

# I. Basic Statistics.

Probability: 확률의 법칙.

Statistics: Data → Probability → 의사결정.  
기록 + 측정, 비교, test.

Bayesian vs. Frequentist

: 바뀌어 가는 목표와 추측하는 양도 다 다를 수 있다.

Properties from Population.

Sample of population

: Useful  
Not perfect: (Caution)

→ Estimate & iid

Q. iid라면 어떻게 봐야?  
A. iid도 상황에 따라 다르다.  
Maximum Likelihood method  
production (CP)로 사용되어  
못 하기 때문.  
And iid라는 가정 자체가  
Sample에 따라 다른 값이 나올  
수 있다는 것. 특히 확률 분포의  
특징을 무시하게 된다.

Distribution of Random Variable → Empirical Distribution (Sample version)

Graph: Histogram/Boxplot  
CDF, Scatter plot

Mean / Variance and the other moments.

Q. 어떤 값에 대한 분포?  
Distribution에 대한  
특징이 있어야.

→ Expectation으로 편하게 생각해 줘.

Sample Mean / Variance

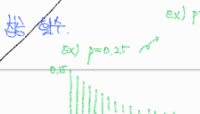
Median, Percentiles, IQR  
Mode, Range

## Distribution

Bernoulli: 1번 시행

Generalization

Binomial: n번 시행, k번 성공 확률.



Geometric (기하)

: 성공 1번

나를 때까지.

Q. 1과 2를 주고 주어진 확률에  
대응하는 값을 구하라.

A. 성공할 확률은 확률 1과 2를 구하라.

Continuous

Exponential

: 성공 1번 나올 때까지

대기 시간

Generalization

Beta

Bayesian의 Prior인 확률.

Continuous

Dirichlet

:  $\frac{\Gamma(\sum \alpha_i)}{\prod \Gamma(\alpha_i)} \prod x_i^{\alpha_i-1}$ ,  $\sum \alpha_i = 1, x_i \geq 0$

: 모든 범위의 개수  
가져오는 비율.

Uniform: 모든 값이 나올 확률이 동일.

Log Normal: Gaussian을 exponential한 것.

Student's t: Gaussian을 따르는 변수에서 Sample 수가 작을 때

Normal:  $p \rightarrow 0$  일 때

Poisson:  $\lambda \rightarrow 0$  일 때

Chi-squared:  $N(0,1)$ 의 제곱합.

2nd moment의

증명을 살피는 더 유용.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Poisson은 성공 4번은 Binomial의  
한 번, Exponential은 성공 4번은  
Geometric의 한 번.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* Scale 이후에 Data를 Standardization하면  $N(0,1)$ 로 변환하고 Data가 잘 분포된 것인지 확인  
주요 2nd moment 즉, Data가 잘 분포된 것  
보여 준다.

\* 분포에 대해 말하는 이유? 우리가 Control할 수 없는 요인들을 control 하려고 하지 말고 (Random이니까)

최소한의 노력으로 최대한의 효과를 기대하는 것. (Expectation을 보는 이유 중 하나)

# C.I & Test.

• Estimator: Mean & Estimator. →

Confidence Intervals:

$$\begin{array}{c} -se \cdot z \quad tse \cdot z \\ \hline \text{Mean} \end{array}$$

z: critical value at  $\alpha$

• CLT

: Mean 1, 2, 3, ... from Sample.

$$\rightarrow N(\mu, \frac{\sigma^2}{n})$$

$$\text{Variable: } \bar{X}_i \rightarrow E(\bar{X}) = \frac{\sum \mu}{n} = \frac{n \cdot \mu}{n} = \mu.$$

$$\text{Var}(\bar{X}) = \frac{\sum \sigma^2}{n^2} = \frac{n \sigma^2}{n^2} = \frac{\sigma^2}{n}$$

\* Sample set # vs. # of one sample set data.  
(Mean) (Variable)

Poisson ver. t-test (Mean = Variance)

$$\begin{array}{l} H_0: \lambda = \theta \\ H_1: \lambda < \theta \end{array} \quad \text{Reject } H_0 \text{ if } P(X \leq x) \leq \alpha \iff |t| = \left| \frac{\bar{x} - \theta}{\sqrt{\bar{x}}} \right| \leq t_{\alpha/2} (n-1) \quad \text{nominal significance level}$$

→ Test at the nominal significance level

\* Why? 정규 분포 가설로 small sample일 때, + dist가 나왔기 때문.

t-test

$$t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}, \quad H_0: \mu = \mu_0$$

if  $\sigma$  is unknown, use  $\hat{\sigma} = \frac{s}{\sqrt{n-2}}$

Hypothesis testing.

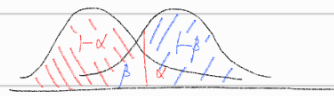
A/B Tests (Two-sample tests)

: pooled estimation으로 sample과 같은 분포 가정에 해당함. (Homoskedastic)

Type I error:  $H_0$  선택해서 틀림.

vs

Type II error:  $H_0$  선택해서 틀림.



$\alpha$ :  $H_0$  참일 때,  $H_1$ 이 틀릴 확률.

$1-\alpha$ :  $H_0$  참일 때,  $H_1$ 이 맞을 확률.

$\beta$ :  $H_1$  참일 때,  $H_0$ 이 틀릴 확률.

$1-\beta$ :  $H_1$  참일 때,  $H_0$ 이 맞을 확률.

→ 대수의 법칙 기억하기

\* Q. 대체  $\alpha$ 가 낮아지면  $1-\beta$ 가 낮아지는 것인가?

