

# Package: HVT (via r-universe)

March 7, 2025

**Type** Package

**Date** 2025-02-03

**Title** Constructing Hierarchical Voronoi Tessellations and Overlay  
Heatmaps for Data Analysis

**Version** 25.2.2

**Description** Facilitates building topology preserving maps for data  
analysis.

**License** Apache License 2.0

**Encoding** UTF-8

**Imports** MASS, grDevices, splancs, stats, dplyr, NbClust, purrr,  
magrittr, ggplot2, tidyr, scales, cluster, reshape2,  
FNN,Rtsne,umap, plyr, gganimate, markovchain, methods,deldir,  
gridExtra

**Depends** R (>= 4.0.0)

**BugReports** <https://github.com/Mu-Sigma/HVT/issues>

**URL** <https://github.com/Mu-Sigma/HVT>

**RoxygenNote** 7.3.2

**Suggests** knitr,rmarkdown,testthat,geozoo,plotly,  
DT,patchwork,sp,Hmisc,data.table,gtable,  
htmlwidgets,installr,skimr,tibble,devtools,gifski,  
tidyverse,DataExplorer,htmltools,corrplot,kableExtra,polyclip,conf.design

**VignetteBuilder** knitr

**NeedsCompilation** no

**Author** Zubin Dowlaty [aut], Mu Sigma, Inc. [cre]

**Maintainer** ``Mu Sigma, Inc." <ird.expericelab@mu-sigma.com>

**Repository** CRAN

**Date/Publication** 2025-02-05 12:50:02 UTC

**Config/pak/sysreqs** libgdal-dev gdal-bin libgeos-dev libglpk-dev make  
libcudnn-dev libpng-dev libxml2-dev libssl-dev libproj-dev  
python3 libsqlite3-dev libudunits2-dev

## Contents

clustHVT . . . . .	2
displayTable . . . . .	4
edaPlots . . . . .	5
getTransitionProbability . . . . .	6
msm . . . . .	7
plotAnimatedFlowmap . . . . .	9
plotHVT . . . . .	11
plotModelDiagnostics . . . . .	13
plotNovelCells . . . . .	14
plotQuantErrorHistogram . . . . .	15
plotStateTransition . . . . .	16
plotZscore . . . . .	17
reconcileTransitionProbability . . . . .	18
removeNovelty . . . . .	20
scoreHVT . . . . .	21
scoreLayeredHVT . . . . .	23
summary . . . . .	25
trainHVT . . . . .	26
<b>Index</b>	<b>30</b>

---

clustHVT

*Performing Hierarchical Clustering Analysis*

---

### Description

This is the main function to perform hierarchical clustering analysis which determines optimal number of clusters, perform AGNES clustering and plot the 2D cluster hvt plot.

### Usage

```
clustHVT(
  data,
  trainHVT_results,
  scoreHVT_results,
  clustering_method = "ward.D2",
  indices,
  clusters_k = "champion",
  type = "default",
  domains.column
)
```

**Arguments**

<code>data</code>	Data frame. A data frame intended for performing hierarchical clustering analysis.
<code>trainHVT_results</code>	List. A list object which is obtained as a result of <code>trainHVT</code> function.
<code>scoreHVT_results</code>	List. A list object which is obtained as a result of <code>scoreHVT</code> function.
<code>clustering_method</code>	Character. The method used for clustering in both <code>NbClust</code> and <code>hclust</code> function. Defaults to 'ward.D2'.
<code>indices</code>	Character. The indices used for determining the optimal number of clusters in <code>NbClust</code> function. By default it uses 20 different indices.
<code>clusters_k</code>	Character. A parameter that specifies the number of clusters for the provided data. The options include "champion," "challenger," or any integer between 1 and 20. Selecting "champion" will use the highest number of clusters recommended by the 'NbClust' function, while "challenger" will use the second-highest recommendation. If a numerical value from 1 to 20 is provided, that exact number will be used as the number of clusters.
<code>type</code>	Character. The type of output required. Default is 'default'. Other option is 'plot' which will return only the clustered heatmap.
<code>domains.column</code>	Character. A vector of cluster names for the clustered heatmap. Used only when type is 'plot'.

**Value**

A list object that contains the hierarchical clustering results.

<code>[[1]]</code>	Summary of k suggested by all indices with plots
<code>[[2]]</code>	A dendrogram plot with the selected number of clusters
<code>[[3]]</code>	A 2D Cluster HVT Plotly visualization that colors cells according to clusters derived from AGNES clustering results. It is interactive, allowing users to view cell contents by hovering over them

**Author(s)**

Vishwavani <[vishwavani@mu-sigma.com](mailto:vishwavani@mu-sigma.com)>

**Examples**

```
data("EuStockMarkets")
dataset <- data.frame(t = as.numeric(time(EuStockMarkets)),
                     DAX = EuStockMarkets[, "DAX"],
                     SMI = EuStockMarkets[, "SMI"],
                     CAC = EuStockMarkets[, "CAC"],
                     FTSE = EuStockMarkets[, "FTSE"])
rownames(EuStockMarkets) <- dataset$t
hvt.results<- trainHVT(dataset[-1],n_cells = 30, depth = 1, quant.err = 0.1,
```

```
distance_metric = "L1_Norm", error_metric = "max",
normalize = TRUE, quant_method = "kmeans")
scoring <- scoreHVT(dataset, hvt.results, analysis.plots = TRUE, names.column = dataset[,1])
centroid_data <- scoring$centroidData
hclust_data_1 <- centroid_data[,2:3]
clust.results <- clustHVT(data = hclust_data_1,
trainHVT_results = hvt.results,
scoreHVT_results = scoring,
clusters_k = 'champion', indices = 'hartigan')
```

---

displayTable	<i>function for displaying table</i>
--------------	--------------------------------------

---

## Description

This is the main function for displaying data in table format

## Usage

```
displayTable(data, scroll = TRUE, limit = 20)
```

## Arguments

data	Data frame. The dataframe to be displayed in table format.
scroll	Logical. A value to have a scroll or not in the table. Default is TRUE.
limit	Numeric. A value to indicate how many rows to display. Default is 20.

