

## magic / % LSmagic dir(lp) dr(R)

%ls -lh working dir with bytes  
 %rm blabla.txt removes blabla.txt  
 variabele met '?' erachter voor info.  
 %prun, %time checkt resources gebruik.  
 %Paste geopieerde zooi met >>> <<< runnen.  
 %run myscript.py use script from some where.  
 acces shell with !  
 pwd = print working directory  
 ls = list working directory  
 cd = change directory of file (geen aanhalingstekens)  
 mv '...' = move de file hierheen.  
 mkdir = make new directory  
 tmp = temporary file place  
 cp = copy file..

regression: RMSE

Classification:  
Recall = TP / (TP + FN)

Precision = TP / (TP + FP)

accuracy how many of all predictions (where correct?)

• plot(kind='bar') df.dropna(subset=['col'])

Data science = Based on patterns in data predicting a value on new or unseen data.

paginas = BeautifulSoup(source\_file, "html.parser")

• find\_all('Page')

newpagina1 = page.get\_text()

" " 2 = re.sub("[ ]+", "", newpagina1)

df['col'] = df['col'].str.replace('15', '')

id = digits \D alles behalve letters \W alles behalve

df[~df.col.str.contains("vragen", na=False)]

alle woorden blijven

## Join combine

df1.append(df2) add rows in df1 to end of df2, (rows should be identical)

pd.concat([df1, df2], axis=1) add cols in df1 to end of df2 (rows should be identical)

df1.join(df2, on='col1, how='inner') joins the cols in df1 with cols in df2 where rows have identical values on col1

## Statistics

df.describe() summary stats for num cols

df.mean() mean of all columns

df.corr() correlation between columns of df.

df.count() number of non-null values per column

df.max() highest value of each column

df.min() lowest value of each column

df.median() median each column

df.std() standard deviation each column

## viewing / inspecting Data

df.head() first n rows

df.tail() last n rows

df.shape rows, columns

df.info() Index, datatype & memory info

df.describe() Summary stats for num columns

s.value\_counts(dropna=False) unique values met counts  
df.apply(pd.Series.value\_counts) For all columns

## Selection

df['col'] return given column as series

df[['col1, col2]] these columns as new df

s.iloc[0] Selection by position

s.iloc['index\_one'] Selection by index

df.iloc[0,:] First row

df.iloc[0,0] First element of first column

## Data cleaning

df.columns = ['a', 'b', 'c'] rename columns

pd.isnull() checks for null return boolean array

pd.isnotnull() checks for not null ""

df.dropna() Drop all rows with null values

df.dropna(axis=1) Drop all columns with null

df.fillna(x) replace nulls with x

s.astype(float) datatype to float.

s.replace(1, 'one') Replace all 1 with 'one'

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

df.rename(columns=lambda x: x+1) mass renaming

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

df.set\_index('column\_one') change index

## Filter, Sort, Groupby

df[df['col'] > 0.5] rows where col > 0.5

df[(df['col'] > 0.5) & (df['col'] < 0.7)] "" ""

df.sort\_values(col1) sort values ascending

df.sort\_values([col1, col2], ascending=[True, False])

df.groupby(col) returns a groupby object for values from one column

df.groupby([col1, col2]) Object for multiple columns

df.groupby(col1)[col2] returns the mean of values in col2 grouped by col1 values. mean can be replaced by any statistic function

df.pivot\_table(index=col1, values=[col2, col3], aggfunc=mean) creates pivot

table that groups col1 and calculates mean of col2 and col3

df.groupby(col1).agg(np.mean) averages across all columns for every unique col1 group.  
df.apply(np.mean) mean for each column  
df.apply(np.max, axis=1) max for each row.







np.arange(1, 10, 2).reshape(3, 3)  
begin ↑ end ↑ step  
 np.concatenate((X, Y)) → join - pd.concat()  
 np.vstack((X, Y)) → vertical  
 np.linspace(0, 6, 4) → 0, 2, 4, 6  
start ↑ end ↑ elements  
 np.add.reduce(x) → sum of all elements  
 np.count\_nonzero(x < 6)  
 np.sum(x < 6)  
 np.add.at(x, i, 1)  
array ↑ index ↑ plus

pd.read\_excel(" ", index\_col = " ")  
 len(set( .index )) == len( .index )  
 df.value\_counts()  
 df.sort\_index()  
 str.lower().str.strip().str.replace(" ", "").str.replace(" ", " ").sort\_values()  
 ~ str.contains('school')  
 m = c.to\_numpy()  
 c.to\_numpy()[m.index[m] - m].tail()  
 df.groupby(c)[i].count().sort\_values().tail()  
 .strip()

np.partition(x, 3) → 3 smallest numbers  
 data = np.zeros(4, dtype = {'names': ('name', 'age'),  
 'formats': ('U10', 'i4')})  
 np.sqrt(((x1 - x1.mean())\*\*2).mean()) = np.std(x1)  
 li = soup.findAll('li', class = re.compile('^interlanguage-line'))  
 {l.a.attrs['lang']: l.a.attrs['href'] for l in li}

html.doc = requests.get(url)  
 soup = BeautifulSoup(html.doc.content)  
 comments = soup.findAll('li', class = 'comment')  
 [soup.title.text for c in comments if c.p]  
 re.findall(b'brewery', line.lower())  
 line.split(b',')[1]  
 np.random.rand(3, 2)

