Name-   Utkarsh.Arvind.Bondade

Gender-Male

Age-   21

City-Nagpur

Pursuing- BTECH from Shri Guru Gobind Singhji Institute of Engineering and Technology(Nanded)

Branch-   Electronics and Telecommunication Engineering(Final Year)

Subject- DATA SCIENCE MAJOR PROJECT

Choose any dataset of your choice ,apply a suitable algorithm(Regression/Classification) and create a model.

Topic- To build predictive model using the Delaney solubility dataset.

Dataset link

https://raw.githubusercontent.com/dataprofessor/data/master/delaney_solubility_with_descriptors.csv

+ Code  + Text

## my first ML **project**

```python
# load data
import  Follow link (ctrl + click)
url = 'https://raw.githubusercontent.com/dataprofessor/data/master/delaney_solubility_with_descriptors.csv'
df = pd.read_csv(url, error_bad_lines=False)
df
```

```
<ipython-input-3-050795025cb2>:4: FutureWarning: The error_bad_lines argument has been deprecated and will be removed in a future version. Use on_bad_lines in the future.

  df = pd.read_csv(url, error_bad_lines=False)
```

| | MolLogP | MolWt | NumRotatableBonds | AromaticProportion | logS |
|---|---|---|---|---|---|
| 0 | 2.59540 | 167.850 | 0.0 | 0.000000 | -2.180 |
| 1 | 2.37650 | 133.405 | 0.0 | 0.000000 | -2.000 |
| 2 | 2.59380 | 167.850 | 1.0 | 0.000000 | -1.740 |
| 3 | 2.02890 | 133.405 | 1.0 | 0.000000 | -1.480 |
| 4 | 2.91890 | 187.375 | 1.0 | 0.000000 | -3.040 |
| ... | ... | ... | ... | ... | ... |
| 1139 | 1.98820 | 287.343 | 8.0 | 0.000000 | 1.144 |
| 1140 | 3.42130 | 286.114 | 2.0 | 0.333333 | -4.925 |
| 1141 | 3.60960 | 308.333 | 4.0 | 0.695652 | -3.893 |
| 1142 | 2.56214 | 354.815 | 3.0 | 0.521739 | -3.790 |
| 1143 | 2.02164 | 179.219 | 1.0 | 0.461538 | -2.581 |

1144 rows × 5 columns

```python
#data preparation
#data seperation in x and y
```

```python
y=df['logS']
y
```

```
0    -2.180
```

---

+ Code  + Text  All changes saved

```python
#data preparation
#data seperation in x and y
```

```python
y=df['logS']
y
```

```
0       -2.180
1       -2.000
2       -1.740
3       -1.480
4       -3.040
         ...
1139     1.144
1140    -4.925
1141    -3.893
1142    -3.790
1143    -2.581
Name: logS, Length: 1144, dtype: float64
```

```python
x=df.drop('logS',axis=1)
x
```

| | MolLogP | MolWt | NumRotatableBonds | AromaticProportion |
|---|---|---|---|---|
| 0 | 2.59540 | 167.850 | 0.0 | 0.000000 |
| 1 | 2.37650 | 133.405 | 0.0 | 0.000000 |
| 2 | 2.59380 | 167.850 | 1.0 | 0.000000 |
| 3 | 2.02890 | 133.405 | 1.0 | 0.000000 |
| 4 | 2.91890 | 187.375 | 1.0 | 0.000000 |
| ... | ... | ... | ... | ... |
| 1139 | 1.98820 | 287.343 | 8.0 | 0.000000 |
| 1140 | 3.42130 | 286.114 | 2.0 | 0.333333 |
| 1141 | 3.60960 | 308.333 | 4.0 | 0.695652 |
| 1142 | 2.56214 | 354.815 | 3.0 | 0.521739 |
| 1143 | 2.02164 | 179.219 | 1.0 | 0.461538 |

