

```
df = df.set_index('A')
```

```
df = df.astype(dtype)
```

```
df = df.T
```

```
top_left_corner_df = df.iloc[:4, :4]
```

```
df.set_index
```

PD.TAB

```
df_copy = pd.concat
```

```
df.merge
```

```
df = df[df['col2'] >= 0.0]
```

```
df = df[df['col1'].isin([1, 2, 5, 7, 11])]
```

```
df = df[df['col1'].str.contains('hello')]
```

```
Value = df.loc['row', 'col']
```

```
df = df.loc[rows = row3, 'col1' : 'col3']
```

```
df = df['col3'] >= 10 | df['col1'] < 5
```

```
df[x] = df[x].str.split(';').str[1]
```

```
for x in list(df.columns):
```

```
df[x] = df[x].str.split(';').str[1].astype(int)
```

```
pwd.index = pd.to_datetime(pwd.index)
```

```
titanic.sex.value_counts()
```

```
titanic.groupby('sex')['survived'].count()
```

```
{s: titanic[titanic.sex == s].sex.count() for s in titanic.sex.unique() }
```

```
titanic.pivot_table(index='sex', values='survived', aggfunc='len')
```

```
pd.crosstab(titanic.sex, titanic.survived).sum(axis=1)
```

```
titanic.pivot_table(index='sex', columns='class', values='survived')
```

df. value_counts() | [x for sublist in L for x in sublist] | [A E I O U Y a e i o u y]

df.groupby('A').B.~~value_counts~~count()

df.pivot_table(index='A', values='B', aggfunc=len)

pd.crosstab(df.A, df.B).sum(axis=1)

df.pivot_table(index='A', columns='B', values='C')

pd.crosstab(df.A, df['B']) #relative

df.pivot_table(index='A', columns='B', values='C', aggfunc=[np.mean, np.median])

datum

html

<[1?]+>

worden

'\b[Aa]\b'

df=pd.read_excel(file, index_col='A')

(df.A.str.findall('abc').str.len() / df.A.str.len()).describe()

df.A.str.contains('abc').sum()

df.join(df2, lsuffix='_suf')

sns.scatterplot(x=df.A, y=df.B)

corr = df.A.corr(df.B)

dict2 = {x: dict[x] for x in dict if dict[x] > 4}

df = pd.read_csv(file, sep='\t', encoding='ISO-8859-1').sort_values(by='A').drop(remove_cols, axis=1)

df = pd.read_csv(file, sep='\t', usecols=use_cols, index_col=0).sort_index()

np.mean, np.median, df.loc[df.A > mean]

df_series = pd.Series(df.A)

df_series[:20].plot(logy=True)

df.A.mean()

df.A.median()

df.A.std()

df2 = df.groupby('A').sum().sort_values(by='B', ascending=F)

intersection = set(a) & set(b)

union = set(a) | set(b)

jaccard = int / un

sns.distplot(df.A)

plt.subplots()

sns.distplot(df.A.apply(np.log))

Pivot-table (welke data, index = x, columns = y)

- SNS - LOAD DataSet ('planets')

- ~~AGG~~ ReCat = multiple functions

make it with GroupBy:

- Filter Groups with filter (FilterFunc)

FilterFunc = x[Data] > y

- Transform By Group Data (lambda x: x - x.mean())

- APPLY APPLIES ARBITRARY function

open Bestand:

- List in GP to Define WHICH Group

with open (Bestand) as f

- {A: Vowel, B: Constant, C: Constant} picks index

for l in f:

- str.lower picks index

for l in list(l)

- 10 // x[year] * 10

- Pivot table me GP = - GP (x1, x2)

make date with datetime()

Regex =

pd.datetimeindex() make a date index

één van: [ABC]

SUBTRACT index from another to get [0 days, 1 days]

niet één van: [^ABC]

- make df with pd.date_range

any = .

