

Data Science and its role in Big Data analytics

WE ALL CAN BE DATA SCIENTISTS NOW!

- I. Properties of Data: What is a Big Data?
- II. What is Data Science?
- III. Data Science Process
- IV. Why Data Science is in a sudden boom?
- V. Who is a Data Scientist?

Kyungwon Kim

Assistant Professor
Department of International Trade
College of Global Political Science and Economics
Incheon National University

WE ALL CAN BE DATA SCIENTISTS NOW! -

- I. Properties of Data: What is a Big Data?
- II. What is Data Science?
- III. Data Science Process
- IV. Why Data Science is in a sudden boom?
- V. Who is a Data Scientist?

Properties of Data: What is Big Data?

Big Data vs. Small Data

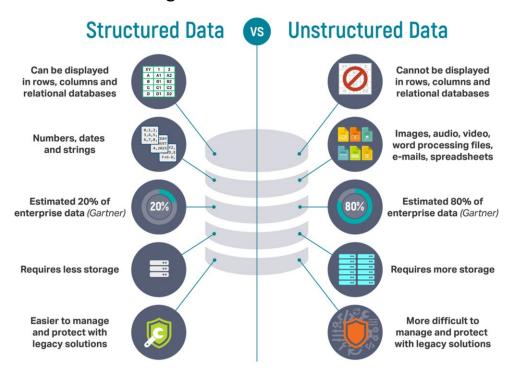
- Small Data: Getting machines to do what humans are good at.
- Big Data: Feeding an algorithm data to learn and predict something.
- => A Powerful Equation: Real Human Insight = Big Data + Small Data



Properties of Data: What is a Big Data?

> Structured vs. Unstructured Data

- Structured Data: the type of data that fits nicely into a relational database. It's highly organized and easily analyzed. Most IT staff are used to working with structured data.
- **Unstructured Data:** It doesn't fit nicely into a spreadsheet or database. It can be textual or non-textual. It can be human- or machine-generated.



Properties of Data: What is a Big Data?

Quantitative vs. Qualitative Data

- Quantitative Data: Numerical calculations and measurements.
- Unstructured Data: Sensations, feelings, and experiences.

Quantitative Data

- money
- time
- speed
- movement
- height
- length
- area
- volume
- weight
- temperature
- humidity
- pressure
- · sound level
- · categories (age, gender, occupation)
- positioning
- status

Qualitative Data

- · verbal and written feedback
- narrative story
 - first-hand (direct experience)
 - second-hand (telling someone else)
 - third-hand (outside story-teller)
- · visual images, drawings, or models
- · experiential sensations
- · descriptions of
 - colors
 - textures
 - smells
 - tastes
 - appearance
- beauty

- feelings
- intuition
- sensations
- choices
- values
- beliefs

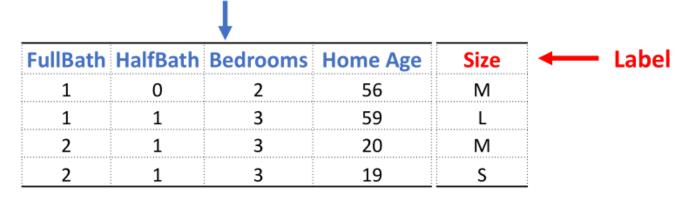


Properties of Data: What is a Big Data?

> Output (Y): Labelled vs Unlabelled

Lets say we want to Classify Houses by Size

Given Features or Feature Set



Unsupervised

SIZE is missing! We need to look for similarities in the data and group them into clusters.

