

④ and ⑤ p-value

$$\begin{aligned}
 \text{p-value} &: P(|t_{\text{stat}}| > t_{(n-1, \alpha)}^*) / H_0 \rightarrow \text{two-sided test} \\
 &= P(t_{\text{stat}} < -t_{(n-1)} | H_0) + P(t_{\text{stat}} > t_{(n-1)} | H_0) \\
 &= pt(q = 42, df = 399, \text{lower.tail} = F) \times 2 \\
 &= pt(q = -42, df = 399, \text{lower.tail} = T) + \\
 &\quad pt(q = 42, df = 399, \text{lower.tail} = F)
 \end{aligned}$$

⑥ Residual = observed - predicted

$$e_i, (i=1, \dots, 400) \quad e_i = y_i - \hat{y}_i$$

$$\approx \text{depths} \$ \text{alt} - 3628.5$$

vector of length 400

$$\text{residual standard error} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n e_i^2}$$

S = residual Std. Error

ONLY in the

intercept only model

$$= \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

$$= \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$$

= Standard deviation
(S)

$$= 1730 = \text{sd}(\text{depths} \$ \text{alt})$$

③ $H_0: \mu_0 = \mu_1$ or $\mu_1 - \mu_0 = 0$

$H_a: \mu_0 \neq \mu_1$ or $\mu_1 - \mu_0 \neq 0$

$$t_{\text{stat}} = \frac{(\bar{y}_1 - \bar{y}_0) - (\mu_1 - \mu_0)}{SE(\bar{y}_1 - \bar{y}_0)}$$

$$= \frac{211}{173} = 1.22$$

n_0 : sample size
in North

n_1 : sample size
in South

④ and ⑤

$$\underline{p\text{-value}} = P(|t_{\text{stat}}| > t^*_{(n_0-1+n_1-1)} | H_0)$$

$df = \overbrace{n-p}$

n : sample size

p : # of determinants + ① Interpret

$$pt(q=1.22 | df=398, \text{lower.tail}=F) \chi^2$$

⑥ Residual std. error.

$$= \sqrt{\frac{1}{n-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = 1730$$

② 95% CI for $\Delta\mu$

$$(\bar{y}_0 - \bar{y}_1) \pm 1.96 \times SE_{\bar{y}_0 - \bar{y}_1}$$

$$\rightarrow 211 \pm 1.96 \times 173$$

Example 3

$$\log(u) = \begin{cases} \log(u_0) & \text{if } \text{South}=0 \\ \log(u_0) + \log(\theta) & \text{if } \text{South}=1 \end{cases}$$

① (Intercept) Estimate = $\log(\bar{y}_0)$

\rightarrow an estimate of the parameter $\log(u_0)$

$$\log(\bar{y}_0) = 8.167 \Rightarrow \bar{y}_0 = \exp(8.167)$$

South Estimate: $\widehat{\log(\theta)} = 0.0581$

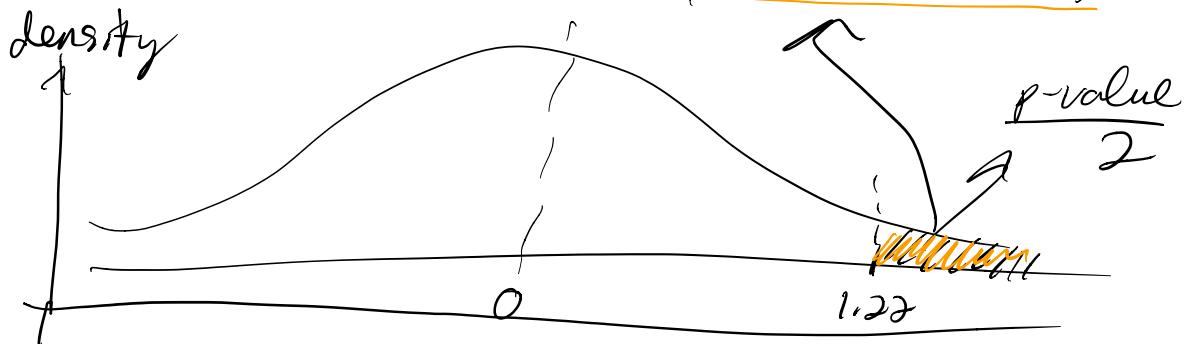
$$\begin{aligned}\widehat{\theta} &= \exp(0.0581) \\ &= 1.06\end{aligned}$$

