

Package ‘lrgs’

August 11, 2020

Type Package

Title Linear Regression by Gibbs Sampling

Version 0.5.4

Date 2020-08-10

Author Adam Mantz

Maintainer Adam Mantz <amantz@slac.stanford.edu>

Description Implements a Gibbs sampler to do linear regression with multiple covariates, multiple responses, Gaussian measurement errors on covariates and responses, Gaussian intrinsic scatter, and a covariate prior distribution which is given by either a Gaussian mixture of specified size or a Dirichlet process with a Gaussian base distribution. Described further in Mantz (2016) <DOI:10.1093/mnras/stv3008>.

License MIT + file LICENSE

LazyLoad yes

Imports mvtnorm

URL <https://github.com/abmantz/lrgs>

NeedsCompilation no

Repository CRAN

Date/Publication 2020-08-11 14:20:03 UTC

R topics documented:

lrgs-package	2
Gibbs.post2dataframe	2
Gibbs.regression	3

Index	8
--------------	----------

lrgs-package

Linear Regression by Gibbs Sampling

Description

Implements a Gibbs sampler to do linear regression with multiple covariates, multiple responses, Gaussian measurement errors on covariates and responses, Gaussian intrinsic scatter, and a covariate prior distribution which is given by either a Gaussian mixture of specified size or a Dirichlet process with a Gaussian base distribution.

Details

Package: lrgs
Type: Package
Version: 0.5.4
Date: 2018-05-17
License: MIT
LazyLoad: yes

See help for Gibbs.regression

Author(s)

Adam Mantz <amantz@slac.stanford.edu>

Gibbs.post2dataframe

Linear Regression by Gibbs Sampling

Description

Transforms a set of posterior samples produced by Gibbs.regression into a data frame for more straightforward analysis.

Usage

```
Gibbs.post2dataframe(p)
```

Arguments

p an object returned from Gibbs.regression.

Value

A data frame with one column corresponding to each parameter stored in post.

Author(s)

Adam Mantz

See Also[Gibbs.regression](#)

Gibbs.regression

*Linear Regression by Gibbs Sampling***Description**

Runs a Gibbs sampler to simulate the posterior distribution of a linear model with (potentially) multiple covariates and response variables. Throughout this help file, we use the following notation: there are n data points, m response variables and p covariates.

Usage

```
Gibbs.regression(x.in, y.in, M, Nsamples, Ngauss = 1, dirichlet=FALSE,
  M.inv = NULL, intercept = TRUE, trace = "bs", fix = "", start = list(),
  B.prior.mean=NULL, B.prior.cov=NULL, Sigma.prior.scale=NULL,
  Sigma.prior.dof=-1, dp.prior.alpha=NULL, dp.prior.beta=NULL,
  mention.every=NA, save.every=NA, save.to=NA)
```

Arguments

<code>x.in</code>	the measured covariates. Either an $n \times p$ matrix, a vector (if $p=1$) or NULL (if $p=0$, i.e. a the model is constant).
<code>y.in</code>	the measured responses. Either an $n \times m$ matrix or a vector (if $m=1$).
<code>M</code>	a $(p+m) \times (p+m) \times n$ array holding the n measurement covariance matrices, with covariates ordered before responses within each matrix. If both <code>M</code> and <code>M.inv</code> are NULL, all measurement errors are taken to be zero.
<code>Nsamples</code>	the number of iterations of the Gibbs sampler to run.
<code>Ngauss</code>	number of Gaussian mixture components describing the distribution of covariates. Pass 0 for the special case of a uniform distribution.
<code>dirichlet</code>	if TRUE, the prior on covariates is given not by a Gaussian mixture of size <code>Ngauss</code> , but by a Dirichlet process with a Gaussian base distribution (see details). This can be thought of as a Gaussian mixture where <code>Ngauss</code> is learned from the data and marginalized over.
<code>M.inv</code>	as an alternative to <code>M</code> , the inverses of the n measurement covariance matrices in the same format. Note that the <code>M</code> argument takes precedence; pass <code>M=NULL</code> to use <code>M.inv</code> . If both <code>M</code> and <code>M.inv</code> are NULL, all measurement errors are taken to be zero.
<code>intercept</code>	if TRUE, the model includes constant (intercept) terms for each response. Otherwise, the intercepts are fixed to zero.

