

Multivariate statistics in R

Hannes PETER Martin BOUTROUX



Recap

- data exploration
 - summary statistics
 - visualization
- transformations
- resemblance
 - dis/similarity, distance
- unsupervised classification
 - cluster analysis



Recap

(Numerical) Classification

Unsupervised

search for main gradients and homogeneous groups in the data.

- No a priori knowledge/assumptions
- Results depend mainly structure of the dataset.
- distance/similarity metric, choice of clustering method
- assignment of samples into groups may change even with slight changes of the dataset (e.g. by adding more samples)
- examples of unsupervised methods are cluster analysis, TWINSPAN

Supervised

use external criteria to classify the dataset

- you supply information/rules about how to classify
- assignment of samples to groups remain the same despite changes in the structure of the dataset
- examples are classification and regression trees (CART), random forest classifier, artificial neural networks (ANN), etc.

(k-means clustering, can either be supervised or unsupervised)



Supervised classification

Classification tree analysis (CT) for qualitative response variables. Classification of new samples using a decision tree.

Regression tree analysis (RT) for <u>quantitative</u> response variables.

Classification and regression tree analysis (CART) combines these two procedures.

Random forest classifier combines bootstrapping (repeated random subsampling) of variables to produce more robust decision trees



Classification Trees

Uses two datasets

- Univariate response
- Multivariate explanatory variables
- Divides the dataset into groups (nodes)
- Applies a logical condition for each division (binary splits)
- Construction of a decision tree
 - Reads from top to bottom (divisive)
 - Allows discriminating among explanatory variables



Decision Trees

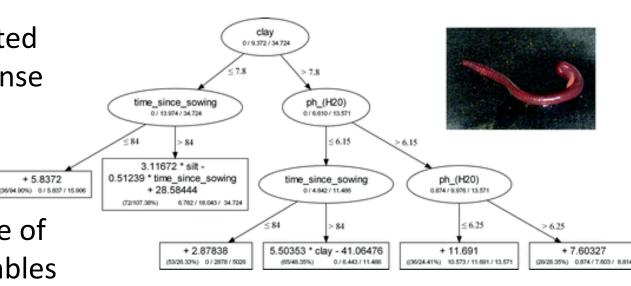
> tree-like diagram resulting from repeated splitting of the response data (dichotomous prevision model)

> shows the influence of the explanatory variables at each split

+ 5.8372

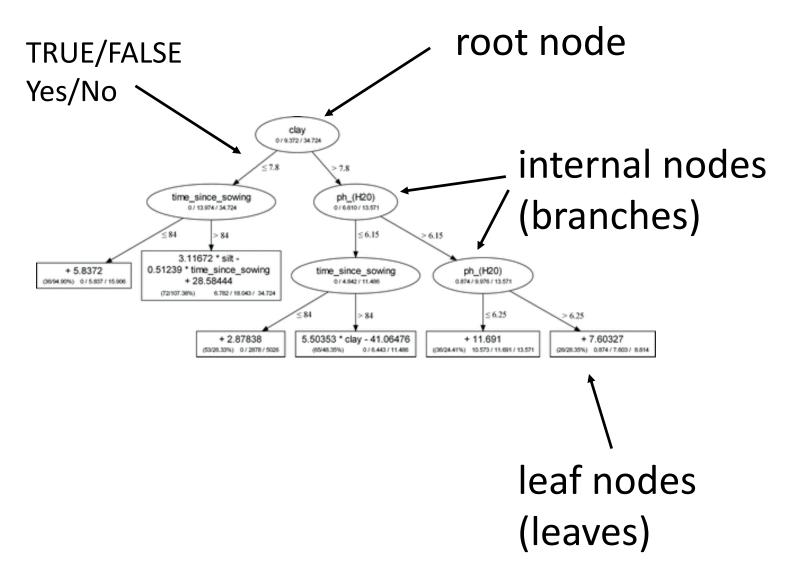
> allows prediction/decisions

Which factors determine the biomass of earthworms?