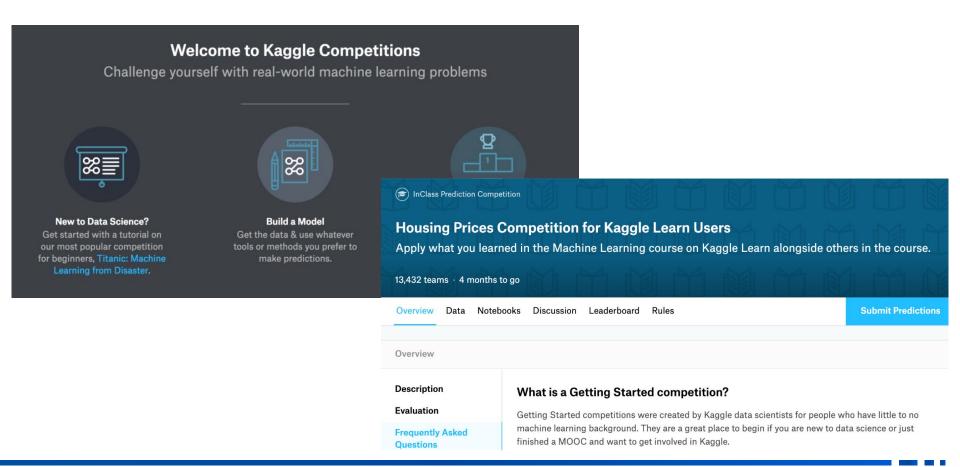
Supervised Learning

Use the labels to build a model. Model used to classify new house size based ONLY on the known feature set.

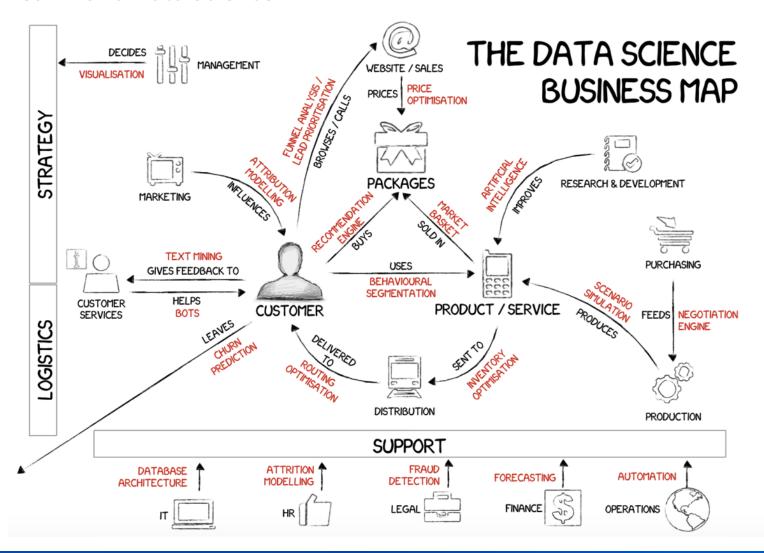
WE ALL CAN BE DATA SCIENTISTS NOW! -

- I. Properties of Data: What is a Big Data?
- II. What is Data Science?
- III. Data Science Process
- IV. Why Data Science is in a sudden boom?
- V. Who is a Data Scientist?

- The Real World Data Science is not a Kaggle Competition
 - It can be worthwhile to step back a little and realize what exactly your ultimate goal is.
 - The best performance might not be equivalent to a model yielding the best score in real.

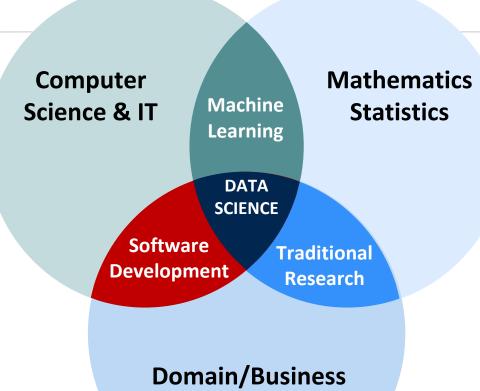


> The Real World Data Science



Data Science: Solving Problems with Data (Real Human Insight)

Computer Science,
Data Engineering,
Data Warehouse,
Pattern Recognition,
High Performance



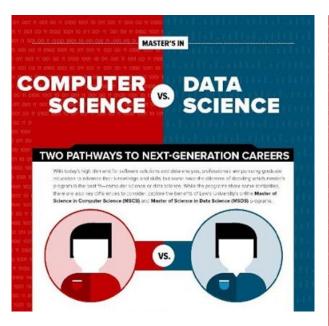
Statistical Learning,
Machine Learning,
Probability,
Numerical Techniques
to derive Insights

Domain/Business
Experience
& Knowledge

Problem and Objectives,
Domain Knowledge,
Business Experience,
Value to the Business

> Computer Science, Science, and Data Science

- Computer Science: The study of the theory and practice of how computers work.
- Science: Focusing on solving problems through the lens of the domain's scientific principles.
- Data Science: Interdisciplinary field involving computer science and statistics.





