

NP. reshape (rows, columns)

NP. Shape

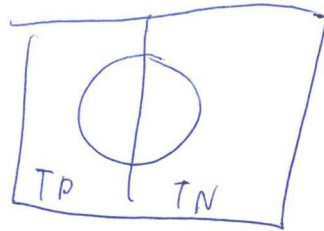
??? (rows, columns)

NP. size

??? (size)

NP. std

??? standard dev.



$TP / (TP - FP)$

pd. df. plot (kind = '...'): box, barh, hist, box, density, area, scatter, pie

pd. sort\_values (ascending = False)

pd. query ('column > 5')

x = 13 'column == @x'

c = Counter (text.split)

pd [index == 'James']

pd. groupby (by = 'column', axis = 0)

example:

df. groupby (['animal']). mean() # takes mean of each animal

# find highest number of 1 kind

df. groupby (0)[1]. count(). sort\_values(). tail(1) # of .aggregate()

pd. merge ()

pd. join

% time it \* returns time the line took on average

% time # quicker

%% time # for whole cell

NP. random. randint (0, 10)

?

% IS - ln

% rm # remove

% prun

x?

Data science is:  $precision = TP / (TP + FP)$   $F1 = 2 * prec * recall / (prec + recall)$   
 $acc = TP + TN + All$   
 $recall = TP / (TP + FN)$

Based on patterns in data predicting a value on new unseen instances of data.

Classification + predicting a class

Regression + predicting a value

Different from traditional statistics that integrate in which you try to explain something.

: Interacting with the outside world

2 preparation 3 transformation

4 modeling and computing 5. presentation

Informatic verwerking is:

- Omzetten ene formaat naar andere

- uitleggen wat er bedoelt is

impliciete informatie expliciet maken = leesbaar computer

onderzoek cyclus:

1. Kom met idee/hypothese 2. Verzamel data

3. Schoon data op 4. herstructureer data 5. analyse

6. Rapporteer.

Waarom geen excel?

1. Schakelt niet

2. Beperkte functionaliteit

3. Integratie met andere tools

Series: one-dimensional labeled

array capable of holding any datatype.

Data frame: two-dimensional labeled data structure

with columns of potentially different types.

```
CS = pd.Series(data = Counter(text.split()), sort_values(ascending=False))
```

```
iris[iris.species.str.len() > 5].sort_values('sepal.length', ascending=False)
```

```
iris.groupby('species')['sepal.length'].max()
```

```
iris.sort_values(['sepal.length', 'sepal.width'], ascending=False, species)
```

```
X, y = iris.drop(columns=['species'], iris.species)
```

```
iris.species.value_counts().plot(kind='bar')
```

```
np.sum(np.square(big_array))
```

```
sum([x**2 for x in big_array])
```

```
X1[X1%3==0] X2[X1%2, [0,1,2], [0,1,2]]
```

```
np.sort([(X1 - X1.mean())**2].mean()) = X1.std()
```

```
X2.reshape(12, 4) np.square(X1 - X1.mean())
```

```
X1[3:] X2[2,1:] X2.shape, X2.size
```

pd. "TAB"

pd.read\_csv?

\*?

help(collections)

dir(scipy) \* stats ~~extend, man,~~

\*\* power x4

% modulo

// floor

L + np.arange(0, L, step)

transpose = .T

L[i:i, i:i]

(Grid 1 = 3).sum() - 10

pd.crosstab(year, party) index.values

Kur[~Kur.party.str.contains('vragen')]

Kur.loc[Kur.party == 'pvv', 'ministerie'].value

index max = idx\_max iloc[[0,1]]

new column = Kur['text'] = ...

