



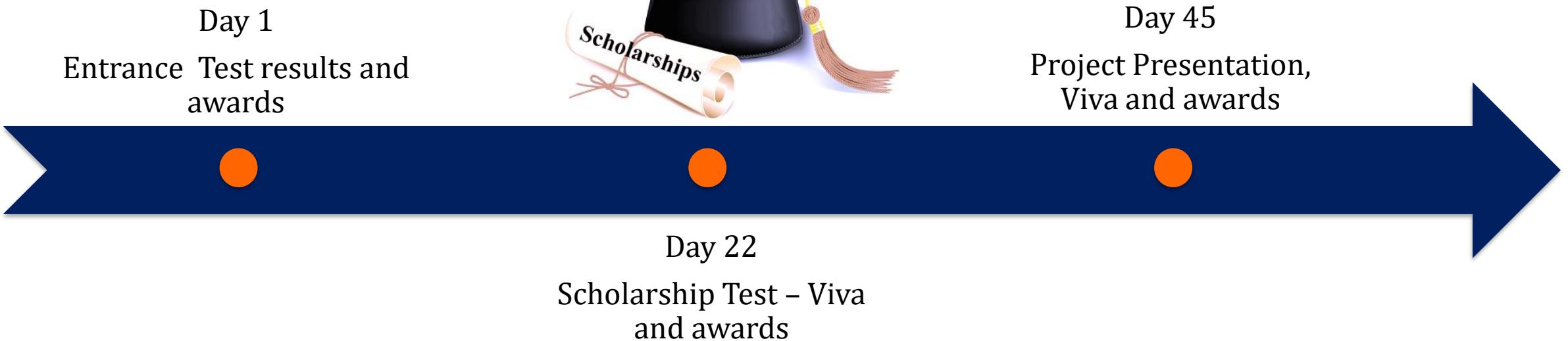
# Inspire...Educate...Transform.

# **Certificate Program in Engineering Excellence**

## **Orientation – Batch 28**

**25<sup>th</sup> March, 2017**

# Scholarship



- 10% of your batch revenue is allocated as the total scholarship amount.
- Number of scholarships and the amount disbursed will depend on the performance at each phase. The academic team in-charge will conduct the evaluations and it is up to their discretion to decide the amount and number of scholarships.
- A student can claim either partial or full fee amount in this process.
- The maximum a student can gain is INR 3,00,000
- Internship announcements as required by INSOFE



# Class Structure

- 4 hours lecture...focus on explaining concepts.
- 4 hours lab...focus on hands-on.
- Hands-on activities using R & Python. Do not underestimate its importance.  
Will have focus in exams.
- Welcome to come to office during weekdays and clarify questions. But please make an appointment as Data Scientists are also involved in CPR activities.



# Course Material and Discussion



## Piazza

**Piazza** CSE 7202c Batch 10: 274

Q & A Resources Statistics Manage Class Question History

**question** 24 views

**Few basic questions....**

Hello sir...  
I had few basic questions ...

> When we have multiple predictors... what are the plots we should look for initially... histogram of single predictor.. or scatter plot of target variable Vs single predictor... or any other?  
> When to deal with outliers before model building or after model building... and also which method to prefer for before(boxplots, normal dist 3sigma,etc)... how to deal with influential observations...  
> Which value to look for comparing different models (AIC or Adjusted R2 or something else)  
> Which error metric should we look for while testing model on test data (mae or mse or rmse or mape)...  
> Which search procedures are mostly used in linear and logistics regression.

**edit** good question | 1 Updated 4 months ago by Sridhar Pappu and amitbabesabhe sasane

**i the instructors' answer**, where instructors collectively construct a single answer

Amit,

The straightforward answer to all these is what I have mentioned a couple of times in the class that in Data Science, there is no set procedure that if you do things in a certain way, the output is defined and guaranteed. You will have to look at all the information and make decisions. For that, you need to understand what each of the metrics and/or methods means and when they are applicable. Having said that, I am glad you asked these because people always get hung up on such things and miss the fundamental understanding. I would like everyone to read my comments below and especially the last part.

Things are not very clear a lot of the times and assumptions made, such as normal distribution, etc. may not be applicable. In that sense, you would need to everything you can do and if outputs vary drastically, you will have to relook the data and the approach.

To answer each of your questions specifically:

> When we have multiple predictors... what are the plots we should look for initially... histogram of single predictor.. or scatter plot of target variable Vs single predictor... or any other?  
Visual information is always very useful and can give you important insights but you should not make decisions in isolation as there could be interactions and correlations may not necessarily mean causation. However, as and when you have to make a decision, it is always good to use your own logic and make it a valid one. Visuals are just tools to help you make a decision. Over a period of time, you...

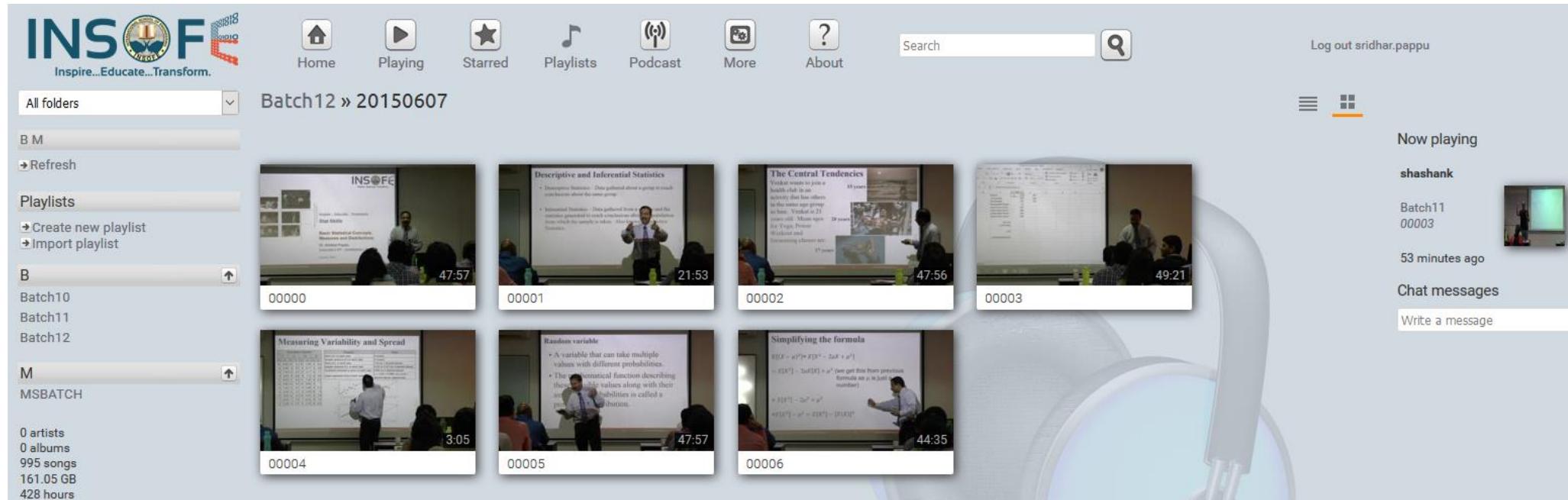
Average Response Time: Special Mentions: N/A Sridhar Pappu answered Step Wise Regression in 15 min. 4 months ago

Online Now | This Week 0 | 15

- Ask a lot of questions
- Post questions, assignments, etc. in the correct module; else, they will be lost.



# Video Recordings



- Do not miss classes. Video  $\neq$  Classroom
- Will be available for 3 months after the program.
- Contact if you need extension.



# Assessment

- ~30 WUQs, 6 GNQs, Project
- Weightage
  - Project: 40%
  - GNQ: 30%
  - WUQ: 10%
  - Feedback: 5%
  - Attendance (85%+): 5%
  - Lab exercises: 10%
- Non-zero cumulative score in each a must



# Feedback

## Zoho Survey

- 5% grade as incentive.
- One per module.
- Provide as soon as you receive the link.

Instructors: Dr. Sridhar Pappu

**PLEASE LIMIT YOUR FEEDBACK TO THIS MODULE, ITS FACULTY AND DATA SCIENTISTS WHO HANDLED LABS FOR THIS MODULE.**

1. Please tell us about yourself. We will ask your permission in the last question on using this information in our collateral. We will NOT use this information without your explicit permission.

Your Name:\*

Your Designation:\*

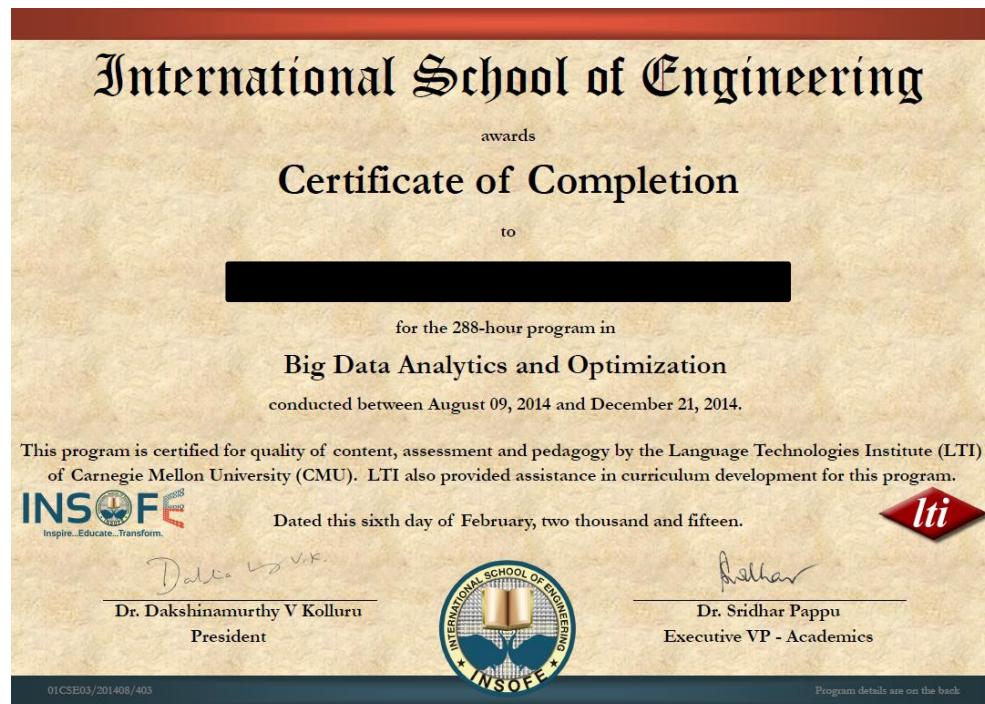
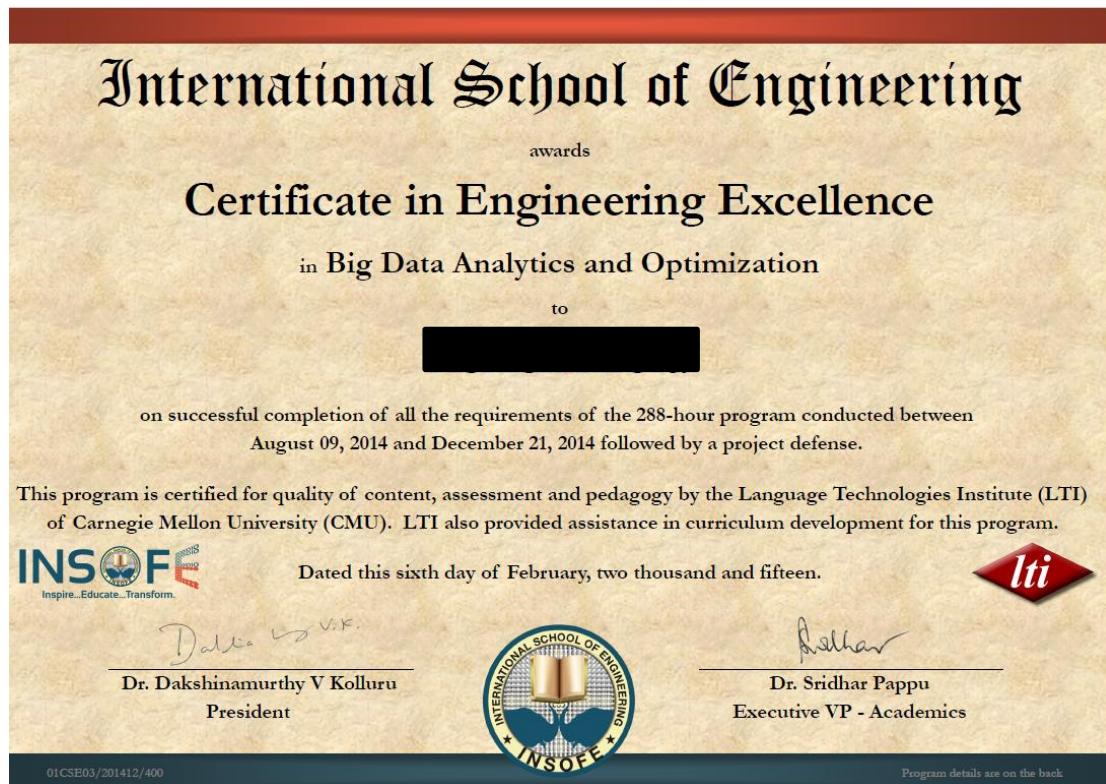
Your Company:\*

2. Content\*

|  | Poor                  | Average               | Good                  | Excellent             | Not Applicable        |
|--|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
| How smooth was the flow of the course material?                  | <input type="radio"/> |
| Were the topics covered in sufficient detail?                    | <input type="radio"/> |
| Did it provide real world experience?                            | <input type="radio"/> |
| How useful is the instructional and reference material provided? | <input type="radio"/> |

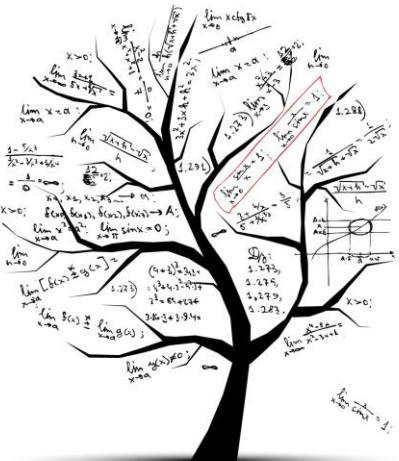
# Certification (or not)

- No certificate if attendance is below 70% and cumulative grade < 50%.
- You have to work really hard to be (in)eligible for this.



