

EST-25134: Aprendizaje Estadístico

Profesor: Alfredo Garbuno Iñigo — Primavera, 2023 — Regresión lineal.

Objetivo. Repasaremos los conceptos de regresión lineal desde un punto de vista de inferencia. Haremos conexiones interesantes con conceptos clave en probabilidad y teoría de la información. Veremos métricas de desempeño para modelos de regresión. Nos llevará a cuestionar el modelo bajo el enfoque de predicción.

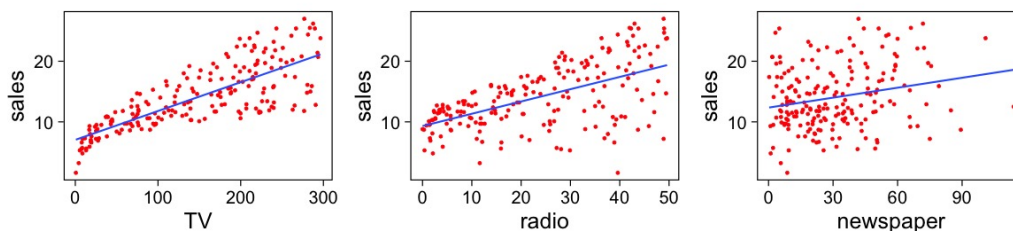
Lectura recomendada: Capítulo 3 de James et al. [1].

1. INTRODUCCIÓN

Aproximamos $f(X)$ como una combinación lineal de las entradas.

1.1. Datos de marketing

¿Qué preguntas serían las que nos interesarían?



¿Hay alguna relación entre datos de entrada? En especial, ¿hay alguna relación entre cuánto se le dedica al presupuesto de *marketing* y las ventas observadas? ¿Qué tipo de *marketing* contribuye mas a las ventas? ¿Qué tan precisos podemos ser con nuestras predicciones? La relación es lineal? ¿Hay algún tipo de sinergia?

2. EL MODELO SIMPLE

$$Y = \beta_0 + \beta_1 X + \varepsilon. \quad (1)$$

Si tuvieramos un estimador $\hat{\beta}$, ¿cómo podemos realizar predicciones?

2.1. Estimación de parámetros

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i. \quad (2)$$

Formulación de problema de optimización. Conexiones interesantes.

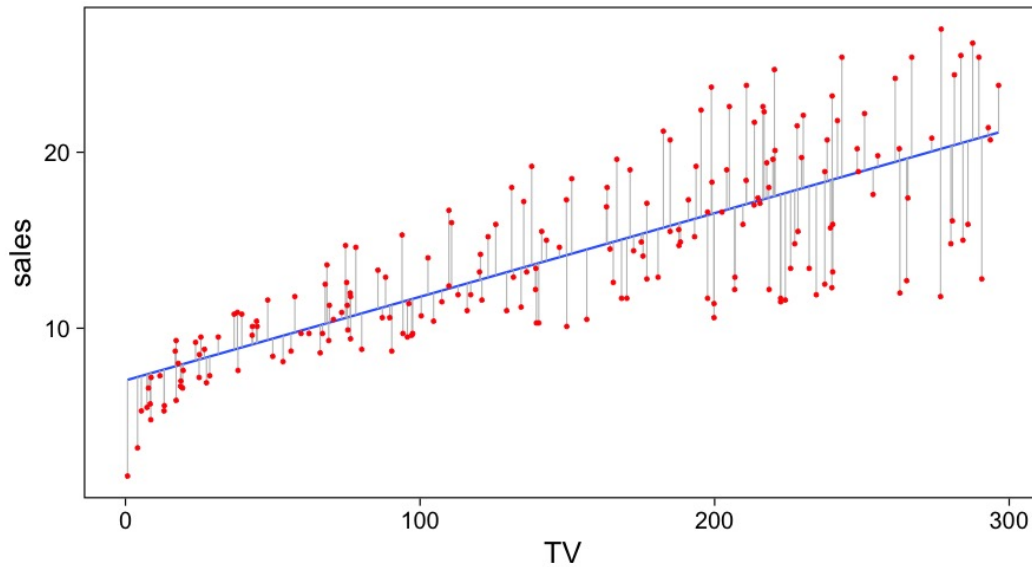


FIGURA 1. Ajuste y residuales a la recta de mínimos cuadrados.

