

# AI-Powered Prediction of Drug-Target Interaction Potency using Molecular Representations

## Abstract

The prediction of drug-target interaction (DTI) is a cornerstone of modern drug discovery.<sup>1</sup> This study explores the use of a machine learning approach to predict the potency of these interactions, quantified by the half-maximal inhibitory concentration (IC<sub>50</sub>). A dataset of IC<sub>50</sub> values was sourced from the ChEMBL database and underwent rigorous preprocessing, including filtering and conversion to a logarithmic scale (pIC<sub>50</sub>). In an initial baseline experiment, a Random Forest Regressor was trained on a set of traditional, handcrafted molecular descriptors (e.g., Molecular Weight, LogP). This model achieved a Root Mean Squared Error (RMSE) of 1.324 and an R<sup>2</sup> score of 0.254. While establishing a baseline, the limited performance highlights the constraints of using low-dimensional descriptors. We discuss the clear path forward, advocating for the use of more sophisticated Graph Neural Network (GNN) architectures, such as Graph Convolutional Networks (GCNs) and Graph Isomorphism Networks (GINs), which can learn directly from the rich topology of molecular graphs and have the potential to deliver significantly more accurate predictions.

## 1. Introduction

The process of discovering and developing new pharmaceuticals is notoriously long and expensive.<sup>2</sup> A critical bottleneck in this pipeline is the efficient identification and optimization of interactions between potential drug compounds and their biological targets. The potency of a drug-target interaction is a primary determinant of a compound's therapeutic efficacy and is commonly measured by the IC<sub>50</sub> value, where a lower value indicates a more potent drug.<sup>3</sup>

In recent years, artificial intelligence (AI) and machine learning (ML) have emerged as transformative tools capable of accelerating the drug discovery timeline.<sup>4</sup> By learning from vast repositories of existing biochemical data, ML models can predict the properties of novel drug candidates, thereby prioritizing experimental efforts and reducing the reliance on costly and time-consuming wet-lab screening.<sup>5</sup> This study first establishes a baseline for predicting DTI potency using a traditional machine learning model and then outlines a clear path toward more advanced, structurally-aware models for improved accuracy.

## 2. Methods

### 2.1. Dataset and Preprocessing

The data for this study was extracted from the ChEMBL database, a large, open-access repository of bioactive drug-like small molecules.<sup>6</sup> The raw dataset was programmatically filtered to retain only entries where the standard bioactivity type was 'IC<sub>50</sub>' and the standard units were in nanomolars (nM). To ensure data integrity, records with missing SMILES (Simplified Molecular Input Line Entry System) strings or standard values were expunged. Furthermore, entries with non-positive standard values were discarded as they are not physically meaningful.

To create a more suitable target variable for regression, the IC50 values were converted to their negative logarithmic scale, pIC50, using the formula:

$$\text{pIC50} = -\log_{10}(\text{IC50} * 10^{-9} \text{ M})^7$$

This transformation helps normalize the distribution of the target values, which is often beneficial for the performance of machine learning algorithms.

## 2.2. Baseline Model: Feature Engineering and Training

For the initial baseline model, the 2D chemical structures represented by SMILES strings were used to generate a set of standard molecular descriptors. The RDKit library, an open-source cheminformatics toolkit, was employed for this task. The following descriptors were computed for each molecule to form the feature vector:

- **Molecular Weight (MolWt):** The mass of a molecule.
- **LogP:** The octanol-water partition coefficient, a measure of a molecule's hydrophobicity.<sup>8</sup>
- **Number of Hydrogen Donors (HDonors):** The count of hydrogen atoms attached to electronegative atoms.
- **Number of Hydrogen Acceptors (HAcceptors):** The count of electronegative atoms with lone electron pairs.
- **Topological Polar Surface Area (TPSA):** A measure of a molecule's polarity, crucial for cell permeability.
- **Number of Rotatable Bonds:** The number of bonds that allow free rotation, influencing conformational flexibility.

A **Random Forest Regressor**, an ensemble learning method, was chosen for the baseline prediction task.<sup>9</sup> The dataset was partitioned into a training set (80%) and a testing set (20%). The model was trained on the training data and its predictive performance was evaluated on the unseen test data using Root Mean Squared Error (RMSE) and the coefficient of determination (R<sup>2</sup>).