1144 rows × 4 columns

```python
#data splitting
from sklearn.model_selection import train_test_split
```

```python
#data splitting
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=100)
x_train
```

|  | MolLogP | MolWt | NumRotatableBonds | AromaticProportion |
|---|---|---|---|---|
| 107 | 3.14280 | 112.216 | 5.0 | 0.000000 |
| 378 | -2.07850 | 142.070 | 0.0 | 0.000000 |
| 529 | -0.47730 | 168.152 | 0.0 | 0.000000 |
| 546 | -0.86740 | 154.125 | 0.0 | 0.000000 |
| 320 | 1.62150 | 100.161 | 2.0 | 0.000000 |
| ... | ... | ... | ... | ... |
| 802 | 3.00254 | 250.301 | 1.0 | 0.842105 |
| 53 | 2.13860 | 82.146 | 3.0 | 0.000000 |
| 380 | 6.76304 | 266.340 | 0.0 | 0.000000 |
| 78 | 3.89960 | 186.339 | 10.0 | 0.000000 |
| 792 | 2.52334 | 310.297 | 3.0 | 0.300000 |

916 rows × 4 columns

```python
x_test
```

|  | MolLogP | MolWt | NumRotatableBonds | AromaticProportion |
|---|---|---|---|---|
| 822 | 2.91000 | 172.268 | 7.0 | 0.000000 |
| 118 | 7.27400 | 360.882 | 1.0 | 0.666667 |
| 347 | 1.94040 | 145.161 | 0.0 | 0.909091 |
| 1123 | 1.98640 | 119.378 | 0.0 | 0.000000 |
| 924 | 1.70062 | 106.140 | 0.0 | 0.750000 |
| ... | ... | ... | ... | ... |
| 1114 | 1.76210 | 478.913 | 4.0 | 0.000000 |
| 427 | 6.32820 | 276.338 | 0.0 | 1.000000 |
| 711 | 0.04430 | 218.205 | 5.0 | 0.000000 |
| 4 | 2.91890 | 187.375 | 1.0 | 0.000000 |

Ds completed at 7:43PM

---

```python
#model building
#linear regression
```

```python
from sklearn.linear_model import LinearRegression
lr=LinearRegression()
lr.fit(x_train,y_train)
```

```
▾ LinearRegression
LinearRegression()
```

```python
#applying the model to make a prediction
y_lr_train_pred=lr.predict(x_train)
y_lr_test_pred=lr.predict(x_test)
```

```python
y_lr_train_pred
```

```
-4.27331505, -0.41608127, -0.92992939, -2.85102130, -2.38452854,
-1.35467797, -5.35250249, -4.65318325, -3.36770115, -6.54753068,
-1.21120552, -1.41277480, -0.62804611, -5.46722691, -2.00332864,
-4.78222762, -4.20128488, -3.75207911, -2.89485114, 0.44470271,
-2.75272546, -6.40483191, -3.43050076, -1.33564461, -7.30980571,
-2.3552617 , -1.97390964, -3.31913336, -0.99788515, -1.40117282,
2.54283347, -2.57933977, -1.5372532 , -1.40969589, -5.50256395,
-2.84523605, -5.79050585, -0.72972899, -4.50417995, -3.41743324,
-1.09701519, -3.04694269, -1.10327803, -1.40988982, -2.81732645,
-6.11716924, -7.16934656, -1.96622906, -3.91516316, -3.61372065,
-6.40483191, -2.05580278, -3.67078501, -5.54561991, -2.93380057,
-3.27180065, -2.72328327, -2.53353851, -1.11363181, -2.90610261,
-1.10853770, -1.49259585, -1.15572753, -4.70049144, 1.99249372,
-1.80643013, -0.46728525, -4.32045753, -1.19196129, -2.25004086,
-1.04336081, -1.87633503, -9.15826626, -3.65922507, -2.29683578,
-0.5325366 , -7.5914365 , -2.15587374, -1.90396018, -1.78665244,
-4.36090015, -1.91112045, -5.28950294, -3.70838271, -3.34373141,
-3.82343254, -3.84290966, -1.75395277, -2.02944093, -1.92755736,
-3.95353481, -2.51949142, -4.16972153, -0.47784193, -2.00094662,
-2.08049675, -1.57387693, -3.70234071, -1.27341971, -1.10845548,
-3.82977779, -1.80911411, -2.73044354, -1.76722676, -2.22965253,
-3.56321848, -6.67044148, -1.51156575, -0.24983049, -1.72672119,
2.14097281, -4.05818764, -3.37140330, -1.097265 , -1.42797042,
-3.08706826, -4.24525450, -6.52588866, -4.68029775, -4.67032068,
-2.03122625, -3.11841945, -5.04235077, 0.00553271, -8.46728525,
-4.23480474, -2.75684478, -2.14339889, -1.5231707 , -2.67587664,
-3.59418029, -5.68837433, -5.02111611, -0.94099440, -3.53200413,
-1.41666505, -1.9765392 , -2.67587664, -5.07753962, -3.73205662,
-2.65590022, -2.46014399, -3.01774254, -8.46728525, -0.43381828,
-3.81993282, -3.25038467, -5.64350181, -4.18419636, -6.47162256,
-2.29452458, -3.70387135, -6.42099768, -1.53182046, -2.45513763,
```

