

# This week

- Ensemble methods
  - Bagging
  - Random forests

# Today

- Bagging
  - code: `ants_bag.R`, `ants_bag.py`
  - inference algorithm (predictive performance)
  - tuning
- Parallel processing

# Ensemble methods

- Train many models (ensemble)
- Average the models to predict
- Averaging reduces prediction variance

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

Takeaway:

variance of the mean of y is  
less than the variance of y

# Bagging

- Bootstrap
  - form many new datasets by resampling from the data
  - sample with replacement
  - train model on each dataset
- Aggregate
  - average over trained models

# 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

## Algorithm pseudocode to R

```
setup {  
  # Bagging algorithm  
  boot_reps <- 500  
  n <- nrow(forest_ants)  
  nx <- nrow(grid_data)  
  boot_preds <- matrix(rep(NA, nx*boot_reps), nrow=nx, ncol=boot_reps)  
  # for many repetitions  
  for ( i in 1:boot_reps ) {  
    # resample the data (rows) with replacement  
    boot_indices <- sample(1:n, n, replace=TRUE)  
    boot_data <- forest_ants[boot_indices,]  
    # train the base model  
    boot_train <- tree(richness ~ latitude, data=boot_data)  
    # record prediction  
    boot_preds[,i] <- predict(boot_train, newdata=grid_data)  
  }  
  bagged_preds <- rowMeans(boot_preds)  
}
```

