

Ciencia de datos

Práctica 4. Visualización de la información con plotly

Alberto Benavides

De los datos correspondientes a los registros obtenidos de los PDFs de la Secretaría de Salud de México

```
In [104]: import pandas as pd
import plotly.plotly as py
import plotly.graph_objs as go

data1 = pd.read_csv("D:/FIME/Epidemia/Data/csvSemanales/enf.csv")
data = data1[data1.estado == "TOTAL"]
data.loc[:, 'cie'] = data['cie'].astype(str).str[0]

# https://plot.ly/python/table/
trace = go.Table(
    header=dict(
        values=list(data),
        fill = dict(color='#a1c3d1'),
    ),
    cells=dict(values=data.sample(10).T)) # https://stackoverflow.com/a/19483025

table = [trace]
py.iplot(table, filename = 'semanalesTodas')
```

Out[104]:

# m	f1	f2	f3	f4	ac1	ac2	ac3	ac4	ac5	ac6	enf	cie	cluster
247537660	83486464	937317943	670563041	684695313	408282325	974057260	935253089	714119233	104939230	169361906	LCERAS ASTRITIS UODENITI	K	9
067739572	693527338	386192015	544808031	409449043	309651978	524507871	639404751	058482943	377191427	293940901	TOXICACI OR ONZONA JIMALES	T	11
579839400	602818565	471689475	124010299	107675162	861322503	446081546	177816572	394855873	662882740	551715190	COCERCO	B	13
274729603	442769226	937317943	397286671	396185233	532406722	177974941	254266928	946332905	270881563	89670227	EATON ESIONADO N CCIDENTE E RANSPORT	V	0
567075220	445000700	126644250	000000000	000000000	000000000	000000000	000000000	000000000	000000000	000000000		A	12

EDIT CHART

y estos mismos preprocesados

```
In [105]: data2 = pd.read_csv("D:/FIME/Epidemia/Data/semanalesTodasKmeans.csv")
trace = go.Table(
    header=dict(
        values=list(data2),
        fill = dict(color='#a1c3d1'),
    ),
    cells=dict(values=data2.sample(10).T)) # https://stackoverflow.com/a/19483025

table = [trace]
py.iplot(table, filename = 'semanalesTodas')
```

Out[105]:

# m	f1	f2	f3	f4	ac1	ac2	ac3	ac4	ac5	ac6	enf	cie	cluster
247537660	83486464	937317943	670563041	684695313	408282325	974057260	935253089	714119233	104939230	169361906	LCERAS ASTRITIS UODENITI	K	9

067739572693527338386192015544808031409449041309651978524507871639404751058482943377191427293940901	TOXICACI OR ONZONA E JIMALES	T	11
5798394006028185654716894751240102991076751628613225034460815461177816572394855879662882740551715190	COCERCO	B	13
27472960344276922693731794339728667139618523353240672217797494125426692394633290527088156389670227	EATON ESIONADO N CCIDENTE E RANSPORT	V	0
5570752231455000704135544750200000767557488508041125015507651541465155451243707080550148914800735005		A	12

EDIT CHART

se pueden mostrar visualizaciones de los resultados de las estadísticas básicas reportadas. Empezaremos por mostrar un diagrama de cantidad de registros por letra inicial de CIE de los datos extraídos de los PDFs

```
In [106]: # https://plot.ly/python/pie-charts/
x = data['cie'].value_counts()

values = x.values

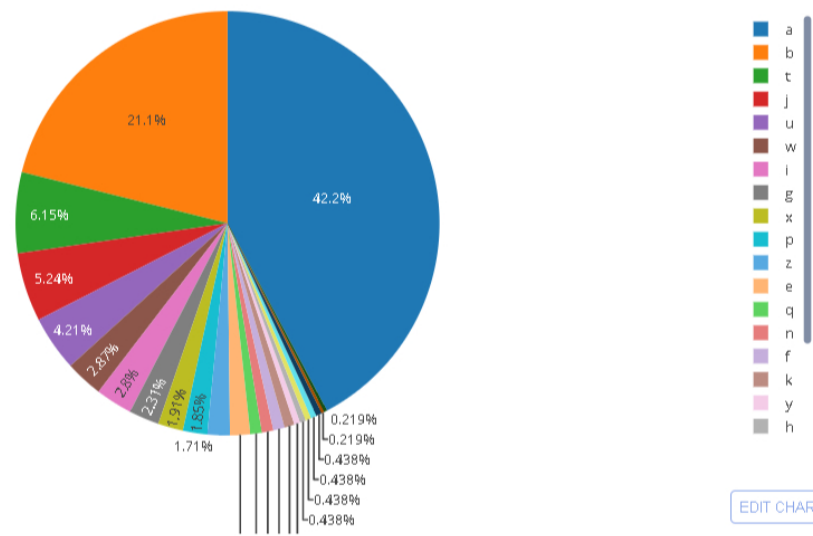
x = x.to_frame().T

labels = x.columns.values

trace = go.Pie(labels = labels, values = values)

py.iplot([trace], filename='conteoCIE')
```

