First Edition

# Data Visualization in R and python

## Leela S. Dodda

# *Contents*

# II R recipies 22

# Introduction

> *This is one of my finer quotations.*
> *–John Smith*

This is a great place to write an introduction or prologue[1].

---

[1] *You can even use a footnote to seem smarter*

# Part I

# PYTHON RECIPIES

*Chapter 1*

# *Linear Regression*

Packages required to run this code

**pandas**  for reading csv files[1] format

**scipy**  for doing linear regression analysis and obtaining the statistics

**matplotlib**  for making the plots

---

[1] *data not shown as tables*

```
1   import matplotlib
2   matplotlib.use('Agg')
3   matplotlib.rc('font', family='serif')
4   import pandas as pd
5   import matplotlib.pyplot as plt
6   from scipy import stats
7   cm5=pd.read_csv("CM5_ARRANGED_DATA_FROM_R.csv")
8   cm1=pd.read_csv("CM1A_ARRANGED_DATA_FROM_R.csv")
9   xcm5=cm5['Expt']
10  ycm5=cm5['G121']
11  xcm1=cm1['Expt']
12  ycm1=cm1['G105']
13  m1,c1,r1,p1,se1=stats.linregress(xcm1,ycm1)
14  m5,c5,r5,p5,se5=stats.linregress(xcm5,ycm5)
15  fig=plt.figure(figsize=(10, 5),dpi=300)
16  ax1 = plt.subplot(121)
17  cm1lab="$"+('y=%2.2fx+%2.2f, r^2=%1.1f'%(m1,c1,r1**2))
        +"$"
18  ax1.plot(xcm1,ycm1,'^',mfc='none',mec='b',mew=1.2)
19  ax1.plot(xcm1, m1*xcm1+c1,'k—',linewidth=2,label=
        cm1lab)
20  plt.grid()
21  plt.ylabel(r'$\Delta G^{GB/SA}_{hyd}~$ 1.05*CM1A (kcal
        /mol)',fontsize=16)
22  plt.xlabel(r'$\Delta G^{Expt}_{hyd}~$ (kcal/mol)',
        fontsize=16)
23  ax1.legend( loc='upper left')
24  ax2 = plt.subplot(122)
25  cm5lab="$"+('y=%2.2fx+%2.2f, r^2=%1.1f'%(m5,c5,r5**2))
        +"$"
26  ax2.plot(xcm5,ycm5,'o',mfc='none',mec='r',mew=1.2)
27  ax2.plot(xcm5, m5*xcm5+c5,'k—',linewidth=2,label=
        cm5lab)
28  ax2.legend( loc='upper left')
29  plt.ylabel(r'$\Delta G^{GB/SA}_{hyd}~$ 1.21*CM5 (kcal/
        mol)',fontsize=16)
```
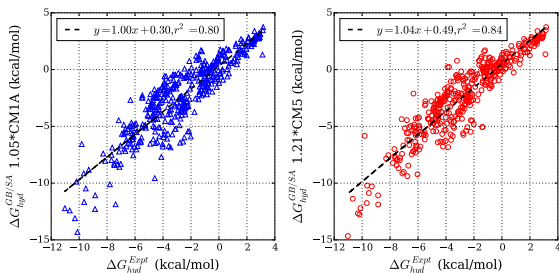
Figure 1.1: *Linear regression analysis has been performed for two sets of data and the resulting model is shown in the legends of each figure*

```
30  plt.xlabel(r'$\Delta G^{Expt}_{hyd}~$ (kcal/mol)',
            fontsize=16)
31  plt.grid()
32  fig.subplots_adjust(left   = 0.15,hspace = .001)
33  fig.tight_layout()
34  plt.savefig('GBSA_comp.pdf')
```

