

I python  
 %time  
 %timeit  
 %ls  
 %prun  
 ?  
 %rm  
 %cd  
 %mkdir

np.reshape (rows, columns)  
 np.mean  
 np.max, np.min  
 np.std  
 np.sum  
 np.size, np.shape  
 np [rows, columns] ~~shape~~  
 np [first: last: step]  
 np.indices  
 np.transpose (array)  
 np.random.shuffle  
 np.invert  
 np.zeros, np.ones  
 np.count\_nonzero (value)  
 np.random.rand

df.argmax  
 df.index  
 df.groupby()  
 df.value\_counts  
 df.sort\_values (ascending=, [ ]  
 df.plot (kind='bar')  
 df.drop  
 df.dropna  
 df.reset\_index()

pd.crosstab (rows, columns)  
 pd.pivot

100 man?  
 10 zieh 100mal g, nicht  
 90 nicht zieh

Series.values() = np.array

\*df.index.weekday / month / day  
 df.merge / df.join?

echt  
 zieh nicht zieh  
 Test zieh g  
 nicht zieh 1 g

\*str.contains  
 \*str.len  
 \*str.replace

1000  
 echt  
 zieh nicht zieh  
 Test zieh g  
 nicht zieh gg  
 990-99 891

% ls magic

Shell met ! (pwd, ls, cd)

? access info

- \* - wildcard options

dot  $\sum_i (s_i \cdot S_i)$

• evaluation

Classification

→ accuracy

recall

$TP / (TP + FN)$

Precision

$TP / (TP + FP)$

RMSE = regression

NP → dir(np)

np.empty/zeros → array

np.shape/ndim/size

→ [i]  
↓  
[j]

np.info(np.ndarray.dtype)

math aggregate

.exp →  $e^x$   
.sqrt →  $\sqrt{x}$   
.sum  
.min  
.std  
.max

np.arange (start, stop, step)

[::-1] → backwards

[::2] → from 0 with step 2

a[2] → el. at 2nd index

a[1,2] → 2nd row, 3rd col

a[0:2] → items at 0 and 1

• boolean indexing

a[a < 2] → alles in a kleiner dan 1

• fancy

b[[1,0,1,2], [0,1,0,1]]

↳ (1,0), (0,1), (1,0), (2,1)

• transpose

np.T

np.transpose()[1]/[2] → add rows/cols

np.concatenate()

arr[i]=4 (maak 4 van index i)

np.add(arr1, arr2) → elementwise

(arr, 4) → 4 by elk element

pd → dir(pd) → help(<sup>pd.</sup>Series, loc)  
pd.read\_csv('file.csv', header, rows)  
df['a'] → get one element  
df[':'] → get subset  
Position → df.iloc[0],[0] + select  
by label → df.loc[0],[country]  
by index → df.ix[2] → single row/sub  
→ df.ix[:, ['Country']]  
bool select  
s[s > 1] → s where s is not  
s[(s < -1) | (s > 2)] → values under -1  
& above 2  
sort\_index()  
sort\_values(by=)

• Statistics

df.describe → Summary

df.mean → mean of all cols

df.corr → corr between cols

df.count → value count col

df.min/max/mean/median/std

• Cleaning

df.columns['a','b','c'] → rename cols

pd.isnull() check null, bool array

df.dropna() → drop row nan values

df.fillna(x) → fill nan with x

(s.mean()) → fill na with mean

s.replace([1,3], [one, three])

↳ replace all 1 → one 3 → three

s['col'].str.lower()

str.replace('vb', '7')

delete all met ...

s[s['col'].str.contains('...')]

