

# Análisis de encuestas de hogares con R

## Módulo 4: Modelos de regresión

CEPAL - Unidad de Estadísticas Sociales

# Tabla de contenidos I

Modelos de regresión bajo diseños de muestreo complejos

Diagnostico del modelo

Modelos de regresión bajo diseños de muestreo complejos

# Introducción

- ▶ Un modelo matemático es una relación funcional entre variables.
- ▶ El objetivo es encontrar modelos que relacionen variables de entrada con una variable de salida.
- ▶ A lo largo de la historia, varios autores han discutido el impacto de los diseños muestrales complejos en las inferencias relacionadas con modelos de regresión.

# Introducción

- ▶ **Kish y Frankel (1974):** Fueron los primeros en abordar, de manera empírica, cómo los diseños muestrales complejos afectan las inferencias en modelos de regresión.
- ▶ **Fuller (1975):** Desarrolló un estimador de varianza que considera ponderaciones desiguales de observaciones, especialmente relevantes en contextos de muestreo complejo de dos etapas.
- ▶ **Sha et al. (1977):** Discutieron las violaciones de supuestos en modelos de regresión lineal y presentaron evaluaciones empíricas del desempeño de estimadores de varianza basados en la linealización para modelos de regresión lineal con datos de encuestas.
- ▶ **Binder (1983):** Se centró en las distribuciones muestrales de estimadores para parámetros de regresión en poblaciones finitas y estimadores de varianza relacionados.

# Introducción

- ▶ **Skinner et al. (1989):** Trabajaron en estimadores de varianza para los coeficientes de regresión que permitieron diseños de muestras complejas, y recomendaron el uso de métodos de linealización u otros métodos para la estimación de la varianza.
- ▶ **Fuller (2002):** Ofreció un resumen de los métodos de estimación para modelos de regresión que involucran información relacionada con muestras complejas.
- ▶ **Pfeffermann (2011):** Discutió enfoques basados en el ajuste de modelos de regresión lineal a datos de encuestas de muestras complejas, respaldando el uso de un método “q-weighted.”

# Modelos de Regresión Lineal Simple y Múltiple

- ▶ Un modelo de regresión lineal simple se define como

$$y = \beta_0 + \beta_1 x + \varepsilon$$

.

- ▶ Los modelos de regresión lineal múltiples extienden este concepto para múltiples variables predictoras:

$$y = X\beta + \varepsilon$$

.

- ▶ El valor esperado de la variable dependiente condicionado a las variables independientes se representa como  $E(y|x)$ .

# Consideraciones en Modelos de Regresión

- ▶  $E(\varepsilon_i|x_i) = 0$ : El valor esperado de los residuos condicionado a las covariables es igual a 0.
- ▶  $Var(\varepsilon_i|x_i) = \sigma_{y,x}^2$ : Homogeneidad de varianza, la varianza de los residuos condicionados es constante.
- ▶  $\varepsilon_i|x_i \sim N(0, \sigma_{y,x}^2)$ : Normalidad en los errores, los residuos condicionados se distribuyen normalmente.
- ▶  $cov(\varepsilon_i, \varepsilon_j|x_i, x_j)$ : Independencia en los residuos, los residuos en diferentes sujetos no están correlacionados con los valores de sus variables predictoras.



## Resultados para el modelo de regresión

Una vez definido el modelo de regresión lineal y sus supuestos, se puede deducir los siguiente:

$$\begin{aligned}\hat{y} &= E(y \mid x) \\ &= E(x\beta) + E(\varepsilon) \\ &= x\beta + 0 \\ &= \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p\end{aligned}$$

y Adicionalmente,

$$\begin{aligned}var(y_i \mid x_i) &= \sigma_{y,x}^2, \\ cov(y_i, y_j \mid x_i, x_j) &= 0 \\ &\text{y} \\ y_i &\sim N(x_i\beta, \sigma_{y,x}^2)\end{aligned}$$

# Estimación de los parámetros en un modelo de regresión simple.

La estimación del coeficiente de regresión  $\beta_1$  en un modelo de regresión simple con muestras complejas involucra el uso de ponderaciones y totales. El estimador  $\hat{\beta}_1$  se calcula como un cociente de totales ponderados.

