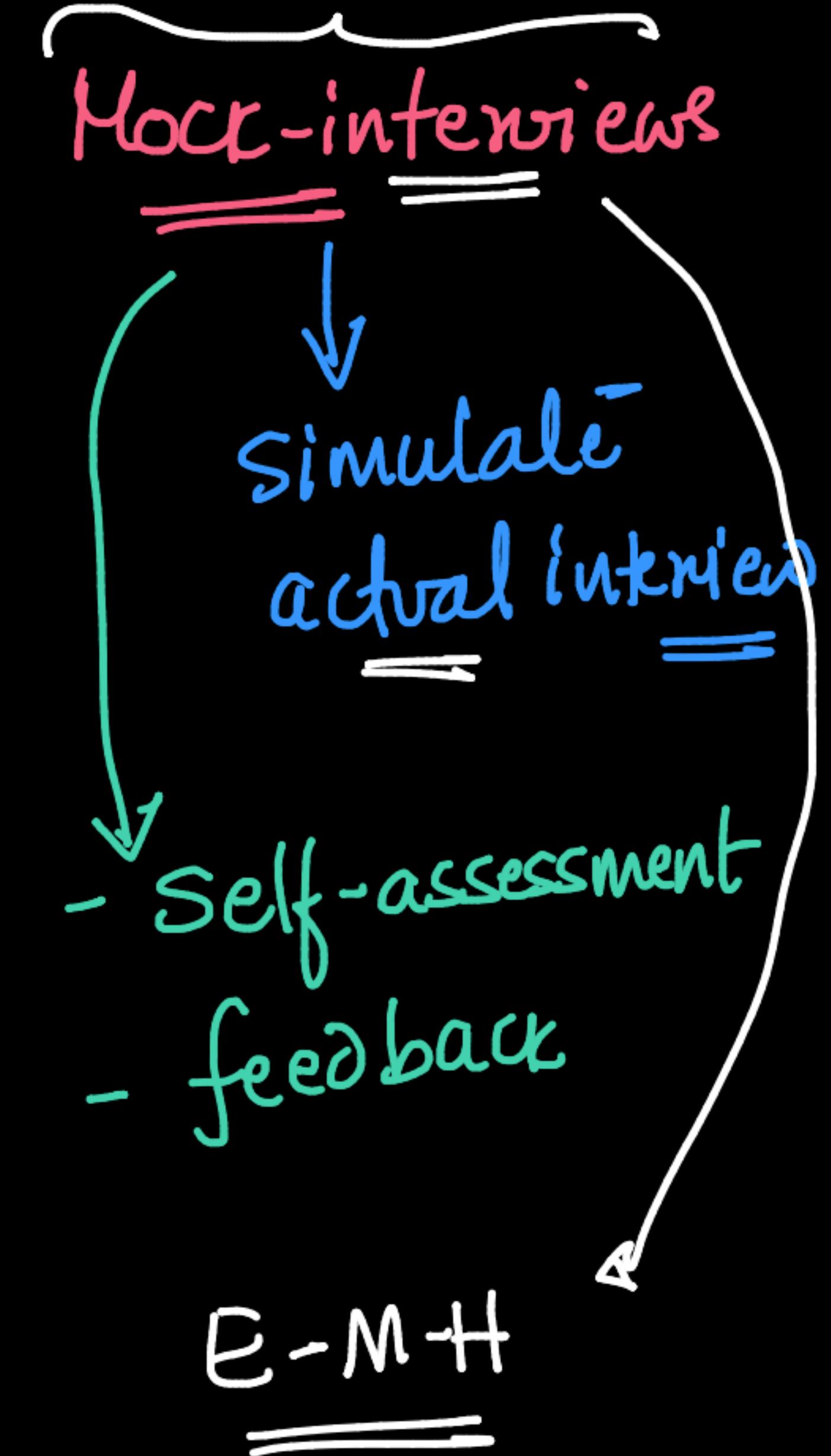
