

Probability and Bayes' Review

Discretely ...

Marginal Distributions $p(A), p(B)$

Joint Distributions $p(A, B) = p(A)p(B)$ if A, B are independent events

Conditional Distributions $p(A|B), p(B|A)$

$$p(A) = \sum_B p(A, B), \quad p(B) = \sum_A p(A, B)$$

$$p(A|B) = \frac{p(A, B)}{p(B)} = \frac{p(A, B)}{\sum_A p(A, B)}, \quad p(B|A) = \frac{p(A, B)}{p(A)} = \frac{p(A, B)}{\sum_B p(A, B)}$$

Given $p(A|B), p(B)$, we can also find $p(B|A)$.

$$p(A|B) = \frac{p(A, B)}{p(B)} \Rightarrow p(A, B) = p(A|B)p(B)$$

Bayes' Rule

$$p(B|A) = \frac{p(A, B)}{p(A)} = \frac{p(A, B)}{\sum_B p(A, B)} = \frac{p(A|B)p(B)}{\sum_B p(A|B)p(B)}$$

Continuously ...

Joint Probability Density $p(x, y)$

Marginal Probability Density $p(x) = \int p(x, y) dy, \quad p(y) = \int p(x, y) dx$

$$p(x|y) = \frac{p(x, y)}{p(y)} = \frac{p(x, y)}{\int p(x, y) dx}, \quad p(y|x) = \frac{p(x, y)}{p(x)} = \frac{p(x, y)}{\int p(x, y) dy}$$

$$= \frac{p(y|x)p(x)}{\int p(y|x)p(x) dx} = \frac{p(x|y)p(y)}{\int p(x|y)p(y) dy}.$$

Bayes' Rule

\rightarrow we can find $p(x|y)$ given that we know $p(y|x)$ and $p(x)$!

Maximum Likelihood Estimation

or pdf for continuous case

parameter that we are trying to solve for our distribution.

Let X be a discrete r.v. with pmf p depending on parameter θ .

$L(\theta|x) = p_\theta(x) = P_\theta(X=x)$ is the likelihood function, given the outcome x of the r.v. X .

For Bernoulli r.v. X ,

$$L(\theta|x) = P_\theta(x) = \theta^x (1-\theta)^{1-x} \text{ s.t. } P_\theta(X=1) = \theta : P_\theta(X=0) = (1-\theta) \\ = 1 - P_\theta(X=1)$$

Suppose we have collected some data ... (Bernoulli coin tosses)
i.i.d.

$$\text{data} = \{x_1, x_2, x_3, \dots, x_N\}$$

Likelihood is: $L(\theta|\text{data}) = P_\theta(\text{data}) = \prod_{i=1}^N P_\theta(x_i) = \prod_{i=1}^N \theta^{x_i} (1-\theta)^{1-x_i}$

↳ Probability of observing the data that we observe assuming each coin toss was i.i.d

$$\text{e.g. } L(\theta|x_1=1, x_2=0, x_3=1) = \theta^1 (1-\theta)^0 \theta^1$$

We want to find the value of θ that maximizes $L(\theta|\text{data})$!

i.e. the value of θ that makes the data that we collect most probable.

To Maximize $L(\theta|\text{data})$, we want to take its derivative w.r.t. θ . st. $\frac{dL}{d\theta} = 0$

$$\hat{\theta} = \arg \max_{\theta} L(\theta). \quad \begin{array}{l} \text{Note that argmax is the operation that finds the argument} \\ \text{that gives the max value from function } L(\theta). \end{array}$$

Most of the time it is better to take the log of the likelihood before differentiating,

- and it usually leads to a simpler expression for the derivative that is easier to set to 0 and solve.
- since log is a monotonous function, i.e. whatever maximizes L also maximizes $\log L$.

$$l(\theta) = L(\theta) = \log \prod_{i=1}^N \theta^{x_i} (1-\theta)^{1-x_i} = \sum_{i=1}^N \{x_i \log \theta + (1-x_i) \log (1-\theta)\}.$$

$$\text{Set } \frac{dl}{d\theta} = \frac{1}{\theta} \sum_{i=1}^N x_i - \frac{1}{1-\theta} \sum_{i=1}^N (1-x_i) = 0$$

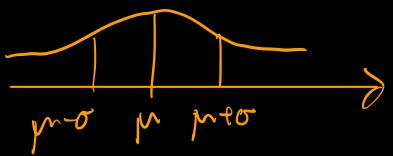
Solve for θ

$$(1/\theta) \sum_{i=1}^N x_i = (1/(1-\theta)) \sum_{i=1}^N (1-x_i) \\ \sum_{i=1}^N x_i - \theta \sum_{i=1}^N x_i = N\theta - \theta \sum_{i=1}^N x_i$$

