

Week 6: Statistics (Basic Math/Concepts for Stats.)

Instructor: Daejin Choi (djchoi@inu.ac.kr)



**INCHEON
NATIONAL
UNIVERSITY**

Before Starting

- In the lecture today, python-based codes are used for better explanation
- You can try if you have already installed, or you can use online python interpreter

A screenshot of a Google search results page. The search query "online python interpreter" is entered in the search bar. The results show several links to online Python compilers and interpreters:

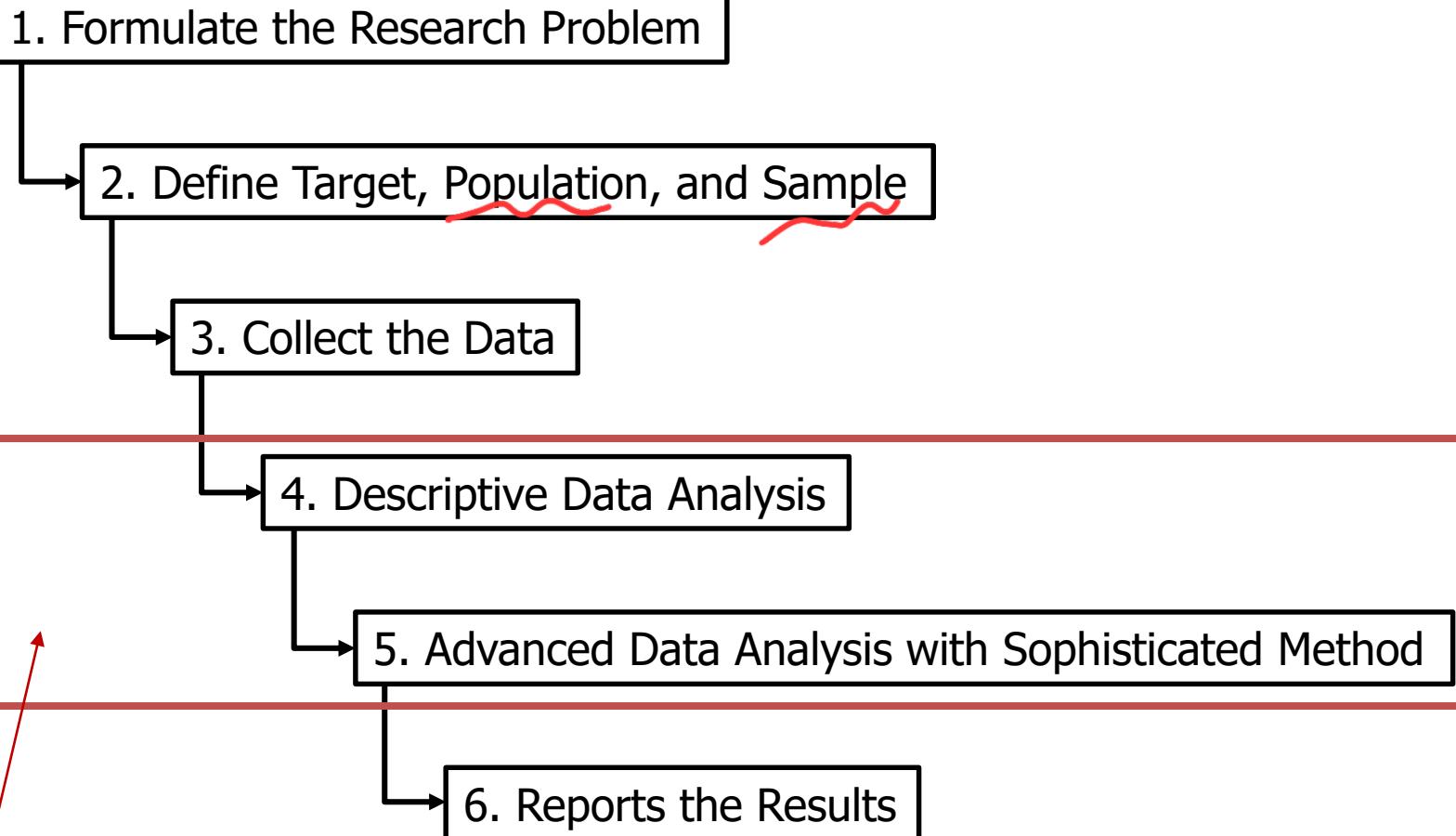
- repl.it** › languages › python3 ▾
[Python Online Compiler and IDE - Fast, Powerful, Free - Repl.it](#)
Repl.it is the world-leading online coding platform where you can collaborate, compile, run, share, and deploy Python online. Code in 50+ programming ...
- [www.onlinegdb.com](#) › online_pyth... ▾ 이 페이지 번역하기
[Online Python Interpreter - online editor - OnlineGDB](#)
OnlineGDB is online IDE with python interpreter. Quick and easy way to run python program online. It supports python3.
- [www.programiz.com](#) › online-com... ▾ 이 페이지 번역하기
[Online Python Compiler \(Interpreter\) - Programiz](#)
Write and run Python code using our online compiler (interpreter). You can use Python Shell like IDLE, and take inputs from the user in our Python compiler.
- [www.tutorialspoint.com](#) › execute_... ▾ 이 페이지 번역하기
[Online Python Compiler - Online Python Editor - Online ...](#)
... Python Coding Online, Practice Python Online, Execute Python Online, Compile Python Online, Run Python Online, [Online Python Interpreter](#), Execute Python ...
- [www.python.org](#) › shell ▾ 이 페이지 번역하기
[Welcome to Python.org](#)
Python is a programming language that lets you work quickly and integrate ... for Python's standard library along with tutorials and guides are available online

On the right side of the search results, there is a sidebar titled "함께 검색한 항목" (Search terms) with the following suggestions:

- 파이썬 온라인 에디터
- 온라인 파이썬 컴파일
- repl.it 파이썬
- 온라인 파이썬 코딩

Why Statistics?

