# Lecture 1: Data Analytics

Pilsung Kang

School of Industrial Management Engineering

Korea University

# AGENDA

# Data Analytics

- Data Analytics by Amazon

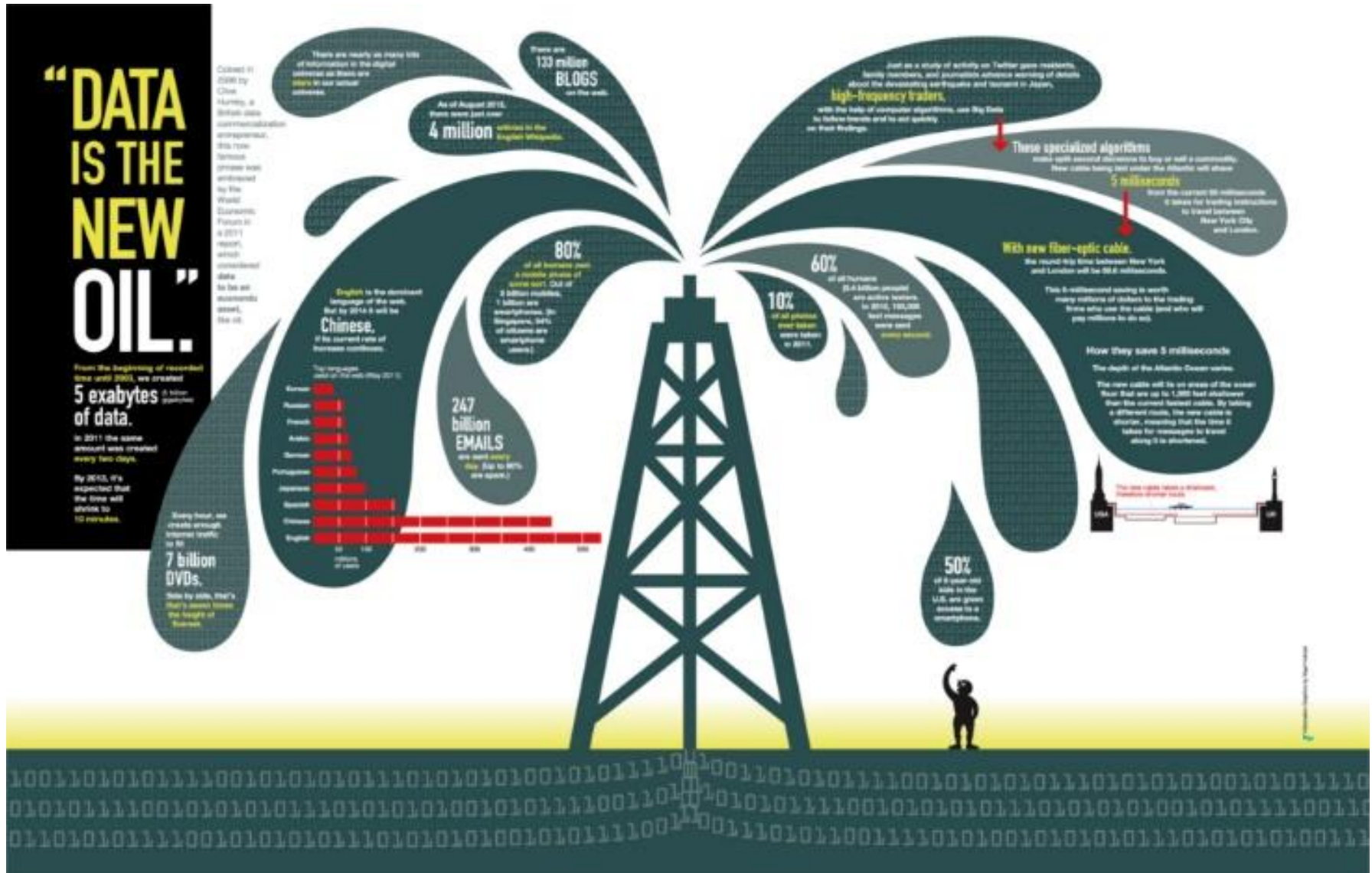# Data Analytics: The Era of "Big Data"

