Análisis de encuestas de hogares con R

Andrés Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

Análisis de encuestas de hogares con R Modulo 6: Modelos lineales generalizados

Andrés Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

CEPAL - Unidad de Estadísticas Sociales

Análisis de encuestas de hogares con R Andrés Gutiérrez,

Ph.D. Stalyn Guerrero M.Sc.

Lectura de la base

Análisis de encuestas de hogares con R

```
encuesta <- readRDS("../Data/encuesta.rds")
data("BigCity", package = "TeachingSampling")</pre>
```

Definir diseño de la muestra con srvyr

```
Análisis de
encuestas de
hogares con R
```

```
library(srvyr)

diseno <- encuesta %>%
  as_survey_design(
    strata = Stratum,
    ids = PSU,
    weights = wk,
    nest = T
)
```

definir nuevas variables

```
Análisis de
encuestas de
hogares con R
```

```
diseno <- diseno %>% mutate(
   pobreza = ifelse(Poverty != "NotPoor", 1, 0),
   desempleo = ifelse(Employment == "Unemployed", 1, 0
```

Tablas de doble entrada para el tamaño

Análisis de

encuestas de hogares con R

Andrés

Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

	Sex	factor(pobreza)0	factor(pobreza)1	se1
Female	Female	48366	30824	2411
Male	Male	43032	28044	2522

30824

28044

(tab <- svytable(~pobreza + Sex, desig)

pobreza/Sex Female Male

0 48366 43032

Tablas de doble entrada para el proporción

Análisis de encuestas de hogares con R

```
(tab_pobreza_sexo <- svyby(~factor(pobreza), ~Sex,
    FUN = svymean, design = as.svrepdesign(diseno),
    se=F, na.rm=T, ci=T, keep.var=TRUE))</pre>
```

	Sex	factor(pobreza)0	factor(pobreza)1	se1	se2
Female	Female	0.6108	0.3892	0.0316	
Male	Male	0.6054	0.3946	0.0366	

```
prop.table(tab, margin = 2)
```

pobreza/Sex	Female	Male
0	0.6108	0.6054
1	0.3892	0.3946

Análisis de encuestas de hogares con R

> Andrés Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

$$\hat{\pi}_{rc} = \frac{n_{r+}}{n_{++}} \times \frac{n_{+c}}{n_{++}}$$

$$\chi^2_{pearsom} = n_{++} \times \sum_r \sum_c \left(\frac{\left(p_{rc} - \hat{\pi}_{rc} \right)^2}{\hat{\pi}_{rc}} \right)$$

$$G^2 = 2 \times n_{++} \times \sum_r \sum_c p_{cr} \times \ln \left(\frac{p_{rc}}{\hat{\pi}_{rc}} \right)$$

donde, R es el número de filas y C representa el número de columnas, la prueba tiene $(R-1)\times (C-1)$ grados de libertad.

Análisis de encuestas de hogares con R Andrés

Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

$$\chi^2_{(R-S)} = \chi^2_{(\textit{Pearson})}/\textit{GDEFF}$$

$$G_{(R-S)}^2 = G^2/GDEFF$$

con GDEFF el efecto generalizado del diseño, esta dado por

GDEFF =
$$\frac{\sum_{r} \sum_{c} (1 - p_{rc}) d^{2}(p_{rc}) - \sum_{r} (1 - p_{r+}) d^{2}(p_{r+}) - \sum_{c} (1 - p_{+c}) d^{2}(p_{r+})}{(R - 1)(C - 1)}$$

Análisis de encuestas de hogares con R Andrés Gutiérrez, Ph.D. Stalvn

> Guerrero M.Sc.

$$F_{R-S,Pearson} = \chi^2_{R-S} / \left[(R-1) (C-1) \right] \sim F_{(R-1)(C-1),(R-1)(C-1)d}$$

$$F_{R-S,LRT} = G_{R-S}^2 / (C-1) \sim F_{(C-1),df}$$

donde C es el número de columnas de la tabla cruzada

Análisis de encuestas de hogares con R

```
summary(tab, statistic = "Chisq")
```

```
## Sex
## pobreza Female Male
## 0 48366 43032
## 1 30824 28044
##
## Pearson's X^2: Rao & Scott adjustment
##
## data: NextMethod()
## X-squared = 0.077, df = 1, p-value = 0.8
```

