



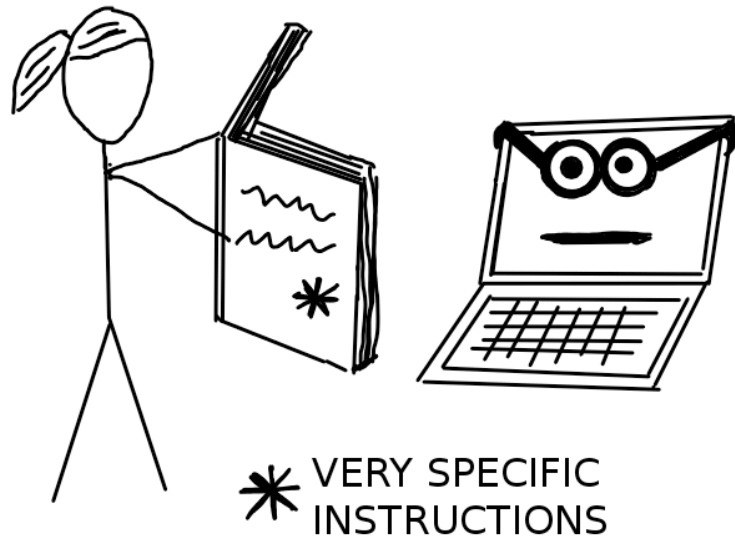
National
Oceanography
Centre

Machine Learning for feature detection

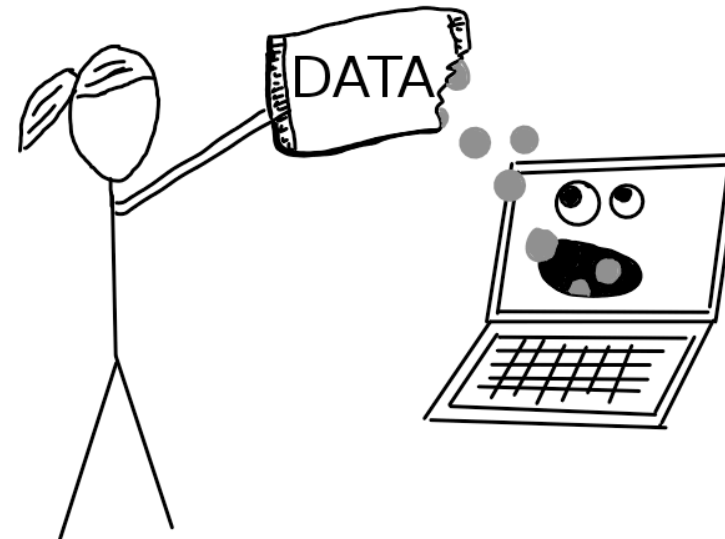
Linking zooplankton diversity and particulate organic carbon (POC)

A Machine Learning definition

Without Machine Learning



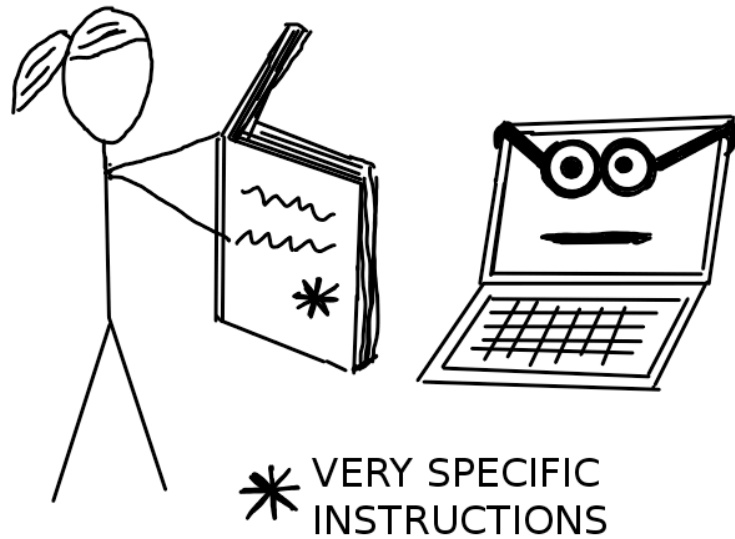
With Machine Learning



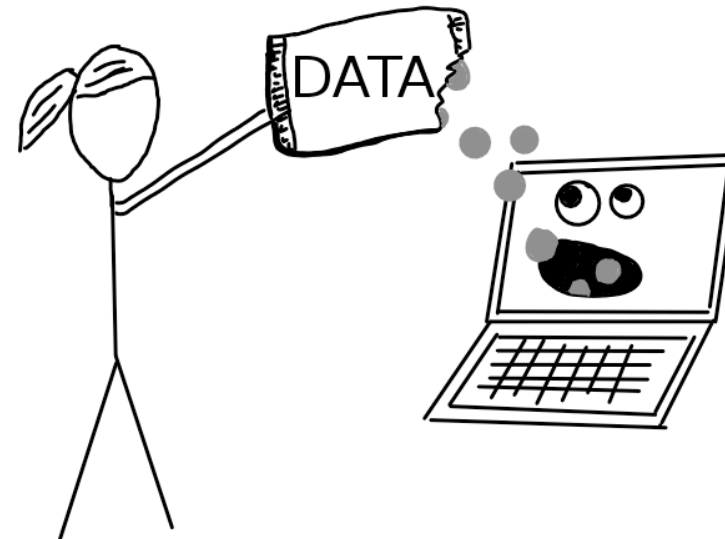
A Machine Learning definition

ML: finding patterns in data, without specific instruction, and possibly predicting outcome for new data

Without Machine Learning



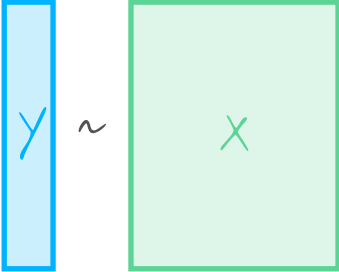
With Machine Learning



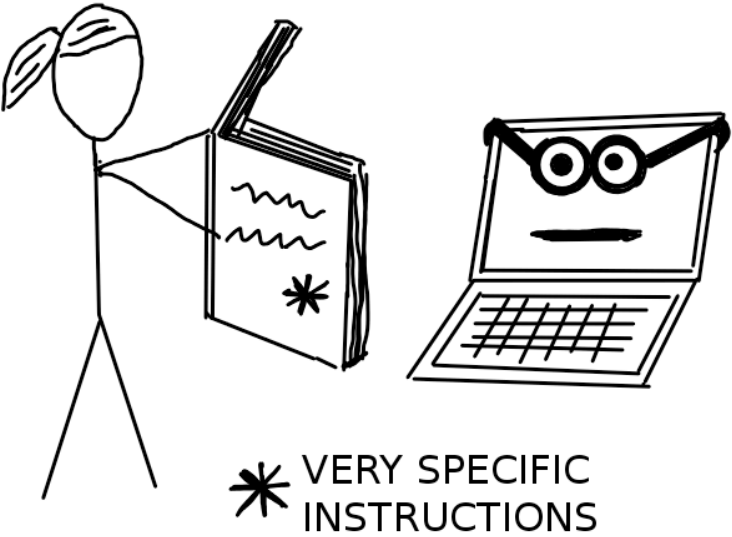
A Machine Learning definition

ML: finding patterns in data, without specific instruction, and possibly predicting outcome for new data

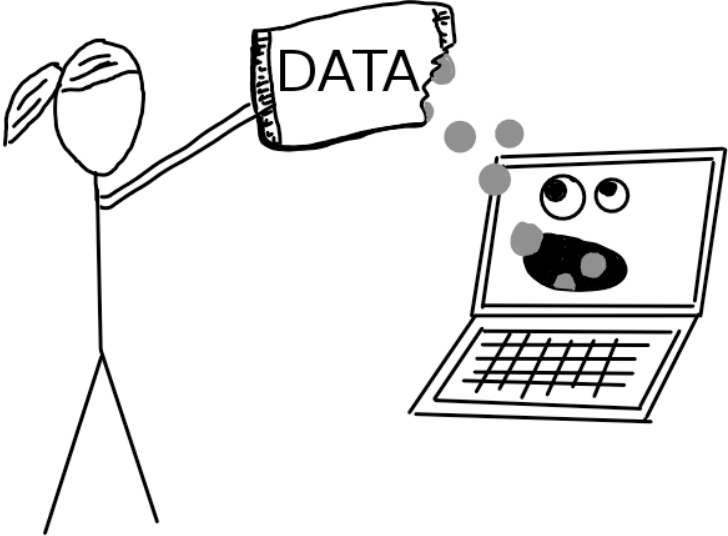
Supervised ML: relating **target response variable(s)** to **features predictors**



Without Machine Learning



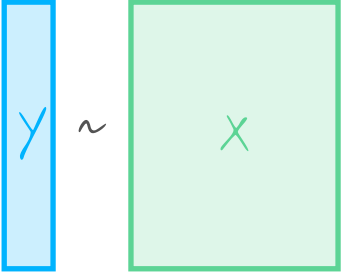
With Machine Learning



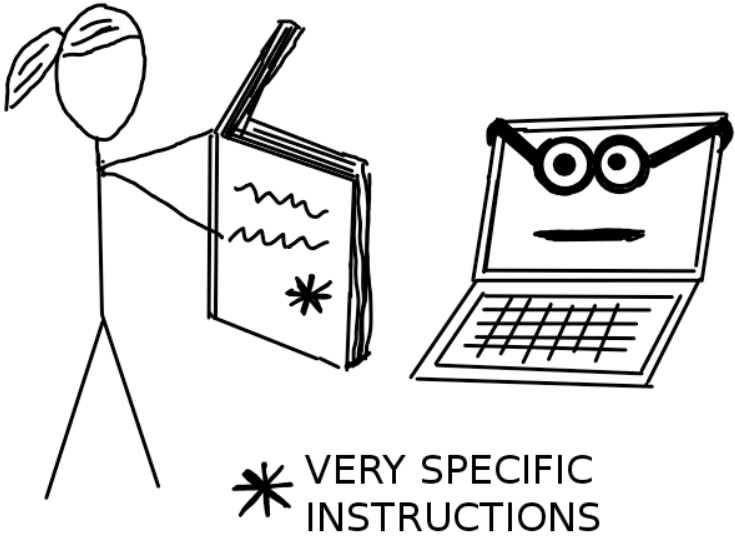
A Machine Learning definition

ML: finding patterns in data, without specific instruction, and possibly predicting outcome for new data

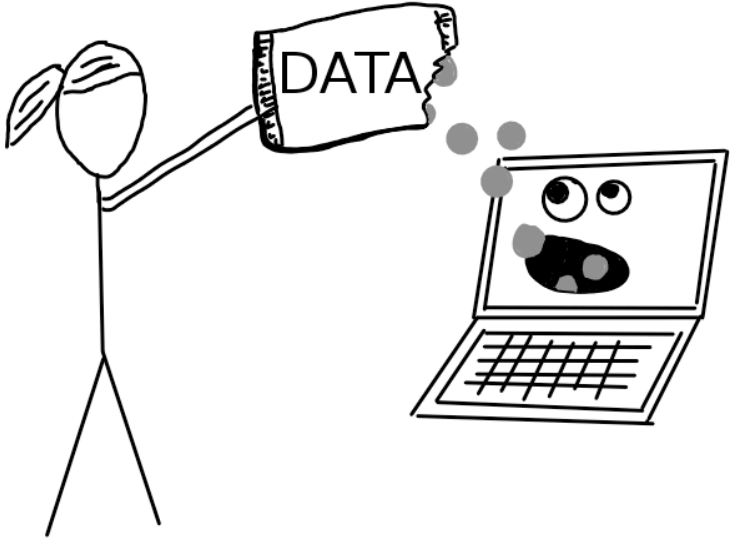
Supervised ML: relating **target response variable(s)** to **features predictors** importance?



Without Machine Learning



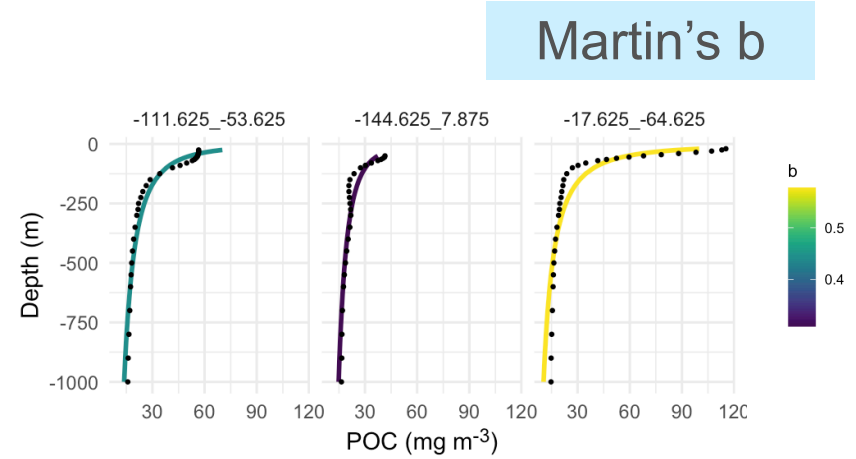
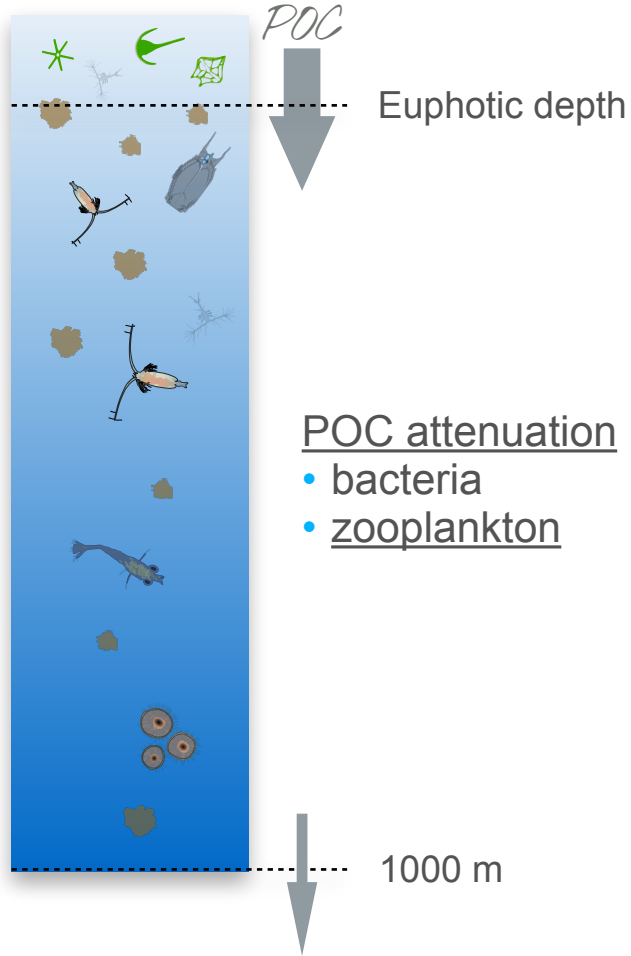
With Machine Learning



Response

Predictors

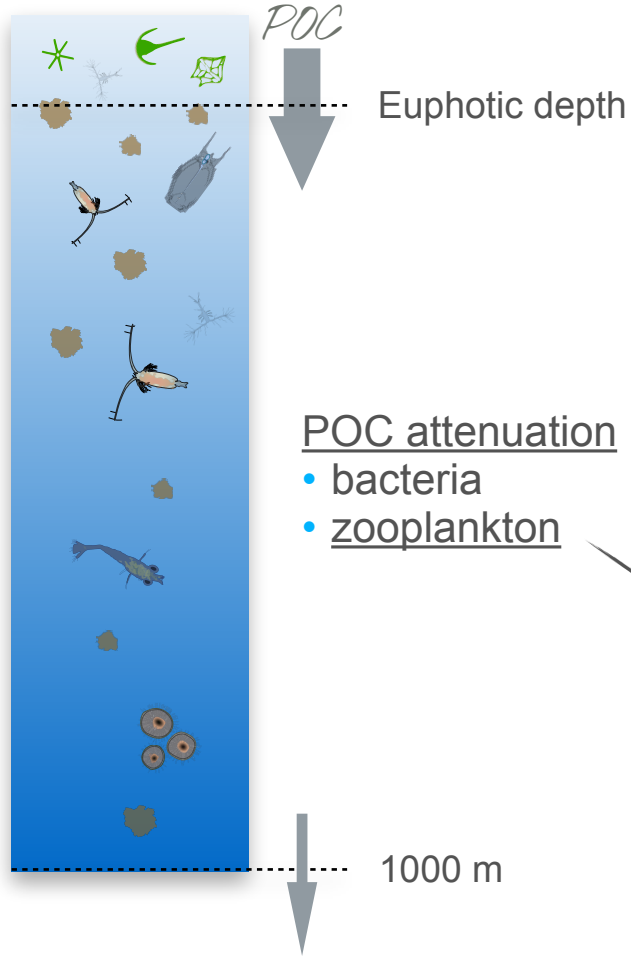
Context: relating POC attenuation to zooplankton diversity



Response

Predictors

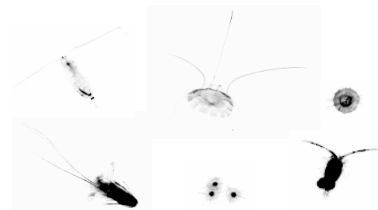
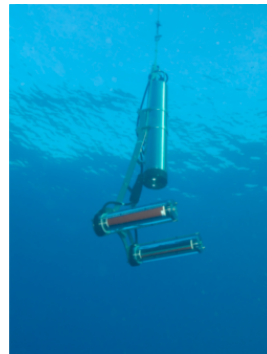
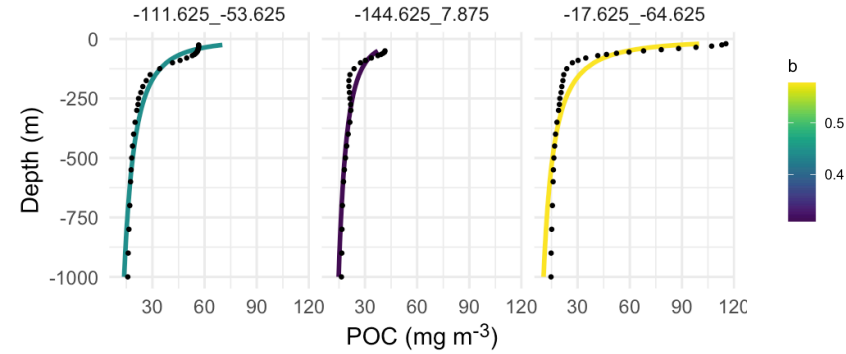
Context: relating POC attenuation to zooplankton diversity