$$\textcircled{2} \quad H_0: \log(\theta) = 0 \quad H_A: \log(\theta) \neq 0$$

$$t_{\text{test}} = \frac{\hat{\log}(\theta) - \log(\theta)}{SE_{\log(\theta)}} = \frac{0.0581 - 0}{0.0477} = 1.22$$

$$\textcircled{3} \text{ and } \textcircled{4} \quad p\text{-value} = P(|t_{\text{test}}| > t_{(398)}(H_A))$$

$$\underline{pt(q=1.22, df=398, \text{lower.tail}=F) \times 2}$$



\textcircled{5} 95\% CI:

$$0.0581 \pm qt(c(0.025, 0.975), df=398) \times 0.0477$$

$$= [\underbrace{\text{Lower}, \text{upper}}]$$

95% CI for the $\log(\theta)$

a 95% CI for θ

$$= [\exp(\text{lower}), \exp(\text{upper})]$$

This is incorrect:

$$\exp(\log(\theta)) \pm 1.96 \cdot \text{SE}_{\log(\theta)}$$

This is correct:

$$\exp \left[\overbrace{\log(\theta)}^{\wedge} \pm 1.96 \cdot \text{SE}_{\log(\theta)} \right]$$

⑦ Pseudo R^2

$$= 1 - \frac{\text{residual deviance} \rightarrow \text{for the model being fit}}{\text{null deviance} \rightarrow \text{measure of fit from an intercept only model}}$$

$\text{lm} \rightarrow$ it's always family = gaussian

$\text{glm} \rightarrow$ family = gaussian(link = identity)

`confint(glm object)`

\hookrightarrow 95% CI for $\log(\theta)$

`exp(confint(glm object))`



Poisson Regression Nov 6, 2020

Example 1

① Coefficient estimate for PT

is $\hat{\lambda}_0 = 3.39$ cases per infant year.

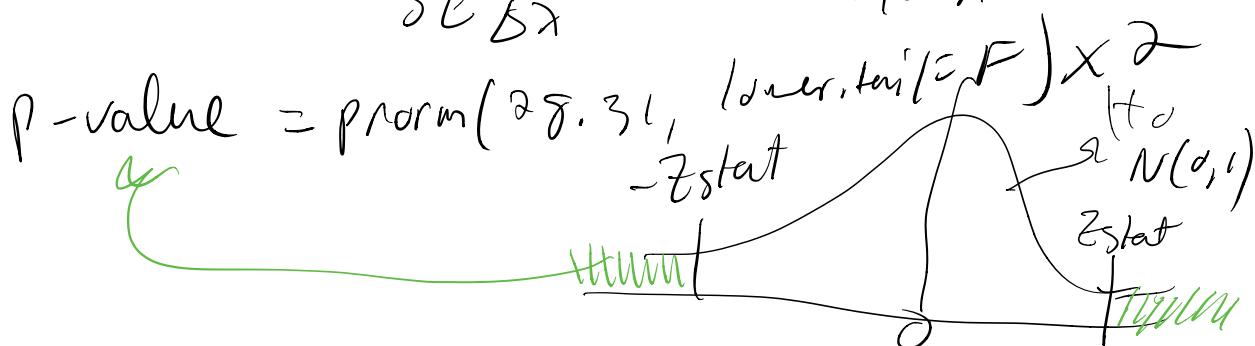
② Estimate for $\Delta\lambda \rightarrow$ rate difference

$$\begin{aligned} &= \hat{\lambda}_1 - \hat{\lambda}_0 \\ &= 0.31 \rightarrow \hat{\lambda}_1 = 3.7 \end{aligned}$$

Cases per
infant year.

