

Lecture #1 | Introduction to big data and data science

SE271 Object-oriented Programming (2017)

Prof. Min-gyu Cho

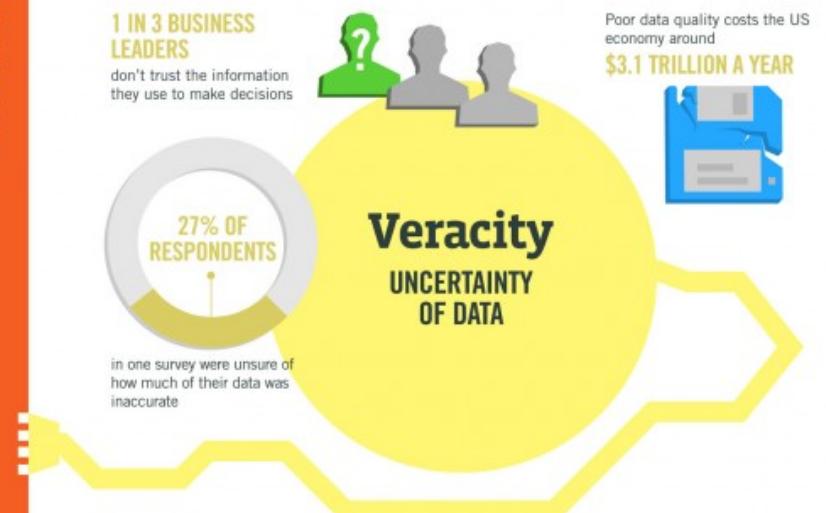
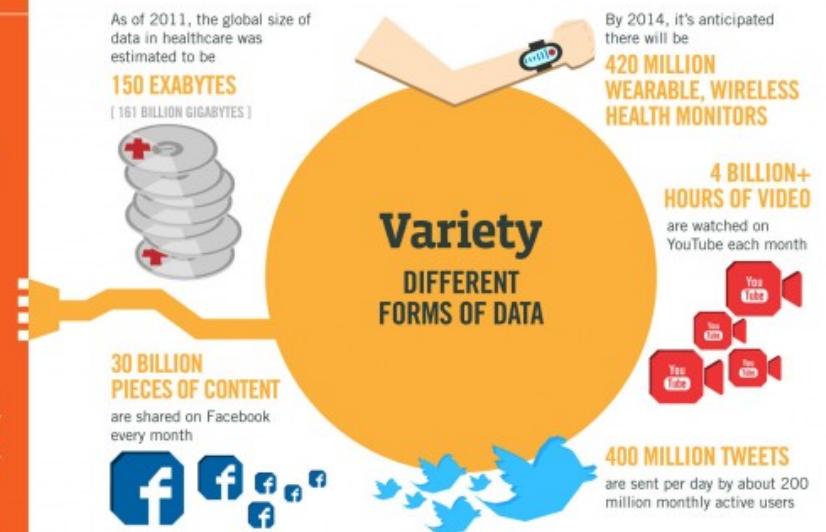
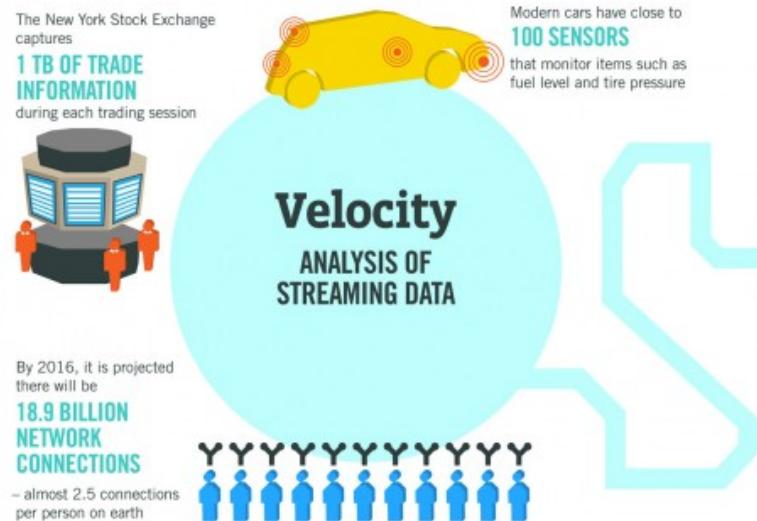
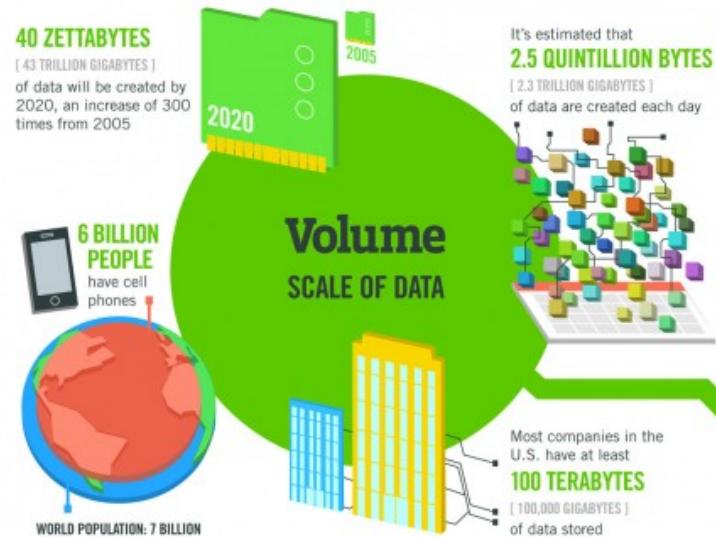
Big Data?