**Index**  
 intersection &  
 union |  
 difference ^

i loc → index  
 loc  
 ix → both

df.query("sepal.length > 5 & ... @a")  
 df.groupby('species')['sepal.length'].max()  
 df.sort\_values([" ", " "], ascending = [False, True])  
 df.drop('species')

**pd.MultiIndex**  
 pd.merge(df1, df2) → how = 'outer'  
 axis = 0 → columns  
 axis = 1 → rows  
 df.set\_index(" ")  
 df.join()  
 % magic

? -> help(C) | % quickref A = and  
help(len), re? | 1 = or

NP.array([...]), dtype='float'  
NP.array([range(i, i+3) for i in [2,4,6]])  
np.zeros | np.full((3,5), '2', dtype=str)  
np.ones  
np.arange(0,20,2) | .ndim .copy()  
np.linspace(0,1,5) | .shape np.sort()  
3D Array: | .size  
| .dtype

np.random.randint(10, size=(3,4,5))

2D Array:  
2d[x,y]  
2d[:, :2] every other column  
2d[:, :-1, :-1] Reversed array

1D Array:  
1d[N], 1d[:N], 1d[N:]  
1d[:, :2] every other element  
1d[:, :-1] all reversed

Van 1D naar 2D:  
np.arange(1,10).reshape((3,3))

Samenvoegen:  
np.concatenate([x,y,z], axis=...)  
np.vstack([x,grid])  
np.hstack([grid,y])

Splitten:  
x1,x2,x3 = np.split(x, [3,5])  
upper, lower = np.vsplit(grid, [2])  
left, right = np.hsplit(grid, [2])

Sum/mult. of all elements:  
np.add.reduce(x)  
np.multiply.reduce(x)

Mult. table:  
np.multiply.outer(x,x)

Hoeveel minder dan 6?  
np.count\_nonzero(x<6)

List(data.items())  
data.values  
data.index  
data.columns  
data['columns/rows']  
2D array -> Df

Pd.DataFrame(np.random.rand(3,2),  
columns=['foo','bar'],  
index=['a','b','c'])

Indexing & Slicing

data.loc[1]  
data.iloc[1]

Ufunc  
column.index

NaN -> 0  
f: fill\_value = 0  
data.fillna(0)  
data.isnull()  
data[data.notnull()]

Drop all if column value = 0  
df.dropna(axis='columns')  
(axis='row', how='all')

Mult. indexed -> DF

data.unstack()

Samenvoegen:

Pd.concat([ser1, ser2])

Merge into one DF:

df3 = Pd.merge(df1, df2)

Pd.merge(df1, df2, left\_on='column1', right\_on='col2').dropna('col2', axis=1)  
Pd.merge(df1, df2, how='inner')

'outer', 'left', 'right'

Van 2D naar 1D Array:  
L.flatten()  
(n,1) -> 1D met n aantal waarde in L:  
np.reshape(L, (len(L.flatten()), 1))  
Laatste kolom 2D Array:  
Return (L[:, -1:])

Array tafel van (n,x):  
np.arange(0, n\*x, x)

Alle elem. uit array L deelbaar door n:  
output =  
[L[x] in range(len(L)) if L[x] % n == 0]  
Return np.asarray(output)

Crosstab:  
Pd.crosstab(df.col1, df.col2)  
Bepaalde value in column:  
df[df['col1'] == 'value']['col2']  
Bepaalde value tellen:  
df['col1'].str.count('?')

(index='\_', columns='col1','col2', values='col',  
aggfunc=np.sum, margins=True, margins\_name='all')

Sorteren op groot -> klein:  
df.nlargest(len(df), 'column')

Reading from the web:

import requests  
f = requests.get(url)  
lines = f.text.split('\n') -> f.close() -> len(lines)  
lines[:2]  
[line.split(';') for line in lines]

Sqrt = int(np.sqrt(len(L)))  
np.reshape(L, (sqrt, sqrt))  
Mean per column:  
np.mean(L, axis=0)  
omdraaien 2D array:  
L.transpose()

df['col'].value\_counts()  
df.str.replace(['WLS'], '')  
df[df.str.contains("-")]

Pivot table:  
df.pivot\_table()



```

df.index = df.index.replace(L['a-z0-9'], , regex=True)
df = df.drop(df[df.index.str.match('.*word.*', na=False)].index)
df.index.value_counts()[df.index.value_counts() > 1]
df.index = df.index.str.lower()
df['RMSE'] = np.sqrt((df.A - df.B)**2)
M = df.mean(); (df[m.index] - M)
df.A.sort_values().plot(kind='barh') (== .plot.barh())
df[df.A > 100].B.value_counts().sort_values.plot(kind='barh')
np.random.randint('int')
    .choice('list')

```

V x z z 4 8 6 7 8 9 10  
 T 2  
 X Z 3 4 5 6 7 8 9 10

```

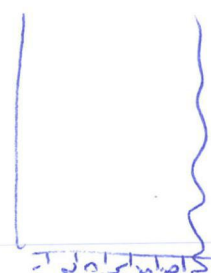
for i, x in df.iterrows()
pd.read_csv(file, sep='t', dtype=str, header=None, names=['...', ...], index_col=...)

