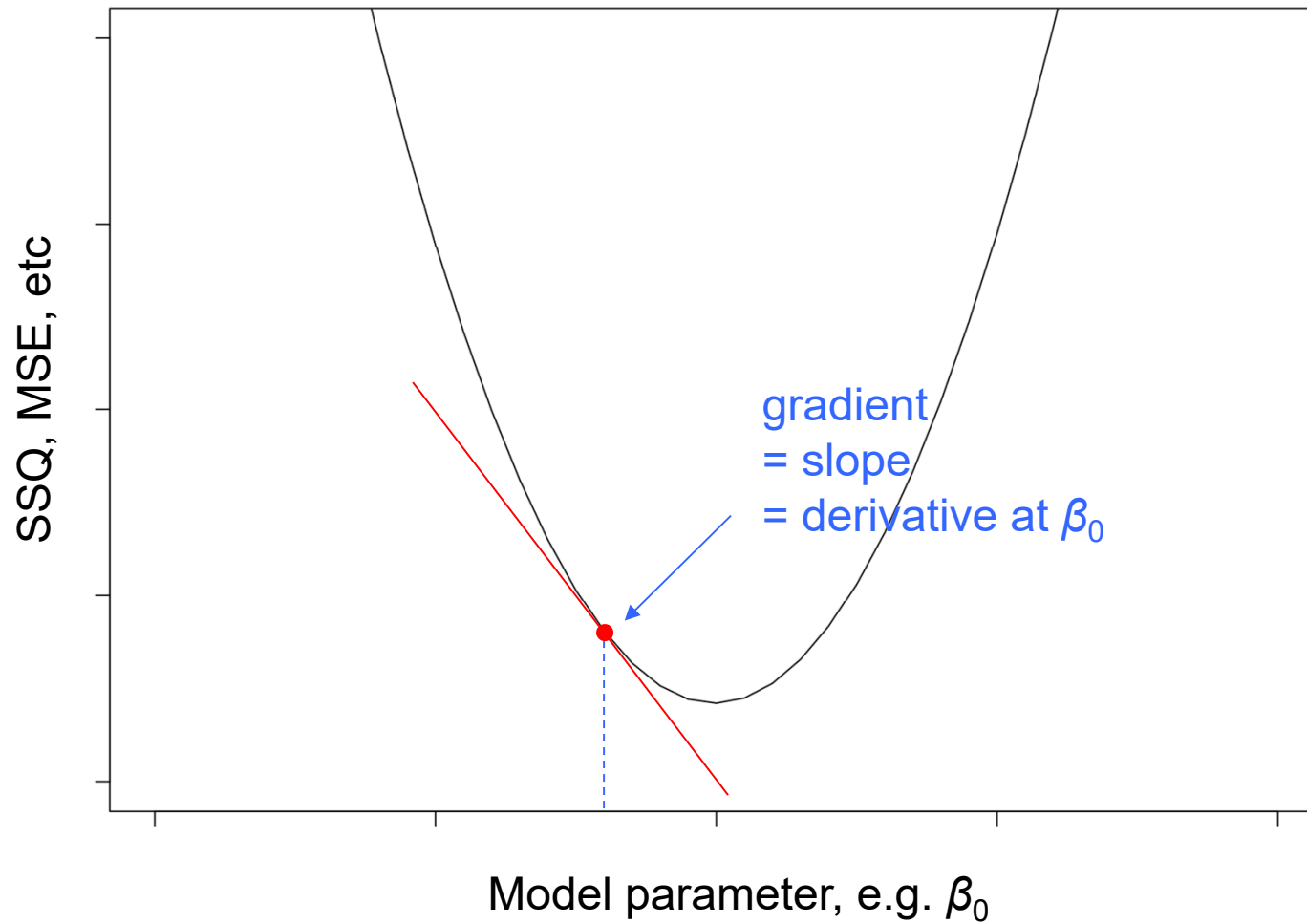


Today

- Ensemble methods
 - Bagging
 - Random forest
 - Boosting
- But first
 - basic gradient descent
 - boosting is a variant