- In general, analysis process is ...



You can easily start with statistical methods!

But, a little knowledge is required. → Today we will cover!

Contents for Today

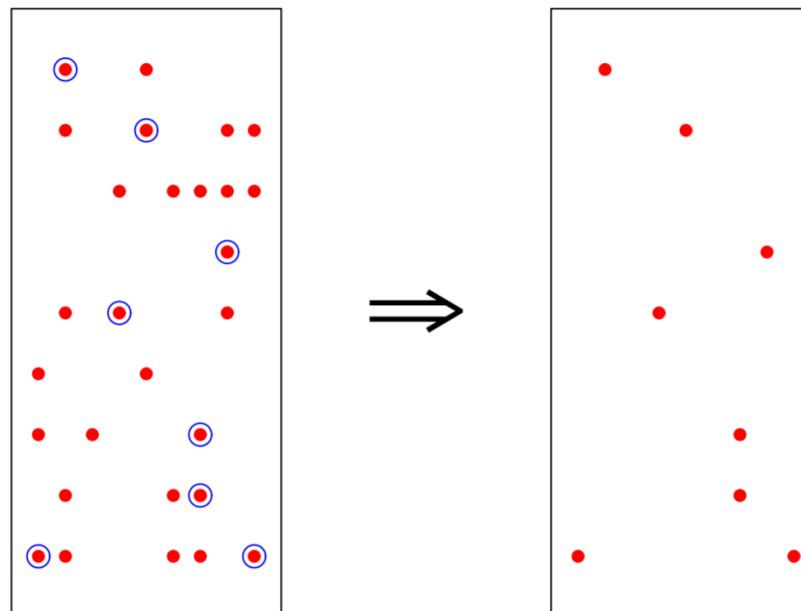
- Goal: Introducing the brief concepts / maths for Statistics.
- Contents
 - Terminologies
 - Central Tendency
 - Probability
 - Distribution & Information Theory

Terminologies

Terminology: Population vs. Sample

- **Population:** A (statistical) population is the set of measurements (or record of some qualitative trait) corresponding to the **entire collection of units** for which inferences are to be made
Target 집단 *수집한 데이터셋*
- **Sample:** A sample from statistical population is the set of measurements that are actually collected in the course of an investigation

Population vs. Sample



Terminology: Descriptive vs. Inferential

- Descriptive Statistics: consists of methods for **organizing and summarizing information** (Weiss, 1999) ⇒ 요약한 정보
- Inferential Statistics: consists of methods for **drawing and measuring the reliability of conclusions** about **신뢰성** population based on information obtained from a **sample** of the population (Weiss, 1999)
유의미하든 줄 나와내기 위한
통계학적 test
- NOTE: Descriptive and inferential statistics are **interrelated**
- Assume that you have just started data analysis. What you will do?
 - Preliminary descriptive analysis
 - (Choice of appropriate inferential method)
 - Inferential analysis

Terminology: Parameters vs. Statistics

- Parameters: an **unknown** numerical summary of the population (Agresti & Finalay, 1997) α, β
 - Statistics: a **known** numerical summary of the sample which can be used to make inference about parameters (Agresti & Finalay, 1997) 알려진 것, 평균 다는 것
↳ 방법 + 방법으로 나타나는 결과도
 - Example
 - The proportion p of 18-30 year-olds who use Twitter at least once a day → Parameter
 - The proportion p' of 18-30 year-olds using Twitter at least once a day, calculated from a sample of 18-30 year-olds
- sample에서 계산
적어도 한 번 사용 → 확률(p')

Terminology: Discrete vs. Continuous

유한하다

- Discrete Variables: indicates a finite numbers of values or as many values as there are integers. E.g.) tweet counts per day int
- Continuous Variables: indicates any real values e.g., time difference between two consecutive tweet postings R (실수)

중점적 경향

Central Tendency

Measure of Central Tendency: Mean

- The sum of observed values in a data divided by the number of observations
- p.s. Often affected by extreme values in the data (i.e., outliers)

```
>>> import numpy  
>>> numpy.mean(numpy.array([1,2,3,1,3,3,54,6]))  
9.125
```

Measure of Central Tendency: Median

- The value of the variable in a data set that divides the set of the observed values in half, so that the observed values in one half are less than or equal to the median value and the observed values in the other half are greater or equal to the median value
- To compute **median**, the observed values of variable in a data should be sorted
 - If the # observation is odd, then the sample median is the observed value exactly in the middle of the ordered list
 - If the # observation is even, then the sample median is the number halfway between two middle observed values in the ordered list

```
>>> import numpy  
>>> numpy.median(numpy.array([1,2,3,1,3,3,54,6]))
```

3

Measure of Central Tendency: Mode

- A discrete variable (or its value) is that value of the variable which occurs with the greatest frequency in a data set

```
>>> from statistics import mode  
>>> mode([1, 1, 2, 3, 3, 3, 3, 4])  
3
```

최빈

- p.s. if the greatest frequency is 1 (i.e., no value occurs more than once), then the variable has no mode

Measure of Variation: Range

- The difference between its maximum and minimum values in a data set
- Formally, $\text{range} = \text{max.} - \text{min.}$
- p.s. range is affected by only max. and min. values, which means that other values may be ignored.
 - E.g., $X = [1, 10, 10, 11, 11, 10, 11, 10, 11, 10, 11, 11, 10, 11, 100]$
 - Can we say the range of X is $[1, 100]$?

