

NORMAL APPROXIMATION TO BINOMIAL
 • When $n \rightarrow \infty$, p remains constant and near 0. Use when $np > 5$ & $n(1-p) > 5$. If $X \sim \text{Bin}(n, p)$, $Z = \frac{X - np}{\sqrt{np}} \sim N(0, 1)$

• Do **CONTINUITY CORRECTION BEFORE NORMALISATION:**
LHS: $\leq \Rightarrow a - 0.5 < \Rightarrow a + 0.5$
RHS: $< \Rightarrow a - 0.5 \leq \Rightarrow a + 0.5$
 $P(X = k) \approx P(k - 0.5 < X < k + 0.5)$
 $P(a < X < b) \approx P(a + 0.5 < X < b - 0.5)$
 $P(X \leq c) = P(0 \leq X \leq c) \approx P(-0.5 < X < c + 0.5)$
 $P(X > c) = P(c < X \leq n) \approx P(c + 0.5 < X < n + 0.5)$

SIMPLE RANDOM SAMPLE Every subset of n observations (members) of the population has same prob of being selected.
 • #possible sample sizes of n from N population = N choose n .

INFINITE POPULATION Equals to sampling from **finite** population with replacement. Prob within a distribution varies, but does not vary across samples $X_1 \dots X_n$ of same distribution.
 • Joint Probability Function: $f_{x_1 \dots x_n}(x_1, \dots, x_n) = f_{x_1}(x_1) \dots f_{x_n}(x_n)$

STATISTICS
Statistic (RANDOM VARIABLE) = function of $(X_1 \dots X_n)$ [random sample of n observations]. Replace with observed \rightarrow Realisation.

SAMPLE MEAN $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ [random variable]
 $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad E(\bar{X}) = E(X)$

SAMPLE VARIANCE $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ [random var]
 $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad E(S^2) = \sigma^2$

MEAN AND VARIANCE OF \bar{X}
 Sampling distribution of the sample mean \bar{X} from a population with mean μ_X and variance σ_X^2 . Note, $\sigma_{\bar{X}}$ is the variance for the sampling distribution for the sample mean, not of the population.

$\mu_{\bar{X}} = E(\bar{X}) = \mu_X$ [valid estimator for population mean]
 $\sigma_{\bar{X}}^2 = V(\bar{X}) = \frac{\sigma_X^2}{n}$

For an inf population, $n \rightarrow \infty$, variance becomes smaller and smaller, accuracy of \bar{X} as an estimator improves

STANDARD ERROR Describes spread of sampling distribution. Is $\sigma_{\bar{X}}$, the standard deviation of sampling distribution of \bar{X} .

Standard error of \bar{X} shows how much \bar{x} differs among samples.

CENTRAL LIMIT THEOREM
 Law of large numbers If $X_1 \dots X_n$ are iid RV with same μ and σ^2 ,
 As $n \rightarrow \infty$, $P(|\bar{X} - \mu| > \epsilon) \rightarrow 0$, $\epsilon \in \mathbb{R}$ i.e. $\bar{X} \rightarrow \mu$
 • When $P(|\bar{X} - \mu| > \frac{\sigma}{\sqrt{n}}) > c$ then $P(Z > 1) > c$

CLT: $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \frac{\bar{X} - \mu}{\sigma_{\bar{X}}} \rightarrow Z \sim N(0, 1)$ (Approximate for large $n \geq 30$, non-normal sample. Exact for normal sample)
 $\bar{X} \sim N(\mu, \frac{\sigma^2}{n}) \lim_{n \rightarrow \infty} P(-\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq x) = \Phi(x)$
 $P(a \leq \bar{X} \leq b) = P(\frac{a - \mu}{\sigma/\sqrt{n}} \leq \bar{X} \leq \frac{b - \mu}{\sigma/\sqrt{n}})$

χ^2 (n) distribution, n degrees of freedom
 χ^2 is a χ^2 distribution with 1 degree of freedom. If Z_1, \dots, Z_n are iid standard normal RV, then $Z_1^2 + \dots + Z_n^2 \sim \chi^2(n)$. All χ^2 density functions have long right tail, determined by degree of freedom.