4V of Big Data

- **Volume:** large amount of data
- **Variety:** various types of data
 - Numbers/text/pictures/videos
 - Web info, SNS, IoT
- **Velocity:** speed of data flow
- **Veracity:** uncertainty or inconsistency of data





Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS

IBM

Example: Search engines

The screenshot shows a Google search results page for the query "dgist". The results include:

- dgist**
https://www.dgist.ac.kr/
사이트의 robots.txt 때문에 검색결과의 설명을 확인할 수 없습니다.
자세히 알아보기
- 디지스트**
기초학부 - 통근 및 셔틀버스 - 프로필 - 채용정보 - ...
- 대구경북과학기술원(DGIST)**
발급 대상자: 신입생: 통합계정 발급 화면 내
"통합계정 발급시 유의사 ..."
- 학사일정 | DGIST**
2017년 01월 학사일정. 2017년 01월 학사
일정을 안내해 드립니다.
- 대구경북과학기술원 - 나무위키**
https://namu.wiki/w/대구경북과학기술원
2017. 8. 16. - DGIST Blue 1, 2, 3[1] DGIST Gray 1, 2, 3[2] ... DGIST는 연구소로 시작한 기관으로서 기존의 연구부들
이 존재하고 석박사과정 교수와 학사과정 ...
- 대구경북과학기술원 - 위키백과, 우리 모두의 백과사전**
https://ko.wikipedia.org/wiki/대구경북과학기술원
재단법인 대구경북과학기술원(大邱慶北科學技術院, Daegu Gyeongbuk Institute of Science and Technology)은
2003년 12월 11일 제정된 대구경북과학기술원법(법률 제699호)에 따라 설립된 미래창조과학부 산하 기타공공기관이
다. 약칭은 DGIST이다. 2004년 국립대학원으로 출범한 DGIST는 2011년 대학원 석박사 과정을 개설했고, 2014년 학
부과정을 개설해 교육을 시작했다. [위키백과](#)

채용정보
600, DGIST 제2차 계약직 연구원 채용공
고, 17.08.07, 1343, 첨부파일.

탑메뉴
종합체육관 시설현황. 홈페이지 :
<https://sport.dgist.ac.kr/index...>

DGIST 정보통신융합공학전공
DGIST 정보통신융합공학전공. 전공소개 ·
전공소개 · 연혁 · 오시는길 ...

대구경북과학기술원
대한민국 대구광역시의 국립대학교

재단법인 대구경북과학기술원은 2003년 12월 11일 제정된 대구경북과학기술원법에
따라 설립된 미래창조과학부 산하 기타공공기관이다. 약칭은 DGIST이다. 2004년 국
립대학원으로 출범한 DGIST는 2011년 대학원 석박사 과정을 개설했고, 2014년 학
부과정을 개설해 교육을 시작했다. [위키백과](#)

주소: 대구광역시 달성군 현풍면 상리 50-1
창당: 2004년 9월 7일
연락처: 053-785-1114
상징색상: 파랑, 회색

수정 제안하기 · 이 비즈니스의 소유주인가요?

리뷰
Google 리뷰 10개

리뷰하기 사진 추가

함께 찾은 검색어

5개 이상 항목 더보기

Example: Netflix recommender system

Everything is a Recommendation



Over 75% of what people watch comes from our recommendations

Source: <https://www.slideshare.net/justinbasilico/recommendation-at-netflix-scale>

Example: 2012 US presidential elections

- Which e-mail is better?

Ready to fight? Please donate \$19 or more today, ahead of the FEC deadline.

Because you've saved your payment information, your donation will go through immediately:

[QUICK DONATE: \\$19](#)

- [QUICK DONATE: \\$35](#)
- [QUICK DONATE: \\$50](#)
- [QUICK DONATE: \\$100](#)
- [QUICK DONATE: \\$250](#)

Or donate another amount:

<https://donate.barackobama.com/June-Deadline>

Ready to fight? Please donate \$19 or more today, ahead of the FEC deadline.

Because you've saved your payment information, your donation will go through immediately:

[QUICK DONATE: \\$19](#)

- [QUICK DONATE: \\$35](#)
- [QUICK DONATE: \\$50](#)
- [QUICK DONATE: \\$100](#)
- [QUICK DONATE: \\$250](#)

Or donate another amount:

