

```

#parse xml
tree = lxml.etree.parse("Kandidatenlijsten_TK2017_Amsterdam.eml.xml")
#haal alle elementen uit xml
lst = []
candidates = tree.findall("//urn:oasis:names:tc:evs:schema:eml:Candidate")
affil = tree.findall("//urn:oasis:names:tc:evs:schema:eml:Affiliation")
for a in affil:
    candidates = a.findall("//urn:oasis:names:tc:evs:schema:eml:Candidate")
    for c in candidates:
        dct = {}
        dct["partij"] = a.find("{urn:oasis:names:tc:evs:schema:eml:AffiliationIdentifier}://{urn:oasis:names:tc:evs:schema:eml:RegisteredName}").text
        dct["initials"] = c.find("{urn:oasis:names:tc:evs:schema:eml:CandidateFullName}://{urn:oasis:names:tc:evs:schema:eml:NameLine}").text
        dct["voornaam"] = c.find("{urn:oasis:names:tc:evs:schema:eml:CandidateFullName}://{urn:oasis:names:tc:evs:schema:eml:FirstName}").text
        dct["gender"] = c.find("{urn:oasis:names:tc:evs:schema:eml:Gender}").text
        dct["achternaam"] = c.find("{urn:oasis:names:tc:evs:schema:eml:CandidateFullName}://{urn:oasis:names:tc:evs:schema:eml:LastName}").text
        dct["id"] = c.find("{urn:oasis:names:tc:evs:schema:eml:CandidateIdentifier}").values()
        dct["woonplaats"] = c.find("{urn:oasis:names:tc:evs:schema:eml:QualifyingAddress}://{urn:oasis:names:tc:evs:schema:eml:Locality}://{urn:oasis:names:tc:evs:schema:eml:LocalityName}").text
        lst.append(dct)

```

```

#neem specifieke partij
cda = df[df.partij == 'CDA']
#kruistabel van alle partijen en hun gender verdeling
geslacht = pd.crosstab(df.partij, df.gender, margins=True)
geslacht.head()

```

```

# man/vrouw verhouding met ratio
geslacht["ratio"] = geslacht["male"] / geslacht["female"]
#ratio = geslacht.apply(np.log2)
#ratio.loc[(ratio["ratio"] > 0.5) & (ratio["ratio"] < 0.5)].sort_values(ascending=False, by="All").head()
geslacht.head().sort_values(ascending=False, by="ratio")
#ratio

```

```

# man en vrouw verhouding, stacked graph kieslijsten
new = pd.crosstab(df.partij, df.gender).sort_values(['female'], ascending=True).plot(kind='bar', figsize=(5, 5), stacked=True)

```

```

#percentage man en vrouw
geslacht["avg_F"] = (geslacht["female"] / geslacht["All"])*100
geslacht["avg_M"] = (geslacht["male"] / geslacht["All"])*100
geslacht.sort_values(['avg_M'], ascending=False).head(10).round(2)

```

```

#aantal procent man en aantal procent vrouw gehele kieslijst
geslacht2 = pd.crosstab(df.partij, df.gender, margins=True)
total = np.sum(geslacht2.ix[['PVV (Partij voor de Vrijheid)'], ['male', 'female']].values)
(geslacht2.ix[['PVV (Partij voor de Vrijheid)'], ['male', 'female']].sum(axis=0)/total * 100).plot(kind='pie')

```

```

# partijleden van de vvd die in rotterdam of amsterdam wonen
vvd_ams_rot = df.loc[(df.partij == 'VVD') & ((df.woonplaats == 'Rotterdam') | (df.woonplaats == 'Amsterdam'))]

```

```

ams = "Kandidatenlijsten_TK2017_Amsterdam.eml.xml"
tree = ET.parse(ams)
root = tree.getroot()
data = []
for item in tree.iter():
    if item.tag == '{urn:oasis:names:tc:evs:schema:eml:Affiliation}':
        partij = item.getchildren()[0][0].text
        for candidate in item.findall("{urn:oasis:names:tc:evs:schema:eml:Candidate}"):
            id = candidate.getchildren()[0].attrib["id"]
            firstname = candidate.getchildren()[1][0][1].text
            lastname = candidate.getchildren()[1][0][2].text
            initialen = candidate.getchildren()[1][0][0].text
            gender = candidate.getchildren()[2].text
            place = candidate.getchildren()[3][0][0].text
            data.append({"initialen":initialen, "Partij":partij, "id":id, "naam": firstname, "achternaam": lastname, "geslacht": gender, "woonplaats":place})

```

1) kandidaten = pd.DataFrame(data)

```

counter = df.woonplaats.value_counts() #Value_counts sorteert ook meteen voor je
counter = counter[:20]
counter.plot(kind='barh', title = '20 meest populaire woonplaatsen voor de kandidaten van Amsterdam')

```

```

#PLOT DE ZETELVERDELING VOLGENS DE PEILINGWIJZER
peil[['Partij', 'Zetels']].plot(x='Partij', y='Zetels', kind='barh', title = 'Zetel verdeling')
#Crosstab Partij en Vrouwen hele kieslijst
vrouwen = df[df['gender'] == 'female']
vrouwen = pd.crosstab(vrouwen.partij, vrouwen.gender).sort_values(['female'], ascending=True).plot(kind='barh', figsize=(10, 10), stacked=True, title='Man/Vrouw')

```

```

#Hoeveel % van totaal kandidaten staat op de lijst voor specifieke partij?
(df.partij == 'VVD').mean()*100
#VIND ALLE KANDIDATEN DIE IN EEN VAN DE VOLGENDE STEDEN WONEN: BREDA, TILBURG, S-HERTOGENBOSCH, EINDHOVEN EN PLA
df_woonplaats = df['woonplaats'].loc[df['woonplaats'].isin(['Breda', 'Tilburg', 'Eindhoven'])]
df_woonplaats.value_counts().plot(kind='bar')

```

```

#s-HERTOGENBOSCH WORDT NIET GEVONDEN DOOR HET STREEPJE (HYPHEN) DUS MOET JE ZO ZOEKEN ALS JE DEZE WILT VINDEN
df['woonplaats'].loc[df['woonplaats'].str.contains('Hertogenbosch')]
#GEEF EEN LIJST VAN DE PARTIJEN DIE VROUWEN IN HUN KANDIDATEN LIJST HEBBEN
df['partij'].loc[(df['geslacht'] == 'vrouw')].unique()

```

```

#pivot table op basis van gender
piv = pd.pivot_table(df, index=['partij', 'gender'], aggfunc='count').dropna()
piv.piv.drop(['achternaam', 'initials', 'voornaam', 'woonplaats'], axis=1).rename(columns={'id': 'aantal'})

```

```

genders = kandidaten.groupby('Partij')['gender'].sum()
dfg = genders.to_frame()
dfg = dfg.reset_index()
compleet = peilingwijzer.merge(dfg, on='Partij')

```

```

kandidaten.gender.replace({'male': 'm', 'female': 'v'}, inplace=True)
vrouwenlaag = []
mannelaag = []
vrouwenhoog = []
mannelaag = []
vrouwen = []
mannelaag = []

```

```

#partijnamen goed zetten
peilingwijzer.Partij[2] = 'Partij van de Arbeid (P.v.d.A.)'
peilingwijzer.Partij[3] = 'PVV (Partij voor de Vrijheid)'
peilingwijzer.Partij[4] = 'SP (Socialistische Partij)'
peilingwijzer.Partij[5] = 'Democraten 66 (D66)'
peilingwijzer.Partij[7] = 'ChristenUnie'
peilingwijzer.Partij[8] = 'GroenLinks'
peilingwijzer.Partij[9] = 'Staatkundig Gereformeerde Partij (SGP)'
peilingwijzer.Partij[10] = 'Partij voor de Dieren'
peilingwijzer.Partij[12] = 'VNL (VoorNederland)'
peilingwijzer.Partij[13] = 'DENK'
peilingwijzer.Partij[14] = 'Forum voor Democratie'
peilingwijzer.Partij[15] = 'PiratesPartij'

```

```

df.c.replace('v', 'V')
df.c.replace('m', 'M')
df.c.replace('v', 'V')
df.c.replace('m', 'M')

```

#read peilingwijzer
peil = pd.read_excel('peilingwijzer.xlsx')
print peil.shape

vr = df2[df2['partij'] == 'PVV (Partij voor de Vrijheid)']
vr = vr[vr['nummer'] < 26]
vr.gender.value_counts()

%Is -Ik "File"
import pandas as pd
import numpy as np
import seaborn as sn
import lxml
from lxml import etree
%matplotlib inline
Kind: str
'line': line plot (default)
'bar': vertical bar plot
'barh': horizontal bar plot
'hist': histogram
'box': boxplot
'kde': Kernel Density Estimation plot
'density': same as 'kde'
'area': area plot
'pie': pie plot
'scatter': scatter plot
'hexbin': hexbin plot

help(pd.DataFrame)
%time

Meerdere files
for file in os.listdir(path):
if not file.endswith('.eml'): continue
with open(path, 'r') as f:
tree = ET.parse(path, file)

Min/Max/Aantal vrouwen/mannen
for item in compleet.iterrows():
vrouwenlaag.append(item[1]['geslacht'][:item[1]['ZetelsLaag']].count("v"))
compleet["vrouwenlaag"] = vrouwenlaag etc...

Import xml.etree.ElementTree as ET

date time

3) %matplotlib inline
ding = compleet[['Partij', 'mannen', 'vrouwen']]
ding[(ding.mannen > 0) | (ding.vrouwen > 0)]

RATIO + TOTAAL
df["total"] = (df["m"] + df["v"])
df["ratio"] = df["m"] / (df["m"] + df["v"])

