

Introduction to Feature Improvement & Data Cleaning

Applications of Machine Learning

Vanessa Gómez Verdejo vanessag@ing.uc3m.es

Scaling

- Linear Regression | Insensitive
- KNN | Sensitive
- Regularization LR | Sensitive
- Trees | Insensitive

In this session, we are going to take the results of our understanding of the data and use them to improve the data. **Feature improvement** is about recognizing areas of issue and improvement in our data and figuring out which cleaning methods will be the most effective. Once we detect problems in our data, instead of immediately dropping rows/columns, we should think about the best ways of fixing these problems. More often than not, our machine learning performance will thank us in the end.

In this notebook we will explore different techniques for preparing/improving our datasets before using them with a ML pipeline. In particular, we will review:

- the usefulness of the **data normalization** and other **data transformations**,
- how to deal with **missing values**,
- and how to clean our data from **outliers**.

```
In [194... from IPython.core.display import Image, display
%matplotlib inline
%config InlineBackend.figure_format = 'retina'
```

```
In [195... import matplotlib.pyplot as plt
from matplotlib import rc
import numpy as np
import pandas as pd
```

```
# Configuración de las figuras matplotlib
plt.rcParams['figure.figsize'] = [8, 6]
plt.rcParams.update({'font.size': 8})
```

1. Data normalization

To analyze the need of the data normalization, let's work with the boston dataset.

```
In [196... boston_df = pd.read_csv("http://www.tsc.uc3m.es/~sevisal/housing.csv",
# boston_df = pd.read_csv('housing.csv', skiprows = 1)
boston_df.rename(columns={"MEDV": "target"}, inplace=True)
```

```
In [197... col_name = boston_df.columns
col_name
```

```
Out[197... Index(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RA
D', 'TAX',
      'PTRATIO', 'B', 'LSTAT', 'target'],
      dtype='object')
```

```
In [198... data = boston_df[col_name[:-1]]
# data = boston_df.iloc[:, :-1]
data
```

```
Out[198...      CRIM  ZN  INDUS  CHAS  NOX  RM  AGE  DIS  RAD  TAX  PTR
0  0.00632  18.0   2.31    0  0.538  6.575  65.2  4.0900   1  296
1  0.02731   0.0   7.07    0  0.469  6.421  78.9  4.9671   2  242
2  0.02729   0.0   7.07    0  0.469  7.185  61.1  4.9671   2  242
3  0.03237   0.0   2.18    0  0.458  6.998  45.8  6.0622   3  222
4  0.06905   0.0   2.18    0  0.458  7.147  54.2  6.0622   3  222
...      ...   ...   ...   ...   ...   ...   ...   ...   ...   ...
501  0.06263   0.0  11.93    0  0.573  6.593  69.1  2.4786   1  273
502  0.04527   0.0  11.93    0  0.573  6.120  76.7  2.2875   1  273
503  0.06076   0.0  11.93    0  0.573  6.976  91.0  2.1675   1  273
504  0.10959   0.0  11.93    0  0.573  6.794  89.3  2.3889   1  273
505  0.04741   0.0  11.93    0  0.573  6.030  80.8  2.5050   1  273
```

506 rows x 13 columns

```
In [199... targets = boston_df[col_name[-1]]
```

```
# targets = boston_df.iloc[:, -1]
targets
```

Out [199...

	target
0	24.0
1	21.6
2	34.7
3	33.4
4	36.2
...	...
501	22.4
502	20.6
503	23.9
504	22.0
505	11.9

506 rows × 1 columns

dtype: float64

In [200...

```
boston_df.describe()
```

Out [200...

	CRIM	ZN	INDUS	CHAS	NOX	
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	3.593761	11.363636	11.136779	0.069170	0.554695	6.284186
std	8.596783	23.322453	6.860353	0.253994	0.115878	0.702919
min	0.006320	0.000000	0.460000	0.000000	0.385000	3.561541
25%	0.082045	0.000000	5.190000	0.000000	0.449000	5.885499
50%	0.256510	0.000000	9.690000	0.000000	0.538000	6.208519
75%	3.647423	12.500000	18.100000	0.000000	0.624000	6.623158
max	88.976200	100.000000	27.740000	1.000000	0.871000	8.780123

As we can see by looking at the mean values and standard deviations of each variable, the variables in this problem have very different ranges. Methods based on distances between observations, such as k-NN or kernel methods, are very

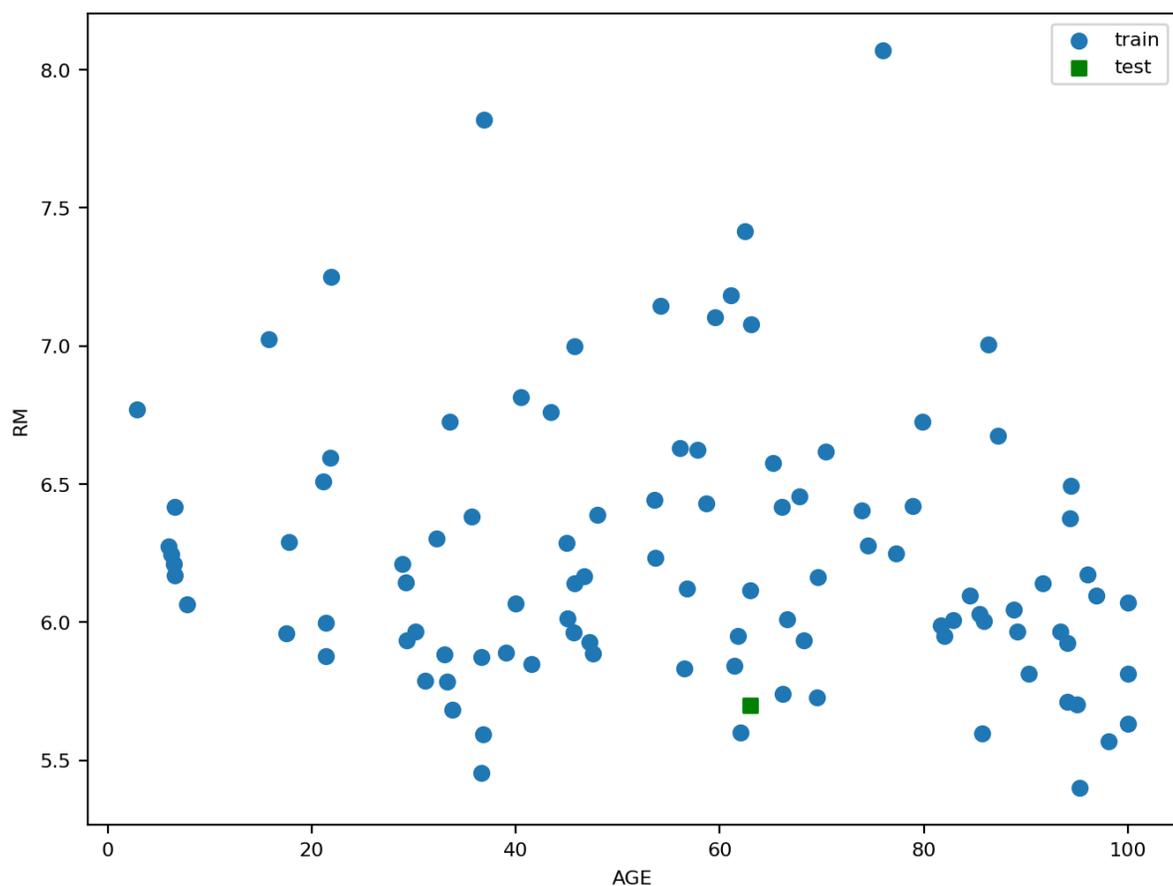
dependent on the dynamic ranges of the variables.

To analyze this, let's solve a mini-problem with only 100 observations and the AGE and RM variables using k-NN.

```
In [201... boston_df.columns
```

```
Out[201... Index(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX',
        'PTRATIO', 'B', 'LSTAT', 'target'],
        dtype='object')
```

```
In [202... mini_boston = boston_df.loc[:100,['AGE','RM','target']]
xt = np.array([63,5.7])
plt.figure()
plt.scatter(mini_boston['AGE'].values,
            mini_boston['RM'].values, label='train')
plt.scatter(xt[0],xt[1],marker='s',color='green', label='test')
plt.legend()
plt.xlabel('AGE')
_ = plt.ylabel('RM')
```



Which data are the 5 closest neighbors to the test observation?

```
In [203... from sklearn.metrics import pairwise_distances
```

```

def plot_example(x,y,xt,ax):
    ax.scatter(x[:,0],x[:,1],color='blue',marker='o')
    ax.scatter(xt[0], xt[1], marker='s', color='green')

def plot_radio(c,x,ax):
    for ii in range(len(x)):
        ax.plot([c[0],x[ii,0]],
                [c[1],x[ii,1]],
                linestyle=':',
                linewidth=2)

x = mini_boston.loc[:,['AGE','RM']].values
y = mini_boston['target'].values

distances = pairwise_distances(xt.reshape(1,-1), x,'euclidean')

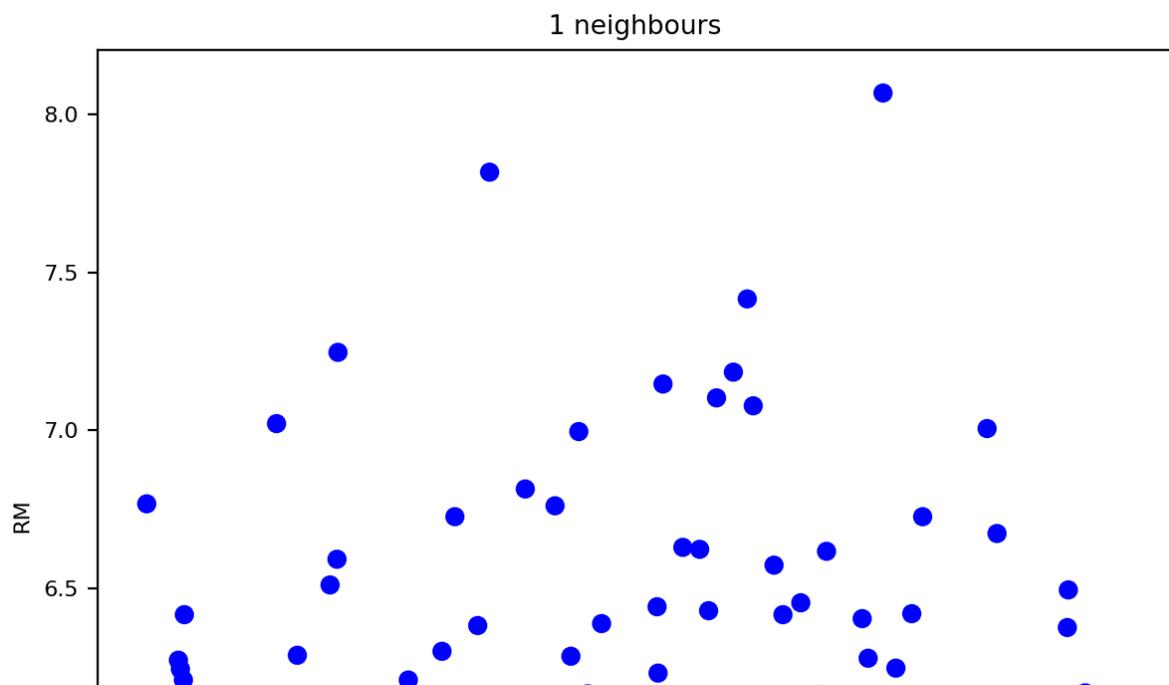
# Sort all neighbours for the test data
id_neigh = np.argsort(distances) # indexes over the training data, whi

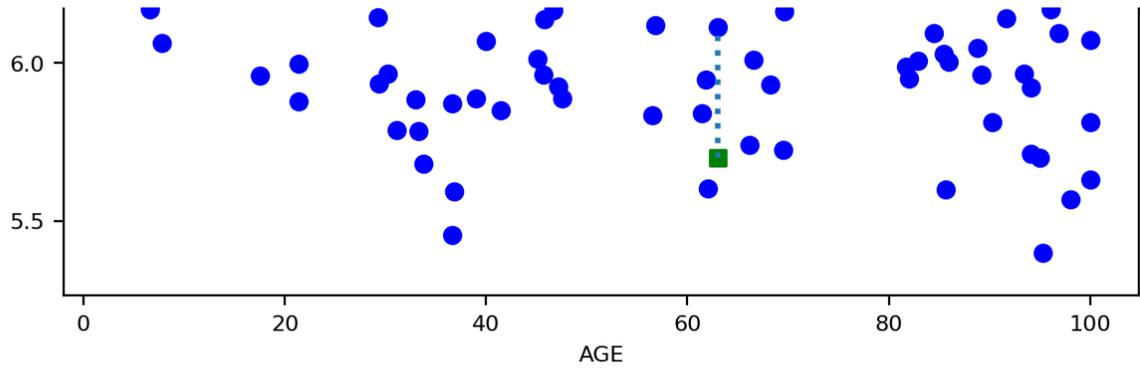
sorted_distances = 1e-6+np.sort(distances) # distances to neighbours h

# class each neighbour votes for
call_neigh = y[id_neigh]

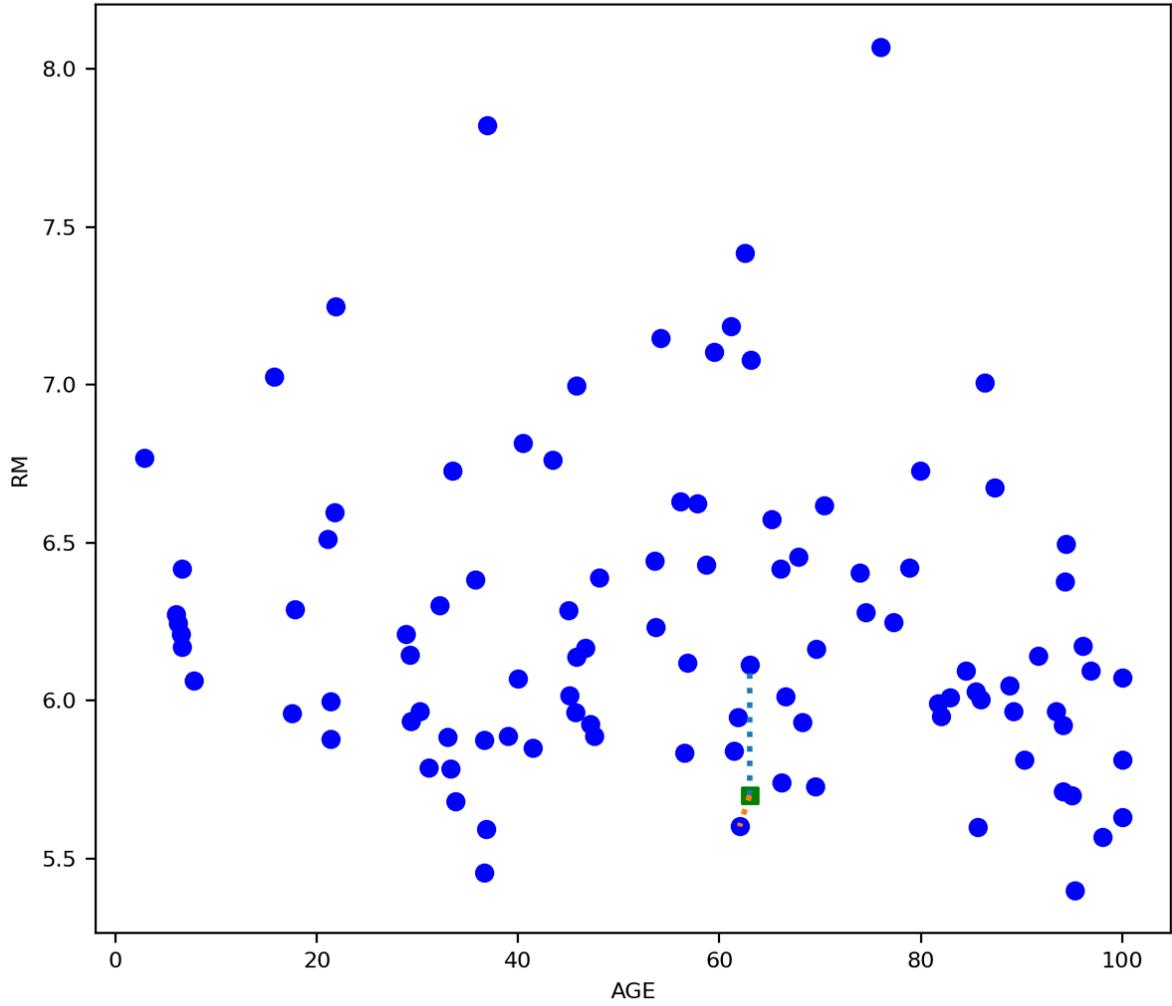
v_k = np.array([1,2,4,8])
ii=7
fx,ax = plt.subplots(len(v_k),1,figsize=(ii,len(v_k)*ii))
for ik,kk in enumerate(v_k):
    plot_example(x,y,xt,ax[ik])
    plot_radio(xt,np.array([x[cc,:] for cc in id_neigh[0,:kk]]),ax[ik])
    ax[ik].set_title('{0:d} neighbours'.format(kk))
    ax[ik].set_xlabel('AGE')
    _=ax[ik].set_ylabel('RM')

