

红酒评论数据集

1. 数据集加载

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 import pandas as pd
4 %matplotlib inline
5
6 path=['./dataset/wine_reviews/winemag-data-130k-v2.csv',
7       './dataset/wine_reviews/winemag-data_first150k.csv']
8 data = pd.concat([pd.read_csv(p) for p in path])
9 data.info()
```

```
1 <class 'pandas.core.frame.DataFrame'>
2 Int64Index: 280901 entries, 0 to 150929
3 Data columns (total 14 columns):
4  #   Column                Non-Null Count  Dtype
5  ---  ---
6  0   Unnamed: 0            280901 non-null  int64
7  1   country               280833 non-null  object
8  2   description           280901 non-null  object
9  3   designation           197701 non-null  object
10  4   points               280901 non-null  int64
11  5   price                258210 non-null  float64
12  6   province              280833 non-null  object
13  7   region_1             234594 non-null  object
14  8   region_2             111464 non-null  object
15  9   taster_name          103727 non-null  object
16  10  taster_twitter_handle  98758 non-null   object
17  11  title                 129971 non-null  object
18  12  variety               280900 non-null  object
19  13  winery                280901 non-null  object
20 dtypes: float64(1), int64(2), object(11)
21 memory usage: 32.1+ MB
```

2. 数据集可视化可摘要

2.1 数据摘要和可视化

- 共13个属性(ID除外),其中:
- 标称属性:
 1. country 国家
 2. designation 葡萄园
 3. province 省份
 4. region_1 葡萄酒生产地
 5. region_2 详细生产地
 6. variety 葡萄酒种类
 7. winery 酿酒厂

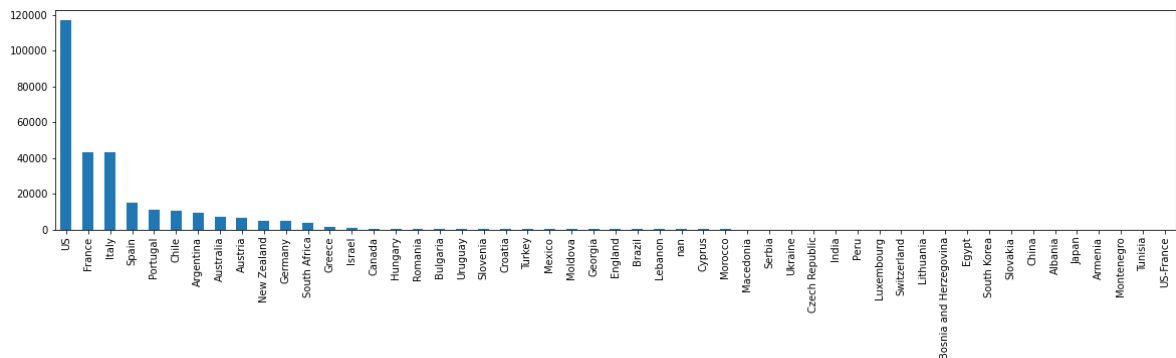
- 数值属性:
 8. points 评分
 9. price 价格
- 其他:
 10. title 评价标题
 11. description 评价内容
 12. taster_name 评价人
 13. taster_twitter_handle 评价人twitter

(1) country属性

```
1 print(data['country'].value_counts(dropna = False).head(10))
2 data['country'].value_counts(dropna = False).plot(kind="bar",figsize=(20,4))
```

```
1 US          116901
2 France      43191
3 Italy       43018
4 Spain       14913
5 Portugal    11013
6 Chile       10288
7 Argentina   9431
8 Australia   7286
9 Austria     6402
10 New Zealand 4739
11 Name: country, dtype: int64
```

