

# Package ‘SCE’

July 25, 2025

**Title** Stepwise Clustered Ensemble

**Version** 1.1.1

**Description** Implementation of Stepwise Clustered Ensemble (SCE) and Stepwise Cluster Analysis (SCA) for multivariate data analysis. The package provides comprehensive tools for feature selection, model training, prediction, and evaluation in hydrological and environmental modeling applications. Key functionalities include recursive feature elimination (RFE), Wilks feature importance analysis, model validation through out-of-bag (OOB) validation, and ensemble prediction capabilities. The package supports both single and multivariate response variables, making it suitable for complex environmental modeling scenarios. For more details see Li et al. (2021) <[doi:10.5194/hess-25-4947-2021](https://doi.org/10.5194/hess-25-4947-2021)>.

**URL** <https://doi.org/10.5194/hess-25-4947-2021>

**License** GPL-3

**Encoding** UTF-8

**RoxxygenNote** 7.2.3

**Depends** R (>= 3.5.0)

**Imports** stats (>= 3.5.0), utils (>= 3.5.0)

**Suggests** testthat (>= 3.0.0), knitr, rmarkdown

**NeedsCompilation** no

**Author** Kailong Li [aut, cre]

**Maintainer** Kailong Li <lk198509509@gmail.com>

**Repository** CRAN

**Date/Publication** 2025-07-25 21:00:02 UTC

## Contents

Air_quality_datasets . . . . .	2
evaluate . . . . .	3
importance . . . . .	4
Plot_RFE . . . . .	4
predict . . . . .	5
print . . . . .	6

RFE_SCE . . . . .	6
SCA . . . . .	7
SCE . . . . .	8
Streamflow_datasets . . . . .	10
<b>Index</b>	<b>12</b>

---



---

*Air\_quality\_datasets Air Quality Dataset*

---

## Description

These datasets contain air quality measurements for training and testing purposes. They include various air pollutant concentrations and meteorological variables measured at different locations and times.

## Usage

```
data("Air_quality_training")
data("Air_quality_testing")
```

## Format

Both datasets are data frames with 8760 rows and 12 variables:

**Date** Date and time of measurement (POSIXct format)

**PM2.5** Particulate matter with diameter less than 2.5 micrometers ( $\mu\text{g}/\text{m}^3$ )

**PM10** Particulate matter with diameter less than 10 micrometers ( $\mu\text{g}/\text{m}^3$ )

**SO2** Sulfur dioxide concentration ( $\mu\text{g}/\text{m}^3$ )

**NO2** Nitrogen dioxide concentration ( $\mu\text{g}/\text{m}^3$ )

**CO** Carbon monoxide concentration ( $\mu\text{g}/\text{m}^3$ )

**O3** Ozone concentration ( $\mu\text{g}/\text{m}^3$ )

**TEMP** Temperature ( $\text{degree C}$ )

**PRES** Atmospheric pressure (hPa)

**DEWP** Dew point temperature ( $\text{degree C}$ )

**RAIN** Precipitation amount (mm)

**WSPM** Wind speed (m/s)

## Details

### Dataset Differences:

- **Air\_quality\_training:** Used for training SCA and SCE models
- **Air\_quality\_testing:** Used for testing trained models

### Variable Descriptions:

- **PM2.5, PM10:** Particulate matter concentrations, important indicators of air quality
- **SO<sub>2</sub>, NO<sub>2</sub>, CO, O<sub>3</sub>:** Major air pollutants regulated by environmental agencies
- **TEMP, PRES, DEWP:** Meteorological variables affecting air quality
- **RAIN, WSPM:** Weather conditions that influence pollutant dispersion

## Source

Air quality monitoring stations

---

evaluate

*Evaluate SCE and SCA Model Performance*

---

## Description

Evaluate model performance for SCE or SCA models.