POC attenuation

- bacteria
- zooplankton

Martin's b



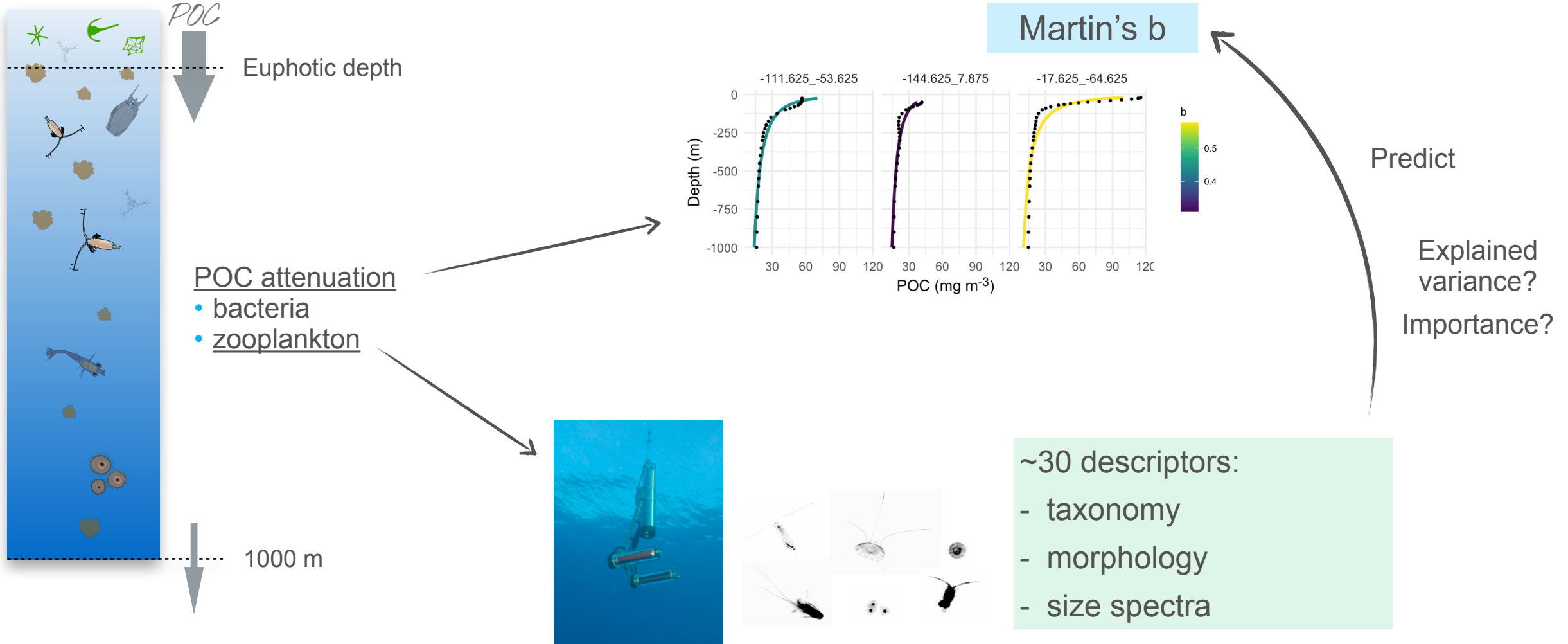
- ~30 descriptors:
- taxonomy
 - morphology
 - size spectra

Response

Predictors



Context: relating POC attenuation to zooplankton diversity



*Response**Predictors*

Context: relating **POC attenuation** to **zooplankton diversity**

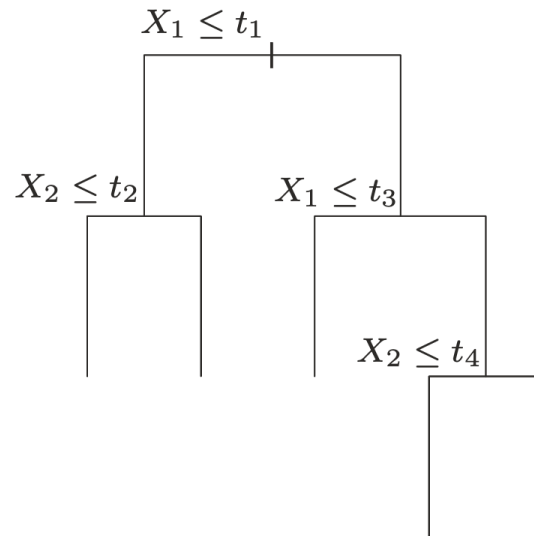


Boosted regression trees

Context: relating *Response* POC attenuation to *Predictors* zooplankton diversity

Boosted regression trees

Binary splits to relate
response to predictors



Context: relating **POC attenuation** to **zooplankton diversity**

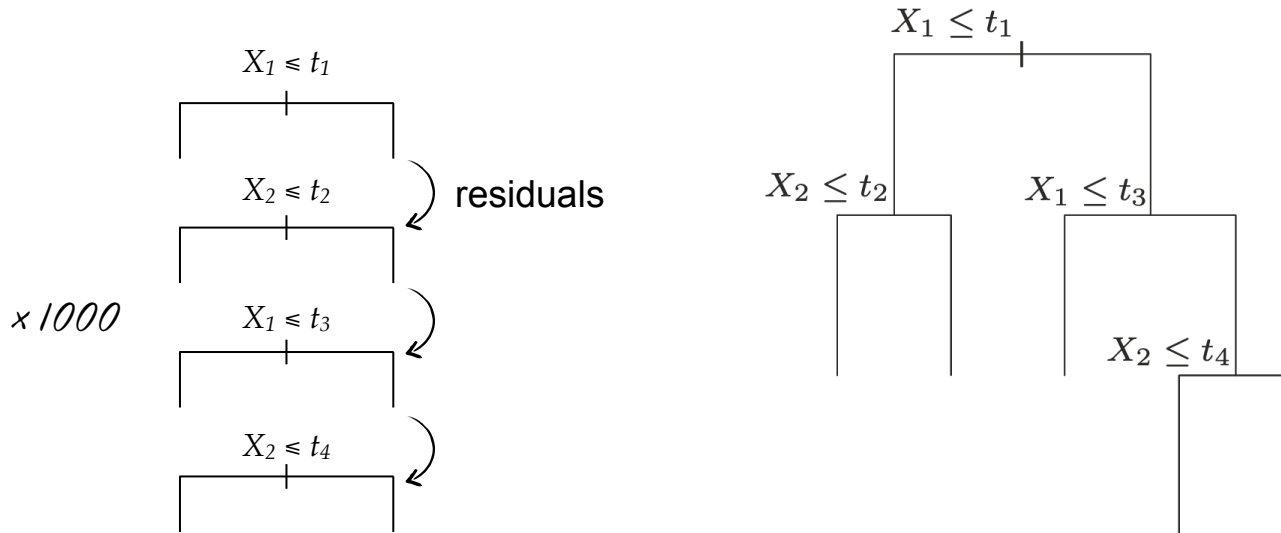
Response

Predictors

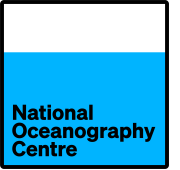
Boosted regression trees

Combining many
small models

Binary splits to relate
response to predictors



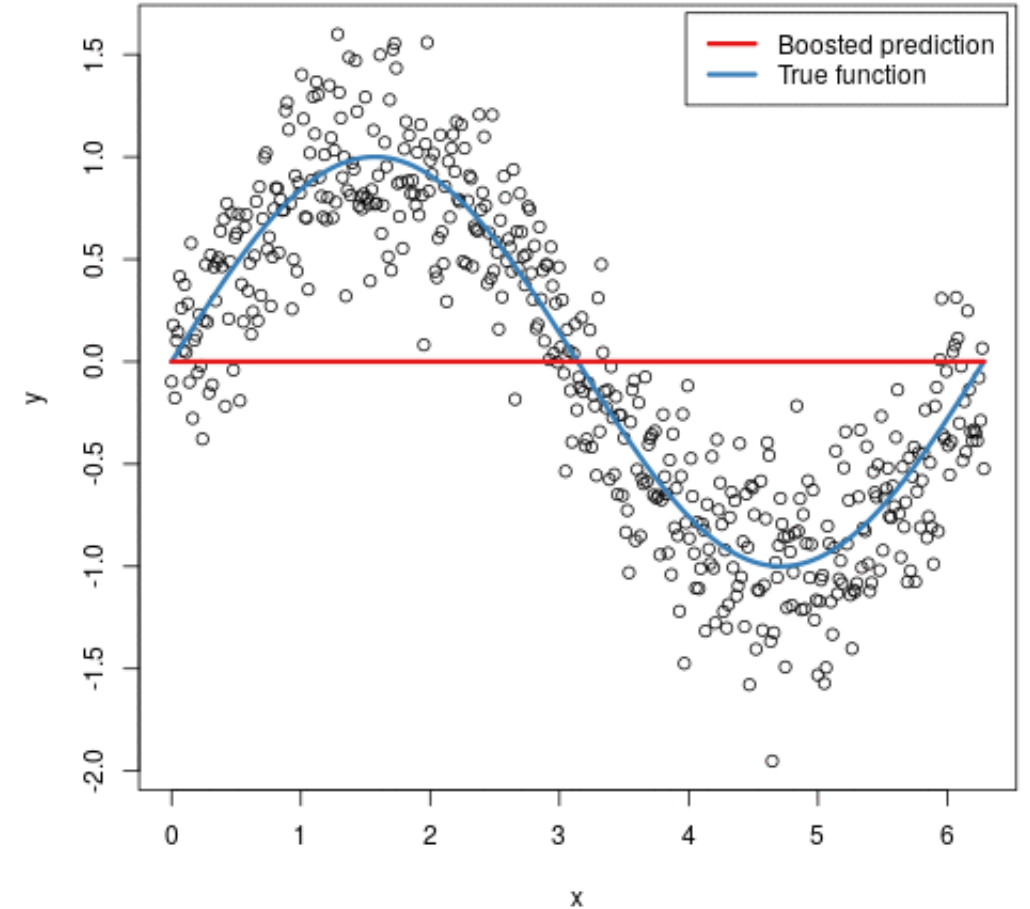
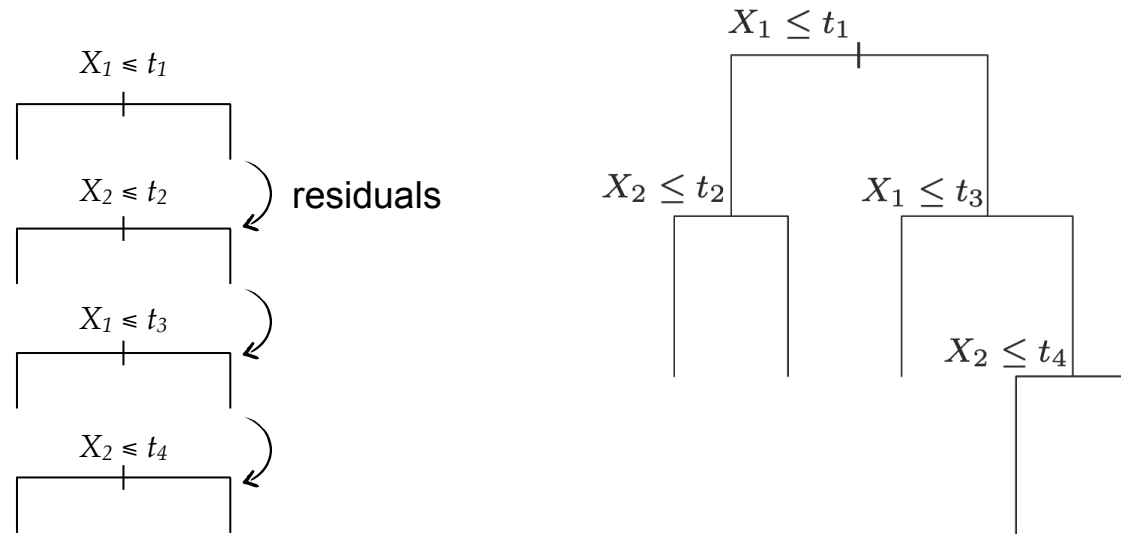
Context: relating POC attenuation to zooplankton diversity



Boosted regression trees

Combining many small models

Binary splits to relate response to predictors



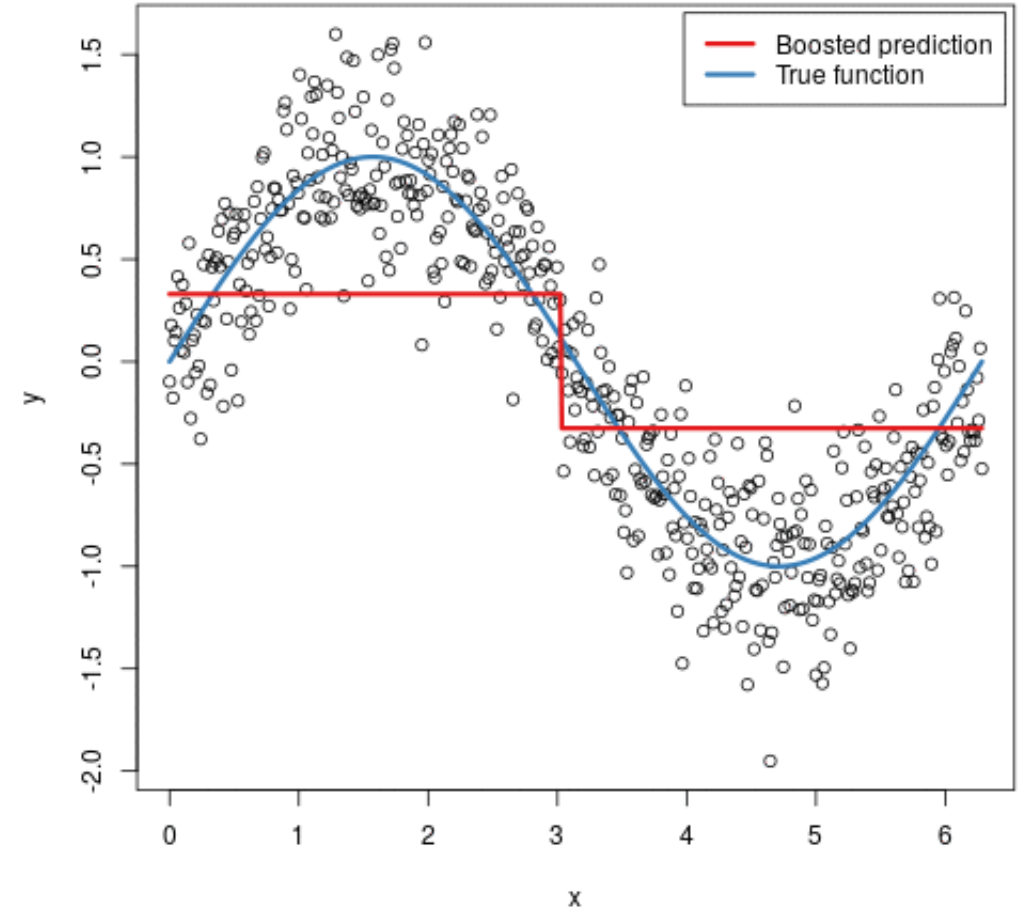
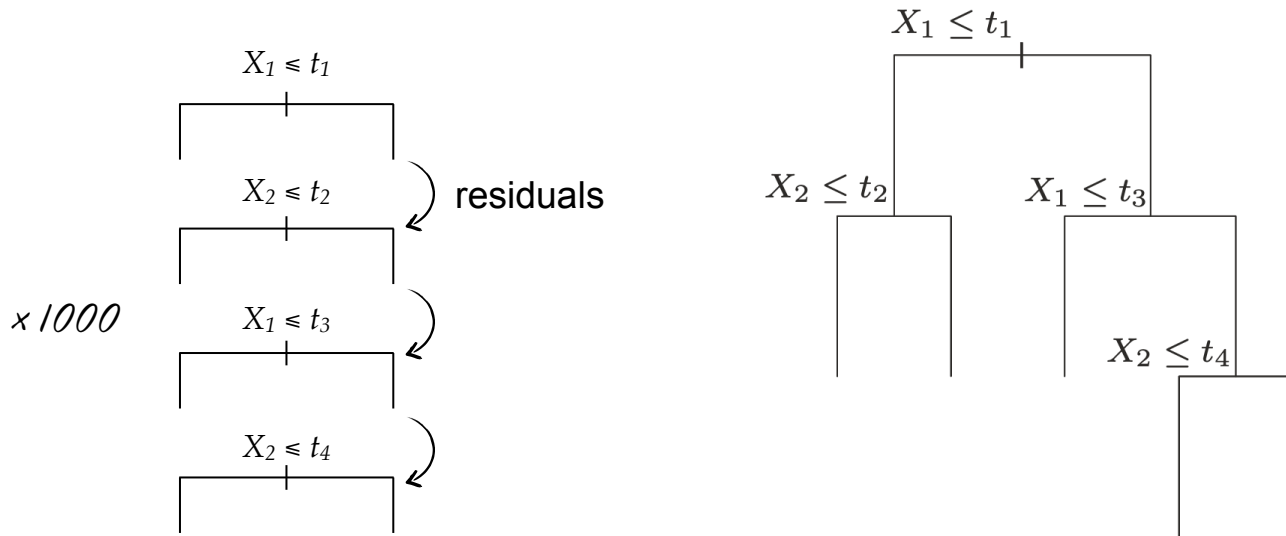


