

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.matsml.org.

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

# Load a dataset
dataset_name='molec_CH3NHOH'
data=Datasets(dataset_name=dataset_name)
data.load_dataset()

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

matsML, v1.0.1
*****
Load requested dataset(s)
Data saved in molecs_CH3NHOH
file_name target
0 CH3NHOH_00001.xyz -940.288539
1 CH3NHOH_00002.xyz -940.580309
2 CH3NHOH_00003.xyz -940.184809
3 CH3NHOH_00004.xyz -940.460977
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 atomization 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

summary=os.path.join(os.getcwd(),'molecs_CH3NHOH/summary.csv')
data_loc=os.path.join(os.getcwd()), 'molecs_CH3NHOH/'
fp_dim=50 # intended fingerprint dimensionality; the final number can be smaller
verbosity=0 # 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/huan/work/matsml/examples/ex2_molecules/molecules_CH3NHOH/summary.csv
data_loc /home/huan/work/matsml/examples/ex2_molecules/molecules_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/huan/work/matsml/examples/ex2_molecules/molecules_CH3NHOH/summary.csv
data_loc /home/huan/work/matsml/examples/ex2_molecules/molecules_CH3NHOH/
fp_type soap_molecules
fp_file fp_soap.csv
fp_dim 50
verbosity 0
Read input num_structs 9999
Computing SOAP fingerprint with DScRibe [=====] 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 saved 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_cols=['target'] # columns for (one or more) target properties
comment_cols=[] # comment columns, anything not counted into ID, fingerprints, and target
n_trains=0.8 # 80% for training, 20% for validating
sampling='random' # method for train/test 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_pcm.csv', 'id_col':id_col, 'y_cols':y_cols, 'comment_cols':comment_cols,
                 'y_scaling':y_scaling, 'x_scaling':x_scaling, 'sampling':sampling, 'n_trains':n_trains}

data_params_soap={'data_file':'fp_soap.csv', 'id_col':id_col, 'y_cols':y_cols, 'comment_cols':comment_cols,
                  'y_scaling':y_scaling, 'x_scaling':x_scaling, 'sampling':sampling, 'n_trains':n_trains}

