# Machine learning for ecology with R
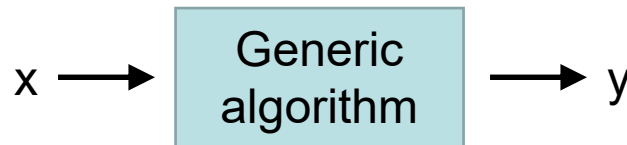
Brett Melbourne
Associate Professor, EBIO
brett.melbourne@colorado.edu
Pronouns: he, him, his

# GitHub

- Code for this presentation
    - github.com/melbourne-lab/ml4e-nutshell


- Machine learning for ecology graduate course
    - github.com/EBIO5460Spring2022

# What is machine learning?

- Using generic algorithms to predict outputs y from inputs x

- Emphasis: predictive skill

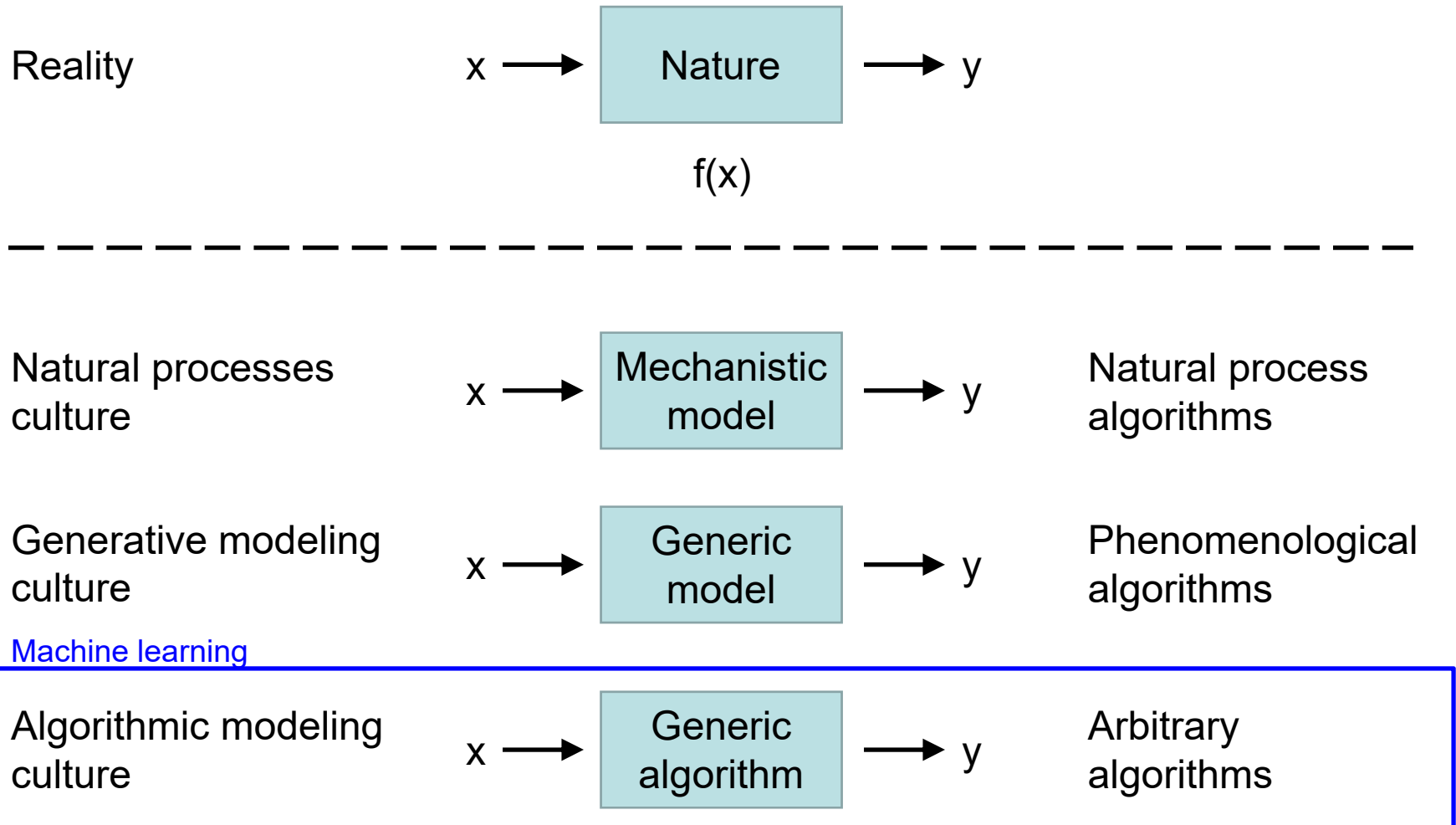x $\longrightarrow$ [Generic algorithm] $\longrightarrow$ y

# Examples in ecology

- Species distribution models (SDMs)
  - predicting the spatial distribution of a species from explanatory variables
- Counting the number of penguins in Antarctica from satellite imagery
- Identifying species in camera trap images
- Identifying bird species from audio recordings

# Trying to learn a function f

Reality

$x \longrightarrow$ Nature $\longrightarrow$ y

f(x)

# Trying to learn a function f



Reality    x ⟶ [Nature] ⟶ y

f(x)

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Natural processes culture    x ⟶ [Mechanistic model] ⟶ y    Natural process algorithms

Generative modeling culture    x ⟶ [Generic model] ⟶ y    Phenomenological algorithms

Machine learning

Algorithmic modeling culture    x ⟶ [Generic algorithm] ⟶ y    Arbitrary algorithms

f can mean different things in different data science cultures

# Goal of prediction

Use data to find a function $\hat{f}$ that has good predictive performance given X

That is, $\hat{f}$ is accurate on new observations

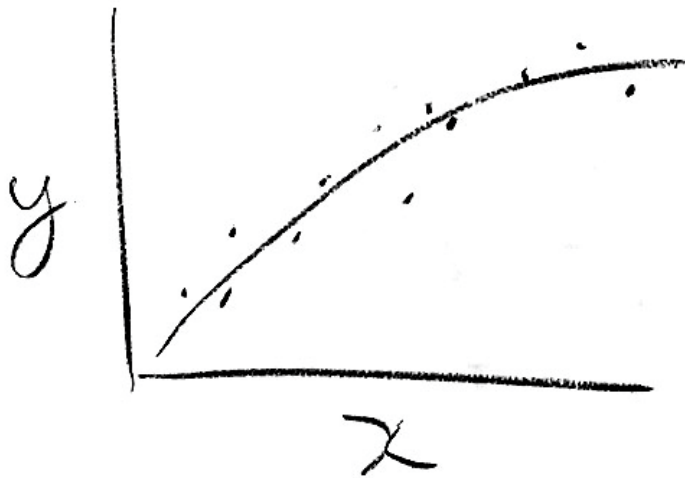# Goal of machine learning

## To predict accurately!

- Species distribution
  - map
  - predict accurately for places we won't visit
- Climate change forecast
  - predict accurately for the future
- Antelopes in camera trap images
  - hand over the identification task to a machine so we don't have to look at images!
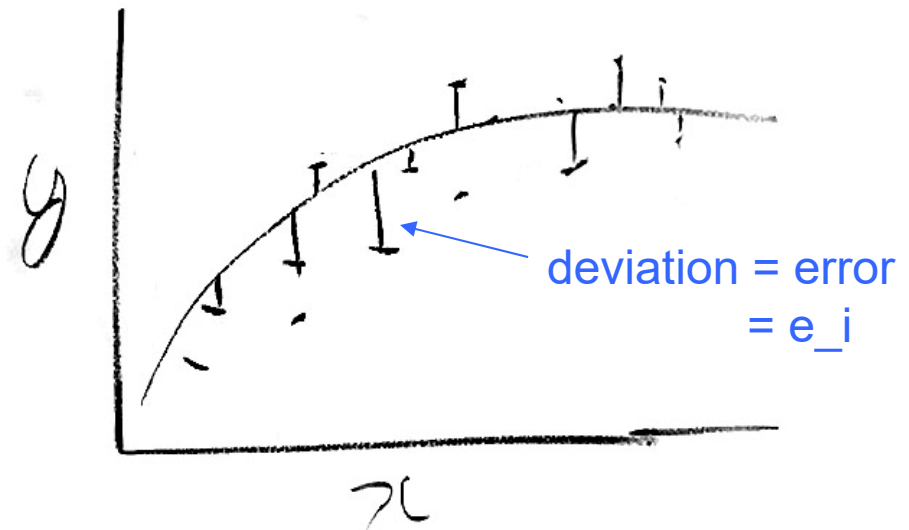  - predict accurately for images that we'll never look at

# Predictive skill

Basic idea: out-of-sample accuracy

$\hat{f}$ fitted on training data

$\hat{f}$ predicting new data



deviation = error
= e_i

e.g. mean square error (MSE) $\dfrac{1}{n}\sum_i^n e_i^2$

# Basic full ML setup

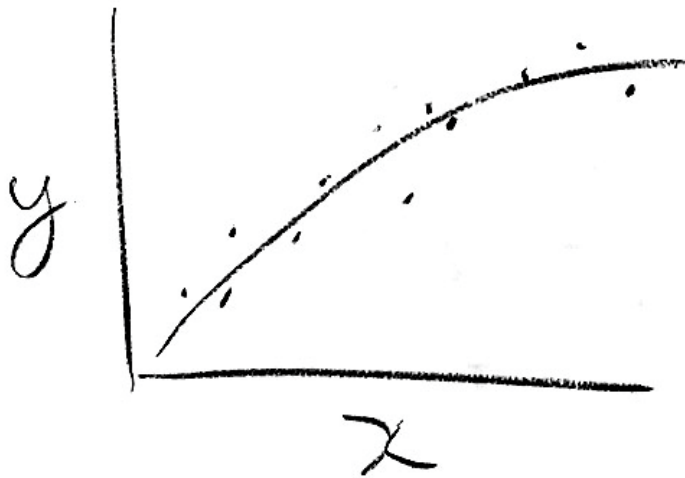<span style="color:blue">Overall algorithm:</span>

1. Create a <span style="color:blue">model algorithm</span> for $\hat{f}(x)$
2. Use a <span style="color:blue">training algorithm</span> to find parameter values of $\hat{f}(x)$
3. Use an <span style="color:blue">inference algorithm</span> to compare predictive skill among models (model families, tuning parameters, $x$ sets, etc).
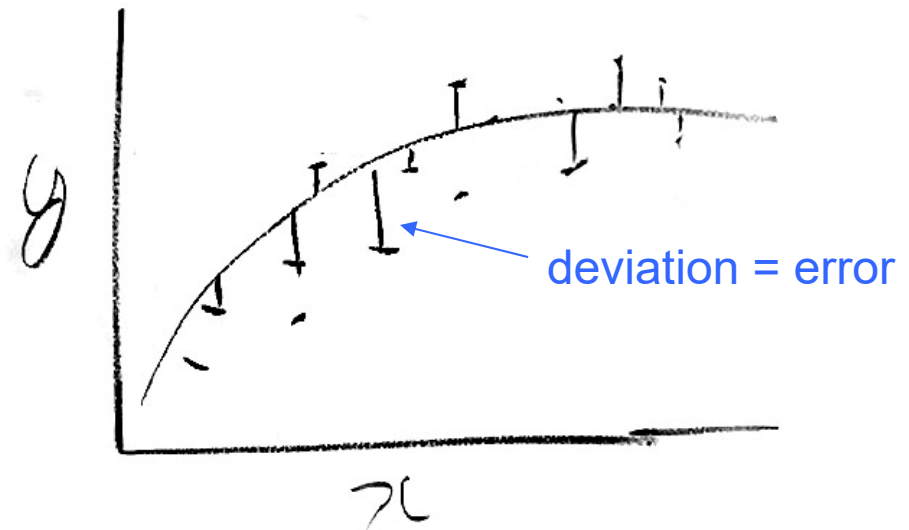
# Inference algorithm

