# Promoting Statistics for the Youth

August 28, 2013

# Byeong U. Park

Department of Statistics
Seoul National University

# What is Statistics?

- A study of the
  - *collection,*
  - *analysis,*
  - *interpretation*

  of **Data**.

- Linked to and applied to various real problems in industry as well as scientific problems in other disciplines, founded on mathematical thinking.
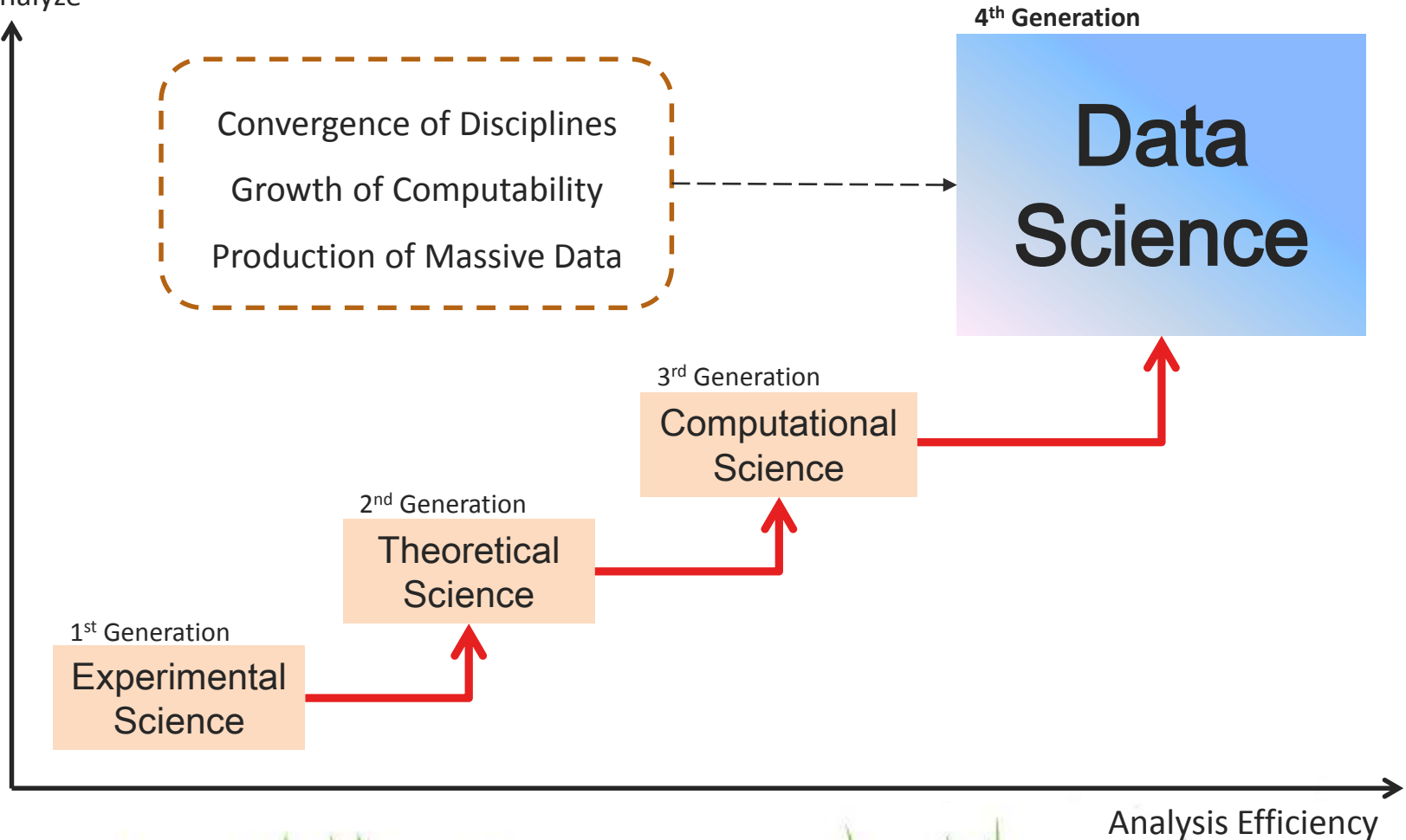
# What is Statistics?