Decision Tree





an example:

What predicts whether or not a person loves Cool as Ice?

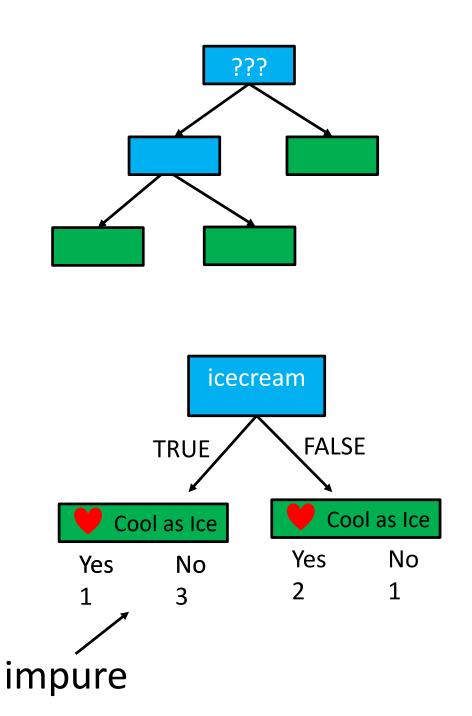
	Loves Icecream	Drinks coke	Age	Loves Cool as Ice?
Person 1	Yes	Yes	7	No
Person 2	Yes	No	12	No
Person 3	No	Yes	18	Yes
Person 4	No	Yes	35	Yes
Person 5	Yes	Yes	38	Yes
Person 6	Yes	No	50	No
Person 7	No	No	88	No



find the root...

icecream, beer, age?

	Loves Icecream	Drinks coke	Age	Loves Cool as Ice?
person1	Yes	Yes	7	No
person2	Yes	No	12	No
person3	No	Yes	18	Yes
person4	No	Yes	35	Yes
person5	Yes	Yes	38	Yes
person6	Yes	No	50	No
person7	No	No	88	No

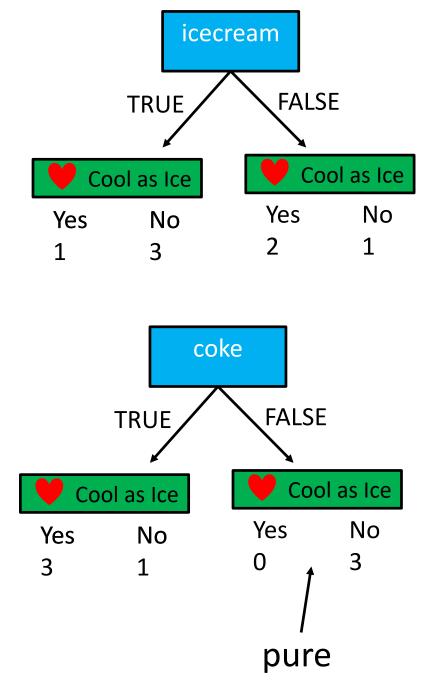




find the root...

icecream, coke, age?

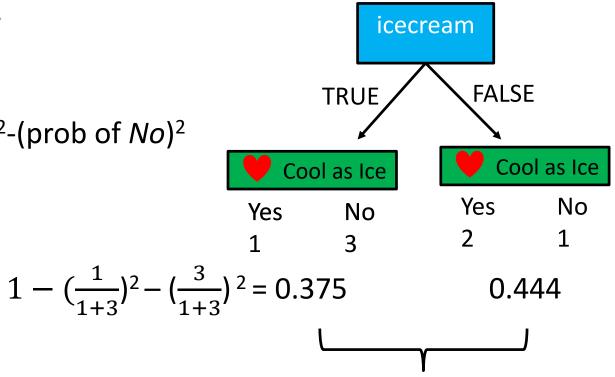
	Loves Icecream	Drinks coke	Age	Loves Cool as Ice?
person1	Yes	Yes	7	No
person2	Yes	No	12	No
person3	No	Yes	18	Yes
person4	No	Yes	35	Yes
person5	Yes	Yes	38	Yes
person6	Yes	No	50	No
person7	No	No	88	No





