

Package ‘BPHO’ documentation

of

April 7, 2008

Version 1.2-5

Title Bayesian Prediction with High-order Interactions

Author Longhai Li <longhai@math.usask.ca>

Maintainer Longhai Li <longhai@math.usask.ca>

Depends R (>= 2.5.1)

Description This software can be used in two situations. The first is to predict the next outcome based on the previous states of a discrete sequence. The second is to classify a discrete response based on a number of discrete covariates. In both situations, we use Bayesian logistic regression models that consider the high-order interactions. The models are trained with slice sampling method, a variant of Markov chain Monte Carlo. The time arising from using high-order interactions is reduced greatly by our compression technique that represents a group of original parameters as a single one in MCMC step.

License GPL (>=2)

URL <http://www.r-project.org>, <http://math.usask.ca/~longhai>

R topics documented:

comp_train_pred	2
compression	7
training	9
prediction	12
gendata	14

Index	16
--------------	-----------

comp_train_pred	<i>User-level functions for compressing parameters, training the models with MCMC, and making predictions for test cases</i>
-----------------	--

Description

The function `comp_train_pred` can be used for all of three tasks: compressing parameter, training the models with MCMC, and making prediction for test cases. When **new_compression=1**, it compresses parameters based on training cases and the information about parameter compression is written to the binary file `ptn_file`. When **new_compression=0**, it uses the existing `ptn_file`. When **iters_mc > 0**, it trains the models with Markov chain Monte Carlo and the Markov chain iterations are written to the binary file `mc_file`. The methods of writing to and reading from the files `ptn_file` and `mc_file` can be found from the documentations [compression](#) and [training](#). When **iters_pred > 0**, it predicts the responses of test cases and the result is written to the file `pred_file` and also returned as a value of this function.

The function `cv_comp_train_pred` is a short-cut function for performing cross-validation with the function `comp_train_pred`.

The argument `is_sequence=1` indicates that a sequence prediction model is fitted to the data, and `is_sequence=0` indicates that a general classification model based on discrete predictor variables is fitted.

Usage

```
comp_train_pred(
##### specify data information #####
test_x,train_x,train_y,no_cls=c(),nos_fth=c(),
##### specify for compression #####
is_sequence=1,order,ptn_file=".ptn.log",new_compression=1,do_comp=1,
##### specify for priors #####
alpha=1,log_sigma_widths=c(),log_sigma_modes=c(),
##### specify for mc sampling #####
mc_file=".mc.log",start_over=FALSE,iters_mc=200,iters_bt=10,
iters_sgm=50,w_bt=5,w_sgm=1,m_bt=20,m_sgm=20,ini_log_sigmas=c(),
##### specify for prediction #####
pred_file=c(),iter_b = 100,forward = 1,iters_pred = 100)

cv_comp_train_pred(
##### Specify data,order,no_fold #####
no_fold=10,train_x,train_y,no_cls=c(),nos_fth=c(),
##### specify for compressing#####
is_sequence=1,order,ptn_file=".ptn.log",new_compression=1,do_comp=1,
##### specify for priors #####
alpha=1,log_sigma_widths=c(),log_sigma_modes=c(),
##### specify for mc sampling #####
mc_file=".mc.log",iters_mc=200,iters_bt=10,iters_sgm=50,
w_bt=5,w_sgm=1,m_bt=20,m_sgm=20,ini_log_sigmas=c(),
```

```
##### specify for prediction #####
pred_file = c(), iter_b = 100, forward = 1, iters_pred = 100)
```

