

Predicting molecular properties with machine learning

Md Naim Hassan Saykat

Petko Petkov

Project objective

- train machine learning models to predict molecular properties from the 3D structure of the molecules
- perform data exploration on QM9 dataset, convert the data into an appropriate format for machine learning models
- compare multiple regression models after cross-validation and hyperparameter tuning

Dataset

- QM9 (quantum chemical properties of approximately 134,000 stable small organic molecules, we're using 1000 random samples)
- includes computed geometric, energetic, electronic, and thermodynamic properties for each molecule
- contains 16 physical/chemical features (we are predicting 2 of them - heat capacity and isotropic polarizability) and Euclidean coordinates of the atoms

Dataset features

I.	Property	Unit	Description
1	tag	-	gdb9; string constant to ease extraction via grep
2	index	-	Consecutive, 1-based integer identifier of molecule
3	A	GHz	Rotational constant A
4	B	GHz	Rotational constant B
5	C	GHz	Rotational constant C
6	mu	Debye	Dipole moment
7	alpha	Bohr^3	Isotropic polarizability
8	homo	Hartree	Energy of Highest occupied molecular orbital (HOMO)
9	lumo	Hartree	Energy of Lowest occupied molecular orbital (LUMO)
10	gap	Hartree	Gap, difference between LUMO and HOMO
11	r2	Bohr^2	Electronic spatial extent
12	zpve	Hartree	Zero point vibrational energy
13	U0	Hartree	Internal energy at 0 K
14	U	Hartree	Internal energy at 298.15 K
15	H	Hartree	Enthalpy at 298.15 K
16	G	Hartree	Free energy at 298.15 K
17	Cv	cal/(mol K)	Heat capacity at 298.15 K

Preprocessing

- the dataset contains XYZ format files (coordinates of atoms and the molecular properties)
- extracting the atoms' coordinates and features from the XYZ files

Example XYZ format file:

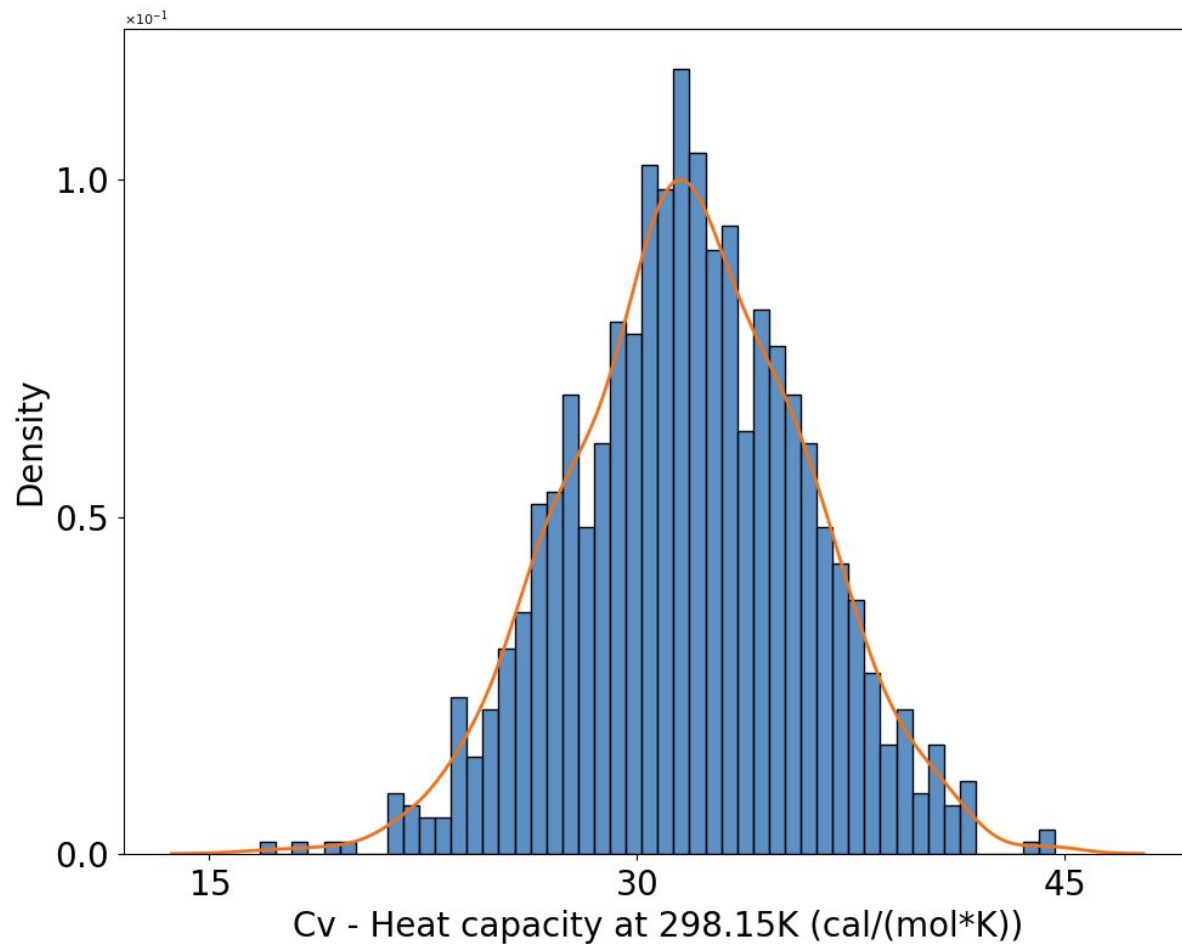
qm9-molecules > data > ≡ qm9_460.xyz

```
1  13
2  gdb 460 6.95109 3.605 2.72215 1.8411 49.49 -0.246 0.0261 0.2721 504.1131 0.109554 -286.510514 -286.504912 -286.503968 -286.539796 19.819
3  C -0.1875361387 1.5483985282 -0.0015224394 -0.396727
4  C 0.07229417 0.0511345788 0.0166833042 0.113356
5  C 0.9452326913 -0.4820790606 -1.1540812291 -0.192434
6  N 2.1643700861 -1.0161526732 -0.5211874978 -0.262577
7  C 1.9974828955 -0.8558646542 0.7241073334 0.122525
8  O 0.8641487372 -0.2873843153 1.1949494463 -0.230837
9  H 0.7540120674 2.1008437961 -0.0814498065 0.129414
10 H -0.8130444335 1.8082348907 -0.8619794537 0.122652
11 H -0.7027820073 1.8688477128 0.9080091098 0.130441
12 H -0.8700079159 -0.4978928383 0.1094746394 0.094257
13 H 0.4474514051 -1.2731170793 -1.7242119199 0.117336
14 H 1.2106716005 0.3121861877 -1.8613440684 0.113784
15 H 2.7073449821 -1.1405681236 1.4932744417 0.138808
```

