

randint (range order, range boven) for x in range(lenge)

random.choice (list)

plt.plot (list1, list2)

re.sub = substitute

re.findall = match

text = re.sub(r'\w = word of a-g, . = 0 of meer, ? = 0 of meer, + = 1 of meer, \d = digit, \ = escape, () = capture inhoud, {3} = exact 3 instances, [1-4] = cijfers 1 t/m 4, \w{3}+\$ = de laatste 3 letters van string)

precision = TP / (TP + FP)
Recall Accuracy = TP / (TP + FN)
Accuracy Recall = TP + TN / Total F1 = 2 \* (Precision \* Recall) / (Precision + Recall)



np.arange (begin range, eind range, grootte van stappen) = array van cijfers

L.reshape (int(math.sqrt(len(L))), int(math.sqrt(len(L)))) = herwin dimension array, vierkant niet geluk

L.reshape (L.size, 1) = 1 dimensionale array maken

L[:, -1] geeft laatste kolom terug L[-1, :] geeft de laatste rij terug. L.mean (axis=0) column

Tel de 0e column 0 op, etc: for row in range(0, len(L)): for column in range(0, len(L[row]))

np.transpose(L) = spiegelen L L[row, column] += column return L

Middelste vierkant bestaande uit 4 cellen vierkant array: def middle\_square(L):

midden = np.empty((4,4), int), midden lijst = [], for row in range(int((len(L)-4)/2), int((len(L)+4)/2))

temp = [], for column in range(...): temp.append(L[row, column])
middenlijst.append(temp), midden = np.array(middenlijst) return midden

np.ones (m \* n, dtype = int) = maak grid met 1'en

Buren vinden: boven = grid[(i-1)%grid.shape[0], j%grid.shape[1]]
rechts = grid[(i-1)%grid.shape[0], (j+1)%grid.shape[1]]

testgrid = buren(buren(grid, i, j)) == grid[i, j] = schrijf 1 voor 1 de waarden in buren die matchen op basis van conditie (vergelijkbaar met DF)

kvr = pd.read\_csv('kvr.csv.gz', compression = 'gzip', sep = '\t', names = ['bestand', 'jaar', 'skipinitialspace = True])

Counter dict. most\_common(10) = 10 meest voorkomende

kvr.loc[kvr['Partij'].isin('poppartje-')] of kvr.loc[kvr['Partij'] == 'd66']

pd.crosstab(kvr.jaar, kvr.Partij, margins = True) = tel aantal combinaties voor beide kolommen

kvr['Partij'].str.lower()

kvr[kvr['Partij'].str.contains('vragen')] = drop alle Null in partij

kvr[kvr['Partij'].str.contains('vragen')] = verwijder alle rijen met vragen in de string

plt.bar (l1, l2) = horizontale bar graph plt.hist = histogram

pivot = df.pivot\_table(values = 'Count', index = 'Name', columns = 'Sex', aggfunc = sum, margins = True)

pivot = pivot.sort\_values(['All', 'Name'], ascending = [False, True])

kvr.Partij.sum = totaal van waarden partij. kvr.Partij.index.values = waarden van index

np.info(), np.arange? df.loc[df['a'] > 10, ['a', 'c']]

df1 = " "

df1 = " "

df1 = " " request.get(url)

5 = BeautifulSoup(soup(html\_doc, content))

comments = s.findAll("li", class="")

### Numpy

numpy.arange() → math array  
reshape()  
np.vstack()  
np.hstack()  
np.concatenate()

### DATAFRAME

df.drop("col", axis=0/1)

df.sort\_index() → sort by label

df.sort\_values() → by "country"

f = lambda x: x \* 2

df.apply(f)

df["col"].plot()

df["col"].value\_counts()

df.loc["row/col"]

→ df.iloc voor integers.

df.dropna(thresh, how)

df.ixrows()

for index, row in ..

with grip.open("sou.qz") as fn:  
for line in fn:

### Numpy slicing ID

x = np.arange(10)

x[:5] → eerste 5

x[5:] → element na 5

x[::2] → every other element

x[1::2] → ↑ starting out 1

x[::1] → all elements, reversed

x[5::-2] → reversed every other from index 5

### 2D

x2[2, :3] → 2 rows, 3 cols

x2[:, :3] → alle rows, eerste 3 cols

x2[:, 0] → eerste col

### Oefentaken

x2 = [3 5 24  
7 6 8 8  
1 6 7 7]

x1 = [5 0 3 3 7 9]

> x1cat5k → kolommen of first 2 rows:

x2[0:2, -3:]

> change each element of x1 by its squared difference from the mean of x1

np.square(x1 - np.mean(x1))

> apply a had mask to x1 and reduce x1 to its divisible by 3

x1[x1 % 3 == 0]

> use fancy indexing to get the diagonal of x2 (3, 6, 7)

x2[0, 1, 2], [0, 1, 2]

slides

precision = TP / (TP + FN)

Data science = based on patterns, predicting a value on new, unseen instances of data.  
steps: interact with outside world, prep - transform - modeling & computation, presentation

Impliciete info exploited maken = leesbaar maken voor pc, format doet computationeel snail is & aansluit bij breedte van die data

Nucleel excel = schuif net (geen heavy berekening), helpelt functionele, integratie andere tools

Onderzoeks cyclus: idee, hypothese, vraag, data verzamelen, opschonen, herstructureren, analyse: rekenen, doe statistieken, steekproef, rapporteer.

