

<https://www.youtube.com/watch?v=3CC4N4z3GJc>

# Gradient Boosting Model

① GBM for regression

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

original data

Weight (kg)
88
76
56
73
77
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.

for continues model, it's gradient boost for regression



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

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

- (A) Compute  $r_m = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for **Part 2** in this series will dive deep into the math behind the Gradient Boost algorithm for Regression,
- (B) Fit a regression tree to the  $r_m$  values and walking through it step-by-step and proving that what we cover to day is correct.
- (C) For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \underset{\gamma}{\operatorname{argmin}} \sum_{i \in R_{jm}} L(y_i, F_{m-1}(x_i) + \gamma)$
- (D) Update  $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

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

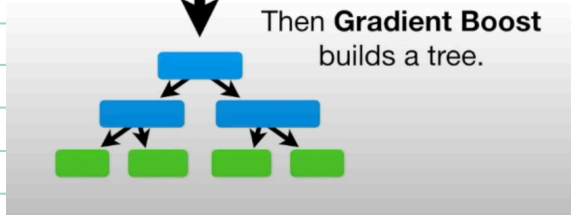
← GB Algo

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

① GB starts with a single leaf.

73.3

which is the average -



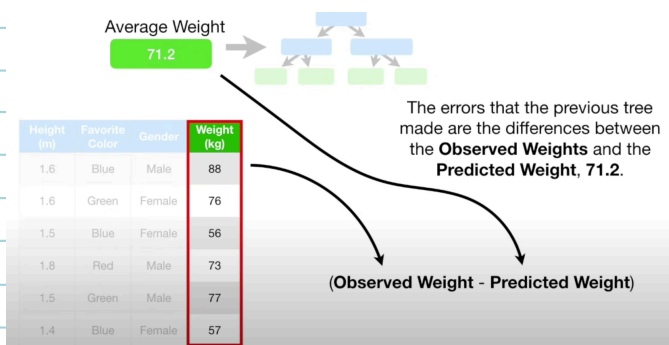
(usually maximum number of tree between 8-22)

Average Weight  
71.2

	$X_1$	$X_2$	$X_3$	$Y$
	Height (m)	Favorite Color	Gender	Weight (kg)
1	1.6	Blue	Male	88
2	1.6	Green	Female	76
3	1.5	Blue	Female	56
4	1.8	Red	Male	73
5	1.5	Green	Male	77
6	1.4	Blue	Female	57

The first thing we do is calculate the average **Weight**.

1. Calculate the avg.  $\bar{Y} = 71.2$



2. we calculate the errors.

took mean as  $\hat{y}$

then error is  $y_i - \hat{y}$

errors are the pseudual residuals.

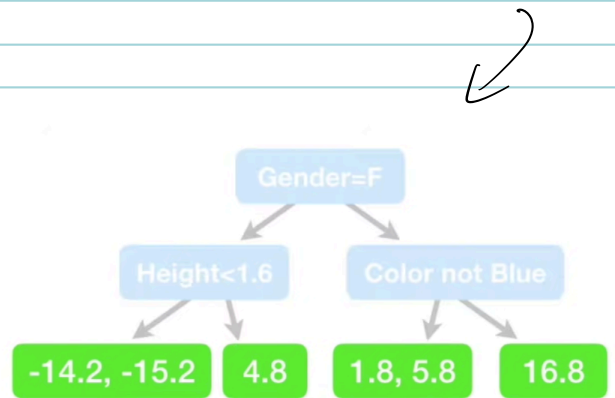
Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	16.8
1.6	Green	Female	76	4.8
1.5	Blue	Female	56	-15.2
1.8	Red	Male	73	1.8
1.5	Green	Male	77	5.8
1.4	Blue	Female	57	-14.2

$X_1$   $X_2$   $X_3$   $Y_R$

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

If it seems strange to Predict the Residuals instead of the original Weights, just bear with me and soon all will become clear.

