

Package ‘gamsel’

February 4, 2022

Type Package

Title Fit Regularization Path for Generalized Additive Models

Version 1.8-2

Date 2022-02-03

Description Using overlap grouped-lasso penalties, 'gamsel' selects whether a term in a 'gam' is nonzero, linear, or a non-linear spline (up to a specified max df per variable). It fits the entire regularization path on a grid of values for the overall penalty lambda, both for gaussian and binomial families.

License GPL-2

Imports foreach, mda, splines

URL <https://arxiv.org/abs/1506.03850>

Encoding UTF-8

RoxygenNote 7.1.2

Depends R (>= 3.6)

NeedsCompilation yes

Author Alexandra Chouldechova [aut],
Trevor Hastie [aut, cre],
Balasubramanian Narasimhan [ctb],
Vitalie Spinu [ctb]

Maintainer Trevor Hastie <hastie@stanford.edu>

Repository CRAN

Date/Publication 2022-02-04 13:50:02 UTC

R topics documented:

gamsel-package	2
basis.gen	2
cv.gamsel	3
gamsel	5
getActive	8

plot.cv.gamsel	9
plot.gamsel	10
predict.gamsel	11
print.gamsel	12
pseudo.bases	13
summary.gamsel	14

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

gamsel-package	<i>Fit Regularization Path for Generalized Additive Models</i>
----------------	--

Description

Using overlap grouped lasso penalties, `gamsel` selects whether a term in a gam is nonzero, linear, or a non-linear spline (up to a specified max df per variable). It fits the entire regularization path on a grid of values for the overall penalty lambda, both for gaussian and binomial families. Key functions are [gamsel](#), [predict.gamsel](#), [plot.gamsel](#), [print.gamsel](#), [summary.gamsel](#), [cv.gamsel](#), [plot.cv.gamsel](#)

Author(s)

Alexandra Chouldechova, Trevor Hastie Maintainer: Trevor Hastie <hastie@stanford.edu>

basis.gen	<i>Generate basis</i>
-----------	-----------------------

Description

Generate basis

Usage

```
basis.gen(x, df = 6, thresh = 0.01, degree = 8, parms = NULL, ...)
```

Arguments

<code>x</code>	A vector of values for <code>basis.gen</code> , or a matrix for <code>pseudo.bases</code>
<code>df</code>	The degrees of freedom of the smoothing spline.
<code>thresh</code>	If the next eigenvector improves the approximation by less than threshold, a truncated bases is returned. For <code>pseudo.bases</code> this can be a single value or a vector of values, which are recycled sequentially for each column of <code>x</code>
<code>degree</code>	The nominal number of basis elements. The basis returned has no more than <code>degree</code> columns. For <code>pseudo.bases</code> this can be a single value or a vector of values, which are recycled sequentially for each column of <code>x</code>
<code>parms</code>	A parameter set. If included in the call, these are used to define the basis. This is used for prediction.
<code>...</code>	other arguments

Value

the basis

 cv.gamsel

Cross-validation Routine for Gamsel

Description

A routine for performing K-fold cross-validation for gamsel.

Usage

```
cv.gamsel(
  x,
  y,
  lambda = NULL,
  family = c("gaussian", "binomial"),
  degrees = rep(10, p),
  dfs = rep(5, p),
  bases = pseudo.bases(x, degrees, dfs, parallel = parallel, ...),
  type.measure = c("mse", "mae", "deviance", "class"),
  nfolds = 10,
  foldid,
  keep = FALSE,
  parallel = FALSE,
  ...
)
```

Arguments

x	x matrix as in gamsel
y	response y as in gamsel
lambda	Optional user-supplied lambda sequence. If NULL, default behaviour is for gamsel routine to automatically select a good lambda sequence.
family	family as in gamsel
degrees	degrees as in gamsel
dfs	dfs as in gamsel
bases	bases as in gamsel
type.measure	Loss function for cross-validated error calculation. Currently there are four options: mse (mean squared error), mae (mean absolute error), deviance (deviance, same as mse for family="gaussian"), class (misclassification error, for use with family="binomial").
nfolds	Numer of folds (default is 10). Maximum value is nobs. Small values of nfolds are recommended for large data sets.

foldid	Optional vector of length nobobs with values between 1 and nfolde specifying what fold each observation is in.
keep	If keep=TRUE, a <i>prevalidated</i> array is returned containing fitted values for each observation and each value of lambda. This means these fits are computed with this observation and the rest of its fold omitted. The folid vector is also returned. Default is keep=FALSE
parallel	If TRUE, use parallel foreach to fit each fold. See the example below for usage details.
...	Other arguments that can be passed to gamsel.