Arguments

test_x	Discrete features (also called inputs, covariates, independent variables, explanatory variables, predictor variables) of test data on which the predictions are based. The row is subject and the columns are inputs, which are coded with 1, 2, ..., with 0 reserved to represent that this input is not considered in a pattern. When the sequence prediction models are fitted, it is assumed that the first column is the state closest to the response. For example, a sequence 'x1, x2, x3, x4' is saved in test_x as 'x4, x3, x2, x1', for predicting the response 'x5'.
train_x	Discrete features of training data of the same format as test_x.
train_y	Discrete response of training data, a vector with length equal to the row of train_x. Assumed to be coded with 1, 2, ... no_cls.
no_cls	the number of possibilities (classes) of the response, default to the maximum value in train_y.
nos_fth	a vector, with each element storing the number of possibilities (classes) of each feature, default to the maximum value of each feature.
is_sequence	is_sequence=1 indicates that sequence prediction models are fitted to the data, and is_sequence=0 indicates that general classification models based on discrete predictor variables are fitted.
no_fold	Number of folders in cross-validation.
order	the order of interactions considered, default to the total number of features, i.e. ncol(train_x).
ptn_file	a character string, the name of the binary file to which the compression result is written. The method of writing to and reading from ptn_file can be found from the documentation for compression .
new_compression	new_compression=1 indicates removing the old file ptn_file if it exists and doing the compression once again. new_compression=0 indicates using the old file ptn_file without doing compression once again. Note that when new_compression=0, the specification related to training cases does not take effect.
do_comp	do_comp=1 indicates doing compression, and do_comp=0 indicates using original parametrization. This is used only to make comparison. In practice, we definitely recommend using our compression technique to reduce the number of parameters.
alpha	alpha=1 indicates that Cauchy prior is used, alpha=2 indicates that Gaussian prior is used.
log_sigma_widths, log_sigma_modes	two vectors of length order+1, which are interpreted as follows: the Gaussian distribution with location log_sigma_modes[o] and standard deviation log_sigma_widths[o] is the prior for 'log(sigmas[o])', which is the hyperparameter (width parameter of Gaussian distribution or Cauchy distribution) for the regression coefficients (i.e. 'beta's) associated with the interactions of order 'o'.

<code>mc_file</code>	A character string, the name of the binary file to which Markov chain is written. The method of writing to and reading from <code>mc_file</code> can be found from the documentation for training .
<code>start_over</code>	<code>start_over=TRUE</code> indicates that the existing file <code>mc_file</code> is deleted before a Markov chain sampling starts, otherwise the Markov chain will continue from the last iteration stored in <code>mc_file</code> .
<code>iters_mc, iters_bt, iters_sgm</code>	<code>iters_mc</code> iterations of super-transition will be run. Each super-transition consists of <code>iters_bt</code> iterations of updating 'beta's, and for each updating of 'beta's, the hyperparameters 'log(sigma)'s are updated <code>iters_sgm</code> times. When <code>iters_mc=0</code> , no Markov chain sampling will be run and other arguments related to Markov chain sampling take no effect.
<code>w_bt, w_sgm, m_bt, m_sgm</code>	<code>w_bt</code> is the amount of stepping-out in updating 'beta' with slice sampling, <code>m_bt</code> is the maximum number of stepping-out in slice sampling for updating 'beta'. <code>w_sgm</code> and <code>m_sgm</code> are interpreted similarly for sampling for 'log(sigma)'.
<code>ini_log_sigmas</code>	Initial values of 'log(sigma)', default to <code>log_sigma_modes</code> .
<code>pred_file</code>	A character string, the name of the file to which the prediction result is written. If <code>pred_file=c()</code> , the prediction result is printed out on screen (or sent to standard output).
<code>iter_b, forward, iters_pred</code>	Starting from <code>iter_b</code> , one of every forward Markov chain samples, with the number of total samples being \leq <code>iters_pred</code> and the maximum usable in the file <code>mc_file</code> , is used to make prediction.

