

PANDAS CHEAT SHEET

ABBREVIATIONS

nv = new variable	nv_df = new dataframe variable
v = variable	df = dataframe variable

READ & WRITE FILES

function	command	explanation, comments, examples
load file	<code>nv = "folder_name/file_name"</code>	Can use indirect path as well (<code>os.path.join()</code>).
read csv file	<code>nv_df = pd.read_csv(v)</code>	When file is read it is also converted into a dataframe.
write csv file	<code>df.to_csv("new_file_name.csv", index=False)</code>	<code>index=False</code> , exporting without index.
write xlsx file	<code>df.to_excel("new_file_name.xlsx", index=False)</code>	

ORGANIZING DATA

function	command	explanation, comments, examples
rename columns	<code>nv_df = df.rename(columns={"original1": "new1", "original2": "new2"})</code>	
change column order	<code>nv_df = df[["column1", "column2", "column3"]]</code>	It can be done with a passing a variable too: <code>nv = ["column1", "column2", "column3"]</code> <code>df = df[nv]</code>
adding a new column	<code>df["new_column"] = nv</code>	<code>nv</code> = condition that we pass in the new variable, for example: <code>membership_weeks = training_data["Membership(Days)"]/7</code> <code>training_data["Membership(Weeks)"] = membership_weeks</code> <code>training_data.head()</code>
sorting data: descending order	<code>nv_df = df.sort_values(["column"], ascending=False)</code>	
sorting data: ascending order	<code>nv_df = df.sort_values(["column"], ascending=True)</code>	
reset index	<code>nv_df = df.reset_index(drop=True)</code>	After sorting, the indexes are not in order. If we want to have them in order we use <code>reset_index()</code> .

INSPECTING DATA		
function	command	explanation, comments, examples
general information	<code>nv = df.info()</code>	nv = optional
count rows in each column	<code>nv = df.count()</code>	Get a count of the rows for all and each column containing data ("null", NaNs, are not included).
count rows in a specific column	<code>nv = df["column"].count()</code>	Count all rows in given column.
unique values in a column	<code>nv = df["column"].value_counts()</code>	How many times unique rows appear in a column. Note: see <code>set_index()</code> method for comparison.
list all columns	<code>df.columns</code>	Returns columns in a list.
unique values	<code>nv = df["column"].unique()</code>	Returns a list of unique values. Use <code>len(df["column"].unique())</code> to count = > nv = <code>len(df["column"].unique())</code>
missing values	<code>df.isnull()</code>	Determine empty rows, returns Boolean.
sum of missing values	<code>df.isnull().sum()</code>	Determine empty rows and return a sum of a missing values.
not-missing values	<code>df.notnull().sum()</code>	Determine not-empty rows and return a sum of a not-missing values.
describe	<code>df.describe()</code>	Returns measures of central tendency.
minimum/maximum	<code>df.min() / df.max()</code>	Returns minimum / maximum
columns	<code>df.columns</code>	Returns name of columns in a list.
data type	<code>df.dtypes</code>	Determine datatype.
data ytppe	<code>df.column.dtype</code>	determine datatype for specific column.
loc/filtering data	<code>nv_df = df.loc[rows, ["column1", "column2"]]</code>	<code>nv_df = df.loc[:, ["column1", "column2"]]</code> = selecting all rows
loc with conditional and selected columns after comma	<code>nv_df = df.loc[(df["column"] > 7), ["column1", "column2"]]</code>	
loc with conditional operators	<code>nv_df = df.loc[(df["column"] == "row/value") & (df["column"] == "row/value") (df["column"] > n)]</code>	
SET NEW INDEX		
function	command	explanation, example
set new index	<code>nv = df.set_index(["column_to_be_new_index"])["column"]</code>	set new index for given column and get a series with a specific column.
Remove the index name.	<code>df.index.name = None</code>	

CLEANING DATA		
function	command	explanation, comments, examples
replace exact given string into new one	<code>df["column"] = df["column"].str.replace("existing", "replacement")</code>	for word in prefixes_suffices: <code>student_data_df["student_name"] = student_data_df["student_name"].str.replace(word, "")</code> (.str = datatype change. We need to convert object to a string before using replace() method!)
replace/rename row names to a new one	<code>df = df.replace({"old_row_name": "new_row_name"})</code>	Use when rows has similar names and can be combined into one.
drop	<code>nv_df = df.dropna(how = 'any')</code>	Drops the row with NaNs.
fill	<code>nv_df = df.fillna(value)</code>	To fill the row with certain value.
to tilst	<code>nv = df["column"].tolist()</code>	Add column to a new list.
split	<code>split()</code>	split string object on the whitespace, or where there is no text.
	<code>set()</code>	The set() method returns all unique items in a list when that list is added inside the parentheses
	<code>strip()</code>	The strip() method removes any combination of letters and words that are inside the parentheses.
delete a column	<code>del df['column']</code>	
MERGE FILES		
function	command	explanation, comments, examples
merge files	<code>nv_df = pd.merge(first_df, second_df, on=["column", "column"])</code>	Note: best if combine on a same column with the same name (rename if neccessary).
merge files and matematical operators		Note: when working with merged dataframe, caution how we use sum() When we apply the sum() method on the <code>merged_df["column"]</code> we get combined output, since rows are repeting. Instead use <code>sum() on single_df["column"]</code> .
setting sufices with merge files	<code>nv_df = pd.merge(first_df, second_df, on=["column", "column"], suffixes = ("sufix1", "sufix2"))</code>	
DISPLAYING DATA		
function	command	explanation, comments, examples
first 5	<code>df.head()</code>	
last 5	<code>df.tail()</code>	
first n	<code>df.head(n)</code>	
last n	<code>df.head(n)</code>	
list desired columns	<code>nv = pd[["column1" "column2", "column3"]]</code>	note: you can do it with loc too

