

# Package ‘hmm.discnp’

August 7, 2019

**Version** 2.1-11

**Date** 2019-08-07

**Title** Hidden Markov Models with Discrete Non-Parametric Observation Distributions

**Author** Rolf Turner

**Maintainer** Rolf Turner <r.turner@auckland.ac.nz>

**Depends** R (>= 2.10)

**Imports** nnet (>= 7.3.12)

**Description** Fits hidden Markov models with discrete non-parametric observation distributions to data sets. The observations may be univariate or bivariate. Simulates data from such models. Finds most probable underlying hidden states, the most probable sequences of such states, and the log likelihood of a collection of observations given the parameters of the model. Auxiliary predictors are accommodated in the univariate setting.

**LazyData** true

**ByteCompile** true

**License** GPL (>= 2)

**NeedsCompilation** yes

**Repository** CRAN

**Date/Publication** 2019-08-07 12:40:02 UTC

## R topics documented:

anova.hmm.discnp . . . . .	2
ccprSim . . . . .	3
cnvrtRho . . . . .	4
fitted.hmm.discnp . . . . .	5
hmm . . . . .	7
hydroDat . . . . .	16

lesionCount . . . . .	18
logLikHmm . . . . .	19
misstify . . . . .	21
mps . . . . .	23
nafracCalc . . . . .	24
pr . . . . .	25
predict.hmm.discnp . . . . .	27
rhmm . . . . .	28
scovmat . . . . .	30
sp . . . . .	32
squantCI . . . . .	33
SydColDisc . . . . .	35
update.hmm.discnp . . . . .	36
viterbi . . . . .	38
weissData . . . . .	40

**Index** **43**

anova.hmm.discnp      *Anova for hmm.discnp models*

**Description**

Performs a likelihood ratio test to compare two discrete non-parametric hidden Markov models.

**Usage**

```
## S3 method for class 'hmm.discnp'
anova(object, ...)
```

**Arguments**

object      An object of class “hmm.discnp” as returned by the function [hmm\(\)](#).  
 ...      A second object of class “hmm.discnp”. There must be only *one* such object.

**Value**

A list with entries

stat      The likelihood ratio statistic.  
 df      The degrees of freedom.  
 pvalue      The p-value of the test.

This list has an attribute “details” which is a vector consisting of the first and second log likelihoods and the associated numbers of parameters, in order of these numbers of parameters. (See **Warning**.)

**Warning**

Hidden Markov models can be numerically delicate and the fitting algorithm can converge to a local maximum of the likelihood surface which is not the global maximum. Thus it is entirely possible for the log likelihood of the model with the greater number of parameters to be *smaller* than that of the model with the lesser number of parameters.

**Author(s)**

Rolf Turner <r.turner@auckland.ac.nz>

**See Also**

[hmm\(\)](#)

**Examples**

```
xxx <- with(SydColDisc, split(y, f=list(locn, depth)))
fit1 <- hmm(xxx, K=1, itmax=10)
fit2 <- hmm(xxx, K=2, itmax=10)
anova(fit1, fit2)
```

---

ccprSim

*Simulated monocyte counts and psychosis symptoms.*

---

**Description**

Discretised values of monocyte counts, and ratings of level of psychosis simulated from a model fitted to a data set consisting of observations made on a number of patients from the Northern District Health Board system. The real data must be kept confidential due to ethics constraints.

**Usage**

```
data("ccprSim")
```

**Format**

The object `ccprSim` is a list of length 1258. Each entry of this list is to be considered to correspond to an individual subject. The entries consist of matrices having two columns named `cellCount` and `psychosisRating`. The number of rows of these matrices varies from entry to entry of the list (i.e. from subject to subject).

Most of the entries of these matrices are NA. The entries are temporally ordered and correspond to the number of weeks from the start of observation. Observations in the real data set were made only when the patient in question visited a physician and so weeks in which no visit was made resulted in an "observation" of NA. The object `ccprSim` was simulated in such a way as to imitate this characteristic. The fraction of missing observations in each variate (i.e. `cellCount` and `psychosisRating`) is roughly commensurate with the corresponding fractions in the real data.

The values in the first column of each matrix (the cellCount column consist of integers from 1 to 5 and are to be interpreted as indicators of cell counts in units of  $10^9$  cells per litre, discretised according to the following scale:

- $0.0 \leq c \leq 0.3 \leftrightarrow 1$
- $0.3 < c \leq 0.5 \leftrightarrow 2$
- $0.5 < c \leq 0.7 \leftrightarrow 3$
- $0.7 < c \leq 1.0 \leftrightarrow 4$
- $1.0 < c \leq 2.0 \leftrightarrow 5$

where  $c$  represents “count”.

The values in the second column of each matrix (the psychosisRating column consist of integers from 0 to 4 and are to be interpreted as indicators of a physician’s assessment of the level of psychosis of the patient. A value of 0 corresponds to “no symptoms”; a value of 4 corresponds to “severe”.

The question of essential interest in respect of the real data was “Is there any association between the cell count values and the psychosis ratings?” More specifically it was “Can the level of psychosis be *predicted* from the cell counts?”

## Source

The real data, on the basis of which these data were simulated, were supplied by Dr. Jonathan Williams, Northern District Health Board.

## Examples

```
## Not run:
# Takes a long time:
  fit <- hmm(ccprSim,K=2,indep=FALSE,itmax=500,verbose=TRUE)
# Throws an error; inadequate data:
  subsetFit <- hmm(ccprsim[1021:1025],K=2,indep=FALSE,itmax=5)

## End(Not run)
```

---

cnvrtRho

*Convert Rho between forms.*

---

## Description

Converts the “old style” (matrix) specification of the emission probabilities to the “new style” (data frame) specification, or vice versa.

## Usage

```
cnvrtRho(Rho)
```

**Arguments**

Rho A specification of the emission probabilities of a discrete valued hidden Markov model. It may be either a matrix of these probabilities, in which case it is converted to a three column data frame, or it may be a three column data frame, in which case it is converted to a matrix of probabilities. See [hmm\(\)](#) for more details about the structure of Rho, in either form.

**Details**

The “new style” specification of Rho allows for predictor variables  $x$ . If Rho is of the “new style” and is designed to allow for predictor variables, then it will have more than three columns and cannot be converted to the “old style”. In such cases `cnvrtRho` will throw an error.

**Value**

A data frame if the argument Rho is a matrix (“old style” specification), or a matrix if the argument Rho is a data frame (“new style” specification).

**Author(s)**

Rolf Turner <r.turner@auckland.ac.nz>

**See Also**

[hmm\(\)](#)

**Examples**

```
Yval <- LETTERS[1:10]
Tpm <- matrix(c(0.75,0.25,0.25,0.75),ncol=2,byrow=TRUE)
Rho <- cbind(c(rep(1,5),rep(0,5)),c(rep(0,5),rep(1,5)))/5
rownames(Rho) <- Yval
newRho <- cnvrtRho(Rho)
oldRho <- cnvrtRho(newRho)
```

---

fitted.hmm.discnp

*Fitted values of a discrete non-parametric hidden Markov model.*

---

**Description**

Calculates the fitted values of a discrete non-parametric hidden Markov model. If the data are numeric these are the conditional expectations of the observations, given the entire observation sequence (and the estimated parameters of the model). If the data are categorical (whence “expectations” make no sense) the “fitted values” are taken to be the probabilities of each of the possible values of the observations, at each time point.

**Usage**

```
## S3 method for class 'hmm.discnp'
fitted(object, warn=TRUE, drop=TRUE, ...)
```

**Arguments**

object	An object of class <code>hmm.discnp</code> as returned by <code>hmm()</code> .
warn	Logical scalar. See the help for <code>sp()</code> .
drop	Logical scalar. If there is a single sequence of observations (i.e. if <code>object[["y"]]</code> consists of a matrix or a list of length 1) and if <code>drop</code> is <code>TRUE</code> then the returned value is a single entity (matrix, list of two matrices, or 3-dimensional array, depending on circumstances. Otherwise the returned value is a list of such entities, one for each observation sequence.
...	Not used.

**Details**

The observation sequence(s) must be present in `object` (which will be the case if `object` was returned by `hmm()` and if the argument `keep.y` was set to `TRUE`). If it is not present an error is thrown.

However, if such an error is thrown, do not despair! You *do not* have to start from scratch when fitting your model with `keep.y==TRUE`. If `fit` is your fitted model that you obtained *without* setting `keep.y==TRUE`, then you can just re-fit the model using `fit` as the starting values:

```
fit2 <- hmm(<whatever>, par0=fit, keep.y=TRUE)
```

This will of course converge instantaneously. You could also do:

```
fit2 <- update(fit, data=<whatever>, keep.y=TRUE)
```

**Value**

If the observations consist of a single sequence and if `drop` is `TRUE` then the returned value consists of a single object (matrix, list of two matrices, or 3-dimensional array, depending on circumstances; see below). Otherwise the returned value is a list of such objects, one for each observation sequence.

If the observations are numeric then the object corresponding to each observation sequence is a matrix. If the model is univariate (see `hmm()`) then matrix has a single column constituting the sequence of fitted values corresponding to the observations in the given sequence. The number of rows is the number of observations and the entry in row `t` is the fitted value (conditional expectation) corresponding to the observation made at time `t`. If the model is bivariate (either independent or dependent) then the matrix has two columns corresponding respectively to the two variables in the bivariate model.

If the observations are categorical then the nature of the object returned changes substantially depending on whether the data are univariate, bivariate independent or bivariate dependent. (See `hmm()`).

In the univariate case the object corresponding to each sequence is a matrix, the number of rows of which is the number of observations and the number of columns of which is the number of unique *possible* values of the observations. The entry of this matrix in row  $t$  and column  $j$  is the conditional probability that an emission, at time  $t$ , is equal to  $u_j$  where  $u_1, \dots, u_m$  are the unique possible values.

In the bivariate independent case the object is a *list* of two matrices, each of which is of the same nature as that produced in the univariate case, corresponding respectively to the first and second of the two variables.

In the bivariate dependent case the object is a 3-dimensional array of dimension  $m_1 \times m_2 \times n$  where  $m_1$  is the number of unique possible values of the first variable,  $m_2$  is the number of unique possible values of the second variable, and  $n$  is the number of observations. The  $(i, j, t)$ -th entry of this array is the conditional probability that an emission, at time  $t$ , is equal to  $(u_i, v_j)$  where the  $u_i$  are the unique possible values of the first variable and the  $v_j$  are the unique possible values of the second variable.

### Author(s)

Rolf Turner <r.turner@auckland.ac.nz>

### See Also

`sp()` `link{predict.hmm.discnp}()`

### Examples

```
P <- matrix(c(0.7,0.3,0.1,0.9),2,2,byrow=TRUE)
R <- matrix(c(0.5,0,0.1,0.1,0.3,
             0.1,0.1,0,0.3,0.5),5,2)
set.seed(42)
l11 <- sample(250:350,20,TRUE)
y <- rhmm(ylengths=l11,nsim=1,drop=TRUE,tpm=P,Rho=R)
fit <- hmm(y,K=2,verb=TRUE,keep.y=TRUE,itmax=10)
fv <- fitted(fit)
```

---

hmm

*Fit a hidden Markov model to discrete data.*

---

### Description

Effects a maximum likelihood fit of a hidden Markov model to discrete data where the observations come from one of a number of finite discrete distributions, depending on the (hidden) state of the Markov chain. These distributions (the “emission probabilities”) are specified non-parametrically. The observations may be univariate, independent bivariate, or dependent bivariate. By default this function uses the EM algorithm. In the univariate setting it may alternatively use a “brute force” method.