Note that for Bernoulli dist.,
 $P(x=k) = \binom{n}{k} \theta^k (1-\theta)^{n-k}$

for Gaussian (Normal) r.v.,

data = $\{x_1, x_2, \dots, x_N\}$ i.i.d.



$$L(\theta | \text{data}) = \prod_{i=1}^N p_\theta(x_i), \text{ where } \theta = \{\mu, \sigma^2\}$$

$$= \frac{1}{N} \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x_i-\mu}{\sigma}\right)^2\right]$$

Find parameters that best describe data collected.

Goal :

These values will maximize the likelihood

$$\hat{\mu}, \hat{\sigma}^2 = \arg \max_{\mu, \sigma^2} L(\mu, \sigma^2 | \text{data})$$

$$\lambda(\mu, \sigma^2) = \log L(\mu, \sigma^2) = \log \frac{1}{N} \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x_i-\mu}{\sigma}\right)^2\right]$$

$$= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x_i-\mu}{\sigma}\right)^2\right] = \sum_{i=1}^N \left[\log \frac{1}{\sqrt{2\pi\sigma^2}} + \log e^{-\frac{1}{2}\left(\frac{x_i-\mu}{\sigma}\right)^2} \right]$$

$$= \sum_{i=1}^N \left[\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2}\left(\frac{x_i-\mu}{\sigma}\right)^2 \right] = -\underbrace{\frac{1}{2} \log(2\pi\sigma^2)}_{\text{Constant}} \sum_{i=1}^N 1 - \frac{1}{2} \sum_{i=1}^N \left(\frac{x_i-\mu}{\sigma}\right)^2$$

$$\frac{\partial \lambda}{\partial \mu} = \sum_{i=1}^N \left(\frac{x_i-\mu}{\sigma}\right) \frac{1}{\sigma}$$

$$\text{Set } \frac{\partial \lambda}{\partial \mu} = 0, \quad \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu) = 0.$$

$$\sum_{i=1}^N x_i - N\mu = 0.$$

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i = f(x_1, x_2, \dots, x_N).$$

$$l(\mu, \sigma^2) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2} \sum_{i=1}^N \left(\frac{x_i - \mu}{\sigma} \right)^2$$

$$l(\mu, v) = -\frac{N}{2} \log(2\pi v) - \frac{1}{2} \frac{1}{v} \sum_{i=1}^N (x_i - \mu)^2, \text{ where } v = \sigma^2$$

$$\frac{\partial l}{\partial v} = -\frac{N}{2} \frac{1}{v} - \frac{1}{2} (-1) \frac{1}{v^2} \sum_{i=1}^N (x_i - \mu)^2$$

$$\text{Set } \frac{\partial l}{\partial v} = 0, \quad N = \frac{1}{v} \sum_{i=1}^N (x_i - \mu)^2$$

$$v = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2 = g(x_1, x_2, \dots, x_N).$$

Our MLR estimates are also random variables!

If they are random variables, then we can ask:

- What is their distribution and expectation?

$$\begin{aligned} \hat{\mu} &\sim ? & E(\hat{\mu}) &=? \\ \hat{\sigma}^2 &\sim ? & E(\hat{\sigma}^2) &=? \end{aligned}$$

Can use characteristic functions to prove!

$$\hat{\mu} \sim N\left(\mu, \frac{\sigma^2}{N}\right)$$

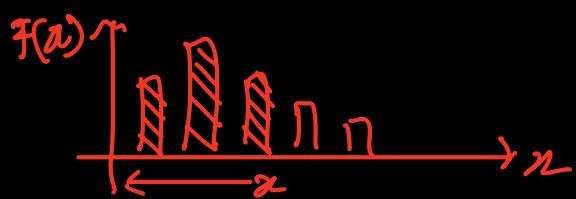
$$E(\hat{\mu}) = \mu$$

$$E(\hat{\sigma}^2) \neq \sigma^2; \quad E(\hat{\sigma}^2) = \frac{N-1}{N} \sigma^2 \rightarrow \sigma^2 \text{ as } N \rightarrow \infty.$$

See 'unbiased estimate of the covariance matrix' for full derivation

CDF and Percentiles

for Discrete rvs:

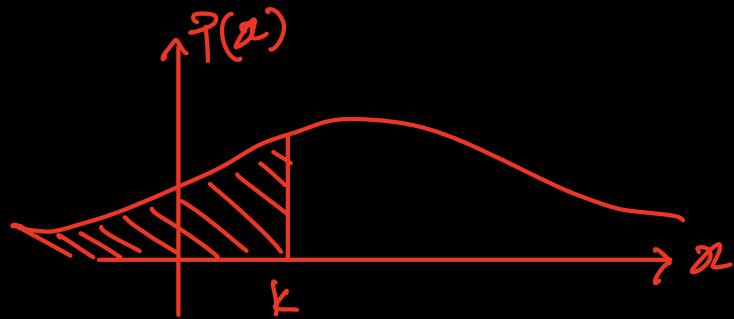


CDF: $F(x) = P(X \leq x) = \sum_{k=-\infty}^x p(k)$, where $p(k) = \text{Prob}(X=k)$

for Continuous rvs:

CDF: $F(k) = P(X \leq k) = \int_{-\infty}^k f(t) dt$ *t is a dummy variable, and disappears after integration.*

PDF: $f(x) = \frac{dF(x)}{dx}$



Inverse CDF (percentile function)

$F^{-1}(p)$ — What value of x would yield a CDF value of p ?

E.g. Heights — $\mu = 170 \text{ cm}$, $\sigma = 7 \text{ cm}$

$$F^{-1}(0.95, \mu = 170, \sigma = 7) \approx 181.5$$

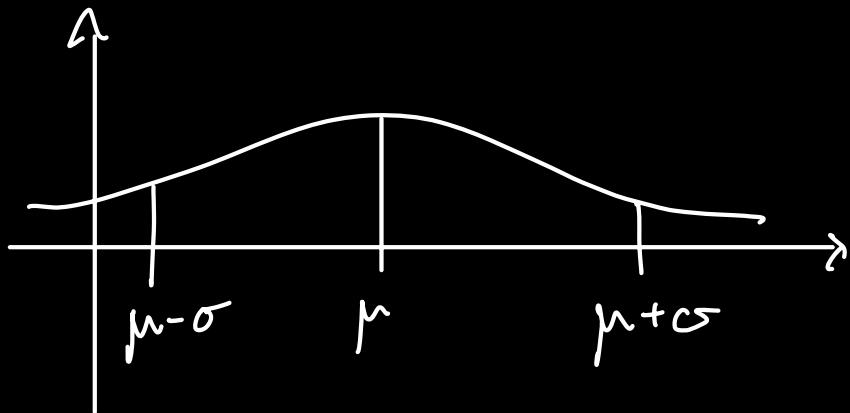
$$F(160, \mu = 170, \sigma = 7) \approx 0.08$$

Trend Mental A/B Testing

all about using point estimates to model a population's distribution

Using height as an example, suppose that height has a normal distribution,

$$X \sim N(\mu, \sigma^2)$$



Note that:

In probability theory, CLT establishes that in many situations for i.i.d samples, the standardized sample mean tends towards the standard normal distribution even if the original variables are not normally distributed.

We can say that we are more confident that the mean of our sample is representative of the entire population if:

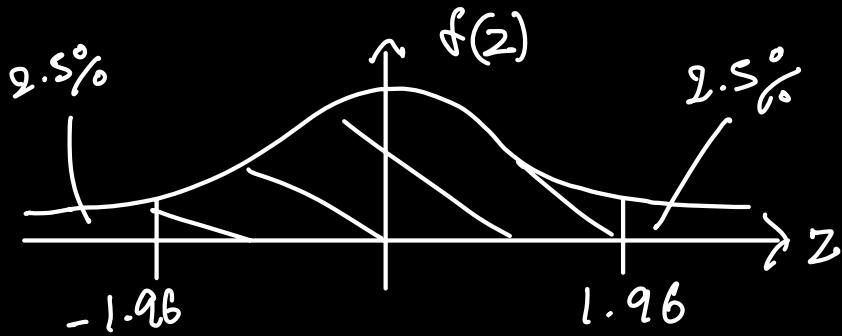
1. we have a large sample
2. the sample variance is small

$$\hat{\mu} \sim N\left(\mu, \frac{\sigma^2}{N}\right) \leftarrow \hat{E}(\hat{\mu}), \text{Var}(\hat{\mu})$$

In statistics, we typically want to standardize our data so we can easily calculate the probability of certain values occurring in our distribution, or to compare data sets with different means and standard deviations.

$$\text{i.e. } Z \sim N(0, 1) \rightarrow Z = \frac{x - \mu}{\sigma}$$

For a 95% confidence interval, the lower and upper endpoints are -1.96 and 1.96.