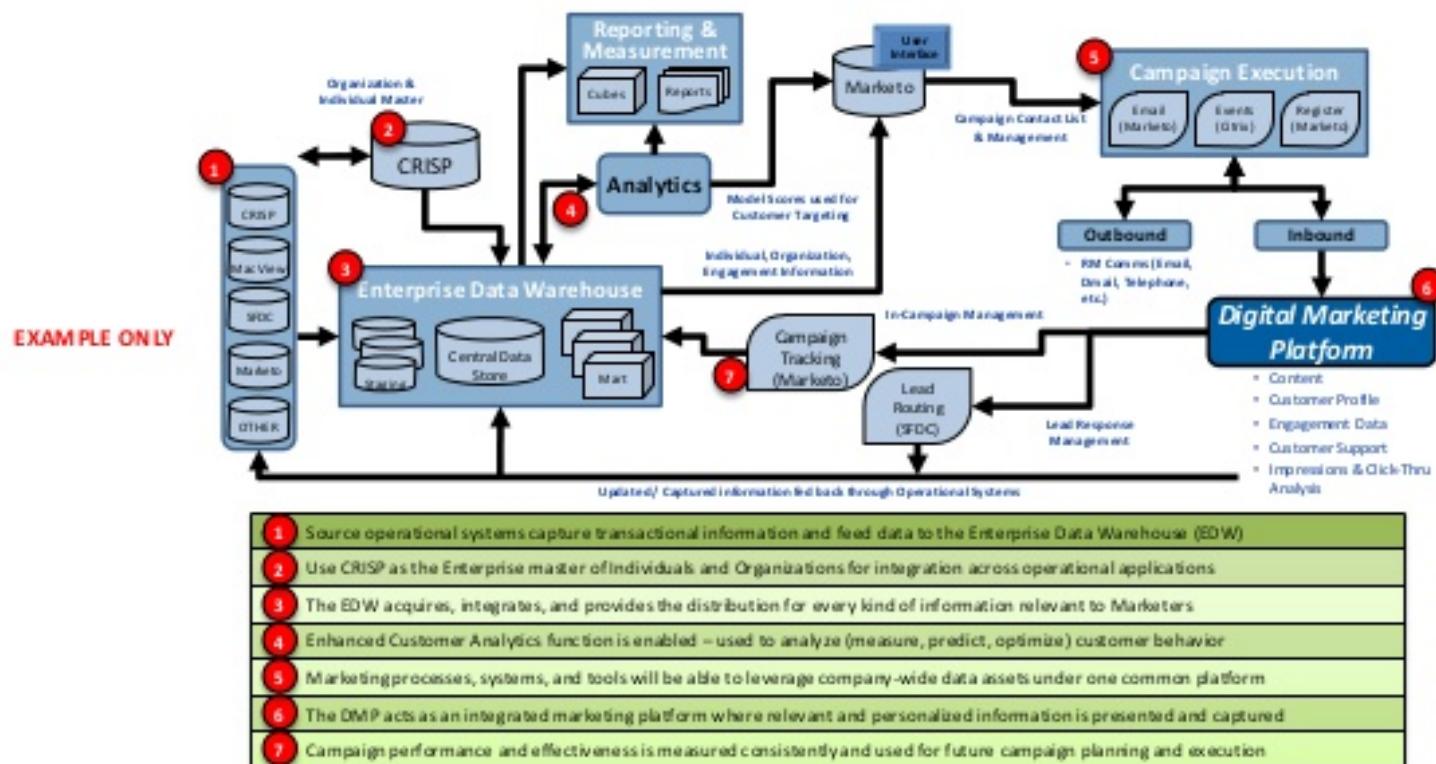
<https://donate.barackobama.com/June-Deadline>

- A/B test and re-test is IMPORTANT!!!

Source: <https://www.slideshare.net/jessday/obama-2012-lessons-from-a-datadriven-campaign>

Example: Digital Marketing

Part of the Future Vision Would be to Map the Enablement of the *Marketing Process* through the Marketing Data Platform.



Source: <https://www.slideshare.net/martinwalsh/example-digital-marketing-platform-strategy>

Example: Deep Learning

- Used 22 mil. sentences as a base and generated simulated data sets to increase speech recognition accuracy of Google Home

Table 1: *Speech recognition Word Error Rates(WERs)*

	Trained with the room simulator	Baseline system
Original Test Set	11.97 %	12.02 %
Simulated Noise Set	19.55 %	47.88 %
Device 1	21.98 %	50.14 %
Device 2	22.23 %	48.65 %
Device 3	22.05 %	56.27 %
Device 3 (Noisy Condition)	34.83 %	76.01 %
Device 3 (Multi-talker Condition)	44.79 %	78.95 %

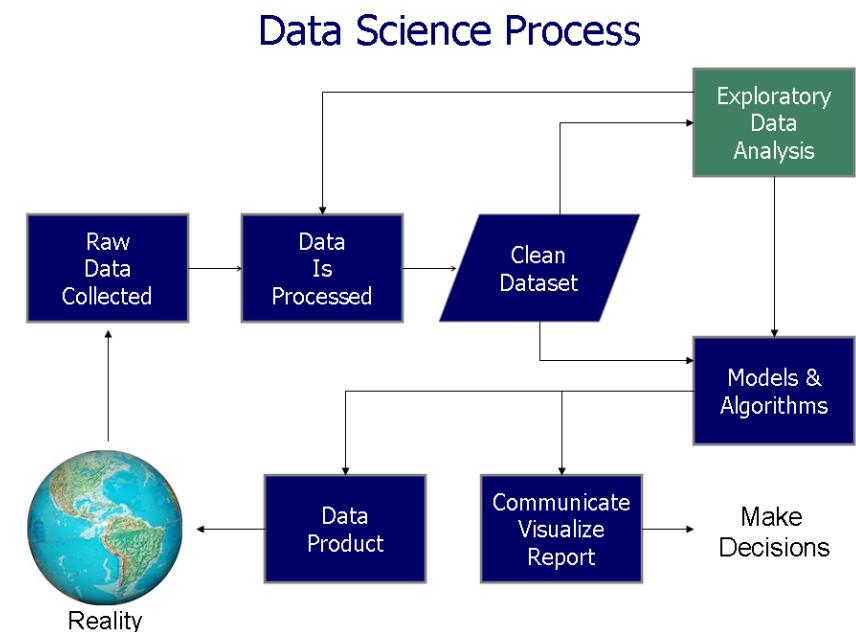
Source: <https://research.google.com/pubs/pub46130.html>

Why is big data analytics important?

