# Election project

## By: Hans Michael Célestin

```
In [47]:   import pandas as pd
           from pandas import Series, DataFrame
           import matplotlib.pyplot as plt
           import seaborn as sns
           sns.set_style('whitegrid')
           %matplotlib inline
           import scipy.stats as pearsonr
           from pandas_datareader import data, wb, DataReader
           import numpy as np
           from datetime import datetime
```

```
In [2]:    import requests # grab info from web
           from io import StringIO
```

```
In [3]:    url = "http://elections.huffingtonpost.com/pollster/api/charts/2012-general-election-romney-vs-obama.csv"

           source =requests.get(url).text
           poll_data =StringIO(source) # avoid string io error
```

```
In [4]:    poll_df = pd.read_csv(poll_data)
```

```
In [5]:    poll_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 428 entries, 0 to 427
Data columns (total 13 columns):
 #   Column                Non-Null Count  Dtype
---  ------                --------------  -----
 0   Obama                 428 non-null    float64
 1   Romney                428 non-null    float64
 2   Undecided             276 non-null    float64
 3   Other                 137 non-null    float64
 4   poll_id               428 non-null    int64
 5   pollster              428 non-null    object
 6   start_date            428 non-null    object
 7   end_date              428 non-null    object
 8   sample_subpopulation  428 non-null    object
 9   sample_size           414 non-null    float64
 10  mode                  428 non-null    object
 11  partisanship          428 non-null    object
 12  partisan_affiliation  428 non-null    object
dtypes: float64(5), int64(1), object(7)
memory usage: 43.6+ KB
```

```
In [6]:    poll_df.head()
```
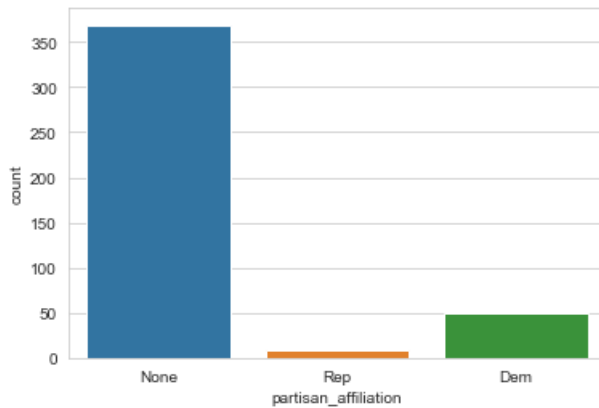
Out[6]:

| | Obama | Romney | Undecided | Other | poll_id | pollster | start_date | end_date | sample_subpopulation | sample_size | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 47.0 | 47.0 | 6.0 | NaN | 16674 | Politico/GWU/Battleground | 2012-11-04 | 2012-11-05 | Likely Voters | 1000.0 | L |
| 1 | 49.0 | 47.0 | 3.0 | NaN | 16733 | YouGov/Economist | 2012-11-03 | 2012-11-05 | Likely Voters | 740.0 | |
| 2 | 48.0 | 48.0 | 4.0 | NaN | 16681 | Gravis Marketing | 2012-11-03 | 2012-11-05 | Likely Voters | 872.0 | A |
| 3 | 50.0 | 49.0 | NaN | 1.0 | 16679 | IBD/TIPP | 2012-11-03 | 2012-11-05 | Likely Voters | 712.0 | L |
| 4 | 48.0 | 49.0 | NaN | NaN | 16677 | Rasmussen | 2012-11-03 | 2012-11-05 | Likely Voters | 1500.0 | A |

```
In [7]:    sns.countplot('partisan_affiliation',data=poll_df)
```

C:\Users\hansm\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
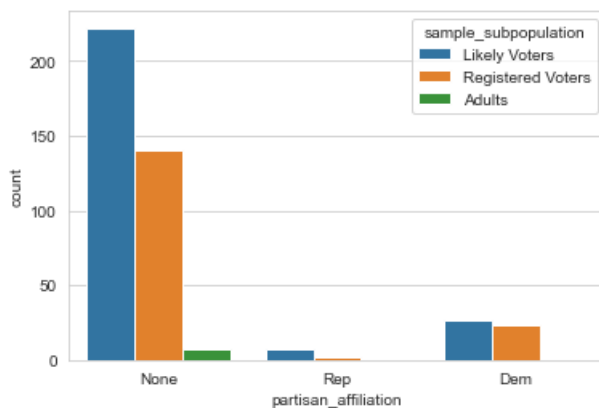  warnings.warn(