Pivot[[F,M]].misp(axis=1).odm()

Pivot[pivot.ratio < 5][F].sum

pd.groupby(0)[1].countsort\_values(1).text(1)

CS.value\_counts().sort\_index(1).plot(kind='bar')

CS[CS = [sum(1) / CS.count(1) (sum)]

CS[CS \* CS.index.str.len() = 24].index

CP, env, man, more

% mv = move file o.o. / textit.1

% mkdir = make new dir

% cd = change direct

% pwd = print work dir

X? % cat = new file

time / time it / paan

% rm = remove

% ls -lh = dir & size

# Theory

precision = positive predicted value.  
 measures how many of the samples predicted positive actually are positive.  $TP / (TP + FP)$

recall = how many of the positive samples are captured by the positive predictions.  $TP / (TP + FN)$

onderzoekscyclus: idee/hypothese / vraag. → data  
 Variabelen / opschonen → herstructureren / analyseren → pandas

Data science: based on patterns in data, & predict a value on new unseen data instances of data.  
 Steps: interacting with outside world, → preparation → (task) formation → modeling & computing → presentation.  
model excel → schakel met, bepaalde functionaliteit, integreer andere tools.

# NumPy

min/max uitbreiden  
 array for each new/col  
 sort - with axis  
 axis = 0: row  
 axis = 1: col

add / subtract / multiply / exp / sqrt / mod  
 vstack → np.vstack((A, B))  
 hstack  
 concat (Eg: A, B, C)  
 axis = 1

# array index:

array (E[1,2,3], [1,2,3]) → row 1 col 1, row 2 col 1 ...

slicing:  
 a[C0:4] → index 0,1,2,3  
 a[C0:4, 3] → row 0,1,2,3 at col 3  
 a[C:2] → index 0,1  
 a[C:, 1] → col 1 all rows

np.average(A[2]).reshape((3,4))  
 row = np.average(C[1,2])  
 col = np.average(C[1,2,3])

X[rows, cols]

table multiply:  
 X = np.average(C[1,6])  
 np.multiply.outer(X, X)

value array:  
 X[C[:, -1], :-1]

X[C[:, 2], :-3] → 2 rows 3 cols  
 X[C[:, 3], :-2] → 2 rows 2 cols  
 X[C[:, 0], :-1] → first col  
 X[C[0, :], :-1] → first row

# Pandas

→ X for x in df.iterrows()

count and plot:  
 df[column].value\_counts().plot.bar()

max value for species:  
 df.groupby('species')['sepal.length'].max()

Show top 10 sorted  
 df.sort\_values(ascending=False).head(10) of E[:10]

lower all in column:  
 df[column].str.lower()

replaces all low words with " ":  
 df[column].replace(to\_replace="\\w", value=" ", regex=True)  
 can also be value np.nan

Series.value\_counts()

df.sort\_index()

get length of index words:

Series.index.str.len()  
 pd.drop(axis=1)  
 - shape  
 - info  
 - count / - sum / - mean  
 loc[[0], [0]] → pos  
 i[0:2] → row  
 i[0, "col"]

3C ~ (5 > 1)  
 5C (5 < -1) (5 > 2)  
 select: 5C[1:5] = 6

df[df.col.isin(list\_of\_options)]  
Counter (df[col]) for df in records if 12 in v most-common.

pd.crosstab(df.col, df.col, margins=True)  
 pd.pivot\_table(df, index=col, values=count, aggfunc="sum", margins=True)

pd.Series.find\_all("class", class="class")  
 first column = 2C[0]  
 print (first column: path(A, B))  
 np.count\_zero(X < 6)  
 np.sum(X < 6, axis=1)



- Any Character
- \\d - (0-9)
- ID - Not
- \\w - (a-z, A-Z, 0-9, -)
- W - Not
- \\s - (space tab, newline)
- S - Not
- \\t - tab
- \\b - white space around words
- B - Not
- ^ - Beginning of str
- \$ - end
- ^[] - Matching characters NOT in []
- [abc] - a, b, or c (range)
- + - Match 1 or more
- ? - Match 0 or 1
- \* - Match 0 or more
- | - either, or
- {x} - expecting 'x' amount

```

p = re.compile('\\d+')
p.findall('string here')
re.split('[\\W]+', 'string')
pandas: kvl ~ not str.contains(), inplace=True
kvr.dropna(subset=['Party']) # removes all rows with NaN values
pd.crosstab(kvr.Jaar, kvr.Party) [Topic value count of party index values]
pd.read_csv('file', header=None, sep='\\t', names=[' '], skipinitialspace=True)
DataFrame.idxmax(axis=0, skipna=True) # for each column
Numpy:
np.random.randint(2, size=n)
np.arange(start, stop, step) # output is array
L.reshape(int1, int2) # 1D to 2D
np.mean(L, axis=0) # gemiddelde, axis 0 is horizontaal, 1 is verticaal
L.sum(axis=0) # vertical axis, 1 = horizontal

