## 1. XGBoost for Regression

1. Similarity Score = $\dfrac{\text{Sum of residuals squared}}{\text{Number of residuals} + \lambda}$ ↳ Regularization

* Note: "$\lambda$" is a regularization parameter & when $\lambda > 0$, it results in more pruning, by shrinking the similarity scores, and it results in smaller output values for the leaves.

2. Gain = Left similarity + Right similarity − Root similarity.

3. Pruning: Gain − $\gamma$ $\begin{cases} > 0 : \text{keep the node} \\ < 0 : \text{drop the node} \end{cases}$

## 2. XGBoost for classification:

# Similarity $= \dfrac{\left(\sum \text{Residuals}\right)^2}{\sum \left[ \text{Prev. Probability}_i * (1 - \text{Prev. Probability}_i) \right] + \lambda}$

## # MATHS:

# XGBoost uses loss functions to build trees by minimizing this equation:

$$\underbrace{\left[ \sum_{i=2}^{n} \mathcal{L}(y_i, P_i) \right]}_{\substack{\text{Loss} \\ \text{function}}} + \underbrace{\frac{1}{2} \lambda \, O_{\text{value}}^2}_{\text{Regularization}}$$

$\begin{cases} \lambda > 0 ; \text{shrinks } O_{\text{value}} \\ \lambda < 0 ; \text{expands } O_{\text{value}} \end{cases}$

→ The goal is to find the output value ($O_{\text{value}}$) for the leaf that minimizes the whole equation.
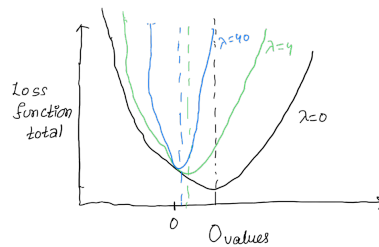
Prediction = $\boxed{0.5}$ + 

(XGBoost tree)



Figure: $O_{\text{values}}$ Vs Loss function total for multiple regularization

* The more emphasis we give the regularization penalty by increasing "$\lambda$" the optimal $O_{\text{value}}$ gets closer to 0.

# XGBoost uses the second order Taylor approximation for both Regression & classification.

i.e.

$$\mathcal{L}(y, P_i + O_{\text{value}}) \approx \mathcal{L}(y, P_i) + \left[ \frac{d}{dP_i} \mathcal{L}(y, P_i) \right] O_{\text{value}} + \frac{1}{2} \left[ \frac{d^2}{dP_i^2} \mathcal{L}(y, P_i) \right] O_{\text{value}}^2$$

$$= L(y, P_i) + \underset{\substack{\downarrow \\ \text{"Gradient"}}}{g} O_{\text{value}} + \frac{1}{2} \underset{\substack{\downarrow \\ \text{"Hessian"}}}{h} O_{\text{value}}^2$$

Now, to minimize:

$$\left[ \sum_{i=1}^{n} \mathcal{L}(y_i, P_i^0 + O_{\text{value}}) \right] + \frac{1}{2} \lambda \, O_{\text{value}}^2$$

to determine optimal $O_{\text{value}}$.

We can write it as:

$$\left.\begin{array}{l} \underline{L\,(y_1,\,P_1^0)} + g_1\,O_{value} + \tfrac{1}{2}\,h_1\,O^2_{value} + \\[4pt] \underline{L\,(y_2,\,P_2^0)} + g_2\,O_{value} + \tfrac{1}{2}\,h_2\,O^2_{value} + \cdots + \\[4pt] \underline{L\,(y_n,\,P_n^0)} + g_n\,O_{value} + \tfrac{1}{2}\,h_n\,O^2_{value} + \tfrac{1}{2}\,\lambda\,O^2_{value} \end{array}\right\} \text{"Taylor Series Expansion"}$$

No contribution in $O_{value}$

So, for optimization we reduce it to:

$$(g_1 + g_2 + g_3 + \cdots + g_n)\,O_{value} +$$
$$\tfrac{1}{2}\,(h_1 + h_2 + h_3 + \cdots + h_n + \lambda)\,O^2_{value}$$

$\dfrac{d}{dO_{value}}$

Now minimize the function by setting its derivative $= 0$.