- **Cost reduction:** Big data technologies such as Hadoop and cloud-based analytics bring significant cost advantages when it comes to storing large amounts of data – plus they can identify more efficient ways of doing business.
- **Faster, better decision making:** With the speed of Hadoop and in-memory analytics, combined with the ability to analyze new sources of data, businesses are able to analyze information immediately – and make decisions based on what they've learned.
- **New products and services:** With the ability to gauge customer needs and satisfaction through analytics comes the power to give customers what they want. Davenport points out that with big data analytics, more companies are creating new products to meet customers' needs.

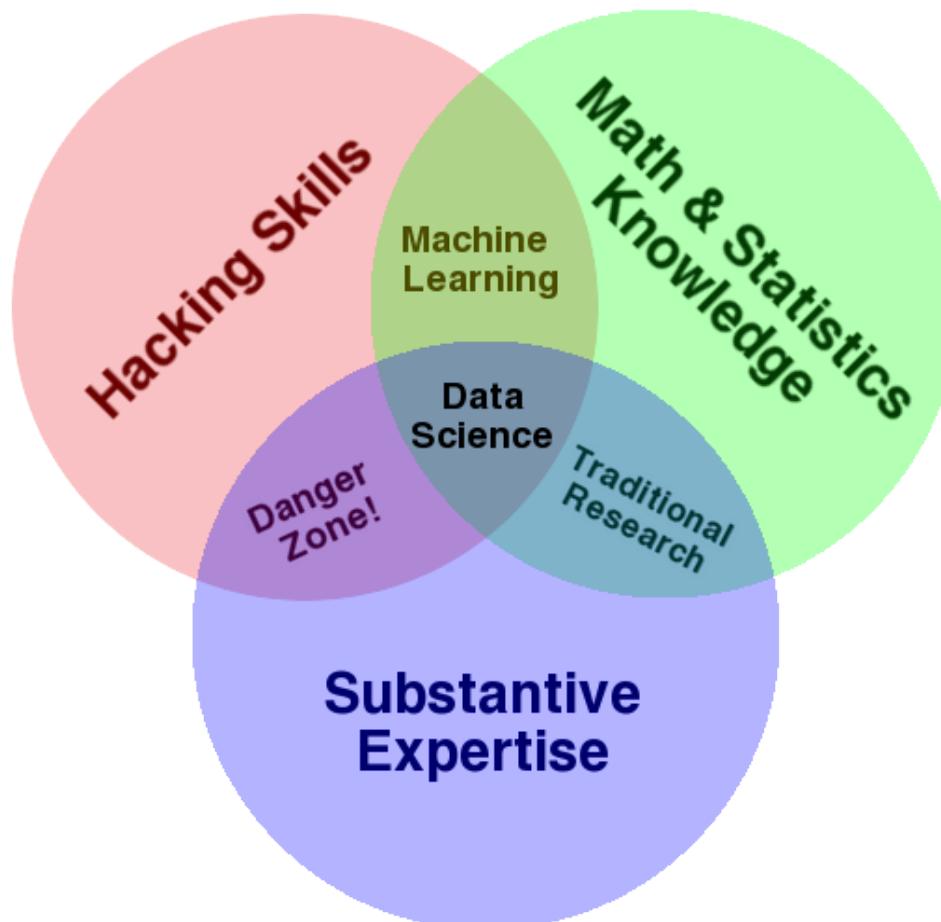
What is data science?

- Data science is an interdisciplinary field about scientific methods, processes, and systems to extract knowledge or insights from data in various forms, either structured or unstructured, similar to data mining for '**decision making**'



