

Mapping the area of applicability - Case Study

Supplementary to the paper ‘Predicting into unknown space? Estimating the area of applicability of spatial prediction models’

Hanna Meyer

1/4/2020

Introduction

This script contains the complete code for the case study within the paper and creates the figures presented there. It can further be used to run experiments under different settings. Note that running the code takes a few minutes (depending on the number of training data and the size of the study area)!

Getting started

Major functionality needed is the ‘aoa’ function from the ‘CAST’ package that is doing the distance based estimation of the area of applicability. ‘Caret’ is needed for model training. The case study uses a simulated prediction task based on the ‘virtualspecies’ package.

```
rm(list=ls())
#install_github("HannaMeyer/CAST")
library(virtualspecies)
library(caret)
library(CAST)
library(viridis)
library(gridExtra)
library(knitr)
library(grid)
library(latticeExtra)
library(hydroGOF)
```

Settings for generating the predictors and response need to be defined, as well as the number of training data points and the seed used for all functions that involve randomness. The settings specified here are used for the case study published in the paper. Feel free to change them to see how things work under different scenarios!

Getting started

Note: To reproduce results used in the paper use the following settings:

- Default Case study with random data: npoints=50,design=“random”,meansPCA= c(3, -1), sdPCA=c(2, 2),simulateResponse=c(“bio2”,“bio5”,“bio10”, “bio13”,“bio14”,“bio19”),studyarea=c(-15, 65, 30, 75), seed=10

- Case Study with clustered data: npoints=500,nclusters=50,design="clustered", maxdist <- 0.6, meansPCA=c(3, -1), sdPCA=c(2, 2), simulateResponse=c("bio2","bio5","bio10", "bio13","bio14", "bio19"), studyarea=c(-15, 65, 30, 75), seed=10

Further interesting settings:

- Case Study biased with outlier: npoints=50,design="biasedWithOutlier",countries = c("Germany","Ireland","France", "Sweden"), countriesOutlier= "Turkmenistan" maxdist <- 0.6, meansPCA=c(3, -1), sdPCA=c(2, 2), simulateResponse=c("bio2","bio5","bio10", "bio13","bio14", "bio19"), studyarea=c(-15, 65, 30, 75), seed=10
- Case Study biased: same as biased with outlier but without the outlier

```
npoints <- 50 # number of training samples

design <- "random" # either clustered,biased,random, biasedWithOutlier

countries <- c("Germany", "Ireland", "France", "Sweden") #if design==biased
countriesOutlier <- "Turkmenistan" #if design==biasedWithOutlier A single point is set here
nclusters <- 50 #number of clusters if design==clustered
maxdist <- 0.6 #maxdist for clustered samples if design==clustered
meansPCA <- c(3, -1) # means of the gaussian response functions to the 2 axes
sdPCA <- c(2, 2) # sd's of the gaussian response functions to the 2 axes
simulateResponse <- c("bio2", "bio5", "bio10", "bio13",
                      "bio14", "bio19") # variables used to simulate the response
studyarea <- c(-15, 65, 30, 75) # extent of study area. Default: Europe
seed <- 10
multiCV <- TRUE
```

Get data

Bioclim data are downloaded and cropped to the study area.

```
predictors_global <- getData('worldclim', var='bio', res=10, path='../data/')
wp <- extent(studyarea)
predictors <- crop(predictors_global,wp)

#create a mask for land area:
mask <- predictors[[1]]
values(mask)[!is.na(values(mask))] <- 1
```

Generate Predictors and Response

The virtual response variable is created based on the PCA of a subset of bioclim predictors. See the virtualspecies package for further information.

```
response_vs <- generateSpFromPCA(predictors[[simulateResponse]],
                                    means = meansPCA, sds = sdPCA, plot=F)

response <- response_vs$suitab.raster
```

