12- Evaluación de

modelos

Calidad modelos regresión

cuadrado de los residuos (los errores que

cometemos al predecir el target de cada

ejemplo) con la dispersión de los targets

R^2=1 sólo si todos los residuos son 0. Por

tanto, **cuanto más cercano a 1, mejor es** 

nuestro modelo.

Positive (P) True positive (TP) False negative (FN

Negative (N) False positive (FP) True negative (TN)

• **Precision**: Precision=1 ~ todos los ejemplos estimados positivos lo eran.

F1-score: Media entre precision y recall
Accuracy: porcentaje aciertos

Curva ROC y AUROC

mayor es AUROC, mejor es el modelo.

model = LogisticRegression()

# Cálculo de las probabilidades

# Cálculo de la curva ROC y AUC

# Visualización de la curva ROC

plt.xlabel('False Positive Rate')

plt.ylabel('True Positive Rate')

y\_score = model.predict\_proba(testX)[:, 1]

fpr, tpr, thresholds = roc\_curve(testY, y\_score)

model.fit(trainX, trainY)

roc\_auc = auc(fpr, tpr)

plt.plot(fpr, tpr)

# Se hace entrenamiento con regresión logística:

from sklearn.metrics import roc\_curve, auc, confusion\_matrix

# Entrenamiento del modelo de regresión logística

• Recall: Recall=1 ~todos los ejemplos positivos han sido estimados correctamente.

Positive (P)

True positive (TP),
hit

(FN),
type II error, miss,
underestimation

Sensitivity (SEN),
probability of detection, hit rate,
power  $= \frac{TP}{P} = 1 - FNR$ False fregative rate (FNR),
miss rate  $= \frac{FN}{P} = 1 - TPR$ 

Negative (N) type I error, false alarm, (TN), probability of false alarm, fall-out specificity (SPC), selectivity

overestimation correct rejection  $= \frac{FP}{N} = 1 - TNR$   $= \frac{TN}{N} = 1 - FPR$ 

 $Prevalence \\ = \frac{P}{P+N} \qquad Positive predictive value (PPV), \\ = \frac{TP}{PP} = 1 - FDR \qquad Positive predictive value (PPV), \\ = \frac{FN}{PN} = 1 - NPV \qquad Positive likelihood ratio (LR+) \\ = \frac{TPR}{FPR} \qquad (LR+) \\ = \frac{FNR}{TNR}$ 

 $\begin{array}{l} \text{Balanced} \\ \text{accuracy (BA)} \\ = \frac{\text{TPR} + \text{TNR}}{2} \end{array} \\ = \frac{2 \text{ PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}} \\ = \frac{2 \text{ TP}}{2 \text{ TP} + \text{FP} + \text{FN}} \\ \end{array} \\ \begin{array}{l} \text{Fowlkes-Mallows} \\ \text{index (FM)} \\ = \sqrt{\text{PPV} \times \text{TPR}} \\ \end{array} \\ \begin{array}{l} \text{Matthews correlation coefficient} \\ \text{(MCC)} \\ = \sqrt{\text{TPR} \times \text{TNR} \times \text{PPV} \times \text{NPV}} \\ - \sqrt{\text{FNR} \times \text{FPR} \times \text{FOR} \times \text{FDR}} \\ \end{array} \\ \begin{array}{l} \text{Jaccard index} \\ = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}} \\ \end{array} \\ \end{array}$ 

Accuracy (ACC) False discovery rate (FDR)  $= \frac{FP}{PP} = 1 - PPV$  False discovery rate (FDR)  $= \frac{FP}{PP} = 1 - PPV$  False discovery rate (FDR)  $= \frac{FP}{PP} = 1 - PPV$  Markedness (MK), deltaP ( $\Delta P$ )  $= \frac{TN}{P} = 1 - FOP$ 

Dado un modelo que estima la probabilidad de p(y=1|x), tal que yhat=1 si p(y=1|x) > U;

la curva ROC es el trazado de pares (FPR Y TPR) cuando hacemos variar el umbral *U*.

La curva ROC es diferente a la mtriz de confusión porque se calcula para todo el rango

de posibles umbrales U, mientras que la matriz de confusión ya asume un *U* fijado.

Midiendo el área bajo la curva ROC (AUROC) podemos comparar modelos: cuanto

False positive (FP), True negative False positive rate (FPR), True negative rate (TNR),

NOTA: También se puede utilizar la propia función de pérdida, el "problema" es que NO

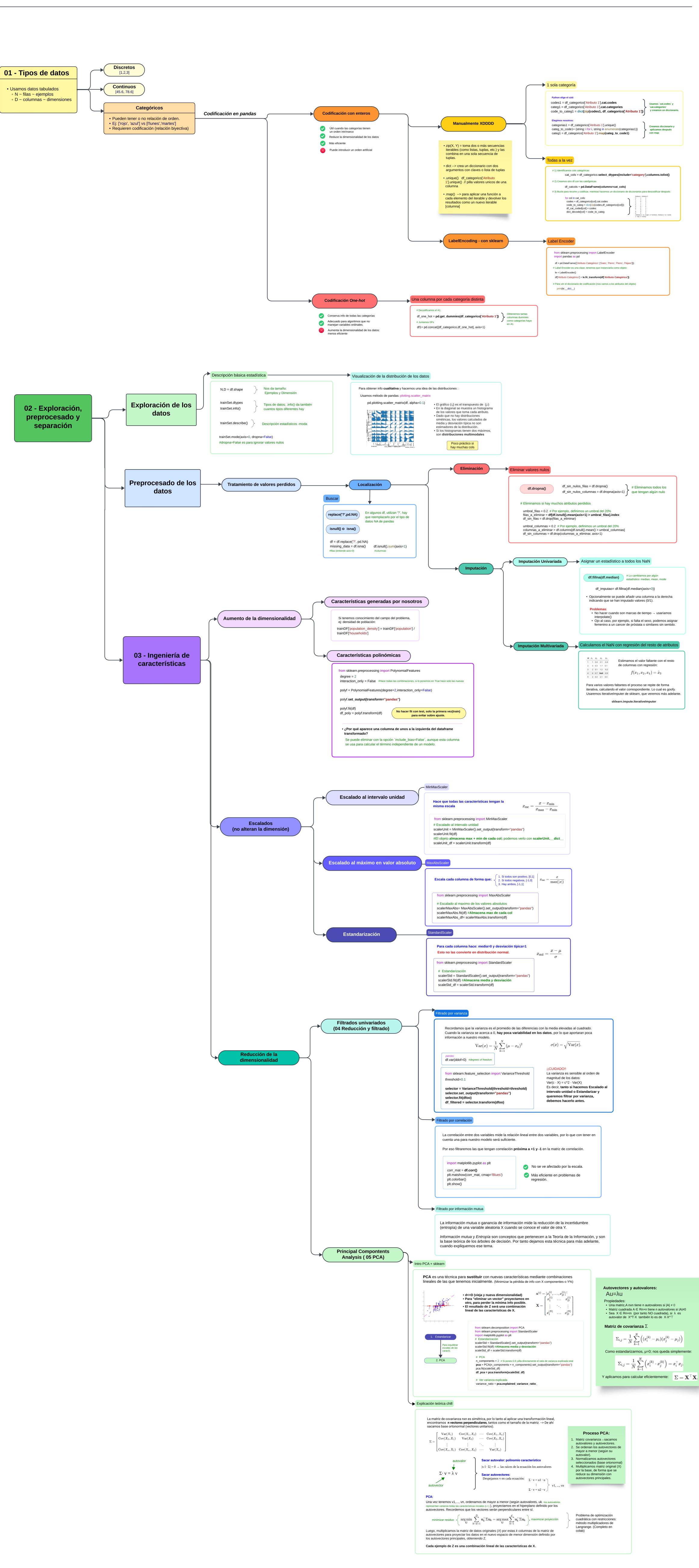
devuelve un valor acotado, dependerá de los datos. R^2 siempre tiene valor máximo 1.