Out[7]: <AxesSubplot:xlabel='partisan_affiliation', ylabel='count'>



In [8]: ```python
sns.countplot('partisan_affiliation',data=poll_df,hue='sample_subpopulation')
```

C:\Users\hansm\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
  warnings.warn(

Out[8]: <AxesSubplot:xlabel='partisan_affiliation', ylabel='count'>



In [9]: ```python
poll_df.head()
```

Out[9]:

| | Obama | Romney | Undecided | Other | poll_id | pollster | start_date | end_date | sample_subpopulation | sample_size | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 47.0 | 47.0 | 6.0 | NaN | 16674 | Politico/GWU/Battleground | 2012-11-04 | 2012-11-05 | Likely Voters | 1000.0 | L |
| 1 | 49.0 | 47.0 | 3.0 | NaN | 16733 | YouGov/Economist | 2012-11-03 | 2012-11-05 | Likely Voters | 740.0 | |
| 2 | 48.0 | 48.0 | 4.0 | NaN | 16681 | Gravis Marketing | 2012-11-03 | 2012-11-05 | Likely Voters | 872.0 | A |
| 3 | 50.0 | 49.0 | NaN | 1.0 | 16679 | IBD/TIPP | 2012-11-03 | 2012-11-05 | Likely Voters | 712.0 | L |
| 4 | 48.0 | 49.0 | NaN | NaN | 16677 | Rasmussen | 2012-11-03 | 2012-11-05 | Likely Voters | 1500.0 | A |

In [10]: ```python
avg =pd.DataFrame(poll_df.mean())
avg.drop(['poll_id','sample_size'],axis =0 ,inplace=True)
```

In [11]: ```python
avg.head()
```

Out[11]:

| | 0 |
|---|---|

|  | 0 |
| --- | --- |
| **Obama** | 47.161215 |
| **Romney** | 45.228972 |
| **Undecided** | 5.615942 |
| **Other** | 2.686131 |

In [12]:
```python
std =pd.DataFrame(poll_df.std())
std.drop(['poll_id','sample_size'],axis =0 ,inplace=True)
```
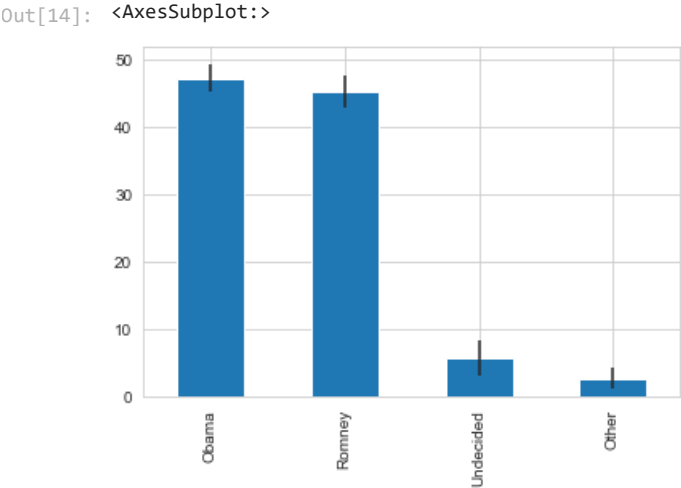
In [13]:
```python
std.head()
```

Out[13]:

|  | 0 |
| --- | --- |
| **Obama** | 2.100449 |
| **Romney** | 2.370565 |
| **Undecided** | 2.629407 |
| **Other** | 1.612232 |

In [14]:
```python
avg.plot(yerr=std,kind='bar',legend=False)
```

Out[14]: <AxesSubplot:>



In [15]:
```python
poll_avg = pd.concat([avg,std],axis=1)
poll_avg.columns= ['Avarage','STD']
poll_avg
```

Out[15]:

|  | Avarage | STD |
| --- | --- | --- |
| **Obama** | 47.161215 | 2.100449 |
| **Romney** | 45.228972 | 2.370565 |
| **Undecided** | 5.615942 | 2.629407 |
| **Other** | 2.686131 | 1.612232 |