<code>trace</code>	determines which variables are returned at the end of the call. See details.
<code>fix</code>	determines which parameters are NOT varied. See details.
<code>start</code>	a list containing starting values for any of the model parameters. These are optional, with the exception of Sigma and Tau if they are fixed. See details.
<code>B.prior.mean</code>	optional vector giving the mean of a Gaussian prior to be applied to the coefficients. The order is B11, B21, ..., B12, B22, ..., where the first index refers to the covariate and the second to the response. Assumed to be zero if it is not specified but B.prior.cov is.
<code>B.prior.cov</code>	optional matrix giving the covariance of a Gaussian prior to be applied to the coefficients. The parameter order is the same as in B.prior.mean. If not specified, an improper uniform prior is used.
<code>Sigma.prior.scale</code>	optional matrix giving the scale parameter of an inverse-Wishart prior to be applied to the intrinsic covariance. Default is zero.
<code>Sigma.prior.dof</code>	number of degrees of freedom for the inverse-Wishart prior on the intrinsic covariance. Default is -1, which corresponds to the Jeffreys (i.e. minimally informative) prior when the default scale matrix is used. Using $-(m+1)$ degrees of freedom with the default scale matrix corresponds to a prior that is uniform in $\det(\text{Sigma})$.
<code>dp.prior.alpha</code>	shape parameter of a Gamma prior to be applied to the Dirichlet process concentration parameter, if <code>dirichlet=TRUE</code> .
<code>dp.prior.beta</code>	rate parameter of a Gamma prior to be applied to the Dirichlet process concentration parameter, if <code>dirichlet=TRUE</code> .
<code>mention.every</code>	if set to a positive integer N, a message will be printed after every N iterations to confirm that something is happening.
<code>save.every</code>	if set to a positive integer N, the result will be saved as an object named "res" to the file named in the <code>save.to</code> argument every N iterations.
<code>save.to</code>	the name of a file to periodically save the results to.

Details

An in-depth description of the model and algorithm can be found in the references below. The full set of parameters is

X: true values of all covariates, arranged as a design matrix

Y: true values of all responses

B: matrix of intercepts and coefficients

Sigma: intrinsic covariance matrix of the responses

G: vector encoding which Gaussian mixture component of the covariate distribution model each data point belongs to

pi: vector of relative proportions of the mixture components

mu: matrix holding the mean of each mixture component

Tau: covariance matrices for each mixture component

mu0: mean of the Gaussian hyper-prior applied to mu

U: covariance matrix of the hyper-prior for mu

W: matrix defining the Wishart hyper-prior for Tau

alpha: the Dirichlet process concentration parameter

Note that mu0, U and W are meaningful only if $N_{\text{gauss}} > 1$. If $N_{\text{gauss}} = 0$ these, as well as G, pi, mu and Tau, are ignored.

If a Dirichlet process prior on the covariates is used instead of a Gaussian mixture, then mu and Tau represent the hyperparameters of the Gaussian base distribution of the process. In this case, pi, mu0, U and W are not sampled, but G can still be interpreted as giving which cluster each data point belongs to. alpha is only meaningful if the Dirichlet process is being used.

In the trace and fix arguments, these parameters (in the order above) are indicated by the characters "xybsgpmtzuwa". So, for example, the default value trace="bs" means that only the coefficients and intrinsic scatter are returned to the user, while fix="x" would fix the values of X (typically to the measured input values, x.in).

The contents of the list optionally passed to start should be similar to those returned by this function, with one fewer dimension (corresponding to the multiple samples returned). That is, if the \$X item returned by this function has dimension $c(n, p+1, N_{\text{samples}})$, a valid value for start\$X has dimensions $c(n, p+1)$.

