# dtwhaclustering

Release 1.0

**Utpal Kumar** 

# **USAGE**

1	Installation	1
	1.1 Requirements	1
	1.2 Install from PyPI	1
2	Dynamic Time Warping based Hierarchical Agglomerative Clustering	3
3	dtwhaclustering.analysis_support	5
4	dtwhaclustering.plot_linear_trend	7
5	dtwhaclustering.leastSquareModeling	9
6	dtwhaclustering.dtw_analysis	11
7	dtwhaclustering.plot_stations	15
8	Indices and tables	17
Ру	ython Module Index	19
In	ndex	21

## **ONE**

## **INSTALLATION**

# 1.1 Requirements

1. dtaidistance: For computing DTW distance

2. pygmt: For plotting high-resolution maps

3. pandas: Analyze tabular data

4. numpy: Computation

5. matplotlib: Plotting time series

6. scipy: Interpolating data

7. xarray: Multilayered data structure

# 1.2 Install from PyPI

This package is available on PyPI (requires Python 3):

pip install dtwhaclustering

**TWO** 

# DYNAMIC TIME WARPING BASED HIERARCHICAL AGGLOMERATIVE CLUSTERING

Codes to perform Dynamic Time Warping Based Hierarchical Agglomerative Clustering of GPS data

author Utpal Kumar

date 2021/08

copyright 2021, Institute of Earth Sciences, Academia Sinica.



### THREE

# DTWHACLUSTERING.ANALYSIS\_SUPPORT

DTW HAC analysis support functions (analysis\_support)

author Utpal Kumar, Institute of Earth Sciences, Academia Sinica

dtwhaclustering.analysis\_support.dec2dt(start)

Convert the decimal type time array to the date-time type array

Parameters start (list) – list or numpy array of decimal year values e.g., [2020.001]

**Returns** date-time type array

Return type list

dtwhaclustering.analysis\_support.dec2dt\_scalar(st)

Convert the decimal type time value to the date-time type

Parameters st – scalar decimal year value e.g., 2020.001

**Returns** time as datetime type

Return type str

dtwhaclustering.analysis\_support.toYearFraction(date)

Convert the date-time type object to decimal year

**Parameters date** – the date-time type object

Returns decimal year

Return type float

**FOUR** 

## DTWHACLUSTERING.PLOT LINEAR TREND

Plot linear trend values on a geographical map

author Utpal Kumar, Institute of Earth Sciences, Academia Sinica

 $\label{linear_trend.compute_interpolation} df, \textit{method='nearest'}, \textit{lonrange=}(120.0, 122.0), \textit{latrange=}(21.8, 25.6), \textit{step=}0.01)$ 

Interpolate linear trend values using Scipy's griddata and returns xarray object

#### **Parameters**

- **df** (pandas.DataFrame) pandas dataframe containing the columns lon, lat and slope
- **method** (*str*) Method of interpolation. One of {'linear', 'nearest', 'cubic'}. For more, see *scipy.interpolate.griddata*
- lonrange (tuple) minimum and maximum values of the longitude to be interpolated
- latrange (tuple) minimum and maximum values of the latitude to be interpolated
- **step** (*float*) stepsize to interpolate data spatially

**Returns** *xarray.DataArray* of dims, and coords ("lat", "long"), maximum of the absolute interpolated array values

```
\label{linear_trend_on_map} \begin{tabular}{llll} dtwhaclustering.plot_linear_trend_on_map(df, maplonrange=(120.0, 122.1), \\ maplatrange=(21.8, 25.6), \\ intrp_lonrange=(120.0, 122.0), \\ intrp_latrange=(21.8, 25.6), \\ outfig='Maps/slope-plot.png', \\ frame=['alf0.25', 'WSen'], cmap='jet', \\ step=0.01, stn_labels=None, \\ justify='left', \\ labelfont='6p, Helvetica-Bold, black', \\ offset='5p/-5p', markerstyle='il0p', \\ defpen='lp, black', \\ stn_labels\_color='black', \\ rand\_justify=False, \\ \end{tabular}
```

Plot the interpolated linear trend values along with the original data points on a geographical map using PyGMT

#### **Parameters**

• **df** (pandas.DataFrame) – Pandas dataframe containing the columns lon, lat and slope

water\_color='skyblue')

- maplonrange (tuple) longitude min/max of the output map
- maplatrange (tuple) latitude min/max of the output map

