# Data Science

An Introduction

Setia Pramana

# Data, data and data everywhere……



**Big Data is affecting people everywhere.**
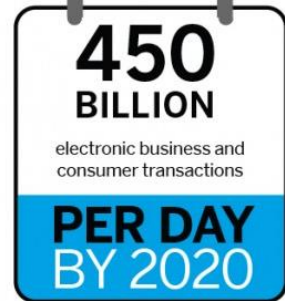
THE WORLD'S TOTAL DATA IS **DOUBLING EVERY** 2 YEARS

There are more mobile phones than people on earth

5,000,000,000 mobile phones are being used to call, text, tweet & browse worldwide each day

Average Google searches per day
YEAR **1998** 9,800
YEAR **2012** 5,134M

% of world's data that was digital
YEAR **1986** 6%
YEAR **2012** More than 99%

**GOOGLE ANSWERS**
1,000,000,000
questions daily from people in **181 countries**

**Big Data is changing business**

36% **ANNUAL INCREASE** in the amount of business data

54% OF C-SUITE EXECUTIVES SAY BOOSTING SALES IS A TOP ADVANTAGE OF BIG DATA

Big data industry estimated to be worth $100 BILLION+

450 BILLION electronic business and consumer transactions **PER DAY BY 2020**

50% of information-intensive businesses will have a **Chief Data Officer** by 2015

80% of the **most** competitive organizations capitalize on data for decision-making.

58% of the **least** competitive organizations capitalize on data for decision-making.

42% of Asia Pacific organisations expect customer service analytics to benefit most from an in-memory data management and analysis technology

**7.9 ZETTABYTES (ZB)** ESTIMATED AMOUNT OF DIGITAL DATA WORLDWIDE **BY 2015** If one dollar bill represented one byte, a zettabyte would stretch from the Sun to Pluto **18,000 times over**

**SAP is helping customers get real value from Big Data**

**MKI performs genome analysis with SAP HANA** "...with this (SAP HANA) we've found a way to shorten the genomic analysis time from several days down to **only 20 minutes.**" YUKIHISA KATO, CTO AND DIRECTOR OF MKI

**eBay uses predictive analytics to gain new insights** "With the **speed of HANA** great people become exceptional at what they do because of the **speed** that they can **interact** with the data. That is truly **awesome.**" DANIEL SCHWARZBACH, VP & CFO EBAY NORTH AMERICA AT EBAY INC.

**Bigpoint solves big data challenges with SAP HANA** Our expectation – and it actually seems to be coming true – is that the use of this technology and the methods behind it helps us realize **sales growth spurts of 10 – 30%.** MICHAEL GUTSMANN, CFO BIGPOINT
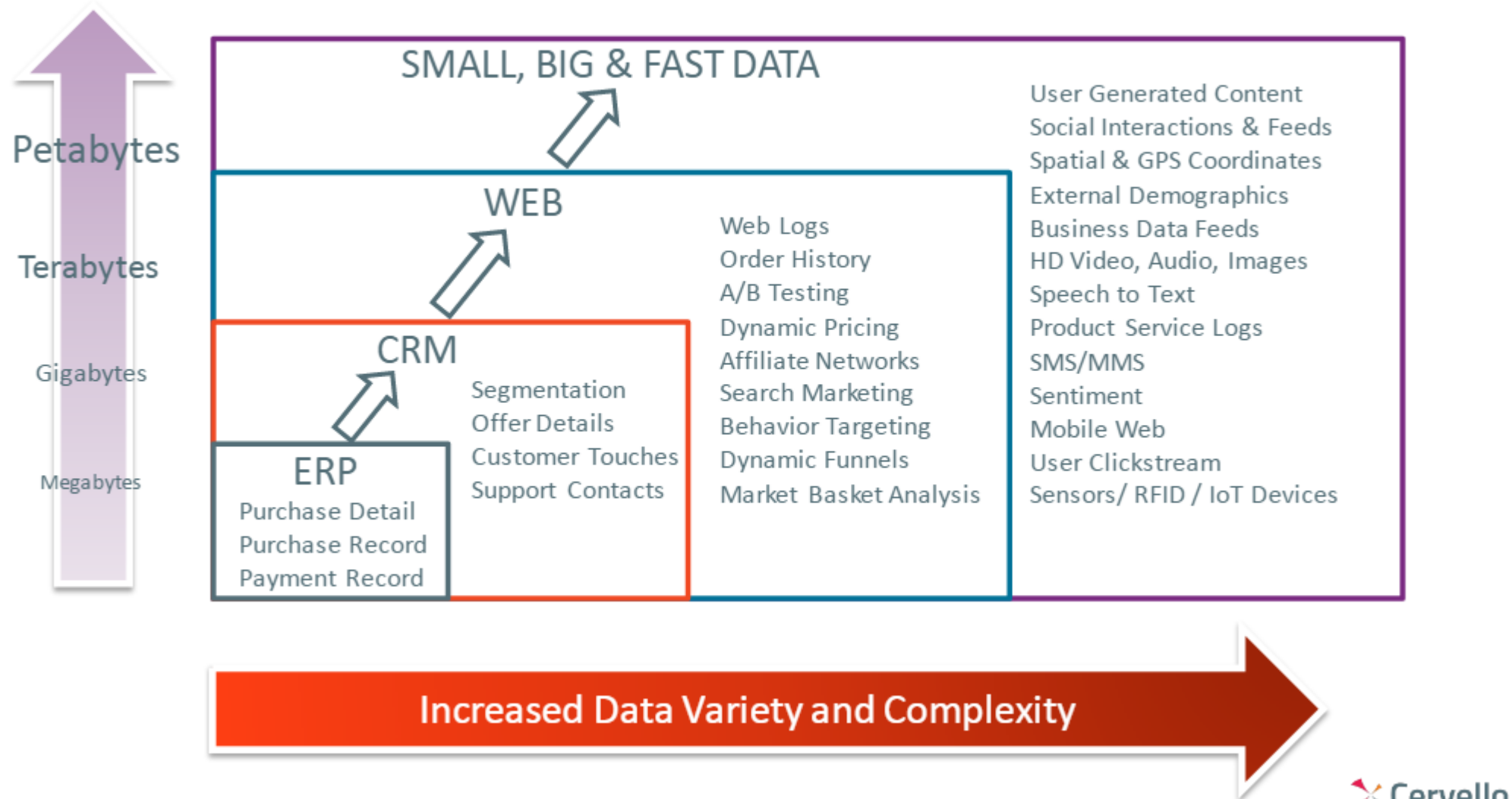
**IDC believes that within five years.....** companies will finally be able to run a real-time enterprise that simultaneously transacts, analyzes and acts on big data.