Expectation: Variance: $2n$ n large $\rightarrow \chi^2 \approx \chi^2(n, m)$
 • Both $Y_1 \sim \chi^2(n)$, $Y_2 \sim \chi^2(m)$ independent $\rightarrow Y_1 + Y_2 \sim \chi^2(n+m)$

• $\chi^2(n, \alpha) = Y - Y(n)$, $P(Y > y(n; \alpha)) = \alpha$

If $X \sim N(\mu, \sigma^2)$ $\forall i$, use $\frac{(x_i - \mu)^2}{\sigma^2} \sim \chi^2(n-1)$ for n iid RV with

$E(X) = \mu$, $\text{Var}(X) = \sigma^2$, sample variance $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$.

• If S^2 is variance of a random sample size n taken from a normal population w variance σ^2 : $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$

---- t distribution, t(n), with n degrees of freedom ----

Independent $Z \sim N(0, 1)$ and $U \sim \chi^2(n) \Rightarrow T = \frac{Z}{\sqrt{U/n}} \sim t(n)$

• $\frac{U}{n} = \frac{Z^2 + U - \chi^2(n)}{n}$ and $E(Z^2) = 1$ when $n \rightarrow \infty \rightarrow T \sim N(0, 1)$ if $n \geq 30$

• Expectation $E(T) = 0$ Variance $V(T) = \frac{n}{n-2} > 2$

• Right tail probability: $P(T > t_{n, \alpha}) = \alpha$ [same shape as N]

• $X_1 \dots X_n$ iid normal RV with mean μ , var $\sigma^2 \Rightarrow \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim t(n-1)$

-- F distribution, $F(m, n)$, with (m, n) degrees of freedom --

Independent $U \sim \chi^2(m)$ and $V \sim \chi^2(n) \Rightarrow F = \frac{U/m}{V/n} \sim F(m, n)$

• $E(F) = \frac{n}{n-2} > 2 \quad \text{Var}(F) = \frac{2n^2(m+n-2)}{m(n-2)^2(n-4)}$, $n > 4$

• $F = \frac{U/(n-1)}{V/(n-1)} = \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} \sim F(n-1, n-2, 1)$

• $F \sim (m, n) \Rightarrow \frac{1}{F} \sim (n, m) \quad F(m, n; 1 - \alpha) = \frac{1}{F(n, m; \alpha)}$

• Upper tail probability of $F(m, n; \alpha)$ s.t. $P(F > F(m, n; \alpha)) = \alpha$

UNBIASED ESTIMATOR

$E(\hat{\theta}) = \theta$ where $\hat{\theta}$ is an estimator of parameter θ , is a RV based on sample, and has mean value equal to true value of parameter.

Confidence Intervals / Test Statistics: Population Mean: • $100(1-\alpha)\%$ confidence interval formulas for population mean μ ,
 • test statistics for the (null) hypothesis: $H_0: \mu = \mu_0$

!! Remember to **SQUARE ROOT VARIANCE σ^2 or s^2 TO GET σ or s !!**

	Population	σ	n	Test Statistic	E	n for desired ϵ, α	Confidence Interval
I	Normal	Known	Any	$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim N(0, 1)$	$z_{\alpha/2} \cdot (\sigma/\sqrt{n})$	$\left(\frac{z_{\alpha/2} \cdot \sigma}{E_0} \right)^2$	$\bar{X} \pm z_{\alpha/2} \cdot (\sigma/\sqrt{n})$
II	Any	Known	Large	$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim N(0, 1)$	$z_{\alpha/2}$	$t < -t_{n-1, \alpha/2}$	$\bar{X} \pm z_{\alpha/2} \cdot (\sigma/\sqrt{n})$
III	Normal	Unknown	Small	$T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \sim t_{n-1}$ n-1 degrees of freedom	$t_{n-1, \alpha/2} \cdot \left(\frac{s}{\sqrt{n}} \right)$	$\left(\frac{t_{n-1, \alpha/2} \cdot s}{E_0} \right)^2$	$\bar{X} \pm t_{n-1, \alpha/2} \cdot (s/\sqrt{n})$ n-1 degrees of freedom
IV	Any	Unknown	Large $n \geq 30$	$Z = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \sim N(0, 1)$	$z_{\alpha/2} \cdot (s/\sqrt{n})$	$\left(\frac{z_{\alpha/2} \cdot s}{E_0} \right)^2$	$\bar{X} \pm z_{\alpha/2} \cdot (s/\sqrt{n})$

