

XGBoost: A Scalable Tree Boosting System

Oct 9, 2023

Seungeun Lee

① | Recap

② | Split Finding Algorithms

③ | H/W Optimization

④ | Details

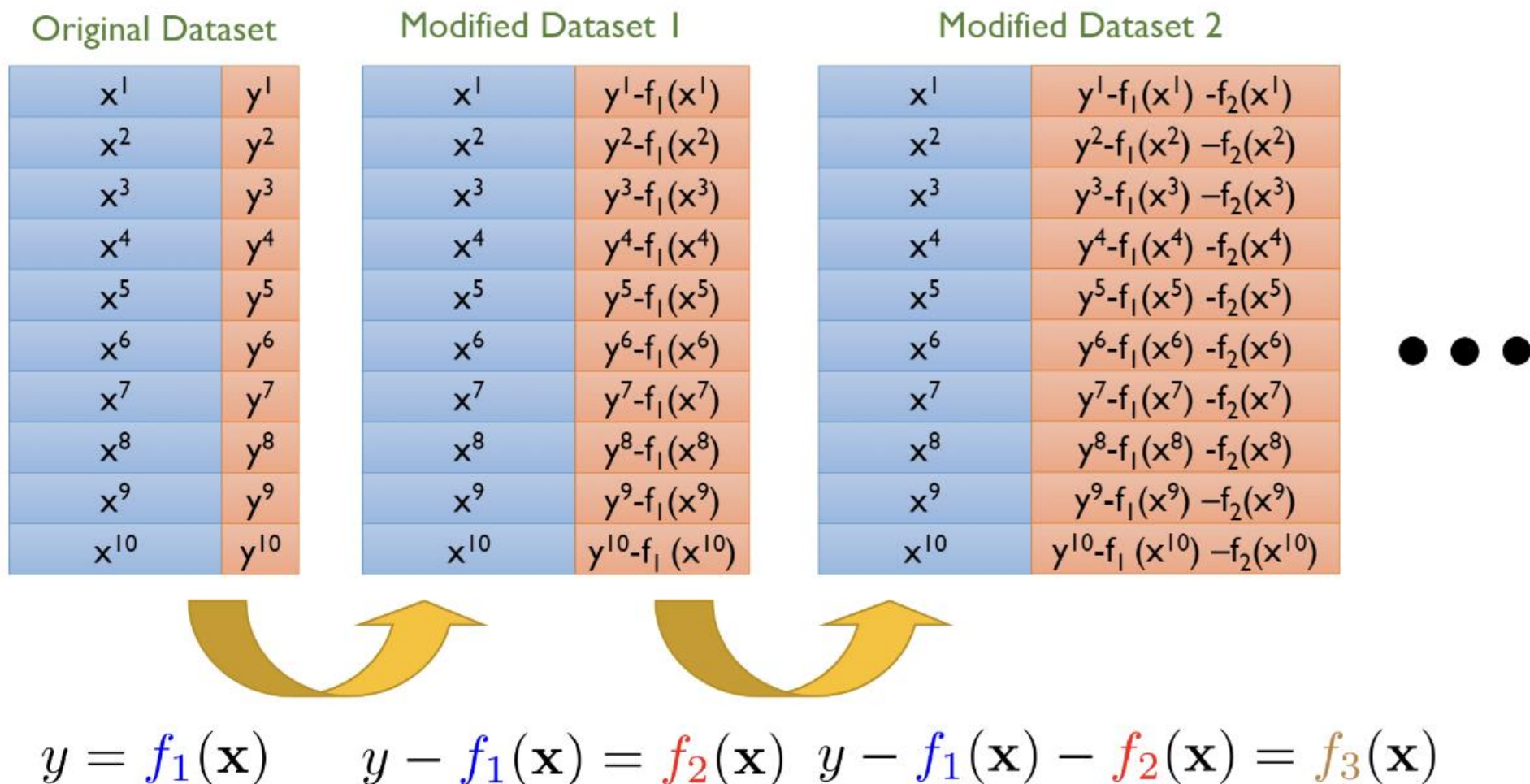
Recap

Gradient Tree Boosting

- expand $\Omega(f_k)$
- f_k : independent tree structure q and leaf weights w
- $I_j = \{i \mid q(x_i) = j\}$
- $$\begin{aligned}\tilde{L}^{(t)} &= \sum_{i=1}^n [g_i f_t(x_i) + \frac{1}{2} h_i (f_t)^2(x_i)] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T (w_j)^2 \\ &= \sum_{j=1}^T [(\sum_{i \in I_j} g_i) w_j + \frac{1}{2} (\sum_{i \in I_j} h_i + \lambda) (w_j)^2] + \gamma T\end{aligned}$$
- Note. $\operatorname{argmin}_x (Gx + \frac{1}{2} Hx^2) = -\frac{G}{H}, H > 0$
- For a fixed structure $q(x)$, we can compute the optimal weight w_j^* of leaf j by
- $$w_j^* = -\frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$$
- (scoring function to measure the quality of a tree structure q)
$$\tilde{L}^{(t)}(q) = -\frac{1}{2} \sum_{j=1}^T \frac{(\sum_{i \in I_j} g_i)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma T$$

Recap

- Why?



Recap

- How is this idea related to the gradient?
 - ✓ Loss function of the ordinary least square (OLS)

