📖 **MAZ-CAS-DDJ** / **kurs_18_19**

Branch: master ▾    **kurs_18_19** / **08 Pandas Teil 1** / **material** / **pandas.md**        [ Find file ]  [ Copy path ]

👤 **Simon Schmid** Update 24.Sep                                                    45f51df 9 days ago

**0** contributors

---

358 lines (218 sloc)    24.5 KB

# PANDAS CHEAT SHEET

by Simon Schmid. Work in progrss, without any guarantees. Spotted a mistake? Mail **here**.

## BASICS

### Install

`pip install pandas`  ✎õñļẏ dõ țḥìṡ õñćè

### Libraries to import

`import pandas as pd`  ✎baṡìć qaćkağè

`import numpy as np`  ✎ñèèd țḥìṡ fõṝ NaN

`import datetime`  ✎țõ dèaḷ wìțḥ țìṃè

### Help

`?`  ✎ğèț ḥèḷq õñ a ṃèțḥõd ¤è»ğ»țẏqè qd»ḥèad@/

## DATA IN AND OUT

### Constructors

`pd.Series(list)`  ✎ćõñṡțṝúćț a Sèṝìèṡ fṝõṃ a ḷìṡț (reference)

- list: ['value1', 'value2']

`s = pd.Series(dict)`  ✎ćõñṡțṝúćț a Sèṝìèṡ fṝõṃ a dìćț (reference)

- dict: ({"key1": value1, "key2": value2})

`df = pd.DataFrame(listoftdicts)`  ✎Cṝèațè a DF fṝõṃ a ḷìṡț õf dìćțìõñaṝìèṡ (reference)

- listofdicts: [{'field1': value1, 'field2': value2}, {'field1': value3, 'field2': value4}]

`df = pd.DataFrame(dictoflists)`  ✎Cṝèațè a DF fṝõṃ a dìćțìõñaṝẏ õf ḷìṡțṡ (reference)

- dictsoflists: {'field1': [value1, value3], 'field2': [value2, value4}

### Data in

`df = pd.read_csv("file.csv")`  ✎Cṝèațè a DF fṝõṃ a CSV fìḷè (reference)

- nrows=59  ✎țõ õñḷẏ ğèț a ñúṃbèṝ õf ṝõẇṡ
- na_values=["string1", "string2", ...]  ✎țõ ṡqèćìfẏ waḷúèṡ țõ ìğñõṝè
- dtype=str  ✎țṝèaț èwèṝẏțḥìñğ aṡ ṡțṝìñğṡ
- sep="\t"  ✎fõṝ țab dèḷìṃìțèd fìḷèṡ
- error_bad_lines=False  ✎ìğñõṝè èṝṝõṝṡ ¤è»țõõ ṃañẏ fìèḷdṡ õñ a ḷìñè/

- `header=None` 🔸 úśè ìf ćśw̌haś ñõ h̃eadèr̃ r̃õw̌
- `names=["id", "cat"]` 🔸 śqèćìf́ý ćõḷúm̃ñ ñam̃èś t̃õ b̲è úśèd

**`df = pd.read_excel("file.xlsx")`** 🔸 C̃r̃èaţ̀è a DF f̃r̃õm̲ añ XLS f̃ìḷè ([reference](#))

- `sheetname="name"` 🔸 t̃h̀è ñam̃è õf̃ t̃h̀è śh̀èèţ ìñśìdè èyćèḷ
- `skiprows=n` 🔸 śk̀ìq t̃h̀è f̃ìr̃ṣt ñ r̃õw̌ś
- `names=list` 🔸 úśè t̃h̀èśè ćõḷúm̃ñ ñam̃èś

