

Sensitivity, recall, hit rate (TPR)

$$TPR = \frac{TP}{P} = \frac{TP}{TP+FN}$$

Precision (PPV)

$$ppv = \frac{TN}{N} = \frac{TN}{TN+FP}$$

Accuracy (ACC)

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

F1 score

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

Nuttige functions numpy

- np.concatenate
- .vstack/hstack
- .r_
- .where() returns indices

select (=>) pandas select

pandas pivot crosstab

marginal = total
 agg = np.[mean, sum, median]

np.array [row, col]

- np.tile
- np.reshape
- np.repeat

from numpy import dot
 np.lookfor('query')

$$FP = \frac{FP}{100}$$

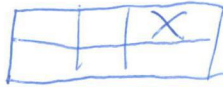
$$\frac{TN}{NN+FP}$$

agg 0.98
0.98

n



NP, array

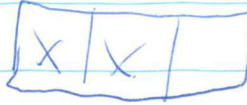


NP .shape

a[0, 2]

.size

a[0:2]



.add(a, b)

.arange(0, 100, 10)

a[a < 2]



a +-*/ b

.reshape (keep)

.resize (delete) - ~~resize()~~

from numpy import doc

Pandas

? np.doc

np.lookfor

df = pd.DataFrame([[1, 2, 3],
[4, 5, 6],
[7, 8, 9]],

pd .query('a >= 100')

index = [1, 2, 3],

.rename('val', 'value')

columns = [a, b, c]

.melt(df) to

~~print(df)~~

df.pivot(col=, row=) to

	a	b	c
1	1	2	3
2	4	5	6
3	7	8	9

Concat ([df1, df2]) axis=0 rows
axis=1 columns

a+b = error?
broadcasting apply

.drop (columns =)

df[df.length > 7]

.dropna

df[df['length'] > 7]

.sort_values

.value_counts()

.sort_index

.unique()

.head

.count

.tail

sum, mean, min, max, med, std

.fillna(x)

.plot, hist

.nlargest (n, 'value')

.read_html (link)

.nsmallest (n, 'value')

.groupby (by = "col")

CHEAT SHEET #1



NP (y, x) np.zeros() - np.ones() - np.full() -> can shape
 ↳ n_j ↳ kolumn
 np.arange() -> cannot shape, is linear & np.linspace() evenly spaced between given values.
 else: np.random.random() - np.random.normal() - np.random.randint()
 => creating arrays
 x2 [0, 0] = 12 -> changing value of that index to 12 - x2.astype('float64') -> got to assign to new variable - x[start:stop:step] - x[::-1] -> all elements reversed - x[:, :2] -> all rows, every other column
 accessing column: x[:, 0] -> first column of x - accessing row: x[0, :] -> first row of x. - if you modify a subarray, the original is changed -> use .copy
 combination of np.arange() - reshape -> can get a multidimensional something
 np.concatenate(array1, array2, axis = 1) : joining two or more arrays
 np.hstack() & np.vstack() -> add the array according to columns
 np.mathematics: $x_1 - x_2$, $x_1 + x_2$, x_1 / x_2 , $x_1 * x_2$, $x_1 ** x_2$, np.sqrt(), np.log(), -> must be the same shape, np.floor_divided(), %
 Python's built-in absolute value abs() -> np.absolute()
 np.reduce() -> applies a given operation to elem of an array until only a single result remains, using np.cumulative() stores all intermediate results of the operation - np.multiply.outer(x, x) (x = np.arange(1, 6)) -> creates multiplication table.
 x3.min(axis = 0) -> min value across each column - np.sum / prod / mean / std / max / min / var / argmin / argmax / median / any / all have nan-safe versions, just put nan in front & nan_safe -> finds index
 Broadcasting: array - array is element by element basis if arrays are same size
 (m, n) + (4, n) -> copies m times to (m, n) use [:, np.newaxis]
 (m, 1) * (m, 1) -> copies n times to (m, n) np.sum([inches > 1 & inches < 1]) -> how many entries
 np.sum(x < 6) -> counts how many are bigger than 6
 np.sum(x[x < 6]) -> counts values in x, if smaller than 6 [:, np.newaxis]
 fancy indexing: shape of result reflects shape of indexing array i.e. shape of array being indexed - np.sort() & np.argsort(x), np.partition(x, 3) -> smallest 3 values to the left
 LET OP AXIS: + 2 [] obj, c looks in masks merge

PD df: index, df: value - df: column = list(series.items()) - df['new column'] = ... string by explicit index series ["a": "c"] -> find index = included -> implicit index = data.loc[1] -> finds explicit index -> data.iloc[1] finds implicit index - add new column: data['density'] = data['pop'] / data['area']
 swap rows & columns: df.T - passing a single index to df. evaluate gives you a row data['area']
 passing a single index to df passes a column, use .loc to slice a df. - the ix, indexer allows you to do both! data.loc[data['density'] > 100, ['area', 'density']] -> columns you want in resulting dataframe. Indexing refers to columns & slicing to rows
 A.add(B, fillvalue = 0) -> makes sure adding doesn't result in NaN, use A.stack.mean() to get the overall mean, ^{row} use: A + B column wise add(B, axis = 0)

isnull(), dropna(), fillna(), notnull()
 ↳ axis = "columns" drops all columns containing null values, how = all -> drops where both row or column is nan, use stack() multi-index -> regular and vice versa
 data.sort_index() - pd.concat(axis = ..., ignore_index = True/False, join = inner, append) / pd.merge(on = ..., left_on = "employee", right_on = "name", df.join -> merge that defaults to joining on index.
 df.set_index("id"), pd.merge(df1, df2, left_index = True, right_index = "name"), merging df's with not much in common? pd.merge(df1, df2, how = "outer") also left and right, merge with overlapping columns
 names? pd.merge(df1, df2, on = "name", suffixes = ["_L", "_R"]), pd.read_csv - df.value_counts
 lambda x: x + 10, isnull().any() -> check if any value is missing, np.argsort() / sort
 df.sort_index(by = 'F') - df.groupby('birth') - df.groupby('opposing') ['aantal kinderen'].sum()

precision: how many of pred of class C were correct? TP / (TP + FP) recall: how many instances of class C predicted to be C? TP / (TP + FN), accuracy: how many of all values were correct TP + TN / (TP + TN + FP + FN)
 sort_values(by = 'F') regex: str.replace(['^a-20-9'], '') / ['^a-2A-20-9'] + str.replace("allen | border", "")
 Multi-index: pop[2, 200]
 neg. pos neg TP; kans op iets & accuracy FP: (1 - kans op iets) & (1 - accuracy)
 pos v TP FP
 neg v FN FN
 from numpy import dot np.dot(a, b)

NP Calc. X[start, stop, step] (-1 for reverse)

mean/mode/median
add/subtract/mod/power
sum/min/max (arg/max)
any/all/percentile
sort/partitioning (arg for k-nearest)

reorder

random/copy

arrange

reshape/concatenate

shape

(V/H) (stack/split)

transpose

roll

linspace

full

test

unique

reduce/accumulate

over

Pd (Series / DataFrame)

read_csv

crossref/pivottable

dropna

value_counts

sort_index / sort_values

mean/mode/median

corr

groupby.fun_

nlargest

STR

replace

lower

strip

contains

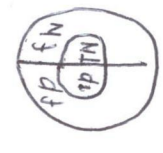
count

len

I

ix (index into elements mix)
add

Pd.Series.sort.IETS?



