

Introduction

Population

The entire group we are concerned with.

Sample

A representative subset of the population.

Quantitative Data

Numerical data. Could be nominal (discrete or continuous), or ordinal (ordered).

Qualitative Data

Categorical data, e.g. colour, model.

Measures of Centrality

Mean

Given a set of n observations x_1, x_2, \dots, x_n , the **arithmetic mean** or **average** is defined as

$$\frac{1}{n} \sum_{i=1}^n x_i$$

The sample mean is denoted \bar{x} . The population mean is denoted μ .

Median

A drawback to the mean is that it can be misleading when the data is skewed. The **median** is the middle value of a set of n observations when arranged from largest to smallest.

If n is odd:

$$\text{median} = x^{(\frac{n+1}{2})}$$

or the $(n+1)/2$ th value of the sorted list. If n is even, the median is the :

$$\text{median} = \frac{x^{(\frac{n}{2})} + x^{(\frac{n}{2}+1)}}{2}$$

Mode

Given discrete data, the mode is defined as the most common value in a set of observations.

Measures of Dispersion

Dispersion refers to how much variation there is in a set of observations.

Range

The range is the difference between the maximum and minimum observation.

Variance

The variance is the average of the squared deviations from the mean.

- Given the observations x_1, x_2, \dots, x_N , from a population of size N with mean μ , the **population variance** is defined as

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

- Given the observations x_1, x_2, \dots, x_n , from a sample of size n with mean \bar{x} , the **sample variance** is defined as

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Standard Deviation

The standard deviation is the square root of the variance. The **population standard deviation** is defined as $\sigma = \sqrt{\sigma^2}$. The **sample standard deviation** is defined as $s = \sqrt{s^2}$.

Theorem 3.3.1 (Chebyshev's Theorem). *Given a set of n observations, at least*

$$1 - \frac{1}{k^2}$$

of them are within k standard deviations of the mean, where $k \geq 1$.

Theorem 3.3.2 (Empirical Rule). *If a histogram of the data is approximately unimodal and symmetric, then,*

- 68% of the data falls within **one** standard deviation of the mean
- 95% of the data falls within **two** standard deviations of the mean
- 99% of the data falls within **three** standard deviations of the mean

Often the standard deviation cannot be computed directly, but can be approximated using the Empirical rule. Here we assume that

$$\text{range} \approx 4s$$

so that

$$s = \frac{\text{range}}{4}.$$

Skew

The **skew** describes the asymmetry of the distribution. For a finite population of size N , the **population skew** is defined as

$$\frac{1}{N} \sum_{i=1}^N \left(\frac{x_i - \mu}{\sigma} \right)^3$$

For a sample of size n , the **sample skew** is defined as

$$\frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\left(\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \right)^3}$$

- When the skew is **positive**, the data is **right-skewed** and the “tail” of the distribution is **longer on the right**
- When the skew is **negative**, the data is **left-skewed** and the “tail” of the distribution is **longer on the left**

Measures of Rank

Z-Score

The Z-score is a unitless quantity and can be used to make comparisons of relative rank between members of a population.

$$Z = \frac{x - \mu}{\sigma} \quad \text{or} \quad \frac{x - \bar{x}}{s}$$

Quantiles

For a set of n observations, x_q is the q -th quantile, if $q\%$ of the observations are less than x_q .

Inter-Quartile Range

The inter-quartile range (IQR) is the difference between the 75th and 25th quantiles, or the range covered by the middle 50% of data.

Covariance and Correlation Coefficients

Covariance is the measure of the linear correlation between variables.

$$s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}.$$

Note that when $x = y$, the formula simplifies to the sample variance of x . The covariance has the following characteristics:

- $s_{xy} > 0$: As x increases, y also increases.
- $s_{xy} < 0$: As x increases, y decreases.
- $s_{xy} \approx 0$: No relationship between x and y .

Correlation Coefficient

$$-1 \leq r_{xy} = \frac{s_{xy}}{s_x s_y} \leq 1$$

Note that a correlation coefficient of 0 indicates **no linear relationship** between the variables, and not necessarily indicative of **no relationship**.