## Value

A table with proper formatting for html notebook

## Author(s)

Vishwavani <[vishwavani@mu-sigma.com](mailto:vishwavani@mu-sigma.com)>

## Examples

```
data <- datasets::EuStockMarkets
dataset <- as.data.frame(data)
displayTable(dataset)
```

---

edaPlots *plots for data analysis*

---

## Description

This is the main function that provides exploratory data analysis plots

## Usage

```
edaPlots(  
  df,  
  time_column,  
  output_type = "summary",  
  n_cols = -1,  
  greyBars = NULL  
)
```

## Arguments

df	Dataframe. A data frame object.
time_column	Character. The name of the time column in the data frame. Can be given only when the data is time series
output_type	Character. The name of the output to be displayed. Options are 'summary', 'histogram', 'boxplot', 'timeseries' & 'correlation'. Default value is summary.
n_cols	Numeric. A value to indicate how many columns to be included in the output.
greyBars	List. A list of timestamps where each list contains two elements: start and end period, which will be highlighted in gray in the time series plot. Default value is NULL.

## Value

Five objects which include time series plots, data distribution plots, box plots, correlation plot and a descriptive statistics table.

## Author(s)

Vishwavani <[vishwavani@mu-sigma.com](mailto:vishwavani@mu-sigma.com)>

## Examples

```
dataset <- data.frame(date = as.numeric(time(EuStockMarkets)),  
  DAX = as.numeric(EuStockMarkets[, "DAX"]),  
  SMI = as.numeric(EuStockMarkets[, "SMI"]),  
  CAC = as.numeric(EuStockMarkets[, "CAC"]),  
  FTSE = as.numeric(EuStockMarkets[, "FTSE"]))  
edaPlots(dataset)  
edaPlots(dataset, time_column = 'date', output_type = 'timeseries', n_cols = 4)
```

---

```
getTransitionProbability
```

*Creating Transition Probability Matrix*

---

### Description

This is the main function to create transition probability matrix. The transition probability matrix quantifies the likelihood of transitioning from one state to another. States: The table includes the current states and the possible next states. Probabilities: For each current state, it lists the probability of transitioning to each of the next possible states.

### Usage

```
getTransitionProbability(
  df,
  cellid_column,
  time_column,
  type = "with_self_state"
)
```

### Arguments

df	Data frame. The input data frame should contain two columns, cell ID from scoreHVT function and time stamp of that dataset.
cellid_column	Character. Name of the column containing cell IDs.
time_column	Character. Name of the column containing time stamps.
type	Character. A character value indicating the type of transition probability table to create. Accepted entries are "with_self_state" and "without_self_state".

### Value

Stores a data frames with transition probabilities.

### Author(s)

PonAnuReka Seenivasan <ponanureka.s@mu-sigma.com>, Vishwavani <vishwavani@mu-sigma.com>

### Examples

```
dataset <- data.frame(t = as.numeric(time(EuStockMarkets)),
  DAX = EuStockMarkets[, "DAX"],
  SMI = EuStockMarkets[, "SMI"],
  CAC = EuStockMarkets[, "CAC"],
  FTSE = EuStockMarkets[, "FTSE"])
hvt.results<- trainHVT(dataset[-1],n_cells = 60, depth = 1, quant.err = 0.1,
  distance_metric = "L1_Norm", error_metric = "max",
  normalize = TRUE,quant_method = "kmeans")
```

```

scoring <- scoreHVT(dataset, hvt.results)
cell_id <- scoring$scoredPredictedData$Cell.ID
time_stamp <- dataset$t
dataset <- data.frame(cell_id, time_stamp)
table <- getTransitionProbability(dataset, cellid_column = "cell_id", time_column = "time_stamp")

```

**Description**

This is the main function to perform Monte Carlo simulations of Markov Chain on the dynamic forecasting of HVT States of a time series dataset. It includes both ex-post and ex-ante analysis offering valuable insights into future trends while resolving state transition challenges through clustering and nearest-neighbor methods to enhance simulation accuracy.

**Usage**

```

msm(
  state_time_data,
  forecast_type = "ex-post",
  initial_state,
  n_ahead_ante,
  transition_probability_matrix,
  num_simulations = 100,
  trainHVT_results,
  scoreHVT_results,
  actual_data = NULL,
  raw_dataset,
  k = 5,
  handle_problematic_states = FALSE,
  n_nearest_neighbor = 1,
  show_simulation = TRUE,
  mae_metric = "median"
)

```

**Arguments**

<code>state_time_data</code>	DataFrame. A dataframe containing state transitions over time(cell id and timestamp)
<code>forecast_type</code>	Character. A character to indicate the type of forecasting. Accepted values are "ex-post" or "ex-ante".
<code>initial_state</code>	Numeric. An integer indicating the state at t0.
<code>n_ahead_ante</code>	Numeric. A vector of n ahead points to be predicted further in ex-ante analyzes.

<code>transition_probability_matrix</code>	DataFrame. A dataframe of transition probabilities/ output of 'getTransition-Probability' function
<code>num_simulations</code>	Integer. A number indicating the total number of simulations to run. Default is 100.
<code>trainHVT_results</code>	List. 'trainHVT' function output
<code>scoreHVT_results</code>	List. 'scoreHVT' function output
<code>actual_data</code>	Dataframe. A dataframe for ex-post prediction period with the actual raw data values
<code>raw_dataset</code>	DataFrame. A dataframe of input raw dataset from the mean and standard deviation will be calculated to scale up the predicted values
<code>k</code>	Integer. A number of optimal clusters when handling problematic states. Default is 5.
<code>handle_problematic_states</code>	Logical. To indicate whether to handle problematic states or not. Default is FALSE.
<code>n_nearest_neighbor</code>	Integer. A number of nearest neighbors to consider when handling problematic states. Default is 1.
<code>show_simulation</code>	Logical. To indicate whether to show the simulation lines in plots or not. Default is TRUE.
<code>mae_metric</code>	Character. A character to indicate which metric to calculate Mean Absolute Error. Accepted entries are "mean", "median", or "mode". Default is "median".

**Value**

A list object that contains the forecasting plots and MAE values.