2.2. Solución

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad (3)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}. \quad (4)$$

¿Qué deficiencia encuentras en el ajuste?

2.3. ¿Precisión en los estimadores?

$$SE(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad (5)$$

$$SE(\hat{\beta}_0)^2 = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]. \quad (6)$$

¿A qué se debe esta variabilidad? Se pueden construir **intervalos de confianza**.

```
1 model >
2   summary()
```

LISTING 1. Resumen del modelo.

```
1
2 Call:
3 lm(formula = sales ~ TV, data = data)
```

```

4
5 Residuals:
6   Min      1Q  Median      3Q      Max
7  -8.386  -1.955  -0.191   2.067   7.212
8
9 Coefficients:
10              Estimate Std. Error t value Pr(>|t|)
11 (Intercept)   7.03259    0.45784   15.4   <2e-16 ***
12 TV            0.04754    0.00269   17.7   <2e-16 ***
13 ---
14 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15
16 Residual standard error: 3.26 on 198 degrees of freedom
17 Multiple R-squared:  0.612,    Adjusted R-squared:  0.61
18 F-statistic: 312 on 1 and 198 DF,  p-value: <2e-16

```

```

1 model >
2   broom::tidy()

```

LISTING 2. Resumen del modelo (tidy).

```

1 # A tibble: 2 × 5
2   term      estimate std.error statistic  p.value
3   <chr>      <dbl>    <dbl>    <dbl>    <dbl>
4 1 (Intercept)   7.03      0.458     15.4 1.41e-35
5 2 TV            0.0475    0.00269    17.7 1.47e-42

```

```

1 genera_datos <- function(id){
2   a <- 1; b <- 0; n <- 100
3   tibble(x = runif(n, -1, 1),
4         y = a * x + b + rnorm(n, sd = 1))
5 }
6 ajusta_modelo <- function(datos){
7   modelo <- lm(y ~ x, datos)
8   modelo
9 }

```

```

1 simulacion <- tibble(id = seq(1, 10)) >
2   mutate(datos = map(id, genera_datos),
3         modelo = map(datos, ajusta_modelo),
4         ajuste = map(modelo, broom::tidy))

```

2.4. Prueba de hipótesis

H_0 : No hay relación entre X y Y , (7)

H_1 : Existe una hay relación entre X y Y . (8)

La prueba de hipótesis se efectúa en el contexto del modelo que estamos proponiendo.

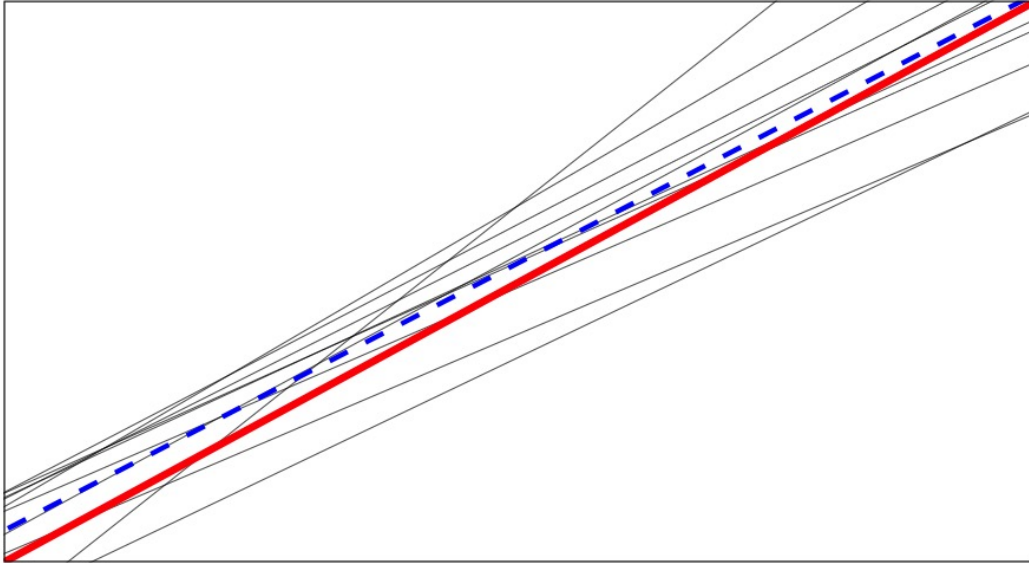


FIGURA 2. Simulación de ajuste (variación en datos).

