

np.reshape(shape, dtype)
 np.ones()
 np.full(shape, fill_value)
 np.arange(0, 20, 2)
 np.linspace(0, 1, 5)
 np.random.random((3, 3))
 Between 0 and 1
 np.random.randint(low, high, (3, 3))
 np.random.randint(10, size=(3, 4))
 .ndim, shape, size

a	0	1	2	3	2e
0	1	2	3	4	
1	5	6	7	8	
2	9	10	11	12	
3	13	14	15	16	

a[:, 3] = 9, 12, 16
 a[:, 1] = 2, 6, 10
 a[0, 1:2] = 2, 3
 a[:, 3, :2] = [1, 5, 9], [3, 7, 11]
 Re = a[:, -1, :-1]
 a[0] = 1, 2, 3, 4
 a[:, 0] = 1, 5, 9, 13
 a[:, :2] = [1, 2, 3, 4], [9, 10, 11, 12]

a[0, 0] = ...
 np.arange(1, 10).reshape((3, 3))
 x[:, np.newaxis] = receive
 np.concatenate([x, y])
 - axis = 1, 0, none
 np.vstack / np.hstack
 np.split(x, [1])
 np.hsplit / np.vsplit
 np.abs(absute)
 np.negative()
 np.add.reduce(x)
 np.add.accumulate
 - timestamp
 np.ctype
 np.min(x, axis=...) of x.min()
 x.mean() - (1) - (1)
 np.count_nonzero(x > b)
 - count True
 np.sum(x > b, axis=1)
 np.any(x > 8) np.all(x > 8)
 np.sum((x <= 6) | (x >= 7))
 x[x < 5]

np.add.at(x, i(locations), 1)
 np.sort(x, axis=0)
 np.argsort
 np.partition(x, 3, kind='split')
 x[x%2 == 1] = -1 of
 x.where(x%2 == 1, -1, x)
 np.intersect1d
 np.setdiff1d
 np.where(a == b)
 x = x[~np.isnan(x)]
 A[np.isnan(A)] = 0
 x.reshape((x.size, 1))

%paste - copy/paste site,
 handle multi-line input
 %cpaste - more code
 %timeit - test time,
 multiple runs
 %a(ls) magic
 pwd - print working dir.
 ls - list work directory
 cd - change dir.
 mv - move file
 use: ! ;
 %prun - in content
 hints

R %timeit - time execution single stat
 R %timeit - time repeated execution of stat
 R %prun - run code with profiler
 %bprun - run line-by-line
 %memit - measure memory use of single stat.
 %mprun - run code with l-b-l mem profiler

Precision = TP / (TP + FP)
 Recall = TP / (TP + FN)
 Accuracy = (TP + TN) / (TP + TN + FP + FN)

$$2 \frac{P \cdot 0.9}{(P \cdot 0.9) + ((1 - P) \cdot 0.1)}$$

$$\sqrt{2 \frac{(1 - P) \cdot 0.9}{(P \cdot 0.1) + ((1 - P) \cdot 0.9)}}$$

np.nanmax(), np.nanmin()

pd.DataFrame(np.random.rand(10, 20, (2, 2)),
 columns = list('AB'))
 bill = A.stack().mean()
 A.add(B, fill_value = bill)

pd.concat([df1, df2]) - concat index
 , axis = 1, ignore_index = True
 , join = 'inner' -> column in both combining
 , join_axes = [dfs.columns]

df1.append(df2) columns / index hetrogene

pd.merge(df1, df2) -> column self name as row

pd.merge(, on = '...', left_on =, right_on =
 (self de woordes onder de kolomnaam)

bij index -> df1.a.join(df2.a)

pd.merge(df1, df2, how = 'outer') 'left' 'right'

pd.plot(kind = 'bar') 'hbar'

df.set_index(['']) title = ...

- dh -> how large in bytes
 for i in range(i.shape[0]):

L[:, 0] [np.isnan(L[:, 0])] = np.nanmedian(L[:, 0])

pd.read_csv("...", index_col = "names")

df.sort_values("names") OR df.sort_index(inplace = True)

df.isnull() -> true/false -> serie data [data.isnull()]

df = df[df.Party.isna()]

df.dropna() / df.fillna() -> axis =, how / thresh

df.fillna(method = 'ffill', axis = ...)

df["new"] = df["..."] + df["..."]

df["..."] > 9, sum()

df = df.sort_values(["...", "..."], ascending = [False, True])

df = df.loc[df.Party == "prodd", "ministerie"]

self.de -> df[df.Party == "prodd", "ministerie"]

df = df.drop([0], axis = 1)

for c in [...]: / df.columns

df[c] = df[c].str... / df.columns

df["..."].value_counts()[0].index.tolist()

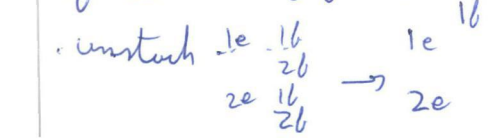
df[df["..."].isin(['...'])]

df.groupby(["...", "..."]).sum()

loc -> value, iloc -> position

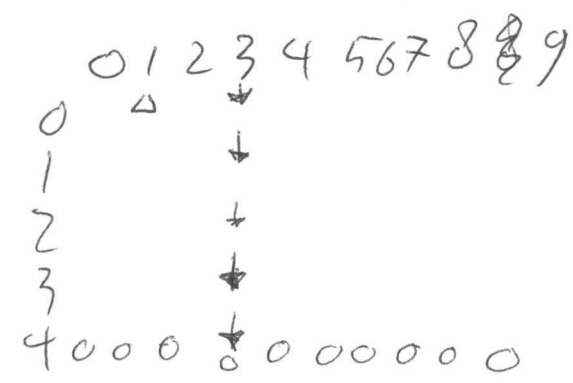
Serie onder index -> A.add(B, fill_value = 0)

df - df.iloc[0] -> 1e row = 0
 df.reset_index()



NP ARRAY (List of tuples)

- ZEROS (3,4) = val 0
- ones " " 1
- arange (start, stop, step) = open/closed
- linspace (1)
- Full (Cshape, Int) volle array
- Random, Random (shape)
- shape
- ndim = dimensions
- info() = help
- axis 0 = Row's
" 1 = Column's



- s[0,1] = Δ
- s[1,3] = ↓
- s[1,:] = 0

- a. sort()
- Reshape (shape)
- concatenate ((a,b) axis=0)
- vstack
- hstack
- Random, randint (start, stop (shape))

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

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

F1 = $\frac{2TP}{2TP+FP+FN}$

ACC = $\frac{TP+TN}{P+N}$

Pandas pd.

DataFrame

- melt (df) columns → Rows
- concat = append (axis 0,1)
- Pivot (columns='var', values='val')
- SORT-values (' ') descending = false
- Rename
- SORT_index
- drop

• Pivot_table aggfunc = functie als 'count', 'sum', ...

- CROSS-tab (col, col)
- index[0]

FILTER (Regex) →

- loc [i, x2, x4] columns van x2 tot x4
- iloc [3, [1,2,3]] 1, 2, 3 de columns
- values_count()
- nunique()
- describe
- Groupby (bv = 'col')

- alles met
- length\$ eindigt met 'length'
- ^start\$ begint met 'start'
- x[1-5]\$ " x, eindigt met 1,2,3,4,5
- ^?!species\$ alles behalve met species

Python magic:
%ls: directory read
%magic: alle magic commands

list comprehension:
new_list = [i+1 for i in old_list if i < 5]
new_list = [i for i in range(10)]

Numpy:

$x[pos]$: get position in 1D array.
 $x[pos, pos]$: get position in 2D array.
~~x.reshape~~ $x.reshape((row, column))$: zo reshape je een numpy array.
 $np.array([1, 2, 3])$: maak een list en daarvoor een numpy array.
 $x - 1$: voer berekening uit op elk item in x .
 $x.reshape((getal, 1))$ maak van een 2d array een 1d array.
 $x + 1 == 2$: apply boolean mask on array.

Pandas:

$pd.read_csv('csv')$: een csv bestand inladen.
 $df.column.value_counts()$: laat de values van de verschillende elements zien.
 $df.drop('column/row', axis=0/1)$: verwijder een column of row van df.
 ~~$df.sort_values('column')$~~
 $df.sort_values(['column1', 'column2'], ascending=[True, False])$
 $df.groupby('column')['Another column'].mean()$
 $indexno = df[(df['column1'] <= 10) \& (df['column2'] >= 5)]$
 $df.drop(indexno, inplace=True)$: drop meerdere dingen tegelijk
 ~~x~~ $Counter(z.split())$: Tokenize text
 $xy = pd.Series(x).sort_values()$: maak een pd series en sorteer de values.
 $xy[xy==1]$: geef alle unieke woorden in een series.
 $(xy==1).mean()$: geef gemiddelde van alle woorden met value 1.
 $df.describe()$: laat de values van df zien.
 $df.index.str.contains('x').head(10)$
 $(df - df.mean()) / df.std()$: z-normalize

Formules:

Precision: $TP / (TP + FP)$

Recall: $TP / (TP + FN)$

accuracy: $(TP + TN) / (P + N)$

Numpy $P = \text{array}$

Diagonal of array:

$d = \text{np. arrange}(P.\text{reshape}[-1])$

$D[d, d]$

String array:

$P[\text{row}, \text{column}]$

: alles, differ = grens (+)

Array Flatten

$D.\text{flatten}()$

~~$\text{np.sort}(\text{axis}=0)$~~

$\text{np.sort}(P, \text{axis} =)$

0 = kolom, 1 = row

np.random.randint

(1, 10, (5, 10))

(problemen)

$\text{np.T} = \text{transpose}$

From numpy import doc

df.join

df.groupby(by='col')

df.merge

~~df~~ ? np.doc

np.lookfor

Axis = 0: vertical/row/* Axis = 1: horizontal/column/*

np.array (rij/kolom), shape (0,1)

Logic operators

(A=xor)

! = or

& = and

~ = not

Numpy

make

np.zeros
np.arange
np.linspace
np.full
np.random.random

inspecten

a.shape
a.ndim
a.dtype
a.isnan
a.nanmedian

math

np.median
np.mean
np.mode
np.sort

Regex

re.match
re.sub
re.findall

\: escape/start

.: any except newline

^: start

?: end

[]: set

|: or

(): group

insert:

- range

^ negate

A: start

Z: end

\b: empty

\d\D: (non) digit

\w\W: (non) alphanum

\s: whitespace

manipuleren

a.flatten
a.transpose
a.reshape
a.resize
a.column_stack
a.unstack
a.hsplit
a.vsplit

overig

a.sort
a.argmax
a.diagonal
np.nonzero(L)

Boolean indexing

bool_arr = ~~arr~~ [cond]

newarr = ~~bool_arr~~ [boolarr]

sort() axis = 1 per rij

axis = 0 per kolom

broadcasting a = [c1], [c2], [c3]

b = [1, 2, 3]

arb

Formules

Z-norm = $\frac{x - \text{mean}}{\text{std}}$

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

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

$$A = \frac{TP + TN}{\text{All}}$$

$$F_1 = 2 \cdot \frac{P \cdot R}{P + R}$$

Actual

P N

P TP FP
N FN TN

Pandas

select

df.loc -> kolom
df.iloc -> index
df.iatmin/iatmax
df.ix
df.where
df.item
df.columns (name)

manipuleren

df.drop
df.dropna
df.fillna
df.apply
df.pivot_table
df.corr()
df.reset_index
df.unstack
df.set_index

plotten

df.plot()
df.bar
df.barh
df.hist
df.scatter
df.T (transpose)

sort

df.sort_values([val1, val2], ~~asc~~ = [True, False])
df.sort_index

overig

df.unique
df.nunique
df.value_counts
df.groupby
df.random.sample
-> i

python

?
* = wildcard
%time
%prun
%ls -lh
%os
%pwd

Numpy

```

np.arange(start, stop, step)
np.linspace(start, stop, steps)
np.zeros((3,3)) / np.ones((shape))
np.full((shape), value)
np.random.random((shape))
a.shape / ndim / size / dtype
a.astype(type) - convert to type
    
```

```

np.transpose(a, (1,0,2)) - kies dimensions
a.ravel() / a.flatten() - make 1D
a.reshape(shape) / a.resize(shape)
np.append(a, value) / np.delete(a, [1])
np.insert(a, value, place)
np.hstack(a, b) / np.vstack(a, b)
np.repeat(a, 3) -> [1,1,1,2,2,2,3,3,3]
np.tile(a, 3) -> [1,2,3,1,2,3,1,2,3]
np.vsplit(a, 3) / np.hsplit(a, 3)
    
```

```

a[row, column]
b[0:2, 1] - rows 0&1 of column 1
a[a<2] - boolean indexing
fancy indexing: b[[1,0,1,0], [0,1,2,0]]
paket = (1,0)(0,1)(1,2)(0,0)
    
```

```

def middle-square(L):
    n = len(L)
    begin = int(n/2 - 2)
    eind = int(n/2 + 2)
    return L[begin:eind, begin:eind]
    
```

```

• Elements in 2 arrays match:
  np.where(a == b)
• Get all numbers between 5 and 10
  np.where((a > 5) & (a <= 10))
  a[(a > 5) & (a <= 10)]
    
```

```

• np.sqrt(((x1 - x1.mean())**2).mean()) == np.std()
• L + np.arange(L.shape[1]) - to +1 +2 etc.
• vcr.Party = kwr.Party.stk.lower().stk.replace(words, "")
  words = re.compile("word|dog")
• greater than 9 = (df[['col1', 'col2', 'col3', ...]].mean().sum() / sum)
  ser1[~ser1.isin(ser2)] - np.percentile(1, 9)
    
```

Pandas

```

s.value_counts(dropna=) values + count of Series
df[col] - return column of df as Series
df.columns['a', 'b', 'c'] - rename columns
df[['col1', 'col2']] - return columns as new df
pd.isnull() / pd.notnull() - boolean array
df.dropna() / df.fillna()
df.set_index('column-name')
df[(df[col] > 0.5) & (df[col] < 0.8)]
df.sort_values(col-name, ascending=)
df.groupby(col) / df.groupby(['col1', 'col2'])
df.groupby(col1)[col2].mean()
df1.append(df2) - add rows of df1 to end of df2
pd.concat([df1, df2], axis=1) add col of 1 to 2
df.pivot_table(index=col1, values=[col2, col3])
pivot = df.pivot_table(values="count", index="Name", columns="Sex", aggfunc="sum")
iris.species.value_counts().plot(kind='bar')
X = iris.drop('species', axis=1)
iris.sort_values(['sepal-length', 'sepal-width'], ascending=[True, False])
iris.groupby('species')['sepal-length'].max()
iris[(iris['sepal-length'] > 5) & (iris['sepal-width'] < sqrt(5))]
    
```

```

txt = text.split()
c = Counter(txt)
CS = pd.Series(c).sort_values(ascending=False)
CS[CS*(CS.index.str.len() == 24)]
CS[CS==1].sum() / CS.count() unique
(CS==1).sum() / CS.sum() of all words

names = pd.read_csv('babynames.csv')
last-char = names.copy()
last-char['Name'] = last-char.Name.str[-1:]
mannen = last-char[last-char.Sex == 'M']
vrouwen = " " " " " " " = 'F'
mannen-sorted = mannen[['Name', 'Count']]
  .groupby('Name').sum()
vrouwen-sorted idem
mannen-sorted.columns = ['mannen']
vrouwen-sorted.columns = ['vrouwen']
combined = mannen-sorted
combined.join(vrouwen-sorted)

Precision: TP / (TP + FP) - (0.9 * 0.01) / (0.9 * 0.01 + 0.1 * 0.95)
%ols %time %timeit
%matplotlib %who-Is
%sm
    
```

from numpy import doc
? np.doc

? df.x

np.locfor('')

Slicing

np.random.randint(start, stop, (rows, columns))
x[row, column]

Array manipulation:

x.flatten() / x.ravel()
x.reshape(row, column)

np.append(x, i)
np.insert(x, i, s)
np.delete(x, [i])
np.concatenate(x, i)
np.vstack((a, b))
np.hstack((a, b))
np.vsplit
np.hsplit

x.shape
x.astype(int)

np.zeros((3, 4))
np.ones('')

Pandas

df.sort_values()
df.rename()
df.sort_index()
df.reset_index()
df.drop()

df.plot(kind = bar/hist)

df['w'].value_counts
df[df.length > 7]

df[['column1', 'column2', 'column3']]

pd.series() ← make pandas object of something

TP = kans • aantal

FP = kansfout • aantal

Precision = TP / (TP + FP)

Recall = TP / (TP + FN)

		predicted	
		negative	positive
actual	negative	TN	FP
	positive	FN	TP

###

import

```

import b22
import codecs
import os
import random
from collections import Counter, defaultdict
import matplotlib.pyplot as plt
import numpy as np
import re
from bs4 import BeautifulSoup
from nose.tools import assert_almost_equal, assert_equal

```

```

from time import time
import tqdm_notebook
%matplotlib inline
pandas as pd

```

```

* irrelevant chars
verwijderen
for i in range(len(page-list)):
    for j in range(len(page-list[i])):
        clear_regex = re.compile('<.*?>')

```

np.nammedian
 np.a.flatten()

1

```

np.random.randint(n, size=n*nm), random.choice([j for j in range(n) if j != guess[i] and previous[i]] for i in range(len(guess)))
plt.plot, plt.xlabel('hoi'), plt.ylabel(), plt.plot(x, list, label='hoi'), b22.open(loadfile(), encoding='utf-8', mode='rt')
as source_file: source_file.readlines(), if '</page>' in line,

```

2

```

linktopdf: if 'file.txt' in os.listdir() return 'file', elif os.path.exists('..1..data/Week 2/'):
return file

```

```

np.arange(0, n*o, o), np.array([1, 2, 3]), np.array([1, 2, 3], [4, 5, 6]), a.shape = (3, 2) of (2, 3)
a.size = 3, a[0,:] := all -> 1, 2, 3, a[1,3] -> 6, a[1,-2] -> 5, [2,2] [0,1:6:-2] -> 2, 4, 6
veranderen = a[0,2] = [1, 2], np.min(a), np.max(a), np.max(a, axis=1),
np.zeros(5) -> [0, 0, 0, 0, 0] np.zeros([2, 3]) [[0, 0, 0], [0, 0, 0]], np.ones, np.full((2, 2), 99)
np.random.rand(4, 2) -> [[0.81, 0.43], [0.69, 0.61], ..., ...]
array a + b, np.sum(a), a.reshape((3, 1)) -> [[1], [2], [3]]
np.sum(a, axis=0) np.arange(p.shape[0])

```

```

np.vstack(a, b), np.hstack(a, b), np.any(a > 50, axis=0), np.all(a > 50, axis=0)
a > 50, a[a > 50] -> [[60, 92, 80]], ((a > 50) & (a < 100)), (~((a > 50) & (a < 100)))
n = int(L.size ** 0.5)
np.reshape(L, (n, n)), L[L % n == 0], np.mean(L, axis=0), L.shape[1], a.T

```

3

```

pd.read_csv(linktofile(), index_col=0, sep='t', names=['jaar'], skip
pd.crosstab(kvr['jaar'], kvr['partij'], plot=a.filter
pd.dropna(subset=['Partij'], inplace=True print-table(0 "count", index="name
str.lower(), str.replace("W", "") ", columns="Sex", aggfunc=sum,
str.replace("|", join(waardenlijst), "") margin=True
kvr[~kvr['Partij']].str.contains("vragen")

```

```

• unique (df - df.mean()) / df.std()
kvr[kvr > 4] pred, yes no
value_counts().index[0], values actual yes TP FN
plot(kind="barh") no FP TN
str.count('i?') pred prec = TP : (TP + FP), Rec = TP : (TP + FN)
str.len() acc = (TP + TN) : (TP + TN + FP + FN)
• mean(), • median(), • mode()
corr = cor(df.columns, method='pearson') • sort_index(), sort_values(by='all'
ascending=True

```

Spielbrief Jop Richtigbaron 11/05/16

np.array

np.arange

np.count_nonzero

np.concatenate((arr1, arr2))

np.ravel np.reshape(L, (x, y))

np.sort

np.ravel(arr)

L[:, np.newaxis]

L [row:column]

arr.transpose

arr.reshape

arr.flatten

arr.shape [0]

$L[L \neq 0]$

arr.size

np.add

np.subtract

np.negative

np.multiply

np.divide

np.floor_divide

np.power

np.abs

np.sqrt

np.sum

prod

mean

min

max

std

var

argmin

argmax

median

percentile

any

all

abs

nonzero

nonzero

end:

np.cumsum(arr)

