

Second Exercise: Wines

In this exercise, We work with data about wines and We want to get some information from that data.

Load both files, the .csv as pandas dataframe and the .json as json

The data is stored in two datasets contained in two files.
One is a json file and the other is a csv file.

As We want to get some information about this data, We proceed to create one dataframe for each datasets. To do it we load the json file and obtain a dataframe and We read the csv as dataframe

```
In [1]: #We import all the modules that We will need to solve the tasks
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import json
```

```
In [2]: # We open the JSON file using the open() function
with open('Files/winemag-data-130k-v2.json') as f:
    #We load the content of the file
    data = json.load(f)

# We convert it in a dataframe using pandas.DataFrame() function
df_1 = pd.DataFrame(data)
```

	points	title	description	taster_name	taster_twitter_handle	price	designation	variety	region_1	region_2
0	87	Nicosia 2013 Vulkà Bianco (Etna)	Aromas tropical fruit, broom, brinestone...	Kerin O'Keefe	@kerinokeefe	NaN	Vulkà Bianco	White Blend	Etna	

1	87	Quinta dos Avidagos 2011 Avidagos Red (Douro)	This is ripe and fruity, a wine that is smooth...	Roger Voss	@vossroger	15.0	Avidagos	Portuguese Red	None	
---	----	---	---	------------	------------	------	----------	----------------	------	--

2	87	Rainstorm 2013 Pinot Gris (Willamette Valley)	Tart and snappy, the flavors of lime flesh and...	Paul Gregutt	@paulgwine	14.0	None	Pinot Gris	Willamette Valley	Willamette Valley
---	----	---	---	--------------	------------	------	------	------------	-------------------	-------------------

3	87	St. Julian 2013 Reserve Late Harvest Riesling... (Alsace)	Pineapple,梨, lemon, path and orange blossom...	Alexander Petreire		None	Reserve Late Harvest	Riesling	Lake Michigan Shore	
---	----	---	--	--------------------	--	------	----------------------	----------	---------------------	--

4	87	Sweet Cheeks 2012 Vintner's Reserve Wild Child... (Ala.)	Much like the regular bottling from 2012, this...	Paul Gregutt	@paulgwine	65.0	Vintner's Reserve Wild Child Block	Pinot Noir	Willamette Valley	Willamette Valley
---	----	--	---	--------------	------------	------	------------------------------------	------------	-------------------	-------------------

...
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

129966	90	Dr. H. Thanisch (Erben Müller-Burggraf) 2013... (Mosel)	Notes of honeysuckle and cantaloupe sweeten this...	Anna Lee C. Iijima		None	Brauneberger Juffer-Sonnenuhr Spätlese	Riesling	None	
--------	----	---	---	--------------------	--	------	--	----------	------	--

129967	90	Citation 2004 Pinot Noir (Oregon)	Citation is given as much as a decade of bottli...	Paul Gregutt	@paulgwine	75.0	None	Pinot Noir	Oregon	Oregon
--------	----	-----------------------------------	--	--------------	------------	------	------	------------	--------	--------

129968	90	Gresser 2013 Kriitt Gewurztraminer (Alsace)	Well-drained gravel soil gives this wine its c...	Roger Voss	@vossroger	30.0	Kriitt	Gewurztraminer	Alsace	
--------	----	---	---	------------	------------	------	--------	----------------	--------	--

129969	90	Domaine Marcel Deiss 2012 Pinot Gris (Alsace)	A dry style of Pinot Gris, this is crisp with ...	Roger Voss	@vossroger	32.0	None	Pinot Gris	Alsace	
--------	----	---	---	------------	------------	------	------	------------	--------	--

129970	90	Schöffel 2013 Lieux-dit Harth Cuvée Caroline (Alsace)	Big, rich and off-dry, this is powered by inte...	Roger Voss	@vossroger	21.0	Lieux-dit Harth Cuvée Caroline	Gewurztraminer	Alsace	
--------	----	---	---	------------	------------	------	--------------------------------	----------------	--------	--