3. then we built a new tree with  $X_1, X_2, X_3$  to predict the  $Y_R$ .  
(distance)  
then we get a tree



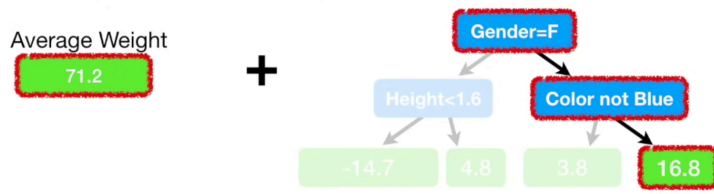
Remember, in this example we are only allowing up to four leaves...

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

So we replace these residuals with their average.

$$\frac{(-14.2 + -15.2)}{2} = -14.7$$

④ now combine original leaf with new tree



to make prediction for individual obs from training data

...so the Predicted Weight =  $71.2 + 16.8 = 88$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88

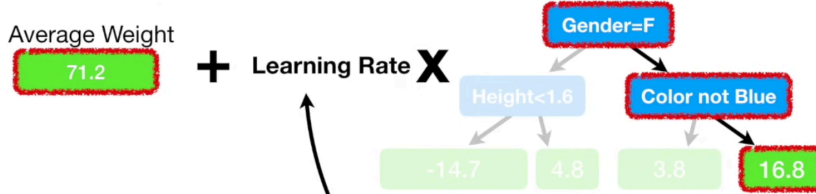
same as observed weight.

Predicted Weight =  $71.2 + 16.8 = 88$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88

No. The model fits the Training Data too well.

⑤ so we will scale the tree with learning rate



Gradient Boost deals with this problem by using a Learning Rate to scale the contribution from the new tree.

The Learning Rate is a value between 0 and 1.



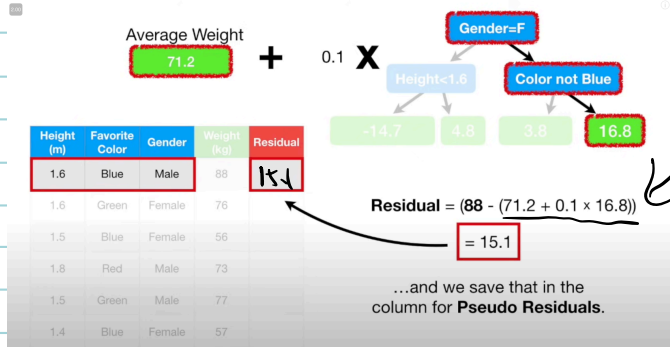
Predicted Weight =  $71.2 + (0.1 \times 16.8) = 72.9$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88

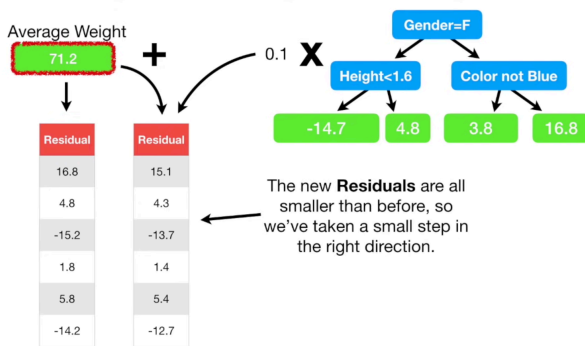
With the Learning Rate set to 0.1, the new Prediction isn't as good as it was before...

(with learning rate, it means a small step to the right direction)

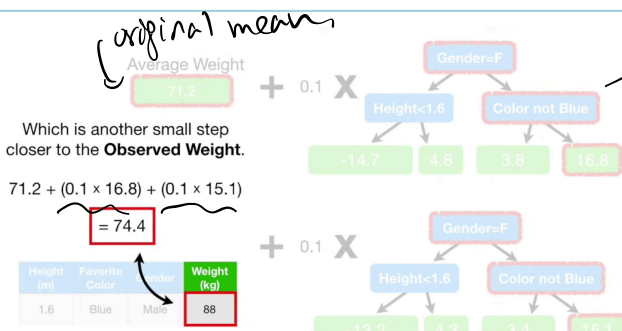
So we calculate another tree:



72.9.  $88 - 72.9 = 15.1$



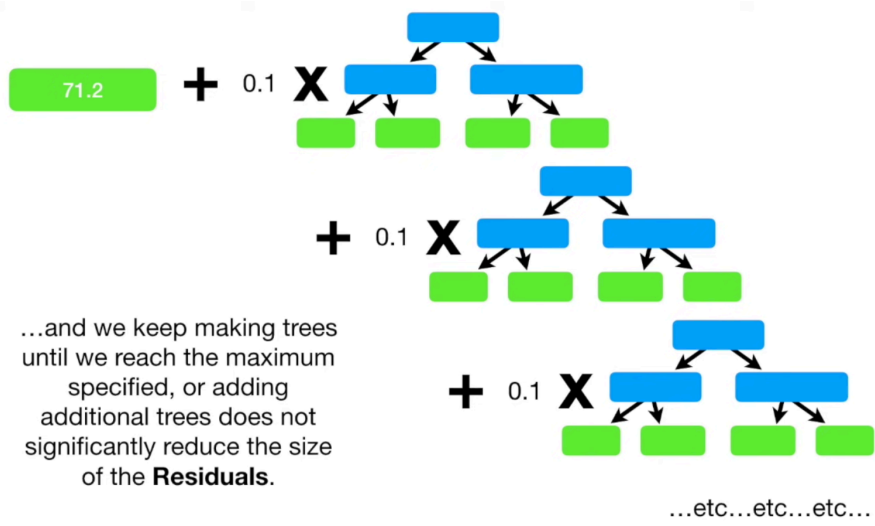
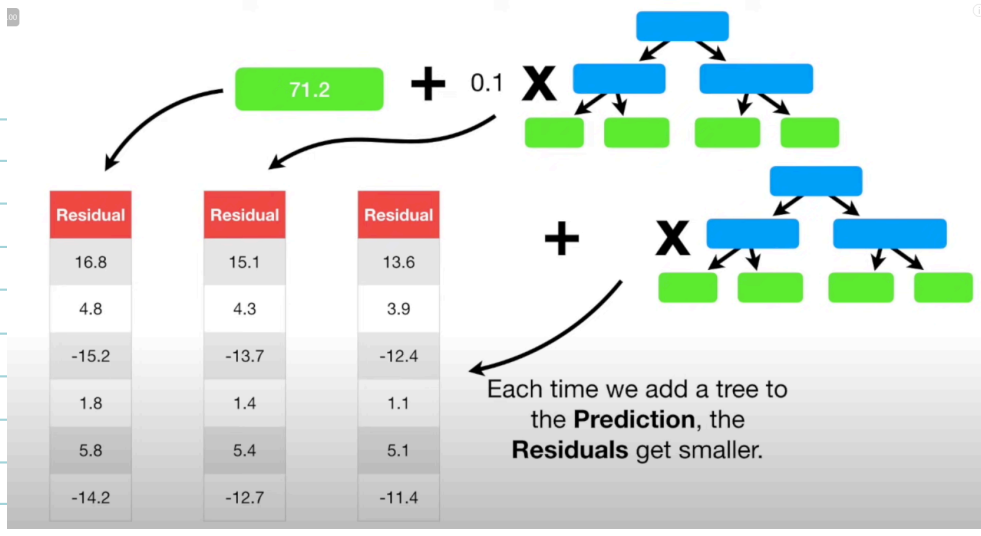
we repeat the process.



