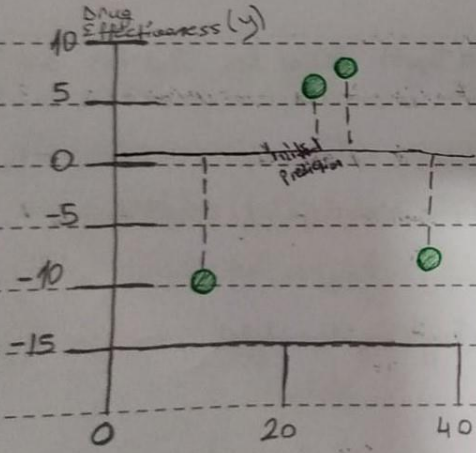# XGBoost

## a) Regression:

Örneği kolay anlatmak için super simple bir training set kullanacağız.



1) The first step in fitting XGBoost to the Training Data is to make initial prediction.

Bu initial herhangi bir sayı olabilir. Ancak Default'da hem Regresyon hem Sınıflandırma için 0.5 alınıyor.

2) Calculate first (initial) residuals. Difference between Observed and Predicted values.

3) Tüm ağaçlar single leaf ile başlar. Bu single leaf'te ilk olarak Residualler olur.

4) Bu residualler için "Similarity Score" hesaplanır.

$$\text{Similarity Score} = \frac{(\text{Sum of the Residuals})^2}{\text{Number of residuals} + \lambda}$$ (Regularization parametre)
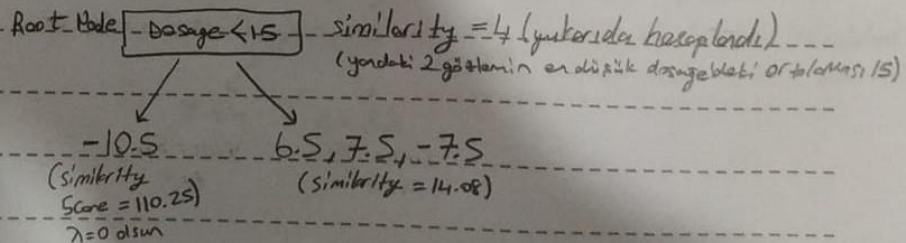
$$\boxed{-10.5, \, 6.5, \, 7.5, \, -7.5} \rightarrow \text{ilk leaf ilk residualler ile başlar demiştik.}$$

$$\text{Sim. score} = \frac{(-10.5 + 6.5 + 7.5 - 7.5)^2}{4 + \lambda_{\rightarrow 0 \text{ olsun şimdilik}}} = 4$$

5) İlk leaf similarity = 4 oldu. Şimdi şunu düşünmeliyiz, Residualleri 2 gruba böler isek daha iyi cluster etmiş olur muyuz.

6) Bunun için ilk leaf node'u en iyi bölecek Dosage (x) karar vermek lazım. En iyi demek en fazla "Gain" sağlayan demek. Bunun için Bağımsız değişkenin en küçük 2 değerinin ortalamasını alarak denemeye başlarız ve her denemede bir sonraki 2 gözlemin ortalamasını alarak leaf denenir.

Root Node $\boxed{\text{Dosage} < 15}$ similarity = 4 (yukarıda hesaplandı)
(yandaki 2 gözlemin en düşük dosageleki ortalaması 15)

$-10.5$        $6.5, 7.5, -7.5$
(Similarity Score = 110.25)    (similarity = 14.08)
$\lambda = 0$ olsun

Gain = Left similarity + Right sim - Root sim $\Rightarrow$ 110.25 + 14.08 - 4 = 120.33

../..../....

Now, we calculated the Gain for the threshold "Dosage <15", ----
we can compare it to the "Gain" calculated for other Thresholds.

↳

Dosage <15 , Dosage <22.5 , Dosage <30 Gain'lerine
bakılır. En yüksek Gain veren Root olarak seçilir. Bu örnekte
en yüksek Gain Dosage <15 iken alıyor.

