
Determination of lattice structure of perovskites using ML

Kaling Vikram Singh
kalingvikram.singh@niser.ac.in
NISER

S Sunil Raja
sunilraja.s@niser.ac.in
NISER

Abstract

1 The crystalline structure is an important notion that controls various physical
2 properties of a matter and thus, it is important to correctly estimate the crystal
3 structure through various techniques. The various techniques are time consuming
4 and costly. Thus, one alternative to correctly estimate the crystal structure is using
5 a machine learning model. Perovskites are a special class of elements with ABO_3
6 type configuration. Several properties such as valency, atomic radii, band gap,
7 etc are used to predict the crystal structure. We have used various classification
8 techniques, after oversampling the data to get the required classification model.
9 We have utilized SVM, Light BGM and XgBoost to get our results. It was found
10 that LBGm gives the best accuracy of 94.32% followed by Xgboost with 94.13%
11 and weighted SVM with 92.03% accuracy. This accuracy is much higher than the
12 results noted in the reference paper.

13 1 Introduction

14 Perovskites are a class of elements that have similar crystal structure as the compound calcium
15 titanium, ABX_3 . These are used in various industries and its main applications include creating
16 solar panels that could be coated on various surfaces. These materials are lightweight and cheap and
17 are used in photovoltaic industry. These have high variations in A,B and can be found in various
18 structures such as cubic, monoclinic, orthorhombic, tetrahedral, hexagonal or rhombohedral. The
19 most significant reasons for the change in shapes include (i) displacement of the cation, (ii) distortion
20 of the octahedra and (iii) tilting of the octahedra. The displacement and distortions are instability
21 driven factors.

22 Good oxide ion conductivity, which is necessary in fuel cell applications, is an important property
23 of a cubic perovskite structure's reduced distortion. In a cubic perovskite, the 3D framework leads
24 to corner sharing of BO_6 octahedra and the A-cation is enclosed within 12 equidistant atoms. The
25 coordination number of O is 2 and is low since the A-O distance is almost 1.4 times the B-O bond
26 distance.

27 There is extensive work in progress that include Density functional Theory (DFT) calculations and
28 tough techniques such as X-ray diffraction (XRD) to get the the crystalline structure of materials.
29 These are power intensive and costly processes which can be reduced if Machine learning (ML)
30 models are used.

31 1.1 Related works

32 Santosh and Taher et al.^[1], provided the basis for this work, although the data they used included only
33 675 data entries and the data was biased towards orthorhombic structures. The methods used include
34 XgBoost, SVM, Light BGM, and Random Forest (RF) with accuracy of 74.8%, 76.6%, 80.3% and
35 62.8% respectively. The paper accounts for tolerance factor (τ), derived from radius of a,b, which has

36 been neglected in current paper and model is built with the available data-set. They reported the best
37 accuracy for RBF (which is explained using the density of states function which is also RBF) kernel
38 in SVM and Light GBM to give the best overall accuracy. Also, the sampled data was cut down to
39 just the ones that have each and every feature value.

40 Jarin et al.^[2] reported various models through which they have found the crystal structures without
41 oversampling with accuracy of 95% using genetic algorithm support vector regression. They also
42 utilized various Neural networks to achieve the best possible accuracy. They have not taken into
43 account the tolerance factor (τ) but also removed some features that are of less importance according
44 to the importance matrix.

45 2 Baselines

46 We have used simple yet fast and reliable models such as Light BGM, XgBoost, SVMs to obtain the
47 crystal structures. Boosting algorithms work on decision trees in which sequential tree growth using
48 gradient boosting improves performance by correcting categorization errors made by earlier trees.
49 This is fast and effective methods. Further, SVMs are simple machines that work on principle of
50 support vectors along with kernel functions. The corresponding kernel functions are used to get the
51 best maximum accuracies. Further, classifiers are used to get the importance matrix through which
52 we can see the corresponding importance of each features.

53 The main aim is to build a classification model that correctly classifies a perovskite. The workflow is
54 done as follows:

- 55 (i) Database collection
- 56 (ii) Feature selection and data-processing
- 57 (iii) Model selection
- 58 (iv) Hyper-parameter optimization
- 59 (v) Testing for accuracy.

