

Confidence Intervals

Estimating Means

For large r , $t_r \rightarrow N(0, 1)$ is a good approximation.

Normal, Single Mean, Known σ

$(\bar{x} \pm c \frac{\sigma}{\sqrt{n}})$; with $c = F^{-1}(1 - \frac{\alpha}{2})$ from pivot $N(0, 1)$.

Normal, Single Mean, Unknown σ

$(\bar{x} \pm c \frac{s}{\sqrt{n}})$; with $c = F^{-1}(1 - \frac{\alpha}{2})$ from pivot t_{n-1} .

Normal, Two Means, Two Known σ 's

$(\bar{x} - \bar{y}) \pm c \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}$; with c from pivot $N(0, 1)$.

Normal, Two Means, Unknown σ 's, Many Samples

$(\bar{x} - \bar{y}) \pm c \sqrt{\frac{s_X^2}{n} + \frac{s_Y^2}{m}}$; with c from pivot $N(0, 1)$.

Normal, Two Means, Unknown σ 's, Common Variance

$(\bar{x} - \bar{y}) \pm c \cdot s_P \sqrt{\frac{1}{n} + \frac{1}{m}}$; with c from pivot t_{n+m-2} .

$$s_P = \sqrt{\frac{(n-1)s_X^2 + (m-1)s_Y^2}{n+m-2}}$$

Normal, Two Means, Unknown σ 's, \neq Variances

$(\bar{x} - \bar{y}) \pm c \sqrt{\frac{s_X^2}{n} + \frac{s_Y^2}{m}}$; with c from pivot t_r .

$$r = \left(\frac{s_X^2}{n} + \frac{s_Y^2}{m} \right)^2 / \left(\frac{s_X^4}{n^2(n-1)} + \frac{s_Y^4}{m^2(m-1)} \right)$$

Normal, Paired Samples

For pairs (X_i, Y_i) , let $D_i = X_i - Y_i$. $D_i \sim N(\mu_D, \sigma_D^2)$.

$(\bar{d} \pm c \frac{s_d}{\sqrt{n}})$; with c from pivot t_{n-1} .

Estimating Variance

Normal, Single Variance - **estimate of σ^2 not σ !**

$(\frac{(n-1)s^2}{b}, \frac{(n-1)s^2}{a})$; with a, b from pivot χ_{n-1}^2 .

Normal, Two Variances

A confidence interval for the ratio of the variances σ_X^2/σ_Y^2 is
 $(a \cdot \frac{s_X^2}{s_Y^2}, b \cdot \frac{s_X^2}{s_Y^2})$; with a, b from pivot $F_{m-1, n-1}$.

Estimating Proportions (Large n - Normal Approx.)

Single Proportion

$\approx (\hat{p} \pm c \sqrt{\frac{\hat{p}(1-\hat{p})}{n}})$; with c from pivot $N(0, 1)$.

Two Proportions

$\approx (\hat{p}_1 - \hat{p}_2 \pm c \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}})$; c from pivot $N(0, 1)$.

Can also derive a confidence interval from the exact distribution of the binomial RV $n \cdot \hat{p} \sim \text{Bi}(n, p)$.

Prediction Intervals

Let X^* be a future realisation of X .

If $X \sim N(\mu, \sigma^2)$, $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$, and $(\bar{X} - X^*) \sim N(0, \sigma^2 + \frac{\sigma^2}{n})$.

$(\bar{x} \pm c \sqrt{s^2 + \frac{s^2}{n}})$ is a prediction interval with c from t_{n-1} .

If σ is known, use pivot $N(0, 1)$. If μ is known use χ_{n-1}^2 .

Regression

Ordinary Least Square (OLS) Estimators

All parameters are unbiased and except $\hat{\sigma}^2$ are normally distributed (Variances as below).

