

Análisis de estadístico sobre jugadores de la MLB

Informe final de Estadística

Realizado por:
Gavazut E. ; Riera L. ; Cordero M.





Resumen

El deporte como actividad social del ser humano no es ajena a la ciencia, en particular, a la matemática y la estadística. En beisbol, por ejemplo, se recopila cada mínimo información de todo lo que sucede durante el juego, aspectos como: tasa de bateo, carreras anotadas o ponches. Y es que estos datos permiten medir cuan bueno o acertado es el desempeño de cada jugador. Debido a la gran cantidad de información, que además se ha incrementado con los años, es necesario recurrir a la ciencia y a modelos computacionales de predicción que ofrezcan un punto de objetividad que permita a los equipos mejorar su competitividad.

En este trabajo, se mostrará como realizar un análisis estadístico sobre una base de datos de jugadores de la Major League Baseball (MLB), relativo a la tasa de bateo, carreras anotadas, triples, dobles y ponches por veces al bate. De este análisis se destaca el comprobar como la tasa de ponches es mayor a la tasa de jonrones de un jugador y que es posible hallar una relación lineal entre las tasas de bateo y las tasas de carreras anotadas, de dobles y ponches por veces al bate. Esto permitirá predecir con un nivel 0.8589 de error cuadrático ajustado, cual será la tasa de carreras anotadas por jugador según su desempeño en el campo. Más aún, un análisis de varianza (ANOVA), permite demostrar que no hay mayor distinción entre jugadores con diferentes tasas de bateo



Planteamiento del problema

En el presente proyecto, el objetivo es tomar una base de datos con diversas métricas que corresponden a jugadores de la MLB para realizar los siguientes estudios sobre ella:

1. Análisis descriptivo.
2. Intervalo de confianza de 97% para la media de cada variable.
3. Probar (a nivel de 0.05) que el promedio de bateo es inferior a 0.300.
4. Estudiar si la tasa de ponches y de jonrones son iguales.
5. Prueba de bondad de ajuste para la tasa de bateo para determinar si tiene distribución normal.
6. Gráfico de dispersión y matriz de correlación para las variables.
7. Modelo de regresión final y predicción para la tasa de bateo.
8. Separar a la tasa de bateo en tres grupos: los que tienen menos de 0.200, los que tienen entre 0.200 y 0.300, y los que tienen más de 0.300, y realización de un análisis de varianza para estudiar si los promedios de tasas de las otras variables son iguales.



Descripción de la base de datos

La base de datos a estudiar cuenta con 45 observaciones de 6 variables, las cuales son:

- **X1 = tasa de bateo (hits/veces al bate).** Entiéndase la conexión efectuada por el bateador que coloca la pelota dentro del terreno de juego, permitiéndole alcanzar al menos una base, sin que se produzca un error de defensa del equipo contrario o algún otro jugador sea declarado como fuera de juego.
- **X2 = tasa de carreras anotadas (carreras anotadas/veces al bate).** Entiéndase carrera por anotación, y se logra al recorrer un corredor la totalidad de las bases volviendo al home, bien de manera continua (por medio de un jonrón) o de forma alternada consecutiva antes de que se realicen 3 outs.
- **X3 = tasa de dobles (dobles/veces al bate).** Entiéndase por doble como un hit en el que el bateador logra llegar a segunda base sin ser puesto out y sin que haya error alguno de la defensiva.



Descripción de la base de datos

- **X4 = tasa de triples (triples/veces al bate).** Entiéndase por triple como un hit en el que el bateador logra llegar satisfactoriamente a tercera base, sin que ocurra ningún error por parte de la defensiva.
- **X5 = tasa de jonrones (jonrones/veces al bate).** Un jonrón se da cuando el bateador hace contacto con la pelota de una manera que le permita recorrer las bases y anotar una carrera (junto con todos los corredores en base) en la misma jugada, sin que se registre ningún out ni error de la defensa.
- **X6 = tasa de ponches (ponches/veces al bate).** Por último, un ponche es la acción de retirar a un bateador con una cuenta de 3 strikes, al que la recibe se le suele llamar ponchao o ponchado.

De esta forma, vemos que cada una de las variables miden números bastante relevantes para cada jugador. Como cada una de estas estadísticas pueden ocurrir una sola vez mientras se está al bate, cada una será un número entre el 0 y el 1



Metodología



Para la realización de esta investigación se hará uso del software estadístico R en el entorno de desarrollo integrado (IDE) RStudio.

En este se iniciará por una descripción de los datos y variables almacenadas en el archivo fuente *Baseball.xlsx*, tales como: mínimo, media, cuantiles y desviación estándar.

Para la media de las variables se obtendrá un intervalo de confianza del 95%. Como se desea estudiar la relación de la tasa de bateo respecto al resto de las variables, se buscará determinar la mejor distribución de probabilidad que se ajuste a esta variable.

Finalmente, se estudiará la eficiencia del mejor modelo lineal de predicción que se ajuste a los datos y permita establecer si en efecto existe tal relación entre las variables y las implicaciones que tendría en las estrategias para futuros juegos de beisbol.



Análisis de los datos



Para la realización de este proyecto se contó con una archivo de excel con la información de algunos jugadores de la Major League Beisbol o MLB, el cual se almacenó en una variable llamada `Baseball`.

De esta archivo podemos realizar el siguiente análisis de datos.





