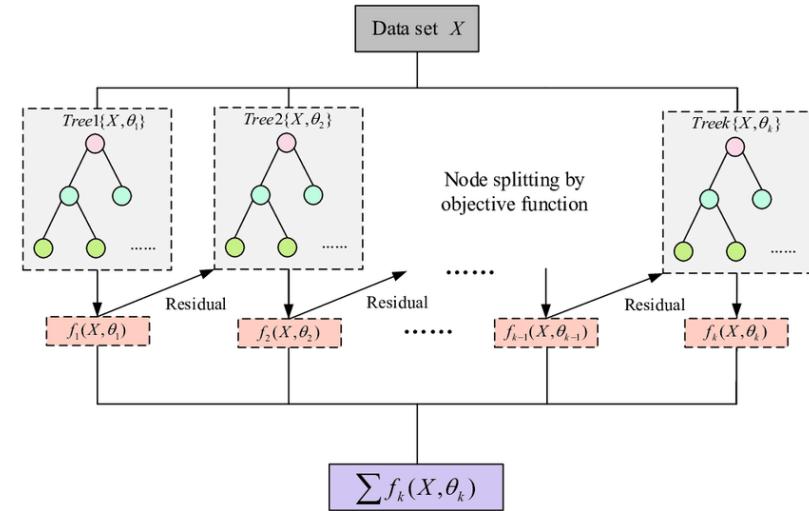
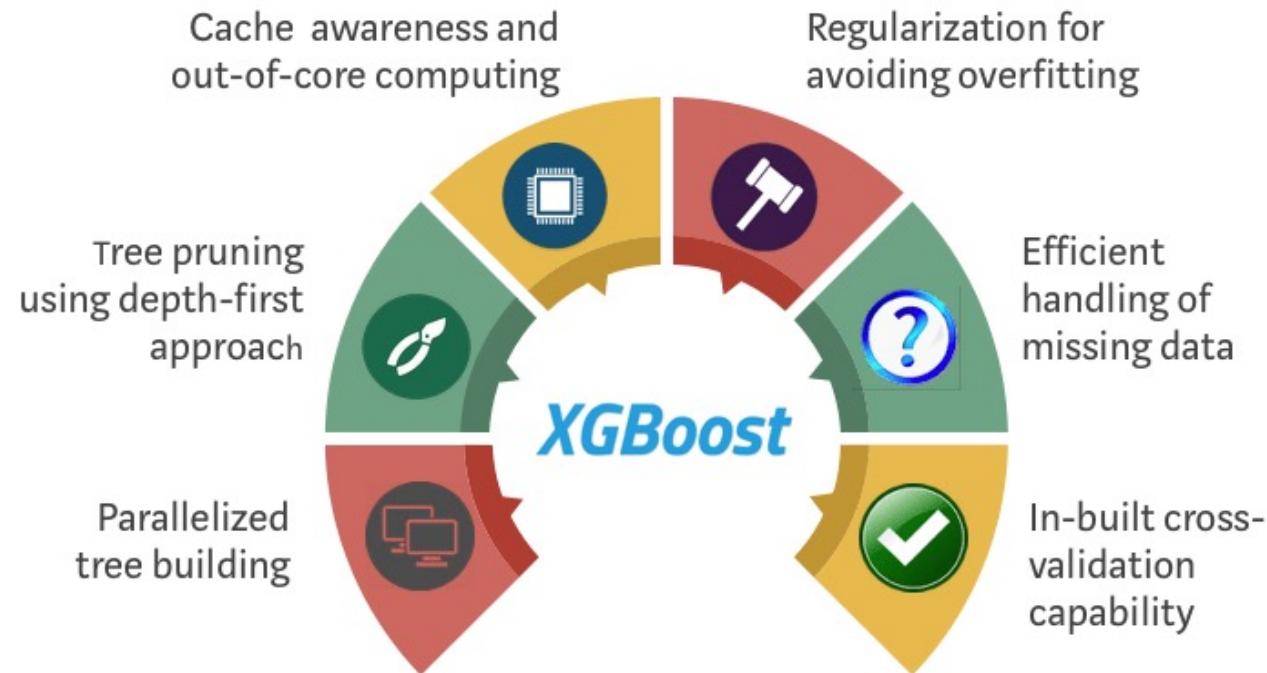


XGBoost



Vinh Dinh Nguyen
PhD in Computer Science

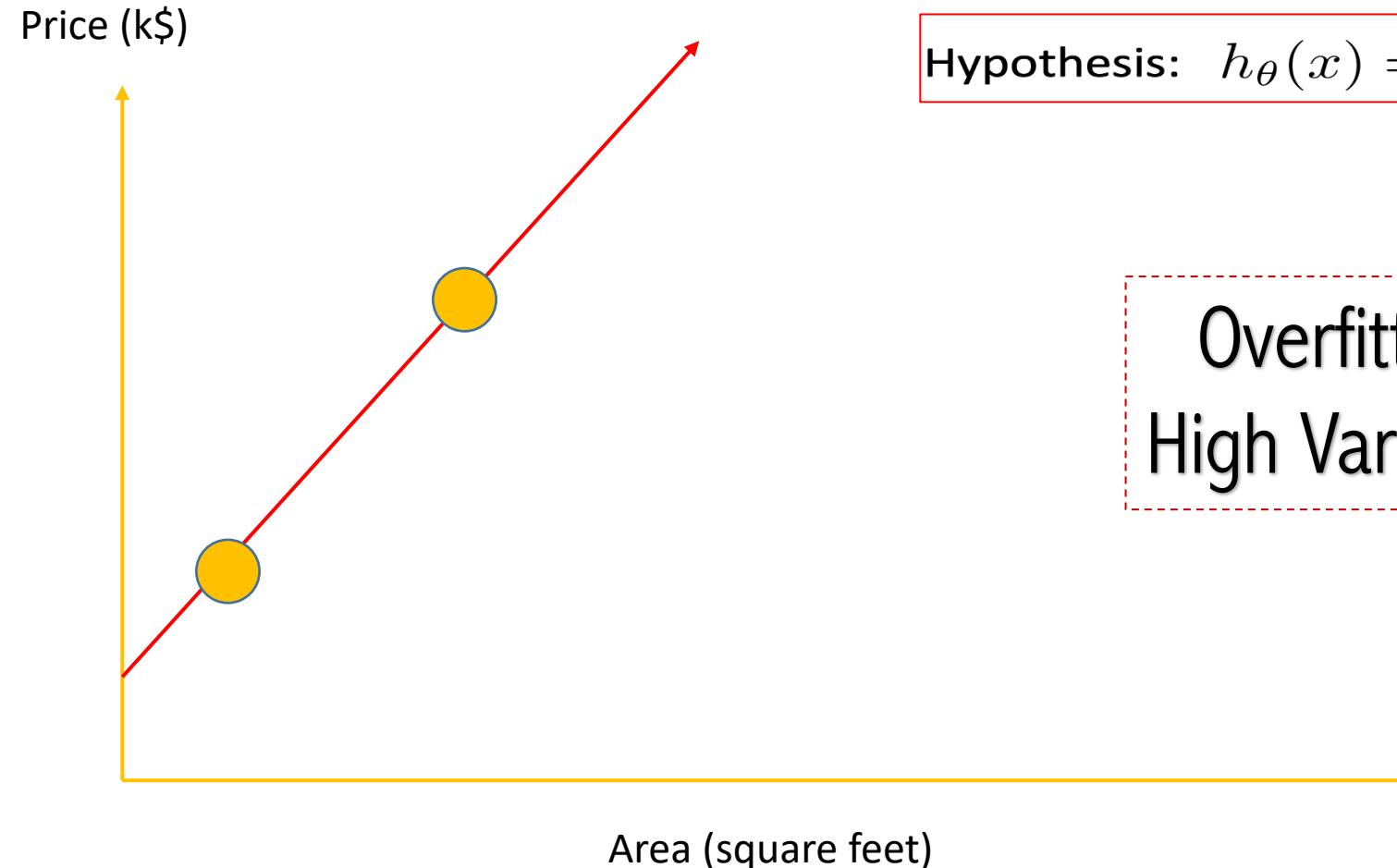
Outline

- **Regularization**
- **XGBoost For Regression**
- **XGBoost For Classification**
- **XGBoost: Mathematical Explanation**
- **How To Fill Missing Values**
- **Example**

Outline

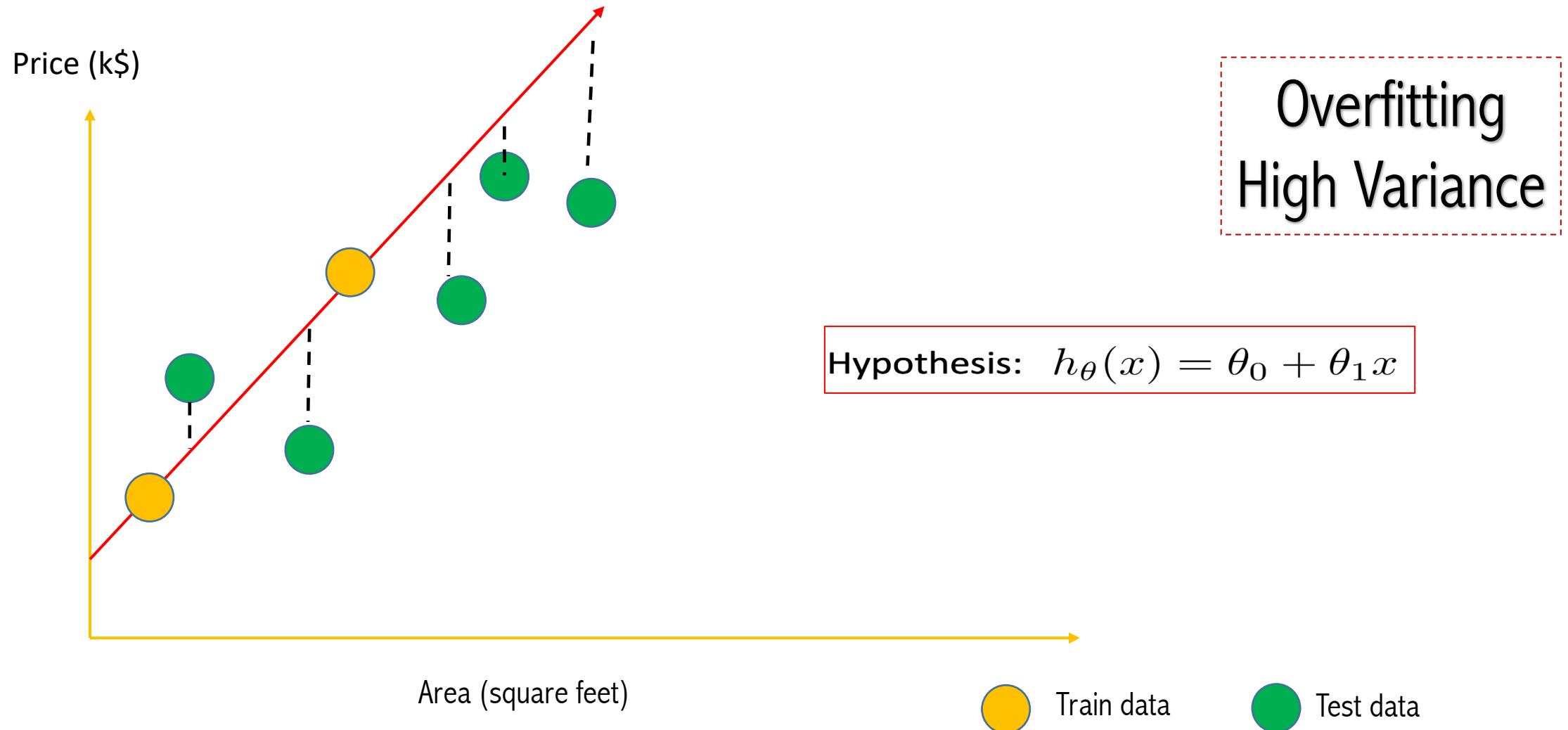
- **Regularization**
- **XGBoost For Regression**
- **XGBoost For Classification**
- **XGBoost: Mathematical Explanation**
- **How To Fill Missing Values**
- **Example**

Regularization

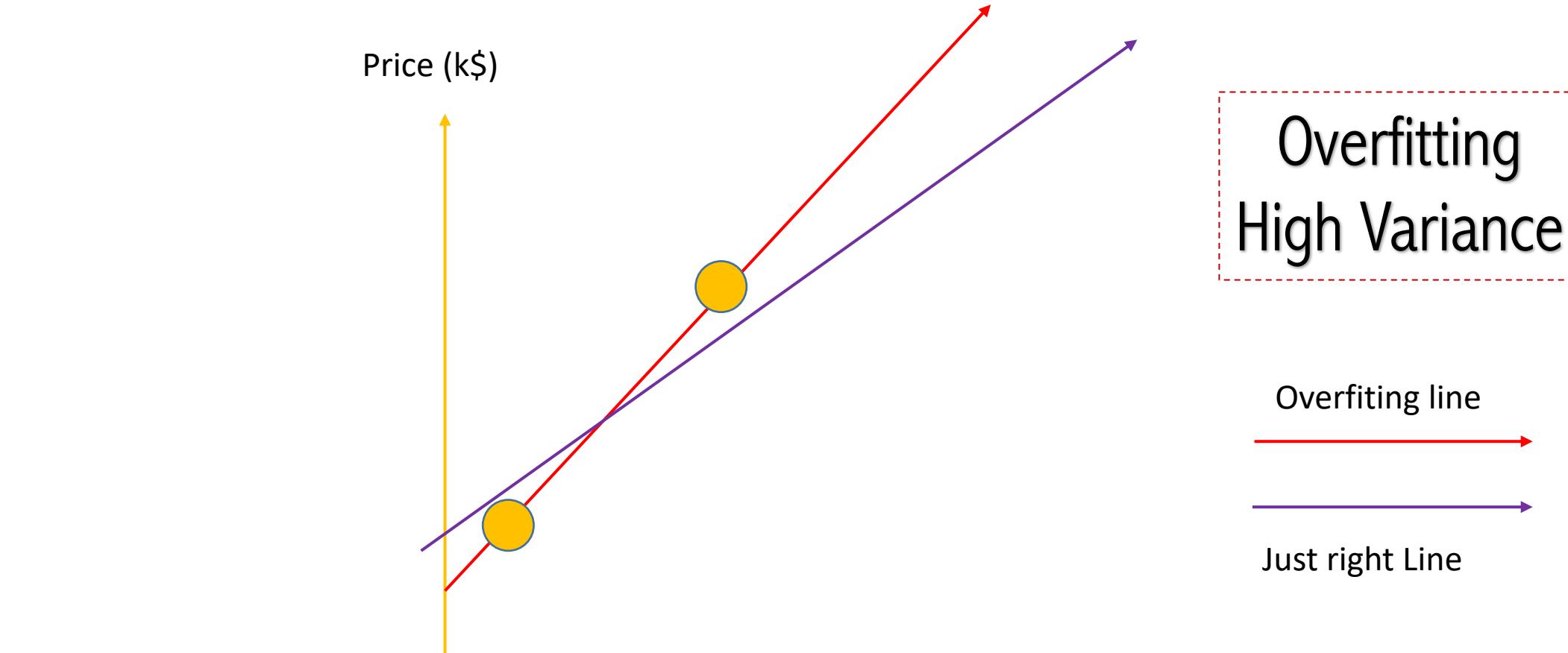


Overfitting
High Variance

Regularization

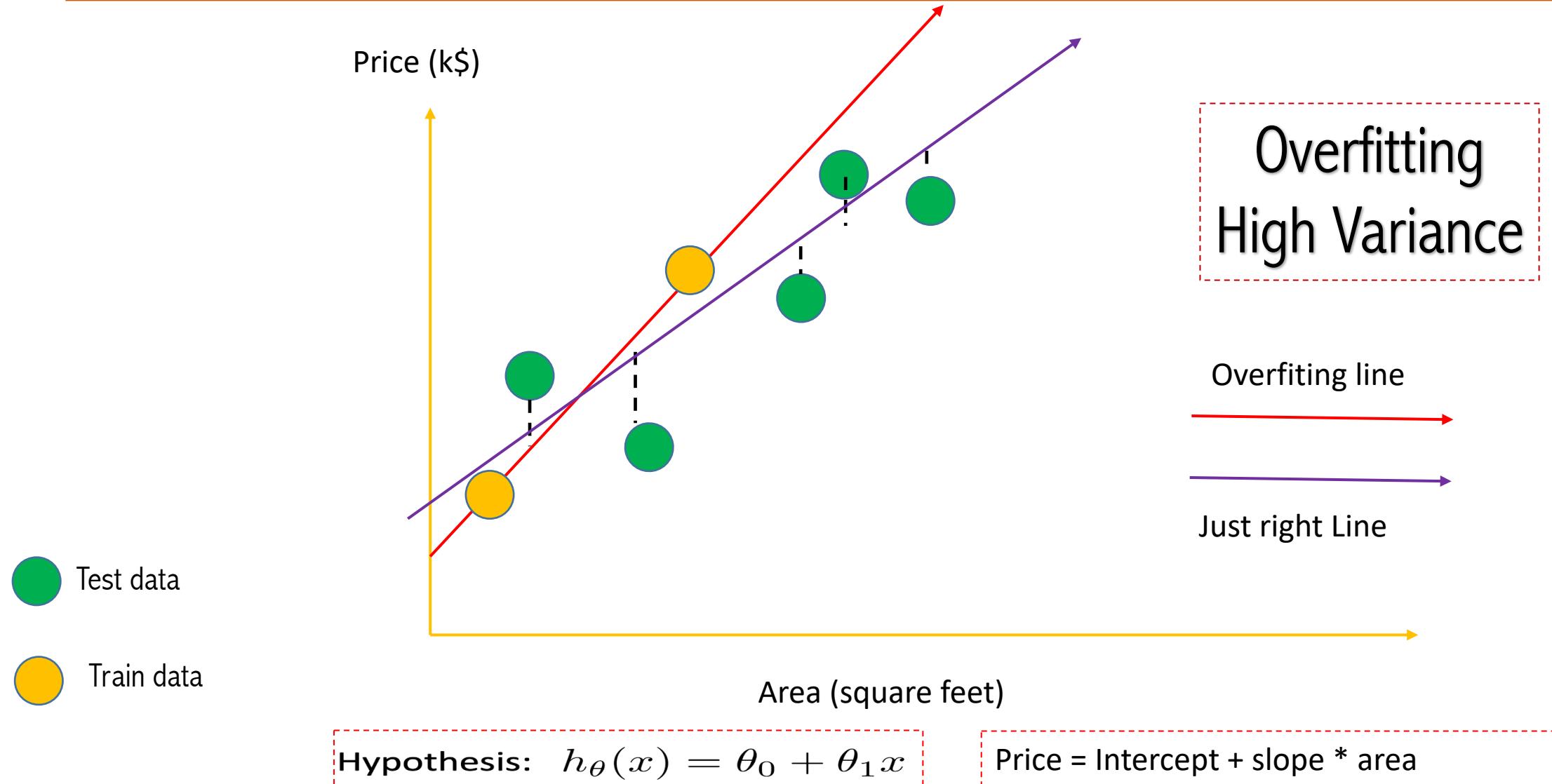


Regularization



$$\text{Hypothesis: } h_{\theta}(x) = \theta_0 + \theta_1 x$$

Regularization



Regularization

Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$

Parameters: θ_0, θ_1

Cost Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Goal: $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$

Regularization.

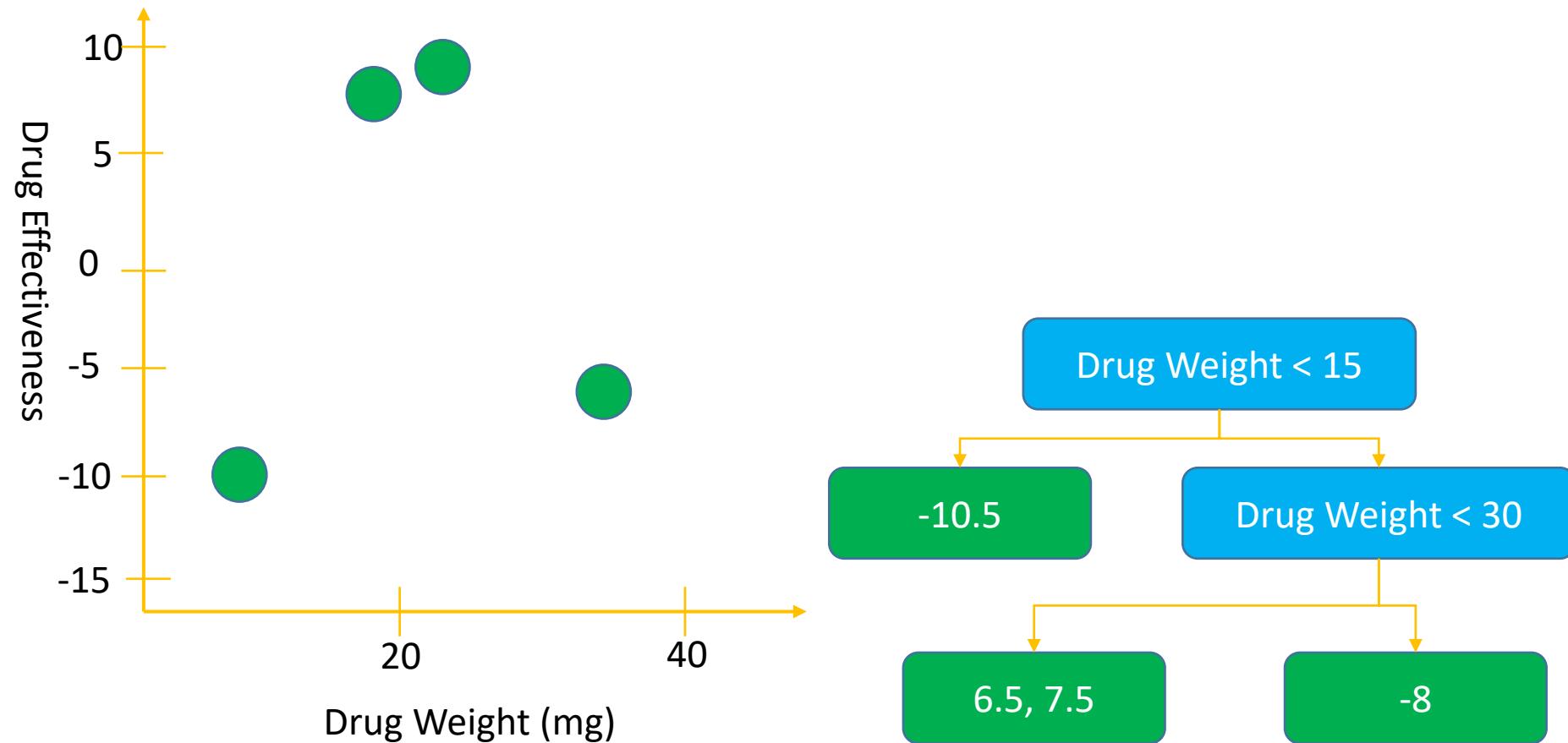
$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

Price = [Intercept + slope * area] + $\lambda * \text{slope}^2$

Outline

- **Regularization**
- **XGBoost For Regression**
- **XGBoost For Classification**
- **XGBoost: Mathematical Explanation**
- **How To Fill Missing Values**
- **Example**

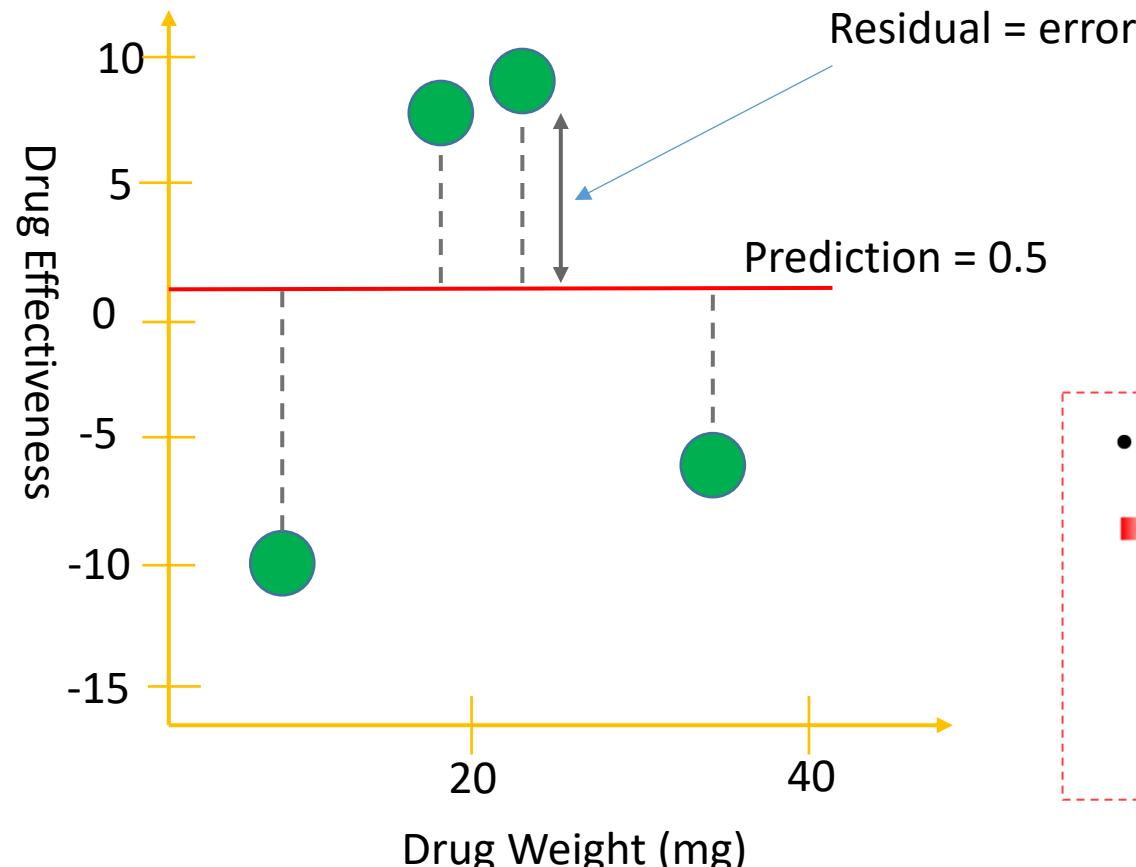
XGBoost For Regression



XGBoost For Regression

Step 1

- Initialize the first prediction for drug effectiveness
- Any number, for default, we set 1st prediction = 0.5



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

• $\{(x_i, y_i)\}_{i=1}^n$

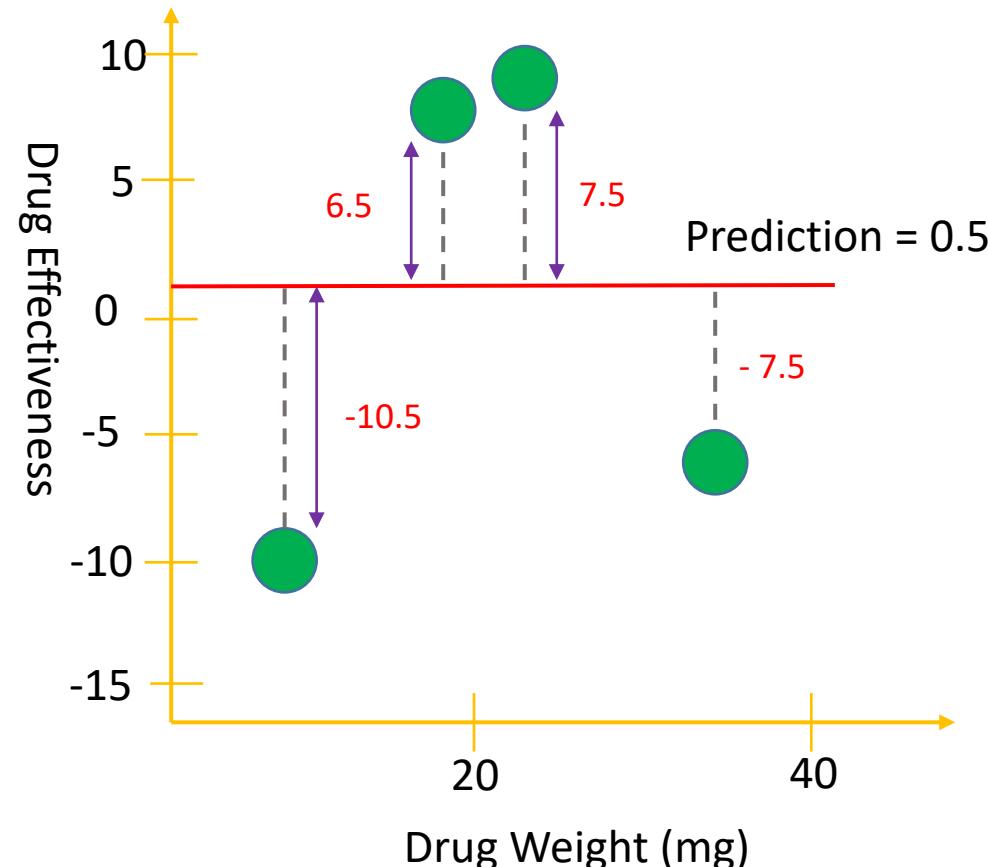
■ Loss function = $L(y_i, F(x)) = 1/2 * (\text{Output} - \text{Predicted})^2$

$$\frac{dL}{d\text{Predicted}} = 2/2 (\text{Output} - \text{Predicted}) * -1 = -(\text{Output} - \text{Predicted})$$

Tricky implementation here

XGBoost For Regression

Step 1



Start with single Leaf of residuals

-10.5, 6.5, 7.5, -7.5

Compute Similarity Score

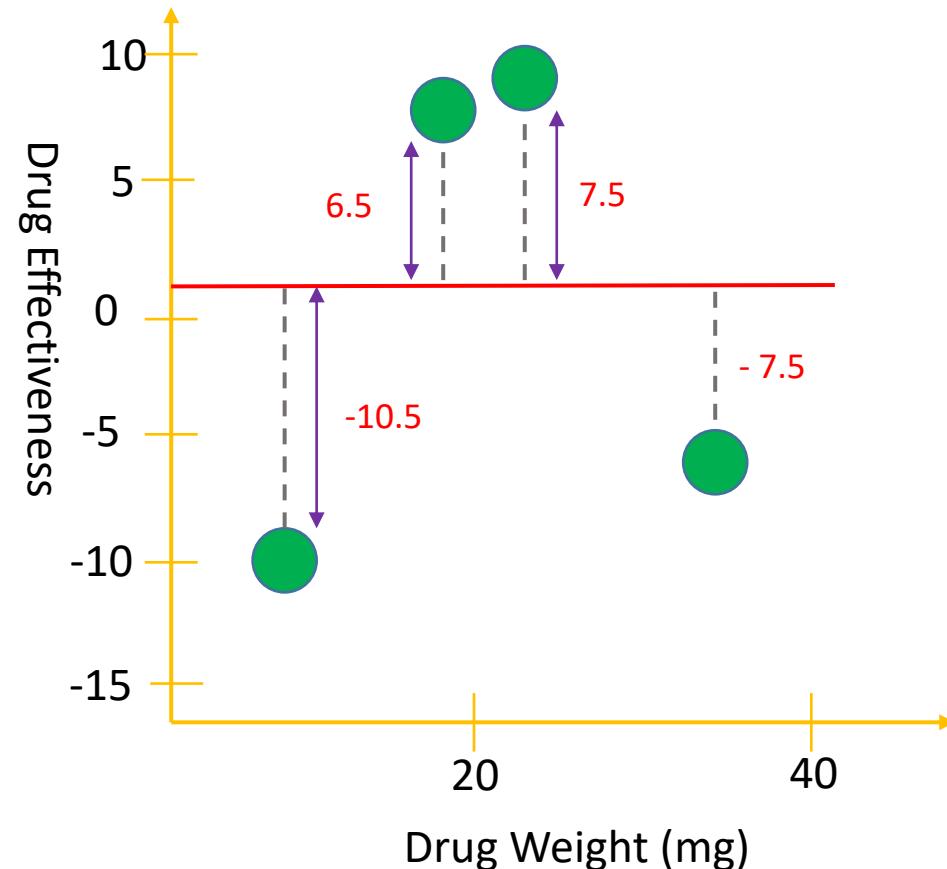
$$SC = \frac{[\sum (\text{output} - \text{predicted})]^2}{m + \lambda}$$

m : number of samples

λ : regularization parameters

XGBoost For Regression

Step 1



Start with single Leaf of residuals

-10.5, 6.5, 7.5, -7.5

$m = 4$
 $\lambda = 0$

Compute Similarity Score

$$SC = \frac{[\sum (output - predicted)]^2}{m + \lambda}$$

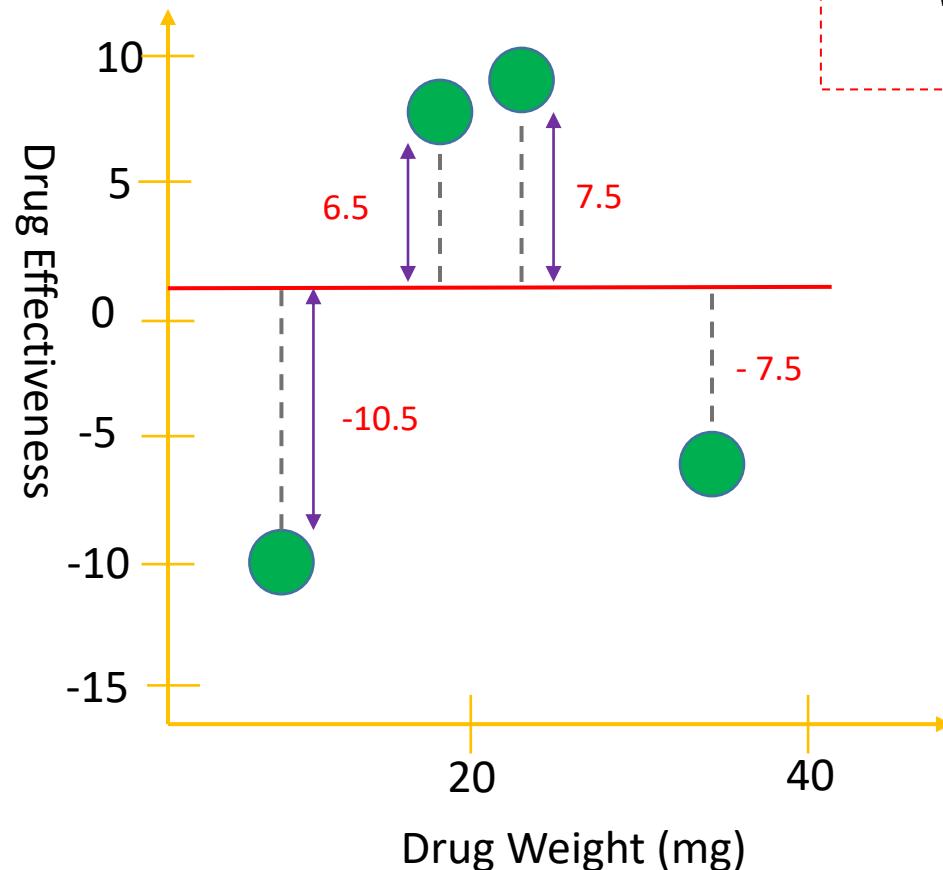
$$SC = \frac{[-10.5 + 7.5 + 6.5 + (-7.5)]^2}{4} = 4$$