③ $H_0: \Delta\lambda = 0$ $H_a: \Delta\lambda \neq 0$

$$z_{\text{stat}} = \frac{\hat{\Delta\lambda} - 0}{SE_{\hat{\Delta\lambda}}} = \frac{0.31 - 0}{0.010951} = 28.31$$



95% CI for $\Delta\lambda$ based on
Normal approximation

$$q_{\text{norm}} \left(p = c(0.025, 0.975) \right),$$

mean = 0.31,
sd = 0.010951

OR

$$0.31 \pm 1.96 \cdot \boxed{0.010951}$$

$\overline{SE}_{\Delta\lambda}$

confint(fit)

Nov 11, 2020

$$\log(A \cdot B) = \log(A) + \log(B)$$

$$\log(A^x) = x \log(A)$$

properties (rules)

$$\ln(\text{depth} \sim \text{South}) = \beta_0 + \beta_1 \cdot \text{South}$$

(Intercept) $\rightarrow \beta_0$

South $\rightarrow \beta_1$

Location	Depth	South
1	1000	0
2	200	1
:	:	0

$$M = \lambda_0 \cdot \theta^{NBF} \cdot PT \quad (1)$$

glm does not know how to fit this model.

Take "log" both sides of (1)

$$\log(M) = \log(\lambda_0 \cdot \theta^{NBF} \cdot PT)$$

property of
logs

$$= \log(\lambda_0) + \log(\theta^{NBF}) + \log(PT)$$

$$= \log(\lambda_0) + NBF \cdot \log(\theta) + \log(PT)$$

$$\log(\hat{\lambda}_d) = 1.220832 \Rightarrow \hat{\lambda}_d = \exp(1.22)$$

$$= 3.39$$

$$\log(\hat{\theta}) = 0.087505 \Rightarrow \hat{\theta} = \exp(0.087)$$

$$= 1.09$$

1.09 is the incidence rate ratio

the rate of respiratory infections among not breastfed infants is 9% higher than those who were breastfed.

