

# Package ‘CAST’

June 15, 2024

**Type** Package

**Title** 'caret' Applications for Spatial-Temporal Models

**Version** 1.0.2

**Author** Hanna Meyer [cre, aut],  
Carles Milà [aut],  
Marvin Ludwig [aut],  
Jan Linnenbrink [aut],  
Fabian Schumacher [aut],  
Philipp Otto [ctb],  
Chris Reudenbach [ctb],  
Thomas Nauss [ctb],  
Edzer Pebesma [ctb]

**Maintainer** Hanna Meyer <hanna.meyer@uni-muenster.de>

**Description** Supporting functionality to run 'caret' with spatial or spatial-temporal data. 'caret' is a frequently used package for model training and prediction using machine learning. CAST includes functions to improve spatial or spatial-temporal modelling tasks using 'caret'. It includes the newly suggested 'Nearest neighbor distance matching' cross-validation to estimate the performance of spatial prediction models and allows for spatial variable selection to select suitable predictor variables in view to their contribution to the spatial model performance. CAST further includes functionality to estimate the (spatial) area of applicability of prediction models. Methods are described in Meyer et al. (2018) <doi:10.1016/j.envsoft.2017.12.001>; Meyer et al. (2019) <doi:10.1016/j.ecolmodel.2019.10210X.13650>; Milà et al. (2022) <doi:10.1111/2041-210X.13851>; Meyer and Pebesma (2022) <doi:10.1038/s41467-022-29838-9>; Linnenbrink et al. (2023) <doi:10.5194/egusphere-2023-1308>. The package is described in detail in Meyer et al. (2024) <doi:10.48550/arXiv.2404.06978>.

**License** GPL (>= 2)

**URL** <https://github.com/HannaMeyer/CAST>,  
<https://hannameyer.github.io/CAST/>

**Encoding** UTF-8

**LazyData** false

**Depends** R (>= 4.1.0)

**BugReports** <https://github.com/HannaMeyer/CAST/issues/>

**Imports** caret, stats, utils, ggplot2, graphics, FNN, plyr, zoo,  
methods, grDevices, data.table, sf, forcats, twosamples, terra,  
sp

**Suggests** doParallel, lubridate, randomForest, knitr, geodata, mapview,  
rmarkdown, scales, parallel, gridExtra, viridis, stars, scam,  
rnaturalearth, MASS, RColorBrewer, tmap, PCAmixdata, gower,  
clustMixType, testthat (>= 3.0.0)

**RoxygenNote** 7.3.1

**VignetteBuilder** knitr

**Config/testthat/edition** 3

**NeedsCompilation** no

**Repository** CRAN

**Date/Publication** 2024-06-15 00:10:02 UTC

## Contents

aoa . . . . .	2
bss . . . . .	7
CAST . . . . .	8
clustered_sample . . . . .	10
cookfarm . . . . .	11
CreateSpacetimeFolds . . . . .	11
errorProfiles . . . . .	13
ffs . . . . .	16
geodist . . . . .	20
global_validation . . . . .	23
knndm . . . . .	24
nndm . . . . .	29
normalize_DI . . . . .	33
plot . . . . .	35
print . . . . .	37
splotdata . . . . .	37
trainDI . . . . .	38
<b>Index</b>	<b>42</b>

## Description

This function estimates the Dissimilarity Index (DI) and the derived Area of Applicability (AOA) of spatial prediction models by considering the distance of new data (i.e. a `SpatRaster` of spatial predictors used in the models) in the predictor variable space to the data used for model training. Predictors can be weighted based on the internal variable importance of the machine learning algorithm used for model training. The AOA is derived by applying a threshold on the DI which is the (outlier-removed) maximum DI of the cross-validated training data. Optionally, the local point density is calculated which indicates the number of similar training data points up to the DI threshold.

## Usage

```
aoa(
  newdata,
  model = NA,
  trainDI = NA,
  train = NULL,
  weight = NA,
  variables = "all",
  CVtest = NULL,
  CVtrain = NULL,
  method = "L2",
  useWeight = TRUE,
  LPD = FALSE,
  maxLPD = 1,
  indices = FALSE,
  verbose = TRUE
)
```

## Arguments

<code>newdata</code>	A <code>SpatRaster</code> , stars object or <code>data.frame</code> containing the data the model was meant to make predictions for.
<code>model</code>	A train object created with <code>caret</code> used to extract weights from (based on variable importance) as well as cross-validation folds. See examples for the case that no model is available or for models trained via e.g. <code>mlr3</code> .
<code>trainDI</code>	A <code>trainDI</code> object. Optional if <code>trainDI</code> was calculated beforehand.
<code>train</code>	A <code>data.frame</code> containing the data used for model training. Optional. Only required when no model is given
<code>weight</code>	A <code>data.frame</code> containing weights for each variable. Optional. Only required if no model is given.
<code>variables</code>	character vector of predictor variables. if "all" then all variables of the model are used or if no model is given then of the train dataset.
<code>CVtest</code>	list or vector. Either a list where each element contains the data points used for testing during the cross validation iteration (i.e. held back data). Or a vector that contains the ID of the fold for each training point. Only required if no model is given.

CVtrain	list. Each element contains the data points used for training during the cross validation iteration (i.e. held back data). Only required if no model is given and only required if CVtrain is not the opposite of CVtest (i.e. if a data point is not used for testing, it is used for training). Relevant if some data points are excluded, e.g. when using <a href="#">nndm</a> .
method	Character. Method used for distance calculation. Currently euclidean distance (L2) and Mahalanobis distance (MD) are implemented but only L2 is tested. Note that MD takes considerably longer.
useWeight	Logical. Only if a model is given. Weight variables according to importance in the model?
LPD	Logical. Indicates whether the local point density should be calculated or not.
maxLPD	numeric or integer. Only if LPD = TRUE. Number of nearest neighbors to be considered for the calculation of the LPD. Either define a number between 0 and 1 to use a percentage of the number of training samples for the LPD calculation or a whole number larger than 1 and smaller than the number of training samples. CAUTION! If not all training samples are considered, a fitted relationship between LPD and error metric will not make sense (@seealso <a href="#">DItoErrormetric</a> )
indices	logical. Calculate indices of the training data points that are responsible for the LPD of a new prediction location? Output is a matrix with the dimensions num(raster_cells) x maxLPD. Each row holds the indices of the training data points that are relevant for the specific LPD value at that location. Can be used in combination with <a href="#">exploreAOA(aoa)</a> function from the <a href="#">CASTvis package</a> for a better visual interpretation of the results. Note that the matrix can be quite big for examples with a high resolution and a larger number of training samples, which can cause memory issues.
verbose	Logical. Print progress or not?

## Details

The Dissimilarity Index (DI), the Local Data Point Density (LPD) and the corresponding Area of Applicability (AOA) are calculated. If variables are factors, dummy variables are created prior to weighting and distance calculation.

Interpretation of results: If a location is very similar to the properties of the training data it will have a low distance in the predictor variable space (DI towards 0) while locations that are very different in their properties will have a high DI. For easier interpretation see [normalize\\_DI](#) See Meyer and Pebesma (2021) for the full documentation of the methodology.

## Value

An object of class `aoa` containing:

parameters	object of class <code>trainDI</code> . see <a href="#">trainDI</a>
DI	SpatRaster, stars object or data frame. Dissimilarity index of newdata
LPD	SpatRaster, stars object or data frame. Local Point Density of newdata.
AOA	SpatRaster, stars object or data frame. Area of Applicability of newdata. AOA has values 0 (outside AOA) and 1 (inside AOA)

**Note**

If classification models are used, currently the variable importance can only be automatically retrieved if models were trained via `train(predictors,response)` and not via the formula-interface. Will be fixed.