¿Qué clase es la base de datos?

Con el comando `class`, se pudo determinar el tipo de base de datos utilizada o lo que es equivalente, la clase de la variable `Baseball`.

El resultado que se obtuvo indica que es del tipo `tbl_df`, que es una subclase de la clase `data.frame`. `tbl_df` cumple con tener propiedades diferentes por defecto y se suele referir a ellas como `tibble`.

Es una clase eficiente para trabajar con bases de datos grandes y su visualización.



Variables en la base de datos

Si se desea saber que tipo de variables están almacenadas en la base de datos, se puede utilizar el comando `str`.

Esta función nos indica que se cuentan con 6 variables denominadas `X1`, `X2`, `X3`, `X4`, `X5`, `X6`, y distribuidas de tal manera que representan la columnas de la base de datos.

Cada una de estas variables tienen 45 valores de tipo `double` o número decimal, que representan las 45 observaciones aleatorias (una por fila) realizadas a jugadores de la (MLB).



Estadísticos

Para obtener los estadísticos de las seis (6) variables de esta base de datos, se inicia por guardar las 45 observaciones en un vector que represente a cada variable.

Y aplicar las siguientes funciones:

- `mean` que permite obtener la media de los datos,
- `median` para obtener la mediana,
- `quantile` para retornar los cuantiles al 0.25%, 0.50% y 0.75% de cada variable,
- `min` para el valor mínimo,
- `max` para el valor máximo,
- `var` para la varianza,
- `sd` que es para la desviación estándar,
- `IQR` es para el rango intercuartil,
- `stad/media` el coeficiente de variación

Estadísticos



	Mínimo	25 %	Media	Mediana (50 %)	75 %	Máximo	RIC	Varianza	Desv. Estándar	Coef. Variación
X1	0.188	0.248	0.2805	0.290	0.308	0.367	0.060	0.0019	0.0440	0.1569
X2	0.064	0.119	0.1509	0.150	0.189	0.259	0.070	0.0018	0.0420	0.2784
X3	0.025	0.039	0.0464	0.045	0.053	0.068	0.014	0.0001	0.0105	0.2255
X4	0.001	0.007	0.0113	0.009	0.016	0.030	0.009	0.0000	0.0070	0.6165
X5	0.000	0.009	0.0243	0.013	0.039	0.085	0.030	0.0005	0.0223	0.9173
X6	0.000	0.062	0.1043	0.095	0.138	0.264	0.076	0.0040	0.0631	0.6044

Tabla 1: Resumen Estadístico de las variables



Estadísticos

De estos resultados hay varios puntos que podemos destacar:

La varianza de los datos es muy baja indicativo que entre los datos hay pocos valores atípicos o datos muy dispersos, lo que se refleja en valores mas cercanos a la media. Misma interpretación se puede extender a la desviación estándar pues es la raíz cuadrada de la varianza.

Una consecuencia de la baja varianza es que la media y la mediana son valores muy cercanos. Esto es particularmente útil al analizar el valor del RIC, que toma como medida central la mediana de los datos. Es decir, nos indica donde se encuentra el 50% de los datos, cuánto mas bajo es el valor del RIC menos dispersos están los datos.

Diagramas e histograma de los datos por cada variable



Histograma: Tasa de hits

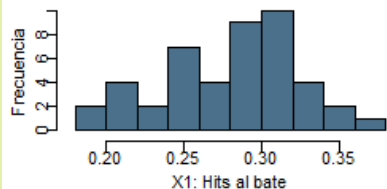
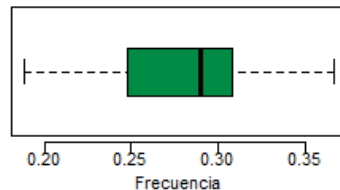


Gráfico de Cajas: Tasa de hits



Histograma: Tasa de triples

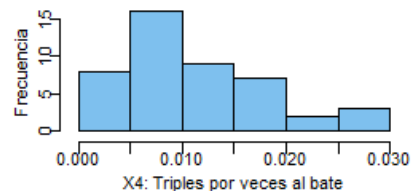
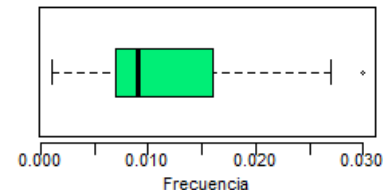


Gráfico de Cajas: Tasa de triples



Histograma: Tasa carreras anotadas

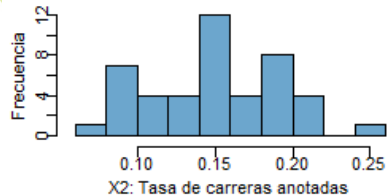
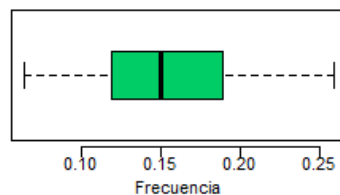


Gráfico de Cajas: Tasa carreras anotadas



Histograma: Tasa de jonrones

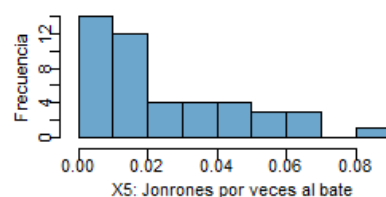
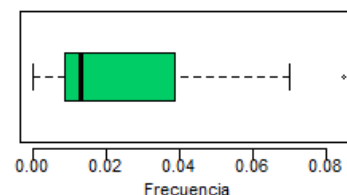


