

Naam	Formule	R-code	Extra
Populatieparameters toetsen			
Onbekende verwachting μ en onbekende populatie-variantie	$\frac{\bar{x} - \mu_0}{s / \sqrt{n}} \sim t_{n-1}$	<p>p-waarde: <code>t.test(x = vector, mu = ..., alternative = "two.sided/less/greater", conf.level = 0.95)</code></p> <p>power: verschil tussen 2 gemiddelden in een steekproef detecteren <code>power.t.test(n/power = ..., delta = $\mu_0 - \mu_a$, sd = ..., alternative = "two.sided/less/greater", sig.level = 0.05, type = "one.sample")</code></p>	Voorwaarden: <ul style="list-style-type: none"> - X is tenminste van intervalniveau - X is normaal verdeeld of de steekproef is groot
Onafhankelijke onbekende verwachtingen μ_1 en μ_2 en onbekende populatie-varianties die niet gelijk zijn aan elkaar	<p>p-waarde: Welch t-toets:</p> $\frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{S_1^2}{n} + \frac{S_2^2}{n}}} \sim t_l$ $l = \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}\right)^2}{\frac{S_1^4}{n_1^2(n_1 - 1)} + \frac{S_2^4}{n_2^2(n_2 - 1)}}$ <p>power: Effectsize d:</p> $d = \frac{\mu_1 - \mu_2}{\sqrt{s^2_{pooled}}}$ $s^2_{pooled} = \frac{SS_X + SS_Y}{n_1 + n_2 - 2}$	<p>p-waarde: <code>t.test(x = vector, y = vector, alternative = "two.sided/less/greater", conf.level = 0.95)</code></p> <p>power: verschil tussen 2 gemiddelden gedeeld door de gepoolde sd detecteren <code>pwr.t2n.test(n1/power = ..., n2/power = ..., d = effectsize, sig.level = 0.05, alternative = "less/greater/two.sided")</code></p> <p>deze functie kan enkel gebruikt worden indien $n_1 = n_2$ (en onafhankelijk): <code>power.t.test(delta = $\mu_1 - \mu_2$, sd = effectsize d, n/power = ..., sig.level = 0.05, alternative = "two.sided/one.sided", type = "two.sample")</code></p> <p>normaliteit van 2 onafhankelijke gemiddeldes nagaan: <code>qqnorm(x1)</code> <code>qqnorm(x2)</code></p>	Voorwaarden: <ul style="list-style-type: none"> - X is tenminste van intervalniveau - X is normaal verdeeld in beide populaties, of beide steekproeven zijn groot - De steekproeven zijn onafhankelijk <p>Benodigde package om power te berekenen: "pwr"</p> <p>d is negatief indien $\mu_1 - \mu_2$ negatief is</p> <p>Equivalent aan lineaire regressie indien de predictor dichotoom is ➔ Bij Welch: 2 vectoren creëren uit de twee levels van de dichotome variabele ➔ de gemiddelden van deze 2 vectoren ga je vergelijken</p>
Afhankelijke onbekende verwachtingen μ_1 en μ_2 en onbekende populatie-varianties	$d = \mu_1 - \mu_2$ $\frac{\bar{d}}{s / \sqrt{n}} \sim t_{n-1}$	<p>p-waarde: 1 steekproef van maken <code>t.test(x = d, mu = 0, alternative = "greater/less/two.sided")</code></p> <p><code>t.test(x = vector, y = vector, alternative = "two.sided/less/greater", paired = TRUE)</code></p> <p>power: <code>power.t.test(n/power = ..., delta = $\mu_1 - \mu_a$, sd = sd van delta, type = "paired", alternative = "two.sided/one.sided")</code></p>	Voorwaarden: <ul style="list-style-type: none"> - X is tenminste van intervalniveau - $X_1 - X_2$ moet normaal verdeeld zijn of de steekproef moet groot zijn <p>Bij afhankelijke steekproeven zijn de steekproefgroottes altijd even groot</p>

