

# DATA SCIENCE CHEATSHEET 2019

ASHWIN

`ROUND(2.04, 2)` Round of up 2 decimals

`Float(2)` mach komma getal

`List.append('item4')` add item to list

`List[:2]` selection of list

`print 'Henk is %s jaar oud' % str(leeftijd)`

`if 'Henk' in members:` check of Henk in members voorkomt

CREATE  
DATAFRAME

`df = pd.DataFrame(data, index=[1,2,3,4],`

`columns=['date', 'D', 'C', 'd'])` maak dataframe van strings

`df = pd.read_excel('http://url.nl.xls')` read file

`df = pd.read_csv('./output.csv', sep=';')`

`skiprows=1, header=None, names=['name', 'F', 'M'])`

SUMMARISE  
DATAFRAME

`df.head(10)`

`len(df)` # of rows in df

`len(df.columns.unique())` no distinct values

`df.columns.value_counts()` count # of rows for each unique value

`df.apply(pd.value_counts)` voer functie uit op alle columns van df

Doing stuff  
with  
DATAFRAME

`df.sort_values('column_name', ascending=False)` sorteer rijen

`df.sort_index()` sorteer index van df

`df.reset_index()` reset index, maakt # van rijen van, move old index to column

`new_df = df.set_index('date')` set column as index

`df.columns=['name', 'F', 'M']` set column names

`df.groupby(by='column_name')` group by values of in columns

`df.columns = df.columns.astype(int)` change datatype of values in column

`df.transpose()` draai df horizontaal, maak van columns rijen

`df.columns.min()` min max van column

`df.sum(axis=2)` van van waarde

`df.count()`

DELETE  
DATA

`df = df.drop(['length', 'height'], axis=1)` delete column

`df.drop_duplicates()` remove duplicates in rows

MISSING  
DATA

`df = df.dropna()` drop rows with any column with NaN

`df = df.dropna(axis=1)` drop rows with NaN in specific column

`df.dropna(subset=['column_name'])`

`df = df.fillna(value)` replace NaN with value

# Data Science cheatsheet

df.loc [voorwaarde] , df.column, iloc [ ]

df."column".value\_counts() Fancy indexing: data [[1,2]]

df."column".sort\_values()

df.index

df.columns

df = pd.read\_csv(data, index\_col=)

pd.Series(dict)

Regex for removing tags: <[^>]\*>

pd.DataFrame.from\_dict(data, orient, columns=)

~~df["column"]~~

.dict.get( )

df["column name"] = list

df[[column1, column2]].plot.bar()

.sort\_index()

.set\_index()

.groupby()

.idxmax(), .dropna(), .astype()

.stack() creates multiindex

.aggregate() can take list of multiple aggregates

.pivot\_table(data, values=, index=, columns=, aggfunc=)

pd.crosstab(titanic.sex, titanic['class']) hoeveel per geslacht per klasse

titanic.pivot\_table(index='sex', columns='class', values='survived')

prince.text.str.contains(r'\b[5s] ex\b').sum()

df ↑ nummers met het woord sex

prince.text.str.starts with ('sex').sum() + " ('sex')

↑ nummers beginnend met sex

axis=0 → rows  
axis=1 = columns

```
df.groupby(['Animal']).mean()  
(['Animals', *'Year', ...])
```

# count woman & man

```
1 df.sex.value_counts()
```

```
2 df.groupby('sex')['survived'].count()
```

```
3 {s: df[df.sex=s].sex.count() for s in df.sex.unique()}
```

```
4 df.pivot_table(index='sex', values='survived', aggfunc='len')
```

```
6 pd.crosstab(df.sex, df.survived).sum(axis=1)
```

# gemiddelde & medicijn een overlevende:

```
df.pivot_table(index='sex', columns='survived',  
values='age', aggfunc=[np.mean, np.median])
```

def teldeletters:

from collections import defaultdict

telling = defaultdict(int)

with open('bestand') as f:

for l in f:

for c in list(l):

telling[c] += 1

tseries = pd.Series(telling).sort\_values(asc=False)

return tseries

# dict comprehension:

```
{k: v**3 for k, v in genlist.items()}
```

```
pd.DataFrame.from_dict(data, orient='index')
```

```
# selectie = df[df['Year'] == '2018']
```

```
df1.join(df2, suffix="-my")
```

```
pd.Series(np.diag(df))
```

```
prince.text = prince.text.astype(str)
```

```
prince.text.contains(r'\b[ss]ex\b').sum()
```

```
prince.text.startswith('sex').sum()
```