129971 rows x 13 columns

```
In [3]: #Don't forget to ask about what our client wants
df_1 = pd.read_json('Files2/winemag-data-130k-v2.json')
```

```
In [4]: df_2 = pd.read_csv('Files/winemag-data-130k-v1.csv')
```

	Unnamed: 0	country	description	designation	points	price	province	region_1	region_2	variety	winery
0	0	US	This tremendous 100% varietal wine hails from ...	Martha's Vineyard	96	235.0	California	Napa Valley	Napa	Cabernet Sauvignon	Heltz

1	1	Spain	Ripe aromas of fig, blackberry and cassis are ...	Caradour Selección Especial Reserva	96	110.0	Northern Spain	Toro	NaN	Tinta de Toro	Bodega Carmen Rodríguez
---	---	-------	---	-------------------------------------	----	-------	----------------	------	-----	---------------	-------------------------

2	2	US	Mac Watson honors the memory of a wine once ...	Selected Late Harvest	96	90.0	California	Knights Valley	Sonoma	Sauvignon Blanc	Macaulay
---	---	----	---	-----------------------	----	------	------------	----------------	--------	-----------------	----------

3	3	US	This spent 20 months in 50% new French oak, an...	Reserve	96	65.0	Oregon	Willamette Valley	Willamette Valley	Pinot Noir	Ponzi
---	---	----	---	---------	----	------	--------	-------------------	-------------------	------------	-------

4	4	France	This is the top wine from La Bégude, named after...	La Brûlée	95	66.0	Provence	Bandol	NaN	Provence red blend	Domaine de la Bégude
---	---	--------	---	-----------	----	------	----------	--------	-----	--------------------	----------------------

...
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

150925	150925	Italy	Many people represent southern Italy.	NaN	91	20.0	Southern Italy	Fiano di Avellino	NaN	White Blend	Feudi di San Gregorio
--------	--------	-------	---------------------------------------	-----	----	------	----------------	-------------------	-----	-------------	-----------------------

150926	150926	France	Offers an intriguing nose with ginger, lime an...	Cuvée Prestige	91	27.0	Champagne	Champagne	NaN	Champagne Blend	H.Germain
--------	--------	--------	---	----------------	----	------	-----------	-----------	-----	-----------------	-----------

150927	150927	Italy	This classic example comes from a cru vineyard...	Terre di Dora	91	20.0	Southern Italy	Fiano di Avellino	NaN	White Blend	Terredora
--------	--------	-------	---	---------------	----	------	----------------	-------------------	-----	-------------	-----------

150928	150928	France	A perfect shade, with scents of peaches...	Grand Brut Rosé	90	52.0	Champagne	Champagne	NaN	Champagne Blend	Gosset
--------	--------	--------	--	-----------------	----	------	-----------	-----------	-----	-----------------	--------

150929	150929	Italy	More Pinot Grigios should taste like this. A...	NaN	90	15.0	Northeastern Italy	Alto Adige	NaN	Pinot Grigio	Alois Lageder
--------	--------	-------	---	-----	----	------	--------------------	------------	-----	--------------	---------------

150930 rows x 11 columns

Once We have both dataframes, We explore the columns of each dataframe

```
In [5]: df.columns
Out[5]: ['points', 'title', 'description', 'taster_name', 'taster_twitter_handle', 'price', 'designation', 'variety', 'region_1', 'region_2', 'province', 'winery']
```

```
In [6]: df_2.columns
Out[6]: Index(['Unnamed: 0', 'country', 'description', 'designation', 'points', 'price', 'province', 'region_1', 'region_2', 'variety', 'winery'], dtype='object')
```

Get a single dataframe that is the union of both files that only has the following columns: country, designation, points, price, state, winery.