Details

This function has the effect of running gamsel nfolde+1 times. The initial run uses all the data and gets the lambda sequence. The remaining runs fit the data with each of the folds omitted in turn. The error is accumulated, and the average error and standard deviation over the folds is computed. Note that cv.gamsel does NOT search for values for gamma. A specific value should be supplied, else gamma=.4 is assumed by default. If users would like to cross-validate gamma as well, they should call cv.gamsel with a pre-computed vector foldid, and then use this same fold vector in separate calls to cv.gamsel with different values of gamma. Note also that the results of cv.gamsel are random, since the folds are selected at random. Users can reduce this randomness by running cv.gamsel many times, and averaging the error curves.

Value

an object of class "cv.gamsel" is returned, which is a list with the ingredients of the cross-validation fit.

lambda	the values of lambda used in the fits.
cvm	The mean cross-validated error - a vector of length length(lambda).
cvsd	estimate of standard error of cvm.
cvup	upper curve = cvm+cvsd.
cvlo	lower curve = cvm-cvsd.
nzero	number of non-zero coefficients at each lambda.
name	a text string indicating type of measure (for plotting purposes).
gamsel.fit	a fitted gamsel object for the full data.
lambda.min	value of lambda that gives minimum cvm.
lambda.1se	largest value of lambda such that error is within 1 standard error of the minimum.
fit.preval	if keep=TRUE, this is the array of prevalidated fits. Some entries can be NA, if that and subsequent values of lambda are not reached for that fold
foldid	if keep=TRUE, the fold assignments used
index.min	the sequence number of the minimum lambda.
index.1se	the sequence number of the 1se lambda value.

Author(s)

Alexandra Chouldechova and Trevor Hastie
Maintainer: Trevor Hastie <hastie@stanford.edu>

References

Chouldechova, A. and Hastie, T. (2015) *Generalized Additive Model Selection*

See Also

gamsel, plot function for cv.gamsel object.

Examples

```
##data=gamsel:::gendata(n=500,p=12,k.lin=3,k.nonlin=3,deg=8,sigma=0.5)
data = readRDS(system.file("extdata/gamsel_example.RDS", package = "gamsel"))
attach(data)
bases=pseudo.bases(X,degree=10,df=6)
# Gaussian gam
gamsel.out=gamsel(X,y,bases=bases)
par(mfrow=c(1,2),mar=c(5,4,3,1))
summary(gamsel.out)
gamsel.cv=cv.gamsel(X,y,bases=bases)
par(mfrow=c(1,1))
plot(gamsel.cv)
par(mfrow=c(3,4))
plot(gamsel.out,newx=X,index=20)
```

gamsel

Fit Regularization Path for Gaussian or Binomial Generalized Additive Model

Description

Using overlap grouped lasso penalties, gamsel selects whether a term in a gam is nonzero, linear, or a non-linear spline (up to a specified max df per variable). It fits the entire regularization path on a grid of values for the overall penalty lambda, both for gaussian and binomial families.

Usage

```
gamsel(  
  x,  
  y,  
  num_lambda = 50,  
  lambda = NULL,  
  family = c("gaussian", "binomial"),
```

```

degrees = rep(10, p),
gamma = 0.4,
dfs = rep(5, p),
bases = pseudo.bases(x, degrees, dfs, parallel = parallel, ...),
tol = 1e-04,
max_iter = 2000,
traceit = FALSE,
parallel = FALSE,
...
)