White sp = \s

Digit = \d

één of meer: +

0 of 1 = ?

n = start

= end

WB: \B

3 VAN A = A{3}

Letters: [A-Z]

- Normalized min max = (df - df.min()) / (df.max() - df.min())

- Normalize z-score = (df - df.mean()) / df.std()

- SNS. PAIRPLOT is used to visually find correlations

- Zipfile - Zipfile() use to read but not open zip

- errorbar(x, y, yerr = 0.5, fmt = '.k')

- fill - Between (xfit, yfit - dyfit, yfit + dyfit, color)

- plt(histo for histogram

- Distplot Combines KDE & Histo

- SNS - FacetGrid (tips df, row = s, column = t, makebottle)

- Grid-map (plt.hist, tip-plot, bins = np.linspace(0, 40, 15))

- SNS - factorplot ("day", "total", "sex", data = tips, kind = "box")

↳ map total against day for FANDOM

`df['status'].astype('cat')`
`df['status'].cat.set_categories(['won', 'lose'], inplace=True)`
`df.pivot_table(index=['M', 'S'], columns='P', values=['Q', 'Pi'], aggfunc={'Q': len, 'Pi': np.sum}, fill_value=0, margins=True)`

test acc = 95%
 scc = 1/1000
 $TP = .95 * 0.001$
 $FP = (1 - .95) * 0.999$
 $TP / (TP + FP)$
 $1 - (TP / (TP + FP))$
 of: $FP / (TP + FP)$

`table.query('Manager == ['Mick'])`
`df.country.str.replace('.', '')`
`df.norm = (df - df.min()) / (df.max() - df.min())`

`df.country.str.replace('.', '')`
`df.norm = (df - df.min()) / (df.max() - df.min())`
 Z-score: $(df - df.mean()) / df.std()$

`titanic.sex.value_counts()`
`titanic.groupby('sex')['survived'].count()`
`tit.pivot_table(index='sex', values='survived', aggfunc=len)`
`pd.crosstab(tit.sex, tit.survived).sum(axis=1)`
`tit.pivot_table(index='sex', columns='class', values='survived')`
`pd.crosstab(titanic.pclass, titanic['class'])`

for l in tqdm_notebook(j):
 rawdata = l.split('it')
 data = []
 totalEdits = 0
 for n in range(1, len(rawdata)):
 currentEdits = int(rawdata[n].split(' ')[1].strip('n'))
 if currentEdits > 1:
 totalEdits += currentEdits

`t = tit.dropna(subset=['age'])`
`t.pivot_table(index='age', columns='sex', values='adult-male', aggfunc=sum)`
`pd.testing.assert_frame_equal(t.query("sex == 'male' and age >= 16"), t['adult-male'])`

`def teldeletters(bostand):`
 from collections import defaultdict
 telling = defaultdict(int)
 with open(bostand) as f:
 for l in f:
 for c in list(l):
 telling[c] += 1
 tellingseries = pd.Series(telling).sort_values()
 return tellingseries

if totalEdits > n and len(data) <= 5:
 for lang in data:
 edits[lang] += lang
 edits[" "] += 1
 if len(data) > 1:
 Edits-by-mult[lang] += lang
 multilin[lang] += 1
 if len(data) >= 2:
 data = sorted(data)

`teldeletters(bostand).plot(kind='bar', figsize=(15,7), logy=True)`
`prince.text.str.contains('1b[SS]ex1b').sum()`
`prince. " " " (r'^[SS]ex1b').sum()`
`u. " " " startswith('Sex').sum()`

`eg = " " " .str.findall(r'[SS]ex1w+')`
`eg[eg.str.len() > 0].head().count()`
`Counter([S for L in eg.values for S in L])`

$a^*, a^+, a^? \rightarrow$
 $0 \geq 1 \geq 017$
 $\{a\} \{1,3\}$

`prince.text.str.findall(r'1b[A-Za-z0-9]{0-9}[A-Za-z][A-Za-z0-9]{1b|1b[A-Za-z0-9]{0-9}[A-Za-z][A-Za-z0-9]{0-9}[A-Za-z0-9]{1b}')`

`(100 * prince.text.str.findall(r'[AETOUyaexuy]').str.len())`
 any char exc newline