Let $K = \sum_{i=1}^n (x_i - \bar{x})^2$ and $D^2 = \sum_{i=1}^n (Y_i - \hat{\alpha}_0 - \hat{\beta}(x_i - \bar{x}))^2$:

$$\hat{\alpha}_0 = \bar{Y}; \quad \hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}; \quad \hat{\alpha} = \hat{\alpha}_0 - \hat{\beta}\bar{x}; \quad \hat{\sigma}^2 = \frac{D^2}{n-2}$$

$$\text{Var}(\hat{\alpha}) = \left(\frac{1}{n} + \frac{\bar{x}^2}{K} \right) \sigma^2; \quad \text{Var}(\hat{\beta}) = \frac{\sigma^2}{K}; \quad \text{Var}(\hat{\alpha}_0) = \frac{\sigma^2}{n}$$

$$\text{Cov}(\hat{\alpha}_0, \hat{\beta}) = 0; \quad \text{Var}(\hat{\mu}(x)) = \left(\frac{1}{n} + \frac{(x - \bar{x})^2}{K} \right) \sigma^2$$

Regression Pivots

Notice that all the t_{n-2} distributed pivots are estimate \div SE.

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2; \quad \frac{\hat{\beta} - \beta}{\hat{\sigma}/\sqrt{K}} \sim t_{n-2}; \quad \frac{\hat{\mu}(x) - \mu(x)}{\hat{\sigma}\sqrt{\frac{1}{n} + \frac{(x-\bar{x})^2}{K}}} \sim t_{n-2}$$

$$\frac{\hat{\alpha} - \alpha}{\hat{\sigma}\sqrt{\frac{1}{n} + \frac{\bar{x}^2}{K}}} \sim t_{n-2}; \quad \frac{\hat{\alpha}_0 - \alpha_0}{\hat{\sigma}/\sqrt{n}} \sim t_{n-2}$$

Distribution Free Methods

Sign Test

Hypotheses: $H_0 : m = m_0$; $H_1 : m \neq m_0$ or 1 sided H_1 .

Let $Y = \sum_{i=1}^n I(X_i - m_0 > 0)$. Under H_0 , $Y \sim \text{Bi}(n, 0.5)$. Can calculate p values or a critical region from $Y \sim \text{Bi}(n, 0.5)$.

Paired Samples: Replace (x_i, y_i) with $\text{sgn}(x_i - y_i)$ and use the same distribution to check for equal medians.

Wilcoxon Signed Rank / One Sample Test

Hypotheses: $H_0 : m = m_0$; $H_1 : m \neq m_0$ or 1 sided H_1 .

Rank the $|X_i - m_0|$ from $1 \rightarrow n$ starting at 1 and replace X_i with $\text{sgn}(X_i - m_0) \cdot \text{rank}(|X_i - m_0|)$. Let W be the sum of the signed ranks. Under H_0 :

$$Z = \frac{W - 0}{\sqrt{n(n+1)(2n+1)/6}} \approx N(0, 1)$$

For paired samples we can take the difference and test for equality of medians. When tied ranks occur give them each the average of the ranks they span.

Wilcoxon Rank Sum / Two Sample Test

Hypotheses: $H_0 : m_X = m_Y$; $H_1 : \bar{H}_0$ or 1 sided H_1 .

Order the combined sample and let W be the sum of the ranks of Y_1, \dots, Y_{n_Y} . W is approximately normal with:

$$\mathbb{E}(W) = \frac{n_Y(n_X + n_Y + 1)}{2}; \quad \text{Var}(W) = \frac{n_X n_Y (n_X + n_Y + 1)}{12}$$

Bayesian Inference

USE THE LIKELIHOOD NOT THE PDF!

Law of Total Probability

$$f(x) = \int_{-\infty}^{\infty} f(x, y) dy = \int_{-\infty}^{\infty} f(x | y) f(y) dy$$

Normal Distribution, Known σ , Inference for μ

$X_1, \dots, X_n \sim N(\mu, \sigma^2)$. Let $Y = \bar{X} \sim N(\mu, \sigma^2/n)$.

Prior: $\mu \sim N(\mu_0, \sigma_0^2)$

Posterior: $f(\mu | y) \propto f(y | \mu) f(\mu) \propto \exp \left[-\frac{(\mu - \mu_1)^2}{2\sigma_1^2} \right]$

Where:

$$\mu_1 = \frac{\frac{\mu_0}{\sigma_0^2} + \frac{y}{\sigma^2/n}}{\frac{1}{\sigma_0^2} + \frac{1}{\sigma^2/n}}; \quad \frac{1}{\sigma_1^2} = \frac{1}{\sigma_0^2} + \frac{1}{\sigma^2/n}$$

So $\mu | y \sim N(\mu_1, \sigma_1^2)$

Binomial Distribution, Beta / Uniform Prior

$X \sim \text{Bi}(n, \theta)$. Prior: $\theta \sim \text{Beta}(\alpha, \beta)$

Posterior: $\theta | x \sim \text{Beta}(\alpha + x, \beta + n - x)$

Pseudodata

Gamma prior: $\gamma(\alpha, \beta)$ where α is the total count over all samples and β is the number of samples taken.