### outer products

x = np.arange(1, 6)

np.multiply.outer(x, x)

> count # of True values in bool array

np.count\_nonzero(x < 6)

> how many < 6 in each row?

np.count\_zeros(x < 6, axis=1)

> any values < 6?

np.any(x < 6) → same with np.all()

> sort each col of x

np.sort(x, axis=0)



\* Plot how many samples there are in the iris dataset in each species → iris[iris.species].value\_counts()

\* Create the df X containing all the numerical data in iris and the series y which contains species information  
 ⇒ X, y = iris.drop('species', axis=1), iris.species ⇒ X, y = iris[iris.columns[0:4]], iris.species

\* Sort iris first by 'sepal.length' and then by 'sepal.width', ⇒ iris.sort(by=['s.length', 's.width'], ascending=[True, False])

\* What are the maximal s-lengths for each species? ⇒ iris.groupby('species').s.length.max()

\* # Groupby('species') zorgt dat de kolom species alle zelfde values achter elkaar komen en alfabetisch geord.

\* Restrict iris to those samples with a s.length > 5 and a s.width < than  $\sqrt{s}$  ⇒ iris[(iris.s.length > 5) & (iris.s.width < np.sqrt(s))]

\* text = 'large string' ⇒ C = Counter(text.split()) CS = pd.Series(C).sort\_values(ascending=False)

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

\* pd.read\_csv(... .csv, sep='\t', compression='gzip', header=0, skipinitialspace=True, name=['id', ...], in)

\* len(cito.index) == len(cito.index.unique()) of len(cito.index) == len(set(cito.index))

\* dubbel = cito.index.value\_counts() dubbel[dubbel >= 2]

\* cito.sort\_index()

\* cito.str.lower(), str.strip(), str.replace(' ', ' '), str.replace("'", "'"), sort\_values()

\* cito ~~index~~ [cito.index.str.lower().str.contains('school')]

\* cito['RSM'] = np.sqrt((cito.quasicyto - cito.verwacht) \*\* 2)

\* Je wilt van elke waarde het gemiddelde van die waarden in die kolom aftrekken, ⇒ M = cito.mean(cito[M.index] - M)

\* Plot de inkomens per school, zorg netjes geordend. ⇒ cito.inkomen.sort\_values().plot()

\* In welke gemeente staan de scholen waar de inkomens meer dan 5000 zijn? Hoeveel scholen heb je dan per gemeente? ⇒ cito[cito['inkomen'] > 5000]['gemeente'].value\_counts().sort\_values(ascending=False).plot(kind='bar')

\* Welke files zitten in deze directory (.ls) en hoe groot zijn ze? ⇒ .ls -lh %lsmagic

\* Remove the file no-good.txt ⇒ rm no-good.txt

\* Hoeveel methodes die beginnen met een i heeft een verzameling in Python? ⇒ Een set aanmaken (a = set()) en daarna naam.iTab en dan tellen

Fancy indexing: x = np.random.randint(100, size=10) [x[3], x[7], x[2]] vs. x[[3, 7, 2]]

x2 = np.arange(12).reshape((3, 4)) x2[[0, 1, 2], [2, 1, 3]] ⇒ output: [2, 5, 11]

x[2, [2, 0, 1]] ⇒ [x[2, 2], x[2, 0], x[2, 1]]

Boolean array as Mask: x > 3 ⇒ x[x > 3]

Sorting arrays: x = np.array([2, 1, 3, 4, 5]) (np.sort(x) ⇒ array([1, 2, 3, 4, 5]) (np.argsort # returns indices ⇒ [1, 0, 3, 2, 4])

+ Regex: r'[a-z\d]'

a ⇒ Match any character that is not in the set

a-z ⇒ Range. Match any char. in the range

\d ⇒ Digit. Match any digit char. (0-9)

(.str.replace(r'[a-z\d]', ''))

↳ dit zorgt dat alle specifieke tekens wegvallen

\* np.where()

\* df.groupby('columns')

\* df.sort\_values(by=[column])

\* Operating on Null values:

- isnull(): generates a boolean mask indicating missing values.
- isnotnull(): opposite of isnull()
- dropna(): return a filtered version of the data #data.dropna(axis=1)
- fillna(): return a copy of the data with missing values filled or imputed

\* pd.Series (data, index=[index])

\* pd.DataFrame (np.array, columns=[...], index=[...])

