

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

Great! Now that we know all about this super simple **Training Dataset**, let's walk through the original **Gradient Descent** algorithm step-by-step.

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$  and a differentiable **Loss Function**  $L(y_i, F(x))$

In this case, a **Loss Function** is just something that evaluates how well we can predict **Weight**.

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

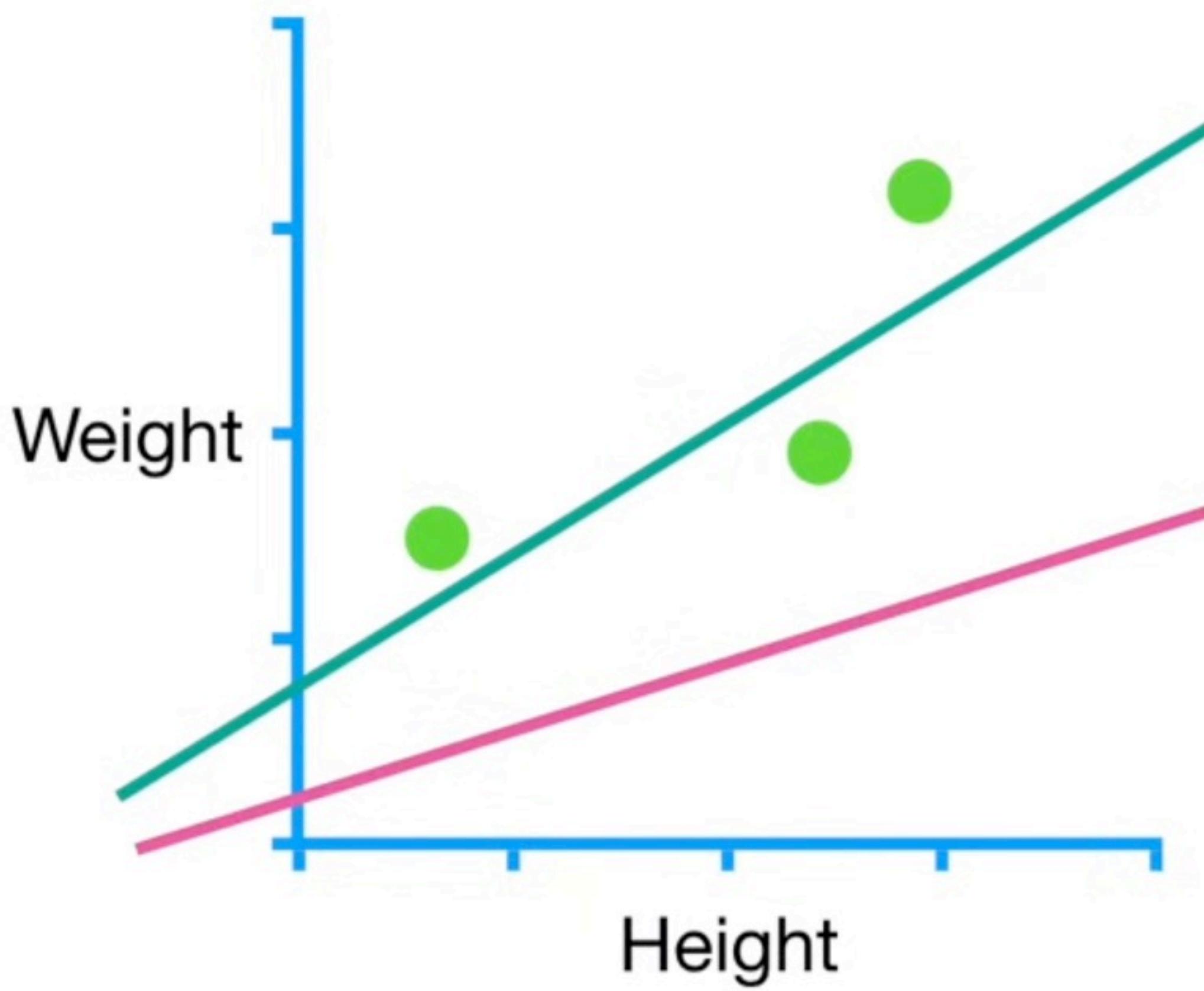
**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$  and a differentiable **Loss Function**  $L(y_i, F(x))$

The **Loss Function** that is most commonly used when doing **Regression with Gradient Boost** is...

$$\frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

In other words, it doesn't matter if  
the **Loss Function** is...



**(Observed - Predicted)<sup>2</sup>**

...of if it's...

$\frac{1}{2} (\text{Observed} - \text{Predicted})^2$

...since both will tell that  
the **Greenish Line** has  
the best fit.

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

The reason why people choose this **Loss Function** for **Gradient Boost** is that when we differentiate it with respect to “**Predicted**”...



$$\frac{d}{d \text{ Predicted}} \frac{1}{2} (\text{Observed} - \text{Predicted})^2$$



5:13 / 26:45

Step 0: The data and the loss function &gt;



**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

...and that leaves you with the  
**Observed** minus the **Predicted**  
multiplied by **-1**.

$$\frac{d}{d \text{ Predicted}} \frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

$$= \frac{2}{2} (\text{Observed} - \text{Predicted}) \times -1$$

$$= -(\text{Observed} - \text{Predicted})$$



**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

$$L(y_i, F(x))$$

The  $y_i$ 's are the **Observed** values...

$$\frac{1}{2} \boxed{\text{Observed}} - \text{Predicted}^2$$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56



**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y, F(x))$

...and  $F(x)$  is a function that gives us the **Predicted** values.

$$\frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56



5:57 / 26:45

Step 0: The data and the loss function &gt;



**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**NOTE:** There are other **Loss Functions** to choose from, but this is the most popular one for **Regression**.

$$\frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:

$$F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$$

We start by initializing  
the model with a  
constant value...

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:

$$F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$$

...and that constant  
value is determined by  
this funky looking thing.

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$

$$\frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

...and that funky symbol, called **gamma**, refers to the **Predicted** values.

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:  $F_0(x) = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, \gamma)$

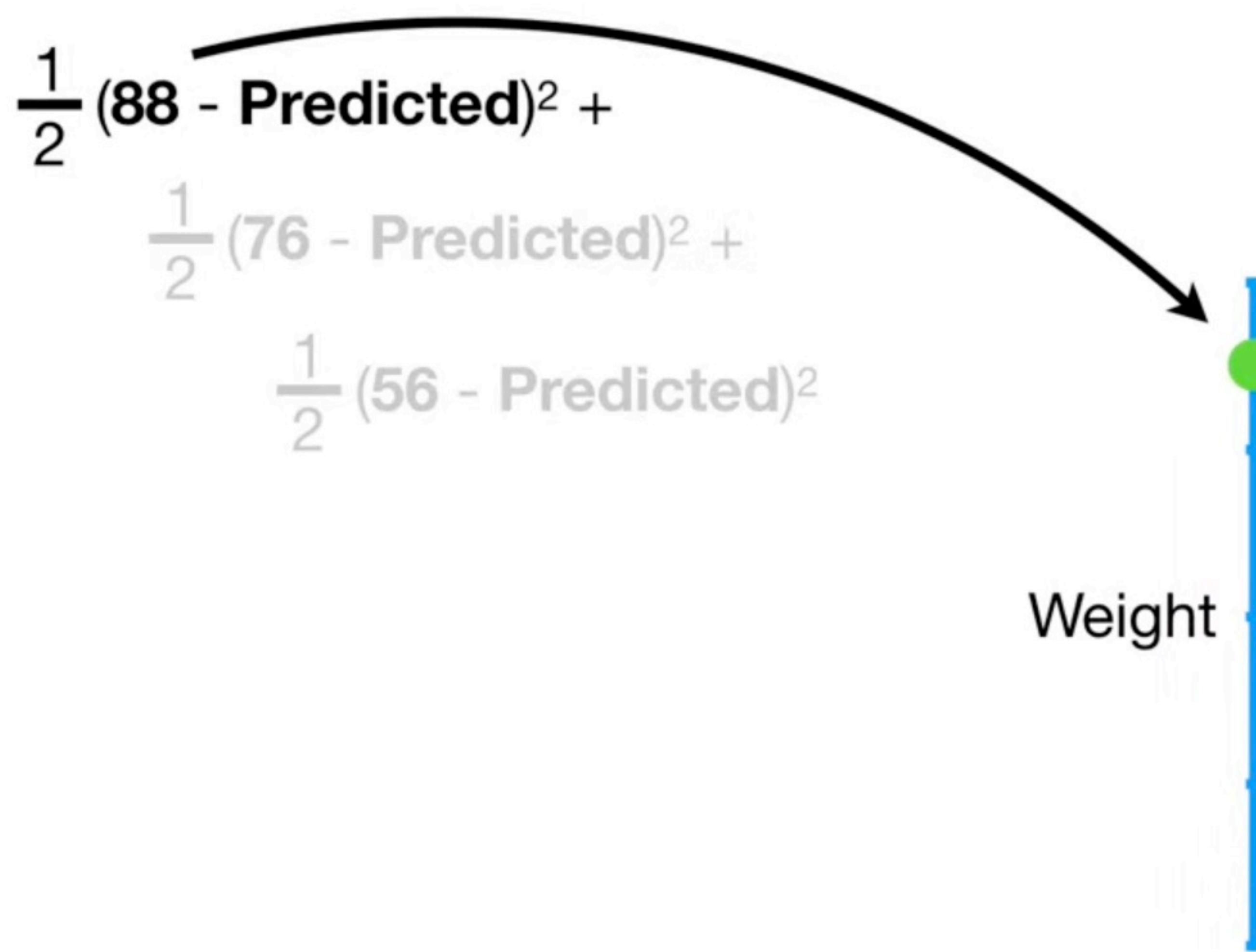
$$\frac{1}{2}(88 - \text{Predicted})^2 +$$
$$\frac{1}{2}(76 - \text{Predicted})^2 +$$
$$\frac{1}{2}(56 - \text{Predicted})^2$$

...and the “**argmin over gamma**” means we need to find a **Predicted** value that minimizes this sum.

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$



In other words, if we plot the  
**Observed Weights** on a  
number line...

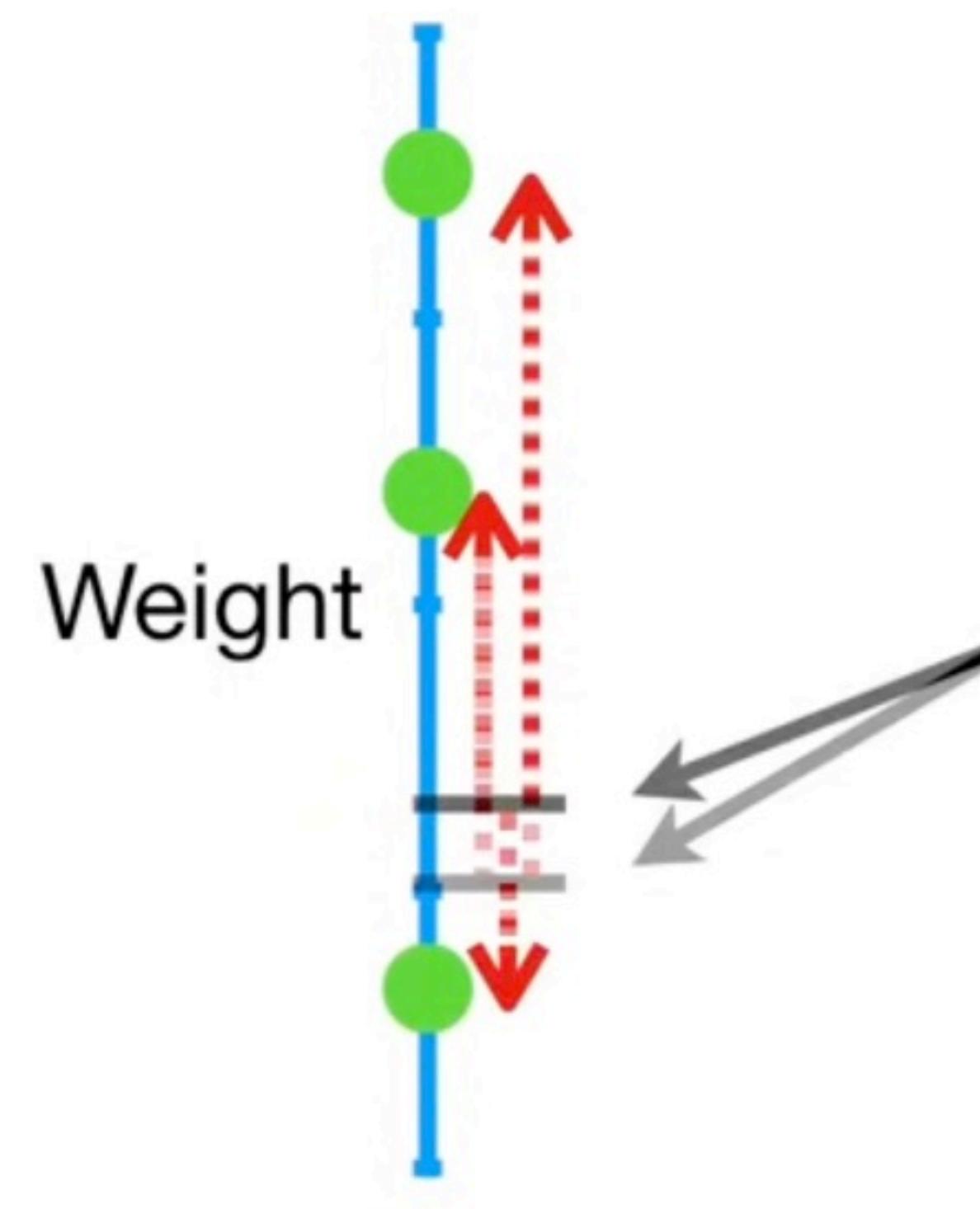
**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$

$$\frac{1}{2} (88 - \text{Predicted})^2 +$$

$$\frac{1}{2} (76 - \text{Predicted})^2 +$$

$$\frac{1}{2} (56 - \text{Predicted})^2$$



...then we want to find the point on the line that minimizes the sum of the squared residuals...