2.5. El valor-p

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}, \quad \text{distribución } t_{n-2}. \quad (9)$$

2.6. Midiendo la precisión del modelo

$$RSE = \sqrt{\frac{1}{n-2} RSS}. \quad (10)$$

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

Es una métrica usual de ajuste. Nos dice qué tan precisos podemos ser al calcular una nueva predicción. Tiene sentido cuando comparamos con las unidades de la variable respuesta. Es decir, no es una métrica absoluta.

$$R^2 = \frac{TSS - RSS}{TSS}. \quad (11)$$

$$\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2.$$

Hay que tener cuidado pues la R^2 es una métrica de correlación lineal, no de ajuste. Esto lo podemos ver en el caso sencillo de una variable respuesta y un modelo lineal. También está sujeta a la variación de la respuesta. Depende la aplicación para determinar cuándo tenemos un buen coeficiente de variación explicada.

3. EL MODELO MULTIVARIADO

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \varepsilon. \quad (12)$$

Cuando tenemos múltiples predictores nos gustaría poder entender la relación de cada uno con la respuesta. ¿Ajustaríamos un modelo independiente con sólo un predictor?

3.1. Interpretación

$$\text{sales} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times \text{newspaper} + \varepsilon. \quad (13)$$

El modelo de regresión lineal usualmente se interpreta como los efectos (¿promedio, esperados?) de cada variable al mantener todas las demás *constantes*. Hay problemas cuando hay correlación entre predictores. Cuidado con datos observacionales.

3.2. Estimación

```
1 model <- lm(sales ~ ., data)
```

```
1 model >
2 broom::tidy()
```

```
1 # A tibble: 4 × 5
2   term      estimate std.error statistic  p.value
3   <chr>      <dbl>    <dbl>    <dbl>    <dbl>
4 1 (Intercept)  2.94      0.312      9.42 1.27e-17
5 2 TV          0.0458   0.00139    32.8 1.51e-81
6 3 radio       0.189    0.00861    21.9 1.51e-54
7 4 newspaper  -0.00104   0.00587    -0.177 8.60e- 1
```

Desarrollo de verosimilitud.

3.3. ¿Existe una relación entre la respuesta y los predictores?

Nos preguntamos si es que existe alguna $\beta_j \neq 0$.

$$F = \frac{(\text{TSS} - \text{RSS})/p}{\text{RSS}/(n - p - 1)} \sim F_{p, n-p-1}. \quad (14)$$

La prueba de hipótesis que formularíamos sería probar contra alguna $\beta_j \neq 0$. Se puede probar que si el supuesto del modelo lineal es correcto y bajo la hipótesis nula el cociente será cercano a 1. En caso de que la hipótesis **alternativa** sea cierta entonces $F > 1$.

```
1 model >
2   broom::glance() >
3   select(statistic, p.value, df, df.residual)
```

LISTING 3. Resumen global del modelo (tidy).

```
1 # A tibble: 1 × 4
2   statistic p.value    df df.residual
3   <dbl>     <dbl> <dbl>    <int>
4 1      570. 1.58e-96     3        196
```

- ¿Por qué tenemos que evaluar en conjunto?

¿Qué pasa en el caso con 100 predictores donde no hay relación?

3.4. ¿Cuáles son los predictores importantes?

Métodos de selección.

La idea mas ingenua es ajustar todas las posibles combinaciones. Pero se pueden construir modelos de manera secuencial. Usualmente ajustando y comparando con respecto a *alguna métrica*. Mas adelante lo estudiaremos.