Measure of Variation: Percentile & Quartiles

- A method dividing a given dataset with the equal portions
 - Quartile (4 parts)
- Let n denote the number of observations in a data set.
Arrange the observed values of variable in a data in increasing order
 - The 1st quartile Q_1 is at position $(n+1)/4$
 - The 2nd quartile Q_2 (the median) is at position $(n+1)/2$
 - The 3rd quartile Q_3 is at position $3(n+1)/4$In the ordered list
- Interquartile range: the difference between 1st and 3rd quartiles of the variable (i.e., $\text{IQR} = Q_3 - Q_1$)
 - Roughly, the IQR gives the range of the middle 50% of the observed values

Measure of Variation: Interquartile Range = IQR

- The sample of interquartile range represents the length of the interval covered by the center half of the observed values of the variable
- This measure of variation is not disturbed if a small fraction the observed values are very large or very small
outlier 제거 효과
- Five-number summary: consists of min., max., and quartiles written in increasing order
 - Min, Q1, Q2, Q3, Max

Measure of Variation: Standard Deviation

- A measure of a distribution's deviation from its mean. (square root of variance) 
- For a variable x , the sample standard deviation is

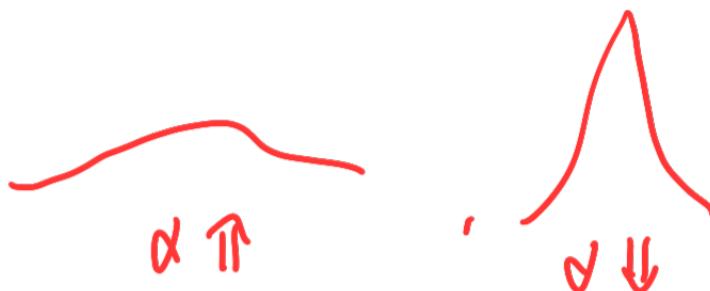
$$\sigma = \sqrt{\frac{1}{N} [(x_1 - \mu)^2 + (x_2 - \mu)^2 + \cdots + (x_N - \mu)^2]}, \text{ where } \mu = \frac{1}{N}(x_1 + \cdots + x_N),$$


- Preferred when the mean is used as the measure of center (e.g., symmetric distribution)

Measure of Variation: Standard Deviation

- Trivially, more variation → larger standard deviation
- However, it can be strongly affected by a few extreme observations

```
>>> import numpy  
>>> numpy.std([0.1,2.8,3.7,2.6,5,3.4])  
1.4851
```



Probability

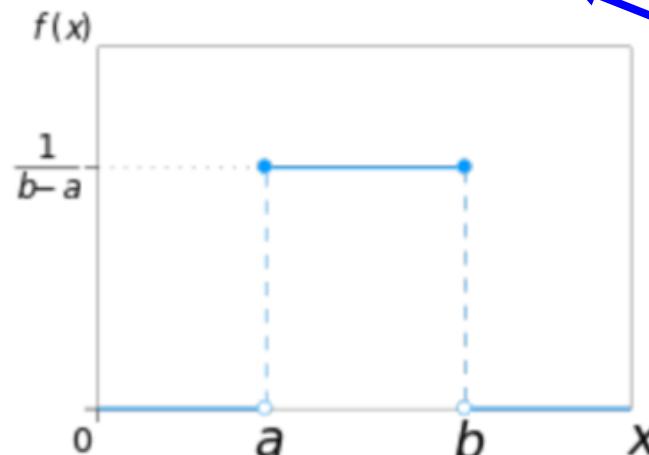
Probability Mass Function (PMF)

- For discrete random variables 이산적 확률
- The domain of P must be the set of all possible states of x
- $\forall x \in X, 0 \leq P(x) \leq 1$. An impossible event has probability 0 and no state can be less probable than that. Likewise, an event that is guaranteed to happen has probability 1, and no state can have a greater chance of occurring
- $\sum_{x \in X} P(x) = 1$. We refer to this property as being normalized. Without this property, we could obtain probabilities greater than one by computing the probability of one of many events occurring
- Uniform distribution among k states: $P(x = x_i) = \frac{1}{k}$

Probability Density Function (PDF)

- For continuous random variables
- The domain of P must be the set of all possible states of x
- $\forall x \in X, p(x) \geq 0$. Note that we do not require $p(x) \leq 1$
- $\int p(x)dx = 1$
- Uniform distribution

$$u(x; a, b) = \frac{1}{b-a}.$$



Parameterized by

Rules associated with Probability

■ Computing marginal probability with the Sum Rule

- $\forall x \in X, P(x = x) = \sum_y P(x = x, y = y).$
- $p(x) = \int p(x, y) dy.$

Sum Rule

■ Conditional Probability

$$P(y = y | x = x) = \frac{P(y = y, x = x)}{P(x = x)} = \frac{p(x \wedge y)}{p(x)}$$

■ Bayes' Rule:

$$\underline{P(x | y) = \frac{P(x) P(y | x)}{P(y)}}$$

$$P(A|B \cap C) = \frac{P(A \cap B \cap C)}{P(B \cap C)}$$

■ Chain Rule of Probability

$$P(x^{(1)}, \dots, x^{(n)}) = P(x^{(1)}) \prod_{i=2}^n P(x^{(i)} | x^{(1)}, \dots, x^{(i-1)}).$$

E.g., $P(a, b, c) = P(a|b,c) P(b|c) P(c)$ $P(A|B \cap C) \times P(B|C) \times P(C)$

Independence

$$\forall x \in X, y \in Y, p(X = x, Y = y) = p(X = x)p(Y = y).$$

Notation: $x \perp \underbrace{y}_{\text{independence}}$ $P(A)P(B) = 0$

■ Conditional Independence

- $\forall x \in X, y \in Y, z \in Z$
- $p(X = x, Y = y | Z = z) = p(X = x | Z = z)p(Y = y | Z = z)$
= conditional Independence
- Equivalently, $p(x|y, z) = p(x|z)$
 - Proof?
- Notation: $x \perp y | z$