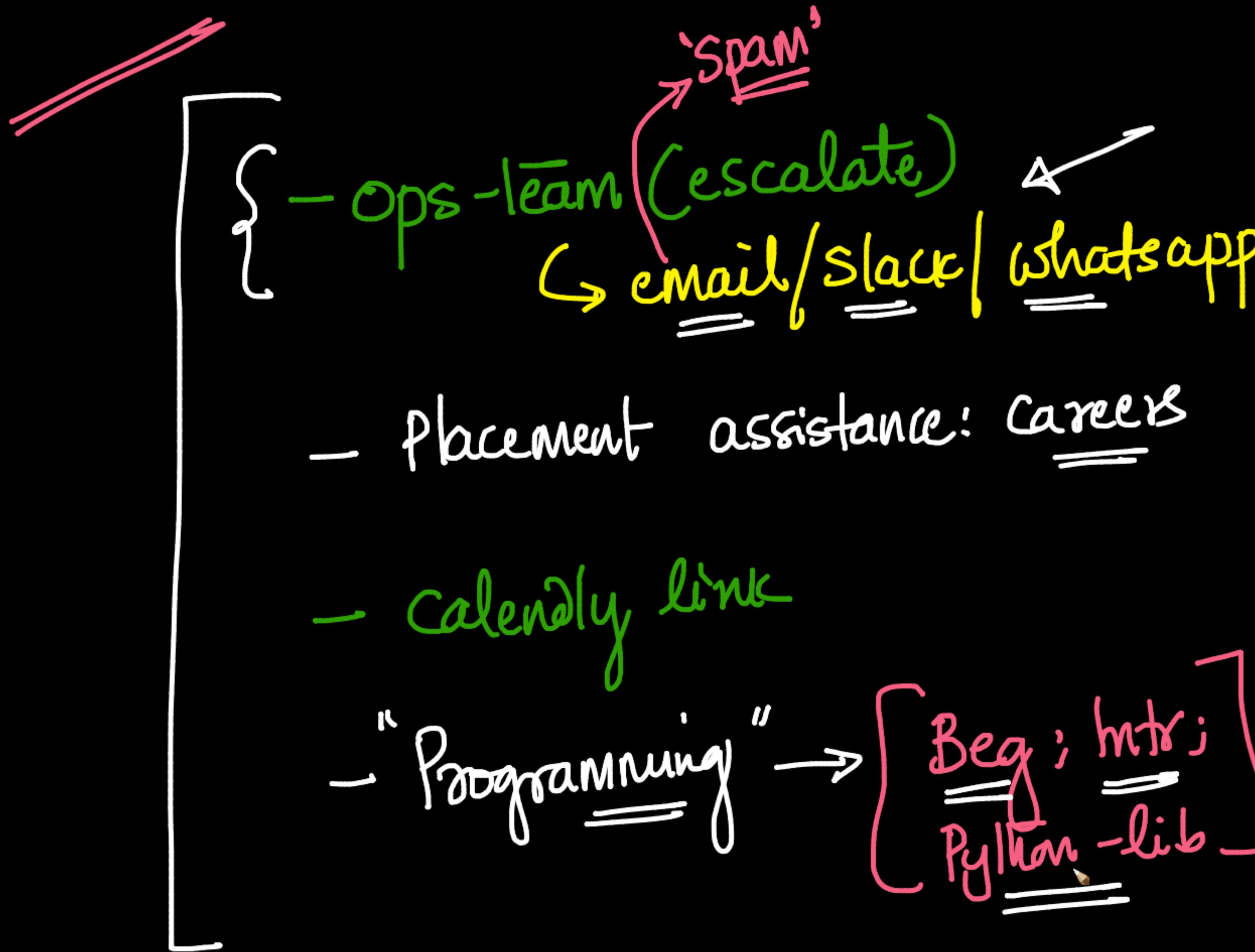


