

Package ‘netropy’

May 9, 2026

Title Statistical Entropy Analysis of Network Data

Version 0.3.0

Description Statistical entropy analysis of network data as introduced by Frank and Shafie (2016) <[doi:10.1177/0759106315615511](https://doi.org/10.1177/0759106315615511)>, and a in text-book which is in progress.

License MIT + file LICENSE

Encoding UTF-8

LazyData true

Imports ggraph, ggplot2, igraph

RoxygenNote 7.3.2

Language en-US

Depends R (>= 3.6)

Suggests testthat (>= 3.0.0), rmarkdown, knitr

URL <https://github.com/termehs/netropy>

Config/testthat/edition 3

VignetteBuilder knitr

NeedsCompilation no

Author Termeh Shafie [aut, cre]

Maintainer Termeh Shafie <termeh.shafie@uni-konstanz.de>

Repository CRAN

Date/Publication 2026-04-24 08:50:02 UTC

Contents

assoc_graph	2
div_gof	3
entropy_bivar	6
entropy_tetravar	8
entropy_trivar	10

get_dyad_var	11
get_triad_var	13
joint_entropy	15
lawdata	16
make_pred_plot	18
prediction_power	19
redundancy	21

Index	23
--------------	-----------

assoc_graph	<i>Association Graphs</i>
-------------	---------------------------

Description

Draws association graphs (graphical models) based on joint entropy values to detect and visualize different dependence structures among the variables in the dataframe.

Usage

```
assoc_graph(dat, cutoff = 0)
```

Arguments

dat	dataframe with rows as observations and columns as variables. Variables must all be observed or transformed categorical with finite range spaces.
cutoff	the cutoff point for the edges to be drawn based on joint entropies. Default is 0 and draws all edges.

Details

Draws association graphs based on given thresholds of joint entropy values between pairs of variables represented as nodes. Thickness of edges between pairs of nodes/variables indicates the strength of dependence between them. Isolated nodes are completely independent and paths through certain nodes/variables indicate conditional dependencies.

Value

A ggraph object with nodes representing all variables in dat and edges representing (the strength of) associations between them based on joint entropies.

Author(s)

Termeh Shafie

References

Frank, O., & Shafie, T. (2016). Multivariate entropy analysis of network data. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique*, 129(1), 45-63.

See Also[joint_entropy](#)**Examples**

```
library(ggraph)
# use internal data set
data(lawdata)
df.att <- lawdata[[4]]

# three steps of data editing:
# 1. categorize variables 'years' and 'age' based on
# approximately three equally size groups (values based on cdf)
# 2. make sure all outcomes start from the value 0 (optional)
# 3. remove variable 'senior' as it consists of only unique values (thus redundant)
df.att.ed <- data.frame(
  status = df.att$status,
  gender = df.att$gender,
  office = df.att$office - 1,
  years = ifelse(df.att$years <= 3, 0,
    ifelse(df.att$years <= 13, 1, 2)
  ),
  age = ifelse(df.att$age <= 35, 0,
    ifelse(df.att$age <= 45, 1, 2)
  ),
  practice = df.att$practice,
  lawschool = df.att$lawschool - 1
)

# association graph based on cutoff 0.15
assoc_graph(df.att.ed, 0.15)
```

div_gof

Divergence Tests of Goodness of Fit

Description

Performs divergence-based goodness-of-fit tests for discrete data, including tests of uniformity, pairwise independence, conditional independence, and nested model comparisons.