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_h^H \sum_{\alpha}^{a_h} \sum_{i=1}^{n_{h\alpha}} \omega_{h\alpha i} (y_{h\alpha i} - \bar{y}_{\omega}) (x_{h\alpha i} - \bar{x}_{\omega})}{\sum_h^H \sum_{\alpha}^{a_h} \sum_{i=1}^{n_{h\alpha}} \omega_{h\alpha i} (x_{h\alpha i} - \bar{x}_{\omega})^2} \\ &= \frac{t_{xy}}{t_x^2}\end{aligned}$$

## Varianza estimada

La varianza del estimador  $\hat{\beta}_1$  se calcula considerando la varianza de los totales ponderados y sus covarianzas. Esta varianza estimada tiene en cuenta el diseño muestral y la estructura de ponderación.

$$var(\hat{\beta}_1) = \frac{var(t_{xy}) + \hat{\beta}_1^2 var(t_{x^2}) - 2\hat{\beta}_1 cov(t_{xy}, t_{x^2})}{(t_{x^2})^2}$$

## Extensión a modelos de regresión múltiple:

Para modelos de regresión múltiple, la estimación de la varianza se generaliza a través de una matriz de varianza-covarianza que involucra los coeficientes de regresión.

$$\text{var}(\hat{\beta}) = \hat{\Sigma}(\hat{\beta}) = \begin{bmatrix} \text{var}(\hat{\beta}_0) & \text{cov}(\hat{\beta}_0, \hat{\beta}_1) & \cdots & \text{cov}(\hat{\beta}_0, \hat{\beta}_p) \\ \text{cov}(\hat{\beta}_0, \hat{\beta}_1) & \text{var}(\hat{\beta}_1) & \cdots & \text{cov}(\hat{\beta}_1, \hat{\beta}_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(\hat{\beta}_0, \hat{\beta}_p) & \text{cov}(\hat{\beta}_1, \hat{\beta}_p) & \cdots & \text{var}(\hat{\beta}_p) \end{bmatrix}$$

Este enfoque de estimación garantiza que se tengan en cuenta las particularidades del diseño muestral en la inferencia sobre los coeficientes de regresión.

# Aplicación en encuestas de hogares

El proceso inicia con la lectura de la muestra y definiendo algunas variables de interés.

- ▶ CANTIDAD\_PERSONAS: Cantidad de miembros pertenecientes al hogar.
- ▶ YDISPONIBLE\_PER: Corresponde al ingreso disponible del hogar, dividido por la cantidad de personas en el hogar.
- ▶ GASTO\_CORRIENTE\_HOGAR: Gasto corriente del Hogar
- ▶ CONSUMO\_FINAL\_HOGAR: Gasto de consumo final del Hogar

# Aplicación en encuestas de hogares

CANTIDAD_PERSONAS	GASTO_CORRIENTE_HOGAR	CONSUMO_FINAL_HOGAR	YDISPONIBLE_PER	ingreso_per	ingreso_hog
4	18193	18011	-2811.51	0.00	0.00
1	28131	18184	-1772.00	0.00	0.00
1	21560	15961	-1702.58	0.00	0.00
2	13479	12968	-255.54	0.00	0.00
1	46836	46579	-189.08	0.00	0.00
2	7079	6789	-53.54	0.00	0.00
1	2652	2652	0.00	0.00	0.00
4	5384	5384	18.75	18.75	75.00
3	5494	5494	20.83	20.83	62.50
2	8967	8967	57.71	57.71	115.42
1	12585	11418	68.06	68.06	68.06
4	2828	2828	111.94	111.94	447.75
3	32166	15307	160.00	160.00	479.99
3	5928	5928	166.67	166.67	500.00
3	4625	4625	180.56	180.56	541.67
2	8213	8213	187.50	187.50	375.00
3	4063	4063	196.06	196.06	588.17
6	14287	14287	226.67	226.67	1360.00
4	6406	6406	232.08	232.08	928.33
3	8294	8294	247.00	247.00	741.00