		<p>normaliteit van 2 afhankelijke gemiddeldes nagaan: 1 steekproef van maken <code>qqnorm(x1 - x2)</code></p>	
Onbekende proportie	$P(B(n, \pi) < k) = p\text{waarde}$ $\frac{n!}{k!(n-k)!} \pi^k (1-\pi)^{(n-k)} \sim B(n, \pi)$	<p>p-waarde: <code>binom.test(x = k, n = n, p = \pi, alternative = "one.sided/two.sided")</code></p> <p>power: <code>powerBinom(n/power = ..., p0 = ..., p1 = ..., sig.level = 0.5, alternative = "one.sided/two.sided")</code></p>	<p>Benodigde package om power te berekenen: "exactci"</p> <p>Als Ha tweezijdig is, dan moeten uitgerekende kansen vermenigvuldigd met 2 worden</p>
Onbekende B1	<p>p-waarde bij enkelvoudig lineair verband: Via t-verdeling:</p> $\frac{B_1}{\sqrt{\frac{SS_{Res}}{(n-2)SS_X}}} \sim t_{n-2}$ <p>Via F-verdeling: nulmodel en lineair model vergelijken</p> $\frac{SS_{Res0} - SS_{Res1}}{SS_{Res1}/(n-2)} \sim F_{1, n-2}$ <p>Waarbij SSRes0 gelijk is aan SSY</p> <p>p-waarde bij meervoudig lineair verband: Via F-verdeling: model A en model B vergelijken</p> $\frac{(SS_{ResA} - SS_{ResB}) / (df_A - df_B)}{SS_{ResB} / df_B} \sim F\left(\frac{df_A - df_B}{df_B}\right)$ <p>Waarbij: df_A: aantal vrijheidsgraden van model A (n-k-1) df_B: aantal vrijheidsgraden van model B (n-p-1) k: subset van predictoren p</p> <p>power bij enkelvoudig lineair verband: correlatie schatten om te gebruiken voor R-functie</p> $\rho = B_1 \frac{\sigma_x}{\sigma_y}$	<p>p-waarde bij enkelvoudig lineair verband: Via t-verdeling: <code>pt(q = ..., df = ..., lower.tail = FALSE)</code></p> <p><code>pf(q = ..., df1 = 1, df2 = n-2, lower.tail = FALSE)</code></p> <p>Via F-verdeling: toetst nulmodel vs lineair model <code>summary(lm)</code>: nagaan of een specifieke coëfficiënt 0 is</p> <p>Voor nominale predictor met <u>dummy/effect codering</u>: enkel F-toets interpreteren, t-toets is betekenisloos (t-toetsen zijn wél betekenisvol bij dichotome r-predictoren!!!) <code>summary(lm)</code></p> <p>p-waarde bij meervoudig lineair verband: Via F-verdeling: model A (k) en model B (p) vergelijken <code>anova(lmA, l_mB)</code></p> <p>Via F-verdeling: nulmodel en lineair model vergelijken <code>summary(lm)</code>: nagaan of de regressiecoëfficiënten van alle predictoren 0 zijn</p> <p>nulmodel creëren om te gebruiken in anova-functie: <code>lm0 <- lm(formula = y ~ NULL) → anova(lm0, lm)</code></p> <p>Voor 1 nominale predictor met <u>dummy/effect codering</u>: enkel F-toets interpreteren!!! <code>summary(lm)</code>: nulmodel vs lm, coëfficiënten en gemiddeldes te weten komen</p>	<p>Voorwaarden:</p> <ul style="list-style-type: none"> - X is tenminste interval, ratio of dichotoom - Y is tenminste interval of ratio - $E(\varepsilon_i) = 0$ voor alle i - $V(\varepsilon_i) = V(\varepsilon_j)$ voor alle i, j - $COV(\varepsilon_i, \varepsilon_j) = 0$ voor alle i, j - ε_i moet normaal verdeeld zijn of de steekproef moet groot zijn <p>F-verdeeld met $df_0 - df_1$ vrijheidsgraden in de teller en df_1 in de noemer</p> <p>F-toets is altijd eenzijdig: alfa niet delen door 2! Ongeacht de alternatieve hypothese</p> <p>Package inlezen om power te berekenen: "pwr"</p> <p>Volgorde van de argumenten in de functie anova: eerst het model met minder predictoren (LM), dan het model met meer predictoren (LM.geslacht)</p> <p>Opgelet voor collineariteit: package "car"</p> <ul style="list-style-type: none"> - VIFs = 1 → goed - VIFs < 3 → oké - VIFs < 10 en één VIF van < 3 → grijze zone - Één VIF > 10 → geen regressie