Regression and Least Squares

A linear relationship between two variables x and y is defined as $y = a + bx$. The least squares best fit determines the coefficients a and b that minimise the

sum of the squares of the residuals

$$b = r \frac{s_y}{s_x} = \frac{s_{xy}}{s_x^2}$$
$$a = \bar{y} - b\bar{x}.$$

Events and Probability

Multiplication Rule

For independent events A and B

$$\Pr(AB) = \Pr(A) \Pr(B).$$

For dependent events A and B

$$\Pr(AB) = \Pr(A|B) \Pr(B)$$

Addition Rule

For independent A and B

$$\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(AB).$$

If $AB = \emptyset$, then $\Pr(AB) = 0$, so that $\Pr(A \cup B) = \Pr(A) + \Pr(B)$.

De Morgan's Laws

$$\overline{A \cup B} = \bar{A} \bar{B}$$
$$\overline{AB} = \bar{A} \cup \bar{B}.$$

$$\Pr(A \cup B) = 1 - \Pr(\bar{A} \bar{B})$$
$$\Pr(AB) = 1 - \Pr(\bar{A} \cup \bar{B})$$

Bayes' Theorem

$$\Pr(A|B) = \frac{\Pr(B|A) \Pr(A)}{\Pr(B)}$$

Random Variables

Measurable variable whose value holds some uncertainty. An event is when a random variable assumes a certain value or range of values.

Probability Distribution

The probability distribution of a random variable X is a function that links all outcomes $x \in \Omega$ to the probability that they will occur $\Pr(X = x)$.

Probability Mass Function

$$\Pr(X = x) = p_x$$

Probability Density Function

$$\Pr(x_1 \leq X \leq x_2) = \int_{x_1}^{x_2} f(x) dx$$

Cumulative Distribution Function

Probability that a random variable is less than or equal to a particular realisation x .

$F(x)$ is a valid CDF if:

1. F is monotonically increasing and continuous
2. $\lim_{x \rightarrow -\infty} F(x) = 0$
3. $\lim_{x \rightarrow \infty} F(x) = 1$

$$\frac{dF(x)}{dx} = \frac{d}{dx} \int_{-\infty}^x f(u) du = f(x)$$

Complementary CDF (Survival Function)

$$\Pr(X > x) = 1 - \Pr(X \leq x) = 1 - F(x)$$

p-Quantiles

$$F(x) = \int_{-\infty}^x f(u) du = p$$

Special Quantiles

$$\text{Lower quartile } q_1: \quad p = \frac{1}{4}$$
$$\text{Median } m: \quad p = \frac{1}{2}$$
$$\text{Upper quartile } q_2: \quad p = \frac{3}{4}$$
$$\text{Interquartile range IQR: } \quad q_2 - q_1$$

Quantile Function

$$x = F^{-1}(p) = Q(p)$$

Expectation (Mean)

Expected value given an infinite number of observations. For $a < c < b$:

$$E(X) = - \int_a^c F(x) dx$$
$$+ \int_c^b (1 - F(x)) dx + c$$

Variance

Measure of spread of the distribution (average squared distance of each value from the mean).

$$\text{Var}(X) = \sigma^2 = E(X^2) - E(X)^2$$

Standard Deviation

$$\sigma = \sqrt{\text{Var}(X)}$$

Central Limit Theorem

For a sample of size n from any random probability distribution with expected value μ and variance σ^2 ,

$$\frac{\sqrt{n}(\bar{x} - \mu)}{\sigma} \xrightarrow{p} N(0, 1)$$

meaning that increasing the sample size will lead to a more normal distribution. In this case, a sample size of $n = 30$ is sufficient to approximate a normal distribution.

Standard Error

$$SE(\bar{x}) = \frac{\sigma^2}{n}$$

Sample Proportion

For a sample of size n let x be the number of members with a particular characteristic. The sample estimate of the population proportion p is

$$\hat{p} = \frac{x}{n}.$$

By assuming that the samples statistic x follows a binomial distribution with probability p and size n , then $E(x) = np$ and $\text{Var}(x) = np(1-p)$. Therefore the expectation is

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

and the standard error is

$$SE(\hat{p}) = \sqrt{\frac{p(1-p)}{n}}.$$

For the above to apply, we must assume that the sample proportion and size are sufficiently large. In general, if $np > 5$ and $n(1-p) > 5$, then we can assume that the sampling distribution of \hat{p} is approximately normal.

