

Towards Interpretable Machine Learning Models for Materials Discovery^[1]

Shreecharan Sundar (CH17B075)
Dept of Chemical engineering, IIT Madras

1. MOTIVATION

Quantitative structure activity relationship (QSAR) modeling is now more than 50 years old, and its utility has been shown in numerous research publications and by the number of new drug and agrochemical entities that have been developed with its aid. The method has evolved very substantially since the seminal linear regression QSAR models. Many new types of molecular descriptors have been developed, new mathematical methods such as neural networks, support vector machines, kernel regression, and random forest have been applied to mapping structure to activity, and QSAR has now incorporated 3D structures, conformation, and chirality. Two main branches of QSAR have evolved. The first of these, where the model is often relatively simple and linear and interpretable in terms of molecular interactions or biological mechanisms. The second type focuses much more on modeling structure–activity relationships in large data sets with high chemical diversity using a variety of regression or classification methods, and often the interpretation of the model is obscure or impossible. A recent publication ^[2] provides detailed discussion on this topic. Data-driven ML models are considerably more useful if more chemically interpretable descriptors are used to train them, as long as these models also accurately recapitulate the properties of materials. Current modelling methods use arcane mathematical descriptors which meant that model predictive power and interpretability remained essentially orthogonal. This paper introduces a particular type of molecular fragment descriptor, the signature descriptor which achieves these joint aims of accuracy and interpretability.

2. INTRODUCTION

2.1 Quantitative structure property relationship (QSPR/QSAR) modeling:

QSPR/QSAR model represent mathematical relationship between the response of chemical compound(activity/property) with their structural and physiochemical information. It works on the principle that structurally similar compounds have similar properties. A QSPR has the form of a mathematical model :

$$\text{Property} = f(D_1, D_2, \dots, D_4) + \text{error}$$

D₁, D₂, ..., D₄ – chemical or structural property (Molecular Descriptors)

A **QSPR model** is composed of dependent and independent variables. In the **QSPR model**, the dependent variable is estimated from the independent variables. Independent variables are **molecular descriptors** computed from the chemical structure of compounds, and the dependent variable is the property under study.

¹ Paulius Mikulskis, Morgan R. Alexander, and David Alan Winkler, ‘Toward Interpretable Machine Learning Models for Materials Discovery’, *Advanced Intelligent Systems*, 1.8 (2019) <<https://doi.org/10.1002/aisy.201900045>>.

² Toshio Fujita and David A. Winkler, ‘Understanding the Roles of the “Two QSARs”’, *Journal of Chemical Information and Modeling*, 56.2 (2016), 269–74 <<https://doi.org/10.1021/acs.jcim.5b00229>>.

CH5650 : Molecular Data Science and Informatics Project report

2.2 Molecular Descriptors

Molecular descriptors are the result of a logical and mathematical procedure which transforms chemical information encoded within a symbolic representation of a molecule into a number. Molecular descriptors are divided into two main categories: **experimental measurements**, such as log P, molar refractivity, dipole moment, polarizability, and **theoretical molecular descriptors**, which are derived from a symbolic representation of the molecule.

The main classes of theoretical molecular descriptors are: 1) **0D-descriptors** (i.e. constitutional descriptors, count descriptors), 2) **1D-descriptors** (i.e. list of structural fragments, fingerprints), 3) **2D-descriptors** (invariants of molecular graphs, e.g., connectivity indices, information indices, counts of paths and walk), 4) **3D-descriptors** (based on spatial properties of the molecules)

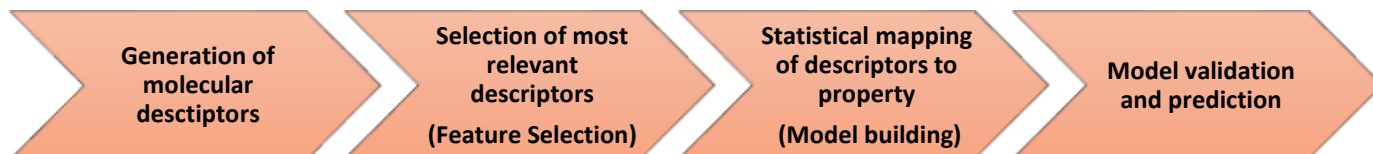
Types of Descriptor		
Descriptor type	Description	Examples
Constitutional Descriptors	They represent properties related to molecular structure	molecular weight, total number of atoms in the molecule, number of aromatic rings
Electrostatic	They represent properties related to the electronic nature of the compound	atomic net partial and charges
Topological Descriptors	They represent properties which can be inferred by treating the structure of the compound as a graph, with atoms as vertices and covalent bonds as edges	total number of bonds in shortest paths between all pairs of non-hydrogen atoms
Geometrical Descriptors	They represent properties related to spatial arrangement of atoms constituting the compound	Vander Waals Area
Fragment based Descriptors	They represent properties related to substructural motifs	MDL Keys and Molecular Fingerprints