```

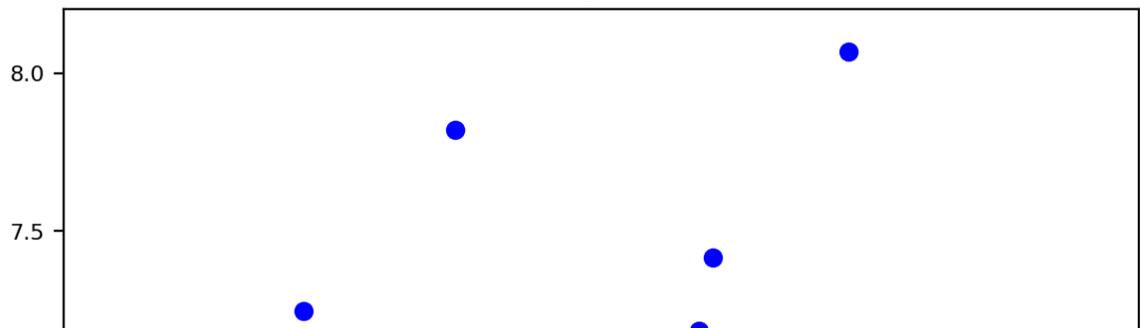


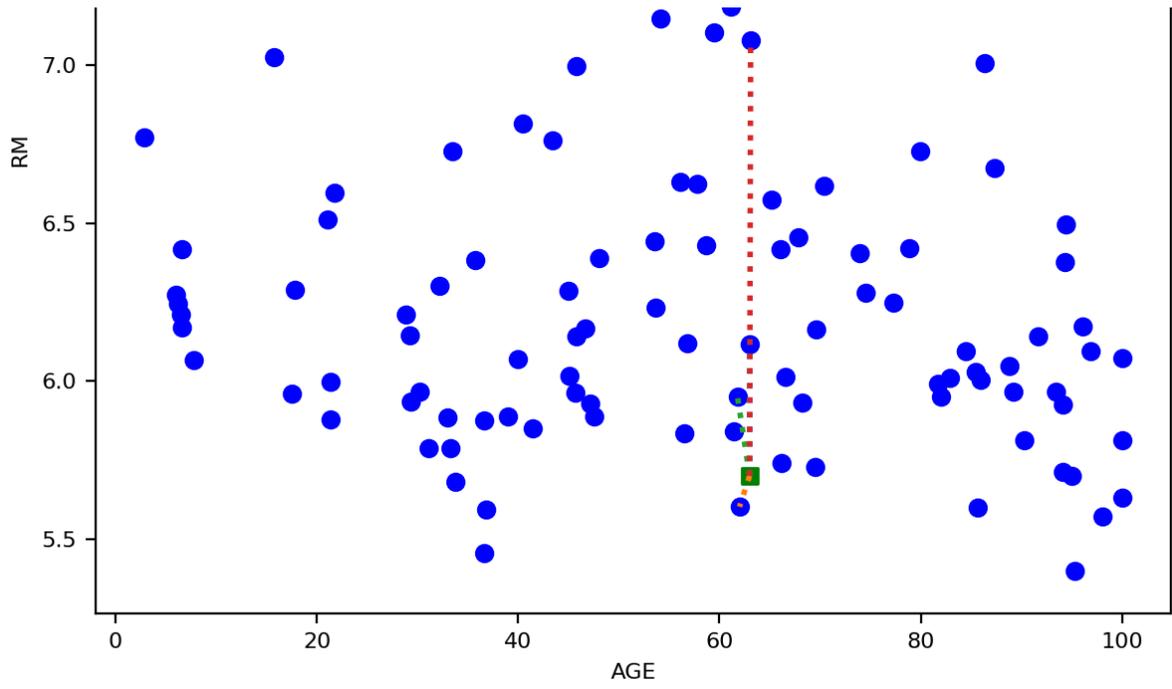


2 neighbours

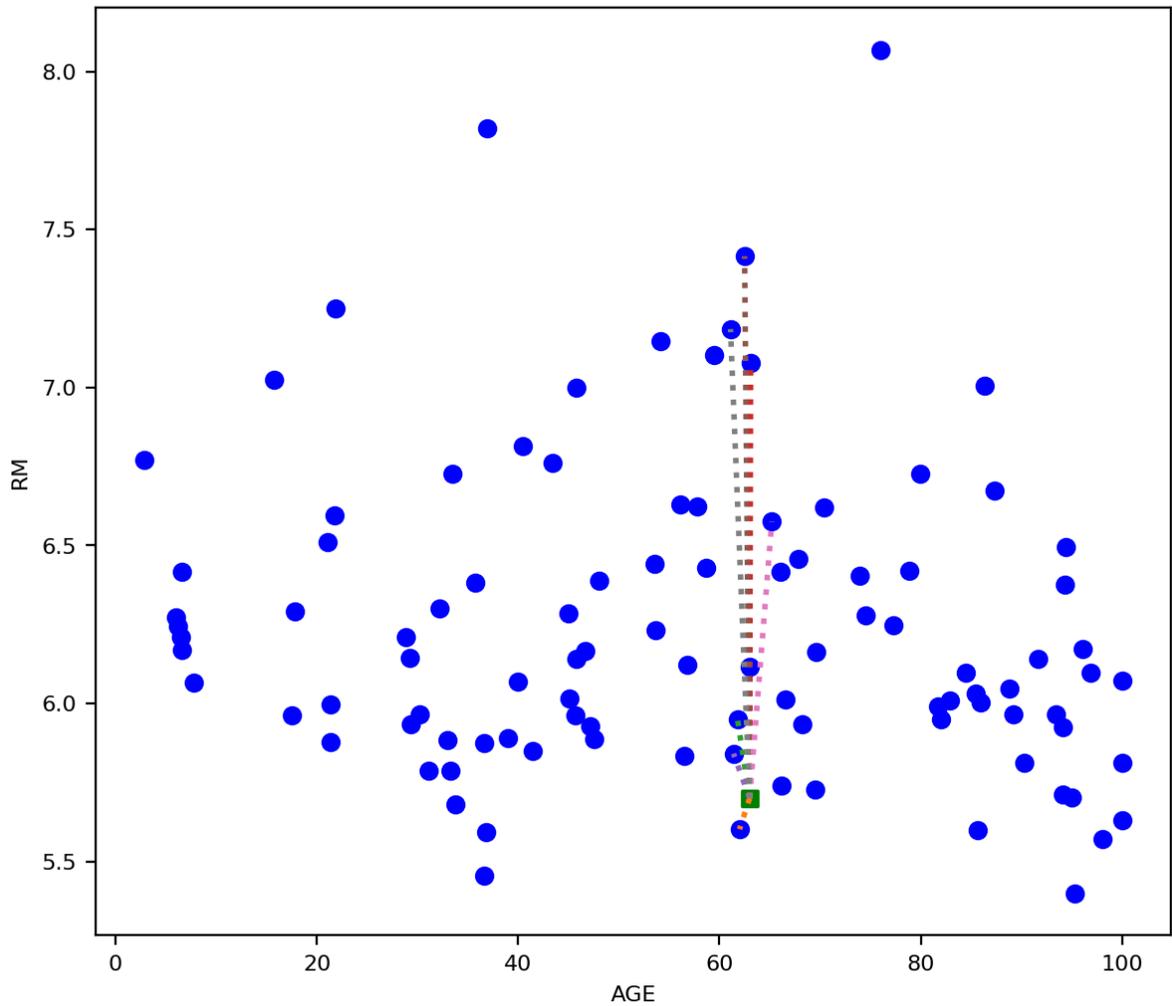


4 neighbours





8 neighbours



Neighbors are found by using only RM! Why do you think that this is happening?

We need to normalize the data so that AGE is also taken into account.

Data normalization

The usual preprocessing in machine learning is **standardize** each variable, that is, transform each column so that it has a mean of 0 and variance of 1.

In sklearn this is achieved with the `StandardScaler`.

```
In [204... from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaler.fit(x)
print("Means of the variables without scaling")
print(x.mean(0))
print("Standard deviations of the variables without scaling")
print(x.std(0))
x_s = scaler.transform(x)
print("Means of the scaled variables")
print(x_s.mean(0))
print("Standard deviations of the scaled variables")
print(x_s.std(0))
xt_s = scaler.transform(xt.reshape(-1,2))[0]
print("Test sample")
print(xt_s)
```

```
Means of the variables without scaling
[56.48613861  6.23928713]
Standard deviations of the variables without scaling
[27.30395478  0.4883953 ]
Means of the scaled variables
[-2.02258452e-16 -2.21494990e-15]
Standard deviations of the scaled variables
[1. 1.]
Test sample
[ 0.23856842 -1.10420213]
```