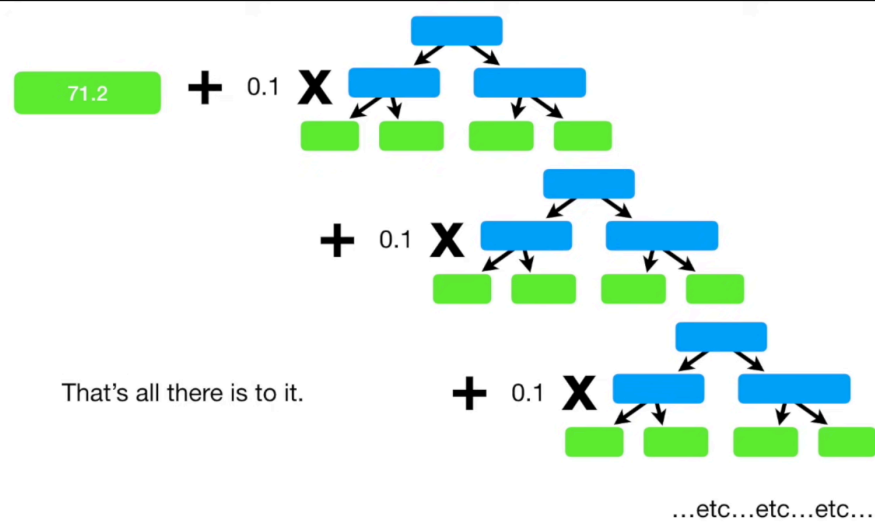
from first tree.

from second trees

another small step.



...and we keep making trees until we reach the maximum specified, or adding additional trees does not significantly reduce the size of the **Residuals**.



GB for classification.

Likes Popcorn	Age	Favorite Color	Loves Troll 2
Yes	12	Blue	Yes
Yes	87	Green	Yes
No	44	Blue	No
Yes	19	Red	No
No	32	Green	Yes
No	14	Blue	Yes

Likes Popcorn	Age	Favorite Color	Loves Troll 2
Yes	12	Blue	Yes
Yes	87	Green	Yes
No	44	Blue	No
Yes	19	Red	No
No	32	Green	Yes
No	14	Blue	Yes

initial prediction

$\log(4/2) = 0.7$

...which we will put into our initial leaf.

$\log(\frac{4}{2}) = 0.7$

2 no log of the odds.

$\log(4/2) = 0.7$   
Probability of Loving Troll 2 = 0.7

NOTE: These two numbers, the log and the Probability are the same because I'm rounding. If I allowed digits passed the decimal place...

Likes Popcorn	Age	Favorite Color	Loves Troll 2
Yes	12	Blue	Yes
Yes	87	Green	Yes
No	44	Blue	No
Yes	19	Red	No
No	32	Green	Yes

$\log(4/2) = 0.7$   
Probability of Loving Troll 2 = 0.7

NOTE: These two numbers, the  $\log(4/2)$  and the Probability are the same only because I'm rounding. If I allowed 4 digits passed the decimal place...

$\log(\frac{4}{2}) = 0.6931$

$\frac{e^{\log(4/2)}}{1 + e^{\log(4/2)}} = 0.6667$

way of turning it to prob.

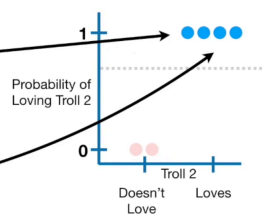
Since majority agree.  $\log(4/2)$  is yes.

so initial is yes

$\log(4/2) = 0.7$   
Probability of Loving Troll 2 = 0.7

...and the Blue Dots, with a Probability of Loving Troll 2 = 1, represent the four people that Love Troll 2.

Likes Popcorn	Age	Favorite Color	Loves Troll 2
Yes	12	Blue	Yes
Yes	87	Green	Yes
No	44	Blue	No
Yes	19	Red	No
No	32	Green	Yes
No	14	Blue	Yes

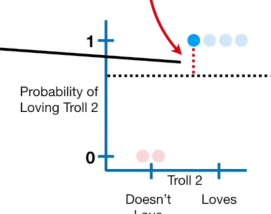


$\log(4/2) = 0.7$   
Probability of Loving Troll 2 = 0.7

...and we save the Residual in a new column.