• Filter, Sort, Group

df[df[col] > 1] → row where col val > 1

df[(df[col] > 1) & (df[col] < 2)]

df.sort\_values(col, ascending)

df.groupby(col)

[col1], [col2]

df.pivot\_table(index=col, values=col2,3, aggf)

→ group by col, aggfunc on 2,3

• Join / combine

df1.append(df2) → add row to end df1

.Concat(df1, df2, axis=1)

→ add cols to end df1

df.join(df2, on=col1, how=inner)

• View, inspect

df.head

df.tail

df.shape

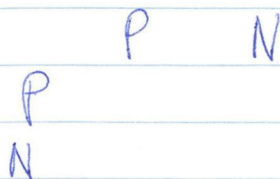
df.info

df.value\_counts

len(.unique) → unique len

vrg → len('Party') == 'pvdd'

- np slices and indices are views not copies so changing them changes original array, use `.copy()` for a copy.
- Broadcasting rules:
  1. If the two arrays differ in their number of dimension the shape of the one with fewer dimensions is padded with on its leading (left) side.
  2. If the shape of the two arrays does not match in a dimension, the array with shape equal to one in that dimension is stretched to match the other shape.
  3. If in any dimension the sizes disagree and neither is equal to 1, an error is raised.
- `x[::-1]` all elements, reversed
- Count amount of True entries in a Boolean array with `np.count_nonzero` (True = 1, False = 0)
- Quickly check any/all values are True with `np.any` / `np.all`
- Data science is: Based on ~~the~~ patterns in data predicting a value or new unseen instances of data (classification/regression). Steps: interaction with outside world, preparation, transformation, modeling and computational presentation. Data Science is the study of the generalizable extraction of knowledge from data.  $\text{Precision} = \frac{TP}{TP+FP}$   $\text{Recall} = \frac{TP}{TP+FN}$   
 $F_1 = 2 \times \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$
- `np.nditer` to iterate through numpy arrays  
`.value_counts()` `.index.values` `.str` (`['\n']`)  
`x[~x.y.str.contains(" " , na=False)]`  
`.value_counts().sort_index()` `x.loc[x['x'] == '']`  
`normalize = True * 100` voor percentages



row column

X[1, 2]

X[2, 3] # two rows, three columns

X[3, :2] # all rows, every other column

X[1, :-1] reverse

its divisible by 3  
X[X % 3 == 0]

X[:, 0] first column

X[0] first row

if  $(X-x) \cdot \text{mean}$   
 $(X-x) \cdot \text{std}(\text{head}) / x \cdot \text{std}(\text{head})$

masking

df[df > 0, 3] (cdf[k, g, 8])

X.std() std() X.std()

df.loc = explicit indexing

df.iloc = python style indexing

# sort rows by height than width

res.sort\_values(['height', 'width', 'ascending=False', True])

res.groupby('species')['length'].max()

X.sort() of X.argsort

sort indices

K nearest 2:

k=2 z=np.argsort(x.shape[0])

for i in range(x.shape[0]):

for j in z[i, h+1]:

plt.plot(2p(xL, j) \* x[i], color='black')

def shape (n=len(x), y):

return L.reshape(np.ma.count(y), 1)

all others much ~~are~~ occupy 2d array

C = pd.Series(C).sort\_values(ascending=False)

C = Counter(x.split())

(S, I) = (S \* C).index.sort\_values(ascending=False)

contains('row') / df.columns.stk

df.columns [df.columns.stk

(m, n)]

[df['varname']] stk: std.cult

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

df[['varname', 'varname2', 'varname3']]

help(len) of len?

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

df.columns

Data science cheat sheet

np

% ls magic = all magic commands

`X [:2, -3, ] = last 3 columns first 2 rows` | `X * + np.arange(len(X))`

`X.reshape(12,)` = make one dimensional | `X [0, 1, 2, 3].copy(1, 2, 3, 3)` - pang index

`(X - X.mean())**2` `sqr dif. each elem from mean of X` hand wr. std vs function

`np.sqrt((X - X.mean()).mean())` = `X.std()` = ~~standard~~

`X [X % 3 == 0] = apply bool. mask`

`X [0, 1, 2, 3].copy(1, 2, 3)` diagonal (pang index)

