

Package ‘bWGR’

January 8, 2020

Type Package

Title Bayesian Whole-Genome Regression

Version 1.6.6

Date 2020-01-07

Author Alencar Xavier, William Muir, Kyle Kocak, Shizhong Xu, Katy Rainey.

Maintainer Alencar Xavier <alencxav@gmail.com>

Description Whole-genome regression methods on Bayesian framework fitted via EM or Gibbs sampling, univariate and multivariate, with optional kernel term and sampling techniques.

License GPL-3

Imports Rcpp

LinkingTo Rcpp

Depends R (>= 3.2.0)

NeedsCompilation yes

Repository CRAN

Date/Publication 2020-01-08 11:50:02 UTC

R topics documented:

| | |
|------------------------|-----------|
| bWGR-package | 2 |
| Dataset | 2 |
| WGR1 (MC) | 3 |
| WGR2 (EM) | 6 |
| WGR3 (MV) | 8 |
| XTRA 1 | 9 |
| XTRA 2 | 11 |
| Index | 12 |

bWGR-package

Bayesian Whole-Genome Regression

Description

Whole-genome regression methods on Bayesian framework fitted via EM or Gibbs sampling, univariate and multivariate, with optional kernel term and sampling techniques.

Details

Package: bWGR
Type: Package
Version: 1.6.6
Date: 2020-01-07
License: GPL-3

Author(s)

Alencar Xavier, William Muir, Kyle Kocak, Shizhong Xu, Katy Rainey.

Maintainer: Alencar Xavier <alencxav@gmail.com>

Examples

```
## Not run:  
data(tpod)  
Fit = wgr(y,gen)  
cor(y,Fit$hat)  
  
## End(Not run)
```

Dataset

Tetra-seed Pods

Description

Two biparental crosses phenotyped for the percentage of pods containing four seeds

Usage

```
data(tpod)
```

Details

Soybean nested association panel with 2 families (*fam*) containing 196 individuals. Genotypic matrix (*gen*) have 376 SNP across 20 chromosome (*chr*). Phenotypic information (*y*) regards the proportion of tetra-seed pods. Data provided by Rainey Lab for Soybean Breeding and Genetics, Purdue University.

Author(s)

Alencar Xavier and Katy Rainey

WGR1 (MC)

MCMC Whole-genome Regression

Description

Univariate model to find breeding values through regression with optional resampling techniques (Xavier et al. 2017) and polygenic term (Kernel). See "Details" for additional standalone functions written in C++.

Usage

```
wgr(y, X, it=1500, bi=500, th=1, bag=1, rp=FALSE, iv=FALSE, de=FALSE,
    pi=0, df=5, R2=0.5, eigK=NULL, VarK=0.95, verb=FALSE)
```

Arguments

| | |
|------------|---|
| <i>y</i> | Numeric vector of observations (<i>n</i>) describing the trait to be analyzed. NA is allowed. |
| <i>X</i> | Numeric matrix containing the genotypic data. A matrix with <i>n</i> rows of observations and (<i>m</i>) columns of molecular markers. |
| <i>it</i> | Integer. Number of iterations or samples to be generated. |
| <i>bi</i> | Integer. Burn-in, the number of iterations or samples to be discarded. |
| <i>th</i> | Integer. Thinning parameter, used to save memory by storing only one every 'th' samples. |
| <i>bag</i> | If different than one, it indicates the proportion of data to be subsampled in each Markov chain. For datasets with moderate number of observations, values of bag from 0.30 to 0.60 may speed up computation without losses in prediction properties. This argument enable users to enhance MCMC through subsampling (Xavier et al. 2017). |
| <i>rp</i> | Logical. Use replacement for bootstrap samples when bag is different than one. |
| <i>iv</i> | Logical. Assign markers independent variance, a T prior from a mixture of normals. If true, turns the default model BLUP into BayesA. |
| <i>de</i> | Logical. Assign markers independent variance through double-exponential prior. If true, turns the default model BLUP into Bayesian LASSO. This argument overrides iv. |

| | |
|------|--|
| pi | Value between 0 and 1. If greater than zero it activates variable selection, where markers have expected probability pi of having null effect. |
| df | Prior degrees of freedom of variance components. |
| R2 | Expected R2, used to calculate the prior shape. |
| eigK | Output of function eigen. Spectral decomposition of the kernel used as a second random effect (eg. pedigree matrix). |
| VarK | Numeric between 0 and 1. For reduction of dimensionality. Indicates the proportion of variance explained by Eigenpairs used to fit second random effect. |
| verb | Logical. If verbose is TRUE, function displays MCMC progress bar. |