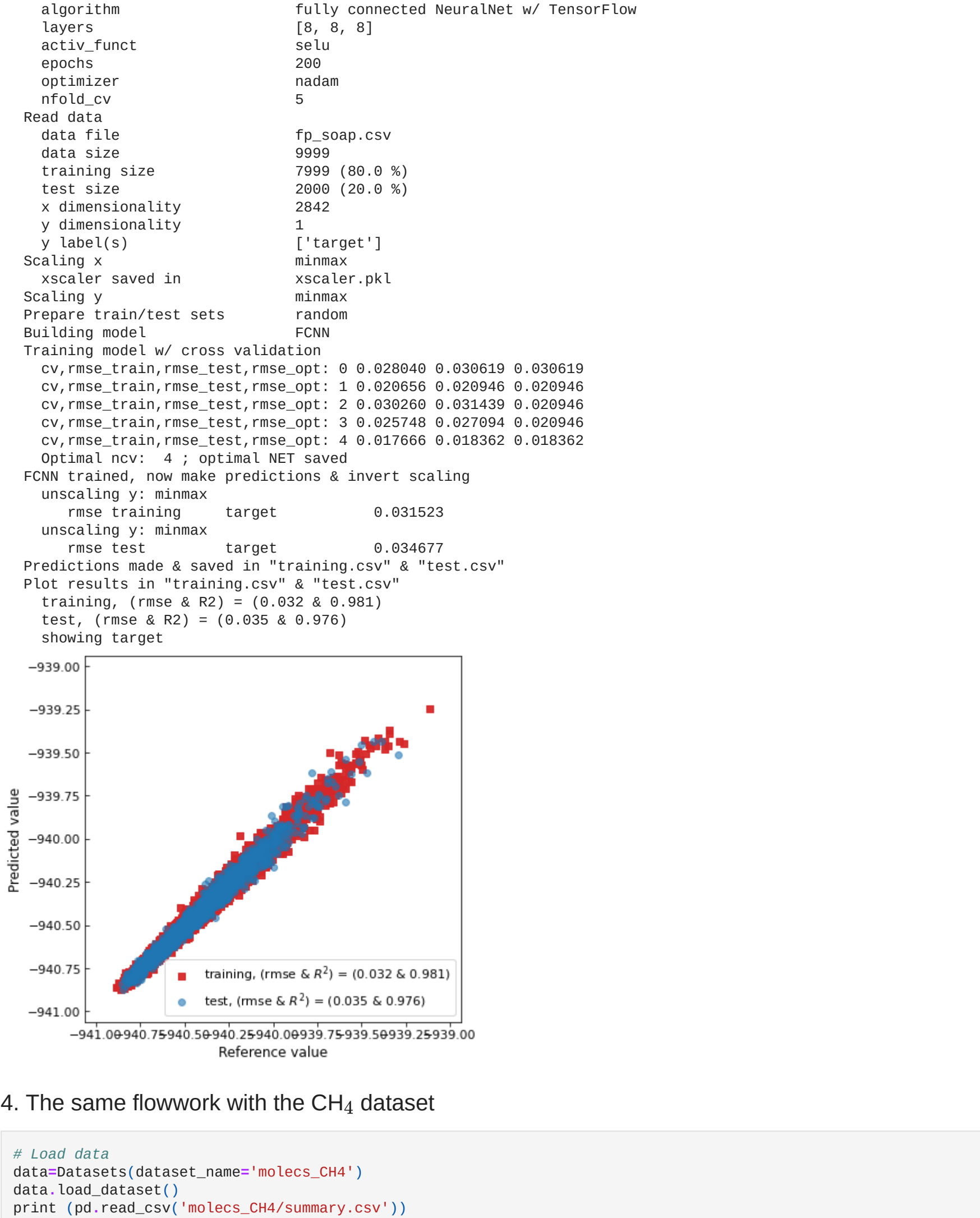
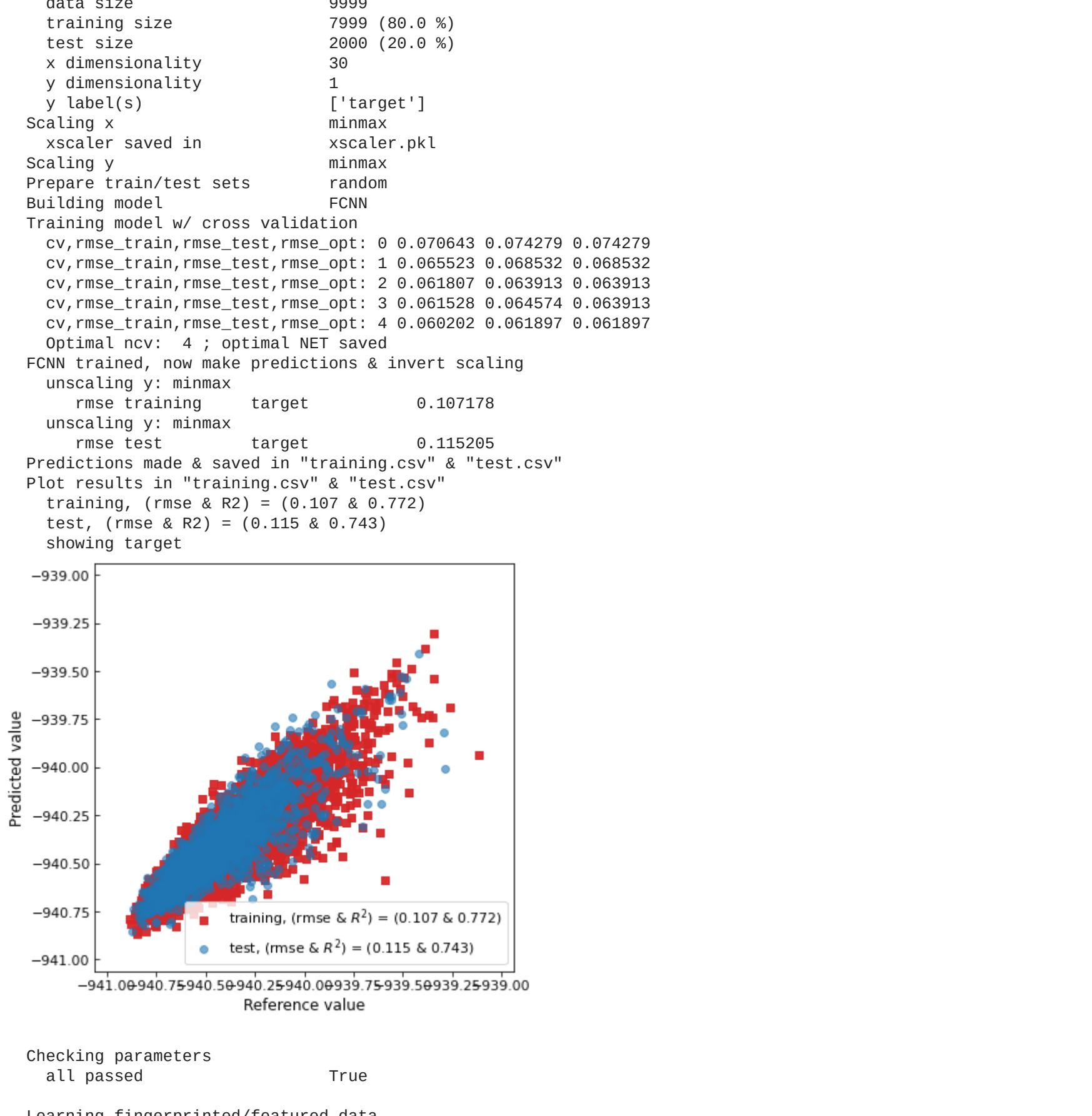
In [5]: from matsml.models import FCNN