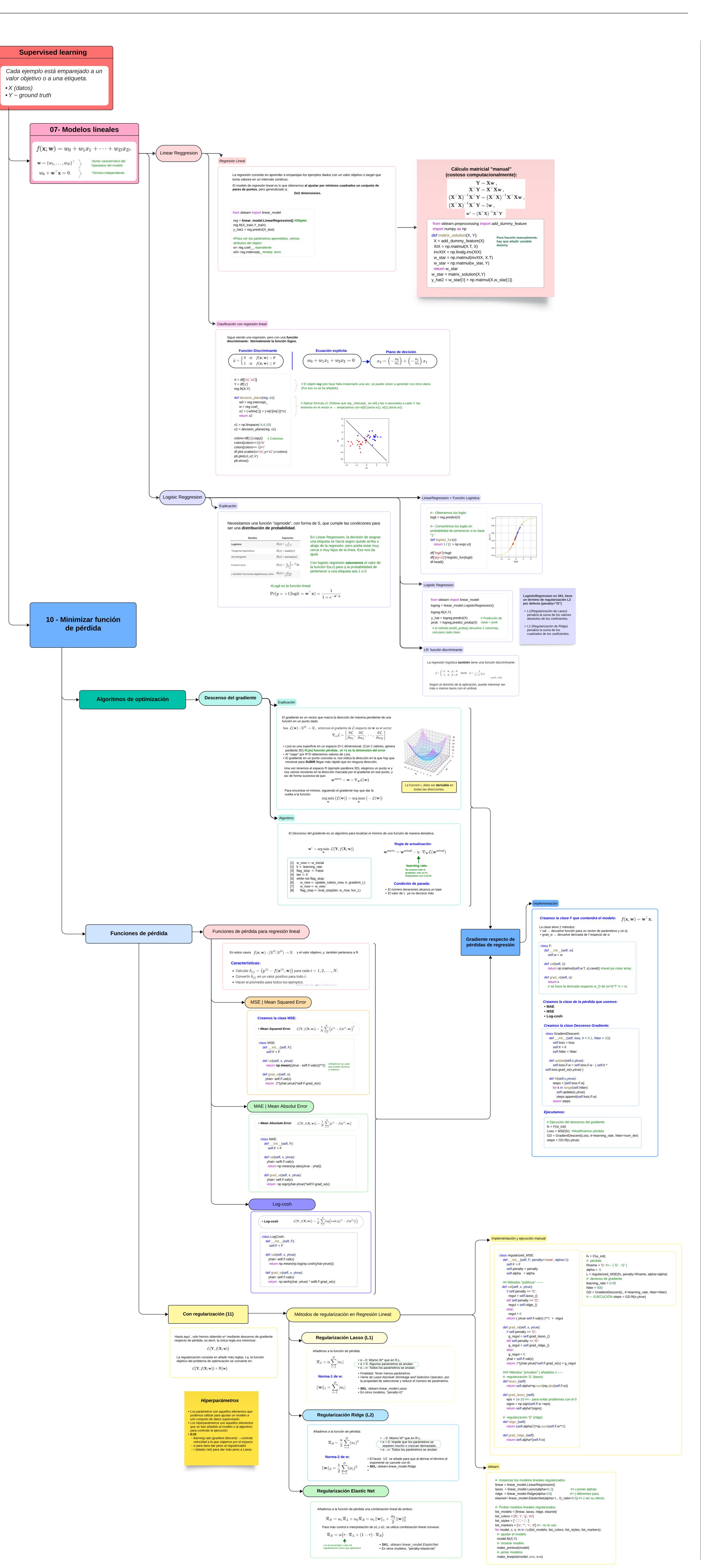
Calidad modelos de clasificación binaria:

La métrica típica en regresión es el valor:

sklearn.metrics.r2\_score

→ Matriz de confusión





## 01-fichero-entrenamiento

May 15, 2024

# 1 Aprendizaje Automático: Práctica sobre Diagnóstico de Enfermedades

## 1.1 Exploración y preprocesado de datos:

#### 1.1.1 Exploración de datos

1. Escribir un informe con una descripción de cada atributo de la columna que sea informativa y útil.

```
[]: import numpy as np
  import pandas as pd
  import matplotlib.pyplot as plt
  import seaborn as sns
  import pickle

from sklearn.preprocessing import LabelEncoder
  from sklearn.preprocessing import PolynomialFeatures, StandardScaler
  from sklearn.decomposition import PCA
  from sklearn.feature_selection import VarianceThreshold
  from sklearn.model_selection import train_test_split
  from sklearn.linear_model import LogisticRegression
  from sklearn.preprocessing import LabelEncoder
  from sklearn.metrics import confusion_matrix
  from sklearn.metrics import roc_curve, auc
```

```
Index(['Unnamed: 0', 'age', 'bp', 'sg', 'al', 'su', 'rbc', 'pc', 'pcc', 'ba',
'bgr', 'bu', 'sc', 'sod', 'pot', 'hemo', 'pcv', 'wbcc', 'rbcc', 'htn', 'dm',
```

```
'cad', 'appet', 'pe', 'ane'], dtype='object')
[]: # Comprobar el tipo de datos de los Data Frames y "qué pinta tienen"
     print("\n Muestra del X_train:\n", X_train.head(), "\n Muestra del Y_train:\n", | 
      →Y_train.head() )
     print("\nInfo X_train: \n")
     X_train.info()
     print("\nInfo Y_train:\n")
     Y_train.info()
     Muestra del Y_train:
        Unnamed: 0
                     age
                                                   cad appet pe ane
                            bр
                                        al ...
                                                dm
                                   sg
    0
              195 70.0 90.0
                                1.02 2.0 ...
                                              yes
                                                   yes
                                                        poor
                                                              no
                                                                   no
    1
              358 47.0 60.0
                                1.02 0.0
                                                         good
    2
              173 17.0 70.0
                               1.015
                                      1.0
                                               no
                                                         good
                                                              no
    3
              306 52.0 80.0
                                1.02 0.0
                                                        good no
                                               no
                                                    no
                                                                   no
              377
                   64.0 70.0
                                1.02 0.0 ...
                                                        good no
                                               no
                                                                   nο
    [5 rows x 25 columns]
     Muestra del Y_train:
        Unnamed: 0
                     class
              195
    0
                      ckd
    1
              358 notckd
              173
                      ckd
    3
              306 notckd
              377
                  notckd
    Info X_train:
    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 320 entries, 0 to 319
    Data columns (total 25 columns):
                     Non-Null Count Dtype
         Column
                     _____
     0
         Unnamed: 0 320 non-null
                                     int64
                     320 non-null
     1
                                     object
         age
     2
                     320 non-null
         bp
                                     object
     3
                     320 non-null
                                     object
         sg
     4
         al
                     320 non-null
                                     object
     5
                     320 non-null
         su
                                     object
     6
                     320 non-null
         rbc
                                     object
```

