

**Numpy functions**

Np.log2()  
 Np.absolute()  
 Np.nan  
 Np.linspace(...)  
 Np.arange(start, stop, step)

**np.mean()  
 np.median()  
 np.sum()**

**Matplotlib**

df.plot()  
 df.plot.set()  
 df.plot.scatter()  
 df.plot.set\_xlabel(...)  
 df.plot.set\_ylabel(...)  
 df.col.hist()

1	stay	switch
3	perc.	perc.
4	perc.	perc.

↳ **key:** stay & switch  
 y-as: perc.  
 x-as: 3,4

**Subsets**

sub = df[df.col == ...]  
 sub2 = df[(df.col1 = ...) & (df.col2 > ...)]  
 (df.col = ...).sum()/mean

**Pandas functions**

Pd.isnull()  
 Pd.notnull()  
 Pd.read\_csv()  
 Pd.read\_excel()  
 Pd.concat  
 ↳ can be used for chunks if chunk size is passed to Pd.read\_csv  
 Pd.crosstab()  
 ↳ Pass series for index, series for cols & margins = True for 'All'  
 ↳ crosstab[...].unstack() for same result as pivot-table()

**Pandas series**

df.col.value\_counts()  
 df.col.dtype  
 df.col.to\_datetime()  
 df.col.str.contains(...)  
 ↳ boolean, dus kan neg.mean()/sum() achter.  
 df.col.str.strip()  
 df.col.idxmax()  
 df.col.astype(type)  
 ↳ change type of col  
 df.col.values  
 ↳ get values from a col as an array  
 df.col.index  
 df.col.str.findall(pattern)  
 df.col.min/max/sum/mean/mode

**Plotly**

df.col.iplot(kind=...)  
 df.col.value\_counts().iplot()  
 df.crosstab.iplot()

% magic  
 % time  
 % timeit  
 % ls

**Linux command**

**Seaborn**

sns.set() → to use seaborn on matplotlib.pyplot as plt  
 sns.boxplot(df)  
 sns.pairplot(df, hue='categorieën', column=...)  
 sns.distplot(df)  
 ↳ combination of hist & kdeplot  
 sns.lmplot()  
 sns.heatmap(df.corr())  
 sns.jointplot()  
 sns.scatterplot  
 sns.kdeplot()  
 sns.FacetGrid(df, row, col)  
 ↳ after this, use grid.map(plt.hist, bins=...)  
 sns.plt.add\_legend()  
 sns.factorplot()

**Counter(list)**

**Strings**

s.len()  
 s.get()  
 s.split('separator')  
 s.count(...)  
 s.find(...)  
 \*' ' join(list)  
 s.lower()  
 s.replace()  
 s.splitlines()  
 s.starts with()  
 s.strip()  
 s.upper()

**regex**

re.compile('...')  
 re.findall('pattern', string)  
 re.sub('pattern', 'replace str', 'search str')

**Pandas dataframe functions**

df.iterrows()  
 df.sort\_index()  
 df.set\_index()  
 df.reset\_index()  
 df.applymap(lambda x: functie wat je met x doet en waar je x mee vervangt\*)  
 df.shape()  
 df.pivot\_table(index='colname', columns='colname')  
 df.groupby('colname')  
 df.col.dropna(subset='...')  
 ↳ subset to define in which col to look for missing values  
 df.sort\_values()  
 df.count() → non na values  
 df.round(int/dict/series)  
 df.mean  
 df.corr()  
 df.to\_pickle('filename')  
 ↳ save df as pickle file  
 df.describe  
 df.fillna()  
 df.head()/df.tail()

**Pandas subsets**

columns: df[['col1', 'col2', ...]]  
 single column: df['col'] as series  
 row: df.loc['row-index']  
 rows: df.loc[['row1', 'row2', ...]]  
 range of rows: df.loc['idx1':'idx4'] → +/m  
 rows & columns: df.loc[[rows][cols]]  
 rows w/ step: df.loc['idx1':'idx4':2]  
 row & column: df.loc['row', 'col']  
 all rows & column: df.loc[:, ['col1', 'col2']]  
 index selecting: df.iloc[slicing as np slicing]

[...] ↳ choice  
 [...] ↳ any except within  
 [a-z] ↳ any characters between  
 • ↳ any single char.  
 \s ↳ any white space  
 \S ↳ any non white sp.  
 \d ↳ any digit  
 \D ↳ any non digit  
 \w ↳ any word char.  
 \W ↳ any non word char  
 (a|b) ↳ match either a or b  
 a? ↳ 0 or one of a  
 a\* ↳ 0 or more of a  
 a+ ↳ one or more of a  
 a{3} ↳ exactly 3 of a  
 \$ ↳ end of string  
 ^ ↳ start of string

**Theory 1**

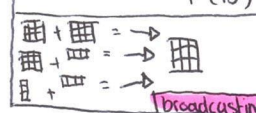
Accuracy = (TP+TN) / alles  
 Precision = TP / (TP+FP)  
 ↳ = P(ja | wit slag ja)  
 recall = TP / (TP+FN)  
 vb) accuracy 95%  
 occurrence 1 op 1000  
 hoeveel mensen met positieve test zijn niet besmet?  
 TP = 0.95 \* 0.001  
 FP = (1-0.95) \* (1-0.001)  
 Precision = TP / (TP+FP)

**Theory 2**

$P(A|B) = \frac{P(A \cap B)}{P(B)}$   
 $= \frac{P(A \cap B)}{P(B)}$   
 ↳ complement van Precision =  $1 - (TP / (TP+FP))$   
 $= FP / (TP+FP)$

**five different ways how many men and women were on boat**

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)



**broadcasting**



str. contains	str. extract	str. startswith	re.sub	re.findall
\d	digit	{3}	exactly 3x	
\w	letter	{2,4}	2 to 4 x	
\s	whitespace	{3,}	3 or more x	
\D	non-digit	*	0 or more x	
\S	non-whitespace	?	once or none	
[ ]	one of	^	start of string	
[a-z]	a to z	\$	end of string	

```
def strcount(file):
    i = defaultdict(int)
    with open(file) as f:
        for line in f:
            for letter in line:
                i[letter] += 1
```

omgekeerd sort = hoog -> laag dus  
ascending = false

nltk.tokenize.wordpunct\_tokenize()

$df[df[x] < df[x].quantile(.99)]$  -> verwijder 1% outliers

		actual values	
		P	N
predicted values	P	TP	FP
	N	FN	TN

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

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

import seaborn as sns.

y = sns.load\_dataset("x")

y.groupby(["x"]).sum/count() → dataframe gegroepeerd

y.sort\_values("z", ascending=False) → geef dataframe gesorteerd hoog-laag o.b.v. z

y["z"].value\_counts() → hoeveel z voorkomt in y

y.loc[y["z"] == "three"] → dataframe met alleen de regels waar 'three' in z voorkomt

y["z"].values → Returns array of values in z

y[["z", "w"]] → make dataframe y with columns z and w

np.mean, np.median, np.std

Pd.DataFrame.from\_dict(x, orient="index")

