

Distributions

Distribution	$E[x]$	$\text{Var}(x)$
$X \sim U(a, b)$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
$X \sim \text{Bin}(n, p)$	np	npq
$X \sim \text{Hg}(N, n, p)$	np	$npq \left(1 - \frac{n-1}{N-1}\right)$
$X \sim \text{Geom}(p)$	$\frac{1}{p}$	$\frac{q}{p^2}$
$X \sim \text{Exp}(\lambda)$	λ^{-1}	λ^{-2}
$X \sim \text{Pois}(\lambda)$	λ	λ
$X \sim N(\mu, \sigma^2)$	μ	σ^2
$X \sim \text{Gamma}(\alpha, \lambda)$	$\frac{\alpha}{\lambda}$	$\frac{\alpha}{\lambda^2}$

Moments

- Mean $\mu = \bar{X} = \frac{1}{n} \sum X_i$
- Variance $\sigma^2 = E[(X - \mu)^2] = \frac{1}{n} \sum (X_i - \bar{X})^2$
- Skewness $\gamma_1 = E\left[\left(\frac{X - \mu}{\sigma}\right)^3\right] = \frac{1}{\sigma^3 n} \sum (X_i - \bar{X})^3$
- Kurtosis $\gamma_2 = \frac{E[(X - \mu)^4]}{E[(X - \mu)^2]^2} = \frac{1}{\sigma^4 n} \sum (X_i - \bar{X})^4$

$$\text{Cov}(X, Y) = E[XY] - E[X]E[Y] = \frac{n}{n-1} (\bar{X}\bar{Y} - \bar{X}\bar{Y})$$

Normal distribution properties

P-value $P(Z \geq z) = 1 - \Phi(z)$ where Φ is the CDF of $N(0, 1)$.

Central limit theorem If X_1, \dots, X_n are IID with $E[X_i] = \mu$, $\text{Var}(X_i) = \sigma^2$, $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$.

Normal approximations

- $\text{Bin}(n, p) \approx N(np, npq)$
- $\text{Pois}(\lambda) \approx N(\lambda, \lambda)$
- $\text{Hg}(N, n, p) \approx N\left(np, npq \frac{N-n}{N-1}\right)$

Point estimates

Unbiased estimates of μ

• Sample mean $\bar{X} = \frac{X_1 + \dots + X_n}{n}$. Variance is $s_{\bar{X}}^2 = \sigma^2/n$ if IID, $s_{\bar{X}}^2 = \frac{\sigma^2}{n} \left(1 - \frac{n-1}{N-1}\right)$ otherwise.

• Stratified sample mean $\bar{X}_s = W_1\bar{X}_1 + \dots + W_L\bar{X}_L$. Variance is $s_{\bar{X}}^2 = (W_1s_{\bar{X}_1}^2) + \dots + (W_Ls_{\bar{X}_L}^2)$ where

$$s_{\bar{X}_i}^2 = \frac{\sigma_i^2}{n_i}$$

Unbiased estimates of σ^2

• Sample variance $s^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2 = \frac{n}{n-1} (\bar{X}^2 - \bar{X}^2)$

Stratification

- Mean $\mu = W_1\mu_1 + \dots + W_L\mu_L$
- Variance $\sigma^2 = \bar{\sigma}^2 + \sum W_i(\mu_i - \mu)^2$
- Avg. variance $\bar{\sigma}^2 = W_1\sigma_1^2 + \dots + W_L\sigma_L^2$
- Avg. stddev $\bar{\sigma} = W_1\sigma_1 + \dots + W_L\sigma_L$
- Pooled sample mean $\bar{X}_p = \frac{1}{n} (n_1\bar{X}_1 + \dots + n_L\bar{X}_L)$ biased with $\sum \left(\frac{n_i}{n} - W_i\right) \mu_i!$

Optimal allocation $n_i = n \frac{W_i\sigma_i}{\bar{\sigma}}$, proportional allocation $n_i = nW_i$.

Estimation

Method of moments

Given $E[X] = f(\theta, \gamma)$ and $E[X^2] = g(\theta, \gamma)$ solve the system $\bar{X} = f(\theta, \gamma)$, $\bar{X}^2 = g(\theta, \gamma)$.

Maximum likelihood

Given $L(\theta) = f(x_1, \dots, x_n | \theta)$ minimize $L(\theta)$ to obtain θ . For IID samples $L(\theta) = f(x_1 | \theta) \dots f(x_n | \theta)$.

