

Development of web-based tools with RDKit & Django

Paul Czodrowski, Merck KGaA Darmstadt
Small Molecule Platform, Computational Chemistry
`paul.czodrowski@merckgroup.com`

RDKit User Group Meeting, London / October 2012

a 3-class aqueous solubility model

- Huuskonen dataset
(<http://dx.doi.org/10.1021/ci9901338> &
<http://www.cheminformatics.org/datasets/huuskonen/index.html>)
- Training set: 1025 compounds
- Test set: 257 compounds
- Classification based on aqueous solubility

aqueous solubility	solubility class	# compounds
≤ -3.0	0 [low]	417
> -3.0 & ≤ -1.0	1 [medium]	402
> -1.0	2 [high]	206

code snippet – calculate the descriptors

fingerprint
calculation

```
descr_list = [x[0] for x in Descriptors._descList]
descr_list.remove('MolecularFormula')
descr_obj=MoleculeDescriptors.MolecularDescriptorCalculator(descr_list)
calc_descrs = [descriptor_object.CalcDescriptors(x) for x in cpds]

pts = []
for i,m in enumerate(ms):
    if m.GetProp('SOL_classification')== '(A) low':
        act=2
    elif m.GetProp('SOL_classification')== '(B) medium':
        act=1
    else:
        act=0
    pts.append([m.GetProp('NAME'),list(calc_descrs[i]),act])
```

compound name

list of descriptors

activity

code snippet – model statistics

```
res = ScreenComposite.ShowVoteResults(range(len(pts)), pts, cmp, \
    3, 0, errorEstimate=True)
```

```

*** Vote Results ***
misclassified: 245/1025 (%23.90)      245/1025 (%23.90)

average correct confidence:    0.9543
average incorrect confidence:  0.8496

Results Table:
-----
      353      57      1      | 84.45
      64     333     111     | 82.63
       0      12      94      | 45.41
-----
      84.65     82.84     45.63
    
```

performance
on the test set

		experiment		
		class0	class1	class2
prediction	class0	89	17	0
	class1	13	90	7
	class2	0	8	33

overall accuracy: 83 %

code snippet – similarity search

fingerprints
of training set

```
fps = [AllChem.GetMorganFingerprintAsBitVect(x,2,2048) for x in ms]  
pts = []  
for i,m in enumerate(ms):  
    pts.append([m,m.GetProp('SOL_classification'),fps[i]])  
cPickle.dump(pts,file('cpds.pkl','wb+'))
```

similarity search

```
sol_trainset = cPickle.load('cpds.pkl', "rb")  
  
top3 = TopNContainer(3)  
refFP = AllChem.GetMorganFingerprintAsBitVect(mol,2,2048)  
hits = top3.GetPts()  
  
for datasetmol in sol_trainset:  
    sol_TL = datasetmol[1]  
    fp = datasetmol[2]  
    tani = DataStructs.TanimotoSimilarity(refFP,fp)  
    top3.Insert((tani,sol_TL,datasetmol[0]))  
  
hits = top3.GetPts()  
hits.reverse()
```

how-to-set-up a simple website

```
urls.py
settings.py
manage.py
media/
templates/start_site.html
templates/results.html
webinter/views.py
```

development environment	productive environment
Comes with django installation	Apache/WSGI

pitfalls Django/WSGI/SLES

- our setup:
 - Apache server 2.2.12 on SLES 11SP1 (64 bit) with a separate virtualHost configured (to have it answering on and with an application-specific hostname)
 - python (compiled by user)
- WSGIPythonHome in `django.wsgi` not properly recognized
- `LD_LIBRARY_PATH` setting of original Apache shall not be modified
- necessary to run *distinct* Apache (`httpd-django.conf` besides `httpd.conf`)

Any experience with virtualenv?

django code snippets

urls.py

```
urlpatterns = patterns('webinter.views',
    (r'^$', "jmecall"),
    (r'^chime4sol/$', "chime4sol"),
    (r'^get_sol/$', "get_sol"),
    (r'^jme/$', "jmecall"),
    (r'^draw_molecule/$', "draw_molecule"),
)
```

settings.py

```
FILE_CMP = cPickle.load(file(os.path.join(MEDIA_ROOT, \
    'cmp.sol_model_descr.pkl'), "rb"))
FILE_CPDS = cPickle.load(file(os.path.join(MEDIA_ROOT, \
    'cpds.MorganFP.pkl'), "rb"))
```

views.py

```
def get_sol(request):
    cmp = settings.FILE_CMP
    sol_trainset = settings.FILE_CPDS
    top3 = TopNContainer(3)
    errors = []
    if not 'q' in request.GET:
        errors.append('No query provided')
    else:
        query=request.GET['q']
        qj = QueryJob(inputtext=query, username=request.user.username)
        qj.setJobid()
        for word in query.split():
            mol = Chem.MolFromSmiles(str(query))
            mol.SetProp('_Name', 'nothing')
            #
            # nearest neighbor
            refFP = AllChem.GetMorganFingerprintAsBitVect(mol, 2, 2048)
            for datasetmol in sol_trainset:
                sol_TL = datasetmol[1]
                fp = datasetmol[2]
                tani = DataStructs.TanimotoSimilarity(refFP, fp)
                top3.Insert((tani, sol_TL, datasetmol[0]))
            hits = top3.GetPts()
            hit_list_dict = {}
            for cnt, sim_cpd in enumerate( hits ):
                skey = "image_query_cpd_%s" % cnt
                request.session[skey] = sim_cpd[2]
                if sim_cpd[1] == "(A) low":
                    sol_TL = "0"
                elif sim_cpd[1] == "(B) medium":
                    sol_TL = "1"
                else:
                    sol_TL = "2"
```

results.html

```
{% for key, value in hit_list_dict.items %}
<tr align=center>
  <td> </td>
  <td {% if value.2 == 0 %} bgcolor="orangered"
    {% else %} bgcolor="lawngreen" {% endif %}> {{ value.0}}</td>
  <td>{{ value.1 }}</td>
</tr>
{% endfor %}
```

Acknowledgments

- Daniel Kuhn
- Michael Krug
- Christian Griebel, Frank Morawietz, Trung Phan An
- Friedrich Rippmann
- Klaus Urbahns

- Greg Landrum
- RDKit mailing list!