Groupby
g = amsterdam.groupby(['partij', 'gender'])
g.size()

df.c * map({'female': 0, 'male': 1})

pd.merge(df1, df2, on=c)

df.c.replace('v', 'V')

```

from bs4 import BeautifulSoup
import pandas as pd
import re
from collections import Counter
import numpy as np
from future import division
import matplotlib inline
import seaborn

hertentamen= pd.read_excel('Cijfers_Peilingwijzer.xlsx')
hertentamen.shape
hertentamen.head()

#Open de XML
soup = BeautifulSoup(open('Kandidatenlijsten_TK2017_Amsterdam.eml.xml').read())
#zoek alle partij info
parties= soup('affiliation')
#zoek alle partijnamen
partynames= [p.text for p in soup.findAll('registeredname')]
#maak de genderlist per partij
genderlist= p.find('registeredname').text : [g.text for g in p.findAll('gender')] for p in parties

# Vertalen van partijnamen
translate={
    '506468': 'SOLUS',
    'CDA': 'CDA',
    '001': 'ChristenUnie',
    '066': 'Democraten 66 (D66)',
    'Denk': 'DENK',
    'Pvd': 'Forum voor Democratie',
    'GL': 'GroenLinks',
    'PVV': 'PVV (Partij voor de Vrijheid)',
    'PvdA': 'Partij van de Arbeid (P.v.d.A.)',
    'PvdD': 'Partij voor de Eieren',
    'SDP': 'Staatkundig Gereformeerde Partij (SGP)',
    'SP': 'SP (Socialistische Partij)',
    'VVD': 'VVD (Voor Nederland)',
    'VVD': 'VVD'
}

AantalVrouwenPerPartijPerAantalGewonnenZetels={p:N: Counter(genderlist[p][:N])[ 'female' ] for N in range(50)}
for p in genderlist)

df1 = pd.DataFrame.from_dict(AantalVrouwenPerPartijPerAantalGewonnenZetels, orient='index')
df1[49].plot(kind='barh')

# Voor de eerste N mensen op de lijst zet N =80 (is maximum) voor de hele lijst
N=50
#Partijnaam met aantal mannen en vrouwen
gencount={p:Counter(genderlist[p][:N]) for p in genderlist}

gdf=pd.DataFrame.from_dict(gencount, orient='index').fillna(0).astype(int)
#min lijst lengte= (gdf.sum(axis=1)).min()
#print min lijst lengte
# Voeg de ratio toe
gdf['Percentage'] = ((gdf.female / gdf.sum(axis=1))*100).round()
gdf.sort_values('Percentage', ascending=False)

ef=pd.DataFrame(translate.values())
ef.index=ef[0]
print ("Aantal vrouwen in de top (met max 50) van elke lijst")
ef.join(gdf).sort_values('Percentage')[ 'male', 'female', 'Percentage' ]

# Peilingwijzer
laatstepeiling="https://s3.eu-central-1.amazonaws.com/louwerse/Public/Peilingwijzer/Last/Cijfers_Peilingwijzer.xlsx"
allepeilingen="https://s3.eu-central-1.amazonaws.com/louwerse/Public/Peilingwijzer/Last/Results_Longitudinal.xlsx"
laatstedef= pd.read_excel(laatstepeiling)
alldf= pd.read_excel(allepeilingen)
print alldf.shape
alldf.tail()
laatstedef.index= laatstePartij
#laatstedef
laatstedef.rename(translate, inplace=True)

# Dit is een handige file met de zetel aantallen
alldf=pd.read_csv("https://dlbje97if6urz.cloudfront.net/Public/Peilingwijzer/Last/Results_DyGraphs_Seats.csv",
index_col='Date')
print alldf.shape
print (alldf.sum(axis=1) == 150).mean() # check if all add to 150
alldf = alldf.applymap(lambda s: int(s.split(':')[1])) # want (ci low; mean; ci high)
# now translate to our names
alldf=alldf.T
alldf.rename(translate, inplace=True)
alldf=alldf.T
alldf.tail()

alldf.rename(columns = {'PP': 'Piratenpartij'}, inplace=True)
for p in alldf.columns:
    alldf[p]= alldf[p].apply(lambda z: AantalVrouwenPerPartijPerAantalGewonnenZetels[p][z])
fem= alldf.plot(figsize=(18,7),
title='Aantal vrouwen in de Tweede Kamer door de tijd, per partij.\n Gel

som= alldf.T.sum()
print som.describe()
fem= alldf.T.sum().plot(figsize=(18,7),
title='Aantal vrou

fig=fem.get_figure()
fig.savefig('FemThroughTime.png')

df=laatstedef.join(gdf)
df

# Percentage naar aantal zetels
def Perc2Zetels(p):
    return round(p * (150/100))

def Ratio2Zetels(p):
    return round(p * 150)

ef= gdf.join(alldf.T.rename(translate)).dropna().drop(['male', 'female', 'Percentage'], axis =1)

ZetelsDoorDeTijd= ef.applymap(Ratio2Zetels)
ZetelsDoorDeTijd.head()

ef= ZetelsDoorDeTijd.T
for p in ef.columns:
    ef[p]= ef[p].apply(lambda z: AantalVrouwenPerPartijPerAantalGewonnenZetels[p][z])

fem= (ef.T.sum() / ZetelsDoorDeTijd.sum()*100).plot(figsize=(18,7),
title='Percentage vrouwen in de Tweede

#Hoeveel vrouwen bij ZetelsLaag CDA
len([1 for mp in genderlist['CDA'][:df.ZetelsLaag['CDA']] if mp=='female' ])

df.rename(index = {'PP': 'Piratenpartij'}, inplace=True)
def aantalvrouwen(p,column,df):
    zetels = df[column][p]
    aantalvrouwen= len([1 for mp in genderlist[p][:zetels ] if mp=='female' ])
    return aantalvrouwen

for c in [c for c in df.columns if c.startswith('Z')]:
    Pc = 'Fem'+c
    # compute aantal vrouwen
    Fem=pd.DataFrame.from_dict({p:aantalvrouwen(p,c,df) for p in df.index},orient='index')
    Fem.columns=[Pc]
    # voeg toe aan df
    df=df.join(Fem)
    # bereken de ratio
    df['FratioPeiling'+c]= (df[Pc] /df[c]).fillna(0)

# Aantal vrouwen in de kamer
VrouwenPeilingwijzer=df.sort_values('Zetels',ascending=True)[[ 'Zetels', 'FemZetels', 'u', 'FratioPeilingZetels' ]]

print "Op basis van deze peiling komen er %s vrouwen in de kamer. Dat is %s procent" % (VrouwenPeilingwijzer.FemZetels.sum(),
round(VrouwenPeilingwijzer.FemZetels.sum() /
VrouwenPeilingwijzer.Zetels.sum()*100))

print VrouwenPeilingwijzer

```


pd.DataFrame.from_dict(acantal_vrouwen['50 Plus'], orient='index')

```
# Open xml bestand
bestand = open('Kandidatenlijsten_TK2017_Amsterdam.eml.xml').read()
soup = bs.BeautifulSoup(bestand, 'xml')

# Namen van de partijen
partij_namen = []
for partij in soup.findAll('RegisteredName'):
    partij_namen.append(partij.text)
print partij_namen
```

```
# Zoek voor elke partij de vrouwen en de mannen
partijen = soup('Affiliation')
genderdict = {}
for partij in partijen:
    genderdict[partij.find('RegisteredName').text]
    = [gender.text for gender in partij.findAll('Gender')]
```

```
# Zoek namen
namendict = {}
for partij in partijen:
    namendict[partij.find('RegisteredName').text]
    = [name.text for name in partij.findAll('FirstName')]
namendict['50PLUS'].count('Henk')
```

```
# Maak dataframe mannen en vrouwen
vrouwen = {}
mannen = {}
for partij in genderdict:
    male_female = genderdict[partij]
    vrouwen[partij] = male_female.count('female')
    mannen[partij] = male_female.count('male')
vrouwen_df = pd.DataFrame.from_dict(vrouwen, orient='index')
mannen_df = pd.DataFrame.from_dict(mannen, orient='index')
```

```
pd.concat([vrouwen_df, mannen_df], axis=1)
# Per aantal zetels komen er x aantal vrouwen in de kamer
aantal_vrouwen = {partij:{N: Counter(genderdict[partij][:N])['female']
    for N in range(50)} for partij in genderdict}
Counter(genderdict[partij voor de Dieren][:25])['female']
gendercount = {p: Counter(genderdict[p][:N]) for p in genderdict}
sqdf = pd.DataFrame.from_dict(gendercount, orient='index', dtype=int)
```

IMPORTING DATA

```
pd.read_csv(filename) - From a CSV file
pd.read_table(filename) - From a delimited text file (like TSV)
pd.read_excel(filename) - From an Excel file
pd.read_sql(query, connection_object) - Read from a SQL table/database
pd.read_json(json_string) - Read from a JSON formatted string, URL or file.
pd.read_html(url) - Parses an HTML URL, string or file and extracts tables to a list of dataframes
pd.read_clipboard() - Takes the contents of your clipboard and passes it to read_table()
pd.DataFrame(dict) - From a dict, keys for columns names, values for data as lists
```

EXPORTING DATA

```
df.to_csv(filename) - Write to a CSV file
df.to_excel(filename) - Write to an Excel file
df.to_sql(table_name, connection_object) - Write to a SQL table
df.to_json(filename) - Write to a file in JSON format
df.to_html(filename) - Save as an HTML table
df.to_clipboard() - Write to the clipboard
```

CREATE TEST OBJECTS

```
Useful for testing
pd.DataFrame(np.random.rand(20,5)) - 5 columns and 20 rows of random floats
pd.Series(my_list) - Create a series from an iterable my_list
df.index = pd.date_range('1980/1/30', periods=df.shape[0]) - Add a date index
```

VIEWING/INSPECTING DATA

```
df.head(n) - First n rows of the DataFrame
df.tail(n) - Last n rows of the DataFrame
df.shape() - Number of rows and columns
df.info() - Index, Datatype and Memory information
df.describe() - Summary statistics for numerical columns
s.value_counts(dropna=False) - View unique values and counts
df.apply(pd.Series.value_counts) - Unique values and counts for all columns
```

value_counts
↳ count freq that value occurs
bv weenplaats hoe vaak

filteren df om plotten

```
peilingen.sort('zetels', ascending=False).plot('Party', 'Zetels', kind='bar')
```

translate = (0,0), (RegisteredName, 'D66', 'Democraten 66', 'D66')

```
laatstedef.index = laatstedef['Partij']
laatstedef.rename(translate, inplace=True)
# Voor alle waarden in kolom apply
zetels = zetels.applymap(lambda s: int(s.split(';')[1]))
nieuwdef = alledef.apply(lambda x: (x*150))
```

```
# Aantal vrouwen in de Tweede Kamer
for p in zetels.columns:
    zetels[p] = zetels[p].apply(lambda z: aantal_vrouwen[p][z])
zetels = dan_nieuwe_met_aantal_vrouwen(zetels.sort(as=1))
# Skippen met loc
alledef = alledef.T
alledef = alledef.iloc[:,3]
```

```
# String contains
barry_voor = hoofdelijke_stemming_df['pro'].str.contains('Barry Madlener')
len(adam_GL[adam_GL['voornaam'].str.startswith('T')])
```

```
# Crosstab
kruistabel = pd.crosstab(rutte2['proposalttype'], rutte2['result'], margins=True)
```

```
# Hoeveel procent stemmen heeft voor gestemd
partijnamen = rutte2.columns[11:50]
stemmen = rutte2[partijnamen].apply(pd.value_counts)
```

```
partij_dict = {}
for partij in partijnamen:
    partij_dict[partij] = float(stemmen[partij].loc[1])/float(stemmen[partij].sum())
```

```
dataframe = pd.DataFrame.from_dict(partij_dict, orient='index')
dataframe = dataframe.dropna()
dataframe = dataframe.sort(0, ascending=False)
```

```
plotje_partijen = dataframe.plot(kind='barh', title='Partijen die voor hebben gestu')
pd.DataFrame(populair).sort(0, ascending=False).plot(kind='barh', title='mk')
```

