

SURP 23b2444 Harsh kumar singh

PREDICTING MOLECULAR PROPERTIES

DATA PROCESSING

In this assignment we start with data processing where we read csv file for SMILE string of molecules and stored it in dataframe. Then we converted the smiles into morgan fingerprints vector using Rdkit (i used 512 bit because most molecules were smaller and we can have compact representations).

EDA(EXPLORATORY DATA ANALYSIS)

We visualized the data (shape, density of 4.11%) and examined fingerprint bit frequencies. Basically it shows **how common each fingerprint feature is across your molecules**.

Tanimoto similarity was also used in code to measure how similar two molecules are based on their fingerprint. I didn't perform exploratory data analysis.

MODEL DEVELOPMENT

Dataset is randomly split into test:train of 25:75 where x is morgan fingerprints and y is log of measured solubility in moles per litre. Then i used fingerprints smiles and performed simple random forest regression but on testing we get R² of 0.67 (which must be near 1.0 for good model) so this model didn't perform very well.

We then use descriptors on smiles of dataset to get more information like molecular weight, no. of h bonds, rotatable bonds etc. These descriptors + morgan fingerprints help in making model which is better than before.

SVM is supervised learning algorithm which i also used as baseline but R² was lesser and RSME was high. (Uses **kernels** to handle nonlinear relationships by mapping input features into higher dimensions.)

I then moved to ensemble models like random forest, xgboost and lightgbm....

—> RANDOM FOREST regressor

i used random forest model which is based on principle that many decision trees working independently and parallelly generating generalized output (bagging based). R² was higher than before now..

—> XGBOOST regressor

Gradient boosting algorithm where decision trees work sequentially and it divides at each level. This model performed well with R² of 0.88 which is nice.

—> LIGHTGBM regressor

As the name suggests the light gradient boosting algorithm uses sequential boosting of decision trees but here tree divides at a leaf node and not every level. This makes it work faster than other models (R² was also around 0.87)

We used morgan fingerprints and descriptors for them which helped in increasing accuracy overall for a model.

```
from rdkit.Chem import Draw
```

We used this package to draw 2D structure of molecules as well for top 10 molecules. Helped in visualization.

```
from torch_geometric.utils import to_networkx
```

This package was used to visualize the molecule as node and edge diagram where each atom is node and edge are bonds.

Graph neural network GNN

GNN outperforms them by learning directly from molecular graph structure, capturing spatial and relational information between atoms. Input: Molecular graphs derived from SMILES.