

```
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
```

```
#downloaden bestand
#lwget https://www.kiesraad.nl/binaries/kiesraad/documenten/rapporten/2017/2/definitieve-ki
url = urllib2.urlopen('https://www.kiesraad.nl/binaries/kiesraad/documenten/rapporten/2017/
with open('map.zip', 'wb') as code:
code.write(uri.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_lijs = []
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_lijs.append(int(event.attrib.values()[0]))
```

```
if tag == '{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}NameLine':
initialen.append(value)
```

```
if tag == '{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}FirstName':
voornaam.append(value)
partij.append(partij_val)
```

```
if tag == '{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}LastName':
achternaam.append(value)
prefix.append(prefix_val)
```

```
if tag == '{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}NamePrefix':
prefix_val = value
else:
prefix_val = ''
```

```
if tag == '{urn:oasis:names:tc:ciq:xsd:schema:xNL: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')
```

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

```
df = df[['partij', 'positie_lijs', '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
--------	-------	------------	----------------	----------------	--------	------------	------------

```
#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

XML

```
#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 ZOEKEN 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-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 = 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_lijs'] <= (pw_df['zetels'].loc[pw_df['Partij'] == a]).values[0]) &
(df['gender'] == 'vrouw')]))
mannen.append(len(df.loc[(df['partij'] == a) &
(df['positie_lijs'] <= (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 WOONPLA
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 UNIEKE 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
#ALLES 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 éenen 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())
```

```

from bs4 import BeautifulSoup
import pandas as pd
# import re
from collections import Counter
import numpy as np
from _future_ import division
import pprint
%matplotlib inline
import matplotlib
import matplotlib.pyplot as plt
matplotlib.style.use('ggplot')
plt.rcParams['figure.figsize'] = (12, 5)

# hoeveel % van VVD is vrouw
len(KandidatenlijstFrame[KandidatenlijstFrame['Partij'] == 'VVD'].values)/len(KandidatenlijstFrame)*100

# dit is een test voor de loop in volgende code block
AantalVrouwenVanEenPartijTest = KandidatenlijstFrame[(KandidatenlijstFrame['Partij'] == 'Partij voor de Dieren')]
AantalVrouwenVanEenPartijTestRangeTest =
AantalVrouwenVanEenPartijTest[(AantalVrouwenVanEenPartijTest['Notering'] < 20)]
AantalVrouwenVanEenPartijTestRangeTestAantalVrouwen =
AantalVrouwenVanEenPartijTestRangeTest[(AantalVrouwenVanEenPartijTestRangeTest['Geslacht'] == 'female').values]
len(AantalVrouwenVanEenPartijTestRangeTestAantalVrouwen)
# Dit is dus de loop waar in vorige code block wordt gerefereerd
PartijDict = {}
for p in PartijenLijst:
    AantalVrouwenVanEenPartijTest = KandidatenlijstFrame
    [(KandidatenlijstFrame['Partij'] == p)
    AantalVrouwenDict = {}
    for i in range(50):
        AantalVrouwenVanEenPartijTestRangeTest =
AantalVrouwenVanEenPartijTest[(AantalVrouwenVanEenPartijTest['Notering'] < i)]
AantalVrouwenVanEenPartijTestRangeTestAantalVrouwen =
len(AantalVrouwenVanEenPartijTestRangeTest[(AantalVrouwenVanEenPartijTestRangeTest['Geslacht'] == 'female').values])
AantalVrouwenDict[i] = AantalVrouwenVanEenPartijTestRangeTestAantalVrouwen
PartijDict[p] = AantalVrouwenDict
PartijDict
# MAAK DATAFRAME VAN DEZE DICT
PartijDictFrame = pd.DataFrame.from_dict(PartijDict, orient='columns')
# plot aantal vrouwen per partij voor bepaalde partij
for p in PartijenLijst:
    AantalVrouwenNavAantalZetelsFrame =
pd.DataFrame.from_dict(PartijDict[p], orient='index')
AantalVrouwenNavAantalZetelsFrame.plot(kind='bar', title='Aantal Vrouwen
N.a.v. Aantal Zetels Voor '+p).set_ylim(0, 40)

# PEILINGWIJZER METHODES OP GOED DATAFRAME TE KRIJGEN
# LAATSTE PEILINGEN
LaatstePeilingWijzer = pd.read_excel('https://s3.eu-central-1.amazonaws.com/louwerse/Public/Peilingwijzer/Last/Cijfers_Peilingwijzer.xlsx')
PartijNamenXML = {u'PvdA':u'Partij van de Arbeid (P.v.d.A.)', u'PVV':u'PVV (Partij voor de Vrijheid)', u'SP':u'SP (Socialistische Partij)', u'D66':u'Democraten 66 (D66)', u'CU':u'ChristenUnie', u'GL':u'GROENLINKS', u'SGP':u'Staatkundig Gereformeerde Partij (SGP)', u'PvdD':u'Partij voor de Dieren', u'VNL':u'VNL (VoorNederland)', u'Denk':u'DENK', u'FvD':u'Forum voor Democratie', u'PP':u'Piratenpartij'}
LaatstePeilingWijzer = LaatstePeilingWijzer.set_index('Partij')
LaatstePeilingWijzer.rename(PartijNamenXML, inplace=True)
LaatstePeilingWijzer = LaatstePeilingWijzer.reset_index()
LaatstePeilingWijzer

# ALLE PEILINGEN
AllePeilingWijzers = pd.read_excel('https://s3.eu-central-1.amazonaws.com/louwerse/Public/Peilingwijzer/Last/Results_Longitudinal.xlsx')
AllePeilingWijzers = AllePeilingWijzers[['VVD', 'PvdA', 'PVV', 'SP', 'CD', 'D66', 'CU', 'GL', 'SGP', 'PvdD', '50PLUS', 'VNL', 'Denk', 'FvD', 'PP']] * 150
AllePeilingWijzers = AllePeilingWijzers.round()
# AllePeilingWijzers.plot(figsize=(18,10)).set_ylim(0,45)
AllePeilingWijzers = AllePeilingWijzers.rename(columns={u'PvdA':u'Partij van de Arbeid (P.v.d.A.)', u'PVV':u'PVV (Partij voor de Vrijheid)', u'SP':u'SP (Socialistische Partij)', u'D66':u'Democraten 66 (D66)', u'CU':u'ChristenUnie', u'GL':u'GROENLINKS', u'SGP':u'Staatkundig Gereformeerde Partij (SGP)', u'PvdD':u'Partij voor de Dieren', u'VNL':u'VNL (VoorNederland)', u'Denk':u'DENK', u'FvD':u'Forum voor Democratie', u'PP':u'Piratenpartij'})
AllePeilingWijzers = AllePeilingWijzers.reset_index()
AllePeilingWijzers = AllePeilingWijzers.rename(columns={'index':'date'})
AllePeilingWijzers = AllePeilingWijzers.set_index('date')
AllePeilingWijzers.head()

# ALLE PEILINGEN SECURDER
alle=pd.read_csv('https://d1b1gq97ifurz.cloudfront.net/Public/Peilingwijzer/Last/Results_DyGraphs_Seats.csv', index_col='Date')
#print (alle.shape)
alle =alle.applymap(lambda s: int(s.split(';')[1])) #if you only want to keep the middle number
#now translate to our names
alle=alle.T
alle.rename(PartijNamenXML, inplace=True)
alle = alle.T
alle

# OPEN FILE
Kandidatenlijsten = BeautifulSoup(open('Kandidatenlijsten_EML_TK2017-20170213/Kandidatenlijsten_TK2017_Amsterdam.eml.xml').read(), 'lxml')

# Merge de dataframes samen
PartijGeslachtDF =
pd.DataFrame.from_dict(PartijGeslachtDict, orient='index')
PartijGeslachtDF =
PartijGeslachtDF.stack().reset_index()
PartijGeslachtDF.columns = ['Partij', 'Notering', 'Geslacht']

KandidatenlijstenFrame =
PartijKandidaatDF.merge(PartijGeslachtDF)
.merge(PartijWoonplaatsDF)
KandidatenlijstenFrame['Notering'] =
KandidatenlijstenFrame[['Notering']] + 1
KandidatenlijstenFrame

# Opmaken tweede kamer op basis van laatste peiling
AllePeilingWijzersLaatstePeiling = AllePeilingWijzers.tail(1)
VerdelingLijst = []
for partij in AllePeilingWijzersLaatstePeiling.columns:
    if len(AllePeilingWijzersLaatstePeiling[AllePeilingWijzersLaatstePeiling[partij]>0]):
        Kandidaten = int(AllePeilingWijzersLaatstePeiling[partij].item())
        PartijFrame = AllePeilingWijzersLaatstePeiling[partij]
        KandidatenPerPartij = KandidatenlijstenFrame[KandidatenlijstenFrame['Partij'] == partij][0:Kandidaten]
        VerdelingLijst.append(KandidatenPerPartij)
tweede_kamer = pd.concat(VerdelingLijst)
tweede_kamer.reset_index(inplace=True)
tweede_kamer

# OOM DATAFRAME TE MAKEN EN PLOTTEN HOEVEEL VROUWEN OVER TIJD IN DE KAMER MAAK COPY
AllePeilingWijzerscopy = AllePeilingWijzers.copy()
AllePeilingWijzerscopy2 = AllePeilingWijzerscopy.copy()
AllePeilingWijzerscopy2

# AANTAL VROUWEN IN DE TWEEDE KAMER OVER DE TIJD (PEILINGWIJZER)
for partij in AllePeilingWijzerscopy.columns:
    for i in AllePeilingWijzerscopy[partij]:
        vrouwen = PartijDict[partij][i]
        AllePeilingWijzerscopy2.loc[AllePeilingWijzerscopy[partij] == i, partij] = vrouwen
AllePeilingWijzerscopy2

# PLOT TOTALE TIJD AAN VROUWEN OVER DE TIJD IN DE KAMER
AllePeilingWijzersTotalen = AllePeilingWijzerscopy2.copy()
AllePeilingWijzersTotalen['Totaal'] = AllePeilingWijzersTotalen.sum(axis=1)
AllePeilingWijzersTotalen['Totaal percentage'] = AllePeilingWijzersTotalen['Totaal']/150
AllePeilingWijzersTotalen

# hoeveel vrouwen in 2e kamer op basis van laatste peiling?
Aantal_vrouwen = AllePeilingWijzersTotalen.tail(1)
print ('Het aantal vrouwen volgens de peiling van 14 Maart =', int(Aantal_vrouwen['Totaal'].item()))

# PLOT PERCENTAGE TOTALE TIJD OVER DE TIJD IN DE KAMER
AllePeilingWijzersTotalen['Totaal percentage'].plot(figsize=(18,7))
# AANTAL VROUWEN LAATSTE PEILING 14/03
AllePeilingWijzers14Maart = AllePeilingWijzers.tail(1)
AllePeilingWijzers14MaartT = AllePeilingWijzers14Maart.T
AllePeilingWijzers14MaartT

# AANTAL VROUWEN LAATSTE PEILING 14/03
AllePeilingWijzerscopy14Maart =
AllePeilingWijzerscopy2.tail(1)
AllePeilingWijzerscopy14MaartT =
AllePeilingWijzerscopy14Maart.T
AllePeilingWijzerscopy14MaartT

# PLOT MAN/VROUW VERDELING PER PARTIJ
t = pd.merge(AllePeilingWijzerscopy14MaartT, AllePeilingWijzers14MaartT, left_index=True, right_index=True, how='inner')
t.columns = ['Vrouwen', 'Zetels']
t['Manner'] = t['Zetels'] - t['Vrouwen']
t[['Manner', 'Zetels']].plot(kind='bar')

# PERCENTAGE VROUWEN
t['Percentage Vrouwen'] = (t['Vrouwen']/t['Zetels'])*100
t[['Percentage Vrouwen']].plot(kind='bar')
t

ascending=[True,False] - Sort values by col1 in ascending order then col2 in descending order
df.groupby(col1) - 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

```

pd.read\_html

```
#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:ciq:xds:schema:xAL:2.0:NameLine]').text
        dct["voornaam"] = c.find('//*[urn:oasis:names:tc:evs:schema:eml:CandidateFullName]//*[urn:oasis:names:tc:ciq:xds:schema:xAL: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:xds:schema:xAL:2.0:LastName]').text
        dct["id"] = c.find('//*[urn:oasis:names:tc:evs:schema:eml:CandidateIdentifier]').text
        dct["woonplaats"] = c.find('//*[urn:oasis:names:tc:evs:schema:eml:QualifyingAddress]//*[urn:oasis:names:tc:ciq:xds:schema:xAL:2.0:Locality]//*[urn:oasis:names:tc:ciq:xds:schema:xAL:2.0:LocalityName]').text
        lst.append(dct)
```

```
#hoeveel unieke partijen
cda = df[df.partij == 'CDA']
len(df.partij.unique())

#ruis tabel 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="ratio").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
            initials = candidate.getchildren()[1][0][0].text
            gender = candidate.getchildren()[2].text
            place = candidate.getchildren()[3][0][0].text
            data.append({'initials':initials, "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')

#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')
```

```
#VIND ALLE KANDIDATEN DIE IN EN VAN DE VOLGENDE STEDEN WONEN: BREDA, TILBURG, S-HERTOGENBOSCH, EINDHOVEN EN PL
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 STREEPJKE (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()
complete = peilingwijzer.merge(dfg, on="Partij")
```

```
ding = complete[["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()
```

```
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 (Voortvanderland)'
peilingwijzer.Partij[13] = 'DENK'
peilingwijzer.Partij[14] = 'Forum voor Democratie'
peilingwijzer.Partij[15] = 'Piratenpartij'
peilingwijzer
```

import xml.etree.ElementTree as ET

1

3

4

5

```
bestand = open(filename).read()
soup = bs.BeautifulSoup(bestand, 'xml')
```

# Data Science Cheat Sheet

Pandas

## KEY

We'll use shorthand in this cheat sheet  
df - A pandas DataFrame object  
s - A pandas Series object

## IMPORTS

Import these to start  
import pandas as pd  
import numpy as np

```
x = {}
for item in soup.findAll('aff.'):
    x[item.find('regn.').text] =
    Ct.text for t in item.findAll('gend')]
print Counter(x['vvd'])[0:25])
```