Sorten eerst dataframe maken oplopend

```
ascending=[True,False]) - Sort values by col1 in ascending order then col2 in descending order
```

```
df.groupby(col) - Return a groupby object for values from one column
```

```
df.groupby([col1,col2]) - Return a groupby object values from multiple columns
```

```
df.groupby(col1)[col2].mean() - Return the mean of the values in col2 grouped by the values in col1 (mean can be replaced with almost any function from the statistics section)
```

```
df.pivot_table(index=col1, values=[col2,col3], aggfunc=max) - Create a pivot table that groups by col1 and calculates the mean of col2 and col3
```

```
df.groupby(col1).agg(np.mean) - find the average across all columns for every unique column
```

```
data.apply(np.mean) - apply a function across each column
```

```
data.apply(np.max, axis=1) - apply a function across each row
```

JOIN/COMBINE

```
df1.append(df2) - Add the rows in df1 to the end of df2 (columns should be identical)
```

```
df.concat([df1, df2], axis=1) - Add the columns in df1 to the end of df2 (rows should be identical)
```

```
df1.join(df2, on=col1, how='inner') - SQL style join the columns in df1 with the columns on df2 where the rows for col have identical values. how can be one of 'left', 'right', 'outer', 'inner'
```

STATISTICS

These can all be applied to a series as well.

```
df.describe() - Summary statistics for numerical columns
```

```
df.mean() - Return the mean of all columns
```

```
df.corr() - finds the correlation between columns in a DataFrame.
```

```
df.count() - counts the number of non-null values in each DataFrame column.
```

```
df.max() - finds the highest value in each column.
```

```
df.min() - finds the lowest value in each column.
```

```
df.median() - finds the median of each column.
```

```
df.std() - finds the standard deviation of each column.
```

names.Count.idxmax()

↳ geeft index

```
names[names.Count == names.Count.median()]
```

pd.read_csv
zetels = zetels.set_index('Da
zetels = applymap.
zetels = zetels.T
zetels.rename
zetels = zetels.T
zetels2 = zetels.
zetels2.describe

zetels2.describe

Teu

mk

Filteren in een dataframe
peilingen.sort('zetels', ascending=True).plot('Partij', 'Zetels', kind =

```
from lxml import etree
from bz2file import BZ2File
import codecs
import urllib2
import urllib
import os
import seaborn as sn
import zipfile
import bs4 as bs
from collections import Count
```

```
import sys
stdout = sys.stdout
reload(sys)
sys.setdefaultencoding('utf-8')
sys.stdout = stdout
```

names.Count.idxmax()

names.Count.idxmax()

names[names.Count == names.Count.median()]

```

tree = ET.parse('/Users/Joel/Downloads/Kandidatenlijsten_EML_TK2017-20170213/
Kandidatenlijsten_TK2017_Amsterdam.eml.xml')
root = tree.getroot()
L = []
for hoi in root:
    for doei in hoi:
        for hee in doei:
            for nee in hee.findall('{urn:oasis:names:tc:evs:schema:eml}Affiliation'):
                for iden in nee.findall('{urn:oasis:names:tc:evs:schema:eml}AffiliationIdentifier'):
                    partij = iden.find('{urn:oasis:names:tc:evs:schema:eml}RegisteredName').text
            for iden in nee.findall('{urn:oasis:names:tc:evs:schema:eml}Candidate'):
                for i in iden.findall('{urn:oasis:names:tc:evs:schema:eml}CandidateIdentifier'):
                    ld = i.get('ld')
                for can in iden:
                    for person in can.findall('{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}PersonName'):
                        initial = person.find('{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}NameLine').text
                        name = person.find('{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}FirstName').text
                        surname = person.find('{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}LastName').text
                        gender = iden.find('{urn:oasis:names:tc:evs:schema:eml}Gender').text
                    for adres in can.findall('{urn:oasis:names:tc:ciq:xsd:schema:xAL:2.0}Locality'):
                        stad = adres.find('{urn:oasis:names:tc:ciq:xsd:schema:xAL:2.0}LocalityName').text
                    L.append({'partij': partij, 'ID': int(ld), 'initial': initial, 'name': name, 'surname': surname, 'gender': gender,
'stad': stad})
Readen pd.read_excel('path') (of read_csv)
Values van allen rutte2.proposaltype.value_counts()
alle votetype met roll call hoofdelijke_stemming_df = rutte2[(rutte2['votetype'] == 'Roll call')] (je kan em & 'en)
hoeveel adopted gemiddeld voorkomt (hoofdelijke_stemming_df['result'] == 'adopted').mean() * 100
Vind alles van Barry (handig als t een lijst heeft binnen in) hoofdelijke_stemming_df.pro.str.contains('Barry Madlener').sum()
Tabel pd.crosstab(rutte2.proposaltype, rutte2.result, margins=True)
2 manieren rutte2.pivot_table(index='proposaltype', columns='result', values='pro', aggfunc=len, margins=True)
Blijkbaar belangrijk im = rutte2[(rutte2.authorparty==rutte2.supporterparties) & (rutte2.authorparty != '[]')]
im.shape Checken hoe tabel eruit ziet
staafdiagram maken im.authorparty.value_counts().sort_values().plot(kind='barh', title='aantal')
Checken of er 568 boys zijn im.authorparty.value_counts().sum()==568
Column names print list(rutte2)
Columnen pakken, gemiddelde pakken en procent, NaN weg, sorten op values en histogram maken - find the
(rutte2[['partijen']].mean()*100).dropna().sort_values(ascending=False).plot(kind='bar')
Partijen stemmen pakken van normal en aantal stemmen aanwezig = rutte2[rutte2.votetype=='Normal']
[['partijnamen']].count()
Niet-stemmers weghalen aanwezig[(aanwezig > 0)]
Denk datum pakken denk = rutte2.dropna(subset=[u'Kuzu/Öztürk'])
denk.date.min()
columnen pakken met 1 of 0 zonder NaN stems = rutte2[rutte2.votetype=='Normal'][['partijen']].dropna(axis=1)
pakt eerste alle 1 PVV, dan de mean, dan sort je em, dan 2 decimaal en procent
stems[stems.PVV==1].mean().sort_values(ascending=False).round(2) *100
Zelfde value vinden hey.loc[hey['partij'].isin(['VVD', 'CDA', '50PLUS'])]

```

```

import pandas as pd
import numpy as np
import seaborn as sn
%matplotlib inline
import matplotlib.pyplot as plt

len(set(df['partij']))
len(peil['Partij'])

counter = PVV['gender'].value_counts()
counter.plot(kind = 'pie', title = 'Verhouding mannen tot
vrouwen binnen de PVV')
plt.show()
Groen = df[(df['partij'] == 'GROENLINKS') & (df['gender']
== 'male')]
len(Groen['name'].str.startswith('T'))
steden = df['stad'].value_counts()[:20]
steden.sort_values(ascending = False).plot(kind='bar',
title='hoi ik ben joel')
plt.show()

cross = pd.crosstab(df.stad, df.gender, margins=True)
#print cross.sort_values(by='All', ascending=False)[:20]
cross[(cross['female'] > cross['male'])]

peil = pd.read_excel('/Users/Joel/Documents/
Informatiekunde/Cijfers_Peilingwijzer.xlsx')
peil.sort_values(by='Zetels', ascending =
False).plot('Partij', 'Zetels', kind='bar', title='hoi ik ben
joel')

len(set(df['partij']))
len(peil['Partij'])

partij_kandidatenlijst = ['partijen']
partij_peilingwijzer = list(set(peil.Partij))
partij_peilingwijzer = sorted(partij_peilingwijzer)
partij_kandidatenlijst = sorted(partij_kandidatenlijst)

dic = {}
for partij in partij_kandidatenlijst:
    index = partij_kandidatenlijst.index(partij)
    dic[partij] = partij_peilingwijzer[index]
df.replace(dic, inplace=True)
dfj = df.join(peil.set_index('Partij'), on = 'partij')
dfj.update(dfj[['columnen']].fillna(0))

