

Package ‘first’

February 16, 2024

Title Factor Importance Ranking and Selection using Total Indices

Version 1.1

Imports stats, parallel, FNN, twinning

Description

A model-independent factor importance ranking and selection procedure that is based on total Sobol' indices. Please see Huang and Joseph (2024) <[arXiv:2401.00800](#)>. This research is supported by U.S. National Science Foundation grants DMS-2310637 and DMREF-1921873.

License GPL (>= 2)

Encoding UTF-8

RoxygenNote 7.1.2

NeedsCompilation no

Author Chaofan Huang [aut, cre],
V. Roshan Joseph [aut]

Maintainer Chaofan Huang <chaofan.huang@gatech.edu>

Repository CRAN

Date/Publication 2024-02-16 04:30:02 UTC

R topics documented:

first	1
totalsobol.knn	4
Index	7

first *Factor Importance Ranking and Selection using Total (Sobol') indices*

Description

`first` implements the **model-independent** factor importance and selection procedure proposed in Huang and Joseph (2024). The importance measure is based on total Sobol' indices from global sensitivity analysis. Factor importance computation and selection are performed directly from the noisy data. Parallel computations are available to accelerate the estimation. For categorical data inputs, please convert them to factor type before calling the function.

Usage

```
first(
  X,
  y,
  n.knn = NULL,
  n.mc = nrow(X),
  twin.mc = FALSE,
  rescale = TRUE,
  n.forward = 2,
  par1 = NULL,
  verbose = FALSE
)
```

Arguments

<code>X</code>	a matrix or data frame for the factors / predictors.
<code>y</code>	a vector for the responses.
<code>n.knn</code>	number of nearest-neighbors for the inner loop conditional variance estimation. <code>n.knn=2</code> is recommended for regression, and <code>n.knn=3</code> for binary classification.
<code>n.mc</code>	number of Monte Carlo samples for the outer loop expectation estimation.
<code>twin.mc</code>	a logical indicating whether to use twinning subsamples, otherwise random subsamples are used. It is supported when the reduction ratio is at least 2.
<code>rescale</code>	a logical logical indicating whether to standardize the factors / predictors.
<code>n.forward</code>	number of times to run the forward selection phase to tradeoff between efficiency and accuracy. <code>n_forward=2</code> is recommended. To run the complete forward selection, please set <code>n_forward</code> to the number of factors / predictors.
<code>par1</code>	number of cores on which to parallelize the computation. If <code>NULL</code> , then no parallelization is done.
<code>verbose</code>	a logical indicating whether to display intermediate results.

Details

`first` provides factor importance ranking and selection directly from scattered data without any model fitting. Factor importance is computed based on total Sobol' indices (Sobol', 2001), which is connected to the approximation error introduced by excluding the factor of interest (Huang and Joseph, 2024). The estimation procedure adapts from the Nearest-Neighbor estimator in Broto et al. (2020) to account for the noisy data. Integrating it with forward selection and backward elimination allows for factor selection.

`first` belongs to the class of forward-backward selection with early dropping algorithm (Borboudakis and Tsamardinos, 2019). In forward selection, each time we find the candidate that maximizes the output variance that can be explained. For candidates that cannot improve the variance explained conditional on the selected factors, they are removed from the candidate set. This forward selection step is run `n_forward` times to tradeoff between accuracy and efficiency. `n_forward=2` is recommended in Yu et al. (2020). To run the complete forward selection, please set `n_forward` to the number of factors / predictors. In backward elimination, we again remove one factor at a time, starting with the factor that can improve the explained variance most, till no factor can further improve.

`n.knn=2` nearest-neighbors is recommended for integer/numeric output, and `n.knn=3` is suggested for binary output. For numeric inputs, it is recommended to standardize them via setting the argument `rescale=TRUE`. Categorical inputs are transformed via one-hot encoding for the nearest-neighbor search. To speed up the nearest-neighbor search, k-d tree from the **FNN** package is used. Also, parallel computation is also supported via the **parallel** package.

For large datasets, we support the use of subsamples for further acceleration. Use argument `n.mc` to specify the number of subsamples. Two options are available for finding the subsamples: random and twinning (Vakayil and Joseph, 2022). Twinning is able to find subsamples that better represent the big data, i.e., providing a more accurate estimation, but at a slightly higher computational cost. For more details, please see the **twinning** package.

Value

A numeric vector of the factor importance, with zero indicating that the factor is not important to the prediction of the response.

Author(s)

Chaofan Huang <chaofan.huang@gatech.edu> and V. Roshan Joseph <roshan@gatech.edu>

References

- Huang, C., & Joseph, V. R. (2024). Factor Importance Ranking and Selection using Total Indices. arXiv preprint arXiv:2401.00800.
- Sobol', I. M. (2001). Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates. *Mathematics and computers in simulation*, 55(1-3), 271-280.
- Broto, B., Bachoc, F., & Depecker, M. (2020). Variance reduction for estimation of Shapley effects and adaptation to unknown input distribution. *SIAM/ASA Journal on Uncertainty Quantification*, 8(2), 693-716.
- Borboudakis, G., & Tsamardinos, I. (2019). Forward-backward selection with early dropping. *The Journal of Machine Learning Research*, 20(1), 276-314.
- Yu, K., Guo, X., Liu, L., Li, J., Wang, H., Ling, Z., & Wu, X. (2020). Causality-based feature selection: Methods and evaluations. *ACM Computing Surveys (CSUR)*, 53(5), 1-36.
- Vakayil, A., & Joseph, V. R. (2022). Data twinning. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, 15(5), 598-610.