- intrp\_lonrange (tuple) longitude min/max for the interpolation of data
- intrp\_latrange (tuple) latitude min/max for the interpolation of data
- **step** (*float*) resolution of the interpolation
- **cmap** (*str*) colormap for the output map
- frame (list) frame of the output map. See PyGMT docs for details
- **outfig** (str) output figure name with extension, e.g., slope-plot.png
- water\_color (str) color of the water, default: skyblue
- stn\_labels (array) label the station names. Only the stations provided will be labeled
- **justify** (*str*) justification of the station labels 'left', 'right'
- rand\_justify (boolean) randomly decide the justification for each labels
- markerstyle (str) marker style and size

Returns None

## DTWHACLUSTERING.LEASTSQUAREMODELING

Least square modeling of GPS displacements for seasonality, tidal, co-seismic jumps.

author Utpal Kumar, Institute of Earth Sciences, Academia Sinica

class dtwhaclustering.leastSquareModeling.LSQmodules(dUU, sel\_eq\_file=None,

station\_loc\_file='helper\_files/stn\_loc.txt', comp='U', figdir='LSQOut', periods=(13.6608, 14.7653, 27.5546, 182.62, 365.26, 6793.836))

Bases: object

Compute the least-squares model using multithreading

#### **Parameters**

- plot\_results (boolean) plot the final results
- remove\_trend (boolean) return the time series after removing the linear trend
- **remove\_seasonality** (*boolean*) return the time series after removing the seasonal signals
- remove\_jumps (boolean) return the time series after removing the co-seismic jumps
- plotformat (str) plot format of the output figure, e.g. "png". "pdf" by default.

jump(t, t0)

heaviside step function

#### **Parameters**

- t (list) time data
- **t0** earthquake origin time

 $\label{ling:leastSquareModeling.lsqmodeling} \begin{subarray}{l} dtwhaclustering.leastSquareModeling.lsqmodeling($dUU$, $dNN$, $dEE$, $stnlocfile$, $plot_results=True$, $remove\_trend=False$, $remove\_seasonality=True$, $remove\_jumps=False$, $sel\_eq\_file='helper\_files/selected\_eqs\_new.txt'$, $true, $tr$ 

figdir='LSQOut')

Least square modeling for the three component time series

- **dUU** (pandas.DataFrame) Vertical component pandas dataframe time series
- **dNN** (pandas. DataFrame) North component pandas dataframe time series
- **dEE** (*pandas.DataFrame*) East component pandas dataframe time series

- plot\_results (boolean) plot the final results
- remove\_trend (boolean) return the time series after removing the linear trend
- **remove\_seasonality** (*boolean*) return the time series after removing the seasonal signals
- **remove\_jumps** (*boolean*) return the time series after removing the co-seismic jumps
- **sel\_eq\_file** (*str*) File containing the earthquake origin times e.g., 2009,1,15 with header info, e.g. *year\_val,month\_val,date\_val*
- **stnlocfile** (*str*) File containing the station location info, e.g., *DAWU*,120.89004,22.34059, with header *stn*,lon,lat

**Returns** Pandas dataframe corresponding to the vertical, north and east components e.g., final\_dU, final\_dN, final\_dE

Return type pandas.DataFrame, pandas.DataFrame, pandas.DataFrame

## DTWHACLUSTERING.DTW ANALYSIS

Classes and functions for the DTW analysis and plotting maps and figures (*dtw\_analysis*) This module is built around the dtaidistance package for the DTW computation and scipy.cluster

author Utpal Kumar, Institute of Earth Sciences, Academia Sinica

note See dtaidistance for details on HierarchicalTree, dtw, dtw\_visualisation

compute\_cluster(clusterMatrix=None, window='constrained', windowfrac=0.25)

Parameters compute\_cluster - data matrix to cluster

Returns model, cluster\_idx

**compute\_cut\_off\_inconsistency** (*t*=None, depth=2, criterion='inconsistent', return\_cluster=False')

Calculate inconsistency statistics on a linkage matrix following *scipy.cluster.hierarchy.inconsistent*. It com-

Calculate inconsistency statistics on a linkage matrix following *scipy.cluster.hierarchy.inconsistent*. It compares each cluster merge's height h to the average avg and normalize it by the standard deviation *std* formed over the depth previous levels

#### **Parameters**

- t (scalar) threshold to apply when forming flat clusters. See scipy.cluster.hierarchy.fcluster for details
- **depth** (*int*) The maximum depth to perform the inconsistency calculation