$$\min L = \frac{1}{2} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$$

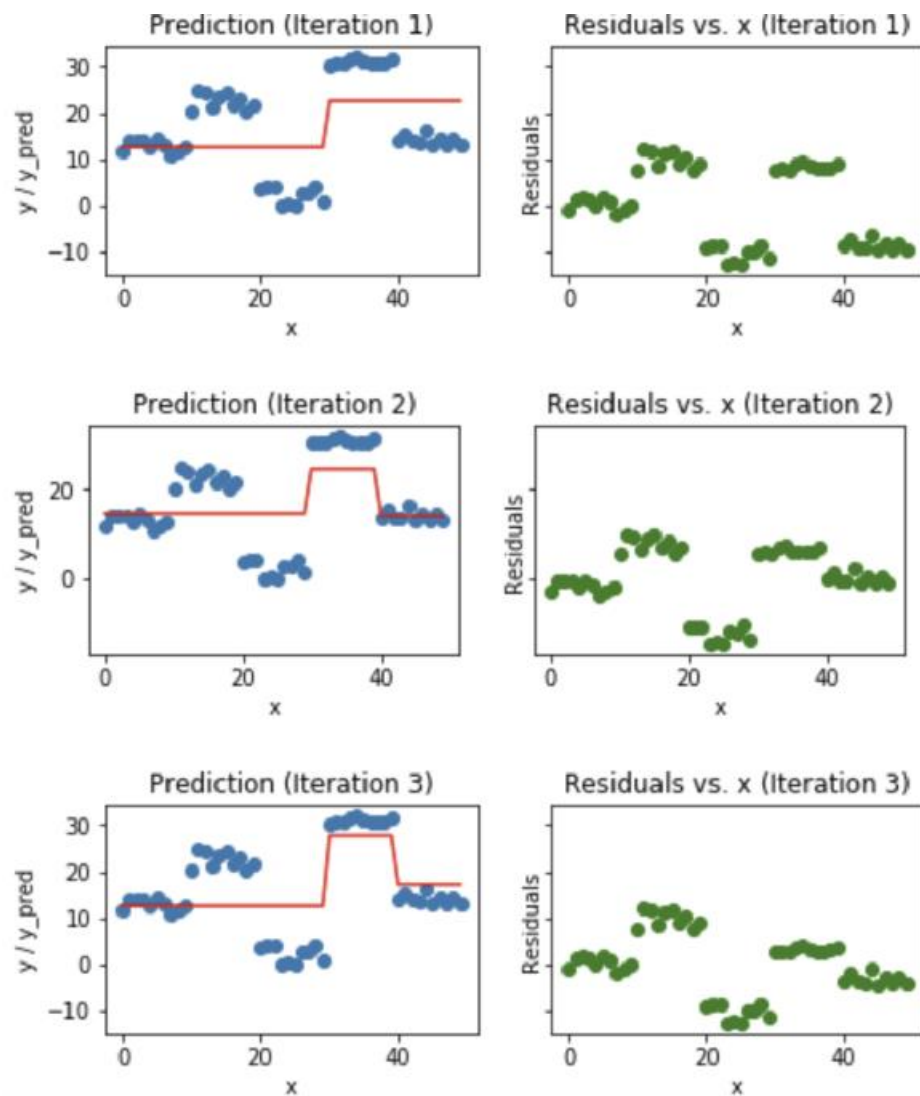
- ✓ Gradient of the Loss function

$$\frac{\partial L}{\partial f(\mathbf{x}_i)} = f(\mathbf{x}_i) - y_i$$

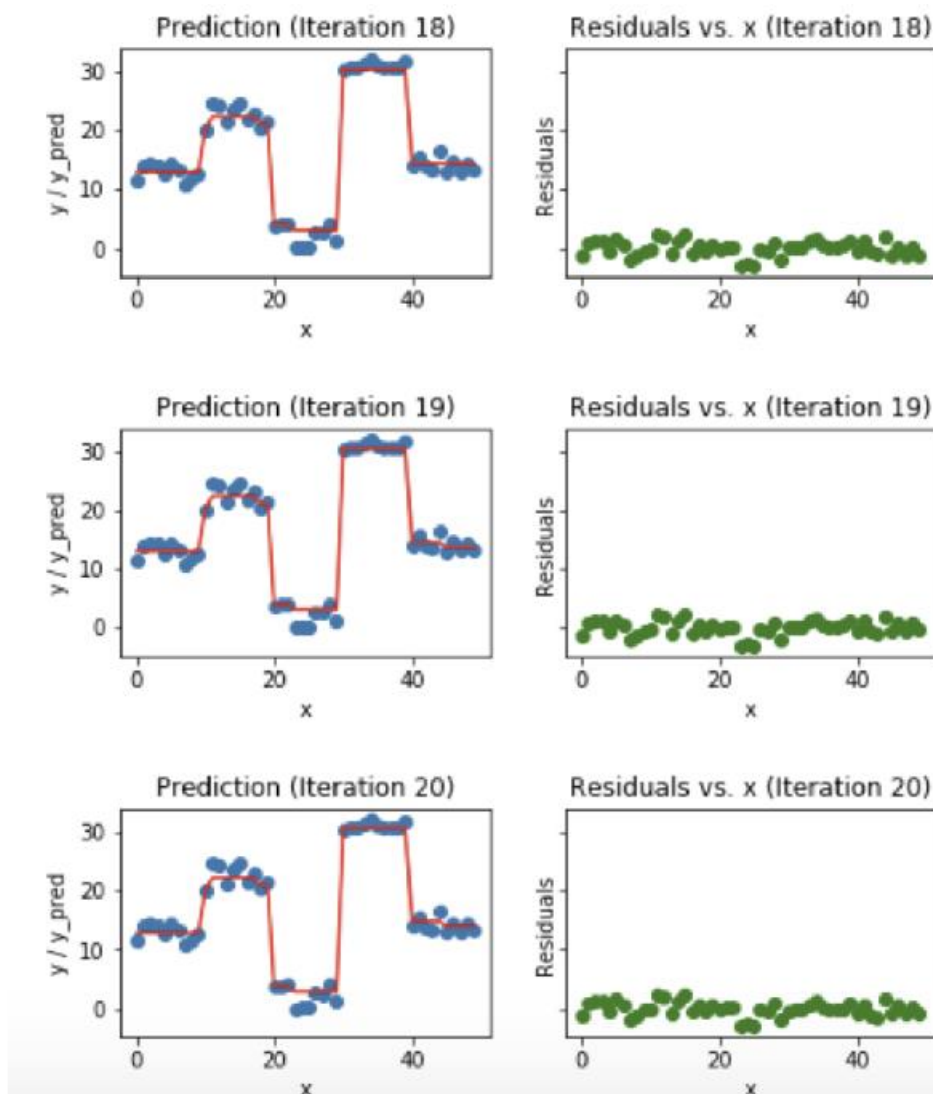
- ✓ Residuals are the negative gradient of the loss function

$$y_i - f(\mathbf{x}_i) = -\frac{\partial L}{\partial f(\mathbf{x}_i)}$$

Recap



Tree depth



Recap

- Gradient Boosting: Algorithm

1. Initialize $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$.

2. For $m = 1$ to M :

2.1 For $i = 1, \dots, N$ compute

$$g_{im} = \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i) = f_{m-1}(x_i)}$$