```

Arguments

<code>x</code>	Input (predictor) matrix of dimension <code>nobs</code> x <code>nvars</code> . Each observation is a row.
<code>y</code>	Response variable. Quantitative for <code>family="gaussian"</code> and with values in $\{0, 1\}$ for <code>family="binomial"</code>
<code>num_lambda</code>	Number of <code>lambda</code> values to use. (Length of <code>lambda</code> sequence.)
<code>lambda</code>	User-supplied <code>lambda</code> sequence. For best performance, leave as <code>NULL</code> and allow the routine to automatically select <code>lambda</code> . Otherwise, supply a (preferably gradually) decreasing sequence.
<code>family</code>	Response type. "gaussian" for linear model (default). "binomial" for logistic model.
<code>degrees</code>	An integer vector of length <code>nvars</code> specifying the maximum number of spline basis functions to use for each variable.
<code>gamma</code>	Penalty mixing parameter $0 \leq \gamma \leq 1$. Values $\gamma < 0.5$ penalize linear fit less than non-linear fit. The default is $\gamma = 0.4$, which encourages a linear term over a nonlinear term.
<code>dfs</code>	Numeric vector of length <code>nvars</code> specifying the maximum (end-of-path) degrees of freedom for each variable.
<code>bases</code>	A list of orthonormal bases for the non-linear terms for each variable. The function <code>pseudo.bases</code> generates these, using the parameters <code>dfs</code> and <code>degrees</code> . See the documentation for pseudo.bases .
<code>tol</code>	Convergence threshold for coordinate descent. The coordinate descent loop continues until the total change in objective after a pass over all variables is less than <code>tol</code> . Default is $1e-4$.
<code>max_iter</code>	Maximum number of coordinate descent iterations over all the variables for each <code>lambda</code> value. Default is 2000.
<code>traceit</code>	If <code>TRUE</code> , various information is printed during the fitting process.
<code>parallel</code>	passed on to the <code>pseudo.bases()</code> function. Uses multiple process if available.
<code>...</code>	additional arguments passed on to <code>pseudo.bases()</code>

Details

The sequence of models along the `lambda` path is fit by (block) coordinate descent. In the case of logistic regression the fitting routine may terminate before all `num_lambda` values of `lambda` have

been used. This occurs when the fraction of null deviance explained by the model gets too close to 1, at which point the fit becomes numerically unstable. Each of the smooth terms is computed using an approximation to the Demmler-Reinsch smoothing spline basis for that variable, and the accompanying diagonal penalty matrix.

Value

An object with S3 class `gamsel`. %% If it is a LIST, use

<code>intercept</code>	Intercept sequence of length <code>num_lambda</code>
<code>alphas</code>	<code>nvars</code> x <code>num_lambda</code> matrix of linear coefficient estimates
<code>betas</code>	<code>sum(degrees)</code> x <code>num_lambda</code> matrix of non-linear coefficient estimates
<code>lambdas</code>	The sequence of lambda values used
<code>degrees</code>	Number of basis functions used for each variable
<code>parms</code>	A set of parameters that capture the bases used. This allows for efficient generation of the bases elements for <code>predict.gamsel</code>

, the `predict` method for this class.

<code>family</code>	"gaussian" or "binomial"
<code>nulldev</code>	Null deviance (deviance of the intercept model)
<code>dev.ratio</code>	Vector of length <code>num_lambda</code> giving fraction of (null) deviance explained by each model along the lambda sequence
<code>call</code>	The call that produced this object
%% ...	

Author(s)

Alexandra Chouldechova and Trevor Hastie
 Maintainer: Trevor Hastie <hastie@stanford.edu>

References

Chouldechova, A. and Hastie, T. (2015) *Generalized Additive Model Selection*, <https://arxiv.org/abs/1506.03850>

See Also

[predict.gamsel](#), [cv.gamsel](#), [plot.gamsel](#), [summary.gamsel](#), [basis.gen](#),

Examples

```
##data=gamsel:::gendata(n=500,p=12,k.lin=3,k.nonlin=3,deg=8,sigma=0.5)
data = readRDS(system.file("extdata/gamsel_example.RDS", package = "gamsel"))
attach(data)
bases=pseudo.bases(X,degree=10,df=6)
# Gaussian gam
```

```

gamsel.out=gamsel(X,y,bases=bases)
par(mfrow=c(1,2),mar=c(5,4,3,1))
summary(gamsel.out)
gamsel.cv=cv.gamsel(X,y,bases=bases)
par(mfrow=c(1,1))
plot(gamsel.cv)
par(mfrow=c(3,4))
plot(gamsel.out,newx=X,index=20)
# Binomial model
gamsel.out=gamsel(X,yb,family="binomial")
par(mfrow=c(1,2),mar=c(5,4,3,1))
summary(gamsel.out)
par(mfrow=c(3,4))
plot(gamsel.out,newx=X,index=30)