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

np.zeros ((3,4))
np.ones ((2,2))
np.full ((2,2), 7)

a.shape = array of dimensions
len(a) = length of array
b.ndim = num array dimensions
e.size = num array elements

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

np.exp(b)
np.sqrt(b)

a.sum()
a.min()

a.sort() → sort an array
c.sort(axis=0) → sort the elements of an array's axis
np.sort(c,axis=0) → sort each column of c

b.max(axis=0)
↳ max of row

a.mean()
np.std(b)

np.mean(x) / np.mean(x[0])

a.sum()

np.mean(x,axis=1)
↳ gem.

• mean(1)
= gerekend oer kolommen

P is een multidimensionale array
P.flatten() → nu 1 dim

np.all → all elements along a given axis must be true
np.any → any

L[np.nonzero(L)] → alle elementen die niet nul zijn uit L

x.namen.value_counts()

naam1 370
naam2 198

b = np.arange(14)[:, np.newaxis]
array([13, [2], [3]])
ipv. [1,2,3]

np.sum(p < n)
= hoeveel waarden kleiner dan p

a[:, [1,0,2]]
kolommen nieuwe volgorde

np.isnan

> bool, returned True / False

df.iloc[2,1]
↳ selecting specific data

df.loc[df['a'] > 10]
↳ selecting rows meeting logical condition

df['W'].value_counts()
↳ count number of rows

df["New Column"]
= (a.4 * 2) + (0.6 * 1)

data[(data > 0.3) & (data < 0.7)]

data[~data.isin(data2)]
↳ data die niet ook in data2 staat

df.groupby(['Animal'], mean)

KVR = KVR[~KVR.Party.str.contains("vragen", na=False)]

KVR = KVR[~KVR.Party.str.contains("vragen", na=False)]

KVR.Party.value_counts().sort_index()

horizontal = columns : axis = 1
vertical = rijen : axis = 0

a[::-1] > reversed array a

a[a < 2] > select elements from a less than 2

i.transpose() → rijen en kolommen versuisselen

g.reshape((3,3)) np.abs()
↳ geen negatieve

np.append(h,g)

np.delete(a,[1])

np.concatenate((x,y),axis=0)
↳ samenvoegen arrays

np.intersect1d(a,b)
↳ get common values

np.setdiff1d(a,b)
↳ from a remove all b

np.where(a == b)
↳ positions where a & b match

df["Ass.mean"] = df[df.columns[0:6]].mean(axis=1)

all = df[df.columns[0:6]]

greater-than-g = all[all > g].count().sum()

df.sort_values(['Name', 'Age'], ascending = [False, True])

df.sort_values(['Name', 'Age'])
df.sort_index()

df.drop(columns = ['Length'])
pd.concatenate([df1, df2], axis=1)
↳ reshaping

df3 = pd.merge(df1, df2)
↳ combine datasets

df.fillna(0)

df.dropna(subset = ['Age'], inplace = True)

df.values() → raw data

df.transpose() → rijen x kolommen om draaien

len(df) → num of rows

str.lower() / str.replace()

x = x.drop(6, axis=1) → kolom 6 weghalen

Np.array([1,2,3])

a.ndim → Dimensions

Np.info

a[2] = [1 2 3] 2e item

a[1,2] =

1	2	3
4	5	6

↑ ↑
Rij kolom

B[0:2,1] =

1	2	3
4	5	6

↑ ↑
Rij kolom

a[0:2] = [1 2 3]
OT 2

a[:, :-1] = Reverse

np.arange(3) = np.array([0,1,2])
Start, Stop, step

np.reshape((3,2)) =

1	2
2	3
3	4

↑ ↑
Rij kolom

a[:, :-1] =

1	2
3	4

↑ ↑
Rij kolom

np.nditer = iter array

np.mean

np.transpose =

1	2
3	4

 =

1	3
2	4

np.ndimmax = [[1,2] [3,4]]
↑ ↑
① ②
↑ ↑
POS NR

pd.read_csv(link, sep=',', na=)

Knr. Value-Counts = 2.0 → 2x etc

vor = knr.loc[knr["room"].min

pd.Ground = [Rij, Kolom]

Knr.dropna(subset=["room"])

Knr["room"] = knr["room"].str.lower()

.str.replace('W', '1').str.replace('1', 'W')

R =

1	2	3
4	5	6

Knr.loc[Knr["p"] > 4]

vor = knr["p"].str.contains('1')

vor2 = vor.where(vor <= 50)
mean(), median(), mode()[0]

vor_corr(knr["p"], knr.loc[...])

pd.pivot_table(df, values, index, columns, margins)

not_values [A, P] or = [T, F]

.between(-1, 1)

idxmax() = id barisan kolom

Pandas

- df.sort_values(['col', 'col'], ascending=[F, T])
- df.fillna('waardoor?', inplace=True)
- pd.read_csv('naam', index_col='naam', names=[...])
- df[] > of < 'waarde'
- pd.concat([df1, df2]) → append rows
- pd.concat([df1, df2], axis=1) → append columns
- df.rename(columns={'y': 'year'})
- df.dropna() → drop rows where at least 1 element is missing
- df.dropna(how='all') → where all elements are missing
- df.loc[:, 'x2': 'x4'] → select all columns between x2 and x4 incl
- df.drop(columns=['a', 'b'])
- pd.isnull()
- pd.merge(adf, bdf, how='left'/'right'/'inner'/'outer', on='x1')
- df['w'].nunique() → # of distinct values in column
- pd.notnull() / notna()
- .str.lower()
- .str.contains('string')
- ~~str~~ series.str.startswith('string')
- series.str.replace('iets', 'door dit')
- series.value_counts(sort=True, ascending=, dropna=)
- df.drop_duplicates([columns], inplace=)
- df.iloc[:, [1, 2, 5]] → select cols in positions 1, 2, 5