Gráfico de Cajas: Tasa de jonrones



Histograma: Tasa de dobles

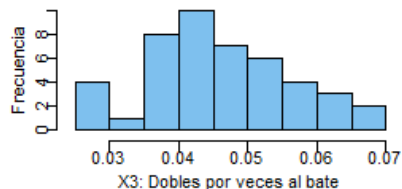
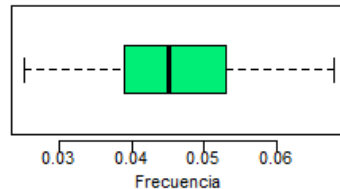


Gráfico de Cajas: Tasa de dobles



Histograma: Tasa de ponches

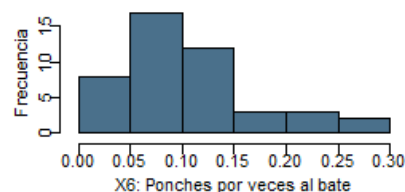
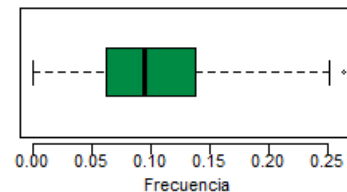


Gráfico de Cajas: Tasa de ponches





Diagramas e histograma de los datos por cada variable

Para la variable X_1 , podemos ver que los valores máximos de los datos se obtienen luego de la media, pero el mayor volumen de ellos se encuentra antes como bien se observa en el diagrama de caja que permite confirmar, además, la ausencia de datos atípicos.

Para la variable X_2 , se puede comprobar que hay simetría de los datos que se podía apreciar en la tabla de estadísticos. Simetría, particular, respecto al valor 0.15 que coincide a su vez con la media de los datos.

El diagrama de caja permite confirmar la ausencia de los valores atípicos.



Diagramas e histograma de los datos por cada variable

Por su parte, para la variable X_3 y X_4 . Vemos que en general, ambos diagramas de caja son bastante parecidos, con la única diferencia siendo que el de triples está 0.03 puntos corrido hacia arriba y los datos desde el primer cuartil hasta la mediana están muchos más dispersos.

Otra diferencia es que el diagrama de cajas para los triples no cuenta con datos atípicos, en cambio los dobles si, que corresponde a 0.3. Todo esto hace que el diagrama de los triples sea casi simétrico, y el de los dobles sea más chato entre el valor mínimo y la mediana, en comparación con lo que tenemos entre la mediana y el máximo valor.



Diagramas e histograma de los datos por cada variable

De la gráfica para la variable X_5 podemos ver como a medida que nos vamos acercando a 1, la frecuencia de jonrones decae rápidamente, mientras que al inicio es muy alta.

De la gráfica para la variable X_6 podemos ver que la mayoría de los jugadores se ponchan menos de un 15% de las veces que estan al bate.



Intervalo de confianza para la media de las variables

Con el uso de la función `t.test()` se puede encontrar el intervalo de confianza con una significancia de 0.03 o (97% de confianza) para las variables estudiadas.

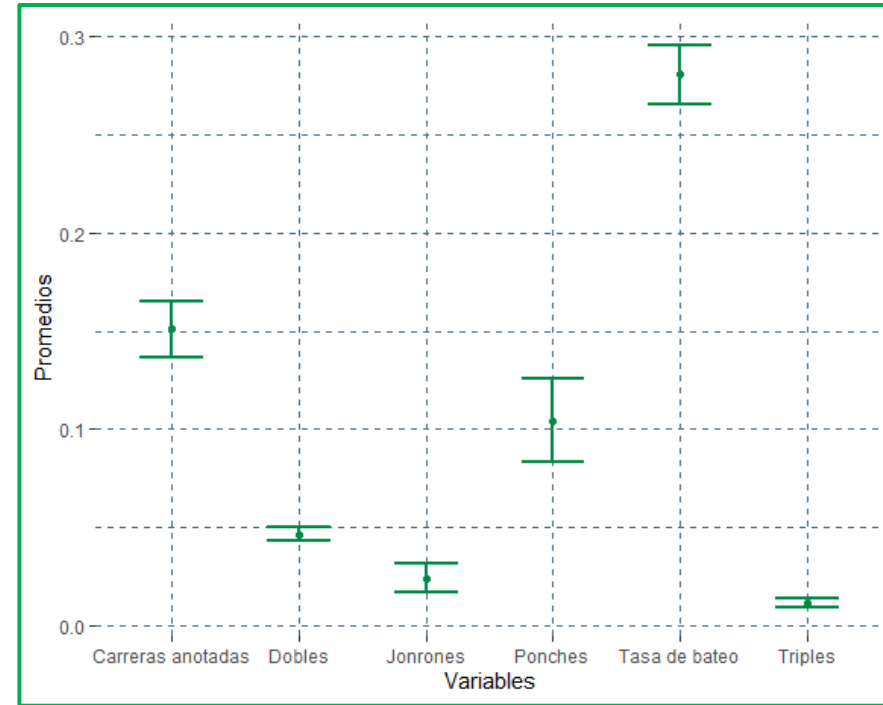
Los resultados de aplicar esta función, se pueden visualizar en la tabla siguiente tabla y gráfico.