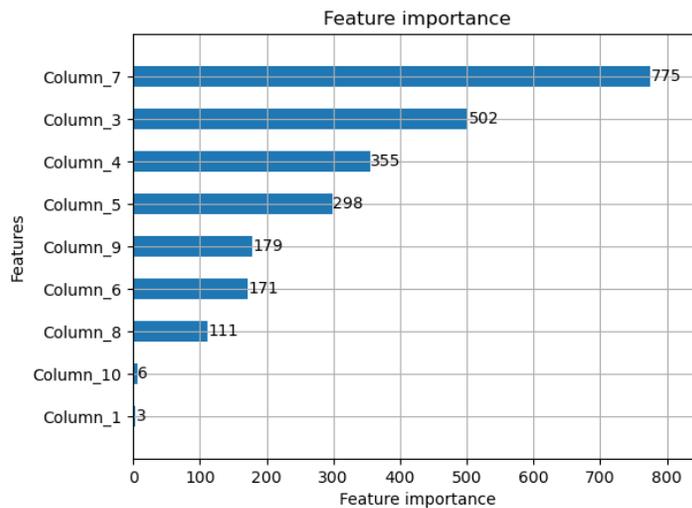


Figure 1: Feature importance

60 2.1 Model Environment

61 Python 3 was used for this project. Scikit-learn, Pandas, Imblearn, Numpy, etc librarians were used
62 for data processing and implementation of different models such as XgBoost, SMOTE , SVM , and
63 Light GBM.

64 2.2 Feature selection

65 Unlike others authors, we planned to give importance to each and every feature even if their importance
66 was low. The lgb.classifier was used to get the relative feature importance (fig 1). The column 7
67 representing the ionic radii of A was most important feature. The correlation matrix was plotted
68 using the data (fig 2). In case of weighted SVM, more importance was given to the instances which
69 are found in nature and which have balanced atoms. Only, the compound names were omitted from
70 the calculations as they did not show any importance. The bond angles and lattice edge length were
71 omitted as they are precursors of crystalline structure.

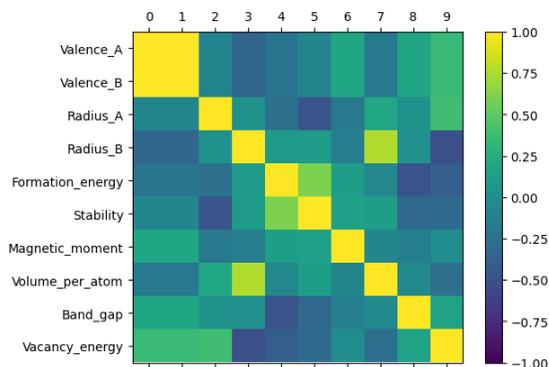


Figure 2: Correlation matrix

72 3 Experiment

73 GitHub repository: <https://github.com/Kalingvikramsingh/Machine-Learning>

74 3.1 Data-Pre processing

75 The data consists of 5329 datasets which include data that cannot be used for our model and thus, such
76 rows were deleted. The data was obtained through DFT calculations and includes various chemical
77 and physical properties of compounds. The dataset is highly imbalanced which can lead to overfitting
78 of particular type of dataset. SMOTE was used to equalize the number of instances of each label.
79 The features were selected beforehand. Since the values were not highly scattered, standardization
80 and normalization techniques were not utilized. The data was divided in training and test before
81 SMOTE to reduce bias in the data. After sampling, the new training data includes 5116 samples and
82 test contains 528 samples.

83 We have used Boosting methods on decision trees that can be used to predict results on the test data.
84 The boosting methods use sequential corrections in order to calculate the loss and get best possible
85 predictions using decision tree. In case of SVM, we have utilized different kernels to study the best
86 fit. Further, we used weighted SVM to calculate the predictions.

87 3.2 Light BGM

88 The parameters used were the height of the decision tree, learning rate. Validation dataset was
89 prepared from the training dataset which was used to get the hyper-parameter values. The optimum
90 height was found to be 7 and learning rate was 0.08. The losses were calculated on a logarithmic loss
91 function. Finally, the hyper parameters were used to get the best accuracy of 94.32% using LBGM.

