

$mean = \sum(x) / \text{len}(x)$
 $var = \sum(x - np.mean(x))^2 / \text{len}(x)$
 $std = \text{math.sqrt}(var)$

$A \cap B = \text{intersection of } A \text{ \& } B$
 $A \cup B = \text{union of } A, \text{ of } B$
 $P(A \cap B) / P(B) = \text{voorwaarde/jh kans op } A$

```

def grad(x, y):
    while y != 0:
        (x, y) = (y, x % y)
    return x
    
```

out: n y v u p y
 return (abs(a*b)/grad(a,b))

from collections import
 defaultdict, Counter

$RMSE = np.sqrt(((predictions - targets)^2).mean())$
 Normalize data = $z = \frac{x - \min(x)}{\max(x) - \min(x)}$

- any char except new line
- \w matches alphanumeric chars: [a-zA-Z0-9_]
- \W matches non-alphanumeric chars: [^ \w]
- \d matches digit: [0-9]
- \D matches non digit: [^\d]
- \s matches whitespace char [t n f r p e]
- \S matches non-whitespace char [^\s]
- [] Character class. Matches any character contained between the square brackets
- ^[] Negated char. matches any char that is not contained between the square brackets
- * matches 0 or more repetitions of preceding symbol
- + matches 1 or more repetitions of preceding symbol
- ? makes the preceding symbol optional
- {n, m} Braces. matches at least 'n' but not more than 'm' repetitions of preceding sym. {n, m}
- (xyz) Char group. matches the chars xyz in that exact order.
- | Alternation. matches either chars before or after the symbol
- \ Escape next character. u can use: [] () * ? ^ \$
- ^ matches the beginning of input
- \$ matches the end of input

$recall = \frac{TP}{TP + FN}$
 $precision = \frac{TP}{TP + FP}$
 pd.cut(df.age, [0, 10, 30])
 of qcut(df.age, 2)
 .str.strip() .str.split()

```

get_dummies()
monte.str.get_dummies('')
monte = ['BICID', 'BID']
  
```

TN = is neg, predic neg
 TP = is pos, predic pos
 FN = is POS predic neg
 FP = is neg, predic pos

- match()
- extract()
- findall()
- replace()
- contains()
- count()
- split()

```

pd.crosstab(df.sex, df.survived, margins=True)
df.pivot_table(c='class', i='sex', v='survived', afunc='count')
Counter(df.sex) OF df.sex.value_counts()
aggfunc = {'survived': sum, 'fare': 'mean'}
  
```

recall = pabhang

SEP

titanic.groupby(['sex', 'class'])['survived'].agg(nyrate=mean).unstack = .pivot('survived', index='sex', columns='class')

df[col], str.count('word').sum()

df.col.value_counts()

df.groupby('col1')['col2'].count()

df.groupby('col1', axis=0).mean()

df.groupby(['col1', 'col2'])['col3'].mean()

{s: df[df.col==s].col.count() for s in titanic.col.unique() }

~~df~~

df.pivot(index='col1', values='col2', aggfunc=len)

pd.crosstab(df.col1, df.col2, sum(axis=1))

value van col1 in value col2 by value col3

df.pivot(index='col1', columns='col2', values='col3', aggfunc=sum/len)

foo	bar	baz	zoo	A	B	C
one	A	1	y			
one	B	2	y			
one	C	3	y	1	2	3
two	A	4	z			
two	B	5	w			
two	C	6	w			

df -> df.pivot(index='foo', columns='bar', values='baz')

Stacked -> stacked.unstack()

multi index -> -1 index

stack.unstack(colname)

df.unstack().stack() = df.sort_index()

pd.crosstab(df.col1, df['col2'])

mean and median of age^{per-sex} for col survived:

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

Relatie tussen 2 columns -> pd.crosstab

pd.crosstab(df['col1'], df['col2'])

pd.testing.assert_frame_equal(df.query('...'), df[df.col])

def tel(f:file):

telling = defaultdict(int)

with open(f) as f:

for letter in f:

telling[letter] += 1

return telling

tel(frequentie van elk letter in file)

mak pd series

tseries = pd.Series(telling).sort_values(ascending=False)

Return tseries

tel(file).plot(kind='bar', figsize=(...), logy=True)

df = pd.read_csv(file)

df.col1 = df.col1.astype(str)

df.col1.str.contains('a|b')

df.col1.str.startswith('sex')

df.col1.str.findall('Bs').tolist()

Counter(Cs for i in df.values for s in P)

pd.read_csv(loadf(k), sep='|')

df.loc[df['col'] == 'value'] -> alleen rijen met bepaalde waarde uit een column

df.sort_values(by='column', ascending=False).drop

set(tuple(zip(df.col1, df.col2)))

df1.join(df2, how='right', rsuffix=('ing'))

pd.read_excel(f:file).set_index('col')