```
1 <AxesSubplot:>
```



(2) designation 属性

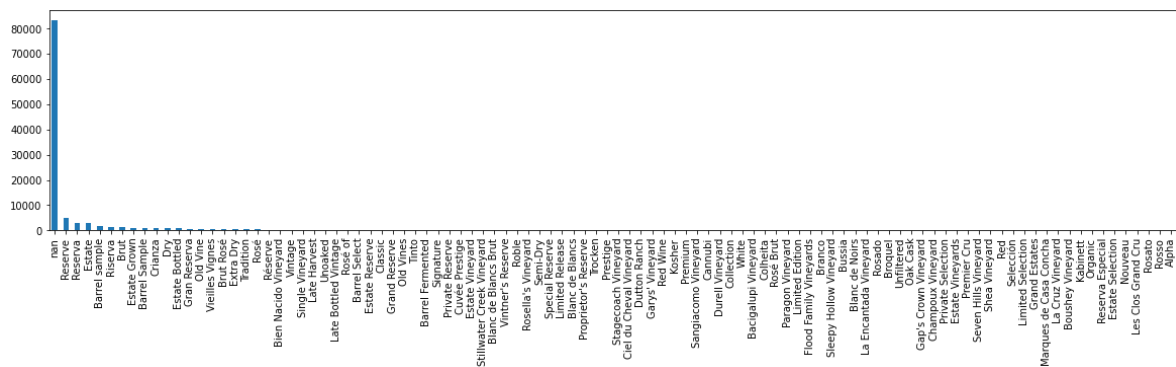
```
1 print(data['designation'].value_counts(dropna = False).head(10))
2 data['designation'].value_counts(dropna = False)
  [:100].plot(kind="bar",figsize=(20,4))
```

```

1  NaN                83200
2  Reserve            4761
3  Reserva            3069
4  Estate             2893
5  Barrel sample      1701
6  Riserva            1452
7  Brut               1137
8  Estate Grown       1070
9  Barrel Sample       891
10 Crianza             846
11 Name: designation, dtype: int64

```

```
1 | <AxesSubplot:>
```



(3) province 属性

```

1  print(data['province'].value_counts(dropna = False).head(10))
2  data['province'].value_counts(dropna = False)[:100].plot(kind="bar",figsize=
(20,4))

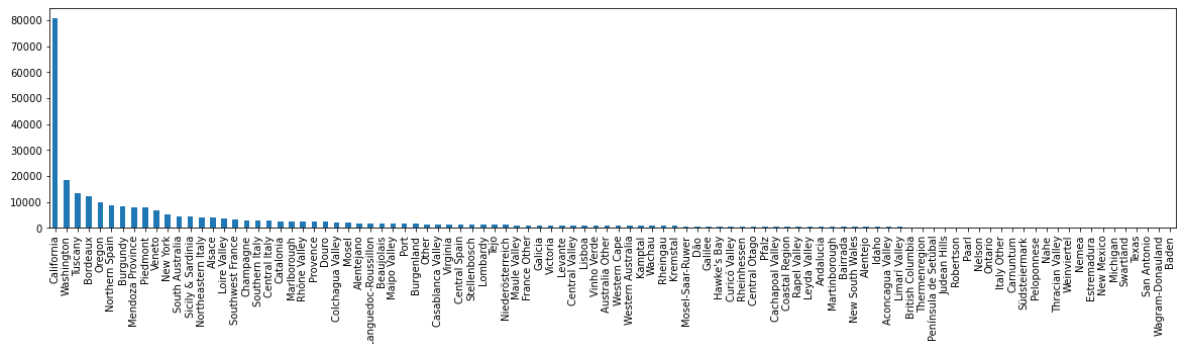
```

```

1  California          80755
2  washington          18389
3  Tuscany             13178
4  Bordeaux           12052
5  Oregon              9962
6  Northern Spain      8743
7  Burgundy            8288
8  Mendoza Province    8006
9  Piedmont            7822
10 Veneto              6678
11 Name: province, dtype: int64

```