# Data Analytics: The Era of "Big Data"
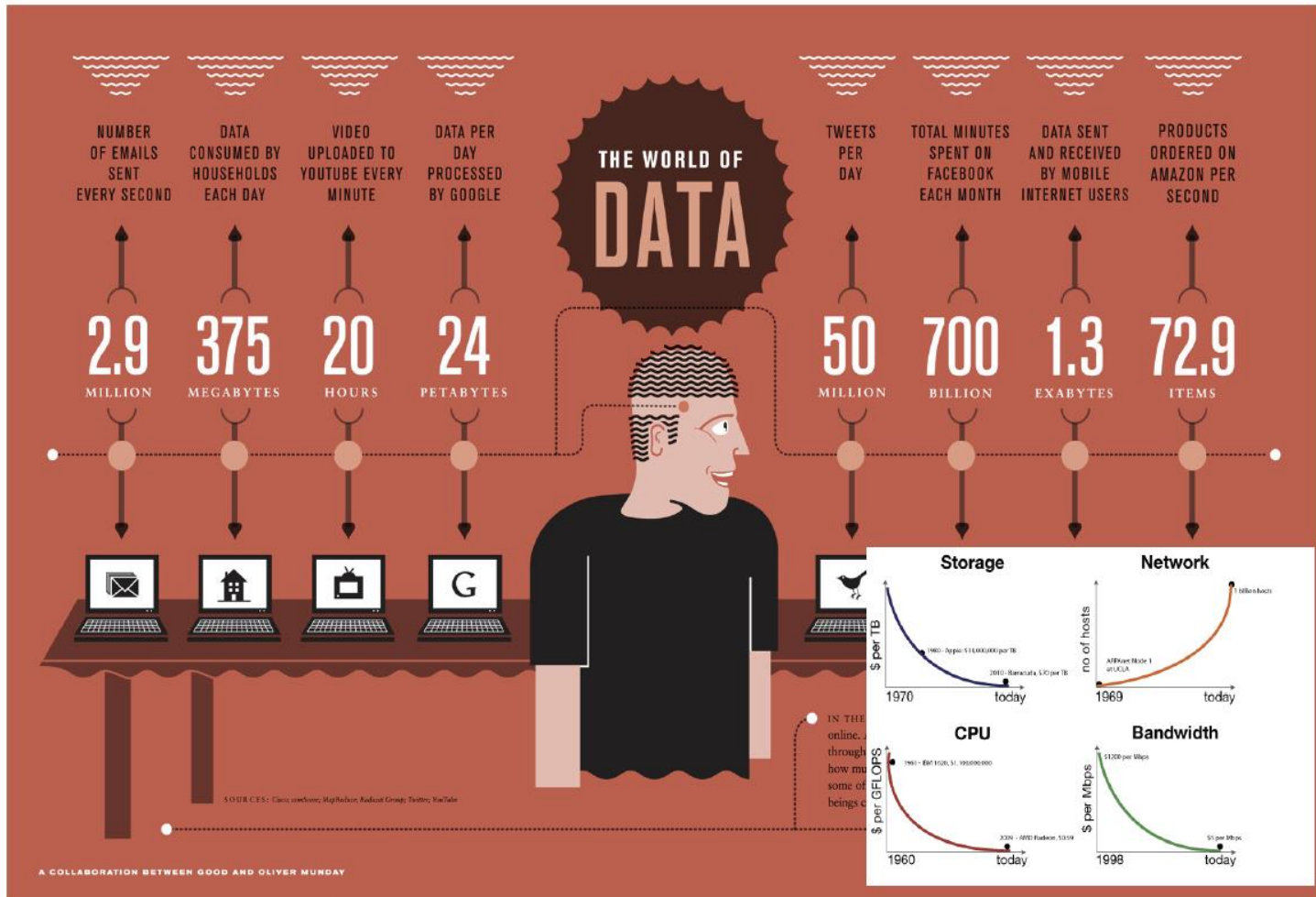
- 4Vs in Big Data



**Volume**



**Velocity**



**Variety**



**Value**

# Data Analytics: The Era of "Big Data"

• Volume of Big Data

# Data Analytics: The Era of "Big Data"

- Velocity of Big Data

# Data Analytics: The Era of "Big Data"

- Variety in Big Data

As of 2011, the global size of data in healthcare was estimated to be

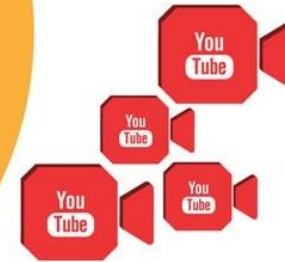## 150 EXABYTES
[ 161 BILLION GIGABYTES ]

By 2014, it's anticipated there will be

## 420 MILLION WEARABLE, WIRELESS HEALTH MONITORS

## 4 BILLION+ HOURS OF VIDEO

are watched on YouTube each month

# Variety
## DIFFERENT FORMS OF DATA

## 30 BILLION PIECES OF CONTENT

are shared on Facebook every month

## 400 MILLION TWEETS

are sent per day by about 200 million monthly active users

# Data Analytics: The Era of "Big Data"

- Value in Big Data

## Big data can generate significant financial value across sectors

**US health care**
- $300 billion value per year
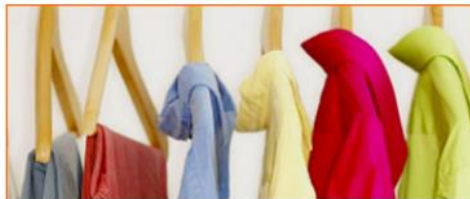- ~0.7 percent annual productivity growth

