

"Machine Learning for Climate Predictions" Workshop

Women in Data Science Livermore - March 7th, 2022

Presented by Gemma Anderson

Jupyter Notebook prepared by Cindy Gonzales

Artificial intelligence (AI) methods have begun to show great promise in climate science applications. This notebook will introduce and describe a dataset and methods to make machine learning predictions of a multiresolution climate model ensemble.

For more information on the dataset and methods used, please reference the paper:

[Machine Learning Predictions of Multi-Resolution Climate Model Ensembles– Anderson, G. J. and Lucas, D. D. \(GRL: in review\). IM release number: LLNL-JRNL-744178.](#)

Auspices Statement: This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.

IM #: LLNL-PRES-832485

Downloading the data

Let's start by downloading the data locally. Uncomment this cell to run the following commands to download the dataset or follow the instructions in the paper to access the data.

PRO TIP: To uncomment a full cell, use `Ctrl + A` to select all lines then `Ctrl + /` to quick uncomment all selected lines.

In []:

```
# ! wget ftp://gdo148.ucllnl.org/multires-random-forest.tgz
# ! tar -xvzf multires-random-forest.tgz
```

```
--2023-01-23 11:02:35-- ftp://gdo148.ucllnl.org/multires-random-forest.tgz
      => 'multires-random-forest.tgz'
Resolving gdo148.ucllnl.org (gdo148.ucllnl.org)... 192.12.137.151
Connecting to gdo148.ucllnl.org (gdo148.ucllnl.org)|192.12.137.151|:21... connected.
Logging in as anonymous ... Logged in!
==> SYST ... done.      ==> PWD ... done.
==> TYPE I ... done.   ==> CWD not needed.
==> SIZE multires-random-forest.tgz ... 983884
```

==> PASV ... done. ==> RETR multires-random-forest.tgz ... done.
Length: 983884 (961K) (unauthoritative)

multires-random-for 100%[=====>] 960.82K --.-KB/s in 0.1s

2023-01-23 11:02:35 (9.11 MB/s) - 'multires-random-forest.tgz' saved [983884]

x multi-resolution/
x multi-resolution/README.pdf
x multi-resolution/fnet_files/
x multi-resolution/fnet_files/energy_balance_random_forest.pkl
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_01.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_02.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_03.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_04.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_05.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_06.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_07.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_08.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_09.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_10.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_100.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_101.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_102.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_103.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_104.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_105.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_106.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_107.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_108.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_109.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_11.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_110.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_111.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_112.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_113.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_114.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_115.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_116.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_117.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_118.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_119.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_12.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_120.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_121.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_122.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_123.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_124.npy
x multi-resolution/fnet_files/energy_balance_random_forest.pkl_125.npy


```
x multi-resolution/p_files/precipitation_random_forest.pkl_77.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_78.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_79.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_80.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_81.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_82.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_83.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_84.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_85.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_86.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_87.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_88.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_89.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_90.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_91.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_92.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_93.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_94.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_95.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_96.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_97.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_98.npy
x multi-resolution/p_files/precipitation_random_forest.pkl_99.npy
x multi-resolution/multiresolution_params_and_metrics.npy
x multi-resolution/predict.py
```

Loading the needed packages and data

Now that we've downloaded the data, let's first import some needed packages that we can use such as numpy, pandas, matplotlib and useful functions from sklearn. We'll also need to read the data into a useable format.

```
In [ ]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.ensemble import (
    ExtraTreesRegressor, GradientBoostingRegressor, RandomForestRegressor
)
from sklearn.neural_network import MLPRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import explained_variance_score
```

```
In [ ]: features = np.load('multi-resolution/multiresolution_params_and_metrics.npy')
```

```
In [ ]: df = pd.DataFrame(
    features,
    columns=['cldfrc_rhminl', 'micromg_dcs', 'eddydiff_a2l', 'uwshcu_criqc', 'uwshcu_rkm',
            'uwshcu_rpen', 'zmconv_alfa', 'zmconv_c0_ocn', 'zmconv_dmpdz', 'zmconv_ke',
            'zmconv_tau', 'resolution', 'TOA_energy_flux', 'precipitation']
)
```

```
In [ ]: df
```

```
Out[ ]:
```