Basic idea: out-of-sample validation

Fit model to training dataset

Test model on validation dataset
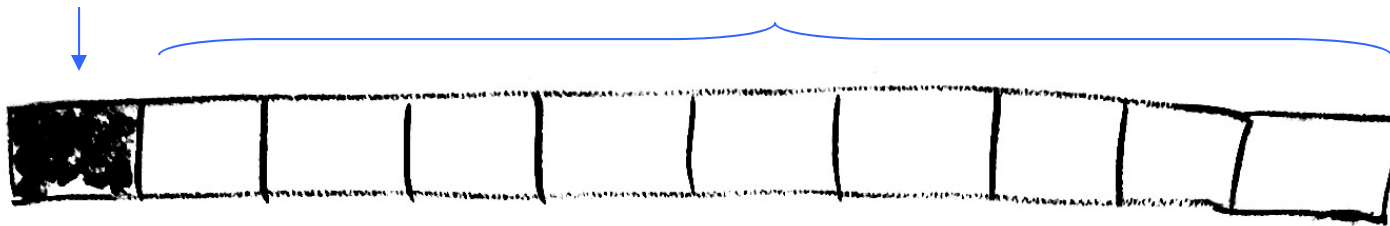


deviation = error

e.g. mean square error (MSE)

# k-fold cross validation (CV)

Divide dataset into k parts (preferably randomly)

test (validation) data
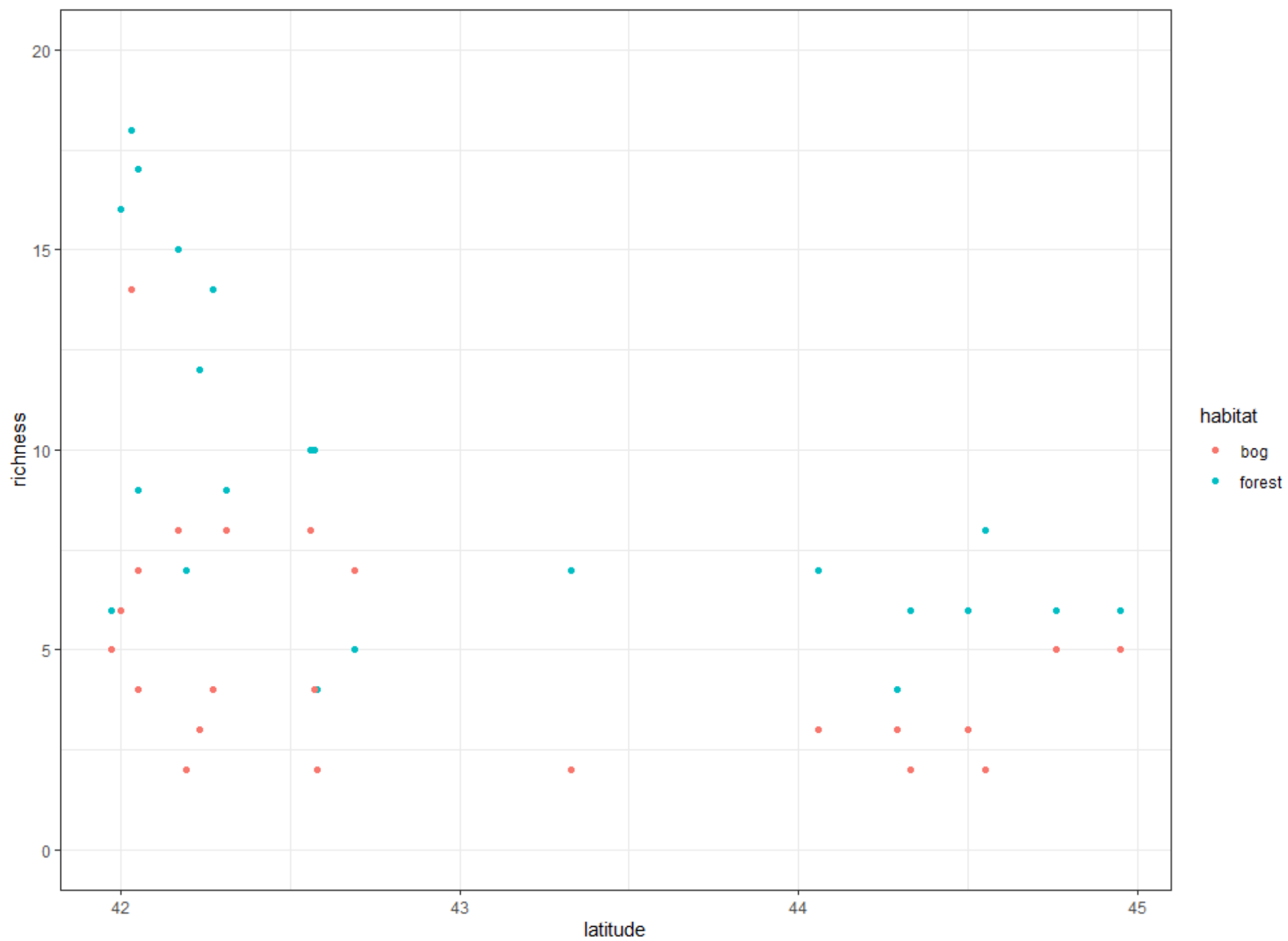
training data

repeat with next test subset

... repeat with each test subset

# Regression & classification

- Regression:
  - numerical response variable
  - predict a numerical value given x
  - e.g. number of species given latitude
- Classification:
  - categorical response variable
  - predict the category given x
  - e.g. is it a bird, deer, tree, or mountain lion?
  - e.g. is it dead or alive?; present or absent?

# Ants data

```
> head(ants)
  site habitat latitude elevation richness
1  TPB  forest    41.97       389        6
2  HBC  forest    42.00         8       16
3  CKB  forest    42.03       152       18
4  SKP  forest    42.05         1       17
5   CB  forest    42.05       210        9
6   RP  forest    42.17        78       15
```
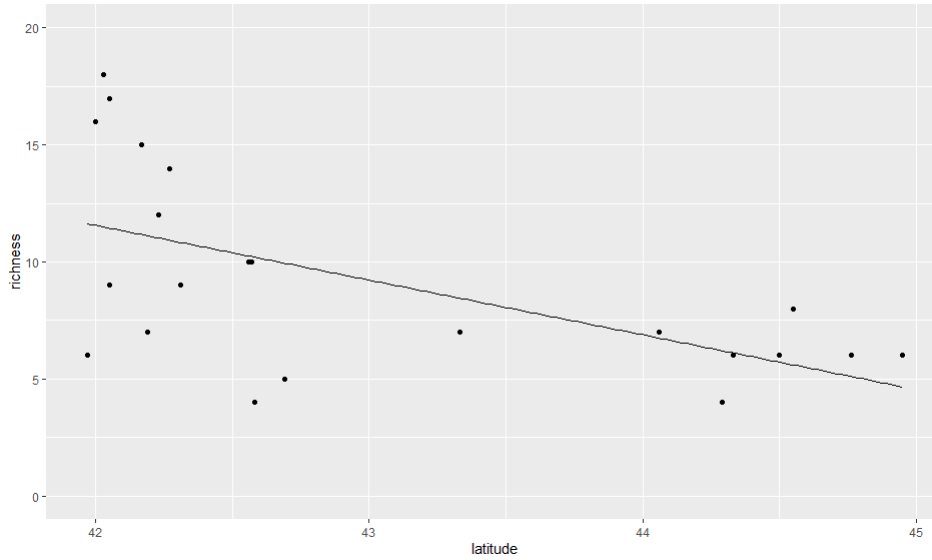
# Basic full ML setup

- Polynomial example, 3 algorithms:
  - model: flexible function $\hat{f}(x)$; polynomial linear model

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \ldots + \beta_m x^m \qquad \text{m=order}$$

# Basic full ML setup
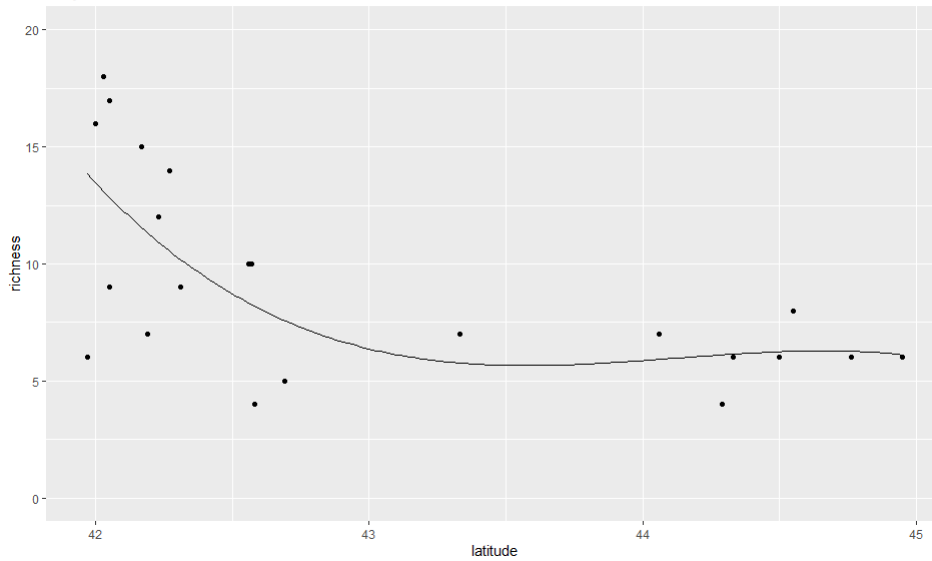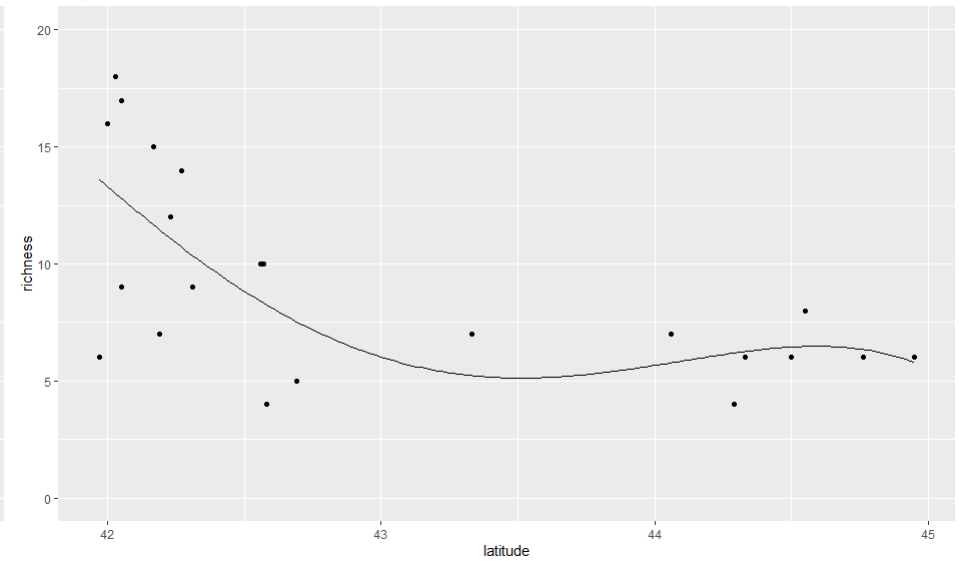
- Polynomial example, 3 algorithms:
  - model: flexible function $\hat{f}(x)$; polynomial linear model

    $$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \ldots + \beta_m x^m \qquad \text{m=order}$$