**Europe public sector administration**
- €250 billion value per year
- ~0.5 percent annual productivity growth

**Global personal location data**
- $100 billion+ revenue for service providers
- Up to $700 billion value to end users

**US retail**
- 60+% increase in net margin possible
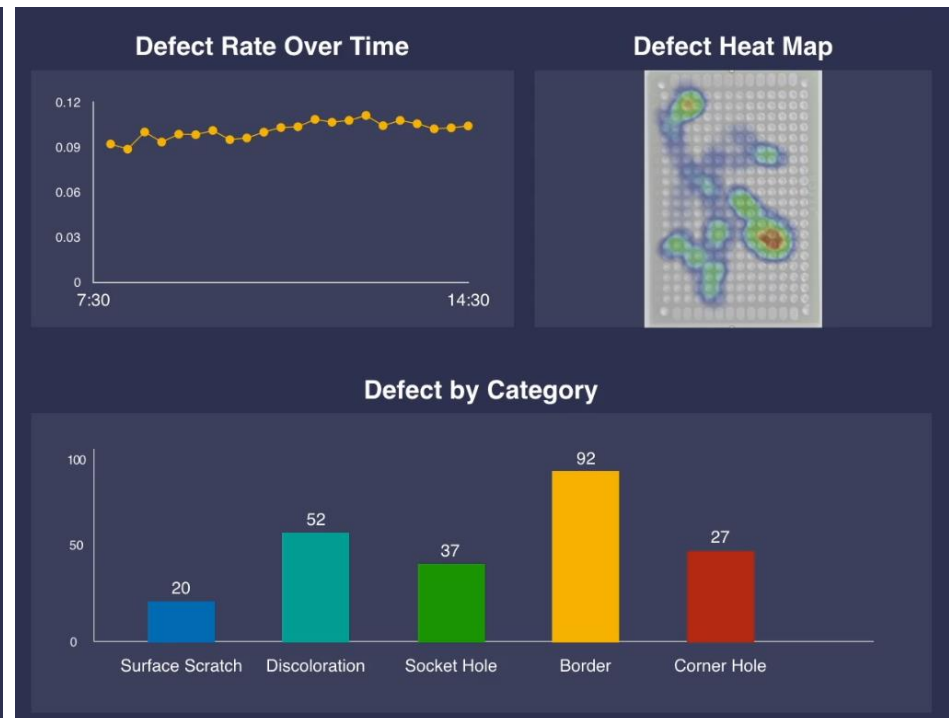- 0.5–1.0 percent annual productivity growth

**Manufacturing**
- Up to 50 percent decrease in product development, assembly costs
- Up to 7 percent reduction in working capital

# Data Analytics

- Data Analytics in Industrial Engineering (Manufacturing)

  ✓ Landing.ai: AI Startup found by Andrew Ng

  ✓ Provide various image/video analytics solutions for fault detection, leak defect detection, etc.

# Understanding Analytics

- Descriptive vs. Predictive vs. Prescriptive Analytics



**Understanding analytics**
Definitions, sample applications and opportunities, and underlying technologies

| | Descriptive | Predictive | Prescriptive |
|---|---|---|---|
| | What **HAS** happened? | What **COULD** happen? | What **SHOULD** happen? |
| **What the user needs to DO** | • **Increase** asset reliability<br>• **Reduce** labor and inventory costs | • **Predict** infrastructure failures<br>• **Forecast** facilities space demands | • **Increase** asset utilization<br>• **Optimize** resource schedules |
| **What the user needs to KNOW** | • The **number and types** of asset failures<br>• Why **maintenance costs** are high<br>• The value of the **materials inventory** | • How to **anticipate failures** for specific asset types<br>• When to **consolidate underutilized** facilities<br>• How to **determine costs** to improve service levels | • How to **increase** asset production<br>• Where to **optimally route** service technicians<br>• Which strategic facilities plan provides the **highest long-term utilization** |
| **How analytics gets ANSWERS** | • **Standard reporting** - What happened?<br>• **Query/drill down** - Where exactly is the problem?<br>• **Ad hoc reporting** - How many, how often, where? | • **Predictive modeling** - What will happen next?<br>• **Forecasting** - What if these trends continue?<br>• **Simulation** - What could happen?<br>• **Alerts** - What actions are needed? | • **Optimization** - What is the best possible outcome?<br>• **Random variable optimization** - What is the best outcome given the variability in specified areas? |
| **What makes this analysis POSSIBLE** | • Alerts, reports, dashboards, **business intelligence** | • Predictive **models**, forecasts, statistical **analysis**, scoring | • Business rules, organization **models**, comparisons, **optimization** |

Business value →

| 데이터분석을 위한 프로그래밍 언어<br>수리통계 및 실습<br>자료구조 및 알고리즘<br>응용통계 및 실습 | 데이터마이닝<br>다변량분석<br>예측애널리틱스<br>영상정보시스템 | OR-I 및 실습<br>OR-II 및 실습<br>최적화이론<br>최적화응용<br>메타휴리스틱 |
|---|---|---|

# Data Scientist

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

  ✓ Harvard Business Review

  ✓ https://hbr.org/2012/10/data-scientist-the-sexiest-job-of-the-21st-century

# Data Science Tasks

- Data Mining

  ✓ The process of exploration and analysis, by automatic or semi-automatic means, of large quantities of data in order to discover meaningful patterns and rules. (Berry and Linoff, 1997, 2000)