<code>[[1]]</code>	Simulation plots and MAE values for state and centroids plot
<code>[[2]]</code>	Summary Table, Dendrogram plot and Clustered Heatmap when <code>handle_problematic_states</code> is TRUE

**Author(s)**

Vishwavani <[vishwavani@mu-sigma.com](mailto:vishwavani@mu-sigma.com)>

**Examples**

```
dataset <- data.frame(t = as.numeric(time(EuStockMarkets)),
  DAX = EuStockMarkets[, "DAX"],
  SMI = EuStockMarkets[, "SMI"],
  CAC = EuStockMarkets[, "CAC"],
  FTSE = EuStockMarkets[, "FTSE"])
hvt.results<- trainHVT(dataset[,-1],n_cells = 60, depth = 1, quant.err = 0.1,
```



```

        distance_metric = "L1_Norm", error_metric = "max",
        normalize = TRUE, quant_method = "kmeans")
scoring <- scoreHVT(dataset, hvt.results)
cell_id <- scoring$scoredPredictedData$Cell.ID
time_stamp <- dataset$t
temporal_data <- data.frame(cell_id, time_stamp)
table <- getTransitionProbability(temporal_data,
cellid_column = "cell_id", time_column = "time_stamp")
colnames(temporal_data) <- c("Cell.ID", "t")
ex_post_forecasting <- dataset[1800:1860,]
ex_post <- msm(state_time_data = temporal_data,
        forecast_type = "ex-post",
        transition_probability_matrix = table,
        initial_state = 2,
        num_simulations = 100,
        scoreHVT_results = scoring,
        trainHVT_results = hvt.results,
        actual_data = ex_post_forecasting,
        raw_dataset = dataset,
        mae_metric = "median",
        show_simulation = FALSE)

```

---

plotAnimatedFlowmap     *Generating flow maps and animations based on transition probabilities*

---

## Description

This is the main function for generating flow maps and animations based on transition probabilities including self states and excluding self states. Flow maps are a type of data visualization used to represent the transition probability of different states. Animations are the gifs used to represent the movement of data through the cells.

## Usage

```

plotAnimatedFlowmap(
  hvt_model_output,
  transition_probability_df,
  df,
  animation = "All",
  flow_map = "All",
  fps_time = 1,
  fps_state = 1,
  time_duration = 2,
  state_duration = 2,
  cellid_column,
  time_column
)

```

**Arguments**

<code>hvt_model_output</code>	List. Output from a trainHVT function.
<code>transition_probability_df</code>	List. Output from getTransitionProbability function
<code>df</code>	Data frame. The input dataframe should contain two columns, cell ID from scoreHVT function and time stamp of that dataset.
<code>animation</code>	Character. Type of animation ('state_based', 'time_based', 'All' or NULL)
<code>flow_map</code>	Character. Type of flow map ('self_state', 'without_self_state', 'All' or NULL)
<code>fps_time</code>	Numeric. A numeric value for the frames per second of the time transition gif. (Must be a numeric value and a factor of 100). Default value is 1.
<code>fps_state</code>	Numeric. A numeric value for the frames per second of the state transition gif. (Must be a numeric value and a factor of 100). Default value is 1.
<code>time_duration</code>	Numeric. A numeric value for the duration of the time transition gif. Default value is 2.
<code>state_duration</code>	Numeric. A numeric value for the duration of the state transition gif. Default value is 2.
<code>cellid_column</code>	Character. Name of the column containing cell IDs.
<code>time_column</code>	Character. Name of the column containing time stamps

**Value**

A list of flow maps and animation gifs.

**Author(s)**

PonAnuReka Seenivasan <ponanureka.s@mu-sigma.com>, Vishwavani <vishwavani@mu-sigma.com>

**See Also**

[trainHVT](#)  
[scoreHVT](#)  
[getTransitionProbability](#)

**Examples**

```
dataset <- data.frame(date = as.numeric(time(EuStockMarkets)),
                     DAX = EuStockMarkets[, "DAX"],
                     SMI = EuStockMarkets[, "SMI"],
                     CAC = EuStockMarkets[, "CAC"],
                     FTSE = EuStockMarkets[, "FTSE"])
hvt.results<- trainHVT(dataset,n_cells = 60, depth = 1, quant.err = 0.1,
                      distance_metric = "L1_Norm", error_metric = "max",
                      normalize = TRUE,quant_method = "kmeans")
scoring <- scoreHVT(dataset, hvt.results)
cell_id <- scoring$scoredPredictedData$Cell.ID
time_stamp <- dataset$date
```

```
dataset <- data.frame(cell_id, time_stamp)
table <- getTransitionProbability(dataset, cellid_column = "cell_id", time_column = "time_stamp")
plots <- plotAnimatedFlowmap(hvt_model_output = hvt.results, transition_probability_df = table,
df = dataset, animation = 'All', flow_map = 'All', fps_time = 1, fps_state = 1, time_duration = 3,
state_duration = 3, cellid_column = "cell_id", time_column = "time_stamp")
```

---

plotHVT

*Plot the hierarchical tessellations.*


---

### Description

This is the main plotting function to construct hierarchical voronoi tessellations in 1D, 2D or Interactive surface plot.

### Usage

```
plotHVT(
  hvt.results,
  line.width = 0.5,
  color.vec = "black",
  centroid.size = 0.6,
  centroid.color = "black",
  child.level = 1,
  hmap.cols,
  separation_width = 7,
  layer_opacity = c(0.5, 0.75, 0.99),
  dim_size = 1000,
  plot.type = "2Dhvt",
  quant.error.hmap = NULL,
  cell_id = FALSE,
  cell_id_position = "bottom",
  cell_id_size = 2.6
)
```

### Arguments

hvt.results	(1D/2DProj/2Dhvt/2Dheatmap/surface_plot) List. A list containing the output of trainHVT function which has the details of the tessellations to be plotted.
line.width	(2Dhvt/2Dheatmap) Numeric Vector. A vector indicating the line widths of the tessellation boundaries for each level.
color.vec	(2Dhvt/2Dheatmap) Vector. A vector indicating the colors of the boundaries of the tessellations at each level.
centroid.size	(2Dhvt/2Dheatmap) Numeric Vector. A vector indicating the size of centroids for each level.
centroid.color	(2Dhvt/2Dheatmap) Numeric Vector. A vector indicating the color of centroids for each level.