Pandas

data[1:3] implicit

data.loc[1:3] explicit

data.iloc[1:3] explicit index

Df.add
Df.subtract
Df.multiply
Df.divide
Df.floor_div
Df.mod
Df.pow

Df.dropna, fillna()
Df.drop
Df.concat
Df.merge
Df.append

pd.crosstab
pd
df.mean
df.mode
df.median
df.std

stack
unstack
kobm = C.loc[ry] kkmemo
uraagio = {ry, kobm}

ry = df.max().idxmax()

df.str (lower / upper / contains / len / strip / replace)

df.sort_values([columns]).ascending([true, false])

df.value_counts().index.tolist()[0]

groupby.apply(unstack)

Perc = (C df ['assignment', 'mean'] < 5, 5) / (df < 5, 5).mean * 100

for i in range(L.shape[0])

L[:, i] [np.isnan(L[:, i])] = np.nanmedian(L[:, i])

Regex

<.*>

[^ | \w | \s]

\w letter

\s whitespace

\n newline

\\ tab

\$ end of string

^ start

\d decimal

P accuracy	(1-P) * (1-acc)
TP	FN
P * (1-acc) FP	TN (1-P) * acc

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

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

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

$$F_1 = \frac{2 \cdot \text{recall} \cdot \text{Precision}}{\text{recall} + \text{Precision}}$$

precision = TP / (TP + FP) recall = TP / (TP + FN) Acc = nT

np.int0 (np.array - dtype) → b.dtype (int)

np.exp(b), np.add(a, b), np.divide(a, b)

b.cumsum(axis=1), a.concatenate() np.hsplit(a, 3)

np.copy(a), = a.copy(), = a.view()

a.sort(), a.sort(axis=0), np.concatenate((a,d), axis=0)

np.flatten() == b.ravel(), .resize(), .append, insert, delete

np.zeros, ones, arange(10, 25, 5), .full(), .random.random(), .empty

diagonal(s) = d = np.arange(s.shape[0]), return s[d, d]

np.isnan(L[:, i]), np.nanmedian(L[:, i]), np.nonzero()

np.arange(1, 4)[:], np.newaxis, np.sort(axis=0)

np.all(), np.any(), np.absolute == np.abs, IPython = help(len)

pd.melt(df), pd.concat([df1, df2], axis=1) data[-1] = laatste

df.sort_index(), df.reset_index(), df.drop(columns=['a', 'b'])

regex: (\), 'Length\$' = eindigt op length, '^sepal' = begint met sepal

df.drop_duplicates(), df.sample(), df.nlargest(n, 'value')

df.loc[:, a:b], df.iloc[:, [1, 2, 5]], df['w'].value_counts()

df['w'].nunique(), df.describe(), df.var(), dropna / fillna()

pd.merge(df1, df2, how="left/right/inner/outer", on="x1")

df1[~df1.x1.isin(df2.x1)] pd.MultiIndex()

websites dt.plot_hist(), dt.plot_scatter(x=a, y=b)

(df[assignments > 9].sum().sum()) np.linspace(0, 1, 5)

c = df[as].corr(), x = (c[!1].max()).idxmax()

y = (c[!1].loc[x].idxmax()), pd.Series(data=values, index=labels)

pivot[["F", "M"], min(axis=1) ~~sum()~~ idxmax()]

groupby → apply(functie) → unstack() gebruik .values ipv tolist()

df[df.x1.astype(str).str.contains("string") → alleen int)

.loc = rows/cols met labels vd index, .iloc = rows/cols met index posities

%ls -lh, %rm <filez, %time (it), %prun

.str.replace("1", join(...))

ggd(g, k) = if k == 0: return k else: return ggd(k, g%k) ^{ligvlarb:} return a*b (ggd(a, b))

$$F1 = 2 \times \frac{\text{Prec.} \times \text{recall}}{\text{Prec.} + \text{recall}}$$

Tips: .TAB, ?, *?, %magic, %time(it), %ls (-lh), %rm, %prun, %cp, %cd, %mkdir

Precisie = $TP / (TP + FP)$, Recall = $TP / (TP + FN)$, Acc. = $(TP + TN) / \text{All}$

PD Pandas: Series = 1D labeled array, DF = 2D

pd.read_'type' (file, ...) | dimensies: .shape | dubbels: $\text{len}(\text{set}(\text{df.index})) == \text{len}(\text{df.index})$ | orden op index: namen en hoe vaak dubbel: $\text{dubbel} = \text{df.index.value_counts}() \rightarrow \text{dubbel}[\text{dubbel} > 1]$ | sort_index(inplace=True)

rm spaties en "" + lowercase: $\text{df.index.str.lower().str.strip().str.replace(" ", "").str.replace("'", "")$

beperk tot substring: $\text{df}[\sim \text{df.index.str.lower().str.contains('sub')}]$ | df.describe() | sort_values()

new kolom $\sqrt{\text{var}}$ van verschil: $\text{df}['n'] = \text{np.sqrt}((\text{df.vers1} - \text{df.vers2})^{**2})$ | trek gem. van waarden in

kolom van elke waarde af / normalisatie: $M = \text{df.mean}()$ | $(\text{df}[\text{M.index}] - M)$ | plot iets op df:

$\text{df.kolomnaam.sort_values().plot}()$ | kolom > iets in gemeente + hor. bar chart: $\text{df}[\text{df.kolom} > 0]$

.gemeente.value_counts().sort_values().plot(kind='barh') | Vervang NaN: $\text{df} = \text{df.fillna}(0)$

rm rows/cols met NaN: $\text{new} = \text{df.dropna}(axis=0/1, how='any')$ | Perc. $\frac{\text{df.mean}()}{\text{df.std}()} * 100$ | Z-normalisatie:

$z\text{-df} = (\text{df} - \text{df.mean}()) / \text{df.std}()$ | Corr. tussen opdr.: $C = \text{df}[[\text{" "}]].\text{corr}() \rightarrow C = C[!:] - 1$

$r_{ij} = C.\text{max}().\text{idxmax}() \rightarrow \text{kolom} = C.\text{loc}[r_{ij}].\text{idxmax}() \rightarrow \text{vr} = \{r_{ij}, \text{kolom}\}$ | KGV: return

NP $\text{int}(\text{abs}(a*b) / \text{ggd}(a,b))$ | (Numpy): normaliseer tot 1: $(P / (P.\text{sum}(axis=0)))$ | np.diagonal(s)

Fill NaN by median: $\text{for } x \text{ in range}(L.\text{shape}[1]): \rightarrow L[:,x][\text{np.isnan}(L[:,x])] = \text{np.nan}$

median(L[:,x]) \rightarrow return L | P. flatten() | return values: $\text{np.sum}(P < n)$ | vmd niet mul:

$L[\text{np.nonzero}(L)]$ | RMSE: $\text{np.sqrt}(\text{np.mean}((L[:,1] - L[:,0])^{**2}))$ | np.arange(1, 4) x 2

$[:, \text{np.newaxis}]$ | Sort array L naar R: $\text{np.sort}(L, axis=1)$ | test waarden: $\text{np.all}(\text{np.abs}(L) > 10)$

| select last 3: $x[-3:]$ | select last 3 cols of first 2 rows: $x[:, 2, -3:]$ | reshape

into 1D: $x.\text{reshape}(x.\text{size},)$ | change elem. by std^{**2} from mean: $(x - x.\text{mean}())^{**2}$ | bereken

std: $\text{np.sqrt}(\text{np.mean}(x - x.\text{mean}())^{**2}) == x.\text{std}()$ | bool mask: $x[x \% 3 == 0]$ | fancy index

diagonal: $d = x.\text{shape}[0] \rightarrow \text{dd} = \text{np.arange}(d) \rightarrow x[\text{dd}, \text{dd}]$ | sort sepal lengths for species

iris. $\text{groupby}('species')['sepal-length'].max()$ | Tokenize text, count tokens, create Series, sort: $C =$

$\text{Counter}(\text{text.split}()) \rightarrow (S = \text{pd.Series}(C).\text{sort_values}(\text{ascending=False}))$ | find tokens:

$C[S * C.\text{index.str.len}() == 24].\text{index}()$ | perc. unique words 1x en perc. of all: $(C == 1).$

.mean(), $(C == 1).\text{sum}() / C.\text{sum}()$ | land met meste: $\text{bd} = \text{pd.read}(\text{file})[1] \rightarrow$

$\text{bd.groupby}(0)[1].\text{count}().\text{sort_values}().\text{tail}(1)$ | np.setdiff1D() | kruistabel: pd.crosstab

df.pivot_table() | ^{2 \rightarrow deelstaten} | With open(text, encoding='utf-8') as f: $\rightarrow t = f.read() \rightarrow t = 2 =$

re.split(r'\W+', t) $\rightarrow \text{pdf} = \text{Counter}(\text{len}(w) \text{ for } w \text{ in } t[2:]) \rightarrow$ return pdf |

! theory slices + indices = views, no copies. Broadcasting: 2 arrays differ in dimensions, kleinere is padded met 1

als shape niet matcht, array met shape==1 is stretched. sizes niet zelfde en niet==2 \rightarrow error.

reversed: $x[:, ::-1, ::-1]$ | $x[:, 0]$ eerste column / $x[0, :]$ eerste row | np.concatenate \rightarrow

join arrays along axis | isnull(), notnull() | $\text{df}[[\text{opdrachten}]] > 9.\text{sum}().\text{sum}()$

$\downarrow (5.5 \mid 5.5)$

Numpy

Data types:

- Checking: arr.dtype
- Changing: arr.astype(...)
- int, float, str

Operations:

- arr1 +, -, /, * arr2
- np.exp, sqrt, log, sin, cos (arr)

Creating arrays:

- 1D: np.array([1, 4, 5])
- 2D: np.array([(1, 4, 5), (2, 3, 6)])
- np.arange(10, 25, 5) → [10, 15, 20]

dtype = ...

rows, axis 0

cols, axis 1

They 'collapse' along this axis

Inspecting:

- arr.shape → (2, 3) row, column
- len(arr) → 2
- arr.size → 6

Comparing:

- Element-wise:
 - arr1 == arr2 → [T, F, T, ...]
 - arr1 > arr2 → ...
- Array-wise: np.array_equal(arr1, arr2)

Slicing:

- select: arr[1, 2] row 1, column 2
- slicing: arr[0:2, 1] row 0/1, column 1
- arr[1] row 1
- arr[:, 1] column 1
- bool indexing: arr[arr > 3] → [...]