**Author(s)**

Hanna Meyer, Fabian Schumacher

**References**

Meyer, H., Pebesma, E. (2021): Predicting into unknown space? Estimating the area of applicability of spatial prediction models. *Methods in Ecology and Evolution* 12: 1620-1633. doi:[10.1111/2041-210X.13650](https://doi.org/10.1111/2041-210X.13650)

**See Also**

[trainDI](#), [normalize\\_DI](#), [errorProfiles](#)

**Examples**

```
## Not run:
library(sf)
library(terra)
library(caret)
library(viridis)

# prepare sample data:
data(cookfarm)
dat <- aggregate(cookfarm[,c("VW", "Easting", "Northing")],
  by=list(as.character(cookfarm$SOURCEID), mean)
pts <- st_as_sf(dat, coords=c("Easting", "Northing"), crs=26911)
pts$ID <- 1:nrow(pts)
set.seed(100)
pts <- pts[1:30,]
studyArea <- rast(system.file("extdata", "predictors_2012-03-25.tif", package="CAST"))[[1:8]]
trainDat <- extract(studyArea, pts, na.rm=FALSE)
trainDat <- merge(trainDat, pts, by.x="ID", by.y="ID")

# visualize data spatially:
plot(studyArea)
plot(studyArea$DEM)
plot(pts[,1], add=TRUE, col="black")

# train a model:
set.seed(100)
variables <- c("DEM", "NDRE.Sd", "TWI")
model <- train(trainDat[,which(names(trainDat)%in%variables)],
  trainDat$VW, method="rf", importance=TRUE, tuneLength=1,
  trControl=trainControl(method="cv", number=5, savePredictions=T))
print(model) #note that this is a quite poor prediction model
```

```

prediction <- predict(studyArea,model,na.rm=TRUE)
plot(varImp(model,scale=FALSE))

#...then calculate the AOA of the trained model for the study area:
AOA <- aoa(studyArea, model)
plot(AOA)
plot(AOA$AOA)
#... or if preferred calculate the aoa and the LPD of the study area:
AOA <- aoa(studyArea, model, LPD = TRUE, maxLPD = 1)
plot(AOA$LPD)

####
#The AOA can also be calculated without a trained model.
#All variables are weighted equally in this case:
####
AOA <- aoa(studyArea,train=trainDat,variables=variables)

####
# The AOA can also be used for models trained via mlr3 (parameters have to be assigned manually):
####

library(mlr3)
library(mlr3learners)
library(mlr3spatial)
library(mlr3spatiotempcv)
library(mlr3extralearners)

# initiate and train model:
train_df <- trainDat[, c("DEM","NDRE.Sd","TWI", "VW")]
backend <- as_data_backend(train_df)
task <- as_task_regr(backend, target = "VW")
lrn <- lrn("regr.randomForest", importance = "mse")
lrn$train(task)

# cross-validation folds
rsmp_cv <- rsmp("cv", folds = 5L)$instantiate(task)

## predict:
prediction <- predict(studyArea,lrn$model,na.rm=TRUE)

### Estimate AOA
AOA <- aoa(studyArea,
           train = as.data.frame(task$data()),
           variables = task$feature_names,
           weight = data.frame(t(lrn$importance())),
           CVtest = rsmp_cv$instance[order(row_id)]$fold)

## End(Not run)

```

**Description**

Evaluate all combinations of predictors during model training

**Usage**

```
bss(  
  predictors,  
  response,  
  method = "rf",  
  metric = ifelse(is.factor(response), "Accuracy", "RMSE"),  
  maximize = ifelse(metric == "RMSE", FALSE, TRUE),  
  globalval = FALSE,  
  trControl = caret::trainControl(),  
  tuneLength = 3,  
  tuneGrid = NULL,  
  seed = 100,  
  verbose = TRUE,  
  ...  
)
```

**Arguments**

predictors	see <a href="#">train</a>
response	see <a href="#">train</a>
method	see <a href="#">train</a>
metric	see <a href="#">train</a>
maximize	see <a href="#">train</a>
globalval	Logical. Should models be evaluated based on 'global' performance? See <a href="#">global_validation</a>
trControl	see <a href="#">train</a>
tuneLength	see <a href="#">train</a>
tuneGrid	see <a href="#">train</a>
seed	A random number
verbose	Logical. Should information about the progress be printed?
...	arguments passed to the classification or regression routine (such as <code>randomForest</code> ).

## Details

bss is an alternative to [ffs](#) and ideal if the training set is small. Models are iteratively fitted using all different combinations of predictor variables. Hence,  $2^X$  models are calculated. Don't try running bss on very large datasets because the computation time is much higher compared to [ffs](#).

The internal cross validation can be run in parallel. See information on parallel processing of caret's train functions for details.

## Value

A list of class train. Beside of the usual train content the object contains the vector "selectedvars" and "selectedvars\_perf" that give the best variables selected as well as their corresponding performance. It also contains "perf\_all" that gives the performance of all model runs.

## Note

This variable selection is particularly suitable for spatial cross validations where variable selection MUST be based on the performance of the model for predicting new spatial units. Note that bss is very slow since all combinations of variables are tested. A more time efficient alternative is the forward feature selection ([ffs](#)).

## Author(s)

Hanna Meyer

## See Also

[train](#), [ffs](#), [trainControl](#), [CreateSpacetimeFolds](#), [nndm](#)

## Examples

```
## Not run:
data(iris)
bssmodel <- bss(iris[,1:4],iris$Species)
bssmodel$perf_all
plot(bssmodel)

## End(Not run)
```



## Description

Supporting functionality to run 'caret' with spatial or spatial-temporal data. 'caret' is a frequently used package for model training and prediction using machine learning. CAST includes functions to improve spatial-temporal modelling tasks using 'caret'. It includes the newly suggested 'Nearest neighbor distance matching' cross-validation to estimate the performance of spatial prediction models and allows for spatial variable selection to select suitable predictor variables in view to their contribution to the spatial model performance. CAST further includes functionality to estimate the (spatial) area of applicability of prediction models by analysing the similarity between new data and training data. Methods are described in Meyer et al. (2018); Meyer et al. (2019); Meyer and Pebesma (2021); Milà et al. (2022); Meyer and Pebesma (2022); Linnenbrink et al. (2023). The package is described in detail in Meyer et al. (2024).

## Details

'caret' Applications for Spatio-Temporal models

## Author(s)

Hanna Meyer, Carles Milà, Marvin Ludwig, Jan Linnenbrink, Fabian Schumacher

## References

- Meyer, H., Ludwig, L., Milà, C., Linnenbrink, J., Schumacher, F. (2024): The CAST package for training and assessment of spatial prediction models in R. arXiv, <https://doi.org/10.48550/arXiv.2404.06978>.
- Linnenbrink, J., Milà, C., Ludwig, M., and Meyer, H.: kNNDM: k-fold Nearest Neighbour Distance Matching Cross-Validation for map accuracy estimation, *EGUsphere* [preprint], <https://doi.org/10.5194/egusphere-2023-1308>, 2023.
- Milà, C., Mateu, J., Pebesma, E., Meyer, H. (2022): Nearest Neighbour Distance Matching Leave-One-Out Cross-Validation for map validation. *Methods in Ecology and Evolution* 00, 1– 13.
- Meyer, H., Pebesma, E. (2022): Machine learning-based global maps of ecological variables and the challenge of assessing them. *Nature Communications*. 13.
- Meyer, H., Pebesma, E. (2021): Predicting into unknown space? Estimating the area of applicability of spatial prediction models. *Methods in Ecology and Evolution*. 12, 1620– 1633.
- Meyer, H., Reudenbach, C., Wöllauer, S., Nauss, T. (2019): Importance of spatial predictor variable selection in machine learning applications - Moving from data reproduction to spatial prediction. *Ecological Modelling*. 411, 108815.
- Meyer, H., Reudenbach, C., Hengl, T., Katurji, M., Nauß, T. (2018): Improving performance of spatio-temporal machine learning models using forward feature selection and target-oriented validation. *Environmental Modelling & Software* 101: 1-9.

## See Also

Useful links:

- <https://github.com/HannaMeyer/CAST>
- <https://hannameyer.github.io/CAST/>
- Report bugs at <https://github.com/HannaMeyer/CAST/issues/>

---

clustered\_sample      *Clustered samples simulation*

---

### Description

A simple procedure to simulate clustered points based on a two-step sampling.

### Usage

```
clustered_sample(sarea, nsamples, nparents, radius)
```

### Arguments

sarea	polygon. Area where samples should be simulated.
nsamples	integer. Number of samples to be simulated.
nparents	integer. Number of parents.
radius	integer. Radius of the buffer around each parent for offspring simulation.

### Details

A simple procedure to simulate clustered points based on a two-step sampling. First, a pre-specified number of parents are simulated using random sampling. For each parent,  $(nsamples - nparents) / nparents$  are simulated within a radius of the parent point using random sampling.

### Value

sf object with the simulated points and the parent to which each point belongs to.

### Author(s)

Carles Milà

### Examples

```
# Simulate 100 points in a 100x100 square with 5 parents and a radius of 10.
library(sf)
library(ggplot2)

set.seed(1234)
simarea <- list(matrix(c(0,0,0,100,100,100,100,0,0,0), ncol=2, byrow=TRUE))
simarea <- sf::st_polygon(simarea)
simpoints <- clustered_sample(simarea, 100, 5, 10)
simpoints$parent <- as.factor(simpoints$parent)
ggplot() +
  geom_sf(data = simarea, alpha = 0) +
  geom_sf(data = simpoints, aes(col = parent))
```

---

cookfarm	<i>Cookfarm soil logger data</i>
----------	----------------------------------

---

### Description

spatio-temporal data of soil properties and associated predictors for the Cookfarm in Washington, USA. The data are a subset of the cookfarm dataset provided with the [GSIF package](#).