Value

<code>times</code>	The time in second for, as this order, compressing parameters, training the model, predicting for test cases
<code>pred_result</code>	a data frame with first <code>no_cls</code> columns being the predictive probability and the next column being the predicted response value is returned.
<code>files</code>	Three character strings: the 1st is the name of the file storing compression information, the 2nd is the name of the file storing Markov chain, and the 3rd one is the name of the file containing the detailed prediction result, i.e., <code>pred_result</code>

Author(s)

Longhai Li, <http://math.usask.ca/~longhai>

References

<http://math.usask.ca/~longhai/doc/seqpred/seqpred.abstract.html>

See Also

[gendata](#), [compression](#), [training](#), [prediction](#)

Examples

```
## loading package
## library("BPHO", lib.loc=~/.rlib")

#####
#####The following are demonstrations of using the whole package
#####

## generate data from a hidden Markov model
data_hmm <- gen_hmm(n=200,p=10,no_h=8,no_o=2,
                    prob_h_stay=0.8,prob_o_stay=0.8)

## compressing parameters, training model, making prediction
comp_train_pred(
  ##### specify data information #####
  test_x=data_hmm$X[1:100,],train_x=data_hmm$X[-(1:100),],
  train_y=data_hmm$y[-(1:100)],no_cls=2,nos_fth=rep(2,10),
  ##### specify for compression #####
  is_sequence=1,order=4,ptn_file=".ptn_file.log",
  new_compression=1,do_comp=1,
  ##### specify for priors #####
  alpha=1,log_sigma_widths=c(),log_sigma_modes=c(),
  ##### specify for mc sampling #####
  mc_file=".mc_file.log",start_over=TRUE,itors_mc=100,
  iters_bt=1,itors_sgm=2,w_bt=5,w_sgm=1,
  m_bt=20,m_sgm=20,ini_log_sigmas=c(),
  ##### specify for prediction #####
  pred_file=".pred_file.csv",iter_b = 10,forward = 1,
  iters_pred = 90)

## display summary information about compression
display_ptn(ptn_file=".ptn_file.log")

## display the pattern information for group 1 and group 2
display_ptn(ptn_file=".ptn_file.log",gid=c(1,2))

## display the general information of Markov chain sampling
display_mc(mc_file=".mc_file.log")

## read Markov chain values of log-likelihood from ".mc_file.log"
read_mc(group="lprobs",ix=0,mc_file=".mc_file.log",
        iter_b=0,forward=1,n=100)

## particularly read `betas' by specifying the group and class id
read_betas(mc_file=".mc_file.log",ix_g=5,ix_cls=2,
           iter_b=0,forward=1,n=100)

## display the information on the pattern related to a `beta'
display_a_beta(mc_file=".mc_file.log",
              ptn_file=".ptn_file.log",id_beta=5)

## calculate the medians of samples of each 'beta'
```

```

calc_medians_betas(mc_file=".mc_file.log",iter_b=10,forward=1,n=90)

## evaluate prediction with true values of the response
evaluate_prediction(
  test_y=data_hmm$y[1:100],
  pred_result=read.csv(".pred_file.csv"),
  file_eval_details="eval_details")

#perform cross-validation with training data only
cv_comp_train_pred(
  ##### specify data information #####
  no_fold=2,train_x=data_hmm$X[-(1:100)],,
  train_y=data_hmm$y[-(1:100)],no_cls=2,nos_fth=rep(2,10),
  ##### specify for compression #####
  is_sequence=1,order=4,ptn_file=".ptn_file.log",
  new_compression=1,do_comp=1,
  ##### specify for priors #####
  alpha=1,log_sigma_widths=c(),log_sigma_modes=c(),
  ##### specify for mc sampling #####
  mc_file=".mc_file.log",iters_mc=100,
  iters_bt=1,iters_sgm=2,w_bt=5,w_sgm=1,
  m_bt=20,m_sgm=20,ini_log_sigmas=c(),
  ##### specify for prediction #####
  pred_file=".pred_file.csv",iter_b = 10,forward = 1,
  iters_pred = 90)

#####
#####

## generating a classification data
data_class <- gen_bin_ho(n=400,p=3,order=3,alpha=1,
  sigmas=c(0.3,0.2,0.1),nos_features=c(4,4,4),beta0=0)

## compressing parameters, training model, making prediction
comp_train_pred(
  ##### specify data information #####
  test_x=data_class$X[1:100],train_x=data_class$X[-(1:100)],,
  train_y=data_class$y[-(1:100)],no_cls=2,nos_fth=rep(4,3),
  ##### specify for compression #####
  is_sequence=0,order=3,ptn_file=".ptn_file.log",
  new_compression=1,do_comp=1,
  ##### specify for priors #####
  alpha=1,log_sigma_widths=c(),log_sigma_modes=c(),
  ##### specify for mc sampling #####
  mc_file=".mc_file.log",start_over=TRUE,iters_mc=500,
  iters_bt=1,iters_sgm=5,w_bt=5,w_sgm=0.5,
  m_bt=20,m_sgm=20,ini_log_sigmas=c(),
  ##### specify for prediction #####
  pred_file=".pred_file.csv",iter_b = 100,forward = 1,
  iters_pred = 400)

## display summary information about compression
display_ptn(ptn_file=".ptn_file.log")

