

EDA Exploratory Data Analysis	ETL: Extract, Transform, Load	Machine Learning: Preparación	Tests estadísticos	Normalización																
<h3>Análisis exploratorio de datos</h3> <p>El Análisis Exploratorio de Datos se refiere al proceso de realizar una serie de investigaciones inciales sobre los datos que tenemos para poder descubrir patrones, detectar anomalías, probar hipótesis y comprobar suposiciones con la ayuda de estadísticas y representaciones gráficas.</p>	<h3>Extraccion</h3> <ul style="list-style-type: none"><li>- obtener datos crudos y almacenarlos<ul style="list-style-type: none"><li>- Tablas de bases de datos SQL o NoSQL</li><li>- Ficheros de texto plano</li><li>- Emails</li><li>- Información de páginas web</li><li>- Hojas de cálculo</li><li>- Ficheros obtenidos de API's</li></ul></li></ul> <h3>Transformación</h3> <ul style="list-style-type: none"><li>- procesar los datos, unificarlos, limpiarlos, validarlos, filtrarlos, etc.<ul style="list-style-type: none"><li>- Formetear fechas</li><li>- Reordenar filas o columnas</li><li>- Unir o separar datos</li><li>- Combinar las fuentes de datos</li><li>- Limpiar y estandarizar los datos</li><li>- Verificar y validar los datos</li><li>- Eliminar duplicados o datos erroneos</li><li>- Filtrado, realización de calculos o agrupaciones</li></ul></li></ul> <h3>Carga</h3> <ul style="list-style-type: none"><li>- cargar los datos en su formato de destino, el tipo de lo cual dependerá de la naturaleza, el tamaño y la complejidad de los datos. Los sistemas más comunes suelen ser:<ul style="list-style-type: none"><li>- Ficheros csv</li><li>- Ficheros json</li><li>- Bases de datos</li><li>- Almacenes de datos (Data Warehouse)</li><li>- Lagos de datos (Data Lakes)</li></ul></li></ul>	<h3>Hipotesis Nula y Errores Tipo I y II</h3> <h4>Hipótesis nula (H0)</h4> <ul style="list-style-type: none"><li>- en general es la afirmación contraria a la que queremos probar</li></ul> <h4>Hipótesis alternativa (H1)</h4> <ul style="list-style-type: none"><li>- en general la afirmación que queremos comprobar</li></ul> <h4>p-valor</h4> <ul style="list-style-type: none"><li>- medida de la probabilidad de que una hipótesis nula sea cierta</li><li>- valor entre 0 y 1</li><li>- si *p-valor* &lt; 0.05 ✗ Rechazamos la hipótesis nula.</li><li>- si *p-valor* &gt; 0.05 ✔ Aceptamos la hipótesis nula.</li></ul> <h4>Error Tipo I:</h4> <ul style="list-style-type: none"><li>- rechazar la hipótesis nula cuando es verdadera</li></ul> <h4>Error Tipo II:</h4> <ul style="list-style-type: none"><li>- aceptar la hipótesis nula cuando es falsa</li></ul>	<h3>Independencia</h3> entre variables predictoras <ul style="list-style-type: none"><li>- las variables predictoras tienen que ser independientes para poder crear un modelo de regresión lineal</li></ul> <h3>Variables numéricas: Correlaciones</h3> <ul style="list-style-type: none"><li>- pairplot</li><li><b>sns.pairplot(df)</b></li><li>- covarianza</li><li><b>df_numéricas.cov()</b></li><li>- correlación de Pearson (relación lineal)</li><li><b>df_numéricas.corr()</b></li><li>- correlación de Spearman (relación no lineal)</li><li><b>df_numéricas.corr(method = 