

ML models for computed band gap of HOIPs

Huan Tran, Georgia Institute of Technology

This notebook is a part of [V. N. Tuo, Nga, T. T. Nguyen, V. Sharma, and T. D. Huan, *Probabilistic deep learning approach for targeted hybrid organic-inorganic perovskites*, 2021], and will also be an example of matsML toolkit. Results obtained here can be found in this work.

The original (raw) dataset containing the computed band gap of 1,346 atomic structures predicted for 192 chemical compositions of hybrid organic-inorganic perovskites (HOIPs) is available at [C. Kim, T.D. Huan, S. Krishnan, and R. Ramprasad, Scientific Data 4, 170057 ("17); <https://www.nature.com/articles/sdata201757>]. Here, three fingerprinted versions of this dataset (S1, S2, and S3) will be fetched from <http://www.matsml.org/> and learned to develop 5 ML models (M1, M2, M3, M4, and M5), which are based on Gaussian Process Regression, fully connected Neural Net, and Probability Neural Net. Computations performed using matsML toolkit, available at <https://github.com/huand/matsml.git>.

Among 5 models developed, M5 demonstrates a reasonable way to handle the aleatoric uncertainty in deep learning of materials data. More details on this topic can be found in "*Probabilistic deep learning approach for targeted hybrid organic-inorganic perovskites*", the reference mentioned above.

1. Download data

Three (fingerprinted) datasets (S1, S2, and S3) used for the work will be obtained. In fact, S2 has 2 versions, one with selector and one not.

```
In [1]: from matsml.data import Datasets

data=Datasets(S1='fp_hoips_S1_1dest.csv.gz',S2a='fp_hoips_S2a_2dest',
              S2b='fp_hoips_S2b_1dest',S3='fp_hoips_S3_4tfp')
data.load_dataset()

matsML, version 1.0.0
****
Load requested dataset(s)
Data saved in fp_hoips_S1_1dest.csv.gz
Data saved in fp_hoips_S2a_2dest.csv.gz
Data saved in fp_hoips_S2b_1dest.csv.gz
Data saved in fp_hoips_S3_4tfp.csv.gz
```

2. Obtained datasets parameters

```
In [2]: # data parameters for learning

n_train=0.9 # 90% for training, 10% for test
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
data1_params={'data_file':'fp_hoips_S1_1dest.csv.gz','id_col':['ID'],'y_cols':['Ymean'],'comment_cols':[],
              'y_scaling':y_scaling,'x_scaling':x_scaling,'sampling':sampling,'n_trains':n_train}
data2a_params={'data_file':'fp_hoips_S2a_2dest.csv.gz','id_col':['ID'],'y_cols':['Ymean'],'ystd':,
               'comment_cols':[],'y_scaling':y_scaling,'x_scaling':x_scaling,'sampling':sampling,'n_trains':n_train}
data2b_params={'data_file':'fp_hoips_S2b_1dest.csv.gz','id_col':['ID'],'y_cols':['prop_value'],
               'comment_cols':['Ymean','Ystd','hid'],'y_scaling':y_scaling,'x_scaling':x_scaling,
               'sampling':sampling,'n_trains':n_trains}
data3_params={'data_file':'fp_hoips_S3_4tfp.csv.gz','id_col':['ID'],'y_cols':['Egap'],'comment_cols':[],
              'x_scaling':x_scaling,'y_scaling':y_scaling, 'sampling':sampling,'n_trains':1.0}
```

3. ML Models

3a. Model M1: GPR on S1

```
In [3]: from matsml.models import GPR

# Model parameters
nfold_cv=5 # Number of folds for cross validation
model_file='M1.pkl' # Name of the model file to be created
verbosity=0
rmse_cv=False
n_restarts_optimizer=100

model_params={'nfold_cv':nfold_cv,'n_restarts_optimizer':n_restarts_optimizer,'model_file':model_file,
              'verbosity':verbosity,'rmse_cv':rmse_cv}

model=GPR(data_params=data1_params,model_params=model_params)
model.train()
model.plot(pdf_output=False)
```

Learning fingerprinted/featured data

algorithm	Gaussian process regression w/ scikit-learn
nfold_cv	5
optimizer	fmin_l_bfgs_b
n_restarts_optimizer	100
rmse_cv	False

Checking parameters

all passed	True
------------	------

Read data

data file	fp_hoips_S1_1dest.csv.gz
data size	192
training size	172 (89.6 %)
test size	20 (10.4 %)
x dimensionality	32
y dimensionality	1
y label(s)	['Ymean']