Assessing Normality

- Histograms: if the data is approximately normal, then the histogram will be approximately symmetric and unimodal.
- Boxplots: boxplots can be useful for showing outliers and skewness. Extreme clusters of an excessive number of outliers can be evidence of non-normality.
- Normal probability plots ($q-q$ plots): these plots are constructed by plotting the sorted data values against their Z -scores. If the data is approximately normal, then the points will lie approximately on a straight line.

Large Sample Estimation

Point Estimation

Method of Moments

The moments of a probability distribution are defined

$$\mu_n = E(X^n) = \int_{-\infty}^{\infty} x^n f(x) dx$$

where $f(x)$ is the probability density function of the distribution. Here $\mu_1 = \text{E}(X)$ and $\text{Var}(X) = \mu_2 - \mu_1^2$. Sample moments are defined similarly

$$m_n = \frac{1}{n} \sum_{i=1}^n x_i^n$$

where $\bar{x} = m_1$.

Method of Maximum Likelihood Estimation

Definition 9.1 (Likelihood function).

$$\mathcal{L}(\theta | \mathbf{x}) = \prod_{i=1}^n f(x_i)$$

Definition 9.2 (Maximum likelihood estimator).

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta | \mathbf{x}).$$

Definition 9.3 (Log-likelihood function). The **log-likelihood function** is defined as

$$\ell(\theta | \mathbf{x}) = \sum_{i=1}^n \log(f(x_i))$$

Due to the monotonicity of the log function, the maximum likelihood estimator is the same as the maximum log-likelihood estimator.

Definition 9.4 (Maximum log-likelihood estimator). The **maximum log-likelihood estimator** is defined as

$$\hat{\theta} = \arg \max_{\theta} \ell(\theta | \mathbf{x}).$$

Properties of Estimators

Definition 9.5 (Bias). The **bias** of an estimator is defined as the difference between the expected value of the estimator $\text{E}(\hat{\theta})$ and the true value of the parameter θ_0 .

$$\text{Bias}(\hat{\theta}) = \text{E}(\hat{\theta}) - \theta$$

An estimator $\hat{\theta}$ is **unbiased** if

$$\text{E}(\hat{\theta}) = \theta$$

so that the bias is zero.

We can also compare the variance of two estimators, to assess which one is more preferable. If the variance of the estimator is small, then the estimator is more precise. Given two estimators $\hat{\theta}_1$ and $\hat{\theta}_2$, we would choose $\hat{\theta}_1$ over $\hat{\theta}_2$ if

$$\text{Var}(\hat{\theta}_1) < \text{Var}(\hat{\theta}_2)$$

Definition 9.6 (Mean square error). Given data x_i with variance σ^2 , the estimators of $\theta = \text{E}\{X\}$ are selected such that they minimise the **mean square error**:

$$\begin{aligned} \text{MSE}(\hat{\theta}) &= \text{E}\left((\hat{\theta} - \theta)^2\right) \\ &= \text{Bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta}). \end{aligned}$$

This quantity is used to determine the **bias-variance trade-off** of an estimator. The **root mean square error** is defined as

$$\text{RMSE}(\hat{\theta}) = \sqrt{\text{MSE}(\hat{\theta})}.$$

Confidence Intervals

This interval ranges from the lower confidence limit (UCL) to the upper confidence limit (LCL)

$$L < \theta < U.$$

This interval has a **confidence coefficient** of $1 - \alpha$, or a **confidence level** of $(1 - \alpha)\%$. The confidence interval is defined as

$$CI_{1-\alpha} = \hat{\theta} \pm Z_{\alpha/2} \text{SE}(\hat{\theta})$$

Confidence Interval for the Mean

$$CI_{1-\alpha} = \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

where $\text{SE}(\bar{x}) = \frac{\sigma}{\sqrt{n}}$.