3.5. ¿Qué tan bien ajusta el modelo?

Podemos usar las métricas típicas como el RSE o la R^2 .

R^2 : Agregar predictores siempre ayuda (en datos de entrenamiento).
RSE: Podemos tener problemas pues mientras mas variables agregemos si el cambio en residuales es pequeño en relación al aumento de p .

3.6. ¿Cómo predecimos y que tan precisa es nuestra predicción?

Podemos utilizar **intervalos confianza**. Mejor aún, podemos utilizar **intervalos de predicción**.

4. EXTENSIONES

4.1. Predictores cualitativos

Modelo con respuestas binarias (1D). ¿Qué tal que tenemos mas categorias?

4.2. Interacciones

Eliminar el supuesto aditivo: *interacciones* y *no-linealidad*.

```
1 model.1 <- lm(sales ~ TV + radio, data)
2 model.2 <- lm(sales ~ TV + radio + TV:radio, data)
```

LISTING 4. Ajuste de modelos sin/con interacciones.

```
1 tibble(modelo = list(model.1, model.2),
2       tipo    = c("lineal", "interaccion")) >
3   mutate(resultados = map(modelo, broom::tidy)) >
4   select(-modelo) >
5   unnest(resultados) >
6   select(tipo, term, estimate, p.value)
```

```
1 # A tibble: 7 × 4
2   tipo      term      estimate  p.value
3   <chr>    <chr>      <dbl>    <dbl>
4 1 lineal  (Intercept)  2.92     4.57e-19
5 2 lineal  TV           0.0458   5.44e-82
6 3 lineal  radio        0.188    9.78e-59
7 4 interaccion (Intercept)  6.75     1.54e-68
8 5 interaccion TV           0.0191   2.36e-27
9 6 interaccion radio      0.0289   1.40e- 3
10 7 interaccion TV:radio    0.00109  2.76e-51
```

LISTING 5. Resúmenes sobre los coeficientes.

```
1 tibble(modelo = list(model.1, model.2)) >
2   mutate(resultados = map(modelo, broom::glance)) >
3   select(-modelo) >
4   unnest(resultados) >
5   select(r.squared, sigma, AIC, deviance)
```

```
1 # A tibble: 2 × 4
2   r.squared sigma    AIC deviance
3   <dbl> <dbl> <dbl>    <dbl>
4 1   0.897 1.68   780.    557.
5 2   0.968 0.944  550.    174.
```

LISTING 6. Resúmenes globales de los modelos.

El efecto de incrementar el presupuesto en un canal de ventas puede aumentar la efectividad de otro.

4.3. Jerarquías

¿Qué pasa cuando un valor- p de una interacción es pequeño, pero de los términos individuales no?