Numpy

- np.intersect1d(a, b)
- np.setdiff1d(a, b) verwijder uit arrA wat in B zit.
- np.arange(start, stop, step)
- np.isnan(array)
- np.nanmedian(array, axis=)
- np.nanmean(array, axis=)
- arr.flatten() → maakt van MA 1D array
- np.nonzero(arr)
- np.sort(arr, axis=)
- np.sqrt(x)
- np.multiply(x1, x2)
- np.vstack([a, b])
- np.hstack([a, b])
- np.argmax(x > 1) → geeft indices van values die voldoen
- np.random.choice(5, 3) → kiest 3 dingen uit np.arange(5)
- np.where(condition)
- np.unique(arr, return_index=, return_counts, axis)
- np.linspace(0, 2, 9)
- arr = arr[~np.isnan(arr)]

Actual

	N	Predicted P
Negative	TN	FP
Positive	FN	TP

- plt.bar(x, height, color, label)
- plt.xlabel('string', fontweight, fontsize)
- plt.xticks(np.arange(x), rotation)
- plt.barh(x='ministerie')

Regex

- \d → matches any decimal digit
- \D → matches any char that isn't dec digit
- \w → matches any word char
- \W → matches any non-word char
- \s → matches whitespace chars
- 'length\$' → matches strings ending with 'length'
- '^sepal' → matches strings beginning
- ? → 0 or 1 reps
- + → 1 or more reps
- * → 0 or more reps

accuracy = $\frac{TP+TN}{total}$

precision = $\frac{TP}{TP+FP} = \frac{(P \cdot acc)}{(P \cdot acc) + (1-P) \cdot (1-acc)}$ = Precision niet

Recall = $\frac{TP}{TP+FN} = \frac{(P \cdot acc)}{(P \cdot acc) + (1-P) \cdot acc}$ = Precision niet

lambda x: x * 10

<u>Numpy</u>	<u>Pandas</u>	<u>Regex</u>
np.setdiff1d	df.columns	1 = match any
np.equal	df.value_counts	a-z = range
np.reshape	df.index	\d = digit
np.transpose	df.values	v = not
np.argwhere	df.crosstab(df.X, df.Y)	\w = word
np.unique	df.copy	= or
np.equal	df.dropna(subset=['x'])	* = 0 or more
np.logical	df.astype(str)	. = any char
np.where	df.contains	\D = non digit
np.shape	df.nunique	r' = regex
np.ndim	df.sort_index	\n = newline
np.size	df.loc	
np.idxmax	df.count	
np.nonzero	df.describe	
np.absolute	df.pivot_table	
	df.concat	

```
pd.MultiIndex.from_tuples([('a', 2000), ('b', 3000)])
      ryen      kolom
      ↓        ↓
      L[:, 0]
```

```
[:, newaxis]
```

P / P.sum(axis=0) = normalise

$$((1-z) * a) / (((1-z) * a) + (z * (1-a))) = nP$$

$$(a * z) / ((a * z) + ((1-z) * (1-a))) = P$$

a = accur
z = zick

$$\text{Precision} = \frac{TP}{TP+FP} \quad \text{Recall} = \frac{TP}{TP+FN}$$

round(, 1)
 %ls
 %cp
 %time
 %timeit
 a= set
 a.iTab

Recall = TP / (TP + FN)
 precision = TP / (TP + FP)

	has dis	no dis
posit	TP	FP
neg ts	FN	TN

95 495 590
 5 9405 9410
 100 9900 10000

acc = -- % reliable
 $0.01 \times 0.95 / (0.01 \times 0.95 + 0.99 \times 0.05) = 0.16$

isin() from collections import Counter
 %matplotlib inline
 value_counts()
 drop() = dropna() => how='all' = drops rows/cols = all nan thresh=3
 unique() subset=['...']
 sort_index() sort_values() (by=[,], asc=[, f])
 kv = kv [kv >= 1]

for i in range(L.shape[0]):
 LE[i][np.isnan(LE[i])] = np.nanmedian(LE, LC?)
 to_datetime

inplace = True
 idxmax
 copy
 print
 print("%s %s" % (per, mis))
 merge
 melt
 query()
 unique()
 clip() bu. orderse en borsew 10

plt.plot(kind='bar' / 'barh')
 plt.hist
 plt.scatter
 lineStyle=steps
 plt.xlabel('...')
 plt.legend('...')
 plt.title('...')
 plt.show()
 pd.MultiIndex?
 pd.merge() on='...' how='inner/outer'

Str
 .len()
 .contains('...')
 .replace
 .count
 .strip()
 .loc [3]
 .iloc [3:2]
 .ix [3, 'pop']
 ascending
 .largest
 .smallest(n, value)
 .filter
 .stack()
 .unstack() level=0/1
 .f.latter()

Set((...) & | ^ (...))
 concat (2 df's samenvoegen); ignore_index=True; join='inner/outer'; append()
 index col => no sort_val => eerste
 groupby('year', 'sex') [count] => by=, level=
 aggregate ('first'). unstack()
 index_col = ...
 reset_index() method='ffill' -> use val
 fillna(0 or 1) 'bfill' -> use val :: 2 = 2 stappen
 mean() etc. maar 0 = col 1 = ry
 isnull() notnull() cumsum min max prod
 apply()

nieuw col = df['...'] = ...
 sort_values('...', axis=1, ascending=[, f]) // floor/div / div / pow / mod
 corr()
 index dtype
 describe
 items()
 any() head() tail() * * pow / % mod
 all() shape()
 transpose() added() level='year'
 fill_value=0

np.array(L...)
 np.arange(10) start, stop, step
 zeros ones
 full(3, 5, 3.14) empty
 linspace() start, stop, hoeveel dus 1, 2, 3 = 1, 1.5, 2
 eye(3) = 1 0 0 / 0 1 0 / 0 0 1 . eye(3, k=1) = 0 0 0 / 0 1 0 / 0 0 1
 add.at(x, i, 1)
 np.random.randint(low=, high=, size=) np.partition argpartition
 random(3, 3) 2d array sorted() gear np!
 normal
 rand

ndim
 shape
 size
 dtype
 setdiff1d
 count_nonzero
 reshape()
 asarray . astype
 np.newaxis
 row/col
 & | ^ ~ => niet => np.bitwise_not
 logical_and
 unique
 isna
 reverse 2d:
 x[::-1, ::-1]
 first row/col
 x2[:, 0] col
 x2[0, :] row
 np.newaxis