Tools to calculate Molecular descriptors		
Name	Type	Availability
Chem Development Kit	Continuous	Free. https://cdk.github.io/
Padel	Continuous /Fingerprint	Free. http://www.yapcwsoft.com/dd/padeldescriptor
RDKit	Continuous /Fingerprint	Free. http://www.rdkit.org
Dragon	Continuous	Commercial. http://www.talete.mi.it/products/dragon_description.htm

2.3 QSPR Workflow

1. calculation of molecular descriptors (mathematical representations of molecular properties)
2. descriptor selection (i.e. choosing the subset of descriptors most relevant to the problem)
3. finding an optimum (usually nonlinear) relationship between the descriptors and response variable
4. validation of the model (i.e. how predictive/ robust is it, and domain of applicability)

CH5650 : Molecular Data Science and Informatics Project report



Feature Selection
Importance of feature selection : <ul style="list-style-type: none">Improves Interpretation<ul style="list-style-type: none">Less features, simpler models.Expert-driven feature selection enhances the mechanistic interpretation of the models.Reduces Overfitting<ul style="list-style-type: none">Less redundant data means lesser decisions based on noise.Reduces Training Time<ul style="list-style-type: none">Less data to learn from ensures quicker model development.
Common feature selection methods: <ul style="list-style-type: none">Filter: e.g.: Chi-squared test, T-test, Correlation coefficientWrapper: e.g.: Genetic search, Recursive feature eliminationEmbedded: e.g. : Decision trees ,Support vector machines, Lasso regularization
Model development
<ul style="list-style-type: none">Multiple linear regression, k-Nearest Neighbors (k-NN), linear discriminant analysis, decision trees, random forests, support vector machines, artificial neural networks and other ensemble methods are some of the most used methods for QSPR modelling
Model validation
<ul style="list-style-type: none">Internal validation<ul style="list-style-type: none">K-fold cross validation: The dataset is split into K parts. K models are developed using (K-1) sets and the Kth set is used as the test set.Leave one out cross-validation: N models are developed each with (N – 1) chemicals as training set and 1 chemical as the test set.External test set validation<ul style="list-style-type: none">Regression metrics: Root-mean-squared-error, Mean Average Error, Coefficient of DeterminationClassification metrics: Accuracy, Sensitivity, Positive Predictivity, Recall, Precision, F1- score

2.4 Signature and Dragon Descriptors

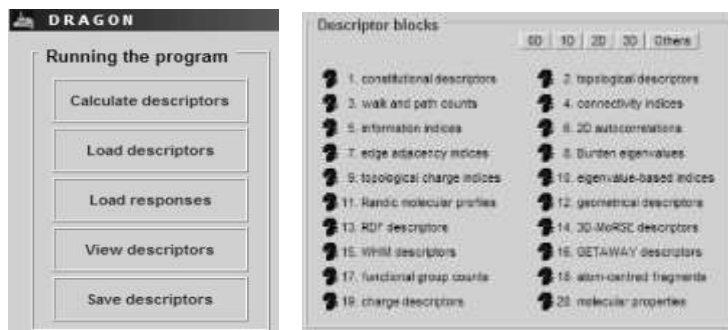
- (1) **Dragon :** Dragon is an application for the calculation of molecular **descriptors**. The Dragon program can calculate 5000 molecular descriptors, some of which are easy to interpret, but the **majority are quite arcane and almost impossible to understand in terms of chemical structure**.

DRAGON was designed as a user-friendly software which performs descriptor calculation using a simple logical sequence (Figure 1) :

1. loading of the molecular files
2. selection of the descriptors
3. calculation of the descriptors
4. saving of the calculated descriptors

CH5650 : Molecular Data Science and Informatics Project report

Together with the molecular descriptors, DRAGON also allows the calculation of the Principal Components (PCs) of descriptors included in the selected block to compress into a small number of variables.



(2) Signature descriptors

The signature of an atom in a polymer is a representation of the atom's connected environment up to a predefined depth d . Signature descriptor is generated from the chemical graph of a molecule or material, and describes the number of atoms, atom types and hybridizations, bond types. Signature descriptors are more chemically interpretable than many other types of descriptors as they have been used successfully in a range of drug discovery applications due to their inherent **ability to be mapped back to chemical structures** to provide clear guidance to which functionalities in molecules contribute to, desired biological responses. A brief description of how signature descriptors are generated is described in figure 1. The algorithms for generating signature descriptors is included in these publications^{[3][4]}.