### Usage

```
data(cookfarm)
```

### Format

A sf data.frame with 128545 rows and 17 columns:

**SOURCEID** ID of the logger

**VW** Response Variable - Soil Moisture

**altitude** Measurement depth of VW

**Date, cdata** Measurement Date, Cumulative Date

**Easting, Northing** Location Coordinates (EPSG:26911)

**DEM, TWI, NDRE.M, NDRE.Sd, Precip\_wrcc, MaxT\_wrcc, MinT\_wrcc, Precip\_cum** Predictor Variables

### References

- Gash et al. 2015 - Spatio-temporal interpolation of soil water, temperature, and electrical conductivity in 3D + T: The Cook Agronomy Farm data set [doi:10.1016/j.spasta.2015.04.001](https://doi.org/10.1016/j.spasta.2015.04.001)
- Meyer et al. 2018 - Improving performance of spatio-temporal machine learning models using forward feature selection and target-oriented validation [doi:10.1016/j.envsoft.2017.12.001](https://doi.org/10.1016/j.envsoft.2017.12.001)

---

CreateSpacetimeFolds	<i>Create Space-time Folds</i>
----------------------	--------------------------------

---

### Description

Create spatial, temporal or spatio-temporal Folds for cross validation based on pre-defined groups

**Usage**

```
CreateSpacetimeFolds(  
  x,  
  spacevar = NA,  
  timevar = NA,  
  k = 10,  
  class = NA,  
  seed = sample(1:1000, 1)  
)
```

**Arguments**

x	data.frame containing spatio-temporal data
spacevar	Character indicating which column of x identifies the spatial units (e.g. ID of weather stations)
timevar	Character indicating which column of x identifies the temporal units (e.g. the day of the year)
k	numeric. Number of folds. If spacevar or timevar is NA and a leave one location out or leave one time step out cv should be performed, set k to the number of unique spatial or temporal units.
class	Character indicating which column of x identifies a class unit (e.g. land cover)
seed	numeric. See ?seed

**Details**

The function creates train and test sets by taking (spatial and/or temporal) groups into account. In contrast to `nndm`, it requires that the groups are already defined (e.g. spatial clusters or blocks or temporal units). Using "class" is helpful in the case that data are clustered in space and are categorical. E.g This is the case for land cover classifications when training data come as training polygons. In this case the data should be split in a way that entire polygons are held back (`spacevar="polygonID"`) but at the same time the distribution of classes should be similar in each fold (`class="LUC"`).

**Value**

A list that contains a list for model training and a list for model validation that can directly be used as "index" and "indexOut" in caret's `trainControl` function

**Note**

Standard k-fold cross-validation can lead to considerable misinterpretation in spatial-temporal modelling tasks. This function can be used to prepare a Leave-Location-Out, Leave-Time-Out or Leave-Location-and-Time-Out cross-validation as target-oriented validation strategies for spatial-temporal prediction tasks. See Meyer et al. (2018) for further information.

**Author(s)**

Hanna Meyer

## References

Meyer, H., Reudenbach, C., Hengl, T., Katurji, M., Nauß, T. (2018): Improving performance of spatio-temporal machine learning models using forward feature selection and target-oriented validation. *Environmental Modelling & Software* 101: 1-9.

## See Also

[trainControl](#), [ffs](#), [nndm](#)

## Examples

```
## Not run:
data(cookfarm)
### Prepare for 10-fold Leave-Location-and-Time-Out cross validation
indices <- CreateSpacetimeFolds(cookfarm, "SOURCEID", "Date")
str(indices)
### Prepare for 10-fold Leave-Location-Out cross validation
indices <- CreateSpacetimeFolds(dat, spacevar="SOURCEID")
str(indices)
### Prepare for leave-One-Location-Out cross validation
indices <- CreateSpacetimeFolds(dat, spacevar="SOURCEID",
  k=length(unique(dat$SOURCEID)))
str(indices)

## End(Not run)
```

---

errorProfiles

*Model and inspect the relationship between the prediction error and measures of dissimilarities and distances*

---

## Description

Performance metrics are calculated for moving windows of dissimilarity values based on cross-validated training data

## Usage

```
errorProfiles(
  model,
  trainDI = NULL,
  locations = NULL,
  variable = "DI",
  multiCV = FALSE,
  length.out = 10,
  window.size = 5,
  calib = "scam",
  method = "L2",
  useWeight = TRUE,
```

```

    k = 6,
    m = 2
  )

```

### Arguments

model	the model used to get the AOA
trainDI	the result of <code>trainDI</code> or <code>aoa</code> object <code>aoa</code>
locations	Optional. <code>sf</code> object for the training data used in model. Only used if <code>variable=="geodist"</code> . Note that they must be in the same order as <code>model\$trainingData</code> .
variable	Character. Which dissimilarity or distance measure to use for the error metric. Current options are "DI" or "LPD"
multiCV	Logical. Re-run model fitting and validation with different CV strategies. See details.
length.out	Numeric. Only used if <code>multiCV=TRUE</code> . Number of cross-validation folds. See details.
window.size	Numeric. Size of the moving window. See <code>rollapply</code> .
calib	Character. Function to model the DI/LPD~performance relationship. Currently <code>lm</code> and <code>scam</code> are supported
method	Character. Method used for distance calculation. Currently euclidean distance (L2) and Mahalanobis distance (MD) are implemented but only L2 is tested. Note that MD takes considerably longer. See <code>?aoa</code> for further explanation
useWeight	Logical. Only if a model is given. Weight variables according to importance in the model?
k	Numeric. See <code>mgcv::s</code>
m	Numeric. See <code>mgcv::s</code>

### Details

If `multiCV=TRUE` the model is re-fitted and validated by `length.out` new cross-validations where the cross-validation folds are defined by clusters in the predictor space, ranging from three clusters to LOOCV. Hence, a large range of dissimilarity values is created during cross-validation. If the AOA threshold based on the calibration data from multiple CV is larger than the original AOA threshold (which is likely if extrapolation situations are created during CV), the AOA threshold changes accordingly. See Meyer and Pebesma (2021) for the full documentation of the methodology.

### Value

A `scam`, linear model or exponential model

### Author(s)

Hanna Meyer, Marvin Ludwig, Fabian Schumacher

## References

Meyer, H., Pebesma, E. (2021): Predicting into unknown space? Estimating the area of applicability of spatial prediction models. [doi:10.1111/2041210X.13650](https://doi.org/10.1111/2041210X.13650)

## See Also

[aoa](#)

## Examples

```
## Not run:
library(CAST)
library(sf)
library(terra)
library(caret)

data(splotdata)
predictors <- terra::rast(system.file("extdata", "predictors_chile.tif", package="CAST"))

model <- caret::train(st_drop_geometry(splotdata)[,6:16], splotdata$Species_richness,
  ntree = 10, trControl = trainControl(method = "cv", savePredictions = TRUE))

AOA <- aoa(predictors, model, LPD = TRUE, maxLPD = 1)

### DI ~ error
errormodel_DI <- errorProfiles(model, AOA, variable = "DI")
plot(errormodel_DI)
summary(errormodel_DI)

expected_error_DI = terra::predict(AOA$DI, errormodel_DI)
plot(expected_error_DI)

### LPD ~ error
errormodel_LPD <- errorProfiles(model, AOA, variable = "LPD")
plot(errormodel_LPD)
summary(errormodel_DI)

expected_error_LPD = terra::predict(AOA$LPD, errormodel_LPD)
plot(expected_error_LPD)

### geodist ~ error
errormodel_geodist = errorProfiles(model, locations=splotdata, variable = "geodist")
plot(errormodel_geodist)
summary(errormodel_DI)

dist <- terra::distance(predictors[[1]], vect(splotdata))
names(dist) <- "geodist"
expected_error_DI <- terra::predict(dist, errormodel_geodist)
plot(expected_error_DI)

### with multiCV = TRUE (for DI ~ error)
```

```

errormodel_DI = errorProfiles(model, AOA, multiCV = TRUE, length.out = 3, variable = "DI")
plot(errormodel_DI)

expected_error_DI = terra::predict(AOA$DI, errormodel_DI)
plot(expected_error_DI)

# mask AOA based on new threshold from multiCV
mask_aoa = terra::mask(expected_error_DI, AOA$DI > attr(errormodel_DI, 'AOA_threshold'),
  maskvalues = 1)
plot(mask_aoa)

## End(Not run)

```

---

ffs

*Forward feature selection*


---

## Description

A simple forward feature selection algorithm

## Usage

```

ffs(
  predictors,
  response,
  method = "rf",
  metric = ifelse(is.factor(response), "Accuracy", "RMSE"),
  maximize = ifelse(metric == "RMSE", FALSE, TRUE),
  globalval = FALSE,
  withinSE = FALSE,
  minVar = 2,
  trControl = caret::trainControl(),
  tuneLength = 3,
  tuneGrid = NULL,
  seed = sample(1:1000, 1),
  verbose = TRUE,
  cores = 1,
  ...
)

```

## Arguments

predictors	see <a href="#">train</a>
response	see <a href="#">train</a>
method	see <a href="#">train</a>



metric	see <a href="#">train</a>
maximize	see <a href="#">train</a>
globalval	Logical. Should models be evaluated based on 'global' performance? See <a href="#">global_validation</a>
withinSE	Logical Models are only selected if they are better than the currently best models Standard error
minVar	Numeric. Number of variables to combine for the first selection. See Details.
trControl	see <a href="#">train</a>
tuneLength	see <a href="#">train</a>
tuneGrid	see <a href="#">train</a>
seed	A random number used for model training
verbose	Logical. Should information about the progress be printed?
cores	Numeric. If > 2, mclapply will be used. see <a href="#">mclapply</a>
...	arguments passed to the classification or regression routine (such as randomForest).

