# Formulas statistics

z-score: $z_i = \dfrac{x_i - \bar{x}}{s}$

$r = \dfrac{\hat{\sum}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\hat{\sum}(x_i - \bar{x})^2 \hat{\sum}(y_i - \bar{y})^2}}$  correlation $[-1, +1]$  $\qquad r^2 = \dfrac{Var(\hat{y})}{Var(y)}$

$\hat{y} = a + bx$  $\hat{y}$: predicted value  $\qquad b = r\dfrac{s_y}{s_x} = \dfrac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}$  $\qquad a = \bar{y} - b\bar{x}$

residual sum of squares: $\hat{\sum}(y_i - \hat{y}_i)^2$

sampling distribution of population proportion: $N \sim \left(p, \left(\sqrt{\dfrac{p(1-p)}{n}}\right)^2\right)$

sampling distribution of sample mean: $N \sim \left(\mu, \left(\dfrac{\sigma}{\sqrt{n}}\right)^2\right)$  $\qquad \sigma$: population std  $\qquad s$: sampling std

$\sim N(\mu, \hat{\sigma}^2)$ by CLT for categorical  
if $n\hat{p}$ and $n(1-\hat{p}) \geqslant 15$  
$\sim N(\mu, \sigma^2)$ by CLT when $n > 30$ for quantitative

$CI$ of population proportion (estimate): $\hat{p} \pm \underset{d/2}{z}(se)$ ; $se = \sqrt{\dfrac{\hat{p}(1-\hat{p})}{n}}$  $\quad (1-d$ confidence level$)$

$CI$ of population mean (estimate): $\bar{x} \pm \underset{d/2}{z}(se)$ ; $se = \dfrac{s}{\sqrt{n}}$  with $n$ large or underlying normally distributed population

$\bar{x} \pm t_{df, d/2}(se)$ ; $se = \dfrac{s}{\sqrt{n}}$  $df = n-p$  with small $n$; normality assumption / with $n$ large: CLT

CI for two means:

two sample pooled t-interval  $\qquad (\bar{x} - \bar{y}) \pm t_{d/2, n+m-2}\, S_p \sqrt{\dfrac{1}{m} + \dfrac{1}{n}}$ ; $S_p^2$ is the pooled sample variance

$\qquad Sp^2 = \dfrac{(n-1)S_x^2 + (m-1)S_y^2}{m+n-2}$  (unbiased estimator of common $\sigma^2$)

two sample t-interval (different $\sigma^2$)  $\qquad (\bar{x} - \bar{y}) \pm z_{d/2} \sqrt{\dfrac{s_x^2}{n} + \dfrac{s_y^2}{m}}$

CI for difference in two population proportions  $\qquad (\hat{p}_1 - \hat{p}_2) \pm z_{d/2} \sqrt{\dfrac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \dfrac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$

## Significance test

| | proportions | means |
|---|---|---|
| 1. assumptions | categorical 0-1, sample size large $\sim$ normal | quantitative, population $\sim$ normal |
| 2. hypothesis | $H_0 : p = p_0$ ; $H_1 : p < p_0 / p \neq p_0 / p > p_0$ | $H_0 : \mu = \mu_0$ , $H_1 : \mu < \mu_0 / \mu \neq \mu_0 / \mu > \mu_0$ |
| 3. test statistic | $z_{obs} = (\hat{p} - p_0) / \sqrt{\dfrac{\hat{p}(1-p_0)}{n}}$ $\left[\dfrac{\hat{p} - p_0}{\sqrt{\hat{p}(1-p_0)}}\sqrt{n}\right]$ | $T_{obs} = (\bar{x} - \mu_0)/\sqrt{\dfrac{s^2}{n}}$ $\left[\dfrac{\bar{x} - \mu_0}{s}\sqrt{n}\right]$ |
| 4. p-value | $H_1 : p < p_0 \Rightarrow P(z < z_{obs})$ | $H_1 : \mu_0 < \mu_0 \Rightarrow P(z < b_{obs})$ |
| 5. conclusion | $H_1 : p \neq p_0 \Rightarrow P(|t| > |z_{obs}|) = 2 \cdot P(z > |z_{obs}|)$ | $H_1 : \mu_0 \neq \mu_0 \Rightarrow P(|t| > |T_{obs}|) = 2 \cdot P(z > |T_{obs}|)$ |
| reject $H_0$ if p-value $< d$ | $H_1 : p > p_0 \Rightarrow P(z > z_{obs})$ | $H_1 : \mu > \mu_0 \Rightarrow P(z > T_{obs})$ |

$\bar{y}$ mean of a sample  
$\mu ; \mu = E(y)$ mean of a population  $\qquad E(y) = \beta_0 + \beta_1 x$  $\qquad b_0 ; b_1$ in the sample

$Y = \beta_0 + \beta_1 x + \varepsilon$ ; $\varepsilon \sim N(0, \sigma^2)$  is the simple linear regression model

fitted model: $\hat{y}_i = b_0 + b_1 x_i$

$SSE = \hat{\sum} (y_i - \hat{y}_i)^2$  sum of squared errors (residuals)