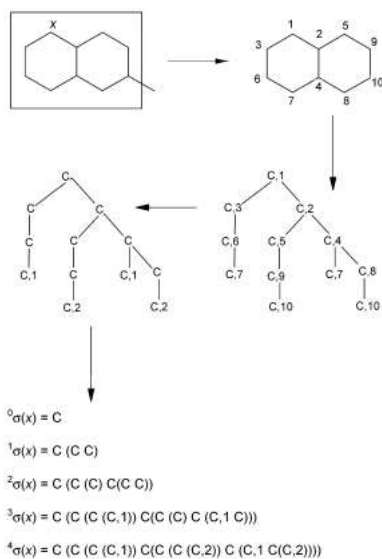


Figure 1. Atomic signature. The figure illustrates the five steps procedure for computing the atomic signature of atom x in methyl-naphthene. (1) The subgraph containing all atoms at distance 4 from atom x is extracted. (2) The subgraph is canonized atom x having label 1. (3) A tree spanning all edges of the subgraph is constructed. (4) All labels appearing only once are removed, and the remaining labels are renumbered in the order they appear. (5) The signature is printed reading the tree in a depth-first order. Here we show atomic signatures from $h = 0$ to $h = 4$.

³ Pablo Carbonell, Lars Carlsson, and Jean Loup Faulon, 'Stereo Signature Molecular Descriptor', *Journal of Chemical Information and Modeling*, 53.4 (2013), 887–97 <<https://doi.org/10.1021/ci300584r>>.

⁴ Jean Loup Faulon, Michael J. Collins, and Robert D. Carr, 'The Signature Molecular Descriptor. 4. Canonizing Molecules Using Extended Valence Sequences', *Journal of Chemical Information and Computer Sciences*, 44.2 (2004), 427–36 <<https://doi.org/10.1021/ci0341823>>.

2.5 Machine Learning frameworks for model building and feature selection

2.5.1 Multiple linear regression with EM algorithm for feature selection (MLREM)

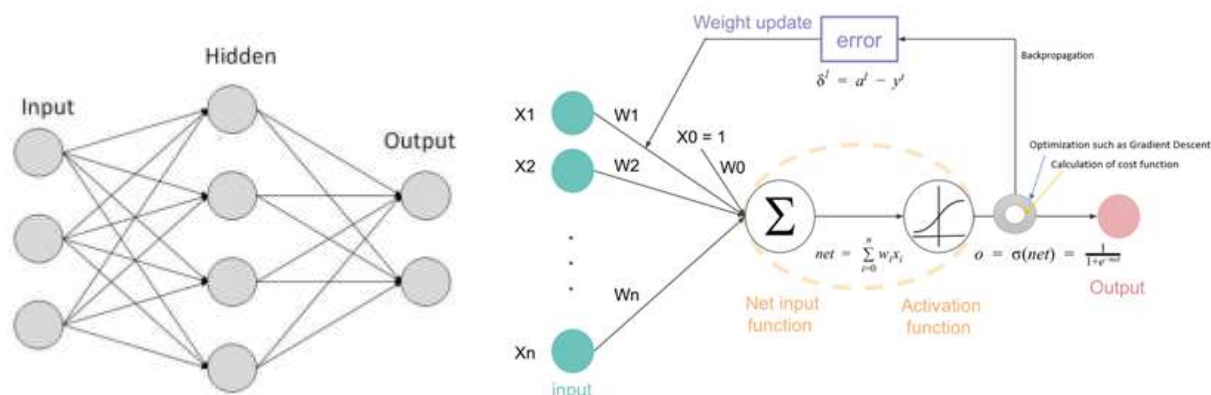
Feature selection represents the most critical step when deriving QSAR/QSPR models from a large parametric space. QSAR belongs to a family of mathematical problems, called grossly underdetermined systems, in which the number of independent variables greatly exceeds the number of dependent variables.

MLREM^[5] is a Bayesian method that uses a sparse prior to eliminate uninformative or low-relevance descriptors by compressing them to zero. It is an L1 feature selection method like the least absolute shrinkage and selection operator method. Feature significance was encoded by sizes and magnitudes of the regression coefficients for each descriptor. Descriptors identified using the MLREM feature selection procedure were used to train nonlinear neural network regression models. In this project, I implemented (L1 feature selection method) Least absolute shrinkage and selection operator that selects the most relevant descriptors and shrinks other coefficients to zero. The chosen descriptors were then used to build the neural network model

2.5.2 Artificial Neural Networks (feed forward and backpropagation)

ANNs are artificial systems simulating the function of the human brain. Three components constitute a neural network: the processing elements or nodes, the topology of the connections between the nodes, and the learning rule by which new information is encoded in the network. While there are a number of different ANN models, the most frequently used type of ANN in QSAR is the three-layered feed-forward network. In this type of networks, the neurons are arranged in layers (an input layer, one hidden layer and an output layer). Each neuron in any layer is fully connected with the neurons of a succeeding layer and no connections are between neurons belonging to the same layer.