Análisis de encuestas de hogares con R

```
summary(tab, statistic = "F")
```

```
## Sex
## pobreza Female Male
## 0 48366 43032
## 1 30824 28044
##
## Pearson's X^2: Rao & Scott adjustment
##
## data: NextMethod()
## F = 0.056, ndf = 1, ddf = 119, p-value = 0.8
```

Estadístico de Wald

Análisis de encuestas de hogares con R

> Andrés Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

$$Q_{wald} = \hat{oldsymbol{Y}}^t \left(oldsymbol{H} \hat{oldsymbol{V}} \left(\hat{oldsymbol{N}}
ight) oldsymbol{H}^t
ight)^{-1} \, \hat{oldsymbol{Y}}$$

donde,

$$\hat{\mathbf{Y}} = (\hat{\mathbf{N}} - \mathbf{E})$$

es un vector de $R \times C$ de diferencias entre los recuentos de celdas observadas y esperadas, esto es, $\hat{N}_{rc} - E_{rc}$

La matriz $\hat{HV}(\hat{N})H^t$, representa la matriz de varianza-covarianza estimada para el vector de diferencias.

Estadístico de Wald

Análisis de encuestas de hogares con R

Andrés Gutiérrez, Ph.D. Stalyn Guerrero M.Sc. La matriz \boldsymbol{H} es la inversa de la matriz \boldsymbol{J} dada por:

$$m{J} = -\left[rac{\delta^2 \ln PL\left(m{B}
ight)}{\delta^2 m{B}}
ight] \mid m{B} = \hat{m{B}}$$

$$\sum_{h}\sum_{a}\sum_{i}x_{hai}^{t}x_{hai}w_{hai}\hat{\pi}_{hai}\left(\boldsymbol{B}\right)\left(1-\hat{\pi}_{hai}\left(\boldsymbol{B}\right)\right)$$

Bajo la hipótesis nula, el estadístico

$$Q_{wald} \sim \chi^2_{(R-1)\times(C-1)}$$

Estadístico de Wald

Análisis de encuestas de hogares con R

$$F_{wald} = Q_{wald} imes rac{df - (R-1)(C-1) + 1}{(R-1)(C-1)df} \sim F_{(R-1)(C-1),df - (R-1)(C-1) + 1}$$

Análisis de encuestas de hogares con R

```
summary(tab, statistic = "Wald")
```

```
## Sex
## pobreza Female Male
## 0 48366 43032
## 1 30824 28044
##
## Design-based Wald test of association
##
## data: NextMethod()
## F = 0.056, ndf = 1, ddf = 119, p-value = 0.8
```

Prueba de independencia adjWald

Análisis de encuestas de hogares con R

```
summary(tab, statistic = "adjWald")
```

```
## Sex
## pobreza Female Male
## 0 48366 43032
## 1 30824 28044
##
## Design-based Wald test of association
##
## data: NextMethod()
## F = 0.056, ndf = 1, ddf = 119, p-value = 0.8
```

Prueba de independencia lincom

Análisis de encuestas de hogares con R

```
summary(tab, statistic = "lincom")
```

```
## Sex
## pobreza Female Male
## 0 48366 43032
## 1 30824 28044
##
## Pearson's X^2: asymptotic exact distribution
##
## data: NextMethod()
## X-squared = 0.077, p-value = 0.8
```

Prueba de independencia saddlepoint

Análisis de encuestas de hogares con R

```
summary(tab, statistic = "saddlepoint")
```

```
## Sex
## pobreza Female Male
## 0 48366 43032
## 1 30824 28044
##
## Pearson's X^2: saddlepoint approximation
##
## data: NextMethod()
## X-squared = 0.077, p-value = 0.8
```

Análisis de encuestas de hogares con R

> Andrés Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

$$\log(p_{ijk}) = \mu + \lambda_i^X + \lambda_j^Y + \lambda_k^Z + \lambda_{ij}^{XY},$$

donde:

lacksquare $p_{ijk}=$ la proporción esperada en la celda bajo el modelo.