Examples

```

ishigami <- function(x) {
  x <- -pi + 2*pi*x
  y <- sin(x[1]) + 7*sin(x[2])^2 + 0.1*x[3]^4*sin(x[1])
  return (y)
}
set.seed(123)
n <- 1000
p <- 6
X <- matrix(runif(n*p), ncol=p)
y <- apply(X,1,ishigami) + rnorm(n)
imp <- first(X, y, n.knn=2, rescale=FALSE, verbose=TRUE)
print(round(imp,3)) # Only first 3 factors are important

```

totalsobol.knn

Estimating Total Sobol' Indices from Data

Description

totalsobol.knn implements the estimation of the total Sobol' indices directly from scattered data. Parallel computation is available to accelerate the estimation. For categorical inputs, please convert them to factor before calling this function. For large datasets, we support the use of subsample to reduce the computational cost.

Usage

```

totalsobol.knn(
  X,
  y,
  noise,
  n.knn = NULL,
  n.mc = nrow(X),
  twin.mc = FALSE,
  rescale = TRUE,
  parl = NULL
)

```

Arguments

X	a matrix or data frame for the factors / predictors.
y	a vector for the responses.
noise	a logical indicating whether the responses are noisy.
n.knn	number of nearest-neighbors for the inner loop conditional variance estimation. n.knn=2 is recommended for regression, and n.knn=3 for binary classification.
n.mc	number of Monte Carlo samples for the outer loop expectation estimation.

<code>twin.mc</code>	a logical indicating whether to use twinning subsamples, otherwise random subsamples are used. It is supported when the reduction ratio is at least 2.
<code>rescale</code>	a logical logical indicating whether to standardize the factors / predictors.
<code>par1</code>	number of cores on which to parallelize the computation. If NULL, then no parallelization is done.

Details

`totalsobol.knn` provides consistent estimation of the total Sobol' indices (Sobol, 1993) from scattered data. When the output is clean/noiseless (`noise=FALSE`), `totalsobol.knn` implements the Nearest-Neighbor estimator proposed in Broto et al. (2020). See `shapleysobol_knn` from the **sensitivity** package for another implementation of the nearest-neighbor estimator. When the output is noisy (`noise=TRUE`), `totalsobol.knn` implements the Noise-Adjusted Nearest-Neighbor (NANNE) estimator (Huang and Joseph, 2024). NANNE estimator can correct the estimation bias of the nearest-neighbor estimator caused by the random noise. Please see Huang and Joseph (2024) for a more detailed discussion and comparison.

For integer/numeric output, `n.knn=2` nearest-neighbors is recommended for the noisy data (Huang and Joseph, 2024), and `n.knn=3` nearest-neighbors is suggested for the clean/noiseless data (Broto et al., 2020). For numeric inputs, it is recommended to standardize them via setting the argument `rescale=TRUE`. Categorical inputs are transformed via one-hot encoding for the nearest-neighbor search. To speed up the nearest-neighbor search, k-d tree from the **FNN** package is used. Also, parallel computation is also supported via the **parallel** package.

Last, for large datasets, we support the use of subsamples for further acceleration. Use argument `n.mc` to specify the number of subsamples. Two options are available for finding the subsamples: `random` and `twinning` (Vakayil and Joseph, 2022). `twinning` is able to find subsamples that better represent the big data, i.e., providing a more accurate estimation, but at a slightly higher computational cost. For more details, please see the **twinning** package.

Value

A numeric vector of the total Sobol' indices estimation.

Author(s)

Chaofan Huang <chaofan.huang@gatech.edu> and V. Roshan Joseph <roshan@gatech.edu>

References

- Huang, C., & Joseph, V. R. (2024). Factor Importance Ranking and Selection using Total Indices. arXiv preprint arXiv:2401.00800.
- Sobol', I. M. (2001). Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates. *Mathematics and computers in simulation*, 55(1-3), 271-280.
- Broto, B., Bachoc, F., & Depecker, M. (2020). Variance reduction for estimation of Shapley effects and adaptation to unknown input distribution. *SIAM/ASA Journal on Uncertainty Quantification*, 8(2), 693-716.
- Vakayil, A., & Joseph, V. R. (2022). Data twinning. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, 15(5), 598-610.

Examples

```
ishigami <- function(x) {  
  x <- -pi + 2*pi*x  
  y <- sin(x[1]) + 7*sin(x[2])^2 + 0.1*x[3]^4*sin(x[1])  
  return (y)  
}  
  
set.seed(123)  
n <- 10000  
p <- 3  
X <- matrix(runif(n*p), ncol=p)  
y <- apply(X,1,ishigami) + rnorm(n)  
tsi <- totalsobol.knn(X, y, noise=TRUE, n.knn=2, rescale=FALSE)  
print(round(tsi,3)) # Analytical Total Sobol' Indices: 0.558, 0.442, 0.244
```

Index

`first`, 1

`totalsobol.knn`, 4