Soup = BeautifulSoup(open(file)).read() features = "x m"

Soup.findall('string')

for key, value in dict.items(): col.append(dict[key][N].count('string'))

df.iloc[0] = my words

np.count_nonzero
df[col].tolist()

$$\text{Precision} = TP / (TP + FP) \quad \text{Recall} = TP / (TP + FN)$$

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

actual	FN	TP
	FP	TN

$$\text{Precision} = TP / (TP + FP) \quad F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

```
titanic.pivot_table(index='sex', columns='class', values='survived', aggfunc='sum')
titanic.pivot_table(index='sex', columns='class', values='survived', aggfunc='len')
pd.crosstab(titanic.sex, titanic['class'])
```

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

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

```
titanic.dropna(subset=['age'])
titanic.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])
```

```
telldelimiters function (prince csv); def telldelimiters (bestand): - from collections import defaultdict - telling = defaultdict(int) - with open (bestand) as f: - for l in f: - for c in list (l): - telling [c] += 1
```

```
return tellingseries = pd.Series (telling, sort_values (ascending=False)) - return telldelimiters (prince.csv), plt.cbook: 'bar', figsize=(15,7), logy=True); || prince = pd.read_csv (bestand) prince.text = prince.text.astype (str) # prince.text.str.contains ('\b[55]ex\b', sum()) # of prince.text.str.startswith ('sex'), sum() + prince.text.str.startswith ('sex').sum() || ef = prince.text.str.findall ('[55]ex\b') - print (ef)
```

```
count() #4: print all words with the 'sex' befallen -> from collections import Counter print (Counter (l for l in ef.values for c in l)) #5: dinge die werden 8 mal befallen; -> ef = prince.text.str.findall ('\b[A-Za-z0-9]*[0-9]+[A-z0-9]*\b[A-Za-z0-9]*\b | \b[A-Za-z0-9]*[A-Za-z0-9]*\b') print (ef[ef.str.len() > 0].head (10)) -
```

```
prince.text.str.findall ('[AEIOUYaeiouy]'), str.len() / prince.text.str.len() describe() Normalise = (x - df.min()) / (df.max() - df.min()) Znormalise = (x - df.mean()) / df.std()
```

```
pd.merge (df1, df2, left_on='employee', right_on='name').drop ('name', axis=1) apply (lambda x: x[0])
```

```
titanic.groupby ('continent').head (1)
```

df.groupby(column)[column].max() True/False
 df.sort_values(column, ascending=False, inplace=True)
 df.column.value_counts() "how many"
 df.groupby(column)[column].count() "how many"
 df.pivot_table(index=column, values=column, aggfunc=len)
 pd.crosstab(df.column, df.column).sum(axis=1)
 df[column].sum() "column sum" Fill_value=0
 df.sum() "series with column name and sum" aggfunc=[np.mean, np.median]
 df.sum(axis=1) "series per row sum"
 df[column]=... "define new column" std() "standard dev"
 df.corr df.quantile(0.99) (df[df < q]).corr()
 df.stack() df.unstack() "per column, per index"
 df.loc[:, (df > 1).any()] "columns with any values > 1"
 df.loc[:, (df > 1).all()] "columns with only values > 1"
 df.set_index([column, column]) "Multi index"

df.values[0] "to get row"
 df[column] "to get column"
 df.iloc[:3, :2] "first 3 rows, 2 columns"
 pd.read_csv("file")
 pd.read_excel("file", index_col=1)
 df1 & df2 "intersection, set files"
 df1 | df2 "union, alles"
 df1 ^ df2 "verschillen"
 df.dropna() "filtered data"
 df.fillna() "filled null data"
 df.merge() "merge by column"
 df.join() "join by row"
 df.set_index(column, inplace=True)
 pd.DataFrame.from_dict(dict, orient=index).fillna(0, dtype=int)
 df.drop([value, value]) "drop from row"
 df.drop(column, axis=1) "drop from column"

		true condition		
		positive	negative	
predicted condition	positive	TP	FP	Precision = $TP / (TP + FP)$ Recall = $TP / (TP + FN)$ Accuracy = $TP + TN / \text{all}$
	negative	FN	TN	

- any char except newline
- character "o"
- string "ab"
- a or b
- 0 or more "a"
- escapes a special char
- 1 or more
- 0 or 1
- exactly 2
- between 2 and 5
- 2 or more
- up to 5
- [a-b-d] char between a/d
- [^a-b-d] char except a/d
- \d one digit
- \s one whitespace
- \w one word