Scaling x

minmax	minmax
--------	--------

Scaling y

xscaler saved in	xscaler.pkl
------------------	-------------

Prepare train/test sets

random	random
--------	--------

Training model w/ cross validation

cv,rmse_train,rmse_test,rmse_opt	0 0.018522 0.038449 0.038449
cv,rmse_train,rmse_test,rmse_opt	1 0.018840 0.051316 0.038449
cv,rmse_train,rmse_test,rmse_opt	2 0.021326 0.037287 0.037287
cv,rmse_train,rmse_test,rmse_opt	3 0.018741 0.028500 0.028500
cv,rmse_train,rmse_test,rmse_opt	4 0.019046 0.043592 0.028500

GPR model trained, now make predictions & invert scaling

rmse training	Ymean	0.088035
unscaling y: minmax		
rmse test	Ymean	0.08475

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

Plot results in "training.csv" & "test.csv"

training, (rmse & R²) = (0.088 & 0.993)

test, (rmse & R²) = (0.085 & 0.990)

showing tmean

3b. Model M2: FCNN on S1

```
In [4]: from matsml.models import FCNN

# Model parameters
layers=[5] # 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='M2.pkl' # Name of the model file to be created
verbosity=0 # Verbosity, 0 or 1
batch_size=32 # Default = 32
loss='mse'
activ_func='selu' # Options: "cnnh", "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}

model=FCNN(data_params=data1_params,model_params=model_params)
model.train()
model.plot(pdf_output=False)
```

Learning fingerprinted/featured data

algorithm	fully connected NeuralNet w/ TensorFlow
layers	[5]
activ_func	selu
epochs	2000
optimizer	nadam
nfold_cv	5

Checking parameters

all passed	True
------------	------

Read data

data file	fp_hoips_S1_1dest.csv.gz
data size	192
training size	172 (89.6 %)
test size	20 (10.4 %)
x dimensionality	32
y dimensionality	1
y label(s)	['Ymean']

Scaling x

minmax	minmax
--------	--------

Scaling y

xscaler saved in	xscaler.pkl
------------------	-------------

Prepare train/test sets

random	random
--------	--------

Building model

FCNN	FCNN
------	------

Training model w/ cross validation

cv,rmse_train,rmse_test,rmse_opt	0 0.025405 0.033866 0.033866
cv,rmse_train,rmse_test,rmse_opt	1 0.018399 0.043445 0.033866
cv,rmse_train,rmse_test,rmse_opt	2 0.019758 0.084789 0.033866
cv,rmse_train,rmse_test,rmse_opt	3 0.020367 0.028783 0.028783
cv,rmse_train,rmse_test,rmse_opt	4 0.018889 0.026327 0.026327

Optimal ncw: 4 / Optimal NET saved

FCNN trained, now make predictions & invert scaling

unscaling y: minmax		
rmse training	Ymean	0.088506
unscaling y: minmax		
rmse test	Ymean	0.137268

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

Plot results in "training.csv" & "test.csv"

training, (rmse & R²) = (0.089 & 0.993)

test, (rmse & R²) = (0.137 & 0.982)

showing Ymean

3c. Model M3: FCNN on S2a

```
In [5]: # Model parameters
layers=[5,5]
epochs=10000
nfold_cv=5
use_bias=True
model_file='M3.pkl'
loss='mse'
verbosity=0
batch_size=32
activ_func='elu'
optimizer='nadam'

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}

model=FCNN(data_params=data2a_params,model_params=model_params)
model.train()
model.plot(pdf_output=False)
```

Learning fingerprinted/featured data

algorithm	fully connected NeuralNet w/ TensorFlow
layers	[5, 5]
activ_func	elu
epochs	10000
optimizer	nadam
nfold_cv	5

Checking parameters

all passed	True
------------	------

Read data

data file	fp_hoips_S2a_2dest.csv.gz
data size	192
training size	172 (89.6 %)
test size	20 (10.4 %)
x dimensionality	31
y dimensionality	2
y label(s)	['Ymean', 'Ystd']

Scaling x

minmax	minmax
--------	--------

Scaling y

xscaler saved in	xscaler.pkl
------------------	-------------

Prepare train/test sets

random	random
--------	--------

Building model

FCNN	FCNN
------	------

Training model w/ cross validation

cv,rmse_train,rmse_test,rmse_opt	0 0.050190 0.086481 0.086481
cv,rmse_train,rmse_test,rmse_opt	1 0.042709 0.078880 0.078880
cv,rmse_train,rmse_test,rmse_opt	2 0.034806 0.087286 0.078880
cv,rmse_train,rmse_test,rmse_opt	3 0.034643 0.078801 0.078801
cv,rmse_train,rmse_test,rmse_opt	4 0.033208 0.087348 0.078801

Optimal ncw: 3 / Optimal NET saved

FCNN trained, now make predictions & invert scaling

unscaling y: minmax		
rmse training	Ymean	0.114854
rmse training	Ystd	0.070648
unscaling y: minmax		
rmse test	Ymean	0.173631
rmse test	Ystd	0.153502

