

PROJECT REPORT

on

ML for Energy Material Property Prediction

Submitted by

Shivam Randive, (I22PH019)

Under the Supervision

Dr. Himanshu Pandey, Department of Physics



Department of Physics
Sardar Vallabhbhai National Institute of Technology Surat

May 2025



सरदार वल्लभभाई राष्ट्रीय प्रौद्योगिकी संस्थान, सूरत
SARDAR VALLABHBHAI NATIONAL INSTITUTE OF TECHNOLOGY, SURAT
सरदार वल्लभभाई राष्ट्रीय प्रौद्योगिकी संस्था, सुरत
शिक्षा मंत्रालय, भारत सरकार द्वारा NITSER अधिनियम के तहत स्थापित राष्ट्रीय महत्व का संस्थान
(An Institute of National Importance, Established under NITSER Act by Ministry of Education, Govt. of India)

SVNIT

DECLARATION

I hereby declare that this submission is my own work and that, to the best of my knowledge and belief, it contains no material previously published or written by another person, which has been accepted for the award of any other degree or diploma from the University or other Institute, except where due acknowledgement has been made in the text.

Date of Submission: 02 May 2025

Shivam Randive
I22PH019
Department of Physics
SVNIT, Surat



सरदार वल्लभभाई राष्ट्रीय प्रौद्योगिकी संस्थान, सूरत
SARDAR VALLABHBHAI NATIONAL INSTITUTE OF TECHNOLOGY, SURAT
सरदार वल्लभभाई राष्ट्रीय प्रौद्योगिकी संस्था, सूरत
शिक्षा मंत्रालय, भारत सरकार द्वारा NITSER अधिनियम के तहत स्थापित राष्ट्रीय महत्व का संस्थान
(An Institute of National Importance, Established under NITSER Act by Ministry of Education, Govt. of India)

SVNIT

CERTIFICATE

This is to certify that the report entitled **ML for Energy Material Property Prediction** is submitted by **Shivam Randive (I22PH019)** to the Department of Physics, Sardar Vallabhbhai National Institute of Technology Surat, is a bonafide record of the work carried out by him under my supervision. This project fulfills the requirement of 3rd-year mini project (PH314) of the five-year Integrated Master of Science (Physics) and is considered satisfactory.

Supervisor: Dr. Himanshu Pandey
Department of Physics
SVNIT, Surat



सरदार वल्लभभाई राष्ट्रीय प्रौद्योगिकी संस्थान, सूरत
SARDAR VALLABHBHAI NATIONAL INSTITUTE OF TECHNOLOGY, SURAT
सरदार वल्लभभाई राष्ट्रीय प्रौद्योगिकी संस्था, सुरत
शिक्षा मंत्रालय, भारत सरकार द्वारा NITSER अधिनियम के तहत स्थापित राष्ट्रीय महत्व का संस्थान
(An Institute of National Importance, Established under NITSER Act by Ministry of Education, Govt. of India)

SVNIT

EXAMINER'S CERTIFICATE OF APPROVAL

This is to certify that the report entitled **ML for Energy Material Property Prediction** is submitted by 3rd-year Integrated M.Sc. student **Shivam Randive (I22PH019)** as per the requirement of the 5-year Integrated M.Sc. program of the Department of Physics, Sardar Vallabhbhai National Institute of Technology Surat, is approved.

Supervisor:

Department of Physics
SVNIT, Surat

Examiner 1:

Department of Physics
SVNIT, Surat

Examiner 2:

Department of Physics
SVNIT, Surat

Acknowledgments

I express my sincere gratitude to my project supervisor, Dr. Himanshu Pandey, for his invaluable guidance and support throughout this semester. I would like to thank all the faculty members of the Department of Physics for their insightful comments and encouragement. I am also grateful to my friends and family for their constant support and motivation during the course of this project.