`Planets.groupby('method')['distance'].max()`
`Planets[" "] .apply(sort_values('distance'))`
`df.index = df.index.droplevel(0)`
`Planets.loc[]`
`Planets. loc []`
`Planets. loc []`

| word | IS -> word idigit | whitespace
 | word | IS -> NOT " " |
 | abc | -> any -> [^ abc] -> not
 | a-g | -> between | b | B -> word bound not word bound

Hoeverel mensen zaten er per geslacht in elke klasse?

```
- (titanic.pivot_table(index='sex', columns='class', values='survived', aggfunc=sum) /
```

↳ hetzelfde ..., aggfunc=len)

Hoeverel heeft de ramp daarvan overleefd?

```
titanic.pivot_table(index='sex', columns='class', values='survived')
```

```
pd.crosstab(titanic.sex, titanic['class'])
```

Wat is de gemiddelde en median leeftijd per geslacht v.d mensen die t hebben overleefd en niet.

```
titanic.pivot_table(index='sex', columns='survived', values='age', aggfunc=[np.mean, np.median])
```

Beschrijf de relatie tussen de variabele class en pclass.

```
pd.crosstab(titanic.pclass, titanic['class'])
```

Wann is de variabele 'adult_male' waar?

```
t = titanic.dropna(subset=['age'])
```

```
t.pivot_table(index='age', columns='sex', values='adult_male', aggfunc=sum)['male']
```

```
pd.testing.assert_frame_equal(t.query('sex="male" and age >= 16'), t[t.adult_male])
```

1. Hoeverel nummers bevatten het woord -sex.
2. Hoeverel beginnen met dat woord
3. Hoeverel nummers bevatten een woord sex, maar is niet gelijk aan -sex.
4. print al die woorden met doe -sex. bevatten
5. ~~print~~ vind alle woorden in zijn langs die bestaan uit gemiddelen letters en andere tekens als (02, 4ut

Apreprocess.

```
prince = pd.read_csv(bestand)
```

```
prince.text = prince.text.astype(str)
```

```
prince.head()
```

```
1&2. print(prince.text.str.contains(r'\b[0-9]sex\b').sum())
```

```
print(prince.text.str.contains(r'^[0-9]sex\b').sum())
```

```
prince.text.str.startswith('sex').sum() + prince.text.str.startswith('sex').sum()
```

```
#3. ef = prince.text.str.findall(r'[0-9]sex\b')
```

```
print(ef[ef.str.len() > 0].head(10))
```

```
print(ef.value_counts())
```

```
#4. from collections import Counter
```

```
print(Counter([s for l in ef.values for s in l]))
```

```
titanic.sex.value_counts()
titanic.groupby('sex')['survived'].count
titanic.pivot_table(values='survived', columns='sex', aggfunc='len')
pd.crosstab(index=titanic.sex, columns=titanic.survived).sum(axis=1)
```

Prob A given D

```
df.groupby('A')['B'].value_counts() / df.groupby('A')['B'].count
```

```
{s: titanic[titanic.sex == s].sex.count() for s in titanic.sex.unique() }
```

```
titanic.pivot_table(ind=sex, col=class, val=survived, aggfunc=sum) / ... agg=len of mean
pd.crosstab(titanic.sex, titanic['class'])
```

```
titanic.pivot_table(sex, survived, age, aggfunc=[np.mean, np.median])
pd.crosstab(titanic.pclass, titanic[class])
```

```
t = titanic.dropna(subset=['age'])
t.pivot_table(ind=age, col=sex, val=adult_male, aggfunc=sum)['male']
pd.testing.assert_frame_equal(t.query("sex == 'male' and age >= 16"), t[['adult_male']])
```

```
from collections import Counter
let_count = defaultdict(int)
with open('best', encoding='utf8') as f:
    for line in f:
        for letter in list(line):
            let_count[letter] += 1
pd.Series(let_count).sort_values(ascending=False)
```