```

```
## display the pattern information for group 1 and group 2
display_ptn(ptn_file=".ptn_file.log",gid=c(1,2))

## display the general information of Markov chain sampling
display_mc(mc_file=".mc_file.log")

## read Markov chain values of log-likelihood from ".mc_file.log"
read_mc(group="lprobs",ix=0,mc_file=".mc_file.log",
        iter_b=0,forward=1,n=500)

## particularly read `betas' by specifying the group and class id
read_betas(mc_file=".mc_file.log",ix_g=5,ix_cls=2,
           iter_b=0,forward=1,n=500)

## display the information on the pattern related to a `beta'
display_a_beta(mc_file=".mc_file.log",ptn_file=".ptn_file.log",
              id_beta=5)

## calculate the medians of samples of each 'beta'
calc_medians_betas(mc_file=".mc_file.log",iter_b=100,forward=1,n=400)

## evaluate prediction with true values of the response
evaluate_prediction(
  test_y=data_class$y[1:100],
  pred_result=read.csv(".pred_file.csv"),
  file_eval_details="eval_details")
```

compression

Functions related to parameter compression

Description

The function `compress` groups the patterns in a way such that the interaction patterns in a group are expressed by the same training cases. In training the models with MCMC, we need to use only one parameter for each group, which represents the sum of all the parameters in this group. The original parameters are seemingly compressed. A large amount of training time is saved by this compression techniques.

The result of this grouping is saved in a binary file in a way such that it can be retrieved as a linked list in C, with each node consisting of a description (an integer vector of fixed length) of the group of patterns and the indice (an integer vector of varying length, with 0 for the first training case) of training cases expressing this group of patterns. This file is needed to train the models with MCMC and to predict the responses of test cases using the function `comp_train_pred`.

The function `display_ptn` displays the summary information about this compression, such as the number of groups and total number of patterns expressed by the training cases. When `gids` is nonempty, it also displays the detailed information about the groups specified by `gids`, such as the pattern description and the indice of training cases associated with this group.

Usage

```
compress(features,nos_fth=c(),no_cases_ign=0,
         ptn_file=".ptn_file.log",quiet=1,
         do_comp=1,sequence=1,order=ncol(features))
display_ptn(ptn_file, gids=c())
```

Arguments

<code>features</code>	Discrete features (also called features,covariates,independent variables, explanatory variables, predictor variables) of training data on which the predictions are based. The row is subject and the columns are inputs, which are coded with 1,2,..., with 0 reserved to represent that this input is not considered in a pattern. When the sequence prediction models are fitted, it is assumed that the first column is the state closest to the response. For example, a sequence 'x1,x2,x3,x4' is saved in <code>test_x</code> as 'x4,x3,x2,x1', for predicting the response 'x5'.
<code>nos_fth</code>	a vector, with each element storing the number of possibilities (classes) of each feature, default to the maximum value of each feature.
<code>order</code>	the order of interactions considered, default to the total number of features, i.e. <code>ncol(features)</code> .
<code>ptn_file</code>	a character string, the name of the binary file to which the compression result is written.
<code>do_comp</code>	<code>do_comp=1</code> indicates doing compression, and <code>do_comp=0</code> indicates using original parametrization. This is used only to make comparison. In practice, we definitely recommend using our compression technique to reduce the number of parameters.
<code>sequence</code>	<code>sequence=1</code> indicates that sequence prediction models are fitted to the data, and <code>sequence=0</code> indicates that general classification models based on discrete predictor variables are fitted.
<code>gids</code>	an integer vector, containing the indice of groups whose information you want to display, with 0 for the first group.
<code>no_cases_ign</code>	When the number of training cases for a pattern is no more than <code>no_cases_ign</code> , this pattern will be ignored, default to 0, i.e. considering all interactions. So far there is no other justification to set it a value greater than 0, except that it can reduce the number of groups.
<code>quiet</code>	If <code>quiet=0</code> , some messages during compression are printed on screen for monitor the compression, if <code>quiet=1</code> the function works silently.