**Find out what opportunities Big Data holds for you.** Follow us on Twitter **@asiarunbetter**  Visit **www.sap.com/bigdata**  Contact **SAP**

Gartner: The Nexus of Forces Boosts Information Governance and MDM, 2012 | IDC: IDC Predictions 2012: Competing for 2020 | IDC: Ingraining Insights into the DNA of People and Process, 2013 | International Telecoms Union, 2013 | Economist, 2013 | SAP, 2013

# Data = Transactions + Interactions + Observations



Petabytes

Terabytes

Gigabytes

Megabytes

**SMALL, BIG & FAST DATA**

**WEB**

**CRM**

**ERP**
Purchase Detail
Purchase Record
Payment Record

Segmentation
Offer Details
Customer Touches
Support Contacts

Web Logs
Order History
A/B Testing
Dynamic Pricing
Affiliate Networks
Search Marketing
Behavior Targeting
Dynamic Funnels
Market Basket Analysis

User Generated Content
Social Interactions & Feeds
Spatial & GPS Coordinates
External Demographics
Business Data Feeds
HD Video, Audio, Images
Speech to Text
Product Service Logs
SMS/MMS
Sentiment
Mobile Web
User Clickstream
Sensors/ RFID / IoT Devices

**Increased Data Variety and Complexity**

Cervello

http://mycervello.com

3

# THE WORLD OF DATA

| NUMBER OF EMAILS SENT EVERY SECOND | DATA CONSUMED BY HOUSEHOLDS EACH DAY | VIDEO UPLOADED TO YOUTUBE EVERY MINUTE | DATA PER DAY PROCESSED BY GOOGLE | TWEETS PER DAY | TOTAL MINUTES SPENT ON FACEBOOK EACH MONTH | DATA SENT AND RECEIVED BY MOBILE INTERNET USERS | PRODUCTS ORDERED ON AMAZON PER SECOND |
|---|---|---|---|---|---|---|---|
| **2.9** MILLION | **375** MEGABYTES | **20** HOURS | **24** PETABYTES | **50** MILLION | **700** BILLION | **1.3** EXABYTES | **72.9** ITEMS |

IN THE 21ST CENTURY, we live a large part of our lives online. Almost everything we do is reduced to bits and sent through cables around the world at light speed. But just how much data are we generating? This is a look at just some of the massive amounts of information that human beings create every single day.

# Type of Data

- Relational Data (Tables/Transaction/Legacy Data)
- Text Data (Web)
- Semi-structured Data (XML)
- Graph Data
- Social Network, Semantic Web (RDF), …
- Streaming Data
- You can afford to scan the data once

# What to do?

- Aggregation and Statistics
  - Data warehousing and OLAP

- Indexing, Searching, and Querying
  - Keyword based search
  - Pattern matching (XML/RDF)

- Knowledge discovery
  - Data Mining
  - Statistical Modeling

# Analytics Approaches

- Descriptive: What happened or what is happening now?
- Diagnostic: Why did it happen or Why is it happening now?
- Predictive: What will happen next? What will happen under various conditions?
- Prescriptive: What are the options to create the most optimal/high value result/outcome?

# Data Science

"Applying advanced statistical tools to existing data to solve problems, generate new insights, improve products/services"

"Everything that has something to do with data: Collecting, analyzing, modeling...... yet the most important part is its applications --- all sorts of application"

# What is Data Science?

- Theories and techniques from many fields and disciplines are used to investigate and analyze a large amount of data to help decision makers in many industries such as science, engineering, economics, politics, finance, and education
  - Computer Science
    - Pattern recognition, visualization, data warehousing, High performance computing, Databases, AI
  - Mathematics
    - Mathematical Modeling
  - Statistics
    - Statistical and Stochastic modeling, Probability.

# Data Science: An Action Plan for Expanding the Technical Areas of the Field of Statistics

**William S. Cleveland**

*Statistics Research, Bell Labs*

**Abstract:**   An action plan to expand the technical areas of statistics focuses on the data analyst. The plan sets out six technical areas of work for a university department and advocates a specific allocation of resources devoted to research in each area and to courses in each area. The value of technical work is judged by the extent to which it benefits the data analyst, either directly or indirectly. The plan is also applicable to government research labs and corporate research organizations.   © 2014 Wiley Periodicals, Inc. Statistical Analysis and Data Mining 7: 414–417, 2014

## 1.   SUMMARY OF THE PLAN

This article describes a plan to broaden the major areas of technical work of the field of statistics. Because the plan is ambitious and implies su... will be called 'data scien...

The focus of the plan is practicing the data analyst. A

- *Models and Methods for Data (20%)*: statistical models; methods of model building; and methods of estimation and distribution based on probabilistic i... ...hardware systems;

# Data Science Is Multidisciplinary

By Brendan Tierney, 2012

# Data Science

- A Mashed Up Discipline
- A multi-disciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from structured and unstructured data

| Math and Theory | • Statistics, Linear Algebra, Optimization, Time Series, etc. |
| --- | --- |
| Applied Algorithms | • Machine Learning, Data Structures, Parallel Algorithms, etc. |
| Engineering and Technologies | • Storage and computing platforms, statistical tools ,etc. |
| Domain Expertise | • Text, Finance, Images, Econometrics etc. |
| Art | • Visualization, Infographics |
| Best practices and hacks | • Handle missed values in data, transform and represent data, etc. |

# Data Science



**Driving the Success of Data Science Solutions:
Skills, Roles and Responsibilities ...**

Ask good questions

Know the constraints
(e.g., legal, ethics, market)

Latency at Execution?

Decision Making

**Business
Skills**

Build, Buy, Outsource

Transparent Versus "black box"

Gauge political
friction

Performance Criteria That Matter
(ROI, accuracy, profitability
versus market gain)

Deployment

"Analytics Leader"?

Data

Feature Engineering

**IT
Skills**

**Data
Science**

Data
Logistics

Recalibration With
New Data?

High-
performance
Computing

Which Analytics
to Choose?