Context: relating POC attenuation to zooplankton diversity

Boosted regression trees

Combining many small models

Binary splits to relate response to predictors



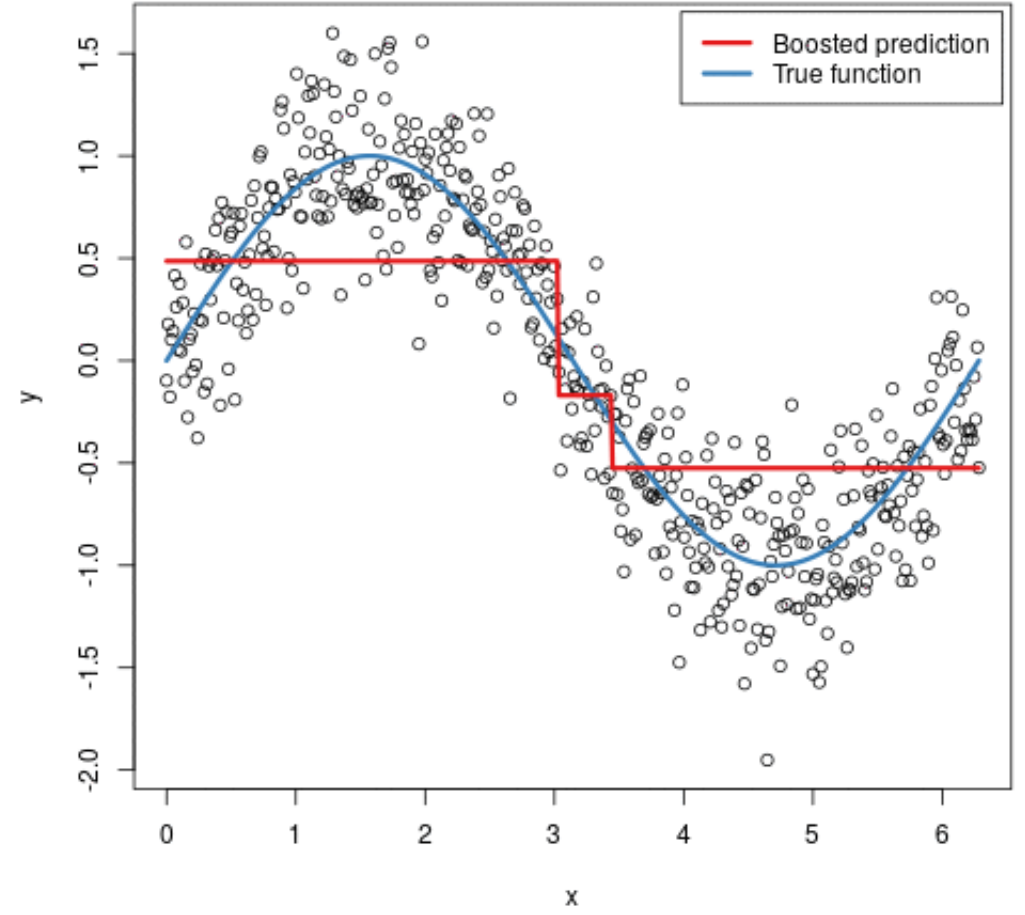
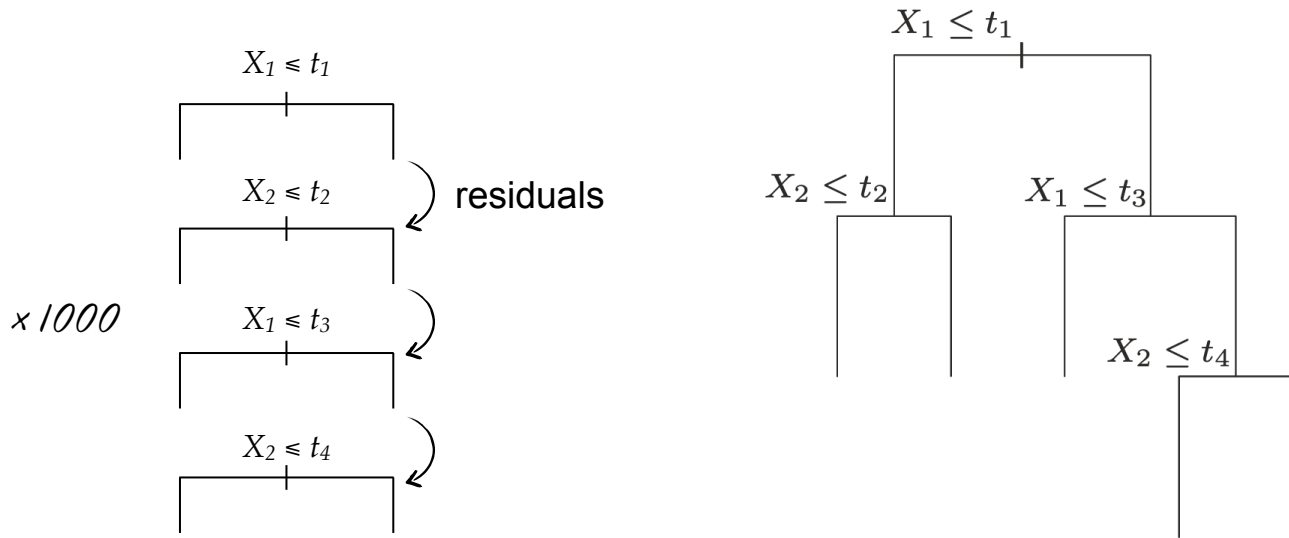


Context: relating POC attenuation to zooplankton diversity

Boosted regression trees

Combining many small models

Binary splits to relate response to predictors



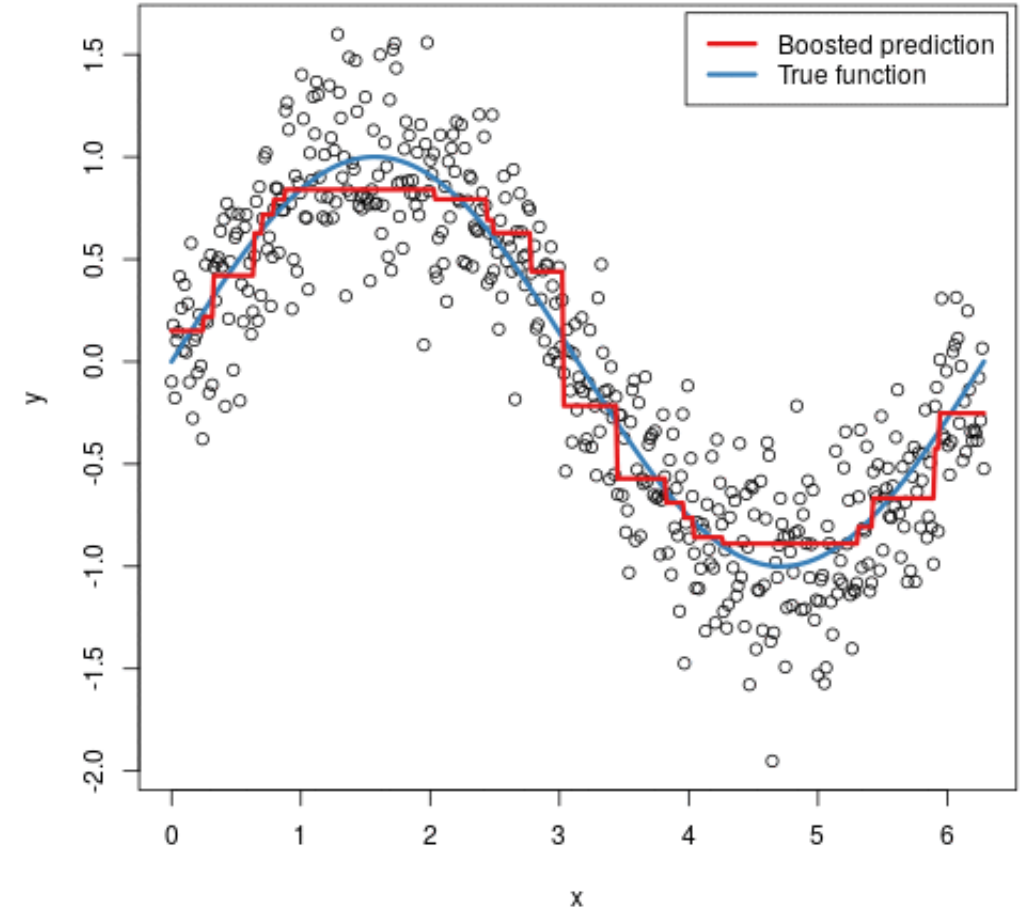
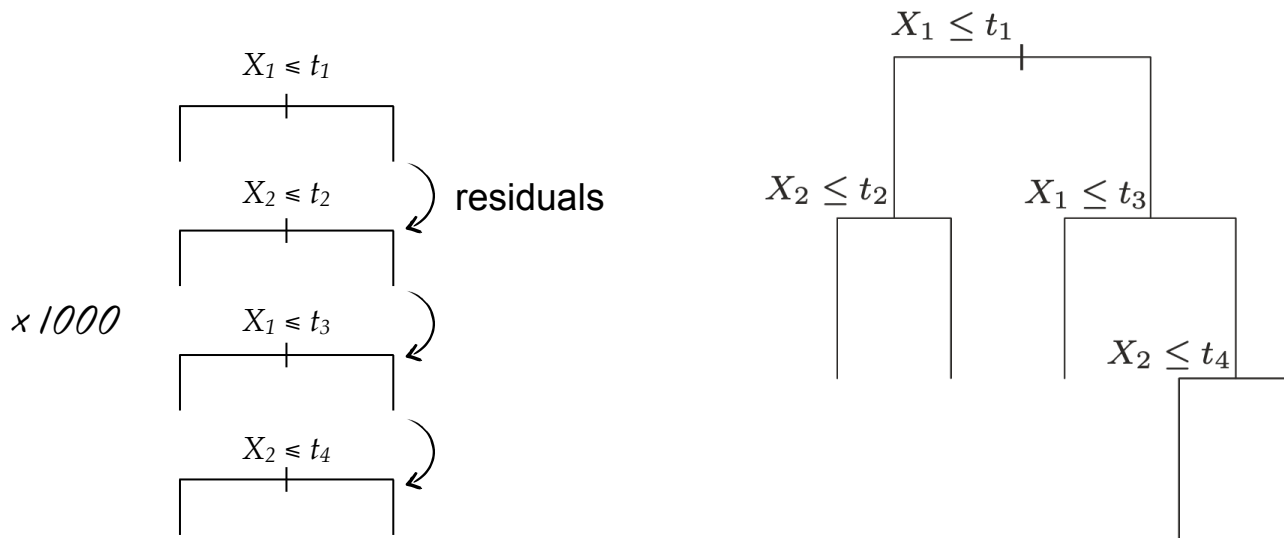
Context: relating POC attenuation to zooplankton diversity



Boosted regression trees

Combining many small models

Binary splits to relate response to predictors



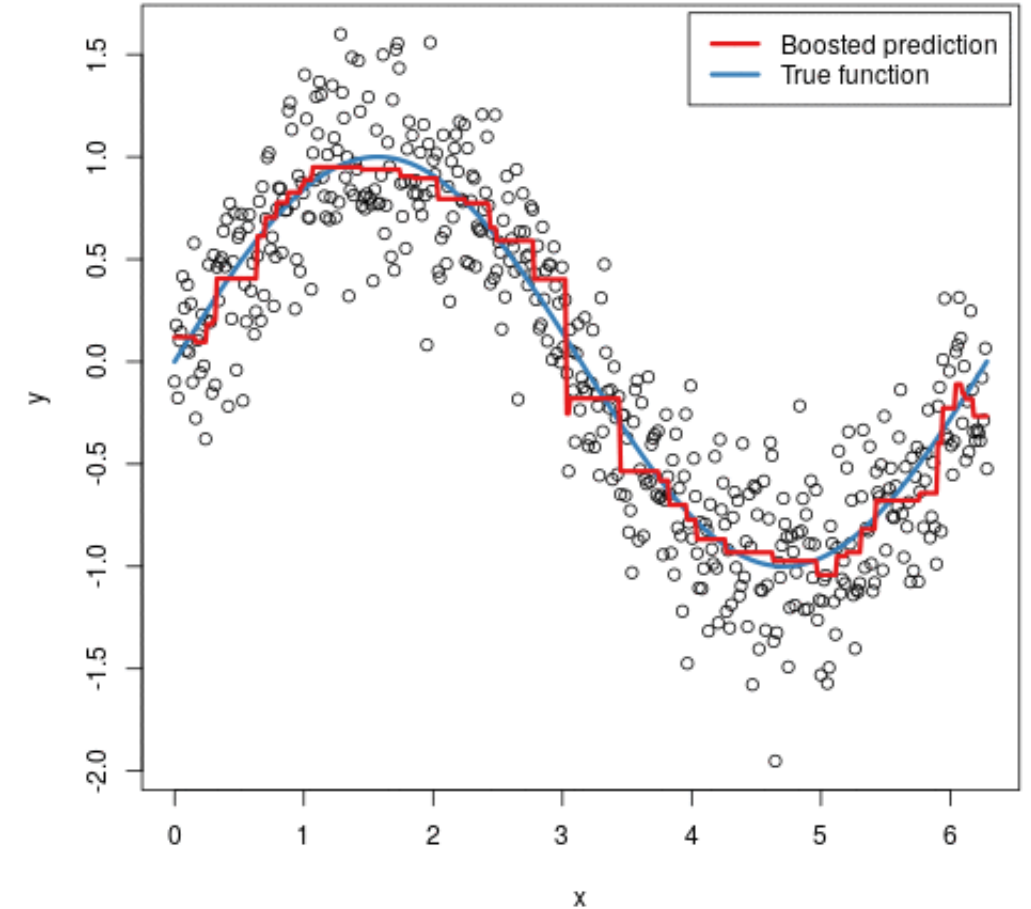
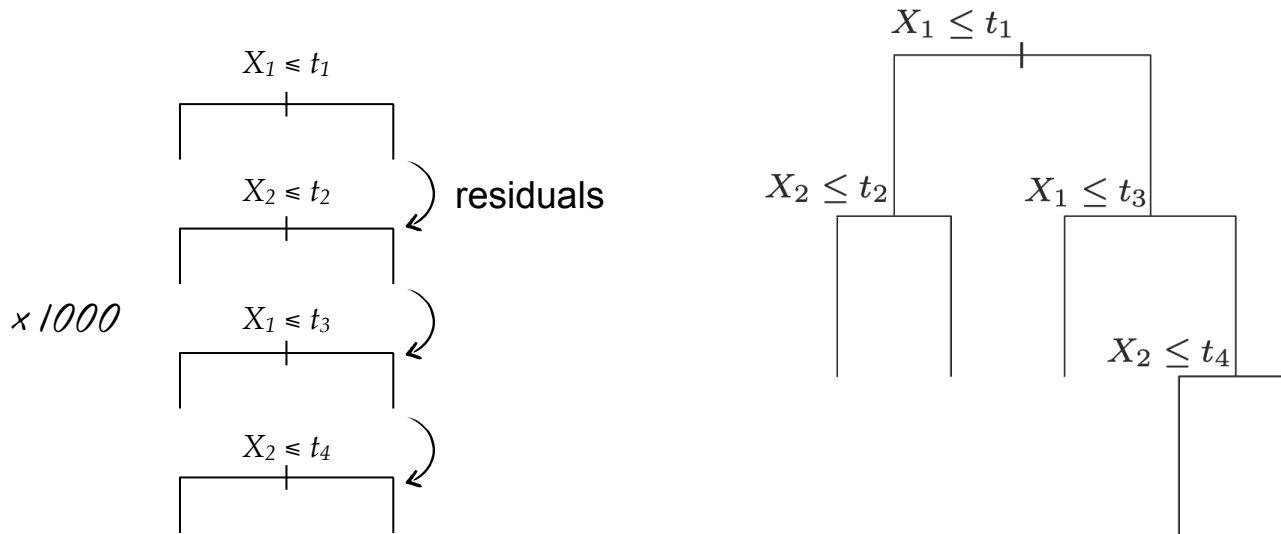
Context: relating POC attenuation to zooplankton diversity



Boosted regression trees

Combining many small models

Binary splits to relate response to predictors

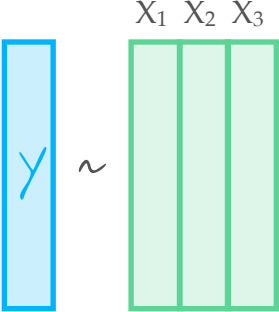


A 3D-rendered maze with a person standing in the center, symbolizing a complex problem or a path to a solution.

ML PRO TIP #1
“CHOOSE AN APPROPRIATE MODEL”

Tree ensembles are fantastic*

*if input is numeric

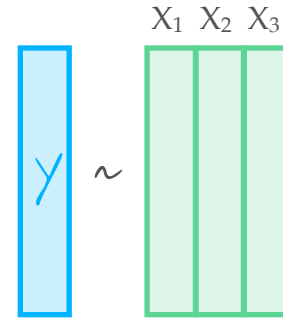


Tree ensembles are fantastic*

Many advantages:

- input flexibility (type, distribution, missing values, relevance)
- complex non-linear relationships + interactions
- good predictive power
- interpretable
- many implementations (R, Python)

*if input is numeric



Elith et al., 2008; Hastie et al., 2009



Tree ensembles are fantastic*

Many advantages:

- input flexibility (type, distribution, missing values, relevance)
- complex non-linear relationships + interactions
- good predictive power
- interpretable
- many implementations (R, Python)

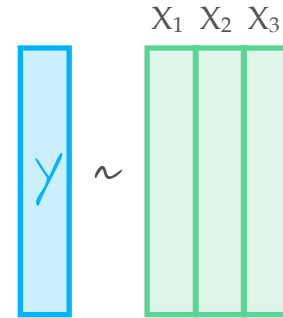
Elith et al., 2008; Hastie et al., 2009

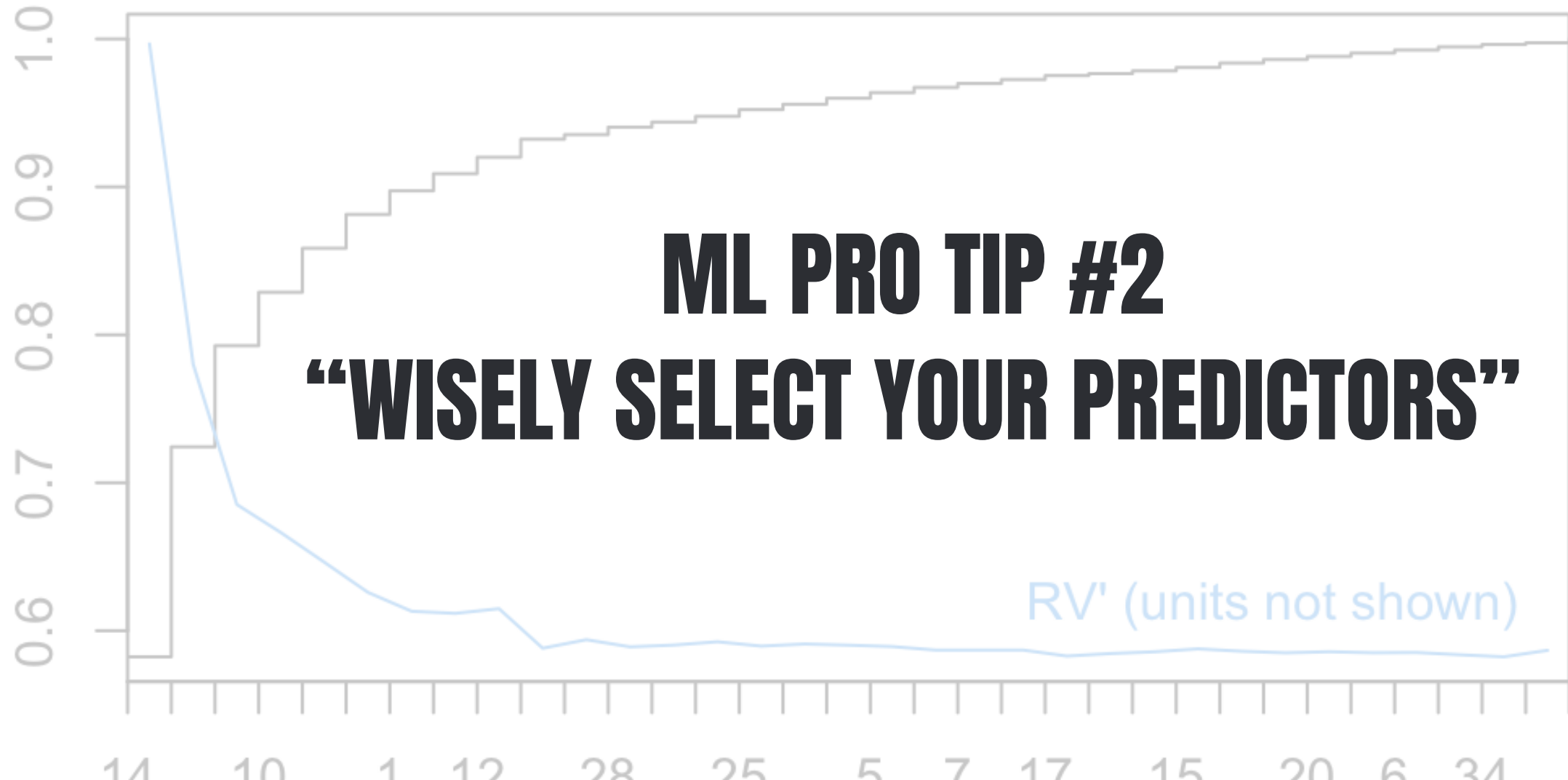
Classification: tree ensembles (RF) > neural network (ANN/MLP)

Regression?

