

# BOOSTING

- AdaBoost
- Gradient Boost
- XGBoost

## GRADIENT BOOST

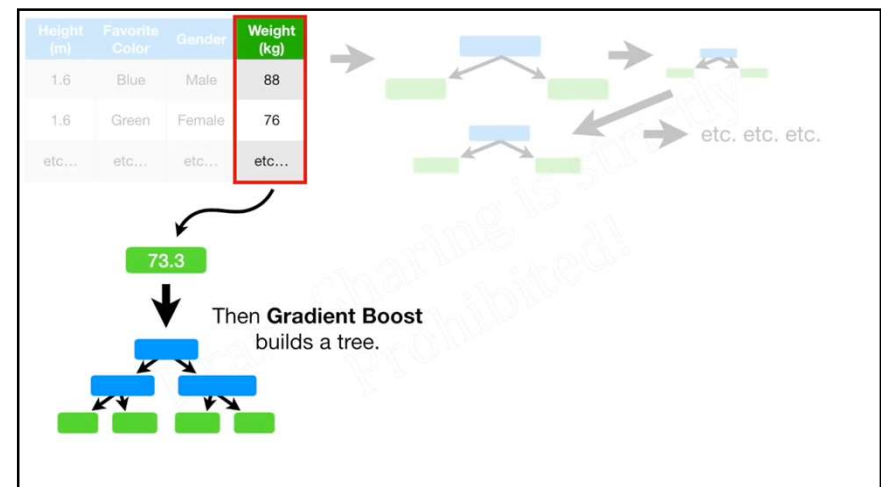
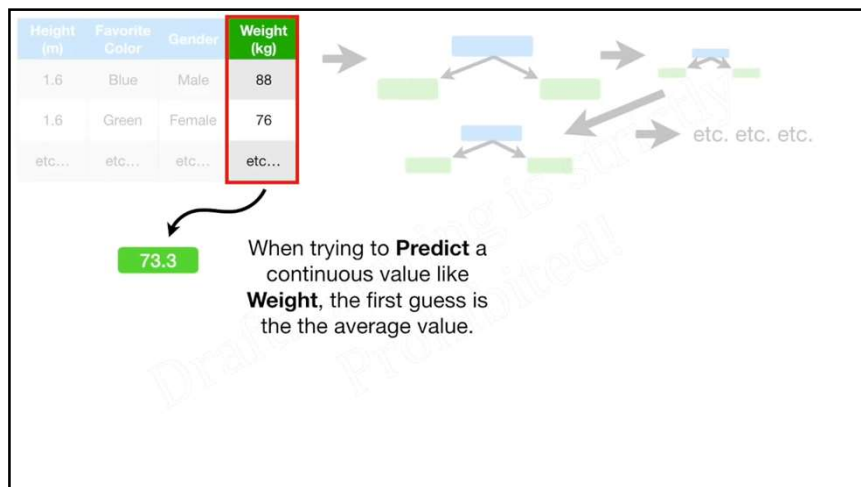
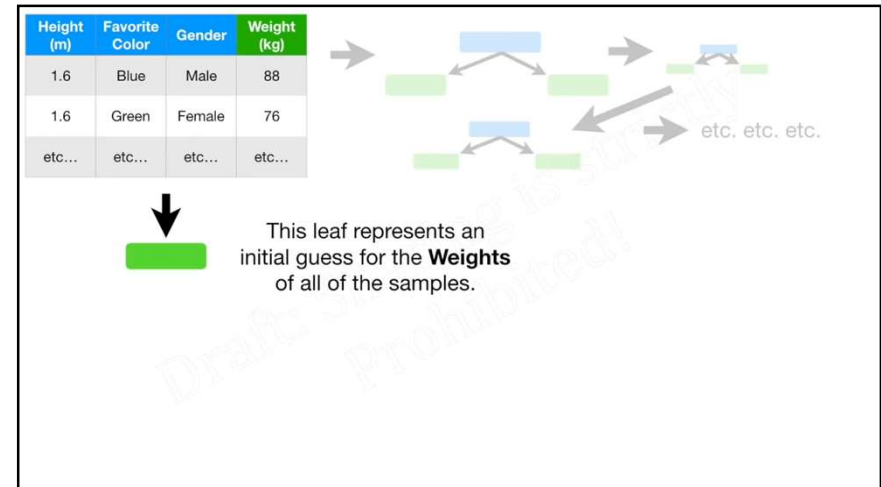
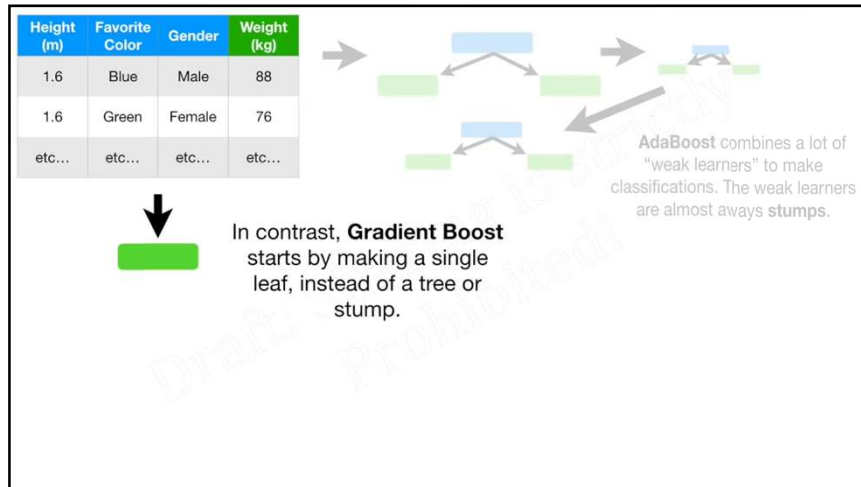
Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
etc...	etc...	etc...	etc...

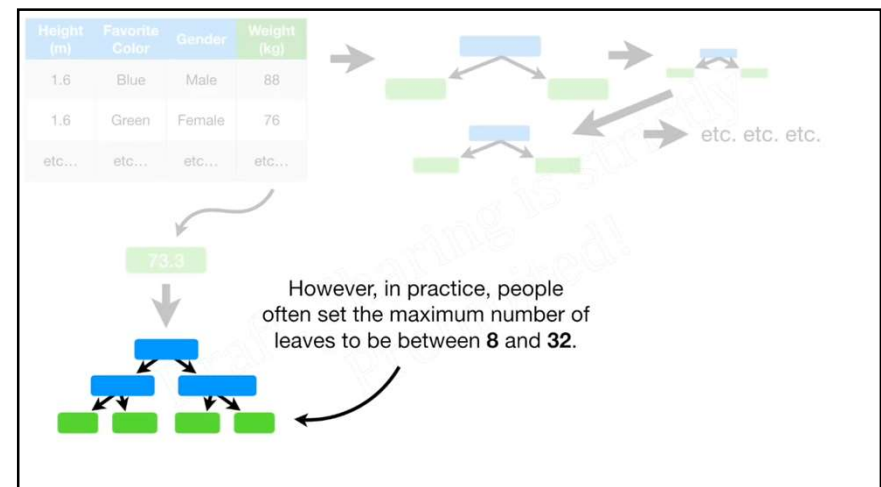
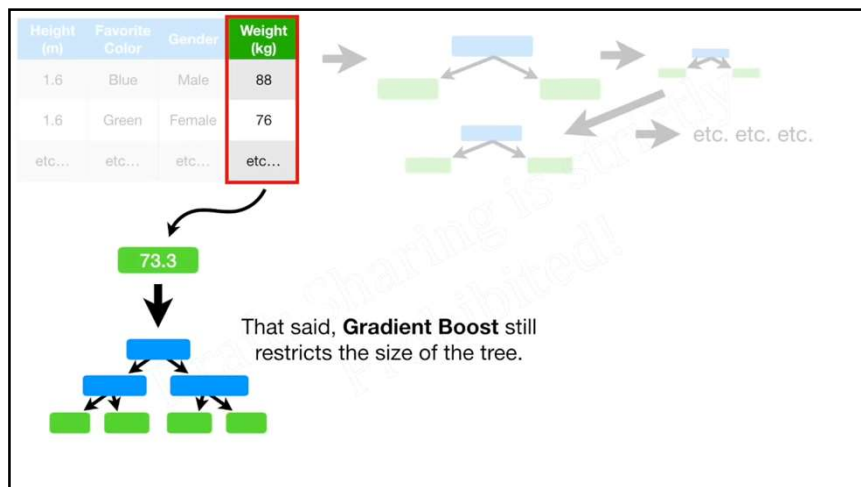
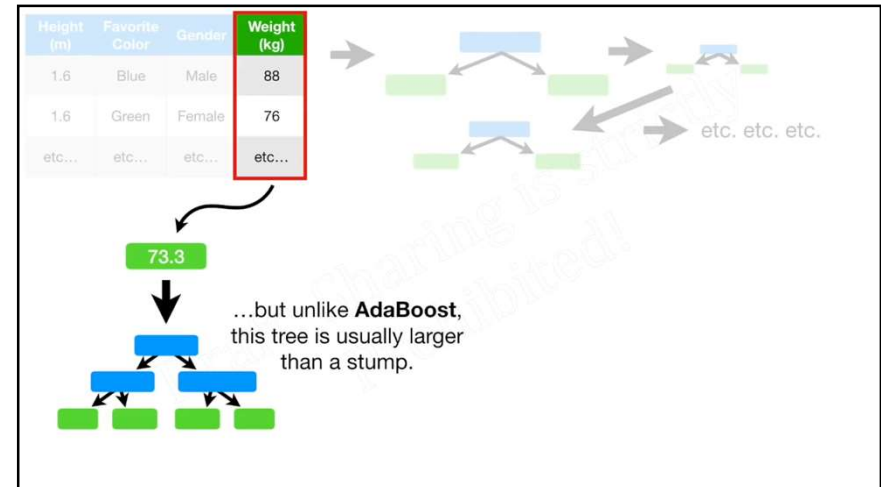
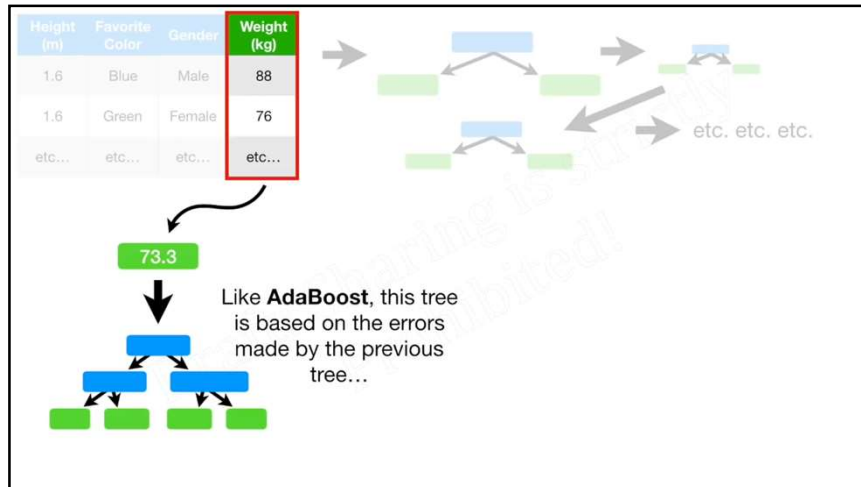
← we want to use these measurements to  
**Predict Weight...**