Project
Execution/Monitoring

Data Exploration

Data Governance

Creativity, Communication

**Gartner**

13

14

THE DATA SCIENCE
**HIERARCHY OF NEEDS**

AI, DEEP LEARNING

A/B TESTING, EXPERIMENTATION, SIMPLE ML ALGORITHMS

ANALYTICS, METRICS, SEGMENTS, AGGREGATES, FEATURES, TRAINING DATA

CLEANING, ANOMALY DETECTION, PREP

RELIABLE DATA FLOW, INFRASTRUCTURE, PIPELINES, ETL, STRUCTURED AND UNSTRUCTURED DATA STORAGE

INSTRUMENTATION, LOGGING, SENSORS, EXTERNAL DATA, USER GENERATED CONTENT

LEARN/OPTIMIZE

AGGREGATE/LABEL

EXPLORE/TRANSFORM

MOVE/STORE

COLLECT

@mrogati

Monica Rogati https://hackernoon.com/the-ai-hierarchy-of-needs-18f111fcc007 15

# THE DATA SCIENCE HIERARCHY OF NEEDS

**Start Up**

**Data Scientist**

LEARN/OPTIMIZE

AGGREGATE/LABEL

EXPLORE/TRANSFORM

MOVE/STORE

COLLECT

AI, DEEP LEARNING

A/B TESTING, EXPERIMENTATION, SIMPLE ML ALGORITHMS

ANALYTICS, METRICS, SEGMENTS, AGGREGATES, FEATURES, TRAINING DATA

CLEANING, ANOMALY DETECTION, PREP

RELIABLE DATA FLOW, INFRASTRUCTURE, PIPELINES, ETL, STRUCTURED AND UNSTRUCTURED DATA STORAGE

INSTRUMENTATION, LOGGING, SENSORS, EXTERNAL DATA, USER GENERATED CONTENT

@mrogati

Monica Rogati https://hackernoon.com/the-ai-hierarchy-of-needs-18f111fcc007 16

THE DATA SCIENCE **HIERARCHY OF NEEDS**

**Medium**

AI, DEEP LEARNING

LEARN/OPTIMIZE

A/B TESTING, EXPERIMENTATION, SIMPLE ML ALGORITHMS

**Data Scientist**

AGGREGATE/LABEL

ANALYTICS, METRICS, SEGMENTS, AGGREGATES, FEATURES, TRAINING DATA

EXPLORE/TRANSFORM

CLEANING, ANOMALY DETECTION, PREP

MOVE/STORE

RELIABLE DATA FLOW, INFRASTRUCTURE, PIPELINES, ETL, STRUCTURED AND UNSTRUCTURED DATA STORAGE

**Data Engineer**

COLLECT

INSTRUMENTATION, LOGGING, SENSORS, EXTERNAL DATA, USER GENERATED CONTENT

@mrogati

**Software Engineer**

Monica Rogati https://hackernoon.com/the-ai-hierarchy-of-needs-18f111fcc007 17

# THE DATA SCIENCE HIERARCHY OF NEEDS

**Large**

LEARN/OPTIMIZE

AGGREGATE/LABEL

EXPLORE/TRANSFORM

MOVE/STORE

COLLECT

AI, DEEP LEARNING

A/B TESTING, EXPERIMENTATION, SIMPLE ML ALGORITHMS

ANALYTICS, METRICS, SEGMENTS, AGGREGATES, FEATURES, TRAINING DATA

CLEANING, ANOMALY DETECTION, PREP

RELIABLE DATA FLOW, INFRASTRUCTURE, PIPELINES, ETL, STRUCTURED AND UNSTRUCTURED DATA STORAGE

INSTRUMENTATION, LOGGING, SENSORS, EXTERNAL DATA, USER GENERATED CONTENT

@mrogati

**Research Scientist, DS Core + ML/AI Engineer**

**Data Science Analytics**

**Data Engineer**

**Software Engineer**

Monica Rogati https://hackernoon.com/the-ai-hierarchy-of-needs-18f111fcc007 18

# Languages, Tools & Software



Source: datacamp

## Data Scientist

also known as Data Managers, statisticians.

A data scientist will be able to take data science projects from end to end. They can help store large amounts of data, create predictive modelling processes and present the findings.

*Skills:* Mathematics, Programming, Communication

*Will use programmes such as:*
SQL, Python, R

## Data Engineers

also known as database administrators and data architects.

They are versatile generalists who use computer science to help process large datasets. They typically focus on coding, cleaning up data sets, and implementing requests that come from data scientists.

*Skills:* Programming, Mathematics, Big data

*Will use programmes such as:*
Hadoop, NoSQL, and Python

## Data Analysts

also known as business Analysts.

They typically help people from across the company understand specific queries with charts.

*Skills:* Statistics, Communication, Business knowledge

*Will use programmes such as:*
Excel, Tableau, SQL

21

# Current Data Scientists Profile



22

Other (9%)

Computer science (20%)

Engineering (9%)

Natural Sciences (11%)

Statistics and Mathmatics (19%)

Data Science and Analysis (13%)

Economics and Social Sciences (19%)
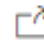
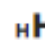Data scientists have heterogeneous academic profiles

91%

365√DataScience

# Data Scientist: The Sexiest Job of the 21st Century

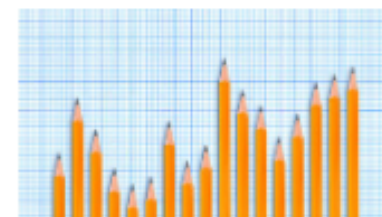by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

Summary    Save    Share    16 Comment    HH Text Size    Print    **$8.95** Buy Copies

**W**hen Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social
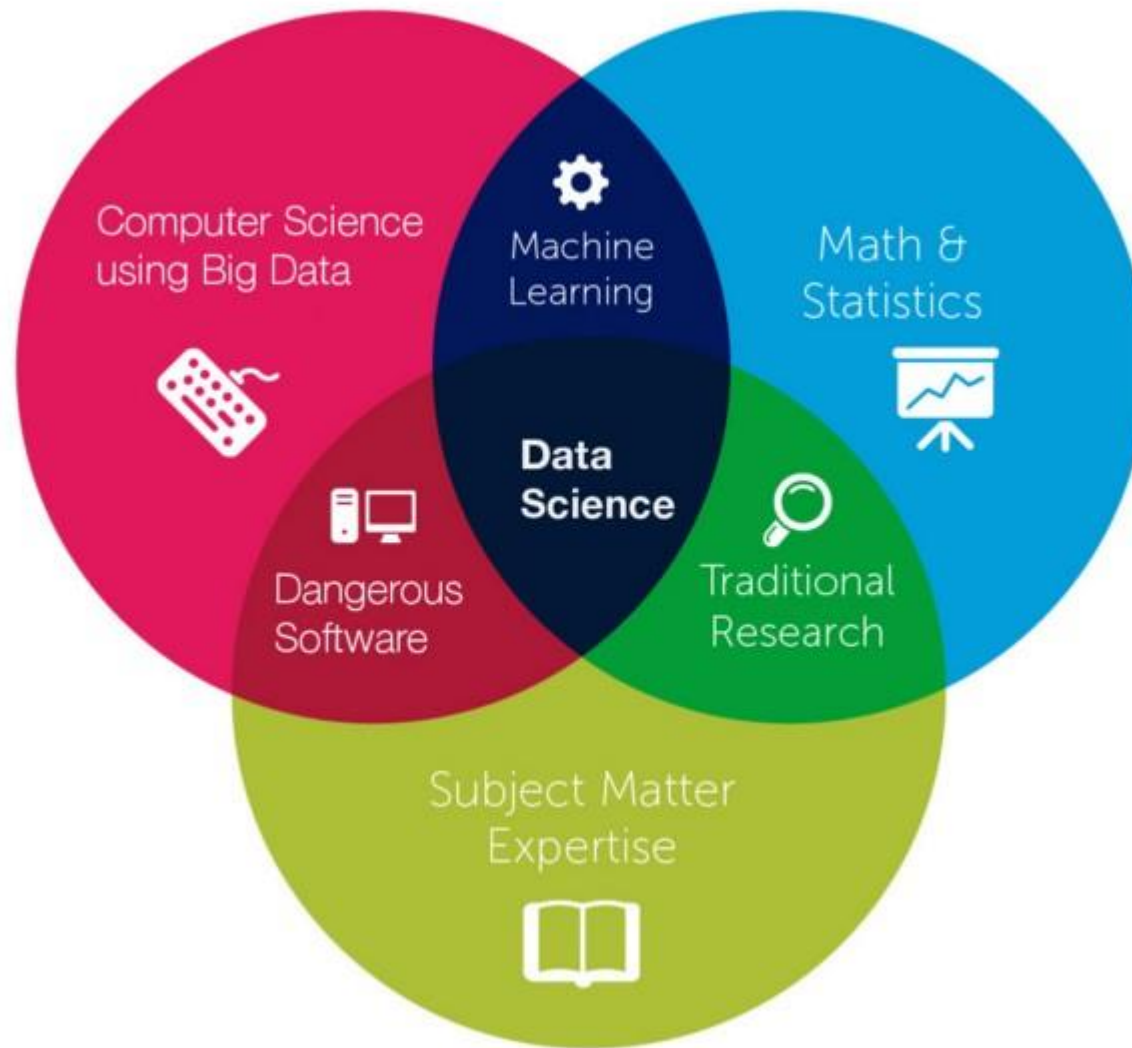
24

# Methods

Sentiment analysis

Time series analysis

Data mining

Multilevel modeling

Missing data imputations

Classification and clustering

Survival analysis

Pattern recognition

Principal component and factor analysis

AB testing

Machine learning

Forecasting

Propensity score matching

Logistic, multinomial and multiple linear regression techniques

Network analysis

# Tools

| Languages | Libraries | Data Engineering | Visualization |
|-----------|-----------|------------------|---------------|
| Python | SciPy | Profiling | D3.js |
| R | Pandas | ETL | Gephi |
| SQL | Scikit-learn | Job notices | R |
| Javascript | GPText | APIs | Leaflet |
| NodeJS | OpenNLP | Optimized data pipelines | PowerBI |
| | Mahout | Optimized data storage/access | ggplot2 |
| | +many others | | shiny |