- . = any char, + = a lot
- \d = digit
- \w = word char
- \s = white space
- \ = escape special char
- ^ = except
- = ... to ...
- {3} = 3 matches
- [^abc] = any char except abc
- ^ = start
- \$ = end
- ^...\$ = exact match
- roar = any string with roar in it

```
prince.text.str.contains(r'\b[5s]ex\b').sum()
prince.text.str.startswith('Sex:').sum()
prince.text.str.findall(r'[5s]<w+').sum()
df[df.str.len() > 0].count()
Counter([s for l in df.values for s in l])
100 * prince.text.str.findall(r'[AEIOUYaeiouy]').str.len() / prince.text.str.len().count()
```

1 of 1000 500

$$TP = 0.95 \times 0.001$$

$$FP = 1 - 0.95 \times 0.999$$

```
carsdf.hp.fillna(carsdf.hp.mean())
carsdf.hp.replace(0, carsdf.hp.mode(), inplace=True)
... = carsdf.hp.apply(lambda x: np.random.choice(carsdf.hp.dropna()))
```

$$TP / (TP + FP) = \text{Recall precision}$$

$$FW = \text{recall}$$

```
z score = (df.col - df.col.mean()) / df.col.std()
min-max normal = df - df.min() / (df.max - df.min)
normalize = P / (P.sum(axis=0))
```

pd.cut = data in bins

```
filter = titan[titan[class] == 'first'].groupby('sex')['survived'].sum()
pd.items(), set(), df.loc[df['n_sec'].idxmax()], title, df.at[1, value]
root waardelijke kans = A & B / B
```

reindex

```
df.apply(pd.value_counts, axis=1).fillna(0).astype(int) / apply(pd.Series).stack().value_counts
```

```

pd.Series([1, 2, 3, 4, 5, 6, 7], index=[2, 5, 3, 2])
pd.Series({'a': 1, 'b': 2})
A.add(B, fill_value=0)
data.deopna() axis=0, how='all', thresh=3
data.fillna(0)
index = pd.MultiIndex.from_tuples(index)
pop[2010]
pop.unstack().level0, 1
pop.index.names = ['state', 'year']
pd.MultiIndex.from_product([['a', 'b'], [1, 2]])
np.concatenate([x, y, z])
pd.concat([s1, s2]) ignore_index=True, join='inner'
df1.append(df2)
pd.merge(df1, df2) on='employee'
df1.join(df2)
merged[merged['population'].isnull()]
df.query("year == 2010 & ages == 'total'")
df.groupby("age").sum().aggregate(["min", "max"])
pd.cut(titanic["age"], [0, 18, 35])
pd.qcut(titanic["fare"], 2)
margins = True -> all
.astype(int)
ser.str.capitalize().extract("(^[A-Z]+)", expand=False)
pd.Series(["AEIOU"] * ["aeiou"])
spl.it(1).str.get(-1)
full_monte["info"].str.get_dummies("1")
with open("file", "r") as f:
    data = [line.strip() for line in f]
recipes.name[ry.argmax(recipes.ingredient.str.len())]
    longest list
recipes.describe.contains(["Bb"])
recipes.name[selection.index]

```

```

with open("file", encoding="utf-8") as f:
    editors = defaultdict(int)
    comblist = C
    D = E3
    for l in tqdm_notebook(F):
        ll = l.split(" ")
        for u, v in [(i, split(" ")[0]), int(i, split(" ")[1], restap(" \n")))] for i in ll[1:]:
            editors[u] += 1
            editors[v] += 1
        D[ll[0]] = E1.split(" ")[0] : int(i, split(" ")[1], restap(" \n")) for i in ll[1:]
    for u, v in D.items():
        if len(v) > 1:
            comblist.append(list(Combinations(sorted(v.keys()), 2)))
            for i, j in v.items():
                multilingual[i] += 1
                Edits_by_mults[i] += j
    for i in comblist:
        for j in i:
            pairs[i] += 1
pd.DataFrame.from_dict(editors, orient='index', columns=['Editors'])
dataset_stats_wikipedia = pd.read_csv(local_stats, sep="t", index_col='Code', usecols=['Code', 'Editors'])
.sort_index()
country_str_replace("USA", "USA")
country.plot.scatter()
sns.lmplot()
min_max
(df - df.min()) / (df.max() - df.min())
z_score
t60z = (t60 - t60.mean()) / t60.std()
sns.distplot(t60z)
sns.pairplot(df[["C"]])
sns.heatmap(chuandf.corr())
chuandf["C"] = chuandf["C"].str.strip() == "True"
pd.crosstab(chuandf["ig"], chuandf["C"]).plot.bar(stacked=True)
sns.lord_datasets("es")
df.sort_values(["", ""], ascending=[False, True])
df["-d"] str.contains("-")