  - training: optimize least squares objective function
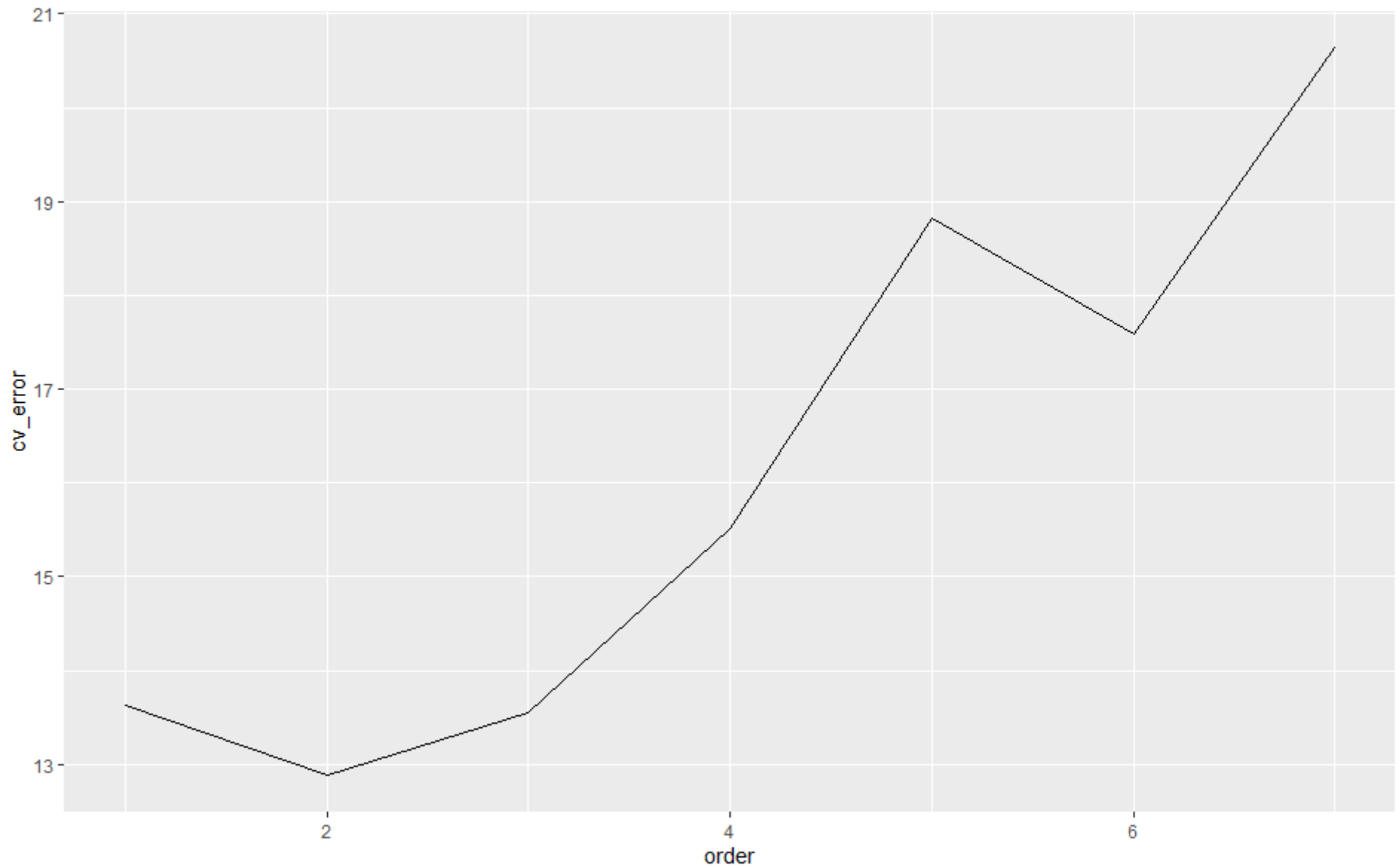  - minimize $SSQ = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2$ for training data

```
lm(richness ~ poly(latitude, order), data=forest_ants)
```

# Basic full ML setup

- Polynomial example, 3 algorithms:
  - model: flexible function $\hat{f}(x)$; polynomial linear model

    $$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \ldots + \beta_m x^m \qquad \text{m=order}$$

  - training: optimize least squares objective function
  - minimize $SSQ = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2$ for training data
  - inference: tuning parameter (order of poly); k-fold cross validation

# Inference algorithm

k-fold cross validation

```r
cv_ants <- function(k, order) {
    forest_ants$fold <- random_folds(nrow(forest_ants), k)
    e <- rep(NA, k)
    for ( i in 1:k ) {
        test_data <- forest_ants %>% filter(fold == i)
        train_data <- forest_ants %>% filter(fold != i)
        poly_trained <- lm(richness ~ poly(latitude, order), data=train_data)
        pred_richness <- predict(poly_trained, newdata=test_data)
        e[i] <- mean((test_data$richness - pred_richness) ^ 2)
    }
    cv_error <- mean(e)
    return(cv_error)
}
```
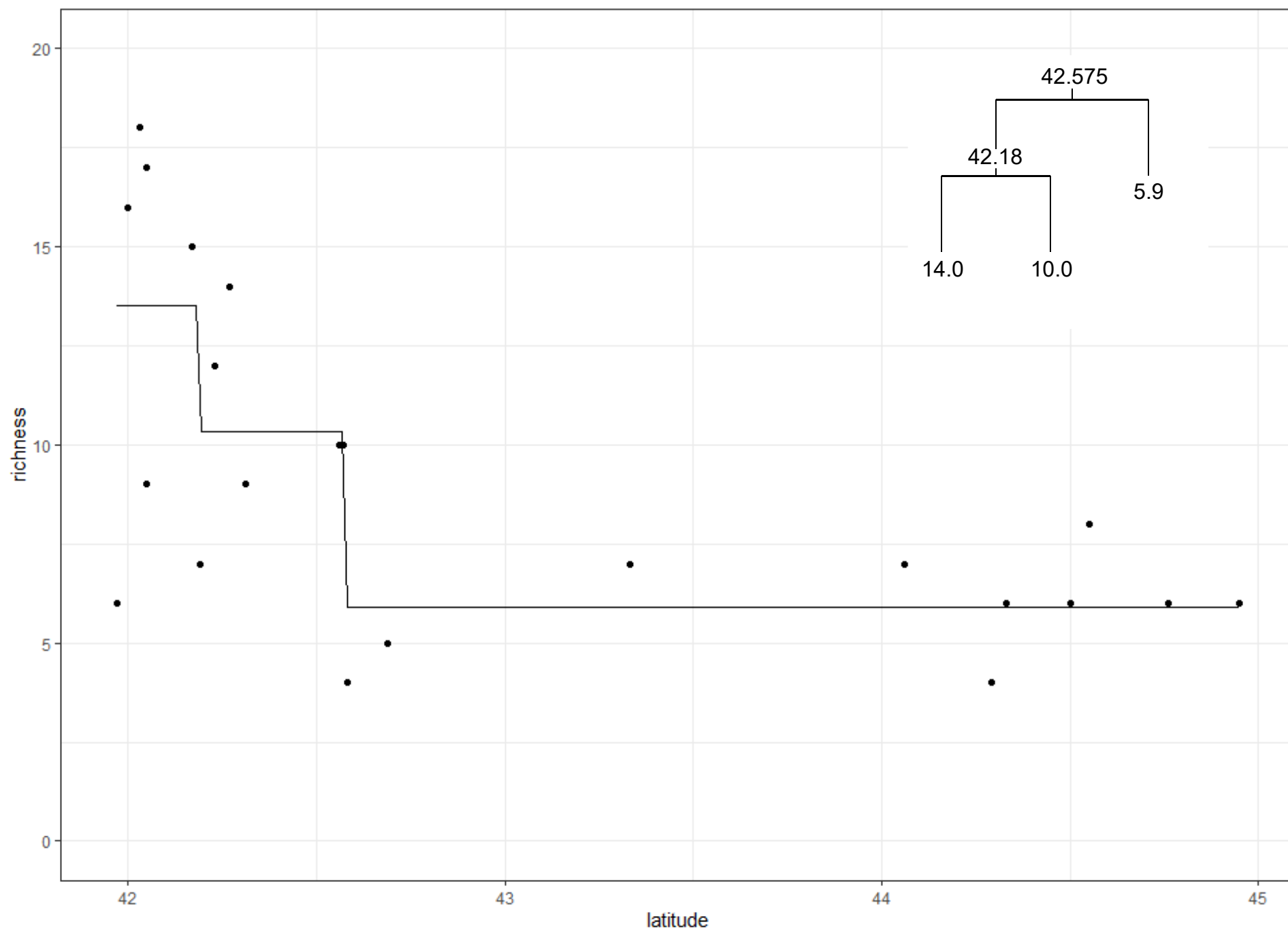
# Leave one out CV (k-fold CV, k=n)
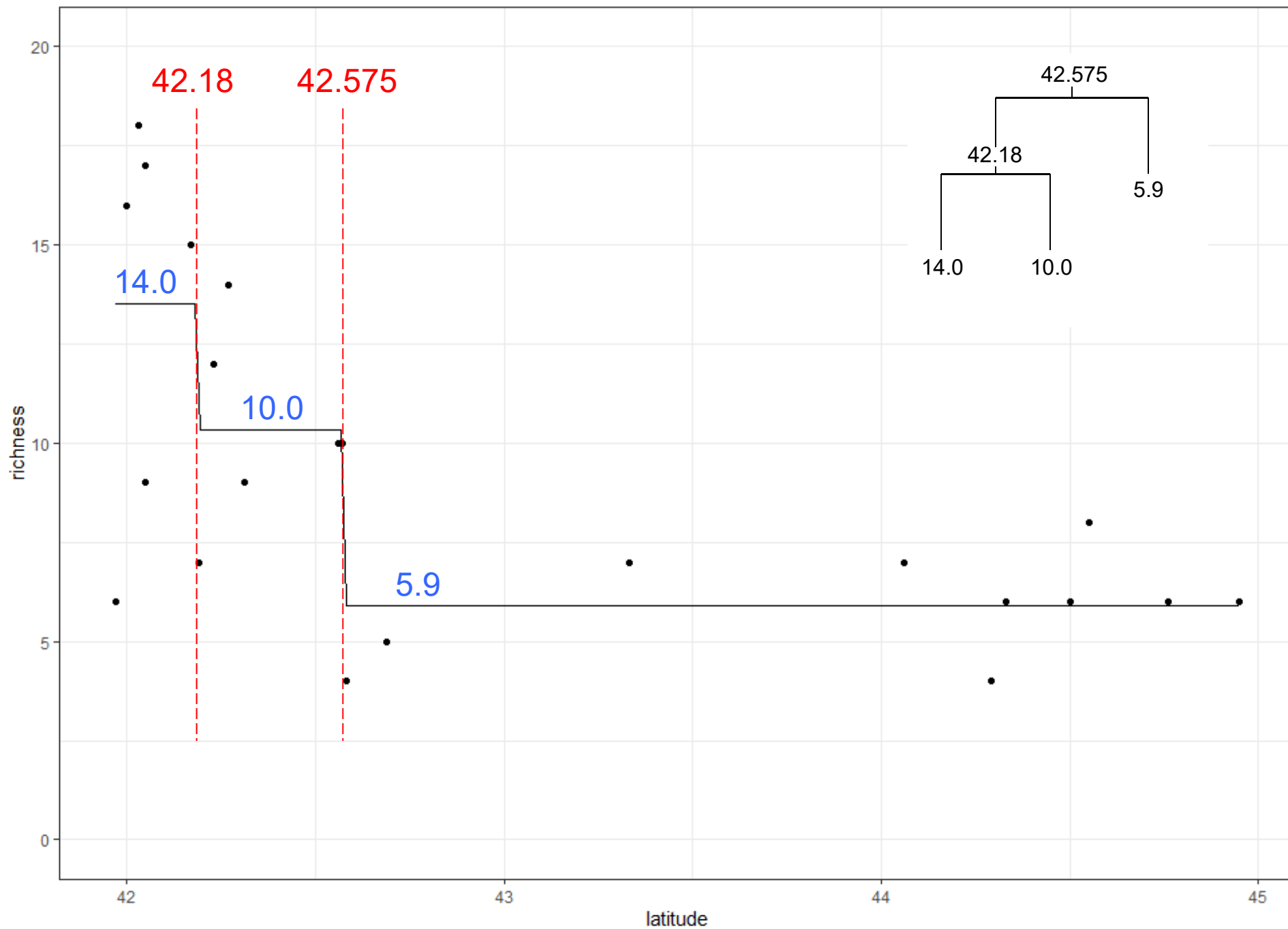
# Basic ML algorithms