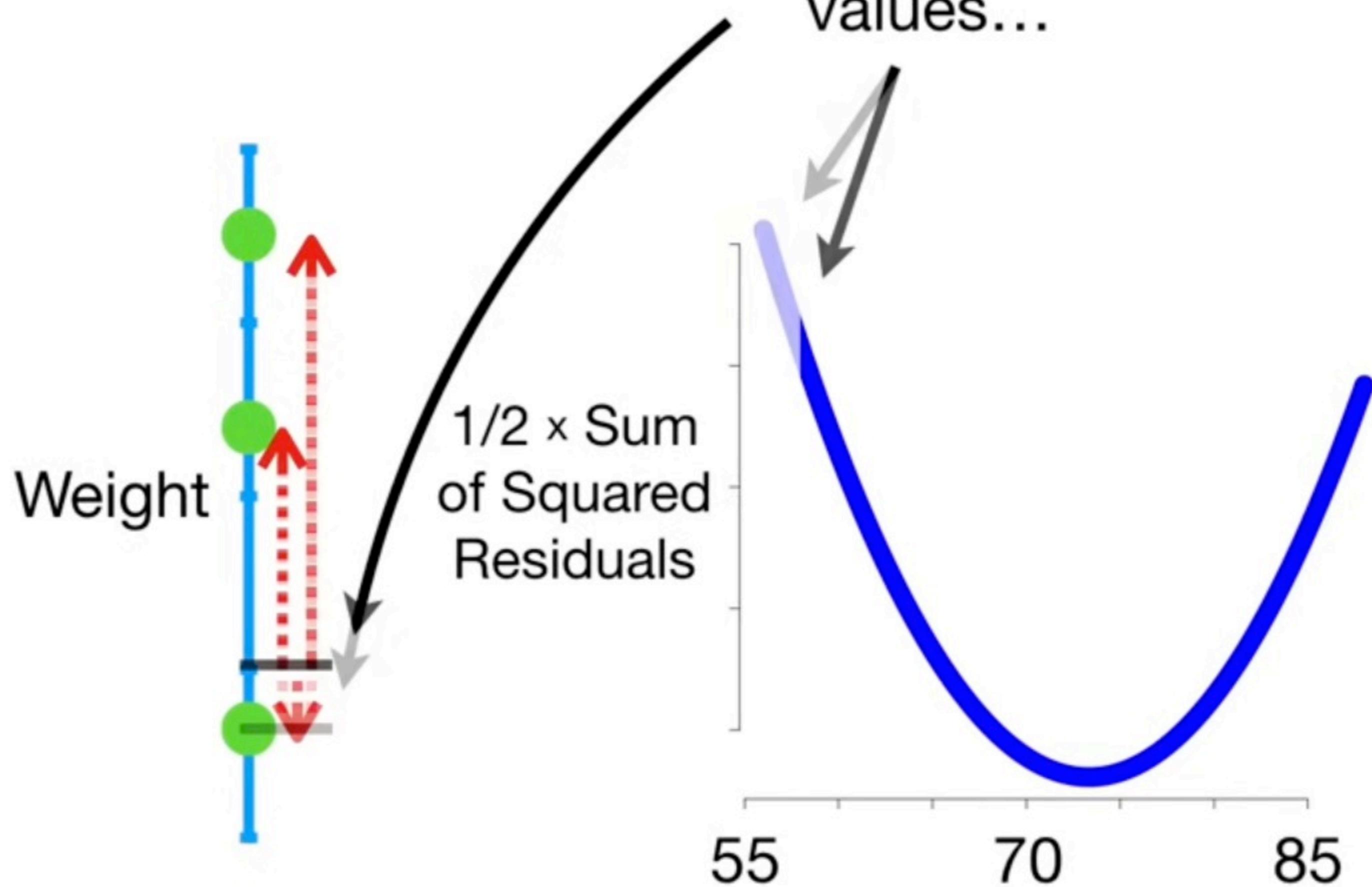
Input: Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Fu**

Step 1: Initialize model with a constant value:  $F_0(x)$

$$\frac{1}{2} (88 - \text{Predicted})^2 +$$

$$\frac{1}{2} (76 - \text{Predicted})^2 +$$

$$\frac{1}{2} (56 - \text{Predicted})^2$$



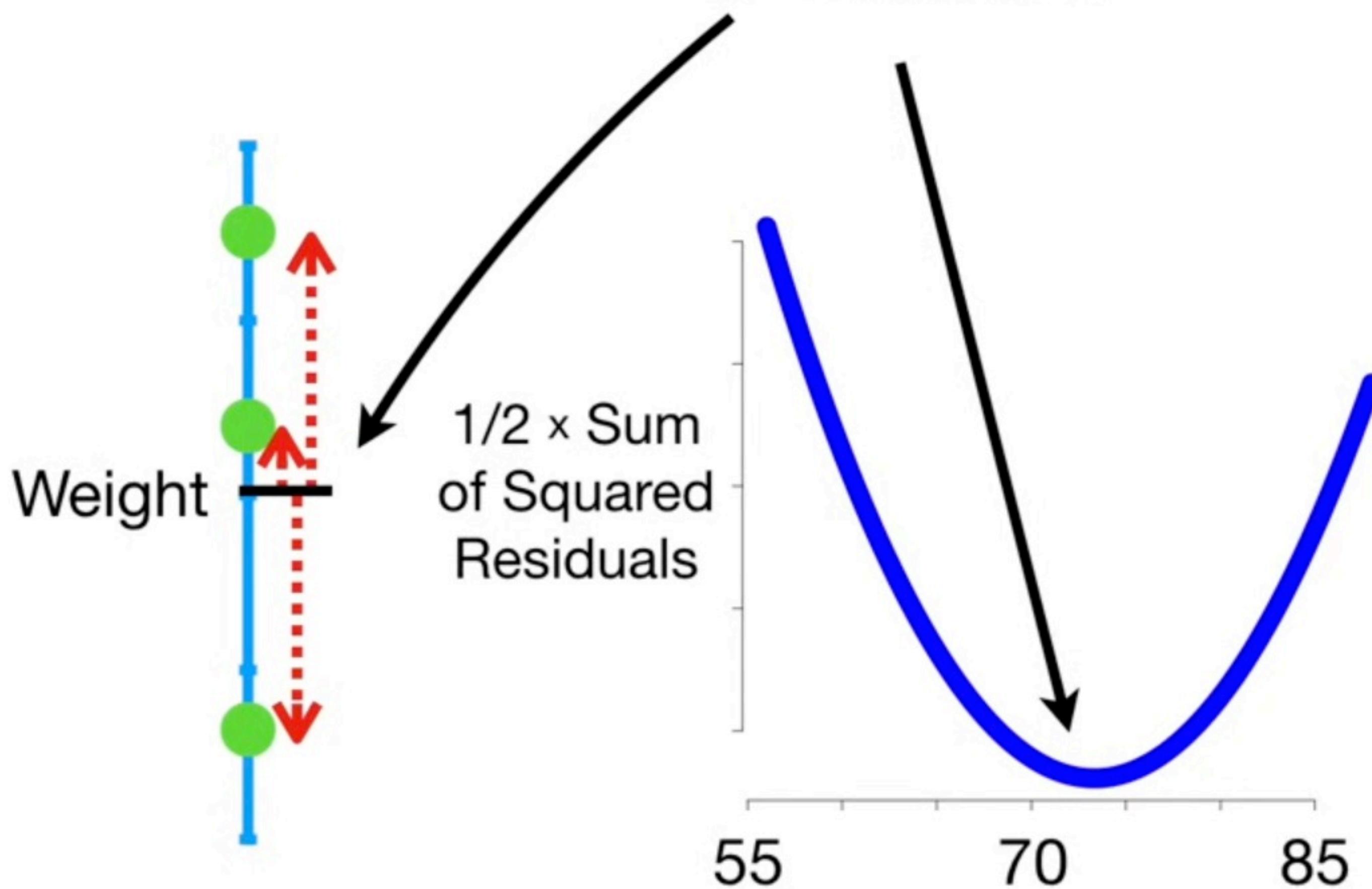
Input: Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Fu**

Step 1: Initialize model with a constant value:  $F_0(x)$

$$\frac{1}{2} (88 - \text{Predicted})^2 +$$

$$\frac{1}{2} (76 - \text{Predicted})^2 +$$

$$\frac{1}{2} (56 - \text{Predicted})^2$$



Input: Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Fu**

Step 1: Initialize model with a constant value:  $F_0(x)$

$$\frac{1}{2}(88 - \text{Predicted})^2 +$$
$$\frac{1}{2}(76 - \text{Predicted})^2 +$$
$$\frac{1}{2}(56 - \text{Predicted})^2$$

...but we can also just solve for it, because the math isn't that hard.

Input: Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

Step 1: Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$

$$\begin{aligned} & \frac{1}{2} (88 - \text{Predicted})^2 + \rightarrow -\mathbf{(88 - Predicted)} + \\ & \frac{1}{2} (76 - \text{Predicted})^2 + \rightarrow -\mathbf{(76 - Predicted)} + \boxed{=} 0 \\ & \frac{1}{2} (56 - \text{Predicted})^2 \rightarrow -\mathbf{(56 - Predicted)} \end{aligned}$$

$\frac{d}{d \text{Predicted}}$

Then we set sum of the derivatives equal to zero...

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

$$\text{Predicted} = \frac{88 + 76 + 56}{3}$$

...and we end up with  
the **Average** of the  
**Observed Weights.**

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:

$$F_0(x) = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$

...the value for **gamma**  
that minimizes this  
sum...

Predicted =  $\frac{88 + 76 + 56}{3}$

$$\frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$

...is the average of the  
**Observed Weights.**



$$\text{Predicted} = \frac{88 + 76 + 56}{3}$$



Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$

We have now created  
the initial predicted  
value,  $F_0(x)$ ...

$$F_0(x) = \frac{88 + 76 + 56}{3}$$

**Input:** Data  $\{(x_i, y_i)\}_{i=1}^n$ , and a differentiable **Loss Function**  $L(y_i, F(x))$

**Step 1:** Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$

That means that the initial predicted value,  $F_0(x)$ , is just a leaf.

$$F_0(x) = \frac{88 + 76 + 56}{3} = 73.3$$



73.3

The leaf predicts that all samples will weigh **73.3**.

**Step 2** is huge, but we'll take it one step at a time. :)

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$

(B) Fit a regression tree to the  $r_{im}$  values and create terminal regions  $R_{jm}$ , for  $j = 1 \dots J_m$

(C) For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_\gamma \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$

(D) Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

**Step 2:** for  $m = 1$  to  $M$ :



In generic terms, we will make  **$M$**  trees, but in practice, most people set  **$M = 100$**  and make **100** trees.

**Step 2:** for  $m = 1$  to  $M$ :



Little  $m$  refers to an individual tree. So when little  $m = 1$ , then we're talking about the first tree.

**Step 2:** for  $m = 1$  to  $M$ :



When little  $m$  = big  $M$ , then  
we're talking about the last  
tree...

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$

**Part A** of **Step 2** looks nasty,  
but it's not.

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = \boxed{-} \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$



This big minus sign tells us to multiply the derivative by -1...



-1 x -(Observed - Predicted)

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$

...and that leaves us with the  
**Observed** value minus the  
**Predicted** value.



**(Observed - Predicted)**

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$

In other words, this  
nasty looking thing...  ...is just a **Residual**.

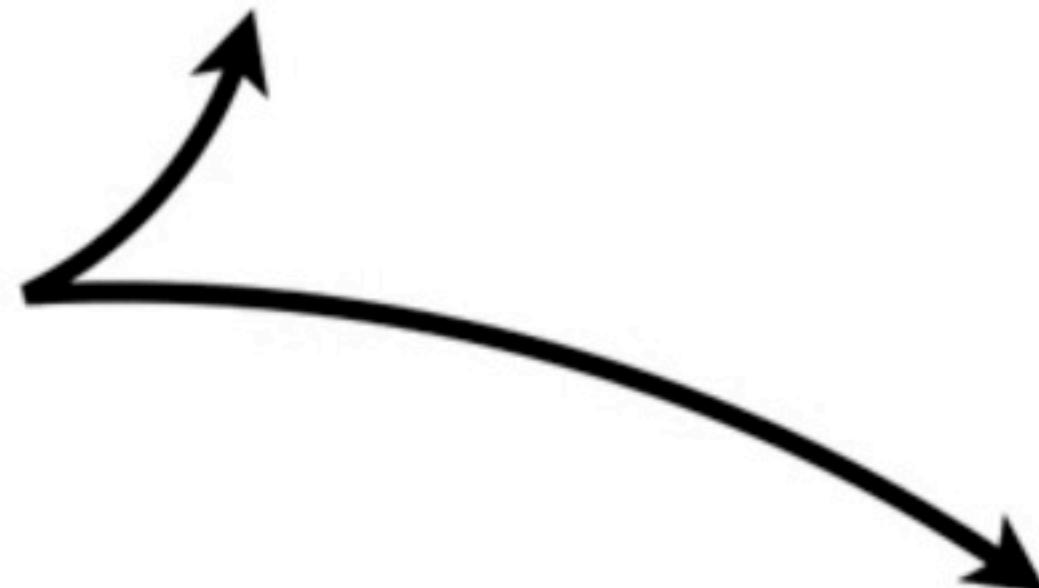
**(Observed - Predicted)**

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]$  for  $i = 1, \dots, n$

$F(x)=F_{m-1}(x)$

Now we plug  $F_{m-1}(x)$  in  
for **Predicted...**



(Observed -  $F_{m-1}(x)$ )

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]$  for  $i = 1, \dots, n$

$$F(x) = F_{m-1}(x)$$

...and since  $m = 1$ , that means we plug in  $F_0(x)$  for  $F_{m-1}(x)$ ...



(Observed -  $F_0(x)$ )

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im}$  =  $-\left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$



Now we can compute  $r_{i,m}$ , where  $r$  is short for **Residual**,  $i$  is the sample number and  $m$  is the tree that we're trying to build.

**(Observed - 73.3)**

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

This tells us to calculate  
**Residuals** for all 3  
samples in the dataset.

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = -\left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17.3

Hooray! We've finished **Part A** of  
**Step 2** by calculating a **Residual** for  
each sample.

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$



**NOTE:** Before we move on, I just want to point out that this derivative is the **Gradient** that **Gradient Boost** is named after.

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$

(B) Fit a regression tree to the  $r_{im}$  values and create terminal regions  $R_{jm}$ , for  $j = 1 \dots J_m$



All this is saying is that we will build a regression tree...

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$

(B) Fit a regression tree to the  $r_{im}$  values and create terminal regions  $R_{jm}$ , for  $j = 1 \dots J_m$

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17.3

...we predict the **Residuals** instead of the **Weights**.



Height < 1.55

-17.3

14.7, 2.7

Here's the new tree.

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17.3



Yes, I know this is just a stump and **Gradient Boost** almost always uses larger trees.

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7



Height < 1.55

-17.3

14.7, 2.7

However, in order to demonstrate details of the **Gradient Boost** algorithm, we need at least one leaf with more than 1 sample in it...

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17.3



However, in order to demonstrate details of the **Gradient Boost** algorithm, we need at least one leaf with more than **1** sample in it...

...and when you only have **3** samples, then you can't have more than **2** leaves.

Height (m)	Favorite Color	Gender		
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17.3

So we're stuck with using stumps, even though they are not typically used with **Gradient Boost**.

Height < 1.55

-17.3

14.7, 2.7

The **Residual** for the third sample,  $x_3$ ,  
goes to the leaf on the left...