# Topics:

[revised]

- ✓ - Mock-interviews
- {
  - GBDT
    - Training algo ✓
    - Bias-variance tradeoff ✓
    - Variations ✓
  - Intuition behind AdaBoost
- ✓ {
  - XgBoost ✓
  - LightGBM ✓



optional  
Module:

Adv DSA for MLE  $\rightarrow$  Greedy; DP; Graphs



pseudo-residual

{already trained  
m-1 classifiers}



$f_{m-1}(x_i)$

$i : l \rightarrow n$  training points

$\overbrace{1 \rightarrow N}$

m<sup>th</sup> classifier

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

Loss: regression  $\Rightarrow$  squared-loss

$$L(y_i, f(x_i)) = \{y_i - f(x_i)\}^2$$

$$\frac{\partial L}{\partial f(x_i)} = -2(y_i - f(x_i)) \quad (\text{P}_{\text{class}})$$

For point  $i$ :

$$\gamma_{im} = \frac{1}{2} [y_i - f_m(\underline{x}_i)]$$

$f_{m-1}(\underline{x}_i)$

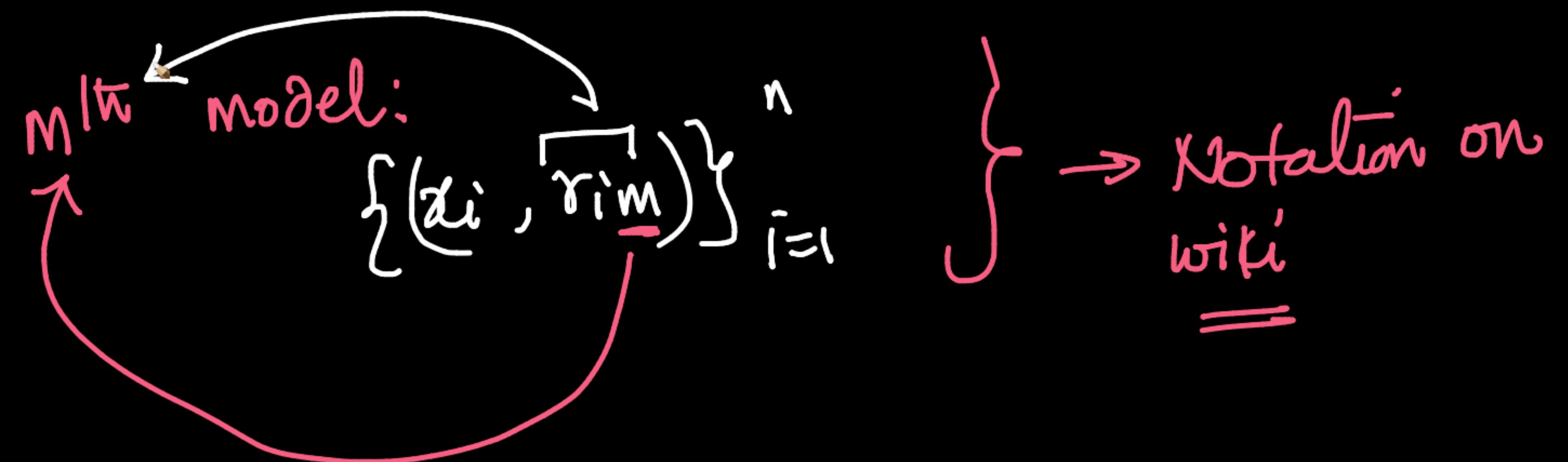
$\gamma_i: l \rightarrow n$

Already have:  $D = \{(x_i, y_i)\}_{i=1}^n$

$\hat{f}_{m-1}(\underline{x}_i)$

equivalent  
to

$$\frac{1}{2} (y_i - \hat{f}_{m-1}(\underline{x}_i))^2$$



$$\mathbf{x}_i \in \mathbb{R}^d$$

$$\text{scalar} \leftarrow y_i \in \mathbb{R}^1$$

Idea:

$$\text{pseudo-residual} = \frac{-\frac{\partial L}{\partial f(x_i)}}{\frac{\partial^2 L}{\partial f(x_i)^2}} \approx \text{residual}$$

any loss function  
that is  
differentiable

$$F_{M-1}(x_i) = \underline{\alpha_0 h_0(x_i) + \alpha_1 h_1(x_i) + \dots + \alpha_{M-1} h_{M-1}(x_i)}$$