Fernández-Delgado et al., 2014

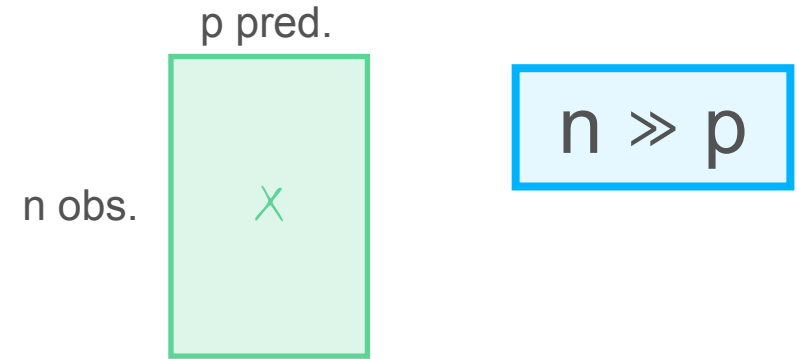
*if input is numeric





Less is more

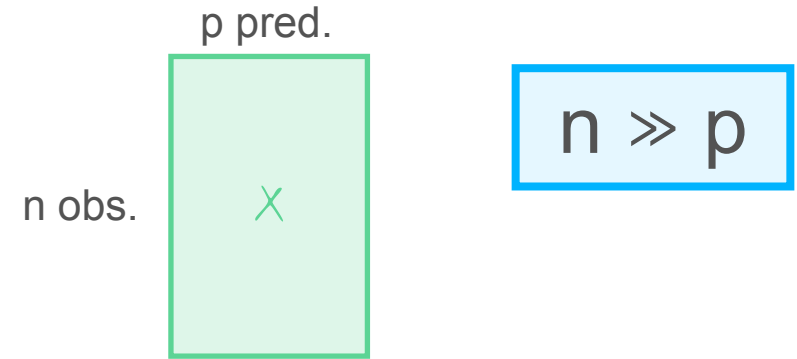
Number of predictors/features VS number of observations.



Less is more

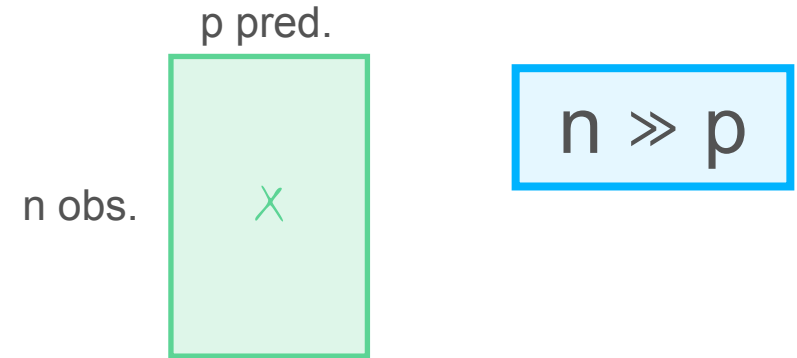
Number of predictors/features VS number of observations.

Trees can ignore non-relevant predictors.



Less is more

Number of predictors/features VS number of observations.



Trees can ignore non-relevant predictors.

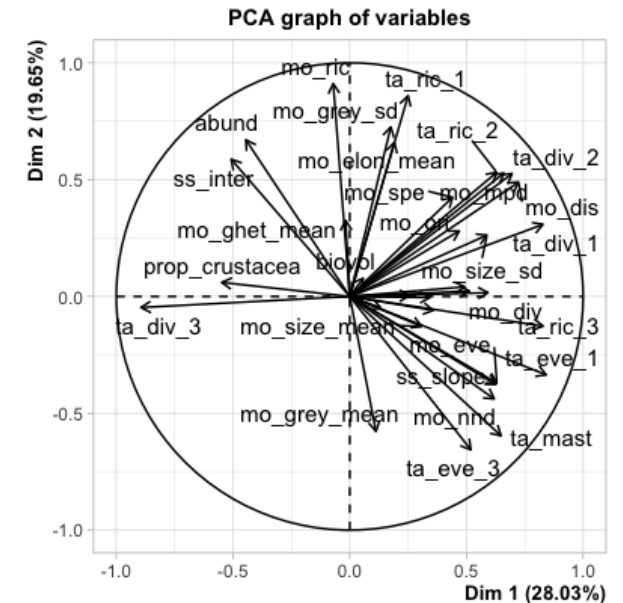
Parsimony

Feature selection

- PCA
- Escoufier's equivalent vectors
- VIF

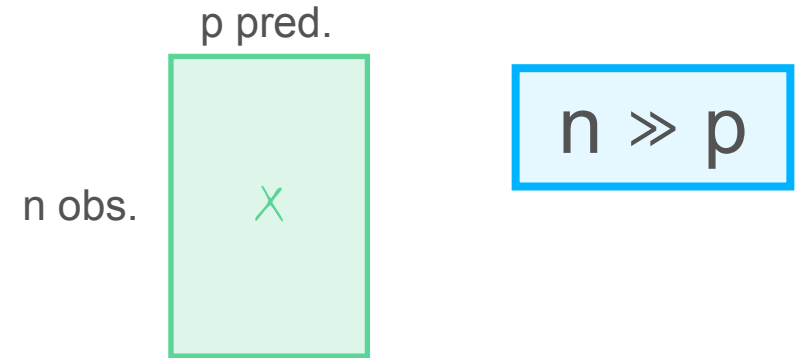
Feature engineering

- PCA: use PCs as predictors



Less is more

Number of predictors/features VS number of observations.



Trees can ignore non-relevant predictors.

Parsimony

Feature selection

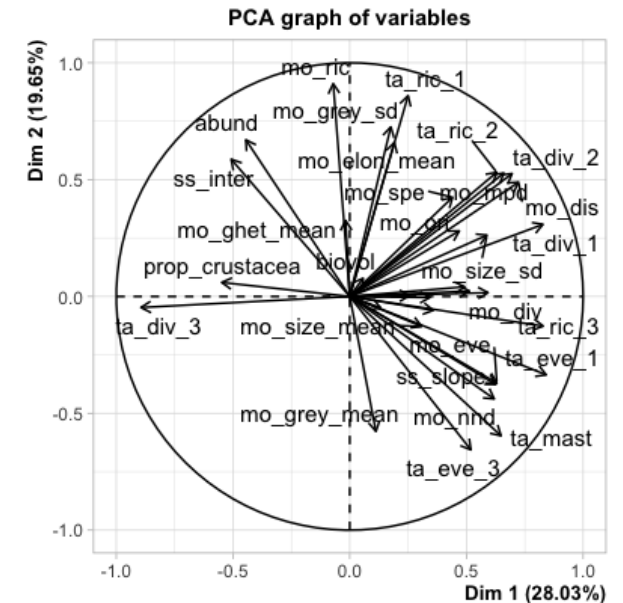
- PCA
- Escoufier's equivalent vectors
- VIF

Feature engineering

- PCA: use PCs as predictors

Correlated features?

Depends on your model. Tree ensembles are fairly robust.



ML PRO TIP #3
“MANAGE YOUR DATA BUDGET”

Need to spend the data

~80%

Train VS Test

~20%

Fit the model

Evaluate model performance, **at the very end, single use.**

How the model will perform with new data?

Regression

- R^2
- RMSE

Classification

- Accuracy
- Precision
- Recall

Need to spend the data



~80%

Train VS Test

~20%

Fit the model

Evaluate model performance, **at the very end, single use.**

How the model will perform with new data?

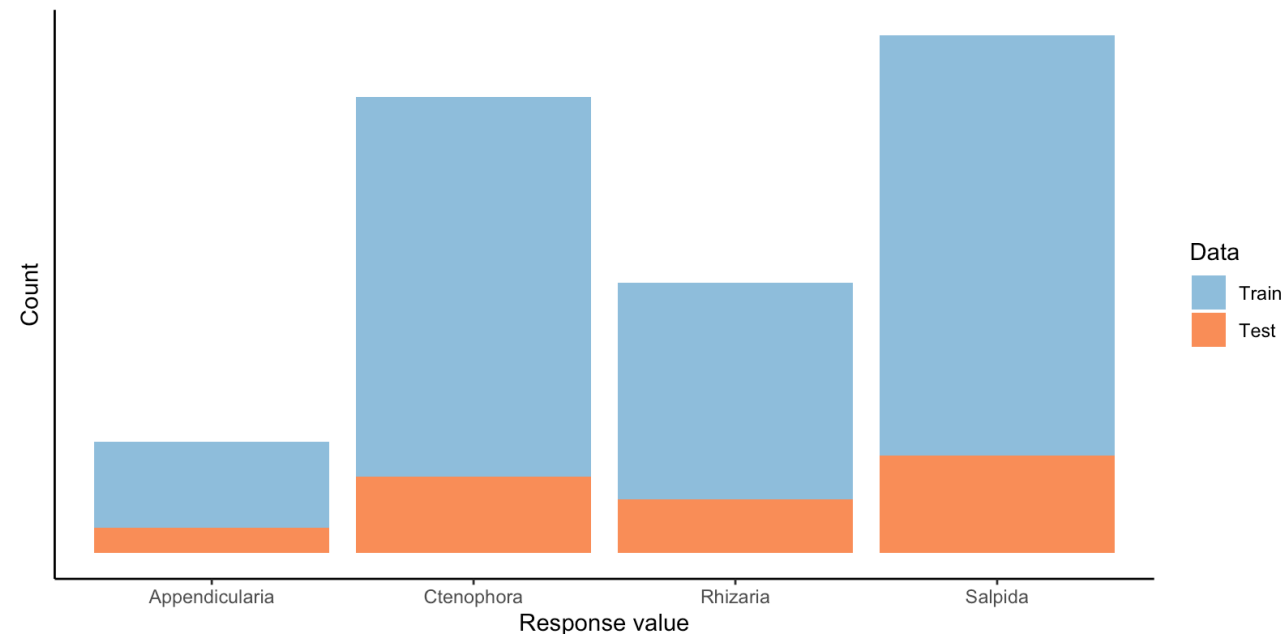
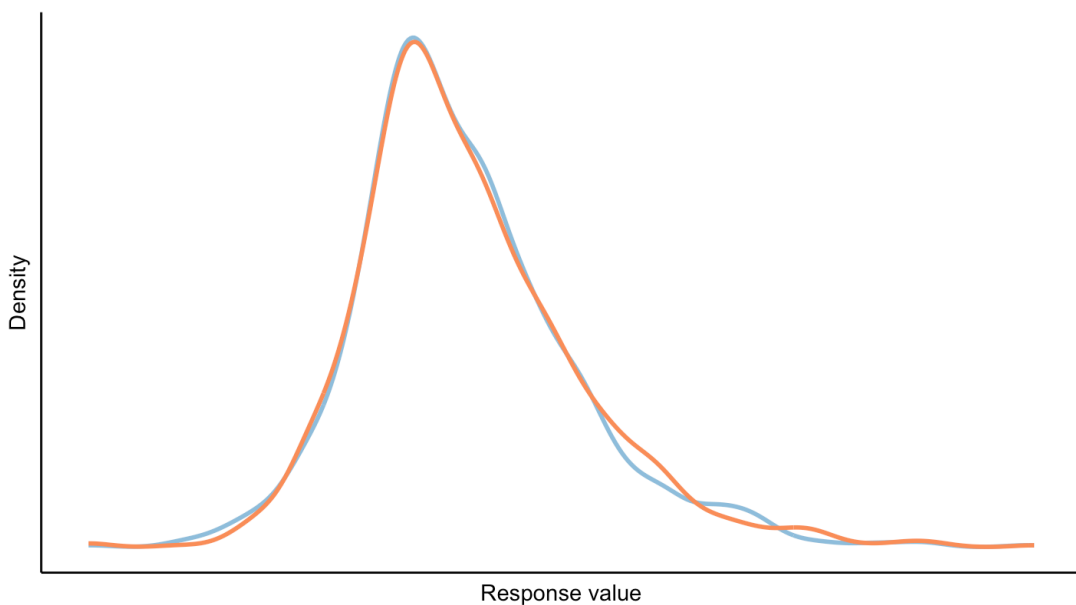
Test set representative of the training set.

Regression

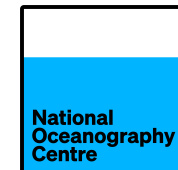
- R^2
- RMSE

Classification

- Accuracy
- Precision
- Recall



Need to spend the data



~80%

Train VS Test

~20%

Fit the model

Evaluate model performance, **at the very end, single use.**

How the model will perform with new data?

Test set representative of the training set.

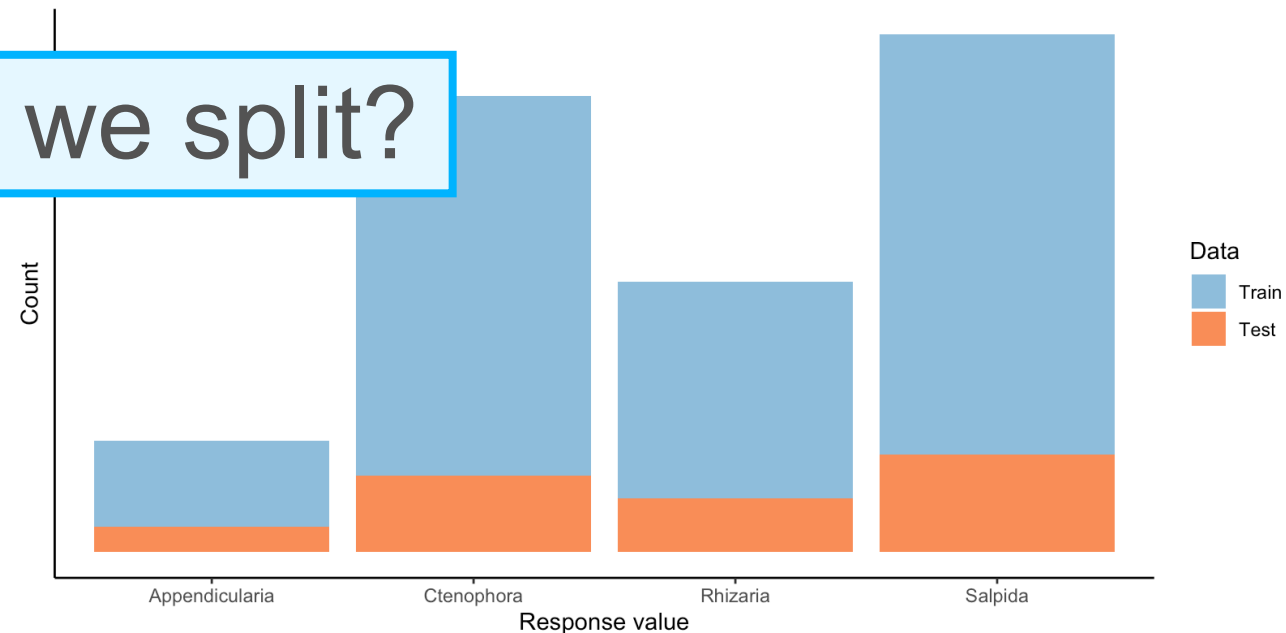
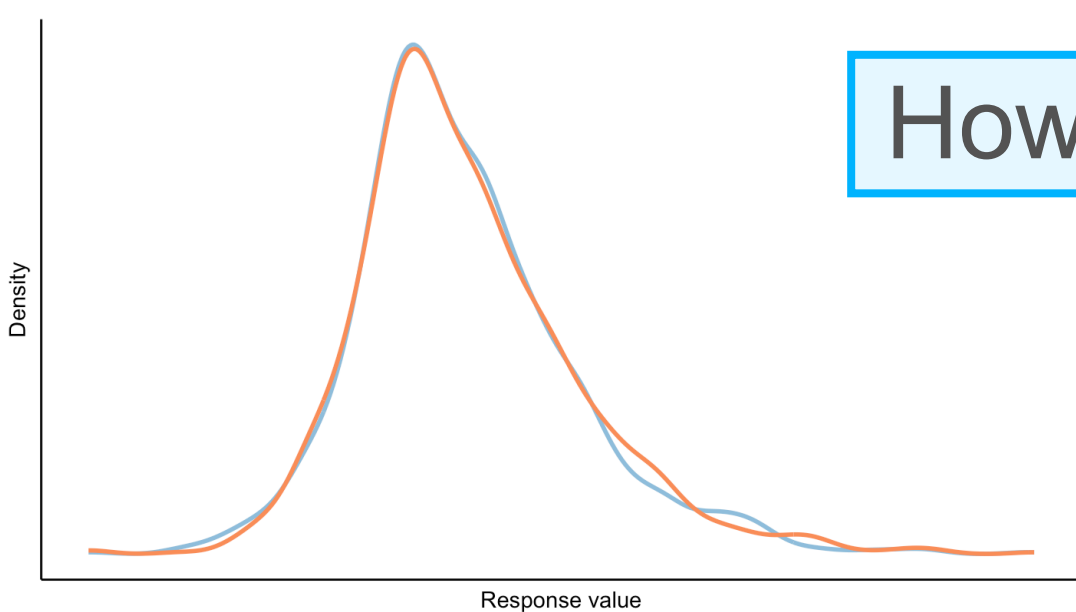
Regression

- R^2
- RMSE

Classification

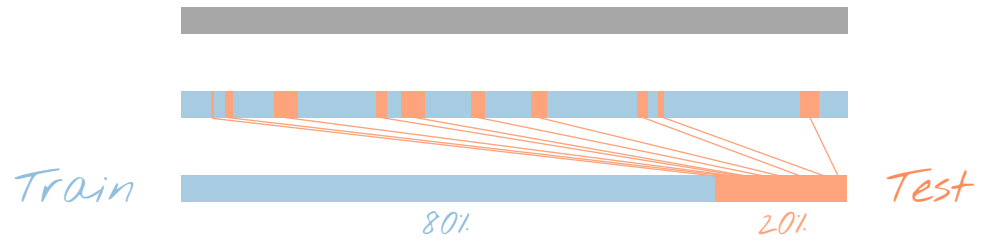
- Accuracy
- Precision
- Recall

How do we split?