According to the supervised learning adopted, the networks are taught by giving them examples of input patterns and the corresponding target outputs. Through an iterative process, the connection weights are modified until the network gives the desired results for the training set of data. A back-propagation algorithm is used to minimize the error function.



⁵ F. R. Burden and Dave A. Winkler, 'Optimal Sparse Descriptor Selection for QSAR Using Bayesian Methods', *QSAR and Combinatorial Science*, 28.6–7 (2009), 645–53 <<https://doi.org/10.1002/qsar.200810173>>.

2.6 Feature significance of the models

Feature significance (parametric sensitivity) for ML models was evaluated by generating the partial derivative of the response variable (**a**) with respect to the descriptor at a position on the response surface close to the response variable maximum. This was achieved by inducing a small change of 0.01 to each normalized descriptor value (**d**) in turn, with the remaining descriptors kept fixed. The difference between original predicted and perturbed predicted property value obtained from neural network (NN) allowed the partial derivative to be calculated

$$\text{partial derivate (delta)} = \frac{\Delta a}{\Delta d} = \frac{\text{change in output}}{\text{change in feature}}$$

3. Implementation

This section explains the different implementation steps followed to reproduce the results obtained in the paper. The paper demonstrates the ability of signature descriptors to generate predictive and interpretable QSPR models by modelling seven diverse materials data sets. I have chosen and modelled “**aqueous solubility**” property and compared the results obtained with those in the paper. The dataset for this is associated with the publication ^[6] and available from the <http://pubs.acs.org>.

3.1 Dataset:

Summary of dataset				
Property	Dataset size	Train-Test split	No of signature descriptors generated	Modelled value
AQSOL	1143	(80-20 %)	3378	Log (octanol-water partition coefficient)

3.2 RDKit Package

RDKit is an opensource cheminformatics toolkit written in C++ and Python. In this project, for converting the SMILES structure to .mol format the RDKit package is utilized. RDKit package can be installed using Anaconda prompt using the following command:

```
conda install -c rdkit rdkit
```

```
# the following code will creat mol files for the 1143 compounds in the data set from their smiles representation
from rdkit import Chem
df=pd.read_csv("aqsol_data.txt",header=None)
# the format of the csv file is (Chemical name of the molecule,Aqsol true, predicted aqsol, smiles representation)
name=df[0]
smiles=df[3]
for i in range(len(smiles)):
    molObj = Chem.MolFromSmiles(smiles[i])
    imported = Chem.MolToMolBlock(molObj)
    with open(str(name[i])+".mol") as newfile:
        newfile.write(imported)
```

⁶ John S. Delaney, ‘ESOL: Estimating Aqueous Solubility Directly from Molecular Structure’, *Journal of Chemical Information and Computer Sciences*, 44.3 (2004), 1000–1005 <<https://doi.org/10.1021/ci034243x>>.

CH5650 : Molecular Data Science and Informatics Project report

Example Mol file generated for 1,1,1,2- Tetrachloroethane:

```
1,1,1,2-Tetrachloroethane - Notepad
File Edit Format View Help

RDKit      2D

6  5  0  0  0  0  0  0  0  0  0999 V2000
-1.2990   0.7500   0.0000 C1  0  0  0  0  0  0  0  0  0  0  0  0
 0.0000   0.0000   0.0000 C  0  0  0  0  0  0  0  0  0  0  0  0
 1.2990   0.7500   0.0000 C  0  0  0  0  0  0  0  0  0  0  0  0
 2.5981   1.5000   0.0000 C1  0  0  0  0  0  0  0  0  0  0  0  0
 2.0490  -0.5490   0.0000 C1  0  0  0  0  0  0  0  0  0  0  0  0
 0.5490   2.0490   0.0000 C1  0  0  0  0  0  0  0  0  0  0  0  0
1  2  1  0
2  3  1  0
3  4  1  0
3  5  1  0
3  6  1  0
M  END
```

3.3 SSCAN package (Linux Executable)

sscan is an open-source package that has the algorithm to generate signature descriptor of compounds from their mol file. It is in Linux executable format.

1. Download the package from <https://sourceforge.net/projects/molSIG/>
2. An executable sscan file will be generated on running make file.

INSTALLATION

1. Go to src folder and run "make".
2. Executable "sscan" will be installed in bin folder.

```
srijoo@LAPTOP-I3EGEPHT:~$ cd /mnt/c/Users/sreen/Desktop/molSIG-code/src
srijoo@LAPTOP-I3EGEPHT:~/molSIG-code/src$ make -f Makefile
/usr/bin/gcc -o ../bin/sscan -m32 \
    allocation_da.o random_da.o structure_da.o structure_nb.o structure_tp.o structure_ut.o twin_da.o hea.o structur
e_he.o allocation_si.o signature.o cansig.o candag.o utility_si.o signature_mo.o in_mo.o out_mo.o utility_mo.o \
    DFConst.o DFManageMemory.o DFFindGeometry.o PCStereoSig.o PCcip.o \
    error.o vector.o main.o -lm
```

3. Signature descriptors can be generated using the sscan file for all the compounds in our dataset up to depth 4 from its mol file (The paper has generated up to depth 7 but due to limited computational capacity I generated up to depth 5).