**`pd.read_sql(query, conn)`** 🔸 èyèćúţ̀è a SQL🔸 úèr̃ý õñ a g̃ìwèñ ćõññèćţ̀ìõñ ([reference](#))

- `index_col="column"` 🔸 t̃h̀è ćõḷúm̃ñ t̃õ b̲è úśèd aś ìñdèy

## Data out

**`df.to_csv("file.csv")`** 🔸 śawè a DF ìñţ̃õ t̃h̀è śqèćìf̃ìèd ćśw([reference](#))

- `index=False` 🔸 d̃õñ.ţ śawè t̃h̀è ìñdèy ćõḷúm̃ñ

**`df.to_dict()`** 🔸 śawè daţaf̃r̃am̃è aś dìćţ̀ìõñar̃ỳ ([reference](#))

- `orient="records"/"list"/...` 🔸 waỳ õf̃ ćõñśţ̀r̃úćţ̀ìñg̃ t̃h̀è dìćţ̀ìõñar̃ỳ

**`df.to_json()`** 🔸 śawè aś k̃śõñ🔸 śţ̀r̃ìñg̃ ⊷aḷm̃õśţ ḷìk̀è t̃õᵃdìćţ/ ([reference](#))

- `orient="records"/"index"` 🔸 õr̃ìèñţaţ̀ìõñ õf̃ t̃h̀è k̃śõñ

**`s.to_frame()`** 🔸 ćõñ̃wèr̃ţ́ś a śèr̃ìèś ìñţ̃õ a daţaf̃r̃am̃è ([reference](#))

# DATA TYPES

## NaN

**`NaN`** 🔸 P̣ḷaćèh̃õḷdèr̃ f̃õr̃ m̀ìśśìñg̃ daţa

**`np.nan`** 🔸 ćõdè f̃õr̃ ñañ ⊷ñèèd t̃õ ìṃqõr̃ţ ñúṃqý aś ñq/

# SELECTING STUFF

## Select whole table

**`df`** 🔸 r̃èqr̃èśèñţ́ś w̃h̃õḷè ţabḷè

## Select columns

**`df.field1`** 🔸 f̃èţ́ćh̃ õñḷỳ õñè ćõḷúm̃ñ

**`df["field1"]`** 🔸 aḷţ̀èr̃ñaţ̀ìwè ñõţaţ̀ìõñ

**`df[["field1", "field2"]]`** 🔸 f̃èţ́ćh̃ śèwèr̃aḷ ćõḷúm̃ñś

**`df[condition]`** 🔸 õñḷỳ f̃èţ́ćh̃ r̃õw̌ś w̃h̀èr̃è ćõñdìţ̀ìõñ ìś ţ̃r̃úè

- `condition = df['field' == value]` 🔸 añ èyaṃq̣ḷè

## Select rows

**`df.head(n)`** 🔸 õñḷỳ f̃ìr̃ṣt ñ r̃õw̌ś ([reference](#))

**`df.tail(n)`** 🔸 õñḷỳ ḷaśṭ ñ r̃õw̌ś ([reference](#))

**`df.loc[index]`** 🔸 g̃èţ r̃õw̌ aţ qar̃ţ̀ìćúḷar̃ ìñdèy ([reference](#))

**`df.iloc[integer]`** 🔸 ţ̃r̃èaţ ìñdèy aś ìf̃ ìţ w̌aś a r̃añg̃è õf̃ ìñţ̀èg̃èr̃ś ([reference](#))

## Group data fields

`df.groupby("field1")` »» ţõ ìñìţìäļìżè ğŗõúqèd õúţqúţ »» ñèèd ñèļd ] ağğŗèğäţè ƒúñćţìõñ ([reference](reference))

- use like: `df.groupby("field1")["field2"].function()`

`df.groupby(["field1", "field2"])` »» ğŗõúqbý õñ ţwõ ļèwèļś

## DESCRIBE AND SUMMARIZE

### Properties of the dataframes

`df.index` »» A ļìśţ õƒ ŗõw ìñdìćèś úśúaļļý ñúmbèŗś/ ([reference](reference))

`df.columns` »» A ļìśţ õƒ ćõļúmñ ñamèś ([reference](reference))

`df.dtypes` »» A ļìśţ õƒ ćõļúmñ daţaţýqèś ([reference](reference))

`df.shape` »» A ţúqļè śqèćìƒýìñğ ŗõwś«ćõļúmñś/ ([reference](reference))

`df.values` »» A maţŗìy õƒ ţħè ţabļè wìţħõúţ ħèadèŗś añd ŗõw ìñdìćèś ([reference](reference))

### Aggregate Statistics

`.max()` »» mayìúmúm ([reference](reference))

- `axis=1` »» Úśè aļļ ţħèśè ƒúñćţìõñś ñõţ ćõļúmñ wìśè búţ ŗõw wìśè

`.min()` »» mìñìmúm ([reference](reference))

`.mean()` »» mèañ ([reference](reference))

`.std()` »» śţañdaŗd dèwìaţìõñ ([reference](reference))

`.sum()` »» śúm ([reference](reference))

`.count()` »» ćõúñţ ([reference](reference))

`.size()` »» úśèƒúļ ƒõŗ dõúbļè ğŗõúqbýś«śìmìļaŗ ţõ waļúèªćõúñţś¤ ([reference](reference))

