

Profesor: Héctor Bahamonde, PhD.

e: hector.bahamonde@uoh.cl

w: www.hectorbahamonde.com

Curso: MLE.

TA: Gonzalo Barria.

I. INFERENCIA E INTERPRETACIÓN

Los GLMs no se pueden interpretar usando la tabla. Los coeficientes no significan lo que significaban en nuestro mundo OLS. Para interpretarlos, debemos usar otras herramientas.

Intervalos de confianza Ya sabemos lo que un intervalo de confianza significa: si nuestra confianza es al 95%, sabemos que el “true value” de la estimación caerá el 95% dentro del margen de los intervalos de confianza. Es en base a esto que antes hablabamos de “significancia estadística”.

Carguemos los datos

```
options(scipen=10000000)
set.seed(2020)
# Data
mydata <- read.csv("https://stats.idre.ucla.edu/stat/data/binary.csv")
head(mydata)

##   admit gre  gpa rank
## 1     0 380 3.61   3
## 2     1 660 3.67   3
## 3     1 800 4.00   1
## 4     1 640 3.19   4
## 5     0 520 2.93   4
## 6     1 760 3.00   2

summary(mydata)

##      admit      gre      gpa      rank
```

```
## Min.      :0.0000   Min.      :220.0   Min.      :2.260   Min.      :1.000
## 1st Qu.:0.0000   1st Qu.:520.0   1st Qu.:3.130   1st Qu.:2.000
## Median :0.0000   Median :580.0   Median :3.395   Median :2.000
## Mean    :0.3175   Mean    :587.7   Mean    :3.390   Mean    :2.485
## 3rd Qu.:1.0000   3rd Qu.:660.0   3rd Qu.:3.670   3rd Qu.:3.000
## Max.    :1.0000   Max.    :800.0   Max.    :4.000   Max.    :4.000
```

```
mydata$rank <- factor(mydata$rank)
```

Estimemos el primer modelo

```
logit.1 <- glm(admit ~ gre + gpa + rank, data = mydata, family = binomial(link = "logit"))
summary(logit.1)
```

```
##
## Call:
## glm(formula = admit ~ gre + gpa + rank, family = binomial(link = "logit"),
##      data = mydata)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.6268  -0.8662  -0.6388   1.1490   2.0790
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -3.989979    1.139951  -3.500  0.000465 ***
## gre          0.002264    0.001094   2.070  0.038465 *
## gpa          0.804038    0.331819   2.423  0.015388 *
## rank2       -0.675443    0.316490  -2.134  0.032829 *
## rank3       -1.340204    0.345306  -3.881  0.000104 ***
## rank4       -1.551464    0.417832  -3.713  0.000205 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 499.98  on 399  degrees of freedom
## Residual deviance: 458.52  on 394  degrees of freedom
## AIC: 470.52
##
## Number of Fisher Scoring iterations: 4
```

Y ahora, estimemos los intervalos de confianza. Los intervalos de confianza son sencillos de calcular. Si queremos un 95% de confianza, miramos la tabla con los valores Z. Para ese porcentaje es 1.960 (para un 90% es 1.645).

$$\hat{\theta} \pm 1.960 \times SE_{\hat{\theta}} \quad (1)$$

donde $\hat{\theta}$ es la estimación (el parámetro), 1.960 es el valor z , DS es la desviación estándar, y N es el largo de la base de datos. Si te fijas, tenemos el signo \pm que implica que debemos restar para obtener el rango mínimo del intervalo de confianza, y sumar para obtener el rango máximo del intervalo de confianza.

Por ejemplo, el coeficiente de “gre” es $\hat{\theta} = 0.0022644$, el $SE = 0.001094$. Entonces, $0.002264426 + 1.96 * 0.001094$ para el intervalo superior, y $0.002264426 - 1.96 * 0.001094$ para el intervalo inferior.

Usemos un paquete.

```
confint(logit.1)

##              2.5 %       97.5 %
## (Intercept) -6.2716202334 -1.792547080
## gre          0.0001375921  0.004435874
## gpa          0.1602959439  1.464142727
## rank2        -1.3008888002 -0.056745722
## rank3        -2.0276713127 -0.670372346
## rank4        -2.4000265384 -0.753542605
```

Odds Ratios Esta manera de interpretar solo funciona con los *link function* tipo logit (logit, multinomial logit, etc.), no probit links (i.e. probit, multinomial probit). Formalmente,

$$\ln \Omega(x) = x\beta$$

$$\Omega(x) = \frac{\Pr(y = 1|x)}{\Pr(y = 0|x)} \quad (2)$$

donde Equation 2 es el *odd ratio* de que y sea 1 dado x , relativo a que y sea 0 dado x . Es un ratio, una fracción. Su interpretación es intuitiva: “cuando x cambia, cuánto cambia el logit estimado ($\hat{\theta}$) manteniendo los otros parametros constantes”?