We get just the columns that the task asks for from the resulting dataframe from the json

```
In [7]: #List with all the columns that We want
chosen_columns = ['country', 'designation', 'points', 'price', 'province', 'winery']

#We select those wanted columns and make a copy() so We still have the original first dataframe untouched
df_3 = df.loc[:, chosen_columns].copy()
```

	country	designation	points	price	province	winery
0	Italy	Vulkà Bianco	87	NaN	Sicily & Sardinia	Nicosia

1	Portugal	Avidagos	87	15.0	Douro	Quinta dos Avidagos
---	----------	----------	----	------	-------	---------------------

2	US	None	87	14.0	Oregon	Rainstorm
---	----	------	----	------	--------	-----------

3	US	Reserve Late Harvest	87	13.0	Michigan	St. Julian
---	----	----------------------	----	------	----------	------------

4	US	Vintner's Reserve Wild Child Block	87	65.0	Oregon	Sweet Cheeks
---	----	------------------------------------	----	------	--------	--------------

...
-----	-----	-----	-----	-----	-----	-----

129967	Germany	Brauneberger Juffer-Sonnenuhr Spätlese	90	28.0	Mosel	Dr. H. Thanisch (Erben Müller-Burggraf)
--------	---------	--	----	------	-------	---

129967	US	None	90	75.0	Oregon	Citation
--------	----	------	----	------	--------	----------

129968	France	Kriitt	90	30.0	Alsace	Domaine Gresser
--------	--------	--------	----	------	--------	-----------------

129969	France	None	90	32.0	Alsace	Domaine Marcel Deiss
--------	--------	------	----	------	--------	----------------------

129970	France	Lieux-dit Harth Cuvée Caroline	90	21.0	Alsace	Domaine Schöffel
--------	--------	--------------------------------	----	------	--------	------------------

129971 rows x 6 columns

We do the same with the second dataframe

```
In [8]: #We select those wanted columns and make a copy() so We still have the original second dataframe untouched
df_4 = df_2.loc[:, chosen_columns].copy()
```

	country	designation	points	price	province	winery
0	US	Martha's Vineyard	96	235.0	California	Heltz

1	Spain	Caradour Selección Especial Reserva	96	110.0	Northern Spain	Bodega Carmen Rodríguez
---	-------	-------------------------------------	----	-------	----------------	-------------------------

2	US	Special Selected Late Harvest	96	90.0	California	Macaulay
---	----	-------------------------------	----	------	------------	----------

3	US	Reserve	96	65.0	Oregon	Ponzi
---	----	---------	----	------	--------	-------

4	France	La Brûlée	95	66.0	Provence	Domaine de la Bégude
---	--------	-----------	----	------	----------	----------------------

...
-----	-----	-----	-----	-----	-----	-----

150925	Italy	NaN	91	20.0	Southern Italy	Feudi di San Gregorio
--------	-------	-----	----	------	----------------	-----------------------

150926	France	Cuvée Prestige	91	27.0	Champagne	H.Germain
--------	--------	----------------	----	------	-----------	-----------

150927	Italy	Terre di Dora	91	20.0	Champagne	Terredora
--------	-------	---------------	----	------	-----------	-----------

150928	France	Grand Brut Rosé	90	52.0	Champagne	Gosset
--------	--------	-----------------	----	------	-----------	--------

150929	Italy	NaN	90	15.0	Northeastern Italy	Alois Lageder
--------	-------	-----	----	------	--------------------	---------------

150930 rows x 6 columns

We compare if the number of rows of the final dataframe is equal to the rows of both previous dataframes

```
In [10]: df_final.shape[0] == df_3.shape[0] + df_4.shape[0]
```

True

Are there duplicates in the dataframe? Delete all duplicate rows

We see if there is duplicate rows

```
In [11]: #This is the final version of this exercise
#This variable will be useful later to see if the elimination of duplicate rows is correct
n_rows_total = df_final.shape[0]

#Shows the duplicate rows of the file
duplicateRows = df_final[df_final.duplicated()]

#To see this more clearly We sort by the values in winery column
duplicateRows.sort_values(by='winery')
```

	country	designation	points	price	province	winery
121401	US	NaN	88	NaN	Washington	'37 Cellars

133867	Spain	Brut	86	13.0	Catalonia	1+1+3
--------	-------	------	----	------	-----------	-------

116517	Spain	Cabernet Sauvignon	82	18.0	Catalonia	1+1+3
--------	-------	--------------------	----	------	-----------	-------