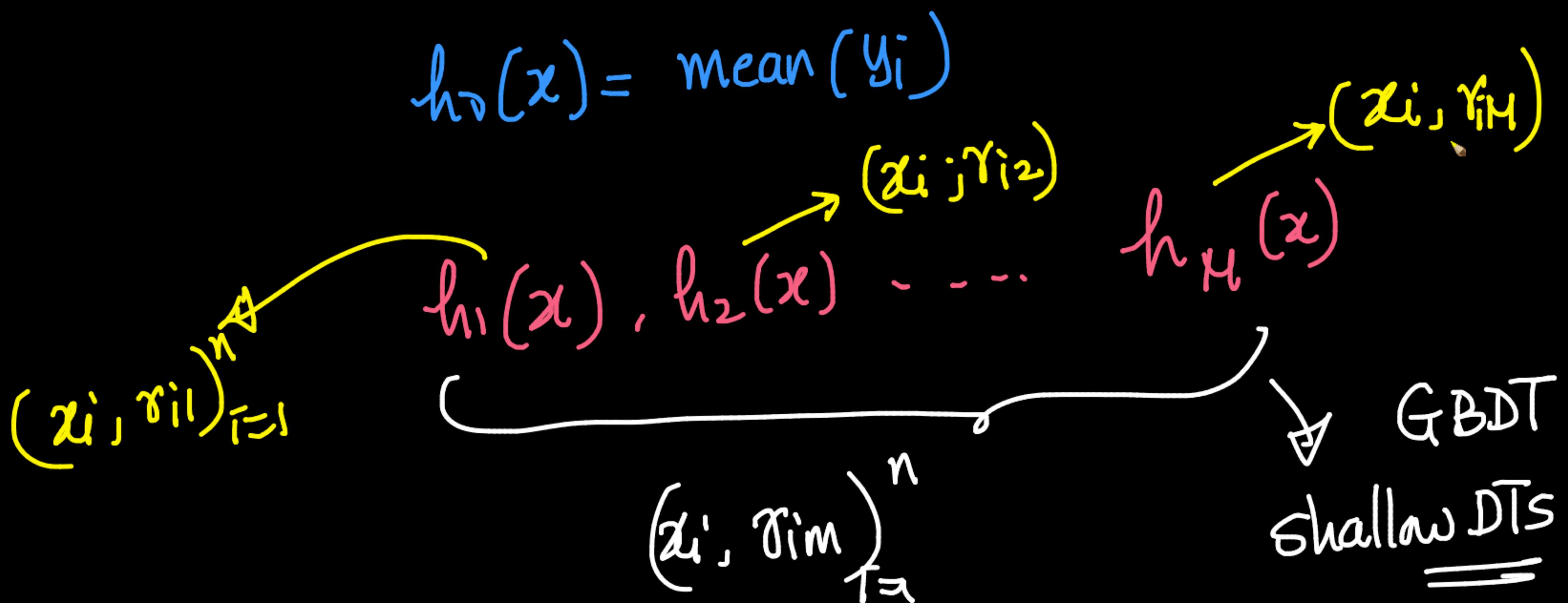
$M-1^{\text{th}}$  iteration  
 $(M-1)^{\text{th}}$  base-model

$$f_m(x_i) = \frac{\gamma_0 h_0(x) + \gamma_1 h_1(x) + \gamma_2 h_2(x) + \dots}{\gamma_{m-1} h_{m-1}(x) + \gamma_m h_m(x)}$$

$m^{\text{th}}$ -iteration

$$\hookrightarrow \left( x_i, \gamma_i m \right)_{i=1}^n \rightarrow \underline{h_m(x)} \quad \checkmark$$

also compute  $\underline{\gamma_m}$  (later)  
1



pseudo residuals are used  
to train

M-base learners

$$f_M(x) = \sum_{i=0}^M \gamma_i \underline{h_i}(x) = \underline{y}$$

$\hookrightarrow$  shallow DT

Test-time:

$$\underline{x}_q \rightarrow \sum_{i=0}^M \gamma_i \hat{h}_i(\underline{x}_q) = \hat{y}_q \quad \checkmark$$

let

$$f_M(x) = \underbrace{\gamma_0 h_0(x)}_{12} + \underbrace{\gamma_1 h_1(x)}_{1} + \underbrace{\gamma_2 h_2(x)}_{1} + \dots + \underbrace{\gamma_M h_M(x)}_{1}$$

$$\boxed{h_0(x_i) = \underline{12 \cdot 0}} \quad \cancel{14 \cdot 2}$$

$$h_1(x_i) = \underline{4 \cdot 2} \quad 10$$

$$h_2(x_i) = \dots = 0$$

$$h_M(x_i) = \dots$$

hyper-param

$$i=10 \text{ (let)} \\ x_i = [1, 2, 1, 0] \\ y_i = \boxed{26 \cdot 2} \quad \checkmark$$

residuals

$$\underline{y_i} \approx \hat{y}_i$$

$\mu \uparrow$  overfit  $\uparrow$

# Overfitting vs Underfitting

Let  $\{$

- $h_0(x) = \text{extremely simple} = \text{mean-model}$
- $h_1(x): \underline{\text{Deep DT}} \Rightarrow \text{overfitting}$
- $\vdots$
- $h_5(x): "$