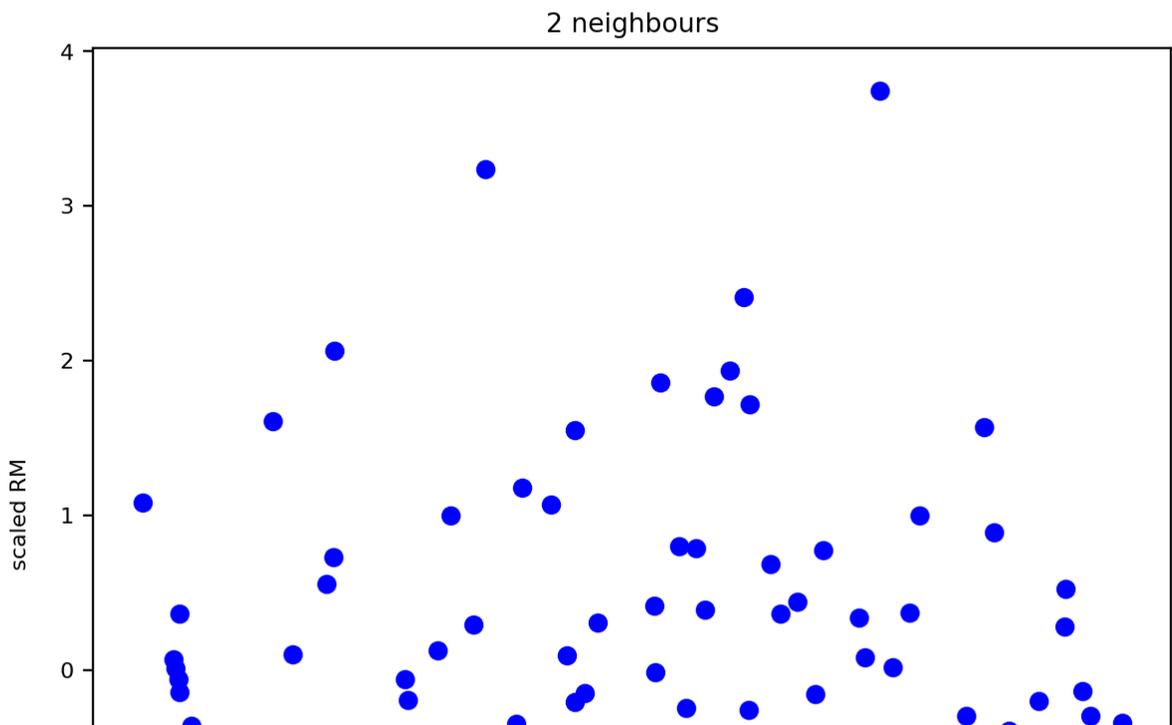
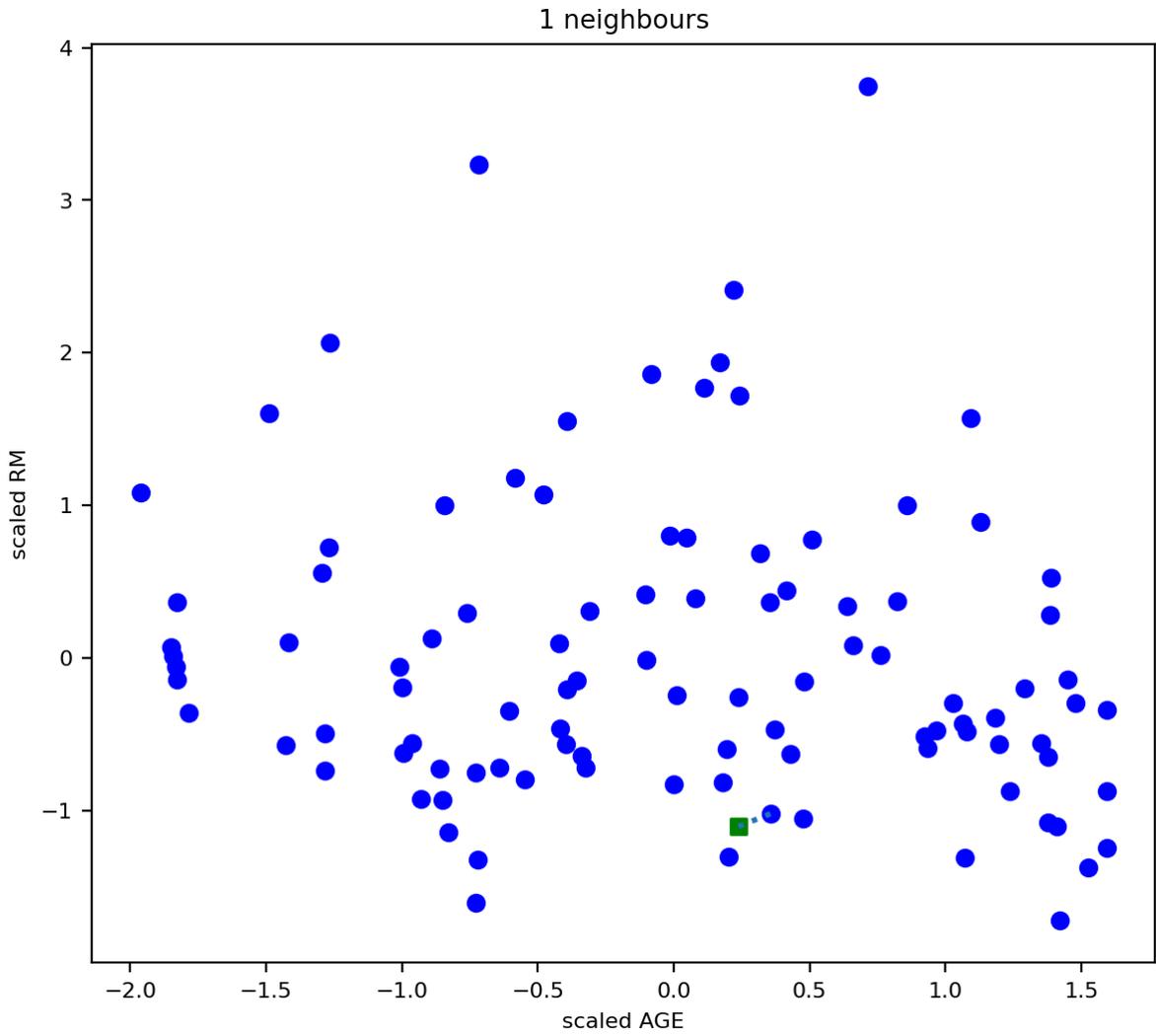
```
In [205... distances = pairwise_distances(xt_s.reshape(1,-1), x_s,'euclidean')

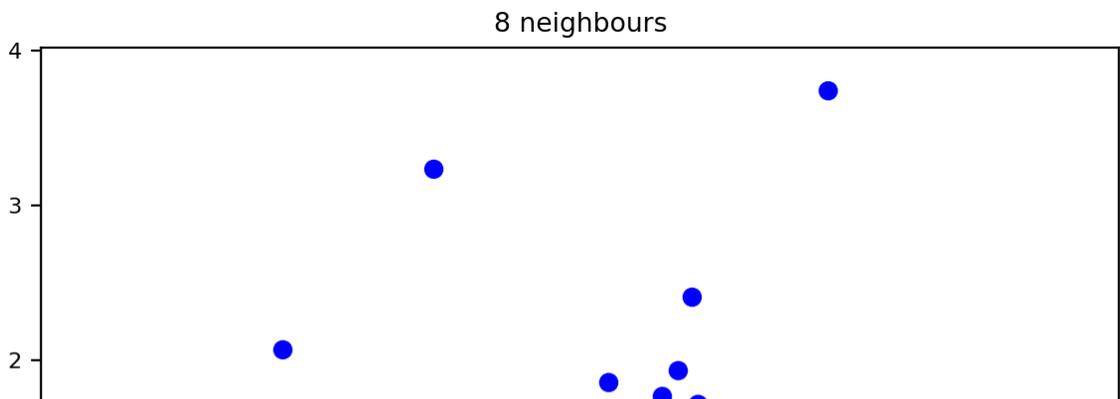
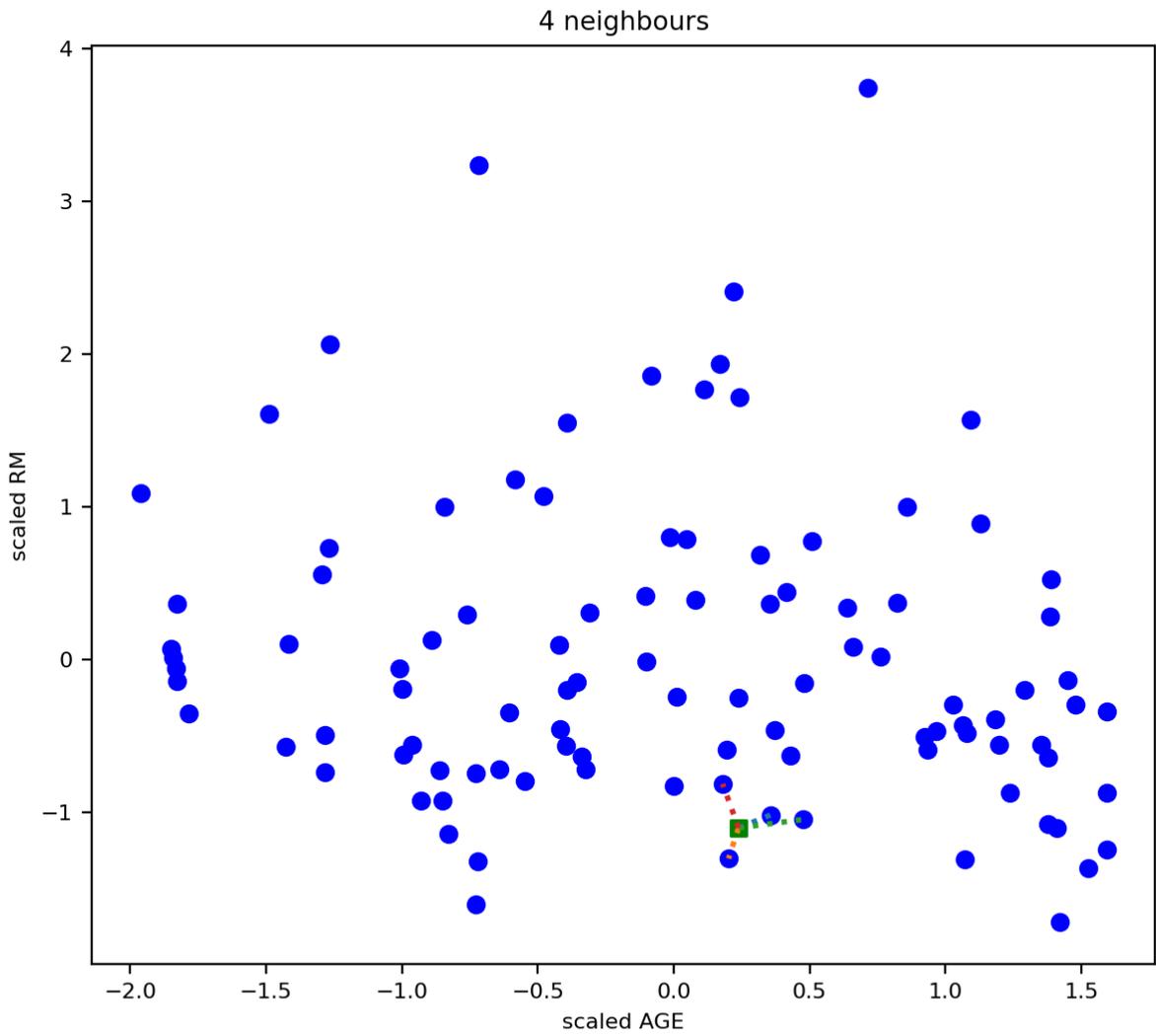
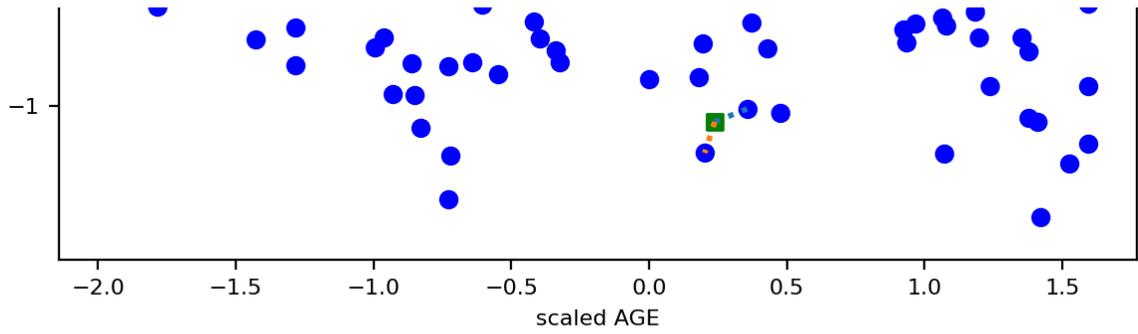
# Sort all neighbours for the test data
id_neigh = np.argsort(distances) # indexes over the training data, whi

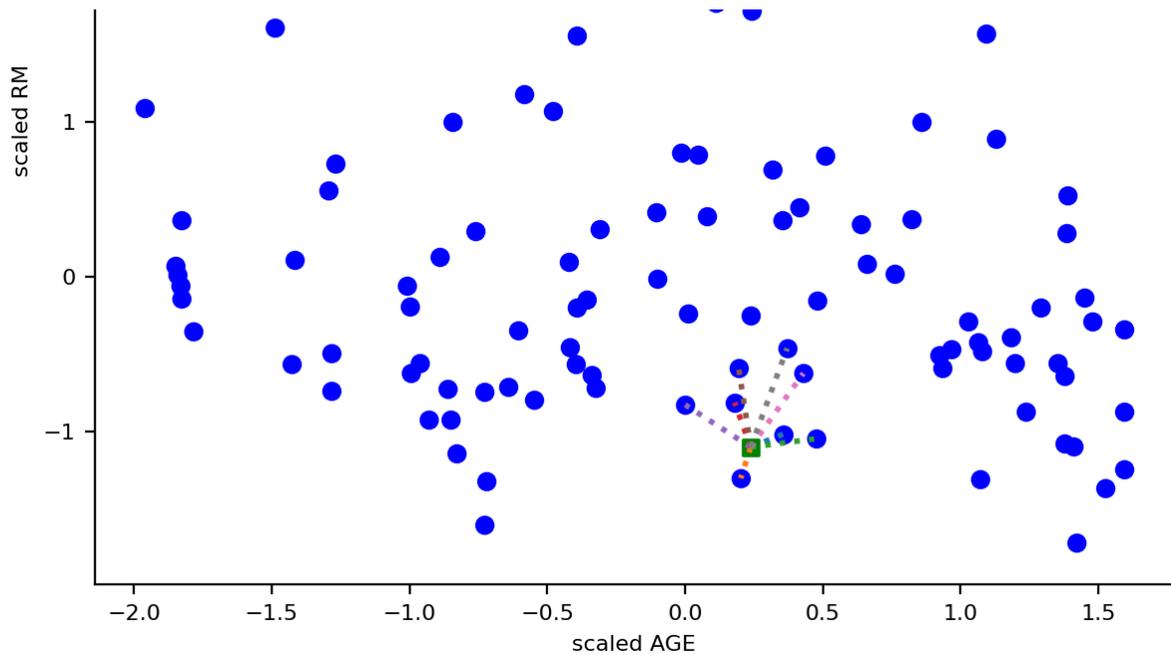
sorted_distances = 1e-6+np.sort(distances) # distances to neighbours h

v_k = np.array([1,2,4,8])
ii=7
fx,ax = plt.subplots(len(v_k),1,figsize=(ii,len(v_k)*ii))
for ik,kk in enumerate(v_k):
    plot_example(x_s,y,xt_s,ax[ik])
    plot_radio(xt_s,np.array([x_s[cc,:] for cc in id_neigh[0,:kk]]),ax[i
ax[ik].set_title('{0:d} neighbours'.format(kk))
ax[ik].set_xlabel('scaled AGE')
```

```
_ = ax[ik].set_ylabel('scaled RM')
```







Influence of the scaling in the ML model result

Now, let's combine scaling and cross-validation to check if the results of k-NN improve using normalized data and the neighbors are calculated with equal impact of the dimensions.

We have to be careful when combining preprocessing with model learning because preprocessing **can only be learned with training samples**, if we use test samples to learn preprocessing, we would be cheating.

1. Results without scaling the variables

```
In [206... from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(data, targets, tes
```

```
In [207... # FILL IN
from sklearn.neighbors import KNeighborsRegressor
for i in range(1,10):
    KNN = KNeighborsRegressor(n_neighbors=i)
    KNN.fit(x_train, y_train)
    R1 = KNN.score(x_train, y_train)
    R2 = KNN.score(x_test, y_test)
    print("R^2 train without scaling: ", R1)
    print("R^2 test without scaling: ", R2)
```

```

R^2 train without scaling: 1.0
R^2 test without scaling: 0.4787149091920948
R^2 train without scaling: 0.831155478417289
R^2 test without scaling: 0.5447309894865939
R^2 train without scaling: 0.7359995089098641
R^2 test without scaling: 0.6202349326037544
R^2 train without scaling: 0.691683878754286
R^2 test without scaling: 0.6060762310692751
R^2 train without scaling: 0.6596756291403296
R^2 test without scaling: 0.5846965270656936
R^2 train without scaling: 0.6258441648192168
R^2 test without scaling: 0.5487428044066404
R^2 train without scaling: 0.6294644970455714
R^2 test without scaling: 0.5324078478249449
R^2 train without scaling: 0.5997214791794261
R^2 test without scaling: 0.5161377811930193
R^2 train without scaling: 0.5894805848606526
R^2 test without scaling: 0.4883436966524203

```

2. Results scaling the variables

```

In [208... # FILL IN
for i in range (1,10):
    scaler = StandardScaler()
    scaler.fit(x_train)
    scaler.fit(x_test)
    x_train_s = scaler.transform(x_train)
    x_test_s = scaler.transform(x_test)

    KNN_s = KNeighborsRegressor(n_neighbors=i)
    KNN_s.fit(x_train_s, y_train)
    R1_s = KNN_s.score(x_train_s, y_train)
    R2_s = KNN_s.score(x_test_s, y_test)
    print("\n ", i, " neighbors")
    print("R^2 train with scaling: ", R1_s)
    print("R^2 test with scaling: ", R2_s)

```

```

1 neighbors
R^2 train with scaling: 1.0
R^2 test with scaling: 0.744718133634135

2 neighbors
R^2 train with scaling: 0.9412916376350793
R^2 test with scaling: 0.7494243445281097

3 neighbors
R^2 train with scaling: 0.9021481139524004
R^2 test with scaling: 0.754632203966948

4 neighbors
R^2 train with scaling: 0.8800547621453377
R^2 test with scaling: 0.75649182792517

```

```
5 neighbors
R^2 train with scaling: 0.8421187671236463
R^2 test with scaling: 0.7467221614449708

6 neighbors
R^2 train with scaling: 0.8040429947249681
R^2 test with scaling: 0.7307623025089145

7 neighbors
R^2 train with scaling: 0.7904798283532792
R^2 test with scaling: 0.7099315019210808

8 neighbors
R^2 train with scaling: 0.7847298213569025
R^2 test with scaling: 0.7072882672118558

9 neighbors
R^2 train with scaling: 0.7718535157913341
R^2 test with scaling: 0.7045531894095389

2 neighbors
R^2 train with scaling: 0.9412916376350793
R^2 test with scaling: 0.7494243445281097

3 neighbors
R^2 train with scaling: 0.9021481139524004
R^2 test with scaling: 0.754632203966948

4 neighbors
R^2 train with scaling: 0.8800547621453377
R^2 test with scaling: 0.75649182792517

5 neighbors
R^2 train with scaling: 0.8421187671236463
R^2 test with scaling: 0.7467221614449708

6 neighbors
R^2 train with scaling: 0.8040429947249681
R^2 test with scaling: 0.7307623025089145

7 neighbors
R^2 train with scaling: 0.7904798283532792
R^2 test with scaling: 0.7099315019210808

8 neighbors
R^2 train with scaling: 0.7847298213569025
R^2 test with scaling: 0.7072882672118558

9 neighbors
R^2 train with scaling: 0.7718535157913341
R^2 test with scaling: 0.7045531894095389
```

In [209... [# Minmax Scaling](#)

```

from sklearn.preprocessing import MinMaxScaler
scaler_mm = MinMaxScaler()
scaler_mm.fit(x_train)
scaler_mm.fit(x_test)
x_train_mm = scaler_mm.transform(x_train)
x_test_mm = scaler_mm.transform(x_test)
KNN_mm = KNeighborsRegressor(n_neighbors=5)
KNN_mm.fit(x_train_mm, y_train)
R1_mm = KNN_mm.score(x_train_mm, y_train)
R2_mm = KNN_mm.score(x_test_mm, y_test)
print("R^2 train with MinMax scaling: ", R1_mm)
print("R^2 test with MinMax scaling: ", R2_mm)

```

R² train with MinMax scaling: 0.8269849352384869
R² test with MinMax scaling: 0.7060914723294562

```

In [210... # robust Scaling
from sklearn.preprocessing import RobustScaler
scaler_rb = RobustScaler()
scaler_rb.fit(x_train)
scaler_rb.fit(x_test)
x_train_rb = scaler_rb.transform(x_train)
x_test_rb = scaler_rb.transform(x_test)
KNN_rb = KNeighborsRegressor(n_neighbors=5)
KNN_rb.fit(x_train_rb, y_train)
R1_rb = KNN_rb.score(x_train_rb, y_train)
R2_rb = KNN_rb.score(x_test_rb, y_test)
print("R^2 train with Robust scaling: ", R1_rb)
print("R^2 test with Robust scaling: ", R2_rb)

```

R² train with Robust scaling: 0.8130610005292604
R² test with Robust scaling: 0.7155938994902362

```

In [211... # max abs scaling
from sklearn.preprocessing import MaxAbsScaler
scaler_ma = MaxAbsScaler()
scaler_ma.fit(x_train)
scaler_ma.fit(x_test)
x_train_ma = scaler_ma.transform(x_train)
x_test_ma = scaler_ma.transform(x_test)
KNN_ma = KNeighborsRegressor(n_neighbors=5)
KNN_ma.fit(x_train_ma, y_train)
R1_ma = KNN_ma.score(x_train_ma, y_train)
R2_ma = KNN_ma.score(x_test_ma, y_test)
print("R^2 train with MaxAbs scaling: ", R1_ma)
print("R^2 test with MaxAbs scaling: ", R2_ma)

```

R² train with MaxAbs scaling: 0.8000387552014049
R² test with MaxAbs scaling: 0.6962324494262893

Although the use of `StandardScaler()` is the most common option for standardizing data, there are other options for this rescaling:

- **MinMaxScaler** : transforms features by scaling each feature to a given range.
- **RobustScaler** : If your data contains many outliers, scaling using the mean and variance of the data is likely to not work very well. In these cases, **RobustScaler** will center the data to the median and component wise scale according to the interquartile range.
- **MaxAbsScaler** : It scales each feature by its maximum absolute value. It is a good choice for sparse data to avoid losing the data sparsity.