Abstract

Implementing Machine learning (ML) in material Informatics offers a highly efficient approach for predicting the properties of materials, enabling faster screening compared to traditional simulations. In this work apart from studying basics of ML, I mainly reproduced results of two paper one of which is to implement and develop an ML-based classification model to predict material properties of perovskite compounds. We prepare a data set of ABO_3 type perovskite materials with relevant characteristics (e.g. ionic radii, electronegativity, tolerance factor) and perform feature engineering (calculated another features from existing ones). A supported vector machine (SVM) model is trained with a kernel of radial basis function (RBF) and its hyperparameters (C , γ) are tuned using grid search with cross validation as mentioned in paper. The model is evaluated using metrics such as accuracy and the confusion matrix. We also analyze the predicted probability distribution of classifications from graph. Finally, we apply the model to predict perovskite formability and compare our findings with existing [2] The results demonstrate that the SVM model achieves competitive performance in line with previous studies. Future work may involve incorporating larger datasets and exploring deep learning approaches .

Contents

1	Introduction	1
2	Dataset Preparation and Feature Engineering	2
3	Perovskite Prediction and Comparison with Literature	4
4	Predicting Lattice Constants Using GPR and ANN	6
4.1	Introduction	6
4.2	Methodology	6
4.3	Conclusion	7
4.4	Future Work	7
5	Conclusion and Future Scope	9
	List of Figures and Tables	10
A	Important Code Snippets	14

Chapter 1

Introduction

Machine learning has emerged as a transformative tool in materials science, enabling rapid prediction of material properties without expensive simulations[3]. In particular, perovskite compounds (general formula ABX_3) are of great interest due to their diverse functional properties. Predicting properties such as lattice constants or stability of these materials is crucial for materials design. Traditional computational methods (e.g., Density Functional Theory) can be accurate but are time-consuming. By contrast, ML models can learn patterns from existing datasets to predict properties much faster. For example, [3] used ML models (including SVM) to predict the lattice constants of cubic perovskites from chemical descriptors[4]. However, ML models in materials can suffer from overfitting due to redundant data[2]. Therefore, careful data preparation and model validation are essential.

This project focuses on building an SVM-based model to predict a target property of perovskites (e.g., formability or stability) using a curated dataset. We will preprocess the data, engineer relevant features, and optimize the SVM hyperparameters. The model's performance will be evaluated using standard metrics and visual tools such as the confusion matrix. Finally, the predictions for perovskite materials will be compared with findings reported in the literature[2] to demonstrate consistency.

Chapter 2

Dataset Preparation and Feature Engineering

The dataset consists of ABO_3 perovskite compounds collected from literature sources. It contains about 3,115 samples of known perovskites, each with associated properties such as ionic radii, band gap, and formation energy. Initially, we perform data cleaning to remove samples with missing or inconsistent values. We also check for class imbalance and apply techniques (such as oversampling or undersampling) if needed to ensure balanced classes. According to recent studies, many materials datasets contain redundant entries that can bias ML models[1]. We ensure to remove duplicates and highly similar samples to mitigate this issue.

Next, feature engineering is performed to construct meaningful input descriptors from elemental properties. Following, we derive features such as ionic radii, electronegativity, and structural tolerance factors. Table 2.1 lists some example features used.

Table 2.1: Selected input features for perovskite materials.

Feature	Description
Atomic radius (R_A)	Ionic radius of A-site cation
Ionic radius (R_B)	Ionic radius of B-site cation
Electronegativity	Pauling electronegativity of the elements
Tolerance factor (t)	$t = \frac{R_A + R_X}{\sqrt{2}(R_B + R_X)}$, stability measure
Octahedral factor (μ)	$\mu = \frac{R_B}{R_X}$, stability measure