Ds completed at 7:43PM

+ Code  + Text   All changes saved

```
y_lr_test_pred
```

```
array([-3.05722870e+00, -7.77781827e+00, -2.55016650e+00, -2.01523582e+00,
       -2.06375090e+00, -9.09672215e-01, -5.04603364e-01, -5.53620003e-01,
       -5.72200956e+00, -3.94006681e+00, -3.95496755e+00, -1.29737000e+00,
       -1.48980354e+00, -1.48988902e+00, -4.64518000e+00, -1.90390018e+00,
       -1.55506313e+00, -3.16424685e+00, -3.70063920e+00, -5.58105660e+00,
       -3.25038467e+00, -5.04233077e+00, -5.09194881e+00, -2.14338049e+00,
       -4.35680341e+00, -5.03964756e+00, -3.10383618e+00, -4.40280964e+00,
       -4.21270272e+00,  5.50508349e-01, -1.45537678e+00, -4.41027396e+00,
       -2.59660773e+00, -1.53336276e+00, -5.55749874e-01, -1.67311795e+00,
       -2.78163675e+00, -3.15393565e+00, -3.27083301e+00, -1.75321446e+00,
       -1.53350725e+00, -2.01255666e+00, -6.57559107e+00, -7.09433046e+00,
       -5.76437127e+00, -4.16420868e+00, -3.43694663e+00,  1.43834212e+00,
       -1.12079105e-02, -2.34521849e+00, -1.86480066e+00, -5.03962756e+00,
        8.55806378e-01, -2.17675292e+00, -5.06764094e+00, -1.99464442e+00,
       -7.77785827e+00, -2.12764093e+00, -9.09541075e-01, -5.04235077e+00,
       -2.43800740e+00, -3.84034056e+00, -2.53403530e+00, -2.36178311e+00,
       -1.63103720e+00, -1.53182046e+00, -3.23911560e+00, -2.88000616e+00,
       -1.88380053e+00, -3.21582220e+00, -3.40245202e+00, -9.01813905e-01,
       -4.82380040e+00, -7.69116343e-01, -7.12894300e+00, -1.05440427e+01,
       -1.95444152e+00, -3.50194744e+00, -7.18167736e+00, -6.01355667e+00,
       -2.08180006e+00, -2.31652280e+00, -3.44556948e+00, -2.05480142e+00,
       -6.01555673e+00, -2.80302999e+00, -4.84067198e+00, -3.51806495e-01,
       -3.54726250e+00, -2.21053919e+00, -4.36058559e+00, -4.21815003e-01,
       -1.63103720e+00, -2.51604201e+00, -2.16707077e+00, -1.48726025e+00,
       -3.20064450e+00, -1.51411241e+00, -1.65013691e+00, -3.66207663e+00,
       -3.26068347e+00, -3.04492313e+00, -4.22580088e+00, -1.68794650e+00,
       -5.98734972e+00, -1.43710934e+00, -1.97033920e+00, -1.85076729e+00,
       -1.14179382e+00, -3.07730028e+00, -2.84867198e+00, -2.19679345e+00,
       -3.68737438e+00, -2.28398218e+00,  1.89408269e+00, -3.61322115e+00,
       -2.79173430e+00, -2.41564138e+00, -7.53910853e-01, -8.54748860e-01,
       -9.20407401e-02, -6.14209981e+00, -3.79386016e+00, -7.77785827e+00,
       -1.79074130e+00, -2.50544035e+00, -3.77102985e+00, -2.25250706e+00,
       -2.57788713e+00, -2.06375990e+00, -3.33843958e+00, -1.03912484e+00,
       -6.68046164e+00, -1.91112045e+00, -2.58735858e+00, -2.19902800e+00,
       -1.90219551e+00, -2.81396751e+00, -4.10972153e+00, -5.72200956e+00,
       -1.60816482e+00, -3.68121117e+00, -4.00929775e+00, -2.45888680e+00,
       -1.13185484e+00, -1.69279652e+00, -7.09025955e+00, -3.79300016e+00,
       -2.99712050e+00, -5.70600137e+00, -2.44845780e+00, -5.20390242e+00,
       -5.29380090e-01, -3.53651118e+00, -3.51200413e+00, -2.02419300e+00,
       -4.47400933e+00, -3.63836536e+00, -4.50596939e+00, -5.57842703e+00,
       -5.30076136e+00, -2.39225449e+00, -4.88290164e+00, -2.61359300e+00,
       -3.11841945e+00, -2.05580278e+00, -3.15002845e+00, -4.91881901e+00,
       -3.03774254e+00, -4.26411548e+00, -3.15002845e+00, -3.49352203e+00,
       -3.81708831e+00, -3.77197358e+00, -2.55016650e+00, -1.97653920e+00,
       -2.59432621e+00, -5.38480406e+00, -5.44932525e+00, -3.04107137e+00,
       -1.87252400e+00, -2.25124657e+00, -2.09215707e+00, -3.40087334e+00,
       -6.06061986e+00, -1.89916369e+00, -1.90035105e+00, -2.45036038e+00,
       -2.79303037e+00, -4.76810415e+00, -1.72379306e+00, -7.00025955e+00,
       -2.86880158e+00, -2.70674744e+00, -4.36825704e+00, -3.11841945e+00])
```