`big-arr = np.random.randint(10**3, size=10**3)`

`% timeit sum(X**2 per x in big-arr)`  
`% timeit sum(X**2 per x in big-arr)`  
`% timeit (big-arr)`

`X [:, -1] = last column 2D array` | `X [X % n == 0] | X.mean(axis=0)`  
`dubbel = df.index.value_counts()`  
`dubbel [dubbel > 1]`

pd

`df.sort_index()`

`df.index.str.lower() .str.strip() .str.replace(' ', '') .str.replace(' ', '')` so

`df [df.index.str.lower().str.contains('string')] .head(10)`

`df.describe`

`df.income.sort_values.plot()`

`df [df.income > 50000].gemean().value_counts().plot()`  
`kind = 'barh'`

= gementes met scholen met inkomen > 50.000, gesorterd op hoeveelheid in

`hor. bar plot`

`top_minis = kur [kur['party'] == 'ppp'] .Ministerie.value_counts()`

`top_minis [top_minis >= 5].plot(kind='barh')`

`df.groupby('party') .count()`

`deelvragen = kur.vraag.str.count('!')`

`no_outliers = deelvragen [deelvragen <= 50]`

`no_outliers.plot(kind='hist')`

`np.mean(no_outliers)`

`median = statistics.median(no_outliers)`

`modus = Counter(no_outliers).most_common(1)[0][0]`

`corr = no_outliers.corr(kur.vraag.str.len())`

`corr = no_outliers.corr(kur.Aanwoord.str.len())`

`df.sort_values(['sepal-len', 'sepal-wid', 'petal-len', 'sepal-len'])`

`df.groupby('species') ['sepal-len'].max()`

`df.query('sepal-len > 5 and sepal-wid < np.sqrt(5)')`

`C = Counter(text.split())`

`CS = pd.Series(C).sort_values(ascending=False)`

`CSCS [CS.index.str.len() == 21].index`

```

import numpy
import matplotlib.pyplot as plt
import random
array
os
math
import matplotlib inline

```

```

np.random.seed(42)
randint(10, size=6) 1d
randint(10, size=(3,7)) 2d
randint(10, size=(4,9,6)) 3d
np.concatenate([X, Y])
np.vstack
np.hstack
np.split([3, 5]) = tot 3/4 5/van 6
np.vsplit([X, [3]]) vanaf rij 3

```

indexing: id = X[0] reverse van's reverse alle rijen tot kolon  
 2d = X[0,0] / X[:, -1] / X[:, -2] / X[:, 2]  
 magic: %ls - the directory  
 %row = remove  
 %col = double line

```

x = np.arange(1, 5)
np.add.reduce(x) = returns it on all elements
np.multiply.outer(x)
np.linspace(start, end, spacing)
np.arange(start, end, step)
a = np.random.random(42)
a.randint(10, size=(3, 9)) -> a < 6 - a > 6 -
np.sort(x, axis=0/1)
np.argsort(x) = values of position
np.partition(x, 3) = 3 smallest values + remaining

```

reshaping: X [np.newaxis, :] = row vector  
 X[:, np.newaxis] = column vector  
 computation =  $2^x - x/2 - x//2 - np.log - np.log 2 - np.log 10(x)$   
 Broadcasting: val i = padded to left  
 2 = 1 dimension array stretched  
 3 = if dimension disagree = error  
 m + A[:, np.newaxis].shape = right padding  
 access simple = X[0]  
 slicing = X[:, 5]  
 masking = X[X > 5]  
 fancy indexing = [X[3], X[5], X[6]]  
 fancy + simple = [X[2, 0, 1]]  
 fancy + fancy = [X[2, 1, 3, 2]]  
 fancy + masking = mask = np.array([0, 1, 0, 1])  
 X[rows[:, np.newaxis], mask]