Residual =  $(1 - 0.7) = 0.3$

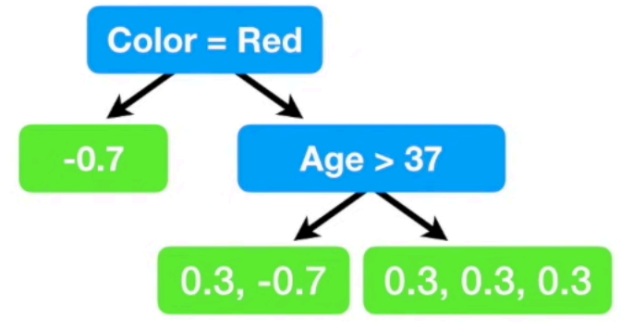
Likes Popcorn	Age	Favorite Color	Loves Troll 2	Residual
Yes	12	Blue	Yes	0.3
Yes	87	Green	Yes	
No	44	Blue	No	
Yes	19	Red	No	
No	32	Green	Yes	
No	14	Blue	Yes	



Residual = (Observed - Predicted)

then we build the tree.

Likes Popcorn	Age	Favorite Color	Loves Troll 2	Residual
Yes	12	Blue	Yes	0.3
Yes	87	Green	Yes	0.3
No	44	Blue	No	-0.7
Yes	19	Red	No	-0.7
No	32	Green	Yes	0.3
No	14	Blue	Yes	0.3



NOTE: The derivation of this formula is quite technical, so I'm saving it for Part 4 of this series when we get into the nitty gritty details of Gradient Boost for Classification.

$$\frac{\sum \text{Residual}_i}{\sum [\text{Previous Probability}_i \times (1 - \text{Previous Probability}_i)]}$$



$\log(\frac{4}{2}) = 0.7$

$\log(4/2) = 0.7$   
Probability of Loving Troll 2 = 0.7

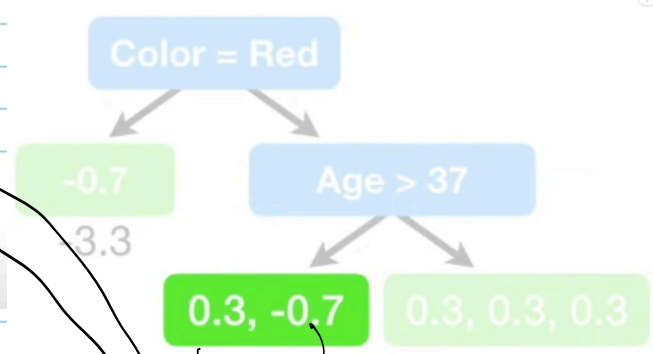
NOTE: These and the Prob because I'm digits pass

$\frac{-0.7}{0.7 \times (1 - 0.7)} = -3.3$

Likes Popcorn	Age	Favorite Color	Loves Troll 2
Yes	12	Blue	Yes
Yes	87	Green	Yes
No	44	Blue	No
Yes	19	Red	No
No	32	Green	Yes

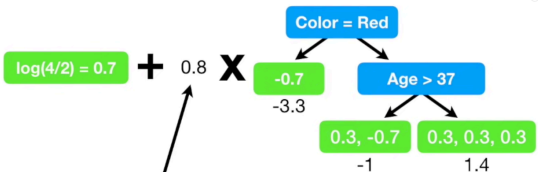
$\log(\frac{4}{2}) = 0.7$

$$\frac{\sum \text{Residual}_i}{\sum [\text{Previous Probability}_i \times (1 - \text{Previous Probability}_i)]}$$