# Model parameters
layers=[8,8,8] # List of nodes in hidden layers
epochs=200 # Epochs
nfold_cv=5 # Number of folds for cross validation
use_bias=True # Use bias term 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'
activ_func='tanh' # Options: "tanh", "relu", and more
optimizer='nadam' # Options: "Nadam", "Adam", and more

# Dict of model parameters
model_params={'layers':layers, 'activ_func':activ_func, '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)
```



4. The same flowwork with the CH₄ dataset

```
In [6]: # Load data
data=Datasets(dataset_name='molecules_CH4')
data.load_dataset()
print (pd.read_csv('molecules_CH4/summary.csv'))

# Fingerprint
from matsml.fingerprint import Fingerprint
summary=os.path.join(os.getcwd(),'molecs_CH4/summary.csv')
data_loc=os.path.join(os.getcwd()), 'molecs_CH4/'
fp_type='pcm_molecules' # projected Coulomb matrix for molecules
fp_dim=100 # intended fingerprint dimensionality; the final number can be smaller
verbosity=0 # 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_cols=['target'] # columns for (one or more) target properties
comment_cols=[] # comment columns, anything not counted into ID, fingerprints, and target
n_trains=0.8 # 80% for training, 20% for validating
sampling='random' # method for train/test splitting
x_scaling='minmax' # method for x scaling
y_scaling='minmax' # method for y scaling

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

data_params_CH4_soap={'data_file':'fp_CH4_soap.csv', 'id_col':id_col, 'y_cols':y_cols, 'comment_cols':comment_cols,
                      'y_scaling':y_scaling, 'x_scaling':x_scaling, 'sampling':sampling, 'n_trains':n_trains}