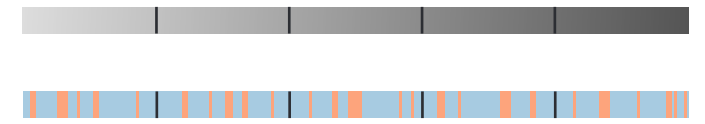


How to split your data

Random

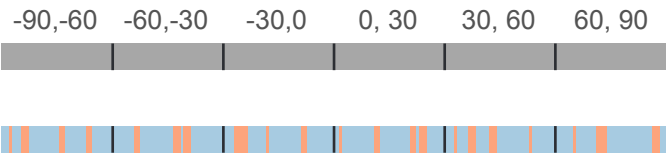


Quantiles of response variable

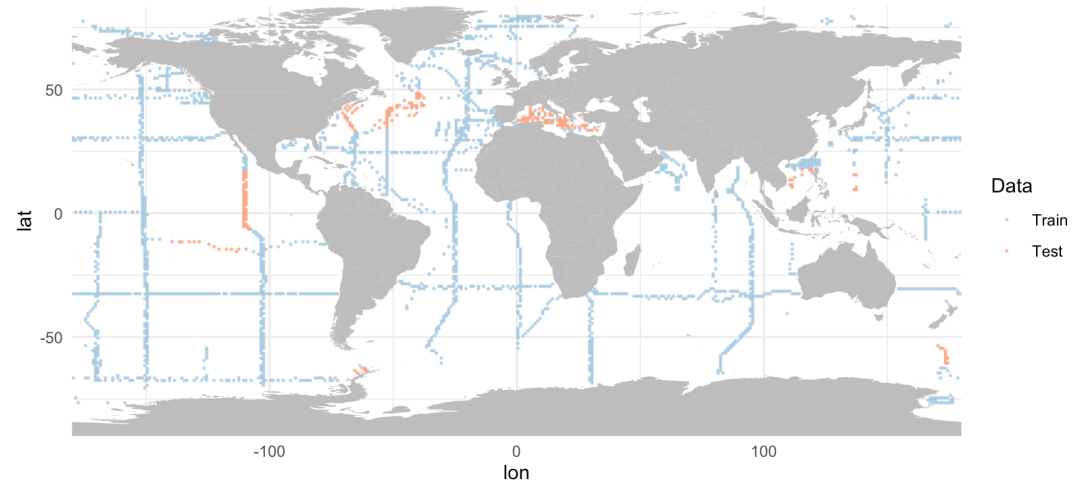


Groups

(e.g. latitude bands)

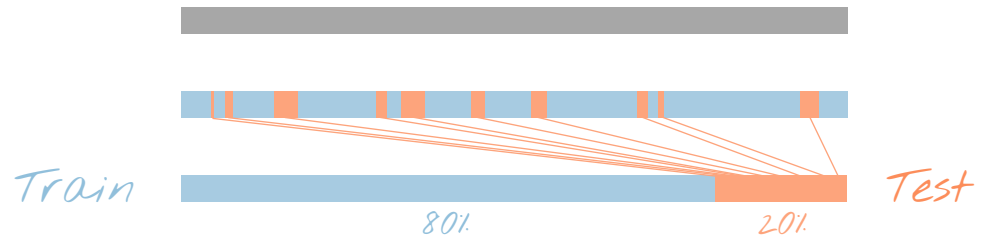


Spatial / temporal

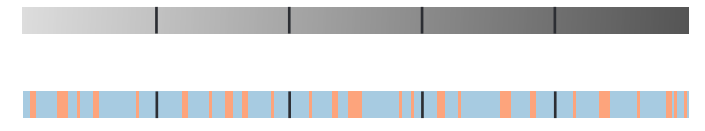


How to split your data

Random



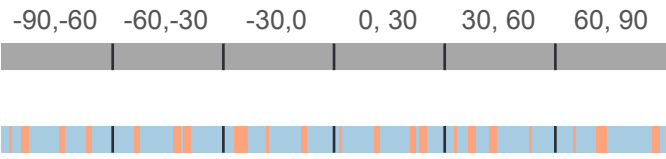
Quantiles of response variable



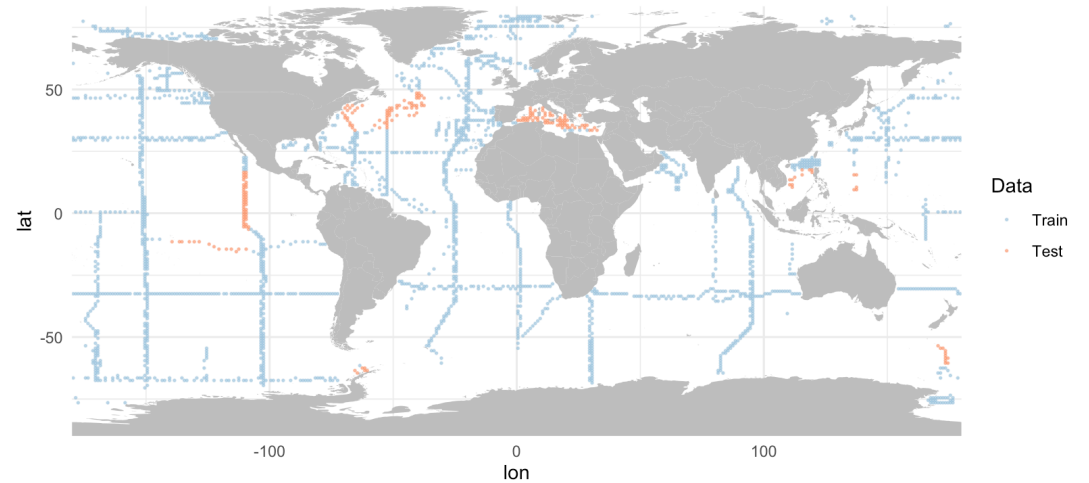
Depends on your data!

Groups

(e.g. latitude bands)



Spatial / temporal



The background of the slide is a dense, overlapping pattern of various-sized, metallic-looking gears. The gears are rendered in shades of grey and white, creating a complex, mechanical texture. Some gears are in sharp focus, while others are blurred in the background, giving a sense of depth and movement.

ML PRO TIP #4
“SOME ML MODELS NEED TO BE TUNED”

Optimising your model hyperparameters

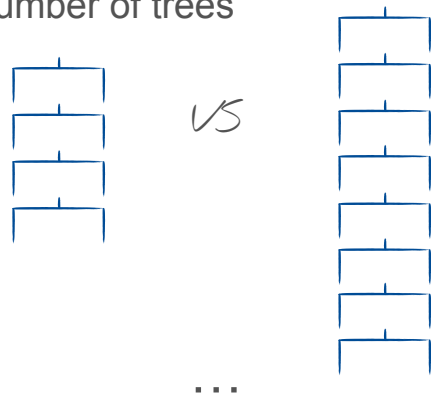
Tree ensembles

Boosted trees

Depth of trees



Number of trees



Optimising your model hyperparameters

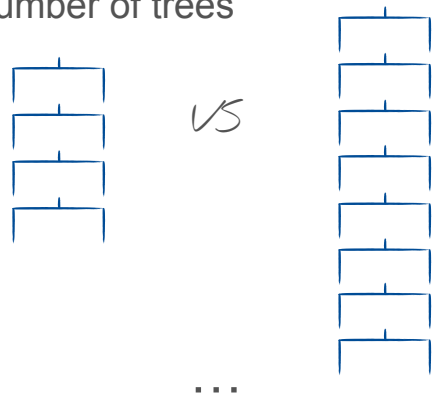
Tree ensembles

Boosted trees

Depth of trees



Number of trees

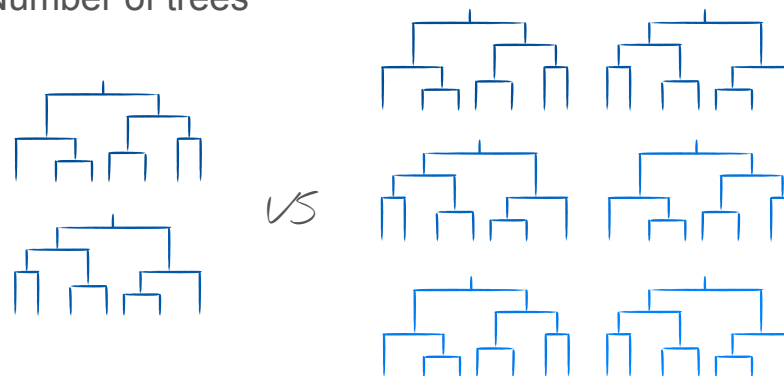


Random Forest

Depth of trees



Number of trees



Optimising your model hyperparameters

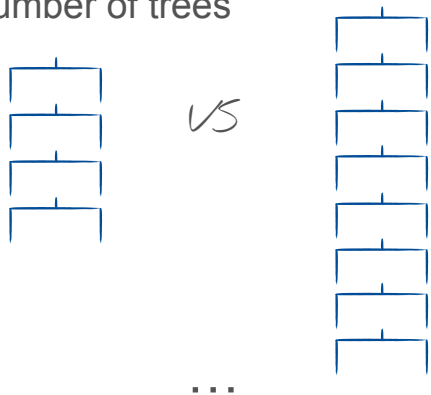
Tree ensembles

Boosted trees

Depth of trees



Number of trees

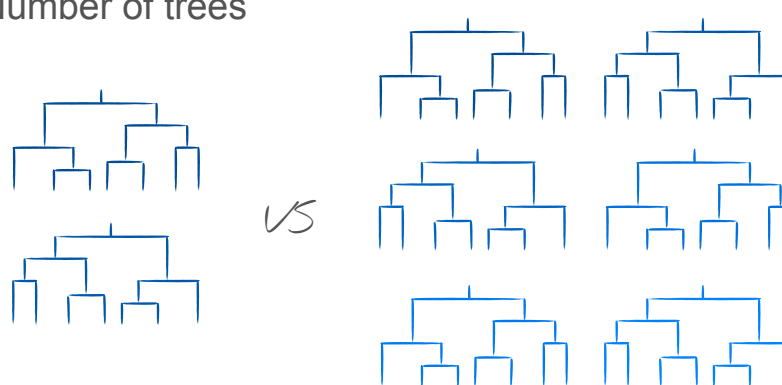


Random Forest

Depth of trees



Number of trees



Neural networks (MLP/ANN)

Size of layers



Number of layers



Optimising your model hyperparameters

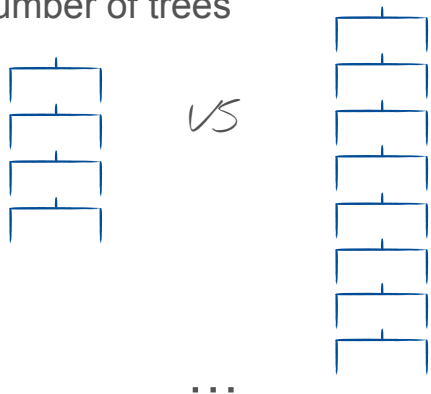
Tree ensembles

Boosted trees

Depth of trees



Number of trees

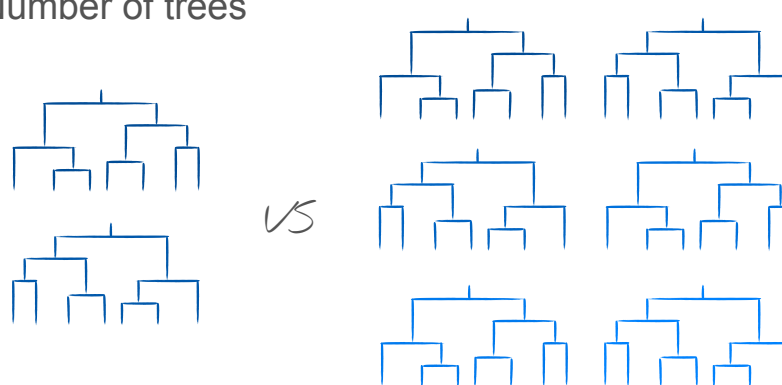


Random Forest

Depth of trees



Number of trees



Neural networks (MLP/ANN)

Size of layers



Number of layers



→ to be tuned: model tuning / gridsearch

That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth

That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth

shallow



Fit

Train



That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth

shallow



Fit

Train



Eval

Test



$R^2 = 74\%$

That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth


shallow



Fit

Train

Eval

Test 
 $R^2 = 74\%$

deep



Fit



That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth

shallow

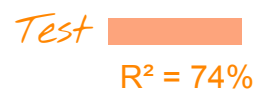


Fit

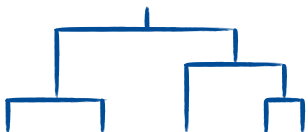


Train

Eval



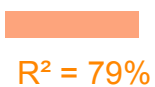
deep



Fit



Eval



That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth


shallow



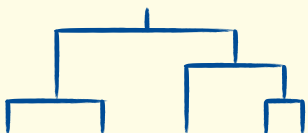
Fit

Train

Eval

Test 
 $R^2 = 74\%$

deep



Fit

Eval

$R^2 = 79\%$

That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth

shallow



Fit

Train

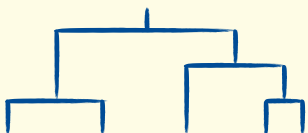


Eval



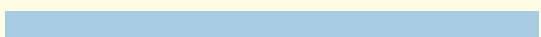
$R^2 = 74\%$

deep



Fit

Eval



$R^2 = 79\%$

Test 

optimising hyperparameters + estimating generalisation error

That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth

shallow

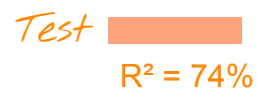


Fit

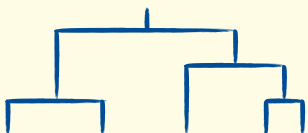


Train

Eval



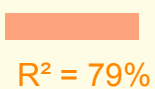
deep



Fit



Eval



optimising hyperparameters ~~Test~~ estimating generalisation error

That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth

shallow

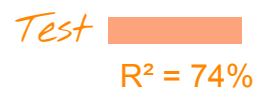


Fit

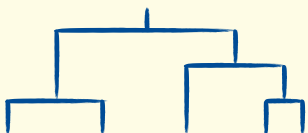


Train

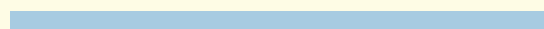
Eval



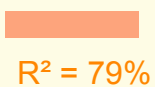
deep



Fit



Eval



optimising hyperparameters ~~Test~~ estimating generalisation error

Test  estimating generalisation error

That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth

shallow



Fit

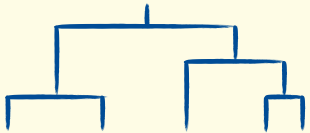
Train

Eval

Test

$R^2 = 74\%$

deep



Fit

Eval

$R^2 = 79\%$

~~optimising hyperparameters~~ ~~estimating generalisation error~~

Test

estimating generalisation error

Validation

optimising hyperparameters

That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth

shallow

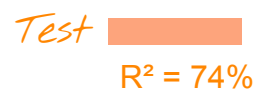


Fit

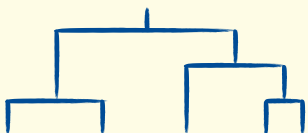


Train

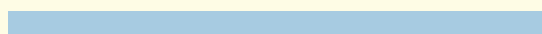
Eval



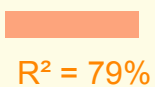
deep



Fit



Eval



~~Test~~
optimising hyperparameters ~~estimating generalisation error~~



Test 

estimating generalisation error

Validation 

optimising hyperparameters

Need 3 splits



That's what the validation set is for

Example for boosted regression trees

Tuning the tree depth

shallow

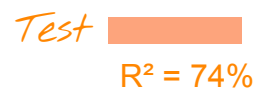


Fit

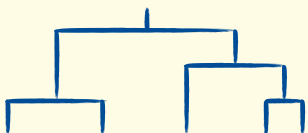


Train

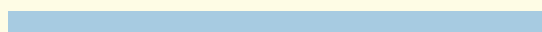
Eval



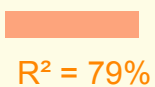
deep



Fit



Eval



~~Test~~
optimising hyperparameters ~~estimating generalisation error~~

Test 

estimating generalisation error

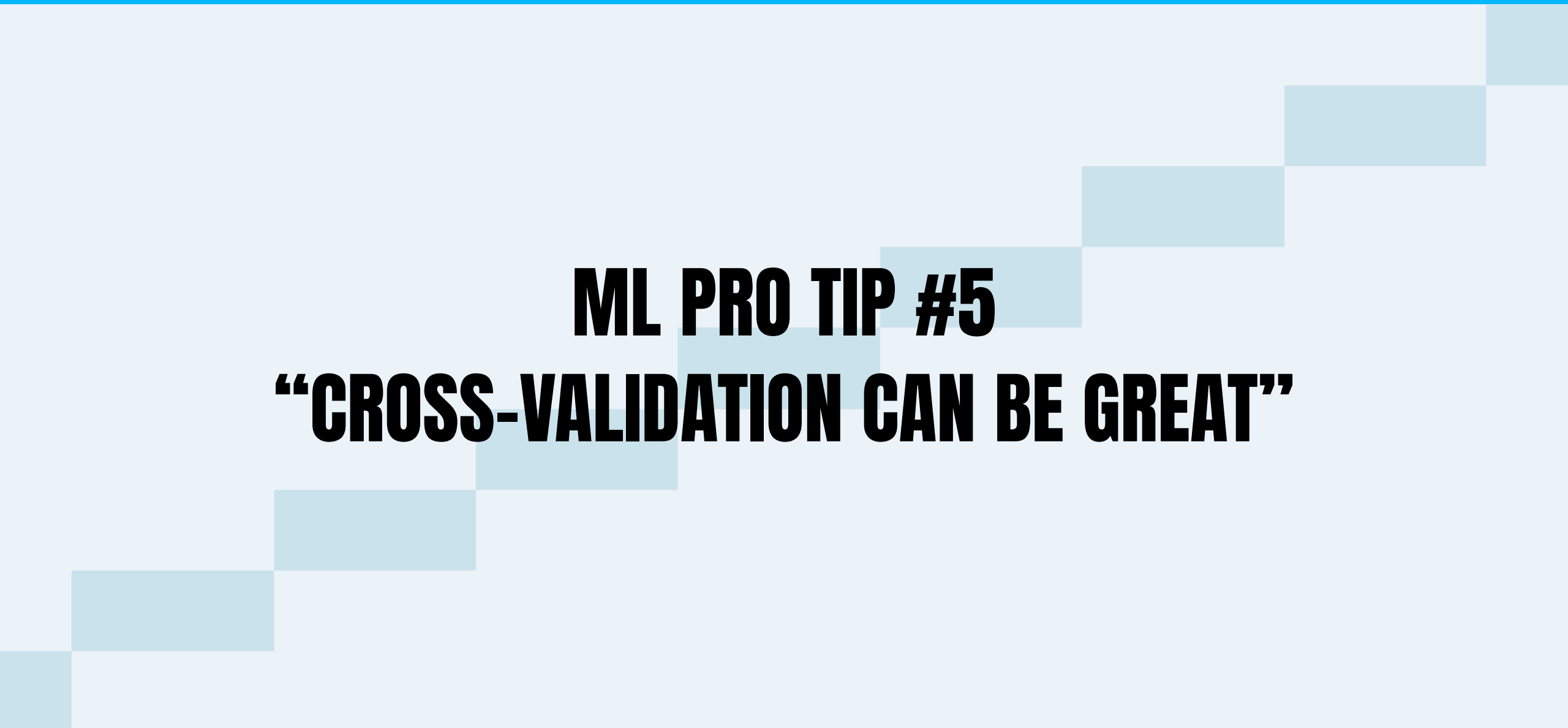
Validation 

optimising hyperparameters

Need 3 splits



Validation set used to tune the model

A decorative graphic consisting of a series of light blue rectangular blocks arranged in a staircase pattern, ascending from the bottom left towards the top right.

