

Band Gap ML part

Anonymous

Paper ID ***

Abstract. Keywords:

1 Introduction

2 Methods

2.1 Machine Learning

In this section we will describe our take on implementing a Machine learning algorithm for Band Gap prediction.

Previous work There were several notable attempts at utilizing Machine Learning to predict physical properties of various materials. Huang et al. reported prediction of band gap properties for ternary metal nitride compounds using ML approach based on calculated hybrid functionals of HSE and DFT PBE and [30]. In their study electronegativity, valence and covalent radius were selected as features for the training of the ML algorithm and prediction. In another study, high accuracy of the prediction was achieved with the ML algorithm trained on a dataset with 3 features of each elements in the compound: ionic radius, electronegativity and number of row associated with position of the element in the periodic table [31].

Data source Firstly, we screened open-access Materials Project and High Performance Computing Center Materials databases to identify stable ternary metal nitride compounds with optical properties that may be suitable for solar energy harvesting. The compounds with predicted band gap values of over 0.5 eV were considered. This is because the DFT method used in the databases for the calculation of the compounds was based on generalized gradient approximation (GGA+U). The latter is known to underestimate the band gap values of the materials relative to the experimentally measured ones [29].

Dataset preparation The algorithm was trained using a dataset of experimentally reported band gap values of about 360 ternary metal nitrides, oxides, sulfides and phosphides compounds.

A lack of data is a common problem in machine learning. This usually imposes limitations on various parts of machine learning pipeline that we will discuss below.

Feature extraction and processing Having low amounts of data also imposes certain limits on the amount of features that can be used without causing the curse of dimensionality [-1]. This suggests us to pick a number of features $d \ll N$, where $N = 360$ is the total number of data samples available in the training split.

The list of characteristic element-specific features used for the machine learning prediction includes detailed electronic configuration, maximum valence, atomic mas, electronegativity, atomic and covalent radius, ionization potential, electron affinity, period, group and block of the constituent elements.

Model Limited dataset size also affects the choice of the model and limits our choices depending on the model's complexity, since a very sophisticated model, e.g. deep neural network with many parameters, will tend to overfit when the number of training samples is too low.

We used support vector regression (SVR) with the nonlinear radial basis function (RBF) kernel as a machine learning model in our experiments. The model of choice was implemented with the use of the Scikit-learn framework and Python 3.x.

Results The ML prediction of the band gap values was used in combination with the data available in the open access databases to reduce the number of potential candidates for the subsequent theoretical simulations with hybrid functionals of HSE. In addition, the remaining compounds were sorted based on the value of effective electron and hole masses available in the following database [26][27][28].

-1) Trunk, G. V. (July 1979). "A Problem of Dimensionality: A Simple Example". IEEE Transactions on Pattern Analysis and Machine Intelligence. PAMI-1 (3): 306–307.

References

090	090
091	091
092	092
093	093
094	094
095	095
096	096
097	097
098	098
099	099
100	100
101	101
102	102
103	103
104	104
105	105
106	106
107	107
108	108
109	109
110	110
111	111
112	112
113	113
114	114
115	115
116	116
117	117
118	118
119	119
120	120
121	121
122	122
123	123
124	124
125	125
126	126
127	127
128	128
129	129
130	130
131	131
132	132
133	133
134	134