Usage

```
div_gof(
  dat,
  var_uniform = NULL,
  var1 = NULL,
  var2 = NULL,
  var_cond = NULL,
  model_full = NULL,
```

```

  model_reduced = NULL,
  alpha = 0.05,
  dec = 3,
  use_approx_cv = TRUE
)

```

Arguments

<code>dat</code>	dataframe with rows as observations and columns as variables. Variables must be categorical with finite range spaces.
<code>var_uniform</code>	character name of a variable in <code>dat</code> to test for uniformity.
<code>var1</code>	character name of the first variable.
<code>var2</code>	character name of the second variable.
<code>var_cond</code>	optional character vector of conditioning variables.
<code>model_full</code>	list containing D and df for the full model.
<code>model_reduced</code>	list containing D and df for the reduced model.
<code>alpha</code>	significance level. Default is 0.05.
<code>dec</code>	number of decimals for rounding. Default is 3.
<code>use_approx_cv</code>	logical; if TRUE, uses the approximate critical value $df + \sqrt{8 * df}$. If FALSE, uses the chi-square quantile.

Details

The function implements four types of tests:

1. Uniformity

$$D = \log r_X - H(X)$$

2. Pairwise Independence

$$D = H(X) + H(Y) - H(X, Y)$$

3. Conditional Independence

$$D = H(X, Z) + H(Y, Z) - H(Z) - H(X, Y, Z)$$

where Z may also represent a vector of conditioning variables.

4. Nested Model Comparison

$$D = D_{reduced} - D_{full}$$

The test statistic is

$$2nD \log(2),$$

since entropies are computed using base 2 logarithms.

Smaller divergence values indicate better model fit.

Value

Dataframe with test type, divergence D , chi-square statistic, degrees of freedom, critical value, and decision.

Author(s)

Termeh Shafie

References

Frank, O., & Shafie, T. (2016). Multivariate entropy analysis of network data. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique*, 129(1), 45-63.

See Also

[joint_entropy](#), [entropy_trivar](#)

Examples

```
data(lawdata)
df_att <- lawdata[[4]]

att_var <- data.frame(
  status = df_att$status - 1,
  gender = df_att$gender,
  office = df_att$office - 1,
  years = ifelse(df_att$years <= 3, 0,
                ifelse(df_att$years <= 13, 1, 2)),
  age = ifelse(df_att$age <= 35, 0,
              ifelse(df_att$age <= 45, 1, 2)),
  practice = df_att$practice,
  lawschool = df_att$lawschool - 1
)

## 1. Test uniformity
div_gof(att_var, var_uniform = "gender")

## 2. Test pairwise independence
div_gof(att_var, var1 = "status", var2 = "gender")

## 3. Test conditional independence

## (a) Conditional independence given a single variable
div_gof(att_var,
        var1 = "status",
        var2 = "gender",
        var_cond = "years")

## (b) Conditional independence given multiple variables
div_gof(att_var,
        var1 = "status",
        var2 = "gender",
        var_cond = c("years", "age"))

## 4. Nested model comparison
## Compare reduced models against the saturated empirical model.
## The saturated model has divergence  $D = 0$  and  $df = 0$ .
```

```

m_full <- list(D = 0, df = 0)

## (a) Pairwise independence model
m_reduced <- div_gof(att_var,
                    var1 = "status",
                    var2 = "gender")

div_gof(att_var,
        model_full = m_full,
        model_reduced = list(D = m_reduced$D, df = m_reduced$df))

## (b) Conditional independence model
m_reduced <- div_gof(att_var,
                    var1 = "status",
                    var2 = "gender",
                    var_cond = "years")

div_gof(att_var,
        model_full = m_full,
        model_reduced = list(D = m_reduced$D, df = m_reduced$df))

## 5. Nested comparison against the saturated empirical model
m_full <- list(D = 0, df = 0)

m_reduced <- div_gof(att_var,
                    var1 = "status",
                    var2 = "gender")

div_gof(att_var,
        model_full = m_full,
        model_reduced = list(D = m_reduced$D, df = m_reduced$df))

```

entropy_bivar

Bivariate Entropy

Description

Computes the bivariate entropies between all pairs of (discrete) variables in a multivariate data set.

Usage

```
entropy_bivar(dat)
```

Arguments

`dat` dataframe with rows as observations and columns as variables. Variables must all be observed or transformed categorical with finite range spaces.