$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{z-\mu}{\sigma})^2}$$

standardizing a variable.

$$\frac{z-\mu}{\sigma} \sim N(0, 1).$$

If $X \sim N(\mu, \sigma^2)$, then $\hat{\mu} \sim N\left(\mu, \frac{\sigma^2}{N}\right)$

$$Z = \frac{\hat{\mu} - \mu}{\sigma/\sqrt{N}}$$

$$\left| \begin{array}{l} \hat{\mu} \sim N\left(\mu, \frac{\sigma^2}{N}\right) \\ \mathbb{E}(\hat{\mu}) = \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N x_i\right] = \frac{1}{N} \left[\sum_{i=1}^N \mathbb{E}(x_i) \right] = \mu \\ \text{Var}(\hat{\mu}) = \mathbb{E}\left[(\hat{\mu} - \mathbb{E}(\hat{\mu}))^2\right] = \mathbb{E}\left[\left(\frac{1}{N} \sum_{i=1}^N x_i - \mu\right)^2\right] \\ = \frac{1}{N^2} \mathbb{E}\left[\left(\sum_{i=1}^N x_i - N\mu\right)^2\right] = \frac{1}{N^2} \text{Var}\left(\sum_{i=1}^N x_i\right) = \frac{1}{N^2} \cdot N\sigma^2 = \frac{\sigma^2}{N} \end{array} \right.$$

Note that $\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$ is the unbiased estimator for sample variance (which we should use to replace σ^2 that is unknown)

for 95% C.I., $-1.96 \leq Z \leq 1.96$

$$-1.96 \leq \frac{\hat{\mu} - \mu}{\sigma/\sqrt{N}} \leq 1.96$$

$$-1.96 \left(\frac{\sigma}{\sqrt{N}} \right) - \hat{\mu} \leq -\mu \leq 1.96 \left(\frac{\sigma}{\sqrt{N}} \right) - \hat{\mu}$$

$$\hat{\mu} - 1.96 \left(\frac{\sigma}{\sqrt{N}} \right) \leq \mu \leq \hat{\mu} + 1.96 \left(\frac{\sigma}{\sqrt{N}} \right)$$

$$\hat{\mu} - 1.96 \frac{\sigma}{\sqrt{N}} \leq \mu \leq \hat{\mu} + 1.96 \frac{\sigma}{\sqrt{N}}$$

replace with $\hat{\sigma}$.

Essentially,

0.95 is the no. of times μ is contained in the 95% C.I. out of the total no. of experiments

For a $\gamma\%$ C.I., $\left[\hat{\mu} + \Phi^{-1}\left(\frac{1-\gamma}{2}\right) \frac{\sigma}{\sqrt{N}}, \hat{\mu} + \Phi^{-1}\left(1 - \frac{1-\gamma}{2}\right) \frac{\sigma}{\sqrt{N}} \right]$

However, we don't know the true variance σ^2 .

Consider a new standardized r.v.

(which is explicitly divided by the estimated sd)

$$t = \frac{\hat{\mu} - \mu}{\hat{\sigma}/\sqrt{N}} = \frac{\hat{\mu} - \mu}{\hat{\sigma}/\sqrt{N}} \times \frac{1/\sigma}{1/\sigma} = \left(\frac{\hat{\mu} - \mu}{\sigma/\sqrt{N}} \right) / \left(\frac{\hat{\sigma}}{\sigma} \right)$$

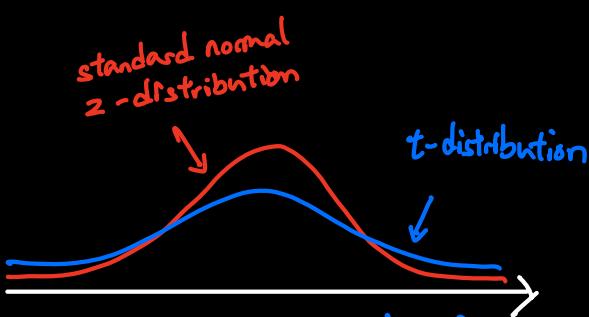
Consider denominator:

$$\left(\frac{\hat{\sigma}}{\sigma} \right)^2 = \frac{1}{N-1} \sum_{i=1}^N \left(\frac{x_i - \bar{x}}{\sigma} \right)^2$$

$$(N-1) \left(\frac{\hat{\sigma}}{\sigma} \right)^2 \sim \chi_{N-1}^2 \quad \text{chi-square distribution with } N-1 \text{ df}$$

