

DATA CLEANING CHEAT SHEET - 1

Data Quality

- 1) Duplicate Values
- 2) Missing Values
- 3) Invalid Values
- 4) Outliers

Duplicates

1. **df.duplicated()** - returns a Series with True and False values that describe which rows in the DataFrame are duplicated and not
2. **df.drop_duplicates()** - return DataFrame with duplicate rows removed
3. **df.reset_index()** - allows you reset the index back to the default 0, 1, 2 etc indexes
4. **df.drop_duplicates(subset=)** - selecting particular rows and columns of data from a DataFrame (or Series)
5. **df.drop_duplicates(keep=)**
 - first : drop duplicates except for the first occurrence.
 - last : drop duplicates except for the last occurrence.
 - False : drop all duplicates.

```
df.duplicated()
df.duplicated(keep = 'first')
df.duplicated(keep = 'last')

a = df.drop_duplicates()
b = df.drop_duplicates(keep = 'first')
c = df.drop_duplicates(keep = 'last')

a.reset_index()
b.reset_index()
c.reset_index()

a.reset_index(drop = True)
b.reset_index(drop = True)
c.reset_index(drop = True)

df.drop_duplicates(subset = ["name", "place"])
```

Missing Values

df.isnull() - returns a DataFrame object where all the values are replaced with a Boolean value True for NULL values, and otherwise False	<pre>df.isnull() df.isnull().sum() df['Country'].isnull() df['Country'].isnull().sum()</pre>
df.notnull() - pandas function that will examine one or multiple values to validate that they are not null	<pre>df.notnull() df.notnull().sum() df['Country'].notnull() df['Country'].notnull().sum()</pre>
df.notna() - returns a DataFrame object where all the values are replaced with a Boolean value True for NOT NA (not-a-number) values, and otherwise False	<pre>df.notna()</pre>
df.isnull().all() - returns one value for each column, True if ALL values in that column are True, otherwise False	<pre>df.isnull().all(axis = 0) df.isnull().sum().all() df['Country'].isnull().all() df['Country'].isnull().sum().all()</pre>
df.isnull().any() - checks whether any value in the caller object (Dataframe or series) is not 0 and returns True for that	<pre>df.isnull().any(axis = 0) df.isnull().sum().any() df['Country'].isnull().any() df['Country'].isnull().sum().any()</pre>

Treating Missing Values

df.dropna() - used to remove missing values	<pre>df.dropna(how = 'all') df.dropna(how = 'all', subset = ['Region', 'Country']) df.dropna(how = 'any') df.dropna(how = 'any', subset = ['Region', 'Country'])</pre>
df.dropna(thresh=) - takes integer value which tells minimum amount of na values to drop	<pre>df.dropna(thresh = 1)</pre>

Replacing Missing Values

<p>fillna - replaces the NULL values with a specified value</p> <ul style="list-style-type: none">• forwardfill - replaces the NULL values with the value from the previous row (or previous column, if the axis parameter is set to 'columns')• backfill - used to backward fill the missing values in the dataset. It will backward fill the NaN values that are present in the pandas dataframe• value - pass in a value into the value= parameter	<pre>df.fillna(method = 'pad') df.fillna(method = 'bfill') df.fillna(value = 'not sure') df['Order_ID'].fillna(value = '9515842735') df['Order_ID'].fillna(df['Order_ID'].mean())</pre>
<p>replace - replaces the specified value with another specified value</p>	<pre>df['Order_ID'] = df['Order_ID'].replace(np.nan , df['Order_ID'].mean())</pre>
<p>interpolate - used to fill NA values in the dataframe or series. But, this is a very powerful function to fill the missing values</p>	<pre>df['Order_ID'].interpolate() df['Order_ID'].interpolate(method = 'polynomial' , order = 2)</pre>