Broadcasting: -> moeten zelfde loc hebben
 np.arange(3) * 5
 0 1 2 -> 5 6 7
 np.ones((3, 3)) + np.arange(3)
 1 1 1 + 0 1 2 1 2 3
 1 1 1 => 1 2 3
 1 1 1
 np.arange(3).reshape((3, 1)) + np.arange(3)
 0 1 2 + 0 1 2 0 1 2
 1 2 3 + 1 2 3 2 3 4

absolute/abs
 sin/cos/tan
 exp . log
 exp2 . log2
 power(3, x) . log10
 reduce()
 accumulate()
 outer(x, x)
 sum() . min() . max()
 mean . prod
 nanmean . std . var
 argmin, max
 median
 percentile
 copy
 from scipy import stats mode
 modus = insert delete

Precision & Recall

- EVALUATING

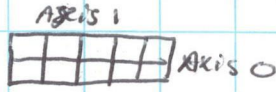
- CLASSIFICATION: ACCURACY, PRECISION, RECALL, F1
- REGRESSION: RMSE

reliability: DOES NOT ACCOUNT FOR ERRORS GOOD

ACCURACY: How many of all predictions correct

RECALL: How many instances of class C predicted to be C $TP / (TP + FN)$

PRECISION: How many predictions of class C are correct $TP / (TP + FP)$



NUMPY

- CREATE:
- NP. ARANGE
 - NP. Linspace
 - NP. Full((2,2), 7) constant
 - NP. RANDOM. RANDOM((2,2))

- INSPECT:
- A. SHAPE
 - LEN(A)
 - A. NDIM
 - A. SIZE

- ARITH:
- NP. SUBTRACT
 - NP. ADD
 - NP. DIVIDE
 - NP. MULTIPLY
 - NP. EXP SORT LOG SIN COS

- A. SUM()
- A. MIN
- A. MAX(axis=1)
- A. MEAN()
- A. CORRCOEF
- NP. STD(A)

- SICES:
- A[1,2]
 - A[[!0,1],[0,1,2]] → (1,0), (0,1), (1,2)

- MANIPUL:
- A.T
 - A. RAVEL (FLATTEN!)
 - A. RESHAPE
 - A. SORT(axis=0)
 - NP. ARGWHLRRE Any char consider W word, digit, whitespace WIDIS Not word, Not digit, not [ABC] Any of ABC [^ABC] Not A, B or C [A-g] char BETWEEN A-g * + ? 0 or more, 1 or more, 0 or 1 a{2,3} 2 of next a's

x matplotlib inline

% MAGIC → explains magic commands.

W = !LS → LEN(w) is level border

PANDAS

DF = PD. DATAFRAME([[...]], INDEX=[1,2,3], COLUMNS=['A','B','C'])

- DF. SORT_VALUES
- DF. SORT_INDEX
- DF. RENAME
- DF. RESET_INDEX
- DF. ILOC[10:20] rows
- DF. LOC[:, x2:x4] row col max colnes
- DF. LOC[df.a > 10, ['a', 'c']]
- DF. W. VALUE_COUNTS
- LEN(DF)
- DF. W. NUNIQUE
- DF. PROPNA
- DF. FILLNA(only)
- DF. GROUPBY(by='cd'). SUM() su.
- PD. PIVOT_TABLE
- DF. IDXMAX is with max val

MISC

- MERGE
- PD. MERGE(a, b, How='left', on='x1')
 - PD. PLOT(title="")
 - DF. PLOT.HIST()
 - DF. PLOT.BARH()
 - DF = DF[DF. NOTNULL()]
 - DF. A. STR. REPLACE(REGEX, "")

$$\text{Precision} = \frac{tp}{tp + fp}$$

$$\text{Recall} = \frac{tp}{tp + fn}$$

$$\text{Accuracy} = \frac{t}{tF}$$

Numpy array indexing.

X[0, 0]
↑ ↑
row column

[start: stop: step]

Slicing

One-dimensional

X[:5] eerste vijf

X[5:] alles na vijf

X[::2] om de twee

X[::-1] omgedraaid.

Multi-dimensional

X[:2, :3] stop bij de tweede rij
stop bij de derde kolom

X[3, ::2] stop bij de derde rij
om de kolom.

X[0, :] rij 0, pau alle volkomen

Filtering pandas

a	0,25
b	0,50
c	0,75

= data[(data > 0,3) & (data < 0,8)]
data.values

pd.read_csv('file', compression='gzip', sep='/', delimiter='\t', header=None, index_col=0)

fillna() → staat niet op

df.iloc[:, -4:-1].sum(axis=1)

(df[[ASS1, ASS2, ASS3]] > 9).sum()

df.sort_values(['deeltvoet_meen', 'names'], ascending=[False, True])

c = df[[ASS1, ...]].corr()

c = c[c!=1]

rij = c.max().idxmax()

wilom = c.loc[rij].idxmax()

vlo = {rij, wilom}

numpy:

df.normalise(p):

return (p / (p.sum(axis=0)))

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

s[d, d], np.diagonal(s)

acc = TP / (TP + FN)

precision = TP / (TP + FP)

recall = TP / (TP + FN)

cit. str.lower(), str.strip(), str.replace(" ", " ").str.replace

regex: str.replace(LR('[^a-z\d]', ''))

pd.crosstab(kvr.partij, kvr.jaar) [most common]. index.array

for loop

maxline = ''

for line in fin:

if len(line) > len(maxline):

maxline = line

df.sort_index(inplace=True)

inplace=True

| = of teksten

perc = ((df[ASS] < 5.5) / (df["ass"] < 5.5)).mean() * 100

normaliseren = (df - df.mean()) / df.std()

df.groupby('columns')['income'].mean()

df.value_counts().index

pd.pivot_table(df, columns='soort')

index='Name', aggfunc='count': np.sum, values='count')

- iloc [0] impliciet
- loc ['a'] expliciet
- idxmax()
- isin()

np.reshape(L, shape)

np.transpose(L)

np.where()

np.flatten() → 1d

np.sum(p < n)

↳ van p, hoeveel daarvan strict kleiner zijn dan n

np.sqrt(np.mean((L[:, 1] - L[:, 0])**2))

np.sort(L, axis=1) → sorteert per rij, van laag → hoog

np.all(np.abs(L) > 10)

L[i, c]

L[np.nonzero(L)]

np.arange(1, 4)

np.arange(1, 4)[:, newaxis]

np.random.randint(0, n)

• append

np.arange(n)

np.reshape(L, i, c)