- Ensemble methods
  - Bagging
  - Random forest
  - Boosting
- Neural networks

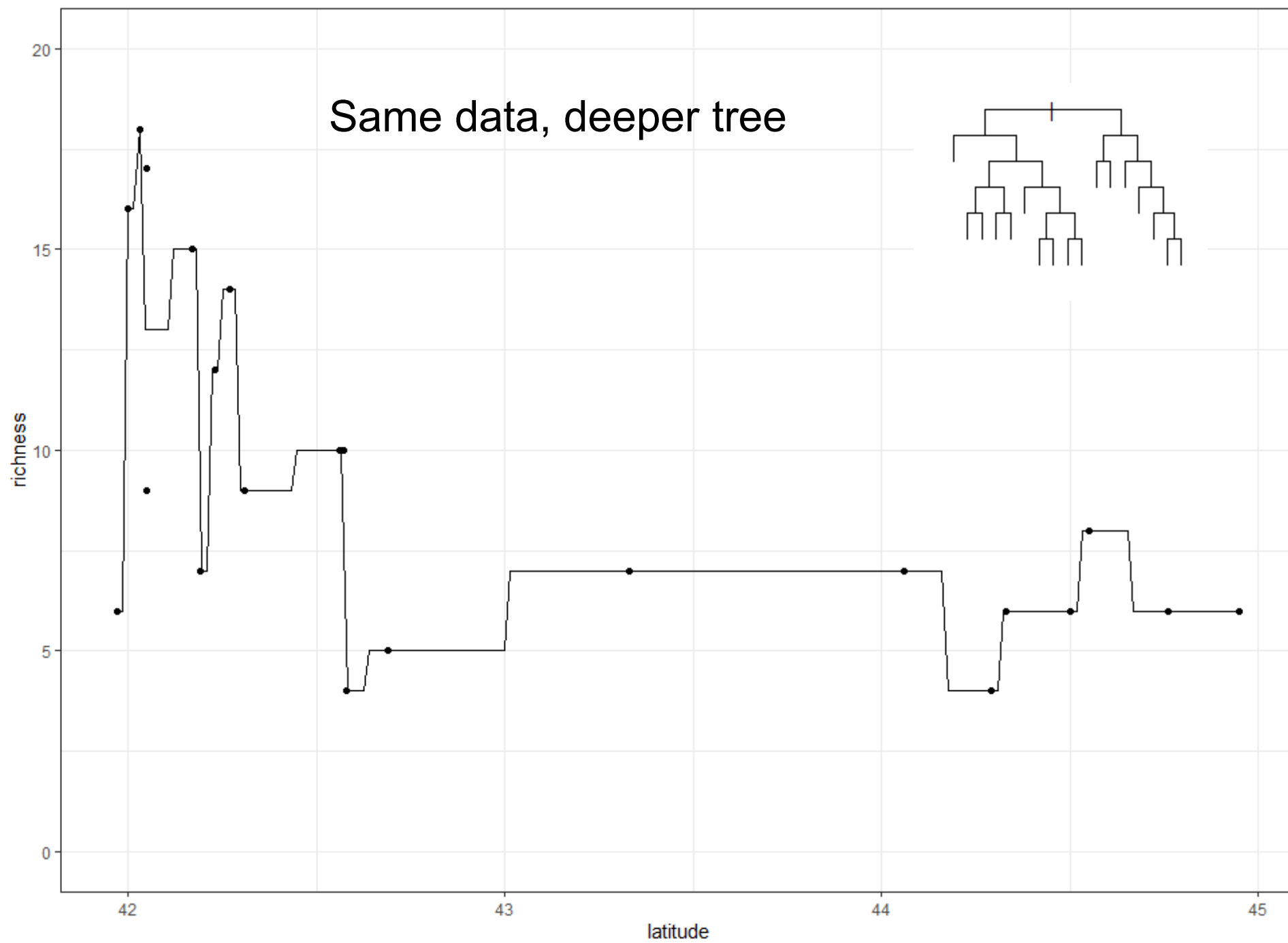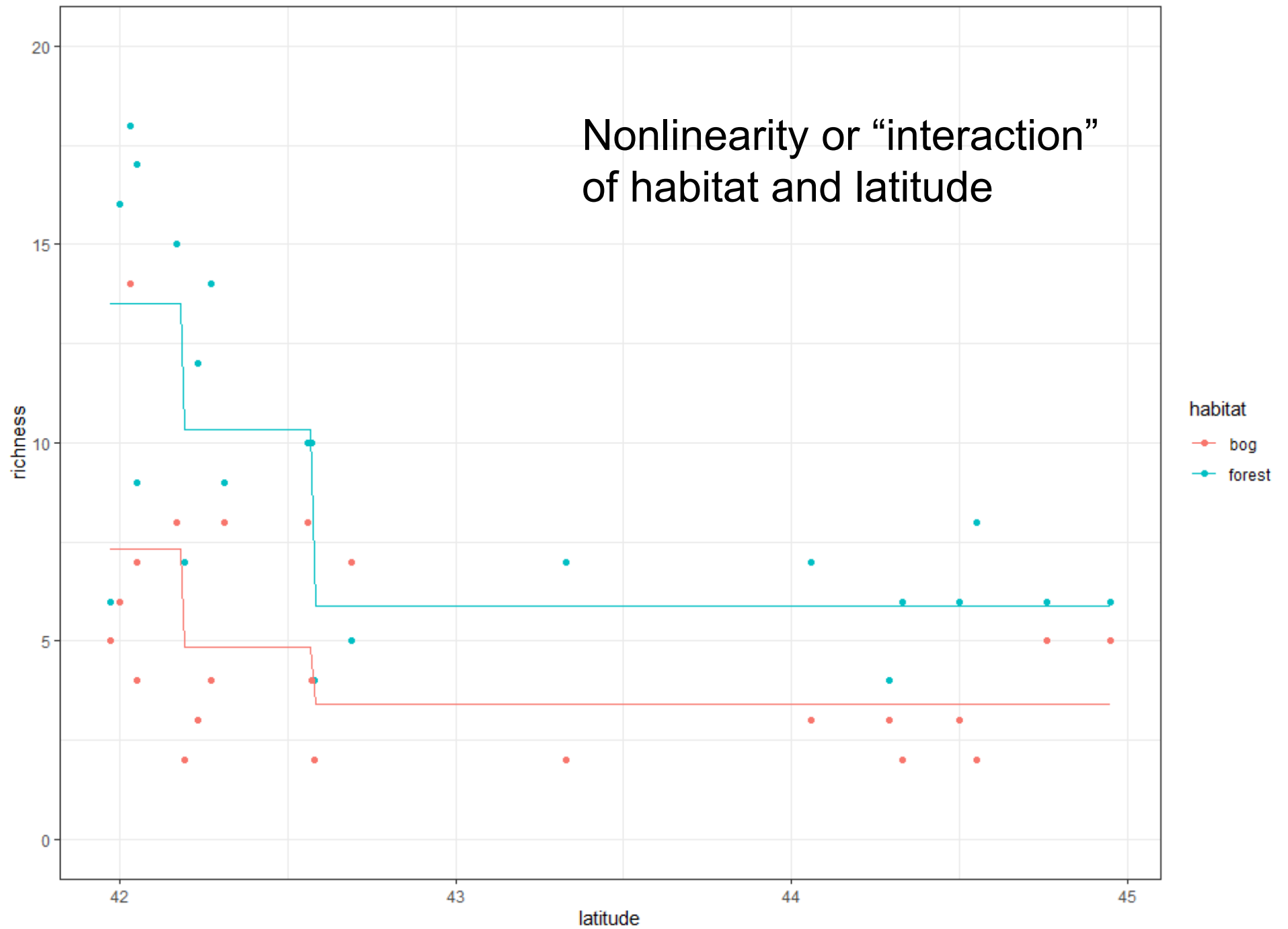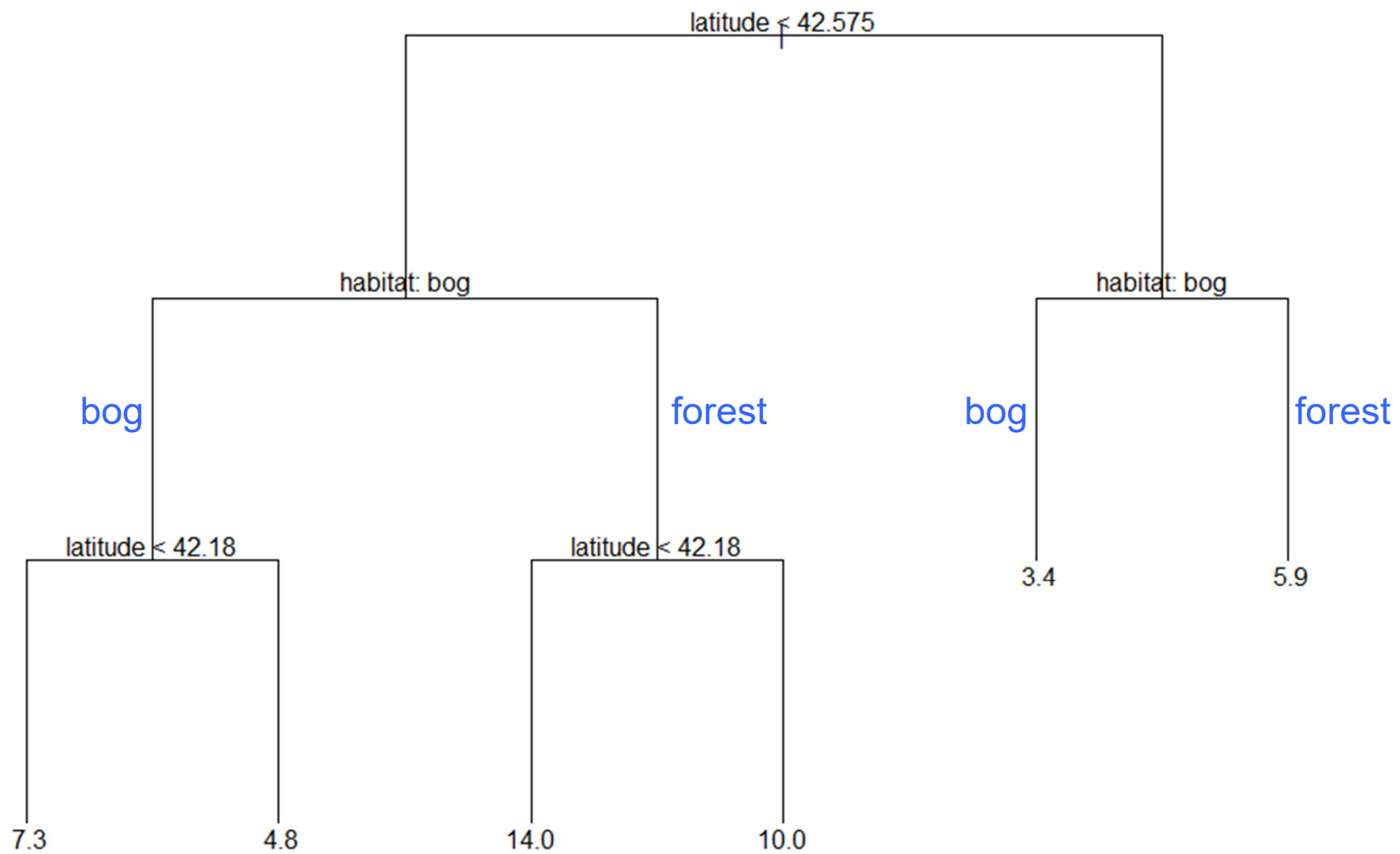# Regression tree base model