```

```

defen
titanic.sex.value_counts()
titanic.groupby("sex")["survived"].count()
s: titanic[titanic.sex == s].sex.count() for s in titanic.sex.unique()
titanic.pivot_table(index="sex", values="survived", aggfunc="len")
pd.crosstab(titanic.sex, titanic.survived).sum(axis=1)
titanic.pivot_table(index="sex", columns="class", values="survived")
aggfunc = sum / aggfunc = len
titanic.pivot_table(index="sex", columns="survived", values="age", aggfunc="np.mean")
pd.crosstab(titanic.pclass, titanic["class"])
t: titanic.describe(subset=["age"])
t.pivot_table(index="age", columns="sex", values="adult_male", aggfunc="sum")
    ["male"]
pd.testing.assert_frame_equal(t.query("sex == 'male' and age >= 16"),
    t["adult_male"])
def kelleletters(bestand):
    from collections import defaultdict - counter
    telling = defaultdict(int)
    with open(bestand) as f:
        for l in f:
            for c in list(l):
                telling[c] += 1
    telling.series = pd.Series(telling).sort_values(asc=False)
    return telling.series
since = pd.read_csv(bestand) - names = [ ] since.rows = 1
since.text = since.text.astype(str)
since.text.str.contains("^[Ss]ex\b").sum() - mid
since.text.str.contains("^[Ss]ex\b").sum() - begin
since.text.str.startswith("Sex ").sum() + - 1 - ("Sex ").sum()
since.text.str.find("^[Ss]ex\b")
['f', str.len() > 0].count()
counter([s for l in f.values for s in l])
100 * since.text.str.find("^[AEIOUaeiou*]").str.len() /
since.text.str.len() = doscat BEL

```

```

TP / (TP + FP) - pos sad test -> besmet
1 - (TP / (TP + FP)) = FP / (TP + FP) - pos sad test -> niet besmet
    = egg
TP = 0.95 * 0.001
    accuracy 1 op 1000
FP = (1 - 0.95) * 0.95
    kans op geen soc
df["-"] corr(df["-"])
jaccard =
A = set(i, sec)
B = set(j, sec)
Return len(A & B) / len(A | B)
    intersection union
df[sections].apply(pd.Series).stack().value_counts()
plt.xticks(1)
chance_dict = dict()
for b in frequent.sections:
    chance_dict[b] = dict()
for a in frequent.sections:
    if dummy_df[b].sum() > 0:
        chance_dict[b][a] = dummy_df[dummy_df[b] == True][a].sum() / dummy_df[b].sum()
    else:
        chance_dict[b][a] = 0

```

```

titanic.sex.value_counts(), titanic.groupby('sex')['survived'].value_counts()
{s: titanic[titanic.sex==s].sex, count()} for s in titanic.sex.unique()
titanic.pivot_table(index='sex', values='survived', aggfunc=Len)
pd.crosstab(titanic.sex, titanic.survived).sum(axis=1)
titanic.pivot_table(index='sex', columns='class', values='survived')
titanic.groupby(('sex', 'class'))['survived'].mean().unstack()
t=titanic.dropna(subset=['age']) #pivot_table(index='age', columns='sex', values='adult_male', aggfunc=sum) ['male']
pd.testing.assert_frame_equal(t.query("sex == 'male' and age >= 16"), t[["adult_male"]])