L.size

last column = L[:, -1]

n = len(L[0])

L[int(n/2)-2:int(n/2)+2, int(n/2)-2:int(n/2)+2]

been:

y-lim, x-lim = grid.shape, (x+xx)%x-lim

grid[(y:yy)%y-lim, (x+xx)%x-lim for xx in (-1, 0, 1) for yy in (-1, 0, 1)]

if not (xx == 0 and yy == 0)

([0, 2])([0, 2])

np. arange(0,0, step) - np. size
 Array. reshape (c,m) np. ndim
 np. astype

~ np. array (L%N, dtype=bool)
 np. transpose (array) $\begin{bmatrix} [1,2], \\ [3,4] \end{bmatrix} \rightarrow \begin{bmatrix} [1,3], \\ [2,4] \end{bmatrix}$

Aggregate Functions
 a. sum() sum array
 a. min()
 b. max(axis=0) max, value per row
 np. sort(axis)

manipulation
 np. concatenate ()
 np. vstack () vertical
 np. r_ [e,f] stack vertical rowwise
 np. hstack
 np. column_stack
 np. hsplit (a,3)
 np. vsplit (a,3) split at index

For y, x in zip (Array, Array)
 np. linspace (0,2,9) $\rightarrow [0, 0.25, 0.5, \dots]$
 np. full ((2,2), 7) $\rightarrow \begin{bmatrix} [7,7] \\ [7,7] \end{bmatrix}$
 np. eye (2,2) \rightarrow identity matrix
 np. random.random ((2,2)) \rightarrow random values
 np. ones ()
 np. zeros ()

slicing
 L[0:3,1] row 0 to 3, Pak value 1
 L[:,1] Pak 1 alles
 L[:1] 1st row
 L[L<2] alle values < 2
 L[L[0,1], L[1,0]] Pak L[0,1] en L[1,0]

broadcasting
 Arange (3)+5
 $\begin{bmatrix} 0 & 1 & 2 \end{bmatrix} + 5 \begin{bmatrix} 3 \end{bmatrix} = \begin{bmatrix} 5 & 6 & 7 \end{bmatrix}$
 np. ones ((3,3)) + Arange (3)
 $\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 1 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 2 & 3 & 4 \\ 3 & 4 & 5 \end{bmatrix}$

Syntax

Index = pd. MultiIndex.from_tuples ([(d,1), ...])

reshape
 pd. melt (df) \rightarrow columns to rows
 df.pivot (table) spread rows to columns
 pd.concat (df1, df2) append rows, axis 1 for column.

df.sort_values ('col')
 high to low: ascending = False
 df.rename (columns = {y: x})
 df.reset_index ()
 df.drop (columns = ['length'])

observations:
 df [df.Length > 7]
 df.head (n)
 df.filter (regex)
 df.drop_duplicates ()
 df.nlargest df.nsmallest (n)
 df.loc [df['a'] > 5, ['a,b']]

df.dropna ()
 df.fillna ()

summarize
 df ['w'], value_counts ()
 df ['w'], nunique ()
 min ()
 max ()
 quantile (n)
 var ()
 apply (np.sort)

df ['Area'] = df.length * df.height
 pd.q (MTE)

df.plot.hist. / plot.bar () / plot.barh
 pd.crosstab ()
 df.pivot_table ()

adF [adF.x1.isin (bdf.x1)]
 df.plot.hist ()

sort_values (by = []) altijd str.lower
 str.replace
 str.contains
 str.len
 str.count

$acc = \frac{TP + TN}{total}$
 Precision = $\frac{TP}{TP + FP}$

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

df.div (sum, axis)

Precision = $\frac{TP}{TP+FP}$
 recall = $\frac{TP}{TP+FN}$
 F1 = $\frac{2 \cdot \text{Prec} \cdot \text{Recall}}{\text{Prec} + \text{Recall}}$
 accuracy = $\frac{TP+TN}{TP+FN+TN+FP}$

Magic
 % ls -lh: files in dir + size
 % rm file x? 1st line for dtype
 % time
 % timeit
 % Prun

a.shape : sort : return none
 a.ndim : sort : return copy
 a.size : argsort : index
 a.type.name
 a.astype = cast to type

np.zeros (len, dtype=int)
 np.multiply.outer (a,b)
 np.add.reduce (a,b) reduce dim
 np.add.accumulate
 np.power
 np.absolute
 np.cumsum (axis=1)

transpose

np.random.seed(a)
 np.random.randint (max, size=amt)
 np.reshape (array, (y,x)) $\begin{bmatrix} 1,2,3 \\ 4,5,6 \end{bmatrix}$
 eye (amt, dtype=int) = $\begin{bmatrix} 1,0 \\ 0,1 \end{bmatrix}$
 resize (= affects og) $\begin{bmatrix} 1,0 \\ 0,1 \end{bmatrix}$

append
 insert
 delete
 concatenate
 vstack = stack vertically
 dropna()
 replace()

Regex
 ^ start string
 \$ end
 \s white space | \S
 \d | \D
 * 0 or more
 + 1 or more
 ? 0 or 1
 {3} 3
 {3,} 3 or more

$\begin{bmatrix} abc \\ a-a \end{bmatrix}$

read_csv (file, skipinitialspace, sep)
 • values
 • index
 • loc / iloc
 • Series ([1,2,3], index=['a','b','c'])
 • head()
 • drop (index or column)
 • size
 • apply (operation)
 • merge

a = np.sort(s) sepal > 5
 df.query (sepal > 5 and width < @a)
 sort_values (['col's'], ascending=[false, true])
 argmax vs max
 sns.heatmap


```
pd.read_csv('linktotfile'), compression='gzip', header=None, names=['Jaar', etc.], sep='\t', quotechar='"', error_bad_lines = False, skipinitialspace=True)
```

```
top10 = kur['Partij'].value_counts()[0:10].index.values
```

```
partij10 = kur[kur['Partij'] == 'PVDA' | etc.]
```

```
cross = pd.crosstab(partij10, Jaar, partij10, partij)
```

```
cross.plot()
```

```
partiesleft = len(set(kur['Partij']))
```

```
kur = kur[kur['Partij'].str.replace('[^w\s]', '')] of str.lower of (' ', '')
```

```
kur = kur[pd.notnull(kur['Partij'])]
```

```
kur = kur[kur['Partij'].str.contains('vragen')]
```

```
kur = kur[kur['Partij'].value_counts().sort_index().index_col of sort_index(inplace=True)]
```