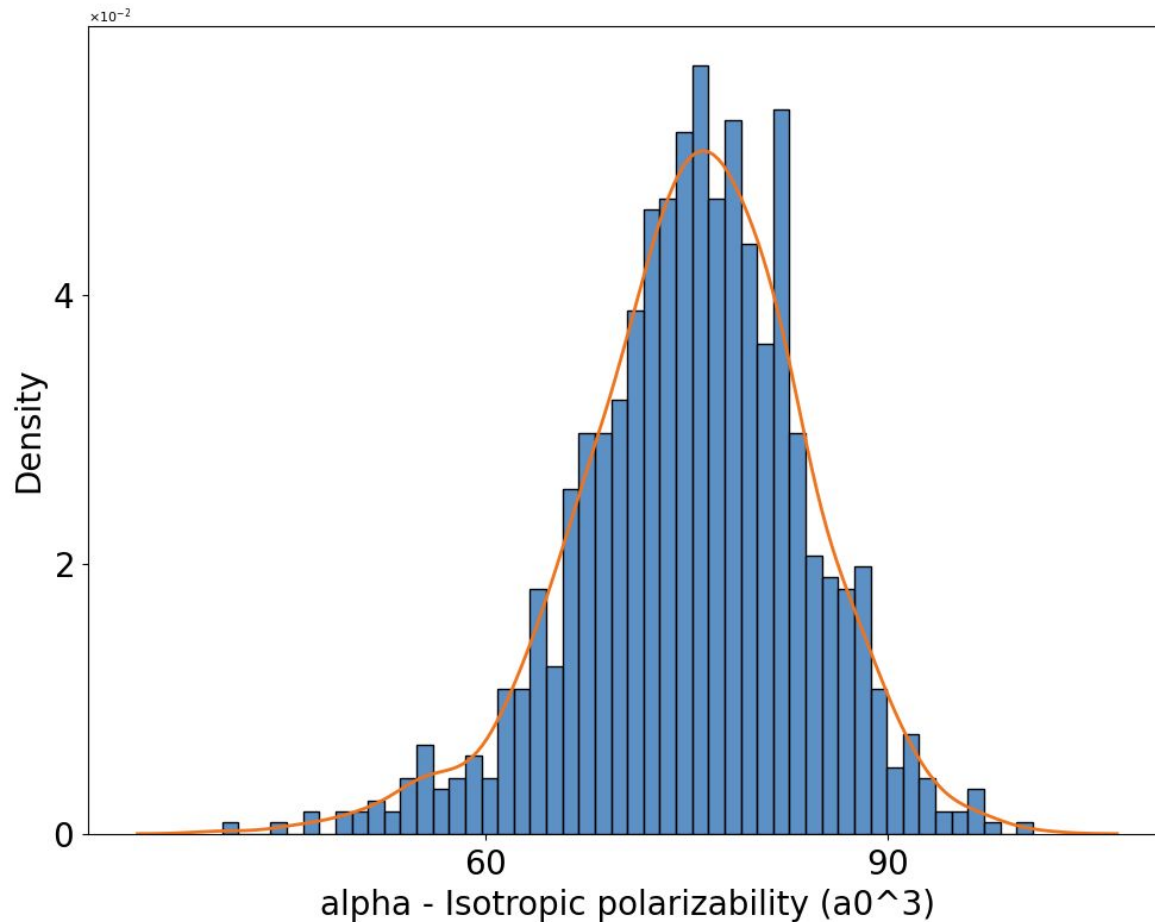
Data exploration

- visualize the data so we can get an idea of the targets we can predict
- make more informed decisions on model selection and evaluation metrics
- identify any potential outliers or skewness in the data that may require special handling during preprocessing