320 non-null

320 non-null

320 non-null

320 non-null

320 non-null

7

8

9

10

рс

рсс

ba

bgr

bu

object

object

object

object

object

```
12
                    320 non-null
                                    object
         sc
     13 sod
                     320 non-null
                                    object
     14
                    320 non-null
                                    object
        pot
                    320 non-null
                                    object
     15 hemo
     16 pcv
                    320 non-null
                                    object
     17
                     320 non-null
                                    object
        wbcc
     18 rbcc
                    320 non-null
                                    object
     19 htn
                    320 non-null
                                    object
     20 dm
                    320 non-null
                                    object
     21 cad
                    320 non-null
                                    object
     22 appet
                    320 non-null
                                    object
                    320 non-null
                                    object
     23
        ре
     24 ane
                    320 non-null
                                    object
    dtypes: int64(1), object(24)
    memory usage: 62.6+ KB
    Info Y_train:
    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 320 entries, 0 to 319
    Data columns (total 2 columns):
         Column
                    Non-Null Count Dtype
                    _____
         Unnamed: 0 320 non-null
                                    int64
         class
                    320 non-null
                                    object
    dtypes: int64(1), object(1)
    memory usage: 5.1+ KB
[]: # Reemplazar "?" con NaN en todas las columnas
    X_train.replace('?', np.nan, inplace=True)
    Y_train.replace('?', np.nan, inplace=True)
    X_test.replace('?', np.nan, inplace=True)
[]: # Columnas categóricas
     columnas_categoricasX = ['rbc', 'pc', 'pcc', 'ba', 'htn', 'dm', 'cad', 'appet', __
      # Convertir columnas categóricas a tipo category en X_train y Y_train
    X_train[columnas_categoricasX] = X_train[columnas_categoricasX].
     ⇔astype('category')
    X_test[columnas_categoricasX] = X_test[columnas_categoricasX].astype('category')
    Y_train['class'] = Y_train['class'].astype('category')
     # Columnas numéricas
```

```
columnas_numericas = ['age', 'bp', 'sg', 'al', 'su', 'bgr', 'bu', 'sc', 'sod', □

→'pot', 'hemo', 'pcv', 'wbcc', 'rbcc']

# Convertir columnas numéricas a tipo float en X_train

X_train[columnas_numericas] = X_train[columnas_numericas].astype(float)

X_test[columnas_numericas] = X_test[columnas_numericas].astype(float)
```

[]: # Comprobar el tipo de datos de los data frames
X\_train.info()
Y\_train.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 320 entries, 0 to 319
Data columns (total 25 columns):

#	Column	Non-Null Count	Dtype
0	Unnamed: 0	320 non-null	int64
1	age	313 non-null	float64
2	bp	311 non-null	float64
3	sg	289 non-null	float64
4	al	289 non-null	float64
5	su	287 non-null	float64
6	rbc	200 non-null	category
7	pc	274 non-null	category
8	pcc	317 non-null	category
9	ba	317 non-null	category
10	bgr	284 non-null	float64
11	bu	304 non-null	float64
12	sc	306 non-null	float64
13	sod	250 non-null	float64
14	pot	249 non-null	float64
15	hemo	281 non-null	float64
16	pcv	263 non-null	float64
17	wbcc	237 non-null	float64
18	rbcc	214 non-null	float64
19	htn	318 non-null	category
20	dm	318 non-null	category
21	cad	318 non-null	category
22	appet	319 non-null	category
23	pe	319 non-null	category
24	ane	319 non-null	category
dtypes: category(10), float64(14), int64(1)			
memory usage: 42.0 KB			
<pre><class 'pandas.core.frame.dataframe'=""></class></pre>			
RangeIndex: 320 entries, 0 to 319			
Data columns (total 2 columns):			
#	Column	Non-Null Count	Dtype

0 Unnamed: 0 320 non-null int64

1 class 320 non-null category

dtypes: category(1), int64(1)

memory usage: 3.1 KB

## []: print(X\_train)

```
Unnamed: 0
                    age
                                                      cad appet
                             bp
                                   sg
                                        al
                                                 dm
                                                                  ре
                                                                       ane
0
             195 70.00
                         90.00 1.02 2.00
                                                      yes
                                                yes
                                                            poor
                                                                  no
                                                                        no
1
             358 47.00
                         60.00 1.02 0.00
                                                 no
                                                       no
                                                            good
                                                                  no
                                                                        no
2
             173 17.00
                         70.00 1.01 1.00
                                                 no
                                                       no
                                                            good
                                                                  no
                                                                        no
3
             306 52.00
                         80.00 1.02 0.00
                                                 no
                                                       no
                                                            good
                                                                  no
                                                                        no
4
             377 64.00
                         70.00 1.02 0.00
                                                 no
                                                       no
                                                            good
                                                                  no
                                                                        no
315
             158 71.00
                         60.00 1.02 3.00
                                                                  no
                                                yes
                                                      yes
                                                            good
                                                                        no
316
               8 52.00 100.00 1.01 3.00
                                                yes
                                                       no
                                                            good
                                                                  no
                                                                       yes
317
             392 57.00
                         80.00 1.02 0.00
                                                 no
                                                       no
                                                            good
                                                                  no
                                                                        no
             317 58.00
                         70.00 1.02 0.00
318
                                                 no
                                                       no
                                                            good
                                                                  no
                                                                        no
319
             300 45.00
                         60.00 1.02 0.00
                                                 no
                                                            good
                                                       no
                                                                  no
                                                                        no
```

[320 rows x 25 columns]

#### #####Variables de X tipo float64:

- 1. age: Representa la edad de los pacientes.
- 2. **bp**: Indica la presión arterial de los pacientes.
- 3. sg: Indica la densidad específica de la orina.
- 4. al: Representa el nivel de albuminuria en la orina.
- 5. su: Significa la cantidad de glucosa en suero.
- 6. bgr: Representa el nivel de glucosa en sangre en ayunas.
- 7. bu: Indica la concentración de urea en sangre.
- 8. sc: Representa la concentración de creatinina sérica en sangre.
- 9. sod: Indica la concentración de sodio en sangre.
- 10. pot: Representa la concentración de potasio en sangre.
- 11. **hemo**: Indica el nivel de hemoglobina en la sangre.
- 12. pcv: Representa el volumen de células empaquetadas en la sangre.
- 13. wbcc: Indica el recuento de glóbulos blancos en la sangre.
- 14. **rbcc**: Representa el recuento de glóbulos rojos en la sangre.