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

```
pivot['ratio'] = np.log2(pivot['Count']['M'] / pivot['Count']['F'])
set(pivot[(pivot['ratio'] > -0.1) & pivot['ratio'] < 0.1])
```

```
df = df.fillna(0)
```

```
df['ass_mean'] = df[['ass_week', etc.]].mean() of .corr()
```

```
greater = (df[['ass', etc.]] > 9).sum().sum() sumi groter dan 9
```

```
df['weighted_mean'] = 0.4 * df.ass_mean + 0.6 * deel_mean * loo
```

```
perc = ((df['ass_mean'] < 5.5) | df['deel_mean'] < 5.5).mean()
```

```
Z-df = (df - df.mean()) / df.std()
```

```
df.sort_values(['deelttoets_mean', 'names'], ascending=[False, True])
```

```
C = C[C != 1]
```

```
rij = C.max().idxmax()
```

```
Kolom = C.loc[rij].idxmax()
```

vraagID = {rij, kolom}

np. setdiff 1d

axis = 0, gaat over de columns ↓
axis = 1, gaat over de rows →

```
top = kvr. Partij.value_counts().head(10)
cross = pd.crosstab(kvr. Jaar, kvr. Partij) [top.keys()]
```

```
remove = ["^ [a-zA-Z0-9]", "\s", "allen", "biden", "group", "fractie", "fr", "lid"]
```

```
kvr. Partij = kvr. Partij.str.lower().replace(regex=remove, value="")
```

```
kvr = kvr[~kvr. Partij.str.contains("vragen", na=False)]
```

```
kvr.dropna(subset=["Partij"], inplace=True)
```

```
parties_left = kvr. Partij.unique().size
```

```
kvr["deelvragen"] = kvr. Vraag.str.count("\?")
```

Pearson correlation

```
numpy = np.corrcoef(x, y)[0,1]
pandas = x.corr(y, method="pearson")
import scipy.stats
```

```
scipy.stats.pearson(x, y)[0]
```

```
df["assignment_mean"] = df[df.columns[0:6]].mean(axis=1)
```

```
all_assignments = df[df.columns[0:6]]
```

```
greater_than_g = all_assignments[all_assignments > g].count().sum()
```

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

```
# vraag 10 = set 0
```

```
c = df[df.columns[0:6]].corr()
```

```
c = c[c!=1]
```

```
row = c.max().idxmax()
```

```
col = c[row].idxmax()
```

```
vraag 10 = (row, col)
```

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

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

RMSE = $\sqrt{\frac{1}{n} \sum (predictions - targets)^2}$
voorkeur - gemiddelen

np.where

np.argmax

normalise = $P/P.sum(axis=0)$

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

Pandas

open files: \rightarrow `pd.read_csv('file.csv', sep=';', encoding='utf-8')`
 \rightarrow with `open('file')` as `f`:
from web \rightarrow `f = requests.get(url)`

- precision accuracy = $\frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$
- recall = $\frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$

Reshaping: `pd.concat`, `pd.melt`
select col: `df.loc`
merge: `pd.merge(df1, df2, how='inner', on='x1')`
`df.groupby()`

Numpy:

inspecting: `a.shape`, `len(a)`, `a.ndim`, `e.size`, `s.dtype (int)`

Sort: `a.sort(axis=0)`

slicing: `a[row][col]`

indexing: `a[acc]` \rightarrow boolean

Fancy Indexing: `[..., ...]` \rightarrow `[..., ...]`

- null NaN \rightarrow `df.fillna(0)`
- pd nouvelle colonne: `df["new-column"] = ...`
- boolean indexing: `df["column"] > 4`
- $Z\text{-df} = (df - df.mean()) / df.std()$ \leftarrow Z-normale

Precision

$P = .1$

$$\text{Precision} = \frac{\text{Fout} \times TP}{((F \times TP) + ((1-F) \times .1))}$$
$$\frac{0.9}{0.9 + 0.1}$$

$$\text{Prec.} = \frac{TP}{TP+FP} \quad \text{Rec.} = \frac{TP}{TP+FN} \quad \text{Acc.} = \frac{TP+TN}{\text{total}}$$

np.int0 (np.array.dtype) → b.dtype(int) ? → np.int64, float32

np.: add, subtract, divide, multiply, sqrt, log, log2

np.: sum, min, max, cumsum, corrcoef, corr, std

np.: zeros((3,4)), ones((2,3,4)), arange(10,25,5), linspace(0,2,9)

np.: any, all, .ravel() == .flatten(), reshape vs. resize

np.: sort(), sort(axis=0), concatenate(a,b,axis=0)

diag.: return L[range(L.shape[0]), range(L.shape[0])]

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

np.arange(1,4)[:, np.newaxis], np.split()

%ls -lh, %rm <file>, %%.time, %time it, <TAB>, ?

df[x==2] vs. df.loc[x==2] vs. df.iloc[1:3]

pd.isnull() vs pd.notnull(), df.sort_values, df.sort_index counts

df.dropna(subset=["x"]), df.nlargest, nsmallest, df[x].value_counts

df.plot(kind="bar/barh"), df.plot.hist, df.plot.scatter, df.fillna()

pd.merge(a,b, how="inner/outer/left/right", on="x")

df.drop_duplicates(), pd.Series(data=values, index=labels)

(df[assignments] > 9).sum().sum(), data[-1], data[:,2]

df.groupby → .apply() → .unstack() → .head()

df[~df.x1.isin(df2.x1)], .index, .values

~~df.str~~ df.x1.str.replace("l".join(["\s","W"], "4"))

df[~df.x1.str.contains("string")], pivot[[F,M]].min(axis=1)

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

y = c[c!=1].loc[x].idxmax(), np.transpose == np.T, np.abs()

pd.merge(), pd.concat(), df.reset_index(), pd.melt

np.copy(a), a.copy(), a.view(), random <TAB>

divisible: L[L % n == 0], round(float, 1)

from numpy import doc → np.doc? → np.lookfor

^ starts, ends \$, (\.) alle., (C\A(?!\Species\$)) *r alles behut species
 \d digit, \W not lw, \s whitespace, {3,4} 3,4, 3+x

Numpy

Maken

`np.random.randint(0, 100, (shape))`
`np.arange(start, stop, steps)`

Informatie

`1.shape` - dimensies
`2.shape[0]` - rij
`2.shape[1]` - kolom
`2.size` - elementen
`2.dtype` - elementen
`2.astype(int)` - convert