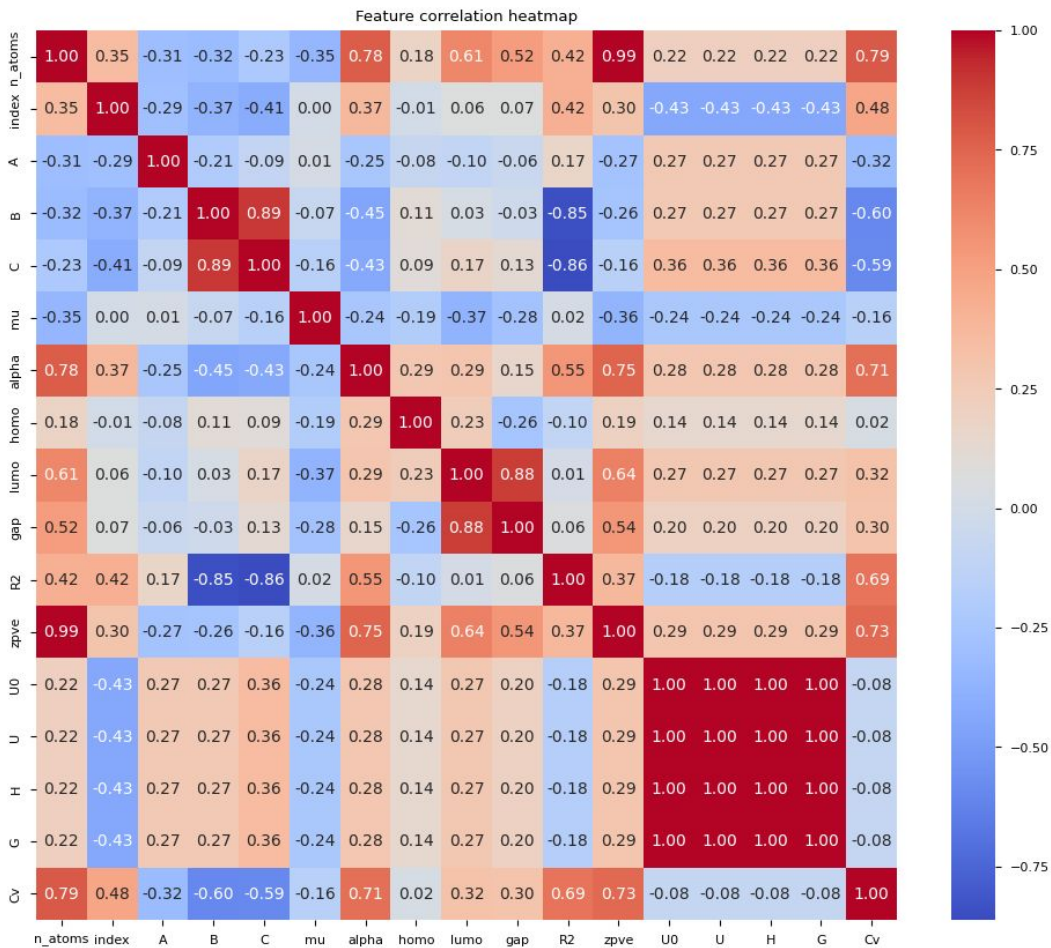
Heat capacity distribution



Isotropic polarizability distribution



Correlation between features



Calculate Coulomb matrices

- simple global descriptor which mimics the electrostatic interaction between nuclei
- N x N dimension where N is the number of atoms. The number of eigenvalues is also equal to N

$$M_{ij} = \begin{cases} \frac{Z_i Z_j}{\|\mathbf{R}_i - \mathbf{R}_j\|}, & i \neq j \\ 0.5 Z_i^{2.4}, & i = j \end{cases}$$