### Customers who viewed this item also viewed these products

| | | | |
|---|---|---|---|
| Dualit Food XL1500 Processor | Kenwood kMix Manual Espresso Machine | Weber One Touch Gold Premium Charcoal Grill-57cm | NoMU Salt Pepper and Spice Grinders |
| $560 | ★★★★☆ $250 | $225 | $3 |
| 🛒 Add to cart | 🔧 Select options | 🛒 Add to cart | 🛒 View options |

# Data Science Tasks

- Machine Learning

  ✓ A computer program is said to learn from <u>experience</u> **E** with respect to some class
  of <u>tasks</u> **T** and <u>performance measure</u> **P**, if its performance at task in T, as measured by
  P, improves with experience E," – Mitchell et al. (2013)

```
┌──────────┐      ┌──────────────┐      ┌──────────┐
│   Data   │ ───▶ │ Methodology  │ ───▶ │ Results  │
│          │      │   (Model)    │      │          │
└──────────┘      └──────────────┘      └──────────┘
```

고려대학교
KOREA UNIVERSITY

DSBA
Data Science & Business Analytics

# Data Science Tasks

- Machine Learning Models in Self-Driving Cars

# Data Science Tasks

- Artificial Intelligence

  ✓ Computers and computer software that are capable of intelligent behavior

  ✓ Intelligent agent perceives its environment and takes actions that maximize its chance of success

# AGENDA

DSBA
Data Science & Business Analytics

# Process Monitoring & Control

| Input Data | Algorithms | Output |
|---|---|---|
| Sensor Data from Production Equipment | Regression + Novelty Detection | Metrological Values With Reliability Scores |



Process equipment

Processed wafers

Sampled wafers

Metrology equipment

Metrology result

UCL

LCL

Metrology data

Prediction Model

Process sensor data

Novelty Detection Model

Outlier    Normal    Outlier

Virtual metrology result

UCL

LCL

◆ : Reliable VM result

◆ : Unreliable VM result

고려대학교
KOREA UNIVERSITY

DSBA
Data Science & Business Analytics

# Process Monitoring & Control

| Input Data | Algorithms | Output |
|---|---|---|
| Sensor Data from Production Equipment | Regression + Novelty Detection | Metrological Values With Reliability Scores |

# Are You a Valid User?

| Input Data | Algorithms | Output |
|---|---|---|
| Time stamps collected during typing | Novelty Detection | Valid user score |



**Through various input devices**

**In any stages**

**Customized (real-time) authenticator**

# Who Looks Happy?

| Input Data | Algorithms | Output |
|---|---|---|
| Image | Convolutional Neural Networks | Emotions Recognized |

https://www.microsoft.com/cognitive-services/en-us/emotion-api

# Favorite Artistic Style

| Input Data | Algorithms | Output |
|---|---|---|
| Image | Convolutional Neural Networks | Image with Preferred Style |

Gatys, L. A., Ecker, A. S., & Bethge, M. (2015). A neural algorithm of artistic style. arXiv preprint arXiv:1508.06576.

고려대학교 KOREA UNIVERSITY

DSBA
Data Science & Business Analytics

# Understanding Emotions

| Input Data | Algorithms | Output |
|---|---|---|
| Movie Review Text | Convolutional Neural Networks | Sentiment (Pos/Neg) & Keyword attention |

# Understanding Emotions

| Input Data | Algorithms | Output |
|---|---|---|
| Movie Review Text | Convolutional Neural Networks | Sentiment (Pos/Neg) & Keyword attention |

| Method | Sentence |
|---|---|
| Raw text | One of the funniest most romantic and most musical movies ever; definitely worth renting/buying especially if you have a taste for older style of cinematography. The animals and the songs alone will make you smile while watching the movie. A definite must for Madonna fans. :o) *(10 / 10 points)* |
| Rand | One of the **the funniest most romantic** and musical movies ever definitely worth renting buying especially if you have a taste for older style cinematography The animals songs alone will make smile while watching movie A definite must Madonna **fans** **Positive** |
| Static | One of the **funniest most romantic and** musical movies ever definitely worth renting buying especially if you have a taste for older style cinematography The animals songs alone will make smile while watching movie **A** definite must Madonna fans **Positive** |
| NStatic | One **of the funniest most** romantic and musical movies ever definitely worth renting buying especially if you have a taste for older style cinematography The animals songs alone will make smile while watching movie A definite must Madonna **fans** **Positive** |
| 2ch | One **of the funniest most** romantic and musical movies ever definitely worth renting buying especially if you have a taste for older style cinematography The animals songs alone will make smile while watching movie A definite must Madonna **fans** **Positive** |