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17.3

Height < 1.55

-17.3

14.7, 2.7

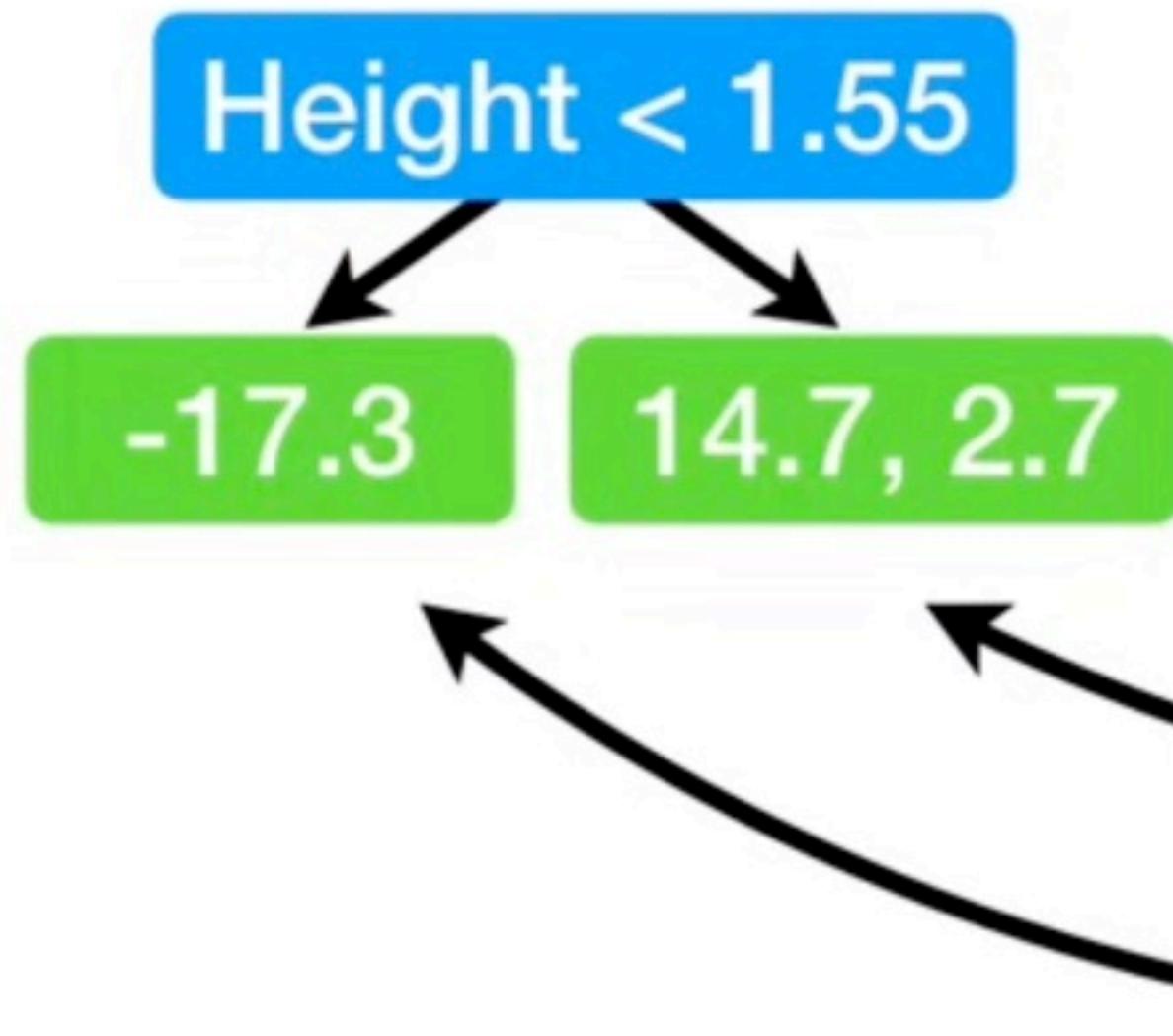
...and the **Residuals** for samples  
 $x_1$  and  $x_2$  go to the leaf on the  
right.

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	58	-17.3



Now we need to “create terminal regions  $R_{j,m}$ ”.

(B) Fit a regression tree to the  $r_{im}$  values and create terminal regions  $R_{jm}$ , for  $j = 1 \dots J_m$



This part is super easy because the **Leaves** are the “terminal regions  $R_{j,m}$ ”.

(B) Fit a regression tree to the  $r_{im}$  values and create terminal regions  $R_{jm}$ , for  $j = 1 \dots J_m$

Height < 1.55

-17.3

14.7, 2.7

This part is super easy because the **Leaves** are the “terminal regions  $R_{j,m}$ ”.

**NOTE:** This little  $m$  is the index for the tree we just made. Since this is the first tree,  $m = 1$ .

...and this little  $j$  is the index for each leaf in the tree.

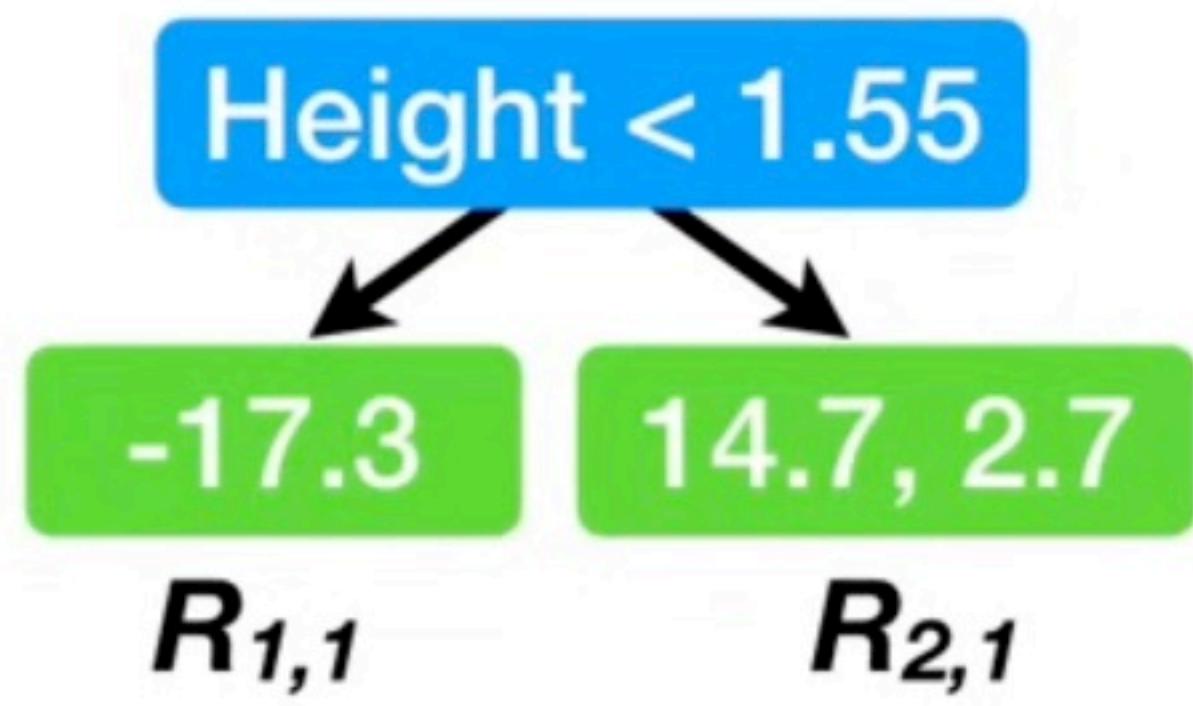
Height < 1.55

-17.3

14.7, 2.7

Since this tree has 2 leaves,  $J_m = 2$ .

(B) Fit a regression tree to the  $r_{im}$  values and create terminal regions  $R_{jm}$ , for  $j = 1 \dots J_m$



**NOTE:** It doesn't matter which leaf gets which label. However, once we give a leaf a label, we need to keep track of it.

- (B) Fit a regression tree to the  $r_{im}$  values and create terminal regions  $R_{jm}$ , for  $j = 1 \dots J_m$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

Hooray!!! We've finished  
**Part B** of **Step 2** by fitting a  
**Regression Tree** to the  
residuals and labeling the  
leaves.

(B) Fit a regression tree to the  $r_{im}$  values and create terminal  
regions  $R_{jm}$ , for  $j = 1 \dots J_m$

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17.3

## Now let's do **Part C.**

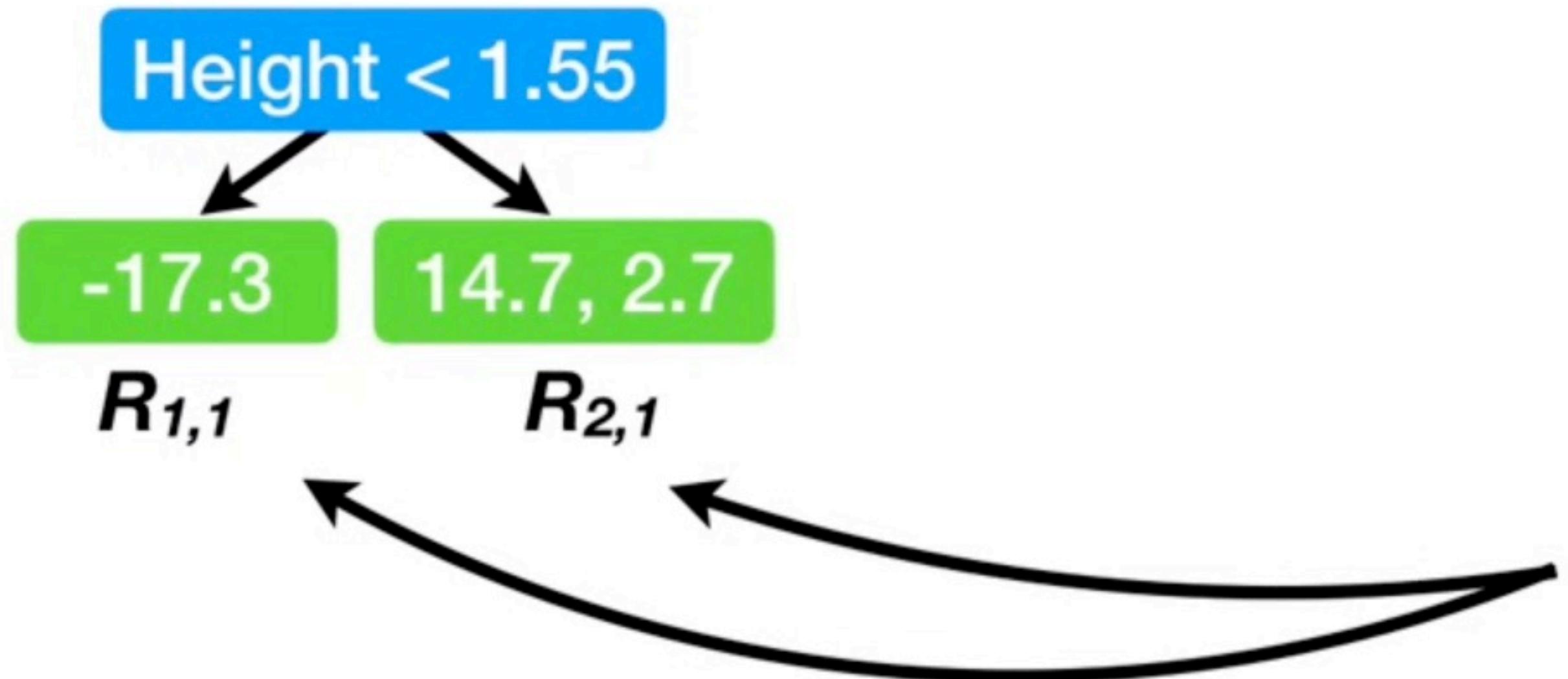
**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$

(B) Fit a regression tree to the  $r_{im}$  values and create terminal regions  $R_{jm}$ , for  $j = 1 \dots J_m$

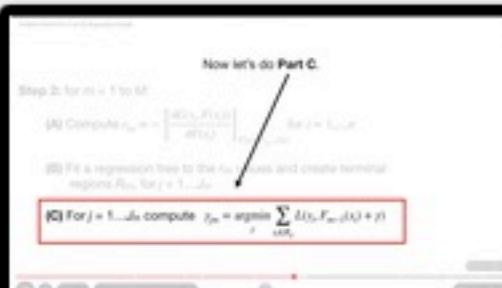
**(C)** For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_\gamma \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$





In this part, we determine the **Output Values** for each leaf.

**(C)** For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$





Specifically, since two residuals ended up in this leaf, it's unclear what its **Output Value** should be.

**(C)** For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$



So for each leaf in  
the new tree...

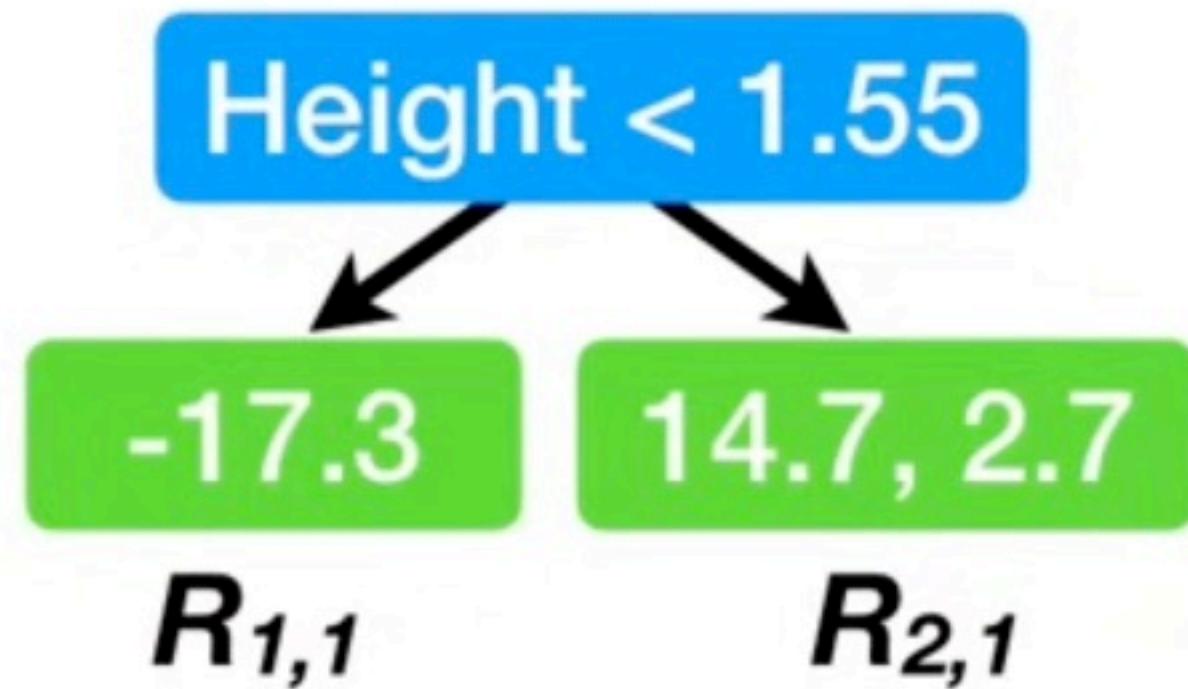
(C) For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$



15:07 / 26:45

Step 2.C: Optimize leaf output values &gt;





...we compute an  
**Output Value,**  
**“gamma<sub>j,m</sub>”.**

(C) For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$



Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

The **Output Value** for each leaf is the value for **gamma** that minimizes this summation.