Análisis de encuestas de hogares con R

> Andrés Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

$$\log(p_{ijk}) = \mu + \lambda_i^X + \lambda_j^Y + \lambda_k^Z + \lambda_{ij}^{XY},$$

donde:

- lacksquare $p_{ijk} =$ la proporción esperada en la celda bajo el modelo.
- $\blacksquare \ \mu = \log(p_0) = \frac{1}{\# \ de \ celdas}$

```
Análisis de
encuestas de
hogares con R
```

```
Andrés
Gutiérrez,
Ph.D.
Stalyn
Guerrero
M.Sc.
```

```
mod1 <- svyloglin(~pobreza+Sex + pobreza:Sex , diseno
(s1 <- summary(mod1))</pre>
```

```
## Loglinear model: svyloglin(~pobreza + Sex + pobrez

## coef se p

## pobreza1 0.219673 0.06778 0.001192

## Sex1 0.052843 0.01625 0.001145

## pobreza1:Sex1 0.005583 0.02350 0.812175
```

```
Análisis de
encuestas de
hogares con R
Andrés
Gutiérrez,
```

```
Ph.D.
Stalyn
Guerrero
M.Sc.

M.Sc.

M.Sc.

Mod2 <- svyloglin(~pobreza+Sex, diseno)
(s2 <- summary(mod2))
```

```
## Loglinear model: svyloglin(~pobreza + Sex, diseno)
## coef se p
## pobreza1 0.21997 0.06752 0.0011230
## Sex1 0.05405 0.01577 0.0006076
```

```
Análisis de
encuestas de
hogares con R
```

```
anova(mod1, mod2)
```

```
## Analysis of Deviance Table
## Model 1: y ~ pobreza + Sex
## Model 2: y ~ pobreza + Sex + pobreza:Sex
## Deviance= 0.07719 p= 0.8126
## Score= 0.07719 p= 0.8126
```

Modelo de regresión logistica

 $g(\pi(x)) = logit(\pi(x)) = z = ln\left(\frac{\pi(x)}{1 - \pi(x)}\right) = B_0 + B_1x_1 + \dots + B_p$

 $\hat{\pi}\left(\boldsymbol{x}\right) = \frac{\exp\left(\boldsymbol{X}\hat{\boldsymbol{B}}\right)}{1 - \exp\left(\boldsymbol{X}\hat{\boldsymbol{B}}\right)} = \frac{\exp\left(\hat{B}_0 + \hat{B}_1x_1 + \dots + \hat{B}x_p\right)}{1 - \exp\left(\hat{B}_0 + \hat{B}_1x_1 + \dots + \hat{B}x_p\right)}$

 $PL(\mathbf{B} \mid X) = \prod_{i=1}^{n} \left\{ \pi(x_{i})^{y_{i}} (1 - \pi(x_{i}))^{1 - y_{i}} \right\}^{w_{i}}$

 $\pi(x_i) = \frac{\exp(x_i \mathbf{B})}{1 - \exp(x_i \mathbf{B})}$

 $\operatorname{var}\left(\hat{oldsymbol{B}}
ight) = oldsymbol{J}^{-1}\operatorname{var}\left(S\left(\hat{oldsymbol{B}}
ight)
ight)oldsymbol{J}^{-1}$

encuestas de hogares con R Andrés

Análisis de

Gutiérrez, Ph.D. Stalvn

Guerrero

M.Sc.

con

Modelo de regresión logistica

Análisis de encuestas de hogares con R Andrés

$$S\left(B\right) = \sum_{h} \sum_{a} \sum_{i} w_{hai} \boldsymbol{D}_{hai}^{t} \left[\left(\pi_{hai}\left(\boldsymbol{B}\right)\right) \left(1 - \pi_{hai}\left(\boldsymbol{B}\right)\right) \right]^{-1} \left(y_{hai} - \pi_{hai}\left(\boldsymbol{B}\right)\right) = 0$$

$$D_{hai} = rac{\delta \left(\pi_{hai} \left(oldsymbol{B}
ight)
ight)}{\delta B_{i}}$$

$$j=0,\ldots,p$$

Prueba de Wald para los parámetros del modelo

Análisis de encuestas de hogares con R

$$G = -2 \ln \left[\frac{L \left(\hat{\beta}_{MLE} \right)_{reduced}}{L \left(\hat{\beta}_{MLE} \right)_{full}} \right]$$

$$\hat{\psi} = \exp\left(\hat{B}_1\right)$$

$$CI(\psi) = \exp\left(\hat{B}_{j} \pm t_{df,1-\frac{\alpha}{2}}se\left(\hat{B}_{j}\right)\right)$$

Tablas de contingencia

Análisis de
encuestas de
hogares con R
Andrés
Gutiérrez,
Ph.D.
Stalyn

Guerrero M.Sc.