power bij meervoudig lineair verband: via F-verdeling het verschil tussen de determinatiecoëfficiënten (voor lm0 vs lm of lmA vs lmB)

$$f^2 = \frac{R_B^2 - R_A^2}{1 - R_B^2} \sim F_{p-k; n-p-k}$$

anova(lm, lmdummy): model zonder en model met dummycodering (die volgorde!), f-toets die nagaat of de extra variantie die het lm met dummy verklaart toevallig is

power bij enkelvoudig lineair verband: correlatiecoëfficiënt detecteren

`pwr.r.test(n/power = ..., r = rho, sig.level = 0.05)`

power bij meervoudig lineair verband: verschil van de determinatiecoëfficiënten van de 2 modellen detecteren (voor lm0 vs lm én lmA vs lmB)

`pwr.f2.test(u = p-k, f2 = f^2, power = ..., sig.level = 0.05)`

waarbij v = n-p-1 (gebruiken om n te berekenen)

predicties:
`fitted(lm)`

residuen:
`residuals(lm)`

SSMod:
`sum((fitted(lm)-mean(Y))^2)`

SSTot:
`var(y)*(n-1)`

Bij df_B bij dummycodering: elke hulpveranderlijke wordt gezien als een aparte predictor, en b0 is geen predictor

B1 voor Ha berekenen: als Y stijgt met 1, dan stijgt X met 4 = ¼

Bij alle soorten modelvergelijkingen of variatievergelijkingen: het model met het MINST predictoren wordt eerst in de formule/functie gestopt, want dit model heeft de grootste variantie (anders negatief)

Bij F- en chikwadraat verdelingen zoeken we altijd kansen in de vorm van $P(X>x)$ (omdat we geïnteresseerd zijn in extreme waarden, en beide verdelingen beginnen pas vanaf 0)

Equivalent aan Welch t-toets indien de predictor dichotoom is

Pearson chikwadraat toets

Zijn de p proporties in de k populaties identiek?

Zijn twee categorische variabelen afhankelijk of niet?

$$X^2 = \sum_{i=1}^k \sum_{j=1}^p \frac{(f_{i,j} - n_i \hat{\pi}_{.j})^2}{n_i \hat{\pi}_{.j}}$$

$$\sim \chi^2_{(p-1)(k-1)}$$

waarbij:
p = aantal categorieën van de categorische variabele
k = aantal populaties
 $\hat{\pi}_{.j}$ = schatting van de verwachte proporties in de categorie j indien de nulhypothese correct is

p-waarde en afhankelijkheid:
`chisq.test(table(variabele 1, variabele 2))`

power:
`power.chisq.test(w = effectgrootte, N = n, df = (p-1)(k-1), sig.level = 0.05)`

effectgrootte w:
`ES.w2(bivariate kansverdeling Ha)`

Altijd een tweezijdige toets: zijn de proporties identiek in de twee populaties?

