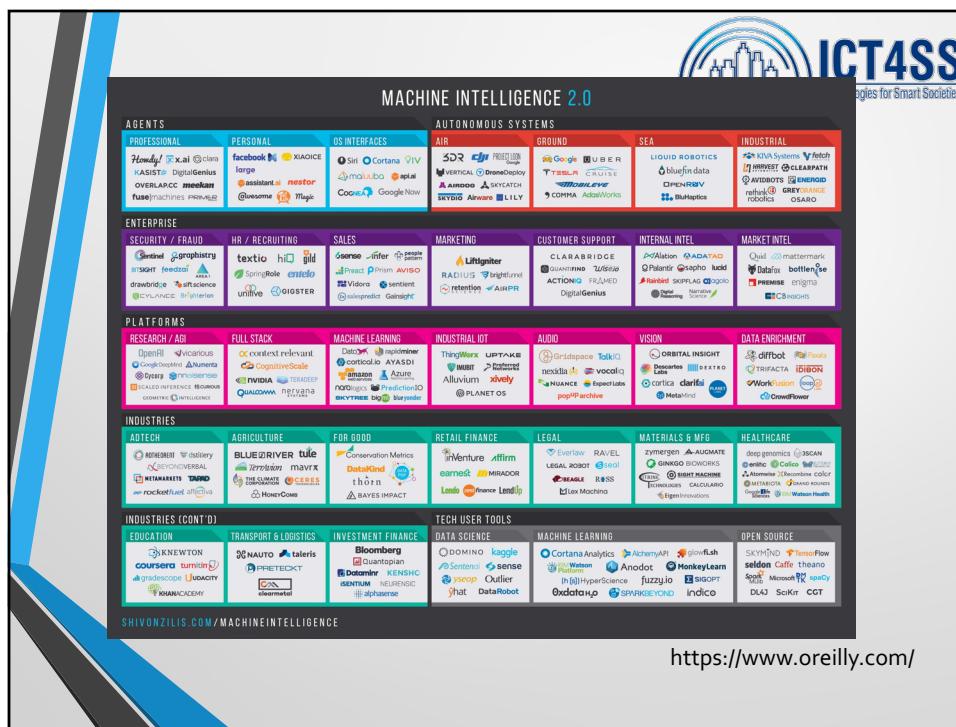
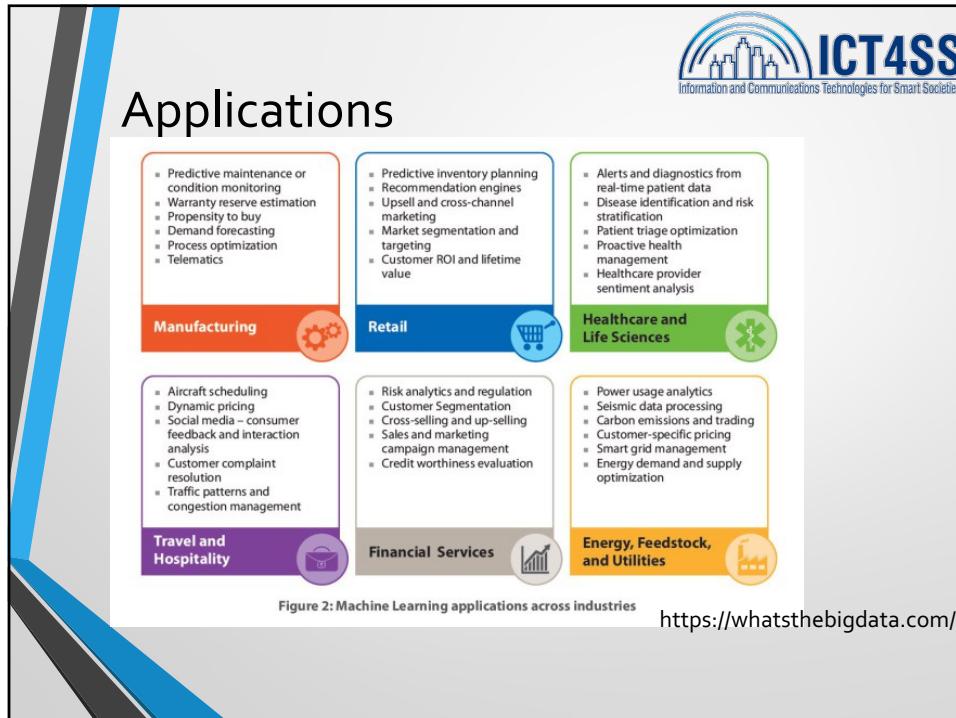


