

Example 2: Fingerprint atomic crystal structures and learning their energy

Huan Tran

The main objective of this example is to get a small dataset of atomic crystal structures and their energy, fingerprint them, and develop some ML models using different learning algorithms.

1. Download data

The dataset contains 329 equilibrium structures of 13 different stoichiometries of Mg and Si, whose energy was computed using DFT. This dataset was reported in IT7. D. Huan, *Pressure-stabilized binary compounds of magnesium and silicon*, Phys. Rev. Materials **2**, 023803 (2018). It will be obtained from www.matsml.org. More information on the available datasets can be found at www.matsml.org as well.

```
In [1]: from matsml.data import Datasets
import os
import pandas as pd

# Load data
data = Datasets(ds_name='MgSi')
data_loader = Datasets(ds_name='MgSi')
data_loader.load_data()

# Have a look at the content
print(pd.read_csv(os.path.join(os.getcwd(), 'strids', name), 'summary.csv'))

matsML_v1.3.0
=====
Load requested dataset(s)
Data saved in crystals_MgSi
=====
file name      target
0  mg2si_struct_01.vasp  -8.5249797
1  mg2si_struct_02.vasp  -34.985707
2  mg2si_struct_03.vasp  -17.246812
3  mg2si_struct_04.vasp  -34.062642
4  mg2si_struct_05.vasp  -34.035175
...
324  mgSi_struct_30.vasp  -40.696471
325  mgSi_struct_31.vasp  -40.598719
326  mgSi_struct_32.vasp  -40.499177
327  mgSi_struct_33.vasp  -6.786034
328  mgSi_struct_34.vasp  -6.32384

[329 rows x 2 columns]
```

2. Fingerprint the obtained data

Two kinds of crystal fingerprints will be used in this example

- Ewald sum matrix [Fe, Faber, A. Lindmaa, O. Anatole von Lilienfeld, and R. Armiento. *Crystal structure representations for machine learning models of formation energies* Int. J. Quantum Chem., 115, 1094 (2015)] is an analogy to the Coulomb matrix for molecules, and its size also depends on the number of atoms of the structure. We used a similar projection on a set of Gaussian. Keyword for this fingerprint is **pesm_crystals**.
- Smooth Overlap of Atomic Positions (SOAP) [S. De, A. P. Bartók, G. Csányi, and M. Ceriotti, *Comparing molecules and solids across structural and alchemical space*, Phys. Chem. Chem. Phys. **18**, 13754 (2016)] is a more sophisticated fingerprint. Different from the Ewald sum matrix which is defined for the whole system, SOAP is defined for each atom. Herein, for simplicity, the atomic fingerprints are added up to make the fingerprint for the whole system. In some ML potential, the SOAP fingerprints are used in a different way, involving the "atomic energy". The keyword for SOAP in matsML is **soap_crystals**.

```
In [2]: import os
import pandas as pd
from matsml.fingerprint import Fingerprint

summary = os.path.join(os.getcwd(), 'crystals_MgSi/summary.csv')
data_loc = os.path.join(os.getcwd(), 'crystals_MgSi')
fp_dim = 20 # intended fingerprint dimensionality
verbosity = 0 # verbosity, 0 or 1

# Ewald sum matrix
data_params_pesm = {
    'summary': summary,
    'data_loc': data_loc,
    'fp_file': 'fp_crystals_MgSi_pesm.csv',
    'fp_type': 'pesm_crystals',
    'fp_dim': fp_dim,
    'verbosity': verbosity,
}

fp_pesm = Fingerprint(data_params_pesm)
fp_pesm.get_fingerprint()

# SOAP
data_params_soap = {
    'summary': summary,
    'data_loc': data_loc,
    'fp_file': 'fp_crystals_MgSi_soap.csv',
    'fp_type': 'soap_crystals',
    'fp_dim': fp_dim,
    'verbosity': verbosity,
}

fp_soap = Fingerprint(data_params_soap)
fp_soap.get_fingerprint()
```

