

EDA y ETL	ETL: Extract, Transform, Load	Machine Learning: Preparación	Tests estadísticos	Normalización
<div>EDA: Análisis exploratorio de datos</div> <div>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.</div>	<div>Extraccion<ul style="list-style-type: none">- obtener datos crudos y almacenarlos<ul style="list-style-type: none">- Tablas de bases de datos SQL o NoSQL- Ficheros de texto plano- Emails- Información de páginas web- Hojas de cálculo- Ficheros obtenidos de API's</div> <div>Transformación<ul style="list-style-type: none">- procesar los datos, unificarlos, limpiarlos, validarlos, filtrarlos, etc.<ul style="list-style-type: none">- Formetear fechas- Reordenar filas o columnas- Unir o separar datos- Combinar las fuentes de datos- Limpiar y estandarizar los datos- Verificar y validar los datos- Eliminar duplicados o datos erroneos- Filtrado, realización de calculos o agrupaciones</div> <div>Carga<ul style="list-style-type: none">- 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">- Ficheros csv- Ficheros json- Bases de datos- Almacenes de datos (Data Warehouse)- Lagos de datos (Data Lakes)</div>	<div>Hipotesis Nula y Errores Tipo I y II</div> <div><div>Hipótesis nula (H0)<ul style="list-style-type: none">- en general es la afirmación contraria a la que queremos probar</div><div>Hipótesis alternativa (H1)<ul style="list-style-type: none">- en general la afirmación que queremos comprobar</div><div>p-valor<ul style="list-style-type: none">- medida de la probabilidad de que una hipótesis nula sea cierta- valor entre 0 y 1- si *p-valor* < 0.05 ✗ Rechazamos la hipótesis nula.- si *p-valor* > 0.05 ✓ Aceptamos la hipótesis nula.</div><div>Error Tipo I:<ul style="list-style-type: none">- rechazar la hipótesis nula cuando es verdadera</div><div>Error Tipo II:<ul style="list-style-type: none">- aceptar la hipótesis nula cuando es falsa</div></div> <div>Tests estadísticos</div> <div><div>Normalidad<ul style="list-style-type: none">- la variable respuesta tiene que tener una distribución normal para poder crear un modelo de regresión lineal</div><div>Visualmente:<ul style="list-style-type: none">- histograma o distribución- grafico de cuantiles teóricos (Q-Q)más alineados están los puntos entorno a la recta, más normales serán nuestros datos<div>import statsmodels.api as sm</div><div>sm.qqplot(datos, line ='45')</div></div><div>Metodos analiticos:</div><div>Asimetría<ul style="list-style-type: none">- distribuciones asimétricas positivas: media > mediana y moda- distribuciones asimétricas negativas: media < mediana y moda<div>from scipy.stats import skew</div><div>skew(datos_normales)</div><div>método de scipy que calcula el sesgo</div><div>df['columna'].skew()</div><div>método de pandas que calcula el sesgo</div></div><div>Curtosis<ul style="list-style-type: none">- leptocurtosis: valor de curtosis mayor que 0 (pico alto)- mesocurtosis: valor de curtosis igual a 0 (pico medio)- platikurtosis: valor de curtosis menor que 0 (plana)<div>from scipy.stats import kurtosistest</div><div>kurtosistest(datos)</div><div>devuelve un p-valor</div><div>- p-valor del test > 0.05: datos normales ✓</div><div>- p-valor del test < 0.05: datos NO normales</div></div><div>Test de Shapiro-Wilk<ul style="list-style-type: none">- para muestras < 5000<div>hipótesis nula: distribución normal</div><div>from scipy import stats</div><div>stats.shapiro(df["datos"])</div><div>- p-valor del test > 0.05: datos normales ✓</div><div>- p-valor del test <) 0.05: datos NO normales</div></div><div>Test de Kolmogorov-Smirnov<ul style="list-style-type: none">- para muestras > 5000<div>hipótesis nula: distribución normal</div><div>from scipy import kstest</div><div>kstest(df["datos"], 'norm')</div><div>- p-valor del test > 0.05: datos normales ✓</div><div>- p-valor del test < p-valor (alfa) 0.05: datos NO normales</div></div></div>	<div>Independencia entre variables predictoras<ul style="list-style-type: none">- las variables predictoras tienen que ser independientes para poder crear un modelo de regresión lineal</div> <div>Variables numéricas: Correlaciones<ul style="list-style-type: none">- pairplot<div>sns.pairplot(df)</div>- covarianza<div>df_numéricas.cov()</div>- correlación de Pearson (relación lineal)<div>df_numéricas.corr()</div>- correlación de Spearman (relación no lineal)<div>df_numéricas.corr(method = 'spearman')</div>- correlación de Kendall (datos numéricos pero categóricos y ordinales)<div>df_numéricas.corr(method = 'kendall')</div></div> <div>Variables categóricas: Chi-cuadrado<ul style="list-style-type: none">- V-Cramer: varía entre 0 y 1<ul style="list-style-type: none">- más cerca a 1 más dependientes- resultado < 0,7 para hacer ML ✓<div>import researchpy as rp</div><div>crosstab, test_results, expected = rp.crosstab</div><div>(df["col1"], df["col2"], test= "chi-square",</div><div>expected_freqs= True, prop= "cell")</div><div>test_results</div><div>devuelve los resultados del test en un dataframe</div></div> <div>Homocedasticidad (homogeneidad de varianzas)<ul style="list-style-type: none">- las variables predictoras tienen que tener homogeneidad de varianzas en comparación con la variable respuesta</div> <div>Visualmente:<ul style="list-style-type: none">- violinplot- boxplot- regplot (columnas numéricas vs variable respuesta)</div> <div>Metodos analiticos:<ul style="list-style-type: none">- test de Levene (más robusto ante falta de normalidad) o Bartlett<div>from scipy import stats</div><div>from scipy.stats import levene</div><div>Variables categóricas:<ul style="list-style-type: none">- hay que crear un dataframe para cada valor único de las columnas categóricas<div>df_valor1 = df[df['col1'] == 'valor1']['col_VR']</div><div>df_valor2 = df[df['col1'] == 'valor2']['col_VR']</div><div>levene_test = stats.levene(df_valor1, df_valor2,</div><div>center='median')</div><div>bartlett_test = stats.bartlett(df_valor1, df_valor2,</div><div>center='median')</div></div><div>Variables numéricas:<ul style="list-style-type: none">- hay que crear un dataframe de las columnas numéricas sin la variable respuesta<div>for col in df_numericas.columns:</div><div>statistic, p_val = levene(df[col], df['col_VR'],</div><div>center='median')</div><div>resultados[col] = p_val</div></div><div>devuelve los p-valores en un diccionario<ul style="list-style-type: none">- p-valor del test > 0.05: varianzas iguales, homocedasticidad ✓- p-valor del test < 0.05: varianzas diferentes, heterocedasticidad</div></div>	<div>Método manual<ul style="list-style-type: none">- 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<div>df["col_norm"] = (df["col_VR"] -</div><div>df["col_VR"].media()) / (df["col_VR"].max() -</div><div>df["col_VR"].min())</div></div> <div>Método logarítmica<ul style="list-style-type: none">*no se puede hacer si algún valor sea 0*<div>df["col_norm"] = df["col_VR"].apply(lambda x:</div><div>np.log(x) if x > 0 else 0)</div></div> <div>Método raiz cuadrada<div>import math</div><div>df["col_norm"] = df["col_VR"].apply(lambda x:</div><div>math.sqrt(x))</div></div> <div>Método stats.boxcox()<div>aplica una transformación logarítmica para los valores positivos y exponencial para valores negativos de nuestra columna</div><div>from scipy import stats</div><div>df["col_norm"] = df["col_VR"].apply(lambda x:</div><div>stats.boxcox(df["col_VR"])</div></div> <div>Método MinMaxScaler<div>from sklearn.preprocessing import MinMaxScaler</div><div>modelo = MinMaxScaler(feature_range=(0,1),</div><div>copy=True)</div><div>modelo.fit(df["col_VR"])</div><div>datos_normalizados = modelo.transform(df["col_VR"])</div><div>df_datos_norm = pd.DataFrame(datos_normalizados,</div><div>columns = ['col_norm'])</div><div>df['col_norm'] = df_datos_norm</div></div> <div>ANOVA<div>import statsmodels.api as sm</div><div>from statsmodels.formula.api import ols</div><div>lm = ols('col_VR ~ col_VP1 + col_VP2 + col_VP3',</div><div>data=df).fit()</div><div>devuelve un dataframe de los resultados:</div><div>df (degrees of freedom): para variables categóricas será el número de valores únicos menos 1; para variables numéricas será siempre 1</div><div>sum_sq: medida de variación/desviación de la media</div><div>mean_sq: es el resultado de dividir la suma de cuadrados entre el número de grados de libertad.</div><div>F: un test que se utiliza para evaluar la capacidad explicativa que tiene la variable predictora sobre la variación de la variable respuestae</div>- PR(>F): si el p-valor < 0.05 es una variable significativa; que puede afectar a la VR<div>lm.summary()</div><div>devuelve una resumen de los resultados:</div>coef: representa los cambios medios en la VR para una unidad de cambio en la VP mientras se mantienen constantes el resto de las VP; los signos nos indican si esta relación es positiva o negativastd err: cuanto menor sea el error estándar, más precisa será la estimaciónt: es el resultado de dividir el coeficiente entre su error estándar</div>