**Returns** maximum inconsistency coefficient for each non-singleton cluster and its children; the inconsistency matrix (matrix with rows of avg, std, count, inconsistency); cluster

```
compute_dendrogram(color_thresh=None)
```

compute\_distance\_accl(clusterMatrix=None)

compute\_distance\_matrix(compact=True, block=None)

get\_linkage(clusterMatrix=None)

optimum\_cluster\_elbow(minmax=False, plotloc=False)

Gives the optimum number of clusters required to express the maximum difference in the similarity using the elbow method

```
plot_cluster_geomap(dtw_distance='optimal', minlon=None, maxlon=None, minlat=None, maxlat=None, figname='dtw_cluster.pdf', colorbar=True, colorbarstep=1, doffset=1, dpi=720, topo_data='@earth_relief_15s', plot_topo=False, markerstyle='c0.3c', cmap_topo='topo', cmap_data='jet', projection='M4i', topo_cpt_range='-8000/8000/1000', landcolor='#6666666')
```

#### **Parameters**

- **figname** (str) output figure name
- colorbar (bool) plot colorbar

Plot the cluster points on a geographical map

- **colorbarstep** (*int*) step for the colorbar
- cmap (str) colormap for the cluster points
- projection (str) map projection
- topo\_data (str) topographic data resolution, see pygmt docs for details
- topo\_cpt\_range min/max/step for topographic color
- landcolor (str) color for land region
- **dpi** (int) output figure resolution

```
plot_cluster_geomap_interpolated(dtw_distance='optimal', lonrange=(120.0, 122.0), latrange=(21.8, 25.6), gridstep=0.01, figname='dtw_cluster_interp.pdf', minlon=None, maxlon=None, minlat=None, maxlat=None, maxlet=Vone, maxlet=style='c0.3c', dpi=720, doffset=1, plot_data=True, plot_intrp=True)
```

#### **Parameters**

- **dtw\_distance** (*str or float*) use *dtw\_distance* value to obtain the cluster division. If *optimal* then the optimal *dtw\_distance* will be calculated
- **lonrange** (tuple) minimum and maximum of longitude values for interpolation
- latrange (tuple) minimum and maximum of latitude values for interpolation
- **gridstep** (*float*) step size for interpolation

Plot the cluster points in a rectangular coordinate system

- **dtw\_distance** (*str or float*) use *dtw\_distance* value to obtain the cluster division. If *optimal* then the optimal *dtw\_distance* will be calculated
- **figname** (str) output figure name
- colorbar (bool) plot colorbar
- colorbarstep (int) step for the colorbar
- scale (int) figure scale
- markersize cluster points size
- edgecolors marker edge color

• cmap (str) – colormap for the cluster points

Returns figure, axes

plot\_dendrogram(figname=None, figsize=(20, 8), xtickfontsize=26, labelfontsize=26, xlabel='3-D Stations', ylabel='DTW Distance', truncate\_p=None, distance\_threshold=None, annotate above=0, plotpdf=True, leaf rotation=0)

**Parameters** truncate\_p – show only last truncate p out of all merged branches

plot\_hac\_iteration(figname=None, figsize=(10, 8), xtickfontsize=26, labelfontsize=26, xlabel='Iteration #', ylabel='DTW Distance', plot\_color='CO', plotpdf=True)

**Parameters xlim** (list) – x limits of the plot e.g., [1,2]

plot\_polar\_dendrogram(figsize=(20, 20), normfactor=None, Nyticks=7, gap=0.05, Nsmooth=100, linewidth=1, xtickfontsize=20, ytickfontsize=20, plotstyle='seaborn', figname=None, plotpdf=True, gridcolor=None, gridstyle='--', gridwidth=1, tickfontweight='bold', distance\_threshold=None)

Plot polar dendrogram of the clustering result

#### **Parameters**

- **figsize** figure size
- **normfactor** (optional) normalization factor for the log spacing between the yticks, np.log(dcoord+normfactor)
- Nyticks number of y ticks
- gap gap for the yticks
- plotstyle matplotlib plot style, *plotstyle* by default
- **distance\_threshold** str, float. threshold for the coloring of branches. If str="optimal", then the optimal number of clusters based on elbow method will be used

plot\_signals(figname=None, figsize=(10, 6), fontsize=26)

#### **Parameters**

- **figname** (str) output figure name
- **figsize** (tuple) output figure size

compute\_distance(pruning=True, best\_path=False)