Gradient descent



Finding the gradient

- It turns out that the gradient for a linear model is a function of the residuals
- See math

$$SSQ = \sum_i^n (y_i - \hat{y}_i)^2$$

$$= \sum_i^n r_i^2$$

$$= \sum_i^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$= \sum_i^n y_i^2 - 2y_i \beta_0 - 2y_i \beta_1 x_i + \beta_0^2 + 2\beta_0 \beta_1 x_i + \beta_1^2 x_i^2$$

$$\hat{y}_i = \beta_0 + \beta_1 x_i$$

$$r_i = y_i - \beta_0 - \beta_1 x_i$$

$$\frac{\partial SSQ}{\partial \beta_1} = \sum_i^n (-2y_i x_i + 2\beta_0 x_i + 2\beta_1 x_i^2)$$

$$= \sum_i^n -2x_i (y_i - \beta_0 - \beta_1 x_i)$$

$$= -2 \sum_i^n r_i x_i$$

$$\frac{\partial SSQ}{\partial \beta_0} = \sum_i^n (-2y_i + 2\beta_0 + 2\beta_1 x_i)$$

$$= -2 \sum_i^n r_i$$

Gradient descent

Gradient descent training algorithm for a linear model

set lambda

make initial guess for β_0, β_1

for many iterations

 find gradient at β_0, β_1

 step down: $\beta = \beta - \text{lambda} * \text{gradient}(\beta)$

print β_0, β_1

Gradient boosting

Gradient boosting algorithm (intuitive version)

set λ

fit a model to the data

calculate the left over variation (\hat{r})

fit a model to \hat{r} , the left over variation

calculate the new left over variation (\hat{r})

fit a model to \hat{r} again

calculate the new left over variation (\hat{r})

...

Keep going until we can no longer explain the variation

Boosted linear model

Algorithm

load y, x, x_{new}

guess parameters:

set $\hat{f}(x_{\text{new}}) = 0$

set $r \leftarrow y$ (residuals equal to the data)

for m in 1 to iterations

train model on r and x

predict residuals, $\hat{r}_b(x)$, from trained model

update residuals: $r \leftarrow r - \lambda \hat{r}_b(x)$

predict y increment, $\hat{f}_b(x_{\text{new}})$, from trained model

update prediction: $\hat{f}(x_{\text{new}}) \leftarrow \hat{f}(x_{\text{new}}) + \lambda \hat{f}_b(x_{\text{new}})$

