

# 데이터 분석



소프트웨어융합대학원  
진혜진

# 1. 데이터의 시대

## 1. 데이터의 비즈니스 활용

- 우리는 데이터의 시대(the age of data)에 살고 있음, 정보화 시대 → 데이터의 시대
- 우리를 둘러싼 모든 것들이 데이터 소스와 연결되고, 우리 삶의 많은 부분이 데이터에 의존하여 영위 ex) 이메일, SNS, 전화사용 기록, 신용카드거래 기록, 병원 치료 기록, 성적, 인터넷, 주민정보, 등기정보, 판매정보, 주식거래 정보 등
- 데이터는 기업 활동에도 중요함, 대형마트들은 소비자의 구매 내역 데이터를 바탕으로 구매 패턴을 분석하고 이를 영업에 활용



맥주를 산 고객이 견과류도 함께 구매하는 비율이 높다고 분석되면



맥주 바로 옆에 견과류를 진열



동반 매출 상승

그림 1-1 판매유통 대형마트 진열대: 구매 패턴 데이터를 분석하여 활용

# 1. 데이터의 시대

- GE 에비에이션은 비행기 엔진에 수많은 센서를 부착하고, 이 센서로부터 수집된 데이터를 활용하여 엔진의 이상 유무나 부품 교체 시기 등을 알려주는 서비스를 제공하여 추가 매출을 올리고 있음

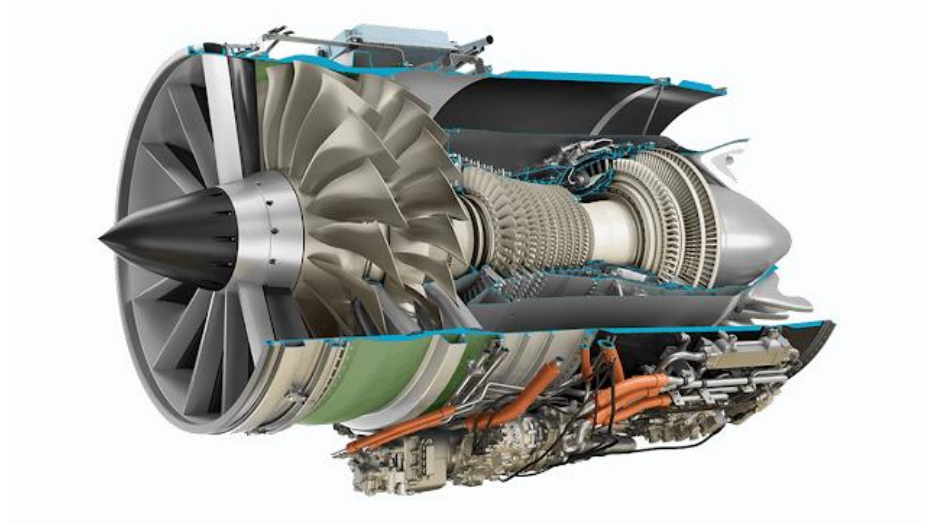


그림 1-2 GE 에비에이션: 비행기 엔진 센서 데이터를 분석하여 활용



# 1. 데이터의 시대

## 2. 4차 산업혁명과 데이터

- 2016년 1월, 스위스 다보스에서 열렸던 세계경제포럼(World Economic Forum)에서 클라우스 슈밥(Klaus Schwab)은 기술 혁명의 새로운 시대가 열렸음을 천명하면서 이를 '**4차 산업혁명(The Fourth Industrial Revolution)**'이라고 명명
- 4차 산업혁명이란 **인공지능**(Artificial Intelligence, AI), **빅데이터**(big data), **로봇**(robot), **사물인터넷**(Internet of Things, IoT), **생명공학기술**(Biotechnology), **3D 프린터**(3D printer) 등 새로운 과학기술이 사회, 경제, 문화 전반에 영향을 미치게 되고, 이러한 변화를 잘 수용하고 가능성을 최대화 하는 시대를 말함
- 인공지능**과 **빅데이터**가 4차 산업혁명의 핵심 기술로 인식