\* A = pd.index([1, 3, 5, 7, 9]) # intersection (A & B) ⇒ Index([3, 5, 7])

B = pd.index([2, 3, 5, 7, 11])

# union (A | B) ⇒ Index([1, 2, 3, 5, 7, 9, 11])

\* data.loc['a'] ⇒ explicit

\* data.iloc[0] ⇒ implicit

\* ix index = hybrid of two (A ^ B) ⇒ Index([1, 2, 9, 11])

np.array (Erange(1, 1+6) For i in n[2, 4, 6])  
 np.zeros  
 np.ones  
 np.full  
 np.arange  
 np.random.random  
 np.random.randn

np.eye  
 np.empty  
 np.dtype  
 x.ndim  
 x.shape  
 x.size  
 x.itemsize  
 x.nbytes  
 x[-1]  
 x[-2]  
 x[2, 2]  
 x[5:]  
 x[2, 2]  
 x[:, 2] *alles behalte*  
 x[1::2]  
 x[:, 3, ::2]  
 all rows, except  
 x[:, 0] first column  
 x[0, :] first row  
 x[0, 0] = 99 *between*  
 x.shape[0] = x2[x2, :2].copy() *grid, reshape*  
 grid.reshape  
 x[np.newaxis, :]  
 concatenate (grid, grid), axis=1  
 np.vstack  
 np.hstack  
 np.sort(x)  
 np.argsort(x)  
 np.sort(x, axis=1)  
 np.partition  
 data[-1][name]

np.arcsin  
 np.exp exp2 Power  
 np.log log2  
 special.gamma  
 np.multiply  
 np.multiply.outer  
 np.random.random(100)  
 x.size  
 min (+) max (+)  
 M.sum() *M, axis=0*  
 M.argmax (axis=0)  
 np.all  
 np.any  
 np.sort\_values (col)  
 data[:, :].sort  
 np.std  
 mean  
 median  
 Percentile  
 plt.plot  
 array + 5

df.head  
 tail  
 shape  
 value\_counts  
 dropna  
 replace  
 rename  
 groupby  
 str.len  
 count

d.values  
 df.columns  
 df[ ]  
 df[1:2]  
 df['b']  
 ['c', 'e']  
 {2:4, 1:6}  
 df.loc  
 df.index  
 df.items  
 df[(data > 3) & data < 6, 3]  
 df[['a', 'e']]  
 df.loc[1]  
 df.iloc  
 df.iloc[5, 12]

pd.concat -> append  
 df.append (df2)  
 pd.merge  
 pd.merge (df, df2 on='hours') how = inner/outer/left/left

x < 3  
 7  
 >=  
 !=  
 --

np.count\_nonzero (x < 6)  
 np.sum (x < 6)  
 np.any ()  
 np.sum (~ (lin < 0.5 | in > 0.5))  
 ~ = not

x [x < 5]  
 in1 = [3, 7, 5]  
 x [in1]  
 row = [1, 2, 3]  
 col = [1, 2, 3]

x [row, col]  
 x [row[:, np.newaxis], met]  
 x [0] == 10

df.head  
 tail  
 shape  
 value\_counts  
 dropna  
 replace  
 rename

groupby  
 str.len  
 count

Recall  
 Precision  
 TP / TP + F  
 TP  
 TP + FP

Zieh nie  
 + 95 495  
 - 5 9405  
 Prec 95  
 95 + 495 = 0,1  
 Recall

- interaction  
 Preparation  
 Transformation  
 modeling and comp  
 Presentation

study of the generalizable  
 extraction of knowledge from  
 data  
 common epistemic requirement in  
 a science whether new knowledge is  
 actionable for decision making is  
 its predictive power, not just its  
 ability to explain the past  
 requires skills in  
 math, ml, AI  
 stat, db, optimized  
 + understand  
 the prob

df.loc [data, dens] loc, ['POP', 'dens']  
 df = pd.DataFrame (data, columns = [''])  
 np.exp, sin, pi  
 A.add (data, fill\_value=0)  
 random  
 stack  
 df.subtract (df['R'], axis=0)  
 df.iloc [0, :2]

isnull, notnull, dropna  
 POP [[i for i in POP.index if i[1] == 2010]]  
 health.loc [1, ('Bob', 'HR')]  
 data.sort\_index()  
 reset\_index

pd.concat -> append  
 df.append (df2)  
 pd.merge  
 pd.merge (df, df2 on='hours') how = inner/outer/left/left



help(function)  
DF = DataFrame  
NA = Numpy Array

%timeit

Recall =  $\frac{True\ P}{True\ P + False\ N}$

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

Numpy => np -> documentatie

Datatypes: Int [8/16/32/64], bool, float [32/64]

np.array([list], dtype="DT")

np.ones/zeros([length/(length, width)], dtype="DT")

np.random.[random/random/random](low, high, (height, width))

np.linspace(low, high, Nsteps)

np.array(low, high, step-size)

NA.[shape/size/ndim]

NA[ row-slice, col-slice]

np.[concatenate/stack/concat/dstack/split/np.split]=adding/splitting

VA.reshape((height, width))

np.[log2/log/log/sqrt]

np.[sum/[idx/arg][min/max],()

NA.[sum/mean/std]()

mask => NA[NA % 3 == 0]

pandas

pd.read\_[csv/html/excel]('file', args...)