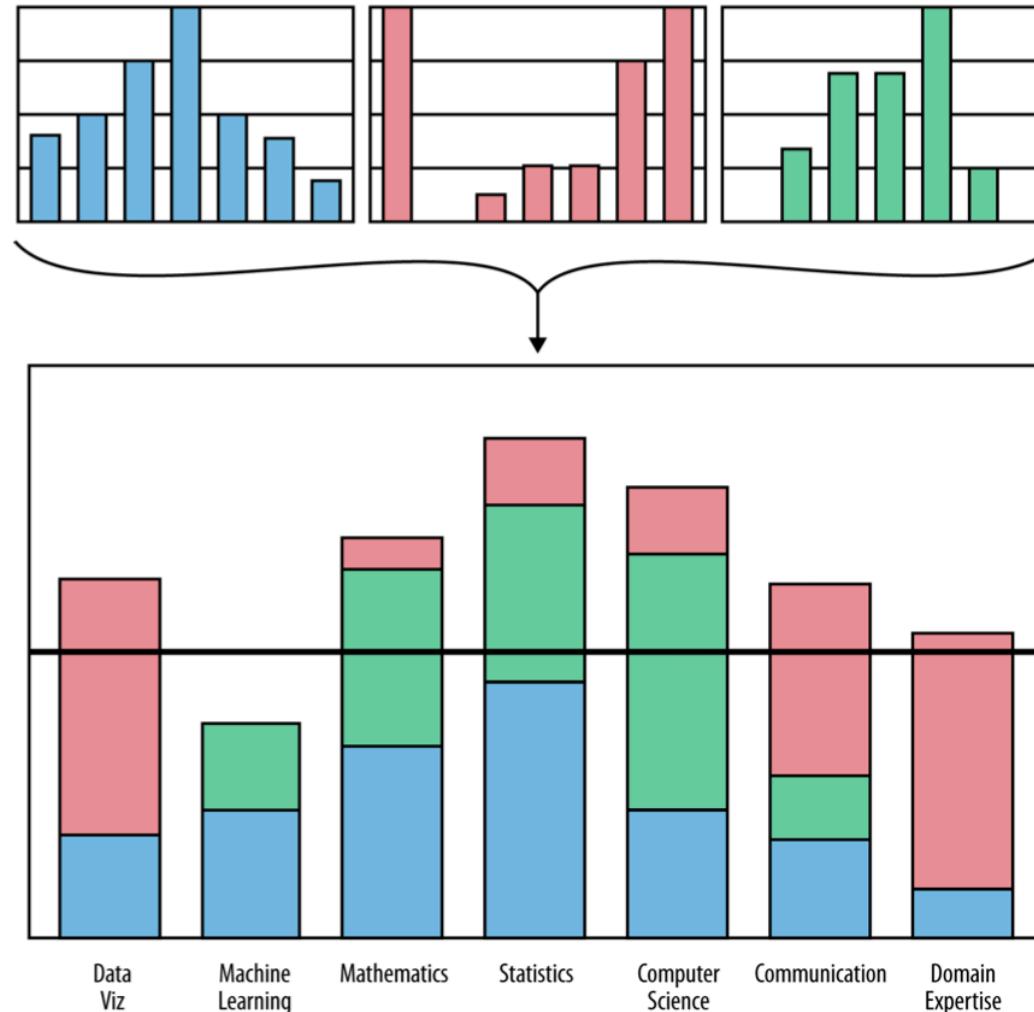
Source: https://en.wikipedia.org/wiki/Data_science

Another definition of data science



Source: <http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>

No one person can be the perfect data scientist, so we need teams.



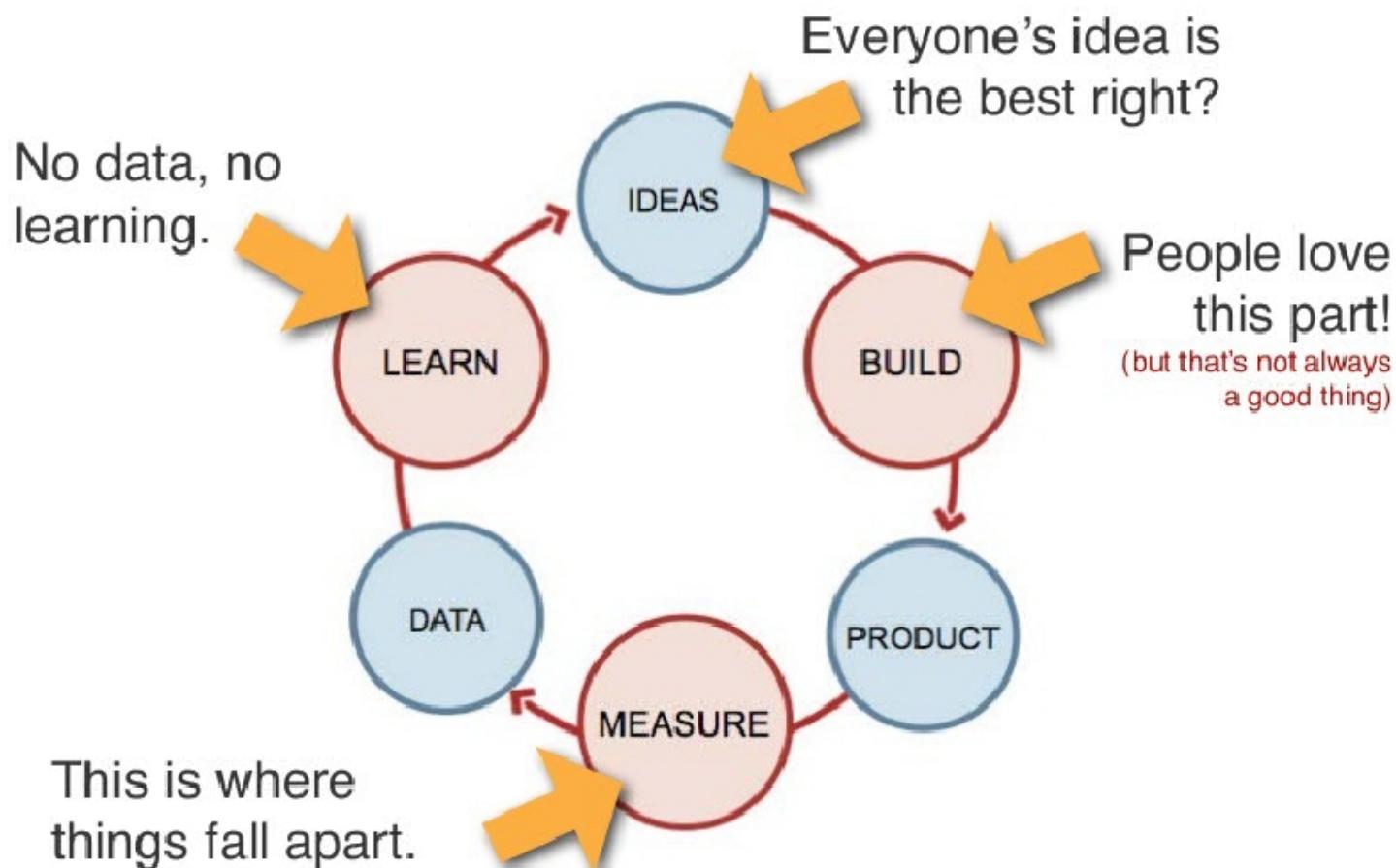
Source: "Doing Data Science"