# 100

# Science Paradigm for 21$^{st}$ Century

Complexity of Data
to Analyze

<Microsoft (2009), "The Fourth Paradigm; Data-Intensive Scientific Discovery". >

**4$^{th}$ Generation**

Data Science

Convergence of Disciplines

Growth of Computability

Production of Massive Data

3$^{rd}$ Generation

Computational Science

2$^{nd}$ Generation

Theoretical Science

1$^{st}$ Generation

Experimental Science

Analysis Efficiency

# Role of Statistics in Data Science

Massive Data
with **Uncertainty**

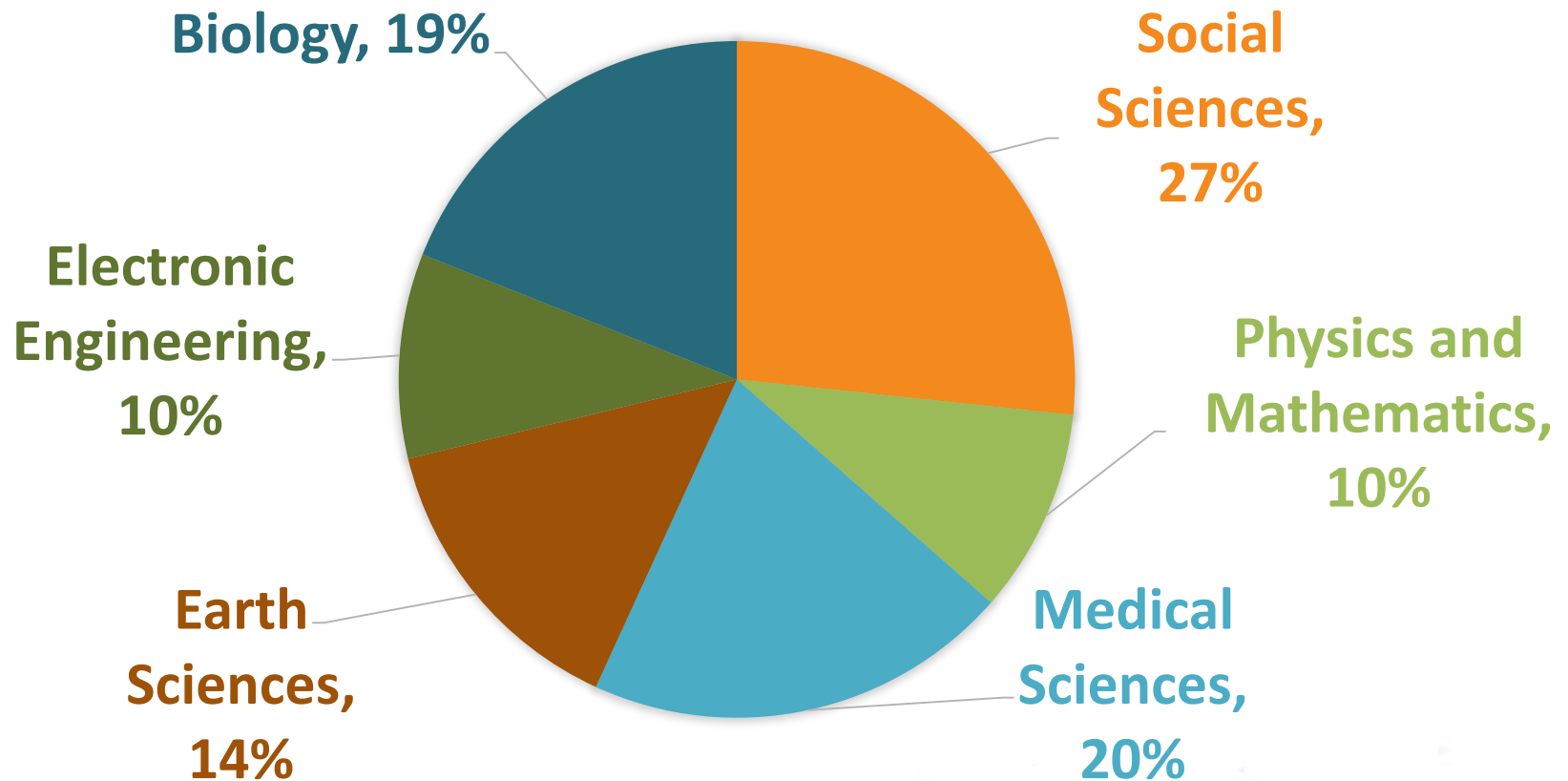**Statistical Modeling
and Data Analysis**

