

Example 2: Fingerprint molecules data and learning energy

Huan Tran

The main objective of this example is to demonstrate a generic workflow of materials, involving (1) obtaining a small dataset of molecules and their energy, (2) fingerprint them, and (3) develop some ML models.

1. Download a dataset

The dataset contains 10,000 non-equilibrium structures of CH₃-NH-OH molecules, whose energy was computed using BigDFT package and HGH norm-conserving pseudopotentials. It is available at www.matssl.org.

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

# Load a dataset
dataset_name = 'molec3s_CH3NHOH'
data = Datasets(dataset_name=dataset_name)
data.load_dataset()

# Have a look at the content
summary_path = os.path.join(os.getcwd(), str(dataset_name), 'summary.csv')
print(pd.read_csv(summary_path))

matsML_v1.3.0
****
Load requested dataset(s)
Data saved in molec3s_CH3NHOH
file name      target
0  CH3NHOH_00001.xyz -940.288539
1  CH3NHOH_00002.xyz -940.580380
2  CH3NHOH_00003.xyz -940.184809
3  CH3NHOH_00004.xyz -940.466977
4  CH3NHOH_00005.xyz -940.579457
...
9994 CH3NHOH_09996.xyz -940.286083
9995 CH3NHOH_09997.xyz -940.744461
9996 CH3NHOH_09998.xyz -940.553979
9997 CH3NHOH_09999.xyz -940.650902
9998 CH3NHOH_10000.xyz -940.059079

[9999 rows x 2 columns]
```

2. Fingerprint the obtained data

Two kinds of fingerprints will be demonstrated here

- Coulomb matrix (CM) [M. Rupp, A. Tkatchenko, K.-R. Müller, and O. Anatole von Lilienfeld, *Fast and accurate modeling of molecular interaction energies with machine learning*, Phys. Rev. Lett., 108, 058301 (2012)] is perhaps one of the earliest fingerprints used in materials informatics. It was defined as an $N \times N$ matrix for a molecule of N atoms. The key advantage of CM is that it is invariant under rotations and translations, required to represent materials structure as a whole. However, its size depends on the molecule size, making it not directly usable for machine learning. Normally, the eigenvalues of these matrices are computed and sorted, and then zero padding is used to make fixed-size vectors. Here, we defined a projection of these Coulomb matrices onto a set of Gaussian functions, covering the entire range of the Coulomb matrix element values. The results are also a set of fixed-size fingerprints, which are ready for learning. Keyword for this fingerprint is **pcm_molecules**.
- 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. Keyword for this fingerprint is **soap_molecules**.

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

summary = os.path.join(os.getcwd(), 'molec3s_CH3NHOH/summary.csv')
data_loc = os.path.join(os.getcwd(), 'molec3s_CH3NHOH/')
fp_dim = 50 # Intended fingerprint dimensionality; the final number can be smaller
# verbosity = 0 or 1

# PCM
data_params_pcm = {
    'fp_type': 'pcm_molecules',
    'summary': summary,
    'data_loc': data_loc,
    'fp_file': 'fp_pcm.csv',
    'fp_dim': fp_dim,
    'verbosity': verbosity
}

fp_pcm = Fingerprint(data_params_pcm)
fp_pcm.get_fingerprint()

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

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

```
Atomic structure fingerprinting
summary /home/huantran/work_local/matssl/examples/ex2_molecules/molec3s_CH3NHOH/summary.csv
v
data_loc /home/huantran/work_local/matssl/examples/ex2_molecules/molec3s_CH3NHOH/
fp_type pcm_molecules
fp_file fp_pcm.csv
fp_dim 50
verbosity 0
Read input
num_structs 9999
Computing Coulomb matrix [=====] 100%
Projecting Coulomb matrix to create fingerprints [=====] 100%
Done fingerprinting, results saved in fp_pcm.csv
Atomic structure fingerprinting
summary /home/huantran/work_local/matssl/examples/ex2_molecules/molec3s_CH3NHOH/summary.csv
v
data_loc /home/huantran/work_local/matssl/examples/ex2_molecules/molec3s_CH3NHOH/
fp_type soap_molecules
fp_file fp_soap.csv
fp_dim 50
verbosity 0
Read input
num_structs 9999
Computing SOAP fingerprint with D5cribe [=====] 100%
Done fingerprinting, results saved in fp_soap.csv
```

The fingerprinting step is slow. A version of fingerprinted data can also be obtained in case you want to skip this step. Pandas can read gzip files for no need to unzip them.

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

# Load data
data = Datasets(ds1='fp_molecules_CH3NHOH_pcm', ds2='fp_molecules_CH3NHOH_soap')
data.load_dataset()

print(os.path.isfile('fp_molecules_CH3NHOH_pcm.csv.gz'))
print(os.path.isfile('fp_molecules_CH3NHOH_soap.csv.gz'))

Load requested dataset(s)
Data saved in fp_molecules_CH3NHOH_pcm.csv.gz
Data label in fp_molecules_CH3NHOH_soap.csv.gz
True
True
```