```

```

np.arange()
np.array.reshape(x,y)
np.array.shape
np.mean()
np.median()

```



```

df.A.mode()
df.A.corr(df.B, method='pearson')
pd.pivot_table(df, index='...', columns='...', values='...', aggfunc=np.sum, fill-
value=0, margins=True).sort_values(by='...', ascending=False,
kind='mergesort')

```

```

pivot.loc[(pivot.A < .1) & (pivot.A > -.1)].index.values

```

```

list.extend(list2)

```

```

plt.plot(x,y)
plt.legend(['...', ...])
plt.show()

```

```

page = re.sub('<[1<]+?>', '', page)

```

Import re  
Import collections - (Counter, most, common)  
Import math - math, sqrt()

Basics  
Random, randint()  
Random, choice()  
L [-1:] last element  
L [-1:] last element

$$g \log(a) = b - g^b = a$$

2 log 16 = 4 - 2 = 16 Pandas

KVR = pd.read\_csv('file', sep = '\t', skipinitialspace=True)

KVR.columns = ["name", "number"]  
KVR\_obj = KVR.select\_dtypes(include="number")  
KVR[KVR\_obj.columns] = KVR\_obj.apply(lambda x: x \* 2)

X.Strip() } .NumDuel()  
KVR["party"] } .NumDuel()

make error = 1/19 x 100  
total error = sum(error) / sum(["all"]) x 100  
Sort-Values (ascending = False)  
Subset = Names[Names.names == 'john']

- Numpy
- info()
- shape()
- size()
- in space()
- fill()
- eyes()
- abstract
- transpose()
- mean()
- reshape
- arrange
- ones
- grid, shape - [-1, 0]
- mode(c)
- groupby
- rel-error = boogstic-error
- most, error = mistake x row

Accuracy = tp + tn / tp + tn + fp + fn

recall = tp / tp + fn  
precision = tp / tp + fp  
f1 = 2 \* precision \* recall / (precision + recall)

re.sub(r'', ' ', text)  
re.findall()

midlen = len(L) - 1  
L[0] any char  
L[-1] last char  
L[1:-1] middle  
L[1:-2] middle  
L[1:-2:-1] reverse middle

PLI

pit, high  
pit, barh  
John, pivot-table(values)  
'births', index: 'year', column: 'sex', aggfunc = sum, p



np.random.randint(low, high, size) **Numpy functions**  
 np.random.choice(a, size) np.log(x)

win-percentage (Prizedoors, guesses):  
 $100 * (np.array(prizedoors) == np.array(guesses).mean())$   
 → == op numpy array kijkt per positie of ze gelijk zijn en geeft True/False terug

np.absolute()  
 np.arange(start, stop, step) → make 1D array  
 np.reshape: [array-naam].reshape(shape) → shortcut for np.reshape()  
 → [1, 2, 3] = a, a.reshape(a.size, 1), a = [[1], [2], [3]]

np.shape(array) → ook array.shape  
 np.size(array) → ook array.size

np.extract(condition, array) → return elements of array that satisf. condition in array + 0.  
 np.mean(array, axis)  
 np.transpose(array)

**Plotje panaja** df = pd.DataFrame(data)  
 Plotje = df.plot(kind=..., title='...', ...);  
 Plotje.set\_xlabel('...');  
 Plotje.set\_ylabel('...');  
 → if you want to plot your own date, e.g. Switched & not switched, consider using a dict comprehension to make the data for the dataframe

n	stay	switch
3	perc.	perc.
4	perc.	perc.

data = {'switch': {percentage}, 'stay': {n: percentage}}  
 → stay switch → lijnen, 3,4 x as, perc

hist = kolom.hist(bins=...)  
 hor.bar = kolom.value\_counts().plot(kind='bar')

**Pandas functi**  
 pd.read\_csv()  
 pd.crosstab(col1, col2)  
 dataframe['column'].value\_counts() → counts how many time elements occur in the colou  
 dataframe[column].value\_counts() → many time elements occur in the colou  
 df[df[column] == '...'].column.value\_counts() → waar rij in deze column gelijk is aan string have  
 df[column1, column2].corr() df.sort\_values()  
 df[column].mean() df.drop() → never subset maken  
 df[column].mode() df[['column1', 'column2']]  
 df[column].median() df[column].min/max/sum  
 df.pivot\_table() df.copy()

