

Package ‘spNNGP’

March 30, 2020

Title Spatial Regression Models for Large Datasets using Nearest Neighbor Gaussian Processes

Version 0.1.4

Date 2020-03-18

Maintainer Andrew Finley <finleya@msu.edu>

Author Andrew Finley [aut, cre],
Abhirup Datta [aut],
Sudipto Banerjee [aut]

Depends R (>= 3.5.0), coda, Formula, RANN

Description Fits univariate Bayesian spatial regression models for large datasets using Nearest Neighbor Gaussian Processes (NNGP) detailed in Finley, Datta, Banerjee (2020) <arXiv:2001.09111>, and Finley, Datta, Cook, Morton, Andersen, and Banerjee (2019) <doi:10.1080/10618600.2018.1537924> and Datta, Banerjee, Finley, and Gelfand (2016) <doi:10.1080/01621459.2015.1044091>.

License GPL (>= 2)

Encoding UTF-8

URL <http://www.finley-lab.com>

Repository CRAN

NeedsCompilation yes

Date/Publication 2020-03-30 15:50:02 UTC

R topics documented:

BCEF	2
MI_TSCA	3
PGLogit	3
predict.NNGP	5
print.spDiag	8
spConjNNGP	9
spDiag	13
spNNGP	15
summary.NNGP	19
summary.PGLogit	20

BCEF	<i>Forest Canopy Height from NASA Goddard's LiDAR Hyperspectral and Thermal (G-LiHT) over Bonanza Creek Experimental Forest</i>
------	---

Description

Forest canopy height (FCH) estimates from NASA Goddard's LiDAR Hyperspectral and Thermal (G-LiHT; Cook et al. 2013) Airborne Imager and percent tree cover (Hansen et al. 2013) over a subset of Bonanza Creek Experimental Forest, AK, collected in Summer 2014.

The BCEF matrix columns are longitude (x), latitude (y), forest canopy height (FCH) in meters from ground, and Landsat derived percent tree cover (PTC) for 188,717 observations. Longitude and latitude are in Albers Equal Area (proj4string "+proj=aea +lat_1=55 +lat_2=65 +lat_0=50 +lon_0=-154 +x_0=0 +y_0=0 +ellps=GRS80 +towgs84=0,0,0,0,0,0 +units=km +no_defs"). The last column (holdout) identifies a subset of data useful for assessing wall-to-wall predictive performance.

Usage

data(BCEF)

Format

A matrix containing 188,717 rows and 6 columns named longitude, latitude, FCH, PTC, and holdout.

Source

G-LiHT data were downloaded from <https://gliht.gsfc.nasa.gov>.

References

Cook, B.D., L.W. Corp, R.F. Nelson, E.M. Middleton, D.C. Morton, J.T. McCorkel, J.G. Masek, K.J. Ranson, and V. Ly. (2013) NASA Goddard's Lidar, Hyperspectral and Thermal (G-LiHT) airborne imager. *Remote Sensing* 5:4045-4066.

Hansen, M.C., Potapov, P.V., Moore, R., Hancher, M., Turubanova, S.A., Tyukavina, A., Thau, D., Stehman, S.V., Goetz, S.J., Loveland, T.R., Kommareddy, A., Egorov, A., Chini, L., Justice, C.O., and Townshend, J.R.G. (2013), High-Resolution Global Maps of 21st-Century Forest Cover Change, *Science*, 342, 850-853.

MI_TSCA

*Occurrence of Tsuga canadensis (Eastern hemlock) in Michigan***Description**

Eastern hemlock (*Tsuga canadensis* L.) analyzed in Lany et al. (2019). The data comprise hemlock occurrence (binomial outcome) on 17,743 forest stands across Michigan, USA. A set of covariates were also observed at each stand and can be used to explain the probability of hemlock occurrence. Covariates included minimum winter temperature (MIN), maximum summer temperature (MAX), total precipitation in the coldest quarter of the year (WIP), total precipitation in the warmest quarter of the year (SUP), annual actual evapotranspiration (AET) and annual climatic water deficit (DEF). Spatial coordinates are recorded in Longitude (long) and latitude (lat) which are Albers Equal Area (proj4string "+proj=aea +lat_1=55 +lat_2=65 +lat_0=50 +lon_0=-154 +x_0=0 +y_0=0 +ellps=GRS80 +towgs84=0,0,0,0,0,0,0 +units=m +no_defs"). See Lany et al. (2019) for details.

Usage

```
data(MI_TSCA)
```

Format

A data frame containing 17,743 rows and 9 columns.

References

Lany, N. K., Zarnetske, P. L., Finley, A. O. and McCullough, D. G. 2019. Complimentary strengths of spatially-explicit and multi-species distribution models. – *Ecography* doi: 10.1111/ecog.04728

PGLogit

*Function for Fitting Logistic Models using Polya-Gamma Latent Variables***Description**

The function PGLogit fits logistic models to binomial data using Polya-Gamma latent variables.