✓ 0s  completed at 7:43PM

---

+ Code  + Text   All changes saved

```
[14] #evaluate model performance
     from sklearn.metrics import mean_squared_error,r2_score
     lr_train_mse=mean_squared_error(y_train,y_lr_train_pred)
     lr_train_r2=r2_score(y_train,y_lr_train_pred)

     lr_test_mse=mean_squared_error(y_test,y_lr_test_pred)
     lr_test_r2=r2_score(y_test,y_lr_test_pred)
```

```
[31] print('LR MSE (train):',lr_train_mse)
     print('LR R@ (train):',lr_train_r2)
     print('LR MSE (test):',lr_test_mse)
     print('LR R2 (test):',lr_test_r2)

     LR MSE (train): 1.0075362951003687
     LR R@ (train): 0.7645051774663391
     LR MSE (test): 1.0206953660061033
     LR R2 (test): 0.7891616188563282
```

```
lr_train_mse

1.0075362951003687
```

```
[17] lr_train_r2

0.7645051774663391
```

```
[18] lr_test_mse

1.0206953660061033
```

```
[19] lr_test_r2

0.7891616188563282
```

```
[20] lr_results=pd.DataFrame(['Linear regression',lr_train_mse,lr_train_r2,lr_test_mse,lr_test_r2]).transpose()
     lr_results.columns=['Method','Training MSE','Training R2','Test MSE','Test R2']
     lr_results
```

| | Method | Training MSE | Training R2 | Test MSE | Test R2 |
|---|---|---|---|---|---|
| 0 | Linear regression | 1.007536 | 0.764505 | 1.020695 | 0.789162 |