4.4. Interacciones y modelos múltiples

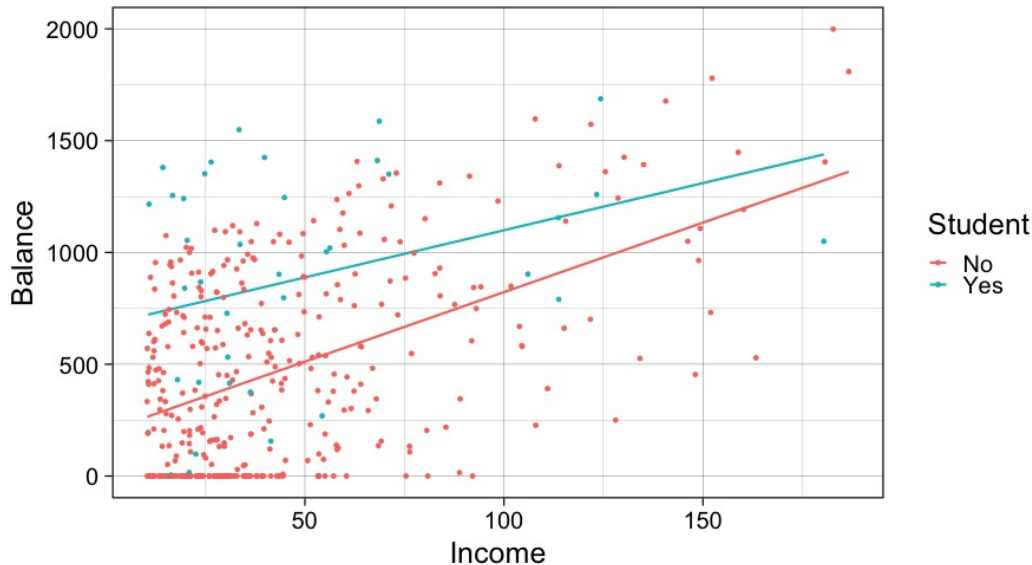


FIGURA 3. Ajuste con interacción cualitativa y cuantitativa.

4.5. Problemas con supuestos.

- No hay una relación lineal.
- Los errores están correlacionados.
- No hay varianza constante.
- Valores atípicos.
- Multicolinealidad.
- Puntos ancla.

5. GENERALIZACIONES

- Problemas de clasificación (siguiente).
- No-linealidad.
- Interacciones.
- Regularización.

6. APLICACIÓN: PREDICCIÓN DE ASISTENCIA EN PARTIDOS DE FOOTBALL

Tomado de esta [liga](#) y es un buen ejemplo de lo que se puede empezar a hacer con [tidymodels](#) [2].

```
1 library(tidyverse)
2 base_url <- "https://raw.githubusercontent.com/rfordatascience/tidytuesday/
  master/data/2020/2020-02-04/"
3
4 attendance <- read_csv(paste(base_url, "attendance.csv", sep = ""),
```



```

5         progress = FALSE, show_col_types = FALSE)
6 standings <- read_csv(paste(base_url, "standings.csv", sep = ""),
7                       progress = FALSE, show_col_types = FALSE)
8
9 attendance_joined <- attendance >
10    left_join(standings, by = c("year", "team_name", "team"))
11 attendance_joined

```

```

1 # A tibble: 10,846 × 20
2   team      team...1n year total   home   away week ...2weekl wins loss
3   <chr>    <chr>    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <
4   1 Arizona ...Cardina 2000 893926 387475 506451 1 77434 3 13
5   2 Arizona ...Cardina 2000 893926 387475 506451 2 66009 3 13
6   3 Arizona ...Cardina 2000 893926 387475 506451 3 NA 3 13
7   4 Arizona ...Cardina 2000 893926 387475 506451 4 71801 3 13
8   5 Arizona ...Cardina 2000 893926 387475 506451 5 66985 3 13
9   6 Arizona ...Cardina 2000 893926 387475 506451 6 44296 3 13
10  7 Arizona ...Cardina 2000 893926 387475 506451 7 38293 3 13
11  8 Arizona ...Cardina 2000 893926 387475 506451 8 62981 3 13
12  9 Arizona ...Cardina 2000 893926 387475 506451 9 35286 3 13
13 10 Arizona ...Cardina 2000 893926 387475 506451 10 52244 3 13
14 # ... with 10,836 more rows, 9 more variables: points_against <dbl>,
15 #   points_differential <dbl>, margin_of_victory <dbl>,
16 #   strength_of_schedule <dbl>, simple_rating <dbl>, offensive_ranking <dbl>,
17 #   defensive_ranking <dbl>, playoffs <chr>, sb_winner <chr>, and abbreviated
18 #   variable names 1team_name, 2weekly_attendance, 3points_for
19 # Use 'print(n = ...)' to see more rows, and 'colnames()' to see all variable
   names

```

```

1 attendance_df <- attendance_joined >
2 filter(!is.na(weekly_attendance)) >
3 select(
4   weekly_attendance, team_name, year, week,
5   margin_of_victory, strength_of_schedule, playoffs
6 )
7
8 attendance_df

```

```

1 # A tibble: 10,208 × 7
2   weekly_attendance team_name year week margin_of_victory strength...1
3   ...2playo
4   <dbl> <chr>    <dbl> <dbl> <dbl> <dbl> <chr>
5   1 77434 Cardinals 2000 1 -14.6 -0.7 No
6   ...Pla