, 독립

Expectation

- Discrete variable: $\mathbb{E}_{\tilde{x} \sim P}[f(x)] = \sum_x P(x)f(x),$
- Continuous variable: $\mathbb{E}_{\tilde{x} \sim p}[f(x)] = \int p(x)f(x)dx.$
- Linearity of Expectations
 - NOTE: this always holds, even when $f(x)$ and $g(x)$ are dependent

$$\mathbb{E}(\alpha x + \beta y) = \alpha \mathbb{E}(x) + \beta \mathbb{E}(y)$$

Covariance

$$\text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))]$$

$$\text{Cov}(f(x), g(y)) = \mathbb{E} [(f(x) - \mathbb{E}[f(x)])(g(y) - \mathbb{E}[g(y)])].$$

Note: Covariance can be negative 음수가 될 수 있음

- Covariance Matrix: $\text{Cov}(\mathbf{x})_{i,j} = \text{Cov}(x_i, x_j)$.
 - Diagonal elements are variance of i-th element

Randomized Algorithms (=Probabilistic Algo.)

	전투경사 Las Vegas	Monte Carlo
Type of Answer	Exact	Random amount of error
Runtime	Random (until answer found)	Chosen by user (longer runtime gives less error)

Estimating sums / integrals with samples

$$s = \sum_{\mathbf{x}} p(\mathbf{x})f(\mathbf{x}) = E_p[f(\mathbf{x})]$$

$$s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x} = E_p[f(\mathbf{x})]$$

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)}).$$

Justification

- Unbiased:
 - The expected value for finite n is equal to the correct value
 - The value for any specific n samples will have random error, but the errors for different sample sets cancel out
- Low variance:
 - Variance is $O(1/n)$
 - For very large n , the error converges “almost surely” to 0

Distribution & Information Theory

Bernoulli Distribution

- PDF

베르누이

$$P(x = 1) = \phi$$

$$P(x = 0) = 1 - \phi$$

$$P(x = x) = \phi^x (1 - \phi)^{1-x}$$

$$\phi^x (1 - \phi)^{1-x}$$

- Expectation: $\mathbb{E}_x[x] = \phi$

$$E(x) = P$$

- Variance: $\text{Var}_x(x) = \phi(1 - \phi)$

$$V(x) = PQ$$

Gaussian Distribution

- Parameterized by variance:

- $\underbrace{E[x]}_{\mu}, \underbrace{\text{Var}[x]}_{\sigma^2}$

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right).$$

정규분포

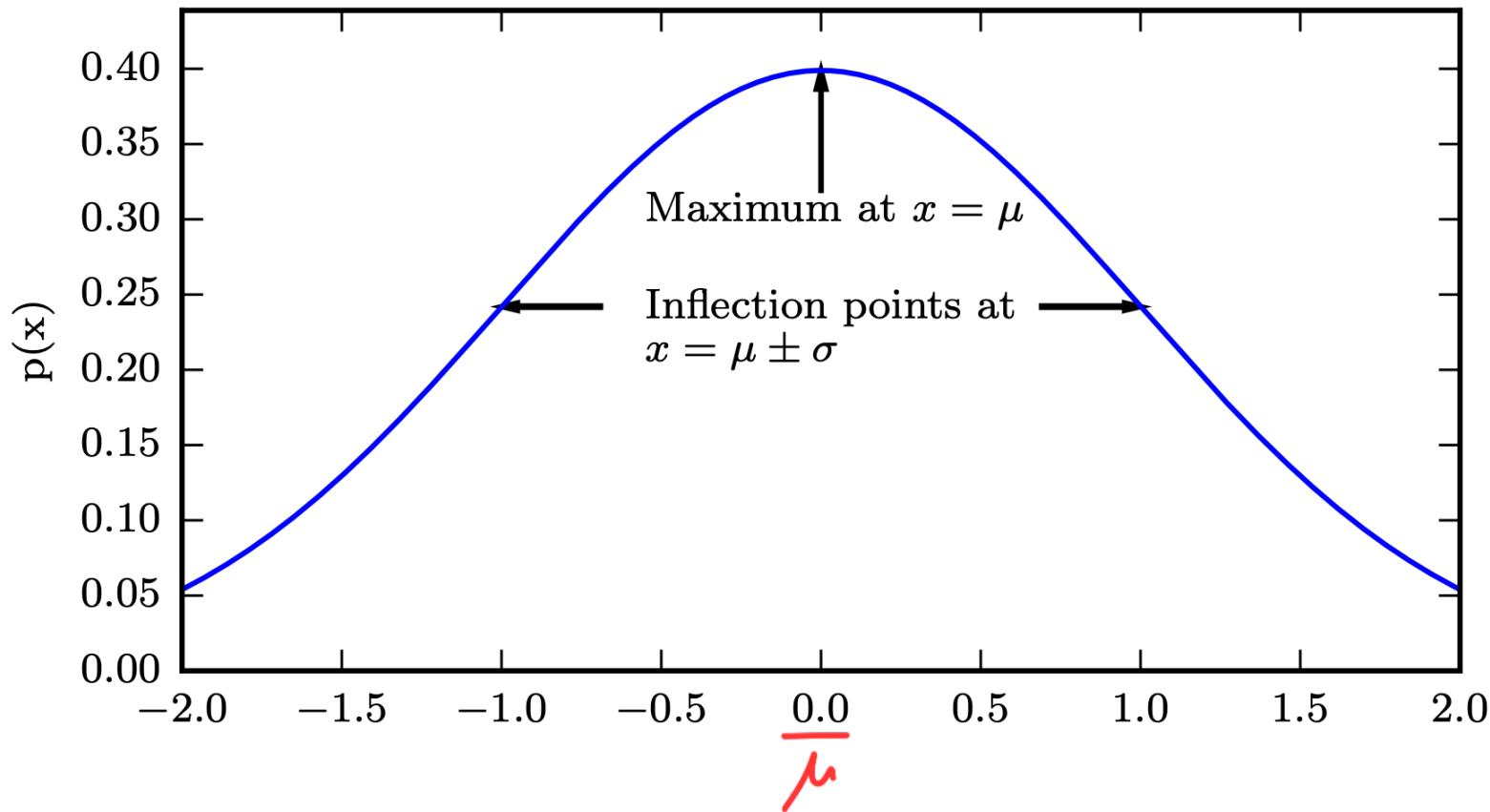


- Parameterized by precision:

$$\mathcal{N}(x; \mu, \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{1}{2}\beta(x - \mu)^2\right).$$

$$\sigma^2 = \frac{1}{\beta}$$

Gaussian Distribution



Theorems related to Gaussian Distribution

- Central limit theorem: the sum of many independent random variables is approximately normally distributed

$$\frac{\sqrt{n}}{\sigma} (\overline{X_n} - \mu) \xrightarrow{\text{Normalize}} N(0,1) \quad \text{as } n \rightarrow \infty$$

대수법칙

- Law of large numbers: the sample average converges to the expectation as the sample size goes to infinity

$$\overline{X_n} \rightarrow \mu \quad \text{as } n \rightarrow \infty, \text{ where } \overline{X_n} = \frac{1}{n} (X_1 + X_2 + \dots + X_n)$$

$$E(x) = \mu$$

기대값이 평균에 수렴한다

Multivariate Gaussian

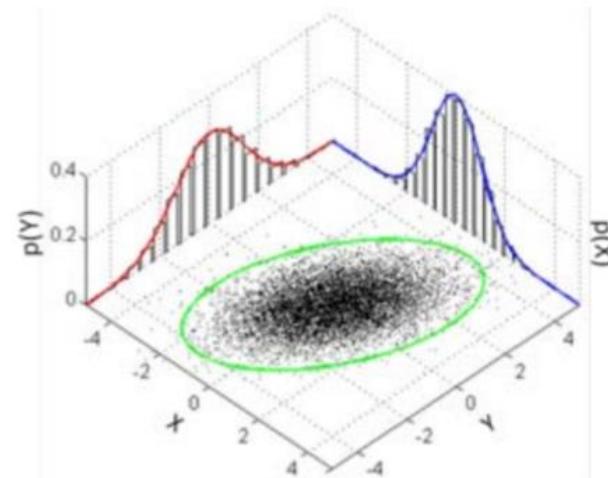
- Parameterized by covariance matrix

$$\mathcal{N}(\underline{x}; \mu, \Sigma) = \sqrt{\frac{1}{(2\pi)^n \det(\Sigma)}} \exp\left(-\frac{1}{2}(\underline{x} - \mu)^\top \Sigma^{-1}(\underline{x} - \mu)\right).$$

\downarrow
vector

ad-

- μ is a vector
- Σ is a covariance matrix



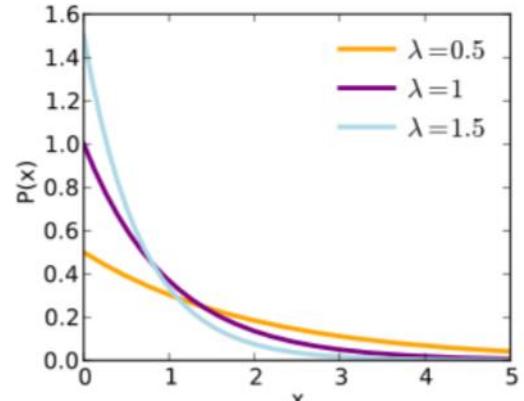
- Parameterized by precision matrix

$$\mathcal{N}(\underline{x}; \mu, \beta^{-1}) = \sqrt{\frac{\det(\beta)}{(2\pi)^n}} \exp\left(-\frac{1}{2}(\underline{x} - \mu)^\top \beta(\underline{x} - \mu)\right).$$

More Distributions

- Exponential:

$$p(x; \lambda) = \lambda \mathbf{1}_{x \geq 0} \exp(-\lambda x).$$

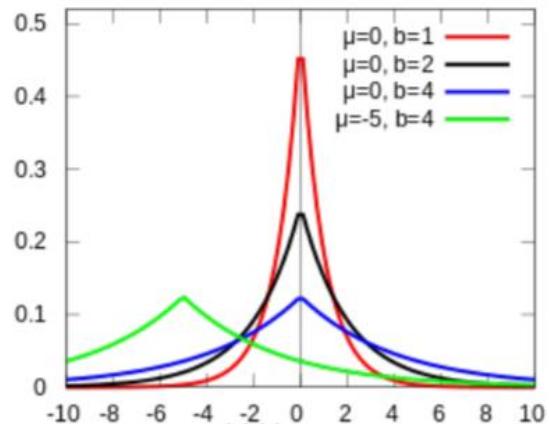


- Laplace:

$$\text{Laplace}(x; \mu, \gamma) = \frac{1}{2\gamma} \exp\left(-\frac{|x - \mu|}{\gamma}\right).$$

- Dirac Delta: It is zero-valued everywhere except at μ , yet integrates to 1

$$p(x) = \delta(x - \mu).$$



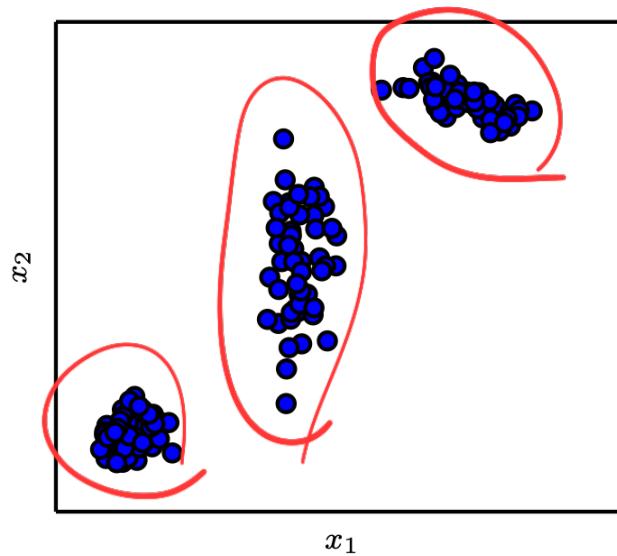
- Empirical Distribution:

$$\hat{p}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x} - \mathbf{x}^{(i)})$$