'spearman')</b></li><li>- correlación de Kendall (datos numéricos pero categóricos y ordinales)</li><li><b>df_numéricas.corr(method = 'kendall')</b></li></ul> <h3>Variables categóricas: Chi-cuadrado</h3> <ul style="list-style-type: none"><li>- V-Cramer: varía entre 0 y 1<ul style="list-style-type: none"><li>- más cerca a 1 más dependientes</li><li>- <b>resultado &lt; 0,7 para hacer ML ✔</b></li></ul></li></ul> <b>import researchpy as rp</b> <b>crosstab, test_results, expected = rp.crosstab</b> <b>(df["col1"], df["col2"], test= "chi-square", expected_freqs= True, prop= "cell")</b> <b>test_results</b> devuelve los resultados del test en un dataframe	<h3>Método manual</h3> <ul style="list-style-type: none"><li>- cogemos el valor que queremos normalizar y restamos la media de la columna, y dividimos el resultado por el maximo restado por el mínimo de la columna</li></ul> <b>df["col_norm"] = (df["col_VR"] - df["col_VR"].media()) / (df["col_VR"].max() - df["col_VR"].min())</b> <h3>Método logarítmica</h3> <ul style="list-style-type: none"><li>*no se puede hacer si algún valor sea 0*</li></ul> <b>df["col_norm"] = df["col_VR"].apply(lambda x: np.log(x) if x &gt; 0 else 0)</b> <h3>Método raiz cuadrada</h3> <b>import math</b> <b>df["col_norm"] = df["col_VR"].apply(lambda x: math.sqrt(x))</b> <h3>Método stats.boxcox()</h3> <p>aplica una transformación logarítmica para los valores positivos y exponencial para valores negativos de nuestra columna</p> <b>from scipy import stats</b> <b>df["col_norm"], lambda ajustada = stats.boxcox(df["col_VR"])</b> <h3>Método MinMaxScaler</h3> <b>from sklearn.preprocessing import MinMaxScaler</b> <b>modelo = MinMaxScaler(feature_range=(0,1), copy=True)</b> <b>modelo.fit(df["col_VR"])</b> <b>datos_normalizados = modelo.transform(df["col_VR"])</b> <b>df_datos_norm = pd.DataFrame(datos_normalizados, columns = ['col_norm'])</b> <b>df['col_norm'] = df_datos_norm</b>																
<h3>1. Entender las variables</h3> <ul style="list-style-type: none"><li>- que variables temenos</li><li>.head(), .tail(), .describe(), .info(), .shape</li><li>- que tipos de datos</li><li>.dtypes(), .info()</li><li>- si temenos nulos o duplicados</li><li>.isnull().sum()</li><li>.duplicated().sum()</li><li>- que valores unicos temenos</li><li>.unique(), .value_counts()</li></ul>																				
<h3>2. Limpiar el dataset</h3> <ul style="list-style-type: none"><li>- quitar duplicados (filas o columnas)</li><li>- cambiar nombres de columnas</li><li>- cambiar tipo de datos de columnas</li><li>- ordenar columnas</li><li>- separar columna en dos con str.split()</li><li>- crear intervalos con pd.cut()</li><li>- crear porcentajes o ratios</li><li>- decidir como tratar outliers: mantenerlos, eliminarlos, o reemplazarlos con la media, mediana o moda; o aplicar una imputacion</li></ul> <ul style="list-style-type: none"><li>- decidir como tratar nulos:<ul style="list-style-type: none"><li>- eliminar filas o columnas con nulos drop.na()</li><li>- imputar valores perdidos:<ul style="list-style-type: none"><li>- reemplazarlos con la media, mediana o moda usando .fillna() o .replace()</li><li>- imputer con metodos de machine learning usando la libreria