In [17]:
```python
poll_df.head()
```

Out[17]:

|  | Obama | Romney | Undecided | Other | poll_id | pollster | start_date | end_date | sample_subpopulation | sample_size |  |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| **0** | 47.0 | 47.0 | 6.0 | NaN | 16674 | Politico/GWU/Battleground | 2012-11-04 | 2012-11-05 | Likely Voters | 1000.0 | L |
| **1** | 49.0 | 47.0 | 3.0 | NaN | 16733 | YouGov/Economist | 2012-11-03 | 2012-11-05 | Likely Voters | 740.0 |  |
| **2** | 48.0 | 48.0 | 4.0 | NaN | 16681 | Gravis Marketing | 2012-11-03 | 2012-11-05 | Likely Voters | 872.0 | A |

| | Obama | Romney | Undecided | Other | poll_id | pollster | start_date | end_date | sample_subpopulation | sample_size | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **3** | 50.0 | 49.0 | NaN | 1.0 | 16679 | IBD/TIPP | 2012-11-03 | 2012-11-05 | Likely Voters | 712.0 | L |
| **4** | 48.0 | 49.0 | NaN | NaN | 16677 | Rasmussen | 2012-11-03 | 2012-11-05 | Likely Voters | 1500.0 | A |

```
In [20]: poll_df.plot(x='end_date',y=['Obama', 'Romney','Undecided'],linestyle='',marker='o')
```

```
Out[20]: <AxesSubplot:xlabel='end_date'>
```



```
In [21]: poll_df['Difference'] = (poll_df.Obama-poll_df.Romney)/100
```

```
In [22]: poll_df
```

```
Out[22]:
```

| | Obama | Romney | Undecided | Other | poll_id | pollster | start_date | end_date | sample_subpopulation | sample_size |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 47.0 | 47.0 | 6.0 | NaN | 16674 | Politico/GWU/Battleground | 2012-11-04 | 2012-11-05 | Likely Voters | 1000.0 |
| **1** | 49.0 | 47.0 | 3.0 | NaN | 16733 | YouGov/Economist | 2012-11-03 | 2012-11-05 | Likely Voters | 740.0 |
| **2** | 48.0 | 48.0 | 4.0 | NaN | 16681 | Gravis Marketing | 2012-11-03 | 2012-11-05 | Likely Voters | 872.0 |
| **3** | 50.0 | 49.0 | NaN | 1.0 | 16679 | IBD/TIPP | 2012-11-03 | 2012-11-05 | Likely Voters | 712.0 |
| **4** | 48.0 | 49.0 | NaN | NaN | 16677 | Rasmussen | 2012-11-03 | 2012-11-05 | Likely Voters | 1500.0 |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| **423** | 44.0 | 44.0 | 5.0 | NaN | 12456 | YouGov/Economist | 2012-01-07 | 2012-01-10 | Registered Voters | 715.0 |
| **424** | 48.0 | 43.0 | 5.0 | 4.0 | 12444 | Ipsos/Reuters | 2012-01-05 | 2012-01-09 | Registered Voters | 900.0 |
| **425** | 45.0 | 47.0 | 8.0 | NaN | 12422 | CBS | 2012-01-04 | 2012-01-08 | Registered Voters | 1247.0 |
| **426** | 42.0 | 42.0 | 8.0 | 8.0 | 12401 | Rasmussen | 2012-01-03 | 2012-01-04 | Likely Voters | 1000.0 |
| **427** | 49.0 | 40.0 | 6.0 | NaN | 12402 | YouGov/Economist | 2011-12-31 | 2012-01-03 | Registered Voters | 715.0 |

428 rows × 14 columns

```
In [24]: poll_df = poll_df.groupby(['start_date'],as_index=False).mean()
         poll_df.head()
```
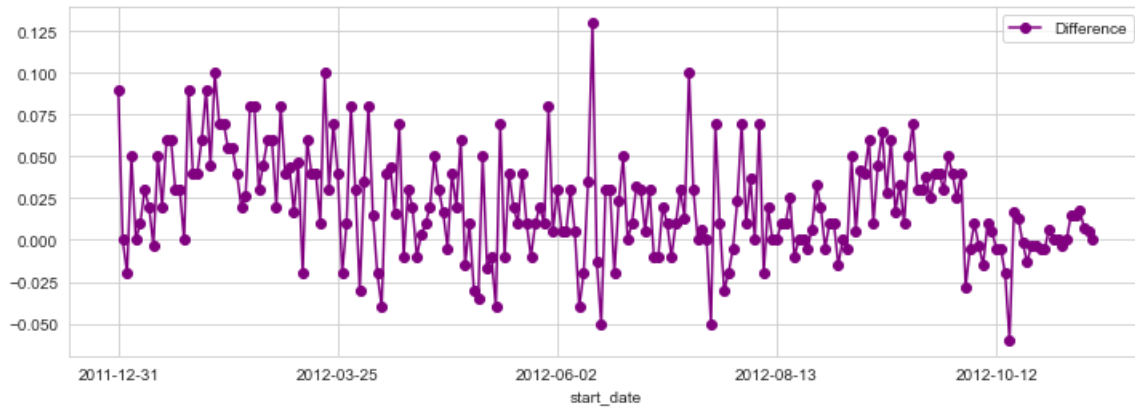
```
Out[24]:
```