Tests

Large-sample proportion

- $H_0 : p = p_0 \implies Z = \frac{\hat{p} - p_0}{\sqrt{p_0q_0/n}}$
- Two-sided rejection region: $\{Z \geq z_{\alpha/2}, Z \leq -z_{\alpha/2}\}$, $\Phi z_{\alpha} = 1 - \alpha$.
- Power function $P_w = P(\text{reject } H_0 | H_1 \text{ is true})$.
- P-value: $P(Z \geq Z_{\text{observed}})$, reject H_0 if $P \leq \alpha$.

Small-sample proportion

- $H_0 : p = p_0 \implies Z = \text{Bin}(n, p)$

Tests for mean

- Large samples: $H_0 : \mu = \mu_0 \implies Z = \frac{\bar{X} - \mu_0}{S_{\bar{X}}} \sim N(0, 1)$
- Small samples: $H_0 : \mu = \mu_0 \implies Z = \frac{\bar{X} - \mu_0}{S_{\bar{X}}} \sim t_{n-1}$
- CI method: reject at $\alpha\%$ if a $(100 - \alpha)\%$ CI doesn't cover μ_0

Likelihood ratio

- Testing $H_0 : \theta = \theta_0$ against $H_1 : \theta = \theta_1$
- Statistic: $\Lambda = \frac{L(\theta_0)}{L(\theta_1)}$
- Reject H_0 if $\Lambda \leq \lambda_{\alpha}$

Pearson's χ^2 test

- Observation belongs in one of J classes
- $H_0 : (p_1, \dots, p_J) = (p_1(\lambda), \dots, p_J(\lambda))$
- Statistic: $X^2 = \sum_j \frac{(O_j - E_j)^2}{E_j}$ with cell counts $E_j = n \cdot p_j(\hat{\lambda})$

Variance analysis

One-way ANOVA

One factor, I levels, I independent IID samples Y_{i1}, \dots, Y_{ij} . H_0 : all treatments have the same effect. Key data:

- $SS_{TOT} = SS_A + SS_E = \sum \sum (Y_{ij} - \bar{Y}_{..})^2$
- $SS_A = J \sum \hat{\alpha}_i^2$
- $SS_E = \sum \sum \hat{\epsilon}_{ij}^2$
- $MS_A = \frac{SS_A}{I-1}$, $E[MS_A] = \sigma^2 + \frac{J}{I-1} \sum \alpha_i^2$
- $MS_E = \frac{SS_E}{I(J-1)}$, $E[MS_E] = \sigma^2$
- $s_p^2 = MS_E = \frac{1}{I(J-1)} \sum \sum (Y_{ij} - \bar{Y}_{i.})^2$

Normal theory model $Y_{ij} \sim N(\mu_i, \sigma^2)$, $Y = \mu + \alpha_i + \epsilon_{ij}$, $\epsilon_{ij} \sim N(0, \sigma^2)$. MLE pooled sample mean $\hat{\mu} = \bar{Y}_{..}$, $\hat{\alpha}_i = \bar{Y}_{i.} - \bar{Y}_{..}$. Reject H_0 for large values of $\frac{MS_A}{MS_E}$ with null distribution $F_{I-1, I(J-1)}$.

Bonferroni method Overall level α in k independent tests if each test has level α/k . Simultaneous CI for $\binom{k}{2}$ pairwise differences is $(\bar{Y}_{u.} - \bar{Y}_{v.}) \pm t_{I(J-1)} \left(\frac{\alpha}{I(I-1)} \right) s_p \sqrt{\frac{2}{J}}$.

Tukey method I independent $N(\mu_i, \sigma^2)$ samples with equal size J gives Tukey's simultaneous CI as $(\bar{Y}_{u.} - \bar{Y}_{v.}) \pm q_{I, I(J-1)}(\alpha) \frac{s_p}{\sqrt{J}}$.

Kruskal-Wallis Nonparametric test for H_0 : equal distributions. Does not assume normality. Pooled sample size $N = J_1 + \dots + J_I$, pooled sample ranking $R_{ij} =$ ranks of Y_{ij} with $\sum \sum R_{ij} = \frac{N(N+1)}{2}$ and $\bar{R}_{..} = \frac{N+1}{2}$. Test statistic becomes $K = \frac{12}{N(N+1)} \sum J_i \left(\bar{R}_{i.} - \frac{N+1}{2} \right)^2$ with null distribution χ_{I-1}^2 .