(C) For  $j = 1 \dots J_m$  compute

$$\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$$

**Step 1:** Initialize model with a constant value:

$$F_0(x) = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$

**NOTE:** This minimization is like what we did in **Step 1**.

(C) For  $j = 1 \dots J_m$  compute

$$\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$$

**Step 1:** Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$

One small difference is that  
now we are taking the  
previous **Prediction** into  
account...



(C) For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$

Step 1: Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$

One small difference is that  
now we are taking the  
previous **Prediction** into  
account...

...while before, since we  
were just starting out,  
there was no “previous  
**Prediction**”.

(C) For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_\gamma \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$

**Step 1:** Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$

The other difference is that  
this summation is picky  
about which samples it  
includes.



(C) For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$

Step 1: Initialize model with a constant value:  $F_0(x) = \operatorname{argmin}_\gamma \sum_{i=1}^n L(y_i, \gamma)$

$$\sum_{i=1}^n$$



...while before, the summation included all of the samples.

The other difference is that this summation is picky about which samples it includes.



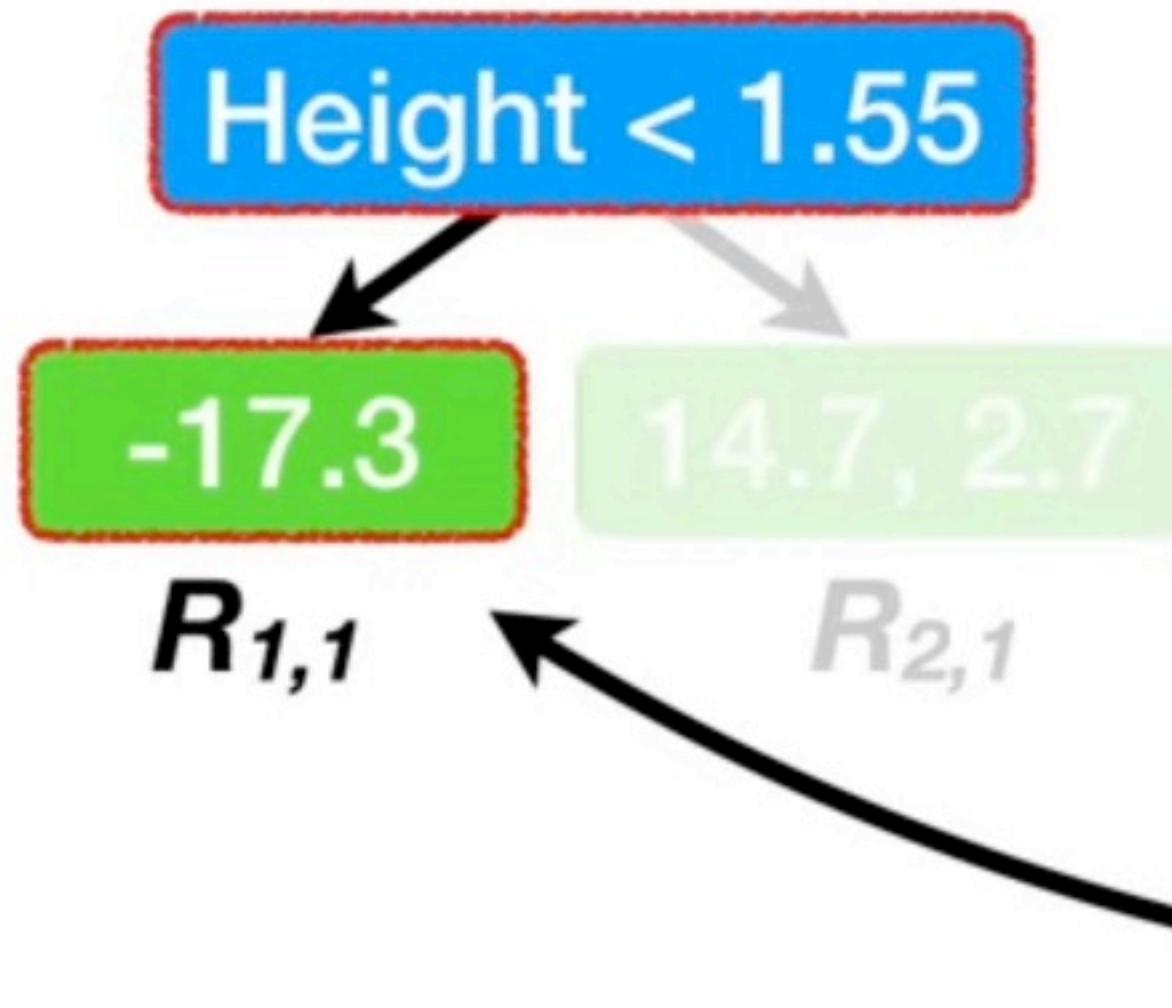
(C) For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_\gamma \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$

$$\sum_{x_i \in R_{ij}}$$

Specifically, the  $x_i$  in  $R_{i,j}$   
means that...



(C) For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \operatorname{argmin}_\gamma \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$



...since only sample,  
 $x_3$ , goes to leaf  $R_{1,1}$ ...

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17

$$\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$$



...then only  $x_3$  is used to calculate the **Output**

**Value** for  $R_{1,1}$ ...

$$\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$$

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17



...and since only two samples,  $x_1$  and  $x_2$ , go to leaf  $R_{1,2} \dots$

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	55	-17

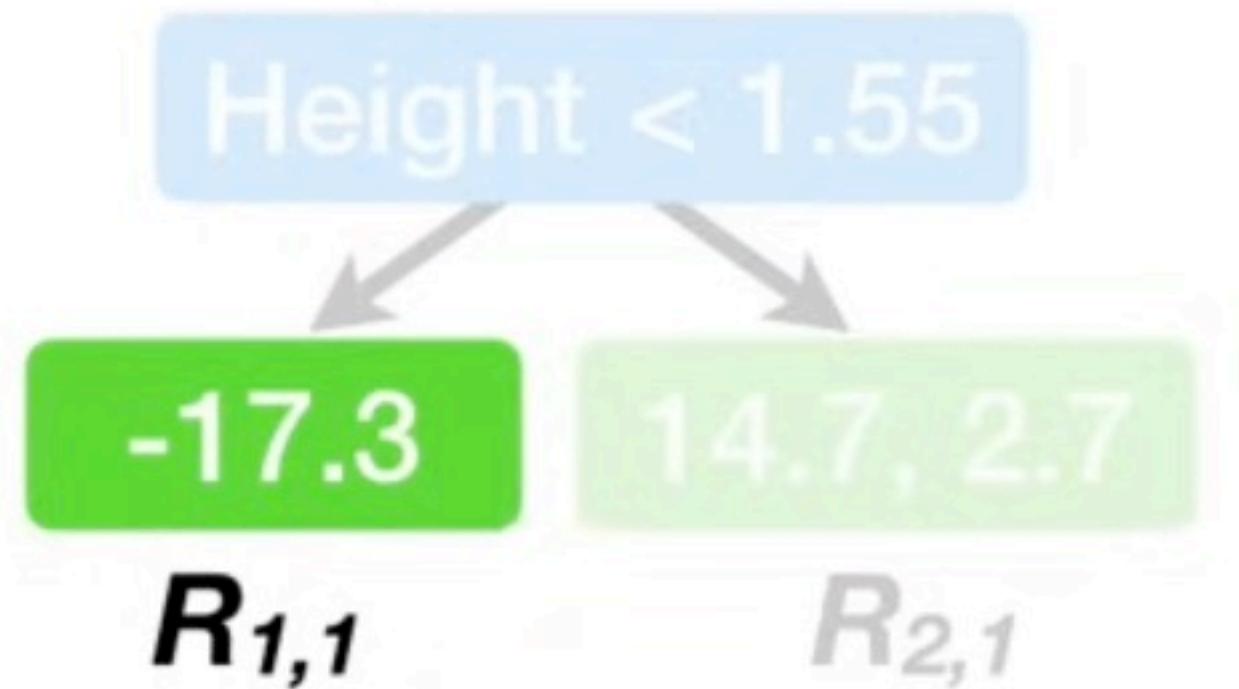
$$\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$$



Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17

...then only  $x_1$  and  $x_2$  are used to calculate the **Output Value** for  $R_{1,2}$ .

$$\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$$



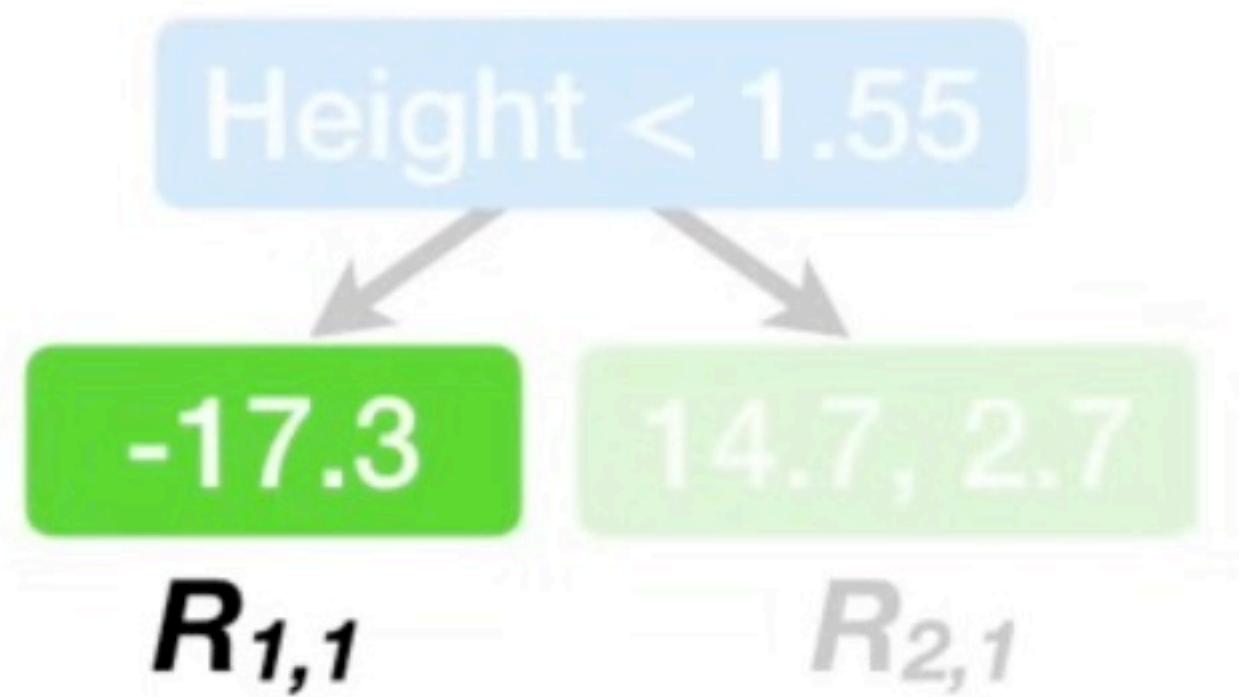
Let's start by calculating  
the **Output Value** for the  
leaf on the left,  $R_{1,1}$ .

$$\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$$



That means  $j = 1$ , since this is the first leaf, and  $m = 1$ , since this is the first tree.

$$\gamma_{1,1} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$$



Now let's replace the generic **Loss Function**...

$$\frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

...with the actual **Loss Function** that we decided to use...

$$\gamma_{1,1} = \operatorname{argmin}_{\gamma} \sum_{x_i \in K_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$$



...and let's expand the summation into individual terms.

$$\gamma_{1,1} = \underset{\gamma}{\operatorname{argmin}} \sum_{x_i \in R_{ij}} \frac{1}{2} (y_i - (F_{m-1}(x_i) + \gamma))^2$$



Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17.3

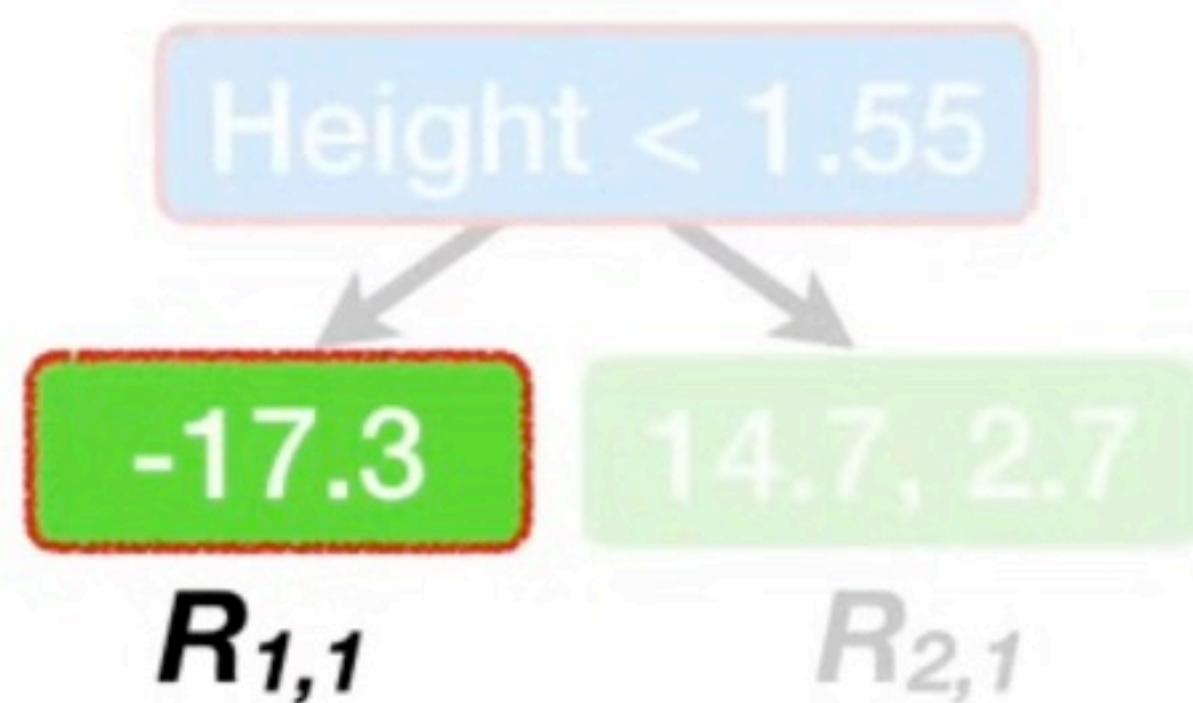
Now we plug in the value for  $y_3$ , the **Observed** value...

$$\gamma_{1,1} = \underset{\gamma}{\operatorname{argmin}} \frac{1}{2}(y_3 - (F_{m-1}(x_3) + \gamma))^2$$



...and the most recent  
**Predicted** value for  $x_3$ .

$$\gamma_{1,1} = \operatorname{argmin}_{\gamma} \frac{1}{2}(56 - (F_{m-1}(x_3) + \gamma))^2$$