Note: $Z \sim N(0,1) \rightarrow \sum_{i=1}^N (Z_i - \bar{Z})^2 \sim \chi_{N-1}^2$

$$t = \frac{z}{\sqrt{V/v}} \sim t_v$$



standard normal
 $\frac{z}{\sqrt{\text{chi-square/deg of freedom}}} \sim t_{\text{deg of freedom}}$

intuitively, t-distribution has fatter tails since we don't know σ , so we expect the confidence interval to be fatter

t-distribution is parameterized by $(N-1)$ degrees of freedom!

$$t_{\text{left}} = F^{-1}(0.025; \text{df} = N-1); \quad t_{\text{right}} = F^{-1}(0.975; \text{df} = N-1).$$

$$t_{\text{left}} \leq t \leq t_{\text{right}}$$

$$\hat{\mu} + t_{\text{left}} \frac{\hat{\sigma}}{\sqrt{N}} \leq \mu \leq \hat{\mu} + t_{\text{right}} \frac{\hat{\sigma}}{\sqrt{N}}$$

For large values of N , t-values \approx z-values
as $N \rightarrow \infty$, $\hat{\sigma} \rightarrow \sigma \Rightarrow t \rightarrow \text{Normal}$

Z-statistic

$Z = \frac{\hat{\mu} - \mu}{\sigma/\sqrt{N}}$, where $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$; $\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$
and σ is replaced with estimator $\hat{\sigma}$.

$$Z_{left} = \Phi^{-1}(0.025) ; Z_{right} = \Phi^{-1}(0.975)$$

$$\text{lower} = \hat{\mu} + Z_{left} \cdot \frac{\hat{\sigma}}{\sqrt{N}} ; \text{upper} = \hat{\mu} + Z_{right} \cdot \frac{\hat{\sigma}}{\sqrt{N}}$$

t-statistic

$$t = \frac{\hat{\mu} - \mu}{\hat{\sigma}/\sqrt{N}} = \frac{\hat{\mu} - \mu}{\hat{\sigma}/\sqrt{N}} \times \left(\frac{1/\sigma}{1/\hat{\sigma}} \right) = \frac{\hat{\mu} - \mu}{\sigma/\sqrt{N}} \left(\frac{\hat{\sigma}}{\sigma} \right)$$

$$\Rightarrow t = \frac{Z}{\sqrt{V/v}} \sim t_v, \text{ where } v = N-1 \text{ degrees of freedom}$$

$$t_{left} = \Phi^{-1}(0.025, v=N-1); t_{right} = \Phi^{-1}(0.975, v=N-1)$$

$$\text{lower} = \hat{\mu} + t_{left} \cdot \frac{\hat{\sigma}}{\sqrt{N}} ; \text{upper} = \hat{\mu} + t_{right} \cdot \frac{\hat{\sigma}}{\sqrt{N}}$$

Hypothesis Testing.

- Groups or Data :

1 → 1-Sample Test e.g. Avg. daily stock return

2 → 2-Sample Test Control vs Treatment Group e.g. Drug efficacy

- 1-sided vs 2-sided test

Consider 2-sample test for new drug.

(2-sided test)

$$H_0: \mu_1 = \mu_2 \quad H_0: \mu_1 - \mu_2 = 0$$

$$H_1: \mu_1 \neq \mu_2 \quad H_1: \mu_1 - \mu_2 \neq 0$$

(1-sided test)

$$H_0: \mu_1 \leq \mu_2$$

$$H_1: \mu_1 > \mu_2$$

→ Output :

Test statistic (does it fall in the rejection region?)

* P-value

tells us whether the difference is statistically significant
(given a significance threshold / level of significance)

Probability to reject the null when the null is true.

How likely that the data observed is to have occurred
under the Null Hypothesis.