where

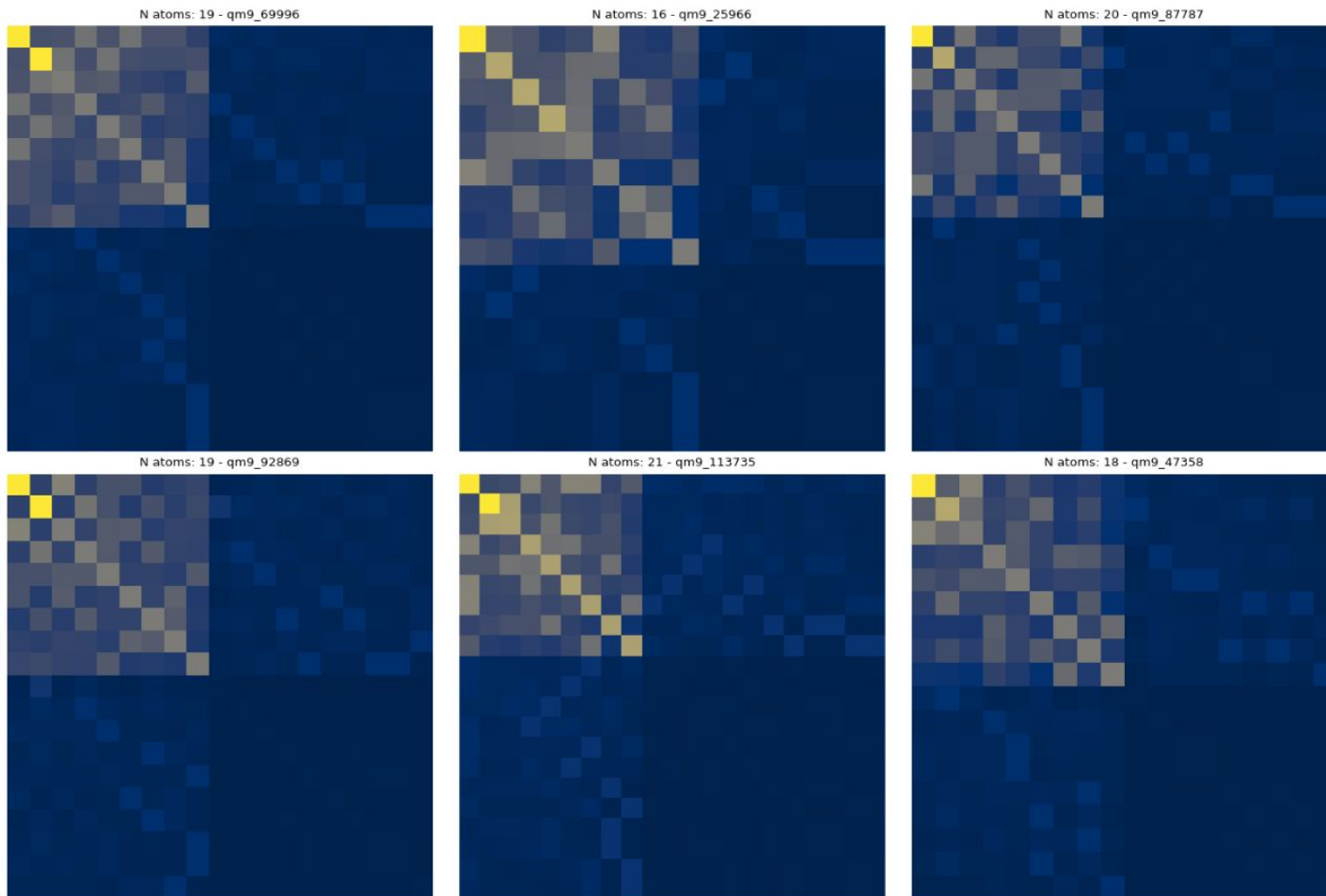
Z_i is the atomic number of atom i ,

\mathbf{R}_i is the position vector of atom i ,

i, j are indices for atoms,

$\|\mathbf{R}_i - \mathbf{R}_j\|$ is the Euclidean distance between atoms i and j .

Example Coulomb matrices



Calculation and padding of eigenvalues

- use eigenvalues of Coulomb matrix instead of the matrix itself as input to the models
- more efficient compressed version of the matrix
- apply padding to the eigenvalues to ensure that all eigenvalues have the same dimension

[illegible]

Models training

- compare multiple regression models - Linear regression, Random forest, SVR (Epsilon-support vector regression), K-nearest neighbors, XGBoost (Extreme gradient boosting) and Ridge regression
- perform hyperparameter optimization for each of the models in order to find the best hyperparameters
- use cross-validation to assess how well the models will generalize on unseen data