## Usage

```
## S3 method for class 'SCE'  
evaluate(object, Testing_data, Training_data, digits = 3, ...)  
## S3 method for class 'SCA'  
evaluate(object, Testing_data, Training_data, digits = 3, ...)
```

## Arguments

object	An SCE or SCA model object
Testing_data	Testing dataset
Training_data	Training dataset
digits	Number of decimal places (default: 3)
...	Additional arguments

## Value

Model performance metrics.

## See Also

[SCE](#), [SCA](#), [predict](#)

**importance***Variable Importance for SCE and SCA Models***Description**

Calculate variable importance for SCE or SCA models.

**Usage**

```
## S3 method for class 'SCE'
importance(object, OOB_weight = TRUE, ...)
## S3 method for class 'SCA'
importance(object, ...)
```

**Arguments**

<code>object</code>	An SCE or SCA model object
<code>OOB_weight</code>	Use out-of-bag weights for importance calculation (SCE only, default: TRUE)
<code>...</code>	Additional arguments

**Value**

Variable importance rankings.

**See Also**

[SCE](#), [SCA](#), [RFE\\_SCE](#)

**Plot\_RFE***Plot Recursive Feature Elimination Results***Description**

Plot Recursive Feature Elimination results.

**Usage**

```
Plot_RFE(rfe_result,
         main = "OOB Validation and Testing R2 vs Number of Predictors",
         col_validation = "blue",
         col_testing = "red",
         pch = 16,
         lwd = 2,
         cex = 1.2,
         legend_pos = "bottomleft",
         ...)
```

**Arguments**

rfe_result	Result object from RFE_SCE function
main	Plot title
col_validation	Color for validation line
col_testing	Color for testing line
pch	Point character
lwd	Line width
cex	Point size
legend_pos	Legend position
...	Additional arguments

**Value**

Plot showing validation and testing R2 vs number of predictors.

**See Also**

[RFE\\_SCE](#)

---

predict

*Predict Using SCE and SCA Models*

---

**Description**

Make predictions on new data using SCE or SCA models.

**Usage**

```
## S3 method for class 'SCE'  
predict(object, newdata, ...)  
## S3 method for class 'SCA'  
predict(object, newdata, ...)
```

**Arguments**

object	An SCE or SCA model object
newdata	New data for prediction
...	Additional arguments

**Value**

Predictions for the new data.

**See Also**

[SCE](#), [SCA](#), [evaluate](#)

<code>print</code>	<i>Print SCE and SCA Model Objects</i>
--------------------	--

### Description

Print information about SCE or SCA model objects.

### Usage

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

### Arguments

x	An SCE or SCA model object
...	Additional arguments (not used)

### Details

For SCE objects, prints ensemble information including number of trees, parameters, predictors, predictants, and OOB performance metrics.

For SCA objects, prints tree structure information including total nodes, leaf nodes, cutting/merging actions, and variable names.

### Value

Prints model information and returns the object invisibly.

### See Also

[SCE](#), [SCA](#), [summary](#)

<code>RFE_SCE</code>	<i>Recursive Feature Elimination for SCE Models</i>
----------------------	---

### Description

Recursive Feature Elimination for SCE models to identify the most important predictors.

### Usage

```
RFE_SCE(Training_data, Testing_data, Predictors, Predictant, Nmin, Ntree,
        alpha = 0.05, resolution = 1000, step = 1, verbose = TRUE,
        parallel = TRUE)
```

**Arguments**

Training_data	Training dataset
Testing_data	Testing dataset
Predictors	Character vector of predictor names
Predictant	Character vector of predictant names
Nmin	Minimum samples per node
Ntree	Number of trees
alpha	Significance level (default: 0.05)
resolution	Resolution for splitting (default: 1000)
step	Number of predictors to remove per iteration (default: 1)
verbose	Print progress (default: TRUE)
parallel	Use parallel processing (default: TRUE)

**Value**

RFE results with performance metrics and importance scores.