## Usage

```
hmm(y, yval=NULL, par0=NULL, K=NULL, rand.start=NULL,
     method=c("EM", "bf", "LM", "SD"), hglmeth=c("fortran", "oraw", "raw"),
     optimiser=c("nlm", "optim"), optimMethod=NULL, stationary=cis,
     mixture=FALSE, cis=TRUE, indep=NULL, tolerance=1e-4, digits=NULL,
     verbose=FALSE, itmax=200, crit=c("PCLL", "L2", "Linf", "ABSGRD"),
     keep.y=FALSE, keep.X=keep.y, newstyle=TRUE, X=NULL,
     addIntercept=TRUE, lmc=10, hessian=FALSE,...)
```

## Arguments

- |            |   |
|------------|---|
| y          | A one or two column matrix of discrete data or a list of such matrices; missing values are allowed. If y is a vector, or list of vectors (of discrete data) these vectors are coerced to one column matrices.   |
| yval       | A vector (of length m, say) of possible values for the data or a list of two such vectors (of lengths m1 and m2, say, one for each of the two variates in the bivariate settings). These vectors default to the sorted unique values of the respective variates as provided in y. If yval is supplied and any value of y does not match some value of yval, then an error is thrown.  |
| par0       | <p>An optional (<i>named</i>) list of starting values for the parameters of the model, with components tpm (transition probability matrix), <i>optionally</i> ispd (initial state probability distribution) and Rho. The object Rho specifies the probability that the observations take on each of the possible values of the variate or variates, given the state of the hidden Markov chain. See <b>Details</b>. Note that in the case of independent bivariate data Rho is a list of two matrices. These matrices may (and in general will) have different row dimensions, but must have identical column dimensions (equal to K, the number of states; see below).</p> <p>If the model is stationary (i.e. if <code>stationary</code> is TRUE) then you should almost surely not specify the <code>ispd</code> component of <code>par0</code>. If you do specify it, it really only makes sense to specify it to be the stationary distribution determined by <code>tpm</code> and this is a waste of time since this is what the code will take <code>ispd</code> to be if you leave it unspecified.</p> <p>If <code>par0</code> is not specified, starting values are created by the function <code>init.all()</code>.</p> |
| K          | <p>The number of states in the hidden Markov chain; if <code>par0</code> is not specified K <i>MUST</i> be; if <code>par0</code> is specified, K is ignored.</p> <p>Note that K=1 is acceptable; if K is 1 then all observations are treated as being independent and the non-parametric estimate of the distribution of the observations is calculated in the dQuoteobvious way.</p>   |
| rand.start | A list consisting of two logical scalars which must be named <code>tpm</code> and <code>Rho</code> , if <code>tpm</code> is TRUE then the function <code>init.all()</code> chooses entries for then starting value of <code>tpm</code> at random; likewise for <code>Rho</code> . This argument defaults to <code>list(tpm=FALSE, Rho=FALSE)</code> .   |
| method     | Character string, either "bf", "EM", "LM" or "SD" (i.e. use numerical maximisation via either <code>nlm()</code> or <code>optim()</code> , the EM algorithm, the Levenberg-Marquardt algorithm, or the method of steepest descent). May be abbreviated. Currently the "bf", "LM" and "SD" methods can be used only in the univariate setting, handle only stationary models (see below) and do not do mixtures.   |

hglmethod	Character string; one of "fortran", "oraw" or "raw". May be abbreviated. This argument determines the procedure by which the hessian, gradient and log likelihood of the model and data are calculated. If this argument is equal to "fortran" (the default) then (obviously!) dynamically loaded fortran sub-routines are used. The other two possibilities effect the calculations in raw R; "oraw" ("o" for "original" uses code that is essentially a direct transcription of the fortran code, do-loops being replaced by for-loops. With method "raw" the for-loops are eliminated and matrix-vector calculations are applied. The "oraw" method is about 25 times slower than the "fortran" method and the "raw" method is more than 30 times slower. The "raw" methods are present mainly for debugging purposes and would not usually be used in practice. This argument is used only if the method is "LM" or "SD" (and is involved only peripherally in the latter instance). It is ignored otherwise.
optimiser	Character string specifying the optimiser to use when the "bf" method of optimisation is chosen. It should be one of "nlm" or "optim", and may be abbreviated. Ignored unless method="bf".
optimMethod	Character string specifying the optimisation method to be used by <code>optim()</code> . Should be one of "Nelder-Mead", "BFGS", "CG", "L-BFGS-B", "SANN", or "Brent". Ignored if the method is not "bf" or if the optimiser is not "optim".
stationary	Logical scalar. If TRUE then the model is fitted under the stationarity assumption, i.e. that the Markov chain was in steady state at the time that observations commenced. In this case the initial state probability distribution is estimated as the stationary distribution determined by the (estimated) transition probability matrix. Otherwise if <code>cis</code> (see below) is TRUE the initial state probability distribution is estimated as the mean of the vectors of conditional probabilities of the states, given the observation sequences, at time $t=1$ . If stationary is TRUE and <code>cis</code> is FALSE an error is thrown. Currently if the method is "bf", "LM" or "SD", and stationary is FALSE, then an error is thrown.
mixture	A logical scalar; if TRUE then a mixture model (all rows of the transition probability matrix are identical) is fitted rather than a general hidden Markov model. Currently an error is thrown if mixture=TRUE and the method is "bf", "LM" or "SD".
cis	A logical scalar specifying whether there should be a <b>constant initial state</b> probability distribution. If stationary is FALSE and <code>cis</code> is FALSE then the initial state probability distribution for a given observation sequence is equal to 1 where the (first) maximum of the vector of conditional probabilities of the states, given the observation sequences, at time $t=1$ , occurs, and is 0 elsewhere. If stationary is TRUE and <code>cis</code> is FALSE an error is given.
indep	Logical scalar. Should the bivariate model be fitted under the assumption that the two variables are (conditionally) independent given the state? If this argument is left as NULL its value is inferred from the structure of <code>Rho</code> in <code>par0</code> if the latter is supplied. If the data are bivariate and neither <code>indep</code> nor <code>par0</code> is supplied, then an error is given. If the data are bivariate and if the value of <code>indep</code> is inconsistent with the structure of <code>par0\$Rho</code> then an error is given. If the data are univariate then <code>indep</code> is ignored.
tolerance	If the value of the quantity used for the stopping criterion is less than tolerance then the algorithm is considered to have converged. Ignored if method="bf".

	Defaults to 1e-4.
digits	Integer scalar. The number of digits to which to print out “progress reports” (when verbose is TRUE). There is a “sensible” default (calculated from tolerance). Not used if the method is “bf”.
verbose	A logical scalar determining whether to print out details of the progress of the algorithm. If the method is “EM”, “LM” or “SD” then when verbose is TRUE information about the convergence criteria is printed out at every step that the algorithm takes. If method=“bf” then the value of verbose determines the value of the argument print.level of <code>nlm()</code> or the value of the argument trace of <code>optim()</code> . In the first case, if verbose is TRUE then print.level is set to 2, otherwise it is set to 0. In the second case, if verbose is TRUE then trace is set to 6, otherwise it is set to 0.
itmax	When the method is “EM”, “LM” or “SD” this is the maximum number of steps that the algorithm takes. If the convergence criterion has not been met by the time itmax steps have been performed, a warning message is printed out, and the function stops. A value is returned by the function anyway, with the logical component “converged” set to FALSE. When method=“bf” the itmax argument is passed to <code>nlm()</code> as the value of iterlim or to <code>optim()</code> as the value of maxit. If the (somewhat obscure) convergence criteria of <code>nlm()</code> or <code>optim()</code> have not been met by the time itmax “iterations” have been performed, the algorithm ceases. In this case, if <code>nlm()</code> is used. the value of code in the object returned set equal to 4 and if <code>optim()</code> is used then the value of convergence returned is set equal to 1. Note that the value of code, respectively convergence is returned as the converged component of the object returned by <code>hmm()</code> . A value of 1 indicates successful completion of the <code>nlm()</code> procedure. A value of 0 indicates successful completion of the <code>optim()</code> procedure.
crit	The name of the stopping criterion used. When method=“EM” it must be one of “PCLL” (percent change in log-likelihood; the default), “L2” (L-2 norm, i.e. square root of sum of squares of change in coefficients), or “Linf” (L-infinity norm, i.e. maximum absolute value of change in coefficients). When method=“LM” or method=“SD” there is a fourth possibility, namely “ABSGRD” the (maximum) absolute value of the gradient. It may not be advisable to use this criterion in the current context (i.e. that of discrete non-parametric distributions). See <b>Warnings</b> . This argument defaults to “PCLL”. It is ignored if method=“bf”. (The <code>nlm()</code> and <code>optim()</code> functions have their own obscure stopping criteria.)
keep.y	Logical scalar; should the observations y be returned as a component of the value of this function?
keep.X	Logical scalar; should the predictors X be returned as a component of the value of this function? Note that the value of keep.X will be silently set equal to FALSE unless it actually “makes sense” to keep X. I.e. unless the observations are <i>univariate</i> , newstyle is TRUE, and X is actually supplied, i.e. is not NULL.
newstyle	Logical scalar; if TRUE then the “new style” of of parameterising Rho (see <b>Details</b> ) is used, in the univariate setting. Note that if X is not NULL or if the method is not “EM” then newstyle <i>cannot</i> be set to FALSE; if it is, then an error is thrown.
X	An optional <i>numeric</i> matrix, or a list of such matrices, of <i>predictors</i> . The use of such predictors is (currently, at least) applicable only in the univariate emissions

setting, and then only if Rho is supplied in the “newstyle” format. If  $X$  is a list it must be of the same length as  $y$  and all entries of this list must have the same number of columns. If the columns of any entry of the list are named, then they must be named for *all* entries, and the column names must be the *same* for all entries. The number of rows of each entry must be equal to the length of the corresponding entry of  $y$ . If  $X$  is a matrix then  $y$  should be a vector or one-column matrix (or a list with a single entry equal to such).

There may be at most one constant column in  $X$  or the components thereof. If there are *any* constant columns there must be precisely one (in all components of  $X$ ), it must be the first column and all of its entries must be equal to 1. If the columns have names, the names of this first column must be “Intercept”.

Note that  $X$  (or its entries) must be a *numeric* matrix (or must be matrices) — no data frames! Factor predictors are not permitted. It may be possible to use factor predictors by supplying  $X$  or its entries as the output of `model.matrix()`; this will depend on circumstances.