Descartamos hacer una representación gráfica de las variables para obtener información cualitativa y hacernos una idea de las distribuciones, ya que hay bastantes columnas. En su lugar mostraremos una tabla con estadísticos básicos para las variables numéricas:

```
[]: desc = X_train[columnas_numericas].describe()
mode = X_train[columnas_numericas].mode().iloc[0]
mode_df = mode.to_frame().T.rename(index={0: 'mode'})
descripcion = pd.concat([desc, mode_df])
print(descripcion)
```

```
su ...
                                                                       wbcc
         age
                 bp
                        sg
                               al
                                                pot
                                                      hemo
                                                              pcv
rbcc
count 313.00 311.00 289.00 289.00 287.00 ... 249.00 281.00 263.00
                                                                     237.00
214.00
                              1.04
                                               4.75 12.54 38.84
mean
       51.45 76.59
                      1.02
                                     0.47 ...
                                                                   8424.47
4.72
std
       17.51
             13.27
                      0.01
                             1.35
                                     1.11 ...
                                               3.55
                                                      2.93
                                                             8.90
                                                                   2956.03
1.03
min
        2.00 50.00
                      1.00
                             0.00
                                     0.00 ...
                                               2.50
                                                      3.10
                                                             9.00
                                                                   2200.00
2.10
25%
       41.00 70.00
                      1.01
                             0.00
                                     0.00 ...
                                               3.90 10.40
                                                            33.00 6500.00
4.00
50%
       55.00 80.00
                      1.02
                             0.00
                                     0.00 ...
                                               4.50
                                                    12.60
                                                                   8100.00
                                                            40.00
4.75
75%
       65.00 80.00
                      1.02
                                               4.90
                                                     15.00
                                                            45.00
                                                                   9800.00
                             2.00
                                     0.00 ...
5.50
max
       90.00 180.00
                      1.02
                             5.00
                                     5.00 ...
                                              47.00 17.80
                                                            54.00 26400.00
8.00
mode
       65.00 80.00
                      1.02
                             0.00
                                     0.00 ...
                                               5.00 15.00 41.00 9800.00
4.50
```

#### [9 rows x 14 columns]

```
age
                 bp
                                al
                                       su
                        sg
count 313.00 311.00 289.00 289.00 287.00
mean
       51.45 76.59
                      1.02
                             1.04
                                     0.47
       17.51 13.27
                             1.35
                                     1.11
std
                      0.01
                             0.00
min
        2.00 50.00
                      1.00
                                     0.00
25%
       41.00 70.00
                      1.01
                             0.00
                                     0.00
50%
       55.00 80.00
                      1.02
                             0.00
                                     0.00
75%
       65.00 80.00
                      1.02
                             2.00
                                     0.00
                                     5.00
max
       90.00 180.00
                      1.02
                             5.00
       65.00 80.00
                      1.02
                             0.00
                                     0.00
mode
```

```
bgr
                                sod
                  bu
                         SC
                                       pot
count 284.00 304.00 306.00 250.00 249.00
      146.32 57.77
                       3.20 137.78
                                      4.75
mean
       76.91
              51.12
                       6.24
                             11.17
                                      3.55
std
min
       22.00
               1.50
                       0.40
                               4.50
                                      2.50
25%
       99.00
              27.00
                       0.90 135.00
                                      3.90
50%
      122.00
              44.00
                       1.30 139.00
                                      4.50
75%
      159.25
              66.00
                       2.88 142.00
                                      4.90
      490.00 391.00
                      76.00 163.00
                                     47.00
max
mode
       99.00 25.00
                       1.20 135.00
                                      5.00
                pcv
        hemo
                         wbcc
                                 rbcc
count 281.00 263.00
                       237.00 214.00
       12.54
              38.84
                      8424.47
                                 4.72
mean
        2.93
               8.90
                      2956.03
                                 1.03
std
        3.10
               9.00
                      2200.00
                                 2.10
min
       10.40
              33.00
                      6500.00
                                 4.00
25%
       12.60
              40.00
                      8100.00
                                 4.75
50%
75%
       15.00
              45.00
                      9800.00
                                 5.50
max
       17.80
              54.00 26400.00
                                 8.00
mode
       15.00
              41.00 9800.00
                                 4.50
```

#### #####Variables de X tipo category:

```
[]: print("\nDominio de las categorías en X_train:")
for columna in columnas_categoricasX:
    print(f"{columna}: {X_train[columna].cat.categories.tolist()}")
```

```
Dominio de las categorías en X_train:
rbc: ['abnormal', 'normal']
pc: ['abnormal', 'normal']
pcc: ['notpresent', 'present']
ba: ['notpresent', 'present']
htn: ['no', 'yes']
dm: ['no', 'yes']
cad: ['no', 'yes']
appet: ['good', 'poor']
pe: ['no', 'yes']
ane: ['no', 'yes']
```

- 1. **rbc**: Este atributo indica el recuento de glóbulos rojos en la sangre. Cuenta con 2 categorías: abnormal y normal.
- 2. **pc**: Indica la apariencia de las células presentes en la orina. Cuenta con 2 categorías: abnormal y normal.
- 3. pcc: Representa la presencia de cilindros patológicos en la orina. Cuenta con 2 categorías:

present y notpresent.

- 4. ba: Indica la presencia de gránulos hialinos en la orina. Cuenta con 2 categorías: present y notpresent.
- 5. htn: Este atributo indica la presencia de hipertensión. Cuenta con 2 categorías: yes y no.
- 6. dm: Representa la presencia de diabetes mellitus. Cuenta con 2 categorías: yes y no.
- 7. cad: Indica la presencia de enfermedad arterial coronaria. Cuenta con 2 categorías: yes y no.
- 8. appet: Representa el apetito del paciente. Cuenta con 2 categorías: good y poor.
- 9. pe: Indica la presencia de edema periférico. Cuenta con 2 categorías: yes y no.
- 10. ane: Representa la presencia de anemia. Cuenta con 2 categorías: yes y no.

Una vez hayamos hecho el tratamiento de valores nulos, codificaremos las variables de tipo categórico con Label Encoding, ya que puede ser conflictivo si lo hacemos antes.

