

Implementing Dataset Enhancements on the THREDDS Data Server

Unidata 2023 Summer Internship

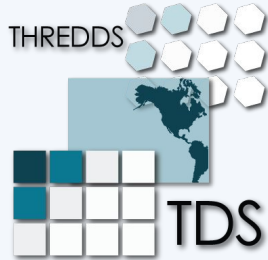
Jessica Souza

Mentor: Tara Drwenski





Background



THREDDS Data Server¹ is a web server that provides catalog, metadata and data access for real-time and archived datasets of environmental data sources at a number of distributed server sites, using a variety of remote data access protocols.

THREDDS Data Server

Welcome to THREDDS Data Server top-level TDS Catalog.
Hosted by Unidata.

Catalog

Dataset	Size	Last Modified
Realtime data from IDD		--
Forecast Model Data		--
Forecast Products and Analyses		--
Radar Data		--
Satellite Data		--
Text Products		--
Other Unidata Data		--
Unidata case studies		--



Motivation

There is an expressive use of machine learning methods in earth sciences research.

AI / ML targeted dataset enhancements

Data preprocessing steps in machine learning generally involves cleaning, **rescaling** and splitting the data.

The goal of rescaling is to transform features to be on a similar range.

This improves the performance and training stability of the model.

Standardizer and Normalizer



Motivation

Standardization (Z- score, Standard Scaler)

$$s = \frac{z - \mu}{\sigma}$$

z	Data point
μ	Mean value in the variable
σ	Standard deviation value in the variable
s	Standardized data point

Obtain a random variable s with mean 0 and standard deviation 1.



Motivation

Normalization (Min-Max Scalar)

$$n = \frac{z - z_{\min}}{z_{\max} - z_{\min}}$$

z Data point

n Normalized data point

Obtain values between 0 and 1.

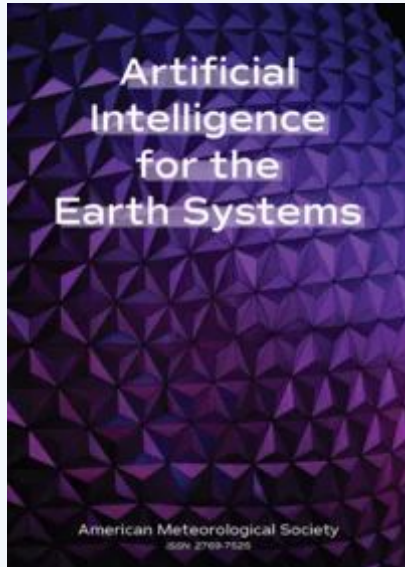
Sensitive to outliers.



Motivation

AMS Journal

Artificial Intelligence for the Earth Systems²



Standardization

Papers

1. Yu et al., 2022
2. White et al., 2022
3. Mamalakis et al., 2022
4. Liu et al 2022
5. Li et al., 2022
6. Miralles et al., 2022
7. Galea et al., 2023
8. Straaten et al., 2023
9. Fulton et al., 2023
10. Connolly et al., 2023

ML

- Regression model
- Neural Network
- Generative Adversarial Network

Dataset

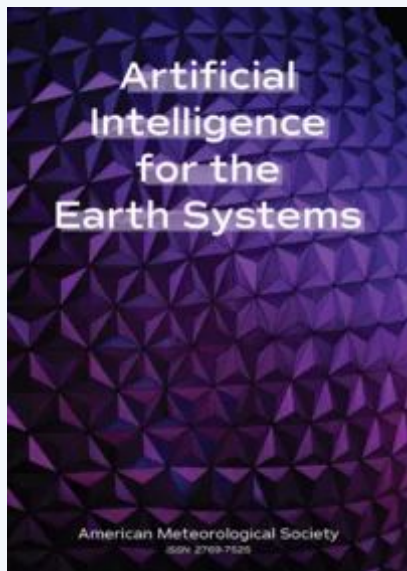
- GOES16 (ABI channels)
- Suomi NPP (VIIRS channels)
- Satellite Altimeter
- ERA-interim Reanalysis



Motivation

AMS Journal

Artificial Intelligence for the Earth Systems²



Normalization

Papers

1. Cheung et al., 2022
2. Osborne et al., 2022
3. Chen et al., 2023

ML

→ Neural Network

Dataset

- GEFS
- Radar
- Met Office Forecast

sklearn.preprocessing: Preprocessing and Normalization



The `sklearn.preprocessing` module includes scaling, centering, normalization, binarization methods.