• sort\_index()

• idxmax() → hoogste index waarde

Pd.Index(a) & Pd.Index(b) → intersect

Pd.Index(a) | Pd.Index(b) → union

Pd.Index(a) ^ Pd.Index(b) → Symmetrische differentie

y[(y > 0.3)] = masking

y.T → transpose

y[y.x > 100]

• fillna()

y.droptna()

y.set\_index("x")

Pd.read\_csv("x", encoding="ISO-8859-1", sep=";" or sep="\t")  
index\_col=1  
• Capitalize() → make a capital letter of the first letter of str

Pivot voorbeeld: Maak pivot die de 10 meest voorkomende materialen geeft van de 10 meest voorkomende kunstenaars. Index moet de naam kunstenaar zijn, kolomnaam materiaal en value aantal keer dat de kunstenaar het gebruikt heeft

Antw: - to 10: kunst [["Naam.kunstenaar", "Id"], groupby(["Naam.kunstenaar"], count).sort\_values("Id", ascending=False): 10, index, values  
• kunst.pivot\_table(["Naam.kunstenaar"], ["Id"], values="Id", index="Naam.kunstenaar", columns="Material")

↳ aggfunc = len is count, sum = values ogetel

Pd.crosstab(y, z, y["w"])

Count how many time each word occurs from collections import defaultdict

with open(bestand as f:

for l in f:

for c in list(l):

telling[c] += 1

telling = defaultdict()

telling = Pd.Series(telling).sort\_values(ascending=False)  
Return telling

axis=1 (cols) = 0 (row)

Precision = TP / (TP + FP)

Recall = TP / (TP + FN)

Rij selecteren DF: ik["..."]

Series naar DF: Pd.DataFrame({"col1": ser1, "col2": ser2})

DF.describe() → count, mean, std, min, max



df = DataFrame  
 se = Series  
 xy = var (int)  
 pd.cats? ?

Weggeven

- pd.Series (data, index = 'index')
- pd.DataFrame (data, columns = [list], index = [list])
- df/se]. [index / columns] => namen van I on C
- df.size / shape / ndim / dtype
- se [se > 0, 3] & (se < 0, 8)] masking
- re.loc [df.col > x] & (df.col < y), [k1:col2] masking
- df.transpose => switch cols and index
- sum(), count(), median(), min(), max()
- df.dropna(). describe()
- count() => totaal aantal
- defaultdict() => makkelijk te toevoegen (int / float / str etc)
- with open (file) as f: regel voor regel
- pd.read\_csv [excel] (file)
- df.value\_counts()
- df.sort\_index / values
- df.query()
- pd.crosstab
- df.binc()
- with B22 File (file) open as f:

Groepen

- df.groupby('key'). sum()
- df.groupby('key'). aggregate ['min', np.median]) meerdere functies
- df.groupby('key'). apply (eigen functie)
- df.pivot\_table ('col', index = 'col', columns = 'col')

Functies op DF

- df.col [i + % \* x] df.col
- df.substract (df.col, axis = 0) => op kolommen
- None / np.nan
- isnull(), notnull() = bool, alles niet null
- dropna(), fillna (method = 'bfill', axis = 0/1)
- Samen voegen
- pd.concat (df1, df2), pd.merge (df1, df2)
- df1.append (df2)

1104

Regel

se.str.contains ('r') => bool the.startwith()

re.findall ('r', data)

- Vd, VD = Any getal, geen getal
- Is, IS = Any spatie, geen spatie
- lw, LW = letter, geen letters
- Ab = alles behalve newline

+ \* ? = It, @+, !oto ~~keer~~  
 {3,3} = 3x  
 {3,3+} = 3 of groep  
 | or ^ = niet iets  
 (....) groep



## Grouping & Pivot

planets.groupby(kolom1).sum()

df.groupby('kolom1')

- aggregate(['min', np.median, max]) → drie methoden
- aggregate({'data1': 'min', 'data2': 'max'}) → methode per kolom voor elke kolom
- filter(filter\_func) → filteren of gem. groter is dan na van een cat.
- transform(lambda x: x - x.mean()) → houdt zelfde shape
- apply(norm\_by\_data2) → alles delen doorsom van kolom data2

df.groupby(mapping)

mapping = {'A': 'rowel', 'B': 'consonant', 'C': 'consonant'}

df.groupby([str.lower(), mapping]).mean() → hierarchical indexing

L = [0, 1, 0, 1, 2, 0]

df.groupby(L) → index 0, 2, 5 samen

titanic.groupby(['sex', 'class'])

['survived'].aggregate('mean').round(2)

= titanic.pivot\_table('survived', index='sex', columns='class')

indexen cols kunnen 5 den zijn voor hierarchical

decade = 10 \* (planets['year'] // 10)

decade.name = 'decade'

planets.groupby(['method', decade])

['number'].sum().stack().fillna(0)

pd.qcut(titanic['fare'] by col of index doet numerieke variabelen in tweeën

pivot\_table(agg\_func = {'survived': sum, 'fare': 'mean', margins = True voor totals

data L '2015'] werkt op dataframe  
str.lower(), str.<TAB>, str.capitalize  
Als 'info': ['BJCID', ...] dan maakt  
df['info'].str.get\_dummies('1') dan  
dummy variabelen

str.contains('Breakfast').sum()

dummy\_df[dummy\_df.isin(frequent\_sections)]

pd.get\_dummies(Categorical Series).dtype('bool')

for file in tqdm\_notebook(namelist):

fs.archive.open(file) root = tree.getroot()

tree = etree.parse(f) root.findall('handelingen')

.set\_index('file', inplace=True)

df['n\_sections'].idxmax()

.value\_counts → histogram

.quantile(0.99)

& → intersection | → union

^ → symmetric diff.

pd.DataFrame.from\_dict(dict, orient='index')

kr.sort\_values(by=['kolom'], inplace=True,

ascending=False)

data[data.notnull()] → Null of None

## Theorie

Omgaan met grote bestanden:

- 1 Niet verspreiden, gebruik file (get voor regx)
- 2 Clean memory als je klaar bent
- 3 redelijk conversie over versch computers.