```

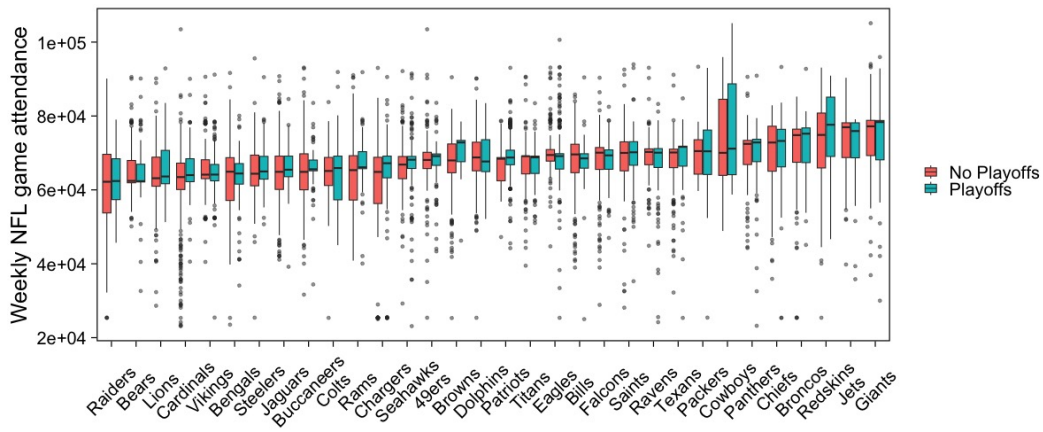


FIGURA 4. Asistencia en estadios por equipo.

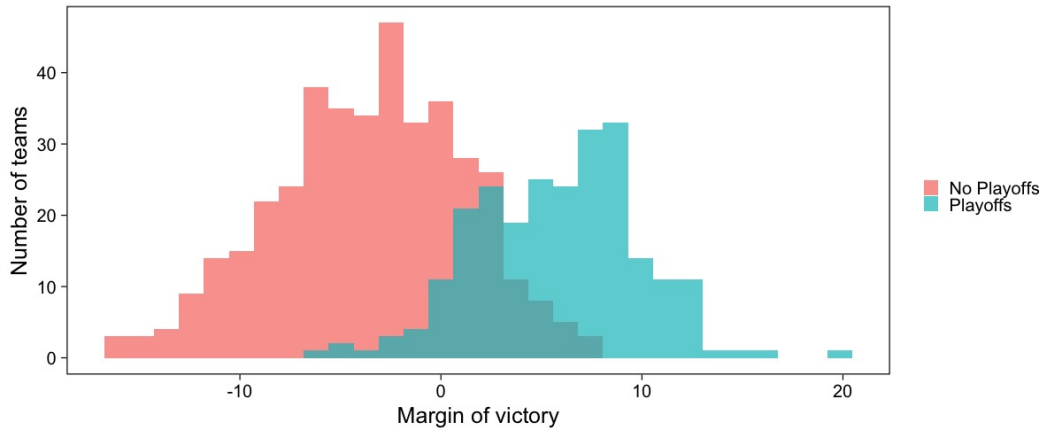


FIGURA 5. Diferencia de puntos en partidos ganados.