```

```

DataFrame.plot(kind='bar')
df.describe()
df.size # amount of elements
df.shape # output (rows, columns)
df.sort
df[['columnname']]
or df.columnname
df.sort_values('columnname', ascending=False)
slice columns: df[[' ', ' ']]
kvr[Party], str.replace('woord', 'iets')
deelbaar = L[L % n == 0]

```

Steps:

1. Interacting with the outside world
2. Preparation
3. Transformation
4. Modeling & computation
- Connecting your data to statistical models, machine learning algorithms or other computational tools.
5. Presentation

- Data science is the study of the generalizable extraction of knowledge from data.
- A ~~common~~ epistemic requirement is predictive power.
- A data scientist requires skill set: Mathematics, machine learning, AI, statistics, databases, optimization, and a deep understanding of the craft of problem formulation to engineer effective solutions.
- \* Evaluate models: Classification & regression

Precision = how many of the predictions of class Cc were correct. (Dus de hoeveelheid false beschuldigingen)

Recall = how many instances of class Cc predicted to be Cc.

↳  $TP / (TP + FN)$  = how many of the "true class" (those with disease) are picked by the test.

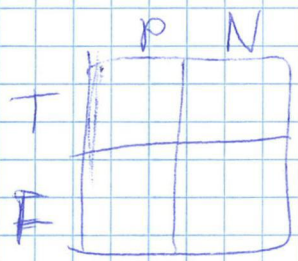
↳  $TP / (TP + FP)$  = how many of those with a positive test do have the disease?

Accuracy = how many of all predictions were correct.

- Pivot table is typisch een 3D Matrix
- Gebruik bij categoriale eigenschappen.
- Groepeer op een eigenschap (bijv. geslacht)
- Kies 1 kolom, maak groepjes voor elke waarde in de groep.
- Kies een kolom en valdies samen.

Waarom niet blijven Excelen?

1. Schaalt niet
2. Beperkte functionaliteit (niet makkelijke schoonmaken)
3. Integratie met andere tools. In python wel!





```

bs4 imp. BeautifulSoup
imp re
1 url = 'https:'
2 html_doc = requests.get(url)
3 soup = BeautifulSoup(html_doc, 'html.parser')
# get li elements
li = soup.findAll('li', class_='re.compile('class 2'))
c. attrs
Dict = {1: a.attr['tag']: 1: a.attr['href']} for 1 in li

```

```

local = './data/name'
doc = open(local).read()

```

```

pd.read_csv
df.loc[0, 0]
- select by row/column
df.loc
- select by row (column name)
df.ix [2] = row
[1, 'capital'] = column
df.drop
sort_index()
sort_values(by='country')
Shape (rows, columns)
columns
count = # non NA values
sum
min/max
idxmin (idx max) = index of min/max
describe = summary
mean/median

```

append (insert) (stable)

```

w. oil companies Comment Harvest
1 -> 2 -> 3
comments = soup.findAll('li', class_='comment')
com_list = [(soup.title.text, c.attrs['data-d'],
             c.find('span', class_='user-name').text, c.p.text)
            for c in comments if c.p]
comment DF = pd.DataFrame(com_list)
comment DF. columns = ['article', 'comment id', 'username', 'rea']

```

```

np.where(a, np)
np.arange(1, 2, 3)
zeros (row, col)
np.average(0, 15, 5)
a = np.array([1, 2, 3])
a. shape #dim int
ndim #dim
size #elem
dtype
astype = convert on type
len(a)
np.subtract
add
divide
multiply
Exp
sqrt

```

```

if both arrays:
a == b
[True, False, False]
[False, True, True]
a. sum
: min (axis = 0)
: cumsum (axis = 0)
: mean (median)
: corrcoeff
a [: :-1] = reversed array
: transpose

```

~~pd.read\_csv~~ np.random.randn(int)

pd.o.c - dropna  
 ↳ missing values replaced by None  
 ↳ drop ... where at least one element is missing (axis)

pd.a. fillna(0)  
 ↳ replaces NA with sub.