2.2 Fit a regression tree to the targets g_{im} giving terminal regions

$$R_{jm}, j = 1, \dots, J_m.$$

2.3 For $j = 1, \dots, J_m$ compute

Ground Truth, accumulated values

$$\gamma_{jm} = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma)$$

2.4 Update $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

3. Output $\hat{f}(x) = f_M(x)$.

Recap

- However, Overfitting
- Memorizes the uncertainty or noise of the data
- Shrinkage (minimizes the power of 'overfitted' models), Subsampling (maintains the size of the dataset), Early Stopping ...

Recap

Gradient Tree Boosting

- expand $\Omega(f_k)$
- f_k : independent tree structure q and leaf weights w
- $I_j = \{i \mid q(x_i) = j\}$
- $$\begin{aligned}\tilde{L}^{(t)} &= \sum_{i=1}^n [g_i f_t(x_i) + \frac{1}{2} h_i (f_t)^2(x_i)] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T (w_j)^2 \\ &= \sum_{j=1}^T [(\sum_{i \in I_j} g_i) w_j + \frac{1}{2} (\sum_{i \in I_j} h_i + \lambda) (w_j)^2] + \gamma T\end{aligned}$$
- Note. $\operatorname{argmin}_x (Gx + \frac{1}{2} Hx^2) = -\frac{G}{H}, H > 0$
- For a fixed structure $q(x)$, we can compute the optimal weight w_j^* of leaf j by
- $$w_j^* = -\frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$$
- (scoring function to measure the quality of a tree structure q)
$$\tilde{L}^{(t)}(q) = -\frac{1}{2} \sum_{j=1}^T \frac{(\sum_{i \in I_j} g_i)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma T$$

Recap

- $w_j^* = -\frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$
- (scoring function to measure the quality of a tree structure q) $\tilde{L}^{(t)}(q) = -\frac{1}{2} \sum_{j=1}^T \frac{(\sum_{i \in I_j} g_i)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma T$
- BUT we can make infinite number of trees
- We have to decide “when to SPLIT the trees” to make an optimal tree structure

Recap

- $L_{split} = \frac{1}{2} \left[\frac{(\sum_{i \in I_L} g_i)^2}{\sum_{i \in I_L} h_i + \lambda} + \frac{(\sum_{i \in I_R} g_i)^2}{\sum_{i \in I_R} h_i + \lambda} - \frac{(\sum_{i \in I} g_i)^2}{\sum_{i \in I} h_i + \lambda} \right] - \gamma$
- (Loss function before split) - (Loss function after split)

- (loss of left node after the split) + (loss of right node after the split) – (loss before the split)
- Choose the split with the maximized loss reduction
- Shrinkage, Column subsampling

① | Recap

② | Split Finding Algorithms

③ | H/W Optimization

④ | Details

Split finding Algorithms

Algorithm 1: Exact Greedy Algorithm for Split Finding

Input: I , instance set of current node

Input: d , feature dimension

$gain \leftarrow 0$

$G \leftarrow \sum_{i \in I} g_i, H \leftarrow \sum_{i \in I} h_i$

for $k = 1$ **to** m **do**

$G_L \leftarrow 0, H_L \leftarrow 0$

for j **in** $sorted(I, \text{by } \mathbf{x}_{jk})$ **do**

$G_L \leftarrow G_L + g_j, H_L \leftarrow H_L + h_j$

$G_R \leftarrow G - G_L, H_R \leftarrow H - H_L$

$score \leftarrow \max(score, \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{G^2}{H + \lambda})$

end

end

Output: Split with max score

- Exact Greedy Algorithm
- Find all possible split point greedily
- OOM Error, cannot be done under distributed settings

Split finding Algorithms

Algorithm 2: Approximate Algorithm for Split Finding

```
for  $k = 1$  to  $m$  do
    | Propose  $S_k = \{s_{k1}, s_{k2}, \dots, s_{kl}\}$  by percentiles on feature  $k$ .
    | Proposal can be done per tree (global), or per split(local).
end
for  $k = 1$  to  $m$  do
    |  $G_{kv} \leftarrow \sum_{j \in \{j | s_{k,v} \geq \mathbf{x}_{jk} > s_{k,v-1}\}} g_j$ 
    |  $H_{kv} \leftarrow \sum_{j \in \{j | s_{k,v} \geq \mathbf{x}_{jk} > s_{k,v-1}\}} h_j$ 
end
Follow same step as in previous section to find max
score only among proposed splits.
```

- Approximation Algorithm
- k : index of the variables (features), l : # of buckets
- 2 methods: per tree (global), per split (local)
- Epsilon as a hyperparameter

Split finding Algorithms

[illegible]

- Compute the gradient for each bucket and find the best split

[illegible]

