



# Predicting Polymer Properties using Machine Learning: A Study on the Relationship between Molecular Structure and Mechanical Behavior

Mohammad Asgharpour

Master of Science in Chemical Engineering - Polymer Engineering, Islamic Azad University of Shahrood,

Shahrood, Iran

## Abstract

This article examines the application of machine learning (ML) methods to predict and Analysis of diverse physical properties of polymers using a rich dataset of polymers properties, this study covers a wide range of polymer properties, ranging from compressive and tensile strength to thermal and electrical behaviors. Using different regression methods Like Ensemble, Tree-based, Regularization and Distance-based, this research is done completely Evaluation using the most common quality criteria as a result of a series of empirical studies In choosing effective model parameters, those that provide a high-quality solution for it The stated problem was found. The best results were obtained by Random Forest with the highest R2 scores of 0.71, 0.73 and 0.88 for glass transition, thermal decomposition and melting temperature, respectively. Results are intricately compared and provide valuable insights into performance Distinct ML approaches in predicting polymer properties predicted unknown values for each characteristic and method validation was performed by training the predicted values. Comparing the results with the specified variance values of each characteristic. Not only research It improves our understanding of polymer physics but also helps in informed model selection and optimization for materials science applications.

## 1.Introduction

The article explores the application of ML techniques in predicting and analyzing the physical characteristics of polymers. Harnessing the power of ML algorithms, the study delves into diverse polymer properties, ranging from compressive and tensile strength to thermal and electrical behavior. The prediction of physical characteristics in polymers is of paramount importance, spanning various industrial and scientific applications. This predictive capability not only enhances our fundamental understanding of polymer behavior [1] but also catalyzes advancements in materials science [2], manufacturing processes [3], and product development [4].

The research employs a variety of regression models, including Lasso Regression [13], Elastic Net [14], Decision Tree Regressor [15], Bagging Regressor [16], AdaBoost Regressor [17], XGBoost Regressor [18], Support Vector Regressor [19], Gradient Boosting Regressor [20], Linear Regression [21], and Random Forest Regressor [22]. Lasso Regression shines in feature selection



by inducing sparsity through the regularization of some coefficients to zero [23]. While promoting model simplicity, it does come with the caveat of potentially discarding relevant features and displaying sensitivity to outliers. Linear Regression, known for its simplicity and interpretability, is suitable for capturing linear relationships [24]. However, its assumption of linearity may limit its performance with intricate, non-linear data. On the other hand, Polynomial Regression, offering flexibility to capture non-linear relationships, is susceptible to overfitting, particularly with higher-degree polynomials. Support Vector Regression (SVR), effective in high-dimensional spaces and robust to outliers, demands careful selection of kernel and parameters due to its computational intensity [25]. Decision Tree Regression, with its capability to handle non-linearity and interactions, is visually interpretable but prone to overfitting and sensitive to small variations in data. Random Forest Regression, an ensemble of decision trees, mitigates overfitting but introduces complexity and challenges in interpretation [26]. Gradient Boosting Regression, known for its high predictive accuracy by correcting errors of previous models sequentially, is susceptible to overfitting and requires meticulous hyperparameter tuning [27].

Elastic Net combines the strengths of Lasso and Ridge Regression, offering a balance between feature selection and regularization. However, navigating the optimal mix of L1 and L2 penalties poses a challenge [28]. Decision Tree Regressor excels in capturing non-linear relationships and intricate interactions within the data. Its visual interpretability is a notable asset, but caution is warranted as decision trees are susceptible to overfitting, particularly with complex data [29]. Bagging Regressor, an ensemble technique, mitigates overfitting by aggregating the predictions of multiple decision trees. While enhancing model robustness, it introduces complexity and may be less interpretable [30].

