

Why the Trees in XGBoost are Better!

Step 1: build a loss function that minimize both Bias and Variance

total loss function

Minimize the Bias

Minimize the Variance

$$\mathcal{L}(\phi) = \sum_i l(y_i, \hat{y}_i) + \sum_k \Omega(f_k)$$

Sum over the samples

Number of leaves

Sum over the trees

where $\Omega(f) = \underbrace{\gamma T}_{\text{Push the number of leaves to be lower}} + \frac{1}{2} \lambda \underbrace{\|w\|^2}_{\text{Push the predictions to be smaller}}$

Prediction at each leaf

Step 2: Approximate the loss function with a 2nd order Taylor expansion

prediction

target

At the t^{th} boosting iteration

Features

$$\underbrace{l(y_i, \hat{y}_i)}_{\text{sample level loss function}} = l(y_i, \underbrace{\hat{y}_i^{(t-1)}}_{\text{prediction of the target at the (t-1)-th iteration}} + \underbrace{f_t(\mathbf{x}_i)}_{\text{prediction of the new tree}})$$

Second order Taylor approximation

$$l(y_i, \hat{y}_i) \simeq l(y_i, \hat{y}_i^{(t-1)}) + \underbrace{g_i}_{\text{Gradient}} f_t(\mathbf{x}_i) + \frac{1}{2} \underbrace{h_i}_{\text{Hessian}} f_t^2(\mathbf{x}_i)$$
$$g_i = \partial_{\hat{y}^{(t-1)}} l(y_i, \hat{y}_i^{(t-1)})$$

Gradient

$$h_i = \partial_{\hat{y}^{(t-1)}}^2 l(y_i, \hat{y}_i^{(t-1)})$$

Hessian

TheAiEdge.io

Example: the cross-entropy loss function

```
1 loss = y * np.log(y_hat) + (1 - y) * np.log(1 - y_hat)
2 grad = -y / y_hat + (1 - y) / (1 - y_hat)
3 hess = y / (y_hat ** 2) + (1 - y) / ((1 - y_hat) ** 2)
```

Step 3: Just use the loss function depending on the new tree

$$\mathcal{L}^{(t)} \simeq \sum_{i=1}^n \left[\underbrace{l(y_i, \hat{y}_i^{(t-1)})}_{\text{Does not depends on the new tree}} + \underbrace{g_i f_t(\mathbf{x}_i)}_{\text{Depends on the new tree}} + \underbrace{\frac{1}{2} h_i f_t^2(\mathbf{x}_i)}_{\text{Depends on the new tree}} \right] + \underbrace{\Omega(f_t)}_{\text{Depends on the new tree}}$$
$$\Rightarrow \tilde{\mathcal{L}}^{(t)} = \sum_{i=1}^n \left[g_i f_t(\mathbf{x}_i) + \frac{1}{2} h_i f_t^2(\mathbf{x}_i) \right] + \Omega(f_t)$$

Step 4: Rearrange summation from samples to leaves and leaf instances

$$\tilde{\mathcal{L}}^{(t)} = \sum_{i=1}^n \left[g_i f_t(\mathbf{x}_i) + \frac{1}{2} h_i f_t^2(\mathbf{x}_i) \right] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

Sum on samples

prediction at leaf j

$$= \sum_{j=1}^T \left[\left(\sum_{i \in I_j} g_i \right) w_j + \frac{1}{2} \left(\sum_{i \in I_j} h_i + \lambda \right) w_j^2 \right] + \gamma T$$

Sum on leafs

Sum on instances in each leaf

Step 5: Find optimal weights by setting the derivative to zero

Find the minimum

$$\partial_{w_j} \tilde{\mathcal{L}}^{(t)} = 0$$
$$\sum_{i \in I_j} g_i$$

Solution

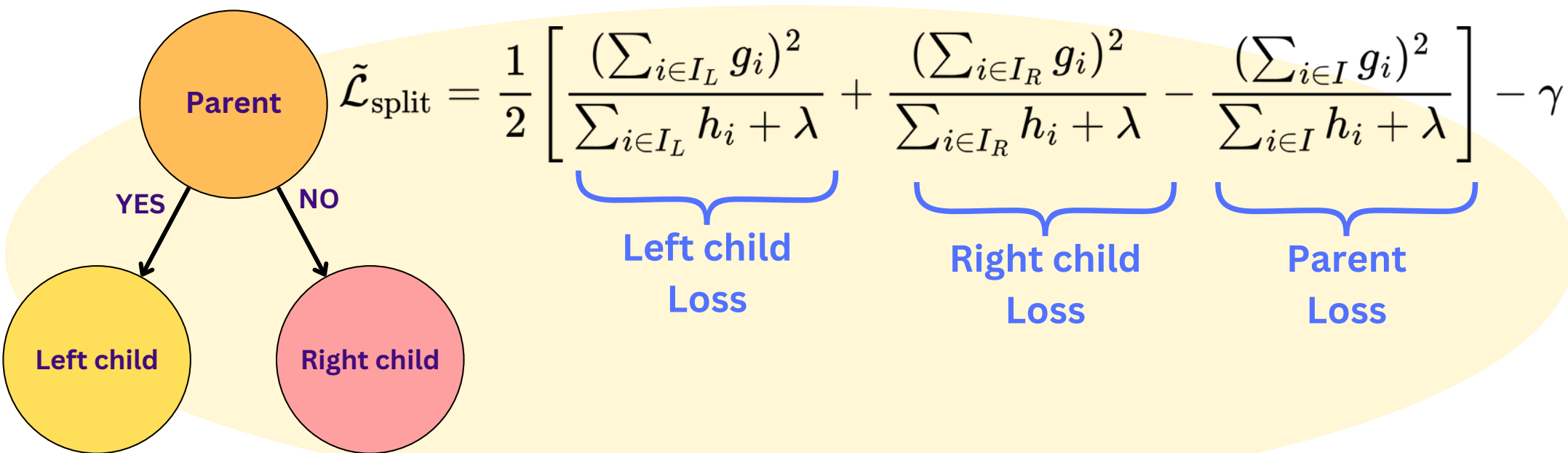
$$w_j^* = \frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$$

Step 6: Replace in the loss function

**Minimum
value of the
Loss**

$$\tilde{\mathcal{L}}^{(t)} = -\frac{1}{2} \sum_{j=1}^T \frac{(\sum_{i \in I_j} g_i)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma T$$

The best splits in the trees are the ones that decrease the loss



Step 7: Iterate and choose the splits that minimize the loss