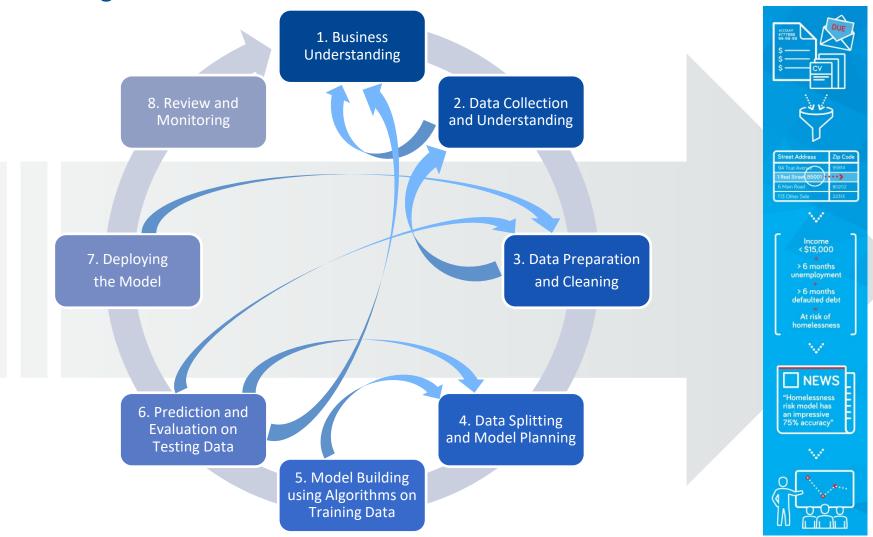


WE ALL CAN BE DATA SCIENTISTS NOW! -

- I. Properties of Data: What is a Big Data?
- II. What is Data Science?
- III. Data Science Process
- IV. Why Data Science is in a sudden boom?
- V. Who is a Data Scientist?

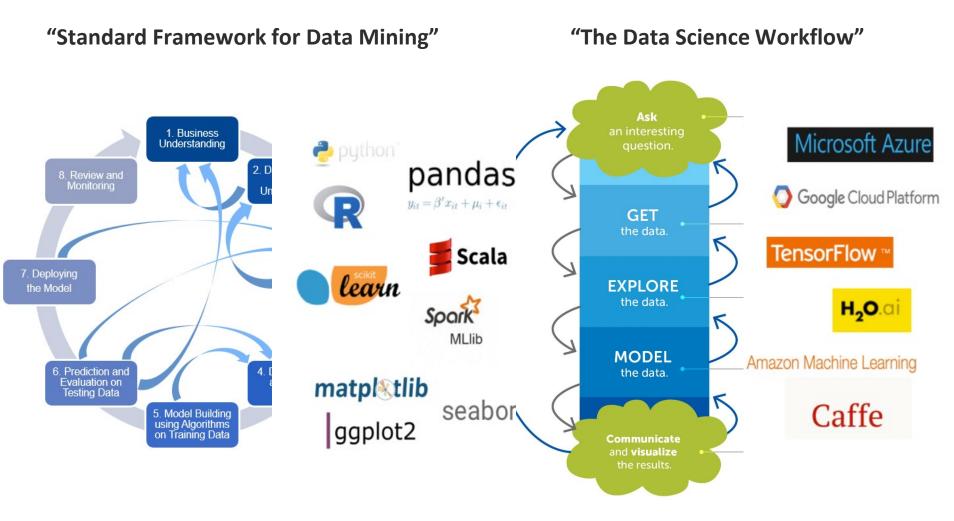
Data Science Process

> Getting from Raw Data to Outcomes



Data Science Process

→ Getting from Standard Framework to Data Science



WE ALL CAN BE DATA SCIENTISTS NOW! -

- I. Properties of Data: What is a Big Data?
- II. What is Data Science?
- III. Data Science Process
- IV. Why Data Science is in a sudden boom?
- V. Who is a Data Scientist?

processes, finding a 'correct' model

> Why use Machine Learning instead of Traditional Statistics?