- .contains("R")
- .findall("R")
- .startswith("R")

Regex

^ - start of string

\$ - end of string

\s - whitespace

\d - digit

\w - word

pd.read_csv()

.contains()

.startswith()

pd.pivot_table()

pd.crosstab()

From collections import Counter

* - 0 or more

+ - 1 or more

? - 0 or 1

any character except

(a|b) - a or b

(...) - group

[abc] - range

[^abc] - not a or not b or not c

[a-z] - lower case

[A-Z] - upper case

[0-7] - digit from 0 to 7

{3} - exact 3

5 Manieren Mannen en vrouwen

Titanic.sex.value_counts()

Titanic.groupby('sex')['survived'].count()

dict(Titanic.loc[:, 'sex'].values)

pd.crosstab(Titanic.sex, Titanic.survived)

.sum(axis=1)

Pandas

df.dropna(inplace=True)
 axis=1 index=1
 sep='|'
 index=col
 rename columns
 axis=0 kolom
 axis=1 rij

Pd. read_csv() df.columns = ['a', 'b', 'c']

Select: df[df.kolom == 'value']
 df.loc[df.kolom == 'value']

Toevoegen: df['naam'] = df.b / df.a

Index: df.set_index('naam').index

Sorteren: df.sort_values(['kolom', 'kolom'], ascending=[false, false])

Tellen: df.value_counts('kolom')
 tellen van kolom, hoe vaak komt iets voor?

Boolean mask: x[x % 3 == 0]

Samenvoegen: pd.merge(x, y) on=how='inner'

Pd.concat([x, y], ignore_index=True, join='inner')

x.append(y) df1.join(df2) ^{rsuffix}

1-1, many-1 : 1/2 heeft duplicaten, many-many

Groupby: df.groupby('key').sum()

key	data
0	A 0
1	B 1
2	C 2
3	A 3
4	B 4
5	C 5

df.describe()

df.groupby('Species')['length'].max()

geeft grootste waarde van kolom lengte met specie ervoor
 What are the max lengths for each specie?

pivot: pd.pivot_table(df, index='sex', columns='class')

class	first	second	third
female	0,9	0,4	0,1
male	0,7	0,5	0,1

plotten: df.plot(kind='barh') x = df.kolom

aanpassen: df.kolom = df.kolom.str.lower()

nan: df.fillna()

df.dropna()
 df.isnull()
 df.notnull() } Boolean

Pd.crosstab(df.kolom, df.kolom) [lijst] margins=True

rename: df = df.rename(translate, axis='index')

5. replace('i', 'one') inplace=True

Literatuur ^{datum}
 datetime -> df.kolom.to_datetime()

import gzip
 with gzip.open('input.gz', 'r') as fin:
 for line in fin:
 print(line)

5 verschillende manieren om aantal mlf te tellen

- tit, sex, value_counts()
- tit.groupby('sex')['survived'].count()
- %s: tit[sex == 'S'].sex.count() tops in tit.sex.unique()
- tit.pivot_table(index='sex', values='survived', aggfunc=len)
- pd.crosstab(tit.sex, tit.survived).sum(axis=1)

Hoeveel nummers bevatten het woord sex?

prince.text.str.contains('sex').sum()

Hoeveel ~~woorden~~ beginnen met dat woord?

prince.text.str.startswith('Sex').sum()

Dataframe van dic: dic = {'Party': party}

df = pd.DataFrame.from_dict(dic, orient='columns')

gemiddelde aantal van kolom

df.kolom.mean()

meer dan median

len(df.loc[df.kolom > df.median])

alle kolommen ~~namen~~
 alle kolommen values een .nalen

df.kolom.str.replace('.', '')

head printen van een boel files

!head ../Data/db/chuen.all

!cat ../Data/db/chuen.names

chuen_df = pd.read_csv('../Data/db/chuen.all', dtype={})

meerdere kolommen

df[['kolom', 'kolom']].sum

type
 df.kolom.dtype = str

soup = BeautifulSoup(open('loadlijst.txt')).read()

print(soup.prettify())
 partijen = [party.text for party in soup.findall('RegisteredName')]
 genders = [[j.text for j in i.findall('gender')] for i in soup.findall('Affiliation')]

aantal mannen = []
 for i in genderlist.values():
 aantalmannen.append(i['30'].count('male'))

plotten van meerdere kolommen op a-z volgorde
 df[['kolom', 'kolom']].sort_index().plot()

selecteer alle kolommen behalve 1

df.loc[:, df.columns != 'kolom']

hoe vaak komt iets voor in een kolom
 (df.kolom == 'string').value_counts()
 iets droppen kolom die index is gemaakt
 df.drop(kolom, inplace=True, axis=1)