## Definición del objeto survey.design

```
diseno <- encuesta %>% as_survey_design(  
  strata = estrato, # Id de los estratos.  
  ids = F1_A0_UPM, # Id para las observaciones.  
  weights = Factor, # Factores de expansión.  
  nest = TRUE # Valida el anidado dentro del estrato  
)
```

### Sub-grupos

Dividir la muestra en sub-grupos de la encuesta.

```
sub_Urbano <- diseno %>% filter(Area == "1. Urbana") #  
sub_Rural <- diseno %>% filter(Area == "2. Rural") #
```

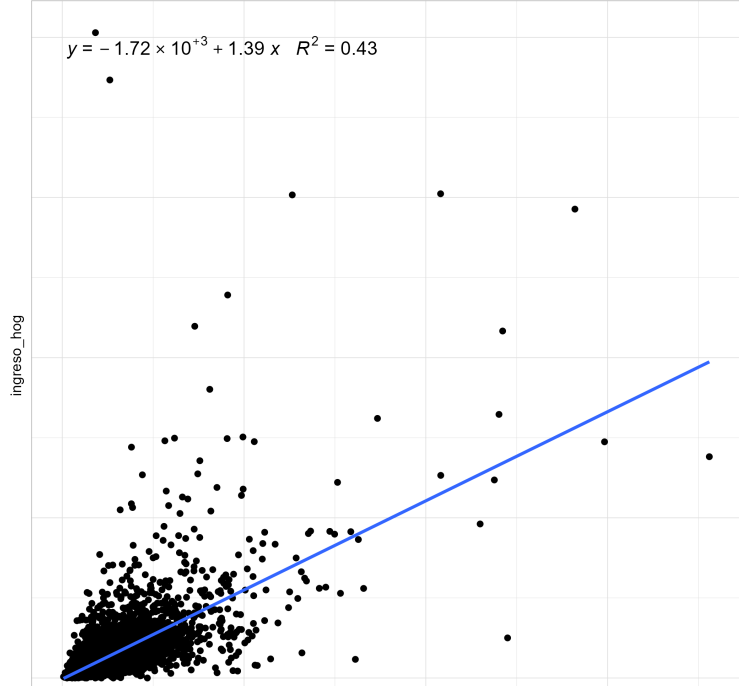
# Scatterplot con los datos encuesta sin ponderar

Una sintaxis similar permite construir el scatterplot en la muestra.

```
plot_sin <-  
  ggplot(data = encuesta,  
    aes(x = GASTO_CORRIENTE_HOGAR, y = ingreso_hog)) +  
  geom_point() +  
  geom_smooth(method = "lm",  
    se = FALSE,  
    formula = y ~ x) +  
  theme_cepal()  
plot_sin <- plot_sin + stat_poly_eq(formula = y~x,  
  aes(label = paste(..eq.label..  
    ..rr.label..., sep = "~~~"), size = 5),  
  parse = TRUE)
```



# Scatterplot con los datos encuesta sin ponderar



## Modelo sin ponderar

El modelo ignorando los factores de expansión quedas así:

```
fit_sinP <- lm(ingreso_hog ~GASTO_CORRIENTE_HOGAR, data = encuesta)
stargazer(fit_sinP, header = FALSE,
          title = "Modelo encuesta Sin ponderar",
          style = "ajps")
```

