

Predicting the critical temp of superconductors

Christopher Carman, Joakim Lajer, Nikolaj Andersson
NBI



All group members have contributed evenly to the project



Outline

- ① Motivation
- ② Approach
- ③ Data
- ④ Data Exploration
- ⑤ Feature selection
- ⑥ Methods
 - Keras Neural Network
 - Randomized Search Tree
 - Xgboost
 - Random Forest Tree
 - K Nearest Neighbors
- ⑦ Evaluation
- ⑧ References
- ⑨ Appendix



Motivation

- ① What is a superconductor
- ② Why is it important
- ③ No (successful) theory for predicting T_c

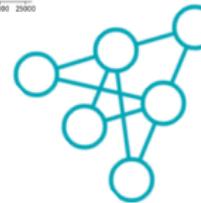
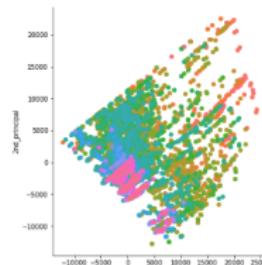


Approach

- Get data
- Kam Hamidieh - [github](#)
- Data exploration
- Feature selection
- Go crazy with ML



National Institute for Materials Science



Data

- ① Data From 2018 paper
- ② Based on Japanese SuperCon Database



Data Processing

- ① Features derived from elemental properties

Variable	units
Atomic mass	AMU
First Ionization Energy	KJ/mol
Atomic Radius	pm
Density	kg/m ³
Electron Affinity	KJ/mol
Fusion Heat	KJ/mol
Thermal Conductivity	(W/(m × k))
Valence	NA

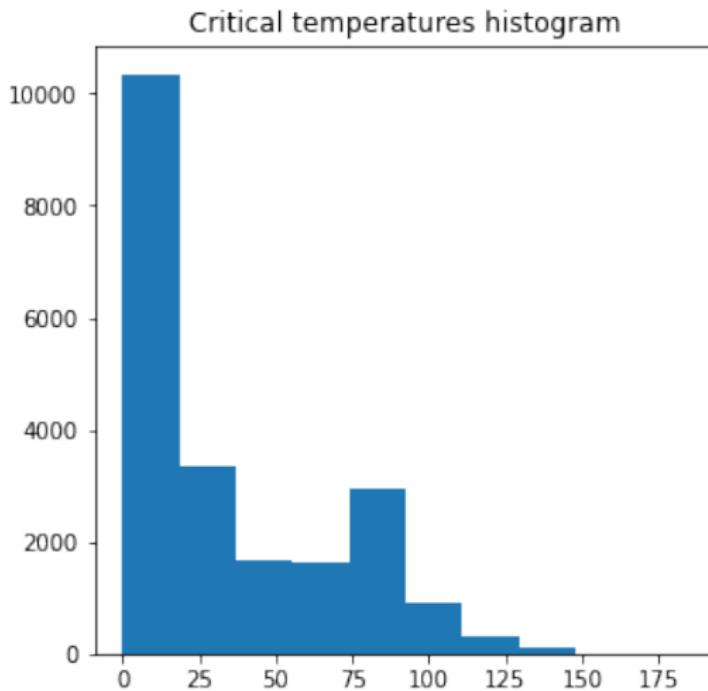


Data Features

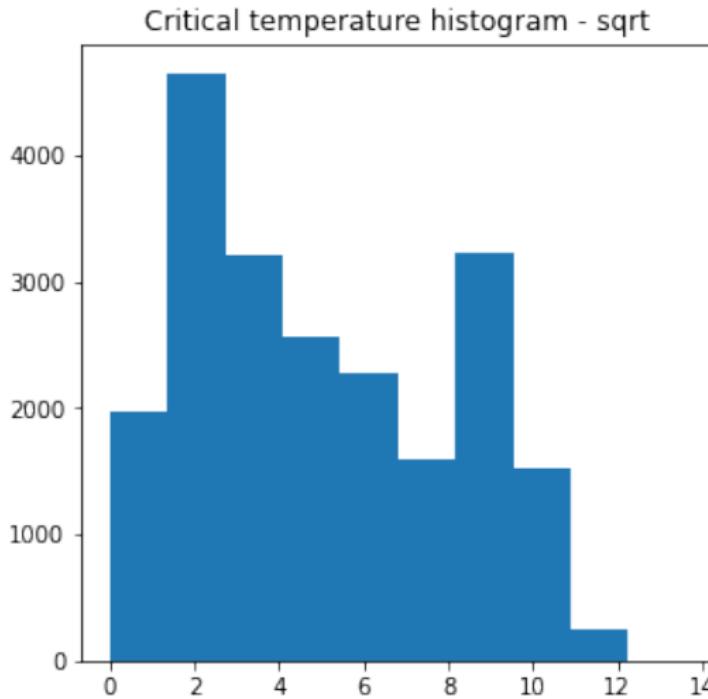
- ① Mean
- ② Weighted mean
- ③ Geometric mean
- ④ Weighted geometric mean
- ⑤ Entropy
- ⑥ Weighted entropy
- ⑦ Range
- ⑧ Weighted range
- ⑨ Standard deviation
- ⑩ Weighted standard deviation