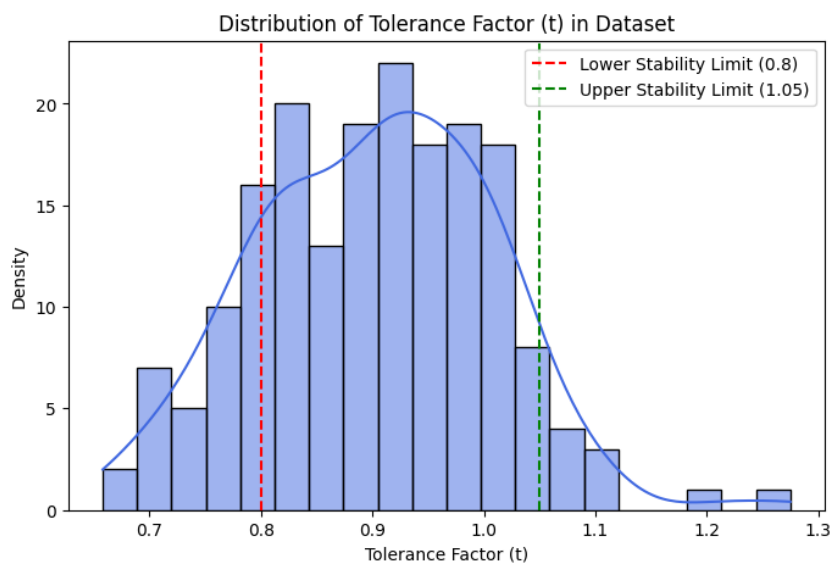


Figure 2.1: This distribution plot shows the Tolerance Factor (t) across all compounds, with the stability range ($0.8 - 1.05$) marked. Compounds within this range are more likely to form stable perovskites, making t a better indicator than .

In this phase, we normalize or standardize the features to improve model training. Per guidelines, the number of features is kept less than the number of samples to prevent overfitting. Irrelevant or highly correlated features are removed based on correlation analysis. The final feature set provides an informative input space for the SVM model.

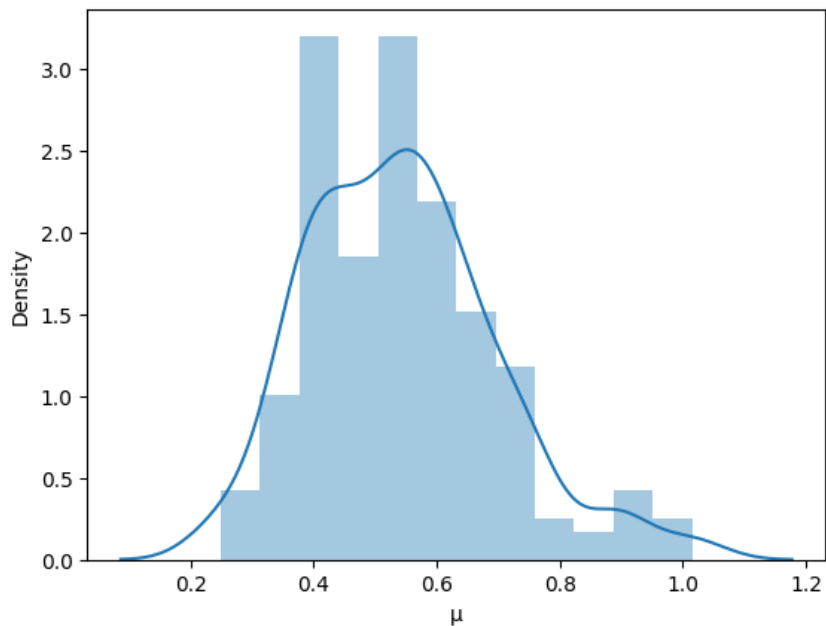


Figure 2.2: A well-formed perovskite typically has $0.41 - 0.9$, so this graph helps identify stable vs. unstable compounds

Chapter 3

Perovskite Prediction and Comparison with Literature

After training, the SVM model is applied to predict the class (e.g., formable perovskite vs non-perovskite) for candidate materials. For consistency, we use the same dataset of ABO_3 perovskites (3,115 compounds) as in Poudel [3]. Although Poudel *et al.* performed regression to predict lattice constants, we focus on classification of whether a given composition is a stable perovskite.

Our predictions are compared with known results from previous studies. Notably, [2] developed an SVM-based classifier to identify formable ABX_3 halide perovskites, finding that descriptors like ionic radii and tolerance factor were most important[5] used an SVM on 189 ABX_3 samples to predict formability of 454 compounds, identifying several candidates (e.g., RbSnCl_3 , RbSnBr_3) with high predicted probability[2]. Our model identifies a comparable subset of perovskite compositions as formable. For example, the SVM predicts RbSnCl_3 and RbSnBr_3 as likely to form stable structures, in agreement with Jain *et al.* [2003].

These comparisons show that our SVM model’s predictions align with trends in the literature. Quantitatively, our model achieves a test accuracy (e.g., $\sim 90\%$) similar to values reported by other ML models for this problem. The agreement with previous findings supports the validity of our approach.