$F_0(x)$   
73.3

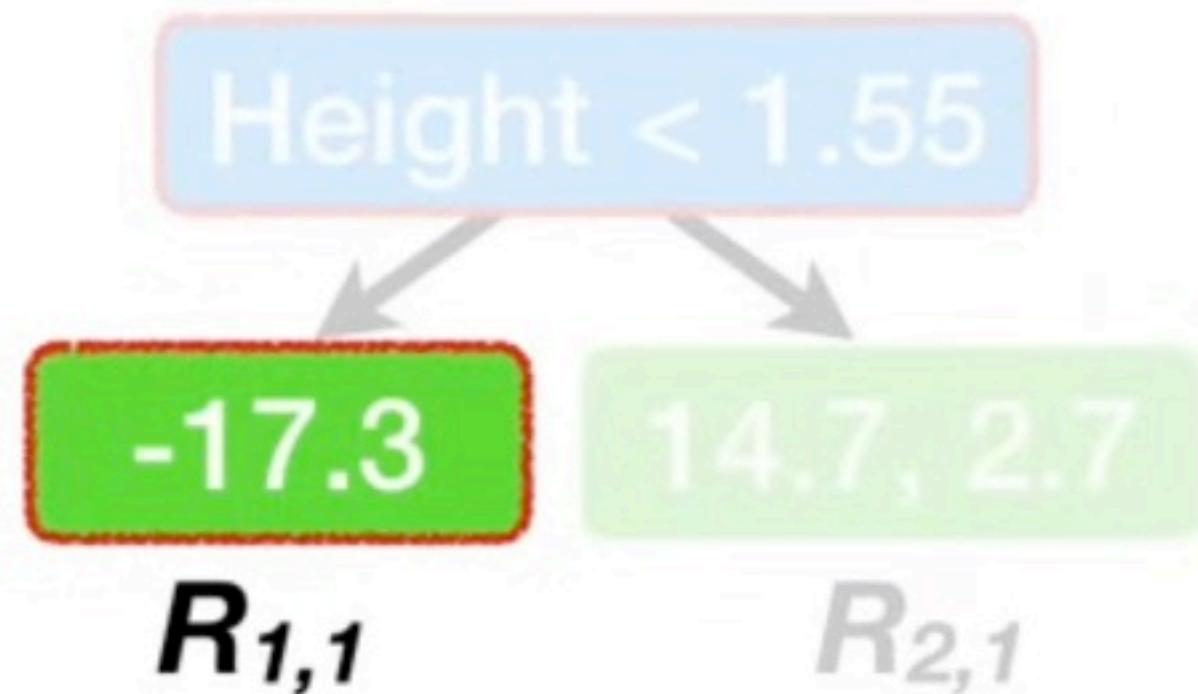
Since  $m = 1$ , the most recent **Prediction** was  $F_0(x)$ , which **Predicted** that all samples weighed 73.3...

$$\gamma_{1,1} = \operatorname{argmin}_{\gamma} \frac{1}{2}(56 - (F_{m-1}(x_3) + \gamma))^2$$



...and simplify what's inside  
the parentheses.

$$\gamma_{1,1} = \operatorname{argmin}_{\gamma} \frac{1}{2} (-17.3 - \gamma)^2$$



Now we need to find the value for **gamma** that minimizes this equation.

$$\gamma_{1,1} = \operatorname{argmin}_{\gamma} \frac{1}{2}(-17.3 - \gamma)^2$$



17:33 / 26:45

Step 2.C: Optimize leaf output values &gt;



**Step 1:** Initialize model with a constant value:

$$F_0(x) = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$

Just like **Step 1**, we can try different values for **gamma** or solve it analytically.



$$\gamma_{1,1} = \operatorname{argmin}_{\gamma} \frac{1}{2}(-17.3 - \gamma)^2$$

$$\frac{d}{d\gamma} \frac{1}{2}(-17.3 - \gamma)^2$$

First, we take the derivative of the **Loss Function** with respect to **gamma**, just like we did at the very start.

$$\gamma_{1,1} = \operatorname{argmin}_{\gamma} \frac{1}{2}(-17.3 - \gamma)^2$$

$$\frac{d}{d\gamma} \frac{1}{2}(-17.3 - \gamma)^2 \rightarrow 17.3 + \gamma = 0$$

Now we set the derivative equal to **0**...

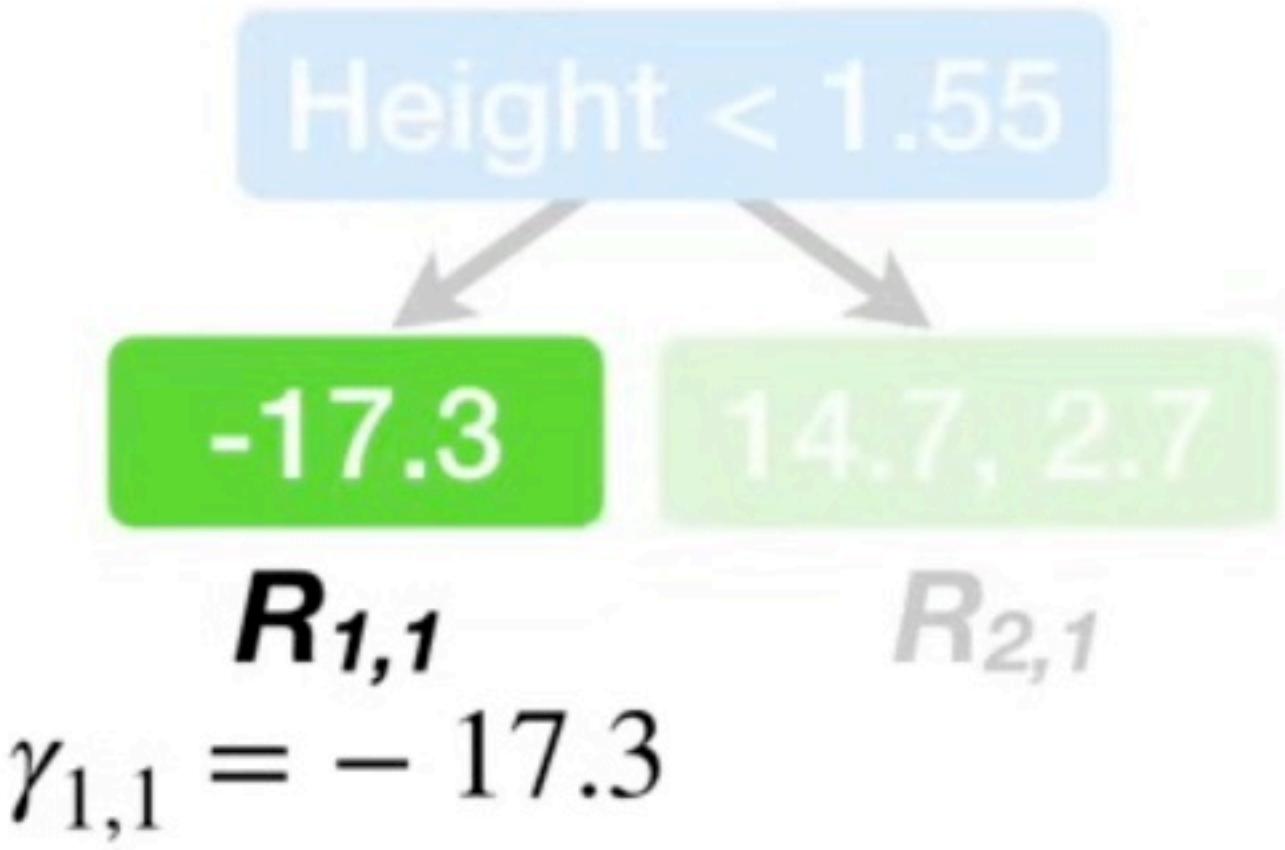
$$\gamma_{1,1} = \operatorname{argmin}_{\gamma} \frac{1}{2}(-17.3 - \gamma)^2$$

$$\gamma = -17.3$$



The value for **gamma** that minimizes this equation is **-17.3**.

$$\gamma_{1,1} = \operatorname{argmin}_{\gamma} \frac{1}{2}(-17.3 - \gamma)^2$$



$$\gamma_{1,1} = -17.3$$

...and ultimately, the leaf,  $R_{1,1}$ , has an **Output Value** of -17.3.

$$\gamma_{1,1} = -17.3$$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

$\gamma_{1,1} = -17.3$

Now let's solve for the  
**Output Value** for  $R_{2,1}$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

$\gamma_{1,1} = -17.3$

That means  $j = 2$ , since this is  
the second leaf, and  $m = 1$ ,  
since this is still the first tree.

$$\boxed{\gamma_{2,1}} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + \gamma)$$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

$$\gamma_{1,1} = -17.3$$

Height (m)	Favorite Color	Gender	Weight (kg)	$r_{i,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17.3

...plug in the **Observed Weights**...

$$\gamma_{2,1} = \operatorname{argmin}_{\gamma} \left[ \frac{1}{2} (y_1 - (F_{m-1}(x_1) + \gamma))^2 + \frac{1}{2} (y_2 - (F_{m-1}(x_2) + \gamma))^2 \right]$$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

$\gamma_{1,1} = -17.3$

...and plug in 73.3  
for  $F_{m-1}(x_1)$  and  
 $F_{m-1}(x_2)$ .

$F_0(x)$

73.3

$$\gamma_{2,1} = \operatorname{argmin}_\gamma \left[ \frac{1}{2} (88 - (F_{m-1}(x_1) + \gamma))^2 + \frac{1}{2} (76 - (F_{m-1}(x_2) + \gamma))^2 \right]$$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

$$\gamma_{1,1} = -17.3$$

Now simplify what's inside the parentheses...

$$\gamma_{2,1} = \operatorname{argmin}_{\gamma} \left[ \frac{1}{2}(14.7 - \gamma)^2 + \frac{1}{2}(2.7 - \gamma)^2 \right]$$

$$\frac{d}{d\gamma} \left[ \frac{1}{2}(14.7 - \gamma)^2 + \frac{1}{2}(2.7 - \gamma)^2 \right] \longrightarrow -14.7 + \gamma + -2.7 + \gamma$$



...and that gives us  
this derivative...

$$\gamma_{2,1} = \operatorname{argmin}_{\gamma} \left[ \frac{1}{2}(14.7 - \gamma)^2 + \frac{1}{2}(2.7 - \gamma)^2 \right]$$

$$\frac{d}{d\gamma} \left[ \frac{1}{2}(14.7 - \gamma)^2 + \frac{1}{2}(2.7 - \gamma)^2 \right] \longrightarrow \quad 2\gamma = 14.7 + 2.7$$

...and then we solve.

$$\gamma_{2,1} = \operatorname*{argmin}_{\gamma} \left[ \frac{1}{2}(14.7 - \gamma)^2 + \frac{1}{2}(2.7 - \gamma)^2 \right]$$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$\gamma_{1,1} = -17.3$

$$\gamma = \frac{14.7 + 2.7}{2}$$

We end up with the average of the  
**Residuals** that ended in leaf  $R_{2,1}$ .

$$\begin{aligned}\gamma_{2,1} = \operatorname{argmin}_{\gamma} & \left[ \frac{1}{2}(14.7 - \gamma)^2 \right. \\ & \left. + \frac{1}{2}(2.7 - \gamma)^2 \right]\end{aligned}$$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

$\gamma_{1,1} = -17.3$

$$\gamma = \frac{14.7 + 2.7}{2} = 8.7$$

So the value for **gamma** that minimizes this equation is **8.7...**

$$\gamma_{2,1} = \operatorname{argmin}_{\gamma} \left[ \frac{1}{2}(14.7 - \gamma)^2 + \frac{1}{2}(2.7 - \gamma)^2 \right]$$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

$\gamma_{1,1} = -17.3$

$\gamma_{2,1} = 8.7$

...and ultimately, the  
leaf,  $R_{2,1}$ , has an  
**Output Value of 8.7.**

$\gamma_{2,1} = 8.7$



We just saw that the **Output Value** for this leaf,  $R_{2,1}$ , is the average of the residuals that ended up here.



Given our choice of **Loss Function**,  
the **Output Values** are *always* the  
average of the **Residuals** that end  
up in the same leaf.

$$\frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

## Now let's do **Part D!!!**

**Step 2:** for  $m = 1$  to  $M$ :

(A) Compute  $r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=L_{m-1}(x)}$  for  $i = 1, \dots, n$

(B) Fit a regression tree to the  $r_{im}$  values and create terminal regions  $R_{jm}$ , for  $j = 1 \dots J_m$

(C) For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \underset{\gamma}{\operatorname{argmin}} \sum_{x \in R_{jm}} L(y_i, F_{m-1}(x_i) + \gamma)$

(D) Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

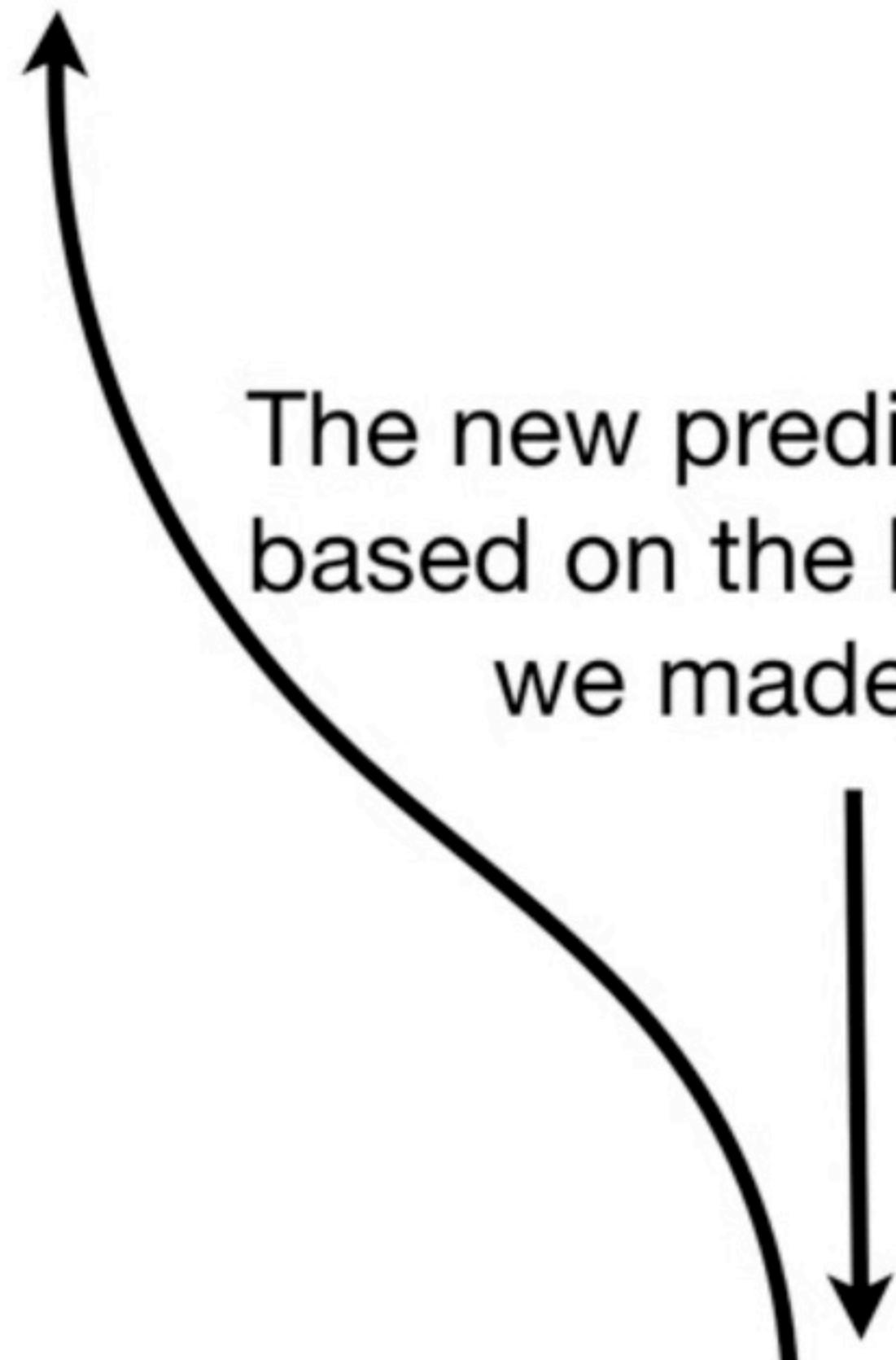
Since this is our first pass through **Step 2** and  $m = 1$ , this new prediction will be called  $F_1(x)$ .