Atomic structure fingerprinting /home/huantran/work_local/matsml/examples/ex3_crystals/crystals_MgSi/summary.csv

```
sv
data_loc      /home/huantran/work_local/matsml/examples/ex3_crystals/crystals_MgSi
fp_type       soap_crystals
fp_file       fp_crystals_MgSi_soap.csv
fp_dim        20
verbosity      0
Read input
num_structs   329
Computing Ewald sum Matrix
[=====] 100%
Projecting Ewald sum matrix to create fingerprints
[=====] 100%
Done fingerprinting, results saved in fp_crystals_MgSi_pesm.csv
Atomic structure fingerprinting
summary
/home/huantran/work_local/matsml/examples/ex3_crystals/crystals_MgSi/summary.csv
sv
data_loc      /home/huantran/work_local/matsml/examples/ex3_crystals/crystals_MgSi
fp_type       soap_crystals
fp_file       fp_crystals_MgSi_soap.csv
fp_dim        20
verbosity      0
Read input
num_structs   329
Computing SOAP fingerprint with D5cribe
[=====] 100%
Done fingerprinting, results saved in fp_crystals_MgSi_soap.csv
```

The fingerprinting step maybe a bit slow for a tutorial because we need to set **n_atoms_max=28**, which results in quite large Ewald sum matrices. A version of fingerprinted data can also be obtained in case you want to skip this step.

```
In [3]: from matsml.data import Datasets
import os
import pandas as pd

# Load data
data = Datasets(ds_name='fp_crystals_MgSi_soap', ds_pesm='fp_crystals_MgSi_pesm')
data_loader = Datasets(ds_name='fp_crystals_MgSi_soap', ds_pesm='fp_crystals_MgSi_pesm')
data_loader.load_data()

print(os.path.isfile('fp_crystals_MgSi_soap.csv.gz'))
print(os.path.isfile('fp_crystals_MgSi_pesm.csv.gz'))

Load requested dataset(s)
Data saved in fp_crystals_MgSi_soap.csv.gz
Data saved in fp_crystals_MgSi_pesm.csv.gz
True
True
```

3. Train several ML models with two fingerprint files just created

```
In [4]: # data parameters for learning

# ID column in the fingerprint data
y_col = ['target']
x_cols = ['x1', 'x2', 'x3', 'x4', 'x5', 'x6', 'x7', 'x8', 'x9', 'x10', 'x11', 'x12', 'x13', 'x14', 'x15', 'x16', 'x17', 'x18', 'x19', 'x20']
comment_cols = []
n_trains = 0.9 # 90% for training, 10% for validating
sampling = 'random' # method for train/test splitting
x_scaling = 'minmax'
y_scaling = 'minmax'
```

3a. Fully-connected NeuralNet

```
In [5]: from matsml.models import FCNN

# Model parameters
layers = [4, 4] # List of nodes in hidden layers
epochs = 2000 # Epochs
nfold_cv = 5 # Number of folds for cross validation
use_bias = True # Use bias term or not
model_file = 'model_fcnn.pkl' # Name of the model file to be created
verbosity = 0 # Verbosity, 0 or 1
batch_size = 32 # Default batch size
loss = 'mse'
activation = 'selu' # Options: "tanh", "relu", and more
optimizer = 'nadam' # Options: "Nadam", "Adam", and more

model_params = {
    'layers': layers,
    'activation': activation,
    'epochs': epochs,
    'nfold_cv': nfold_cv,
    'optimizer': optimizer,
    'use_bias': use_bias,
    'model_file': model_file,
    'loss': loss,
    'batch_size': batch_size,
    'verbosity': verbosity,
    'rmse_cv': False,
}

# PESM
model = FCNN(data_params=data_params_pesm, model_params=model_params)
model.train()
model.plot(pdf_output=False)
```

```
# SOAP
model = FCNN(data_params=data_params_soap, model_params=model_params)
model.train()
model.plot(pdf_output=False)
```