고려대학교
KOREA UNIVERSITY

DSBA
Data Science & Business Analytics

# Understanding Emotions

| Input Data | Algorithms | Output |
|---|---|---|
| Movie Review Text | Convolutional Neural Networks | Sentiment (Pos/Neg) & Keyword attention |

| Method | Sentence |
|---|---|
| **Raw text** | This is one of the most boring films I've ever seen. The three main cast members just didn't seem to click well. Giovanni Ribisi's character was quite annoying. For some reason he seems to like repeating what he says. If he was the Rain Man it would've been fine but he's not. *(3 / 10 points)* |
| **Rand** | This is one of the most boring films I ve ever seen The three main cast members just didn t seem to click w ell Giovanni Ribisi s character was quite annoing For some reason he seems like repeating what says If Rain Man it would been fine but he s not Negative |
| **Static** | This is one of the most boring films I ve ever seen The three main cast members just didn t seem to click w ell Giovanni Ribisi s character was quite annoying For some reason he seems like repeating what says If Rai n Man it would been fine but he s not Negative |
| **NStatic** | This is one of the most boring films I ve ever seen The three main cast members just didn t seem to click w ell Giovanni Ribisi s character was quite annoying For some reason he seems like repeating what says If Rai n Man it wouldbeen fine but not Negative |
| **2ch** | This is one of the most boring films I ve ever seen The three main cast members just didn t seem to click w ell Giovanni Ribisi s character was quite annoying For some reason he seems like repeating what says If Rai n Man it would beenfine but not Negative |
| **4ch** | This is one of the most boring films I ve ever seen The three main cast members just didn t seem to click w ell Giovanni Ribisi s character was quite annoying For some reason he seems like repeating what says If Rai n Man it would been fine but not Negative |

고려대학교 KOREA UNIVERSITY

DSBA
Data Science & Business Analytics

# What are the Main Topics in Anonymous Posts of University Students?

| Input Data | Algorithms | Output |
|---|---|---|
| Facebook Posts (Bamboo forest) | Latent Dirichlet Allocation (Topic Modeling) | Topic Distribution |



[1] 28개 대학 페이스북 대나무숲 주제 분포

총 36만개 DATAx

| | | |
|---|---|---|
| 0 | TOPIC | 사회편견 / 학벌 |
| 1 | TOPIC | 시험 / 학업 / 성적 |
| 2 | TOPIC | 연애 / 그리움 / 이별 |
| 3 | TOPIC | 먹는 것 / 학교주변 생활 |
| 4 | TOPIC | 재정 / 용돈 / 알바 |
| 5 | TOPIC | 제보 / 부탁 / 분실물 |
| 6 | TOPIC | 동아리 / 공연 / 축제 |
| 7 | TOPIC | 썸 / 이상형 |
| 8 | TOPIC | 회상 / 그리움 / 감성글 |
| 9 | TOPIC | 대인관계 / 가정환경 |
| 10 | TOPIC | 학교생활 / 선후배 / 입학 / 동기 |
| 11 | TOPIC | 의견 / 칼럼글 |