141687	Spain	Brut	86	13.0	Catalonia	1+1+3
--------	-------	------	----	------	-----------	-------

93375	Spain	Brut	87	16.0	Catalonia	1+1+3
-------	-------	------	----	------	-----------	-------

...
-----	-----	-----	-----	-----	-----	-----

15282	US	NaN	89	28.0	Washington	lMaurice
-------	----	-----	----	------	------------	----------

120292	US	NaN	92	34.0	Washington	lMaurice
--------	----	-----	----	------	------------	----------

81557	US	NaN	89	34.0	Washington	lMaurice
-------	----	-----	----	------	------------	----------

21901	US	Fred Estate	89	45.0	Washington	lMaurice
-------	----	-------------	----	------	------------	----------

91274	Slovenia	Izbrani	88	20.0	Kras	Štoka
-------	----------	---------	----	------	------	-------

77782 rows x 6 columns

```
In [12]: #Deleting duplicated rows
df_final.drop_duplicates(inplace=True)
```

	country	designation	points	price	province	winery
0	Italy	Vulkà Bianco	87	NaN	Sicily & Sardinia	Nicosia

1	Portugal	Avidagos	87	15.0	Douro	Quinta dos Avidagos
---	----------	----------	----	------	-------	---------------------

2	US	None	87	14.0	Oregon	Rainstorm
---	----	------	----	------	--------	-----------

3	US	Reserve Late Harvest	87	13.0	Michigan	St. Julian
---	----	----------------------	----	------	----------	------------

4	US	Vintner's Reserve Wild Child Block	87	65.0	Oregon	Sweet Cheeks
---	----	------------------------------------	----	------	--------	--------------

...
-----	-----	-----	-----	-----	-----	-----

149633	US	NaN	84	40.0	New York	Wliffier
--------	----	-----	----	------	----------	----------

149634	France	NaN	84	15.0	Alsace	W. Gisselbrecht
--------	--------	-----	----	------	--------	-----------------

149635	US	Bungalow Red	84	15.0	California	Casa Barranca
--------	----	--------------	----	------	------------	---------------

149636	Portugal	30-year old tawny	84	0.0	Port	Casa Santa Eufemia
--------	----------	-------------------	----	-----	------	--------------------

149638	Argentina	NaN	84	9.0	Mendoza Province	Finca El Portillo
--------	-----------	-----	----	-----	------------------	-------------------

203119 rows x 6 columns

```
In [13]: #If the number of rows in the dataframe with just the duplicate rows plus the number of rows without duplicates is equal to the original dataframe number of rows?
n_rows_total == df_final.shape[0] + duplicateRows.shape[0]
```

True

Make a study of the null values of the dataset, assign a zero to all the rows that have a null value in the price, after that, delete the rest of the null values of the dataframe. What is the resulting number of rows?

```
In [14]: #We use a list because the criteria to order a list is Num values, alpha uppercase from A-Z,
#alpha lowercase from a-z, and numbers from the littlest one to the biggest one
#If there is an string or a None value, sort() will raise an error
unique_prices = list(df_final.price.unique())
unique_prices.sort()
```

nan,

4.0,

5.0,

6.0,

7.0,

8.0,

9.0,

10.0,

11.0,

12.0,

13.0,

14.0,

15.0,

16.0,

17.0,

18.0,

19.0,

20.0,

21.0,

22.0,

23.0,

24.0,

25.0,

26.0,

27.0,

28.0,

29.0,

30.0,

31.0,

32.0,

33.0,

34.0,

35.0,

36.0,

37.0,

38.0,

39.0,

40.0,

41.0,

42.0,

43.0,

44.0,

45.0,

46.0,

47.0,

48.0,

49.0,

50.0,

51.0,

52.0,

53.0,

54.0,

55.0,

56.0,

57.0,

58.0,

59.0,

60.0,

61.0,

62.0,

63.0,

64.0,

65.0,

66.0,

67.0,

68.0,

69.0,

70.0,

71.0,

```
In [20]: #Select the rows that has as country Spain and any designation
#With the word "reserva"
mask_reserva_spain = (df_final.country=="Spain") & (df_final.designation.str.contains(pat = "reserva", case = False))
df_reserva = df_final[mask_reserva_spain].copy()
df_reserva
```