**Numpy slicing** fixed type arrays

X1 = [1, 2, 3, 4, 5, 6]  
 X1[0] ⇒ 1  
 X1[-1] ⇒ 6  
 X1[:3] ⇒ [1, 2, 3]  
 X1[3:] ⇒ [4, 5, 6]  
 X1[::2] ⇒ [1, 3, 5]  
 X1[1::2] ⇒ [2, 4, 6]  
 ↳ every other element starting at index 1

X1[::-1] ⇒ [6, 5, 4, 3, 2, 1]  
 X1[3::-2] ⇒ [4, 2]  
 ↳ every other element reversed, starting at index 3

Y1 = [98, 99]  
 np.concatenate([X1, Y1])  
 ⇒ [1, 2, 3, 4, 5, 6, 98, 99]  
 ↳ also possible with more than 2, or with grids (2D)  
 default: add 2nd grid as rows, with axis=1 add as columns

X2 = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]  
 X2[2, 0] ⇒ 7 (X2[row, column])  
 X2[0, 0] = 7 (assigning value 7 to row 0 column 0)  
 X2[:2, :3] ⇒ [[1, 2, 3], [4, 5, 6]] two rows, three columns  
 X2[:, :2] ⇒ [[1, 3], [4, 6], [7, 9]] all rows, every column  
 X2[:, :-1] ⇒ [[9, 8, 7], [6, 5, 4], [3, 2, 1]] first column of x2  
 X2[0, :] ⇒ [1, 2, 3] first row of x2  
 ↳ zelfde: X2[0]  
 X2[:2, :2] ⇒ [[1, 2], [4, 5]] two rows, two columns

**Pandas subsets**  
 columns: df[['col1', 'col2']]  
 single column: df['col1'] as series  
 row: df.loc['index-of-row']  
 rows: df.loc['index1', 'index2']  
 range of rows: df.loc[idx1:idx4] (1/m)  
 range of rows w/ steps: df.loc[idx1:idx4:2]  
 rows & columns: df.loc[rows][cols]  
 row & column: df.loc[row, col]  
 all rows & column: df.loc[:, col]  
 index selecting: df.iloc [numpy slicing]

**cleaning** from b34, import Beach  
 soup = Beautiful Soup('string', HTML parser)  
 soup.findAll('html-tag')  
 re.sub('regex', 'subst', 'stri')  
 list-of-patterns = ['is', 'w', 'word', ...]  
 compiled re = re.compile("".join(list-of-pattern))  
 df[column] = df[column].lower().str.replace(compiled-re, "")  
 df = df.dropna(subset=...)

**file** with b22, open file, encoding='utf8', mode='rt') as file:  
 for line in tqdm\_notebook(file):  
 ... progress-bar

**Pd series** value\_counts(), min(), max(), sum(), mean(), mode(), median(), abs()

**Numpy Broadcasting**  
 Precision = TP / (TP + FP)  
 recall = T / (TP + FN)

0	0	0
10	10	10
20	20	20
30	30	30

0	1	2
0	1	2
0	1	2
0	1	2

=

0	1	2
10	11	12
20	21	22
30	31	32

**Boolean indexing**  
 arr = [1, 2, 3, 4]  
 arr[True, False, True, False] = [1, 3]  
 arr[arr < 3] = [1, 2]  
 arr[(arr > 1) & (arr < 4)] = [2, 3]

**Splitting of Arrays**  
 grid = [[0, 1, 2, 3], [4, 5, 6, 7], [8, 9, 10, 11], [12, 13, 14, 15]]  
 X = [1, 2, 3, 99, 99, 3, 2, 1]  
 X1, X2, X3 = np.split(X, 3)  
 ⇒ [1, 2, 3] [99, 99] [3, 2, 1]

upper = np.vsplit(grid, [2]) ⇒ upper: [[0 1 2 3], [4 5 6 7]]  
 lower: [[8 9 10 11], [12 13 14 15]]

left, right = np.hsplit(grid, [2]) ⇒ left: [[0 1], [4 5], [8 9], [12 13]]  
 right: [[2 3], [6 7], [10 11], [14 15]]

Grid = np.concatenate([grid, grid]) np.concatenate([grid, grid], axis=1) np.vstack([x, grid]) np.hstack([grid, x])  
 [1, 2, 3], [4, 5, 6] [1, 2, 3, 1, 2, 3], [4, 5, 6, 4, 5, 6] [[7, 8, 9], [1, 2, 3], [4, 5, 6]] [[-1, 2, 3, 99], [4, 5, 6, 99]]

x = [7, 8, 9] [1, 2, 3] [4, 5, 6]  
 y = [99, 99, 99] [4, 5, 6]

% magic % time it % CP  
 % Lim % IC

**Concatenation of Arrays**



Pandas

pd.read\_csv() / xml

Select: df[df.kolom == 'value'], kolom[0]  
df.loc[df.kolom == 'value']

toevoegen: df['naam'] = df.b / df.a

oortoen: df.sort\_values('kolom', ascending=False)

value\_counts() = waarden tellen van kolom. Hoe vaak komt iets voor?