Checking parameters
model_file not ends with ".weights.h5", renamed

Learning fingerprinted/featured data
algorithm fully connected NeuralNet w/ TensorFlow
layers [4, 4]
activation selu
epochs 2000
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_crystals_MgSi_soap.csv
data size 329
training size 90.0 %
test size 10.0 %
x dimensionality 735
y dimensionality 1
y label(s) ['target']

Scaling x minmax
xscaler saved in xscaler.pkl

Scaling y minmax

Prepare train/test sets random

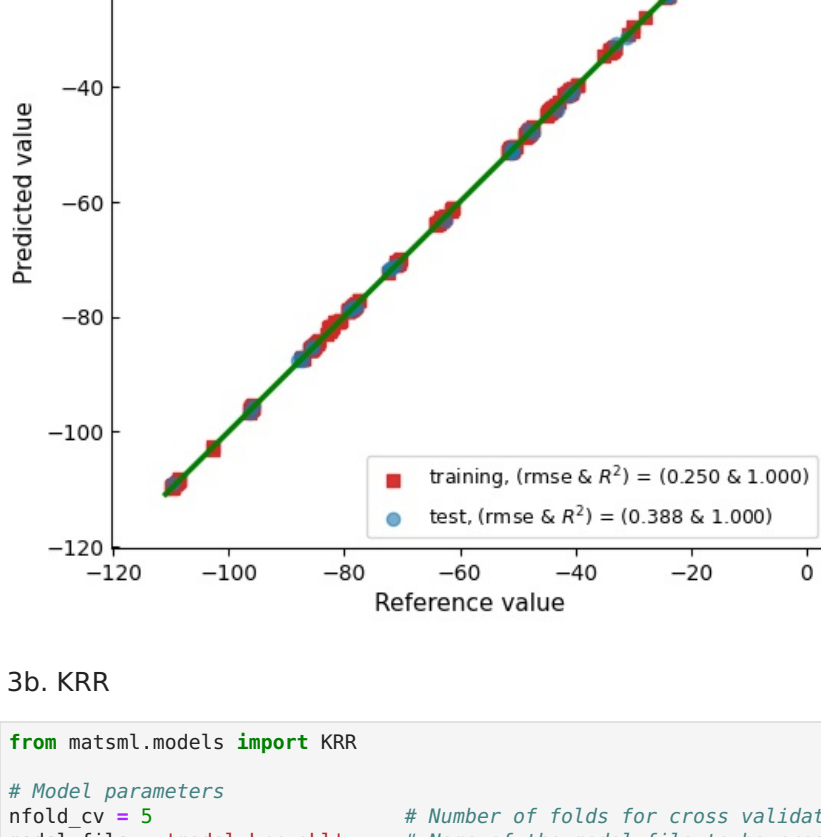
Building model FCNN

Training model w/ cross validation
8/8 [=====] - 0s 506us/step
2/2 [=====] - 0s 880us/step
cv_rmse_train,rmse_test,rmse_opt: 0 0.060827 0.088732 0.088732
2/2 [=====] - 0s 751us/step
cv_rmse_train,rmse_test,rmse_opt: 1 0.056679 0.084836 0.084836
8/8 [=====] - 0s 520us/step
2/2 [=====] - 0s 799us/step
cv_rmse_train,rmse_test,rmse_opt: 2 0.055376 0.071329 0.071329
8/8 [=====] - 0s 511us/step
2/2 [=====] - 0s 784us/step
cv_rmse_train,rmse_test,rmse_opt: 3 0.054338 0.072815 0.071329
8/8 [=====] - 0s 460us/step
2/2 [=====] - 0s 717us/step
cv_rmse_train,rmse_test,rmse_opt: 4 0.055881 0.068306 0.068306
Optimal ncv: 4 ; optimal NET saved

FCNN trained, now make predictions & invert scaling
10/10 [=====] - 0s 506us/step
unscaled y: minmax target 0.249655
r mse training target 0.249655
2/2 [=====] - 0s 1ms/step
unscaled y: minmax target 0.388191
r mse test target 0.388191

Predictions made & saved in "training.csv" & "test.csv"

