

The challenge of missing or duplicate data

The pros and cons of data outliers

Change categorical data to numerical data

Value: Not numbers versus names

Value: Label encoding in Python

Handling: Other approaches to data visualization

Handling: reference guide: Data cleaning reference

Practice Quiz: Test your knowledge

Changing: original data to cleaned data

Input: validation

Review: Clean your data

Reference guide: Data cleaning in Python

This reference guide contains common functions and methods that data professionals use to clean data. The reference guide contains three different tables of useful tools, each grouped by cleaning category: missing data, outliers, and label encoding.

Save this course item

You may want to save a copy of this guide for future reference. You can use it as a resource for additional practice or in your future professional projects. To access a downloadable version of this course item, click the link below and select "Use Template."

Reference guide: [Data cleaning in Python](#)

OR

If you don't have a Google account, you can download the item directly from the attachment below.

[Reference guide: Data cleaning in Python](#)
DOCX File

Missing data

The following pandas functions and methods are helpful when dealing with missing data.

isnull()

- Description:** A `dataframe` method that returns a concise summary of the dataframe, including a "non-null count", which helps you know the number of missing values.

Example:

```
1 print(df)
2 print()
3 df.isnull()
```

Run
Reset

notnull().isnull()

- Description:** `pd.isna()` is a pandas function that returns a same-sized Boolean array indicating whether each value is null (you can also use `pd.isnull()` as an alias). Note that this function also exists as a `dataframe` method.

Example:

```
1 print(df)
2 print('is After: pd.isnull(): ' + ')')
3
4 pd.isnull(df)
```

Run
Reset

notnull().isnotnull()

- Description:** A pandas function that returns a same-sized Boolean array indicating whether each value is NOT null (you can also use `pd.notnull()` as an alias). Note that this function also exists as a `dataframe` method.

Example:

```
1 print(df)
2 print('is After: notnull(): ' + ')')
3 pd.notnull(df)
```

Run
Reset

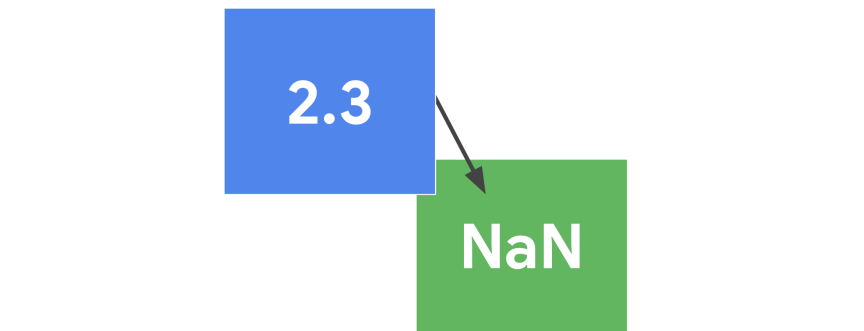
fillna()

- Description:** A `dataframe` method that fills in missing values using specified method.

Example:

```
1 print(df)
2 print('is After: fillna(): ' + ')')
3
4 df.fillna(2)
```

Run
Reset



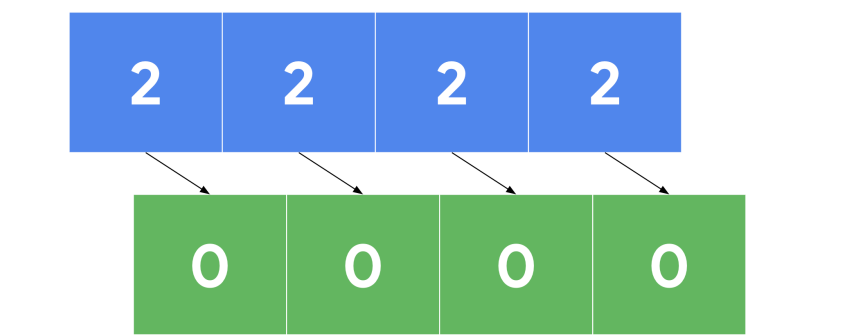
replace()

- Description:** A `dataframe` method that replaces specified values with other specified values. Can also be applied to pandas `Series`.

Example:

```
1 print(df)
2 print('is After: replace(): ' + ')')
3
4 df.replace('A', 'B')
```

Run
Reset



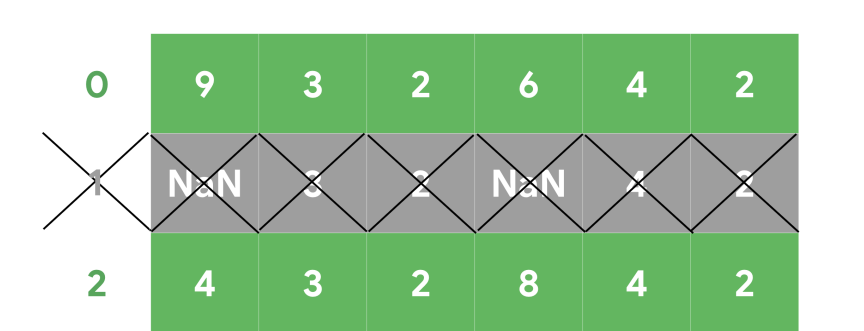
dropna()

- Description:** A `dataframe` method that removes rows or columns that contain missing values, depending on the axis you specify.