```

getActive	<i>Returns active variables</i>
-----------	---------------------------------

Description

Extract active variables of different kinds from a gamsel object

Usage

```

getActive(
  object,
  index = NULL,
  type = c("nonzero", "linear", "nonlinear"),
  EPS = 0
)

```

Arguments

object	gamsel object
index	index or vector of indices at which to obtain active information. NULL returns all.
type	type of active variables to report. One of c("nonzero", "linear", "nonlinear")
EPS	threshold for what is nonzero; default is 0

Details

Returns a vector of variables indices of variables having the desired properties.

Value

vector of indices

`plot.cv.gamsel`*Plotting Routine for Gamsel Cross-Validation Object*

Description

Produces a cross-validation curve with standard errors for a fitted gamsel objecty.

Usage

```
## S3 method for class 'cv.gamsel'  
plot(x, sign.lambda = 1, ...)
```

Arguments

<code>x</code>	cv.gamsel object
<code>sign.lambda</code>	Either plot against $\log(\lambda)$ (default) against $-\lambda$ if <code>sign.lambda=-1</code> .
<code>...</code>	Optional graphical parameters to plot.

Details

A plot showing cross-validation error is produced. Nothing is returned.

Author(s)

Alexandra Chouldechova and Trevor Hastie
Maintainer: Trevor Hastie <hastie@stanford.edu>

References

Chouldechova, A. and Hastie, T. (2015) *Generalized Additive Model Selection*

Examples

```
##data=gamsel:::gendata(n=500,p=12,k.lin=3,k.nonlin=3,deg=8,sigma=0.5)  
data = readRDS(system.file("extdata/gamsel_example.RDS", package = "gamsel"))  
attach(data)  
bases=pseudo.bases(X,degree=10,df=6)  
# Gaussian gam  
gamsel.out=gamsel(X,y,bases=bases)  
gamsel.cv=cv.gamsel(X,y,bases=bases)  
par(mfrow=c(1,1))  
plot(gamsel.cv)
```

plot.gamsel

Plotting Routine gamsel Object

Description

Produces plots of the estimated functions for specified variables at a given value of lambda.

Usage

```
## S3 method for class 'gamsel'
plot(x, newx, index, which = 1:p, rugplot = TRUE, ylims, ...)
```

Arguments

x	Fitted gamsel object.
newx	nobs_new x p matrix giving values of each predictor at which to plot.
index	Index of lambda value (i.e., model) for which plotting is desired.
which	Which values to plot. Default is all variables, i.e. {1,2,...,nvars}. Besides indices, which can take two special values: "nonzero" will plot only the nonzero functions, and "nonlinear" only the nonlinear functions.
rugplot	If TRUE, a rugplot showing values of x is shown at the bottom of each fitted function plot.
ylims	ylim argument for plotting each curve, which overrides the default which is the range of all the functions.
...	Optional graphical parameters to plot.

Details

A plot of the specified fitted functions is produced. Nothing is returned.

Author(s)

Alexandra Chouldechova and Trevor Hastie
 Maintainer: Trevor Hastie <hastie@stanford.edu>

References

Chouldechova, A. and Hastie, T. (2015) *Generalized Additive Model Selection*

See Also

gamsel, and print.gamsel, summary.gamsel

Examples

```
##set.seed(1211)
##data=gamсел::gendata(n=500,p=12,k.lin=3,k.nonlin=3,deg=8,sigma=0.5)
data = readRDS(system.file("extdata/gamсел_example.RDS", package = "gamсел"))
attach(data)
bases=pseudo.bases(X,degree=10,df=6)
# Gaussian gam
gamсел.out=gamсел(X,y,bases=bases)
par(mfrow=c(3,4))
plot(gamсел.out,newx=X,index=20)
```

predict.gamсел

Gamсел Prediction Routine

Description

Make predictions from a gamсел object.