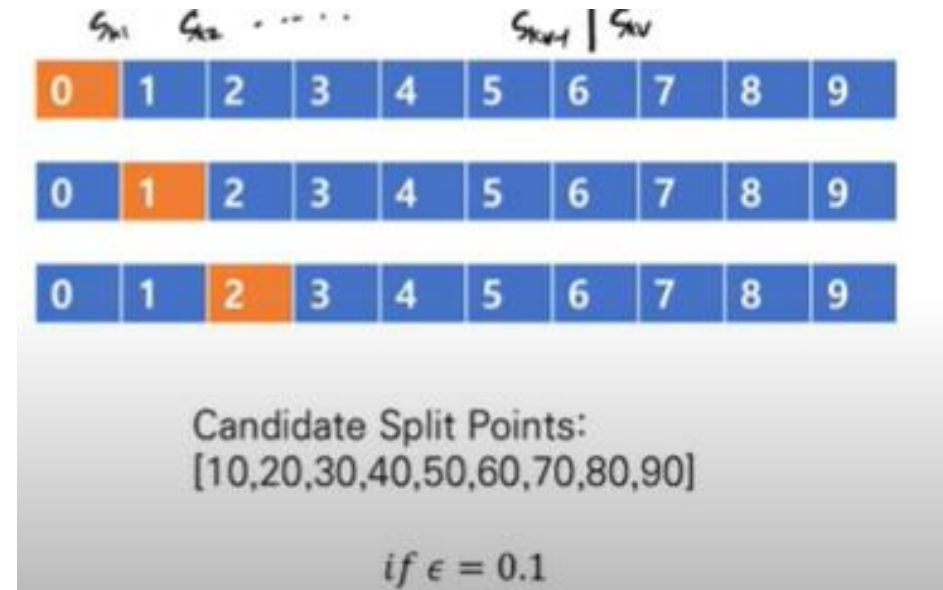
Best split point

- Ascending order
- # of buckets = 10 i.e. $\epsilon = 0.1$
- Exact greedy: 39 / approximation: $3 \times 10 = 30$

Split finding Algorithms

- Construct candidate split points w/ the percentile of feature distribution (epsilon)
- Update w/ G, H

Split finding Algorithms



- For data points that are only “inside” the split points
- Enables “Parallelization”

Split finding Algorithms

Algorithm 2: Approximate Algorithm for Split Finding

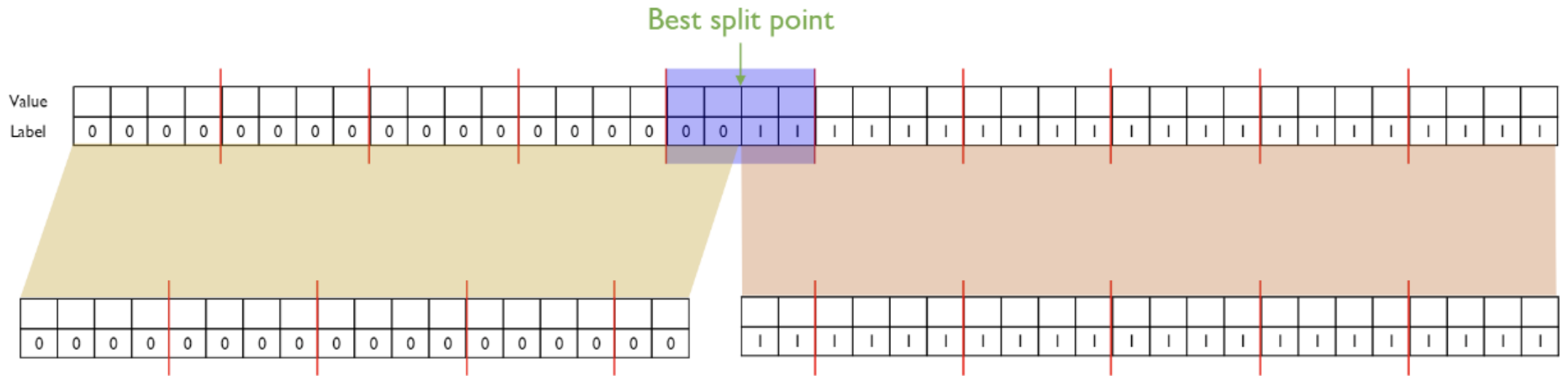
```
for  $k = 1$  to  $m$  do  
    | Propose  $S_k = \{s_{k1}, s_{k2}, \dots, s_{kl}\}$  by percentiles on feature  $k$ .  
    | Proposal can be done per tree (global), or per split(local).  
end  
for  $k = 1$  to  $m$  do  
    |  $G_{kv} \leftarrow \sum_{j \in \{j | s_{k,v} \geq \mathbf{x}_{jk} > s_{k,v-1}\}} g_j$   
    |  $H_{kv} \leftarrow \sum_{j \in \{j | s_{k,v} \geq \mathbf{x}_{jk} > s_{k,v-1}\}} h_j$   
end
```

Follow same step as in previous section to find max score only among proposed splits.

Split finding Algorithms

- Global variant (per tree)