Intervalo de confianza para la media de las variables



	Limite inferior	Promedio	Limite Superior
Tasa de bateo	0.2658	0.2805	0.2952
Carreras anotadas	0.1368	0.1509	0.1649
Dobles	0.0429	0.0464	0.0498
Triples	0.0090	0.0113	0.0136
Jonrones	0.0168	0.0243	0.0317
Ponches	0.0833	0.1043	0.1254

Tabla 2: Intervalos de confianza para las medias de las variables





Intervalo de confianza para la media de las variables

Note que general, los intervalos de confianza más estrechos son los de dobles y triples, lo que nos indica que, con una probabilidad del 97%, podemos asegurar que los jugadores de la MLB tendrán un promedio de triples y dobles que puede ser estimado con bastante certeza.

Pero vemos que las carreras anotadas, los ponches y la tasa de bateo tienen un intervalo de confianza mucho más grande, por lo que no podemos asegurar que el promedio será estimado de forma tan certera.



Promedio de bateo

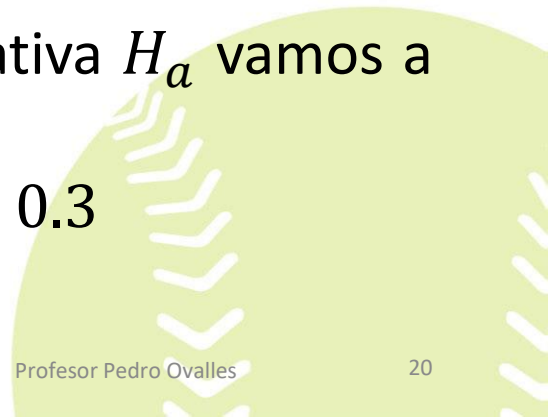


Con lo obtenido en los intervalos de confianza del apartado anterior se tiene que la tasa de bateo toma valores por debajo de 0.300.

Para corroborar este resultado, se realizará un prueba de hipótesis con un nivel de significancia de $\alpha = 0.05$.

Como hipótesis nula H_0 y como hipótesis alternativa H_a vamos a suponer que

$$H_0: \mu_{\text{bateo}} \leq 0.3, \quad H_a: \mu_{\text{bateo}} > 0.3$$





Promedio de bateo



Si suponemos que los datos presentan una distribución normal, podemos aplicar el comando `t.test` de `R`, que permite realizar pruebas de hipótesis sobre las medias de los datos cuando se trabaja con una sola variable.





Promedio de bateo



Con esta función, se obtuvo que el valor para el **estadístico t** es -23.811 , con 44 grados libertad.

El **p – valor** es bastante alto, es igual $0,9976$ (que representa un 99.76%). Y se cumple que $\alpha = 0.05 < 99.76$.

Por lo tanto, la hipótesis alternativa se rechaza, mas aún, se rechaza para todo nivel de significancia porque se necesita un valor para α más alto que el *p – valor*, para rechazar la hipótesis nula.

Se afirma entonces, con seguridad, que la tasa de bateo es inferior a 0.300 , tal como se podía apreciar con el intervalo de confianza.



Comparación entre las tasas de ponches y las de jonrones

Ahora, deseamos comparar las tasas de ponches y de jonrones para determinar si son o no parecidas. Como no tenemos conocimiento acerca de las varianzas poblacionales, usaremos el test de Welch tal y como es explicado en *Heumann, Schomaker (2017)* para comparar las medias.

En este caso, haremos una prueba de hipótesis, donde tomaremos como hipótesis nula, H_0 e hipótesis alternativa H_a las dadas por:

$$H_0: \mu_{\text{jonrones}} - \mu_{\text{ponches}} = 0 \quad \text{vs.} \quad H_a: \mu_{\text{jonrones}} - \mu_{\text{ponches}} \neq 0$$



Comparación entre las tasas de ponches y las de jonrones

Es decir, queremos determinar si las tasas de jonrones y ponches son distintas. Ahora, con apoyo del comando anterior `t.test()`, pero esta vez para comparar dos variables, podremos determinar cuál de estas hipótesis es aceptada.

Como resultado se obtuvo que el $p\text{-valor} = 1.112 \times 10^{-8}$, que es extremadamente pequeño, mucho más que el nivel de significancia $\alpha = 0.01$ que es razonable utilizar para nuestra prueba de hipótesis.

Adicionalmente, el intervalo de confianza que se obtuvo fue de $(-0.1068, -0.0593)$ que no incluye el cero.



Comparación entre las tasas de ponches y las de jonrones

Otra cosa que podemos hacer es evaluar el estadístico de prueba con el comando `qt()` (vemos por lo anterior que $dt = 55$ y $\alpha = 0.05$). Por lo que $t = -8$, vemos que el estadístico cae en la región de rechazo (porque es de cola doble).

Para cualquiera de estos casos, podemos concluir que la hipótesis nula se rechaza, es decir que hay suficiente evidencia para creer que $\mu_{\text{jonrones}} - \mu_{\text{ponches}} \neq 0$.

Y además, como el intervalo de confianza es negativo, concluimos que $\mu_{\text{ponches}} > \mu_{\text{jonrones}}$ con un nivel de confianza del 95%, como se podía apreciar en los intervalos de confianza de la media.