3. Train some ML models with "fp_pcm.csv" and "fp_soap.csv" just created

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

id_col = ['id'] # column for data ID
y_col = ['target'] # columns for (one or more) target properties
comment_cols = [] # comment columns, anything not counted into ID, fingerprints, and target
n_train = 0.8 # 80% for training, 20% for validating
sampling = 'random' # method for training & splitting
x_scaling = 'minmax' # method for x scaling
y_scaling = 'minmax' # method for y scaling

# Dict of data parameters
data_params_pcm = {
    'data_file': 'fp_molecules_CH3NHOH_pcm.csv.gz',
    'id_col': id_col,
    'y_cols': y_col,
    'comment_cols': comment_cols,
    'y_scaling': y_scaling,
    'x_scaling': x_scaling,
    'sampling': sampling,
    'n_train': n_train
}

data_params_soap = {
    'data_file': 'fp_molecules_CH3NHOH_soap.csv.gz',
    'id_col': id_col,
    'y_cols': y_col,
    'comment_cols': comment_cols,
    'y_scaling': y_scaling,
    'x_scaling': x_scaling,
    'sampling': sampling,
    'n_train': n_train
}
```

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

# Model parameters
layers = [8, 8, 0] # list of nodes in hidden layers
epochs = 200 # Epochs
nfold_cv = 5 # Number of folds for cross validation
use_bias = True # Use bias or not
model_file = 'model_nn.pkl' # Name of the model file to be created
verbosity = 0 # Verbosity, 0 or 1
batch_size = 32 # Default = 32
loss = 'mse'
activation = 'selu' # Options: "tanh", "relu", and more
optimizer = 'nadam' # Options: "nadam", "Adam", and more

# Dict of model parameters
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
}
```

```
# PCM
model = FCNN(data_params=data_params_pcm, 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 [8, 8, 8]
activation selu
epochs 200
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_molecules_CH3NHOH_pcm.csv.gz