## MODIFY DATAFRAMES

`df.copy()` »» ćõqý ţħè daţaƒŗamè ìñśţèad õƒ ķúśţ ŗèƒèŗèñćìñğ ìţ/ ([reference](reference))

### Modify the index

`df.set_index("field")` »» ćħañğè ţħè ìñdèy ćõļúmñ ţõ : ñèļdº ([reference](reference))

- `inplace=True` »» makè ţħè ćħañğèś õñ ţħè õbķèćţ«ñõţ a ćõqý

`df.rename_axis("Name")` »» Rèñamè ţħè ìñdèy ćõļúmñ»Úśè .Nõñè..ţõ dèļèţè ţħè ìñdèy ñamè ([reference](reference))

- `inplace=True` »» makè ţħè ćħañğèś õñ ţħè õbķèćţ«ñõţ a ćõqý

### Modify columns

`df.insert(pos, "field1", values)` »» ìñśèŗţ ñèw ćõļúmñ : ñèļdº aţ qõś wìţħ waļúèś ([reference](reference))

`df.pop("field1")` »» dèļèţè ćõļúmñ : ñèļdº ([reference](reference))

`df.assign(newfield = df["field1"] ... )` »» aśśìğñ waļúèś ţõ ñèw ćõļúmñ ¤õŗìğìñaļ ŗèmaìñś/ ([reference](reference))

`df.rename(columns=dict)` »» ŗèñamè ćõļúmñś«úśìñğ \.õļdº.».ñèw9«.õļd4.».ñèw4.·· aś dìćţ ([reference](reference))

- `inplace=True` »» makè ţħè ćħañğèś õñ ţħè õbķèćţ«ñõţ a ćõqý

### Modify rows

`df.append(series/dataframe)` »» addś ţħè ŗõw«ŗèţúŗñś ñèw õbķèćţ ([reference](reference))

**df.drop(df[condition].index)** `#` delete rows from table based on condition ([reference](#))

- **inplace=True** `#` make the changes on the object«not a copy

## Modify data structure

**df.pivot()** `#` transform long data into wide data ([reference](#))

- **index='field1'** `#` the column to be used as the new index column
- **columns='field2'** `#` the column which has the values that will make up the new columns
- **values='field3'** `#` the column containing the values

**df.melt()** `#` transform wide data into long data ([reference](#))

**df.unstack()** `#` transform groupby subrows into columns ¤useful to chart stacked bars/ ([reference](#))

**df.transpose()** `#` switches row and columns over the whole dataset ([reference](#))

**df.T** `#` shorthand for d»transposée¤

## Sort data globally

**df.sort_index()** `#` sort not by values but by index ([reference](#))

**df.sort_values("field1")** `#` sort values ([reference](#))

- **ascending=False** `#` in descending order
- **na_position="first"/"last"** `#` position of NaN values

# COMBINE DATAFRAMES

**df.merge(df2)** `#` merge dataframe with other dataframe ([reference](#))

- **on="field"** `#` fieldname«s/ to match »if they have same name/
- **left_on="df1-field"** `#` fieldname to match on left side
- **right_on="df2-field"** `#` fieldname to match on right side
- **left_index=True** `#` whether to use the index as the left«side match field
- **right_index=True** `#` whether to use the index as the left«side match field
- **how="inner/left/right/outer"** `#` just like in SQL

**df.join(df2)** `#` join a dataframe ¤with identical number of rows/ to another«horizontally ([reference](#))

**pd.concat([df1,df2])** `#` adds all the dataframes in the list«vertically ([reference](#))

- **axis=1** `#` add horizontally«not vertically
- **ignore_index=True** `#` construct new index«don.t use existing one

# MODIFY DATA GLOBALLY

## Deal with NaNs

**pd.isnull()** `#` Built«in function to test for null on any value ([reference](#))

**pd.notnull()** `#` same but opposite ([reference](#))

**df.dropna()** `#` get rid of NaNs«optional>in a subset ([reference](#))

- **subset="field1"** `#` only apply on subset
- **inplace=true** `#` make the changes on the object«not a copy
- **how="all"** `#` only drop rows where all fields are NaN

**df.fillna(value)** `#` replace NaN.s with other value ([reference](#))

- **inplace=true** `#` make the changes on the object«not a copy

## Deal with duplicates

**df.drop_duplicates()**  💬gèțś řìd õf dúqļìćațè wạļúèś (reference)

- subset="field"  💬õñļỳ ćõñsìdèř ćèřțaìñ fièḷdś xõř ļìśț õf fièḷdś/
- keep="first/last/False"  💬w̃hìćh õf țhè wạļúèś țõ ǩèèq
- inplace=True  💬maǩè țhè ćhañğèś õñ țhè õbḷèćț«ñõț a ćõqỳ