**See Also**

[Plot\\_RFE](#), [SCE](#), [importance](#)

---

SCA

*Stepwise Cluster Analysis (SCA)*

---

**Description**

Builds a single Stepwise Cluster Analysis (SCA) tree model that recursively partitions the data space based on Wilks' Lambda statistic.

**Usage**

```
SCA(Training_data, X, Y, Nmin, alpha = 0.05, resolution = 1000, verbose = FALSE)
```

**Arguments**

Training_data	A data.frame containing the training data
X	Character vector of predictor variable names
Y	Character vector of predictant variable names
Nmin	Minimum number of samples in a leaf node
alpha	Significance level for clustering (default: 0.05)
resolution	Resolution for splitting (default: 1000)
verbose	Print progress information (default: FALSE)

**Value**

An S3 object of class "SCA" containing the tree model.

**See Also**

[SCE](#), [predict](#), [importance](#), [evaluate](#)

**Examples**

```
# Load example data
data(Streamflow_training_10var)
data(Streamflow_testing_10var)

# Define variables
Predictors <- c("Prcp", "SRad", "Tmax", "Tmin", "VP", "smlt", "swvl1", "swvl2", "swvl3", "swvl4")
Predictants <- c("Flow")

# Build SCA model
sca_model <- SCA(
  Training_data = Streamflow_training_10var,
  X = Predictors,
  Y = Predictants,
  Nmin = 5,
  alpha = 0.05,
  resolution = 1000
)

# Use S3 methods
print(sca_model)
summary(sca_model)
sca_predictions <- predict(sca_model, Streamflow_testing_10var)
sca_importance <- importance(sca_model)
sca_evaluation <- evaluate(sca_model, Streamflow_testing_10var, Streamflow_training_10var)
```

**Description**

Builds a Stepwise Clustered Ensemble (SCE) model, which is an ensemble of SCA trees built using bootstrap samples and random feature selection, providing improved prediction accuracy and robustness.

**Usage**

```
SCE(Training_data, X, Y, mfeature, Nmin, Ntree, alpha = 0.05,
  resolution = 1000, verbose = FALSE, parallel = TRUE)
```

## Arguments

Training_data	A data.frame containing the training data
X	Character vector of predictor variable names
Y	Character vector of predictant variable names
mfeature	Number of features to randomly select for each tree
Nmin	Minimum number of samples in a leaf node
Ntree	Number of trees in the ensemble
alpha	Significance level for clustering (default: 0.05)
resolution	Resolution for splitting (default: 1000)
verbose	Print progress information (default: FALSE)
parallel	Use parallel processing (default: TRUE)

## Value

An S3 object of class "SCE" containing the ensemble model.

## See Also

[SCA](#), [predict](#), [importance](#), [evaluate](#)

## Examples

```
# Load example data
data(Streamflow_training_10var)
data(Streamflow_testing_10var)

# Define variables
Predictors <- c("Prcp", "SRad", "Tmax", "Tmin", "VP", "smlt", "swvl1", "swvl2", "swvl3", "swvl4")
Predictants <- c("Flow")

# Build SCE model
sce_model <- SCE(
  Training_data = Streamflow_training_10var,
  X = Predictors,
  Y = Predictants,
  mfeature = round(0.5 * length(Predictors)),
  Nmin = 5,
  Ntree = 48,
  alpha = 0.05,
  resolution = 1000,
  parallel = FALSE
)

# Use S3 methods
print(sce_model)
summary(sce_model)
sce_predictions <- predict(sce_model, Streamflow_testing_10var)
sce_importance <- importance(sce_model)
```

```
sce_evaluation <- evaluate(sce_model, Streamflow_testing_10var, Streamflow_training_10var)
```

**Streamflow\_datasets**      *Streamflow Dataset*

### Description

These datasets contain streamflow and related environmental variables for training and testing purposes. They are used in examples to demonstrate the SCE package functionality with different levels of complexity.