Simulate training points

To simulate field locations that are typically used as training data, “npoints” locations are randomly selected. If a clustered design is used, the “npoints” are distributed over “nclusters” with a maximum distance between each point of a cluster (maxdist, in degrees).

```
mask <- rasterToPolygons(mask,dissolve=TRUE)
set.seed(seed)
if (design=="clustered"){
  samplepoints <- csample(mask,npoints,nclusters,maxdist=maxdist,seed=seed)
}
if (design=="biased"){
  countryboundaries <- getData("countries", path='..../data/')
  countryboundaries <- countryboundaries[countryboundaries$NAME_ENGLISH%in%c(countries),]
  samplepoints <- spsample(countryboundaries,npoly, "random")
}
if (design=="biasedWithOutlier"){
  countryboundaries <- getData("countries", path='..../data/')
  countryboundariesOut <- countryboundaries[countryboundaries$NAME_ENGLISH%in%c(countriesOutlier),]
  countryboundaries <- countryboundaries[countryboundaries$NAME_ENGLISH%in%c(countries),]
  samplepoints <- spsample(countryboundaries,npoly, "random")
  samplepoints <- rbind(samplepoints,spsample(countryboundariesOut,1, "random"))
}
if (design=="random"){
  samplepoints <- spsample(mask,npoly, "random")
}
```

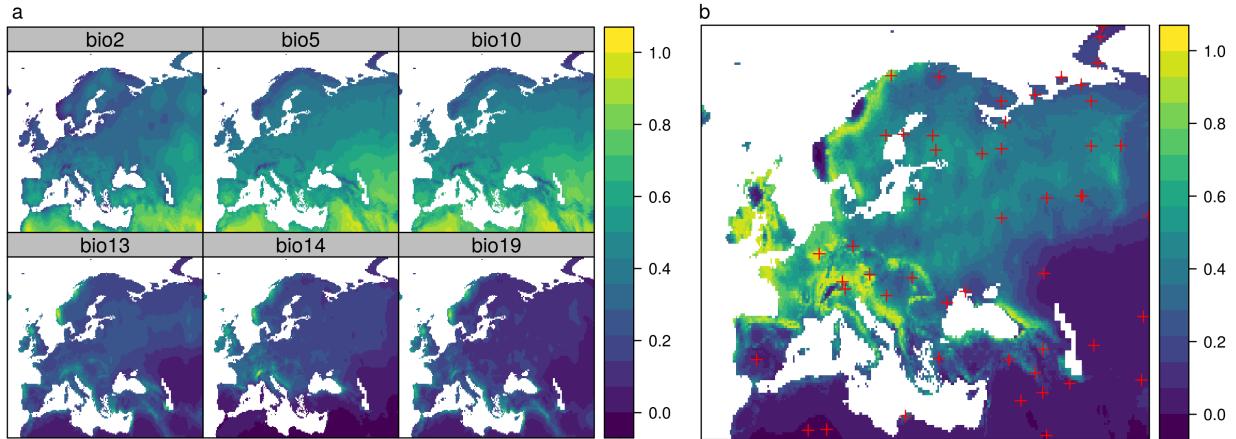


Figure 1: Subset of the predictors as well as the response variable and the selected training data points

Model training and prediction

Preparation

To prepare model training, predictor variables are extracted for the location of the selected sample data locations.

```

trainDat <- extract(predictors,samplepoints,df=TRUE)
trainDat$response <- extract (response,samplepoints)

if (design=="clustered"){
  trainDat <- merge(trainDat,samplepoints,by.x="ID",by.y="ID")
}

trainDat <- trainDat[complete.cases(trainDat),]

#pca_pred <- predict(response_vs$details$pca,trainDat[,response_vs$details$variables])

#axes <- response_vs$details$axes
#axes.to.plot <- axes[1:2]
#xmin <- min(response_vs$details$pca$li[, axes.to.plot[1]]) - 0.3 * diff(range(response_vs$details$pca$li[, axes.to.plot[1]]))
#xmax <- max(response_vs$details$pca$li[, axes.to.plot[1]])
#
# pdf("../figures/plotresponse.pdf",width=6,height=6)
# plotResponse(response_vs,no.plot.reset=T)
# points((pca_pred[,1]+10)/21.5,
#         pca_pred[,2],pch=16)
# dev.off()

```

Training and cross-validation

Model training is then done using the caret package. Note that other packages work as well as long as variable importance can be derived. The model output gives information on the general estimated model performance based on random cross validation.

```

set.seed(seed)
if(design!="clustered"){
model <- train(trainDat[,names(predictors)],
               trainDat$response,
               method="rf",
               importance=TRUE,
               tuneGrid = expand.grid(mtry = c(2:length(names(predictors)))),
               trControl = trainControl(method="cv",savePredictions = TRUE))

print("random cross-validation performance:")
print(model)
}

## [1] "random cross-validation performance:"
## Random Forest
##
## 50 samples
## 19 predictors
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 46, 45, 45, 45, 46, 46, ...
## Resampling results across tuning parameters:
##