## 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('1900/1/30', 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

## SELECTION

df[col] - Return column with label col as Series  
df[[col1, col2]] - Return Columns as a new DataFrame  
s.iloc[0] - selection by position  
s.loc[0] - selection by index  
df.iloc[0, :] - first row  
df.iloc[0,0] - first element of first column

## DATA CLEANING

df.columns = ['a', 'b', 'c'] - Rename columns  
pd.isnull() - Checks for null Values, Returns Boolean Array  
pd.notnull() - Opposite of s.isnull()  
df.dropna() - Drop all rows that contain null values  
df.dropna(axis=1) - Drop all columns that contain null values  
df.dropna(axis=1, thresh=n) - Drop all rows have have less than n non null values  
df.fillna(x) - Replace all null values with x  
s.fillna(s.mean()) - Replace all null values with the mean (mean can be replaced with almost any function from the statistics section)  
s.astype(float) - Convert the datatype of the series to float  
s.replace(1, 'one') - Replace all values equal to 1 with 'one'  
s.replace([1,3], ['one', 'three']) - Replace all 1 with 'one' and 3 with 'three'  
df.rename(columns=lambda x: x + 1) - mass renaming of columns  
df.rename(columns={'old\_name': 'new\_name'}) - selective renaming  
df.set\_index('column\_one') - change the index  
df.rename(index=lambda x: x + 1) - mass renaming of index

## FILTER, SORT, & GROUPBY

df[df[col] > 0.5] - Rows where the col column is greater than 0.5  
df[(df[col] > 0.5) & (df[col] < 0.7)] - Rows where 0.7 > col > 0.5  
df.sort\_values(col1) - Sort values by col1 in ascending order  
df.sort\_values(col2, ascending=False) - Sort values by col2 in descending order  
df.sort\_values([col1, col2],

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 1 group  
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 col1 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.

```
..future.. division
collections Counter
bs4 as bs
```

```
.value_counts()
plot.barh()
pd.crosstab(col1, col2, margins=True)
df[col].str.contains('string')
```

```
import seaborn as sn
%matplotlib inline
from IPython import pyplot as plt
plt.figure(figsize=(x,y))
plt.bar(y=legend, y_lim=(10,100), title)
```

```

# LIJST AAN KANDIDATEN INLEZEN XML
import bs4 as bs
bestand = open('Kandidatenlijsten_TK2017_Amsterdam.eml2.xml').read()
soup = bs.BeautifulSoup(bestand, 'xml')

# NAMEN VAN DE PARTIJEN
party_names = []
for x in soup.findAll('RegisteredName'):
    party_names.append(x.text)
print party_names

# PLOT HOEVEEL MENSEN ANNE HETEN PER PARTIJ
namendict = {}
for partij in soup('Affiliation'):
    namendict[partij.find('RegisteredName').text] = \
        [name.text for name in partij.findAll('FirstName')].count('Anne')
y = pd.DataFrame.from_dict(namendict, orient='index')
y.plot(kind='barh')

# DICT PARTIJNAAM: [male, female, male]
manvrouwdict = {}
for p in soup('Affiliation'):
    manvrouwdict[p.find('RegisteredName').text] = [g.text for g in p.findAll('Gender')]

# Stel dat de PVV 25 zetels haalt, hoeveel vrouwelijke PVV'ers komen er dan in de Kamer?
from collections import Counter
print Counter(manvrouwdict['PVV (Partij voor de Vrijheid)'][:25])['female']

# PLOT VROUWEN VAN PVV IN DE KAMER MET N GEWONNEN ZETELS
vrouwenperzetel = {p: {N: Counter(manvrouwdict[p][:N])['female']} \
    for N in range(50)} for p in manvrouwdict}
pd.DataFrame.from_dict(vrouwenperzetel['PVV (Partij voor de Vrijheid)'], orient='index') \
    .plot(kind='bar', title='vrouwen van PVV in de kamer met N aantal zetels');

# DATAFRAME AANTAL MENSEN OP DE LIJST, PERCENTAGE MAN/VROUW
vrouwen = {}
mannen = {}
for x in manvrouwdict:
    manofvrouw = manvrouwdict[x]
    vrouwen[x] = manofvrouw.count('female')
    mannen[x] = manofvrouw.count('male')
df = pd.DataFrame.from_dict(vrouwen, orient='index')
df2 = pd.DataFrame.from_dict(mannen, orient='index')
df.columns = ['vrouwen']
df2.columns = ['mannen']
new = pd.concat([df, df2], axis=1)
new['vrouwenperc'] = new.vrouwen / (new.vrouwen + new.mannen) * 100
new['mannenperc'] = new.mannen / (new.vrouwen + new.mannen) * 100
new.sort_values('vrouwenperc', ascending=False)

# PLOT VROUWENPERCENTAGES PER PARTIJ
new['vrouwenperc'].sort_values().plot(kind='barh')

# MINIMAAL AANTAL MENSEN VAN PARTIJ
min_lijst_lengte = (new.filter(['vrouwen', 'mannen']).sum(axis=1)).min()
print min_liist_lengte

# LEES PEILINGWIJZER IN, VERANDER NAMEN OM SAMEN TE VOEGEN
peilingen = pd.read_excel('Cijfers_Peilingwijzer.xlsx')
translate = {u'50PLUS': u'50PLUS',
             u'CDA': u'CDA',
             u'CU': u'ChristenUnie',
             u'D66': u'Democraten 66 (D66)',
             u'Denk': u'DENK',
             u'Fvd': u'Forum voor Democratie',
             u'GL': u'GROENLINKS',
             u'PVV': u'PVV (Partij voor de Vrijheid)',
             u'PvdA': u'Partij van de Arbeid (P.v.d.A.)',
             u'PvdD': u'Partij voor de Dieren',
             u'SGP': u'Staatkundig Gereformeerde Partij (SGP)',
             u'SP': u'SP (Socialistische Partij)',
             u'VNL': u'VNL (VoorNederland)',
             u'VVD': u'VVD',
             u'PP': u'Piratenpartij'}
peilingen.index = peilingen.Partij
peilingen.rename(translate, inplace=True)
del peilingen['Partij']

# JOIN DE TWEDE DATAFRAMES
peilingen.join(new)

# PEILINGWIJZER DOOR DE TIJD OPGESCHOOND
allepeilingen = pd.read_excel('Results_Longitudinal.xlsx')
x = allepeilingen.T.iloc[:, 3]
x.rename(translate, inplace=True)
x.T.tail(3)

# ANDERE OPTIE
zetelstijd = pd.read_csv('https://d1bjgg97if6urz.cloudfront.net/Public/Peilingwijzer/Last/Results_Dy_
    , index_col='Date')
# check if all add to 150
print (zetelstijd.sum(axis=1) == 150).mean()
# pak alleen de middelste
zetelstijd = zetelstijd.applymap(lambda s: int(s.split(':')[1]))

# WEER NAMEN AANPASSEN
zetelstijd = zetelstijd.T
zetelstijd.rename(translate, inplace=True)
zetelstijd = zetelstijd.T
zetelstijd.tail()

# ALLEEN VROUWEN IN DE KAMER DOOR DE TIJD
for p in zetelstijd.columns:
    zetelstijd[p] = zetelstijd[p].apply(lambda z: vrouwenperzetel[p][z])

# PLOT VROUWEN IN DE KAMER PER PARTIJ DOOR DE TIJD
zetelstijd.plot(figsize=(20,10), title='Aantal vrouwen in de Tweede Kamer door de tijd, per partij');

# TOTAAL AANTAL VROUWEN IN DE KAMER DOOR DE TIJD
zetelstijd.T.sum().plot(figsize=(18,7), title='Aantal vrouwen in de Tweede Kamer door de tijd');

# VERDELING ZETELS VOLGENS PEILINGWIJZER
df_peilingwijzer.sort_values(by='Zetels', ascending=False) \
    .plot('Partij', 'Zetels', kind='bar', title='Verdeling Zetels volgens

```

```

from future import division
# HOEVEEL PROCENT VAN DE STEMMINGEN ZIJN AANGENOMEN?
hoofdelijke_stemming_df = rutte2[rutte2['votetype'] == 'Roll call']
print len(hoofdelijke_stemming_df[hoofdelijke_stemming_df['result'] == 'adopted']) \
    / len(hoofdelijke_stemming_df) * 100

# HOE VAAK STEMDE BARRY MEE?
barry_stemt_voor = hoofdelijke_stemming_df['pro'].str.contains('Barry Madlener')
print 'Barry stemde', len(hoofdelijke_stemming_df[barry_stemt_voor]), 'keer voor'

# KRUISTABEL
pd.crosstab(rutte2['proposaltype'], rutte2['result'], margins=True)

# HOEVEEL STUKKEN ELKE PARTIJ INGEDIEND
im = rutte2[(rutte2['authorparty'] != '') \
    & (rutte2['authorparty'] == rutte2['supporterparties'])]
per_party_counts = im['authorparty'].str[1:-1].value_counts()
per_party_counts.sort_values().plot.barh()

# PLOT PER PARTIJ HOEVEEL PROCENT VD STEMMINGEN PARTIJ VOOR HEEFT GESTEMD
parties = rutte2.columns[11:50]
vote_distribution = rutte2[parties].apply(pd.value_counts)
vote_distribution = vote_distribution.T
vote_distribution.dropna(inplace=True)
vote_distribution.columns = ['contra', 'pro']
vote_distribution['pro_percentage'] = (vote_distribution['pro'] \
    / vote_distribution.sum(axis='columns')) * 100
percentage_sorted = vote_distribution \
    .sort_values(by='pro_percentage', ascending=False)
percentage_sorted.plot.bar(y='pro_percentage', ylim=(0, 100), \
    title='Percentage voor-stemmingen per partij')

# PARTIJEN DIE MINSTENS 1 KEER GESTEMD HEBBEN, HOE VAAK?
im2 = rutte2[rutte2['votetype'] == 'Normal']
party_vote_counts = im2[parties].count()
sorted_vote_counts = party_vote_counts[party_vote_counts > 0] \
    .sort_values(ascending=False)

# WANNEER ONTSTOND DENK?
denk_index = parties[18]
denk_is_nan = np.isnan(rutte2[denk_index])
# All items where DENK voted
denk_votes = rutte2[~denk_is_nan]
denk_votes['date'].min()

# HOE VAAK STEMMEN PARTIJEN MEE MET PVV?
pvv_voted_pro = rutte2[rutte2['PVV'] == 1]
other_parties = pvv_voted_pro[parties].dropna(axis='columns')
other_party_votes = other_parties.sum()
relative_votes = other_party_votes / len(pvv_voted_pro)
relative_votes.sort_values(ascending=False)

```

#### FILTER, SORT, & GROUPBY

```

df[df[col] > 0.5] - Rows where the col column is greater than 0.5
df[(df[col] > 0.5) & (df[col] < 0.7)] - Rows where 0.7 > col > 0.5
df.sort_values(col1) - Sort values by col1 in ascending order
df.sort_values(col2, ascending=False) - Sort values by col2 in descending order
df.sort_values([col1, col2],

```

```

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

```

#### 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

```

#### SELECTION

```

df[col] - Return column with label col as Series
df[[col1, col2]] - Return Columns as a new DataFrame
s.iloc[0] - selection by position
s.loc[0] - selection by index
df.iloc[0, :] - first row
df.iloc[0, 0] - first element of first column

```

#### DATA CLEANING

```

df.columns = ['a', 'b', 'c'] - Rename columns
pd.isnull() - Checks for null Values, Returns Boolean Array
pd.notnull() - Opposite of s.isnull()
df.dropna() - Drop all rows that contain null values
df.dropna(axis=1) - Drop all columns that contain null values
df.dropna(axis=0, thresh=n) - Drop all rows have less than n non null values
df.fillna(x) - Replace all null values with x
s.fillna(s.mean()) - Replace all null values with the mean (mean can be replaced with almost any function from the statistics section)
s.astype(float) - Convert the datatype of the series to float
s.replace(1, 'one') - Replace all values equal to 1 with 'one'
s.replace([1, 3], ['one', 'three']) - Replace all 1 with 'one' and 3 with 'three'
df.rename(columns=lambda x: x + 1) - mass renaming of columns
df.rename(columns={'old_name': 'new_name'}) - selective renaming
df.set_index('column_one') - change the index
df.rename(index=lambda x: x + 1) - mass renaming of index

```

#### 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

```

```

import pandas as pd
import os
from lxml import etree
import xml.etree.ElementTree as ET
%matplotlib inline
import pandas as pd
from lxml import etree
from math import sqrt
import matplotlib.pyplot as plt

```

```

peiling = pd.read_excel('Cijfers_Peilingwijzer.xlsx')
peiling.head()
results = pd.read_excel('Results_Longitudinal.xlsx')
results.head(2)
print peiling.describe()
print peiling.shape

```

```

#ALLE VALUES OPTELLEN
print peiling.sum()
#OPTELLEN SPECIFIEKE KOLOM
peiling['PercentageLaag'].sum()

```

```

#VERSCHIL ZETELSHOOG, ZETELSLAAG IN NIEUWE
KOLOM
peiling['Verschil'] = (peiling['ZetelsHoog'] -
(peiling['ZetelsLaag']))
peiling.head(5)

```

```

#KOLOMEN UIT DATAFRAME EN INDEXCIJFERS WEG
peiling[['Partij', 'Datum',
'Percentage']].set_index('Partij').head(2)

```

```

#WELKE PARTIJ BOVEN 20 ZETELS?
zetels = peiling.loc[peiling['Zetels'] > 20]
zetels

```

```

#PLOT HOEVEEL ZETELS PER PARTIJ
df_peiling.sort_values(by='Zetels', ascending =
False).plot('Partij', 'Zetels', kind = 'bar', title =
'Hoeveelheid zetels per partij')

```

```

#GETAL UIT DATAFRAME
getal = peiling['Verschil'].loc[peiling['Partij'] ==
'CDA'].values[0]
getal

```

```

#SORTEER DATAFRAME OP BEPAALDE KOLOM
peiling1 = peiling.sort_values(by='Zetels',
ascending=False)
peiling1.head(2)