# Seeking Help

- Contact relevant person(s) only.
- All contact details provided in the Guidelines document.



For any subject-matter help

Maheshkumar Duvvarapu | Pavan  
Srungaram | Vivek Bakaraju  
| Shilpa Kadam



All payments and loan related

Ravindra Karanam



For any in-class issues

Yugandhar Reddy  
Chakradhar Reddy



# Some Good Practices

- Classes start at 9 AM. Please respect others' time.
- Please switch off phones or keep them in silent mode. Please do not distract others.
- Please do not browse internet unless instructed to do so. We want you to learn.
- We take attendance in the last hour of the day. We do want you to learn.
- Please do not ask questions on topics not yet covered. No ransom will be paid for hijacking class discussions.
- Refer to the Guidelines document in your welcome kit.
- Write all your exams. In case you are likely to miss any exam, please post a note in piazza in **advance**. Before Wednesday of the following week, take the exam.

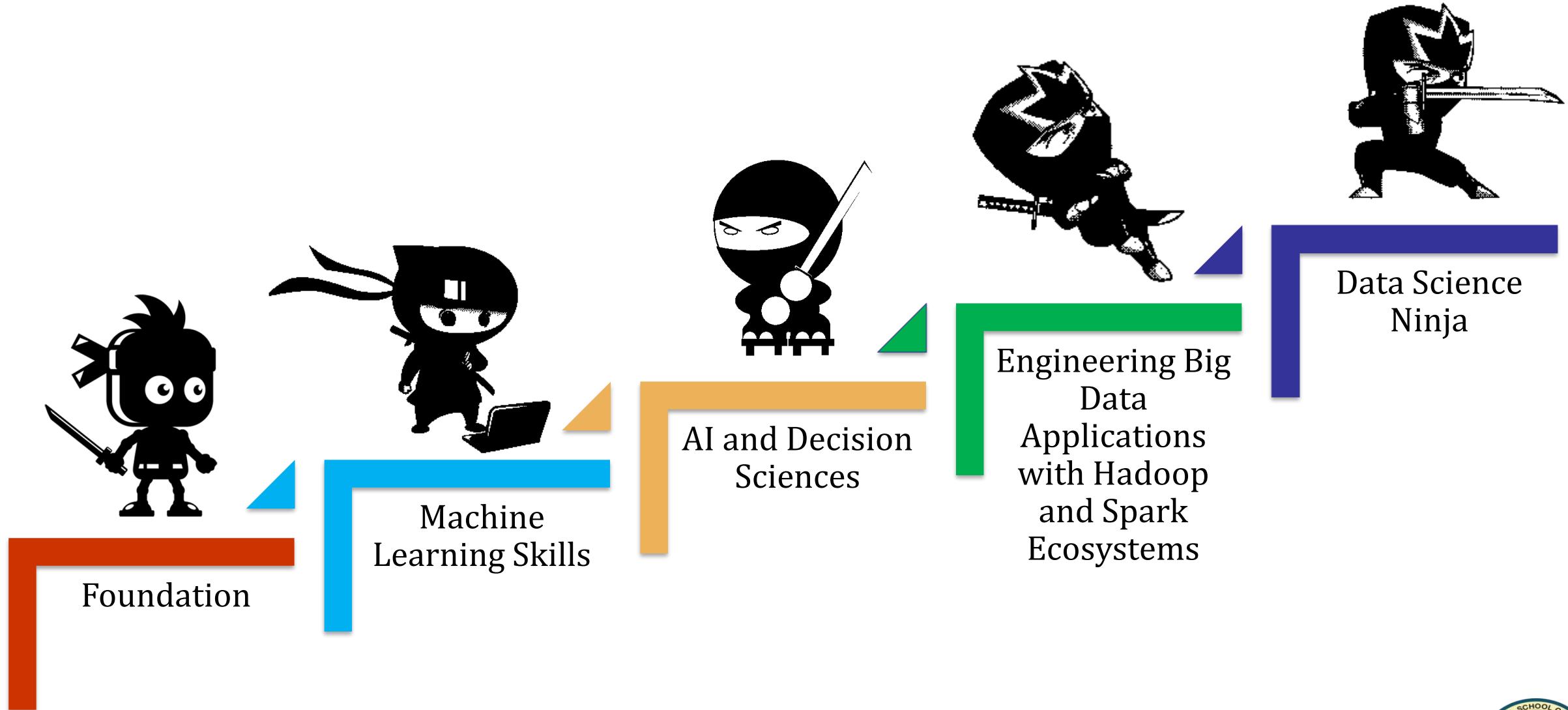






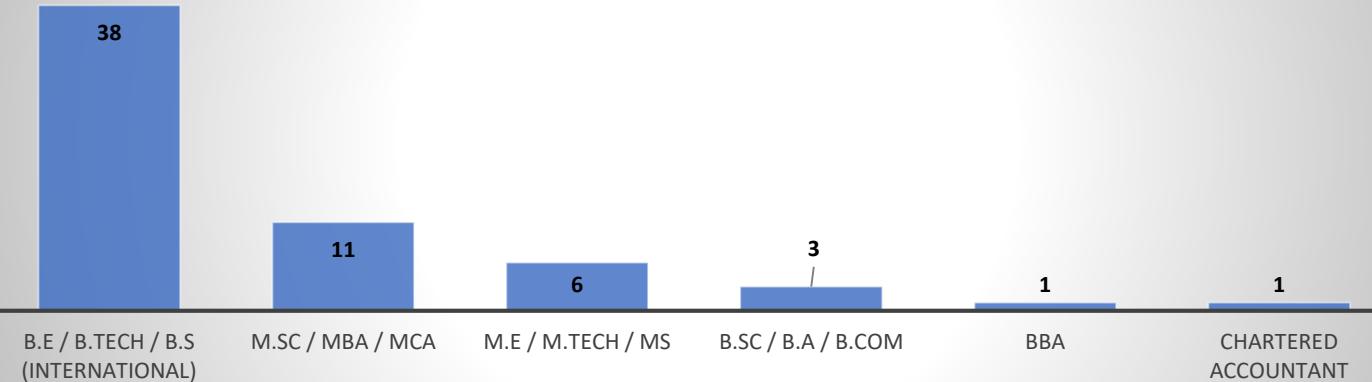


# Curriculum structure – how are we heading?

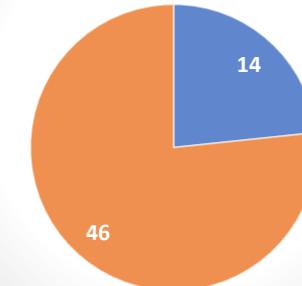


# Background

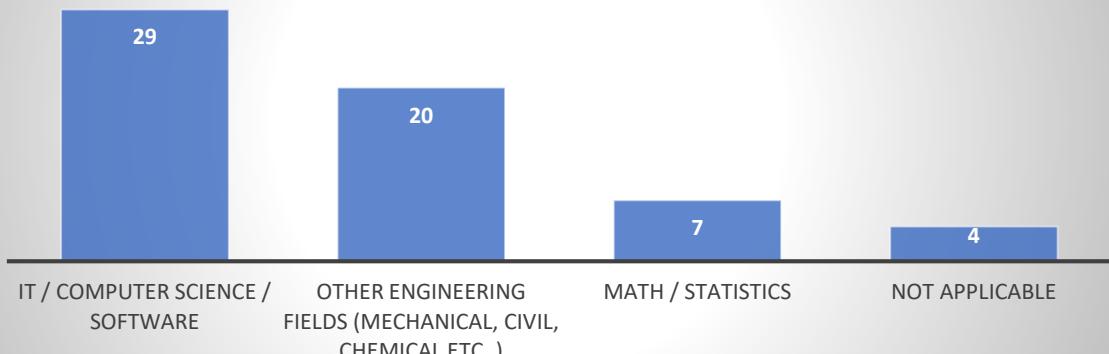
## Educational Background



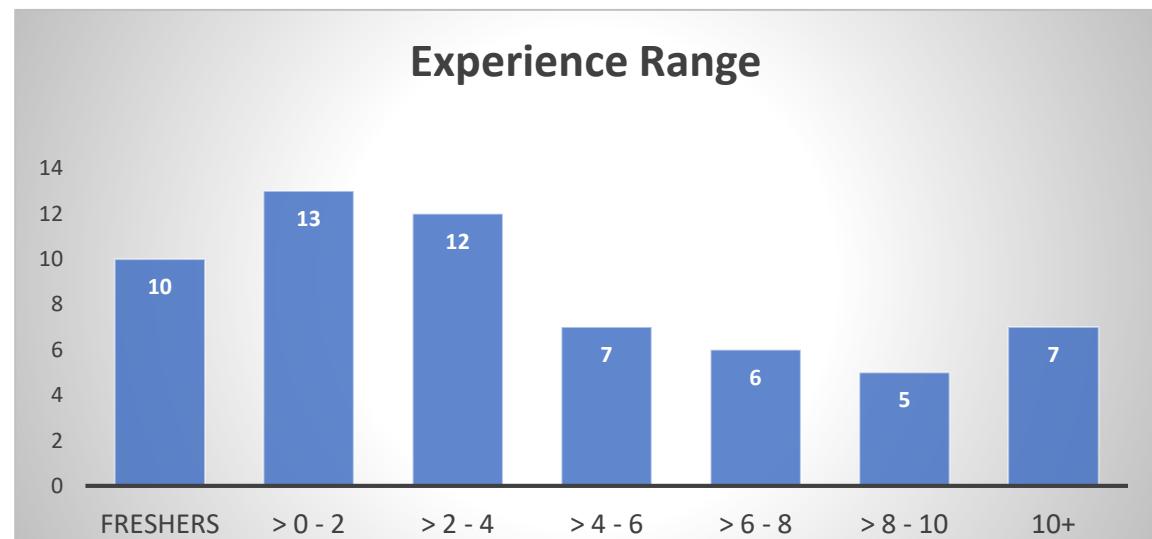
## Gender Distribution



## Educational Background - Specialization



## Experience Range



# Break

# The General Tasks of a Data Scientist

## **20,000-FOOT VIEW OF DATA SCIENCE**

# Data Science

*“Goal is in **extracting meaning from data** and **creating data products** and seeks to **use all available and relevant data** to effectively **tell a story** that can be easily **understood by non-practitioners.**”*

- From Wikipedia (not on latest version of the page)



# Clarifying the Jargon



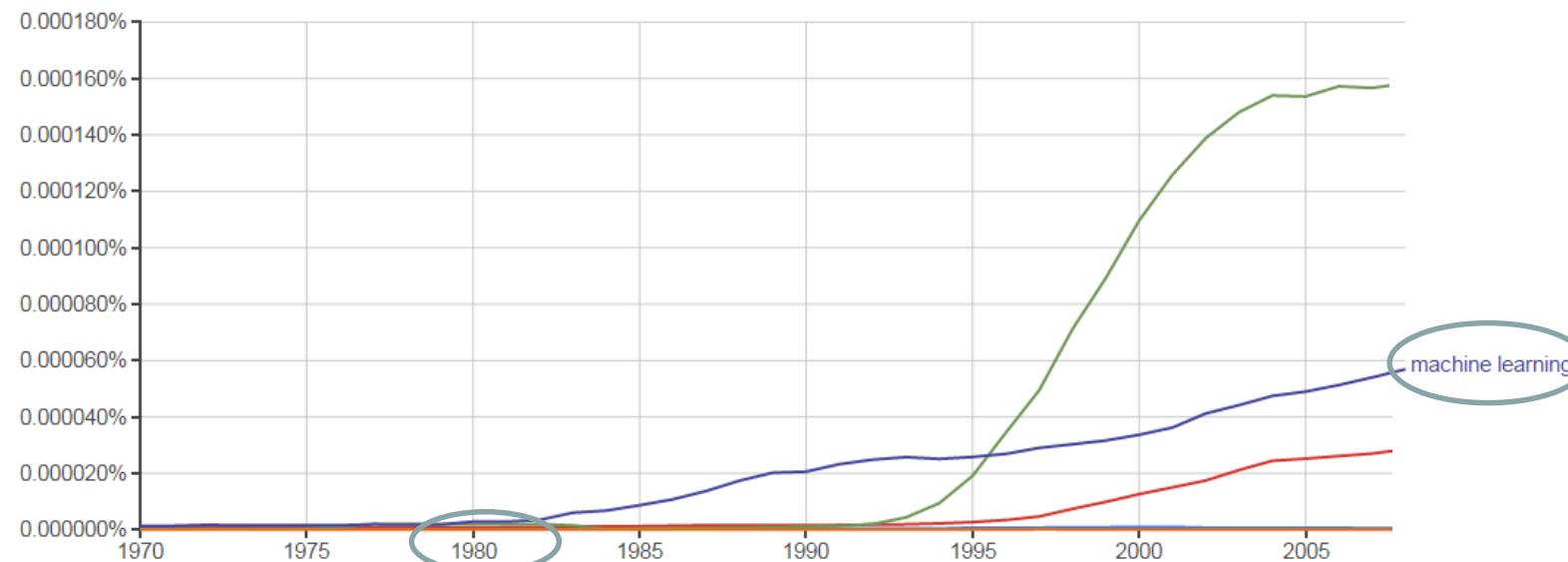
- Machine Learning
- Predictive Analytics/Data Mining
- Big Data Analytics
- Data Science