Other data transformations

Discretization of continuous variables

Another transformation that may be of interest when working with continuous variables is to change their distribution to a uniform or Gaussian.

The transformation to **uniform** tends to separate the most frequent values of that variable and concentrate the less frequent ones, thus contributing to reduce the impact of *outliers*. As it is a **non-linear** transformation, linear correlations between variables can be distorted.

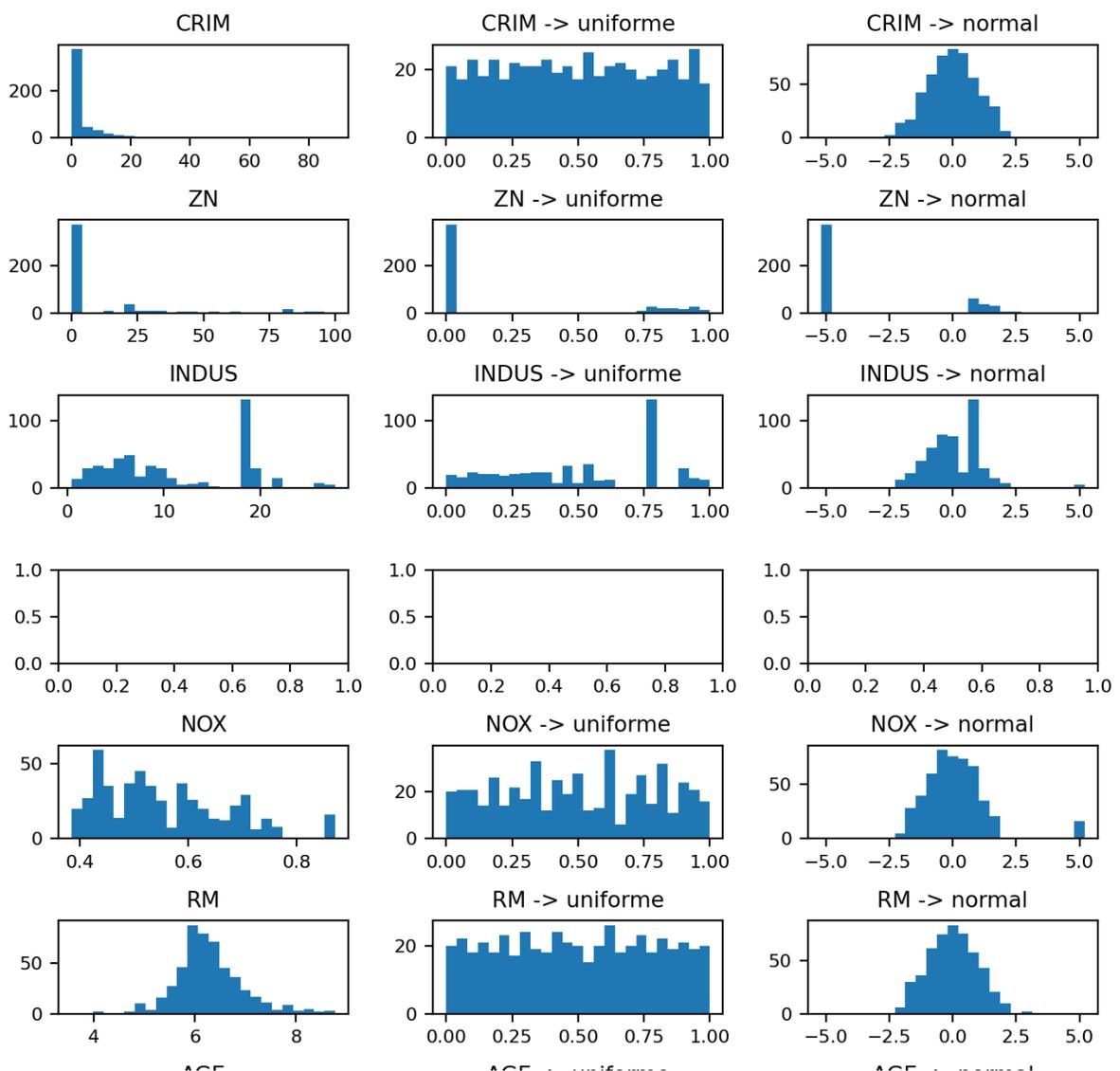
The transformation to **Gaussian** tends to concentrate values around their average. This is very useful (sometimes mandatory) when the subsequent learning model is designed to work with gaussian data distributions (linear regression models, principal component analysis, euclidean distance based models/kernels, ...).

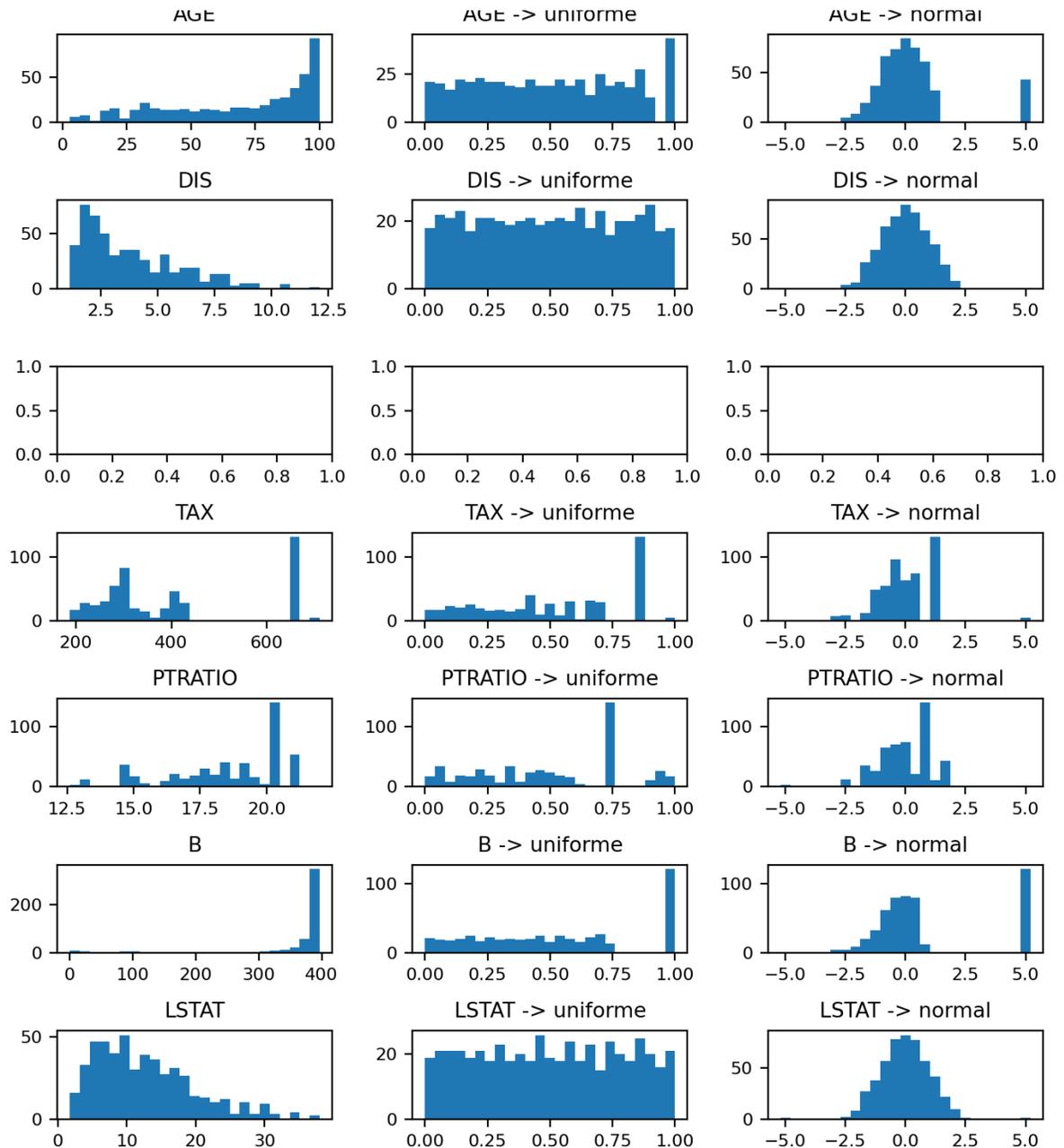
There are several options to get this transformation:

- When the data values expand by several orders of magnitude, but most of the values are around the first order, we can just apply a logarithmic transformation of the data (for example, see the behavior of variable CRIM).
- In case we are working with pandas dataframes, we can use the **quantile** function in combination with the numpy function **digitize** to get a uniform distribution.
- Or, in a more general way, we can use the function **sklearn.preprocessing.quantile_transform** . This function uses an estimate of the cumulative distribution function of each feature to map the original values to the desired output distribution (uniform or gaussian).

Let's see how this last transformation and how it affects the results in the Boston problem.

```
In [212... from sklearn.preprocessing import quantile_transform
ndim = data.shape[1]
fx, ax = plt.subplots(ndim, 3, figsize=(7,14))
nbins = 25
for jj in range(ndim):
    xjj = data.values[:,jj]
    if len(np.unique(xjj)) < nbins:
        continue
    ax[jj][0].hist(xjj,nbins)
    ax[jj][0].set_title('{0}'.format(boston_df.columns[jj]))
    xjj_u = quantile_transform(xjj.reshape(-1,1), n_quantiles=nbins, ran
    ax[jj][1].hist(xjj_u, nbins)
    ax[jj][1].set_title('{0} -> uniforme'.format(boston_df.columns[jj]))
    xjj_n = quantile_transform(xjj.reshape(-1,1), n_quantiles=nbins, ran
    ax[jj][2].hist(xjj_n, nbins)
    ax[jj][2].set_title('{0} -> normal'.format(boston_df.columns[jj]))
fx.tight_layout()
```





We can explore with `GridSearchCV` the impact of these transformations on the final result.

```
In [213... from sklearn.preprocessing import QuantileTransformer
from sklearn.pipeline import Pipeline
from sklearn.model_selection import GridSearchCV
v_nbins = [5,10,20,25,100,200,x_train.shape[1]]
n_neighbours = [1,2,3,4,5,6,7,8,9,10,20,30,40,50,60,70,80,90,100]

# FILL IN
```

One-Hot Encoding Transformation

The **One-Hot encoding** transformation aims to be able to use categorical variables within a model designed for continuous variables.

Let's consider a categorical variable with M possible values, for example the variable `group` that can take 3 values `A`, `B` or `C`. One-hot encoding transforms each categorical variable into M binary variables (they only take the value 1 or 0) which are arranged like this: If an observation has the value m in the categorical variable, all binary variables take the value 0, except the m -th, which takes the value 1.

In the example, the variable `group` would be transformed into 3 variables, [`group_A`, `group_B` and `group_C`] which would code the observations as follows:

- `group = A` → [`group_A = 1`, `group_B = 0`, `group_C = 0`]
- `group = B` → [`group_A = 0`, `group_B = 1`, `group_C = 0`]
- `group = C` → [`group_A = 0`, `group_B = 0`, `group_C = 1`]

In the Boston Housing database, the variable `CHAS`, which measures whether the property faces the Charles River, is coded one-hot, but with the option `drop = "if_binary"`.

To see how this process works in a real problem, we will work with a database of a credit problem ([german dataset](#)). Each observation has data of the clients and of the concrete operation. The `target` indicates if the operation went well (`target=1`) or not (`target=2`). When loading the data we have transformed the `target=2` into zeros to follow the notation usually used. Before applying any transformation, we are going to divide the variables (except the `target`) in two lists:

- `v_cat`, with the names of the variables that are categorical
- `v_num`, with the names of the variables being numeric

```
In [214... def load_data():
    data = pd.read_csv('https://archive.ics.uci.edu/ml/machine-learning-
                      delimiter=" ",
                      header=None)
    data.columns=['existingchecking',
                 'duration',
                 'credithistory',
                 'purpose',
                 'creditamount',
                 'savings',
                 'employmentsince',
                 'installmentrate',
                 'statussex',
                 'otherdebtors',
```

```

        'residencesince',
        'property',
        'age',
        'otherinstallmentplans',
        'housing',
        'existingcredits',
        'job',
        'peopleliable',
        'telephone',
        'foreignworker',
        'target'
    ]
    data.loc[:, 'target'] = data['target'].replace([1,2], [1,0])
    return data

```

```

In [215... data = load_data()
X = data[data.columns[:20]].values
Y = data['target'].values
data.head()

```

```

Out [215...

```

	existingchecking	duration	credithistory	purpose	creditamount	savings	e
0	A11	6	A34	A43	1169	A65	
1	A12	48	A32	A43	5951	A61	
2	A14	12	A34	A46	2096	A61	
3	A11	42	A32	A42	7882	A61	
4	A11	24	A33	A40	4870	A61	

5 rows x 21 columns

```

In [216... # Divide variables into categorial and continuous
v_cat = [element for element in data.dtypes[(data.dtypes=='0')].index.
v_num = [element for element in data.dtypes[(data.dtypes!='0')].index.
target = ['target']

print(v_cat)
print(v_num)

```

```

['existingchecking', 'credithistory', 'purpose', 'savings', 'employment
since', 'statusex', 'otherdebtors', 'property', 'otherinstallmentplan
s', 'housing', 'job', 'telephone', 'foreignworker']
['duration', 'creditamount', 'installmentrate', 'residencesince', 'ag
e', 'existingcredits', 'peopleliable']

```

We will now use sklearn's *one hot encoder* method to code the categorical variables. In addition, we will use the `get_feature_names_out` method to assign a name to each transformed variable so that it is related to the original variable.

```
In [217... from sklearn.preprocessing import OneHotEncoder
onehot = OneHotEncoder()

v_cat_data = onehot.fit_transform(data[v_cat])
v_cat_data = v_cat_data.toarray()

column_names = onehot.get_feature_names_out(v_cat)

one_hot_encoded_frame = pd.DataFrame(v_cat_data, columns=column_names)
```

```
In [218... # We join categorial and continious variables
data_ohe = pd.concat([one_hot_encoded_frame, data[v_num], data[target]],
data_ohe.head()
```

```
Out [218... existingchecking_A11 existingchecking_A12 existingchecking_A13 existingcl
```

	existingchecking_A11	existingchecking_A12	existingchecking_A13	existingcl
0	1.0	0.0	0.0	
1	0.0	1.0	0.0	
2	0.0	0.0	0.0	
3	1.0	0.0	0.0	
4	1.0	0.0	0.0	

5 rows × 62 columns

One hot encoding is quite intuitive but it has **collinearity problems**. If we use linear methods because the columns it generates are not linearly independent, since the sum of all the elements in each row gives 1. To learn a linear model we can add noise to the observations to eliminate collinearity.

Dummy coding corrects the collinearity defect of *one-hot encoding* by removing one of the features. This feature corresponds to a category called **reference category** which is represented with zeros in all features corresponding to the other categories.

In *scikit learn* we use the same class as for *one-hot encoding* but with the parameter `drop` we choose the reference category.

Hashing of features

When the number of feature categories is very large the above techniques produce unwieldy results. In this case, we can use a **Hashing of features** where each category is transformed into a code of m components, which in principle

can be dense. In the *one-hot encoding* variants the code has a length equal to the number of categories (minus one) and with the *hashing* function we reduce this code length to a value of m .

The *hashing* functions are designed to preserve statistics of the original scalar products, so that the *hashed* variables can be used within linear models.

The price paid for this reduction in dimension is that the m categories resulting from *hashing* are **not interpretable** within the a priori knowledge defined by the problem variables.

We can employ the *hashing functionality* of *scikit learn* for this purpose.