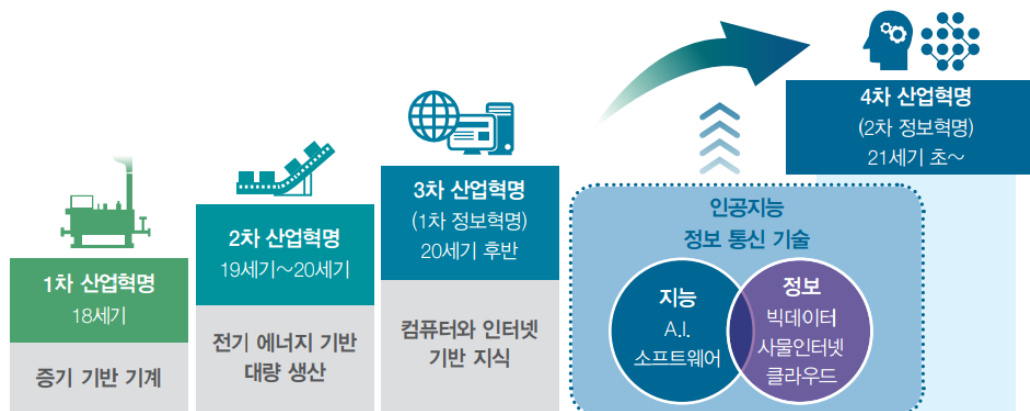


그림 1-4 4차 산업혁명까지의 과정과 핵심 기술

## 전문가들이 표현한 데이터의 중요성

1. 데이터는 비즈니스를 위한 새로운 원천 재료가 되어가고 있다."

- MS 부회장, 크레이그 먼디(Craig Mundie)

2. "데이터가 쏟아지는 수도꼭지가 틀어졌고, 다시 잠기는 일은 없을 것이다."

- 액티언 CTO, 마이크 호스킨(Mike Hoskins, Actian)

3. "데이터는 새로운 석유다."

- 데이터 과학자, 클라이브 험비(Clive Humby)

4. "당신이 정보를 포함하지 않은 데이터를 가질 수는 있겠지만, 데이터에 의하지 않은 정보는 가질 수 없다."

- 프로그래머/데이터 과학자, 다니엘 키즈 모런 (Daniel Keys Moran)

## 2. 빅데이터

### 2. 빅데이터의 성공 사례

#### 2.1 국내 활용 사례: 아파트 관리비 적정성 평가

- 경기도는 국토교통부의 공동주택관리정보시스템에 의무적으로 등록하는 각 아파트 관리사무소의 관리비 내역과 관리비를 구성하는 37개 세부항목의 원천데이터를 비교·분석하는 방식으로 관리비 과다 청구 여부를 분석
- 분석 결과를 가지고 아파트 관리비 산출 표준 모델 및 **아파트관리비부당지수** 개발
- 556개 단지를 샘플로 조사하여 2년간 152억원의 관리비가 부당하게 징수된 사실이 적발 전국적으로 적용될 경우 **연간 1조 1000억원 정도의 관리비를 절감**할 수 있을 것으로 예상

## 2. 빅데이터

### 2. 빅데이터의 성공 사례

#### 2.2 해외 활용 사례: 타깃의 맞춤형 광고



STEP 1

여학생에게 온  
타깃 광고  
메일 내용이  
유아용품?



STEP 2

빅데이터 전문가들이  
여학생 고객의 구매 분석  
기본 로션 →  
무향 로션 구매  
영양제 비구매 →  
미네랄 영양제 구매



STEP 3

타깃은 고객  
데이터베이스에 적용 →  
전국적으로 수만 명의  
임신 추정 여성들을  
가려내 관련 할인 쿠폰 발송



### 3. 데이터 분석 과정



그림 1-7 데이터 분석의 과정

#### 1. 1단계: 문제 정의 및 계획

- 문제가 명확해야 그 문제를 해결하기 위한 데이터가 어떤 것인지를 추정할 수 있고, 어떤 분석기법을 적용해야 할지도 계획할 수 있음

#### 2. 2단계: 데이터 수집

- 기존 시스템의 데이터베이스, 엑셀파일, 종이 문서, 장비내의 파일, 인터넷 등에서 필요한 자료를 수집

## 3. 데이터 분석 과정

### 3. 3단계: 데이터 정제 및 전처리

- 수집된 데이터는 바로 분석에 사용할수 없는 경우가 대부분
- 단위의 차이, 결측값, 오류 데이터 등의 보정 필요
- 수집된 데이터를 분석이 가능한 형태로 정돈하는 과정을 데이터 정제 혹은 전처리 과정