XGBoost For Regression

Step 1

SC = 4

-10.5, 6.5, 7.5, -7.5



What happens if we try to split residuals into two groups =>
measure the similarity score

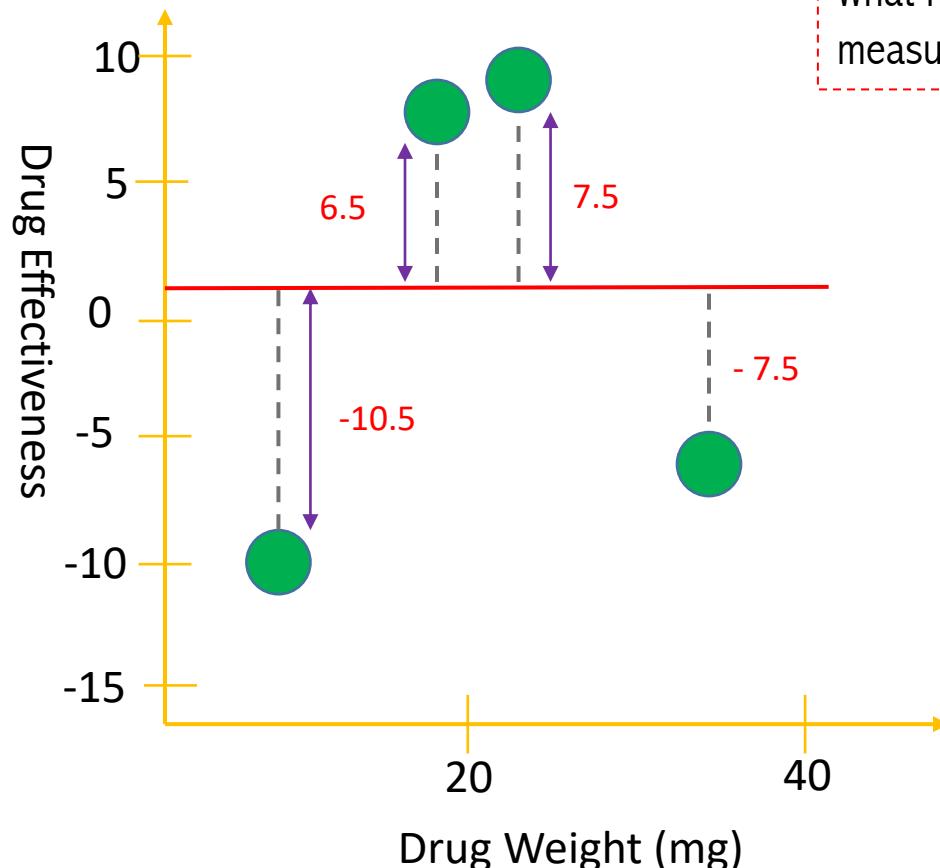


XGBoost For Regression

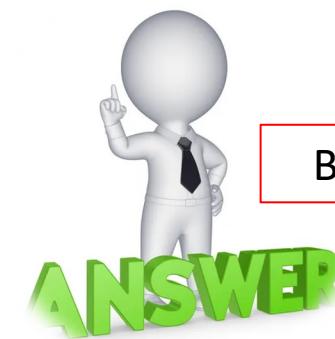
Step 1

SC = 4

-10.5, 6.5, 7.5, -7.5



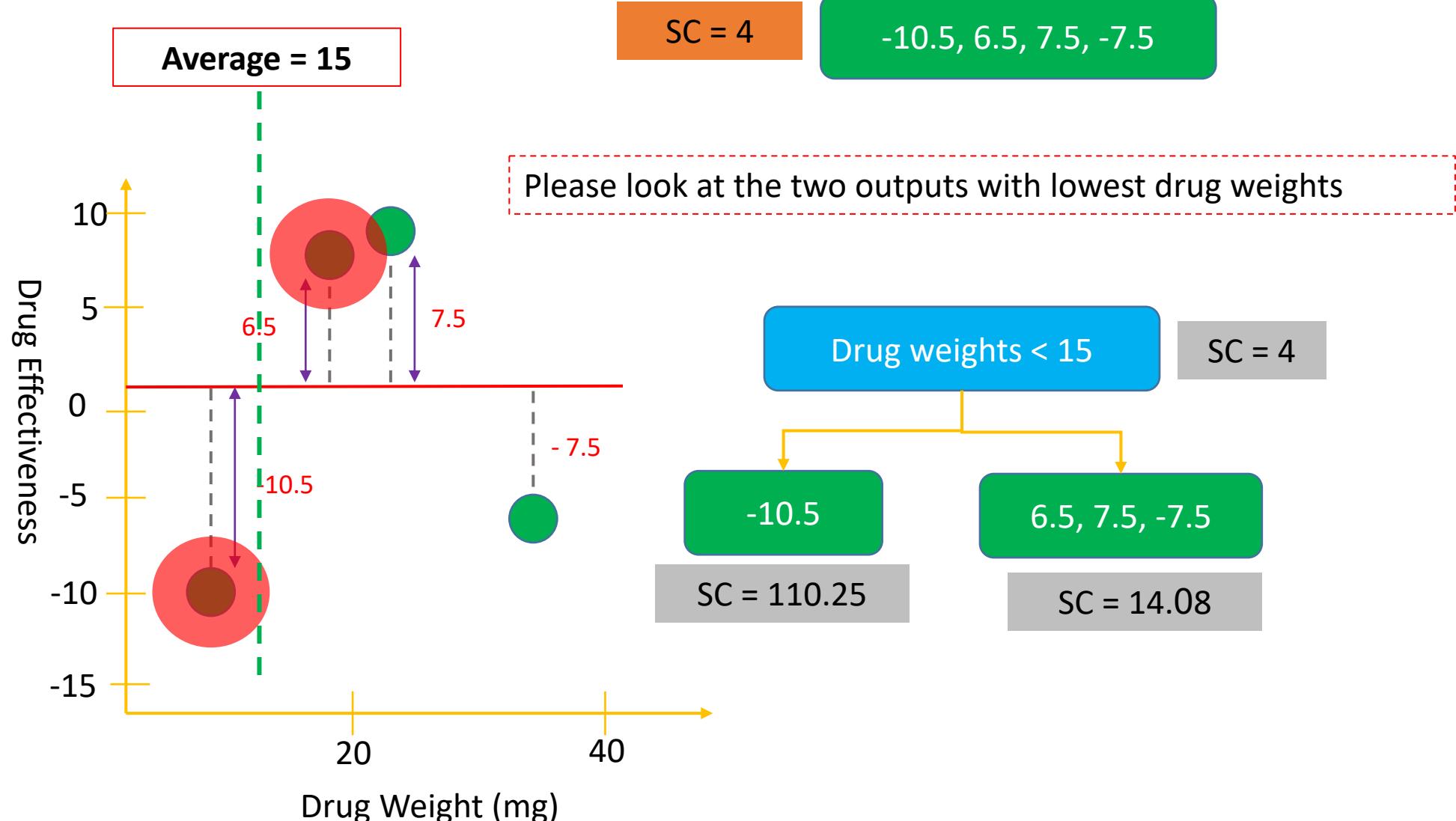
What happens if we try to split residuals into two groups =>
measure the similarity score



Build a tree on it

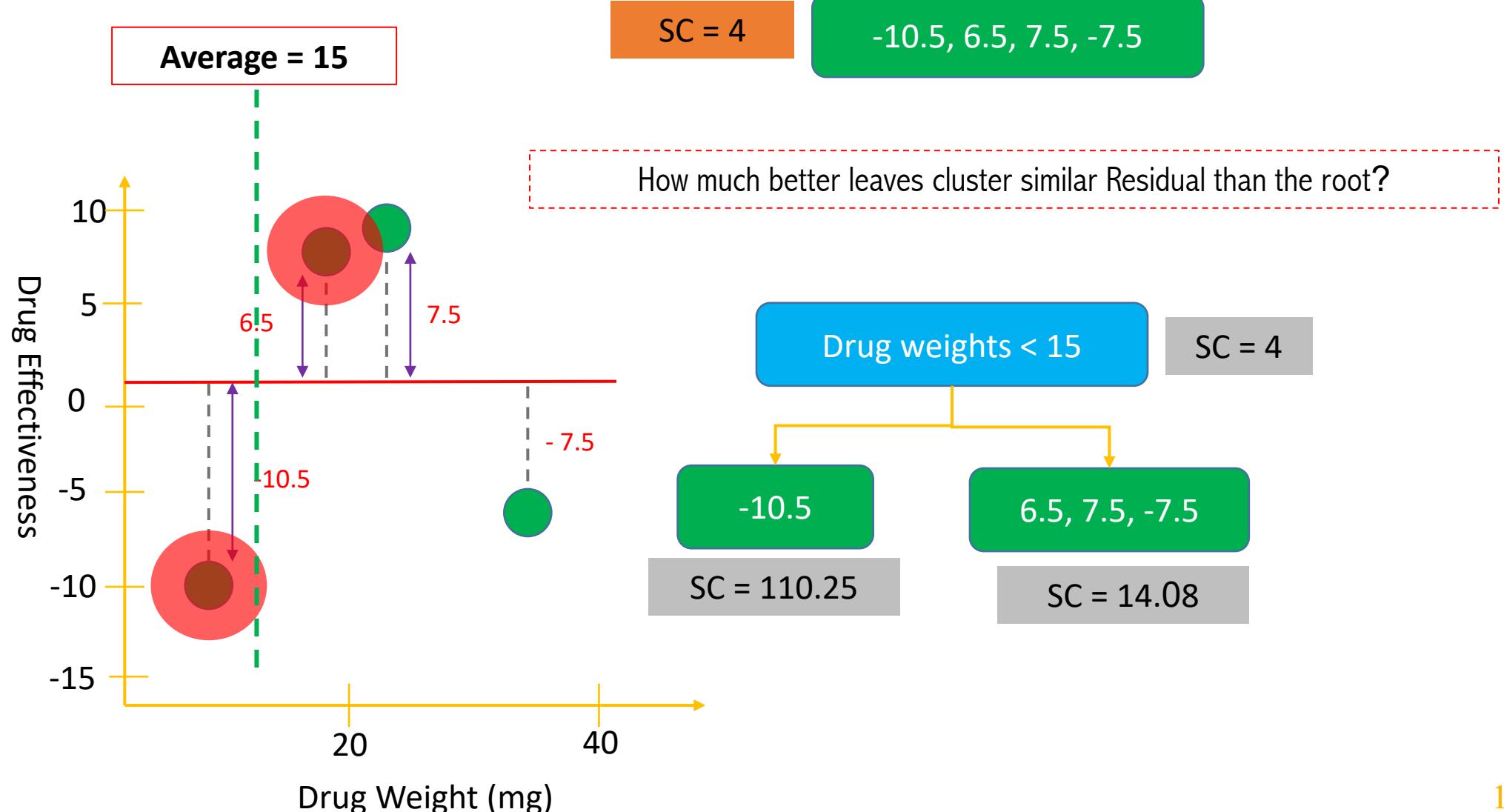
XGBoost For Regression

Step 1



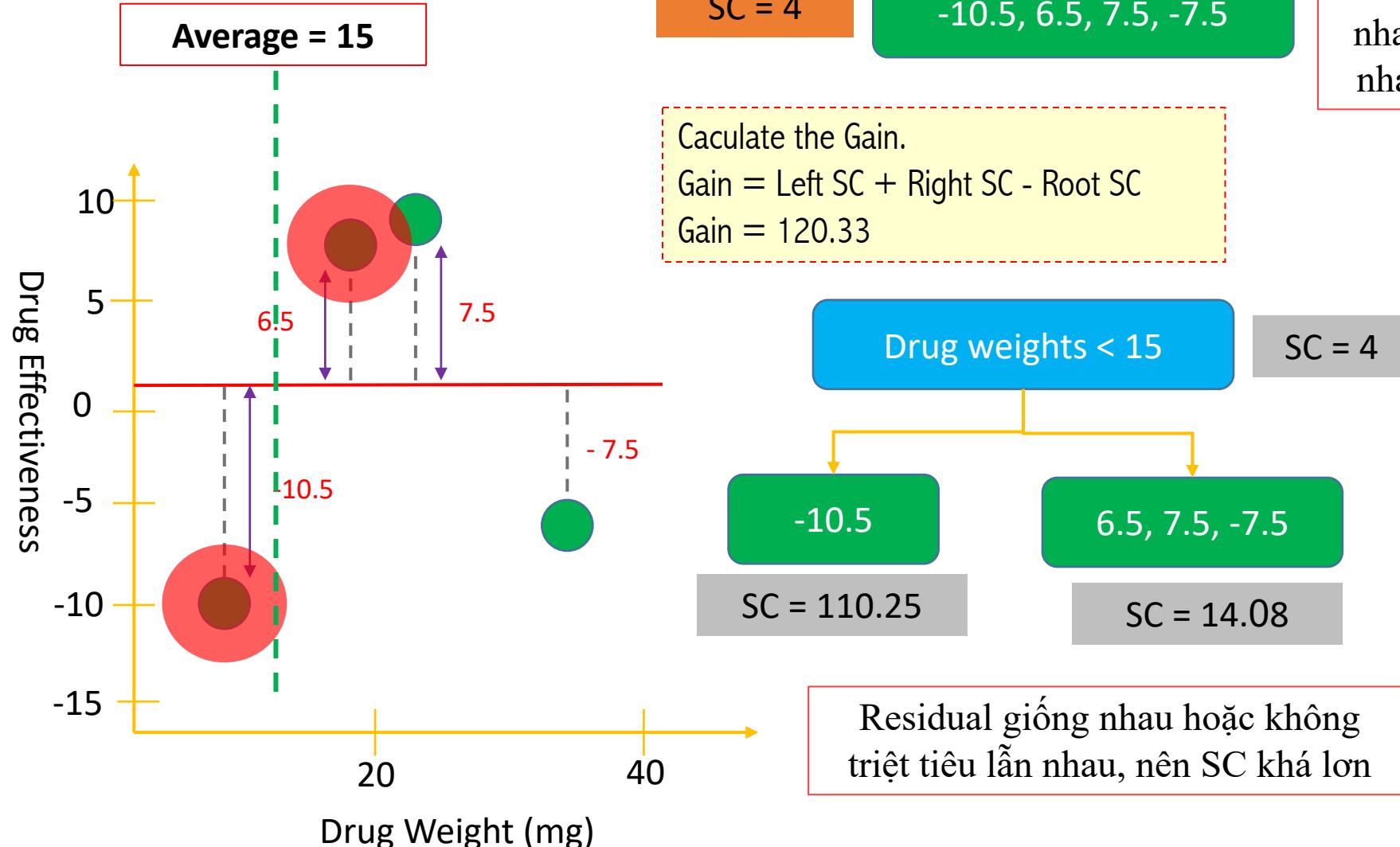
XGBoost For Regression

Step 1



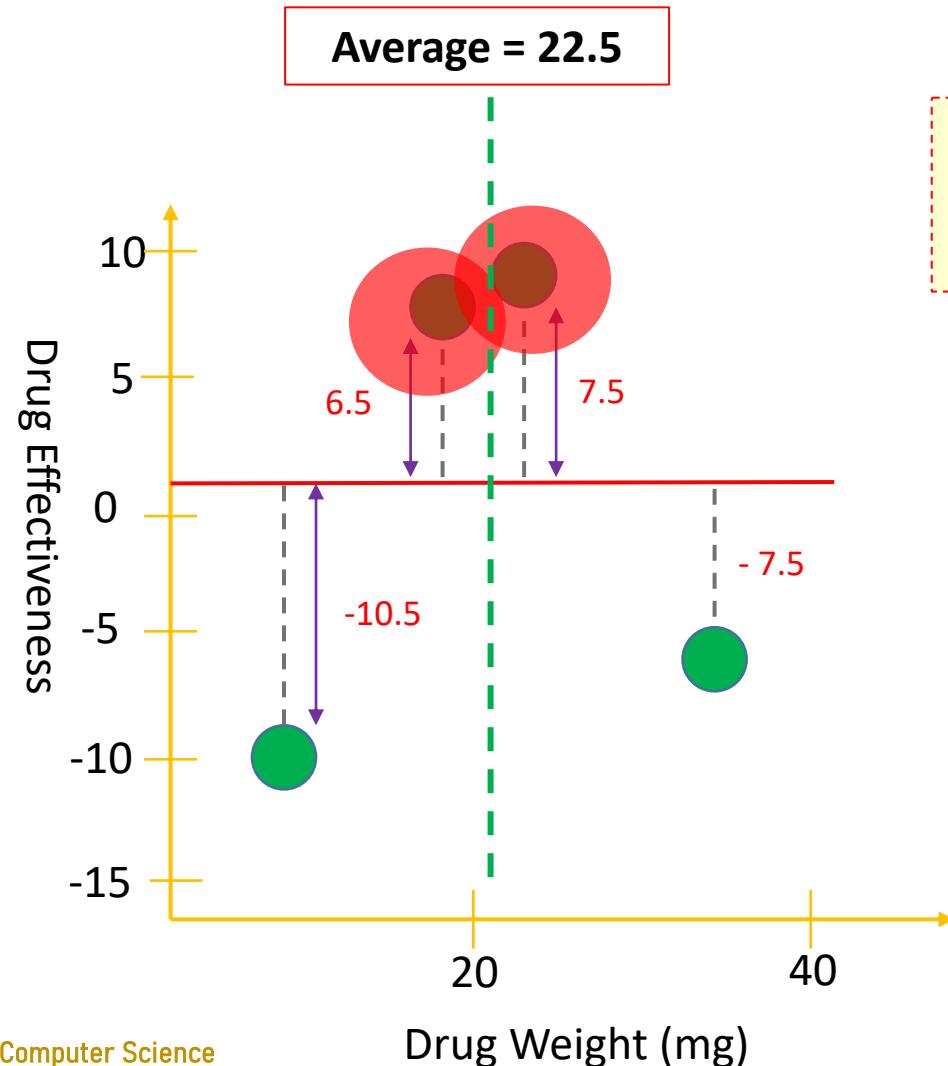
XGBoost For Regression

Step 1



XGBoost For Regression

Step 1



SC = 4

-10.5, 6.5, 7.5, -7.5

Residual rất khác nhau, triệt tiêu lẫn nhau, nên SC nhỏ

Caculate the Gain.

Gain = Left SC + Right SC - Root SC

Gain = 4.0

Drug weights < 22.5

SC = 4

-10.5, 6.5

SC = 8

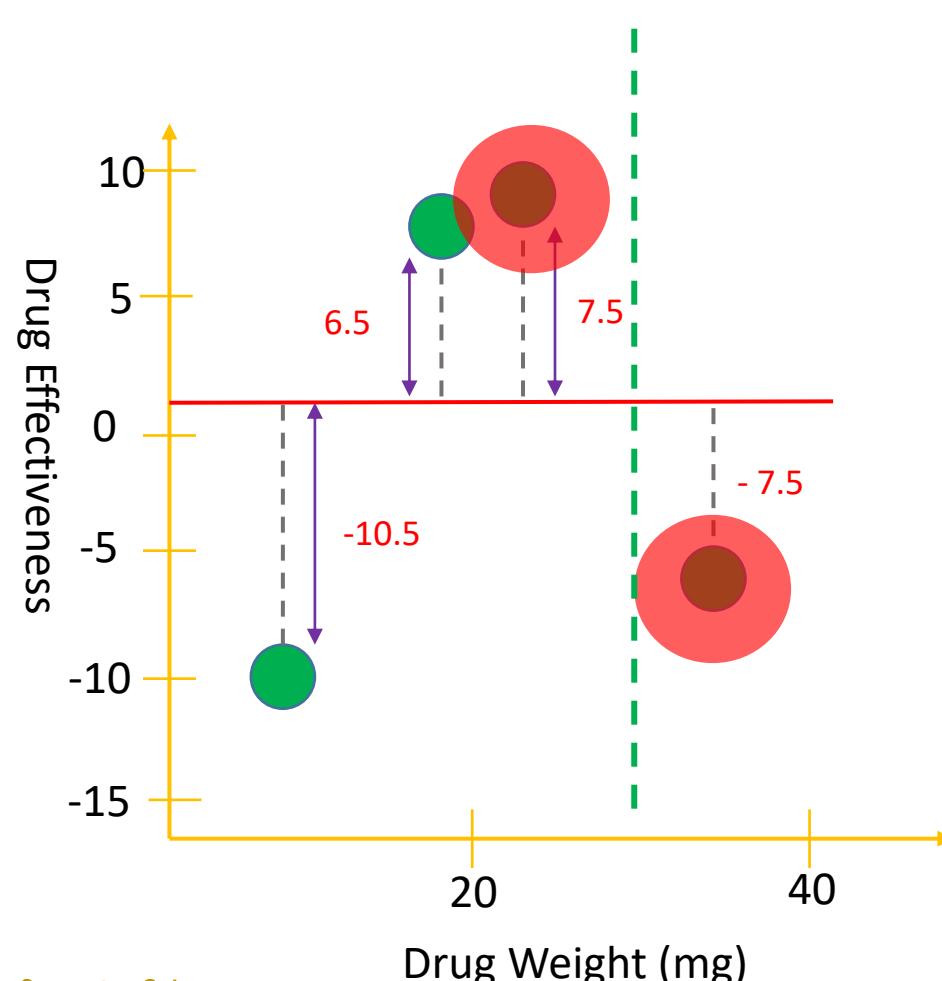
7.5, -7.5

SC = 0

Residual giống nhau hoặc không triệt tiêu lẫn nhau, nên SC khá lớn

XGBoost For Regression

Step 1



Average = 30

SC = 4

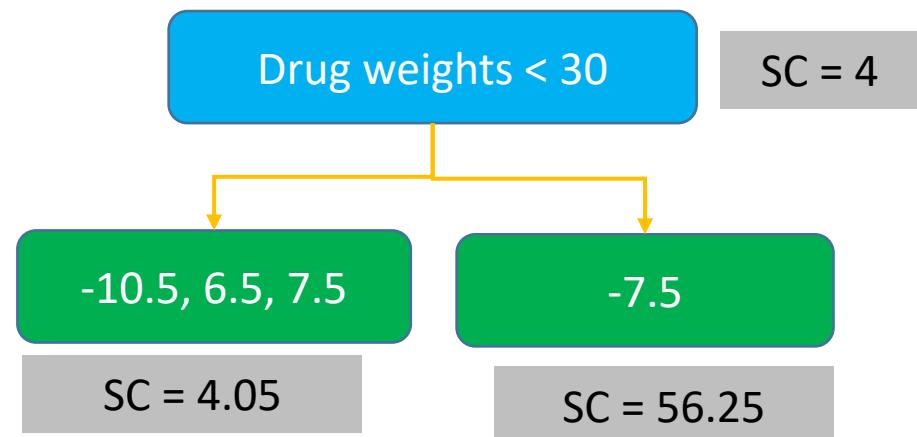
-10.5, 6.5, 7.5, -7.5

Residual rất khác nhau, triệt tiêu lẫn nhau, nên SC nhỏ

Caculate the Gain.