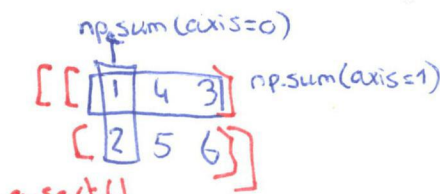
Math

`a.sum()`
`a.min()`
`a.max()`
`a.cumsum()`
`a.mean`
`a.median`
`np.std(a)`

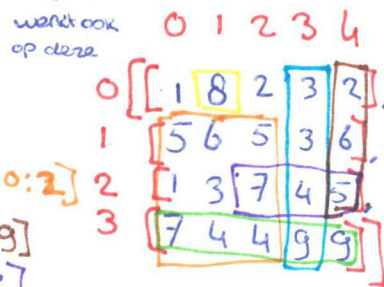
over hele array

Index Slicing

`a[0,1]` `a[1:3, 0:2]`
`a[:,3]` `a[:,-2,9]`
`a[3,:]` `a[2,2:]`



`a.sort()`
`a.sort(axis=0)` kolommen
`(axis=1)` rijen



Manipulate

`np.transpose(a)`
`np.insert(a, 1, 5)`
`np.append()`
`np.delete()`
`np.reshape()`

`a.reshape(size,)` → wordt 1D, nieuwe axis
`np.arange(1, 4)[:, np.newaxis]` - 1D in kolommen

< -10 of > 10: maak er absoluut van! `np.abs`

$Precision = \frac{TP}{TP+FP}$
 $recall = \frac{TP}{TP+FN}$ (actual TP)
 $accuracy = \frac{TP+TN}{TP+TN+FP+FN}$

`np.var(a)`
`np.std(a)`

$$RMSE = np.sqrt(np.mean((L[:,0] - L[:,0])^{**2}))$$

waarden kolommen

Jupyter handig

function? - documentatie
`help(function)` - documentatie
function?? - source code

Regular expressions

1 of `df.str.replace(re, m)`

lw words

W - not-words

[0-9] - numbers

[a-z] - lowercase

[A-Z] - upper

^hoi - begint met hoi

ls space

? alles met ?

@[-] niet.

! is magic

Pandas

Inlezen

`Pd.read_csv('naambest', sep, index_col, skipinitialspace)`

`df.sort_index()`

`df.sort_values('kolomnaam', ascending=false)`

`df.rename(columns={'old': 'new'})`

Displays

`df[df.length > 7]` - rows that meet criteria

`df.drop_duplicates()`

`df.drop(columns=['x', 'y'])`

`df.sample(frac=0.5)`

n=10

`df.nlargest(n, 'values')`

Subsets

`df[['naam', 'mean']]`

`df.filter(regex)`

`df.loc[:, [1, 2, 5]]`

Summarizing

`df[kolom].value_counts`

`df.sum()`

`df.count()`

`df.describe()`

`df.agg(['count', 'min', 'max'])`

for each in

what is total

`df.groupby('kolom').kolom.sum()`
voor totalen bij elkaar

`df.pivot_table('values', index, columns=aggfunc='mean', fill_value, margins=false)`

`df.dropna(subset=['kolom1', 'kolom2'])`

Naam moet in beide sets voorkomen.

`pd.merge(df1, df2, how='left' on='kolom1')`
intersection ↳ optional

`pd.crosstab(df['kolom1'], df['kolom2'])`

↳ altijd count

Reze

`df[kolom1], [kolom2], cross()`

`pd.set_option('display.max_columns', 10)`

`df.isna(), isin(), contains()` → all boolean

lambda x, y: x*y → shortcut for function



numpy

Arrays are created with `np.array([])`
`np.arange(10, 25, 5)` creates an array of evenly spaced values (5 is the step value)
 Inspecting your array is done by
 a. `shape` array dimensions
`len(a)` length of array
 b. `ndim` number of array dimensions
 e. `size` number of array elements
 b. `dtype` datatype of array elements.

`a[2]` subsetting, select element on 2nd index
`b[1, 2]` equivalent to `b[1][2]` or `[row][column]`

$$\text{precision} = \frac{\text{precision} \cdot 0.9}{(1-p)^{0.9}}$$

`a[0:2]` slicing for a/ in a row $(p^{0.9}) + ((1-p)^{0.9})$
`a[0:2, 1]` slicing for a/ in a column
`a[:, -1]` get the last column

$$\text{precision} = \frac{p \cdot \text{precision}}{(p^{0.9}) + ((1-p)^{0.9})}$$

`a[a < 2]` boolean indexing
`a[["fill in columns, row, etc']]` select a subset of the matrix's rows and columns. = fancy indexing

`i = np.transpose(b)` } permute array
`i.T` } dimensions
`np.intersect1d()` or `np.union1d()` or `np.setdiff1d()`
 Flatten the array `np.flatten()` or `b.ravel()`
 Reshape the array `g.reshape(3, 2)`
 concatenate arrays `np.concatenate((a, d), axis=0)`
 axis has to be specified.
`np.arange(start, stop, step)` optional `[: , np.newaxis]`

$g = a - b$ subtraction
`np.subtract(a, b)`
 $g = a + b$ addition
`np.add(a, b)`
 $g = a / b$ division
`np.divide(a, b)`
 $g = a * b$ multiplication
`np.multiply(a, b)`

$$\text{log} = \frac{\log(a) \cdot \log(b)}{\log(a) + \log(b)}$$

`np.exp(b)` Exponentiation
`np.sqrt(b)` square root
`np.log(b)` logarithm
`e. dot(f)` dot product

`a == b` returns True or False
`a < 2` returns True or False
`a != 2` returns True or False
`np.array_equal(a, b)`
`np.all()` or `np.any()`
`a.sum()` sum of values
`a.min()` `b.max()` returns max of min
`a.cumsum()` cumulative sum
`a.mean()` mean
`b.median()` median
`a.corr` correlations
`np.std()` standard deviation
`np.where()` replace values in array. or to see indices of matrices & arrays
`a.sort()` sort an array
`c.sort(axis=0)` sort an array on axis

make a dataframe, a multiIndex dataframe is created
`pd.DataFrame([[4, 7, 10], [5, 8, 11], [6, 9, 12]], index=[2, 2, 3], columns=['a', 'b', 'c'])`

$$\text{mse} = \sqrt{\frac{\sum (\text{Predicted} - \text{Actual})^2}{N}}$$

subset variables (columns)
`df[['width', 'length', 'species']]` select multiple columns & specific names.
`df.width` or `df['width']`
`df.filter(regex='regex')`

Reshaping data is done with the following

`pd.melt(df)` (columns -> to rows)
`pd.concat([df1, df2])` (df1 + df2)
`pd.concat([df1, df2], axis=1)` (df1 + df2)
`df.pivot(columns='var', values='val')` (rows -> columns)

`df.sort_values('mpg')` order rows by values in a column (optional arguments apply)
`df.rename()` rename columns.
`df.sort_index()` sort the index of a dataframe

min

`df['w'].value_counts()` Count number of rows with each unique value as variable
`len(df)` # of rows in df
`df['w'].unique()` # of unique values in a column.
`df.describe()` describes the df.

subset observations (rows)
`df[df.Length > 7]` extract rows that meet logic
`df.iloc[10:20]` select rows by position
`df.nlargest(n, 'value')` select and order top n entries
`df.nsmallest(n, 'value')` select and order bottom n entries

`df.groupby(by='col')` return a groupby object grouped by values in a column.