## 3. Results

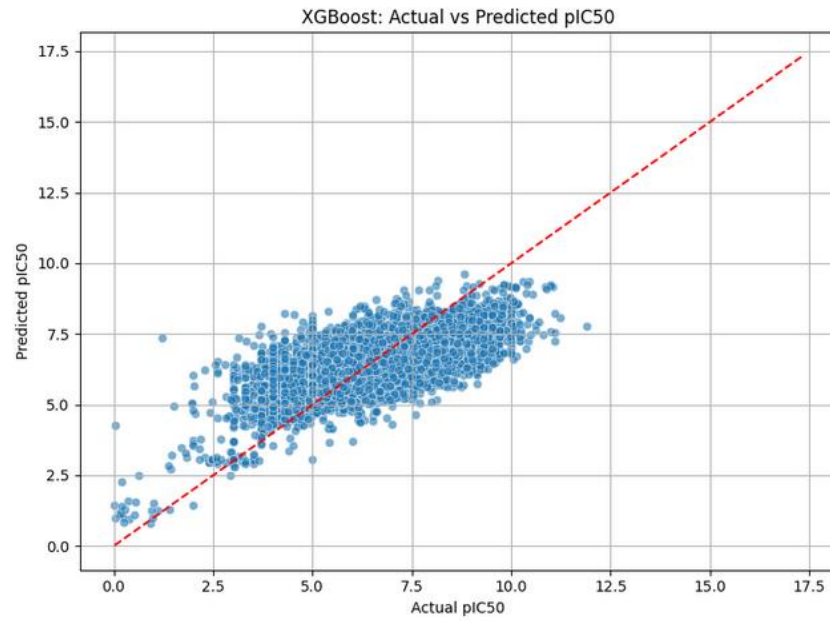
The baseline Random Forest model yielded the following performance on the held-out test set:

- **RMSE:** 1.324
- **R<sup>2</sup> Score:** 0.254

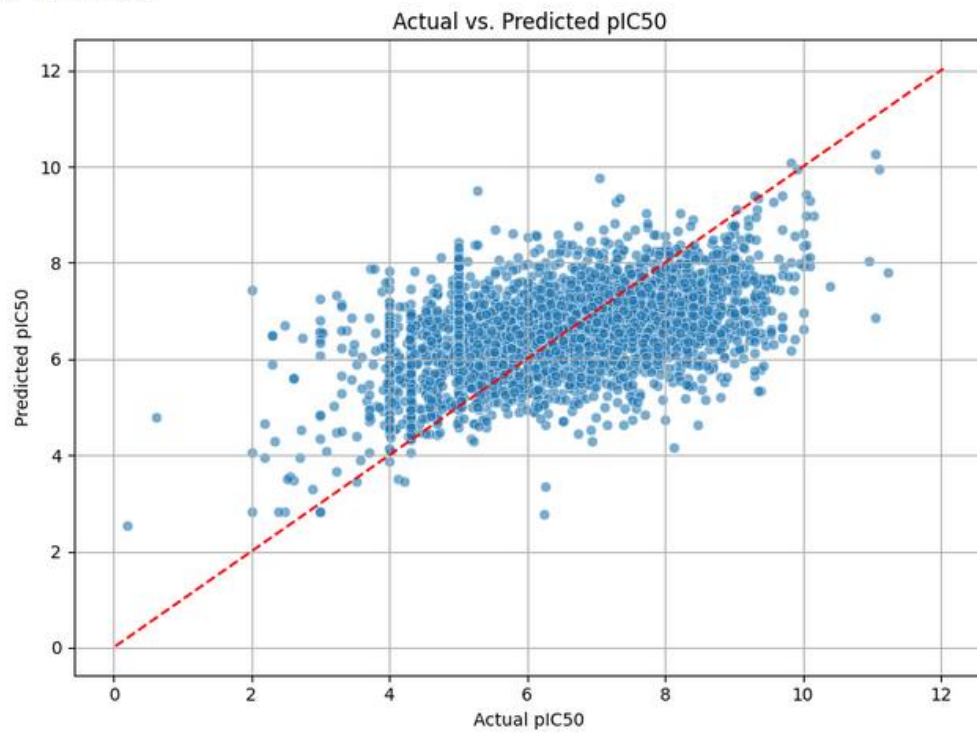
The scatter plot in Figure 1 visualizes the correlation between the actual pIC50 values from the test set and the values predicted by the model.