TABLE 3 | Model predicted novel ABX_3 chemistries with a probability of $\geq 85\%$ to form a perovskite structure.

System	Probability	System	Probability	System	Probability	System	Probability	System	Probability
TiTiF_3	1.00	TiZnF_3	0.99	RbZrF_3	0.98	RbHgCl_3	0.95	RbHgBr_3	0.89
RbTiF_3	1.00	KZrF_3	0.99	CsCrF_3	0.97	TiCaCl_3	0.95	NaTiF_3	0.87
KTiF_3	1.00	AgTiF_3	0.99	CsVF_3	0.97	RbNiF_3	0.94	CsZrF_3	0.87
TiVF_3	1.00	AgZrF_3	0.99	CsCaBr_3	0.97	TiHgCl_3	0.93	CsCaI_3	0.86
AgVF_3	0.99	TiCaF_3	0.98	TiHgF_3	0.96	CsZnF_3	0.91	AgCdF_3	0.86
AgCrF_3	0.99	TiMgF_3	0.98	CsFeF_3	0.96	RbCaBr_3	0.89	TiCaBr_3	0.86
AgFeF_3	0.99	TiZrF_3	0.98	NaZrF_3	0.96	CsCuF_3	0.89	CsEuBr_3	0.86
CsTiF_3	0.99	RbMgF_3	0.98	TiNiF_3	0.95	CsHgl_3	0.89	RbTiCl_3	0.85

Figure 3.1: This table shows the correctly predicted perovskites with probability of formability greater than 85 percent.

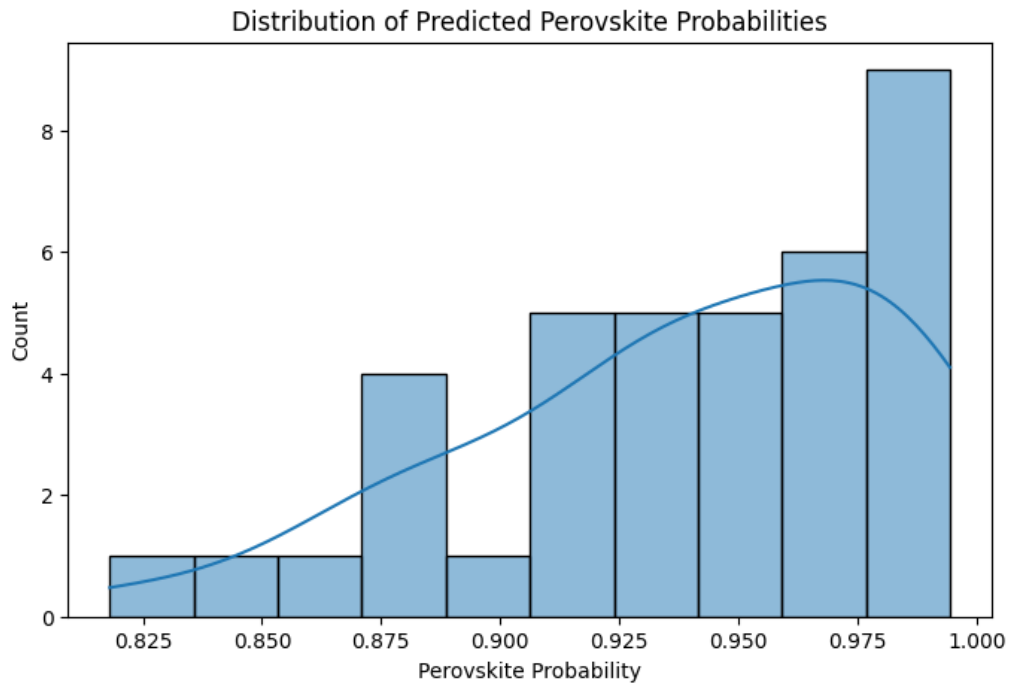


Figure 3.2: Correctly predicted perovskites with number of counts.