# Targeting: Find the needle in the haystack

**What to target?**

**Data Science**

**Service Change**

Target areas

Target categories

Target individuals

**Service Issue:**
Difficult to identify targets in a population

**Data Science Process:**
Use existing data and predictive modeling to identify targets

**Service Change:**
Engage with target subset of population

**Result:** Department resources are spent where most needed

# Prioritizing

**What to prioritize?**     **Data Science**     **Service Change**

**Service Issue:**
Backlog is tackled via first in, first out (FIFO)

**Data Science Process:**
Create a model to categorize and group past and current cases

**Service Change:**
Prioritize cases based on categories in order of risk, need or opportunity

**Result:** Department addresses high priority cases first

# Predictions



**How to detect?**  **Data Science**  **Service Change**

**Service Issue:**
Hard to predict future condition which leads to reactive services
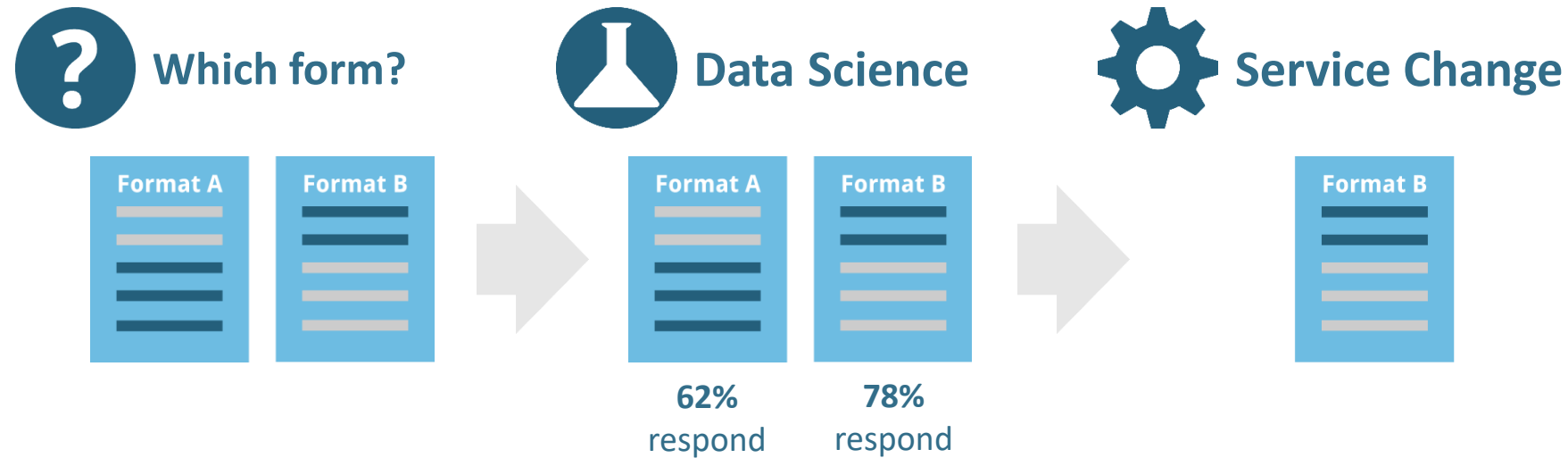
**Data Science Process:**
Use historical and current data to create estimate ranges for potential outcomes

**Service Change:**
Use estimates to change and tailor intervention points

**Result:** Department provides pro-active early interventions

# A/B test

**? Which form?**

**🧪 Data Science**

**⚙ Service Change**

Format A    Format B → Format A    Format B → Format B

62% respond    78% respond

**Service Issue:**
Costly outreach methods are not tested before implementation

**Data Science Process:**
Statistical testing on outreach methods to identify which, when, and to whom to send

**Service Change:**
Use statistically validated outreach method

**Result:** Department increases response rates
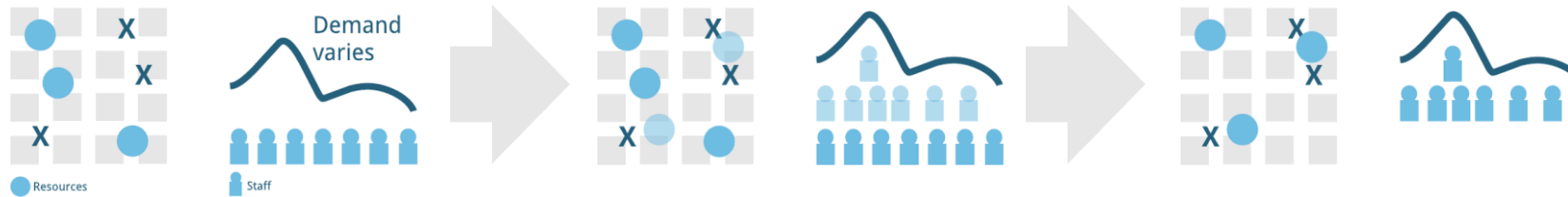
# Optimization

**How to distribute?** **Data Science** **Service Change**

**Service Issue:**
Difficult to identify where to place or distribute resources to be most effective

**Data Science Process:**
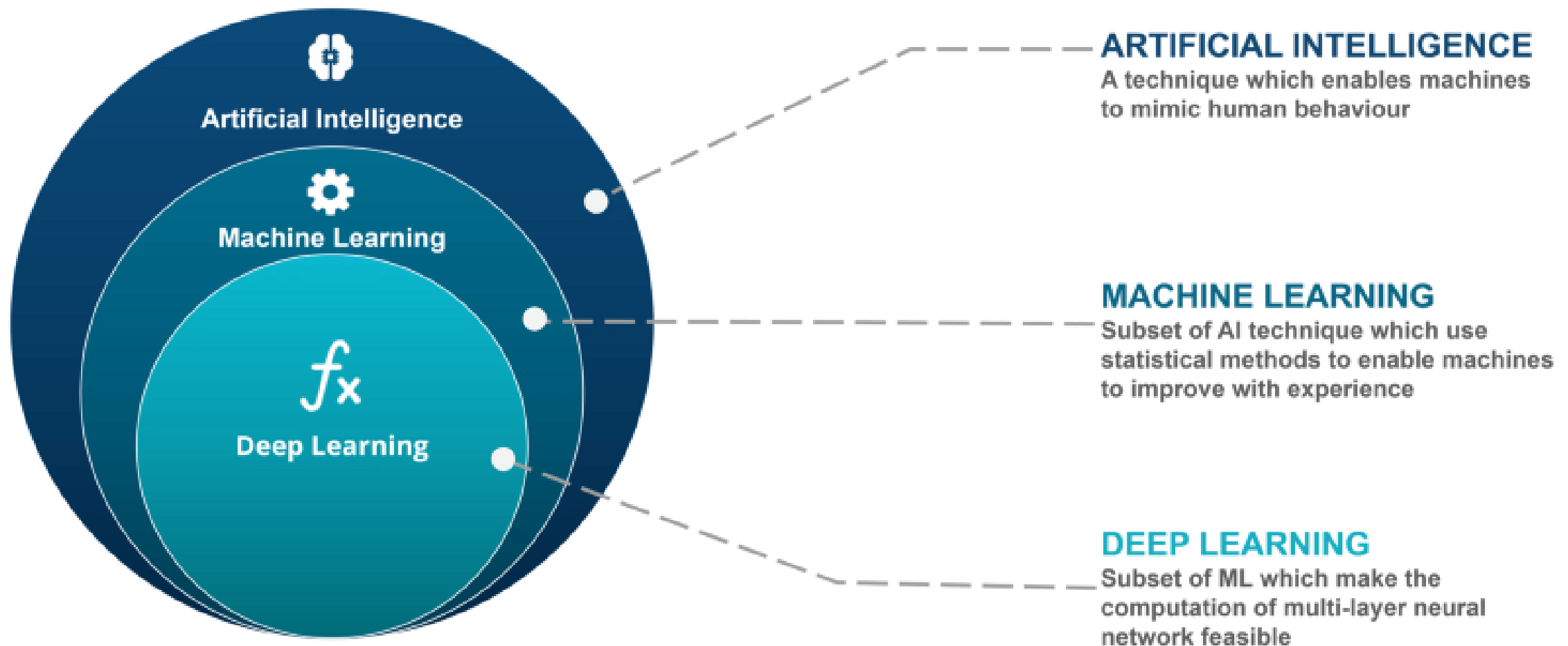Use geospatial and/or other data to identify optimal distribution of resources