```
1 | <AxesSubplot:>
```

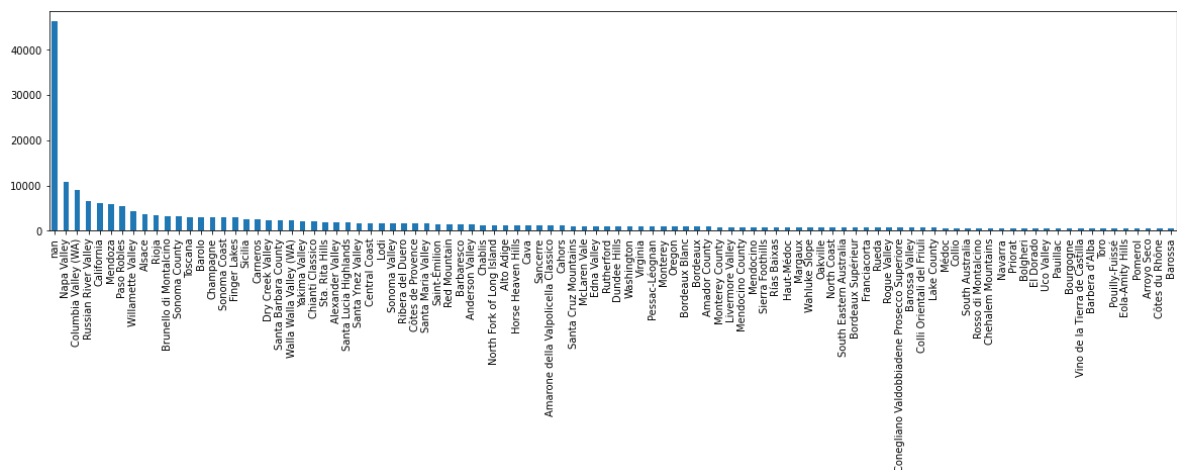


(4) region_1 属性

```
1 print(data['region_1'].value_counts(dropna = False).head(10))
2 data['region_1'].value_counts(dropna = False)[:100].plot(kind="bar",figsize=(20,4))
```

```
1 NaN 46307
2 Napa Valley 10689
3 Columbia Valley (WA) 9099
4 Russian River Valley 6662
5 California 6091
6 Mendoza 5887
7 Paso Robles 5403
8 willamette valley 4397
9 Alsace 3737
10 Rioja 3362
11 Name: region_1, dtype: int64
```

```
1 <AxesSubplot:>
```



(5) region_2 属性

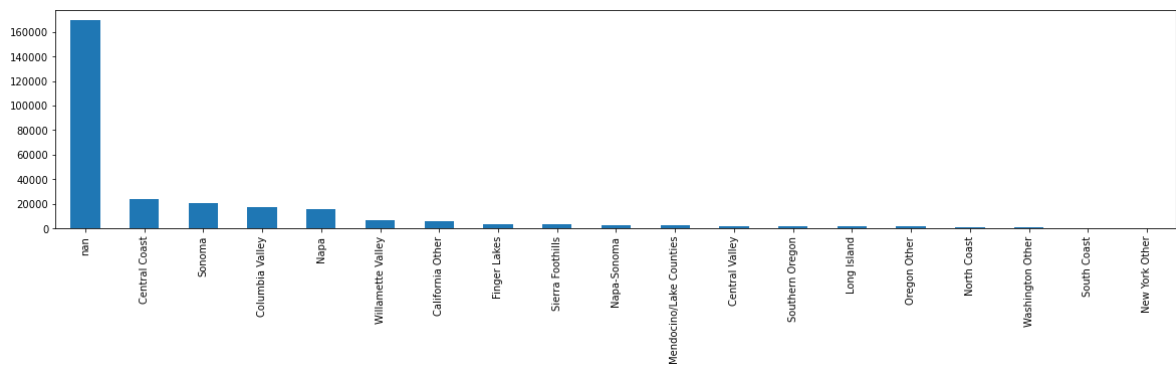
```
1 print(data['region_2'].value_counts(dropna = False).head(10))
2 data['region_2'].value_counts(dropna = False).plot(kind="bar",figsize=(20,4))
```