Details

The model for the whole-genome regression is as follows:

$$y = \mu + Xb + u + e$$

where y is the response variable, μ is the intercept, X is the genotypic matrix, b is the regression coefficient or effect of an allele substitution, with d probability of being included into the model, u is the polygenic term if a kernel is used, and e is the residual term.

Users can obtain four WGR methods out of this function: BRR (pi=0,iv=F), BayesA (pi=0,iv=T), BayesB (pi=0.95,iv=T), BayesC (pi=0.95,iv=F) and Bayesian LASSO or BayesL (pi=0,de=T). Theoretical basis of each model is described by de los Campos et al. (2013).

Gibbs sampler that updates regression coefficients is adapted from GSRU algorithm (Legarra and Misztal 2008). The variable selection of functions *wgr*, *BayesB* and *BayesC* works through the unconditional prior algorithm proposed by Kuo and Mallick (1998), whereas *BayesCpi* and *BayesDpi* are based on Metropolis-Hastings. Prior shape estimates are computed as $S_b = R^2 * df * var(y) / MS_x$ and $S_e = (1 - R^2) * df * var(y)$, with an exception for *BayesC* and *BayesCpi* where the prior shape is $S_b = R^2 * df * var(y) / MS_x / (1 - pi)$. The polygenic term is solved by Bayesian algorithm of reproducing kernel Hilbert Spaces proposed by de los Campos et al. (2010).

In addition to *wgr*, standalone C++ functions available include:

- 01) BayesA(y, X, it=1500, bi=500, df=5, R2=0.5)
- 02) BayesB(y, X, it=1500, bi=500, pi=0.95, df=5, R2=0.5)
- 03) BayesC(y, X, it=1500, bi=500, pi=0.95, df=5, R2=0.5)
- 04) BayesCpi(y, X, it=1500, bi=500, df=5, R2=0.5)
- 05) BayesDpi(y, X, it=1500, bi=500, df=5, R2=0.5)
- 06) BayesL(y, X, it=1500, bi=500, df=5, R2=0.5)
- 07) BayesRR(y, X, it=1500, bi=500, df=5, R2=0.5)

The implementations that support two random effects include:

- 08) BayesA2(y, X1, X2, it=1500, bi=500, df=5, R2=0.5)
- 09) BayesB2(y, X1, X2, it=1500, bi=500, pi=0.95, df=5, R2=0.5)
- 10) BayesRR2(y, X1, X2, it=1500, bi=500, df=5, R2=0.5)

And the cross-validation for the C++ implementations, with arguments analogous to *emCV*.

mcmcCV(y, gen, k=5, n=5, it=1500, bi=500, pi=0.95, df=5, R2=0.5, avg=T, llo=NULL, tbv=NULL)

Value

The function `wgr` returns a list with expected value from the marker effect (b), probability of marker being in the model (d), regression coefficient (g), variance of each marker (Vb), the intercept (mu), the polygene (u) and polygenic variance (Vk), residual variance (Ve) and the fitted value (hat).

Author(s)

Alencar Xavier

References

de los Campos, G., Hickey, J. M., Pong-Wong, R., Daetwyler, H. D., and Calus, M. P. (2013). Whole-genome regression and prediction methods applied to plant and animal breeding. *Genetics*, 193(2), 327-345.

de los Campos, G., Gianola, D., Rosa, G. J., Weigel, K. A., & Crossa, J. (2010). Semi-parametric genomic-enabled prediction of genetic values using reproducing kernel Hilbert spaces methods. *Genetics Research*, 92(04), 295-308.

Kuo, L., & Mallick, B. (1998). Variable selection for regression models. *Sankhya: The Indian Journal of Statistics, Series B*, 65-81.

Legarra, A., & Misztal, I. (2008). Technical note: Computing strategies in genome-wide selection. *Journal of dairy science*, 91(1), 360-366.

