# Data Analysis with Python

## Cheat Sheet: Data Wrangling

| Package/Method | Description | Code Example |
|---|---|---|
| | | ```<br>1. 1<br>2. 2<br>``` |
| Replace missing data with frequency | Replace the missing values of the data set attribute with the mode common occurring entry in the column. | ```<br>1. MostFrequentEntry = df['attribute_name'].value_counts().idxmax()<br>2. df['attribute_name'].replace(np.nan,MostFrequentEntry,>df['attribute_name'].replace(np.nan,MostFrequentEntry, inplace=True)<br>```<br><br>Copied! |
| | | ```<br>1. 1<br>2. 2<br>``` |
| Replace missing data with mean | Replace the missing values of the data set attribute with the mean of all the entries in the column. | ```<br>1. AverageValue=df['attribute_name'].astype(<data_type>).mean(axis=0)<br>2. df['attribute_name'].replace(np.nan, AverageValue, inplace=True)<br>```<br><br>Copied! |
| | | ```<br>1. 1<br>2. 2<br>3. 3<br>``` |
| Fix the data types | Fix the data types of the columns in the dataframe. | ```<br>1. df[['attribute1_name', 'attribute2_name', ...]] =<br>2. df[['attribute1_name', 'attribute2_name', ...]].astype('data_type')<br>3. #data_type is int, float, char, etc.<br>```<br><br>Copied! |
| | | ```<br>1. 1<br>``` |
| Data Normalization | Normalize the data in a column such that the values are restricted between 0 and 1. | ```<br>1. df['attribute_name'] =<br>   df['attribute_name']/df['attribute_name'].max()<br>```<br><br>Copied! |
| | | ```<br>1. 1<br>2. 2<br>3. 3<br>4. 4<br>5. 5<br>6. 6<br>``` |
| Binning | Create bins of data for better analysis and visualization. | ```<br>1. bins = np.linspace(min(df['attribute_name']),<br>2. max(df['attribute_name'],n)<br>3. # n is the number of bins needed<br>4. GroupNames = ['Group1','Group2','Group3',...]<br>5. df['binned_attribute_name'] =<br>6. pd.cut(df['attribute_name'], bins, labels=GroupNames, include_lowest=True)<br>```<br><br>Copied! |
| | | ```<br>1. 1<br>``` |
| Change column name | Change the label name of a dataframe column. | ```<br>1. df.rename(columns={'old_name':\'new_name'}, inplace=True)<br>```<br><br>Copied! |
| | | ```<br>1. 1<br>2. 2<br>``` |
| Indicator Variables | Create indicator variables for categorical data. | ```<br>1. dummy_variable = pd.get_dummies(df['attribute_name'])<br>2. df = pd.concat([df, dummy_variable],axis = 1)<br>```<br><br>Copied! |