```

1 NaN 169437
2 Central Coast 24122
3 Sonoma 20286
4 Columbia Valley 17260
5 Napa 15615
6 Willamette Valley 6604
7 California Other 6179
8 Finger Lakes 3287
9 Sierra Foothills 3122
10 Napa-Sonoma 2814
11 Name: region_2, dtype: int64

```

```
1 <AxesSubplot:>
```



(6) variety 属性

```

1 print(data['variety'].value_counts(dropna = False).head(10))
2 data['variety'].value_counts(dropna = False)[:100].plot(kind="bar",figsize=
(20,4))

```

```

1 Pinot Noir 27563
2 Chardonnay 26235
3 Cabernet Sauvignon 22272
4 Red Blend 19008
5 Bordeaux-style Red Blend 14262
6 Sauvignon Blanc 11287
7 Riesling 10713
8 Syrah 9967
9 Merlot 8172
10 Zinfandel 6513
11 Name: variety, dtype: int64

```

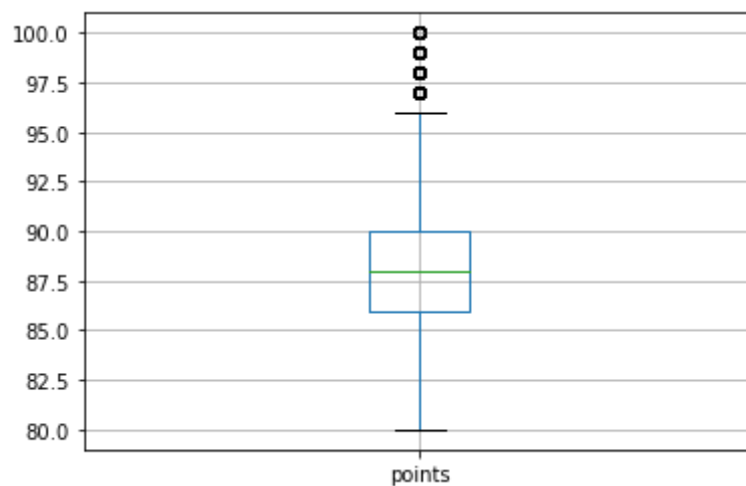
```
1 <AxesSubplot:>
```



```
1 最小值 80
2 第一四分位数 86.0
3 中位数 88.0
4 第三四分位数 90.0
5 最大值 100
6 缺失值个数 0
```

```
1 data.boxplot(['points'],return_type='dict')
```

```
1 {'whiskers': [<matplotlib.lines.Line2D at 0x2b839845fc8>,
2              <matplotlib.lines.Line2D at 0x2b8398614c8>],
3  'caps': [<matplotlib.lines.Line2D at 0x2b839868908>,
4           <matplotlib.lines.Line2D at 0x2b839868d48>],
5  'boxes': [<matplotlib.lines.Line2D at 0x2b839861408>],
6  'medians': [<matplotlib.lines.Line2D at 0x2b839868d88>],
7  'fliers': [<matplotlib.lines.Line2D at 0x2b839868e48>],
8  'means': []}
```



```
1 def detect_outliers(sr):
2     q1 = sr.quantile(0.25)
3     q3 = sr.quantile(0.75)
4     iqr = q3-q1 #Interquartile range
5     fence_low = q1-1.5*iqr
6     fence_high = q3+1.5*iqr
7     outliers = sr.loc[(sr < fence_low) | (sr > fence_high)]
8     return outliers
9 detect_outliers(data['points'])
10
```