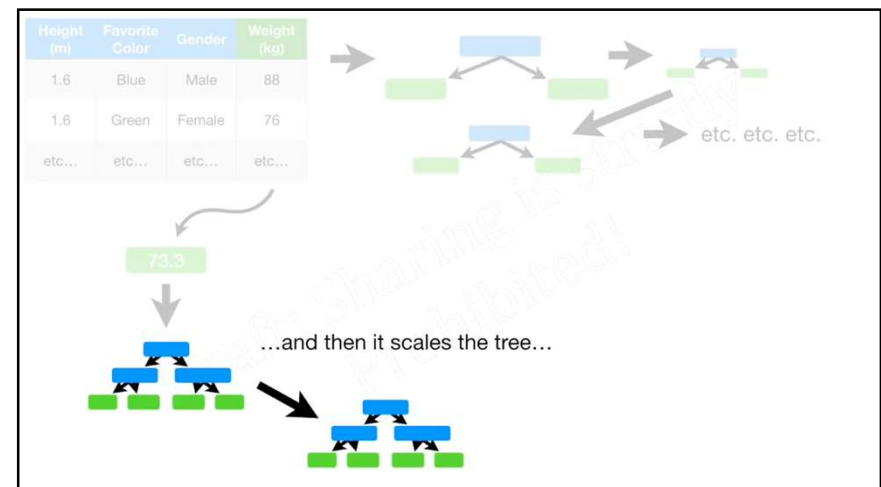
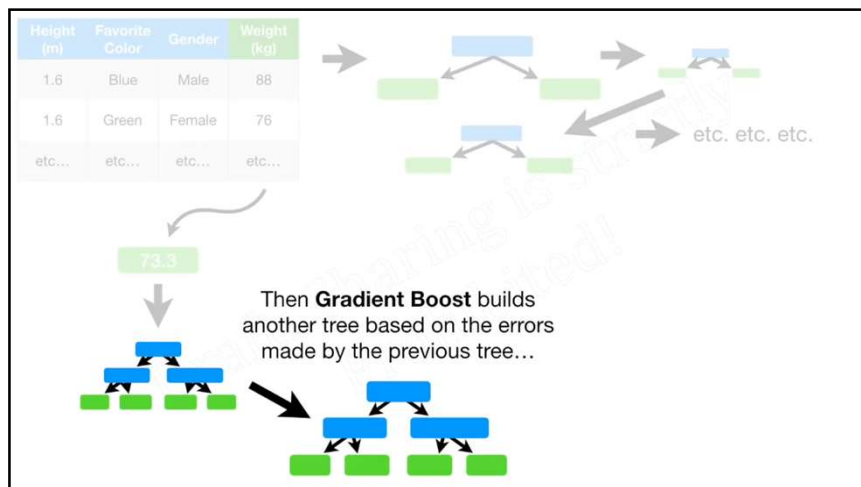
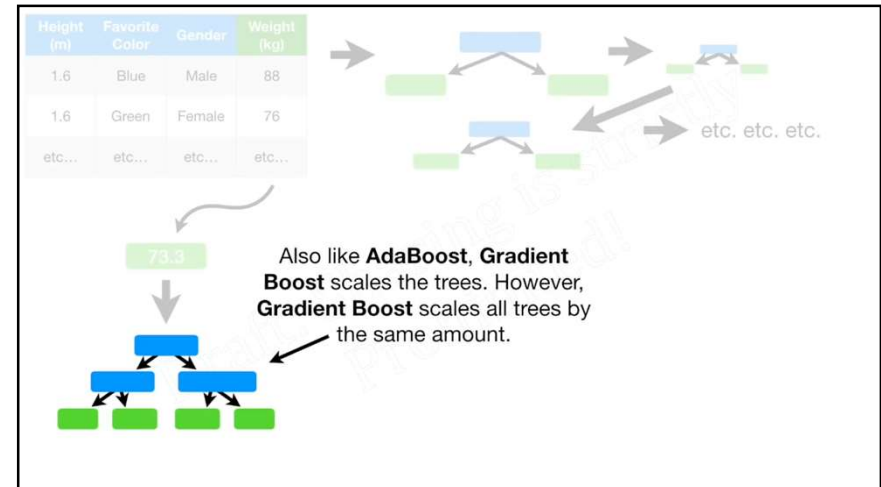
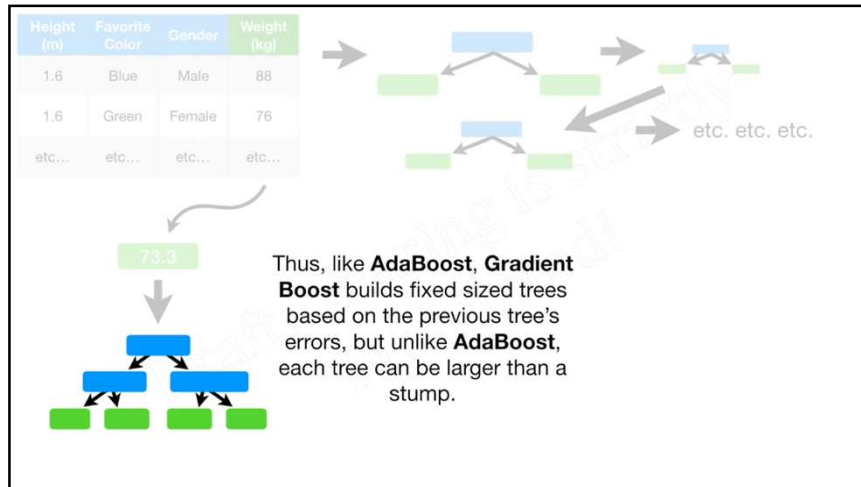
Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	Male	73
1.5	Green	Male	77
1.4	Blue	Female	57

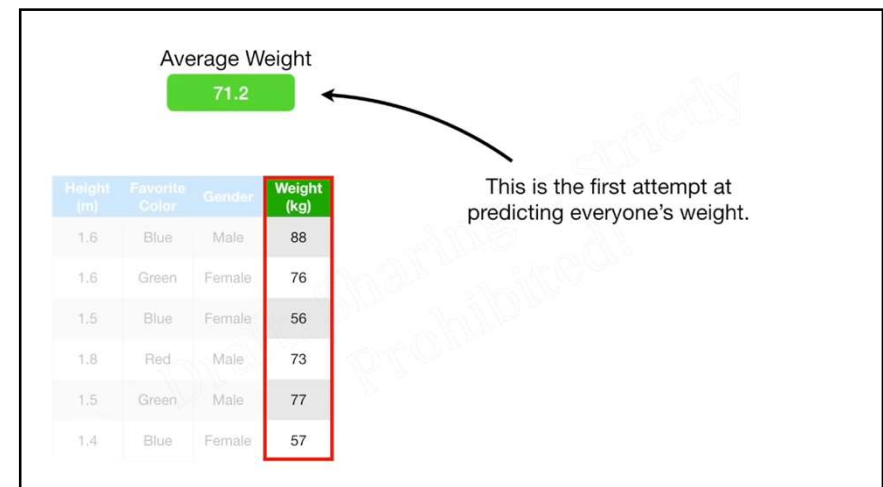
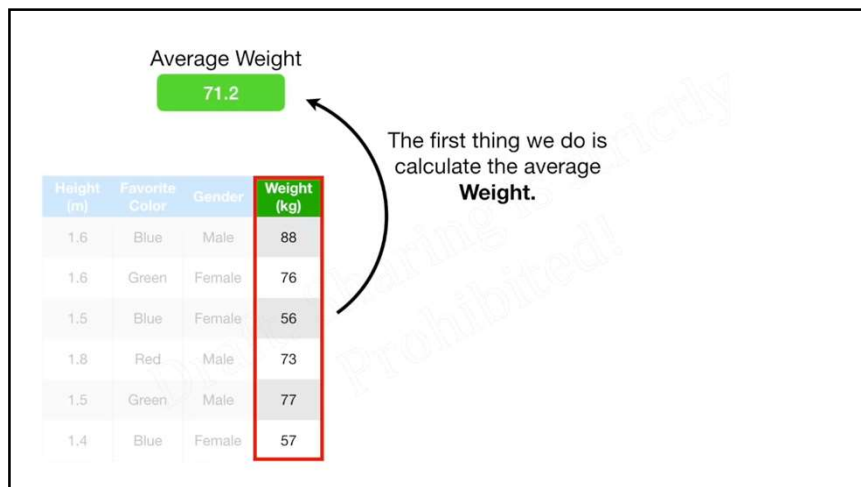
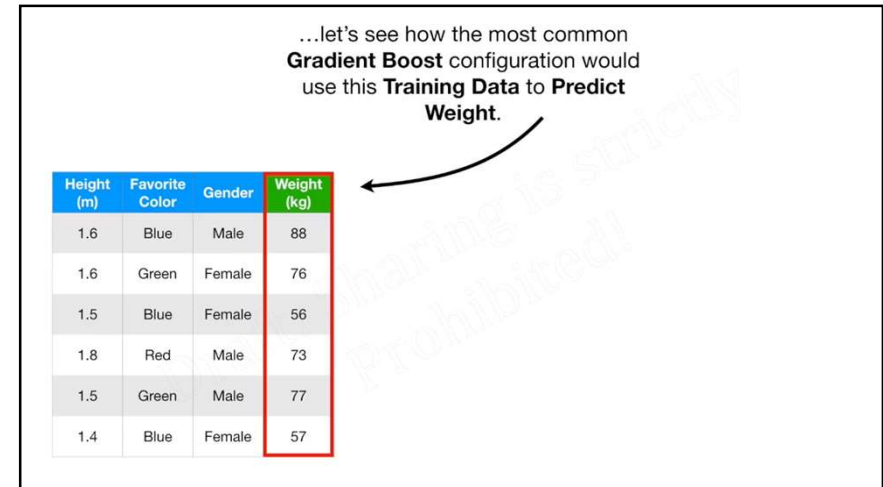
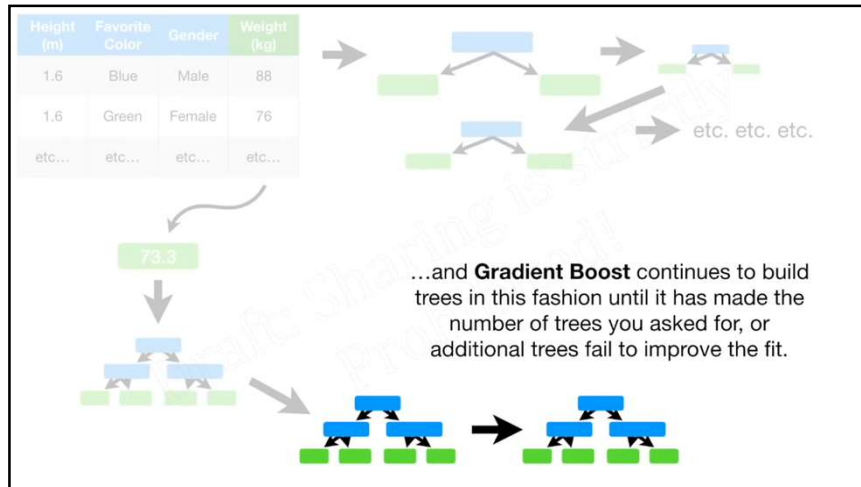
**NOTE:** When **Gradient Boost** is used to **Predict** a continuous value, like **Weight**, we say that we are using **Gradient Boost for Regression**.

Using **Gradient Boost for Regression** is different from doing **Linear Regression**, so while the two methods are related, don't get them confused with each other.









Average Weight

71.2

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	Male	73
1.5	Green	Male	77
1.4	Blue	Female	57

In other words, if we stopped right now, we would predict that everyone **Weighed 71.2 kg**.

However, **Gradient Boost** doesn't stop here.

Average Weight

71.2

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	Male	73
1.5	Green	Male	77
1.4	Blue	Female	57

The next thing we do is build a tree based on the errors from the first tree.

Average Weight

71.2

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	Male	73
1.5	Green	Male	77
1.4	Blue	Female	57

The errors that the previous tree made are the differences between the **Observed Weights** and the **Predicted Weight, 71.2**.

(Observed Weight - Predicted Weight)

Average Weight

71.2

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	Male	73
1.5	Green	Male	77
1.4	Blue	Female	57

So let's start by plugging in **71.2** for the **Predicted Weight**...

(Observed Weight - Predicted Weight)

Average Weight

71.2

Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	16.8
1.6	Green	Female	76	
1.5	Blue	Female	56	
1.8	Red	Male	73	
1.5	Green	Male	77	
1.4	Blue	Female	57	

...and save the difference, which is called a **Pseudo Residual**, in a new column.

$$(88 - 71.2) = 16.8$$

Average Weight

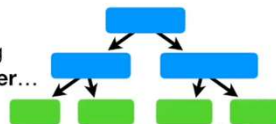
71.2

Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	16.8
1.6	Green	Female	76	4.8
1.5	Blue	Female	56	-15.2
1.8	Red	Male	73	1.8
1.5	Green	Male	77	5.8
1.4	Blue	Female	57	-14.2

Now do the same thing for the remaining **Weights**...

$$(57 - 71.2) = -14.2$$

Now we will build a **Tree**, using **Height, Favorite Color and Gender**...



Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	16.8
1.6	Green	Female	76	4.8
1.5	Blue	Female	56	-15.2
1.8	Red	Male	73	1.8
1.5	Green	Male	77	5.8
1.4	Blue	Female	57	-14.2

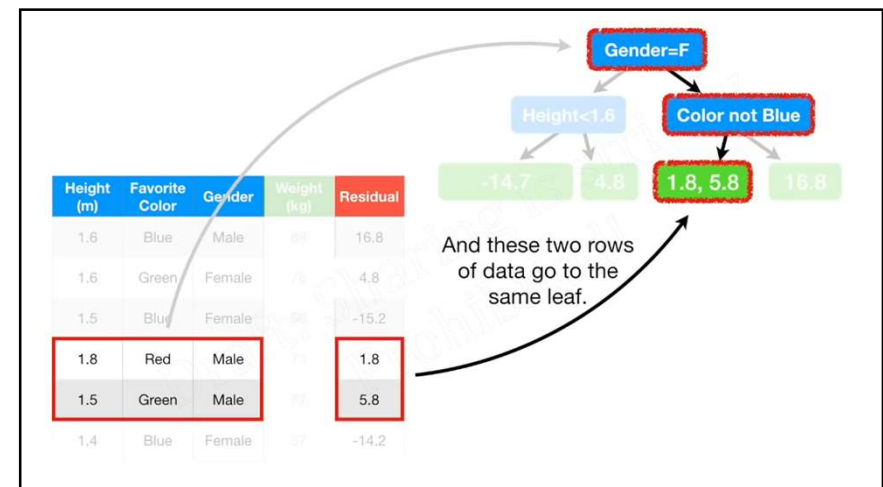
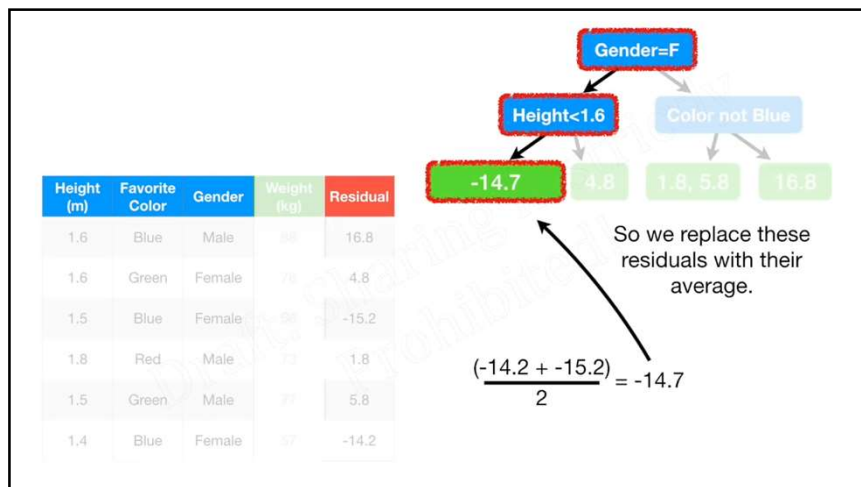
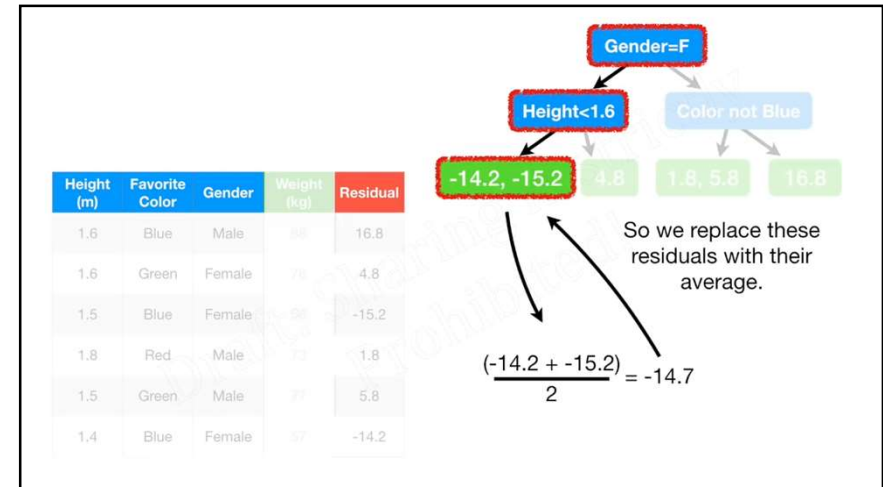
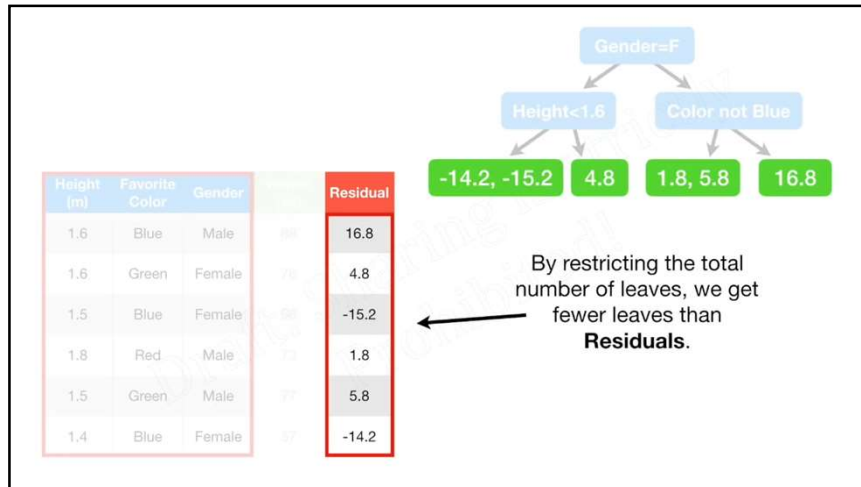
...to Predict the Residuals.



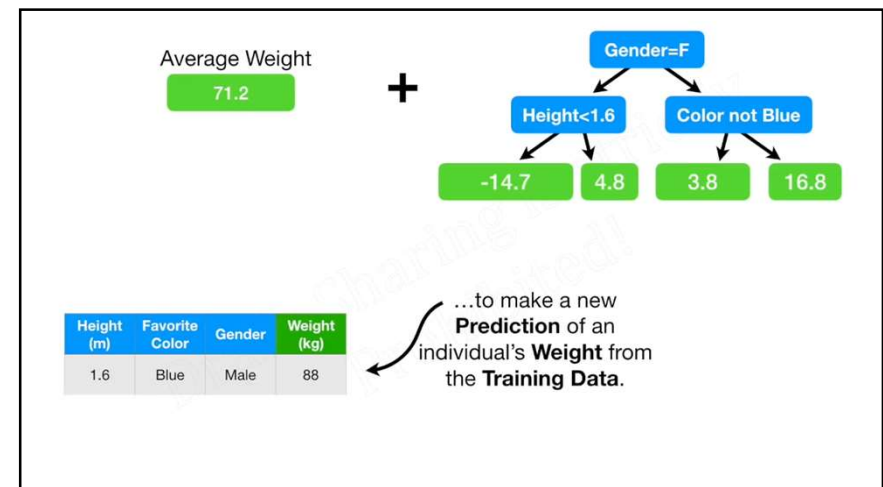
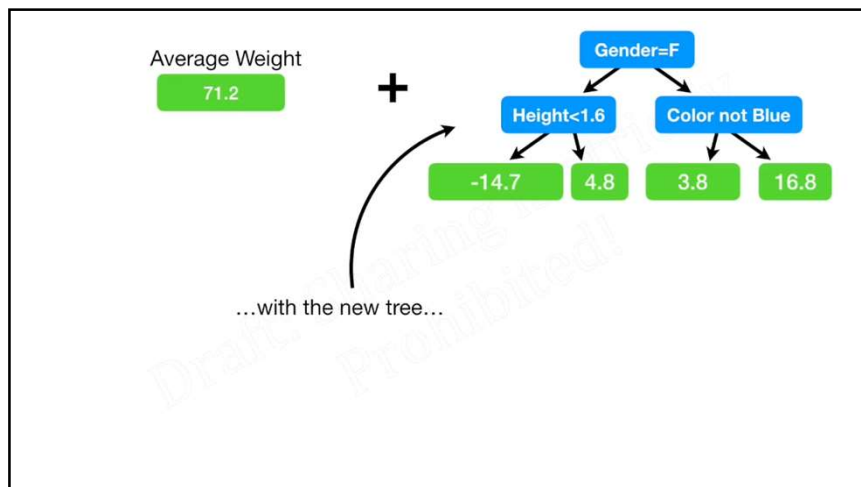
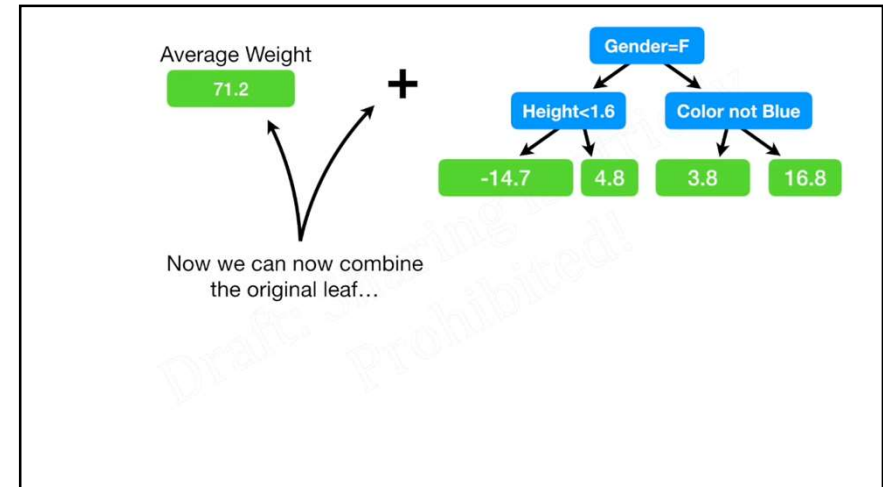
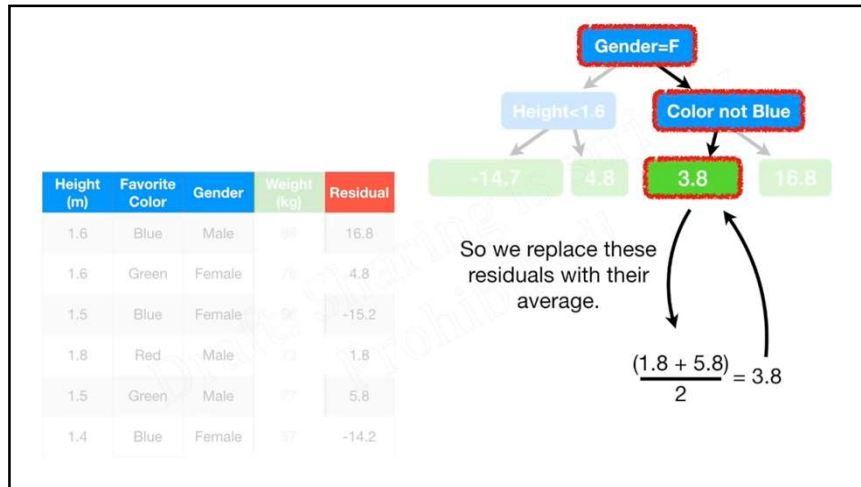
Remember, in this example we are only allowing up to four leaves...