Functions

Or in the form: np.sum(arr, axis=...)

arr.sum(axis=...)

arr.min(...)

arr.max(...)

arr.mean(...)

arr.median(...)

arr.sort(...)

x = arr.view()

x = arr.copy()

Manipulation:

- arr.T → rows become cols, cols become rows
- arr.reshape(..., ...) → new arr, same data
- arr.resize(..., ...) → same arr, new shape

np.append(arr, x, axis=...)

np.insert(arr, index, x, axis=...)

np.delete(arr, indexes, axis=...)

np.concatenate((a, b), axis=a)

np.vstack((a, b))

np.hstack((a, b))

Creating dataframe:

df = pd.DataFrame([4, 7, 10], [5, 8, 11], [6, 9, 12], index=[1, 2, 3], columns=['a', 'b', 'c'])

Pandas

Functions: high to low

df.sort_values('price', ascending=False)

df.rename(columns={'y': 'year', ...})

df.drop(columns=['length', ...])

df.sort_index() / df.set_index('A')

df.reset_index() → index becomes 0, 1, 2, ... again

- df['A'].value_counts()

- df['A'].unique()

- df.dropna() / df.fillna(value)

- df.sum, count, median, ['A'] mean, min, max()

- len(df) → # of rows

New column / row:

- col: df['new-col'] = [...]
- row: df.append({'a': 4, 'b': 7, ...}, ignore_index=True)

Logic: df.Column.isin(values)

pd.isnull() / pd.notnull()

operators: &, |, ~ (not)

df.any(), df.all()

Subset (rows):

df[df.length > 7]

df.iloc[10:20] (by position)

df.nlargest(n, 'value') } selects & orders

df.nsmallest(n, 'value')

df.drop_duplicates()

df.head(n) / df.tail(n)

Subset (columns):

- single col: df['length'] or df.length

- multiple cols: df[['length', ...]]

df.iloc[:, [1, 2, 5]]

df.iloc[df['a'] > 10, ['a', 'c']]

df.loc[:, ['a', 'c']]

df.filter(['a', 'c', 'd'])

Regex: * = anything

'a': containing 'a'

'a\$': ending in 'a'

^a': starting with 'a'

Plot: df.plot.barh, line, etc.

Combining: on top of each other

pd.concat([df1, df2])

next to each other: axis=1

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

Filtering: df1[df1.A.isin(df2.A)]

Grouping: df.groupby(col)[other_col].count()

size()

or .max(), .min(), etc.

Pivot table: df.pivot_table(index='', columns='', values='', aggfunc=np.sum)

Reading Files:

pd.read_csv(parameters!)

String operations:

S.str.count(r'(^F*)').sum()

S.str.match(r.) → returns bool

S.replace('...', '', inplace=True)

S.str.findall(r.)

S.str.contains(r.) → returns bool

S can be df[df['A']...]

For example

Python

% + tab → all magic commands

Numpy cheatheet

- np.arange (start, stop, step) → make numbers between start and stop with step space between
- np.linspace (start, stop, number) → make x numbers between start and stop
- np.array([(1,2,3), (4,5,6)])
- a.shape & a.size + a.ndim + a.dtype + a.astype(int) + np.info
- np.concatenate(a) - np.**? → all functions

Reading files in python → CSV, XML, JSON, PDF, HTML

- general logic → 1. open file → 2. read using special purpose module → 3. convert to dict or list
- spreadsheet to list of lists or list of dicts using columns and indexes + headers
- with open ('./Data/...csv') as f:
 - all_lines = []
 - for line in f:
 - all_lines.append(line.split(';'))
- dict = json.loads(test)
- titles = tree.xpath('//item/label/text()')
- magic commands
 - ls ... /Data/ → mijn_files
 - fort in mijn_files:
 - lwc -w "\$f"
- CSV → pandas
- XML → lxml library and xpath
- JSON → json.load + json.dump
- PDF → pdfminer
- HTML → BeautifulSoup and requests
- Read from disk
 - %ls .. /Data/ → what files are in folder?
 - f = open ('.. /Data/ Monumentale - Bomen.csv')
 - f.readline(), split(';')

np practice

- np.full((3,3), True, dtype=bool)
- np.where(arr % 2 == 1, -1, arr) → copy
- np.r_[np.repeat(a,3), np.tile(a,3)]
- np.intersect1d(a,b)
- np.setdiff1d(a,b) → $-(L[(L<10)|(L>10)])$
- P/P.sum (axB=0) → normalize columns

Numpy

- return P/P.sum (axB=0)
 - ↳ normalize columns to 1
- np.diag
- def rmse(L):
 - test = (L[:,1] - L[:,0])**2
 - return np.sqrt(test.mean())
- np.arange(1,4)[:, np.newaxis]
- np.sort(L)

- (mse: test = (L[:,1] - L[:,0])**2)
 - return np.sqrt(test.mean())
- np.arange(1,4)[:, np.newaxis] b = np.arange(1,4) → broadcasting
- Q: From array return values < -10 and values > 10
- A: return L[(L < -10) | (L > 10)]

Theorie

- precision = TP / (TP + FP) → how many positives are correct? / reliability of test
- recall = TP / (TP + FN) → how much is detected?
- accuracy = TP + TN / (TP + FP + FN + TN) → reliability of test
- F1 = 2 * (recall * precision) / (recall + precision)

	True	False	
pos test	TP 90	FP 990	1080
neg test	FN 10	TN 8910	8920
	100	9900	10000

N = 10000
 P = 0,01
 Precision = 90 / (90 + 990) = 0,083
 accuracy bad for shewel test (close to 1)

Theorie

```
def qgd(q, h):
    if h == 0:
        return q
    else:
        return qgd(h, q%h)
def kv(a, b):
    return abs(a*b) / qgd(a,b)
```

Pandas

- pd.crosstab (col1, col2)
- cross.plot (kind='bar')
- cross.plot.barh()
- kvr.str.lower(), str.replace('[^\w]|allen', '')
- kvr = kvr[~kvr.str.contains('vragen', na=False)]
- vragen['pvad'].sort_values(ascending=False)
- kvr[['Antwoord']].applymap(lambda x: len(x))
- kvr[['vraag']].applymap(lambda x: str.count(x, '?'))
- vraag.hist (bins=50)
- vraag.corr().iloc['kar'][0]

Pandas 2

- ax = df.plot (kind='line')
- ax.set_xlabel('Year')
- ax.set_title('Title')
- ax.legend(['F', 'M'])
- str = 'ride... %os' % (var)

- ind = [3,7,4]
- x[ind] → fancy index
- x[xrowF:np.newaxis, col]

```
df = pd.DataFrame(
    {'a': [4,5,6],
     'b': [7,8,9],
     'c': [10,11,12]},
    index=[1,2,3])
```

- df.sort_values(['All', 'Name'], ascending=[False, True])
- df.query('ratio_col and ratio > -0,1')
- df.index.values
- error['rel:error'].idxmax()
- .sum(skipna=True)

- df.groupby(by='col')
- df.groupby(level='ind')
- pd.merge(a, b, how='outer', on='')

robert9561#

pd.read_csv (File, sep=';', index_col=0, header=None, skipinitialspace=True, compression='gzip')

pd.crosstab (Kvr. partij, Kvr. partij) [partijen]

str.replace / lower / contains

KVR = KVR[KVR.partij.str.contains('bragen')] prec = TP / (TP + FP) recall = TP / (TP + FN)

KVR['deelv'] = KVR.vraag.str.count('\?')

pd.pivot_table (df, index, values, columns, aggfunc, margins=True)

sort values (by = ['All', 'name'], ascending = [F, T])

pivot. table (pivot. ratio < 0.1). index. value_counts ()

L[:i] kolom L[k%n == 0] [3:7, 3:7] vierkant L[i:i+1] = i

np.arange (0, 7, 7x144) Start stop step

reshape (-1, i) zoekt zelt uit

pivot. values ['D.', 'E'] ... aggfunc = {'D': np.mean, 'E': np.min}

np.random.randint (low, high, size)

plt.plot (x-lijst, lijst values, label) plt.show ()

split ('str')

DF. groupby (kolom1, kolom2)

apply (lambda x: x['n'].count, sum)

unstack ()

• <TAB> = mogelijkheden van functie

np.random.randint(begin, eind, size=(kolom, lengte))

ndim = number dimensions

shape = lengte per regel (dimension)

size = Totale lengte (len)

slicing: x [start:stop:step]

[row, column]

Reshape = (rijen, lengte)

joining arrays = np.concatenate([x, y], axis)

↑ stack & Hstack

split arrays: np.split

ufunc:

- add, - subtract
- negative, - multiply
- divide, - power (2)
- mod

np.count_nonzero(x) → hoeveel van iets

np.where(x) → select alles waar voorwaarde

np.nan...() → geeft nan values

np.all/any(array)

np.random.randn(x, x)

↳ array w random values

np.set_printoptions() → hoe je weer geeft

p.d. DataFrame() → series → DF

df[~df.isin(x)] → where functie

p.d.concat([ser1, ser2], axis=1)
↳ Hstack

p.d.Series([list comprehension])

p.d.<TAB> - mogelijkheden

index_col → index by col

df[[kolom]] → waarde van alleen die kolommen

~~df~~.isnull - .notnull

Multi Index → SF 2000 3 mil
2010 5 mil

• unstack → data frame

p.d.concat(ignore_index)

↳ verint nieuwe index voor overlappend

p.d.concat(keys=[""])

↳ laat zien welke waarde van welke df is

p.d.merge() → merge df's

↳ Left/right_on = kolom waarop moet linken

: groupby.count() = hoeveelheid

p.d.pivot_table(index, kolom, waarden, aggfunc)

p.d.crosstab() → laat relatie zien tussen klassen

↓
functie

Accuracy = How many of all predictions were correct $\frac{TP+TN}{TP+TN+FP+FN}$

Precision = how many of the predictions were correct $\frac{TP}{TP+FP}$

Recall = how many of the instances were picked up $\frac{TP}{TP+FN}$

$$RMSE = \sqrt{\frac{\sum(\text{predicted} - \text{actual})^2}{n}}$$

predicted
P F

	P	F
P	TP	FN
F	FP	TN

#Kijk GOED naar de vraag

timeit 'rest' to time rest

Numpy

Create np. zeros
 ones
 Full
 random . randint
 normal
 . arange (start, stop, step)

