

Data Analysis with Python

Cheat Sheet: Data Wrangling

Package/Method	Description	Code Example
Replace missing data with frequency	Replace the missing values of the data set attribute with the mode (common occurring entry) in the column.	<pre>1. pd.isnull(df).sum() 2. df['attribute_name'].fillna(df['attribute_name'].mode()[0])</pre> <div>Copy</div>
Replace missing data with mean	Replace the missing values of the data set attribute with the mean of all the entries in the column.	<pre>1. df['attribute_name'].fillna(df['attribute_name'].mean()) 2. df['attribute_name'].fillna(df['attribute_name'].median())</pre> <div>Copy</div>
Fix the data types	Fix the data types of the columns in the dataframe.	<pre>1. df['attribute_name'].astype('float') 2. df['attribute_name'].astype('int') 3. df['attribute_name'].astype('object')</pre> <div>Copy</div>
Data Normalization	Normalize the data in a column such that the values are restricted between 0 and 1.	<pre>1. df['attribute_name'].min() 2. df['attribute_name'].max() 3. df['attribute_name'].apply(lambda x: (x - df['attribute_name'].min()) / (df['attribute_name'].max() - df['attribute_name'].min()))</pre> <div>Copy</div>
Binning	Create bins of data for better analysis and visualization.	<pre>1. df['attribute_name'].min() 2. df['attribute_name'].max() 3. df['attribute_name'].apply(lambda x: (x - df['attribute_name'].min()) / (df['attribute_name'].max() - df['attribute_name'].min()))</pre> <div>Copy</div>
Change column name	Change the label name of a dataframe column.	<pre>1. df.rename(columns={'old_name': 'new_name'}, inplace=True)</pre> <div>Copy</div>
Indicator Variables	Create indicator variables for categorical data.	<pre>1. df['attribute_name'].value_counts() 2. df['attribute_name'].apply(lambda x: (x == 'category').astype(int))</pre> <div>Copy</div>