*Chapter 2*

# Heat Maps

*Chapter 3*

# *Barplots*

Packages required to run this code

**pandas**  for reading "Hvap.csv[1]" format

**numpy**  for creating and manipulating vectors

**matplotlib**  for making the plots

---

[1] *contains both the raw and devation data required for plot*

Table 3.1: *Data to be plotted using bar plots*

| Molecules | OPLS | CM1A | CM5 | Expt |
|---|---|---|---|---|
| Acetic acid | 12.26 | 13.52 | 14.46 | 12.49 |
| Acetone | 7.23 | 7.74 | 8.92 | 7.48 |
| Acetonitrile | 7.57 | 7.63 | 9.76 | 8.01 |
| Aniline | 11.88 | 16.41 | 14.61 | 12.60 |
| Benzonitrile | 12.52 | 14.45 | 15.49 | 12.54 |
| Cyclohexane | 7.56 | 7.64 | 7.61 | 7.86 |
| Diethylamine | 7.68 | 7.54 | 7.46 | 7.48 |
| Diethyl ether | 6.90 | 7.01 | 7.22 | 6.56 |
| N,N-dimethylacetamide | 13.44 | 14.34 | 15.57 | 11.75 |
| Ethanethiol | 6.67 | 6.48 | 6.68 | 6.58 |
| Ethanol | 10.29 | 9.06 | 10.19 | 10.11 |
| Furan | 6.91 | 8.01 | 7.17 | 6.56 |
| Hexane | 7.54 | 7.48 | 7.34 | 7.54 |
| Methanol | 9.00 | 7.60 | 8.84 | 8.95 |
| Methyl acetate | 7.99 | 10.00 | 10.12 | 7.72 |
| Nitroethane | 9.78 | 14.16 | 11.72 | 9.94 |
| N-methylacetamide | 13.87 | 16.12 | 19.06 | 13.30 |
| Phenol | 14.58 | 14.63 | 14.30 | 13.82 |
| Propylamine | 7.90 | 8.93 | 7.23 | 7.47 |
| Pyridine | 9.76 | 11.16 | 11.16 | 9.61 |
| Pyrrole | 10.32 | 13.81 | 12.37 | 10.80 |
| Tetrahydrofuran | 7.52 | 7.66 | 8.08 | 7.61 |

```python
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  hvap = pd.read_csv("Hvap.csv")
5  n_groups = len(hvap.D_OPLS)
6
7  opls = list(hvap.D_OPLS)
8  x_lab = list(hvap.Molecules)
9  cm5 = list(hvap.D_CM5)
10 cm1a = list(hvap.D_CM1A)
11 fig, ax = plt.subplots()
12 index = np.arange(n_groups)
13 bar_width = 0.33
14 opacity = 0.5
15 rects1 = plt.bar(index, opls, bar_width,
16                  alpha=opacity, color='r', label='OPLS')
17 rects2 = plt.bar(index + bar_width, cm5, bar_width,
18                  alpha=opacity, color='g', label='1.27*
    CM5')
19 rects3 = plt.bar(index + 2 * bar_width, cm1a,
    bar_width,
20                  alpha=opacity, color='b', label='1.14*
    CM1A')
21 plt.ylabel(r'$\Delta H_{vap}^{expt}-\Delta H_{vap}^{
    calc}$ (kcal/mol)')
22 plt.xticks(index + bar_width, x_lab, rotation=90)
23 plt.grid()
24 plt.xlim(-0.5, n_groups + 0.5)
25 plt.legend(loc='lower left', ncol=3)
26 plt.tight_layout()
27 plt.savefig("Tesh_hvap.pdf")
```

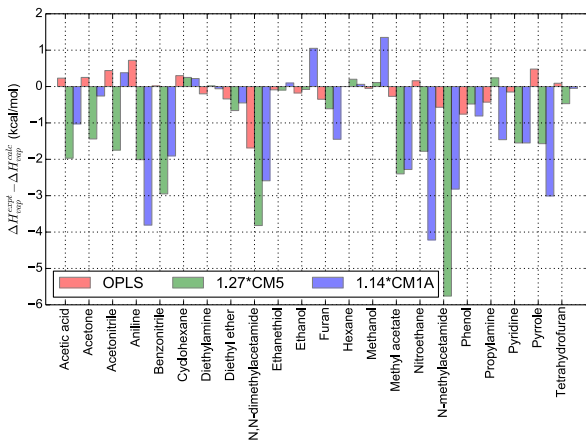Listing 3.1: *Bar plot of the data shown in Table above*