child.level	(2Dheatmap/surface_plot) Numeric. Indicating the level for which the plot should be displayed
hmap.cols	(2Dheatmap/surface_plot) Numeric or Character. The column number or column name from the dataset indicating the variables for which the heat map is to be plotted.
separation_width	(surface_plot) Numeric. An integer indicating the width between hierarchical levels in surface plot
layer_opacity	(surface_plot) Numeric. A vector indicating the opacity of each hierarchical levels in surface plot
dim_size	(surface_plot) Numeric. An integer controls the resolution or granularity of the 3D surface grid
plot.type	Character. An option to indicate which type of plot should be generated. Accepted entries are '1D', '2Dproj', '2Dhvt', '2Dheatmap' and 'surface_plot'. Default value is '2Dhvt'.
quant.error.hmap	(2Dheatmap) Numeric. A number representing the quantization error threshold to be highlighted in the heatmap. When a value is provided, it will emphasize cells with quantization errors equal or less than the specified threshold, indicating that these cells cannot be further subdivided in the next depth layer. The default value is NULL, meaning all cells will be colored in the heatmap across various depths.
cell_id	(2Dhvt/2Dheatmap) Logical. A logical indicating whether the cell IDs should be displayed
cell_id_position	(2Dhvt/2Dheatmap) Character. A character indicating the position of the cell IDs. Accepted entries are 'top', 'bottom', 'left' and 'right'.
cell_id_size	(2Dhvt/2Dheatmap) Numeric. A numeric vector indicating the size of the cell IDs for all levels.

**Value**

plot object containing the visualizations of reduced dimension(1D/2D) for the given dataset.

**Author(s)**

Shubhra Prakash <shubhra.prakash@mu-sigma.com>, Sangeet Moy Das <sangeet.das@mu-sigma.com>, Vishwavani <vishwavani@mu-sigma.com>

**See Also**

[trainHVT](#)

**Examples**

```
data("EuStockMarkets")
hvt.results <- trainHVT(EuStockMarkets, n_cells = 60, depth = 1, quant.err = 0.1,
```

```

distance_metric = "L1_Norm", error_metric = "max",
normalize = TRUE, quant_method="kmeans")

#change the 'plot.type' argument to '2Dproj' or '2DHVT' to visualize respective plots.
plotHVT(hvt.results, plot.type='1D')

#change the 'plot.type' argument to 'surface_plot' to visualize the Interactive surface plot
plotHVT(hvt.results, child.level = 1,
hmap.cols = "DAX", plot.type = '2Dheatmap')

```

---

plotModelDiagnostics *Make the diagnostic plots for hierarchical voronoi tessellations*

---

## Description

This is the main function that generates diagnostic plots for hierarchical voronoi tessellations models and scoring.

## Usage

```
plotModelDiagnostics(model_obj)
```

## Arguments

model\_obj      List. A list obtained from the trainHVT function or scoreHVT function

## Value

For trainHVT, Minimum Intra-DataPoint Distance Plot, Minimum Intra-Centroid Distance Plot, Mean Absolute Deviation Plot, Distribution of Number of Observations in Cells, for Training Data and Mean Absolute Deviation Plot for Validation Data are plotted. For scoreHVT Mean Absolute Deviation Plot for Training Data and Validation Data are plotted

## Author(s)

Shubhra Prakash <shubhra.prakash@mu-sigma.com>

## See Also

[plotHVT](#)

## Examples

```

data("EuStockMarkets")
hvt.results <- trainHVT(EuStockMarkets, n_cells = 60, depth = 1, quant.err = 0.1,
distance_metric = "L1_Norm", error_metric = "max",
normalize = TRUE, quant_method="kmeans", diagnose = TRUE,
hvt_validation = TRUE)
plotModelDiagnostics(hvt.results)

```

---

plotNovelCells      *Plot the identified outlier cells in the voronoi tessellation map.*

---

### Description

This is the main plotting function to construct hierarchical voronoi tessellations and highlight the outlier cells

### Usage

```
plotNovelCells(
  plot.cells,
  hvt.map,
  line.width = c(0.6),
  color.vec = c("#141B41"),
  pch = 21,
  centroid.size = 0.5,
  title = NULL,
  maxDepth = 1
)
```

### Arguments

plot.cells	Vector. A vector indicating the cells to be highlighted in the map
hvt.map	List. A list containing the output of trainHVT function which has the details of the tessellations to be plotted
line.width	Numeric Vector. A vector indicating the line widths of the tessellation boundaries for each level
color.vec	Vector. A vector indicating the colors of the boundaries of the tessellations at each level
pch	Numeric. Symbol of the centroids of the tessellations (parent levels) Default value is 21.
centroid.size	Numeric. Size of centroids of first level tessellations. Default value is 0.5
title	String. Set a title for the plot. (default = NULL)
maxDepth	Numeric. An integer indicating the number of levels. (default = NULL)

### Value

Returns a ggplot object containing hierarchical voronoi tessellation plot highlighting the outlier cells

### Author(s)

Shantanu Vaidya <shantanu.vaidya@mu-sigma.com>

**See Also**

[trainHVT](#)  
[plotHVT](#)

**Examples**

```
data("EuStockMarkets")
hvt.results <- trainHVT(EuStockMarkets, n_cells = 60, depth = 1, quant.err = 0.1,
                        distance_metric = "L1_Norm", error_metric = "max",
                        normalize = TRUE, quant_method="kmeans")
#selected 55,58 are for demo purpose
plotNovelCells(c(55,58),hvt.results)
```

---

plotQuantErrorHistogram

*Make the quantization error plots for training and scoring.*

---

**Description**

This is the function that produces histograms displaying the distribution of Quantization Error (QE) values for both train and test datasets, highlighting mean values with dashed lines for quick evaluation.

**Usage**

```
plotQuantErrorHistogram(hvt.results, hvt.scoring)
```

**Arguments**

`hvt.results` List. A list of `hvt.results` obtained from the `trainHVT` function.  
`hvt.scoring` List. A list of `hvt.scoring` obtained from the `scoreHVT` function.