Theorie:  
 Precision = # of predictions for class C covered  
 Recall = # instances of C predicted to be C - TP / (TP + FN)  
 accuracy = # of all predictions being correct - TP / (TP + FP)  
 accuracy = 95% n=10000  
 prec + 95 495 590  
 rep + 5 905 910  
 precision 95 / (95 + 95) = 0,16

Pandas - good for data analysis and modeling without use of R  
 df = pd.read\_csv('...'), names = ['kolom', 'human', 'runder']  
 df = pd.read\_html('...')[1]  
 %ls - ./Data/\*.csv  
 %cat - ./Data/foe.csv -> df = pd.read\_csv('.. /Data/foe.csv')

- plot samples in dataset = iris.species.value\_counts().plot(kind='bar');
- create dataframe numerical data = x and y = X, y = iris.drop('species', axis=1).iris.species
- z-normalize by subtracting mean by std = ((X - X.mean()) / X.std()).head(1), ((X - X.mean()) / X.std()).std()
- Sort iris by column and row inverse = iris.sort\_values(['...'], ascending=[False, True])
- maximal lengths for each species = iris.groupby('species')['length'].max()
- larger than 5 and smaller than sqrt of 5 = a = np.sqrt(5) - iris.query('x > 5 & y < @a')
- tokenize the text by split, count taken, show head = C = Counter(text.split())
- Find all tokens which have >= 2 = CS = pd.Series(C).sort\_values(ascending=False)
- what percentage of all unique words occur once = CS[CS == 1].sum() / CS.sum()
- make histogram which indicates for each number i - occur how many tokens in text occur i times  
 pd.Series(Counter(CS.values)).plot(kind='bar'); CS.value\_counts().sort\_index().plot(kind='bar')

```

data = pd.read_html('...')[1], data.groupby('...').count().sort_values().tail(1)
data = pd.Series([1, 2, 3]), index=['a', 'b', 'c'], dict = {'ams': 100, 'utc': 50}
pop = pd.Series(dict)

```

Series: pd.Series(5, index=[3, 6, 9])  
 pd.Series({'a': 1, 'b': 2}) -> index = [3, 1]  
 states = pd.DataFrame({'pop': pop, 'area': area})  
 states

```

indexing: data['b'], 'a' in data, list(data.items())
slicing: data['a': 'b'], data[0:3], data[(data > 0) & (data < 3)], data[['d', 'b']]
data.values = alle waarden
data.T = transpose
data.loc[data.x > 100, ['pop', 'area']] = masking + fancy
A.add(B, fill_value=0)

```

slicing in 2d = pop[['ams', 2010], ['utc', 2000]]  
 hard = select all values from 2010 -> pop[[i for i in pop.index if i[1] == 2010]] -> pd.MultiIndex.from\_tuples(index)

concatenate: pd.concat([ser1, ser2], axis=0/1)  
 merge: df3 = pd.merge(df1, df2)  
 categories of join = 1-to-1, many-to-one, many-to-many  
 need key column  
 choose by pd.merge(x, y, on='...', left\_on='...', right\_on='...')  
 pd.merge(x, y, how='outer', 'inner', 'left', 'right')

Lees de vragen goed! Kijk of je iets moet plotten!

iris = ... heeft niet verzameling set() a.tab → kkl

### NUMPY

1D array: `x2.reshape((12,7))`

`x1.std()` =  $\text{np.sqrt}(\frac{\sum (x1 - \text{mean}(x1))^2}{n-1})$   
Squared difference  
variance

diagonaal: `x2[[0,1,2], [0,1,2]]`

`%timeit sum([x**2 for x in bigarray])`

`%timeit sum(bigarray**2)`

- sum of squares in an array -

homogeen vs heterogeen. size = elementen ndim

### PANDAS

`index_col = 'kolomnaam'` die je als index wilt

• `unique()` of `set()`

`str.strip()` haakt whitespace aan begin/ende weg

`str.contains('str')`

`df[naam]`

of `df.describe()` → statistiek data

→ in elke waarde het gem. uitre waarde in die kolom afkruken (normaaliseren)

