

Data Imputation

Learn about data imputation and the various methods to accomplish it.

Chapter Goals:

- Learn different methods for imputing data

A. Data imputation methods

In real life, we often have to deal with data that contains missing values. Sometimes, if the dataset is missing too many values, we just don't use it. However, if only a few of the values are missing, we can perform data imputation ([https://en.wikipedia.org/wiki/Imputation_\(statistics\)](https://en.wikipedia.org/wiki/Imputation_(statistics))) to substitute the missing data with some other value(s).

There are many different methods for data imputation. In scikit-learn, the `SimpleImputer` (<https://scikit-learn.org/stable/modules/generated/sklearn.impute.SimpleImputer.html#sklearn.impute.SimpleImputer>) transformer performs four different data imputation methods.

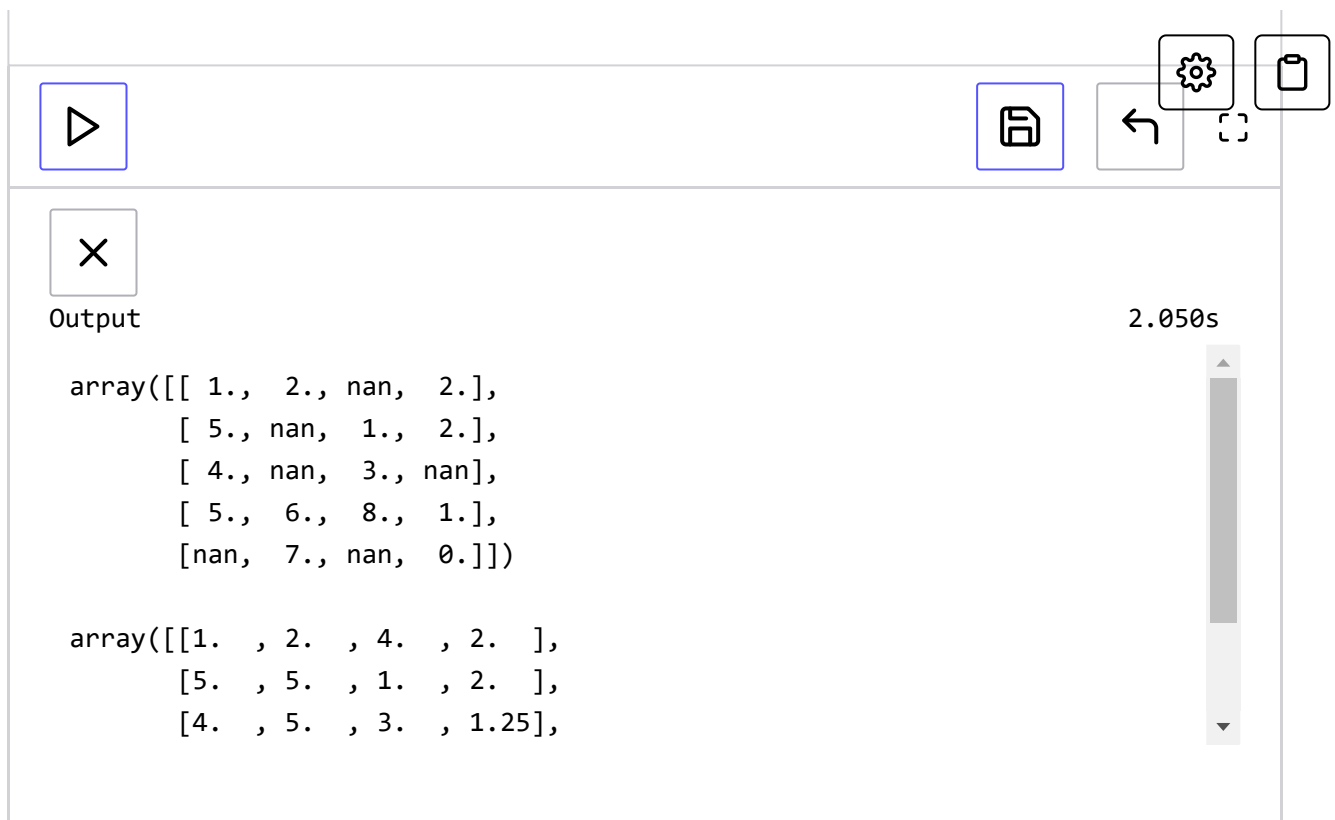
The four methods are:

- Using the mean value
- Using the median value
- Using the most frequent value
- Filling in missing values with a constant