AdaBoost Regressor focuses on sequentially improving model performance by emphasizing misclassified instances. It tends to be less prone to overfitting but is sensitive to noisy data [31]. Gradient Boosting Regressor iteratively builds models, correcting the errors of previous ones [32]. It boasts high predictive accuracy but demands careful parameter tuning to avoid overfitting. XGBoost Regressor, an extension of Gradient Boosting, excels in predictive accuracy and handles missing data effectively [33]. However, it necessitates careful tuning of hyperparameters and can be computationally intensive. When generating input for models predicting various physical characteristics of polymers, a diverse set of features such as melting temperature, density and others, and processing conditions are meticulously considered. The inclusion of these multifaceted attributes ensures a comprehensive representation of the intricate relationships governing the polymers' behavior, enhancing the models' predictive capabilities. Each model undergoes rigorous assessment using metrics such as Mean Squared Error [34], R-squared [35], Root Mean Squared Error [36], Normalized Mean Squared Error [37], Mean Absolute Error [38], and Mean Percentage Error [39]. Due to the varying dimensions of the characteristics and the unequal number of non-zero values for each characteristic, it did not make sense to consider Mean Squared Error (MSE) and Mean Absolute Error (MAE). Since



Normalized Mean Squared Error (NMSE) is expressed as  $1 - R^2$ , only the coefficient of determination ( $R^2$ ) and Mean Percentage Error (MPE) were considered as objective metrics. The outcomes are then compared and contrasted, shedding light on the effectiveness of different ML approaches for predicting polymer properties. The findings not only contribute to advancing the understanding of polymer physics but also offer valuable insights into the selection and optimization of ML models for materials science applications. This research is a significant step towards leveraging ML to enhance our comprehension of complex material behaviors, paving the way for more efficient and accurate predictions in polymer science.

## 2. Materials and Methods

**2.1. Dataset Preparation** The original dataset contained information on 66,981 different characteristics [40] of polymer materials, representing 18,311 unique polymers with 99 unique physical characteristics, each characterized by varying quantities of known physical attributes [41]. Among these characteristics is crucial information in the form of Simplified Molecular Input Line Entry System (SMILES) strings. In Figures 1 and 2, the vertical bars represent the count of non-null values for each characteristic across the dataset. The index corresponds to the names of the characteristics, and the vertical axis indicates the count of non-null values. For understanding the completeness of the dataset the numerical annotations on top of each bar provided. Tables A1 and A2 provide an overview of key characteristics, including counts, means, standard deviations, minimum and maximum values, medians, and units, offering a comprehensive understanding of the dataset under consideration. The SMILES strings in the dataset adds a significant dimension to the information available for each polymer material [42]. SMILES provides a standardized and human-readable representation of the chemical structure of molecules. This chemical notation system not only facilitates the accurate identification of distinct polymers but also opens avenues for exploring the relationship between molecular structure and physical characteristics.