```

```

##   mtry    RMSE      Rsquared     MAE
##   2    0.08312237  0.9449585  0.06038504
##   3    0.08134313  0.9491834  0.05831224
##   4    0.08270354  0.9500045  0.05912983
##   5    0.08255321  0.9448124  0.05969416
##   6    0.08354799  0.9466461  0.06076259
##   7    0.08258087  0.9483848  0.05983775
##   8    0.08358856  0.9414773  0.06077974
##   9    0.08285660  0.9471475  0.06024775
##  10   0.08549180  0.9417758  0.06160476
##  11   0.08714468  0.9394383  0.06346486
##  12   0.08701203  0.9336757  0.06361495
##  13   0.08738857  0.9354565  0.06383382
##  14   0.08733814  0.9365332  0.06300664
##  15   0.08859858  0.9324766  0.06452674
##  16   0.08794901  0.9354402  0.06460736
##  17   0.08750610  0.9333760  0.06407203
##  18   0.08907661  0.9320242  0.06453999
##  19   0.08606493  0.9362805  0.06282772
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was mtry = 3.
#if data are clustered, clustered CV is used:
if(design=="clustered"){
  folds <- CreateSpacetimeFolds(trainDat, spacevar="clstrID",k=nclusters)
  model <- train(trainDat[,names(predictors)],
                  trainDat$response,
                  method="rf",
                  importance=TRUE,
                  tuneGrid = expand.grid(mtry = c(2:length(names(predictors)))),
                  trControl = trainControl(method="cv",index=folds$index,savePredictions = TRUE))
  print("leave-cluster-out cross-validation performance:")
  print(model)

  model_random <- train(trainDat[,names(predictors)],
                         trainDat$response,
                         method="rf",
                         importance=TRUE,
                         tuneGrid = expand.grid(mtry = c(2:length(names(predictors)))),
                         trControl = trainControl(method="cv",savePredictions = TRUE))
  print("random cross-validation performance:")
  print(model_random)
}

```

Prediction and error calculation

The trained model is used to make predictions for the entire study area. The absolute error between prediction and reference is calculated for later comparison with the dissimilarity index.

```

prediction <- predict(predictors,model)
truediff <- abs(prediction-response)

```

Estimating the area of applicability

Variable importance to get weights

The variable importance from model training can be visualized to get information on how variable weighting will be done during to estimation of the dissimilarity index

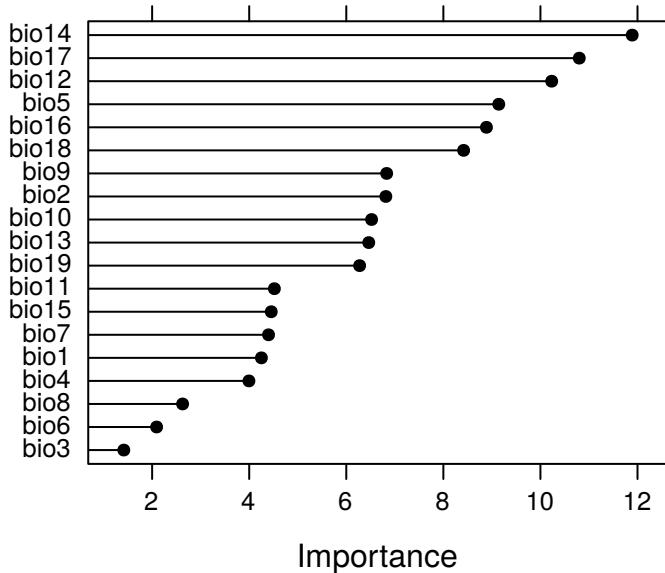


Figure 2: Variable importance

The area of applicability and the dissimilarity index are then calculated. First using weighted variables, and second (for comparison) without weighting. Everything is run in parallel to speed things up.

```
#with variable weighting:  
AOA <- aoa(predictors, model=model, returnTrainDI = TRUE)  
#without weighting:  
AOA_noWeights <- aoa(predictors, train=trainDat, variables = names(predictors))  
  
if (design=="clustered"){  
  AOA_random<- aoa(predictors, model_random)  
}  
  
preds <- model$pred[model$pred$mtry==model$bestTune$mtry,]  
absError <- abs(preds$pred[order(preds$rowIndex)]-preds$obs[order(preds$rowIndex)])  
  
predtrain <- predict(model$finalModel, trainDat, predict.all=TRUE)  
sdpred <- apply(predtrain$individual, 1, sd)
```