Kolom => df.colname of df['colname']

antallen => df.value\_counts()[.plot(kind="")] => plotten

sorteren => df.sort\_[index/values](by='col', ascending=[false/true])

nee maken => df.[col/index].str.[strip/replace/lower]()

kolom delete => df.drop('col', axis=1)

groeperen => df.groupby('col')

crosstabel => df.crosstab(df1, df2)

entry in df => df.isin([df/list])

functie over => df.apply(func)

pivot => pd.pivot\_table(df, index='colname', columns=[col], values=)

aggfunc = np.sum, margins = True

query = pd.df.query('query') => Ub query 'col > 5 and col2 < @var'

var = 2

# Python

\* key word ? (TAB)

% paste Bij planken code

% Run external code

% !smagic !

% history = prev output

Print (out[0])

Print (-) = prev output

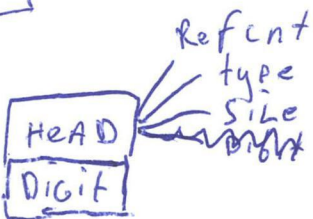
; = Suppress out

% +mode = Debug mode (VERBOSE)

% !python = where most time plain

% memit = memory

## Numpy



step = [::2]

Reverse = [::-1]

ReverseCol = [::2, ::-1]

Fcol = [::, 0]

slice = view

(out=) = where to save

REDUCE = func on ALL VALUES

(axis=0) columns

(axis=1) rows

x.reshape(-1, 3) = x.reshape(1, 3)

x[[3, 2, 9]] = x[3] ...

Shape of index ARRAY

argsort = indices of sort items

R1 -> Left +1

R2 -> 1 set to match

R3 -> any DIS AND no 1 = error

x[[i, l, p]] = 5 !

np.zeros(4, dtype=[(x, type)])

## PANDAS







# Text parsen

" ".join(list) → string van elementen van lijst  
 string.split() → lijst met woorden, split op spatie  
 BeautifulSoup(string) → string zonder HTML, + text alleen text  
 soup.findAll("li", class\_=re.compile("1 interlanguage-link"))  
 df.kolom.str.lower()  
 df.kolom.str.replace(item, " ")

df.fillna(0) *W:W* np.log2(?)  
 log2(16) = 4 np.absolute(?)  
 2<sup>4</sup> = 16 *Kolo*

# Numpy

%.ls - LH grootste file directory  
 Broadc. shapes equal no equal shape 1  
 %.rm remove tek  
 %.time var? je bin datatypen

pandas  
 pandas

df[[kolom1, kolom2]].min(axis=1)  
 df.drop(kolom, axis=1)  
 df.mean()  
 df.std()

(Array 1 + Array 2) / 2 gemiddelde tussen elementen in kolom

list(range(100))[1::2] oneven getalle in range  
 "[2::2] even

Array[rij][kolom]

Accuracy: TP / totaal 1  
 precision: TP / (TP + FP) 2  
 recall: TP / (TP + FN) 3  
 F<sub>1</sub> = 2 \* (precision \* recall) / (precision + recall)  
 Wilkehenij

df.groupby(k)[k2].max  
 Max van groep of k2  
 df.groupby(kolom).count()  
 Counter() → unieke item; lijst tellen

int(float) → geheel getal

Array + np.arange(0, Array.shape[1]) → index + het item

Array.mean(axis=1) → gemiddeldes kolommen

Array[Array % 10 == 0] → alleen waarden deelbaar door 10

Array.size → aantal elementen, len(array) → aantal rijen, Array.shape[0] aantal rijen

Array.reshape(rijen, kolommen) reshape(12, 1) id [12]

random.randint(0, 5), random.choice([ ]) → kan je ook exclude

# Pandas

link tot file

pd.read\_csv() \* sep = , skipinitialspace = , names = , index\_col

df[df.columns[0]] → eerste kolom

df[['string', 'string2']] → selecteer twee kolommen

pd.pivot\_table(values = koloms, index = unieke, columns = sex(F/M), margin = True)

df.kolom.value\_counts() → unieke waarden, count \* values → count \* index → alfabetisch  
 unieke waarden. \* sort\_index() → sorteren (by = [kolom, kolom2], ascending = [T, F])

df[~pd.isnull(df.kolom)] → bepaalde kolom rans uit halen

df[~df.kolom.str.contains(string)] → kolommen waarden in kolommen die bepaalde string bevatten verwijderen

df.kolom.loc[df.kolom >= 5] → rijen waar count >= 5, \* idxmax() → max index, \* max() → waarde



```

M.shape (2,3) -> (2,3) - (2,3)
a.shape (3,) -> (1,3) -> (2,3)
np.sqrt((x1-x1.mean())**2).mean()
(x1-x1.mean())**2
x1.std()
x1[x1%3==0]
x2[x2>1.23, x2<2.3]
for x in bgr:

```

```

pd.cite = pd.read_excel('...', index_col = ...)
cite.shape, cite.columns
dub = cite.index.value_counts()
dub[dub>1]
cite.sort_index()
cite.index.str.lower().str.strip().str.replace(' ','').sort_values()
cite.index.str.lower().str.contains('saloo')
cite.describe()

```

```

x,y = dg.drop('spe', axis=1), dg.spe
((x-x.mean())/x.std()).head()
iris.sort_values(['sepal', 'petal'], asc=[F,T])
iris.groupby('sepal')['petal'].max()
iris[iris.sepal>5] & (iris.sepal<(np.sqrt(5)))

```

```

cite['ren'] = np.sqrt((cite['u-cit.v']*2))
M = cite.mean(), (cite[M.index-1].tail())
cite.income.sort_values().plot()
cite[cite.income>50000].genre.value_counts().sort_val

```

```

a = np.sqrt(5)
iris.query('sepal>5 and sepal<@a')
Counter(text.split()) = C
pd.Series(C).sort_values(asc=False) = CS
CS[CS.index.str.len()==24].index
t = CS.index.str.len()*CS.values == 243
np.where(t) [1]
CS[CS.index==np.where(t)[1]]
CS[CS==1].sum()/CS.sum()
(CS==1).mean()
pd.Series(Counter(CS.values)).plot()
CS.value_counts().sort_index().plot()

```

```

url = 'wires'
html_doc = requests.get(url)
Soup = BeautifulSoup(html_doc.content)
Soup.pretty()
Soup.find_all('li', class_=re.compile('.*interlink.*')) [5]
a = Soup()[0].a
a.attrs

```

```

ii =
dict = {'a.attr': 'lang', 'a.attrs': 'href'} for i in 1:3
comments = s.jA('li', class_='comment')
C_list = [Soup.title, text, C.attrs['data-id'], C.find('span', class_='uname')]
Comm_df = pd.DataFrame(C_list)
Comm_df.columns = [' ', ' ', ' ', ' ']

```

```

x[0:2] = [0,2,4,6]
x[1:2] = [1,3,5,7]
x[0:-1] = [0,5,7,6,3]
x[0:-2] = [5,3,1]
x2[0:-1, 0:-1] -> reverse
x = np.newaxis, :3
x1, x2, x3 = np.split(x, [1,2])

```

```

bd = pd.read_html('...')
bd.groupby(0)[1].count().sort_values().tail(1)
w = bd.groupby(0)[1].count()
w.argmax(), w.max()
bd.groupby(0)[2].value_counts().sort_values().tail(1)

```

```

np.linspace(0,1,5)
np.random.random((3,3))
np.random.normal(3,3), np.random.randint(0,10)(3,3)
np.eye(3), np.empty, np.zeros
sum = np.add.reduce(x)
np.multiply.reduce(x)
" " . accumulate(x)

```

```

np.multiply.outer(x,x)
x = np.random.random((10,3))
x.mean = x.mean(0)
x.center = x - x.mean
np.sort(x)(x, axis=0,1)
np.argsort for index

```

```

np.sum(x < 6, axis=1)
np.any(x > 8)
np.all(x == 6)
data.isnull()
" " . notnull()
fillna(0)
pd.merge(df1, df2)
left-on
right-on
on

```

```

Regex: r'[a-zA-Z\d]'
^ = match any char. not inset.
a-z = range match ' ' in range
\d = digit Match any digit (0-9)
.str.replace(r'[a-zA-Z\d]', '')

```

```

data.keys()
list(data.items())
data[data>0.3]
8(data<0.8)

```

```

np.multiply.reduce(x)
" " . accumulate(x)

```

```

x [3], x [4], x [2]
x [2], [2,0:1]
x [1], [2,0:1]

```

```

np.sum(x < 6, axis=1)
np.any(x > 8)
np.all(x == 6)

```

```

data.isnull()
" " . notnull()
fillna(0)

```

Regex: r'[a-zA-Z\d]'  
 ^ = match any char. not inset.  
 a-z = range match ' ' in range  
 \d = digit Match any digit (0-9)  
 .str.replace(r'[a-zA-Z\d]', '')



```

x1 = array([5,0,3,3,7,9]) x2 = array([3,5,2,4],
                                         [7,6,8,8],
                                         [1,6,7,7])
last3: x1[-3:] = [3,7,9]
last3 col, first 2 rows: x2[:2,-3:]
x2.shape, x2.size: ((3,4),12) Reshape idim: x2.reshape
Squared diff each element: Q = (x1 - x1.mean())**2 (x2.size)
Std deviation: x1.std = np.sqrt(Q.mean())
Adding different shapes: "Broadcast regels worden eekel"
ints div by 3: x1[x1 % 3 == 0] = [0,3,3,9]
fancy index: x2[[0,1,2],[0,1,2]] = [3,6,7]
Sum of squares big array: ba = np.random.randint(10**3, size=10**6)
%timeit sum([x**2 for x in ba])
%timeit sum(ba**2)

```

```

tabel van 7 rux: np.arange(0,(14+7),7)
Reshape to square: length = int(math.sqrt(array.size))
array.reshape(length, length)
Reshape to idim: array.reshape(array.size,)
last row: array[-1:] | last column: array[:, -1]
mean per column: array.mean(axis=0)
middle square: Afreken = int(((array.shape[-1]-4)/2))
% array[Afreken:Afreken+4, Afreken:Afreken+4]

```

```

working file directory: %ls -lh | grep %.Rm no-good.txt W1
inspect algorithm: %timeit %timeit %timeit finddata.py x: x?