92 3.3 XgBoost

93 The parameters used were height of decision tree, learning rate, and number of epochs. After hyper
94 parameter tuning, 20 epochs with 1.25 learning rate and tree height 2 was found to be the best fit.
95 The loss was again logarithmic. The accuracy found through this method was 94.13%

96 **3.4 SVM**

97 The parameters used were C(Penalty parameter) and kernels(Linear,RBF,Sigmoid and Polynomial).
98 After Hyper parameter tuning, RBF (fig 3) was chosen to be the kernel and the value of C as 209.5.
99 Upon further tuning the value of gamma for RBF was chosen to be 0.01.

100 **3.5 Weighted SVM**

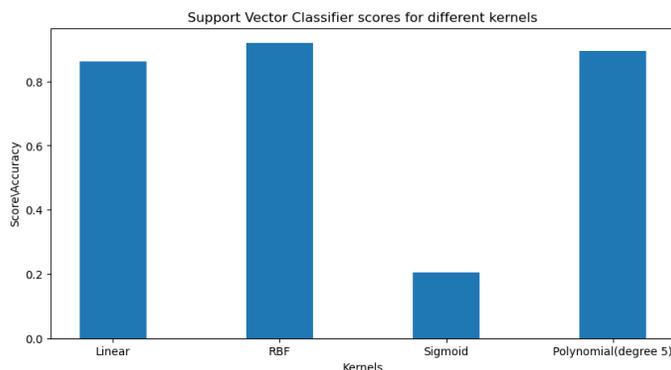


Figure 3: SVM kernels vs accuracy

101 The parameters used were C(Penalty Parameter), kernels and the weights of the instances. It was
102 decided that instances which occur in nature have higher weights compared to those which were
103 produced using SMOTE. The final parameter which were obtained are 208 for penalty parameter,
104 kernel as RBF(with gamma as 0.1) and the value of weights as 5,1,0.5.

Table 1: Accuracy of different models

Accuracy vs Model	
Model	Accuracy
SVM	91.33%
Weighted SVM	92.03%
Light GBM	94.32%
XGBoost	94.13%

105 **4 Conclusion**

106 Out of the sample data we analyzed, the overall accuracy was found to be 91.33% for SVM and
107 92.03% for weighted SVM,94.32% for LBGm and 94.13% for XgBoost. The model ,thus can be
108 used to classify the perovskites fairly accurately. The prediction accuracy was found to be as high as
109 94.32% which is fairly better than 80.3%^[1] than one reported by Santost et al. Thus, this model can
110 be fairly used for crystal structure predictions.

111 **Future Scope**

112 We utilized 4 models to predict the crystal structures of perovskites. Neural networks can be utilized
113 to get the accuracy better than benchmark of 95%^[2] reported in Jarin et al. Further, this model can be
114 used to identify crystal structures of different halides and not necessarily oxide perovskites.

115 **Dataset**

116 The dataset was a part of a research by Emery et al.^[3], licensed by MIT, and can be found at: https://figshare.com/articles/dataset/Wolverton_Oxides_Data/7250417?file=13354619
117

118 **References**

119 [1] Santosh Behara a, et al. “Crystal Structure Classification in ABO₃ Perovskites via Machine Learning.”
120 Computational Materials Science, Elsevier, 1 Dec. 2020.

121 [2] Jarin, S.; Yuan, Y.; Zhang, M.; Hu, M.; Rana, M.; Wang, S.; Knibbe, R. Predicting the Crystal Structure and
122 Lattice Parameters of the Perovskite Materials via Different Machine Learning Models Based on Basic Atom
123 Properties. Crystals 2022, 12, 1570. <https://doi.org/10.3390/cryst12111570>

124 [3] Emery, Antoine Wolverton, Chris. (2017). High-Throughput DFT calculations of formation en-
125 ergy, stability and oxygen vacancy formation energy of ABO₃ perovskites. Scientific Data. 4. 170153.
126 10.1038/sdata.2017.153.