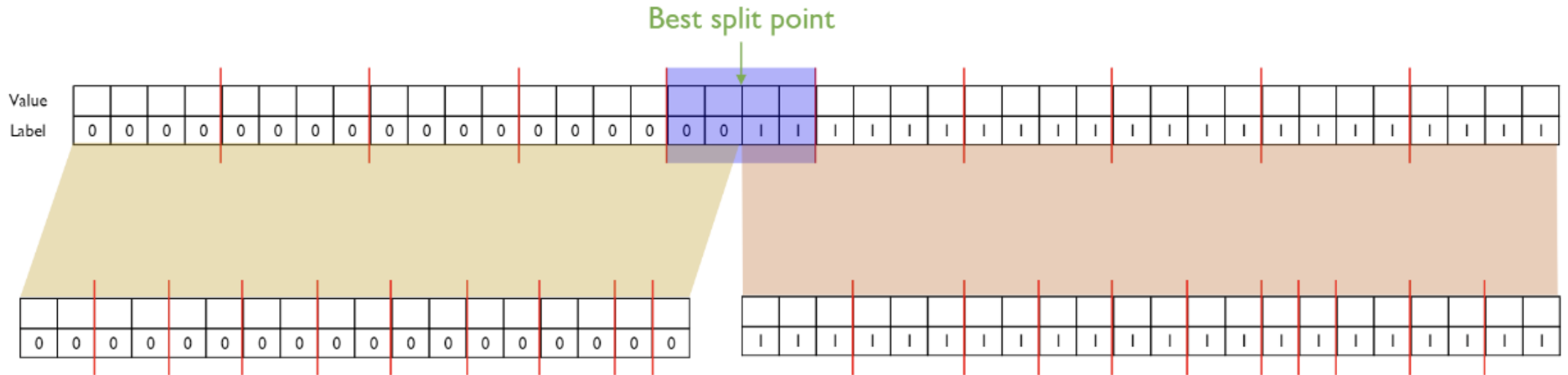
- Global variant vs. Local variant



Split finding Algorithms

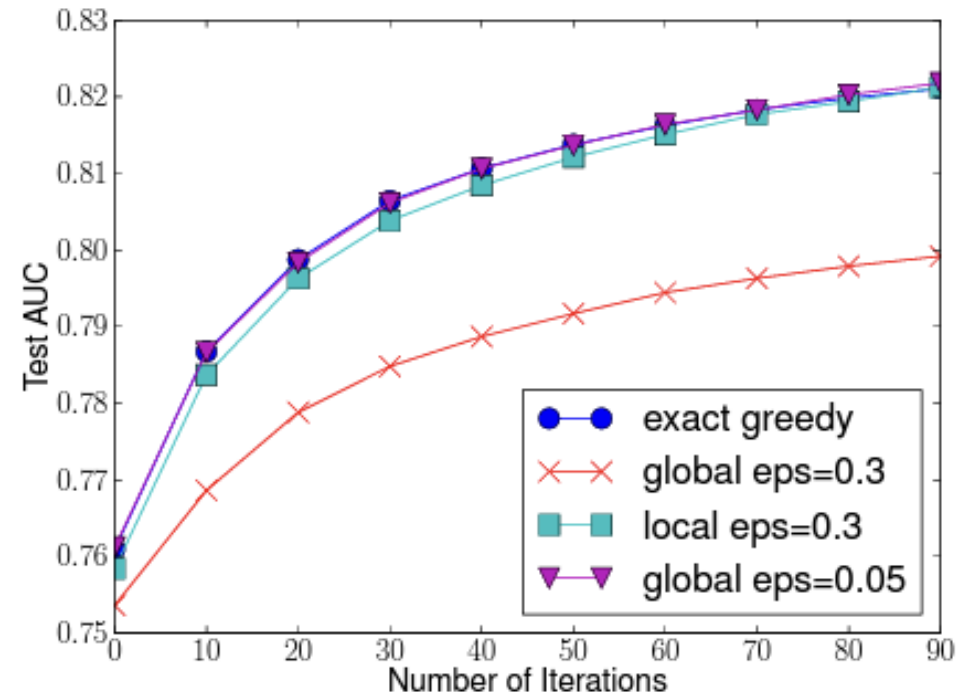
- Local variant (per split)
- Maintains the # buckets

- Global variant vs. Local variant



Split finding Algorithms

- Does not mention which (global or local) algorithm is better
- Suggest the way to choose an appropriate epsilon
- Global: large epsilon does not work
- Local: large epsilon does work
- Epsilon: percentile parameter
- $1/\text{epsilon} \sim \#$ of candidate split points



Split finding Algorithms

- Sketch Algorithm

Get a scheme of the Original Data Distribution w/ sample data sketch

- Quantile Sketch Algorithm

Get a scheme of the Original Data Distribution w/ sample data sketch & quantile

- Weighted Quantile Sketch Algorithm

Normal quantile: each quantile has the same # of data

Weighted quantile: each quantile has the same sum of weights (h_i)

Split finding Algorithms

- Sketch Algorithm, Quantile Sketch Algorithm, Weighted Quantile Sketch Algorithm
- These algorithms enable XGBoost to make a parallel computation

Sparsity-Aware Split Finding

- Missing values

① | Recap

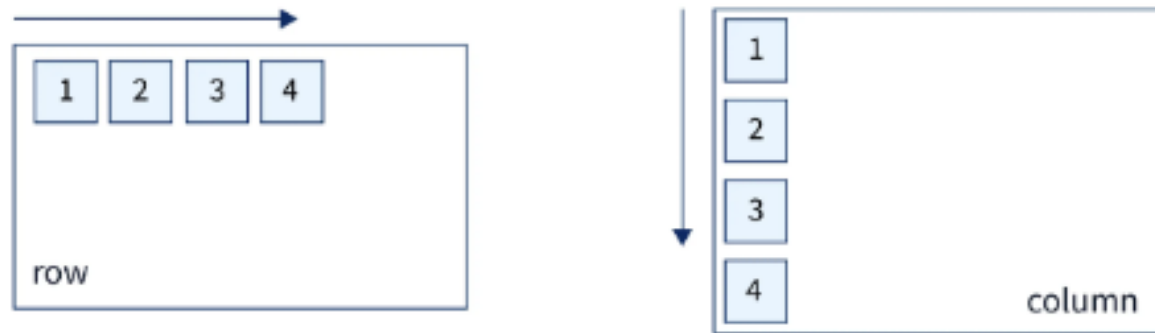
② | Split Finding Algorithms

③ | H/W Optimization

④ | Details

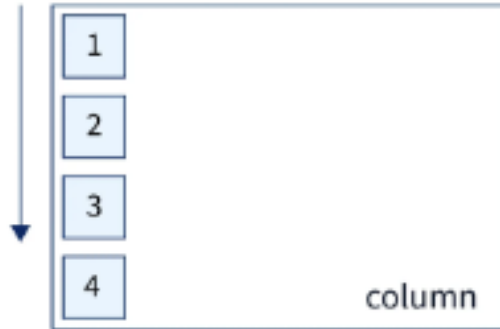
Hardware Optimization

- System Design for Efficient Computing
- Row Orientation vs. Column Orientation



- Row: Transactional Processing, ALL the columns are required
- Column: Only relevant columns are required

Hardware Optimization



- Data in each block is stored in the compressed column (CSC) format, with each column sorted by the corresponding feature value
- This input data layout only needs to be computed “only once” before training and can be reused in later iterations. (by blocks)
- BLOCKS?