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

Plot results in "training.csv" & "test.csv"

training, (rmse & R²) = (0.115 & 0.988)

test, (rmse & R²) = (0.174 & 0.971)

showing Ystd

training, (rmse & R²) = (0.071 & 0.929)

test, (rmse & R²) = (0.154 & 0.585)

showing Ystd

training, (rmse & R²) = (0.071 & 0.829)

test, (rmse & R²) = (0.154 & 0.585)

showing Ystd

3d. Model M4: FCNN on S2b

```
In [6]: # Model parameters
layers=[4,4]
epochs=5000
nfold_cv=5
use_bias=True
model_file='M4.pkl'
loss='mse'
verbosity=0
batch_size=32
activ_func='tanh'
optimizer='nadam'

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}

model=FCNN(data_params=data2b_params,model_params=model_params)
model.train()
model.plot(pdf_output=False)
```

Learning fingerprinted/featured data

algorithm	fully connected NeuralNet w/ TensorFlow
layers	[4, 4]
activ_func	tanh
epochs	5000
optimizer	nadam
nfold_cv	5

Checking parameters

all passed	True
------------	------

Read data

data file	fp_hoips_S2b_1dest.csv.gz
data size	384
training size	345 (89.8 %)
test size	39 (10.2 %)
x dimensionality	53
y dimensionality	1
y label(s)	['prop_value']

Scaling x

minmax	minmax
--------	--------

Scaling y

xscaler saved in	xscaler.pkl
------------------	-------------

Prepare train/test sets

random	random
--------	--------

Building model

FCNN	FCNN
------	------

Training model w/ cross validation

cv,rmse_train,rmse_test,rmse_opt	0 0.045158 0.112265 0.112265
cv,rmse_train,rmse_test,rmse_opt	1 0.048173 0.090974 0.090974
cv,rmse_train,rmse_test,rmse_opt	2 0.047811 0.080596 0.080596
cv,rmse_train,rmse_test,rmse_opt	3 0.050388 0.059932 0.059932
cv,rmse_train,rmse_test,rmse_opt	4 0.045543 0.106781 0.059932

Optimal ncw: 3 / Optimal NET saved

FCNN trained, now make predictions & invert scaling

unscaling y: minmax		
rmse training	selector1 prop_value	0.148816
rmse training	selector2 prop_value	0.07692
unscaling y: minmax		
rmse test	selector1 prop_value	0.130914
rmse test	selector2 prop_value	0.172489

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

Plot results in "training.csv" & "test.csv"

training, (rmse & R²) = (0.119 & 0.996)

test, (rmse & R²) = (0.156 & 0.993)

showing prop_value

3e. Model M5: Probabilistic Neural Net on S3

```
In [9]: from matsml.models import PrFCNN

layers=[5]
epochs=200
nfold_cv=5
model_file='M5.pkl'
loss='mse'
verbosity=0
batch_size=32
activ_func='selu'
optimizer='nadam'

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}

model=PrFCNN(data_params=data3_params,model_params=model_params)
model.train()
model.plot(pdf_output=False)
```

Learning fingerprinted/featured data

algorithm	Probabilistic NeuralNet w/ TensorFlow-Probability
layers	[5]
activ_func	selu
epochs	200
optimizer	nadam
nfold_cv	5

Checking parameters

all passed	True
------------	------

Read data

data file	fp_hoips_S3_4tfp.csv.gz
data size	1346
training size	1346 (100.0 %)
test size	0 (0.0 %)
x dimensionality	221
y dimensionality	1
y label(s)	['Egap']

Scaling x

minmax	minmax
--------	--------

Scaling y

xscaler saved in	xscaler.pkl
------------------	-------------

Prepare train/test sets

random	random
--------	--------

Building model

PrFCNN	PrFCNN
--------	--------

Training PrFCNN w/ cross validation

cv,rmse_train,rmse_test,rmse_opt	0 0.620635 0.582243 0.582243
cv,rmse_train,rmse_test,rmse_opt	1 0.603726 0.649592 0.582243
cv,rmse_train,rmse_test,rmse_opt	2 0.580510 0.597396 0.582243
cv,rmse_train,rmse_test,rmse_opt	3 0.599846 0.621374 0.582243
cv,rmse_train,rmse_test,rmse_opt	4 0.622166 0.654577 0.582243

Optimal ncw: 3

PrFCNN trained, now make predictions & invert scaling

unscaling y: none		
rmse training	Egap	0.428167

Predictions made & saved in "training.csv"