User guide: See the [Preprocessing data](#) section for further details.

<code>preprocessing.Binarizer(*[, threshold, copy])</code>	Binarize data (set feature values to 0 or 1) according to a threshold.
<code>preprocessing.FunctionTransformer([func, ...])</code>	Constructs a transformer from an arbitrary callable.
<code>preprocessing.KBinsDiscretizer([n_bins, ...])</code>	Bin continuous data into intervals.
<code>preprocessing.KernelCenterer()</code>	Center an arbitrary kernel matrix K .
<code>preprocessing.LabelBinarizer(*[, neg_label, ...])</code>	Binarize labels in a one-vs-all fashion.
<code>preprocessing.LabelEncoder()</code>	Encode target labels with value between 0 and <code>n_classes-1</code> .
<code>preprocessing.MultiLabelBinarizer(*[, ...])</code>	Transform between iterable of iterables and a multilabel format.
<code>preprocessing.MaxAbsScaler(*[, copy])</code>	Scale each feature by its maximum absolute value.
<code>preprocessing.MinMaxScaler([feature_range, ...])</code>	Transform features by scaling each feature to a given range.
<code>preprocessing.Normalizer([norm, copy])</code>	Normalize samples individually to unit norm.
<code>preprocessing.OneHotEncoder(*[, categories, ...])</code>	Encode categorical features as a one-hot numeric array.
<code>preprocessing.OrdinalEncoder(*[, ...])</code>	Encode categorical features as an integer array.
<code>preprocessing.PolynomialFeatures([degree, ...])</code>	Generate polynomial and interaction features.
<code>preprocessing.PowerTransformer([method, ...])</code>	Apply a power transform featurewise to make data more Gaussian-like.
<code>preprocessing.QuantileTransformer(*[, ...])</code>	Transform features using quantiles information.
<code>preprocessing.RobustScaler(*[, ...])</code>	Scale features using statistics that are robust to outliers.
<code>preprocessing.SplineTransformer([n_knots, ...])</code>	Generate univariate B-spline bases for features.
<code>preprocessing.StandardScaler(*[, copy, ...])</code>	Standardize features by removing the mean and scaling to unit variance.
<code>preprocessing.TargetEncoder([categories, ...])</code>	Target Encoder for regression and classification targets.

<code>preprocessing.add_dummy_feature(X[, value])</code>	Augment dataset with an additional dummy feature.
<code>preprocessing.binarize(X, *[, threshold, copy])</code>	Boolean thresholding of array-like or <code>scipy.sparse</code> matrix.
<code>preprocessing.label_binarize(y, *, classes)</code>	Binarize labels in a one-vs-all fashion.
<code>preprocessing.maxabs_scale(X, *[, axis, copy])</code>	Scale each feature to the <code>[-1, 1]</code> range without breaking the sparsity.
<code>preprocessing.minmax_scale(X[, ...])</code>	Transform features by scaling each feature to a given range.
<code>preprocessing.normalize(X[, norm, axis, ...])</code>	Scale input vectors individually to unit norm (vector length).
<code>preprocessing.quantile_transform(X, *[, ...])</code>	Transform features using quantiles information.
<code>preprocessing.robust_scale(X, *[, axis, ...])</code>	Standardize a dataset along any axis.
<code>preprocessing.scale(X, *[, axis, with_mean, ...])</code>	Standardize a dataset along any axis.
<code>preprocessing.power_transform(X[, method, ...])</code>	Parametric, monotonic transformation to make data more Gaussian-like.