**Service Change:**
Re-allocates resources to optimal distribution

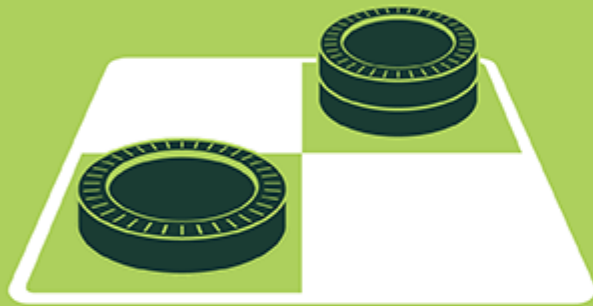**Result:** Department decreases response times; increases volume

**ARTIFICIAL INTELLIGENCE**
A technique which enables machines to mimic human behaviour

**MACHINE LEARNING**
Subset of AI technique which use statistical methods to enable machines to improve with experience

**DEEP LEARNING**
Subset of ML which make the computation of multi-layer neural network feasible

https://medium.com/@StepUpAnalytics/ai-vs-machine-learning-vs-deep-learning-vs-data-science-572b34452c3

ARTIFICIAL INTELLIGENCE
Early artificial intelligence stirs excitement.

MACHINE LEARNING
Machine learning begins to flourish.

DEEP LEARNING
Deep learning breakthroughs drive AI boom.

1950's   1960's   1970's   1980's   1990's   2000's   2010's

Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

https://www.datasciencecentral.com 35

# Data Mining, AI and Machine Learning

- **Data Mining:** extract existing information to highlight patterns, and serves as foundation for AI and machine learning.

- **Artificial Intelligence:** creating machines that perform functions that require intelligence when performed by people.

- **Machine Learning:** Offers data necessary for a machine to learn & adapt. The machine must automatically learn the parameters of models from the data. It uses self-learning algorithms to improve its performance at a task with experience over time

# AI in Sci-Fi Movies

- Terminator



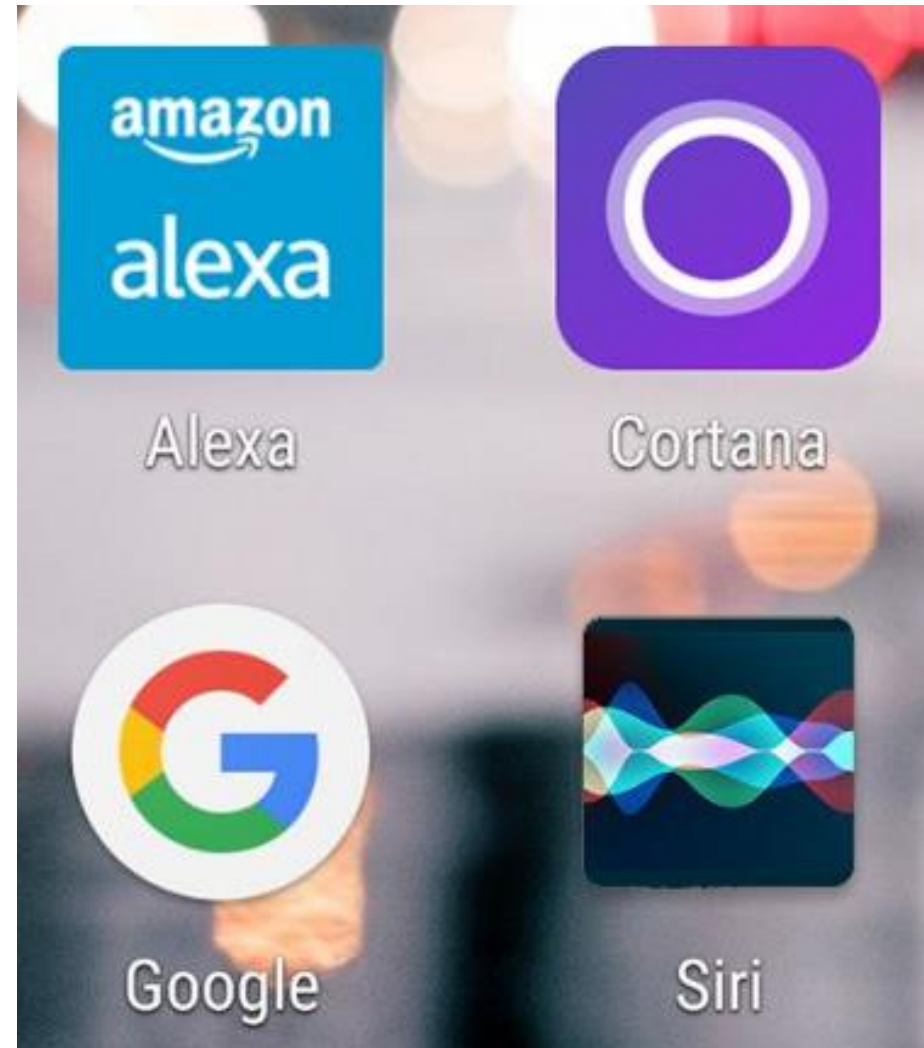http://starwars.com/

- Iron Man Marvel



Just A Rather Very Intelligent System

# AI in Life

## Vacuum Cleaning Robot



## AI assistants

# AI in Life

- Kiva warehouse robot

# What is Artificial Intelligence ?

- The art of creating machines that perform functions that require intelligence when performed by people (Kurzweil, 1990)

- The study of how to make computers do things at which, at the moment, people are better (Rich and Knight, 1991)

- AI: acting humanly

# Machine Learning
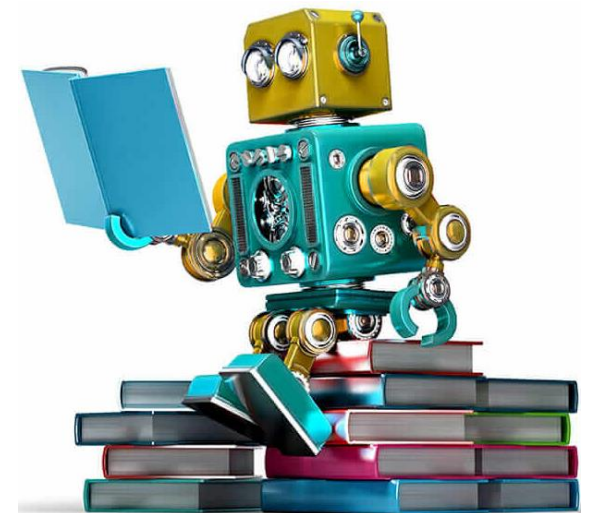
**Learning from Experience**

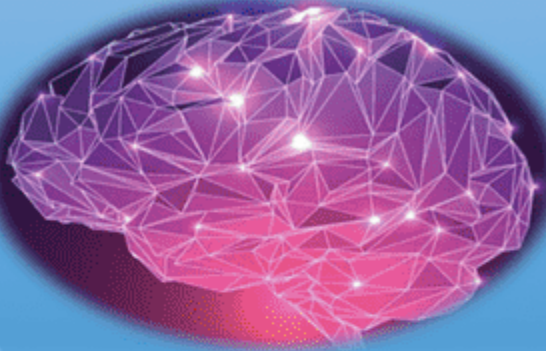**Learning from Data**

**Follow Instructions**

# Machine Learning

- Machine learning is aimed to optimize a certain task using example data or past experience

- The extraction of knowledge from data

- Machine learning is preferred approach to
  - Business Intelligence
  - Speech recognition, Natural language processing
  - Computer vision
  - Robot control
  - Computational biology
  - Crime predictions
  - Etc..