Plot results in "training.csv" & "test.csv"

showing Egap

```
In [10]: import io, requests
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
import matplotlib

# read the trained data
pred=pd.read_csv('training.csv')

# Trained data provides ID of 1346 cases, but we need the name of the organic cations A,
# cation B, and anions X also. They can be obtained here
sum_url='http://www.matsml.org/data/hoips201.csv'
resp=requests.get(sum_url)
pred=pd.read_csv(io.StringIO(resp.get('sum_url').content.decode('utf-8'))))

# We will plot 16 compositions ASnI3 made of 16 cations A and Sn for B and I for X.
organics=['Acetamidinium','Ammonium','Azetidinium','Butylammonium','Dimethylammonium','Ethylammonium',
          'Formamidinium','Guanidinium','Hydrazinium','Hydroxylammonium','Indazolium','Isoniazolium',
          'Methylammonium','Propylammonium','Tetramethylammonium','Trimethylammonium']
anions = ['Sn']

comps = [(organic,cation,anion) for organic in organics for cation in cations\
          for anion in anions]

# 2 DataFrame will be extracted from pred for plotting
bandgap_strs=pd.DataFrame(columns=['cid','id','organic_cat','bandgap'])
bandgap_comp=pd.DataFrame(columns=['cid','mean_comput','std_comput','mean_pred','std_pred'])

# For each of 16 compositions, extract needed data and store in "bandgap_strs" and "bandgap_comp"
for cid in range(len(comps)):
    comp = comps[cid]
    sel_ids = mapping([mapping('organic')==comp[0]] & \
                      (mapping('cation')==comp[1]) & (mapping('anion')==comp[2]))
    sel_id = list(sel_rows['id']).isin(sel_ids)
    sel_pred = pred[pred['id'].isin(sel_ids)]
    sel_pred.reset_index(drop=True, inplace=True)
    bandgap_strs.loc[len(bandgap_strs)]=(cid,np.mean(sel_pred['Egap']), \
    np.std(sel_pred['Egap']),sel_pred.at[0,'nd_Egap'], sel_pred.at[0,'md_Egap_err'])
    for idx,eg in zip(list(sel_id['id']), list(sel_pred['Egap']))):
        bandgap_strs.loc[len(bandgap_strs)]=(cid,idx,comp[0],eg)

# Make figure
fig,ax=plt.subplots(figsize=(8,6),frameon=True)
plt.subplots_adjust(left=0.12, bottom=0.21, right=0.98, top=0.98, wspace=0,
                    hspace=0)

plt.rc('font',size=16)

plt.box(True)
plt.tick_params(axis='x',which='both',bottom=True,top=True,labelbottom=True)
plt.tick_params(axis='y',which='both',right=True,left=True,direction='in',labelleft=True,length=5)
ax.set_ylim([0,4,0])

plt.tick_params(axis='x',which='both',direction='in',labelsize=12,top=True)
plt.tick_params(axis='y',which='both',direction='in',labelsize=12,right=True)
ax.set_xticks(np.arange(0,16,1))
ax.set_xticklabels(organics, rotation=35, ha='right')
plt.ylabel('Eg (eV)',color='black',fontsize=18)
ax.scatter(bandgap_strs['cid'], bandgap_strs['bandgap'],color='royalblue',
           alpha=0.75, zorder=3, label='computed data')
ax.errorbar(bandgap_comp['cid'],bandgap_comp['mean_comput'], yerr=bandgap_comp['std_comput'],
           color='darkgoldenrod', alpha=0.5, zorder=3, markersize=5,
           fsize=12, label='predicted Eg (eV)')
ax.plot(bandgap_comp['cid'], bandgap_comp['mean_pred'],color='tab:red',linewidth=2,
       label='predicted Eg (eV)')

ax.fill_between(bandgap_comp['cid'], bandgap_comp['mean_pred'] - 2* bandgap_comp['std_pred'],
               bandgap_comp['mean_pred'] + 2* bandgap_comp['std_pred'],color='0.15',
               label='predicted Eg (eV)')

handles, labels = ax.get_legend_handles_labels()
labels, handles = zip(*sorted(zip(labels, handles), key=lambda t: t[0]))
ax.legend(handles, labels, loc='lower right', fontsize=15)

Out[10]: <matplotlib.legend.Legend at 0x7f4987b5970>
```



Fig. 1. Electronic band gap E_g (circles) computed for the predicted atomic structures of ASnI₃, 16 HOIP formulas corresponding to 16 organic/anions A, for each formula, the mean and standard deviation of E_g , i.e., E_g^{mean} and E_g^{std} , are given by dark golden squares and associated errors. Predicted $E_g^{\text{mean}} \pm 2E_g^{\text{std}}$ is given in red while the shaded area indicates the 95-percent confidence interval ($E_g \pm 2E_g^{\text{std}}$) of the predictions using the probabilistic model developed in this work

```
In [1]:
```