Value

The function `compress` returns no value. Instead, it saves the result of compression in the file `ptn_file`.

The function `display_ptn` returns a vector of 6 numbers. Their meanings are as follows: `is.sequence` – indicator whether a sequence model is fitted, `order` – the maximum order of interactions considered, `#groups` – the number of groups found, `#patterns` – the number of interaction patterns expressed by the training cases, `#cases` – the number of training cases, `#features` – the number of features.

When `gids` is nonempty, it also displays the details about the queried groups. The information printed on screen for each group is read as follows. Under **superpatterns**, it displays a compact description of the pattern group, which is in a special format defined in the references associated with this software. Under **expression**, it displays the indice of training cases that express this group of patterns. Under **sigmas**, it displays the number of patterns with a certain order, starting from order 0. This information is needed to compute the width parameter of the regression coefficient associated with this group from the values of hyperparameters ‘sigma’s.

See Also

`comp_train_pred`, `training`, `prediction`

Examples

```
## generate features
features <- gen_X(50,5,2)

## compressing the parameter based on 'features'
compress(features,nos_fth=rep(2,5),no_cases_ign=0,
          ptn_file=".ptn_file.log",quiet=1,do_comp=1,
          sequence=1,order=4)

## display the summary information in the file ".ptn_file.log"
display_ptn(".ptn_file.log")

## display the information for group #2 and #3
display_ptn(".ptn_file.log",gids=c(2,3))
```

Description

The models are trained with Markov chain Monte Carlo (MCMC) methods. Slice sampling is used to update ‘beta’s, the regression coefficients for groups, and ‘log(sigma)’, where ‘sigma’ is the width parameter of the prior for ‘beta’.

The function `training` carries out the Markov chain sampling, saving the Markov chain samples in a binary file `mc_file`.

The function `display_mc` displays the summary information in the file `mc_file`.

The function `read_mc` reads the Markov chain samples from the file `mc_file` at given iterations.

The function `read_betas` is based on the function `read_mc`. It specifically reads the ‘beta’ for given group and class identities.

The function `display_a_beta` displays both the pattern information for the group associated with the ‘beta’ specified by `id_beta`, and also return the full Markov chain samples of this ‘beta’.

The function `calc_medians_betas` returns the medians of the Markov chain samples for all ‘beta’s at specified iterations. This function is for discovering important interaction patterns. An interaction pattern with large absolute medians is highly suspected to be an important pattern for predicting the response.

Usage

```
display_mc(mc_file)
read_mc(mc_file, group, ix, iter_b=0, forward=1, n=c(), quiet=1)
read_betas(mc_file, ix_g, ix_cls, iter_b=0, forward=1, n=c(), quiet=1)
display_a_beta(id_beta, mc_file, ptn_file)
calc_medians_betas(mc_file, iter_b=0, forward=1, n=c())
training(mc_file, ptn_file, train_y, no_cls,
         alpha, log_sigma_widths,
         log_sigma_modes, ini_log_sigmas,
         iters_mc, iters_bt, iters_sgm,
         w_bt, w_sgm, m_bt, m_sgm)
```