	country	designation	points	price	province	winery
490	Spain	Reserva Brut	87	19.0	Catalonia	V&N; Cellars
604	Spain	Reserva 1423	87	35.0	Northern Spain	Príncipe de Viana
836	Spain	Cinta Púrpura Brut Reserva	86	13.0	Catalonia	Juvé y Camps
837	Spain	Sweet Reserva	86	15.0	Catalonia	Juvé y Camps
1053	Spain	Reserva	85	30.0	Northern Spain	Finca La Emperatriz
...
145885	Spain	Reserva	88	11.0	Central Spain	Castillo de Almansa
147286	Spain	Ceremonia Reserva De Autor	85	15.0	Levante	Vicente Gandia
147460	Spain	Hoya De Cadenas Reserva	85	9.0	Levante	Gandia
148030	Spain	Reserva	81	18.0	Northern Spain	Señorio de Sarria
148036	Spain	Reserva	81	20.0	Northern Spain	Vilafiesta

1485 rows x 6 columns

Of the previous sub-set with the Spanish reserve wines, do they have a higher average score than all the wines or a lower one?

```
In [21]: df_final.dtypes
Out[21]:
country      object
designation   object
points        object
price        float64
province      object
winery        object
dtype: object

To be able to calculate the mean of the points column We have to change the data type to a numeric one
```

```
In [22]: df_final.points = pd.to_numeric(df_final.points)
df_reserva.points = pd.to_numeric(df_reserva.points)

We get the mean points from all the world
```

```
In [23]: #Gets the mean of the points with just 2 decimal places
mean_total=round(df_final.points.mean(),2)
mean_total
Out[23]:
88.55

We get the mean points from the Spanish wines with a "reserva" designation
```

```
In [24]: mean_reserva_spain = round(df_reserva.points.mean(),2)
mean_reserva_spain
Out[24]:
88.22

In [25]: print( "Are the Spanish Reserva wines better than the average of the world?: ",mean_total<mean_reserva_spain)
Are the Spanish Reserva wines better than the average of the world?: False

The average from "reserva" Spanish wines is lower than the world average
```

Take again the sub-set of Spanish reserve wines, in general, in Spain the wines from Riojaja (a great province in the north of Spain, North Spain will appear in the dataset) have grave fame, therefore , it is to be expected that this province or region is the one with the best wines in this sub-set, check whether this statement is true graphically.