$\}$

✓  $h_0(x)$ : Mean-Model  
↑ train data only  
No CV or test data

# classification

$$(x_i, y_i)_{i=1}^n$$

$$y_i = 0 \text{ or } 1$$

$$p_i = P(y_i=1)$$

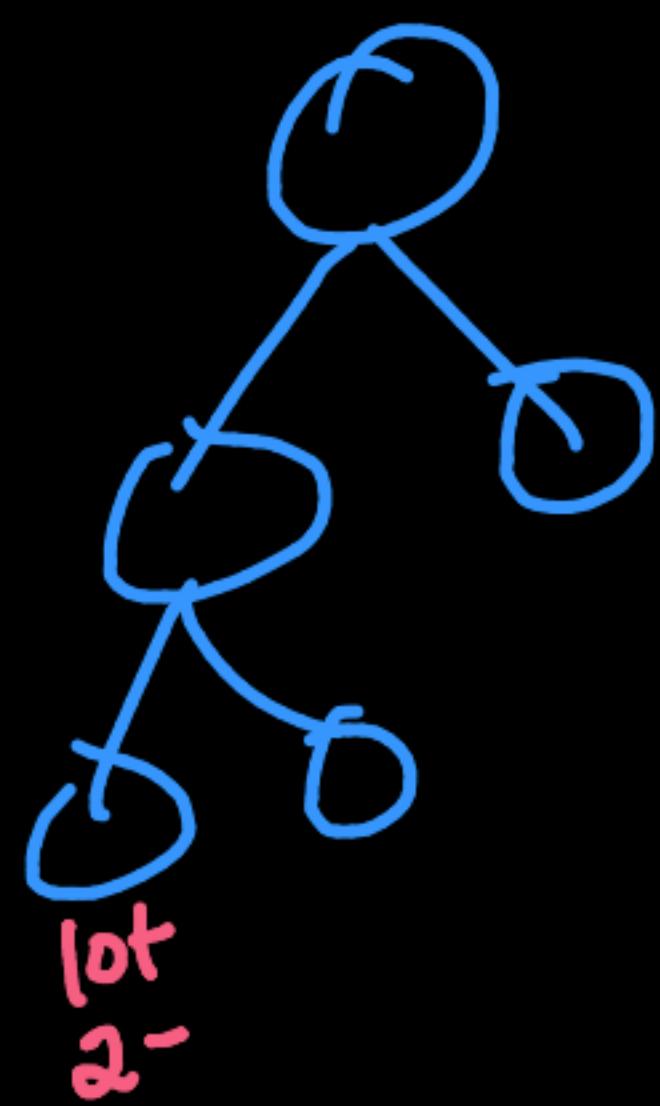
$$f_M(x_i) \rightarrow p_i$$

If  $x_i \in \text{+ve}$   $p_i \rightarrow 1$   
else  $p_i \rightarrow 0$

Loss-function for classfn: log-loss

$$\mathcal{L}(y_i, f(x_i)) = \text{log-loss}$$

$$-\frac{\partial \mathcal{L}}{\partial F(x)} = p_i - y_i \quad (\text{proved in post-read video})$$