<code>addIntercept</code>	Logical scalar. Should a column of ones, corresponding to an intercept term, be prepended to each of the matrices in the list $X$ ? (The user should remember that this argument defaults to TRUE.)
<code>lmc</code>	Numeric scalar. The (initial) “Levenberg-Marquardt correction” parameter. Used only if <code>method="LM"</code> , otherwise ignored.
<code>hessian</code>	Logical scalar. Should the hessian matrix be returned? This argument is relevant only if <code>method="bf"</code> (in which case it is passed along to <code>hmmNumOpt()</code> ) and is ignored otherwise. This argument should be set to TRUE only if you <i>really</i> want the hessian matrix. Setting it to TRUE causes a substantial delay between the time when <code>hmm()</code> finishes its iterations and when it actually returns a value.
<code>...</code>	Additional arguments passed to <code>hmmNumOpt()</code> . There is one noteworthy argument <code>useAnalGrad</code> which is used “directly” by <code>hmmNumOpt()</code> . This argument is a logical scalar and if it is TRUE then calls to <code>nlm()</code> or <code>optim()</code> are structured so that an analytic calculation of the gradient vector (implemented by the internal function <code>get.gl()</code> ) is applied. If it is FALSE then finite difference methods are used to calculate the gradient vector. If this argument is not specified it defaults to FALSE. Note that the name of this argument <b>cannot be abbreviated</b> .  Other “additional arguments” may be supplied for the control of <code>nlm()</code> and are passed on appropriately to <code>nlm()</code> . These are used only if <code>method="bf"</code> and if <code>optimiser="nlm"</code> . These “...” arguments might typically include <code>gradtol</code> , <code>stepmax</code> and <code>steptol</code> . They should <b>NOT</b> include <code>print.level</code> or <code>iterlim</code> . The former argument is automatically passed to <code>nlm()</code> as 0 if <code>verbose</code> is FALSE and as 2 if <code>verbose</code> is TRUE. The latter argument is automatically passed to <code>nlm()</code> with the value of <code>itmax</code> .

## Details

In the univariate case the emission probabilities are, in the case in which `newstyle` is TRUE (the default), specified by means of a data frame `Rho`. The first column of `Rho`, named `y`, is a factor consisting of the possible values of the emissions, repeated  $K$  times (where  $K$  is the number of states). The second column, named `states`, is a factor consisting of integer values  $1, 2, \dots, K$ . Each of these values is repeated  $m$  times where  $m$  is the length of `yval`. Further columns of `Rho` are

numeric and consist of coefficients of the linear predictor of the probabilities of the various values of  $y$ . If  $X$  is NULL then Rho has only one further column named Intercept.

If  $X$  is not NULL then the Intercept column is present only if addIntercept is TRUE. There as many (other, in addition to the possible Intercept column) numeric columns as there are columns in  $X$  or in the matrices in the list  $X$ . The names of these columns are taken to be the column names of  $X$  or the *first* entry of  $X$  if such column names are present. Otherwise the names default to V1, V2 ...

The probabilities of the emissions taking on their various possible values are given by

$$\Pr(Y = y_i | \mathbf{x}, \text{state} = S) = \ell_i / \sum_{j=1}^m \ell_j$$

where  $\ell_j$  is the  $j$ th entry of  $\beta^\top \mathbf{x}$  and where in turn  $\mathbf{x}$  is the vector of predictors and  $\beta$  is the coefficient vector in the linear predictor that corresponds to  $y_i$  and the hidden state  $S$ . For identifiability the vectors  $\beta$  corresponding to the first value of  $Y$  (the first level of Rho\$y) are set equal to the zero vector for all values of the state  $S$ .

If newstyle is FALSE Rho is specified by a matrix  $Rho = [\rho_{ij}]$  where  $\rho_{ij} = \Pr(Y = y_i | S = j)$ . If newstyle is FALSE and if  $X$  is not NULL then an error is thrown.

In the independent bivariate case the emission probabilities are specified by a list of two such matrices. In this setting  $P(Y_1, Y_2) = (y_{i1}, y_{i2}) | S = j) = \rho_{i1,j}^{(1)} \rho_{i2,j}^{(2)}$  where  $R^{(k)} = [\rho_{ij}^{(k)}]$  ( $k = 1, 2$ ) are the two emission probability matrices.

In the independent bivariate case the emission probabilities are specified by a list of two such matrices. In this setting  $P(Y_1, Y_2) = (y_{i1}, y_{i2}) | S = j) = \rho_{i1,j}^{(1)} \rho_{i2,j}^{(2)}$  where  $R^{(k)} = [\rho_{ij}^{(k)}]$  ( $k = 1, 2$ ) are the two emission probability matrices. In the dependent bivariate case the emission probabilities are specified by a three dimensional array. In this setting  $P((Y_1, Y_2) = (y_{i1}, y_{i2}) | S = j) = \rho_{i1,i2,j}^{(k)}$  where  $R_k = [\rho_{ij}^{(k)}]$  is the emission probability array.

The hard work of calculating the recursive probabilities used to fit the model is done by a Fortran subroutine "recurse" (actually coded in Ratfor) which is dynamically loaded. In the univariate case, when  $X$  is provided, the estimation of the "linear predictor" vectors  $\beta$  is handled by the function multinom() from the nnet package. Note that this is a "Recommended" package and is thereby automatically available (i.e. does not have to be installed).

## Value

A list with components:

Rho	The fitted value of the matrix, data frame, list of two matrices, or array Rho specifying the distributions of the observations (the "emission" probabilities).
tpm	The fitted value of the transition probability matrix tpm.
ispd	The fitted initial state probability distribution, or a matrix of (trivial or "deterministic") initial state probability distributions, one (column) of ispd for each observation sequence. If stationary is TRUE then ispd is assumed to be the (unique) stationary distribution for the chain, and thereby determined by the transition probability matrix tpm. If stationary is FALSE and cis is TRUE then ispd is estimated as the

mean of the vectors of conditional probabilities of the states, given the observation sequences, at time  $t=1$ .

If `cis` is FALSE then `ispd` is a matrix whose columns are trivial probability vectors, as described above.

<code>log.like</code>	The final (maximal, we hope!) value of the log likelihood, as determined by the maximisation procedure.
<code>grad</code>	The gradient of the log likelihood. Present only if the method is "LM" or "bf" and in the latter case then only if the optimiser is <code>nlm()</code> .
<code>hessian</code>	The hessian of the log likelihood. Present only if the method is "LM" or "bf".
<code>stopCrit</code>	A vector of the (final) values of the stopping criteria, with names "PCLL", "L2", "Linf" unless the method is "LM" or "SD" in which case this vector has a fourth entry named "ABSGRD".
<code>par0</code>	The starting values used by the algorithms. Either the argument <code>par0</code> , or a similar object with either or both components ( <code>tpm</code> and <code>Rho</code> ) being created by <code>rand.start()</code> .
<code>npar</code>	The number of parameters in the fitted model. Equal to <code>nispar + ntpmpar + nrhopar</code> where (1) <code>nispar</code> is 0 if <code>stationary</code> is TRUE and is $K-1$ otherwise; (2) <code>ntpmpar</code> is $K*(K-1)$ (3) <code>nrhopar</code> is <ul style="list-style-type: none"> <li>• <math>K*(nrow(Rho)-1)</math> for univariate models with <code>newstyle</code> equal to FALSE</li> <li>• <math>(nrow(Rho) - K)*(ncol(Rho)-2)</math> for univariate models with <code>newstyle</code> equal to TRUE</li> <li>• <math>K*(sum(sapply(Rho, nrow))-K)</math> for bivariate independent models</li> <li>• <math>prod(dim(Rho))-K</math> for bivariate dependent models.</li> </ul>
<code>bicm</code>	Numeric scalar. The number by which <code>npar</code> is multiplied to form the BIC criterion. It is essentially the log of the number of observations. See the code of <code>hmm()</code> for details.
<code>converged</code>	A logical scalar indicating whether the algorithm converged. If the EM, LM or steepest descent algorithm was used it simply indicates whether the stopping criterion was met before the maximum number ( <code>itmax</code> ) of steps was exceeded. If <code>method="bf"</code> then <code>converged</code> is based on the <code>code</code> component of the object returned by the optimiser when <code>nlm()</code> was used, or on the <code>convergence</code> component when <code>optim()</code> was used. In these cases <code>converged</code> has an <i>attribute</i> ( <code>code</code> or <code>convergence</code> respectively) giving the (integer) value of the relevant component.  Note that in the <code>nlm()</code> case a value of <code>code</code> equal to 2 indicates "probable" convergence, and a value of 3 indicates "possible" convergence. However in this context <code>converged</code> is set equal to TRUE <i>only</i> if <code>code</code> is 1.
<code>nstep</code>	The number of steps performed by the algorithm if the method was "EM", "LM" or "SD". The value of <code>nstep</code> is set equal to the <code>iterations</code> component of the value returned by <code>nlm()</code> if <code>method="bf"</code> .
<code>prior.emsteps</code>	The number of EM steps that were taken before the method was switched from "EM" to "bf" or to "LM". Present only in values returned under the "bf" or "LM" methods after a switch from "EM" and is equal to 0 if either of these methods was specified in the initial call (rather than arising as the result of a switch).

ylengths	Integer vector of the lengths of the observation sequences (number of rows if the observations are in the form of one or two column matrices).
nafrac	A real number between 0 and 1 or a pair (two dimensional vector) of such numbers. Each number is the the fraction of missing values if the corresponding components of the observations.
y	An object of class "tidyList". It is a tidied up version of the observations; i.e. the observations y after the application of the undocumented function tidyList(). Present only if keep.y is TRUE.
X	An object of class "tidyList". It is tidied up version of the predictor matrix or list of predictor matrices; i.e. the argument X after the application of tidyList() (with argument rp set to "predictor". Present only if X is supplied, is an appropriate argument, and if keep.X is TRUE.
parity	Character string; "univar" if the data were univariate, "bivar" if they were bivariate.
numeric	Logical scalar; TRUE if the (original) data were numeric, FALSE otherwise.
AIC	The value of $AIC = -2 \cdot \log .like + 2 \cdot npar$ for the fitted model.
BIC	The value of $BIC = -2 \cdot \log .like + \log(nobs) \cdot npar$ for the fitted model. In the forgoing nobs is the number of observations. This is the number of <i>non-missing</i> values in unlist(y) in the univariate setting and one half of this number in the bivariate setting.
args	A list of argument values supplied. This component is returned in the interest of making results reproducible. It is also needed to facilitate the updating of a model via the update method for the class hmm.discnp, <a href="#">update.hmm.discnp()</a> . It has components: <ul style="list-style-type: none"> <li>• newstyle</li> <li>• method</li> <li>• optimiser</li> <li>• optimMethod</li> <li>• stationary</li> <li>• mixture</li> <li>• cis</li> <li>• tolerance</li> <li>• itmax</li> <li>• crit</li> <li>• addIntercept</li> </ul>

### Warnings

The ordering of the (hidden) states can be arbitrary. What the estimation procedure decides to call "state 1" may not be what *you* think of as being state number 1. The ordering of the states will be affected by the starting values used.

Some experiences with using the "ABSGRD" stopping criterion indicate that it may be problematic in the context of discrete non-parametric distributions. For example a value of 1854.955 was returned after 200 LM steps in one (non-convergent, of course!) attempt at fitting a model. The stopping criterion "PCLL" in this example was took the "reasonable" value of 0.03193748 when iterations ceased.

## Notes

The package *used* to require the argument *y* to be a *matrix* in the case of multiple observed sequences. If the series were of unequal length the user was expected to pad them out with NAs to equalize the lengths.

The old matrix format for multiple observation sequences was permitted for a while (and the matrix was internally changed into a list) but this is no longer allowed. If *y* is indeed given as a matrix then this corresponds to a single observation sequence and it must have one (univariate setting) or two (bivariate setting) columns which constitute the observations of the respective variates.

If  $K=1$  then *tpm*, *ispd*, *converged*, and *nstep* are all set equal to NA in the list returned by this function.

The estimate of *ispd* in the non-stationary setting is inevitably very poor, unless the number of sequences of observations (the length of the list *y*) is very large. We have in effect “less than one” relevant observation for each such sequence.