data size 9999
training size 80.0 %
test size 20.0 %
x dimensionality 38
y dimensionality 1
y label(s) ['target']

Scaling x minmax
xscaler saved in xscaler.pkl

Scaling y minmax

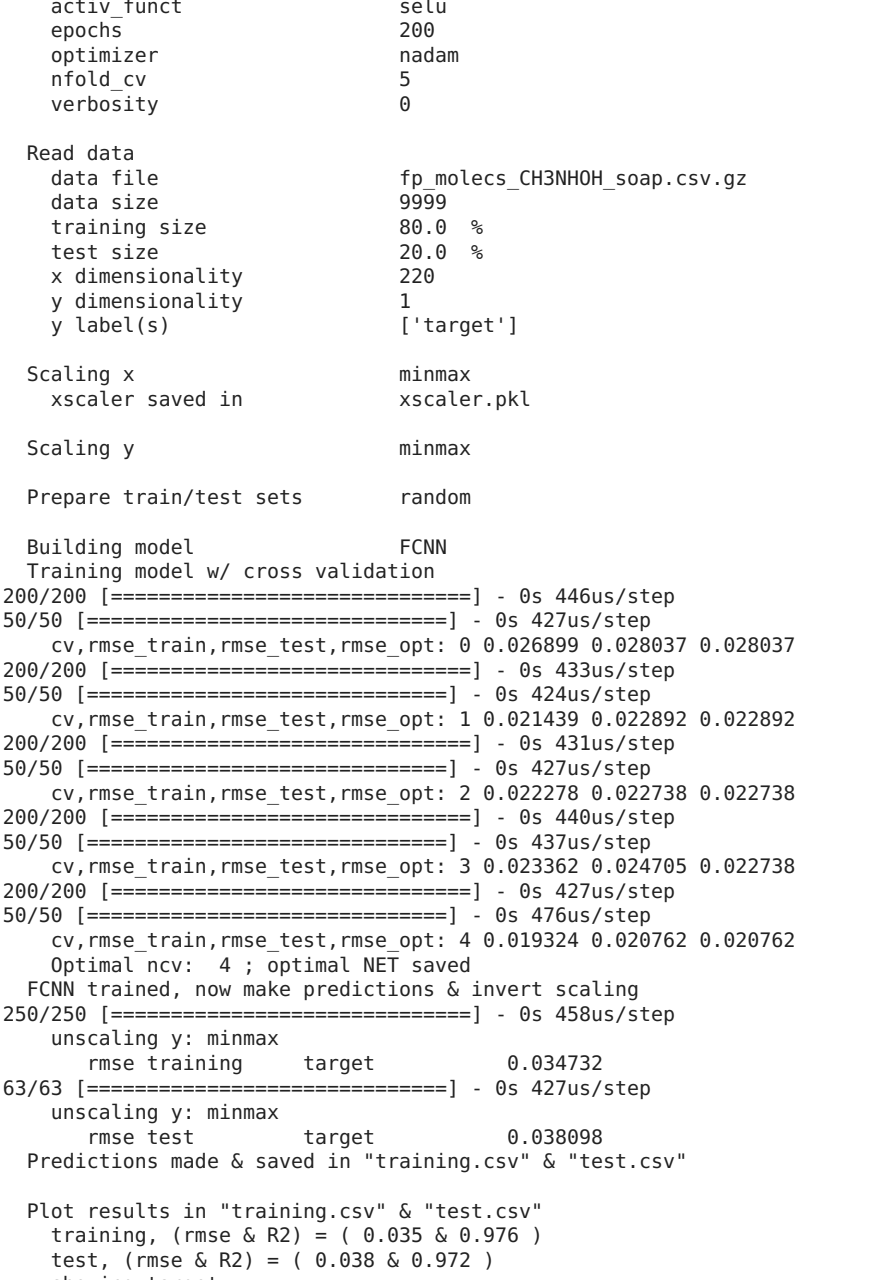
Prepare train/test sets random

Building model FCNN
Training model w/ cross validation
```

```
200/200 [=====] - 0s 359us/step
50/50 [=====] - 0s 374us/step
cv_rmse_train_rmse_test_rmse_opt: 0 0.05092 0.049035 0.049035
200/200 [=====] - 0s 367us/step
50/50 [=====] - 0s 363us/step
cv_rmse_train_rmse_test_rmse_opt: 1 0.040676 0.040751 0.040751
200/200 [=====] - 0s 397us/step
50/50 [=====] - 0s 376us/step
cv_rmse_train_rmse_test_rmse_opt: 2 0.045175 0.040834 0.040751
200/200 [=====] - 0s 360us/step
50/50 [=====] - 0s 359us/step
cv_rmse_train_rmse_test_rmse_opt: 3 0.047309 0.040480 0.040480
200/200 [=====] - 0s 365us/step
50/50 [=====] - 0s 362us/step
cv_rmse_train_rmse_test_rmse_opt: 4 0.047316 0.049786 0.040480
Optimal ncv: 3 ; optimal NET saved - 0s 476us/step
FCNN trained, now make predictions & invert scaling
250/250 [=====] - 0s 418us/step
```

```
unsclayng y: minmax
rmse training target 0.004167
63/63 [=====] - 0s 401us/step
unsclayng y: minmax
rmse test target 0.009707
Predictions made & saved in "training.csv" & "test.csv"
```

```
Plot results in "training.csv" & "test.csv"
training, (rmse & R2) = ( 0.084 & 0.859 )
test, (rmse & R2) = ( 0.090 & 0.846 )
showing target
```



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

Learning fingerprinted/featured data
algorithm fully connected NeuralNet w/ TensorFlow
layers [8, 8, 8]
activation selu
epochs 200
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_molecules_CH3NHOH_soap.csv.gz
data size 9999
training size 80.0 %
test size 20.0 %
x dimensionality 38
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
```