Exploratory Data Analysis (EDA)

- Exploratory Data Analysis (EDA) is an approach/philosophy for data analysis that employs a variety of techniques (mostly graphical) to maximize insight into a data set;
 - Uncover underlying structure
 - Extract important variables
 - Detect outliers and anomalies
 - Test underlying assumptions
 - Develop parsimonious models
 - Determine optimal factor settings
- The EDA approach is precisely that – an approach – not a set of techniques, but an attitude/philosophy about how a data analysis should be carried out
- c.f., “Exploratory Data Analysis”, John W. Tukey (1977)

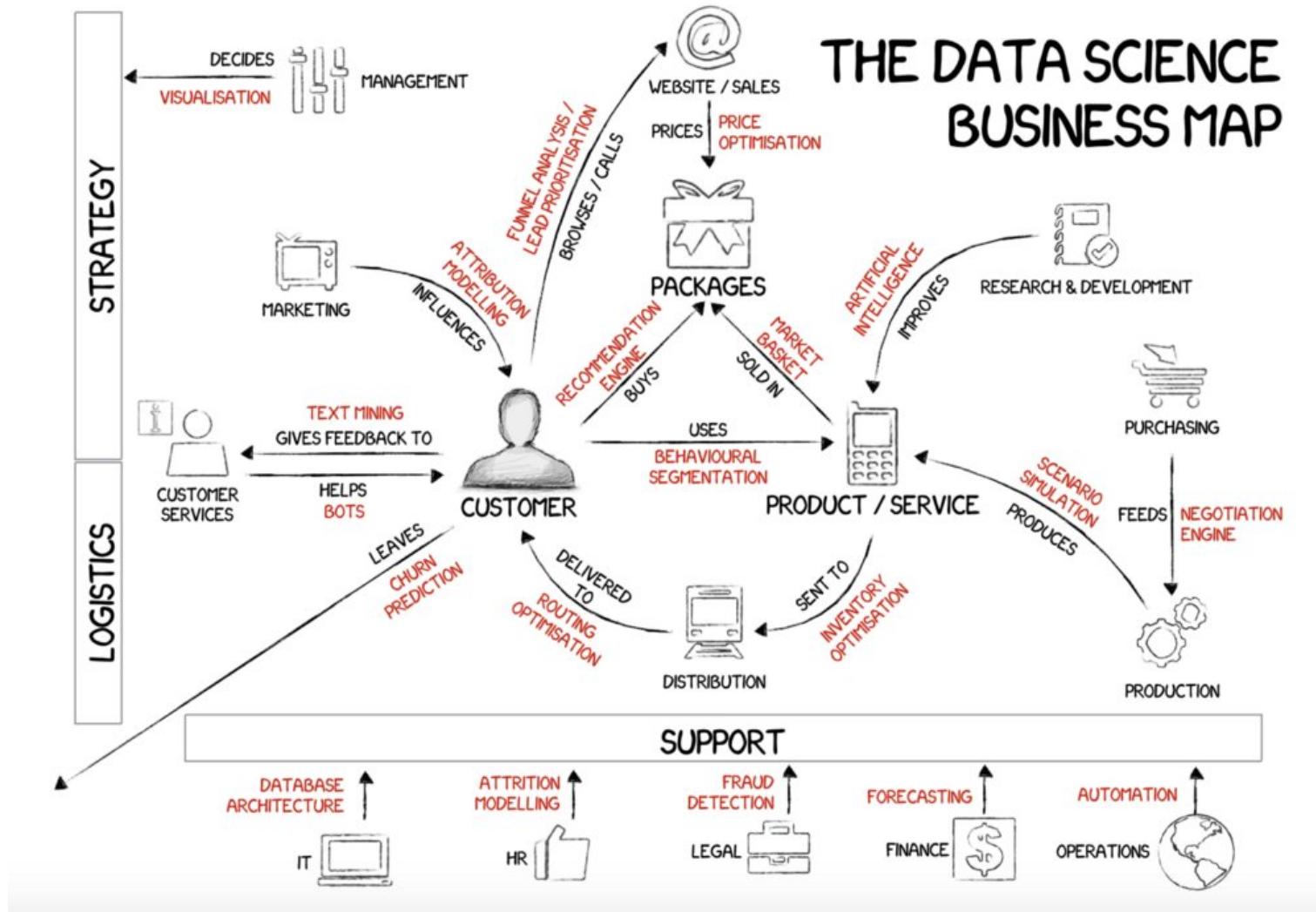
Source: <http://www.itl.nist.gov/div898/handbook/eda/section1/eda11.htm>