Z-score:

$$z = (x - \mu) / \sigma$$

Missing values:

- 1 Invullen met constant (NaN of None?)
- 2 Vervangen door gem (numeriek) of modus (categoriek)
- 3 Vervangen door random waarde
- 4 Vervangen door inferred waarde

Outliers herkennen:

- 1 histogram
- 2 scatterplot
- 3 |z-score| < 3

## Regex

soup.findall('Affiliation')

port.Regex(Name, text)

[a-zA-Z] een van a-z

+ twee taken

\* een of meer

? nul of meer

. nul of een

{3} {3,6} 3, 3-6

^ niet of start string

\$ einde string

. elk karakter

\s, \s (geen) whitespace

\d, \d (geen) cijfer

(...) wat je wilt captureren

re.findall, split, sub

\w, \w (geen) sfer of letter

\b word boundary

df1.intersection(df2)

pd.merge(df1, df2, how='inner')

'outer' ↓ intersection  
union



# Pandas

Lees csv file in

```
pd.read_csv(bestand, sep = ';', usecols = ['...', '...']) = df
```

sort op index of kolom

```
df.sort_index(axis=0, ascending=T/F, inplace=True)
df.sort_values('kolom', ascending=True...)
```

kolom als index

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

wieke waarden tellen

```
df['kolom'].value_counts()
```

Dataframes aan elkaar plakken

```
s = pd.concat([df1, df2], join='inner', axis=1)
```

specifieke kolom van specifieke rij

```
peiling.loc['D66']['zetels']
```

```
peiling['zetels'].loc['D66']
```

rij met meer dan 25 zetels

```
zetels = pdf.loc[df['zetels'] > 25]
```

## Data Cleaning

Naam columns

```
df.columns = ['...']
```

check voor null values

```
df.isnull()
```

Drop rows with null values

```
df.dropna(axis=1)
```

replace null values with mean

```
df.fillna(s.mean())
```

## Groupby

selecteer specifieke waarde in kolom

```
nl = df[df['original-lang'] == 'nld']
```

5 meest vertaalde boeken uit nl

```
meest = nl.groupby('target-lang').max()
```

```
['num-transl'].nlargest(5)
```

naar hier types van

```
nst = set(*for x in meest.items())
```

hoeveel mannen en vrouwen waren aan boord

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

```
titanic['sex'].value_counts()
```

## Pivot tables

```
pd.pivot_table(df, index=['sex'], columns=...
```

```
titanic.pivot_table()
```

beveel mensen zaken er per geslacht

```
titanic.pivot_table(index='sex', columns=...
```

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

gemiddelde leeftijd van mensen die het wel

of niet hebben overleefd

```
titanic.pivot_table(index='sex', columns=...
```

```
np.median])
```

```
import pandas as pd
import seaborn as sns
import re
```

```
from tqdm import tqdm
tqdm.notebook
```

Hoeveel nummers beginnen met sex?

```
print(prince.text.str.contains(r'^[Ss]ex\b').sum() of
prince.text.str.startswith('Sex').sum() + prince.text.startswith('sex')
```

Hoeveel nummers bevatten sex maar niet letterlijk sex als woord

```
ef = prince.text.str.findall(r'[Ss]ex\b')
print([ef[ef.str.len() > 0].count()])
```

Regex \* = 0 or more + = 1 or more ? = 0 or 1 {2,3} = exactly two

{2,5} = between 2 and 5 {2,} = 2 or more [ab-d] = one character of a,b,c,d

[b] = backspace |d = one digit |s = one space |w = one word | = start of string

\$ = end of string |b = word boundary . = any character except newline a = character a

| = escapes special character. monte.str.lower() monte.str.len()

String operations: str.capitalize() = hoofdletter monte.str.startswith('T')

extract() = matched groups, findall() = for each element, replace() = replace pattern with str

contains() = boolean, count() = count occurrences of patterns, rsplit() excepts regex

extract first name: monte.str.extract('([A-Za-z]+)') expand = False. find the names

extract last name monte.str.split().str.get(-1) |w = [a-zA-Z0-9-]

re.search = match pattern anywhere in string (pattern, string)

Wikipedia bestand per regel inlezen with open(bestand, encoding='utf-8') as f:

```
for line in f:
    for l in tqdm.notebook(df)
    split regels: gesplit = l.split('\t')[1:]
    dictje = {}
    for x in gesplit:
        weer = x.split(' ')
        dictje[weer[0]] = int(weer[1])
    for x, y in dictje.items():
        iets[x] += 1
        iets[x] += y
    for iets in dictje.keys():
```

handige overzichten: dir(str) / dir(pd) / dir(np) dict, DF uit dictionary: pd.DataFrame.from\_dict(orient='index')

from Collections import Counter

correlatie = corr1 = df['k'].corr(df['i'])

maak dict van lijsten lege-dict {key: lijst} = [lijst, lijst]

reset index df.reset\_index() Series = 1 kolom uit df

Voorbeeld inlezen en letters tellen bestand = princelyrics.csv def lettertel(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
```

tellingseries = pd.Series(telling).sort\_values(ascending=False)

return tellingseries

50%

wanner precies is de variabele 'adult-male'

waar titanic1 = titanic [titanic ['adult-male']]

titanic1.pivot\_table(index='age', values=...

['adult-male'], aggfunc = sum)

preprocess: prince = pd.read\_csv(bestand)

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

regulair expressions bestand inladen:

bestand = open('prince\_lyrics.csv', encoding='utf-8')

Hoeveel nummers bevatten het woord sex. (met of zonder capital)

print(prince.text.str.contains(r'[Ss]ex\b').sum())

print(prince.text.str.startswith('Sex').sum() + prince.text.startswith('sex')

print([ef[ef.str.len() > 0].count()])

from collections import Counter

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

Regex \* = 0 or more + = 1 or more ? = 0 or 1 {2,3} = exactly two

{2,5} = between 2 and 5 {2,} = 2 or more [ab-d] = one character of a,b,c,d

[b] = backspace |d = one digit |s = one space |w = one word | = start of string

\$ = end of string |b = word boundary . = any character except newline a = character a

| = escapes special character. monte.str.lower() monte.str.len()

String operations: str.capitalize() = hoofdletter monte.str.startswith('T')

extract() = matched groups, findall() = for each element, replace() = replace pattern with str

contains() = boolean, count() = count occurrences of patterns, rsplit() excepts regex

extract first name: monte.str.extract('([A-Za-z]+)') expand = False. find the names

extract last name monte.str.split().str.get(-1) |w = [a-zA-Z0-9-]

re.search = match pattern anywhere in string (pattern, string)

Wikipedia bestand per regel inlezen with open(bestand, encoding='utf-8') as f:

```
for line in f:
    for l in tqdm.notebook(df)
    split regels: gesplit = l.split('\t')[1:]
    dictje = {}
    for x in gesplit:
        weer = x.split(' ')
        dictje[weer[0]] = int(weer[1])
    for x, y in dictje.items():
        iets[x] += 1
        iets[x] += y
    for iets in dictje.keys():
```

handige overzichten: dir(str) / dir(pd) / dir(np) dict, DF uit dictionary: pd.DataFrame.from\_dict(orient='index')

from Collections import Counter

correlatie = corr1 = df['k'].corr(df['i'])

maak dict van lijsten lege-dict {key: lijst} = [lijst, lijst]

reset index df.reset\_index() Series = 1 kolom uit df

Voorbeeld inlezen en letters tellen bestand = princelyrics.csv def lettertel(bestand):

key	Data
A	3
B	5
C	7



**magic** dir(pd), dir(np)  $\approx$  magic  
 with open('Bestand') as file:  
 for page in file:

read\_csv()  
 df.sort\_index()  
 df.sort\_values('kolom')  
 df.set\_index('kolom')  
 df.value\_counts()  
 pd.concat, pd.merge  $\rightarrow$  df1.merge(df2)

**Selection**  
 df[col]  $\rightarrow$  returns series  
 df[['col1', 'col2']]  $\rightarrow$  nieuwe df.  
 df.iloc[0,:]  $\rightarrow$  first row

**Cleaning**  
 df.columns = ['a', 'b', 'c']  $\rightarrow$  rename  
 pd.isnull(), pd.notnull() df.dropna()  
 df.dropna(axis=1)  $\rightarrow$  drop columns  
 df.fillna(x), so.astype(float)  
 s.replace(1, 'one')

**Filter**  
 df[df['col'] > 0.5]  
 df[(df['col'] > 0.5) & (df['col'] < 0.7)]  
 df.sort\_values(['col1', 'col2'], ascending=[True, False])  
 df.groupby(col)  
 df.groupby(['col1', 'col2'])  
 df.groupby('col1')['col2'].mean()/sum/etc.  
 df.groupby('col1').agg(np.mean)  
 df.apply(np.mean), df.apply(np.max, axis=1)  
 pd.pivot\_table(index=col1, values=[col2], aggfunc='mean')

**pd.crosstab?**  
 pd.DataFrame.from\_dict?  
 $\rightarrow$  Dict of Dicts to dataframe  
 pd.DataFrame(dict)  
 pd.Series(dict)  
 titanic['sex'].value\_counts()  
 titanic.pivot\_table(index='sex', columns='class', values='survived')  
 pd.crosstab(titanic['sex'], titanic['class'])  
 index.name = 'naam'  
 titanic.groupby('sex')['survived'].mean()  
 de dubbele namen geven een group frame  
 titanic.groupby(['sex', 'class'])['survived'].agg(gate('mean')).unstack()  $\rightarrow$  zorgt voor mooi frame  
 ge = pd.cut(titanic['age'], [0, 18, 80])  
 titanic.pivot\_table('survived', ['sex', 'age'], 'class')  
 pd.ccut(titanic['age'], 2)

**Regex** dir(str):  
 df[kolom].str.findall(r'blabla').str.len()  
 len(), lower(), extract(), contains(), split()  
 Prince['text'].str.contains(r'sex').sum()  
 $\rightarrow$  alles behalve '\n' (newline)  
 \w  $\rightarrow$  [a-zA-Z0-9\_]  $\rightarrow$  begin overgestelde  
 \b boundary between words and not words  
 \s single white space  
 \S non white space character.  
 \t, \n, \r  $\rightarrow$  tab, newline, return.  
 \d  $\rightarrow$  [0-9] \D  $\rightarrow$  alles behalve  
 ^ = start, \$ = end  
 + = one or more occurrences  
 \* = zero or more  
 ? = match 0 or 1 of thing to the left  
 {abc} = a or b or c [^abc] alles behalve  
 re.sub(regex, vervanging, string)  
 endswith(), str.startswith()

**Plotting**  
 df.plot(x=y, kind='line, bar, hist, barh, box, scatter, pie), subplots=True or False  
 $\rightarrow$  per kolom subplot  
 use\_index=True  $\rightarrow$  sticks for x axis  
 title='string', grid=False, legend=False  
 logx=False, logy=False, loglog=False,  
 xticks=[values for x axis], yticks=[]  
 sort\_columns=False  
 frequent\_sections = [lijst met labels]  
 dummy\_df = df.sections.apply(pd.Series).stack()  
 dummy\_df = pd.get\_dummies(dummy\_df)  
 dummy\_df = dummy\_df[frequent\_sec].groupby('page-id').sum()

**magic**  
 chance\_dict = dict()  
 for b in frequent\_sections:  
 chance\_dict[b] = dict()  
 for a in frequent\_sections:  
 if a == b:  
 chance\_dict[b][a] = 1  
 else:  
 kans = sum(dummy\_df[b] == 1)  
 chance\_dict[b][a] = len(dummy\_df[(dummy\_df[a] == 1) & (dummy\_df[b] == 1)]) / kans  
 conditionele kans: kans op A en B als je weet dat B er is.  $A+B/B$   
 recall = TP / (TP + FN) je bent ziek + test goed  
 precision = TP / (TP + FP) test is positief + je bent ziek  
 languages = [l2.split(',')[0].int(l2.split(',')[1])  
 for l2 in l1.split(' ')]







**PIVOT TABLES**

df.pivot-table(index='...', columns='...', values='...', aggfunc='...')

- je kan ook twee indexen meegeven
- aggfunc = mean by default

**Hoever mannen en vrouwen aan boord?**

titanic.pivot-table(index='sex', columns='class', values='survived', aggfunc=Sum)

**Hoever mensen de ramp hebben overleefd per geslacht, per klasse.**

① titanic.pivot-table(index='sex', columns='class', values='survived', aggfunc=Sum)  
 titanic.pivot-table(index='sex', columns='class', values='survived', aggfunc=len)

**Gemiddelde én mediaan per geslacht van overlevende (niet-overlevende)**

titanic.pivot-table(index='sex', columns='survived', values='age', aggfunc=[np.mean, np.median])

**Aantal 'woorden' in 'original' kolom**

raw['original'].str.count('woord').sum()

**Selecteren van waarde op basis van 2 kolommen**

df.loc['kolom 1']['kolom 2'] of df.loc[kolom1, kolom2]

**Series maken**

pd.Series(data='df[...]')

df.sum(axis=0) → som per kolom

df.sum(axis=1) → som per ry

**Pandas subsets**

- df[['col1', 'col2']] columns
- df[['col1']] single column as series
- df.loc['index-of-row'] row
- df.loc[['index1', 'index2']] rows
- df.loc[['rows']][['cols']] rows & columns
- df.loc['row1', 'col1'] row & column

**Bestand per regel inladen**

with open(bestand) as f:  
 for l in f:

**Detective werk**

Prince = pd.read\_csv(bestand)  
 Prince['text'] = prince['text'].astype(str)  
 prince['text'].str.contains('...').sum()  
 prince['text'].str.starts with('...')  
 prince['text'].str.findall('...')

**5 ways to count men and women**

titanic.sex.value\_counts()  
 titanic.groupby('sex')['survived'].count()  
 {s: titanic[titanic.sex==s].sex.count() for s in titanic.unique()}  
 titanic.pivot-table(index='sex', values='survived', aggfunc=len)  
 pd.crosstab(titanic.sex, titanic.survived).sum(axis=1)

**Regex**

- \n Newline
- [...] Range or character class
- [^...] not " "
- . Any character except newline
- \w word character
- \W nonword character
- \b word boundary
- \B non word boundary
- i case insensitive matching
- (...) group subpattern into 1 2

- \d digit character
- \D non digit character
- \s whitespace character
- \S non whitespace character
- ^ the start of the line of text
- \$ the end of line of text
- \* match 0 or more items
- + match 1 or more items
- ? match 1 or 0 items
- {n} match exactly n items

**Regex methods**

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

**Defentie**

df.teldeletters(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  
 t\_series = pd.Series(telling).sort\_values(ascending=False)  