Example:

```
1 print('Original df: ' + df)
2 print('is After: dropna(): ' + ')')
3 print(df.dropna(inplace=True))
4
5 print('is After: dropna(inplace=True): ' + ')')
6 print(df.dropna(inplace=True))
```

Run
Reset



Outliers

The following tools are helpful when dealing with outliers in a dataset.

describe()

- Description:** A `dataframe` method that returns general statistics about the dataframe which can help determine outliers.

Example:

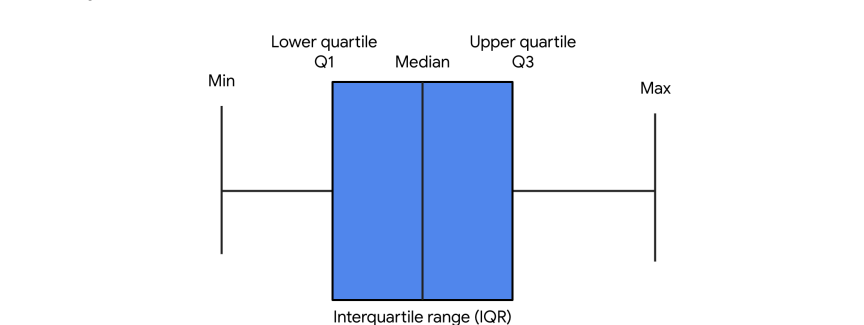
```
1 print(df)
2 print()
3 df.describe()
```

Run
Reset

boxplot()

- Description:** A seaborn function that generates a box plot. Data points beyond 1.5x the interquartile range are considered outliers.

Example:



Label encoding

The following tools are helpful when performing label encoding.

astype()

- Description:** A `dataframe` method that allows you to encode its data as a specified dtype. Note that this method can also be used on `Series` objects.

Example:

```
1 print(df)
2 print('is Original dtype of df: ' + ')')
3
4 print(df.dtypes)
```

```
5 print('is df.dtypes after casting \'class\' column as categorical: ' + ')')
6
7 df['class'] = df['class'].astype('category')
8
9 print(df.dtypes)
```

Run
Reset

Series.cat.codes

- Description:** A `Series` attribute that returns the numeric category codes of the series.

Example:

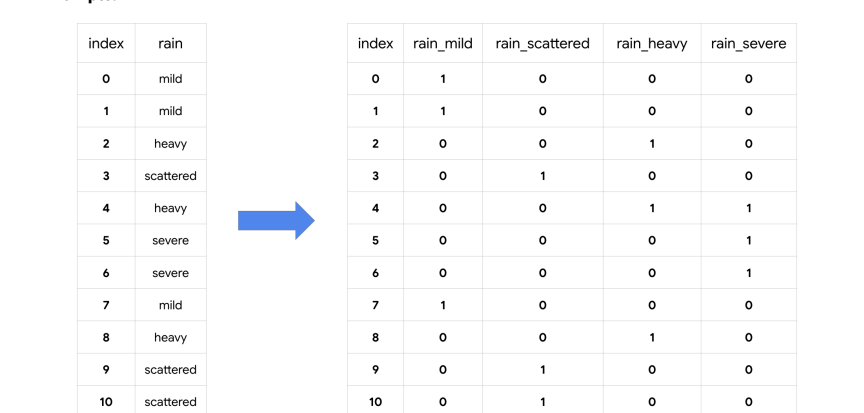
```
1 # Get 'class' column as categorical
2 df['class'] = df['class'].astype('category')
3
4 print('is \'class\' column: ' + ')')
5 print(df['class'])
6
7 print('is Category codes of \'class\' column: ' + ')')
8
9 df['class'].cat.codes
```

Run
Reset

LabelEncoder()

- Description:** A function that converts categorical values into new binary columns—one for each different category.

Example:



LabelEncoder()

- Description:** A transformer from `sklearn.preprocessing` that encodes specified categories or labels with numeric codes. Note that when building predictive models it should only be used on target variables (x, y data).

Example:

```
1 # Import sklearn.preprocessing: LabelEncoder
2
3 # Instantiate LabelEncoder()
4 encoder = LabelEncoder()
5
6 data = ['a', 'b', 'c', 'd']
7
8 # Fit to the data
9 encoder.fit(data)
10
11 # Transform the data
12 transformed = encoder.transform(data)
13
14 # Reverse the transformation
15 inverse = encoder.inverse_transform(transformed)
16
17 print('Data: ' + data)
18 print('is Classes: ' + encoder.classes_)
19 print('is Encoded (transformed): ' + ')')
20 print('is New data: ' + ')')
21 print('is Reverse From encoded classes to original: ' + ')')
22
```

Run
Reset

```
1 # Import sklearn.preprocessing: LabelEncoder
2
3 # Instantiate LabelEncoder()
4 encoder = LabelEncoder()
5
6 data = ['a', 'b', 'c', 'd']
7
8 # Fit to the data
9 encoder.fit(data)
10
11 # Transform the data
12 transformed = encoder.transform(data)
13
14 # New data
15 new_data = ['a', 'b', 'c', 'd']
16
17 # Get classes of new data
18 inverse = encoder.inverse_transform(new_data)
19
20 print('Data: ' + data)
21 print('is Classes: ' + encoder.classes_)
22 print('is Encoded (transformed): ' + ')')
23 print('is New data: ' + ')')
24 print('is Convert new_data to original: ' + ')')
25
```

Run
Reset

Key Takeaways

There are many tools that data professionals can use to perform data cleaning on a wide range of data. The information you learn from missing data, outliers, and transforming categorical to numeric data will help you prepare datasets for further analysis throughout your career.

Mark as completed

[Like](#) [Dislike](#) [Report an issue](#)