95% CI for $\log(\theta)$

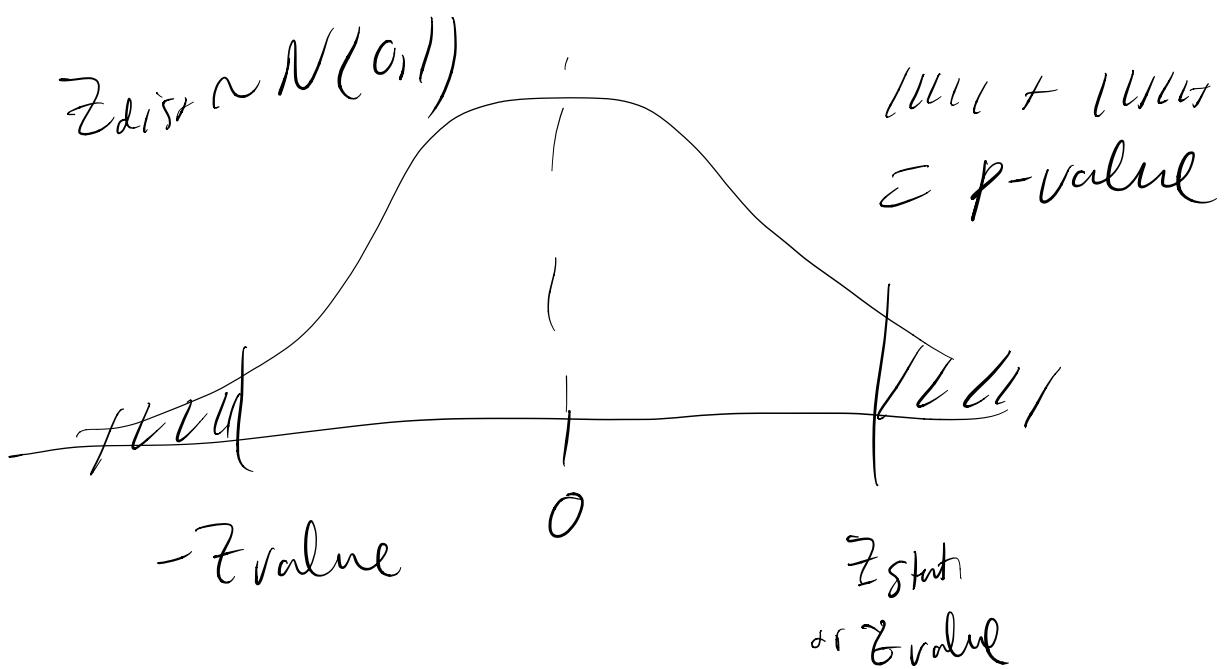
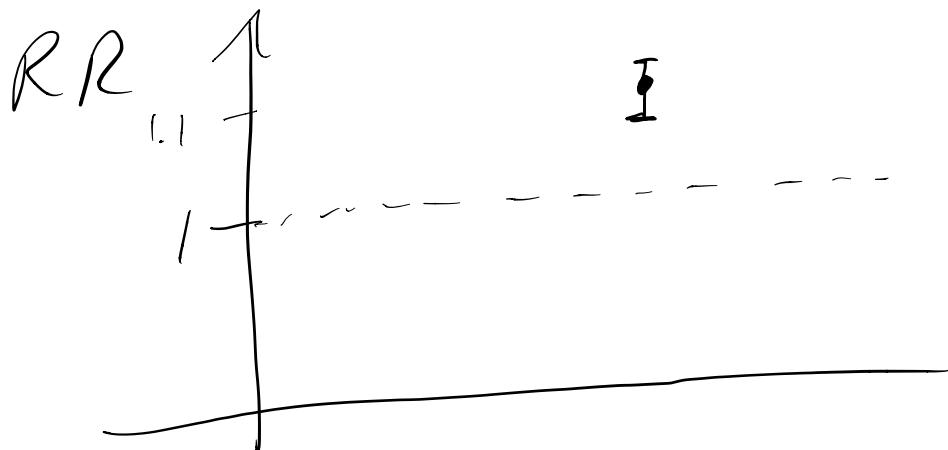
$$0.08 \pm 1.96 \cdot 0.003012$$

95% CI for θ

$$\exp(0.08 \pm 1.96 \cdot 0.003012)$$

$$\exp(\text{confint}(f_i))$$

$$= [1.085023, 1.0979]$$



$$= \text{pnorm}(\text{q} = \alpha \cdot q, \text{lower.tail} = \text{F}) \times 2$$

Bednets

expected # of cases of malaria = Rate \times PT

$$\mu = \boxed{\lambda} \times PT$$

λ model for rate

$$\lambda = \lambda_0 \cdot \theta^{\text{exposure}}$$

$$\theta = \frac{\lambda_1}{\lambda_0} \rightarrow \text{rate ratio}$$

Regression Equation

$$\mu = \lambda_0 \cdot \theta^{\text{exposure}} \times PT$$

$$\log(\mu) = \log(\lambda_0) + \text{exposure} \cdot \log(\theta) + \log(PT)$$

Estimates

$$\hat{\log}(\lambda_0) = 0.683 \Rightarrow \hat{\lambda}_0 = 1.98 \text{ cases/child year}$$

$$\hat{\log}(\theta) = -0.266 \Rightarrow \hat{\theta} = 0.765$$

Goodness of fit

we need to compare observed # of cases to expected # of cases.

$$\hat{M} = \hat{\lambda} \times PT = \hat{\lambda}_0 \cdot \hat{\theta}^{\text{exposure}} \times PT$$

$$\hat{M} = \begin{cases} 1.98 \times PT & \text{if exposed} = 0 \\ 1.98 \times 0.765 \times PT & \text{if exposed} = 1 \end{cases}$$