`df.loc[:, 'x2': 'x4']` select all columns between x2 and x4 (in)
`df.iloc[:, [1, 2, 5]]` select columns in positions 1, 2 & 5
`df.loc[df['a'] > 10, ['a', 'c']]` select rows meeting a logical condition, and only specific to the column.

additional groupby functions =
`size()` * `agg(function)`
`df.groupby(by='')["col"]`
 new column? `df['volume'] = 'value'`

`.sum()` `.mean()` `.min()` `.max()` `.count()` `.median()`
`.var()` `.std()` `.apply(function)` `df.dropna()`
`df.fillna(value)` `pd.merge()` `df.intersect1d()` `df.union1d()`
`df.diff()`

Pandas

Read File: pd.read_csv (filenames, compression, sep)
df.set_index ("column name")
Sorting: df = df.sort_values (by = [columnname], ascending = [bool])
missing values: df.fillna (value) df.dropna
mean of columns: df["column"] = df.loc[:, "von column": "tot column"].mean
 $= df["column"]$
 $= df["column"]$
Filters df = df[df["column to filter"] == 1] df[df["loc"]]
ilvening: df.loc[:, "von": "tot"]
Select more columns: df[["name", "age"]]
df.sort_index
df.sort_values ("name", ascending = False)
df.rename (columns = {"name": "newname"})
df.drop (columns = ["name", "age"])
df["name"].value_counts()
pd.merge (data, data, how, on)
how: left, right, inner, outer on: column
df.groupby (by = "name", as_index = False)
df.pivot_table (index = "name", columns = "age", aggfunc = {"column": sum})
pd.concat (data, data)

logic operators
normal ones
 column.isn
 pd.isnull (obj)
 pd.notnull (obj)
 df.any ()
 df.all ()
describe data
 sum () min ()
 max () count ()
 mean () median ()
 qseries (i) (j)
 std () var ()
 apply (function)
 abs ()

Numpy

np.shape → (row, column)
np.where (array + condition) → gives coords
np.isnan → ~~bool~~ boolean array
np.reshape (row, column)
Filter → array [array condition]
np.asarray → make np array von list
Sorting = ~~array~~ ~~sort~~ (axis) np.sort (array, axis)
np.all = ! np.all (condition) & (condition) → gives bool
np.ndim = array dimensions
array.sum ()
 • min ()
 • max (axis)
 • cumsum ()
 • mean
 • median
 • correlate
np.std (array)
Slicing array [start: end, column: end]
np.random.randint (start, end)
np.arange (start, end, stepsize)
from numpy import doc
np.random
np.lookfor
np.transpose → turn array
np.append
np.insert (array, coord)
np.delete
np.hstack / np.vstack
np.hstack (array, index = split)
precision = HP / (P)
Recall = FP / (FP + FN)
acc = TP + TN / (all)
F1 = 2 * (Precision) / (Precision + Recall)

acc: $\frac{TP+TN}{Total}$ prec: $\frac{TP}{TP+FP}$ rec: $\frac{TP}{TP+FN}$

// = floor div

np.array => np functie: np.sum, np.max, etc...

np.zeros (length/dims, dtype) or reshape (dims) np.concatenate ([arr, arr2]) np.split (grid, [indices])
 [0:5] first 5 np.count_nonzero (statement) counts booleans np.ravel (arr) ↳ hsplit, vsplit
 [5:] all after 5 np.sort (arr), np.argsort (arr) np.arange (n)
 [4:7] subarray arr.transpose () ar.flatten () / np.ravel (arr)
 [::2] step 2

arith: np. add/subtract/negative/multiply/divide/floor-divide/power/mod/abs/sqrt
 agg: np. sum/prod/mean/std/var/min/max/argmin/argmax/median/percentile/any/all (let op axis)
 ↳ nan error en nan's te handlen.
 np.cumsum (axis),

data[1:3] = implicit, data.loc[1:3] = explicit, data.iloc[1:3] explicit index
 df. add/subtract/multiply/divide/floordiv/mod/pow (doe ff? voor args)

dropna (axis, how, thresh) fillna (method (ffill, bfill))

pd.concat () (check? voor args)

df.append (df2)

pd.merge () (check? voor args)

~~pd.value_counts~~ values_count, is_in, cross-tab

.str.lower/upper/contains/startswith/replace

mean/median/mode/std

.T, .unstack, .stack

idxmax, idxmin, nunique, nlargest

groupby.apply.unstack

\d digit	+ one or more	[] select group
\w word char	{3} 3 times	^ start string
\s white space char	{2,4} 2-4 times	[^...] not
\D not \d	{3,} 3 or more	\$ end string
\W not \w	*	
\S not \s	?	
	.	

any except line break

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

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

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

dropna(axis=1)

fillna()

isnull() → boolean

isnotnull → boolean

bevat Niet

→ df[df['party'].str.contains('vliegen')]

→ len(df['party'].unique())

replace tokens + lowercase

→ df['party'].str.lower().str.replace("W", "w")

remove specific list of words

→ df['party'].str.replace(" ".join(list_of_words), "")

→ drop rows with nan

→ df.dropna(subset=['party'])

→ count = len } agrfunction

→ sum = sum

→ df[column].str.count('\?') = tel aantal deelvragen

d.crosstab(df['column1'], df['column2']).filter(...)

d.to_csv(skipinitialspace=True, sep='\t', names='list of columns')

d.merge(df1, df2, how='left/right/inner/outer', on='x1')
df1[x1].isin(df2[x1])

Numpy

np.arange(start, stop, step)

np.reshape(..., ...)

L [L % n = 0] → deelbaar door n

add_per_column = np.arange(L.shape[1]) + 1

add_per_row = np.arange(L.shape[0]) + 1

random.choice(i for i in range(5))

with open bestand as f:

for i in f:

for j in i:

print(j)