Plot results in "training.csv" & "test.csv"
training, (rmse & R²) = (5.980 & 0.932)
test, (rmse & R²) = (17.524 & 0.531)
showing target



Checking parameters
model_file not ends with ".weights.h5", renamed

Learning fingerprinted/featured data
algorithm fully connected NeuralNet w/ TensorFlow
layers [4, 4]
activation selu
epochs 2000
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_crystals_MgSi_pesm.csv
data size 329
training size 90.0 %
test size 10.0 %
x dimensionality 735
y dimensionality 1
y label(s) ['target']

Scaling x minmax
xscaler saved in xscaler.pkl

Scaling y minmax

Prepare train/test sets random

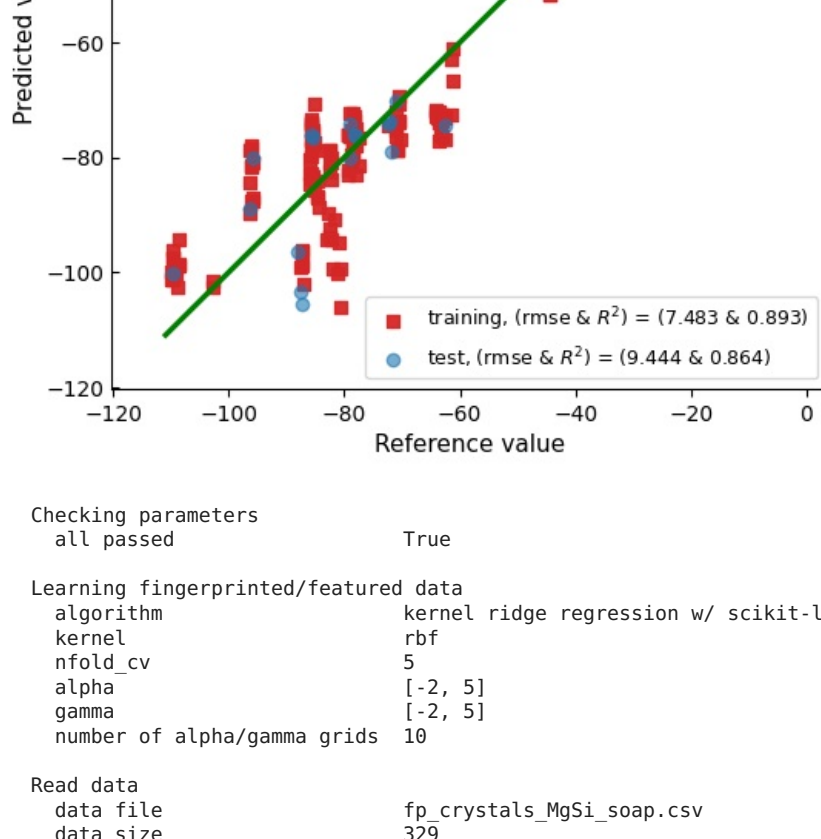
Building model FCNN

Training model w/ cross validation
8/8 [=====] - 0s 558us/step
2/2 [=====] - 0s 880us/step
cv_rmse_train,rmse_test,rmse_opt: 0 0.040827 0.086279 0.086279
8/8 [=====] - 0s 545us/step
2/2 [=====] - 0s 830us/step
cv_rmse_train,rmse_test,rmse_opt: 1 0.002263 0.044627 0.044627
8/8 [=====] - 0s 942us/step
2/2 [=====] - 0s 803us/step
cv_rmse_train,rmse_test,rmse_opt: 2 0.003888 0.044913 0.044627
8/8 [=====] - 0s 544us/step
2/2 [=====] - 0s 890us/step
cv_rmse_train,rmse_test,rmse_opt: 3 0.003208 0.044591 0.044591
8/8 [=====] - 0s 508us/step
2/2 [=====] - 0s 902us/step
cv_rmse_train,rmse_test,rmse_opt: 4 0.002134 0.003330 0.003330
Optimal ncv: 4 ; optimal NET saved

