

Question 5.

cf.

● In the materials taught before, we usually assume, in the statistical modeling, that the data follows a particular joint distribution which contains some unknown parameters of finite dimension. — e.g., normal with 3 parameters μ_x, μ_y, σ^2

- The statistical inferences, estimation and testing, are then based on a formulation of these parameters. — e.g., $\Delta = \mu_x - \mu_y$

Q: What if we do not have any knowledge about the particular form of the joint distribution of data?

Consider the problem of 2-sample comparison.

- Let Ω be the collection of all continuous distributions

F, G : any continuous distribution

- Only assume that $F, G \in \Omega$

Recall.
LNp. 2
case (a)

- Thus, the statistical model is:

parameter space/
model space

1st sample: $X_1, \dots, X_n \sim \text{i.i.d. from } F$
2nd sample: $Y_1, \dots, Y_m \sim \text{i.i.d. from } G$ } \Leftarrow independent (\square)

● This model contains parameters of infinitely many dimension because

$$\dim(\Omega) = \infty \quad (\text{why?})$$

pdf: $\int_{-\infty}^{\infty} f(x) dx = 1$

cdf: nondecreasing $F \rightarrow$

mgf: $M(t) \Rightarrow k\text{th moment} = M^{(k)}(0), k=1,2,\dots$
(chf)

check
LNp. 5
case (a)

● Under this model, a 2-sample comparison examines the null and alternative hypotheses:

$$H_0: F = G \quad \text{vs.} \quad H_A: F \neq G.$$

no difference

Definition 2 (nonparametric models and nonparametric methods)

such as normal, exponential, Poisson, ...

- Nonparametric models do not assume any particular distributional form.
Nonparametric models can be viewed as having infinitely many parameters.

The methods
can work for
data from any
distribution

($\xleftrightarrow{\text{cf.}}$ parametric models: parameters are of finite dimension)

- Statistical methods developed under nonparametric models are called nonparametric methods.

Q: What statistics (transformation) often appear in nonparametric method?

Review 3 (order statistics and ranks)

empirical
cdf
(TBp. 378)

- Let X_1, X_2, \dots, X_n be random variables. We sort the X_i 's and denote by $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ the order statistics. Using the notation,

$$X_{(1)} = \min(X_1, X_2, \dots, X_n) \quad \text{is the minimum,}$$

$$X_{(n)} = \max(X_1, X_2, \dots, X_n) \quad \text{is the maximum.}$$

percentile
median is a
function of
order statistics

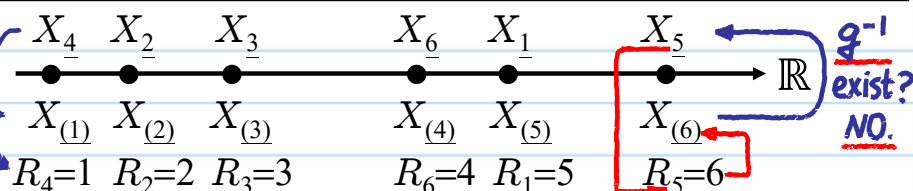
true
cdf
 y_n
 $X_{(1)} X_{(2)} \dots X_{(n)}$

- Let $R(X_1, X_2, \dots, X_n) = (R_1, R_2, \dots, R_n)$ such that $X_i = X_{(R_i)}$, $i = 1, \dots, n$. Then, (R_1, R_2, \dots, R_n) is called the ranks of X_1, X_2, \dots, X_n . Notice that

$$\# \text{ of observations } \leq X_i \rightarrow R_i = \sum_{j=1}^n \delta(X_i - X_j), \quad \text{where } \delta(t) = \begin{cases} 1, & \text{if } t \geq 0, \\ 0, & \text{if } t < 0. \end{cases}$$

complete
information
in X_1, \dots, X_n

data:
transformation g
order statistics:
ranks:



Theorem 8 (sufficient and complete statistics for nonparametric models)

$$\underline{T} = g(\underline{x})$$

$$\begin{matrix} x_1 \\ \vdots \\ x_n \end{matrix} \leftarrow \begin{matrix} x_{(1)} \\ \vdots \\ x_{(n)} \end{matrix}$$

for $x_1 < x_2 < \dots < x_n$ and zero, otherwise.

The proof of completeness is omitted (out of the scope of this course).

Note 7 (Some notes about order statistics and ranks)

- Order statistics and ranks are defined precisely, i.e., no ties, under the condition $P(X_i = X_j) = 0, i \neq j$ (**Note.** this condition holds when $X_1, \dots, X_n \sim$ i.i.d. from F and F is a continuous distribution).
- Under Ω , the dimension of data (i.e., n) cannot be reduced without losing the information about F ($\in \Omega$). \downarrow **combine information**
- Under 1-sample model, ranks + order statistics = complete data
- Order statistics are intuitive estimator of quantiles, e.g., median. \xrightarrow{e} **50% quantile**

- Ranks are invariant under any monotonic transformation of data, i.e.,

$$R(X_1, \dots, X_n) = R(H(X_1), \dots, H(X_n)),$$

if H is a monotone increasing function and

$$R(X_1, \dots, X_n) = (n + 1) - R(H(X_1), \dots, H(X_n)),$$

if H is a monotone decreasing function. ($\xleftrightarrow{\text{cf.}}$ z - or t -tests may change significantly under monotonic transformations of data).

🌀 Replacing the data by their ranks also has the effect of moderating the influence of outliers.

- Many nonparametric methods are based on order statistics and/or ranks.
- **Q:** Why are many nonparametric methods based on replacement of the data by ranks? What information of data are contained in their ranks?

(e) – **Recall.** Let X_1, \dots, X_n be i.i.d. from a continuous cdf F , and let $U_i = \overline{F}(X_i)$, $i = 1, \dots, n$. Then, U_1, \dots, U_n are i.i.d. from $U(0, 1)$.

(e) – **Recall.** If $U_1, \dots, U_n \sim \text{i.i.d.}$ $U(0,1)$, the pdf of the i th-order statistic $U_{(i)}$ is

$$f_{U(\underline{i})}(\underline{u}) = \frac{\underline{n}!}{(\underline{i}-1)!(\underline{n}-i)!} \frac{\underline{u}^{\underline{i}-1}}{\underline{u}} \frac{(1-\underline{u})^{\underline{n}-i}}{\underline{u}},$$

for $0 < u < 1$ and zero, otherwise.

Note that $E(U_{(i)}) = i/(n+1)$.