### Details

Models with two predictors are first trained using all possible pairs of predictor variables. The best model of these initial models is kept. On the basis of this best model the predictor variables are iteratively increased and each of the remaining variables is tested for its improvement of the currently best model. The process stops if none of the remaining variables increases the model performance when added to the current best model.

The forward feature selection can be run in parallel with forking on Linux systems ([mclapply](#)). Each fork computes a model, which drastically speeds up the runtime - especially of the initial predictor search. The internal cross validation can be run in parallel on all systems. See information on parallel processing of caret's train functions for details.

Using `withinSE` will favour models with less variables and probably shorten the calculation time

Per Default, the `ffs` starts with all possible 2-pair combinations. `minVar` allows to start the selection with more than 2 variables, e.g. `minVar=3` starts the `ffs` testing all combinations of 3 (instead of 2) variables first and then increasing the number. This is important for e.g. neural networks that often cannot make sense of only two variables. It is also relevant if it is assumed that the optimal variables can only be found if more than 2 are considered at the same time.

### Value

A list of class `train`. Beside of the usual `train` content the object contains the vector `"selectedvars"` and `"selectedvars_perf"` that give the order of the best variables selected as well as their corresponding performance (starting from the first two variables). It also contains `"perf_all"` that gives the performance of all model runs.

### Note

This variable selection is particularly suitable for spatial cross validations where variable selection **MUST** be based on the performance of the model for predicting new spatial units. See Meyer et al. (2018) and Meyer et al. (2019) for further details.

**Author(s)**

Hanna Meyer

**References**

- Gasch, C.K., Hengl, T., Gräler, B., Meyer, H., Magney, T., Brown, D.J. (2015): Spatio-temporal interpolation of soil water, temperature, and electrical conductivity in 3D+T: the Cook Agronomy Farm data set. *Spatial Statistics* 14: 70-90.
- Meyer, H., Reudenbach, C., Hengl, T., Katurji, M., Nauß, T. (2018): Improving performance of spatio-temporal machine learning models using forward feature selection and target-oriented validation. *Environmental Modelling & Software* 101: 1-9. doi:10.1016/j.envsoft.2017.12.001
- Meyer, H., Reudenbach, C., Wöllauer, S., Nauss, T. (2019): Importance of spatial predictor variable selection in machine learning applications - Moving from data reproduction to spatial prediction. *Ecological Modelling*. 411, 108815. doi:10.1016/j.ecolmodel.2019.108815.
- Ludwig, M., Moreno-Martinez, A., Hölzel, N., Pebesma, E., Meyer, H. (2023): Assessing and improving the transferability of current global spatial prediction models. *Global Ecology and Biogeography*. doi:10.1111/geb.13635.

**See Also**

[train.bss](#), [trainControl](#), [CreateSpacetimeFolds](#), [ndm](#)

**Examples**

```
## Not run:
data(splotdata)
ffsmodel <- ffs(splotdata[,6:12], splotdata$Species_richness, ntree = 20)

ffsmodel$selectedvars
ffsmodel$selectedvars_perf
plot(ffsmodel)
#or only selected variables:
plot(ffsmodel,plotType="selected")

## End(Not run)

# or perform model with target-oriented validation (LLO CV)
#the example is described in Gasch et al. (2015). The ffs approach for this dataset is described in
#Meyer et al. (2018). Due to high computation time needed, only a small and thus not robust example
#is shown here.

## Not run:
# run the model on three cores (see vignette for details):
library(doParallel)
library(lubridate)
cl <- makeCluster(3)
registerDoParallel(cl)

#load and prepare dataset:
data(cookfarm)
```

```

trainDat <- cookfarm[cookfarm$altitude==0.3&
  year(cookfarm$Date)==2012&week(cookfarm$Date)%in%c(13:14),]

#visualize dataset:
ggplot(data = trainDat, aes(x=Date, y=VW)) + geom_line(aes(colour=SOURCEID))

#create folds for Leave Location Out Cross Validation:
set.seed(10)
indices <- CreateSpacetimeFolds(trainDat,spacevar = "SOURCEID",k=3)
ctrl <- trainControl(method="cv",index = indices$index)

#define potential predictors:
predictors <- c("DEM","TWI","BLD","Precip_cum","cday","MaxT_wrcc",
  "Precip_wrcc","NDRE.M","Bt","MinT_wrcc","Northing","Easting")

#run ffs model with Leave Location out CV
set.seed(10)
ffsmodel <- ffs(trainDat[,predictors],trainDat$VW,method="rf",
  tuneLength=1,trControl=ctrl)
ffsmodel
plot(ffsmodel)
#or only selected variables:
plot(ffsmodel,plotType="selected")

#compare to model without ffs:
model <- train(trainDat[,predictors],trainDat$VW,method="rf",
  tuneLength=1, trControl=ctrl)
model
stopCluster(c1)

## End(Not run)

## Not run:
## on linux machines, you can also run the ffs in parallel with forks:
data("splotdata")
spatial_cv = CreateSpacetimeFolds(splotdata, spacevar = "Biome", k = 5)
ctrl <- trainControl(method="cv",index = spatial_cv$index)

ffsmodel <- ffs(predictors = splotdata[,6:16],
  response = splotdata$Species_richness,
  tuneLength = 1,
  method = "rf",
  trControl = ctrl,
  ntree = 20,
  seed = 1,
  cores = 4)

## End(Not run)

```

---

geodist	<i>Calculate euclidean nearest neighbor distances in geographic space or feature space</i>
---------	--

---

## Description

Calculates nearest neighbor distances in geographic space or feature space between training data as well as between training data and prediction locations. Optional, the nearest neighbor distances between training data and test data or between training data and CV iterations is computed.

## Usage

```
geodist(
  x,
  modeldomain = NULL,
  type = "geo",
  cvfolds = NULL,
  cvtrain = NULL,
  testdata = NULL,
  preddata = NULL,
  samplesize = 2000,
  sampling = "regular",
  variables = NULL,
  timevar = NULL,
  time_unit = "auto"
)
```

## Arguments

x	object of class sf, training data locations
modeldomain	SpatRaster, stars or sf object defining the prediction area (see Details)
type	"geo" or "feature". Should the distance be computed in geographic space or in the normalized multivariate predictor space (see Details)
cvfolds	optional. list or vector. Either a list where each element contains the data points used for testing during the cross validation iteration (i.e. held back data). Or a vector that contains the ID of the fold for each training point. See e.g. ?create-Folds or ?CreateSpacetimeFolds or ?nndm
cvtrain	optional. List of row indices of x to fit the model to in each CV iteration. If cvtrain is null but cvfolds is not, all samples but those included in cvfolds are used as training data
testdata	optional. object of class sf: Point data used for independent validation
preddata	optional. object of class sf: Point data indicating the locations within the modeldomain to be used as target prediction points. Useful when the prediction objective is a subset of locations within the modeldomain rather than the whole area.

samplesize	numeric. How many prediction samples should be used?
sampling	character. How to draw prediction samples? See <a href="#">spsample</a> . Use sampling = "Fibonacci" for global applications.
variables	character vector defining the predictor variables used if type="feature". If not provided all variables included in modeldomain are used.
timevar	optional. character. Column that indicates the date. Only used if type="time".
time_unit	optional. Character. Unit for temporal distances See <code>?difftime</code> . Only used if type="time".

### Details

The modeldomain is a sf polygon or a raster that defines the prediction area. The function takes a regular point sample (amount defined by samplesize) from the spatial extent. If type = "feature", the argument modeldomain (and if provided then also the testdata and/or preddata) has to include predictors. Predictor values for x, testdata and preddata are optional if modeldomain is a raster. If not provided they are extracted from the modeldomain rasterStack. If some predictors are categorical (i.e., of class factor or character), gower distances will be used. W statistic describes the match between the distributions. See Linnenbrink et al (2023) for further details.

### Value

A data.frame containing the distances. Unit of returned geographic distances is meters. attributes contain W statistic between prediction area and either sample data, CV folds or test data. See details.