ML PRO TIP #5

“CROSS-VALIDATION CAN BE GREAT”

What about cross validation?

Overcome the effect of randomness in your splits

What about cross validation?

Overcome the effect of randomness in your splits

Dataset



5 folds CV

What about cross validation?

Overcome the effect of randomness in your splits

5 folds CV

Dataset



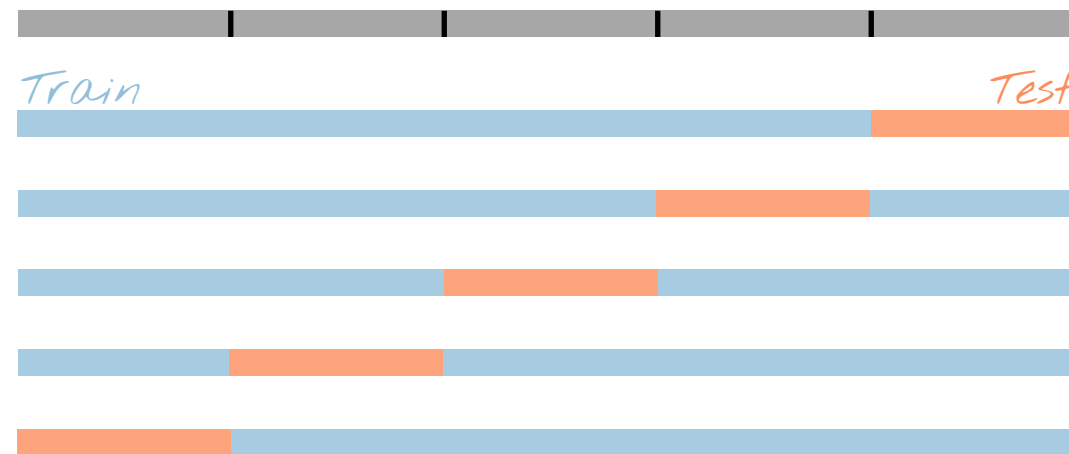
What about cross validation?

Overcome the effect of randomness in your splits

5 folds CV

- 5 iterations
- 5 models fitted
- 5 performance estimates

Dataset



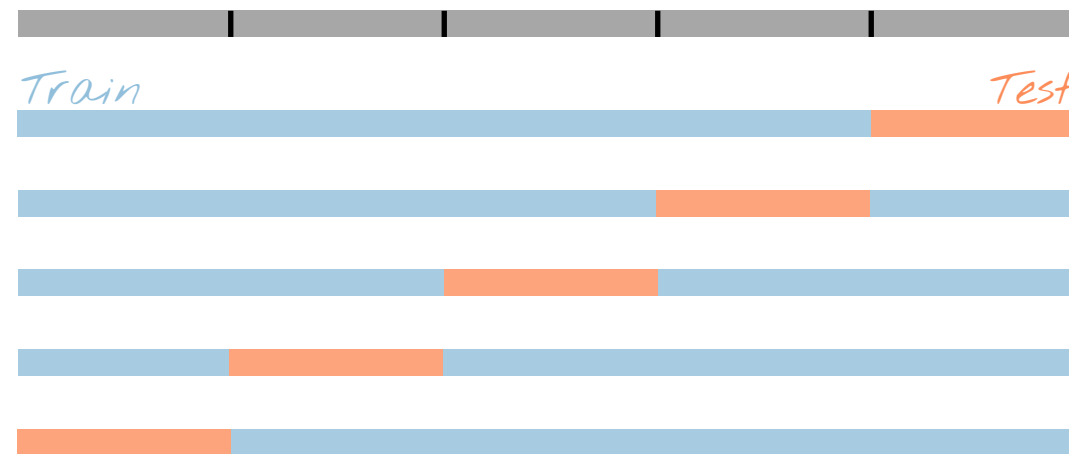
What about cross validation?

Overcome the effect of randomness in your splits

5 folds CV

- 5 iterations
- 5 models fitted
- 5 performance estimates

Dataset



Validation and model tuning?

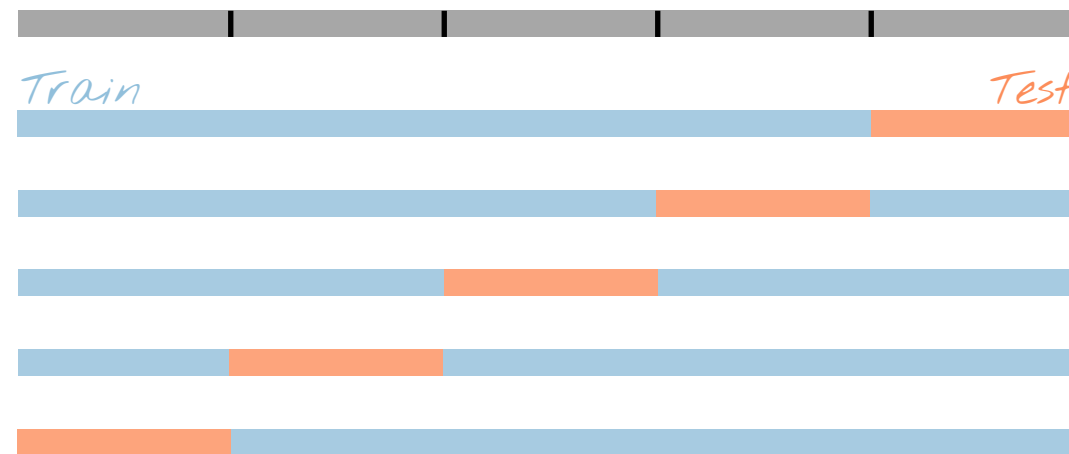
What about cross validation?

Overcome the effect of randomness in your splits

5 folds CV

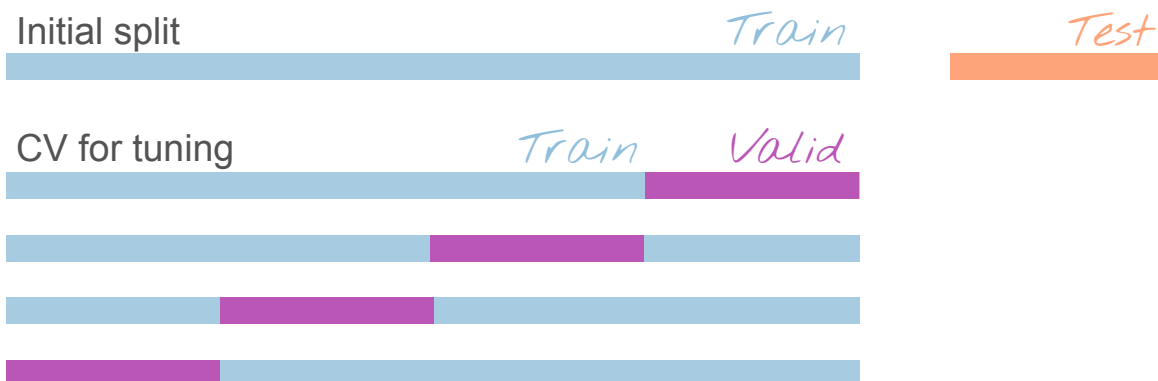
- 5 iterations
- 5 models fitted
- 5 performance estimates

Dataset

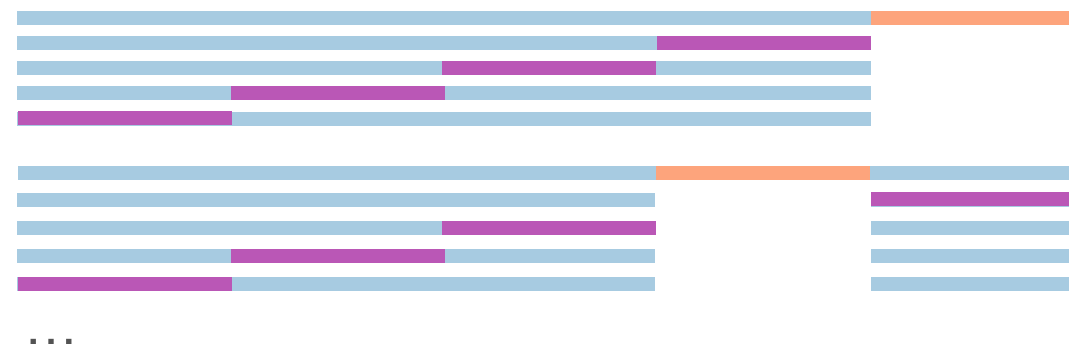


Validation and model tuning?

Initial split + CV



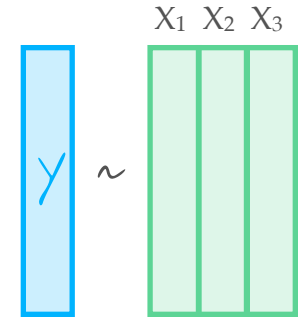
Nested CV



The background of the slide is a 3D rendering of a large number of white, semi-transparent cubes arranged in a grid-like pattern, receding into the distance. One cube in the foreground is highlighted with a bright yellow glow.

ML PRO TIP #6
“ML MODELS CAN BE INTERPRETED”

ML models are not black boxes

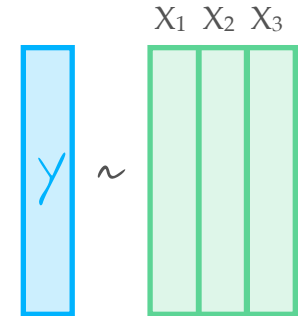
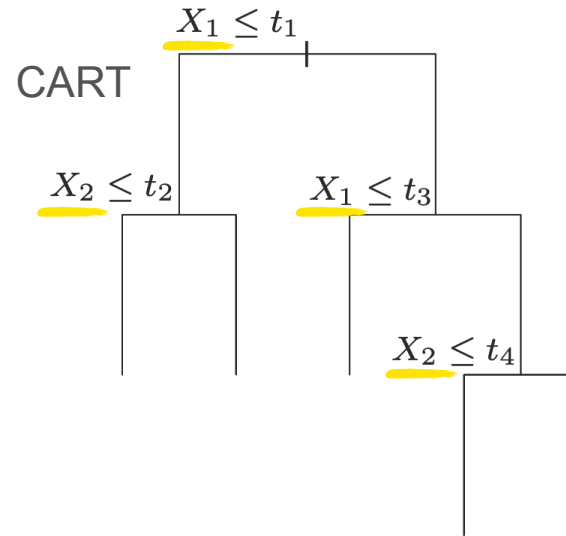


ML models are not black boxes

- Some models are easy to interpret

linear regression

$$Y = b + \underline{3.7} \times X_1 + 0.01 \times X_2 + \underline{1.6} \times X_3$$



ML models are not black boxes

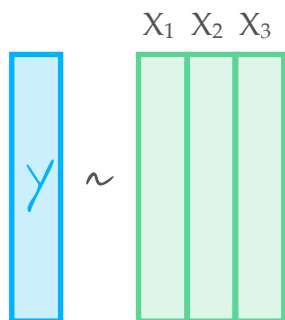
- Some models are easy to interpret

linear regression

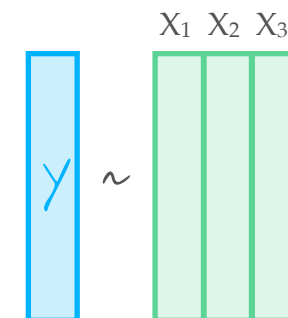
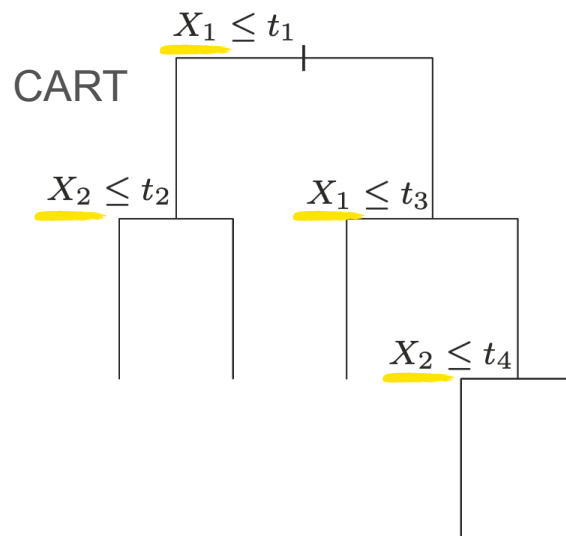
$$Y = b + \underline{3.7} \times X_1 + 0.01 \times X_2 + \underline{1.6} \times X_3$$

- For other ones, there are workarounds

Orig. model



Orig. performance

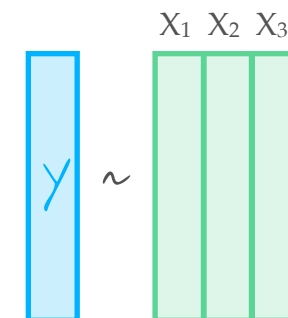
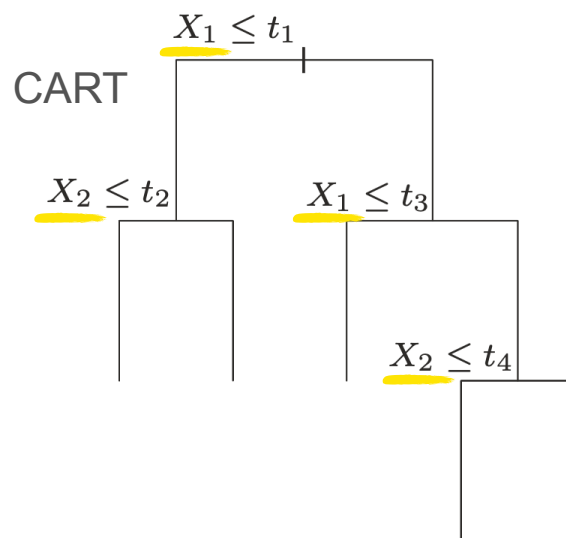


ML models are not black boxes

- Some models are easy to interpret

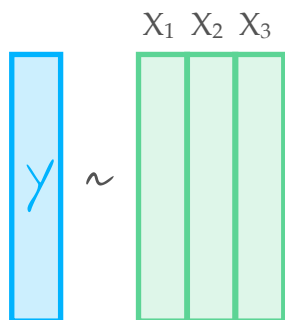
linear regression

$$Y = b + 3.7 \times X_1 + 0.01 \times X_2 + 1.6 \times X_3$$



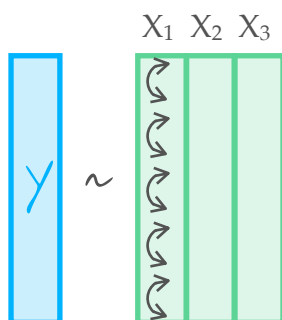
- For other ones, there are workarounds

Orig. model



Orig. performance

Shuffle X_1



~ same performance

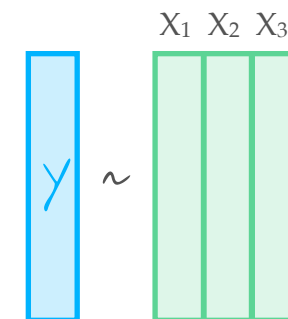
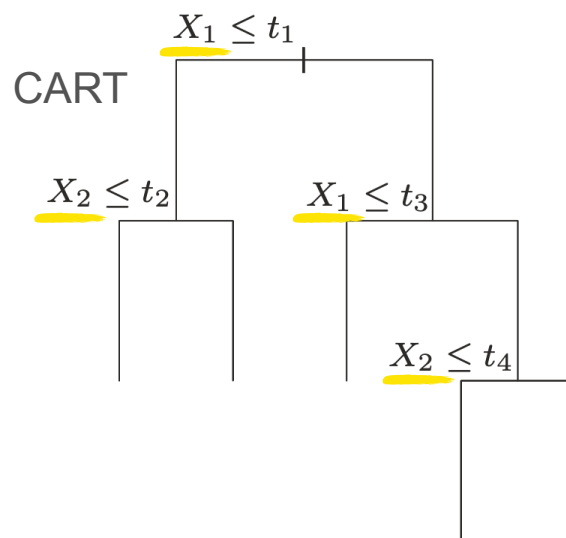
→ X_1 is not important

ML models are not black boxes

- Some models are easy to interpret

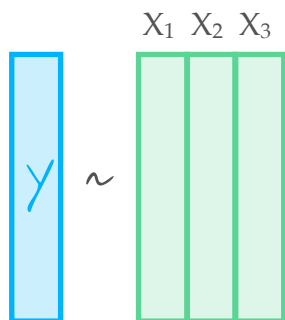
linear regression

$$Y = b + \underline{3.7} \times X_1 + 0.01 \times X_2 + \underline{1.6} \times X_3$$



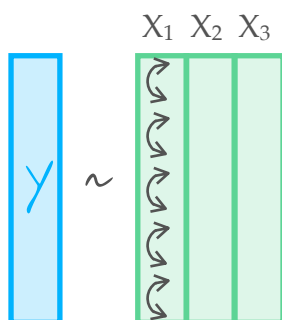
- For other ones, there are workarounds

Orig. model



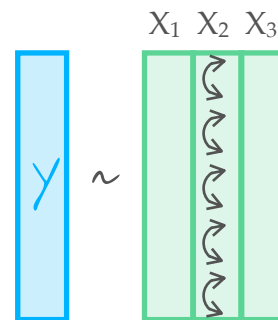
Orig. performance

Shuffle X_1



~ same performance
→ X_1 is not important

Shuffle X_2



Lower performance
→ X_2 is important

...

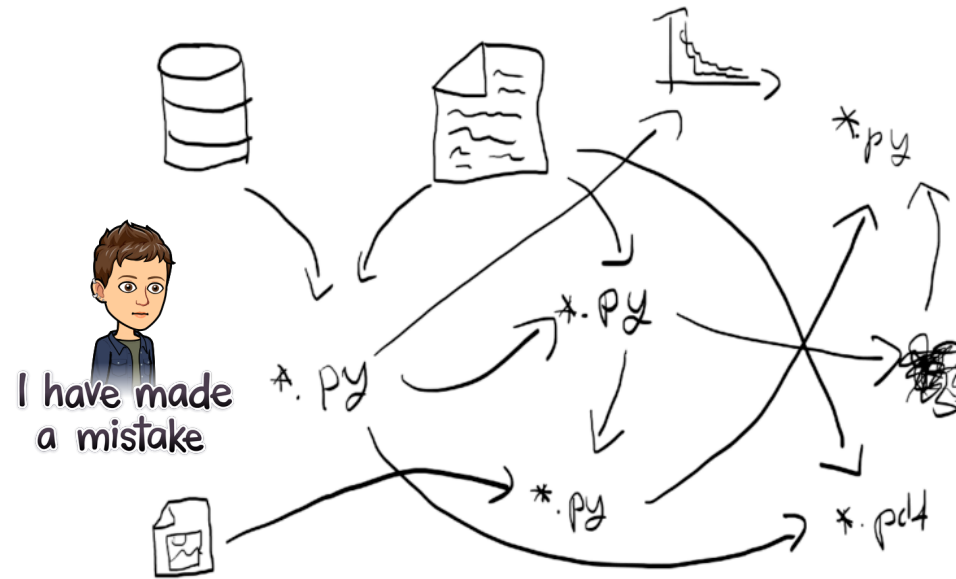
Importance of
each predictor

The background of the slide is a stylized landscape of rolling hills and mountains. The entire scene is rendered in shades of light blue and white, with a dense overlay of binary code (0s and 1s) that creates a digital, data-driven atmosphere.