DF3 = pd.merge(df1, df2)

444: 0.9 x 0.9  
 99 x 0.9 + 1 x 0.9

```
import re | os | numpy as np | matplotlib.pyplot as plt | seaborn as sns
%matplotlib inline
```

```
from nose.tools import assert assert_equal, assert_almost_equal
from pandas.testing import assert_frame_equal
```

```
col_Names = ["Jaar", "Partij", "Titel..."]
kvr = pd.read_csv('linktofile.csv', sep='t', header=None, skipinitialspace=True,
names=col_Names)
```

```
1. cito = pd.read_excel('http...xls', index_col='naam')
```

```
2. print(cito.shape, cito.columns)
   cito.tail()
```

2. laatste 5 rijen

3. dubbele index?

```
(len(set(cito.index)) == len(cito.index))
```

4. ~~dubbel~~ welke zijn dubbel index?

```
dubbel = cito.index.value_counts()
dubbel[dubbel > 1]
```

5. orden.

```
cito.sort_index
```

6. maak schoon

```
cito.index.str.lower().str.replace(' ', '')
```

```
cito.index.str.lower().str.strip().str.replace(' ', '').str.replace(' ', '').
```

7. haal cito de scholen zonder 'school' in naam

```
sort_values()
```

```
cito[~cito.index.str.lower().str.contains('school')].head(5)
```

8. add column 'RMSE' ~~toen~~ (column 'quasicito' en 'verwacht')  
(kwadrant van verschil tussen)

```
cito['RMSE'] = np.sqrt((cito.quasicito - cito.verwacht)**2)
```

9. elke waarde, gemiddeld verande kolom aftrekken, op die manier alle waarden normaliseren

```
M = cito.mean()
```

```
(cito[M.index] - M).tail()
```

10. plot inkomens per school

```
cito.inkomen.sort_values().plot()
```

```
cito[cito.inkomen > 50000].gemeente.value_counts().sort_values().plot(kind='bar');
```

11. welke gemeentes inkomen > 50.000, hoeveel per gemeente, horizontale barplot.



# numpy

create array linear sequence  
np.arange(0, 20, 2)  
[0, 2, 4, 6... 18]

geef 1 dimensionale array van eerste n getallen getallen die deelbaar zijn door o.

return np.arange(n) \* o

accepteert 1 dimensionaal array met lengte van  $\sqrt{n}$  voor een geheel getal n, verandert in vierkant numpy array

return L.reshape((int(np.sqrt(len(L))), int(np.sqrt(len(L))))

L.shape = array dimensions  
L.size = number of elements in array

maak array met shape n, 1

return L.reshape(int(L.size), 1)

laatste kolom in array  
return np.array(L[:, -1]) for column in L

geef 1D array met alle getallen deelbaar door n

return L[L % n == 0] ← boolean mask

RMSE = np.sqrt((x-y)\*\*2)

gemiddelde per kolom  
return np.mean(L, axis=0)

add-per-column voor elk getal + hoeveelheid van index  
return L \* [0] for x in enumerate(L) + L

[10, 15], [20, 30] → [10, 15], [21, 31]  
return add-per-column(L.transpose()).transpose((1, 0))

shift tab by

change each element in x1 by its squared difference from the mean of x1  
 $x1 - x1.mean() ** 2 = variance$

np.sqrt((variance).mean(1)) == x1.std()

np.diagonal() = diagonaal

# Python noteboek vragen

lijst van magic functies = % ls magic  
% magic voor uitgebreide beschrijving

kopieer file naar andere file % cp  
verzameling ? = set()

Dir(np), Dir(pd) = lijst met functies  
create the dataframe ~~is~~ X, y = iris.drop('species', axis=1), iris.species