## ####Variables de Y tipo category:

```
[]: print("\nDominio de las categorías en Y_train:")
print(Y_train['class'].cat.categories.tolist())
```

```
Dominio de las categorías en Y_train: ['ckd', 'notckd']
```

1. class. Representa si el paciente tiene o no la afección. Hay dos categorías: nockd paciente sin afección y ckd: paciente con afección.

```
Unnamed: 0
                     class
191
              126
                       ckd
81
              296
                   notckd
217
              333
                   notckd
174
               42
                       ckd
309
              237
                       ckd
              239
                       ckd
201
51
               86
                       ckd
16
              319
                   notckd
268
              172
                       ckd
188
              139
                       ckd
```

[64 rows x 2 columns]

```
[]: x_train.shape x_val.shape
```

#### []: (64, 25)

#### 1.1.2 Tratamiento de valores faltantes

- 2. Procesar los datos para tratar los valores perdidos. Para ello se deben considerar las siguientes cuestiones:
- ¿Eliminar ese ejemplo porque tiene valores perdidos?
- ¿Eliminar ese atributo **porque** tiene valores perdidos?
- ¿Imputar los valores perdidos?

A la salida de este paso consideramos que los atributos ya son características

**Eliminar valores faltantes** Para eliminar columnas y filas con valores perdidos, hemos establecido un umbral del 20% para eliminarlas:

Problema: Simplemente con este criterio nos quedaba train con 23 y test con 24 columnas... Si eliminabamos la intersección, no eliminabamos ninguna por lo que si después imputabamos los datos con los métodos conocidos hasta ahora, podríamos alterar la distribución de los datos. Por ello, hemos optado por eliminar tanto las columnas q tienen 20% de valores faltantes en train y en test, quedándonos 23 cols.

```
# Eliminar filas con valores faltantes en x_train

# Eliminar filas con valores faltantes en x_train

umbral_filas = 0.2  # Umbral del 20% para eliminar filas

x_train.dropna(axis=0, thresh=int(x_train.shape[1] * (1 - umbral_filas)),u

inplace=True)

# Eliminar columnas con valores faltantes en x_train

umbral_columnas = 0.2  # Umbral del 20% para eliminar columnas

columnas_a_eliminar_train = x_train.columns[x_train.isnull().mean() >u

umbral_columnas]

x_train.drop(columnas_a_eliminar_train, axis=1, inplace=True)

# Eliminar filas en y_train que no tienen correspondencia en x_train

y_train = y_train.loc[x_train.index]

# Eliminar las mismas columnas en x_val

x_val.drop(columnas_a_eliminar_train, axis=1, inplace=True)
```

```
[]: print(x_train.shape, y_train.shape)
print(x_val.shape)
```

```
(219, 23) (219, 2) (64, 23)
```

```
[]: columnas_con_nulos = x_train.columns[x_train.isnull().any()].tolist()
     valores_nulos_por_columna = x_train.isnull().sum()
     print("Columnas con valores nulos:")
     print(columnas_con_nulos)
     print(valores_nulos_por_columna)
    Columnas con valores nulos:
    ['age', 'bp', 'sg', 'al', 'su', 'pc', 'pcc', 'ba', 'bgr', 'bu', 'sc', 'sod',
    'pot', 'hemo', 'pcv', 'wbcc', 'htn', 'dm', 'cad', 'appet', 'pe', 'ane']
    Unnamed: 0
    age
                    3
                    5
    bр
                    6
    sg
    al
                    5
                    5
    su
    рс
                   14
                    3
    рсс
                    3
    ba
                   15
    bgr
                    5
    bu
                    3
    sc
    sod
                   31
                   31
    pot
                   16
    hemo
                   21
    pcv
                   36
    wbcc
    htn
                    2
                    2
    dm
                    2
    cad
                    1
    appet
    ре
    ane
    dtype: int64
```

## Imputar valores faltantes

```
[]: # Visualizar histogramas de columnas numéricas
print("Representación columnas tipo float:")

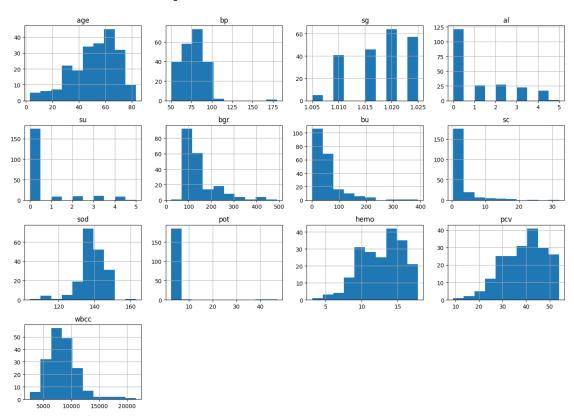
tipos_de_datos = x_train.dtypes

columnas_numericas = tipos_de_datos[tipos_de_datos != 'category'].index.tolist()
columnas_numericas = columnas_numericas[1:] #Quitamos el índice

x_train[columnas_numericas].hist(figsize=(14, 10))
plt.tight_layout()
```

#### plt.show()

Representación columnas tipo float:



#####Imputación univariada variables numéricas Observamos que no son distribuciones normales, por lo que escogemos utilizar la mediana para hacer imputación univariada.

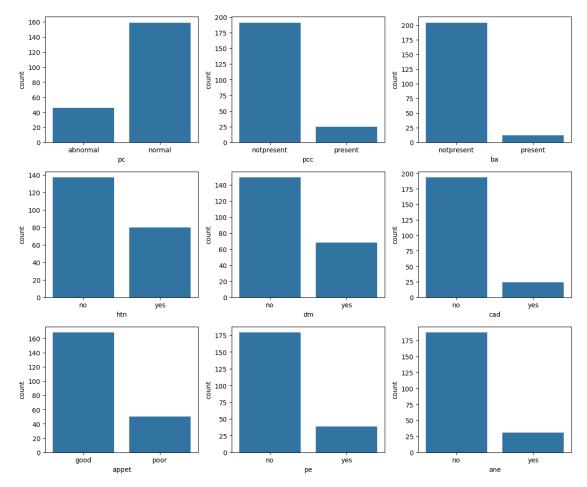
```
columnas_numericas = x_train.select_dtypes(include=['float64']).columns.tolist()
diccionario_medianas = {}

for column in columnas_numericas:
    median_train = x_train[column].median()
    x_train[column].fillna(median_train, inplace=True)
    x_val[column].fillna(median_train, inplace=True)
    diccionario_medianas[column] = median_train

print("Diccionario de medianas en X_train:", diccionario_medianas)
```

Diccionario de medianas en X\_train: {'age': 55.0, 'bp': 80.0, 'sg': 1.02, 'al': 0.0, 'su': 0.0, 'bgr': 119.5, 'bu': 41.5, 'sc': 1.2, 'sod': 139.0, 'pot': 4.5, 'hemo': 13.2, 'pcv': 40.0, 'wbcc': 8100.0}