Hardware Optimization

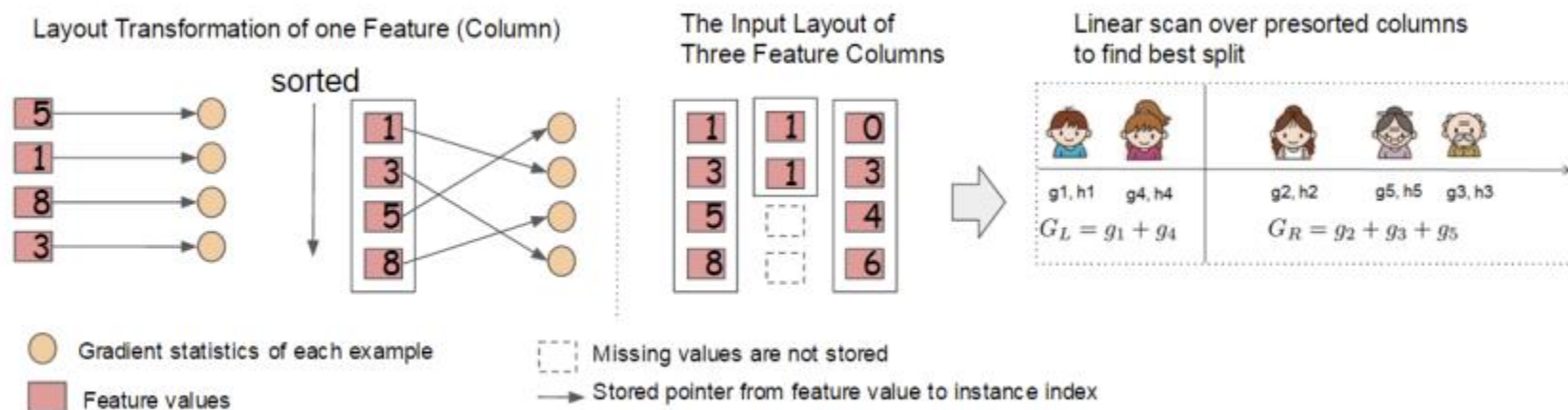
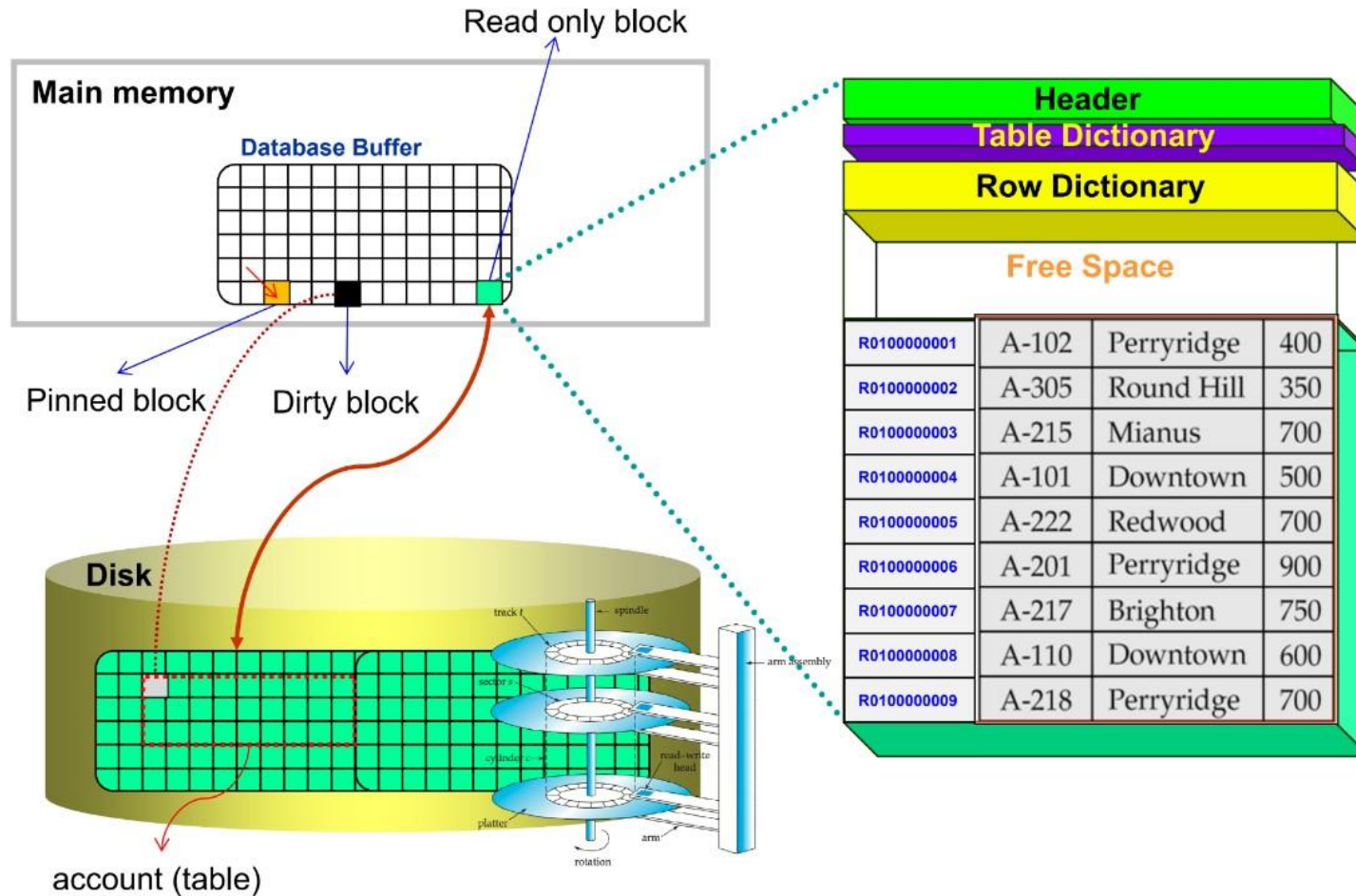


Figure 6: Block structure for parallel learning. Each column in a block is sorted by the corresponding feature value. A linear scan over one column in the block is sufficient to enumerate all the split points.