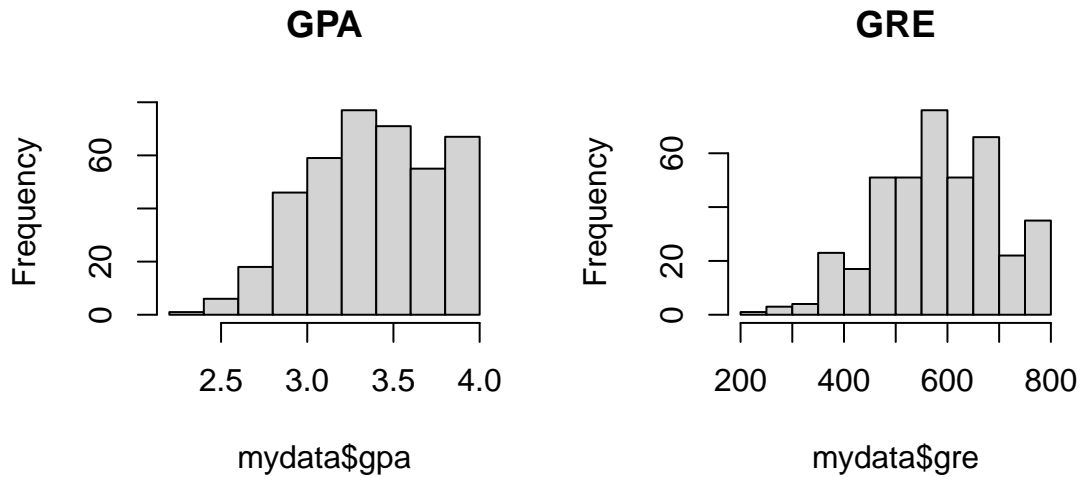
Lo bueno de esta interpretacion es que podemos observar el cambio de la variable dependiente. Ten en cuenta que los *odds ratios* son comparables sólo con la misma variable. **Lo malo** es que seguimos en unidades de la escala logit (que sigue siendo poco interpretable).

```
if (!require("pacman")) install.packages("pacman"); library(pacman)
p_load(oddsratio)
#
or_glm(data = mydata,
        model = logit.1,
        incr = list(gre = 380, gpa = 5))
```

##	predictor	oddsratio	ci_low (2.5)	ci_high (97.5)	increment
## 1	gre	2.364	1.054	5.396	380
## 2	gpa	55.712	2.229	1511.282	5
## 3	rank2	0.509	0.272	0.945	Indicator variable
## 4	rank3	0.262	0.132	0.512	Indicator variable
## 5	rank4	0.212	0.091	0.471	Indicator variable

Aquí vemos que la variable dependiente (**admit**) es 55 veces más probable de ocurrir cuando la variable independiente **gpa** incrementa de su media (3.3899) a 5. Como ves, ganamos interpretatibilidad del modelo.

```
par(mfrow=c(1,2))
hist(mydata$gpa, main = "GPA")
hist(mydata$gre, main = "GRE")
```



Como interpretamos GRE?

```
mean(mydata$gpa)

## [1] 3.3899

mean(mydata$gre)

## [1] 587.7
```

Si bien es cierto que ganamos interpretación del modelo, lo malo es seguimos hablando con poca precisión. Decir que algo es 55 veces más probable es interesante, pero no sé si tan “científico”.

Partial/Marginal Changes in y Los “cambios parciales” son muy parecidos a los *odds ratios*. Son parecidos porque nos muestra “cuando cambia \hat{y} cuando cambia x ”. Formalmente,

$$\frac{\partial \hat{y}}{\partial x_k} = \beta_k \quad (3)$$

lo que significa “por un cambio en x_k , se espera que \hat{y} cambie β_k ” (manteniendo todas las otras variables constantes en su media). Sin embargo, debido a que la varianza de \hat{y} es desconocida (la varianza es un *population parameter*), entonces, esto complica la interpretación de β_k . Para resolver esto, lo que hacemos es pensar en coeficientes β_k estandarizados. Siguiendo [Equation I](#), lo que pensamos es en “por un cambio en x_k , se espera que \hat{y} cambie β_k **desviaciones estándar**”. Formalmente,

$$\beta_k^S = \frac{\sigma_k \beta_k}{\sigma_{\hat{y}}} = \sigma_k \beta_k^{S_{\hat{y}}} \quad (4)$$

Nota que también estandariza \hat{y} . Debido a que $\hat{\sigma}_{\hat{y}} = \sigma_{\hat{y}}$ (i.e. podemos estimar la varianza estandarizada de los datos),

$$\hat{\sigma}_{\hat{y}} = \beta^\top \hat{\sigma}(x) \beta + \sigma(\epsilon) \quad (5)$$

Calculemos dos escenarios. Uno donde el estudiante le iba muy mal en el colegio, pero tuvo un buen puntaje en la prueba de admisión de doctorado (GRE).

```
if (!require("pacman")) install.packages("pacman"); library(pacman)
p_load(margins)
summary(margins(logit.1,
  at = list(
    gre = min(mydata$gre),
    gpa = max(mydata$gpa))
))
```