```



```

import zipfile
import pandas as pd
import xml.etree.ElementTree as ET
import numpy as np
import seaborn as sns
import urllib2 #voor als je geen wget gebruikt
import matplotlib inline

#downloaden bestand
#wget https://www.kiesraad.nl/binaries/kiesraad/documenten/rapporten/2017/2/definitieve-ks
url = urllib2.urlopen('https://www.kiesraad.nl/binaries/kiesraad/documenten/rapporten/2017/
with open('map.zip', 'wb') as code:
code.write(url.read())

#VIND ALLE MANNEN VAN PARTIJ EN ZET ZE ONDERELKAAR
df.loc[df['partij'] == 'VVD' & (df['geslacht'] == 'man')]

#TEL HET AANTAL MANNEN VAN EEN SPECIEFIEKE PARTIJ
len(df.loc[(df['partij'] == 'VVD') & (df['geslacht'] == 'man')])

#VIND ALLE KANDIDATEN DIE IN EEN BEPAALDE WOONPLAATS WONEN
df.loc[(df['woonplaats'] == 'Amsterdam')]

#VIND ALLE KANDIDATEN VAN EEN SPECIEFIEKE PARTIJ DIE IN EEN BEPAALDE WOONPLAATS WONEN
df.loc[(df['partij'] == 'CDA') & (df['woonplaats'] == 'Amsterdam')]

#VIND AANTAL KANDIDATEN DIE PER STAD WONEN EN PLOT DIT IN EEN pie PLOT
df['woonplaats'].value_counts().plot(kind='pie')

#PLOT DE VERHOUDING VAN WOONPLAATS VAN EEN SPECIEFIEKE PARTIJ
df['woonplaats'].loc[(df['partij'] == 'PVV (Partij voor de Vrijheid)')].value_counts().plot(kind='pie', figsize=(10,10))

#PLOT DE VERHOUDING VAN DE VERHOUDING MAN VROUW VAN EEN SPECIEFIEKE PARTIJ IN EEN PIE-GRAPH
df['geslacht'].loc[(df['partij'] == 'PVV (Partij voor de Vrijheid)')].value_counts().plot(kind='pie')

HERTOGENBOSCH WORDT NIET GEVONDEN DOOR HET STREEPJE (HYPHEN) DUS MORT JE ZO ROEKEN ALS JE DEZE WILT VINDEN
'woonplaats'].loc[df['woonplaats'].str.contains('Hertogenbosch')]

#VIND ALLE KANDIDATEN DIE IN EEN VAN DE VOLGENDE STEDEN WONEN: BREDA, TILBURG, S-HERTOGENBOSCH, EINDHOVEN EN PLOT
df_woonplaats = df['woonplaats'].loc[df['woonplaats'].isin(['Breda', 'Tilburg', 'Eindhoven'])]
df_woonplaats.value_counts().plot(kind='bar')

#PARTIJ MET HET AANTAL KANDIDATEN OP DE LIJST PLOT DIT IN EEN BAR GRAPH
df['partij'].value_counts().plot(kind='barh')

#Hoeveel kandidaten staan er per partij op de lijst plot dit
partijen = []
aantal = []
vrouwen = []
mannen = []

#FOR LOOP OM DE DATA TE GENEREN VOOR DE LIJSTEN
for a in df['partij'].unique():
partijen.append(a)
aantal.append(len(df.loc[(df['partij'] == a)]))
vrouwen.append(len(df.loc[(df['partij'] == a) & (df['gender'] == 'vrouw')]))
mannen.append(len(df.loc[(df['partij'] == a) & (df['gender'] == 'man')]))

#MAAK EEN DATAFRAME VAN DE LIJSTEN
kandidaten_per_partij = pd.DataFrame({'partijen':partijen, 'aantal':aantal, 'vrouwen':vrouwen, 'mannen':mannen})
kandidaten_per_partij

#GEEF HET PERCENTAGE MANNEN EN VROUWEN EN VOEG DIT TOE AAN HET NIEUWE DATAFRAME
kandidaten_per_partij['perc_vrouw'] = ((kandidaten_per_partij['vrouwen'] / kandidaten_per_partij['aantal']) * 100).round(2)
kandidaten_per_partij['perc_man'] = ((kandidaten_per_partij['mannen'] / kandidaten_per_partij['aantal']) * 100).round(2)

#PLOT DE VERHOUDING MAN VROUW PER PARTIJ
kandidaten_per_partij[['partijen', 'mannen', 'vrouwen']].plot(x='partijen', kind='barh', figsize=(10,10))

#Vind het aantal vrouwen en mannen die in de tweede kamer komen volgens de laatste peiling
partijen = []
vrouwen = []
mannen = []

for a in pw_df['Partij']:
partijen.append(a)
vrouwen.append(len(df.loc[(df['partij'] == a) & (df['positie_lijst'] <= (pw_df['zetels'].loc[pw_df['Partij'] == a]).values[0]) & (df['gender'] == 'vrouw')]))
mannen.append(len(df.loc[(df['partij'] == a) & (df['positie_lijst'] <= (pw_df['zetels'].loc[pw_df['Partij'] == a]).values[0]) & (df['gender'] == 'man')]))

man_vrouw_in_kamer = pd.DataFrame({'partijen':partijen, 'vrouwen':vrouwen, 'mannen':mannen})
man_vrouw_in_kamer[['partijen', 'mannen', 'vrouwen']]
man_vrouw_in_kamer.plot(kind='barh', x='partijen', figsize=(10,10))

#VINDEN ALLE KANDIDATEN WAARVAN HUN VOORNAAM MET EEN B BEGINT WONEN IN ADAM OF DIEMEN EN GEEF VOORNAAM EN PARTIJ WOONPL
df[['voornaam', 'partij', 'woonplaats']].loc[(df['voornaam'].str.startswith('B') & df['woonplaats'].isin(['Amsterdam',

#PLOT DE ZETELVERDELING VOLGENS DE PEILINGWIJZER
peilingwijzer[['Partij', 'Zetels']].plot(x='Partij', y='Zetels', kind='barh', title='Zetel verdeling')

#VIND ALLE UNIEKEN PARTIJEN VAN DE TWEE VERSCHILLENDE DATAFRAMES
df['partij'].unique(), peilingwijzer['Partij'].unique())

#COMBINEER ALLE LIJSTEN TOT EEN DATAFRAME
df = pd.DataFrame({'partij':partij, 'positie_lijst':positie_lijst, 'prefix':prefix, 'initialen':initialen, 'voornaam':voornaam, 'achternaam':achternaam, 'woonplaats':woonplaats, 'gender':gender})

df = df[['partij', 'positie_lijst', 'initialen', 'voornaam', 'prefix', 'achternaam', 'gender', 'woonplaats']]

df.head()
# vind alle tags: set(tags) voor kopiëren hierboven

#HEEST RECENTE PEILINGWIJZER INLADEN

#MAAK EEN DATAFRAME
peilingwijzer = pd.read_excel('cijfers_peilingwijzer.xlsx')

#GEEF DE MATEN EN DE DESCRIPTION VAN HET DATAFRAME
peilingwijzer.shape, peilingwijzer.describe()

peilingwijzer.head()

Partij Datum Percentage PercentageLaag PercentageHoog Zetels ZetelsLaag ZetelsHoog
1 VVD 2017-03-14 17,2 16,5 17,9 26 24 28

#Stel de partij namen van de twee dataframes aan elkaar gelijk
pw_df_namen = ['VVD', 'Partij van de Arbeid (P.v.d.A.)', 'PVV (Partij voor de Vrijheid)',
'SP (Socialistische Partij)', 'CDA', 'Democraten 66 (D66)', 'ChristenUnie',
'GROENLINKS', 'Staatkundig Gereformeerde Partij (SGP)',
'Partij voor de Dieren', '50PLUS', 'VNL (VoorNederland)',
'DENK', 'Forum voor Democratie', 'Piratenpartij']

pw_df['Partij'] = pw_df_namen
pw_df.head(10)

#VIND ALLE KANDIDATEN DIE OP DE LIJST STAAN VOOR DE PVDA
df.loc[(df['partij'] == 'Partij van de Arbeid (P.v.d.A.)').head()]

#GEF EEN SPECIEFIEK GETAL UIT HET DATAFRAME ALS INTEGER (ALS VARIABLEN GEBRUIKEN VOOR FUNCTIES)
integer = peilingwijzer['PercentageLaag'][(peilingwijzer.Partij == 'PVV (Partij voor de Vrijheid)')].values[0]
integer
13.0

#TEL ALLE VALUES IN EEN KOLON OP
peilingwijzer['Zetels'].sum()

#BEREKEN HET VERSCHIL VAN TWEE KOLONEN EN VOEG DIT TOE IN EEN EXTRA KOLON
peilingwijzer['VerschilHoogLaag'] = (peilingwijzer['ZetelsHoog'] - peilingwijzer['ZetelsLaag'])
peilingwijzer.head()

Partij Datum Percentage PercentageLaag PercentageHoog Zetels ZetelsLaag ZetelsHoog VerschilHoogLaag
1 VVD 2017-03-14 17,2 16,5 17,9 26 24 28 4

#VIND ALLE VROUWEN VAN DE VVD EN HET CDA, MET HET OF STATEMENT
df.loc[(df['partij'] == 'VVD') | (df['partij'] == 'CDA')] & (df['gender'] == 'vrouw')

#VIND ALLE MENSEN OP DE LIJST DIE ALS VOORNAAM GEERT HEBBEN
df.loc[(df['voornaam'] == 'Geert')]

#VIND ALLE MENSEN OP DE LIJST VAN DE VVD MET HET PREFIX: VAN
df.loc[(df['prefix'] == 'van') & (df['partij'] == 'VVD')]

#VIND ALLE VROUWEN DIE ALS DE VVD 26 ZETELS KRIJGEN IN DE KAMER KOMEN
df.loc[(df['partij'] == 'VVD') & (df['nummer'] <= 26) & (df['geslacht'] == 'vrouw')]

#MAAK EEN LIJST VAN ALLE UNIEKE PARTIJEN OP DE LIJST
df['partij'].unique()

#AANTAL PARTIJEN
len(df['partij'].unique())

```



```
import zipfile
import pandas as pd
import xml.etree.ElementTree as ET
import numpy as np
import seaborn as sns
import urllib2 #voor als je geen lwget gebruikt
import matplotlib inline
```

```
#downloaden bestand
#wget https://www.kiesraad.nl/binaries/kiesraad/documenten/rapporten/2017/2/definitieve-kr
url = urllib2.urlopen('https://www.kiesraad.nl/binaries/kiesraad/documenten/rapporten/2017/
with open('map.zip', 'wb') as code:
    code.write(url.read())
```

```
zfile = zipfile.ZipFile('map.zip') #wanneer je wget gebruikt, naam van url ipv map.zip
zfile.extractall()
```

```
file = 'Kandidatenlijsten_EML_TK2017-20170213/Kandidatenlijsten_TK2017_Amsterdam.eml.xml'
#MAAK EEN LIJST VOOR ELKE COLUMN DIE JE WILT HEBBEN
partij = []
positie_lijst = []
initialen = []
voornaam = []
prefix = []
achternaam = []
woonplaats = []
gender = []
prefix_val = ''
tags = []
```

```
#PARSE DE FILE MET DE ET.ITERPARSE
for elem,event in ET.iterparse(file):
    tag = event.tag
    value = event.text
    tags.append(tag)
    if tag == '{urn:oasis:names:tc:evs:schema:eml}RegisteredName':
        partij_val = value
    if tag == '{urn:oasis:names:tc:evs:schema:eml}CandidateIdentifier':
        positie_lijst.append(int(event.attrib.values()[0]))
    if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XML:2.0}NameLine':
        initialen.append(value)
    if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XML:2.0}FirstName':
        voornaam.append(value)
        partij.append(partij_val)
    if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XML:2.0}LastName':
        achternaam.append(value)
        prefix.append(prefix_val)
    if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XML:2.0}NamePrefix':
        prefix_val = value
    else:
        prefix_val = ''
    if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XML:2.0}LocalityName':
        woonplaats.append(value)
    if tag == '{urn:oasis:names:tc:evs:schema:eml}Gender':
        if value == 'male':
            gender.append('man')
        elif value == 'female':
            gender.append('vrouw')
```

