

GEA1000 Quant. Reasoning with Data

AY 24/25 Sem 2 — github/omgeta

1. Studying Data

Population is the entire group of interest.

Population parameter is a population's numerical fact.

Census is an attempted survey of full population.

Sample is a subset of a population from a sampling frame.

Sample statistic is a numeric fact of the sample.

Estimates infer pop. parameters from sample statistics.

Selection bias is caused by flawed sampling frame or non-probability sampling. Non-response bias is caused by systematic exclusion of subjects by unwillingness.

Probability sampling:

- Simple random.**
- Systematic:** k^{th} subject of each size- r component.
- Stratified:** Divide into strata sharing a characteristic, then SRS within each stratum.
- Cluster:** Divide into natural clusters, then SRS including all subjects within selected clusters.

Non-probability sampling:

- Convenience sampling:** subjects chosen by convenience; selection bias.
- Volunteer sampling:** self-selected sample, usually with subjects off strong opinions; selection bias.

Study types:

- Experimental study:** observe dependent variable after direct manipulation of independent variable. Random treatment and control groups are similar. Shows cause-effect relationship.
- Observational study:** observe variable of interest without manipulation of variables. Shows association, not necessarily cause-effect.

Generalizability: frame size \geq population, probability sampling, large sample size and minimal bias.

2. Categorical Data Analysis

Categorical variables are ordinal (naturally ordered) or nominal (no natural order).

Rates

When variables A, B are not associated:

$$\text{i. } \text{rate}(A | B) = \text{rate}(A | B')$$

When variables A, B are associated:

$$\text{i. } \text{rate}(A | B) > \text{rate}(A | B') \text{ and } \text{rate}(A' | B') > \text{rate}(A' | B) \quad (+\text{ve})$$

$$\text{ii. } \text{rate}(A | B) < \text{rate}(A | B') \text{ and } \text{rate}(A' | B') < \text{rate}(A' | B) \quad (-\text{ve})$$

Symmetry Rules:

$$\text{i. } \text{rate}(A | B) > \text{rate}(A | B') \iff \text{rate}(B | A) > \text{rate}(B | A')$$

$$\text{ii. } \text{rate}(A | B) < \text{rate}(A | B') \iff \text{rate}(B | A) < \text{rate}(B | A')$$

$$\text{iii. } \text{rate}(A | B) = \text{rate}(A | B') \iff \text{rate}(B | A) = \text{rate}(B | A')$$

Basic Rule on Rates:

$$\text{i. } \text{rate}(A) \text{ lies between } \text{rate}(A | B) \text{ and } \text{rate}(A | B')$$

$$\text{ii. As } \text{rate}(B) \rightarrow 100\%, \text{rate}(A) \rightarrow \text{rate}(A | B)$$

$$\text{iii. } \text{rate}(B) = 50\% \implies \text{rate}(A) = \frac{1}{2}[\text{rate}(A | B) + \text{rate}(A | B')]$$

$$\text{iv. } \text{rate}(A | B) = \text{rate}(A | B') \implies \text{rate}(A) = \text{rate}(A | B) = \text{rate}(A | B')$$

Simpson's Paradox

Simpson's paradox is the observation that a trend appearing in majority of the groups of the data disappears/reverses when the groups are combined.

Confounders

Confounder is a third variable associated with both the independent and dependent variable being investigated. Randomised assignment can help to remove associations, removing the confounder in experimental studies.

3. Numerical Data Analysis

Numerical variables are discrete or continuous.

Summary Statistics

Mean, \bar{x} , is the average of variable x .

Mode is the most common element in variable x .

Q_1 , Median, Q_3 are the ordered 1st, 2nd, 3rd quarter element of variable x .

Sample variance, Var, and standard deviation, s_x , of variable x are given by:

$$\text{Var} = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$
$$s_x = \sqrt{\text{Var}}$$

Coefficient of variance, $\frac{s_x}{\bar{x}}$, measures spread relative to mean between different variables and has no units.

Median with $IQR = Q_3 - Q_1$ is preferred for asymmetrical distributions or when there are outliers.

Outliers satisfy one of the conditions:

$$\text{i. } x > Q_3 + 1.5 \times IQR$$

$$\text{ii. } x < Q_1 - 1.5 \times IQR$$

Univariate EDA

Histograms

Histograms show data distribution, are better at greater frequencies and represent data points better. Distributions with n peaks are called n -modal.