Arguments

<code>mc_file</code>	A character string, the name of the binary file to which Markov chain is written.
<code>group</code>	A character string giving the group name of values. It can be one of ‘lprobs’, ‘lsigmas’, ‘betas’, ‘evals’. Group ‘lprobs’ contains: the values of log probabilities of data given the values of ‘beta’s (identified by <code>ix=0</code>), the value of log prior of ‘beta’s given ‘sigma’s (identified by <code>ix=1</code>), the value of log prior of ‘log(sigma)’s (identified by <code>ix=2</code>), and the value of log posterior (identified by <code>ix=3</code>), which is the sum of the previous three values. Group ‘lsigmas’ contains: the values of hyperparameters ‘log(sigma)’, with <code>ix</code> indicating the order, starting from 0. Group ‘betas’ contains: the values of ‘betas’, with <code>ix</code> indicating the index of ‘beta’. The ‘beta’s in each iteration is placed as that the <code>no_cls</code> values of ‘beta’s for pattern group ‘i’ are followed by the next <code>no_cls</code> values for pattern group ‘i+1’. The smallest index is 0. Group ‘evals’ contains: the average times of evaluating the posterior distribution in updating each ‘beta’ using slice sampling (identified by <code>ix=0</code>), and the average rejection rate of updating each ‘log(sigma)’ with Metropolis sampling (identified by <code>ix=1</code>).
<code>ix</code>	index of parameters inside each group, as discussed for <code>group</code> above.
<code>ix_g</code>	index of pattern group, starting from 0.
<code>ix_cls</code>	index of class, ranging from 1 to <code>no_cls</code> .
<code>id_beta</code>	index of ‘beta’, starting from 0.

<code>iter_b, forward, n</code>	Starting from <code>iter_b</code> , one of every <code>forward</code> Markov chain samples, with the number of total samples being $\leq n$ and the maximum usable in the file <code>mc_file</code> , is read.
<code>train_y</code>	Discrete response of training data. Assumed to be coded with 1,2,... <code>no_cls</code> .
<code>no_cls</code>	the number of possibilities (classes) of the response, default to the maximum value in <code>train_y</code> .
<code>alpha</code>	<code>alpha=1</code> indicates that Cauchy prior is used, <code>alpha=2</code> indicates that Gaussian prior is used.
<code>log_sigma_widths, log_sigma_modes</code>	two vectors of length <code>order+1</code> , which are interpreted as follows: the Gaussian distribution with location <code>log_sigma_modes[o]</code> and standard deviation <code>log_sigma_widths[o]</code> is the prior for ' <code>log(sigmas[o])</code> ', which is the hyperparameter (width parameter of Gaussian distribution or Cauchy distribution) for the regression coefficients (i.e. ' <code>beta</code> 's) associated with the interactions of order ' <code>o</code> '.
<code>ptn_file</code>	a character string, the name of the binary file where the compression result is saved. The method of writing to and reading from <code>ptn_file</code> can be found from the documentation for compression .
<code>iters_mc, iters_bt, iters_sgm</code>	<code>iters_mc</code> iterations of super-transition will be run. Each super-transition consists of <code>iters_bt</code> iterations of updating ' <code>beta</code> 's, and for each updating of ' <code>beta</code> 's, the hyperparameters ' <code>log(sigma)</code> 's are updated <code>iters_sgm</code> times. When <code>iters_mc=0</code> , no Markov chain sampling will be run and other arguments related to Markov chain sampling take no effect.
<code>w_bt, w_sgm, m_bt, m_sgm</code>	<code>w_bt</code> is the amount of stepping-out in updating ' <code>beta</code> ' with slice sampling, <code>m_bt</code> is the maximum number of stepping-out in slice sampling for updating ' <code>beta</code> '. <code>w_sgm</code> and <code>m_sgm</code> are interpreted similarly for sampling for ' <code>log(sigma)</code> '.
<code>ini_log_sigmas</code>	Initial values of ' <code>log(sigma)</code> ', default to <code>log_sigma_mode</code> .
<code>quiet</code>	<code>quiet=1</code> suppresses the messages printed during reading the file <code>mc_file</code> .

Value

The function `display_mc` returns a vector with names as `#iters, #class, #groups, order, alpha`.