Traditional Statistics Machine Learning A Data Science Continuum White-box modelling **Black-box modelling** simpler computation, emphasis on high computational complexity, emphasis introspection, form, causal effects and

Bayes Theorem:

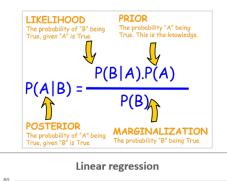
Thomas Bayes mid 1700's

Regression:

Legendre, Gauss and Galton early 1800's

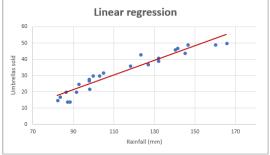
Neural Networks:

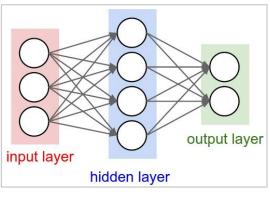
McCulloch and Pitts early 1940s



on speed and quality of prediction,

finding a 'performant' model



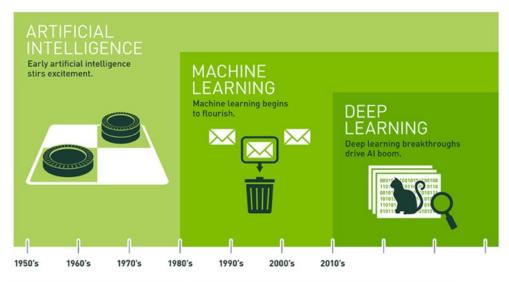


➤ Why use Machine Learning instead of Traditional Statistics?



- AI: Getting machines to do what humans are good at
- Machine Learning:
 Feeding an algorithm data
 to learn and predict something
- Deep Learning:

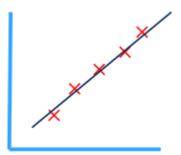
A type of machine learning



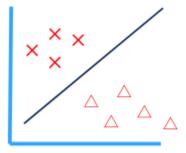
Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

> Solution directions to the black box problem

How much is the stock of Samsung Electronics tomorrow?



<u>Regression</u> – Looking for a statistical relationship across variables that may give us an estimate of a particular outcome. Will Samsung Electronics' stocks rise or fall tomorrow?



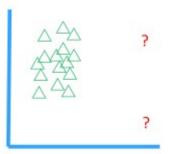
<u>Classification</u> – Similar to regression but looking for separations in the data given predefined classes. (Supervised)

 Are Samsung Electronics and Naver similar business companies?



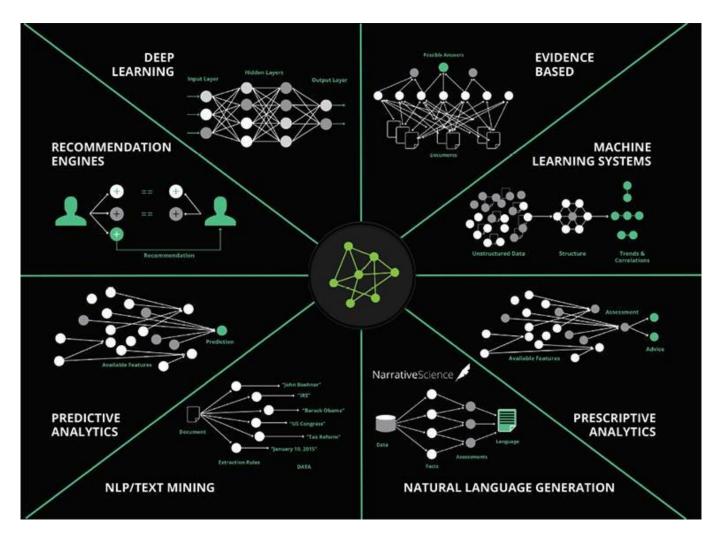
<u>Clustering</u> – Do not have predefined classes but trying to find groups or sets based upon data at hand. (Unsupervised)

Is the recent sharp drop in KOSPI outliers or normal?