```

```

bestand = 'prince_lyrics.csv' # prince = pd.read_csv('prince_lyrics.csv')
from collections import defaultdict # prince.text = prince.text.astype(str)
def teldeletters(bestand): # prince.text.contains(r'\b[Ss]ex\b').sum()
    aantal = defaultdict(int) # prince.text.startswith('sex').sum()
    with open(bestand) as f: # ef = prince.text.str.findall(r'[Ss]ex\w+', print(ef[ef.str.len() > 0].count()))
        for l in f: # print(Counter([s for l in ef.values for s in l]))
            for c in list(l): # # Alleen woorden die zowel een letter als een cijfer bevatten en of alleen daaruit besta
                aantal[c] += 1 # ef = prince.text.str.findall(r'\b[A-Za-z0-9]*[0-9]+[A-Za-z0-9]*[A-Za-z][A-z0-9]*\b|b[A-Za-z0-9]*[A-Za-z][A-Za-z0-9]*[0-9][A-Za-z0-9]*\b')
    aantal.series = pd.Series(aantal).sort_values() # print(Counter([s for l in ef.values for s in l]))
    return aantal.series

```

```

teldeletters(bestand).plot(kind='bar', figsize=(15,7), logy=True)
#percentage klinkers => (100 * prince.text.str.findall(r'[AEIOUYaeiouy]').str.len() / prince.text.str.len()).describe()
with open(load_userlang(), encoding='utf-8') as f: # Hierarchical indexing
    c = 0 # df.index, droplevel()

```

```

from collections import defaultdict # index = [(('New York', 2000), ('New York', 2010),
editors = defaultdict(int) # ('Texas', 2000), ('Texas', 2010))
edits = defaultdict(int) # population = [1, 2, 3, 4]
multiLingual = defaultdict(int) # index = pd.MultiIndex.from_tuples(index)
Edits-by-multiLinguals = defaultdict(int) # pop = pd.Series(population, index=index)
pairs = defaultdict(int) # pop.unstack() => New York 2000 2010
for f in tqdm_notebook(f): # Texas 2000 3 4
    line = l.strip() # * unstack(level=0/1)
    list_1 = l.split("\t")[1:]
    list_2 = [word.split(',') for word in list_1]
    dict_1 = {x[0]: int(x[1]) for x in list_2}
    data = dict_1.keys()
    data = sorted(data)
    if len(data) > 1:
        pair_list = []
        for p in data:
            pair_list.append(p)
        combi = combinations(pair_list, 2)
        for sp in combi:
            if sp not in pairs:
                pairs[sp] = 1
            else:
                pairs[sp] += 1
        if sum(dict_1.values()) >= 5 and len(dict_1) < 6:
            for x in dict_1.keys():
                if dict_1[x] > 1:
                    editors[x] += 1
                    edits[x] += int(dict_1[x])
                if len(dict_1) > 1:
                    multiLingual[x] += 1
                    Edits-by-multiLinguals[x] += int(dict_1[x])
    c += 1

```

```

df3 = pd.merge(df1, df2)
df1.join(df2.set_index('employee'), on='employee')
dfnom = (df - df.min()) / (df.max() - df.min())
z-score: (df - df.mean()) / df.std().describe()

```

```

with open('... .json', 'r') as f:
    data = [line.strip() for line in f]
    data_json = "[{}]" .format(', '.join(data))

```

```

c += 1

```

precision = $\frac{TP}{TP+FP}$

recall = $\frac{TP}{TP+FN}$

		predicted	
		-	+
actual	-	TN	FP
	+	FN	TP

- * standaarddeviatie
 - bereken het gemiddelde
 - neem van elk getal de afstand tot gemiddelde
 - neem het kwadraat van die afstanden
 - bereken het gemiddelde van die kwadraten
 - neem de wortel van de uitkomst
- * variance:

$$(x_1 - x_1, \text{mean}()) ** 2 \text{ summed}$$
- * standaarddeviatie:

$$np.sqrt(((x_1 - x_1, \text{mean}()) ** 2), \text{mean}()) == x_1, \text{std}()$$