Unimodal distribution shapes can be:

- Symmetrical (mean = mode = median)
- Left-skewed (mean < mode < median)
- Right-skewed (mean > mode > median)

Bell distributions are symmetrical with spread:

- 68% of data within 1 S.D.
- 95% of data within 2 S.D.

Boxplots

Boxplots side-by-side help compare distributions of different data sets, and are better to identify outliers. They consist of minimum, Q_1 , median, Q_3 and maximum.

Boxplot shapes can be:

- Symmetrical (Q_1, Q_3 equidistant to median)
- Left-skewed (Q_1 closer to median)
- Right-skewed (Q_3 closer to median)

Boxplot spread for middle 50% is given by IQR .

Bivariate EDA

Deterministic relationships determine exactly a variable given the value of the other variable.

Association is a statistical relation describing average value of a variable given the value of the other variables

Correlation coefficient, r , is given by:

$$r = \frac{\text{Pop. covariance}}{\text{Pop. SD}_x \times \text{Pop. SD}_y} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \cdot \sum(y_i - \bar{y})^2}}$$

*unaffected by swapping x, y or adding/scaling by constant

Direction, form and magnitude can be summarized by r :

- $r > 0$ (+ve direction)
- $r < 0$ (-ve direction)
- $r = 0$ (Non-linear form)
- $0 < |r| < 0.3$ (Weak association)
- $0.3 < |r| < 0.7$ (Moderate association)
- $0.7 < |r| < 1$ (Strong association)

Linear Regression

Linear regression between variables believed to be linearly associated predicts the average value of the dependent variable given the independent variable.

Least squares regression line for predicting variable Y given X (and not vice versa) is given by:

$$Y = mX + b, \quad m = \frac{s_Y}{s_X}r$$

4. Statistical Inference

Probability of event E in sample space S , $P(E)$, is given by:

- $P(E) = \frac{|E|}{|S|}$, where $0 \leq P(E) \leq 1$
- $P(E') = 1 - P(E)$ (Complement)

Conditional probability of B given A is given by:

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

Mutually exclusive events A, B have special results:

- $P(A \cap B) = 0$ (Intersection)
- $P(A \cup B) = P(A) + P(B)$ (Union)
- $A \cup B = S$ (Total probability)
 $\implies P(C) = P(C | A)P(A) + P(C | B)P(B)$

Independent events A, B have special results:

- $P(A \cap B) = P(A) \cdot P(B)$ (Intersection)
- $P(A | B) = P(A)$ (Conditional)

$$\text{Sensitivity} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

$$\text{Specificity} = \frac{\text{True Negatives}}{\text{True Negatives} + \text{False Positives}}$$

Fallacies

Distribution fallacies:

- Ecological fallacy:** when we generalise group-level correlation to individuals.
- Atomistic fallacy:** when we generalise individual-level relations to groups.

Probability fallacies:

- Conjunction fallacy:** probability of two events occurring together is always less than of either event occurring alone.
- Base rate fallacy:** ignoring the base rate of an event when calculating its probability.

Relation between sample statistic and population parameter is given by:

Sample statistic = pop. parameter + bias + random error

Confidence Intervals

Confidence interval is a range of values likely to contain a population parameter at a certain confidence level.

Given a sample proportion p^* and sample size n , confidence interval for population proportion is given by:

$$p^* \pm z^* \times \sqrt{\frac{p^*(1-p^*)}{n}}$$

where z^* is the z -value for desired confidence level.

Given a sample mean \bar{x} , sample SD s_x and sample size n , confidence interval for population mean is given by:

$$\bar{x} \pm t^* \times \frac{s_x}{\sqrt{n}}$$

where t^* is the t -value for desired confidence level.

Hypothesis Testing

Hypothesis tests can be used for population proportion, population mean, and association, given a null hypothesis H_0 , alternative hypothesis H_1 , and significance value α . For hypothesis test on association, we take:

- H_0 there is no association
- H_1 : there is an association.

p -value can be defined as:

- Probability of obtaining a sample statistic as extreme or more extreme than the observed statistic, assuming H_0 is true.
- Smallest level of significance at which H_0 is rejected, assuming H_0 is true

where we reject H_0 in favour of H_1 when $p\text{-value} < \alpha$ or not reject H_0 (doesn't imply H_0 true) when $p\text{-value} \geq \alpha$