```

1 | 345      100
2 | 346      98
3 | 347      97
4 | 348      97
5 | 349      97
6 | ...
7 | 143626   97
8 | 144494   97
9 | 144652   97
10 | 148047   97
11 | 149172   97
12 | Name: points, Length: 928, dtype: int64

```

其中points>97的均视为离群点

(9) price 属性

```

1 | print('最小值',data['price'].min())
2 | print('第一四分位数',data['price'].quantile(0.25))
3 | print('中位数',data['price'].median())
4 | print('第三四分位数',data['price'].quantile(0.75))
5 | print('最大值',data['price'].max())
6 | print('缺失值个数',data['price'].isnull().sum())

```

```

1 | 最小值 4.0
2 | 第一四分位数 16.0
3 | 中位数 25.0
4 | 第三四分位数 40.0
5 | 最大值 3300.0
6 | 缺失值个数 22691

```

```

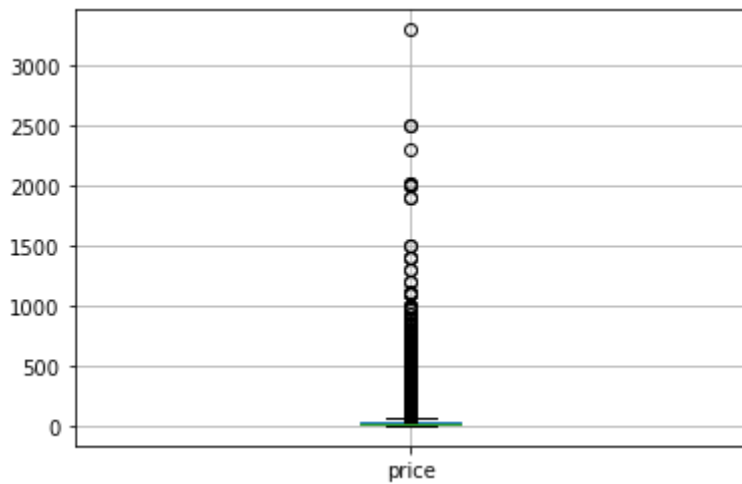
1 | data.boxplot(['price'],return_type='dict')

```

```

1 | {'whiskers': [<matplotlib.lines.Line2D at 0x1c667b47d08>,
2 | <matplotlib.lines.Line2D at 0x1c667b36ac8>],
3 | 'caps': [<matplotlib.lines.Line2D at 0x1c667b36a48>,
4 | <matplotlib.lines.Line2D at 0x1c667b341c8>],
5 | 'boxes': [<matplotlib.lines.Line2D at 0x1c667b44388>],
6 | 'medians': [<matplotlib.lines.Line2D at 0x1c667b34d08>],
7 | 'fliers': [<matplotlib.lines.Line2D at 0x1c667b32308>],
8 | 'means': []}

```

```

1 def detect_outliers(sr):
2     q1 = sr.quantile(0.25)
3     q3 = sr.quantile(0.75)
4     iqr = q3-q1 #Interquartile range
5     fence_low = q1-1.5*iqr
6     fence_high = q3+1.5*iqr
7     outliers = sr.loc[(sr < fence_low) | (sr > fence_high)]
8     return outliers
9 detect_outliers(data['price'])

```

```

1    60      100.0
2   111      85.0
3   118      80.0
4   119      80.0
5   134      78.0
6      ...
7  150570      90.0
8  150613      80.0
9  150727      83.0
10 150762     100.0
11 150765      87.0
12 Name: price, Length: 15128, dtype: float64

```

3. 数据缺失的处理

统计所有属性的数据缺失情况:

```

1 print(data.isnull().sum(axis=0))

```

```

1 Unnamed: 0          0
2 country            68
3 description         0
4 designation       83200
5 points             0

```