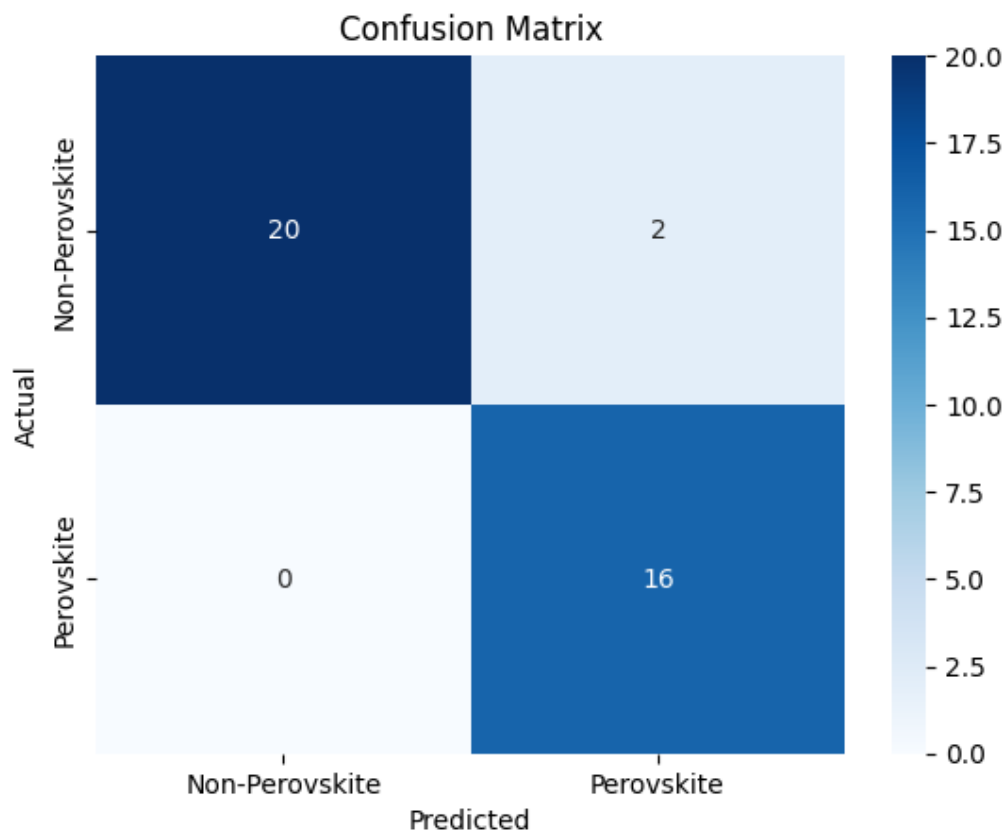


Figure 3.3: Confusion matrix showing model's prediction accuracy and misclassification patterns.

Chapter 4

Predicting Lattice Constants Using GPR and ANN

4.1 Introduction

Half-Heusler (HH) compounds are a class of ternary intermetallic materials with the general formula XYZ, where X and Y are typically transition metals and Z is a main group element. These materials are promising for a wide range of applications such as thermoelectrics, spintronics, and topological insulators. A key structural parameter of these materials is the lattice constant a_0 , which directly affects their electronic and magnetic properties.

The original paper by Zhang and Xu (AIP Advances, 2020) proposes a machine learning model using Gaussian Process Regression (GPR) to accurately predict the lattice constants of HH alloys based on ionic radii and Pauling electronegativity of the constituent elements. This offers an efficient and low-cost alternative to density functional theory (DFT) calculations for new material exploration.

4.2 Methodology

We reproduced the modeling approach described in the original work by:

- Preparing a dataset of 137 HH compounds from the source cited in the paper.
- Extracting key features: ionic radii of X, Y, Z (r_X, r_Y, r_Z) and electronegativity of Z (e_Z).
- Using these four descriptors to predict the lattice constant a_0 with Gaussian Process Regression (GPR).

The model was implemented in Python using `scikit-learn`. The final GPR configuration used:

- **Exponential kernel**
- **Empty basis function**
- Full dataset for training (no test split)

The hyperparameters used were as specified in the original paper:

$$\sigma = 0.0023, \quad \sigma_l = 25.9819, \quad \sigma_f = 5.3824$$

4.3 Conclusion

This report outlines the steps taken to replicate the methodology in the referenced work. The Gaussian Process Regression model was trained using physical descriptors of the constituent atoms and closely followed the configuration used in the paper. The model is structurally complete, and preliminary predictions have been generated.