### 4. 4단계: 데이터 탐색

- 가벼운 데이터 분석
- 전반적인 데이터의 내용을 파악하는 단계

### 5. 5단계: 데이터 분석

- 데이터 탐색 단계에서 파악한 정보를 바탕으로 보다 심화된 분석을 수행하는 단계
- 전통적인 통계분석을 포함하여 고급 분석 기법들이 사용됨
- 머신러닝 기술도 적용됨

### 3. 데이터 분석 과정

#### 6. 6단계: 결과 보고

- 데이터의 분석과 해석이 마무리 되면 그 내용이 정리되고, 보고 되어야 함
- 결과보고 작성단계에서 중요한 기술이 바로 데이터 시각화(visualization)
- 데이터 시각화란 분석된 결과를 단순 숫자의 나열이 아니라 다양한 그래프나 그림을 통해서 결과를 쉽게 이해할 수 있도록 표현하는 것

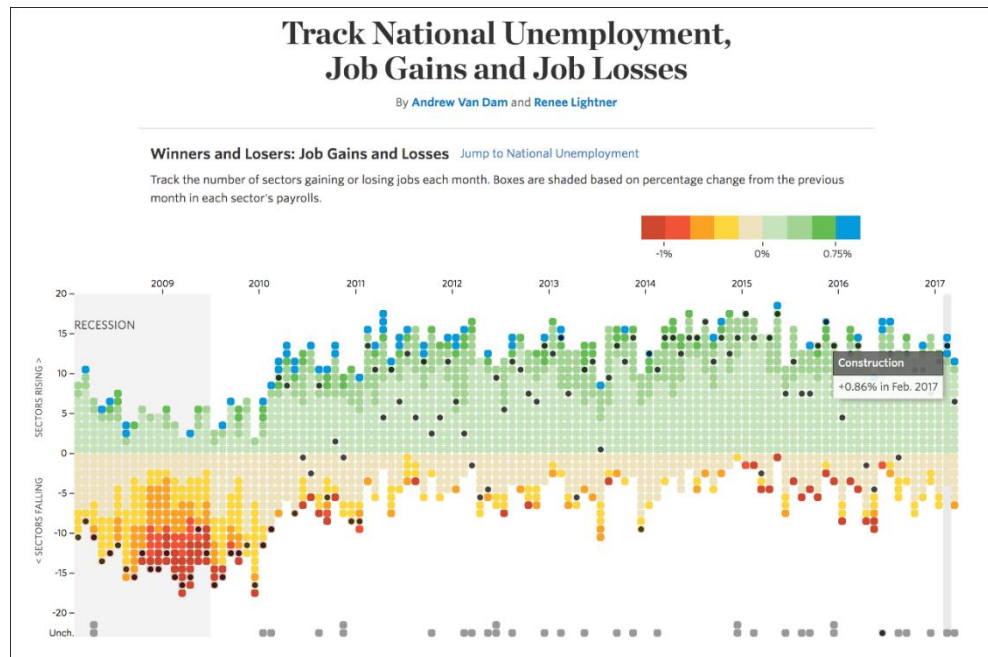


그림 1-8 데이터 시각화의 사례: 미국의 연도별 취업자와 실업자 통계

## 여기서 잠깐! 데이터 분석의 소요 시간

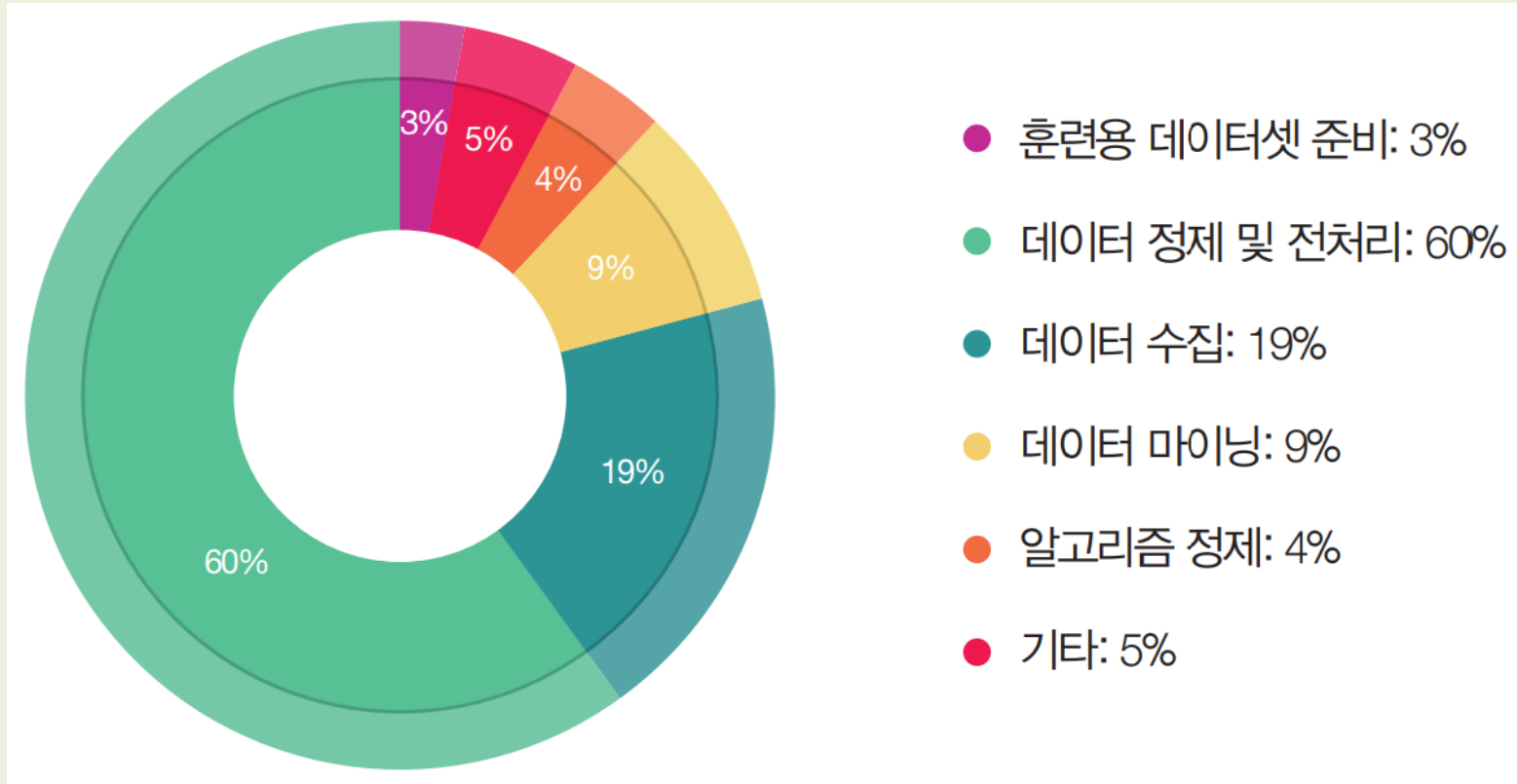


그림 1-9 데이터 분석 작업에 소요되는 시간

1. 데이터를 수집하는 일에 19%, 데이터를 정제하고 전처리하는 데 60%의 시간을 사용 → 즉, 전체 분석 과정에서 약 80%의 시간이 분석을 위한 데이터 준비에 사용
2. 이러한 시간을 얼마나 줄이느냐가 전체 분석 시간을 줄이는 관건