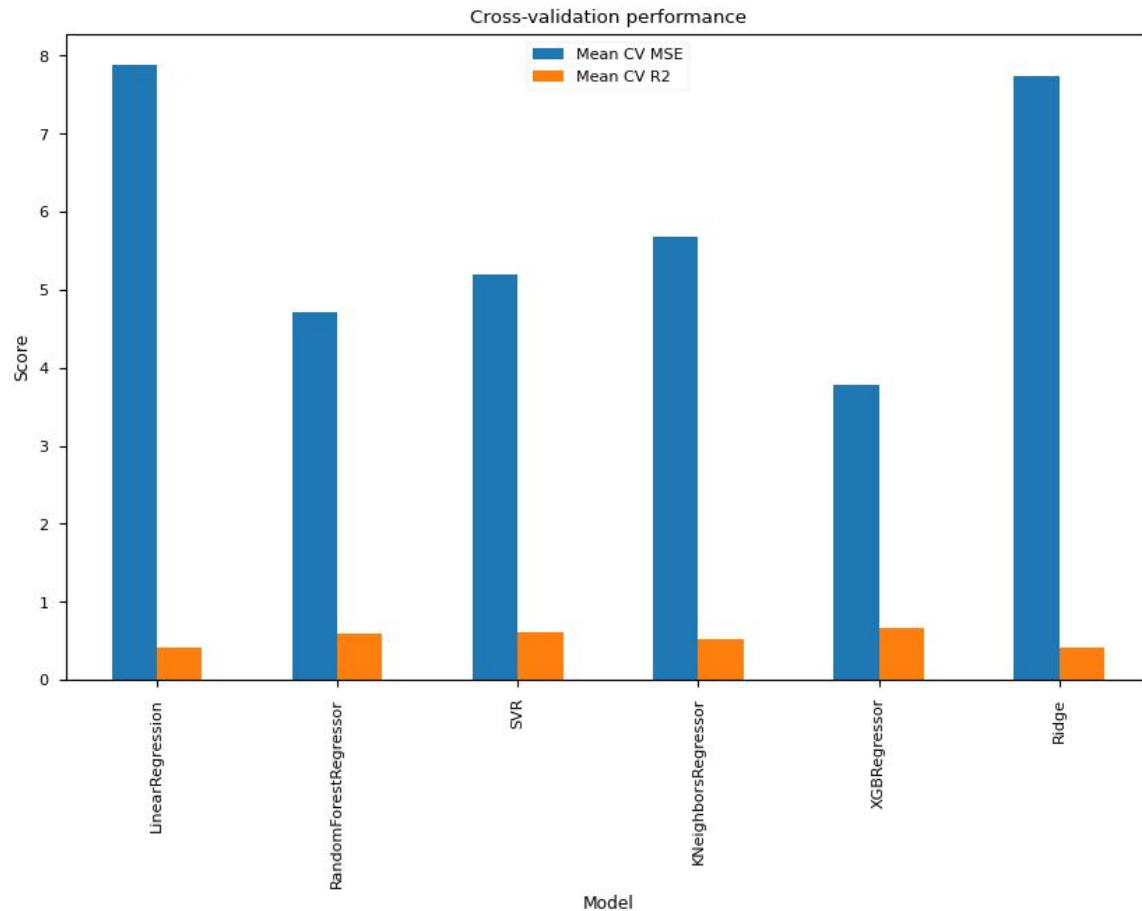
Mean squared error and R² loss functions

Compare models' performance with MSE and R² loss functions:

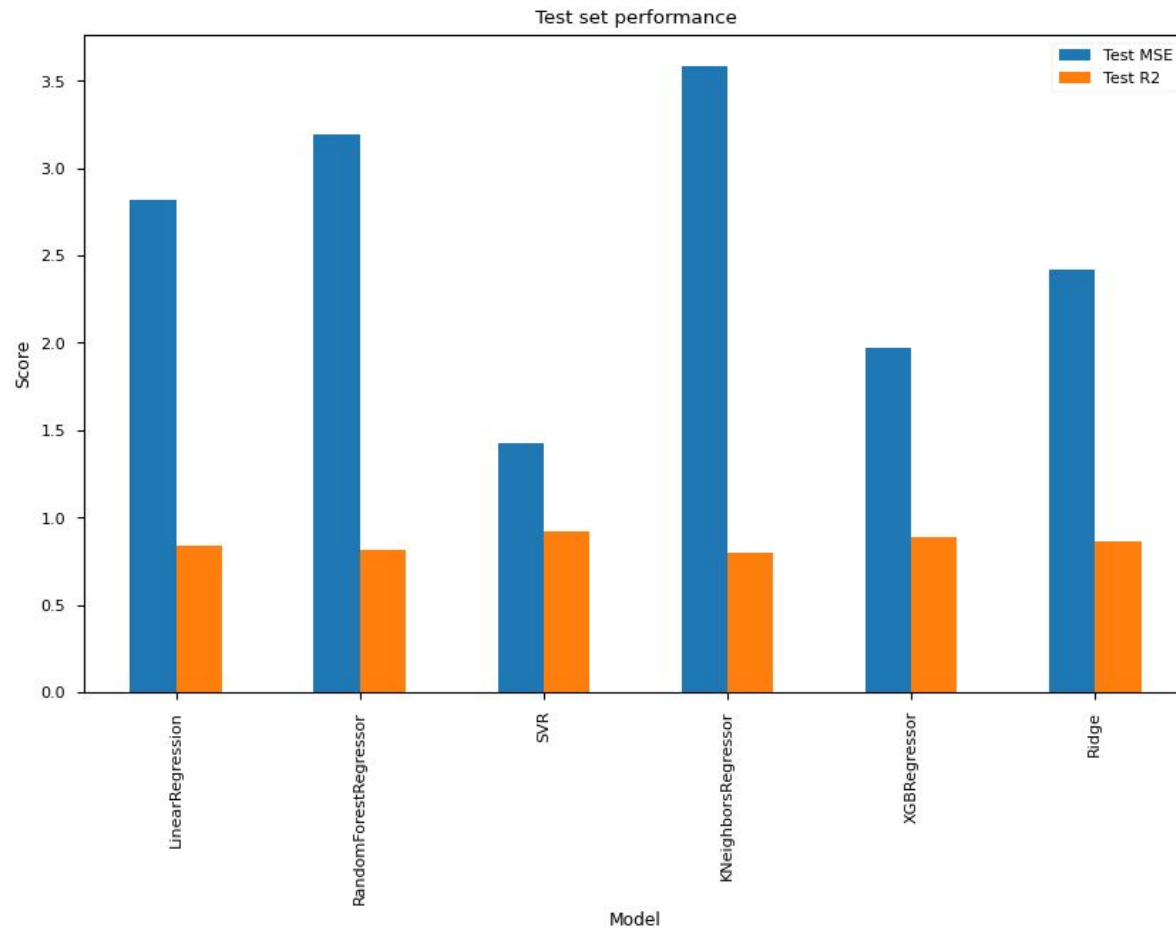
$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Heat capacity models cross-validation performance