Gain = Left SC + Right SC - Root SC

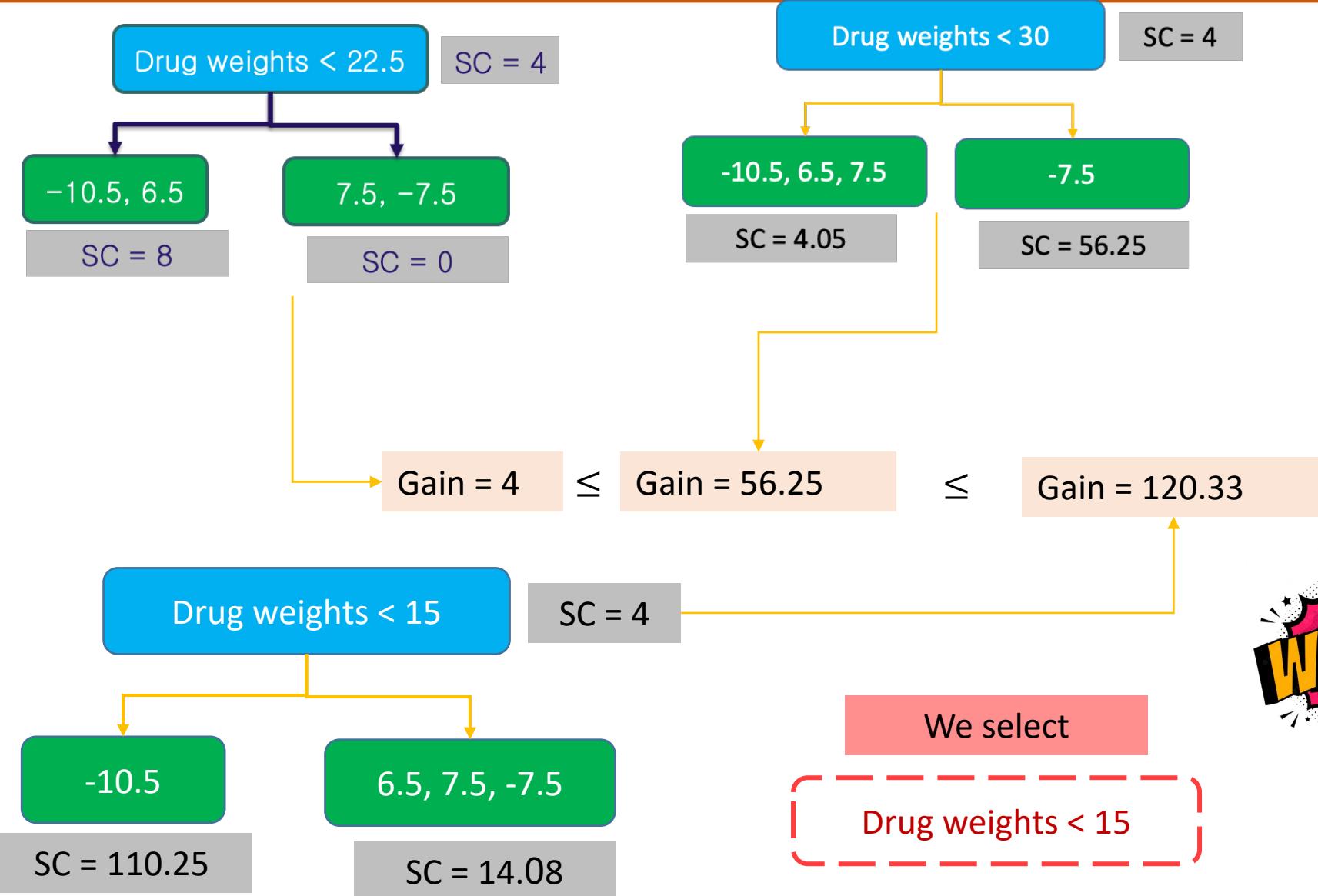
Gain = 56.33



Residual giống nhau hoặc không triệt tiêu lẫn nhau, nên SC khá lớn

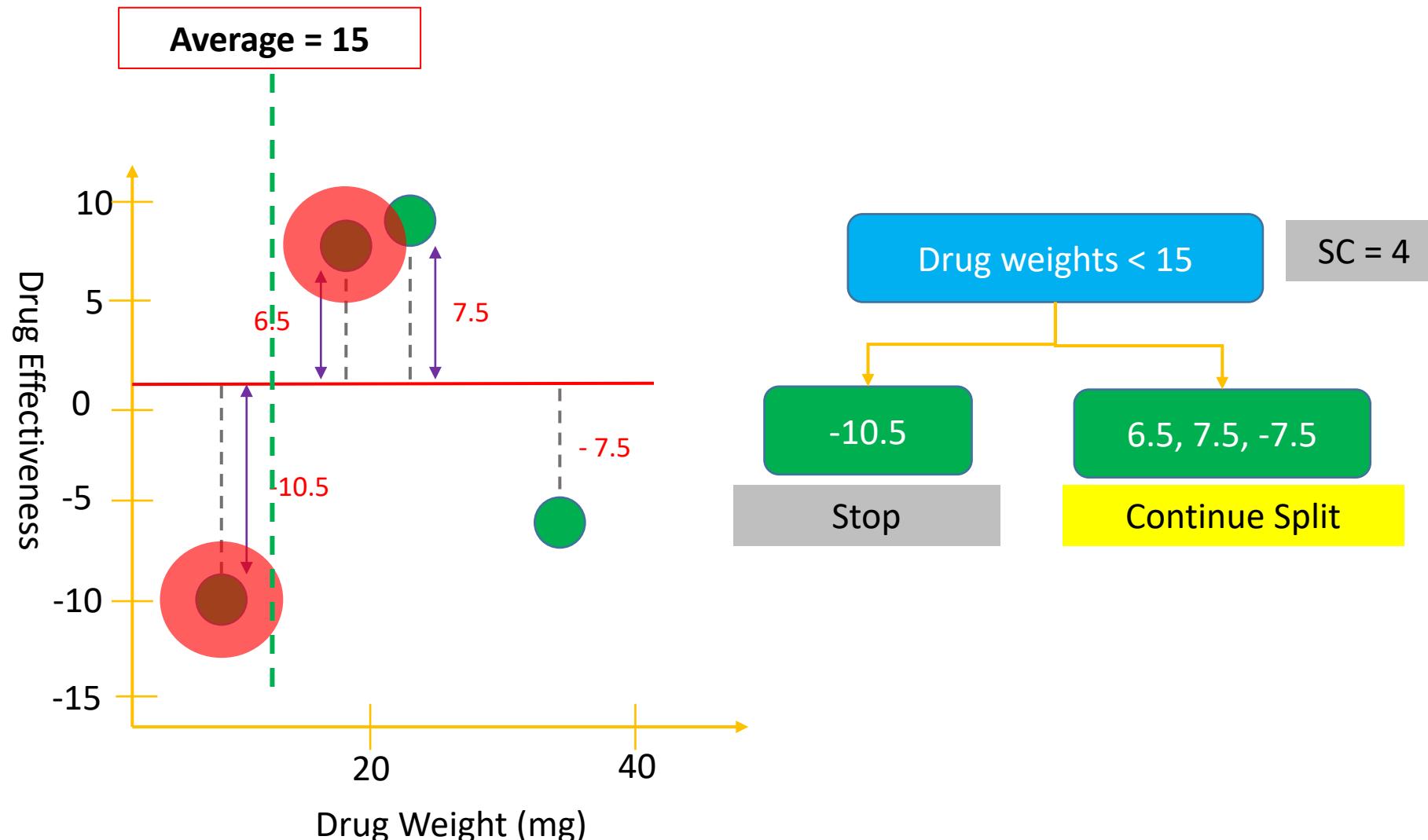
XGBoost For Regression

Step 1



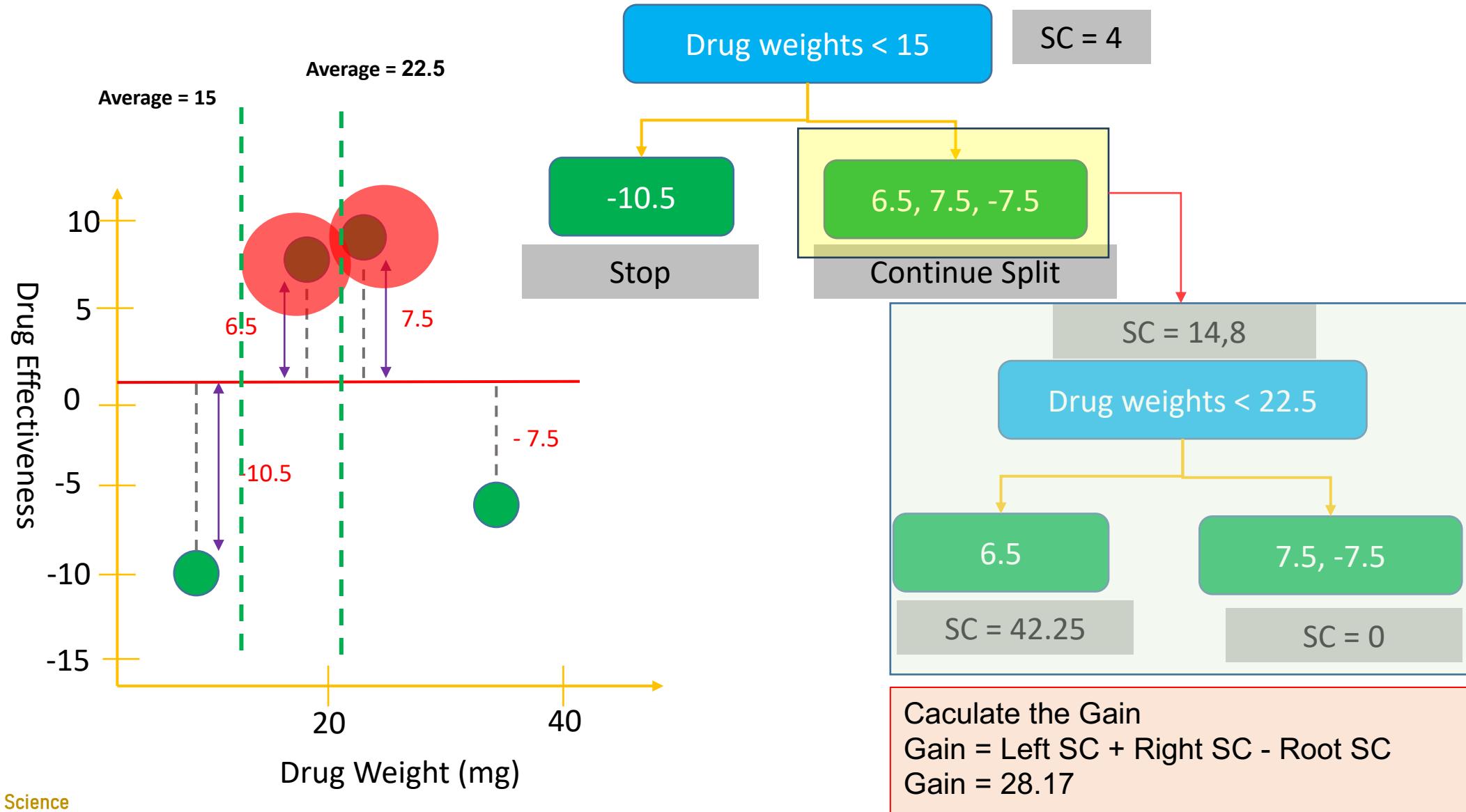
XGBoost For Regression

Step 1



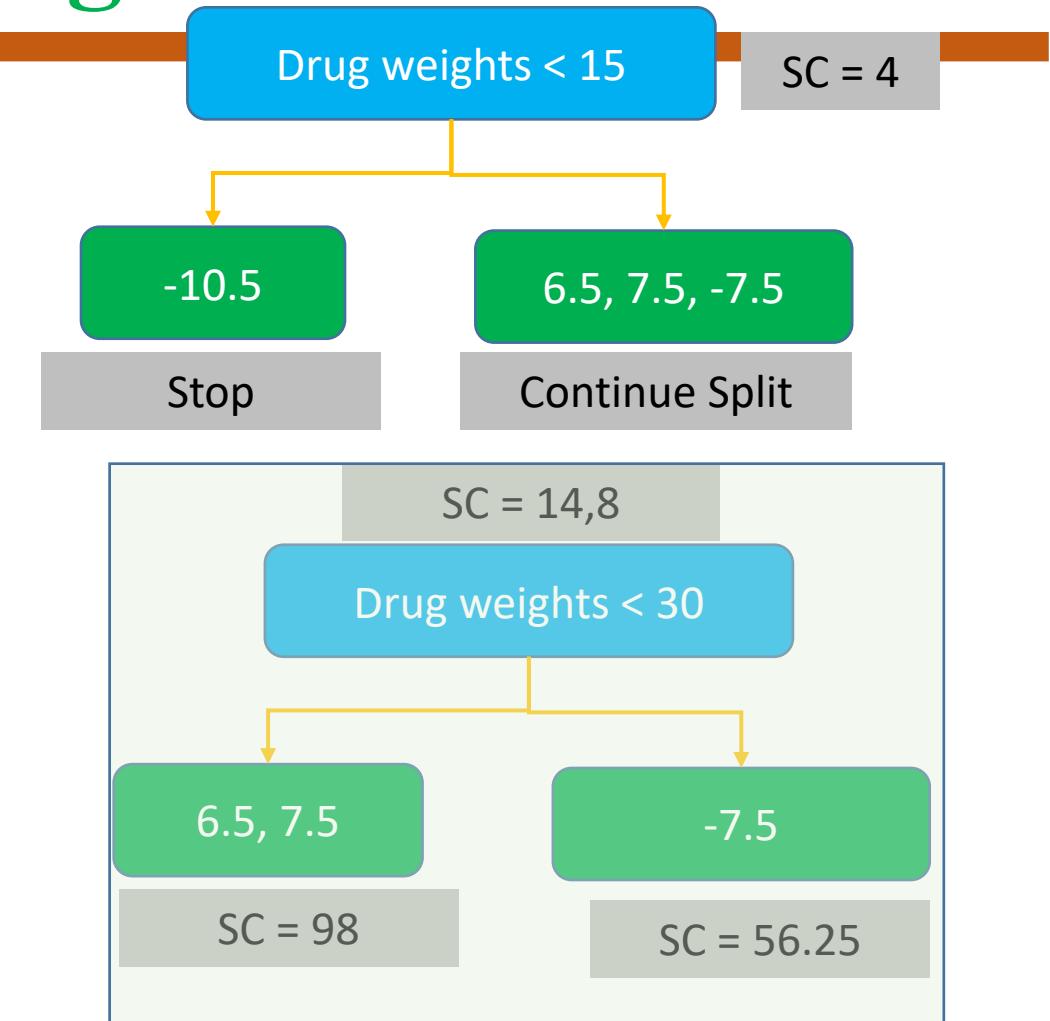
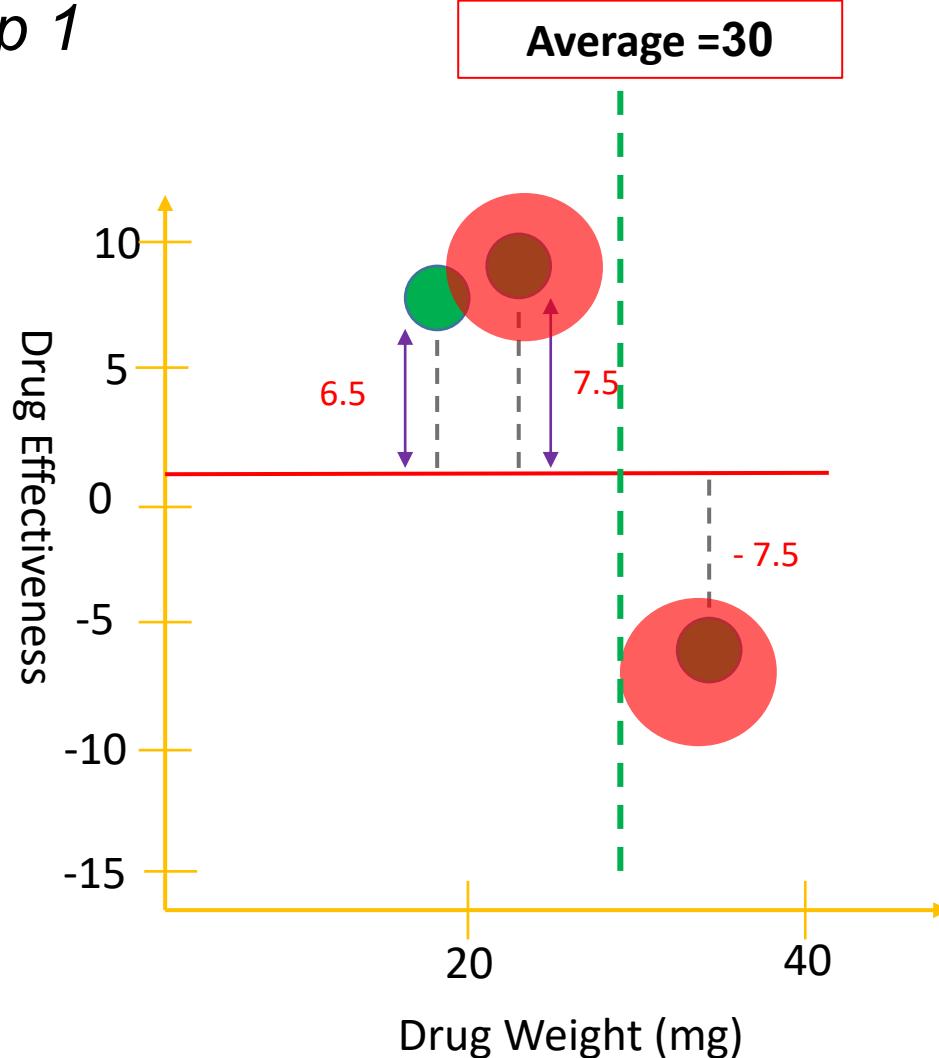
XGBoost For Regression

Step 1



XGBoost For Regression

Step 1

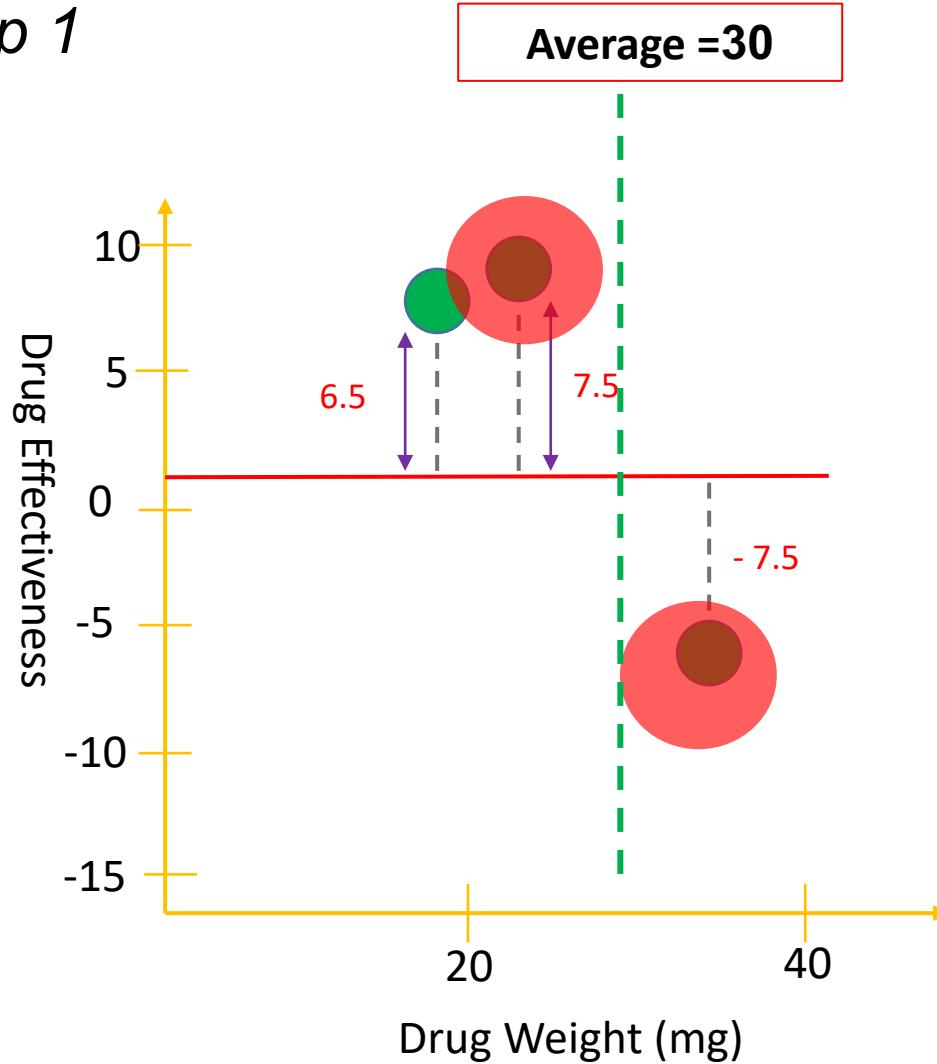


Calculate the Gain
Gain = Left SC + Right SC - Root SC
Gain = 140.17

XGBoost For Regression

Gain = 120.33

Step 1



Drug weights < 15

-10.5

Stop

SC = 140.7

Drug weights < 30

6.5, 7.5

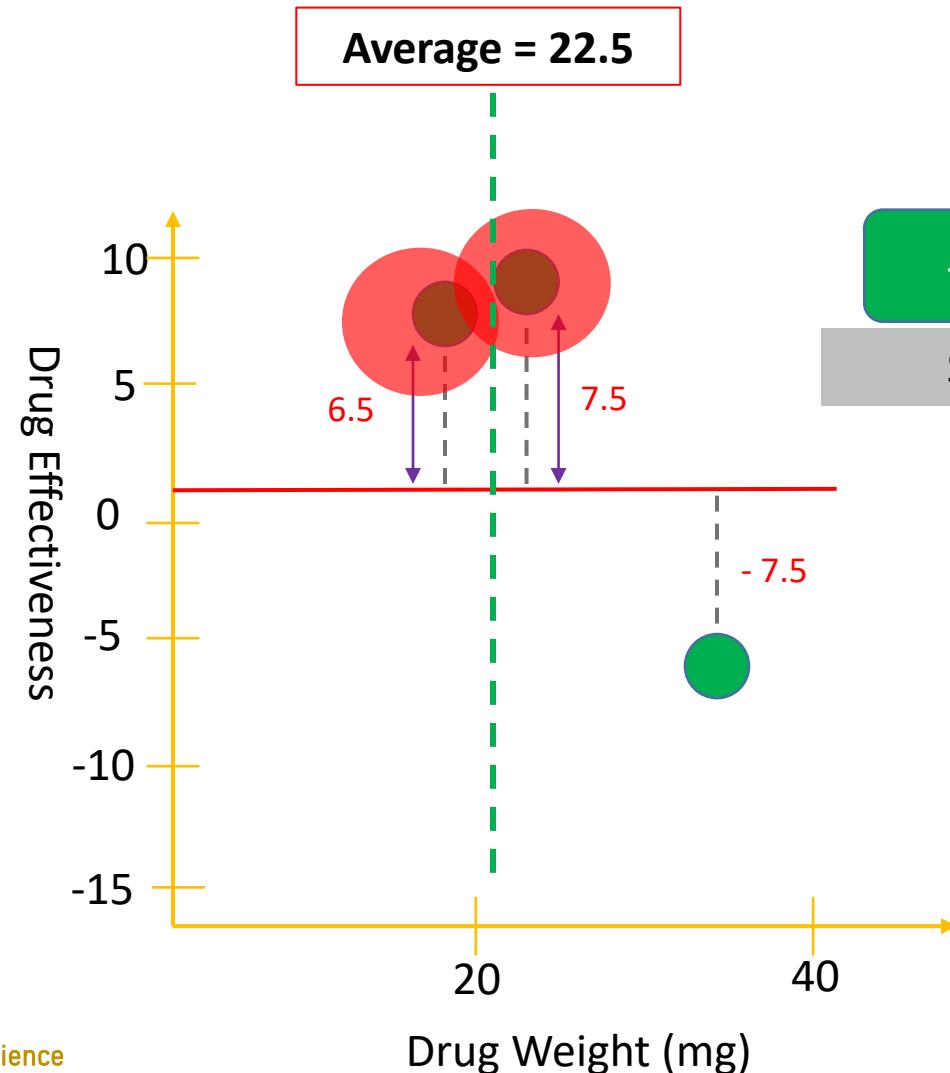
-7.5

How to prune the tree to prevent
Overfitting ? Gain information

$\gamma = 130$

XGBoost For Regression

Step 1

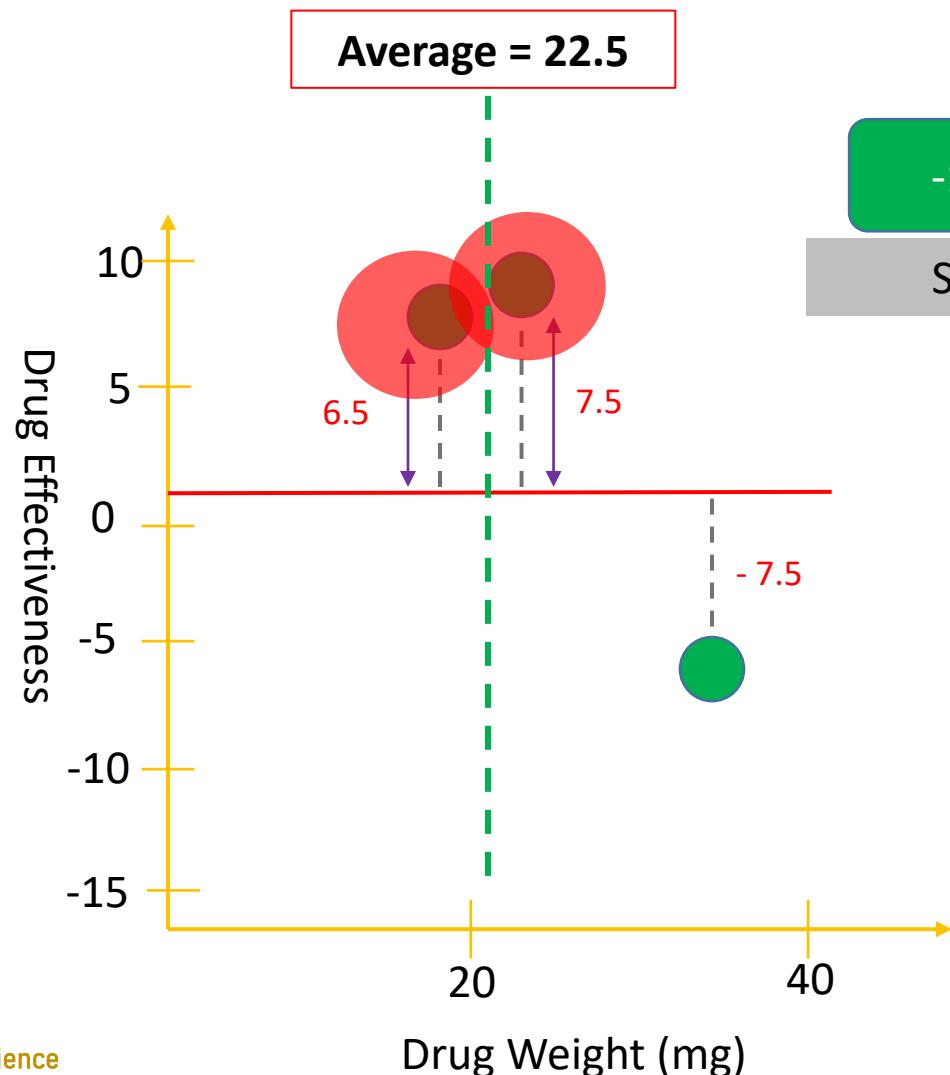


How to prune the tree to prevent Overfitting ? Gain information

$\gamma = 130$
Difference = Gain - γ
If difference > 0, do not remove branch
If difference < 0, remove branch

XGBoost For Regression

Step 1



Gain = 120.33

Drug weights < 15

SC = 140.7

-10.5

Stop

6.5, 7.5, -7.5

How to prune the tree to prevent Overfitting ? Gain information

$\gamma = 150$

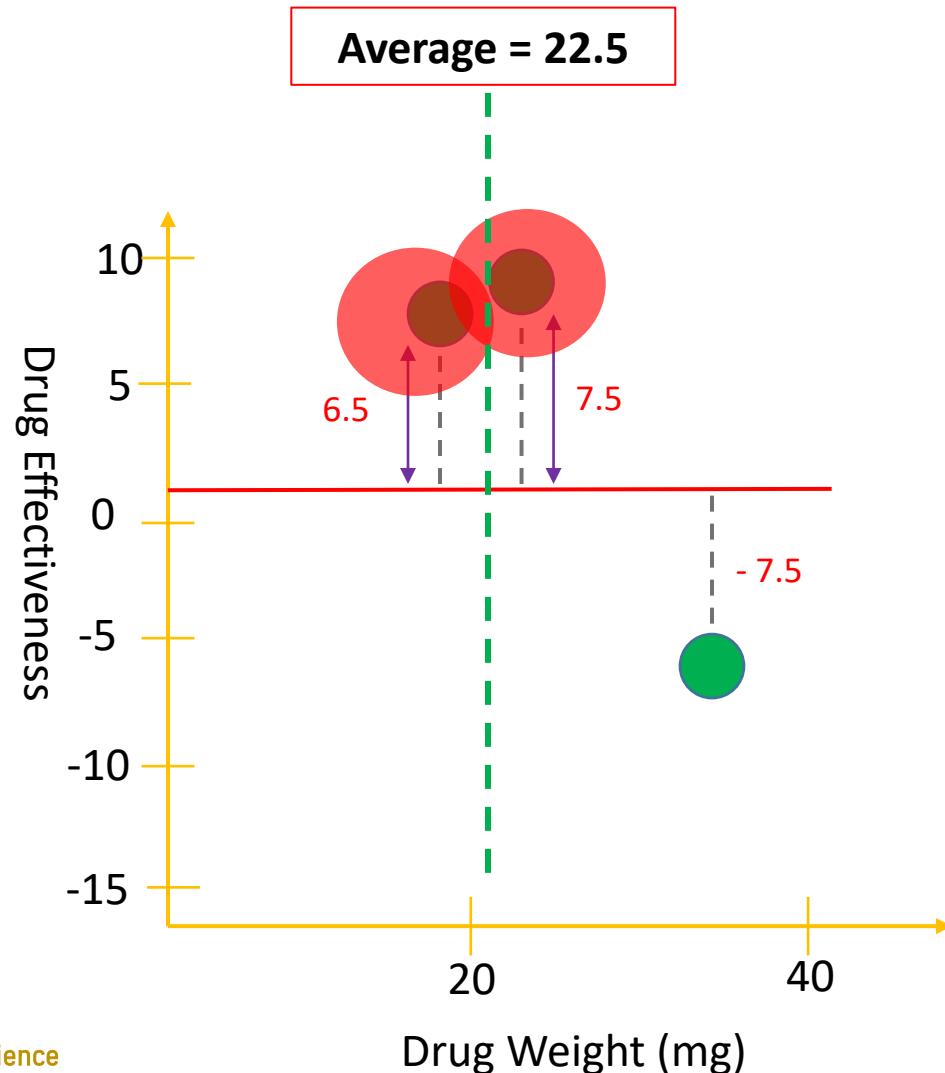
Difference = Gain - γ

If difference > 0, do not remove branch

If difference < 0, remove branch

XGBoost For Regression

Step 1

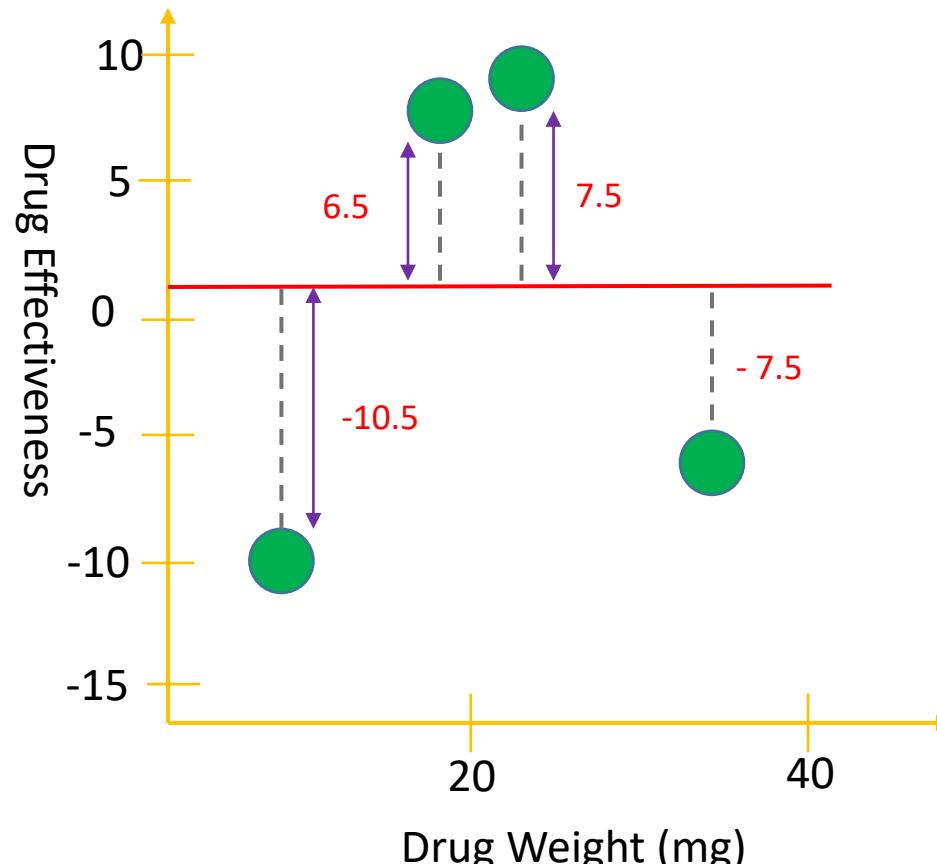


How to prune the tree to prevent Overfitting ? Gain information

$\gamma = 150$
Difference = Gain - γ
If difference > 0, do not remove branch
If difference < 0, remove branch