colkan mask: X[X[:,3]==0]

shape.size . describe() . ndim

samenvoegen: pd.merge(x,y)

many-1 many-many many-1 = 1/2 vd kolommen heeft duplicaten

pd.concat (x,y) of x.append(y)

pd.fillna() df.dropna

np.nansum (aeray) df[df.nnull()]

any indexing 3 rijen, 4 kolommen

row = array (1,2,3,4) kolom  
col = array (1,2,3,4) rij

pd.diagonal (3,4,7) x [0,1,2,3,4,5,6,7]

latten: .plot(kind='barh')

anpassen: df.kolom = df.kolom.str.lower()

resetcol: pd.resetcol(df.kolom, df.kolom) (lijst)

ivot: df = pd.ivot\_csv(title)

h.pivot-table (df, index='name', columns='sex', values='count', aggfunc=np.sum, margins=True)

pivotstijle os Dataframe

index.values

pd.groupby('species')['length'].max  
jeet grootste waarde van kolom  
jeet lengte met specie eluore  
setosa 5.8  
virginica 6.1  
for each specie?

numpy

create: np.arange(start, stop, step)

np.random.randint (5, size=10)

slicing: x[x[:,2]] elk een geel

x[:, : -1] alle elementen, omgekeerd

x[:, 2, :2] geeft 2x2 array

reshape: x.reshape (3,4)

samenvoegen: np.concatenate (x,y)

axis=0 verticaal toevoegen (stack)

axis=1 horizontaal toevoegen (stack)

peken: np.min(axis=0)

x.argmax / index max value

np.count-nonce (x,y)

hoeveel values zijn kleiner dan b?

np.any (x>8) zijn er values getreden?

np.all (x>8) np.argmax (l)

sorteerd np.sort (x, axis=0)

toevoegen np.insert (l, index, object)

literature

fluoray: positive 95 is sick 495 590  
negative 5 9405 9410

recall = TP / (TP + FP) = 95 / (95 + 5)

precision = TP / (TP + FP) = 95 / (95 + 5)

Regression = RMSE, rest = classification

Variance = E(x -  $\bar{x}$ )<sup>2</sup> / N = 2.5, 6.9 = 6.93

STD =  $\sqrt{\text{variance}}$  df, std

a = array (1,2,3,4,5)

b = array (1,2,3)

a(3,1) b(3,1) →

a(3,1) b(3,1)

a+b = array (1,2,3,4,5, 1,2,3,4,5, 1,2,3,4,5)

functions

dataframe van lijst met dicts

df = pd.DataFrame (geom, dict (list))

df['id'] = df['id'].str.get(0)

meest pop wooneplekken

df = df.woonplats.value\_counts() [2:20]

df.plot (kind='barh')

Praktien met vreuren

df.reset\_index (df.reset\_index() = 'vreuren':unique)

tekst tokenize, most occurring tokens

C = Counter (text.split())

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

CS.head (10)

find tokens in text which occur

in total 24 characters

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

check voor dubbel in text

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

keuze namen komen dubbel oore?

Plot inkomens per school

Cito.inkomen.sort\_values().plot()

in welke gemeente is inkomen > 50000

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

sort\_values().plot (kind='barh')

maak lijst met tuples van items

list (data.items())

veranderde kolomen pakken

df[['kolom1', 'kolom2']]

de ie rij pakken

df.iloc [1]

kolom als index

df.set\_index('kolom')

Regex

p = re.compile ('wtf')

p.findall ('string')

re.split ('[w]+', 'string')

re.sub ('[w]+', '', 'string')

w = (a-z, A-Z, 0-9)

w = not \*

\d = (0-9)

\s = (space, tab, newline)

BeautifulSoup import BeautifulSoup

from bs4 import BeautifulSoup (html-doc, content)

BeautifulSoup (html-doc, content)

soup.findall (l)

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]

laatste kolom L[-1:]



NumPy

- 1D → row or col
- min/max axis 0 for cols
- argmin (index)
- np.newaxis, : ] or r ]
- np.concatenate([grid, grid], axis=0/1)
- broadcast casting
- value - counts().plot(kind='bar(h)')
- M.shape = (2,3) → (2,3)
- a.shape = (3,) → (1,3) → (2,3)
- values less than 6?
- np.count\_nonzero(x < 6)
- also in agg. functions
- Boolean operators & array mask
- fancy ind → shape of result reflects shape of index arrays
- np.multiply.outer(x, x)
- reduce array for binary - o. any/all
- np.add.reduce(x)
- store intermediate with accumulate - np.add, subtract, ... (a, b) or ex. cumsum
- np.linspace
- np.eye (3) identity matrix 3x3
- itemsize, nbytes
- slicing
- x[begin:end(:step), begin:end(:step)]
- x[0] = x[0, :]
- np.split(x, [pos1, pos2, n])
- n+1 variables
- vhsplit(x, [i])
- dsplit → 3d axis
- multiplication table
- np.multiply.outer(x, x)
- reduce array for binary - o. any/all
- np.add.reduce(x)
- store intermediate with accumulate - np.add, subtract, ... (a, b) or ex. cumsum
- np.linspace
- np.eye (3) identity matrix 3x3
- itemsize, nbytes
- slicing
- x[begin:end(:step), begin:end(:step)]
- x[0] = x[0, :]