	cldfrc_rhminl	micromg_dcs	eddydiff_a2l	uwshcu_criqc	uwshcu_rkm	uwshcu_rpen	zmconv_alfa	zmconv_c0_ocn	zmconv_dmpdz	zmconv_
0	0.513158	0.750000	0.500000	0.200000	0.750000	1.000000	0.278943	0.826606	0.698970	0.2313
1	0.000000	0.750000	0.500000	0.200000	0.750000	1.000000	0.278943	0.826606	0.698970	0.2313
2	1.000000	0.750000	0.500000	0.200000	0.750000	1.000000	0.278943	0.826606	0.698970	0.2313
3	0.513158	0.000000	0.500000	0.200000	0.750000	1.000000	0.278943	0.826606	0.698970	0.2313
4	0.513158	1.000000	0.500000	0.200000	0.750000	1.000000	0.278943	0.826606	0.698970	0.2313
...
901	0.841311	0.550732	0.862211	0.771593	0.626219	0.712388	0.473620	0.937692	0.952484	0.7954
902	0.420755	0.150635	0.879170	0.933833	0.374344	0.412366	0.253196	0.537597	0.920912	0.7194
903	0.016315	0.200030	0.430349	0.025959	0.347203	0.921252	0.008454	0.543158	0.114437	0.5660
904	0.522297	0.052848	0.662738	0.006282	0.277326	0.519419	0.704073	0.868487	0.069545	0.2549
905	0.876733	0.351345	0.615204	0.457433	0.027828	0.303024	0.687878	0.609783	0.016450	0.0720

906 rows × 14 columns



Looking at the data

Overall, the data has 906 rows and 14 columns. The columns correspond with features and the rows correspond with data points collected. We will use a subset of the data points to train a predictive model.

First, let's briefly discuss what features are included. We collected these features by running simulations in a climate model repeatedly. We have 11 features that have been identified as important that span five physical schemes related to clouds, cloud microphysics, turbulent

mixing, and deep and shallow convection. Below is a table describing the features and their expert-provided ranges of variation:

Description of features:

Parameter	Scheme	Description	Range
cldfrc_rhminl	Cloud processes	Humidity threshold for low clouds	[0.8, 0.99]
micromg_dcs	Microphysics	Autoconversion size threshold (ice to snow)	[1.0E-4, 5.0E-4]
eddydiff_a2l	Eddy diffusion	Moist entrainment enhancement	[10.0, 50.0]
uwshcu_criqc	Shallow convection	Maximum updraft condensate	[0.0005, 0.0015]
uwshcu_rkm	Shallow convection	Fractional updraft mixing efficiency	[8.0, 16.0]
uwshcu_rpen	Shallow convection	Penetrative entrainment efficiency	[1.0, 10.0]
zmconv_alfa	Deep convection	Initial cloud downdraft mass flux	[0.05, 0.6]
zmconv_c0_ocn	Deep convection	Precipitation efficiency over ocean	[0.001, 0.1]
zmconv_dmpdz	Deep convection	Mass entrainment rate	[0.0002, 0.002]
zmconv_ke	Deep convection	Evaporation efficiency	[0.5E-6, 10.0E-6]
zmconv_tau	Deep convection	Convective timescale	[300, 28800]

Splitting the data

For this workshop, let's use the first 12 features (the 11 features listed in the table above plus an added feature for horizontal resolution of the climate models.) Let's also select a value to predict. We can start by predicting the global annual mean net TOA energy flux, which describes the extent to which the Earth's energy budget is in balance. This is an important climate quantity for model developers. We can label these X and Y, respectively.

```
In [ ]: X = features[:, :12]
        Y = features[:, 12]
```

This will be a regression problem, since we want to predict the value of TOA energy flux.