Sex	pobreza	se	ci l	ci u
 Female	0.3892	0.0316	0.3273	0.4512
Male	0.3946	0.0366	0.3228	0.4664
Zone	pobreza	se	ci_l	ci_u
Rural	0.4485	0.0561	0.3386	0.5585
Urban	0.3394	0.0320	0.2766	0.4022

Region	pobreza	se	ci_l	ci_u
Norte	0.3590	0.0555	0.2502	0.4677
Sur	0.3438	0.0435	0.2586	0.4291
Centro	0.3654	0.0786	0.2113	0.5195
Occidente	0.4008	0.0467	0.3092	0.4924
Orionto	0.4519	0.0006	0.2791	0.6255

Análisis de encuestas de hogares con R

```
## Pearson's X^2: Rao & Scott adjustment
pobreza_sex <- svychisq(
   formula = ~pobreza + Sex, design = diseno)
tidy( pobreza_sex) %>% select(-method)
```

ndf	ddf	statistic	p.value
1	119	0.0565	0.8126

```
Análisis de
encuestas de
hogares con R
```

```
pobreza_Zona <- svychisq(
   formula = ~pobreza + Zone, design = diseno)
tidy(pobreza_Zona) %>% select(-method)
```

ndf	ddf	statistic	p.value
1	119	2.954	0.0883

```
Análisis de
encuestas de
hogares con R
```

```
pobreza_Region <- svychisq(
  formula = ~pobreza + Region, design = diseno)
  tidy(pobreza Region) %>% select(-method)
```

ndf	ddf	statistic	p.value
3.008	358	0.4879	0.6914

Modelo log lineal ajustado

Análisis de encuestas de hogares con R

```
mod_loglin <- svyglm(
  pobreza ~ Sex + Zone + Region,
    family=quasibinomial, design=diseno)
tidy(mod_loglin)</pre>
```

term	estimate	std.error	statistic	p.value
(Intercept)	-0.4082	0.2640	-1.5464	0.1248
SexMale	0.0086	0.0915	0.0945	0.9249
ZoneUrban	-0.4378	0.2418	-1.8106	0.0729
RegionSur	0.0063	0.3140	0.0201	0.9840
RegionCentro	0.1915	0.4279	0.4476	0.6553
RegionOccidente	0.2319	0.3144	0.7377	0.4622
RegionOriente	0.3699	0.4259	0.8686	0.3869

Plot de la distribución de los betas

Análisis de encuestas de hogares con R Andrés Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

```
plot_summs(mod_loglin,
                scale = TRUE, plot.distributions = TRUE)
   SeyMale
  ZoneUrban
  RegionSur
  RegionCentro
 RegionOccidente
  RegionOriente
```

Modelo log lineal ajustado

Análisis de encuestas de hogares con R

```
bind_cols(
data.frame(exp_estimado = exp(coef(mod_loglin))),
as.data.frame(exp(confint(mod_loglin)))
)
```

	exp_estimado	2.5 %	97.5 %
(Intercept)	0.6648	0.3941	1.122
SexMale	1.0087	0.8414	1.209
ZoneUrban	0.6454	0.3997	1.042
RegionSur	1.0063	0.5402	1.875
RegionCentro	1.2111	0.5188	2.827
RegionOccidente	1.2611	0.6764	2.351
RegionOriente	1.4476	0.6226	3.366

Estadístico de Wald sobre los parámetros

```
Análisis de
          regTermTest(model = mod_loglin, ~Sex)
encuestas de
hogares con R
 Andrés
          ## Wald test for Sex
 Gutiérrez,
 Ph.D.
          ## in svyglm(formula = pobreza ~ Sex + Zone + Region
 Stalvn
 Guerrero
          ##
                  family = quasibinomial)
 M Sc
          ## F = 0.00893 on 1 and 113 df: p = 0.92
          regTermTest(model = mod_loglin, ~Zone)
          ## Wald test for Zone
          ## in svyglm(formula = pobreza ~ Sex + Zone + Region
          ##
                  family = quasibinomial)
          ## F = 3.278 on 1 and 113 df: p = 0.073
          regTermTest(model = mod_loglin, ~Region)
```

Wald test for Region

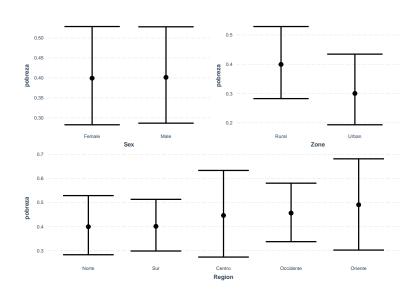
Efecto del modelo.

```
Andrés
Gutiérrez,
Ph.D.
Stalvn
```

