

Pandas Cheat Sheet :

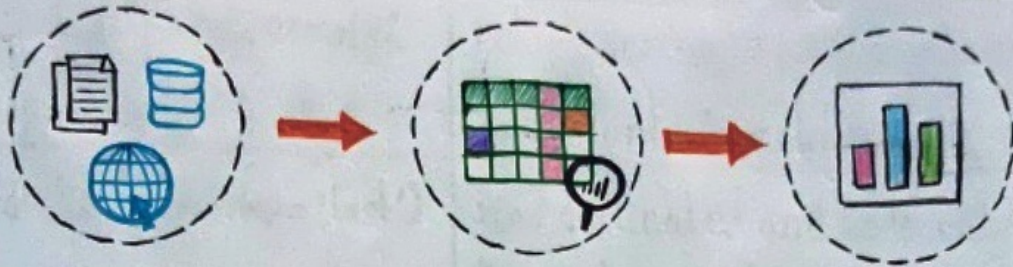
Data Cleaning

Data Cleaning useful for everyday working with data. This Pandas cheat sheet contains ready-to-use codes and steps for data cleaning.



EDA

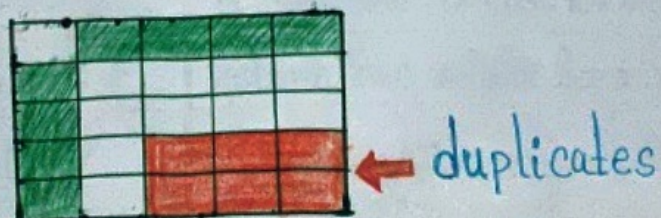
Exploratory Data Analysis.



<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

Duplicates

Detect and Remove duplicates.

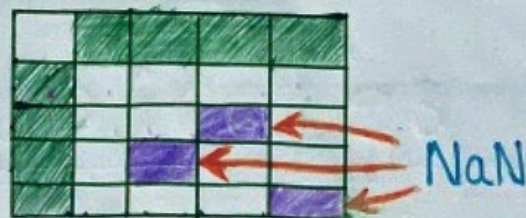


Duplicates

<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:

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() s[s == True]</code>	Columns which contains only NaN values

Missing 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

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'].quantile(0.01)</code> <code>q_low = df['col'].quantile(0.99)</code> <code>df[(df['col'] < q_hi) & (df['col'] > q_low)]</code>	remove outliers with quantiles

Outliers

```
import numpy as np
```

```
ab = np.abs(df['col'] - df['col'].mean())
```

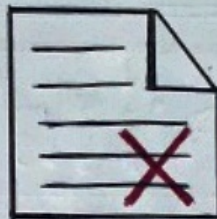
```
std = (3 * df['col'].std())
```

```
df[ab <= std]
```

remove outliers with
standant deviation

Wrong Data

Detect wrong data.



```
df[df['col_1'].str.contains(r'[@#&$%+-/!'])]
```

Detect special symbols

```
df[df['col_1'].map(lambda x: x.isascii())]
```

Detect (non) ascii characters

```
df['col'] \
```

```
.loc[~df['col'].str.match(r'[0-9.]+')]
```

find pattern with regex

```
import numpy as np
```

```
np.where(df['col'] == ' ', df['col2'], df['col'])
```

detect empty spaces

```
df[df['col_1'].str.contains('[A-Za-z]')]
```

Detect latin symbols

```
df.applymap(np.isreal)
```

detect non numeric rows

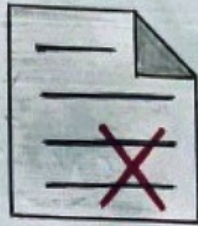
```
df['city'].str.len().value_counts()
```

Count values by lenght

Wrong Format

Detect wrong format.

numeric?
date?



```
df.apply(pd.to_numeric, errors='coerce')  
.isna().any()
```

detect wrong numeric
format

```
pd.to_datetime(df['date_col'],  
errors='coerce')
```

detect wrong datetime
format

```
import pandas_dedupe  
dd_df = pandas_dedupe.dedupe_dataframe(df, field_properties=['col1',  
'col2'], canonicalize=['col1'], sample_size=0.8)
```

Find typos and miss-
pelling - Deduplication
and canonicalization
with pandas_dedupe

```
from difflib import get_close_matches  
w = ['apes', 'apple', 'peach', 'puppy']  
get_close_matches('ape', w, n=3,  
cutoff=0.8)
```

use difflib to find
close matches

Fix Errors

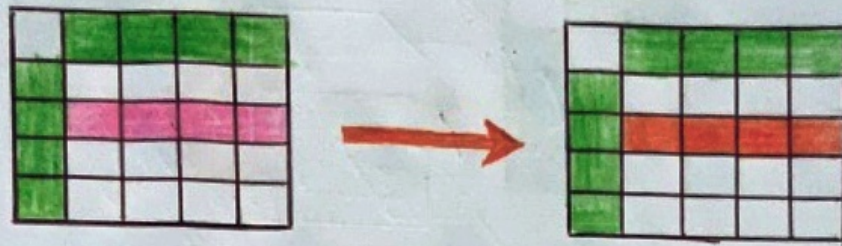
Fix errors in Pandas.



<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
<code>df.fill(axis=0)</code>	Fill the missing values row wise

Replacing

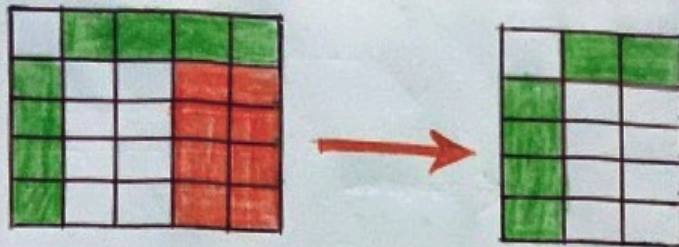
Replace data in DataFrame.



<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('AT', 'T', regex = False)</code>	Replace values in column - no regex
<code>df.replace(r'\n+ \n+ \t+', '', 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('\n\n')</code>	Replace line breaks from the right
<code>p = r'<[<>]*>'</code> <code>df['col'].str.replace(p, '', regex = True)</code>	Replace HTML tags

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 column 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