```
6 price 22691
7 province 68
8 region_1 46307
9 region_2 169437
10 taster_name 177174
11 taster_twitter_handle 182143
12 title 150930
13 variety 1
14 winery 0
15 dtype: int64
```

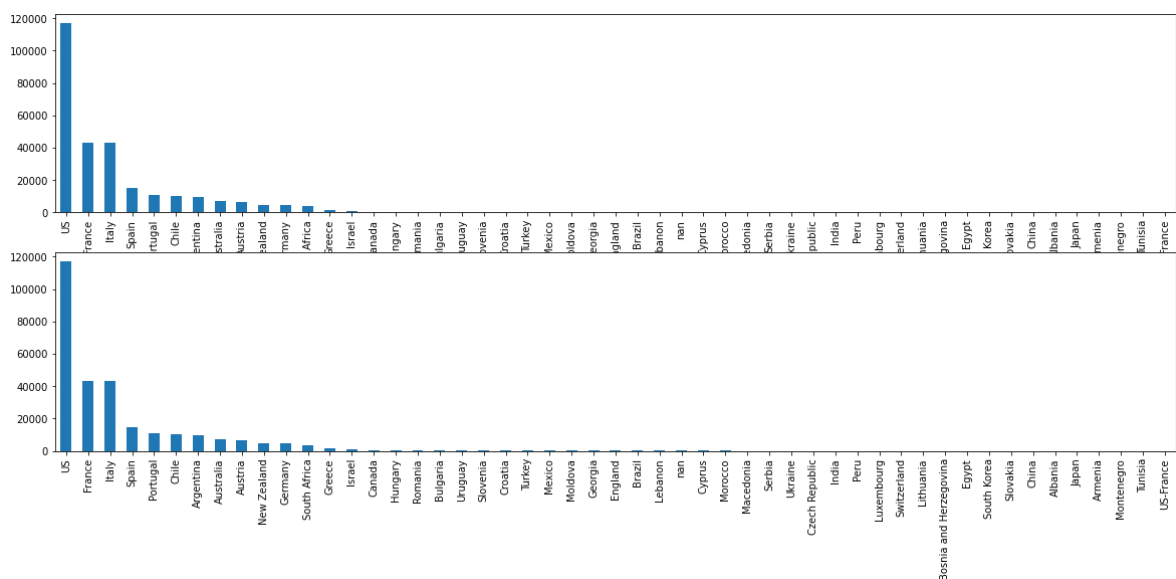
3.1 处理country属性缺失

缺失原因：统计失误，可以通过属性的相关关系来填补缺失值，查找country属性缺失的葡萄园所在国家并进行填充

```
1 import copy
2 designation2country = {
3     "Askitikos": "Greece",
4     "Shah": "US",
5     "Piedra Feliz": "Chile",
6 }
7 data_country = copy.deepcopy(data)
8 for index, row in data_country.iterrows():
9     if not row['country']:
10         row['country'] = designation2country[row['designation']]
```

```
1 import matplotlib.pyplot as plt
2 plt.subplot(2,1,1)
3 data["country"].value_counts(dropna = False).plot(kind='bar',figsize=(20,8))
4 plt.subplot(2,1,2)
5 data_country["country"].value_counts(dropna = False).plot(kind='bar',figsize=(20,8))
```