# Modelo sin ponderar

Tabla 2: Modelo encuesta Sin ponderar

	<b>ingreso_hog</b>
GASTO_CORRIENTE_HOGAR	1.391*** (0.017)
Constant	−1723.000*** (479.200)
N	8746
R-squared	0.434
Adj. R-squared	0.434
Residual Std. Error	28819.000 (df = 8744)
F Statistic	6713.000*** (df = 1; 8744)

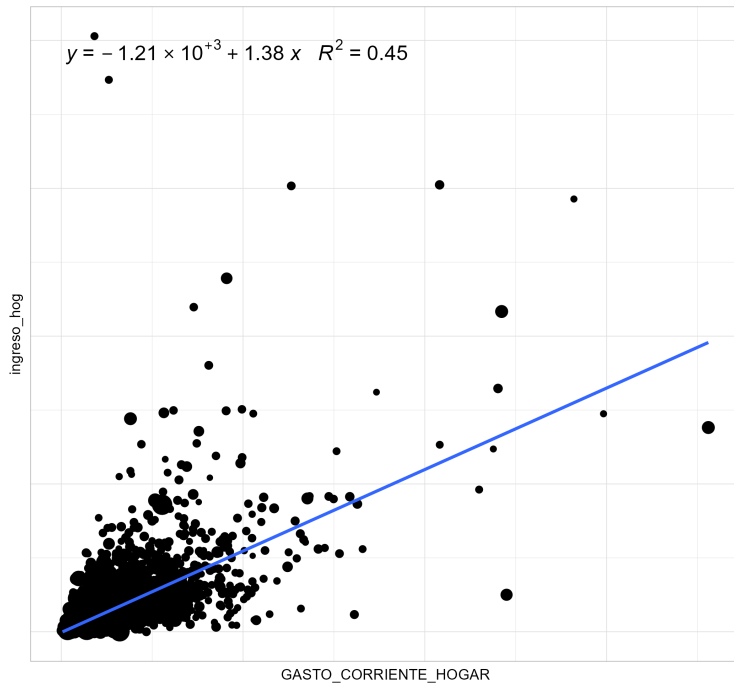
\*\*\*  $p < .01$ ; \*\*  $p < .05$ ; \*  $p < .1$

## Scatterplot con los datos encuesta ponderado

Para que el gráfico tenga en cuenta las ponderaciones debe agregar `mapping = aes(weight = wk)` en la función `geom_smooth`.

```
plot_Ponde <-  
  ggplot(data = encuesta,  
    aes(x = GASTO_CORRIENTE_HOGAR, y = ingreso_hog)) +  
  geom_point(aes(size = Factor)) +  
  geom_smooth(method = "lm",  
    se = FALSE,  
    formula = y ~ x,  
    mapping = aes(weight = Factor)) +  
  theme_cepala()  
plot_Ponde <- plot_Ponde + stat_poly_eq(formula = y~x,  
  aes(weight = Factor,  
    label = paste(..eq.label..,  
      ..rr.label.., sep = "~~~")),  
  parse = TRUE, size = 5)
```

## Scatterplot con los datos encuesta sin ponderar



## Modelo ponderado lm

La función `lm` permite incluir los `weights` en la estimación de los coeficientes.

```
fit_Ponde <- lm(ingreso_hog ~ GASTO_CORRIENTE_HOGAR,  
               data = encuesta, weights = Factor)  
stargazer(fit_Ponde, header = FALSE,  
          title = "Modelo encuesta ponderada",  
          style = "ajps")
```

# Modelo ponderado lm

Tabla 3: Modelo encuesta ponderada

	<b>ingreso_hog</b>
GASTO_CORRIENTE_HOGAR	1.378*** (0.016)
Constant	−1209.000*** (415.200)
N	8746
R-squared	0.452
Adj. R-squared	0.452
Residual Std. Error	448570.000 (df = 8744)
F Statistic	7213.000*** (df = 1; 8744)