The returned values of *tpm* and *Rho* (or the entries of *Rho* when *Rho* is a list) have dimension names. These are the same as the dimension names of the starting values of these objects (as provided in *par0*) if these exist. Otherwise they are formed from the sorted unique values of the observations in *y* or are taken to be the appropriate sequences of integers. Likewise the returned value of *ispd* is a named vector, the names being the same as the row (and column) names of *tpm* and the corresponding dimension names of *Rho* or of its entries.

If *method* is equal to “EM” it *may* get switched to “bf” at some EM step if there is a decrease in the log likelihood. This is “theoretically impossible” but can occur in practice due to an intricacy in the way that the EM algorithm treats *ispd* when *stationary* is TRUE. It turns out to be effectively impossible to maximise the expected log likelihood unless the term in that quantity corresponding to *ispd* is ignored (whence it *is* ignored). Ignoring this term is “asymptotically negligible” but can have the unfortunate effect of occasionally leading to a decrease in the log likelihood.

If such a decrease is detected, then the algorithm switches over to using the “bf” method of maximisation (which is not beset by this difficulty). The current estimates of *ispd*, *tpm* and *Rho* are used as starting values for the “bf” method.

It seems to me that it *should* be the case that such switching can occur only if *stationary* is TRUE. However I have encountered instances in which the switch occurs when *stationary* is FALSE. I have yet to figure out/track down what is going on here.

## Author(s)

Rolf Turner <r.turner@auckland.ac.nz>

## References

- Rabiner, L. R., "A tutorial on hidden Markov models and selected applications in speech recognition," Proc. IEEE vol. 77, pp. 257 – 286, 1989.
- Zucchini, W. and Guttorp, P., "A hidden Markov model for space-time precipitation," Water Resources Research vol. 27, pp. 1917-1923, 1991.
- MacDonald, I. L., and Zucchini, W., "Hidden Markov and Other Models for Discrete-valued Time Series", Chapman & Hall, London, 1997.

Liu, Limin, "Hidden Markov Models for Precipitation in a Region of Atlantic Canada", Master's Report, University of New Brunswick, 1997.

### See Also

[rhmm\(\)](#), [mps\(\)](#), [viterbi\(\)](#)

### Examples

```
# TO DO: Create one or more bivariate examples.
#
# The value of itmax in the following examples is so much
# too small as to be risible. This is just to speed up the
# R CMD check process.
# 1.
Yval <- LETTERS[1:10]
Tpm <- matrix(c(0.75,0.25,0.25,0.75),ncol=2,byrow=TRUE)
Rho <- cbind(c(rep(1,5),rep(0,5)),c(rep(0,5),rep(1,5)))/5
rownames(Rho) <- Yval
set.seed(42)
xxx <- rhmm(ylengths=rep(1000,5),nsim=1,tpm=Tpm,Rho=Rho,yval=Yval,drop=TRUE)
fit1 <- hmm(xxx,par0=list(tpm=Tpm,Rho=Rho),newstyle=FALSE,itmax=10)
print(fit1$Rho)
newRho <- cnvrtRho(Rho)
fit2 <- hmm(xxx,par0=list(tpm=Tpm,Rho=newRho),itmax=10)
print(round(cnvrtRho(fit2$Rho),6))
# Agrees with fit1$Rho to 5 decimal places; the difference
# is due to the fitting procedure used with the "new style"
# Rho, which calls upon multinom() from the nnet package.

# 2.
# See the help for rhmm() for how to generate y.num.
## Not run:
fit.num <- hmm(y.num,K=2,verb=TRUE,itmax=10)
fit.num.mix <- hmm(y.num,K=2,verb=TRUE,mixture=TRUE,itmax=10)
print(fit.num[c("tpm","Rho")])

## End(Not run)
# Note that states 1 and 2 get swapped.

# 3.
xxx <- with(SydColDisc,split(y,f=list(locn,depth)))
fitSydCol <- hmm(xxx,K=2,itmax=10) # Two states: above and below the thermocline.

# 4.
fitLesCount <- hmm(lesionCount,K=2,itmax=10) # Two states: relapse and remission.
```

## Description

Five data sets obtained from the “HYDAT” database, Environment and Climate Change Canada’s database of historical hydrometric data. The data were obtained using the `tidyhydro` package. The data have been trimmed so that there are no gaps in the observation dates and are presented in “raw” form and in discretised form as deciles of the residuals (difference between raw values and the daily mean over years).

## Usage

```
data("linLandFlows")
data("ftLiardFlows")
data("portMannFlows")
data("portMannSedLoads")
data("portMannSedCon")
```

## Format

Data frames with observations on the following 3 variables.

`Date` Dates on which observations were made.

`Value` Numeric vector of observation values.

`mean` The mean over years of `Value`.

`resid` The difference `Value - mean`.

`deciles` A factor with levels `d1, ..., d10`, which are the deciles of the variable `resid`

## Details

The variable `mean` was calculated as follows:

```
yday <- as.POSIXlt(X$Date)$yday
mn <- tapply(X$Value, yday, mean, na.rm=TRUE)
mean <- mn[as.character(yday)]
```

where `X` is the data set being processed.

The data set `linLandFlows` originally consisted of 2008 observations; there were 1980 observations after “trimming”. The data set `ftLiardFlows` originally consisted of 22364 observations; there were 11932 observations after “trimming”. The data set `portMannFlows` originally consisted of 6455 observations; there were 3653 observations after “trimming”. The data set `portMannSedLoads` consists of 2771 observations; no observations were trimmed. The data set `portMannSedCon` consists of 4597 observations; no observations were trimmed.

The units of the “Flows” variables are cubic metres per second ( $m^3/s$ ); the units of “`portMannSedLoads`” are tonnes; the units of “`portMannSedCon`” are milligrams per litre (mg/l).

The “`linLandFlows`” data were obtained at the Lindberg Landing hydrometric station on the Liard River in the Northwest Territories of Canada. The “`ftLiardFlows`” data were obtained at the Fort Liard hydrometric station on the Liard River in the Northwest Territories of Canada. The “`portMann`” data were obtained at the hydrometric station located at the Port Mann pumping station on the Fraser River in the Province of British Columbia in Canada.

**Source**

Environment and Climate Change Canada’s database “HYDAT”, a database of historical hydrometric data. The data were obtained via the tidyhydat package, which is available from “CRAN”, <https://cran.r-project.org>

**Examples**

```
fit <- hmm(linLandFlows$deciles,K=4,itmax=10)
```

---

lesionCount

*Multiple sclerosis lesion counts for three patients.*

---

**Description**

Lesion counts for three multiple sclerosis patients. The counts were obtained by magnetic resonance imaging, and were observed at monthly intervals.

**Usage**

```
lesionCount
```

**Format**

A list with three components each component being the sequence of counts for a given patient and consisting of a vector with non-negative integer entries.

**Modelling**

The hidden Markov models applied to these data by Albert et al. and by MacKay and Petkau were much more complex and elaborate than those fitted in the examples in this package. See the references for details.

**Source**

The data were originally studied by Albert et al., (1994). They were also analyzed by Altman and Petkau (2005). The data were kindly provided by Prof. Altman.

**References**

Albert, P. S., McFarland, H. F., Smith, M. E., and Frank, J. A. Time series for modelling counts from a relapsing-remitting disease: application to modelling disease activity in multiple sclerosis. *Statistics in Medicine* **13** (1994) 453–466.

Altman, Rachel MacKay, and Petkau, A. John. Application of hidden Markov models to multiple sclerosis lesion count data. *Statistics in Medicine* **24** (2005) 2335–2344.

logLikHmm

*Log likelihood of a hidden Markov model***Description**

Calculate the log likelihood of a hidden Markov model with discrete non-parametric observation distributions.

**Usage**

```
logLikHmm(y, model=NULL, tpm=NULL, ispd=NULL, Rho=NULL, X=NULL,
          addIntercept=NULL, warn=TRUE)
```

**Arguments**

y	A vector, or list of vectors, or a two column matrix or a a list of such matrices, whose entries consist of observations from a hidden Markov model with discrete non-parametric observation distributions.
model	An object specifying a hidden Markov model, usually returned by <code>hmm()</code> .
tpm	The transition probability matrix of the Markov chain. Ignored (and extracted from <code>model</code> ) if <code>model</code> is non-NULL.
ispd	The vector of probabilities specifying the initial state probability distribution, or a matrix each of whose columns is a trivial (“delta function”) vector specifying the “most probable” initial state for each observation sequence. If <code>ispd</code> is missing then <code>ispd</code> is calculated as the stationary distribution determined by <code>tpm</code> . Ignored (and extracted from <code>model</code> ) if <code>model</code> is non-NULL.
Rho	An object specifying the “emission” probabilities of the observations. (See the <b>Details</b> in the help for <code>hmm()</code> .) Ignored (and extracted from <code>model</code> ) if <code>model</code> is non-NULL.
X	An optional <i>numeric</i> matrix, or a list of such matrices, of <i>predictors</i> . The use of such predictors is (currently, at least) applicable only in the univariate emissions setting, and then only if <code>Rho</code> is supplied in the “newstyle” format. If <code>X</code> is a list it must be of the same length as <code>y</code> and all entries of this list must have the same number of columns. The number of rows of each entry must be equal to the length of the corresponding entry of <code>y</code> . If <code>X</code> is a matrix then <code>y</code> should be a vector or one-column matrix (or a list with a single entry equal to such).
addIntercept	Logical scalar. See the documentation of <code>hmm()</code> . If this argument is not specified, and if <code>model</code> is NULL then an error is thrown.
warn	Logical scalar; should a warning be issued if <code>Rho</code> hasn’t got relevant dimension names? (Note that if this is so, then the corresponding dimension names are formed from the sorted unique values of <code>y</code> or of the appropriate column(s) of <code>y</code> . And if <i>this</i> is so, then the user should be sure that the ordering of the entries of <code>Rho</code> corresponds properly to the the sorted unique values of <code>y</code> .) This argument is passed to the utility function <code>check.yval()</code> which actually issues the warning if <code>warn=TRUE</code> .

**Details**

If `y` is not provided the function simply returns the `log.like` component of `model` (which could be `NULL` if `model` was not produced by `hmm()`).

The observation values (the entries of the vector or matrix `y`, or of the vectors or matrices which constitute the entries of `y` if `y` is a list) must be consistent with the appropriate dimension names of `Rho` or of its entries when `Rho` is a list. More specifically, if `Rho` has dimension names (or its entries have dimension names) then the observation values must all be found as entries of the appropriate dimension name vector. If a vector of dimension names is `NULL` then the corresponding dimension must be equal to the number of unique observations of the appropriate variate. integers between 1 and `nrow(Rho)`.

**Value**

The loglikelihood of `y` given the parameter values specified in `par`.

**Author(s)**

Rolf Turner <[r.turner@auckland.ac.nz](mailto:r.turner@auckland.ac.nz)>

**References**

See [hmm\(\)](#) for references.