Gini Impurity

Gini = 1-(prob of Yes)²-(prob of No)²

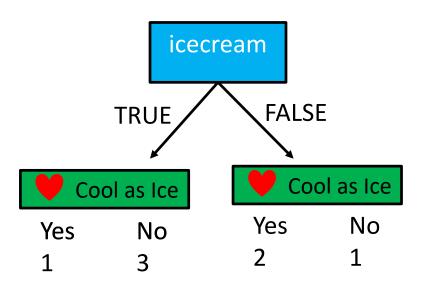


Gini Impurity: weighted average of leaf impurities

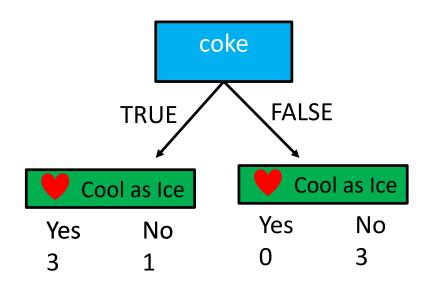
$$\left(\frac{4}{4+3}\right)*0.375 + \left(\frac{3}{4+3}\right)*0.44 = 0.405$$



find the root...

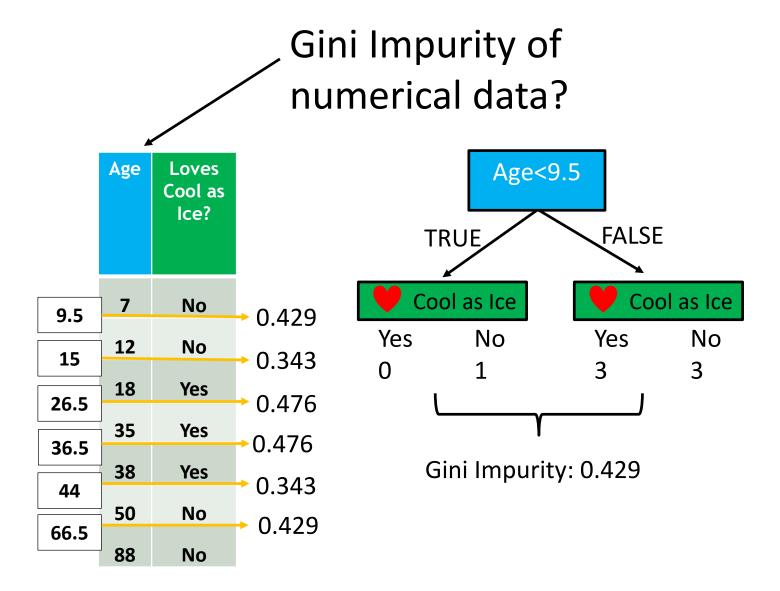


Gini Impurity: 0.405



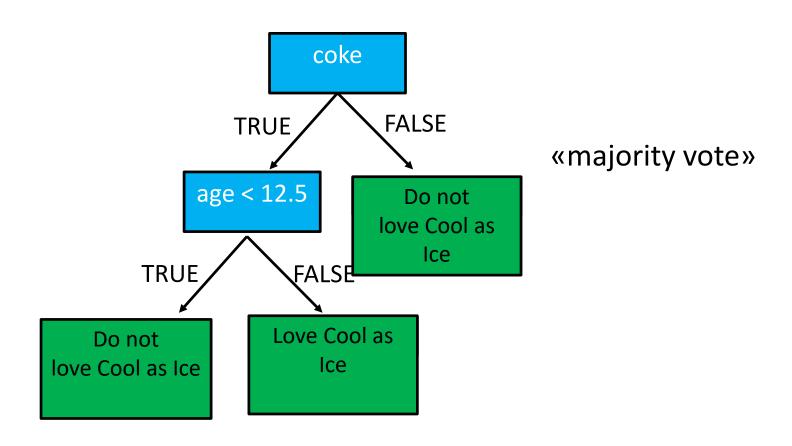
Gini Impurity: 0.214







keep splitting impure nodes...