```

```

len(peiling[peiling.Partij=='VVD'])

```

```

#HOELANG IS KOLOM 'PARTIJ'?
print len(peiling.Partij)
#HOEVEEL UNIEKE PARTIEN ZITTEN IN PEILINGWIJZER
print len(set(peiling.Partij))
#HOEVEEL UNIEKE PARTIEN KANDIDATENLIJST?
len(set(df.partij))
#WIE HEEFT HOOGSTE AANTAL ZETELS?
peiling[peiling['ZetelsHoog']== peiling.ZetelsHoog.max()]

```

```

#PARTIEN SORTEREN IN DATAFRAME MET HOOGSTE
PERCENTAGE
peiling = peiling.set_index('Partij')
percentage1 = peiling.sort_values(by='Percentage',
ascending=False)
percentage1.head(2)

```

```

#PLOTTEN VERSCH MANIEREN
percentage1['Percentage'].plot(kind='bar', title =
'Hoogste Percentage')
percentage1['Percentage'].plot(kind='pie', title =
'Hoogste Percentage')
percentage1['Percentage'].plot(title = 'Hoogste
Percentage')

```

```

#XML FILES
import xml.etree.ElementTree as ET
tree =
ET.parse('Kandidatenlijsten_TK2017_Amsterdam.eml.xml')
root = tree.getroot()
root.tag
root.attrib

```

```

#ALLE TAGS UIT XML FILE HALEN
for tag in tree.iter():
print tag

```

```

#HOEVEEL M-V ZITTEN IN KANDIDATENLIJST VAN
AMSTERDAM
mannen = 0
vrouwen = 0
for event, elem in
ET.iterparse('Kandidatenlijsten_TK2017_Amsterdam.eml.
xml'):
value = elem.text
tag = elem.tag
if tag ==
"urn:oasis:names:tc:evs:schema:eml:Gender" :
if value == 'female':
vrouwen += 1
else:
mannen += 1
elem.clear()
print "TEKST" mannen
print "TEKST" vrouwen
print "TEKST TOTAAL" (mannen + vrouwen)

```

```

#EEN DATAFRAME UIT XML FILE HALEN
partij = [], voornaam = [], initialen = [], achternaam = [],
geslacht = [], woonplaats = []

```

```

for event, elem in
ET.iterparse('Kandidatenlijsten_TK2017_Amsterdam.eml.
xml'):
tag = elem.tag
value = elem.text (ZOEKEN IN ALLE TAGS)
if tag ==
"urn:oasis:names:tc:evs:schema:eml:RegisteredName":
partij.value = value
if tag ==
"urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0:FirstName":
partij.append(partij.value)
voornaam.append(value)
if tag ==
'urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0:NameLine':
initialen.append(value)
if tag ==
"urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0:LastName":
achternaam.append(value)
if tag ==
"urn:oasis:names:tc:evs:schema:eml:Gender":
if value == 'male':
geslacht.value = 'man'
else:
geslacht.value = 'vrouw'
geslacht.append(geslacht.value)
if tag ==
'urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0:LocalityName':
woonplaats.append(value)

```

```

df = pd.DataFrame({'partij':partij, 'voornaam':voornaam,
'initialen':initialen, 'achternaam':achternaam,
'geslacht':geslacht, 'woonplaats':woonplaats})
df.head(5)

```

```

print len(df[df.geslacht=='man'])
print len(df[df.geslacht=='vrouw'])

```

```

#WAT IS PERCENTAGE KANDIDATEN VAN DENK VAN
TOTAAL KANDIDATEN
(df.partij == 'DENK').mean()*100

```

```

#WAT IS LIJST VAN ALLE UNIEKE PARTIEN
print (df['partij']).unique()
len(set(df['partij']))

```

```

#HOEVEEL KANDIDATEN ZIJN ER PER PARTIJ + PLOT?
totaalaantalkandidaten = df['partij'].value_counts()
totaalaantalkandidaten.plot(kind='bar', title =
'Hoeveelheid kandidaten per partij')

```

```

#VIND ALLE MANNEN OF VROUWEN UIT DENK
df.loc[df['partij'] == 'DENK'] & (df['geslacht'] ==
'man').head()
df.loc[(df['partij'] == 'DENK') & (df['geslacht'] ==
'vrouw')].head()

```

```

#HOEVEEL PROCENT IS MAN VAN DENK?
#TOTAAL MANNEN UIT DENK
mannen_denk = len(df.loc[(df['partij'] == 'DENK') &
(df['geslacht']=='man')])
#TOTAAL VROUWEN UIT DENK
vrouwen_denk = len(df.loc[(df['partij'] == 'DENK') &
(df['geslacht']=='vrouw')])
totaal_denk = len(df[df.partij == 'DENK'])
procentman = (mannen_denk * 100) / totaal_denk
print procentman
#WAT IS RELATIEVE VERHOUDING M-V? ROND AF 2
DECIMALEN.
np.round(vrouwen_denk/float(mannen_denk), 2)

```

```

#VERHOUDING M-V DENK + PLOT
manvrouw = df['geslacht'].loc[(df['partij'] ==
'DENK')].value_counts()
manvrouw.plot(kind='bar', title='Verhouding man vrouw
in partij Denk')

```

```

#VERHOUDING WOONPLAATSEN KANDIDATEN UIT
DENK + PLOT
verhouding = df['woonplaats'].loc[(df['partij'] ==
'DENK')].value_counts()
verhouding.plot(kind='pie', title='Verdeling
woonplaatsen kandidaten DENK')

```

```

#KANDIDAAT UIT WOONPLAATS?
df.loc[(df['woonplaats'] == 'Wassenaar')]
#KANDIDAAT UIT WOONPLAATS EN PARTIJ?
df.loc[(df.woonplaats == 'Wassenaar') & (df.partij ==
'VVD')]
# of optie df.loc[(df['woonplaats'] == 'Wassenaar') &
(df['partij'] == 'VVD')]
#LIJST KANDIDATEN MEERDERE WOONPLAATSEN?
df.loc[(df['woonplaats'].isin(['Wassenaar', 'Assen',
'Utrecht']))].head()
df.loc[(df['partij'].isin(['VVD', 'SP']))].head()

```

```

#GEEF EEN LIJST VAN DE PARTIEN DIE EEN BAS ALS
KANDIDAAT HEBBEN
df['partij'].loc[(df['voornaam'] == 'Bas')].unique()
#WEERGEGEVEN IN DATAFRAME
df.loc[(df['voornaam'] == 'Bas')]

```

```

#KANDIDATEN ACHTERNAAM MET W BEGINT
df[['voornaam', 'achternaam',
'partij']].loc[(df['achternaam'].str.startswith('W'))].head()

```

```

#PVV 21 ZETELS, HOEVEEL VROUWELIJKE PVVERS IN DE
KAMER?
1. DATAFRAME ALLE KANDIDATEN PVV
df_vanpvv = df.loc[(df['partij'] == 'PVV (Partij voor de
Vrijheid)')].head(peiling['Zetels']).loc[peiling['Partij'] ==
'PVV'].values.item()
2. HOEVEEL VROUWEN ZITTEN HIER IN?

```

```

len(df_vanpvv.loc[(df_vanpvv['geslacht'] == 'vrouw')])
3.PLOTTEN
df_vanpvv['geslacht'].value_counts().plot(kind='bar')
OFFFF OPTIE:
df_kandidatenlijst_adam_pv20 =
df_kandidatenlijst_adam[df_kandidatenlijst_adam.partij
== 'PVV (Partij voor de Vrijheid)'].head(20)
print
len(df_kandidatenlijst_adam_pv20[df_kandidatenlijst_
adam_pv20.gender == 'V']), 'vrouwen.'

```

```

#VERHOUDING M/V + PLOT
vrouwen_partij = []
mannen_partij = []
for a in df['partij'].unique():
try:
vrouw1 = df['geslacht'].loc[(df['partij'] == a) &
(df['geslacht'] == 'vrouw')].value_counts()[0]
except:
vrouw1 = 0

```

```

man1 = df['geslacht'].loc[(df['partij'] == a) &
(df['geslacht'] == 'man')].value_counts()[0]
vrouwen_partij.append(vrouw1)
mannen_partij.append(man1)

```

```

#DATAFRAME MAKEN
mannenvrouwen_df =
pd.DataFrame({'partij':df['partij'].unique(),
'mannen':mannen_partij, 'vrouwen':vrouwen_partij})
mannenvrouwen_df

```

```

#PLOTTEN ZELFDE GRAFIEK, TWEE KLEUREN .3
plotverhouding =
mannenvrouwen_df.plot(x='partij', y='mannen',
kind = 'bar')
mannenvrouwen_df.plot(x='partij', y='vrouwen',
kind = 'bar', color='red', ax=ax)

```

```

#HOOGSTE PERCENTAGE M-V PARTIJ?
mannenvrouwen_df['percentagevanvrouw'] =
((mannenvrouwen_df['vrouwen'] /
(mannenvrouwen_df['mannen'] +
mannenvrouwen_df['vrouwen'])) * 100).round(2)
mannenvrouwen_df['percentagevanman'] =
((mannenvrouwen_df['mannen'] /
(mannenvrouwen_df['mannen'] +
mannenvrouwen_df['vrouwen'])) * 100).round(2)
mannenvrouwen_df
#PLOT DEZE PERCENTAGES IN EEN DATAFRAME
Zelfde als .3 hierboven alleen dan met y =
'percentagevanman' en y = 'percentagevanvrouw'

```

```

# DOOR ALLE LIJSTEN LOOPEN NAAR ALLE MANNEN EN
VROUWEN (DUS DUBBELE)
totaalmannen = 0
totaalvrouwen = 0

```

```

path = 'mappie/'
#print os.listdir(path)
for file in os.listdir(path):
if not file.endswith('.xml'):continue
# print_file
aantalmannen = 0
aantalvrouwen = 0
with open(path+_file):
for event, elem in ET.iterparse(path+_file):
tag = elem.tag
value = elem.text
if tag ==
"urn:oasis:names:tc:evs:schema:eml:Gender" :
if value == 'male':
aantalmannen += 1
totaalmannen += 1
else:
aantalvrouwen += 1

```

```

elem.clear() # discard the element
print_file.strip('eml.xml'), aantalmannen,
aantalvrouwen, (aantalmannen + aantalvrouwen)

```

```

#20 POPULAIRE WOONPLAATSEN
woonplaats = df.woonplaats.value_counts()
woonplaats = woonplaats[:20]
woonplaats.plot(kind='barh', title = 'blablalbl')
(barh = horizontaal)

```

```

#MAAK KRUISTABEL WAARIN JE GENDER TEGEN DE
WOONPLAATS PLAATST
cross = pd.crosstab(df.woonplaats, df.gender, margins =
True)
print cross.sort_values(by='All', ascending =
False).head()

```

```

#NaN WAARDEN OPlossen
df_join.update(df_join[['Percentage', u'PercentageLaag',
u'PercentageHoog', u'Zetels', u'ZetelsLaag',
u'ZetelsHoog']].fillna(0))
df_join.head(5)

```

## VIEWING/INSPECTING DATA

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

## SELECTION

df[col] - Return column with label col as Series  
df[[col1, col2]] - Return Columns as a new DataFrame  
s.iloc[0] - selection by position  
s.loc[0] - selection by index df.iloc[0,] - first row  
df.iloc[0,0] - first element of first column

## DATA CLEANING

df.columns = ['a','b','c'] - Rename columns  
pd.isnull() - Checks for null Values, Returns Boolean Array  
pd.notnull() - Opposite of s.isnull() df.dropna() - Drop all rows that contain null values  
df.dropna(axis=1) - Drop all columns that contain null values  
df.dropna(axis=1, thresh=n) - Drop all rows have less than n non null values df.fillna(x) - Replace all null values with x  
s.fillna(s.mean()) - Replace all null values with the mean (mean can be replaced with almost any function from the statistics section)  
s.astype(float) - Convert the datatype of the series to float  
s.replace(1,'one') - Replace all values equal to 1 with 'one'  
s.replace([1,3],[one,'three']) - Replace all 1 with 'one' and 3 with 'three'  
df.rename(columns=lambda x: x + 1) - mass renaming of columns df.rename(columns={'old\_name': 'new\_name'}) - selective renaming df.set\_index('column\_one') - change the index df.rename(index=lambda x: x + 1) - mass renaming of index

## FILTER, SORT, & GROUPBY

df[df[col] > 0.5] - Rows where the col column is greater than 0.5  
df[(df[col] > 0.5) & (df[col] < 0.7)] - Rows where 0.7 > col > 0.5  
df.sort\_values(col1) - Sort values by col1 in ascending order  
df.sort\_values(col2, ascending=False) - Sort values by col2 in descending order  
df.sort\_values([col1, col2], 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 1 group  
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