Next, we'll want to split these arrays into a training set and a testing set. The training set will be used to train the model and the testing set will be used to evaluate the model's performance. `sklearn` has a handy function called `train_test_split` which will split the arrays provided into a "train" and a "test" set based on the percentage of data indicated for the test set.

```
In [ ]: x_train, x_test, y_train, y_test = train_test_split(X, Y, test_size=0.2, random_state=12345)
```

Great, now that we've split the dataset...

Train your own predictive model!

In the paper, the model presented is an Extra Trees Regressor and achieved a goodness-of-fit metric (R^2) of **0.86**. Can you beat this?

You can use the helper function below to also help you visualize how well your model is fitting the data. It will generate a plot of the predictions vs. ground truth.

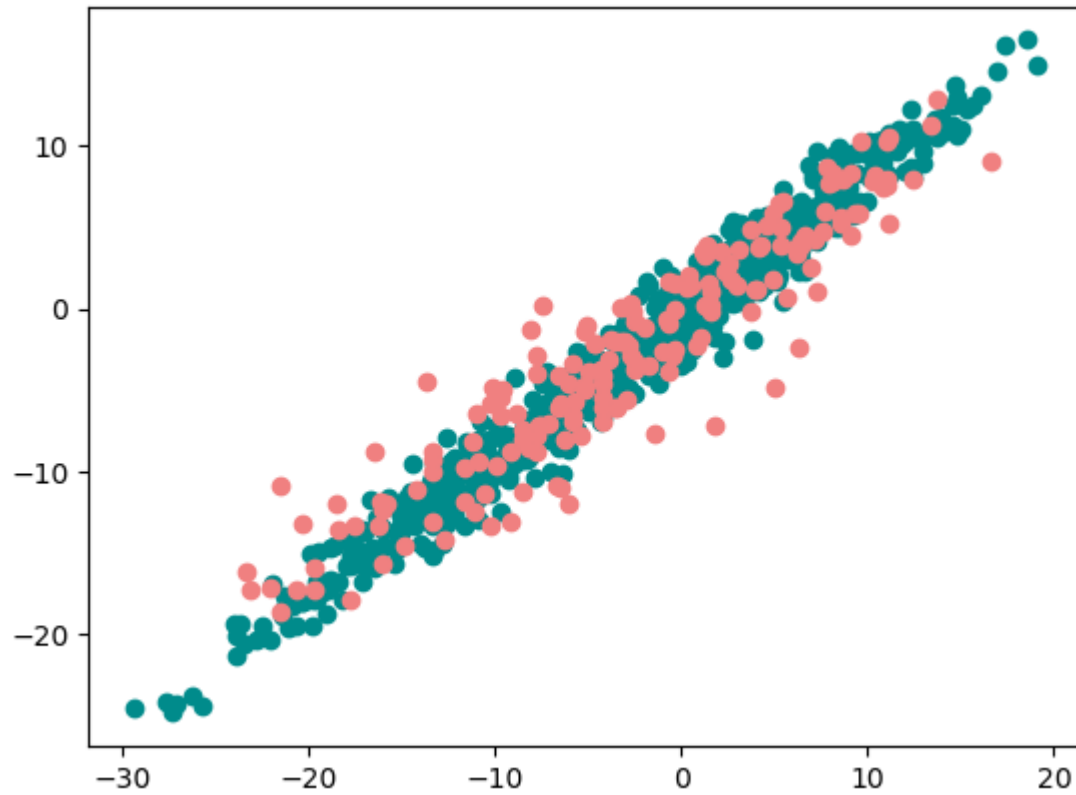
```
In [ ]: def show_plot(train_pred, train_actual, test_pred, test_actual):
        plt.scatter(train_actual, train_pred, color='darkcyan')
        plt.scatter(test_actual, test_pred, color='lightcoral')
        plt.show()
        plt.close()
```

We've included some codes below to help you get started. The first cell contains the code used to train the model presented in the paper. The next few are some thoughts from Gemma and Cindy for methods that might achieve similar (or better) performance. Can you think of any other methods that can be used on this data?