(D) Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

$F_0(x)$

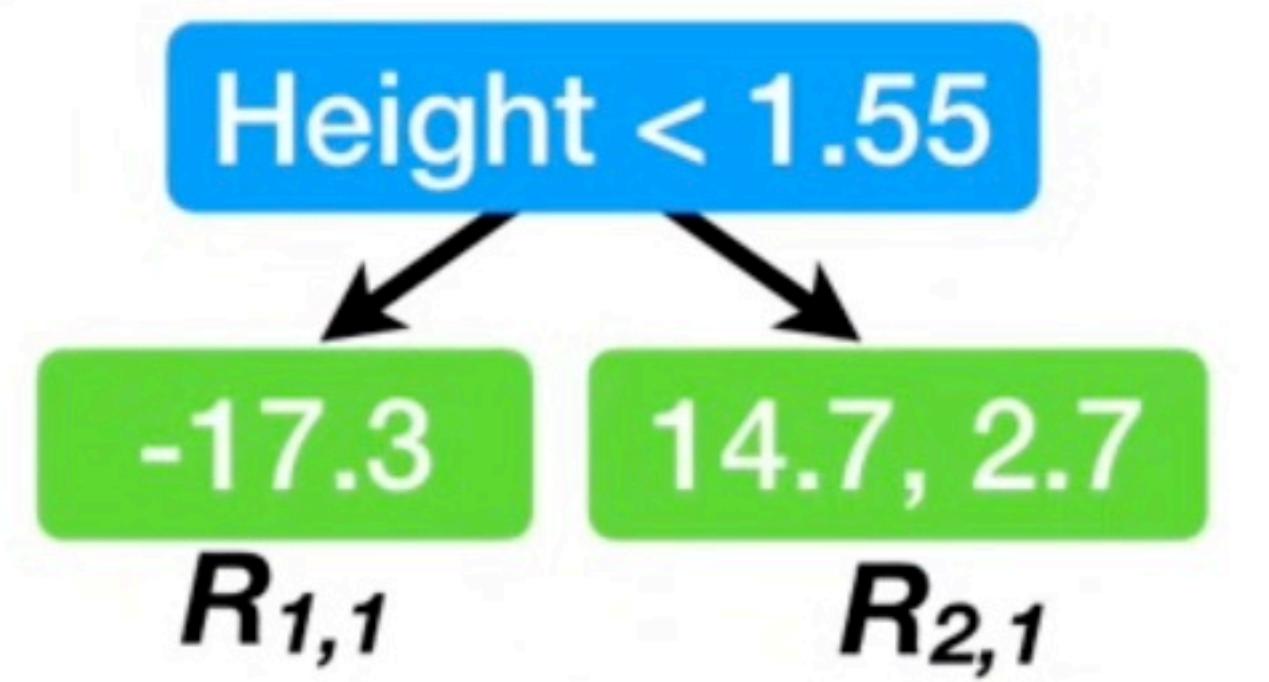
$F_1(x) =$  73.3



The new prediction,  $F_1(x)$ , is  
based on the last prediction  
we made,  $F_0(x)$ ...

**(D)** Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

$$F_0(x) \\ F_1(x) = 73.3 +$$

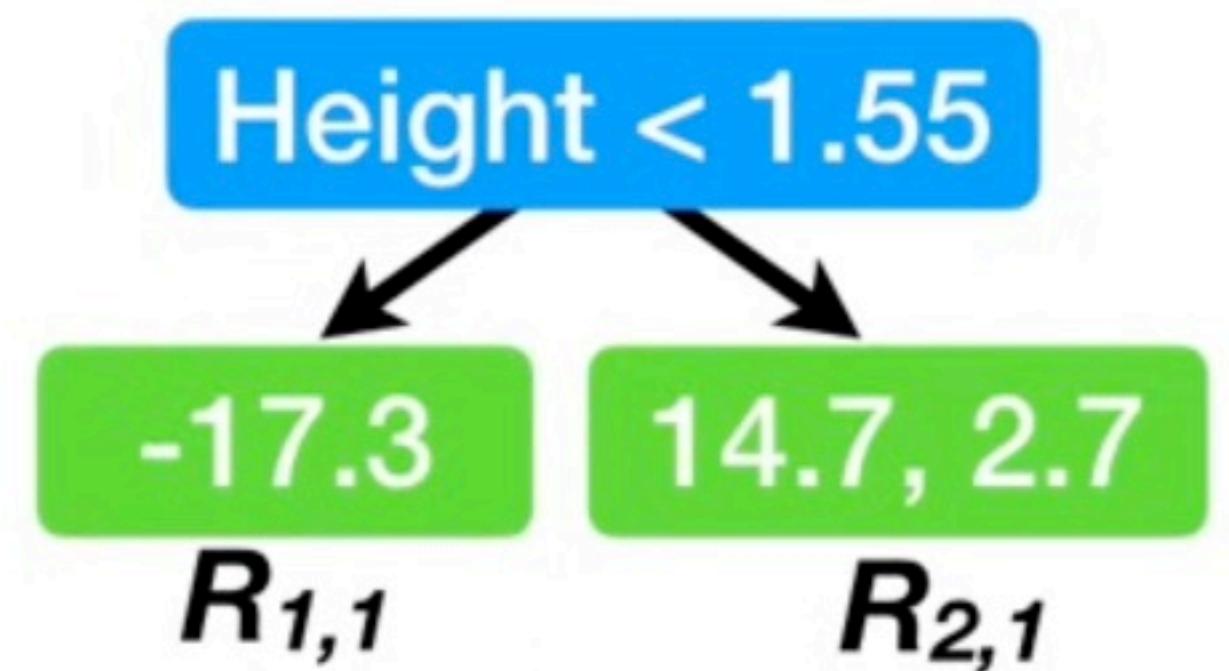


$$\gamma_{1,1} = -17.3 \quad \gamma_{2,1} = 8.7$$

...and the tree that we  
just finished making.

(D) Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

$$F_0(x) \\ F_1(x) = 73.3 +$$



$$\gamma_{1,1} = -17.3 \quad \gamma_{2,1} = 8.7$$

The summation says we should add up the **Output Values,  $\gamma_{j,m}$ 's**, for all the leaves,  $R_{j,m}$ , that a sample,  $x$ , can be found in.



(D) Update  $F_m(x) = F_{m-1}(x) + \gamma_{jm} I(x \in R_{jm})$

$$\sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$$

$$F_0(x) \\ F_1(x) = 73.3 + \begin{array}{c} \text{Height} < 1.55 \\ \downarrow \quad \downarrow \\ -17.3 \quad 14.7, 2.7 \\ R_{1,1} \quad R_{2,1} \\ \gamma_{1,1} = -17.3 \quad \gamma_{2,1} = 8.7 \end{array}$$

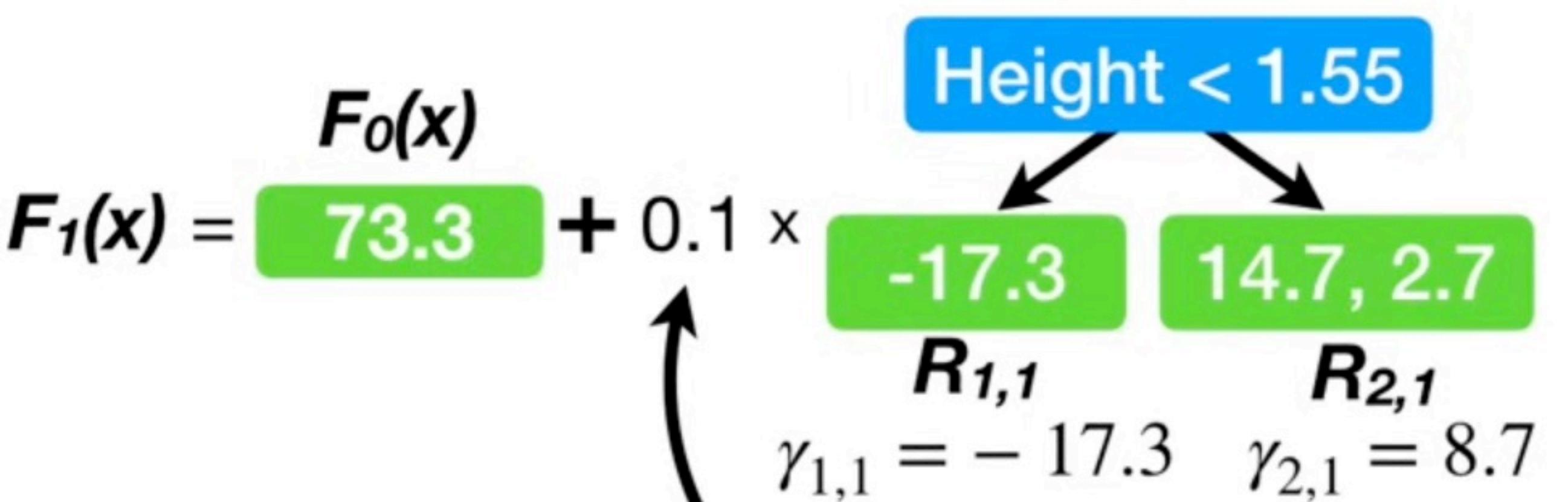
The last thing in this equation is  
this Greek character “nu”.


$$\mathbf{(D)} \text{ Update } F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$$

$$F_0(x) + F_1(x) = 73.3 + \begin{cases} -17.3 & \text{Height} < 1.55 \\ 14.7, 2.7 & \end{cases}$$
$$R_{1,1} \quad R_{2,1}$$
$$\gamma_{1,1} = -17.3 \quad \gamma_{2,1} = 8.7$$

A small **Learning Rate** reduces the effect each tree has on the final prediction, and this improves accuracy in the long run.

(D) Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$



In this example, we'll set **nu** to **0.1**.

**(D) Update**  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

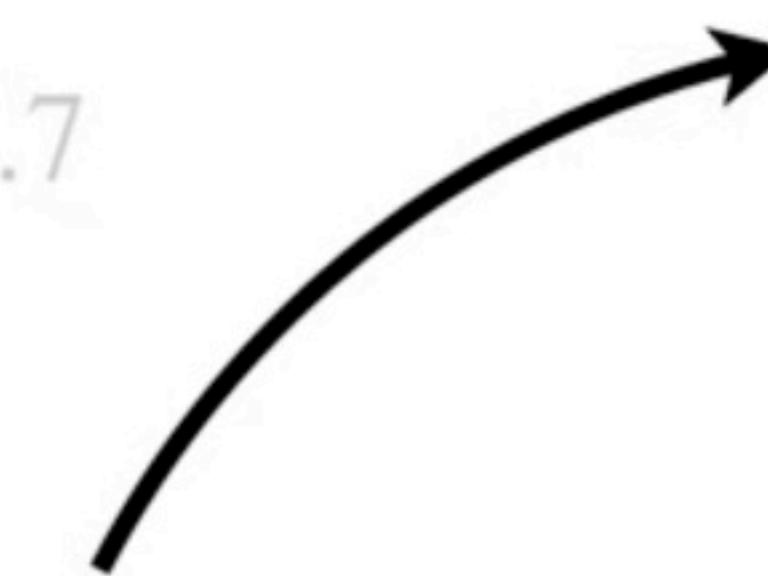
$$F_0(x)$$

$$F_1(x) = 73.3 + 0.1 \times$$

$R_{1,1}$        $R_{2,1}$

$\gamma_{1,1} = -17.3$      $\gamma_{2,1} = 8.7$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56



Now we will use  $F_1(x)$  to make new **Predictions** for each sample.

**(D) Update**  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$



$F_0(x)$

$F_1(x) = 73.3 + 0.1 \times$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

$\gamma_{1,1} = -17.3$

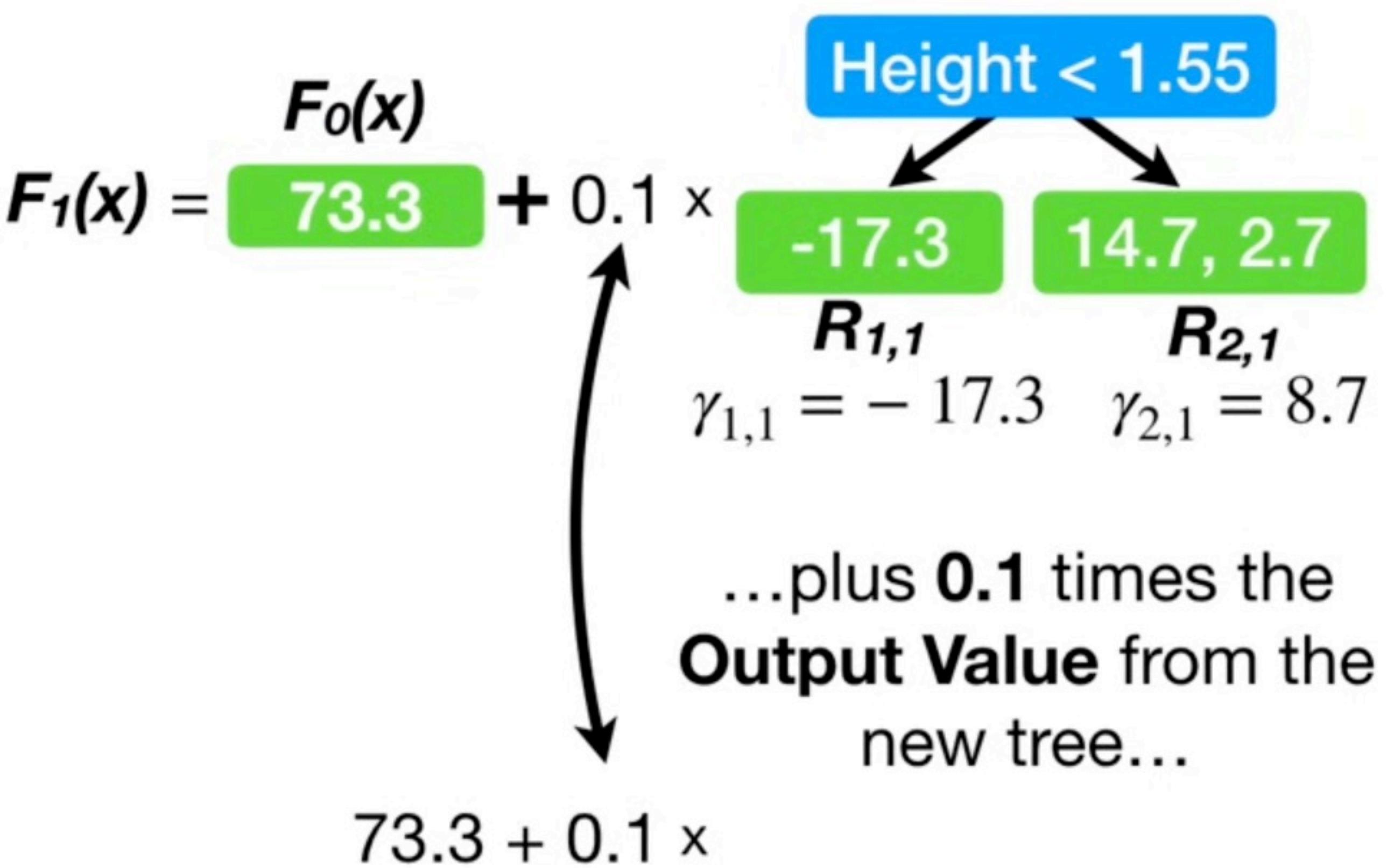
$\gamma_{2,1} = 8.7$

The new **Prediction** for  $x_1$   
starts with the last **Prediction**,  
 $F_0(x)$ , which is 73.3...