FCNN trained, now make predictions & invert scaling
10/10 [=====] - 0s 506us/step
unscaled y: minmax target 0.249655
r mse training target 0.249655
2/2 [=====] - 0s 1ms/step
unscaled y: minmax target 0.388191
r mse test target 0.388191

Predictions made & saved in "training.csv" & "test.csv"

Plot results in "training.csv" & "test.csv"
training, (rmse & R²) = (0.250 & 0.999)
test, (rmse & R²) = (0.388 & 1.000)
showing target



Checking parameters
model_file not ends with ".weights.h5", renamed

Learning fingerprinted/featured data
algorithm fully connected NeuralNet w/ TensorFlow
layers [4, 4]
activation selu
epochs 2000
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_crystals_MgSi_soap.csv
data size 329
training size 90.0 %
test size 10.0 %
x dimensionality 735
y dimensionality 1
y label(s) ['target']

Scaling x minmax
xscaler saved in xscaler.pkl

Scaling y minmax

Prepare train/test sets random

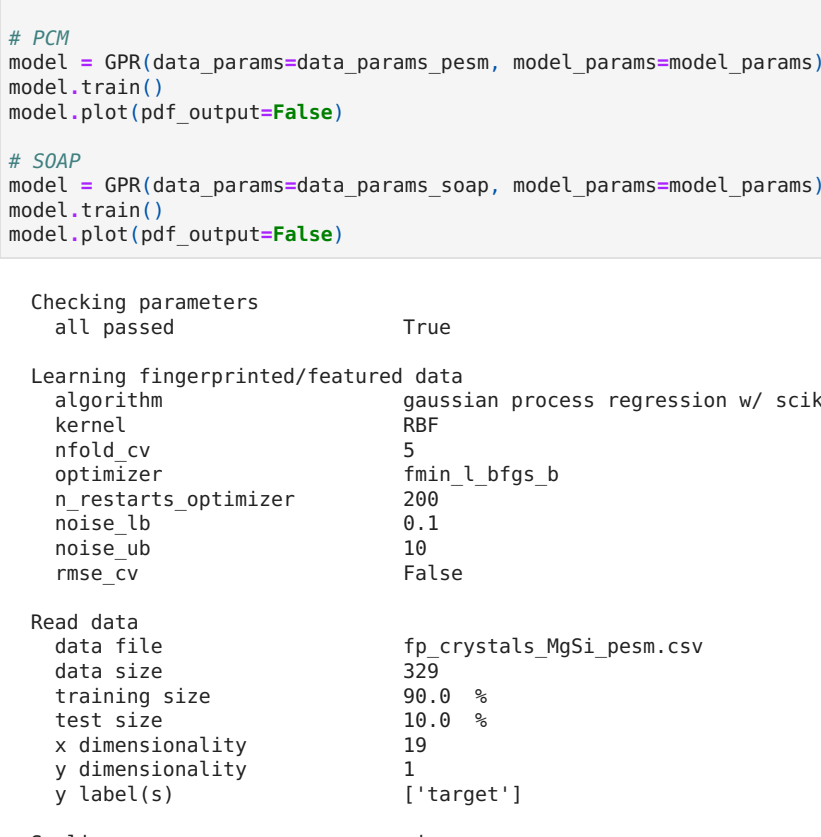
Building model FCNN

Training model w/ cross validation
8/8 [=====] - 0s 558us/step
2/2 [=====] - 0s 880us/step
cv_rmse_train,rmse_test,rmse_opt: 0 0.040827 0.086279 0.086279
8/8 [=====] - 0s 545us/step
2/2 [=====] - 0s 830us/step
cv_rmse_train,rmse_test,rmse_opt: 1 0.002263 0.044627 0.044627
8/8 [=====] - 0s 942us/step
2/2 [=====] - 0s 803us/step
cv_rmse_train,rmse_test,rmse_opt: 2 0.003888 0.044913 0.044627
8/8 [=====] - 0s 544us/step
2/2 [=====] - 0s 890us/step
cv_rmse_train,rmse_test,rmse_opt: 3 0.003208 0.044591 0.044591
8/8 [=====] - 0s 508us/step
2/2 [=====] - 0s 902us/step
cv_rmse_train,rmse_test,rmse_opt: 4 0.002134 0.003330 0.003330
Optimal ncv: 4 ; optimal NET saved