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.

```
import pandas as pd
import os
from lxml import etree
import xml.etree.ElementTree as ET
%matplotlib inline
from math import sqrt
import matplotlib.pyplot as plt

#DF van excel
pwijzer = pd.read_excel('naam.xlsx')

#informatie van je DF (eerst aantal tabellen)
print (pwijzer.shape, pwijzer.describe)

pwijzer.head()

# geeft kolom
pwijzer['Partij']

pwijzer['Partij'].sum() #telt alle waardes op
pwijzer['Partij'].count() #telt aantal rijen

#kolom toevoegen met nieuwe waarde, bijv. gem.
pwijzer['GemZetels'] = (pwijzer['ZetelsLaag'] +
pwijzer['ZetelsHoog']) / 2
pwijzer

#getal uit DF (als variabele voor functie)
inter = pwijzer['Percentage'].loc[[pwijzer['Partij']
== 'PVV']].values[0] #kan ook met 'Partij'

#kolom weg en index op partij
pwijzer[['Partij',
'Percentage']].set_index('Partij').head()

#partij meer dan 25 zetels
zetels = pwijzer.loc[pwijzer['Zetels'] > 25]

#Sorteer op kolom
pwijzer2 = pwijzer2.sort_values(by='Zetels',
ascending=False)

#XML
#TAGS
import xml.etree.ElementTree as ET
tree = ET.parse('xml')
root = tree.getroot()
root.tag
root.attrib

# tags:
for tag in tree.iter():
    print (tag)

#OF
tags = [ ]
for event, elem in ET.iterparse('naam.xml'):
    tag = elem.tag
    tags.append(tag)
    elem.clear()

tags

#aantal V/M per kandidaatlijst (loop door alle .xml)
totalman = 0
totalvrouw = 0
path = 'naam-map'
#Kandidatenlijsten_EML_TK2017-20170213/
#print os.listdir(path) (test)
for _file in os.listdir(path):
    if not _file.endswith('xml'):continue
    #print _file(test)
    countman = 0
    countvrouw = 0
    with open(path+_file):
        for event, elem in ET.iterparse(path + _file):
            tag = elem.tag
            value = elem.text
            if tag ==
"({urn:oasis:names:tc:evs:schema:eml}Gender)":
                if value == 'male':
                    countman += 1
                    totalman += 1
                else:
                    countvrouw += 1

            elem.clear()
            print (_file.strip('.eml.xml'), countman,
countvrouw, (countman + countvrouw))

#hoeveel m/v+tot. Per kandidaatlijst
countman2 = 0
countvrouw2 = 0
for event, elem in
ET.iterparse('amsterdam.xml'):
    tag = elem.tag
    value = elem.text
    if tag ==
"({urn:oasis:names:tc:evs:schema:eml}Gender)":
        if value == 'male':
            countman2 += 1
        else:
            countvrouw2 += 1
    elem.clear()
    print ('aantal man: = ', countman2)
    print ('aantal vrouw: = ', countvrouw2)
    print ('total: = ', (countvrouw2 + countman2))
```

```
#list partijen op kandidaatlijst
partijen = [ ]
for event, elem in
ET.iterparse('amsterdam.xml'):
    tag = elem.tag
    value = elem.text
    if tag ==
"({urn:oasis:names:tc:evs:schema:eml}Registere
dName)":
        partijen.append(value)
    elem.clear()
partijen

#MAAK DF
partij = [ ]
voornaam = [ ]
achternaam = [ ]
gender = [ ]

for event, elem in
ET.iterparse('Kandidatenlijsten_TK2017_Amster
dam.eml.xml'):
    tag = elem.tag
    value = elem.text
    if tag ==
"({urn:oasis:names:tc:evs:schema:eml}Registere
dName)":
        partij_val = value
        if tag ==
"({urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}Firs
tName)":
            partij.append(partij_val)
            voornaam.append(value)
        if tag ==
"({urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}Las
tName)":
            achternaam.append(value)
        if tag ==
"({urn:oasis:names:tc:evs:schema:eml}Gender)":
            if value == 'male':
                gender_val = 'man'
            else:
                gender_val = 'vrouw'
            gender.append(gender_val)

df = pd.DataFrame({'partij':partij,
'voornaam':voornaam,
'achternaam':achternaam, 'geslacht':gender})
df.head()

#aantal unique partijen
len(df['partij'].unique())

#values kolom weergeven
df['partij'].head()

#alle M van partij
df.loc[(df['partij'] == 'X') & (df['geslacht'] ==
'man')]

len(df.loc[(df['partij'] == 'X') & (df['geslacht'] ==
'man')])

#wie zit in partij x?
df.loc[df['partij'].isin(['X','Y'])]

#Vrouw in partij?
df['partij'].loc[(df['geslacht'] ==
'vrouw')].unique()

len(df['partij'].loc[(df['geslacht'] ==
'vrouw')].unique())

#alle voornamen met Mo
df[['voornaam',
'partij']].loc[(df['voornaam'].str.startswith('Mo'))
]

PLOTTEN

#M/V van een partij
verhouding = df['geslacht'].loc[(df['partij'] ==
'X')].value_counts()

verhouding.plot(kind='pie')

#woonplaats binnen een partij
verhouding = df['woonplaats'].loc[(df['partij'] ==
'X')].value_counts()

verhouding.plot(kind='pie')

#partij met aantal kandidaten
partij_aantal = df['partij'].value_counts()
partij_aantal.plot(kind='bar')

#M-V van partij in BAR
vrouwen = [ ]
mannen = [ ]
for a in df['partij'].unique():
    try:
        V = df['geslacht'].loc[(df['partij'] == a) &
(df['geslacht'] == 'vrouw')].value_counts()[0]
    except:
        V = 0

M = df['geslacht'].loc[(df['partij'] == a) &
(df['geslacht'] == 'man')].value_counts()[0]
vrouwen.append(V)
mannen.append(M)
```

```
#DF
mvd =
pd.DataFrame({'partij':df['partij'].unique(),
'mannen':mannen, 'vrouwen':vrouwen})
mvd

#BARPLOT
ax = manvroudf.plot(x='partijen', y='mannen',
kind = 'bar')
manvroudf.plot(x='partijen', y='vrouwen', kind
= 'bar', color='red', ax=ax)

#HOOGSTE PERCENTAGE MANNEN EN
VROUWEN PER PARTIJ OP DE
KANDIDATENLIJST?
manvroudf['percentagevrouwen'] =
((manvroudf['vrouwen'] /
(manvroudf['mannen'] +
manvroudf['vrouwen'])) * 100).round(2)

manvroudf['percentageman'] =
((manvroudf['mannen'] /
(manvroudf['mannen'] +
manvroudf['vrouwen'])) * 100).round(2)

manvroudf

#PLOT DATAFRAME
ax = manvroudf.plot(x='partijen',
y='percentageman', kind = 'bar')
manvroudf.plot(x='partijen',
y='percentagevrouwen', kind = 'bar', color='red',
ax=ax)

manvroudf

#PVV 21 zetels, Hoeveel vrouwelijk PVV in de
kamer?

1. DF met X aantal namen van kandidatenlijst
in partij
df_pvv = df.loc[(df['partij'] == 'PVV (Partij voor
de
Vrijheid)')].head(pwijzer['Zetels']).loc[[pwijzer['P
artij'] == 'PVV']].values.item()

#OF
df_adam_pvv_20 =
df_kandidatenlijst_adam[df_kandidatenlijst_adam
.m partij == 'PVV (Partij voor de
Vrijheid)'].head(20)
print
(len(df_kandidatenlijst_adam_pvv_20[df_kandid
atenlijst_adam_pvv_20.gender == 'v'],
'vrouwen.').

2. Tel conditie in DF
len(df_pvv.loc[(df_pvv['geslacht'] == 'vrouw')])

3. plot
df_pvv['geslacht'].value_counts().plot(kind='bar'
)

#kruistabel: woonplaats x gender
cross_gender_ woonplaats =
pd.crosstab(df_kandidatenlijst_adam_ woonplaats
, df_kandidatenlijst_adam.gender, margins
= True)
print cross_gender_ woonplaats.sort_values(by
= 'All', ascending = False).head(5)

#Maak een horizontale staafdiagram met de
verdeling van de 20 meest populaire
woonplaatsen van alle kandidaten van de
kandidatenlijst van Amsterdam. Doe dit ook
netjes op volgorde en geef de grafiek een
passende titel.

counter =
df_kandidatenlijst_adam.woonplaats.value_cou
nts()
counter = counter[:20]
counter.plot(kind = 'barh', title = '20 meest
populaire woonplaatsen voor de kandidaten van
Amsterdam')

#Wat is de verdeling van de zetels volgens de
peilingwijzer? Neem hierbij de kolom 'zetels'
en plot dit in een staafdiagram tegen de
verschillende partijen op volgorde van de
hoeveelheid zetels met een passende titel.

df_peilingwijzer.sort_values(by='Zetels',
ascending = False).plot('Partij', 'Zetels', kind =
'bar', title = 'Verdeling Zetels volgens
Peilingwijzer')

#Vervang de NaN waarden met 0 voor alle
kolommen behalve de kolom 'Datum' in een
onliner en laat de eerste 5 regels zien.
df_join.update(df_join[['Percentage',
u'PercentageLaag', u'PercentageHoog', u'Zetels',
u'ZetelsLaag', u'ZetelsHoog']].fillna(0))

df_join.head(5)
```

```
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 col- umns
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
pd.DataFrame(np.random.rand(20,5)) - 5 col- umns and
20 rows of random floats
pd.Series(my_list) - Create a
series from an iterable my_list
df.index = pd.date_range('1900/1/30',
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 informa-
tion
df.describe() - Summary statistics for numerical
columns
s.value_counts(dropna=False) - View unique
values and counts of apply(pd.Series, value_counts) -
Unique values and counts for all columns

SELECTION
df[col] - Return column with label col as Series
df[[col1,
col2]] - Return Columns as a new DataFrame
s.iloc[0] - selection by position
s.loc[0] - selection by index
df.iloc[0,0] - first
row
df.iloc[0,0] - first element of first column

DATA CLEANING
df.columns = ['a','b','c'] - Rename columns
pd.isnull() -
Checks for null Values, Returns Boolean
Array
pd.notnull() - Opposite of s.isnull()
df.dropna() -
Drop all rows that contain null values
df.dropna(axis=1) - Drop all columns that con-
tain null values
df.dropna(axis=1, thresh=n) - Drop all rows have
less than n non null values
df.fillna(x) - Replace all
null values with x
df.fillna(s.mean()) - Replace all null values
with the mean (mean can be replaced with almost any
function from the statistics section)
s.astype(float) -
Convert the datatype of the series to float
s.replace(1,'one') - Replace all values equal to 1 with
'one'
s.replace(1,3,['one','three']) - Replace all 1 with
'one' and
3 with 'three'
df.rename(columns=lambda x: x + 1) - mass
renaming of columns
df.rename(columns='old_name':
'new_name') - selective renaming
df.set_index('column_one') - change the index
df.rename(index=lambda x: x + 1) - mass renaming of
index

FILTER, SORT, & GROUPBY
df[df[col] > 0.5] - Rows where the col column is
greater than 0.5
df[(df[col] > 0.5) & (df[col] < 0.7)] - Rows where
0.7 > col > 0.5
df.sort_values(col1) - Sort values by col1 in
ascending order
df.sort_values(col2, ascending=False) -
Sort values by col2 in descending order
df.sort_values(col1,col2),

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 1
group
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
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.
```

```
150 mensen in TK
```

```
%matplotlib inline
```

```
from __future__ import division
```

### Element letterlijk printen

```
print etree.tostring(xml OF element,  
pretty_print=True)
```

### Boxplotje maken

```
df.boxplot(by = 'Party_name', column =  
'Fraction', figsize = (16,10), rot =  
30)
```

### Bepaalde kolommen (met bepaalde inhoud) niet krijgen:

```
party_columns =  
df.columns[~(df.columns.str.contains('low') |  
df.columns.str.contains('high'))]
```

### Selecteren op inhoud van index:

```
index.str.contains ('....')
```

### Specifieke kolom pakken

```
df['kolomnaam']
```

### Meerdere kolommen pakken

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

### Specifieke rij pakken met index

```
df.loc['index-waarde']
```

### De i-de rij pakken

```
df.iloc[i]  
df.iloc[3]  
df.iloc[-1]
```

### Kolom als index instellen

```
df = df.set_index('kolom')
```

### Specifieke kolom van specifieke rij pakken

```
peiling.loc['D66']['Zetels']  
peiling['Zetels'].loc['D66']
```

### Unieke waarden in kolom

```
df['kolom'].unique()
```

### Unieke waarden tellen

```
df['kolom'].value_counts()
```

### Groupby element unstacken

```
groupby element.size().unstack()
```

### dataframe unstacken + index resetten

```
df.unstack().reset_index()
```

### Capslock van kandidaten vermijden

```
candidates_df['Party_name'].str.lower(  
) .str.contains(party_in_lijst.lower())
```

### Sorteren op kolom

```
df.sort_values('kolom',  
ascending=True)
```

### List van DataFrames aan elkaar plakken

```
pd.concat
```

### XML bestand importen

```
from lxml import etree
```

```
with open
```

```
('Kandidatenlijsten_TK2017_Amsterdam.e  
ml.xml') as f:
```

```
parser =  
etree.XMLParser(ns_clean=True)  
xml = etree.parse(f, parser)
```

```
ns = { ns: '{'+ url +}'' for ns, url  
in xml.getroot().nsmap.items() }
```

```
ns
```

```
contest = xml.find('//'+ ns[None]  
+'Contest')
```

```
all_candidates = []
```

```
affiliations =  
contest.findall(ns[None] +  
'Affiliation')
```

```
for party in affiliations:  
> party.find(affiliation_identifler)  
> affiliation_identifler.find  
(registered_name).text  
> candidates: party.findall
```

```
for candidate in candidates:
```

```
candidate_name,  
candidate_adress, candidate_country,  
candidate_gender
```

```
if candidate_country is not  
None:
```

```
candidate_country_initials,  
candidate_locality_notNL,  
candidate_adress  
else: candidate_country_initials  
= 'NL', candidate_locality,  
candidate_adress
```

### Attribuut van element

```
Candidate_identifler =  
candidate.find(ns[None] +  
'CandidateIdentifler')
```

```
candidate_position =  
candidate_identifler.attrib['Id']
```

```
if candidate_prefix is not None: ...
```