We obtain the mean for each province in Spain and represent the result so We can more easily view the results.

```
In [26]: mean_points_by_province =df_final[df_final.country=="Spain"].groupby('province').points.mean()
mean_points_by_province.sorted = mean_points_by_province.sort_values()
mean_points_by_province.sorted
Out[26]:
province
Central Spain    85.145242
Levante          86.338104
Spain Other      86.366667
Catalonia        87.255867
Northern Spain   87.360326
Galicia          87.434203
Spanish Islands  88.400000
Andalucia        89.455479
Name: points, dtype: float64

In [27]: #We obtain both arrays to be used as variables in the barplot
names_spain_province = np.array(mean_points_by_province.sorted.keys())
points_spain_province = np.array(mean_points_by_province.sorted)

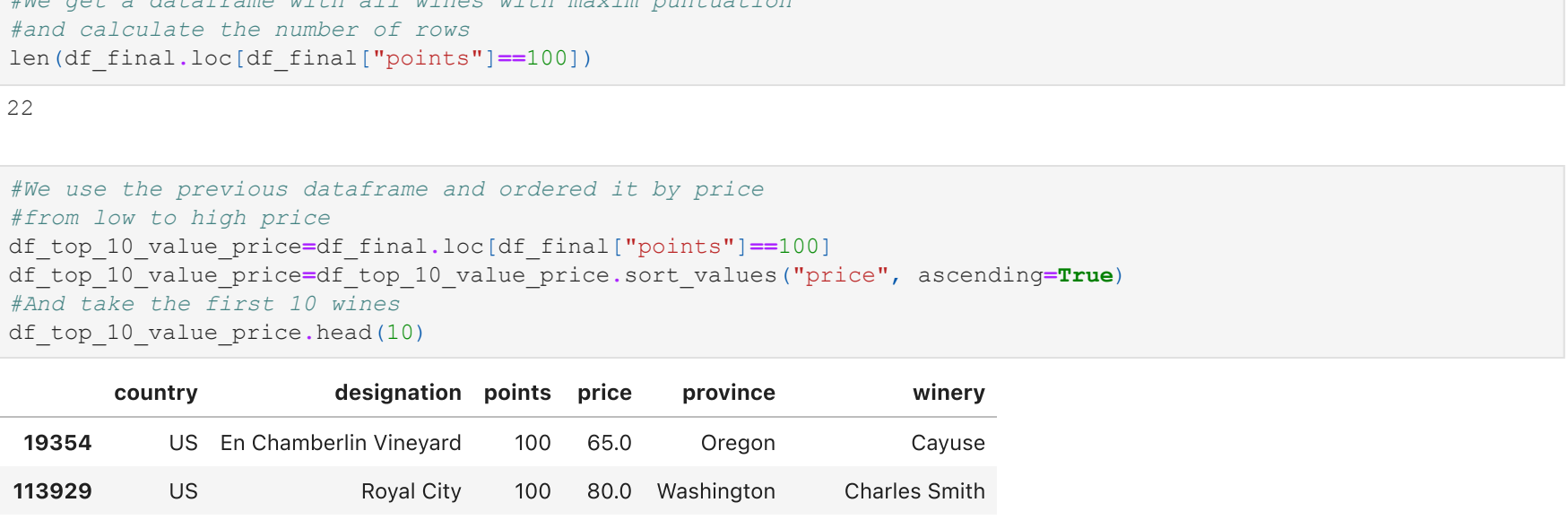
#The values are rounded.
points_spain_province = np.around(points_spain_province , decimals=2)

#Determining the size of the figure
fig, ax = plt.subplots(1, 1, figsize=(16,4))

#A barplot is created
sns.barplot(x=names_spain_province, y=points_spain_province)

#The values of each bar are set just above them
for i in range(1,len(points_spain_province)):
    ax.text(i, points_spain_province[i]*1, points_spain_province[i], horizontalalignment='center', size=20)

#The limits of the axes are set for a better view
ax.set(xlim=(-1,len(points_spain_province)), ylim=(0,120))
plt.show()
```



We observe that Rioja wines can be more popular but they are not the best valued in Spain.

Show the top 10 of the best valued wines.

We get the 10 most valued wines:

```
In [28]: #The top 10 of the best valued wines
df_top_10_wines = df_final.sort_values("points", ascending=False).head(10)
df_top_10_wines
Out[28]:
country      designation  points  price  province  winery
111754  Italy      Cerreto    100    270.0  Tuscany   Casanova di Neri
45781   Italy      Riserva    100    550.0  Tuscany   Biondi Santi
42197   Portugal   Barca-Velha  100    450.0  Douro     Casa Ferreira
114972  Portugal   Nacional Vintage  100    650.0  Port      Quinta do Noval
98647   US         Litton Estate Vineyard  100    100.0  California  Williams Selyem
89729   France  Le Meunil Blanc de Blancs Brut  100    617.0  Champagne  Salon
89728   France  Cristal Vintage Brut    100    250.0  Champagne  Louis Roederer
24151   Italy      Masseto    100    460.0  Tuscany   Tenuta dell'Ormeaia
113929  US         Royal City    100    80.0  Washington  Charles Smith
39286   Italy      Masseto    100    460.0  Tuscany   Tenuta dell'Ormeaia

As We noticed there are more than 10 wines but the maxim calification (100 points). We decided to choose the ones with lower price, as an added good feature.
```

```
In [29]: #We get a dataframe with all wines with maxim puntuation
#and calculate the number of rows
len(df_final.loc[df_final["points"]==100])
Out[29]:
22

In [30]: #We use the previous dataframe and ordered it by price
#from low to high price
df_top_10_value_price=df_final.loc[df_final["points"]==100]
df_top_10_value_price=df_top_10_value_price.sort_values("price", ascending=True)
#And take the first 10 wines
df_top_10_value_price.head(10)
```