Figure 3.1: *Data in table above is plotted where instead of raw data, deviations from experiments for each method is plotted*

```
 1  import matplotlib
 2  matplotlib.use('Agg')
 3  matplotlib.rc('font', family='serif')
 4  from matplotlib import pylab
 5  from pylab import rcParams
 6  rcParams['figure.figsize'] = 11, 7
 7  import pandas as pd
 8  import matplotlib.pyplot as plt
 9  import numpy as np
10  import matplotlib.cm as cm
11  dat = pd.read_csv('all_cm5_dat.csv')
12  method=list(dat['Molecules'])
13  ################################
14  legend= (dat.columns.values)[1:]
15  #colors = cm.Greens(np.linspace(0, 1, len(legend)))
16  colors = cm.Spectral(np.linspace(0, 1, len(legend)))
17  index = np.arange(len(method))
18  bar_width = 1.0/len(legend)
19  opacity = 1.0
20  for i,c in zip(range(0,len(legend)),colors):
21    plt.bar(index+bar_width*i,dat[legend[i]] , bar_width
       ,
22                      alpha=opacity,
23                      color=c,
24                      label=legend[i])
25  plt.ylabel(r'$\Delta H_{vap}^{expt}-\Delta H_{vap}^{
       calc}~$ (kcal/mol)',fontsize=14)
26  plt.xticks(index + bar_width*len(legend)/2, method,
       rotation=90,fontsize=12)
27  plt.grid()
28  plt.xlim(-0.1, len(method) + 0.0)
29  #plt.legend(bbox_to_anchor=(1.0, 1.01),fontsize=9,loc
       =0,frameon=False)
30  plt.legend(loc='upper left', ncol=6,fontsize=10,
       frameon=False)
31  plt.tight_layout(rect=[0.0,0.01,0.99,1])
```
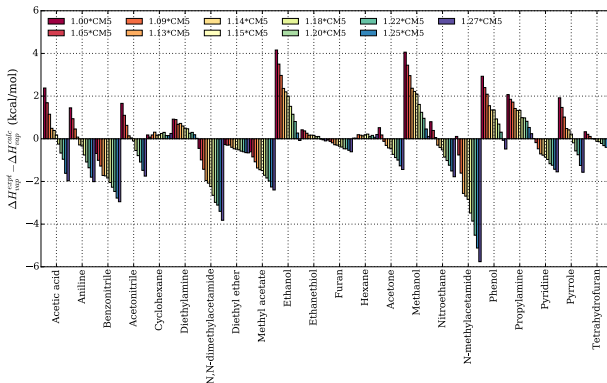
Figure 3.2: *Data in table above is plotted where instead of raw data, deviations from experiments for each method is plotted*

```
32  plt.savefig("Thh.pdf")
```