Details

The bivariate entropy $H(X,Y)$ of two discrete random variables X and Y can be used to check for functional relationships and stochastic independence between pairs of variables. The bivariate entropy is bounded according to

$$H(X) \leq H(X,Y) \leq H(X) + H(Y)$$

where $H(X)$ and $H(Y)$ are the univariate entropies.

Value

Upper triangular matrix giving bivariate entropies between pairs of variables given as rows and columns of the matrix. The univariate entropies are given in the diagonal.

Author(s)

Termeh Shafie

References

Frank, O., & Shafie, T. (2016). Multivariate entropy analysis of network data. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique*, 129(1), 45-63.

See Also

[joint_entropy](#), [entropy_trivar](#), [redundancy](#), [prediction_power](#)

Examples

```
# use internal data set
data(lawdata)
df.att <- lawdata[[4]]

# three steps of data editing:
# 1. categorize variables 'years' and 'age' based on
# approximately three equally size groups (values based on cdf)
# 2. make sure all outcomes start from the value 0 (optional)
# 3. remove variable 'senior' as it consists of only unique values (thus redundant)
df.att.ed <- data.frame(
  status = df.att$status,
  gender = df.att$gender,
  office = df.att$office - 1,
  years = ifelse(df.att$years <= 3, 0,
    ifelse(df.att$years <= 13, 1, 2)
  ),
  age = ifelse(df.att$age <= 35, 0,
    ifelse(df.att$age <= 45, 1, 2)
  ),
  practice = df.att$practice,
  lawschool = df.att$lawschool - 1
```

```

)

# calculate bivariate entropies
H.biv <- entropy_bivar(df.att.ed)
# univariate entropies are then given as
diag(H.biv)

```

entropy_tetravar *Tetrivariate Entropy*

Description

Computes tetrivariate entropies, expected conditional entropies, and expected conditional joint entropies for all quadruples of variables in a multivariate discrete data set.

Usage

```
entropy_tetravar(dat, dec = 2)
```

Arguments

dat dataframe with rows as observations and columns as variables. Variables must be categorical with finite range spaces.

dec number of decimals used for rounding the entropy values. Default is 2.

Details

For four variables X , Y , Z , and U , the tetrivariate entropy is denoted $H(X,Y,Z,U)$. The expected conditional entropies are computed as

$$EH(U|X, Y, Z) = H(X, Y, Z, U) - H(X, Y, Z)$$

and

$$EH(Z|X, Y, U) = H(X, Y, Z, U) - H(X, Y, U).$$

The expected conditional joint entropy is computed as

$$EJ(X, Y|Z, U) = H(X, Z, U) + H(Y, Z, U) - H(Z, U) - H(X, Y, Z, U).$$

This quantity measures deviation from conditional independence of the form $X \perp Y | Z, U$. Smaller values indicate weaker conditional dependence.

Value

A dataframe with one row for each ordered decomposition of four variables into predictors and conditioning variables. The columns are:

X first variable in the pair of interest.

Y second variable in the pair of interest.

Z	first conditioning variable.
U	second conditioning variable.
H_XYZU	tetrivariate entropy $H(X,Y,Z,U)$.
EH_U_XYZ	expected conditional entropy $EH(U X,Y,Z)$.
EH_Z_XYU	expected conditional entropy $EH(Z X,Y,U)$.
EJ_XY_ZU	expected conditional joint entropy $EJ(X,Y Z,U)$.

Author(s)

Termeh Shafie

References

Frank, O., & Shafie, T. (2016). Multivariate entropy analysis of network data. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique*, 129(1), 45-63.