i.e.

$$\rightarrow (g_1 + g_2 + g_3 + \cdots + g_n) + (h_1 + h_2 + h_3 + \cdots + h_n + \lambda)\,O_{value} = 0$$

Or,

$$\boxed{O_{value} = \dfrac{-(g_1 + g_2 + g_3 + \cdots + g_n)}{(h_1 + h_2 + h_3 + \cdots + h_n + \lambda)}}$$

(A) For Regression:

* $L(y_i, P_i) = \tfrac{1}{2}(y_i - P_i)^2$

$\Rightarrow g_i = \dfrac{d}{dP_i}\,\tfrac{1}{2}(y_i - P_i)^2 = -(y_i - P_i) = -\text{Residual}$

lly, $h_i = \dfrac{d^2}{dP_i^2}\,\tfrac{1}{2}(y_i - P_i)^2 = \dfrac{d}{dP_i}-(y_i - P_i) = 1$

So,
$$\boxed{O_{value} = \dfrac{\text{sum of residuals}}{\text{Number of residuals} + \boldsymbol{\lambda}}}$$

This is the specific formula for the output value for a leaf when using XGBoost for regression.

(B) For Classification:

* $L(y_i, P_i) = -\big[\,y_i\,\log(P_i) + (1 - y_i)\,\log(1 - P_i)\,\big]$

$\rightarrow g_i = -(y_i - P_i) = -\text{Residuals}$
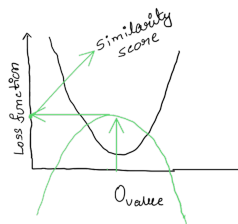
$\rightarrow h_i = P_i * (1 - P_i)$

So,
$$\boxed{O_{value} = \dfrac{\sum \text{Residuals}_i}{\sum \big[\text{Prev. Probability}_i * (1 - \text{Prev. Probability}_i)\big] + \lambda}}$$

(C) Similarity score for tree growing:

The first thing XGBoost does is multiply the optimization eqn by "$-1$" to flip the parabola upside down.

i.e.



$$-1 * \Big[\,(g_1 + g_2 + \cdots + g_n)\,O_{value} + \tfrac{1}{2}\,(h_1 + h_2 + h_3 + \cdots + h_n + \lambda)\,O^2_{value}\,\Big]$$

use the $O_{value}$ from above & simplify:

$$\boxed{\text{Similarity score} = \dfrac{1}{2}\,\dfrac{(g_1 + g_2 + \cdots + g_n)^2}{(h_1 + h_2 + \cdots + h_n + \lambda)}}$$

* For regression

$$\text{Similarity score} = \frac{\text{sum of residuals, squared}}{\text{Number of residuals} + \lambda}$$

& cover = Number of residuals

Ands

* for **C**lassification

$$\text{Similarity score} = \frac{\left(\sum \text{Residuals}_i\right)^2}{\sum\left[\text{Prev. Probability}_i * (1 - \text{Prev. Probability}_i)\right] + \lambda}$$

& cover = $\sum\left[P_i^o * (1 - P_i)\right]$

#Extra:

```
  ┌──────────────┐         ┌──────────────────┐
  │   Initial    │──────▶  │   Residual       │───▶ Build model to predict
  │   Model      │         │  calculations    │     those residuals.
  └──────────────┘         └──────────────────┘          │
                                   ▲                      │
                                   │    ┌──────────────┐  │
                                   └────│ Add last model to │◀─┘
                                        │  Ensemble    │
                                        └──────────────┘
```
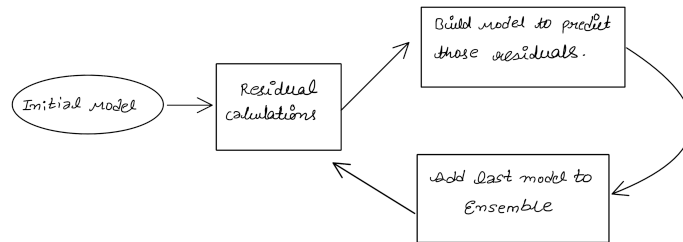
fig: General xgBoost Algorithm