

# Reporte: Modelo de Regresión Lineal

## Introducción

Este reporte documenta el proceso de implementación de un modelo de regresión lineal basado en el conjunto de datos "beisbol.csv".

El objetivo del modelo es predecir la cantidad de carreras (`runs`) utilizando como principales variables predictoras el número de bateos (`bateos`) y el equipo al que pertenece.

## Justificación del Algoritmo

La regresión lineal fue seleccionada como algoritmo base debido a su simplicidad, interpretabilidad y adecuación para problemas que implican relaciones lineales entre variables. El análisis inicial de los datos indicó una correlación moderada entre el número de bateos y las carreras, lo que sugiere que un modelo lineal podría ser un punto de partida razonable. Además, la facilidad de implementación y optimización de la regresión lineal la hace ideal para un primer enfoque a este problema.

## Descripción del Diseño del Modelo

El diseño del modelo incluye un pipeline que combina preprocesamiento de datos y entrenamiento del modelo de regresión lineal.

El preprocesamiento se encargó de:

- Escalar la variable numérica `bateos` mediante estandarización.
- Codificar la variable categórica `equipos` utilizando OneHotEncoding.

El modelo final se entrenó utilizando un conjunto de entrenamiento (80% de los datos) y se evaluó con un conjunto de prueba (20% restante).

## Modelo de Regresión Lineal con Pipeline

```
[5]: # Separar variables predictoras y objetivo
X_filtered = filtered_data[["equipos", "bateos"]]
y_filtered = filtered_data["runs"]

# Configurar preprocesador y pipeline
preprocessor = ColumnTransformer(
    transformers=[
        ("num", StandardScaler(), ["bateos"]),
        ("cat", OneHotEncoder(drop="first", handle_unknown="ignore"), ["equipos"]),
    ]
)

pipeline = Pipeline(
    steps=[
        ("preprocessor", preprocessor),
        ("regressor", LinearRegression())
    ]
)

# Dividir los datos en conjuntos de entrenamiento y prueba
X_train_f, X_test_f, y_train_f, y_test_f = train_test_split(X_filtered, y_filtered, test_size=0.2, random_state=42)

# Entrenar el modelo
pipeline.fit(X_train_f, y_train_f)

# Evaluar el modelo
y_pred_f = pipeline.predict(X_test_f)
mse_f = mean_squared_error(y_test_f, y_pred_f)
r2_f = r2_score(y_test_f, y_pred_f)

print(f"Mean Squared Error (MSE): {mse_f}")
print(f"R^2 Score: {r2_f}")
```

## Evaluación y Optimización del Modelo

El modelo inicial mostró un bajo desempeño, con un coeficiente de determinación ( $R^2$ ) negativo y un error cuadrático medio (MSE) elevado.

Se implementaron las siguientes optimizaciones:

- Filtrado de valores atípicos para las variables `bateos` y `runs`.
- Manejo de categorías desconocidas en la codificación de `equipos` mediante OneHotEncoder.

A pesar de estas optimizaciones, el modelo mantuvo un desempeño limitado debido a la ausencia de características predictoras adicionales en el conjunto de datos.

```
Mean Squared Error (MSE): 7973.666667125901
R^2 Score: -1.3180012275633288
```

## Gráfica Personalizada e Interpretación de Resultados

A continuación se presenta una gráfica que ilustra la relación entre los valores predichos por el modelo y los valores reales observados en el conjunto de prueba.

La interpretación de esta gráfica sugiere que el modelo lineal no es capaz de capturar completamente las dinámicas de los datos disponibles, lo que refuerza la necesidad de características adicionales o el uso de modelos más complejos.

```
# Gráfica de valores reales vs predicciones
plt.scatter(y_test_f, y_pred_f, alpha=0.7, label="Datos")
plt.plot(
    [y_test_f.min(), y_test_f.max()],
    [y_test_f.min(), y_test_f.max()],
    'k--', lw=2, label="Línea Ideal"
)
plt.xlabel("Valores Reales")
plt.ylabel("Predicciones")
plt.title("Regresión Lineal: Predicciones vs Valores Reales")
plt.legend()
plt.show()
```