Python Scikit-learn library³

Python Scikit-learn library³

sklearn.preprocessing: Preprocessing and Normalization

The `sklearn.preprocessing` module includes scaling, centering, normalization, binarization methods.

User guide: See the [Preprocessing data](#) section for further details.

<code>preprocessing.Binarizer(*[, threshold, copy])</code>	Binarize data (set feature values to 0 or 1) according to a threshold.
<code>preprocessing.FunctionTransformer((func, ...))</code>	Constructs a transformer from an arbitrary callable.
<code>preprocessing.KBinsDiscretizer([n_bins, ...])</code>	Bin continuous data into intervals.
<code>preprocessing.KernelCenterer()</code>	Center an arbitrary kernel matrix K .
<code>preprocessing.LabelBinarizer(*[, neg_label, ...])</code>	Binarize labels in a one-vs-all fashion.
<code>preprocessing.LabelEncoder()</code>	Encode target labels with value between 0 and <code>n_classes-1</code> .
<code>preprocessing.MultiLabelBinarizer(*[, ...])</code>	Transform between iterable of iterables and a multilabel format.
<code>preprocessing.MaxAbsScaler(*[, copy])</code>	Scale each feature by its maximum absolute value.
<code>preprocessing.MinMaxScaler((feature_range, ...))</code>	Transform features by scaling each feature to a given range.
<code>preprocessing.Normalizer([norm, copy])</code>	Normalize samples individually to unit norm.
<code>preprocessing.OneHotEncoder(*[, categories, ...])</code>	Encode categorical features as a one-hot numeric array.
<code>preprocessing.OrdinalEncoder(*[, ...])</code>	Encode categorical features as an integer array.
<code>preprocessing.PolynomialFeatures((degree, ...))</code>	Generate polynomial and interaction features.
<code>preprocessing.PowerTransformer([method, ...])</code>	Apply a power transform featurewise to make data more Gaussian-like.
<code>preprocessing.QuantileTransformer(*[, ...])</code>	Transform features using quantiles information.
<code>preprocessing.RobustScaler(*[, ...])</code>	Scale features using statistics that are robust to outliers.
<code>preprocessing.SplineTransformer([n_knots, ...])</code>	Generate univariate B-spline bases for features.
<code>preprocessing.StandardScaler(*[, copy, ...])</code>	Standardize features by removing the mean and scaling to unit variance.
<code>preprocessing.TargetEncoder([categories, ...])</code>	Target Encoder for regression and classification targets.

<code>preprocessing.add_dummy_feature(X[, value])</code>	Augment dataset with an additional dummy feature.
<code>preprocessing.binarize(X, *[, threshold, copy])</code>	Boolean thresholding of array-like or <code>scipy.sparse</code> matrix.
<code>preprocessing.label_binarize(y, *, classes)</code>	Binarize labels in a one-vs-all fashion.
<code>preprocessing.maxabs_scale(X, *[, axis, copy])</code>	Scale each feature to the <code>[-1, 1]</code> range without breaking the sparsity.
<code>preprocessing.minmax_scale(X[, ...])</code>	Transform features by scaling each feature to a given range.
<code>preprocessing.normalize(X[, norm, axis, ...])</code>	Scale input vectors individually to unit norm (vector length).
<code>preprocessing.quantile_transform(X, *[, ...])</code>	Transform features using quantiles information.
<code>preprocessing.robust_scale(X, *[, axis, ...])</code>	Standardize a dataset along any axis.
<code>preprocessing.scale(X, *[, axis, with_mean, ...])</code>	Standardize a dataset along any axis.
<code>preprocessing.power_transform(X[, method, ...])</code>	Parametric, monotonic transformation to make data more Gaussian-like.





Code

External library implementations (Apache Commons Mathematics Library 4)

```
import org.apache.commons.math.stat.descriptive.SummaryStatistics;
```

Computes summary statistics for very large data streams.