| | start_date | Obama | Romney | Undecided | Other | poll_id | sample_size | Difference |
|---|---|---|---|---|---|---|---|---|

|   | start_date | Obama | Romney | Undecided | Other | poll_id | sample_size | Difference |
|---|-----------|-------|--------|-----------|-------|---------|-------------|------------|
| 0 | 2011-12-31 | 49.0 | 40.0 | 6.0 | NaN | 12402.0 | 715.0 | 0.09 |
| 1 | 2012-01-03 | 42.0 | 42.0 | 8.0 | 8.0 | 12401.0 | 1000.0 | 0.00 |
| 2 | 2012-01-04 | 45.0 | 47.0 | 8.0 | NaN | 12422.0 | 1247.0 | -0.02 |
| 3 | 2012-01-05 | 48.0 | 43.0 | 5.0 | 4.0 | 12444.0 | 900.0 | 0.05 |
| 4 | 2012-01-07 | 44.0 | 44.0 | 5.0 | NaN | 12456.0 | 715.0 | 0.00 |

```
In [27]: poll_df.plot('start_date','Difference',figsize=(12,4),marker='o',linestyle='-',color='purple')
```

```
Out[27]: <AxesSubplot:xlabel='start_date'>
```



```
In [29]: row_in =0
         xlimit= []
         for date in poll_df['start_date']:
             if date[0:7] == '2012-10' :
                 xlimit.append(row_in)
                 row_in +=1
             else:
                 row_in +=1
         print( min(xlimit))
         print( max(xlimit))
```
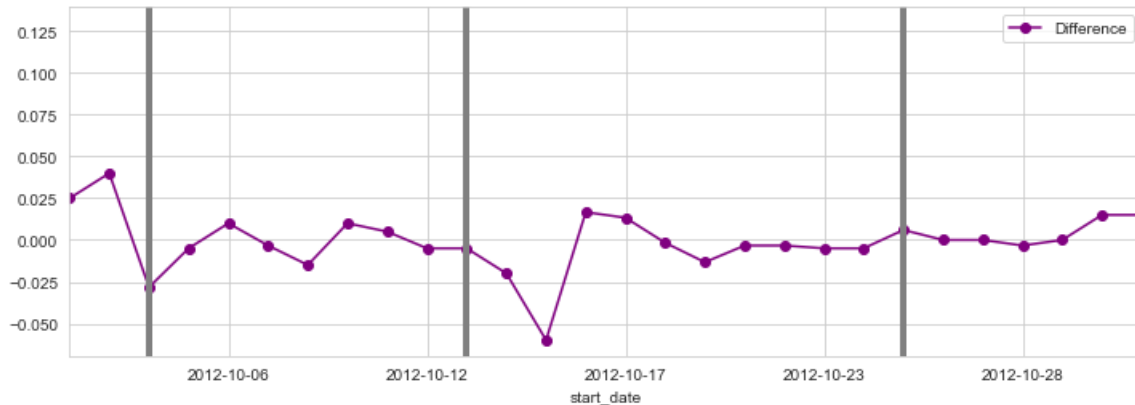
```
191
218
```

```
In [35]: poll_df.plot('start_date','Difference',figsize=(12,4),marker='o',linestyle='-',color='purple',xlim=(191,218))
         plt.axvline(x=191+2,linewidth =4, color='grey')
         plt.axvline(x=191+10,linewidth =4, color='grey')
         plt.axvline(x=191+21,linewidth =4, color='grey')
```

```
Out[35]: <matplotlib.lines.Line2D at 0x1ace00e02b0>
```



DONOR DATA

```
In [36]: donor_df =pd.read_csv('Election_Donor_Data.csv')
```

```
C:\Users\hansm\anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3146: DtypeWarning: Columns (6) have
```

```
mixed types.Specify dtype option on import or set low_memory=False.
  has_raised = await self.run_ast_nodes(code_ast.body, cell_name,
```