- np.concatenate([grid, grid], axis=0/1)
- broadcast casting
- M.shape = (2,3) → (2,3)
- a.shape = (3,) → (1,3) → (2,3)
- values less than 6?
- np.count\_nonzero(x < 6)
- also in agg. functions
- Boolean operators & array mask
- fancy ind → shape of result reflects shape of index arrays
- np.multiply.outer(x, x)
- reduce array for binary - o. any/all
- np.add.reduce(x)
- store intermediate with accumulate - np.add, subtract, ... (a, b) or ex. cumsum
- np.linspace
- np.eye (3) identity matrix 3x3
- itemsize, nbytes
- slicing
- x[begin:end(:step), begin:end(:step)]
- x[0] = x[0, :]

Pandas

- fx in df.iterrows
- df.iloc[[0],[0]] → single value by row & col
- df.iat[0][0]
- df.loc[[0], ['country']] → row & col label
- df.ix[2] → single row of subset of cols
- df.ix[:, 'capital'] → single col of subset of cols
- df.ix[1, 'capital'] → select rows & cols
- df.sort\_values(by=...)
- df.rank() → obsing ranks
- df.idxmax() → index
- get name where age < 30
- data[data['age'] < 30]['name']
- pd.pivot-table
- pd.pivot-table

Data Science → based on patterns in data. Predictions are made on unseen data

Steps: interacting with outside world → preparation → transformation → modeling & computing → presentation

Study of the generalizable extraction of knowledge

make implicit knowledge explicit & in format that is computationally fast for computer and still corresponds to human's view of that information

Theorie

Precision → how many positive actually positive TP/TP+FP

Recall → how many positive samples are captured by the pos. prediction TP/TP+FP

np.sort(a, axis=-1, kind='order) sorted copy

- np.add, subtract, ... (a, b)



LL start: even: stop ] → even: LC:: 2]
   
Dot product:  $\sum x_i y_i \rightarrow s.dot(S)$ 
  
DS: Predicting value on unseen instances based on patterns in data → regn./class
   
evaluatie: class: Accuracy, Precision, recall, FI
   
regn: RMSE

Precision: how many pred. were corr?  $\frac{TP}{TP+FP}$ 
  
Recall: How many C's predicted as C?  $\frac{TP}{TP+FN}$ 
  
Accuracy: how many of all pred. were corr?

	Has D	No D	
Posit	TP	FP	= alles P voorspeld = alles N voorspeld
Negit	FN	TN	
	= alle zieken		= alle gezonden

Impliciete inform expliciet maken
   
- lees baan voor computer in een formaat dat computationeel snel is, en aansluit op het beeld vd mens mbt die informatie
   
2015-02-09: alfabetisch ordenen.
   
spread sheets: taken:
   
delete, insert, update: rijen/kols weghalen/toev..
   
= nieuw element / nieuw eigenschap

Andere view: sorteren, slicing
   
Aggregatie: nieuw veld = functie
   
Combineren van sheets
   
excell/open off/Google: lastig mergen
   
Google fusion = goed in grote sheets
   
excell: schooft niet naar grote sets
   
- steeds dezelfde handmatige GML commands met veel zware berekeningen
   
- Functionaliteit beperkt: opschonen moeilyk
   
- Geen integratie met andere tools/libery's.

Stappen DS: 1. Interacting outside world: reading/wrt. data formats
   
2. Prep.: cleaning, combining, normaliz. slicing, etc.
   
3. transformation: applying maths/stats
   
4. Model./computation: connect data to ML alg,
   
5. Presentation: create visualizations
   
MM: Verzamel, clean, herstructureer, analyse, rapporteer.

Wagen - PvdD. ministerie. value\_counts().idxmax()
   
banchan = 4, banchan[banchan >= 5].plot.kmh()
   
Pivot (values = waarden in cell, index, cols, aggf: operatie →
   
→ voor values, margins)
   
Pivot(CF, M).min(axis=1).sum()
   
No. of HL = Kvr. vraag [Kvr. vraag. str.count('!') <= 50].str.count
   
add\_per\_row = add\_per\_column(L.T).T

NP: (Nj, kolom): 2D, (kolommen/waarden): 1
   
Axis 0: within each column
   
1: within each row
   
np.concatenate voor joinen 2 am's in NP
   
↳ vstack/hstack
   
np.abs = absolute
   
• linspace = start, end, aantal elementen

Broadc. 1. Shape with fewer dimen + 1 left side
   
2. am. with 1 stretch to match
   
PD: • iloc = rijen pakken op indexen (getalle)
   
• loc = rijen pakken op label

~~voorbeeld~~ voorbeeld TT:
   
Boolean mask:  $x1(x1 \% 3 = 0)$ 
  
Fancy index:  $x2[[0, 1, 2]]$ 
  
• Plot(kind = 'barh')
   
mis.drop('species', axis=1)
   
mis.sort\_values(['col', 'col'], ascending: [T, F])
   
mis.groupby('species')['sepal.length'].max()
   
Counter(List) → pd.Series(Counter(list))
   
CS[CS.index.str.len() = 24].index

bd.groupby(['party', 'land']).count().sort\_values...
   