Value

A list containing samples of the model parameters specified by the trace argument. The parameters are defined above; the dimensionality of the result will be

X	$n, p+1, N_{\text{samples}}$
Y	n, m, N_{samples}
B	$p+1, m, N_{\text{samples}}$
Sigma	m, m, N_{samples}
G	n, N_{samples}
pi	$N_{\text{gauss}}, N_{\text{samples}}$
mu	$N_{\text{gauss}}, p, N_{\text{samples}}$
Tau	$p, p, N_{\text{gauss}}, N_{\text{samples}}$
mu0	p, N_{samples}
U	p, p, N_{samples}
W	p, p, N_{samples}
alpha	N_{samples}

This assumes intercept=TRUE; otherwise, replace p+1 with p everywhere. This list can be transformed into a data frame using the Gibbs.post2dataframe function.

Note

The output is a Monte Carlo Markov Chain, and therefore may take some time to converge to the true posterior distribution. Traces of the parameters of interest should be examined in order to identify and remove this "burn-in" period.

If your data set is missing the occasional covariate or response measurement, it should be sufficient to set the corresponding variance in M to a very large number. (Alternatively, this case can be handled exactly by passing the inverse-covariances through $M.inv$, and setting the appropriate elements to zero.) Similarly, if a subset of values are measured with no error, use a very small number for the variance (but do not attempt to pass zero covariance or Inf in the inverse-covariance).

Author(s)

Adam Mantz

References

Mantz (2016; MNRAS 457:1279; arXiv:1509.00908; doi:10.1093/mnras/stv3008) for this function, Kelly (2007, ApJ 665:1489; arXiv:0705.2774; doi:10.1086/519947) for the same approach with a single response variable and a Gaussian mixture of covariates, Neal (2000, Journal of Computational and Graphical Statistics 9:249) for the Dirichlet process algorithm employed.

See Also

[lm](#) for classical linear regression.

Examples

```
## example using the default Ngauss=1 with no measurement errors
x <- rnorm(500, 0, 5)
y <- pi*x + rnorm(length(x), 0, 0.1)
post <- Gibbs.regression(x, y, NULL, 50, trace='bsmt', fix='xy')
m <- lm(y~x)
plot(post$B[1,1,-(1:10)], col=4); abline(h=m$coefficients[1], lty=2, col=2)
plot(post$B[2,1,-(1:10)], col=4); abline(h=m$coefficients[2], lty=2, col=2)
plot(post$Sigma[1,1,-(1:10)], col=4); abline(h=var(m$residuals), lty=2, col=2)
plot(post$mu[1,1,-(1:10)], col=4); abline(h=mean(x), lty=2, col=2)
plot(post$Tau[1,1,1,-(1:10)], col=4); abline(h=var(x), lty=2, col=2)

## verbose example using a Dirichlet process, including measurement errors
## in practice, you would want a longer chain, i.e. larger nmc
xx <- rnorm(100, c(-15,0,15), 1)
yy <- xx + rnorm(length(xx)) + rnorm(length(xx), 0, 3)
xx <- xx + rnorm(length(xx))
M <- array(0, dim=c(2,2,length(xx)))
M[1,1,] <- 1
M[2,2,] <- 1
nmc <- 10
post <- Gibbs.regression(xx, yy, M, nmc, dirichlet=TRUE, trace='bsgmta', mention.every=1)
plot(xx, yy, col=post$G[,nmc]) # plot clusters at the last iteration
m <- lm(yy~xx)
```

```
plot(post$B[1,1,-1], col=4); abline(h=m$coefficients[1], lty=2, col=2)
plot(post$B[2,1,-1], col=4); abline(h=m$coefficients[2], lty=2, col=2)
plot(post$Sigma[1,1,-1], col=4); abline(h=var(m$residuals), lty=2, col=2)
plot(post$mu[1,1,-1], col=4); abline(h=mean(xx), lty=2, col=2)
plot(post$Tau[1,1,1,-1], col=4); abline(h=var(xx), lty=2, col=2)
plot(post$alpha[-1], col=4)
```

Index

* **multivariate**

Gibbs.regression, 3

* **package**

lrgs-package, 2

* **regression**

Gibbs.regression, 3

Gibbs.post2dataframe, 2

Gibbs.regression, 3, 3

lm, 6

lrgs (lrgs-package), 2

lrgs-package, 2