XGBoost For Regression

Step 1



Start with single Leaf
of residuals

-10.5, 6.5, 7.5, -7.5

$m = 4$
 $\lambda = 1$

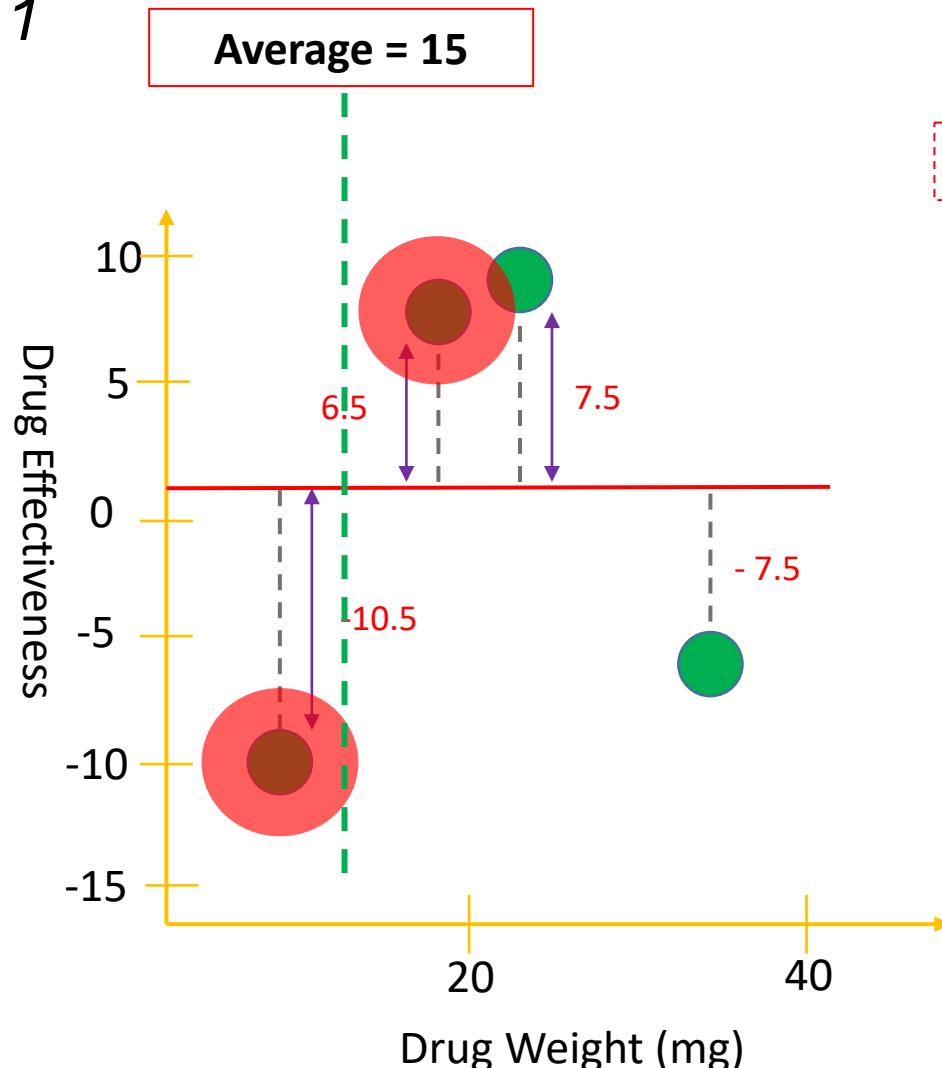
Compute Similarity Score

$$SC = \frac{\sum (output - predicted)^2}{m + \lambda}$$

$$SC = \frac{[-10.5 + 7.5 + 6.5 + (-7.5)]^2}{4+1} = 3.2$$

XGBoost For Regression

Step 1



SC = 3.2

-10.5, 6.5, 7.5, -7.5

Please look at the two outputs with lowest drug weights

SC = 3.2

SC = 4

Drug weights < 15

-10.5

SC = 55.12

SC = 110.25

6.5, 7.5, -7.5

SC = 10.56

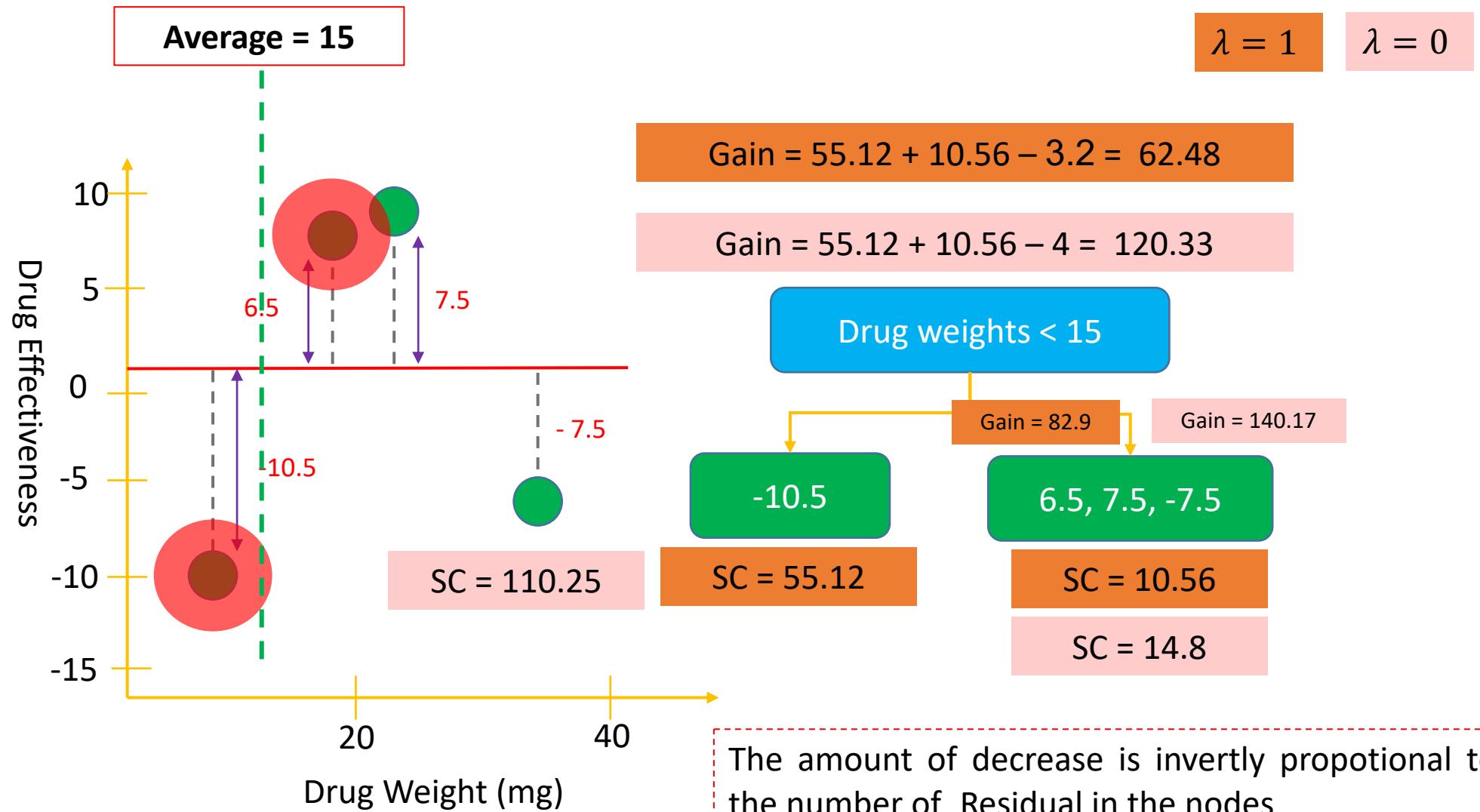
SC = 14.8

$\lambda = 1$

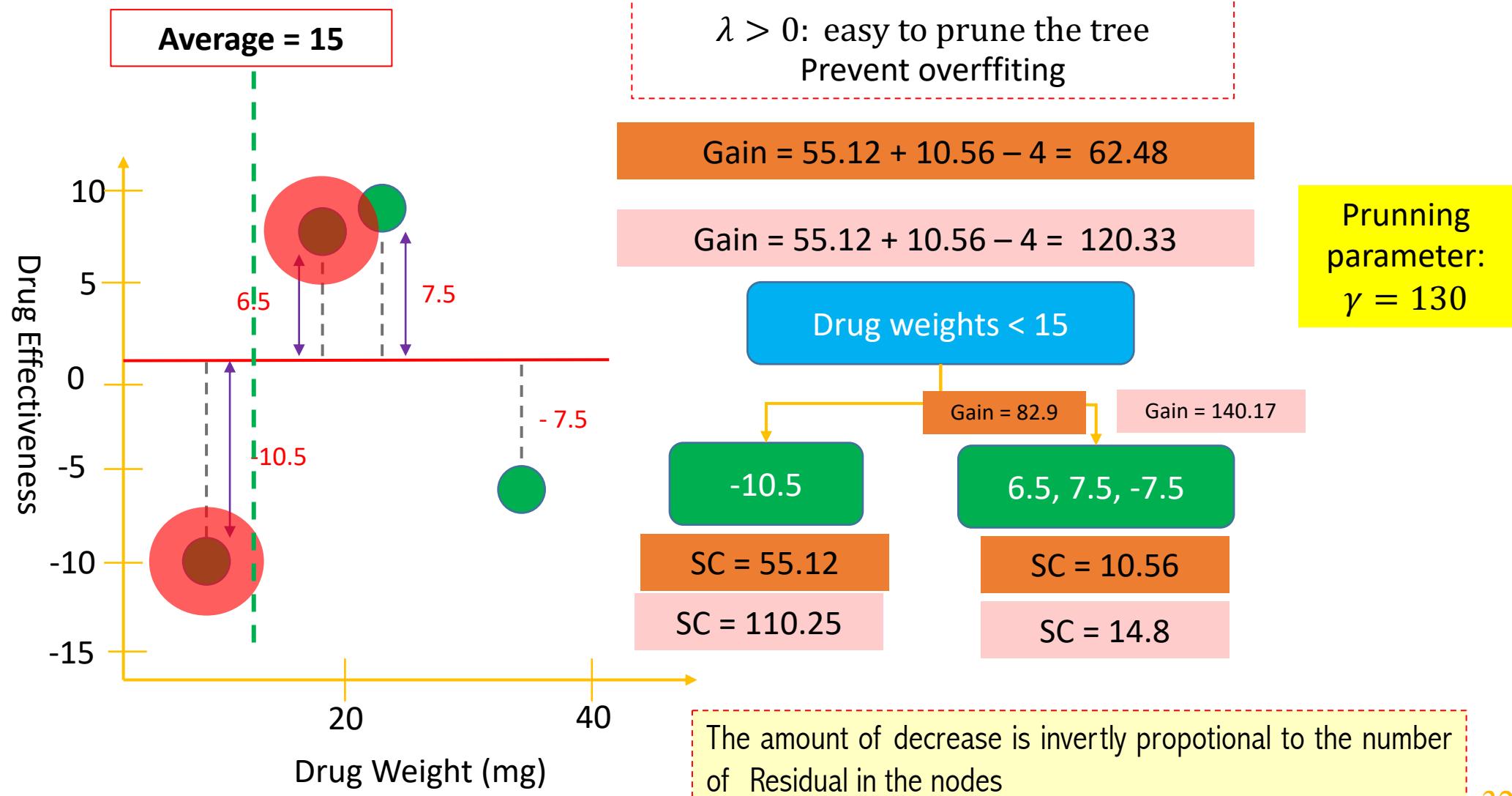
$\lambda = 0$

When $\lambda > 0$, the similarity score are smaller
Inversely proportional to the number of residuals

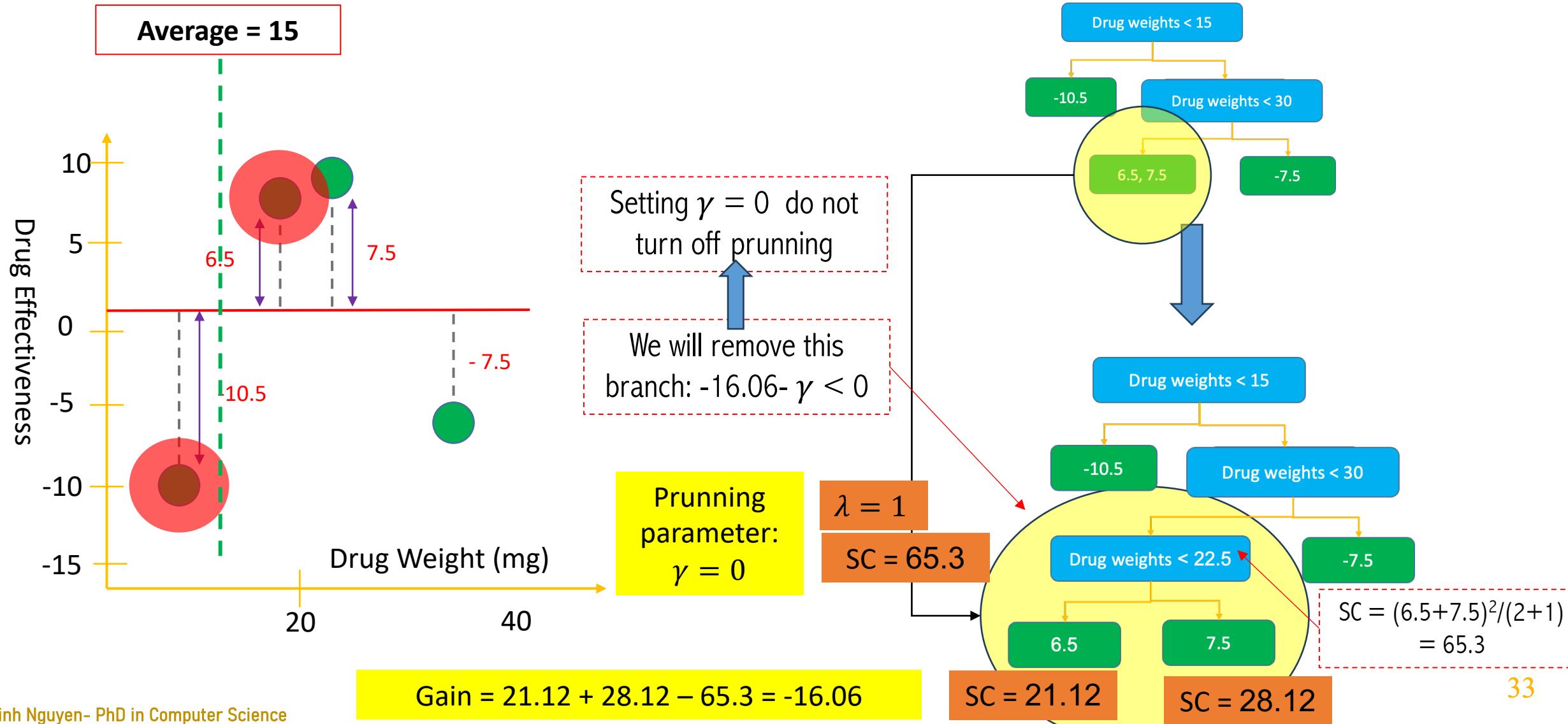
XGBoost For Regression



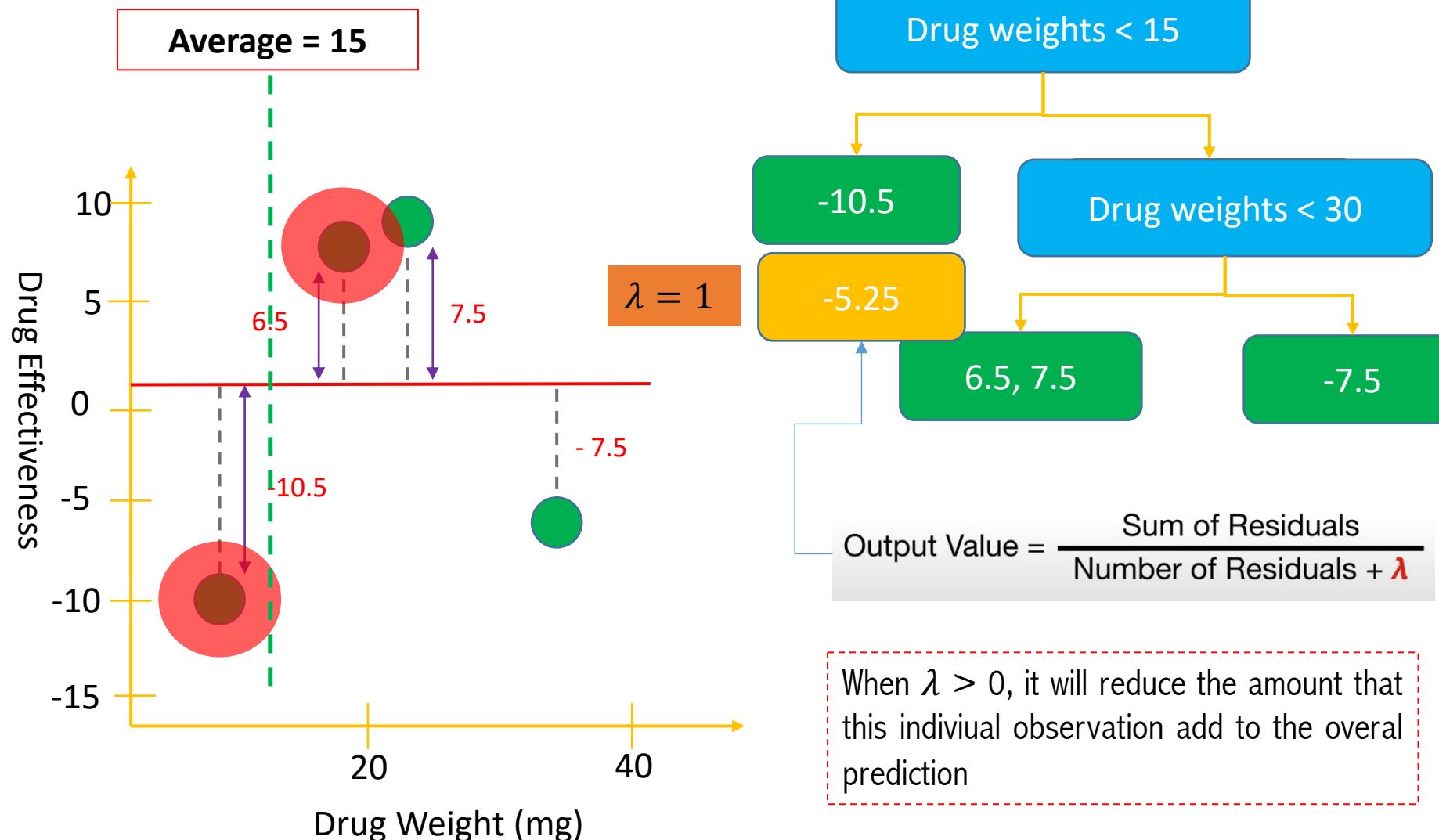
XGBoost For Regression



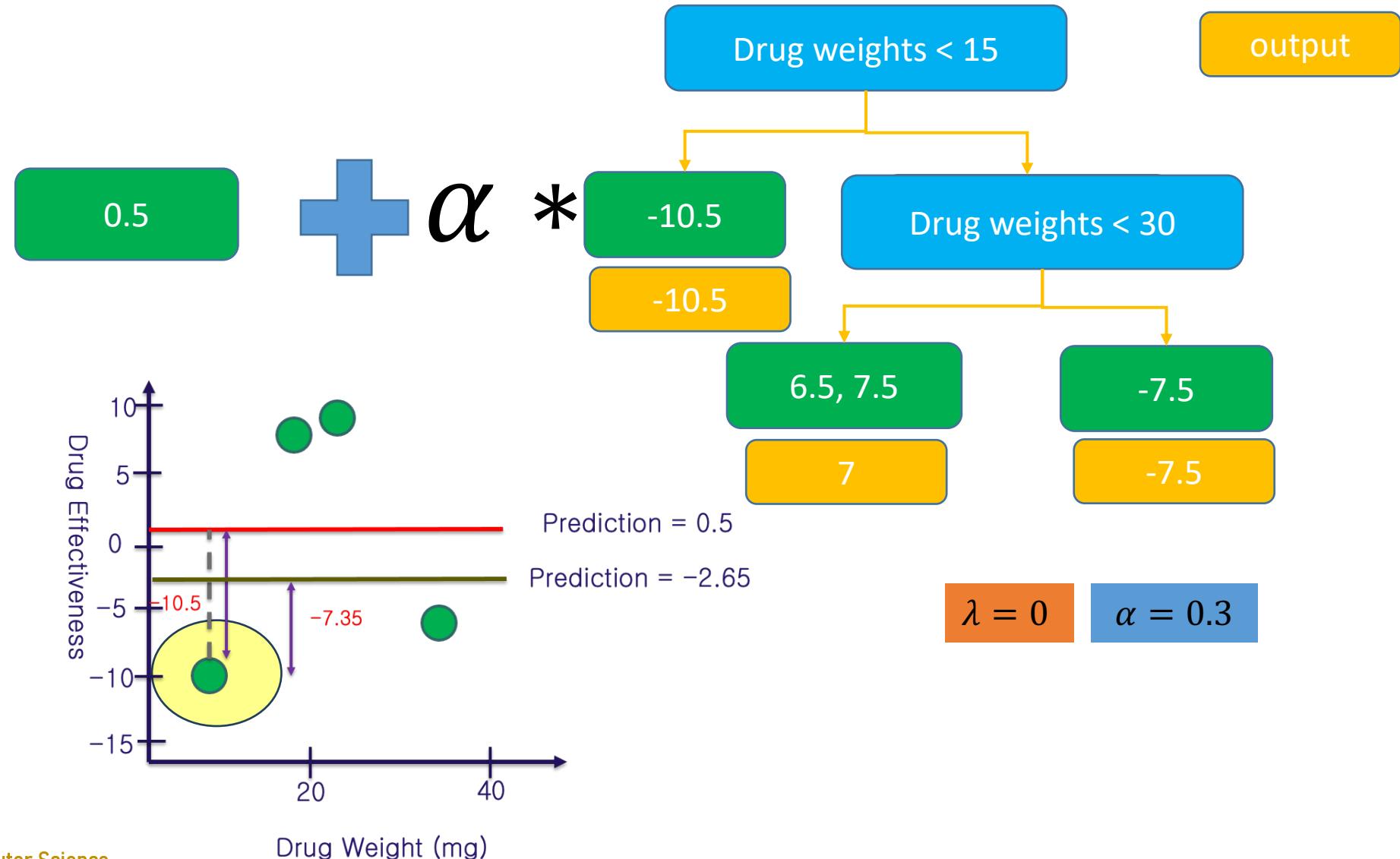
XGBoost For Regression



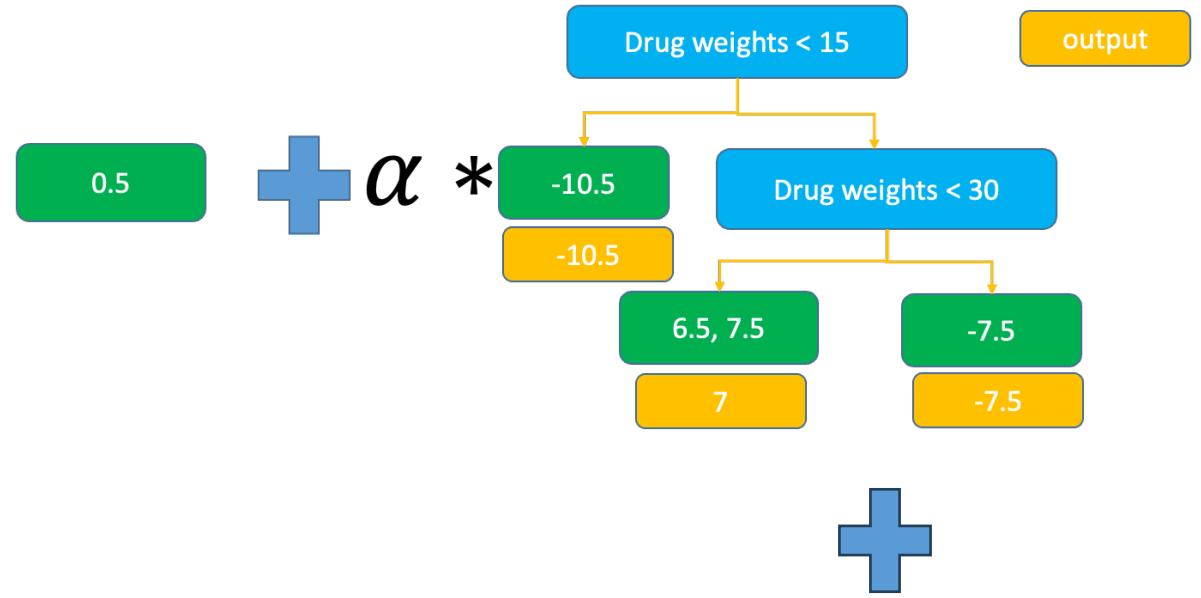
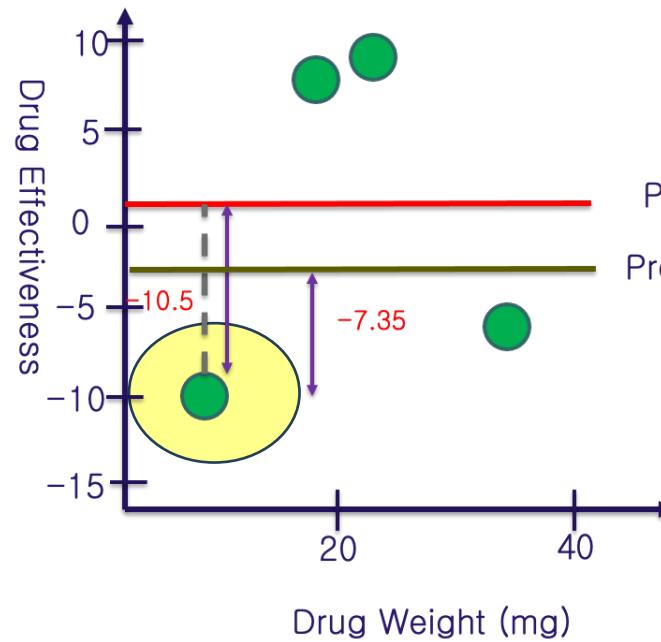
How to Predict Value



How to Predict Value



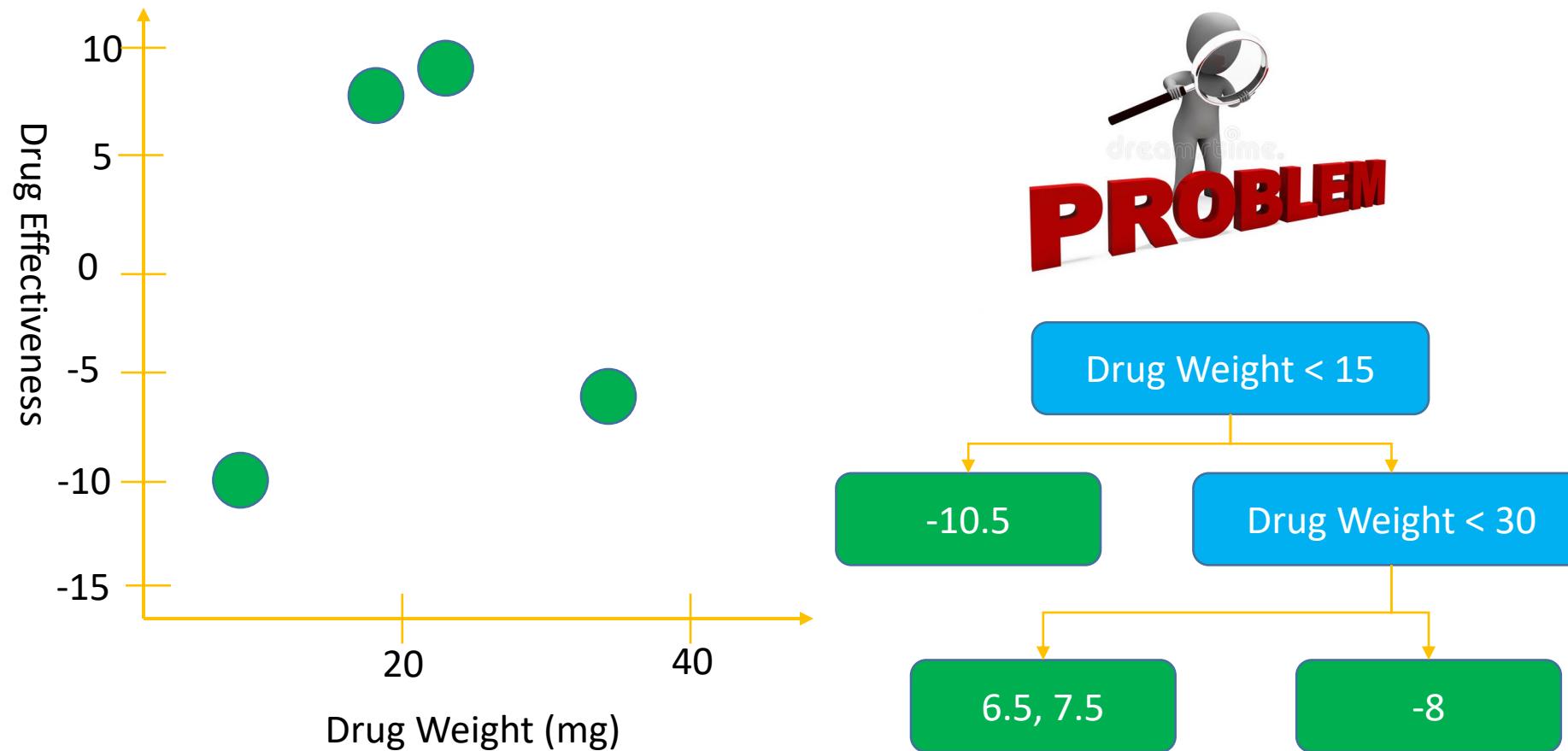
Building the Next Tree



$\alpha * \text{Next Tree Result}$

Keep bulding the Tree until the Residual are reach the predefined threshold. Or we reach to the maximum number of Tree

XGBoost for Regression

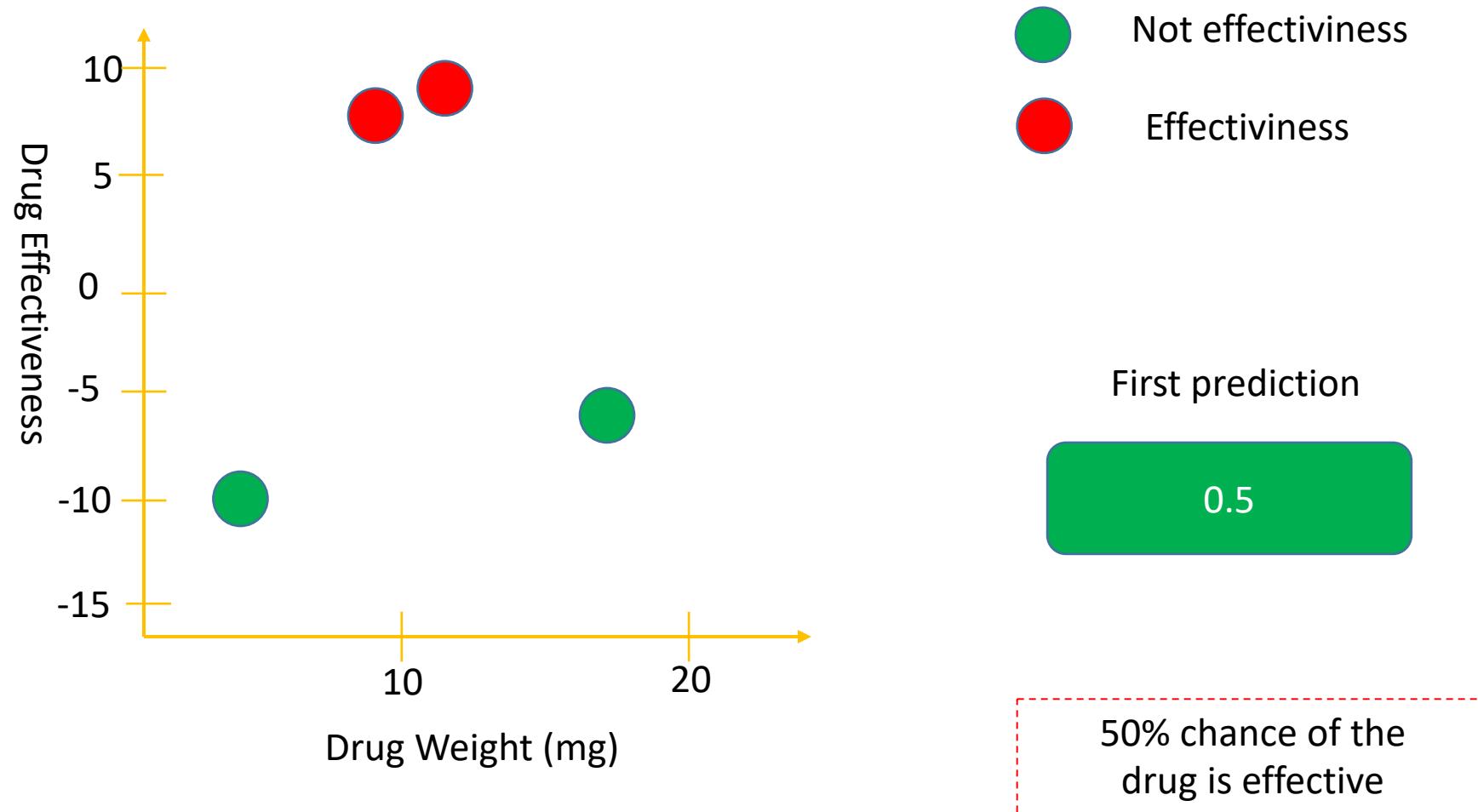


HOW TO FIND QUANTILES? => QUANTILE SKETCH APPROXIMATE SOLUTION

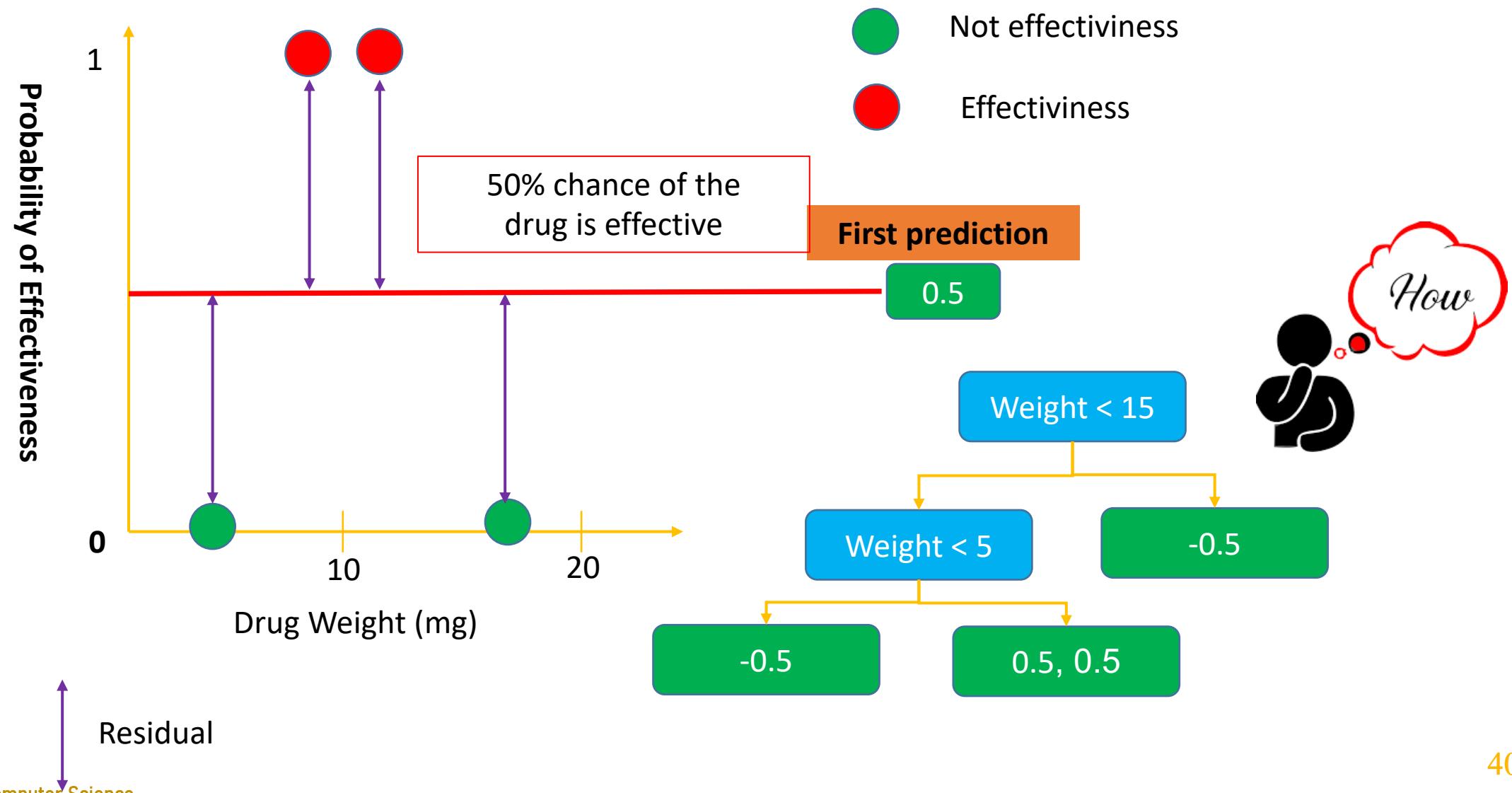
Outline

- **Regularization**
- **XGBoost For Regression**
- **XGBoost For Classification**
- **XGBoost: Mathematical Explanation**
- **How To Fill Missing Values**
- **Example**

XGBoost For Classification



XGBoost For Classification



XGBoost For Classification

Similarity Score for Classification:

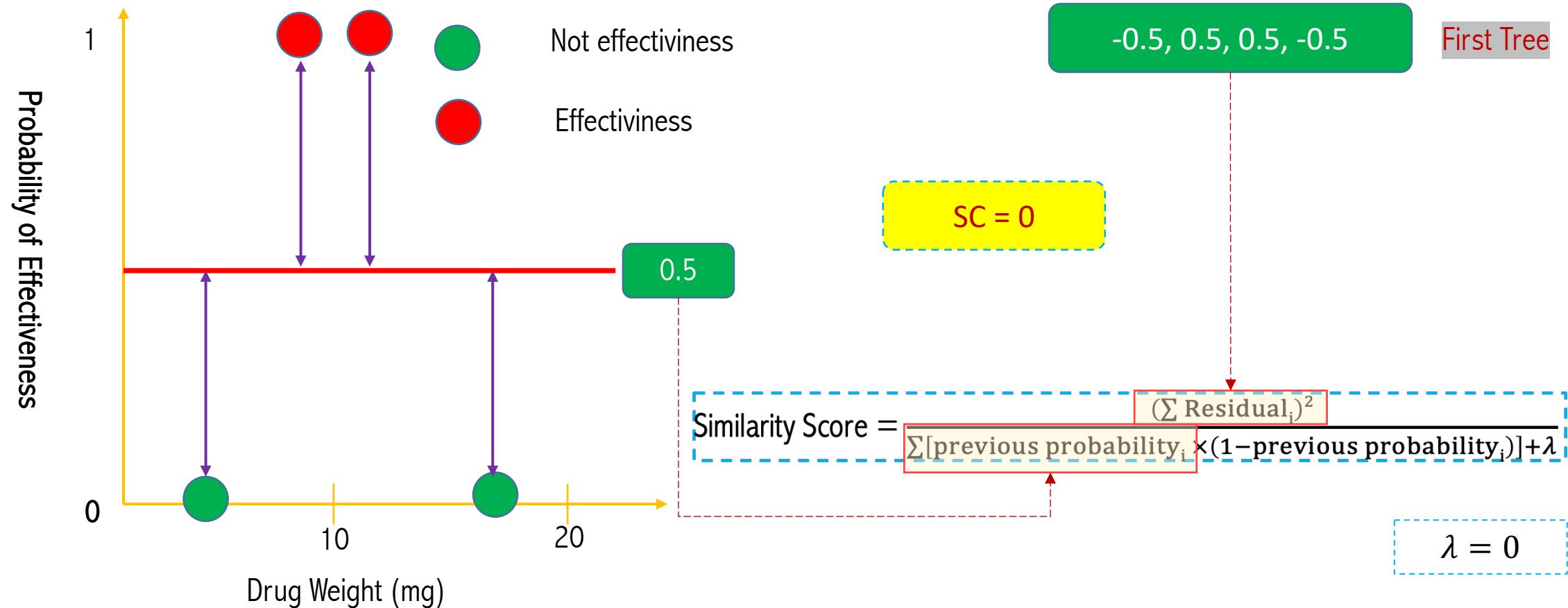
$$\text{Similarity Score} = \frac{(\sum \text{Residual}_i)^2}{\sum [\text{previous probability}_i \times (1 - \text{previous probability}_i)] + \lambda}$$