<center>Dosage <15</center>

```
        ↙              ↘
    -10.5            6.5, 7.5, -7.5  → similarity = 14.08
  (Burada tek Residual      (Burası split edilir)
   olduğu için, we can't
   split it any further).
        λ=0
```

22.5 için :   Bu durumda; split edilecek node için yine yukarıdaki
             gibi threshold'lar denemeye başlanır.

<center>Dosage <15</center>

```
        ↙              ↘
    -10.5          Dosage <22.5 (ilk olarak 22.5 denensin)
                       ↙          ↘
                     6.5         7.5, -7.5
                (similarity =   (similarity = 0)
                  42.25)
```

Gain = 42.25 + 0 - 14.08 = 28.17

30 için :        Dosage <15

```
        ↙              ↘
    -10.5          Dosage <30 → sim = 14.08
                       ↙          ↘
                    6.5, 7.5      -7.5
                    sim = 98     sim = 56.25
                     λ=0          λ=0
```

Gain = 98 + 56.25 - 14.08 = 140.17  } Dosage 30 seçildi.

7) Tree'yi Prune etme i → γ (gamma) = 130 olsun diyelim

★ Gain - γ (gamma) < 0 ise , we remove the Branch.

Mesela → Dosage < 30 olduğu durumda Gain = 140.17 idi
0 node'u remove edemeyiz çünkü 140.17 > 130.

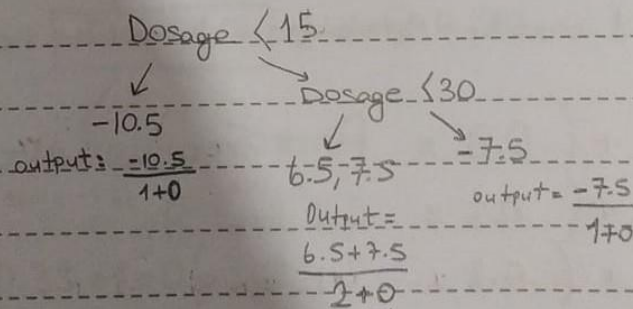8) λ = 0 aldık hep işlem kolaylığı için. λ → Regularization parameter and reduce the prediction's sensitivity to individual observations. ①"When λ > 0 , the similarity scores are smaller." ②" When λ > 0 , the Gain values are smaller also."

λ → Training Data'da overfit'i engeller.

9) Leaf Output Value = $\dfrac{\text{Sum of Residuals}}{\text{Number of Residuals} + \lambda}$

Yani ;

Dosage < 15

↓
-10.5
output: $\dfrac{-10.5}{1+0}$

→ Dosage < 30

↓                    ↘
6.5, 7.5              = 7.5
Output =             output = $\dfrac{-7.5}{1+0}$
$\dfrac{6.5+7.5}{2+0}$

↳ λ > 0 olduğunda , it reduce the amount that this individual observation adds to the overall prediction.

..../..../....

10) Bu durumda, örneğin;

**New Predicted value**, for observation with Dosage = 10
is calculated as follows: ──────────── Rate

Learning

Original Prediction + ($ε$ (eta) × -10.5)

$$0.5 + 0.3 \times \underset{-10.5}{\overset{\text{Dosage} < \text{is}}{\downarrow}} = -2.65$$

(new prediction)

output = -10.5
$λ = 0$ iken

★ New residual is smaller than the previous (initial) residual.

✗ Aynı şekilde; New prediction for observation with Dosage = 20

$$0.5 + (0.3 \times 7) = 2.6$$
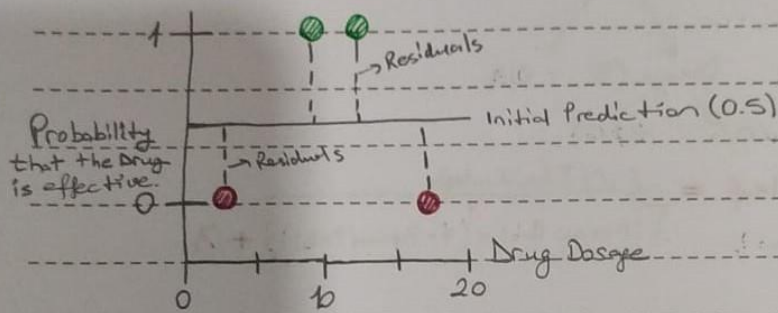
↳ New residual is smaller then before again.