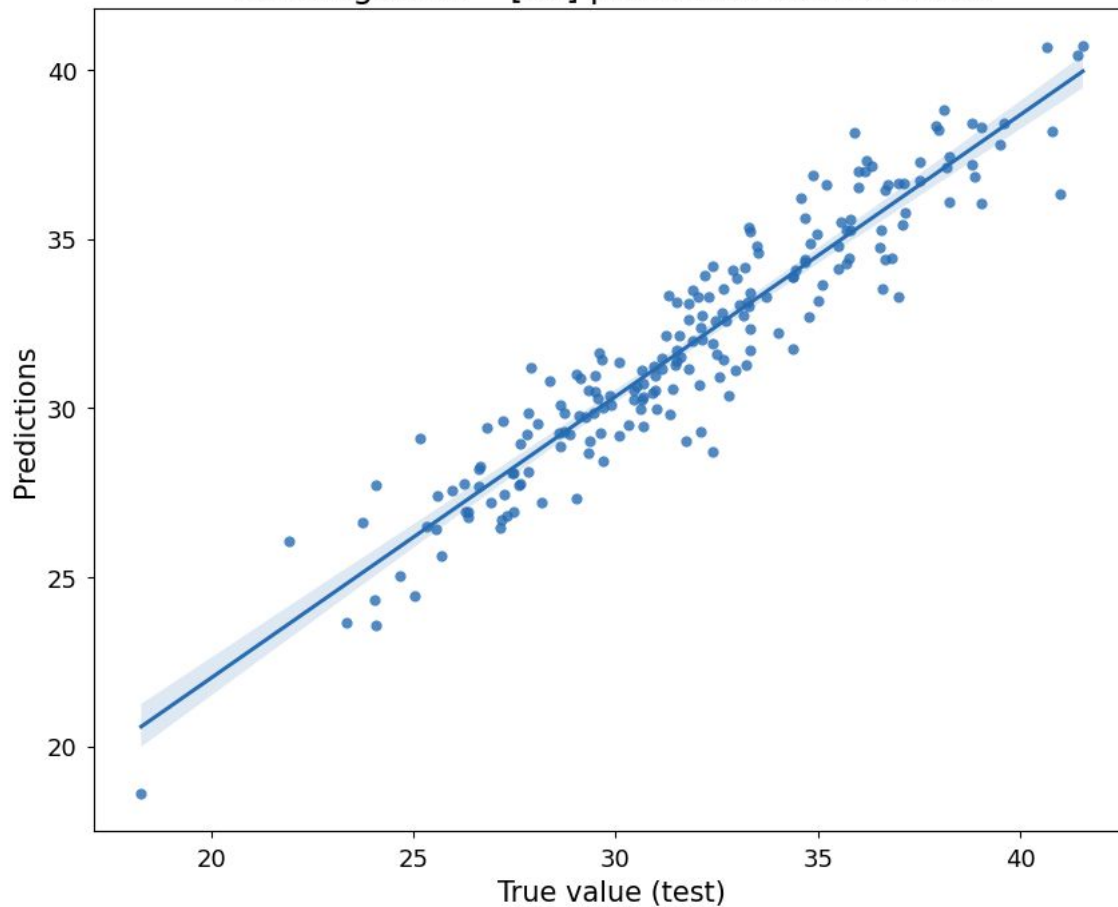


Heat capacity models test set performance

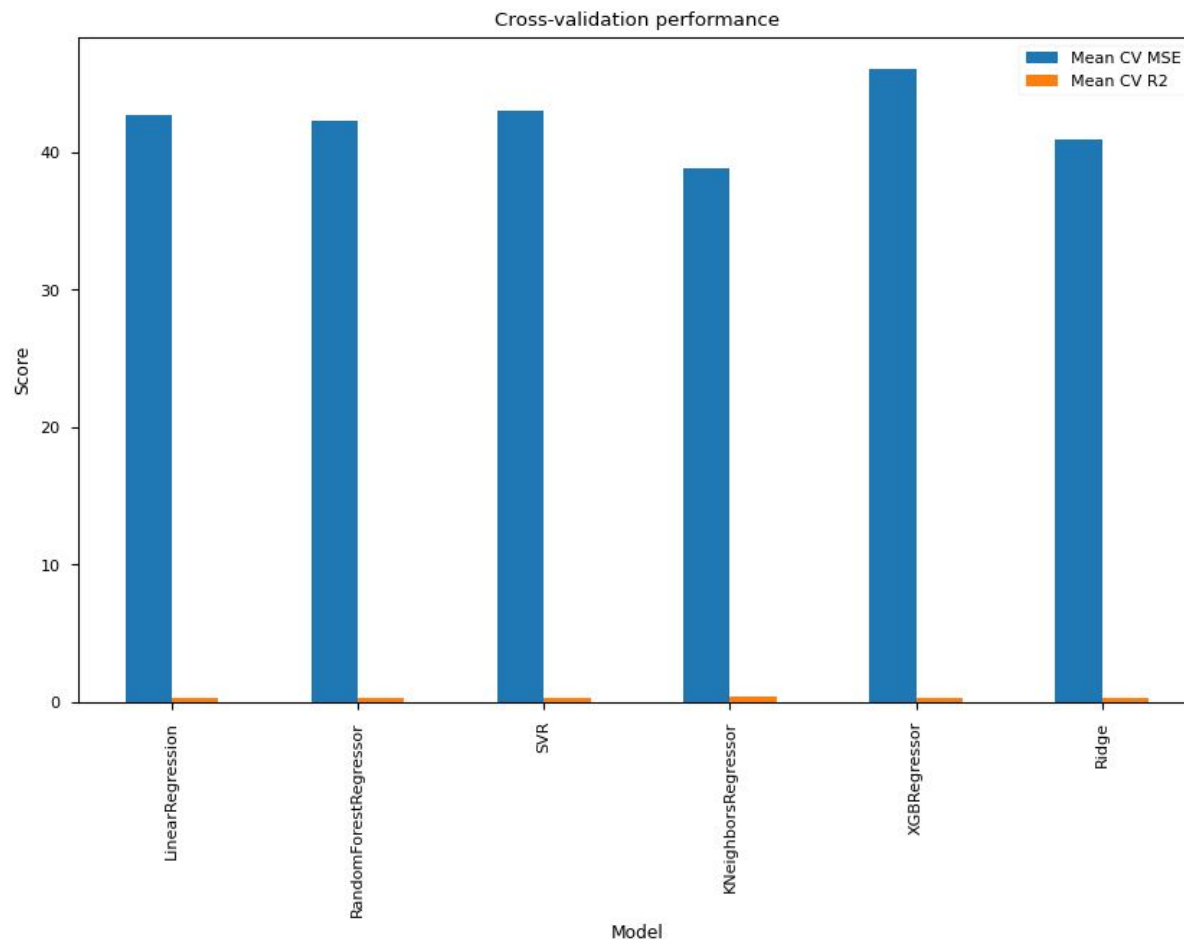


