

② Jueves

- Python

> Anaconda > Jupyter

> Importar bibliotecas

> Conexión con SQL

> Query donde se hace join STATS con goles_casa y goles_vis

> Eliminar columnas duplicadas (team y season)

> Ver la correlación de las variables con el target "wins"

> Nuevo DF con variables ~~correlacionadas~~ con mayor correlación (goles_casa, goles_visita, total_pacs, corner_tallen, touches, season).

> ~~ver minimos, maximos, nulos, outliers~~

> Limpieza datos, minimos, maximos, nulos, outliers

> Revisar la distribución de los datos

> Train y Test del Modelo Lineal

> DF 2 donde este la primera a la penultima temporada

1° Train

> Dropear target "wins"

> Dividir categoricos y numericos

> Normalizar numericos

> Estandarizar numericos

> Tratar categoricos "season" con OneHot Encoding

↓

> Concatenar con X_norm o X-standarized con Encoding

> ~~Split~~ Split train test

2° Test

> Crear Modelo

> Entrenar modelo

> Predecir data

> R^2

↓

normalizadas (X_norm)

$R^2 = .75$

estandarizadas (X-standarized)

$R^2 = .73$

3° cross-validation

> Elegir modelo \rightarrow X_norm

> Hacer validación del modelo con los datos de la última temporal. (Cross-Validation).

Comparar Resultados de Predicción de "wins" con las verdaderas "wins"