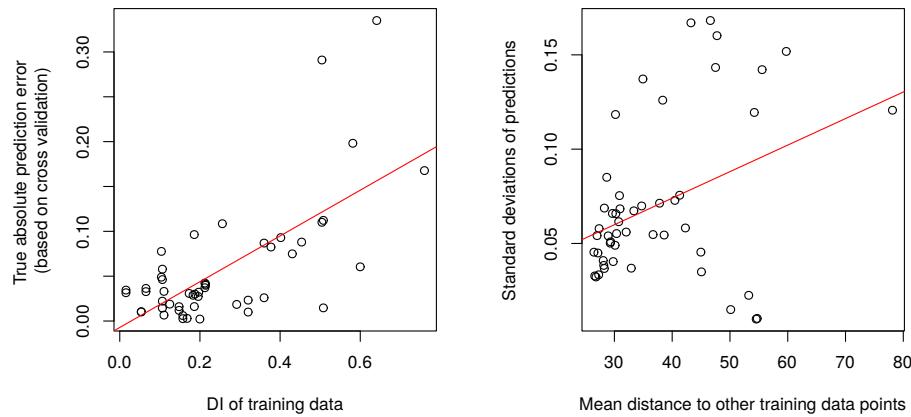


Figure 3: Relationship between the DI of the training data and the true prediction error (a) as well as the relationship between the mean distance to other data points (as a indicator for data point density) and the standard deviations of predictions (b).

Standard deviation from individual trees for comparison

For comparison to what is often used as uncertainty information, the standard deviations of the individual predictions from the 500 developed trees within the Random Forest model are calculated.

```
predsd <- RFsd(predictors,model)
```

Comparison

The dissimilarity index, as well as the standard deviations can then be compared to the true error.

```
compare <- stack(response,prediction,
                  predsd,trueendiff,
                  AOA$DI)
names(compare) <- c("response","prediction", "sd","true_diff","DI")
summary(values(compare))
```

```
##      response      prediction        sd      true_diff
##  Min.   :0.00   Min.   :0.00   Min.   :0.01   Min.   :0.00
##  1st Qu.:0.07   1st Qu.:0.14   1st Qu.:0.06   1st Qu.:0.01
##  Median :0.32   Median :0.34   Median :0.08   Median :0.04
##  Mean   :0.31   Mean   :0.32   Mean   :0.10   Mean   :0.06
##  3rd Qu.:0.45   3rd Qu.:0.44   3rd Qu.:0.15   3rd Qu.:0.07
##  Max.   :1.00   Max.   :0.81   Max.   :0.31   Max.   :0.76
##  NA's    :47804  NA's    :47804  NA's    :47804  NA's    :47804
##          DI
##  Min.   :0.00
##  1st Qu.:0.12
##  Median :0.21
##  Mean   :0.25
##  3rd Qu.:0.34
##  Max.   :2.90
##  NA's   :47804
```