The code below shows how to perform data imputation using mean values from each column.

```
1 # predefined data
2 print('{}\n'.format(repr(data)))
3
4 from sklearn.impute import SimpleImputer
5 imp_mean = SimpleImputer()
6 transformed = imp_mean.fit_transform(data)
7 print('{}\n'.format(repr(transformed)))
```





The image shows a Jupyter Notebook interface. At the top, there is a toolbar with icons for running (a play button), saving (a floppy disk), undo (a curved arrow), redo (a curved arrow), settings (a gear), and a clipboard. Below the toolbar is a code cell with a close button (an 'X' in a square). The code cell contains two NumPy array definitions. The first array is a 5x4 matrix with some missing values (nan). The second array is a 3x4 matrix with numerical values. To the right of the code cell, the word 'Output' is displayed, followed by the execution time '2.050s'. Below this, the output of the code is shown, displaying the two arrays. A vertical scrollbar is on the right side of the output area.

```
array([[ 1.,  2., nan,  2.],
       [ 5., nan,  1.,  2.],
       [ 4., nan,  3., nan],
       [ 5.,  6.,  8.,  1.],
       [nan,  7., nan,  0.]])


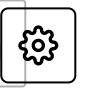
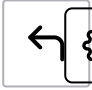


array([[1. , 2. , 4. , 2. ],
       [5. , 5. , 1. , 2. ],
       [4. , 5. , 3. , 1.25],
```


In NumPy arrays, missing data is represented by the `np.nan` value. In the above example, we replaced each missing value with the mean of the values in its column.

The default imputation method for `SimpleImputer` is using the column means. By using the `strategy` keyword argument when initializing a `SimpleImputer` object, we can specify a different imputation method.

The code below demonstrates various initialization strategies for `SimpleImputer`.

```
1 # predefined data
2 print('{}\n'.format(repr(data)))
3
4 from sklearn.impute import SimpleImputer
5 imp_median = SimpleImputer(strategy='median')
6 transformed = imp_median.fit_transform(data)
7 print('{}\n'.format(repr(transformed)))
8
9 imp_frequent = SimpleImputer(strategy='most_frequent')
10 transformed = imp_frequent.fit_transform(data)
11 print('{}\n'.format(repr(transformed)))
```





Output2.096s

```
array([[ 1.,  2., nan,  2.],
       [ 5., nan,  1.,  2.],
       [ 4., nan,  3., nan],
       [ 5.,  6.,  8.,  1.],
       [nan,  7., nan,  0.]])

array([[1. , 2. , 3. , 2. ],
       [5. , 6. , 1. , 2. ],
       [4. , 6. , 3. , 1.5],
```

The 'median' strategy fills in missing data with the median from each column, while the 'most_frequent' strategy uses the value that appears the most for each column.

The final imputation method that `SimpleImputer` provides is to fill in missing values with a specified constant. This can be useful if there is already a suitable substitute for missing data (e.g. 0 or -1).

The code below demonstrates how to fill in missing data with a specific constant. The `fill_value` keyword argument is used when initializing the `SimpleImputer` object, to specify the constant.

```
1 # predefined data
2 print('{}\n'.format(repr(data)))
3
4 from sklearn.impute import SimpleImputer
5 imp_constant = SimpleImputer(strategy='constant',
6                               fill_value=-1)
7 transformed = imp_constant.fit_transform(data)
8 print('{}\n'.format(repr(transformed)))
```

Output

```
array([[ 1.,  2., nan,  2.],
       [ 5., nan,  1.,  2.],
       [ 4., nan,  3., nan],
       [ 5.,  6.,  8.,  1.],
       [nan,  7., nan,  0.]])

array([[ 1.,  2., -1.,  2.],
       [ 5., -1.,  1.,  2.],
       [ 4., -1.,  3., -1.]])
```

1.575



B. Other imputation methods

The `SimpleImputer` object only implements the four imputation methods shown in section A. However, data imputation is not limited to those four methods.

There are also more advanced imputation methods such as k-Nearest Neighbors (https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm) (filling in missing values based on similarity scores from the kNN algorithm) and MICE (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3074241/>) (applying multiple chained imputations, assuming the missing values are randomly distributed across observations).

In most industry cases these advanced methods are not required, since the data is either perfectly cleaned or the missing values are scarce. Nevertheless, the advanced methods could be useful when dealing with open source datasets, since these tend to be more incomplete.

← Back

Next →

Normalizing Data

PCA

☒ Mark as Completed

35% completed, meet the criteria and claim your course certificate!



Report
an Issue



Ask a Question

(https://discuss.educative.io/tag/data-imputation__data-preprocessing-with-scikit-learn__machine-learning-for-software-engineers)