Listing 3.2: *Barplot liquid properties using CM5 charges with different scale factors*

```
1  import matplotlib
2  matplotlib.use('Agg')
3  matplotlib.rc('font', family='serif')
4  import pandas as pd
5  import matplotlib.pyplot as plt
6  import numpy as np
7  import matplotlib.cm as cm
8  from matplotlib.ticker import FuncFormatter
9  hvap = pd.read_csv('Hvap.csv')
10 den = pd.read_csv('Den.csv')
11 hvap.drop(hvap.columns[[1,2,3,4,5,7]], axis=1, inplace
       =True)
12 den.drop(den.columns[[1,2,3,4,5,7]], axis=1, inplace=
       True)
13 m1=list(hvap['Molecules'])
14 m2=list(den['Molecules'])
15 ###############################
16 def millions(x, pos):
17     'The two args are the value and tick position'
18     return '%2.2f' % (x)
19 formatter = FuncFormatter(millions)
20 ###############################
21 l1 = (hvap.columns.values)[1:]
22 l2 = (den.columns.values)[1:]
23 colors = cm.Spectral(np.linspace(0, 1, len(l1)))
24 fig,(ax1, ax2) = plt.subplots(2, sharex=True)
25 index = np.arange(len(m1))
26 bar_width = 1.0/len(l1)
27 opacity = 1.0
28 patterns = [ "*", "o", "."]
29 for i,c in zip(range(0,len(l1)),colors):
30   ax1.bar(index+bar_width*i,hvap[l1[i]] , bar_width,
31                     alpha=opacity,
32                     color=c,
33       hatch=patterns[i],
34                     label=l1[i][2:])
```

```
35  ax1.legend(fontsize=10,loc=9, bbox_to_anchor=(0.5,
        1.2),ncol=3,frameon=False)
36  ax1.set_ylabel(r'$\Delta H_{vap}^{expt}-\Delta H_{vap
        }^{calc}~$ (kcal/mol)')
37  ax1.yaxis.set_major_formatter(formatter)
38  ax1.yaxis.set_ticks(np.arange(-6,3,2))
39  for i,c in zip(range(0,len(l2)),colors):
40    ax2.bar(index+bar_width*i,den[l1[i]] , bar_width,
41                    alpha=opacity,
42                    color=c,
43        hatch=patterns[i],
44                    label=l2[i][2:])
45  ax2.set_ylabel(r'$\Delta \rho^{expt}-\Delta \rho^{calc
        }~$ (g/cc)')
46  ax2.yaxis.set_ticks(np.arange(-0.1,0.08,0.04))
47  plt.xticks(index + bar_width*len(l2)/2, m2, rotation
        =90,fontsize=10)
48  ax1.xaxis.grid()
49  ax2.xaxis.grid()
50  ax1.yaxis.grid()
51  ax2.yaxis.grid()
52  ax1.set_xlim(-0.1, len(m1) + 0.0)
53  ax2.set_xlim(-0.1, len(m2) + 0.0)
54  plt.tight_layout(rect=[0.0,0.01,0.89,0.99])
55  plt.savefig("Multi_bar.pdf")
```

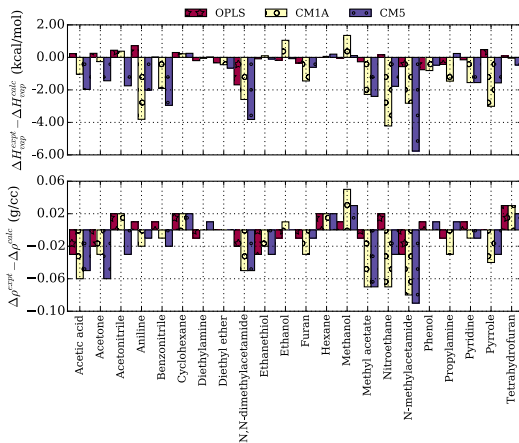Listing 3.3: *Multiple bar plots in matplotlib*