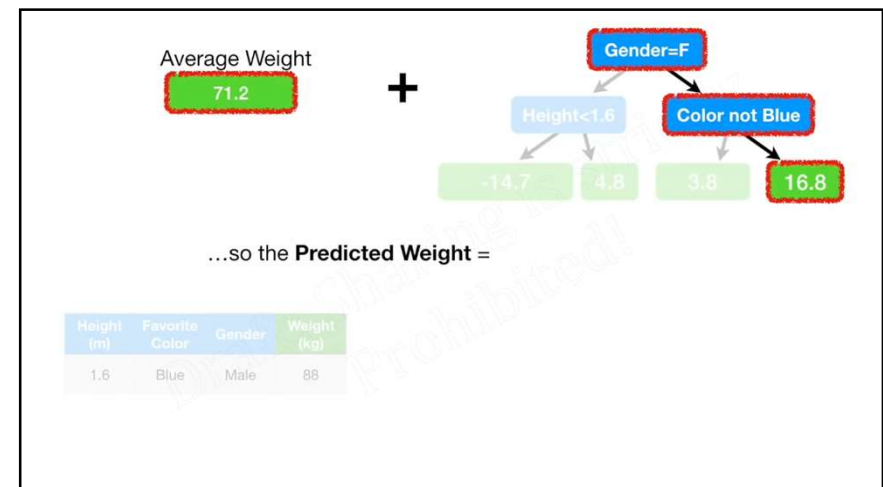
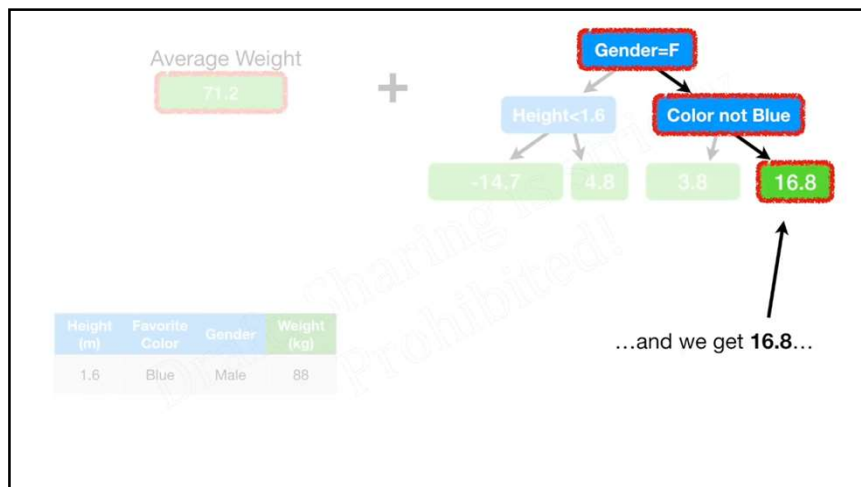
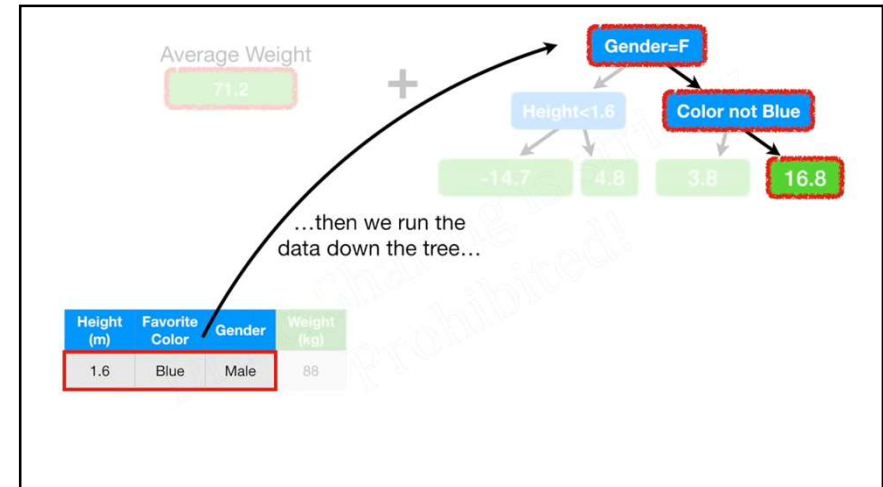
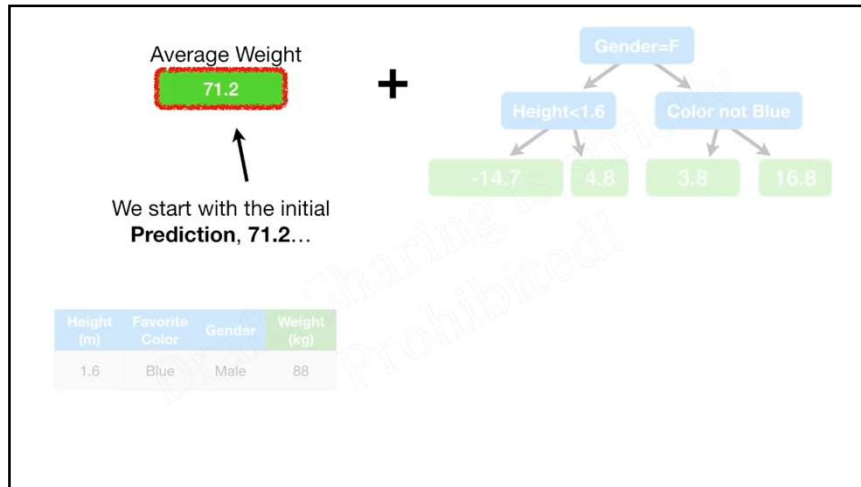
...but when using a larger dataset, it is common to allow anywhere from 8 to 32.

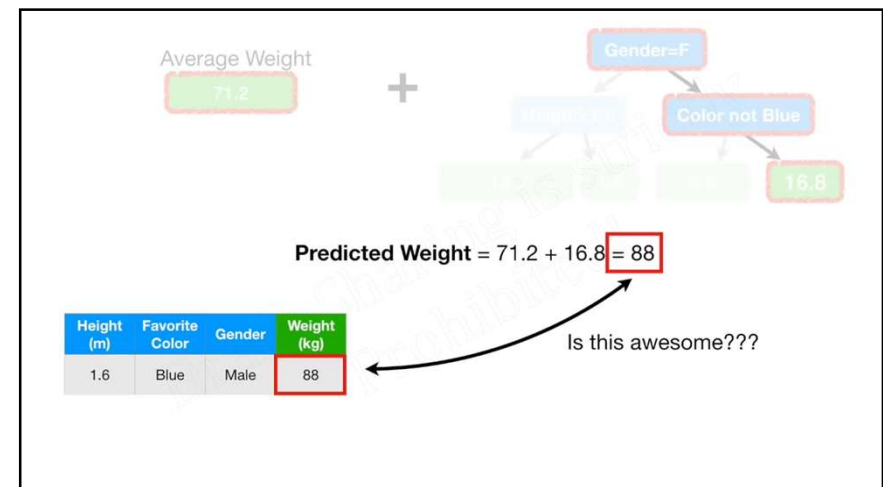
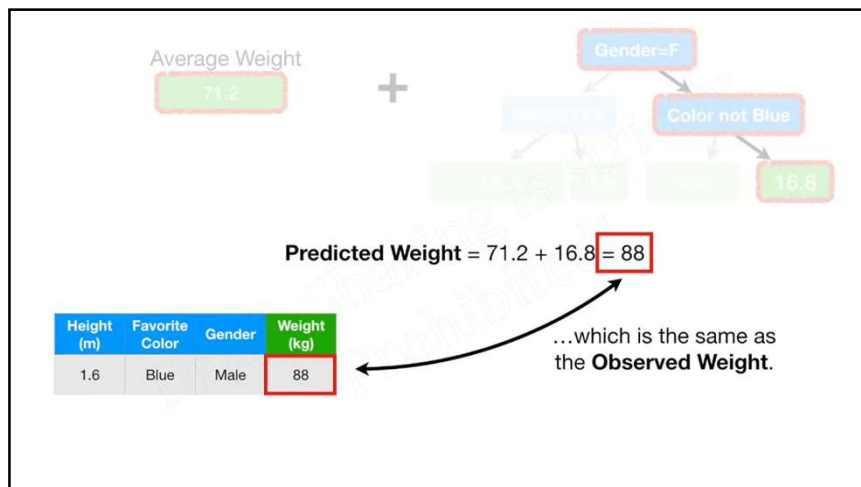
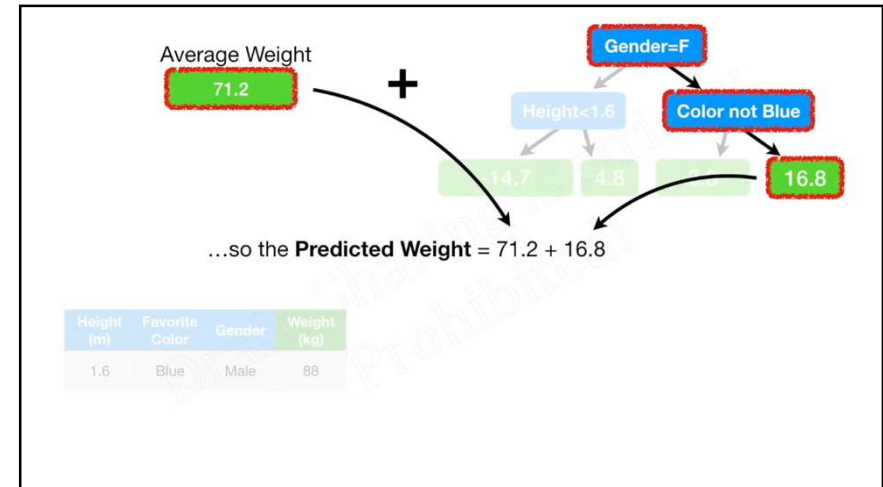
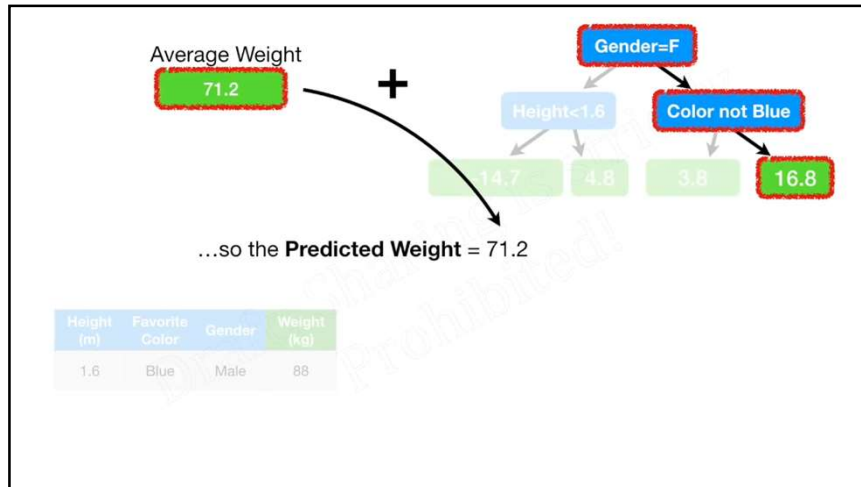


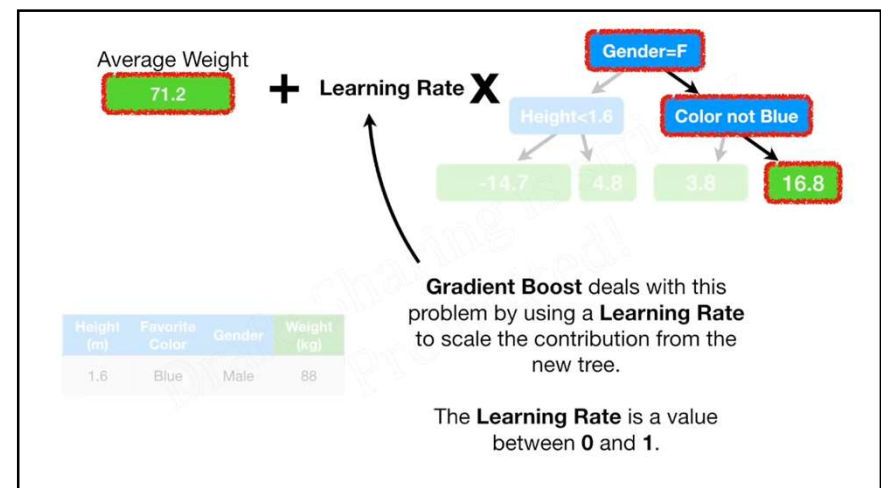
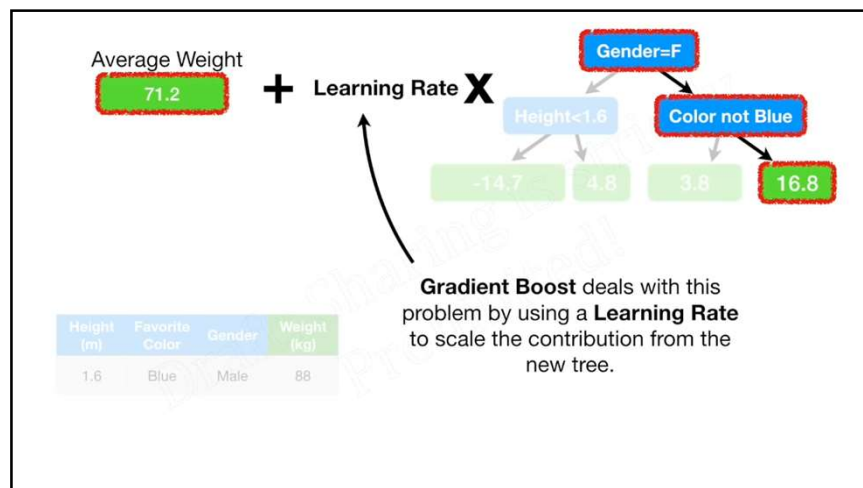
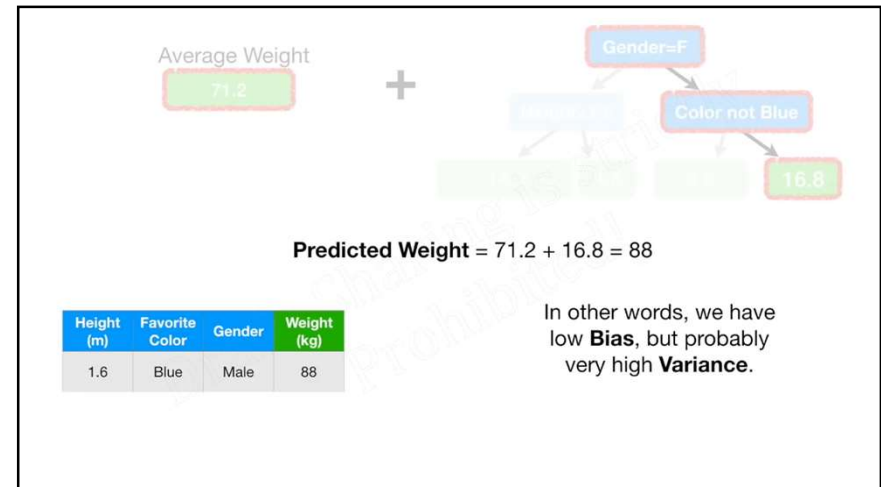
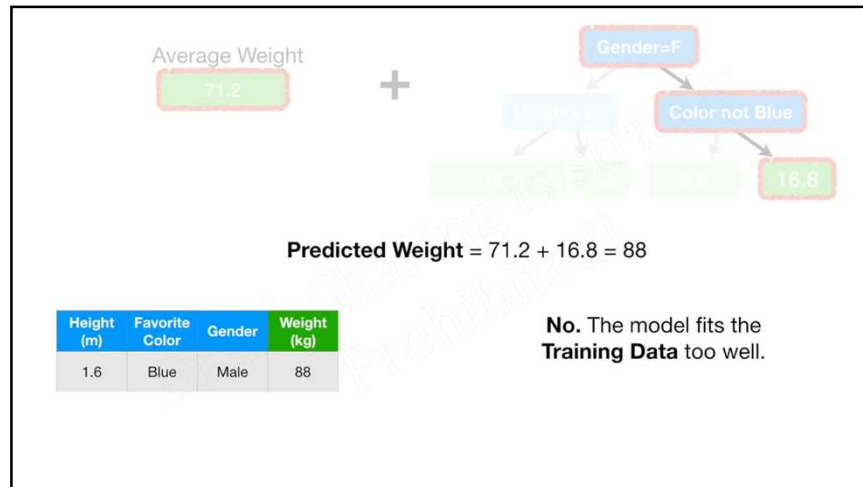