Usage

```
PGLogit(formula, weights = 1, data = parent.frame(), n.samples,
         n.omp.threads = 1, fit.rep = FALSE, sub.sample, verbose = TRUE, ...)
```

Arguments

<code>formula</code>	a symbolic description of the regression model to be fit. See example below.
<code>weights</code>	specifies the number of trials for each observation. The default is 1 trial for each observation. Valid arguments are a scalar value that specifies the number of trials used if all observations have the same number of trials, and a vector of length n that specifies the number of trials for each observation when there are differences in the number of trials.
<code>data</code>	an optional data frame containing the variables in the model. If not found in data, the variables are taken from <code>environment(formula)</code> , typically the environment from which PGLo <i>git</i> is called.
<code>n.samples</code>	the number of posterior samples to collect.
<code>n.omp.threads</code>	a positive integer indicating the number of threads to use for SMP parallel processing. The package must be compiled for OpenMP support. For most Intel-based machines, we recommend setting <code>n.omp.threads</code> up to the number of hyperthreaded cores. Note, <code>n.omp.threads > 1</code> might not work on all systems.
<code>fit.rep</code>	if TRUE, regression fitted and replicate data will be returned. The argument <code>sub.sample</code> controls which MCMC samples are used to generate the fitted and replicated data.
<code>sub.sample</code>	an optional list that specifies the samples used for <code>fit.rep</code> . Valid tags are <code>start</code> , <code>end</code> , and <code>thin</code> . Given the value associated with the tags, the sample subset is selected using <code>seq(as.integer(start), as.integer(end), by=as.integer(thin))</code> . The default values are <code>start=floor(0.5*n.samples)</code> , <code>end=n.samples</code> and <code>thin=1</code> .
<code>verbose</code>	if TRUE, model specification and progress of the sampler is printed to the screen. Otherwise, nothing is printed to the screen.
<code>...</code>	currently no additional arguments.

Value

An object of class PGLo*git* which is a list comprising:

<code>p.beta.samples</code>	a coda object of posterior samples for the regression coefficients.
<code>y.hat.samples</code>	if <code>fit.rep=TRUE</code> , regression fitted values from posterior samples specified using <code>sub.sample</code> .
<code>y.hat.quantiles</code>	if <code>fit.rep=TRUE</code> , 0.5, 0.025, and 0.975 quantiles of the <code>y.hat.samples</code> .
<code>y.rep.samples</code>	if <code>fit.rep=TRUE</code> , replicated outcome from posterior samples specified using <code>sub.sample</code> .
<code>y.rep.quantiles</code>	if <code>fit.rep=TRUE</code> , 0.5, 0.025, and 0.975 quantiles of the <code>y.rep.samples</code> .
<code>s.indx</code>	if <code>fit.rep=TRUE</code> , the subset index specified with <code>sub.sample</code> .
<code>run.time</code>	MCMC sampler execution time reported using <code>proc.time()</code> .

The return object will include additional objects used for subsequent prediction and/or model fit evaluation.

Note

Some of the underlying code used for generating random number from the Polya-Gamma distribution is taken from the **pgdraw** package written by Daniel F. Schmidt and Enes Makalic. Their code implements Algorithm 6 in PhD thesis of Jesse Bennett Windle (2013) <https://repositories.lib.utexas.edu/handle/2152/21842>.

Author(s)

Andrew O. Finley <finleya@msu.edu>,
Abhirup Datta <abhidatta@jhu.edu>,
Sudipto Banerjee <sudipto@ucla.edu>

References

Polson, N.G., J.G. Scott, and J. Windle. (2013) Bayesian Inference for Logistic Models Using Polya-Gamma Latent Variables. *Journal of the American Statistical Association*, 108:1339-1349.

Examples

```
##Generate binary data
set.seed(1)
n <- 100

x <- cbind(1, rnorm(n), runif(n,0,1))
beta <- c(0.1,-5, 5)
p <- 1/(1+exp(-(x%*%beta)))

##Assume 5 trials per outcome
weights <- rep(5, n)

y <- rbinom(n, size=weights, prob=p)

m <- PGLogit(y~x-1, weights = rep(5, n), n.samples = 1000)

summary(m)
```

predict.NNGP

Function for prediction at new locations using NNGP models.

Description

The function `predict` collects posterior predictive samples for a set of new locations given an object of class `NNGP`.

Usage

```
## S3 method for class 'NNGP'
predict(object, X.0, coords.0, sub.sample,
        n.omp.threads = 1, verbose=TRUE, n.report=100, ...)
```

Arguments

object	an object of class NNGP.
X.0	the design matrix for prediction locations. An intercept should be provided in the first column if one is specified in sp.obj model.
coords.0	the spatial coordinates corresponding to X.0.
sub.sample	an optional list that specifies the samples to included in the composition sampling a non-Conjugate model. Valid tags are start, end, and thin. Given the value associated with the tags, the sample subset is selected using seq(as.integer(start), as.integer(end), by=thin). The default values are start=floor(0.5*n.samples), end=n.samples and thin=1.
n.omp.threads	a positive integer indicating the number of threads to use for SMP parallel processing. The package must be compiled for OpenMP support. For most Intel-based machines, we recommend setting n.omp.threads up to the number of hyperthreaded cores. Note, n.omp.threads > 1 might not work on some systems.
verbose	if TRUE, model specification and progress of the sampler is printed to the screen. Otherwise, nothing is printed to the screen.
n.report	the interval to report sampling progress.
...	currently no additional arguments.

Value

An object of class predict.NNGP which is a list comprising:

p.y.0	a matrix that holds the response variable posterior predictive samples where rows are locations corresponding to coords.0 and columns are samples.
p.w.0	a matrix that holds the random effect posterior predictive samples where rows are locations corresponding to coords.0 and columns are samples. This is only returned if the input class has method = "latent".
run.time	execution time reported using proc.time().