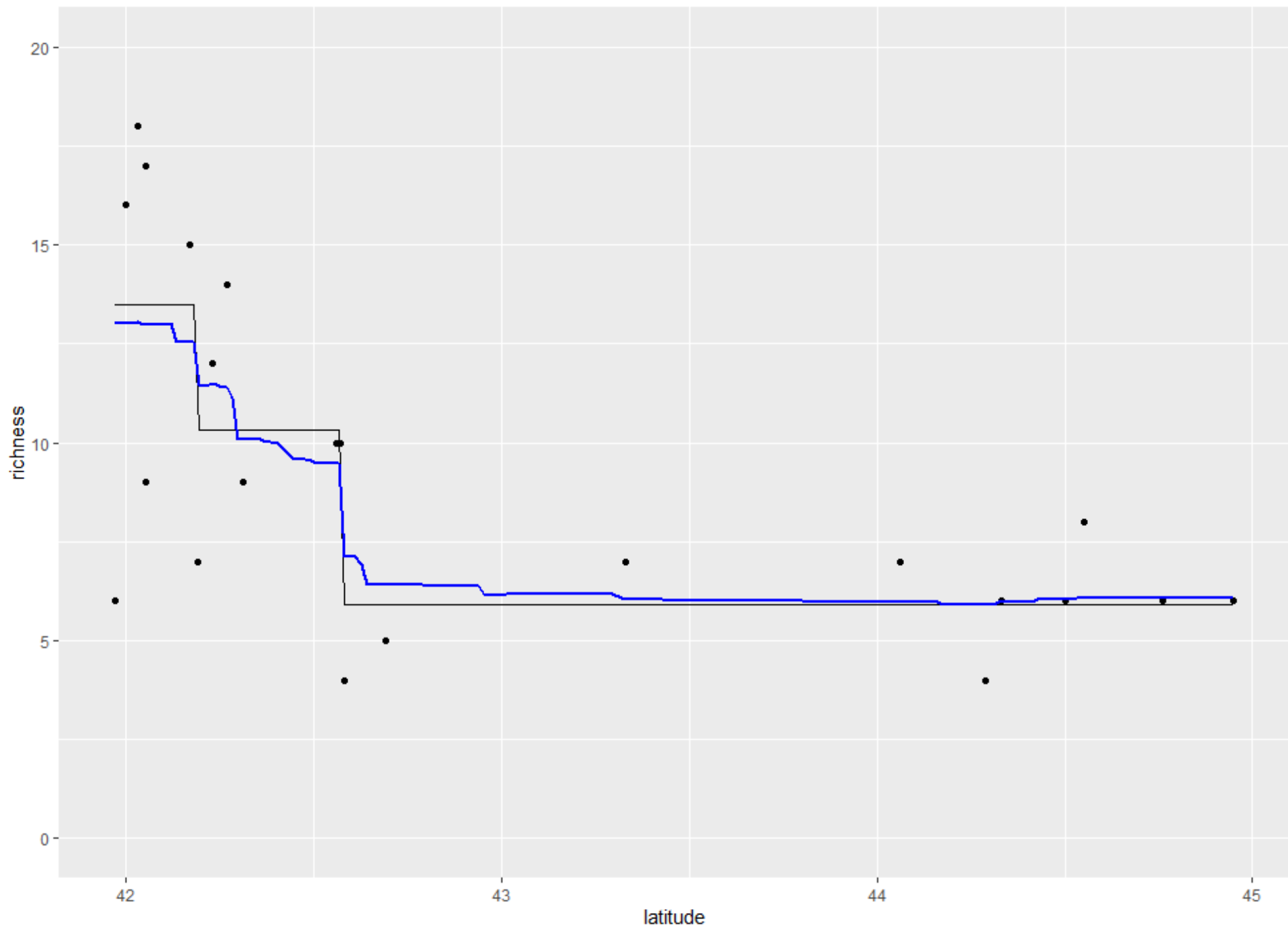
# Bagged regression tree

## Algorithm pseudocode to Python

setup {

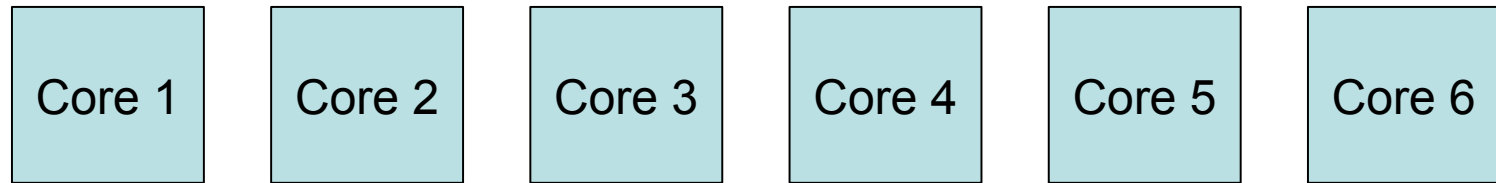
```
# Bagging algorithm
boot_reps = 500
dt = tree.DecisionTreeRegressor(max_depth=2) #define base model
n = len(forest_ants)
nx = len(grid_data)
boot_preds = np.full((nx, boot_reps), np.nan)
# for many repetitions
for i in range(boot_reps):
    # resample the data (rows) with replacement
    boot_indices = rng.choice(range(n), n, replace=True)
    boot_data = forest_ants.iloc[boot_indices]
    # train the base model
    boot_train = dt.fit(boot_data[["latitude"]], boot_data["richness"])
    # record prediction
    boot_preds[:,i] = boot_train.predict(grid_data)
# mean of predictions
bagged_preds = np.mean(boot_preds, axis=1)
```