Similarity Score for Prediction (regression):

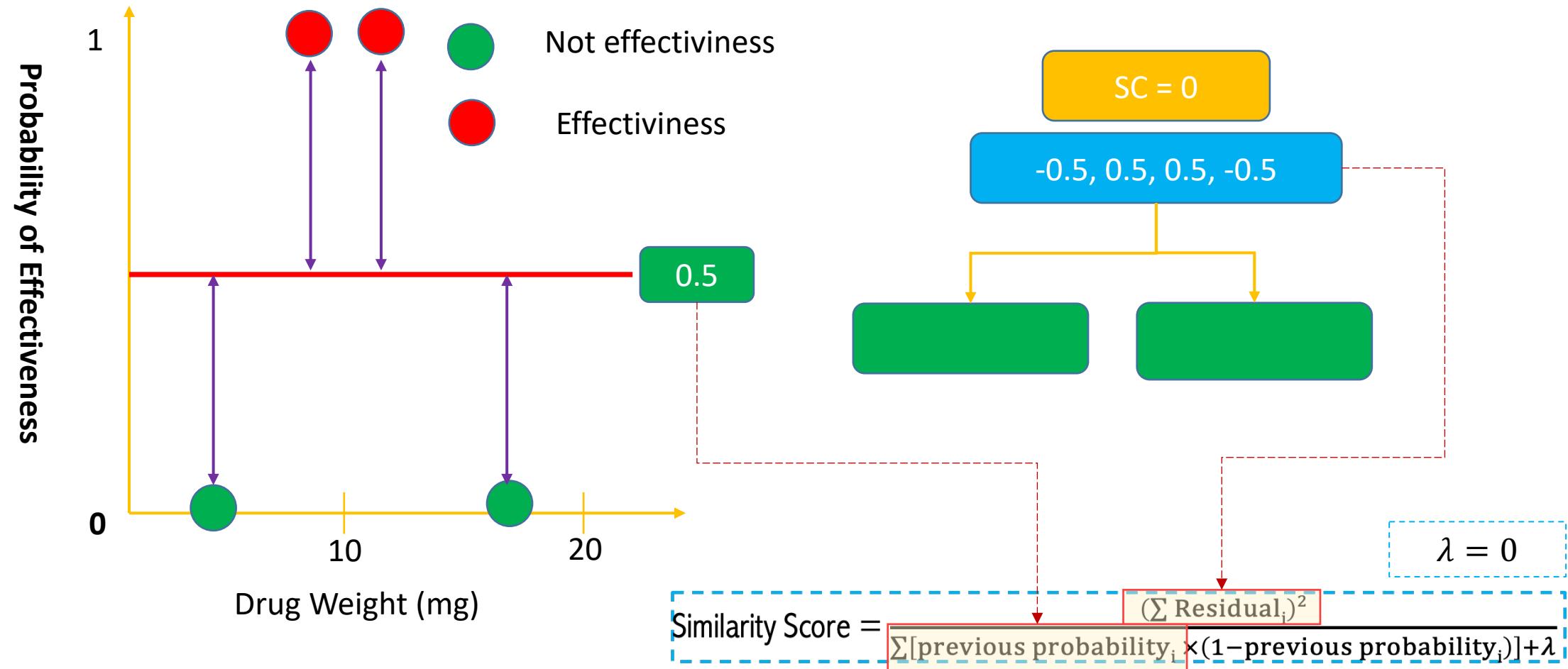
$$\text{Similarity Score} = \frac{(\sum \text{Residual}_i)^2}{\text{number of residual} + \lambda}$$



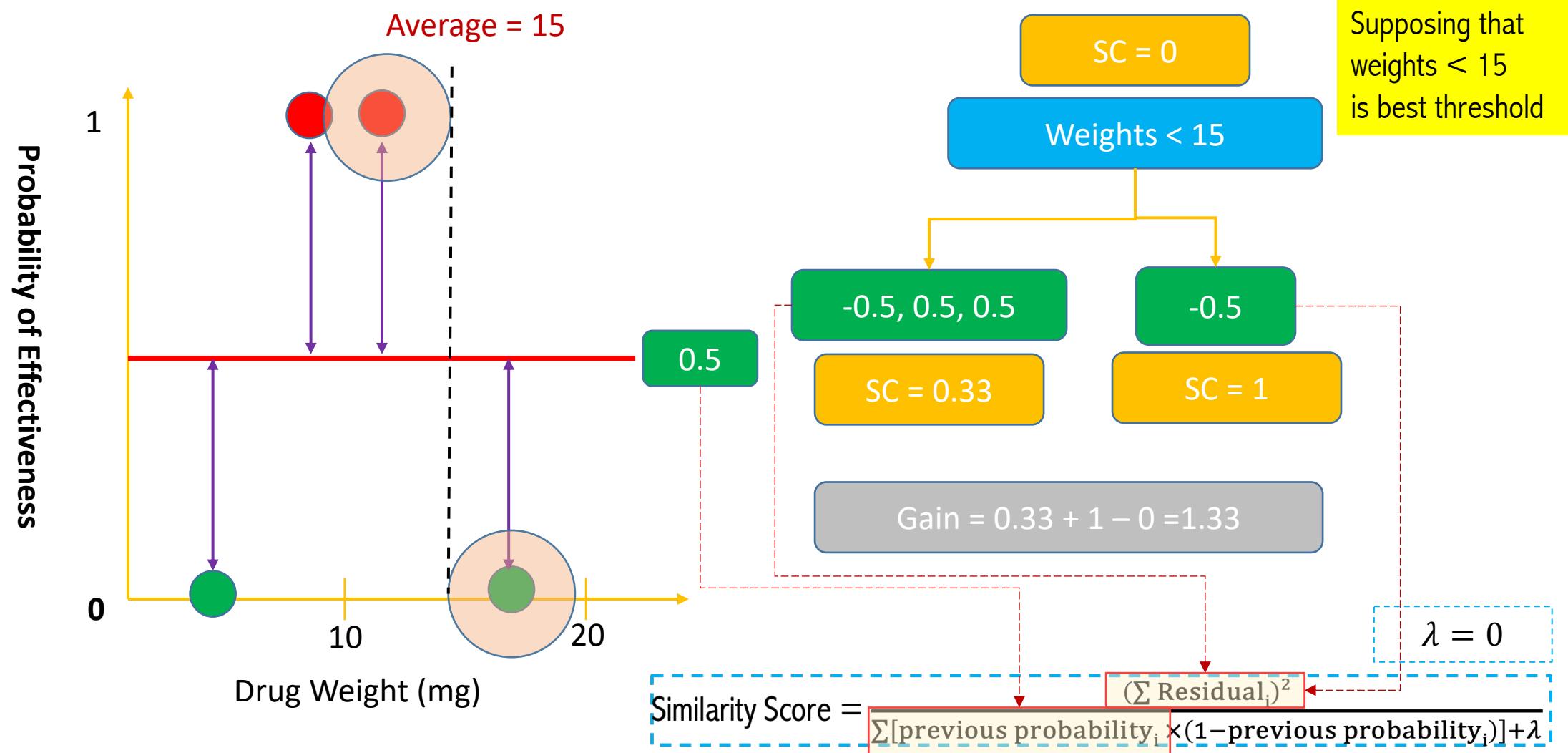
XGBoost for Classification



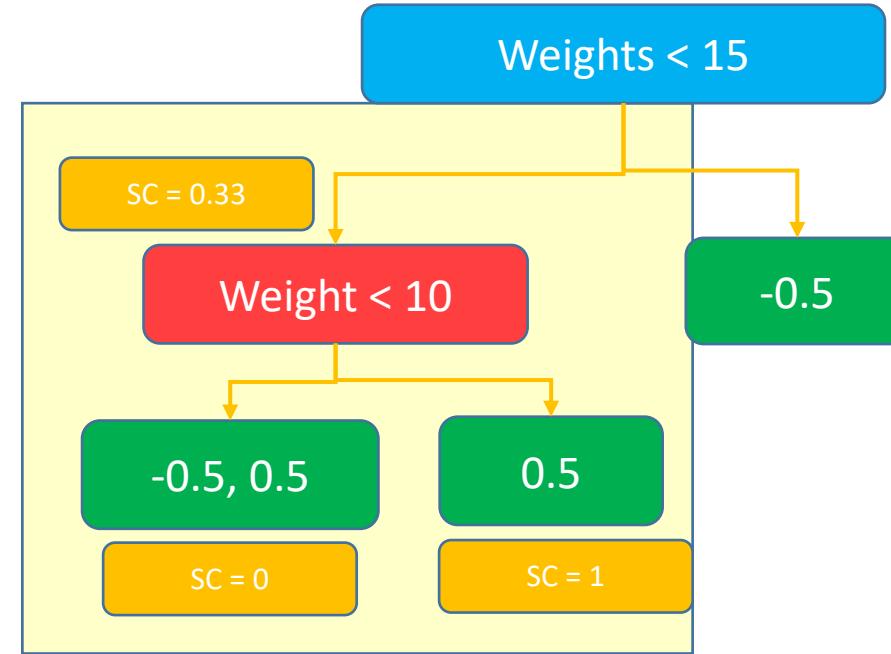
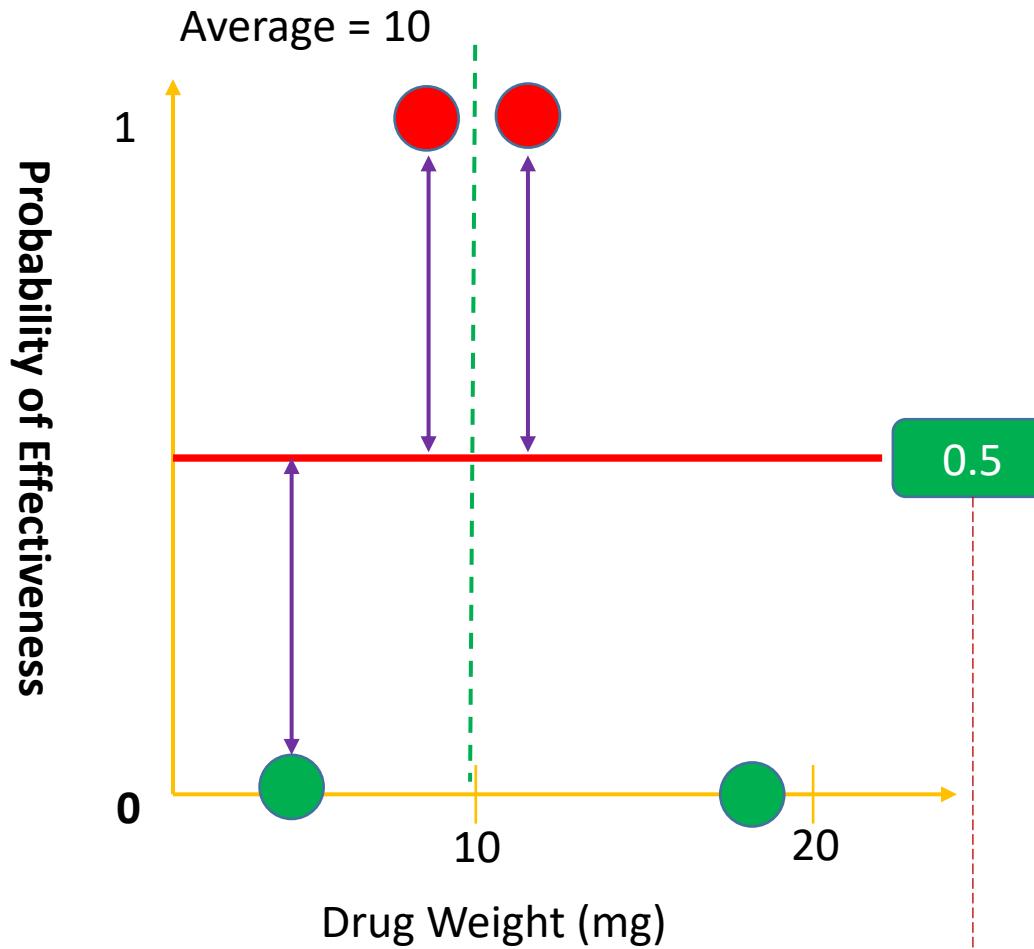
XGBoost For Classification



XGBoost For Classification



XGBoost For Classification

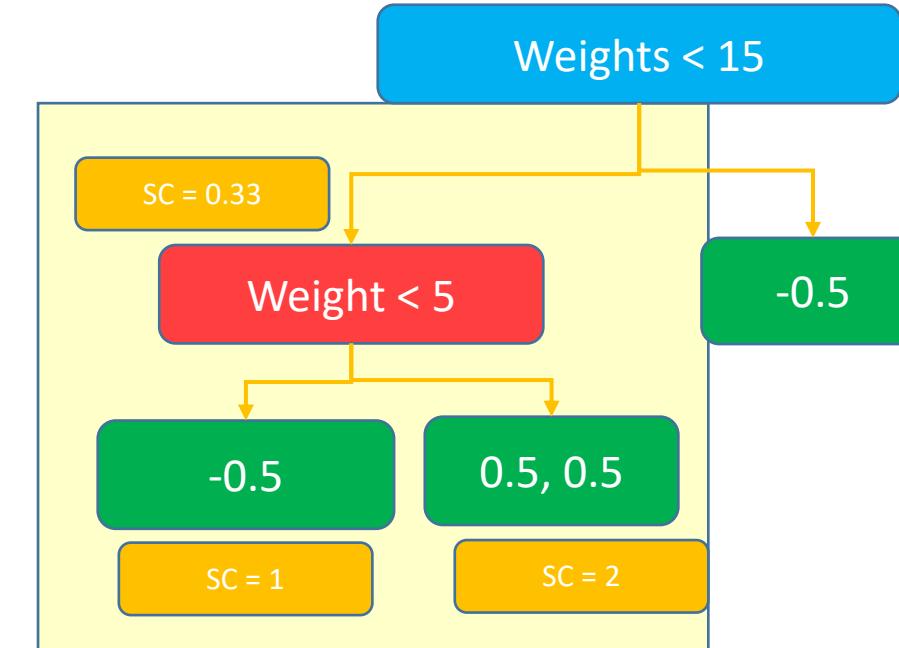
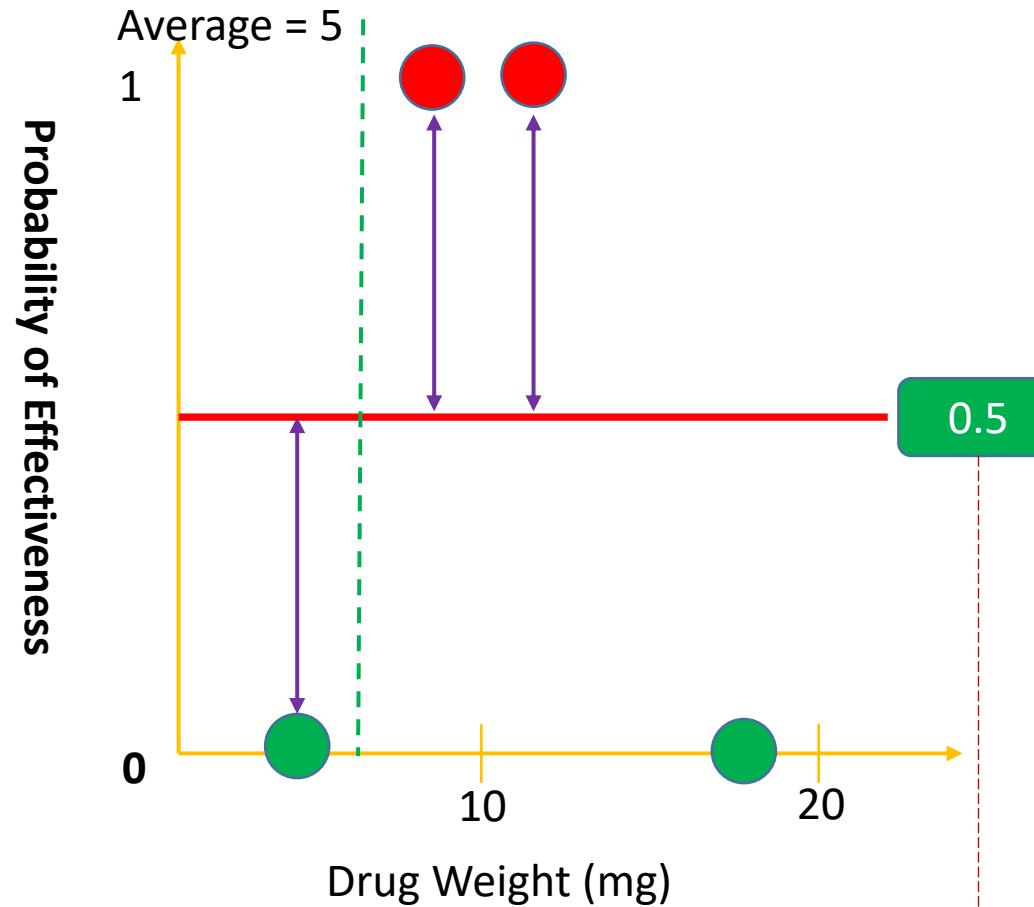


$$\text{Gain} = 0 + 1 - 0.33 = 0.66$$

$\lambda = 0$

Similarity Score =
$$\frac{(\sum \text{Residual}_i)^2}{\sum [\text{previous probability}_i \times (1 - \text{previous probability}_i)] + \lambda}$$

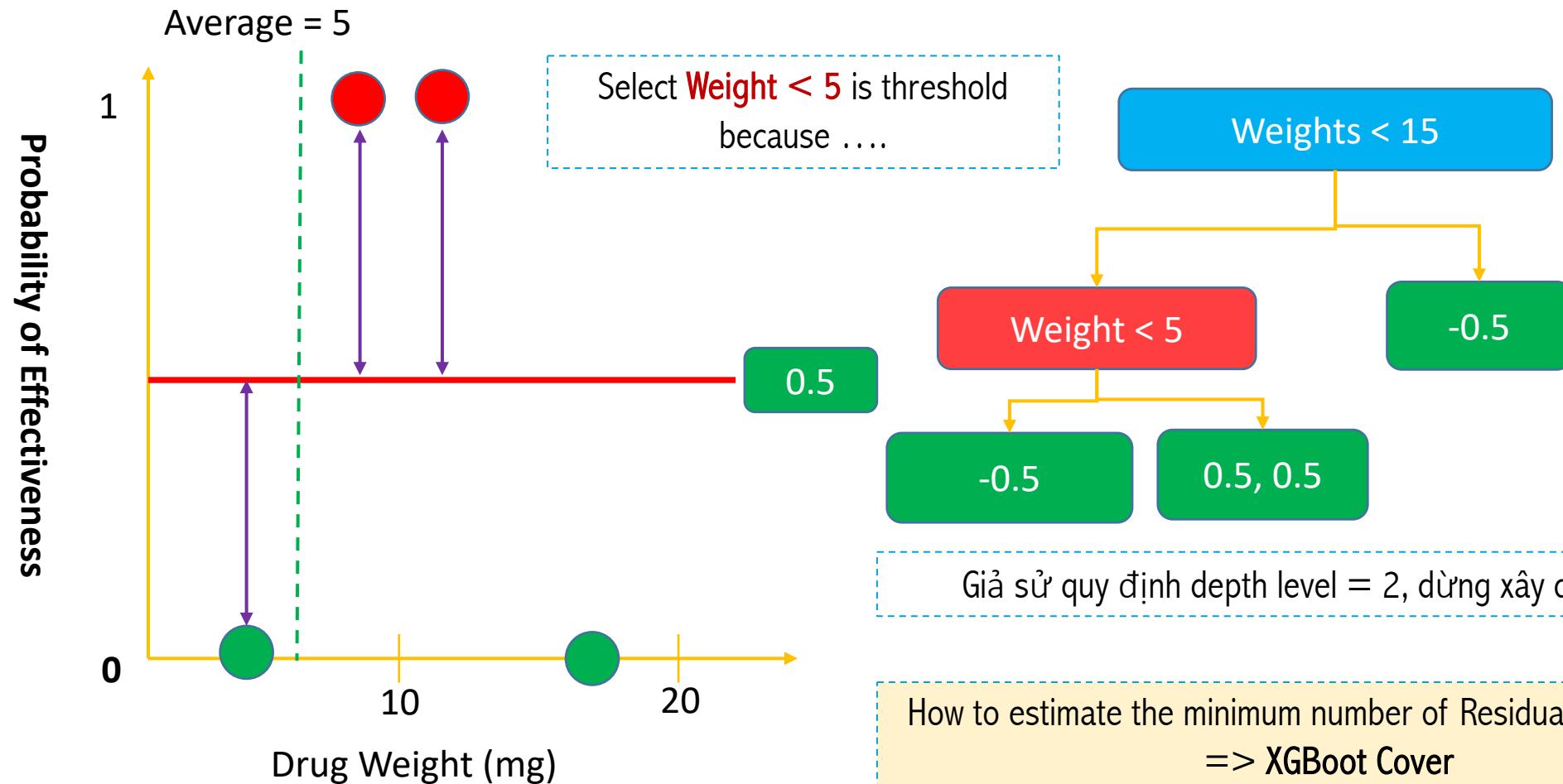
XGBoost for Classification



$$\text{Gain} = 1 + 2 - 0.33 = 2.66$$

Similarity Score = $\frac{(\sum \text{Residual}_i)^2}{\sum [\text{previous probability}_i \times (1 - \text{previous probability}_i)] + \lambda}$

XGBoost for Classification



What is a Cover

Similarity Score for Classification:

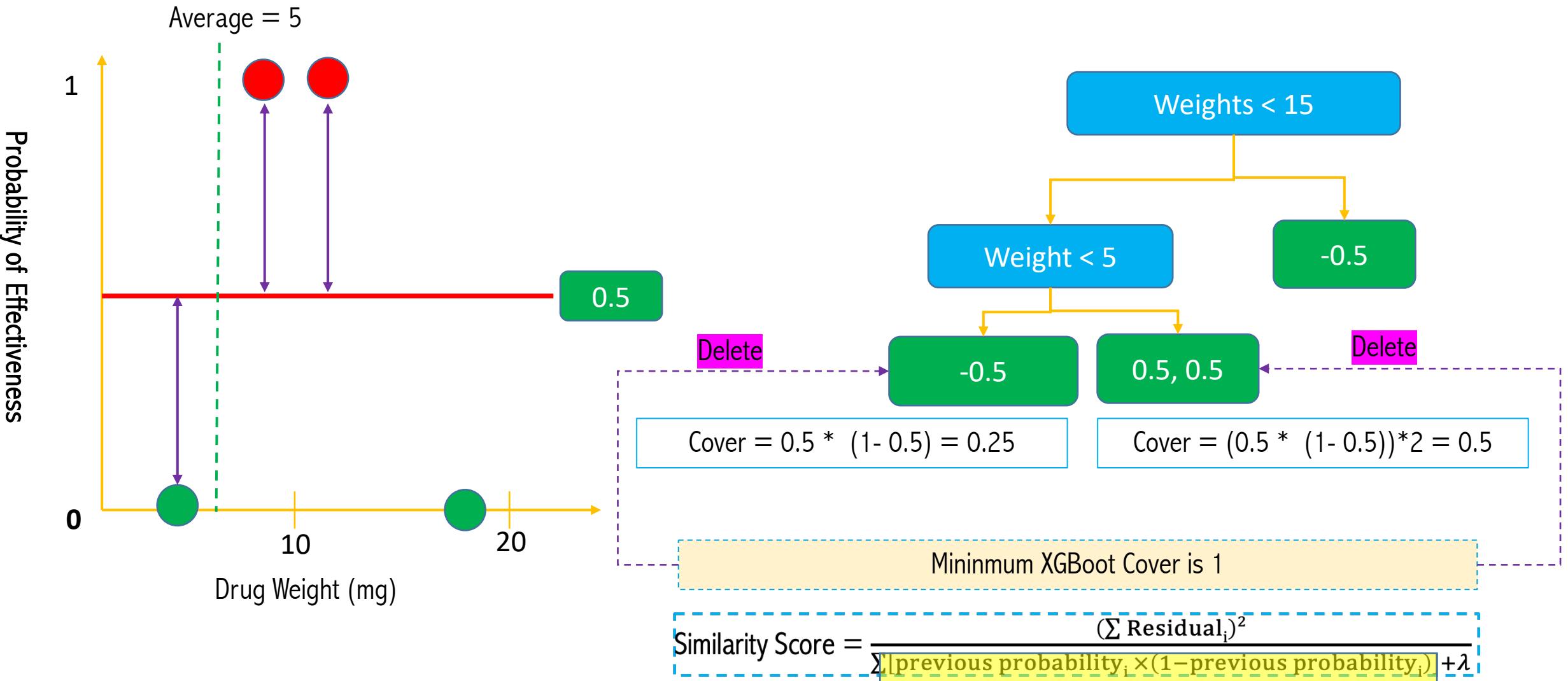
$$\text{Similarity Score} = \frac{(\sum \text{Residual}_i)^2}{\sum [\text{previous probability}_i \times (1 - \text{previous probability}_i)] + \lambda}$$

Cover

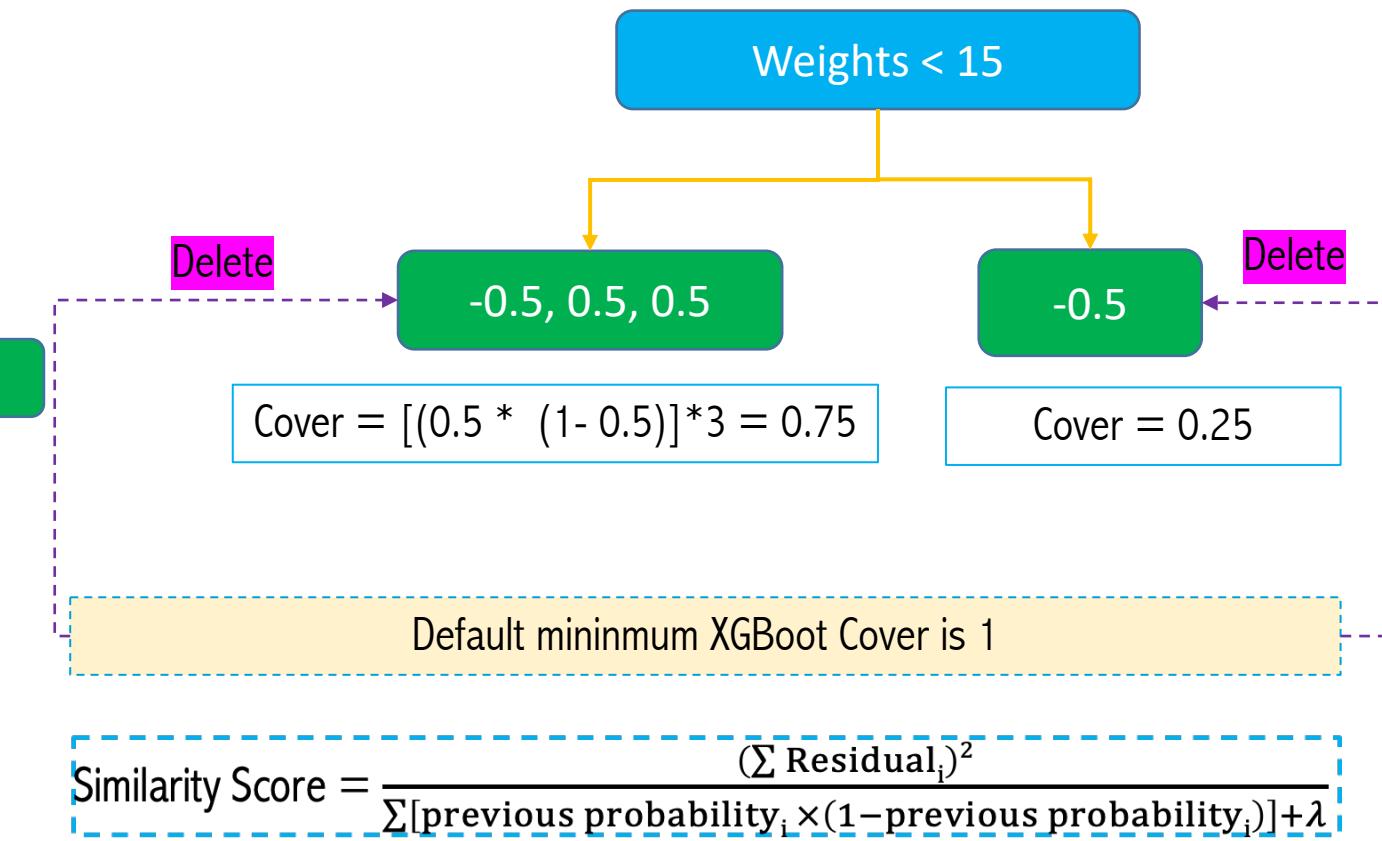
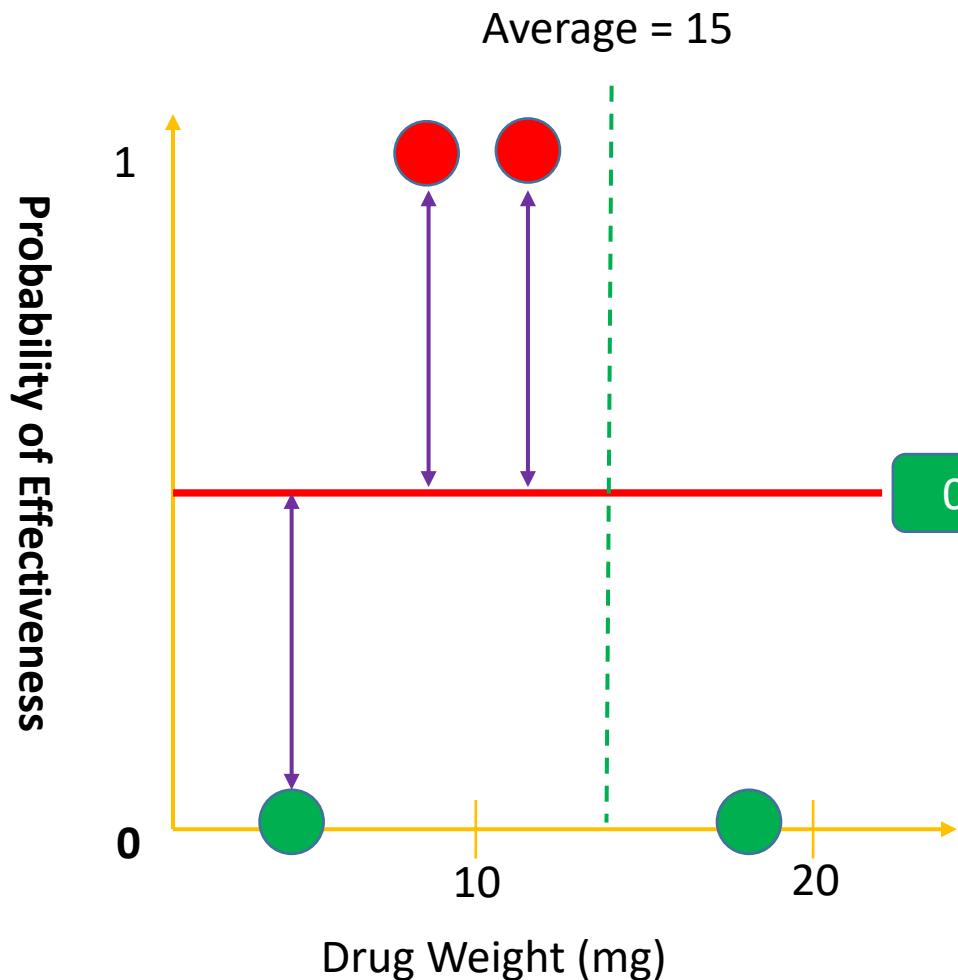
Similarity Score for Prediction:

$$\text{Similarity Score} = \frac{(\sum \text{Residual}_i)^2}{\text{number of residual} + \lambda}$$