Guerrero M Sc

Análisis de encuestas de

Análisis de encuestas de hogares con R



Modelo log lineal ajustado con interacciones

Andrés Gutiérrez, Ph.D. Stalyn Guerrero

M Sc

term

Análisis de encuestas de

				'
ZoneUrban	-0.4248	0.2562	-1.6580	0.100
(Intercept)	-0.4289	0.2849	-1.5055	0.135
SexMale:RegionSur	0.2871	0.2774	1.0348	0.303
RegionOriente	0.3843	0.4279	0.8980	0.371
RegionOccidente	0.3342	0.3783	0.8835	0.379
SexMale:RegionOccidente	-0.2302	0.2868	-0.8026	0.424

estimate

std.error statistic

p.valu

Modelo log lineal ajustado con interacciones

Análisis de encuestas de hogares con R

> Andrés Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

tab_mod %>% slice(7:12)

term	estimate	std.error	statistic	p.value
RegionCentro	0.2466	0.4560	0.5408	0.5897
SexMale:RegionCentro	-0.1162	0.2791	-0.4162	0.6781
RegionSur	-0.1325	0.3464	-0.3825	0.7028
SexMale	0.0478	0.1994	0.2399	0.8109
SexMale:RegionOriente	-0.0304	0.2878	-0.1057	0.9161
SexMale:ZoneUrban	-0.0154	0.1872	-0.0824	0.9345

Plot de la distribución de los betas

Análisis de plot_summs(mod_loglin_int, mod_loglin, scale = TRUE, encuestas de hogares con R Andrés Gutiérrez, SexMale Ph.D. Stalyn Guerrero Zonel Irban M.Sc. RegionSur RegionCentro RegionOccidente RegionOriente SexMale:ZoneUrban SexMale:RegionSur SexMale:RegionCentro SexMale:RegionOccidente

SexMale:RegionOriente

Model Model 1

Model 2

Modelo log lineal ajustado

Análisis de encuestas de hogares con R

	exp_estimado	2.5 %	97.5 %
(Intercept)	0.6512	0.3702	1.145
SexMale	1.0490	0.7065	1.557
ZoneUrban	0.6539	0.3935	1.087
RegionSur	0.8759	0.4408	1.740
RegionCentro	1.2797	0.5183	3.160
RegionOccidente	1.3968	0.6599	2.957
RegionOriente	1.4685	0.6288	3.430
SexMale:ZoneUrban	0.9847	0.6795	1.427
SexMale:RegionSur	1.3325	0.7689	2.309
SexMale:RegionCentro	0.8903	0.5120	1.548
SexMale:RegionOccidente	0.7944	0.4499	1.403
SexMale:RegionOriente	0.9701	0.5484	1.716

Estadístico de Wald sobre los parámetros

```
regTermTest(model = mod_loglin_int, ~Sex)
## Wald test for Sex
##
       design = diseno, family = quasibinomial)
## F = 0.05753 on 1 and 108 df: p = 0.81
regTermTest(model = mod_loglin_int, ~Zone)
## Wald test for Zone
```

```
##
```

F = 2.749 on 1 and 108 df: p = 0.1

Wald test for Region

regTermTest(model = mod_loglin_int, ~Region)

```
in svyglm(formula = pobreza ~ Sex + Zone + Region
##
       design = diseno, family = quasibinomial)
```

in svyglm(formula = pobreza ~ Sex + Zone + Region

Análisis de

encuestas de hogares con R Andrés

> Gutiérrez, Ph.D.