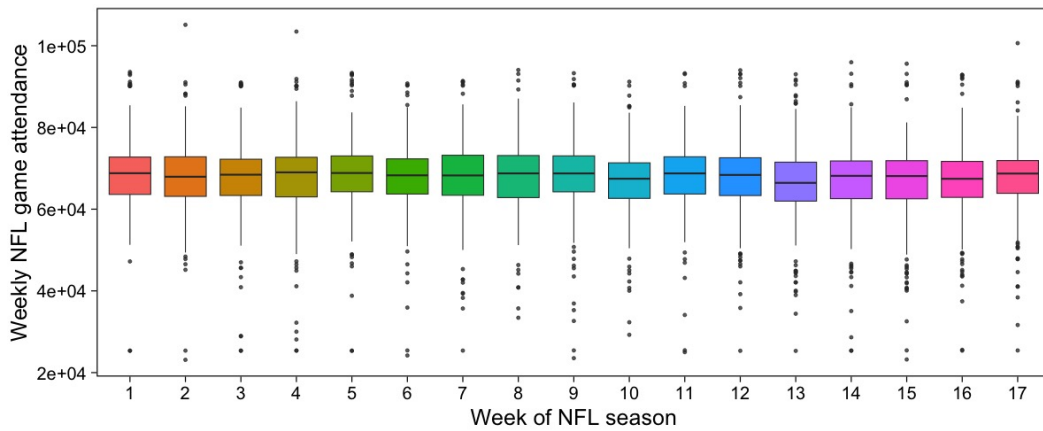


FIGURA 6. Asistencia a lo largo de la temporada.

```

5  2      ...Pla      66009 Cardinals  2000      2      -14.6      -0.7 No
6  3      ...Pla      71801 Cardinals  2000      4      -14.6      -0.7 No
7  4      ...Pla      66985 Cardinals  2000      5      -14.6      -0.7 No
8  5      ...Pla      44296 Cardinals  2000      6      -14.6      -0.7 No
9  6      ...Pla      38293 Cardinals  2000      7      -14.6      -0.7 No
10 7      ...Pla      62981 Cardinals  2000      8      -14.6      -0.7 No
11 8      ...Pla      35286 Cardinals  2000      9      -14.6      -0.7 No
12 9      ...Pla      52244 Cardinals  2000     10      -14.6      -0.7 No
13 10     ...Pla      64223 Cardinals  2000     11      -14.6      -0.7 No
14 # ... with 10,198 more rows, and abbreviated variable names
15 # 1 strength_of_schedule, 2playoffs
16 # Use 'print(n = ...)' to see more rows

```

```

1 library(tidymodels)
2
3 set.seed(108727)
4 attendance_split <- attendance_df >
5   initial_split(strata = playoffs)
6
7 nfl_train <- training(attendance_split)
8 nfl_test <- testing(attendance_split)

```

```

1 lm_spec <- linear_reg() >
2   set_engine(engine = "lm")
3
4 lm_spec

```

```

1 Linear Regression Model Specification (regression)
2
3 Computational engine: lm

```

```

1 lm_fit <- lm_spec >
2   fit(weekly_attendance ~ .,
3     data = nfl_train
4   )
5
6 lm_fit > broom::tidy()

```

```

1 # A tibble: 37 × 5
2   term                estimate std.error statistic  p.value
3   <chr>                <dbl>    <dbl>    <dbl>    <dbl>
4 1 (Intercept)        -79125.    33598.    -2.36  1.85e- 2
5 2 team_nameBears      -2613.     767.     -3.41  6.64e- 4
6 3 team_nameBengals    -4757.     771.     -6.17  7.08e-10

```

```

7 4 team_nameBills      106.      766.      0.138  8.90e- 1
8 5 team_nameBroncos   2889.      775.      3.73   1.93e- 4
9 6 team_nameBrowns    -21.6      775.     -0.0278 9.78e- 1
10 7 team_nameBuccaneers -2796.      752.     -3.72   2.04e- 4
11 8 team_nameCardinals -5905.      767.     -7.70   1.53e-14
12 9 team_nameChargers  -5098.      774.     -6.59   4.76e-11
13 10 team_nameChiefs   1802.      763.      2.36   1.83e- 2
14 # ... with 27 more rows
15 # Use 'print(n = ...)' to see more rows

```

```

1 lm_fit > broom::tidy(conf.int = TRUE)