### Note

See Meyer and Pebesma (2022) for an application of this plotting function

### Author(s)

Hanna Meyer, Edzer Pebesma, Marvin Ludwig, Jan Linnenbrink

### See Also

[nndm knndm](#)

### Examples

```
## Not run:
library(CAST)
library(sf)
library(terra)
library(caret)
library(rnaturalearth)
library(ggplot2)

data(splotdata)
studyArea <- rnaturalearth::ne_countries(continent = "South America", returnclass = "sf")
```

```

##### Distance between training data and new data:
dist <- geodist(splotdata, studyArea)
# With density functions
plot(dist)
# Or ECDFs (relevant for nndm and knnmd methods)
plot(dist, stat="ecdf")

##### Distance between training data, new data and test data (here Chile):
plot(splotdata[, "Country"])
dist <- geodist(splotdata[splotdata$Country != "Chile",], studyArea,
               testdata = splotdata[splotdata$Country == "Chile",])
plot(dist)

##### Distance between training data, new data and CV folds:
folds <- createFolds(1:nrow(splotdata), k=3, returnTrain=FALSE)
dist <- geodist(x=splotdata, modeldomain=studyArea, cvfolds=folds)
# Using density functions
plot(dist)
# Using ECDFs (relevant for nndm and knnmd methods)
plot(dist, stat="ecdf")

##### Distances in the feature space:
predictors <- terra::rast(system.file("extdata", "predictors_chile.tif", package="CAST"))
dist <- geodist(x = splotdata,
               modeldomain = predictors,
               type = "feature",
               variables = c("bio_1", "bio_12", "elev"))
plot(dist)

dist <- geodist(x = splotdata[splotdata$Country != "Chile",],
               modeldomain = predictors, cvfolds = folds,
               testdata = splotdata[splotdata$Country == "Chile",],
               type = "feature",
               variables=c("bio_1", "bio_12", "elev"))
plot(dist)

##### Distances in temporal space
library(lubridate)
library(ggplot2)
data(cookfarm)
dat <- st_as_sf(cookfarm, coords=c("Easting", "Northing"))
st_crs(dat) <- 26911
trainDat <- dat[dat$altitude== -0.3 & lubridate::year(dat$date)==2010,]
predictionDat <- dat[dat$altitude== -0.3 & lubridate::year(dat$date)==2011,]
trainDat$week <- lubridate::week(trainDat$date)
cvfolds <- CreateSpacetimeFolds(trainDat, timevar = "week")

dist <- geodist(trainDat, preddata = predictionDat, cvfolds = cvfolds$indexOut,
               type="time", time_unit="days")
plot(dist)+ xlim(0,10)

```

```
##### Example for a random global dataset
##### (refer to figure in Meyer and Pebesma 2022)

### Define prediction area (here: global):
ee <- st_crs("+proj=eqearth")
co <- ne_countries(returnclass = "sf")
co.ee <- st_transform(co, ee)

### Simulate a spatial random sample
### (alternatively replace pts_random by a real sampling dataset (see Meyer and Pebesma 2022):
sf_use_s2(FALSE)
pts_random <- st_sample(co.ee, 2000, exact=FALSE)

### See points on the map:
ggplot() + geom_sf(data = co.ee, fill="#00BFC4",col="#00BFC4") +
  geom_sf(data = pts_random, color = "#F8766D",size=0.5, shape=3) +
  guides(fill = "none", col = "none") +
  labs(x = NULL, y = NULL)

### plot distances:
dist <- geodist(pts_random,co.ee)
plot(dist) + scale_x_log10(labels=round)

## End(Not run)
```

---

global\_validation      *Evaluate 'global' cross-validation*

---

## Description

Calculate validation metric using all held back predictions at once

## Usage

```
global_validation(model)
```

## Arguments

model                  an object of class `train`

## Details

Relevant when folds are not representative for the entire area of interest. In this case, metrics like  $R^2$  are not meaningful since it doesn't reflect the general ability of the model to explain the entire gradient of the response. Comparable to LOOCV, predictions from all held back folds are used here together to calculate validation statistics.

**Value**

regression ([postResample](#)) or classification ([confusionMatrix](#)) statistics

**Author(s)**

Hanna Meyer

**See Also**

[CreateSpacetimeFolds](#)

**Examples**

```
## Not run:
library(caret)
data(cookfarm)
dat <- cookfarm[sample(1:nrow(cookfarm),500),]
indices <- CreateSpacetimeFolds(dat,"SOURCEID","Date")
ctrl <- caret::trainControl(method="cv",index = indices$index,savePredictions="final")
model <- caret::train(dat[,c("DEM","TWI","BLD")],dat$VW, method="rf", trControl=ctrl, ntree=10)
global_validation(model)

## End(Not run)
```

---

knndm

*K-fold Nearest Neighbour Distance Matching*

---

**Description**

This function implements the kNNDM algorithm and returns the necessary indices to perform a k-fold NNDM CV for map validation.

**Usage**

```
knndm(
  tpoints,
  modeldomain = NULL,
  predpoints = NULL,
  space = "geographical",
  k = 10,
  maxp = 0.5,
  clustering = "hierarchical",
  linkf = "ward.D2",
  samplesize = 1000,
  sampling = "regular",
  useMD = FALSE
)
```



**Arguments**

tpoints	sf or sfc point object, or data.frame if space = "feature". Contains the training points samples.
modeldomain	sf polygon object or SpatRaster defining the prediction area. Optional; alternative to predpoints (see Details).
predpoints	sf or sfc point object, or data.frame if space = "feature". Contains the target prediction points. Optional; alternative to modeldomain (see Details).
space	character. Either "geographical" or "feature".
k	integer. Number of folds desired for CV. Defaults to 10.
maxp	numeric. Maximum fold size allowed, defaults to 0.5, i.e. a single fold can hold a maximum of half of the training points.
clustering	character. Possible values include "hierarchical" and "kmeans". See details.
linkf	character. Only relevant if clustering = "hierarchical". Link function for agglomerative hierarchical clustering. Defaults to "ward.D2". Check 'stats::hclust' for other options.
samplesize	numeric. How many points in the modeldomain should be sampled as prediction points? Only required if modeldomain is used instead of predpoints.
sampling	character. How to draw prediction points from the modeldomain? See 'sf::st_sample'. Only required if modeldomain is used instead of predpoints.
useMD	boolean. Only for 'space'='feature': shall the Mahalanobis distance be calculated instead of Euclidean? Only works with numerical variables.

**Details**

knndm is a k-fold version of NNDM LOO CV for medium and large datasets. Briefly, the algorithm tries to find a k-fold configuration such that the integral of the absolute differences (Wasserstein  $W$  statistic) between the empirical nearest neighbour distance distribution function between the test and training data during CV ( $G_j^*$ ), and the empirical nearest neighbour distance distribution function between the prediction and training points ( $G_{ij}$ ), is minimised. It does so by performing clustering of the training points' coordinates for different numbers of clusters that range from  $k$  to  $N$  (number of observations), merging them into  $k$  final folds, and selecting the configuration with the lowest  $W$ .

Using a projected CRS in 'knndm' has large computational advantages since fast nearest neighbour search can be done via the 'FNN' package, while working with geographic coordinates requires computing the full spherical distance matrices. As a clustering algorithm, 'kmeans' can only be used for projected CRS while 'hierarchical' can work with both projected and geographical coordinates, though it requires calculating the full distance matrix of the training points even for a projected CRS.

In order to select between clustering algorithms and number of folds 'k', different 'knndm' configurations can be run and compared, being the one with a lower  $W$  statistic the one that offers a better match.  $W$  statistics between 'knndm' runs are comparable as long as 'tpoints' and 'predpoints' or 'modeldomain' stay the same.

Map validation using 'knndm' should be used using 'CAST::global\_validation', i.e. by stacking all out-of-sample predictions and evaluating them all at once. The reasons behind this are 1) The

resulting folds can be unbalanced and 2) nearest neighbour functions are constructed and matched using all CV folds simultaneously.

If training data points are very clustered with respect to the prediction area and the presented ‘knndm’ configuration still show signs of  $G_j^* > G_{ij}$ , there are several things that can be tried. First, increase the ‘maxp’ parameter; this may help to control for strong clustering (at the cost of having unbalanced folds). Secondly, decrease the number of final folds ‘k’, which may help to have larger clusters.

The ‘modeldomain’ is either a sf polygon that defines the prediction area, or alternatively a SpatRaster out of which a polygon, transformed into the CRS of the training points, is defined as the outline of all non-NA cells. Then, the function takes a regular point sample (amount defined by ‘samplesize’) from the spatial extent. As an alternative use ‘predpoints’ instead of ‘modeldomain’, if you have already defined the prediction locations (e.g. raster pixel centroids). When using either ‘modeldomain’ or ‘predpoints’, we advise to plot the study area polygon and the training/prediction points as a previous step to ensure they are aligned.