Noemt ook de homogeniteitstoets

De vraag of twee categorische variabelen identiek zijn is equivalent aan het vragen of twee variabelen afhankelijk zijn

Benodigde package om de effectgrootte en de power te berekenen: "pwr"

Bij F- en chikwadraat verdelingen zoeken we altijd kansen in de vorm van $P(X>x)$ (omdat we

geïnteresseerd zijn in extreme waarden, en beide verdelingen beginnen pas vanaf 0)

Puntschattingen (incorrecte geheugensteuntjes)

Gemiddelde of verwachting	$E(X) = \bar{x}$ $SE_{E(X)} = \frac{\sigma}{\sqrt{n}}$	mean(...)	Voorwaarden: - X is tenminste interval, ratio of absoluut
Variantie	$V(X) = s_x^2$ $E(S_x^2) = \sigma_x^2$ $E(SN_x^2) = \frac{n-1}{n} \sigma_x^2$ $V(X) = \frac{SS_{Res}}{n-1}$	var(...) sd(...)	Voorwaarden: - X is tenminste interval, ratio of absoluut
Covariantie	$COV_{X,Y} = cov_{x,y}$	cov(..., ...)	Voorwaarden: - X is tenminste interval, ratio of absoluut
Correlatiecoëfficiënt	$\rho_{X,Y} = r_{x,y} = \frac{cov_x}{s_x s_y}$	cor(..., ...)	Voorwaarden: - X is tenminste interval, ratio of absoluut
B1	<p>Voor een enkelvoudig en meervoudig lineair model:</p> $\beta_1 = b_1$ $E(B_1) = \beta_1$ <p>Voor een enkelvoudig lineair model:</p> $V(B_1) = \frac{\sigma_\varepsilon^2}{SS_X} = \frac{\sigma_\varepsilon^2}{(n-1)s_x^2}$	lm(formula = y ~ x) summary(lm)	<p>Voorwaarden:</p> <ul style="list-style-type: none"> - X is tenminste interval, ratio of absoluut - $E(\varepsilon_i) = 0$ voor alle i - $V(\varepsilon_i) = V(\varepsilon_j)$ voor alle i,j - $COV(\varepsilon_i, \varepsilon_j) = 0$ voor alle i,j <p>Variantie van B1 verkleinen:</p> <ul style="list-style-type: none"> - n zo groot mogelijk - s_x^2 zo groot mogelijk - σ_ε^2 zo klein mogelijk <p>Opgelet voor collineariteit: package "car"</p> <ul style="list-style-type: none"> - VIFs = 1 → goed - VIFs < 3 → oké - VIFs < 10 en één VIF van < 3 → grijze zone - Één VIF > 10 → geen regressie
B0	<p>Voor een enkelvoudig en meervoudig lineair model:</p> $\beta_0 = b_0$	lm(formula = y ~ x) lm(formula = y ~ x1 + x2 + ...)	Voorwaarden: - X is tenminste interval, ratio of absoluut - $E(\varepsilon_i) = 0$ voor alle i