Beta prior: $\text{Beta}(\alpha, \beta)$ where α is the number of successes and β is the number of failures.

Pivots

Sampling from a normal distribution:

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1); \quad \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t_{n-1}; \quad \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$$

ANOVA

All populations assumed to have equal variance.

F-Statistics / F-test

To test any H_0 in ANOVA, use the F statistic. Since the MS terms are χ^2 distributed, $F = MS(X) \div MS(E)$ is distributed as F_{df_X, df_E} and the rejection region is $F > c$.

eg. To test $H_{0AB} : \gamma_{ij} = 0 \forall i, j$

Single Factor (One Way) ANOVA

Total Sample Size: $n = n_1 + \dots + n_k$

Means:

$$\bar{X}_{i.} = \frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij}$$

$$\bar{X}_{..} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} X_{ij} = \frac{1}{n} \sum_{i=1}^k n_i \bar{X}_{i.}$$

Hypotheses: $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$ vs. $H_1 : \bar{H}_0$

Sum of Squares:

Treatment SS:

$$SS(T) = \sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{X}_{i.} - \bar{X}_{..})^2 = \sum_{i=1}^k n_i (\bar{X}_{i.} - \bar{X}_{..})^2$$

Error SS:

$$SS(E) = \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_{i.})^2 = \sum_{i=1}^k (n_i - 1) S_i^2$$

Total SS:

$$SS(TO) = \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_{..})^2 = SS(T) + SS(E)$$

Degrees of Freedom:

$$\text{Treatment : } k - 1; \quad \text{Error : } n - k$$

Null Distributions of Sum of Squares Terms:

$$\frac{SS(T)}{\sigma^2} \sim \chi_{k-1}^2; \quad \frac{SS(E)}{\sigma^2} \sim \chi_{n-k}^2; \quad \frac{SS(TO)}{\sigma^2} \sim \chi_{n-1}^2$$

Two Factor (Two Way) ANOVA

Assume $X_{ij} \sim N(\mu_{ij}, \sigma^2)$. $\mu_{ij} = \mu + \alpha_i + \beta_j$. $\sum_{i=1}^a \alpha_i = 0$ and $\sum_{j=1}^b \beta_j = 0$. Take one observation from each combination of a and b ($n = ab$).

Hypotheses: $H_{0A} : \alpha_1 = \alpha_2 = \dots = \alpha_k = 0$ vs. $H_1 : \bar{H}_0$

or $H_{0B} : \beta_1 = \beta_2 = \dots = \beta_k = 0$ vs. $H_1 : \bar{H}_0$

Means:

$$\bar{X}_{..} = \frac{1}{ab} \sum_{i=1}^a \sum_{j=1}^b X_{ij}; \quad \bar{X}_{i.} = \frac{1}{b} \sum_{j=1}^b X_{ij}; \quad \bar{X}_{.j} = \frac{1}{a} \sum_{i=1}^a X_{ij}$$

Sum of Squares:

$$SS(A) = b \sum_{i=1}^a (\bar{X}_{i.} - \bar{X}_{..})^2$$

$$SS(B) = a \sum_{j=1}^b (\bar{X}_{.j} - \bar{X}_{..})^2$$

Error SS:

$$SS(E) = \sum_{i=1}^a \sum_{j=1}^b (X_{ij} - \bar{X}_{i.} - \bar{X}_{.j} + \bar{X}_{..})^2$$

Total SS:

$$SS(TO) = SS(A) + SS(B) + SS(E)$$

Degrees of Freedom:

$$A : a - 1; \quad B : b - 1; \quad \text{Error : } (a - 1)(b - 1)$$

Two Factor ANOVA with Interaction Terms

Take samples over two factors, but c samples for each factor pair. $\mu_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij}$.