Knowledge Production
and Scientific Discovery
from **Data**

# Far-reaching Influence of Statistics

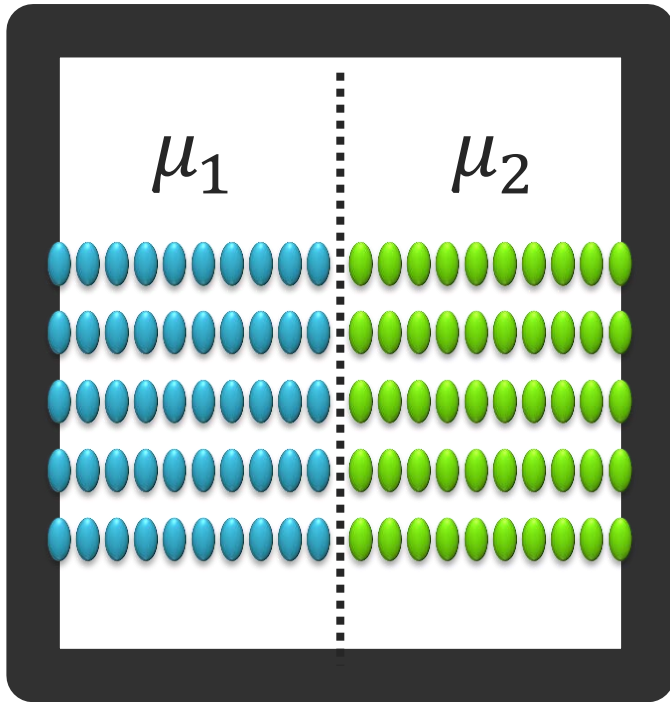- Citation distribution of Bootstrap paper (Efron, 1979)

# Far-reaching Influence of Statistics

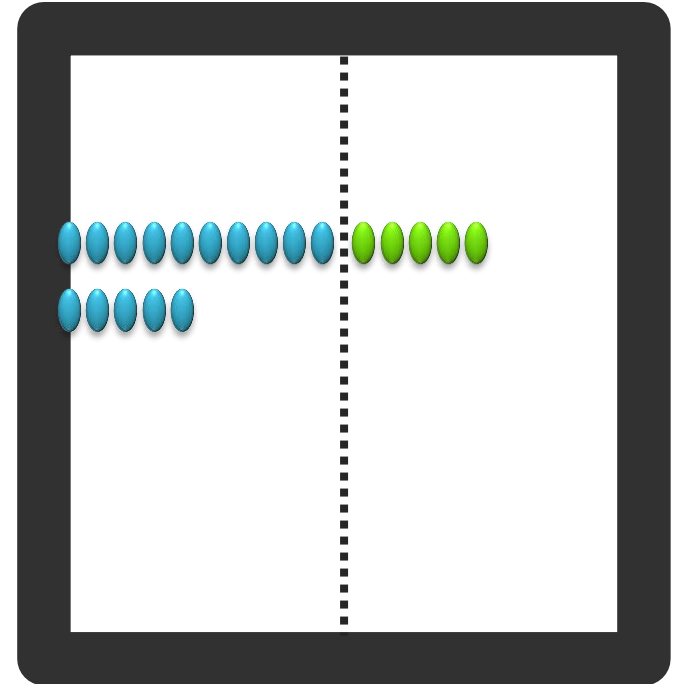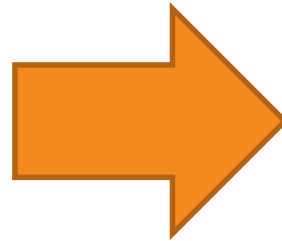| Statistical Topic | Author | Journal Title | Citation* |
|---|---|---|---|
| EM Algorithm | Dempster et al. (1977) | Journal of the Royal Statistical Society Series B | 33,645 |
| False Discovery Rate | Benjamini and Hochberg (1995) | Journal of the Royal Statistical Society Series B | 17,964 |
| Co-integration** | Engle and Granger (1987) | Econometrica | 20,249 |