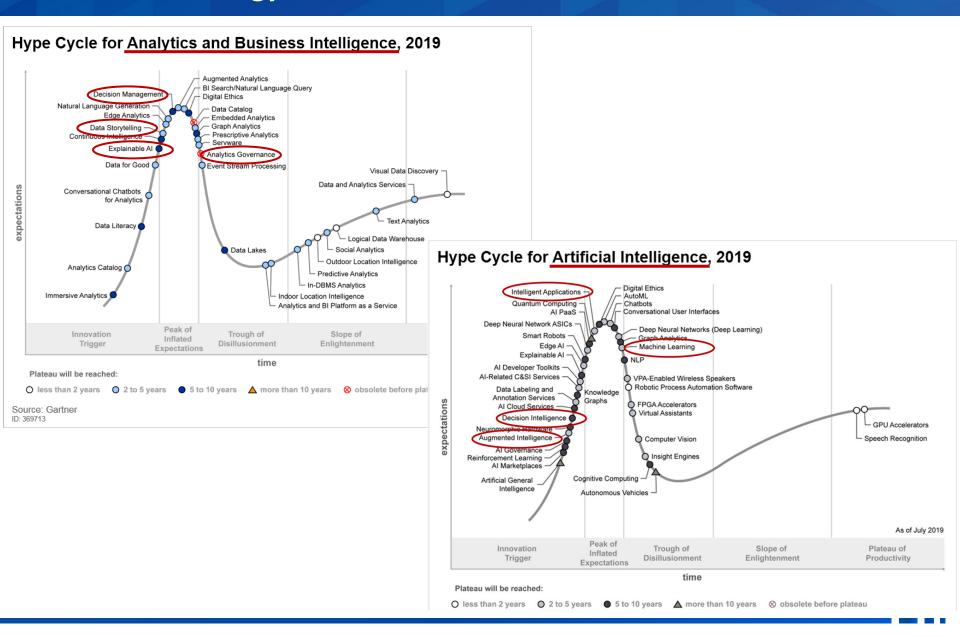


Anomaly Detection – Identification of outliers based upon expected ranges of data.

Different types of Machine Learning algorithms explained



Global "Technology" Trends



WE ALL CAN BE DATA SCIENTISTS NOW! -

- I. Properties of Data: What is a Big Data?
- II. What is Data Science?
- III. Data Science Process
- IV. Why Data Science is in a sudden boom?
- V. Who is a Data Scientist?

> What should we do?





MarketingDistillery.com is a group of practitioners in the area of e-commerce marketing. Our fields of expertise includemarketing strategy and optimization customer bracking and on site analytics predictive analytics and econometrics: data wavehouring and big data systems: marketing channel insights in Paid Search, ESO, Social, CRM and brand.



> Traditional Specialist of Data Science Team



Data Analyst (DA)

: Assist DS with domain understanding, data preprocessing and problem defining.



Data Scientist (DS)

: Prepares data, engineers features, most valuable skill: training models.



Data Engineer (DE)

: Data acquisition focus. Build data pipelines. Not uncommon to have 5:1 ratio DE:DS



Data Application Architect (DAA)

: Design complete solution; deploy and maintain models in production

> Typical Collaboration of Data Science Project

- Makes data science teams more productive
- Broad support for open source libraries in various languages





















Understand Business Objectives

ID Mapping Procure Training Data

Prepare Data and Build **Features**

Train, Tune, and Test **Models**

Deploy and Operationalize Models

Update Models

> What is a Project of Data Science? (Q) (Q) (Q)









TASKS

: In addition to advanced analytic skills, this individuals are also proficient at integrating and preparing large, varied datasets, architecting specialized database and computing environments, and communicating results.

MISSION

: A data scientist may or may not have specialized industry knowledge to aid in modeling business problems and with understanding and preparing data.