Heat capacity predictions

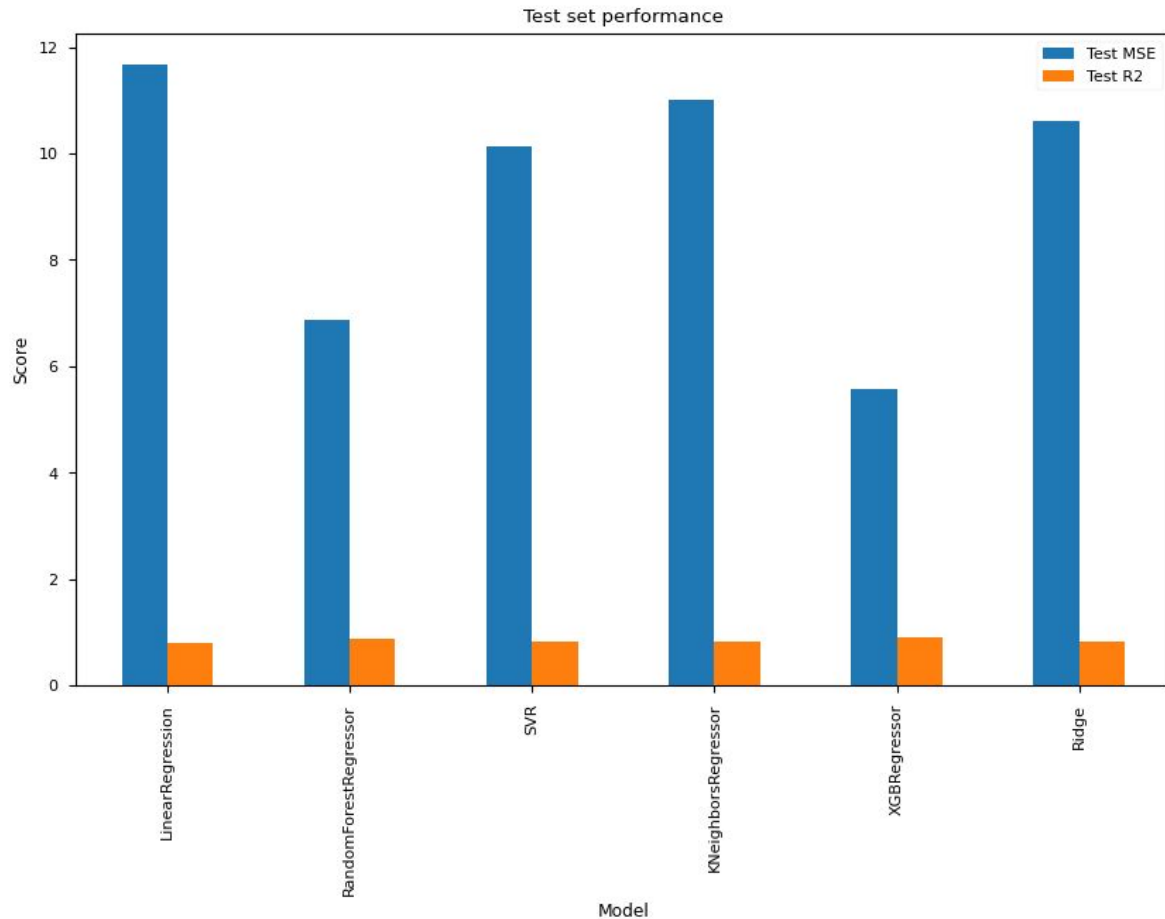
XGBRegressor - [Cv] prediction vs true value