	$E(B_0) = \beta_0$ <p>Voor een enkelvoudig lineair model:</p> $V(B_0) = \sigma_\varepsilon^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2} \right)$ $= \sigma_\varepsilon^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{SS_X} \right)$	<p>coef(lm)</p> <p>summary(lm)</p>	<ul style="list-style-type: none"> - $V(\varepsilon_i) = V(\varepsilon_j)$ voor alle i, j - $COV(\varepsilon_i, \varepsilon_j) = 0$ voor alle i, j <p>Variantie van B_0 verkleinen:</p> <ul style="list-style-type: none"> - n zo groot mogelijk - s_x^2 zo groot mogelijk - σ_ε^2 zo klein mogelijk
Predictie Y_i	<p>Voor een enkelvoudig en meervoudig lineair model:</p> $\hat{Y}_i = b_0 + b_1 x_i$ $E(Y_i X_i) = b_0 + b_1 x_i$ <p>Voor een enkelvoudig lineair model:</p> $V(\hat{Y}_i) = \sigma_\varepsilon^2 \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{(n-1)s_x^2} \right)$ $= \sigma_\varepsilon^2 \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{SS_X} \right)$		<p>Voorwaarden:</p> <ul style="list-style-type: none"> - X is tenminste interval, ratio of absoluut - $E(\varepsilon_i) = 0$ voor alle i - $V(\varepsilon_i) = V(\varepsilon_j)$ voor alle i, j - $COV(\varepsilon_i, \varepsilon_j) = 0$ voor alle i, j <p>De variantie van de predictie van Y_i stijgt naarmate de predictie verder van het gemiddelde ligt → predicties zijn beter voor punten die dicht bij het gemiddelde liggen</p> <p>Variantie van Y_i verkleinen:</p> <ul style="list-style-type: none"> - n zo groot mogelijk - s_x^2 zo groot mogelijk - σ_ε^2 zo klein mogelijk
Variantie van de foutterm (residuen) σ_ε^2	<p>Voor een enkelvoudig lineair model:</p> $\sigma_\varepsilon^2 = s_\varepsilon^2 = \frac{1}{n-2} \sum_i^n (y_i - \hat{y}_i)^2 = \frac{SS_{Res}}{n-2}$ <p>Voor een meervoudig lineair model:</p> $S_\varepsilon^2 = \hat{\sigma}_\varepsilon^2 = \frac{1}{n-p-1} \sum_{i=1}^n (Y_i - B_0 - B_1 x_{i1} - \dots - B_p x_{ip})^2$ $= \frac{1}{n-p-1} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ $= \frac{SS_{Res}}{n-p-1}$	$((\text{sum}(\text{residuals}(\text{lm})^2))/n-1)$	<p>Voorwaarden:</p> <ul style="list-style-type: none"> - X is tenminste interval, ratio of absoluut - $E(\varepsilon_i) = 0$ voor alle i - $V(\varepsilon_i) = V(\varepsilon_j)$ voor alle i, j - $COV(\varepsilon_i, \varepsilon_j) = 0$ voor alle i, j
Betrouwbaarheidsintervallen			
BI voor het gemiddelde	$\left[\bar{x} \pm t_{n-1; \frac{\alpha}{2}} \frac{Sx}{\sqrt{n}} \right] \sim t_{n-1}$	<p>qt(p = 0.025, df = ..., lower.tail = FALSE)</p> <p>qt(p = 0.975, df = ...)</p>	<p>Voorwaarden:</p> <ul style="list-style-type: none"> - X is normaal verdeeld of de steekproef is groot

