

# 3\_MoleculeSetsAndSimilarityExample

October 23, 2019

## 1 Molecule sets and similarity example

Andrea Volkamer

### 1.1 Working with Tables and csv Files

```
[1]: import pandas as pd
df = pd.read_csv('./data/EGFR-course.csv', delimiter=',', names=['Smiles', 'Name'], header=None)
df.head()
```

```
[1]:
```

	Smiles	Name
0	<chem>C0c1cc2ncnc(Nc3ccc(F)c(Cl)c3)c2cc1OCCCN1CCOCC1</chem>	Gefitinib
1	<chem>C#Cc1cccc(Nc2ncnc3cc(OCCOC)c(OCCOC)cc23)c1</chem>	Erlotinib
2	<chem>CS(=O)(=O)CCNCc1ccc(-c2ccc3ncnc(Nc4ccc(OCc5ccc...</chem>	Lapatinib
3	<chem>CN(C)C/C=C/C(=O)Nc1cc2c(Nc3ccc(F)c(Cl)c3)ncnc2...</chem>	Afatinib
4	<chem>C=CC(=O)Nc1cc(Nc2nccc(-c3cn(C)c4cccc34)n2)c(O...</chem>	Osimertinib

```
[2]: df.shape
```

```
[2]: (5, 2)
```

```
[3]: df.columns
```

```
[3]: Index(['Smiles', 'Name'], dtype='object')
```

```
[4]: from rdkit import Chem
from rdkit.Chem.Draw import IPythonConsole
from rdkit.Chem import Draw, PandasTools

PandasTools.AddMoleculeColumnToFrame(df, smilesCol='Smiles')

# Draw molecules
df[['Name', 'ROMol']]
```

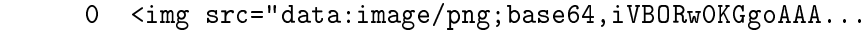
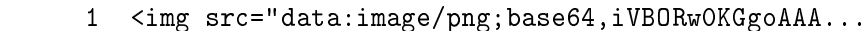
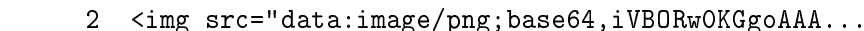
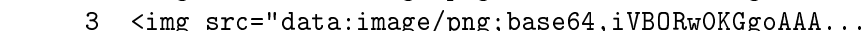
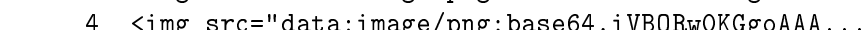
```
[4]:
```

	Name	ROMol	
0	Gefitinib	COc1cc2ncnc(Nc3ccc(F)c(Cl)c3)c2cc1OCCCN1CCOCC1</chem>	Gefitinib
1	<chem>C#Cc1cccc(Nc2ncnc3cc(OCCOC)c(OCCOC)cc23)c1</chem>	Erlotinib	
2	<chem>CS(=O)(=O)CCNCc1ccc(-c2ccc3ncnc(Nc4ccc(OCc5ccc...</chem>	Lapatinib	
3	<chem>CN(C)C/C=C/C(=O)Nc1cc2c(Nc3ccc(F)c(Cl)c3)ncnc2...</chem>	Afatinib	
4	<chem>C=CC(=O)Nc1cc(Nc2nccc(-c3cn(C)c4cccc34)n2)c(O...</chem>	Osimertinib	

	ROMol	HeavyAtoms
0	COc1cc2ncnc(Nc3ccc(F)c(Cl)c3)c2cc1OCCCN1CCOCC1</chem>	Gefitinib
1	<chem>C#Cc1cccc(Nc2ncnc3cc(OCCOC)c(OCCOC)cc23)c1</chem>	Erlotinib
2	<chem>CS(=O)(=O)CCNCc1ccc(-c2ccc3ncnc(Nc4ccc(OCc5ccc...</chem>	Lapatinib
3	<chem>CN(C)C/C=C/C(=O)Nc1cc2c(Nc3ccc(F)c(Cl)c3)ncnc2...</chem>	Afatinib
4	<chem>C=CC(=O)Nc1cc(Nc2nccc(-c3cn(C)c4cccc34)n2)c(O...</chem>	Osimertinib

	ROMol	HeavyAtoms \
0		31
1		29
2		40
3		34
4		37

	sim2Gefitinib
0	1.000000
1	0.609195
2	0.517073
3	0.641711
4	0.371134

```
[11]: df.sort_values(['sim2Gefitinib'], inplace=True, ascending=False)
```

```
[12]: df.head()
```

```
[12]:
```

	Smiles	Name \
0	<chem>COc1cc2ncnc(Nc3ccc(F)c(Cl)c3)c2cc1OCCCN1CCOCC1</chem>	Gefitinib
3	<chem>CN(C)C/C=C/C(=O)Nc1cc2c(Nc3ccc(F)c(Cl)c3)ncnc2...</chem>	Afatinib
1	<chem>C#Cc1cccc(Nc2ncnc3cc(OCCOC)c(OCCOC)cc23)c1</chem>	Erlotinib

```

2  CS(=O)(=O)CCNCc1ccc(-c2ccc3ncnc(Nc4ccc(OCc5ccc...   Lapatinib
4  C=CC(=O)Nc1cc(Nc2nccc(-c3cn(C)c4cccc34)n2)c(0...   Osimertinib

```

	ROMol	HeavyAtoms	\
0	<img src="data:image/png;base64,iVBORwOKGgoAAA..."	31	
3	<img src="data:image/png;base64,iVBORwOKGgoAAA..."	34	
1	<img src="data:image/png;base64,iVBORwOKGgoAAA..."	29	
2	<img src="data:image/png;base64,iVBORwOKGgoAAA..."	40	
4	<img src="data:image/png;base64,iVBORwOKGgoAAA..."	37	

```

sim2Gefitinib
0      1.000000
3      0.641711
1      0.609195
2      0.517073
4      0.371134

```

```

[13]: # Save file
df.drop('ROMol', axis=1).to_csv('./data/mytest_csvFile.csv')

```

```

[:

```