**See Also**

[hmm\(\)](#), [pr\(\)](#), [sp\(\)](#)

**Examples**

```
# TO DO: One or more bivariate examples.
P <- matrix(c(0.7,0.3,0.1,0.9),2,2,byrow=TRUE)
R <- matrix(c(0.5,0,0.1,0.1,0.3,
             0.1,0.1,0,0.3,0.5),5,2)

set.seed(42)
l11 <- sample(250:350,20,TRUE)
set.seed(909)
y.num <- rhmm(ylengths=l11,nsim=1,tpm=P,Rho=R,drop=TRUE)
set.seed(909)
y.let <- rhmm(ylengths=l11,nsim=1,tpm=P,Rho=R,yval=letters[1:5],drop=TRUE)
row.names(R) <- 1:5
l11 <- logLikHmm(y.num,tpm=P,Rho=R)
row.names(R) <- letters[1:5]
l12 <- logLikHmm(y.let,tpm=P,Rho=R)
l13 <- logLikHmm(y.let,tpm=P,Rho=R,ispd=c(0.5,0.5))
fit <- hmm(y.num,K=2,itmax=10)
l14 <- logLikHmm(y.num,fit) # Use the fitted rather than the "true" parameters.
```

---

misstify                      *Insert missing values.*

---

### Description

Insert missing values into data simulated by rhmm.

### Usage

```
misstify(y, nafrac, fep = NULL)
```

### Arguments

- |        |  |
|--------|--|
| y      | A data set (list of matrices with one or two columns) or a list of such data sets (objects of class "multipleHmmDataSets" such as might be generated by <a href="#">rhmm()</a> )   |
| nafrac | <p>A numeric vector, some entries of which could be ignored. (See below.) Those which do not get ignored must be probabilities <i>strictly</i> less than 1. (Having <i>everything</i> missing makes no sense!)</p> <p>The vector nafrac will be replicated to have an "appropriate" length. If y is of class "multipleHmmDataSets" then this length is length(y) if the data are univariate and is 2*length(y) if the data are bivariate. In the former case the entries of the replicated vector form the fraction of missing values in the corresponding data set. In the latter case the odd numbered entries form the fraction of missing values for the first variable and the even numbered entries the fraction for the second variable. If y is not of class "multipleHmmDataSets" then this length is either 1 (univariate case) or 2 (bivariate case).</p> <p>Note that replication discards entries that are not needed to make up the required length, and such entries are thereby ignored. E.g. <code>rep(c(0.2, 0.7, 1.6), length=2)</code> yields <code>[1] 0.2 0.7</code>, i.e. the entry 1.6 is ignored.</p> <p>The fraction(s) of missing values in a given data set may be determined by <a href="#">nafracCalc()</a>.</p> |
| fep    | <p>"First entry present". A list with one or two entries, the first being a logical scalar (which might be named "present". If there is a second entry it should be a scalar probability (which might be named "p2"). In an application of interest, observation sequences always begin at an observed event, i.e. at a time point at which the "emission" has at least one non-missing value. If <code>fep[[1]]</code> is TRUE the NAs will be inserted in such a way that the resulting data have this characteristic. If fep is left NULL then its first (possibly only) entry is set to TRUE.</p> <p>For <i>bivariate</i> data, <code>fep[[2]]</code> specifies the probability that <i>both</i> values of the initial pair of observations are non-missing. In this case one of the entries of the initial pair is chosen to be "potentially" missing, with probabilities <code>nafrac/sum(nafrac)</code>. This entry is left non-missing with probability <code>fep[[2]]</code>. (The other entry is always left non-missing.)</p>   |

If the data are univariate or if `fep[[1]]` is FALSE, then `fep[[2]]` is ignored. If the data are bivariate and `fep[[2]]` is not specified, it defaults to the (estimated) conditional probability that both entries of the initial pair of observations are present given that at least one is present, under the assumption of independence of these events. I.e. it is set equal to  $\text{prod}(1-\text{nafrac})/(1-\text{prod}(1-\text{nafrac}))$ .

## Value

An object with a structure similar to that of `y`, containing the same data as `y` but with some of these data having been replaced by missing values (NA). In particular, if `y` is of class "multipleHmDataSets" then so is the returned value.

Note that `rhmm()` calls upon `misstify()` to effect the replacement of a certain fraction of the simulated observations by missing values. If `rhmm()` is applied to a fitted model, then by default, this "certain fraction" is determined, using `nafracCalc()`, from the data set to which the model was fitted.

## Author(s)

Rolf Turner <r.turner@auckland.ac.nz>

## See Also

[rhmm\(\)](#) [nafracCalc\(\)](#)

## Examples

```
P <- matrix(c(0.7,0.3,0.1,0.9),2,2,byrow=TRUE)
R <- matrix(c(0.5,0,0.1,0.1,0.3,
             0.1,0.1,0,0.3,0.5),5,2)
set.seed(42)
l11 <- sample(250:350,20,TRUE)
y1 <- rhmm(ylengths=l11,nsim=1,tpm=P,Rho=R)
y1m <- misstify(y1,nafrac=0.5,fep=list(TRUE))
y2 <- rhmm(ylengths=l11,nsim=5,tpm=P,Rho=R)
set.seed(127)
y2m <- misstify(y2,nafrac=0.5,fep=list(TRUE))
nafracCalc(y2m) # A list all of whose entries are close to 0.5.
set.seed(127)
y2ma <- lapply(y2,misstify,nafrac=0.5,fep=list(TRUE))
## Not run:
nafracCalc(y2ma) # Throws an error.

## End(Not run)
sapply(y2ma,nafracCalc) # Effectively the same as nafracCalc(y2m).
```

---

mps *Most probable states.*

---

### Description

Calculates the most probable hidden state underlying each observation.

### Usage

```
mps(y, model = NULL, tpm, Rho, ispd=NULL, warn=TRUE)
```

### Arguments

y	The observations for which the underlying most probable hidden states are required. May be a sequence of observations, or a list each component of which constitutes a (replicate) sequence of observations. If y is missing it is set equal to the y component of model, given that that object and that component exist. Otherwise an error is given.
model	An object describing a fitted hidden Markov model, as returned by <code>hmm()</code> . In order to make any kind of sense, model should bear some reasonable relationship to y.
tpm	The transition probability matrix for a hidden Markov model; ignored if model is non-null. Should bear some reasonable relationship to y.
Rho	An object specifying the probability distributions of the observations (“emission” probabilities) for a hidden Markov model. See <code>hmm()</code> . Ignored if model is non-null. Should bear some reasonable relationship to y.
ispd	A vector specifying the initial state probability distribution for a hidden Markov model, or a matrix each of whose columns are trivial (“delta function”) vectors specifying the “most probable” initial state for each observation sequence. This argument is ignored if model is non-null. It should bear some reasonable relationship to y. If both ispd and model are NULL then ispd is taken to be the stationary distribution of the chain, calculated from tpm.
warn	Logical scalar; should a warning be issued if Rho hasn’t got relevant dimension names? (Note that if this is so, then the corresponding dimension names are formed from the sorted unique values of y or of the appropriate column(s) of y. And if <i>this</i> is so, then the user should be sure that the ordering of the entries of Rho corresponds properly to the the sorted unique values of y.) This argument is passed to the utility function <code>check.yval()</code> which actually issues the warning if warn=TRUE.

### Details

For each  $t$  the maximum value of  $\gamma_t(i)$ , i.e. of the (estimated) probability that the state at time  $t$  is equal to  $i$ , is calculated, and the value of the state with the corresponding index is returned.

**Value**

If  $y$  is a single observation sequence, then the value is a vector of corresponding most probable states.

If  $y$  is a list of replicate sequences, then the value is a list, the  $j$ -th entry of which constitutes the vector of most probable states underlying the  $j$ -th replicate sequence.

**Warning**

The *sequence of most probable states* as calculated by this function will not in general be the *most probable sequence of states*. It may not even be a *possible* sequence of states. This function looks at the state probabilities separately for each time  $t$ , and not at the states in their sequential context.

To obtain the most probable sequence of states use `viterbi()`.

**Author(s)**

Rolf Turner <r.turner@auckland.ac.nz>

**References**

Rabiner, L. R., "A tutorial on hidden Markov models and selected applications in speech recognition," Proc. IEEE vol. 77, pp. 257 – 286, 1989.

**See Also**

`hmm()`, `rhmm()`, `viterbi()`

**Examples**

```
# See the help for rhmm() for how to generate y.num.
## Not run:
fit.num <- hmm(y.num,K=2,verb=TRUE)
s.1 <- mps(y.num,fit.num)
s.2 <- mps(y.num,tpm=P,ispd=c(0.25,0.75),Rho=R) # P and R as in the help
                                                # for rhmm().
# The order of the states has gotten swapped; (3 - s.1[,1]) is much
# more similar to s.2[,1] than is s.1[,1].

## End(Not run)
```

---

nafracCalc

*Calculate fractions of missing values.*

---

**Description**

Calculate the fraction (univariate case) or fractions (bivariate case) of missing values in the data or in each component of the data.

**Usage**

```
nafracCalc(y, drop=TRUE)
```

**Arguments**

**y** A vector or a one or two column matrix of discrete data or a list of such vectors or matrices, or a *list* of such objects (an object of class "multipleHmmDataSets" such as might be produced by `rhmm()`).

**drop** Logical scalar. If *y* is of class "multipleHmmDataSets" but actually consists of a single data set, and if `drop` is TRUE, then the returned value is not a list but rather the single component that such a list "would have had" were `drop` equal to FALSE. This argument is ignored if *y* is not of class "multipleHmmDataSets" or has length greater than 1.

**Value**

If *y* is *not* of class "multipleHmmDataSets", then the returned value is a scalar (between 0 and 1) if the data are univariate or a pair (2-vector) of such scalars if the data are bivariate. The values are equal to the ratios of the total count of missing values in the appropriate column to the total number of observations.

If *y* is of class "multipleHmmDataSets", and if *y* has length greater than 1 or `drop` is FALSE, then the returned value is a *list* of such scalars or 2-vectors, each corresponding to one of the data sets constituting *y*. If *y* has length equal to 1 and `drop` is TRUE, then the returned value is the same as if *y* were not of class "multipleHmmDataSets".

**Author(s)**

Rolf Turner <r.turner@auckland.ac.nz>

**See Also**

`rhmm()` `misstify()`

**Examples**

```
xxx <- with(SydColDisc, split(y, f=list(locn, depth)))
nafracCalc(xxx) # 0.7185199
```

---

pr

*Probability of state sequences.*

---

**Description**

Calculates the conditional probability of one or more state sequences, given the corresponding observations sequences (and the model parameters).

**Usage**

```
pr(s, y, model=NULL, tpm, Rho, ispd=NULL, warn=TRUE)
```

**Arguments**

s	A sequence of states of the underlying Markov chain, or a list of such sequences.
y	A one or two column matrix of observations from a hidden Markov model, corresponding to the state sequence s, or a list of such matrices corresponding to the state sequences in the list s. If y consists of a single matrix, it is assumed to correspond to each of the state sequences in s in turn. Otherwise the length of the list y must be the same as the length of the list s. (If not, then an error is given). If y is missing, it is extracted from model provided that model and its y component are not NULL. Otherwise an error is given.
model	An object of class <code>hmm.discnp</code> as returned by <code>hmm()</code> .
tpm	The transition probability matrix of the chain. Ignored (and extracted from model instead) if model is not NULL.
Rho	An object specifying the “emission” probabilities of observations, given the underlying state. See <code>hmm()</code> . Ignored (and extracted from model instead) if model is not NULL.
ispd	The vector specifying the initial state probability distribution of the Markov chain. Ignored (and extracted from model instead) if model is not NULL. If both ispd and model are NULL then ispd is taken to be the stationary distribution of the chain, calculated from tpm.
warn	Logical scalar; should a warning be issued if Rho hasn’t got relevant dimension names? (Note that if this is so, then the corresponding dimension names are formed from the sorted unique values of y or of the appropriate column(s) of y. And if <i>this</i> is so, then the user should be sure that the ordering of the entries of Rho corresponds properly to the the sorted unique values of y.) This argument is passed to the utility function <code>check.yval()</code> which actually issues the warning if <code>warn=TRUE</code> .

