

# 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.	<div>1</div> <div>2</div> <ul style="list-style-type: none"><li>• <code>MostFrequentEntry = df['attribute_name'].value_counts().idxmax()</code></li><li>• <code>df['attribute_name'].replace(np.nan, MostFrequentEntry, &gt;df['attribute_name'].replace(np.nan, MostFrequentEntry, inplace=True)</code></li></ul> <div>Copied!</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.	<div>1</div> <div>2</div> <ul style="list-style-type: none"><li>• <code>AverageValue=df['attribute_name'].astype(&lt;data_type&gt;).mean(axis=0)</code></li><li>• <code>df['attribute_name'].replace(np.nan, AverageValue, inplace=True)</code></li></ul> <div>Copied!</div>

Fix the data types	Fix the data types of the columns in the dataframe.	<div>1</div> <div>2</div> <div>3</div> <ul style="list-style-type: none"> <li>• <code>df[['attribute1_name', 'attribute2_name', ...]] =</code></li> <li>• <code>df[['attribute1_name', 'attribute2_name', ...]].astype('data_type')</code></li> <li>• <code>#data_type is int, float, char, etc.</code></li> </ul> <div>Copied!</div>
Data Normalization	Normalize the data in a column such that the values are restricted between 0 and 1.	<div>1</div> <ul style="list-style-type: none"> <li>• <code>df['attribute_name'] =</code></li> <li><code>df['attribute_name']/df['attribute_name'].max()</code></li> </ul> <div>Copied!</div>
Binning	Create bins of data for better analysis and visualization.	<div>1</div> <div>2</div> <div>3</div> <div>4</div> <div>5</div> <div>6</div> <ul style="list-style-type: none"> <li>• <code>bins =</code></li> <li><code>np.linspace(min(df['attribute_name']),</code></li> <li><code>max(df['attribute_name']),n)</code></li> <li>• <code># n is the number of bins needed</code></li> </ul>

		<div> <ul style="list-style-type: none"> <li>• <code>GroupNames =</code>  <code>['Group1','Group2','Group3,...']</code></li> <li>• <code>df['binned_attribute_name'] =</code></li> <li>• <code>pd.cut(df['attribute_name'], bins,</code>  <code>labels=GroupNames, include_lowest=True)</code></li> </ul> </div> <p>Copied!</p>
Change column name	Change the label name of a dataframe column.	<div> 1 <ul style="list-style-type: none"> <li>• <code>df.rename(columns={'old_name':'new_name'}, inplace=True)</code></li> </ul> </div> <p>Copied!</p>
Indicator Variables	Create indicator variables for categorical data.	<div> 1 2 <ul style="list-style-type: none"> <li>• <code>dummy_variable =</code>  <code>pd.get_dummies(df['attribute_name'])</code></li> <li>• <code>df = pd.concat([df, dummy_variable],axis</code>  <code>= 1)</code></li> </ul> </div> <p>Copied!</p>