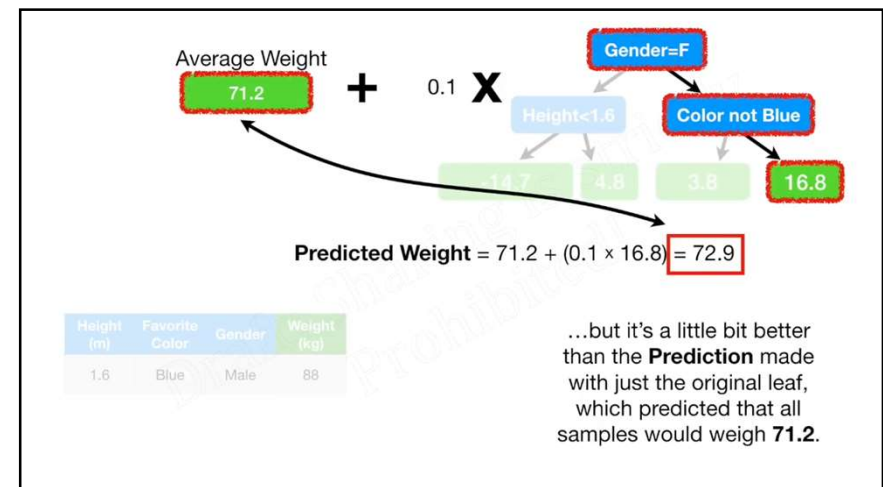
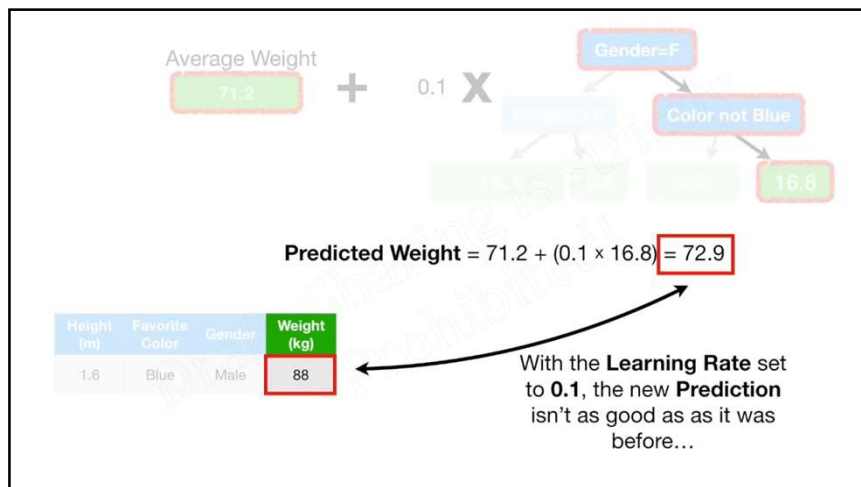
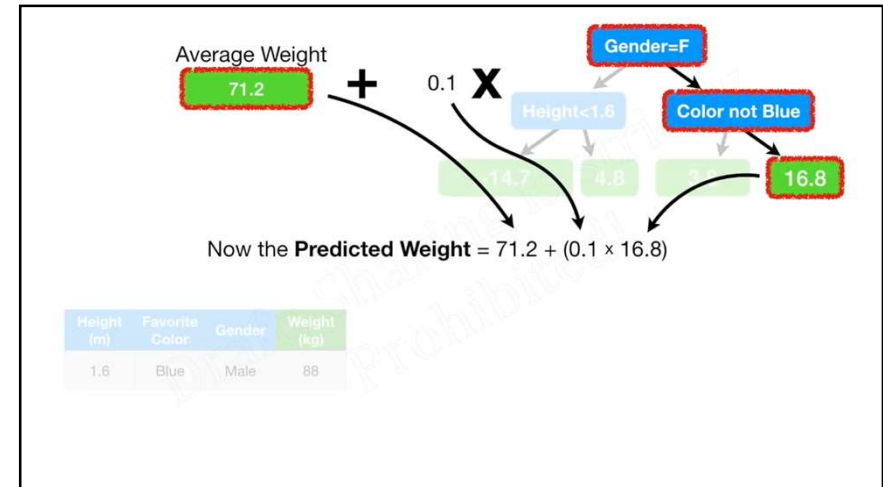
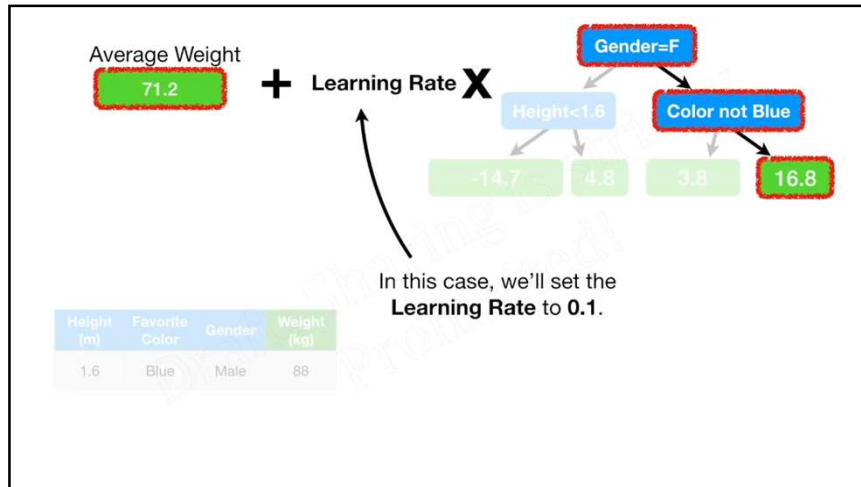














So let's build another tree so we can take another small step in the right direction.

Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	
1.6	Green	Female	76	
1.5	Blue	Female	56	
1.8	Red	Male	73	
1.5	Green	Male	77	
1.4	Blue	Female	57	

Just like before, we calculate the **Pseudo Residuals**, the difference between the **Observed Weights** and our latest **Predictions**.

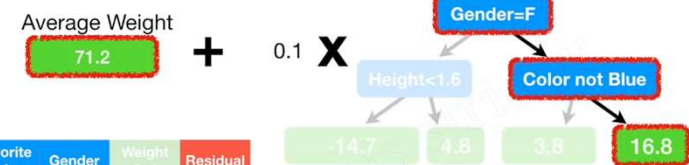
←  $\text{Residual} = (\text{Observed} - \text{Predicted})$



Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	
1.6	Green	Female	76	
1.5	Blue	Female	56	
1.8	Red	Male	73	
1.5	Green	Male	77	
1.4	Blue	Female	57	

$$\text{Residual} = (88 - (71.2 + 0.1 \times 16.8)) = 15.1$$

...and we get **15.1**...

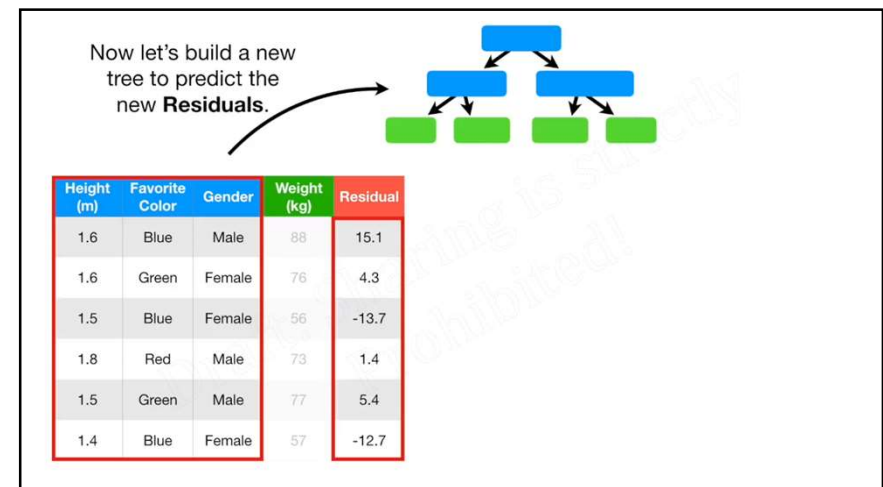
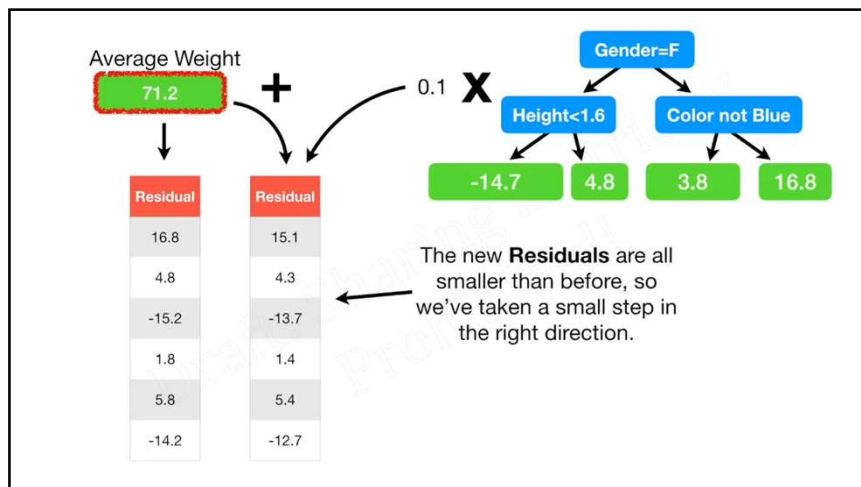
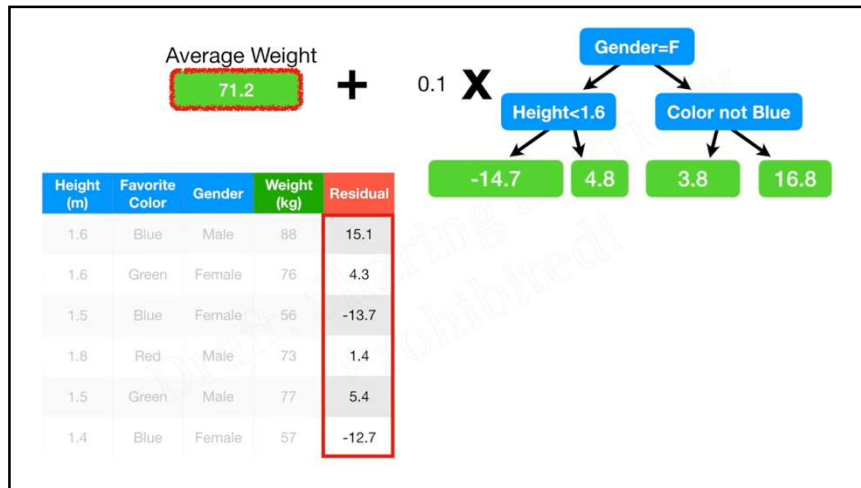


Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	15.1
1.6	Green	Female	76	
1.5	Blue	Female	56	
1.8	Red	Male	73	
1.5	Green	Male	77	
1.4	Blue	Female	57	