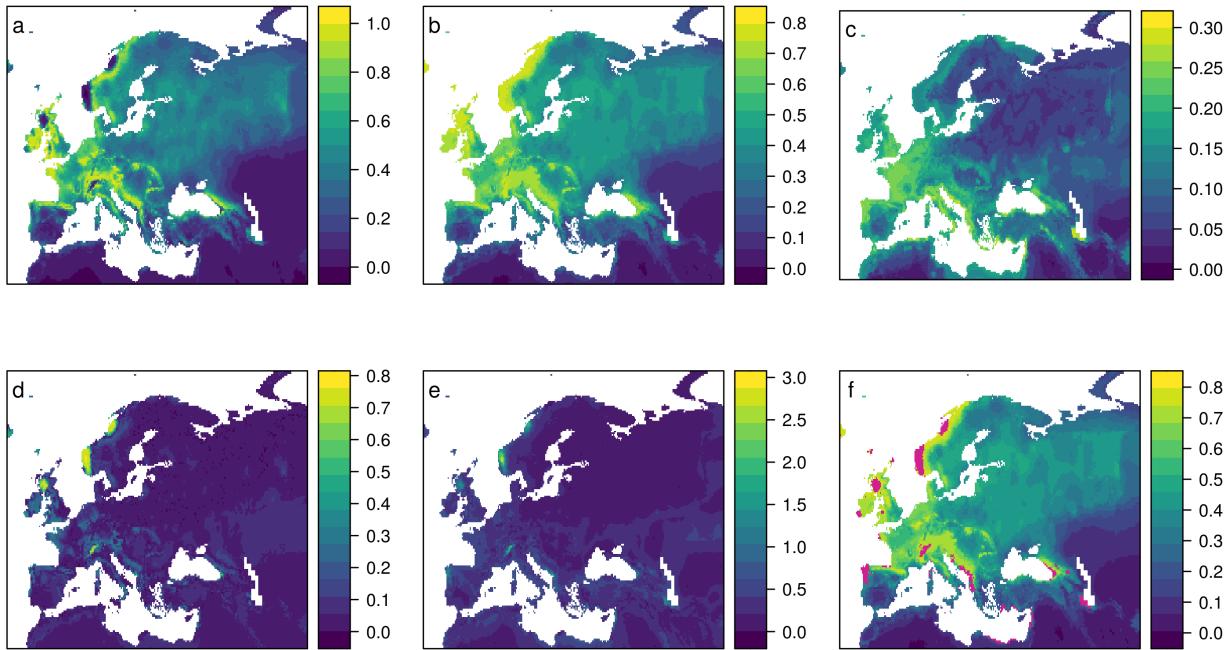


Figure 4: Comparison between reference (a), prediction (b), standard deviation of predictions (c), the true error (d), DI based on weighted variables (e), masked predictions (f)

Comparison between prediction, error, DI and cross validation

```
#DI with weights:
cor(values(truediff),values(AOA$DI),use="complete.obs")

## [1] 0.7149642

#DI no weights:
cor(values(truediff),values(AOA_noWeights$DI),use="complete.obs")

## [1] 0.6182329

#comparison prediction~ref
summary(lm(values(response)~values(prediction)))$r.squared

## [1] 0.8564091

cor(values(response),values(prediction),use = "complete.obs")

## [1] 0.9254238

rmse(values(response),values(prediction))

## [1] 0.09536494

#comparison prediction for the AOA-ref
print(attributes(AOA)$aoa_stats)

## $Mean_train
## [1] 37.39012
```

```

## $threshold_stats
##      25%      50%      75%      90%      95%      99%     100%
## 0.1102882 0.1916736 0.3599672 0.5079589 0.5915618 0.7017908 0.7601471
##
## $threshold
## [1] 0.6410527
##
## $lower_threshold
## [1] 0.01540481

predictionAOI <- prediction
values(predictionAOI)[values(AOA$AOA)==0] <- NA
paste0("R^2 AOA= ",summary(lm(values(response)~values(predictionAOI)))$r.squared)

## [1] "R^2 AOA= 0.938743450608877"
paste0("r AOA= ",cor(values(response),values(predictionAOI),use="complete.obs"))

## [1] "r AOA= 0.968887738909352"
paste0("RMSE AOA= ",rmse(values(response),values(predictionAOI)))

## [1] "RMSE AOA= 0.0711663273926504"
# ...and for outside the AOA
predictionNOTAOI <- prediction
values(predictionNOTAOI)[values(AOA$AOA)==1] <- NA
paste0("R^2 outside AOA= ",summary(lm(values(response)~values(predictionNOTAOI)))$r.squared)

## [1] "R^2 outside AOA= 1.7919287941121e-06"
paste0("r outside AOA= ",cor(values(response),values(predictionNOTAOI),use="complete.obs"))

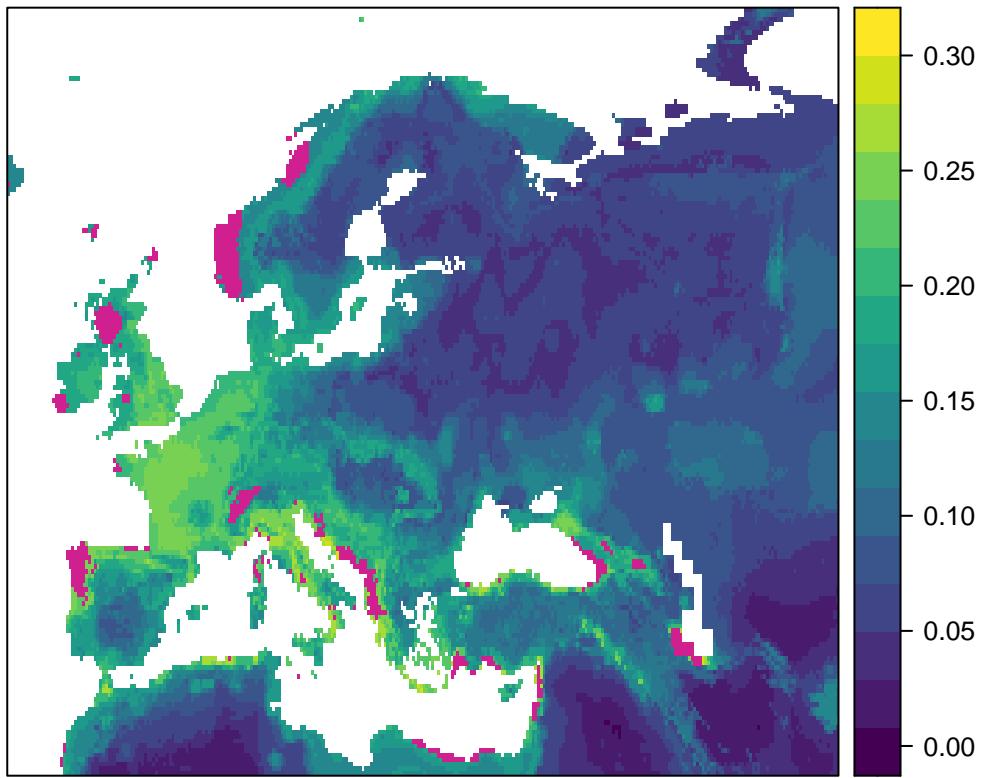
## [1] "RMSE outside AOA= -0.00133862944615461"
paste0("RMSE outside AOA= ",rmse(values(response),values(predictionNOTAOI)))

## [1] "RMSE outside AOA= 0.436650330543256"