$b_1 = \dfrac{\hat{\sum}(x_i - \bar{x})(y_i - \bar{y}_i)}{\hat{\sum}(x_i - \bar{x})^2}$  $\quad b_0 = \bar{y} - b_1\bar{x}$  estimates of $\beta_1$ and $\beta_0$  $\qquad b_1 \sim N(\beta_1, Var(b_1))$  $Var(b_1) = \dfrac{\sigma^2}{\hat{\sum}(x_i - \bar{x})^2}$

$SSE = \hat{\sum}(y_i - \hat{y}_i)^2$  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad b_0 \sim N(\beta_0, Var(b_0))$  $Var(b_0) = \sigma^2\left[\dfrac{1}{n} + \dfrac{\bar{x}^2}{\hat{\sum}(x_i - \bar{x})^2}\right]$

$MSE = \dfrac{SSE}{n-p}$

$S = \sqrt{MSE}$  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \dfrac{b_1 - \beta_1}{\sqrt{Var(b_1)}} \sim N(0,1)$  $\dfrac{b_1 - \beta_1}{S(b_1)} \sim t(n-p)$

$R^2 = \dfrac{SST - SSE}{SST} = \dfrac{SSR}{SST}$  $\quad r = \sqrt{R^2}$  $\qquad\qquad\qquad\qquad \dfrac{b_0 - \beta_0}{\sqrt{Var(b_0)}} \sim N(0,1)$  $\dfrac{b_0 - \beta_0}{S(b_0)} \sim t(n-p)$

$SST = \hat{\sum}(y_i - \bar{y})^2$

$SST = SSR + SSE \Rightarrow SSR = SST - SSE$

$SSR = \hat{\sum}(\hat{y}_i - \bar{y})^2$  $\qquad MSR = \dfrac{SSR}{p-1}$

Hypothesis testing

1. state null and alternative Hypothesis
2. test statistic
3. conclusion : critical value approach
           p-value

two-sided test

$H_0 : \beta = 0$ ; $H_1 : \beta \neq 0$

$$T = \frac{b}{S(b)} \sim t(n-p) \qquad \left[ \frac{b-\beta}{S(b)} \sim t(n-p) \right]$$

crit. val:
     if $|T_0| > t_{\alpha/2}(n-p)$ reject $H_0$

p-val:
     if p-value $< \alpha$ reject $H_0$
     p-val = $P(|T| > |T_0|)$

     (one sided test) : reject $H_0$ when $T_0 \gtrless t_{\alpha/2}(n-p)$ or $P(T \gtrless T_0) < \alpha$
     $H_1 : \beta \gtrless 0$

note: $p = 2$
(only $\beta_0$ and $\beta_1$)
simple linear regression

$(t\text{-value})^2 = (F\text{-value})$   simple linear
$F(1, n-1)$     regression

Confidence interval for slope $\beta_1$

$$b_1 \pm t_{\alpha/2}(n-2) \cdot S(b_1)$$
$$\hookrightarrow \sqrt{\frac{MSE}{\hat{\sum}(x_i - \bar{x})^2}}$$

CI for $E(y)$ :

$$\hat{y} \pm t_{(1-\alpha)/2}(n-2) \cdot S\{\hat{y}\}$$
$\hat{y} = b_0 + b_1 x$
$$S\{\hat{y}\} = \sqrt{MSE \cdot \left( \frac{1}{n} + \frac{(x-\bar{x})^2}{\hat{\sum}(x_i - \bar{x})^2} \right)}$$

PI for $y$ :   (prediction interval)

$$\hat{y} \pm t_{\alpha/2}(n-2) \sqrt{S^2\{\hat{y}\} + MSE}$$

Multiple linear regression

$E(Y) = X\beta$   regression equation
estimate of vector $\beta$ : $b = (X^T X)^{-1} X^T Y \sim MVN(\beta, \sigma^2(X^T X)^{-1})$
$\hat{y} = Xb$   fitted model

                      ┌→ residual
$Y = E(Y) + \varepsilon$   regression model

in hypothesis testing : $T = \frac{\text{sample coefficient}}{\text{std error of the coefficient}} \sim t(n-p)$ under $H_0$

$SSE = e^T e \left[ = (y - xb)^T (y - xb) \right] \quad \hat{\sum}(y_i - \hat{y}_i)^2$
$MSE = \frac{SSE}{n-p}$
$S = \sqrt{MSE}$
$SST = \sum (y_i - \bar{y})^2$
$SSR = SST - SSE$
$MSR = \frac{SSR}{p-1}$      $R^2 = \frac{SSR}{SSE}$

ANOVA table   (decomposition of variance)

| source | DF | SS | MS | F |
|---|---|---|---|---|
| regression | $p-1$ | SST-SSE | MSR | $\frac{MSR}{MSE}$ |
| error | $n-p$ | SSE | MSE | |
| total | $n-1$ | SST | | |

(significance test)

F-test   $H_0 : \beta_1 = \beta_2 = \ldots = \beta_{p-1} = 0$
        $H_1 :$ at least one $\beta_i \neq 0$ for $i = 1, \ldots, p-1$

$F = \frac{MSR}{MSE} \sim F(p-1, n-p)$ if $F_0 > F_\alpha(p-1, n-p)$ reject $H_0$

                  if $P(F > F_0) < \alpha$ reject $H_0$

t-test tests linear relation between $y$ and a certain $x_i$ while all other x-variables are in the model
F-test tests linear relation between $y$ and all x-variables together