


Pandas Cheatsheet: Data Cleaning

EDA

Exploratory Data Analysis



<code>df.info()</code>	DataFrame columns, dtypes and memory
<code>df.describe()</code>	Returns columns coverage and types
<code>df.head(7)</code>	Returns first N rows
<code>df.sample(2)</code>	return random samples
<code>df.shape</code>	return DataFrame dimensions
<code>df.columns</code>	returns DataFrame columns

Missing values

Working with missing data



<code>df.isna()</code>	return True or False for missing values
<code>df['col_1'].notna()</code>	return True or False for non-NA data
<code>df.isna().all()</code> <code>s[s == True]</code>	Columns which contains only NaN values
<code>df.isna().any()</code>	Detect columns with NaN values
<code>df['col_1'].fillna(0)</code>	Fill NaN with string or 0
<code>import seaborn as sns</code> <code>sns.heatmap(df.isna(), cmap = 'Greens')</code>	plot missing values
<code>s.loc[0] = None</code> <code>s.loc[0] = np.nan</code>	Insert missing data
<code>df.dropna(axis=0)</code>	dropping rows with missing data
<code>df.dropna(axis=1, how='any')</code>	Drop columns with NaN values

Wrong data


Detect wrong data



<code>df[df['col_1'].str.contains(r'[@#%\$%^&*~!'])]</code>	Detect special symbols
<code>df[df['col_1'].map(lambda x: x.isascii())]</code>	Detect (non) ascii characters
<code>df['col_1'].loc[~df['col_1'].str.match(r'[0-9.]+')]</code>	find pattern with regex
<code>import numpy as np</code> <code>np.where(df['col_1']=='', df['col_2'], df['col_1'])</code>	detect empty spaces
<code>df[df['col_1'].str.contains('[A-Za-z]')]</code>	Detect latin symbols
<code>df.applymap(np.isreal)</code>	detect non numeric rows
<code>df['city'].str.len().value_counts()</code>	Count values by length

Drop


Drop rows, columns, index, condition



<code>df.drop('col_1', axis=1, inplace=True)</code>	Drop one column by name
<code>df.drop(['col1', 'col2'], axis=1)</code>	Drop multiple columns by name
<code>df.dropna(axis=1, how='any')</code>	Drop columns with NaN values
<code>df.drop(0)</code>	Drop rows by index - 0
<code>df.drop([0, 2, 4])</code>	drop multiple rows
<code>df[(df['col1'] > 0) & (df['col2'] != 'open')]</code>	drop rows by condition
<code>df.reset_index()</code>	drop index

Duplicates


Detect and Remove duplicates



<code>df.diet.nunique()</code>	number of unique values in column
<code>df.diet.unique()</code>	unique values in column
<code>df['col_1'].value_counts(dropna=False)</code>	return series of unique values and counts in column
<code>df.duplicated(keep='last')</code>	find duplicates and keep only the last record
<code>df.drop_duplicates(subset=['col_1'])</code>	drop duplicates from column(s)
<code>df[df.duplicated(keep=False)].index</code>	get indexes of all detected duplications:

Replacing


Replace data in DataFrame



<code>df['col'] = df['col'].str.replace(' M', '')</code>	replace string from column
<code>df['col'].str.replace(' M', '')\n.fillna(0).astype(int)</code>	replace and convert column to integer
<code>df['col'].str.replace('A7', '7', regex=False)</code>	Replace values in column - no regex
<code>df.replace(r'[^a-zA-Z]+', '', regex=True)</code>	Find and replace line breaks - new line, tab - regex
<code>df['col'].str.replace('\s+', '', regex=True)</code>	Replace multiple white spaces
<code>df['col'].str.rstrip('\r\n')</code>	Replace line breaks from the right
<code>p = r'<[^>]*>'</code> <code>df['col1'].str.replace(p, '', regex=True)</code>	Replace HTML tags

Outliers


Detect and remove outliers



<code>df['col_1'].describe()</code>	detecting outliers with describe()
<code>import seaborn as sns</code> <code>sns.boxplot(data=df[['col_1', 'col_2']])</code>	detect outliers with boxplot
<code>q_low = df['col_1'].quantile(0.01)</code> <code>q_hi = df['col_1'].quantile(0.99)</code> <code>df[(df['col_1'] < q_low) & (df['col_1'] > q_hi)]</code>	remove outliers with quantiles
<code>import numpy as np</code> <code>ab = np.abs(df['col_1'] - df['col_1'].mean())</code> <code>std = (3*df['col_1'].std())</code> <code>df[ab <= std]</code>	remove outliers with standart deviation

Wrong format


Detect wrong format



<code>df.apply(pd.to_numeric, errors='coerce')\n.isna().any()</code>	detect wrong numeric format
<code>pd.to_datetime(df['date_col'], errors='coerce')</code>	detect wrong datetime format
<code>import pandas_dedupe</code> <code>dd_df = pandas_dedupe.dedupe_dataframe(df, field_properties=['col1', 'col2'], canonicalize=['col1'], sample_size=0.8)</code>	Find typos and misspelling - Deduplication and canonicalization with pandas_dedupe
<code>from difflib import get_close_matches</code> <code>w = ['apes', 'apple', 'peach', 'puppy']</code> <code>get_close_matches('ape', w, n=3, cutoff=0.8)</code>	Use difflib to find close matches

Fix errors

Fix errors in Pandas



<code>df.convert_dtypes()</code>	Convert the DataFrame to use best possible dtypes
<code>df.astype({'col_1': 'int32'})</code>	Cast col_1 to int32 using a dictionary
<code>df.fillna(method='ffill')</code>	Propagate non-null values forward or backward
<code>values = {'A': 0, 'B': 1, 'C': 2, 'D': 3}</code> <code>df.fillna(value=values)</code>	Replace all NaN elements in column with dict