Figure 3.3: *Data in table above is plotted where instead of raw data, deviations from experiments for each method is plotted*

```
1  import matplotlib
2  matplotlib.use('Agg')
3  matplotlib.rc('font', family='serif')
4  import pandas as pd
5  import matplotlib.pyplot as plt
6  import numpy as np
7  import matplotlib.cm as cm
8  import seaborn as sns
9  sns.set()
10 from matplotlib.ticker import FuncFormatter
11 hvap = pd.read_csv('Hvap.csv')
12 den = pd.read_csv('Den.csv')
13 d_std=pd.read_csv('STD_Den.csv')
14 h_std=pd.read_csv('STD_Hvap.csv')
15 h_std['CM5']=h_std.CM5_127
16 hvap.drop(hvap.columns[[1,2,3,4,5,7]], axis=1, inplace
       =True)
17 den.drop(den.columns[[1,2,3,4,5,7]], axis=1, inplace=
       True)
18 m1=list(hvap['Molecules'])
19 m2=list(den['Molecules'])
20 ################################
21 def millions(x, pos):
22     'The two args are the value and tick position'
23     return '%2.2f' % (x)
24 formatter = FuncFormatter(millions)
25 ################################
26 l1= (hvap.columns.values)[1:]
27 l2= (den.columns.values)[1:]
28 colors = cm.rainbow(np.linspace(0, 1, len(l1)))
29 fig ,(ax1, ax2) = plt.subplots(2, sharex=True)
30 index = np.arange(len(m1))
31 bar_width = 1.0/len(l1)
32 opacity = 1.0
33 patterns = [ "*", "o", "."]
34 for i,c in zip(range(0,len(l1)),colors):
```

```
35    ax1.bar(index+bar_width*i,hvap[l1[i]]  ,bar_width,
36                        alpha=opacity,
37                        facecolor=c,
38          edgecolor = c,
39 #        hatch=patterns[i],
40                        label=l1[i][2:],yerr=h_std[l1[i
       ][2:]],ecolor='k',capsize=1)
41 ax1.legend(fontsize=10,loc=9, bbox_to_anchor=(0.5,
       1.2),ncol=3,frameon=False)
42 ax1.set_ylabel(r'$\Delta H_{vap}^{expt}−\Delta H_{vap
       }^{calc}~$ (kcal/mol)')
43 ax1.yaxis.set_major_formatter(formatter)
44 ax1.yaxis.set_ticks(np.arange(−6,3,2))
45 for i,c in zip(range(0,len(l2)),colors):
46   ax2.bar(index+bar_width*i,den[l1[i]]  ,bar_width,
47                        alpha=opacity,
48                        facecolor=c,
49          edgecolor=c,
50          #hatch=patterns[i],
51                        label=l2[i][2:],yerr=d_std[l1[i
       ][2:]],ecolor='k',capsize=1)
52 ax2.set_ylabel(r'$\Delta \rho^{expt}−\Delta \rho^{calc
       }~$ (g/cc)')
53 ax2.yaxis.set_ticks(np.arange(−0.1,0.08,0.04))
54 plt.xticks(index + bar_width*len(l2)/2, m2, rotation
       =90,fontsize=10)
55 #ax1.xaxis.grid()
56 #ax2.xaxis.grid()
57 #ax1.yaxis.grid()
58 #ax2.yaxis.grid()
59 ax1.set_xlim(−0.1, len(m1) + 0.0)
60 ax2.set_xlim(−0.1, len(m2) + 0.0)
61 plt.tight_layout(rect=[0.0,0.01,0.89,0.99])
62 plt.savefig("Ers_Multi_bar.pdf")
```

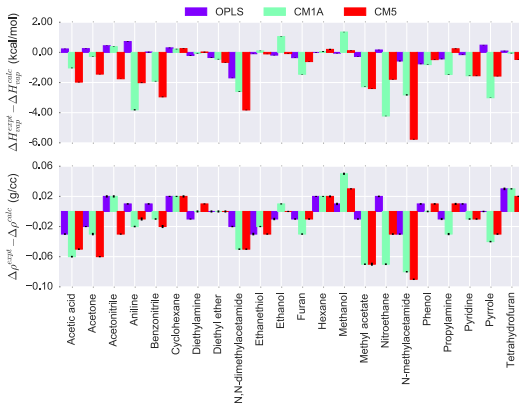Listing 3.4: *Multiple bar plots in matplotlib*

Figure 3.4: *Data in table above is plotted where instead of raw data, deviations from experiments for each method is plotted*

# Part II

# R RECIPIES