### Alle variabelen toevoegen aan een lijstje

```
.append(( ))
```

```
Kruistabel: pd.crosstab(df[kolom],  
df[kolom], margins = True (voor  
total))
```

```
apply voor elk element in het hele  
dataframe: applymap
```

```
partijen = peilingen.columns[1:]
```

```
for partij in partijen:
```

```
zetels_gesplit =  
peilingen[partij].str.split(';').str
```

```
peilingen[partij + '.low'],  
peilingen[partij], peilingen[partij  
+ '.high'] = zetels_gesplit
```

```
def alleen_de_middelste(x):  
return x.split(';')[1]
```

```
peilingen[partijen] =  
peilingen[partijen].applymap(alleen_de  
_middelste)
```



```

#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}").text
        dct["woonplaats"] = c.find("{urn:oasis:names:tc:evs:schema:eml:QualifyingAddress}/{urn:oasis:names:tc:evs:schema:eml:Localities}/{urn:oasis:names:tc:evs:schema:eml:Localities}").text
        lst.append(dct)

#neem specifieke partij
cda = df[df.partij == 'CDA']
#kruistabel van alle partijen op 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})

#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')

#TEL ALLE VALUES IN EEN KOLOM OP
print peil['Zetels'].sum()
#BEREKEN HET VERSCHIL VAN TWEE KOLOMMEN EN VOEG DIT TOE IN EEN EXTRA KOLOM
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()
complete = 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] = 'VNL (VoorNederland)'
peilingwijzer.Partij[13] = 'DENK'
peilingwijzer.Partij[14] = 'Forum voor Democratie'
peilingwijzer.Partij[15] = 'Piratenpartij'
peilingwijzer

vrouwenlaag = {}
mannenlaag = {}
vrouwenhoog = {}
mannenhoog = {}
vrouwen = {}
mannen = {}

for item in complete.iterrows():
    vrouwenlaag.append(item["gender"] + item["ZetelsLaag"].count("v"))
    vrouwenhoog.append(item["gender"] + item["ZetelsHoog"].count("v"))
    mannenlaag.append(item["gender"] + item["ZetelsLaag"].count("m"))
    mannenhoog.append(item["gender"] + item["ZetelsHoog"].count("m"))
    vrouwen.append(item["gender"] + item["Zetels"].count("v"))
    mannen.append(item["gender"] + item["Zetels"].count("m"))

complete["vrouwenlaag"] = vrouwenlaag
complete["vrouwenhoog"] = vrouwenhoog
complete["mannenlaag"] = mannenlaag
complete["mannenhoog"] = mannenhoog
complete["mannen"] = mannen
complete["vrouwen"] = vrouwen

#matplotlib inline
ding = complete[["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.mmap(C{'female':0, 'male':1})
pd.merge(df1, df2, on=c)
df.c.replace(v, av)
df.nivo6-table

```

#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()

%s -lh "file"  
import pandas as pd  
import numpy as np  
import seaborn as sns  
import lxml  
from lxml import etree  
import matplotlib inline

ax.is = 0 = 1;  
↓  
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('xml'): continue  
with open(path + file):  
tree = ET.parse(path + file)

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

Import xml.etree.ElementTree as ET

1 kandidaten = pd.DataFrame(data)

datafilem66

2

4

5

6

7

3

df.c.replace(v, av)

df.nivo6-table

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

c='sex'

```

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
partynamen= [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={'SOPLOS': 'SOPLOS',
           'CDA': 'CDA',
           'CU': 'ChristenUnie',
           'D66': 'Democraten 66 (D66)',
           'Denk': 'DENK',
           'FVD': 'Forum voor Democratie',
           'GL': 'GroenLinks',
           'PVV': 'PVV (Partij voor de Vrijheid)',
           'PvdA': 'Partij van de Arbeid (P.v.d.A.)',
           'PvdD': 'Partij voor de Dieren',
           'SGP': 'Staatkundig Gereformeerde Partij (SGP)',
           'SP': 'SP (Socialistische Partij)',
           'VVD': 'VVD (VoorNederland)',
           '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"
laatstedf= pd.read_excel(laatstepeiling)
alldf= pd.read_excel(allepeilingen)
print alldf.shape
alldf.tail()
laatstedf.index= laatstePartij
#laatstedf
laatstedf.rename(translate, inplace=True)

# Dit is een handige file met de zetel aantallen
allezdf=pd.read_csv('https://d1b3gg97if6urz.cloudfront.net/Public/Peilingwijzer/Last/Results_DyGraphs_Seats.csv',
                   index_col='Date')
print allezdf.shape
print (allezdf.sum(axis=1) == 150).mean() # check if all add to 150
allezdf=allezdf.applymap(lambda s: int(s.split(':')[1])) # want (ci low; mean; ci high)
# now translate to our names
allezdf=allezdf.T
allezdf.rename(translate, inplace=True)
allezdf=allezdf.T
allezdf.tail()

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

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

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

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

# Aantal vrouwen in de kamer
VrouwenPeilingwijzer=df.sort_values('Zetels',ascending=True)[[ u'Zetels',u'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(acental_vocuen ['50 plus'], orient='index')
plot(kind='barh', figsize=(10,8))
```

```
# 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']}
Counter(genderdict[partij voor de Dieren][:25])['female']
gendercount = {p: Counter(genderdict[p][:N]) for p in genderdict}
gdf = pd.DataFrame.from_dict(gendercount, orient='index', fillna(0).astype(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  
SEP = ;  
**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('1900/1/30', 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

**SELECTION**  
df[col] - Return column with label col as Series  
df[[col1, col2]] - Return Columns as a new DataFrame  
s.iloc[0] - selection by position  
s.loc[0] - selection by index  
df.iloc[0,:] - first row  
df.iloc[0,0] - first element of first column

**DATA CLEANING**  
df.columns = ['a', 'b', 'c'] - Rename columns  
pd.isnull() - Checks for null Values, Returns Boolean Array  
pd.notnull() - Opposite of s.isnull()  
df.dropna() - Drop all rows that contain null values  
df.dropna(axis=1) - Drop all columns that contain null values  
df.dropna(axis=1, thresh=n) - Drop all rows have less than n non null values  
df.fillna(x) - Replace all null values with x  
s.fillna(s.mean()) - Replace all null values with the mean (mean can be replaced with almost any function from the statistics section)  
s.astype(float) - Convert the datatype of the series to float  
s.replace(1, 'one') - Replace all values equal to 1 with 'one'  
s.replace([1,3], ['one', 'three']) - Replace all 1 with 'one' and 3 with 'three'  
df.rename(columns=lambda x: x + 1) - mass renaming of columns  
df.rename(columns={'old\_name': 'new\_name'}) - selective renaming  
df.set\_index('column\_name') - change the index  
df.rename(index=lambda x: x + 1) - mass renaming of index

**FILTER, SORT, & GROUPBY**  
df[df[col] > 0.5] - Rows where the col column is greater than 0.5  
df[(df[col] > 0.5) & (df[col] < 0.7)] - Rows where 0.7 > col > 0.5  
df.sort\_values(col1) - Sort values by col1 in ascending order  
df.sort\_values(col2, ascending=False) - Sort values by col2 in descending order  
df.sort\_values([col1, col2],

~~translate = lambda s: s.split(' ')[0]~~

```
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_vocuen(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 stemmingen 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='m')
```

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

```
zetels = pd.read_csv
zetels = zetels.set_index('Da
zetels = applymap
zetels = zetels.T
zetels = zetels.join
zetels = zetels.zetels2 = zetels.SU
```

```
zetels2.cleasa
zetels.sort(as=1)
```

```
str.startswith('...')
```

```
Teu
```

```
pd.value_counts
```

```
float(stemmen[partij].loc[1])/float(stemmen[partij].sum())
```

```
pd.DataFrame(populair).sort(0, ascending=False).plot(kind='barh', title='m')
```

```
# Filteren in een dataframe
peilingen.sort('zetels', ascending=False)
.plot('Partij', 'zetels', kind =
```

```
#matplotlib inline
import pandas as pd
import re
from lxml import etree
from bz2file import BZ2File
import codecs
import urllib2
import urllib
import os
import seaborn as sns
import zipfile
import bs4 as bs
from collections import Count
```

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

~~Handwritten scribbles and notes in blue ink, including 'zetels', 'partij', 'zetels2', 'names', 'Count', 'idxmax', 'geeft index', 'names.Count', 'median', 'names.Count == names.Count.median()'.~~

value-counts  
↳ count freq that value occurs  
bv woonplaats hoe vaak  
# filteren df om plotten  
peilingen.sort('zetels', ascending=False).plot('Party', 'zetels', kind='barh')

```

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 df.groupby(col1).agg(np.mean)
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-ki
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_TR2017-20170213/Kandidatenlijsten_TR2017_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:XNL:2.0}NameLine':
initialen.append(value)

if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XNL:2.0}FirstName':
voornaam.append(value)
partij.append(partij_val)

if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XNL:2.0}LastName':
achternaam.append(value)
prefix.append(prefix_val)

if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XNL:2.0}NamePrefix':
prefix_val = value
else:
prefix_val = ''

if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XNL: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')
```

```
#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
```

#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
1	VVD	2017-03-14	17.2	16.5	17.9	26	24

```
#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

```
#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 ZOEKEN 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-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 = 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 WOONPLAATS

```
df[['voornaam', 'partij', 'woonplaats']].loc[(df['voornaam'].str.startswith('B')) & df['woonplaats'].isin(['Amsterdam', 'Hertogenbosch'])]
```

#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()
```

*woonplacem*

#MAAK EEN CROSSTAB VAN DE WOONPLAATS EN HIER DE MAN VROUW VERHOUDING, PLOT DIT IN EEN STACKED PLOT

```
#ALLES 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())
```

```
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-ki
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_list = []
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_list.append(int(event.attrib.values()[0]))

    if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XNL:2.0}NameLine':
        initialen.append(value)

    if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XNL:2.0}FirstName':
        voornaam.append(value)
        partij.append(partij_val)

    if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XNL:2.0}LastName':
        achternaam.append(value)
        prefix.append(prefix_val)

    if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XNL:2.0}NamePrefix':
        prefix_val = value
    else:
        prefix_val = ''

    if tag == '{urn:oasis:names:tc:ciq:xsd:schema:XAL: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')
```

```
#COMBINEER ALLE LIJSTEN TOT EEN DATAFRAME
```

```
df = pd.DataFrame({'partij':partij, 'positie_list':positie_list, 'prefix':prefix, 'initialen':initialen, 'voornaam':v,
                  'achternaam':achternaam, 'woonplaats':woonplaats, 'gender':gender})
```

```
df = df[['partij', 'positie_list', '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()
```

```
df['gender_nieuw'] = df.gender.apply(gender_to_num)
```

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

```
#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 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 KOLOM OP
peilingwijzer['zetels'].sum()
```

```
#BEREKEN HET VERSCHIL VAN DRIE 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 ZOEKEN 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-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 = 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_list'] <= (pw_df['zetels'].loc[pw_df['Partij'] == a]).values[0]) &
                              (df['gender'] == 'vrouw')]))
    mannen.append(len(df.loc[(df['partij'] == a) &
                              (df['positie_list'] <= (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 WOONPLAATS
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 UNIEKE 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
```

```
#ALLES 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())
```

XML inladen

DF

pw\_df

Mannen

```
#parse xml
tree = lxml.etree.parse("Kandidatenlijst_2017_Amsterdam.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}First").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)

#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
```

```
#neem specifieke partij #aantal unieke partijen
cda = df[df.partij == 'CDA'] len(df.partij.unique())
''' Pak het aantal zetels uit de peilingen
dataset voor d66 en kijk hoeveel mannen en
vrouwen er uiteindelijk in de kamer komen voor PVV'''
vr = df[df2["partij"] == 'PVV (Partij voor de Vrijheid)']
vr = vr[vr["nummer"] < 26]
vr.gender.value_counts()
```

```
#kruis tabel 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 = "Kandidatenlijst_2017_Amsterdam.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})
```

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')
```

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

```
#HOEVEEL % VAN TOTAAL KANDIDATEN STAAT OP DE LIJST VOOR SPECIFIEKE PARTIJEN
(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'})
```

④

```
genders = kandidaten.groupby("Partij")["gender"].sum()
dfg = genders.to_frame()
dfg = dfg.reset_index()
```

②

```
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] = 'Staatskundig 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] = 'Piratenpartij'
```

⑤

⑤

⑥

```
compleet = peilingwijzer.merge(dfg, on="Partij")
vrouwenlaag = {}
mannelaag = {}
vrouwenhoog = {}
mannelaaghoog = {}
vrouwen = {}
mannen = {}
for item in compleet.iterrows():
    vrouwenlaag.append(item[1]["gender"][:item[1]["ZetelsLaag"]].count("v"))
    vrouwenhoog.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"))
    vrouwen.append(item[1]["gender"][:item[1]["Zetels"]].count("v"))
    mannen.append(item[1]["gender"][:item[1]["Zetels"]].count("m"))
compleet["vrouwenlaag"] = vrouwenlaag
compleet["vrouwenhoog"] = vrouwenhoog
compleet["mannelaag"] = mannelaag
compleet["mannelaaghoog"] = mannelaaghoog
compleet["vrouwen"] = vrouwen
compleet["mannen"] = mannen
```

```
#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()
```

#read peilingwijzer  
peil = pd.read\_excel("peilingwijzer.xlsx")  
print peil.shape

%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('.xml'):  
with open(path+file):  
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





```
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_peilingen.head()

#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.xml'):
    tag = elem.tag
    value = elem.text
    if tag == '{urn:oasis:names:tc:evs:schema:eml}Gender':
        if value == 'female':
            vrouwen += 1
        else:
            mannen += 1
    elem.clear()
print ('aantal mannen: ', mannen)
print ('aantal vrouwen: ', vrouwen)
print ('totaal: ', (mannen + vrouwen))

#Haal alle variabelen uit het dataframe:
ijfers_peilingen.describe()

#Haal het aantal rijen en kolommen zien:
ijfers_peilingen.shape

#Haal alle waarden op van een specifieke kolom (zetels, zetelslaag, zetelshoog, percentage laag, percentage hoog):
ijfers_peilingen.sum()
ijfers_peilingen['kolom'].sum()

#Verschil tussen twee kolommen in nieuwe kolom (tussen percentage laag/percentage hoog - zetelslaag/zetelshoog)
ijfers_peilingen['Verschil'] = (cijfers_peilingen['kolom'] - cijfers_peilingen['kolom'])
ijfers_peilingen.head(10)

#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

#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()

#Welke partijen hebben boven de 25 zetels of onder bijvoorbeeld (<?)
etels = cijfers_peilingen.loc[(cijfers_peilingen['Zetels'] > 25)
etels

#Sorteer dataframe op bepaalde column - ascending = 'alse (boven naar onder)
ijfers_peilingen_new = cijfers_peiling.sort_values(by='Zetels', ascending=False)
ijfers_peiling_new.head()

#Hoeveel partijleden zitten er in de dataset van de PVV?
en(cijfers_peilingen[cijfers_peilingen.Partij=='PVV'])

#Hoe lang is de kolom partij?
en(cijfers_peilingen.Partij)

#Hoe veel unieke/dubbele partijen zitten er in deze kolom?
en(set(cijfers_peilingen.Partij))

#Hoogste aantal zetels
ijfers_peilingen[cijfers_peilingen['ZetelsHoog']==ijfers_peilingen.ZetelsHoog.max()]

#Sorteer de partijen met het hoogste percentage?
ijfers_peilingen = cijfers_peilingen.set_index('partij')

#percentage_new = cijfers_peilingen.sort_values(by='Percentage', ascending=False)
percentage_new.head()

#lot de partij en het percentage - bar
percentage_new['Percentage'].plot(kind='bar', title='Hoogste percentage')

#lot de partij en het percentage - pie
percentage_new['Percentage'].plot(kind='pie', title='Hoogste percentage plot')

#lot de partij en het percentage - pie
percentage_new['Percentage'].plot(kind='pie', title='Hoogste percentage plot')

#lot de partij en het percentage - gragic
percentage_new['Percentage'].plot(title='Hoogste Percentage plot')

#open een excel bestand en laat de eerste 5 zien?
results_longitudinal = pd.read_excel('Results_Longitudinal.xlsx')
results_longitudinal.head()

#ML Files openen.
import xml.etree.ElementTree as ET
tree_new = ET.parse('bestandnaam')
root = tree_new.getroot()

#ML Files openen.
root.tag

#ML Files openen.
root.attrib

#ML Files openen.
for tag in tree_new.iter():
    print (tag)

#Hoeveel mannen en vrouwen zijn in kandidatenlijst van Amsterdam?
mannen = 0
vrouwen = 0
for event, elem in ET.iterparse('Kandidatenlijsten_TK2017_Amsterdam.xml'):
    tag = elem.tag
    value = elem.text
    if tag == '{urn:oasis:names:tc:evs:schema:eml}Gender':
        if value == 'female':
            vrouwen += 1
        else:
            mannen += 1
    elem.clear()
print ('aantal mannen: ', mannen)
print ('aantal vrouwen: ', vrouwen)
print ('totaal: ', (mannen + vrouwen))

#Hoeveel mannen en vrouwen zijn in kandidatenlijst van Amsterdam?
STATISTICS
These can all be applied to a series as well.
df.describe() - Summary statistics for numerical columns df.mean() - Return t
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.