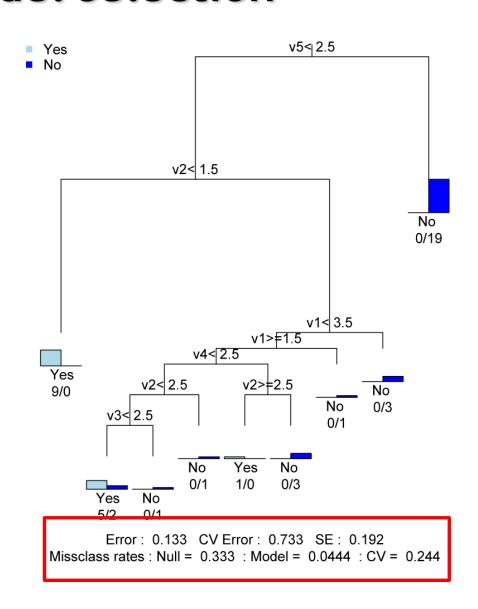




model error and model selection

Model error (Error)

- misclassifications (errors) of the terminal leaves
- trends to 0 with increasing tree size
- Prediction error, Relative error or Cross Validation Error (CV)
 - sum of the errors of the terminal branches achieved by cross-validation
 - measures uncertainty of forecasting
 - reaches a minimum for an optimal tree size (=> pruning)





Regression trees

- Quantitative response
 - Synthetic variable representing a global feature
 - Ex.: diversity index of the community, density of an invasive species, etc...
 - Aim: evaluate the response variable according to the explanatory variables and forecast the values of this quantitative descriptor from explanatory variables



dune vegetation dataset (R::vegan)

- cover (Braun-Blanquet classes, 0 9) of 30 plant species at 20 dune meadow sites (2 x 2 m) in the Netherlands Batterink & Wijfels 1983
 - species names are abbreviated (4 genus + 4 species letters, ex: Achimill = Achillea millefolium)

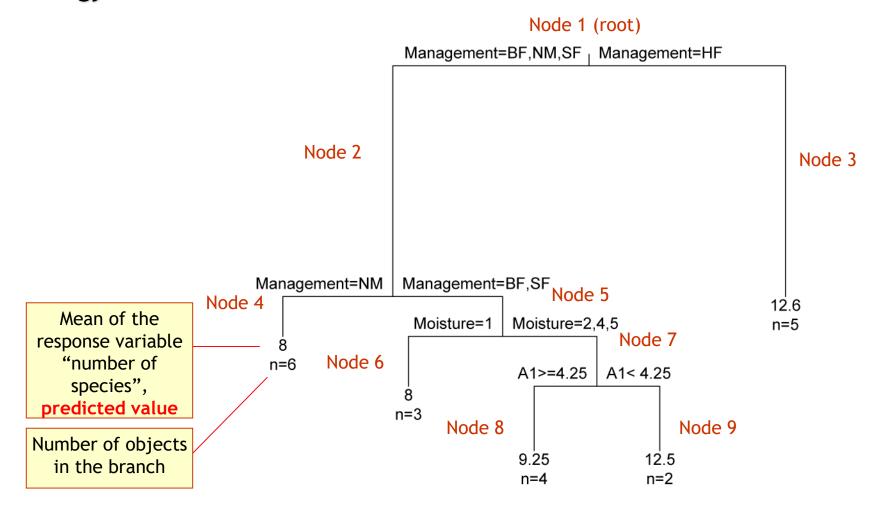
dune.env

- □ 20 observations of environmental data
 - A1: thickness of soil A1 horizon (numeric)
 - ☐ Moisture: soil moisture with levels 1-5
 - Management: biological farming (BF), Hobby Farming (HF), NM (Nature and Conservation Management), SF (Standard farming) (factors)
 - Use: land-use with levels Hayfield, Haypastu, Pasture
 - ☐ Manure: factor with levels 0-4



Univariate regression tree

How does the number of plant species depend on management strategy?