‘knndm’ can also be performed in the feature space by setting ‘space’ to “feature”. Euclidean distances or Mahalanobis distances can be used for distance calculation, but only Euclidean are tested. In this case, nearest neighbour distances are calculated in n-dimensional feature space rather than in geographical space. ‘tpoints’ and ‘predpoints’ can be data frames or sf objects containing the values of the features. Note that the names of ‘tpoints’ and ‘predpoints’ must be the same. ‘predpoints’ can also be missing, if ‘modeldomain’ is of class SpatRaster. In this case, the values of the SpatRaster will be extracted to the ‘predpoints’. In the case of any categorical features, Gower distances will be used to calculate the Nearest Neighbour distances [Experimental]. If categorical features are present, and ‘clustering’ = “kmeans”, K-Prototype clustering will be performed instead.

### Value

An object of class *knndm* consisting of a list of eight elements: `indx_train`, `indx_test` (indices of the observations to use as training/test data in each kNNDM CV iteration), `Gij` (distances for G function construction between prediction and target points), `Gj` (distances for G function construction during LOO CV), `Gjstar` (distances for modified G function during kNNDM CV), `clusters` (list of cluster IDs), `W` (Wasserstein statistic), and `space` (stated by the user in the function call).

### Author(s)

Carles Milà and Jan Linnenbrink

### References

- Linnenbrink, J., Milà, C., Ludwig, M., and Meyer, H.: kNNDM: k-fold Nearest Neighbour Distance Matching Cross-Validation for map accuracy estimation, EGUSphere [preprint], <https://doi.org/10.5194/egusphere-2023-1308>, 2023.
- Milà, C., Mateu, J., Pebesma, E., Meyer, H. (2022): Nearest Neighbour Distance Matching Leave-One-Out Cross-Validation for map validation. *Methods in Ecology and Evolution* 00, 1– 13.

### See Also

[geodist](#), [nndm](#)

## Examples

```
#####
# Example 1: Simulated data - Randomly-distributed training points
#####

library(sf)
library(ggplot2)

# Simulate 1000 random training points in a 100x100 square
set.seed(1234)
simarea <- list(matrix(c(0,0,0,100,100,100,100,0,0,0), ncol=2, byrow=TRUE))
simarea <- sf::st_polygon(simarea)
train_points <- sf::st_sample(simarea, 1000, type = "random")
pred_points <- sf::st_sample(simarea, 1000, type = "regular")
plot(simarea)
plot(pred_points, add = TRUE, col = "blue")
plot(train_points, add = TRUE, col = "red")

# Run kNNDM for the whole domain, here the prediction points are known.
knndm_folds <- knndm(train_points, predpoints = pred_points, k = 5)
knndm_folds
plot(knndm_folds)
plot(knndm_folds, type = "simple") # For more accessible legend labels
plot(knndm_folds, type = "simple", stat = "density") # To visualize densities rather than ECDFs
folds <- as.character(knndm_folds$clusters)
ggplot() +
  geom_sf(data = simarea, alpha = 0) +
  geom_sf(data = train_points, aes(col = folds))

#####
# Example 2: Simulated data - Clustered training points
#####
## Not run:
library(sf)
library(ggplot2)

# Simulate 1000 clustered training points in a 100x100 square
set.seed(1234)
simarea <- list(matrix(c(0,0,0,100,100,100,100,0,0,0), ncol=2, byrow=TRUE))
simarea <- sf::st_polygon(simarea)
train_points <- clustered_sample(simarea, 1000, 50, 5)
pred_points <- sf::st_sample(simarea, 1000, type = "regular")
plot(simarea)
plot(pred_points, add = TRUE, col = "blue")
plot(train_points, add = TRUE, col = "red")

# Run kNNDM for the whole domain, here the prediction points are known.
knndm_folds <- knndm(train_points, predpoints = pred_points, k = 5)
knndm_folds
plot(knndm_folds)
plot(knndm_folds, type = "simple") # For more accessible legend labels
plot(knndm_folds, type = "simple", stat = "density") # To visualize densities rather than ECDFs
```

```

folds <- as.character(knndm_folds$clusters)
ggplot() +
  geom_sf(data = simarea, alpha = 0) +
  geom_sf(data = train_points, aes(col = folds))

## End(Not run)
#####
# Example 3: Real- world example; using a modeldomain instead of previously
# sampled prediction locations
#####
## Not run:
library(sf)
library(terra)
library(ggplot2)

### prepare sample data:
data(cookfarm)
dat <- aggregate(cookfarm[,c("DEM", "TWI", "NDRE.M", "Easting", "Northing", "VW")],
  by=list(as.character(cookfarm$SOURCEID)), mean)
pts <- dat[,-1]
pts <- st_as_sf(pts, coords=c("Easting", "Northing"))
st_crs(pts) <- 26911
studyArea <- rast(system.file("extdata", "predictors_2012-03-25.tif", package="CAST"))
pts <- st_transform(pts, crs = st_crs(studyArea))
terra::plot(studyArea[["DEM"]])
terra::plot(vect(pts), add = T)

knndm_folds <- knndm(pts, modeldomain=studyArea, k = 5)
knndm_folds
plot(knndm_folds)
folds <- as.character(knndm_folds$clusters)
ggplot() +
  geom_sf(data = pts, aes(col = folds))

#use for cross-validation:
library(caret)
ctrl <- trainControl(method="cv",
  index=knndm_folds$indx_train,
  savePredictions='final')
model_knndm <- train(dat[,c("DEM", "TWI", "NDRE.M")],
  dat$VW,
  method="rf",
  trControl = ctrl)
global_validation(model_knndm)

## End(Not run)
#####
# Example 4: Real- world example; kNNDM in feature space
#####
## Not run:
library(sf)
library(terra)
library(ggplot2)

```

```

data(splotdata)
splotdata <- splotdata[splotdata$Country == "Chile",]

predictors <- c("bio_1", "bio_4", "bio_5", "bio_6",
               "bio_8", "bio_9", "bio_12", "bio_13",
               "bio_14", "bio_15", "elev")

trainDat <- sf::st_drop_geometry(splotdata)
predictors_sp <- terra::rast(system.file("extdata", "predictors_chile.tif", package="CAST"))

terra::plot(predictors_sp[["bio_1"]])
terra::plot(vect(splotdata), add = T)

knndm_folds <- knndm(trainDat[,predictors], modeldomain = predictors_sp, space = "feature",
                    clustering="kmeans", k=4, maxp=0.8)
plot(knndm_folds)

## End(Not run)

```

---

nndm

*Nearest Neighbour Distance Matching (NNDM) algorithm*


---

### Description

This function implements the NNDM algorithm and returns the necessary indices to perform a NNDM LOO CV for map validation.

### Usage

```

nndm(
  tpoints,
  modeldomain = NULL,
  predpoints = NULL,
  space = "geographical",
  samplesize = 1000,
  sampling = "regular",
  phi = "max",
  min_train = 0.5
)

```

### Arguments

tpoints	sf or sfc point object, or data.frame if space = "feature". Contains the training points samples.
modeldomain	sf polygon object or SpatRaster defining the prediction area. Optional; alternative to predpoints (see Details).

predpoints	sf or sfc point object, or data.frame if space = "feature". Contains the target prediction points. Optional; alternative to modeldomain (see Details).
space	character. Either "geographical" or "feature". Feature space is still experimental, so use with caution.
samplesize	numeric. How many points in the modeldomain should be sampled as prediction points? Only required if modeldomain is used instead of predpoints.
sampling	character. How to draw prediction points from the modeldomain? See 'sf::st_sample'. Only required if modeldomain is used instead of predpoints.
phi	Numeric. Estimate of the landscape autocorrelation range in the same units as the tpoints and predpoints for projected CRS, in meters for geographic CRS. Per default (phi="max"), the maximum distance found in the training and prediction points is used. See Details.
min_train	Numeric between 0 and 1. Minimum proportion of training data that must be used in each CV fold. Defaults to 0.5 (i.e. half of the training points).

### Details

NNDM proposes a LOO CV scheme such that the nearest neighbour distance distribution function between the test and training data during the CV process is matched to the nearest neighbour distance distribution function between the prediction and training points. Details of the method can be found in Milà et al. (2022).

Specifying *phi* allows limiting distance matching to the area where this is assumed to be relevant due to spatial autocorrelation. Distances are only matched up to *phi*. Beyond that range, all data points are used for training, without exclusions. When *phi* is set to "max", nearest neighbor distance matching is performed for the entire prediction area. Euclidean distances are used for projected and non-defined CRS, great circle distances are used for geographic CRS (units in meters).

The *modeldomain* is either a sf polygon that defines the prediction area, or alternatively a SpatRaster out of which a polygon, transformed into the CRS of the training points, is defined as the outline of all non-NA cells. Then, the function takes a regular point sample (amount defined by *samplesize*) from the spatial extent. As an alternative use *predpoints* instead of *modeldomain*, if you have already defined the prediction locations (e.g. raster pixel centroids). When using either *modeldomain* or *predpoints*, we advise to plot the study area polygon and the training/prediction points as a previous step to ensure they are aligned.