XML inladen

```
#PLOT DE VERHOUDING VAN WOONPLAATS VAN EEN SPECIFIEKE PARTIJ
df[woonplaats].loc[df['partij'] == 'PVV (Partij voor de Vrijheid)'].value_counts().plot(kind='pie', figsize=(10,10))
#PLOT DE VERHOUDING VAN DE VERHOUDING MAN VROUW VAN EEN SPECIFIEKE PARTIJ IN EEN PIE-GRAPH
df['geslacht'].loc[df['partij'] == 'PVV (Partij voor de Vrijheid)'].value_counts().plot(kind='pie')
```

```
HERTOGENBOSCH WORDT NIET GEVONDEN DOOR HET STREEPJE (HYPHEN) DUS MOET JE ZO SOEKEN ALS JE DEZE WILT VINDEN
df[woonplaats].loc[df[woonplaats].str.contains('Hertogenbosch')]
#VIND ALLE KANDIDATEN DIE IN EEN VAN DE VOLGENDE STEDEN WONEN: BREDA, TILBURG, S-BERTOGENBOSCH, SINDHOVEN EN PLOT
df[woonplaats = df[woonplaats'].isin(['breda', 'tilburg', 'Eindhoven'])]
df[woonplaats.value_counts().plot(kind='bar')
#PARTIJ MET HET AANTAL KANDIDATEN OP DE LIJST PLOT DIT IN EEN BAR GRAPH
df[partij].value_counts().plot(kind='barh')
```

```
#hoeveel kandidaten staan er per partij op de lijst plot dit
partijen = []
aantal = []
vrouwen = []
mannen = []
#FOR LOOP OM DE DATA TE GENEREN VOOR DE LIJSTEN
for a in df[partij].unique():
    partijen.append(a)
    aantal.append(len(df.loc[df['partij'] == a]))
    vrouwen.append(len(df.loc[(df['partij'] == a) & (df['gender'] == 'vrouw')]))
    mannen.append(len(df.loc[(df['partij'] == a) & (df['gender'] == 'man')]))
```

```
#MAAK EEN DATAFRAME VAN DE LIJSTEN
kandidaten_per_partij = pd.DataFrame({'partijen':partijen, 'aantal':aantal, 'vrouwen':vrouwen, 'mannen':mannen})
kandidaten_per_partij = kandidaten_per_partij[['partijen', 'aantal', 'vrouwen', 'mannen']]
#GEEF HET PERCENTAGE MANNEN EN VROUWEN EN VOEG DIT TOE AAN HET NIEUWE DATAFRAME
kandidaten_per_partij['perc_vrouw'] = (kandidaten_per_partij['vrouwen'] /
kandidaten_per_partij['aantal']) * 100).round(2)
kandidaten_per_partij['perc_man'] = ((kandidaten_per_partij['mannen'] /
kandidaten_per_partij['aantal']) * 100).round(2)
```

```
#PLOT DE VERHOUDING MAN VROUW PER PARTIJ
kandidaten_per_partij[['partijen', 'mannen', 'vrouwen']].plot(x='partijen', kind='barh', figsize=(10,10))
#vind het aantal vrouwen en mannen die in de tweede kamer komen volgens de laatste peiling
partijen = []
vrouwen = []
mannen = []
```

```
for a in pw_df['Partij']:
    partijen.append(a)
    vrouwen.append(len(df.loc[(df['partij'] == a) &
(df['positie_lijst'] == (pw_df['zetels'].loc[pw_df['Partij'] == a]).values[0]) &
(df['gender'] == 'vrouw')]))
    mannen.append(len(df.loc[(df['partij'] == a) &
(df['positie_lijst'] == (pw_df['zetels'].loc[pw_df['Partij'] == a]).values[0]) &
(df['gender'] == 'man')]))
man_vrouw_in_kamer = pd.DataFrame({'partijen':partijen, 'vrouwen':vrouwen, 'mannen':mannen})
man_vrouw_in_kamer[['partijen', 'mannen', 'vrouwen']]
man_vrouw_in_kamer.plot(kind='barh', x=partijen, figsize=(10,10))
```

```
#VINDEN ALLE KANDIDATEN WAARVAN HUN VOORNAAM MET EEN B BEGINT WONEN IN ADAM OF DIEMEN EN GEEF VOORNAAM EN PARTIJ WOONPL.
df[['voornaam', 'partij', 'woonplaats']].loc[df['voornaam'].str.startswith('B')] & df[woonplaats'].isin(['Amsterdam',
#PLOT DE ZETELVERDELING VOLGENS DE PEILINGWIJZER
peilingwijzer[['Partij', 'Zetels']].plot(x='Partij', y='Zetels', kind='barh', title='Zetel verdeling')
```

```
#VIND ALLE UNIEKEN PARTIJEN VAN DE TWEE VERSCHILLENDE DATAFRAMES
df[partij].unique(), peilingwijzer[Partij].unique()
```

```
#COMBINEER ALLE LIJSTEN TOT EEN DATAFRAME
df = pd.DataFrame({'partij':partij, 'positie_lijst':positie_lijst, 'prefix':prefix, 'initialen':initialen, 'voornaam':v,
'achternaam':achternaam, 'woonplaats':woonplaats, 'gender':gender})
```

```
df = df[['partij', 'positie_lijst', 'initialen', 'voornaam', 'prefix', 'achternaam', 'gender', 'woonplaats']]
df.head()
# vind alle tags: set(tags) voor kopiëren hierboven
```

```
#MEEST RECENTE PEILINGWIJZER INLADEN
#MAAK EEN DATAFRAME
peilingwijzer = pd.read_excel('Cijfers_Peilingwijzer.xlsx')
#GEEF DE MATEN EN DE DESCRIPTION VAN HET DATAFRAME
peilingwijzer.shape, peilingwijzer.describe()
peilingwijzer.head()
```

Partij	Datum	Percentage	PercentageLaag	PercentageHoog	Zetels	ZetelsLaag	ZetelsHoog
--------	-------	------------	----------------	----------------	--------	------------	------------

```
#Stel de partij namen van de twee dataframes aan elkaar gelijk
pw_df_namen = ['VVD', 'Partij van de Arbeid (P.v.d.A.)', 'PVV (Partij voor de Vrijheid)',
'SP (Socialistische Partij)', 'CDA', 'Democraten 66 (D66)', 'ChristenUnie',
'GROENLINKS', 'Staatkundig Gereformeerde Partij (SGP)',
'Partij voor de Dieren', '50PLUS', 'VNL (VoorNederland)',
'DENK', 'Forum voor Democratie', 'Piratenpartij']
pw_df['Partij'] = pw_df_namen
pw_df.head(10)
```

```
#VIND ALLE KANDIDATEN DIE OP DE LIJST STAAN VOOR DE PVDA
df.loc[df['partij'] == 'Partij van de Arbeid (P.v.d.A.)'].head()
```

```
#GEEF EEN SPECIFIEK GETAL UIT HET DATAFRAME ALS INTEGER (ALS VARIABLEN GEBRUIKEN VOOR FUNCTIES)
integer = peilingwijzer['PercentageLaag'][peilingwijzer.Partij == 'PVV (Partij voor de Vrijheid)'].values[0]
integer
13.0
```

```
#TEL ALLE VALUES IN EEN KOLOM OP
peilingwijzer['zetels'].sum()
```

```
#BEREKEN HET VERSCHIL VAN TWEE KOLOMMEN EN VOEG DIT TOE IN EEN EXTRA KOLOM
peilingwijzer['verschilHoogLaag'] = (peilingwijzer['zetelsHoog'] - peilingwijzer['zetelsLaag'])
peilingwijzer.head()
```

Partij	Datum	Percentage	PercentageLaag	PercentageHoog	Zetels	ZetelsLaag	ZetelsHoog	VerschilHoogLaag	
1	VVD	2017-03-14	17.2	16.5	17.9	26	24	28	4

```
#VIND ALLE MANNEN VAN PARTIJ EN ZET ZE ONDERELKAAR
df.loc[(df['partij'] == 'VVD') & (df['geslacht'] == 'man')]
```

```
#TEL HET AANTAL MANNEN VAN EEN SPECIFIEKE PARTIJ
len(df.loc[(df['partij'] == 'VVD') & (df['geslacht'] == 'man')])
```

```
#VIND ALLE KANDIDATEN DIE IN EEN BEPAALDE WOONPLAATS WONEN
df.loc[df[woonplaats] == 'Amsterdam']
```

```
#VIND ALLE KANDIDATEN VAN EEN SPECIFIEKE PARTIJ DIE IN EEN BEPAALDE WOONPLAATS WONEN
df.loc[(df['partij'] == 'CDA') & (df[woonplaats] == 'Amsterdam')]
```

```
#VIND AANTAL KANDIDATEN DIE PER STAD WONEN EN PLOT DIT IN EEN pie PLOT
df[woonplaats].value_counts().plot(kind='pie')
```

```
#PLOT DE VERHOUDING VAN WOONPLAATS VAN EEN SPECIFIEKE PARTIJ
df[woonplaats].loc[df['partij'] == 'PVV (Partij voor de Vrijheid)'].value_counts().plot(kind='pie', figsize=(10,10))
```

```
#PLOT DE VERHOUDING VAN DE VERHOUDING MAN VROUW VAN EEN SPECIFIEKE PARTIJ IN EEN PIE-GRAPH
df['geslacht'].loc[df['partij'] == 'PVV (Partij voor de Vrijheid)'].value_counts().plot(kind='pie')
```

```
HERTOGENBOSCH WORDT NIET GEVONDEN DOOR HET STREEPJE (HYPHEN) DUS MOET JE ZO SOEKEN ALS JE DEZE WILT VINDEN
df[woonplaats].loc[df[woonplaats].str.contains('Hertogenbosch')]
#VIND ALLE KANDIDATEN DIE IN EEN VAN DE VOLGENDE STEDEN WONEN: BREDA, TILBURG, S-BERTOGENBOSCH, SINDHOVEN EN PLOT
df[woonplaats = df[woonplaats'].isin(['breda', 'tilburg', 'Eindhoven'])]
df[woonplaats.value_counts().plot(kind='bar')
#PARTIJ MET HET AANTAL KANDIDATEN OP DE LIJST PLOT DIT IN EEN BAR GRAPH
df[partij].value_counts().plot(kind='barh')
```

```
#hoeveel kandidaten staan er per partij op de lijst plot dit
partijen = []
aantal = []
vrouwen = []
mannen = []
#FOR LOOP OM DE DATA TE GENEREN VOOR DE LIJSTEN
for a in df[partij].unique():
    partijen.append(a)
    aantal.append(len(df.loc[df['partij'] == a]))
    vrouwen.append(len(df.loc[(df['partij'] == a) & (df['gender'] == 'vrouw')]))
    mannen.append(len(df.loc[(df['partij'] == a) & (df['gender'] == 'man')]))
```

```
#MAAK EEN DATAFRAME VAN DE LIJSTEN
kandidaten_per_partij = pd.DataFrame({'partijen':partijen, 'aantal':aantal, 'vrouwen':vrouwen, 'mannen':mannen})
kandidaten_per_partij = kandidaten_per_partij[['partijen', 'aantal', 'vrouwen', 'mannen']]
#GEEF HET PERCENTAGE MANNEN EN VROUWEN EN VOEG DIT TOE AAN HET NIEUWE DATAFRAME
kandidaten_per_partij['perc_vrouw'] = (kandidaten_per_partij['vrouwen'] /
kandidaten_per_partij['aantal']) * 100).round(2)
kandidaten_per_partij['perc_man'] = ((kandidaten_per_partij['mannen'] /
kandidaten_per_partij['aantal']) * 100).round(2)
```

```
#PLOT DE VERHOUDING MAN VROUW PER PARTIJ
kandidaten_per_partij[['partijen', 'mannen', 'vrouwen']].plot(x='partijen', kind='barh', figsize=(10,10))
#vind het aantal vrouwen en mannen die in de tweede kamer komen volgens de laatste peiling
partijen = []
vrouwen = []
mannen = []
```

```
for a in pw_df['Partij']:
    partijen.append(a)
    vrouwen.append(len(df.loc[(df['partij'] == a) &
(df['positie_lijst'] == (pw_df['zetels'].loc[pw_df['Partij'] == a]).values[0]) &
(df['gender'] == 'vrouw')]))
    mannen.append(len(df.loc[(df['partij'] == a) &
(df['positie_lijst'] == (pw_df['zetels'].loc[pw_df['Partij'] == a]).values[0]) &
(df['gender'] == 'man')]))
man_vrouw_in_kamer = pd.DataFrame({'partijen':partijen, 'vrouwen':vrouwen, 'mannen':mannen})
man_vrouw_in_kamer[['partijen', 'mannen', 'vrouwen']]
man_vrouw_in_kamer.plot(kind='barh', x=partijen, figsize=(10,10))
```

```
#VINDEN ALLE KANDIDATEN WAARVAN HUN VOORNAAM MET EEN B BEGINT WONEN IN ADAM OF DIEMEN EN GEEF VOORNAAM EN PARTIJ WOONPL.
df[['voornaam', 'partij', 'woonplaats']].loc[df['voornaam'].str.startswith('B')] & df[woonplaats'].isin(['Amsterdam',
#PLOT DE ZETELVERDELING VOLGENS DE PEILINGWIJZER
peilingwijzer[['Partij', 'Zetels']].plot(x='Partij', y='Zetels', kind='barh', title='Zetel verdeling')
```

```
#VIND ALLE UNIEKEN PARTIJEN VAN DE TWEE VERSCHILLENDE DATAFRAMES
df[partij].unique(), peilingwijzer[Partij].unique()
```

```
#MAAK EEN CROSSTAB VAN DE WOONPLAATS EN HIER DE MAN VROUW VERHOUDING, PLOT DIT IN EEN STACKED PLOT
#ALLEZ ACHTER ELKAAR
pd.crosstab(df.woonplaats,
df.gender).sort_values(by='man',
ascending=False).head(5).plot(stacked=True,
colormap='Greens',
kind='barh',
title='Man vrouw verhouding top 5 grootste steden')
```

```
#Maak een formule die de genders naar nullen en eënen maakt, vervolgens maak je hier een nieuwe kolom voor. Hierna deel je
alle vrouwen door het totaal aantal leden van de partij zodat je erachter komt hoeveel procent vrouw is.
def gender_to_num(x):
    if x == 'man':
        return 0
    if x == 'vrouw':
        return 1
df['gender_nieuw'] = df.gender.apply(gender_to_num)
```

```
a = df.groupby('partij').gender_nieuw.sum() / df.partij.value_counts() * 100
a.round(2)
```

```
#VIND DE SPECIFIEKE NUMMER VAN EEN PERSOON AAN DE HAND VAN VOORNAAM EN ACHTERNAAM
df['nummer'].loc[(df['voornaam'] == 'Mark') & (df['achternaam'] == 'Rutte')]
```

```
#TOTAAL AANTAL KANDIDATEN VOOR SPECIFIEKE LIJST
len(df.loc[(df['partij'] == 'VVD')])
#HOEVEEL % VAN TOTAAL KANDIDATEN STAAT OP DE LIJST VOOR SPECIFIEKE PARTIJ?
(df.partij == 'VVD').mean()*100
```

```
#vind alle vrouwen van de vvd en het cda, met het of statement
df.loc[(df['partij'] == 'VVD') | (df['partij'] == 'CDA')] & (df['gender'] == 'vrouw')
```

```
#VIND ALLE MENSEN OP DE LIJST DIE ALS VOORNAAM GEERT HEBBEN
df.loc[df['voornaam'] == 'Geert']
```

```
#VIND ALLE MENSEN OP DE LIJST VAN DE VVD MET HET PREFIX: VAN
df.loc[(df['prefix'] == 'van') & (df['partij'] == 'VVD')]
```

```
#VIND ALLE VROUWEN DIE ALS DE VVD 26 ZETELS KRIJGEN IN DE KAMER KOMEN
df.loc[(df['partij'] == 'VVD') & (df['nummer'] <= 26) & (df['geslacht'] == 'vrouw')]
```

```
#MAAK EEN LIJST VAN ALLE UNIEKE PARTIJEN OP DE LIJST
df[partij].unique()
#AANTAL PARTIJEN
len(df[partij].unique())
```