Hypotheses: $H_{0A} : \alpha_1 = \alpha_2 = \dots = \alpha_k = 0$ vs. $H_1 : \bar{H}_0$

or $H_{0B} : \beta_1 = \beta_2 = \dots = \beta_k = 0$ vs. $H_1 : \bar{H}_0$

or $H_{0AB} : \gamma_{ij} = 0 \forall i, j$ vs. $H_1 : \bar{H}_0$

Means:

$$\bar{X}_{ij.} = \frac{1}{c} \sum_{k=1}^c X_{ijk}; \quad \bar{X}_{i..} = \frac{1}{bc} \sum_{j=1}^b \sum_{k=1}^c X_{ijk}$$

$$\bar{X}_{.j.} = \frac{1}{ac} \sum_{i=1}^a \sum_{k=1}^c X_{ijk}; \quad \bar{X}_{...} = \frac{1}{abc} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^c X_{ijk}$$

Sum of Squares:

$$SS(A) = bc \sum_{i=1}^a (\bar{X}_{i..} - \bar{X}_{...})^2$$

$$SS(B) = ac \sum_{j=1}^b (\bar{X}_{.j.} - \bar{X}_{...})^2$$

$$SS(AB) = c \sum_{i=1}^a \sum_{j=1}^b (\bar{X}_{ij.} - \bar{X}_{i..} - \bar{X}_{.j.} + \bar{X}_{...})^2$$

Error SS:

$$SS(E) = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^c (X_{ijk} - \bar{X}_{ij.})^2$$

Total SS:

$$SS(TO) = SS(A) + SS(B) + SS(AB) + SS(E)$$

Degrees of Freedom:

$$A : a - 1; \quad B : b - 1; \quad AB : (a - 1)(b - 1); \quad \text{Error : } ab(c - 1)$$

Mean Squares / Mean Square Error

MSE is found by dividing the $SS(E)$ by its degrees of freedom.

$\hat{\sigma}^2 = MS(E)$ is an unbiased estimator for the variance.

Distribution of Order Statistics

CDF of $X_{(k)}$

$$G_k(x) = \Pr(X_k \leq x) = \sum_{i=k}^n \binom{n}{i} F(x)^i (1 - F(x))^{n-i}$$

PDF of $X_{(k)}$

$$g_k(x) = k \binom{n}{k} F(x)^{k-1} (1 - F(x))^{n-k} f(x)$$

Goodness-of-fit Test (χ^2)

O_i is the Observed count of data in class i , E_i (> 5) is the Expected count under H_0 . Be careful to use the count and not the proportion. For k classes and e estimated parameters:

$$Q_{k-1} = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \approx \chi_{k-1}^2$$

Since the numerator $(O_i - E_i)^2$ will be larger when H_0 is false, we reject H_0 for $Q_{k-1} > c$, with c from χ_{k-1-e}^2 .

Two Classes ($k = 2$)

Testing $H_0 : p = p_1$ vs. $H_1 : p \neq p_1$

More than Two Classes ($k > 2$)

Let p_1, \dots, p_k define the proportions of a categorical distribution. $H_0 : "p_1, \dots, p_k \text{ do define the distribution}"$ vs. $H_1 : \bar{H}_0$

Sufficient Statistics

Exponential Family

$$f(x | \theta) = \exp\{K(x)p(\theta) + S(x) + q(\theta)\}$$

Has $\sum_{i=1}^n K(X_i)$ as a sufficient statistic for θ .

Factorisation Theorem

$Y = g(x_1, \dots, x_n)$ is sufficient for θ iff:

$$f(x_1, \dots, x_n | \theta) = \phi\{g(x_1, \dots, x_n) | \theta\} h(x_1, \dots, x_n)$$

Other Formulae

Sample Variance

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

Cramér-Rao Lower Bound

The CR Lower Bound is the minimum possible variance of an estimator $\hat{\theta}$ of θ . It is the asymptotic variance of the MLE.

$$\ell(\theta) = \ln L(\theta); \quad U(\theta) = \frac{\partial \ell}{\partial \theta} \quad (\text{Score Function})$$

$$V(\theta) = -\frac{\partial U}{\partial \theta}; \quad I(\theta) = \mathbb{E}(V(\theta)) \quad (\text{Fisher Information})$$

$$\text{Var}(\hat{\theta}) \geq \frac{1}{I(\theta)} \quad (\text{CR Lower Bound})$$

In the Fisher Information, take $\mathbb{E}(V(\theta))$ over all x , not θ .

Gamma Function in Beta + γ Distributions

$$\Gamma(n) = (n - 1)!$$