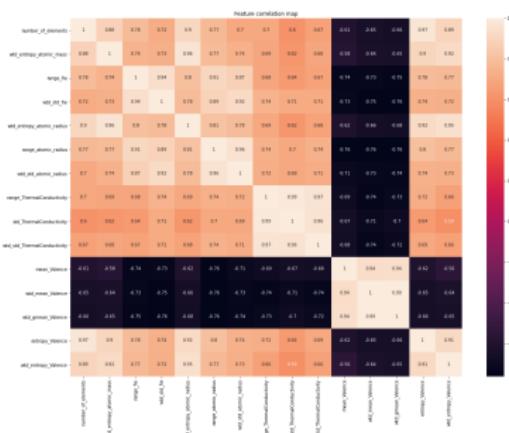
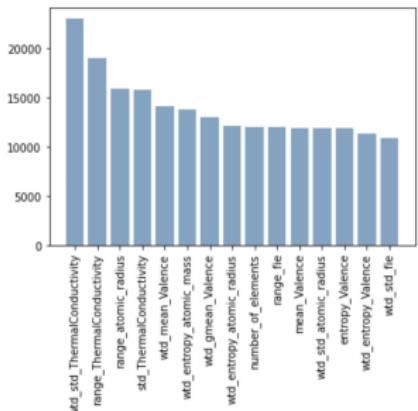
Data Exploration



Transform



Feature selection



Methods

- ① Keras Neural Network
- ② Randomized Search Tree
- ③ Random Forest Tree
- ④ Xgboost
- ⑤ K Nearest Neighbors



Keras Neural Network

① Model/Dimensions

Model: "sequential"

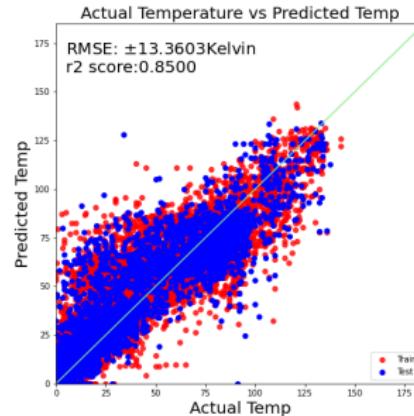
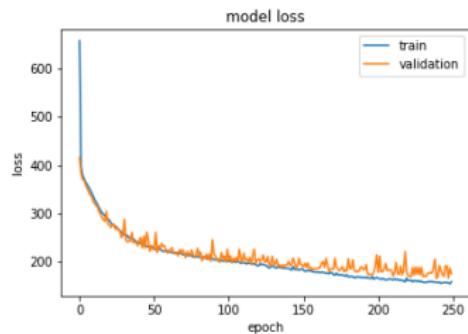
Layer (type)	Output Shape	Param #
dense (Dense)	(None, 24)	384
dense_1 (Dense)	(None, 30)	750
dense_2 (Dense)	(None, 24)	744
dense_3 (Dense)	(None, 1)	25