Two-way ANOVA

Two factors A with I rows and B with J columns, and K observations per cell. Key data:

- $SS_{TOT} = SS_A + SS_B + SS_{AB} + SS_E = \sum \sum \sum (Y_{ijk} - \bar{Y}_{...})^2$
- $SS_A = JK \sum \hat{\alpha}_i^2$
- $SS_B = IK \sum \hat{\beta}_j^2$
- $SS_{AB} = K \sum \sum \hat{\delta}_{ij}^2$
- $SS_E = \sum \sum \sum \hat{\epsilon}_{ijk}^2$
- $MS_A = \frac{SS_A}{I-1}$, $E[MS_A] = \sigma^2 + \frac{JK}{I-1} \sum \alpha_i^2$
- $MS_B = \frac{SS_B}{J-1}$, $E[MS_B] = \sigma^2 + \frac{IK}{J-1} \sum \beta_j^2$
- $MS_{AB} = \frac{SS_{AB}}{(I-1)(J-1)}$, $E[MS_{AB}] = \sigma^2 + \frac{K}{(I-1)(J-1)} \sum \sum \delta_{ij}^2$
- $MS_E = \frac{SS_E}{IJ(K-1)}$, $E[MS_E] = \sigma^2$
- $s_p^2 = MS_E = \frac{1}{IJ(K-1)} \sum \sum \sum (Y_{ijk} - \bar{Y}_{ij.})^2$

Normal theory model $Y_{ijk} = \mu + \alpha_i + \beta_j + \delta_{ij} + \epsilon_{ijk}$, $\epsilon_{ijk} \sim N(0, \sigma^2)$. MLEs $\hat{\mu} = \bar{Y}_{...}$, $\hat{\alpha}_i = \bar{Y}_{i..} - \bar{Y}_{...}$, $\hat{\beta}_j = \bar{Y}_{.j.} - \bar{Y}_{...}$, $\hat{\delta}_{ij} = \bar{Y}_{ij.} - \bar{Y}_{i..} - \hat{\alpha}_i - \hat{\beta}_j = \bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j.} + \bar{Y}_{...}$.

Tukey method Tukey's simultaneous CI: $(\bar{Y}_{u..} - \bar{Y}_{v..}) \pm q_{I, IJ(K-1)}(\alpha) \frac{s_p}{\sqrt{J}}$

Additive model For $K = 1$ no interaction ($\delta_{ij}^2 = 0$). Statistics $\frac{MS_A}{MS_E} \sim F_{I-1, (I-1)(J-1)}$ and $\frac{MS_B}{MS_E} \sim F_{J-1, (I-1)(J-1)}$.

Randomized block design

Experimental design with I treatments randomly assigned within J blocks. H_0 : no treatment effects. Parametric uses two-way ANOVA.

Friedman's test Ranking within j th block $(R_{1j}, \dots, R_{Ij}) =$ ranks of (Y_{1j}, \dots, Y_{Ij}) giving $R_{1j} + \dots + R_{Ij} = \frac{I(I+1)}{2}$, implying $\frac{1}{I}(R_{1j} + \dots + R_{Ij}) = \frac{I+1}{2}$ and $\bar{R}_{.j} = \frac{I+1}{2}$. Test statistic $Q = \frac{12J}{I(I+1)} \sum (\bar{R}_{i.} - \frac{I+1}{2})^2$ with $Q \sim \chi_{I-1}^2$.

Bayesian inference

Conjugate priors

Data	Prior	Posterior
$X \sim N(\theta, \sigma^2)$	$\mu \sim N(m, v^2)$	$N(\gamma m + (1 - \gamma)\bar{x}, \gamma v^2)$
$X \sim \text{Bin}(n, p)$	$p \sim \text{Beta}(a, b)$	$\text{Beta}(a + x, b + n - x)$
$\text{Mn}(n; p_1, \dots, p_r)$	$D(\alpha_1, \dots, \alpha_r)$	$D(\alpha_1 + x_1, \dots, \alpha_r + x_r)$
$X \sim \text{Pois}(\mu)$	$\mu \sim \Gamma(\alpha, \lambda)$	$\Gamma(\alpha + x, \lambda + 1)$
$X \sim \text{Exp}(\rho)$	$\rho \sim \Gamma(\alpha, \lambda)$	$\Gamma(\alpha + 1, \lambda + x)$

with $\gamma = \frac{\sigma^2}{\sigma^2 + nv^2}$.