Author(s)

Andrew O. Finley <finleya@msu.edu>,
 Abhirup Datta <abhidatta@jhu.edu>,
 Sudipto Banerjee <sudipto@ucla.edu>

References

Datta, A., S. Banerjee, A.O. Finley, and A.E. Gelfand. (2016) Hierarchical Nearest-Neighbor Gaussian process models for large geostatistical datasets. *Journal of the American Statistical Association*, 111:800-812.

Finley, A.O., A. Datta, B.D. Cook, D.C. Morton, H.E. Andersen, and S. Banerjee. (2019) Efficient algorithms for Bayesian Nearest Neighbor Gaussian Processes. *Journal of Computational and Graphical Statistics*, <https://doi.org/10.1080/10618600.2018.1537924>.

Examples

```
rmvn <- function(n, mu=0, V = matrix(1)){
  p <- length(mu)
  if(any(is.na(match(dim(V),p))))
    stop("Dimension problem!")
  D <- chol(V)
  t(matrix(rnorm(n*p), ncol=p)%*%D + rep(mu,rep(n,p)))
}

##Make some data
set.seed(1)
n <- 100
coords <- cbind(runif(n,0,1), runif(n,0,1))

x <- cbind(1, rnorm(n))

B <- as.matrix(c(1,5))

sigma.sq <- 5
tau.sq <- 1
phi <- 3/0.5

D <- as.matrix(dist(coords))
R <- exp(-phi*D)
w <- rmvn(1, rep(0,n), sigma.sq*R)
y <- rnorm(n, x%*%B + w, sqrt(tau.sq))

ho <- sample(1:n, 50)

y.ho <- y[ho]
x.ho <- x[ho,,drop=FALSE]
w.ho <- w[ho]
coords.ho <- coords[ho,]

y <- y[-ho]
x <- x[-ho,,drop=FALSE]
w <- w[-ho,,drop=FALSE]
coords <- coords[-ho,]

##Fit a Response, Latent, and Conjugate NNGP model
n.samples <- 500
```

```

starting <- list("phi"=phi, "sigma.sq"=5, "tau.sq"=1)

tuning <- list("phi"=0.5, "sigma.sq"=0.5, "tau.sq"=0.5)

priors <- list("phi.Unif"=c(3/1, 3/0.01), "sigma.sq.IG"=c(2, 5), "tau.sq.IG"=c(2, 1))

cov.model <- "exponential"

n.report <- 500

##Latent
m.s <- spNNGP(y~x-1, coords=coords, starting=starting, method="latent", n.neighbors=10,
             tuning=tuning, priors=priors, cov.model=cov.model,
             n.samples=n.samples, n.omp.threads=1, n.report=n.report)

p.s <- predict(m.s, X.0 = x.ho, coords.0 = coords.ho, n.omp.threads=1)

plot(apply(p.s$p.w.0, 1, mean), w.ho)
plot(apply(p.s$p.y.0, 1, mean), y.ho)

##Response
m.r <- spNNGP(y~x-1, coords=coords, starting=starting, method="response", n.neighbors=10,
             tuning=tuning, priors=priors, cov.model=cov.model,
             n.samples=n.samples, n.omp.threads=1, n.report=n.report)

p.r <- predict(m.r, X.0 = x.ho, coords.0 = coords.ho, n.omp.threads=1)

points(apply(p.r$p.y.0, 1, mean), y.ho, pch=19, col="blue")

##Conjugate
theta.alpha <- c(phi, tau.sq/sigma.sq)
names(theta.alpha) <- c("phi", "alpha")

m.c <- spConjNNGP(y~x-1, coords=coords, n.neighbors=10,
                theta.alpha=theta.alpha, sigma.sq.IG=c(2, sigma.sq),
                cov.model=cov.model, n.omp.threads=1)

p.c <- predict(m.c, X.0 = x.ho, coords.0 = coords.ho, n.omp.threads=1)

points(p.c$y.0.hat, y.ho, pch=19, col="orange")

```

print.spDiag

Methods for spDiag Object

Description

Methods for extracting information from spDiag.

Usage

```
## S3 method for class 'spDiag'
print(x, ...)
```

Arguments

x object of class spDiag.
... currently not used.

Details

A standard extractor function for printing objects of class spDiag.

spConjNNGP	<i>Function for Fitting Univariate Bayesian Conjugate Spatial Regression Models</i>
------------	---

Description

The function spConjNNGP fits Gaussian univariate Bayesian conjugate spatial regression models using Nearest Neighbor Gaussian Processes (NNGP).

Usage

```
spConjNNGP(formula, data = parent.frame(), coords, knots, n.neighbors = 15,
            theta.alpha, sigma.sq.IG, cov.model = "exponential",
            k.fold = 5, score.rule = "crps",
            X.0, coords.0, n.omp.threads = 1, search.type = "cb",
            ord, return.neighbor.info = TRUE,
            neighbor.info, fit.rep = FALSE, n.samples, verbose = TRUE, ...)
```

Arguments

formula a symbolic description of the regression model to be fit. See example below.
data an optional data frame containing the variables in the model. If not found in data, the variables are taken from environment(formula), typically the environment from which spConjNNGP is called.
coords an $n \times 2$ matrix of the observation coordinates in R^2 (e.g., easting and northing), or if data is a data frame then coords can be a vector of length two comprising coordinate column names or indices. There can be no duplicate locations.
knots an $r \times 2$ matrix of the observation coordinates in R^2 (e.g., easting and northing). Adding the knots argument invokes SLGP, see Shin et al. (2019) below.
n.neighbors number of neighbors used in the NNGP.