The function `read_mc` returns the Markov chain samples for a variable at specified iterations.

The function `read_betas` returns the Markov chain samples for a '`beta`' at specified iterations.

The function `display_a_beta` displays the pattern group information for the group associated with the queried '`beta`', and also returns the Markov chain samples of this '`beta`'. The method of reading the on-screen messages about a pattern group is documented in [compression](#).

The function `calc_medians_betas` returns the medians of Markov chain samples of all ‘beta’s at given iterations.

The function `training` returns no value. Instead, the Markov chain samples are written to the binary file `mc_file`.

See Also

[comp_train_pred](#), [compression](#), [prediction](#)

Examples

```
## examples are given in comp_train_pred.
```

prediction

Functions related to prediction

Description

The function `predict_bpho` predicts the response of test cases.

The function `evaluate_prediction` evaluates the performance of the prediction in terms of average minus log probabilities and error rate. The function `split_cauchy` draws samples from a Cauchy distribution of two variables constraint to that their sum is fixed.

Usage

```
predict_bpho(test_x, no_cls, mc_file, ptn_file, iter_b, forward,
             iters_pred)
evaluate_prediction(test_y, pred_result, file_eval_details=c())
split_cauchy(n, s, signal, sigmasum, debug=1)
```

Arguments

<code>test_x</code>	Discrete features (also called inputs, covariates, independent variables, explanatory variables, predictor variables) of test data on which the predictions are based. The row is subject and the columns are inputs, which are coded with 1,2,..., with 0 reserved to represent that this input is not considered in a pattern. When the sequence prediction models are fitted, it is assumed that the first column is the state closest to the response. For example, a sequence ‘x1,x2,x3,x4’ is saved in <code>test_x</code> as ‘x4,x3,x2,x1’, for predicting the response ‘x5’.
<code>test_y</code>	Discrete responses of test data, a vector with length equal to the row of <code>test_x</code> . Assumed to be coded with 1,2,... <code>no_cls</code> .
<code>no_cls</code>	the number of possibilities (classes) of the response.
<code>ptn_file</code>	a character string, the name of the binary file to which the compression result is saved. The method of writing to and reading from <code>ptn_file</code> can be found from the documentation compression .

<code>mc_file</code>	A character string, the name of the binary file to which Markov chain is written. The method of writing to and reading from <code>mc_file</code> can be found from the documentation training .
<code>iter_b, forward, iters_pred</code>	Starting from <code>iter_b</code> , one of every forward Markov chain samples, with the number of total samples being \leq <code>iters_pred</code> and the maximum usable in the file <code>mc_file</code> , is used to make prediction.
<code>pred_result</code>	the value returned from the function <code>predict_bpho</code> .
<code>file_eval_details</code>	the details of evaluation is sent to the file <code>file_eval_details</code> .
<code>n</code>	number of samples one wishes to obtain.
<code>s</code>	sum of two Cauchy random variables.
<code>signal</code>	scale parameter for the first Cauchy random variable.
<code>sigmasum</code>	the sum of scale parameters for two Cauchy random variables.
<code>debug</code>	indicator whether you are debugging the C program.

Value

The function `predict_bpho` returns a data frame, with the first `no_cls` columns storing the predictive probabilities for each class, and the last column is the guess for the response by choosing the label of the class with largest predictive probability.

The function `evaluate_prediction` returns the following values:

<code>eval_details</code>	a data frame. The first column is the true response, the second is the guessed value by taking the label of class with largest predictive probability, the third is indicator whether a wrong decision is made, the last column is the predictive probability at the true class.
<code>error_rate</code>	the proportion of wrong prediction.
<code>aml1</code>	the average of minus log probabilities at true class, i.e. the average of the logarithms of the last column of <code>eval_details</code> .

The function `split_cauchy` returns a vector of `n` random numbers.