## Clustering

- Objectives
  - detecting groups of similar data objects
  - identifying exceptions and outliers







**ICT4SS**  
Information and Communications Technologies for Smart Societies

## Conclusions and open issues

- We are leaving in a data deluge epoch
  - Getting data is easier and easier
- **Extracting value from the data is still a challenge**
  - Don't blindly trust the magician – no Harry Potter stick here

The Spells of Harry Potter scatter plot shows various magical spells plotted against parameters like Complexity, Power, and Duration. Harry Potter is shown in the center of the plot.

**ICT4SS**  
Information and Communications Technologies for Smart Societies

## Is it for free?

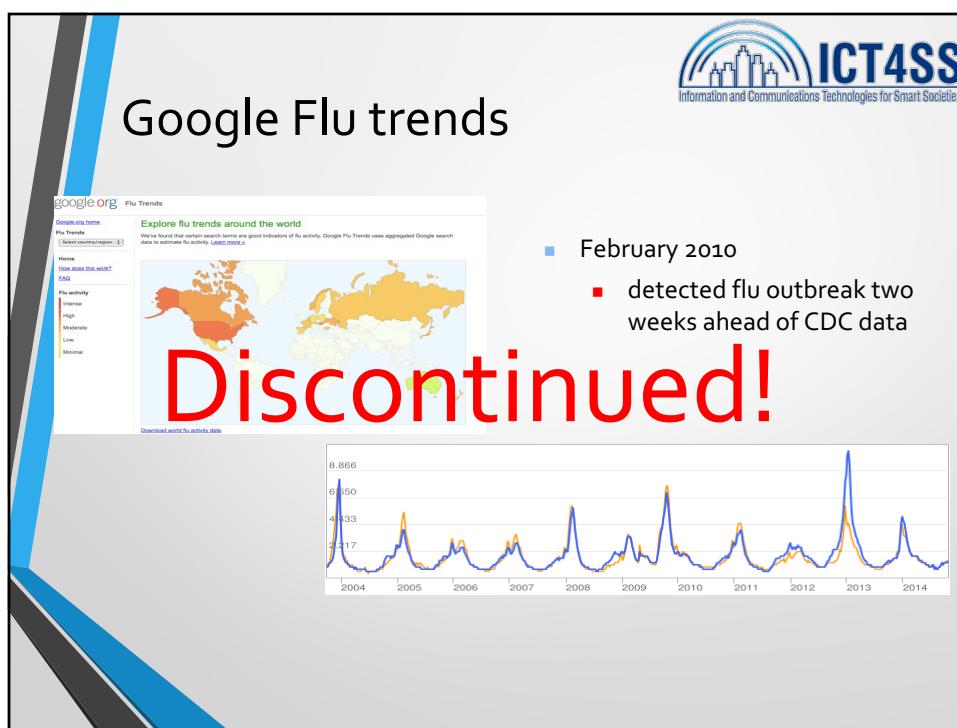
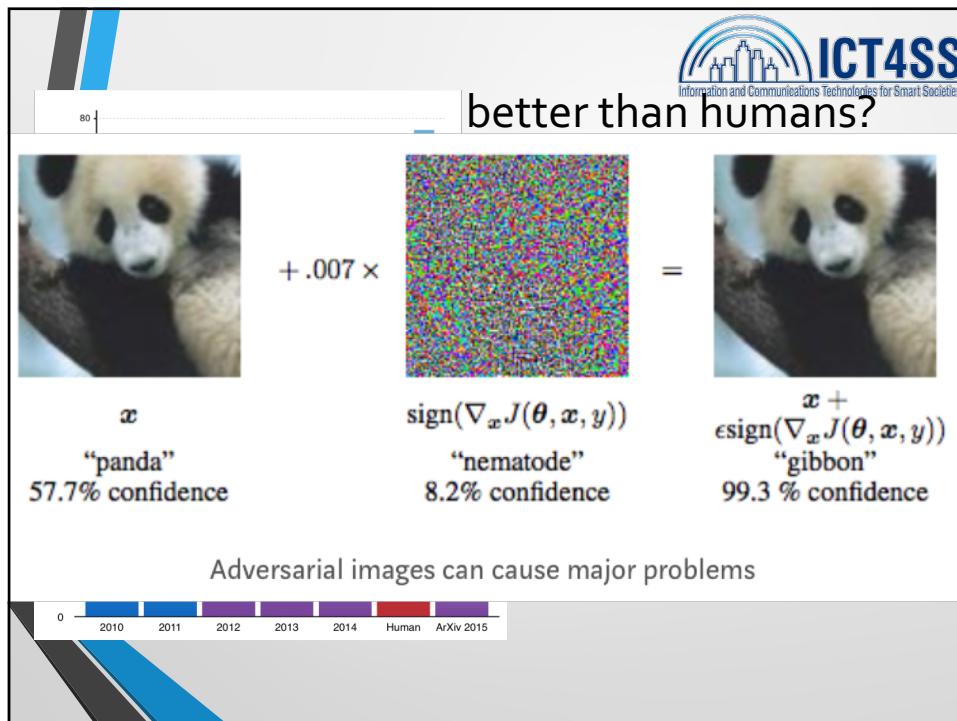
**AlphaGo**