2 normalize the data in x  
(X - X.mean()) / X.std().head()

sort values

iris.sort\_values(['sepal-length', 'sepal-width'], ascending=[False, True])

What are the maximal sepal lengths for each species  
iris.groupby('species')['sepal-length'].max()

C = counter(text.split())  
CS = pd.Series(C).sort\_values(ascending=False)

what percentage of unique words occurrence.  
CS[CS == 1].sum() / CS.count()

# Pandas

1) Lees de file in een df en maak de naam van de school de index. Df = cito

cito = pd.read\_excel('... .xls', index\_col='naam')

2) geef de dimensies en de namen van kolommen  
print(cito.shape, cito.columns) | cito.tail()/head

3 Zitten er dubbel in de index? laat zien  
len(set(cito.index)) == len(cito.index)

4 laat zien welke namen dubbel voorkomen  
dubbel = cito.index.value\_counts()  
dubbel[dubbel > 1]

5 orden cito op index  
cito.sort\_index()

6 verwijders spaties aan begin en eind, telkens als " en andere rommel van de namen van index.  
cito.index.str.lower().str.strip().str.replace(' ', '').str.replace('"', '').sort\_values()

7 Beperk cito tot de scholen zonder de substring school in hun naam  
cito[~ cito.index.str.lower().str.contains('school')].head(10)

8 Voeg een nieuwe kolom 'RMSE' toe waarin de wortel... tussen de kolommen quasico en verwacht staat  
cito.describe() → handige lijst met berekeningen  
cito['RMSE'] = np.sqrt((cito.quasico - cito.verwacht)\*\*2)

10 plot de inkomens  
cito.inhomen.sort\_values().plot()

" In welke gemeentes staan de scholen met een inkomens van meer dan 50000. Hoeveel scholen zijn er per gemeente. Maak een horizontale barplot. Gemeente met meeste scholen bovenaan.  
cito[cito.inhomen > 50000].gemeente.value\_counts().sort\_values().plot(kind='barh');  
winner.argmax() → winner.max() → led.

groupby(0) [1].count().sort\_values() = winner land naam bondleider

Tokenize the text in variable text using split() count each token create a Series with token counts and sort with most occurring tokens on top



# PANDAS

`len(set(...)) = unique` | `soms()` achter functie | `shape = dimensies`  
`dubbel = cito.index.value_counts()` | `cito[~cito.index.str.lower().str.contains('school')]`  
`dubbel[dubbel > 1] < BM` | `cito[cito.mean().index] - cito.mean()`  
 iets(x) voor iets(y) = `Groupby('x')[y]` | `10..7..3 = .ascending[False, True]` | `CS[CS...]`  
`iris.query('column > 5')` | `RMSE = np.sqrt((x-y)**2)` | keren dat het voor komt \* len voor  
`C = Counter(text.split())` | `CS[CS==1].sum / CS.count()` = just once  
`CS = pd.Series(C).sort_values(ascending=False)` | `CS[CS==1].sum() / CS.sum()` = all words | `float(...)`  
`kvr[kvr.Portij == 'pvdd'].Ministerie.U.C()` | `kvr.drop(kvr[kvr >= 5].index)`

# Numpy

`Tafel van x = np.arange(n) * x` | `Last_c = np.array([C[-1] for c in L])`  
`Varray L = int(np.sqrt(len(L)))` | `% = L[L % N = 0]` | `diag = x.diagonal()`  
`verticaal = np.reshape(L, np.size(L), 1)` | `mean_p_c = np.mean(L, axis=0)`  
`add 0,1,2 = [x[0] for x in enumerate(L)] + L` | `std = np.sqrt(((x-x.mean())**2).mean)`  
 middle square:  
`r = int((L.shape[0]-4)/2)` | `sqrt diff from mean = (x1 - x1.mean())**2` | `np.array([1,3,6], dtype='float32')`  
`c = int((L.shape[1]-4)/2)` | `ndim = n. of dimensions`  
`return = (L[r:-r, c:-c])` | `n.shape = size each dimension` | `n.size = size total array`  
`big = np.random.randint(10**3, size=10**2)` | `every other element = x[::2] of x[1::2]`  
`comer = sum(x**2 for x in big)` | `x1, x2, x3 = np.split(x, [3, 5])` | `np.multiply(x, x)`  
`np = sum(big**2)` | `((x >= 1) | (x < 0.5))` | `np.argsort = sort indices`  
`↑ sum of squares`  
`np.partition(x, n)` | `2 log(8) = 3, 2^3 = 8, 2 log(16) = 4 = 2^4 = 16`