* 2013.06.06    ** Nobel Prize in Economic Sciences (2003)
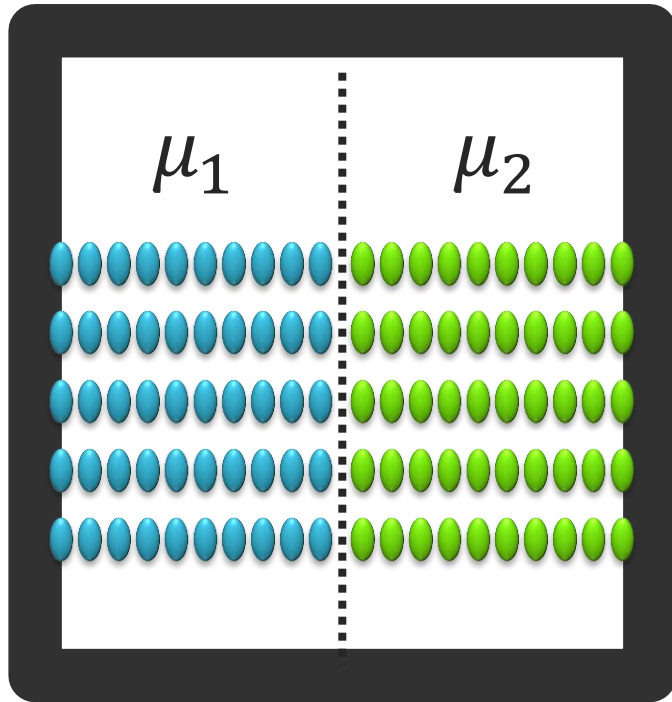
# Ball Sampling



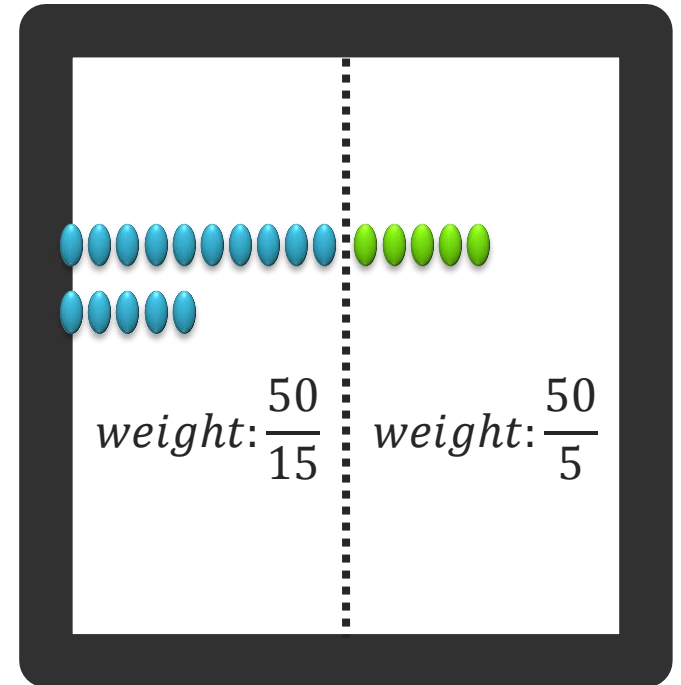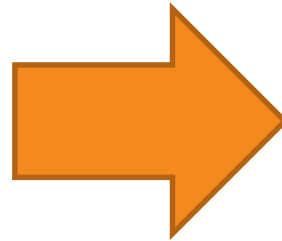$Average: \dfrac{\mu_1 + \mu_2}{2}$

$Average: \dfrac{15\mu_1 + 5\mu_2}{20} \left( \neq \dfrac{\mu_1 + \mu_2}{2} \right)$

# Ball Sampling



$Average: \dfrac{\mu_1 + \mu_2}{2}$

$Weighted\ Average: \dfrac{15\left(\dfrac{50}{15}\mu_1\right) + 5\left(\dfrac{50}{5}\mu_2\right)}{15 \times \dfrac{50}{15} + 5 \times \dfrac{50}{5}} \left(= \dfrac{\mu_1 + \mu_2}{2}\right)$

# Credit Loan Survey

- Investigated the relation between overdue status of credit loan and employment status to establish a strategy for allowing loan or not.