# Machine Learning & Some Use Cases

## Machine Learning
Where business and experience meet emerging technology and decides to work together".

FB use it at extreme level for Spam control, discovering new content, recommendations and Ad sales

Google use it for many many reasons i.e. for maps, route calculations, data collection, translations, email spam and many more.
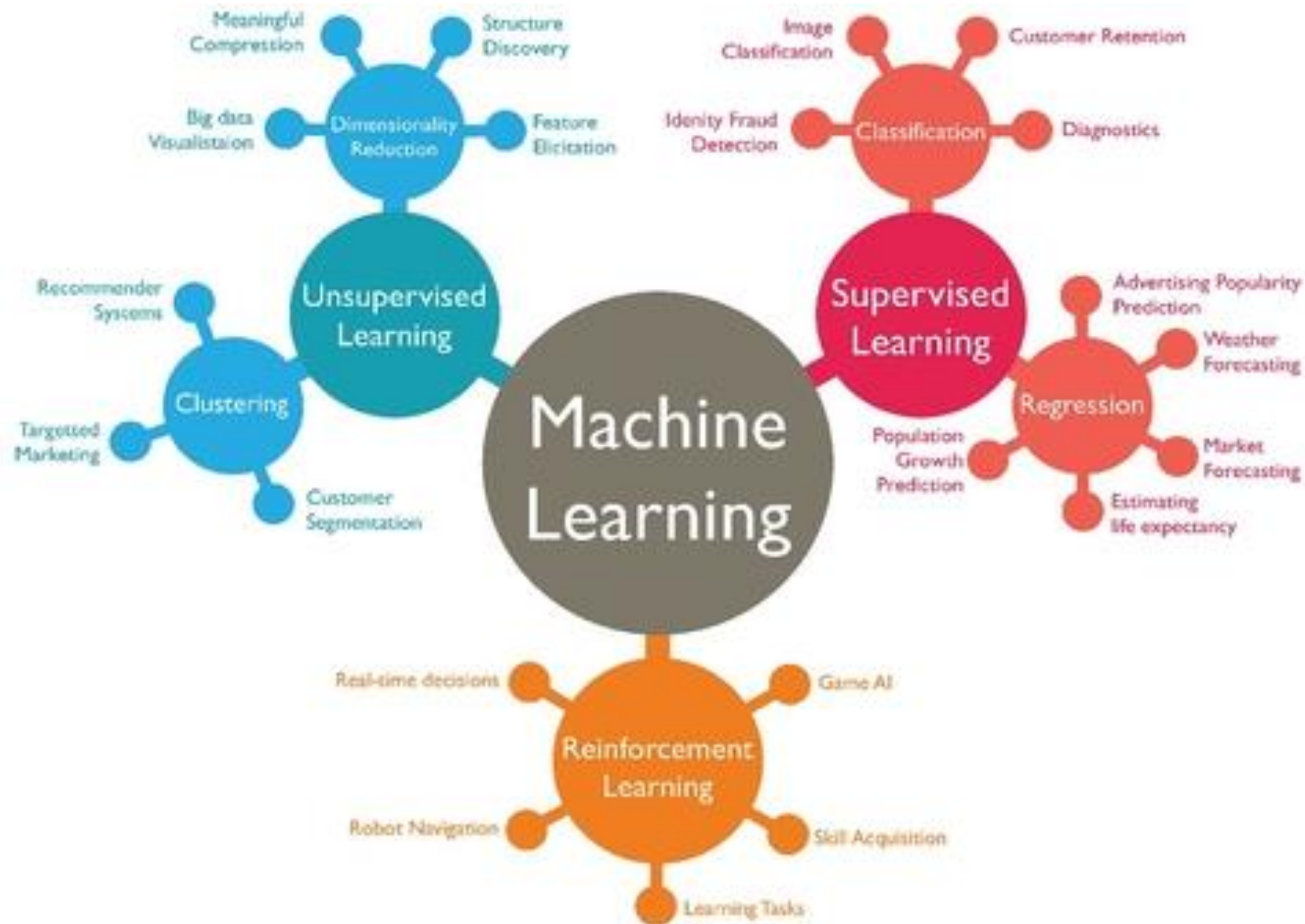
Credit Card Companies getting this now deeper and deeper to minimize the frauds and give safer transaction experience for customers

# Types of Learning

- **Supervised (inductive) learning**
  - Training data includes desired outputs

- **Unsupervised learning**
  - Training data does not include desired outputs

- **Reinforcement learning**
  - Rewards from sequence of actions

# Methods

- **Supervised learning**
  - Decision tree induction
  - Rule induction
  - Naïve Bayes
  - Neural networks
  - Support vector machines
  - Model ensembles
  - Etc.

- **Unsupervised learning**
  - Clustering
  - Dimensionality reduction

- **Reinforcement learning**
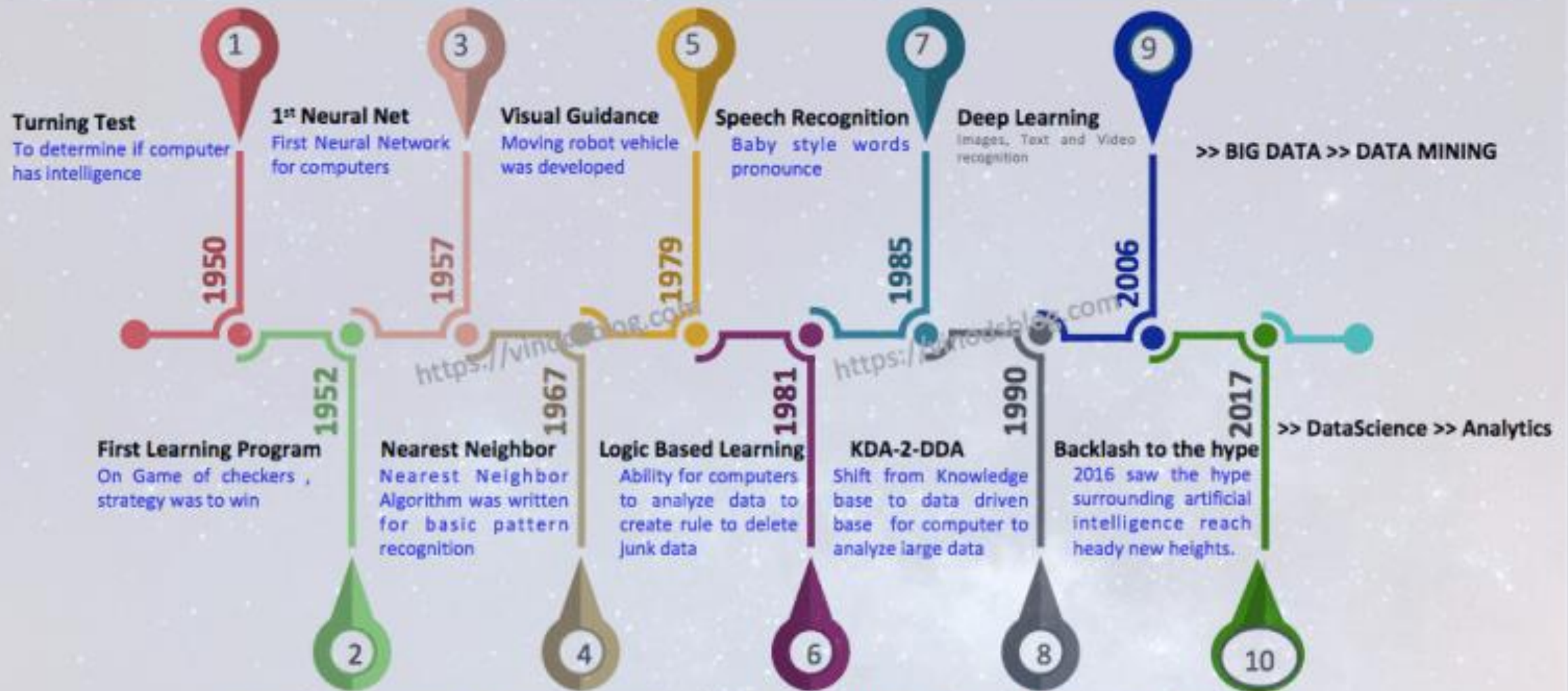  - Decision making (robot, chess machine)