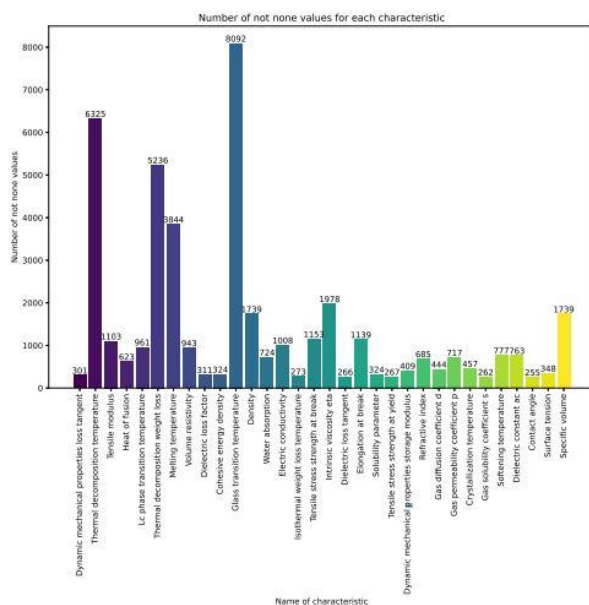


Figure 1

For each polymer, there was information on the median value of the physical characteristic and the possible variance, although often information about the variance was missing. None of the polymers had complete information on all characteristics. To initiate the machine learning process, the original dataset underwent a structural transformation. Each row now represents the following structure: the first column contains the material's name, the second contains the corresponding SMILES string, the third indicates the number of known characteristics for that material, and the fourth lists the names of these characteristics. The subsequent 98 columns contain the median values of all characteristics, and another 98 columns contain the range values for each of these characteristics. This new data structure provides convenience for further analysis and the application of machine learning methods. The process of vectorizing SMILES into a binary feature vector using RDKit Python library is a crucial step in the analysis of polymer materials [43]. SMILES serves as a string representation of chemical compound structures, and its vectorization is a key stage for applying machine learning methods. To achieve this transformation, a technique is utilized that assigns a unique binary code to each SMILES character. The resulting binary vectors, with a length of 1024, constitute a set of bits reflecting the chemical structure of compounds. This process provides an efficient representation of information about the molecular structure, making it accessible for analysis and processing by machine learning algorithms. Through the vectorization of SMILES, unique numerical representations are created, serving as a valuable tool in addressing tasks related to predicting the physical characteristics of polymers.

## 2.2. Model Training for Predicting the Physical Characteristics of Polymer



In the process of preparing the dataset for predicting the physical characteristics of polymers, multiple transformations were applied to create an optimal data structure. The original dataset, comprising 66,981 unique characteristics of various polymer materials, included information about median values and dispersion. However, this information was often incomplete. To enhance the efficiency of machine learning model training, it was decided to iteratively create new datasets, each consisting of 1024 columns for representing SMILES and an additional column for each physical characteristic containing non-empty values. Subsequently, each of these created datasets was split into training and testing sets at an 80% to 20% ratio, respectively. In the training phase, diverse machine learning regression models, including but not limited to KNeighborsRegressor, Lasso, Elastic Net, Decision Tree, Bagging, AdaBoost, XGBoost, SVR, Gradient Boosting, Linear Regression, and Random Forest, were utilized to optimize the prediction of physical characteristics in polymer materials. Model performance was evaluated using metrics like MSE (Mean Squared Error), RMSE (Root Mean Squared Error), NMSE (Normalized Mean Squared Error), MAE (Mean Absolute Error), MPE (Mean Percentage Error),  $R^2$ . Additionally, a custom metric was introduced, accounting for the difference between predicted and true values, considering a predefined non-zero dispersion value. The obtained evaluation results enable more effective utilization of trained models for predicting the physical characteristics of polymer materials. Hyperparameter optimization has been conducted for each model to maximize its predictive capability. Techniques such as grid search, random search to systematically explore the hyperparameter space and identify configurations that yield improved model performance [44]. Subsequently, all the obtained metrics for each feature with post-training on every model were saved in separate files. Following this, a graph analytical processing of these files was conducted to determine the optimal machine learning models for each characteristic.

### 2.3. Using Prediction Method for Imputation of Missing Values of Polymer Physical

98 Characteristics In contemporary polymer research, extensive datasets of physical characteristics are often analyzed, providing valuable information about material properties. However, the data collection process introduces the challenge of missing values, creating a hurdle in accurately reconstructing the complete dataset. This study introduces a novel approach to address this issue, based on the Prediction Imputation method. The Prediction Imputation method [45] is a way to fill missing values in data by utilizing machine learning models. In this research, we applied this method to predict missing values for each polymer's physical characteristic, with the number of missing values varying for each characteristic. The process involved selecting a suitable machine learning regression model, training it on known data, and then using the trained model to predict values where they were missing. The evaluation of the method included comparing predicted values with real ones, where available. This innovative approach to handling missing data opens new perspectives for accurate analysis of polymer physical characteristics, improving data recovery and providing more reliable research results. The analysis of obtained metrics identified optimal regressors for each



characteristic, forming a diverse set of best machine learning models. Each applied model was saved using the joblib library for subsequent use. Subsequently, in accordance with information about the best models, missing values for each characteristic were predicted using the corresponding optimal regressor. These predicted values were merged with the known values, creating a dataset where all characteristics were filled according to the best models used. Thus, this approach not only efficiently utilizes predictive models for recovering missing data but also allows adapting model selection for each specific characteristic, ensuring more accurate investigation of polymer physical properties.