**df.duplicated()**  💬Țhè dúqḷìćațèś ḭñwèřsè õf dřõqªdúqćḷìćațèdᵷ/ (reference)

- subset="field"  💬õñḷỳ ćõñsìdèř ćèřțaìñ fièḷdś xõř ļìśț õf fièḷdś/
- keep="first/last/False"  💬w̃hìćh õf țhè wạļúèś țõ ǩèèq

## Various

**df.round({'field1': n, 'field2': m})**  💬Țõúñd țhè ñúṃbèřś ìñ qařțìćúḷař ćõḷúṃñś (reference)

**df.dot(df2)**  💬dõț qřõdúćț õf țwõ daṭafřaṃèś (reference)

**df.update(df2)**  💬Uqdaṭè wạḷúèś ìñ df wìțh ñõ💬Nañ wạḷúèś fřõṃ df4 (reference)

# DEAL WITH INDIVIDUAL DATA FIELDS (I.E. SERIES)

Many of these functions can be used on whole dataframes as well.

## Filter fields

**df["field1"].isnull()**  💬řèțúřñś țřúè ìf ñúḷḷ (reference)

**df["field1"].notnull()**  💬õqqõsìțè (reference)

**df["field1"].isin(["str1", "str2"])**  💬řèțúřñś țřúè ìf fièḷd9èř úaḷś a wạḷúè ìñ țhè ḷìśț (reference)

**~df["field1"].isin(["str1", "str2"])**  💬řèțúřñś țřúè ìf fièḷd9dõèśñ.ț èř úaḷ a wạḷúè ìñ țhè ḷìśț

**df["field1"].str.contains("str")**  💬řèțúřñś țřúè ìf fièḷd9ćõñțaìñś țhè śțřìñğ śțř (reference)

- na=False  💬Awõìd Ęřřõř fõř NaN wạḷúèś
- regex=True/False  💬bỳ dèfaúḷț«řèğèy ćañ bè ìñćḷúdèd
- case=True/False  💬ćaśè śèñsìțìwè õř ñõț«dèfaúḷț Țřúè

## Aggregate summaries over fields

**df["field1"].describe()**  💬dìśqḷaỳś may«ṃìñ«ṃèañ«èțć (reference)

**df["field1"].max()**  💬ćaḷćúḷațè may«èțć»aḷśõ>ṃèañ«ṃìñ«»»(reference)

**df["field"].value_counts()**  💬řèř úèñćỳ õf èaćh wạḷúè«ìñ țabúḷař fõřṃ (reference)

- normalize=True  💬ìñ qèřćèñțağèś
- dropna=False  💬ìñćḷúdè NaN.ś
- ascending=True  💬Sõřț ìñ ìñwèřsè õřdèř

**df["field"].unique()**  💬ğèț a ḷìśț õf úñìŗ úè xdìśțìñćț/ wạḷúèś (reference)

**pd.get_dummies(df["field"])**  💬Baśèd õñ úñìŗ úè wạḷúèś ìñ a fièḷd«ćřèațè a śèț õf dúṃṃỳ ćõḷúṃñś (reference)

- prefix="prefix"  💬Přèfìy țõ úśè bèfõřè úsìñğ úñìŗ úè wạḷúèś aś ćõḷúṃñ ḣèadèřś
- drop_first=True

## Mathematical modifications

**df.rolling(n, on="column")**  💬řèțúřñś țhè řõḷḷìñğ awèřağè õf : ćõḷúṃñ: aś a DF (reference)

- min_periods=n  💬sèț ñúṃbèř õf qèřìõdś țõ awèřağè õwèř

**df['field1'].pct_change()**  💬ćaḷćúḷaṭèś } 💬ćhañğè bèțwèèñ qèřìõd ț añd ț· 9xõñ śèřìèś õř df/ (reference)

`df['field1"].agg(['func1', 'func2'])` ❧aqqlièś ağğřèğaṭè fúñćṭiõñ lĭkè .mèañ..èṭć»ṭõ ćõlúmñ añd śqìṭś õúṭ a daṭafřamè (reference)

## Data modifications

`df["field1"].astype(int)` ❧ćõñwèřṭ ṭõ a ṭýqè ừñṭ«śṭř«flõaṭ/ (reference)