# From Data Mining to Knowledge Discovery in Databases

*Usama Fayyad, Gregory Piatetsky-Shapiro, and Padhraic Smyth*

■ Data mining and knowledge discovery in databases have been attracting a significant amount of research, industry, and media attention of late. What is all the excitement about? This article provides an overview of this emerging field, clarifying how data mining and knowledge discovery in databases are related both to each other and to related fields, such as machine learning, statistics, and databases. The article mentions particular real-world applications, specific data-mining techniques, challenges involved in real-world applications of knowledge

This article begins by discussing the historical context of KDD and data mining and their intersection with other related fields. A brief summary of recent KDD real-world applications is provided. Definitions of KDD and data mining are provided, and the general multistep KDD process is outlined. This multistep process has the application of data-mining algorithms as one particular step in the process. The data-mining step is discussed in more detail in the context of specific data-mining al-

Machine Learning Evolution Over The Years

# Business Intelligence in Banking

- Customer account data and demographics
- Core banking data
- Transactional data at every level of detail
- Wire and payment data
- Trade and position data
- General ledger data including accounts payable, accounts receivable, cash management, purchasing information
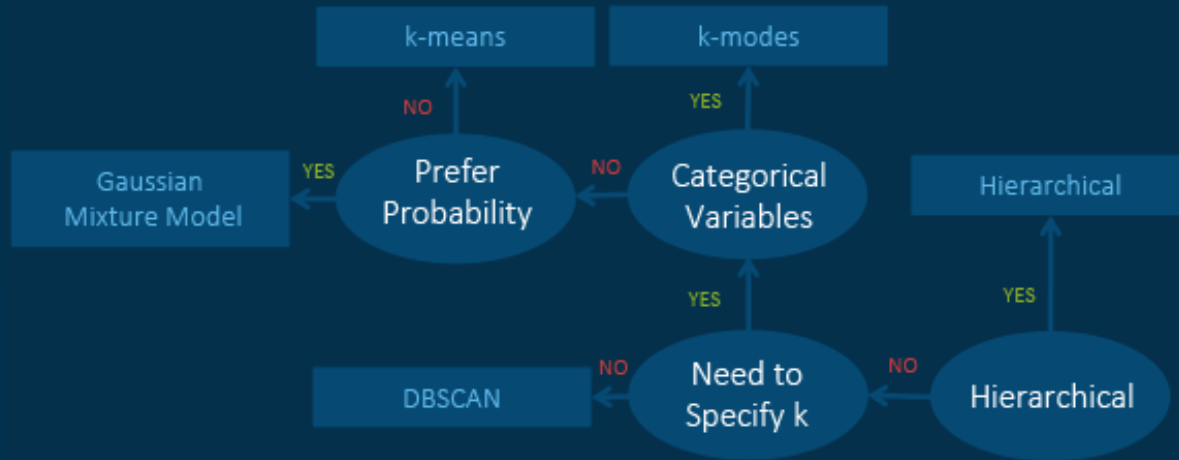- Support data from banking reporting

# Machine Learning in Finance

- Fraud prevention
- Portfolio and Risk Management
- Investment predictions
- Customer service
- Digital assistants
- Marketing
- Network security
- Loan underwriting
- Algorithmic trading
- Customer Service (Chatbot)