DF

Manen


```
#parse xml
tree = lxml.etree.parse("Kandidatenlijsten_TK2017_Amsterdam.eml.xml")
#haal alle elementen uit xml
lst = []
candidates = tree.findall('//*[urn:oasis:names:tc:evs:schema:eml]Candidate')
affil = tree.findall('//*[urn:oasis:names:tc:evs:schema:eml]Affiliation')
for a in affil:
    candidates = a.findall('//*[urn:oasis:names:tc:evs:schema:eml]Candidate')
    for c in candidates:
        dct = {}
        dct['partij'] = a.find('//*[urn:oasis:names:tc:evs:schema:eml]AffiliationIdentifier//*[urn:oasis:names:tc:evs:schema:eml]RegisteredName').text
        dct['initials'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]CandidateFullName//*[urn:oasis:names:tc:evs:schema:eml]NameLine').text
        dct['voornaam'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]CandidateFullName//*[urn:oasis:names:tc:evs:schema:eml]FirstName').text
        dct['gender'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]Gender').text
        dct['achternaam'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]CandidateFullName//*[urn:oasis:names:tc:evs:schema:eml]LastName').text
        dct['id'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]CandidateIdentifier').values()
        dct['woonplaats'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]QualifyingAddress//*[urn:oasis:names:tc:evs:schema:eml]Locality//*[urn:oasis:names:tc:evs:schema:eml]LocalityName').text
        lst.append(dct)
```

```
#neem specifieke partij
cda = df[df.partij == 'CDA']
#kruistabel van alle partijen en hun gender verdeling
geslacht = pd.crosstab(df.partij, df.gender, margins=True)
# man/vrouw verhouding met ratio
geslacht["ratio"] = geslacht["male"] / geslacht["female"]
#ratio = geslacht.apply(np.log2)
#ratio.loc[(ratio["ratio"] > -0.5) & (ratio["ratio"] < 0.5)].sort_values(ascending=False, by="All").head()
geslacht.head().sort_values(ascending=False, by="ratio")
#ratio
```

```
# man en vrouw verhouding, stacked graph kieslijsten
new = pd.crosstab(df.partij, df.gender).sort_values(['female'], ascending=True).plot(kind='bar', figsize=(5, 5), stacked=True)
```

```
#percentage man en vrouw
geslacht["avg_F"] = (geslacht["female"] / geslacht["All"]) * 100
geslacht["avg_M"] = (geslacht["male"] / geslacht["All"]) * 100
geslacht.sort_values(['avg_M'], ascending=False).head(10).round(2)
```

```
#aantal procent man en aantal procent vrouw gehele kieslijst
geslacht2 = pd.crosstab(df.partij, df.gender, margins=True)
total = np.sum(geslacht2.ix[['PVV (Partij voor de Vrijheid)'], ['male', 'female']].values)
(geslacht2.ix[['PVV (Partij voor de Vrijheid)'], ['male', 'female']].sum(axis=0) / total * 100).plot(kind='pie')
```

```
# partijleden van de vvd die in rotterdam of amsterdam wonen
vvd_ams_rot = df.loc[(df.partij == 'VVD') & ((df.woonplaats == 'Rotterdam') | (df.woonplaats == 'Amsterdam'))]
```

```
ams = "Kandidatenlijsten_TK2017_Amsterdam.eml.xml"
tree = ET.parse(ams)
root = tree.getroot()
data = []
for item in tree.iter():
    if item.tag == '//*[urn:oasis:names:tc:evs:schema:eml]Affiliation':
        partij = item.getchildren()[0][0].text
        for candidate in item.findall('//*[urn:oasis:names:tc:evs:schema:eml]Candidate'):
            id = candidate.getchildren()[0].attrib["id"]
            first_name = candidate.getchildren()[1][0][1].text
            last_name = candidate.getchildren()[1][0][2].text
            initialen = candidate.getchildren()[1][0][0].text
            gender = candidate.getchildren()[2].text
            place = candidate.getchildren()[3][0][0].text
            data.append({"initialen": initialen, "Partij": partij, "id": id, "naam": first_name, "achternaam": last_name, "geslacht": gender, "woonplaats": place})
```

```
counter = df.woonplaats.value_counts() #Value_counts sorteert ook meteen voor je
counter = counter[:20]
counter.plot(kind = 'barh', title = '20 meest populaire woonplaatsen voor de kandidaten van Amsterdam')
```

```
#PLOT DE ZETELVERDELING VOLGENS DE PEILINGWIJZER
peil[['Partij', 'Zetels']].plot(x='Partij', y='Zetels', kind='barh', title='Zetel verdeling')
```

```
#BEREKEN HET VERSCHIL VAN TWEE KOLONNEN EN VOEG DIT TOE IN EEN EXTRA KOLM
peil['VerschilHoogLaag'] = (peil['ZetelsHoog'] - peil['ZetelsLaag'])
peil.head()
```

```
#s-HERTOGENBOSCH WORDT NIET GEVONDEN DOOR HET STREEPJE (HYPHEN) DUS MOET JE ZO ZOEKEN ALS JE DEZE WILT VINDEN
df['woonplaats'].loc[df['woonplaats'].str.contains('Hertogenbosch')]
#GEEF EEN LIJST VAN DE PARTIJEN DIE VROUWEN IN HUN KANDIDATEN LIJST HEBBEN
df['partij'].loc[(df['geslacht'] == 'vrouw')].unique()
```

```
#pivot table op basis van gender
piv = pd.pivot_table(df, index=['partij', 'gender'], aggfunc='count').dropna()
piv = piv.drop(['achternaam', 'initials', 'voornaam', 'woonplaats'], axis=1).rename(columns={'id': 'aantal'})
```

```
genders = kandidaten.groupby('Partij')['gender'].sum()
dfg = genders.to_frame()
dfg = dfg.reset_index()
compleet = peilingwijzer.merge(dfg, on='Partij')
```

```
kandidaten.gender.replace({'male': 'm', 'female': 'v'}, inplace=True)
peilingwijzer = pd.read_excel("peilingwijzer.xlsx")
#partijnamen goed zetten
peilingwijzer.Partij[2] = 'Partij van de Arbeid (P.v.d.A.)'
peilingwijzer.Partij[3] = 'PVV (Partij voor de Vrijheid)'
peilingwijzer.Partij[4] = 'SP (Socialistische Partij)'
peilingwijzer.Partij[6] = 'Democraten 66 (D66)'
peilingwijzer.Partij[7] = 'ChristenUnie'
peilingwijzer.Partij[8] = 'GROENLINKS'
peilingwijzer.Partij[9] = 'Staatkundig Gereformeerde Partij (SGP)'
peilingwijzer.Partij[10] = 'Partij voor de Dieren'
peilingwijzer.Partij[12] = 'VVD (Voor Nederland)'
peilingwijzer.Partij[13] = 'DENK'
peilingwijzer.Partij[14] = 'Forum voor Democratie'
peilingwijzer.Partij[15] = 'Piratenpartij'
```

```
wrouwenlaag = {}
mannelaag = {}
wrouwenhoog = {}
mannelaaghoog = {}
wrouwen = {}
mannen = {}
for item in compleet.iterrows():
    wrouwenlaag.append(item[1]['gender'] + item[1]['ZetelsLaag'].count("v"))
    wrouwenhoog.append(item[1]['gender'] + item[1]['ZetelsHoog'].count("v"))
    mannelaag.append(item[1]['gender'] + item[1]['ZetelsLaag'].count("m"))
    mannelaaghoog.append(item[1]['gender'] + item[1]['ZetelsHoog'].count("m"))
    wrouwen.append(item[1]['gender'] + item[1]['ZetelsLaag'].count("v"))
    mannen.append(item[1]['gender'] + item[1]['ZetelsLaag'].count("m"))
compleet['wrouwenlaag'] = wrouwenlaag
compleet['wrouwenhoog'] = wrouwenhoog
compleet['mannelaag'] = mannelaag
compleet['mannelaaghoog'] = mannelaaghoog
compleet['mannen'] = mannen
compleet['vrouwen'] = wrouwen
```

```
ding = compleet[['Partij', 'mannen', 'vrouwen']]
ding[ding.mannen > 0 | (ding.vrouwen > 0)]
RATIO + TOTAAL
dfl["total"] = (dfl["m"] + dfl["v"])
dfl["ratio"] = dfl["m"] / (dfl["m"] + dfl["v"])
Groupby
g = amsterdam.groupby(['partij', 'gender'])
g.size()
```

