

Types of variables

regression-type analysis (LNp.1-4)

response-explanatory distinction
(dependent-independent, response-predictor)

- ① values determined before collecting Y data, e.g., DOE
- ② if random, analysis conditional on the observed values of X_1, \dots, X_m .

➤ response variables Y : regarded as random. ← assign distribution.

check LNp.1-4 ➤ explanatory variables X_1, \dots, X_m : regarded as deterministic.

➤ causal relationship? not necessary.

- continuous-discrete distinction → Recall. continuous vs. discrete r.v.
(uncountable)(pdf) ← (countable)(pmf)

➤ whether a variable can take any values within an interval
(if yes ⇒ continuous; if no, only countable values ⇒ discrete)

e.g.,
height,
weight,
...

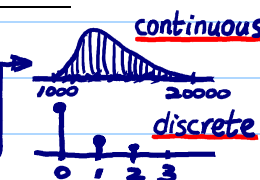
➤ Q: does continuous data really exist in the real world?

← No, because of the precision limitation of measuring instrument.

↳ any continuous data can be regarded as an interval data (LNp.2-3)

➤ a better distinction approach from the viewpoint of data analysis: according the number of values a variable can take within the range in which most data fall

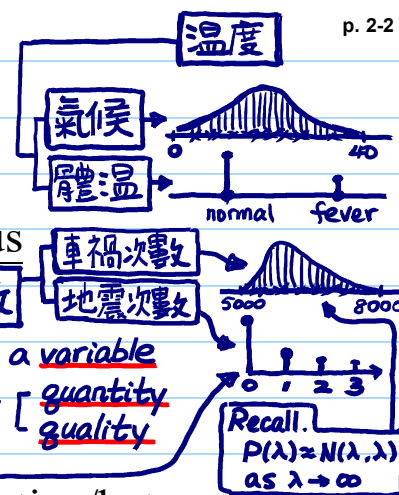
any pdf is an approximation to dense interval data



a trick:
check
histogram
↳ pdf
↳ pmf

variable taking lots of values (dense) ⇒ continuous; few values ⇒ discrete

■ Q: Should Poisson data (infinite possible values) be treated as discrete or continuous from this data analysis viewpoint?



- quantitative-qualitative distinction → 2 outcomes of a variable differ in quantity vs. quality
量性的 → 質性的
- continuous variable must be quantitative
- discrete variable could be quantitative/qualitative/between

- categorical (qualitative or between) variables can be further classified into:

➤ nominal variable: no natural ordering between categories (e.g., religious affiliation, mode of transportation, favorite type of music, ...)

→ or even if some order exists btwn categories, not interested in the effect of this order

what order?

r categories
denoted by
{1, 2, 3, ..., r}

- values that represent categories have no numeric meaning
- no value exist between categories
- analysis irrelevant to the order of listing the categories

In analysis, cannot arbitrarily change the order of categories as in nominal case

ordinal variable: there exist some ordering between categories (e.g., size of automobile, social class, political philosophy, patient condition, ...)

e.g., size of car - a (latent) continuous variable hidden behind

discrete (relatively few intervals)

compact $\xrightarrow{\text{distance}}$ large \geq middle $\xleftarrow{\text{distance}}$ large

compact $\xrightarrow{\text{distance}}$ large \geq middle $\xleftarrow{\text{distance}}$ large

exact distance between ordered categories are unknown

\rightarrow discrete interval or continuous interval

interval variable: categories are non-overlapping intervals (e.g., functional life length of television set, length of prison term, ...)

cf. continuous variable (an interval variable with many and very dense intervals)

0.5 1 2 3 4 5 6 7.5 10 length (years)

known

can have numerical distances between 2 categories

sometimes, possible to compare the ratio of 2 categories

Sometimes, it is the way that a variable is measured determined its classification, e.g., education:

nominal when measured as public/private school

$X_1 \equiv$ ordinal when measured as none/high school/bachelor/...