Lean Startup Cycle



Source: "Lean Analytics" by Alistair Croll

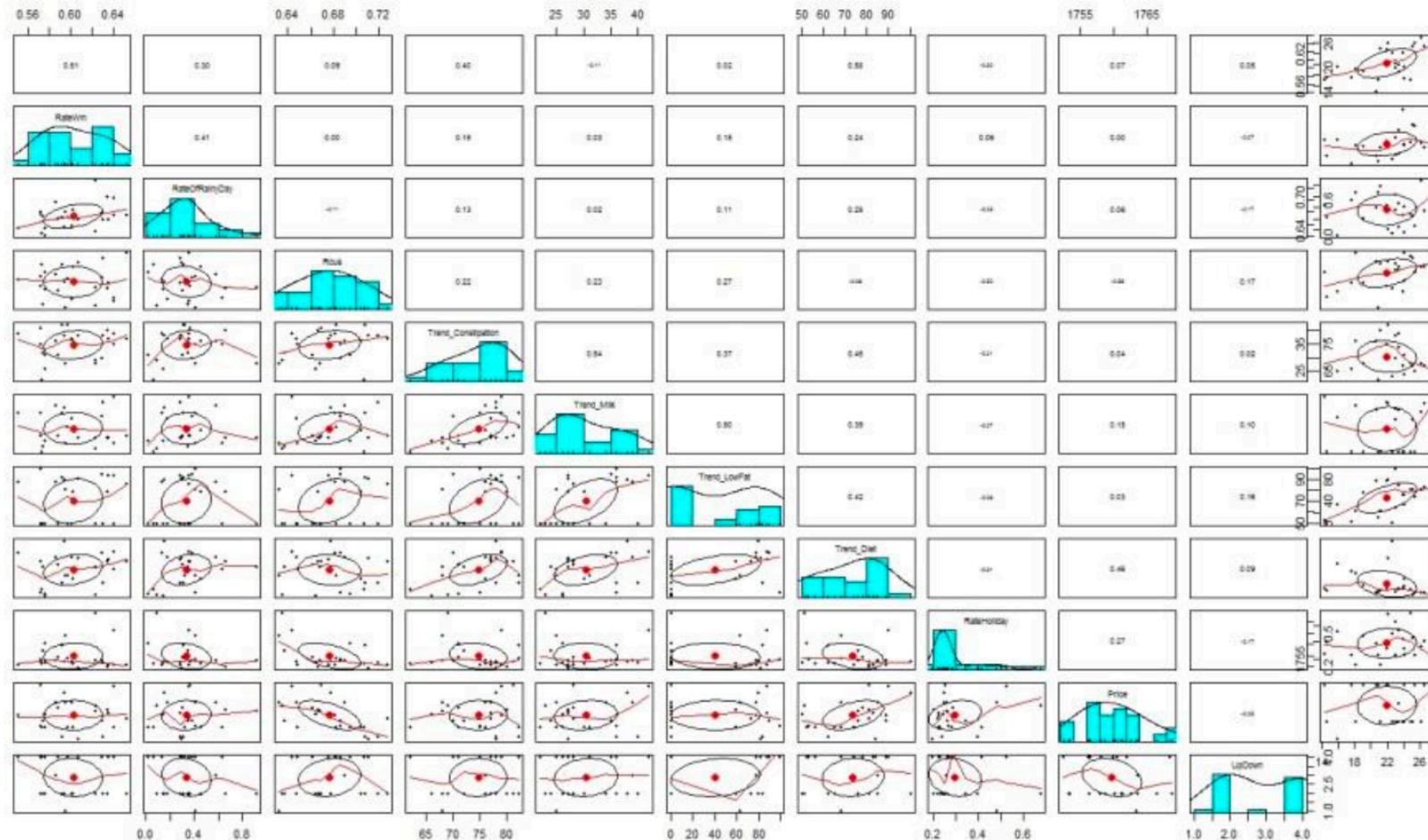
15



Data Visualization

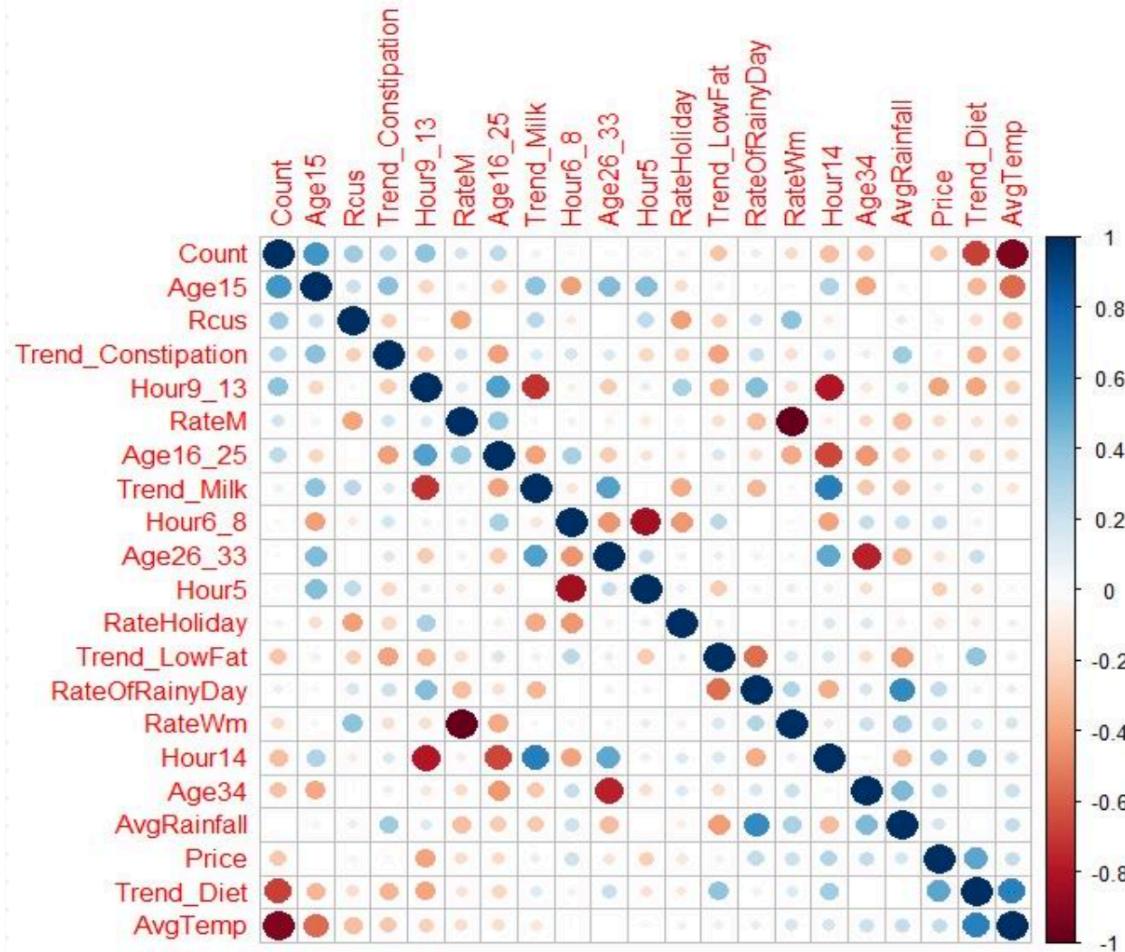
- Data visualization is the creation and study of the visual representation of data, meaning "information that has been abstracted in some schematic form, including attributes or variables for the units of information"
- Two objectives of data visualization
 - Visualization for analysis: figure out hidden relationship between several features of given data, typically during EDA process
 - Visualization for communication: deliver insights to the targeted audience efficiently and effectively without any distortion of the facts, possibly using infographics or interactive visualization

Example of data visualization during EDA



Source: Sales prediction of dairy products by Hyunlim Yang

Example of data visualization during EDA (cont.)



Source: Sales prediction of dairy products by Hyunlim Yang

Examples of Data Visualization

- 축구 패스 분포도: <http://news.joins.com/Digitalspecial/124>
- 선거의 득표 수를 어떻게 지도에 표현할 것인가?
 - <http://www.vw-lab.com/40>
- 2017년 19대 대통령 선거개표결과, 중앙일보
 - <http://news.joins.com/DigitalSpecial/179>
- Curated list of data visualization
 - <https://www.maptive.com/17-impressive-data-visualization-examples-need-see/>
 - <https://www.webdesignerdepot.com/2009/06/50-great-examples-of-data-visualization/>

Data Science v.s. Data Engineering

- **Data Engineering:** manage automatic and scalable data flow (close to CS)
 - **Extract:** The process of reading data from a database