In [39]: `donor_df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1001731 entries, 0 to 1001730
Data columns (total 16 columns):
 #   Column            Non-Null Count    Dtype
---  ------            --------------    -----
 0   cmte_id           1001731 non-null  object
 1   cand_id           1001731 non-null  object
 2   cand_nm           1001731 non-null  object
 3   contbr_nm         1001731 non-null  object
 4   contbr_city       1001712 non-null  object
 5   contbr_st         1001727 non-null  object
 6   contbr_zip        1001620 non-null  object
 7   contbr_employer   988002 non-null   object
 8   contbr_occupation 993301 non-null   object
 9   contb_receipt_amt 1001731 non-null  float64
 10  contb_receipt_dt  1001731 non-null  object
 11  receipt_desc      14166 non-null    object
 12  memo_cd           92482 non-null    object
 13  memo_text         97770 non-null    object
 14  form_tp           1001731 non-null  object
 15  file_num          1001731 non-null  int64
dtypes: float64(1), int64(1), object(14)
memory usage: 122.3+ MB
```

In [40]: `donor_df.head()`

Out[40]:

| | cmte_id | cand_id | cand_nm | contbr_nm | contbr_city | contbr_st | contbr_zip | contbr_employer | contbr_occupation | contb_ |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | C00410118 | P20002978 | Bachmann, Michelle | HARVEY, WILLIAM | MOBILE | AL | 3.6601e+08 | RETIRED | RETIRED | |
| 1 | C00410118 | P20002978 | Bachmann, Michelle | HARVEY, WILLIAM | MOBILE | AL | 3.6601e+08 | RETIRED | RETIRED | |
| 2 | C00410118 | P20002978 | Bachmann, Michelle | SMITH, LANIER | LANETT | AL | 3.68633e+08 | INFORMATION REQUESTED | INFORMATION REQUESTED | |
| 3 | C00410118 | P20002978 | Bachmann, Michelle | BLEVINS, DARONDA | PIGGOTT | AR | 7.24548e+08 | NONE | RETIRED | |
| 4 | C00410118 | P20002978 | Bachmann, Michelle | WARDENBURG, HAROLD | HOT SPRINGS NATION | AR | 7.19016e+08 | NONE | RETIRED | |

In [42]: `donor_df['contb_receipt_amt'].value_counts()`

```
Out[42]: 100.00    178188
50.00     137584
25.00     110345
250.00     91182
500.00     57984
           ...
97.15          1
122.32         1
188.65         1
122.40         1
132.12         1
Name: contb_receipt_amt, Length: 8079, dtype: int64
```

In [44]: 
```
don_mean =donor_df['contb_receipt_amt'].mean()
don_std =donor_df['contb_receipt_amt'].std()
print( 'the avrg donation was %.2f with a std %.2f' %(don_mean,don_std))
```

the avrg donation was 298.24 with a std 3749.67

In [56]: 
```
top_donor= donor_df['contb_receipt_amt'].copy()
top_donor.sort_values(ascending=False)
top_donor
```

```
Out[56]: 0          250.0
1           50.0
2          250.0
3          250.0
```

```
4           300.0
            ...
1001726    5000.0
1001727    2500.0
1001728     500.0
1001729     500.0
1001730    2500.0
Name: contb_receipt_amt, Length: 1001731, dtype: float64
```

In [60]:
```python
top_donor =top_donor[top_donor > 0]
top_donor.sort_values(ascending=False)
top_donor.value_counts().head(10)
```

Out[60]:
```
100.0     178188
50.0      137584
25.0      110345
250.0      91182
500.0      57984
2500.0     49005
35.0       37237
1000.0     36494
10.0       33986
200.0      27813
Name: contb_receipt_amt, dtype: int64
```
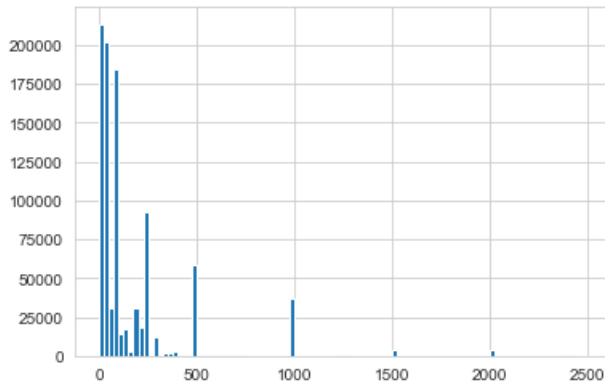
In [63]:
```python
com_don = top_donor[top_donor<2500]
com_don.hist(bins=100)
```

Out[63]: &lt;AxesSubplot:&gt;



In [68]:
```python
candidate = donor_df.cand_nm.unique()
candidate
```

Out[68]:
```
array(['Bachmann, Michelle', 'Romney, Mitt', 'Obama, Barack',
       "Roemer, Charles E. 'Buddy' III", 'Pawlenty, Timothy',
       'Johnson, Gary Earl', 'Paul, Ron', 'Santorum, Rick',
       'Cain, Herman', 'Gingrich, Newt', 'McCotter, Thaddeus G',
       'Huntsman, Jon', 'Perry, Rick'], dtype=object)
```

In [70]:
```python
party_map = {'Bachmann, Michelle': 'Republican',
             'Cain, Herman': 'Republican',
             'Gingrich, Newt': 'Republican',
             'Huntsman, Jon': 'Republican',
             'Johnson, Gary Earl': 'Republican',
             'McCotter, Thaddeus G': 'Republican',
             'Obama, Barack': 'Democrat',
             'Paul, Ron': 'Republican',
             'Pawlenty, Timothy': 'Republican',
             'Perry, Rick': 'Republican',
             "Roemer, Charles E. 'Buddy' III": 'Republican',
             'Romney, Mitt': 'Republican',
             'Santorum, Rick': 'Republican'}