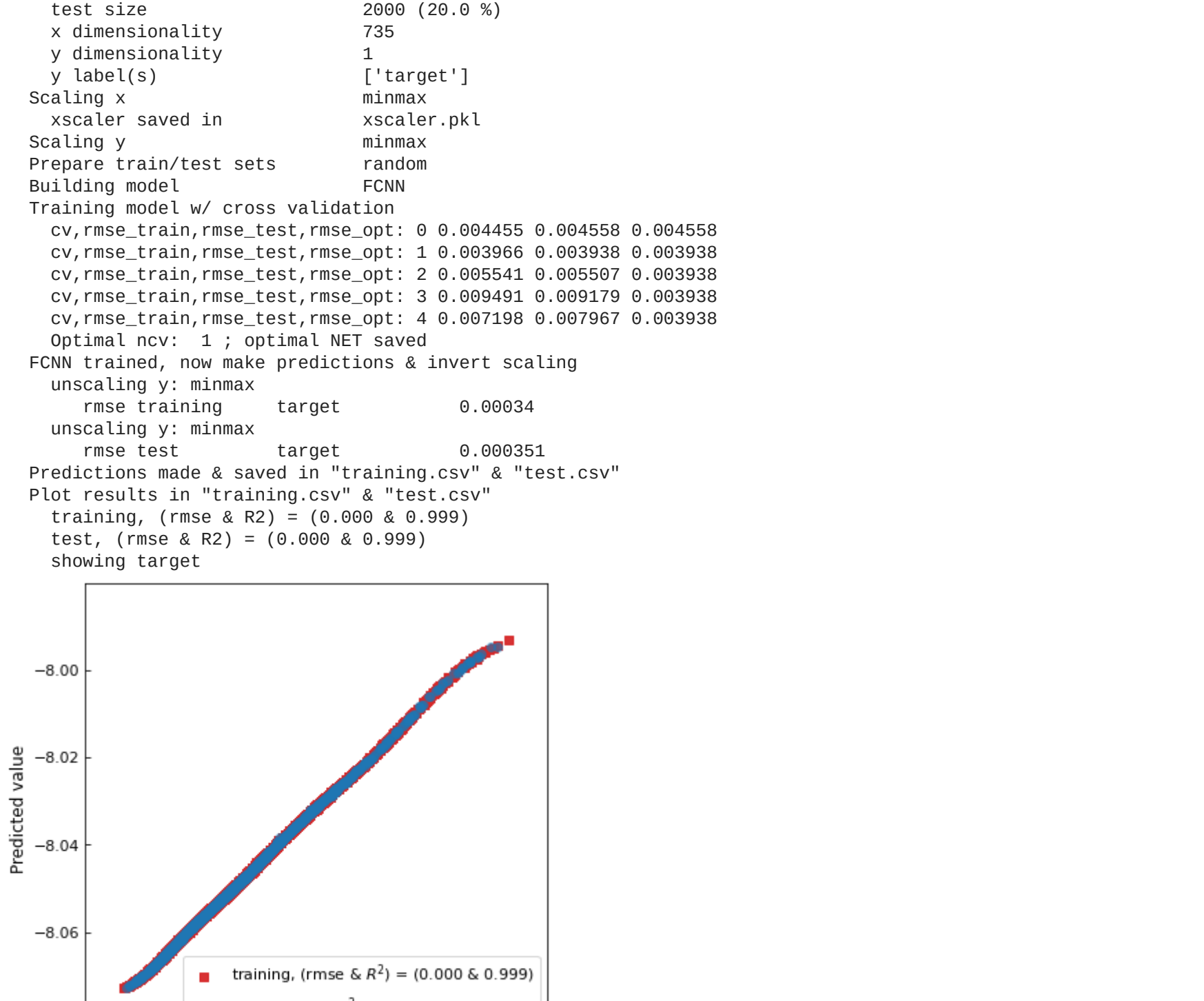
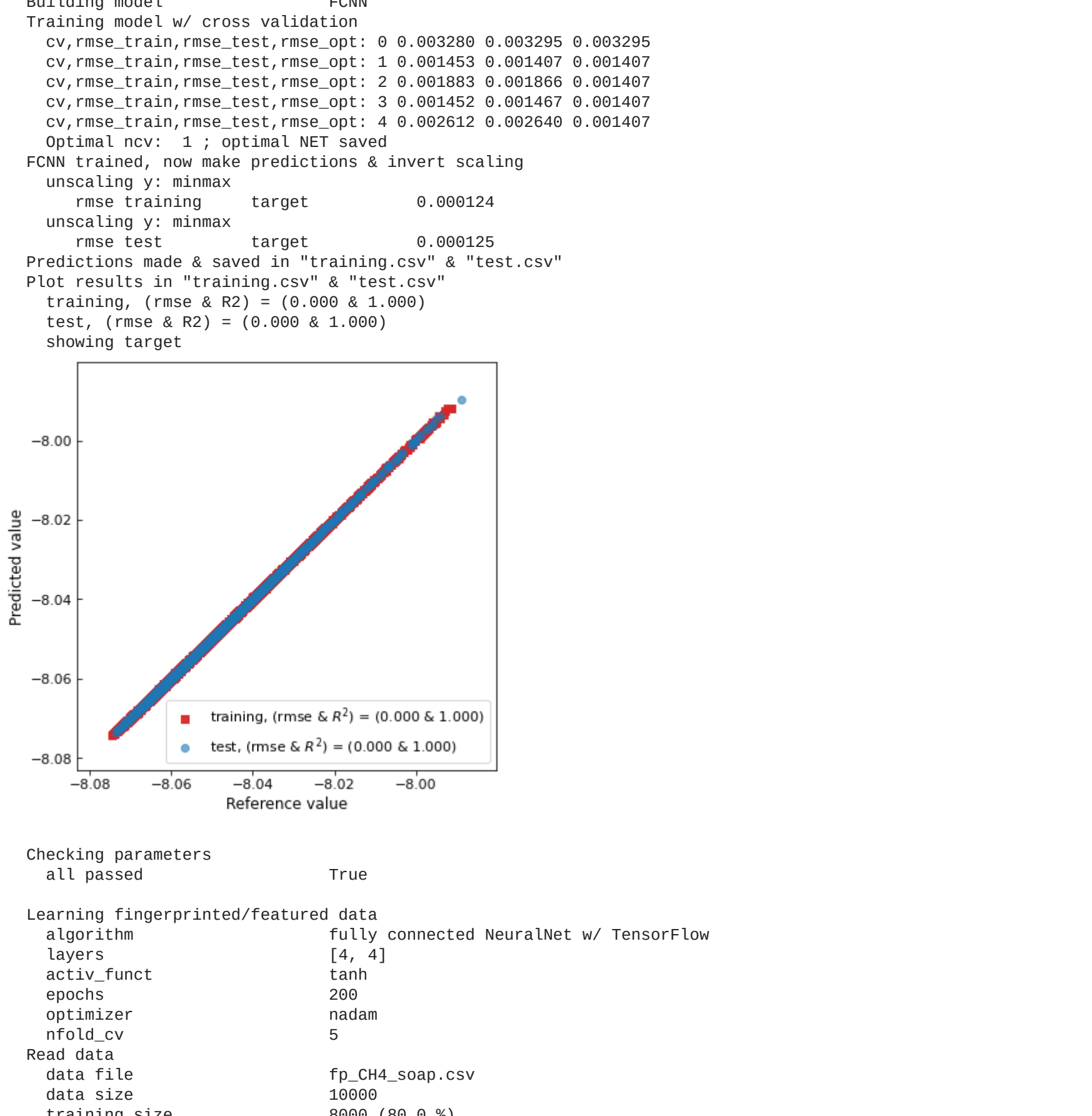
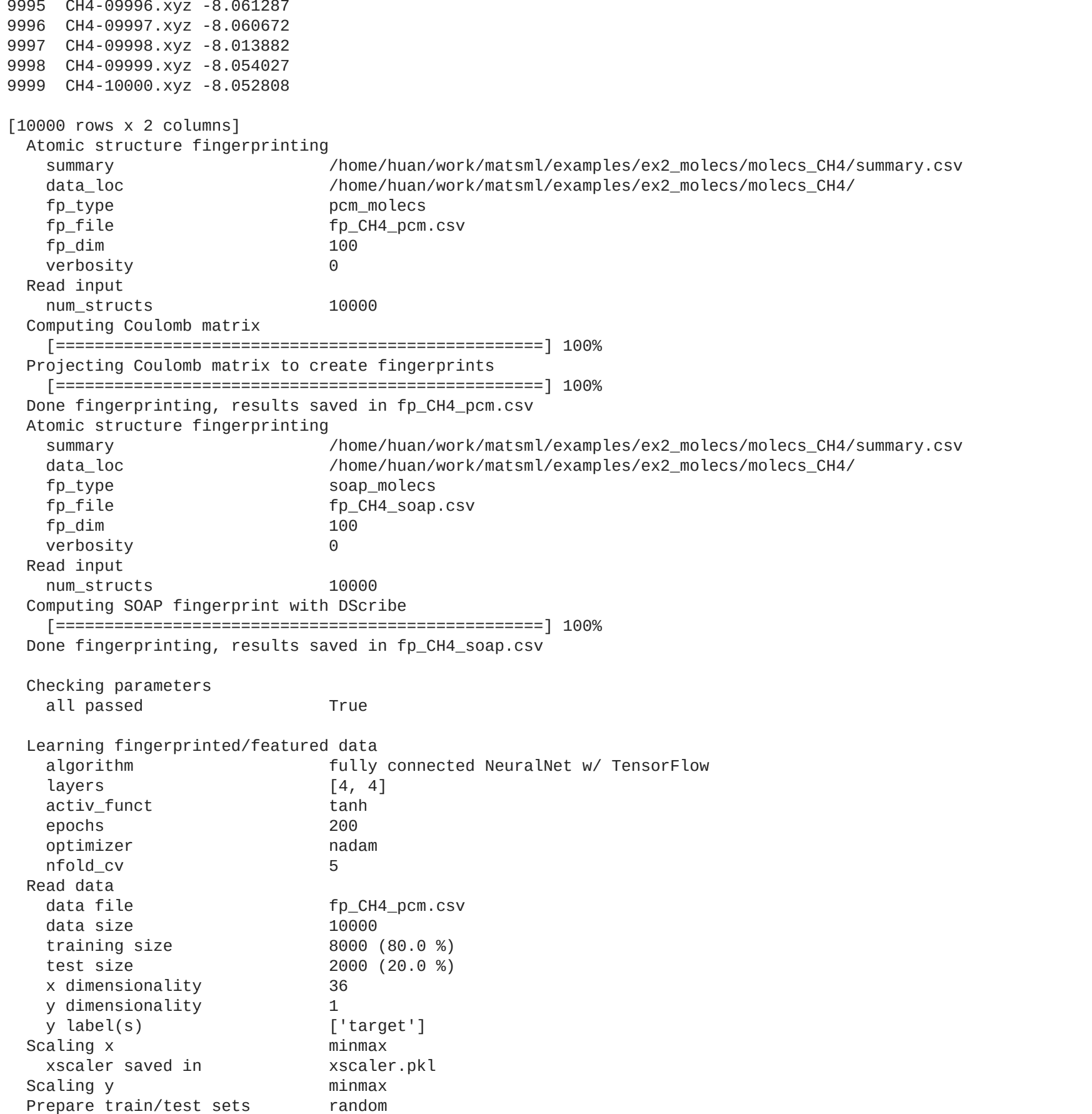
# Models with FCNN
from matsml.models import FCNN

# Model parameters
layers=[4,4] # List of nodes in hidden layers
epochs=200 # Epochs
nfold_cv=5 # Number of folds for cross validation
use_bias=True # Use bias term 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'
activ_func='tanh' # Options: "tanh", "relu", and more
optimizer='nadam' # Options: "Nadam", "Adam", and more

# Dict of model parameters
model_params={'layers':layers, 'activ_func':activ_func, '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_CH4_pcm, model_params=model_params)
model.train()
model.plot(pdf_output=False)

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



```
In [ ]: 
```