 - **Transform:** the process of converting the extracted data from its previous form into the form it needs to be in so that it can be placed into another database. Transformation occurs by using rules or lookup tables or by combining the data with other data
 - **Load:** the process of writing the data into the target database
- **Data Science:** find insights from data (interdisciplinary domain)
 - **Discover:** Find, identify the sources of good data, and the metrics. Sometimes request the data to be created (work with data engineers and business analysts)
 - **Access:** Access the data. Sometimes via an API, a web crawler, an Internet download, a database access or sometimes in-memory within a database
 - **Distill:** Extract essence from data, the stuff that leads to decisions, increased ROI, and actions (such as determining optimum bid prices in an automated bidding system)

Assignment #0

- Due: 9/11 Monday
- Register to elice
 - Register dgist.elice.io
 - Use dgist email and Korean first/last name
 - Subscribe to “Introduction to Big Data Analysis and Visualization”
 - Subscription password: dgist_2017_ds
- Set up anaconda on your notebook

Reading List

- 데이터과학을 시작할 때 도움되는 것들
 - <https://brunch.co.kr/@amangkim/37>
 - <https://brunch.co.kr/@amangkim/38>
- 스타트업은 데이터를 어떻게 바라봐야 할까?
 - <https://www.slideshare.net/yongho/ss-32267675>
- Data Scientist versus Data Engineer
 - <http://www.datasciencecentral.com/profiles/blogs/data-scientist-versus-data-engineer>
- A very short history of data science
 - <https://www.forbes.com/sites/gilpress/2013/05/28/a-very-short-history-of-data-science/>



ANY QUESTIONS?