 return t\_series

carsdf.year.value\_counts().plot(kind='bar').sort\_index

By value-counts plot → sort-index

**min-max normalization**  
 (df.min()) (df.max - df.min)  
 z-score  
 (df.col - df.col2).mean() / df.col.std()

**Plots**

outliers vinden  
 df.plot.scatter(x=col1, y=col2, title=...)  
 Regressie lyn → data=df, fit-reg=True  
 sns.displot(df)  
 df.plot(kind='bar')

**Missing Values**

Replace by:  
 - Constant → creation of info  
 - mean / mode of field → overoptimistic confidence levels  
 - Random variable → spread should remain closer to original  
 → spread will be reduced

**Big files**

with import gzip  
 with gzip.open(bestand) as bla:  
 for line in bla:

**Filling missing values**

df[col].fillna(df[col].mean(axis=0)) → mean  
 df[col].fillna(random.random())

→ explicit index  
 data.loc[0]  
 → implicit index  
 Python index.

pd.read\_csv(bestand, index\_col='...')

.set\_index()  
 .sort\_values(by=...)

df.columns = [...]  
 → rename columns

- .value\_counts()
- .min()
- .max()
- .sum()
- .mean()
- .mode()
- .median()
- .abs()

pd.series

df[['col1', 'col2']].corr()

df[col].mean()  
 df[col].mode()  
 df[col].median()

df.sort\_values  
 df[col].min/max/sur

.astype(str/float/lin)

df['kolom'].nlargest(r, columns=...)

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

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

pd.crosstab(col1, col2)

**Min-max normalization**

$X^* = \frac{x - \min(x)}{\text{range}(x)}$   
 $\frac{x - \min(x)}{\max(x) - \min(x)}$

→ is by min waarde gelyk aan 0

**Z-score**

$X^* = \frac{x - \text{mean}(x)}{SD(x)}$

→ is by mean waarde gelyk aan 0.

Z-score is niet geschikt om outliers te identificeren omdat mean & SD worden beïnvloed door outliers.

Be careful with using correlated variables

- at best → overemphasize one data component

- at worst → unstable & unreliable results.

OF series  
 .count()

**Str methods**

- len()
- lower()
- upper()
- index()
- strip()
- startswith()
- endswith()
- isupper()
- isnumeric()
- split()





HEZELBURCHT

# Data Science CheatSheet

kolom selecteren

`DF[["...", "..."]] , DF["..."]`

rij selecteren

`DF.ix[["..."], ...]`

} `DF.ix[["...", "..."], ["...", "..."]]`  
→ kan ook zonder list als enk.  
→ put this in `pd.Series()` to get a series object

Series naar DF => `pd.DataFrame({"col1": seri, "col2": ser2})`

dict naar DF => `pd.DataFrame(my_dict, key = index)`

`DF.set_index("name")`

`DF.sort_values(by = "name")` | `DF.sort_index()`

`DF.add_suffix(str)`

`DF.value_counts()` => telt hoe vaak waarde voorkomt in Series

`DF.describe()` => count, mean, std, min, 25, 50, 75, max

↳ categorical series => count, unique, top, freq

`DF.plot(x = column, y = column, kind = "...")` ← plotten!!!

`pd.pivot_table()`, `pd.crosstab(row, column)`, `DF.groupby()`

↳ margins = True => column & row with "All"

aggfunc = len => ipv count

get string from Series => `values[index]`

remove items from serie1 not in serie2 => `serie1[~serie1.isin(serie2)]`

get index with max value in Series => `serie.idxmax()` (also column)

loop door bestand => with `open(file) as f:`

dataset seaborn inladen => `sns.load_dataset("...")`

excel file lezen => `pd.read_excel`

csv file lezen => `pd.read_csv(".", sep = "...")`

`dict.items()` => k, v

`dict(zip(list1, list2))` => dict van 2 lijsten

precision =  
 $TP / (TP + FP)$

Als index tuple is, naar multi ↴

`pd.MultiIndex.from_tuples(index)`

recall =  
 $TP / (TP + FN)$

join 2 DF's `DF1.join(DF2)`

intersection = &  
union = |  
symmetric diff = ^

str → change to string  
[ss] = one or the other  
\b = blank space  
list to string = ".join(list)"

DF.column.contains('r').sum()  
starts with, ends with  
findall  
list to string = ".join(list)"



df.ndim / shape / size / dtype

x[2,0]=12  
 x[:,2]  
 x[1:2]  
 x[1:-1]  
 x[:,0] first col  
 x[0,:] first row

pd.read\_csv(loadfile(), sep='\t')  
 with open(file) as f:  
 sns.load\_dataset() usecols=[]  
 pd.read\_excel()

import re - default dict  
 import tqdm import tqdm\_notebook  
 from

np.array(x, z=1)  
 # extract odd nrs.  
 np.concatenate([x, y, m])  
 np.vstack([x])  
 np.hstack()  
 np.array + 5 = ufuncs  
 np.abs(x) / np.absolute()  
 np.sum() / min() / max()  
 np.any() / np.all()  
 np.argsort()  
 np.sort(axis=0/1)  
 np.log2(x)  
 2\*log(8) = 2^3=8  
 np.percentile(series, q=[25, 50, 95])  
 np.argwhere(series > 0)

view under outliers  
 df.loc[df[x] < df[x].quantile(0.95)]

plot  
 plt.hist(name)  
 df.head(10).plot(kind='barh')  
 df.plot.scatter(x, y, fit\_reg=True, style='o')  
 sns.distplot(np.log(x))  
 df.plot(kind='line', loglog=True, use\_index=False)  
 histogram => value\_counts()

create set of tuples  
 set(zip(df.x, df.y))

REGEX.  
 re.match() contains()  
 re.extract() count()  
 re.findall() split()  
 re.replace() rsplit()  
 word.split() .title  
 .startswith .capitalize  
 .endswith .swapcase  
 .replace  
 strip() .istitle()  
 (strip()) .isspace()  
 rstrip() .replace('W', '')

TRANSFORM  
 pd.DataFrame(data)  
 df.to\_frame()  
 df.squeeze  
 df = pd.DataFrame({'col1': ser1, 'col2': ser2})  
 df = pd.DataFrame(ser.values.reshape())  
 ind A & ind B = intersection  
 A | B = union  
 A ^ B = symmetric difference

pd.crosstab(df.x, df.y).sum(axis=1)  
 df = pd.DataFrame.from\_dict(x, orient='index')  
 df.name.unique.size  
 df.pivot\_table(index=, values=, columns=, aggfunc= [np.mean, np.median])

m = re.search('<p>.\*?</p>')  
 m.group  
 re.sub(x, y, str)  
 re.split  
 [] set of chars \* zero or more  
 . any char + one or more  
 ^ startswith {2} exactly 2  
 \$ endswith | either..or

data['a': 'c']  
 data[['a', 'c']]  
 data.loc[data.x > 2, ['p', 'c']]  
 index[0]  
 df.dropna(axis='columns')  
 df.fillna(0)  
 pd.concat([x, y])  
 df1.append(df2)  
 pd.merge(df1, df2, on='employee', left\_on='emp', right\_on='emp', left\_index=True, suffixes=(['b'], 'a'))  
 df1.join(df2)

# relative tessen x en y  
 pd.crosstab(df.x, df.y)  
 with open(file) as f:  
 c = 0  
 e = defaultdict(int)  
 for l in tqdm\_notebook(f):  
 ....  
 c += 1

precision = TP / (TP + FP)  
 recall = TP / (TP + FN)  
 df.column.contains(r'').sum  
 aggfunc(ipv count)  
 string from series values [index]  
 serie.idxmax()

AGGREGATIONS na.groupby  
 count() min() max()  
 first() std() var()  
 ast() mad() prod()  
 mean()  
 median() sum()

df.groupby('key').aggregate(['min', 'max'])  
 (data1 = min, data2 = max)  
 margins = true  
 first row & iloc[0]

df.column.dropna(subset=['age'])  
 r' \b[A-ZA-Z0-9-]'  
 prince\_text.str.contains(r'\b[ss]ex\b')

df.groupby('key').aggregate(['min', 'max'])  
 margins = true  
 first row & iloc[0]

percentage winners:  
 prince\_text.str.findall(r'[AEIOUaeiou]').st

prince\_text.str.contains(r'\b[ss]ex\b')