Xavier, A., Xu, S., Muir, W., & Rainey, K. M. (2017). Genomic prediction using subsampling. *BMC bioinformatics*, 18(1), 191.

Examples

```
## Not run:

# Load data
data(tpod)

# BLUP
fit_BRR = wgr(y,gen,iv=FALSE,pi=0)
cor(y,fit_BRR$hat)

# BayesA
fit_BayesA = wgr(y,gen,iv=TRUE,pi=0)
cor(y,fit_BayesA$hat)

# BayesB
fit_BayesB = wgr(y,gen,iv=TRUE,pi=.95)
cor(y,fit_BayesB$hat)

# BayesC
fit_BayesC = wgr(y,gen,iv=FALSE,pi=.95)
cor(y,fit_BayesC$hat)

# BayesCpi
fit_BayesCpi = BayesCpi(y,gen)
```

```

cor(y, fit_BayesCpi$hat)

# BayesDpi
fit_BayesDpi = BayesDpi(y, gen)
cor(y, fit_BayesDpi$hat)

# BayesL
fit_BayesL = wgr(y, gen, de=TRUE)
cor(y, fit_BayesL$hat)

# Bagging BLUP
fit_Bag = wgr(y, gen, bag=0.5)
cor(y, fit_Bag$hat)

## End(Not run)

```

WGR2 (EM)

Expectation-Maximization WGR

Description

Univariate models to find breeding values through regression fitted via expectation-maximization implemented in C++.

Usage

```

emRR(y, gen, df = 10, R2 = 0.5)
emBA(y, gen, df = 10, R2 = 0.5)
emBB(y, gen, df = 10, R2 = 0.5, Pi = 0.75)
emBC(y, gen, df = 10, R2 = 0.5, Pi = 0.75)
emBL(y, gen, R2 = 0.5, alpha = 0.02)
emEN(y, gen, R2 = 0.5, alpha = 0.02)
emDE(y, gen, R2 = 0.5)
emML(y, gen, D = NULL)
emCV(y, gen, k = 5, n = 5, Pi = 0.75, alpha = 0.02,
      df = 10, R2 = 0.5, avg=TRUE, llo=NULL, tbv=NULL)

```

Arguments

| | |
|-----|--|
| y | Numeric vector of response variable (n). NA is not allowed. |
| gen | Numeric matrix containing the genotypic data. A matrix with n rows of observations and m columns of molecular markers. |
| df | Hyperprior degrees of freedom of variance components. |
| R2 | Expected R2, used to calculate the prior shape (de los Campos et al. 2013). |
| Pi | Value between 0 and 1. Expected probability pi of having null effect (or $1-\text{Pi}$ if $\text{Pi}>0.5$). |

| | |
|-------|--|
| alpha | Value between 0 and 1. Intensity of L1 variable selection. |
| D | NULL or numeric vector with length p. Vector of weights for markers. |
| k | Integer. Folding of a k-fold cross-validation. |
| n | Integer. Number of cross-validation to perform. |
| avg | Logical. Return average across CV, or correlations within CV. |
| llo | NULL or a vector (numeric or factor) with the same length as y. If provided, the cross-validations are performed as Leave a Level Out (LLO). This argument allows the user to predefine the splits. This argument overrides k and n. |
| tbv | NULL or numeric vector of 'true breeding values' (n) to use to compare cross-validations to. If NULL, the cross-validations will have the phenotypes as prediction target. |

Details

The model for the whole-genome regression is as follows:

$$y = \mu + Xb + e$$

where y is the response variable, μ is the intercept, X is the genotypic matrix, b is the effect of an allele substitution (or regression coefficient) and e is the residual term. A k-fold cross-validation for model evaluation is provided by *emCV*.

Value

The EM functions returns a list with the intercept (μ), the regression coefficient (b), the fitted value (\hat{y}), and the estimated intraclass-correlation (h^2).

The function *emCV* returns the predictive ability of each model, that is, the correlation between the predicted and observed values from k -fold cross-validations repeated n times.

Author(s)

Alencar Xavier

Examples

```
## Not run:

data(tpod)
emCV(y, gen, 3, 3)

## End(Not run)
```

WGR3 (MV)

*Multivariate Regression***Description**

Multivariate model to find breeding values.