Domain Topic: Topic2&Topic9
대학생 : 관계에 대한 고민

고려대학교 KOREA UNIVERSITY

DSBA Data Science & Business Analytics

# What are the Main Topics in Anonymous Posts of University Students?
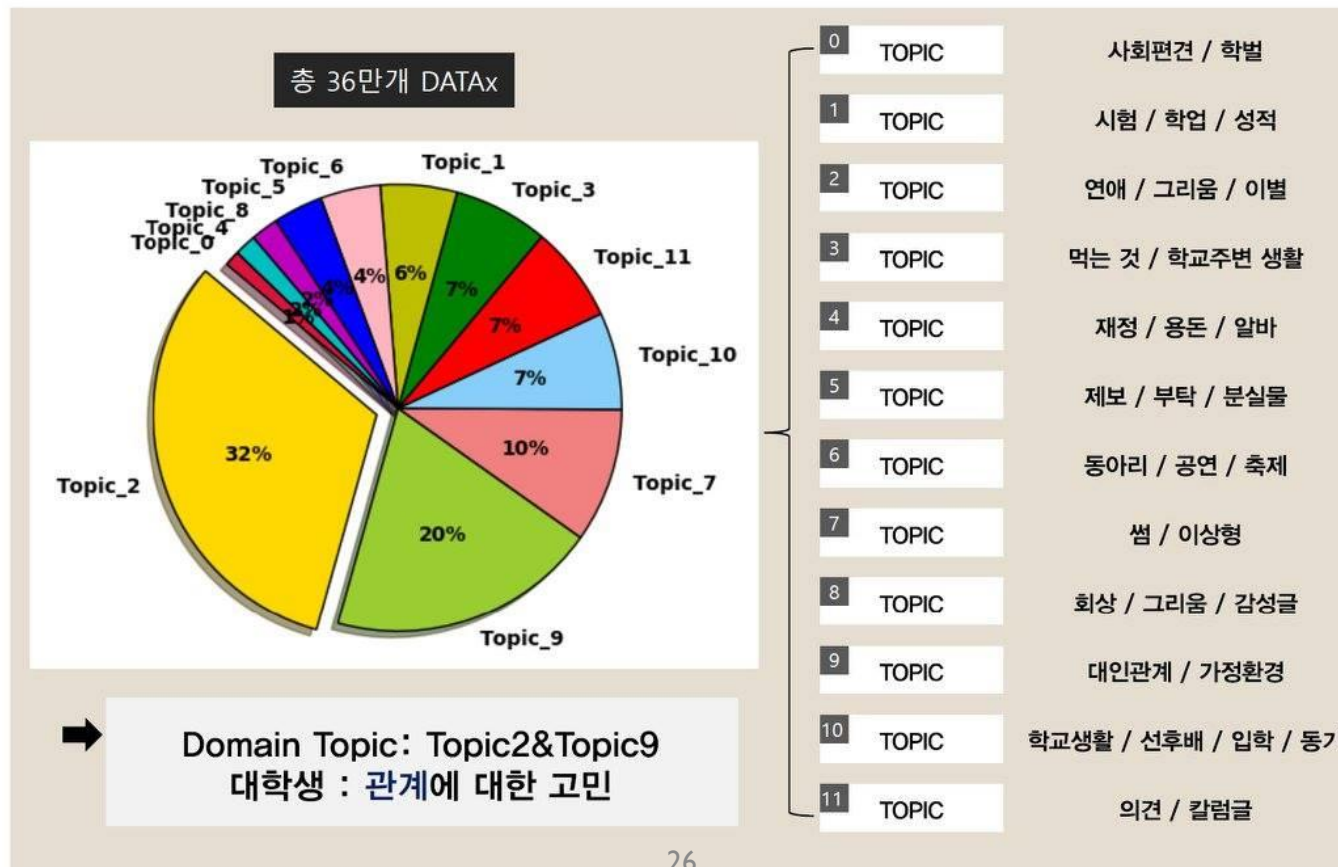
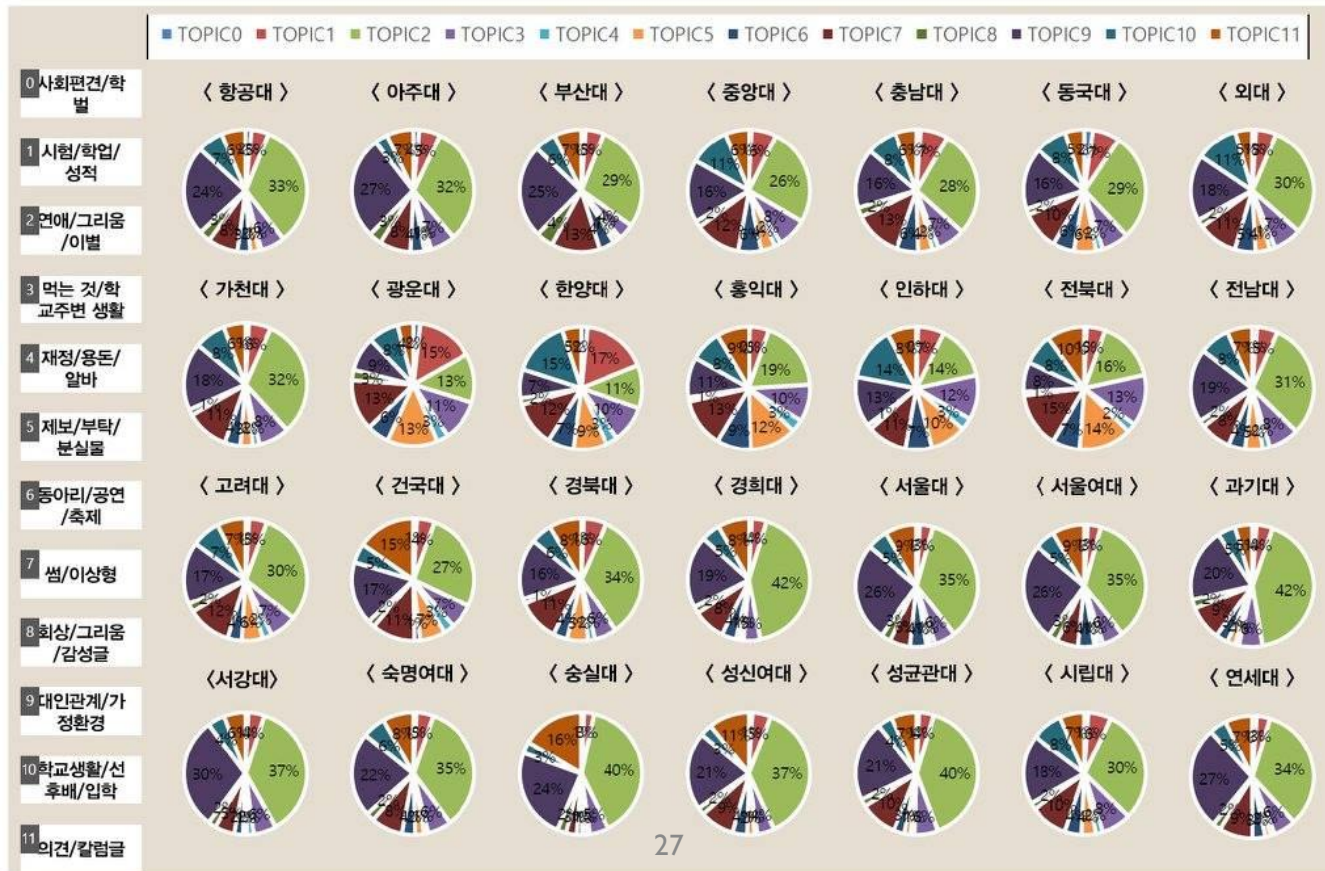| Input Data | Algorithms | Output |
|---|---|---|
| Facebook Posts (Bamboo forest) | Latent Dirichlet Allocation (Topic Modeling) | Topic Distribution |