$X_2 \equiv$ interval variable when measured by # of year intervals

$X_1 = f(X_2)$, f : non-invertible \Rightarrow information lost in X_1

"0" must be well-defined

X_2 has higher "resolution" than X_1 . Q: Which is better in analysis? It depends.

more information \leftarrow

hierarchy of measurement scale: interval variable (highest) > ordinal > nominal (lowest)

$X_{ord} = f(X_{int})$
 $X_{nom} = g(X_{ord})$

ignore or disregard some information in the higher level

statistical methods for variables of one type can be used for variables at higher level, but not at lower levels

nominal variable \Rightarrow qualitative; interval variable \Rightarrow (close to) quantitative; ordinal variable \Rightarrow between (fuzzy)

choice of statistical method/model for different types of variables - a rough classification:

only response variables, no explanatory variable

- 1 response variable \Rightarrow uni-variate analysis
e.g., $Y_1, \dots, Y_k \text{ i.i.d. } N(\mu, \sigma^2) \text{ or } B(n, p) \text{ or } \text{Exp}(\lambda) \dots$
estimate or test the parameters.
- more than 1 response \Rightarrow multi-variate analysis
e.g., $Y_1, \dots, Y_k \text{ i.i.d. multivariate normal } (\underline{\mu}, \underline{\Sigma})$
information contained in the var-cov matrix

principal component
factor analysis
canonical correlation

both response and explanatory variables \rightarrow regression-type

	Resp.	Expl.
Cont.		
Disc.	counts (LNp2-8)	
Quan.		
Disc. Inte.		
Ordi.		
Nomi.		

check in LNp2-6

assign distribution

use base functions

categorical data

- response: (1) regarded as random variable; (2) modeling depends on the types of variable \uparrow assign distribution

assigning probabilities on counts:

- binomial
- multinomial
- Poisson
- negative binomial
- hypergeometric

- continuous & normal \Rightarrow linear model (LM)
- continuous but not normal (including exponential family, such as Weibull, gamma, ...) \Rightarrow generalized linear model

- discrete \Rightarrow generalized linear model (GLM) \leftarrow main focus of this course

- explanatory: (1) regarded as deterministic; (2) same treatment (sum of base functions multiplied by their coefficients) for any types of explanatory variables in modeling

Recall $\mathbf{X}\beta$ in LM

- LNp.1-23,24 \leftarrow □ quantitative: base functions like polynomial or other continuous transformations (e.g., log, exp, sin, cos, ...)

- LNp.1-25~32 \leftarrow □ qualitative: base function like dummy variables \leftarrow e.g., treatment or sum coding for nominal, Helmert coding for ordinal

Sufficient statistics of categorical responses

- Q: how to convert categorical observations (symbols or notations) into numerical and computable data without losing any important information? Q: what important information? meaningless if taking difference or ratio or numeric calculation on

Example 1.

➤ scenario

- observe $Y_{x,i} \in r$ categories $\{1, 2, \dots, r\}$, $i=1, \dots, k_x$. \leftarrow could be nominal, ordinal, or discrete interval
- category j observed with probability $p_{x,j}$, $j=1, \dots, r$, and \leftarrow a fixed known value

$$p_{x,1} + p_{x,2} + \dots + p_{x,r} = 1. \quad \leftarrow \text{Note. } Y_x \text{ is random}$$

- statistical modeling: $Y_{x,1}, Y_{x,2}, \dots, Y_{x,k_x}$ are independent and identically distributed from multinomial(1; $p_{x,1}, \dots, p_{x,r}$).

➤ sufficient statistics \leftarrow categorical responses, might not be meaningful if do calculation like

- the joint pmf of $Y_{x,1}, Y_{x,2}, \dots, Y_{x,k_x}$ is

$$\frac{N_{x,1}!}{p_{x,1}^{N_{x,1}}} \cdot \frac{N_{x,2}!}{p_{x,2}^{N_{x,2}}} \cdots \frac{N_{x,r}!}{p_{x,r}^{N_{x,r}}}$$

$$\sum_i Y_{x,i} / k_x$$

\leftarrow Linear model

* can be treated as a predictor with r levels

where $N_{x,j}$ = number of category j in the k_x trials, $j=1, \dots, r$, and $k_x = N_{x,1} + \dots + N_{x,r}$. \leftarrow count

\leftarrow numerical data, difference, sum, ratio, ... are meaningful

- by factorization thm, $(N_{x,1}, \dots, N_{x,r})$ are the sufficient

can understand how X influence $p_{x,i}$ in analysis.

statistics for the parameters $(p_{x,1}, \dots, p_{x,r})$ and $(N_{x,1}, \dots, N_{x,r}) \sim \text{multinomial}(k_x; p_{x,1}, \dots, p_{x,r})$. \leftarrow important information (dim=r-1)

- note: two categories and binomial is a special case of $r=2$.