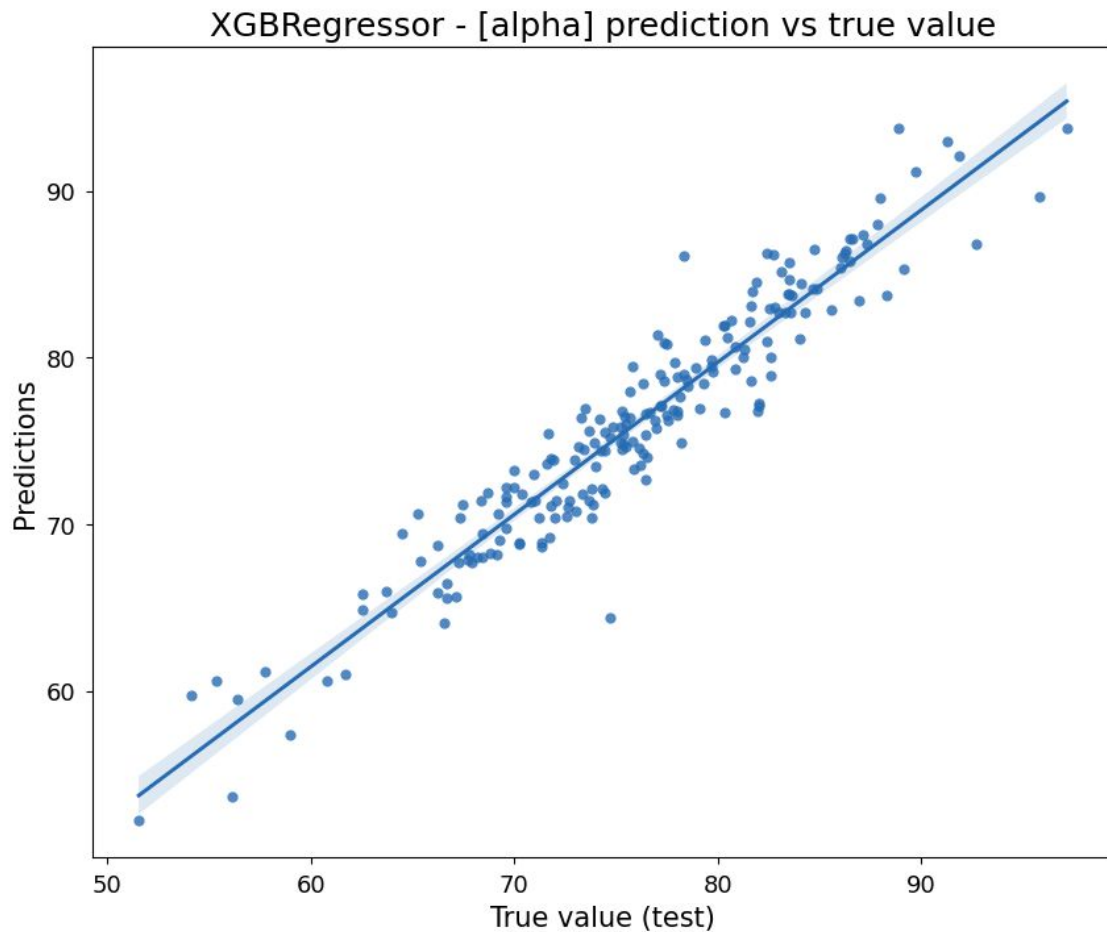
Isotropic polarizability cross-validation performance



Isotropic polarizability test set performance



Isotropic polarizability predictions



Conclusion

- data preprocessing and exploration on the QM9 dataset
- trained and compared multiple regression models on the eigenvalues of the Coulomb matrices to predict heat capacity and isotropic polarizability, performed hyperparameter tuning and cross-validation
- XGBoost is the best model for this task
- if we had more time: extend the dataset by adding molecules with more atoms and compare how these models would perform when the size of the molecules is significantly increased