### Usage

```
data("Streamflow_training_10var")
data("Streamflow_training_22var")
data("Streamflow_testing_10var")
data("Streamflow_testing_22var")
```

### Format

**Streamflow\_training\_10var:** Basic environmental variables (12 columns):

**Date** Date and time of measurement  
**Prep** Monthly mean daily precipitation (mm)  
**SRad** Monthly mean daily solar radiation (W/m<sup>2</sup>)  
**Tmax** Monthly mean daily maximum temperature (°C)  
**Tmin** Monthly mean daily minimum temperature (°C)  
**VP** Monthly mean daily vapor pressure (Pa)  
**smlt** Monthly snowmelt (m)  
**swvl1** Soil water content layer 1 (m<sup>3</sup>/m<sup>3</sup>)  
**swvl2** Soil water content layer 2 (m<sup>3</sup>/m<sup>3</sup>)  
**swvl3** Soil water content layer 3 (m<sup>3</sup>/m<sup>3</sup>)  
**swvl4** Soil water content layer 4 (m<sup>3</sup>/m<sup>3</sup>)  
**Flow** Monthly mean daily streamflow (cfs)

**Streamflow\_training\_22var:** Extended variables with climate indices (24 columns):

**Flow** Streamflow measurements  
**IPO** Interdecadal Pacific Oscillation  
**IPO\_lag1** IPO with 1-month lag  
**IPO\_lag2** IPO with 2-month lag  
**Nino3.4** Nino 3.4 index

**Nino3.4\_lag1** Nino 3.4 with 1-month lag  
**Nino3.4\_lag2** Nino 3.4 with 2-month lag  
**PDO** Pacific Decadal Oscillation  
**PDO\_lag1** PDO with 1-month lag  
**PDO\_lag2** PDO with 2-month lag  
**PNA** Pacific North American pattern  
**PNA\_lag1** PNA with 1-month lag  
**PNA\_lag2** PNA with 2-month lag  
**Precipitation** Monthly precipitation  
**Precipitation\_2Mon** 2-month precipitation  
**Radiation** Solar radiation  
**Radiation\_2Mon** 2-month solar radiation  
**Tmax** Maximum temperature  
**Tmax\_2Mon** 2-month maximum temperature  
**Tmin** Minimum temperature  
**Tmin\_2Mon** 2-month minimum temperature  
**VP** Vapor pressure  
**VP\_2Mon** 2-month vapor pressure

**Testing datasets:** Same structure as corresponding training datasets.

## Details

### Dataset Structure:

- **10var datasets:** Basic environmental variables (12 columns)
- **22var datasets:** Extended variables with climate indices (24 columns)
- **Training datasets:** Used for model building
- **Testing datasets:** Used for model evaluation

**Climate Indices:** IPO (Interdecadal Pacific Oscillation), Nino3.4 (El Niño), PDO (Pacific Decadal Oscillation), PNA (Pacific North American pattern)

**Data Sources:** ERA5 Land, Daymet, USGS, and climate indices databases

## Source

Environmental monitoring stations, climate indices databases, ERA5 Land, Daymet, and USGS

# Index

Air\_quality\_datasets, 2  
Air\_quality\_testing  
    (Air\_quality\_datasets), 2  
Air\_quality\_training  
    (Air\_quality\_datasets), 2  
  
evaluate, 3, 5, 8, 9  
  
importance, 4, 7–9  
  
Plot\_RFE, 4, 7  
predict, 3, 5, 8, 9  
print, 6  
  
RFE\_SCE, 4, 5, 6  
  
SCA, 3–6, 7, 9  
SCE, 3–8, 8  
Streamflow\_datasets, 10  
Streamflow\_testing\_10var  
    (Streamflow\_datasets), 10  
Streamflow\_testing\_22var  
    (Streamflow\_datasets), 10  
Streamflow\_training\_10var  
    (Streamflow\_datasets), 10  
Streamflow\_training\_22var  
    (Streamflow\_datasets), 10  
summary, 6