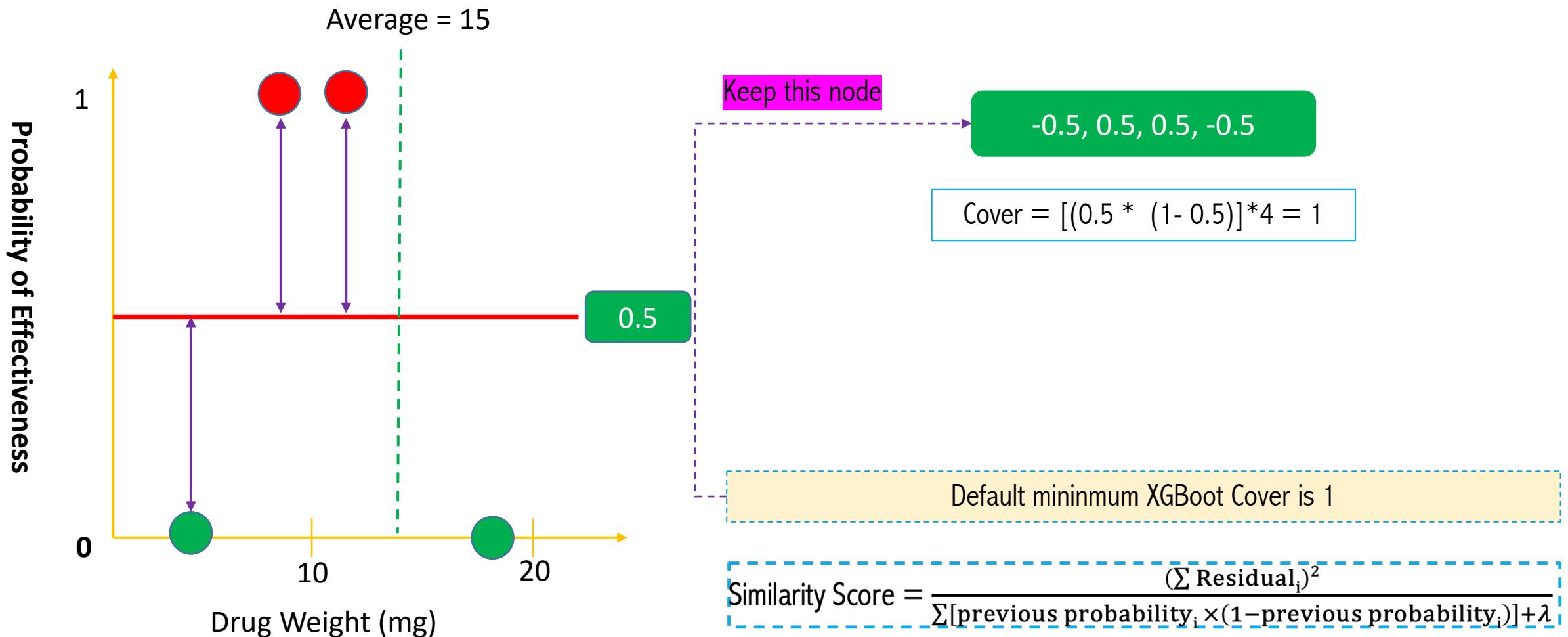
What is a Cover



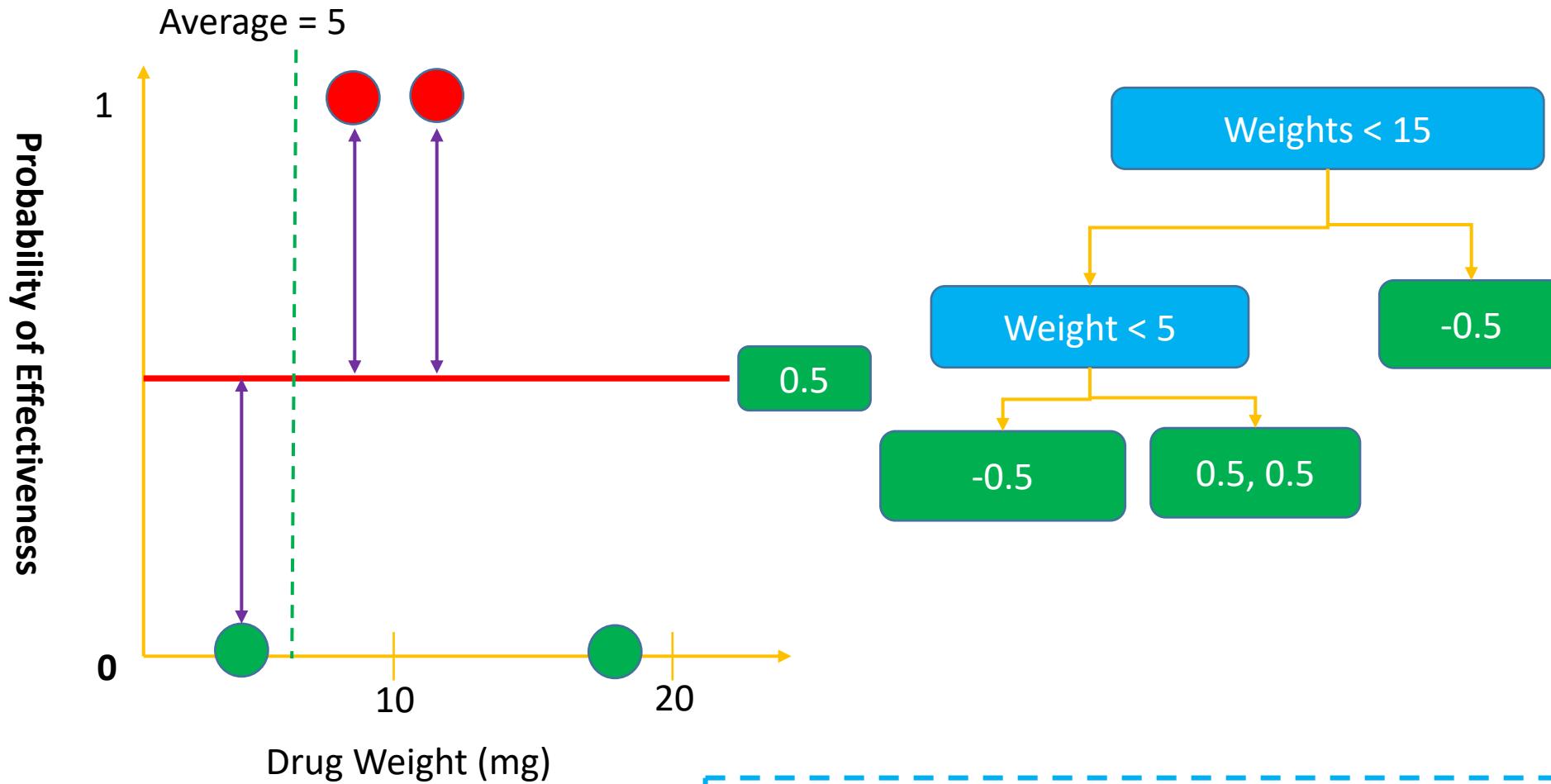
Xgboost for Classification



Xgboost for Classification



How to predict the value



$$\text{Output Value} = \frac{(\sum \text{Residual}_i)}{\sum [\text{previous probability}_i \times (1 - \text{previous probability}_i)] + \lambda}$$

How to predict the value

Drug Weight	Drug Effectiveness
Green	No
Red	Yes
Red	Yes
Green	No

Initial prediction is that the probability of drug effective is 50%

2 Yes and 2 No => Probability Yes = $2/4 = 1/2 = 0.5$

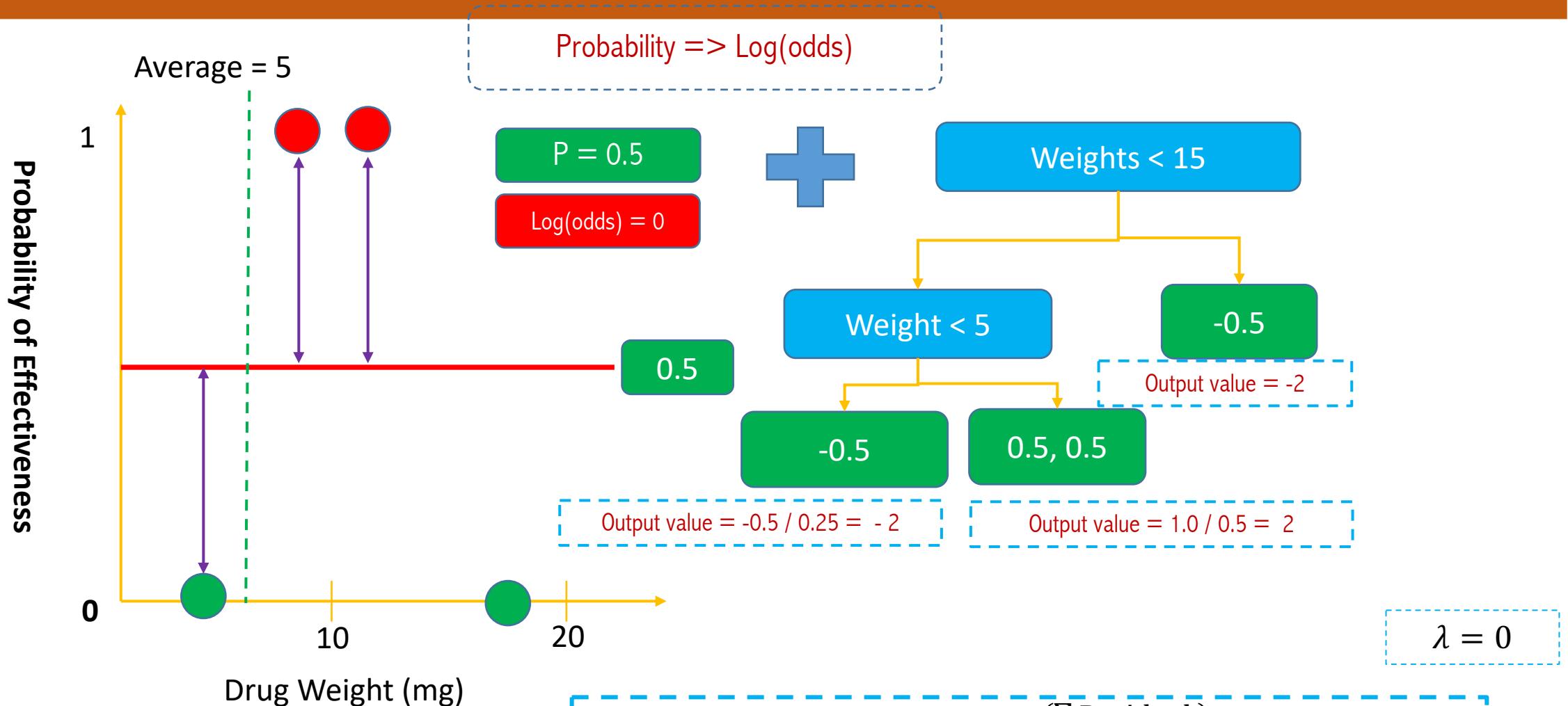
$$\text{Log(Odds)} = \log \left(\frac{\text{Probability Yes}}{\text{Probability No}} \right) = 0$$

In XGBoost (or Gradient Boost), the initial prediction is that the log(Odds)

$$\text{Probability of Drug Effectiveness} = \frac{e^{\text{log(Odds)}}}{1 + e^{\text{log(Odds)}}}$$

$$\text{Probability of Drug Effectiveness} = \frac{e^0}{1 + e^0} = 0.5$$

How to predict the value



$$\text{Output Value} = \frac{(\sum \text{Residual}_i)}{\sum [\text{previous probability}_i \times (1 - \text{previous probability}_i)] + \lambda}$$

Transformation formula for getting value at a leaf.

How to predict the value

Probability => Log(odds)

$$P = 0.5 \quad +$$

$$\text{Log}(odds) = 0$$

α

Weights < 15

Weight < 5

-0.5

-0.5

0.5, 0.5

Output value = -2

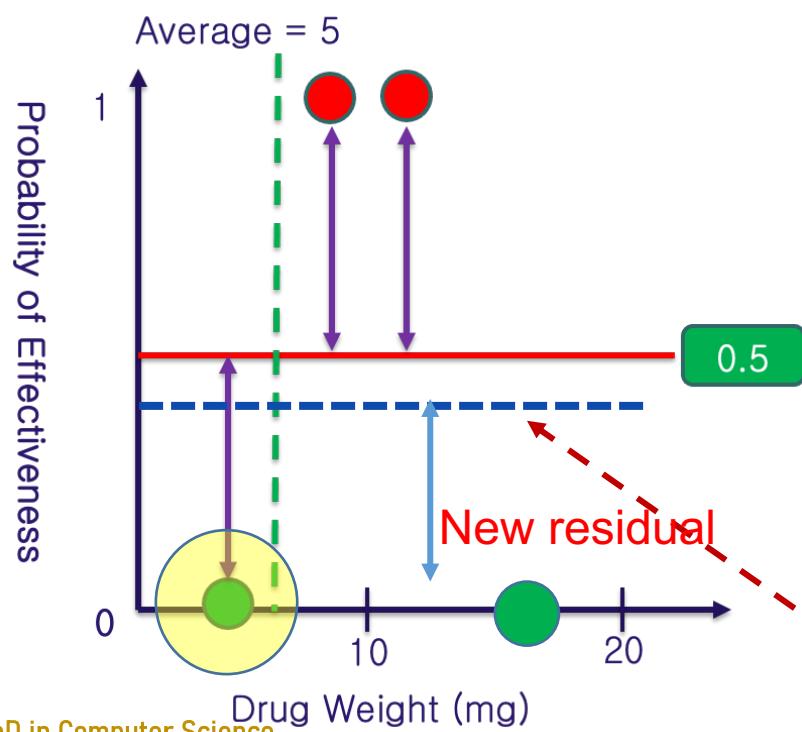
Output value = -2

Output value = 2

$$\frac{p}{1-p} = \text{odds}$$

$$\text{Log}\left(\frac{p}{1-p}\right) = \text{log(odds)}$$

$$\alpha = 0.3$$



$$\text{Prediction} = 0 + 0.3 * (-2) = -0.6$$

$$\text{Probability} = \frac{e^{\text{log(odds)}}}{1 + e^{\text{log(odds)}}}$$

$$\text{Probability} = \frac{e^{-0.6}}{1+e^{-0.6}} = 0.35$$

How to predict the value

Probability => Log(odds)

Can we
change P?

P = 0.5

Log(odds) = 0



α *

Weights < 15

Weight < 5

-0.5

Output value = -2

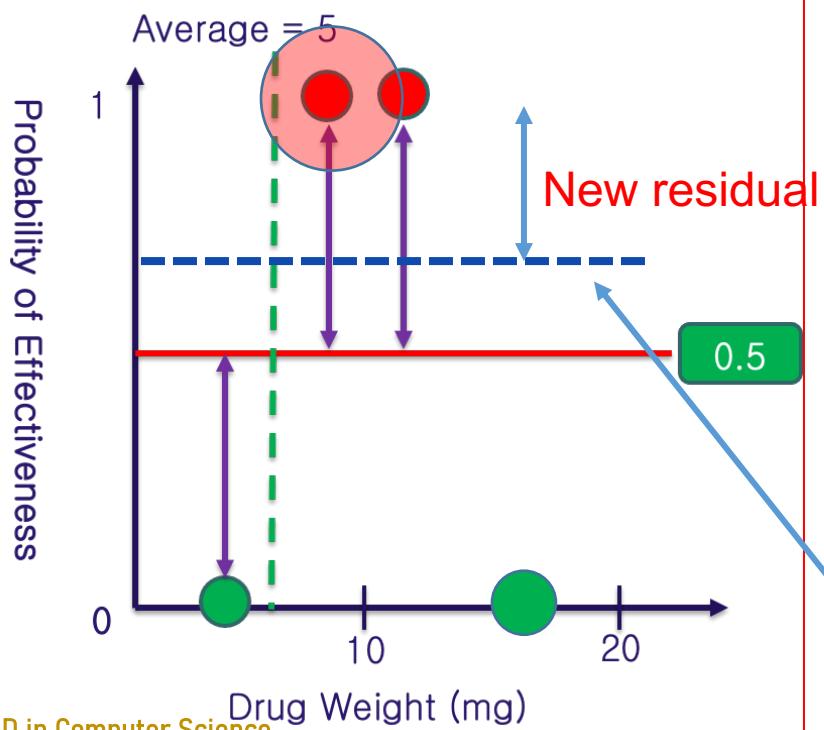
0.5, 0.5

Output value = 2

$$\frac{p}{1-p} = \text{odds}$$

$$\log\left(\frac{p}{1-p}\right) = \log(\text{odds})$$

$$\alpha = 0.3$$



$$\text{Log(odds)} = \text{Prediction} = 0 + 0.3 * (2) = 0.6$$

$$\text{Probability} = \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}}$$

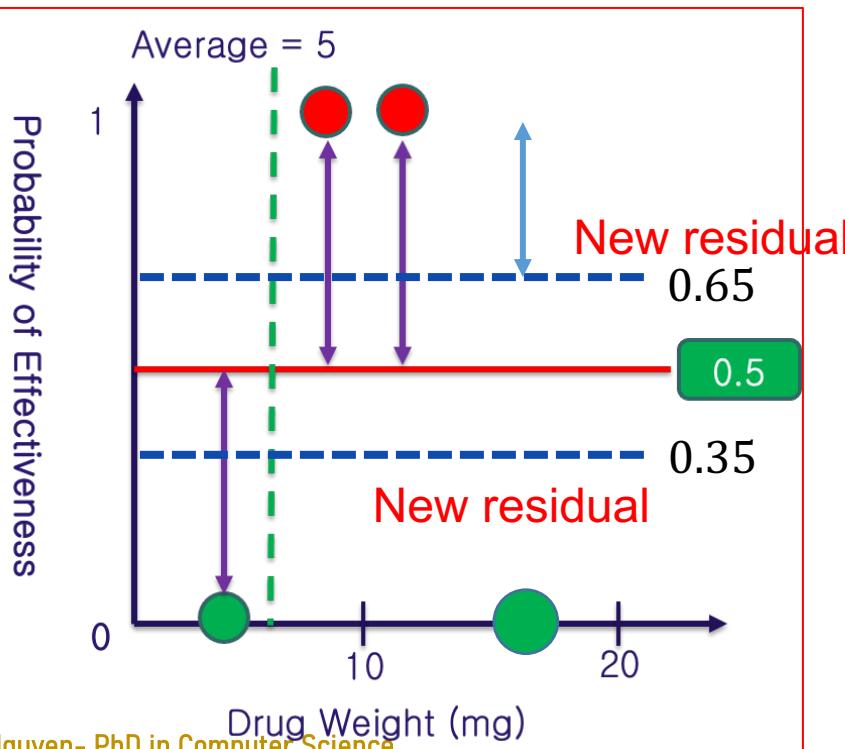
$$\text{Probability} = \frac{e^{0.6}}{1 + e^{0.6}} = 0.65$$

Build 2nd Tree

Probability => Log(odds)

P = 0.5

Log(odds) = 0

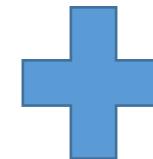


Weights < 15

Weight < 5

-0.5
0.5, 0.5

-0.5
Output value = -2
Output value = 2



* *

-0.35, 0.35, 0.35, -0.35

$$\text{Similarity Score} = \frac{(-0.35 + 0.35 + 0.35 - 0.35)^2}{0.35 \times (1-0.35) + 0.65 \times (1-0.65) + 0.65 \times (1-0.65) + 0.35 \times (1-0.35)}$$

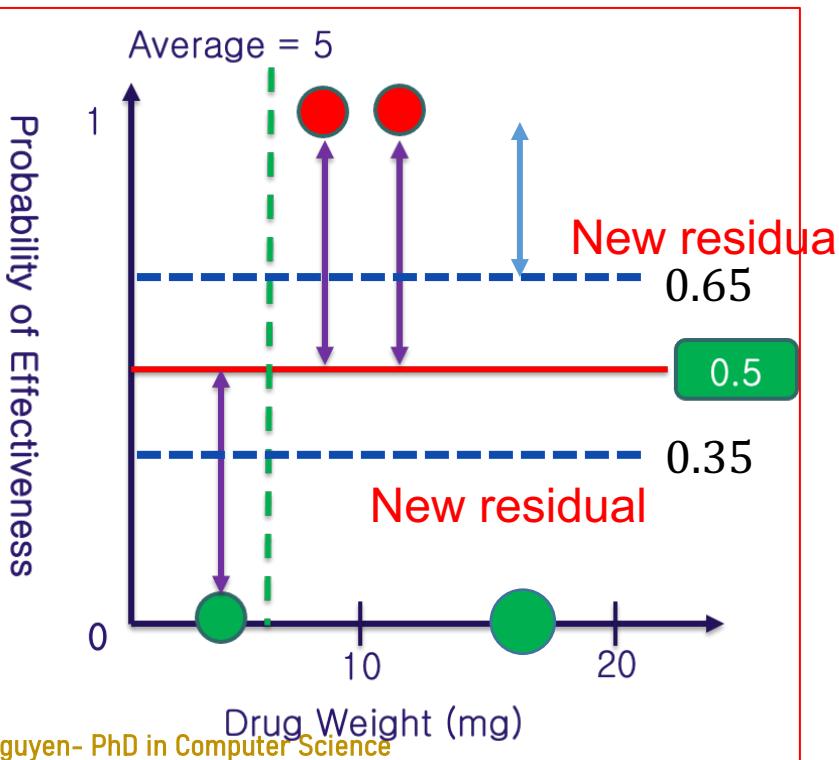
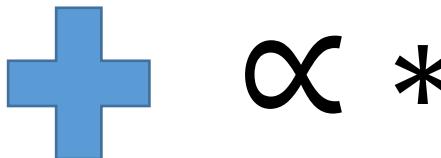
$$\text{Similarity Score} = \frac{(\sum \text{Residual}_i)^2}{\sum [\text{previous probability}_i \times (1-\text{previous probability}_i)] + \lambda}$$

Build 2nd Tree

Probability => Log(odds)

P = 0.5

Log(odds) = 0



Weights < 15

Weight < 5

-0.5
0.5, 0.5

Output value = -2

Output value = 2



* *

-0.35, 0.35, 0.35, -0.35

$$\text{Output Score} = \frac{(-0.35 + 0.35 + 0.35 - 0.35)}{0.35 \times (1-0.35) + 0.65 \times (1-0.65) + 0.65 \times (1-0.65) + 0.35 \times (1-0.35) + \lambda}$$

$$\text{Output Score} = \frac{(\sum \text{Residual}_i)}{\sum [\text{previous probability}_i \times (1-\text{previous probability}_i)] + \lambda}$$

Build 2nd Tree

Probability => Log(odds)

$P = 0.5$

$\text{Log(odds)} = 0$

$$+ \alpha * *$$

Weights < 15

Weight < 5

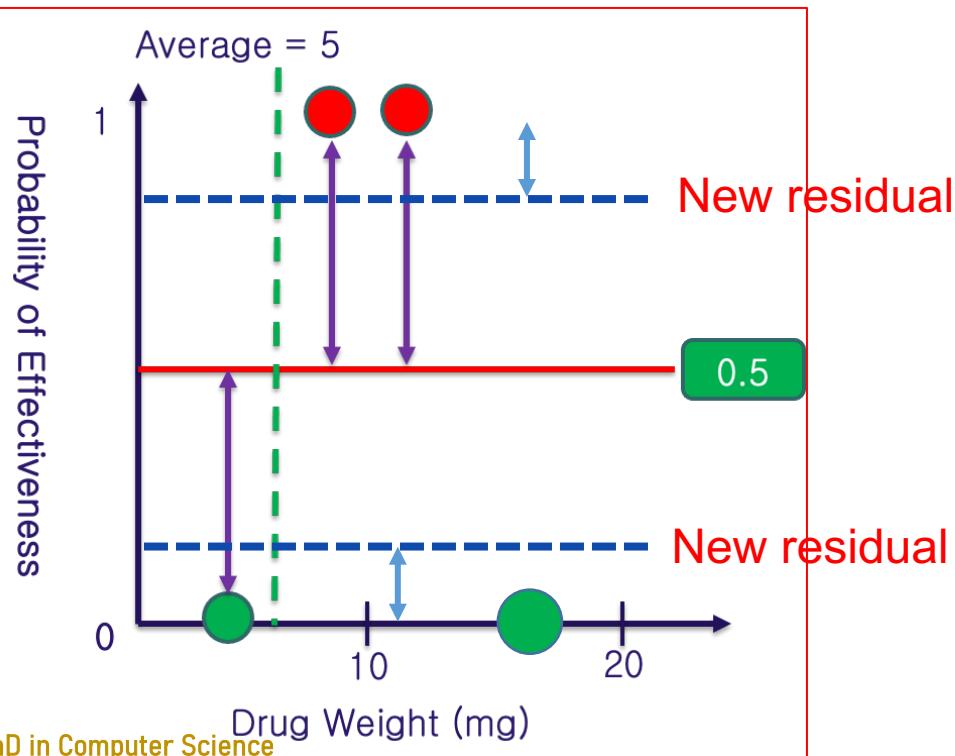
-0.5

-0.5

0.5, 0.5

$$+$$

Output value = -2 Output value = 2



$$\alpha *$$

Weights < 5

-0.35

Weight < 15

0.35, 0.35

-0.35

Review Questions

1. When do you stop to build the Tree

2. What's happen when $\lambda > 0$



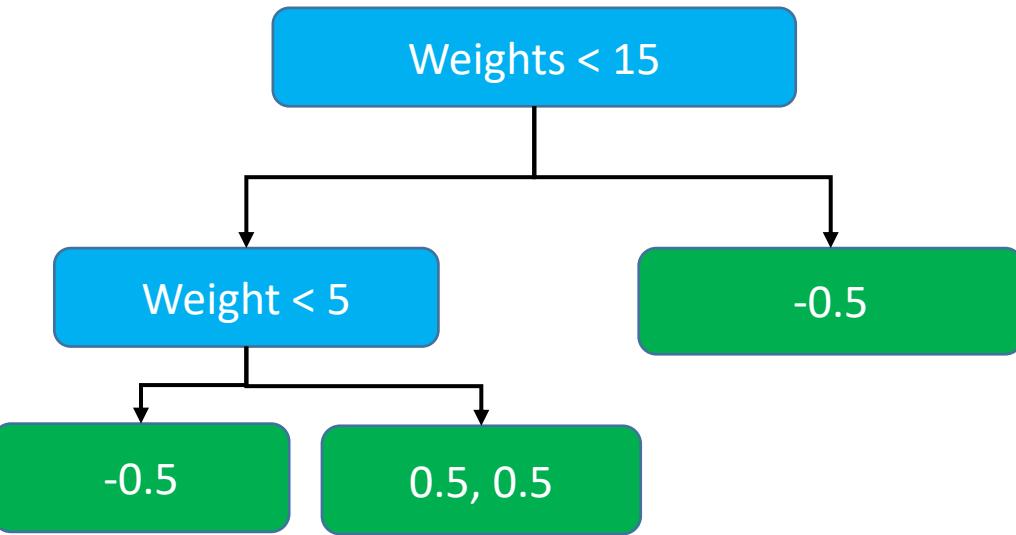
$$\text{Similarity Score} = \frac{(\sum \text{Residual}_i)^2}{\sum [\text{previous probability}_i \times (1 - \text{previous probability}_i)] + \lambda}$$

Outline

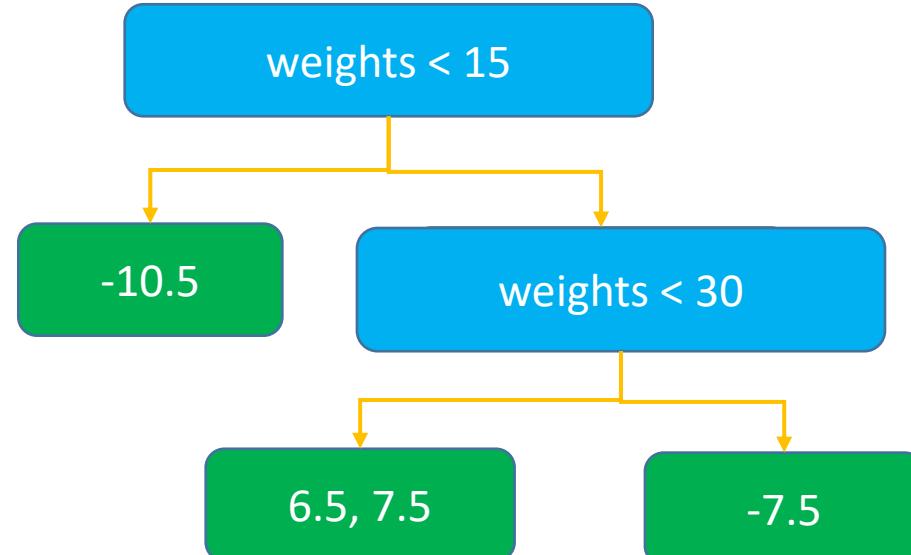
- **Regularization**
- **XGBoost For Regression**
- **XGBoost For Classification**
- **XGBoost: Mathematical Explanation**
- **How To Fill Missing Values**
- **Example**

XGBoost: Behind The Scenes

Classification



Regression



$$\text{Similarity Score} = \frac{(\sum \text{Residual})^2}{\sum \bar{y}_i \times (1 - \bar{y}_i) + \lambda}$$

$$\text{Output value} = \frac{(\sum \text{Residual})}{\sum \bar{y}_i \times (1 - \bar{y}_i) + \lambda}$$

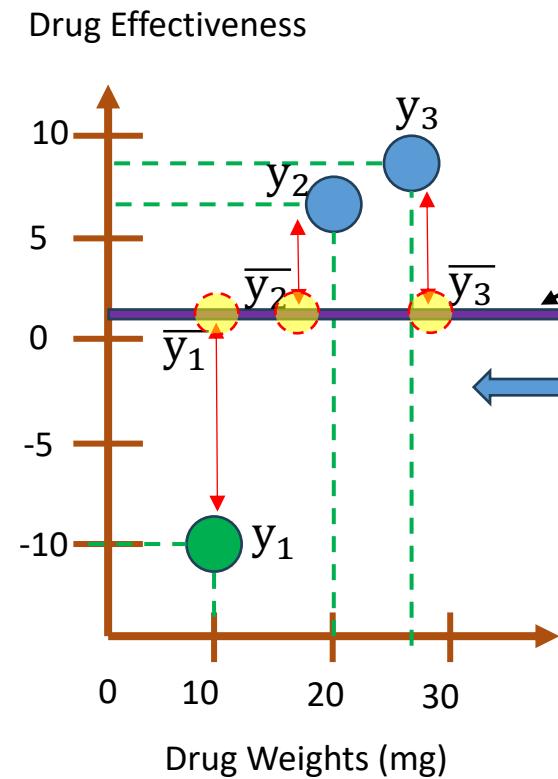


$$\text{Similarity Score} = \frac{(\sum \text{Residual})^2}{\text{Number of Residual} + \lambda}$$

$$\text{Output Value} = \frac{(\sum \text{Residual})}{\text{Number of Residual} + \lambda}$$