`mean = cito.mean()`

`cito[mean.index] = mean`

`cito.gemeente[cito.inkomen > 50000].value_counts()`

• `sort_values()`, `plot(kind='bar')`

→ in welke gem. staan scholen waar inkomen meer dan 50000  
`plot` aantal scholen per gemeente

• `drop(columns='naam')` of `(naam, axis=1)`

→ Z-normalize data by subtracting mean of column from each cell and divide by standard dev of each column  
 $(X - \text{mean}(X)) / \text{std}(X)$ , `head()`

→ also show std of each col of the Z-normalized data • `std()` → moet gelijk zijn aan 1

→ Sort iris first by S-L & then by S-W. Sort in reverse order (multindex op species)  
`iris.sort_values(['S-L', 'S-W'], ascending=[False, True])`

→ max S-L for each species  
`iris.groupby('species').sepal.length.max()`

→ Find all tokens in text that occupy 24 char in total  
`CS[CS.index.str.len() * CS.values == 24].index`  
(`CS = pd.Series(Counter(text.split()))`)

→ What % of all unique words in text occurs once

$\frac{CS[CS==1].sum()}{CS.count()}$   
of  $\text{len}(CS.index.unique()[CS.values==1]) / \text{len}(CS.index.unique())$

→ What % of all words in text occurrence

$\frac{CS[CS==1].sum()}{CS.sum()}$   
of  $\text{len}(CS.index.unique()[CS.values==1]) / \text{len}(text.split())$

→ Make a histogram (bar plot) which indicates for each number of occurrences how many tokens occur i times

`CS.value_counts().plot(kind='bar')`

→ Land met meeste bendidagleden + aantal leden

`bdEO.value_counts().head(1)`

of `bd.groupby(0)[1].count().sort_values().tail(1)`

of `winnar` → `bd.groupby(0)[1].count()`  
`winnar.argmax()`, `winnar.max()`

→ welke partij heeft grootste aantal leden 1 deelstaat

`bd.groupby(0)[2].value_counts().sort_values().tail(1)`

### IPYTHON

`%ls` files in directory } %magic  
`%cp` file file2 } %magic  
`%rm` no-good.txt

plod print work - array  
cd naam / change dir  
mkdir naam - make dir  
mv .. / . / move one dir up

`inspect resources: %timeit`  
• `berschil` → %time ? ! %time?

`datatype of var. x` → `x?` inspect 1st line of output

`get documentation: help(...)`

`print(...)`

Supress output;

% debug interactief

% mode hoe exceptions zien

### THEORY

precision:  $\frac{TP}{TP+FP}$  accuracy:  $\frac{TP+TN}{\text{total}}$  recall:  $\frac{TP}{TP+FN}$

	T	F	
T	95	495	590
F	5	9405	9410
	100	9900	10000

accuracy 95%  
zieke 1/100

precision =  $(0.01 \times 0.95) / ((0.01 \times 0.95) + (0.99 \times 0.05)) = 0.16$

xml/json: cel + lijst (als er niet evenveel kolommen zijn → spreadsheet)

teken spreadsheet:

- insert / delete / update

- andere view (string / ding)

- combineren v. sheets

- aggregate / correlate

dot prod: `A.dot(B)`

$\sum(a[i] * b[i])$  for i in range(len(a))

no births p/year p/gender

`names.groupby(['year', 'sex'])['births'].sum().unstack()`

### numpy hw wk2

• create array (n, 0):  
`return np.arange(0, n*0, 0)` eerste n getallen  
deelbaar door 0

• add-per-column(L)  
for i in range(L.shape[1]):  
`L[:, i] = L[:, i] + i` tel i op voor item  
in ide kolom  
return L

• add-per-row(L)  
`tp = np.transpose(L)`  
return `add-per-column(tp).transpose()`