See Also

[entropy_trivar](#), [entropy_bivar](#), [prediction_power](#)

Examples

```
# use internal data set
data(lawdata)

# extract node attributes
df_att <- lawdata[[4]]

# data editing:
# 1. discretize 'years' and 'age' into three approximately balanced groups
# 2. recode selected variables so categories start at 0
att_var <- data.frame(
  status = df_att$status - 1,
  gender = df_att$gender,
  office = df_att$office - 1,
  years = ifelse(df_att$years <= 3, 0,
                 ifelse(df_att$years <= 13, 1, 2)),
  age = ifelse(df_att$age <= 35, 0,
               ifelse(df_att$age <= 45, 1, 2)),
  practice = df_att$practice,
  lawschool = df_att$lawschool - 1
)

# compute tetrivariate entropy quantities for five selected variables
entropy_tetravar(
  dat = att_var[, c("gender", "years", "age", "office", "practice")]
)
```

entropy_trivar	<i>Trivariate Entropy</i>
----------------	---------------------------

Description

Computes trivariate entropies of all triples of discrete variables in a multivariate data set.

Usage

```
entropy_trivar(dat)
```

Arguments

`dat` dataframe with rows as observations and columns as variables. Variables must all be observed or transformed categorical with finite range spaces.

Details

Trivariate entropies can be used to check for functional relationships and stochastic independence between triples of variables.

The trivariate entropy $H(X,Y,Z)$ of three discrete random variables X , Y , and Z is bounded according to

$$H(X, Y) \leq H(X, Y, Z) \leq H(X, Z) + H(Y, Z) - H(Z).$$

The increment between the trivariate entropy and its lower bound is equal to the expected conditional entropy.

Value

Dataframe with the first three columns representing possible triples of variables (X , Y , Z) and the fourth column giving trivariate entropies $H(X, Y, Z)$.

Author(s)

Termeh Shafie

References

Frank, O., & Shafie, T. (2016). Multivariate entropy analysis of network data. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique*, 129(1), 45-63.

See Also

[entropy_bivar](#), [prediction_power](#)

Examples

```

# use internal data set
data(lawdata)
df_att <- lawdata[[4]]

# data editing:
# 1. categorize variables 'years' and 'age' into approximately
# equally sized groups
# 2. recode selected variables so categories start at 0
att_var <- data.frame(
  status   = df_att$status - 1,
  gender   = df_att$gender,
  office   = df_att$office - 1,
  years    = ifelse(df_att$years <= 3, 0,
                    ifelse(df_att$years <= 13, 1, 2)),
  age      = ifelse(df_att$age <= 35, 0,
                    ifelse(df_att$age <= 45, 1, 2)),
  practice = df_att$practice,
  lawschool = df_att$lawschool - 1
)

# calculate trivariate entropies
h_trivar <- entropy_trivar(att_var)

```

get_dyad_var

Get Dyad Variables

Description

Transforms vertex variables or observed directed/undirected ties into dyad variables.

Usage

```
get_dyad_var(var, type = "att")
```

Arguments

var	variable vector (actor attribute) or adjacency matrix (ties) to be transformed to a dyad variable.
type	either 'att' for actor attribute (default) or 'tie' for relations.

Details

Dyad variables are given as pairs of incident vertex variables or actor attributes. Here, unique pairs of original attribute values constitute the outcome space. Note that the actor attributes need to be categorical with finite range spaces. For example, binary attribute yields outcome space (0,0), (0,1), (1,0), (1,1) coded as (0),(1),(2),(3). Warning message is shown if actor attribute has too many unique outcomes as it will yield too many possible outcomes once converted in to a dyad variable.

For directed relations, pairs of indicators from the adjacency matrix constitute the four outcomes representing possible combinations of sending and receiving ties: (0,0), (0,1), (1,0), (1,1) coded as (0),(1),(2),(3).

For undirected relations, an indicator variable which is directly read from the adjacency matrix represents the dyadic variable.

Value

Dataframe with three columns: first two columns show the vertex pairs u and v where $u < v$, and the third column gives the value of the transformed dyadic variable var .

Author(s)

Termeh Shafie

References

Frank, O., & Shafie, T. (2016). Multivariate entropy analysis of network data. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique*, 129(1), 45-63.