return $\hat{f}(x_{\text{new}})$

Can be any
model



Gradient
descent



Gradient descent

predict residuals, $\hat{r}_b(x)$, from trained model

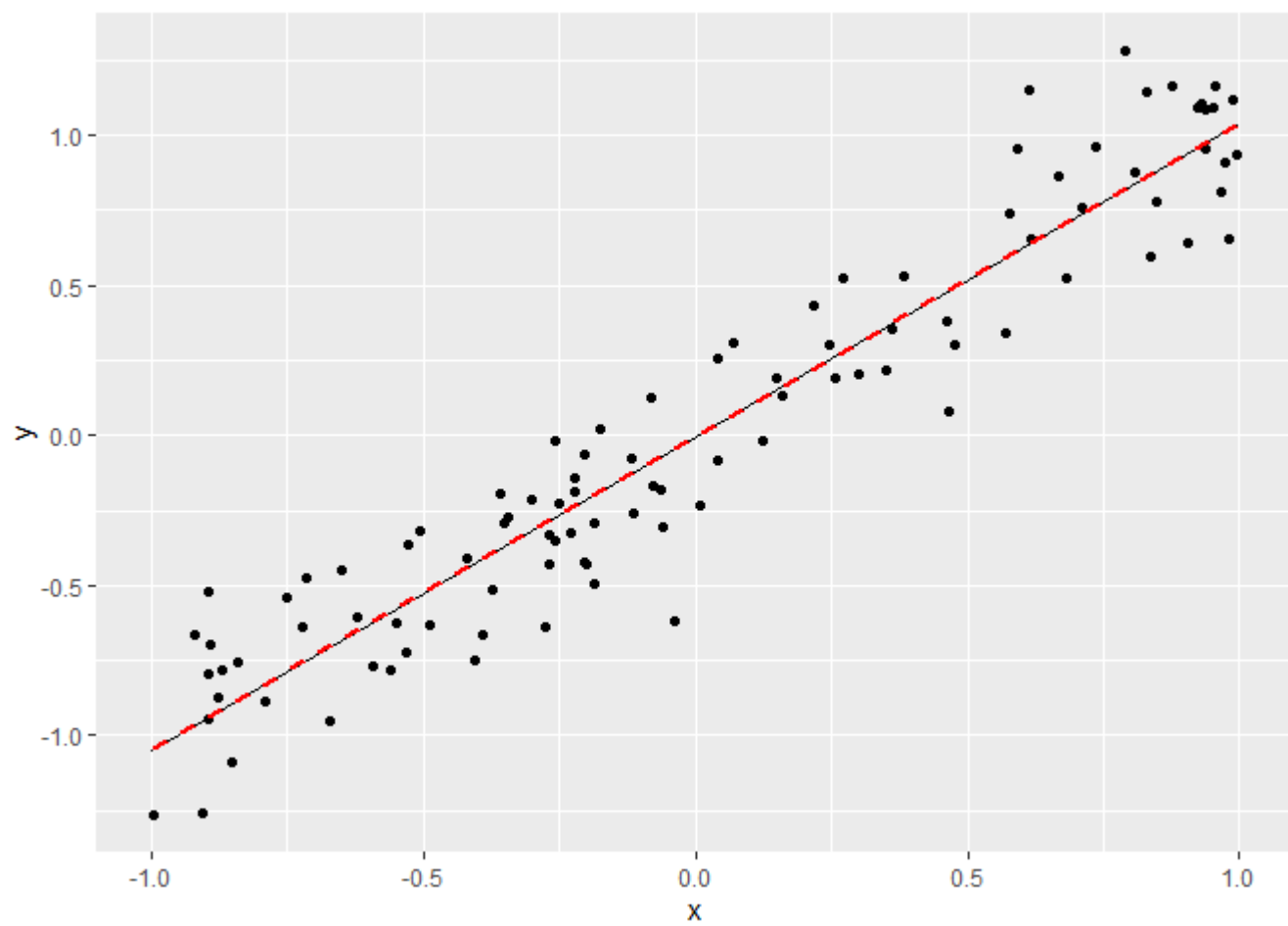
update residuals: $r \leftarrow r - \lambda \hat{r}_b(x)$

Loss function is MSE and we are descending its surface

Estimated gradient at x is the predicted residual $\hat{r}_b(x)$

λ is the **increment** taken down the gradient

r gets closer to 0 at each step, so MSE goes down



Boosted regression tree

Algorithm

load y, x, x_{new}

set parameters: mincut, ntrees, λ

set $\hat{f}(x_{\text{new}}) = 0$

set $r \leftarrow y$ (residuals equal to the data)