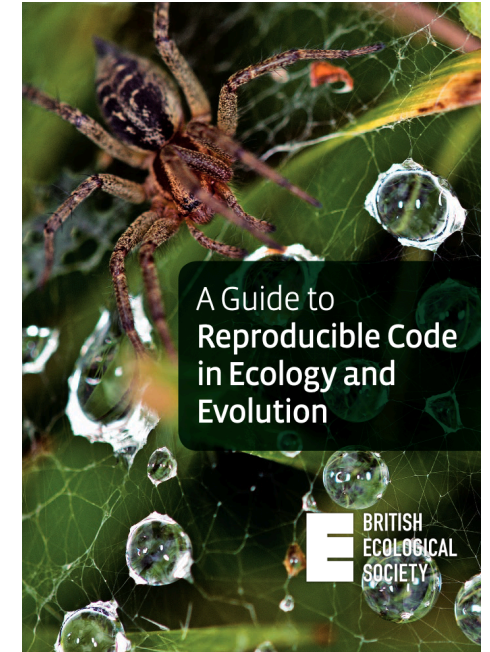
ML PRO TIP #7
**“SHARE WITH OTHERS,
LEARN FROM OTHERS”**

The importance of version control and sharing

- Version control
 - no more mess
 - go back in time



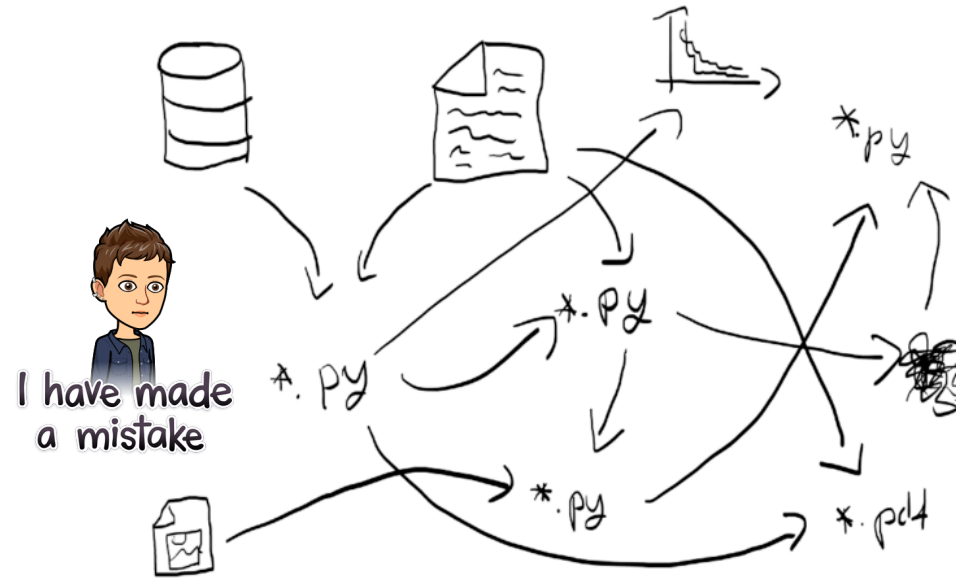
Broad Research Communication Lab



BES & Cooper, N. A Guide to Reproducible Code in Ecology and Evolution. (2017).

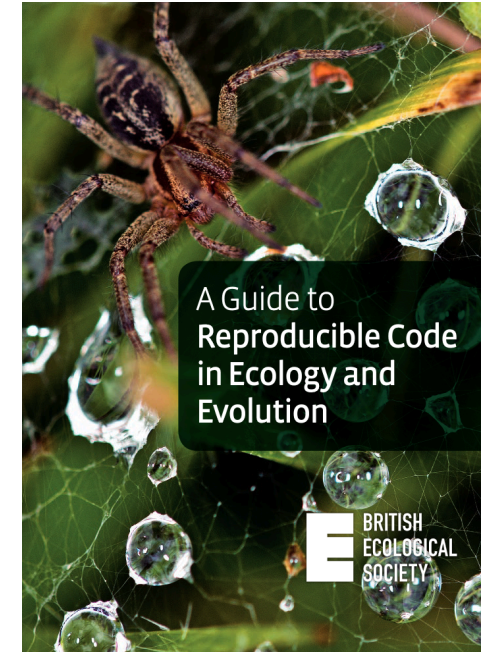
The importance of version control and sharing

- Version control
 - no more mess
 - go back in time



Broad Research Communication Lab

- Sharing
 - improve yourself
 - reproducibility

BES & Cooper, N. A Guide to Reproducible Code in Ecology and Evolution. (2017).

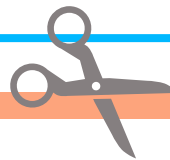
Conclusions

- Supervised ML relates response variable(s) to predictors
- ML is neither black magic, nor a black box
- **ML models can be interpreted**
- A few checks are essential to do it ~~right~~ *not wrong*
- Multiple choices are possible, depends on your data

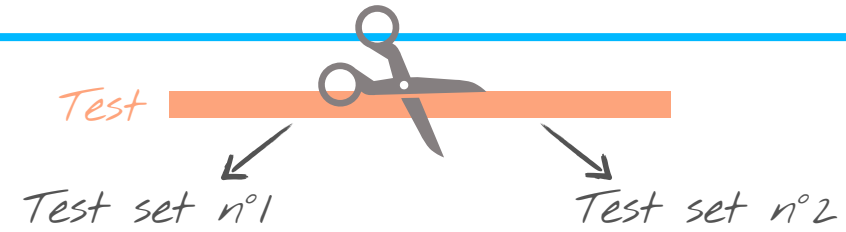
Thank you

On the importance of the validation set

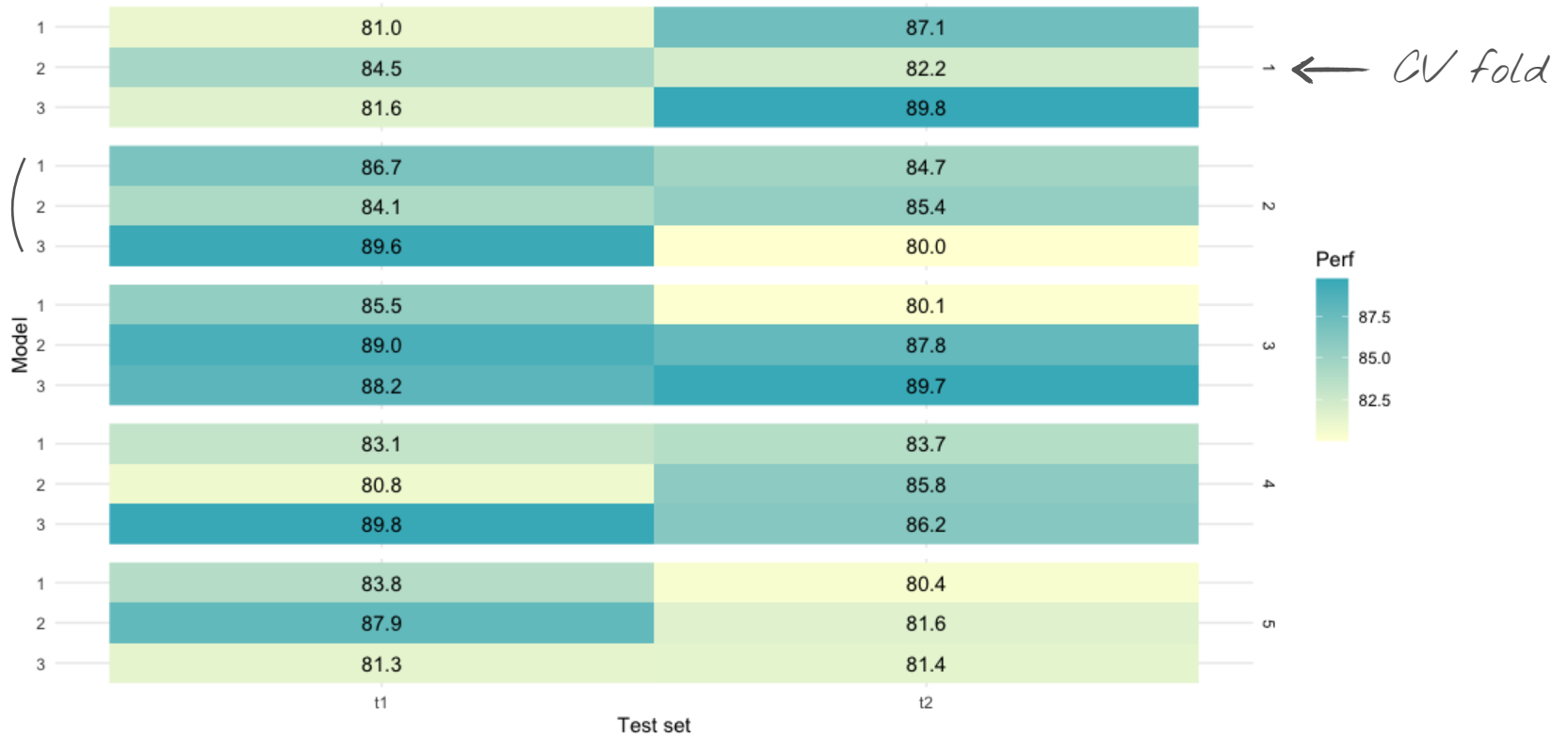
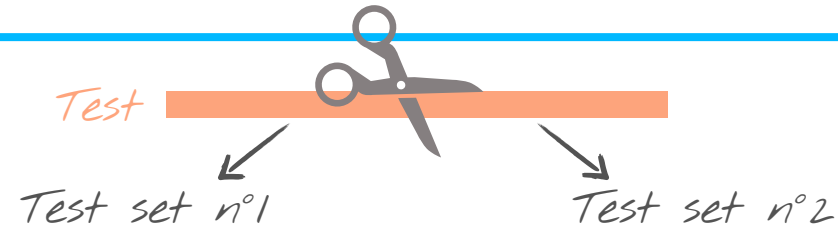
Test