✓ 0s  completed at 7:43PM

+ Code + Text   All changes saved

```
[19] lr_test_r2
```

```
0.789163618563282
```

```
[20] lr_results=pd.DataFrame(['Linear regression',lr_train_mse,lr_train_r2,lr_test_mse,lr_test_r2]).transpose()
     lr_results.columns=['Method','Training MSE','Training R2','Test MSE','Test R2']
     lr_results
```

| | Method | Training MSE | Training R2 | Test MSE | Test R2 |
|---|---|---|---|---|---|
| 0 | Linear regression | 1.007536 | 0.764505 | 1.020695 | 0.789162 |

```
##random forest
#training the model
from sklearn.ensemble import RandomForestRegressor
rf=RandomForestRegressor(max_depth=2,random_state=100)
rf.fit(x_train,y_train)
```

```
          RandomForestRegressor
RandomForestRegressor(max_depth=2, random_state=100)
```

```
[23] #applying the model to make prediction
     y_rf_train_pred=rf.predict(x_train)
     y_rf_test_pred=rf.predict(x_test)
```

```
[25] #evaluate model performance
     from sklearn.metrics import mean_squared_error,r2_score
     rf_train_mse=mean_squared_error(y_train,y_rf_train_pred)
     rf_train_r2=r2_score(y_train,y_rf_train_pred)

     rf_test_mse=mean_squared_error(y_test,y_rf_test_pred)
     rf_test_r2=r2_score(y_test,y_rf_test_pred)
```

```
[26] rf_results=pd.DataFrame(['Random Forest',rf_train_mse,rf_train_r2,rf_test_mse,rf_test_r2]).transpose()
     rf_results.columns=['Method','Training MSE','Training R2','Test MSE','Test R2']
     rf_results
```

| | Method | Training MSE | Training R2 | Test MSE | Test R2 |
|---|---|---|---|---|---|
| 0 | Random Forest | 1.028228 | 0.759669 | 1.407688 | 0.709223 |

---

+ Code + Text   All changes saved

```
[25] from sklearn.metrics import mean_squared_error,r2_score
     rf_train_mse=mean_squared_error(y_train,y_rf_train_pred)
     rf_train_r2=r2_score(y_train,y_rf_train_pred)

     rf_test_mse=mean_squared_error(y_test,y_rf_test_pred)
     rf_test_r2=r2_score(y_test,y_rf_test_pred)
```

```
[26] rf_results=pd.DataFrame(['Random Forest',rf_train_mse,rf_train_r2,rf_test_mse,rf_test_r2]).transpose()
     rf_results.columns=['Method','Training MSE','Training R2','Test MSE','Test R2']
     rf_results
```

| | Method | Training MSE | Training R2 | Test MSE | Test R2 |
|---|---|---|---|---|---|
| 0 | Random Forest | 1.028228 | 0.759669 | 1.407688 | 0.709223 |

```
[27] #model comparison
     df_models=pd.concat([lr_results,rf_results],axis=0)
     df_models
```

| | Method | Training MSE | Training R2 | Test MSE | Test R2 |
|---|---|---|---|---|---|
| 0 | Linear regression | 1.007536 | 0.764505 | 1.020695 | 0.789162 |
| 0 | Random Forest | 1.028228 | 0.759669 | 1.407688 | 0.709223 |

```
df_models.reset_index(drop=True)
```

| | Method | Training MSE | Training R2 | Test MSE | Test R2 |
|---|---|---|---|---|---|
| 0 | Linear regression | 1.007536 | 0.764505 | 1.020695 | 0.789162 |
| 1 | Random Forest | 1.028228 | 0.759669 | 1.407688 | 0.709223 |

```
[29] #data visulization of prediction results
```

```
[30] import matplotlib.pyplot as plt
     import numpy as np
     plt.figure(figsize=(5,5))
     plt.scatter(x=y_train,y=y_lr_train_pred,c='#7CAE00',alpha=0.3)

     z=np.polyfit(y_train,y_lr_train_pred,1)
     p=np.poly1d(z)

     plt.plot(y_train,p(y_train),'#F8766D')
```