Stalvn Guerrero

M Sc

Estadístico de Wald sobre los parámetros

Análisis de encuestas de

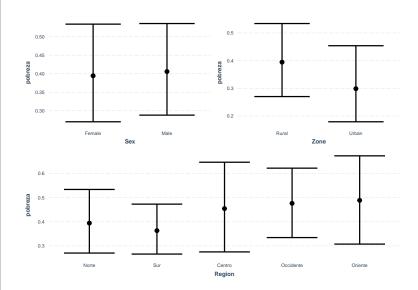
```
regTermTest(model = mod_loglin_int, ~Sex:Zone)
hogares con R
 Andrés
 Gutiérrez,
 Ph D
          ## Wald test for Sex:Zone
 Stalvn
 Guerrero
              in svyglm(formula = pobreza ~ Sex + Zone + Region
          ##
 M Sc
          ##
                 design = diseno, family = quasibinomial)
          ## F = 0.006789 on 1 and 108 df: p= 0.93
          regTermTest(model = mod_loglin_int, ~Sex:Region)
          ## Wald test for Sex:Region
              in svyglm(formula = pobreza ~ Sex + Zone + Region
          ##
                 design = diseno, family = quasibinomial)
          ##
          ## F = 1.058 on 4 and 108 df: p = 0.38
```

```
Andrés
Gutiérrez,
Ph.D.
Stalyn
Guerrero
```

M Sc

Análisis de encuestas de

Análisis de encuestas de hogares con R



Modelo log lineal ajustado con Q_Weighting

```
encuestas de
hogares con R
Andrés
Gutiérrez,
Ph.D.
```

Análisis de

weights = wk2,

nest = T

Modelo log lineal ajustado con Q_Weighting

```
Análisis de
encuestas de
hogares con R
```

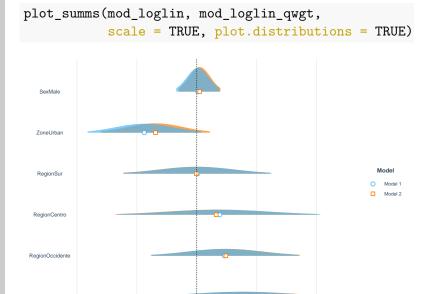
Modelo log lineal ajustado con Q_Weighting

Análisis de encuestas de hogares con R

term	estimate	std.error	statistic	p.value
(Intercept)	-0.4644	0.2630	-1.7656	0.0802
SexMale	0.0241	0.0883	0.2726	0.7857
ZoneUrban	-0.3445	0.2311	-1.4903	0.1389
RegionSur	-0.0041	0.3116	-0.0130	0.9896
RegionCentro	0.1613	0.4270	0.3778	0.7063
RegionOccidente	0.2424	0.3147	0.7705	0.4426
RegionOriente	0.3937	0.4319	0.9115	0.3639

Plot de la distribución de los betas

Análisis de encuestas de hogares con R



Modelo log lineal ajustado

Análisis de encuestas de hogares con R

	exp_estimado	2.5 %	97.5 %
(Intercept)	0.6285	0.3732	1.058
SexMale	1.0244	0.8600	1.220
ZoneUrban	0.7086	0.4482	1.120
RegionSur	0.9960	0.5371	1.847
RegionCentro	1.1750	0.5043	2.738
RegionOccidente	1.2744	0.6832	2.377
RegionOriente	1.4824	0.6301	3.488

Estadístico de Wald sobre los parámetros

Análisis de

encuestas de

Andrés

Gutiérrez, Ph.D.

Stalvn Guerrero

M Sc

```
regTermTest(model = mod_loglin_qwgt, ~Sex)
hogares con R
         ## Wald test for Sex
         ## in svyglm(formula = pobreza ~ Sex + Zone + Region
         ##
                family = quasibinomial)
         ## F = 0.0743 on 1 and 113 df: p = 0.79
         regTermTest(model = mod_loglin_qwgt, ~Zone)
         ## Wald test for Zone
         ## in svyglm(formula = pobreza ~ Sex + Zone + Region
         ##
                family = quasibinomial)
```

F = 2.221 on 1 and 113 df: p = 0.14

Wald test for Region

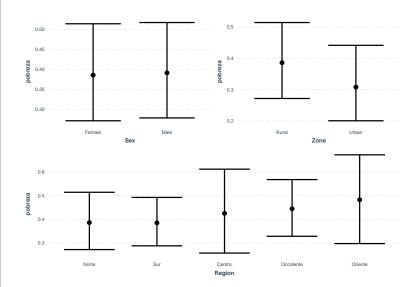
regTermTest(model = mod_loglin_qwgt, ~Region)

```
Andrés
Gutiérrez,
Ph.D.
Stalyn
Guerrero
```

M Sc

Análisis de encuestas de

Análisis de encuestas de hogares con R



¡Gracias!

Análisis de encuestas de hogares con R

> Andrés Gutiérrez, Ph.D. Stalyn Guerrero M.Sc.

> > Email: andres.gutierrez@cepal.org