

Statistics

Marie Biolková

Useful Properties

Always:

$$\begin{aligned}\mathbb{E}(aX + b) &= a\mathbb{E}(X) + b \\ \mathbb{E}(X + Y) &= \mathbb{E}(X) + \mathbb{E}(Y) \\ \text{Var}(aX + b) &= a^2\text{Var}(X)\end{aligned}$$

Only if *independent*:

$$\begin{aligned}\mathbb{E}(XY) &= \mathbb{E}(X)\mathbb{E}(Y) \\ \text{Var}(X + Y) &= \text{Var}(X) + \text{Var}(Y) \\ \text{Var}(X - Y) &= \text{Var}(X) + \text{Var}(Y) \\ \text{Var}\left(\sum_{i=1}^n X_i\right) &= \sum_{i=1}^n \text{Var}(X_i)\end{aligned}$$

Sample Mean

Unbiased and consistent estimator of μ .

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

Sample Variance

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - \frac{1}{n} \left(\sum_{i=1}^n X_i \right)^2 \right)$$

- Unbiased and consistent estimator of σ^2 .
- Since $\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2$ we have that $\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$ so we can estimate $\text{Var}(\bar{X})$ by $\frac{S^2}{n}$.

Sample Covariance and Correlation

$$\text{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y) \quad (\text{covariance})$$

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \quad (\text{correlation})$$

Sample Covariance:

$$\begin{aligned}S_{xy} &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n X_i Y_i - \frac{1}{n} \left(\sum_{i=1}^n X_i \right) \left(\sum_{i=1}^n Y_i \right) \right]\end{aligned}$$

- Unbiased and consistent estimator of $\text{Cov}(X, Y)$
- S_{xx} and S_{yy} are the sample variances for X and Y , recall $\text{Cov}(X, X) = \text{Var}(X)$.

Sample Correlation:

$$R_{xy} = \frac{S_{xy}}{S_x S_y}$$

Maximum Likelihood Estimators (MLEs)

Assuming the data are independent, the *likelihood function* is

$$L(\theta; x_1, \dots, x_n) = \prod_{i=1}^n f(x_i; \theta).$$

The *log-likelihood* is therefore

$$l(\theta; x_1, \dots, x_n) = \sum_{i=1}^n \log f(x_i; \theta).$$

- $\hat{\sigma}^2$ is not the sample variance S^2 .
- In general MLEs are biased estimators.
- Consistent estimators.

Invariance Property of MLEs: Let $\hat{\theta}$ be the MLE of θ and g be any function of θ . Then the MLE of $g(\theta)$ is $g(\hat{\theta})$.

Properties of the Sample Mean and Variance for the Normal Distribution

Let X_1, \dots, X_n be independent $N(\mu, \sigma^2)$ rvs, then

- $\frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \sim N(0, 1)$
- $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$
- \bar{X} and S^2 are independent.

Normal Distribution with Known Variance

Assume $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ are independent rvs, σ^2 known. Recall $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ then the linear transform

$$Z = \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}}$$

is such that $Z \sim N(0, 1)$.

The $(1 - \alpha)\%$ *confidence interval* for μ is given by

$$\bar{x} \pm \frac{z_{\alpha/2}\sigma}{\sqrt{n}}.$$

- To calculate $z_{\alpha/2}$ in R use `qnorm(1-alpha/2, 0, 1)`, e.g. `qnorm(0.975, 0, 1)` for 95% CI.
- CI is larger for smaller sample size.
- Higher % confidence interval results in wider interval.

Normal Distribution with Unknown Variance

Assume $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ are independent rvs, σ^2 unknown. Consider

$$T = \frac{\bar{X} - \mu}{\sqrt{S^2/n}}.$$

χ^2 Distribution

Let Z_1, \dots, Z_n be independent $N(0, 1)$ rvs and $X = \sum_{i=1}^n Z_i^2$. Then X has chi-squared distribution with n degrees of freedom, $X \sim \chi_n^2$.