donor_df['Party'] = donor_df.cand_nm.map(party_map)
```

In [72]:
```python
donor_df =donor_df[donor_df.contb_receipt_amt > 0]
```

In [74]:
```python
donor_df.head()
```

Out[74]:

| | cmte_id | cand_id | cand_nm | contbr_nm | contbr_city | contbr_st | contbr_zip | contbr_employer | contbr_occupation | contb_ |
|---|---|---|---|---|---|---|---|---|---|---|

| | cmte_id | cand_id | cand_nm | contbr_nm | contbr_city | contbr_st | contbr_zip | contbr_employer | contbr_occupation | contb_ |
|---|---------|---------|---------|-----------|-------------|-----------|------------|-----------------|-------------------|--------|
| 0 | C00410118 | P20002978 | Bachmann, Michelle | HARVEY, WILLIAM | MOBILE | AL | 3.6601e+08 | RETIRED | RETIRED | |
| 1 | C00410118 | P20002978 | Bachmann, Michelle | HARVEY, WILLIAM | MOBILE | AL | 3.6601e+08 | RETIRED | RETIRED | |
| 2 | C00410118 | P20002978 | Bachmann, Michelle | SMITH, LANIER | LANETT | AL | 3.68633e+08 | INFORMATION REQUESTED | INFORMATION REQUESTED | |
| 3 | C00410118 | P20002978 | Bachmann, Michelle | BLEVINS, DARONDA | PIGGOTT | AR | 7.24548e+08 | NONE | RETIRED | |
| 4 | C00410118 | P20002978 | Bachmann, Michelle | WARDENBURG, HAROLD | HOT SPRINGS NATION | AR | 7.19016e+08 | NONE | RETIRED | |

In [76]:
```python
donor_df.groupby('cand_nm')['contb_receipt_amt'].count()
```

Out[76]:
```
cand_nm
Bachmann, Michelle              13082
Cain, Herman                    20052
Gingrich, Newt                  46883
Huntsman, Jon                    4066
Johnson, Gary Earl               1234
McCotter, Thaddeus G               73
Obama, Barack                  589127
Paul, Ron                      143161
Pawlenty, Timothy                3844
Perry, Rick                     12709
Roemer, Charles E. 'Buddy' III   5844
Romney, Mitt                   105155
Santorum, Rick                  46245
Name: contb_receipt_amt, dtype: int64
```

In [78]:
```python
donor_df.groupby('cand_nm')['contb_receipt_amt'].sum()
```

Out[78]:
```
cand_nm