[2] 대학별 페이스북 대나무숲 주제 분포

# What are the Main Topics in Anonymous Posts of University Students?

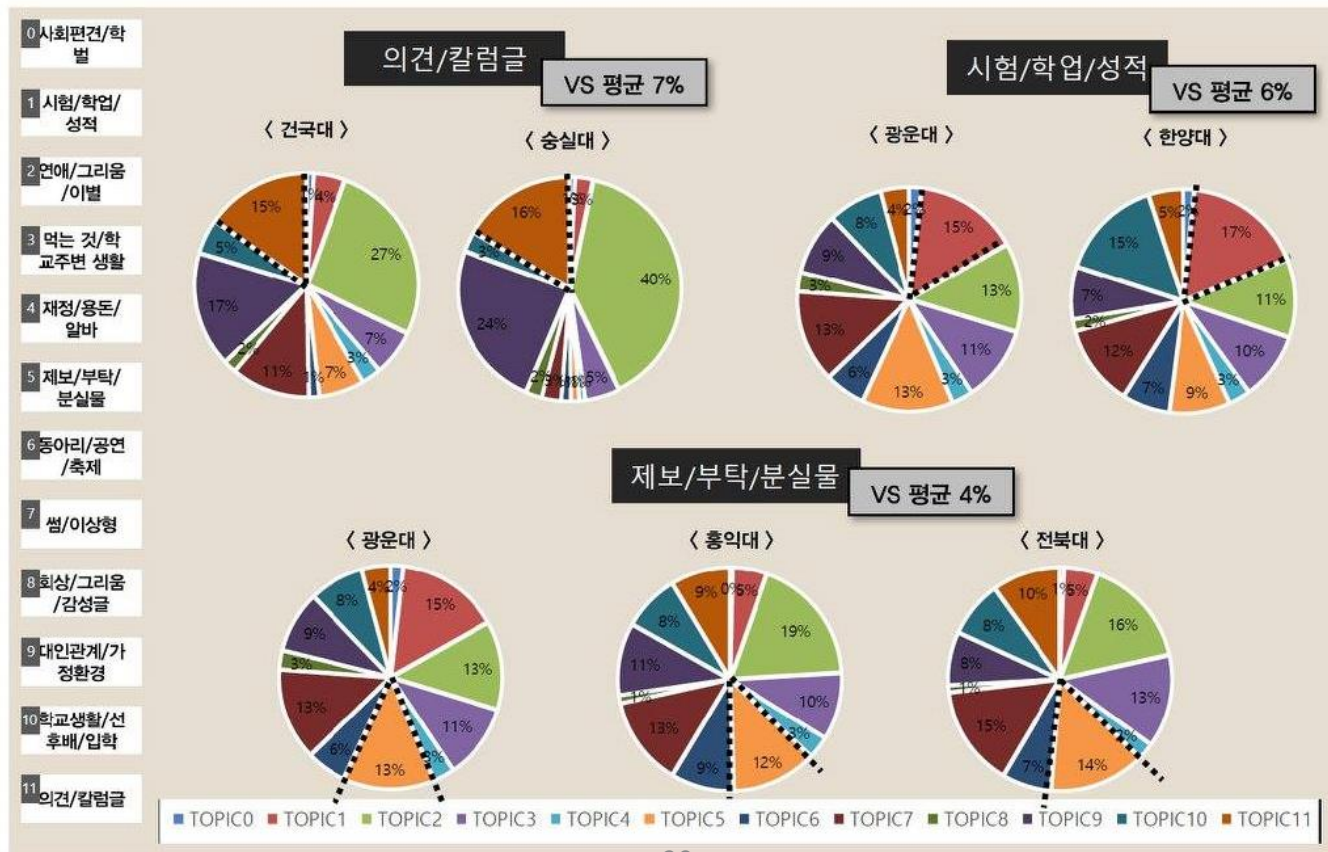| Input Data | Algorithms | Output |
|---|---|---|
| Facebook Posts (Bamboo forest) | Latent Dirichlet Allocation (Topic Modeling) | Topic Distribution |



[3] 대학별 특이 사항

# Connecting the Dots

# Connecting the Dots

You can't connect the dots looking forward; you can only connect them looking backwards.

So you have to trust that the dots will somehow connect in your future.