Data values are not stored in memory.

```
public class Standardizer
```

```
statistics.getMean()
```

```
statistics.getStandardDeviation()
```

```
public class Normalizer
```

```
statistics.getMin()
```

```
statistics.getMax() - statistics.getMin()
```



Code

Integrating with netcdf-java

- CDM.java
 - ◆ Create constants / attributes in the Common Data Model class of constants.

- NetcdfDataset.java
 - ◆ Include Standardizer / Normalizer in the set of enhancements (ApplyStandardizer, ApplyNormalizer)

- VariableDS.java
 - ◆ Check if attribute "standardizer" "normalizer" is present and data is floating point
 - ◆ Apply enhancements to the data



Code

Using it in the TDS

The preprocessing can be achieved in the THREDDS catalog.xml through NetCDF Markup Language (NcML)⁵.

NcML creates a virtual dataset without changing the original data.

```
<netcdf xmlns="http://www.unidata.ucar.edu/namespaces/netcdf/ncml-2.2" enhance="all">  
  <variable name="scaledVar">  
    <attribute name="standardize"/>  
  </variable>  
</netcdf>
```

```
<netcdf xmlns="http://www.unidata.ucar.edu/namespaces/netcdf/ncml-2.2" enhance="all">  
  <variable name="scaledVar">  
    <attribute name="normalize"/>  
  </variable>  
</netcdf>
```



Testing

- Test Convert
 - ◆ Doubles, Floats
 - ◆ **NaNs, Equal values**

Standardizer

- Test Calculate Mean
 - ◆ Doubles, Floats
- Test Calculate Standard Deviation
 - ◆ Doubles, Floats

Normalizer

- Test Calculate Minimum
 - ◆ Doubles, Floats
- Test Calculate Range
 - ◆ Doubles, Floats



Testing

- Read variables from NcML file
- Doubles, Floats
- Ints
- Equal values

```
@Test
public void testEnhanceStandardizer() throws IOException {
    try (NetcdfFile ncfile = NetcdfDatasets.openDataset( location: dataDir + "testStandardizer.ncml", enhance: true, cancelTask: null)) {
        Variable doubleVar = ncfile.findVariable( fullNameEscaped: "doublevar");
        assertNotNull(doubleVar);
        assertEquals(doubleVar.getDataType(), DataType.DOUBLE);
        assertTrue(doubleVar.attributes().hasAttribute( attName: "standardize"));
        Array dataDoubles = doubleVar.read();
        assertTrue(nearlyEquals(dataDoubles, DATA_DOUBLES));

        Variable sameDoubleVar = ncfile.findVariable( fullNameEscaped: "samedoublevar");
        assertNotNull(sameDoubleVar);
        assertEquals(sameDoubleVar.getDataType(), DataType.DOUBLE);
        assertTrue(sameDoubleVar.attributes().hasAttribute( attName: "standardize"));
        Array dataSameDoubles = sameDoubleVar.read();
        assertTrue(nearlyEquals(dataSameDoubles, DATA_SAMEDOUBLES)); // The enhancement doesn't apply if all the
                                                                    // values are the equal, so it returns the
                                                                    // same data


        Variable floatVar = ncfile.findVariable( fullNameEscaped: "floatvar");
        assertNotNull(floatVar);
        assertEquals(floatVar.getDataType(), DataType.FLOAT);
        assertTrue(floatVar.attributes().hasAttribute( attName: "standardize"));
        Array dataFloats = floatVar.read();
        assertTrue(nearlyEquals(dataFloats, DATA_FLOATS));

        Variable intVar = ncfile.findVariable( fullNameEscaped: "intvar");
        assertNotNull(intVar);
        assertEquals(intVar.getDataType(), DataType.INT);
        assertTrue(intVar.attributes().hasAttribute( attName: "standardize"));
        Array data = intVar.read();
        assertTrue(nearlyEquals(data, DATA_INTS)); // The enhancement doesn't apply to ints, so the data should
                                                                    // be equal to the input array
    }
}
```









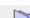


Results

Thredds-test server






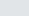
 THREDDS Data Server

Welcome to THREDDS Data Server top-level TDS Catalog.
Hosted by Unidata.