Mixture Distributions

$$P(\mathbf{x}) = \sum_i P(c = i)P(\mathbf{x} | c = i)$$

- Gaussian mixture: $P(x|c = i)$ is Gaussian

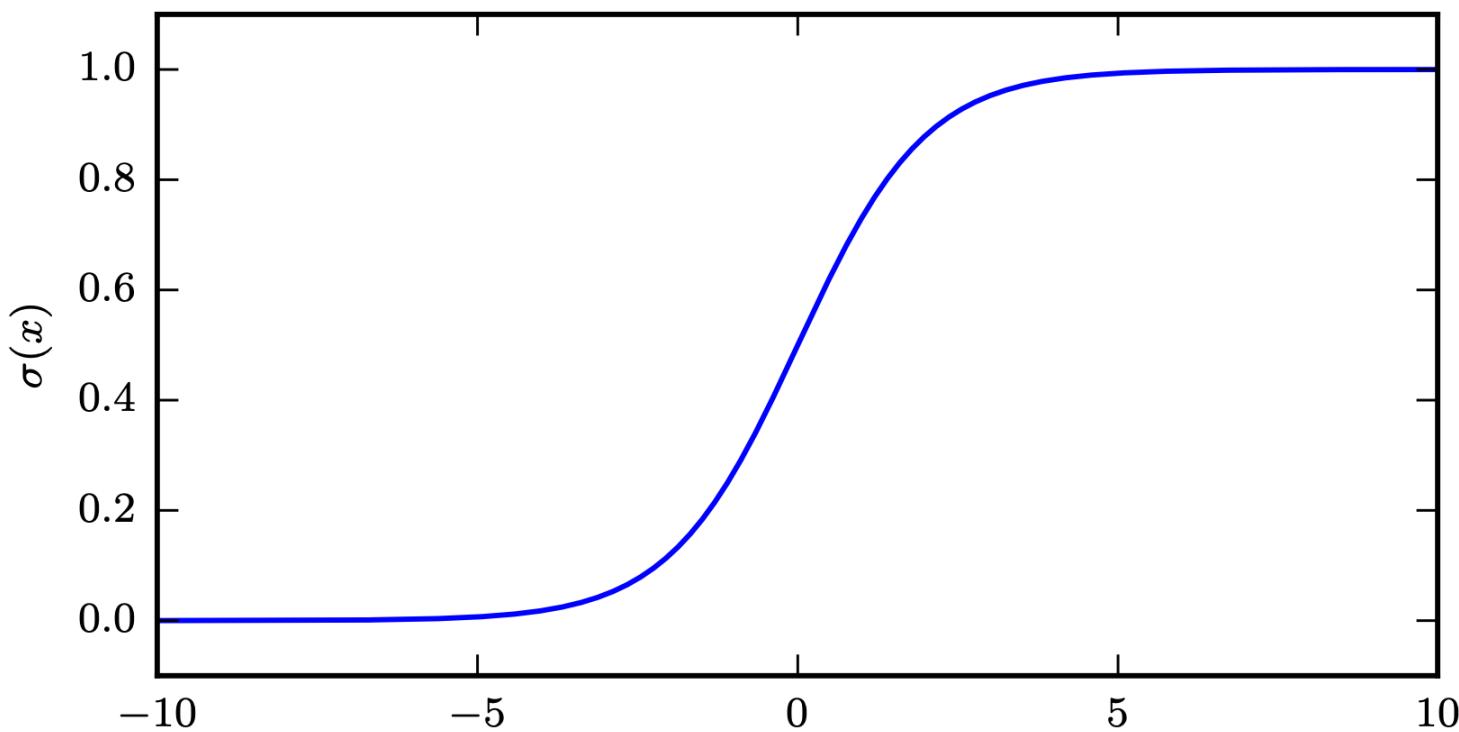


복분 복분이
가우시안 분포 $\sim \mathcal{N}$

Gaussian mixture with three components

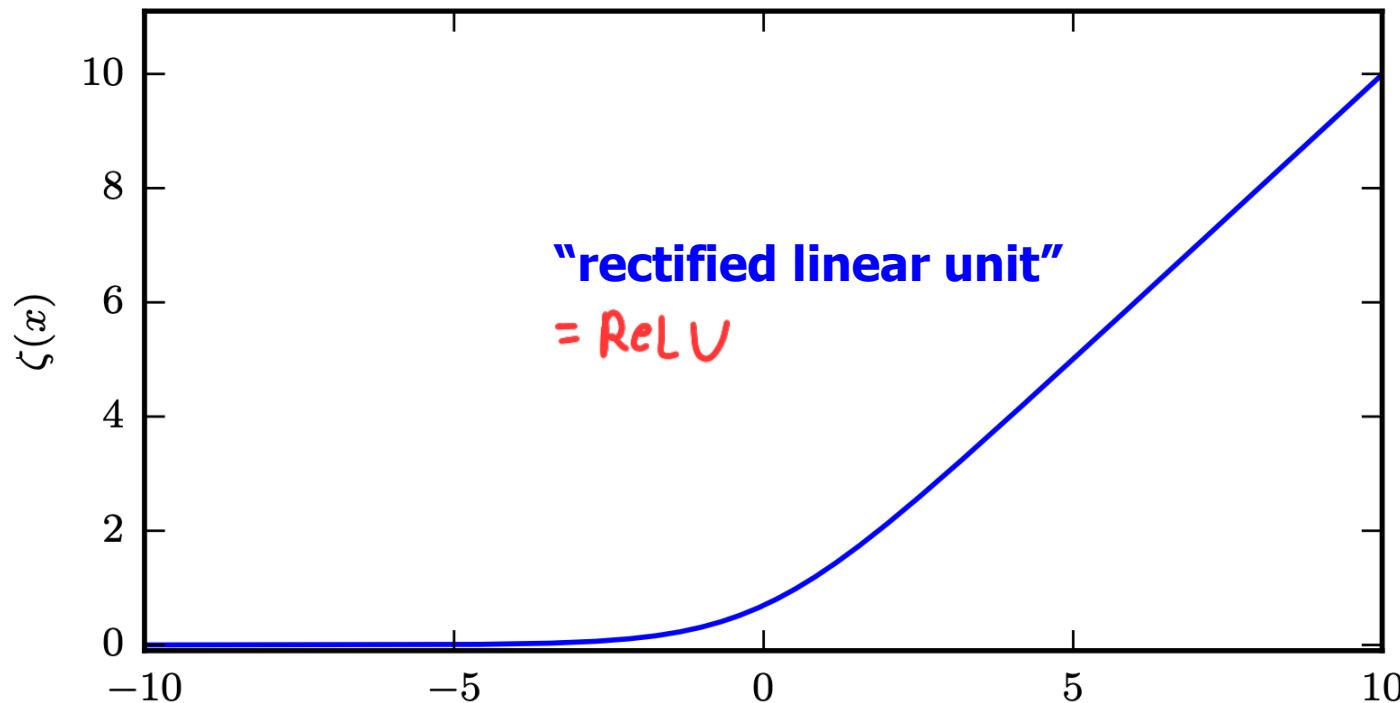
Logistic Sigmoid

- $\sigma(x) = \frac{1}{1+\exp(-x)}$



Softplus Function

- $\zeta(x) = \log(1 + \exp(x))$



- “softened” version of $x^+ = \max(0, x)$



Information Theory

- Information theory: quantifying how much information is present in a signal
정보의 양을 측정