# percentage winners

```
(100 * prince.text.findall(r'BEU(Stacougt)').str.len() /  
prince.text.len()) describe()
```

```
freq_serie = serie.values().nlargest(10)  
index
```

```
dummydf = pd.Series.apply(pd.Series).stack
```

```
dict = dict()
```

```
for b in freq_serie:
```

```
dict[b] = dict()
```

```
for a in freq_serie:
```

```
if dummydf[b].sum() > 0:
```

```
dict[b][a] = dummydf[dummydf
```

```
[b] == True][a].sum() /  
dummydf[b].sum()
```

```
else:
```

```
dict[b][a] = 0
```

```
def selectie(df):
```

```
df.dropna(inplace=True)
```

```
x = pd.DataFrame.from_dict(counter
```

```
(list(it.chain.from_iterable
```

```
(df.sections)), orient='index')
```

```
y = pd.Series(x[0]).sort_values()
```

```
return y[y > 5]
```

```
% is -th row column
```

```
% smagic x[1,2]
```

```
x[1,2,3] two rows three columns
```

```
x[1,3::2] rows, every other col
```

```
x[::1,::-1] reverse
```

```
x[:,0] first column
```

```
x[0] first row
```

```
x[x % 3 == 0] del base door 3
```

```
df.sort_values(['height', 'width', ascending=False, True])
```

```
df.groupby('Year')['length'].max()
```

```
df.index.str.lower().str.strip().str
```

```
.replace("x", "y").sort_values()
```

```
df.apply(np.mean) # mean of each col.
```

```
df.isin(['Breda', 'Tilburg'])
```

```
df['ratio'] = df['m'] / (df['m'] + df['v'])
```

counting:  
find unique

df.column.value\_counts()

df.groupby(column)[column2].count()

Ex: df[df.column == 5].column.count() for 5 in df.column.unique()

df.pivot\_table(index='column', values=~~column~~ column, aggfunc='kw')

p2. crosstab Titanic sex, Titanic survived). sum(axis=1)

p2. crosstab()

p2. pivot\_table(index, column, values, aggfunc='')

with open chestand, as f:

for l in f:  
code

• plot(kind='bar', 'hbar' etc.)

# Pandas

selecting:  
 df.loc[:, df > 1]. any  
 df.loc[:, df > 1]. all  
 df.loc[:, df.isnull().any]

indexing  
 df[df.country.isin(df2.type)]  
 df.filter(items=['a', 'b'])  
 df.select(lambda x: x \* 2)  
 df.where(s > 0)  
 df.query('second > first')

df.ix[:, 'capita']  
 df.ix[1, 'country']  
 df[(s > 1)]  
 df[(s < -1) | (s > 2)]  
 df[df['popul'] > 12000]  
 df['logic'] = np.where(df['A'] > 5, 'high', 'low')

ult' index  
 ip = pd.Series(populations, index=index)  
 det = pd.MultiIndex.from\_tuples(index)  
 p = pop.reindex(index)  
 ip\_df = pop.unstack()

## TITANIC 5 different ways

titanic.sex.value\_counts()  
 titanic.groupby('sex')['survived'].count()  
 titanic[titanic.sex == 'M'].sex.count() for s in titanic.sex.unique()  
 titanic.pivot\_table(index='sex', values='survived', aggfunc=len)  
 .crosstab([tit.sex, tit.survived], sum(axis=1))  
 titanic.pivot\_table(index='sex', columns='survived', values='age', aggfunc=[np.mean, np.median])  
 schrijf relatie class en pclass  
 x = crosstab(titanic.pclass, titanic['class'])  
 wanneer is adult male waar?  
 t = titanic.dropna(subset='age')  
 b = pivot\_table(index='age', columns='sex', values='adult\_male', aggfunc=sum)['male']