#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.xml'):
    tag = elem.tag
    value = elem.text
    if tag == '{urn:oasis:names:tc:evs:schema:eml}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:eml}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)
df = pd.DataFrame({'partij':partij, 'voornaam':voornaam, 'initialen':initialen,
'achternaam':achternaam,
'geslacht':gender, 'woonplaats':woonplaats})
df.head(5)

#Print het aantal mannen en het aantal vrouwen?
print (len(df[df.geslacht=='man']))
print (len(df[df.geslacht=='vrouw']))

#Print een lijst van alle unieke partijen op de kandidatenlijst?
print (df['partij'].unique())
len(set(df['partij']))

#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')

#Bepaalde partijen met een bepaald geslacht
df.loc[(df['partij'] == 'VVD') & (df['geslacht'] == 'man')].head()

#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'])
aantal_procenten_man = (mannen_VVD * 100) / totaal_VVD
print (aantal_procenten_man)

#Plot de verhouding man vrouw van de VVD?
verhouding = df['geslacht'].loc[(df['partij'] == 'VVD')].value_counts()
verhouding.plot(title='Verhouding man vrouw VVD')

#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')

#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')

#Zoek een Kandidaat op uit een bepaalde woonplaats
df.loc[(df['woonplaats'] == 'Amsterdam')]

#Kandidaat uit Amsterdam die in DENK zit
df.loc[(df.woonplaats == 'Amsterdam') & (df.partij == 'DENK')]

#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'.

#Geef een lijst met partijen die een 'naam' (Peter) als kandidaat hebben?
df['partij'].loc[(df['voornaam'] == 'naam')].unique()

#Zoek naar kandidaten waar achternaam met een K begint
df[(df['voornaam', 'achternaam', 'partij']).loc[(df['achternaam'].str.startswith('K'))].head()

#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

#Tel vervolgens de gegeven conditie binnen dit nieuwe Dataframe
len(df_vvd.loc[(df_vvd['geslacht'] == 'man')])

#Plotten van het dataframe
df_vvd['geslacht'].value_counts().plot(kind='pie' title='aantal mannen en vrouwen in de kamer VVD')

#Wat is de verdeling van de zetels volgens de peilingwijzer? Neem hierbij de kolom 'zetels' en plot dit in een staafdiagram tegen de verschillende partijen op volgorde van de hoeveelheid zetels met een passende titel.
df_peilingwijzer.sort_values(by='Zetels', ascending = False).plot('Partij', 'Zetels', kind = 'bar', title = 'Verdeling Zetels volgens Peilingwijzer')

#Maak een kruistabel waarin je het gender tegen de woonplaats plaatst. Laat de totalen zien zet en laat de eerste 5 regels van deze kruistabel zien. Het resultaat moet er als volgt uitzien:
cross_gender_woonplaats = pd.crosstab(df_kandidatenlijst_adam.woonplaats, df_kandidatenlijst_adam.gender, margins = True)
print cross_gender_woonplaats.sort values(by = 'All', ascending = False).head(5)

#Maak een kruistabel waarin je het gender tegen de woonplaats plaatst. Laat de totalen zien zet en laat de eerste 5 regels van deze kruistabel zien. Het resultaat moet er als volgt uitzien:
cross_gender_woonplaats = pd.crosstab(df_kandidatenlijst_adam.woonplaats, df_kandidatenlijst_adam.gender, margins = True)
print cross_gender_woonplaats.sort values(by = 'All', ascending = False).head(5)

#Maak een horizontale staafdiagram met de verdeling van de 20 meest populaire woonplaatsen van alle kandidaten van de kandidatenlijst van Amsterdam. Doe dit ook netjes op volgorde en geef de grafiek een passend titel.
counter = df_kandidatenlijst_adam.woonplaats.value_counts()
Value_counts sorteert
counter = counter[:20]
Maak van deze verhouding mannen op vrouwen ook een pie chart met een passende titel.
counter = df_kandidatenlijst_adam.gender.value_counts()
counter.plot(kind = 'pie', title = 'Verhouding mannen tot vrouwen binnen de PVV')

#En hoeveel mannen? Wat is de relatieve verhouding van mannen tot vrouwen?
Rond dit af op 2 decimalen achter de komma.
vrouwen = len(df_kandidatenlijst_adam[(df_kandidatenlijst_adam.partij == 'PVV (Partij voor de Vrijheid)') & (df_kandidatenlijst_adam.gender == 'v')])
mannen = len(df_kandidatenlijst_adam[(df_kandidatenlijst_adam.partij == 'PVV (Partij voor de Vrijheid)') & (df_kandidatenlijst_adam.gender == 'm')])
print vrouwen, 'vrouwen en', mannen, 'mannen op de kandidatenlijst van Amsterdam.

#Beperk je bij de kandidatenlijst tot de mensen die kiesbaar zijn voor de kiesring Amsterdam. Als de PVV 20 zetels zouden behalen, hoeveel vrouwen zouden er dan in de kamer komen volgens de kandidatenlijst van Amsterdam?
df_kandidatenlijst_adam = df_kandidatenlijst(df_kandidatenlijst.kieskrin == 'Amsterdam')
df_kandidatenlijst_adam_pv20 = df_kandidatenlijst_adam[(df_kandidatenlijst_adam.partij == 'PVV (Partij voor de Vrijheid)')].head(20)
print len(df_kandidatenlijst_adam_pv20[df_kandidatenlijst_adam_pv20.geslacht == 'v'])
index == 'v'. Vrouwen
```

```
#Hoeveel mannen en vrouwen zijn in kandidatenlijst van Amsterdam?
mannen = 0
vrouwen = 0
for event, elem in ET.iterparse('Kandidatenlijsten_TK2017_Amsterdam.xml'):
    tag = elem.tag
    value = elem.text
    if tag == '{urn:oasis:names:tc:evs:schema:eml}Gender':
        if value == 'female':
            vrouwen += 1
        else:
            mannen += 1
    elem.clear()
print ('aantal mannen: ', mannen)
print ('aantal vrouwen: ', vrouwen)
print ('totaal: ', (mannen + vrouwen))

#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.xml'):
    tag = elem.tag
    value = elem.text
    if tag == '{urn:oasis:names:tc:evs:schema:eml}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:eml}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)
df = pd.DataFrame({'partij':partij, 'voornaam':voornaam, 'initialen':initialen,
'achternaam':achternaam,
'geslacht':gender, 'woonplaats':woonplaats})
df.head(5)

#Print het aantal mannen en het aantal vrouwen?
print (len(df[df.geslacht=='man']))
print (len(df[df.geslacht=='vrouw']))

#Print een lijst van alle unieke partijen op de kandidatenlijst?
print (df['partij'].unique())
len(set(df['partij']))

#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')

#Bepaalde partijen met een bepaald geslacht
df.loc[(df['partij'] == 'VVD') & (df['geslacht'] == 'man')].head()

#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'])
aantal_procenten_man = (mannen_VVD * 100) / totaal_VVD
print (aantal_procenten_man)

#Plot de verhouding man vrouw van de VVD?
verhouding = df['geslacht'].loc[(df['partij'] == 'VVD')].value_counts()
verhouding.plot(title='Verhouding man vrouw VVD')

#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')

#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')

#Zoek een Kandidaat op uit een bepaalde woonplaats
df.loc[(df['woonplaats'] == 'Amsterdam')]

#Kandidaat uit Amsterdam die in DENK zit
df.loc[(df.woonplaats == 'Amsterdam') & (df.partij == 'DENK')]

#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'.

#Geef een lijst met partijen die een 'naam' (Peter) als kandidaat hebben?
df['partij'].loc[(df['voornaam'] == 'naam')].unique()

#Zoek naar kandidaten waar achternaam met een K begint
df[(df['voornaam', 'achternaam', 'partij']).loc[(df['achternaam'].str.startswith('K'))].head()

#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

#Tel vervolgens de gegeven conditie binnen dit nieuwe Dataframe
len(df_vvd.loc[(df_vvd['geslacht'] == 'man')])

#Plotten van het dataframe
df_vvd['geslacht'].value_counts().plot(kind='pie' title='aantal mannen en vrouwen in de kamer VVD')
```

#Hoeveel mannen en vrouwen zijn in kandidatenlijst van Amsterdam?

STATISTICS

These can all be applied to a series as well.

df.describe() - Summary statistics for numerical columns df.mean() - Return t  
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.

Wat is de verdeling van de zetels volgens de peilingwijzer? Neem hierbij de kolom 'zetels' en plot dit in een staafdiagram tegen de verschillende partijen op volgorde van de hoeveelheid zetels met een passende titel.

df\_peilingwijzer.sort\_values(by='Zetels', ascending = False).plot('Partij', 'Zetels', kind = 'bar', title = 'Verdeling Zetels volgens Peilingwijzer')

Maak een kruistabel waarin je het gender tegen de woonplaats plaatst. Laat de totalen zien zet en laat de eerste 5 regels van deze kruistabel zien. Het resultaat moet er als volgt uitzien:

```
cross_gender_woonplaats = pd.crosstab(df_kandidatenlijst_adam.woonplaats, df_kandidatenlijst_adam.gender, margins = True)
print cross_gender_woonplaats.sort values(by = 'All', ascending = False).head(5)
```

Maak een kruistabel waarin je het gender tegen de woonplaats plaatst. Laat de totalen zien zet en laat de eerste 5 regels van deze kruistabel zien. Het resultaat moet er als volgt uitzien:

```
cross_gender_woonplaats = pd.crosstab(df_kandidatenlijst_adam.woonplaats, df_kandidatenlijst_adam.gender, margins = True)
print cross_gender_woonplaats.sort values(by = 'All', ascending = False).head(5)
```

Maak een horizontale staafdiagram met de verdeling van de 20 meest populaire woonplaatsen van alle kandidaten van de kandidatenlijst van Amsterdam. Doe dit ook netjes op volgorde en geef de grafiek een passend titel.

```
counter = df_kandidatenlijst_adam.woonplaats.value_counts()
Value_counts sorteert
counter = counter[:20]
Maak van deze verhouding mannen op vrouwen ook een pie chart met een passende titel.
counter = df_kandidatenlijst_adam.gender.value_counts()
counter.plot(kind = 'pie', title = 'Verhouding mannen tot vrouwen binnen de PVV')
```

En hoeveel mannen? Wat is de relatieve verhouding van mannen tot vrouwen?

Rond dit af op 2 decimalen achter de komma.

```
vrouwen = len(df_kandidatenlijst_adam[(df_kandidatenlijst_adam.partij == 'PVV (Partij voor de Vrijheid)') & (df_kandidatenlijst_adam.gender == 'v')])
mannen = len(df_kandidatenlijst_adam[(df_kandidatenlijst_adam.partij == 'PVV (Partij voor de Vrijheid)') & (df_kandidatenlijst_adam.gender == 'm')])
print vrouwen, 'vrouwen en', mannen, 'mannen op de kandidatenlijst van Amsterdam.
```

Beperk je bij de kandidatenlijst tot de mensen die kiesbaar zijn voor de kiesring Amsterdam. Als de PVV 20 zetels zouden behalen, hoeveel vrouwen zouden er dan in de kamer komen volgens de kandidatenlijst van Amsterdam?