Total params: 1,903
Trainable params: 1,903
Non-trainable params: 0



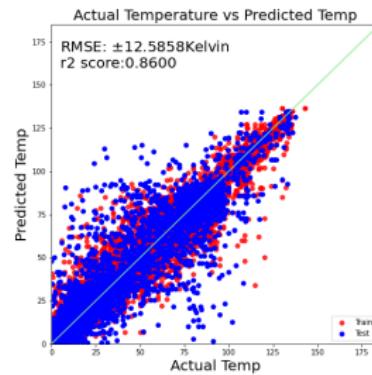
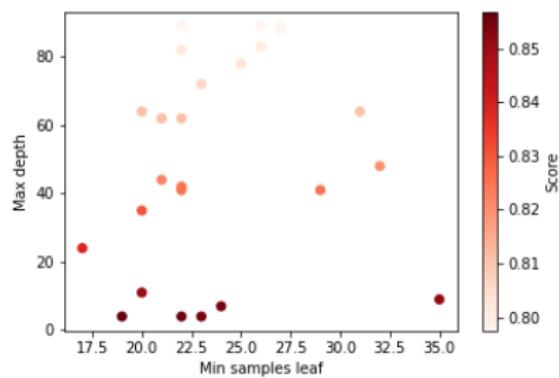
Keras Neural Network

① Efficiency/Results

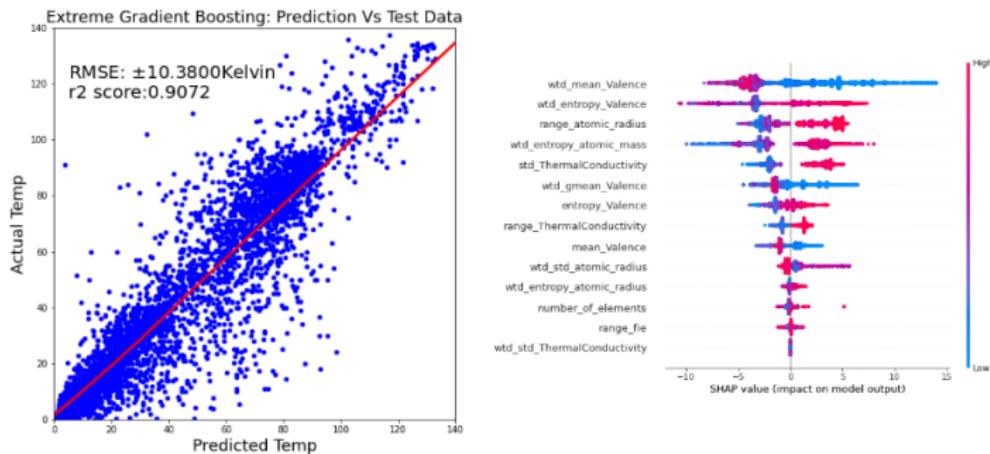


Randomized Search Tree

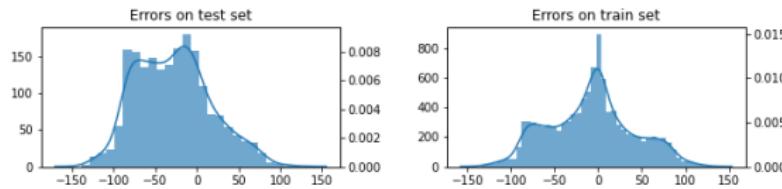
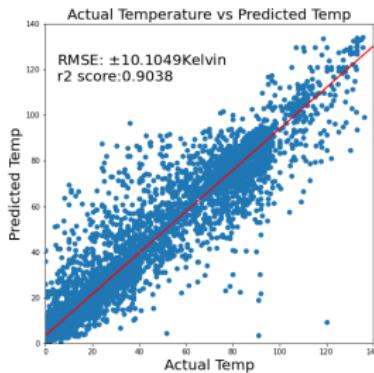
Top 5			
Max depth	Min samples leaf	Score	Rank
19	4	0.8569	1
22	4	0.8561	2
24	7	0.8556	3
23	4	0.8548	4
35	9	0.8517	5