## prince

def teldeletters(bestand):  
 telling = defaultdict(int)  
 with open(bestand) as f:  
 for l in f:  
 for c in list(l):  
 telling[c] += 1  
 tseries = pd.Series(telling).sort\_values(ascending=False)  
 return tseries

lezen:  
 prince = pd.read\_csv(bestand)  
 prince.text = prince.text.astype(str)  
 prince.text.str.contains(r'\b(S|s|F|f)\b').sum()  
 prince.text.str.contains(r'\b(S|s|F|f)\b').sum()  
 prince.text.str.startswith('S').sum() + "" ('s').sum() / df2 = df[n-words] < df[n-words].quantile(.99)

df = pd.Series(''.join('S' if s else 's' for s in df))  
 means = (df.n-words / df.n-sections)  
 y = [mean\_words] = mean  
 df = df.dropna(subset=['mean\_words'])  
 mean\_words.describe()

df.replace('a', 'f')  
 df.fillna(0) (inplace = true for columns)  
 merge pd.merge(data1, data2, how='left', on='x1')  
 - left, right, inner, outer  
 data1.join(data2, how='right')  
 data1.append(data2)  
 pd.concat([s1, s2], axis=1, keys=['1', '2'])  
 pd.concat([data1, data2], axis=1, join='inner')

df.drop(['a', 'c']) = from rows  
 df.drop('country', axis=1) = from columns  
 f = lambda x: x \* 2  
 df.apply(f)  
 df.applymap(f)

pd.read\_csv(bestand, etc)  
 pd.read\_excel( )  
 sep = '\t', usecols =  
 test = pd.DataFrame(random.rand(20, 5))  
 pd.DataFrame.from\_dict(dict, orient='index')

def jaccard(i, j, df):  
 i = list(df.loc[df['title'] == i], 'sections')[0]  
 j = list(df.loc[df['title'] == j], 'sections')[0]  
 s1 = set(i)  
 s2 = set(j)  
 return len(s1.intersection(s2)) / len(s1.union(s2))

## (cito pandas)

> cito.index.str.lower().strip().str.replace(' ', '').str.replace(';', '').sort\_values()  
 - cito['isme'] = np.sqrt((cito.quas cito - cito.verwacht) \* 2)  
 - dubbel = cito.index.value\_counts()  
 dubbel[dubbel > 2]  
 len(wet(cito.index) == len(cito.index))

% is // churn\*  
 !head // churn.all  
 !cat 'churn.names'

other examples:  
 • nl-trans = raw.loc[crawl['orig\_lang'] == 'nl']  
 most-nl = nl-trans.sort\_values(by=['num\_trans'], ascending=True).head()  
 most-trans = set(most-nl['target\_lang', 'num\_trans']).apply(tuple, axis=1)  
 • total-trans = pd.pivot(raw, index='org', columns='target', values='num\_trans', aggfunc=sum, margins=True, fill\_value=0).reset\_index().rename\_axis(None, axis=1).sort\_values(by=['all'], ascending=False).iloc[1:]  
 total-trans = total-trans.groupby('org', axis=1)  
 total-trans = total-trans.reset\_index('org')  
 total-trans.head(20).plot(kind='bar', logy=True)

• TO series: total-trans.squeeze()  
 • corr = df.n-sections.corr(df.n-words)  
 df2 = df[n-words] < df[n-words].quantile(.99)  
 corr2 = df2.n-sections.corr(df2.n-words)

log = np.log(df.mean-words)  
 normal = df.mean-words  
 sns.distplot(logs)  
 sns.distplot(normal)

col. Value-Counts()

groupby(col)[col].count()

Pivot-table(index=col, value=col, aggfunc=)

resstab(df.col, df.col).sum(axis=1)

df[df.col==x].col.count() for x in df.col.unique()

## pandas vragen

- 5 manieren hoeveel m/v t. sex = male/female
  - 1) titanic.sex.value\_counts()
  - 2) titanic.groupby('sex')['survived'].count()
  - 3) Es: titanic[titanic.sex == S].sex.count()
   
for s in titanic.sex.unique():
  - 4) titanic.pivot\_table(index='sex', values='survived', aggfunc='len')
  - 5) pd.crosstab(titanic.sex, titanic.survived).sum(axis=1)
- Hoeveel per geslacht in elke klasse, welk deel survived?
  - 1) titanic.pivot\_table(index='sex', columns='class', values='survived', aggfunc='sum', aggfunc=len)
  - 1) pd.crosstab(titanic.sex, titanic['class'])
- gem & med. leeftyd per geslacht per overlevd/niet overlevd
   
titanic.pivot\_table(index='sex', columns='survived', values='age', aggfunc=[np.mean, np.median])
- beziehung relatie 2 var.
   