ALL THESE TERMS ARE PART OF THE SAME FIELD USED AT DIFFERENT TIMES BY PEOPLE WITH  
DIFFERENT BACKGROUNDS



# Clarifying jargon: Chronology

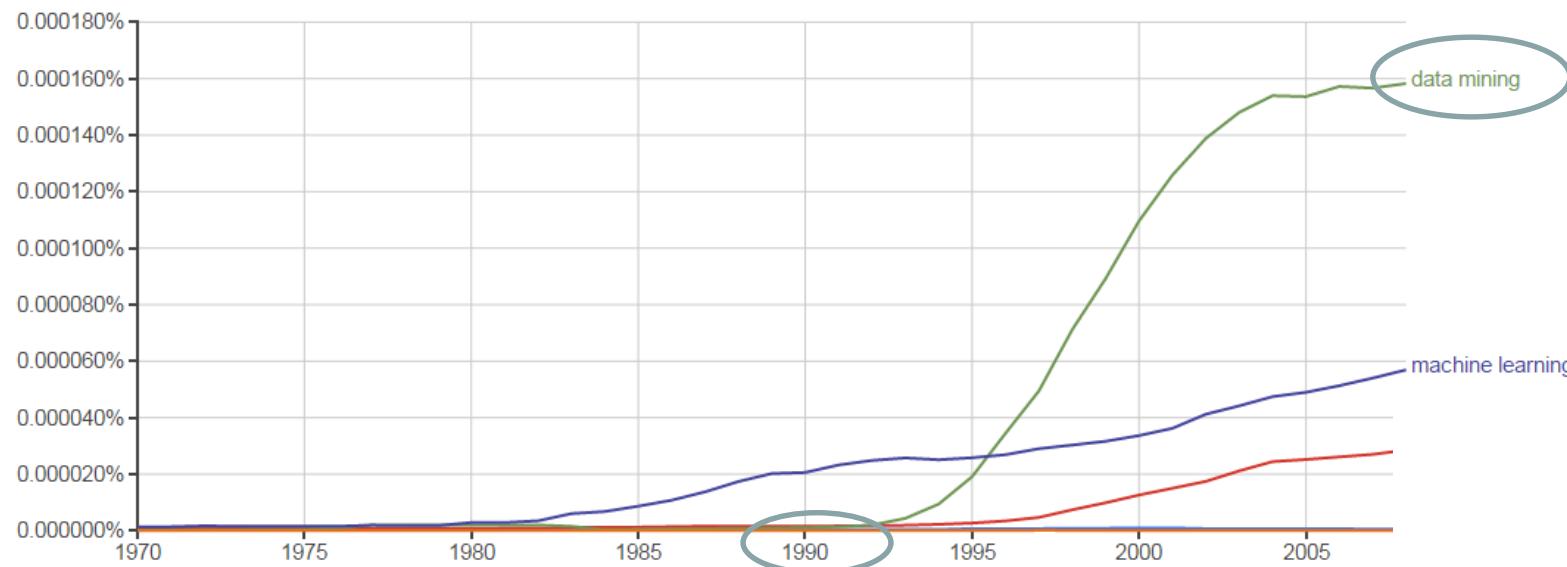
- Machine learning – 1980's
  - Computing Departments called it
  - Focus was on algorithm and the amount of data is limited



**Source:** Google Books chart comparing the frequency of occurrence of "big data", "business intelligence", "data mining", "data science" y "machine learning" in the historical records of this service.  
[www.mikelnino.com](http://www.mikelnino.com)

# Clarifying jargon: Chronology

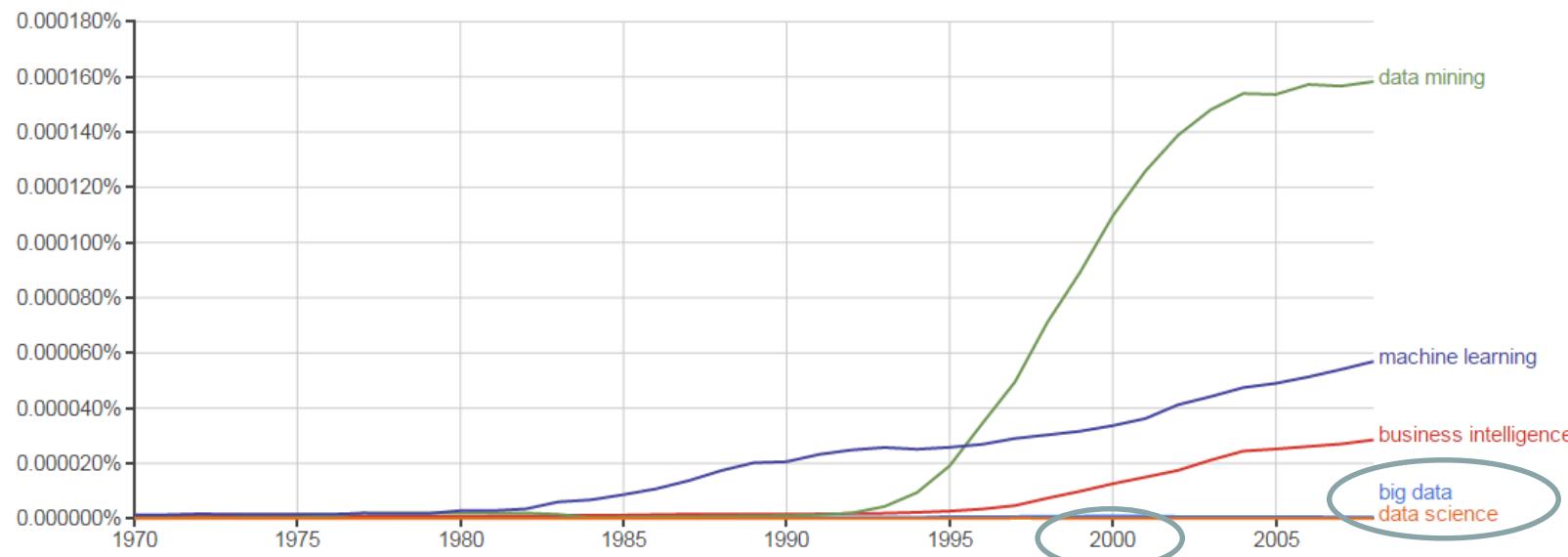
- Predictive analytics | Data Mining – 1990's
  - Business world started adopting Data Analytics
  - Used algorithms that are developed and applied on large amount of data



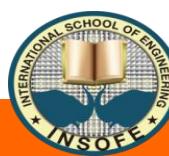
Source: Google Books chart comparing the frequency of occurrence of "big data", "business intelligence", "data mining", "data science" y "machine learning" in the historical records of this service.  
[www.mikelnino.com](http://www.mikelnino.com)

# Clarifying jargon: Chronology

- Big Data Analytics – 2000's
  - Focus was on computing on big volume of data
  - Google, Yahoo, Facebook, Twitter etc. used to apply algorithms on much large amount of data

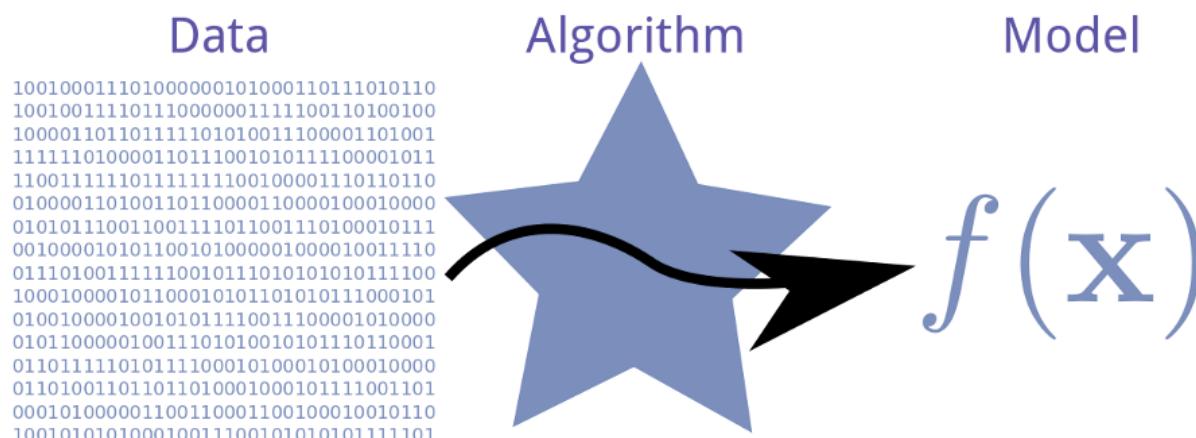


**Source:** Google Books chart comparing the frequency of occurrence of "big data", "business intelligence", "data mining", "data science" y "machine learning" in the historical records of this service.  
[www.mikelnino.com](http://www.mikelnino.com)

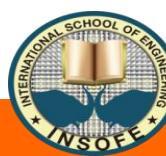


# Clarifying jargon: Chronology

- Data Science – over the last decade
  - Term is used to indicate a field where complex algorithms work on large volumes of data to solve important business problems to non-practitioners
  - A lot of emphasis on Visualization and story telling



Source: Google Books chart comparing the frequency of occurrence of "big data", "business intelligence", "data mining", "data science" y "machine learning" in the historical records of this service. [www.mikelnino.com](http://www.mikelnino.com)



# What are various sources of data?

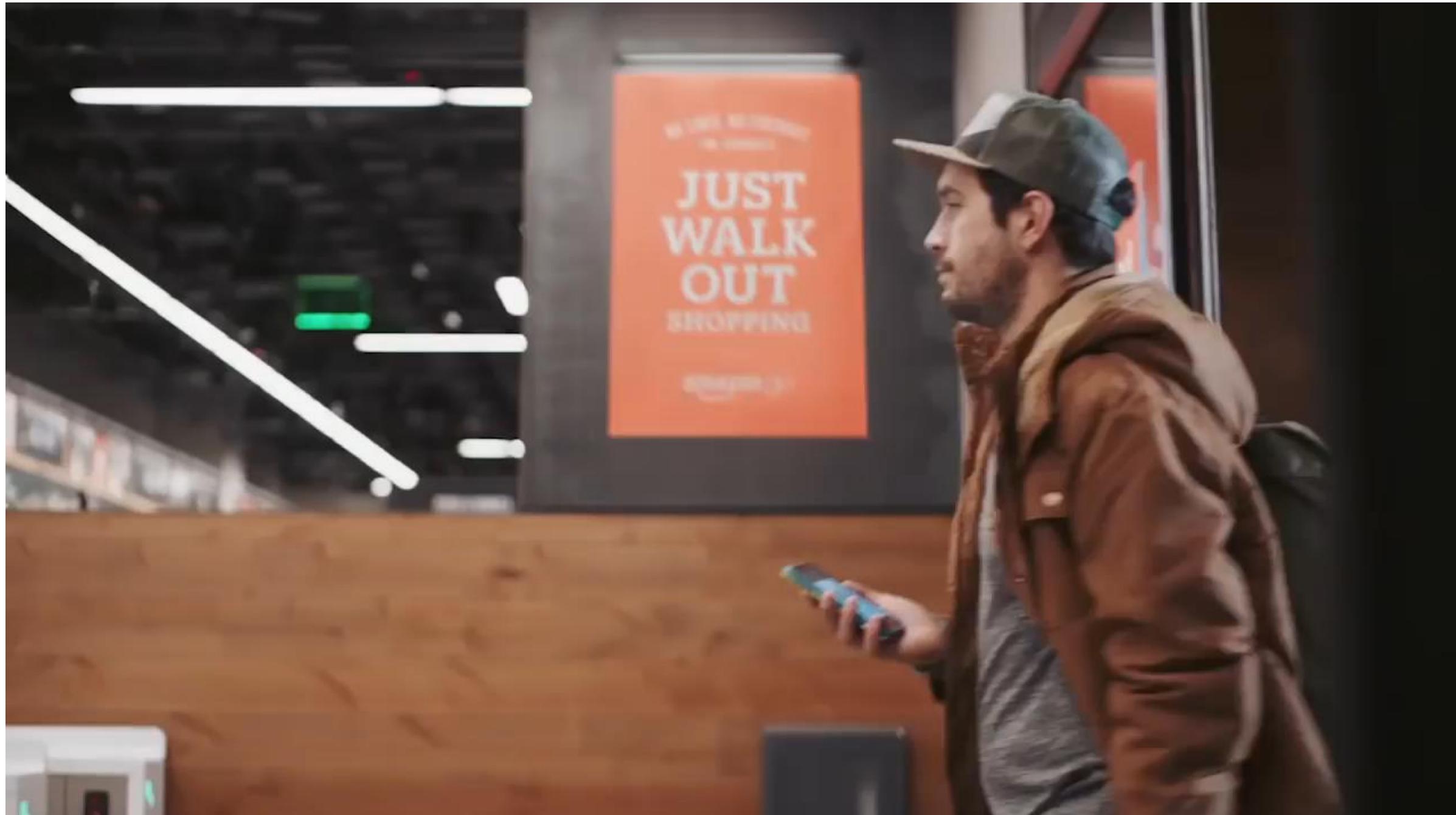


Conventional data sources



Evolution with internet of things

<https://www.youtube.com/watch?v=Q3ur8wzzhBU>



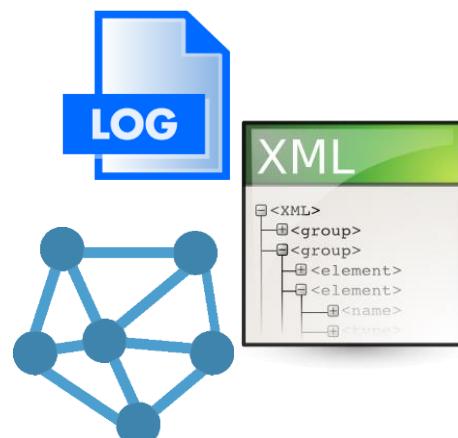
# What are various Types of data?