theta.alpha	a vector or matrix of parameter values for phi, nu, and alpha, where $\alpha = \tau^2/\sigma^2$ and nu is only required if cov.model="matern". A vector is passed to run the model using one set of parameters. The vector elements must be named and hold values for phi, nu, and alpha. If a matrix is passed, columns must be named and hold values for phi, nu, and alpha. Each row in the matrix defines a set of parameters for which the model will be run.
sigma.sq.IG	a vector of length two that holds the hyperparameters, <i>shape</i> and <i>scale</i> respectively, for the inverse-Gamma prior on σ^2 .
cov.model	a quoted keyword that specifies the covariance function used to model the spatial dependence structure among the observations. Supported covariance model key words are: "exponential", "matern", "spherical", and "gaussian". See below for details.
k.fold	specifies the number of <i>k</i> folds for cross-validation. If theta.alpha is a vector then cross-validation is not performed and k-fold and score.rule are ignored. In <i>k</i> -fold cross-validation, the data specified in model is randomly partitioned into <i>k</i> equal sized subsamples. Of the <i>k</i> subsamples, <i>k</i> -1 subsamples are used to fit the model and the remaining <i>k</i> samples are used for prediction. The cross-validation process is repeated <i>k</i> times (the folds). Root mean squared prediction error (RMSPE) and continuous ranked probability score (CRPS; Gneiting and Raftery, 2007) rules are averaged over the <i>k</i> fold prediction results and reported for the parameter sets defined by theta.alpha. The parameter set that yields the <i>best</i> performance based on the scoring rule defined by score.rule is used to fit the final model that uses all the data and make predictions if X.0 and coords.0 are supplied. Results from the <i>k</i> -fold cross-validation are returned in the k.fold.scores matrix.
score.rule	a quoted keyword "rmspe" or "crps" that specifies the scoring rule used to select the <i>best</i> parameter set, see argument definition for k.fold for more details.
X.0	the design matrix for prediction locations. An intercept should be provided in the first column if one is specified in model.
coords.0	the spatial coordinates corresponding to X.0.
n.omp.threads	a positive integer indicating the number of threads to use for SMP parallel processing. The package must be compiled for OpenMP support. For most Intel-based machines, we recommend setting n.omp.threads up to the number of hyperthreaded cores. Note, n.omp.threads > 1 might not work on some systems.
search.type	a quoted keyword that specifies type of nearest neighbor search algorithm. Supported method key words are: "cb" and "brute". The "cb" should generally be much faster. If locations do not have identical coordinate values on the axis used for the nearest neighbor ordering (see ord argument) then "cb" and "brute" should produce identical neighbor sets. However, if there are identical coordinate values on the axis used for nearest neighbor ordering, then "cb" and "brute" might produce different, but equally valid, neighbor sets, e.g., if data are on a grid.
ord	an index vector of length <i>n</i> used for the nearest neighbor search. Internally, this vector is used to order coords, i.e., coords[ord,], and associated data. Nearest neighbor candidates for the <i>i</i> -th row in the ordered coords are rows 1:(i-1),

with the `n.neighbors` nearest neighbors being those with the minimum euclidean distance to the location defined by ordered `coords[i,]`. The default used when `ord` is not specified is x-axis ordering, i.e., `order(coords[,1])`. This argument should typically be left blank. This argument will be ignored if the `neighbor.info` argument is used.

<code>return.neighbor.info</code>	if TRUE, a list called <code>neighbor.info</code> containing several data structures used for fitting the NNGP model is returned. If there is no change in input data or <code>n.neighbors</code> , this list can be passed to subsequent <code>spNNGP</code> calls via the <code>neighbor.info</code> argument to avoid the neighbor search, which can be time consuming if n is large. In addition to the several cryptic data structures in <code>neighbor.info</code> there is a list called <code>n.indx</code> that is of length n . The i -th element in <code>n.indx</code> corresponds to the i -th row in <code>coords[ord,]</code> and holds the vector of that location's nearest neighbor indices. This list can be useful for plotting the neighbor graph if desired.
<code>neighbor.info</code>	see the <code>return.neighbor.info</code> argument description above.
<code>fit.rep</code>	if TRUE, regression fitted and replicate data will be returned. The argument <code>n.samples</code> controls the number of fitted and replicated data samples.
<code>n.samples</code>	gives the number of posterior samples returned. Note, point and associated variance estimates for model parameters are not based on posterior samples. Only specify <code>n.samples</code> if you wish to generate samples from parameters' posteriors (this is an exact sampling algorithm). If <code>fit.rep</code> is TRUE, then <code>n.samples</code> also controls the number of fitted and replicated data samples.
<code>verbose</code>	if TRUE, model specification and progress is printed to the screen. Otherwise, nothing is printed to the screen.
<code>...</code>	currently no additional arguments.