In [219...

data[v_cat]

Out [219...

	existingchecking	credithistory	purpose	savings	employmentsince	statu
0	A11	A34	A43	A65	A75	
1	A12	A32	A43	A61	A73	
2	A14	A34	A46	A61	A74	
3	A11	A32	A42	A61	A74	
4	A11	A33	A40	A61	A73	
...
995	A14	A32	A42	A61	A74	
996	A11	A32	A41	A61	A73	
997	A14	A32	A43	A61	A75	
998	A11	A32	A43	A61	A73	
999	A12	A34	A41	A62	A71	

1000 rows x 13 columns

In [220...

```
from sklearn.feature_extraction import FeatureHasher

h =FeatureHasher(n_features=len(v_cat), # number of columns of the out
                 input_type='string')
v_cat_transf_hashing = h.transform(data[v_cat].values) # notice you do

print("Show the output of the hash")
v_cat_transf_hashing.toarray()[:20,:]
```

Show the output of the hash

```

Out[220...] array([[ 0., -1.,  1.,  0., -1.,  0., -1.,  1.,  1.,  0., -1., -1., -
1.],
      [ 0.,  0.,  0., -1., -2.,  0.,  0.,  0.,  1.,  0., -2.,  0., -
1.],
      [ 1.,  0.,  1.,  1., -2.,  1., -1.,  0.,  1.,  0., -1.,  0.,
0.],
      [ 1.,  1.,  1.,  0., -2.,  0., -1.,  0.,  1.,  0.,  0.,  0.,
0.],
      [ 0.,  1.,  1.,  0.,  0.,  0., -1.,  1.,  2., -1.,  0.,  0., -
2.],
      [ 0.,  0.,  1.,  1.,  0.,  1., -1., -1.,  1.,  0.,  0.,  0., -
1.],
      [ 0.,  0.,  1., -1., -2.,  0., -1., -1.,  1.,  0.,  0., -1.,
1.],
      [ 0.,  0.,  1.,  0., -1., -1., -2., -1.,  2.,  0., -1., -1., -
1.],
      [ 1.,  0.,  0., -1., -2.,  1.,  0., -1.,  1., -1., -1.,  0.,
0.],
      [ 0.,  0., -1.,  0.,  0., -1., -1.,  0.,  2.,  0., -1., -1.,
0.],
      [ 1.,  0.,  0.,  0., -1., -1., -1.,  0.,  2.,  0., -1.,  0.,
0.],
      [ 1.,  1.,  0.,  0., -1.,  0., -1.,  0.,  1.,  0.,  0.,  0.,
0.],
      [ 0.,  0.,  0., -1., -2., -1.,  0.,  0.,  1.,  0., -1.,  0., -
1.],
      [ 0.,  0.,  1.,  0., -1.,  0., -1.,  0.,  2.,  0.,  0., -1., -
1.],
      [ 0.,  0.,  0.,  0., -1., -1., -1.,  0.,  2.,  0.,  0.,  0., -
2.],
      [ 0.,  0.,  0., -1., -1.,  0.,  0., -2.,  1.,  0.,  0.,  0., -
2.],
      [ 0., -1.,  1.,  0., -2.,  0., -1.,  1.,  1.,  0.,  0., -1.,
1.],
      [ 2.,  0.,  1., -1.,  0., -1., -1.,  2.,  1., -1.,  0.,  0., -
1.],
      [ 0.,  1.,  0.,  0.,  0.,  0.,  0., -1.,  2.,  0., -1., -2.,
0.],
      [ 0.,  0.,  1., -1., -2., -1., -1., -1.,  1.,  0.,  0., -1.,
0.]])

```

Missing values and data imputation

Part of this section is adapted from this online [tutorial](#)

When we have to face real problems is quite common find that our datasets contain missing values, often encoded as `?`, `nan`, `N/A`, blank cell, or sometimes `-999`, `inf`, `-inf`. The use of these datasets by our scikit-learn estimators is not straightforward since they assume that all values in an array are

numerical.

A basic strategy to use these datasets is to discard entire rows and/or columns containing missing values. However, this comes at the price of losing data which may be valuable (even though incomplete). A better strategy is to impute the missing values, i.e., to infer them from the known part of the data.

The aim of this section is to provide an introduction of missing data and describe some basic methods on how to handle them. For this purpose, let's use the seaborn dataset `tips` and include at random missing values in some of its variables.

```
In [221... import seaborn as sb
from sklearn.model_selection import train_test_split

tips = sb.load_dataset('tips')

# Use a subset of the data for this example
df = tips.loc[:,['total_bill', 'size', 'tip']]
df = df.sample(frac=1).reset_index(drop=True) # Sort the rows

# Split in train/test partitions
data = df[['total_bill', 'size']]
targets = df['tip']
x_train, x_test, y_train, y_test = train_test_split(data, targets, tes

# Add missing values (NaNs) in variables size and total_bill
df_train_all = pd.DataFrame(pd.concat((x_train, y_train), axis=1), colu
df_train = df_train_all.copy()
df_train.loc[:20, 'size'] = np.nan
df_train.loc[150:, 'total_bill'] = np.nan
df_train
```

```
Out [221...

```

	total_bill	size	tip
234	15.98	NaN	3.00
227	14.15	NaN	2.00
180	30.46	NaN	2.00
5	30.40	NaN	5.60
56	28.97	NaN	3.00
...
106	NaN	2.0	3.51
14	NaN	4.0	5.00
92	NaN	2.0	2.00
179	NaN	4.0	4.50
102	NaN	2.0	4.00

170 rows × 3 columns

As we can imagine, the simplest thing is to ignore the missing values and consider only the observations in which all the variables are known.

To drop entries with missing values in any column in pandas, we can use:

```
In [222...
idx = df_train[['total_bill', 'size']].isnull().sum(1)<1
idx
```

Out [222... 0

234 False

227 False

180 False

5 False

56 False

... ...

106 False

14 False

92 False

179 False

102 False

170 rows × 1 columns

dtype: bool

```
In [223... x_train = df_train.loc[df_train[['total_bill', 'size']].isnull().sum(1)
x_train
```

Out [223... total_bill size

```
In [224... y_train = df_train.loc[df_train[['total_bill', 'size']].isnull().sum(1)
y_train
```

Out [224... tip**dtype:** float64

```
In [225... # Remove entries with missing values
df_train_dropna = df_train.dropna(axis=0)
x_train_dropna = df_train_dropna[['total_bill', 'size']]
y_train_dropna = df_train_dropna['tip']

print(f"Original training data shape: {df_train.shape}")
print(f"After dropping missing values: {df_train_dropna.shape}")
print(f"Number of complete rows: {len(x_train_dropna)}")
x_train
y_train
```

```
Original training data shape: (170, 3)
After dropping missing values: (0, 3)
Number of complete rows: 0
```

Out[225... **tip**

dtype: float64

```
In [226... from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsRegressor
from sklearn.model_selection import GridSearchCV

# Use the complete data (dropna approach)
# Since dropna resulted in empty dataset, we'll show why this approach
print(f"Available complete samples after dropna: {len(x_train_dropna)}")

if len(x_train_dropna) > 5: # Only run if we have enough data for CV
    param_grid = {
        'knn__n_neighbors': range(1,10),
        'knn__weights':['uniform','distance'],
    }
    pipe = Pipeline([('scaler', StandardScaler()), ('knn', KNeighborsR
    grid_knn = GridSearchCV(pipe, param_grid, cv=5)
    grid_knn.fit(x_train_dropna, y_train_dropna)
    R2_train = grid_knn.score(x_train_dropna, y_train_dropna)
    R2_test = grid_knn.score(x_test, y_test)
    print("R^2 en el conjunto de entrenamiento: {0:.2f}".format(R2_tra
    print("R^2 en el conjunto de test: {0:.2f}".format(R2_test))
    print(grid_knn.best_params_)
else:
    print("ERROR: Not enough complete samples to perform cross-validat
    print("This demonstrates why dropping all rows with missing values
    print(f"Original training samples: {len(df_train)}")
    print(f"Complete samples after dropna: {len(x_train_dropna)}")
    print("We need imputation methods instead of dropping rows.")
```

```
Available complete samples after dropna: 0
ERROR: Not enough complete samples to perform cross-validation!
This demonstrates why dropping all rows with missing values is problema
tic.
Original training samples: 170
Complete samples after dropna: 0
We need imputation methods instead of dropping rows.
```

In general, this method should not be used unless the proportion of missing values is very small (<5%) since this approach has the cost of having less data and the result is highly likely to be biased.

So, to avoid apply this scheme, here we will introduce some **imputation** schemes able to replace the missing values with some guessed/estimated ones.

Mean, median, mode imputation

A simple guess of a missing value is the mean, median, or mode (most frequently appeared value) of that variable. This strategy is called **univariate**, since it imputes values in the i -th feature dimension using only non-missing values in that feature dimension.

In pandas, `.fillna` can be used to replace `NA`'s with a specified value. For instance, next cell uses the mean to impute the missing values in `size` and `total_bill`.

```
In [227... # Mean imputation for size and total_bill
mean_size = df_train['size'].mean()
mean_total_bill = df_train['total_bill'].mean()
df_mean = df_train.fillna(value={'size':mean_size, 'total_bill':mean_t
df_mean
```

```
Out [227...
   total_bill  size  tip
234  15.980000   3.0  3.00
227  14.150000   3.0  2.00
180  30.460000   3.0  2.00
  5  30.400000   3.0  5.60
 56  28.970000   3.0  3.00
...         ...   ...   ...
106  20.225932   2.0  3.51
 14  20.225932   4.0  5.00
 92  20.225932   2.0  2.00
179  20.225932   4.0  4.50
102  20.225932   2.0  4.00
```

170 rows × 3 columns

Sklearn also includes imputation facilities. In particular the `SimpleImputer` class provides basic strategies for imputing missing values. Missing values can be imputed with a provided constant value, or using the statistics (mean, median or most frequent) of each column in which the missing values are located. This class also allows for different missing values encodings and when using the

most_frequent or constant strategy it supports categorical data imputation. For example...

```
In [228... from sklearn.impute import SimpleImputer
imp = SimpleImputer(missing_values=np.nan, strategy='mean')
imp.fit(df_train[['total_bill', 'size']])
imp_values=imp.transform(df_train[['total_bill', 'size']])
pd.DataFrame(imp_values,columns=['total_bill', 'size'])
```

```
Out [228...
   total_bill  size
0  15.980000   3.0
1  14.150000   3.0
2  30.460000   3.0
3  30.400000   3.0
4  28.970000   3.0
...         ...   ...
165 20.225932   2.0
166 20.225932   4.0
167 20.225932   2.0
168 20.225932   4.0
169 20.225932   2.0
```

170 rows × 2 columns

Besides, the fact of using a sklearn class allows us integrate the missing value imputation in our pipeline....

```
In [229... # Extract x_train and y_train
x_train = df_train[['total_bill', 'size']]
y_train = df_train['tip']
```

```
In [230... from sklearn.impute import SimpleImputer
param_grid = {
    'imp_strategy' : ['mean', 'most_frequent', 'median'],
    'kNN_n_neighbors': range(1,10),
    'kNN_weights':['uniform', 'distance'],
}
pipe = Pipeline([('imp', SimpleImputer(missing_values=np.nan)), ('scal
grid_knn = GridSearchCV(pipe, param_grid, cv=5)
grid_knn.fit(x_train, y_train)
R2_train = grid_knn.score(x_train, y_train)
```



```
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='mea
n'.
  warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='mea
n'.
  warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='mea
n'.
  warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='mea
n'.
  warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='mos
t_frequent'.
  warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='mos
t_frequent'.
  warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='mos
t_frequent'.
  warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='mos
t_frequent'.
  warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='mea
n'.
  warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='mea
```



```
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='med
ian'.
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='med
ian'.
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='med
ian'.
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='med
ian'.
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='med
ian'.
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='med
ian'.
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='med
ian'.
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='med
ian'.
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/impute/_base.py:635: Us
erWarning: Skipping features without any observed values: ['size']. At
least one non-missing value is needed for imputation with strategy='med
ian'.
R^2 en el conjunto de entrenamiento: 0.10
R^2 en el conjunto de test: 0.34
{'imp__strategy': 'median', 'kNN__n_neighbors': 5, 'kNN__weights': 'uni
form'}
```

Linear Regression imputation

Univariate imputation only look at the distribution of the values of the variable with missing entries. If we know there is a correlation between the missing value and other variables, we can often get better guesses by regressing the missing variable on other variables. This strategy is called **multivariate imputation**.

Let's check the correlation among our variables.

```
In [231... ## Regression imputation for total_bill

# Check correlation among variables
corr = df_train_dropna.corr()
corr
```

```
Out [231...

```

	total_bill	size	tip
total_bill	NaN	NaN	NaN
size	NaN	NaN	NaN
tip	NaN	NaN	NaN

As we can see, `tip` is the most correlated variable with `total_bill`. Thus, we can use a simple linear model regressing `total_bill` on `tip` to fill the missing values in `total_bill`.

```
In [232... from sklearn.linear_model import LinearRegression

# 1. Create training data with a subset of data where there are no missing
df_bill_tip = df_train.dropna(axis=0, subset=['total_bill', 'tip'])
# Create training data: X (tip) and y (total_bill)
X = df_bill_tip[['tip']]
y = df_bill_tip['total_bill']

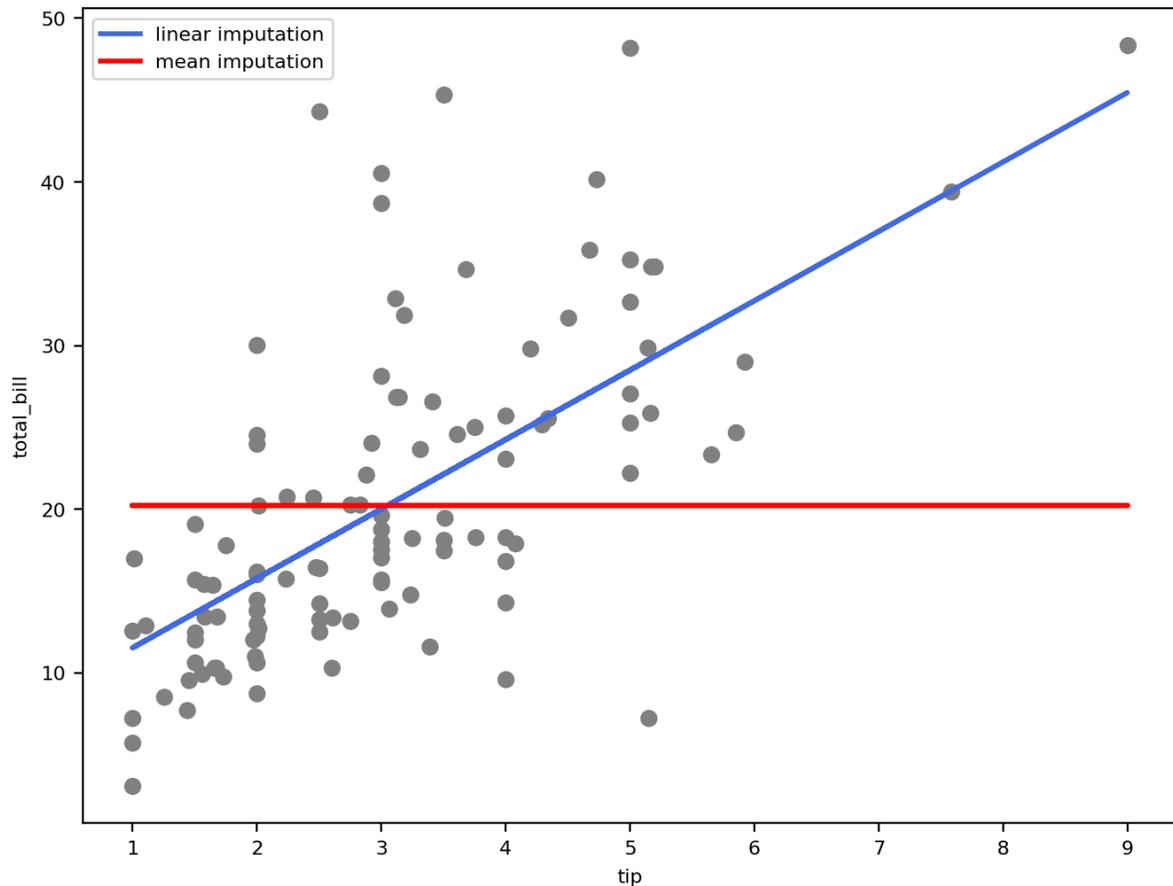
# 2. Create test data to impute missing values
# find the entries with total_bill missing
missing_bill = df_train['total_bill'].isnull()
# extract the tips of observations with total_bill missing
tip_misbill = pd.DataFrame(df_train['tip'][missing_bill])

## 3. Train a linear regression model to estimate total bill from tip

# fit a linear model
lm = LinearRegression().fit(X, y)
# use fitted model and tip values to predict missing total_bill
bill_pred = lm.predict(tip_misbill)
```

```
In [233... # Visualize the imputed total_bill values comparing to its true missing
plt.scatter(tip_misbill, df_train_all['total_bill'][missing_bill], color='gray', linewidth=1)
plt.plot(tip_misbill, bill_pred, color='royalblue', linewidth=2, label='linear imputation')
plt.plot(tip_misbill, mean_total_bill*np.ones(tip_misbill.shape), color='red', linewidth=2, label='mean imputation')

plt.xlabel("tip")
plt.ylabel("total_bill")
plt.legend()
plt.show()
```



As we can see, the imputed `total_bill` from a simple linear model from `tips` does not exactly recover the truth but capture the general trend (and it is better than a single value imputation such as mean imputation). We can, of course, use more variables in the regression model to get better imputation.