Code->

# Load data

import pandas as pd

url = 'https://raw.githubusercontent.com/dataprofessor/data/master/delaney_solubility_with_descriptors.csv'

df = pd.read_csv(url, error_bad_lines=False)

df

#data preparation

#data seperation in x and y

y=df['logS']

```python
y

x=df.drop('logS',axis=1)

x
```

```python
#data splitting
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=100)
x_train
x_test
```

```python
#model building
#linear regression
from sklearn.linear_model import LinearRegression
lr=LinearRegression()
lr.fit(x_train,y_train)
#applying the model to make a prediction
y_lr_train_pred=lr.predict(x_train)
y_lr_test_pred=lr.predict(x_test)
y_lr_train_pred
y_lr_test_pred
```

```python
#evaluate model performance

from sklearn.metrics import mean_squared_error,r2_score

lr_train_mse=mean_squared_error(y_train,y_lr_train_pred)

lr_train_r2=r2_score(y_train,y_lr_train_pred)


lr_test_mse=mean_squared_error(y_test,y_lr_test_pred)

lr_test_r2=r2_score(y_test,y_lr_test_pred)

print('LR MSE (train):',lr_train_mse)

print('LR R@ (train):',lr_train_r2)

print('LR MSE (test):',lr_test_mse)

print('LR R2 (test):',lr_test_r2)

lr_train_mse

lr_train_r2

lr_test_mse

lr_test_r2

lr_results=pd.DataFrame(('Linear
regression',lr_train_mse,lr_train_r2,lr_test_mse,lr_test_r2)).transpose()

lr_results.columns=['Method','Training MSE','Training R2','Test MSE','Test R2']

lr_results
```

```python
##random forest

#training the model

from sklearn.ensemble import RandomForestRegressor

rf=RandomForestRegressor(max_depth=2,random_state=100)

rf.fit(x_train,y_train)



#applying the model to make prediction

y_rf_train_pred=rf.predict(x_train)

y_rf_test_pred=rf.predict(x_test)

#evaluate model performance

from sklearn.metrics import mean_squared_error,r2_score

rf_train_mse=mean_squared_error(y_train,y_rf_train_pred)

rf_train_r2=r2_score(y_train,y_rf_train_pred)


rf_test_mse=mean_squared_error(y_test,y_rf_test_pred)

rf_test_r2=r2_score(y_test,y_rf_test_pred)

#model comparison

df_models=pd.concat((lr_results,rf_results),axis=0)

df_models

df_models.reset_index()
```

```python
#data visulization of prediction results

import matplotlib.pyplot as plt

import numpy as np

plt.figure(figsize=(5,5))

plt.scatter(x=y_train,y=y_lr_train_pred,c='#7CAE00',alpha=0.3)

z=np.polyfit(y_train,y_lr_train_pred,1)

p=np.poly1d(z)

plt.plot(y_train,p(y_train),'#FB766D')

plt.ylabel('Predict LogS')

plt.xlabel('Experimental LogS')
```
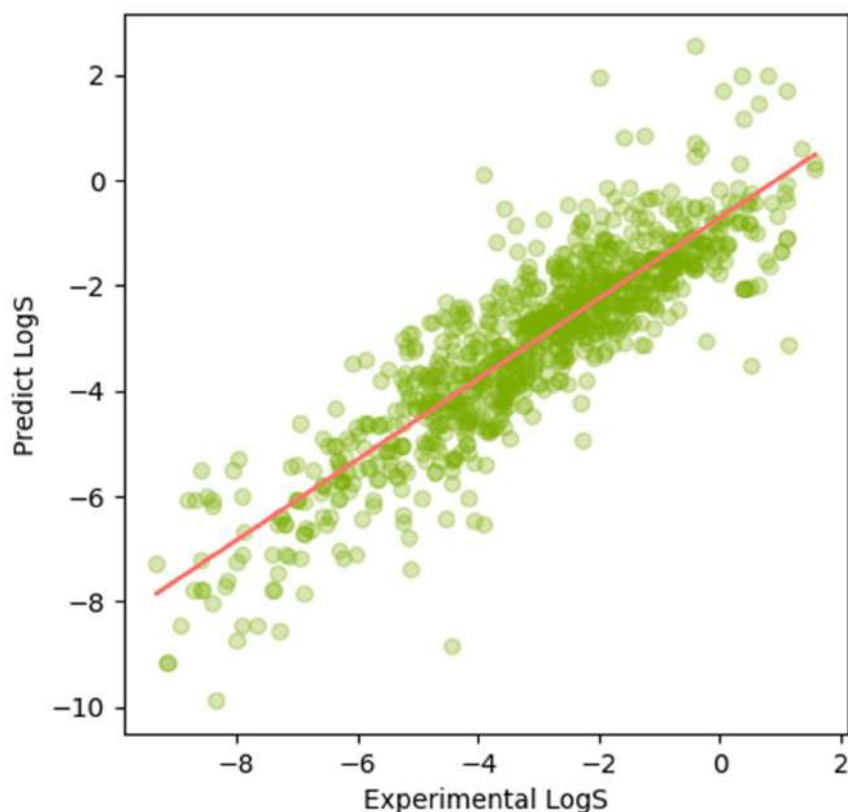
Output->