Value

An object of class NNGP and conjugate, and, if `knots` is provided, SLGP. Among other elements, the object contains:

<code>theta.alpha</code>	the input <code>theta.alpha</code> vector, or <i>best</i> (according to the selected scoring rule) set of parameters in the <code>theta.alpha</code> matrix. All subsequent parameter estimates are based on this parameter set.
<code>beta.hat</code>	a matrix of regression coefficient estimates corresponding to the returned <code>theta.alpha</code> .
<code>beta.var</code>	<code>beta.hat</code> variance-covariance matrix.
<code>sigma.sq.hat</code>	estimate of σ^2 corresponding to the returned <code>theta.alpha</code> .
<code>sigma.sq.var</code>	<code>sigma.sq.hat</code> variance.
<code>k.fold.scores</code>	results from the k-fold cross-validation if <code>theta.alpha</code> is a matrix.
<code>y.0.hat</code>	prediction if <code>X.0</code> and <code>coords.0</code> are specified.
<code>y.0.var.hat</code>	<code>y.0.hat</code> variance.
<code>n.neighbors</code>	number of neighbors used in the NNGP.
<code>neighbor.info</code>	returned if <code>return.neighbor.info=TRUE</code> see the <code>return.neighbor.info</code> argument description above.

`run.time` execution time for parameter estimation reported using `proc.time()`. This time does not include nearest neighbor search time for building the neighbor set.

Author(s)

Andrew O. Finley <finleya@msu.edu>,
Abhirup Datta <abhidatta@jhu.edu>,
Sudipto Banerjee <sudipto@ucla.edu>

References

Datta, A., S. Banerjee, A.O. Finley, and A.E. Gelfand. (2016) Hierarchical Nearest-Neighbor Gaussian process models for large geostatistical datasets. *Journal of the American Statistical Association*, 111:800-812.

Finley, A.O., A. Datta, B.D. Cook, D.C. Morton, H.E. Andersen, and S. Banerjee. (2019) Efficient algorithms for Bayesian Nearest Neighbor Gaussian Processes. *Journal of Computational and Graphical Statistics*, <https://doi.org/10.1080/10618600.2018.1537924>.

Gneiting, T and A.E. Raftery. (2007) Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102:359-378.

Shirota, S., A.O. Finley, B.D. Cook, and S. Banerjee (2019) Conjugate Nearest Neighbor Gaussian Process models for efficient statistical interpolation of large spatial data. <https://arxiv.org/abs/1907.10109>.

Examples

```
rmvn <- function(n, mu=0, V = matrix(1)){
  p <- length(mu)
  if(any(is.na(match(dim(V),p))))
    stop("Dimension problem!")
  D <- chol(V)
  t(matrix(rnorm(n*p), ncol=p)%*%D + rep(mu,rep(n,p)))
}

##Make some data
set.seed(1)
n <- 2000
coords <- cbind(runif(n,0,1), runif(n,0,1))

x <- cbind(1, rnorm(n))

B <- as.matrix(c(1,5))

sigma.sq <- 5
tau.sq <- 1
phi <- 3/0.5

D <- as.matrix(dist(coords))
R <- exp(-phi*D)
w <- rmvn(1, rep(0,n), sigma.sq*R)
```

```

y <- rnorm(n, x%*%B + w, sqrt(tau.sq))

ho <- sample(1:n, 1000)

y.ho <- y[ho]
x.ho <- x[ho,,drop=FALSE]
w.ho <- w[ho]
coords.ho <- coords[ho,]

y <- y[-ho]
x <- x[-ho,,drop=FALSE]
w <- w[-ho,,drop=FALSE]
coords <- coords[-ho,]

##Fit a Conjugate NNGP model and predict for the holdout
sigma.sq.IG <- c(2, sigma.sq)

cov.model <- "exponential"

g <- 10
theta.alpha <- cbind(seq(phi,30,length.out=g), seq(tau.sq/sigma.sq,5,length.out=g))

colnames(theta.alpha) <- c("phi", "alpha")

m.c <- spConjNNGP(y~x-1, coords=coords, n.neighbors = 10,
                 X.0 = x.ho, coords.0 = coords.ho,
                 k.fold = 5, score.rule = "crps",
                 n.omp.threads = 1,
                 theta.alpha = theta.alpha, sigma.sq.IG = sigma.sq.IG, cov.model = cov.model)

m.c

```

spDiag

Model fit diagnostics

Description

The function spDiag calculates measurements of model fit for objects of class NNGP and PGLogit.

Usage

```
spDiag(object, sub.sample, ...)
```

Arguments

object an object of class NNGP or PGLogit.

sub.sample	an optional list that specifies the samples to included in the computations. Valid tags are start, end, and thin. Given the value associated with the tags, the sample subset is selected using <code>seq(as.integer(start), as.integer(end), by=as.integer(thin))</code> . The default values are <code>start=floor(0.5*n.samples)</code> , <code>end=n.samples</code> and <code>thin=1</code> . If sub.samples is not specified, then it is taken from object, or, if not aviable in object the default values of start, end, and thin are used. Note, if the object is a NNGP response model and n is large, then computing the replicated data needed for GPD and GRS can take a long time.
...	currently no additional arguments.

Value

A list with the following tags:

DIC	a data frame holding Deviance information criterion (DIC) and associated values. Values in DIC include DIC the criterion (lower is better), D a goodness of fit, and pD the effective number of parameters, see Spiegelhalter et al. (2002) for details.
GPD	a data frame holding $D=G+P$ and associated values. Values in GPD include G a goodness of fit, P a penalty term, and D the criterion (lower is better), see Gelfand and Ghosh (1998) for details.
GRS	a scoring rule, see Equation 27 in Gneiting and Raftery (2007) for details.
WAIC	a data frame hold Watanabe-Akaike information criteria (WAIC) and associated values. Values in WAIC include LPPD log pointwise predictive density, P.1 penalty term defined in unnumbered equation above Equation (11) in Gelman et al. (2014), P.2 an alternative penalty term defined in Equation (11), and the criteria WAIC.1 and WAIC.2 (lower is better) computed using P.1 and P.2, respectively.
y.rep.samples	if y.rep.samples in object were not used (or not available), then the newly computed y.rep.samples is returned.
y.fit.samples	if y.fit.samples in object were not used (or not available), then the newly computed y.fit.samples is returned.
s.indx	the index of samples used for the computations.