- Learning that an unlikely event has occurred is more informative than learning that a likely event has occurred 확률적으로 적은 것들을 갖고오는데 더 유효하다
- Self-Information of x
 - $I(x) = -\log P(x)$
 $P(x) \uparrow \rightarrow I(x) \downarrow, P(x) \downarrow \rightarrow I(x) \uparrow$
 - Intuition: minimum # of bits to express (encode) an event with probability $P(x)$ $\log \frac{1}{P} \rightarrow x$ 값을 나타내기 위한 최소 bit
 - Rare event has a large information content

Information Theory: Entropy

불확실성

■ Entropy: expectation of self-information

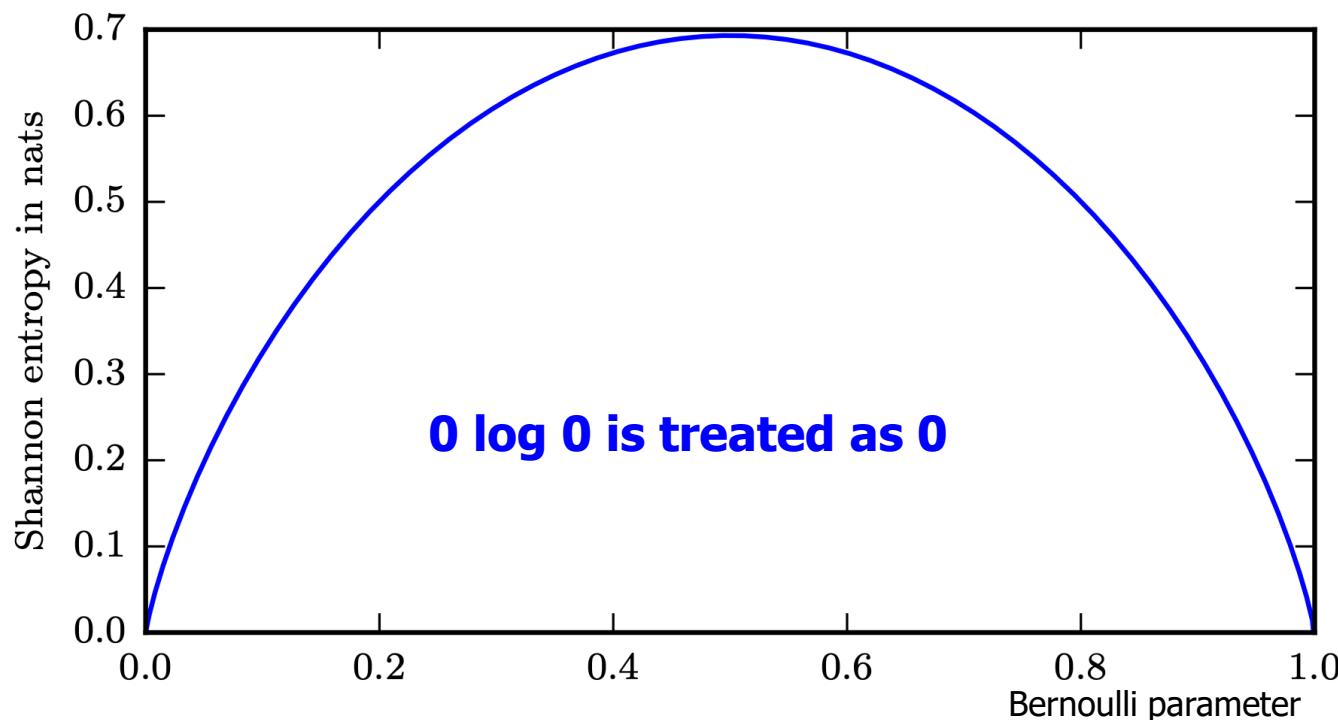
- Minimum expected # of bits to express a distribution