## Examples

Structured Data



Semi-structured Data



Unstructured Data



# What happens in an INTERNET MINUTE?



## Volume

## Variety

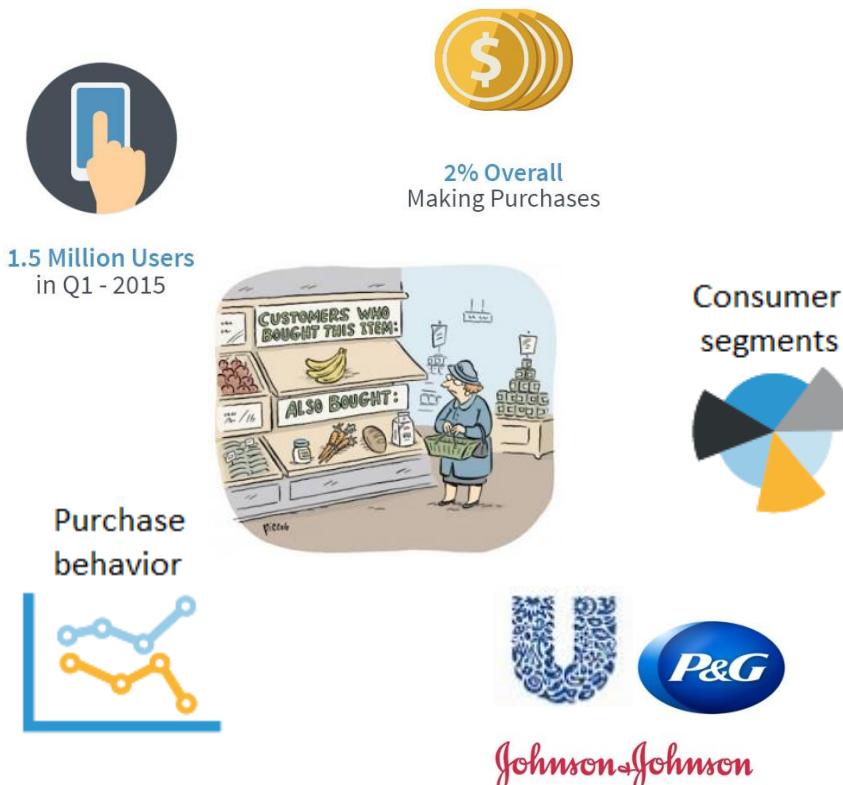
## Velocity

# Big versus Nano data

- Extremely small data sets (often found in healthcare, national calamities) are challenging
  - Most techniques fail to catch subtle trends
  - Deep learning, ensembles are able to tackle these problems better than before

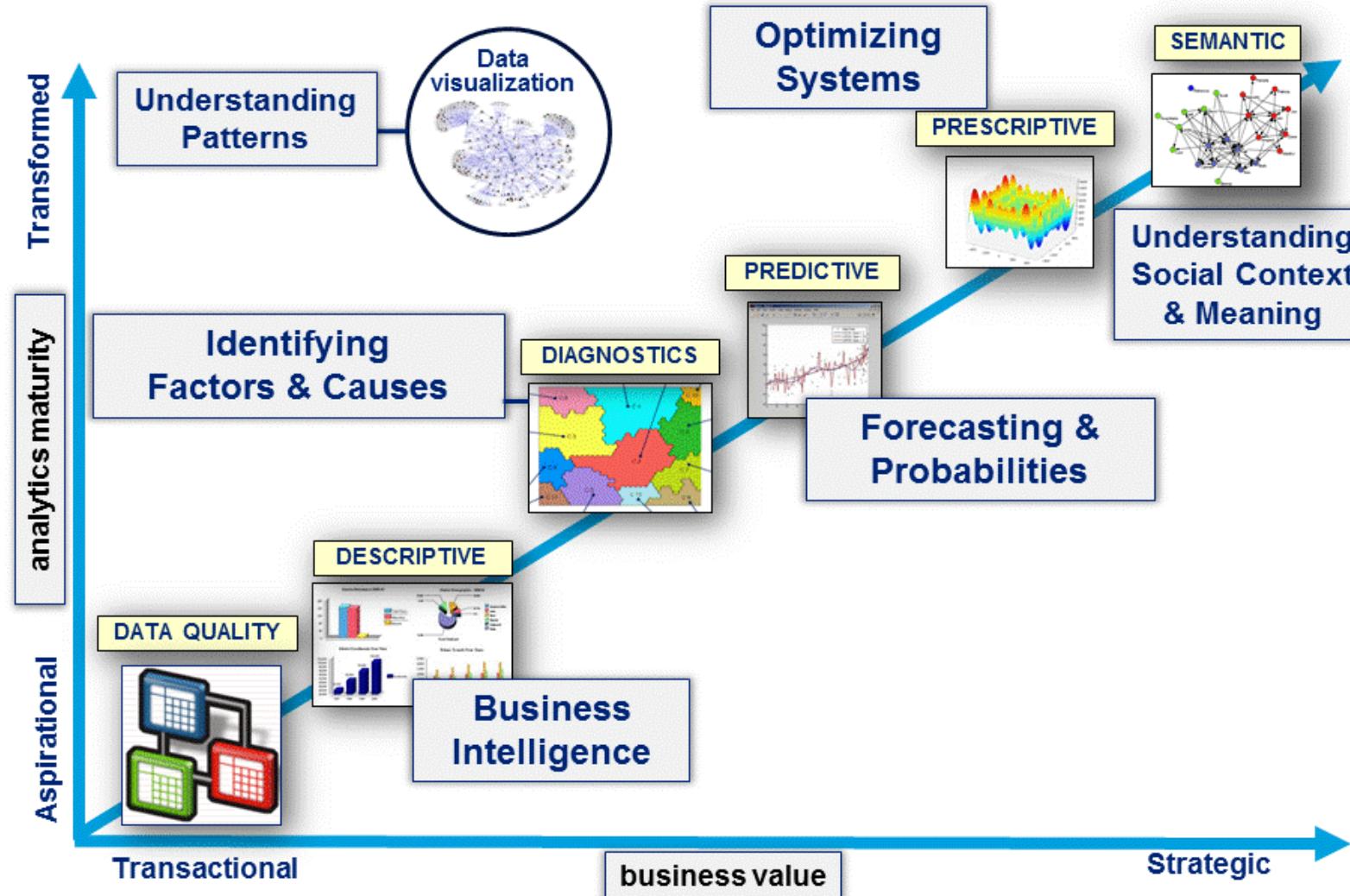
# What does the business need?

Quickly leverage data using data analytical methods in strategic decision making



- FMCG is a consumer driven business and enterprises continuously look out for ways to:
  - Innovate new products and services
  - Get visibility on trade promotions
  - Get Quick insights on consumer behaviour
  - Counter fierce competition, etc.

# BI to Predictive Analytics & Beyond



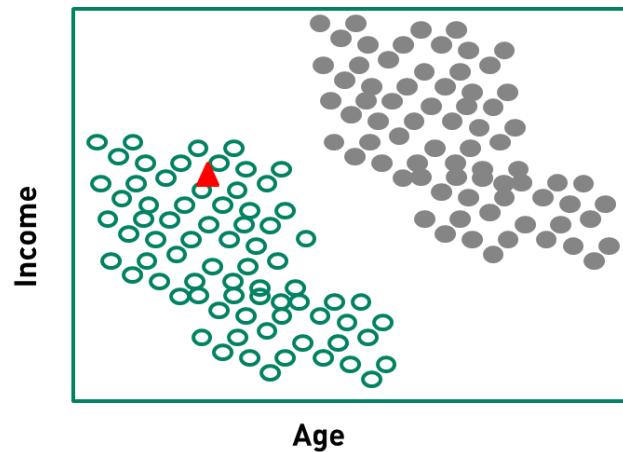
# Types of problems solved in Data Science

# Machine Learning Tasks

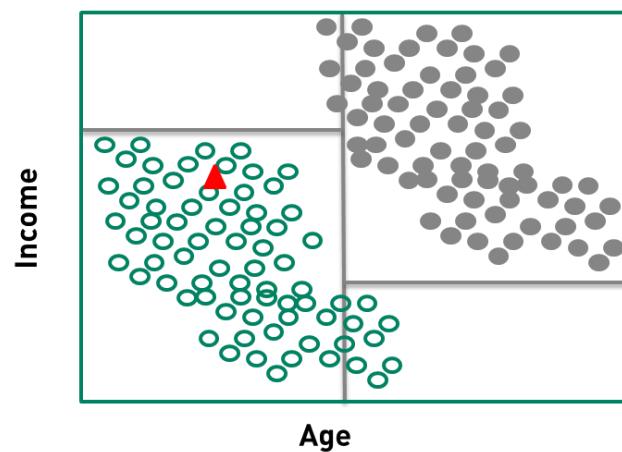
- Classification
- Regression
- Clustering
- Optimization