See Also

[get_triad_var](#)

Examples

```
# use internal data set
data(lawdata)
adj.advice <- lawdata[[1]]
adj.cowork <- lawdata[[3]]
df.att <- lawdata[[4]]

# three steps of data editing of attribute dataframe:
# 1. categorize variables 'years' and 'age' based on
# approximately three equally size groups (values based on cdf)
# 2. make sure all outcomes start from the value 0 (optional)
# 3. remove variable 'senior' as it consists of only unique values (thus redundant)
df.att.ed <- data.frame(
  status = df.att$status,
  gender = df.att$gender,
  office = df.att$office - 1,
  years = ifelse(df.att$years <= 3, 0,
    ifelse(df.att$years <= 13, 1, 2)
  ),
  age = ifelse(df.att$age <= 35, 0,
    ifelse(df.att$age <= 45, 1, 2)
  ),
  practice = df.att$practice,
  lawschool = df.att$lawschool - 1
```

```

)

# actor attribute converted to dyad variable
dyad.gend <- get_dyad_var(df.att.ed$gender, "att")

# directed tie converted to dyad variable
dyad.adv <- get_dyad_var(adj.advice, "tie")

# undirected tie converted to dyad variable
dyad.cwk <- get_dyad_var(adj.cowork, "tie")

```

get_triad_var

Get Triad Variables

Description

Transforms vertex variables or observed directed/undirected ties into triad variables.

Usage

```
get_triad_var(var, type = "att")
```

Arguments

var	variable vector (actor attribute) or adjacency matrix (ties) to be transformed to a triad variable.
type	either 'att' for actor attribute (default) or 'tie' for relations.

Details

For actor attributes, unique triples of original attribute values constitute the outcome space. Note that the actor attributes need to be categorical with finite range spaces. For example, binary attributes have 8 possible triadic outcomes (0,0,0),(1,0,0),(0,1,0),(1,1,0),(0,0,1),(1,0,1),(0,1,1),(1,1,1) which are coded 0-7. Warning message is shown if actor attribute has too many unique outcomes as it will yield too many possible outcomes once converted in to a triad variable.

For directed relations, a sequence of indicators of length 6 created from the adjacency matrix constitutes the 64 outcomes representing possible combinations of sending and receiving ties.

For undirected relations, triples of indicators are created from the adjacency matrix.

Value

Dataframe with four columns: first three columns show the vertex triad u, v, w , and the fourth column gives the value of the transformed triadic variable var.

Author(s)

Termeh Shafie

References

Frank, O., & Shafie, T. (2016). Multivariate entropy analysis of network data. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique*, 129(1), 45-63.

See Also

[get_dyad_var](#)

Examples

```
# use internal data set
data(lawdata)
adj.advice <- lawdata[[1]]
adj.cowork <- lawdata[[3]]
df.att <- lawdata[[4]]

# three steps of data editing:
# 1. categorize variables 'years' and 'age' based on
# approximately three equally size groups (values based on cdf)
# 2. make sure all outcomes start from the value 0 (optional)
# 3. remove variable 'senior' as it consists of only unique values (thus redundant)
df.att.ed <- data.frame(
  status = df.att$status,
  gender = df.att$gender,
  office = df.att$office - 1,
  years = ifelse(df.att$years <= 3, 0,
    ifelse(df.att$years <= 13, 1, 2)
  ),
  age = ifelse(df.att$age <= 35, 0,
    ifelse(df.att$age <= 45, 1, 2)
  ),
  practice = df.att$practice,
  lawschool = df.att$lawschool - 1
)

# actor attribute converted to triad variable
triad.gend <- get_triad_var(df.att.ed$gender, "att")

# directed tie converted to triad variable
triad.adv <- get_triad_var(adj.advice, type = "tie")

# undirected tie converted to triad variable
triad.cwk <- get_triad_var(adj.cowork, type = "tie")
```

joint_entropy	<i>Joint Entropy</i>
---------------	----------------------

Description

Computes the joint entropies between all pairs of (discrete) variables in a multivariate data set.