K-nearest neighbour (KNN) imputation

Besides model-based imputation like regression imputation, neighbor-based imputation can also be used. K-nearest neighbor (KNN) imputation is an example of neighbor-based imputation. For a discrete variable, KNN imputer uses the most frequent value among the k nearest neighbors and, for a continuous variable, use the mean.

To use KNN for imputation, first, a KNN model is trained using complete data. For continuous data, commonly the Euclidean is used distance metric, whereas for discrete data, hamming distance is a frequent choice.

In the next example we use `total_bill` and `tip` to impute `size` values.

```
In [234... from sklearn.neighbors import KNeighborsRegressor

# For the sake of simplicity, we set K=3
knn = KNeighborsRegressor(n_neighbors=3, weights = 'distance')
knn.fit(df_train_dropna.loc[:,['total_bill', 'tip']], df_train_dropna.

# Find the missing values in size
missing_size = df_train['size'].isnull()
# extract the tips of observations with total_bill missing
df_missing_size = pd.DataFrame(df_train[['total_bill', 'tip']][missing

# used trained K-NN to predict missing sizes
imputed_size = knn.predict(df_missing_size)
```

```
-----
ValueError                                Traceback (most recent call l
ast)
/tmp/ipython-input-3266564550.py in <cell line: 0>()
      3 # For the sake of simplicity, we set K=3
      4 knn = KNeighborsRegressor(n_neighbors=3, weights = 'distance')
----> 5 knn.fit(df_train_dropna.loc[:,['total_bill', 'tip']], df_train_
dropna.loc[:, 'size'])
      6
      7 # Find the missing values in size

/usr/local/lib/python3.12/dist-packages/sklearn/base.py in wrapper(esti
mator, *args, **kwargs)
    1387         )
    1388         ):
-> 1389         return fit_method(estimator, *args, **kwargs)
    1390
    1391         return wrapper

/usr/local/lib/python3.12/dist-packages/sklearn/neighbors/_regression.p
y in fit(self, X, y)
    220         The fitted k-nearest neighbors regressor.
    221         """
--> 222         return self._fit(X, y)
    223
    224         def predict(self, X):

/usr/local/lib/python3.12/dist-packages/sklearn/neighbors/_base.py in _
fit(self, X, y)
    476         if self.__sklearn_tags__().target_tags.required:
```

```

477         if not isinstance(X, (KDTree, BallTree, NeighborsBa
se)):
--> 478             X, y = validate_data(
479                 self,
480                 X,

/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py in
validate_data(_estimator, X, y, reset, validate_separately, skip_check_
array, **check_params)
2959         y = check_array(y, input_name="y", **check_y_param
s)
2960     else:
-> 2961         X, y = check_X_y(X, y, **check_params)
2962         out = X, y
2963

/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py in
check_X_y(X, y, accept_sparse, accept_large_sparse, dtype, order, copy,
force_writeable, force_all_finite, ensure_all_finite, ensure_2d, allow_
nd, multi_output, ensure_min_samples, ensure_min_features, y_numeric, e
stimator)
1368     ensure_all_finite = _deprecate_force_all_finite(force_all_f
inite, ensure_all_finite)
1369
-> 1370     X = check_array(
1371         X,
1372         accept_sparse=accept_sparse,

/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py in
check_array(array, accept_sparse, accept_large_sparse, dtype, order, co
py, force_writeable, force_all_finite, ensure_all_finite, ensure_non_ne
gative, ensure_2d, allow_nd, ensure_min_samples, ensure_min_features, e
stimator, input_name)
1128         n_samples = _num_samples(array)
1129         if n_samples < ensure_min_samples:
-> 1130             raise ValueError(
1131                 "Found array with %d sample(s) (shape=%s) while
a"
1132                 " minimum of %d is required%s."

ValueError: Found array with 0 sample(s) (shape=(0, 2)) while a minimum
of 1 is required by KNeighborsRegressor.

```

```

In [235... # Compare the imputed size distribution to its true missing values
size_obs = df_train['size'].dropna(axis=0)
size_mis = df_train_all[['size']][missing_size]
size_knn_imp = pd.DataFrame(imputed_size)
size_mean_imp = df_mean[['size']][missing_size]

legend_1 = ['observed_sizes', 'missing sizes']
legend_2 = ['observed_sizes', 'knn imputed sizes']
legend_3 = ['observed_sizes', 'mean imputed sizes']

```

```

plt.figure(figsize=(15,5))
plt.subplot(1,3,1)
plt.hist([size_obs, np.squeeze(size_mis)], color=['gold', 'limegreen'])
plt.xlabel("size")
plt.ylabel("Frequency")
plt.legend(legend_1)

plt.subplot(1,3,2)
plt.hist([size_obs, np.squeeze(size_knn_imp)], color=['gold', 'royalbl
plt.xlabel("size")
plt.ylabel("Frequency")
plt.legend(legend_2)

plt.subplot(1,3,3)
plt.hist([size_obs, np.squeeze(size_mean_imp)], color=['gold', 'viole
plt.xlabel("size")
plt.ylabel("Frequency")
plt.legend(legend_3)

plt.show()

```

```

-----
NameError                                Traceback (most recent call l
ast)
/tmp/ipython-input-3757783298.py in <cell line: 0>()
      1 # Compare the imputed size distribution to its true missing val
ues
      2 size_obs = df_train['size'].dropna(axis=0)
----> 3 size_mis = df_train_all[['size']][missing_size]
      4 size_knn_imp = pd.DataFrame(imputed_size)
      5 size_mean_imp = df_mean[['size']][missing_size]

NameError: name 'missing_size' is not defined

```

In the plot above, we compared the missing `sizes` and imputed `sizes` using both 3NN imputer and mode imputation. As we can see, KNN imputer gives much better imputation than ad-hoc methods like mode imputation.

In general, KNN imputer is simple, flexible (can be used to any type of data), and easy to interpret. However, if the dataset is large, using a KNN imputer could be slow.

Besides, sklearn has a `KNNImputer` class which provides imputation for filling in missing values using the K-NN approach. Besides, this implementation includes several ad-hoc functionalities:

- It includes a euclidean distance metric that supports missing values, `nan_euclidean_distances`.
- By default, each missing feature is imputed using values from

`n_neighbors` nearest neighbors that have a value for the feature.

- If a sample has more than one feature missing, then the neighbors for that sample can be different depending on the particular feature being imputed.
- When the number of available neighbors is less than `n_neighbors`, the mean value for that feature is used during imputation.
- If a feature is always missing in training, it is removed during transform.

The following cell shows how to use this class:

```
In [ ]: from sklearn.impute import KNNImputer
imp = KNNImputer(n_neighbors=2, weights="uniform")
imp.fit(df_train[['total_bill', 'size']])
imp_values=imp.transform(df_train[['total_bill', 'size']])
pd.DataFrame(imp_values,columns=['total_bill', 'size'])
```

Note that in the above code, we may want to learn the imputation model with a training data partition and later predict values (`transform`) over both train and test data; besides, we can also use the training labels to learn the model and improve our imputation.

Anyway, this can be done for us integrating the imputation with a ML pipeline...

```
In [ ]: from sklearn.impute import KNNImputer
param_grid = {
    'imp__n_neighbors' : range(1,10,2),
    'imp__weights': ['uniform', 'distance'],
    'knn__n_neighbors': range(1,10,2),
    'knn__weights': ['uniform', 'distance'],
}
pipe = Pipeline([('imp', KNNImputer()), ('scaler', StandardScaler()),
grid_knn = GridSearchCV(pipe, param_grid, cv=5)
grid_knn.fit(x_train, y_train)
R2_train = grid_knn.score(x_train, y_train)
R2_test = grid_knn.score(x_test, y_test)
print("R^2 en el conjunto de entrenamiento: {:.2f}".format(R2_train))
print("R^2 en el conjunto de test: {:.2f}".format(R2_test))
print(grid_knn.best_params_)
```

Iterative Imputation

A more sophisticated approach is to use the `IterativeImputer` class of sklearn, which models each feature with missing values as a function of other features, and uses that estimate for imputation. It does so in an iterated fashion: at each step, a feature column is designated as output y and the other feature columns are treated as inputs X . A regressor is fit on (X, y) for known y . Then, the regressor is used to predict the missing values of y . This is done for

each feature in an iterative fashion, and then is repeated for `max_iter` imputation rounds. Finally, the model returns the values of the final imputation round.

```
In [236... # This estimator is still experimental for now
# explicitly require this experimental feature
from sklearn.experimental import enable_iterative_imputer
from sklearn.impute import IterativeImputer
imp = IterativeImputer(max_iter=5, random_state=0)
imp.fit(df_train[['total_bill', 'size']])
imp_values=imp.transform(df_train[['total_bill', 'size']])
pd.DataFrame(imp_values, columns=['total_bill', 'size'])
```

```
Out [236...
   total_bill  size
0  15.980000  3.0
1  14.150000  3.0
2  30.460000  3.0
3  30.400000  3.0
4  28.970000  3.0
...         ...   ...
165  20.225932  2.0
166  20.225932  4.0
167  20.225932  2.0
168  20.225932  4.0
169  20.225932  2.0
```

170 rows × 2 columns

```
In [237... from sklearn.impute import IterativeImputer
from warnings import filterwarnings
filterwarnings('ignore')
param_grid = {
    'imp_estimator' : [KNeighborsRegressor(), LinearRegression()],
    'knn_n_neighbors': range(1,10,2),
    'knn_weights': ['uniform', 'distance'],
}
pipe = Pipeline([('imp', IterativeImputer(max_iter=20, tol=0.1, random
grid_knn = GridSearchCV(pipe, param_grid, cv=5)
grid_knn.fit(x_train, y_train)
R2_train = grid_knn.score(x_train, y_train)
R2_test = grid_knn.score(x_test, y_test)
print("R^2 en el conjunto de entrenamiento: {:.2f}".format(R2_train))
```

```
print("R^2 en el conjunto de test: {:.2f}".format(R2_test))
print(grid_knn.best_params_)
```

R² en el conjunto de entrenamiento: 0.20

R² en el conjunto de test: 0.34

```
{'imp__estimator': LinearRegression(), 'kNN__n_neighbors': 5, 'kNN__weights': 'uniform'}
```

Other imputation methods

So far, we have talked about some common methods that can be used for missing data imputation. Depending on the nature of the data or data type, some other imputation methods may be more appropriate. For example:

- For **longitudinal data**, such as patients' weights over a period of visits, it might make sense to use last valid observation to fill the NA's. In pandas, this can be done using the `ffill` method in `.fillna`.

```
df_l = df.fillna(method='ffill')
```

- For **time-series data**, it might make sense to use interpolation of observed values before and after a timestamp for missing values. In pandas, various interpolation methods (e.g. polynomial, splines) can be implemented using `.interpolation`.

```
df_interp = df.interpolate(method = 'linear',
limit_direction = 'forward', axis=0)
```

The mean, median, mode imputation, regression imputation, KNN imputer are all methods that create a single replacement value for each missing entry. **Multiple Imputation (MI)**, rather than a different method, is more like a general approach/framework of doing the imputation procedure multiple times to create different plausible imputed datasets. The key motivation to use MI is that a single imputation cannot reflect sampling variability from both sample data and missing values, whereas the final analysis results (with different imputations) allow the data scientist to obtain understanding of how analytic results may differ as a consequence of the inherent uncertainty caused by the missing values.

However, when the final user is not interested in measuring this uncertainty, we need to properly combine this information to provide a single imputation value. In this sense, the iterative imputer is considered a MI approach with Chained Equations to combine all the imputation values.

Outlier detection

Outliers are rare values that deviate from other observations on data; in other words, an outlier is an observation that is far from the rest of the observations or the center of mass of observations. Outlier detection estimators, thus, try to fit the regions where the training data is the most concentrated, ignoring the deviant observations.

Outliers can be caused human, instrumental or processing errors, for example, during the data acquisition, recording or processing or, measurement or experimental errors.

In machine learning is important to identify and remove outliers from data when training algorithms for predictive modeling. Outliers can skew statistical measures and data distributions, providing a misleading representation of the underlying data and relationships. Removing outliers from training data prior to modeling can result in a better fit of the data and, in turn, more skillful predictions.

As outlier detection consists in detecting the observations that are far from the rest of the observations, most of the methods are based in probability density estimation, such us using the Gaussian Mixture Model (GMM) or One-class SVM (1-SVM). Besides, sklearn includes other three specific approaches:

- Isolation Forest
- Minimum Covariance Determinant
- Local Outlier Factor

As in previous sections, let's consider a dataset, in this case a synthetic one that we will contaminate with outliers, and then we will see method by method analyzing its capacity to detect outliers.

```
In [238... from sklearn.datasets import make_blobs
# Generate synthetic data
X, y_true = make_blobs(n_samples=450, centers=4,
                      cluster_std=0.60, random_state=0)

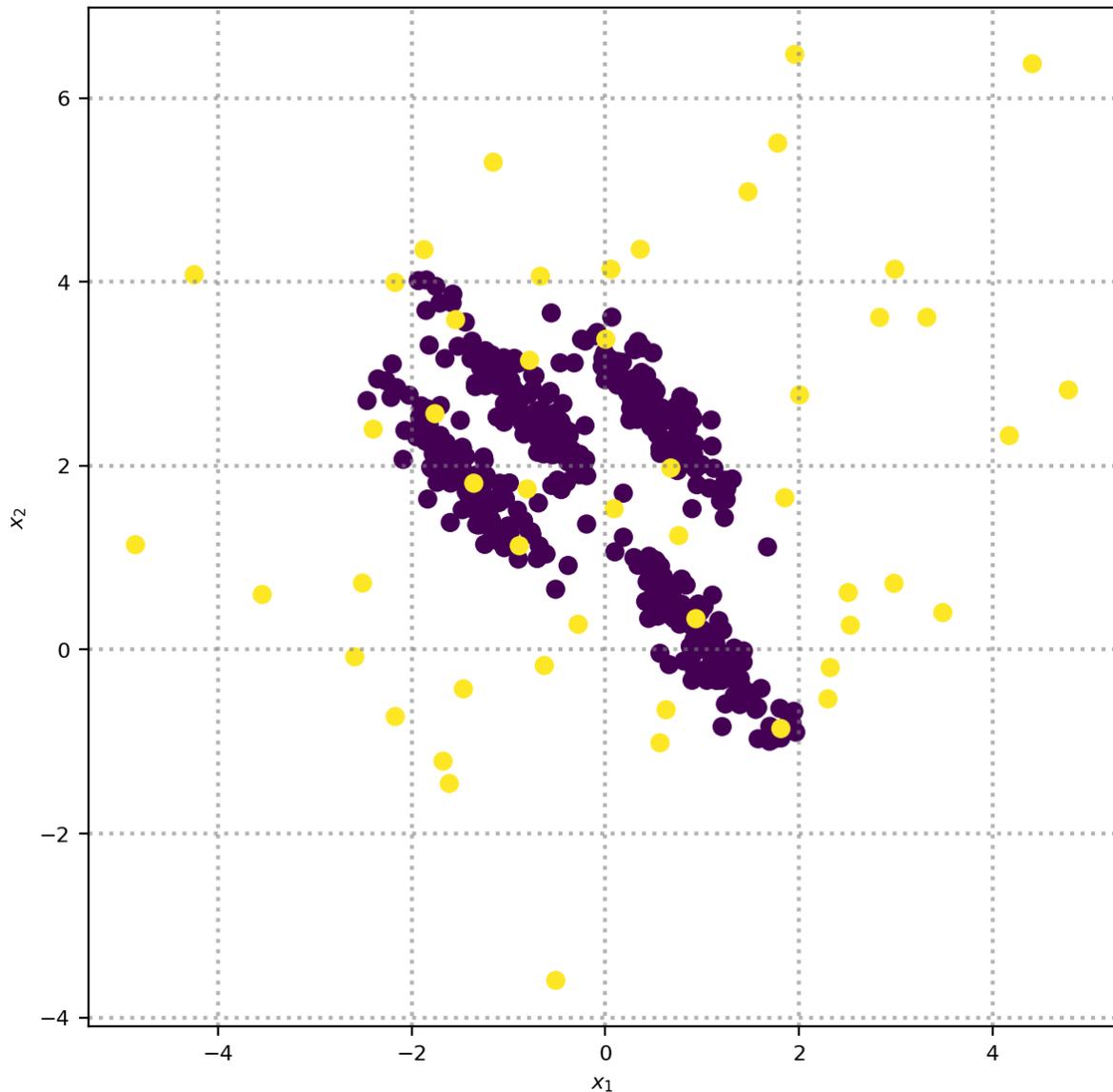
rng = np.random.RandomState(13)
X_stretched = np.dot(X, rng.randn(2, 2))

# Add outliers
center_outliers = np.array([0,2]).reshape([1,-1])
X_outlier,_ = make_blobs(n_samples=50,centers=center_outliers, cluster
X_out_example = np.vstack([X_stretched,X_outlier])
y_outlier = np.zeros([X_out_example.shape[0],])
y_outlier[-50:] = 1 # Vector binario que vale 1 en los outliers introd
```

```

fig, ax = plt.subplots(figsize=(7, 7))
ax.scatter(X_out_example[:, 0], X_out_example[:, 1], c=y_outlier, s=40)
ax.set_xlabel('$x_1$')
ax.set_ylabel('$x_2$')
ax.grid(which='major', color='gray', alpha=0.6, linestyle='dotted', lw

```



When generating the outliers, we have entered data in random positions, so it is possible that some of them will fall on the inliers distribution and logically they will not be outliers and we will not be able to detect them.