```

Data Science: Based on patterns in data prediction:  
 - a value on new unseen instances of data (classification class) - regression (grade)  
 - traditional statistics you explain  
 Steps: Interacting with outside world - Prediction - transformation  
 - modeling & computation - presentation  
 def: the study of the generalizable extraction of knowl.  
 Recall: TP/(TP+FN) | Precision: TP/(TP+FP)  
 A column: Mozilla/5.0  
 frame.a.dropna(), str.split('/'), str.get(0), value\_counts, head()  
 for mean, max, min, std: dataframe.describe()  
 number of unique names: dataframe.names.unique().size  
 Pivot table with year as index, gender as column, small bits as value  
 names.pivot\_table(values='bieth', index='year', columns='sex', aggfunc=Sum)  
 john = dataframe[dataframe.name='john']  
 femjohn = dataframe[dataframe.name='john'] && (names.sex='F')

impliciete informatie expliciet maken:  
 - leesbare maken voor een computer  
 - in een formaat dat computationeel snel is en aansluit bij het beeld dat de mens heeft van die informatie

nodulen excel: 1. schraaft niet (geen grote datasets, kan geen heavy berekeningen aan) 2. beperkte functionaliteit (opschonen van data is lastig) 3. integratie met andere tools (kan natuurlijk heel fijn in Python)

Onderzoeks cyclus: kom met idee/hypot/ vraag - Verzamel data - schoon je data op - heestruktuur je data - analyse: ga rekenen, doe statistics, toets hypot - Rapporteer  
 Pandas datastructurees: Series & DataFrame

```

Series: one dimensional labeled array (dict)
s = pd.Series([a,b,c,d,e], index=['a','b','c','d','e'])
S['a'] = 0.448597 | S.mean() = -0.337068
S[2] = -0.992218 | S.max() = 0.20144
S.idxmax() = 'a'
The dot product / cosine similarity:
[1,2,3] · [5,6,7] = (1·5) + (2·6) + (3·7)
S.dot(s)

```

```

DataFrame = df = pd.read_csv('x.csv', index_col='date')
eerste kolom: df.date | df.index
kolom a & b: df[['A','B']]
df.C > 3 (true/false)
df.B.median = 3 | df.loc['2009-01-02'] | row
df[df['B'] > 3] = rows that give true
df * 5 = mult everything by 5

```

date	A	B	C
2009-01-01	a	1	2
2009-01-02	b	3	4
2009-01-03	c	4	5

```

iris = pd.read_csv('bestand.csv')
iris: sepal.length sepal.width petal.l petal.w species
0 5.1 3.5 1.4 0.2 setosa
1 4.9 3.2 5.7 2.3 virginica
aantal samples each species + plot: iris.species.value_counts().plot(kind='bar')
create df x with all numerical: X = iris.drop('species', axis=1)
y with species: iris.species
2 normalize data: ((X - X.mean()) / X.std()).std() (overall)
sort by sepal.lw: iris.sort_values(['sepal.l', 'sepal.w'], ascending=[True, False])
max sepal length per species: iris.groupby('species')['sepal.l'].max()
sepal.l > 5 & sepal.w < 5 sqrt(s): x = np.sqrt(s)

```

```

Tokenize text: C = Counter(text.split())
CS = pd.Series(C).sort_values(ascending=False)
find all tokens which occupy 2 characters:
CS[CS.index.str.len() == 2].index
Percentage unique words: (CS == 1).mean()
Percentage all words: (CS == 1).sum() / CS.sum()
plot no of words how many times:
pd.Series(Counter(s.values)).plot(kind='bar'); CS.value_counts().sort_index().plot(kind='bar');

```

```

kvr = pd.read_csv('linktofile', sep='|', header=None, skipinitialspace=True)
kvr.columns = ['id', 'jaer', 'poetij', 'onderwerp', 'vraag', 'reactie', 'minis']
cross table jaer x poetij en top 10 poetij:
meest actief = kvr.poetij.value_counts()[0:10]
CROSS = pd.crosstab(kvr.poetij, kvr.jaer)
vraag counte = kvr.vraag.value_counts()
vraag counte = kvr.vraag.value_counts().str.count('i?') <= 50