sklearn: Simple-Imputer, Iterative-Imputer, o KNN Imputer</li></ul></li></ul></li></ul>	<h3>APIs</h3> <b>import requests</b> libreria para realizar petitions HTTP a una URL, para hacer web scraping <b>url = 'enlace'</b> el enlace de la que queremos extraer datos <b>header = {}</b> opcional; contiene informacion sobre las peticiones realizadas (tipo de ficheros, credenciales) <b>response = requests.get(url=url, header = header)</b> pedimos a la API que nos de los datos <b>variables = {'parametro1': 'valor1', 'parametro2': 'valor2'}</b> <b>response = request.get(url=url, params=variables)</b> pedimos a la API que nos de los datos con los parametros segun el diccionario de parametros que le pasamos <b>response.status_code</b> devuelve el status de la peticion <b>response.reason</b> devuelve el motive de codigo de estado <b>response.text</b> devuelve los datos en formato string <b>response.json()</b> devuelve los datos en formato json <b>df = pd.json_normalize(response.json)</b> devuelve los datos en un dataframe																			
<h3>3. Analizar relaciones entre variables</h3> <p>Analizar relaciones entre las variables</p> <ul style="list-style-type: none"><li>- para encontrar patrones, relaciones o anomalías</li></ul> <p>Relaciones entre dos variables numéricas:</p> <ul style="list-style-type: none"><li>- scatterplot</li><li>- regplot - scatterplot con línea de regresion</li><li>- matriz de correlación y heatmap</li><li>- joinplot - permite emparejar dos gráficas - una histograma con scatter o reg plot por ejemplo</li></ul> <p>Relaciones entre dos variables categóricas:</p> <ul style="list-style-type: none"><li>- countplot</li></ul> <p>Relaciones entre variables numéricas y categóricas:</p> <ul style="list-style-type: none"><li>- swarmplot</li><li>- violinplot</li><li>- pointplot</li><li>- boxplot</li></ul>	<h3>Codigos de respuesta de HTTP</h3> <table><tr><td>1XX informa de una respuesta correcta</td><td>4XX error durante peticion</td></tr><tr><td>2XX codigo de exito</td><td>401 peticion incorrecta</td></tr><tr><td>200 OK</td><td>402 sin autorizacion</td></tr><tr><td>201 creado</td><td>403 prohibido</td></tr><tr><td>202 aceptado</td><td>404 no encontrado</td></tr><tr><td>204 sin contenido</td><td>5XX error del servidor</td></tr><tr><td>3XX redireccion</td><td>501 error interno del servidor</td></tr><tr><td></td><td>503 servicio no disponible</td></tr></table>	1XX informa de una respuesta correcta	4XX error durante peticion	2XX codigo de exito	401 peticion incorrecta	200 OK	402 sin autorizacion	201 creado	403 prohibido	202 aceptado	404 no encontrado	204 sin contenido	5XX error del servidor	3XX redireccion	501 error interno del servidor		503 servicio no disponible			
1XX informa de una respuesta correcta	4XX error durante peticion																			
2XX codigo de exito	401 peticion incorrecta																			
200 OK	402 sin autorizacion																			
201 creado	403 prohibido																			
202 aceptado	404 no encontrado																			
204 sin contenido	5XX error del servidor																			
3XX redireccion	501 error interno del servidor																			
	503 servicio no disponible																			