\*\*\*  $p < .01$ ; \*\*  $p < .05$ ; \*  $p < .1$

## Modelo ponderado svyglm

Ahora, emplee la función `svyglm` de `survey`

```
fit_svy <- svyglm(ingreso_hog ~ GASTO_CORRIENTE_HOGAR,  
                  design = diseno)
```



# Resumen del Modelo

Tabla 4: Modelo encuesta ponderada, svyglm

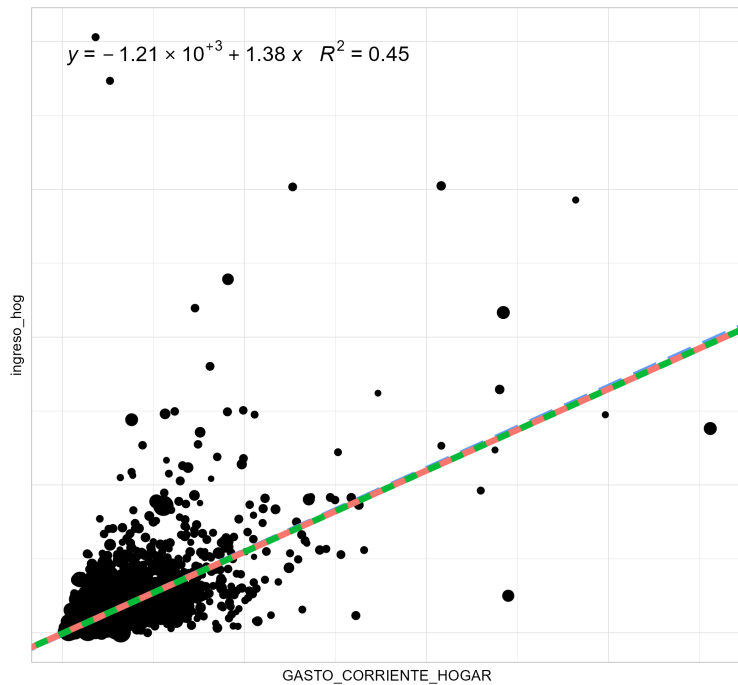
	<b>ingreso_hog</b>
GASTO_CORRIENTE_HOGAR	1.378*** (0.081)
Constant	−1209.000 (1361.000)
N	8746
AIC	205813.000

\*\*\*  $p < .01$ ; \*\*  $p < .05$ ; \*  $p < .1$

## Comparando los resultados

```
df_model <- data.frame(  
  intercept = c(coefficients(fit_sinP)[1],  
                 coefficients(fit_Ponde)[1],  
                 coefficients(fit_svy)[1]),  
  slope = c( coefficients(fit_sinP)[2],  
            coefficients(fit_Ponde)[2],  
            coefficients(fit_svy)[2]),  
  Modelo = c("Sin ponderar",  
            "Ponderado(lm)", "Ponderado(svyglm)"))  
  
plot_Ponde2 <- plot_Ponde + geom_abline( data = df_model,  
  mapping = aes( slope = slope,  
                intercept = intercept, linetype = Modelo,  
                color = Modelo ), size = 2  
)
```

## Comparando los resultados



## Comparando los resultados

Variable	Sin Pond	Ponde(lm)	Ponde(svyglm)
(Intercept)	-1722.608	-1209.015	-1209.015
p-value	(<0.001)	(0.004)	(0.375)
GASTO_CORRIENTE_HOGAR	1.391	1.378	1.378
p-value	(<0.001)	(<0.001)	(<0.001)
Num.Obs.	8746	8746	8746
R2	0.434	0.452	0.452
R2 Adj.	0.434	0.452	-5.619
AIC	204445.8	205813.2	202638.9
F	6712.633	7213.379	288.839
RMSE	28815.76	28817.66	28817.66

## Diagnostico del modelo

# Diagnostico del modelo

**Adecuado Ajuste del Modelo:** - Verificar que el modelo se ajuste adecuadamente a los datos recopilados en la encuesta. - Evaluar si la relación funcional especificada es apropiada para representar las variables de interés.