**Value**

Returns the `ggplot` object containing the quantization error distribution plots for the given HVT results of training and scoring

**Author(s)**

Shubhra Prakash <shubhra.prakash@mu-sigma.com>

**See Also**

[plotHVT](#)

**Examples**

```

data("EuStockMarkets")
dataset <- data.frame(date = as.numeric(time(EuStockMarkets)),
                     DAX = EuStockMarkets[, "DAX"],
                     SMI = EuStockMarkets[, "SMI"],
                     CAC = EuStockMarkets[, "CAC"],
                     FTSE = EuStockMarkets[, "FTSE"])
rownames(EuStockMarkets) <- dataset$date
#Split in train and test
train <- EuStockMarkets[1:1302, ]
test <- EuStockMarkets[1303:1860, ]

hvt.results<- trainHVT(train,n_cells = 60, depth = 1, quant.err = 0.1,
                      distance_metric = "L1_Norm", error_metric = "max",
                      normalize = TRUE, quant_method = "kmeans")
scoring <- scoreHVT(test, hvt.results)
plotQuantErrorHistogram(hvt.results, scoring)

```

---

plotStateTransition    *Creating State Transition Plot*

---

**Description**

This is the main function to create a state transition plot from a data frame. A state transition plot is a type of data visualization used to represent the changes or transitions in states over time for a given system. State refers to a particular condition or status of a cell at a specific point in time. Transition refers to the change of state for a cell from one condition to another over time.

**Usage**

```

plotStateTransition(
  df,
  sample_size = NULL,
  line_plot = NULL,
  cellid_column,
  time_column,
  v_intercept = NULL,
  time_periods = NULL
)

```

**Arguments**

df	Data frame. The Input data frame should contain two columns. Cell ID from scoreHVT function and time stamp of that dataset.
sample_size	Numeric. An integer indicating the fraction of the data frame to visualize in the plot. Default value is 0.2
line_plot	Logical. A logical value indicating to create a line plot. Default value is NULL.



cellid_column	Character. Name of the column containing cell IDs.
time_column	Character. Name of the column containing time stamps.
v_intercept	Numeric. A numeric value indicating the time stamp to draw a vertical line on the plot.
time_periods	List. A list of vectors, each containing start and end times for highlighting time periods.

**Value**

A plotly object representing the state transition plot for the given data frame.

**Author(s)**

PonAnuReka Seenivasan <ponanureka.s@mu-sigma.com>, Vishwavani <vishwavani@mu-sigma.com>

**Examples**

```
dataset <- data.frame(date = as.numeric(time(EuStockMarkets)),
  DAX = EuStockMarkets[, "DAX"],
  SMI = EuStockMarkets[, "SMI"],
  CAC = EuStockMarkets[, "CAC"],
  FTSE = EuStockMarkets[, "FTSE"])

hvt.results<- trainHVT(dataset,n_cells = 60, depth = 1, quant.err = 0.1,
  distance_metric = "L1_Norm", error_metric = "max",
  normalize = TRUE,quant_method = "kmeans")

scoring <- scoreHVT(dataset, hvt.results)
cell_id <- scoring$scoredPredictedData$Cell.ID
time_stamp <- dataset$date
dataset <- data.frame(cell_id, time_stamp)
plotStateTransition(dataset, sample_size = 1, cellid_column = "cell_id",time_column = "time_stamp")
```

---

plotZscore

*Plots of z scores*

---

**Description**

This is the main function to plot the z scores against cell ids.

**Usage**

```
plotZscore(
  data,
  cell_range = NULL,
  segment_size = 2,
  reference_lines = c(-1.65, 1.65)
)
```

**Arguments**

<code>data</code>	Data frame. A data frame of cell id and features.
<code>cell_range</code>	Vector. A numeric vector of cell id range for which the plot should be displayed. Default is NULL, which plots all the cells.
<code>segment_size</code>	Integer. A numeric value to indicate the size of the bars in the plot. Default is 2.
<code>reference_lines</code>	Vector. A numeric vector of confidence interval values for the reference lines in the plot. Default is <code>c(-1.65, 1.65)</code> .

**Value**

A grid of plots of z score against cell id of teh given features.

**Author(s)**

Vishwavani <[vishwavani@mu-sigma.com](mailto:vishwavani@mu-sigma.com)>

**Examples**

```
data("EuStockMarkets")
dataset <- data.frame(t = as.numeric(time(EuStockMarkets)),
                     DAX = EuStockMarkets[, "DAX"],
                     SMI = EuStockMarkets[, "SMI"],
                     CAC = EuStockMarkets[, "CAC"],
                     FTSE = EuStockMarkets[, "FTSE"])
rownames(EuStockMarkets) <- dataset$t
hvt.results<- trainHVT(dataset[-1],n_cells = 60, depth = 1, quant.err = 0.1,
                      distance_metric = "L1_Norm", error_metric = "max",
                      normalize = TRUE,quant_method = "kmeans")
col_names <- c("Cell.ID", "DAX", "SMI", "CAC", "FTSE")
data <- dplyr::arrange(dplyr::select(hvt.results[[3]][["summary"]], col_names), Cell.ID)
data <- round(data, 2)
plotZscore(data)
```

---

reconcileTransitionProbability

*Reconciliation of Transition Probability*

---

**Description**

This is the main function for creating reconciliation plots and tables which helps in comparing the transition probabilities calculated manually and from markovchain function

**Usage**

```
reconcileTransitionProbability(
  df,
  hmap_type = NULL,
  cellid_column,
  time_column
)
```

**Arguments**

df	Data frame. The input data frame should contain two columns, cell ID from scoreHVT function and timestamp of that dataset.
hmap_type	Character. ('self_state', 'without_self_state', or 'All')
cellid_column	Character. Name of the column containing cell IDs.
time_column	Character. Name of the column containing timestamps

**Value**

A list of plotly heatmap objects and tables representing the transition probability heatmaps.