Stel 95% accuracy & 1/1000 mensen heeft SOA

TP = (test ja, SOA) = $0,95 * 0,001$

FP = (test ja, geen SOA) = $(1 - 0,95) * 0,999$

Wat is de kans geen SOA maar test positief

$1 - (TP / (TP + FP))$ of $FP / (TP + FP)$

Series

pd.Series([...], index=[...]) or dict
series.values / index / keys() / items()

DataFrame

df.index / columns / items() / iter^{rows} / iter^{columns} ()
pd.Index -> &, ^ pd.DataFrame (data, index=, columns=)

Slicing, masking & fancy

series[0:2] / ['a':'c'] excl / incl
series[(series > 0.3) & (series < 0.8)]
series[['a', 'e']]
series.loc[1] explicit index
series.iloc[1] python-style
df.T swap rows & cols
df.values[0] row df[0] col
df.iloc[:3, :2] python
df.loc[: 'Illinois', : 'pop'] incl
df.ix[:3, : 'pop'] hybrid (excl & incl)
df.loc[df.density > 100, ['pop', 'density']]
df.iloc[3, 2] = 80 assign
df[1:3] -> rows
df[df.density > 100]

Ufuncs

np.exp(series) / np.sin(df * np.pi / 4)
A+B or A.add(B, fill_value=num)
f.stack.mean() all values
df.subtract(df['R'], axis=0) column-wise

Missing values

np.isnan(vals)
isnull() boolean
notnull() opposite
dropna (axis=1) for cols
how='any/all'
thresh=3 inplace=True
fillna (method='ffill')
axis=1

Multi index

pd.MultiIndex.from_tuples() / from_arrays()
df.reindex(), level=0/1/...
df.unstack() series to df
df.stack() df to series
df.index.names = ['ind1', 'ind2']
df.sort_index()

from_product() * 0 or more + 1 or more ? 0 or 1 {2}
(levels = [['a', 'b'], [1, 2]], [2, 5] {2, 3} {5} up to 5
labels = [[0, 1, 1, 0], [1, 1, 0, 0]]) 1 & 1n
pd.cut(df['age'], [0, 18, 80]) 1 A start of string
re.compile match, search, split, sub

Concat, merge

pd.concat([ob1, ob2], join='outer', ...)
df1.append(df2)
pd.merge(df1, df2)
df1.join(df1b)

df.col.value_counts()

Query

df.query("col==2010 & other-col== 'total'")
@ var pd.eval('df2.T[0] + df3.iloc[3]')

Aggregation

df.mean (axis='columns')
df.describe(), count(), first/last(), mod(), prod()

Groupby

[...].aggregate(['min', np.median])
df.groupby('col')['col2'].sum()
def filter_func(x):
 return x['data'].std() > 4
df.groupby('col').filter(filter_func)
df.groupby('col').transform(lambda x: x - x.mean())
df.groupby('col').apply(norm-by-data)

Pivot tables

df.pivot_table('value-col', index=, columns=)
df.pivot_table(data, values=, index=, columns=, aggfunc=, fill_value=, ...)
1 or list 1 or list

Strings

col.str.match() boolean
extract() matched groups as str
findall each elem
replace()
contains() boolean
count() occ. of pattern
r.split() excepts regex
get_dummies()
get

Regex

^ start, \$ end, \b word boundary, \B non-word bound
\d one digit, \D non-digit, \w, \W, \s whitespace, \S
[abc] one char [^abc] not
* 0 or more + 1 or more ? 0 or 1 {2}
(levels = [['a', 'b'], [1, 2]], [2, 5] {2, 3} {5} up to 5
labels = [[0, 1, 1, 0], [1, 1, 0, 0]]) 1 & 1n
1 A start of string

String operations

decade = 10 * (year // 10)
 decade = decade.astype(str) + 's'
 decade.name = 'decade'
 Groupen by decade

Series

s.unstack() → multiindex → kolommen ~~en~~ index
 s.unstack(0) → → kolom 0 index
 pd.Series(np.diag(df))
 df[df != np.diag(df)]
 s.to_frame().reset_index() → df met index series als kolom
 ser1[~ser1.isin(ser2)]

str.split().str[0]
 str.contains(r'\b[SS]ex\b').sum()
 str.contains(r'\^[SS]ex\w+').sum()
 R = str.findall(r'[SS]ex\w+')
 R[R.str.len() > 0]
 from collections import Counter
 Counter([s for L in R, values for s in L])