Same data, deeper tree

Nonlinearity or "interaction" of habitat and latitude

# Ensemble methods

- Train many models
- Average the models to predict
- Averaging reduces variance

$$\text{e.g.} \qquad \text{Var}(\bar{y}) = \frac{\sigma_y^2}{n}$$

# Bagging

- Bootstrap
  - form new datasets by resampling from the data

- Aggregate
  - average over bootstrapped model fits

# Bagging algorithm

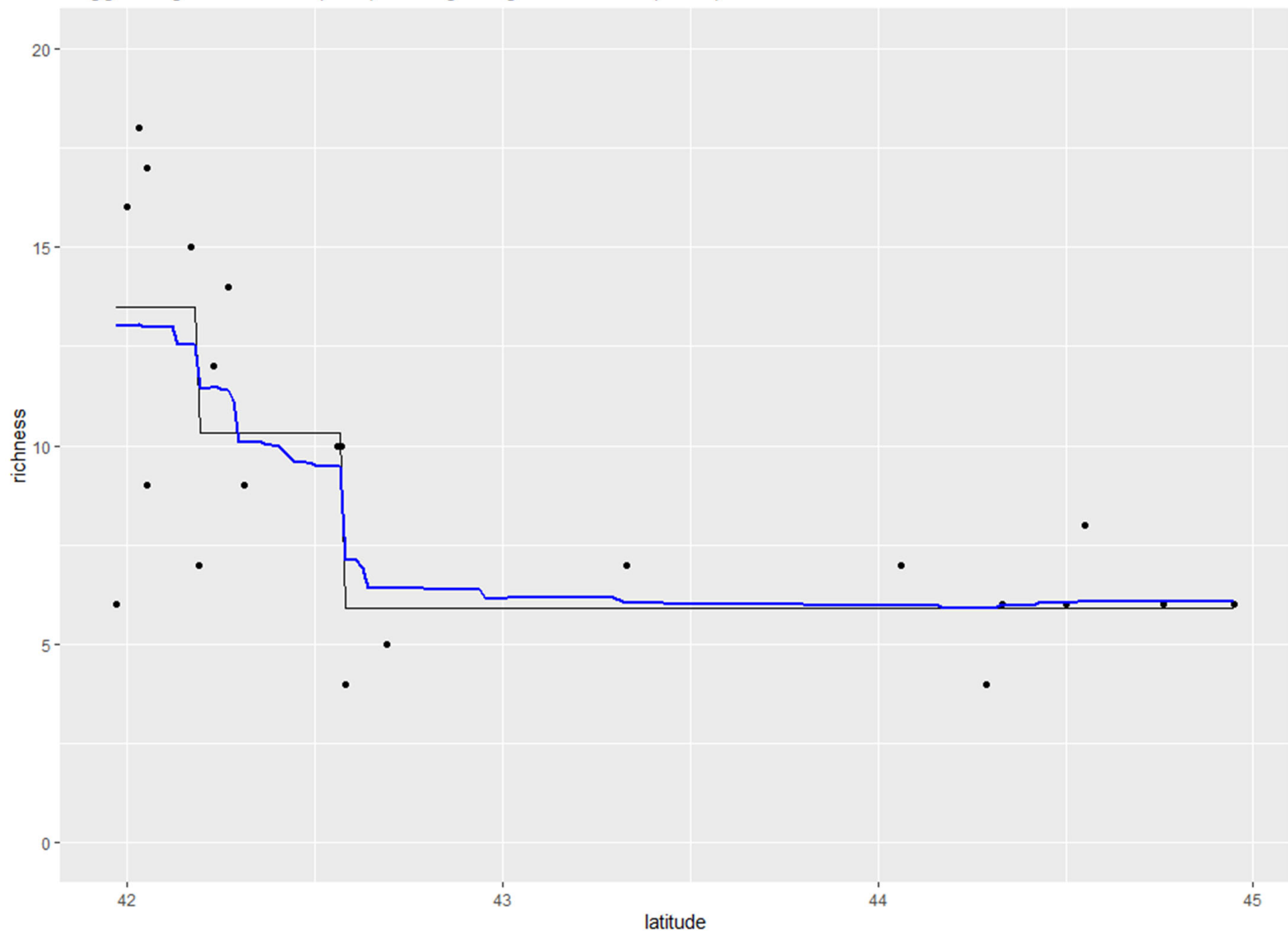for many repetitions
    resample the data with replacement
    train the base model
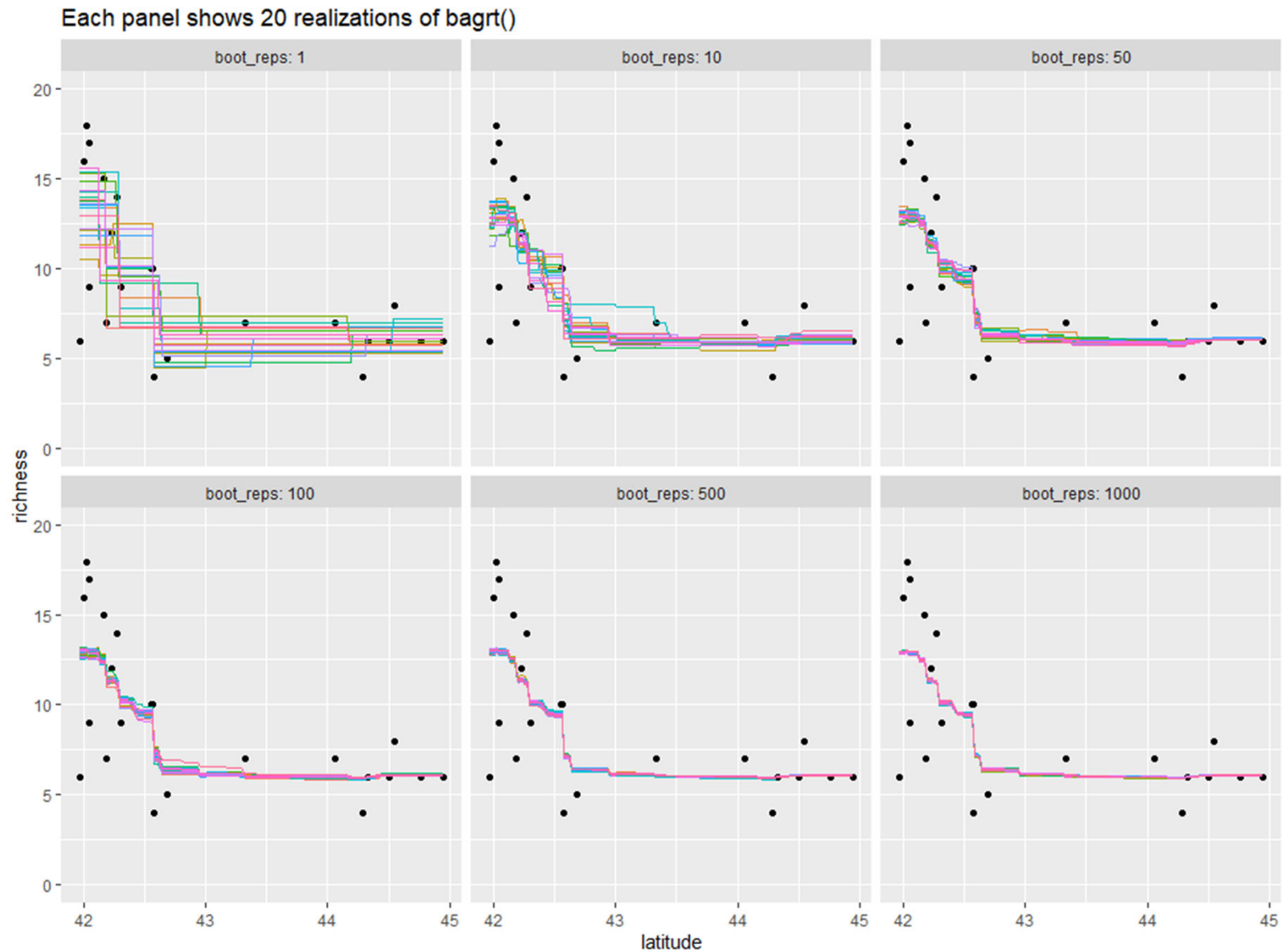    record prediction
final prediction = mean of predictions

Base model: can be any type of model

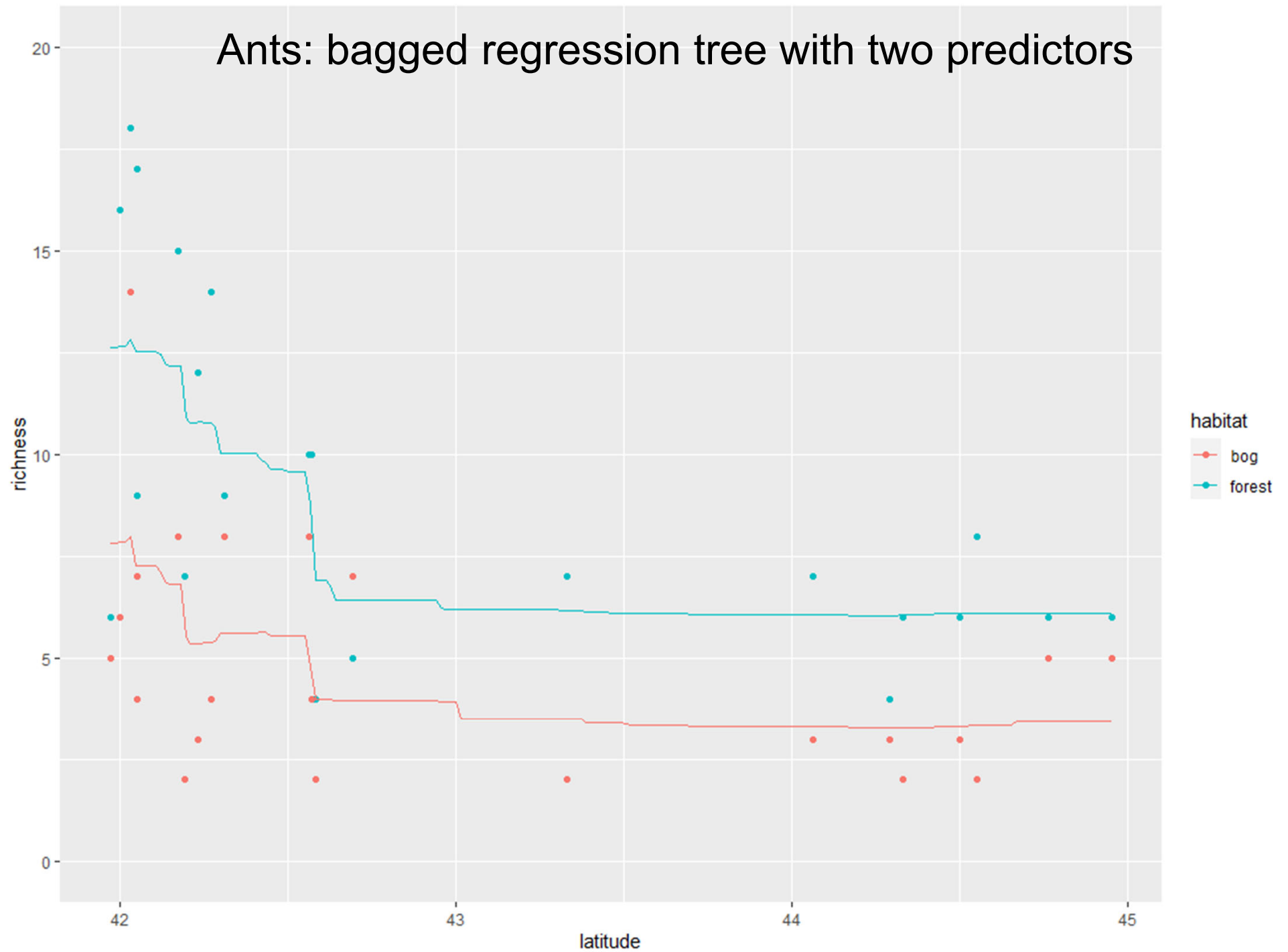Bagged regression tree (blue) vs single regression tree (black)

# Bagging reduces prediction variance



Each panel shows 20 realizations of bagrt()