See Also

[comp_train_pred](#), [compression](#), [training](#)

Examples

```
## the function `predict_bpho' is demonstrated with the function
## `comp_train_pred' which calls `predict_bpho' inside.

## examples of 'evaluate_prediction' can be found from
## the documentation for comp_train_pred.

## testing the function split_cauchy
split_cauchy(100,10,1,5)
```

Description

`gen_hmm` generates sequences using hidden Markov models. `gen_bin_ho` generates general discrete data using logistic models, with high-order interactions considered; the response is binary. `text_to_3number` converts an English text file into sequence of 1 (special symbols such as space, symbol), 2 (vowel), 3 (consonant). `text_to_number` converts an English text into sequence of 1 - 27, 1-26 for letter a-z, and 27 for all other symbols.

Usage

```
gen_hmm(n, p, no_h, no_o, prob_h_stay, prob_o_stay)
gen_bin_ho(n, p, order, alpha, sigmas, nos_features, beta0)
text_to_number(p, file)
text_to_3number(p, file)
gen_X(n, p, K)
```

Arguments

<code>n</code>	number of cases.
<code>p</code>	number of features, or length of sequence.
<code>K</code>	number of possibilities for each feature.
<code>no_h</code>	number of states of hidden Markov chain.
<code>no_o</code>	number of states of output in hidden Markov model.
<code>prob_h_stay</code>	In simulating the hidden Markov chain, a chain will stay in its previous state with probability <code>prob_h_stay</code> , and move to other states with some minor probabilities adding up to $1 - \text{prob_h_stay}$.
<code>prob_o_stay</code>	In simulating the output state of hidden Markov model, the "output" is equal to ("hidden state" mod <code>no_o</code>)+1 with probability <code>prob_o_stay</code> and equally likely other states.
<code>order</code>	the order of interactions considered in simulating data from general classification models.
<code>alpha</code>	<code>alpha=2</code> indicates that Gaussian distributions are used to generate the "beta"s and <code>alpha=1</code> indicates that Cauchy distributions are used.
<code>sigmas</code>	hyperparameters in generating "beta"s, a vector of length <code>order</code> .
<code>nos_features</code>	number of states for each feature, i.e., the number of possibilities for each feature. A vector of length <code>p</code> .
<code>beta0</code>	intercept of linear function in generating classification data.
<code>file</code>	name of the file containing text file, a character string.

Value

X	values of predictors, a matrix. Each row is a case. For sequence, the data for each case (a row) is placed in the reverse order of time. For example, sequence "x1,x2,x3" is represented with a row of X: x3,x2,x1. The values of predictor X are coded by 1,2,3,...,nos_features. The function gen_X generates only this matrix.
y	values of the response, a vector, coded by 1,2,...
betas	a matrix of two columns saving the values of "betas" used in generating classification data. The first column is the absolute identity of this beta, and the 2nd column is the value. The total number of "betas" is saved in no_betas.

See Also

[comp_train_pred](#)

Examples

```
data_hmm <- gen_hmm(100,10,8,2,0.8,0.8)
data_bin_ho <- gen_bin_ho(100,3,2,1,c(5,2),c(3,3,3),0)
X <- gen_X(100,5,3)
```

Index

*Topic **classif**

- comp_train_pred, 1
- compression, 7
- prediction, 12
- training, 9

*Topic **datagen**

- gendata, 13

begin.BPHO(*comp_train_pred*), 1

calc_medians_betas(*training*), 9

comp_train_pred, 1, 7, 9, 11, 13, 14

compress(*compression*), 7

compression, 2–4, 7, 11–13

cv_comp_train_pred
(*comp_train_pred*), 1

display_a_beta(*training*), 9

display_mc(*training*), 9

display_ptn(*compression*), 7

evaluate_prediction(*prediction*),
12

gen_bin_ho(*gendata*), 13

gen_hmm(*gendata*), 13

gen_X(*gendata*), 13

gendata, 4, 13

predict_bpho(*prediction*), 12

prediction, 4, 9, 11, 12

read_betas(*training*), 9

read_mc(*training*), 9

split_cauchy(*prediction*), 12

text_to_3number(*gendata*), 13

text_to_number(*gendata*), 13

training, 2–4, 9, 9, 12, 13