```

```

1 # A tibble: 37 × 7
2   term                estimate std.error statistic  p.value conf.low conf.
3   <chr>                <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>
4 1 (Intercept)        -79125.    33598.    -2.36   1.85e- 2 -144985.
5 2 team_nameBears      -2613.      767.     -3.41   6.64e- 4  -4117.
6 3 team_nameBengals    -4757.      771.     -6.17   7.08e-10  -6268.
7 4 team_nameBills       106.      766.      0.138  8.90e- 1  -1396.
8 5 team_nameBroncos    2889.      775.      3.73   1.93e- 4   1371.
9 6 team_nameBrowns     -21.6      775.     -0.0278 9.78e- 1 -1541.
10 7 team_nameBuccaneers -2796.      752.     -3.72   2.04e- 4  -4271.
11 8 team_nameCardinals  -5905.      767.     -7.70   1.53e-14  -7408.
12 9 team_nameChargers   -5098.      774.     -6.59   4.76e-11  -6615.
13 10 team_nameChiefs    1802.      763.      2.36   1.83e- 2    306.
14 # ... with 27 more rows
15 # Use 'print(n = ...)' to see more rows

```

```

1 lm_fit > broom::glance()

```

```

1 # A tibble: 1 × 12
2   r.squared adj.r...1. sigma ...2stati  p.value  df  logLik  AIC  BIC
3   <dbl>    <dbl> <dbl>    <dbl>    <dbl> <dbl>    <dbl> <dbl> <dbl> <
4 1 0.152    0.148 8379.    38.1 7.30e-242    36 -80005. 1.60e5 1.60e5 5.35
5 # ... with 2 more variables: df.residual <int>, nobs <int>, and abbreviated
6 # variable names 1adj.r.squared, 2statistic, 3deviance
7 # Use 'colnames()' to see all variable names

```

```

1 results_train <- lm_fit >
2   predict(new_data = nfl_train) >
3   mutate(truth = nfl_train$weekly_attendance,
4          model = "lm")
5 results_train >
6   rmse(truth = truth, estimate = .pred)

```

```

1 # A tibble: 1 × 3
2   .metric .estimator .estimate
3   <chr>   <chr>       <dbl>
4 1 rmse    standard      8358.

```

```

1 results_test <- lm_fit >
2   predict(new_data = nfl_test) >
3   mutate(truth = nfl_test$weekly_attendance,
4          model = "lm")
5 results_test >
6   rmse(truth = truth, estimate = .pred)

```

```

1 # A tibble: 1 × 3
2   .metric .estimator .estimate
3   <chr>   <chr>       <dbl>
4 1 rmse    standard      8192.

```

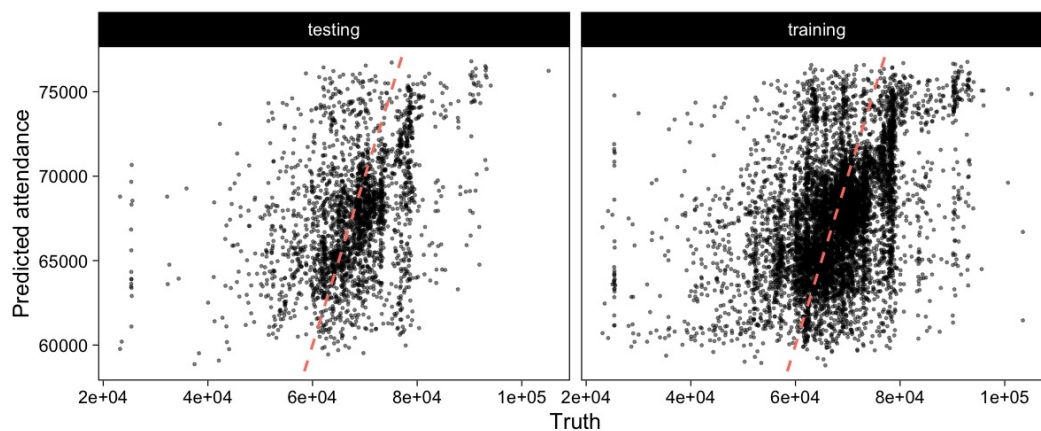


FIGURA 7. Comparativo de predicciones contra valores reales.

REFERENCIAS

- [1] G. James, D. Witten, T. Hastie, and R. Tibshirani. *An Introduction to Statistical Learning: With Applications in R*. Springer Texts in Statistics. Springer US, New York, NY, 2021. [1](#)
- [2] M. Kuhn and J. Silge. *Tidy Modeling with R*. O'Reilly Media, Inc., 2022. [8](#)