• mean-per-column(L)  
return `np.mean(L, axis=0)`

alles reversed: `x[::-1, ::-1]`

1D → 2D `x.reshape((1,3))` of `x[0].newaxis, :]`

`x1, x2 = np.split(x, [3,5])` → `[:3]` `[3:5]` `[5:]`

`np.sum of L.sum()`

broadcasting - ele regels verbreden het om x1 naar de  
shape van x2 te krijgen.

how many vals less than zero: `np.sum(x < 0)`  
in each row: `np.sum(x < 0, axis=1)`  
`np.any()` `np.all()`

combined indexing: `x[2, [2, 0, 1]] = x[2,2], x[2,0], x[2,1]`

get name from last row `data[-1]['name']`

`data.loc` explicit → indices

`data.iloc` implicit

`data.ix` hybride vorm

`data.T` transpose

→ `zip` zegt dat ik heb  
heb ik het ook

## Numpy

- 1 Last 3 columns of first 2 rows:  $x[2, -3:]$
- 2 Change each element in  $x_1$  by its squared difference from the mean of  $x_1$ .  
 $(x_1 - x_1.mean())**2 = \text{variance}$
- 3 standard deviation = ~~np~~ sqrt of variance ( $np.sqrt(\text{variance}) == x_1.std()$ )
- 4 1D en 2D arrays vanwege broadcast regels niet samen te voegen, 1D kan niet naar 2D
- 5 boolean mask to  $x_1$  and reduce  $x$  to ints divisible by 3:  $x_1[x_1 \% 3 == 0]$
- 6 Fancy indexing for diagonal of  $x_2$   
 $x_2[[0,1,2], [0,1,2]]$

TP / TP + FP  
TN / TN + FN

## JPython

% ls -lh a = working directory ~~size~~ b = size

% rm = remove

% time, timeit, prun

% cp = copy previous, new

Pandas  
pd.read\_excel(" ", index\_col = " ")

dimensions = .shape

test voor dubbel = len(set(" ")) = len(" ")

cito.index.value\_counts() telt hoeraak index elk item in cito.index voor komt

>1 betekent dubbel (dubbel = cito.index.value\_counts() | dubbel [dubbel > 1])

sort\_index()

str.contains

~ = negation

↗ (cito[M.index] - M.tail())

↘ M = cito.mean()

cito[cito.inkomen > 50000].gemeente.value\_counts() 7,2

.sort\_values() ascending = [ 6, 8

.groupby .max() / .min()

~~np~~ a = np.sqrt(5)

iris.query('sepal\_length > 7.5 and sepal\_width < @a')

C = Counter(text.split()) = text splitten en tellen 6,2 + 0,7c  
pd.Series(C)

10, 10: 10, 60, 10, 10 60 0,6 60 14,9 7,9 7,45  
10 70  
5,5 10 70/2 = 0,9



① import re  
 import collections → (counter, most-common())  
 import math → math.sqrt()

**BASICS**

random.random() → rand  
 random.choice() → sel()  
 [start: stop: step]  
 L[-1:] → laatste rij  
 L[:, -1] → laatste kolom  
 \*  $g \log(a) = b \rightarrow g^b = a$   
 2  $\log(16) = 4 \rightarrow 2^4 = 16$

**Numpy**

**Pandas**

.log2(), .info(), .shape(), .inplace(), .size(), .full(), .eye(), .subtract(), .multiply(), .sin(), .cos(), .sqrt(), .corrcoef(), .transpose(), .mean(), .reshape(), .ravel(), .ones(), df[[]], min(axis=)