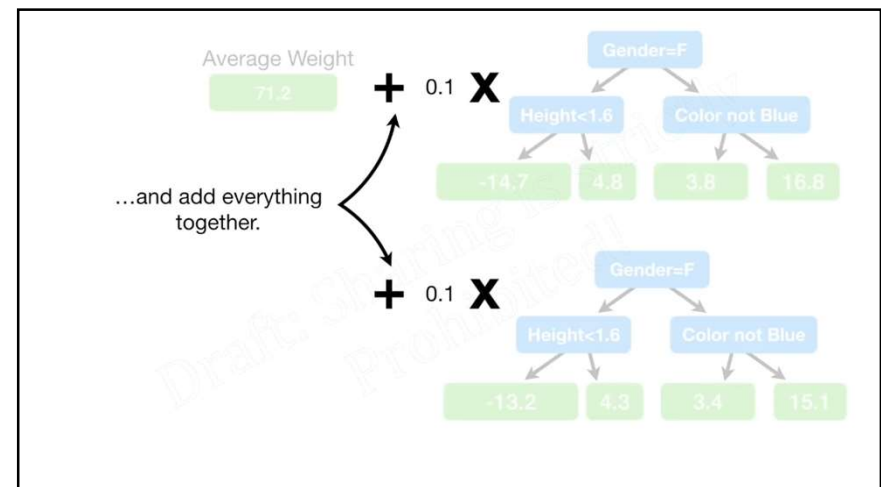
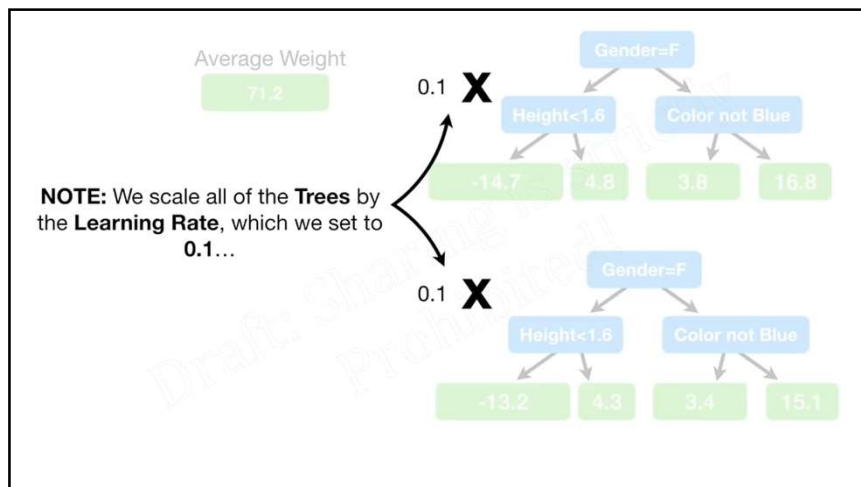
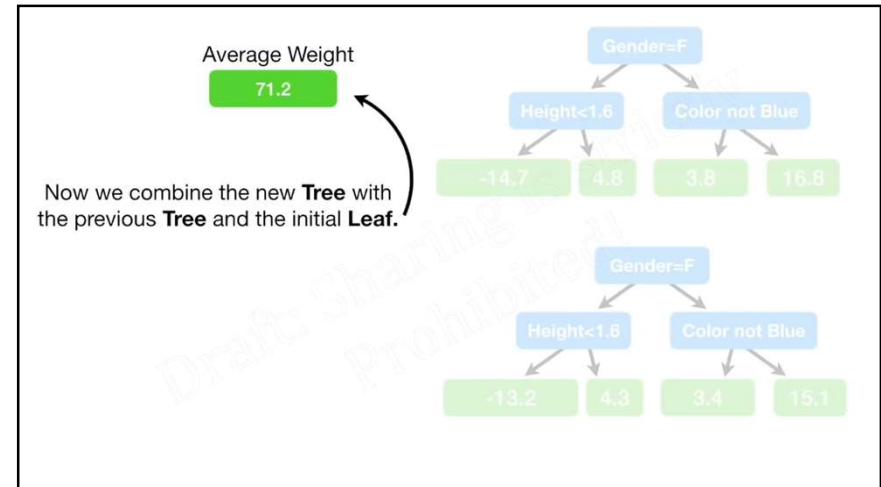
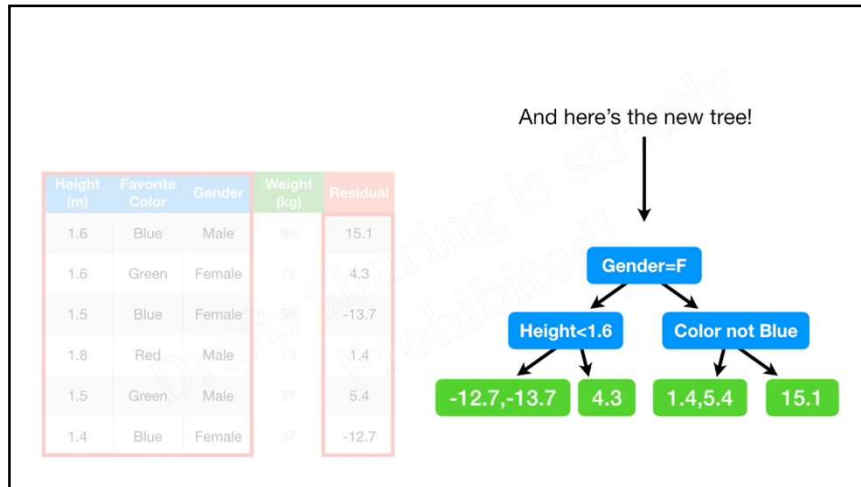
$$\text{Residual} = (88 - (71.2 + 0.1 \times 16.8)) = 15.1$$

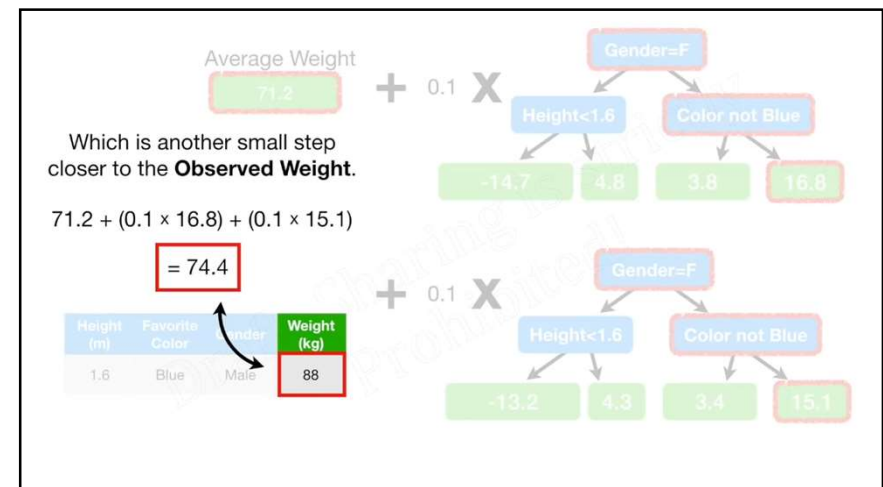
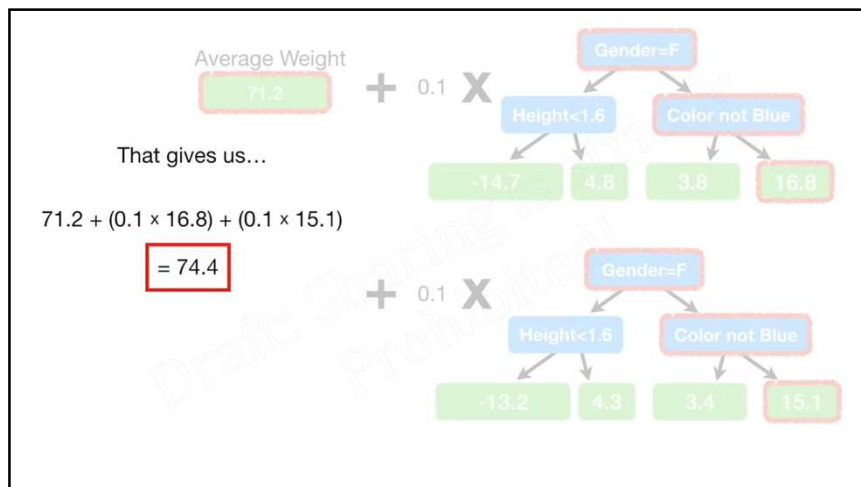
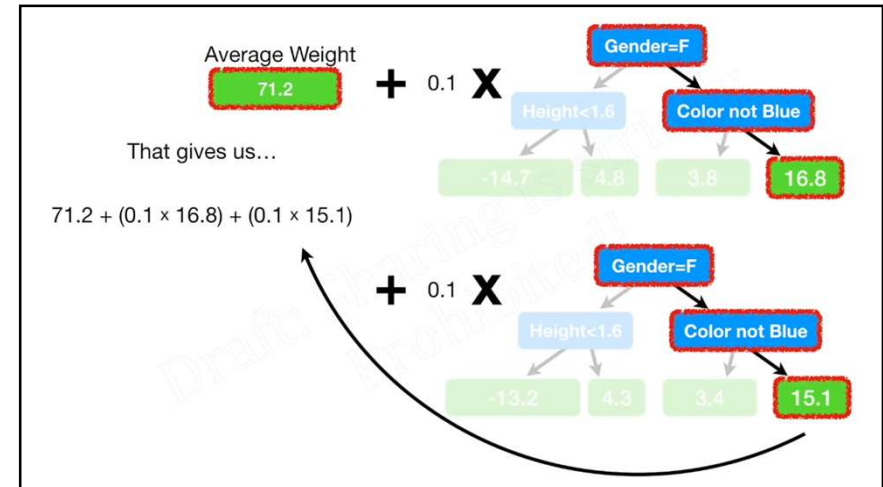
...and we save that in the column for **Pseudo Residuals**.

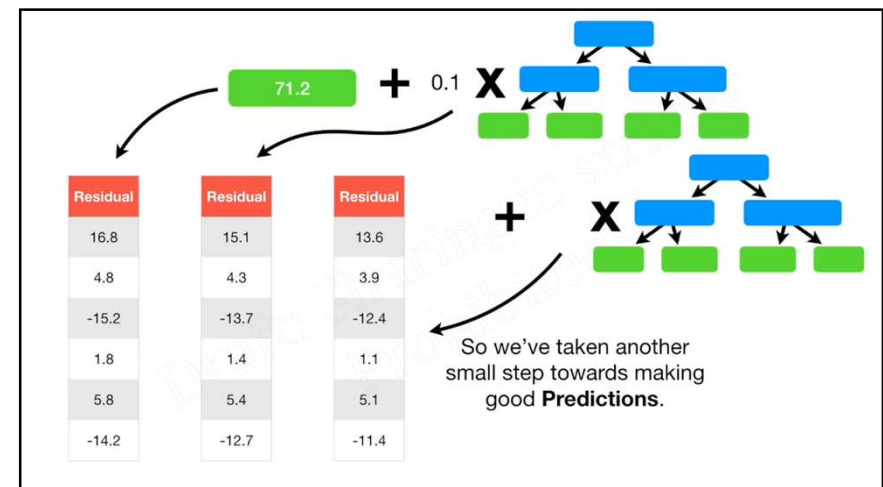
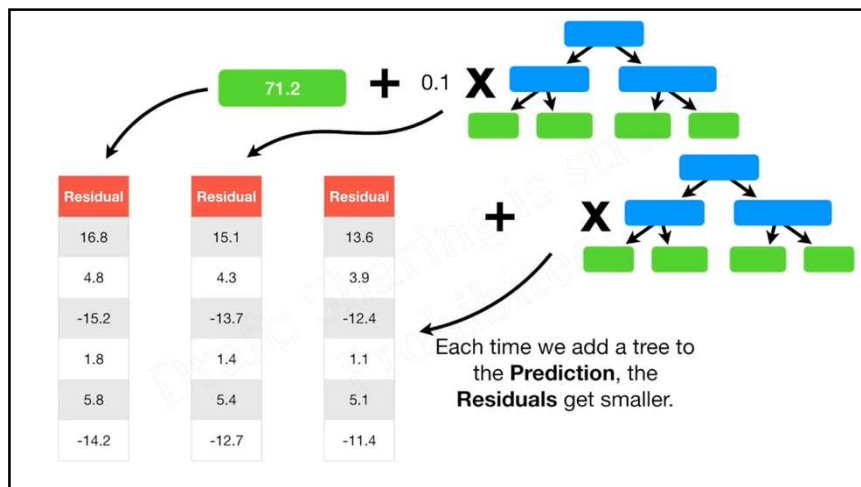
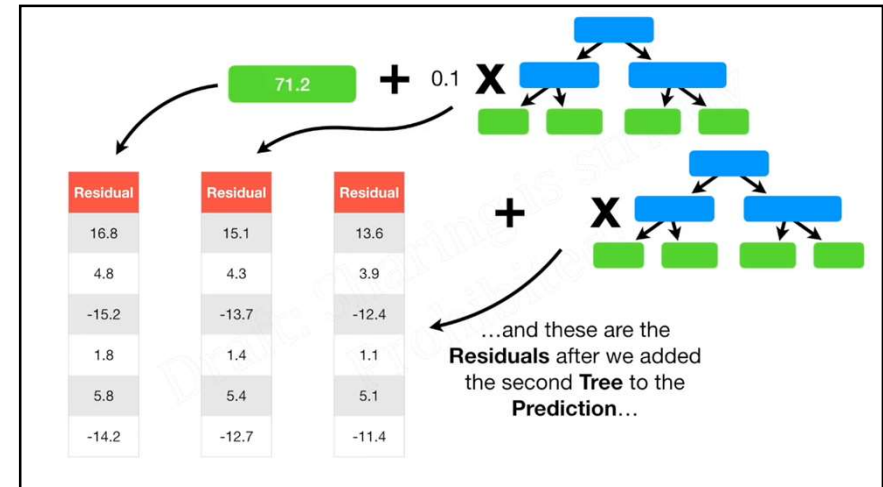
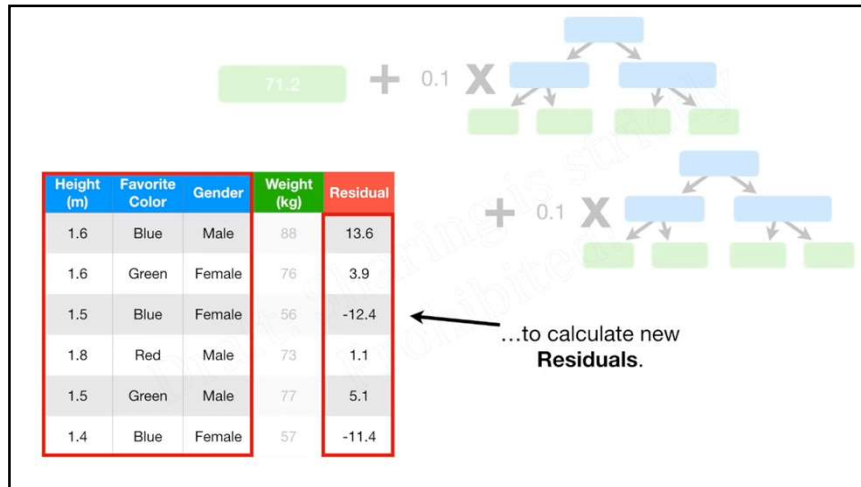