genderlist = {}

for i in range(len(partijen)):
 if partijen[i] in genderlist:
 genderlist[partijen[i]].append(genders[i])
 else:
 genderlist[partijen[i]] = genders[i]

kolom toevoegen hoeveel vrouwen in kamer worden verwacht

dicv = {}

for party in df.index:
 for zetel, vrouw in secret[party].items():
 zetels = int(df.loc[df.index == party].zetels.values)
 if zetel == zetels:
 dicv[party] = vrouw
 df['vrouwen'] = dicv.values()

p = re.compile('lw+')

p.findall('string')

re.split('lw+', 'string')

re.sub('lw+', '', 'string')

lw = (a-z, A-Z, 0-9)

lw = not ^

ld = 0-g

ls = space, tab, newline

274975 f
 200960 M

~/5.../Data/db/chuen*

`.findall(r')`

`.contains()`

`.startswith()`

`^` = start of a string

`\d` = one digit 0 to 9

`\w` = word character

`\s` = whitespace

`\D` = not a digit

`\W` = not a word

`A-b-1`

`ABC`

`"+=)`

`+` = one or more

`{3}` = exactly three times

`*` = zero or more times

`|` = or

`[^abc]` = not a b or c

`(...)` = group

`[A-Z]` = uppercase

`\t` = tab

`[a-z]` = lowercase

`\n` = any except newline

`pd.crosstab(titanic.sex, titanic.class)`

`df.pivot(table) = pivot_table()`

`pd.pivot_table(index='sex', columns='survived',`

`values='age', aggfunc=[np.mean, np.median])`

`pd.crosstab()`

`pd['sex'].str.contains`

`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)`

Pivot → gem & median leeftijd per geslacht van overlevende & niet overlevende

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

df.pivot_table(index = 'sex', values = 'survived', aggfunc = len) → df [...]

df.pivot_table(index = [col1, col2], values = 'col3', margins = True, margins_name = 'Total')

df[df.col.str.startswith('G').values]

df[df.col.isin(['Eng', 'NL']).values]

Groupby

df.groupby('col')[waarde].agg(aggfunct [min, len ...])

axis = 0 (rows) or 1 (cols)

Filter

df.groupby('col').filter(lambda x: len(x) > 3)

transform df.groupby('col').transform(lambda x: (x - x.mean()) / x.std() * 10)

regeer

Series.str.extract('([A-Za-z]+)') ← losse woorden

findall('in[AETIOU]*[Aerou]\$', 'series') → species → [AETIOU]

A = start string
 \$ = end string
 [] = set of chars
 * = any char except \n
 + = 1, meer
 ? = 0, meer

mean absolute deviation

plot

df.plot(kind = 'bar', x = 'col1', y = 'col2', title = "...")

↳ logy = True

Inladen

var = pd.read_csv(file, sep = '\t', index_col = 0, usecols = [...])

with gzip.open('file.gz', 'r') as f: for line in f: print(line)

Pandas

df.drop_duplicates([col1, col2 ...])

df.iloc → integer based zoeken, weight of positie

df.loc → index based

df.query("...", inplace = True) → @ var

df.loc[data.col < 2000], [col1, col2]

nieuwe col maken

↳ df["newcol"] = ...

↳ df.eval("D = (A+B) / C", inplace = True)

df.dropna() → how = "all", axis = "columns"

df.fillna(value) == df.replace(np.nan, value)

df.transform(lambda x: x + 1)

Quantal kolonnen selecteren

df.loc[:, 'col1': 'col2'] . apply(np.argmax, axis = 1)

absolute value ↑

variance ↑

standard error ↑

df.apply(func)

↳ def func(x): return x

apply → row/col wise df

applymap → element-wise df

map → element-wise series

change col order

df = df[['col1, col2 ...]]

find dtype's → df.dtypes

last element df[-1:]

Regex .str.findall(r'...')

\d	- digit	[]	- Matches char in brackets
\w	- word char	[^]	- negates
\s	- whitespace		- either or
\b	- word boundary	()	- group
^	- beginning of string	quantifiers:	
\$	- end of string	*	- 0 or more
.	- Anything	+	- 1 or more
-	- symbol range	?	- 0 or 1
(? <! ignore me)		{3}	- exact number
		{3, 5}	- range

Pandas

df.pivot_table(index=' ', values=' ', columns=' ', aggfunc=[len, 'mean', 'median'])

Pd. ~~cross~~ cross tab (df.x, df.y). Sum (axis=1)

Pd. column.value_counts()

Pd. groupby('column1')['column2'].count()

↑
np.mean

Extra

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

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

$$A = \frac{TP + TN}{TP + TN + FP + FN}$$