Ants: bagged regression tree with two predictors

# Random forest

Algorithm

for many repetitions

    randomly select m predictor variables

    resample the data (rows) with replacement
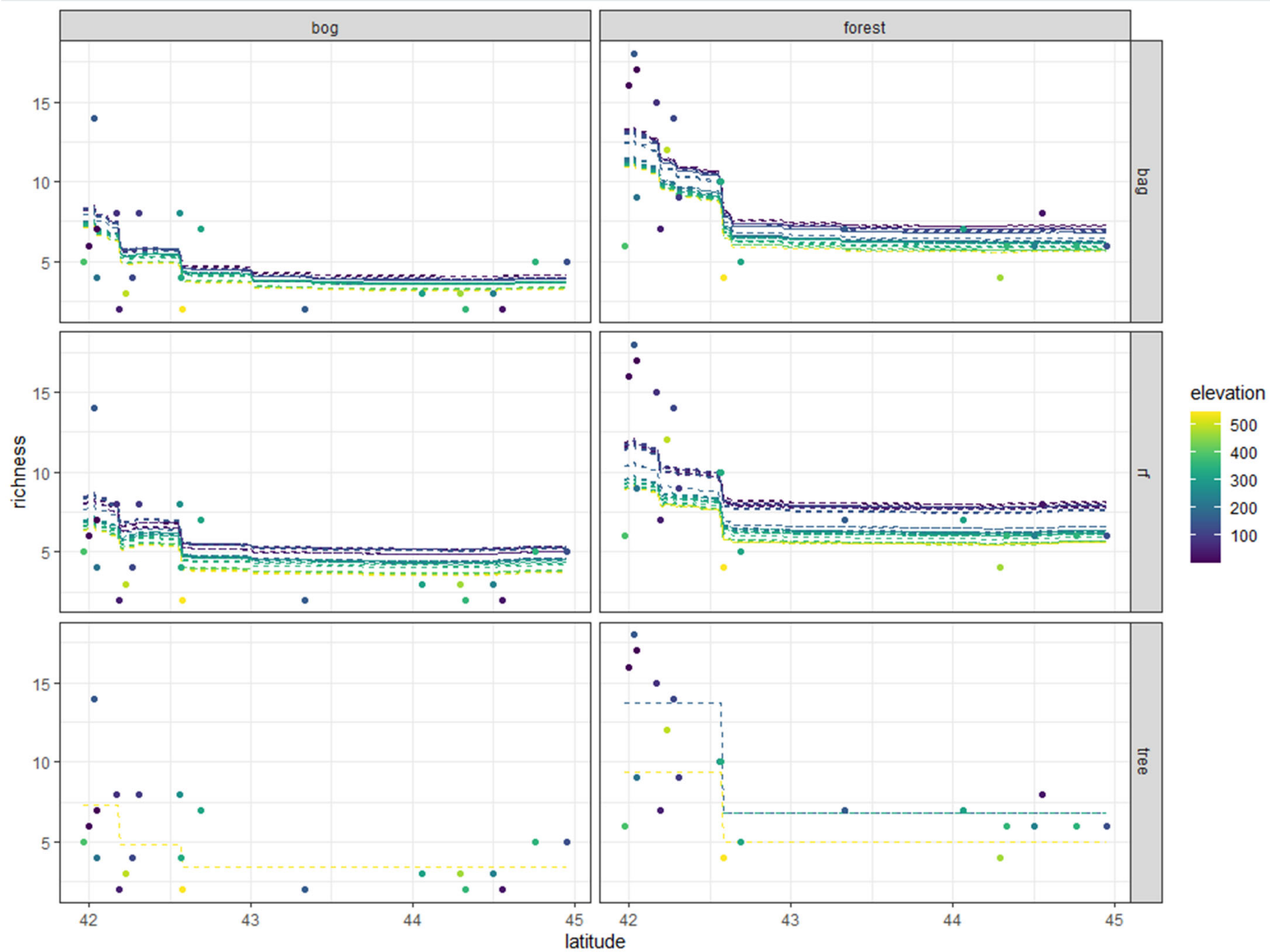
    train the tree model

    record prediction

final prediction = mean of predictions

# R packages

- randomForest
  - original Breimen (2001) algorithm
  - Fortran

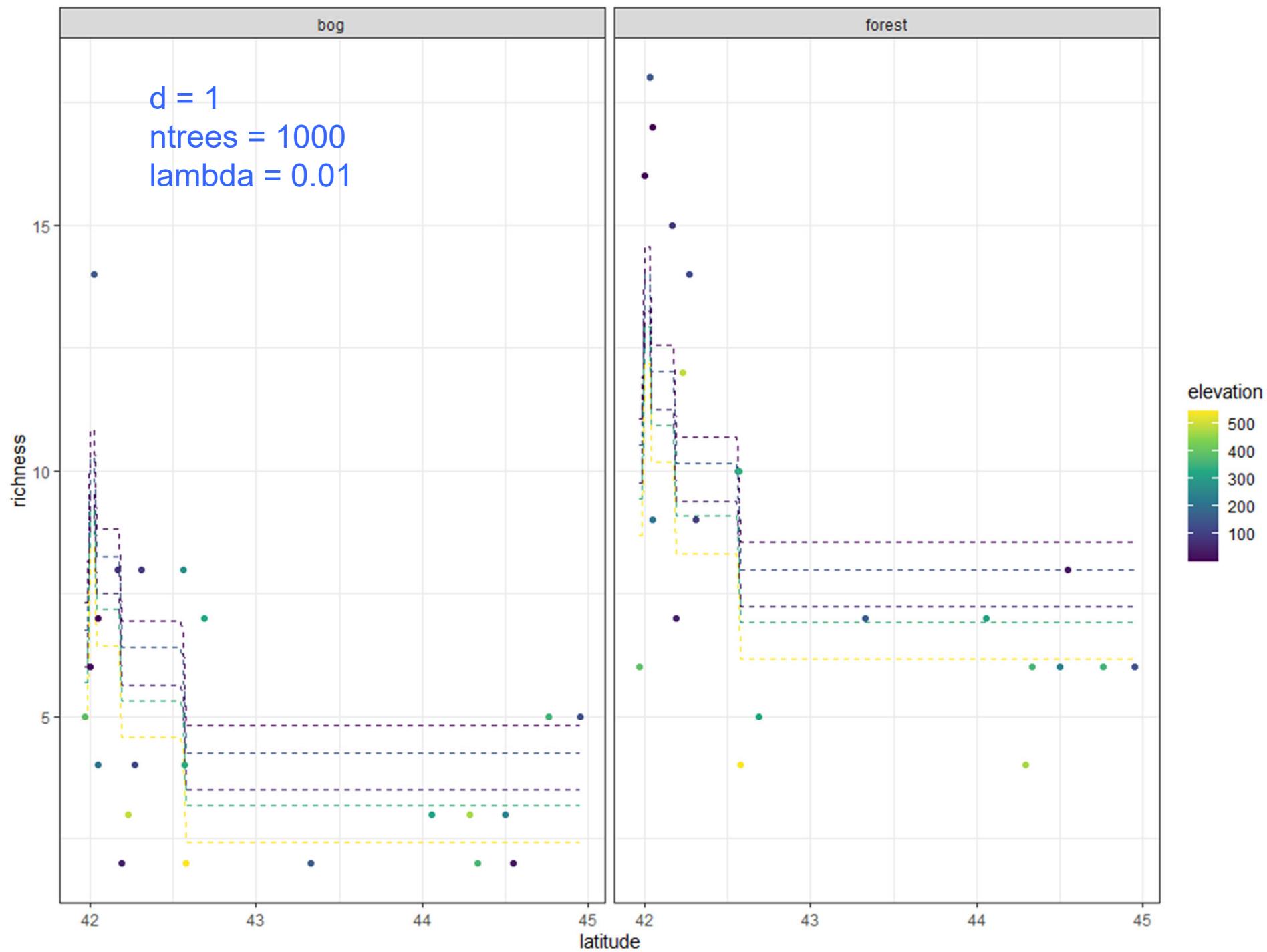- ranger
  - fast implementation
  - C++

# Boosting algorithm

- Too complex for this quick overview
- Basic idea: learn slowly by building up an ensemble iteratively from many models, each with small weight
- Key tuning parameter: learning rate
- Can include aspects of bagging and RF
- Training algorithm: gradient descent

# Boosting packages in R

- **gbm:** gradient boosting machines
  - boosted decision trees
  - good, stable, maintained
  - retired (no new features)
- **xgboost:** extreme gradient boosting
  - R interface to very fast C++ library
  - additional algorithm innovations
  - current industry standard
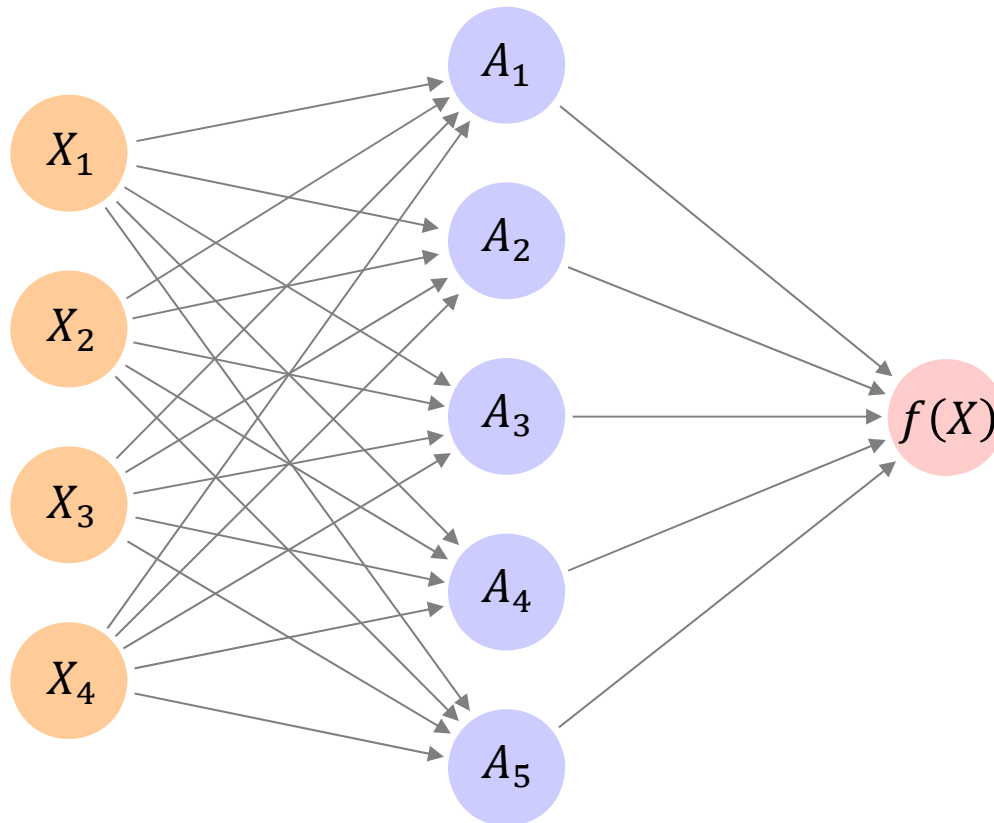
# Neural Networks in R

- keras
  - currently easiest to use
  - tensorflow (Google)
- torch
  - coming along (v 0.9)
  - PyTorch (Facebook, now Linux Foundation)