FCNN trained, now make predictions & invert scaling
10/10 [=====] - 0s 506us/step
unscaled y: minmax target 0.249655
r mse training target 0.249655
2/2 [=====] - 0s 1ms/step
unscaled y: minmax target 0.388191
r mse test target 0.388191

Predictions made & saved in "training.csv" & "test.csv"

Plot results in "training.csv" & "test.csv"
training, (rmse & R²) = (0.250 & 0.999)
test, (rmse & R²) = (0.388 & 1.000)
showing target



Checking parameters
model_file not ends with ".weights.h5", renamed

Learning fingerprinted/featured data
algorithm fully connected NeuralNet w/ TensorFlow
layers [4, 4]
activation selu
epochs 2000
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_crystals_MgSi_pesm.csv
data size 329
training size 90.0 %
test size 10.0 %
x dimensionality 735
y dimensionality 1
y label(s) ['target']

Scaling x minmax
xscaler saved in xscaler.pkl

Scaling y minmax

Prepare train/test sets random

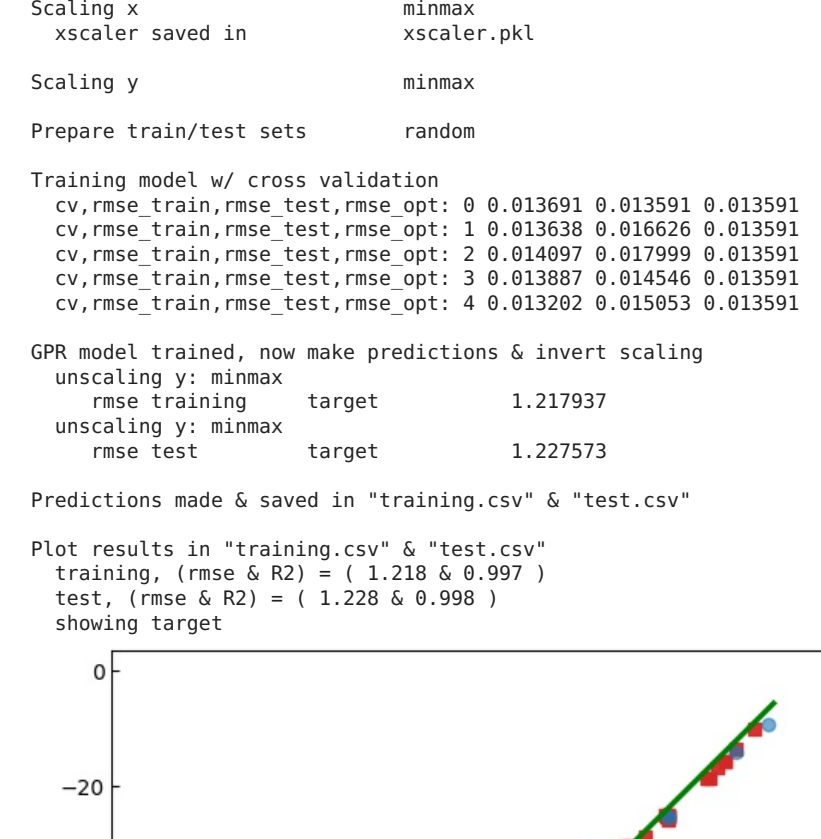
Building model FCNN

Training model w/ cross validation
8/8 [=====] - 0s 558us/step
2/2 [=====] - 0s 880us/step
cv_rmse_train,rmse_test,rmse_opt: 0 0.040827 0.086279 0.086279
8/8 [=====] - 0s 545us/step
2/2 [=====] - 0s 830us/step
cv_rmse_train,rmse_test,rmse_opt: 1 0.002263 0.044627 0.044627
8/8 [=====] - 0s 942us/step
2/2 [=====] - 0s 803us/step
cv_rmse_train,rmse_test,rmse_opt: 2 0.003888 0.044913 0.044627
8/8 [=====] - 0s 544us/step
2/2 [=====] - 0s 890us/step
cv_rmse_train,rmse_test,rmse_opt: 3 0.003208 0.044591 0.044591
8/8 [=====] - 0s 508us/step
2/2 [=====] - 0s 902us/step
cv_rmse_train,rmse_test,rmse_opt: 4 0.002134 0.003330 0.003330
Optimal ncv: 4 ; optimal NET saved

FCNN trained, now make predictions & invert scaling
10/10 [=====] - 0s 506us/step
unscaled y: minmax target 0.249655
r mse training target 0.249655
2/2 [=====] - 0s 1ms/step
unscaled y: minmax target 0.388191
r mse test target 0.388191

Predictions made & saved in "training.csv" & "test.csv"

Plot results in "training.csv" & "test.csv"
training, (rmse & R²) = (0.250 & 0.999)