→ Either reject the null hypothesis,
or fail to reject the null hypothesis

$$\hat{\mu} \sim N\left(\mu, \frac{\sigma^2}{N}\right)$$

Consider 1-Sample Test:

For 2-sided test: $\begin{cases} H_0: \mu = \mu_0 \\ H_1: \mu \neq \mu_0 \end{cases}$

Rewrite the null hypothesis

$$H_0: \hat{\mu} \sim N(\mu_0, \frac{\sigma^2}{N}) \longleftrightarrow H_0: Z = \frac{\hat{\mu} - \mu_0}{\sigma/\sqrt{N}} \sim N(0, 1)$$

To calculate P-value,

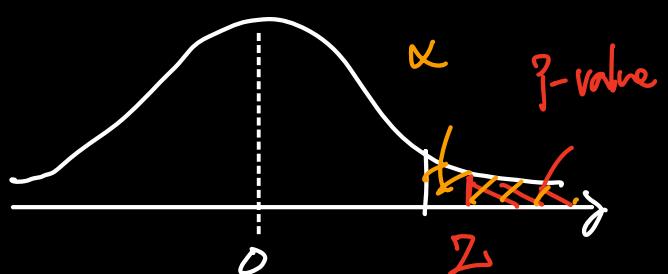
$$P_{\text{left}} = \Phi(-|Z|); P_{\text{right}} = 1 - \Phi(|Z|)$$

$$P = P_{\text{left}} + P_{\text{right}} = P_{\text{left}} \times 2 \quad (\text{since distribution is symmetric})$$

For 1-sided test: $\begin{cases} H_0: \mu \leq \mu_0 \\ H_1: \mu > \mu_0 \end{cases}$

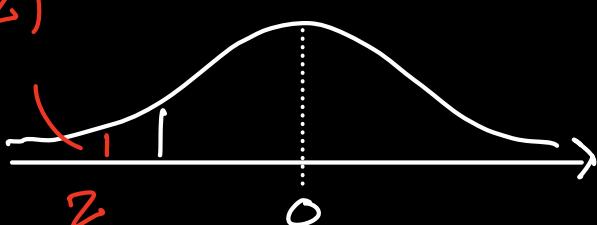
(for the greater than hypothesis)

$$P = 1 - \Phi(Z)$$



Similarly, for the lesser than hypothesis,

$$P = \Phi(Z)$$



For 2-sample test ,

For 2-sided test : $\left\{ \begin{array}{l} H_0: \mu_1 = \mu_2 \\ H_1: \mu_1 \neq \mu_2 \end{array} \right. \leftrightarrow \left\{ \begin{array}{l} H_0: \mu_1 - \mu_2 = 0 \\ H_1: \mu_1 - \mu_2 \neq 0 \end{array} \right.$

Let $y = \mu_1 - \mu_2$, $\left\{ \begin{array}{l} H_0: y = 0 \\ H_1: y \neq 0 \end{array} \right.$

$$\hat{y} = \hat{\mu}_1 - \hat{\mu}_2 = \frac{\sum_{i=1}^{N_1} x_i}{N_1} - \frac{\sum_{i=1}^{N_2} x_i}{N_2}$$

$$\begin{aligned} \text{Var}(\hat{y}) &= \text{Var}\left(\frac{x_1}{N_1}\right) + \dots + \text{Var}\left(\frac{x_{N_1}}{N_1}\right) + \text{Var}\left(\frac{x'_1}{N_2}\right) + \text{Var}\left(\frac{x'_{N_2}}{N_2}\right) \\ &= \frac{\sigma_1^2}{N_1^2} + \dots + \frac{\sigma_1^2}{N_1^2} + \frac{\sigma_2^2}{N_2^2} + \dots + \frac{\sigma_2^2}{N_2^2} \\ &= \frac{\sigma_1^2}{N_1} + \frac{\sigma_2^2}{N_2} \end{aligned}$$

$$\sigma_{\hat{Y}} = \sqrt{\frac{\sigma_1^2}{N_1} + \frac{\sigma_2^2}{N_2}}$$

$$Z = \frac{\hat{y} - \mu_0}{\sigma_{\hat{Y}}} = \frac{\hat{y} - \mu_0}{\sqrt{\frac{\sigma_1^2}{N_1} + \frac{\sigma_2^2}{N_2}}}$$

For P-value , 2-sided test :

$$P_{\text{right}} = 1 - \Phi(|Z|) ; P_{\text{left}} = \Phi(-|Z|) ; P = P_{\text{left}} + P_{\text{right}}$$

1-sided test :

$$(>) : P = 1 - \Phi(Z) \quad \text{Graph: A bell curve on a coordinate system with a vertical axis and a horizontal axis. A red vertical line is drawn at } Z \text{ on the positive side of the horizontal axis. The area under the curve to the right of this line is shaded red.}$$

$$(<) : P = \Phi(Z) \quad \text{Graph: A bell curve on a coordinate system with a vertical axis and a horizontal axis. A red vertical line is drawn at } Z \text{ on the negative side of the horizontal axis. The area under the curve to the left of this line is shaded red.}$$