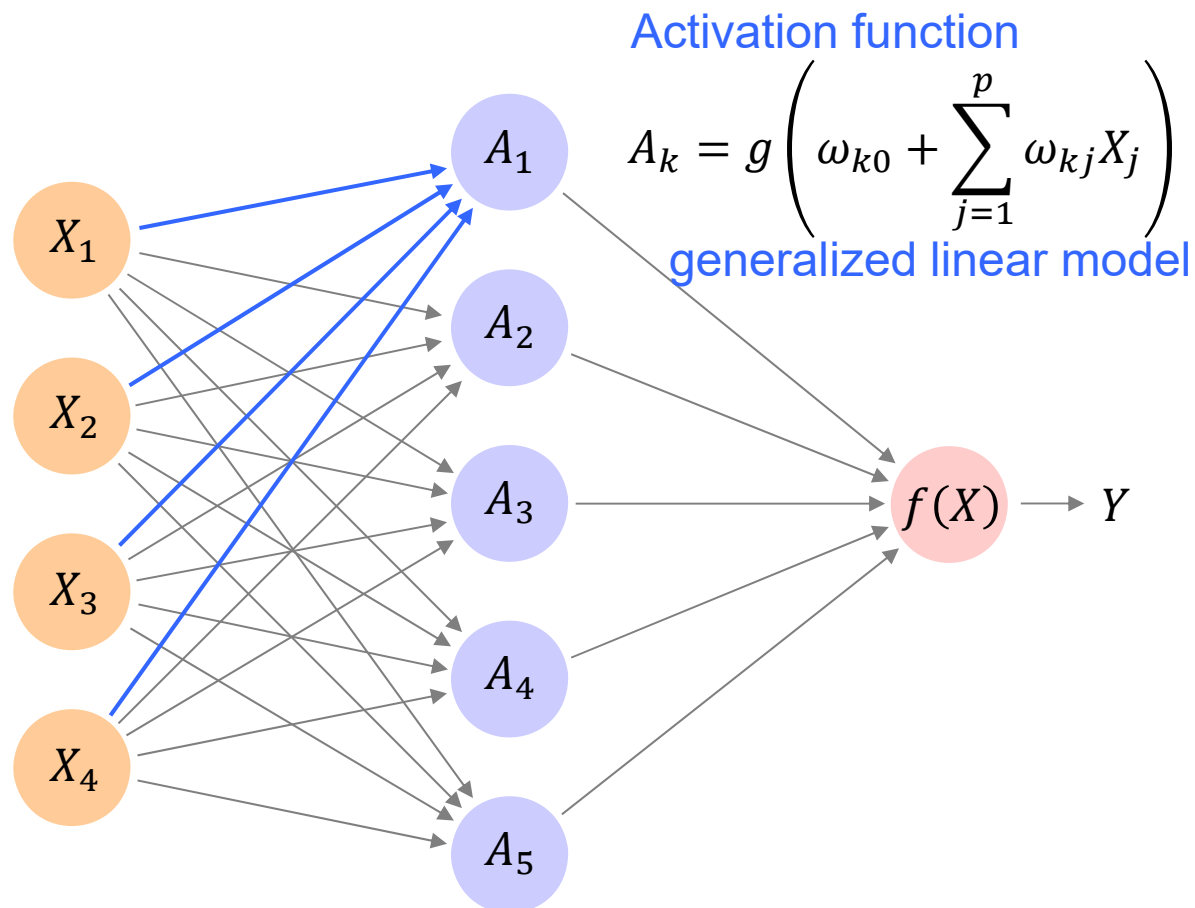# Single layer NN



Input layer     Hidden layer     Output layer

$X_1$, $X_2$, $X_3$, $X_4$

$A_1$, $A_2$, $A_3$, $A_4$, $A_5$

$f(X)$

# Single layer NN

Input layer

Hidden layer

Output layer

Activation function

$$A_k = g\left(\omega_{k0} + \sum_{j=1}^{p} \omega_{kj} X_j\right)$$

generalized linear model

Transform from X to A

# Single layer NN

Input
layer

Hidden
layer

Output
layer

Activation function

$$A_k = g\left(\omega_{k0} + \sum_{j=1}^{p} \omega_{kj}X_j\right)$$

generalized linear model

Transform
from X to A

$X_1$

$X_2$

$X_3$

$X_4$

$A_1$

$A_2$

$A_3$

$A_4$

$A_5$

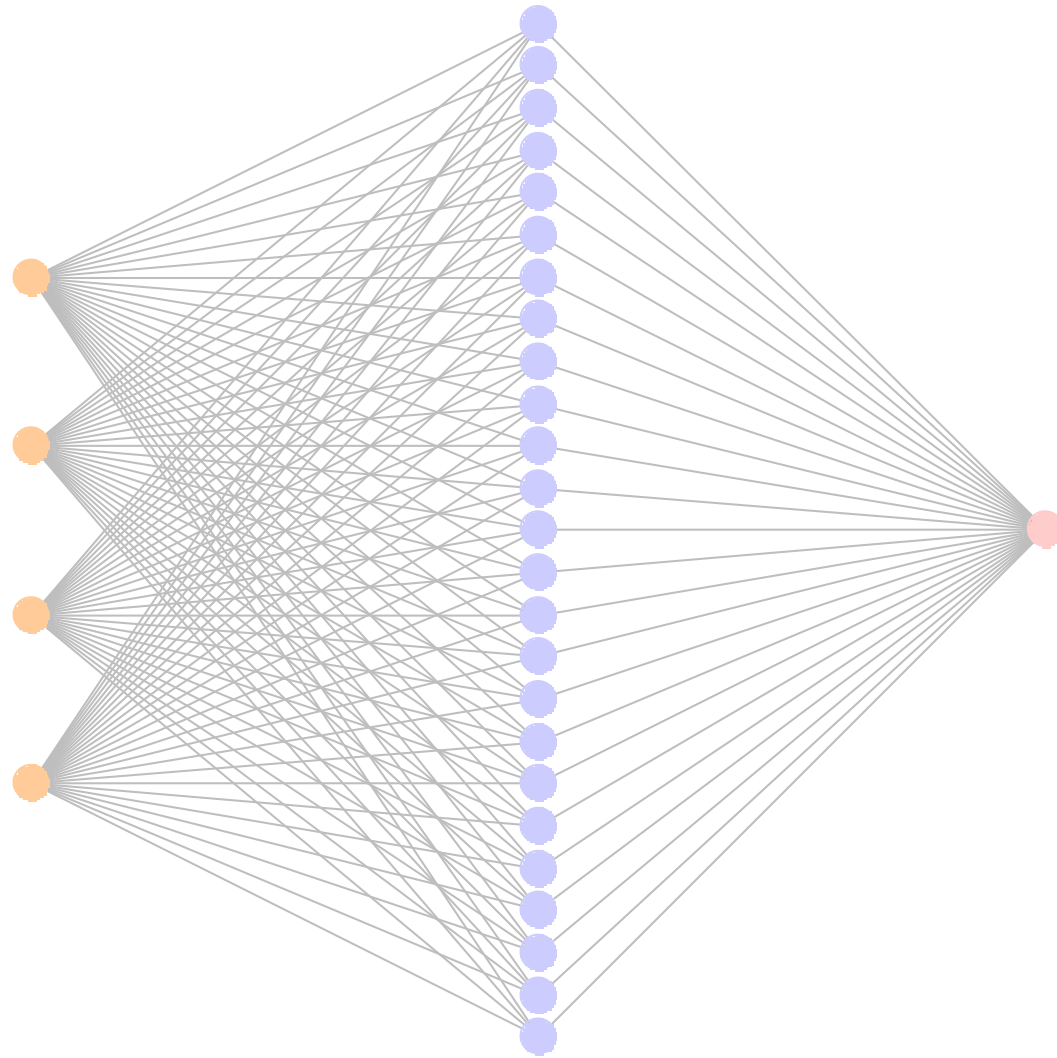$f(X)$ → $Y$

$$f(X) = \beta_0 + \sum_{k=1}^{K} \beta_k A_k$$

multiple linear
regression

# Training algorithm

- Stochastic gradient descent
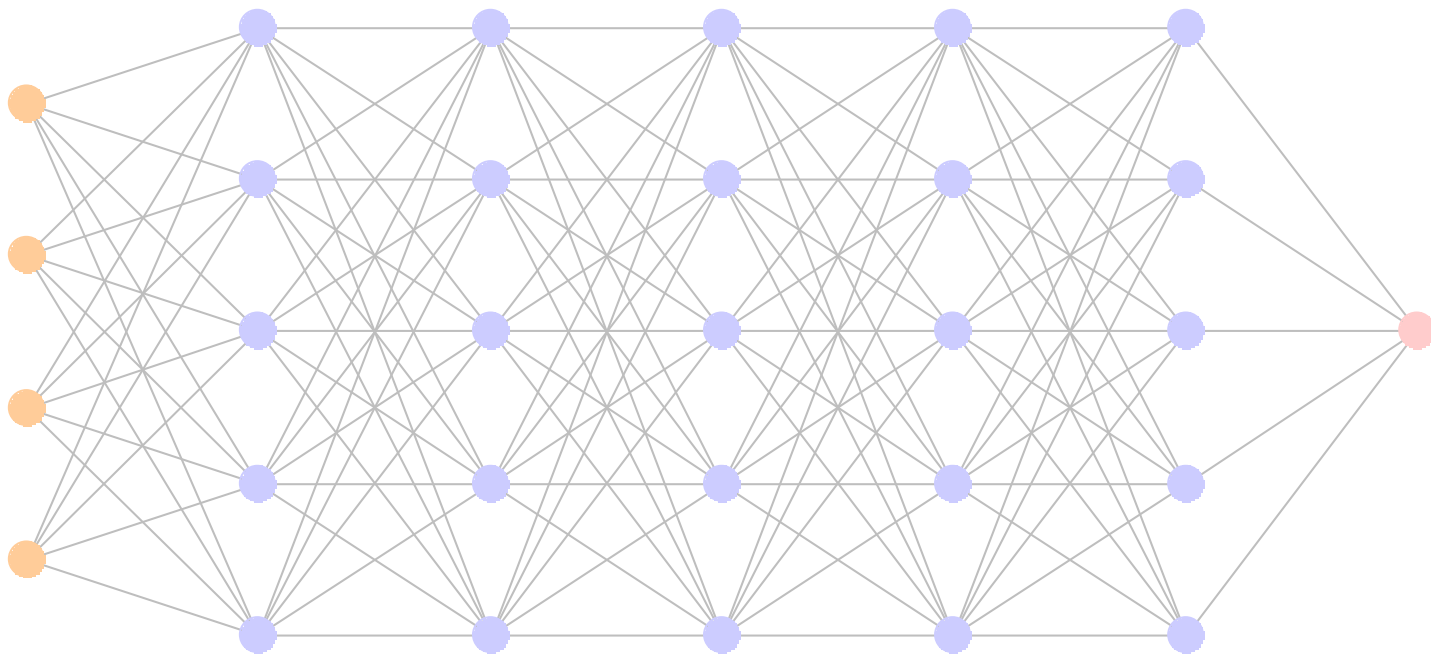  - with back propagation