## Machine Learning: Preprocesamiento

### Estandarización

- cambiar los valores de nuestras columnas de manera que la desviación estándar de la distribución sea igual a 1 y la media igual a 0; para que las VP sean comparables

#### Método manual

```
df["col_esta"] = (df ["col_VR"] - df ["col_VR"].media()) / (df ["col_VR"].std())
```

#### Sklearn StandardScaler

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaler.fit(df_num_sin_VR)
datos_estandarizados = scaler.transform (df_num_sin_VR)
df_datos_esta = pd.DataFrame(datos_estandarizados, columns = df_num_sin_VR.columns)
```

#### Sklearn RobustScaler

```
from sklearn.preprocessing import RobustScaler
scaler = RobustScaler()
scaler.fit(df_num_sin_VR)
datos_estandarizados = scaler.transform (df_num_sin_VR)
df_datos_esta = pd.DataFrame(datos_estandarizados, columns = df_num_sin_VR.columns)
```

### Encoding

#### Variables categóricas

**Ordinaria:** no requiere números pero si consta de un orden o un puesto; diferencias de medianas entre categorías

**Nominal:** variable que no es representada por números, no tiene algún tipo de orden, y por lo tanto es matemáticamente menos precisa; no habrá grandes diferencias de medianas entre categorías

**Binaria:** dos posibilidades; puede tener orden o no

**Variables sin orden:** creamos una columna nueva por valor único, asignando unos y zeros

#### One-Hot Encoding

```
from sklearn.preprocessing import OneHotEncoder
oh = OneHotEncoder()
df_transformados = oh.fit_transform(df[['columna']])
oh_df = pd.DataFrame(df_transformados.toarray())
oh_df.columns = oh.get_feature_names_out()
df_final = pd.concat([df, oh_df], axis=1)

get_dummies
df_dum = pd.get_dummies(df['col'], prefix='prefijo', dtype=int)
df[df_dum.columns] = df.dum
df.drop('col', axis=1, inplace=True)
```

#### Variables que tienen orden:

Label Encoding asigna un número a cada valor único de una variable

```
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df['col_VR_le'] = le.fit_transform(df[col_VR'])

map() asigna el valor que queramos según el mapa que creamos
df['col_VR_map'] = df[col_VR'].map(diccionario)

Ordinal-Encoding asignamos etiquetas basadas en un orden o jerarquía
from sklearn.preprocessing import OrdinalEncoder
```

### Regresión Lineal: Métricas

```
from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
```

**R2:** representa la proporción de la varianza que puede ser explicada por las VP del modelo; mayor R2=mejor modelo

```
r2_score(y_train,y_predict_train)
r2_score(y_test,y_predict_test)
```

**MAE (Mean absolute error):** medida de la diferencia entre los valores predichos vs los reales; menor MAE=mejor modelo

```
mean_absolute_error(y_train,y_predict_train)
mean_absolute_error(y_test,y_predict_test)
```

**MSE (Mean Squared Error):** mide el promedio(media) de los errores al cuadrado; menor MSE=mejor modelo

```
mean_squared_error(y_train,y_predict_train)
mean_squared_error(y_test,y_predict_test)
```

**RMSE (Root Mean Squared Error):** distancia promedio entre los valores predichos y los reales; menor RMSE=mejor modelo

```
np.sqrt(mean_squared_error(y_train,y_predict_train))
np.sqrt(mean_squared_error(y_test,y_predict_test))
```

### Linear Regression: Modelo

- separar los datos de las variables predictoras (x) de la variable respuesta (y)
- dividimos los datos en datos de entrenamiento y datos de test con train\_test\_split()

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
```

- Ajustamos el modelo

```
from sklearn.linear_model import LinearRegression
lr = LinearRegression(n_jobs=-1)
lr.fit(x_train, y_train)
```

- Hacemos las predicciones

```
y_predict_train = lr.predict(x_train)
y_predict_test = lr.predict(x_test)
```

5. Guardamos los resultados en dataframes y los concatenamos

```
train_df = pd.DataFrame({'Real': y_train, 'Predicted': y_predict_train, 'Set': ['Train']*len(y_train)})
test_df = pd.DataFrame({'Real': y_test, 'Predicted': y_predict_test, 'Set': ['Test']*len(y_test)})
resultados = pd.concat([train_df,test_df], axis = 0)
```

6. creamos una columna de los residuos: la diferencia entre los valores observados y los de la predicción

```
resultados['residuos'] = resultados['Real'] - resultados['Predicted']
```

#### Cross-validation

```
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import cross_validate
```

**cv\_scores = cross\_val\_score(estimator = LinearRegression(), X = X, y = y, scoring = 'neg\_root\_mean\_squared\_error', cv = 10)**

**cv\_scores.mean()**

calcula la media de los resultados de CV de una métrica

**cv\_scores = cross\_validate(estimator = LinearRegression(), X = X, y = y, scoring = 'r2', 'neg\_root\_mean\_squared\_error', cv = 10)**