XGBoost: Behind The Scenes

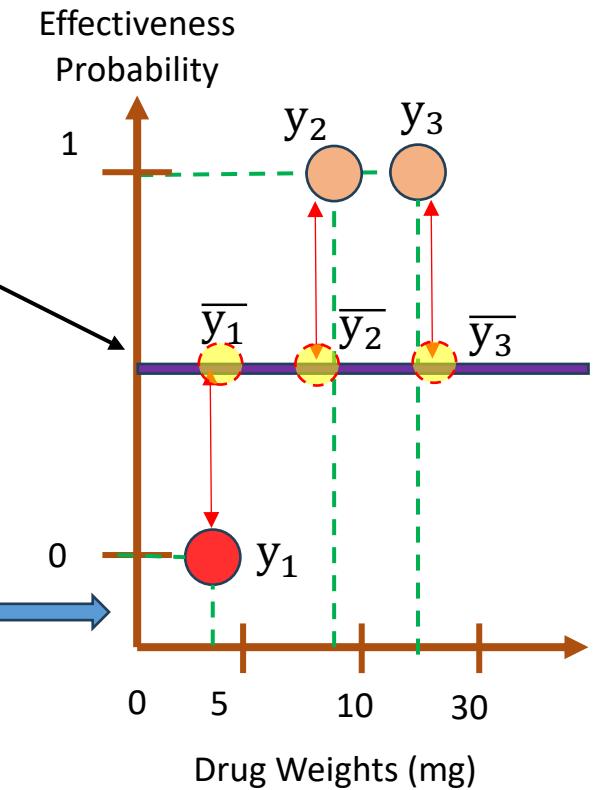
Regression



Dự đoán ban đầu hiệu quả thuốc

0.5

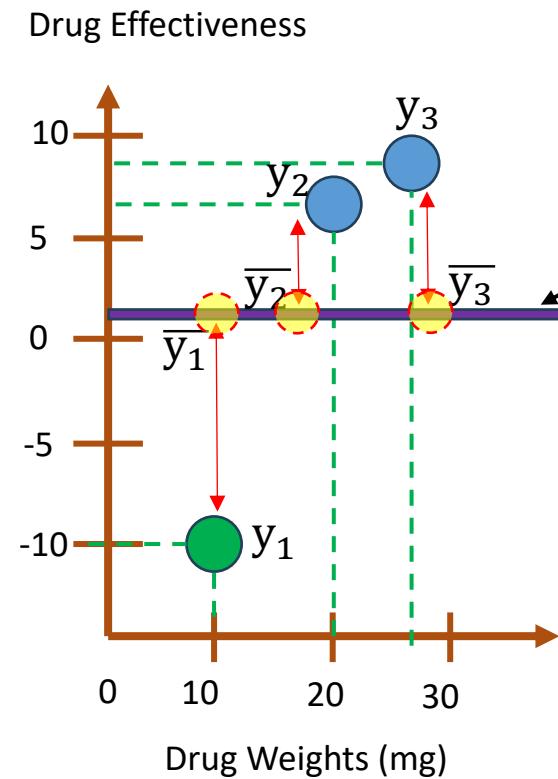
Classification



Sử dụng loss functions xây dựng cây

XGBoost: Behind The Scenes

Regression



Dự đoán ban đầu hiệu quả thuốc

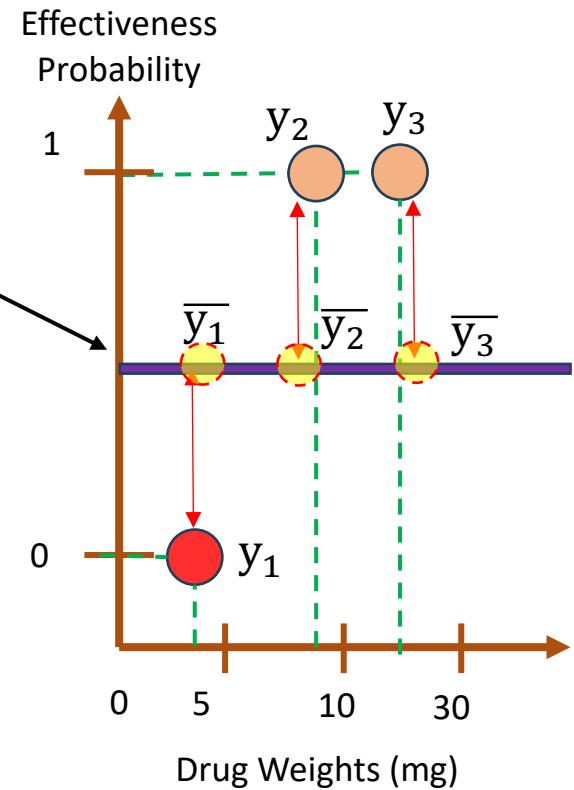
0.5

$$\sum_{i=1}^n \mathcal{L}(y_i, \bar{y}_i) + \gamma T + \lambda P^2$$

γ is a user definable penalty to encourage pruning

XGBoost can prune even when $\gamma = 0$

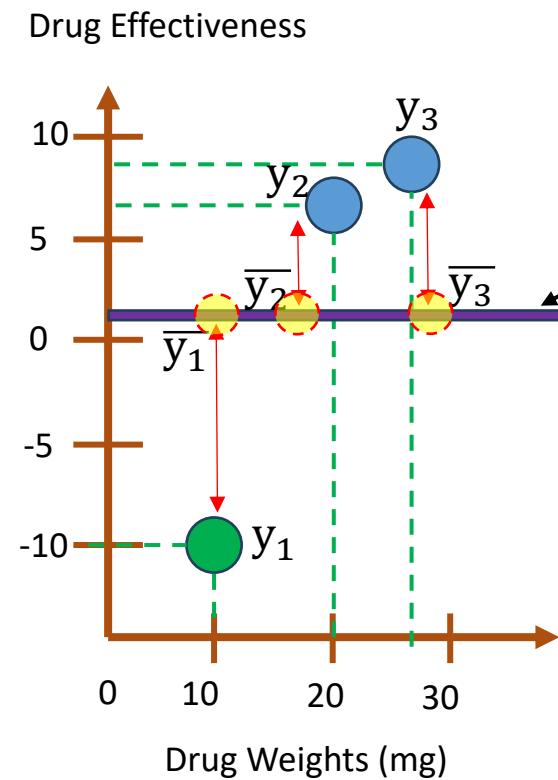
Classification



Pruning is executed after the full tree built => It plays no role in deriving the Optimal Output Values

XGBoost: Behind The Scenes

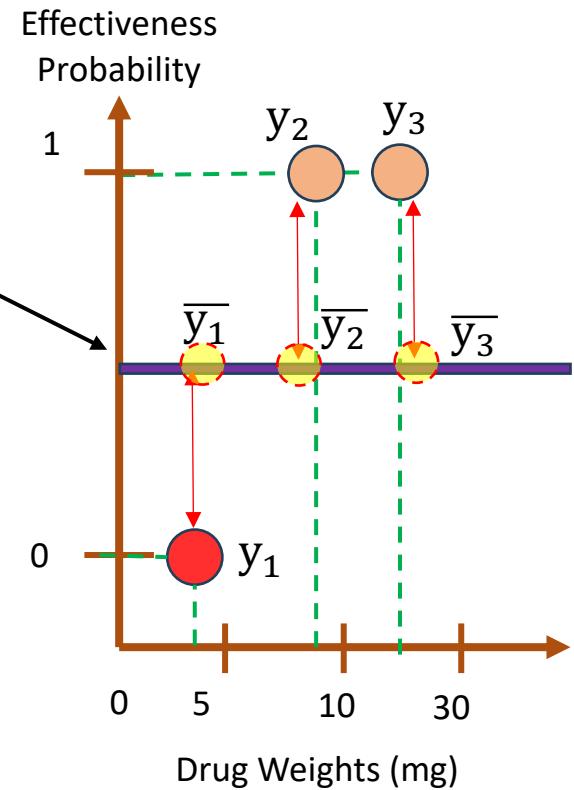
Regression



Dự đoán ban đầu hiệu quả thuốc

0.5

Classification



Bỏ qua γ

$$\sum_{i=1}^n \mathcal{L}(y_i, \bar{y}_i) + \cancel{\gamma T} + \lambda P^2$$

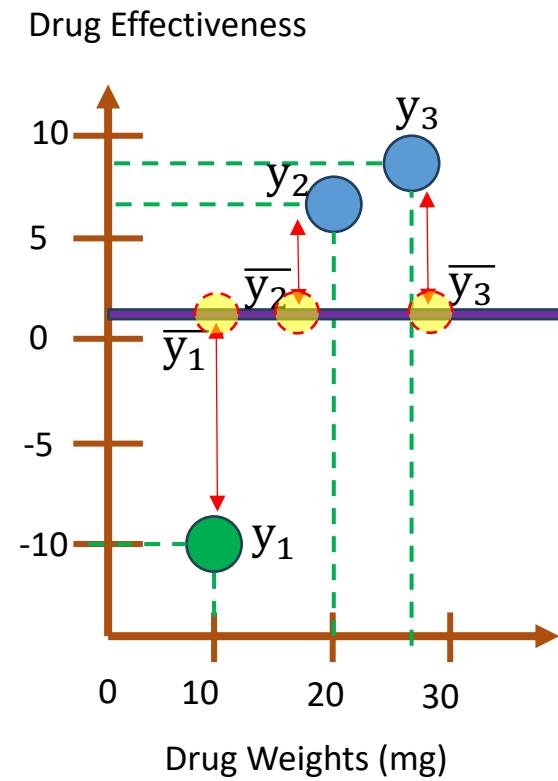
γ is a user definable penalty to encourage pruning

XGBoost can prune even when $\gamma = 0$

Pruning is executed after the full tree built => It plays no role in deriving the Optimal Output Values

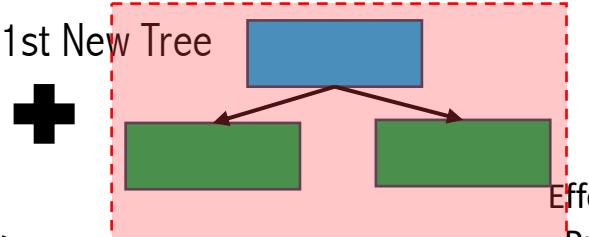
XGBoost: Behind The Scenes

Regression

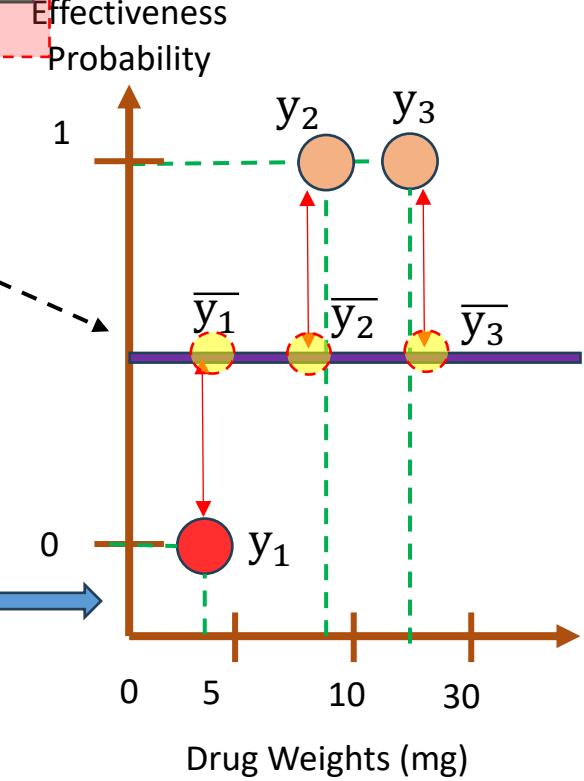


Dự đoán ban đầu
hiệu quả thuốc

1st New Tree



Classification



XGBoost builds the new tree based on the loss function:

$$\sum_{i=1}^n \mathcal{L}(y_i, \hat{y}_i^0 + P) + \frac{1}{2} \lambda P^2$$

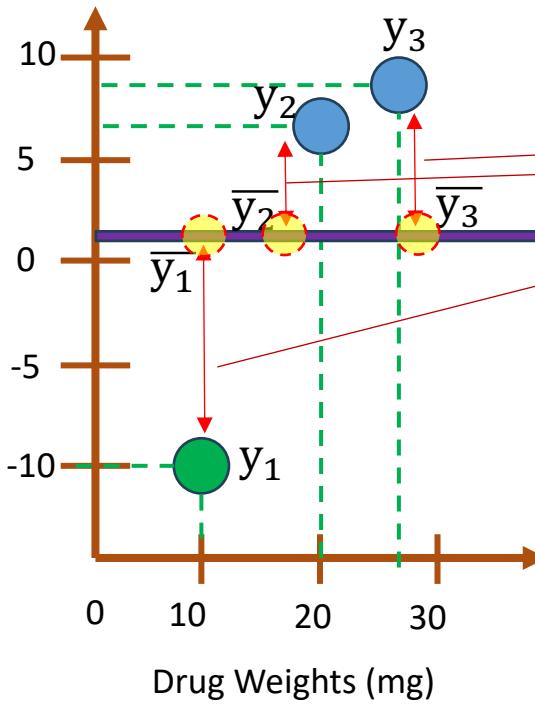
Ridge Regression Regularization term

Mục tiêu: tìm giá trị dự đoán cho mỗi leaf (P) của cây mới nhằm minimize hàm loss.

XGBoost Regression: Behind The Scenes

Regression

Drug Effectiveness



Dự đoán ban đầu
hiệu quả thuốc

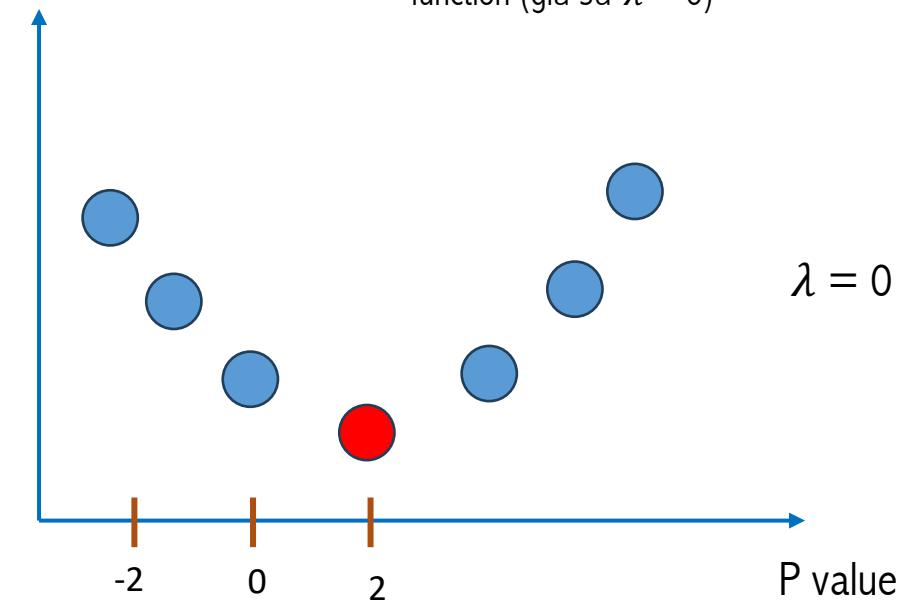
0.5

-10.5, 6.5, 7.5

Loss function

Giá trị P cần tìm là giá trị ứng với đạo hàm
của loss theo P bằng 0

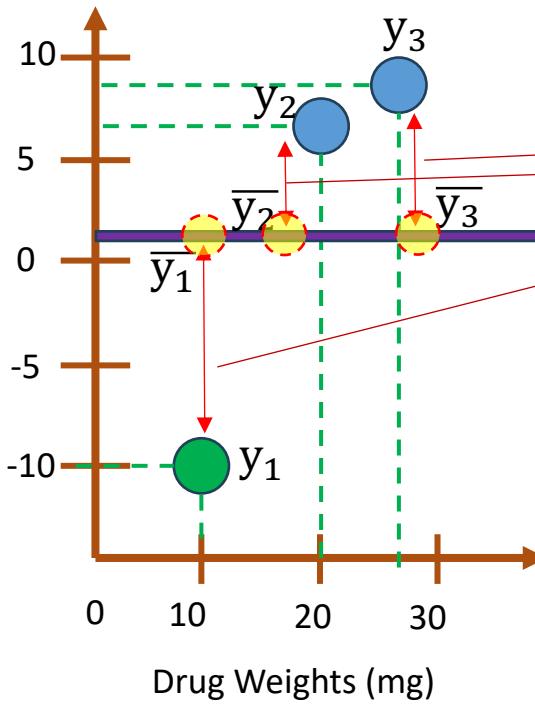
Chúng cần tìm giá trị đầu ra của nút lá này (giá trị P) bằng cách minimize loss
function (giả sử $\lambda = 0$)



XGBoost Regression: Behind The Scenes

Regression

Drug Effectiveness

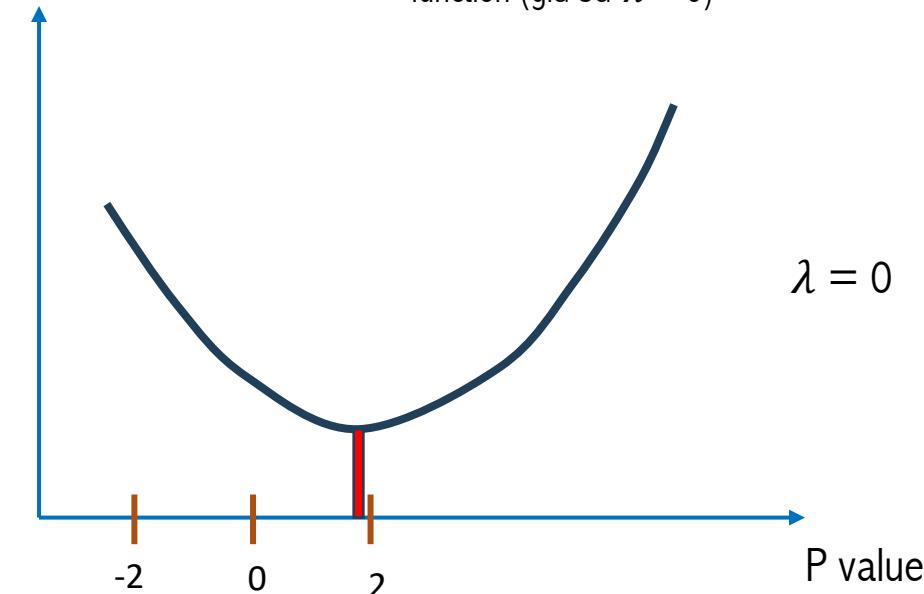


Giá trị P cần tìm là giá trị ứng với đạo hàm của loss theo P bằng 0

Loss function

$\sum_{i=1}^n \mathcal{L}(y_i, \bar{y}_i + P) + \frac{1}{2} \lambda P^2$

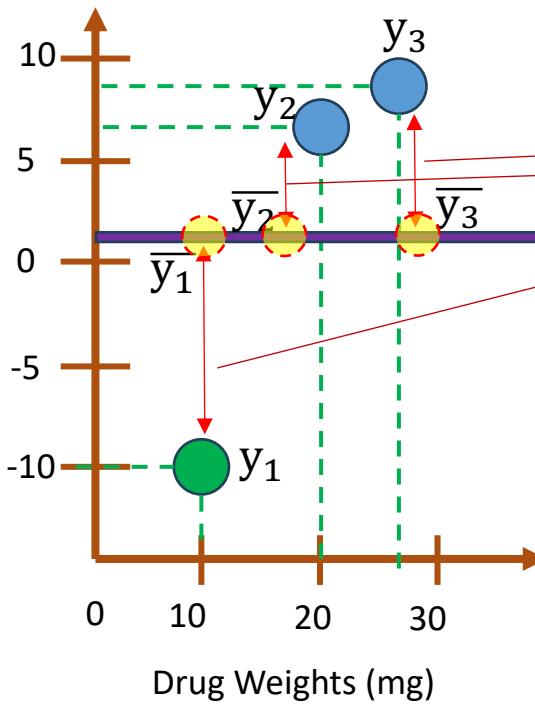
Chúng cần tìm giá trị đầu ra của nút lá này (giá trị P) bằng cách minimize loss function (giả sử $\lambda = 0$)



XGBoost: Behind The Scenes

Regression

Drug Effectiveness



Dự đoán ban đầu
hiệu quả thuốc

0.5

Residual
-10.5, 6.5, 7.5

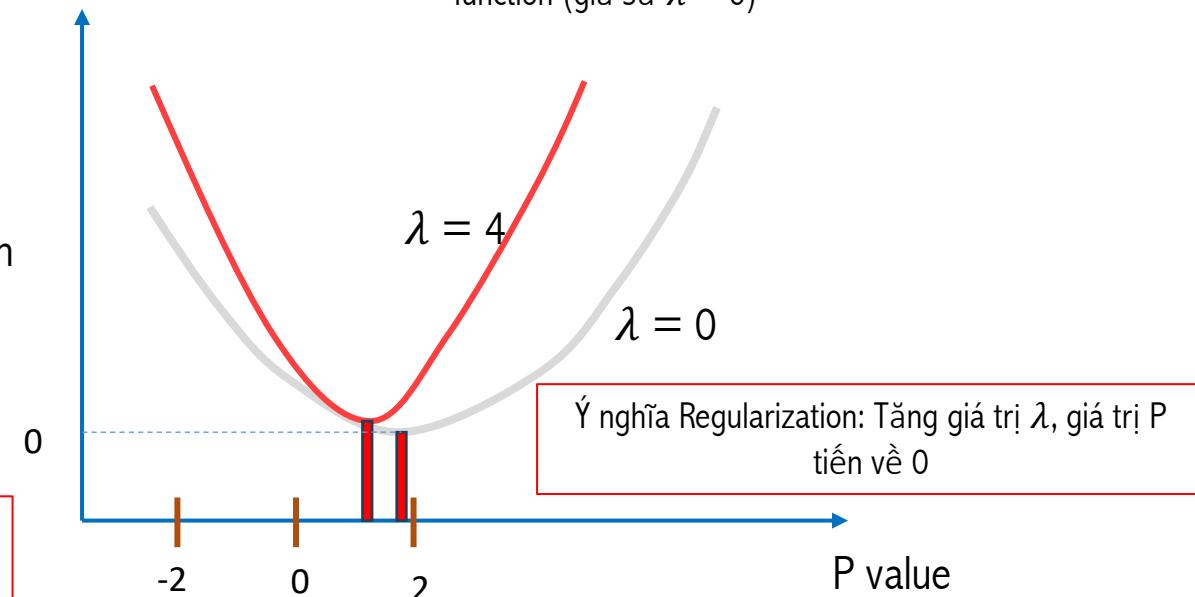
Loss function

Giá trị P cần tìm là giá trị ứng với đạo hàm
của loss theo P bằng 0

What's happen if λ is
very large?

$$\sum_{i=1}^n \mathcal{L}(y_i, \bar{y}_i + P) + \frac{1}{2} \lambda P^2$$

Chúng cần tìm giá trị đầu ra của nút lá này (giá trị P) bằng cách minimize loss function (giả sử $\lambda = 0$)



Ý nghĩa Regularization: Tăng giá trị λ , giá trị P
tiến về 0

XGBoost Regression: Behind The Scenes

$$\sum_{i=1}^n \mathcal{L}(y_i, \bar{y}_i^0 + P) + \frac{1}{2} \lambda P^2$$

Rất khó để tìm optimization, nên cũng ta sẽ sắp sỉ hàm loss bằng Second Order Tayler Appriximation

$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + \left[\frac{d}{d\bar{y}_i} \mathcal{L}(y_i, \bar{y}_i) \right] P + \frac{1}{2} \left[\frac{d^2}{d\bar{y}_i^2} \mathcal{L}(y_i, \bar{y}_i) \right] P^2$$

$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + g_i P + \frac{1}{2} h_i P^2$$

g (gradient) presents the first derivative of the loss function
h (hessian) presents the second derivative of the loss function

$$\mathcal{L}(y_1, \bar{y}_1^0) + g_1 P + \frac{1}{2} h_1 P^2 + \mathcal{L}(y_2, \bar{y}_2^0) + g_2 P + \frac{1}{2} h_2 P^2 + \dots + \mathcal{L}(y_n, \bar{y}_n^0) + g_n P + \frac{1}{2} h_n P^2 + \frac{1}{2} \lambda P^2$$

Tìm giá trị P cần tìm sao cho đạo
hàm của loss function theo P bằng 0

$$\frac{d}{dP} \left[(g_1 + g_2 + \dots + g_n)P + \frac{1}{2} (h_1 + h_2 + \dots + h_n + \lambda)P^2 \right] = 0$$

XGBoost Regression: Output Value

1

$$\frac{d}{dP} \left[(g_1 + g_2 + \dots + g_n)P + \frac{1}{2}(h_1 + h + \dots + h_n + \lambda)P^2 \right] = 0$$

$$(g_1 + g_2 + \dots + g_n) + (h_1 + h + \dots + h_n\lambda)P = 0$$

2

$$g_i = \frac{d}{dy_i} \frac{1}{2} (y_i - \bar{y}_i)^2 = -(y_i - \bar{y}_i)$$

$$h_i = \frac{d^2}{dy_i^2} \frac{1}{2} (y_i - \bar{y}_i)^2 = 1$$

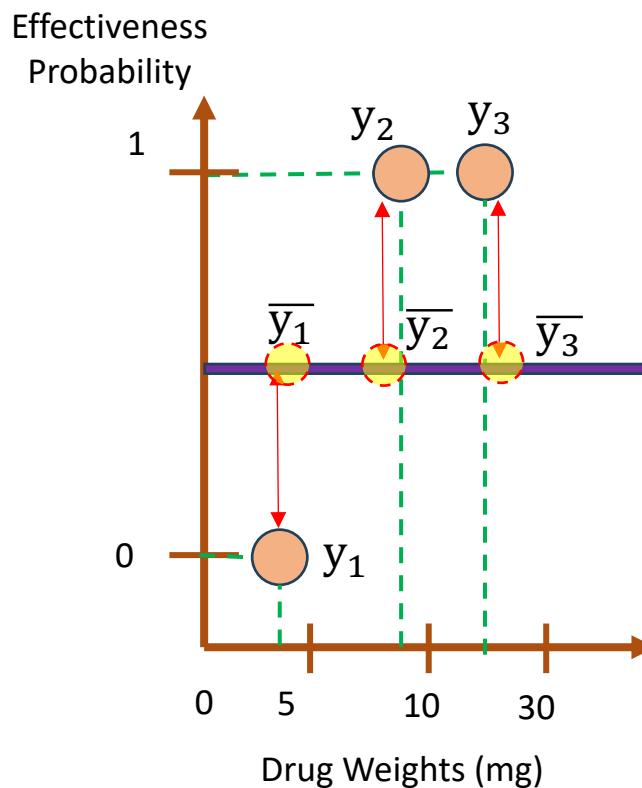
$$P = \frac{-(g_1 + g_2 + \dots + g_n)}{h_1 + h_2 + \dots + h_n + \lambda} = \frac{-(-(y_1 - \bar{y}_1) + -(y_2 - \bar{y}_2) + \dots + -(y_n - \bar{y}_n))}{1+1+\dots+1+\lambda} = \frac{\text{sum of residual}}{\text{num of sum residual} + \lambda}$$



Output value of the leaf
(or terminal node)

XGBoost For Classification: Output Value

Classification



$$\mathcal{L}(y_i, \bar{y}_i) = -[y_i \log(\bar{y}_i) + (1 - y_i) \log(1 - \bar{y}_i)]$$

Convert probability to Log(odds)

$$\mathcal{L}(y_i, \log(\text{odds})) = -y_i \log(\text{odds}) + \log(1 + e^{\log(\text{odds})})$$

$$g_i = \frac{d}{d \log(\text{odds})} \mathcal{L}(y_i, \log(\text{odds})) = -y_i + \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}} = -(y - \bar{y}_i)$$

$$h_i = \frac{d^2}{d \log(\text{odds})^2} \mathcal{L}(y_i, \log(\text{odds})) = \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}} \times \frac{1}{1 + e^{\log(\text{odds})}} = \bar{y}_i \times (1 - \bar{y}_i)$$

$$P = \frac{-(g_1 + g_2 + \dots + g_n)}{h_1 + h_2 + \dots + h_n + \lambda} = \frac{\text{sum of residual}}{\bar{y}_1 \times (1 - \bar{y}_1) + \bar{y}_2 \times (1 - \bar{y}_2) + \dots + \bar{y}_n \times (1 - \bar{y}_n) + \lambda} = \frac{(\Sigma \text{Residual})}{\sum \bar{y}_i \times (1 - \bar{y}_i) + \lambda}$$



XGBoost: Similarity Score

1

$$\sum_{i=1}^n \mathcal{L}(y_i, \bar{y}_i + P) + \frac{1}{2} \lambda P^2$$

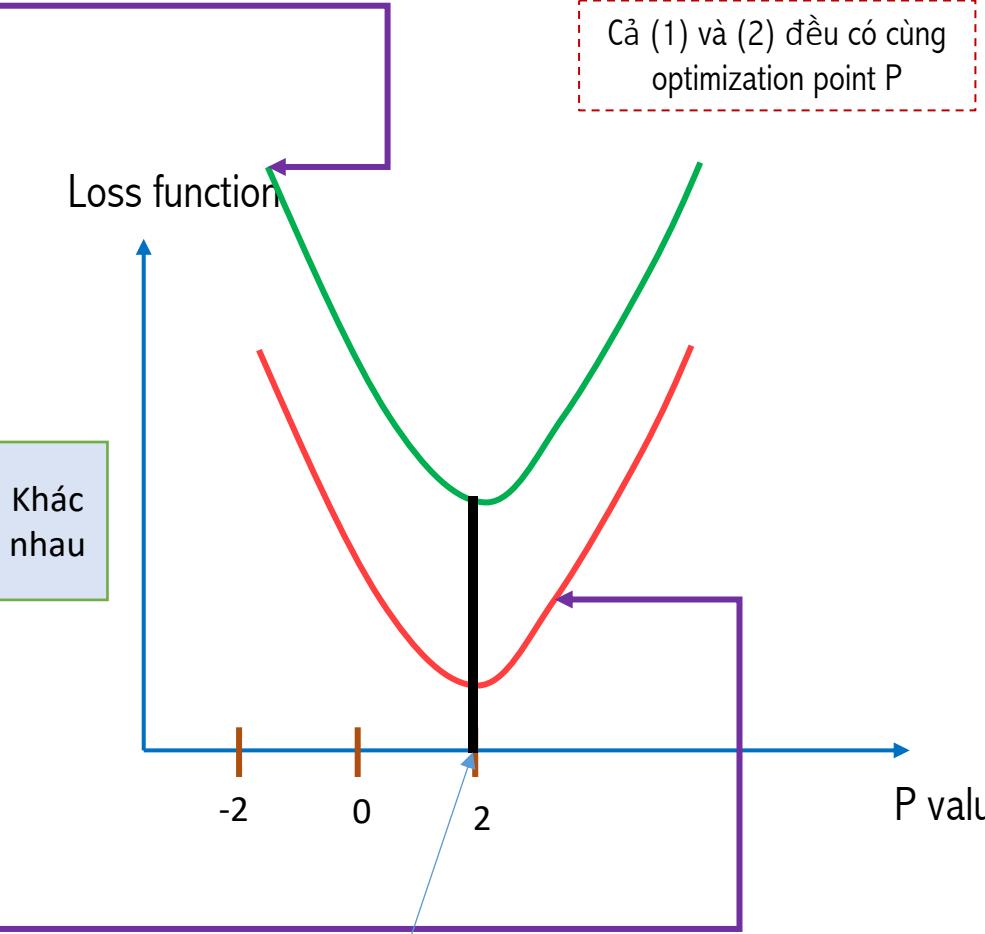
Rất khó để tìm optimization, nên cũng ta sẽ sắp sỉ hàm loss bằng Second Order Tayler Appriximation

$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + \left[\frac{d}{dy_i} \mathcal{L}(y_i, \bar{y}_i) \right] P + \frac{1}{2} \left[\frac{d^2}{dy_i^2} \mathcal{L}(y_i, \bar{y}_i) \right] P^2$$

$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + gP + \frac{1}{2} hP^2$$

2

$$(g_1 + g_2 + \dots + g_n)P + \frac{1}{2} (h_1 + h_2 + \dots + h_n + \lambda)P^2$$



$$P = \frac{-(g_1 + g_2 + \dots + g_n)}{h_1 + h_2 + \dots + h_n + \lambda}$$

XGBoost: Similarity Score

1

$$\sum_{i=1}^n \mathcal{L}(y_i, \bar{y}_i + P) + \frac{1}{2} \lambda P^2$$

Rất khó để tìm optimization, nên cũng ta sẽ sắp sỉ hàm loss bằng Second Order Tayler Approximation

$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + \left[\frac{d}{dy_i} \mathcal{L}(y_i, \bar{y}_i) \right] P + \frac{1}{2} \left[\frac{d^2}{dy_i^2} \mathcal{L}(y_i, \bar{y}_i) \right] P^2$$

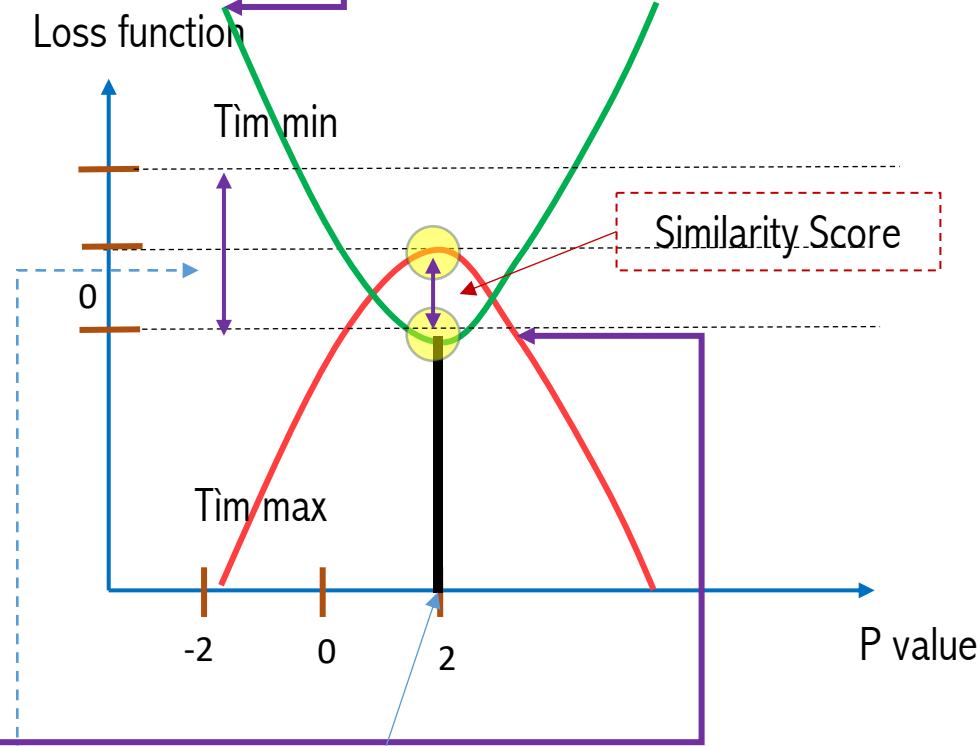
$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + gP + \frac{1}{2} hP^2$$