# Magic shell

Files in working directory, and how large = `%.Ls -Lh` | Remove = `?.rm`  
 datatype of x = `x?` | `!` = access shell

# Beautiful soup

`x = requests.get(url)` | `f.text.split('\n')` | `html_text`  
`line.split(s) for line in lines` | `".join` | `BeautifulSoup(url, 'xml')` | `re.sub(ok, oke, text)`  
`Spec tekens = [^a-zA-Z0-9]` | `BFS(Source file, 'HTML.Parser').find_all('page')`  
`Recall = TP / (TP + FN)` | `\w = alle woorden`  
`Precision = TP / (TP + FP)`

# Extra Pandas

`S.value_counts(dropna=False)` | `count = nr of rows. 3+3=2r.`  
`df.apply(pd.Series.value_counts)` | `sum = 3+3=6`  
`rename c = df.columns = ['a', 'b']` | `Party_10 = [ ... ]`  
`dropna(axis=1) = for columns` | `cross = pd.crosstab(kvr, jaar, kvr, p.) [Party_10]`  
`% = vragen.max() / sum(vragen) * 100`  
`Pivot_table(df, values='count', columns=['sex'], index=['name'], fill_value=0, margins=True, aggfunc=np.sum)`



# Numpy

**Slicing** `x[start:stop:step]`

10

`x[:5]` → first 5 elements  
`x[5:]` → elements after index 5  
`x[4:7]` → middle sub-array  
`x::x]` → every other element  
`x[1::2]` → " " starting at index 1  
 when step element is negative → reverse  
`x[::-1]` → All elements reversed  
`x[5::-2]` → reversed every other element from index 5

**Reshape** → andere vorm  
**concatenate** → arrays platten  
`axis=0` → kolommen  
`M.min(axis=0)` → kleinste waarde per kolom  
**fancy indexing**: passing an array of indices to access multiple array elements at once.

**Multi dimensional**  
`x2[2,:3]` → two rows, three columns  
`x2[3,::2]` → three rows, every other column  
`x2[:,1,::-1]` → rjen én kolommen omdraaien

**Accessing**  
`x2[:,0]` → first column  
`x2[0,:]` → first row  
`x2[2,2]` → extract 2x2

**Broadcasting**  
`a = np.array([0,1,2])`  
`b = np.array([5,5,5])`  
`a+b`  
 → array([5,6,7])

2 kolommen selecteren  
`df[['kolom1', 'kolom2']]`  
**Broadcasting**:

$\begin{bmatrix} 0 & 1 & 2 \end{bmatrix} + \begin{bmatrix} 5 & 5 & 5 \end{bmatrix} = \begin{bmatrix} 5 & 6 & 7 \end{bmatrix}$  `np.arange(3)+5`

$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 6 & 7 & 2 \\ 0 & 7 & 2 \\ 0 & 7 & 2 \end{bmatrix} = \begin{bmatrix} 7 & 8 & 3 \\ 1 & 8 & 3 \\ 1 & 8 & 3 \end{bmatrix}$  `ones((3,3))+np.arange(3)`

$\begin{bmatrix} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 2 & 1 & 2 \end{bmatrix} + \begin{bmatrix} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 0 & 2 & 4 \\ 0 & 2 & 4 \\ 2 & 2 & 4 \end{bmatrix} + \text{np.arange}(3)$  `np.arange(3).reshape((3,1))+np.arange(3)`

`np.sum()`  
`np.prod`  
`np.mean`  
`np.var`  
`np.std`  
`np.min`  
`np.max`  
`np.argmin`  
`np.argmax`  
`np.median`  
`np.any` → any True?  
`np.all` → all True?  
`np.size`  
`np.log`  
`np.split` → 1d  
`np.vsplit` → 2d  
`np.hsplit` → 2d

# Pandas

Series als dictionaries  
`pd.Series([lyste], index=[lyste], index=[lyste])`  
`loc`: allows indexing and slicing that always references the explicit index → key value achtig  
`iloc`: allows indexing and slicing that allows references the implicit index → echt indexen!