Author(s)

Andrew O. Finley <finleya@msu.edu>
Sudipto Banerjee <sudipto@ucla.edu>

References

- Spiegelhalter, D.J., Best, N.G., Carlin, B.P., van der Linde, A. (2002). Bayesian measures of model complexity and fit (with discussion). *Journal of the Royal Statistical Society, Series B.*, 64:583-639.
- Gelfand A.E. and Ghosh, S.K. (1998). Model choice: a minimum posterior predictive loss approach. *Biometrika*, 85:1-11.
- Gelman, A., Hwang, J., and Vehtari, A. (2014). Understanding predictive information criteria for Bayesian models. *Statistics and Computing*, 24:997-1016.

Gneiting, T. and Raftery, A.E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102:359-378.

 spNNGP

Function for Fitting Univariate Bayesian Spatial Regression Models

Description

The function spNNGP fits Gaussian and non-Gaussian univariate Bayesian spatial regression models using Nearest Neighbor Gaussian Processes (NNGP).

Usage

```
spNNGP(formula, data = parent.frame(), coords, method = "response",
        family="gaussian", weights, n.neighbors = 15,
        starting, tuning, priors, cov.model = "exponential",
        n.samples, n.omp.threads = 1, search.type = "cb", ord,
        return.neighbor.info = FALSE, neighbor.info,
        fit.rep = FALSE, sub.sample, verbose = TRUE, n.report = 100, ...)
```

Arguments

formula	a symbolic description of the regression model to be fit. See example below.
data	an optional data frame containing the variables in the model. If not found in data, the variables are taken from environment(formula), typically the environment from which spNNGP is called.
coords	an $n \times 2$ matrix of the observation coordinates in R^2 (e.g., easting and northing), or if data is a data frame then coords can be a vector of length two comprising coordinate column names or indices. There can be no duplicate locations.
method	a quoted keyword that specifies the NNGP sampling algorithm. Supported method keywords are: "response" and "latent". When n is large, the "response" algorithm should be faster and provide finer control over Metropolis acceptance rate for covariance parameters. In general, unless estimates of spatial random effects are needed, the "response" algorithm should be used.
family	a quoted keyword that specifies the data likelihood. Choices are "gaussian" for continuous outcome and "binomial" for discrete outcome which assumes a logistic link modeled using Polya-Gamma latent variables.
weights	specifies the number of trials for each observation when family="binomial". The default is 1 trial for each observation. Valid arguments are a scalar value that specifies the number of trials used if all observations have the same number of trials, and a vector of length n that specifies the number of trials for each observation used there are differences in the number of trials.
n.neighbors	number of neighbors used in the NNGP.
starting	a list with each tag corresponding to a parameter name. Valid tags are beta, sigma.sq, tau.sq, phi, and nu. nu is only specified if cov.model="matern". The value portion of each tag is the parameter's starting value.

tuning	a list with each tag corresponding to a parameter name. Valid tags are <code>sigma.sq</code> , <code>tau.sq</code> , <code>phi</code> , and <code>nu</code> . If <code>method="latent"</code> then only <code>phi</code> and <code>nu</code> need to be specified. The value portion of each tag defines the variance of the Metropolis sampler Normal proposal distribution.
priors	a list with each tag corresponding to a parameter name. Valid tags are <code>sigma.sq.ig</code> , <code>tau.sq.ig</code> , <code>phi.unif</code> , and <code>nu.unif</code> . Variance parameters, <code>sigma.sq</code> and <code>tau.sq</code> , are assumed to follow an inverse-Gamma distribution, whereas the spatial decay <code>phi</code> and smoothness <code>nu</code> parameters are assumed to follow Uniform distributions. The hyperparameters of the inverse-Gamma are passed as a vector of length two, with the first and second elements corresponding to the <i>shape</i> and <i>scale</i> , respectively. The hyperparameters of the Uniform are also passed as a vector of length two with the first and second elements corresponding to the lower and upper support, respectively.
cov.model	a quoted keyword that specifies the covariance function used to model the spatial dependence structure among the observations. Supported covariance model key words are: "exponential", "matern", "spherical", and "gaussian". See below for details.
n.samples	the number of posterior samples to collect.
n.omp.threads	a positive integer indicating the number of threads to use for SMP parallel processing. The package must be compiled for OpenMP support. For most Intel-based machines, we recommend setting <code>n.omp.threads</code> up to the number of hyperthreaded cores. Note, <code>n.omp.threads > 1</code> might not work on some systems.
fit.rep	if TRUE, regression fitted and replicate data will be returned. The argument <code>sub.sample</code> controls which MCMC samples are used to generate the fitted and replicated data.
sub.sample	an optional list that specifies the samples used for <code>fit.rep</code> . Valid tags are <code>start</code> , <code>end</code> , and <code>thin</code> . Given the value associated with the tags, the sample subset is selected using <code>seq(as.integer(start), as.integer(end), by=as.integer(thin))</code> . The default values are <code>start=floor(0.5*n.samples)</code> , <code>end=n.samples</code> and <code>thin=1</code> .
verbose	if TRUE, model specification and progress of the sampler is printed to the screen. Otherwise, nothing is printed to the screen.
search.type	a quoted keyword that specifies type of nearest neighbor search algorithm. Supported method key words are: "cb" and "brute". The "cb" should generally be much faster. If locations do not have identical coordinate values on the axis used for the nearest neighbor ordering (see <code>ord</code> argument) then "cb" and "brute" should produce identical neighbor sets. However, if there are identical coordinate values on the axis used for nearest neighbor ordering, then "cb" and "brute" might produce different, but equally valid, neighbor sets, e.g., if data are on a grid.
ord	an index vector of length n used for the nearest neighbor search. Internally, this vector is used to order <code>coords</code> , i.e., <code>coords[ord,]</code> , and associated data. Nearest neighbor candidates for the i -th row in the ordered <code>coords</code> are rows $1:(i-1)$, with the <code>n.neighbors</code> nearest neighbors being those with the minimum euclidean distance to the location defined by <code>ordered coords[i,]</code> . The default