- X is a continuous rv and $x \geq 0$.
- Let $Z \sim N(0, 1)$ and $Y = Z^2$. Then $Y \sim \chi_1^2$.
- Let $X \sim \chi_n^2$ and $Y \sim \chi_m^2$, independently. Then $X + Y \sim \chi_{n+m}^2$.
- If $X \sim \chi_n^2$ then $\mathbb{E}(X) = n$ and $\text{Var}(X) = 2n$.

t Distribution

Let X and Y be independent rvs such that $Z \sim N(0, 1)$ and $Y \sim \chi_n^2$. Let $T = \frac{Z}{\sqrt{Y/n}}$, then T has a t -distribution with n degrees of freedom, i.e. $T \sim t_n$.

- T is a continuous rv, $t \in \mathbb{R}$.
- As $n \rightarrow \infty$, $t_n \rightarrow N(0, 1)$.
- If $T \sim t_n$ the $\mathbb{E}(T) = 0$ and $\text{Var}(T) = n/(n-1)$ for $n > 2$.
- Denote $t_{n;\alpha}$ the *upper α quantile*, i.e. $\mathbb{P}(T \geq t_{n;\alpha}) = \alpha$.
- Symmetrical about 0.

$$\frac{\bar{X} - \mu}{\sqrt{S^2/n}} \sim t_{n-1}$$

The $(1 - \alpha)\%$ *confidence interval* for μ is

$$\bar{x} \pm t_{n-1;\alpha/2} \frac{s}{\sqrt{n}}.$$

- To calculate $t_{n-1;\alpha/2}$ in R use `qt(1-alpha/2, n-1)`.
- The CI is larger when the variance is unknown.

Hypothesis Testing

- *Type I error:* Reject H_0 when it is in fact true.
- *Type II error:* Fail to reject H_0 when it is false.
- *Significance level α :* Probability that we reject H_0 when it is true, $\mathbb{P}(\text{Type I error}) = \alpha$.
- *Power β :* Probability that we reject H_0 when it is false, $\mathbb{P}(\text{Type II error}) = 1 - \beta$.
- *Power function:* $\beta(\theta) = \mathbb{P}(\text{reject } H_0 : \theta = \theta_0 \text{ when the true value is } \theta)$.
- *Test statistic:* Function of the data chosen, is expected to take a different range of values when H_0 is true than when it is false.
- *Critical region C* The set of values of t that lead us to reject H_0 .
- *p-value* is the probability of observing a result at least as extreme as t if H_0 is true.
 - *p-value small ($< \alpha$):* reject H_0 .
 - *p-value large ($\geq \alpha$):* no evidence to reject H_0 .

Increasing sample size means we are more likely to reject H_0 if it is false.

z-test

X_1, \dots, X_n independent $N(\mu, \sigma^2)$ rvs, σ^2 known.

- $H_0 : \mu = \mu_0$ vs $H_1 : \mu \neq \mu_0$
- Test statistic: $T = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$, then under H_0 , $T \sim N(0, 1)$.
- Critical region: $|T| = \left| \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \right| \geq z_{\alpha/2}$.
- p-value: $\mathbb{P}(|T| \geq t_0) = 2\mathbb{P}(T \geq t_0) = 2\mathbb{P}(T \leq -t_0)$ ¹.

One Sample t-test

X_1, \dots, X_n independent $N(\mu, \sigma^2)$ rvs, σ^2 unknown.

- $H_0 : \mu = \mu_0$ vs $H_1 : \mu \neq \mu_0$
- Test statistic: $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$, then under H_0 , $T \sim t_{n-1}$.
- Critical region: reject H_0 if $|T| \geq t_0 = t_{n-1, \alpha/2}$.
- p-value: $\mathbb{P}(|T| \geq t_0) = 2\mathbb{P}(T \geq t_0) = 2\mathbb{P}(T \leq -t_0)$

Paired t-test

Paired data $(X_1, Y_1), \dots, (X_n, Y_n)$ where the two measurements are dependent. Consider the difference such that $D_i = Y_i - X_i$ for $i = 1, \dots, n$.