# SOME CLASSIFICATION METHODS

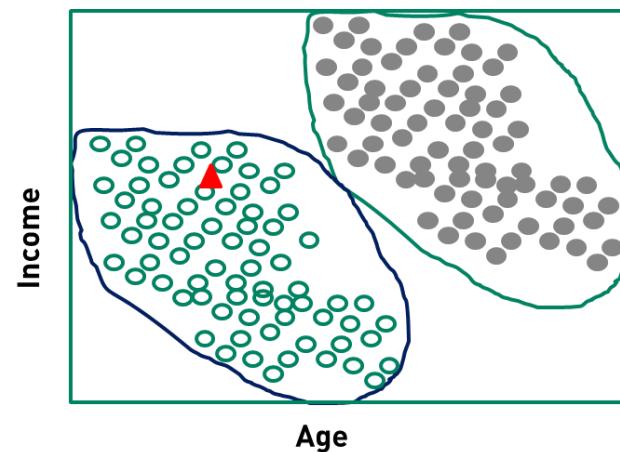
Binary Classification



Decision Trees

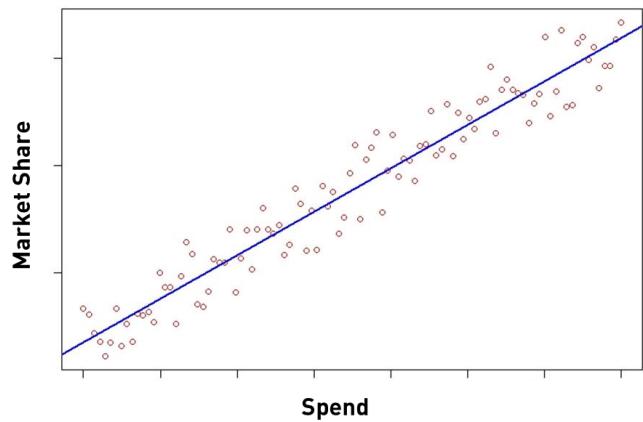


Naïve Bayes

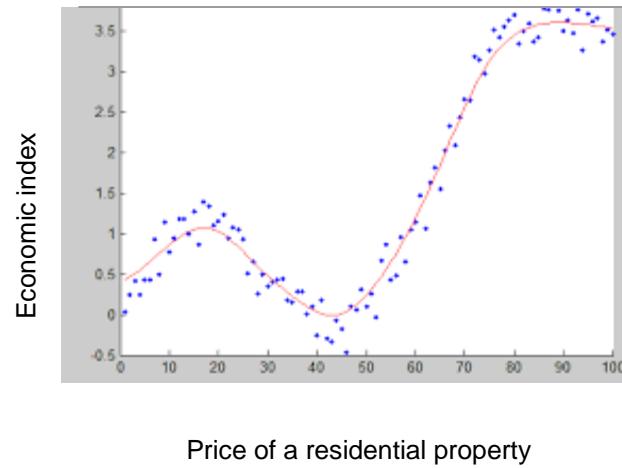


# SOME PREDICTION METHODS

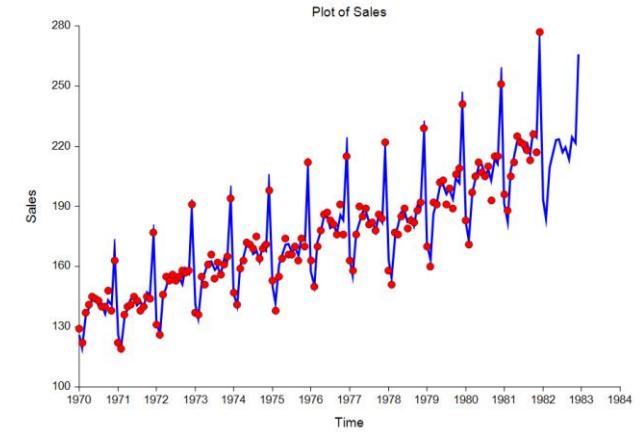
Linear Regression



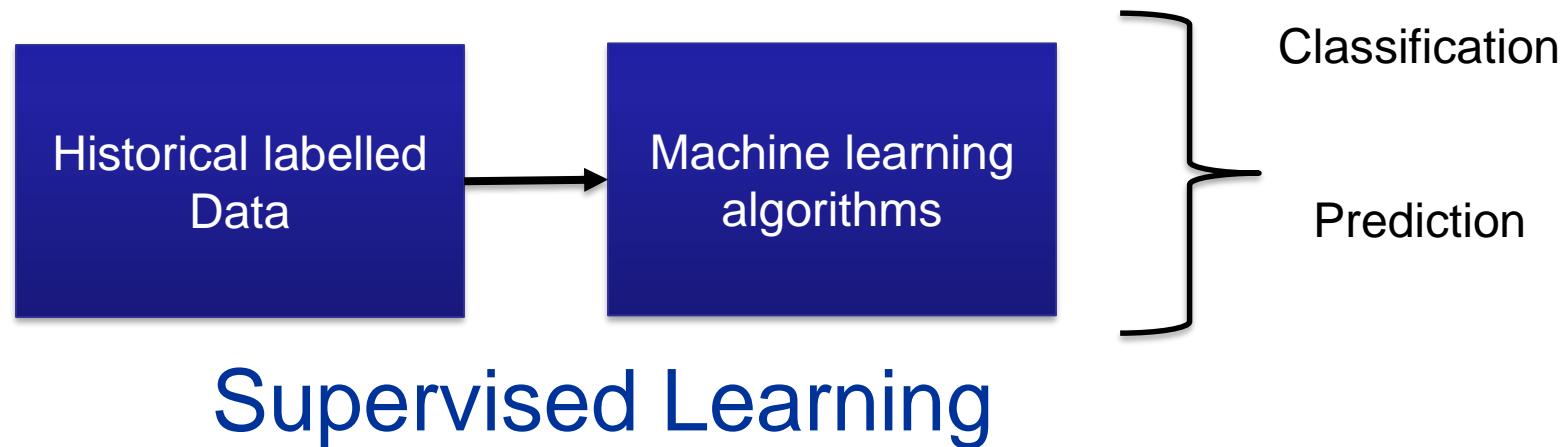
Polynomial Regression



Time Series Forecasting

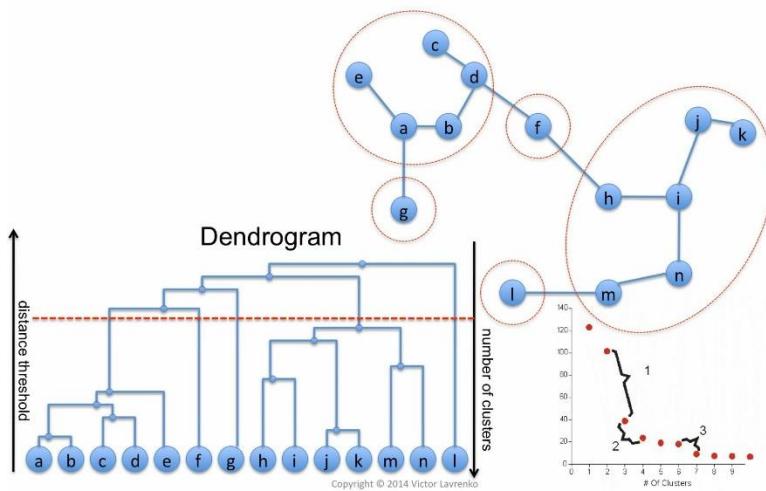


# Learning Method

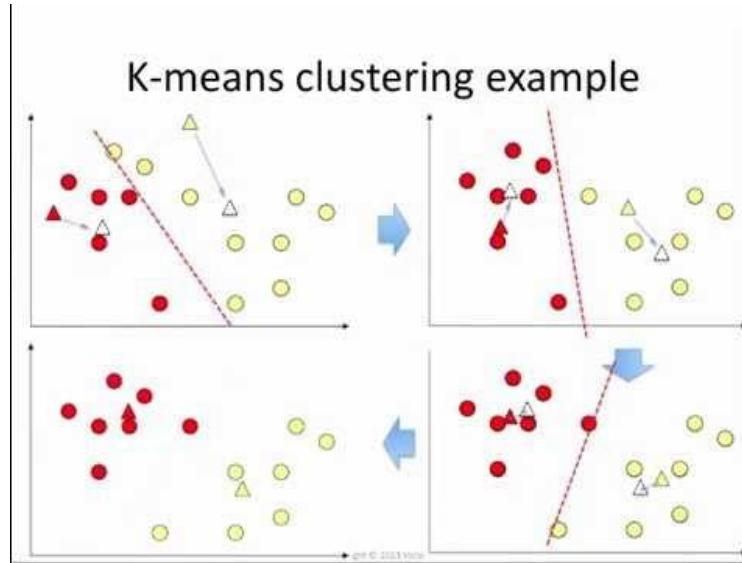


# Some Clustering methods

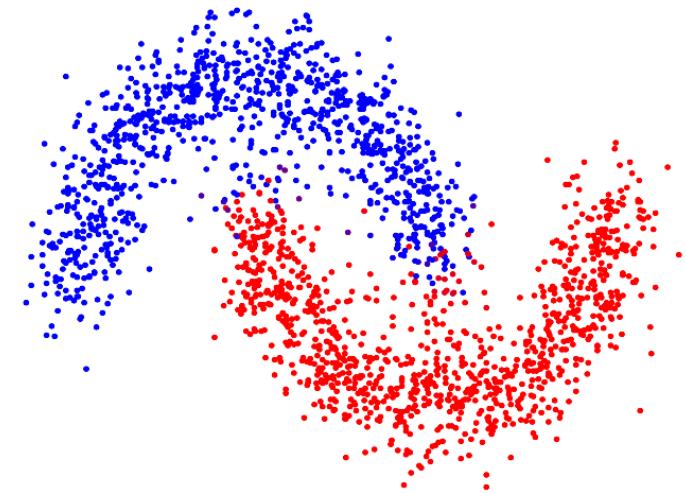
Agglomerative clustering: example



K-means clustering example



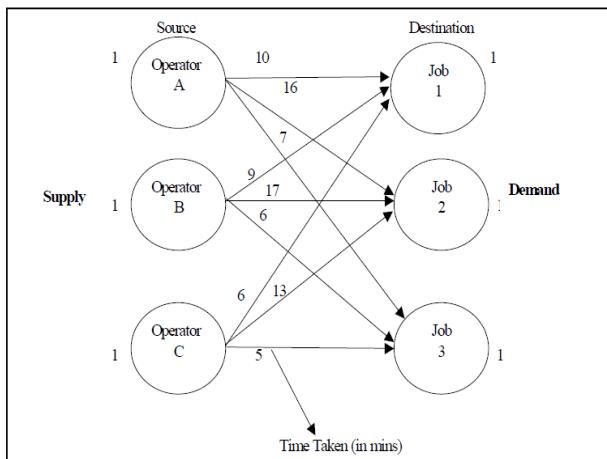
Spectral clustering example



## Unsupervised Learning

# SOME OPTIMIZATION METHODS

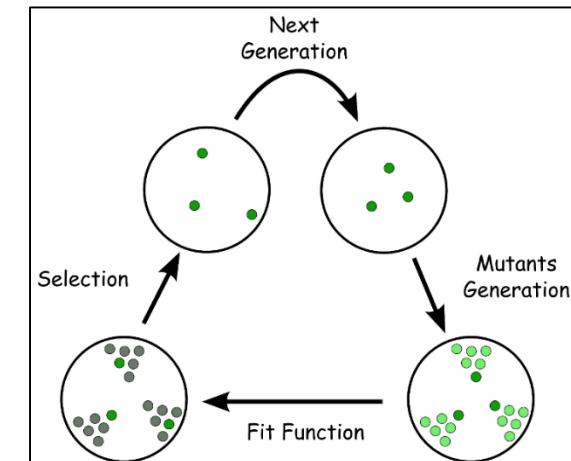
## Assignment Problem



## Transportation Problem

|        |                | Destination    |                |                |                | Supply |
|--------|----------------|----------------|----------------|----------------|----------------|--------|
|        |                | D <sub>1</sub> | D <sub>2</sub> | D <sub>3</sub> | D <sub>4</sub> |        |
| Source | S <sub>1</sub> | 19             | 30             | 50             | 10             | 7      |
|        | S <sub>2</sub> | 70             | 30             | 40             | 60             | 9      |
| Source | S <sub>3</sub> | 40             | 8              | 70             | 20             | 18     |
|        | Demand         | 5              | 8              | 7              | 14             |        |

## Genetic Algorithm



Key terms: Minimize / Maximize

# A Catalog of methods

## Descriptive statistics

- Mean, Median, Mode
- Correlations
- Sampling & Distributions
- T Test, F Test
- Normal distribution
- Poisson distribution. Etc.

## Prediction (Supervised)

- Multiple linear regression
- Neural nets
- Support vector machines
- Gradient Boosting
- Etc.

## Classification (Supervised)

- Logistic regression
- Decision trees
- Bayesian analysis
- Random Forest
- Etc.

## Classification (Unsupervised)

- Clustering
  - K-means
  - Hierarchical
- Association Rules, etc.



## Optimization

- Operations Research
- Linear Programming
- Genetic Algorithm, etc.



# WHY SO MANY METHODS?

There are many different methods for prediction and classification.

Each method has its advantages and disadvantages.

The usefulness of a particular method can depend on factors such as –

- the particular goal of the analysis
- size of the data set
- the types of patterns that exist in the data
- whether the data meet some underlying assumptions of the method

Different methods can lead to different results, and their performance can vary.

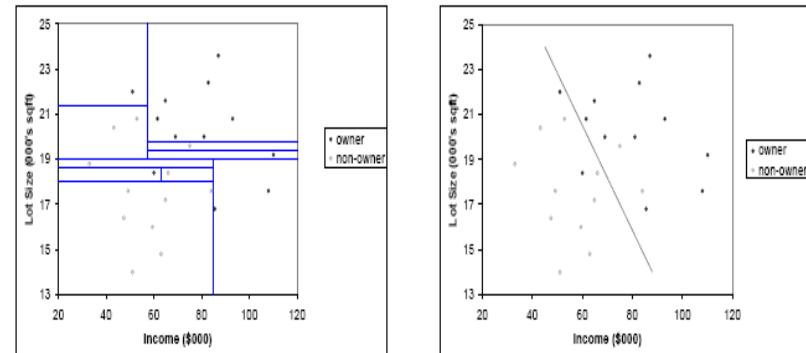
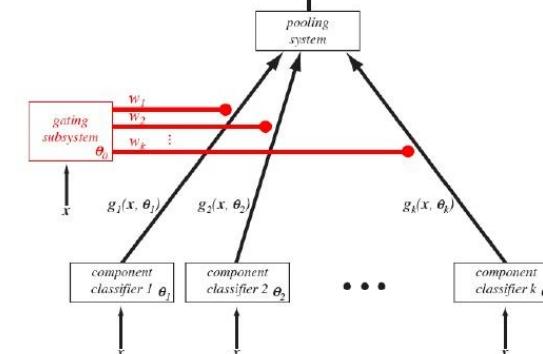


Figure 1.1: Two different methods for separating buyers from non-buyers

Mixtures of Experts



# The General Tasks of a Data Scientist

1. Get a little domain understanding
2. Define the problem statement well



# Understanding the Problem



# The Approach



# Results - The Streak!



# The General Tasks of a Data Scientist

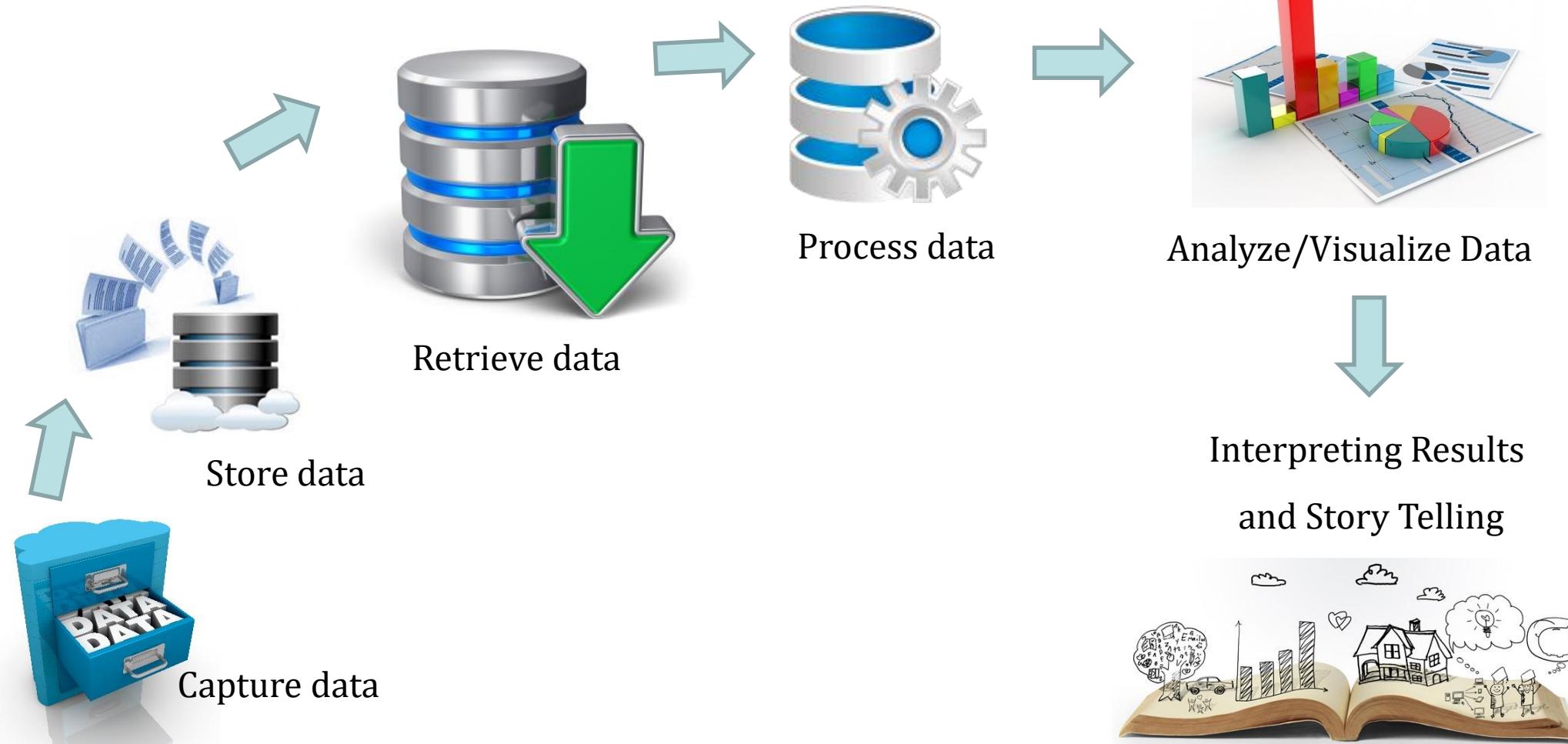
3. Pre-process data to fix data issues like duplicates, missing values, etc.
4. Visualize data to the extent possible for better understanding and to see basic patterns
5. Identify what kind of a problem it is (Prediction/Forecasting, Classification, Optimization and/or Managing Big Data)
6. Identify appropriate modeling techniques and build models
7. Analyze results and iterate, as needed; DO NOT trust software outputs blindly – Remember: Garbage In, Garbage Out
8. Visualize outputs and Communicate