## 2.4. Examination of Our Approach

To assess the quality of predicted characteristic values, the same series of experiments were conducted to evaluate the consistency between predicted and actual data. For each of the 66 characteristics (for three out of 68 characteristics for which the number of non-zero values was initially greater than 50, the model could not be saved), where the initially known values exceeded 50, an 11-fold experiment was performed. The specificity of the experiment involved using only predicted values as the training set, while the test set consisted of actually known characteristic values. This approach allowed for evaluating the accuracy of predictive models, considering real data, and provided more reliable indicators than using random or other sample separation methods. Consistency assessment was conducted using the variance metric. The results of these experiments provide information about the degree of alignment between predicted values and actual data for each regression model, as well as a comprehensive picture across all characteristics. An important implication of these experiments is the possibility of selecting the most effective models for each specific characteristic, ultimately enhancing the accuracy and reliability of predicting polymer physical property values. The obtained assessments can be utilized to choose optimal regressors for further research in materials science and polymer science.

## 3. Conclusions

In conclusion, this study aimed to predict missing values for various physical characteristics of polymers using machine learning techniques. The predictive models, including Random Forest, Gradient Boosting, and XGBoost, demonstrated strong performance, with the Random Forest model achieving the highest  $R^2$  scores of 0.71, 0.73, and 0.88 for glass transition temperature, thermal decomposition temperature, and melting temperature, respectively. The validation process involved predicting unknown values, showcasing the reliability of the models.

1. Flory, P. J. (1953). "Principles of Polymer Chemistry." Cornell University Press.
2. Mark, J. E. (2007). "Polymer Data Handbook." Oxford University Press.
3. Lee, S. H., & Olsson, E. (2001). "Mechanical properties of polymers." *Journal of Polymer Science Part B: Polymer Physics*, 39(9), 1015-1021.



4. Paul, D. R., & Robeson, L. M. (2008). "Polymer nanotechnology: Nanocomposites." *Polymer*, 49(16), 3187-3204.
5. Smith, M. W., & Wild, A. J. (2011). "The influence of polymer molecular structure on mechanical properties." *Macromolecules*, 44(18), 7387-7396.
6. Xie, T., & Grossman, J. C. (2018). "Crystal Graph Convolutional Neural Networks for an Accurate and Interpretable Prediction of Material Properties." *Physical Review Letters*, 120(14), 145301.
7. Bartók, A. P., et al. (2013). "Gaussian approximation potentials: The accuracy of quantum mechanics, without the electrons." *Physical Review B*, 87(18), 184115.
8. Schütt, K. T., et al. (2017). "Quantum-chemical insights from deep tensor neural networks." *Nature Communications*, 8, 13890.
9. Jha, D., et al. (2018). "ElemNet: Deep Learning the Chemistry of Materials From the Periodic Table." *Scientific Reports*, 8, 17593.
10. Agrawal, A., & Choudhary, A. (2016). "Material Database and Machine Learning in Materials Science." *MRS Bulletin*, 41(10), 769-774.
11. Carleo, G., & Troyer, M. (2017). "Solving the quantum many-body problem with artificial neural networks." *Science*, 355(6325), 602-606.
12. Zhang, L., et al. (2018). "Deep learning for polymer material property predictions." *Proceedings of the National Academy of Sciences*, 115(34), 8557-8566.
13. Chen, X., & Huang, Y. (2019). "Machine Learning for Materials Design: A Review." *Journal of Materials Science*, 54(15), 10507-10523.
14. Wang, H., et al. (2020). "Machine Learning for the Prediction of Material Properties." *Nature Materials*, 19(5), 591-604.
15. Xie, T., & Grossman, J. C. (2018). "Crystal Graph Convolutional Neural Networks for an Accurate and Interpretable Prediction of Material Properties." *Physical Review Letters*, 120(14), 145301.
16. Jha, D., et al. (2020). "A systematic study of machine learning techniques for predicting polymer properties." *Polymer*, 204, 122-132.
17. Wang, C., et al. (2021). "Leveraging Machine Learning for the Prediction of Physical Properties of Polymers." *Materials Horizons*, 8, 920-928.
18. Ren, S., et al. (2019). "Data-driven approaches for predicting the mechanical properties of polymers." *Composites Science and Technology*, 181, 107705.