#####Imputación univariada variables categóricas



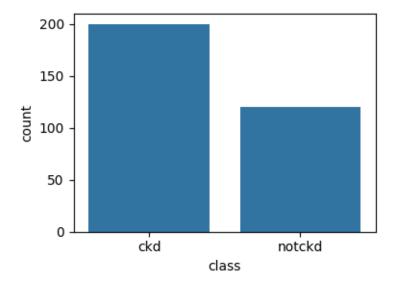
```
x_val[column].fillna(moda_train, inplace=True)

diccionario_modas[column] = moda_train

print("Diccionario de modas en X_train:", diccionario_modas)
```

```
Diccionario de modas en X_train: {'pc': 'normal', 'pcc': 'notpresent', 'ba': 'notpresent', 'htn': 'no', 'dm': 'no', 'cad': 'no', 'appet': 'good', 'pe': 'no', 'ane': 'no'}
```

```
[]: # Representar distribución variable objetivo
plt.figure(figsize=(4, 3))
sns.countplot(x='class', data=Y_train)
plt.tight_layout()
plt.show()
```



## 1.1.3 Codificación categorica

```
# Variable objetivo en Y
y_train['class'] = le.fit_transform(y_train['class'])
y_val['class'] = le.fit_transform(y_val['class'])
```

### 1.2 Ingeniería de características

3. Procesar las características para decidir si se elimina alguna de ellas.

## Funciones modulares para ingeniería de características

```
[]: def polyf_numerico (df_train, degree, interaction_only):
         # Crear características polinómicas de las numéricas
         cols_numericas = df_train.select_dtypes(include=float).columns.tolist()
         cols_categoricas = df_train.select_dtypes(include='category').columns.
      →tolist()
         polyf= PolynomialFeatures(degree=degree, interaction_only=interaction_only)
         polyf.set_output(transform="pandas")
         polyf.fit(df_train[cols_numericas])
         df_polyf_num = polyf.transform(df_train[cols_numericas])
         df_polyf_num = pd.concat([df_polyf_num, df_train[cols_categoricas]], axis=1)
         return df_polyf_num, polyf
     def polyf (df_train, degree, interaction_only):
         # Crear características polinómicas de todos
         polyf= PolynomialFeatures(degree=degree, interaction_only=interaction_only)
         polyf.set_output(transform="pandas")
         polyf.fit(df_train)
         df_polyf= polyf.transform(df_train)
         return df_polyf, polyf
     def filtrar_var_train(df_train, threshold):
         selector = VarianceThreshold(threshold=threshold)
         selector.set_output(transform="pandas")
         selector.fit(df_train)
         X_train_fil_var = selector.transform(df_train)
         return X_train_fil_var, selector
```

```
def filtrar_corr_train (df_train, umbral_correlacion):
    correlation_matrix = df_train.corr().abs()
    # Identificar las características altamente correlacionadas
   upper_tri = correlation_matrix.where(np.triu(np.ones(correlation_matrix.
 ⇒shape), k=1).astype(bool))
   high_correlation_features = [column for column in upper_tri.columns if_
 →any(upper_tri[column] > umbral_correlacion)]
    # Eliminar características altamente correlacionadas
   df train = df train.drop(columns=high correlation features)
   return df_train, high_correlation_features
def std_pca_train (df_train, n_components):
    # Estandarizar
    scalerStd = StandardScaler().set_output(transform="pandas")
    scalerStd.fit(df_train)
   df_train_scl = scalerStd.transform(df_train)
   # Aplicar PCA
   pca = PCA(n_components=n_components).set_output(transform="pandas")
   pca.fit(df_train_scl)
   df train pca = pca.transform(df train scl)
   return df_train_pca, scalerStd, pca
# Procesado 1: filtrado por varianza, estandarizacion y PCA (los datos seu
 ⇔modifican al invocar)
def procesado_train_1 (df_train, threshold, n_components):
   df_train, selector= filtrar_var_train(df_train, threshold)
   df_train, scaler, pca =std_pca_train(df_train, n_components)
   return df_train, selector, scaler, pca
def procesado_test_1(df_test, selector, scaler, pca):
   df_test_fil_var =selector.transform(df_test)
   df_test_scl= scaler.transform(df_test_fil_var)
   df_test_pca= pca.transform(df_test_scl)
   return df_test_pca
# Procesado 2: polynomialfeatures, estandarizacion, PCA
```

```
def procesado train 2 (df_train, pf_degree, pf_interaction_only,
      →pf_numeric_only, pca_n_components):
         if pf_numeric_only:
          df_train, polyf = polyf_numerico(df_train, pf_degree, pf_interaction_only)
        else:
          df train, polyf = polyf(df train, pf degree, pf interaction only)
        df train, scaler, pca = std pca train(df train, pca n components)
        return df_train, polyf, scaler, pca
    def procesado_test_2(df_test, pf_degree, pf_interaction_only, pf_numeric_only,_
      →polyf, scaler, pca):
        if pf_numeric_only:
            df_test, polyf = polyf_numerico(df_test, pf_degree, pf_interaction_only)
        else:
            df_test, polyf = polyf(df_test, pf_degree, pf_interaction_only)
        df_test = polyf.transform(df_test)
        df_test = scaler.transform(df_test)
        df_test = pca.transform(df_test)
        return df_test
[]: | # # Configuraciones distintas para procesado 1 (filtrado por varianza, ___
     ⇔estandarización, PCA):
     # Configuración modelo 1_1
    x_train_1_1, selector_1_1, scaler_1_1, pca1_1= procesado_train_1(x_train,_u
      →threshold=0.1, n_components=0.85)
    x_val_1_1= procesado_test_1(x_val, selector_1_1, scaler_1_1, pca1_1)
    # Configuración modelo 1_2
    x_train_1_2, selector_1_2, scaler_1_2, pca1_2= procesado_train_1(x_train,_
     x_val_1_2= procesado_test_1(x_val, selector_1_2, scaler_1_2, pca1_2)
    # Configuración modelo 1 3
    x_train_1_3, selector_1_3, scaler_1_3, pca1_3= procesado_train_1(x_train,_
      →threshold=0.1, n_components=2)
    x_val_1_3= procesado_test_1(x_val, selector_1_3, scaler_1_3, pca1_3)
    # # Configuraciones distintas para procesado 2 (polynomialfeatures, ___
     ⇔estandarización, PCA):
    x_train_2_1, polyf_2_1, scaler_2_1, pca2_1 = procesado_train_2(x_train,__
```

→pf\_degree=2,

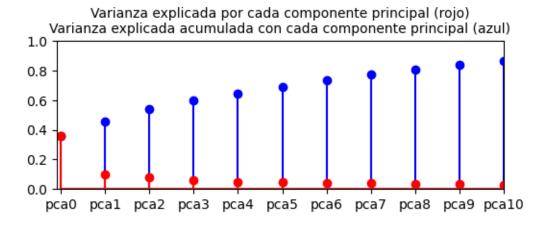
```
opf_interaction_only=True,

opf_numeric_only=True,

opca_n_components=0.85)
```

```
[]: # Función para generar el gráfico de explicación de varianza por número de
      \hookrightarrow componentes
     def plot_PCA(pca):
       plt.stem(pca.explained_variance_ratio_.cumsum(),'b')
       plt.stem(pca.explained_variance_ratio_,'r')
       titleStr = 'Varianza explicada por cada componente principal (rojo)'
       titleStr = titleStr+'\n'
       titleStr = titleStr+'Varianza explicada acumulada con cada componenteu
      ⇔principal (azul)'
       plt.title(titleStr, fontsize=10)
       ax = plt.gca()
       ax.axis([-0.1,1.1,0,1])
       ax.set_xticks([i for i in range(pca.n_components_)])
       ax.set_xticklabels(["pca"+str(i) for i in range(pca.n_components_)])
       fig = plt.gcf()
       fig.set_size_inches(6,2)
       plt.show()
```

```
[]: plot_PCA(pca1_1)
   pca1_1.n_components_
   sum(pca1_2.explained_variance_ratio_)
   sum(pca1_3.explained_variance_ratio_)
```