Isolation Forest

The `IsolationForest` *isolates* observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature.

Since recursive partitioning can be represented by a tree structure, the number of splittings required to isolate a sample is equivalent to the path length from the root node to the terminating node. The most typical data will need more iterations (splittings) or longer path lengths to be isolated; so this path length, averaged over a forest of such random trees, is a measure of normality and our decision function for outlier detection. Note that random partitioning produces noticeably shorter paths for anomalies. Hence, when a forest of random trees collectively produce shorter path lengths for particular samples, they are highly likely to be anomalies.

Perhaps the most important hyperparameter in the model is the `contamination` argument, which is used to help estimate the number of outliers in the dataset. This is a value between 0.0 and 0.5 and by default is set to 0.1.

```
In [250... from sklearn.ensemble import IsolationForest
# identify outliers in the training dataset
iForest = IsolationForest(contamination=0.1)
youtliers = iForest.fit_predict(X_out_example)

idx_outliers = np.where(youtliers == -1)[0]

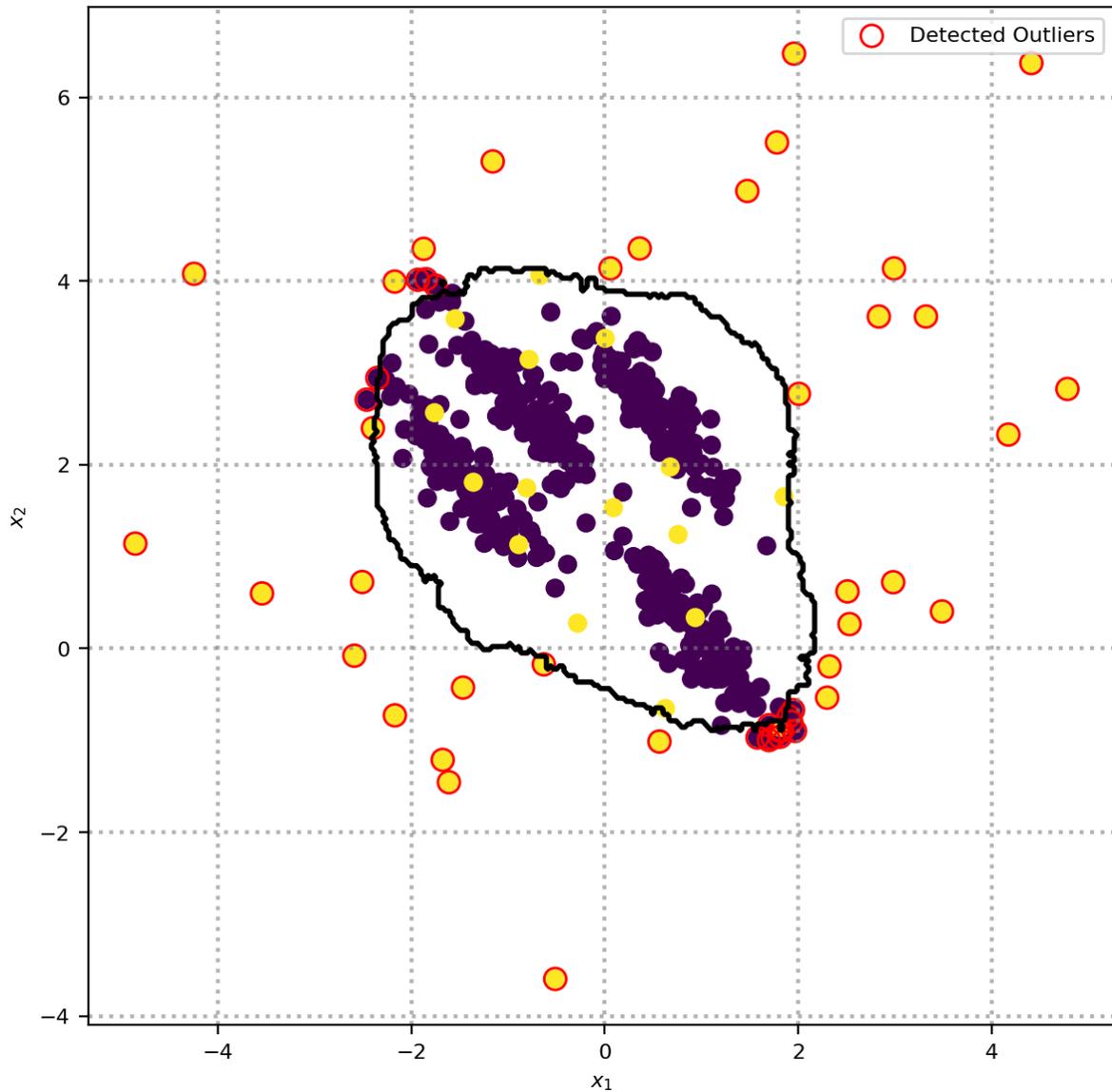
fig, ax = plt.subplots(figsize=(7, 7))

ax.scatter(X_out_example[:, 0], X_out_example[:, 1], c=y_outlier, s=40)
ax.scatter(X_out_example[idx_outliers, 0], X_out_example[idx_outliers, 1],
plt.legend()
ax.set_xlabel('$x_1$')
ax.set_ylabel('$x_2$')
ax.grid(which='major', color='gray', alpha=0.6, linestyle='dotted', lw

# Plot boundary
# We created a grid
intervals = 200
x = np.linspace(-3, 3, intervals)
y = np.linspace(-2, 5, intervals)
X,Y = np.meshgrid(x,y)
xys = np.vstack([X.ravel(), Y.ravel()]).T
ygrid = iForest.predict(xys)
ygrid = ygrid.reshape([intervals,intervals])

ax.contour(X, Y, ygrid, levels=[0], linewidths=2, colors='black')
```

```
Out[250... <matplotlib.contour.QuadContourSet at 0x78a283c4d580>
```



```
In [240... df_out = pd.DataFrame(X_out_example, columns=['x1', 'x2'])
df_out['anomaly_score'] = iForest.decision_function(X_out_example)
df_out['anomaly'] = youtliers
df_out
```

```
Out [240...
```

	x1	x2	anomaly_score	anomaly
0	0.768658	2.364070	0.118761	1
1	0.848419	2.539324	0.111895	1
2	0.898920	-0.063674	0.100955	1
3	0.896827	1.533960	0.090779	1
4	-0.701717	2.601242	0.126777	1
...
495	4.784731	2.825824	-0.186368	-1
496	1.957472	6.476287	-0.217820	-1
497	1.854925	1.652729	-0.001544	-1
498	-4.853358	1.142175	-0.179058	-1
499	0.090980	1.533816	0.090761	1

500 rows x 4 columns

Minimum Covariance Determinant or Elliptic Envelope

If the input variables have a Gaussian distribution, then simple statistical methods can be used to detect outliers. For example, if the dataset has two input variables and both are Gaussian, then the feature space forms a multi-dimensional Gaussian and knowledge of this distribution can be used to identify values far from the distribution. This approach can be generalized by defining a hypersphere (ellipsoid) that covers the normal data, and data that falls outside this shape is considered an outlier.

The `EllipticEnvelope` class of `sklearn` implements this approach and it provides the `contamination` argument that defines the expected ratio of outliers to be observed.

```
In [241... from sklearn.covariance import EllipticEnvelope

# identify outliers in the training dataset
ee = EllipticEnvelope(contamination=0.1)
youtliers = ee.fit_predict(X_out_example)

idx_outliers = np.where(youtliers == -1)[0]
```

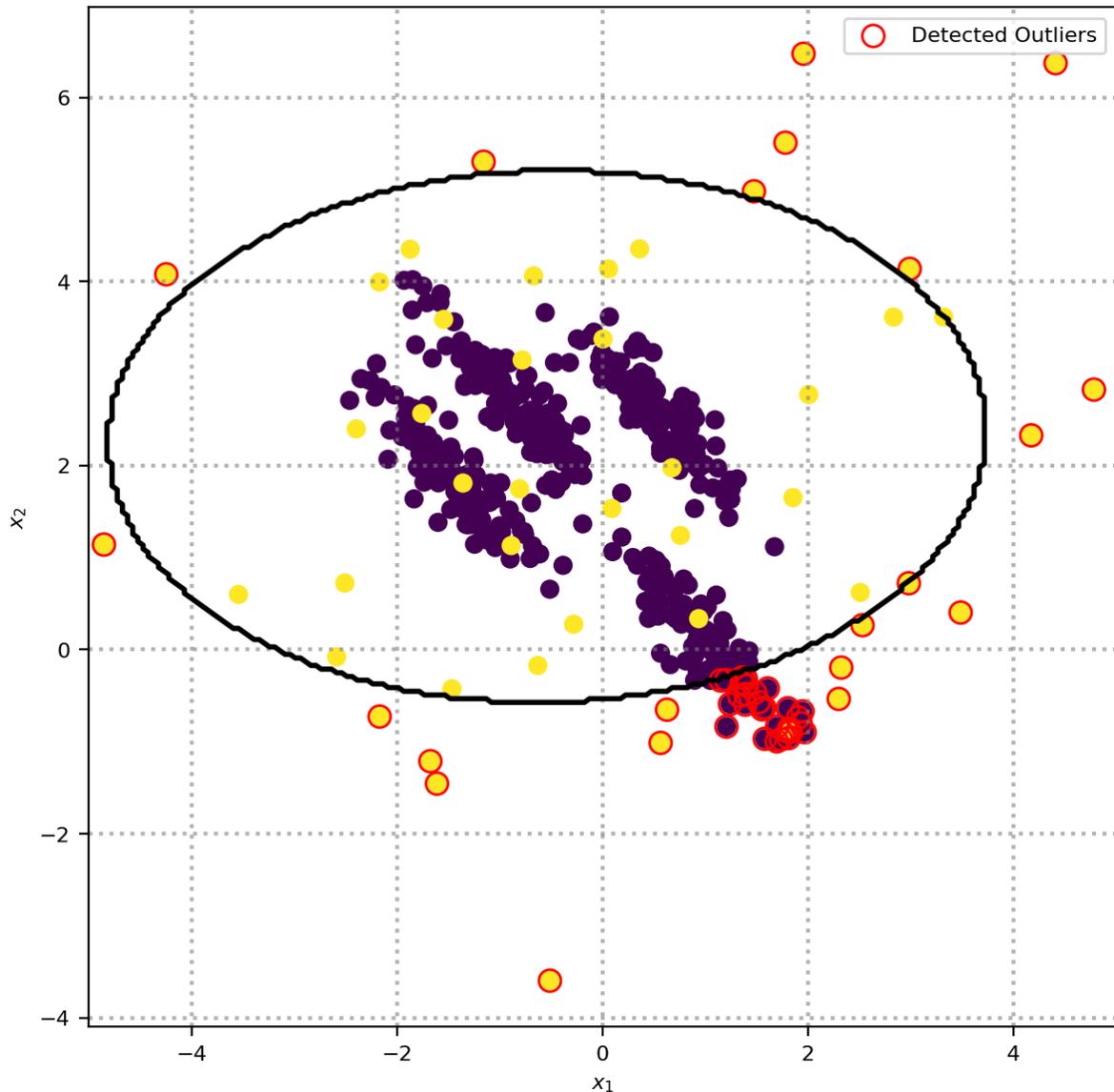
```
fig, ax = plt.subplots(figsize=(7, 7))

ax.scatter(X_out_example[:, 0], X_out_example[:, 1], c=y_outlier, s=40)
ax.scatter(X_out_example[idx_outliers, 0], X_out_example[idx_outliers, 1], c='red', s=40)
plt.legend()
ax.set_xlabel('$x_1$')
ax.set_ylabel('$x_2$')
ax.grid(which='major', color='gray', alpha=0.6, linestyle='dotted', lw=1)

# Plot boundary
# We created a grid
intervals = 200
x = np.linspace(-5, 5, intervals)
y = np.linspace(-2, 6, intervals)
X, Y = np.meshgrid(x, y)
xys = np.vstack([X.ravel(), Y.ravel()]).T
ygrid = ee.predict(xys)
ygrid = ygrid.reshape([intervals, intervals])

ax.contour(X, Y, ygrid, levels=[0], linewidths=2, colors='black')
```

Out[241]... <matplotlib.contour.QuadContourSet at 0x78a28aa14bc0>



Local Outlier Factor

The **local outlier factor**, or LOF for short, measures the local deviation of density of a given sample with respect to its neighbors. It is local in that the anomaly score depends on how isolated the object is with respect to the surrounding neighborhood. More precisely, locality is given by a KNN, whose distance is used to estimate the local density. By comparing the local density of a sample to the local densities of its neighbors, one can identify samples that have a substantially lower density than their neighbors.

This can work well for feature spaces with low dimensionality (few features), although it can become less reliable as the number of features is increased, referred to as the *curse of dimensionality*.

Here we find two parameters to adjust, the `contamination` level and the

number of neighbors (set to 20 by default).