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





# Parallel processing

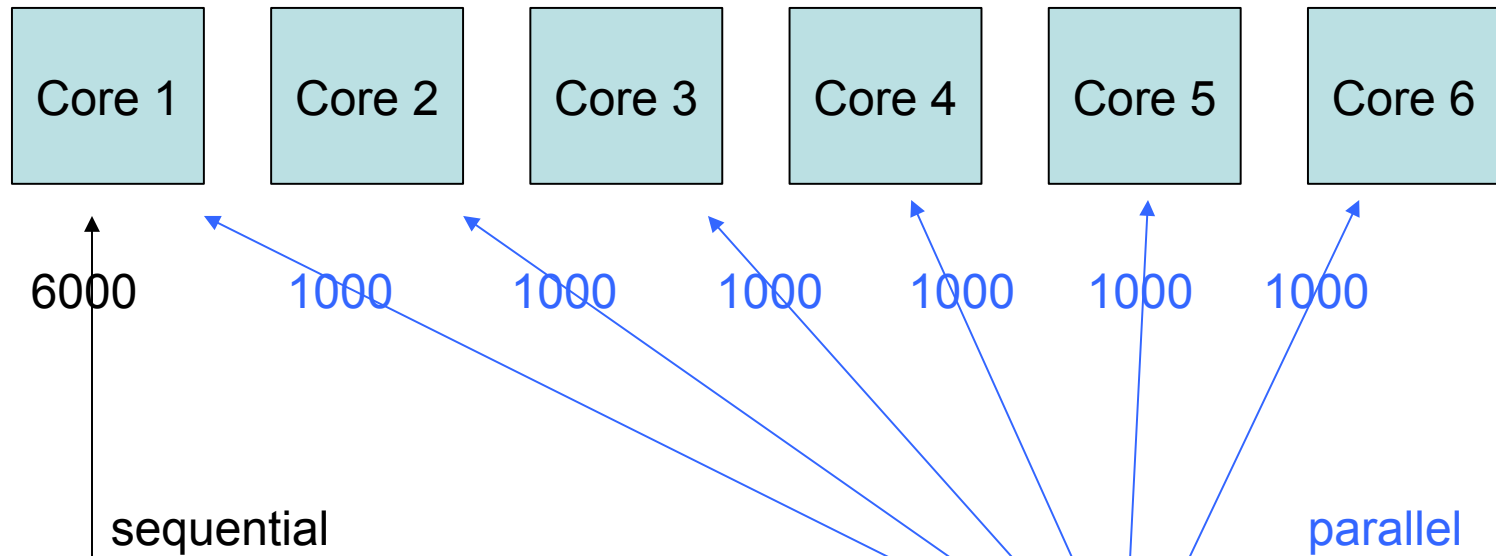


6000

sequential

```
for ( i in 1:6000 ) {  
  y[i] <- long_job()  
}
```

# Parallel processing



```
for ( i in 1:6000 ) {  
  y[i] <- long_job()  
}
```

```
y <- foreach ( i=1:6000 ) %doRNG% {  
  long_job()  
}
```

# Setup

```
library(doFuture)      → library(future)
library(doRNG)          library(foreach)
registerDoFuture()

availableCores()
plan(multisession, workers=8)

# Handy timing function:
system.time( some_function() )

# Saving/loading long jobs
save(myresult1, myresult2, file="/saved/myresult.Rdata")
load("/saved/myresult.Rdata")
```

# Inference algorithm

- 5-fold CV, 500 splits
- Ants: mean prediction error (MSE)

12.93 +/- 0.07

Model	LOOCV	5-fold CV
Polynomial 2	12.88	13.51
Single reg tree	12.68	13.15
KNN 7	12.63	13.03
KNN 6	12.95	13.01
Bagged reg tree	13.23*	12.93
Smoothing spline 3	12.52	12.77

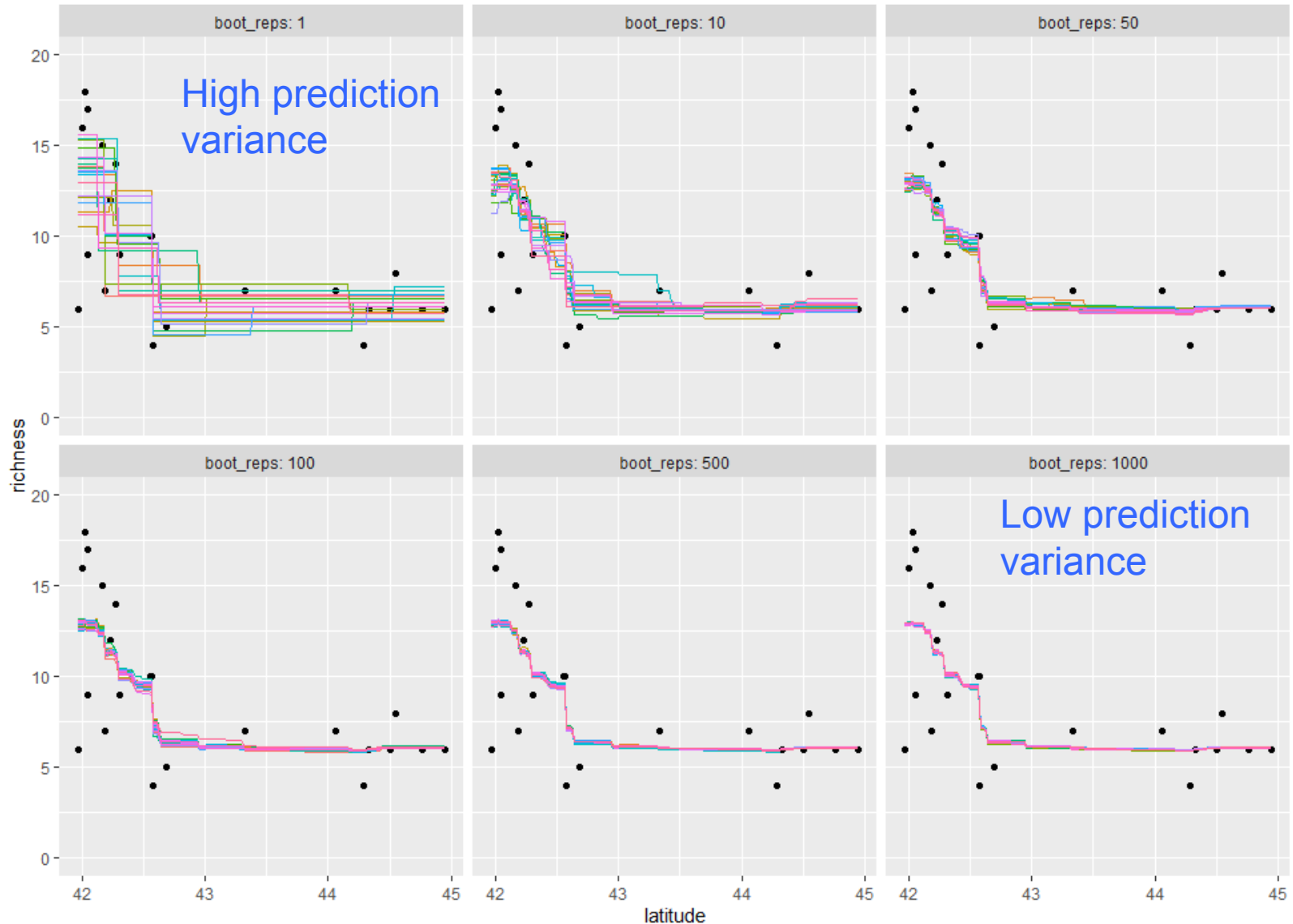
2<sup>nd</sup> best on 5-fold CV:  
generally good predictive function