#### []: 0.4576948478897684

##Modelo lineal: Regresión Logística

logreg3.fit(x\_train\_1\_3, target)

```
[]: # # Modelo 1:
     logreg = LogisticRegression(penalty="11", solver='liblinear') # No se aprecian
     ⇔diferencias con l2
     target = y_train['class']
     # Ajustar el modelo con los datos de entrenamiento
     logreg.fit(x_train_1_1, target)
     # Realizar predicciones en el conjunto de validación
     y_hat = logreg.predict(x_val_1_1)
     # Calcular las probabilidades de las predicciones en el conjunto x_train_1_1
     prob = logreg.predict_proba(x_train_1_1)
     # print(y_hat, prob)
[]: # # Modelo 1_2:
     logreg2 = LogisticRegression(penalty="11", solver='liblinear')
     target = y_train['class']
     # Ajustar el modelo con los datos de entrenamiento
     logreg2.fit(x_train_1_2, target)
     # Realizar predicciones en el conjunto de validación
     y_hat2 = logreg2.predict(x_val_1_2)
     # Calcular las probabilidades de las predicciones en el conjunto x_train_1_2
     prob2 = logreg2.predict_proba(x_train_1_2)
     # print(y_hat, prob)
[]: # # Modelo 1_3:
     target = y_train['class']
     logreg3 = LogisticRegression(penalty="11", solver='liblinear')
     # Ajustar el modelo con los datos de entrenamiento
```

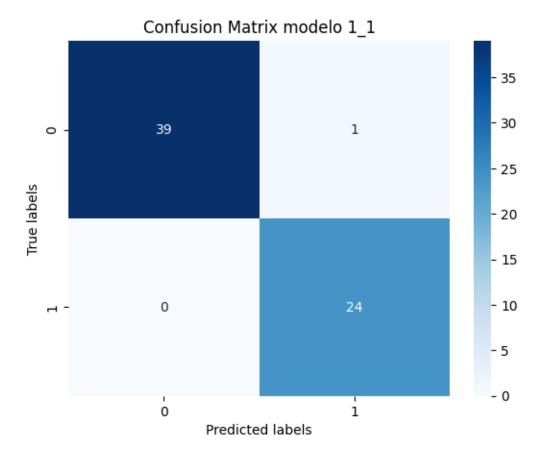
```
# Realizar predicciones en el conjunto de validación
y_hat3 = logreg3.predict(x_val_1_3)

# Calcular las probabilidades de las predicciones en el conjunto x_train_1_3
prob3 = logreg3.predict_proba(x_train_1_3)

# print(y_hat, prob)
```

```
[]: # Calcular la matriz de confusión
y_val_array = y_val['class'].values
cm1_1 = confusion_matrix(y_val_array, y_hat)

# Visualizar la matriz de confusión
sns.heatmap(cm1_1, annot=True, cmap='Blues', fmt='g')
plt.xlabel('Predicted labels')
plt.ylabel('True labels')
plt.title('Confusion Matrix modelo 1_1')
plt.show()
```



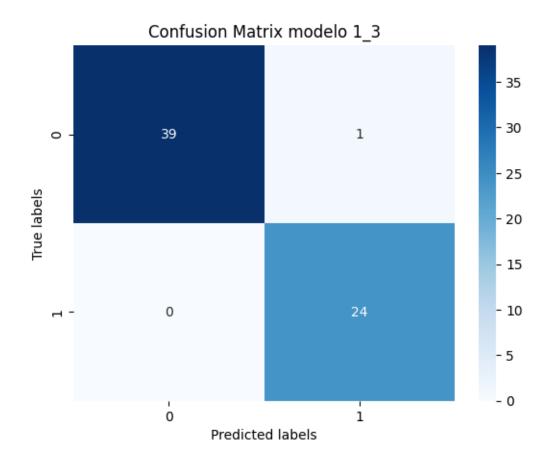
```
[]: y_val_array = y_val['class'].values
cm1_2 = confusion_matrix(y_val_array, y_hat2)

sns.heatmap(cm1_2, annot=True, cmap='Blues', fmt='g')
plt.xlabel('Predicted labels')
plt.ylabel('True labels')
plt.title('Confusion Matrix modelo 1_2')
plt.show()
```

# 

```
[]: cm1_3 = confusion_matrix(y_val_array, y_hat3)

sns.heatmap(cm1_3, annot=True, cmap='Blues', fmt='g')
plt.xlabel('Predicted labels')
plt.ylabel('True labels')
plt.title('Confusion Matrix modelo 1_3')
plt.show()
```



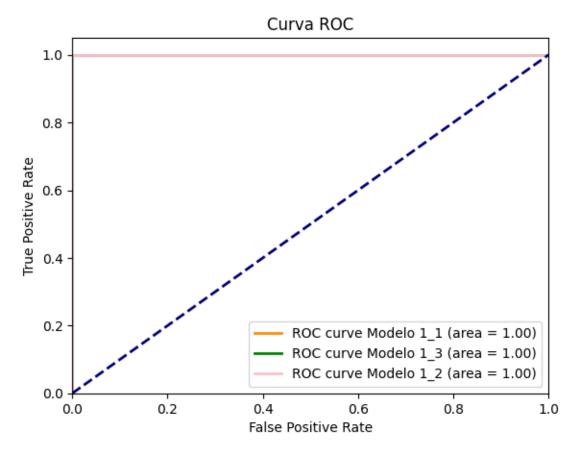
```
[]: # Calcular las probabilidades de las predicciones en el conjunto de validación prob_val1_1 = logreg.predict_proba(x_val_1_1) prob_val1_2 = logreg2.predict_proba(x_val_1_2) prob_val1_3 = logreg3.predict_proba(x_val_1_3)

# Calcular la curva ROC y el área bajo la curva (AUC) fpr1_1, tpr1_1, thresholds1_1 = roc_curve(y_val_array, prob_val1_1[:, 1]) roc_auc1_1 = auc(fpr1_1, tpr1_1)

fpr1_2, tpr1_2, thresholds1_2 = roc_curve(y_val_array, prob_val1_2[:, 1]) roc_auc1_2= auc(fpr1_2, tpr1_2)

fpr1_3, tpr1_3, thresholds1_3= roc_curve(y_val_array, prob_val1_3[:, 1]) roc_auc1_3 = auc(fpr1_3, tpr1_3)

# Visualizar la curva ROC para ambos modelos en la misma gráfica plt.figure()
```



```
[]: # Calcular los valores necesarios para el cálculo manual del F1-Score listay_hat=[y_hat, y_hat3]
```

```
def f1_score (y_hat, y_val_array):
    TP = np.sum((y_hat == 1) & (y_val_array == 1))
    FP = np.sum((y_hat == 1) & (y_val_array == 0))
    FN = np.sum((y_hat == 0) & (y_val_array == 1))

# Calcular la precisión y el recall
    precision = TP / (TP + FP)
    recall = TP / (TP + FN)

# Calcular el F1-Score
    f1_score = 2 * (precision * recall) / (precision + recall)
    return f1_score

f1_score1= f1_score(y_hat,y_val_array)
f1_score3= f1_score(y_hat2,y_val_array)
f1_score2=f1_score(y_hat2,y_val_array)
```

```
[]: print(f1_score1, f1_score2, f1_score3)
```