Bachmann, Michelle             2.711439e+06
Cain, Herman                   7.101082e+06
Gingrich, Newt                 1.283277e+07
Huntsman, Jon                  3.330373e+06
Johnson, Gary Earl             5.669616e+05
McCotter, Thaddeus G           3.903000e+04
Obama, Barack                  1.358774e+08
Paul, Ron                      2.100962e+07
Pawlenty, Timothy              6.004819e+06
Perry, Rick                    2.030575e+07
Roemer, Charles E. 'Buddy' III 3.730099e+05
Romney, Mitt                   8.833591e+07
Santorum, Rick                 1.104316e+07
Name: contb_receipt_amt, dtype: float64
```

In [81]:
```python
cand_amount = donor_df.groupby('cand_nm')['contb_receipt_amt'].sum()
i=0
for don in cand_amount:
    print('the candidate %s raise %.0f dollars' %(cand_amount.index[i],don))
    print('\n')
    i+=1
```

the candidate Bachmann, Michelle raise 2711439 dollars


the candidate Cain, Herman raise 7101082 dollars


the candidate Gingrich, Newt raise 12832770 dollars


the candidate Huntsman, Jon raise 3330373 dollars


the candidate Johnson, Gary Earl raise 566962 dollars


the candidate McCotter, Thaddeus G raise 39030 dollars

the candidate Obama, Barack raise 135877427 dollars

the candidate Paul, Ron raise 21009620 dollars

the candidate Pawlenty, Timothy raise 6004819 dollars

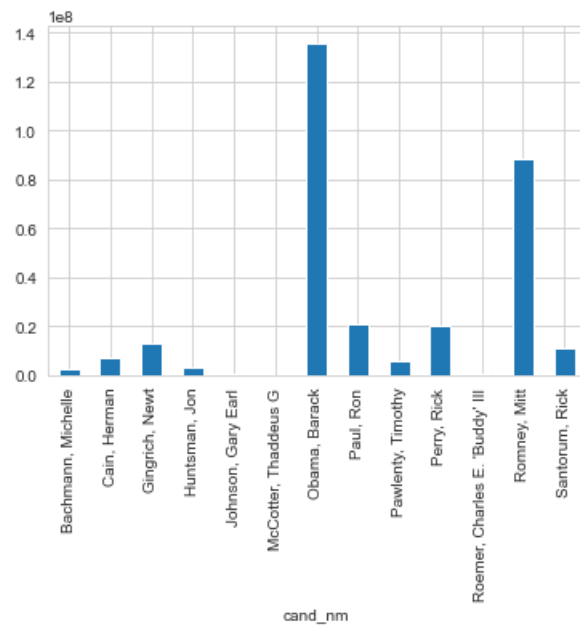the candidate Perry, Rick raise 20305754 dollars

the candidate Roemer, Charles E. 'Buddy' III raise 373010 dollars

the candidate Romney, Mitt raise 88335908 dollars
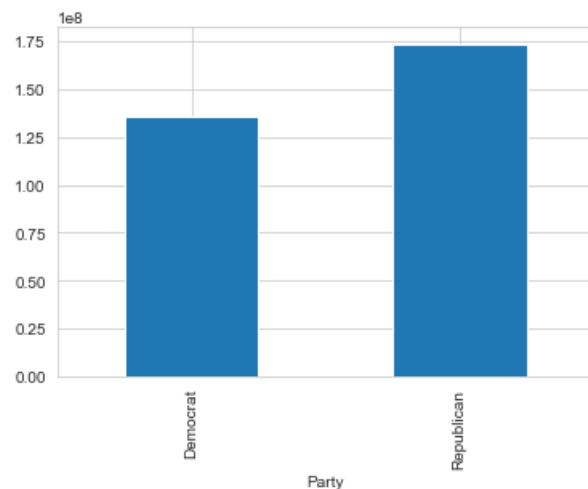
the candidate Santorum, Rick raise 11043159 dollars

In [83]:
```python
cand_amount.plot(kind='bar')
```

Out[83]: <AxesSubplot:xlabel='cand_nm'>



In [85]:
```python
donor_df.groupby('Party')['contb_receipt_amt'].sum().plot(kind='bar')
```

Out[85]: <AxesSubplot:xlabel='Party'>



In [88]:
```python
occupation_df =donor_df.pivot_table('contb_receipt_amt', index='contbr_occupation', columns ='Party', aggfunc='su
```