$\alpha$ 는  $H_0$ 가 참일 때  $H_1$ 을 선택할 확률이 낮아지는 것이므로

$H_1$ 가 참일 때 관측으로 보면  $H_0$ 의 틀릴 확률( $\beta$ )이 높아

지는 것이므로 선택할 수 있고( $\beta \uparrow$ ) 맞을 확률이 높아진다고

생각할 수 있다( $1-\beta \downarrow$ )

\* 결론  $H_0$ 과  $H_1$  둘 중 하나는 맞아야 하는 서로의 가정

속에서 누가 맞을 확률이 높아지면 반대쪽은 틀릴 확률이 높아지는 것.

$H_0$ 를 바탕으로 우의 estimation이 진행됐고 test-statistics이 나옴.

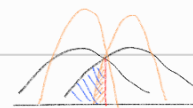
\* Data 1000개 짜리 Data set 하나.

.. 30개 짜리 .. 30개

들은 뭐가 다를까?

반정 cases: ① 분산 &  $\alpha$ 가 변할 때,

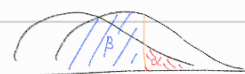
$$1-\beta \propto \frac{\alpha}{\text{분산의 크기}}$$



⇒  $\frac{\sigma^2}{n}$ 이 커지면 데이터가 커지면 됨 (검정하는 data set)

\* Tip.  $\alpha$ 와  $\beta$ 는 trade off 관계니까 model의 power를 높이기 위해서는 data를 많이 모아줘야 함을 알 수 있다.

② 분산 함수의 모양이 달라질 때,



정규 분포를 바탕으로 한 가설이 깨지면  $1-\beta$  또한 바뀌게 됨.

(3rd, 4th moment까지 고려해야 함.)

# QnA

Q.  $y \leftarrow x$  으로 설명하는 게 부족. 여러 개의 variable로 설명되나?

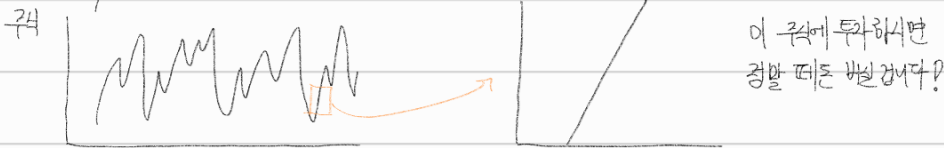
A. ① 설명이 부족하다.  $R^2$  or  $R^2_{adj}$

②  $x$ 으로 설명하고 남은 error 부분이  $x$ 에 따라 움직인다.

↳ Regression 작동 수  $E(\epsilon) = 0$  이 되어야 한다. ( $y = x\beta + \epsilon$ ) But expectation을 쓴다는 것 자체가 이미 error가 random이라는 것.

그런데  $x$ 에 따라 error가 결정되는 것 같다? error가 random이라는 기본적임 가정이 깨진 것.

③ Population의 특정 부분만 본 것 (Selection Bias)



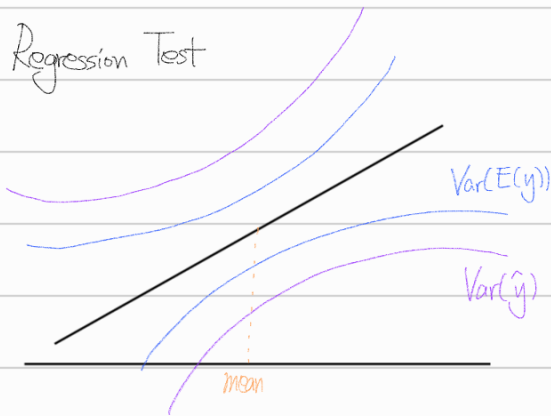
④ 시간이 지나면서 Data shape이 바뀌는 경우.

④ - 1. 단선이 평선, 분포만 바뀌는 경우. → 시간이라는 변수를 통해 control과거나 test를 따로 간별. \*대? 설명력이 떨어지지 않아서.

④ - 2. 분포가 달라짐. → 평균은 CLT 개념하에 t-test 간별.

평선 아니면? test 자체가 바뀌어야 함. 같은  $y$ 를  $x$ 으로 설명하고 싶은 상황.

## Regression Test



\*Variance는 왜 더 커질까?

Mean(expectation)에서 멀어지면서

$\epsilon$ 이 커진다. Regression은 평균에 대한 계산이니. 평균에서 멀어질수록 편차가 커져서 분산이 커짐.

$\hat{y}$ 를 예측하는 보라  $E(y)$ 를 예측하는  
검은 색은  $\epsilon$ 의 유무다.

$E(y|x) = E(x\beta|x) + E(\epsilon|x) = E(x\beta)$ : 검은 선은 위아. 검은 선만 올라간 나옴.

$Var(y|x) = Var(x\beta|x) + Var(\epsilon|x)$

$E(x\beta|x)$ 의 variance ↓ 매가 여기에 추가된 것.

: 단위  $\epsilon$ 의 유무로 분산에 차이가 생긴 것.

\*추가 설명 degrees of freedom

Data가 구어져 있을 때 아무 값이나 될 수 있는 특변. 그러는 평균이 고정되어 있기 때문에 Data들의 자유도도 I 감소 ex)  $\bar{x} = 2$   $x_1 = 1$   $x_2 = 2$   $x_3 = ?$

가능한 아무 값이나 될 수 있는가?