```
1 | <AxesSubplot:>
```

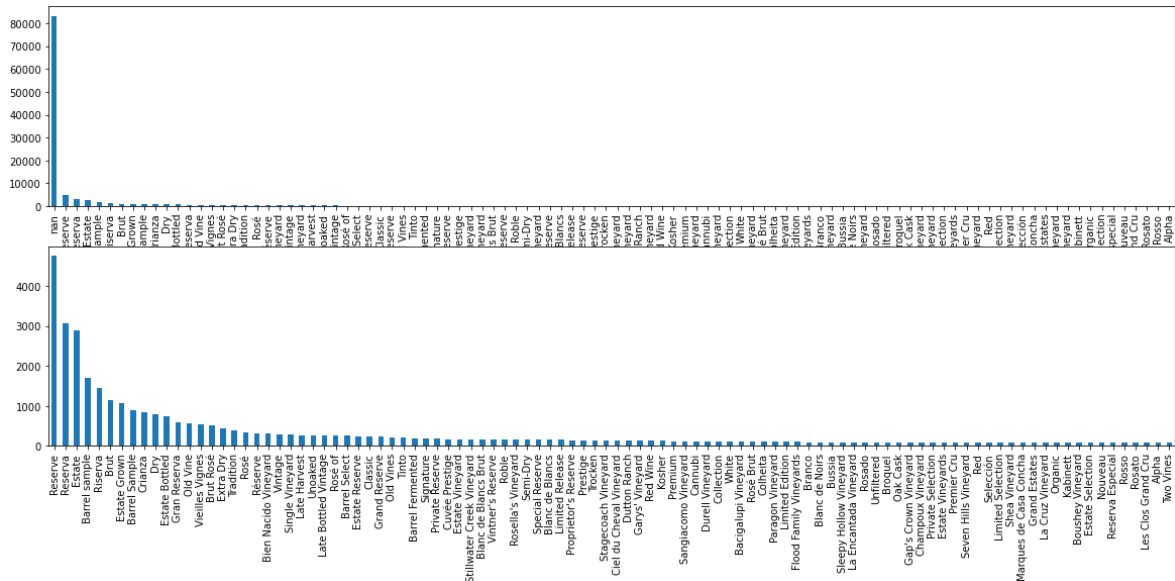


3.2 处理designation属性缺失

缺失原因：葡萄园可能较小或葡萄来源不明，将缺失部分剔除

```
1 data_designation = data.dropna(subset=['designation'])
2 plt.subplot(2,1,1)
3 data["designation"].value_counts(dropna = False)
4 [:100].plot(kind='bar',figsize=(20,8))
5 plt.subplot(2,1,2)
6 data_designation["designation"].value_counts(dropna = False)
7 [:100].plot(kind='bar',figsize=(20,8))
```

1 <AxesSubplot:>

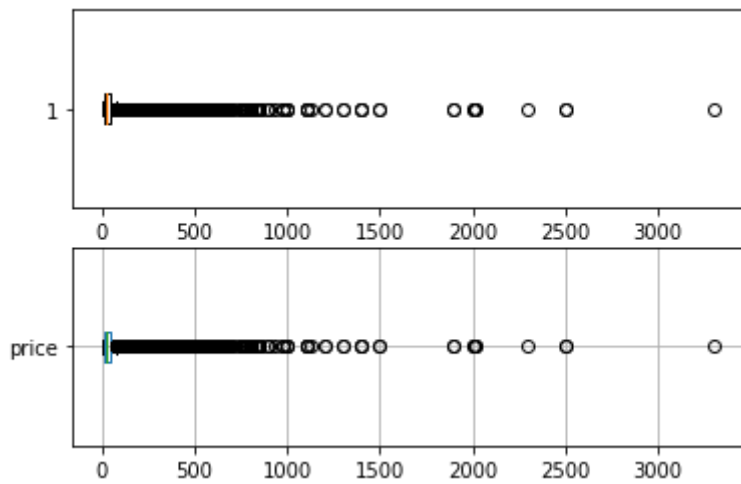


3.3 处理price属性缺失

缺失原因：通过数据平均值来填补缺失值

```
1 data_price= copy.deepcopy(data)
2 mean_price = data_price['price'].mean()
3 data_price = data_price['price'].fillna(mean_price)
4
5 plt.subplot(2,1,1)
6 plt.boxplot(data_price,vert=False)
7 plt.subplot(2,1,2)
8 data.boxplot(['price'],return_type='dict',vert=False)
```