Prueba de bondad de ajuste para la distribución de X_1

Para continuar con el análisis a un nivel más profundo, resulta conveniente determinar si los datos en la variable X_1 , sobre la tasa de bateos, sigue una distribución normal.

Para esto, primero note que en el histograma para la variable, se obtuvo que si se subdivide en intervalos de longitud 0.02, las frecuencias son como las descritas en las siguientes tabla:

Prueba de bondad de ajuste para la distribución de X_1

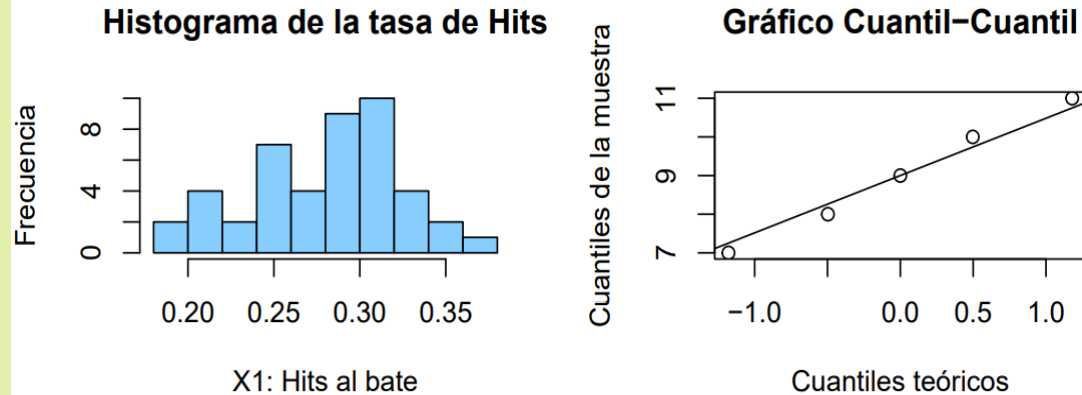


Figura 1: Histograma y Gráfico cuantil-cuantil de la variable X_1

<i>Intervalo</i>	0,18- 0,20	0,22- 0,24	0,26- 0,28	0,30-0,32	0,34-0,36	0,38
<i>Frecuencia</i>	2	4	7	10	2	1

Tabla 3: Tabla de clases y frecuencias

Prueba de bondad de ajuste para la distribución de X_1



Ahora, agruparemos los datos en categorías de frecuencia mayor o igual a 5 (para poder aplicar el método de bondad de ajuste) tal como puede apreciarse en la tabla siguiente:

<i>Clases</i>	[0,18 0,24)	[0,24 0,28)	[0,28 0,30)	[0,30 0,32)	[0,32 0,38)
<i>Frecuencia</i>	8	11	9	10	7

Tabla 4: Nueva agrupación en clases con frecuencia mayor o igual a 5

Con la gráfica cuantil-cuantil de la figura anterior, podemos ver que esta agrupación se ajusta bien a una distribución normal (representada por la recta).



Prueba de bondad de ajuste para la distribución de X_1

Vamos a proceder a realizar una prueba χ^2 , que es una prueba de hipótesis que compara la distribución observada de los datos con la distribución esperada de los datos.

Para este tipo de pruebas, el estadístico de χ^2 cuantifica que tanto varía la distribución respecto a la distribución hipotética.

La hipótesis nula H_0 y la hipótesis alternativa H_a vienen dadas por:

H_0 : Los datos siguen una distribución normal

H_a : Los datos no siguen una distribución normal



Prueba de bondad de ajuste para la distribución de X_1

Como estadístico χ^2 tenemos:

$$\chi^2 = \sum_{i=1}^k \frac{[n_i - E(n_i)]^2}{E(n_i)}$$

con $k = 5$ el número de clases o categorías, n_i las frecuencias de cada categoría, $E(n_i) = n * p_i$, el valor esperado con n el número total de datos y p_i la probabilidad de cada clase n_i .



Prueba de bondad de ajuste para la distribución de X_1

Para calcular las probabilidades p_i se obtuvo la media y la desviación estándar de los datos agrupados como

$$\bar{x} = 0.2822 \text{ y } \sigma = 0.045,$$

respectivamente. Con \bar{x} y σ se obtuvieron las siguientes probabilidades para cada clase:

$$p_1 = 0.172, \quad p_2 = 0.3082, \quad p_3 = 0.1747, \quad p_4 = 0.1466, \\ p_5 = 0.1986.$$



Prueba de bondad de ajuste para la distribución de X_1

Sustituyendo los datos en el estadístico tenemos que:

$$\chi^2 = 2.9421,$$

y el p – *valor* viene dado por $1 - P(\chi^2 < 2.9421) = 0.2297$.

El p – *valor* es bastante alto por lo que la hipótesis nula no se rechaza para ningún nivel de significancia.

Por tanto, los datos siguen una distribución normal con media 0.2822 y desviación estándar 0.045.