### Value

An object of class *nndm* consisting of a list of six elements: *indx\_train*, *indx\_test*, and *indx\_exclude* (indices of the observations to use as training/test/excluded data in each NNDM LOO CV iteration), *Gij* (distances for G function construction between prediction and target points), *Gj* (distances for G function construction during LOO CV), *Gjstar* (distances for modified G function during NNDM LOO CV), *phi* (landscape autocorrelation range). *indx\_train* and *indx\_test* can directly be used as "index" and "indexOut" in caret's [trainControl](#) function or used to initiate a custom validation strategy in mlr3.

### Note

NNDM is a variation of LOOCV and therefore may take a long time for large training data sets. See [knndm](#) for a more efficient k-fold variant of the method.

**Author(s)**

Carles Milà

**References**

- Milà, C., Mateu, J., Pebesma, E., Meyer, H. (2022): Nearest Neighbour Distance Matching Leave-One-Out Cross-Validation for map validation. *Methods in Ecology and Evolution* 00, 1– 13.
- Meyer, H., Pebesma, E. (2022): Machine learning-based global maps of ecological variables and the challenge of assessing them. *Nature Communications*. 13.

**See Also**[geodist](#), [knndm](#)**Examples**

```
#####
# Example 1: Simulated data - Randomly-distributed training points
#####

library(sf)

# Simulate 100 random training points in a 100x100 square
set.seed(123)
poly <- list(matrix(c(0,0,0,100,100,100,100,0,0,0), ncol=2, byrow=TRUE))
sample_poly <- sf::st_polygon(poly)
train_points <- sf::st_sample(sample_poly, 100, type = "random")
pred_points <- sf::st_sample(sample_poly, 100, type = "regular")
plot(sample_poly)
plot(pred_points, add = TRUE, col = "blue")
plot(train_points, add = TRUE, col = "red")

# Run NNDM for the whole domain, here the prediction points are known
nndm_pred <- nndm(train_points, predpoints=pred_points)
nndm_pred
plot(nndm_pred)
plot(nndm_pred, type = "simple") # For more accessible legend labels

# ...or run NNDM with a known autocorrelation range of 10
# to restrict the matching to distances lower than that.
nndm_pred <- nndm(train_points, predpoints=pred_points, phi = 10)
nndm_pred
plot(nndm_pred)

#####
# Example 2: Simulated data - Clustered training points
#####

library(sf)
```

```

# Simulate 100 clustered training points in a 100x100 square
set.seed(123)
poly <- list(matrix(c(0,0,0,100,100,100,100,0,0,0), ncol=2, byrow=TRUE))
sample_poly <- sf::st_polygon(poly)
train_points <- clustered_sample(sample_poly, 100, 10, 5)
pred_points <- sf::st_sample(sample_poly, 100, type = "regular")
plot(sample_poly)
plot(pred_points, add = TRUE, col = "blue")
plot(train_points, add = TRUE, col = "red")

# Run NNDM for the whole domain
nndm_pred <- nndm(train_points, predpoints=pred_points)
nndm_pred
plot(nndm_pred)
plot(nndm_pred, type = "simple") # For more accessible legend labels

#####
# Example 3: Real- world example; using a SpatRast modeldomain instead
# of previously sampled prediction locations
#####
## Not run:
library(sf)
library(terra)

### prepare sample data:
data(cookfarm)
dat <- aggregate(cookfarm[,c("DEM", "TWI", "NDRE.M", "Easting", "Northing", "VW")],
  by=list(as.character(cookfarm$SOURCEID)), mean)
pts <- dat[,-1]
pts <- st_as_sf(pts, coords=c("Easting", "Northing"))
st_crs(pts) <- 26911
studyArea <- rast(system.file("extdata", "predictors_2012-03-25.tif", package="CAST"))
pts <- st_transform(pts, crs = st_crs(studyArea))
terra::plot(studyArea[["DEM"]])
terra::plot(vect(pts), add = T)

nndm_folds <- nndm(pts, modeldomain = studyArea)
plot(nndm_folds)

#use for cross-validation:
library(caret)
ctrl <- trainControl(method="cv",
  index=nndm_folds$indx_train,
  indexOut=nndm_folds$indx_test,
  savePredictions='final')
model_nndm <- train(dat[,c("DEM", "TWI", "NDRE.M")],
  dat$VW,
  method="rf",
  trControl = ctrl)
global_validation(model_nndm)

## End(Not run)

```



```
#####
# Example 4: Real- world example; nndm in feature space
#####
## Not run:
library(sf)
library(terra)
library(ggplot2)

# Prepare the plot dataset for Chile
data(splotdata)
splotdata <- splotdata[splotdata$Country == "Chile",]

# Select a series of bioclimatic predictors
predictors <- c("bio_1", "bio_4", "bio_5", "bio_6",
               "bio_8", "bio_9", "bio_12", "bio_13",
               "bio_14", "bio_15", "elev")

predictors_sp <- terra::rast(system.file("extdata", "predictors_chile.tif", package="CAST"))

# Data visualization
terra::plot(predictors_sp[["bio_1"]])
terra::plot(vect(splotdata), add = T)

# Run and visualise the nndm results
nndm_folds <- nndm(splotdata[,predictors], modeldomain = predictors_sp, space = "feature")
plot(nndm_folds)

#use for cross-validation:
library(caret)
ctrl <- trainControl(method="cv",
                     index=nndm_folds$indx_train,
                     indexOut=nndm_folds$indx_test,
                     savePredictions='final')
model_nndm <- train(st_drop_geometry(splotdata[,predictors]),
                   splotdata$Species_richness,
                   method="rf",
                   trControl = ctrl)
global_validation(model_nndm)

## End(Not run)
```

---

normalize\_DI

*Normalize DI values*


---

### Description

The DI is normalized by the DI threshold to allow for a more straightforward interpretation. A value in the resulting DI larger 1 means that the data are more dissimilar than what has been observed during cross-validation. The returned threshold is adjusted accordingly and is, as a consequence, 1.

**Usage**

```
normalize_DI(AOA)
```

**Arguments**

AOA                    An AOA object

**Value**

An object of class `aoa`

**See Also**

[aoa](#)

**Examples**

```
## Not run:
library(sf)
library(terra)
library(caret)

# prepare sample data:
data(cookfarm)
dat <- aggregate(cookfarm[,c("VW", "Easting", "Northing")],
  by=list(as.character(cookfarm$SOURCEID), mean)
)
pts <- st_as_sf(dat, coords=c("Easting", "Northing"))
pts$ID <- 1:nrow(pts)
set.seed(100)
pts <- pts[1:30,]
studyArea <- rast(system.file("extdata", "predictors_2012-03-25.tif", package="CAST"))[[1:8]]
trainDat <- extract(studyArea, pts, na.rm=FALSE)
trainDat <- merge(trainDat, pts, by.x="ID", by.y="ID")

# train a model:
set.seed(100)
variables <- c("DEM", "NDRE.Sd", "TWI")
model <- train(trainDat[,which(names(trainDat)%in%variables)],
  trainDat$VW, method="rf", importance=TRUE, tuneLength=1,
  trControl=trainControl(method="cv", number=5, savePredictions=T))

#...then calculate the AOA of the trained model for the study area:
AOA <- aoa(studyArea, model)
plot(AOA)
plot(AOA$DI)

#... then normalize the DI
DI_norm <- normalize_DI(AOA)
plot(DI_norm)
plot(DI_norm$DI)
```

```
## End(Not run)
```

---

```
plot Plot CAST classes
```

---

### Description

Generic plot function for CAST Classes

A plotting function for a forward feature selection result. Each point is the mean performance of a model run. Error bars represent the standard errors from cross validation. Marked points show the best model from each number of variables until a further variable could not improve the results. If `type=="selected"`, the contribution of the selected variables to the model performance is shown.

Density plot of nearest neighbor distances in geographic space or feature space between training data as well as between training data and prediction locations. Optional, the nearest neighbor distances between training data and test data or between training data and CV iterations is shown. The plot can be used to check the suitability of a chosen CV method to be representative to estimate map accuracy.