A.size
 A.shape
 A.start
 A.argmax
 np.argmax

Index 2D [row, col] slice = 2D [start; stop: step, ""]
 A[:, 2] = every other val A[:, :3] first 3 cols
 A[::-1] = reverse First 3 row

Computation $A_{n1} + A_{n2} \rightarrow$ calculates element wise

+ np.add
 - Subtract
 / divide
 // Floor-divide
 ** negative
 ** Power
 * multiply
 % mod
 exp
 Power
 log / log2 / log10

operators (can be given along different axis / series etc.)
 np.sum
 • min $\left| \begin{array}{l} \cdot \text{argmin} \\ \cdot \text{argmax} \end{array} \right|$
 • max
 • mean \rightarrow gives indexes
 • median
 • std
 • var
 • Percentile
 • nan 'operator'

Min, Max, etc \uparrow

Reshape, concat, split A.flatten $\rightarrow (n, 1)$
 A.reshape(3, 3) = 1D len(9) \rightarrow 2D len(3, 3)
 A[:, newaxis] = ~~axis~~ makes array (10,)

np.concatenate([A, A]) = 2x as long
 " ([A, A], axis=1) = ~~axis~~ extra column

np.vstack - np.hstack

A1, A2, A3 = np.split(A, [3, 6]) np.vsplit
 \rightarrow [1, 2, 3], [4, 5, 6], [7, 8, 9] np.hsplit

Broadcasting works when arrays:
 1. are equal (dimensions)
 2. one of them == 1
 so (9, 1) and (9, 2) will not work

$SD^2 = (x - \text{mean})^2$
 $SD = \sqrt{SD^2 \cdot \text{mean}()}$
 Precision = $\frac{TP}{TP+FP}$
 Recall = $\frac{TP}{TP+FN}$
 Acc = $\frac{TP+TN}{All}$

ziekte kont voor 1 op 10
 Accuracy = 90%
 $P = (0.1 \times 0.9) / ((0.1 \times 0.9) + ((1 - 0.1) \times 0.1))$

Numpy

np.array() / np.arange()
np.ones() / np.zeros() / np.full() / np.linspace()
np.random.random() / choice() / randint() / seed() / shuffle()
ndim / shape / size
reshape() / x [np.newaxis, :] / np.flatten() / np.transpose()
np.concatenate() / np.vstack() / np.hstack() / np.repeat()
np.split() / np.hsplit() / np.vsplit()
np.add() / subtract() / multiply() / divide() / mod()
np.negative() / abs() / exp() / log() / loge() / power()
np.reduce() / np.accumulate() / np.cumsum() = (all vs steps)
np.outer() (multiply table)
np.sum() / prod() / mean() / mode() / std() / var() / min() / max()
np.argmax() / ~~np.argmax()~~ / median() / percentile() (pleur non error)
np.any() / all() -- b, 1, n (not) / isin()
np.sort(axis=...) / np.argsort() -- axis=0 = v / axis=1 = h
np.dtypes (E names: 'i', 'f', 'o', 'u', 'S', 'O', 'U', 'V', 'M', 'O', 'F', 'D', 'C', 'S', 'b', 'B', 'i', 'f', 'o', 'u', 'S', 'O', 'U', 'V', 'M', 'O', 'F', 'D', 'C', 'S', 'b', 'B')
np.diagonal() / np.setdiff1d()
x[y, x] : = help ry -1 = achterste x[::-1] = reversed
np.where()

pandas

pd.Series() / pd.DataFrame()
values / index / columns / size / shape / ndim
keys() / items()
mean() / mode() / median() / quantile() / corr() / logz()
pd.notnull() / isnull() / dropna() / fillna()
str.len() / lower() / upper() / startswith() / find() / replace() / etc
value_counts() / nlargest()
max() / min() / idxmax() / idxmin()
groupby() / pivot_table() / crosstab()
div() / apply()
plot() / histogram()
pd.read_csv()
sort_values()
loc() / iloc() / ix() = explicit / implicit / hybrid
pd.MultiIndex / reindex() / unstack() / stack()
reset_index() / set_index()
pd.concat() / append() / transpose()
& = intersection | = union ^ = symmetric difference
df[Index], df.index df[df>3]
pd.merge() / join()
tail() / head()

Re

findall() / sub() / replace()
W = letters / cyfers \d = cyfer
o = of mee + = 1 of mee
{ } = exact times
^ = match anything not in set
" ".join(list) < . + ? >
based on patterns in data predicting a value on new unseen instances of data / classification / regression
unlike statistics -> no explaining
selection = [min, min, etc]
bogo = random shuffle
breadfirst = knop, den knoden etc

help
shift+tab
?
.tab
% mode verbose
" context
+ word?
help(function)

accuracy = $\frac{TP+TN}{All}$
precision = $\frac{TP}{TP+FP}$
recall = $\frac{TP}{TP+FN}$
F1 = $2 \times (\text{Rec} \times \text{pres}) / (\text{rec} + \text{pres})$
TP = R.A (ratio/accuracy)
FN = R.(1-A) P
TN = (1-R).A = type 2
FP = (1-R).(1-A) = type 1

- 1 interacting with outside world
- 2 preparation
- 3 transformation
- 4 modeling & computation
- 5 presentation

Numpy dir (np)
 (start, stop, step)
 np.arange generate array
 np.concatenate join arrays
 vstack
 hstack
 h/v split split points as list of indices
 nan
 np.sum
 prod
 mean
 std
 var
 min/max
 argmin/max index of value
 median
 percentile rank-based stats
 any*
 all*
 .ndim nr of dims
 .shape size per dim
 .size total size
 .reshape in shape (x,y)
 .flatten 1D array
 .T transpose

Numpy slicing
 np.newaxis
 np.sort
 np.where
 x[::2] every other el.
 x[::-1] reverse array
 x[:, 2, :3] row, cols
 x[:, 0] first col.
 x[0, :] first row
 x[::-1, ::-1] reversed

Numpy Indexing
 x[y-as, x-as]
 row, column
 axis=0, axis=1

ufuncs
 np.add +
 np.subtract -
 np.negative -
 np.multiply *
 np.divide /
 np.floor-divide //
 np.power **
 np.mod %

Boolean & Masks
 & and
 ^ xor
 | or
 ~ not

x[x > 5] select values that meet requirement

Fancy Indexing
 pass array of indices
 (2, 1) (1, 0) x[[2, 0, 1], [1, 2, 0]]
 → [2, 1] x [1, 0]

Broadcasting
 Apply ufuncs to arrays of different sizes

Pandas dir (pd) of pd.Series
 pd.read_csv inlesen bestand
 .values
 .index
 pd.Series (data, index = index)
 [a:b] incl a, excl. b selection
 pd.DataFrame (input, columns, index = [1])
 serie ['index'] >>> value
 'a' in serie T/F boolean
 serie.keys() get indices
 list(serie.items()) returns list of tuples
 serie ['index'] = 3 change value of 'index'
 serie [(serie > x) & (serie < y)] masking
 serie [['a', 'b']] fancy indexing
 df.columnname > series of column, index + value
 df['columnname']
 df.values display df as 2D array
 df.T Transpose
 df.values[0] returns row
 df.values["index"] returns column
 !df[df[0]] get index for value
 for index, row in df.iterrows()
 fill_value = 0 replace nan

.isnull()/.notnull() boolean mask
 .dropna() drops r/c where nan is present (all, any) thresh
 .fillna()
 .sort_index()
 .sort_values() pass list of columns on which to sort
 .unstack
 .stack
 .set_index/reset_index
 .join
 .append
 pd.concat standard row-wise

testp testn TP FP
 FP FN TN
 Precision TP / (TP + FP)
 Recall TP / (TP + FN)
 nicht zieh zieh

Predicting:
 Classification spam/no spam
 Regression value eg. for course
 name where age < 10
 df[df['age'] < 10][name]

% ls magic = list met alles wat bestaat
 % magic
 % ls
 % cd change dir
 % rm remove
 % line f % lineit (test)
 value - counts

& and
 | or
 ~ not
 .loc Explicit index
 .iloc Implicit index
 .ix combination
 Expl. index: last index INcluded
 Implicit: last index EXcluded.
 direct masking is interpreted row-wise!

NB. Indexing refers to columns [... , ...]
 Slicing refers to rows [... ; ...]

mapers row-wise
 add() +
 sub() -
 mult() *
 div() /
 floordiv() //
 mod() %
 pow() **

Steps of DS
 1. Interaction with the world
 2. Preparation
 3. Transformation
 4. Modeling & Computation
 5. Presentation

Predicting:
 Classification spam/no spam
 Regression value eg. for course
 name where age < 10
 df[df['age'] < 10][name]

NB. Broadcasting!
 a b c
 0
 1
 2
 axis=0
 axis=1
 row, column

np.abs
 np.newaxis voor komma: row
 np.sort na komma column
 rb.array.shape >> 4,
 [:, np.newaxis] >> 4, 1
 [newaxis, :] >> 1, 4

Accuracy: $(TP+TN)/\text{total}$
 Precision: $TP/(TP+FP)$
 Recall: $TP/(TP+FN)$
 F1: $2 \cdot \frac{pr \cdot re}{pr+re}$

	actual	
	ziek	niet ziek
predicted	ziek	TP
	niet ziek	FN
	FP	TN

Ziekte met 95% accuracy, 1/100 personen
 N = 10000 P = 0,01
 TP $0,01 \cdot 0,95 \cdot 10000 = 95$
 FN $0,01 \cdot 0,05 \cdot 10000 = 5$
 TN $(1-0,01) \cdot 0,95 \cdot 10000 = 9405$
 FP $(1-0,01) \cdot 0,05 \cdot 10000 = 495$
 Ziekte met 90% accuracy, 1/10 personen
 P = 0,1

Precisionziek $(P \cdot 0,9) / ((P \cdot 0,9) + ((1-P) \cdot 0,1))$
 Precisionnietziek $((1-P) \cdot 0,9) / ((P \cdot 0,1) + ((1-P) \cdot 0,9))$