Gráfico de dispersión y matrix de correlación

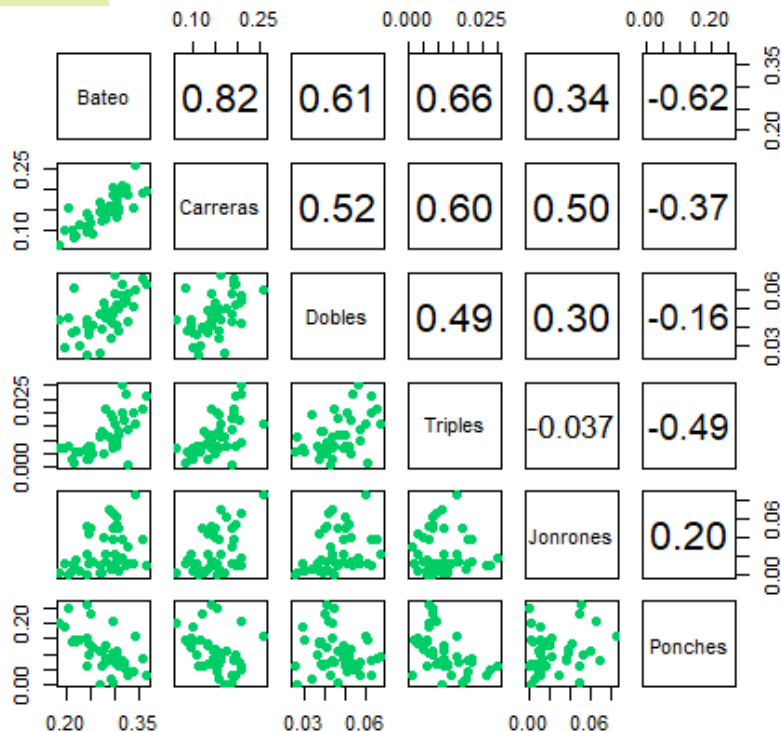




Gráfico de dispersión y matrix de correlación



Es ahora, de nuestro interés estudiar la relación entre las variables de la base de datos. Esto lo podemos observar en la figura anterior.

Note que de las gráficas de dispersión de la mitad inferior de la figura se puede apreciar que para **carreras anotadas, dobles y triples** tenemos algo que se asemeja a una relación lineal positiva.

Mientras que para los ponches, estos disminuyen a medida que la tasa de bateo aumenta.

La única variable que no parece **tener ninguna relación clara** con la tasa de bateo es la **tasa de jonrones**, por lo que es una variable que probablemente no nos ofrezca mayor información si queremos establecer un modelo lineal que relacione a las variables.

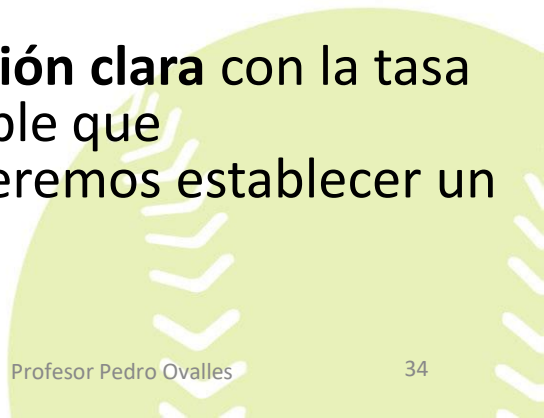




Gráfico de dispersión y matrix de correlación



Por otro lado, con la parte superior de la figura se tienen los coeficientes de correlación por pares de variables.

Estos coeficientes nos indican que, efectivamente, para las **carreras, dobles y triples**, tenemos una correlación positiva (siendo las carreras la que tiene mayor correlación, y los triples la menor).

Además, para los ponches tenemos una correlación negativa bastante significativa, y entre todas las variables, los jonrones tienen la menor correlación.



Muestreo 80%-20%

Por lo visto en la matriz de correlación, parece existir una relación lineal entre las variables, particularmente vamos a estar interesados en ver como se relaciona cada campo de información (carreras, dobles, etc.) con la variable X_1 que es la tasa de bateo.

Con R tenemos la posibilidad de obtener un modelo de regresión lineal con la función `lm`.



Muestreo 80%-20%

Pero para asegurarnos que el modelo sea el más adecuado, primero necesitamos extraer una muestra que permita entrenar al modelo de predicción, y con los datos restantes probar que tan eficiente es el modelo.

Con este objetivo, se dividen los datos en un 80% para el entrenamiento y en un 20% para las pruebas.



Muestreo 80%-20%

Como la base de datos consta de 45 observaciones por variable, el 80% representa tomar una muestra aleatoria de **36 observaciones**, por lo que el 20% restante serán las **9 observaciones** no tomadas en la muestra.

Vale la pena resaltar que se habla de observaciones, o las filas de la base de datos y no de las entradas particulares de cada variable porque se busca estudiar la relación por jugador, de su tasa de bateo, respecto a su tasa de carreras, dobles, triples, jonrones y ponches.

En otras palabras, **las filas son independientes** entre sí y por eso se pueden tomar muestras al azar, pero **las columnas no lo son** por ser datos relativos a un jugador en particular.



Modelo de regresión lineal para la variable X_1

Ahora, teniendo seleccionado nuestros datos, podemos pasar a realizar el modelo.

La mejor manera de realizar un modelo de regresión lineal es seguir el método de **regresión paso a paso**, y de esta manera determinar cuáles variables son significativas o no al tomar en cuenta la tasa de bateo.