Out[106]:



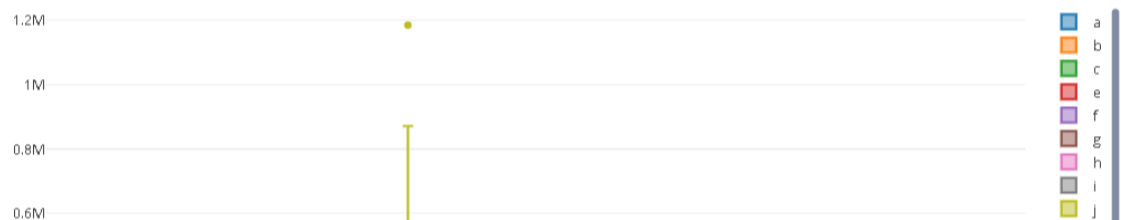
EDIT CHART

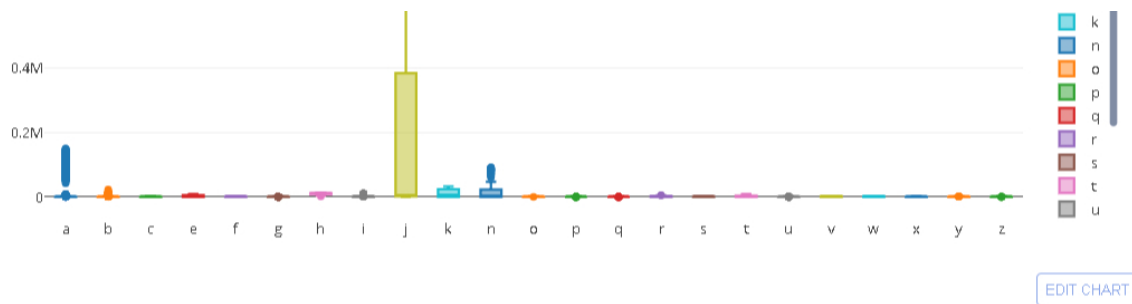
De modo que las enfermedades cuya CIE inicia con la letra A ocupan el 42.2% de los registros.

Ahora bien, se puede obtener una descripción de los datos agrupados por número de casos registrados y letra inicial de CIE

```
In [107]: # https://plot.ly/python/box-plots/
boxes = []
for cie in data.groupby(['cie']):
    trace = go.Box(y=cie[1]['casos'], name=str(cie[0]))
    boxes.append(trace)
py.iplot(boxes)
```

Out[107]:

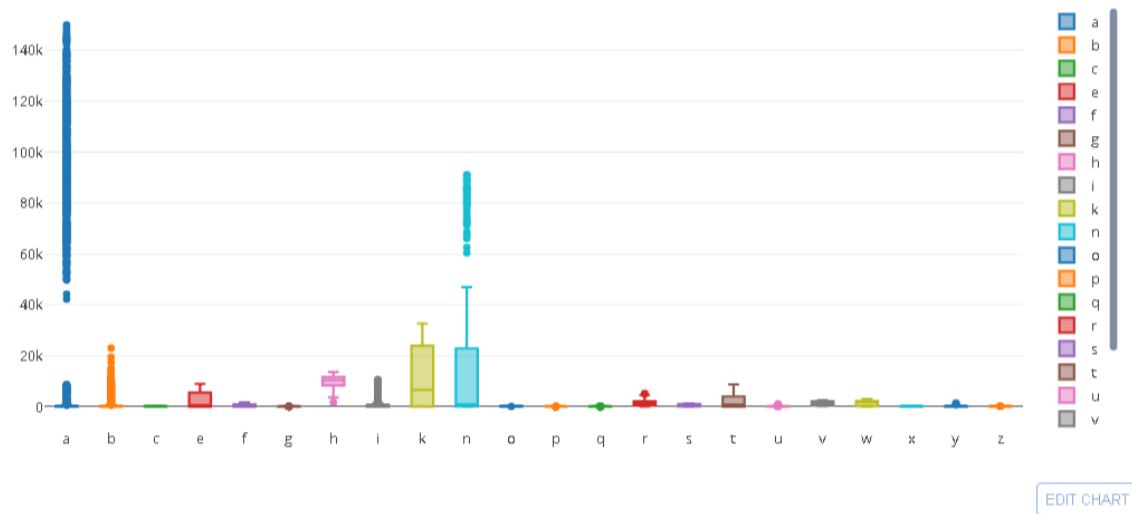




Quizás una mejor visualización del resto de enfermedades, se podría obtener al remover las que inician con j

```
In [108]: boxes = []
for cie in data.groupby(['cie']):
    if cie[0] != 'j':
        trace = go.Box(y=cie[1]['casos'], name=str(cie[0]))
        boxes.append(trace)
py.iplot(boxes)
```

Out[108]:



Al preprocesar los datos, se agrupan las enfermedades por primera letra de la CIE y se descubre que los grupos A y B contienen la mayoría de los registros, contando un 31.19% y un 19.57% respectivamente.

```
In [109]: # https://stackoverflow.com/a/51453257
pd.options.display.max_columns
data = data2

x = data['cie'].value_counts()

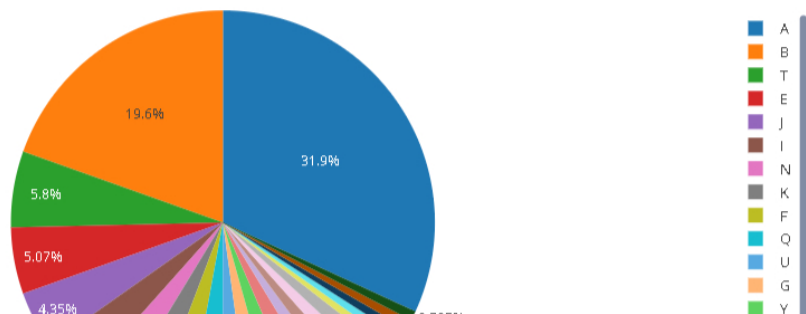
values = x.values

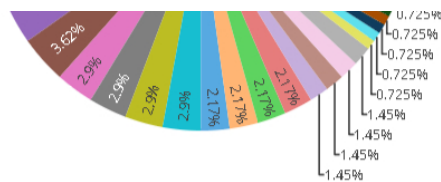
x = x.to_frame().T

labels = x.columns.values

trace = go.Pie(labels = labels, values = values)
py.iplot([trace], filename='conteoCIECluster')
```

Out[109]:





EDIT CHART

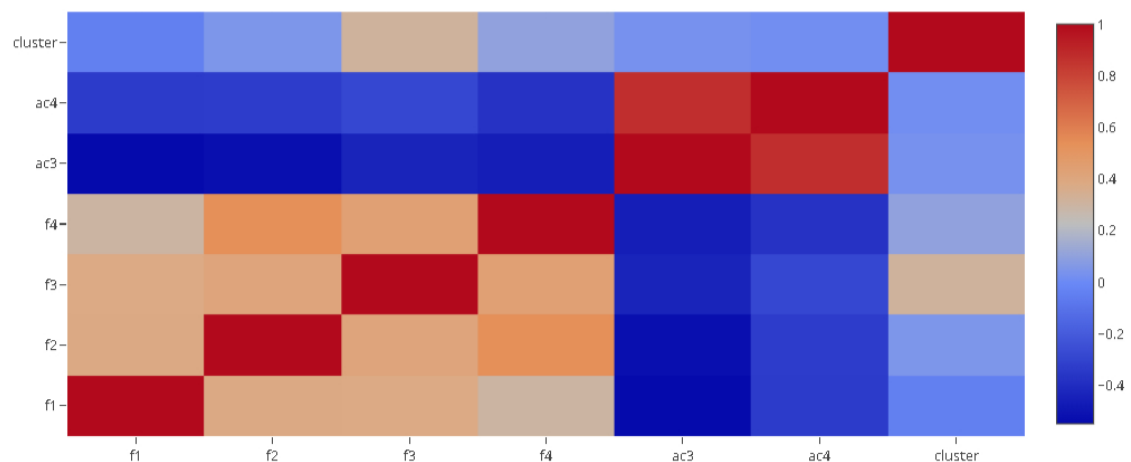
También se obtuvieron las correlaciones existentes entre las características de los datos preprocesados, siendo de interés aquellas que guardan correlación con el tipo de cluster asignado por k -medias

```
In [110]: # https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.factorize.html
data['cle'], uniques = pd.factorize(data['cle'])

#https://stackoverflow.com/a/19483025
#print(list(data.corr()))

trace = go.Heatmap(z=data.corr().values, x = list(data.corr()), y = list(data.corr()))
corr=[trace]
py.iplot(corr, filename='basic-heatmap')
```