**Pivot table**  
`Pivot = pd.Pivot_table(df, index=..., margins=True)`  
 ↓ All

`df = pd.read_csv(..., index_col=...)`  
`df = pd.read_excel(...)`  
`lito[lito.index.str.lower().str.contains('achthoek')]`

`Pd.ien`  
`Pd.lstrip`  
`Pd.startswith`  
`Pd.split`  
`pd.index`  
`pd.replace`  
`pd.value_counts`  
`pd.plot(kind=...)`  
`pd.sort_values()`  
 ↓  
`ascending = [True, False]`  
`str.contains`  
`.sort_index()`  
`.describe()`  
`Pd.groupby`  
`pd.index_col = ...`  
`Pd.tail`  
`Pd.columns`  
`Pd.rows`  
`Pd.shape`

# Python notebooks

`% magic`      `% ls` → welke files in dir  
`% time`        `% cp` kopieer file  
`% timeit`

**DS steps**

- interacting with outside world
- Preperation
- Transformation
- Modeling & computation
- Presentation

**precision** = how many of the prediction of class C TP / (TP+FP) were correct  
**Recall** = how many instances of class C predicted to be C TP / (TP+FN)  
**Accuracy** = how many of all predictions were correct. TP / (TP+FP+FN)  
 werkt niet met:  
 - skewed distribution  
 - you can get a high reliability by the law of large numbers on the False class.

**Waarom geen excel?**  
 - schaal niet  
 - beperkte functionaliteit  
 - integratie met andere tools is moeilijk.  
**Series**: 1d labeled array capable of holding any datatype  
**Dataframe**: 2d labeled data structures



True positive	False Negative
False positive	True Negative

Accuracy

Accuracy =  $TP + TN / all$   
 Precision =  $TP / TP + FP$   
 Recall =  $TP / TP + FN$

$F_1 = 2 * (pre * re) / (pre + re)$

variance =  $\sum ((x - mean)^2 \text{ for } x \text{ in list}) / (len(list) - 1)$   
 standard dev =  $np.sqrt(\sum ((x1 - x1.mean())^2) / len(x1))$   
 Boolean mask =  $x < 5$  voor array of  $x[x < 5]$  voor list values  
 df.fillna(), df.dropna(), df.replace("a", "b")  
 df.unique(), df.duplicated("type"), df.drop\_duplicates()

% matplotlib inline  
 import matplotlib.pyplot as plt  
 plt.show()

specifieke kolom = df["kolom"] of df[["kolom", "kolom2"]]

specifieke rij = df.iloc[i]

kolom van rij = pd.iov.loc["Dbb"]["Zetels"]

nieuwe waarden tellen = df["kolom"].value\_counts()

lowercase = df["PartyName"].str.lower().str.contains("party in lijst lower")

pd.read\_csv('file')  
 df.to\_csv('file.csv')  
 df.shape, df.count()

np.zeros, np.ones, np.arange, np.linspace, np.random.random, np.loadtxt, np.info

ld digit  
 LD non digit  
 lw word  
 LW non word  
 \s space  
 \S non space

90 100 60 5  
 5 30

45 | 5  
 5 | 45

90 | 10  
 10 | 90

		1000		
		10	9	1
	100	990	99	891
	—			
Ziel	10			
niet Ziel	90			
				9 81

810  
 891



np.ones((3,5), dtype=float) → array([1., 1., 1., 1., 1.])

np.arange(0,20,2) → array([0, 2, 4, 6, ..., 18])

X[::2] → every other element

X[:2, :3] → first two rows & three columns

np.arange(1,10).reshape((3,3))

np.any(x>5) True

np.all(x>5) False

X = rand.randint(100, size=10)

np.sort

5 40%

np.argsort

8 60%

df.values  
df.index

? 100%

5·4 + 8·6 = 68/10

data.items()

data[(data > 3) & (data < 5)]