```

```

pivot = df.pivot_table(values='count', index='name', columns='sex', margins=True, aggfunc=sum, sort_values=[False, True])
pivot['ratio'] = pd.Series(np.log2(pivot['m'] / pivot['f']), index=pivot.index)
pivot = pivot[(pivot['ratio'] <= 4) | (pivot['ratio'] >= 4)]
act_error = (pivot['error'].sum() / pivot['all'].sum()) * 100
pivot['kleinst'] = pivot[['m', 'f']].min(axis=1)
most_errors = pivot.kleinst.idxmax()
kvr.poetij.replace(to_replace='calle', value='', regex=True)

```

```

aantal bomen elke soort: boom.df.Boomsoort.str.bwee().value_counts()
aantal bomen per jaar: df.Plantjaar.dropna().astype(int).value_counts().sort_index().plot(kind='bar')
laait 2 kolommen zien: df[df['Plantjaar'] >= 2000][['Plantjaar', 'Boomsoort']]
verzoeken: stammen['stamdiameter'].str.replace('cm', '')

```

```

magic: % linemagic %% cellmagic
Pasting code blocks: % paste % cpaste
Run external code: % run
Timing: % time % timeit % !perm
magic func: % magic -> list: % !s magic % ! history
working dir: % pwd
Content size: % ls
change dir: % cd
make dir: % mkdir
move file: % mv file loc
Remove: % rm
all data var: % var?
mem usage:

```



a science = processing values on raw data  
 1. interact w/ outside world (read/write variety of file formats/datab.)  
 2. prep data  
 3. transformation (apply math/stats)  
 4. modelling + computation  
 5. presentation (visualization)

precision = how many of pred correct (TP/TP+FP)  
 recall = how many instances of true CC pred. to be CC?  
 accuracy = how many actually correct. (TP/TP+FN)

df = 2D labeled array any data type  
 df = 2D labeled data struct.

drop:  $kvr = kvr [kvr.index[kvr.party]] == row$   
 df.pivot\_table(values='count', index='name', columns='sex', aggfunc='sum')

kvr.stabel: pd.cross tab(x, y) [kvr.party-value-count]

np.array([1,2,3], dtype='float64')  
 np.zeros(4, dtype='int')  
 np.ones((3,5), dtype='float')  
 np.full((3,5), 7)  
 np.arange(0,20,2)  
 np.linspace(0,1,5)  
 np.random.random((3,3))  
 np.random.randint(0,10,(3,3))

pd.DataFrame(data)  
 # transform dict to df  
 # index  
 ind.size  
 ind.shape  
 ind.ndim  
 ind.dtype

cities = [r['cy'] for r in records if 'cy' in r]  
 from collections import Counter  
 cities = Counter(cities)  
 cities.most\_common(10)

**Amejet Rojgn**  
 11269111

df.ndim  
 df.shape  
 df.size  
 df.dtype

data.values  
 data.columns  
 data.index

df[2,0] = 12 # index en modify  
 df[::2] # every other element  
 x[1::2] # from index 1  
 x[:: -1] # all elem. reversed.  
 x[5::-2] # reverse from 5  
 x[2, :3] # two rows, three columns  
 x[3, ::2] # all rows, every other column  
 x[:, 0] first column  
 x[0, :] first row

# intersection  
 A & B  
 # union  
 A | B  
 # sym. difference  
 A ^ B

cit = pd.DataFrame.from\_dict(cities, orient='index')  
 cit.columns = ['aantal'] # geef naam  
 cit.sort\_values('aantal', ascending=False, inplace=True)

grid = np.arange(16).reshape((4,4))  
 np.concatenate([x,y,m])  
 np.vstack([x,grid]) vertically  
 np.hstack([grid,x]) stack array  
 np.vsplit(grid,[2]) split hor/  
 np.hsplit(grid,[2]) split vel.  
 np.absolute/abs(x)  
 np.add(x,2) of subtract, multiply, ma

data['a':'c']  
 data[['a','c']]  
 Fancy  
 # masking  
 data[data > 0.3] & (data < 4)

short\_records = [{'country\_code': r['c'], 'timezone': r['t2']} for r in records if 't2' in r and 'c' in r]  
 # how often time zone occurs  
 tz\_counts = frame.dropna().str.split().frame.tz.value\_counts()

ufuncs: array + s  
 np.sum(array)  
 np.min(a)  
 np.max(a)  
 df.min(axis=0)  
 df.min(axis=1)

# locate:  
 data.loc[1:3]  
 data.iloc[[1,3]]  
 for index  
 data.ix[3, 'pop']  
 data.loc[data.dens > 100, ['pop', 'dens']]  
 # transpose  
 data.T

( = frame.a.dropna().str.split('/').str.get(0).value\_counts().head(0)  
 # unique names  
 df.name.unique().size  
 # births per year per gender:  
 total\_births = names.pivot\_table(values='births', index='year', columns='sex', aggfunc='sum')

np.argmax/max → find index of minimum.  
 np.any() any/all elements  
 np.all() true?  
 np.sort(array)  
 np.argsort() # indexes of sorted el.  
 np.sort(x, axis=0) - sort column  
 axis=1 - sort row  
 np.partition(x, axis=1) smallest values left from index 3.