Terminologies-

The Delaney solubility dataset is a well-known dataset in the field of quantitative structure-property relationship (QSPR) modeling and cheminformatics. It is often used to develop predictive models for estimating the aqueous solubility of chemical compounds based on their molecular descriptors.

1. **Chemical Compounds**: These are individual molecules or compounds for which the solubility data is collected. Each compound has a unique structure that can be represented by its molecular formula, which specifies the types and counts of atoms in the molecule.

2. **Molecular Descriptors:** Molecular descriptors are numerical representations of the structural and physicochemical properties of chemical compounds. They provide information about a compound's size, shape, bonding patterns, and other characteristics. Examples of molecular descriptors include:

   - Molecular weight: The sum of the atomic weights of all atoms in the molecule.

   - Number of atoms: Count of different types of atoms (e.g., carbon, hydrogen, oxygen, etc.).

   - Number of bonds: Count of different types of chemical bonds (e.g., single, double, triple bonds).

   - LogP (Partition Coefficient): A measure of a compound's lipophilicity (ability to dissolve in lipids) and hydrophilicity (ability to dissolve in water).

   - Polar Surface Area (PSA): A measure of a compound's surface area that is capable of hydrogen bonding.

   - etc.

3. **Aqueous Solubility:** This is the experimental measurement of how much of a given compound can dissolve in water. It is typically reported as a concentration value (e.g., in mg/L or mol/L) at which the compound reaches its saturation point in water under specific conditions.

## Working Model:

To develop a predictive model for estimating aqueous solubility based on molecular descriptors, a machine learning approach can be used. Here's a general outline of how the working model could be constructed:

1.  **Data Preprocessing:**

    *   Load the Delaney solubility dataset, which includes compound structures, molecular descriptors, and experimental solubility values.

    *   Perform any necessary data cleaning, such as handling missing values or outliers.

2.  **Feature Selection/Engineering:**

    *   Select relevant molecular descriptors that are likely to influence solubility.

    *   Optionally, create new features by combining or transforming existing descriptors if it enhances model performance.

3.  **Splitting the Data:**

    *   Divide the dataset into training and testing subsets. The training set is used to build the model, while the testing set is used to evaluate its performance.

4.  **Model Selection:**

    *   Choose a suitable machine learning algorithm for regression, as solubility is a continuous property prediction task. Common algorithms include linear regression, random forests, support vector machines, and neural networks.

5.  **Model Training:**

- Train the chosen model using the training data and their corresponding solubility values.

6. **Model Evaluation:**

    - Evaluate the model's performance on the testing data using appropriate metrics such as mean squared error (MSE), root mean squared error (RMSE), or coefficient of determination (R-squared).

7. **Prediction and Interpretation:**

    - Use the trained model to predict solubility for new compounds based on their molecular descriptors.

    - Interpret the model's predictions and understand which descriptors have the most significant impact on solubility predictions.

8. **Model Refinement (Optional):**

    - Fine-tune the model parameters or try different algorithms to improve performance.