```
occupation_df
```

Out[88]:

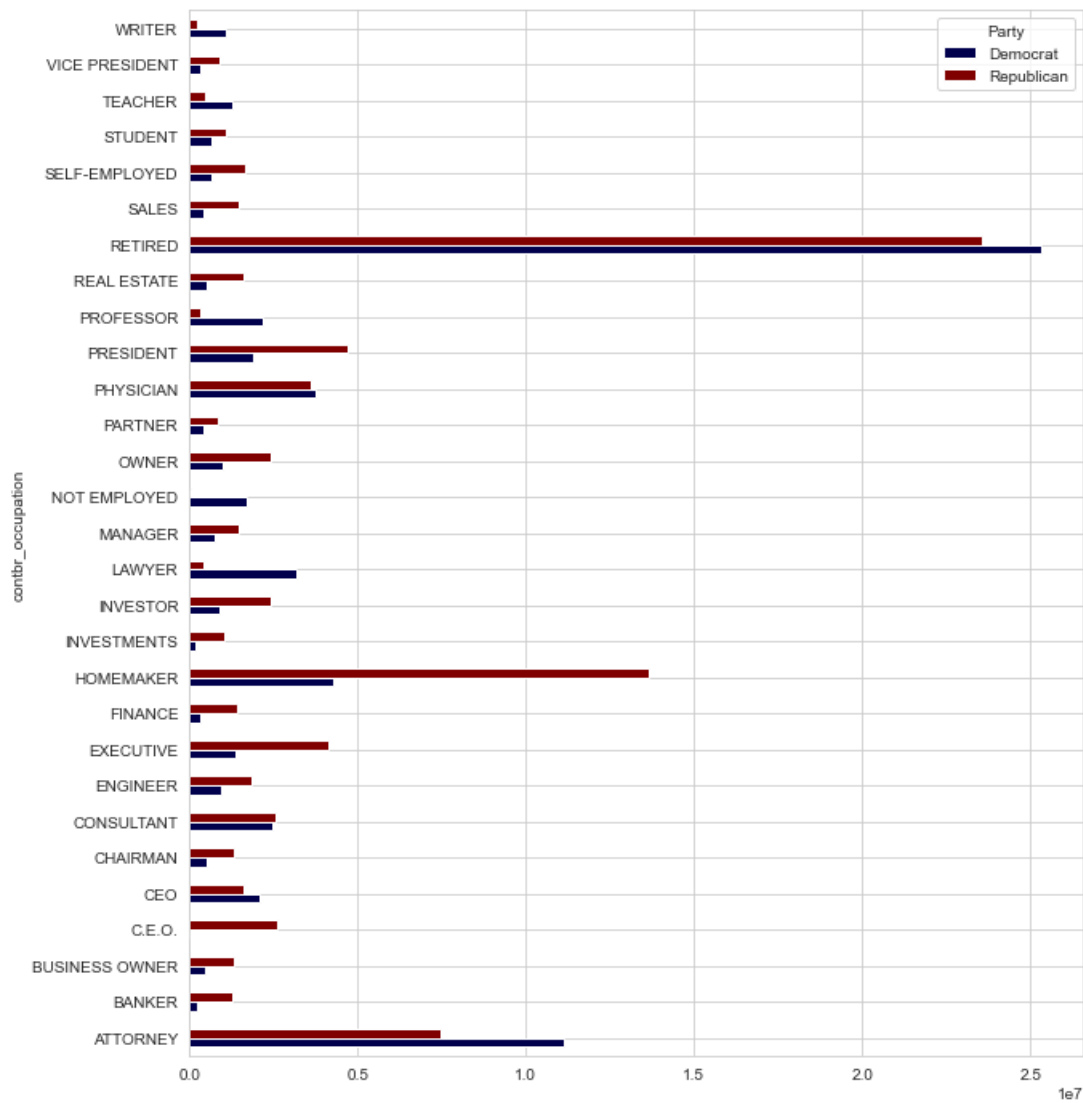| Party | Democrat | Republican |
|---|---|---|
| **contbr_occupation** | | |
| **MIXED-MEDIA ARTIST / STORYTELLER** | 100.0 | NaN |
| **AREA VICE PRESIDENT** | 250.0 | NaN |
| **RESEARCH ASSOCIATE** | 100.0 | NaN |
| **TEACHER** | 500.0 | NaN |
| **THERAPIST** | 3900.0 | NaN |
| **...** | ... | ... |
| **ZOOKEEPER** | 35.0 | NaN |
| **ZOOLOGIST** | 400.0 | NaN |
| **ZOOLOGY EDUCATION** | 25.0 | NaN |
| **\NONE\** | NaN | 250.0 |
| **~** | NaN | 75.0 |

45067 rows × 2 columns

In [92]:
```python
occupation_df.shape
```

Out[92]: (45067, 2)

In [106…
```python
occupation_df = occupation_df[occupation_df.sum(1) > 1000000]
occupation_df.shape
occupation_df.plot(kind='barh',figsize=(10,12),cmap='seismic')
```

Out[106… <AxesSubplot:ylabel='contbr_occupation'>

```python
# occupation_df.drop(['INFORMATION REQUESTED PER BEST EFFORTS','INFORMATION REQUESTED'],axis=0,inplace=True)
occupation_df.loc['CEO'] =occupation_df.loc['CEO']+ occupation_df.loc['C.E.O.']
occupation_df.drop('C.E.O.',inplace=True)
occupation_df.plot(kind='barh',figsize=(10,12),cmap='seismic')
```

Out[108... <AxesSubplot:ylabel='contbr_occupation'>