Z-score: $(df - df.mean()) / df.std()$
 RMSE: $np.sqrt(np.mean((L[:,1] - L[:,0]) ** 2))$

def GGD(g,k):
 if k == 0:
 return g
 else:
 return GGD(k, g%k)

middlesquare (4x4)
 k = L.shape[1]
 st = int((k/2)-2)
 ei = int((k/2)+2)
 L[st:ei, st:ei]

reshape square
 array 144 → 12 x 12
 L.reshape(int(math.sqrt(L.size)), int(math.sqrt(L.size)))

tafel
 eerste n geballen,
 deelbaar door k
 np.arange(0, n*k, k)

euclidean distance
 np.linalg.norm(a-b)

Arr = np.array([0,9,1,4,64])
 Arr[1:3] → 9, 1
 Arr[:2] → 0, 9
 Arr[-2:] → 4, 64
 Arr[-2] = 50 → 0, 9, 1, 50, 50
 Arr[arr < 3] → 0, 1
 Arr[arr < 3] = -1 → -1, 9, -1, 4, 64
 Arr2 = [[1,2,3], [4,5,6]]
 Arr2[0] → 1, 2, 3
 Arr2[0,1] → 2
 Arr2[:,0] → 1, 4
 Arr2[0,0] → 1

Arr2[1,0] → 4
 Arr2[1,-1] → 6
 Arr2[:2, :3] → hele arr2
 Arr2[:2, :2] → [[1,3], [4,6]]
 Arr2[:, :-1, ::-1] → [[6,5,4], [3,2,1]]
 Arr2[2,[2,0,1]] → 6,4,5

x = np.arange(10)
 x[:,2] → tafel 2, start = 0
 x[1::2] → tafel 2, start = 1
 x[:, :-1] → andersom g H m o
 x[5::-2] → 5, 3, 1

def kgv(a,b):
 return int(abs(a*b) / ggd(a,b))

[0,1,2 etc.] optellen bij L
 z = np.arange(0, np.size(L,1), 1)
 L+z
 vertical reshape
 L.reshape(L.size, 1)

se herhaling van 1 in array x
 np.where(x==1)[0][5-1]

wag numbers only
 geen = df.kolom.str.contains('^[0-9.-]')
 df = df.loc[~geen]

SLICING

create [[2,3,4], [3,4,5], [4,5,6]]
 a = np.arange(1,4)
 b = np.arange(1,4)[::-1]
 a+b

rij
 kolom
 np.sort(L, axis=1)
 np.sort(L, axis=0)
 np.all((L < -10) | (L > 10))
 ↳ True of False

np.random.randint(n, size=...)
 size = 6
 = (3,4) 6 getallen
 = (3,4,5) 3 hoog, 4 breed
 = (3,4,5) 3 keer 4 hoog, 5 breed
 np.random.random((3,3)) 3x3 kolom

np.subtract(a,b) a-b
 np.add(a,b) a+b
 np.divide(a,b) a/b
 np.multiply(a,b) a*b

np.mean(arr, axis=0)
 arr.reshape(x,y)
 x - lengte
 y - breedte

np.sum(x < 6, axis=1)
 ↳ in each row
 np.transpose(L)
 10,15 → 10, 33
 33,34 → 15, 34

reshape x2 → x1
 x2.reshape(x2.size)
 variance
 (x1 - x1.mean()) ** 2 = v
 std
 np.sqrt(v.mean())
 == x1.std()

divisible by 3
 x1[x1 % 3 == 0]
 diagonal
 d = x2.shape[0]
 dd = np.arange(d)
 dia = x2[dd, dd]

pd.read_csv('file', sep='\t', index_col=0, header=None, names=[''], compression='gzip', skipinitialspace=True)

df.kolom.value_counts() hoogste bovenaan
 df.sort_values(['kolom1', 'kolom2'], ascending=[False, True]) alfabetisch hoog naar laag

df['kolom'].corr(df['kolom2'])
 df['kolom'].str.count(r'\?')
 str.len(), str.lower(), str.replace(' ','')

df.dropna(subset=[''], inplace=True)
 df.drop(columns=[''])
 pd.concat → samenvoegen df's
 pd.merge

len(set(df.kolom)) = df.kolom.unique().size
 aantal unieke in kolom
 df.groupby('kolom')['lengte'].max()
 ↳ max lengte van kolom

restrict df[df['kolom'] > 5]
 c = counter(text.split())
 cs = pd.Series(c).sort_values(ascending=False)
 cs[cs.index.str.len() == 24].index
 ↳ tokens which occupy 24 char in total

unique words percentage
 cs[cs == 1].sum() / cs.count() of (cs == 1).mean()
 percentage all words
 cs[cs == 1].sum() / cs.sum()

df.pivot_table(index='Name', columns='sex', values='count', aggfunc='sum', margins=True).sort_values(by=['All', 'Name'], ascending=[False, True])

pd.crosstab(df['kolom'], df['index'])
 str.replace(r"[^A-Za-z0-9]+", " ")
 df['kolom'].str.len() of mean/mode etc.

df[df['vr'].str.contains('vr') == False] = df
 ↳ verwijder rijen waar vr in voor komt

PANDAS

df[df['vr'].str.contains('vr') == False].head()

TP
 $precision = \frac{TP}{TP+FP}$
 $recall = \frac{TP}{TP+FN}$
 $F1 = 2 * \frac{precision * recall}{precision + recall}$

$dir()$
 $help()$
 $?$
 $??$
 $.<TAB>$
 $*?$

$np.array([], dtype=)$
 $np.zeros()$
 $np.ones()$
 $np.arange()$
 $np.linspace()$
 $np.full()$
 $np.random.random()$
 $np.empty()$
 $np.arange()[i:, np.newaxis]$
 $np.random.randint(n, size=...)$
 $random.choice(a)$

$a.shape$
 $len(a)$
 $a.ndim$
 $a.size$
 $a.dtype$
 $a.dtype.name$
 $a.astype(int)$

$a-b$
 $np.add(a,b)$
 a/b
 $a*b$
 $a=b$
 $np.subtract(a,b)$
 $np.add(a,b)$
 $np.divide(a,b)$
 $np.multiply(a,b)$
 $a.sum()$
 $a.min()$
 $a.max()$
 $a.mean()$
 $a.median()$
 $a.corrcoef()$
 $np.std(a)$
 $np.log(a)$
 $np.argmax()$
 $np.any()$
 $np.all()$
 $np.exp()$

$a.sort(axis=)$
 $np.sort(a)$
 $np.argsort(a)$
 $np.partition(a,2)$

		actual	
		has disease	no disease
predicted	has disease	TP	FP
	no disease	FN	TN

	0	1	2	3	4	5	6	7	8	9	X[row, column]
0	11	8	2	3	2	9	9	1	9	1	ex[1,1]
1	5	6	5	3	6	8	2	6	8	1	ex[1,1]
2	1	3	7	4	5	3	8	7	3	3	ex[1,1]
3	7	4	4	9	9	6	6	1	4	9	ex[1,-2,9]
4	6	8	9	6	8	4	3	1	6	4	ex[1,4]

$axis=1$ row
 $axis=0$ column
 \cap intersection
 \cup union
 \wedge symm diff

kolomnormalisatie: $a/a.sum(axis=0)$
 rijnormalisatie: $a/a.sum(axis=1)$
 Z-normalisatie: $(df - df.mean()) / df.std()$
 RMSE: $1. np.mean((y - echt - y - gemeten)**2)$
 $2. np.sqrt(\text{antwoord } 1)$

```

def gcd(g, k):
    if k == 0:
        return g
    else:
        return (k, g % k)
    
```

```

def kgv(a, b):
    return int(abs(a*b) / gcd(a, b))
    
```

$df['w'].value_counts()$
 $len(df)$
 $df['w'].nunique()$
 $df.describe()$
 $sum()$
 $count()$
 $median()$
 $apply(function)$
 $min()$
 $max()/idxmax()$
 $mean()$
 $var()$
 $std()$
 $mode()$
 $df['w1'].corr(df['w2'])$

$df['w'].str.count(r'\?')$
 $str.len()$
 $str.contains()$
 $str.lower()$
 $str.upper()$
 $str.replace()$
 $df['w'].isin(values/column)$
 $pd.isnull(obj)$
 $pd.notnull(obj)$
 $2, 1, ~, ^, df.any(), df.all()$
 $df.dropna(subset=[])$
 $df.fillna(value)$
 $str.split('\')$
 $str.get(i)$
 $pd.concat()$

$df.values$
 $df.index$
 $df.columns$
 $df.sort_values('', ascending=)$
 $df.rename(columns = {'y': 'year'})$
 $df.sort_index()$
 $df.reset_index()$
 $df.drop(columns=[])$
 $df.shape()$
 $df.tail()$
 $df.head()$
 $df.plot(kind='barh') \rightarrow \log x = True$
 $list(df.items())$

$html: r'<[^>]*>'$
 $leestekens: r'\W\|s'$

reverse columns: $a[:, [2, 1, 0]] / a[:, :-1]$

reverse rows: $a[[2, 1, 0], :] / a[:, :-1, :]$

max length for each species: $iris.groupby('species')['length'].max()$

middlesquare:
 $k = L.shape[1]$
 $st = int((k/2) - 2)$
 $ei = int((k/2) + 2)$
 $L[st:ei, st:ei]$

vertical reshape
 $L.reshape(L.size, 1)$

euclidean distance
 $np.linalg.norm(a-b)$

notzero(L):
 $return L[np.nonzero(L)]$

diagonal
 $d = x2.shape[0]$
 $old = np.arange(d)$
 $dia = x2[dd, dd]$

$pd.read_csv('file', sep='\t', index_col=0, header=None, names=[""], compression='gzip', skipinitialspace=True)$
 $os.listdir()$

$df.pivot_table(index='Name', columns='Sex', values='Count', aggfunc=sum, margins=True).sort_values(by=['All', 'Name'], ascending=[False, True])$