① kandidaten = pd.DataFrame (data)
②
③
④
⑤
⑥

#maak dataframe van lijst met dictionary, maak van id een nummer
import numpy as np
df = pd.DataFrame.from_dict(lst)
df['id'] = df['id'].str.get(0)
df.id = pd.to_numeric(df.id, errors="coerce").fillna(0).astype(np.int64)
print df.shape
df.head()
#df.partij.value_counts() #aantal kandidaten lijst per partij
#read peilingwijzer
peil = pd.read_excel('peilingwijzer.xlsx')
print peil.shape
vr = df2[df2['partij'] == 'PVV (Partij voor de Vrijheid)']
vr = vr[vr['nummer'] < 26]
vr.gender.value_counts()
%ls -lh "file"
import pandas as pd
import numpy as np
import seaborn as sn
import lxml
from lxml import etree
%matplotlib inline
kind: str
'line': line plot (default)
'bar': vertical bar plot
'barh': horizontal bar plot
'hist': histogram
'box': boxplot
'kde': Kernel Density Estimation plot
'density': same as 'kde'
'area': area plot
'pie': pie plot
'scatter': scatter plot
'hexbin': hexbin plot
help (pd.DataFrame)
%time
Meerdere files
for file in os.listdir(path):
if not file.endswith('mb') or not os.path.isfile(file):
tree = ET.parse(path + file)


```
#parse xml
tree = lxml.etree.parse("Kandidatenlijsten_TK2017_Amsterdam.eml.xml")
#haal alle elementen uit xml
candidates = tree.findall('//*[urn:oasis:names:tc:evs:schema:eml]Candidate')
affil = tree.findall('//*[urn:oasis:names:tc:evs:schema:eml]Affiliation')
for a in affil:
    candidates = a.findall('//*[urn:oasis:names:tc:evs:schema:eml]Candidate')
    for c in candidates:
        dct = {}
        dct['partij'] = a.find('//*[urn:oasis:names:tc:evs:schema:eml]AffiliationIdentifier//{urn:oasis:names:tc:evs:schema:eml}RegisteredName').text
        dct['initials'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]CandidateFullName//{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}NameLine').text
        dct['voornaam'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]CandidateFullName//{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}FirstName').text
        dct['gender'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]Gender').text
        dct['achternaam'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]CandidateFullName//{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}LastName').text
        dct['id'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]CandidateIdentifier').values()
        dct['woonplaats'] = c.find('//*[urn:oasis:names:tc:evs:schema:eml]QualifyingAddress//{urn:oasis:names:tc:ciq:xsd:schema:xAL:2.0}Locality//{urn:oasis:names:tc:ciq:xsd:schema:xAL:2.0}LocalityName').text
        lst.append(dct)
```

```
#neem specifieke partij
cda = df[df.partij == 'CDA']
#kruistabel van alle partijen en hun gender verdeling
geslacht = pd.crosstab(df.partij, df.gender, margins=True)
geslacht.head()
#man/vrouw verhouding met ratio
geslacht["ratio"] = geslacht["male"] / geslacht["female"]
#ratio = geslacht.apply(np.log2)
#ratio.loc[(ratio["ratio"] > 0.5) & (ratio["ratio"] < 0.5)].sort_values(ascending=False, by="All").head()
geslacht.head().sort_values(ascending=False, by="ratio")
#ratio
#man en vrouw verhouding, stacked graph kieslijsten
new = pd.crosstab(df.partij, df.gender).sort_values(['female'], ascending=True).plot(kind='bar', figsize=(5, 5), stacked=True)
```

```
#percentage man en vrouw
geslacht["avg_F"] = (geslacht["female"] / geslacht["All"]) * 100
geslacht["avg_M"] = (geslacht["male"] / geslacht["All"]) * 100
geslacht.sort_values(['avg_M'], ascending=False).head(10).round(2)
#aantal procent man en aantal procent vrouw gehele kieslijst
geslacht2 = pd.crosstab(df.partij, df.gender, margins=True)
total = np.sum(geslacht2.ix[['PVV (Partij voor de Vrijheid)'], ['male', 'female']].values)
(geslacht2.ix[['PVV (Partij voor de Vrijheid)'], ['male', 'female']].sum(axis=0) / total * 100).plot(kind='pie')
```

```
#partijleden van de vvd die in rotterdam of amsterdam wonen
vvd_ams_rot = df.loc[(df.partij == 'VVD') & ((df.woonplaats == 'Rotterdam') | (df.woonplaats == 'Amsterdam'))]
ams = "Kandidatenlijsten_TK2017_Amsterdam.eml.xml"
tree = ET.parse(ams)
root = tree.getroot()
data = []
for item in tree.iter():
    if item.tag == '{urn:oasis:names:tc:evs:schema:eml}Affiliation':
        partij = item.getchildren()[0][0].text
        for candidate in item.findall('//*[urn:oasis:names:tc:evs:schema:eml]Candidate'):
            id = candidate.getchildren()[0].attrib["id"]
            firstname = candidate.getchildren()[1][0][1].text
            lastname = candidate.getchildren()[1][0][2].text
            initialen = candidate.getchildren()[1][0][0].text
            gender = candidate.getchildren()[2].text
            place = candidate.getchildren()[3][0][0].text
            data.append({"initialen": initialen, "partij": partij, "id": id, "naam": firstname, "achternaam": lastname, "geslacht": gender, "woonplaats": place})
```

```
counter = df.woonplaats.value_counts() #value_counts sorteert ook meteen voor je
counter = counter[:20]
counter.plot(kind='barh', title = '20 meest populaire woonplaatsen voor de kandidaten van Amsterdam')
#PLOT DE ZETELVERDELING VOLGENS DE PEILINGWIJZER
peil[['Partij', 'Zetels']].plot(x = 'Partij', y = 'Zetels', kind='barh', title = 'Zetel verdeling')
#BEREKEN HET VERSCHIL VAN TWEE KOLOMMEN EN VOEG DIT TOE IN EEN EXTRA KOL
peil['VerschilHoogLaag'] = (peil['ZetelsHoog'] - peil['ZetelsLaag'])
peil.head()
#HOEVEEL % VAN TOTAAL KANDIDATEN STAAT OP DE LIJST VOOR SPECIFIEKE PARTIJ?
(df.partij == 'VVD').mean() * 100
#S-HERTOGENBOSCH WORDT NIET GEVONDEN DOOR HET STREEPJE (HYPHEN) DUS MOET JE ZO ZOEKEN ALS JE DEZE WILT VINDEN
df['woonplaats'].loc[df['woonplaats'].str.contains('Hertogenbosch')]
#GEEF EEN LIJST VAN DE PARTIJEN DIE VROUWEN IN HUN KANDIDATEN LIJST HEBBEN
df['partij'].loc[(df['geslacht'] == 'vrouw')].unique()
```

```
#pivot table op basis van gender
piv = pd.pivot_table(df, index=['partij', 'gender'], aggfunc='count').dropna()
piv = piv.drop(['achternaam', 'initials', 'voornaam', 'woonplaats'], axis=1).rename(columns={'id': 'aantal'})
genders = kandidaten.groupby('Partij')['gender'].sum()
dfg = genders.to_frame()
dfg = dfg.reset_index()
kandidaten.gender.replace({'male': 'm', 'female': 'v', inplace=True)
peilingwijzer = pd.read_excel("peilingwijzer.xlsx")
#partijnamen goed zetten
peilingwijzer.Partij[2] = 'Partij van de Arbeid (P.v.d.A.)'
peilingwijzer.Partij[3] = 'PVV (Partij voor de Vrijheid)'
peilingwijzer.Partij[4] = 'SP (Socialistische Partij)'
peilingwijzer.Partij[6] = 'Democraten 66 (D66)'
peilingwijzer.Partij[7] = 'ChristenUnie'
peilingwijzer.Partij[8] = 'GroenLinks'
peilingwijzer.Partij[9] = 'Staatkundig Gereformeerde Partij (SGP)'
peilingwijzer.Partij[10] = 'Partij voor de Dieren'
peilingwijzer.Partij[12] = 'VNL (VoorNederland)'
peilingwijzer.Partij[13] = 'DEK'
peilingwijzer.Partij[14] = 'Forum voor Democratie'
peilingwijzer.Partij[15] = 'Piratenpartij'
peilingwijzer
#min. dft.columns
dfLL = ... JJ.mein() ...
```

```
compleet = peilingwijzer.merge(dfg, on='Partij')
vrouwenlaag = []
mannenlaag = []
vrouwenhoog = []
mannenhoog = []
vrouwen = []
mannen = []
for item in compleet.iterrows():
    vrouwenlaag.append(item[['gender', 'ZetelsLaag']].count("v"))
    vrouwenhoog.append(item[['gender', 'ZetelsHoog']].count("v"))
    mannenlaag.append(item[['gender', 'ZetelsLaag']].count("m"))
    mannenhoog.append(item[['gender', 'ZetelsHoog']].count("m"))
    vrouwen.append(item[['gender', 'ZetelsLaag']].count("v"))
    mannen.append(item[['gender', 'ZetelsLaag']].count("m"))
compleet['vrouwenlaag'] = vrouwenlaag
compleet['vrouwenhoog'] = vrouwenhoog
compleet['mannenlaag'] = mannenlaag
compleet['mannenhoog'] = mannenhoog
compleet['vrouwen'] = vrouwen
compleet['mannen'] = mannen
```

```
#matplotlib inline
ding = compleet[['Partij', 'mannen', 'vrouwen']]
ding[ding.mannen > 0] | (ding.vrouwen > 0)
RATIO & TOTAAL
df1["total"] = (df1["m"] + df1["v"])
df1["ratio"] = df1["m"] / (df1["m"] + df1["v"])
Groupby
g = amsterdam.groupby(['partij', 'gender'])
g.size()
#partijnamen goed zetten
peilingwijzer.Partij[2] = 'Partij van de Arbeid (P.v.d.A.)'
peilingwijzer.Partij[3] = 'PVV (Partij voor de Vrijheid)'
peilingwijzer.Partij[4] = 'SP (Socialistische Partij)'
peilingwijzer.Partij[6] = 'Democraten 66 (D66)'
peilingwijzer.Partij[7] = 'ChristenUnie'
peilingwijzer.Partij[8] = 'GroenLinks'
peilingwijzer.Partij[9] = 'Staatkundig Gereformeerde Partij (SGP)'
peilingwijzer.Partij[10] = 'Partij voor de Dieren'
peilingwijzer.Partij[12] = 'VNL (VoorNederland)'
peilingwijzer.Partij[13] = 'DEK'
peilingwijzer.Partij[14] = 'Forum voor Democratie'
peilingwijzer.Partij[15] = 'Piratenpartij'
peilingwijzer
```

```
peilingwijzer
#min. dft.columns
dfLL = ... JJ.mein() ...
```

```
peilingwijzer
#min. dft.columns
dfLL = ... JJ.mein() ...
```

read peilingwijzer.xlsx
peil = pd.read_excel('peilingwijzer.xlsx')
print peil.shape

vr = df2[df2['partij'] == 'PVV (Partij voor de Vrijheid)']
vr = vr[vr['nummer'] < 26]
vr.gender.value_counts()

Y-Is - 1h "file"

import pandas as pd
import numpy as np
import seaborn as sns
import lxml
from lxml import etree
%matplotlib inline

kind: str
'line': line plot (default)
'bar': vertical bar plot
'barh': horizontal bar plot
'hist': histogram
'box': boxplot
'kde': Kernel Density Estimation plot
'density': same as 'kde'
'area': area plot
'pie': pie plot
'scatter': scatter plot
'hexbin': hexbin plot

help(pd.DataFrame)
%time

Maak deze files
for file in os.listdir(path):
if not file.endswith('xml') or not os.path.isfile(file):
tree = ET.parse(path + file)

import lxml.etree.ElementTree as ET
pd.read("...")