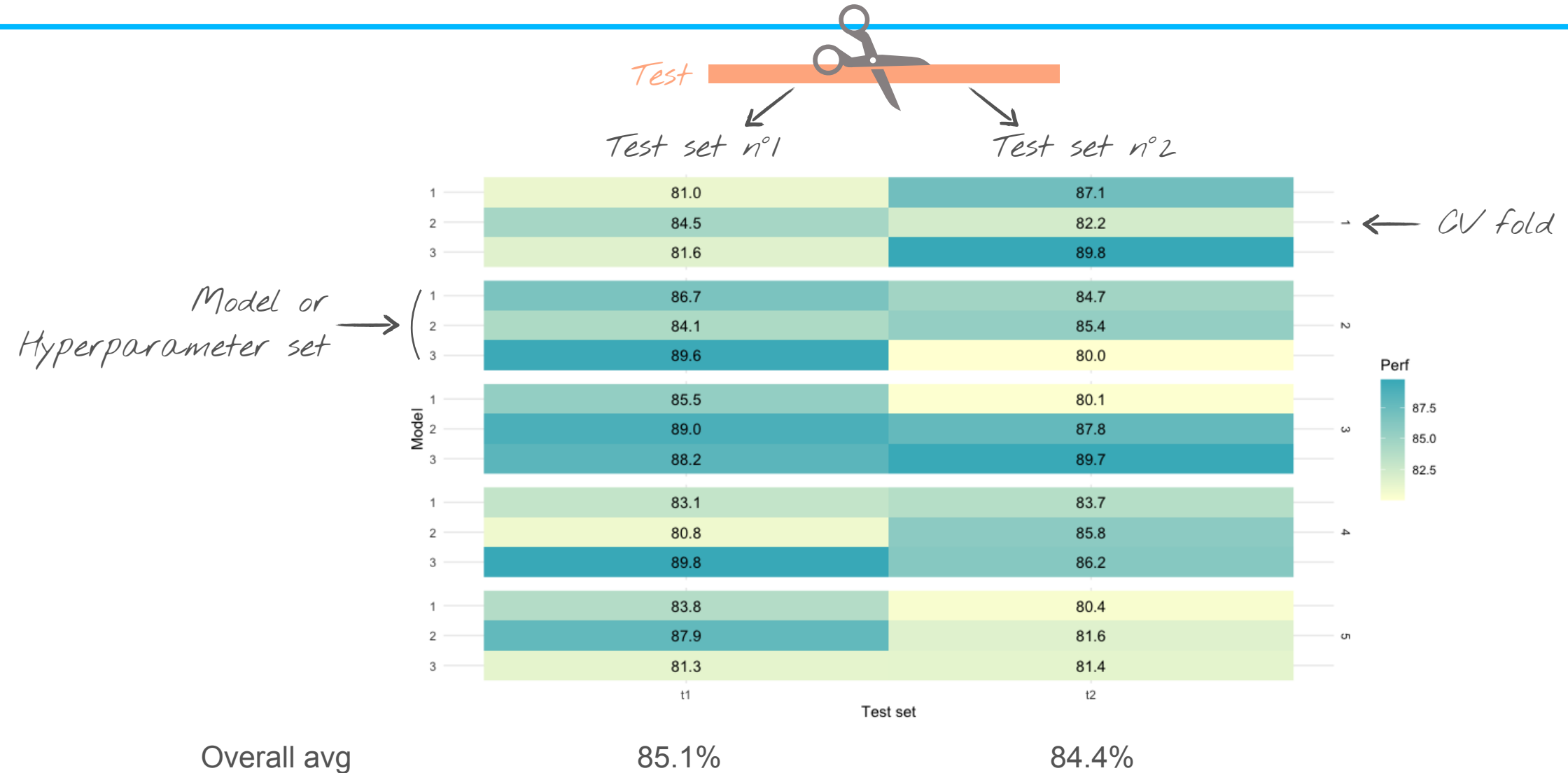
On the importance of the validation set



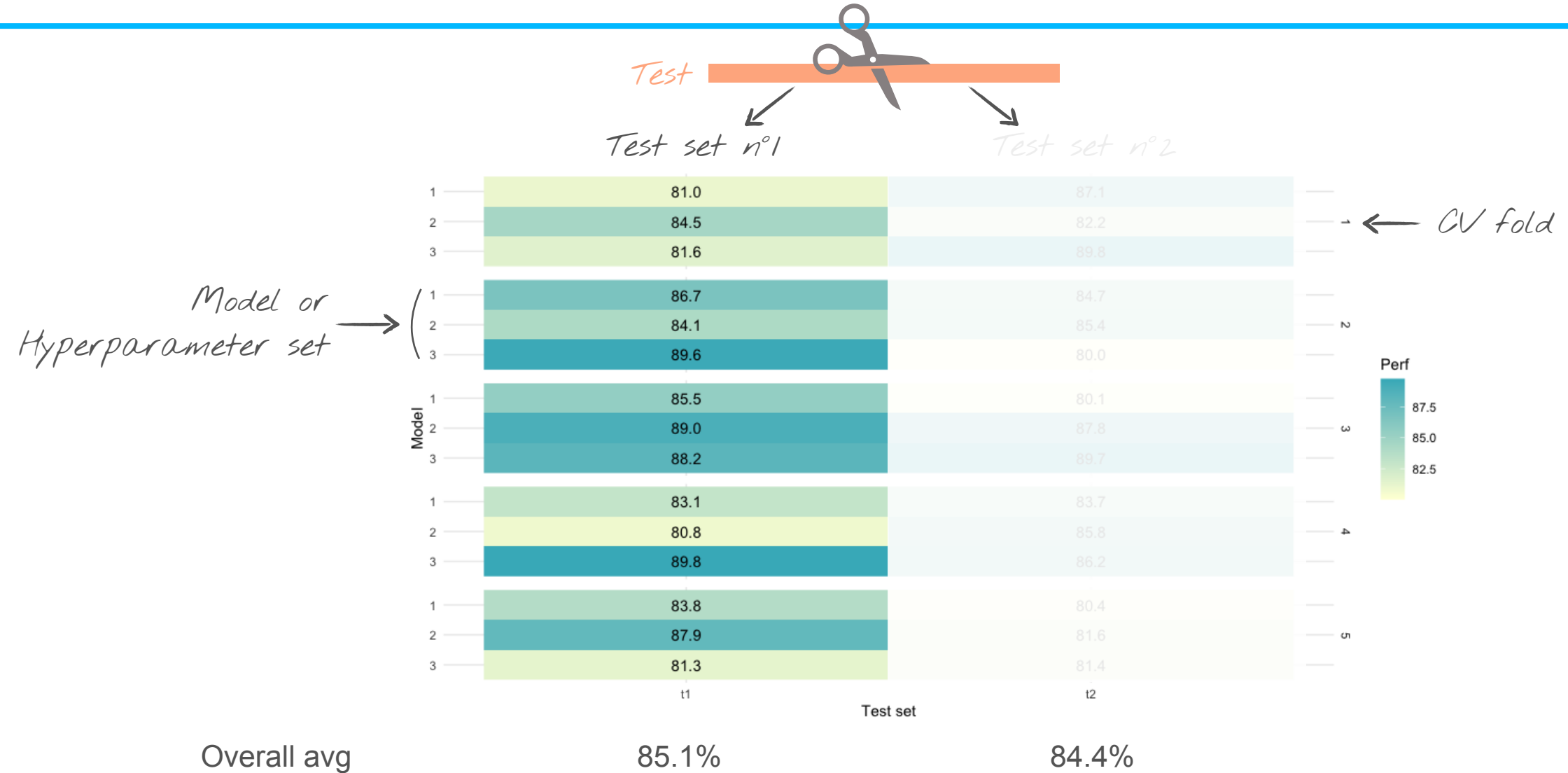
On the importance of the validation set



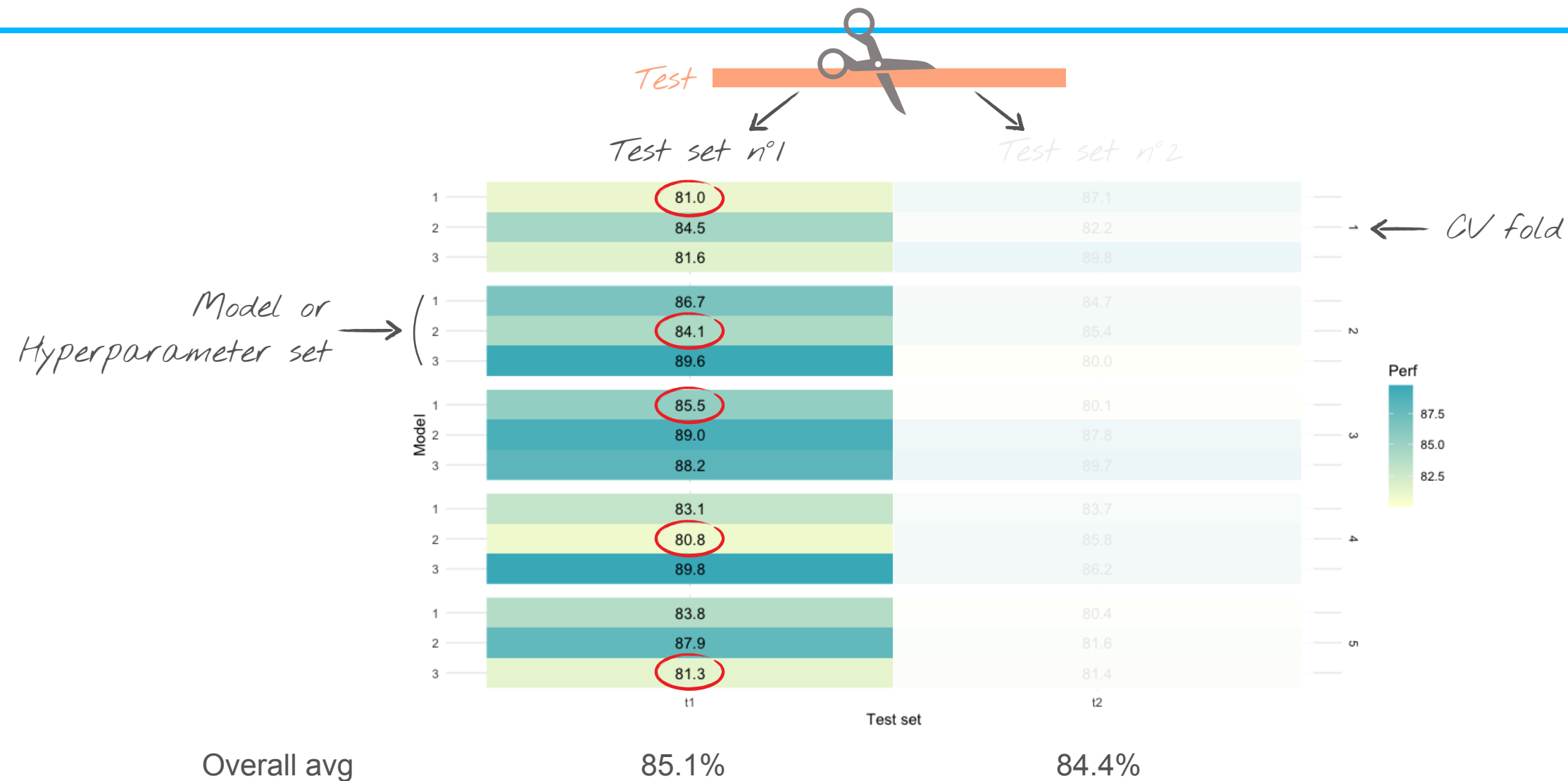
On the importance of the validation set



On the importance of the validation set



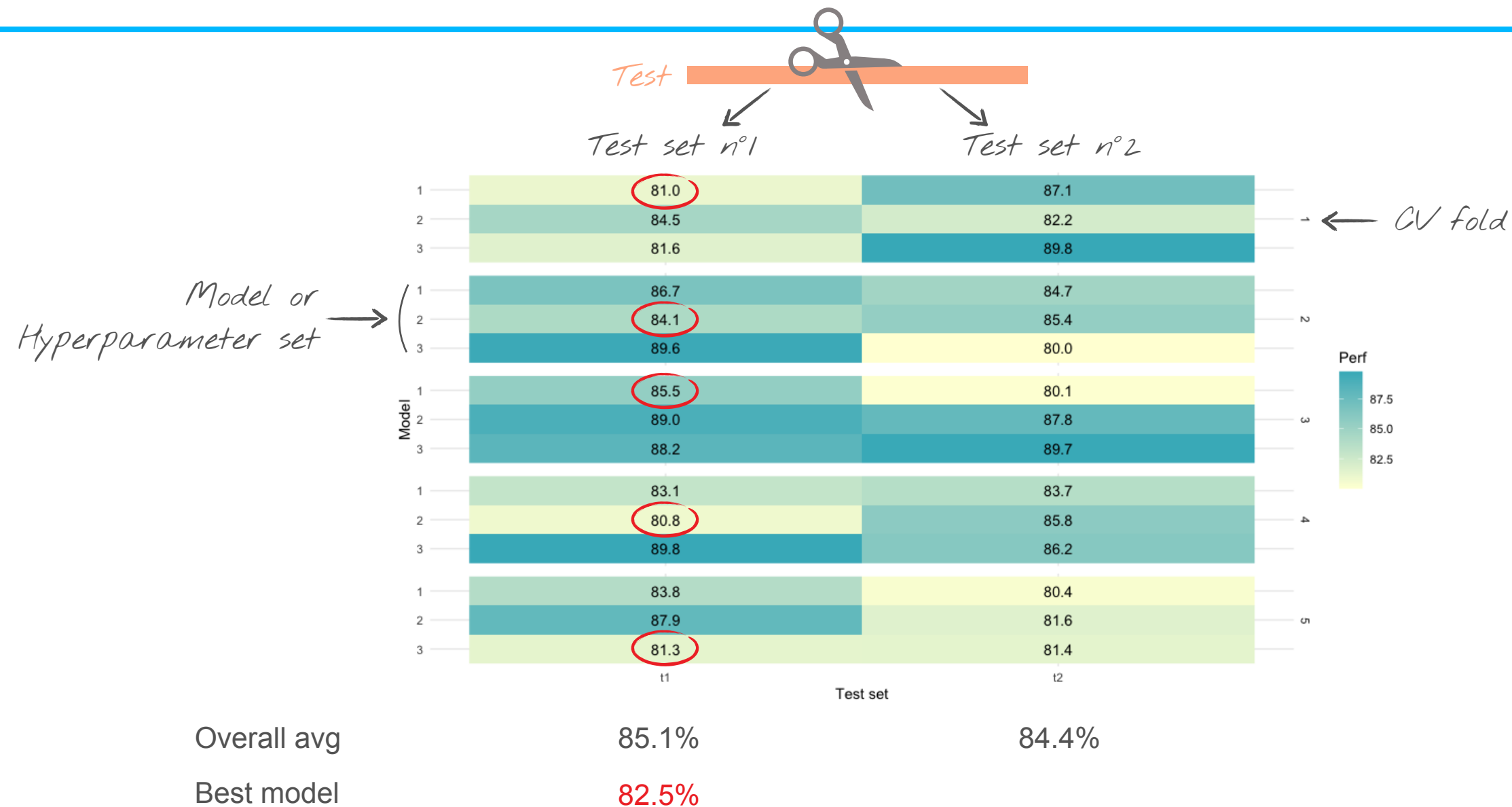
On the importance of the validation set



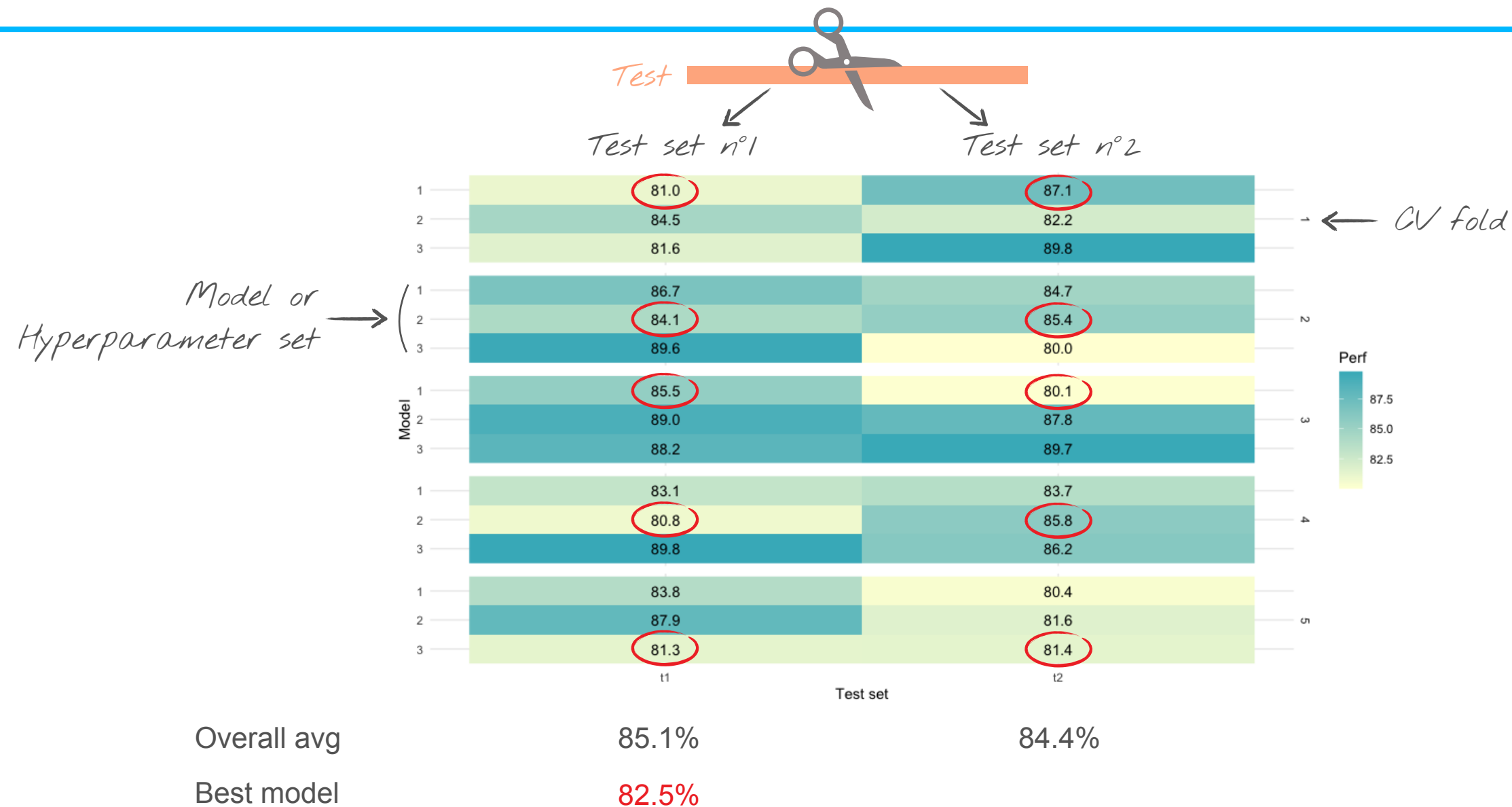
On the importance of the validation set



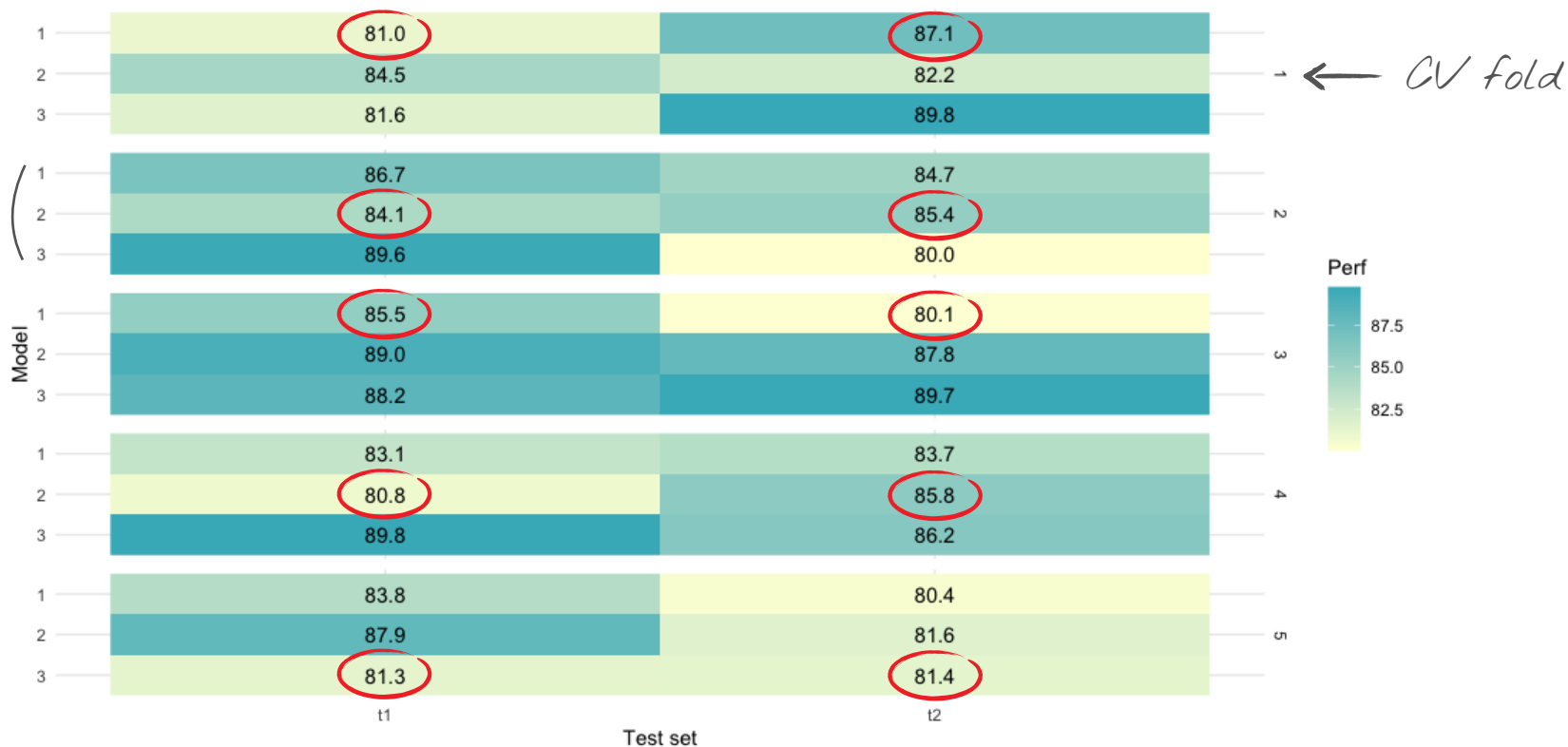
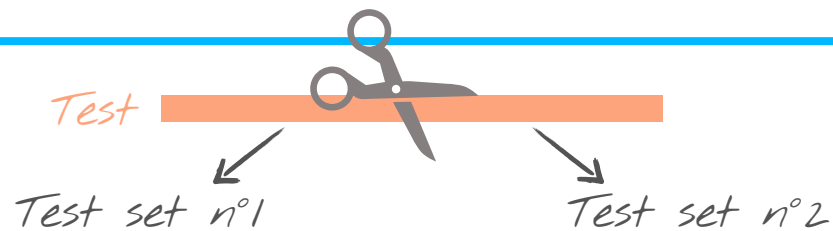
On the importance of the validation set



On the importance of the validation set



On the importance of the validation set



Overall avg

85.1%

84.4%

Best model

82.5%

84.0%

On the importance of the validation set

Validation 

optimising hyperparameters

Test 

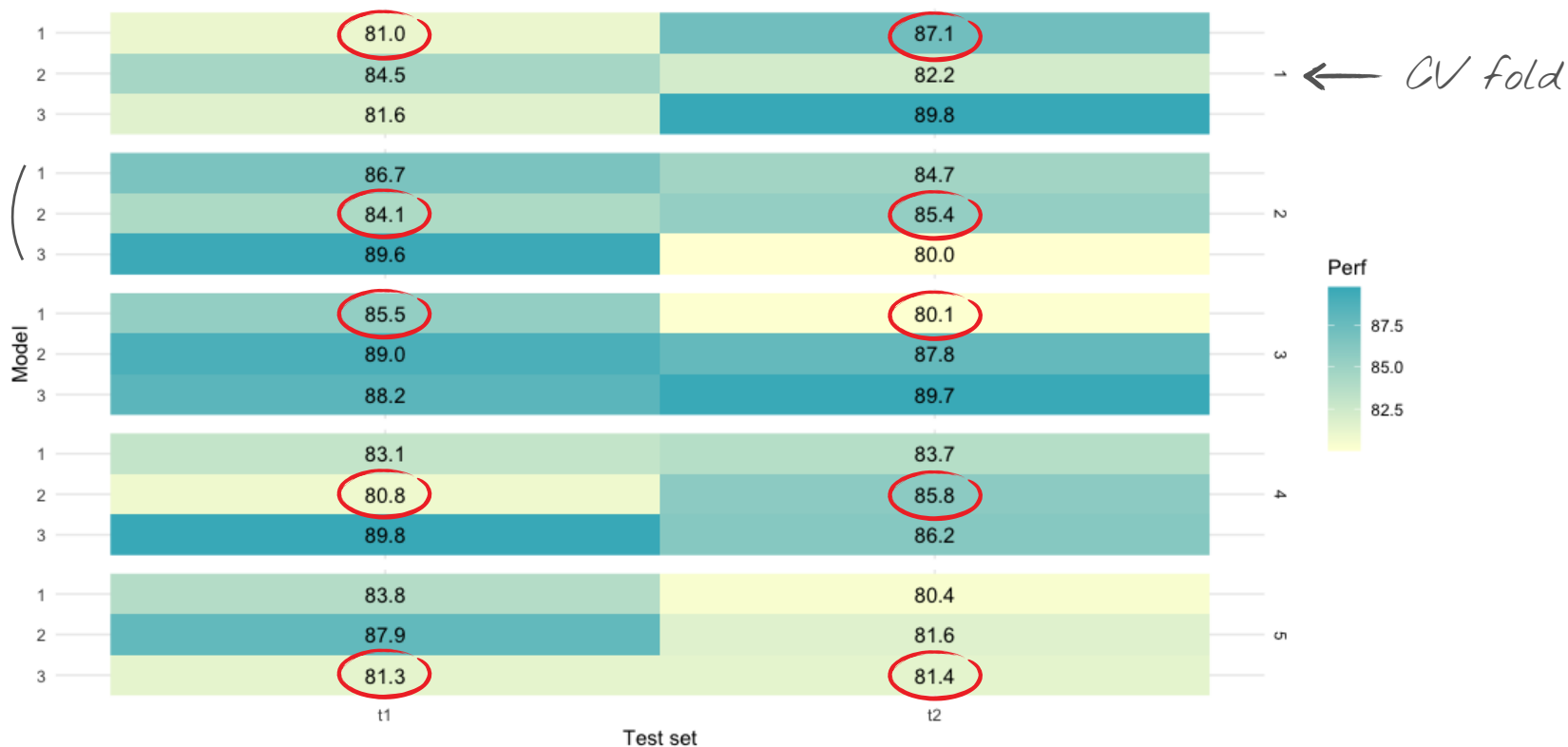
Test set n°1

Test set n°2

Test 

estimating generalisation error

Model or
Hyperparameter set →



Overall avg

85.1%

84.4%

Best model

82.5%

84.0%