test, (rmse & R²) = (0.388 & 1.000)
showing target



Checking parameters
model_file not ends with ".weights.h5", renamed

Learning fingerprinted/featured data
algorithm fully connected NeuralNet w/ TensorFlow
layers [4, 4]
activation selu
epochs 2000
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_crystals_MgSi_soap.csv
data size 329
training size 90.0 %
test size 10.0 %
x dimensionality 735
y dimensionality 1
y label(s) ['target']

Scaling x minmax
xscaler saved in xscaler.pkl

Scaling y minmax

Prepare train/test sets random

Building model FCNN

Training model w/ cross validation
8/8 [=====] - 0s 558us/step
2/2 [=====] - 0s 880us/step
cv_rmse_train,rmse_test,rmse_opt: 0 0.040827 0.086279 0.086279
8/8 [=====] - 0s 545us/step
2/2 [=====] - 0s 830us/step
cv_rmse_train,rmse_test,rmse_opt: 1 0.002263 0.044627 0.044627
8/8 [=====] - 0s 942us/step
2/2 [=====] - 0s 803us/step
cv_rmse_train,rmse_test,rmse_opt: 2 0.003888 0.044913 0.044627
8/8 [=====] - 0s 544us/step
2/2 [=====] - 0s 890us/step
cv_rmse_train,rmse_test,rmse_opt: 3 0.003208 0.044591 0.044591
8/8 [=====] - 0s 508us/step
2/2 [=====] - 0s 902us/step
cv_rmse_train,rmse_test,rmse_opt: 4 0.002134 0.003330 0.003330
Optimal ncv: 4 ; optimal NET saved

FCNN trained, now make predictions & invert scaling
10/10 [=====] - 0s 506us/step
unscaled y: minmax target 0.249655
r mse training target 0.249655
2/2 [=====] - 0s 1ms/step
unscaled y: minmax target 0.388191
r mse test target 0.388191

Predictions made & saved in "training.csv" & "test.csv"

Plot results in "training.csv" & "test.csv"
training, (rmse & R²) = (0.250 & 0.999)
test, (rmse & R²) = (0.388 & 1.000)
showing target

Checking parameters
model_file not ends with ".weights.h5", renamed

Learning fingerprinted/featured data
algorithm fully connected NeuralNet w/ TensorFlow
layers [4, 4]
activation selu
epochs 2000
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_crystals_MgSi_soap.csv
data size 329
training size 90.0 %
test size 10.0 %
x dimensionality 735
y dimensionality 1
y label(s) ['target']