used when `ord` is not specified is x-axis ordering, i.e., `order(coords[,1])`. This argument should typically be left blank. This argument will be ignored if the `neighbor.info` argument is used.

`return.neighbor.info` if TRUE, a list called `neighbor.info` containing several data structures used for fitting the NNGP model is returned. If there is no change in input data or `n.neighbors`, this list can be passed to subsequent `spNNGP` calls via the `neighbor.info` argument to avoid the neighbor search, which can be time consuming if n is large. In addition to the several cryptic data structures in `neighbor.info` there is a list called `n.indx` that is of length n . The i -th element in `n.indx` corresponds to the i -th row in `coords[ord,]` and holds the vector of that location's nearest neighbor indices. This list can be useful for plotting the neighbor graph if desired.

`neighbor.info` see the `return.neighbor.info` argument description above.

`n.report` the interval to report Metropolis sampler acceptance and MCMC progress.

... currently no additional arguments.

Details

Model parameters can be fixed at their starting values by setting their tuning values to zero.

The *no nugget* model is specified by setting `tau.sq` to zero in the starting and tuning lists.

Value

An object of class `NNGP` with additional class designations for `method` and `family`. The return object is a list comprising:

`p.beta.samples` a coda object of posterior samples for the regression coefficients.

`p.theta.samples` a coda object of posterior samples for covariance parameters.

`p.w.samples` is a matrix of posterior samples for the spatial random effects, where rows correspond to locations in `coords` and columns hold the `n.samples` posterior samples. This is only returned if `method="latent"`.

`y.hat.samples` if `fit.rep=TRUE`, regression fitted values from posterior samples specified using `sub.sample`. See additional details below.

`y.hat.quant`s if `fit.rep=TRUE`, 0.5, 0.025, and 0.975 quantiles of the `y.hat.samples`.

`y.rep.samples` if `fit.rep=TRUE`, replicated outcome from posterior samples specified using `sub.sample`. See additional details below.

`y.rep.quant`s if `fit.rep=TRUE`, 0.5, 0.025, and 0.975 quantiles of the `y.rep.samples`.

`s.indx` if `fit.rep=TRUE`, the subset index specified with `sub.sample`.

`neighbor.info` returned if `return.neighbor.info=TRUE` see the `return.neighbor.info` argument description above.

`run.time` execution time for parameter estimation reported using `proc.time()`. This time does not include nearest neighbor search time for building the neighbor set.

The return object will include additional objects used for subsequent prediction and/or model fit evaluation.

Author(s)

Andrew O. Finley <finleya@msu.edu>
 Abhirup Datta <abhidatta@jhu.edu>
 Sudipto Banerjee <sudipto@ucla.edu>

References

Datta, A., S. Banerjee, A.O. Finley, and A.E. Gelfand. (2016) Hierarchical Nearest-Neighbor Gaussian process models for large geostatistical datasets. *Journal of the American Statistical Association*, 111:800-812.

Finley, A.O., A. Datta, B.D. Cook, D.C. Morton, H.E. Andersen, and S. Banerjee. (2019) Efficient algorithms for Bayesian Nearest Neighbor Gaussian Processes. *Journal of Computational and Graphical Statistics*, <https://doi.org/10.1080/10618600.2018.1537924>.