```

Standard deviations of predictions within the AOA

```
spplot(compare$sd,col.regions=viridis(100))+spplot(AOA$AOA,col.regions=c("violetred","transparent"))
```



If applicable: Comparison random vs spatial CV and AOA estimation

#...if clustered compare to default AOA without taking #clusteredness into account #should be comparable to random CV error

```

if(design=="clustered"){
  prediction_random <- predict(predictors,model_random)

  predictionAOI <- prediction_random
  values(predictionAOI)[values(AOA_random$AOA)==0] <- NA
  print(paste0("R^2 AOA= ",summary(lm(values(response)~values(predictionAOI)))$r.squared))
  print(paste0("r AOA= ",cor(values(response),values(predictionAOI),use="complete.obs"))

  print(paste0("RMSE AOA= ", rmse(values(response),values(predictionAOI)))

  predictionNOTAOI <- prediction_random
  values(predictionNOTAOI)[values(AOA_random$AOA)==1] <- NA
  print(paste0("R^2 outside AOA= ",summary(lm(values(response)~values(predictionNOTAOI)))$r.squared))
  print(paste0("r outside AOA= ",cor(values(response),values(predictionNOTAOI),use="complete.obs")))
  print(paste0("RMSE outside AOA= ",rmse(values(response),values(predictionNOTAOI))))
}

```

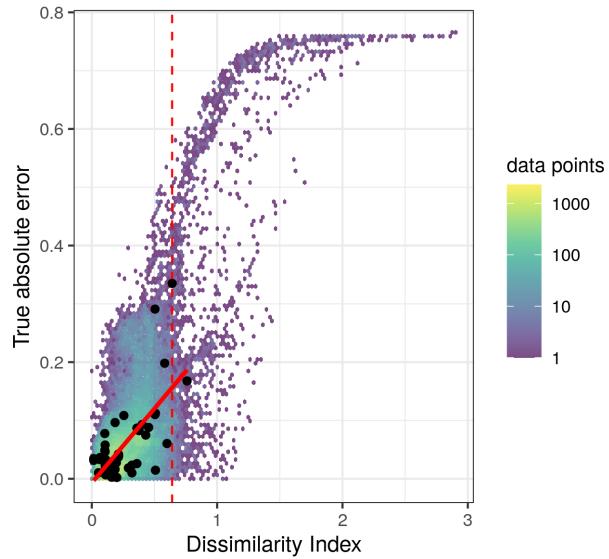


Figure 5: Relationship between the DI and the true error

DI as a quantitative uncertainty measure?

Relationship between the DI and the absolute error on a data-point-level

Relationship between the DI and the requested performance measure (here RMSE)

And the same but including re-fitting and validation with leave-cluster-in-predictorspace-out CV.