$df.groupby(["", ""]).apply().(sum().)unstack()$

as column $L[:, -1]$ $shape = (r, c), ndim$ recall = $TP / (TP + FN)$ axis = 0 \rightarrow kolom $S.max()$ $idxmax()$
 $L[:, :2]$ size = data cellen precision = $TP / (TP + FP)$ axis = 1 \rightarrow rij \downarrow max waarde \downarrow locatie max waarde
 $L[-1, :]$ reverse

df = pd.read_csv('file', index_col='names')
 df.sort_index(inplace=True) of df = df.sort_values(['-', 'names'], ascending (pivot.sort) = [False, True])
 $a[(a <= 3) | (a == 5)]$

file[(file[-1] == 'Puro') | (file[-1] == '2010')] ['Vragen']

groupby all boys/jaar file[file['sex'] == 'M'].groupby('Year')['births'].sum() | groupby(['Year', 'sex']) kans

pivot df.pivot_table(values='count', columns='Sex', index='name', aggfunc=sum, margins=True)

cross tab x = file['party'].value_counts()[10] \rightarrow pd.crosstab(nvr['jaar'], nvr['party']) [x.index]
 nvr.dropna(subset=['-'], inplace=True) | nvr['-'].str.lower() | .str.replace('-', '') | .str.count('1?')
 df.fillna(0, inplace=True)

z-norm $(df - df.mean()) / df.std()$ RMSE \rightarrow np.sqrt(np.mean((L[:, 1] / L[:, 0]) ** 2))

mean df[['-', '-']].mean() | (df[['-', '-']] > 9).sum().sum() pd.Series

% kans ((df['-'] > 5.5) | (df['3'] > 5.5)).mean() * 100 | # narrow's per woord serie = serie.map(lambda x: len(x))

file[file['sex'] == 'F'].loc['John'] ['births'].plot() | file[~file.index.str.lower().str.contains('school')]

df.isnull().values.any()

accuracy = $\frac{TP}{TP + FN}$ precision = $\frac{TP}{TP + FP}$ recall = $\frac{TP}{TP + FN}$ precision = $(0.01 * 0.95) / ((0.01 * 0.95) + (0.99 * 0.05))$

numpy \rightarrow arange, reshape, transpose, sort, intersect1d, np.all, size, shape
 L[L > 5.5].mean() | kans op veldwaarde (L > 5.5).mean() | np.reshape(L, (np.size(L), 1))

+1 +2 +3 +4 L + np.arange(L.shape[1]) | np.arange(0, 0*n, 1*0) | normalise \rightarrow L / L.sum(axis=0)

loop kolommen for x in range(L.shape[1]): L[:, x][np.isnan(L[:, x])] = np.nanmedian(L[:, x])

fest all values Return np.all(L[(L > 10) | (L < -10)]) | %3 L[L % 3 == 0] | np.where[a == b]

reshape ((4, -1)) -1 doet automatisch # kolommen

Correlatie \rightarrow df[['-', '-']].corr() dan c = c[c != 1] rj = c.max(1).idxmax(1) ndim = c.loc[rj].idxmax(1)

Swap \rightarrow x = np.arange(9).reshape(3, 3) return np.warray(x[[1, 0, 2], :])

median \rightarrow float(np.median(df['-'])) | ratio \rightarrow pivot['-'] = np.log2(pivot[M] / pivot[F])

reshape Square \rightarrow x = np.sqrt(len(L)).astype(int) return np.reshape(L, (x, x))

reshape paren \rightarrow np.reshape(x, (int(np.size(x)/2), 2)) | np.concatenate([x, y, z]) merge arrays
 vertically stack \rightarrow np.vstack([x, grid]) | np.vsplit(grid, [2]) | np.ones((3, 3)) | np.arange(3)[:, np.newaxis]

Random \rightarrow np.random.randint(0, 10, (2, 3)) | x = array mag dan x[i] = 99 of x[i] = 10
 distance points \rightarrow i = x[:, np.newaxis, :] - x[np.newaxis, :, :] i.shape

oneer gehalten \rightarrow x[x % 2 == 1] = -1 | [np.repeat(a, 3), np.tile(a, 3)] | 2d array floats tussen 5-10 np.random.uniform(5, 10, size=(5, 3))

items of A that are not in B \rightarrow A[~A.isin(B)] | pd.Series(np.union1d(a, b)) | file['name'].value_counts()

convert array of \rightarrow df = pd.DataFrame(file.values.reshape(7, 5)) | df = a.append(b) | pd.concat([a, b], axis=1)

hoofletter size \rightarrow a.map(lambda x: x.title()) | iris.groupby('species')['sepal-length'].max() max per species

tokenize \rightarrow C = Counter(text.split()) CS = pd.Series(C).sort_values(ascending=False)

column with min by max per row \rightarrow df.apply(lambda x: np.min(x) / np.max(x), axis=1) %15 %CP

encl distance \rightarrow np.linalg.norm(a-b) | column met max # missing values x = df.apply(lambda x: x.isnull().sum()) x.argmax() 'Luggage.recon'

	predicted negative	predicted positive
negative	TN	FP
positive	FN	TP

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

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

$$F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

$$STD = \sqrt{\sum \frac{(x_i - \mu)^2}{n}}$$

square = x^2
 rooted = $\sqrt{\quad}$
 square rooted = $\sqrt{\quad}$

*.values → Series to Numpy
 *.T → transposes a Dataframe

% ls magic = list of magic
 % magic = handbook of magic
 % ls = list of files in directory
 % cd = change current directory
 % cp = keep file nor order file
 % rm = remove file
 % time, % timeit, % prun

+ add
 - subtract
 * multiply
 / divide
 // floor divide

np.random.randint
 np.newaxis
 np.where
 np.repeat
 np.concatenate
 np.sort

df.sort_index
 df.sort_values

pd.Series.value_counts

Regex

\ escape character
 \d any number
 \s space
 \w any character
 \W any except new line
 \b white space around word
 \n new line
 \s space
 \t tab
 \r return

create 1D Array = `np.arange(10)`
create bool array = `np.full((3,3), True, dtype=Bool)`
extract odd numbers = `arr[arr % 2 == 1]`
create new array with condition = `np.where(..., ..., arr)`
reshape array = `np.reshape(2, -1)` # -1 → automatisch
stack arrays vertical = `np.concatenate([a,b], axis=0)`
stack arrays horizontal = `np.concatenate([a,b], axis=1)`
repeat element 3 times = `np.repeat(a, 3)` # (111222333)
repeat array 3 times = `np.tile(a, 3)` # (123123123)
build arrays quickly = `np.R_[a,b]`
get common items = `np.intersect1d(a,b)`
from a remove all b = `np.setdiff1d(a,b)`
get all matches = `np.where(a==b)`
get between range = `a[(a>5) & (a<=10)]`
check values between 2 arrays = `np.vectorize(func)(types=[float])`
Swap columns = `arr[[2,0,1], arr[:, [2,0,1]]]` # swap rows
reverse rows = `arr[::-1]`
reverse columns = `arr[:, ::-1]`
random floats = `np.random.uniform(5,10, size=(5,3))`
limit decimal = `np.set_printoptions(precision=3)`
correlation = `np.corrcoef`
check for any nan = `np.isnan(A).any()`
get unique values in array = `np.unique(a)`
sort 2D array by column = `A[A.argsort()]`
max value each row = `np.amax(a, axis=1)`
only columns from read = `pd.read_csv(..., usecols=[...])`
reverse rows of df = `df.iloc[::-1, :]`
combine 2 df = `pd.merge(df1, df2, how='inner', left_on='fruit', right_on='peunds', suffixes=['left', 'right'])`

PANDAS

dataframe with 2 series:
`pd.concat([A,B], axis=1)`
give name to series
`A.name = "..."`
from ser1 remove items in ser2:
`ser1[~ser1.isin(ser2)]`
get percentile:
`np.percentile(A, q=[0, 25, 75])`
calculate frequency:
`A.value_counts`
reshape series to df with
`shape: pd.DataFrame(ser.values, reshape(7,5))`
get all values for condition:
`np.argwhere(ser % 3 == 0)` # of undere condition
take items from positions:
`ser.take([0, 2, 5, 8])`
stack two series:
`pd.concat([ser1, ser2], axis=0 or 1)`
get loc from items in ser2 in ser1:
`[pd.Index(ser1).get_loc(i) for i in ser2]`
map every element in series:
`ser.map(lambda x: func(x))`
make datetime: `pd.to_datetime(a)`
files in loader = `pd.read_csv(..., index_col=...)`
unique values = `a.index.value_counts`
order by index = `a.sort_index`

check duplicates: `len(set(cite.index)) == len(cite.index)`
kwadrat = `** 2`, wurzel = `np.sqrt(a)`
plot = `plot()` of `plot(kind='bar')`
pivot-table = `A.pivot(index="...", columns="...", values="...", aggfunc=[np.mean, np.median])`
cross-tab = `pd.crosstab(a, b)`

Parse HTML with `gzip.open(..., 'rt')`
`url = "..."`
`html_doc = request.get(url)`
`soup = BeautifulSoup(html_doc.content)`
find all comments
`comments = soup.findAll('li', class_='comment')`
beih struct. ... `prettify()`
make tuples list `[soup.title, c.p.text etc.]`
`df = pd.DataFrame(tuples)`
`df.columns = ['A', 'B']` `index=False`
near excel of csv
`df.to_excel, df.to_csv`

regex `findall('i')`
`" = r '\b[A-Za-z0-9_]{1,10}\b'`

PRECISION + RECALL

THERE IS A TEST FOR A DISEASE WHICH HAS AN ACCURACY OF 95%. THE DISEASE OCCURS IN ONE IN EVERY 100 PERSONS. THIS ACCURACY MEANS THAT THE TEST IS 95% RELIABLE.

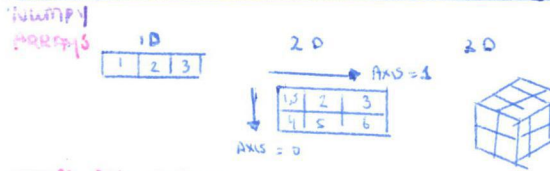
WITH ACCURACY (OR RELIABILITY) DOESN'T SAY MUCH WHEN:
 PREDICTED POSITIVE ACCURACY → 90%
 PREDICTED NEGATIVE ACCURACY → 90%
 PREDICTION ERROR = $CP \times .9 / (CP \times .9) + ((1-P) \times .1)$
 PRECISION/RECALL = $((1-PA) \times .9) / ((1-P) \times .9)$
 PREDICTION ERROR = $CP \times .9 / ((1-P) \times .9) + ((1-P) \times .1)$

WITH HEAVILY SKEWED DISTRIBUTIONS (E.G. DISEASE →)
 YOU CAN GET A HIGH RELIABILITY BY THE LAW OF LARGE NUMBERS
 * EVERY TIME YOU PREDICT SOMEBODY DOES NOT HAVE THE DISEASE → ACCURACY GOES UP
 * SAYING THAT NO ONE HAS THE DISEASE GIVES AN ACCURACY OF 95%
 BUT ALSO A PRECISION OF 0 AND A RECALL OF 0 TOO.