Usage

```
## S3 method for class 'gamсел'
predict(
  object,
  newdata,
  index = NULL,
  type = c("link", "response", "terms", "nonzero"),
  ...
)
```

Arguments

object	Fitted gamсел object.
newdata	nobs_new x p matrix of new data values at which to predict.
index	Index of model in the sequence for which plotting is desired. Note, this is NOT a lambda value.
type	Type of prediction desired. Type link gives the linear predictors for "binomial", and fitted values for "gaussian". Type response gives fitted probabilities for "binomial" and fitted values for "gaussian". Type "terms" returns a matrix of fitted functions, with as many columns as there are variables. Type nonzero returns a list of the indices of nonzero coefficients at the given lambda index.
...	Not used

Value

Either a vector aor a matrix is returned, depending on type.

Author(s)

Alexandra Chouldechova and Trevor Hastie
 Maintainer: Trevor Hastie <hastie@stanford.edu>

References

Chouldechova, A. and Hastie, T. (2015) *Generalized Additive Model Selection*

See Also

[gamsel](#), [cv.gamsel](#), [summary.gamsel](#), [basis.gen](#)

Examples

```
##data=gamsel:::gendata(n=500,p=12,k.lin=3,k.nonlin=3,deg=8,sigma=0.5)
data = readRDS(system.file("extdata/gamsel_example.RDS", package = "gamsel"))
attach(data)
bases=pseudo.bases(X,degree=10,df=6)
# Gaussian gam
gamsel.out=gamsel(X,y,bases=bases)
preds=predict(gamsel.out,X,index=20,type="terms")
```

print.gamsel	<i>print a gamsel object</i>
--------------	------------------------------

Description

Print a summary of the gamsel path at each step along the path

Usage

```
## S3 method for class 'gamsel'
print(x, digits = max(3, getOption("digits") - 3), ...)
```

Arguments

x	fitted gamsel object
digits	significant digits in printout
...	additional print arguments

Details

The call that produced the object x is printed, followed by a five-column matrix with columns NonZero, Lin, NonLin, %Dev and Lambda. The first three columns say how many nonzero, linear and nonlinear terms there are. %Dev is the percent deviance explained (relative to the null deviance).

Value

The matrix above is silently returned

Author(s)

Alexandra Chouldechova and Trevor Hastie
 Maintainer: Trevor Hastie <hastie@stanford.edu>

References

Chouldechova, A. and Hastie, T. (2015) *Generalized Additive Model Selection*

See Also

[predict.gamsel](#), [cv.gamsel](#), [plot.gamsel](#), [summary.gamsel](#), [basis.gen](#)

pseudo.bases	<i>Generate pseudo-spline bases</i>
--------------	-------------------------------------

Description

Generate an approximation to the Demmler-Reinsch orthonormal bases for smoothing splines, using orthogonal polynomials. `basis.gen` generates a basis for a single `x`, and `pseudo.bases` generates a list of bases for each column of the matrix `x`.

Usage

```
pseudo.bases(x, degree = 8, df = 6, parallel = FALSE, ...)
```

Arguments

<code>x</code>	A vector of values for <code>basis.gen</code> , or a matrix for <code>pseudo.bases</code>
<code>degree</code>	The nominal number of basis elements. The basis returned has no more than <code>degree</code> columns. For <code>pseudo.bases</code> this can be a single value or a vector of values, which are recycled sequentially for each column of <code>x</code>
<code>df</code>	The degrees of freedom of the smoothing spline.
<code>parallel</code>	if TRUE, parallelize
<code>...</code>	other arguments for <code>basis.gen</code> can be passed through to <code>[basis.gen]</code>

Details

`basis.gen` starts with a basis of orthogonal polynomials of total degree. These are each smoothed using a smoothing spline, which allows for a one-step approximation to the Demmler-Reinsch basis for a smoothing spline of rank equal to the degree. See the reference for details. The function also approximates the appropriate diagonal penalty matrix for this basis, so that the a approximate smoothing spline (generalized ridge regression) has the target `df`.