Magic: %: line magic %%: cell magic
   
Pasting code blocks: % paste cpaste
   
Run external code % run
   
timing: %timeit, %time, %run
   
magic functions: % magic → list: %lsmas
   
Print (-): Print Previous output
   
Out[2] ← 2 ... % history
   
working dir: %pwd, content+size %ls
   
change dir: %cd, make dir %mkdir
   
move file: %mv file loc, remove: %rm
   
all data/info from van: van?
   
memory usage: %memit %mprun

PD. MultiIndex.from\_tuples([(c), (c), (c)])
   
or. df(data, index=[(c1), (c2)], cols=...)
   
NP.concatenate((x, y, z)) ... axis
   
PD.concat([df1, df2], axis) inner
   
Pandas preserves indices
   
PD.merge: intelligent → zoekt naar zelfde columns
   
PD.merge(CON='keyword')
   
on indices df1.join(df2).
   
• is null



A: species, value, counts(). plot(kind='barh')  
 dataframe x numerical data, series y species inf.  
 X, y = A.drop('species', axis=1), A.species  
 Z-normalize data in each x-column by subtracting mean of the column from each cell and divide by the standard deviation of each column.  
 Also show each column of the normalized data ((x-x.mean())/x.std()). head(), (x-x.mean())/x.std().std()  
 sort by Name and ALL in reverse  
 pivot: sort-values (by='ALL', Name), ascending=[False, True]  
 maximal length each species  
 A.groupby('species')['length'].max()  
 length larger than 5 and width smaller than 5  
 R = np.sqrt(5)  
 A.query('length > 5 and width < @R')  
 tokenize with split and count each token.  
 create serie with this dict. and sort it  
 C = Counter(text.split())  
 cs = A.dseries(C).sort\_values(ascending=False)  
 tokens that occur in total 24 characters  
 eg: if 'ida' occurs 5 times it occ. 15 charact.  
 what % of all unique words occurs just once?  
 and what % of all words in text?  
 cs[cs==1].sum()/cs.sum(), (cs==1).mean(),  
 cs[cs==1].sum()/cs.sum(), (cs==1).sum()/cs.sum()  
 Hist for each number: i of occurrences how many  
 tokens in text occur i many times (so if there are  
 8 words which occur each 4 times, we have an  
 8-high bar at x-value 4)  
 pd.Series(Counter(cs.values)).plot(kind='bar');  
 cs.value\_counts().sort\_index().plot(kind='bar')  
 meeste leden met aantal leden van dat land  
 pd.groupby('country').count().sort\_values().tail(1)  
 grootste aantal leden met 1 deelstaat  
 pd.groupby('country').value\_counts().sort\_values().tail(1)

x1[-3:], x1[3:]  
 last 3 columns of first two rows  
 x2[:, -3:], x2[:, 2:]  
 reshape into 1-dimens. array  
 x2.reshape((172))  
 change each element to its squared difference from the mean  
 (x1 - x1.mean())\*\*2  
 np.sqrt((x1 - x1.mean())\*\*2).mean() = x1.std()  
 boolean mask on x1 and reduce x1 to int  
 divisible by 3  
 x1[x1%3==0]  
 fancy indexing to get the diagonal of x2  
 x2[0,1,2], [0,1,2]  
 use comprehension to compute sum of squares  
 sum(x\*\*2 for x in big\_array)  
 universal functions  
 sum(big\_array\*\*2)  
 files in directory, how large  
 ls -lh  
 remove  
 rm no-good.txt  
 what resources does the algorithm use?  
 %time, %timeit, %prun  
 find out datatype using magic command  
 x?  
 id array deelbaar door 0  
 def... bla  
 np.arange(0)\*0  
 make float van de array, float(32)  
 np.array([1,2,3], dtype=float32)  
 np.array([range(i, i+8) for i in [2,4,6]])  
 np.linspace(0,1,5)  
 np.random.randint(0,10,(3,3)) #array of 3x3

alle elements reversed  
 x[::-1] # 1d array  
 every other from  
 index 5  
 x[5::-2]  
 two rows, 3 columns  
 x2[:, :2], x2[:, 2:]  
 all rows, every other c  
 x2[:, ::2], x2[:, 1::2]  
 completely reversed  
 x2[::-1, :]  
 first column of x2  
 x2[:, 0]  
 3x3 grid  
 .reshape(3,3)  
 x[:, np.newaxis]  
 np.multiply(x, 10, out=y)  
 np.add.reduce(x)  
 sum of all elements  
 min value of column  
 M.min(axis=0)  
 in each row how many?  
 np.sum(x < 6, axis=1)  
 sort each row  
 np.sort(x, axis=1)  
 np.sum((inch < 0.5) | (inch > 1))  
 all non zero in python  
 are True, and is over  
 an entire object.  
 [a tm n]  
 ind = np.array([0, 2], [4, 5])  
 x[ind] #  
 [[d, e], [e, f]]  
 fancy indexing