	HAS DISEASE	NO DISEASE	
POS TEST	95	495	590
NEG TEST	5	945	950
	100	990	

	ACTUAL	
	POS	NEG
PREDICTED POS	TRUE POS	FALSE POS
PREDICTED NEG	FALSE NEG	TRUE NEG

	PRECISION	RECALL
	TRUE POS	TRUE POS
	TRUE POS + FALSE POS	TRUE POS + FALSE NEG
	TRUE POS	TRUE POS
	TOTAL ADJUSTED POS	TOTAL ACTUAL POS
		$\frac{P \times R}{P + R}$
		F1 = 2 x



CREATING ARRAYS
 a = np.array([1, 2, 3])
 b = np.array([[1, 2, 3], [4, 5, 6]], dtype=float)

INITIAL PLACEHOLDERS
 np.zeros((3, 4))
 np.ones((2, 3, 4), dtype=np.int8)
 np.arange(10, 25, 5)
 np.linspace(0, 2, 5)
 np.full((2, 2), 7)
 np.eye(2)
 np.random.randn(2, 2)
 np.empty((3, 2))

INSPECTING YOUR ARRAY
 a.shape
 len(a)
 b.ndim
 c.size
 d.dtype
 e.dtype.name
 f.dtype.kind

ASKING FOR HELP
 np.info(np.ndarray.dtype)

ARRAY DIMENSIONS
 LENGTH OF ARRAY
 NR OF ARRAY DIM
 NR OF ARRAY ELEMENTS
 DATA TYPE OF ARRAY ELEMENTS
 NAME OF DATA TYPE
 CONVERT AN ARRAY TO A DIFFERENT TYPE

ARRAY MATHS MARKS
 G = A - B (SUBTRACTION)
 G = B + A (ADDITION)
 A / B (DIVISION)
 A * B (MULTIPLICATION)
 np.exp(B)
 np.sqrt(B)
 np.sin(A)
 np.cos(B)
 np.log2(A)
 np.dot(C, F)

COMPARISON
 A == B (ELEMENT-COMP)
 A < B (ELEMENT-COMP)
 np.array_equal(A, B) (ARRAY-COMP)
AGGREGATE FUNCTIONS
 a.sum()
 a.min()
 b.max(axis=0)
 a.cumsum(axis=1)
 a.mean()
 b.median()
 a.var(axis=0)
 np.std(B)

CUBBETTING, SLICING, INDEXING
 a[2]
 b[1, 2]
 a[0:2]
 b[0:2, 1]
 c[2, 0:2]
 b[:, 1]
 c[:, 2, 3]

ARRAY MANIPULATION
 np.transpose(C)
 b.ravel()
 G.reshape(3, 2)
 h.resize((2, 4))
 np.append(A, B)
 np.insert(A, 1, 5)
 np.delete(A, [1])
 np.concatenate(A, B)
 np.vstack((A, B))
 np.r_[C, E]
 np.hstack((C, E))
 np.column_stack((A, B))
 np.r_[A, B]
 np.split(A, 2)
 np.all
 np.any
 np.argmax
 np.argmin

SORTING ARRAYS
 a.sort()
 c.sort(axis=0)

2D ARRAY LOOP LOGIC
 row-len = len(shape[0])
 col-len = len(shape[1])
 for row in range(row-len):
 for col in range(col-len):
 if condition (e.g. x > y):
 L[x][y] = np.median(L, axis=0)[y]

DATA ALIGNMENT
 s3 = pd.Series([F, -2, 3], index=['a', 'c', 'd'])
 s + s3
 a 10.0
 b NaN
 c 5.0
 d 7.0

PANDAS SERIES
 s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])

DATAFRAME

	Country	cap	pop
0	Belgium	Brussels	10
1	India	New Delhi	100
2	Brazil	Brasilia	200

data = {'Country': ['Belgium', 'India', 'Brazil'], 'Capital': ['Brussels', 'ND', 'Brasilia'], 'Population': [10, 100, 200]}
 df = pd.DataFrame(data, columns=['Country', 'Capital', 'Population'])

READING
 pd.read_csv('file.csv', header=None, names=['a', 'b', 'c', 'd'])
 pd.read_excel('file.xlsx')
 pd.read_html('file.html')

SELECTION
 df[['a', 'b']]
 df[['a', 'b', 'c', 'd']]
 df[0:2]
 df[['Country', 'Cap', 'Pop']]
 df[['India', 'New Delhi', 100]]
 df[['Brazil', 'Brasilia', 200]]

DROPPING
 df.drop(['a', 'c'])
 df.drop(['Country'], axis=1)

PROF VALUES FROM ROWS (axis=0)
PROF VALUES FROM COLS (axis=1)

DATA ALIGNMENT
 s3 = pd.Series([F, -2, 3], index=['a', 'c', 'd'])
 s + s3
 a 10.0
 b NaN
 c 5.0
 d 7.0

BY POSITION
 df.iloc[0], [0]
 'Belgium'
 df.iat[0], [0]
 'Belgium'

BY LABEL
 df.loc[0], ['Country']
 'Belgium'
 df.at[0], ['Country']
 'Belgium'

BY LABEL / POSITION
 df.ix[2]
 'Brazil'
 df.ix[0], ['Capital']
 'Brussels'
 df.ix[1], ['Capital']
 'New Delhi'

FULL METHODS
 s.add(53, fill_value=0)
 a 10.0
 b -5.0
 c 5.0
 d 7.0
 s.sub(53, fill_value=2)
 s.div(53, fill_value=4)
 s.mul(53, fill_value=3)
 df.fillna(0, inplace=True)

SETTING
 s['a'] = 6

APPLYING FUNCTIONS
 f = lambda x: x * 2
 df.apply(f)
 df.applymap(f)

ASKING FOR HELP
 help(pd.Series.loc)

2 ITEM CORR
 c = df[['a', 'b']].corr()
 c = c[0, 0]
 z1 = c.max(1).idxmax()
 bottom = c.loc[z1].idxmax()
 wday10 = z1, bottom3

SETTING
 s['a'] = 6

ASKING FOR HELP
 help(pd.Series.loc)

ASKING FOR HELP
 z1 = (df - df.mean(axis=0)) / df.std()

ASKING FOR HELP
 e = (df[['a', 'b']] < 5) | (df[['a', 'b']] > 10)
 w perc (...) mean() * 100

precisie/recall
 pr = hoeveel correct in voorsp. C = $\frac{TP}{TP+FP}$
 rc = hoeveel daadwerk C correct = $\frac{TP}{TP+FN}$
 acc = hoeveel correct totaal = $\frac{TP+TN}{\text{totaal}}$

acc = 0.9 ziek = 0.1
 pziek = $\frac{0.1 \cdot 0.9}{(0.1 \cdot 0.9) + (0.9 \cdot 0.1)}$
 pnziek = $\frac{0.9 \cdot 0.9}{(0.9 \cdot 0.9) + (0.1 \cdot 0.1)}$

%ls -th = directory + size %paste / paste
 %rm file = verwijder %run = extern
 %time %timeit %prun = resources
 %variabele? = type pr(-) / out [2] / -2
 %history -n 1-4
 !pwd / !ls / !cd / !mkdir / !mv
 %% = modulo remainder
 // = floordivision (cuts low round)

NUMPY

D1. Shape (values)
 D2. Shape (row, column)

• Shape = vorm → shape [0] / [1] = rijen / kolommen
 • size = grootte / aantal values
 • reshape () → plat maken 2D = 1 waarde met de size of • flatten ()
 • arange (start, stop, step, int) alleen necessary

Veranderen = 1D (-1 etc) → v.b. X1 - X1.mean

np.sqrt () = $\sqrt{\quad}$
 X1.std () = stand.

Boolean mask v.b. X1[X1 % 2 == 0] → geeft False / True
 X1 < 3 = False en True

np.sum () geeft aantal Trues

np.nonzero () : geeft locatie non zero → array [n np.nonzero (a)]
 = de waarden

NORMALISEREN
 v.b. $\frac{2}{12}, \frac{4}{12}, \frac{6}{12} = \frac{p}{p.\text{sum}(\text{axis}=0)}$
 opvall. ↓

broadcasting

- matching dimensies kleinste wordt groter
- shape van array met shape wordt gestretched
- geen shape is 1? → kan niet

np.all > geeft True of False
 np.any

np.sort (a, axis=1)
 = sort de rijen op grootte

v.b. toets a = np.arange (1, 4) b = np.arange (1, 4) [:, np.newaxis]
 a = [0, 1, 2, 3] b = [[0], [1], [2], [3]]
 a = [[1, 2, 3]]
 b = [[1, 1, 1] [2, 2, 2] [3, 3, 3]]
 a + b = [[2, 3, 4] [3, 4, 5] [4, 5, 6]]

RMSE = $\sqrt{\frac{\sum (y - \hat{y})^2}{T}}$
 y = gem aantal waarden
 T = totaal waarden

Pandas

pd.read_csv ('file') → bekijk opties voor naam index etc.
 df.sort_index (ascending)
 voorwaarde (df[['lijst k']]) = / < / etc.

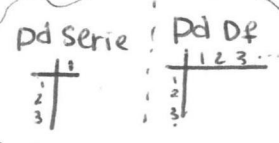
fillna (value)
 dropna (zie condities) → NaN

nieuwe kolom

df[knaam] v.b. DF[['lijst k']].mean (1) geeft gem rij van de kolommen

v.b. len (df[['k']]) / len (df) = % → & doeper stuk

df.sort_values (lijst, bool lijst)



idxmax return index van max val. of k
 corr uit DT
 c = df[['kol']].corr ()
 c = c [c != 1] ; voorwaarde
 rij = c.idxmax ()
 kol = c.loc [rij].idxmax ()
 coord = [rij, kol]

df[kol].value_counts () : aantal v. in kolom

z.normalise = $\frac{x - x.\text{mean}}{x.\text{std}}$

voor kolom df['k']
 voor index df.index
 voor values df.values

df[df['k'].str.contains('w')] groupby?
 df.groupby (kolom) [k] :
 • max / sum etc
 ↳ geeft max van k1
 df[k] = [re.sub (r'...', '', s) for i, s in df[k]]
 allen | - | = woorden
 [^A-Za-z0-9]+ = alleen letterse cyfers

pd.merge
 → corrects index
 on = zelfde kolom
 left-on / right-on
 np.concatenate
 pd.concat