Confidence Interval for the Proportion

Given the sample size n and sample proportion \hat{p} ,

$$\hat{p} \sim N\left(p, \frac{p(1-p)}{n}\right).$$

The confidence interval for the population proportion is

$$CI_{1-\alpha} = \hat{p} \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

where the standard error is given by

$$\text{SE}(\hat{p}) = \sqrt{\frac{p(1-p)}{n}}.$$

with the approximation $p = \hat{p}$. Note that $n\hat{p} > 5$ and $n(1-\hat{p}) > 5$ are required for the approximation to be valid.

Confidence Interval for the Difference of Two Means

$$CI_{1-\alpha} = \bar{x}_1 - \bar{x}_2 \pm Z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}.$$

If the two populations follow a normal distribution, then the sampling distribution is exactly normal. If the two populations are not normal, then the sampling distribution is approximately normal, for $n_1 > 30$ and $n_2 > 30$.

Confidence Interval for the Difference of Two Proportions

$$CI_{1-\alpha} = \hat{p}_1 - \hat{p}_2 \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

Note that the following constraints must be satisfied:

- $n_1\hat{p}_1 > 5$
- $n_1(1-\hat{p}_1) > 5$
- $n_2\hat{p}_2 > 5$
- $n_2(1-\hat{p}_2) > 5$

Hypothesis Testing

Hypothesis Testing for the Population Mean

Given the sample statistic \bar{x} ,

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

the test statistic is defined

$$T(\mathbf{x}) = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \sim N(0, 1).$$

Hypothesis Testing for the Population Proportion

Given the sample statistic \hat{p} ,

$$\hat{p} \sim N\left(p, \frac{p(1-p)}{n}\right)$$

for $n\hat{p} > 5$ and $n(1-\hat{p}) > 5$, the test statistic is defined

$$T(\mathbf{x}) = \frac{\sqrt{n}(\hat{p} - p_0)}{\sqrt{p_0(1-p_0)}}.$$

Hypothesis Testing with Differences

The rejection regions for the difference between two parameters is defined:

Null Hypothesis H_0	Rejection Region
$\theta_1 - \theta_0 = 0$	$ T(\mathbf{x}) > Z_{\alpha/2}$
$\theta_1 - \theta_2 \leq 0$	$T(\mathbf{x}) > Z_{\alpha}$
$\theta_1 - \theta_2 \geq 0$	$T(\mathbf{x}) < -Z_{\alpha}$

Hypothesis Testing for the Difference in Population Means

The point estimator of $\mu_1 - \mu_2$ is given by

$$\bar{x}_1 - \bar{x}_2$$

and the standard error is given by

$$\text{SE}_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}.$$

The test statistic is defined

$$T(\mathbf{x}) = \frac{(\bar{x}_1 - \bar{x}_2) - \Delta_0}{\text{SE}_{\bar{x}_1 - \bar{x}_2}}.$$

where $\Delta_0 = \mu_1 - \mu_2$ is the hypothesized difference between the two population means.

Hypothesis Testing for the Difference in Population Proportions

The point estimator of the difference in proportions where $p_1 = p_2$ is given by

$$\hat{p}_1 - \hat{p}_2$$

and the standard error is defined

$$SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{p_0(1 - p_0) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

where

$$p_0 = \frac{x_1 + x_2}{n_1 + n_2}$$
$$p_0 = \frac{\hat{p}_1 n_1 + \hat{p}_2 n_2}{n_1 + n_2}.$$

so that $p_0 = p_1 = p_2$. The resulting test statistic is defined:

$$T(\mathbf{x}) = \frac{(\hat{p}_1 - \hat{p}_2)}{SE_{\hat{p}_1 - \hat{p}_2}}.$$

When the hypothesised difference is not 0, i.e., $p_1 - p_2 = \Delta_0$, the test statistic is defined:

$$T(\mathbf{x}) = \frac{(\hat{p}_1 - \hat{p}_2) - \Delta_0}{\sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}}.$$

Significance of Results

When interpreting the results from a test statistic, the test can only be used to reject the null hypothesis.