PRO TIP: To uncomment a full cell, use `Ctrl + A` to select all lines then `Ctrl + /` to quick uncomment all selected lines.

```
In [ ]: clf = ExtraTreesRegressor(n_estimators=50, max_depth=8, n_jobs=4, random_state=12345, bootstrap=True, oob_score=True)
        clf.fit(x_train, y_train)
        print(f'The goodness-of-fit is {clf.score(x_test, y_test)}.')
        test_preds = clf.predict(x_test)
        train_preds = clf.predict(x_train)
        show_plot(train_preds, y_train, test_preds, y_test)
```

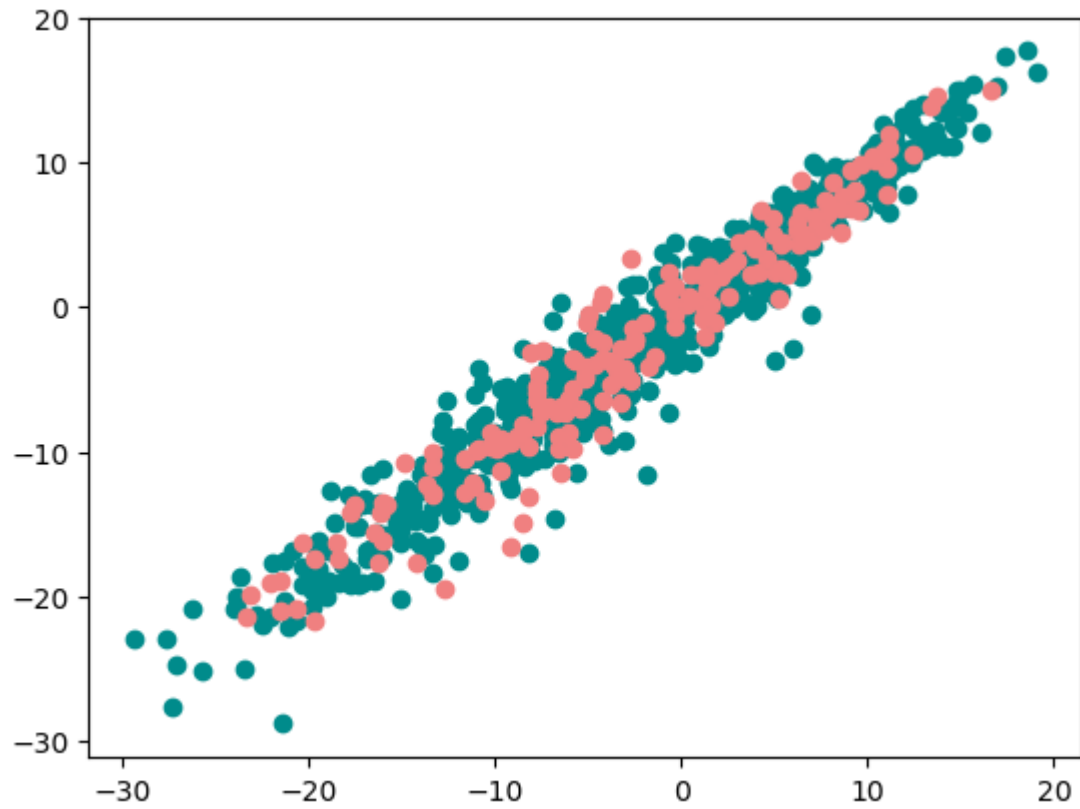
The goodness-of-fit is 0.8586828355957432.



```
In [ ]: clf = MLPRegressor()  
clf.fit(x_train, y_train)  
print(f'The goodness-of-fit is {clf.score(x_test, y_test)}.')  
test_preds = clf.predict(x_test)  
train_preds = clf.predict(x_train)  
show_plot(train_preds, y_train, test_preds, y_test)
```

The goodness-of-fit is 0.9364659626155306.