Modelo de regresión lineal para la variable X1



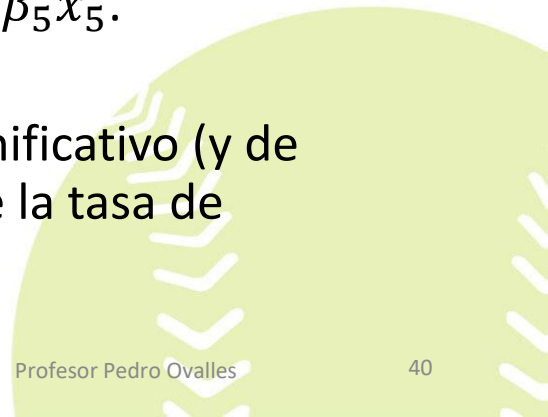
Ahora, pasemos a realizar el modelo lineal utilizando el comando `lm()` de R. Se desarrolla primero el modelo dado por

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5 + \epsilon.$$

Suponiendo que $E(\epsilon) = 0$, buscamos estimar los parámetro β_i para los cuales

$$E[Y] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5.$$

Con el comando `lm()` se obtuvo que el único valor no significativo (y de hecho el p-valor más alto) fue la tasa de triples, seguido de la tasa de jonrones que era significativa a nivel 0.05.





Modelo de regresión lineal para la variable X_1

PRUEBA 2:

De esta forma, realicemos de nuevo el modelo pero sin la variable X_4 correspondiente a los triples. Es decir, el modelo a estimar es:

$$E[Y] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_4 x_4 + \beta_5 x_5.$$

Con \mathbb{R} se obtuvo en esta prueba, que la tasa de jonrones es la variable con p-valor mas alto, con 0.0611.

A pesar, de ser significativa a nivel de 0.1 procedemos a realizar una nueva prueba, esta vez sin la tasa de jonrones.



Modelo de regresión lineal para la variable X1

PRUEBA 3:

El nuevo modelo, consiste en estimar

$$E[Y] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_5 x_5 .$$

Ahora, todas nuestras variables son bastante significativas, por lo que sus p-valores son bastante pequeños, significativos a nivel 0.001.

Los valores estimados fueron:

$$\beta_0 = 0.1630, \beta_1 = 0.5192, \beta_2 = 1.3650 \text{ y } \beta_5 = -0.2451.$$

Como medida del error, tenemos el R^2_{ajus} , con valor 0.8489, indicando que hay un buen ajuste de los datos al modelo.



Modelo de regresión lineal para la variable X1

Además, tenemos que:

- Para los estimadores, los dobles es el mayor de todos, y este nos indica que por cada aumento del 1 % en la tasa de dobles, hay un aumento correlacionado del 136 % en la tasa de bateo. Es interesante ver que este estimador es muchísimo mayor que el de las carreras.
- La varianza es estimada como $\sigma^2 = 0,017052$.
- Para el error estándar (Std. Error), podemos construir los intervalos de confianza para las variables. Primero, tenemos que $t_{\{32,0.975\}} = 2.0369$:
 - $I_{\{dobles\}} = 1.3650 \pm 2.0369 * 0.3471 = (0.6580, 2.0720)$
 - $I_{\{ponches\}} = -0.2451 \pm 2.0369 * 0.0460 = (-0.3388, -0.1514)$

Modelo de regresión lineal para la variable X_1



Ahora, veamos que efectivamente se cumple con las característica de un buen modelo apoyándonos en las gráficas.

Modelo de regresión lineal para la variable X1

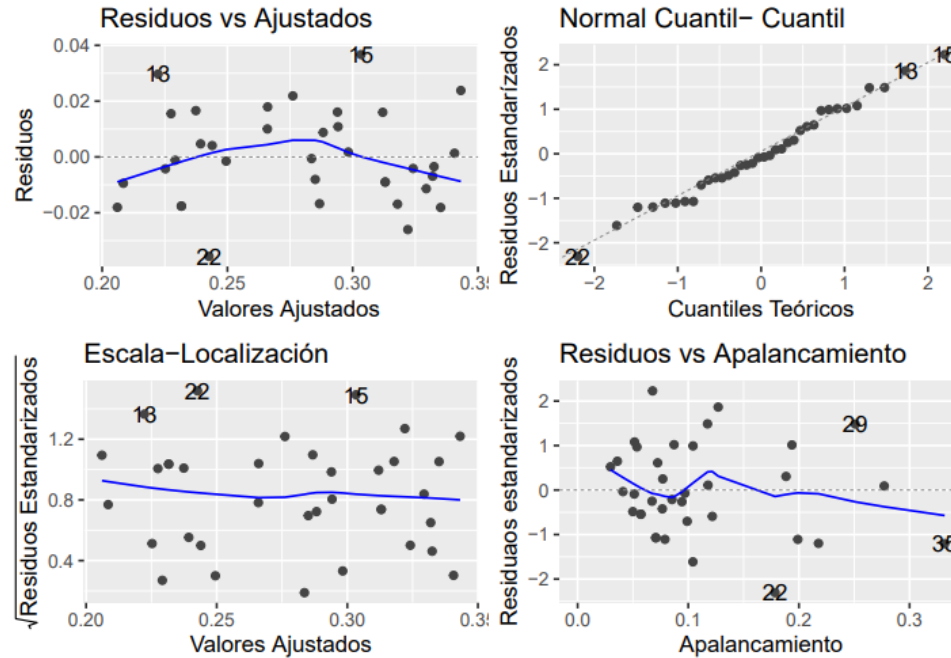


Figura 3: Graficos descriptivos del modelo



Modelo de regresión lineal para la variable X_1

- Cuando vemos la gráfica de "**Residuos vs Ajustados**", nos damos cuenta de que la línea azul es bastante horizontal, y esta además está centrada alrededor del cero, es decir que podemos asumir que no hay independencia entre las variables y la tasa de bateo.
- Al ver el gráfico "**Normal Cuantil-Cuantil**", vemos que todos los valores están bastante cercanos a la recta, lo que nos confirma la normalidad.
- En "**Escala-Localización**" no vemos ningún patrón, lo que nos indica que los valores presentan homocedasticidad.
- Y por último, en "**Residuos vs Apalancamiento**", no hay ningún valor que esté fuera de las líneas rayadas, por lo que no parece haber valores que generen apalancamiento.