Stereo Signature Molecular Descriptor

Usage: sscan <file.mol> <output-type>

<file> is mol file

<output-type>

output types are (x is the depth, any integer between 0 and 100000):

scanx: signature DAG

sscanx: stereo signature DAG (CIP rules)

fscanx: fast stereo signature mode

```
srijoo@LAPTOP-I3EGEPHT:~/molSIG-code/bin$ ./sscan 1,1,1,2-Tetrachloroethane.mol fscan5
```

4. Signature descriptors for all the molecules are generated and compiled

3.4 Implementation approach

3.4.1 Feature set preparation

Compiling all the descriptors and removal of duplicates; creation of [X] matrix

```
1 # Compiling all the generated signature descriptors
2 for i in range(len(filename)):
3     for j in range(5):
4         new=pd.read_csv(str(filename[i])+".fscan"+str(j),header=None,sep=" ")
5         new=new[:-1]
6         for k in range(len(new)):
7             feature.append(str(new[i][k]))
```

```
1 len(feature)
```

37546

```
1 #removal of duplicates
2 final_feature = []
3 for i in feature:
4     if i not in final_feature:
5         final_feature.append(i)
```

```
1 len(final_feature)
```

3378

```
1 # Mapping each compound to its corresponding features for model building
2 X= np.zeros((1143,3379))
3 blank_df = pd.DataFrame(X,columns=final_feature+['Chemical name'])
4 map_dict1 = {}
5 map_dict2 = {}
6 for i in range(len(filename)):
7     for j in range(5):
8         new=pd.read_csv(str(filename[i])+".fscan"+str(j),header=None,sep=" ")
9         new=new[:-1]
10        for k in range(len(new)):
11            feature.append(new[i][k])
12            if new[i][k] not in list(map_dict2.keys()):
13                map_dict2.update([(new[i][k],[new[0][k]])])
14        act_feature =list(map_dict2.keys())
15        feature_df = pd.DataFrame.from_dict(map_dict2)
16        blank_df.loc[i,act_feature] = feature_df.loc[0,act_feature]
17        blank_df.loc[i,'Chemical name'] = filename[i]
18        map_dict2={}
19 #saving the generated matrix to a csv to be exported for model building
20 blank_df.to_csv("file1.csv")
```

3.4.2 Data Matrix preparation, preprocessing and normalization

Generating X (feature set) and Y (prediction) matrix for model building .Removal of highly correlated variables (higher than 90 percent correlation)

```
1 # X and y matrix generation
2 y_vec=pd.read_csv("aqsol data.txt",header=None)
3 Y=y_vec[1] #Target matrix ()
4 Y=np.array(Y)
5 x_vec=pd.read_csv("file1.csv",index_col=0)
6 X=x_vec.values
7 X=X[:, :-1] #feature matrix
```

```
1 print("total no of generated signature descriptors",len(x_vec.columns))
2 new_vec=x_vec.iloc[:, :-1]
```

total no of generated signature descriptors 3379

CH5650 : Molecular Data Science and Informatics Project report

```
1 # function to remove features with more than 90 percent correlation
2 def correlation(dataset, threshold):
3     col_corr = set() # Set of all the names of deleted columns
4     corr_matrix = dataset.corr()
5     for i in range(len(corr_matrix.columns)):
6         for j in range(i):
7             if (corr_matrix.iloc[i, j] >= threshold) and (corr_matrix.columns[j] not in col_corr):
8                 colname = corr_matrix.columns[i] # getting the name of column
9                 col_corr.add(colname)
10                if colname in dataset.columns:
11                    del dataset[colname] # deleting the column from the dataset
12
13    return(dataset)
14 uncor=correlation(new_vec,0.9)
15 X=uncor.values
16 X=X.astype(int)

1 print("total no of signature descriptors after removing feature with more than 90% correlation",len(uncor.columns))
total no of signature descriptors after removing feature with more than 90% correlation 1571
```

Data Normalisation

```
1 Scaler= MinMaxScaler()
2 X_scaled=Scaler.fit_transform(X)
```

3.4.3 L1 feature selection (hyperparameter tuning for the best model)

```
1 param_grid = {'alpha': uniform()}
2 model = Lasso()
3 rand_search = RandomizedSearchCV(estimator=model,
4                                   param_distributions=param_grid,random_state=40,
5                                   n_iter=200)
6
7 rand_search.fit(X_scaled, Y)
8
9 print(rand_search.best_estimator_.alpha)
10 print(rand_search.best_score_)
11 best_alpha=rand_search.best_estimator_.alpha

0.0042243686391818525
0.8037845591191454
```