```
Readfiles | help(Pd.Series.loc); dir(pd/np); dir(str);
```

```
df = Pd.read_csv('file.csv')
Pd.read_excel('file.xlsx')
with open('file') as file:
    for page in file:
df.sort_values(by='col')
df.sort_index()
with open('file') as f:
    for line in f:
```

```
Regex findall; match; extract;
df.text = Prince.text.astype(str)
bevat sex -> prince.text.str.contains(r'\b[ss]ex\b')
begin met sex -> prince.text.str.startswith('sex')
of... startswith('sex').sum(1) -> startswith('sex').sum()
songs met woorden waar sex in voorkomt
et = prince.text.str.findall(r'[ss]ex\w+')
et df.str.len() > 0 -> laet zien of count()
Counter([s for l in et.values for s in l])
woorden met letter encyfer erin
... findall(r'\b[A-Za-z0-9]*[0-9]+[A-Za-z0-9]*[A-Za-z][A-Za-z0-9]*\b
```

- df.replace(1, one)
- df.dropna(subset=)
- df.describe()
- df.corr(col1, col2)
- df.min/max/std
- S['col'].str.lower()
- df.head str.replace
- df.tail
- df.shape
- df.info

```
regels = line.split('\t')
languages = {i.split(',')[0]: int(i.split(',')[1])
              for i in regels}
if sum(languages.values()) >= 5 and len(languages) >= 5:
    languages = {language: value for
                  language, value in languages.items()
                  if int(value) >= 2}
for language, value in languages.items():
    editors[language] += 1
    edits[language] += value
```

```
klinters vinden: findall(r'[AEIOUYaeiouy]')
-> any char except line break
\w -> word char, digit, \_ -> period
\s -> white space
\d -> digit 0-9
+ -> one or more
* -> zero or more
? -> once or none
\t|\n|\r -> tab, newline, return
[abc] -> a/b/c [^ABC] -> alles behalve
re.sub(reg, vervanging, string)
```

```
Pd.DataFrame.from_dict(dict, orient='index')
Zamenwaejen = Pd.concat([dataf1, dataf2], join='inner', axis=1)
merge = df1.merge(df2, how='inner', ...)
selecteren
```

```
from pandas.util.testing
import
test if equal -> assert_frame_equal(df1, df2)
stack() -> level 1 -> stack col to row
level 0 -> upper cols stack to row
unstack() -> row to col
cross tab
Pd.crosstab(df.col1, df.col2) zet de als
groupby([col1, col2])[col3].count().unstack()
grouping -> Pd.crosstab(df.col, [df.col1, col2])
query -> df.query("col1 < col2 < col3")
-> and, or; col1 in [val, val, val]
```

```
df[df['col'] = 'voorwaarde'].sort_values('col1', ascending=True)
df.loc[row, col] -> df[:, col] of gewoon df['col']
integer: df.iloc[0] -> df.iloc[[0]]
by label -> df.loc[[0], 'colname']
by label/posit -> df.ix[[index], 'colname']
bool index -> s[~(s > 1)] = Series where val is not >
s[(s < -1) | (s > 2)] = s where val is < -1 or > 2
df[df['population'] > 1200000] -> Filter df
df.loc[:, (df.any())] -> sel col any
.all() -> sel col all
```

```
tellen | len(df[df['col'] = 'value']) | df.groupby('value').size()
count_values() | df.pivot_table(index, val, aggfunc='len')
dictcom = {x: df[df.col = x].col.count() for x in df.col.unique()}
Pd.crosstab(df.col1, df['col2'])
sex class
df.groupby(['sex', 'class']).count().unstack()
hoeveel deel overleefd van sex per class
Titanic.pivot_table(index='sex', columns='sex', values='survived')
gem en medi van geslacht van overlevers - niet overlevers
Titanic.pivot_table(index='sex', columns='survived', values='age', aggfunc
rijen met nan eruit:
df = Titanic.dropna(subset=['age'])
```

```
concat(df1, df2), axis=0, 1
append: df1.append(df2)
Pd.merge(df1, df2)
merge on index -> left_index=True
error plot vrg 4 w 7
error -> df.Percentage - df.Percentage (aa:
Percs = df.Percentage
Percs.plot, bar(yerr=error)
gemiddelde aantal woorden per sectie
mean = (df['n-words']/df['n-sec'])
log = np.log2(mean)
pbt = sns.distplot(log)
```

```
telde letters (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
tellingseries = Pd.Series(telling).sort_values()
return tellingseries
telde letters (bestand).plot(kind='bar', figsize=(15, 7), logy=True);
```

```
Pivot fill-value, aggfunc=np
Pd.pivot_table(df, values, index, columns)
-> spread rows into columns
aggfunc applies to value, columns
are optional -> segment
```