Assume $D_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ - observed differences are independent of each other and observations are from normal distribution with mean μ and unknown variance σ^2 .

Reduces to a one-sample t-test.

- $H_0 : \mu = 0$ vs $H_1 : \mu \neq 0$
- Test statistic: $T = \frac{\bar{D}}{S/\sqrt{n}}$, then under H_0 , $T \sim t_{n-1}$.

Two Sample t-test

Suppose we have two sets of independent rvs X_1, \dots, X_n and Y_1, \dots, Y_m such that $X_i \sim N(\mu_X, \sigma^2)$, $Y_i \sim N(\mu_Y, \sigma^2)$.

$$\frac{(n-1)S_X^2}{\sigma^2} \sim \chi_{n-1}^2 \quad \frac{(m-1)S_Y^2}{\sigma^2} \sim \chi_{m-1}^2$$

Pooled sample variance:

$$S_p^2 = \frac{(n-1)S_X^2 + (m-1)S_Y^2}{m+n-2}$$

- $H_0 : \mu_X = \mu_Y$ vs $H_1 : \mu_X \neq \mu_Y$
- Test statistic: $T = \frac{\bar{X} - \bar{Y}}{S_p \sqrt{\frac{1}{m} + \frac{1}{n}}}$, then under H_0 , $T \sim t_{m+n-2}$.
- Critical region: reject H_0 if $|T| \geq t_0 = t_{m+n-2, \alpha/2}$.
- p-value: $\mathbb{P}(|T| \geq t_0) = 2\mathbb{P}(T \geq t_0) = 2\mathbb{P}(T \leq -t_0)$

¹ t_0 is the upper quantile, $-t_0$ is the lower quantile.

F-test for Equality of Variance

Suppose we have two independent normal rvs X_1, \dots, X_n and Y_1, \dots, Y_m with variances σ_X^2, σ_Y^2 .

- $H_0 : \sigma_X^2 = \sigma_Y^2$ vs $H_1 : \sigma_X^2 \neq \sigma_Y^2$
- Test statistic: $T = \frac{S_X^2}{S_Y^2}$, then under H_0 , $F \sim F_{n-1, m-1}$.

F Distribution

$U \sim \chi_m^2, V \sim \chi_n^2$ independent rvs. Then $X = \frac{U/m}{V/n}$ has an F distribution with m, n degrees of freedom ($X \sim F_{m, n}$).

- $1/X \sim F_{n, m}$
- Upper α quantile $F_{m, n; \alpha}$ is such that $\mathbb{P}(X \geq F_{m, n; \alpha}) = \alpha$, lower quantile $F_{m, n; 1-\alpha} = 1/F_{n, m; \alpha}$.
- pf and qf commands in R

One-sided Tests

$H_0 : \theta = \theta_0$ vs $H_1 : \theta > \theta_0$
 $H_0 : \theta = \theta_0$ vs $H_1 : \theta < \theta_0$

Linear Regression

$$\mathbb{E}(Y) = \alpha + \beta x$$

Least-Squares Estimation

Want to find $\hat{\alpha}, \hat{\beta}$ that minimise the sum of squares

$$S(\alpha, \beta) = \sum_{i=1}^n [y_i - (\alpha + \beta x_i)]^2 = \sum_{i=1}^n \epsilon_i^2$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \quad \hat{\beta} = \frac{S_{XY}}{S_{XX}}$$

- Requires no assumptions about the distribution.
- $\hat{\alpha}, \hat{\beta}$ are rvs, unbiased and consistent estimators of α, β .