Usage

```
mkr(Y,K)
mrr(Y,X)
```

Arguments

| | |
|---|--|
| Y | Numeric matrix of observations x trait. NA is allowed. |
| K | Numeric matrix containing the relationship matrix. |
| X | Numeric matrix containing the genotyping matrix. |

Details

The model for the kernel regression (mkr) is as follows:

$$Y = Mu + UB + E$$

where Y is a matrix of response variables, Mu represents the intercepts, U is the matrix of Eigenvector of K , b is a vector of regression coefficients and E is the residual matrix.

The model for the ridge regression (mrr) is as follows:

$$Y = Mu + XB + E$$

where Y is a matrix of response variables, Mu represents the intercepts, X is the matrix of genotypic information, B is the matrix of marker effects, and E is the residual matrix.

Algorithm: Residuals are assumed to be independent among traits. Regression coefficients are solved via a multivariate adaptation of Gauss-Seidel Residual Update. Variance and covariance components are solved with an efficient variation of EM-REML.

Other related implementations:

- 01) mkr2X(Y, K1, K2): Solves multi-trait kernel regressions with two random effects.
- 02) mrr2X(Y, X1, X2): Solves multi-trait ridge regressions with two random effects.
- 03) mrrR(Y, X): Variation of mrr that assumes correlated residuals.

Value

Returns a list with the random effect covariances (Vb), residual variances (Ve), genetic correlations (GC), matrix with marker effects (b) or eigenvector effects (if mkr), intercepts (mu), heritabilities (h2), and a matrix with fitted values (hat).

Author(s)

Alencar Xavier

Examples

```

## Not run:

# Load data and compute G matrix
data(tpod)
gen = CNT(gen)
K = tcrossprod(gen)
K = K/mean(diag(K))

# Phenotypes: 3 traits correlated r=0.5
G0 = 0.5+diag(0.5,3)
G = kronecker(G0,K)
diag(G)=diag(G)+0.001
L = chol(G)
TBV = crossprod(L,rnorm(196*3))
Y = rnorm(196*3,10+TBV,sd(TBV))
Phe = matrix(Y,ncol=3)
TBV = matrix(TBV,ncol=3)

# Fit kernel and regression models
test1 = mkr(Phe,K)
test2 = mrr(Phe,gen)

# Genetic correlation
test1$GC
test2$GC

# Heritabilities
test1$h2
test2$h2

# Goodness of fit
diag(cor(TBV,test1$hat))
diag(cor(TBV,test2$hat))

## End(Not run)

```

Description

Function to solve univariate mixed models with or without the usage of omic information. This function allows single-step modeling of replicated observations with marker information available

through the usage of a linkage function to connect to a whole-genome regression method. Genomic estimated values can be optionally deregressed (no shrinkage) while fitting the model.

Usage

```
mixed(y, random=NULL, fixed=NULL, data=NULL, X=list(),
      alg=emML, maxit=10, Deregress=FALSE, ...)
```

Arguments

| | |
|------------------------|---|
| <code>y</code> | Response variable from the data frame containing the dataset. |
| <code>random</code> | Formula. Right-hand side formula of random effects. |
| <code>fixed</code> | Formula. Right-hand side formula of fixed effects. |
| <code>data</code> | Data frame containing the response variable, random and fixed terms. |
| <code>X</code> | List of omic incidence matrix. Row names of these matrices connect the omic information to the levels of the indicated random terms (eg. <code>X=list("ID"=gen)</code>). |
| <code>alg</code> | Function. Whole-genome regression algorithm utilized to solve link functions. These include MCMC (<code>wgr</code> , <code>BayesB</code> , etc) and EM (<code>emEN</code> , <code>emDE</code> , etc) algorithms. By default, it runs maximum likelihood <code>emML</code> . |
| <code>maxit</code> | Integer. Maximum number of iterations. |
| <code>Deregress</code> | Logical. Deregress (unshrink) coefficients while fitting the model? |
| <code>...</code> | Additional arguments to be passed to the whole-genome regression algorithms specified on <code>alg</code> . |

Details

The model for the whole-genome regression is as follows:

$$y = Xb + Zu + Wa + e$$

where y is the response variable, Xb corresponds to the fixed effect term, Zu corresponds to one or more random effect terms, W is the incidence matrix of terms with omic information and a is omic values by $a = Mg$, where M is the genotypic matrix and g are marker effects. Here, e is the residual term. An example is provided using the data from the NAM package with: `demo(mixedmodel)`.