	country	designation	points	price	province	winery
19354	US	En Chamberlin Vineyard	100	85.0	Oregon	Cayuse
113929	US	Royal City	100	80.0	Washington	Charles Smith
123545	US	Bionic Frog	100	80.0	Washington	Cayuse
98647	US	Litton Estate Vineyard	100	100.0	California	Williams Selyem
28954	Italy	Guado de' Gemoli	100	195.0	Tuscany	Giovanni Chiappini
7335	Italy	Occhio di Pernice	100	210.0	Tuscany	Avignonesi
111087	Italy	Occhio di Pernice	100	210.0	Tuscany	Avignonesi
92916	US	Hillside Select	100	215.0	California	Shaffer
114272	US	Red Wine	100	245.0	California	Sloan
89728	France	Cristal Vintage Brut	100	250.0	Champagne	Louis Roederer

Using this criteria it is observed that Portugal doesn't make it to the top 10.

```
In [31]: #The designation and winery are shown alone
df_top_10_value_price.loc[:,["winery", "designation"]].head(10)
Out[31]:
winery      designation
19354      Cayuse      En Chamberlin Vineyard
113929    Charles Smith  Royal City
123545      Cayuse      Bionic Frog
98647      Williams Selyem  Litton Estate Vineyard
28954      Giovanni Chappini  Guado de' Gemoli
7335       Avignonesi      Occhio di Pernice
111087     Avignonesi      Occhio di Pernice
12916      Shaffer       Hillside Select
114272     Sloan         Red Wine
89728     Louis Roederer  Cristal Vintage Brut
```

Show the average price of wines in each country. Which country has the highest average price? Does this country appear in the top 10 above?

```
In [32]: df_final.groupby("country").price.mean().head()
Out[32]:
country
Argentina    26.269494
Armenia      14.500000
Australia    38.229518
Austria      27.653320
Bosnia and Herzegovina  13.000000
Name: price, dtype: float64

We get the mean average price for each country. We order from higher to lower average price. And finally We get the first 5 country with the highest average price.
```

```
In [33]: #Get average price for each country
mean_price_by_country= df_final.groupby("country").price.mean()
#We sort that average price in descending order
mean_price_by_country.sorted = mean_price_by_country.sort_values(ascending=False)
#And with round the result with 2 decimal places
round(mean_price_by_country.sorted,2).head()
Out[33]:
country
Hungary      50.13
England      48.51
Germany      43.41
Canada       39.73
US            39.13
Name: price, dtype: float64

In [34]: print("Country with the highest average price:", mean_price_by_country.sorted.keys()[0], " Average price: ", round(mean_price_by_country.sorted.keys()[0] in df_top_10_wines.loc[:,["country"]])
Country with the highest average price: Hungary Average price: 50.13

To check if this country is in the top 10 list done above, a conditional is done. Is Hungary in the countries of the top 10 best wines?
```

```
In [35]: mean_price_by_country.sorted.keys()[0] in df_top_10_wines.loc[:,["country"]]
Out[35]:
False

This means even wines in Hungary are more expensive in average there is no a top 10 wine. Maybe the average price is higher for other reasons. Maybe is just more expensive to make wine in those countries.
```

Graphically show the price over the score based on the province of the subset of Spanish reserve wines. Which province has the most expensive wine? Is the north of Spain still the province with the best valued reserve wine?