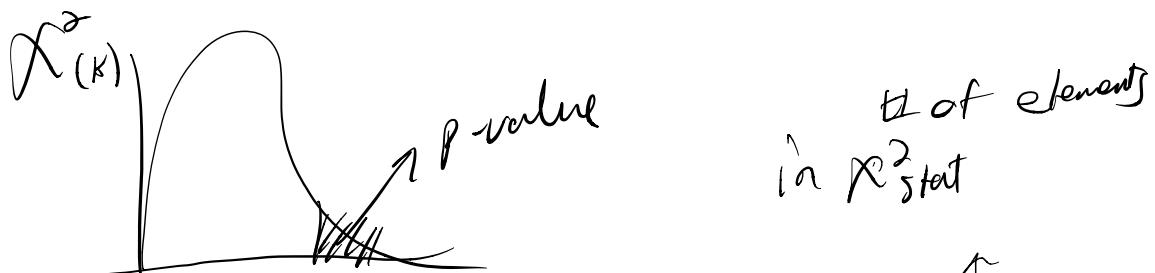
month	exposure	cases	years	expected
June 2014	0	33	79	1.98×79
July 2014	0	454	123	1.98×123
Aug 2014	0	204	103	1.98×103
Aug 2014	1	43	23	$(1.98 \times 0.765 \times 23)$
Sep 2014	0	177	79	
Sep 2014	1	66	39	

H_0 : no lack of fit H_A : lack of fit
 observed \approx expected

$$\chi^2_{\text{stat}} = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{(k-1)}$$

$$= \frac{(33 - 156.43)^2}{156.43} + \frac{(454 - 243)^2}{243} + \dots$$

compare to a chisq. distribution



$$p\text{-value} = p\text{chisq}\left(q = \chi^2_{\text{stat}}, df = k-1, \text{lower, tail of } F\right)$$

vectorized

Cases is a vector
 expect " " "

$$\text{chisq.stat} \leftarrow \text{sum} \left(\frac{(\text{obs} - \text{expected})^2}{\text{expected}} \right)$$

Nov 13, 2020

Age multipliers (M_{75}, M_{80}, M_{85})

Years	Age	Female	Male
2000-2004	70-74	1	1
	75-79	$0.0468/0.027 = 1.69$	$0.0822/0.052 = 1.58$
	80-84	$0.0808/0.027 = 2.97$	2.33
	85-89	$0.137/0.027 = 5.00$	3.49

2000-	70-74	1	1
2004	75-79	1.57	1.66
	80-84	2.58	2.66
	85-89	4.49	4.22

A best estimate for M_{75} might be
 mean $(1.69, 1.57, 1.58, 1.66) \rightarrow$ equal weight
 or to all cells.
 median $(1.69, 1.57, 1.58, 1.66)$

$$\lambda = \boxed{\lambda} \times PT$$

we need a model for lambda

$$Y: \text{count} \quad Y \sim \text{Poisson} (\lambda = \lambda \times PT)$$

$$\lambda = \lambda_0 \cdot M_{75}^{I_{75}} \times M_{80}^{I_{80}} \times M_{85}^{I_{85}} \times M_{90}^{I_{90}} \times M_m^{I_m}$$

I : binary indicator variable.

$$\log(\lambda) = \log(\lambda_0) + I_{75} \cdot \log(M_{75}) + \dots +$$

Age

"[70, 74]"

"[75-79]"

age \rightarrow factor variable

=

factor (age, levels = c("70-74", "75-79", ...))
labels =

obs	I_{75}	I_{80}	I_{90}	I_m	I_{20y}
1	0	0	0	0	0
2	1	0	0	0	0
3	1	1	0	0	

$\text{glm}(\text{cases} \sim I_{75} + I_{80} + I_{90} +$
 $I_m + I_{204} +$
 $\text{offset}(\log(\text{PT})) \mid$
 $\text{family} = \text{poisson}(\text{link} = \text{"log"})$)

$\text{glm}(\text{cases} \sim \text{age} + \text{gender} + \text{period}$
 $+ \text{offset}(\log(\text{PT})) \mid$)

levels(NHANES\$variable)