Value

An orthonormal basis is returned (a list for `pseudo.bases`). This has an attribute `parms`, which has elements `coefs` Coefficients needed to generate the orthogonal polynomials `rotate` Transformation matrix for transforming the polynomial basis `dp` penalty values for the diagonal penalty `dfdf` used `deg` number of columns

Author(s)

Alexandra Chouldechova and Trevor Hastie
Maintainer: Trevor Hastie <hastie@stanford.edu>

References

T. Hastie *Pseudosplines*. (1996) JRSSB 58(2), 379-396.
Chouldechova, A. and Hastie, T. (2015) *Generalized Additive Model Selection*

Examples

```
##data=gamsel:::gendata(n=500,p=12,k.lin=3,k.nonlin=3,deg=8,sigma=0.5)
data = readRDS(system.file("extdata/gamsel_example.RDS", package = "gamsel"))
attach(data)
bases=pseudo.bases(X,degree=10,df=6)
## Not run:
  require(doMC)
  registerDoMC(cores=4)
  bases=pseudo.bases(X,degree=10,df=6,parallel=TRUE)

## End(Not run)
```

summary.gamsel

Gamsel summary routine

Description

This makes a two-panel plot of the `gamsel` object.

Usage

```
## S3 method for class 'gamsel'
summary(object, label = FALSE, ...)
```

Arguments

<code>object</code>	<code>gamsel</code> object
<code>label</code>	if TRUE, annotate the plot with variable labels. Default is FALSE
<code>...</code>	additional arguments to <code>summary</code>

Details

A two panel plot is produced, that summarizes the linear components and the nonlinear components, as a function of lambda. For the linear components, it is the coefficient for each variable. For the nonlinear, we see the norm of the nonlinear coefficients.

Value

Nothing is returned.

Author(s)

Alexandra Chouldechova and Trevor Hastie
Maintainer: Trevor Hastie <hastie@stanford.edu>

References

Chouldechova, A. and Hastie, T. (2015) *Generalized Additive Model Selection*

See Also

gamselect, and methods plot, print and predict for cv.gamselect object.

Examples

```
##data=gamselect::gendata(n=500,p=12,k.lin=3,k.nonlin=3,deg=8,sigma=0.5)
data = readRDS(system.file("extdata/gamselect_example.RDS", package = "gamselect"))
attach(data)
bases=pseudo.bases(X,degree=10,df=6)
# Gaussian gam
gamselect.out=gamselect(X,y,bases=bases)
par(mfrow=c(1,2),mar=c(5,4,3,1))
summary(gamselect.out)
```

Index

- * **additive**
 - [gamsel-package](#), 2
 - * **generalized**
 - [gamsel-package](#), 2
 - * **models**
 - [gamsel-package](#), 2
 - * **nonparametric**
 - [cv.gamsel](#), 3
 - [gamsel](#), 5
 - [plot.cv.gamsel](#), 9
 - [plot.gamsel](#), 10
 - [predict.gamsel](#), 11
 - [print.gamsel](#), 12
 - [pseudo.bases](#), 13
 - [summary.gamsel](#), 14
 - * **package**
 - [gamsel-package](#), 2
 - * **regression**
 - [cv.gamsel](#), 3
 - [gamsel](#), 5
 - [gamsel-package](#), 2
 - [plot.cv.gamsel](#), 9
 - [plot.gamsel](#), 10
 - [predict.gamsel](#), 11
 - [print.gamsel](#), 12
 - [pseudo.bases](#), 13
 - [summary.gamsel](#), 14
 - * **smooth**
 - [cv.gamsel](#), 3
 - [gamsel](#), 5
 - [plot.cv.gamsel](#), 9
 - [plot.gamsel](#), 10
 - [predict.gamsel](#), 11
 - [print.gamsel](#), 12
 - [pseudo.bases](#), 13
 - [summary.gamsel](#), 14
- [basis.gen](#), 2, 7, 12, 13
- [cv.gamsel](#), 2, 3, 7, 12, 13
- [gamsel](#), 2, 5, 12
- [gamsel-package](#), 2
- [getActive](#), 8
- [plot.cv.gamsel](#), 2, 9
- [plot.gamsel](#), 2, 7, 10, 13
- [predict.gamsel](#), 2, 7, 11, 13
- [print.gamsel](#), 2, 12
- [pseudo.bases](#), 6, 13
- [summary.gamsel](#), 2, 7, 12, 13, 14