Usage

```
joint_entropy(dat, dec = 3)
```

Arguments

dat	dataframe with rows as observations and columns as variables. Variables must all be observed or transformed categorical with finite range spaces.
dec	the precision given in number of decimals for which the frequency distribution of unique entropy values is created. Default is 3.

Details

The joint entropy $J(X,Y)$ of discrete variables X and Y is a measure of dependence or association between them, defined as

$$J(X,Y) = H(X) + H(Y) - H(X,Y).$$

Two variables are independent if their joint entropy, i.e. their mutual information, is equal to zero. The frequency distributions can be used to decide upon convenient thresholds for constructing association graphs.

Value

List with	
matrix	an upper triangular joint entropy matrix (univariate entropies in the diagonal).
freq	a dataframe giving the frequency distributions of unique joint entropy values.

Author(s)

Termeh Shafie

References

Frank, O., & Shafie, T. (2016). Multivariate entropy analysis of network data. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique*, 129(1), 45-63.

See Also

[assoc_graph](#), [entropy_bivar](#)

Examples

```
# use internal data set
data(lawdata)
df.att <- lawdata[[4]]

# three steps of data editing:
# 1. categorize variables 'years' and 'age' based on
# approximately three equally size groups (values based on cdf)
# 2. make sure all outcomes start from the value 0 (optional)
# 3. remove variable 'senior' as it consists of only unique values (thus redundant)
df.att.ed <- data.frame(
  status = df.att$status,
  gender = df.att$gender,
  office = df.att$office - 1,
  years = ifelse(df.att$years <= 3, 0,
    ifelse(df.att$years <= 13, 1, 2)
  ),
  age = ifelse(df.att$age <= 35, 0,
    ifelse(df.att$age <= 45, 1, 2)
  ),
  practice = df.att$practice,
  lawschool = df.att$lawschool - 1
)

# calculate joint entropies
J <- joint_entropy(df.att.ed)
# joint entropy matrix
J$matrix
# frequency distribution of joint entropy values
J$freq
```

lawdata

Law Firm

Description

This data set comes from a network study of corporate law partnership that was carried out in a Northeastern US corporate law firm, referred to as SG&R, 1988-1991 in New England. It includes (among others) measurements of networks among the 71 attorneys (partners and associates) of this firm, i.e. their strong- co-worker network, advice network, friendship network, and indirect control networks. Various members' attributes are also part of the data set, including seniority, formal status, office in which they work, gender, law school attended. The ethnography, organizational and network analyses of this case are available in Lazega (2001).

Basic advice network: "Think back over the past year, consider all the lawyers in your Firm. To whom did you go for basic professional advice? For instance, you want to make sure that you are handling a case right, making a proper decision, and you want to consult someone whose professional opinions are in general of great value to you. By advice I do not mean simply technical advice."

Friendship network: "Would you go through this list, and check the names of those you socialize with outside work. You know their family, they know yours, for instance. I do not mean all the people you are simply on a friendly level with, or people you happen to meet at Firm functions."

Strong coworkers network: "Because most firms like yours are also organized very informally, it is difficult to get a clear idea of how the members really work together. Think back over the past year, consider all the lawyers in your Firm. Would you go through this list and check the names of those with whom you have worked with. (By "worked with" I mean that you have spent time together on at least one case, that you have been assigned to the same case, that they read or used your work product or that you have read or used their work product; this includes professional work done within the Firm like Bar association work, administration, etc.)"