HYPOTHESIS TESTING & REJECTION OF H_0

Can use Z-Test in GC Can use T-Test in GC

Given $H_0: \mu = \mu_0$		p-value (note: z = computed test statistic)	Rejection region (Reject H_0) (NOTE α , or $\alpha/2$)
Two-sided	$H_1: \mu \neq \mu_0$	$P(Z > z) = 2P(Z > z) = 2P(Z < - z)$	$z < -z_{\alpha/2}$ or $z > z_{\alpha/2}$ or $t < -t_{n-1, \alpha/2}$ or $t > t_{n-1, \alpha/2}$
One-sided	$H_1: \mu < \mu_0$	$P(Z < - z)$ (only left tail)	$z < -z_{\alpha}$ (α , not $\alpha/2$!) n-1 degrees of freedom
	$H_1: \mu > \mu_0$	$P(Z > z)$ (only right tail)	$z > z_{\alpha}$ (α , not $\alpha/2$!) n-1 degrees of freedom

p-value & significance level α or CONFIDENCE INTERVAL \Rightarrow Reject H_0 ?

p-value $< \alpha$ Confidence interval does NOT contain μ_0 \Rightarrow Reject H_0

p-value $\geq \alpha$ Confidence interval contains μ_0 \Rightarrow DO NOT reject H_0

MAXIMUM ERROR OF ESTIMATE $E = z_{\alpha/2} \cdot (\sigma/\sqrt{n})$

Minimum sample size so that with

$$\text{probability } 1 - \alpha, \text{ the error is at most } E_0: \quad n \geq \left(\frac{z_{\alpha/2} \cdot \sigma}{E_0} \right)^2$$

$(1 - \alpha) \text{ CI} \Rightarrow z_{\alpha/2} \text{ distr} \rightarrow \text{invNorm}(1 - \frac{\alpha}{2}, \dots)$

$\alpha/2$ or α

90% CI, $\alpha = 0.10$: $z_{0.05} = 1.645$ $\text{invNorm}(0.95..)$

95% CI, $\alpha = 0.05$: $z_{0.025} = 1.96$ $\text{invNorm}(0.975..)$

96% CI, $\alpha = 0.04$: $z_{0.02} = 2.05$ $\text{invNorm}(0.98..)$

98% CI, $\alpha = 0.02$: $z_{0.01} = 2.326$ $\text{invNorm}(0.99..)$

99% CI, $\alpha = 0.01$: $z_{0.005} = 2.576$ $\text{invNorm}(0.995..)$

0.995 quantile of standard normal = $z_{0.005} = 2.576$

$(1 - \alpha) \text{ CI} \Rightarrow t_{n-1, \alpha/2} \text{ distr} \rightarrow \text{invT}(1 - \frac{\alpha}{2}, \text{degree of freedom})$

$\alpha/2$

90% CI: $\text{invT}(0.95, n)$ for $t_{n-1, 0.05} \quad ** 0.95 = (1 - 0.10/2)$

95% CI, $\alpha = 0.05 \Rightarrow \text{invT}(0.975, n-1), t_{10, 0.025} = 2.228$

98% CI, $\alpha = 0.02 \Rightarrow \text{invT}(0.99, n-1), t_{9, 0.01} = 2.821$

99% CI, $\alpha = 0.01 \Rightarrow \text{invT}(0.995, n-1), t_{8, 0.005} = 4.604$

$t_{n-1, \alpha} \text{ distr} \rightarrow \text{invT}(1 - \alpha, \text{degree of freedom}) \quad \alpha$

95% CI, $\alpha = 0.05 \Rightarrow \text{invT}(0.95, n-1), t_{7, 0.05} = 1.895$

CONFIDENCE INTERVAL

• Interval estimator: Rule for calculating, from the sample, an interval (ab) where you are fairly certain the parameter lies in.