3.4.4 MLR model with best hyperparameters chosen from gridsearch

```
1 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.20, random_state=40)
2 scaling= MinMaxScaler()
3 X_train=scaling.fit_transform(X_train)
4 X_test=scaling.transform(X_test)

1 model_lasso = lasso(alpha=best_alpha)
2 model_lasso.fit(X_train,y_train)
3 y_pred_train=model_lasso.predict(X_train)
4 y_pred_test=model_lasso.predict(X_test)
5 r2train= r2_score(y_train,y_pred_train)
6 r2test=r2_score(y_test,y_pred_test)

1 r2train,r2test

(0.8777022860095753, 0.8276239373798526)
```

CH5650 : Molecular Data Science and Informatics Project report

3.4.5 Features chosen from L1 feature selection technique and feature significance

```
1 coeff = model_lasso.coef_  
2 a=[]  
3 for i in range(len(coeff)):  
4     if(coeff[i]!=0):  
5         a.append(i)  
6 XNN=X[:,a]  
7 print(" No of Signature descriptors selected using L1 feature selection method:",len(XNN[0]))
```

No of Signature descriptors selected using L1 feature selection method: 83

```
1 dict={'weights':coeff[a].astype(float),'descriptor':uncor.columns[a]}  
2 new_df=pd.DataFrame(data=dict)  
3
```

```
1 new_df.sort_values(by= [new_df.columns[0]],axis=0, ascending=True, inplace=True)  
2 print(new_df.head(12))  
3 print("\n")  
4 print(new_df.tail(12))
```

3.4.6 Regularized Neural Network model building with chosen features

Baseline model and performance:

```
1 X_train, X_test, y_train, y_test = train_test_split(XNN, Y, test_size=0.20, random_state=40)  
2 scaling= MinMaxScaler()  
3 X_train=scaling.fit_transform(X_train)  
4 X_test=scaling.transform(X_test)
```

```
1 #baseline model  
2 baseline_model = Sequential([  
3     Dense(len(X_train[0]), activation='relu'),  
4     Dense(8, activation='relu'),  
5     Dense(1,activation='linear')])  
6 # Compile the model.  
7 baseline_model.compile(optimizer='adam',loss='mean_squared_error')  
8 # Train the model.  
9 history = baseline_model.fit(  
10     X_train,  
11     y_train,  
12     epochs=100,  
13     validation_data = (X_test,y_test)  
14 )
```

Hyperparameter tuning of the neural network model to improve model performance

```
1 #defining callback to save the best model during training  
2 def model_callbacks(weight_file='Weights_v3(epoch:02d).h5'):  
3     checkpoint = ModelCheckpoint(weight_file, monitor = 'val_loss', verbose=1, save_best_only= True,save_weights_only= True,  
4     #checkpoint = ModelCheckpoint(weight_file, monitor = 'val_loss', verbose=1)  
5     return [checkpoint]  
6 #defining function to grid search best hyperparameter for learning rate and batch_size  
7 def my_nn_model(alpha=0.001):  
8     model = Sequential([  
9         Dense(len(X_train[0]), activation='relu',activity_regularizer = regularizers.l2(1e-3)),  
10        Dropout(0.2),  
11        Dense(8, activation='relu',activity_regularizer = regularizers.l2(1e-4)),  
12        Dropout(0.01),  
13        Dense(1,activation='linear')])  
14    model.compile(optimizer=Adam(learning_rate=alpha),loss='mean_squared_error')  
15    return model
```

CH5650 : Molecular Data Science and Informatics Project report

```
1 #searching for the best hyperparameter and saving the weights of the best model
2 step_size_list = [0.0001,0.001,1e-1,1e-2]
3 batch_size_list = [16,32,64,128]
4 param_grid = {'step_size': step_size_list,'batch_size': batch_size_list}
5 grid_search = list(ParameterGrid(param_grid))
6 val_loss = []
7 for index,item in enumerate(grid_search):
8     test_model = my_nn_model(item['step_size'])
9     weight_file='Weights_gs'+str(index)+'(epoch:020).h5'
10    history = test_model.fit(X_train,y_train,epochs=500,batch_size=item['batch_size'],validation_data = (X_test,y_test),call
11    val_loss.append(min(history.history['val_loss']))
12
```

Performance of the tuned neural network model:

```
1 Y_pred=test_model.predict(X_test)
2 from sklearn.metrics import r2_score
3 print('R squared (test):',r2_score(y_test,Y_pred))
4 from sklearn.metrics import mean_squared_error,mean_absolute_error
5 print('Root mean squared error (test):',np.sqrt(mean_squared_error(y_test,Y_pred)))
6 print('Mean absolute error (test):',mean_absolute_error(y_test,Y_pred))
```