# AGENDA

**DSBA**
Data Science & Business Analytics

# R vs. Python

- Learn Data Science, not Programming

  - ✓ R vs. Python, different brushes

  - ✓ Do not choose between R & Python, learn both

  - ✓ It strengthens your data science communication skills

  - ✓ It boosts your data science career

  - ✓ It is not that hard

# R vs. Python

- Some interesting polls (2018)

## Table 1: Top Analytics/Data Science/ML Software in 2018 KDnuggets Poll

| Software | 2018 % share | % change 2018 vs 2017 |
|---|---|---|
| Python | 65.6% | 11% |
| RapidMiner | 52.7% | 65% |
| R | 48.5% | -14% |
| SQL | 39.6% | 1% |
| Excel | 39.1% | 24% |
| Anaconda | 33.4% | 37% |
| Tensorflow | 29.9% | 32% |
| Tableau | 26.4% | 21% |
| scikit-learn | 24.4% | 11% |
| Keras | 22.2% | 108% |



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

# R vs. Python
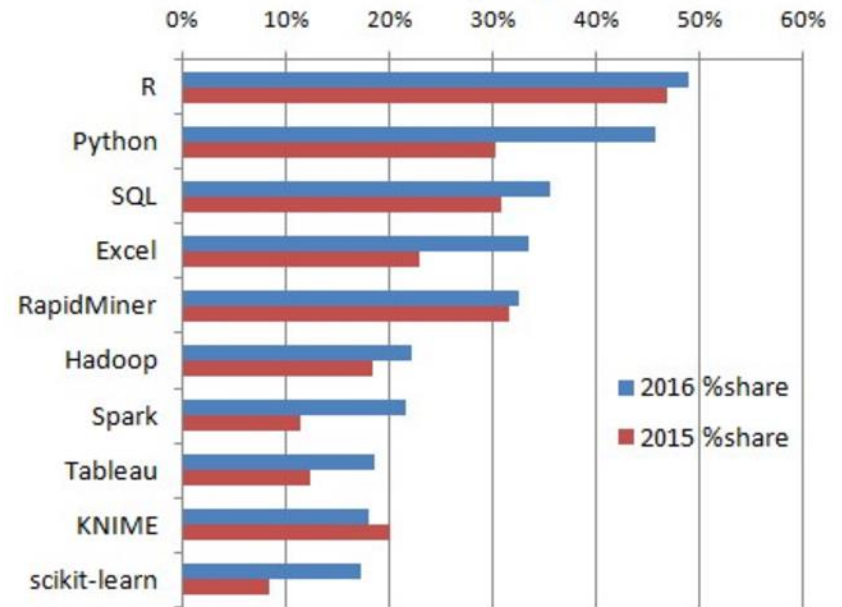
- Some interesting polls (2016)

Next table has the top 10 most popular tools in 2016 poll

| Tool | 2016 % share | % change | % alone |
|------|-------------|----------|---------|
| R | 49% | +4.5% | 1.4% |
| Python | 45.8% | +51% | 0.1% |
| SQL | 35.5% | +15% | 0% |
| Excel | 33.6% | +47% | 0.2% |
| RapidMiner | 32.6% | +3.5% | 11.7% |
| Hadoop | 22.1% | +20% | 0% |
| Spark | 21.6% | +91% | 0.2% |
| Tableau | 18.5% | +49% | 0.2% |
| KNIME | 18.0% | -10% | 4.4% |
| scikit-learn | 17.2% | +107% | 0% |

**KDnuggets Analytics/Data Science 2016 Software Poll, top 10 tools**

# R vs. Python

- Some interesting polls

# R vs. Python

- Difference between R and Python

## Difference between R and Python

| Parameter | R | Python |
|---|---|---|
| Objective | Data analysis and statistics | Deployment and production |
| Primary Users | Scholar and R&D | Programmers and developers |
| Flexibility | Easy to use available library | Easy to construct new models from scratch. I.e., matrix computation and optimization |
| Learning curve | Difficult at the beginning | Linear and smooth |
| Popularity of Programming Language. Percentage change | 4.23% in 2018 | 21.69% in 2018 |
| Average Salary | $99.000 | $100.000 |
| Integration | Run locally | Well-integrated with app |
| Task | Easy to get primary results | Good to deploy algorithm |
| Database size | Handle huge size | Handle huge size |
| IDE | Rstudio | Spyder, Ipthon Notebook |
| Important Packages and library | tydiverse, ggplot2, caret, zoo | pandas, scipy, scikit-learn, TensorFlow, caret |
| Disadvantages | Slow High Learning curve Dependencies between library | Not as many libraries as R |
| Advantages | • Graphs are made to talk. R makes it beautiful<br>• Large catalog for data analysis<br>• GitHub interface<br>• RMarkdown<br>• Shiny | • Jupyter notebook: Notebooks help to share data with colleagues<br>• Mathematical computation<br>• Deployment<br>• Code Readability<br>• Speed<br>• Function in Python |

고려대학교
KOREA UNIVERSITY

DSBA
Data Science & Business Analytics

# R vs. Python

- Data Science Wars: R vs. Python

  ✓ https://101.datascience.community/2015/05/12/data-science-wars-r-vs-python/