200

data.loc  
data.iloc

100 100

TP	TN
FP	FN

np.concatenate → appends

90	90
10	10

0.9 99.0 99.0 99.0



import numpy as np  
 np.zeros/ones/full(n, fill, value)  
 → (x,y,z) waarden in array  
 np.arange(x,y) → als y leeg is, is x de range  
 np.empty(n) → random array met n verschillende waarden.  
 x[:,::n] every other column  
 .reshape(k,y) → reshape into x rows, y columns  
 np.vsplit(x,n) → Split array x in n arrays  
 h.split → spreid over een  
 np.hstack(x) → make x an horizontal array  
 np.vstack(x) → make x a vertical array.  
 x>5 → array vol met booleans  
 x[x>5] → waarden terug die hoger zijn  
 np.sort(x) → sorteert de array

X.std() = Standard deviatie  
 machine = df[df>5] & [df<10]  
 fancy indexing = df[['A','C']]  
 np.sqrt(((x1-x1.mean())\*\*2).mean()) → standaard deviatie

pd.Series values → on array  
 pd.DataFrame(x) → maak van de array x een dataframe  
 X.loc[n] → haal de n rij eruit (print)  
 X.iloc[n] → haal de rij index n eruit (print)  
 X.dropna or X.fillna(0) → vul leeg met 0 of drop dan  
 X.rename(['x','y']) → geef x de nieuwe name  
 X.concat([df1, df2]) → voeg dfs samen  
 X.join() → merge dfs of index  
 ↳ X.set\_index('key').join(other set\_index 'key')  
 X.mean() → per column (n) = columns ↳ per row  
 X.groupby(columnlist) → groepeer + de  
 X.pivot\_table(value, index=col, columns=col)  
 X.date\_range(start='x', end='y')  
 X.date\_range(start='x', periods='y')

X = pd.read\_csv('r.csv')  
 X.species.value\_counts().plot(kind='x') → plots dit naar op x axis  
 X.drop('y', axis=1) → drop de rij y  
 (X-X.mean()/X.std()) → z-score, 2 darts  
 X.sort\_values(['y']), ascending=False → sorteer op y  
 X.groupby('y')[z].max() → Max y for each z  
 X.Counter(text.split())  
 CS = pd.Series(X).sort\_values(ascending=False) →  
 tokenize text and create dict with counts  
 X[X\*X.index.str.len()==24].index → vind tokens waarin 24 karakteres staan  
 CS[CS==1].sum()/CS.count() → percentage unieke woorden die maar 1 versie hebben  
 CS[CS>1].sum()/CS.sum() → percentage woorden op alle woorden  
 pd.Series(Counter(CS.values)).plot(kind='bar') →  
 maak een bar plot met woorden die worden verspreid  
 X["y"] = (len(str(x)) for x in X['z']) →  
 maak nieuwe rij in df met de len



np.zeros, np.ones, np.full, np.arange, np.random, np.random.randn, np.random.normal, np.random.randint, np.random.rand, np.random.choice, np.random.choice, np.random.choice, np.random.choice

np.zeros  
 np.ones  
 np.full (3,5), dtype=float  
 np.arange (start and, stop)  
 np.random.random (3,3)  
 np.random.randn (3,3)  
 np.random.normal (mean, sd, (shape))  
 np.random.randint (0,10, (3,3))  
 np.random.rand (3,3)  
 np.random.choice (1,10), reshaped (3,3)

np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)

np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)

np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)

np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)

np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)

np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)

np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)

np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)

np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)  
 np.random.randn (3,3)

np.random.randn (3,3)

np.random.randn (3,3)

np.random.randn (3,3)

np.random.randn (3,3)

np.random.randn (3,3)



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

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

$$\text{accuracy} = \frac{TP + TN}{\text{Total}}$$

$$F_1 = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

predicted	TP	FP <sup>Type 1 error</sup>
	FN <sub>Type 2 error</sub>	TN
	actual	