Hardware Optimization



Storage Access



Hardware Optimization

- Cache-aware access

Cache -> Main Memory (M/M) -> Disk (SSD, HDD)

I/O Speed: Cache > M/M > Disk

- Block: a virtual unit of data
- It is important to choose the block size adequately
- Bigger the better -> No! (slow access)
- Smaller the better -> No! (slow computation in total)

Hardware Optimization

- Out-of-core computing

Utilize disk space to handle data that does not fit into M/M (block)

Reduce overhead and increase the throughput of disk I/O

- Block Compression (CSC)
- Block Sharding (partition)

When we can utilize more than 1 disks, save the data in each disk in different orders

Can increase the reading (throughput) of disk I/O

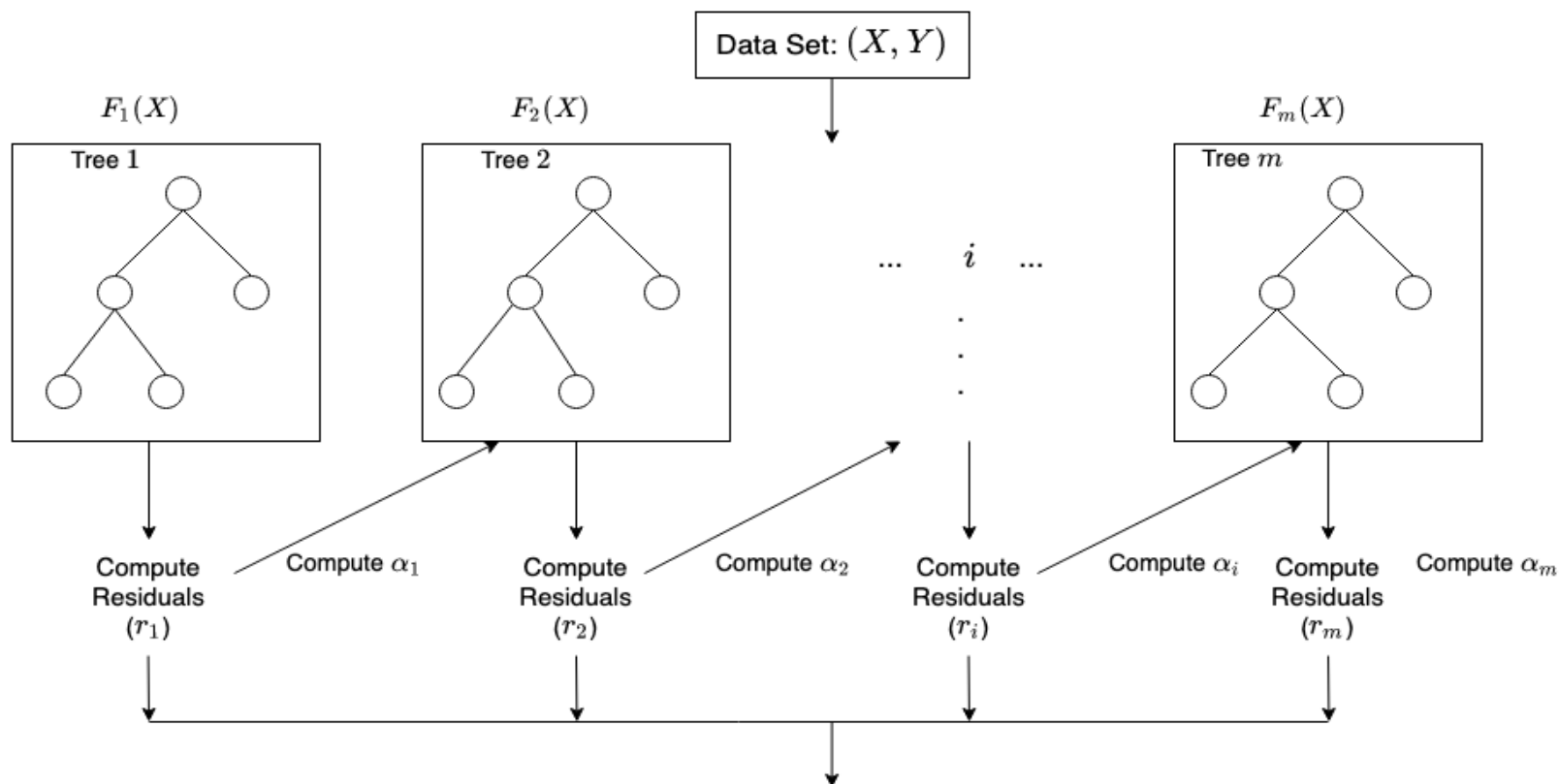
① | Recap

② | Split Finding Algorithms

③ | H/W Optimization

④ | Details

Details



$$F_m(X) = F_{m-1}(X) + \alpha_m h_m(X, r_{m-1}),$$

where α_i , and r_i are the regularization parameters and residuals computed with the i^{th} tree respectively, and h_i is a function that is trained to predict residuals, r_i using X for the i^{th} tree. To compute α_i we use the residuals

computed, r_i and compute the following: $\arg \min_{\alpha} = \sum_{i=1}^m L(Y_i, F_{i-1}(X_i) + \alpha h_i(X_i, r_{i-1}))$ where

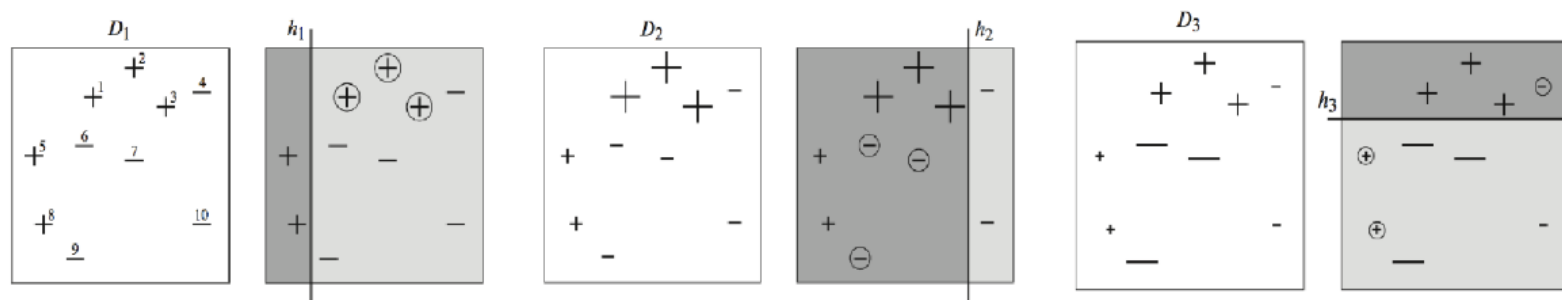
$L(Y, F(X))$ is a differentiable loss function.