TALENT

: Creating value from data requires a range of talents from data integration and preparation, to architecting specialized computing/database environments, to data mining and intelligent algorithms.

RESPONSIBILITY

: An individual responsible for modeling complex business problems, discovering business insights and identifying opportunities through the use of statistical, algorithmic, mining and visualization techniques.

Data Science Roadmap

> We can be the best Data Science team









| Data Analysis | Data Analysis Cycle | | | | |
|----------------------------|--------------------------------------|---|---|---|---|
| | Data Visualization and Communication | | | | • |
| | Data Wrangling and Intuition | | | | |
| Mathematics | Linear Algebra | • | | | |
| | Numerical Analysis | • | | | |
| | Optimization | • | | | • |
| | Multivariate Calculus | • | | • | |
| Statistics | Probability and Statistics | | | • | • |
| | Experimental Design | | | | • |
| | Statistical Thinking and Algorithms | | | • | • |
| Artificial Intelligence | Machine Learning | • | | | • |
| | Deep Learning | • | | | |
| Computation | Databases and Distributed Systems | • | | | |
| | Programming Tools | | | | |
| | Algorithmic and Programming | • | | | |
| | Software Engineering | • | • | | |
| | Platform Understanding | • | • | | |

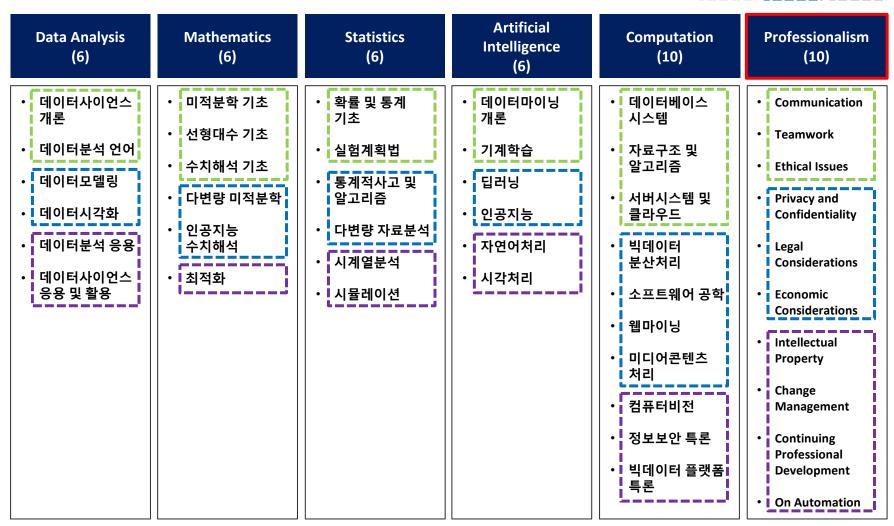
Not that



Data Science Curriculum

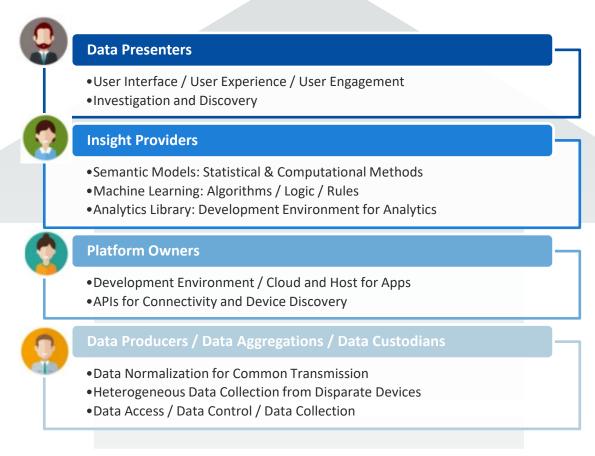
Related 44 Lectures (132 Credits)





> The Insights Revolution?

• A **Data Economy** is a global digital ecosystem in which data is gathered, organized, and exchanged by a network of vendors for the purpose of deriving value from the accumulated information.