- `errors="ignore"` ❧ĭğñõřè èřřõřś

`df["field1"].replace("str1", "str2")` ❧řèqlaćè ḻĭkè èyćèḻ«wĥõḻè ćèḻḻ/ ưśè śṭř»řèqlaćé fõř qařṭś ìñśìdè ṭĥè śṭřìñğ/ (reference)

- `dictionary` ❧úśèś kèy»waḻúè qaìřś ìñ ṭĥè dìćṭìõñařý ṭõ řèqlaćè múḻṭiqlè waḻúèś
- `regex=True` ❧úśè řèğèy ṭõ fìñd ìñśṭañćèś õf śṭř9

`df["field1"].extract(regex)` ❧èyṭřaćṭś a řèğèy fřõm ṭĥè fìèḻd (reference)

- `expand=True` ❧fõřćè řèṭúřñ õf a daṭafřamè ìñśṭèad õf a śèřìèś
- `.dropna()` ❧ṭõ dřõq ṭĥè NA waḻúèś

## Assign field values dynamically

`df.loc[cond, "field"] = "value"` ❧śèṭś : fìèḻd: ~ : waḻúè: ìñ aḻḻ řõwś wĥèřè : ćõñd: ìś Třúè (reference)

`df.apply(function)` ❧aqqlièś śõmè fúñćṭìõñ ṭõ ṭĥè daṭafřamè (reference)

- `axis = 1` ❧ṭèḻḻś ṭĥè fúñćṭìõñ ṭõ ṭřèaṭ ṭĥè df ìñ ROWS»dèfaúḻṭ>COLUMNS

`def function(x):` ❧ñèèd ṭõ dèfìñè ṭĥè fúñćṭìõñ

`df["field1"] = df.apply(function)` ❧śawè ṭĥè řèśúḻṭ ìñ a ñèw ćõḻúmñ

# DEAL WITH TIME

## Data Conversion

`pd.to_datetime(df.column)` ❧Túřñ a śṭřìñğ ìñṭõ a daṭèṭìmè»wìṭĥõúṭ ařğś«ḻèawèś fõřmaṭ úñćĥañğèd (reference)

- `format="format"` ❧śqèćìfý ṭĥè fõřmaṭ ✷è✷»:} Y✷} m✷} d: / (formats)

`df['field1'].apply(lambda (t): t.strftime('format'))` ❧ṭõ ṭřañśfõřm a daṭèṭìmè ćõḻ ìñṭõ a śṭřìñğ (reference)

## Extract datetime info

`df['field1'].year` ❧ğèṭ ṭĥè ýèař fřõm a daṭè fõřmaṭṭèd ćõḻúmñ (reference)

`df['field1'].month` ❧❧❧❧❧ṭĥè mõñṭĥ

`df['field1'].day` ❧❧❧❧❧ṭĥè daý õf ṭĥè mõñṭĥ

`df['field1'].dayofweek` ❧❧❧❧❧ṭĥè daý õf ṭĥè wèèk

## Timedelta

`td = datetime.timedelta()` ❧ğèṭ a řèqřèśèñṭaṭìõñ fõř a ṭìmè ìñṭèřwaḻ (reference)

- `days=n` ❧śqèćìfý ñúmbèř õf daýś

`td.days` ❧řèṭúřñ ñúmbèř õf daýś ìñ a ṭìmèdèḻṭa õbḻèćṭ

`td.years` ❧řèṭúřñ ñúmbèř õf daýś ìñ a ṭìmèdèḻṭa õbḻèćṭ

## Filter date columns

`df['YYYY']` ❧śèḻèćṭ řõwś fřõm ṭĥaṭ ýèař

`df['YYYY':'YYYY']` ❧śèḻèćṭ řañğè õf ýèařś

## Mathematical modifications

`df.resample('rule')` ▸ aggregate data to some specific time interval ([reference](#)), ([rules](#))

`df.rolling(n)` ▸ aggregate with n neighbors»chain with »mean«/«sum«/ or other ([reference](#))

# DISPLAY OPTIONS FOR JUPYTER NOTEBOOKS

`pd.set_option("optionName", value)` ▸ Change the behavior of displayed content in notebooks ([reference](#))

- `"display.max_rows"` ▸ The number of rows of a DataFrame
- `"display.max_columns"` ▸ The number of columns of a DataFrame
- `"display.max_colwidth"` ▸ The number of chars inside a column
- `"display.float_format"` ▸ something like :\×8f-:»format
- etc.

`pd.options.display.max_rows` ▸ to display the current settings