# High-level steps

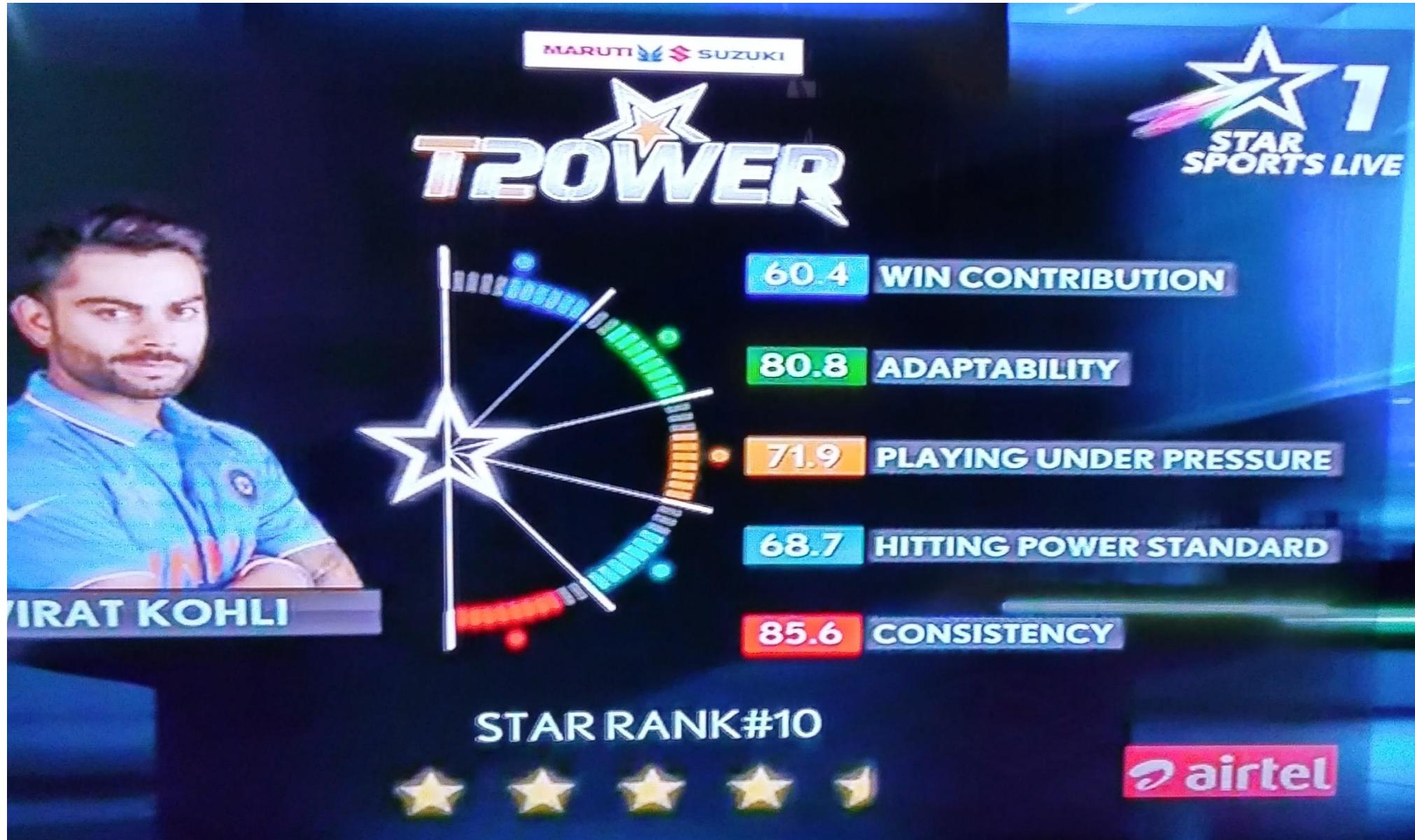
Insight : frequency of shoppers across marks

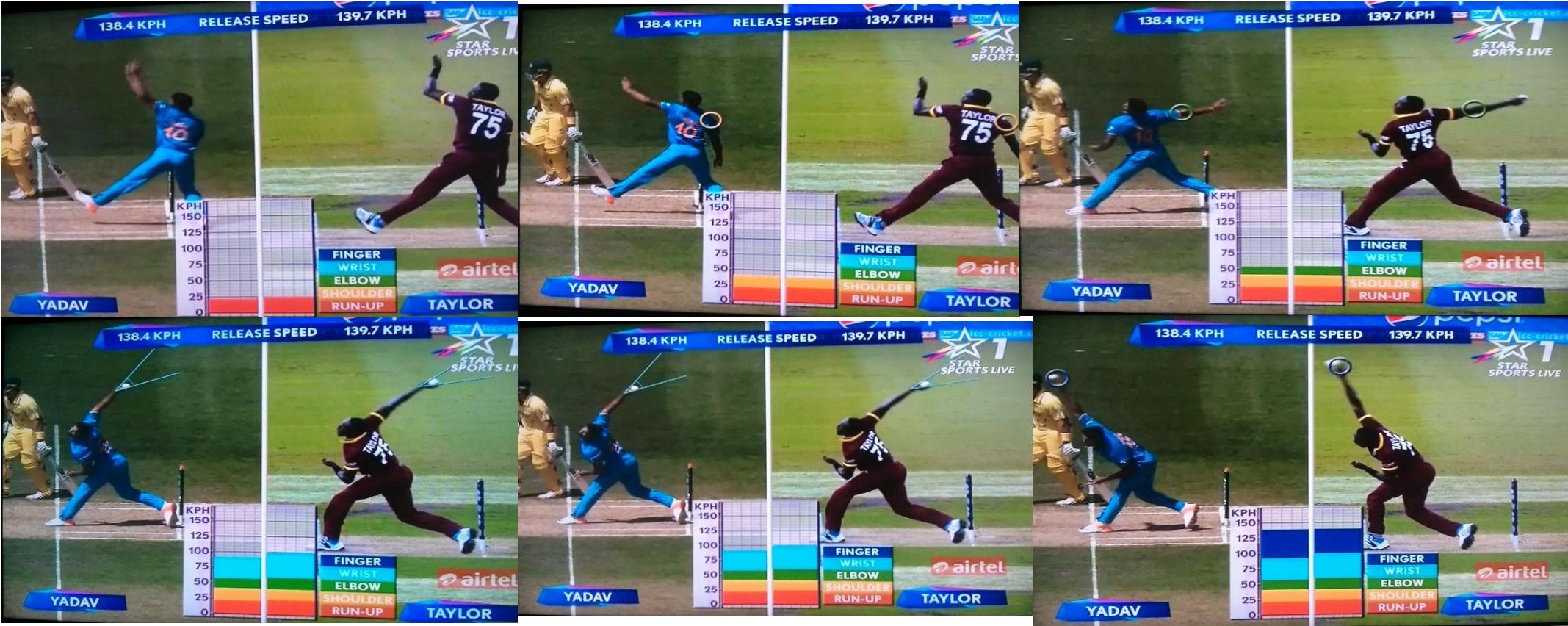
Broadly, here are the steps



# Applications and Significance of Data Science









# Descriptive Statistics



World T20 Finals 2014:  
Kohli scored 77;  
Spinners picked 3  
wickets

# Stock Market Trends

## Google searches can predict stock market crashes: Study

**London, July 29:** A rise in Google searches for terms relating to business and politics can predict a future stock market crash, researchers have claimed.

A team of researchers from Warwick Business School in the UK and Boston University in the US has developed a method to automatically identify topics that people search for on Google before subsequent stock market falls. Applied to data between 2004 and 2012, the method shows

that increases in searches for business and politics preceded falls in the stock market.

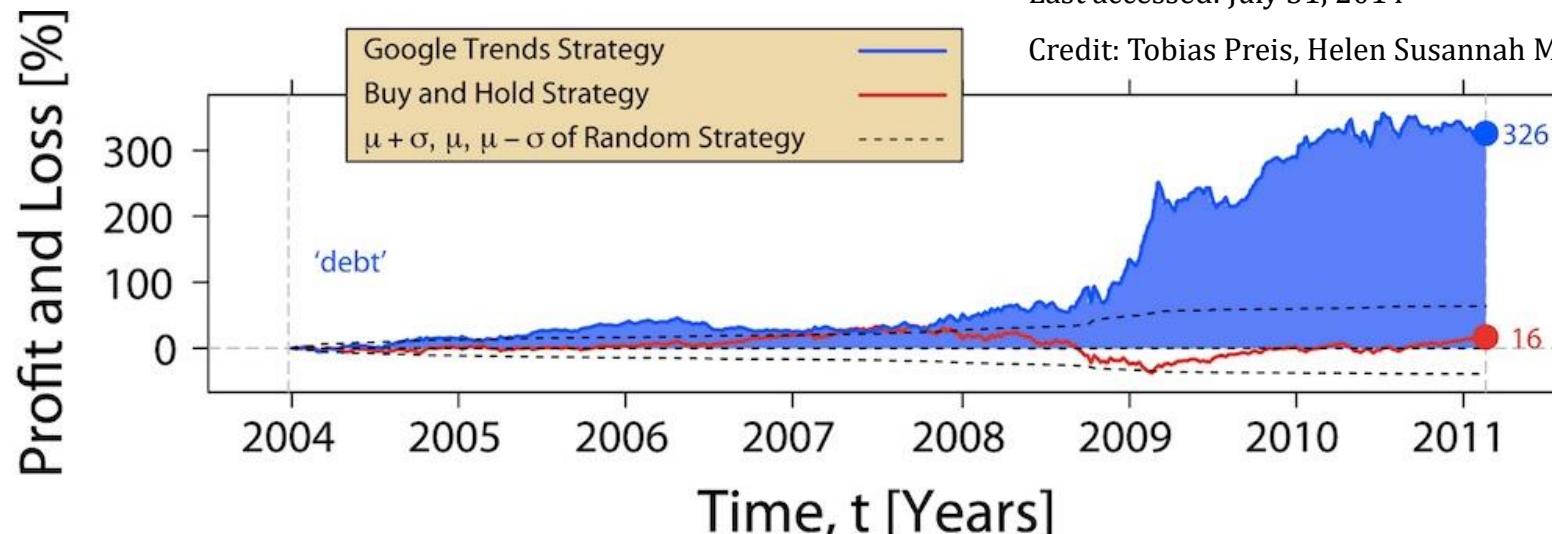
The researchers suggest that this method could be applied to help identify warning signs in search data before a range of real world events. "Search engines, such as Google, record almost everything we search for," said Chester Curme, research fellow at Warwick Business School and lead author of the study. "Records of these search queries allow us to learn about how people gather information online before making decisions in the real world. So there's potential to use these search data to anticipate what large groups of people may do," Mr Curme said.

In previous studies, Mr Curme and his colleagues, Tobias Preis and Suzy Moat of WBS, and H Eugene Stanley of BU, have demonstrated that usage data from Google and Wikipedia may contain early warning signs of stock market moves.