Returns the DTW distance

**Parameters pruning** (*boolean*) – prunes computations by setting max\_dist to the Euclidean upper bound

#### **compute\_warping\_path**(windowfrac=None, psi=None, fullmatrix=False)

Returns the DTW path :param windowfrac: Fraction of the signal length. Only allow for shifts up to this amount away from the two diagonals. :param psi: Up to psi number of start and end points of a sequence can be ignored if this would lead to a lower distance :param full\_matrix: The full matrix of all warping paths (or accumulated cost matrix) is built

**plot\_matrix**(windowfrac=0.2, psi=None, figname=None, shownumbers=False, showlegend=True) Plot the signals with the DTW matrix

**Parameters figname** (str) – output figure name

**plot\_signals** (*figname=None*, *figsize=*(12, 6))

Plot the signals

#### Parameters

- **figname** (str) output figure name
- **figsize** (tuple) output figure size

**Returns** figure and axes object

plot\_warping\_path(figname=None, figsize=(12, 6))

Plot the signals with the warping paths

#### **Parameters**

- **figname** (str) output figure name
- **figsize** (tuple) output figure size

 $\verb|dtwhaclustering.dtw\_analysis.noise\_robustness\_test(|clean\_df|, |scale=0.1|)|$ 

 $\label{loss_equation} $$ dtwhaclustering.dtw\_analysis.plot\_cluster(lons, lats, figname=None, figsize=(10, 10), plotpdf=True, \\ labels=[], markersize=20)$ 

#### **Parameters**

- **figname** (str) output figure name
- **figsize** (tuple) output figure size

dtwhaclustering.dtw\_analysis.plot\_signals(matrix, labels=[], figname=None, plotpdf=True, figsize=(10, 8), ylabelsize=26, color=None)

Plot signals

#### **Parameters**

- color list of colors. If None then matplotlib defaults color sequence will be used
- **figname** (str) output figure name
- **figsize** (tuple) output figure size

 $\label{lem:dtwhaclustering.dtw_analysis.shuffle_signals} (\textit{matrix}, labels = [], \textit{plot\_signals} = \textit{False}, \textit{figsize} = (10, 12), \\ \textit{figname} = \textit{None}, \textit{plotpdf} = \textit{True})$ 

- **figname** (str) output figure name
- **figsize** (tuple) output figure size

### SEVEN

# DTWHACLUSTERING.PLOT\_STATIONS

Plot topographic station map

author Utpal Kumar, Institute of Earth Sciences, Academia Sinica

dtwhaclustering.plot\_stations.plot\_station\_map(station\_data, minlon=None, maxlon=None,

minlat=None, maxlat=None, outfig='station\_map.png', datacolor='blue', topo\_data='@earth\_relief\_15s', cmap='etopo1', projection='M4i', datalabel='Stations', markerstyle='i10p', random\_station\_label=False, stn\_labels=None, justify='left', labelfont='6p,Helvetica-Bold,black', offset='5p/-5p', stn\_labels\_color='red', rand\_justify=False)

Plot topographic station map using PyGMT

- station\_data Pandas dataframe containing columns lon, lat
- minlon Minimum longitude of the map (optional)
- maxlon Maximum longitude of the map (optional)
- **minlat** Minimum latitude of the map (optional)
- maxlat Maximum latitude of the map (optional)
- datacolor color of the data point to plot
- topo\_data etopo data file
- **cmap** colormap for the output topography
- **projection** projection of the map. Mercator of width 4 inch by default
- datalabel Label for the data
- markerstyle Style of the marker. Inverted triangle of size 10p by default.
- outfig Output figure path
- random\_station\_label int. label randomly selected random\_station\_label of stations

# **EIGHT**

# **INDICES AND TABLES**

- genindex
- modindex
- search

## **PYTHON MODULE INDEX**

## d

```
dtwhaclustering.1
dtwhaclustering.analysis_support,5
dtwhaclustering.dtw_analysis,11
dtwhaclustering.leastSquareModeling,9
dtwhaclustering.plot_linear_trend,7
dtwhaclustering.plot_stations,15
```