**Value**

The probability of s given y, or a vector of such probabilities if s and y are lists.

**Warning**

The conditional probabilities will be tiny if the sequences involved are of any substantial length. Underflow may be a problem. The implementation of the calculations is not sophisticated.

**Author(s)**

Rolf Turner <r.turner@auckland.ac.nz>

**See Also**

`hmm()`, `mps()`, `viterbi()`, `sp()`, `fitted.hmm.discnp()`

## Examples

```
# See the help for rhmm() for how to generate y.num.
## Not run:
fit.num <- hmm(y.num,K=2,verb=TRUE)
# Using fitted parameters.
s.vit.1 <- viterbi(y.num,fit.num)
pr.vit.1 <- pr(s.vit.1,model=fit.num)
# Using true parameters from which y.num was generated.
s.vit.2 <- viterbi(y.num,tpm=P,Rho=R)
pr.vit.2 <- pr(s.vit.2,y.num,tpm=P,Rho=R)

## End(Not run)
```

---

predict.hmm.discnp      *Predicted values of a discrete non-parametric hidden Markov model.*

---

## Description

Calculates predicted values given a specification of a discrete non-parametric hidden Markov model. The specification may be provided in the form of a `hmm.discnp` object as returned by `hmm()` or in the form of “components” of such a model: the data `y`, the transition probability matrix `tpm`, the emission probabilities `Rho`, etc. If the data are numeric then these predicted values are the conditional expectations of the observations, given the entire observation sequence (and the — possibly estimated — parameters of the model). If the data are categorical (whence “expectations” make no sense) the “predicted values” are taken to be the probabilities of each of the possible values of the observations, at each time point.

## Usage

```
## S3 method for class 'hmm.discnp'
predict(object, y = NULL, tpm=NULL, Rho=NULL,
        ispd=NULL, X=NULL,addIntercept=NULL,
        warn=TRUE, drop=TRUE, ...)
```

## Arguments

<code>object</code>	If not <code>NULL</code> , an object of class <code>hmm.discnp</code> as returned by <code>hmm()</code> .
<code>y</code>	A data structure from which the fitted model object <i>could</i> have been fitted. If <code>y</code> is <code>NULL</code> , an attempt is made to extract <code>y</code> from <code>model</code> .
<code>tpm,Rho,ispd,X,addIntercept,warn</code>	See the help for <code>sp()</code> .
<code>drop</code>	Logical scalar. See the help for <code>fitted.hmm.discnp()</code> .
<code>...</code>	Not used.

## Details

This function is essentially the same as `fitted.hmm.discnp()`. The main difference is that it allows the calculation of fitted/predicted values for a data object `y` possibly different from that to which the model was fitted. Note that if both the argument `y` and `object[["y"]]` are present, the “argument” value takes precedence. This function also allows the model to be specified in terms of individual components rather than as a fitted model of class `"hmm.discnp"`. These components, (`tpm`, `Rho`, `ispd`, `X`, `addIntercept`) if supplied, *take precedence* over the corresponding components of `object`. The opposite applies with `sp()`. The function `fitted.hmm.discnp()` makes use *only* of the components of `object`.

## Value

See the help for `fitted.hmm.discnp()`.

## Author(s)

Rolf Turner <r.turner@auckland.ac.nz>

## See Also

`sp()` `link{fitted.hmm.discnp}()`

## Examples

```
P <- matrix(c(0.7,0.3,0.1,0.9),2,2,byrow=TRUE)
R <- matrix(c(0.5,0,0.1,0.1,0.3,
             0.1,0.1,0,0.3,0.5),5,2)
set.seed(42)
l11 <- sample(250:350,20,TRUE)
y1 <- rhmm(ylengths=l11,nsim=1,tpm=P,Rho=R,drop=TRUE)
fit <- hmm(y1,K=2,verb=TRUE,keep.y=TRUE,itmax=10)
fv <- fitted(fit)
set.seed(176)
l12 <- sample(250:350,20,TRUE)
y2 <- rhmm(ylengths=l12,nsim=1,tpm=P,Rho=R,drop=TRUE)
pv <- predict(fit,y=y2)
yval <- letters[1:5]
set.seed(171)
y3 <- rhmm(ylengths=l12,yval=yval,nsim=1,tpm=P,Rho=R,drop=TRUE)
fit3 <- hmm(y3,K=2,verb=TRUE,keep.y=TRUE,itmax=10)
pv3 <- predict(fit3) # Same as fitted(fit3).
```

## Description

Simulates one or more replicates of discrete data from a model such as is fitted by the function `hmm()`.

## Usage

```
rhmm(model, ..., nsim, verbose=FALSE)
## Default S3 method:
rhmm(model, ..., nsim=1, verbose=FALSE, ylengths,
      nafrac=NULL, fep=NULL, tpm, Rho, ispd=NULL, yval=NULL,
      drop=TRUE)
## S3 method for class 'hmm.discnp'
rhmm(model, ..., nsim=1, verbose=FALSE, inMiss=TRUE,
      fep=NULL, drop=TRUE)
```

## Arguments

<code>model</code>	An object of class <code>hmm.discnp</code> . This will have the form of a list specifying a hidden Markov model with discrete emissions and emission probabilities specified non-parametrically, i.e. by means of a table or tables. Usually this will be an object returned by <code>hmm()</code> . This argument is ignored by the default method.
<code>...</code>	Not used.
<code>nsim</code>	Integer scalar; the number of data sets to be simulated.
<code>verbose</code>	Logical scalar. If TRUE then the overall index of the simulated value that has been reached is printed out every 1000 iterations. Useful for reassurance when very “large” simulations are undertaken.
<code>ylengths</code>	Integer values vector specify the lengths (or number of rows in the bivariate setting) of the individual observation sequences constituting a data set.
<code>nafrac</code>	See <code>misstify()</code> for an explanation of this argument. If specified a fraction <code>nafrac[[j]]</code> of column <code>j</code> of the data will be randomly set equal to NA.
<code>fep</code>	“First entry present”. See <code>misstify()</code> for an explanation of this argument.
<code>tpm</code>	The transition probability matrix for the underlying hidden Markov chain(s). Note that the rows of <code>tpm</code> must sum to 1. Ignored if <code>ncol(Rho)==1</code> . Ignored by the <code>hmm.discnp</code> method and extracted from <code>model</code> .
<code>Rho</code>	An object specifying the probability distribution of the observations, given the state of the underlying hidden Markov chain. (I.e. the “emission” probabilities.) See <code>hmm()</code> . Note that <code>Rho</code> can be such that the number of states is 1, in which case the simulated data are i.i.d. from the single distribution specified by <code>Rho</code> . Ignored by the <code>hmm.discnp</code> method and extracted from <code>model</code> .
<code>ispd</code>	A vector specifying the initial state probability distribution of the chain. If this is not specified it is taken to be the stationary distribution of the chain, calculated from <code>tpm</code> . Ignored by the <code>hmm.discnp</code> method and extracted from <code>model</code> .
<code>yval</code>	Vector of possible values of the observations, or (in the bivariate setting) a list of two such vectors. If not supplied it is formed from the appropriate dimension names of <code>Rho</code> or of the entries of <code>Rho</code> if <code>Rho</code> is a list. Ignored by the <code>hmm.discnp</code> method.

drop	Logical scalar; if nsim is 1 and if drop is TRUE then the list to be returned by this function is replaced by its first and only entry. (See below.)
inMiss	Logical scalar; if TRUE then missing values will be randomly inserted into the data in the fraction nafrac determined from object.

### Value

If nsim>1 or drop is FALSE then the value returned is a list of length nsim. Each entry of this list is in turn a list of the same length as ylengths, each component of which is an independent vector or matrix of simulated observations. The length or number of rows of component i of this list is equal to ylengths[i]. The values of the observations are entries of yval or of its entries when yval is a list.

If nsim=1 and drop is TRUE then the (“outer”) list described above is replaced by its first and only entry

### Note

You may find it useful to avail yourself of the function [nafracCalc\(\)](#) to determine the fraction of missing values in a given existing (presumably “real”) data set.

### Author(s)

Rolf Turner <r.turner@auckland.ac.nz>

### See Also

[hmm\(\)](#) [nafracCalc\(\)](#) [misstify\(\)](#)

### Examples

```
# To do: one or more bivariate examples.
## Not run:
  y <- list(linLandFlows$deciles,ftLiardFlows$deciles)
  fit <- hmm(y,K=3)
  simX <- rhmm(fit)

## End(Not run)
```

---

scovmat

*Simulation based covariance matrix.*

---

### Description

Produces an estimate of the covariance matrix of the parameter estimates in a model fitted by `hmm.discnp`. Uses a method based on simulation (or “parametric bootstrapping”).

**Usage**

```
scovmat(object, seed = NULL, nsim=100, verbose = TRUE)
```

**Arguments**

object	An object of class <code>hmm.discnp</code> as returned by <code>hmm()</code> .
seed	Integer scalar serving as a seed for the random number generator. If left <code>NULL</code> the seed itself is chosen randomly from the set of integers between 1 and $10^5$ .
nsim	A positive integer. The number of simulations upon which the covariance matrix estimate will be based.
verbose	Logical scalar; if <code>TRUE</code> , iteration counts will be printed out during each of the simulation and model-fitting stages.

**Details**

This function is currently applicable only to models fitted to univariate data. The covariance matrix produced is for the “raw” parameters (entries of `tpm` with the last column dropped — since the rows sum to 1, and the entries of `Rho` with the last row dropped — since the columns sum to 1).

**Value**

A (positive definite) matrix which is an estimate of the covariance of the parameter estimates from the fitted model specified by `object`. It has row and column labels which indicate the parameters to which its entries pertain, in a reasonably perspicuous manner.

This matrix has an attribute `seed` (the random number generation seed that was used) so that the calculations can be reproduced.

**Author(s)**

Rolf Turner <r.turner@auckland.ac.nz>

**See Also**

[squantCI\(\)](#) [link{rhmm}\(\)](#) [link{hmm}\(\)](#)

**Examples**

```
## Not run:
y <- list(lindLandFlows$deciles, ftLiardFlows$deciles)
fit <- hmm(y, K=3)
ccc <- scovmat(fit, nsim=100)

## End(Not run)
```

---

sp *Calculate the conditional state probabilities.*

---

### Description

Returns the probabilities that the underlying hidden state is equal to each of the possible state values, at each time point, given the observation sequence.