the shape of the one with fewer dimensions is padded with ones on its leading (left) side  
 M.shape = (9, 3)  
 N.shape = (3, 3)  
 if the shape of the arrays does not match in any dimension, the array with shape (2, 3) is padded if in any dimension the sizes disagree

pd.merge(a, b, how='inner')  
 pd.merge(a, b, how='left', left\_on='read', right\_on='help').drop('help', axis=1)  
 Niet-deelbaar = np.array(candidates) % d, dtype=bool  
 ↳ Erasthostenes: lijst 2 tot max → kleinste getal → streep alle veelvouden door (incl getal zelf)  
 ↳ toevoegen aan priemgetallen (600, 136)  
 Grootste gemene deler (Euclides)  
 600 = 4 \* 136 + 56  
 136 = 2 \* 56 + 24  
 56 = 2 \* 24 + 8  
 while A != 0 and B != 0:  
 if A > B:  
 A = A % B  
 else:  
 B = B % A  
 return max(A, B)  
 pivot['ratio'] = np.log2(pivot['M'] / pivot['F'])  
 subset = pivot[pivot['ratio'] > -0.1 & (pivot['ratio'] < 0.1)]  
 ambi-names = set(subset)  
 act\_error = 100 \* (pivot['F', 'M']).min(axis=1).sum() / pivot['ALL', 'sum']

krv = pd.read\_csv(linktofile(), compression='gzip', sep='\\t', index\_col=0) # eerste colom = index  
 header=None, names=['joepit', 'ja']  
 skipinitialspace=True  
 dot product:  
 sum(a[i]\*b[i] for i in range(len(a)))  
 sum(i\*i for i in baby)  
 df[df.A == 'b'][[ 'A', 'B']] # kolom A en B  
 projection:  
 kvr[['jaar', 'party']]  
 s.argmax() # naam hoogste cijfer  
 s.argmax() # index of first occurrence of max  
 pivot['col from index'] = pivot.index  
 pivot.drop('col from index', axis=1, inplace=True)  
 ↳ axis=1 dan drop je de kolom  
 lemma = onhandig om definitives in spreadsheet  
 beter in XML of JSON

ersnel je programma door:  
 - vervangen for-loops door vector berekeningen  
 - np.arrays ipv lists en sets  
 - sleep geen onnodige ballast mee.  
 breed-first search: afstand tussen node n in netwerk G als in de lagen bipartiet = netwerk met twee sets.  
 9 feb != ambigu  
 impliciete info: wilders zijn uitspraken

data science = classification (spam), regressie (predicting, e.g. grade) (RMSE) → niet als normale statistiek, je legt niks uit.  
 interacting with outside world, preparation, transformation, modeling and computing, presentation  
 ↳ verschillende file formats  
 ↳ schoonmaken, ↳ berekeningen om nieuwe sets te krijgen  
 ↳ data mining  
 ↳ statistieke modellen  
 ↳ data visualisatie

precision = hoeveel in voorspelde klas zijn correct (hoeveel voorspeld als spam, echt spam).  
 Recall = hoeveel instances van klasse (voorspeld als) (hoeveel van de spammail zie ik dat spam was)  
 Accuracy = hoeveel voorspellingen zijn correct?  
 test op ziekte = 95% correct, 1 op 100 heeft het:

precision = TP / (TP + FP) = 95 / (95 + 495) = .16  
 ↳ om de hoge pakkans te krijgen betaal je de prijs van veel valse beschuldigingen, accuracy is niet de juiste maat als je zegt dat iedereen zwanger is, je precision en recall zijn dan 0  
 Recall = TP / (TP + FN)

	ziek (T)	niet (F)	
posttest(T)	95 TP	495 FP	590
neg-test(F)	5 FN	995 TN	990
	100	9900	10.000

Logaritme:  
 product = log<sub>b</sub> mn = log<sub>b</sub> m + log<sub>b</sub> n  
 ratio = log<sub>b</sub> (m/n) = log<sub>b</sub> m - log<sub>b</sub> n  
 powers = log<sub>b</sub> n<sup>p</sup> = p log<sub>b</sub> n  
 Roots = log<sub>b</sub> √<sup>q</sup>n = 1/q log<sub>b</sub> n  
 change of bases = log<sub>b</sub> n = log<sub>a</sub> n / log<sub>a</sub> b

if y = b<sup>x</sup> then x = log<sub>b</sub> y b ≠ 1  
 y = 2<sup>x</sup> ↔ x = log<sub>2</sub> y  
 2<sup>x</sup> = 16 ↔ log<sub>2</sub> 16 = x = 4  
 3<sup>x</sup> = 81 ↔ log<sub>3</sub> 81

Re.findall('b'brewery', line.lower())  
 line.split('b',')[1]  
 lstrip()  
 index\_col = " "