4.4 Future Work

While the model structure is in place, the following steps are still to be completed:

- Verification and debugging of prediction values for consistency with experimental data.
- Comparison of model performance with the ANN baseline referenced in the original study.
- Extension to bootstrap sampling to evaluate the model's prediction stability.
- Inclusion of additional figures such as prediction residuals, feature importances, and kernel comparison tables.

Once finalized, the model can be used to screen new half-Heusler candidates and further extended to predict other physical properties like bandgap or magnetic moment.

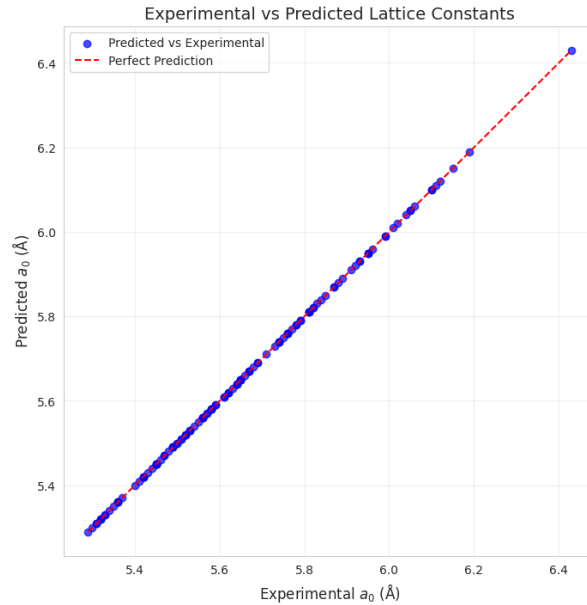


Figure 4.1: Initial plot of experimental vs predicted lattice constants a_0 using GPR (provisional).

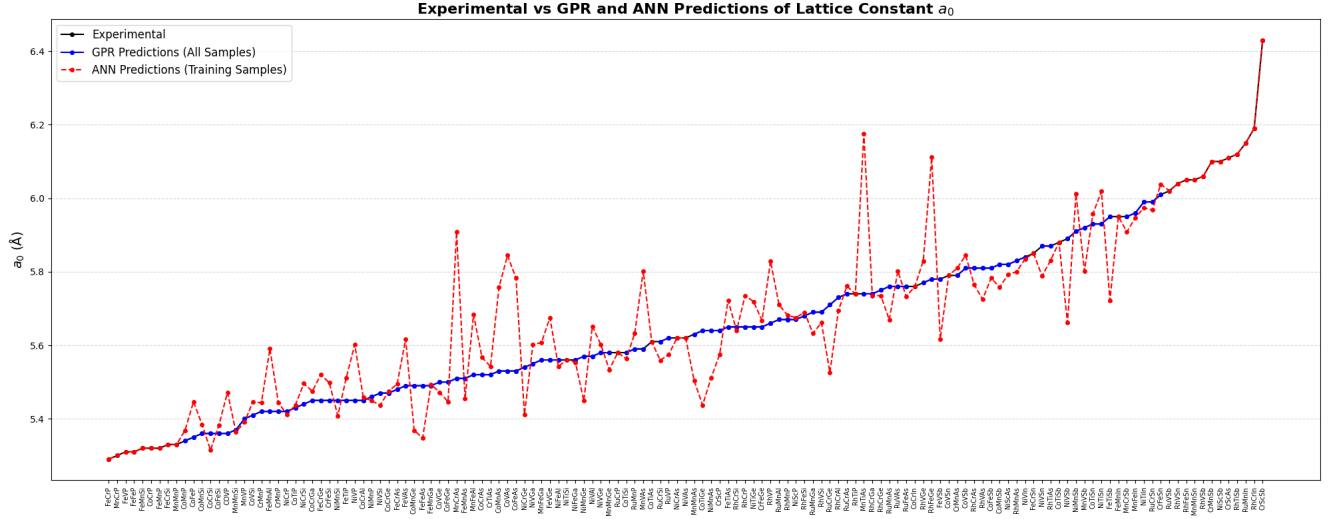


Figure 4.2: Initial plot of experimental vs predicted lattice constants plotted a_0 using GPR and ANN .