### Usage

```
sp(y, model = NULL, tpm=NULL, Rho=NULL, ispd=NULL, X=NULL,
   addIntercept=NULL, warn=TRUE, drop=TRUE)
```

### Arguments

y	The observations on the basis of which the probabilities of the underlying hidden states are to be calculated. May be a (one or two column) matrix of observations, or a list each component of which is such a matrix. If y is missing it is set equal to the y component of model, given that that argument is non-NULL and that that component exist. Otherwise an error is given.
model	An object of class <code>hmm.discnp</code> as returned by <code>hmm()</code> .
tpm	The transition probability matrix for the underlying hidden Markov chain. Ignored if model is not NULL (in which case tpm is extracted from model).
Rho	An object specifying the distribution of the observations, given the underlying state. I.e. the “emission” probabilities. See <code>hmm()</code> . Ignored if model is not NULL (in which case Rho is extracted from model).
ispd	Vector specifying the initial state probability distribution of the underlying hidden Markov chain. Ignored if model is not NULL (in which case ispd is extracted from model). If both <code>model[["ispd"]]</code> and ispd are NULL then ispd is calculated to be the stationary distribution of the chain as determined by tpm.
X	An optional <i>numeric</i> matrix, or a list of such matrices, of <i>predictors</i> . Ignored if model is not NULL (in which case X is extracted from model). The use of such predictors is (currently, at least) applicable only in the univariate emissions setting, and then only if Rho is supplied in the “newstyle” format. If X is a list it must be of the same length as y and all entries of this list must have the same number of columns. The number of rows of each entry must be equal to the length of the corresponding entry of y. If X is a matrix then y should be a vector or one-column matrix (or a list with a single entry equal to such).
addIntercept	Logical scalar. See the documentation of <code>hmm()</code> . Ignored if model is not NULL (in which case addIntercept is extracted from model).
warn	Logical scalar; should a warning be issued if Rho hasn’t got relevant dimension names? (Note that if this is so, then the corresponding dimension names are formed from the sorted unique values of y or of the appropriate column(s) of y. And if <i>this</i> is so, then the user should be sure that the ordering of the entries of Rho corresponds properly to the the sorted unique values of y.) This argument is

passed to the utility function `check.yval()` which actually issues the warning if `warn=TRUE`.

**drop** Logical scalar. If `y` is a matrix, or a list of length 1, and if `drop` is `FALSE` then the returned value is a list whose sole entry is the matrix that would have been returned were `drop` equal to `TRUE`. The argument `drop` is ignored if `y` is a list of length greater than 1.

### Details

Note that in contrast to `predict.hmm.discnp()`, components in `model` take precedence over individually supplied components (`tpm`, `Rho`, `ispd`, `X` and `addIntercept`).

### Value

If `y` is a single matrix of observations or a list of length 1, and if `drop` is `TRUE` then the returned value is a matrix whose rows correspond to the states of the hidden Markov chain, and whose columns correspond to the observation times. Otherwise `probs` is a list of such matrices, one for each matrix of observations.

### Author(s)

Rolf Turner <r.turner@auckland.ac.nz>

### See Also

[hmm\(\)](#), [mps\(\)](#), [viterbi\(\)](#), [pr\(\)](#), [fitted.hmm.discnp\(\)](#)

### Examples

```
P <- matrix(c(0.7,0.3,0.1,0.9),2,2,byrow=TRUE)
R <- matrix(c(0.5,0,0.1,0.1,0.3,
             0.1,0.1,0,0.3,0.5),5,2)

set.seed(42)
y <- rhmm(ylengths=rep(300,20),nsim=1,tpm=P,Rho=R,drop=TRUE)
fit <- hmm(y,K=2,verb=TRUE,keep.y=TRUE,itmax=10)
cpe1 <- sp(model=fit) # Using the estimated parameters.
cpe2 <- sp(y,tpm=P,Rho=R,warn=FALSE) # Using the ``true'' parameters.
# The foregoing would issue a warning that Rho had no row names
# were it not for the fact that "warn" has been set to FALSE.
```

---

squantCI

*Simulation-quantile based confidence intervals.*

---

### Description

Calculates estimates of confidence intervals for the parameters of a model fitted by `hmm.discnp`. Uses a method based quantiles of estimates produced by simulation (or “parametric bootstrapping”).

**Usage**

```
squantCI(object, seed = NULL, alpha = 0.05, nsim=100, verbose = TRUE)
```

**Arguments**

object	An object of class <code>hmm.discnp</code> as returned by <code>hmm()</code> .
seed	Integer scalar serving as a seed for the random number generator. If left <code>NULL</code> the seed itself is chosen randomly from the set of integers between 1 and $10^5$ .
alpha	Positive real number strictly between 0 and 1. A set of $100*(1-\alpha)\%$ confidence intervals will be produced.
nsim	A positive integer. The number of simulations upon which the confidence interval estimates will be based.
verbose	Logical scalar; if <code>TRUE</code> , iteration counts will be printed out during each of the simulation and model-fitting stages.

**Details**

This function is currently applicable only to models fitted to univariate data. The confidence intervals calculated are for the “raw” parameters (entries of `tpm` with the last column dropped — since the rows sum to 1, and the entries of `Rho` with the last row dropped — since the columns sum to 1.

**Value**

A 2-by-`npar` matrix (where `npar` is the number of “independent” parameters in the model) whose columns form the estimated confidence intervals. The column labels indicate the parameters to which each column pertains, in a reasonably perspicuous manner. The row labels indicate the relevant quantiles in percentages.

This matrix has an attribute `seed` (the random number generation seed that was used) so that the calculations can be reproduced.

**Author(s)**

Rolf Turner <r.turner@auckland.ac.nz>

**See Also**

`scovmat()` `link{rhmm}()` `link{hmm}()`

**Examples**

```
## Not run:
y <- list(lindLandFlows$deciles, ftLiardFlows$deciles)
fit <- hmm(y, K=3)
CIs <- squantCI(fit, nsim=100)

## End(Not run)
```

**Description**

Discretised version of counts of faecal coliform bacteria in sea water samples collected at seven locations near Sydney NSW, Australia. There were four “controls”: Longreef, Bondi East, Port Hacking “50”, and Port Hacking “100” and three “outfalls”: Bondi Offshore, Malabar Offshore and North Head Offshore. At each location measurements were made at four depths: 0, 20, 40, and 60 meters. A large fraction of the counts are missing values.

**Usage**

SydColDisc

**Format**

A data frame with 5432 observations on the following 6 variables.

`y` A factor consisting of a discretisation of counts of faecal coliform count bacteria in sea water samples. The original measures were obtained by a repeated dilution process. The data were discretised using the `cut()` function with breaks given by `c(0, 1, 5, 25, 200, Inf)` and labels equal to `c("lo", "mlo", "m", "mhi", "hi")`.

`locn` a factor with levels Longreef, Bondi East, Port Hacking 50, Port Hacking 100, Bondi Offshore, Malabar Offshore and North Head Offshore.

`depth` a factor with levels 0 (0 metres), 20 (20 metres), 40 (40 metres), 60 (60 metres).

`ma.com` A factor with levels no and yes, indicating whether the Malabar sewage outfall had been commissioned.

`nh.com` A factor with levels no and yes, indicating whether the North Head sewage outfall had been commissioned.

`bo.com` A factor with levels no and yes, indicating whether the Bondi Offshore sewage outfall had been commissioned.

**Details**

The observations corresponding to each location-depth combination constitute a (discrete valued) time series. The sampling interval is ostensibly 1 week; distinct time series are ostensibly synchronous. The measurements were made over a 194 week period. Due to exigencies of weather, the unreliability of boats and other factors the collection times were actually highly irregular and have been rounded to the nearest week. Often no sample was obtained at a given site within a week of the putative collection time, in which the observed count is given as a missing value. In fact over **75%** of the counts are missing. See Turner et al. (1998) for more detail.

## Modelling

The hidden Markov models applied in the paper Turner et al. (1998) and in the paper Turner (2008) used a numeric version of the response in this data set. The numeric response was essentially a square root transformation of the original data, and the resulting values were modelled in terms of a Poisson distribution. See the references for details.

## Source

The original data were kindly supplied by Geoff Coade, of the New South Wales Environment Protection Authority (Australia)

## References

T. Rolf Turner, Murray A. Cameron, and Peter J. Thomson. Hidden Markov chains in generalized linear models. *Canadian J. Statist.* **26** (1998) 107 – 125.

Rolf Turner. Direct maximization of the likelihood of a hidden Markov model. *Comp. Statist. Data Anal.* **52** (2008) 4147–4160.

---

update.hmm.discnp	<i>Update a fitted hmm.discnp model.</i>
-------------------	--

---

## Description

An update() method for objects of class `hmm.discnp`.

## Usage

```
## S3 method for class 'hmm.discnp'
update(object,..., data, Kplus1=FALSE,
        tpm2=NULL, verbose=FALSE, method=NULL, optimiser=NULL,
        stationary=NULL, mixture=NULL, cis=NULL, tolerance=NULL,
        itmax=NULL, crit=NULL, newstyle=NULL, X=NULL, addIntercept=NULL)
```

## Arguments

object	An object of class <code>hmm.discnp</code> as returned by <code>hmm()</code> .
...	Not used.
data	The data set to which the (updated) model is to be fitted. See the description of the <code>y</code> argument of <code>hmm()</code> for more detail.
Kplus1	Logical scalar. Should the number of states be incremented by 1? If so then <code>tpm</code> (the transition probability matrix) is re-formed by <code>rbind()</code> -ing on a row all of whose entries are $1/K$ (where $K$ is the “old” number of states) and then <code>cbind()</code> -ing on a column of zeroes. The emission probability matrix <code>Rho</code> is reformed by <code>cbind()</code> -ing on a column all of whose entries are $1/m$ where $m$ is the number of discrete values of the emissions.

Note that the initial likelihood of the “new” model with  $K+1$  states will (should?) be exactly the same as that of the “old” fitted  $K$ -state model.

The `Kplus1` argument is provided mainly so as to provide a set of starting values for the fitting process which will guarantee the log likelihood of a  $K+1$ -state model will be at least as large as that of a  $K$ -state model fitted to the same data set.

Experience indicates that when `Kplus1=TRUE` is used, the fitting process does not “move very far” from the maximum log likelihood found for the  $K$ -state model. It is then advisable to try (many) random starting values so as to (try to) find the “true” maximum for the  $K+1$ -state model.

<code>tpm2</code>	The transition probability matrix to use when updating a model fitted with $K=1$ and <code>Kplus1=TRUE</code> . This argument is ignored otherwise. The default value of this argument is <code>matrix(0.5,2,2)</code> . The value of <code>tpm2</code> makes no difference to the <i>initial</i> value of the likelihood of the $K=2$ model (which will be identical to the likelihood of the fitted $K=1$ model that is being updated). Any two-by-two transition probability matrix “will do”. However the value of <code>tpm2</code> could conceivably have an impact on the final likelihood of the $K=2$ model to which the fitting procedure converges. This is particularly true if the method is (or is switched to) “LM”.
<code>verbose</code>	See the help for <code>hmm()</code> .
<code>method</code>	See the help for <code>hmm()</code> .
<code>optimiser</code>	See the help for <code>hmm()</code> .
<code>stationary</code>	See the help for <code>hmm()</code> .
<code>mixture</code>	See the help for <code>hmm()</code> .
<code>cis</code>	See the help for <code>hmm()</code> .
<code>tolerance</code>	See the help for <code>hmm()</code> .
<code>itmax</code>	See the help for <code>hmm()</code> .
<code>crit</code>	See the help for <code>hmm()</code> .
<code>newstyle</code>	See the help for <code>hmm()</code> .
<code>X</code>	See the help for <code>hmm()</code> .
<code>addIntercept</code>	See the help for <code>hmm()</code> .

### Details

Except for argument `X`, any arguments that are left `NULL` have their values supplied from the `args` component of `object`.