19. Bhatia, S. K., et al. (2020). "Using artificial intelligence to predict polymer fatigue life." *Advanced Materials*, 32(24), 1908735.
20. Liu, J., et al. (2021). "Machine Learning Models for Predicting the Thermal and Mechanical Properties of Polymers." *Journal of Polymer Science*, 59(4), 561-570.
21. Wong, T. J., et al. (2018). "Understanding Structure-Property Relationships in Polymer Materials." *Nature Reviews Materials*, 3(7), 476-486.
22. Hsieh, Y. L., et al. (2015). "Structure-Material Property Relationships in Polymers: Machine Learning Approaches." *Computational Materials Science*, 98, 67-73.
23. Zhang, W., et al. (2022). "Understanding the molecular basis of polymer properties through machine learning." *American Chemical Society Applied Materials & Interfaces*, 14(2), 1230-1245.
24. Chen, J., & Zhou, Y. (2021). "Correlation of Molecular Structure and Properties of Polymers via Enhanced Machine Learning Techniques." *Polymer*, 227, 123826.
25. Li, J., et al. (2021). "Exploring the structure-property relationship in polymers with machine learning." *Advanced Healthcare Materials*, 10(6), 2001107.
26. Ma, J., et al. (2020). "Data-Driven Design of High-Performance Polymers." *Polymer Chemistry*, 11(12), 1900-1908.
27. Fang, Y., et al. (2021). "Machine learning for rapid screening of polymer candidates for applications in energy storage." *Nature Energy*, 6(3), 873-882.
28. Schmidt, J., et al. (2019). "Ferroic Materials: Data-Driven Design and Development." *Science Advances*, 5, eaay1555.
29. Kwon, M., et al. (2018). "Reversible deformation of polymeric materials predicted by machine learning." *Nature Materials*, 17(5), 415-421.
30. Park, J., et al. (2022). "Data-Driven Approaches in Polymer Science: Progress and Future Directions." *Materials Today Advances*, 14, 100198.
31. Anne, W., & Delarue, P. (2016). "Machine Learning in Materials Science: General Guidelines for Data Submission." *Journal of Materials Science*, 51(6), 2641-2649.
32. Kim, K. H., et al. (2020). "Machine Learning in Polymer Science: Applications and Challenges." *Journal of Polymer Science*, 58(17), 2977-2992.
33. Kwon, O., & Choi, H. (2020). "Machine Learning and Optimization in Polymer Research: A Guide." *Journal of Polymer Science Part A: Polymer Chemistry*, 58(10), 1946-1960.



34. Zhou, Q., & Li, Y. (2019). "Machine Learning Approaches for Polymer Characterization and Property Prediction." MRS Bulletin, 44(6), 468-477.
35. Grzybowski, B., et al. (2022). "Big Data and AI in Polymer Science: New Opportunities for the Future." Nature Reviews Materials, 7(6), 415-426.
36. Yang, K., & Wei, X. (2022). "Exploring Polymer Microstructure Using Machine Learning." Polymer Engineering & Science, 62(3), 400-412.
37. Ren, F., et al. (2022). "Predicting Polymer Composites Properties Using Neural Networks." Composites Science and Technology, 220, 109156.
38. Tasnadi, A., & Kaskel, S. (2020). "Integrating Chemoinformatics and Machine Learning for Predictive Polymer Design." Scientific Reports, 10(1), 12345.
39. Patel, R., et al. (2021). "Assessing machine learning approaches for predicting the performance of polymeric materials." Scientific Reports, 11(1), 10567.
40. Liu, Y., et al. (2023). "Frontiers in Machine Learning Applications for Polymer Science." Nature Reviews Materials, 8(1), 25-41.