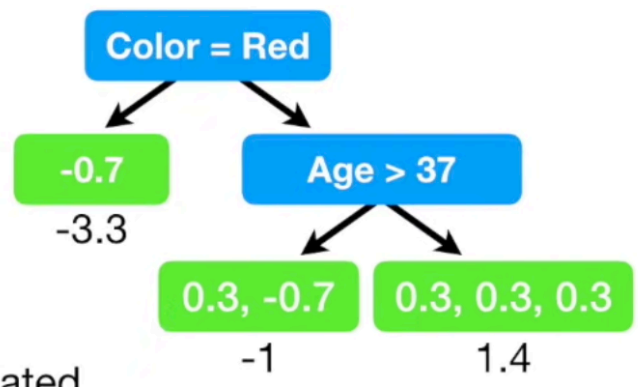
$$\frac{(0.3 - 0.7)}{(0.7 \times (1 - 0.7)) + (0.7 \times (1 - 0.7))} = -1$$





**NOTE:** Just like before, the new tree is scaled by a **Learning Rate**.

This example uses a relatively large **Learning Rate** for illustrative purposes. However, **0.1** is more common.



lated

$\log(4/2) = 0.7$   
Initial Probability of **Loving Troll 2** = 0.7

...so we are taking a small step in the right direction since this person **Loves Troll 2**.

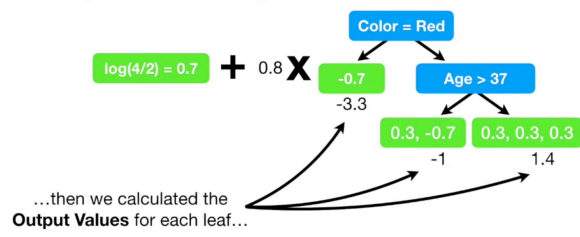
Likes Popcorn	Age	Favorite Color	Loves Troll 2
Yes	12	Blue	Yes
Yes	87	Green	Yes
No	44	Blue	No
Yes	19	Red	No
No	32	Green	Yes
No	14	Blue	Yes

$$\text{Probability} = \frac{e^{1.8}}{1 + e^{1.8}} = 0.9$$

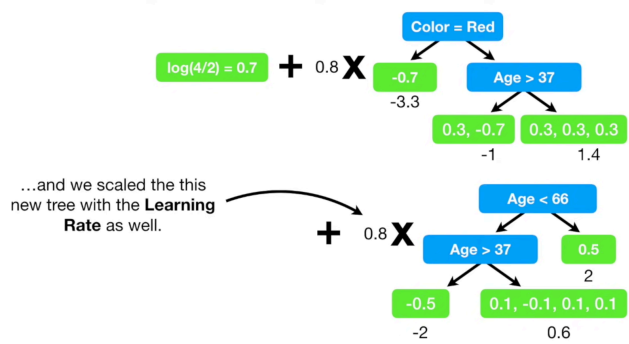
$$\log(\text{odds}) \text{ Prediction} = 0.7 + (0.8 \times 1.4) = 1.8$$

Likes Popcorn	Age	Favorite Color	Loves Troll 2	Predicted Prob.	Residual
Yes	12	Blue	Yes	0.9	0.1
Yes	87	Green	Yes	0.5	0.5
No	44	Blue	No	0.5	-0.5
Yes	19	Red	No	0.1	-0.1
No	32	Green	Yes	0.9	0.1
No	14	Blue	Yes	0.9	0.1

**BAM!**



...then we calculated the **Output Values** for each leaf...



...and we scaled the this new tree with the **Learning Rate** as well.

**Log(odds) Prediction** that someone **Loves Troll 2:**

...and get **2.3** as the **Log(odds) Prediction** that this person **Loves Troll 2**.

$$= 0.7 + (0.8 \times 1.4) + (0.8 \times 0.6) = 2.3$$

**Log(odds) Prediction** that someone **Loves Troll 2:**

$$= 0.7 + (0.8 \times 1.4) + (0.8 \times 0.6) = 2.3$$

...and the **Predicted Probability** that this individual will **Love Troll 2** is **0.9**.

$$\text{Probability} = \frac{e^{2.3}}{1 + e^{2.3}} = 0.9$$

Likes Popcorn	Age	Favorite Color	Loves Troll 2
Yes	25	Green	???

Likes Popcorn	Age	Favorite Color	Loves Troll 2
Yes	25	Green	???