```
boven 20
df.loc[col['x'] > 20]
gem zonder nan df.mean(skipna=True)
gender count
{p: Counter(genderlist[p][:n]) for p
in genderlist}
Pv25 zet is, hoeveel vrouw
Counter(mvdict('pv25')['female
```



	Index as ordered set:	REGEX
Series: index	x = 1, 3, 5, 7, 9	^ start
Series: values	y = 2, 3, 5, 7, 11	\$ end
DataFrame: columns	→ x & y : 3, 5, 7	\w \d \s
data.items() → tuples	→ x   y : 1, 2, 3, 5, 7, 9, 11	\w \d \s
data['a': 'b'] → includes b.	→ x ^ y : 1, 2, 9, 11	[abc] [^abc]

~ → not  
 defaultdict  
 ↳ from collections

data.loc[1:3] → index-based  
 data.iloc[1:3] → python-based  
 data.values[0] → first col.  
 data.loc[data.dew > 100, [x, y]]  
 a.add(b, fill\_value=0)  
 → sub(), mul(), div(), mod(), pow()  
 • notnull() → all non nan values.  
 • dropna / fillna()  
 MultiIndex.unstack() → back to df.  
 df.set\_index([arr]) → multi index  
 df.append(df2)  
 pd.merge(df1, df2, on = '...')

elem.findtext  
 f = archive.open(file)  
 root = etree.parse(f).getroot()  
 node.find  
 node.findall

BS.find(" ").text  
 BS.findall

• query("query") → @var  
 • sort\_values() by =  
 • mean(axis='columns') → mean over rows  
 • count() → number of items.  
 pivot-table('survived', index='sex', columns='class')  
 pivot-table(index='sex', columns='class', aggfunc='sum')

5 ways:  
 • value\_counts()  
 • groupby('').count()  
 • pivot-table aggfunc=len  
 • crosstab().sum(axis=1)

col.str.match → bool  
 • extract → matched groups  
 • findall → list  
 • replace  
 • contains → bool  
 • count → occurrences of pattern.  
 • split  
 • rsplit  
 • get\_dummies → indicator var split.

df.col.str.contains().sum()  
 str.replace(old, new)  
 Counter.most\_common()

• sort\_index()  
 • sort\_values()  
 • crosstab  
 hist  
 df.col.value\_counts().plot.barh()

index\_col sep.  
 filteren met bools:  
 df[filters.values]

df.iterrows()  
 pd.merge  
 pd.read\_excel(file, usecols='A:B')  
 df.col.quantile  
 pd.DataFrame.from\_dict(dict, orient='index', columns=[])



```

(df.col1 == 'test').sum() True False
df[(col1 == 'test') & (df.col2 > 5)] df
df[['col1', 'col2']] DF of 2 cols
df.values DF in list without index
list.items for error to many values to unpack
df.sort_values('col1')
pd.Series(df.col1, index=df.index)
df.drop(columns=['col1', 'col2'], axis=1)
pd.DataFrame.from_dict(dict1, orient='index', columns=['col1'])
pd.concat([dict1, dict2], axis=1, join='inner')
df1.join(df2)
((abs(df[col1] - df[col2])) / df[col1]).mean() * 100 mistake %

```

```

df.mean()
df.median()
df.mode()
df.count()
df.first()
df.last()
df.min()
df.max()
df.std()
df.var()
df.mean()
df.sum()

```

week 6

```

df.loc?
df.col1.quantile(q=0.99)
df.no.na = df.dropna()

```

week 7

```

df.fillna(0, inplace=True)
df.col1 = df.col1.astype(int)
df['percentage col'] = round((df.col1 / (df.col1 + df.col2) * 100), 0)
df['manner'] = df.zetels - df.vrouwen

```

Tentamen 1

- df = pd.read\_csv('xxx.csv', index\_col='names')
  - df.sort\_index(inplace=True)
- df = df.fillna(0)
- df['ass mean'] = df[['ass1', 'ass2']].mean(axis=1)
- (df[['ass1', 'ass2']] > 9).sum().sum()
- df['w mean'] = 0.4 \* df[['bas mean']] + 0.6 \* df[['deelman']]
- z.df = (df - df.mean()) / df.std()
- df.sort\_values(['col1', 'col2'], ascending=(false, True))
- corr = df[['col1', 'col2']].corr()
  - corr = corr[corr != 1]
  - rij = corr.max().idxmax()
  - kolom = corr.loc[rij].idxmax()
  - veroop = (rij, kolom)

socia geensoa

voors TP	FP
voors FN	TN

Hoofdstuk 3

```

df.loc[3, :2]
df.isnull()
pd.merge(df1, df2)
df.unstack()
df.groupby('col1').aggregate([min, np.median, max])

```

```

titanic.groupby(['sex', 'class'])['survived'].aggregate('mean').unstack()
titanic.pivot_table('survived', index='sex', columns='class', aggfunc='mean')
titanic.pivot_table('survived', index='sex', columns='class', aggfunc='mean')
age = pd.cut(titanic['age'], [0, 18, 80])
titanic.pivot_table('survived', index=['sex', 'age'], columns='class', margins=True)
births.pivot_table('births', [births.index.month, births.index.day])
see = pd.Series(['Jan smit', 'John hoas'])
see.str.lower.len.str.startswith(str).get.capitalize

```

```

df.Percentage.plot('bar', yerr=(df.Percentage - df.PercBar).tolist())
r = (Q3 - Q1) / 1.5
Q1 - x en Q3 + x = data
df.col1.fillna((df.col1.mean()), inplace=True)
for bin in freqsec:
    chancedict[b] = dict()
    valuedict = {}
    for a in freqsec:
        numab = len(dummysdf[dummysdf[b] == 1] & (dummysdf[a] == 1))
        valuedict[a] = numab / len(dummysdf[dummysdf[b] == 1])
    chancedict[b] = valuedict
    condchance = pd.DataFrame(chancedict).sort_index(axis=1)

```

```

titanic.groupby(['sex'])['survived'].count()
titanic.groupby(['sex', 'class'])['survived'].count()
titanic.pivot_table(index='sex', values='survived', aggfunc='mean')
pd.crosstab(titanic.sex, titanic.survived).sum(axis=1)
titanic.groupby(['sex', 'class'])['survived'].mean().unstack()
titanic.pivot_table(index='sex', values='survived', columns='class', aggfunc='mean')
pd.crosstab(titanic.sex, titanic.class).aggregate([np.mean, np.std])
test = titanic.dropna(subset=['age'])
test = test[test.sex == 'male']
pd.crosstab(test.age, test.adult_male)

```

```

bestand = 'prince_lyrics.asu'
def teldeletters(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
    tellings = pd.Series(telling).sort_values(ascending=False)
    return tellings

```

```

teldeletters(bestand).plot(kind='bar', logy=True)
prince = pd.read_csv(bestand)
prince.text = prince.text.astype(str)
A: hoeveel nummers bevatten "sex"
prince.text.str.contains('b[ss]ex\b').sum()
B: hoeveel nummers beginnen met 'sex'
prince.text.str.startswith('sex').sum()
C: hoeveel nummers bevatten sexy etc?
sex = prince.text.str.findall(r'[ss]ex\b')
print(sex[sex.str.len() > 0].head(10))
print(sex[sex.str.len() > 0].count())
D: Print al die woorden uit
from collections import Counter
print(Counter([s for l in sex.values for s in l]))
E: U2 lu etc
ef = prince.text.str.findall(r'b[A-Za-z0-9]{0,9}[A-Za-z0-9]')
print(ef[ef.str.len() > 0].head(10))
print(Counter([s for l in ef.values for s in l]))
F: % links per track
(100 * prince.text.str.findall(r'[A-Za-z0-9]{1,4}')).str.len() / prince.text.str.len().describe()

```

```

[abc] anything abc | [^abc] anything not abc | a-z A-Z 0-9
\d whitespace | \d digit | \w word | \b word boundary
() capture | ^ start of string | $ end of string | Regex
a? zero or one | a* zero or more | a+ one or more

```

A = alcoholist B = man	$\frac{x - \min(x)}{\max(x) - \min(x)}$	min-max normalization
$P(A B) = 0,0225$	$\frac{x - \text{mean}(x)}{\text{std}(x)}$	z-score normalization
$P(A B) = \frac{P(A \cap B)}{P(B)}$		
Conditional P(B)		
Kans = $\frac{0,0225}{0,5} = 0,045$		