Worst on LOOCV:  
perhaps an influential data point

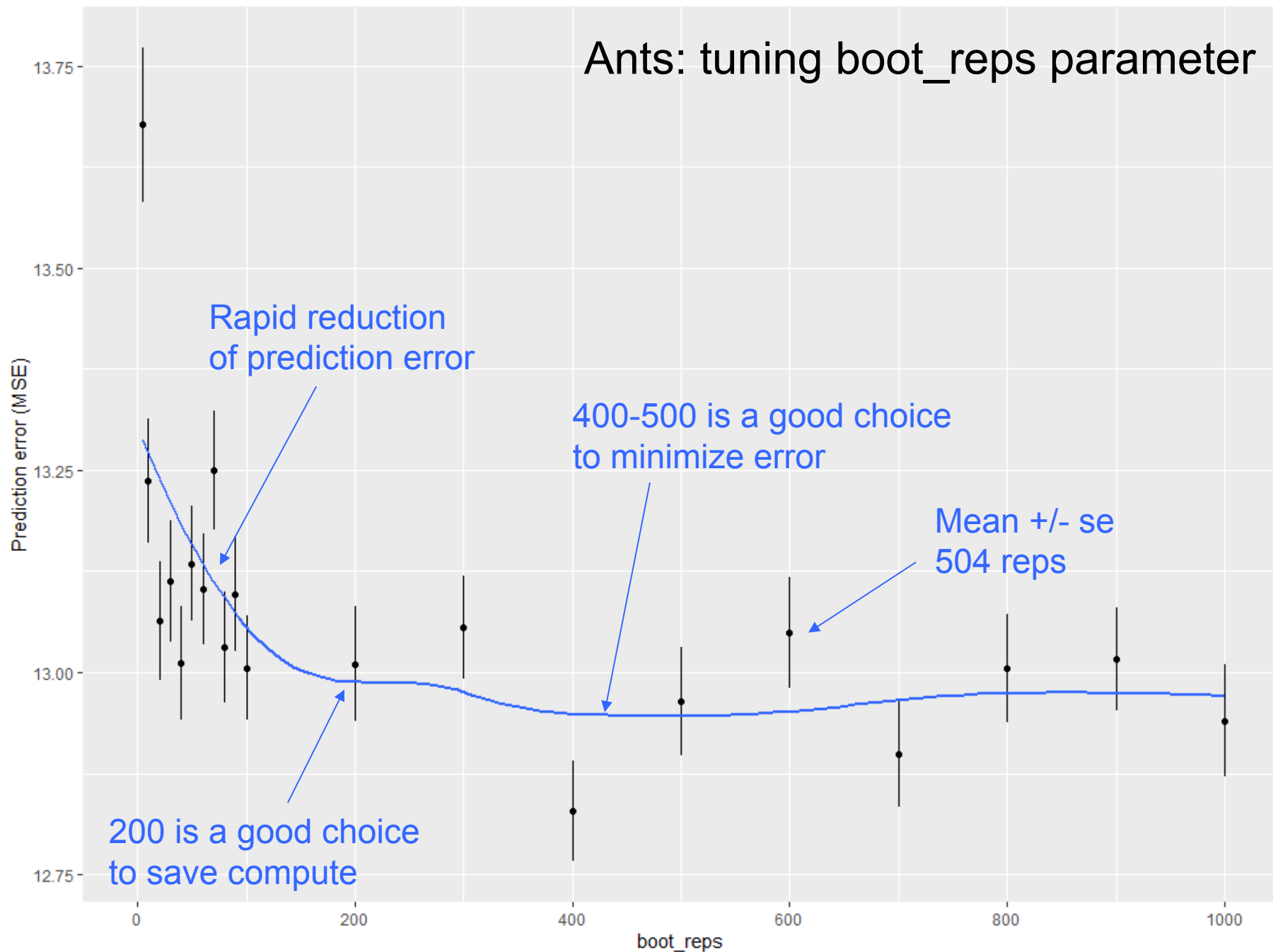
Prediction is always tenuous  
with a small dataset