pd.crosstab(t['class'], t['class']) → in elkaar te vertalen
- wanneer adult = male wcar?
   
t = titanic.dropna(subset=['age'])
   
t.pivot\_table(index='age', columns='sex', values='adult\_male', aggfunc='sum')['male']
- pd series hoe vaak elke letter voorkomt
 

```
def teldeletters(bestand):
    from collections import Counter
    count = Counter()
    with open(bestand) as f:
        encoding = 'utf-8'
        for L in f:
            for c in list(L):
                count[c] += 1
    count_series = pd.Series(count).sort_values(ascending=False)
    return count_series
```
- prince.text = prince.text.astype(str) p.t.s. = prince.text.str
  - 1) bevat 'sex' of 'Sex'
   
p.t.s.contains(r'\b[3s]ex\b').sum()
  - 2) begint met 'sex' of 'Sex'
   
p.t.s.contains(r'^[3s]ex\b').sum()
   
of p.t.s.startswith('sex').sum() + p.t.s.startswith('Sex').sum()
  - 3) hoeveel bevat sex maar != 'sex', welke woorden
   
ef = p.t.s.findall(r'[3s]ex\w+')
   
ef[ef.str.len() > 0].count()
   
+ head(10)
  - 4) alle woorden printen
   
Counter([s for L in ef.values for s in L])
   
[a-zA-Z]
  - 5) % linkers
   
p.t.s.findall(r'[AEIOUaeiou]').str.len() / p.t.s.str.len() \* 100

- index\_col = 'naam' of set\_index('naam')
- dubbele namen index + hoe vaak
   
cito.index.value\_counts() [cito.i.v.c().values != 1]
- verwijder spaces begin/end, quotes
   
cito.index = cito.index.str.lower().str.replace(r'^\s+|\s+\$', '')
   
+ str.strip('')
- geen "school" in naam
   
cito[cito.index.str.contains('school')]
- cito['RMSE'] = np.sqrt((cito.quasi - cito.verw)\*\*2)
- normaliseren door v-ellie waarde gem v. col of te halen
   
mean = cito.mean()
   
cito[mean.index] - mean
   
'sklepprows = 1)
   
sns.distplot(df)
   
df.isnotnull().mean()

- max sep\_length for each species
   
ins.groupby('species')['sep\_length'].max()
- CS = pd series of word occurrence no
   
len word x occurrence == 24
   
CS[CS \* CS.index.str.len() == 24].index
- % of unique words in txt occurs 1x
   
CS[CS == 1].sum() / CS.count()
- % of all words in txt occurs 1x
   
CS[CS == 1].sum() / CS.sum()

## Theorie

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

accuracy 95% zwelke 1/1000

$$TP = 0,95 \times 0,001 \quad TN = 0,95 \times 0,999$$

$$FP = 0,05 \times 0,999 \quad FN = 0,05 \times 0,001$$

precision - de test zegt het, heb ik het ook  
 recall (pakkans - ik heb het, zegt de test het ook)  
 P(test positief, hebt het niet) = 1 - precisie

## voorwaardelijke kans

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{aantal } A \text{ \& B voorkomt} / \text{aantal } B \text{ voorkomt}$$

Min max norm  $\frac{X - \min(x)}{\max(x) - \min(x)}$  0 - 1  
 0.5 gem

Z-score  $\frac{X - \text{mean}(x)}{SD(x)}$  -4 to +4  
 0 = gem

IQE  
 so % tussen 70 & 80 IQE = 80 - 70 = 10  
 outlier als 80 + 1,5 \* 10 = 95 of 70 - 1,5 \* 10 = 55

streamend algoritme / top-k  
 bgv 5 kleinste lln:  
 heap = E3 # ln toe + ln len S  
 tot nu toe langste # langste in heap  
 # nieuw daarmee vgl.  
 with gap.open(...)

meeste  
 # bondsdag leden meer v. well land  
 $w = \text{bd.groupby}(0)[1].count()$   
 $w = \text{id} \times \text{max}(C), w = \text{max}(C)$   
 # welke party grootste aantal leden uit 1 decht  
 $\text{bd.groupby}(0)[2].count().argmax()$