for b in 1 to ntrees

 train tree model on r and x

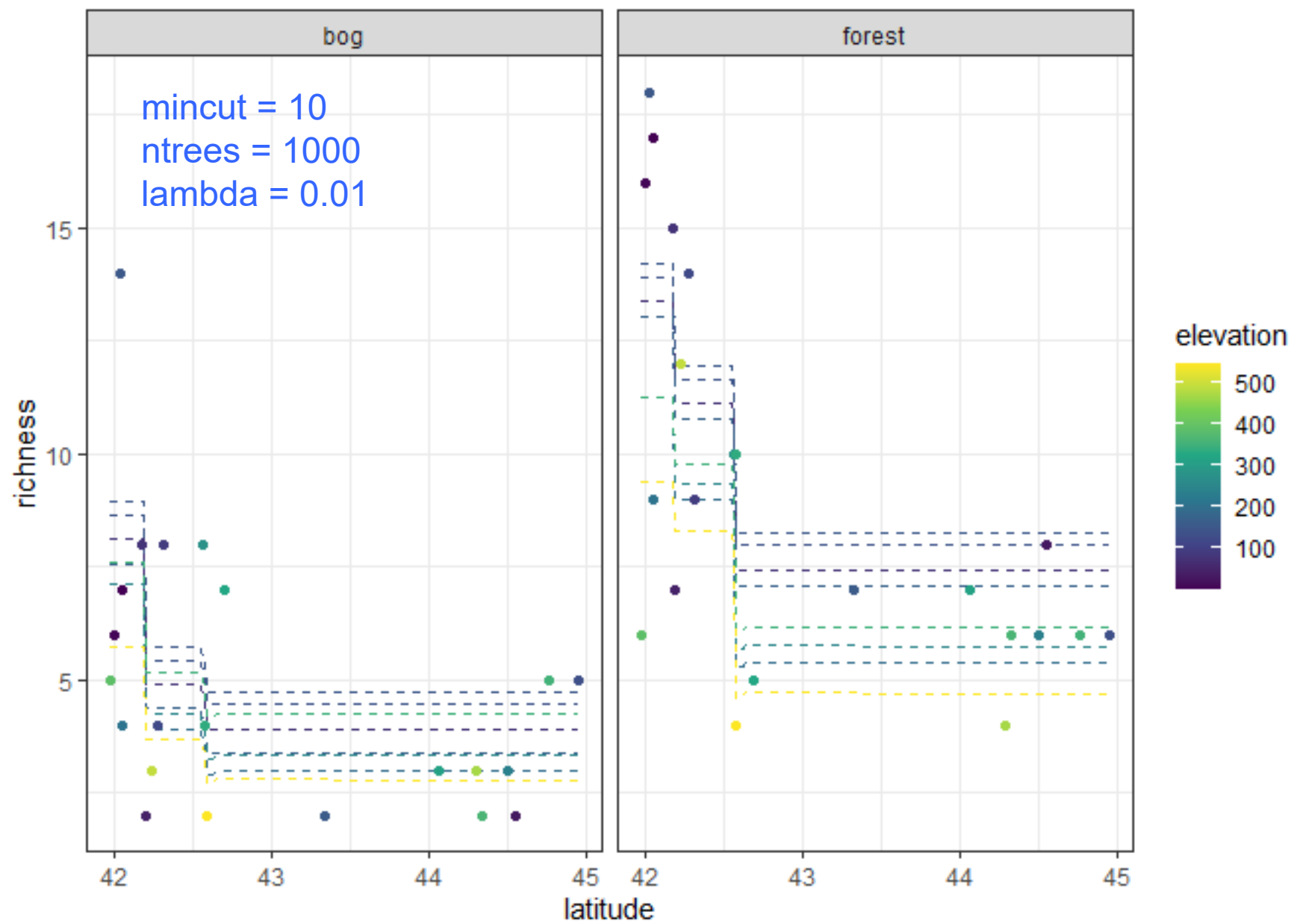
 predict residuals, $\hat{r}_b(x)$, from trained tree

 update residuals: $r \leftarrow r - \lambda \hat{r}_b(x)$

 predict y increment, $\hat{f}_b(x_{\text{new}})$, from trained tree

 update prediction: $\hat{f}(x_{\text{new}}) \leftarrow \hat{f}(x_{\text{new}}) + \lambda \hat{f}_b(x_{\text{new}})$

return $\hat{f}(x_{\text{new}})$



Gradient descent