**Normalidad de Errores:** - Examinar si los errores del modelo siguen una distribución normal. - Esto es crucial para realizar pruebas de hipótesis precisas y estimar intervalos de confianza confiables.

**Varianza Constante de Errores:** - Asegurarse de que la varianza de los errores sea constante en todos los niveles de las variables independientes. - La heterocedasticidad puede impactar en las pruebas y la interpretación de coeficientes.

# Diagnostico del modelo

**Errores No Correlacionados:** - Evaluar si los errores pueden considerarse no correlacionados entre sí. - La autocorrelación de errores puede afectar la eficiencia de las estimaciones.

**Datos Influyentes:** - Identificar valores atípicos o datos influyentes que tienen un efecto desproporcionadamente grande en el modelo de regresión. - Estos datos deben tratarse con precaución y su impacto debe ser evaluado.

**Valores Atípicos (Outliers):**

## Estimación del $R^2$ y $R_{adj}^2$

- ▶ En análisis de regresión, el coeficiente de determinación ( $R^2$ ) mide la variabilidad explicada por el modelo.
- ▶ El  $R_{\omega}^2$  ajusta  $R^2$  para muestras complejas, considerando ponderaciones de la muestra.
- ▶  $R_{\omega}^2$  se basa en la suma de cuadrados totales ponderada (WSST) y la suma de cuadrados del error ponderada (WSSE).
- ▶ La fórmula de  $R^2$  es  $1 - \frac{SSE}{SST}$ , donde SSE es la suma de cuadrados del error y SST es la suma de cuadrados totales.



## Estimación del $R^2$ y $R_{adj}^2$

- Para  $R_{\omega}^2$ , la fórmula es  $1 - \frac{WSSE}{W SST}$ , considerando las ponderaciones de la muestra.

$$\widehat{WSSE}_{\omega} = \sum_h^H \sum_{\alpha}^{a_h} \sum_{i=1}^{n_{h\alpha}} \omega_{h\alpha i} \left( y_{h\alpha i} - x_{h\alpha i} \hat{\beta} \right)^2$$

- Se utiliza el coeficiente de determinación ajustado ( $R_{adj}^2$ ) para tener en cuenta el tamaño de la muestra y el número de predictores en el modelo.
- $R_{adj}^2$  se calcula como  $1 - \frac{(n-1)}{(n-p)} R_{\omega}^2$ , donde  $n$  es el tamaño de la muestra y  $p$  es el número de predictores.

Estimación del  $R^2$  para el modelo del ingreso.

```
fit_svy <- svyglm(ingreso_hog ~ GASTO_CORRIENTE_HOGAR ,
                  design = diseno,family=stats::gaussian())

medY <- diseno %>% summarise(medY = survey_mean(ingreso_hog))

diseno %<>% mutate(
  ypred = fitted(fit_svy, type = "response"),
  medY = medY,
  sst = (ingreso_hog - medY$medY)^2,
  sse = (ypred - medY$medY)^2
)

diseno %>% summarise(WSST = survey_total(sst),
                    WSSE = survey_total(sse)) %>%
  transmute(WSST, WSSE, R2 = WSSE/WSST)
```

## Estimación del $R^2$ para el modelo del ingreso

El resultado para el  $R^2$  es

WSST	WSSE	R2
3.211e+15	1.451e+15	0.452

De forma alternativa es:

```
modNul <- svyglm(ingreso_hog ~ 1, design = diseno)
s1 <- summary(fit_svy)
s0 <-summary(modNul)

WSST<- s0$dispersion
WSSE<- s1$dispersion
R2 = 1- WSSE/WSST
R2
```

```
      variance      SE
[1,]      0.452 1.07e+08
```