The Future of Go Summit, Match Three: Ke Jie & AlphaGo  
中国乌镇围棋峰会  
Google  
浙江卫视直播  
AlphaGo Wins

ALPHAGO 01:55:46  
Lee Sedol 01:55:41  
Google DeepMind Challenge Match  
YouTube Home

**AlphaGO** 1202 CPUs, 176 GPUs, 100+ Scientists.  
**Lee Se-dol** 1 Human Brain, 1 Coffee.

1 MW      30W



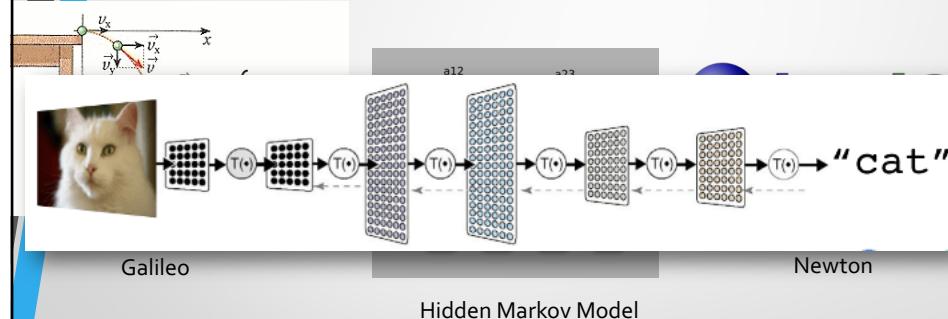


## Conclusions and open issues

- We are leaving in a data deluge epoch
  - Getting data is easier and easier
- Extracting value from the data is still a challenge
  - Don't blindly trust the magician – no Harry Potter stick here
- **Humans need to understand a process**



## Causality and models



**ICT4SS**  
Information and Communications Technologies for Smart Societies

## Conclusions and open issues

- We are leaving in a data deluge epoch
  - Getting data is easier and easier
- Extracting value from the data is still a challenge
  - Don't blindly trust the magician – no Harry Potter stick here
- Humans need to understand a process
- **Social impact of analysis is very important**
  - Interpretability and transparency of the analysis process
  - Privacy preservation
  - Modification of users'habits

**ICT4SS**  
Information and Communications Technologies for Smart Societies

## Privacy

The collage illustrates privacy concerns related to fitness tracking and location sharing. It shows how Strava users share their GPS data, which can reveal sensitive locations like military bases. The BBC article highlights a specific incident where Strava users shared routes around Bagram air base in Afghanistan, raising security concerns.

**Impact on our habits**

- Listen to this 
- 10 most popular questions
  - People cannot read a map
- People become harsh
  - OK google, take me home
    - Ok, here is your path
  - Hey Siri, call mom
    - Ok, calling mom
  - Alexa, buy my pills,
    - Ok, done
- And Kids ?