CALCULATIONS		
function	command	explanation, comments, examples
sum	<code>nv = df["column"].sum()</code>	calculates sum for specific column
horizontal sum	<code>df["new_column"] = df.sum(axis = 1)</code>	horizontally sum up values (same row, different columns) + adding results in a new column with <code>df["new_column"]</code>
other calculations	<code>nv = df["column"].mean() .max(), min(), etc,...</code>	<code>average_reading_score = school_data_complete_df["reading_score"].mean()</code> calculates average for specific column.
operations between columns + creating a new column	<code>nv = df["column1"] / df["column2"] df["new_column"] = nv</code>	
matematical operation for each row in selected column	<code>nv = df["column"] matematical_operator n</code>	<code>nv = df["column"] / n</code>
CONDITIONALS		
function	command	explanation, comments, examples
conditionals	<code>nv = df["column"] >= n</code>	returns boolean
conditionals cont.	<code>nv = df[df["column"] >= n]</code>	returns values (data/dataframe) with this condition
combine conditionals	<code>nv = df[(df["column"] >= n) & (df["column"] >= n)]</code>	
DATA FRAMES		
function	command	explanation, comments, examples
DATA FRAME : list of dictionaries	<code>df = pd.DataFrame({"column1": value1, "column2": value2, "column3":value3}))</code>	
DATA FRAME : list of dictionaries	<code>df = pd.DataFrame([{"column1" : "value1", "column2": "value2" } {"column1" : "value1.1", "column2": "value2.2" }])</code>	
DATA FRAME : dictionary of lists	<code>df = pd.DataFrame({ "column1": ["value1", "value2", "value3"], "column2": ["value1", "value2", "value3"] })</code>	
DATA SERIES	<code>nv = pd.Series(v)</code>	
FORMATTING		
function	command	explanation, comments, examples
map & formatting & changing data type	<code>df["column"] = df["column"].astype(type).map("{:,0f} ".format)</code>	<code>strongest_df["HP"] = strongest_df["HP"].map("{:2f} millions".format)</code>
new column order	<code>df = df[["column1", "column2", "column3"]]</code>	Note: also with passing a variable: Assign district summary df the new column order. <code>district_summary_df = district_summary_df[new_column_order]</code>

CHANGE DATA TYPE		
function	command	explanation, comments, examples
change data type	<code>df["column"] = df["column"].loc[:, "column"].astype(float)</code>	converted_ufo["duration (seconds)"] = converted_ufo.loc[:, "duration (seconds)"].astype(float)
CREATING BINS & GROUPING		
function	command	explanation, comments, examples
CREATING BINS (1)creating bins: range and labels (2)cutting (3)cutting with creating a new row	(1a)bins_variable=[1st_point(i), 2nd_point(ni),3rd_point, etc] (1b)labels_variable=["", "", "", ""] one less than element then in bins (2)pd.cut(df["column"], bins_variable, labels=labels_variable) (3)df["new_column"] = pd.cut(ted_talks_df["views"], views_bins, labels=views_labels)	creating bins/categories, based on selected column. (1a)spending_bins = [0, 585, 630, 645, 675] (1b)group_names = ["<\$584", "\$585-629", "\$630-644", "\$645-675"]
GROUPING (1)create a group based on a column	<code>nv = df.groupby(["column_to_be_grouped_by"]).matemathical_operator()["column_matemathical_operator_apply_to"]</code>	<code>nv = df.groupby(["column"]).mean()["column1"]</code> OR w/o matemathical operator: <code>nv = df.groupby("column") => use this when establishing bins and you want to check how many elements in each bin. print(ted_group["views"].count())</code>
combinations	<code>nv_df = df.groupby(["column_to_be_grouped_by"]).matemathical_operator()["column(s)_matemathical_operator_apply_to"]</code>	<code>ted_group = ted_talks_df.groupby(["View Group"]).mean()[["comments", "duration", "languages"]]</code> <code>ted_group</code> *NOTE: if you skip columns the matemathical operator will apply to all.