**cv\_scores["test\_r2"].mean()**

**cv\_scores["test\_neg\_root\_mean\_squared\_error"].mean()** calcula las medias de los resultados de validación de múltiples métricas

### Regresión Logística: Métricas

#### Matriz de confusión

Matriz de confusión		Predicción	
		Positivo	Negativo
Realidad	Positivo	Verdadero positivo	Falso negativo
	Negativo	Falso positivo	Verdadero negativo

para crear un heatmap de una matriz de confusión:

```
from sklearn.metrics import confusion_matrix
mat_lr = confusion_matrix(y_test, y_pred_test)
plt.figure(figsize = (n,m))
sns.heatmap(mat_lr, square=True, annot=True= plt.xlabel('valor predicho')
plt.ylabel('valor real')
plt.show()
```

#### Métricas

```
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score , cohen_kappa_score, roc_curve,roc_auc_score
```

**Accuracy (exactitud):** porcentaje de los valores predichos están bien predichos

**accuracy\_score(y\_train,y\_predict\_train)**

**accuracy\_score(y\_test,y\_predict\_test)**

**Recall:** porcentaje de casos positivos capturados

\*si preferimos FP, queremos recall alta\*

**recall\_score(y\_train,y\_predict\_train)**

**recall\_score(y\_test,y\_predict\_test)**

**Precisión (sensibilidad):** porcentaje de predicciones positivas correctas

\*si preferimos FN, queremos precisión alta\*

**precision\_score(y\_train,y\_predict\_train)**

**presicion\_score(y\_test,y\_predict\_test)**

**Especificidad:** porcentaje de los casos negativos capturados

**F1:** la media de la precisión y el recall

**f1\_score(y\_train,y\_predict\_train)**

**f1\_score(y\_test,y\_predict\_test)**

**kappa:** una medida de concordancia que se basa en comparar la concordancia observada en un conjunto de datos, respecto a la que podría ocurrir por mero azar

- <0 No acuerdo
- 0.0-0.2 Insignificante
- 0.2-0.4 Bajo
- 0.4-0.6 Moderado
- 0.6-0.8 Bueno
- 0.8-1.0 Muy bueno

**cohen\_kappa\_score(y\_train,y\_predict\_train)**

**cohen\_kappa\_score(y\_test,y\_predict\_test)**

**curva ROC:** forma gráfica de ver la kappa; la sensibilidad vs. la especificidad

**AUC (área under curve):** la área bajo la curva ROC; cuanto más cerca a 1, mejor será nuestro modelo clasificando los VP

### Balanceo para Regresión Logística

#### Downsampling

ajustar la cantidad de datos de la categoría mayoritaria a la minoritaria

#### Método manual

**df\_minoritaria = df[df['col'] == valor\_min]**

**df\_muestra = df[df['col'] == valor\_max].sample (num\_minoritarios, random\_state = 42)**

**df\_balanceado = pd.concat([df\_minoritaria, df\_muestra],axis = 0)**

#### Método RandomUnderSample

```
import imblearn
X = df.drop('col_VR', axis=1)
y = df['col_VR']
down_sampler = RandomUnderSampler()
X_down, y_down = down_sampler.fit_resample(X,y)
df_balanceado = pd.concat([X_down, y_down], axis = 1)
```

#### Método Tomek

**x = df.drop('col\_VR', axis=1)**

**y = df['col\_VR']**

**x\_train, x\_test, y\_train, y\_test = train\_test\_split(x, y, test\_size = 0.2, random\_state = 42)**

**tomek\_sampler = SMOTETomek()**

**X\_train\_res, y\_train\_res =**

**tomek\_sampler.fit\_resample(X\_train, y\_train)**

#### Upsampling

ajustar la cantidad de datos de la categoría minoritaria a la mayoritaria

#### Método manual

**df\_mayoritaria = df[df['col'] == valor\_may]**

**df\_muestra = df[df['col'] == valor\_min].sample (num\_mayoritarias, random\_state = 42)**

**df\_balanceado = pd.concat([df\_mayoritaria, df\_muestra],axis = 0)**

#### Método RandomOverSample

```
import imblearn
X = df.drop('col_VR', axis=1)
y = df['col_VR']
down_sampler = RandomUnderSampler()
X_down, y_down = down_sampler.fit_resample(X,y)
df_balanceado = pd.concat([X_down, y_down], axis = 1)
```