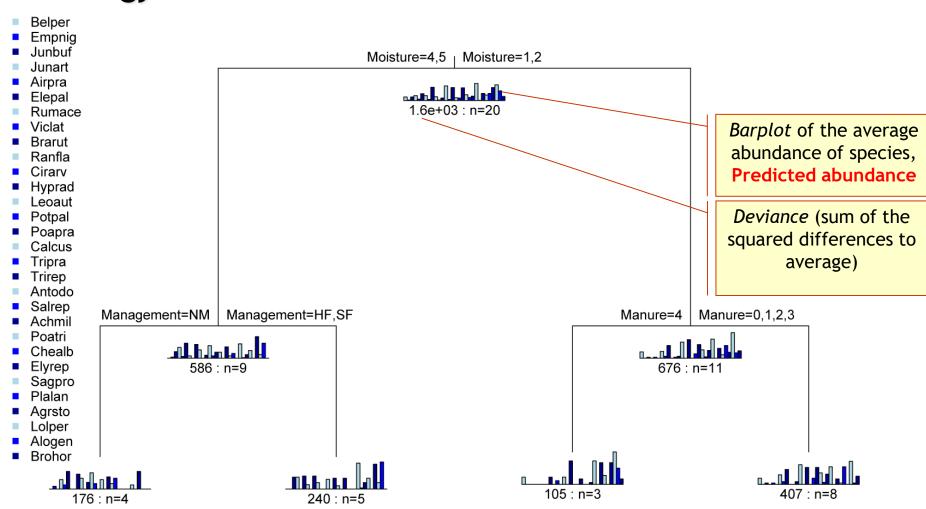


Error: 0.544 CV Error: 0.86 SE: 0.268



multivariate regression tree How do species abundances depend on ma

How do species abundances depend on management strategy?

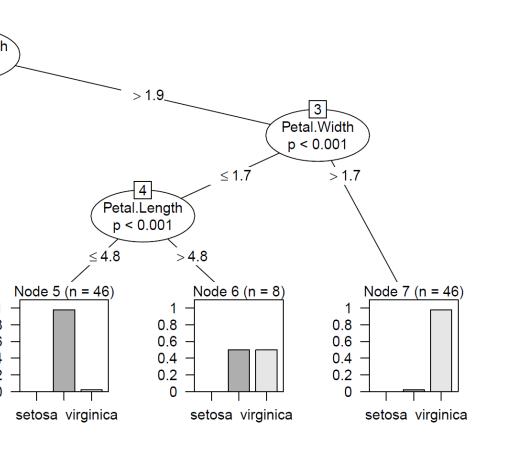


Error: 0.58 CV Error: 1.07 SE: 0.0915



Implementation in R

partykit::ctree (Conditional Inference Trees)



But see also: rpart

- Combines recursive binary partitioning with permutation tests (Bonferroni-adjusted pvalue for each node)
- Stops when no significant associations (regression relationship) between any predictor and the response can be found (no tree pruning necessary)
- Can handle weighted predictors



Advantages of CART methods

- powerful tool for data mining
- easy to build and interpret decision trees
- robust and flexible technique
 - All sorts of variables (binary, multi-class, ordered)
 - Accepts missing data
 - No assumptions on the variable distribution and relationships

BUT: Trees have one aspect that prevents them from being the ideal tool for predictive learning, namely inaccuracy.

=> Random forest classifiers



Random Forest Classifiers

- ▶ Use a large number of decision trees (=forest) each built with a different, random sub-sample of the dataset (bootstraping) and only a subset of explanatory variables
- ▶ Uses all decision trees for making a prediction. Criteria: the most frequent pattern (majority vote).