En conclusión, podemos ver que este es un buen modelo, cuyas variables son todas significativas, no tiene datos que generen apalancamiento y cumple con homocedasticidad.

En resumen, nuestro modelo es:

$$\hat{Y} = 0.1630 + (0.5192)x_2 + (1.3605)x_3 - (0.2451)x_5$$



Prueba y predicción del modelo lineal

Ahora, haremos uso del comando `predict` para hacer la predicción de la variable X_1 (tasa de hits), utilizando las 9 observaciones que se seleccionaron previamente.

Luego calculamos la diferencia entre los valores reales y los valores estimados por el modelo.

Los resultados se muestran en la siguiente tabla.

Prueba y predicción del modelo lineal



Tabla 5: Hits reales vs Hits predichos

Tasa de hits real	Tasa de hits predicha	Diferencia
0.281	0.2817371	-0.0007371
0.290	0.2638287	0.0261713
0.269	0.3130776	-0.0440776
0.307	0.3310601	-0.0240601
0.308	0.2455406	0.0624594
0.358	0.2242784	0.1337216
0.245	0.3030698	-0.0580698
0.294	0.2396552	0.0543448
0.218	0.2968785	-0.0788785

Es claro que los residuos son bastante pequeños, así que se considera que el modelo es suficientemente bueno para predecir la tasa de hits.



ANOVA para la Tasa de Bateo

Para finalizar, estamos interesados en realizar un análisis de varianza sobre la variable X_1 o la tasa de bateo, para compararla con el resto de las variables.

Particularmente queremos realizar, el estudio sobre 3 categorías o grupos:

- Grupo 1: los bateadores con una tasa de bateo igual a ($X_1 < 0.200$).
- Grupo 2: los bateadores con una tasa de bateo igual a ($0.200 \leq X_1 < 0.300$).
- Grupo 3: los bateadores con una tasa de bateo igual a ($0.300 \leq X_1$)

Con esta agrupación se opta por realizar un análisis de varianza con bloques aleatorizados, donde los bloques serán los grupos y los tratamientos o métodos serán las distintas variables de la base de datos.

ANOVA para la Tasa de Bateo



Tabla 6: Tabla a dos factores para las medias de las variables

	X1	X2	X3	X4	X5	X6
Grupo1	0.1935000	0.0820000	0.0365000	0.0070000	0.00650	0.1965000
Grupo2	0.2580400	0.1327600	0.0414400	0.0087200	0.02228	0.1204800
Grupo3	0.3212778	0.1837222	0.0542778	0.0153333	0.02900	0.0716667



ANOVA para la Tasa de Bateo

Con la tabla anterior podemos apreciar las medias de los valores agrupados.

Con estos valores, se puede aplicar el comando `anova` de R para obtener la tabla ANDEVA, tal y como se detalla en la tabla siguiente.

ANOVA para la Tasa de Bateo



Tabla 7: Tabla ANDEVA

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
grupos	2	0.0019826	0.0009913	0.5020692	0.6197546
variables	5	0.1334024	0.0266805	13.5132135	0.0003524
Residuals	10	0.0197440	0.0019744	NA	NA



ANOVA para la Tasa de Bateo

En esta tabla, se aprecia que el p-valor para los grupos es de 0.6198 que es alto, indicando que **la hipótesis nula para los grupos no se puede rechazar**, es decir, que las medias por grupos son iguales.

Sin embargo, para las medias clasificadas por variable o método se obtuvo un p-valor de 0.0004 que es bastante bajo, incluso significativo indicando que las medias son distintas tal como se esperaba por los datos analizados.

Con esto podemos afirmar con seguridad, que los promedios de las tasas son iguales por cada grupo.

Conclusiones



De todo el análisis anterior se deduce que:

1. No hay demasiada variabilidad entre las diferentes tasas de bateo, por lo que en general los jugadores de la MLB proyectan rendimientos similares (aunque esto depende del grado de exactitud con el que se quiera medir).
2. La tasa de bateo es en media al menos mas del doble que la tasa de ponches para cualquier jugador (esto se sigue de la tabla 2).
3. La tasa de hits sigue aproximadamente un distribución normal centrada en 0.2822 y con desviación estándar de 0.045.
4. Las variables más significativas (entre las estudiadas), para predecir la tasa de hits o bateos son la tasa de carreras, la tasa de dobles, y la tasa de ponches, con estas se puede lograr un buen modelo lineal.

An aerial view of a large baseball stadium at night, filled with spectators. The field is illuminated, showing the green grass and brown dirt base paths. A large green graphic overlay is on the left side, featuring a white baseball bat and a white baseball. The text "¡Gracias!" is centered over the field.

¡Gracias!