# Deep learning

- Multilayer neural networks
  - aka deep feedforward networks
  - aka multilayer perceptrons (MLP)
- Model algorithm
  - expressiveness
    - ability to approximate complex nonlinearity
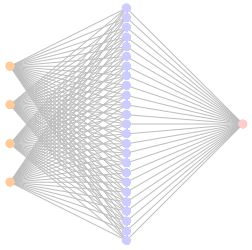  - architecture: width versus depth
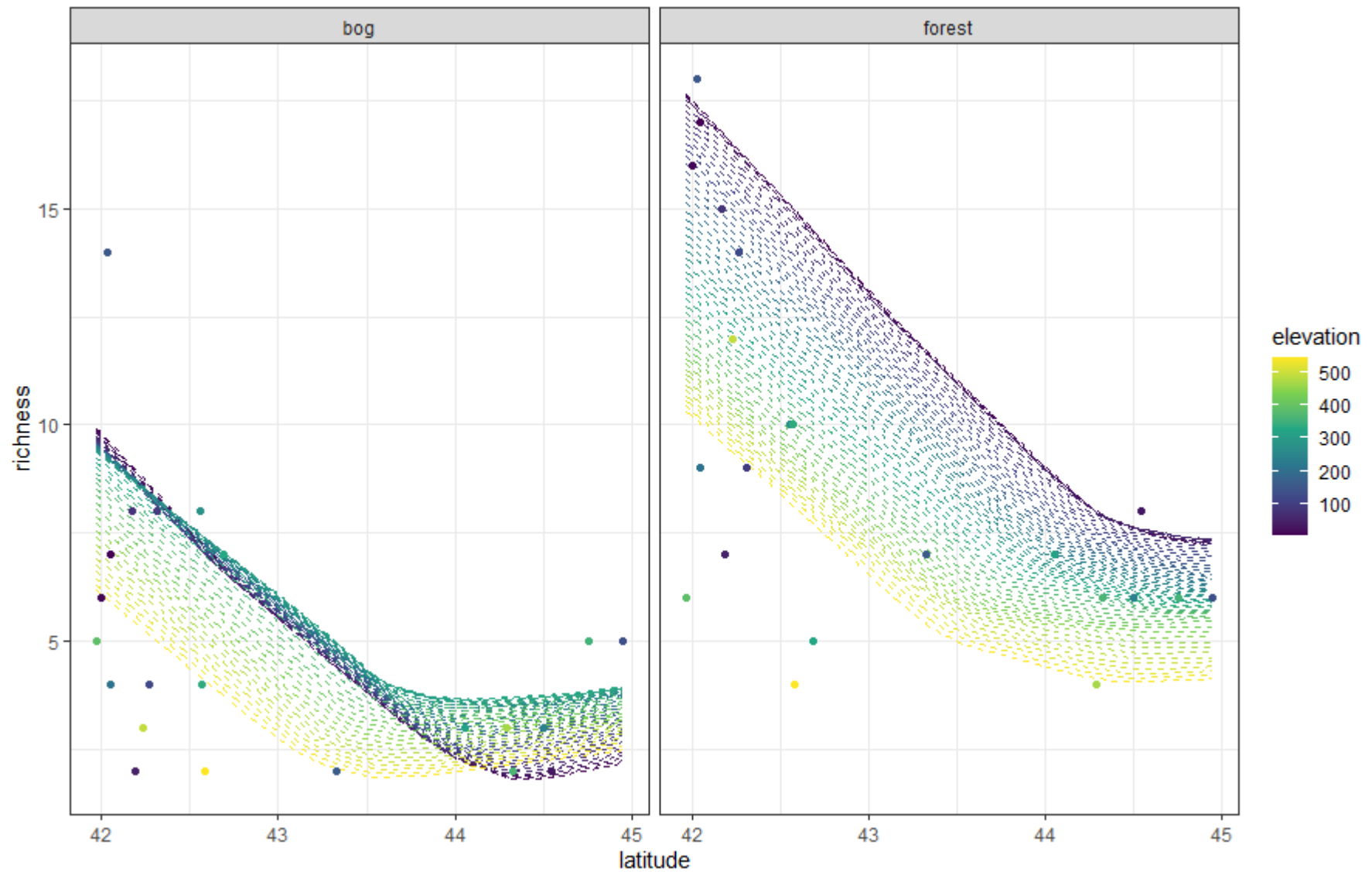
# Wide: 25 hidden units
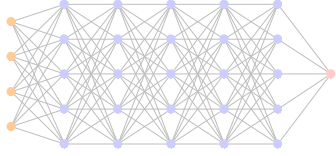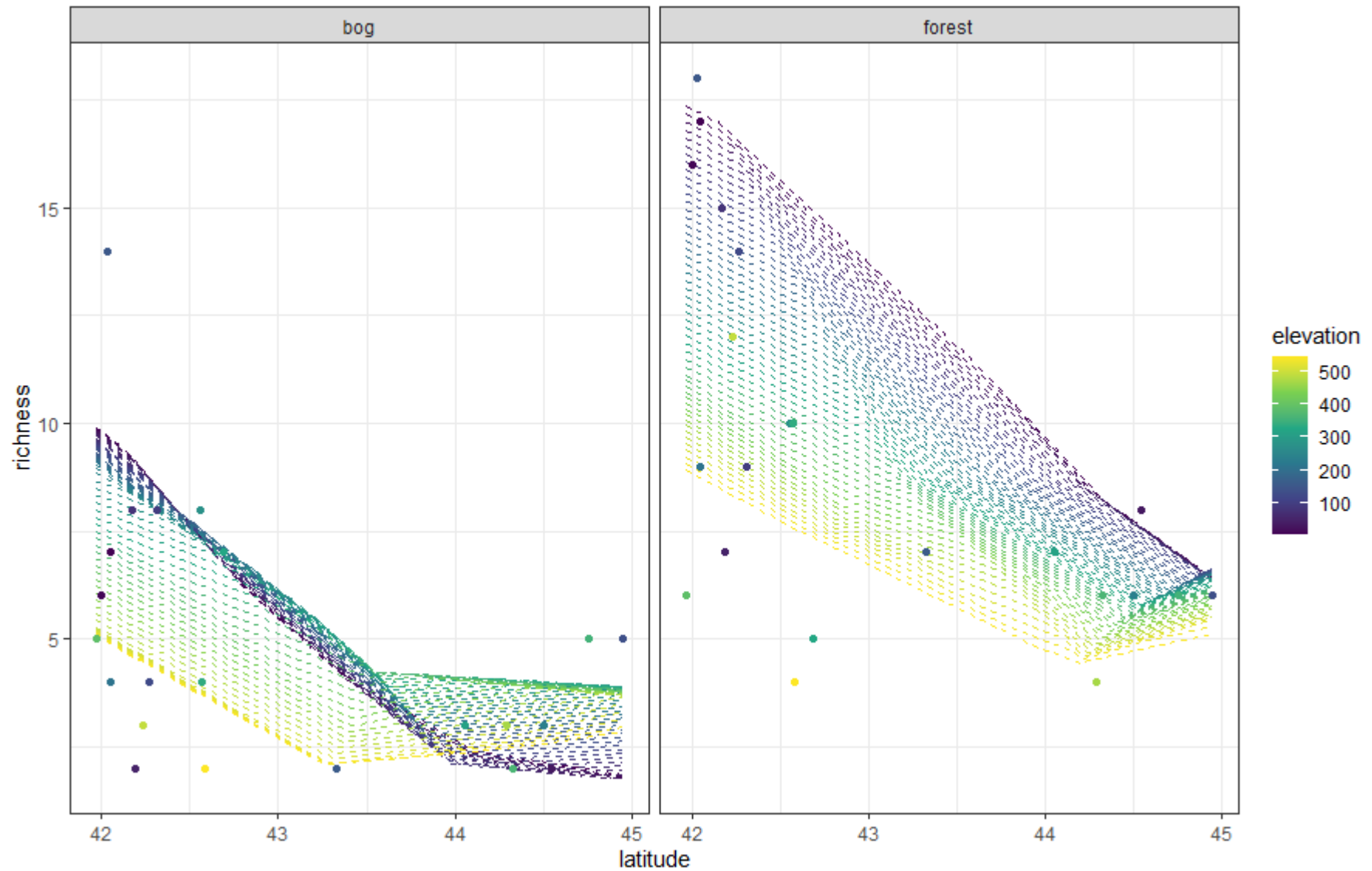
# Deep: 25 hidden units
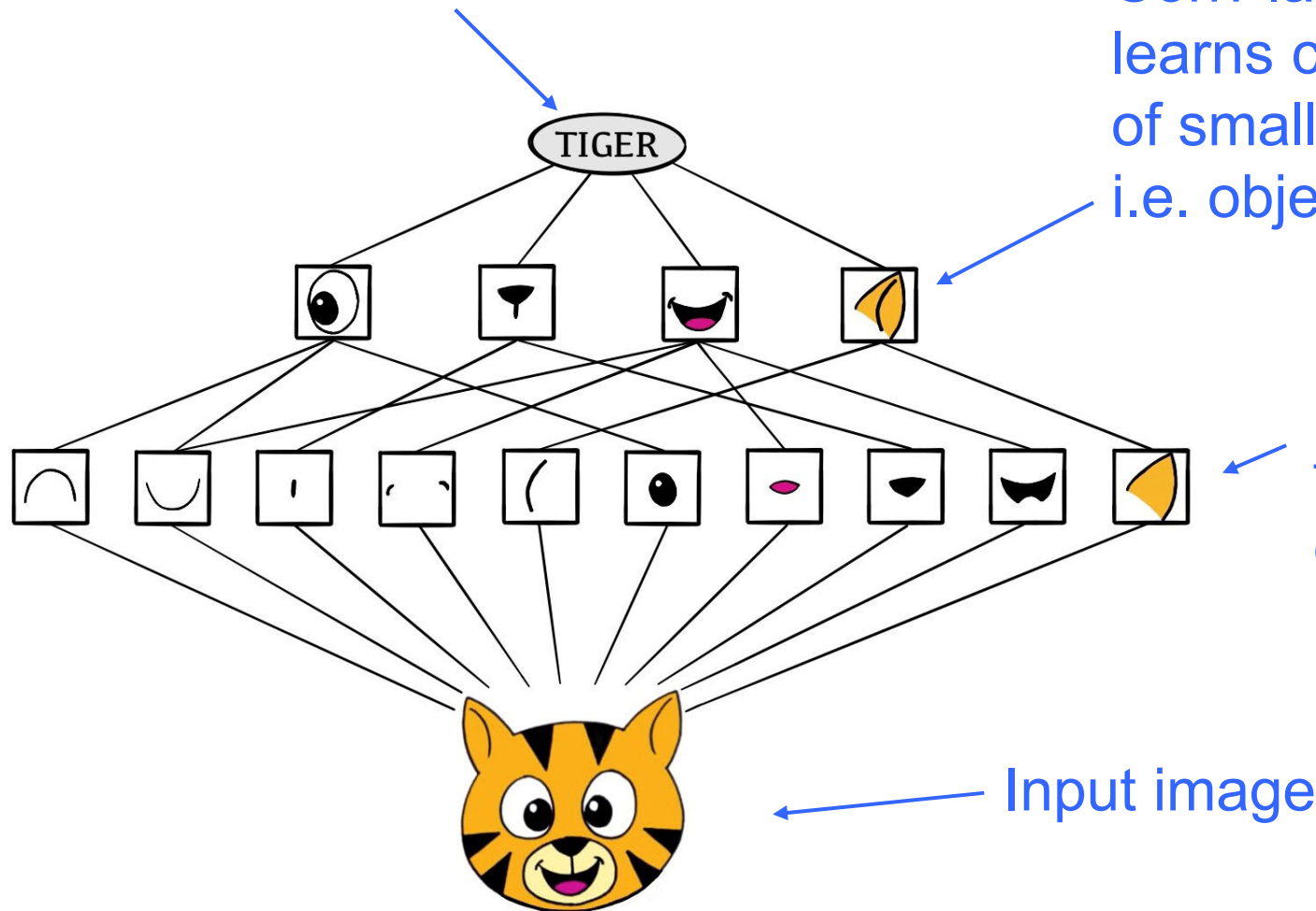
# Ants data: wide

# Ants data: deep

# Convolutional NNs



Output is high-level concept

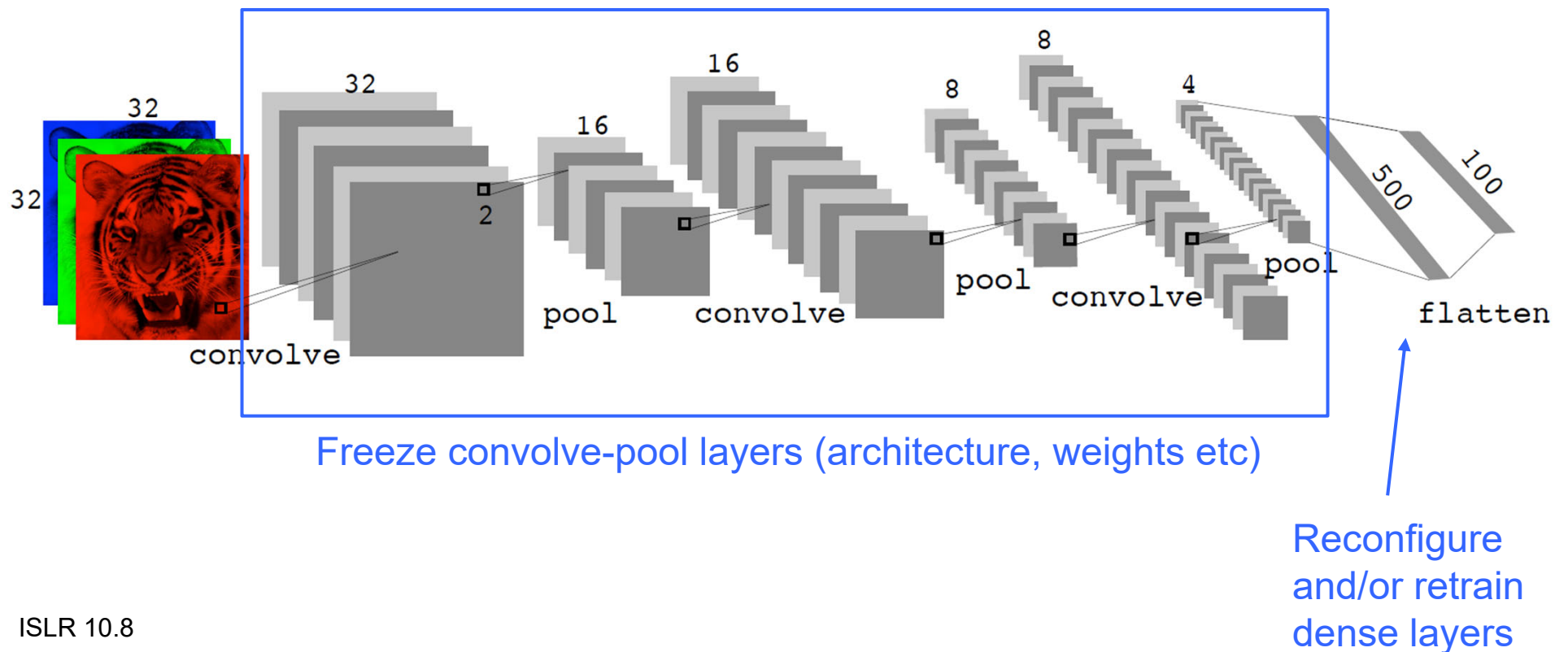Conv layer 2
learns combinations
of small features
i.e. objects

Conv layer 1
learns small
features
e.g. edges

Input image

ISLR 10.6

# Transfer learning

Pretrained model on related big data
Retrain last 1-2 layers on specialized little data



Freeze convolve-pool layers (architecture, weights etc)

Reconfigure and/or retrain dense layers

ISLR 10.8

# Rapid innovation

- Architectures
- Algorithms
- Recent example: transformer

# Outlook

Automated data collection
+
machine learning
=
revolutionary