— PTI

Source: <http://epaper.deccanchronicle.com/articledetailpage.aspx?id=732781>

Last accessed: July 31, 2014



Source: <http://www.livescience.com/29016-google-predicts-stock-market.html>

Last accessed: July 31, 2014

Credit: Tobias Preis, Helen Susannah Moat, H. Eugene Stanley



# Medical non-adherence is likened to be an invisible epidemic

Non-adherence = # of doses not taken or taken incorrectly that jeopardizes patient's therapeutic outcome

Note: Lack of adherence to non-pharmacologic treatments, such as recommended lifestyle changes, is not included



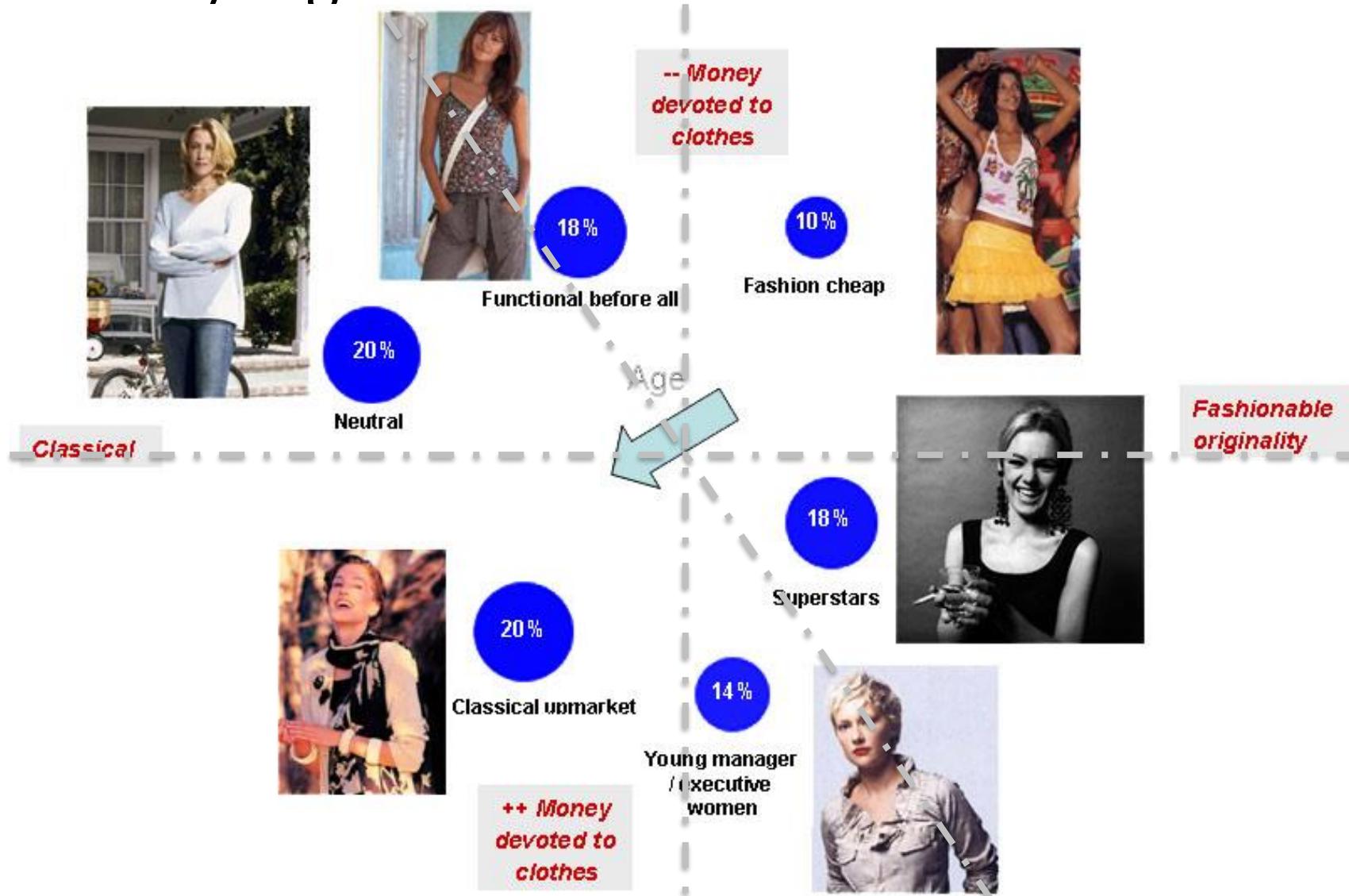
| Stakeholders  | Medical non-adherence impact  | Est. \$\$\$  | INDUSTRY ACTIONS   |
|---|---|--|--|
| <ul style="list-style-type: none"> <li>Insurers</li> <li>Employers</li> <li>Patients</li> </ul>                                   | <ul style="list-style-type: none"> <li>Increases healthcare costs due to disease-related complications</li> </ul>                         | <ul style="list-style-type: none"> <li>\$ 290 B</li> </ul> <p>Source: New England Healthcare Institute</p> | <ul style="list-style-type: none"> <li>Free drugs (<i>Netherlands</i>)</li> <li>Lower co-pays (<i>Aetna</i>)</li> <li>Reminders</li> <li>Reward points (<i>HealthPrize</i>)</li> </ul> |
| <ul style="list-style-type: none"> <li>Pharmaceutical companies</li> <li>Pharmacies</li> <li>Pharmacy benefit managers</li> </ul> | <ul style="list-style-type: none"> <li>Erodes profits due to prescriptions never filled and medications not taken often enough</li> </ul> | <ul style="list-style-type: none"> <li>\$ 188 B</li> </ul> <p>Source: Capgemini</p>                        | <ul style="list-style-type: none"> <li>More targeted reminders (<i>RxAnte</i> uses algorithms)</li> </ul>  |

The entire medical industry is scrambling to address this issue !!!

Advanced analytics may provide the insights to help solve this puzzle

Source: [http://pharmaceuticalcommerce.com/business\\_finance?articleid=26718](http://pharmaceuticalcommerce.com/business_finance?articleid=26718)

# Classification/Segmentation



Source: <http://www.bayesia.com/en/applications/marketing.php>

Last accessed: August 4, 2014



# Churn can be heartbreaking



Source: <http://princesswithapen.hubpages.com/hub/Is-my-girlfriend-cheating-on-me-Signs-of-a-cheating-girlfriend>

Last accessed: August 04, 2014

# Recommendation Engine



Write a Review 

Select Color 

Select Size (UK) 

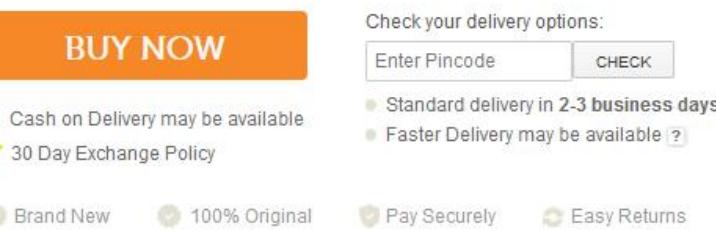
 SIZE CHART

 **1 OFFER**

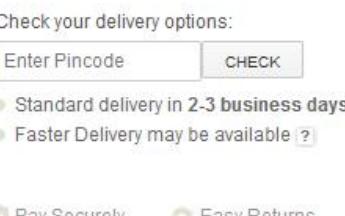
SAVE MORE: Shop for Rs.1199 or more and get 25% Off on Women's Clothing. See final price in cart. [View T&C](#)

[View details](#)

**BUY NOW**



Check your delivery options:  
Enter Pincode [CHECK](#)



**CUSTOMERS WHO VIEWED THIS PRODUCT ALSO VIEWED**

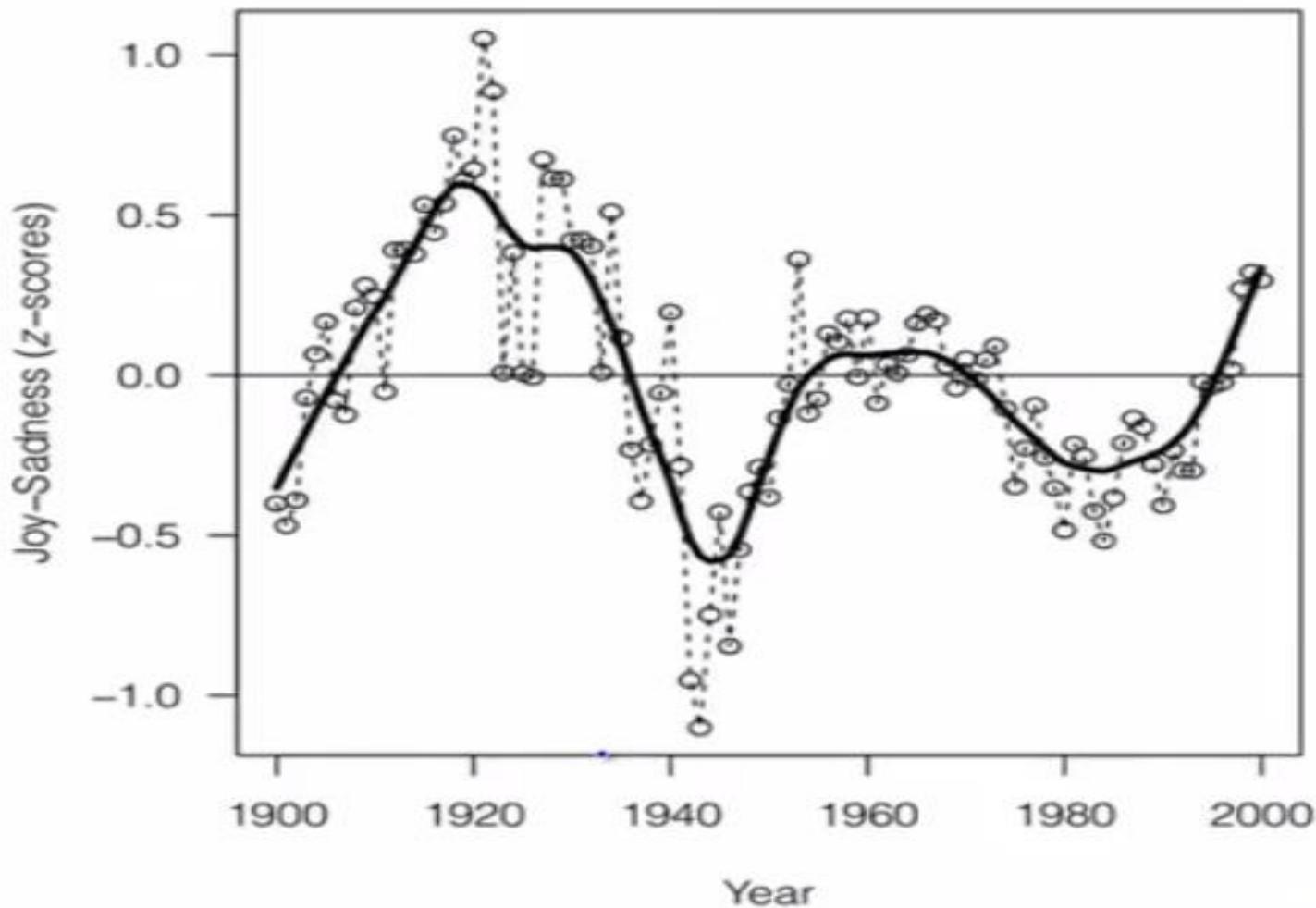


| French Connection | Flying Machine | Flying Machine | People   | People   |
|-------------------|----------------|----------------|----------|----------|
| Rs. 3499          | Rs. 1699       | Rs. 1699       | Rs. 1199 | Rs. 1199 |



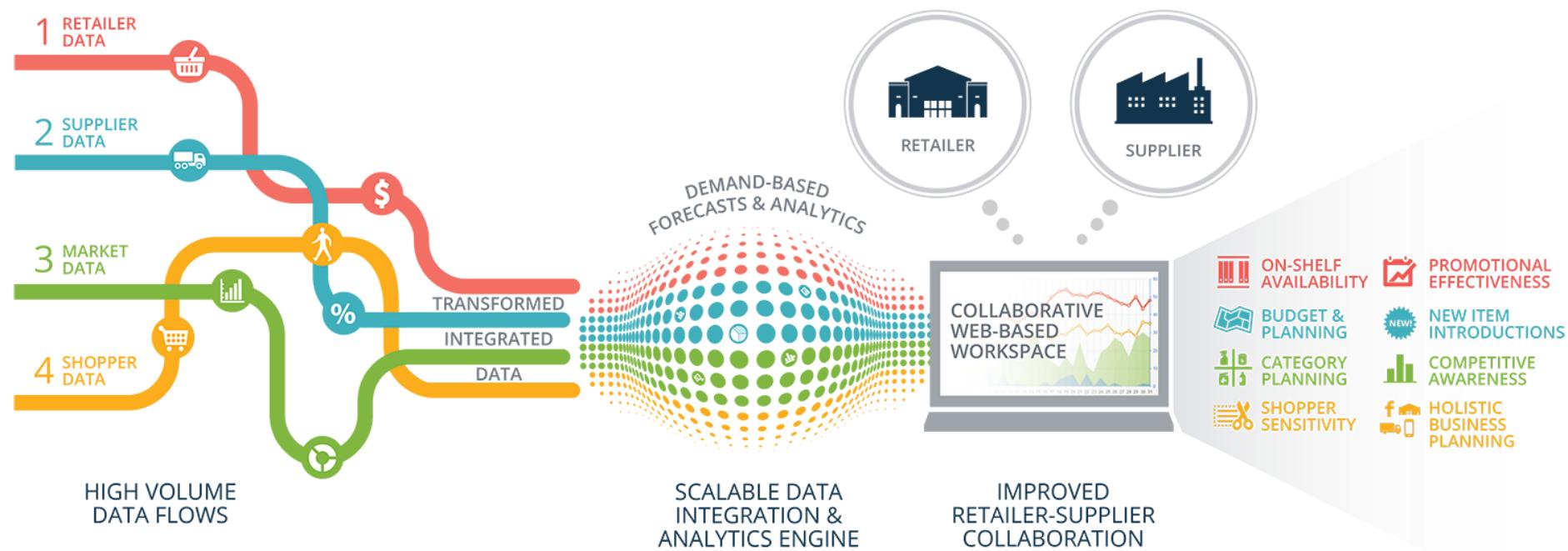
# Text Mining: Are Americans happy or unhappy?