• Degree of confidence / Confidence Level = $1 - \alpha$

• $P(a < \mu < b) = 1 - \alpha$ (a, b) is the $(1 - \alpha)$ confidence interval

• CI at the 1 - α level of confidence is the $(1 - \alpha)$ confidence that a population parameter is found between the interval. Also the % of samples that will have population parameter within interval.

• $(1 - \alpha) \text{ CI} = \bar{X} \pm E$ where E = maximum error of estimate

→ Since confidence interval (a, b) is $\bar{X} \pm E$,
 $\text{Reject } H_0 \Leftrightarrow \bar{X} \notin (a, b)$ **DO NOT reject $H_0 \Leftrightarrow \mu_0 \in (a, b)$**

Confidence intervals can be used to perform two-sided tests.

• $\bar{X} \pm E$ has probability $(1 - \alpha)$ of containing μ

POOLED ESTIMATOR S_p^2

• For small samples with equal variance

• Assume equal variance when $\frac{1}{2} \leq \frac{s_1^2}{s_2^2} \leq 2$

• If s is unknown → Estimate it. Under the equal variance assumption, s_1^2 and s_2^2 are unbiased estimators of σ^2 .

• POOLED ESTIMATOR, $S_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}$

PAIRED DATA (X_i, Y_i)

• Matched pairs (X_i, Y_i) , where $X_1 \dots X_n$ is a random sample from population 1 and $Y_1 \dots Y_n$ is a random sample from population 2.

• X_i and Y_i are dependent. (X_i, Y_i) and (X_j, Y_j) are independent for any $i \neq j$.

• $D = X_i - Y_i$ and $\mu_D = \mu_1 - \mu_2$. Treat $D_1 \dots D_n$ as a random sample from a population with mean μ_D and variance σ_D^2 .

HYPOTHESIS TESTING

1. Set null and alternative hypothesis. Assume H_0 unless there is overwhelming evidence against it. Alternative hypothesis H_1 is the hypothesis we want to prove.

2. Set significance level $\alpha = P(\text{Type I error})$

= $P(\text{Reject } H_0 | H_0 \text{ is true})$

Given Do not reject H_0 Reject H_0

H_0 is true Correct Type I error [serious]

H_0 is false Type II error Correct

• $\beta = P(\text{Type II error}) = P(\text{Do not reject } H_0 | H_0 \text{ is false})$

• Power of the test = $1 - \beta = P(\text{Reject } H_0 | H_0 \text{ is false})$

• To decrease both errors simultaneously, reduce sample size.

3. Identify test statistics (→ How unlikely it is to observe the sample, assuming H_0 is true), distribution, and rejection criteria.

• α Divide possible values of test statistics into REJECTION REGION (CRITICAL REGION) and acceptance region.

• E.g. $N-(\mu, 2^2)$, $X_1 \dots X_n$ RVs from that distribution, if $H_0: \mu = 0$ vs $H_1: \mu \neq 0$, and reject when $|z_{obs}| > 2$ where $Z = \bar{X}/(\sigma/\sqrt{n})$. Suppose we actually reject H_0 , what is the probability that we fail to reject H_0 given $\mu = 2$? **Type II error.**

→ $P(|Z| \leq 2 \mid \mu = 2) = P(-2 \leq \frac{Z - \mu}{\sigma/\sqrt{n}} \leq 2 \mid \mu = 2)$. Standardise by subtracting 2/(2-5) from all terms (represents $\mu = 2$ part)

→ $P(-7 \leq \frac{Z - \mu}{\sigma/\sqrt{n}} \leq -3 \mid \mu = 2) = P(-7 \leq \bar{Z} \leq -3 \mid \mu = 2)$

4. Compute observed test statistic. (Refer to first table)

5. Reject or do not reject the null hypothesis.

P-VALUE / OBSERVED LEVEL OF SIGNIFICANCE

Probability of obtaining a test statistic at least as extreme as (\leq or \geq) than the observed sample value, given H_0 is true.

• **p-value < α ⇒ reject H_0** (test statistic in rejection region)

• **p-value $\geq \alpha$ ⇒ DO NOT reject H_0**

• If asked for $\Phi(\alpha)$ values, and if RIGHT tail ($\mu > \mu_0$) was tested, negate the z or t value from the test (NOT the α value)