### Value

An object of class `hmm.discnp` with an additional component `init.log.like` which is the initial log likelihood calculated at the starting values of the parameters (which may be modified from the parameters returned in the object being updated, if `Kplus1` is `TRUE`). The calculation is done by the function `logLikHmm()`. Barring the strange and unforeseen, `init.log.like` should be (reassuringly) equal to `object$log.like`. See `hmm()` for details of the other components of the returned value.

**Author(s)**

Rolf Turner <r.turner@auckland.ac.nz>

**See Also**

[hmm\(\)](#) [rhmm.hmm.discnp\(\)](#)

**Examples**

```
set.seed(294)
fit <- hmm(WoodPeweeSong,K=2,rand.start=list(tpm=TRUE,Rho=TRUE),itmax=10)
xxx <- rhmm(fit,nsim=1)
sfit <- update(fit,data=xxx,itmax=10)
yyy <- with(SydColDisc,split(y,f=list(locn,depth)))
f1 <- hmm(yyy,K=1)
f2 <- update(f1,data=yyy,Kplus1=TRUE) # Big improvement, but ...
## Not run:
g2 <- hmm(yyy,K=2) # Substantially better than f2.

## End(Not run)
```

---

viterbi

*Most probable state sequence.*

---

**Description**

Calculates “the” most probable state sequence underlying each of one or more replicate observation sequences.

**Usage**

```
viterbi(y, model = NULL, tpm, Rho, ispd=NULL,log=FALSE, warn=TRUE)
```

**Arguments**

y	The observations for which the underlying most probable hidden states are required. May be a sequence of observations, or a list each entry of which constitutes an independent sequence of observations. If y is missing (and if model is not NULL) then y is extracted from model, provided that the y component of model is present. Otherwise an error is given.
model	An object describing a hidden Markov model, as fitted to the data set y by <a href="#">hmm()</a> .
tpm	The transition probability matrix for a hidden Markov model; ignored if model is non-null.

Rho	<p>An object specifying the probability distributions of the observations for a hidden Markov model. See <code>hmm()</code>. Ignored if <code>model</code> is non-null. Should bear some reasonable relationship to <code>y</code>.</p> <p>If Rho has dimension names (or if its entries have dimension names in the case where Rho is a list) then the appropriate dimension names must include all corresponding values of the observations. If a relevant vector of dimension names is NULL then it is formed as the sort unique values of the appropriate columns of the observation matrices. In this case the corresponding dimensions must match the number of unique values.</p>
ispd	<p>The initial state probability distribution for a hidden Markov model; ignored if <code>model</code> is non-null. Should bear some reasonable relationship to <code>y</code>. If <code>model</code> and <code>ispd</code> are both NULL then <code>ispd</code> is set equal to the stationary distribution calculated from <code>tpm</code>.</p>
log	<p>Logical scalar. Should logarithms be used in the recursive calculations of the probabilities involved in the Viterbi algorithm, so as to avoid underflow? If <code>log</code> is FALSE then underflow is avoided instead by a normalization procedure. The quantity <code>delta</code> (see Rabiner 1989, page 264) is replaced by <code>delta/sum(delta)</code> at each step. It should actually make no difference whether <code>log</code> is set to TRUE. I just included the option because I could. Also the HMM package uses the logarithm approach so setting <code>log=TRUE</code> might be of interest if comparisons are to be made between the results of the two packages.</p>
warn	<p>Logical scalar; should a warning be issued if Rho hasn't got relevant dimension names? (Note that if this is so, then the corresponding dimension names are formed from the sorted unique values of <code>y</code> or of the appropriate column(s) of <code>y</code>. And if <i>this</i> is so, then the user should be sure that the ordering of the entries of Rho corresponds properly to the the sorted unique values of <code>y</code>.) This argument is passed to the utility function <code>check.yval()</code> which actually issues the warning if <code>warn=TRUE</code>.</p>

### Details

Applies the Viterbi algorithm to calculate “the” most probable robable state sequence underlying each observation sequences.

### Value

If `y` consists of a single observation sequence, the value is the underlying most probable observation sequence, or a matrix whose columns consist of such sequences if there is more than one (equally) most probable sequence.

If `y` consists of a list of observation sequences, the value is a list each entry of which is of the form described above.

### Warning

There *may* be more than one equally most probable state sequence underlying a given observation sequence. This phenomenon can occur but appears to be unlikely to do so in practice.

**Thanks**

The correction made to the code so as to avoid underflow problems was made due to an inquiry and suggestion from Owen Marshall.

**Author(s)**

Rolf Turner <r.turner@auckland.ac.nz>

**References**

Rabiner, L. R., "A tutorial on hidden Markov models and selected applications in speech recognition," Proc. IEEE vol. 77, pp. 257 – 286, 1989.

**See Also**

[hmm\(\)](#), [rhmm\(\)](#), [mps\(\)](#), [pr\(\)](#), [viterbi\(\)](#)

**Examples**

```
# See the help for logLikHmm() for how to generate y.num and y.let.
## Not run:
fit.num      <- hmm(y.num,K=2,verb=TRUE,keep.y=TRUE)
v.1          <- viterbi(model=fit.num)
rownames(R) <- 1:5 # Avoids a (harmless) warning.
v.2          <- viterbi(y.num,tpm=P,Rho=R)
# P and R as in the help for logLikHmm() and for sp().

# Note that the order of the states has gotten swapped; 3-v.1[[1]]
# is identical to v.2[[1]]; for other k = 2, ..., 20, 3-v.1[[k]]
# is much more similar to v.2[[k]] than is v.1[[k]].

fit.let <- hmm(y.let,K=2,verb=TRUE,keep.y=TRUE)
v.3 <- viterbi(model=fit.let)
rownames(R) <- letters[1:5]
v.4 <- viterbi(y.let,tpm=P,Rho=R)

## End(Not run)
```

---

weissData

*Data from "An Introduction to Discrete-Valued Time Series"*

---

**Description**

Data sets from the book "An Introduction to Discrete-Valued Time Series" by Christian H. Weiß.

**Usage**

```
data(Bovine)
data(Cryptosporidiosis)
data(Downloads)
data(EricssonB_Jul2)
data(FattyLiver)
data(FattyLiver2)
data(goldparticle380)
data(Hanta)
data(InfantEEGsleepstates)
data(IPs)
data(LegionnairesDisease)
data(OffshoreRigcountsAlaska)
data(PriceStability)
data(Strikes)
data(WoodPeweeSong)
```

**Format**

- Bovine A character vector of length 8419.
- Cryptosporidiosis A numeric (integer) vector of length 365.
- Downloads A numeric (integer) vector of length 267.
- EricssonB\_Jul2 A numeric (integer) vector of length 460.
- FattyLiver2 A numeric (integer) vector of length 449.
- FattyLiver A numeric (integer) vector of length 928.
- goldparticle380 A numeric (integer) vector of length 380.
- Hanta A numeric (integer) vector of length 52.
- InfantEEGsleepstates A character vector of length 107.
- IPs A numeric (integer) vector of length 241.
- LegionnairesDisease A numeric (integer) vector of length 365.
- OffshoreRigcountsAlaska A numeric (integer) vector of length 417.
- PriceStability A numeric (integer) vector of length 152.
- Strikes A numeric (integer) vector of length 108.
- WoodPeweeSong A numeric (integer) vector of length 1327.

**Details**

For detailed information about each of these data sets, see the book cited in the **References**.

Note that the data sets `Cryptosporidiosis` and `LegionnairesDisease` are actually called `Cryptosporidiosis_02-08` and `LegionnairesDisease_02-08` in the given reference. The “suffixes” were removed since the minus sign causes problems in a variable name in R.

## Source

These data sets were kindly provided by Prof. Christian H. Weiß. The package author is also pleased to acknowledge the kind permission granted by Prof. Kurt Brännäs (Professor Emeritus of Economics at Umeå University) to include the Ericsson time series data set (EricssonB\_Ju12).

## References

Christian H. Weiß (2018). *An Introduction to Discrete-Valued Time Series*. Chichester: John Wiley & Sons.

## Examples

```
## Not run:
fit1 <- hmm(WoodPeweeSong,K=2,verbose=TRUE)
# EM converges in 6 steps --- suspicious.
set.seed(321)
fit2 <- hmm(WoodPeweeSong,K=2,verbose=TRUE,rand.start=list(tpm=TRUE,Rho=TRUE))
# 52 steps --- note the huge difference between fit1$log.like and fit2$log.like!
set.seed(321)
fit3 <- hmm(WoodPeweeSong,K=2,verbose=TRUE,method="bf",
            rand.start=list(tpm=TRUE,Rho=TRUE))
# log likelihood essentially the same as for fit2

## End(Not run)
```

# Index

- \*Topic **datagen**
  - misstify, 21
  - rhmm, 28
- \*Topic **datasets**
  - ccprSim, 3
  - hydroDat, 16
  - lesionCount, 18
  - SydColDisc, 35
  - weissData, 40
- \*Topic **methods**
  - anova.hmm.discnp, 2
  - update.hmm.discnp, 36
- \*Topic **models**
  - anova.hmm.discnp, 2
  - fitted.hmm.discnp, 5
  - hmm, 7
  - logLikHmm, 19
  - mps, 23
  - pr, 25
  - predict.hmm.discnp, 27
  - sp, 32
  - update.hmm.discnp, 36
  - viterbi, 38
- \*Topic **utilities**
  - cnvrtRho, 4
  - nafracCalc, 24
- \*Topic **utility**
  - scovmat, 30
  - squantCI, 33
- anova.hmm.discnp, 2
- Bovine (weissData), 40
- ccprSim, 3
- cnvrtRho, 4
- Cryptosporidiosis (weissData), 40
- cut, 35
- Downloads (weissData), 40
- EricssonB\_Jul2 (weissData), 40
- FattyLiver (weissData), 40
- FattyLiver2 (weissData), 40
- fitted.hmm.discnp, 5, 26–28, 33
- ftLiardFlows (hydroDat), 16
- goldparticle380 (weissData), 40
- Hanta (weissData), 40
- hmm, 2, 3, 5, 6, 7, 19, 20, 23, 24, 26, 27, 29–34, 36–40
- hydroDat, 16
- InfantEEGsleestates (weissData), 40
- IPs (weissData), 40
- LegionnairesDisease (weissData), 40
- lesionCount, 18
- linLandFlows (hydroDat), 16
- logLikHmm, 19
- misstify, 21, 25, 29, 30
- model.matrix, 11
- mps, 16, 23, 26, 33, 40
- nafracCalc, 21, 22, 24, 30
- nlm, 10, 11, 13
- OffshoreRigcountsAlaska (weissData), 40
- optim, 9, 10, 13
- portMannFlows (hydroDat), 16
- portMannSedCon (hydroDat), 16
- portMannSedLoads (hydroDat), 16
- pr, 20, 25, 33, 40
- predict.hmm.discnp, 27
- PriceStability (weissData), 40
- rhmm, 16, 21, 22, 24, 25, 28, 40
- rhmm.hmm.discnp, 38

scovmat, [30](#), [34](#)  
sp, [6](#), [7](#), [20](#), [26–28](#), [32](#)  
squantCI, [31](#), [33](#)  
Strikes (weissData), [40](#)  
SydColDisc, [35](#)  
  
update.hmm.discnp, [14](#), [36](#)  
  
viterbi, [16](#), [24](#), [26](#), [33](#), [38](#), [40](#)  
  
weissData, [40](#)  
WoodPeweeSong (weissData), [40](#)