# Customer Analytics

## Retail and Telecommunications

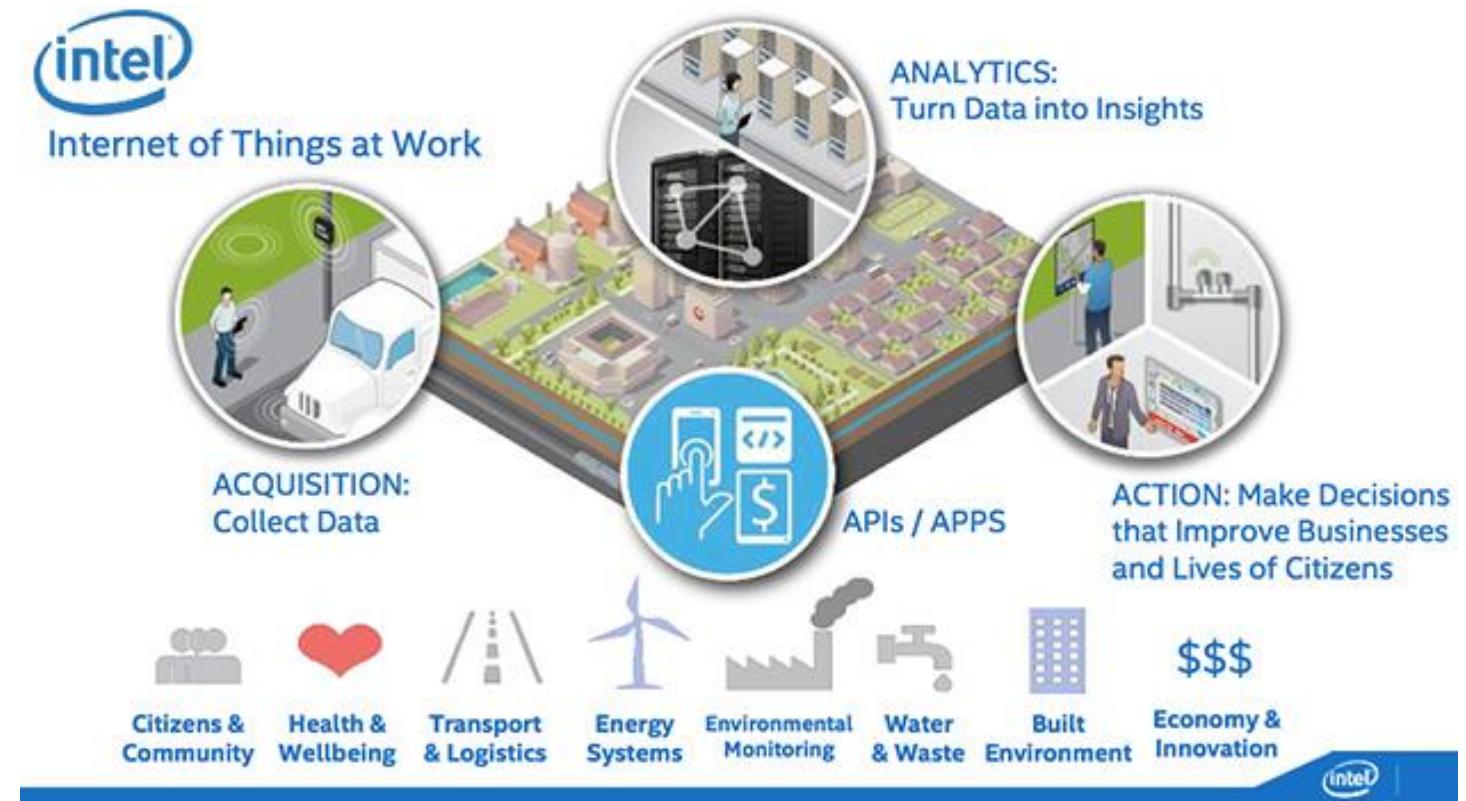
- Customer retention; Enhancing supply chain efficiency; Improving customer service quality; Planning store locations; Cross-selling and upselling; Recommendation systems; Sales forecasting



# Other important applications

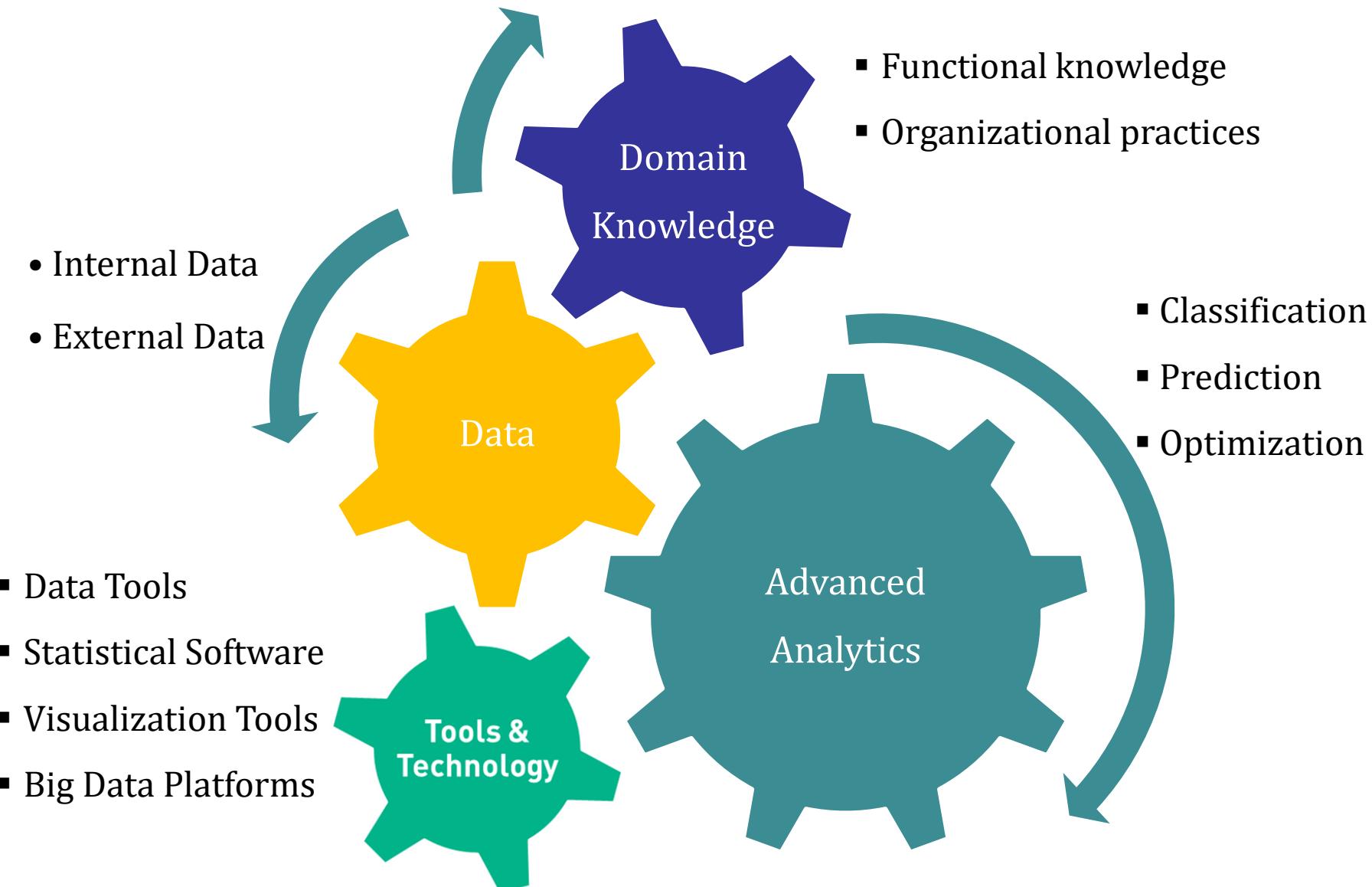
## Government

- Policy planning; Effective use of resources; Security against terrorist attacks; Effective policing by understanding crime patterns; Weather predictions; Calamity predictions

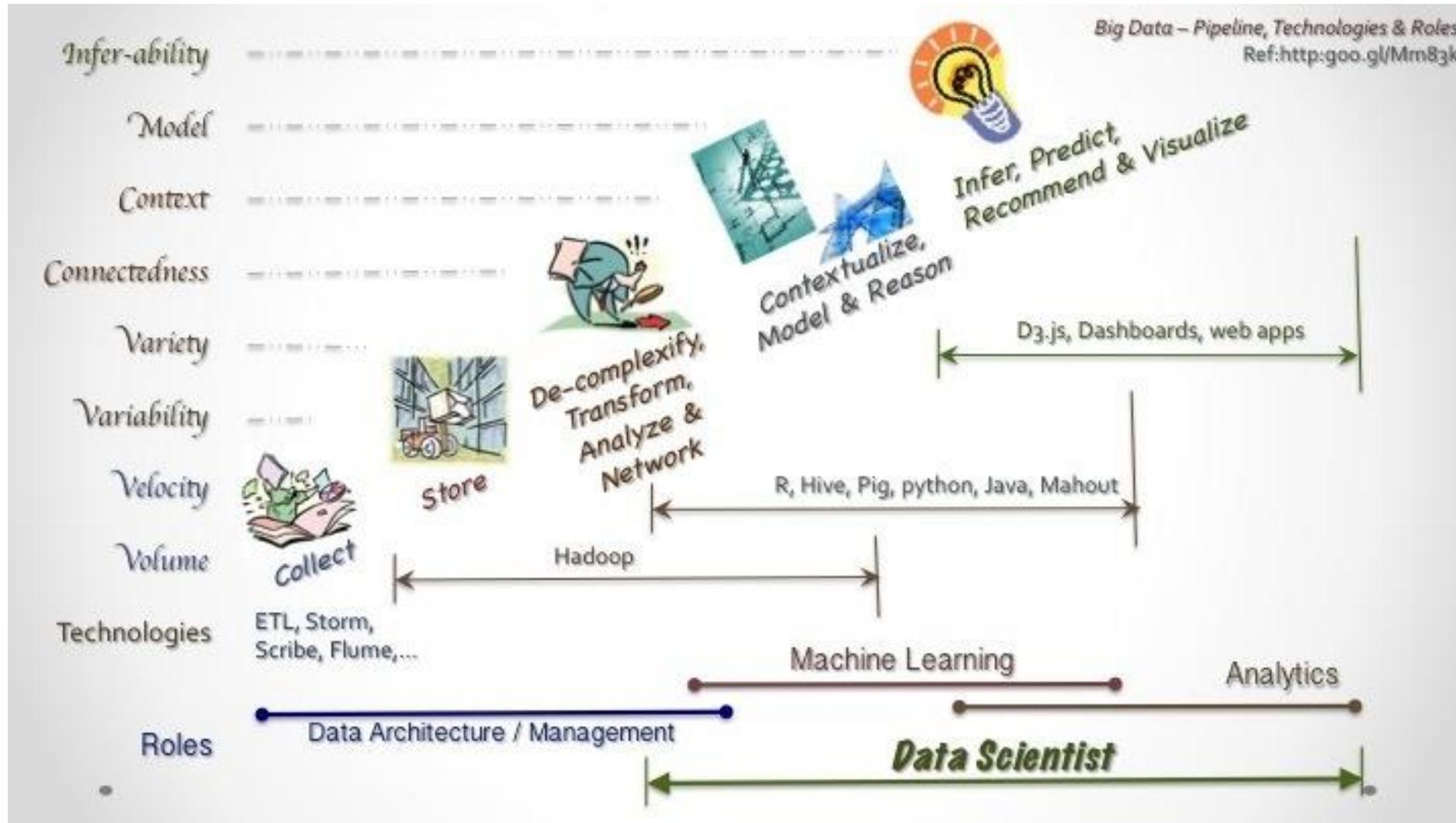




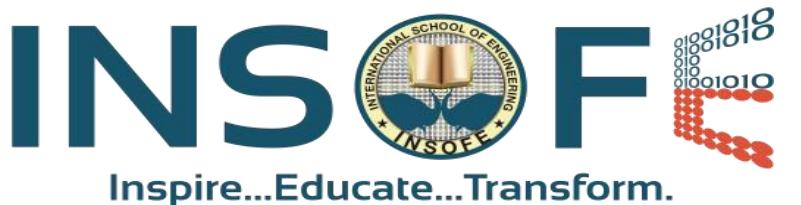
# You Can Make Better Decisions Using Data Science / Big Data Analytics



# Roles



<https://www.quora.com/What-are-the-most-valuable-skills-to-learn-for-a-data-scientist-now>



# HYDERABAD

## Office and Classrooms

Plot 63/A, Floors 1&2, Road # 13, Film Nagar, Jubilee Hills,  
Hyderabad - 500 033  
+91-9701685511 (Individuals)  
+91-9618483483 (Corporates)

## Social Media

- |             |   |              |
|-------------|---|--------------|
| Web:        | <a href="http://www.insofe.edu.in">http://www.insofe.edu.in</a>   | KnowledgeHut |
| Facebook:   | <a href="https://www.facebook.com/insofe">https://www.facebook.com/insofe</a>   | Main Road, T |
| Twitter:    | <a href="https://twitter.com/Insofeedu">https://twitter.com/Insofeedu</a>   | Layout, Beng |
| YouTube:    | <a href="http://www.youtube.com/InsofeVideos">http://www.youtube.com/InsofeVideos</a>   |              |
| SlideShare: | <a href="http://www.slideshare.net/INSOFE">http://www.slideshare.net/INSOFE</a>   |              |
| LinkedIn:   | <a href="http://www.linkedin.com/company/international-school-of-engineering">http://www.linkedin.com/company/international-school-of-engineering</a> |              |

*This presentation may contain references to findings of various reports available in the public domain. INSOFF makes no representation as to their accuracy or that the organization subscribes to those findings.*

# BENGALURU

## Office

Incubex, #728, Grace Platina, 4th Floor, CMH Road, Indira Nagar,  
1st Stage, Bengaluru – 560038  
+91-9502334561 (Individuals)  
+91-9502799088 (Corporates)

## Classroom

KnowledgeHut Solutions Pvt. Ltd., Reliable Plaza, Jakkasandra Main Road, Teacher's Colony, 14th Main Road, Sector – 5, HSR Layout, Bengaluru - 560102