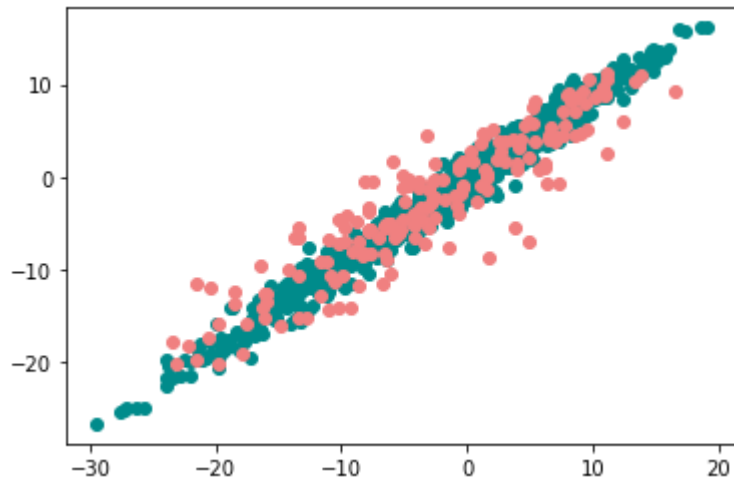
/Users/anderson276/miniconda3/lib/python3.9/site-packages/sklearn/neural_network/_multilayer_perceptron.py:679: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.
warnings.warn(



In []:

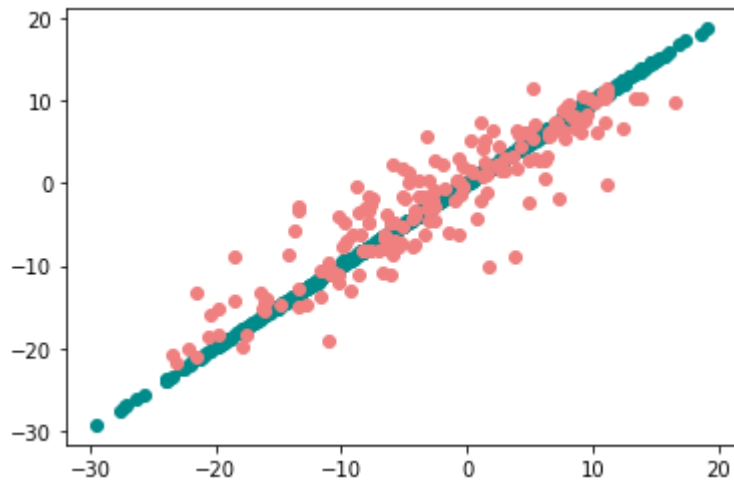
```
# clf = RandomForestRegressor(n_estimators=50, max_depth=8, random_state=12345)
# clf.fit(x_train, y_train)
# print(f'The goodness-of-fit is {clf.score(x_test, y_test)}.')
# test_preds = clf.predict(x_test)
# train_preds = clf.predict(x_train)
# show_plot(train_preds, y_train, test_preds, y_test)
```

The goodness-of-fit is 0.8341147640130518.



```
In [ ]: # clf = GradientBoostingRegressor(n_estimators=50, max_depth=8, random_state=12345)
# clf.fit(x_train, y_train)
# print(f'The goodness-of-fit is {clf.score(x_test, y_test)}.')
# test_preds = clf.predict(x_test)
# train_preds = clf.predict(x_train)
# show_plot(train_preds, y_train, test_preds, y_test)
```

The goodness-of-fit is 0.8158775609361482.



Want to make these models better? Consider hyperparameter tuning! Each of the methods above have keyword arguments for their hyperparameters that you can change the value of to see how the model's performance improves or declines. Consider looking over the source documents to learn more about the hyperparameters you can tune:

- Extra Trees Regressor
- MLP Regressor
- Random Forest Regressor
- Gradient Boosting Regressor

IM #: LLNL-PRES-832485