# Bagging reduces prediction variance

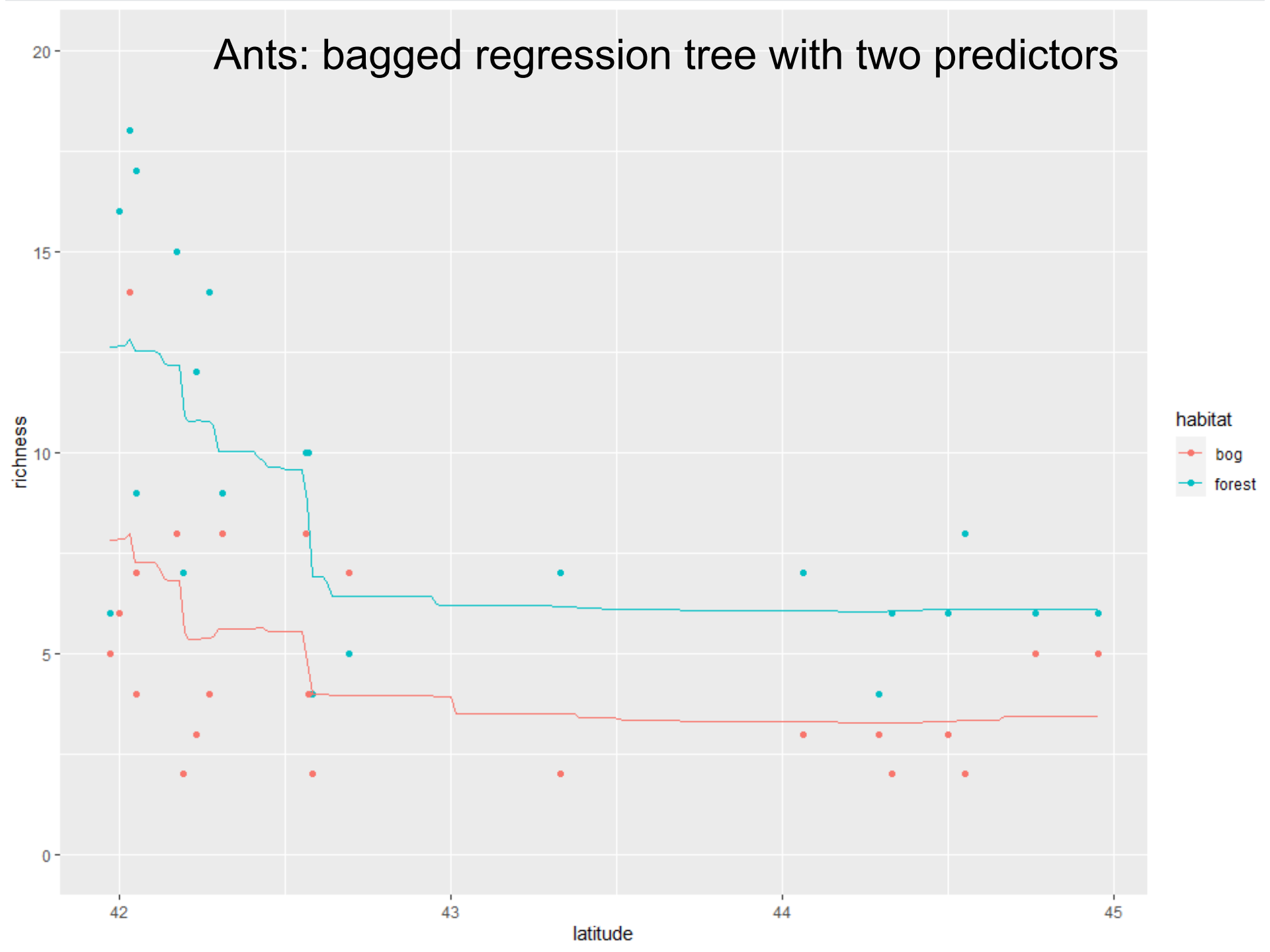
Each panel shows 20 realizations of `bagrt()`



# Ants: tuning boot\_reps parameter



# Ants: bagged regression tree with two predictors



Bagged KNN 7 (blue) compared to single KNN 7 (black)