Simple Linear Regression

Assume Y_1, \dots, Y_n are independent, normally distributed rvs with common variance, and have a mean that is a linear function of the explanatory variable, i.e

$Y_i \stackrel{iid}{\sim} N(\alpha + \beta x_i, \sigma^2)$ $i = 1, \dots, n$.

$$\hat{\alpha} \sim N\left(\alpha, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{XX}}\right)\right) \quad \hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{S_{XX}}\right)$$

$$S^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{y}_i)^2 \quad \hat{y}_i = \hat{\alpha} + \hat{\beta} x_i \text{ (fitted value)}$$

- S^2 is an unbiased estimator of σ^2 with $\frac{(n-2)S^2}{\sigma^2} \sim \chi_{n-2}^2$.
- S^2 is independent of $\hat{\alpha}, \hat{\beta}$ (but $\hat{\alpha}, \hat{\beta}$ are not independent!)
- Standard errors: $s.e.(\hat{\alpha}) = \sqrt{\text{Var}(\hat{\alpha})}$, $s.e.(\hat{\beta}) = \sqrt{\text{Var}(\hat{\beta})}$

- Confidence intervals:

$$\hat{\alpha} \pm t_{n-2; 0.025} \times s.e.(\hat{\alpha})$$

$$\hat{\beta} \pm t_{n-2; 0.025} \times s.e.(\hat{\beta})$$

Regression using R

Command `lm(y~x)`, and `lm(y~x - 1)` for regression through the origin.

Confidence Interval for $\mathbb{E}(Y_0)$

$$\hat{\alpha} + \hat{\beta} x_0 \pm t_{n-2; 0.025} \sqrt{s^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{XX}} \right)}$$

- Interval for the predicted expectation - they reflect uncertainty in our estimates of average observation.

Prediction Interval for Y_0

$$\hat{\alpha} + \hat{\beta} x_0 \pm t_{n-2; 0.025} \sqrt{s^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{XX}} \right)}$$

- Prediction for a single observation as a function of the explanatory variable - we would expect 95% of observations to lie within this interval.
- Prediction intervals for Y_0 are wider than confidence intervals for $\mathbb{E}(Y_0)$ as they take into account uncertainty relating to the expected value and individual variability.
- Confidence and prediction intervals become wider as x_0 moves away from \bar{x} .
- Do not extrapolate beyond the range of data as this is might be very inaccurate.

Multiple Regression

Assume $Y_i \sim N(\alpha + \beta_1 x_{i1} + \dots + \beta_k x_{ki}, \sigma^2)$ for $i = 1, \dots, n$ with Y_1, \dots, Y_n independent, i.e. the observations are independent, normally distributed, have constant variance and the expectations are linearly related to explanatory variables.

$$\mathbb{E}(Y) = \alpha + \beta_1 x_1 + \dots + \beta_k x_k$$

The least-squares estimates $\hat{\alpha}, \hat{\beta}_1, \dots, \hat{\beta}_k$ are values that minimise

$$S(\alpha, \beta_1, \dots, \beta_k) = \sum_{i=1}^n [y_i - (\alpha + \beta_1 x_{i1} + \dots + \beta_k x_{ki})]^2$$

$$S^2 = \frac{1}{n - (k+1)} \sum_{i=1}^n [Y_i - (\hat{\alpha} + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_k x_{ki})]^2$$

Confidence intervals:

$$\hat{\alpha} \pm t_{n-(k+1); 0.025} \times s.e.(\hat{\alpha})$$

$$\hat{\beta}_j \pm t_{n-(k+1); 0.025} \times s.e.(\hat{\beta}_j)$$

Residual sum of squares (rss): $\sum_{i=1}^n (Y_i - \hat{y}_i)^2$

F-test for Model Comparison

Used to see whether or not the full model gives a significantly better fit than a submodel.

H_0 : the specified regression coefficients are zero

H_1 : there is no restriction on the regression coefficients

Analysis of Variance

One-way ANOVA

Assume $Y_{ij} \sim N(\mu_i, \sigma^2)$ for $i = 1, \dots, k$ and $j = 1, \dots, n_i$ independently for all Y_{ij} , i.e. the observations are from a normal distribution, independent, have a common variance and a mean only dependent on the group they are member of.