**Author(s)**

PonAnuReka Seenivasan <ponanureka.s@mu-sigma.com>, Vishwavani <vishwavani@mu-sigma.com>

**Examples**

```
dataset <- data.frame(date = as.numeric(time(EuStockMarkets)),
  DAX = EuStockMarkets[, "DAX"],
  SMI = EuStockMarkets[, "SMI"],
  CAC = EuStockMarkets[, "CAC"],
  FTSE = EuStockMarkets[, "FTSE"])
hvt.results<- trainHVT(dataset,n_cells = 60, depth = 1, quant.err = 0.1,
  distance_metric = "L1_Norm", error_metric = "max",
  normalize = TRUE,quant_method = "kmeans")
scoring <- scoreHVT(dataset, hvt.results)
cell_id <- scoring$scoredPredictedData$Cell.ID
time_stamp <- dataset$date
dataset <- data.frame(cell_id, time_stamp)
reconcileTransitionProbability(dataset, hmap_type = "All",
cellid_column = "cell_id", time_column = "time_stamp")
```

---

removeNovelty	<i>Remove identified novelty cell(s)</i>
---------------	--

---

**Description**

This function is used to remove the identified novelty cells.

**Usage**

```
removeNovelty(outlier_cells, hvt_results)
```

**Arguments**

`outlier_cells` Vector. A vector with the cell number of the identified novelty  
`hvt_results` List. A list having the results of the compressed map i.e. output of `trainHVT` function

**Value**

A list of two items

[[1]] Dataframe of novelty cell(s)  
[[2]] Dataframe without the novelty cell(s) from the dataset used in model training

**Author(s)**

Shantanu Vaidya <shantanu.vaidya@mu-sigma.com>

**See Also**

[trainHVT](#)  
[scoreLayeredHVT](#)

**Examples**

```
data("EuStockMarkets")
hvt.results <- trainHVT(EuStockMarkets, n_cells = 60, depth = 1, quant.err = 0.1,
                       distance_metric = "L1_Norm", error_metric = "max",
                       normalize = TRUE, quant_method="kmeans")
identified_Novelty_cells <- c(2, 10)
output_list <- removeNovelty(identified_Novelty_cells, hvt.results)
data_with_novelty <- output_list[[1]]
data_without_novelty <- output_list[[2]]
```

---

scoreHVT	<i>Score which cell each point in the test dataset belongs to.</i>
----------	--

---

### Description

This function scores each data point in the test dataset based on a trained hierarchical Voronoi tessellations model.

### Usage

```
scoreHVT(
  dataset,
  hvt.results.model,
  child.level = 1,
  mad.threshold = 0.2,
  line.width = 0.6,
  color.vec = c("navyblue", "slateblue", "lavender"),
  normalize = TRUE,
  distance_metric = "L1_Norm",
  error_metric = "max",
  yVar = NULL,
  analysis.plots = FALSE,
  names.column = NULL
)
```

### Arguments

dataset	Data frame. A data frame which to be scored. Can have categorical columns if 'analysis.plots' are required.
hvt.results.model	List. A list obtained from the trainHVT function
child.level	Numeric. A number indicating the depth for which the heat map is to be plotted.
mad.threshold	Numeric. A numeric value indicating the permissible Mean Absolute Deviation.
line.width	Vector. A vector indicating the line widths of the tessellation boundaries for each layer.
color.vec	Vector. A vector indicating the colors of the tessellation boundaries at each layer.
normalize	Logical. A logical value indicating if the dataset should be normalized. When set to TRUE, the data (testing dataset) is standardized by 'mean' and 'sd' of the training dataset referred from the trainHVT(). When set to FALSE, the data is used as such without any changes.
distance_metric	Character. The distance metric can be L1_Norm(Manhattan) or L2_Norm(Eucledian). L1_Norm is selected by default. The distance metric is used to calculate the distance between an n dimensional point and centroid. The distance metric can be different from the one used during training.

error_metric	Character. The error metric can be mean or max. max is selected by default. max will return the max of m values and mean will take mean of m values where each value is a distance between a point and centroid of the cell.
yVar	Character. A character or a vector representing the name of the dependent variable(s)
analysis.plots	Logical. A logical value indicating that the scored plot should be plotted or not. If TRUE, the identifier column(character column) name should be supplied in 'names.column' argument. The output will be a 2D heatmap plotly which gives info on the cell id and the observations of a cell.
names.column	Character. A character or a vector representing the name of the identifier column/character column.

### Value

Dataframe containing scored data, plots and summary

### Author(s)

Shubhra Prakash <shubhra.prakash@mu-sigma.com>, Sangeet Moy Das <sangeet.das@mu-sigma.com>, Vishwavani <vishwavani@mu-sigma.com>

### See Also

[trainHVT](#)  
[plotHVT](#)

### Examples

```
data("EuStockMarkets")
dataset <- data.frame(date = as.numeric(time(EuStockMarkets)),
                      DAX = EuStockMarkets[, "DAX"],
                      SMI = EuStockMarkets[, "SMI"],
                      CAC = EuStockMarkets[, "CAC"],
                      FTSE = EuStockMarkets[, "FTSE"])
rownames(EuStockMarkets) <- dataset$date
# Split in train and test
train <- EuStockMarkets[1:1302, ]
test <- EuStockMarkets[1303:1860, ]
#model training
hvt.results<- trainHVT(train,n_cells = 60, depth = 1, quant.err = 0.1,
                      distance_metric = "L1_Norm", error_metric = "max",
                      normalize = TRUE,quant_method = "kmeans")
scoring <- scoreHVT(test, hvt.results)
data_scored <- scoring$scoredPredictedData
```

---

scoreLayeredHVT	<i>Score which cell and what layer each data point in the test dataset belongs to</i>
-----------------	---

---

## Description

This function that scores the cell and corresponding layer for each data point in a test dataset using three hierarchical vector quantization (HVT) models (Map A, Map B, Map C) and returns a data frame containing the scored layer output. The function incorporates the scored results from each map and merges them to provide a comprehensive result.