```

freqsec = sectietitels(df).nlargest(20).index
dummydf = df.set_index(freqsec).stack()
freqsec = [x for x in freqsec.values]
dummycols = pd.DataFrame(columns=freqsec)
for col in freqsec:
    dummycols[col] = dummydf
dummydf = dummycols == freqsec
dummydf = dummydf.groupby(level=0).sum() == 1
chancedict = dict()
Re.findall Re.sub
Re.split Re.match if beginning with (rf)

```



NaNs: replace with constant  
 \* replace with mean/median with pandas imputed value  
 \* changes center and spread

Miscellaneous: value counts  
 outliers: histogram, boxplot, scatterplot  
 center measures: mean, mode, median  
 Variability: range (max-min), STD, MAD, IQR  
 Normalization (reduce variation):

• min-max:  $X^* = \frac{X - \min(x)}{\text{range}(x)}$   
 alles 0-1  
 • z-score:  $Z = \frac{(x - \text{mean}(x))}{SD(x)}$   
 < mean - > mean +  
 Kom beiden tot n-verd, if skewed.

Decimal:  $X^* = \frac{X}{10^d}$  d = decim. max(x)  
 • Skewness:  $3 \cdot \frac{(\text{mean} - \text{median})}{SD}$   
 right mean > median  
 Transformation (reduce skewness)  
 • log:  $\sqrt{x}, 1/\sqrt{x}$

outliers: if  $z_{sc} > 3$  < -3  
 • z-score: ~~z-score~~  
 • IQR: Q1: 25% Q2: 50 (med)  
 Q3: 75%  
 IQR = Q3 - Q1 → Q1 - 15 IQR  
 Q3 + 15 IQR

outlier if x located  
 1.5 IQR below Q1, above Q3  
 num → cat: binning (k-means e.g.)  
 cat → num: dummy var.

unsupv.: clustering/grouping  
 reduce complex  
 supv: regressie, classificatie  
 3: reinforcement learning  
 evaluatie: vergl. met gouden std.  
 class.: Acc, Prec, Rec, MRR ...  
 Regn.: RMSE, AIME ...

Acc: how many predicted were correct?  
 Prec: how many pos. pred were con TP/TP+FN  
 Rec: how m. pos pred. as pos? TP/TP+FN

Test +	TP	FP
-	FN	TN

alle ziek gezond

Train, Val, test  
 weinig data: k-fold  
 • men/women titanic:  
 + sex, value\_counts()  
 + groupby('sex')(survived).count()  
 + pivot\_table(i=sex, v=surv., agg=len)  
 + crossstab(+sex, +surv).sum(axis=1)  
 • per geslacht / per klasse (doel surv.)  
 + p\_t(1:sex, v=sum, c:class)  
 + crosstab(+sex, +class)(survived).agg('mean')

• gem./men. vestijn sum = 1  
 + p\_t(1:sex, ~~agg='sum'~~  
 c = sum\_val: age, agg=(np.mean, median)  
 • tel de letters:  
 from coll.: import defaultdict  
 with open('file', encoding='utf8') as f:  
 for l in f:  
 for c in l:  
 tel[c] += 1  
 tel = pd.Series(tel).sort\_values ...  
 Return tel

• Regex: str.match → bool  
 • extract → strings  
 • findall - replace  
 • contains → returns bool  
 • count → count occ. of part.  
 [abc]: inv: [^ABC]  
 \s: whitespace |d: digit  
 \w word (a|b) a n b  
 a?: 0/1 a  
 a\*: 0+ a a+: 1+ a  
 a{3}: 3 start of str  
 \$ end of str \b wordbound.  
 \* ? non greedy

• n numbers - sex -  
 Prince.txt.str.contains('b[Es]ex|b').sum()  
 • n numbers starts with - sex -  
 • contains('^[Es]ex|b')  
 • n numbers bevat - sex -  
 • contains('ESs]ex|w+).sum()

• Print alle woorden die - sex -  
 P(counter(x for v in P for x in Y))  
 • Perct. klinkers:  
 findall('aeiouy').str.len()/pr.txt.str.len()

• (Ris)kolom • Axis 0: within each col  
 1: each row.  
 langste str: 1. lang = str()  
 2. pak regel, meet len() vergl. len(lang)  
 3. if >, lang = 2 → next

• Age = pd.cut(titanic.age, [0, 18, 35])  
 + p\_t(sum, ['sex', 'age', 'class'])  
 • con: df1.con(df2)  
 • Jaccard  $\frac{A \cap B}{A \cup B} \rightarrow A \cup B = \frac{(A+B) - A \cap B}{A \cap B}$   
 • conditionele kans  $P(A|B)$   
 =  $\frac{P(A \cap B)}{P(B)} \rightarrow$  waarschijnlijk A als B is  
 P(B) → totaal aantal B's

temp\_dict[a] = len(df[df[a]=a])  
 df[df[b]=1] & (df[df[a]=a])  
 len(df[df[a]=a] & df[df[b]=1])  
 chance\_dict[b] = temp\_dict[a] / len(df[df[a]=a])  
 cond\_chance = pd.DataFrame(chance\_dict)  
 cond\_chance.sort\_index(axis=1)

use rows dat n plaats vindt, gey.  
 dat Book ...  
 P(choc sal.) \* P(womw)  
 P(mouw)

regressie: sns.lmplot(x, y, data, fitreg=T)  
 df → multi. indx → set\_index(cod1, cod2)  
 • Concat: row-wise default  
 short: df1.append(df2)  
 • Merge: join autom. on key or cols.  
 merge(df1, df2, on='col1')  
 left on, right on ...  
 default on indexes → df1.join(df2)

• Pivot-tables (values, index, cols, agg)  
 • groupby. aggregate (col1: mean, col2: no)  
 pd.qcut → quantiles  
 Pivot-tables can be multi leveled.  
 Regex: re.findall(regex, text) → list  
 • match → determine if begin string  
 • split → list  
 • sub → replaces

f-5 = sectionals (df) n largest rows  
 dummy\_df = df.sectionals.agg(pd.Series).stack  
 cols = [x for x in frequent\_sectionals.value  
 df2 = pd.DataFrame(columns=cols)  
 for x in f-5:  
 df2[x] = d\_df  
 Pcol = df2 = temp - sec  
 dum\_df = pd.DataFrame(glob(level=0/1)  
 chance\_dict = dict() == 1  
 for b in frequent\_sectionals:  
 temp\_dict[a] = len(df[df[a]=a])  
 df[df[b]=1] & (df[df[a]=a])  
 len(df[df[a]=a] & df[df[b]=1])  
 chance\_dict[b] = temp\_dict[a] / len(df[df[a]=a])  
 cond\_chance = pd.DataFrame(chance\_dict)  
 cond\_chance.sort\_index(axis=1)