Value

The function `wgr` returns a list with Fitness values (`Fitness`) containing observation `obs`, fitted values `hat`, residuals `res`, and fitted values by model term `fits`; Estimated variance components (`VarComp`) containing the variance components per se (`VarComponents`) and variance explained by each model term (`VarExplained`), regression coefficients by model term (`Coefficients`), and the effects of structured terms (`Structure`) containing the marker effects of each model term where markers were provided.

Author(s)

Alencar Xavier

References

Xavier, A. (2019). Efficient Estimation of Marker Effects in Plant Breeding. G3: Genes, Genomes, Genetics, DOI: 10.1534/g3.119.400728

Examples

```
## Not run:  
demo(mixedmodel)  
  
## End(Not run)
```

XTRA 2

Additional tools

Description

Complementary functions that may help with handling parameters and routine operations.

Details

```
emGWA(y, gen) # Simple MLM for association analysis  
markov(gen, chr=NULL) # Markovian imputation of genotypes coded as 012  
IMP(X) # Imputes genotypes with SNP expectation (column average)  
CNT(X) # Recodes SNPs by centralizing columns in a matrix  
GAU(X) # Creates Gaussian kernel as exp(-Dist2/mean(Dist2))  
GRM(X, Code012=FALSE) # Creates additive kinship matrix VanRaden (2008)  
SPC(y, blk, row, col, rN=3, cN=1) # Spatial covariate  
SPM(blk, row, col, rN=3, cN=1) # Spatial design matrix  
SibZ(id, p1, p2) # Pedigree design matrix compatible to regression methods  
Hmat(ped, gen=NULL) # Kinship combining pedigree and genomics
```

Author(s)

Alencar Xavier

Index

BayesA (WGR1 (MC)), 3
BayesA2 (WGR1 (MC)), 3
BayesB (WGR1 (MC)), 3
BayesB2 (WGR1 (MC)), 3
BayesC (WGR1 (MC)), 3
BayesCpi (WGR1 (MC)), 3
BayesDpi (WGR1 (MC)), 3
BayesL (WGR1 (MC)), 3
BayesRR (WGR1 (MC)), 3
BayesRR2 (WGR1 (MC)), 3
bWGR (bWGR-package), 2
bWGR-package, 2

chr (Dataset), 2
CNT (XTRA 2), 11

Dataset, 2

emBA (WGR2 (EM)), 6
emBB (WGR2 (EM)), 6
emBC (WGR2 (EM)), 6
emBL (WGR2 (EM)), 6
emCV (WGR2 (EM)), 6
emDE (WGR2 (EM)), 6
emEN (WGR2 (EM)), 6
emGWA (XTRA 2), 11
emML (WGR2 (EM)), 6
emML2 (WGR2 (EM)), 6
emRR (WGR2 (EM)), 6

fam (Dataset), 2

GAU (XTRA 2), 11
gen (Dataset), 2
GRM (XTRA 2), 11
GSEN (WGR1 (MC)), 3

Hmat (XTRA 2), 11

IMP (XTRA 2), 11

KMUP (WGR1 (MC)), 3
KMUP2 (WGR1 (MC)), 3

markov (XTRA 2), 11
mcmcCV (WGR1 (MC)), 3
mixed (XTRA 1), 9
mkr (WGR3 (MV)), 8
mkr2X (WGR3 (MV)), 8
mrr (WGR3 (MV)), 8
mrr2X (WGR3 (MV)), 8
mrrFast (WGR3 (MV)), 8
mrrR (WGR3 (MV)), 8
mtgsru (XTRA 1), 9
mtmixed (XTRA 1), 9

SibZ (XTRA 2), 11
SPC (XTRA 2), 11
SPM (XTRA 2), 11

tpod (Dataset), 2

wgr (WGR1 (MC)), 3
WGR1 (MC), 3
WGR2 (EM), 6
WGR3 (MV), 8

XTRA 1, 9
XTRA 2, 11

y (Dataset), 2