R squared (test): 0.8686955540327765
Root mean squared error (test): 0.7180816030369687
Mean absolute error (test): 0.5187011401268072

```
1 Y_pred2=test_model.predict(X_train)
2 from sklearn.metrics import r2_score
3 print('R squared (train):',r2_score(y_train,Y_pred2))
4 from sklearn.metrics import mean_squared_error,mean_absolute_error
5 print('Root mean squared error (train):',np.sqrt(mean_squared_error(y_train,Y_pred2)))
6 print('Mean absolute error (train):',mean_absolute_error(y_train,Y_pred2))
```

R squared (train): 0.9743088772087087
Root mean squared error (train): 0.3403829653076838
Mean absolute error (train): 0.24099120221265816

4. Results obtained from implementation vs results obtained in paper

4.1 Multiple Linear regression model

Performance of Multiple Linear regression model from my implementation					
Property	No of Sig descriptors	R ² test	R ² train	RMSE test	RMSE train
AQSOL	83	0.827	0.877	0.8227	0.7426

Performance of Multiple Linear regression model in the review paper					
Property	No of Sig descriptors	R ² test	R ² train	RMSE test	RMSE train
AQSOL	258	0.88	0.91	0.72	0.62

Summary of AQSOL models in literature trained on tradition descriptors (e.g. Dragon)

study	data set size	training set		test set	
		r ²	SEE	r ²	SEP
Duchowicz et al. ⁵⁴	97 train	0.87	0.903	0.85	
	48 test				0.899
	21 ext test				1.202
Wang et al. ⁸	3664	0.83	0.840		
	LOO	0.84		0.84	0.823
	85/15 crossval.			0.83	0.841
Delaney ¹⁵	2874	0.72	0.97		
Votano et al. ²³	3343/772	0.88	0.74	0.77	0.91
	1674/166	0.88	0.61	0.84	0.75
Johnson et al. ²⁸	581	0.88	0.61	0.93	0.50
this work	4376	0.90 ± 0.00	0.645 ± 0.001	0.90 ± 0.00	0.665 ± 0.003

Clearly, the AQSOL model derived from signature descriptors could predict the properties of the

CH5650 : Molecular Data Science and Informatics Project report

training and test set with good fidelity with r^2 values of 0.877 and 0.827 and standard errors of 0.82 and 0.74 logS for the training and test sets, respectively. This is comparable with many comprehensive and accurate AQSO models available in the literature, trained on Dragon descriptors. There is a small difference between the results in the paper and obtained from my implementation but the accuracy of the model is still very similar. The small difference arises mostly due to generating signature descriptors to a lower depth (5) compared to 7 in the paper. (due to limited computational power)

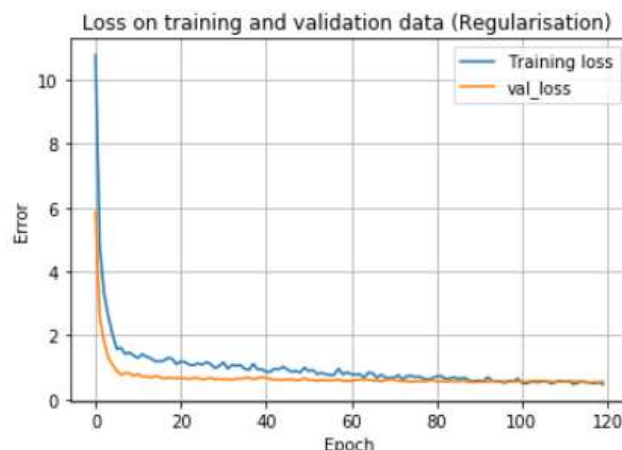
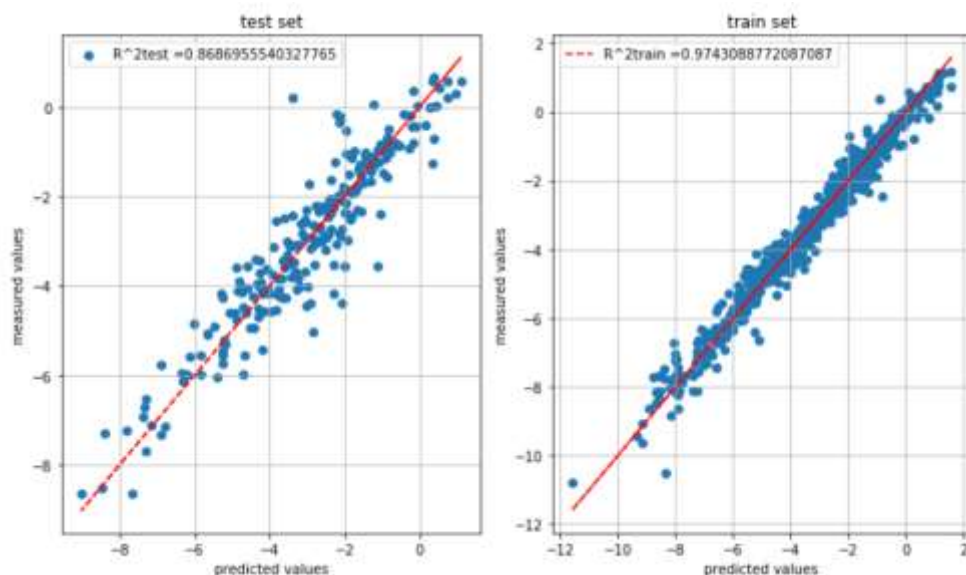
4.2 Neural Network model:

Performance of neural network model (Non linear) from my implementation

# of hidden layers	# of hidden layer neuron	R^2 test	R^2 train	RMSE test	RMSE train
1	8	0.868	0.974	0.718	0.34

Performance of neural network model (Non linear) from the paper

# of hidden layers	# of hidden layer neuron	R^2 test	R^2 train	RMSE test	RMSE train
1	2	0.91	0.97	0.65	0.39



4.3 Descriptors with most influence on Aqsol property

Most relevant features from my implementation (negative MLR coefficients reduce solubility, and positive MLR coefficients increase solubility). It matches to a good accuracy with the result in paper

weights	descriptor	Results from the paper:
-7.288989	[C]	
-5.702998	[Cl]	
-3.243925	[C]([C]([C])([C])([C]))	
-3.157147	[C]=[C]	
-2.818981	[Br]	
-2.283436	[C]([C]=[C])([C])([C]=[C])	
-1.813920	[C]([C])([C])	
-1.813215	[P]=[S]	
-1.653292	[N]([C][O]=[O])	
-1.627357	[I]	
-1.532294	[S]	
-1.293141	[C]([C]=[C])([Cl])	

weights	descriptor
0.536141	[C]=[O]
0.541980	[C]([C])([C])([C][C][O])
0.546411	[C]([C])([C])([O])
0.555033	[C](\=[C])([C])([O])
0.706211	[O]=[P]
0.733317	[C]([C])([C])([C])([C][O]))
2.160589	[O]([C])
2.564086	[O]

Signature descriptor	MLR weights	Signature descriptor	MLR weights
C=C	-0.52	O=P	+0.07
C(CC)	-0.44	C(C=CS)	+0.07
Cl	-0.39	C=O	+0.08
S	-0.21	O(C)	+0.22
C(C=C)C(C=C)	-0.15		
C(C)	-0.14		

5. Conclusion

We have been able to demonstrate that signature descriptors generate robust and predictive models for polymer properties relevant to their biological applications. ML-based models trained using traditional but hard-to-interpret molecular descriptors and molecular signature descriptors exhibited similar prediction accuracies for the structure–property relationships models studied. The signature descriptors have the advantage of more direct molecular interpretability as the most important **signatures can be easily mapped back onto exemplar materials**

6. References

- Burden, F. R., and Dave A. Winkler, ‘Optimal Sparse Descriptor Selection for QSAR Using Bayesian Methods’, *QSAR and Combinatorial Science*, 28.6–7 (2009), 645–53
<https://doi.org/10.1002/qsar.200810173>
- Carbonell, Pablo, Lars Carlsson, and Jean Loup Faulon, ‘Stereo Signature Molecular Descriptor’, *Journal of Chemical Information and Modeling*, 53.4 (2013), 887–97
<https://doi.org/10.1021/ci300584r>
- Delaney, John S., ‘ESOL: Estimating Aqueous Solubility Directly from Molecular Structure’, *Journal of Chemical Information and Computer Sciences*, 44.3 (2004), 1000–1005
<https://doi.org/10.1021/ci034243x>

CH5650 : Molecular Data Science and Informatics Project report

Faulon, Jean Loup, Michael J. Collins, and Robert D. Carr, 'The Signature Molecular Descriptor. 4. Canonizing Molecules Using Extended Valence Sequences', *Journal of Chemical Information and Computer Sciences*, 44.2 (2004), 427–36 <<https://doi.org/10.1021/ci0341823>>

Fujita, Toshio, and David A. Winkler, 'Understanding the Roles of the “Two QSARs”', *Journal of Chemical Information and Modeling*, 56.2 (2016), 269–74 <<https://doi.org/10.1021/acs.jcim.5b00229>>

Mikulskis, Paulius, Morgan R. Alexander, and David Alan Winkler, 'Toward Interpretable Machine Learning Models for Materials Discovery', *Advanced Intelligent Systems*, 1.8 (2019) <<https://doi.org/10.1002/aisy.201900045>>