11) Build another tree based on the newest "Residuals".

$$0.5 + ( 0.3 \times \text{First Tree} )$$
$$+ ( 0.3 \times \text{2nd Tree} )$$
$$± ( 0.3 \times \text{3rd Tree} )$$

## b) Classification:

Regresyona çok benziyor.



$$\text{Similarity Score for Residuals} = \frac{\left(\sum \text{Residuals}_i\right)^2}{\sum \left[\text{Previous Prob}_i \times (1 - \text{Previous Prob}_i)\right] + \lambda}$$

Örneğin

$$\boxed{-0.5, 0.5, 0.5, -0.5} \rightsquigarrow \text{similarity} = 0$$

$$\boxed{\text{Dosage} < 15} \rightsquigarrow \text{similarity} = 0 \text{ idi.}$$

$$-0.5, 0.5, 0.5 \qquad -0.5$$
$$\text{sim} = 0.33 \qquad \text{sim} = 1$$

$$\text{Gain} = 0.33 + 1 - 0 = 1.33$$

Oluşan ilk ağacınız aşağıdaki gibi olsun :

Dosage < 15

Dosage < 5   = 0.5

-0.5   0.5, 0.5

$$\text{Output value for a leaf} = \frac{(\sum \text{Residuals}_i)}{\sum [\text{Previous Prob}_i \times (1 - \text{Previous Prob}_i)] + \lambda}$$

Dosage < 15

Dosage < 5

-0.5
output = -2
$\lambda$=0

-0.5
output= -2
$\lambda$=0 iken

0.5, 0.5
output = 2
$\lambda$=0

-0.5
output =-2
$\lambda$=0

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

**örneğin**

"New Predicted" value for observation whose Dosage=2 :

Predicted Drug Effectiveness

Initial = 0.5   +   ( 0.3 × -2 )  = - 0.6

(Learning rate, output)

↳ log (odds) = 0

↳ Log (odds) = -0.6 | Convert log (odds) into probability.

Use Logistic func.

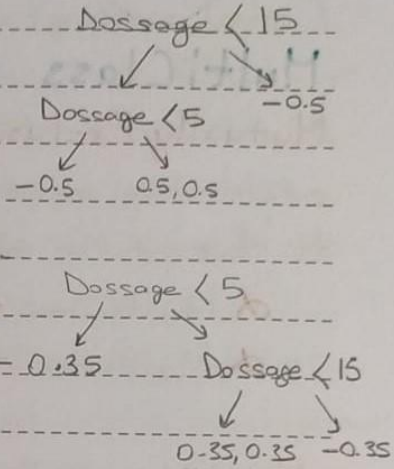$$\text{Probability} = \frac{e^{-0.6}}{1+e^{-0.6}} = 0.35$$

→ Original residual was = 0.5 , with new prediction it becomes
-0.35.

Initial Prediction = 0.5

log(odds) = 0 + 0.3 × [Dossage < 15 tree diagram]

Dossage < 15
↙ ↘
Dossage < 5        −0.5
↙ ↘
−0.5    0.5, 0.5

+ 0.3 × [tree diagram]

Dossage < 5
↙ ↘
−0.35    Dossage < 15
↙ ↘
0.35, 0.35    −0.35

## OPTIMIZATIONS 8

Large dataset olduğunda bir node karar vermek için, her farklı değişkenin her threshold değerine bakmak çok maliyetli.
Yani, Değişken X için tek tek arka arkaya gelen her 2 gözlemin ortalamasına bakıp en iyi threshold buluyor idik ve bunu tüm değişkenler (Y, Z ...) için yapıyorduk. Large dataset olduğunda thresholdları quartile noktalardan seçer XGBoost, böylece bakması gereken threshold sayısını azaltmış olur.

## How Xgboost Deal with Missing Values?

Training kümesinde missing değerlerin sağ veya sol node'dan hangisine gideceğine loss'u minimize etmesi bakımından karar verir ve test kümesindeki missing değerler o node'a gider default olarak. Training kümede NAN yok ise bu durumda her node için Default missing yönü seçer ve Test için bu yön kullanılır.