Credibility interval CI interval for the posterior distribution $(P(\theta_0(x) < \theta < \theta_1(x)) = 1 - \alpha$ for random θ).

Summarizing data

Survival function $S(t) = P(T > t) = 1 - F(t)$ for lifelength T . Hazard function $h(t) = f(t)/S(t) = -\frac{d}{dt} \log(S(t))$.

Measure of location

- Median M . $H_0 : M = M_0$
- Statistic: $Z = \sum I(X_i \leq M_0)$ (number of observations below M_0)
- Reject H_0 if M_0 is not in $(X_{(k)}, X_{(n-k+1)})$ where $k = k_\alpha$ such that $P(Y < k_\alpha) = \frac{\alpha}{2}$

Measures of dispersion

Measures of σ in $N(\mu, \sigma^2)$:

- Sample standard deviation s
- Interquartile range $\frac{x_{0.75} - x_{0.25}}{1.35}$
- Median of absolute deviance $\frac{\text{median}(|X_i - \hat{M}|)}{0.675}$

Comparing samples

Comparing two independent samples

- Large samples: normal approximation $\bar{X} - \bar{Y} \sim N(\mu_x - \mu_y, s_x^2 + s_y^2)$.
- P-value: reject H_0 if $P(Z \geq z | H_0) \leq \alpha$ where z is the observed test statistic. Two-sided P-value is two times the one-sided P-value.

Wilcoxon rank sum test

- Pool samples, replace data by ranks
- Statistic: either $R_x = \sum \text{ranks of } X$ or $R_y = \binom{n+m+1}{2} - R_x$
- Null distributions in table, for large samples apply normal approximation

Paired samples

Paired IID samples $(X_1, Y_1), \dots, (X_n, Y_n)$

- Transform to $D_i = X_i - Y_i$ estimating $\mu_x - \mu_y = \bar{D} = \bar{X} - \bar{Y}$
- Correlation coefficient $\rho = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y} > 0$ for paired observations
- If $\rho > 0$, $\text{Var}(\bar{D}) = \text{Var}(\bar{X}) + \text{Var}(\bar{Y}) - 2\sigma_{\bar{x}}\sigma_{\bar{y}}\rho$

Sign test Test $H_0 : M_D = 0$ with statistics $Y_+ = \sum \{D_i > 0\}$ or $Y_- = \sum \{D_i < 0\}$, null distribution $\text{Bin}(n, 0.5)$.

Wilcoxon signed rank test H_0 : distribution of D is symmetric around $M_D = 0$. Statistic $W_+ = \sum \text{rank}(|D_i|) \cdot (D_i > 0)$ or corresponding W_- . Normal approximation of null distribution has $\mu_W = \frac{n(n+1)}{4}$, $\sigma_W^2 = \frac{n(n+1)(2n+1)}{24}$.

Categorical data

Fisher's exact test

$H_0 : \pi_{11} = \pi_{12}, \pi_{21} = \pi_{22}$. Use n_{11} as a test statistic, null distribution $n_{11} \sim \text{Hg}(N, n, p)$ with parameters $N = n_{..}, n = n_{.1}, Np = n_{1.}, Nq = n_{2.}$.

χ^2 -test of homogeneity

I categories, J populations, H_0 all J distributions are equal. Use sample counts and test statistic $X^2 = \sum_i \sum_j \frac{(n_{ij} - n_{i.}n_{.j}/n_{..})^2}{n_{i.}n_{.j}/n_{..}}$. Reject H_0 for large X^2 , null distribution $X^2 \sim \chi_{df}^2$ with $df = (I-1)(J-1)$.

χ^2 -test of independence

H_0 all pairs of column/row are independent. Use homogeneity test (is equivalent).

McNemar's test

$H_0 : \pi_{12} = \pi_{21}$. Use statistic $X^2 = \frac{(n_{12} - n_{21})^2}{n_{12} + n_{21}}$ with null distribution χ_1^2 . Use normal approximation with 2-sided P-value $2(1 - \Phi(\sqrt{X^2}))$, reject H_0 if $X^2 \leq \alpha$.