$73.3 +$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

(D) Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$



(D) Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

$$F_0(x)$$

$$F_1(x) = 73.3 + 0.1 \times$$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

$\gamma_{1,1} = -17.3$

$\gamma_{2,1} = 8.7$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

...which is **8.7** because  
 $x_1$ 's **Height** is  $> 1.55$ .

$$73.3 + 0.1 \times 8.7$$

(D) Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

$$F_0(x)$$

$$F_1(x) = 73.3 + 0.1 \times$$

Height < 1.55

-17.3

14.7, 2.7

$R_{1,1}$

$R_{2,1}$

$\gamma_{1,1} = -17.3$

$\gamma_{2,1} = 8.7$

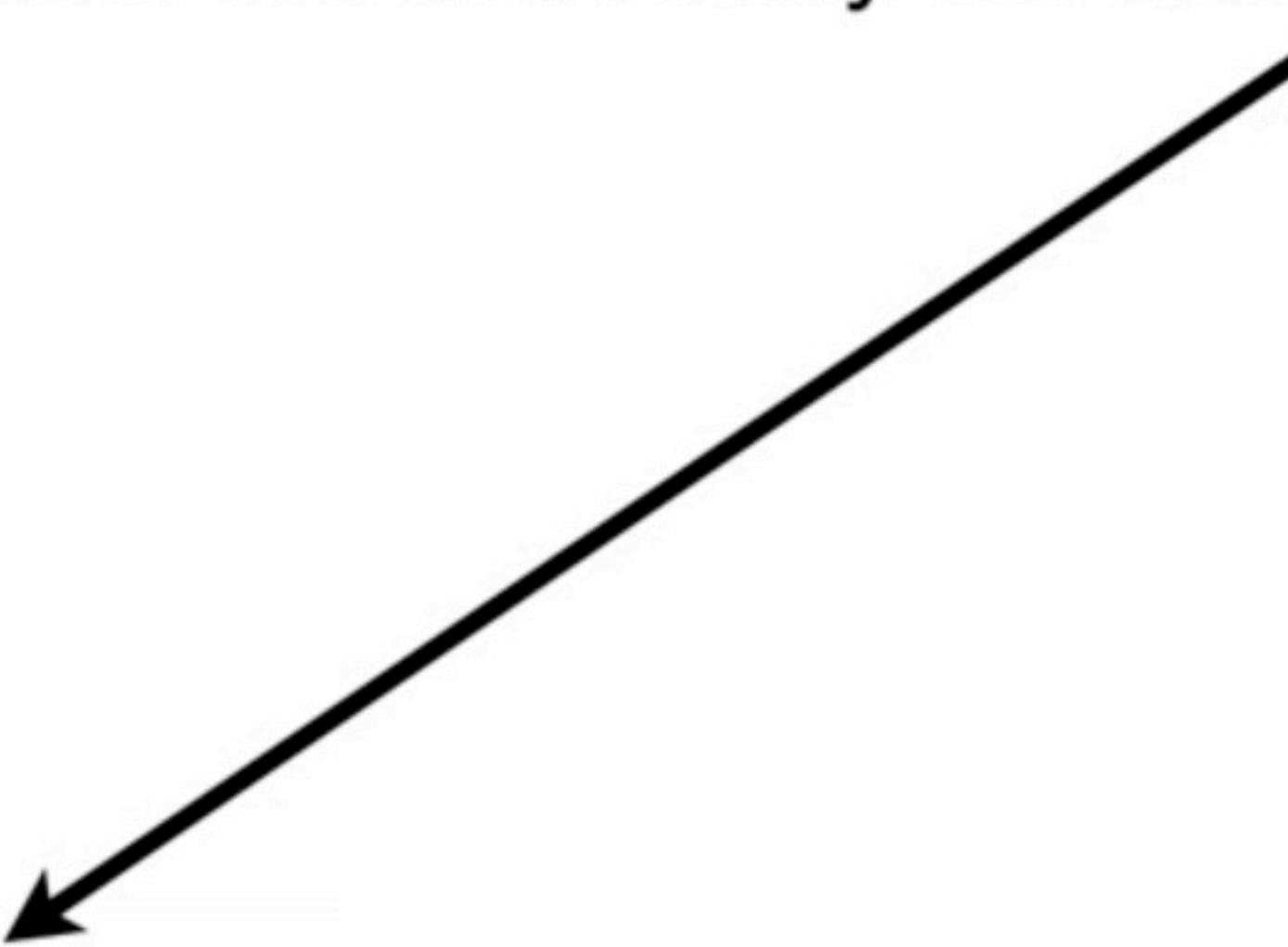
$$73.3 + 0.1 \times 8.7 = 74.2$$

The new **Prediction** for the first sample is **74.2**, which is slightly closer to the **Observed Weight, 88**, than the first **Prediction, 73.3**.

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56

(D) Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

Now we are ready for **Gradient Boost's** 3rd and final step!



**Step 3:** Output  $F_M(x)$

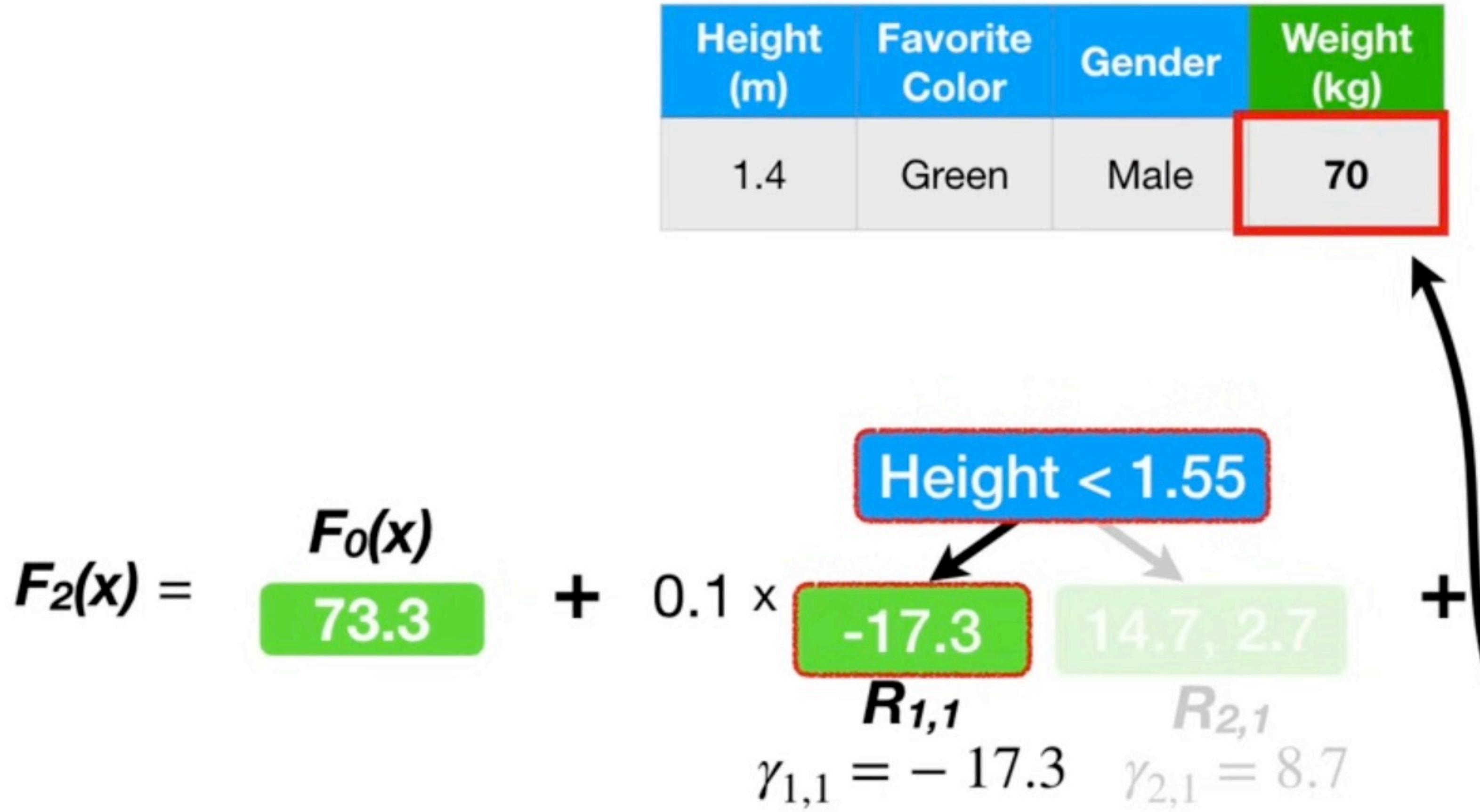
If  $M = 2$ , then  $F_2(x)$  is the output from the **Gradient Boost** algorithm.

$$F_2(x) = F_0(x) + 0.1 \times \begin{array}{c} \text{Height} < 1.55 \\ \downarrow \quad \downarrow \\ -17.3 \quad 14.7, 2.7 \\ R_{1,1} \quad R_{2,1} \\ \gamma_{1,1} = -17.3 \quad \gamma_{2,1} = 8.7 \end{array} + 0.1 \times \begin{array}{c} \text{Height} < 1.55 \\ \downarrow \quad \downarrow \\ -15.6 \quad 13.8, 1.8 \\ R_{1,2} \quad R_{2,2} \\ \gamma_{1,2} = -15.6 \quad \gamma_{2,2} = 7.8 \end{array}$$