- Process automation
- Document interpretation
- Content creation
- Trade settlements
- Money-laundering prevention
- Custom machine learning solutions
- Sales/Recommendations of Financial Products
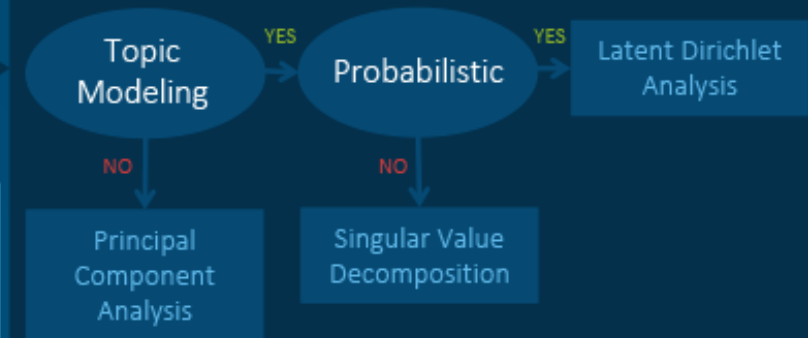- Sentiment/News Analysis
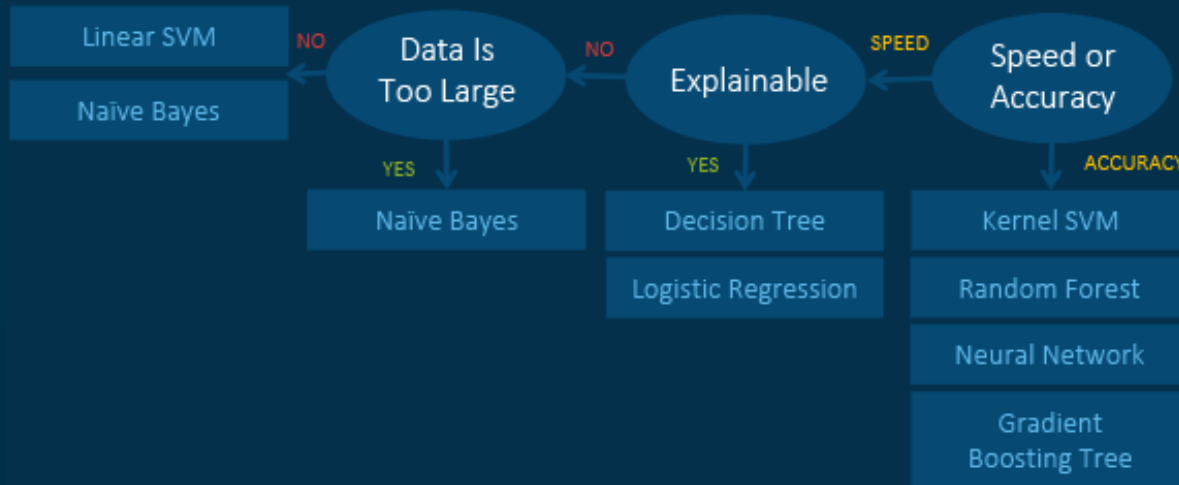
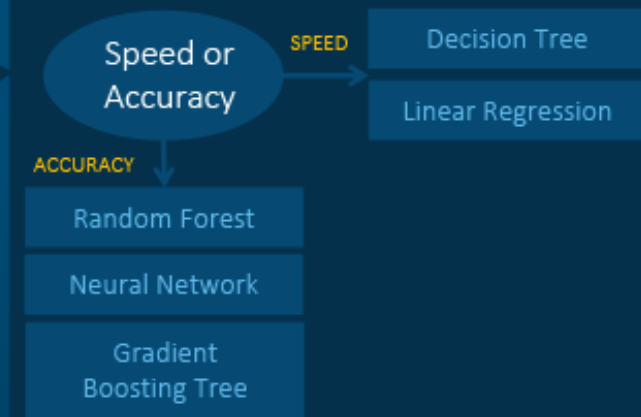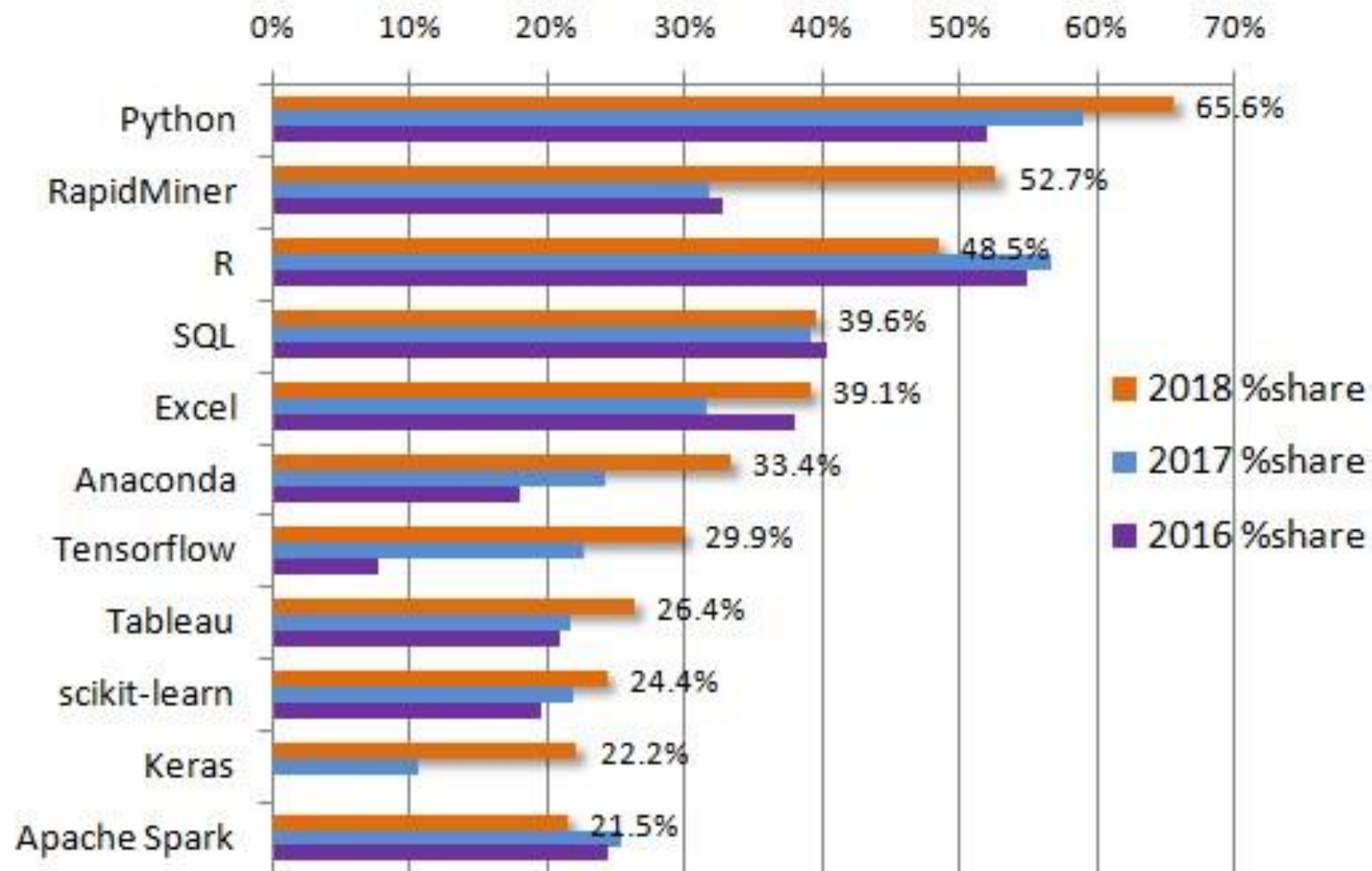Machine Learning Algorithms Cheat Sheet

# Numerous (New) Algorithms

# Numerous (New) Algorithms

Classification algorithms considered in the benchmarking study.

| | BM selection | Classification algorithm | Acronym |
|---|---|---|---|
| **Individual classifier** | n.a. | Bayesian Network | B-Net |
| | | CART | CART |
| | | Extreme learning machine | ELM |
| | | Kernalized ELM | ELM-K |
| | | k-nearest neighbor | kNN |
| | | J4.8 | J4.8 |
| | | Linear discriminant analysis[a] | LDA |
| | | Linear support vector machine | SVM-L |
| | | Logistic regression[a] | LR |
| | | Multilayer perceptron artificial neural network | ANN |
| | | Naive Bayes | NB |
| | | Quadratic discriminant analysis[a] | QDA |
| | | Radial basis function neural network | RbfNN |
| | | Regularized logistic regression | LR-R |
| | | SVM with radial basis kernel function | SVM-Rbf |
| | | Voted perceptron | VP |
| **Classification models from individual classifiers** | | | **16** |
| **Homogenous ensembles** | n.a. | Alternating decision tree | ADT |
| | | Bagged decision trees | Bag |
| | | Bagged MLP | BagNN |
| | | Boosted decision trees | Boost |
| | | Logistic model tree | LMT |
| | | Random forest | RF |
| | | Rotation forest | RotFor |
| | | Stochastic gradient boosting | SGB |

| | | | |
|---|---|---|---|
| | n.a. | Simple average ensemble | AvgS |
| | | Weighted average ensemble | AvgW |
| | | Stacking | Stack |
| **Heterogeneous ensembles** | Static direct | Complementary measure | CompM |
| | | Ensemble pruning via reinforcement learning | EPVRL |
| | | GASEN | GASEN |
| | | Hill-climbing ensemble selection | HCES |
| | | HCES with bootstrap sampling | HCES-B |
| | | Matching pursuit optimization ensemble | MPOE |
| | | Top-T ensemble | Top-T |
| | Static indirect | Clustering using compound error | CuCE |
| | | k-Means clustering | k-Mean |
| | | Kappa pruning | KaPru |
| | | Margin distance minimization | MDM |
| | | Uncertainty weighted accuracy | UWA |
| | Dynamic | Probabilistic model for classifier competence | PMCC |
| | | k-nearest oracle | kNORA |
| **Classification models from heterogeneous ensembles** | | | **17** |

KDnuggets Analytics, Data Science, Machine Learning Software Poll, 2016-2018

| Software | 2018 %share |
|----------|-------------|
| Python | 65.6% |
| RapidMiner | 52.7% |
| R | 48.5% |
| SQL | 39.6% |
| Excel | 39.1% |
| Anaconda | 33.4% |
| Tensorflow | 29.9% |
| Tableau | 26.4% |
| scikit-learn | 24.4% |
| Keras | 22.2% |
| Apache Spark | 21.5% |

Legend: 2018 %share, 2017 %share, 2016 %share

# Next, we learn R… and Python….