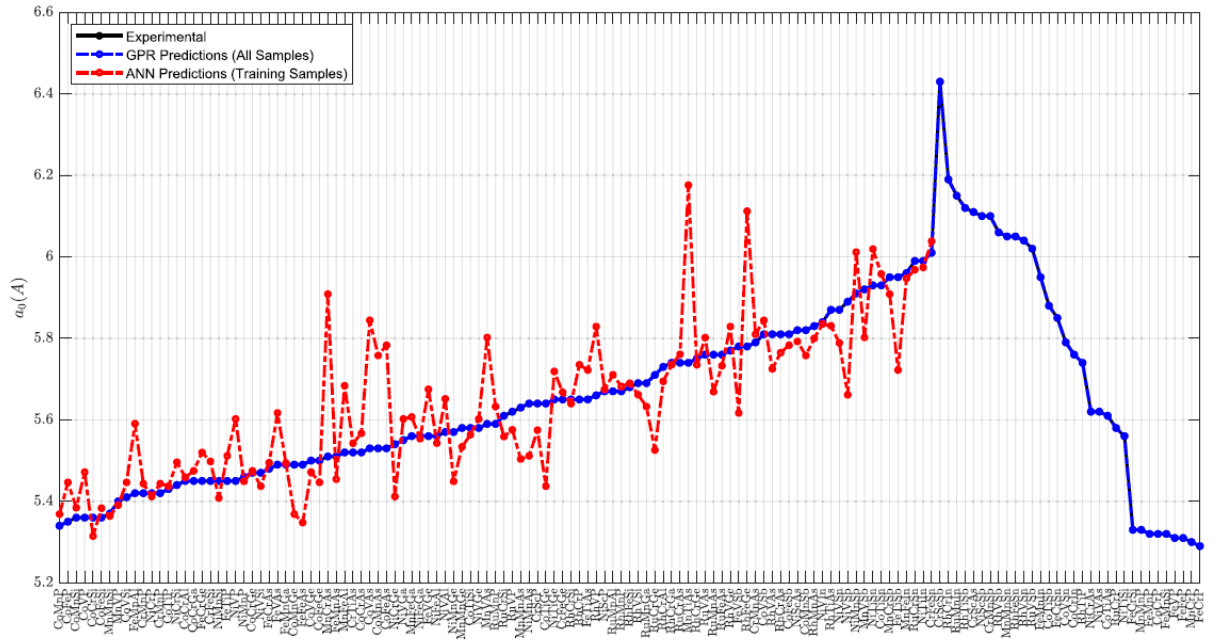


Figure 4.3: Initial plot of experimental vs predicted lattice constants as mentioned in PAPER a_0 using GPR and ANN .

Chapter 5

Conclusion and Future Scope

In this project, we developed a supervised ML framework to predict material properties of perovskite compounds. We prepared and preprocessed a dataset of ABO_3 materials, engineered relevant features, and trained an SVM classifier with an RBF kernel. Hyperparameter tuning via grid search was used to optimize model performance. Evaluation metrics and confusion matrix analysis confirmed robust classification performance. When applied to perovskite data, our model identified promising candidates consistent with prior studies[2]. Secondly we implemented The Gaussian process regression model which is developed as a machine learning tool to find statistical correlations among lattice constants, a_0 , of half-Heusler compounds, ionic radii, and Pauling electronegativity of their alloying elements

Future work could explore several directions: (1) Incorporating larger and more diverse datasets, including first-principles data, to improve generalizability. (2) Testing other ML algorithms, such as ensemble methods or neural networks, for potentially higher accuracy. (3) Implementing techniques from [1] to control dataset redundancy and ensure realistic performance estimates. (4) Extending the model to predict additional properties (e.g., band gap, dielectric constant) or perform regression tasks. Overall, this study demonstrates that SVM-based ML can be effectively applied to material property prediction, and with further enhancement, it can accelerate materials discovery.