(100 * df.kolom.str.findall(r'[AEIOUYaei...]').str.len() / df.kolom.str.len())

Capitalize ~~at~~ list with strings and a None → make a series of the list and do Series.str.capitalize()

str.replace(' ', '').str.lower().str.strip.str.replace("''", "").str.replace("'''", "")

Dataframe operations

df.mean(), mean(1)
 df.count() → aantal regels per kolom
 df.dropna()

↑ spaces verwijden
 From collections import defaultdict
 telling = defaultdict(int)
 with open(bestand, encoding='utf8') as f:
 for L in f:
 for c in L:
 telling[c] += 1

as F:

Groupby & pivot etc

df.loc[indices].kolom → waarde uit kolom bij bepaalde index
 df.loc[df.kolom.isin([list])] → rijen waar kolom waarde is in lijst

alles maar in 0 & 1 R
 [A B C] alles Beltaal A, B of C

(df[df.mean().index] - df)
 len(set(df.index)) == len(df.index) → false → duplicate index etc

A = begin string
 B = einde string
 [] → op → [Aa]
 ? = 0 of 1 x * 0, 1 meer
 + 1 of meer

- * aggregate([min, max, count]) verschillende berekeningen op Groupby
- * aggregate({kolom: functie}) functie alleen op bepaalde kolom
- * df.groupby(kolom).transform(lambda x: x - x.mean())
 Haal van elke cell de mean van de groep waarbij de cell hoort af
- * df.groupby(List) of (dict, nieuwe namen aan indexes geven en daar of groepen
- * df.groupby(str.lower).mean() lowercased indexname en mean groep

duplex_index = df.index.value_counts()
 duplex_index[duplex_index > 1]
 pd.cut(df.kolom, [0, 10, 30, 60, 80])

skiprows = 1 → 1ste rij / skiprows = [1] → index 1 → 2e rij
 [1, 2] → index 1 & 2 → 2e en 3e

df.sort_index(axis=1) → kolommen sorteren

	predicted	
	neg	pos
neg	TN	FP
pos	FN	TP

Theorie & formules

Conditionele kans = kans (A als A ook B bevat)

(A ∩ B) / B
 Intersection (∩) = die beide systemen
 Union (∪) = alle unieke n s/w te had
 sym. dif. (Δ) = in de ene maar niet in de and

Z score = (df - df.mean()) / df.std()
 precision = TP / (TP + FP)
 Recall = TP / (TP + FN)
 Accuracy = (TP + TN) / totaal

- * df.groupby([str.lower, list]) nieuwe namen indexes en oude naam lower multi index
- * pivot_table(df, values = waardes die je wilt tellen, index = index die je wilt, columns = aggfunct = functie over waardes)
 of aggfunct = {kolom: functie} values
 Hoeve je dan niet aan te geven

- * df.groupby(kolom)[[kolom]] → dataframe
 [kolom] → serie
 .idxmax() → vindt alle → max. indexes
- * df.iloc[[df.groupby(...).idxmax()]]
 ↑ vindt de rijen
 def norm_by_dataz(x)
 x['data1'] /= x['data2'].sum()
 return x
 .apply(norm_by_dataz)

TN = true neg → is neg, predic neg

TP = true pos → is pos, predic pos

FN = false neg → is pos, predic neg

FP = false pos → is neg, predic pos

$$\text{recall} = \frac{TP}{TP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

.aggregate

match()

extract()

findall() -

replace() -

contains() -

count()

split() -

rsplit() -

cat() -

df.value_counts()

Counter(df.sex)

pd.crosstab(df.sex, df.gender, margins=True)

df.pivot_table(c='class', i='sex', v='survived', aggfunc='count')

pd.cut(df['age'], [0, 10, 30])

of qcut

aggfunc = {'survived': 'sum', 'fare': 'mean'}