Small Sample Inference

Student's t-distribution:

$$T(\mathbf{x}) \sim t_\nu$$

where the degrees of freedom ν is equal following: to $n - 1$.

$$E(X) = 0$$
$$\text{Var}(X) = \frac{\nu}{\nu - 2}$$

Inferencing

$$T(\mathbf{x}) = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim t_{\nu, \alpha/2}.$$

Hypothesis Testing for the Population Mean

$$T(\mathbf{x}) = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim t_{\nu, \alpha/2}.$$

Hypothesis Testing for the Difference in Population Means

$$T(\mathbf{x}) = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim t_{\nu, \alpha/2}.$$

If the sample variances s_1^2 and s_2^2 are not equal, then we need to determine the common or *pooled* variance s_p^2 .

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{\nu}.$$

where $\nu = n_1 + n_2 - 2$ for the two-sample *t*-test. This results in the following test statistic:

$$T(\mathbf{x}) = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}.$$

The population variances between two samples vary *greatly*, if they satisfy the

$$\frac{\max(s_1^2, s_2^2)}{\min(s_1^2, s_2^2)} > 3.$$

when this is the case, we must modify the test statistic to account for the different variances:

$$T(\mathbf{x}) = \frac{(\bar{x}_1 - \bar{x}_2) - \Delta_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

noting that Δ_0 is typically zero. The degrees of freedom are given by

$$\nu = \left\lfloor \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2}{\frac{(s_1^2/n_1)^2}{n_1 - 1} + \frac{(s_2^2/n_2)^2}{n_2 - 1}} \right\rfloor$$

where the value is truncated (towards zero). Note that the above only applies for independent samples.

Decision	H_0	$\neg H_0$
Reject H_0	α	$1 - \beta$
Fail to reject	$1 - \alpha$	β

H_0	Rejection Region R
$\theta = \theta_0$	$ T(\mathbf{x}) > Z_{\alpha/2}$
$\theta \leq \theta_0$	$T(\mathbf{x}) > Z_\alpha$
$\theta \geq \theta_0$	$T(\mathbf{x}) < -Z_\alpha$

	Discrete	Continuous
$E(X)$	$\sum_{\Omega} x p_x$	$\int_{\Omega} x f(x) \, dx$
$E(g(X))$	$\sum_{\Omega} g(x) p_x$	$\int_{\Omega} g(x) f(x) \, dx$
$\text{Var}(X)$	$\sum_{\Omega} (x - \mu)^2 p_x$	$\int_{\Omega} (x - \mu)^2 f(x) \, dx$

Distribution	Restrictions	PMF	CDF	$E(X)$	$\text{Var}(X)$
$X \sim \text{Uniform}(a, b)$	$x \in \{a, \dots, b\}$	$\frac{1}{b-a+1}$	$\frac{x-a+1}{b-a+1}$	$\frac{a+b}{2}$	$\frac{(b-a+1)^2-1}{12}$
$X \sim \text{Bernoulli}(p)$	$p \in [0, 1], x \in \{0, 1\}$	$p^x (1-p)^{1-x}$	$1-p$	p	$p(1-p)$
$X \sim \text{Binomial}(n, p)$	$x \in \{0, \dots, n\}$	$\binom{n}{x} p^x (1-p)^{n-x}$	$\sum_{u=0}^x \binom{n}{u} p^u (1-p)^{n-u}$	np	$np(1-p)$
$N \sim \text{Poisson}(\lambda)$	$n \geq 0$	$\frac{\lambda^n e^{-\lambda}}{n!}$	$e^{-\lambda} \sum_{u=0}^n \frac{\lambda^u}{u!}$	λ	λ

Distribution	Restrictions	PDF	CDF	$E(X)$	$\text{Var}(X)$
$X \sim \text{Uniform}(a, b)$	$a < x < b$	$\frac{1}{b-a}$	$\frac{x-a}{b-a}$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
$T \sim \text{Exp}(\eta)$	$t > 0$	$\eta e^{-\eta t}$	$1 - e^{-\eta t}$	$1/\eta$	$1/\eta^2$
$X \sim N(\mu, \sigma^2)$	$x \in \{0, \dots, n\}$	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	$\frac{1}{2} \left(1 + \text{erf} \left(\frac{x-\mu}{\sigma\sqrt{2}} \right) \right)$	μ	σ^2