Usage

```
data(lawdata)
```

Format

List containing the following objects as numbered

1. adjacency matrix for advice seeking (directed)
2. adjacency matrix for friendship (directed)
3. adjacency matrix for cowork (undirected)
4. dataframe with the following attributes on each lawyer:
 - 'senior': seniority (ranked from most to least senior)
 - 'status': 1=partner; 2=associate
 - 'gender': 1=man; 2=woman
 - 'office': 1=Boston; 2=Hartford; 3=Providence
 - 'years': years with the firm
 - 'age': age of attorney
 - 'practice': 1=litigation; 2=corporate
 - 'lawschool': 1=harvard/yale; 2=ucon; 3= other

Note: the first 36 out of 71 respondents are the partners in the firm.

Source

https://www.stats.ox.ac.uk/~snijders/siena/Lazega_lawyers_data.htm

References

Emmanuel Lazega, *The Collegial Phenomenon: The Social Mechanisms of Cooperation Among Peers in a Corporate Law Partnership*, Oxford University Press (2001).

Tom A.B. Snijders, Philippa E. Pattison, Garry L. Robins, and Mark S. Handcock. New specifications for exponential random graph models. *Sociological Methodology* (2006), 99-153.

Examples

```
data(lawdata)
## assign the correct names to the objects in the list
adj.advice <- lawdata[[1]]
adj.friend <- lawdata[[2]]
adj.cowork <- lawdata[[3]]
df.att <- lawdata[[4]]
```

make_pred_plot

Prediction Power Heatmap

Description

Creates a heatmap for visualizing prediction power from a prediction power matrix.

Usage

```
make_pred_plot(mat, title, low = "azure4", high = "white", text_size = 2.5)
```

Arguments

mat	matrix returned by prediction_power . Entries should contain expected conditional entropies $EH(Z X, Y)$.
title	character string giving the plot title.
low	color for low expected conditional entropy values. Default is "steelblue".
high	color for high expected conditional entropy values. Default is "white".
text_size	numeric value controlling the size of the cell labels. Default is 2.5.

Details

The plot visualizes expected conditional entropies

$$EH(Z|X, Y)$$

where Z is the target variable and X and Y are predictors. Diagonal entries correspond to prediction using a single predictor, $EH(Z|X)$, while off-diagonal entries correspond to prediction using pairs of predictors, $EH(Z|X, Y)$. Lower values indicate stronger predictive power.

Value

A ggplot object showing a heatmap of expected conditional entropy values. Darker cells indicate lower prediction uncertainty and therefore higher prediction power.

See Also

[prediction_power](#), [entropy_trivar](#)

Examples

```
# use internal data set
data(lawdata)

# extract node attributes
df_att <- lawdata[[4]]

# data editing:
# 1. discretize 'years' and 'age' into 3 categories
# 2. ensure values start at 0 where needed
att_var <- data.frame(
  status = df_att$status - 1,
  gender = df_att$gender,
  office = df_att$office - 1,
  years = ifelse(df_att$years <= 3, 0,
                ifelse(df_att$years <= 13, 1, 2)),
  age = ifelse(df_att$age <= 35, 0,
              ifelse(df_att$age <= 45, 1, 2)),
  practice = df_att$practice,
  lawschool = df_att$lawschool - 1
)

# compute prediction power matrix for 'status'
pred_mat <- prediction_power("status", att_var)

# visualize prediction power
make_pred_plot(pred_mat, "Prediction Power for Status")
```

prediction_power	<i>Prediction Power</i>
------------------	-------------------------

Description

Computes prediction power when pairs of variables in a given dataframe are used to predict a third variable from the same dataframe. The prediction strength is measured by expected conditional entropies.

Usage

```
prediction_power(var, dat)
```

Arguments

var	character string representing the variable in dataframe dat to be predicted by pairs of other variables in the dataframe dat.
dat	dataframe with rows as observations and columns as variables. Variables must all be observed or transformed categorical with finite range spaces.

Details

The expected conditional entropy given by

$$EH(Z|X,Y) = H(X,Y,Z) - H(X, Y)$$

measures the prediction uncertainty when pairs of variables X and Y are used to predict variable Z . The lower the value of EH given different pairs of variables X and Y , the stronger is the prediction of Z .