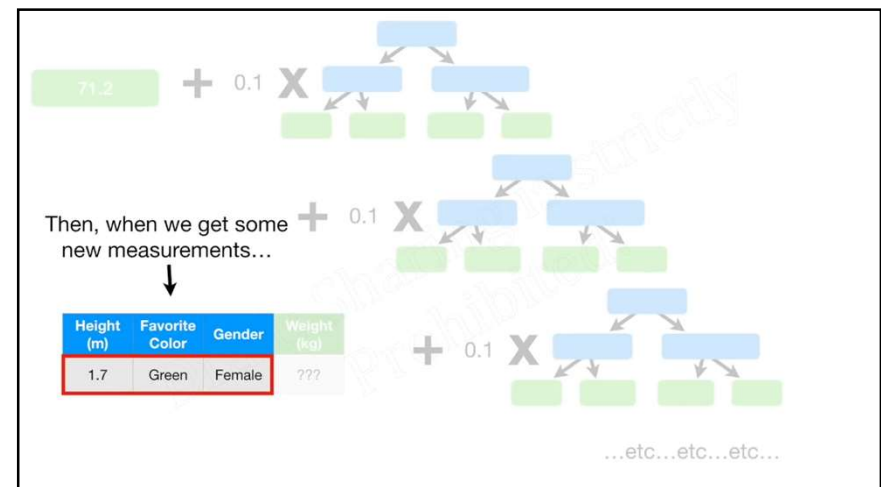
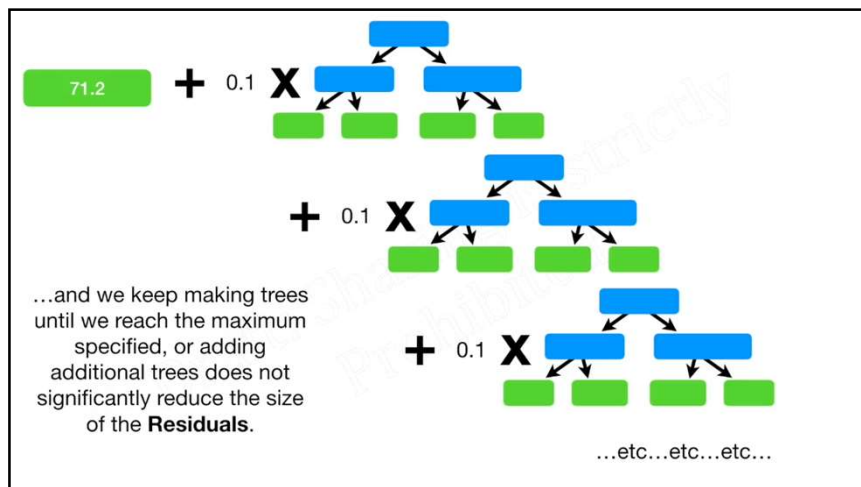
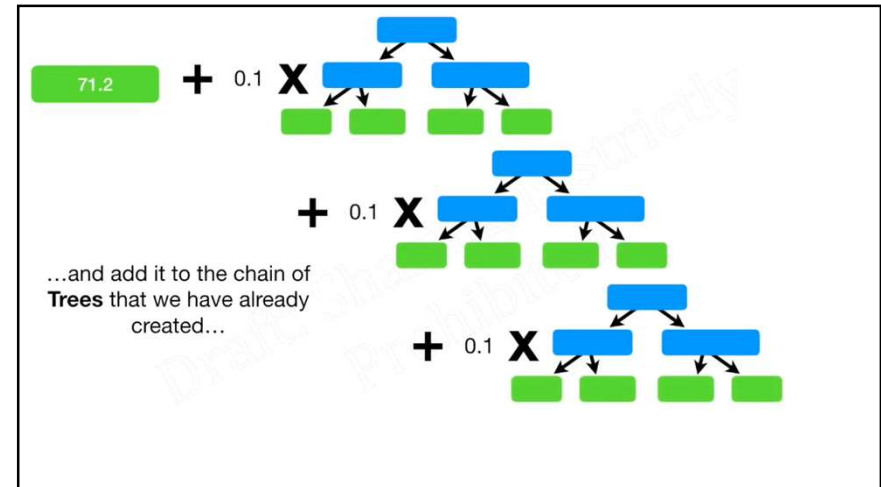
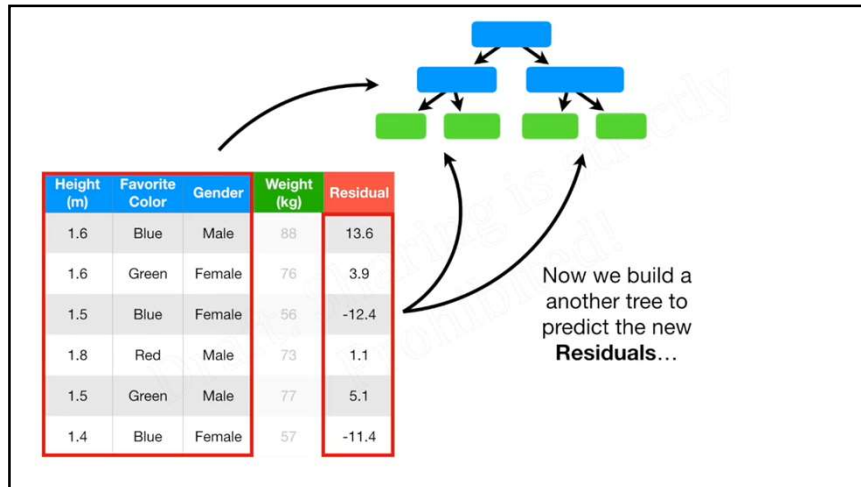


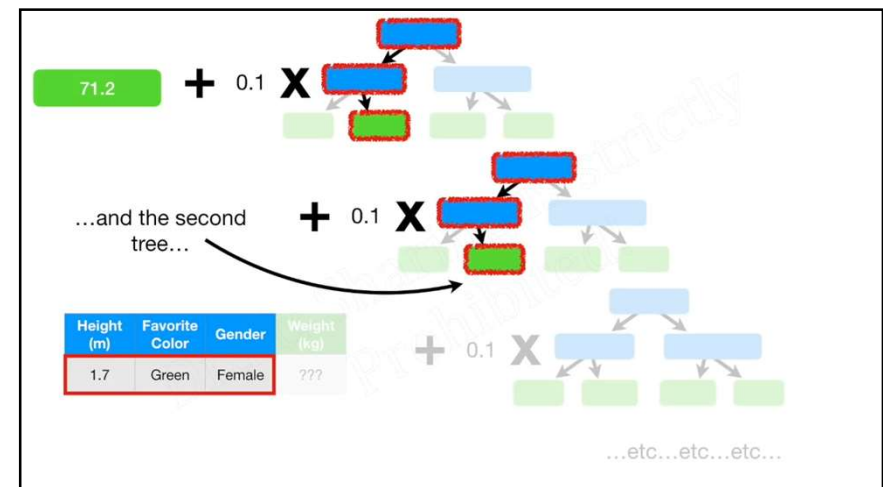
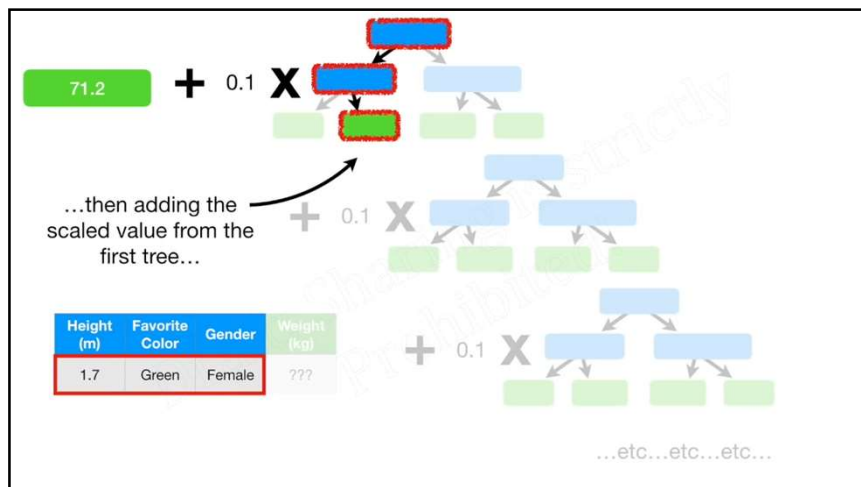
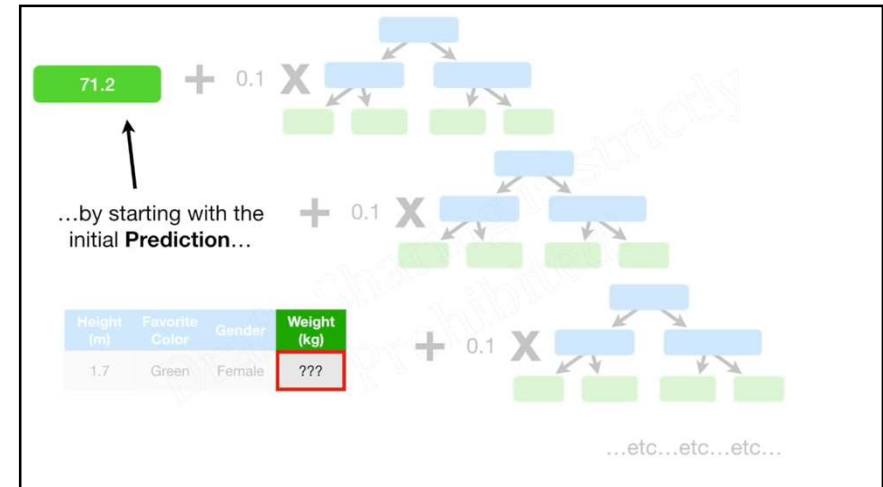
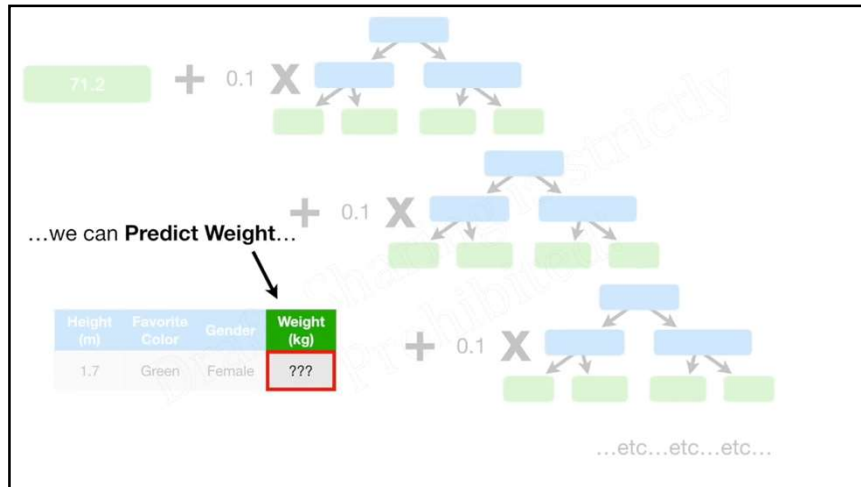