**Step 3:** Output  $F_M(x)$

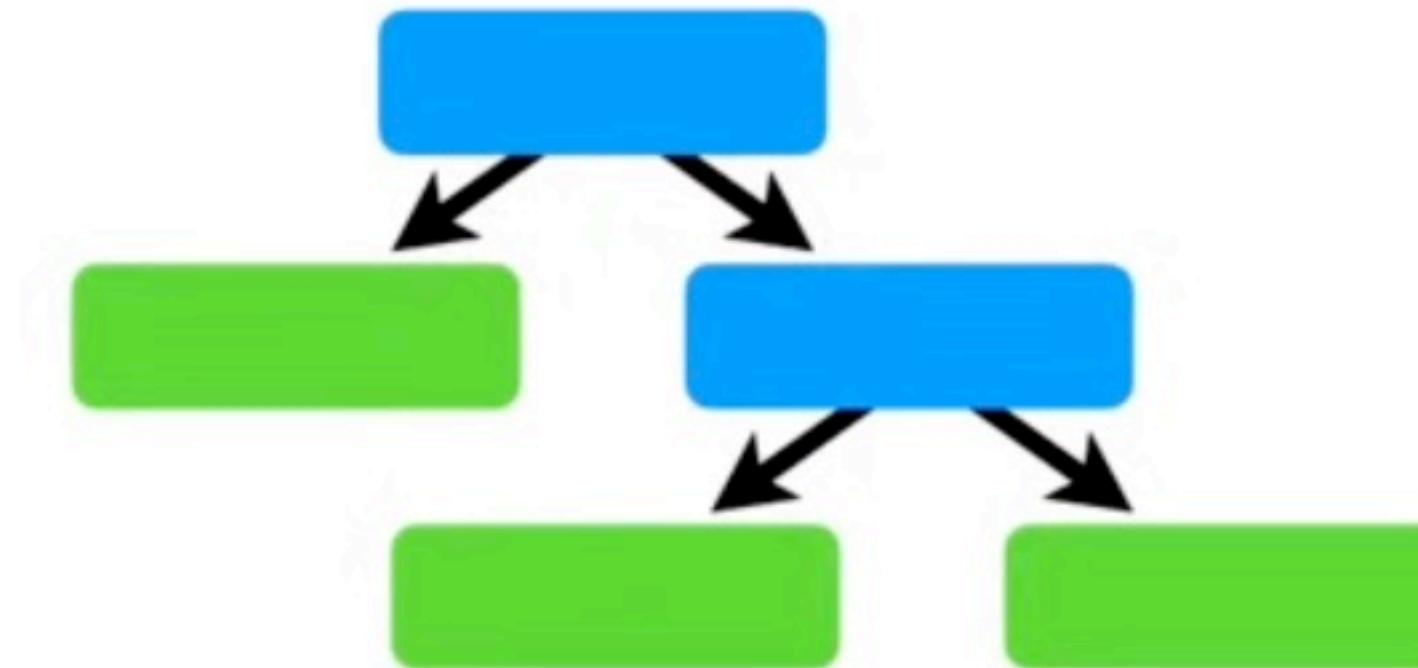
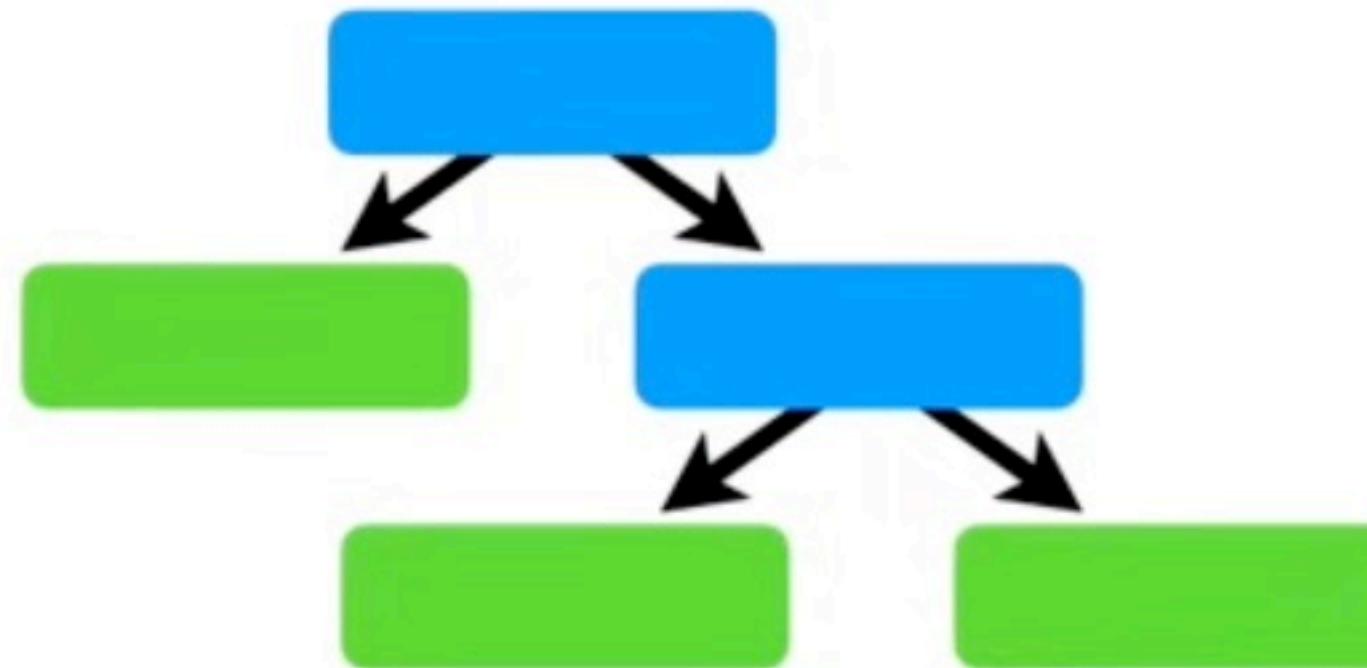
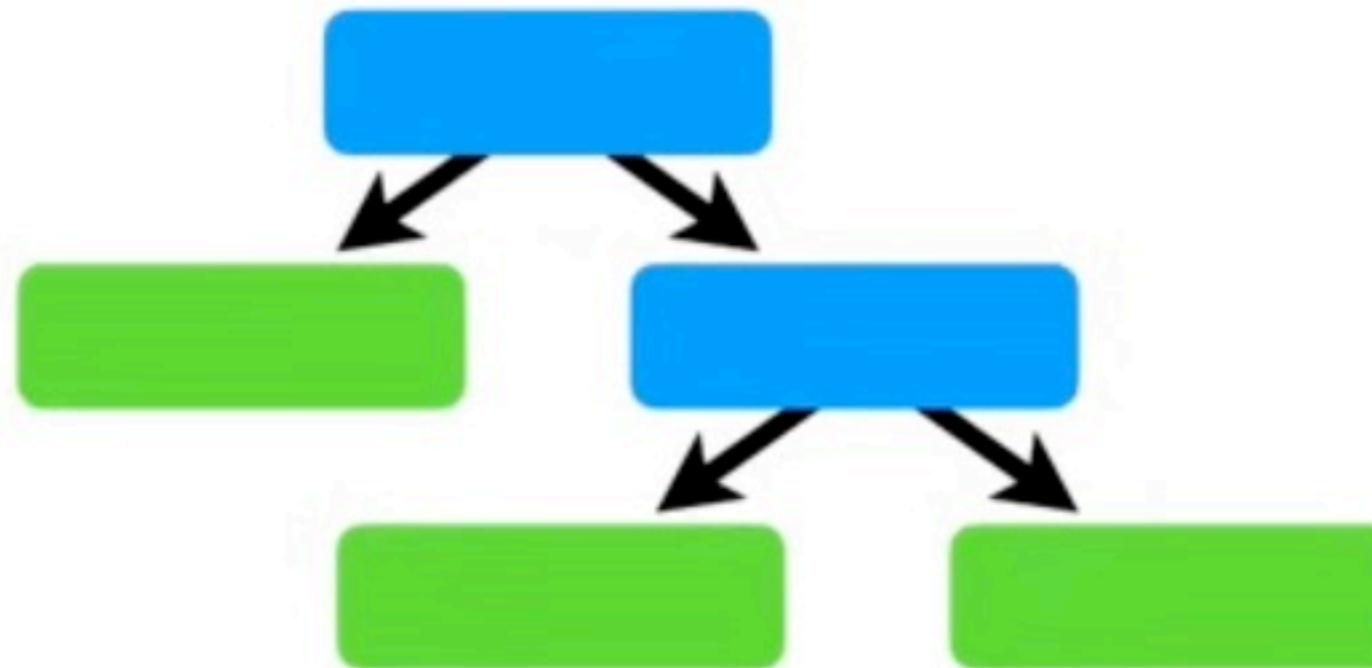
## Gradient Boost

Predicts that this person  
**Weighs 70 kg.**



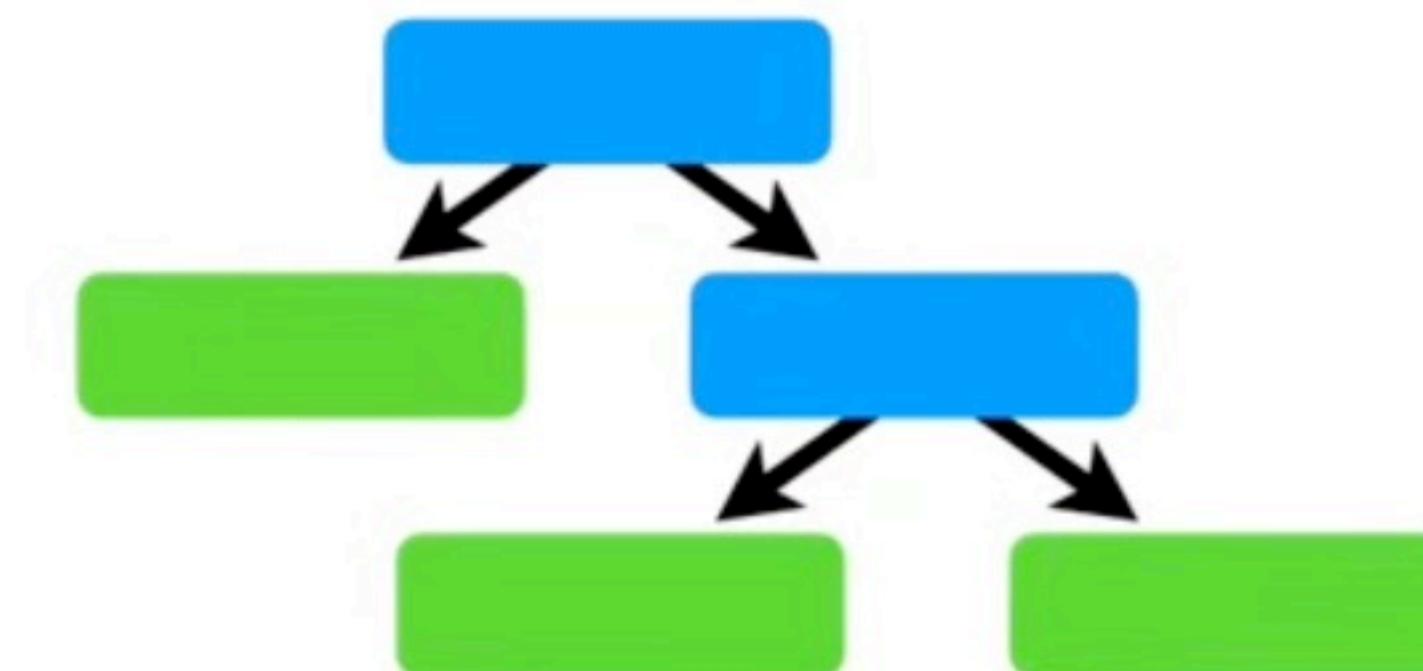
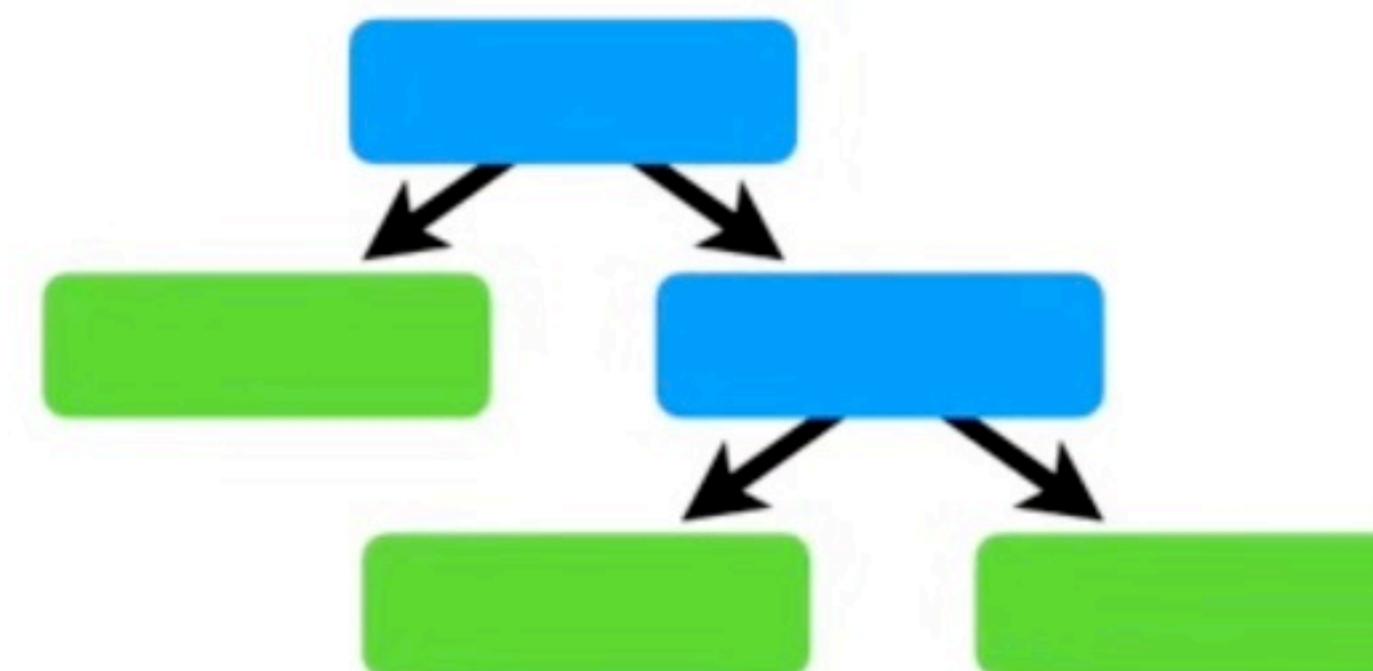
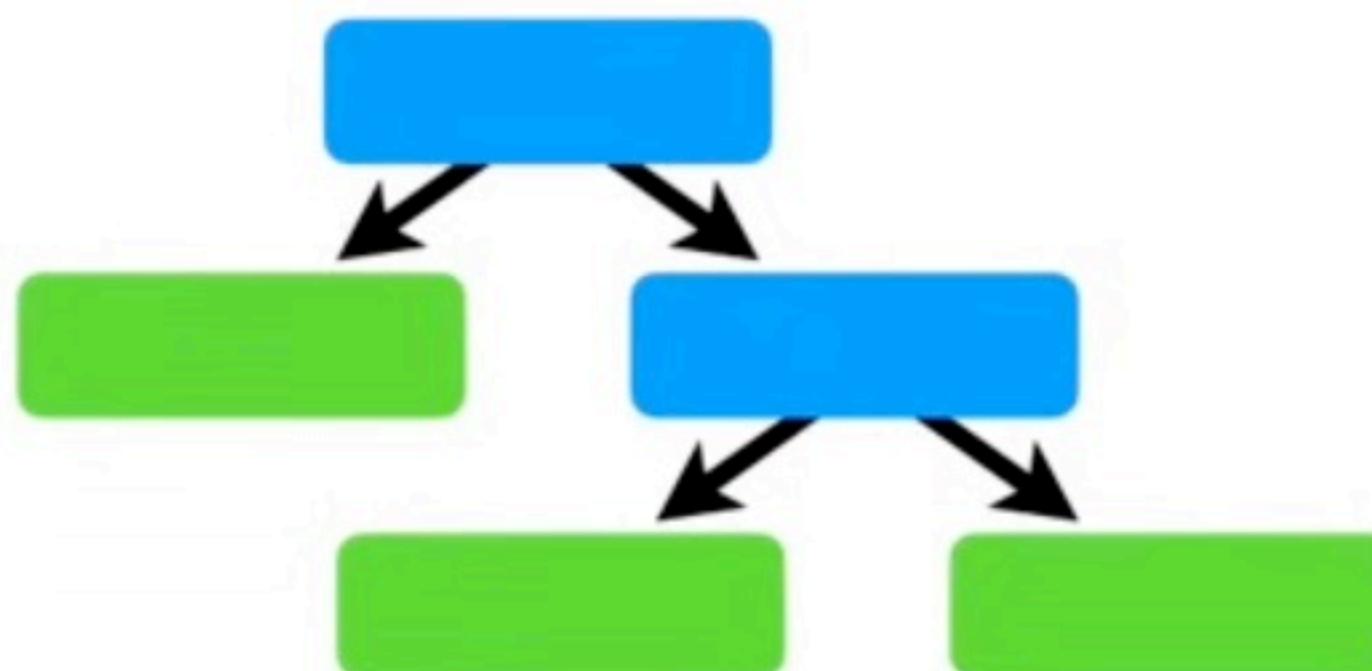
The Predicted Weight =  $73.3 + (0.1 \times -17.3) + (0.1 \times -15.6) = 70$

**NOTE:** Before we go, I want to remind you  
that **Gradient Boost** usually uses trees larger  
than stumps.



**NOTE:** Before we go, I want to remind you  
that **Gradient Boost** usually uses trees larger  
than stumps.

I only used stumps in this tutorial because our  
**Training Dataset** was so darn small.



## ★ SUPER INTUITIVE ANALOGY (you'll love this)

Think of Gradient Boosting as a shooter trying to hit a target:

- First shot misses → that's your initial model
- Second shot is made by looking at **where you missed** → **direction of correction**
- Third shot corrects remaining error
- Each new shot corrects remaining error (gradient)

The "direction of correction" is the **gradient**.

Each correction is **learned by a decision tree**.

## ★ FINAL ANSWER (interview-quality)

"We compute the gradient because it tells us the direction in which predictions must change to reduce the loss.

A decision tree is then trained to approximate that gradient — meaning the tree learns how much correction is needed in different regions of feature space.

Adding this tree to the model moves predictions in the direction that reduces loss, exactly like gradient descent but in function space."