Value

Upper triangular matrix giving the expected conditional entropies of pairs of variables given as rows and columns of the matrix. The diagonal gives $EH(Z|X) = H(X,Z) - H(X)$, that is when only one variable is used to predict var. Note that NA's are in the entire row and column representing the variable being predicted.

Author(s)

Termeh Shafie

References

Frank, O., & Shafie, T. (2016). Multivariate entropy analysis of network data. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique*, 129(1), 45-63.

See Also

[entropy_trivar](#), [entropy_bivar](#)

Examples

```
# use internal data set
data(lawdata)
df.att <- lawdata[[4]]

# three steps of data editing:
# 1. categorize variables 'years' and 'age' based on
# approximately three equally size groups (values based on cdf)
# 2. make sure all outcomes start from the value 0 (optional)
# 3. remove variable 'senior' as it consists of only unique values (thus redundant)
df.att.ed <- data.frame(
  status = df.att$status,
  gender = df.att$gender,
  office = df.att$office - 1,
  years = ifelse(df.att$years <= 3, 0,
    ifelse(df.att$years <= 13, 1, 2)
  ),
  age = ifelse(df.att$age <= 35, 0,
    ifelse(df.att$age <= 45, 1, 2)
  ),
```

```
practice = df.att$practice,  
lawschool = df.att$lawschool - 1  
)  
  
# power of predicting 'status' using pairs of other variables  
prediction_power("status", df.att.ed)
```

redundancy

Redundant Variables & Dimensionality Reduction

Description

Finds redundant variables in a dataframe consisting of discrete variables.

Usage

```
redundancy(dat, dec = 3)
```

Arguments

dat	dataframe with rows as observations and columns as variables. Variables must all be observed or transformed categorical with finite range spaces.
dec	the precision given as number of decimals used to round bivariate entropies in order to find redundant variables (the more decimals, the harder to detect redundancy). Default is 3.

Details

Redundancy is defined as two variables holding the same information (bivariate entropies) as at least one of the variable alone (univariate entropies). Consider removing one of these two variable from the dataframe for further analysis.

Value

Binary matrix indicating which row and column variables hold the same information.

Author(s)

Termeh Shafie

References

Frank, O., & Shafie, T. (2016). Multivariate entropy analysis of network data. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique*, 129(1), 45-63.

See Also

[entropy_bivar](#),

Examples

```
# use internal data set
data(lawdata)
df.att <- lawdata[[4]]

# two steps of data editing:
# 1. categorize variables 'years' and 'age' based on
# approximately three equally size groups (values based on cdf)
# 2. make sure all outcomes start from the value 0 (optional)
df.att.ed <- data.frame(
  senior = df.att$senior,
  status = df.att$status,
  gender = df.att$gender,
  office = df.att$office - 1,
  years = ifelse(df.att$years <= 3, 0,
    ifelse(df.att$years <= 13, 1, 2)
  ),
  age = ifelse(df.att$age <= 35, 0,
    ifelse(df.att$age <= 45, 1, 2)
  ),
  practice = df.att$practice,
  lawschool = df.att$lawschool - 1
)

# find redundant variables in dataframe
redundancy(df.att.ed) # variable 'senior' should be omitted
```

Index

* datasets

lawdata, [16](#)

assoc_graph, [2](#), [15](#)

div_gof, [3](#)

entropy_bivar, [6](#), [9](#), [10](#), [15](#), [20](#), [21](#)

entropy_tetravar, [8](#)

entropy_trivar, [5](#), [7](#), [9](#), [10](#), [18](#), [20](#)

get_dyad_var, [11](#), [14](#)

get_triad_var, [12](#), [13](#)

joint_entropy, [3](#), [5](#), [7](#), [15](#)

lawdata, [16](#)

make_pred_plot, [18](#)

prediction_power, [7](#), [9](#), [10](#), [18](#), [19](#)

redundancy, [7](#), [21](#)