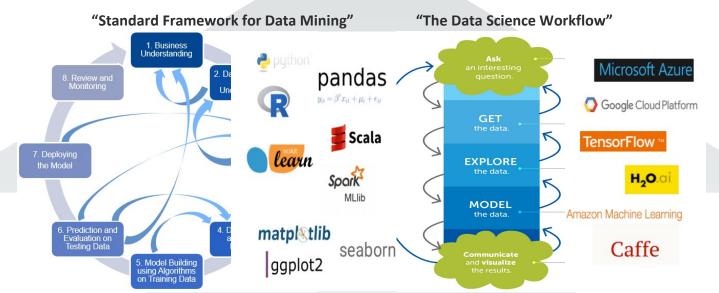


> The Insights Revolution?



Data Presenters

- User Interface / User Experience / User Engagement
- Investigation and Discovery





Insight Providers



Platform Owners



Data Producers / Data Aggregations / Data Custodians

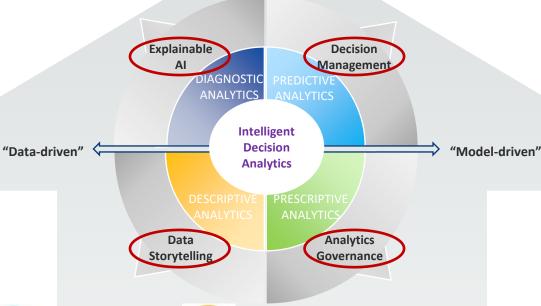
- Machine Learning: Algorithm APIs for Connectivity and D
- Analytics Library: Development Environment for Ana
- Semantic Models: Statistical Development Environment /• Data Normalization for Common Transmission
 - Heterogeneous Data Collection from Disparate Devices
 - Data Access / Data Control / Data Collection

> The Insights Revolution?



Data Presenters

- User Interface / User Experience / User Engagement
- Investigation and Discovery





Insight Providers

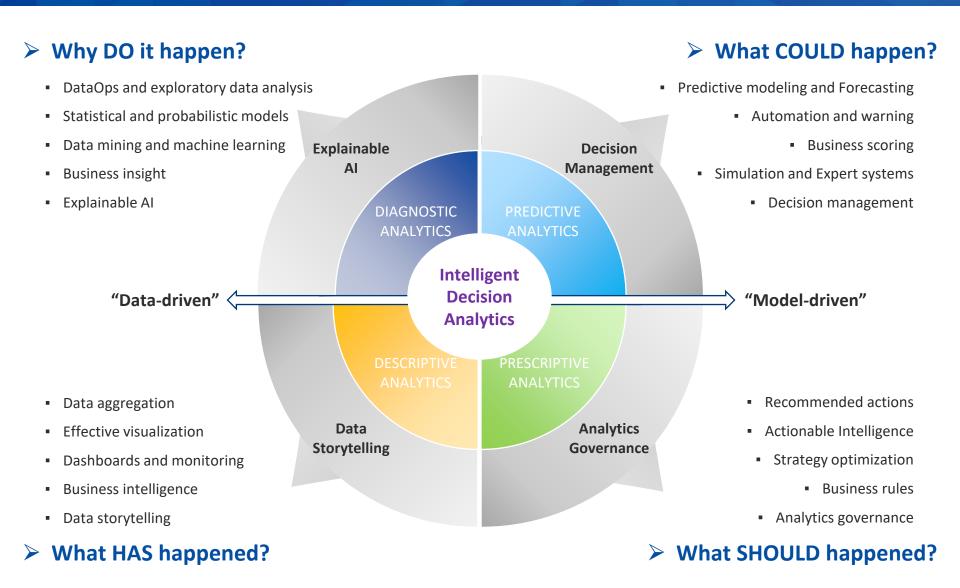


Platform Owners



- Semantic Models: Statistical Development Environment /• Data Normalization for Common Transmission
- Machine Learning: Algorithm APIs for Connectivity and December 1. APIs for Connectivity and December 1.
- Analytics Library: Development Environment for Analytics
- Heterogeneous Data Collection from Disparate Devices
- Data Access / Data Control / Data Collection

Hype Cycle for Analytics and Business Intelligence, 2019



THANK YOU