##	factor	gre	gpa	AME	SE	z	p	lower	upper
##	gpa	220.0000	4.0000	0.1408	0.0839	1.6786	0.0932	-0.0236	0.3052
##	gre	220.0000	4.0000	0.0004	0.0001	3.5999	0.0003	0.0002	0.0006
##	rank2	220.0000	4.0000	-0.1529	0.0768	-1.9898	0.0466	-0.3034	-0.0023
##	rank3	220.0000	4.0000	-0.2657	0.0927	-2.8657	0.0042	-0.4475	-0.0840
##	rank4	220.0000	4.0000	-0.2929	0.1002	-2.9234	0.0035	-0.4893	-0.0965

Ahora calculemos el opuesto:

```
summary(margins(logit.1,
  at = list(
    gre = max(mydata$gre),
    gpa = min(mydata$gpa):mean(mydata$gpa))
  ))
```

##	factor	gre	gpa	AME	SE	z	p	lower	upper
##	gpa	800.0000	2.2600	0.1352	0.0293	4.6134	0.0000	0.0777	0.1926
##	gpa	800.0000	3.2600	0.1792	0.0695	2.5770	0.0100	0.0429	0.3155
##	gre	800.0000	2.2600	0.0004	0.0003	1.4657	0.1427	-0.0001	0.0009
##	gre	800.0000	3.2600	0.0005	0.0003	1.8747	0.0608	-0.0000	0.0010
##	rank2	800.0000	2.2600	-0.1489	0.0771	-1.9320	0.0534	-0.2999	0.0022
##	rank2	800.0000	3.2600	-0.1668	0.0762	-2.1897	0.0285	-0.3161	-0.0175
##	rank3	800.0000	2.2600	-0.2564	0.0876	-2.9273	0.0034	-0.4280	-0.0847
##	rank3	800.0000	3.2600	-0.3193	0.0777	-4.1088	0.0000	-0.4717	-0.1670
##	rank4	800.0000	2.2600	-0.2820	0.0988	-2.8541	0.0043	-0.4756	-0.0883
##	rank4	800.0000	3.2600	-0.3608	0.0879	-4.1047	0.0000	-0.5331	-0.1885

Para entender mejor esto, grafiquemos:

```
par(mfrow=c(1,2))
cplot(logit.1, "gre")
```

##	xvals	yvals	upper	lower
## 1	220.0000	0.1912913	0.3293332	0.05324935
## 2	244.1667	0.1999003	0.3349924	0.06480810
## 3	268.3333	0.2087967	0.3405397	0.07705355
## 4	292.5000	0.2179811	0.3459866	0.08997559
## 5	316.6667	0.2274534	0.3513519	0.10355486
## 6	340.8333	0.2372125	0.3566644	0.11776057
## 7	365.0000	0.2472562	0.3619647	0.13254773
## 8	389.1667	0.2575817	0.3673097	0.14785363
## 9	413.3333	0.2681847	0.3727760	0.16359328

```
## 10 437.5000 0.2790601 0.3784663 0.17965382
## 11 461.6667 0.2902016 0.3845153 0.19588791
## 12 485.8333 0.3016019 0.3910967 0.21210700
## 13 510.0000 0.3132523 0.3984280 0.22807664
## 14 534.1667 0.3251433 0.4067686 0.24351805
## 15 558.3333 0.3372640 0.4164053 0.25812272
## 16 582.5000 0.3496025 0.4276189 0.27158598
## 17 606.6667 0.3621457 0.4406338 0.28365750
## 18 630.8333 0.3748795 0.4555665 0.29419252
## 19 655.0000 0.3877889 0.4723996 0.30317809
## 20 679.1667 0.4008577 0.4909949 0.31072049
```

```
cplot(logit.1, "gpa")
```

```
##      xvals      yvals      upper      lower
## 1  2.2600 0.1798308 0.2984691 0.06119251
## 2  2.3325 0.1885895 0.3051858 0.07199319
## 3  2.4050 0.1976720 0.3118808 0.08346312
## 4  2.4775 0.2070802 0.3185673 0.09559314
## 5  2.5500 0.2168152 0.3252662 0.10836423
## 6  2.6225 0.2268769 0.3320085 0.12174533
## 7  2.6950 0.2372641 0.3388377 0.13569055
## 8  2.7675 0.2479744 0.3458129 0.15013578
## 9  2.8400 0.2590039 0.3530130 0.16499472
## 10 2.9125 0.2703476 0.3605409 0.18015436
## 11 2.9850 0.2819991 0.3685276 0.19547065
## 12 3.0575 0.2939505 0.3771355 0.21076543
## 13 3.1300 0.3061923 0.3865574 0.22582718
## 14 3.2025 0.3187139 0.3970088 0.24041895
## 15 3.2750 0.3315028 0.4087086 0.25429694
## 16 3.3475 0.3445454 0.4218506 0.26724013
## 17 3.4200 0.3578264 0.4365681 0.27908480
```



```
## 18 3.4925 0.3713294 0.4529077 0.28975110
## 19 3.5650 0.3850364 0.4708228 0.29925010
## 20 3.6375 0.3989283 0.4901876 0.30766903
```