### Logistic Regression: Modelo

seguir los mismos pasos como para la Regresión Lineal pero con **LogisticRegression()**

```
from sklearn.linear_model import LogisticRegression
```

### Decision Tree: Modelo

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.ensemble import DecisionTreeRegressor
```

**from sklearn import tree**

seguir los mismos pasos como para la Regresión Lineal pero con **DecisionTreeRegressor()**

**arbol = DecisionTreeRegressor(random\_state=42)**

Para dibujar el árbol:

**fig = plt.figure(figsize = (10,6))**

**tree.plot\_tree(arbol, feature\_names = x\_train.columns, filled = True)**

**plt.show()**

### GridSearch y best\_estimator\_

Despues de hacer las predicciones de un modelo Decision Tree, examinamos lás métricas de los resultados:

- si temenos overfitting hay que reducir la profundidad del modelo

- si temenos underfitting hay que aumentar la profundidad del modelo

**max\_features = np.sqrt(len(x\_train.columns))**

podemos calcular el valor de max\_features siendo la raíz cuadrada del número de variables predictoras

**arbol.tree\_.max\_depth** nos muestra el max depth usado por defecto, para poder ajustarlo; deberíamos usar la mitad como mucho

- GridSearch ejecuta todas las posibles combinaciones de hiperparámetros que le damos con el parámetro 'param' y best\_estimator\_ devuelve la major combinacion encontrado

- Definimos un diccionario de los hiperparametros

```
param = {"max_depth": [n,m,l], "max_features": [a,b,c,d], "min_samples_split": [x,y,z], "min_samples_leaf": [r,s,t]}
from sklearn.model_selection import GridSearchCV
```

- Iniciamos el modelo con GridSearch

```
gs = GridSearchCV(estimator = DecisionTreeRegressor(), param_grid = param, cv=10, verbose=-1, return_train_score = True, scoring = "neg_mean_squared_error")
```

- Ajustamos el modelo en el GridSearch

**gs.fit(x\_train, y\_train)**

- Aplicamos el método de best\_estimator\_

**mejor\_modelo = gs.best\_estimator\_**

devuelve la mejor combinación de hiperparámetros

- Volvemos a sacar las predicciones

**y\_pred\_test\_dt2 = mejor\_modelo.predict(x\_test)**

**y\_pred\_train\_dt2 = mejor\_modelo.predict(x\_train)**

#### Importancia de los predictores

```
importancia_predictores = pd.DataFrame({'predictor': x_train.columns, 'importancia': mejor_modelo.feature_importances_})
```

**importancia\_predictores.sort\_values(by=["importancia"], ascending=False, inplace = True)** crea un dataframe con la relativa importancia de cada VP

- para los variables categóricas nominales a los cuales se ha aplicado encoding, hay que sumar los resultados de las columnas divididas:

**df\_sum = importancia\_predictores.esta.iloc[[n, m]]**

**importancia\_predictores.esta.drop(df\_sum.index, inplace = True)**

**importancia\_predictores.esta.loc[n] =**

**["nombre\_col", df\_sum["importancia"].sum()]**

### Random Forest: Modelo

seguir los mismos pasos como para el Decision Tree pero con **RandomForestRegressor()**

**from sklearn.ensemble import RandomForestRegressor**

- se puede usar los mismos hiperparámetros del best\_estimator\_ o volver a ejecutar el GridSearch