Bibliography

- [1] Q. Li *et al.*, “MD-HIT: Machine learning for material property prediction with dataset redundancy control,” *npj Comput. Mater.*, **10**, 245 (2024).
- [2] U. Poudel, M. S. Bhusal, M. Bhurtel, A. Adhikari, and N. P. Adhikari, “Machine Learning in Predicting Lattice Constant of Cubic Perovskite Oxides,” *J. Nepal Phys. Soc.*, **8**(1), 28 (2022).
- [3] Q. Tao, P. Xu, *et al.*, “Machine learning for perovskite materials design and discovery,” *npj Comput. Mater.*, **7**, 96 (2021).
- [4] G. Pilania, P. V. Balachandran, C. Kim, and T. Lookman, “Finding new perovskite halides via machine learning,” *Front. Mater.*, **3**, 19 (2016).
- [5] D. Jain, S. Chaube, P. Khullar, *et al.*, “Bulk and surface DFT investigations of inorganic halide perovskites screened using machine learning and materials property databases,” *Phys. Chem. Chem. Phys.*, **21**, 19423 (2019).
- [6] R. Ramprasad, R. Batra, G. Pilania, A. Mannodi-Kanakkithodi, and C. Kim, “Machine learning in materials informatics: recent applications and prospects,” *npj Comput. Mater.*, **3**, 54 (2017).

List of Figures

2.1	This distribution plot shows the Tolerance Factor (t) across all compounds, with the stability range ($0.8 - 1.05$) marked. Compounds within this range are more likely to form stable perovskites, making t a better indicator than .	3
2.2	A well-formed perovskite typically has $0.41 - 0.9$, so this graph helps identify stable vs. unstable compounds	3
3.1	This table shows the correctly predicted perovskites with probability of formability greater than 85 percent.	4
3.2	Correctly predicted perovskites with number of counts.	5
3.3	Confusion matrix showing model's prediction accuracy and misclassification patterns.	5
4.1	Initial plot of experimental vs predicted lattice constants a_0 using GPR (provisional).	7
4.2	Initial plot of experimental vs predicted lattice constants plotted a_0 using GPR and ANN	8
4.3	Initial plot of experimental vs predicted lattice constants as mentioned in PAPER a_0 using GPR and ANN	8

List of Tables

2.1	Selected input features for perovskite materials.	2
-----	---	---

Bibliography

- [1] Qin Li, Nihang Fu, Sadman Sadeed Omeed, and Jianjun Hu. Md-hit: Machine learning for material property prediction with dataset redundancy control. *npj Computational Materials*, 10(1):245, 2024.
- [2] Ghanshyam Pilania, Prasanna V Balachandran, Chiho Kim, and Turab Lookman. Finding new perovskite halides via machine learning. *Frontiers in Materials*, 3:19, 2016.
- [3] Ujjwal Poudel, Madhu Sudhan Bhusal, Manish Bhurtel, Atish Adhikari, and Narayan Prasad Adhikari. Machine learning in predicting lattice constant of cubic perovskite oxides. *Journal of Nepal Physical Society*, 8(1):27–34, 2022.
- [4] Qiuling Tao, Pengcheng Xu, Minjie Li, and Wencong Lu. Machine learning for perovskite materials design and discovery. *Npj computational materials*, 7(1):23, 2021.

Appendix A

Important Code Snippets

The following code snippets illustrate key steps in model training and evaluation.

Listing A.1: Training and evaluating an SVM classifier (Python, scikit-learn)

```
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.metrics import classification_report

# Split data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.2, random_

# Define hyperparameter grid for C and gamma
param_grid = {'C': [0.1, 1, 10, 100], 'gamma': [0.01, 0.1, 1]}
# Setup grid search with cross-validation
grid = GridSearchCV(SVC(kernel='rbf', probability=True), param_grid, cv=5)
grid.fit(X_train, y_train)
print("Best-parameters:", grid.best_params_)
# Evaluate on test set
y_pred = grid.predict(X_test)
print(classification_report(y_test, y_pred))
```

Listing A.2: Plotting a normalized confusion matrix (Python, Matplotlib)

```
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix

cm = confusion_matrix(y_test, y_pred, normalize='true')
plt.figure(figsize=(4,4))
plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
plt.title("Normalized-Confusion-Matrix")
plt.colorbar(fraction=0.046, pad=0.04)
plt.xlabel("Predicted-label")
plt.ylabel("True-label")
plt.tight_layout()
plt.show()
```