Xgboost

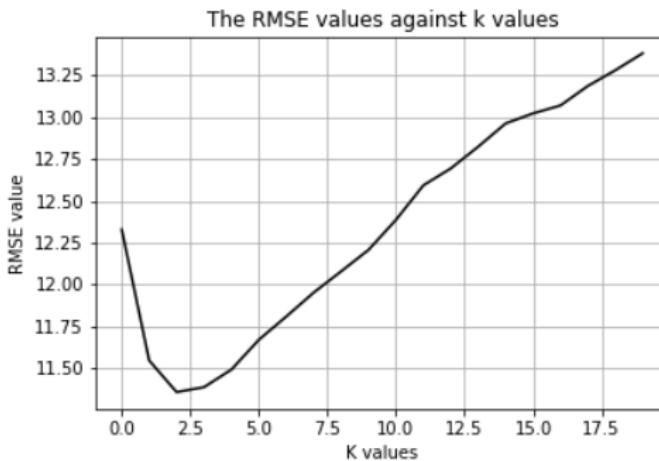


Random Forest Tree

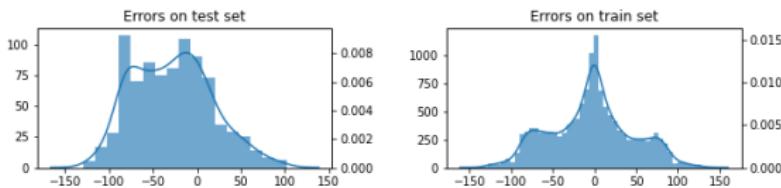
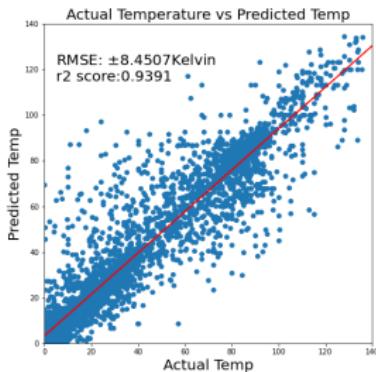


K Nearest Neighbors

① Finding the K value

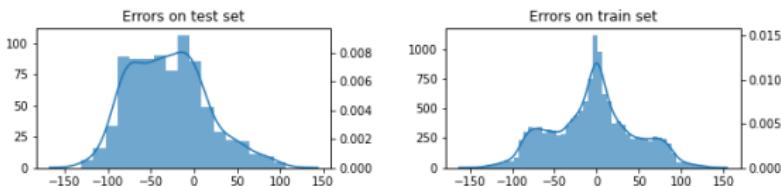
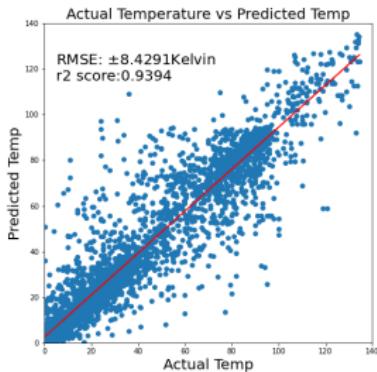


K Nearest Neighbors



K Nearest Neighbors

① Optimizing the algorithm



Evaluation

① Performances of methods

② Ranking

Method:	r2	RMSE	Speed
Keras Neural Network	0.8500	±13.3603	Medium
Randomized Search	0.8600	±12.5858	Fast
Xgboost	0.9072	±10.3800	Slow
Random Forest	0.9038	±10.1049	Fast
K Nearest Neighbors	0.9391	±8.4507	Medium
K Nearest Neighbors Optimized	0.9394	±8.4291	Semi-Slow



Thank you for your attention



References

Links to all Ipython notebooks (google colab) used in this project:

- ① Data exploration
- ② Feature selection
- ③ Keras Neural Network
- ④ Random Search Tree
- ⑤ XGBoost
- ⑥ Random Forrest Tree
- ⑦ kNN
- ⑧ kNN v2



Appendix - Feature selection

```
selector_f=SelectPercentile(f_regression,percentile=25)
selector_f.fit(X_var,y_var)

SelectPercentile(percentile=25,
               score_func=<function f_regression at 0x7f6ac79c9378>

for n,s in zip(X_var,selector_f.scores_):
    print('F-score: %3f\t%3' % (s,n))
```

```
import numpy as np
from sklearn.feature_selection import SelectKBest

# feature extraction
test = SelectKBest(f_regression, k=15)
fit = test.fit(X_var, y_var)

from itertools import cycle, islice
# Get the indices sorted by most important to least important
indices = np.argsort(fit.scores_)[::-1]

# To get top 15 feature names
features = []
for i in range(15):
    features.append(df.columns[indices[i]])

# plot
plt.figure()
plt.xticks(rotation=90)
plt.bar(features, fit.scores_[indices[range(15)]], color=(0.2, 0.4, 0.6, 0.6), align='center')
plt.savefig('feature_selection_rankings.png', bbox_inches = "tight")
plt.show()
```