## Estimación del $R_{adj}^2$ para el modelo del ingreso

Calculamos el  $R_{adj}^2$  utilizando la fórmula adecuada. Asegúrate de definir los valores de  $n$  y  $p$  de acuerdo a tu modelo.

```
n = nrow(encuesta)
p = 2
(R2Adj = 1 - ((n-1)/(n-p)) * R2)
```

	variance	SE
[1,]	0.548	1.07e+08

# Metodología de los Q\_Weighting de pfefferman

Cuando trabajamos con datos de encuestas que siguen un diseño muestral complejo y es posible aplicar la metodología de los q-weights (**Pffeferman, 2011**).,

1. **Ajuste del Modelo de Regresión a los Q-Weights:** Inicialmente, ajustamos un modelo de regresión lineal a los q-weights en R. Esto se hace utilizando la función `lm()`.

```
fit_wgt <- lm(1/Factor ~ GASTO_CORRIENTE_HOGAR, data = encuesta)
```

2. **Obtención de Predicciones de Q-Weights:** A continuación, calculamos las predicciones de los q-weights para cada caso, utilizando las variables predictoras del modelo de regresión.

```
qw <- predict(fit_wgt)  
summary(qw)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.0058	0.0061	0.0062	0.0064	0.0065	0.0151

# Metodología de los Q\_Weighting de pfefferman

3. **Creación de Nuevos Q-Weights:** Para obtener los q-weights ajustados, dividimos los weights originales por las predicciones calculadas en el paso anterior.

```
encuesta <- encuesta %>% mutate(wk1 = Factor/qw)
```

4. **Definición de un Diseño Muestral con Q-Weights:** Usamos los nuevos q-weights obtenidos para definir un diseño muestral que refleje estos pesos.

```
diseno_qwgt <- encuesta %>%  
  as_survey_design(  
    strata = estrato, # Id de los estratos.  
    ids = F1_A0_UPM, # Id para las observaciones.  
    weights = wk1, # Factores de expansión.  
    nest = TRUE # Valida el anidado dentro del estrato  
  )
```

# Modelos empleando los Q\_Weighting

Estimando los coeficientes del modelo con los Q\_Weighting de pfefferman

```
library(tidyr)
fit_svy_qwgt <- svyglm(ingreso_hog ~ GASTO_CORRIENTE_HOGAR,
                      design = diseno_qwgt)
s1_qwgt <- summary(fit_svy_qwgt)
tidy(fit_svy_qwgt)
```

term	estimate	std.error	statistic	p.value
(Intercept)	-1306.548	1109.1754	-1.178	0.2392
GASTO_CORRIENTE_HOGAR	1.383	0.0696	19.872	0.0000

## Calculo del $R^2$ y $R^2_{adj}$

Obtenido el  $R^2$

```
WSST<- s0$dispersion  
WSSE<- s1_qwgt$dispersion  
(R2 = 1- WSSE/WSST)
```

```
      variance      SE  
[1,]    0.531 92109176
```

Obtenido el  $R^2_{adj}$

```
n = nrow(encuesta)  
p = 2  
(R2Adj = 1-(((1-R2)*(n-1))/(n-1-1)))
```

```
      variance      SE  
[1,]    0.531 92109176
```



## Modelos empleando los Q\_Weighting

Tabla 9: Comprando Modelos con Q Weighting

Variable	svyglm(wgt)	svyglm(qwgt)
(Intercept)	-1209.015	-1306.548
p-value	(0.375)	(0.239)
GASTO_CORRIENTE_HOGAR	1.378	1.383
p-value	(<0.001)	(<0.001)
Num.Obs.	8746	8746
R2	0.452	0.423
AIC	202638.9	201262.1
F	288.839	394.884
RMSE	28817.66	28817.13

¡Gracias!

*Email:* [andres.gutierrez@cepal.org](mailto:andres.gutierrez@cepal.org)