## Usage

```
scoreLayeredHVT(
  data,
  hvt_mapA,
  hvt_mapB,
  hvt_mapC,
  mad.threshold = 0.2,
  normalize = TRUE,
  seed = 300,
  distance_metric = "L1_Norm",
  error_metric = "max",
  child.level = 1,
  yVar = NULL
)
```

## Arguments

data	Data Frame. A data frame containing test dataset. The data frame should have all the variable(features) used for training.
hvt_mapA	A list of hvt.results.model obtained from trainHVT function while performing 'trainHVT()' on train data
hvt_mapB	A list of hvt.results.model obtained from trainHVT function while performing 'trainHVT()' on data with novelty(s)
hvt_mapC	A list of hvt.results.model obtained from trainHVT function while performing 'trainHVT()' on data without novelty(s)
mad.threshold	Numeric. A number indicating the permissible Mean Absolute Deviation
normalize	Logical. A logical value indicating if the dataset should be normalized. When set to TRUE, the data (testing dataset) is standardized by 'mean' and 'sd' of the training dataset referred from the trainHVT(). When set to FALSE, the data is used as such without any changes. (Default value is TRUE).
seed	Numeric. Random Seed.

distance_metric	Character. The distance metric can be L1_Norm(Manhattan) or L2_Norm(Eucledian). L1_Norm is selected by default. The distance metric is used to calculate the distance between an n dimensional point and centroid. The distance metric can be different from the one used during training.
error_metric	Character. The error metric can be mean or max. max is selected by default. max will return the max of m values and mean will take mean of m values where each value is a distance between a point and centroid of the cell.
child.level	Numeric. A number indicating the level for which the heat map is to be plotted.
yVar	Character. A character or a vector representing the name of the dependent variable(s)

**Value**

Dataframe containing scored layer output

**Author(s)**

Shubhra Prakash <shubhra.prakash@mu-sigma.com>, Sangeet Moy Das <sangeet.das@mu-sigma.com>, Shantanu Vaidya <shantanu.vaidya@mu-sigma.com>, Somya Shambhawi <somya.shambhawi@mu-sigma.com>

**See Also**

[trainHVT](#)  
[plotHVT](#)

**Examples**

```
data("EuStockMarkets")
dataset <- data.frame(date = as.numeric(time(EuStockMarkets)),
                      DAX = EuStockMarkets[, "DAX"],
                      SMI = EuStockMarkets[, "SMI"],
                      CAC = EuStockMarkets[, "CAC"],
                      FTSE = EuStockMarkets[, "FTSE"])
rownames(EuStockMarkets) <- dataset$date

train <- EuStockMarkets[1:1302, ]
test <- EuStockMarkets[1303:1860, ]

###MAP-A
hvt_mapA <- trainHVT(train, n_cells = 150, depth = 1, quant.err = 0.1,
                    distance_metric = "L1_Norm", error_metric = "max",
                    normalize = TRUE, quant_method = "kmeans")

identified_Novelty_cells <- c(127,55,83,61,44,35,27,77)
output_list <- removeNovelty(identified_Novelty_cells, hvt_mapA)
data_with_novelty <- output_list[[1]]
data_with_novelty <- data_with_novelty[, -c(1,2)]
```



```

### MAP-B
hvt_mapB <- trainHVT(data_with_novelty,n_cells = 10, depth = 1, quant.err = 0.1,
                    distance_metric = "L1_Norm", error_metric = "max",
                    normalize = TRUE,quant_method = "kmeans")
data_without_novelty <- output_list[[2]]

### MAP-C
hvt_mapC <- trainHVT(data_without_novelty,n_cells = 135,
                    depth = 1, quant.err = 0.1, distance_metric = "L1_Norm",
                    error_metric = "max", quant_method = "kmeans",
                    normalize = TRUE)

##SCORE LAYERED
data_scored <- scoreLayeredHVT(test, hvt_mapA, hvt_mapB, hvt_mapC)

```

summary

*Table for displaying summary***Description**

This is the main function for displaying summary from model training and scoring

**Usage**

```
summary(data, limit = 20, scroll = TRUE)
```

**Arguments**

data	List. A listed object from trainHVT or scoreHVT
limit	Numeric. A value to indicate how many rows to display.
scroll	Logical. A value to indicate whether to display scroll bar or not. Default value is TRUE.

**Value**

A consolidated table of summary for training, scoring and forecasting

**Author(s)**

Vishwavani <[vishwavani@mu-sigma.com](mailto:vishwavani@mu-sigma.com)>, Alimpan Dey <[alimpan.dey@mu-sigma.com](mailto:alimpan.dey@mu-sigma.com)>

**Examples**

```

data <- datasets::EuStockMarkets
dataset <- as.data.frame(data)
#model training
hvt.results <- trainHVT(dataset, n_cells = 60, depth = 1, quant.err = 0.1,
                      distance_metric = "L1_Norm", error_metric = "max",
                      normalize = TRUE, quant_method = "kmeans", dim_reduction_method = 'sammon')
summary(data = hvt.results)

```

**Description**

This is the main function to construct hierarchical voronoi tessellations. This is done using hierarchical vector quantization(hvq). The data is represented in 2D coordinates and the tessellations are plotted using these coordinates as centroids. For subsequent levels, transformation is performed on the 2D coordinates to get all the points within its parent tile. Tessellations are plotted using these transformed points as centroids.

**Usage**

```
trainHVT(  
  dataset,  
  min_compression_perc = NA,  
  n_cells = NA,  
  depth = 1,  
  quant.err = 0.2,  
  normalize = FALSE,  
  distance_metric = "L1_Norm",  
  error_metric = "max",  
  quant_method = "kmeans",  
  scale_summary = NA,  
  diagnose = FALSE,  
  hvt_validation = FALSE,  
  train_validation_split_ratio = 0.8,  
  dim_reduction_method = "sammon",  
  tsne_theta = 0.2,  
  tsne_eta = 200,  
  tsne_perplexity = 30,  
  tsne_verbose = TRUE,  
  tsne_max_iter = 500,  
  umap_n_neighbors = 60,  
  umap_n_components = 2,  
  umap_min_dist = 0.1  
)
```

**Arguments**

dataset	Data frame. A data frame, with numeric columns (features) will be used for training the model.
min_compression_perc	Numeric. An integer, indicating the minimum compression percentage to be achieved for the dataset. It indicates the desired level of reduction in dataset size compared to its original size.
n_cells	Numeric. An integer, indicating the number of cells per hierarchy (level).