20 Python Module Index

# **INDEX**

C	module, 11	
compute_cluster() (dtwhacluster-	dtwhaclustering.leastSquareModeling	
ing.dtw_analysis.dtw_clustering method),	module, 9	
11	<pre>dtwhaclustering.plot_linear_trend   module, 7</pre>	
compute_cut_off_inconsistency() (dtwhaclus-	dtwhaclustering.plot_stations	
tering.dtw_analysis.dtw_clustering method),	module, 15	
compute_dendrogram() (dtwhacluster-		
ing.dtw_analysis.dtw_clustering method),	G	
11	<pre>get_linkage() (dtwhacluster-</pre>	
compute_distance() (dtwhacluster-	ing.dtw_analysis.dtw_clustering method),	
ing.dtw_analysis.dtw_signal_pairs method),	11	
compute_distance_accl() (dtwhacluster-	J	
ing.dtw_analysis.dtw_clustering method),	jump() (dtwhaclustering.leastSquareModeling.LSQmodules	
11	method), 9	
<pre>compute_distance_matrix() (dtwhacluster-</pre>		
ing.dtw_analysis.dtw_clustering method),	L	
11	lsqmodeling() (in module dtwhacluster-	
<pre>compute_interpolation() (in module dtwhacluster- ing.plot_linear_trend), 7</pre>	ing.leastSquareModeling), 9 LSOmodules (class in dtwhacluster-	
compute_lsq() (dtwhacluster-	LSQmodules (class in dtwhacluster- ing.leastSquareModeling), 9	
ing.leastSquareModeling.LSQmodules	ing.icusisquaremodeting), )	
method), 9	M	
compute_warping_path() (dtwhacluster-	module	
ing.dtw_analysis.dtw_signal_pairs method),	${\tt dtwhaclustering}, 1$	
13	dtwhaclustering.analysis_support,5	
D	<pre>dtwhaclustering.dtw_analysis, 11 dtwhaclustering.leastSquareModeling, 9</pre>	
<pre>dec2dt() (in module dtwhaclustering.analysis_support),</pre>	dtwhaclustering.plot_linear_trend,7	
5	dtwhaclustering.plot_stations, 15	
dec2dt_scalar() (in module dtwhacluster-		
ing.analysis_support), 5	N	
<pre>dtw_clustering (class in dtwhacluster- ing.dtw_analysis), 11</pre>	<pre>noise_robustness_test() (in module dtwhacluster-</pre>	
dtw_signal_pairs (class in dtwhacluster-	ing.dtw_analysis), 14	
ing.dtw_analysis), 13	0	
dtwhaclustering		
module, 1	<pre>optimum_cluster_elbow() (dtwhacluster- ing.dtw_analysis.dtw_clustering method),</pre>	
dtwhaclustering.analysis_support	11	
<pre>module, 5 dtwhaclustering.dtw_analysis</pre>		
acwinactus certing. acw_anaty 515		

#### Р plot\_cluster() (in module dtwhaclustering.dtw\_analysis), 14 plot\_cluster\_geomap() (dtwhaclustering.dtw\_analysis.dtw\_clustering method), plot\_cluster\_geomap\_interpolated() (dtwhaclustering.dtw\_analysis.dtw\_clustering method), 12 plot\_cluster\_xymap() (dtwhaclustering.dtw\_analysis.dtw\_clustering method), plot\_dendrogram() (dtwhaclustering.dtw\_analysis.dtw\_clustering method), plot\_hac\_iteration() (dtwhaclustering.dtw\_analysis.dtw\_clustering method), plot\_linear\_trend\_on\_map() (in module dtwhaclustering.plot\_linear\_trend), 7 plot\_matrix() (dtwhaclustering.dtw\_analysis.dtw\_signal\_pairs method), 14 plot\_optimum\_cluster() (dtwhaclustering.dtw\_analysis.dtw\_clustering method), plot\_polar\_dendrogram() (dtwhaclustering.dtw\_analysis.dtw\_clustering method), 13 plot\_signals() (dtwhaclustering.dtw\_analysis.dtw\_clustering method), plot\_signals() (dtwhaclustering.dtw\_analysis.dtw\_signal\_pairs method), dtwhaclusterplot\_signals() (in module ing.dtw\_analysis), 14 plot\_station\_map() (in module dtwhaclustering.plot\_stations), 15 plot\_warping\_path() (dtwhaclustering.dtw\_analysis.dtw\_signal\_pairs *method*), S dtwhaclustershuffle\_signals() module (in ing.dtw\_analysis), 14 significance\_test() (dtwhaclustering.dtw\_analysis.dtw\_clustering method), Т toYearFraction() module dtwhacluster-(in ing.analysis\_support), 5

22 Index