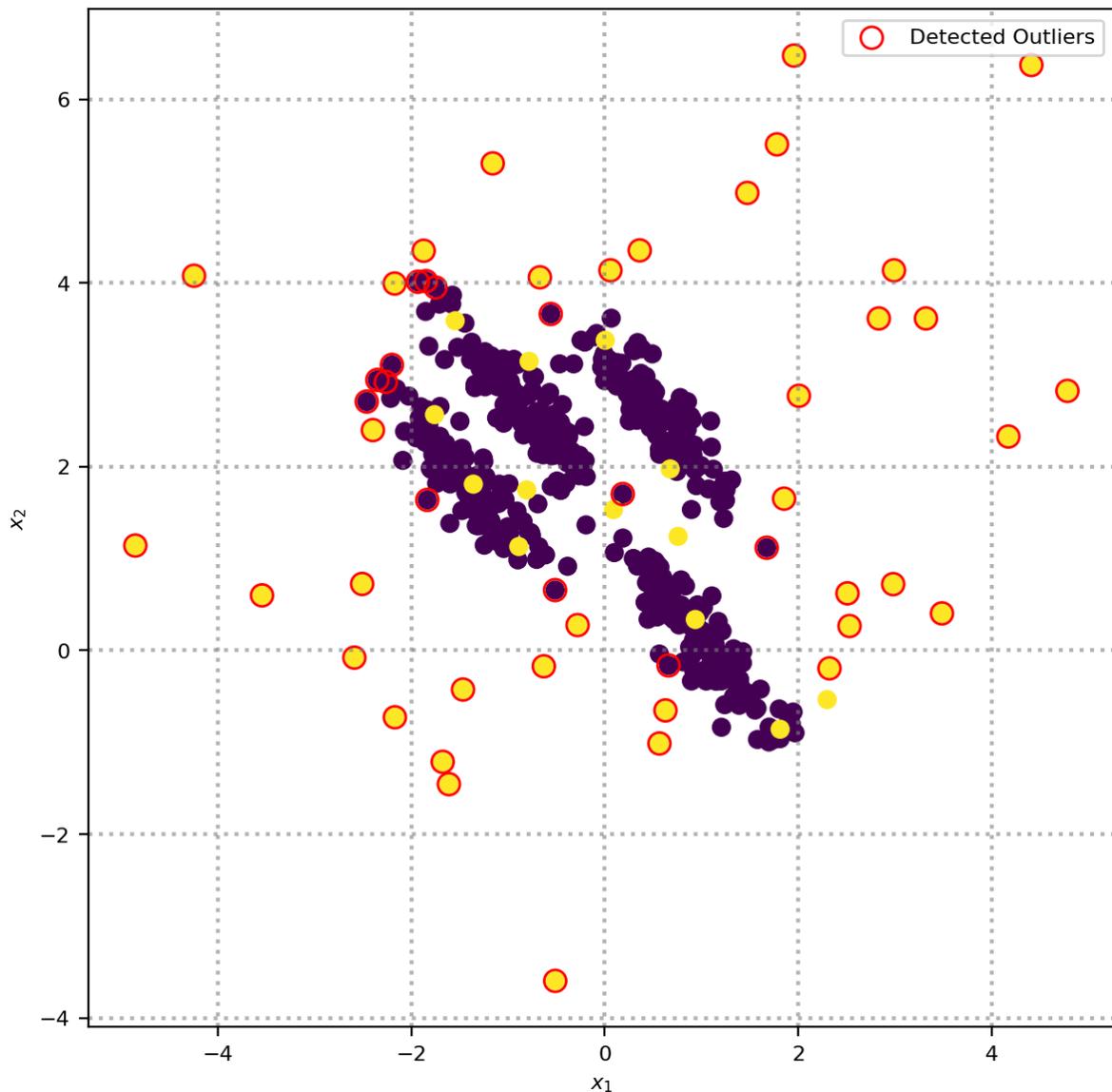
```
In [242... from sklearn.neighbors import LocalOutlierFactor

# identify outliers in the training dataset
lof = LocalOutlierFactor(contamination=0.1)
youtliers = lof.fit_predict(X_out_example)

idx_outliers = np.where(youtliers == -1)[0]

fig, ax = plt.subplots(figsize=(7, 7))

ax.scatter(X_out_example[:, 0], X_out_example[:, 1], c=y_outlier, s=40)
ax.scatter(X_out_example[idx_outliers, 0], X_out_example[idx_outliers,
plt.legend()
ax.set_xlabel('$x_1$')
ax.set_ylabel('$x_2$')
ax.grid(which='major', color='gray', alpha=0.6, linestyle='dotted', lw
```



One-Class SVM

The support vector machine, or SVM, algorithm developed initially for binary classification can be used for one-class classification, that is, it is able to establishing the borders of the data support. This way, the 1-SVM captures the density of the majority class and classifies examples on the extremes of the density function as outliers.

The scikit-learn library provides an implementation of 1-SVM in the `OneClassSVM` class. This model requires the choice of a kernel and a contamination parameter to define a boundary. Normally a kernel RBF is chosen although there is no exact formula or algorithm to establish its bandwidth parameter. The contamination parameter known as the SVM boundary of a class is called `nu` and corresponds to the approximate proportion of outliers in the data set. The 'OneClassSVM' class sets the default value of `nu` to 0.1.

```
In [252... from sklearn.svm import OneClassSVM
from matplotlib import cm

# identify outliers in the training dataset
oneSVM = OneClassSVM(nu=0.1, kernel="rbf", gamma=0.1)
youtliers = oneSVM.fit_predict(X_out_example)

idx_outliers = np.where(youtliers == -1)[0]

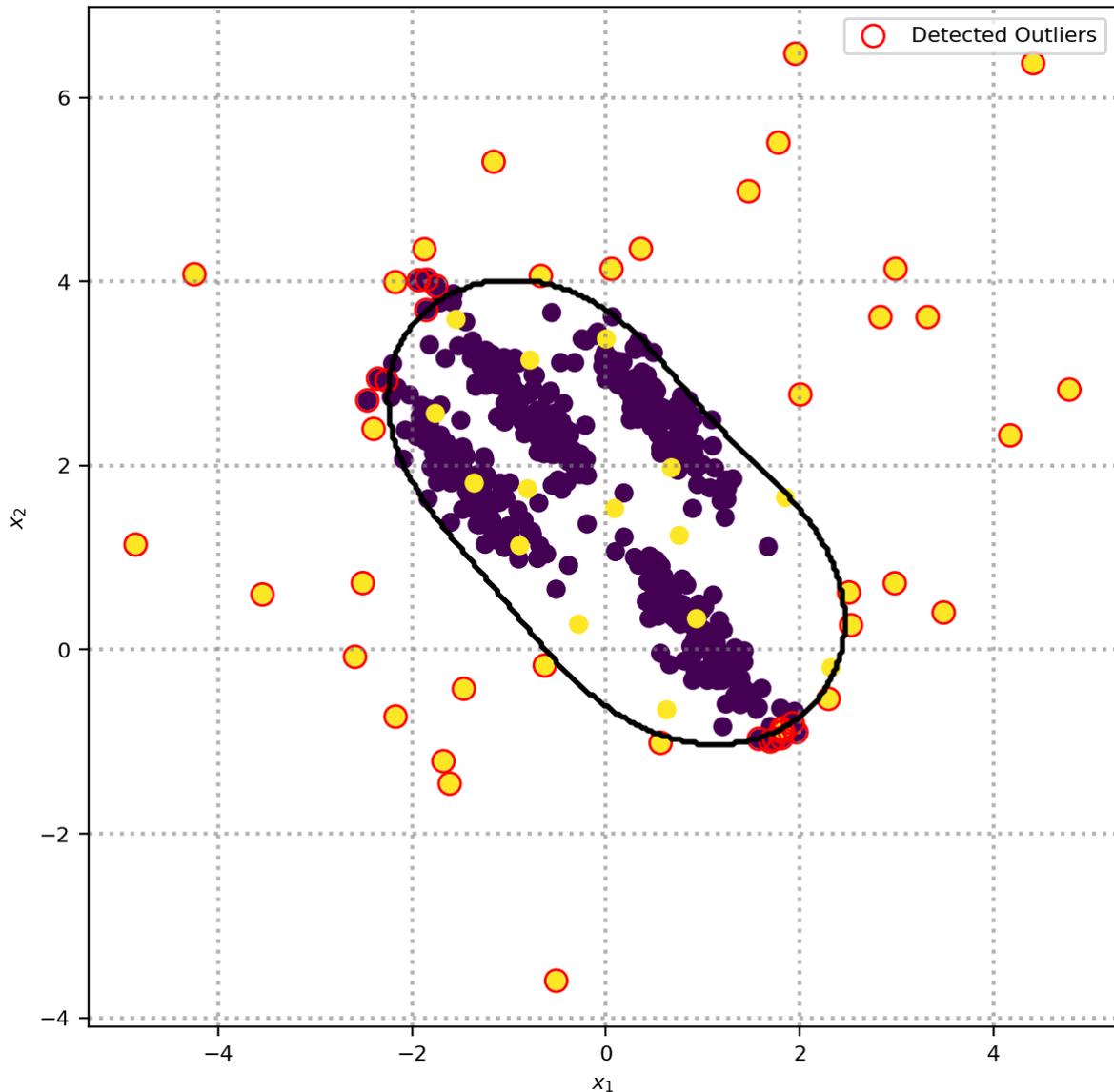
fig, ax = plt.subplots(figsize=(7, 7))

ax.scatter(X_out_example[:, 0], X_out_example[:, 1], c=y_outlier, s=40)
ax.scatter(X_out_example[idx_outliers, 0], X_out_example[idx_outliers, 1],
plt.legend()
ax.set_xlabel('$x_1$')
ax.set_ylabel('$x_2$')
ax.grid(which='major', color='gray', alpha=0.6, linestyle='dotted', lw

# Plot boundary
# We created a grid
intervals = 200
x = np.linspace(-3, 3, intervals)
y = np.linspace(-1.5, 4.5, intervals)
X,Y = np.meshgrid(x,y)
xys = np.vstack([X.ravel(), Y.ravel()]).T
ygrid = oneSVM.predict(xys)
ygrid = ygrid.reshape([intervals,intervals])

ax.contour(X, Y, ygrid, levels=[0], linewidths=2, colors='black')
```

```
Out[252... <matplotlib.contour.QuadContourSet at 0x78a28a57b2f0>
```



Gaussian Mixture Models (GMM)

The [mixing models](#) are a type of probabilistic data model that allows approximate the probability density of the data and, thus, use this information to detect values far from the distribution. A GMM model will adjust a parameterized probability density function to our data as follows:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k) \quad (1)$$

where

- π_k is the probability that the data comes from the k -th gaussian.
- We assume that all data assigned to the k -th gaussian are distributed

according to a Gaussian mean of μ_k and a σ_k covariance matrix.

- $(\pi_1, \dots, \pi_K), (\mu_1, \dots, \mu_K), (\Sigma_1, \dots, \Sigma_K)$ are the **model parameters**.

These model parameters are chosen to maximize the **probability of already observed data or evidence**:

$$\max_{(\pi_1, \dots, \pi_K), (\mu_1, \dots, \mu_K), (\Sigma_1, \dots, \Sigma_K)} \sum_{n=1}^N \log p(\mathbf{x}_n)$$

This problem is solved numerically by an iterative algorithm known as **EM** (Expectation-Maximization).

Although this approach is not included in sklearn as a novelty detection technique, we can use its [GaussianMixture](#) implementation for data modeling and clustering, as it is commonly used.

Note that this approach can be considered an extension of Elliptic Envelope approach where instead of considering a single gaussian to model our data, we use a mixture of gaussian to be able to model more complex data distributions.

```
In [254... from sklearn.mixture import GaussianMixture

# As we know the dataset, we set K=4 (otherwise this parameters has to
gmm = GaussianMixture(n_components=4, covariance_type='full', n_init=200

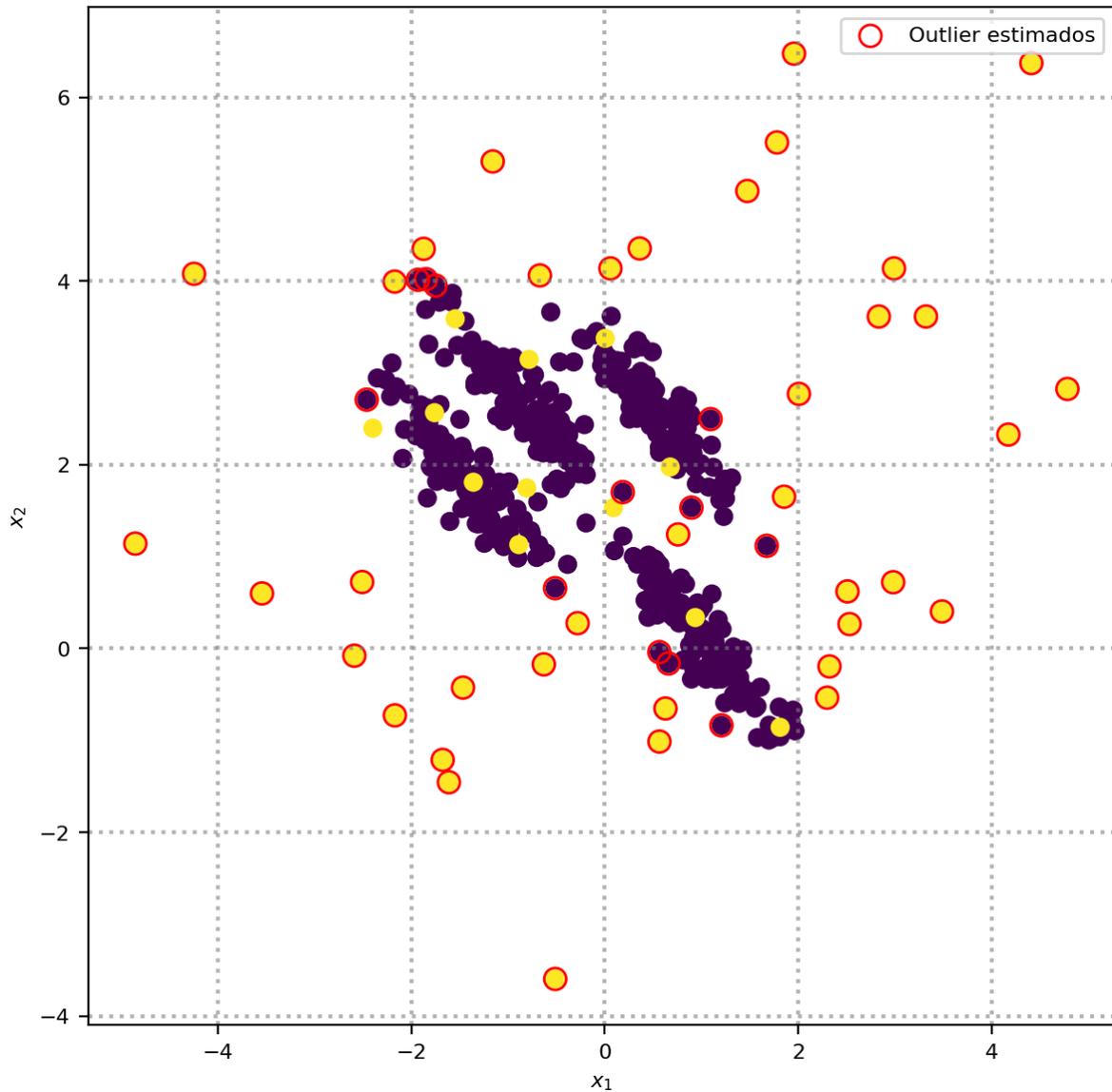
gmm.fit(X_out_example)

log_probs = gmm.score_samples(X_out_example)

y_outlier_pred = np.zeros([X_out_example.shape[0],])
frac = 0.1 # We set the number of contaminated data
idx_outliers = np.argsort(log_probs)[:int(np.round(X_out_example.shape

fig, ax = plt.subplots(figsize=(7, 7))

ax.scatter(X_out_example[:, 0], X_out_example[:, 1], c=y_outlier, s=40
ax.scatter(X_out_example[idx_outliers, 0], X_out_example[idx_outliers,
plt.legend()
ax.set_xlabel('$x_1$')
ax.set_ylabel('$x_2$')
ax.grid(which='major', color='gray', alpha=0.6, linestyle='dotted', lw
```



For this model we can plot the contour curves, that is, the curves over space (x_1, x_2) with equal probability $p(\mathbf{x})$. Also the learned $p(\mathbf{x})$ in a 3D figure. In this way, we can analyze the data distribution.

```
In [255... from scipy.stats import multivariate_normal as mvn #Multivariate normal
from mpl_toolkits.mplot3d import Axes3D
from matplotlib import cm

# Lets plot the pdf contour plot
intervals = 200

# We created a grid
x = np.linspace(-3, 3, intervals)
y = np.linspace(-1.5, 4.5, intervals)

X,Y = np.meshgrid(x,y)

xys = np.vstack([X.ravel(), Y.ravel()]).T
```

```

# We evaluate  $p(x)$  for each point of the grid
K=4
Zgmm = np.zeros(len(xys))
for k in range(K):
    Zgmm += gmm.weights_[k]*mvn(gmm.means_[k,:], gmm.covariances_[k]).

fig = plt.figure(figsize=(14, 7))
ax = fig.add_subplot(1, 2, 1)

Zgmm = Zgmm.reshape([intervals,intervals])
ax.contour(X, Y, Zgmm, 20, cmap=cm.coolwarm)
ax.scatter(X_out_example[:, 0], X_out_example[:, 1], c=y_outlier, s=40)
ax.scatter(gmm.means_[:, 0], gmm.means_[:, 1], s=40, marker='^', facec

ax.grid(which='major', color='gray', alpha=0.6, linestyle='dotted', lw
ax.set_xlabel('$x_1$')
ax.set_ylabel('$x_2$')

ax = fig.add_subplot(1, 2, 2, projection='3d')
surf = ax.plot_surface(X, Y, Zgmm, cmap=cm.coolwarm)
ax.grid(which='major', color='gray', alpha=0.6, linestyle='dotted', lw
fig.colorbar(surf, shrink=0.5, aspect=5)

plt.show()

```