Random Forest Classifiers

- Random Forest trees are generated using bootstraping Bootstrap: resampling the dataset with replacement for constructing the decision tree
- ► For each node, only a subset of explanatory variables are taken randomly
 - iterate using different numbers of explanatory variables
- Generate many trees
- run data through all trees and measure the outcomes, aggregate all outcomes to make a decision (prediction).
- Bootstrap and aggregate => "bagging"
- use out-of-bag samples to evaluate random forest classifier performance



example: Use V1, V2, BMI, Sport to predict health state

bootstrap samples

					•		ii ap 3	arripic.		•
Health state	V1	V2	BMI	Sport?		Health state	V1	V2		ВМІ
III	0	1	25	No		Ill	0	1		25
Well	1	1	45	No		Ill	0	1		32
Well	1	0	65	Yes		Well	1	0		65
III	0	1	32	Yes						
					•	Well	1	0		65 I
								elect n	е	explan
_				l C.		Health state	V1 _{Va}	r PAble s	;	
predict first node						Ill	0	25		
			\neg			Ill	0	22		\ \ \
	Well 1							Health state		V2
						Well	1	Ill		1
	PR	EDICTION		pred	cond nod	le		III		1
PFL		ı						Well		0

RANDOM FOREST CLASSIFIER **DATASET DECISION TREE PREDICTION PREDICTION PREDICTION MAJORITY VOTE TAKEN** FINAL PREDICTION MADE



Out-of-bag samples (OOB)

Health state	V1	V2	ВМІ	Sport?	
Ill	0	1	25	No	
Well	1	1	45	No	\
Well	1	0	65	Yes	
III	0	1	32	Yes	

Health state	V1	V2	ВМІ	Sport?
Ill	0	1	25	No
Ill	0	1	32	Yes
Well	1	0	65	Yes
Well	1	0	65	Yes

Typically ~1/3 of samples are not included by bootstrapping into random forest classifier tree generation

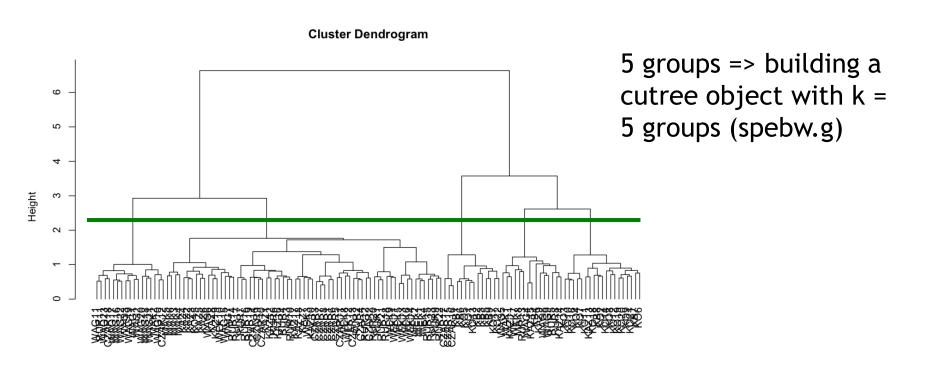
⇒ Use these samples to evaluate random forest classifier performance: (OOB error = fraction correctly classified OOB samples)



Example: microbes

Dataset: microbes and environmental variables (nutrient availability) in a wetland

132 observations, 70 species (spe) and 15 environmental variables (env)







Example: microbes

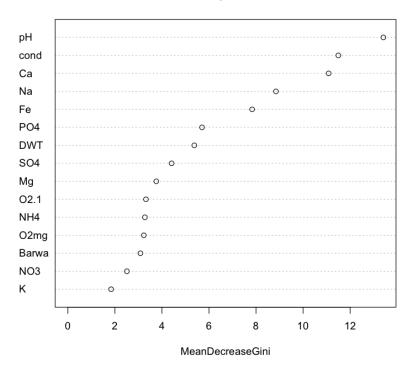
rf.env = randomForest(spebw.g~., env, ntree=500, mtry=10, proximity=T)

ntree: number of decision trees

mtry: number of variables (species) selected for each node in the tree

proximity: keep estimates of closeness of pairs of samples

rf.spebw



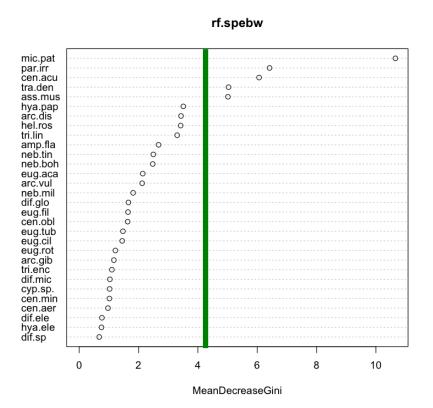
Classification of the most important variables. Gives some important explanatory variables for the given classification.