```
200/200 [=====] - 0s 446us/step
50/50 [=====] - 0s 427us/step
cv_rmse_train_rmse_test_rmse_opt: 0 0.026899 0.028037 0.028037
200/200 [=====] - 0s 433us/step
50/50 [=====] - 0s 424us/step
cv_rmse_train_rmse_test_rmse_opt: 1 0.021439 0.022892 0.022892
200/200 [=====] - 0s 430us/step
50/50 [=====] - 0s 427us/step
cv_rmse_train_rmse_test_rmse_opt: 2 0.022278 0.022738 0.022738
200/200 [=====] - 0s 437us/step
50/50 [=====] - 0s 425us/step
cv_rmse_train_rmse_test_rmse_opt: 3 0.023362 0.024705 0.022738
200/200 [=====] - 0s 427us/step
50/50 [=====] - 0s 476us/step
cv_rmse_train_rmse_test_rmse_opt: 4 0.019324 0.020762 0.020762
Optimal ncv: 4 ; optimal NET saved
FCNN trained, now make predictions & invert scaling
250/250 [=====] - 0s 458us/step
```

```
unsclayng y: minmax
rmse training target 0.034732
63/63 [=====] - 0s 427us/step
unsclayng y: minmax
rmse test target 0.030808
Predictions made & saved in "training.csv" & "test.csv"
```

```
Plot results in "training.csv" & "test.csv"
training, (rmse & R2) = ( 0.035 & 0.976 )
test, (rmse & R2) = ( 0.038 & 0.972 )
showing target
```



4. The same flowwork with the CH₄ dataset

```
In [6]: # Load data
from matsml.data import Datasets
import pandas as pd
import os

data = Datasets(dataset_name='molec3s_CH4')
data.load_dataset()
print(pd.read_csv('molec3s_CH4/summary.csv'))

# Fingerprint
from matsml.fingerprint import Fingerprint

summary = os.path.join(os.getcwd(), 'molec3s_CH4/summary.csv')
data_loc = os.path.join(os.getcwd(), 'molec3s_CH4/')
fp_type = 'pcm_molecules' # Intended fingerprint dimensionality; the final number can be smaller
# verbosity = 0 or 1

# PCM
data_params_pcm = {
    'fp_type': 'pcm_molecules',
    'summary': summary,
    'data_loc': data_loc,
    'fp_file': 'fp_CH4_pcm.csv',
    'fp_dim': fp_dim,
    'verbosity': verbosity
}

fp_pcm = Fingerprint(data_params_pcm)
fp_pcm.get_fingerprint()

# SOAP
data_params_soap = {
    'fp_type': 'soap_molecules',
    'summary': summary,
    'data_loc': data_loc,
    'fp_file': 'fp_CH4_soap.csv',
    'fp_dim': fp_dim,
    'verbosity': verbosity
}

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

```
# Data params
id_col = ['id'] # column for data ID
y_col = ['target'] # target property column(s)
comment_cols = [] # comment columns
n_train = 0.8 # 80% training, 20% validation
sampling = 'random' # train/test split method
x_scaling = 'minmax' # x scaling method
y_scaling = 'minmax' # y scaling method

# Dict of data parameters
data_params_CH4_pcm = {
    'data_file': 'fp_CH4_pcm.csv',
    'id_col': id_col,
    'y_cols': y_col,
    'comment_cols': comment_cols,
    'y_scaling': y_scaling,
    'x_scaling': x_scaling,
    'sampling': sampling,
    'n_train': n_train
}
```

```
# Dict of model parameters
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
}
```

```
# PCM Model
model = FCNN(data_params=data_params_CH4_pcm, model_params=model_params)
model.train()
model.plot(pdf_output=False)
```

```
# SOAP Model
model = FCNN(data_params=data_params_CH4_soap, model_params=model_params)
model.train()
model.plot(pdf_output=False)
```

```
Load requested dataset(s)
Data saved in molec3s_CH4
file name      target
0  CH4-00001.xyz -8.067043
1  CH4-00002.xyz -8.052410
2  CH4-00003.xyz -8.062079
3  CH4-00004.xyz -8.053099
4  CH4-00005.xyz -8.054724
...
9995 CH4-09996.xyz -8.061287
9996 CH4-09997.xyz -8.060672
9997 CH4-09998.xyz -8.013882
9998 CH4-09999.xyz -8.054027
9999 CH4-10000.xyz -8.052008

[10000 rows x 2 columns]
```

```
Atomic structure fingerprinting
summary /home/huantran/work_local/matssl/examples/ex2_molecules/molec3s_CH4/summary.csv
data_loc /home/huantran/work_local/matssl/examples/ex2_molecules/molec3s_CH4/
fp_type fp_CH4_pcm.csv
fp_file fp_CH4_pcm.csv
fp_dim 100
verbosity 0
Read input
num_structs 10000
Computing Coulomb matrix [=====] 100%
Projecting Coulomb matrix to create fingerprints [=====] 100%
Done fingerprinting, results saved in fp_CH4_pcm.csv
Atomic structure fingerprinting
summary /home/huantran/work_local/matssl/examples/ex2_molecules/molec3s_CH4/summary.csv
data_loc /home/huantran/work_local/matssl/examples/ex2_molecules/molec3s_CH4/
fp_type fp_CH4_pcm.csv
fp_file fp_CH4_pcm.csv
fp_dim 100
verbosity 0
Read input
num_structs 10000
Computing SOAP fingerprint with D5cribe [=====] 100%
Done fingerprinting, results saved in fp_CH4_soap.csv
```