- Found that the group of people in their employment has a much higher overdue rate than the group of unemployed.

- **Why?**

# Chain-Ladder Data

| $i$ \\ $j$ | 1 | 2 | ... | ... | ... | ... | $n$ |
|---|---|---|---|---|---|---|---|
| 1 | $Z_{1,1}$ | $Z_{1,2}$ | ... | ... | ... | ... | $Z_{1,n}$ |
| 2 | $Z_{2,1}$ | $Z_{2,2}$ | ... | ... | ... | $Z_{2,n-1}$ | |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | | | |
| $i$ | $Z_{i,1}$ | $Z_{i,2}$ | ... | $Z_{i,n-i+1}$ | | | |
| ⋮ | ⋮ | ⋮ | | | | | |
| ⋮ | ⋮ | $Z_{n-1,2}$ | | | | | |
| $n$ | $Z_{n,1}$ | | | | | | |

(e.g. 1) Future liability for insurance
 i = Accident year,   j= development year

(e.g. 2) Divorce rate
 i = Marriage year,   j=Divorce year

# Body Fat Data

- Body Fat (%) = $\dfrac{\text{Total Weight of Fat}}{\text{Weight}} \times 100$

- Typical Body Fat Amounts

| Description | Men | Women |
|---|---|---|
| Exceptionally lean | 6 ~ 10 % | 10 ~ 15 % |
| Very lean | 11 ~ 14 % | 16 ~ 20 % |
| Lean* | 15 ~ 18 % | 21 ~ 25 % |
| Moderate | 19 ~ 24 % | 26 ~ 29 % |
| Obese | $\geq$25 % | $\geq$30 % |

* Thin but looks strong and healthy

# Body Fat Data



Men

| 3 - 4% | 6 - 7% | 10 - 12% |
| 15% | 20% | 25% |
| 30% | 35% | 40% |

Women

| 10-12% | 15-17% | 20-22% |
| 25% | 30% | 35% |
| 40% | 45% | 50% |

BUILTLEAN

# Body Fat Data

- Body fat percentage is a good health indicator.
- How to measure body fat percentage?
- MRI, CT: Accurate, but **expensive**



**Cheap** but less accurate equipment

| Body Fat Required Data Entry | | |
|---|---|---|
| Select Your Gender | Enter Your Weight In Pounds | |
| ◉ Male<br>◯ Female | Enter Your Waist Size In Inches | |

**Free** but fairly accurate calculator (statistical modeling and prediction)

http://www.csgnetwork.com/bodyfatcalc.html
http://extremebodyfit.com/fat-loss/how-to-calculate-your-body-fat-percentage

# Body Fat Data

- Nonparametric function estimation
  - Quantile additive model
  - Response(Y) = Body Fat percentage
  - Covariate(**X**) = (Age, Height, Abdomen, Hip)
  - $(Body\ Fat) = f_1(Age) + f_2(Height) + f_3(Abdomen) + f_4(Hip)$
    $+(measurement\ error)$

quantile(τ)=0.5
(median estimation)

# Era of Big Data

- **Production of Massive Data**
  - Over 95% of human-made data have been generated just during last 2 years.
  - 1.8 zettabytes in 2011 ($2 \times 10^{11}$ HD movies that can be watched for 47 million years)
  - Expect 50-fold increase by 2020.