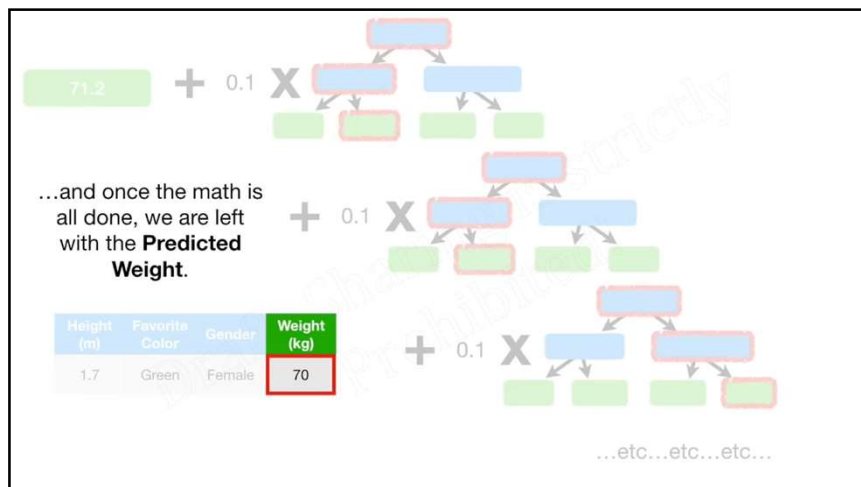
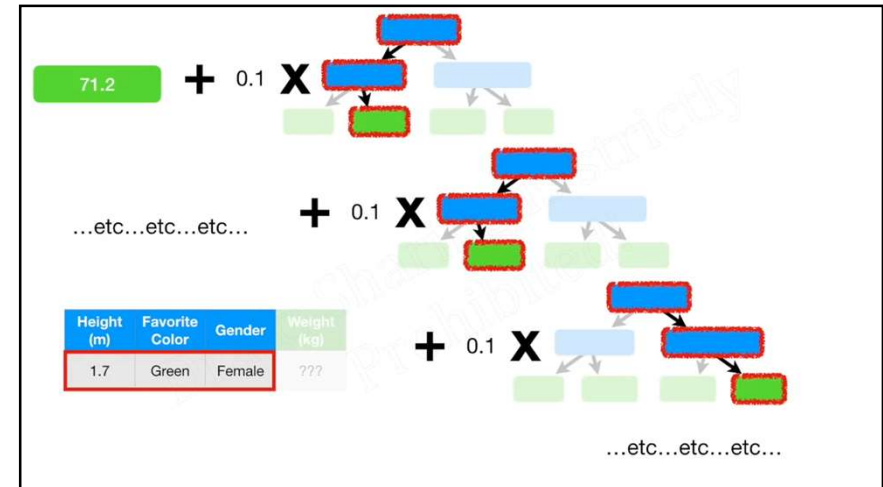
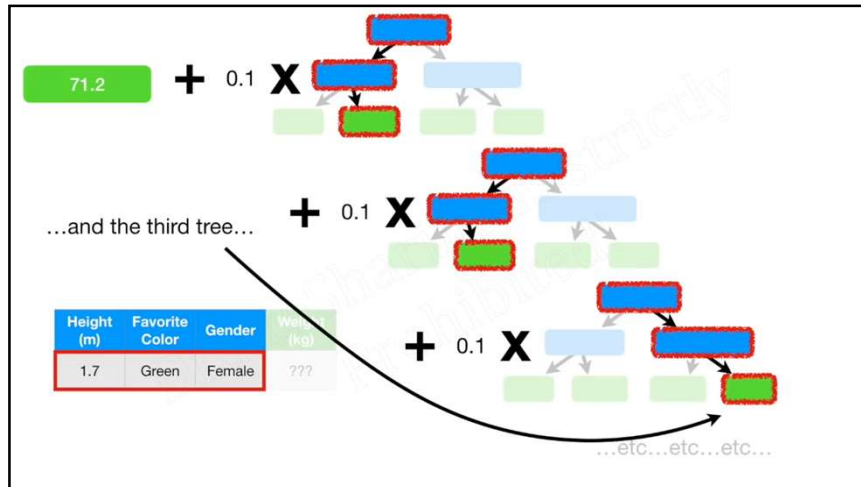




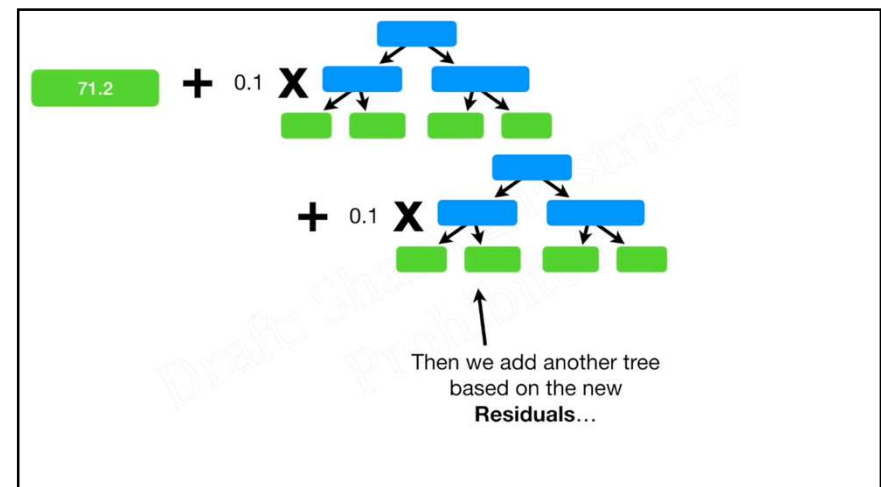
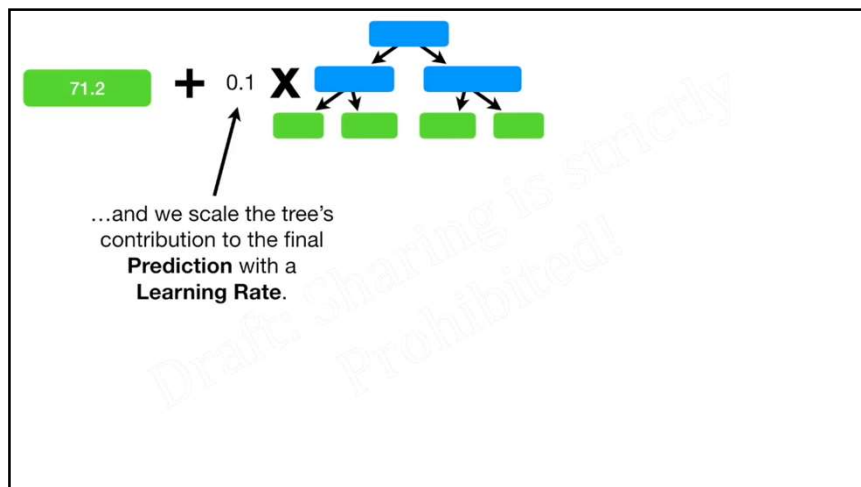
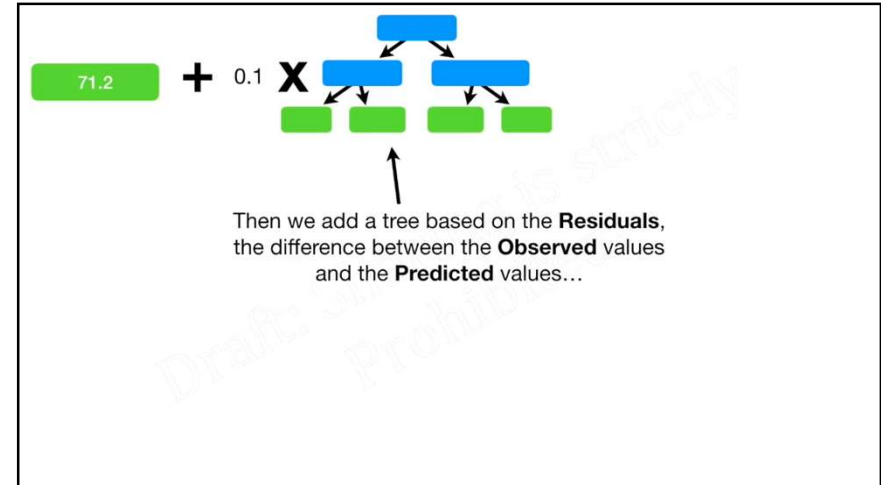
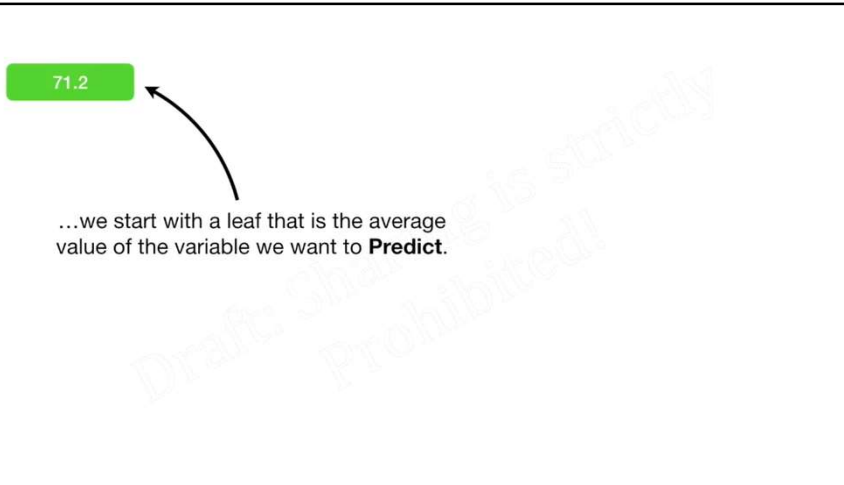




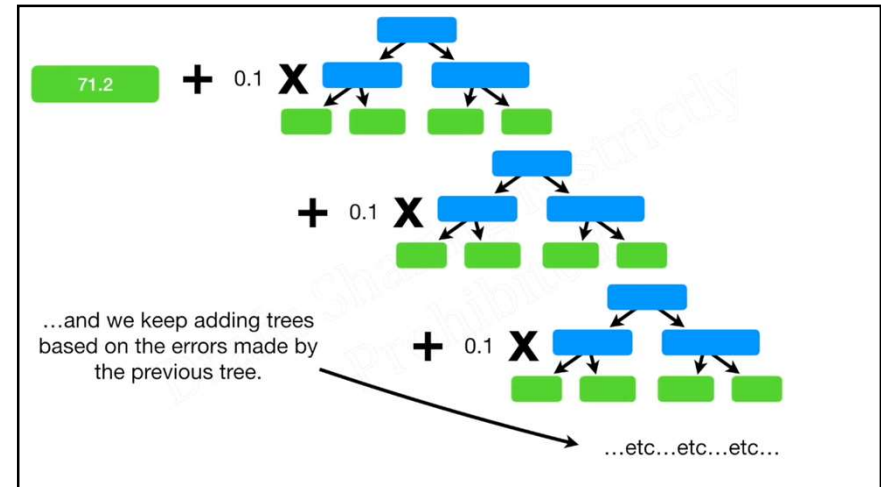
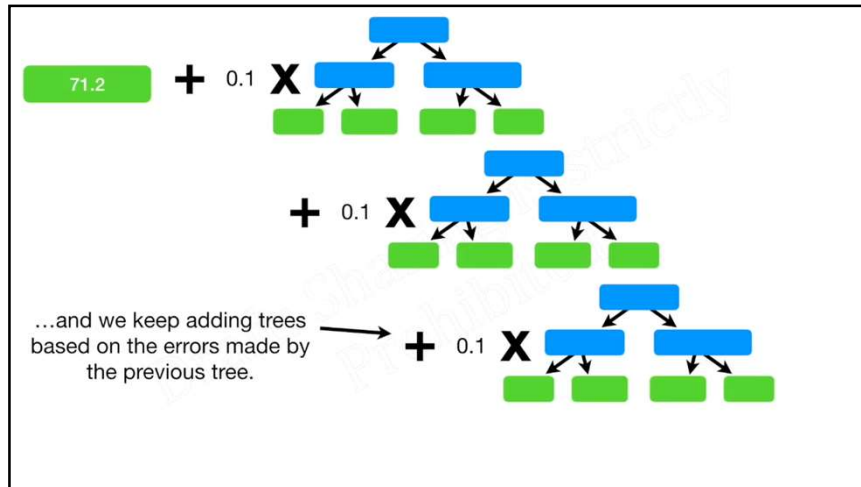




In summary, when **Gradient Boost** is used for **Regression**...








---

THANK YOU

---

StatQuest with Josh Starmer  
<https://www.youtube.com/watch?v=9CC#t=423GJo>