Catalog

Dataset	Size	Last Modified
 Realtime data from IDD		--
 Forecast Model Data		--
 Forecast Products and Analyses		
 Radar Data		
 Satellite Data		
 Text Products		
 Other Unidata Data		
 Unidata case studies		
 Preprocessed data		

Catalog

Dataset	Size	Last Modified
 Standardized data		--
 Standardized GOES West Products / CloudTopTemperature		--
 Standardized GFS Quarter Degree Analysis		--
 Normalized data		--
 Normalized GOES West Products / CloudTopTemperature		--
 Normalized GFS Quarter Degree Analysis		--



Results

Jupyter Notebooks

THREDDS Data Server


Access Preview

Viewers:

Viewer	Type	Description
Godiva3	Browser	
default_viewer.ipynb	Jupyter Notebook	The TDS default viewer attempts to plot any Variable contained in the Dataset.
preprocessed_dataset.ipynb	Jupyter Notebook	Preprocessed Dataset Visualization

Siphon THREDDS Jupyter Notebook - Visualizing Preprocessed Data

Dataset: GFS_Global_0p25deg_ana_20230725_1200.grib2

Powered By 

```
import xarray as xr
from siphon.catalog import TDSCatalog
```

[1] ✓ 1.9s Python

Standardized Data

z: Data point
μ: Mean value in the variable
σ: Standard deviation value in the variable
s: Standardized data point

$$s = \frac{z - \mu}{\sigma}$$



Results

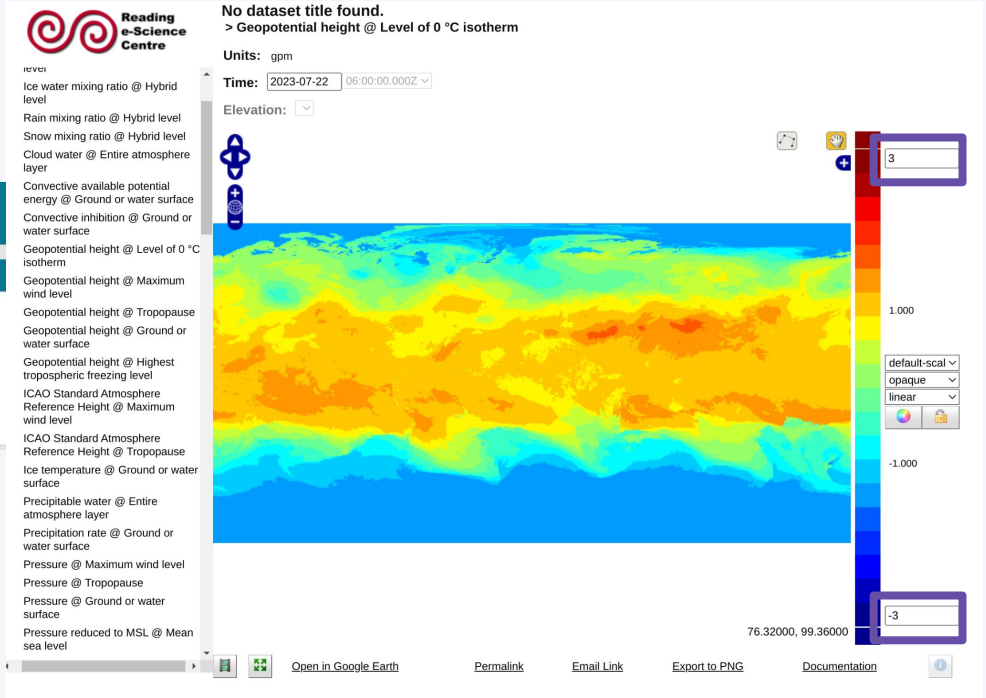
Godiva3

THREDDS Data Server

Access Preview

Viewers:

Viewer	Type	Description
Godiva3	Browser	
default_viewer.ipynb	Jupyter Notebook	The TDS default viewer attempts to plot any Variable contained in the Dataset.
preprocessed_dataset.ipynb	Jupyter Notebook	Preprocessed Dataset Visualization

