

```
In [ ]: import pandas as pd
import numpy as np
```

```
In [ ]: # Leer el archivo nombres.csv y convertirlo en un df
df_base = pd.read_csv('nombres.csv')
# Ver los primeros registros del df
df_base.head()
```

```
Out[ ]:
```

	state	sex	year	name	quantity
0	MO	F	1910	Mary	611
1	MO	F	1910	Helen	313
2	MO	F	1910	Dorothy	270
3	MO	F	1910	Mildred	267
4	MO	F	1910	Ruth	237

```
In [ ]: # Mostrar información sobre el df creado
df_base.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5933561 entries, 0 to 5933560
Data columns (total 5 columns):
#   Column      Dtype
---  -
0   state       object
1   sex         object
2   year        int64
3   name        object
4   quantity    int64
dtypes: int64(2), object(3)
memory usage: 226.3+ MB
```

```
In [ ]: #Filtrar todas los hombres con nombre Mary
#Crear un nuevo dataframe
df_mary = df_base[(df_base["sex"] == "M") & (df_base["name"] == "Mary")]
df_mary
```

Out[ ]:

	state	sex	year	name	quantity
<b>76848</b>	MO	M	1913	Mary	6
<b>78179</b>	MO	M	1916	Mary	5
<b>79914</b>	MO	M	1920	Mary	9
<b>80430</b>	MO	M	1921	Mary	7
<b>81976</b>	MO	M	1924	Mary	5
...	...	...	...	...	...
<b>5815488</b>	CT	M	1955	Mary	5
<b>5816465</b>	CT	M	1958	Mary	5
<b>5816806</b>	CT	M	1959	Mary	5
<b>5817816</b>	CT	M	1962	Mary	6
<b>5899574</b>	KS	M	1936	Mary	5

1054 rows × 5 columns

```
In [ ]: #Ya no necesitamos el sexo y el nombre, hacemos una copia
nuevo_df = df_mary[["state", "quantity"]]
nuevo_df.head()
```

Out[ ]:

	state	quantity
<b>76848</b>	MO	6
<b>78179</b>	MO	5
<b>79914</b>	MO	9
<b>80430</b>	MO	7
<b>81976</b>	MO	5

```
In [ ]: #Obtener la suma de la columna quantity
nuevo_df[["quantity"]].sum()
```

quantity 9638  
dtype: int64

```
In [ ]: estado_df = nuevo_df.groupby("state").agg("sum")
estado_df
```

Out[ ]:           **quantity**

<b>state</b>	
<b>AL</b>	519
<b>AR</b>	159
<b>AZ</b>	6
<b>CA</b>	540
<b>CO</b>	18
<b>CT</b>	26
<b>DC</b>	51
<b>FL</b>	127
<b>GA</b>	482
<b>IA</b>	35
<b>IL</b>	570
<b>IN</b>	213
<b>KS</b>	5
<b>KY</b>	343
<b>LA</b>	237
<b>MA</b>	220
<b>MD</b>	92
<b>MI</b>	293
<b>MN</b>	59
<b>MO</b>	236
<b>MS</b>	392
<b>NC</b>	576
<b>NE</b>	10
<b>NJ</b>	137
<b>NM</b>	69
<b>NY</b>	785
<b>OH</b>	543
<b>OK</b>	92
<b>PA</b>	668
<b>SC</b>	361
<b>TN</b>	316
<b>TX</b>	969
<b>UT</b>	5
<b>VA</b>	263
<b>WA</b>	10

quantity	
state	
WI	92
WV	119

```
In [ ]: # Reto:  
# Obtener la cantidad promedio de Mary por estado usando estados_df  
promedio_df = nuevo_df.groupby("state").agg("mean")  
promedio_df
```

Out[ ]:

	quantity
state	
AL	10.380000
AR	6.625000
AZ	6.000000
CA	8.852459
CO	6.000000
CT	5.200000
DC	12.750000
FL	6.350000
GA	10.041667
IA	5.833333
IL	10.363636
IN	6.454545
KS	5.000000
KY	9.026316
LA	6.771429
MA	6.875000
MD	6.133333
MI	7.918919
MN	5.900000
MO	7.375000
MS	9.333333
NC	12.521739
NE	5.000000
NJ	6.227273
NM	6.272727
NY	11.544118
OH	8.901639
OK	6.133333
PA	11.133333
SC	9.025000
TN	9.028571
TX	13.842857
UT	5.000000
VA	7.514286
WA	5.000000

quantity	
state	
WI	6.133333
WV	6.263158

```
In [ ]: # Ordenar el set de datos por la mayor cantidad
estado_df.sort_values('quantity', ascending=False)
```

Out[ ]:           **quantity**

<b>state</b>	
<b>TX</b>	969
<b>NY</b>	785
<b>PA</b>	668
<b>NC</b>	576
<b>IL</b>	570
<b>OH</b>	543
<b>CA</b>	540
<b>AL</b>	519
<b>GA</b>	482
<b>MS</b>	392
<b>SC</b>	361
<b>KY</b>	343
<b>TN</b>	316
<b>MI</b>	293
<b>VA</b>	263
<b>LA</b>	237
<b>MO</b>	236
<b>MA</b>	220
<b>IN</b>	213
<b>AR</b>	159
<b>NJ</b>	137
<b>FL</b>	127
<b>WV</b>	119
<b>MD</b>	92
<b>OK</b>	92
<b>WI</b>	92
<b>NM</b>	69
<b>MN</b>	59
<b>DC</b>	51
<b>IA</b>	35
<b>CT</b>	26
<b>CO</b>	18
<b>NE</b>	10
<b>WA</b>	10
<b>AZ</b>	6

quantity	
state	
KS	5
UT	5

```
In [ ]: # Reto
# Obtener los 10 registros con la mayor cantidad de Mary
estado_df.sort_values('quantity', ascending=False).head(10)
```

```
Out[ ]:      quantity
state
TX      969
NY      785
PA      668
NC      576
IL      570
OH      543
CA      540
AL      519
GA      482
MS      392
```

```
In [ ]: # Para convertir el índice en una columna del df, se puede generar uno nuevo aleato
# Nota al uso de inplace=True para modificar el df
estado_df.reset_index(inplace=True)
estado_df
```



Out[ ]:

	state	quantity
0	AL	519
1	AR	159
2	AZ	6
3	CA	540
4	CO	18
5	CT	26
6	DC	51
7	FL	127
8	GA	482
9	IA	35
10	IL	570
11	IN	213
12	KS	5
13	KY	343
14	LA	237
15	MA	220
16	MD	92
17	MI	293
18	MN	59
19	MO	236
20	MS	392
21	NC	576
22	NE	10
23	NJ	137
24	NM	69
25	NY	785
26	OH	543
27	OK	92
28	PA	668
29	SC	361
30	TN	316
31	TX	969
32	UT	5
33	VA	263
34	WA	10
35	WI	92

	state	quantity
36	WV	119

```
In [ ]: estado_df.to_excel('estados_mary.xlsx', index=None)
```