- **Type of Data**
  - Conventional: number
  - New: images, sounds, texts, web logs (Facebook, Twitter, …)

# Simpson's Paradox

- ## **Berkeley Gender Bias**

  - Lawsuit for gender bias against women for admission to graduate schools in the fall 1973

  - Men were more likely than women to be admitted.

| Gender | Total Applicants | Admitted Percentage |
|--------|------------------|---------------------|
| Men | 8,442 | **44 %** |
| Women | 4,321 | **35 %** |

*https://en.wikipedia.org/wiki/Simpson's_paradox*

# Simpson's Paradox

— But, disaggregation of the data showed:

| Department | Men | | Women | |
|:---:|:---:|:---:|:---:|:---:|
| | **Applicants** | **Admitted (%)** | **Applicants** | **Admitted (%)** |
| **A** | 825 | **62 %** | 108 | **82 %** |
| **B** | 560 | **63 %** | 25 | **68 %** |
| **C** | 325 | **37 %** | 593 | **34 %** |
| **D** | 417 | **33 %** | 375 | **35 %** |
| **E** | 191 | **28 %** | 393 | **24 %** |
| **F** | 272 | **6 %** | 341 | **7 %** |

*https://en.wikipedia.org/wiki/Simpson's_paradox*

# Simpson's Paradox



- **Data Analysis Paper in *Science***
  - Bickel, Hammel and O'Connell (1975), "Sex Bias in Graduate Admissions: Data from Berkeley", *Science*, Vol. 187, p. 498-404.

- Size of box indicates relative number of applicants to the department.

# PageRank

- Since web pages are extremely diverse, the role of search engines are very important.

- Earlier search engines showed web pages related to the query only based on contents in the page.

    (e.g.) There are too many web pages containing the word "university".
        → Is it reasonable to rank pages just in the order of appearance
            counts of the word "university"?

- We need to consider which page is of the most importance.

# PageRank

- Measure of Importance
  - Page, et al. (1999) "The PageRank Citation Ranking: Bringing Order to the Web", *Technical Report*. Stanford InfoLab.
  - Importance based on the webpages linked to it

# PageRank

- Comparison of query for "university"

# PageRank

# Cluster Analysis on Facebook

- Salter-Townshend (2012), "Analysing My Facebook Friends", *Significance*, p. 40-42.

- A statistical analysis of the link pattern (grey lines) reveals information on the Facebook friends (circles).

- The algorithm picks up 8 groups as different colors.

- In fact Facebook uses closeness measures to make their friend suggestions.

  Cf) Yellow(family), Blue(girl friends), Red(dormitory friends), Green(ski club members), …

# Netflix Prize

- **Netflix**
  - Online DVD rental service company
  - (data) 100,480,507 ratings that 480,189 users gave to 17,770 movies

- **Objective**
  - Three-year (2007-2009) contest for movie recommendation system
  - Improving prediction algorithm for user ratings based on the **data**
  - At least 10% better performance than *Cinematch,* Netflix's algorithm

- **Competition Result**
  - $ 100,000,000 prize open competition
  - Winning to *BellKor* team, AT&T **statistician** group
  - The winning algorithm is also applicable to other marketing area.

# Higgs Boson Data

- ## Large Hadron Collider (LHC)

# Higgs Boson Data

- **Standard Model in Physics**
  - A theory concerning the electromagnetic, weak and strong nuclear interactions
  - There exists 25 adjustable parameters.

- **Higgs Boson**
  - An unconfirmed piece of the theory
  - Exist theoretically but no evidence
  - Strong *belief* for most physicists

# Higgs Boson Data

- **Key Goal of LHC**
  - Establish whether Higgs boson actually exists or not.
  - Measure, if so, properties of Higgs boson.

- **LHC Data**
  - LHC produce close to a billion events per seconds. ($\approx 10^{15}$ bytes)
  - (1 DVD $\approx$ 5 GB) 200,000 DVDs are needed for a second.
  - But, only a tiny fraction of those are of potential interest.

# Statistics in Other Disciplines

- **Bioinformatics**

  – **Bio**logy + **Informatics**

  – Genome sequence analysis and its management

  – Biologically novel discovery from genomic data



- **Medical Science**

  – Finding optimal therapy based on massive medical records combining health insurance data

# Statistics in Other Disciplines

- ## **Brain Science**

  – Complex understanding on brain function

  – Statistical analysis of signals on brain with high-dimensional images

  (e.g.) Functional Magnetic Resonance Image (fMRI) data



Thesen et al. (2012), "Sequential then interactive processing of letters and words in the left fusiform gyrus, *Nature Communications*, Vol. 3.