0.9795918367346939 0.960000000000000 0.9795918367346939

```
[]: # Guardar los objetos en un archivo usando pickle
with open('modelo_1_1_objetos.pkl', 'wb') as f:
    pickle.dump((selector_1_1, scaler_1_1, pca1_1,logreg), f)
```

## 02-fichero-test

May 15, 2024

```
[]: from google.colab import drive drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force\_remount=True).

```
[]: import pandas as pd
import numpy as np
import pickle

from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import PolynomialFeatures, StandardScaler
from sklearn.decomposition import PCA
from sklearn.feature_selection import VarianceThreshold
from sklearn.feature_selection import VarianceThreshold
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.linear_model import LogisticRegression
```

```
[]: # Modificar ruta según ubicación:
X_test = pd.read_csv('/content/retoML1_X_test.csv', delimiter=';')

# Ruta al archivo en tu Google Drive
ruta_al_archivo = '/content/modelo_1_1_objetos.pkl'

# Abrir el archivo y cargar los objetos
with open(ruta_al_archivo, 'rb') as archivo:
    objetos_cargados = pickle.load(archivo)
```

##Funciones preprocesado y procesado

```
[]: def preprocesado_atributos (df, columnas_categoricas, columnas_numericas):
    df.replace('?', np.nan, inplace=True)
# Convertir tipos
    df[columnas_categoricas]=df[columnas_categoricas].astype('category')
    df[columnas_numericas]=df[columnas_numericas].astype(float)
    return df
```

```
def tratar_valores_faltantes(df, columnas_eliminadas_train, __

→dict_median_numeric_vars, dict_modes_category_vars):
  # Eliminar columnas que hemos eliminado en test
  df.drop(columnas_eliminadas_train, axis=1, inplace=True)
  # Imputar valores faltantes
  # Para variables numéricas
  for var, median_value in dict_median_numeric_vars.items():
      df[var].fillna(median_value, inplace=True)
      df[var] = df[var].astype(float)
  # Para variables categóricas
  for var, mode_value in dict_modes_category_vars.items():
      df[var].fillna(mode_value, inplace=True)
      df[var] = df[var].astype('category')
  return df
def encoding(df):
  columnas_categoricas = df.select_dtypes(include=['category']).columns.tolist()
  le = LabelEncoder()
  for columna in columnas_categoricas:
      df[columna] = le.fit_transform(df[columna])
  return df
```

```
# Eliminar características altamente correlacionadas
   df_train = df_train.drop(columns=high_correlation_features)
   return df_train, high_correlation_features
def std_pca_train (df_train, n_components):
   # Estandarizar
   scalerStd = StandardScaler().set_output(transform="pandas")
   scalerStd.fit(df train)
   df_train_scl = scalerStd.transform(df_train)
   # Aplicar PCA
   pca = PCA(n_components=n_components).set_output(transform="pandas")
   pca.fit(df_train_scl)
   df_train_pca = pca.transform(df_train_scl)
   return df_train_pca, scalerStd, pca
# Procesado 1: filtrado por varianza, estandarización y PCA (los datos seu
 ⇔modifican al invocar)
def procesado_train_1 (df_train, threshold, n_components):
   df_train, selector= filtrar_var_train(df_train, threshold)
   df_train, scaler, pca =std_pca_train(df_train, n_components)
   return df_train, selector, scaler, pca
def procesado_test_1(df_test, selector, scaler, pca):
   df_test_fil_var =selector.transform(df_test)
   df_test_scl= scaler.transform(df_test_fil_var)
   df_test_pca= pca.transform(df_test_scl)
   return df_test_pca
```

#### 0.1 Datos extraídos

```
columnas_eliminadas_train = pd.Index(['rbc', 'rbcc'])
    dict_median_numeric_vars = {'age': 55.0, 'bp': 80.0, 'sg': 1.02, 'al': 0.0,
     y'su': 0.0, 'bgr': 119.5, 'bu': 41.5, 'sc': 1.2, 'sod': 139.0, 'pot': 4.5, □
     ⇔'hemo': 13.2, 'pcv': 40.0, 'wbcc': 8100.0}
    dict modes category vars = {'pc': 'normal', 'pcc': 'notpresent', 'ba':
     o'notpresent', 'htn': 'no', 'dm': 'no', 'cad': 'no', 'appet': 'good', 'pe':⊔

    'no', 'ane': 'no'}

   ##Preprocesado
[]: X_test= preprocesado_atributos(X_test, columnas_categoricas, columnas_numericas)
    X test=
     dict_modes_category_vars=dict_modes_category_vars)
    X test= encoding(X test)
   0.2 Ingenería características
[]: # Crear y asignar nombres a cada objeto
    selector = objetos_cargados[0]
    scaler = objetos_cargados[1]
    pca = objetos_cargados[2]
    log_reg = objetos_cargados[3]
    # Usar los objetos como sea necesario
[]: # Procesar
    X_test_procesado= procesado_test_1(X_test, selector, scaler, pca)
   ##Predicción
[]: # Predecir
    predicciones = log_reg.predict(X_test_procesado)
    print(predicciones)
    [0\ 1\ 0\ 0\ 1\ 0\ 0\ 1\ 0\ 1\ 0\ 1\ 0\ 0\ 1\ 0\ 0\ 1\ 0\ 0\ 1\ 0\ 0\ 0\ 0\ 1\ 1\ 1\ 1\ 1\ 1\ 0
    1 1 0 1 0 0]
   #Postprocesado y exportación
[]: mapeo = {0: 'ckd', 1: 'notckd'}
    predicciones= np.array([mapeo[pred] for pred in predicciones])
[]: # Crear un data frame con la lista de predicciones
    df = pd.DataFrame({'class': predicciones})
```

```
# Nombre del archivo CSV de salida
nombre_archivo = 'retoML1_Y_test.csv'

# Exportar el data fFrame a un archivo CSV
df.to_csv(nombre_archivo, index=False)
```