```
df_kandidatenlijst_adam = df_kandidatenlijst(df_kandidatenlijst.kieskrin == 'Amsterdam')
df_kandidatenlijst_adam_pv20 = df_kandidatenlijst_adam[(df_kandidatenlijst_adam.partij == 'PVV (Partij voor de Vrijheid)')].head(20)
print len(df_kandidatenlijst_adam_pv20[df_kandidatenlijst_adam_pv20.geslacht == 'v'])
index == 'v'. Vrouwen
```

at. shape, at. into, at. newserie pa. omruu  
df.iloc[:, ] voor rij, df[col] voor kolom, df.sort\_values(col1) of [col1, [col2]]

```
soup = BeautifulSoup(open('Kandidatenlijst_TK2017_Amsterdam.eml.xml').read())
parties = soup['affiliation']
party_names = [p.text for p in soup.findAll('registeredname')]
genderlist = {p.find('registeredname').text : [g.text for g in p.findAll('gender')] for p in parties}
```

```
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')
```

max aantal vrouwen p partij

```
#PVV= AantalVrouwenPerPartijPerAantalGewonnenZetels['PVV (Partij voor de Vrijheid)']
for p in genderlist:
    titel = p + ': aantal vrouwen per aantal gewonnen zetels'
```

```
pd.DataFrame.from_dict(AantalVrouwenPerPartijPerAantalGewonnenZetels[p], orient='index').plot(kind='bar', title=titel);
```

# translate the party names so that we can join

```
translate = {u'50PLUS':u'50PLUS',
             u'CDA':u'CDA',
             u'CU':u'ChristenUnie',
             u'D66':u'Democraten 66 (D66)',
             u'Denk':u'DENK',
             u'Fvd':u'Forum voor Democratie',
             u'GL':u'GROENLINKS',
             u'PVV':u'PVV (Partij voor de Vrijheid)',
             u'PvdA':u'Partij van de Arbeid (P.v.d.A.)',
             u'PvdD':u'Partij voor de Dieren',
             u'SGP':u'Staatkundig Gereformeerde Partij (SGP)',
             u'SP':u'SP (Socialistische Partij)',
             u'VNL':u'VNL (VoorNederland)',
             u'VVD':u'VVD'}
```

# Of bereken tot welke peiling = pd.read\_excel()
peiling2 = part.no [peiling.no. Partij == 'VVD']
# kolom toevoegen
peiling2['zetels'] = peiling2['zetels'] + peiling2['zetels']

# Voor de eerste N mensen op de lijst zet N=80 (is maximum) voor de hele lijst

```
N=50
gendercount = {p: Counter(genderlist[p][:N]) for p in genderlist}
# print gendercount
gdf = pd.DataFrame.from_dict(gendercount, orient='index').fillna(0).astype(int)
# min_lijst_lengte = (gdf.sum(axis=1)).min()
# print min_lijst_lengte
# Voeg de ratio toe
gdf['Fratio'] = ((gdf.female / gdf.sum(axis=1)) * 100).round()
gdf.sort_values('Fratio', 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('Fratio')[['male', 'female', 'Fratio']]
```

aantal m/v p partij

hoeveel vrouwen bij zetels langs COA

```
pvda = genderlist['Partij van de Arbeid (P.v.d.A.)']
sum([(1/np.log2(2+i)) for i in range(len(pvda)) if pvda[i]=='female'])
```

# Peilingwijzer

```
laatstepeiling = Cijfers_Peilingwijzer.xlsx
allepeilingen = Results_Longitudinal.xlsx
laatstedef = pd.read_excel(laatstepeiling)
alldf = pd.read_excel(allepeilingen)
print alldf.shape
alldf.tail()
laatstedef.index = laatstedef.Partij
laatstedef.rename(translate, inplace=True)
```

# Dit is een handige file met de zetel aantallen
alleZdf = pd.read\_csv('https://d1bjgq97if6urz.cloudfront.net/Public/Peilingwijzer/Last/Results\_DyGraphs\_Seats.csv', index\_col='Date')

```
print alleZdf.shape
print (alleZdf.sum(axis=1) == 150).mean() # check if all add to 150
alleZdf = alleZdf.applymap(lambda s: int(s.split(':')[1])) # want (ci low; mean; ci high)
# now translate to our names
alleZdf = alleZdf.T
alleZdf.rename(translate, inplace=True)
alleZdf = alleZdf.T
alleZdf.tail()
```

```
alleZdf.rename(columns = {'PP': 'Piratenpartij'}, inplace=True)
for p in alleZdf.columns:
    alleZdf[p] = alleZdf[p].apply(lambda z: AantalVrouwenPerPartijPerAantalGewonnenZetels[p][z])
```

```
fem = alleZdf.plot(figsize=(18,7), title='Aantal vrouwen in de Tweede Kamer door de tijd, per partij.\n Gebaseerd op de Peilingwijzer.');
```

```
som = alleZdf.T.sum()
print som.describe()
fem = alleZdf.T.sum().plot(figsize=(18,7), title='Aantal vrouwen in de Tweede Kamer door de tijd.\n Gebaseerd op de Peilingwijzer.');
```

```
df = laatstedef.join(gdf)
df
# Turn percentages into number of zetels
def Perc2Zetels(p):
    return round(p * (150/100))
def Ratio2Zetels(p):
    return round(p * 150)
print df.Percentage.apply(Perc2Zetels).sum()
ef = gdf.join(alldf.T.rename(translate)).dropna().drop(['u'male', 'u'female', 'u'Fratio'], axis=1)
ZetelsDoorDeTijd = ef.applymap(Ratio2Zetels)
# ZetelsDoorDeTijd.columns = pd.to_datetime(ZetelsDoorDeTijd.columns)
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()).plot(figsize=(18,7), title='Percentage vrouwen in de Tweede Kamer door de tijd.\n Gebaseerd op de Peilingwijzer.');
```

```
len([1 for mp in genderlist['VVD'][:df.ZetelsLaag['VVD']] 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')]:
    Fc = 'Fem'+c
# compute aantal vrouwen
Fem = pd.DataFrame.from_dict({p: aantalvrouwen(p,c,df) for p in df.index}, orient='index')
Fem.columns = [Fc]
# voeg toe aan df
df = df.join(Fem)
# bereken de ratio
df['FratioPeiling'+c] = (df[Fc] / df[c]).fillna(0)
```

```
# Aantal vrouwen in de kamer
VrouwenPeilingwijzer = df.sort_values('Zetels', ascending=False)[['u'Zetels', 'u'FemZetels', 'u'FratioPeilingZetels']]
print "Op basis van deze peiling komen er %s vrouwen in de kamer. Dat is %s procent" % (VrouwenPeilingwijzer.FemZetels.sum(),
```

round((vrouwen peilingwijzer.FemZetels.sum() / vrouwen peilingwijzer.Zetels.sum()) \* 100)
print vrouwen peilingwijzer
df.sort\_values(ascending=False, axis=1, columns=['Partij', 'Zetels', 'kind='bar', title='verdeling zetels peilingwijzer'])

```

import pandas as pd, import numpy as np, import seaborn as sn, %matplotlib inline, from lxml import etree
from collections import Counter, import matplotlib.pyplot as plt, from zipfile import ZipFile
from urllib import urlretrieve, from tempfile import mktemp
#Download Data and open zip
filename = 'Kandidatenlijsten.zip'
theurl =
'https://www.kiesraad.nl/binaries/kiesraad/documenten/rapporten/2017/2/definitieve-kandidatenlijsten-tk-2017/definitieve-kandidatenlijsten-tk-2017/Kandidatenlijsten_EML_TK2017-20170213.zip'
name, hdrs = urlretrieve(theurl, filename)
thefile=ZipFile(filename)
thefile.extractall()
thefile.close()
df_peilingwijzer =
pd.read_excel('https://d1bjgq97if6urz.cloudfront.net/Public/Peilingwijzer/Last/Cijfers_Peilingwijzer.xlsx')
xml_data = 'Kandidatenlijsten_EML_TK2017-20170213/Kandidatenlijsten_TK2017_Amsterdam.eml.xml'
#Parse XML
def xml2df(xml_data):
    s1 = '{urn:oasis:names:tc:evs:schema:eml}'
    s2 = '{urn:oasis:names:tc:ciq:xsd:schema:xNL:2.0}'
    s3 = '{urn:oasis:names:tc:ciq:xsd:schema:xAL:2.0}'
    columns = ['partij', 'nummer', 'naam kandidaat', 'voorletters', 'voornaam', 'gender', 'woonplaats']
    L = []
    doc = etree.parse(xml_data)
    root = doc.getroot().find(s1 + 'CandidateList/'+s1+'Election/'+s1+'Contest')
    for child in root.iter(s1+'Affiliation'):
        partij = child.find(s1+'AffiliationIdentifier/'+s1+'RegisteredName').text
        for child_2 in child:
            for candidate in child_2.iter(s1+'Candidate'):
                ID = candidate.find(s1+'CandidateIdentifier').get('Id')
                nameinfo = candidate.find(s1+'CandidateFullName/'+s2+'PersonName')
                initial = nameinfo.find(s2+'NameLine').text
                firstname = nameinfo.find(s2+'FirstName').text
                try:
                    prefix = nameinfo.find(s2+'NamePrefix').text + ' '
                except:
                    prefix = ''
                lastname = prefix + nameinfo.find(s2+'LastName').text
                gender = candidate.find(s1+'Gender').text[0]
                try:
                    living = candidate.find(s1+'QualifyingAddress/'+s3+'Locality/'+s3+'LocalityName').text
                except:
                    country = candidate.find(s1+'QualifyingAddress/'+s3+'Country/'+s3+'CountryNameCode').text
                    living = candidate.find(s1+'QualifyingAddress/'+s3+'Country/'+s3+'Locality/'+s3+'LocalityName').text+'('+country+')'
                L.append([partij, ID, lastname, initial, firstname, gender, living])
    df_kandidatenlijst = pd.DataFrame(L)
    df_kandidatenlijst.columns = columns
    return df_kandidatenlijst

df_kandidatenlijst_adam.partij.replace(dic, inplace=True)
df_join = df_kandidatenlijst_adam.join(df_peilingwijzer.set_index('Partij'), on = 'partij')
df_join.update(df_join[['Percentage', u'PercentageLaag', u'PercentageHoog', u'Zetels', u'ZetelsLaag',
u'ZetelsHoog']].fillna(0))
namesamb.loc[namesamb.F > namesamb.M, 'winning gender'] = 'F'

df_peilingwijzer.sort_values(by='Zetels', ascending = False).plot('Partij', 'Zetels', kind = 'bar', title =
'Verdeling Zetels volgens Peilingwijzer')

df=df_kandidatenlijst_adam.pivot_table(index=["partij", "gender"],aggfunc='count').dropna()
df=df.drop(['naam kandidaat', 'voorletters', 'voornaam', 'woonplaats'],axis=1).rename(columns={'nummer': 'aantal'})
df.unstack('gender').fillna(0).plot(kind='barh',title= 'verhouding man/vrouw' ,stacked = True, figsize=(20,10))

cross_gender_woonplaats = pd.crosstab(df_kandidatenlijst_adam.woonplaats, df_kandidatenlijst_adam.gender, margins =
True)

for f in os.listdir('Directory')

df_peilingwijzer.shape() df_peilingwijzer.describe() alist.abs() alist.sorted()

for partij in df_peilingwijzer.Partij:
    frame = df_peilingwijzer_long[['Datum', partij, partij+'.low', partij+'.high']] #df_peilingwijzer_Long['Datum'] =
df_peilingwijzer_Long.index
    frame.columns = ['Datum', 'zetels', 'zetels_laag', 'zetels_hoog']
    frame['Partij'] = partij # Add the year column with constant value
    pieces.append(frame)
df_peilingwijzer_long_nieuw = pd.concat(pieces, ignore_index=True)

```

## Spiekbrief Datascience hertentamen

```
import pandas as pd
import os
from lxml import etree
import xml.etree.ElementTree as ET
import seaborn as sns
%matplotlib inline
import pandas as pd
import re
import numpy as np
from lxml import etree
#from bz2file import BZ2File
import codecs
import nltk
from collections import defaultdict
from itertools import combinations
# ideal for creating all possible
pairs that out can make out of a set
from __future__ import division
from math import sqrt
import collections
import matplotlib.pyplot as plt
import os
from bs4 import BeautifulSoup
import requests
import matplotlib.pyplot as plt
```

### Bestanden openen:

```
pd.read_excel('peilingwijzer.xlsx')
```

### Select:

```
df[df['column'] == 'value']
```

### Sorteren:

```
df.sort_values('column')
value.counts() -> bij een kolom
aantal waarden tellen
```

**loc:** gebruik .loc[(voorwaarden) & / | (voorwaarden)] vb:  
df.loc[(df['partij'] == 'VVD')] en  
df.loc[(df['partij'] == 'VVD') &  
(df['geslacht'] == 'man')]

*Print out the top 20 most ambiguous names. Take those names with a "log-value" between -0.1-0.1 and 0.10.1 and sort them reversely on the total number of babies with that name.*

```
ambiguous = ef[(ef.logratio < 0.1) &
(ef.logratio > -
0.1)].sort_values('All',
ascending=False)
ambiguous.head(20)
OF
names = total_baby.query('log_2
<= 0.1 & -0.1 <= log_2')[20]
```

### maten + wiskundige dingen df

```
= .shape & .describe
```

### Unieke tags xml:

```
tags = []
for event, elem in -
ET.iterparse('Kandidatenlijsten_TK2
017_Amsterdam.eml.xml'):
    tag = elem.tag
    tags.append(tag)
    elem.clear()
set(tags) #als je set weghaalt heb
je alle tags uit het file
hoeveel mannen/vrouwen in alle
files:
path =
'Kandidatenlijsten_EML_TK2017-20
170213/'
for _file in os.listdir(path):
    if not -
_file.endswith('.xml'):continue
    malecount = 0
    femalecount = 0
    with open(path+_file): ←
        for event, elem in
ET.iterparse(path +_file):
            tag = elem.tag
            value = elem.text
            if tag ==
"{urn:oasis:names:tc:evs:schema:e
ml}Gender" :
                if value == 'male':
                    malecount += 1
                totalmalecount += 1
            else:
                femalecount += 1
            elem.clear() # discard the
element
            print _file.strip('.eml.xml'),
malecount, femalecount,
(malecount+femalecount)
```

### Dataframe van XML:

```
partij = []
voornaam = []
initialen = []
achternaam = []
gender = []
woonplaats = []
nummer = []

for event, elem in
ET.iterparse('Kandidatenlijsten_TK2
017_Amsterdam.eml.xml'):
    tag = elem.tag
    value = elem.text
    if tag ==
"{urn:oasis:names:tc:evs:schema:e
ml}RegisteredName":
        partij_val = value
        if tag ==
"{urn:oasis:names:tc:ciq:xsdscem
a:xNL:2.0}NameLine":
            initialen.append(value)
        if tag ==
"{urn:oasis:names:tc:ciq:xsdscem
a:xNL:2.0}FirstName":
            partij.append(partij_val)
            voornaam.append(value)
```