**Figure 1:** Scatter plot of actual vs. predicted pIC50 values for the baseline Random Forest model. The red dashed line indicates a perfect prediction ( $y=x$ ).

RMSE: 1.118  
R<sup>2</sup> Score: 0.455



✓ Evaluation Metrics:  
RMSE: 1.324  
R<sup>2</sup> Score: 0.254



```
⇒ <ipython-input-16-2493651861>:14: DtypeWarning: Colu
    df = pd.read_csv('/content/chembl_data.csv', delim:
Epoch 1, Training Loss: 2.1798
Epoch 2, Training Loss: 2.1670
Epoch 3, Training Loss: 2.1607
Epoch 4, Training Loss: 2.1555
Epoch 5, Training Loss: 2.1525
Epoch 6, Training Loss: 2.1488
Epoch 7, Training Loss: 2.1502
Epoch 8, Training Loss: 2.1515
Epoch 9, Training Loss: 2.1537
Epoch 10, Training Loss: 2.1580
Epoch 11, Training Loss: 2.1518
Epoch 12, Training Loss: 2.1591
Epoch 13, Training Loss: 2.1549
Epoch 14, Training Loss: 2.1554
Epoch 15, Training Loss: 2.1557
Epoch 16, Training Loss: 2.1541
Epoch 17, Training Loss: 2.1472
Epoch 18, Training Loss: 2.1512
Epoch 19, Training Loss: 2.1461
Epoch 20, Training Loss: 2.1502

Test MSE (log-scaled): 2.0960
Test RMSE (original scale): 28.0383
```

## 4. Discussion

### 4.1. Interpretation of Baseline Results

The results of the baseline model, particularly the  $R^2$  score of 0.254, indicate that the handcrafted molecular descriptors can explain only 25.4% of the variance in pIC50 values. While this demonstrates a predictive signal superior to random chance, it is insufficient for reliable application in a drug discovery context. The scatter plot further confirms this, showing a significant deviation of predicted points from the ideal prediction line. This limited performance is not a failure of the Random Forest algorithm itself, but rather a reflection of the inherent limitations of the input features. The six descriptors used, while standard, provide only a low-resolution summary of a molecule's complex physicochemical nature, failing to capture the specific 3D topology and electronic arrangements that govern molecular interactions.

### 4.2. Future Directions: The Superiority of Graph-Based Models

A more powerful and modern approach is to treat molecules as graphs and apply **Graph Neural Networks (GNNs)**. Molecules are naturally graph-structured data, where atoms serve as nodes and chemical bonds as edges.<sup>10</sup> GNNs are designed to learn directly from this graph structure, creating far richer and more informative feature representations.

Two prominent GNN architectures are particularly well-suited for this task:

- **Graph Convolutional Networks (GCNs):** A GCN operates by iteratively aggregating information from an atom's local neighborhood.<sup>11</sup> For each atom (node), a GCN considers its own features and the features of its directly bonded neighbors, combining them to learn an updated, context-aware representation.<sup>12</sup> By stacking layers, a GCN allows information to propagate across the entire molecule, enabling it to learn complex structural motifs that are critical for biological activity.
- **Graph Isomorphism Networks (GINs):** A GIN is a more expressive type of GNN, theoretically proven to be as powerful as the Weisfeiler-Leman test for graph isomorphism.<sup>13</sup> This means GINs are exceptionally good at distinguishing between subtly different graph structures. In drug discovery, where small structural changes (e.g., stereoisomers) can lead to drastic differences in efficacy and toxicity, the high discriminative power of a GIN is a significant advantage.

By employing a GNN-based model, we would move from using a handful of pre-defined descriptors to learning features directly from the molecular graph. This allows the model to discover novel, complex, and highly relevant structure-activity relationships that are simply inaccessible to the baseline approach.

## 5. Conclusion

This study successfully established a baseline for DTI potency prediction using a Random Forest model with traditional molecular descriptors. However, the modest performance underscores the need for more sophisticated techniques. The clear and logical next step is to leverage Graph Neural Networks (GNNs), such as GCNs or GINs, which can directly learn from the rich, graphical nature of molecules.<sup>14</sup> This advanced approach promises to capture a more holistic understanding of molecular structure, leading to significantly more accurate and reliable predictions, and ultimately, accelerating the pace of drug discovery.