Out[110]:



EDIT CHART

Finalmente, se seleccionan las características de los datos

```
In [111]: features = ["# m", "f1", "f2", "f3", "f4", "ac1", "ac2", "ac3", "ac4", "ac5", "ac6"]
x = data.loc[:, features].values
```

se normalizan

```
In [112]: # https://scikit-Learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html
from sklearn.preprocessing import MinMaxScaler
x = MinMaxScaler().fit(x).transform(x)
```

y con estas características normalizadas se puede hacer una selección a partir del umbral de varianza

```
In [113]: # https://stackoverflow.com/a/7670325
print("Columnas iniciales = {}".format(x.shape[1]))

# https://scikit-Learn.org/stable/modules/generated/sklearn.feature_selection.VarianceThreshold.html#sklearn.feature_selection.VarianceThreshold
from sklearn.feature_selection import VarianceThreshold
th = 0.05 # .8 * (1 - .8)
print("Umbral de varianza = {}".format(th))
sel = VarianceThreshold(threshold=th)
x = sel.fit_transform(x)
print("Columnas finales")
# https://stackoverflow.com/a/39812885
dataSelected = data[data.columns[sel.get_support(indices=True)]]
```

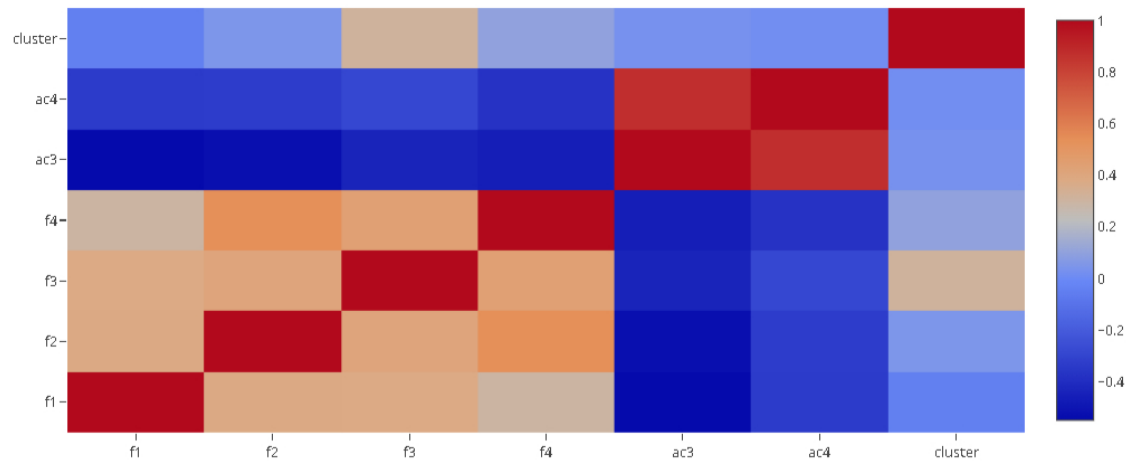
Columnas iniciales = 11
Umbral de varianza = 0.05
Columnas finales

y mostrar sus correlaciones

```
In [114]: # https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html
#print(data[['cle']])
dataSelected = dataSelected.assign(cluster=data[['cle']])
```

```
trace = go.Heatmap(z=dataSelected.corr().values, x = list(dataSelected.corr()), y= list(dataSelected.corr()))  
corr=[trace]  
py.iplot(corr, filename='basic-heatmap')
```

Out[114]:



[EDIT CHART](#)