Appendix - Keras Neural Network

```
model = Sequential()
model.add(Dense(24, input_dim=X_train.shape[1], kernel_initializer='normal', activation='relu'))
model.add(Dense(30, activation='relu'))
model.add(Dense(24, activation='relu'))
model.add(Dense(1, activation='relu'))
model.summary()

model.compile(loss='mse', optimizer='adam', metrics=['mse','mae'])

r2 = round(sm.r2_score(y_test, y_pred), 2)
RMSE = sqrt(mean_squared_error(y_test,y_pred))
```



Appendix - Random Search Tree

```
clf_DecisionTree = DecisionTreeRegressor(random_state=42)

parameters_RandomSearch = {'max_depth': poisson(25),
                           'min_samples_leaf': randint(1, 100)}

n_iter_search = 25
RandomSearch = RandomizedSearchCV(clf_DecisionTree,
                                   param_distributions=parameters_RandomSearch,
                                   n_iter=n_iter_search,
                                   cv=5,
                                   iid=True,
                                   return_train_score=True,
                                   random_state=42)

clf_RandomSearch = RandomSearch.best_estimator_

y_pred = clf_RandomSearch.predict(X_test)

r2 = round(sm.r2_score(y_test, y_pred), 2)
RMSE = sqrt(mean_squared_error(y_test,y_pred))
```



Appendix - XGBoost

```
DM_train = xgb.DMatrix(data = X_train,
                        label = y_train)

DM_test = xgb.DMatrix(data = X_test,
                      label = y_test)

gbm_param_grid = {
    'colsample_bytree': np.linspace(0.5, 0.9, 5),
    'n_estimators':[100, 200],
    'max_depth': [10, 15, 20, 25]
}

grid_mse = GridSearchCV(estimator = gbm, param_grid = gbm_param_grid,
                        scoring = 'neg_mean_squared_error', cv = 5, verbose = 1)
scores = cross_val_score(grid_mse, X, y, scoring='r2')

grid_mse.fit(X_train, y_train)

print("Best parameters found: ",grid_mse.best_params_)
print("Lowest RMSE found: ", np.sqrt(np.abs(grid_mse.best_score_)))
```



Appendix - Random Forrest Tree

```
scaler = MinMaxScaler()
X = df[variables]
X = pd.DataFrame(scaler.fit_transform(X), columns=X.columns)
y = df['critical_temp']

rf = RandomForestRegressor(n_estimators=200,max_features='log2')
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=40)
rf.fit(X_train,y_train)

y_pred_test = pd.Series(rf.predict(X_test))
y_pred_train = pd.Series(rf.predict(X_train))

min_rmse = round(np.sqrt(mean_squared_error(y_test,y_pred_test)),4)
r2 = r2_score(y_pred_test,y_test)
```



Appendix - kNN

```
scaler = StandardScaler()
X = df[variables]
X = pd.DataFrame(scaler.fit_transform(X), columns=X.columns)
y = df['critical_temp']

rmse_val = [] #to store rmse values for different k
for K in range(20):
    K = K+1
    knn = neighbors.KNeighborsRegressor(n_neighbors = K)

    knn.fit(X_train, y_train) #fit the model
    pred=knn.predict(X_test) #make prediction on test set
    error = sqrt(mean_squared_error(y_test,pred)) #calculate rmse
    rmse_val.append(error) #store rmse values
    print('RMSE value for k= ', K , 'is:', error)

#Create KNN regressor
knn = KNeighborsRegressor(n_neighbors=3)

#Train the model using the training sets
knn.fit(X_train, y_train)

#Predict the response for test dataset
y_pred = knn.predict(X_test)
```



Appendix - KNN v2

```
scaler = MinMaxScaler()
X = df[variables]
X2 = df[variables]
X = pd.DataFrame(scaler.fit_transform(X), columns=X.columns)
X_d = pd.DataFrame(scaler.fit_transform(X), columns=X.columns)
#print(X1.head())
y = df['critical_temp']

X = X.copy().apply(np.sqrt)
y = y.copy().apply(np.sqrt)

no_elements = df['number_of_elements']
X = pd.concat([no_elements.to_frame(), X], axis = 1)

from sklearn.model_selection import GridSearchCV
params = {'n_neighbors':[2,3,4,5],
          'leaf_size':np.arange(1,40)}

knn = neighbors.KNeighborsRegressor()

knn = GridSearchCV(knn, params, cv=10)
knn.fit(X_train,y_train)
knn.best_params_
```