```
if tag ==
"{urn:oasis:names:tc:ciq:xsdscem
a:xNL:2.0}LastName":
    achternaam.append(value)
    if tag ==
"{urn:oasis:names:tc:evs:schema:e
ml}Gender":
        if value == 'male':
            gender_val = 'man'
        else:
            gender_val = 'vrouw'
        gender.append(gender_val)
    if tag ==
"{urn:oasis:names:tc:ciq:xsdscem
a:xAL:2.0}LocalityName":
        woonplaats.append(value)
```

```
df_amsterdam =
pd.DataFrame({'partij':partij,
'voornaam':voornaam,
'initialen':initialen,
'achternaam':achternaam,
'geslacht':gender,
'woonplaats':woonplaats})
df_amsterdam =
df_amsterdam[['partij', 'voornaam',
'initialen', 'achternaam', 'geslacht',
'woonplaats']]
```

*stel dat de PVV 21 zetels haalt,
hoeveel vrouwelijke PVV'ers komen
er dan in de Kamer? .*

### #ONELINER

```
len(df.loc[(df['partij'] == 'PVV (Partij
voor de
Vrijheid)']).head(peilingwijzer['Zetel
s'].loc[(peilingwijzer['Partij'] ==
'PVV']).values.item()).loc[(df['geslac
ht'] == 'vrouw')])
#MAAK EEN DATAFRAME MET X
AANTAL NAMEN VAN DE
KANDIDATENLIJST VAN EEN
SPECIFIEKE PARTIJ
df_pvv = df.loc[(df['partij'] == 'PVV
(Partij voor de
Vrijheid)']).head(peilingwijzer['Zetel
s'].loc[(peilingwijzer['Partij'] ==
'PVV']).values.item())
```

#VIND ALLE KANDIDATEN DIE IN.  
df\_woonplaats =  
df['woonplaats'].loc[df['woonplaats']
.isin(['Breda', 'Tilburg',
'Eindhoven'])]

### Random

```
crosstab (wat , tegenover wat,
margins = true) = kruistabel
axis = is kolommen
#PLOT HET DATAFRAME IN EEN
STACKED BARPLOT
ax = manvrouwdf.plot(x='partij',
y='mannen', kind = 'bar')
manvrouwdf.plot(x='partij',
y='vrouwen', kind = 'bar',
color='red', ax=ax)
```

```

# GENERAL #####
# Is gelijk aan: == is niet gelijk aan: !=
round(2.05467, 2) # rond een getal af op 2 decimalen
float(2) # maak kommagetal, dus 2.00 ipv 2
list.append('item 4') # add item to List
list[:2] # selection of List
print 'Henk is %s jaar oud' % str(leeftijd)
if 'Henk' in mensen: # check of Henk voorkomt in mensen

# CREATE DATAFRAME #####
# van string
df = pd.DataFrame(data, index=[1,2,3,4],
columns=['date', 'b', 'c', 'd'])

# van file (read_excel/table/csv etc.)
df = pd.read_excel('http://url.nl/file.xls')
df = pd.read_csv('./output.csv',
sep=',', # of bijvoorbeeld it voor tab
skiprows=1, # sla de eerste rij van file over
header=None, # neem niet header van file over
names=['name', 'F', 'M'] # voeg kolomnamen toe
)

# SUMMARIZE DATAFRAME #####
df.head(10) # of df.tail(10)
df.describe() # algemene info over df
df.shape # aantal kolommen en aantal rijen
len(df) # number of rows in df
len(df.clnm_name.unique()) # no distinct values

df.clnm_name.value_counts() # count number of rows for each
unique value (voor een specifieke kolom)
df.apply(pd.value_counts) # voer de functie van hierboven uit
op alle kolommen van de df

# DOING STUFF WITH DATAFRAME #####
df.sort_values('clnm_name', ascending=False) # sorteer rijen
df.sort_index() # sort index of a df
df.reset_index() # reset index of df to row numbers, moving
the existing index to a column
new_df = df.set_index('date') # set column as index
df.columns=['name', 'F', 'M'] # set column names
df.groupby(by='clnm_name') # return groupby object, grouped
by values in column
df.clnm_name = df.clnm_name.astype(int) # change datatype of
values in a column
df.transpose() # draai df kwartslag/maak columns van rijen
df.clnm_name.min() df.clnm_name.max() # min of max van column
df.sum(axis=1) # som van waarden (kan ook axis=0)
df.count()

# DELETE DATA #####
df = df.drop(['Length', 'Height'], axis=1) # delete columns
df.drop_duplicates() # remove duplicate rows

```

```

# STUFF WITH MISSING DATA #####
df = df.dropna() # drop rows with any column having NA/null
df = df.dropna(axis=1) # drop columns with any column having
NA/null
df.dropna(subset=['column_name']) # drop rows with NA/null in
specific column
df = df.fillna(value) # replace all NA/null with value

# SELECTING DATA FROM DATAFRAME #####
# select rows
df[(df.c >= 3) & (df.b < 5)] # extract rows with condition
df.loc[df.clnm_name == 'John'] # extract rows with condition
df[1:3] # select rows based on position
df.loc[df['a'] > 10, ['a', 'c']] # select rows with condition
and only in specified columns

# select columns
df['width'] # select single column
df[['width', 'length', 'species']] # select multiple columns
df.loc[:, 'width': 'length'] # select all columns between those
df.iloc[:, [1,2,5]] # select columns in positions 1, 2 and 5

df.columns # get the column names (df.column[12:50] voor een
selectie van een deel van de kolommen)
list(df.columns.values)

# CREATE A NEW COLUMN #####
df['Ratio'] = df.M + df.F

# MAKING PLOTS #####
df.plot(kind=kind) # create plot (kind: bar, barh, etc.)
# HISTOGRAM - Maak een "horizontaal histogram" waarin je voor
elke partij aangeeft hoeveel van dit soort stukken die partij
heeft ingediend.
df.clnm_name.value_counts().sort_values(ascending=True).plot(
'barh')
# SCATTER PLOT
df.plot.scatter(x='column 1', y='column 2', title='')
# KRUISTABEL
pd.crosstab(df.proposaltype, df.votetype, margins=True)

# OPDRACHTEN #####
# Maak een nieuwe df (im) van de selectie waarbij de partij
van indiener (kolom: authorparty) gelijk is aan die van
medeindiener(s) (kolom: supporterparties).
im = rutte2.loc[(rutte2.authorparty ==
rutte2.supporterparties) & (rutte2.authorparty != '[]')]

# Hoe vaak heeft Barry Madlener voor gestemd bij hoofdelijke
stemmingen? Kijk hiervoor in de pro kolom.
len([item[1] for item in df.pro.iteritems() if 'Barry
Madlener' in item[1]])

# Check of je resultaat gelijk is aan wat het moet zijn
print (string == 100) # dit print True of False

```

```

# PRETTIFY SOUP #####
print soup.prettify()

# XML #####
import urllib
import zipfile
import os
from collections import defaultdict
from bs4 import BeautifulSoup

urllib.urlretrieve
("https://data.openstate.eu/dataset/85aaef1d-9084-4d3f-a6d1-
1566a9538b6/resource/dd639f9f-4035-4e17-b09b-33fc388e9780/dow
nload/kandidatenlijsteneimtk2017-20170213.zip", "files.zip")
with zipfile.ZipFile("files.zip", "r") as zipf:
zipf.extractall("")

directory = "Kandidatenlijsten_EML_TK2017-20170213/"
candidates_list = []

# Loop through all the xml files in the directory
for filename in os.listdir(directory):
if filename.endswith(".xml"):
f = open(directory + filename, "r")
soup = BeautifulSoup(f.read(), 'lxml')

kieskring = soup.find('contestname').contents[0]

# Loop through all affiliations (partijen)
for affiliation in soup.findAll('affiliation'):
partij = affiliation.find('registeredname').contents[0]

# Loop through all candidates
for c in affiliation.findAll('candidate'):
c_id = c.find('candidateidentifier').attrs['id']
c_initials = c.find('ns3:nameline').contents[0]
c_wp = c.find('ns2:localityname').contents[0]

try: c_lastname=c.find('ns3:nameprefix').contents[0] +
c.find('ns3:lastname').contents[0]
except: c_lastname=c.find('ns3:lastname').contents[0]

c_firstname = c.find('ns3:firstname')
try: c_firstname = c_firstname.contents[0]
except: print 'No firstname'

c_gender = c.find('gender')
try: c_gender = c_gender.contents[0]
except: print 'No gender'

candidates_list.append([kieskring, partij, c_id,
c_lastname, c_initials, c_firstname, c_gender, c_wp])

f.close()

df = pd.DataFrame(candidates_list, columns=['kieskring',
'partij', 'nummer', 'achternaam', 'voorletters', 'voornaam', 'gend
er', 'woonplaats'])

```

### Peilingwijzer namen gelijk 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)'
```

### XML lezen

```
from lxml import etree  
import xml.etree.ElementTree as ET  
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"):  
            nummer = candidate.getchildren()[0].attrib["Id"]  
            voornaam = candidate.getchildren()[1][0][1].text  
            achternaam = 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 .....})  
kandidaten = pd.DataFrame(data)  
kandidaten.geslacht.replace({"male": "m", "female": "v"}, inplace=True)
```

### Dataframes samenvoegen

```
genders = kandidaten.groupby("Partij")["geslacht"].sum().to_frame().reset_index()  
compleet = peilingwijzer.merge(genders, on="Partij")
```

### Min/Max/Aantal vrouwen/mannen

```
for item in compleet.iterrows():  
    vrouwenlaag.append(item[1]["geslacht"][item[1]["ZetelsLaag"]].count("v"))  
compleet["vrouwenlaag"] = vrouwenlaag etc..
```

### Pivot table kandidatenlijst

```
pivot = kandidaten.pivot_table(index=["partij", "gender"], aggfunc="count", values="woonplaats")  
dfl = pd.DataFrame(pivot.unstack("gender"))  
PLOT: dfl.plot(kind="barh", figsize=(20,15))
```

### RATIO + TOTAAL

```
dfl["total"] =(dfl["m"]+dfl["v"])  
dfl["ratio"] =dfl["m"]/(dfl["m"]+dfl["v"])
```

### Groupby

```
g = amsterdam.groupby(["partij", "gender"])  
g.size()
```

### Meerdere files

```
for _file in os.listdir(path):  
    if not _file.endswith('xml'):continue  
    with open(path+_file):  
        tree = ET.parse(path+_file)
```

```

import pandas as pd
import os
from lxml import etree
import xml.etree.ElementTree as ET

XML naar DataFrame
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"):
            nummer = int(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, "nummer":nummer, "naam": firstname, "achternaam": lastname,
"geslacht": gender, "woonplaats":place})

kandidaten = pd.DataFrame(data)
kandidaten.geslacht.replace({"male": "m", "female": "v"}, inplace=True)

Aantal vrouwen/mannen per partij kandidatenlijst
%matplotlib inline
hi= kandidaten.pivot_table(index=["partij","geslacht"], aggfunc='count', values='woonplaats').dropna()
#hi=hi.drop(['naam','initialen','achternaam','woonplaats'],axis=1).rename(columns={'nummer':'aantal'})
hi.unstack("geslacht").sort_values('v').plot(kind="barh", figsize=(20,10))

Ratio man/vrouw
df = hi.unstack("geslacht")
df["ratio"] = df["aantal","v"]/(df["aantal","m"]+df["aantal","v"])*100
df.fillna(0).sort_values("ratio")["aantal"].plot(kind="barh", figsize=(20,10))

Aantal zetels per partij peilingenwijzer 2017
peilingenwijzer=peilingwijzer[["Partij","Zetels"]].set_index("Partij").sort_values('Zetels').plot(kind='barh')

Geslacht op volgorde per partij mergen met pijlingenwijzer en plotten
genders = kandidaten.groupby("Partij")["geslacht"].sum()
dfg = genders.to_frame()
dfg = dfg.reset_index()
compleet = peilingwijzer.merge(dfg, on="Partij")
vrouwenlaag = []
vrouwenhoog = []
vrouwen = []

for item in compleet.iterrows():
    vrouwenlaag.append(item[1]["geslacht"][:item[1]["ZetelsLaag"]].count("v"))
    vrouwenhoog.append(item[1]["geslacht"][:item[1]["ZetelsHoog"]].count("v"))
    vrouwen.append(item[1]["geslacht"][:item[1]["Zetels"]].count("v"))

compleet["vrouwenlaag"] = vrouwenlaag
compleet["vrouwenhoog"] = vrouwenhoog
compleet["vrouwen"] = vrouwen
%matplotlib inline
hoi = compleet[["Partij","vrouwen","mannen"]]
hoi = hoi[(hoi.vrouwen > 0) | (hoi.mannen > 0)]
hoi.set_index('Partij').sort_values('vrouwen').plot(kind='barh').set(xlabel='aantal m/v', ylabel='partijen')

Partijen over de jaren plotten
peilingwijzer = pd.read_excel("Results Longitudinal-2.xlsx")
[col for col in peilingwijzer if not col.endswith("low") and not col.endswith("high")]
%matplotlib inline
(peilingwijzer[[u'VVD',
u'PvdA', u'PVV', u'SP', u'CDA', u'D66', u'CU', u'GL', u'SGP', u'PvdD', u'50PLUS',
u'VNL', u'Denk', u'FvD',
u'PP']]*150).plot(figsize=(20,10)).set(xlabel='datum', ylabel='aantal zetels')

Stel VVD 26 zetels hoeveel vrouwen in kamer?
df_vvd = kandidaten[(kandidaten.Partij == 'VVD') & (kandidaten.geslacht == 'v') & (kandidaten.nummer < 26)]

Top 20 mannen namen
mannen.voornaam.value_counts()[:20].sort_values().plot('barh')

Pijlingenwijzer namen gelijk zetten
peilingwijzer = pd.read_excel("peilingwijzer.xlsx")
peilingwijzer.Partij[2] = 'Partij van de Arbeid (P.v.d.A.)'
peilingwijzer.Partij[3] = 'PVV (Partij voor de Vrijheid)'

peilingwijzer = peilingwijzer.rename(columns={'Partij':'partij'})

```

kolom.apply(lambda x:

)