- For Bernoulli variable,

$$H(x) = -p \log p - (1 - p) \log(1 - p)$$

$H(x) \uparrow$ = 무규칙
 \downarrow = 규칙

$E(x) = p(x) \cdot f(x)$



KL Divergence

- Measure the difference of two distributions $P(x)$ and $Q(x)$ 두 분포 $P(x)$ 와 $Q(x)$ 가 얼마나 차이가 있는지

$$D_{\text{KL}}(P\|Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)].$$

- Properties
 - Always nonnegative: 0 if and only if P and Q are the same
 - Intuition: If $x \sim P$, the best (minimal) encoding is given by assigning $\log P(x)$ bits for each x
 - Not symmetric: 교환법칙 성립X

$$\Leftrightarrow D_{\text{KL}}(P\|Q) \neq D_{\text{KL}}(Q\|P)$$

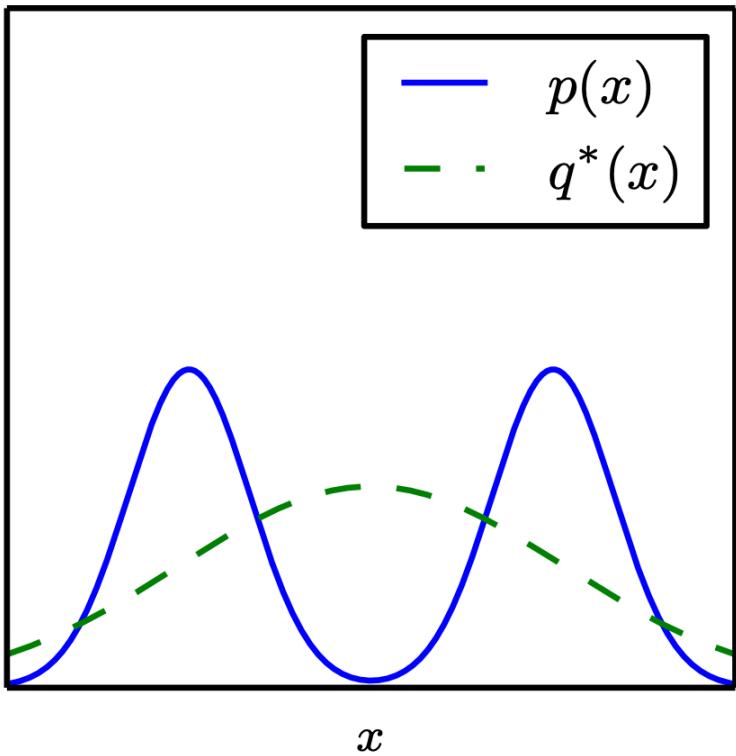
KLD \downarrow = 서로 수렴하는

KLD \uparrow \Rightarrow 서로 많이 떨어짐

KL is Asymmetric

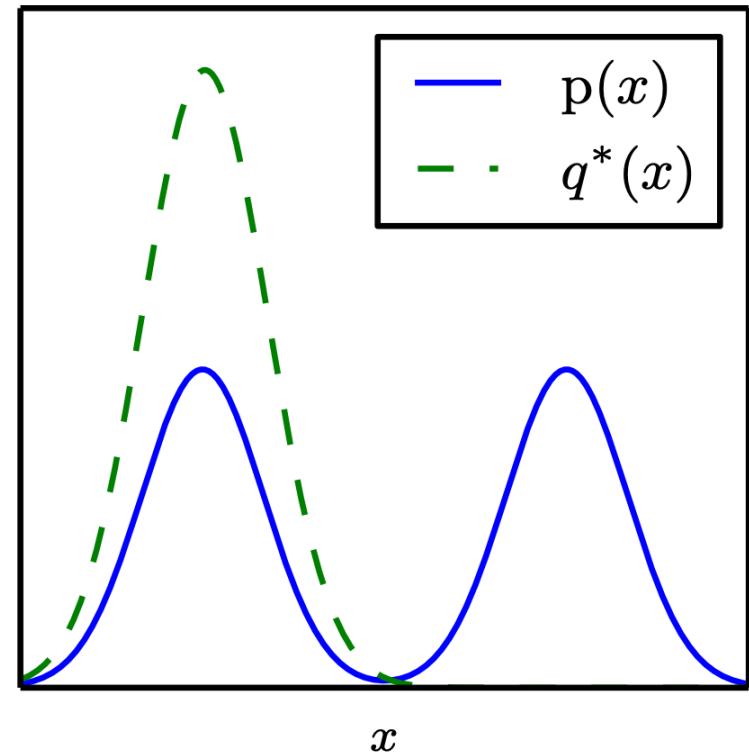
$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p\|q)$$

Probability Density



$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q\|p)$$

Probability Density



Cross-Entropy

- Average # of bits needed to identify an event from the true distribution P , if we use a coding scheme optimized for unnatural distribution Q

- $$\underline{H(P, Q)} = H(P) + \underline{D_{KL}(P||Q)} = -E_{x \sim P} \log Q(x)$$

Entropy + KLD

P 의 불포함 확률
 $\log Q$ 의 평균

- Minimizing the cross-entropy w.r.t. Q is equivalent to minimizing the KL divergence

$KLD \uparrow \rightarrow CE \uparrow$

linear 함

Summary & Next

- “Remember” the key terminology / concepts that popularly used in Stats.
 - Population, Sample, Descriptive, Inferential, Parameters, Statistics, ...
 - Mean, Median, Mode, Range/Quartile, Standard Dev.
 - Probability, Expectation, Covariance
 - Distributions
 - Self-Information, Entropy, KL-Divergence, Cross-Entropy
- Next?
 - Statistical methods (Pearson, ...)

Thank you!

Instructor: Daejin Choi (djchoi@inu.ac.kr)