# Statistics in Other Disciplines

- ## Signal Processing
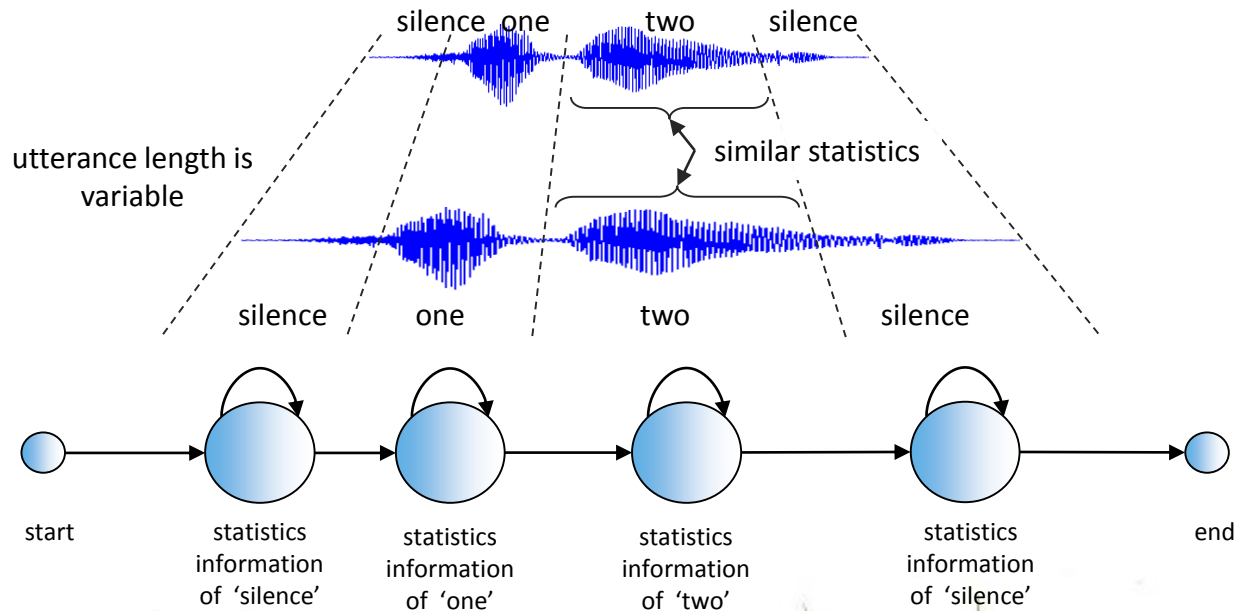  - Voice/picture recognition, Image denoising algorithm, …
  - Conventional: Physical pattern analysis
  - New: Statistical/Machine learning

# Statistics in Other Disciplines

- ## Earth Science

  – Meteorology, Climatology, Oceanography, Seismology, …

  – Probabilistic/Statistical models perform more efficient than computationally heavy physical-dynamic models

  (e.g.) Ensemble method

  – 3$^{rd}$ generation → 4$^{th}$ generation science

- ## Astronomy

  – Spatial-temporal distribution analysis for cosmic evolution

  – Measurement error or observation truncation on magnitude of stars

# And MANY others…!

# Prospects

- *For Today's Graduate, Just One Word: Statistics*
  - New York Times, August 5, 2009
  - "People think of field Archaeology as Indiana Jones, but much of what you really do is **data analysis**." – Carrie Grimes, Google statistician
  - "We are rapidly entering a world where everything is monitored and measured but the big problem is going to be **the ability of human to use, analyze and make sense of the data**" – Erik Brynjolfsson, director of MIT center for digital business

- *Statistics - Dream Job of the Next Decade*
  - "People can make the data tell a story, and everybody has data, but the problem is **how to utilize the data** more effectively." – Hal Varian, Google chief economist

# Thank You!

bupark@stats.snu.ac.kr