OOB estimate of error rate: 27.27% (36 missclassified out of 132 samples)



Example: microbes

rf.spe = randomForest(spebw.g~., spe, ntree=500, mtry=10, proximity=T)



Classification of the most important variables.

Gives some important response variables (species) for the given classification

OOB estimate of error rate: 9.85% (13 misclassified out of 132 samples)



Classification or regression trees

- + Useful for qualitative and quantitative prediction
- + Fast
- + accepts missing data
- + Visual discrimination
- + Tree easy to interpret
- +/- Pruning of the tree
- No option for selecting variables
- Requires a low CV Error for prediction

Random Forest

- + Useful for qualitative and quantitative prediction
- + Provides variable importance (ranking)
- + Robust (bootstraping)
- + accepts missing data (in explanatory variables and in new predictions)
- individual decision trees are not directly interpretable ('Black Box')
- sometimes slow



https://doi.org/10.1038/s41558-018-0393-5

Widespread loss of lake ice around the Northern Hemisphere in a warming world

Sapna Sharma (1)11,111, Kevin Blagrave (1)11, John J. Magnuson (2)11, Catherine M. O'Reilly (1)3,11, Samantha Oliver (4, Ryan D. Batt (1)5, Madeline R. Magee (2,6, Dietmar Straile (7, Gesa A. Weyhenmeyer (1)8, Luke Winslow (1)9 and R. Iestyn Woolway (1)

Ice provides a range of ecosystem services—including fish harvest1, cultural traditions2, transportation3, recreation4 and regulation of the hydrological cycle5—to more than half of the world's 117 million lakes. One of the earliest observed impacts of climatic warming has been the loss of freshwater ice6, with corresponding climatic and ecological consequences7. However, while trends in ice cover phenology have been widely documented^{2,6,8,9}, a comprehensive large-scale assessment of lake ice loss is absent. Here, using observations from 513 lakes around the Northern Hemisphere, we identify lakes vulnerable to ice-free winters. Our analyses reveal the importance of air temperature, lake depth, elevation and shoreline complexity in governing ice cover. We estimate that 14,800 lakes currently experience intermittent winter ice cover, increasing to 35,300 and 230,400 at 2 and 8 °C, respectively, and impacting up to 394 and 656 million people. Our study illustrates that an extensive loss of lake ice will occur within the next generation, stressing the importance of climate mitigation strategies to preserve ecosystem structure and function, as well as local winter cultural heritage.

in some winters^{2,11}. This transitional period from annual winter ice to permanent loss of ice cover may endure for decades². The factors influencing whether or not ice forms are well known; previous research has indicated that air temperature, wind speed, and lake size are essential components to ensure that vertical heat transfer is sufficient to cool surface water temperatures to 0°C^{12,13}. Precipitation¹², snow cover¹⁴, cloud cover, solar radiation¹⁴, distance to coastline¹⁵ and regional differences^{7,16} can govern the timing of ice formation and ice growth during the winter season. However, previous research has not identified how the interactions between features such as climate and lake shape (area and depth) will dictate when and where the threat of lake ice loss is greatest. We provide the first global estimate of how many lakes are likely to lose annual winter ice cover as the climate warms.

We used updated lake ice cover records for 346 lakes in North America, 136 lakes in Europe, and 32 lakes in Asia to evaluate the threat of lake ice loss¹⁷ (Supplementary Fig. 1). Lakes were designated as annual or intermittent winter ice-covered lakes. Annual ice-covered lakes experienced complete ice cover every winter, whereas intermittent ice-covered lakes had one or more winters