pd.read\_csv() → compression, skip initial space, sep, names=K]  
 df.loc[" "].isin() of df[df[" "] == " "].unique()  
 .crosstab(L, E) / .dropna(subset=L) / .notnull()  
 df[df[" "] == " "].str.contains(" ") / .str.lower()  
 df.pivot\_table(values, index, columns, aggfunc, margins=True)  
 df.pivot\_sort\_values(" ", ascending)  
 df.column.sum() / df.column.index.values()  
 pivot[pivot[" "] >= 10.000] @ ((pivot["ratio"] >= 4) & (max\_error = 1/15 \* 100))  
 act = sum(error...) / sum(["all"] \* 100)

plt.hist **PIT**  
 plt.parch  
 \* plt.pcol(K, E), df.plot.scatter  
 re.sub(r',', '', text) **Regex**  
 re.findall()  
 \w → word → [A-Z a-z]  
 \d → digit → [0-9]  
 \ → Escape spec. char.  
 \* ? → 0 of meer, 1 of meer  
 \ → any char  
 ( ) → capture { } → wantat

grid, shape() → [i-1 % " ], ...  
 midden → L [int(len(L)/2)-2: "+2, "-2: "+2]

rel\_error = hoogste error  
 most\_errors = meeste x post

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

f\_error = f.sum() + m\_error \* pivot.m.sum() / pivot.all.sum  
 → = act\_error

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

pd.head(), groupby(), .tail(), .pcol()

```
np.linspace(0, 1, 5) - 0.2, 0.4
```

```
Linear(axis=0)
```

```
np.concatenate([a, b], axis=1)
```

```
np.full((3, 5), 3.14)
```

```
np.all(x==6)
```

```
x[[2, 1, 0]] == 10
```

```
np.argmax / np.argmin index of min value
```

Rule 3: if any dimension the sizes disagree and neither is equal to 1, an error is raised

```
data.loc[data.density > 100, ['pop', 'density']]
```

```
np.nansum(), pd.merge(df1, df2, on='employee')
```

```
df.dropna(), pd.merge(df1, df2, left_on='employee', right_on='name')
```

```
.drop('name', axis=1)
```

```
pd.join is same as merge, alleen autom-op index
```

```
df.pivot_table(values='births', columns='name', index='year',
```

```
aggfunc=sum)
```

Madeken excel: schouwt niet, beperkte functionaliteit, integratie met andere tools

```
cross=crossstab(df.jaar, df.party)[party in [0, 10].index.values]
```

1. Panda helps filling in the gap with Python concerning data analysis and modeling, without having to switch to a more domain specific language like R

2. Python data analyze excels in performance, productivity and the ability to collaborate

```
np.arange(0, 14, 2), create_array(14, 7), np.reshape(L, (y, x))
```

```
laatste kolom = L[:, -1], L[begin: eind]
```

```
df.columns, len(df.index) geeft hoeveelheid rijen
```

```
index.value_counts
```

```
cito.index.str.lower(), str.strip(), str.replace(" ", "").sort_values()
```

```
cito[n.index.str.contains('school')], cito.describe()
```

```
cito['RMSE'] = np.sqrt((cito - cito.verwacht)**2), arr[:, -1]
```

```
np.random(0, 10, (3, 3))
```

```
x[1::2] every other elem, x[5::-2]
```

```
np.split(x, [3, 5]), splitten op 3 en 5
```

```
x < 3 geeft array, delta(i) loc [ ]
```

```
data.nansum, data.fillna(0), data.reset()
```

```
df.groupby('key').sum(), pd.read_csv/html
```

```
data.plot(x='', y='', kind='bar')
```

Kom met hypothese  
verzamel data  
schoon op  
structureer  
analyseer  
rapporteer

2de ingesloten  
gise- pps en klopt niet zelf en klopt  
ziet in font niet zelf en font

$$\frac{g}{g_0} \Big/ \frac{1}{1} \quad acc = \frac{TP + TN}{TP + FP + TN + FN}$$

g <sub>0</sub>	g <sub>1</sub>		
10	10	TP	TN
g <sub>0</sub>	g <sub>1</sub>	FP	FN

$$gg \frac{g}{g_0} \Big/ \frac{1}{1}$$