Scaling x minmax
xscaler saved in xscaler.pkl

Scaling y minmax

Prepare train/test sets random

Building model FCNN

Training model w/ cross validation
8/8 [=====] - 0s 558us/step
2/2 [=====] - 0s 880us/step
cv_rmse_train,rmse_test,rmse_opt: 0 0.040827 0.086279 0.086279
8/8 [=====] - 0s 545us/step
2/2 [=====] - 0s 830us/step
cv_rmse_train,rmse_test,rmse_opt: 1 0.002263 0.044627 0.044627
8/8 [=====] - 0s 942us/step
2/2 [=====] - 0s 803us/step
cv_rmse_train,rmse_test,rmse_opt: 2 0.003888 0.044913 0.044627
8/8 [=====] - 0s 544us/step
2/2 [=====] - 0s 890us/step
cv_rmse_train,rmse_test,rmse_opt: 3 0.003208 0.044591 0.044591
8/8 [=====] - 0s 508us/step
2/2 [=====] - 0s 902us/step
cv_rmse_train,rmse_test,rmse_opt: 4 0.002134 0.003330 0.003330
Optimal ncv: 4 ; optimal NET saved

FCNN trained, now make predictions & invert scaling
10/10 [=====] - 0s 506us/step
unscaled y: minmax target 0.249655
r mse training target 0.249655
2/2 [=====] - 0s 1ms/step
unscaled y: minmax target 0.388191
r mse test target 0.388191

Predictions made & saved in "training.csv" & "test.csv"

Plot results in "training.csv" & "test.csv"
training, (rmse & R²) = (0.250 & 0.999)
test, (rmse & R²) = (0.388 & 1.000)
showing target

Checking parameters
model_file not ends with ".weights.h5", renamed

Learning fingerprinted/featured data
algorithm fully connected NeuralNet w/ TensorFlow
layers [4, 4]
activation selu
epochs 2000
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_crystals_MgSi_soap.csv
data size 329
training size 90.0 %
test size 10.0 %
x dimensionality 735
y dimensionality 1
y label(s) ['target']

Scaling x minmax
xscaler saved in xscaler.pkl

Scaling y minmax

Prepare train/test sets random

Building model FCNN

Training model w/ cross validation
8/8 [=====] - 0s 558us/step
2/2 [=====] - 0s 880us/step
cv_rmse_train,rmse_test,rmse_opt: 0 0.040827 0.086279 0.086279
8/8 [=====] - 0s 545us/step
2/2 [=====] - 0s 830us/step
cv_rmse_train,rmse_test,rmse_opt: 1 0.002263 0.044627 0.044627
8/8 [=====] - 0s 942us/step
2/2 [=====] - 0s 803us/step
cv_rmse_train,rmse_test,rmse_opt: 2 0.003888 0.044913 0.044627
8/8 [=====] - 0s 544us/step
2/2 [=====] - 0s 890us/step
cv_rmse_train,rmse_test,rmse_opt: 3 0.003208 0.044591 0.044591
8/8 [=====] - 0s 508us/step
2/2 [=====] - 0s 902us/step
cv_rmse_train,rmse_test,rmse_opt: 4 0.002134 0.003330 0.003330
Optimal ncv: 4 ; optimal NET saved

FCNN trained, now make predictions & invert scaling
10/10 [=====] - 0s 506us/step
unscaled y: minmax target 0.249655
r mse training target 0.249655
2/2 [=====] - 0s 1ms/step
unscaled y: minmax target 0.388191
r mse test target 0.388191

Predictions made & saved in "training.csv" & "test.csv"

Plot results in "training.csv" & "test.csv"
training, (rmse & R²) = (0.250 & 0.999)
test, (rmse & R²) = (0.388 & 1.000)
showing target

Checking parameters
model_file not ends with ".weights.h5", renamed

Learning fingerprinted/featured data
algorithm fully connected NeuralNet w/ TensorFlow
layers [4, 4]
activation selu
epochs 2000
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_crystals_MgSi_soap.csv
data size 329
training size 90.0 %
test size 10.0 %
x dimensionality 735
y dimensionality 1
y label(s) ['target']

Scaling x minmax
xscaler saved in xscaler.pkl

Scaling y minmax

Prepare train/test sets random

Building model FCNN

Training model w/ cross validation
8/8 [=====] - 0s 558us/step
2/2 [=====] - 0s 880us/step
cv