2

$$-1 \times (g_1 + g_2 + \dots + g_n)P + -1 \times \frac{1}{2} (h_1 + h_2 + \dots + h_n + \lambda)P^2$$

Implementation Similarity Score

Cả (1) và (2) đều có cùng optimization point P



$$P = \frac{-(g_1 + g_2 + \dots + g_n)}{h_1 + h_2 + \dots + h_n + \lambda}$$

XGBoost: Similarity Score

1

$$\sum_{i=1}^n \mathcal{L}(y_i, \bar{y}_i + P) + \frac{1}{2} \lambda P^2$$

Rất khó để tìm optimization, nên cũng ta sẽ sắp sỉ hàm loss bằng Second Order Tayler Approximation

$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + \left[\frac{d}{d\bar{y}_i} \mathcal{L}(y_i, \bar{y}_i) \right] P + \frac{1}{2} \left[\frac{d^2}{d\bar{y}_i^2} \mathcal{L}(y_i, \bar{y}_i) \right] P^2$$

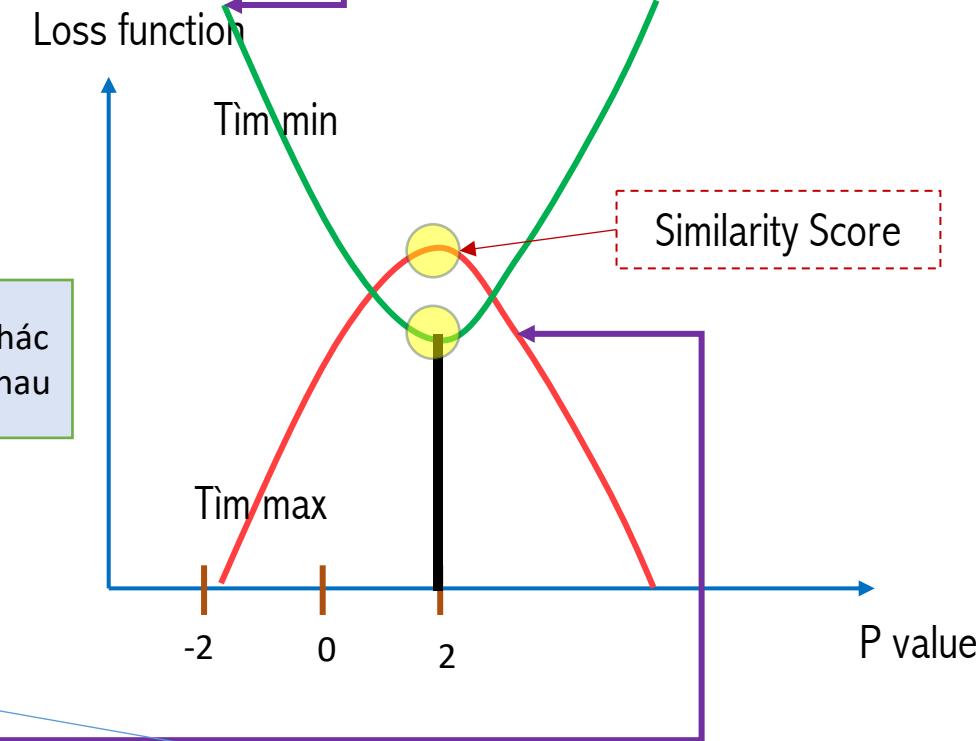
$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + gP + \frac{1}{2} hP^2$$

2

$$-1 \times (g_1 + g_2 + \dots + g_n)P - 1 \times \frac{1}{2} (h_1 + h_2 + \dots + h_n + \lambda)P^2$$

$$\text{Similarity Score} = \frac{1}{2} \frac{(g_1 + g_2 + \dots + g_n)^2}{(h_1 + h_2 + \dots + h_n + \lambda)}$$

Cả (1) và (2) đều có cùng optimization point P



$$P = \frac{-(g_1 + g_2 + \dots + g_n)}{h_1 + h_2 + \dots + h_n + \lambda}$$

XGBoost Regression: Similarity Score

1

$$\sum_{i=1}^n \mathcal{L}(y_i, \bar{y}_i + P) + \frac{1}{2} \lambda P^2$$

Rất khó để tìm optimization, nên cũng ta sẽ sắp sỉ hàm loss bằng Second Order Tayler Appriximation

$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + \left[\frac{d}{d\bar{y}_i} \mathcal{L}(y_i, \bar{y}_i) \right] P + \frac{1}{2} \left[\frac{d^2}{d\bar{y}_i^2} \mathcal{L}(y_i, \bar{y}_i) \right] P^2$$

$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + gP + \frac{1}{2} hP^2$$

2

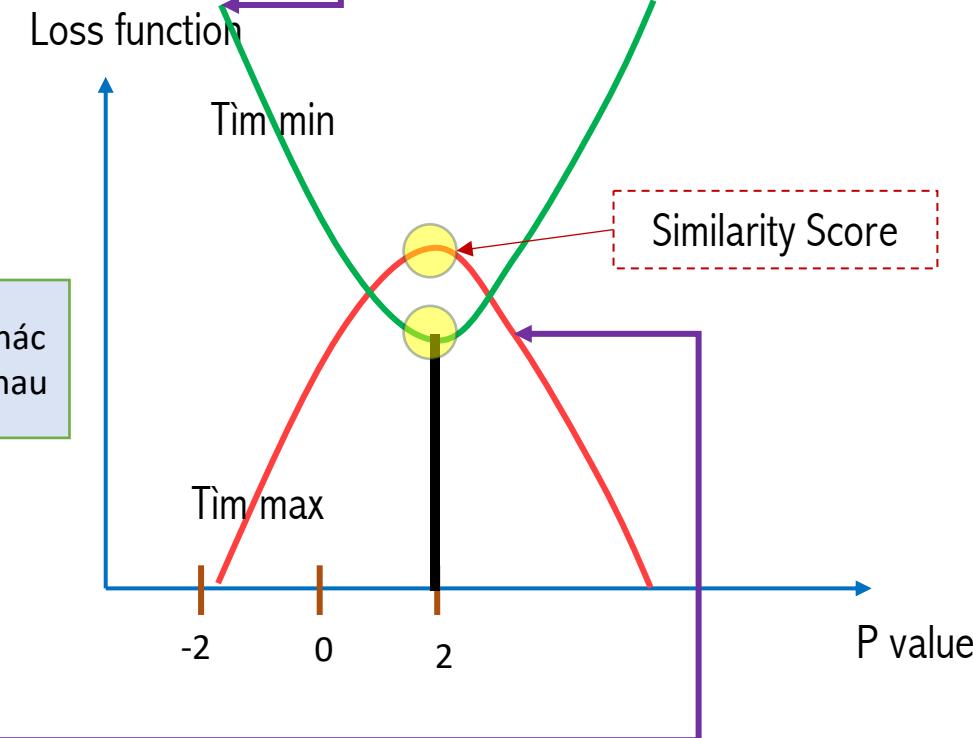
$$-1 \times (g_1 + g_2 + \dots + g_n)P + -1 \times \frac{1}{2} (h_1 + h_2 + \dots + h_n + \lambda)P^2$$

$$h_i = \frac{d}{d\bar{y}_i^2} \frac{1}{2} (y_i - \bar{y}_i)^2 = 1$$

$$g_i = \frac{d}{dy_i} \frac{1}{2} (y_i - \bar{y}_i)^2 = -(y_i - \bar{y}_i)$$

Science

$$\text{Similarity Score} = \frac{1}{2} \frac{(g_1 + g_2 + \dots + g_n)^2}{(h_1 + h_2 + \dots + h_n + \lambda)}$$



Cả (1) và (2) đều có cùng optimization point P

$$\text{Similarity Score} = \frac{(\sum \text{Residual})^2}{\text{Number of Residual} + \lambda}$$

XGBoost Classification: Similarity Score

1

$$\sum_{i=1}^n \mathcal{L}(y_i, \bar{y}_i + P) + \frac{1}{2} \lambda P^2$$

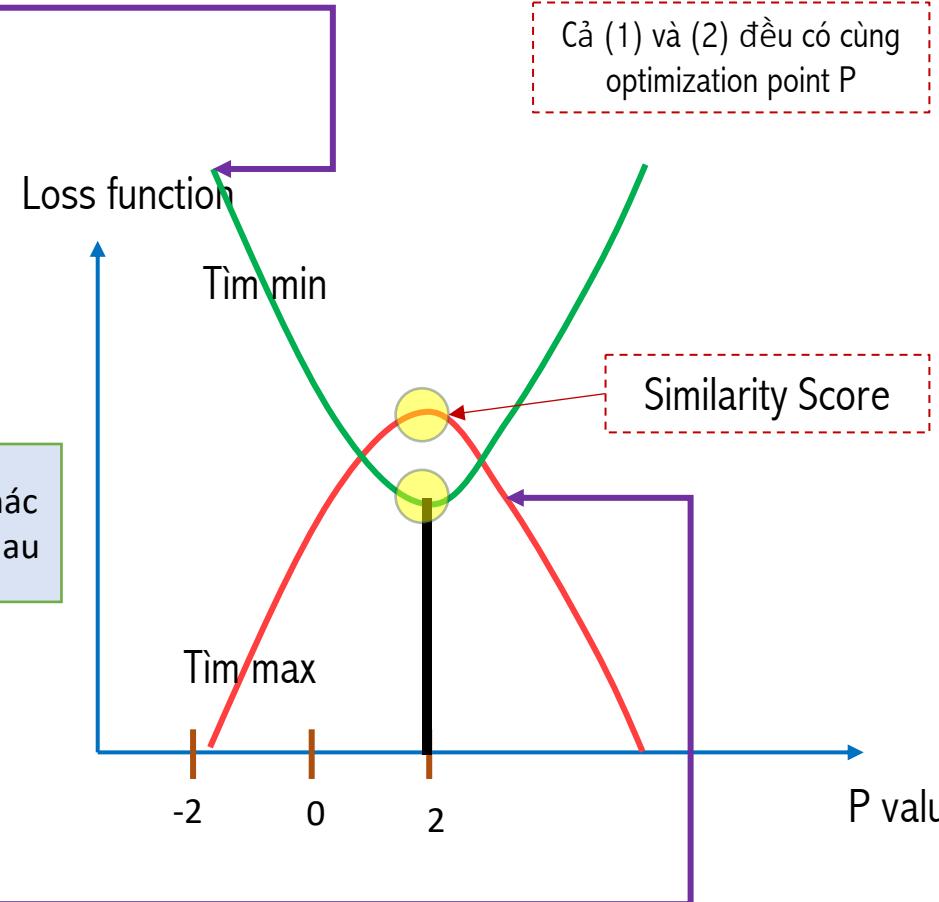
Rất khó để tìm optimization, nên cũng ta sẽ sắp sỉ hàm loss bằng Second Order Tayler Approximation

$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + \left[\frac{d}{dy_i} \mathcal{L}(y_i, \bar{y}_i) \right] P + \frac{1}{2} \left[\frac{d^2}{dy_i^2} \mathcal{L}(y_i, \bar{y}_i) \right] P^2$$

$$\mathcal{L}(y_i, \bar{y}_i + P) \approx \mathcal{L}(y_i, \bar{y}_i) + gP + \frac{1}{2} hP^2$$

2

$$-1 \times (g_1 + g_2 + \dots + g_n)P + -1 \times \frac{1}{2} (h_1 + h_2 + \dots + h_n + \lambda)P^2$$



$$g_i = -(y_i - \bar{y}_i)$$

$$h_i = \bar{y}_i \times (1 - \bar{y}_i)$$

Science

$$\text{Similarity Score} = \frac{1}{2} \frac{-(g_1 + g_2 + \dots + g_n)^2}{(h_1 + h_2 + \dots + h_n + \lambda)}$$

$$\text{Similarity Score} = \frac{(\sum \text{Residual})^2}{\sum \bar{y}_i \times (1 - \bar{y}_i) + \lambda}$$

Outline

- Regularization
- XGBoost For Regression
- XGBoost For Classification
- XGBoost: Mathematical Explanation
- XGBoost: Optimization
- Example

How to Handle Missing Value

Sparse-Aware Split Finding

Initial Prediction for Drug Effectiveness

0.5

Dosage	Drug Effectiveness
10	-7
???	-3
21	7
25	8
5	-5
???	-2



Dosage	Drug Effectiveness	Residual
10	-7	-7.5
???	-3	-3.5
21	7	6.5
25	8	7.5
5	-5	-5.5
???	-2	-2.5

How to Handle Missing Value

Table without missing values

Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5

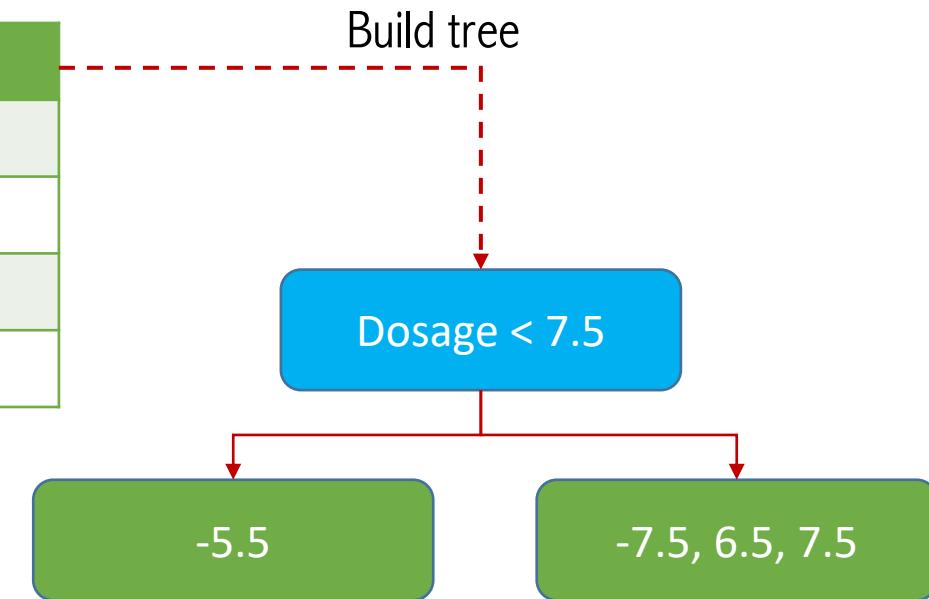


Table with missing values

Dosage	Drug Effectiveness	Residual
???	-3	-3.5
???	-2	-2.5

How to Handle Missing Value

Table without missing values

Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5

Compute Gain
Information Left

Dosage < 7.5

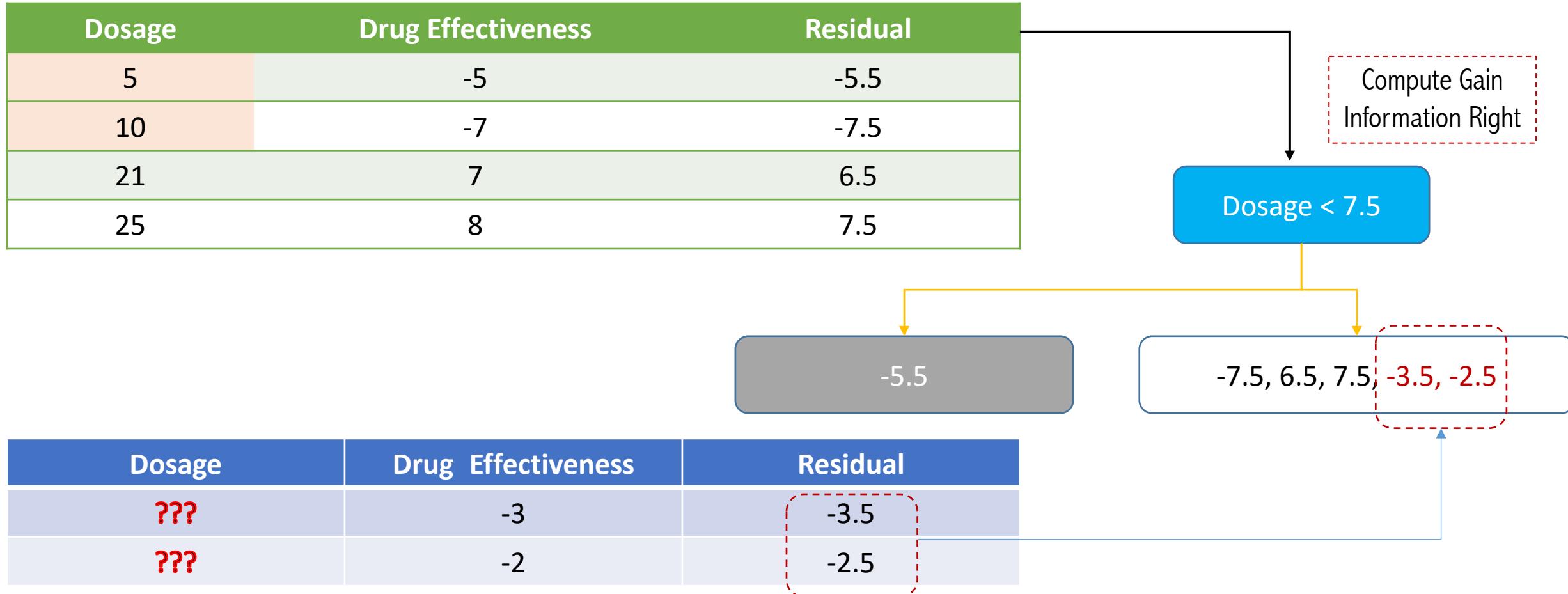
-5.5, -3.5, -2.5

-7.5, 6.5, 7.5

Table with missing values

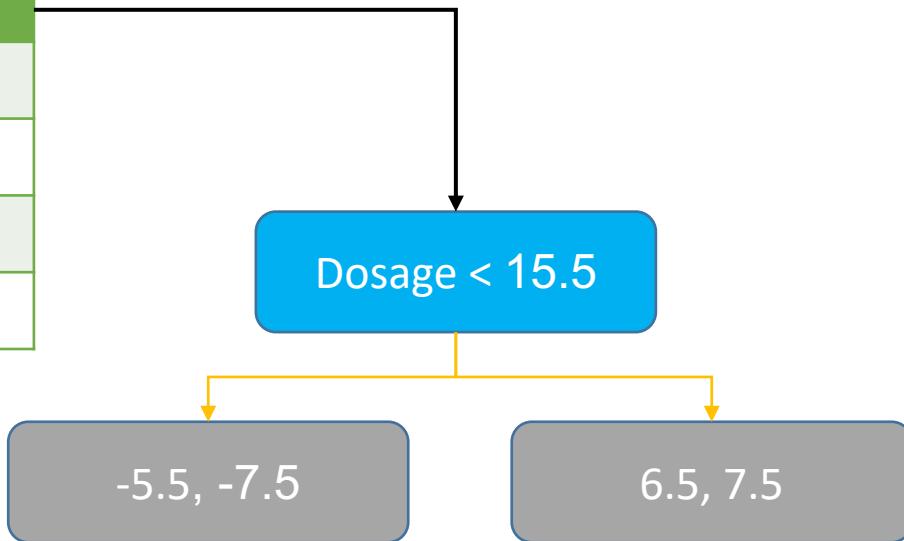
Dosage	Drug Effectiveness	Residual
???	-3	-3.5
???	-2	-2.5

How to Handle Missing Value



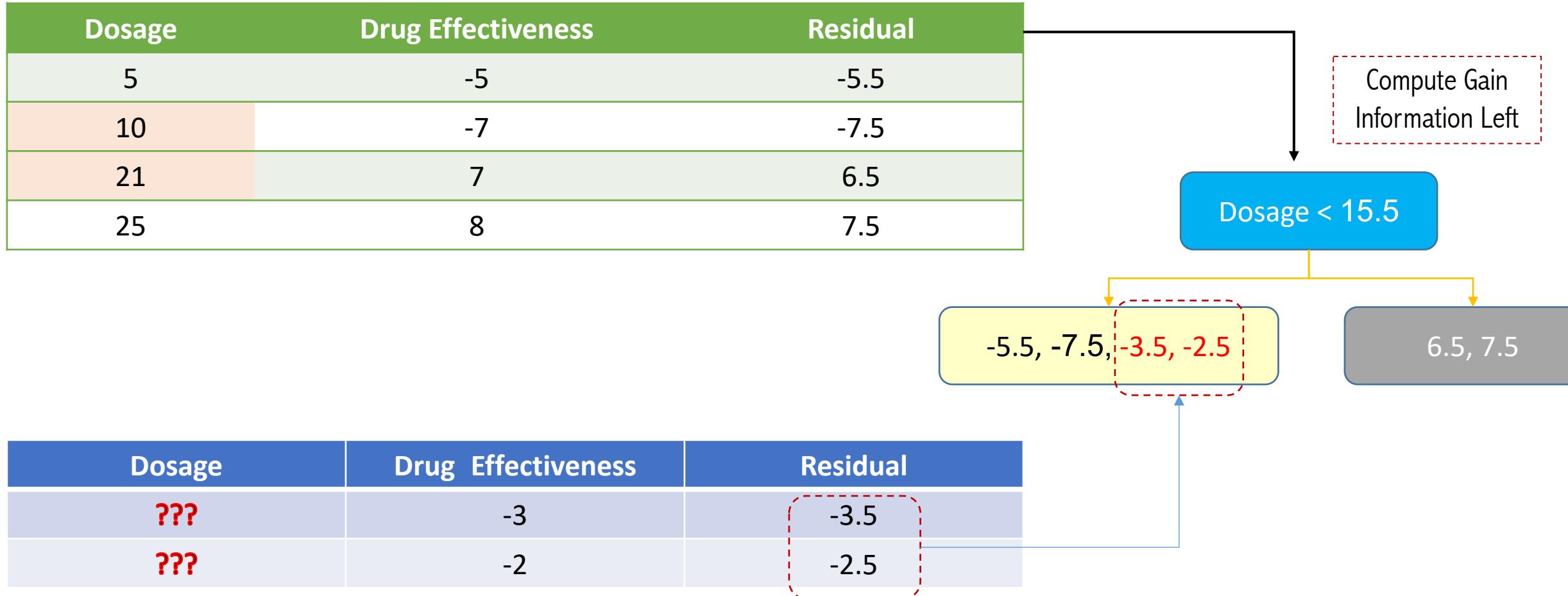
How to Handle Missing Value

Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5



Dosage	Drug Effectiveness	Residual
???	-3	-3.5
???	-2	-2.5

How to Handle Missing Value



How to Handle Missing Value

Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5

Compute Gain
Information Right

Dosage < 15.5

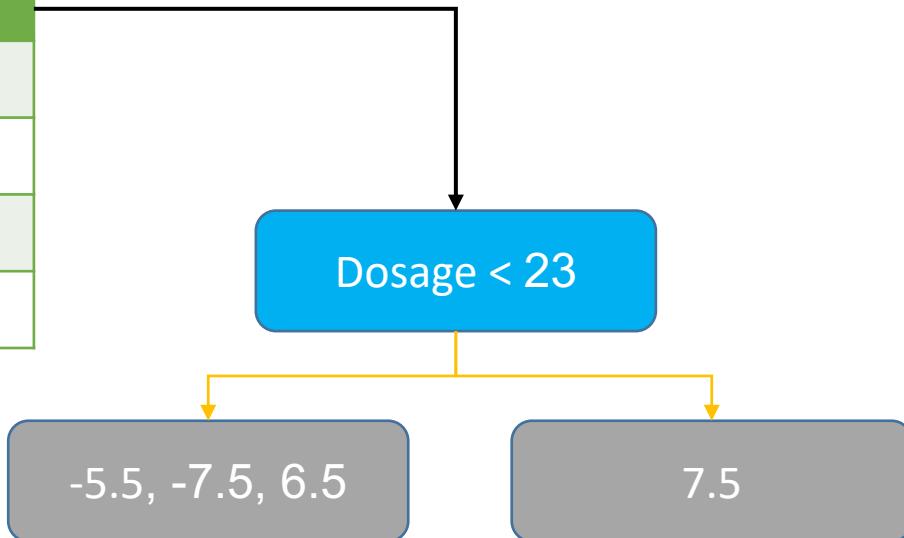
-5.5, -7.5

6.5, 7.5, 3.5, -2.5

Dosage	Drug Effectiveness	Residual
???	-3	-3.5
???	-2	-2.5

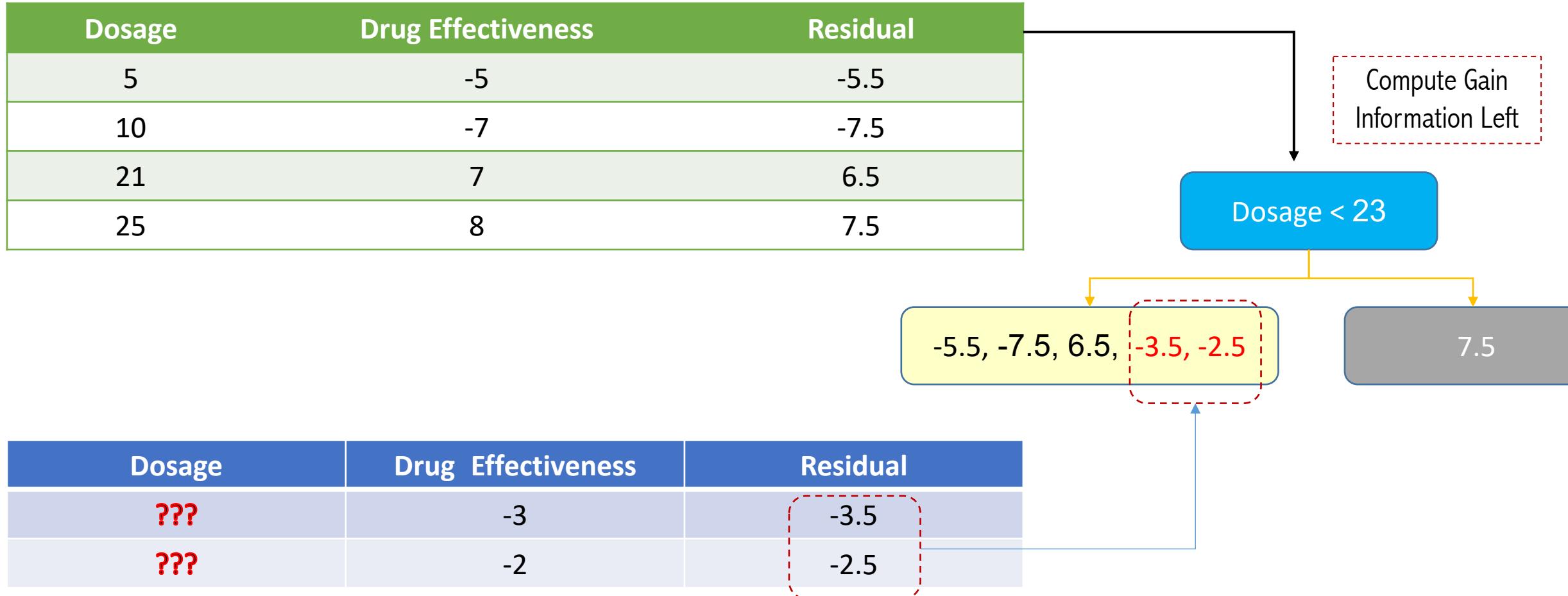
How to Handle Missing Value

Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5



Dosage	Drug Effectiveness	Residual
???	-3	-3.5
???	-2	-2.5

How to Handle Missing Value



How to Handle Missing Value

Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5

Compute Gain
Information Right

Dosage < 23

-5.5, -7.5, 6.5

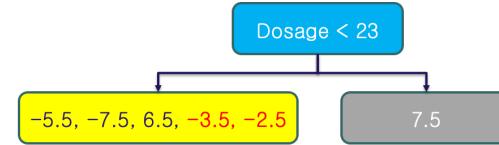
7.5, 3.5, -2.5

Dosage	Drug Effectiveness	Residual
???	-3	-3.5
???	-2	-2.5

Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5

Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5

Compute Gain value



Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5

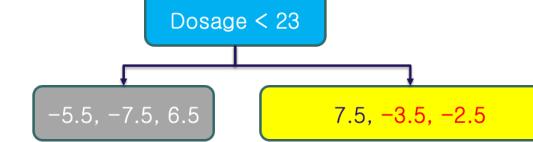
Compute Gain value



Dosage	Drug Effectiveness	Residual
10	-7	-7.5
21	7	6.5
25	8	7.5
5	-5	-5.5

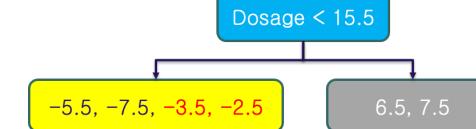
MAXIMUM GAIN VALUE

Compute Gain value



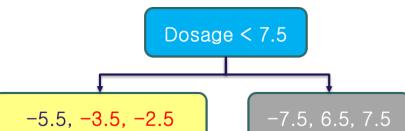
Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5

Compute Gain value



Dosage	Drug Effectiveness	Residual
21	7	-7.5
25	8	6.5
5	-5	7.5
10	-7	-5.5

Compute Gain value



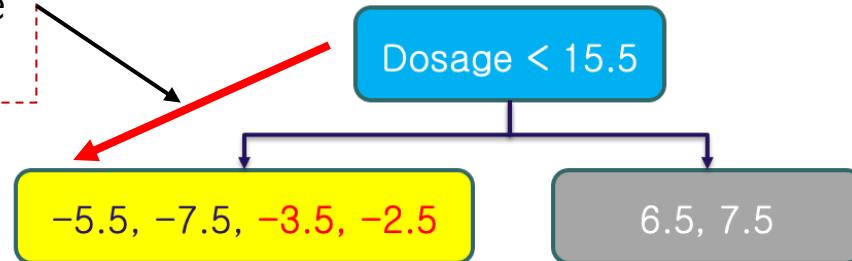
Dosage	Drug Effectiveness	Residual
-5.5	-3.5, -2.5	-7.5
-7.5, 6.5, 7.5	-3.5, -2.5	6.5
5	-5	7.5
10	-7	-5.5

How to Handle Missing Value

Dosage	Drug Effectiveness	Residual
5	-5	-5.5
10	-7	-7.5
21	7	6.5
25	8	7.5

Rẻ nhánh mặc định cho tất cả các missing value trong Dosage

Compute Gain value



Dosage	Drug Effectiveness	Residual
???	-3	-3.5
???	-2	-2.5

Outline

- **Regularization**
- **XGBoost For Regression**
- **XGBoost For Classification**
- **XGBoost: Mathematical Explanation**
- **How To Fill Missing Values**
- **Example**

Sales Prediction

TV	Radio	Newspaper	Sales
230.1	37.8	69.2	22.1
44.5	39.3	45.1	10.4
17.2	45.9	69.3	12
151.5	41.3	58.5	16.5
180.8	10.8	58.4	17.9
8.7	48.9	75	7.2
57.5	32.8	23.5	11.8
120.2	19.6	11.6	13.2
8.6	2.1	1	4.8
199.8	2.6	21.2	15.6



CUSTOM CODE

```
class AIVN_XBoostRegressor():
    '''XBoost from Scratch
    ...
    def __init__(self, params, random_seed=None):
        import numpy as np

        class SquaredErrorObjective():
            def loss(self, y, pred): return np.mean((y - pred)**2)
            def gradient(self, y, pred): return pred - y
            def hessian(self, y, pred): return np.ones(len(y))

    params = {
        'learning_rate': 0.1,
        'max_depth': 5,
        'subsample': 0.8,
        'reg_lambda': 1.5,
        'gamma': 0.0,
        'min_child_weight': 25,
        'base_score': 0.0,
        'tree_method': 'exact',
    }
    num_boost_round = 50

    # train the from-scratch XGBoost model
    model_scratch = AIVN_XBoostRegressor(params, random_seed=42)
    model_scratch.fit(X_train, y_train, SquaredErrorObjective(), num_boost_round)

    pred_scratch = model_scratch.predict(X_test)
    print(f'Our implementation: {SquaredErrorObjective().loss(y_test, pred_scratch)}')
```

Vinh scratch score: 1.7017064684759806

Sales Prediction



```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

#READ DATA
data = pd.read_csv("/content/advertising.csv")

#X,y
X = data.iloc[:, :-1]
y = data.iloc[:, -1]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=100)

from xgboost import XGBRegressor

model = XGBRegressor()
model.fit(X_train, y_train)

pred_scratch = model.predict(X_test)
print(f'XGBRegressor Lib: {SquaredErrorObjective().loss(y_test, pred_scratch)}')

XGBRegressor Lib: 1.7209355894486629
```