1 kandidaten = pd.DataFrame(data)

2

3

4

5

6

test "Femal". sort_values(ascending=True).dropna
.plot(kind='barh')
Calculum 7. apply(lambda x: x.count())

<pre> import pandas as pd import numpy as np import matplotlib.pyplot as plt import matplotlib as mpl import os from lxml import etree #open Bestand csv & xlsx: open kandidatenlijstentk2017.csv en ijfers_Peilingwijzer.xlsx in aparte Dataframes. if kandidaten = pd.read_csv('bestandnaam') if peiling = pd.read_excel('bestandnaam') open bestand excel: ijfers_peilingen = pd.read_excel('Bestandnaam') ijfers_neelinen.html </pre>	<pre> #hoeveel mannen en vrouwen zijn in kandidatenlijst van Amsterdam? mannen = 0 vrouwen = 0 for event, elem in ET.iterparse('Kandidatenlijsten_TK2017_Amsterdam.eml.xml'): tag = elem.tag value = elem.text if tag == '(urn:oasis:names:tc:evs:schema:em)Gender': if value == 'female': vrouwen += 1 else: mannen += 1 elem.clear() print('aantal mannen: ', mannen) print('aantal vrouwen: ', vrouwen) print('totaal: ', (mannen + vrouwen)) </pre>	<p>#Hoeveel mannen en vrouwen zijn in kandidatenlijst van Amsterdam?</p> <p>STATISTICS</p> <p>These can all be applied to a series as well.</p> <p>df.describe() - Summary statistics for numerical columns df.mean() - Return t</p> <p>mean of all columns df.corr() - finds the correlation between columns in a DataFrame.</p> <p>df.count() - counts the number of non-null values in each DataFrame</p> <p>column df.max() - finds the highest value in each column. df.min() - finds the lowest value in each column. df.median() - finds the median of each column.</p> <p>df.std() - finds the standard deviation of each column.</p>
<pre> Haal alle variabele zien binnen het dataframe: ijfers_peilingen.describe() </pre>	<pre> #Maak een dataframe met voornaam, initialen, achternaam, man/vrouw, woonplaats partij = [] voornaam = [] initialen = [] achternaam = [] gender = [] woonplaats = [] for event, elem in ET.iterparse('Kandidatenlijsten_TK2017_Amsterdam.eml.xml'): tag = elem.tag value = elem.text if tag == '(urn:oasis:names:tc:evs:schema:em)RegisteredName': (partij) partijvalue = value if tag == '(urn:oasis:names:tc:ciq:xdschema:xNL:2.0)FirstName': partij.append(partijvalue) voornaam.append(value) if tag == '(urn:oasis:names:tc:ciq:xdschema:xNL:2.0)NameLine': initialen.append(value) if tag == '(urn:oasis:names:tc:ciq:xdschema:xNL:2.0)LastName': achternaam.append(value) if tag == '(urn:oasis:names:tc:evs:schema:em)Gender': if value == 'male': gender_val = 'man' else: gender_val = 'vrouw' gender.append(gender_val) if tag == '(urn:oasis:names:tc:ciq:xdschema:xNL:2.0)LocalityName': woonplaats.append(value) </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> Haal een variabele uit het dataframe uit een kolom - voorbeeld 'Verschil' 'Partij' 'PVV' getal = ijfers_peilingen['Verschil'].loc[(cijfers_peilingen['Partij'] == 'PVV')].values[0] getal </pre>	<pre> df = pd.DataFrame({'partij':partij, 'voornaam':voornaam, 'initialen':initialen, 'achternaam':achternaam, 'geslacht':gender, 'woonplaats':woonplaats}) df.head(5) </pre>	<p>df.count() - counts the number of non-null values in each DataFrame</p> <p>column df.max() - finds the highest value in each column. df.min() - finds the lowest value in each column. df.median() - finds the median of each column.</p> <p>df.std() - finds the standard deviation of each column.</p>
<pre> Haal verschillende kolommen eruit en haal de indexcijfers weg en zet daar een andere kolom voor terug. ijfers_peilingen[['Zetels', 'Datum', 'Percentage']].set_index('Zetels').head() </pre>	<pre> df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels', df.join.head(5) </pre>	<p>df.count() - counts the number of non-null values in each DataFrame</p> <p>column df.max() - finds the highest value in each column. df.min() - finds the lowest value in each column. df.median() - finds the median of each column.</p> <p>df.std() - finds the standard deviation of each column.</p>
<pre> Velke partijen hebben boven de 25 zetels of onder ijvoorbeeld (<)? etels = cijfers_peilingen.loc[cijfers_peilingen['Zetels'] > 25] etels </pre>	<pre> Print het aantal mannen en het aantal vrouwen? print (len(df[df.geslacht=='man'])) print (len(df[df.geslacht=='vrouw'])) </pre>	<p>df.count() - counts the number of non-null values in each DataFrame</p> <p>column df.max() - finds the highest value in each column. df.min() - finds the lowest value in each column. df.median() - finds the median of each column.</p> <p>df.std() - finds the standard deviation of each column.</p>
<pre> Sorteer dataframe op bepaalde column - ascending = 'alse (boven naar onder) ijfers_peilingen_new = cijfers_peiling.sort_values(by='Zetels', scending=False) ijfers_peiling_new.head() </pre>	<pre> Print een lijst van alle unieke partijen op de kandidatenlijst? print (df['partij'].unique()) len(set(df['partij'])) </pre>	<p>Hoeveel partijen heb je op de kandidatenlijst van Amsterdam? En hoeveel partijen worden genoemd door de peilingwijzer?</p> <p>print len(set(df.peiling)), 'genoemd door de peilingwijzer en', len(set(df.kandidaten)), 'op de kandidatenlijst van Amsterdam'</p>
<pre> Hoeveel partijen zitten er in de dataset van de PVV? en(cijfers_peilingen[cijfers_peilingen.Partij=='PVV']) </pre>	<pre> Hoeveel kandidaten per partij en maak een plot? kandidaten_per_partij = df['partij'].value_counts() kandidaten_per_partij.plot(kind='bar', title='kandidaten per partij') </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> Hoe lang is de kolom partij? en(cijfers_peilingen.Partij) </pre>	<pre> Bepaalde partijen met een bepaald geslacht df.loc[(df['partij'] == 'VVD') & (df['geslacht'] == 'man')].head() </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> Hoe veel unieke/dubbele partijen zitten er in deze kolom? en(set(cijfers_peilingen.Partij)) </pre>	<pre> Hoeveel procent is man van de VVD? totaal_aantal_mannen mannen_VVD = len(df.loc[(df['partij'] == 'VVD') & (df['geslacht'] == 'man')]) totaal_aantal_vrouwen vrouwen_VVD = len(df.loc[(df['partij'] == 'VVD') & (df['geslacht'] == 'vrouw')]) totaal_VVD = len(df[df.partij == 'VVD']) </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> Sorteer de partijen met het hoogste percentage ? ijfers_peilingen = cijfers_peilingen.set_index ('partij') </pre>	<pre> aantal_procenten_man = (mannen_VVD * 100) / totaal_VVD print (aantal_procenten_man) </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> percentage_new = cijfers_peilingen.sort_values(by='Percentage', scending=False) percentage_new.head() </pre>	<pre> Plot de verhouding man vrouw van de VVD? verhouding = df['geslacht'].loc[(df['partij'] == 'VVD')].value_counts() verhouding.plot(title='Verhouding man vrouw VVD') </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> lot de partij en het percentage - bar percentage_new['Percentage'].plot(kind='bar', title = 'Hoogste percentage') </pre>	<pre> Plot de verhouding man vrouw van de VVD? verhouding = df['geslacht'].loc[(df['partij'] == 'VVD')].value_counts() verhouding.plot(kind='bar', title='Verhouding man vrouw VVD') </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> lot de partij en het percentage - pie percentage_new['Percentage'].plot(kind='pie', title = 'Hoogste percentage plot') </pre>	<pre> Plot de verhouding man vrouw van de VVD? verhouding = df['geslacht'].loc[(df['partij'] == 'VVD')].value_counts() verhouding.plot(kind='pie', title='Verhouding man vrouw VVD') </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> lot de partij en het percentage - grapi percentage_new['Percentage'].plot(title = 'Hoogste Percentage slot') </pre>	<pre> Zoek een Kandidaat op uit een bepaald woonplaats df.loc[(df['woonplaats'] == 'Amsterdam')] </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> open een excel bestand en laat de eerste 5 zien? results_longitudinal = pd.read_excel('Results_Longitudinal.xlsx') results_longitudinal.head() </pre>	<pre> Kandidaat uit Amsterdam die in DENK zit df.loc[(df.woonplaats == 'Amsterdam') & (df.partij == 'DENK')] </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> XML Files openen. import xml.etree.ElementTree as ET tree_new = ET.parse('bestandnaam') root = tree_new.getroot() </pre>	<pre> Welke van de kandidaten wonen in Rotterdam, Amsterdam en Wassenaar? df.loc[df['woonplaats'].isin(['Rotterdam', 'Amsterdam', 'Wassenaar'])].head() je kan dit ook doen met 'partij' en dan bijvoorbeeld 'DENK', 'VVD'. </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> XML Files openen. root.tag </pre>	<pre> Geef een lijst met partijen die een 'naam' (Peter) als kandidaat hebben? df['partij'].loc[(df['voornaam'] == 'naam').unique() </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> XML Files openen. root.attrib </pre>	<pre> Zoek naar kandidaten waar achternaam met een K begint df[['voornaam', 'achternaam', 'partij']].loc[(df['achternaam'].str.startswith('K'))].head() </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>
<pre> XML Files openen. for tag in tree_new.iter(): print (tag) </pre>	<pre> Stel dat de VVD 20 zetels haalt, hoeveel mannelijke VVD'ers komen er dan in de Kamer? Maak een dataframe met X aantal namen van de kandidatenlijst van een specifieke partij df_vvd = df.loc[(df['partij'] == 'VVD')].head(peiling['Zetels']).loc[(peiling['Partij'] == 'VVD')].values.item() df_vvd </pre>	<p>df.join.update(df.join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels',</p>

index == 'v' // Vrouwen