Examples

```
rmvn <- function(n, mu=0, V = matrix(1)){
  p <- length(mu)
  if(any(is.na(match(dim(V),p))))
    stop("Dimension problem!")
  D <- chol(V)
  t(matrix(rnorm(n*p), ncol=p)%*%D + rep(mu,rep(n,p)))
}

##Make some data
set.seed(1)
n <- 100
coords <- cbind(runif(n,0,1), runif(n,0,1))

x <- cbind(1, rnorm(n))

B <- as.matrix(c(1,5))

sigma.sq <- 5
tau.sq <- 1
phi <- 3/0.5

D <- as.matrix(dist(coords))
R <- exp(-phi*D)
w <- rmvn(1, rep(0,n), sigma.sq*R)
y <- rnorm(n, x%*%B + w, sqrt(tau.sq))

##Fit a Response and Latent NNGP model
n.samples <- 500

starting <- list("phi"=phi, "sigma.sq"=5, "tau.sq"=1)

tuning <- list("phi"=0.5, "sigma.sq"=0.5, "tau.sq"=0.5)

priors <- list("phi.Unif"=c(3/1, 3/0.01), "sigma.sq.IG"=c(2, 5), "tau.sq.IG"=c(2, 1))
```

```

cov.model <- "exponential"

m.s <- spNNGP(y~x-1, coords=coords, starting=starting, method="latent", n.neighbors=10,
             tuning=tuning, priors=priors, cov.model=cov.model,
             n.samples=n.samples, n.omp.threads=1)

summary(m.s)
plot(apply(m.s$p.w.samples, 1, median), w)

m.r <- spNNGP(y~x-1, coords=coords, starting=starting, method="response", n.neighbors=10,
             tuning=tuning, priors=priors, cov.model=cov.model,
             n.samples=n.samples, n.omp.threads=1)

summary(m.r)

##Fit with some user defined neighbor ordering

##ord <- order(coords[,2]) ##y-axis
ord <- order(coords[,1]+coords[,2]) ##x+y-axis
##ord <- sample(1:n, n) ##random

m.r.xy <- spNNGP(y~x-1, coords=coords, starting=starting, method="response", n.neighbors=10,
               tuning=tuning, priors=priors, cov.model=cov.model,
               ord=ord, return.neighbor.info=TRUE,
               n.samples=n.samples, n.omp.threads=1)

summary(m.r.xy)

## Not run:
##Visualize the neighbor sets and ordering constraint
n.indx <- m.r.xy$neighbor.info$n.indx
ord <- m.r.xy$neighbor.info$ord

##This is how the data are ordered internally for model fitting
coords.ord <- coords[ord,]

for(i in 1:n){

  plot(coords.ord, cex=1, xlab="Easting", ylab="Northing")
  points(coords.ord[i,,drop=FALSE], col="blue", pch=19, cex=1)
  points(coords.ord[n.indx[[i]],,drop=FALSE], col="red", pch=19, cex=1)

  readline(prompt = "Pause. Press <Enter> to continue...")
}

## End(Not run)

```

Description

Methods for extracting information from fitted NNGP model of class NNGP and predict.NNGP objects from predict.

Usage

```
## S3 method for class 'NNGP'
summary(object, sub.sample, quantiles = c(0.025, 0.25,
0.5, 0.75, 0.975), digits = max(3L, getOption("digits") - 3L), ...)
## S3 method for class 'NNGP'
print(x, ...)
## S3 method for class 'predict.NNGP'
print(x, ...)
```

Arguments

object, x	object of class NNGP or predict.NNGP.
sub.sample	an optional list that specifies the samples to included in the summary or composition sampling. Valid tags are start, end, and thin. Given the value associated with the tags, the sample subset is selected using <code>seq(as.integer(start), as.integer(end), by=as.i</code> . The default values are <code>start=floor(0.5*n.samples)</code> , <code>end=n.samples</code> and <code>thin=1</code> .
quantiles	for summary, posterior distribution quantiles to compute.
digits	for summary, number of digits to report in summary.
...	currently not used.

Details

A set of standard extractor functions for fitted model objects of class NNGP and prediction object of class predict.NNGP, including methods to the generic functions `print` and `summary`.

summary.PGLogit

Methods for PGLogit Object

Description

Methods for extracting information from fitted PGLogit model.

Usage

```
## S3 method for class 'PGLogit'
summary(object, sub.sample, quantiles
=c(0.025, 0.25, 0.5, 0.75, 0.975), digits = max(3L, getOption("digits")
- 3L), ...)
## S3 method for class 'PGLogit'
print(x, ...)
```

Arguments

object, x	object of class PGLoFit.
sub.sample	an optional list that specifies the samples to included in the summary or composition sampling. Valid tags are start, end, and thin. Given the value associated with the tags, the sample subset is selected using <code>seq(as.integer(start), as.integer(end), by=as.i</code> . The default values are <code>start=floor(0.5*n.samples)</code> , <code>end=n.samples</code> and <code>thin=1</code> .
quantiles	for summary, posterior distribution quantiles to compute.
digits	for summary, number of digits to report.
...	currently not used.

Details

A set of standard extractor functions for fitted model objects of class PGLoFit, including methods to the generic functions `print` and `summary`.

Index

*Topic **datasets**

BCEF, [2](#)

MI_TSCA, [3](#)

*Topic **model**

PGLogit, [3](#)

predict.NNGP, [5](#)

print.spDiag, [8](#)

spConjNNGP, [9](#)

spNNGP, [15](#)

summary.NNGP, [20](#)

summary.PGLogit, [20](#)

*Topic **utilities**

spDiag, [13](#)

BCEF, [2](#)

fitted.NNGP (summary.NNGP), [20](#)

fitted.PGLogit (summary.PGLogit), [20](#)

MI_TSCA, [3](#)

PGLogit, [3](#)

predict.NNGP, [5](#)

print, [20](#), [21](#)

print.NNGP (summary.NNGP), [20](#)

print.PGLogit (summary.PGLogit), [20](#)

print.predict.NNGP (summary.NNGP), [20](#)

print.spDiag, [8](#)

residuals.NNGP (summary.NNGP), [20](#)

residuals.PGLogit (summary.PGLogit), [20](#)

spConjNNGP, [9](#)

spDiag, [13](#)

spNNGP, [15](#)

summary, [20](#), [21](#)

summary.NNGP, [19](#)

summary.PGLogit, [20](#)