```
1 {'whiskers': [<matplotlib.lines.Line2D at 0x7feddd199940>,
2              <matplotlib.lines.Line2D at 0x7feddd199c10>],
3  'caps': [<matplotlib.lines.Line2D at 0x7feddd199ee0>,
4           <matplotlib.lines.Line2D at 0x7feddd1371f0>],
5  'boxes': [<matplotlib.lines.Line2D at 0x7feddd199670>],
6  'medians': [<matplotlib.lines.Line2D at 0x7feddd1374c0>],
7  'fliers': [<matplotlib.lines.Line2D at 0x7feddd137790>],
8  'means': []}
```

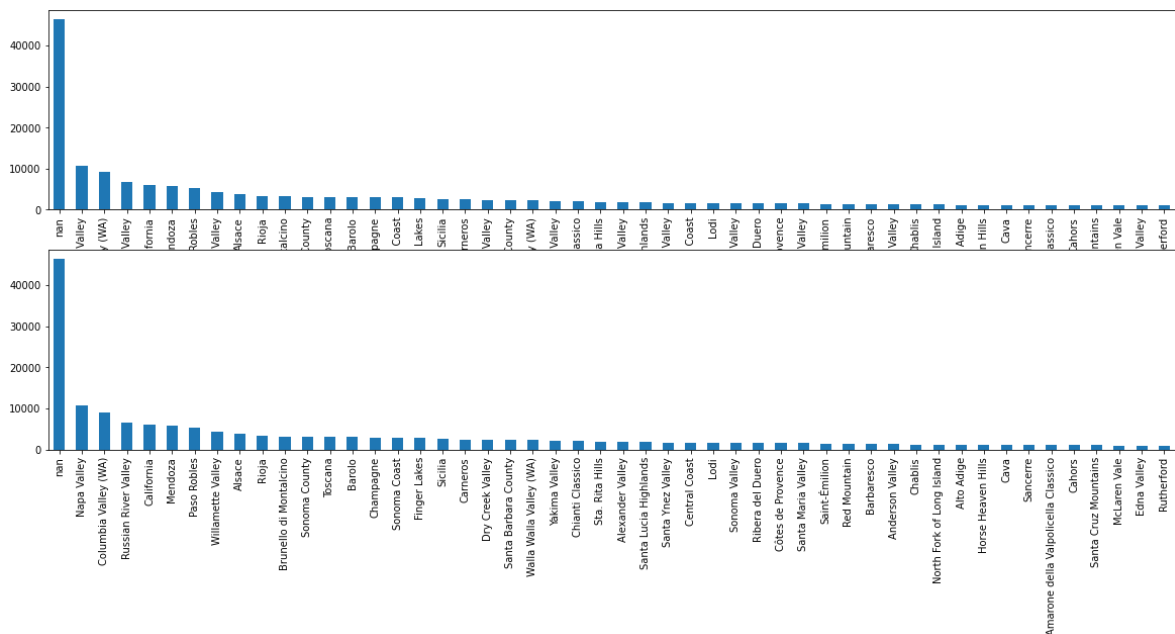


3.4 处理region_1属性缺失

缺失原因：通过相似对象的最高频率来填补缺失值

```
1
2 mode = data['region_1'].mode()
3 data_region1 = data['region_1'].fillna(mode)
4
5 plt.subplot(2,1,1)
6 data['region_1'].value_counts(dropna = False)[:50].plot(kind='bar',figsize=
7 (20,8))
8 plt.subplot(2,1,2)
9 data_region1.value_counts(dropna = False)[:50].plot(kind='bar',figsize=
10 (20,8))
```

1 <AxesSubplot:>



3.5 处理region_2属性缺失

缺失原因：将缺失部分剔除

```

1 data_region2 = data.dropna(subset=['region_2'])
2 plt.subplot(2,1,1)
3 data["region_2"].value_counts(dropna = False)[:50].plot(kind='bar',figsize=
  (20,8))
4 plt.subplot(2,1,2)
5 data_region2["region_2"].value_counts(dropna = False)
  [:50].plot(kind='bar',figsize=(20,8))

```

1 <AxesSubplot:>