np.nan  
 isnull()  
 notnull()  
 df.dropna()  
 axis='columns'  
 how='all'  
 thresh=3  
 df.fillna(0) # fill missing values with 0's

# find col with name:  
 john = names[names.name == 'john']  
 subset = john['mary']  
 # percentage of all kids  
 totals = names.groupby('year')['births'].sum()  
 (johnandmary[john]/totals \* 100).round(2)  
 # percentage of all boys  
 totalsboys = names[names.sex == 'M'].groupby('year')['births'].sum()  
 (johnandmary[john]/totalsboys \* 100).round(2).head(5)

ind = [3,4,7]  
 x[ind] # returns these values  
 x[2,[2,0,1]] # from 3rd row get those indexes.

pd.concat([x,y])  
 of df1.append(df2)  
 pd.concat(x,y, ignore\_index=True)  
 df.merge(d1,d2)  
 df1.join(df2)

# n getallen deelbaar door 0  
 np.arange(0,n\*(0,0))  
 # variant array  
 n = int(math.sqrt(len(L)))  
 np.reshape(L,(n,n))  
 # reshape vertical  
 np.reshape(L,(L.size,1))  
 # laatste kolom vs row  
 L[-1, :] vs L[:, -1]  
 # array deelbaar door n  
 L[np.where(L % n == 0)]  
 # gem per kolom  
 (L.mean(axis=0))  
 # elue waarde + 1  
 L + np.arange(len(L[0]))

# replace special char  
 .replace('\W', '')  
 # remove words  
 kvr[party] = [re.sub(r'a|beiden+', '', str(p)) for p, i in kvr.items()]  
 # For row in kvr.party:  
 if row == "":  
 i = kvr.index[kvr[party] == row]  
 kvr.drop(i)

# name from lastraw:  
 data[-1]['name']  
 # name where age < 30  
 data[data['age'] < 30]['name']

# square:  
 x = int(L.shape[0]/2)  
 L[x-2:x+2, x-2:x+2]  
 # index  
 index[0]

# max value  
 Male.loc[Male['F'].idxmax()]



## NUMPY

2 range: generate array  
 reshape: change shape  
 sort: sorted copy  
 zeros: generate filled 0  
 add: add two arr  
 subtract: - two arr  
 multiply: \* two arr  
 divide: / two arr  
 exp:  $2^x$  two arr  
 sqrt:  $\sqrt{x}$  one arr  
 sum: sum arr  
 min: min of arr  
 max(a=0): max of arr  
 mean: mean of arr  
 mode: mode of arr  
 median: median of arr  
 vstack: 2 arr to matrix  
 concatenate: add 2 arr together  
 count\_nonzero: counts True  
 any: bool over any  
 all: bool over all  
 array: turn list into up  
 log2: log2 over arr  
 transpose: flip arr  
 where: find all to bool

## NUMPY SLICING

a[0:4, 3] on col 3  
 a[:, 1] whole col 1  
 a[::2] every other el  
 a[::-1] all, reversed  
 a[:: -2] every other reversed  
 a[2, :3] two rows, three cols  
 a[3, ::2] every other col  
 a[:, 0] first col  
 a[0, :] first row  
 a[:, -1, :: -1] mirror

## PANDAS

DataFrame: create df  
 iat/iloc: element at index  
 at/loc: element at name  
 ix: select row  
 drop: drop values  
 sort-index: sort along axis  
 sort-values: sort by values  
 rank: assign ranks  
 shape: (x,y) shape  
 index: arr of index col  
 columns: arr of column names  
 info: datatype info  
 count: number of non NA values  
 sum: sum of df  
 cumsum: accum. sum of df  
 apply: apply func to df  
 pivot table: pivot of df  
 crosstab: crosstable of df  
 isin: bool check  
 hist: histogram of df

## STEPS OF DATA SCIENCE

1. interacting with <sup>outside</sup> world
2. preparation
3. transformation
4. modeling & computing
5. presentation

## DEFINITION

Study of the generalizable extraction of knowledge

## CON EXCEL

Schaalt niet, beperkte function.

~ → not

| → or

& → and

for index, row in df.iterrows()

Str.split( delimiter )

del. join (list)

Str.replace(x, y)

## FORMULAS

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

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

## SNIPPETS

x[x % 3 == 0]

x[[0, 1, 2], [0, 1, 2]] → fancy index

df.col.value\_counts().plot.bar()

plt.plot

## ONDERZOEKSCYCLUS

hypothese / vraag → verzamel data → schoon  
 data op → herstructureer data → analyse (Statistiek, hypothesetoets) → rapporteer

pandas / Series: one dimen. labeled arr.  
 pandas / dataframe: two dimen. labeled

## REGEX

|w |d |s

|w |D |S

[abc] ↔ [^abc]

[a-z]

|t |n

(?=abc)

(?!abc)

\* 0+

+ 1+

? 0/1

? lazy.

re.sub

re.findall

re.match

re.search

## RANDOM

random.random

random.choice