The dataframe showing the reserve wines in Spain is used again

```
In [36]: df_reserva
Out[36]:
country      designation  points  price  province  winery
490  Spain      Reserva Brut    87    19.0  Catalonia  V&N; Cellars
604  Spain      Reserva 1423    87    35.0  Northern Spain  Príncipe de Viana
836  Spain  Cinta Púrpura Brut Reserva  86    13.0  Catalonia  Juvé y Camps
837  Spain      Sweet Reserva    86    15.0  Catalonia  Juvé y Camps
1053 Spain      Reserva        85    30.0  Northern Spain  Finca La Emperatriz
...  ...  ...  ...  ...  ...  ...
145885 Spain      Reserva        88    11.0  Central Spain  Castillo de Almansa
147286 Spain  Ceremonia Reserva De Autor  85    15.0  Levante    Vicente Gandia
147460 Spain  Hoya De Cadenas Reserva    85    9.0  Levante    Gandia
148030 Spain      Reserva        81    18.0  Northern Spain  Señorío de Sarria
148036 Spain      Reserva        81    20.0  Northern Spain  Vilafiesta

1485 rows x 6 columns

In [37]: #To show the mean of the points each province has
df_reserva.groupby("province").points.mean()
Out[37]:
province
Andalucia    91.800000
Catalonia    87.786407
Central Spain  85.000000
Levante      85.567568
Northern Spain  86.612225
Spain Other   89.000000
Name: points, dtype: float64

To create the array to get the column of points of the dataframe that will be used to plot:
```

```
In [38]: points = np.array(df_reserva.groupby("province").points.mean())
points
Out[38]:
array([91.8, 85.56756757, 86.61222541, 87.78640719, 85.56756757, 88.61222541, 89.00000000])

To create the array to get the column of price of the dataframe that will be used to plot:
```

```
In [39]: df_reserva.groupby("province").price.mean()
Out[39]:
province
Andalucia    115.800000
Catalonia    26.086826
Central Spain  24.883333
Levante      20.783784
Northern Spain  36.237822
Spain Other   24.000000
Name: price, dtype: float64

To create the array to get the column of provinces of the dataframe that will be used to plot.
```

```
In [40]: price=np.array(df_reserva.groupby("province").price.mean())
price = np.array(df_reserva.groupby("province").price.mean().keys())
province
Out[40]:
array(['Andalucia', 'Catalonia', 'Central Spain', 'Levante', 'Northern Spain', 'Spain Other'], dtype=object)

The dataframe to plot is created:
```

```
In [41]: df_reserva_price_value =pd.DataFrame({"province": province, "price": price , "points" : points})
df_reserva_price_value
Out[41]:
province  price  points
0  Andalucia  115.800000  91.800000
1  Catalonia  26.086826  87.796407
2  Central Spain  14.883333  85.000000
3  Levante    20.783784  85.567568
4  Northern Spain  36.237822  86.612225
5  Spain Other  24.000000  89.000000
```

Finally We make the plotting. We plot for each region from Spain the relationship between the mean price and the points given to his wines.

```
In [42]: fig, ax = plt.subplots(1, 1, figsize=(10,3))

# To create the scatter plot
ax.scatterplot(data = df_reserva_price_value , x= price , y= points , hue=province, marker="p" , s = 400, legend=True)

#To make the fontsize of the x-axis names in the graph bigger
ax.set_context("notebook", font_size=14, cml="font-size:15")

# To show the values of the x-axis near each point
for i, val in enumerate(price):
    plt.annotate(province[i], (price[i]*1.5, points[i]))

#Labels
ax.set_xlabel('PRICE', fontsize=20, fontweight='bold')
ax.set_ylabel('POINTS', fontsize=20, fontweight='bold')
#To set a grid
plt.grid(color='g', linestyle='--', linewidth=0.5)
#To set a color at the background
ax.set_facecolor('lightgrey')

plt.show()
```



The result tell us that "Spain Other" region has the best relationship, quality-price. Wines from Andalucia are slightly better but they are far more expensive

This results are quite accurate, because the prices without a value were set as zero at the beginning was just around 7%. So there is little bias.

```
In [43]: #The number of nulls in the original merged dataframe
nulos=df_3.price.isna().sum()
Out[43]:
0

In [44]: #percentage on NaNs at the price column in the original merged dataframe
nulos/(len(df_3)*100)
Out[44]:
6.921544036746659

Which is the most expensive wine?
```

```
In [45]: df_reserva.price.max()
Out[45]:
600.0

In [46]: max_row_index = df_reserva.price.idxmax()
df_reserva.loc[max_row_index]
```

```
Out[46]:
country      Spain
designation  Unico Reserva Especial
points        92
price        600.0
province     Northern Spain
winery       Vega Sicilia
Name: 97421, dtype: object
```

As we can see, even though this wine is the most expensive, it is not one of the most valued ones.

Questions?