	<p>Breedte:</p> $2 \left(t_{n-1; \frac{\alpha}{2}} \frac{Sx}{\sqrt{n}} \right)$		<ul style="list-style-type: none"> - X is tenminste van intervalniveau <p>Kritieke t-waarde moet altijd positief zijn, dus kijken naar de rechterkant van de verdeling</p> <p>Breedte:</p> <ul style="list-style-type: none"> - Hoe groter sigma, hoe breder - Hoe groter n, hoe smaller - Hoe groter alfa, hoe smaller - Hoe kleiner de kritieke waarde, hoe breder
BI voor b1	<p>Voor enkelvoudige lineaire regressie:</p> $\left[b_1 \pm t_{n-2; \alpha/2} \sqrt{V(B_1)} \right]$ $= \left[b_1 \pm t_{n-2; \alpha/2} \sqrt{\frac{\widehat{\sigma}_\varepsilon^2}{(n-2) s_x^2}} \right]$ $= \left[b_1 \pm t_{n-2; \alpha/2} \sqrt{\frac{\widehat{\sigma}_\varepsilon^2}{SS_X}} \right]$ <p>Waarbij</p> $\widehat{\sigma}_\varepsilon^2 = \frac{SS_{Res}}{n - p - 1}$	<p>Voor enkelvoudige en meervoudige lineaire regressie: confint(lm, level = 0.95)</p>	<p>Voorwaarden:</p> <ul style="list-style-type: none"> - X is normaal verdeeld of de steekproef is groot - X is tenminste van intervalniveau - ε_i is normaal verdeeld - $E(\varepsilon_i) = 0$ voor alle i - $V(\varepsilon_i) = V(\varepsilon_j)$ voor alle i,j - $COV(\varepsilon_i, \varepsilon_j) = 0$ voor alle i,j <p>Smaller BI:</p> <ul style="list-style-type: none"> - n zo groot mogelijk - s_x^2 zo groot mogelijk - σ_ε^2 zo klein mogelijk
BI voor b0	<p>Voor enkelvoudige lineaire regressie:</p> $\left[b_0 \pm t_{n-2; \alpha/2} \sqrt{V(B_0)} \right]$ $= \left[b_0 \pm t_{n-2; \alpha/2} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SS_X}} \right]$	<p>Voor enkelvoudige en meervoudige lineaire regressie: confint(lm, level = 0.95)</p>	<p>Voorwaarden:</p> <ul style="list-style-type: none"> - X is normaal verdeeld of de steekproef is groot - X is tenminste van intervalniveau - ε_i is normaal verdeeld - $E(\varepsilon_i) = 0$ voor alle i - $V(\varepsilon_i) = V(\varepsilon_j)$ voor alle i,j - $COV(\varepsilon_i, \varepsilon_j) = 0$ voor alle i,j <p>Smaller BI:</p> <ul style="list-style-type: none"> - n zo groot mogelijk - s_x^2 zo groot mogelijk - σ_ε^2 zo klein mogelijk
Andere			
Determinatie-coëfficiënt R ²	Voor enkelvoudige en meervoudige lineaire regressie:	summary(lm)	Eigenschappen van R ² : <ul style="list-style-type: none"> - R² is nooit negatief

	$R^2 = \frac{SS_{Mod}}{SS_{Tot}} = \frac{SS_{Tot} - SS_{Res}}{SS_{Tot}}$ <p>waarbij:</p> $SS_Y = SS_{Tot} = \sum_{i=1}^n (y_i - \bar{y})^2$ $SS_{Mod} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ $SS_{Res} = \sum_{i=1}^n (\hat{y}_i - y_i)^2$ <p>Enkelvoudige lineaire regressie: populatie niveau</p> $\bar{R}^2 = 1 - (1 - R^2) \frac{n-1}{n-2}$ <p>Meervoudige lineaire regressie: populatie niveau</p> $\bar{R}^2 = 1 - (1 - R^2) \frac{n-1}{n-p-1}$	<p>predicties: fitted(lm)</p> <p>residuen: residuals(lm)</p> <p>SSMod: sum((fitted(lm)-mean(Y))^2)</p> <p>SSTot: var(y)*(n-1)</p>	<ul style="list-style-type: none"> - R² ligt altijd tussen 0 en 1 (het is een proportie) - Hoe hoger R², hoe sterker het lineair verband (geeft niet de richting van het verband aan) <p>R² = 1 → SSMod = SSTot</p> <p>R² = 0 → SSRes = SSTot</p>
--	--	---	---

Wanneer welke methode gebruiken?

1 variabele	<ol style="list-style-type: none"> 1) Dichotome variabele: exacte binomiale toets 2) Variabele met meerdere proporties: Pearson chikwadraat toets (homogeniteitstoets) 3) Variabele van minstens intervalniveau: one-sample t-toets
2 variabelen	<ol style="list-style-type: none"> 1) 2 nominale variabelen: Pearson chikwadraat toets (homogeniteitstoets) 2) 2 intervalniveau variabelen: lineaire regressie 3) 1 intervalniveau variabele en 1 dichotome variabele: lineaire regressie of Welch t-toets 4) 1 intervalniveau variabele en 1 nominale variabele: lineaire regressie (of ANOVA) 5) 1 variabele gemeten in 2 populaties: Welch t-toets (indien onafhankelijk) of two-sample paired t-toets