Plot the DI/LPD and errormetric from Cross-Validation with the modeled relationship

### Usage

```
## S3 method for class 'trainDI'
plot(x, ...)

## S3 method for class 'aoa'
plot(x, samplesize = 1000, variable = "DI", ...)

## S3 method for class 'nndm'
plot(x, type = "strict", stat = "ecdf", ...)

## S3 method for class 'knndm'
plot(x, type = "strict", stat = "ecdf", ...)

## S3 method for class 'ffs'
plot(
  x,
  plotType = "all",
  palette = rainbow,
  reverse = FALSE,
  marker = "black",
  size = 1.5,
  lwd = 0.5,
  pch = 21,
  ...
)
```

```
## S3 method for class 'geodist'
plot(x, unit = "m", stat = "density", ...)
```

```
## S3 method for class 'errorModel'
plot(x, ...)
```

## Arguments

x	errorModel, see <a href="#">DitoErrormetric</a>
...	other params
samplesize	numeric. How many prediction samples should be plotted?
variable	character. Variable for which to generate the density plot. 'DI' or 'LPD'
type	String, defaults to "strict" to show the original nearest neighbour distance definitions in the legend. Alternatively, set to "simple" to have more intuitive labels.
stat	"density" for density plot or "ecdf" for empirical cumulative distribution function plot.
plotType	character. Either "all" or "selected"
palette	A color palette
reverse	Character. Should the palette be reversed?
marker	Character. Color to mark the best models
size	Numeric. Size of the points
lwd	Numeric. Width of the error bars
pch	Numeric. Type of point marking the best models
unit	character. Only if type=="geo" and only applied to the plot. Supported: "m" or "km".

## Value

```
a ggplot
a ggplot
```

## Author(s)

Marvin Ludwig, Hanna Meyer  
Carles Milà

## Examples

```
## Not run:
data(splotdata)
splotdata <- st_drop_geometry(splotdata)
ffsmodel <- ffs(splotdata[,6:16], splotdata$Species_richness, ntree = 10)
plot(ffsmodel)
#plot performance of selected variables only:
plot(ffsmodel,plotType="selected")

## End(Not run)
```

---

print	<i>Print CAST classes</i>
-------	---------------------------

---

**Description**

Generic print function for trainDI and aoa

**Usage**

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

show.trainDI(x, ...)

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

show.aoa(x, ...)

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

show.nndm(x, ...)

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

show.knndm(x, ...)

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

show.ffs(x, ...)
```

**Arguments**

x	An object of type <i>ffs</i>
...	other arguments.

---

splotdata	<i>sPlotOpen Data of Species Richness</i>
-----------	---

---

**Description**

sPlotOpen Species Richness for South America with associated predictors

**Usage**

```
data(splotdata)
```

**Format**

A sf points / data.frame with 703 rows and 17 columns:

**PlotObservationID, GIVD\_ID, Country, Biome** sPlotOpen Metadata

**Species\_richness** Response Variable - Plant species richness from sPlotOpen

**bio\_x, elev** Predictor Variables - Worldclim and SRTM elevation

**geometry** Lat/Lon

**Source**

- Plot with Species\_richness from [sPlotOpen](#)
- predictors acquired via R package [geodata](#)

**References**

- Sabatini, F. M. et al. sPlotOpen – An environmentally balanced, open-access, global dataset of vegetation plots. (2021). [doi:10.1111/geb.13346](https://doi.org/10.1111/geb.13346)
- Lopez-Gonzalez, G. et al. ForestPlots.net: a web application and research tool to manage and analyse tropical forest plot data: ForestPlots.net. *Journal of Vegetation Science* (2011).
- Pauchard, A. et al. Alien Plants Homogenise Protected Areas: Evidence from the Landscape and Regional Scales in South Central Chile. in *Plant Invasions in Protected Areas* (2013).
- Peyre, G. et al. VegPáramo, a flora and vegetation database for the Andean páramo. *phyto-coenologia* (2015).
- Vibrans, A. C. et al. Insights from a large-scale inventory in the southern Brazilian Atlantic Forest. *Scientia Agricola* (2020).

---

trainDI

*Calculate Dissimilarity Index of training data*

---

**Description**

This function estimates the Dissimilarity Index (DI) within the training data set used for a prediction model. Optionally, the local point density can also be calculated. Predictors can be weighted based on the internal variable importance of the machine learning algorithm used for model training.

**Usage**

```
trainDI(
  model = NA,
  train = NULL,
  variables = "all",
  weight = NA,
  CVtest = NULL,
  CVtrain = NULL,
  method = "L2",
  useWeight = TRUE,
  LPD = FALSE,
  verbose = TRUE
)
```

**Arguments**

model	A train object created with caret used to extract weights from (based on variable importance) as well as cross-validation folds
train	A data.frame containing the data used for model training. Only required when no model is given
variables	character vector of predictor variables. if "all" then all variables of the model are used or if no model is given then of the train dataset.
weight	A data.frame containing weights for each variable. Only required if no model is given.
CVtest	list or vector. Either a list where each element contains the data points used for testing during the cross validation iteration (i.e. held back data). Or a vector that contains the ID of the fold for each training point. Only required if no model is given.
CVtrain	list. Each element contains the data points used for training during the cross validation iteration (i.e. held back data). Only required if no model is given and only required if CVtrain is not the opposite of CVtest (i.e. if a data point is not used for testing, it is used for training). Relevant if some data points are excluded, e.g. when using <a href="#">nndm</a> .
method	Character. Method used for distance calculation. Currently euclidean distance (L2) and Mahalanobis distance (MD) are implemented but only L2 is tested. Note that MD takes considerably longer.
useWeight	Logical. Only if a model is given. Weight variables according to importance in the model?
LPD	Logical. Indicates whether the local point density should be calculated or not.
verbose	Logical. Print progress or not?

**Value**

A list of class trainDI containing:

train            A data frame containing the training data

weight	A data frame with weights based on the variable importance.
variables	Names of the used variables
catvars	Which variables are categorical
scaleparam	Scaling parameters. Output from scale
trainDist_avrg	A data frame with the average distance of each training point to every other point
trainDist_avrgmean	The mean of trainDist_avrg. Used for normalizing the DI
trainDI	Dissimilarity Index of the training data
threshold	The DI threshold used for inside/outside AOA
trainLPD	LPD of the training data
avrgLPD	Average LPD of the training data

**Note**

This function is called within [aoa](#) to estimate the DI and AOA of new data. However, it may also be used on its own if only the DI of training data is of interest, or to facilitate a parallelization of [aoa](#) by avoiding a repeated calculation of the DI within the training data.

**Author(s)**

Hanna Meyer, Marvin Ludwig, Fabian Schumacher

**References**

Meyer, H., Pebesma, E. (2021): Predicting into unknown space? Estimating the area of applicability of spatial prediction models. doi:[10.1111/2041210X.13650](https://doi.org/10.1111/2041210X.13650)

**See Also**

[aoa](#)

**Examples**

```
## Not run:
library(sf)
library(terra)
library(caret)
library(CAST)

# prepare sample data:
data("splotdata")
splotdata = st_drop_geometry(splotdata)

# train a model:
set.seed(100)
model <- caret::train(splotdata[,6:16],
                      splotdata$Species_richness,
                      importance=TRUE, tuneLength=1, ntree = 15, method = "rf",
```



```
trControl = trainControl(method="cv", number=5, savePredictions=T))
# variable importance is used for scaling predictors
plot(varImp(model, scale=FALSE))

# calculate the DI of the trained model:
DI = trainDI(model=model)
plot(DI)

#...or calculate the DI and LPD of the trained model:
# DI = trainDI(model=model, LPD = TRUE)

# the DI can now be used to compute the AOA (here with LPD):
studyArea = rast(system.file("extdata/predictors_chile.tif", package = "CAST"))
AOA = aoa(studyArea, model = model, trainDI = DI, LPD = TRUE, maxLPD = 1)
print(AOA)
plot(AOA)
plot(AOA$AOA)
plot(AOA$LPD)

## End(Not run)
```

# Index

- \* **datasets**
  - cookfarm, [11](#)
  - splotdata, [37](#)
- \* **package**
  - CAST, [8](#)
- aoa, [2](#), [14](#), [15](#), [34](#), [40](#)
- bss, [7](#), [18](#)
- CAST, [8](#)
- CAST-package (CAST), [8](#)
- clustered\_sample, [10](#)
- confusionMatrix, [24](#)
- cookfarm, [11](#)
- CreateSpacetimeFolds, [8](#), [11](#), [18](#), [24](#)
- DItoErrormetric, [4](#), [36](#)
- DItoErrormetric (errorProfiles), [13](#)
- errorProfiles, [5](#), [13](#)
- ffs, [8](#), [13](#), [16](#)
- geodist, [20](#), [26](#), [31](#)
- global\_validation, [7](#), [17](#), [23](#)
- knndm, [21](#), [24](#), [30](#), [31](#)
- mclapply, [17](#)
- nndm, [4](#), [8](#), [12](#), [13](#), [18](#), [21](#), [26](#), [29](#), [39](#)
- normalize\_DI, [4](#), [5](#), [33](#)
- plot, [35](#)
- postResample, [24](#)
- print, [37](#)
- rollapply, [14](#)
- show.aoa (print), [37](#)
- show.ffs (print), [37](#)
- show.knndm (print), [37](#)
- show.nndm (print), [37](#)
- show.trainDI (print), [37](#)
- splotdata, [37](#)
- spsample, [21](#)
- train, [7](#), [8](#), [16–18](#), [23](#)
- trainControl, [8](#), [13](#), [18](#), [30](#)
- trainDI, [3–5](#), [14](#), [38](#)