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

Learning fingerprinted/featured data
algorithm fully connected NeuralNet w/ TensorFlow
layers [4, 4]
activation tanh
epochs 200
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_CH4_pcm.csv
data size 10000
training size 80.0 %
test size 20.0 %
x dimensionality 36
y dimensionality 1
y label(s) ['target']

Scaling x minmax
xscaler saved in xscaler.pkl

Scaling y minmax

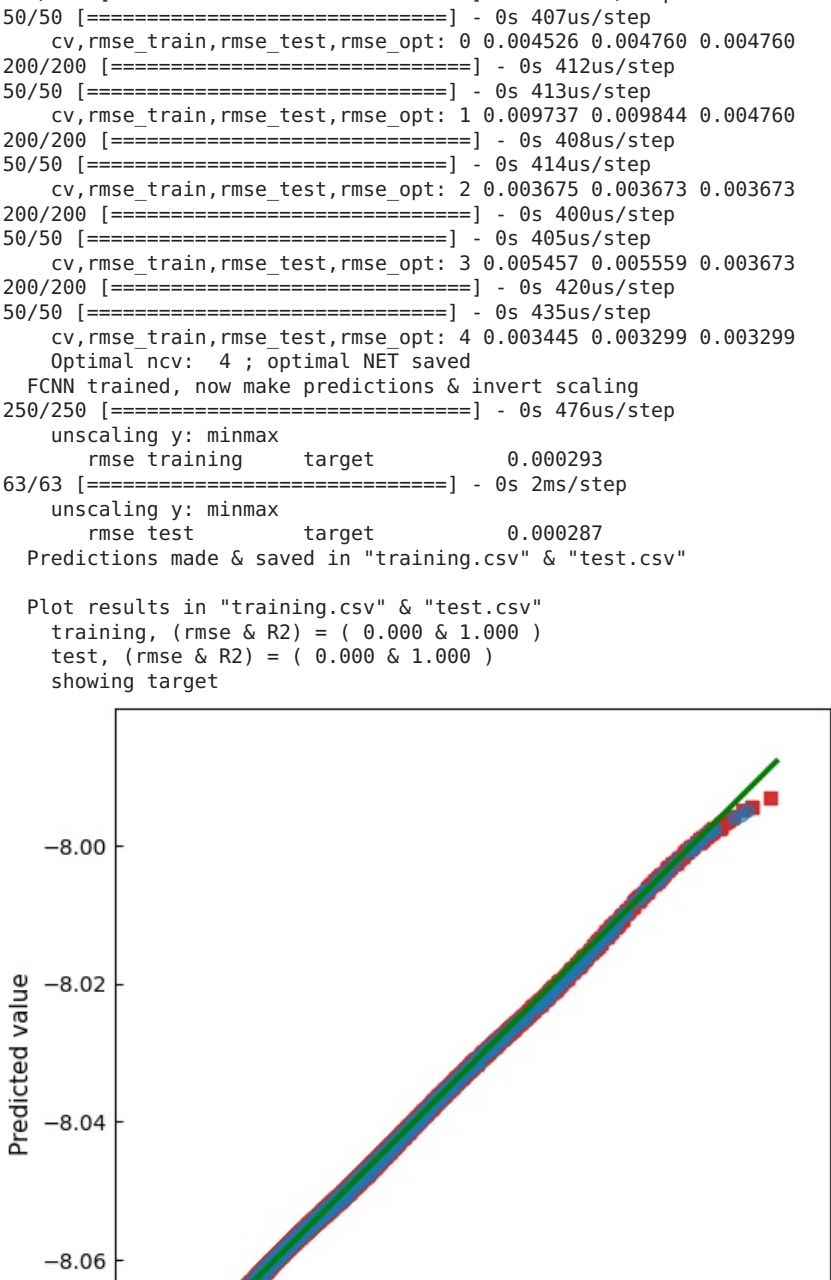
Prepare train/test sets random

Building model FCNN
Training model w/ cross validation
```

```
200/200 [=====] - 0s 301us/step
50/50 [=====] - 0s 349us/step
cv_rmse_train_rmse_test_rmse_opt: 0 0.004526 0.004760 0.004760
200/200 [=====] - 0s 413us/step
50/50 [=====] - 0s 344us/step
cv_rmse_train_rmse_test_rmse_opt: 1 0.009737 0.009844 0.004760
200/200 [=====] - 0s 400us/step
50/50 [=====] - 0s 414us/step
cv_rmse_train_rmse_test_rmse_opt: 2 0.003675 0.003673 0.003673
200/200 [=====] - 0s 420us/step
50/50 [=====] - 0s 435us/step
cv_rmse_train_rmse_test_rmse_opt: 3 0.005457 0.005559 0.003673
200/200 [=====] - 0s 420us/step
50/50 [=====] - 0s 420us/step
cv_rmse_train_rmse_test_rmse_opt: 4 0.003445 0.003299 0.003299
Optimal ncv: 4 ; optimal NET saved
FCNN trained, now make predictions & invert scaling
250/250 [=====] - 0s 387us/step
```

```
unsclayng y: minmax
rmse training target 0.000293
63/63 [=====] - 0s 380us/step
unsclayng y: minmax
rmse test target 0.000287
Predictions made & saved in "training.csv" & "test.csv"
```

```