Details

- Models based on GBM: XGBoost, LightGBM (Microsoft), CatBoost (Categorical dataset)
- AdaBoost – w/o stacked trees, make a complex hyperplane

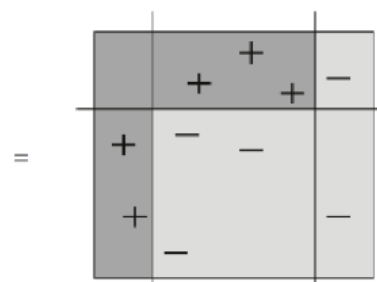
Details

- Adaboost



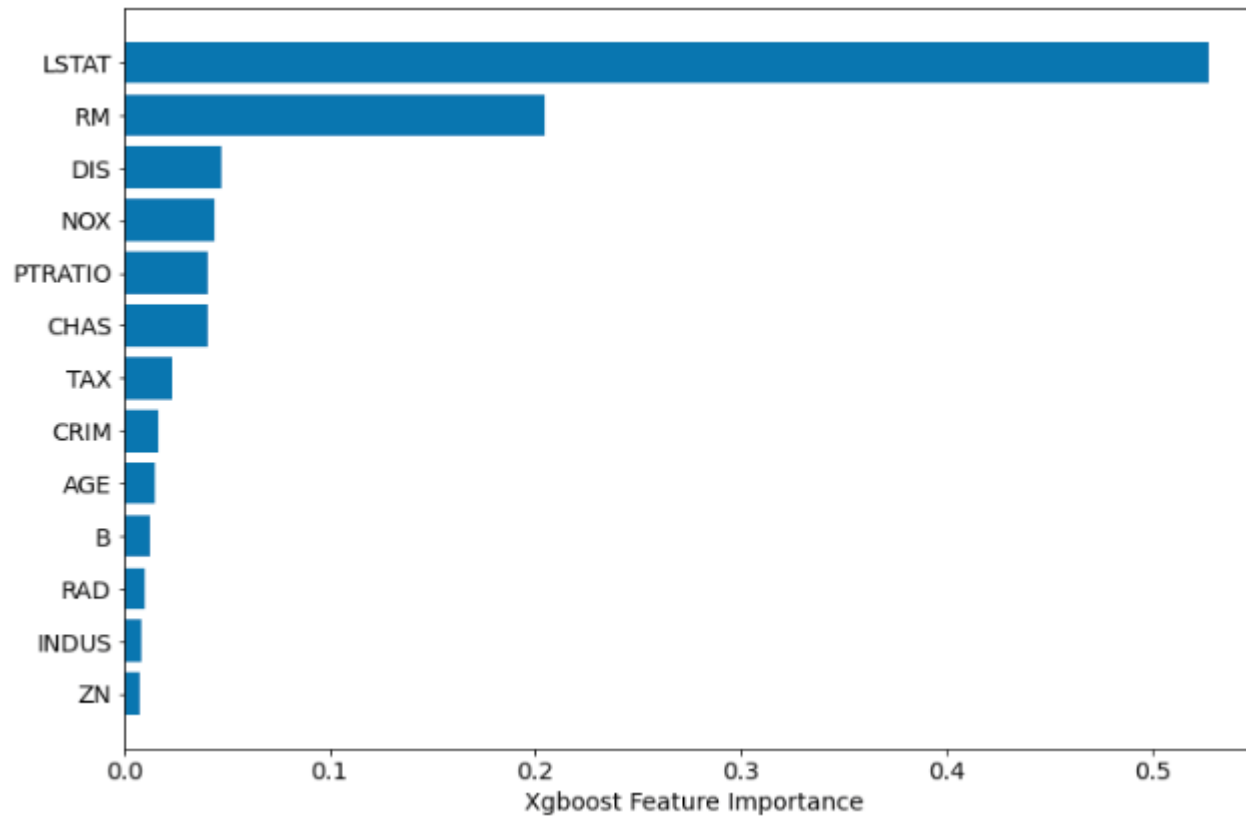
$$H = \text{sign} \left(0.42 \begin{array}{|c|} \hline \text{shaded region} \\ \hline \end{array} + 0.65 \begin{array}{|c|} \hline \text{shaded region} \\ \hline \end{array} + 0.92 \begin{array}{|c|} \hline \text{shaded region} \\ \hline \end{array} \right)$$

$$H(x) = \sum_t \rho_t h_t(x)$$



Details

- Plot_importance



Details

- Plot_importance

- Information gain (default):
$$L_{split} = \frac{1}{2} \left[\frac{(\sum_{i \in I_L} g_i)^2}{\sum_{i \in I_L} h_i + \lambda} + \frac{(\sum_{i \in I_R} g_i)^2}{\sum_{i \in I_R} h_i + \lambda} - \frac{(\sum_{i \in I} g_i)^2}{\sum_{i \in I} h_i + \lambda} \right] - \gamma$$

(Loss function before split) - (Loss function after split)

- get_score: # of the advent of variables (F score)
- Num_rounds: # of boosting
- Max_depth: depth of a tree
- Subsample [0, 1], ETA (learning rate), gamma (regularization, avoid overfitting)

Reference

- <https://github.com/pilsung-kang/Business-Analytics-IME654-/tree/master/04%20Ensemble%20Learning>
- XGBoost: A Scalable Tree Boosting System

Future Work

- Deep Learning for Tabular Data
- Machine Learning based models (in general), usually outperforms the deep learning model 'only' in terms of tabular data analysis
- But what if we need Deep Learning in tabular data analysis?

Happens when multi-modal data analysis is held

- Then let's make Deep Tabular Learning model that utilizes the pros of Machine Learning techniques

Future Work

- But just simple Deep Learning? (DL usually requires a bunch of time to optimize)
- Then Meta Learning (Few-shot, fast meta-learner)
- When Meta Learning is combined w/ Tree-boosted ML techniques...?
- How to make ML techniques differentiable? ... Gumbel-(soft)max techniques?