depth	Numeric. An integer, indicating the number of levels. A depth of 1 means no hierarchy (single level), while higher values indicate multiple levels (hierarchy).
quant.err	Numeric. A number indicating the quantization error threshold. A cell will only breakdown into further cells if the quantization error of the cell is above the defined quantization error threshold.
normalize	Logical. A logical value indicating if the dataset should be normalized. When set to TRUE, scales the values of all features to have a mean of 0 and a standard deviation of 1 (Z-score).
distance_metric	Character. The distance metric can be L1_Norm(Manhattan) or L2_Norm(Eucledian). L1_Norm is selected by default. The distance metric is used to calculate the distance between an n dimensional point and centroid.
error_metric	Character. The error metric can be mean or max. max is selected by default. max will return the max of m values and mean will take mean of m values where each value is a distance between a point and centroid of the cell.
quant_method	Character. The quantization method can be kmeans or kmedoids. Kmeans uses means (centroids) as cluster centers while Kmedoids uses actual data points (medoids) as cluster centers. kmeans is selected by default.
scale_summary	List. A list with user-defined mean and standard deviation values for all the features in the dataset. Pass the scale summary when normalize is set to FALSE.
diagnose	Logical. A logical value indicating whether user wants to perform diagnostics on the model. Default value is FALSE.
hvt_validation	Logical. A logical value indicating whether user wants to holdout a validation set and find mean absolute deviation of the validation points from the centroid. Default value is FALSE.
train_validation_split_ratio	Numeric. A numeric value indicating train validation split ratio. This argument is only used when hvt_validation has been set to TRUE. Default value for the argument is 0.8.
dim_reduction_method	Character. The dim_reduction_method can be one of "tsne", "umap", "sammon".
tsne_theta	Numeric. The tsne_theta is only used when dim_reduction_method is set to "tsne". Default value is 0.5 and common values are between 0.2 and 0.5.
tsne_eta	Numeric. The tsne_eta are used only when dim_reduction method is set to "tsne". Default value is 200.
tsne_perplexity	Numeric. The tsne_perplexity is only used when dim_reduction_method is set to "tsne". Default value is 30 and common values are between between 30 and 50.
tsne_verbose	Logical. A logical value which indicates the t-SNE algorithm to print detailed information about its progress to the console.
tsne_max_iter	Numeric. The tsne_max_iter is used only when dim_reduction_method is set to "tsne". Default value is 1000. More iterations can improve results but increase computation time.

umap_n_neighbors	Integer. The <code>umap_n_neighbors</code> is used only when <code>dim_reduction_method</code> is set to "umap". Default value is 15. Controls the balance between local and global structure in data.
umap_n_components	Integer. The <code>umap_n_components</code> is used only when <code>dim_reduction_method</code> is set to "umap". Default value is 2. Indicates the number of dimensions for embedding.
umap_min_dist	Numeric. The <code>umap_min_dist</code> is used only when <code>dim_reduction_method</code> is set to "umap". Default value is 0.1. Controls how tightly UMAP packs points together.

### Value

A Nested list that contains the hierarchical tessellation information. This list has to be given as input argument to plot the tessellations.

[[1]]	A list containing information related to plotting tessellations. This information will include coordinates, boundaries, and other details necessary for visualizing the tessellations
[[2]]	A list containing information related to Sammon's projection coordinates of the data points in the reduced-dimensional space.
[[3]]	A list containing detailed information about the hierarchical vector quantized data along with a summary section containing no of points, Quantization Error and the centroids for each cell.
[[4]]	A list that contains all the diagnostics information of the model when <code>diagnose</code> is set to TRUE. Otherwise NA.
[[5]]	A list that contains all the information required to generate a Mean Absolute Deviation (MAD) plot, if <code>hvt_validation</code> is set to TRUE. Otherwise NA
[[6]]	A list containing detailed information about the hierarchical vector quantized data along with a summary section containing no of points, Quantization Error and the centroids for each cell which is the output of 'hvq'
[[7]]	model info: A list that contains model-generated timestamp, input parameters passed to the model, the validation results and the dimensionality reduction evaluation metrics table.

### Author(s)

Shubhra Prakash <shubhra.prakash@mu-sigma.com>, Sangeet Moy Das <sangeet.das@mu-sigma.com>, Shantanu Vaidya <shantanu.vaidya@mu-sigma.com>, Bidesh Ghosh <bidesh.gosh@mu-sigma.com>, Alimpan Dey <alimpan.dey@mu-sigma.com>

### See Also

[plotHVT](#)

**Examples**

```
data("EuStockMarkets")
hvt.results <- trainHVT(EuStockMarkets, n_cells = 60, depth = 1, quant.err = 0.1,
                        distance_metric = "L1_Norm", error_metric = "max",
                        normalize = TRUE, quant_method="kmeans")
```

# Index

- \* **Clustering\_Analysis**
    - clustHVT, 2
  - \* **Diagnostics\_or\_Validation**
    - plotModelDiagnostics, 13
    - plotQuantErrorHistogram, 15
  - \* **EDA**
    - edaPlots, 5
    - plotZscore, 17
  - \* **Novelty\_or\_Outliers**
    - plotNovelCells, 14
    - removeNovelty, 20
  - \* **Scoring**
    - scoreHVT, 21
    - scoreLayeredHVT, 23
  - \* **Table\_Formatting**
    - displayTable, 4
    - summary, 25
  - \* **Tessellation\_and\_Heatmap**
    - plotHVT, 11
  - \* **Timeseries\_Analysis**
    - getTransitionProbability, 6
    - msm, 7
    - plotAnimatedFlowmap, 9
    - plotStateTransition, 16
    - reconcileTransitionProbability, 18
  - \* **Training\_or\_Compression**
    - trainHVT, 26
- clustHVT, 2
- displayTable, 4
- edaPlots, 5
- getTransitionProbability, 6, 10
- msm, 7
- plotAnimatedFlowmap, 9
- plotHVT, 11, 13, 15, 22, 24, 28
- plotModelDiagnostics, 13
- plotNovelCells, 14
- plotQuantErrorHistogram, 15
- plotStateTransition, 16
- plotZscore, 17
- reconcileTransitionProbability, 18
- removeNovelty, 20
- scoreHVT, 10, 21
- scoreLayeredHVT, 20, 23
- summary, 25
- trainHVT, 10, 12, 15, 20, 22, 24, 26