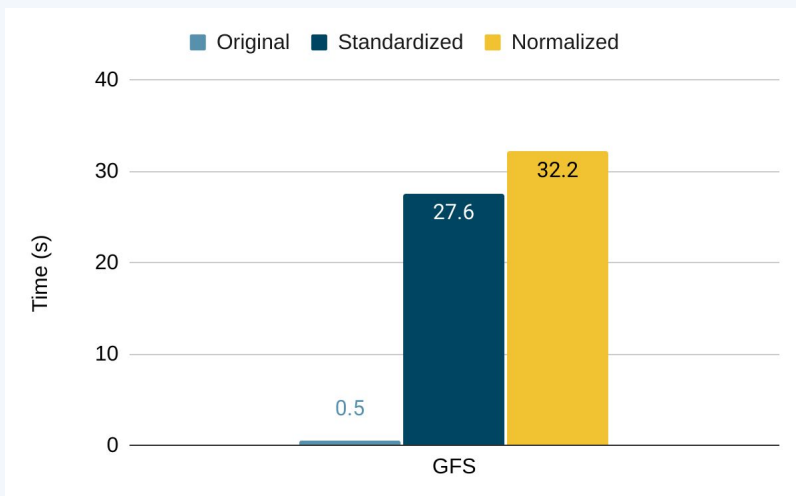




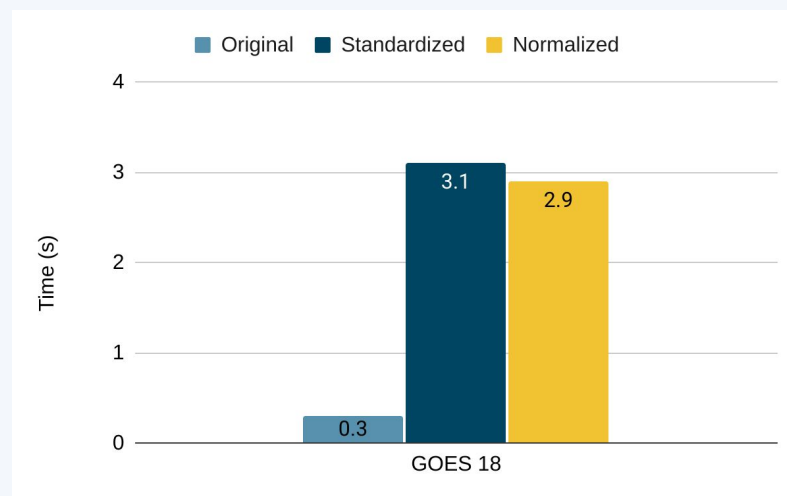
Results

Performance test - Apache HTTP server benchmarking tool⁶
Simple load tests on thredds-test server

GFS 0.25 Degree - Forecast Model Data
+80 variables



GOES 18 Product - Full Conus
1 variable





Final Remarks

Dataset enhancements targeting machine learning applications.

Two most common types of rescaling data as part of preprocessing step.

Available original + standardized & normalized dataset on thredds-test server.

Future work:

- Improvements on performance
 - ◆ Caching for example
- Provide more datasets relevant to the users.

References

¹ THREDDS Data Server

<https://www.unidata.ucar.edu/software/tds/>

² Artificial Intelligence for the Earth Systems

<https://www.ametsoc.org/index.cfm/ams/publications/journals/artificial-intelligence-for-the-earth-systems/>

³ Scikit-learn library

<https://scikit-learn.org/>

⁴ Commons Math: The Apache Commons Mathematics Library

<https://commons.apache.org/proper/commons-math/>

⁵ NetCDF Markup Language (NcML)

https://docs.unidata.ucar.edu/netcdf-java/current/userguide/ncml_overview.html

⁶ Apache HTTP server benchmarking tool

<https://httpd.apache.org/docs/2.4/programs/ab.html>



Acknowledgements