$H_0 : \mu_1 = \dots = \mu_k$ vs $H_1 : \mu_1, \dots, \mu_k$ are not all equal.

| Source | d.f. | SS | MS | F | p |
|---------|---------|------------|--------|-----|-----|
| Between | $k - 1$ | SS_B | MS_B | F | p |
| Error | $n - k$ | SS_W | MS_W | | |
| Total | $n - 1$ | SS_{Tot} | | | |

In R, use `anova(lm())`. Need to express the explanatory variable using `as.factor`.

$$SS_{Tot} = SS_B + SS_W$$

Between groups mean square:

$$MS_B = \frac{SS_B}{k - 1}$$

Within groups mean square:

$$MS_W = \frac{SS_W}{n - k} = s^2 \quad (\text{residual mean square}),$$

where s is the *residual standard error*.

If H_0 is true then $F = \frac{MS_B}{MS_W} \sim F_{k-1, n-k}$.

Least Significant Differences (LSD)

$$t_{n-k; \alpha/2} \sqrt{s^2 \left(\frac{1}{n_i} + \frac{1}{n_j} \right)} \text{ or } t_{n-k; \alpha/2} \sqrt{\frac{2s^2}{m}}$$

if the samples are of equal size.

Two-way ANOVA

Assume $Y_{ij} \sim N(\mu_{ij}, \sigma^2)$ where $\mu_{ij} = \alpha_i + \beta_j$, i.e. the observations are from a normal distribution, independent, have a common variance and a mean that is a function of effect of each group.

Consider b blocks, k treatments, $n = bk$.

Test 1 (block effect):

$H_0 : \alpha_1 = \dots = \alpha_b$ vs $H_1 : \alpha_1, \dots, \alpha_b$ are not all equal

Test 2 (treatment effect):

$H_0 : \beta_1 = \dots = \beta_k$ vs $H_1 : \beta_1, \dots, \beta_k$ are not all equal

| Source | d.f. | SS | MS | F | p |
|-----------|------------------|------------|--------|-------|-------|
| Blocks | $b - 1$ | SS_B | MS_B | F_B | p_B |
| Treatment | $k - 1$ | SS_T | MS_T | F_T | p_T |
| Error | $(b - 1)(k - 1)$ | SS_W | MS_W | | |
| Total | $bk - 1$ | SS_{Tot} | | | |

$$SS_{Tot} = SS_B + SS_T + SS_W$$

$$MS_B = \frac{SS_B}{b - 1}$$

$$MS_W = \frac{SS_W}{(b - 1)(k - 1)}$$

$$F_B = \frac{MS_B}{MS_W}$$

$$MS_T = \frac{SS_T}{k - 1}$$

$$F_T = \frac{MS_T}{MS_W}$$

LSD for two-way ANOVA

Block effect:

$$t_{(b-1)(k-1); \alpha/2} \sqrt{2 \frac{s^2}{k}}$$

Treatment effect:

$$t_{(b-1)(k-1); \alpha/2} \sqrt{2 \frac{s^2}{b}}$$

Two-way ANOVA with r replications

| Source | d.f. | SS | MS | F | p |
|-----------|-------------------|------------|--------|-------|-------|
| Blocks | $b - 1$ | SS_B | MS_B | F_B | p_B |
| Treatment | $k - 1$ | SS_T | MS_T | F_T | p_T |
| Error | $rbk - b - k + 1$ | SS_W | MS_W | | |
| Total | $rbk - 1$ | SS_{Tot} | | | |

LSD for two-way ANOVA with replications

Block effect:

$$t_{rbk-b-k+1; \alpha/2} \sqrt{2 \frac{s^2}{rk}}$$

Treatment effect:

$$t_{rbk-b-k+1; \alpha/2} \sqrt{2 \frac{s^2}{rb}}$$