Plot results in "training.csv" & "test.csv"
training, (rmse & R2) = ( 0.000 & 1.000 )
test, (rmse & R2) = ( 0.000 & 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 tanh
epochs 200
optimizer nadam
nfold_cv 5
verbosity 0

Read data
data file fp_CH4_soap.csv
data size 10000
training size 80.0 %
test size 20.0 %
x dimensionality 36
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
```

```
200/200 [=====] - 0s 420us/step
50/50 [=====] - 0s 420us/step
cv_rmse_train_rmse_test_rmse_opt: 0 0.004526 0.004760 0.004760
200/200 [=====] - 0s 413us/step
50/50 [=====] - 0s 344us/step
cv_rmse_train_rmse_test_rmse_opt: 1 0.009737 0.009844 0.004760
200/200 [=====] - 0s 400us/step
50/50 [=====] - 0s 414us/step
cv_rmse_train_rmse_test_rmse_opt: 2 0.003675 0.003673 0.003673
200/200 [=====] - 0s 420us/step
50/50 [=====] - 0s 435us/step
cv_rmse_train_rmse_test_rmse_opt: 3 0.005457 0.005559 0.003673
200/200 [=====] - 0s 420us/step
50/50 [=====] - 0s 420us/step
cv_rmse_train_rmse_test_rmse_opt: 4 0.003445 0.003299 0.003299
Optimal ncv: 4 ; optimal NET saved
FCNN trained, now make predictions & invert scaling
250/250 [=====] - 0s 387us/step
```

```
unsclayng y: minmax
rmse training target 0.000293
63/63 [=====] - 0s 380us/step
unsclayng y: minmax
rmse test target 0.000287
Predictions made & saved in "training.csv" & "test.csv"
```

```
Plot results in "training.csv" & "test.csv"
training, (rmse & R2) = ( 0.000 & 1.000 )
test, (rmse & R2) = ( 0.000 & 1.000 )
showing target
```


In [] : Processing math: 100%