$$h_0(x_i) = \frac{0.4}{0.6}$$

$$h_1(x_i) = \frac{0.2}{0.4}$$

$$h_2(x_i) = 0.31$$

⋮

$$x_i = [0, 1, 3, 1]$$

$$y_i = [1]$$

$$h_0(x) = 0.4$$

$$-0.4$$

$$h_1(x) = -0.1$$

$$-0.3$$

$$h_2(x) = -0.21$$

$$-0.09$$

$$h_3(x) =$$

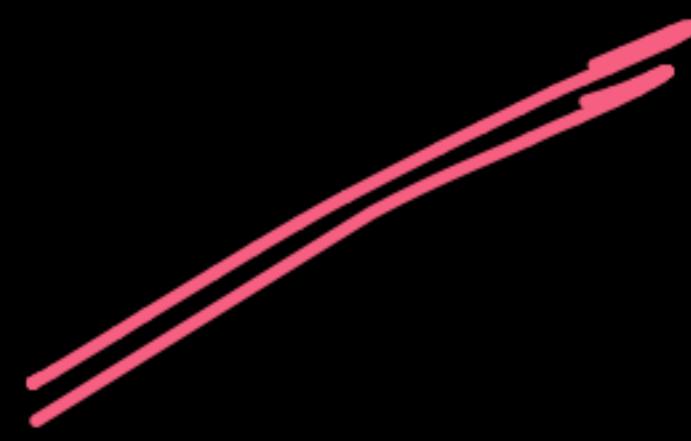
1

2

3

$$x_i^- = [1, 2, 0, 1]$$

$$y_i^- = 0$$



Let  $h_0(x)$  Deep DT err: small

$$\sum_{i=0}^M \gamma_i h_i(x)$$

(prob)

GBDT  $\rightarrow$  Stochastic GBDT (later)

↳ shallow trees

row & col-sampling

Fancy-loss:

$$\frac{-\partial L}{\partial F(x)} \quad \checkmark$$

hinge-loss  
(later SVM)

Sq-loss:

$$L_2 \rightarrow 2(y_i - \hat{y}_i)$$

pseudo residual

↓ ≈  
residual

log-loss

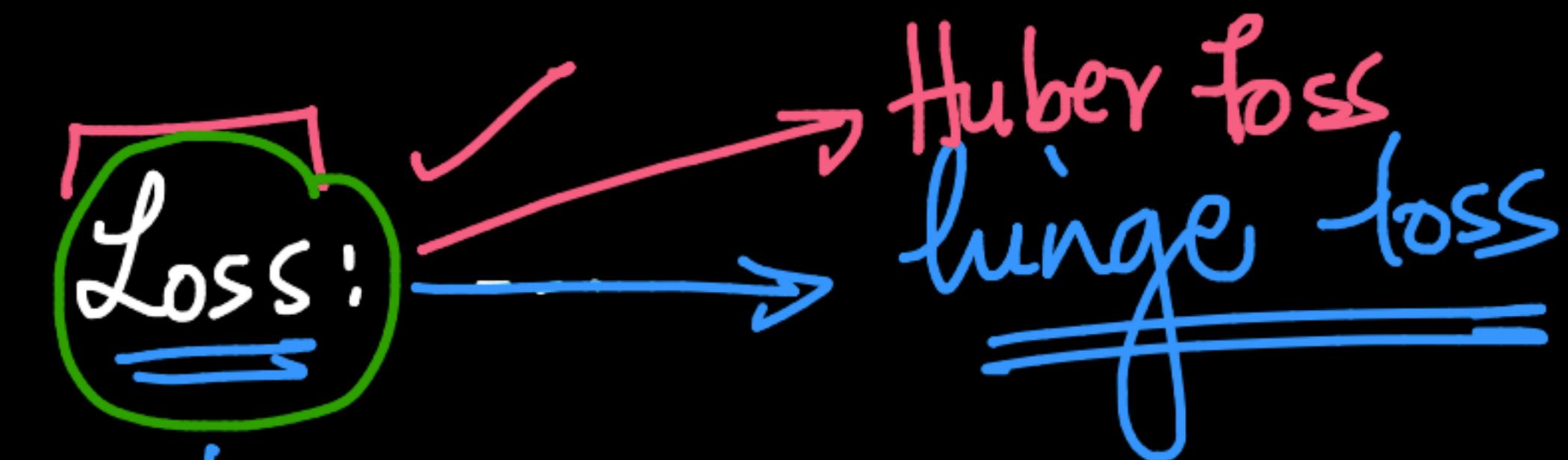
Loss

$$\pi_i - y_i$$

(pseudo residual)

Real-world:

Data j



A blue arrow points downwards from the "Loss:" oval to a mathematical equation. The equation is  $\frac{-\partial L}{\partial f(x_i)} = \text{residual}$ . The word "residual" is written in pink at the end of the arrow. The entire equation is underlined.

$$\frac{-\partial L}{\partial f(x_i)} = \text{residual}$$



Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ . ③

Algorithm: ①

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

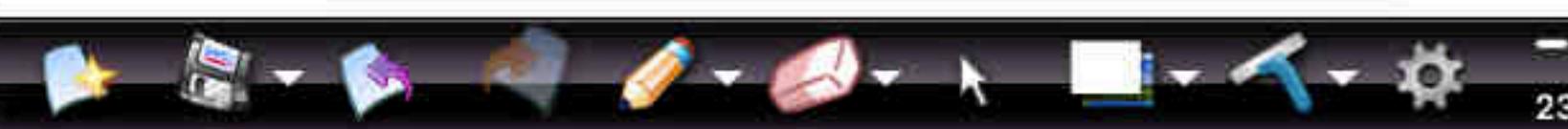
3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output  $F_M(x)$ .



Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

$h_0(x)$  = mean

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

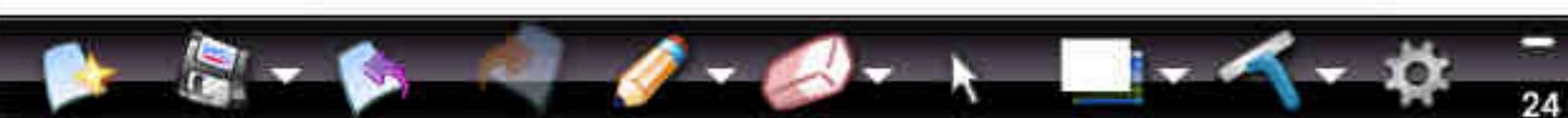
3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output  $F_M(x)$ .



Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

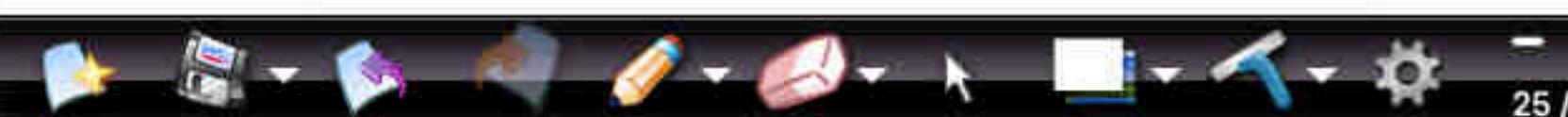
3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output  $F_M(x)$ .



Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

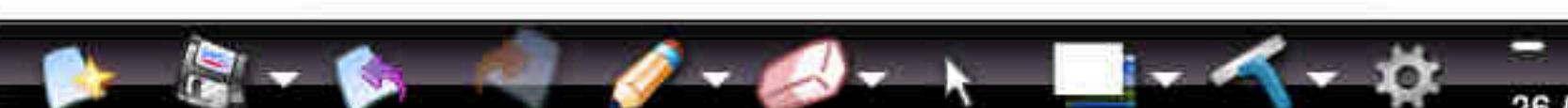
1. Initialize model with a constant value:  

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \underline{\gamma}).$$
mean for SQ-loss
→  $n/3$  +ve  
 $2n/3$  -ve
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:  

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$$
for  $i = 1, \dots, n$ .
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following *one-dimensional optimization problem*:  

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
  4. Update the model:  

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
  3. Output  $F_M(x)$ .



Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

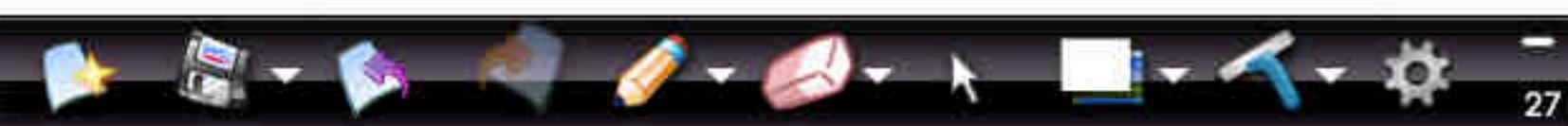
3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output  $F_M(x)$ .



$$f_0(x) = \arg \min_{\gamma} \sum_{i=1}^n \ell(y_i, \gamma)$$
$$= \arg \min_{\gamma} \left[ \sum_{i=1}^n (y_i - \gamma)^2 \right]$$

$g(\gamma)$

↓  
Scalar

$$\frac{\partial g(\gamma)}{\partial \gamma} = \sum_{i=1}^n -\cancel{2}(y_i - \gamma) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i = n \gamma$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n y_i = \gamma$$

ANSWER

Pseudo-residuals for Gradient boosting - Wikipedia AdaBoost - Wikipedia talk.dvi 1603.02754.pdf Python API Reference Random Forests(TM) LightGBM: A Highly Efficient C++ Library for Machine Learning lightgbm.LGBMClass

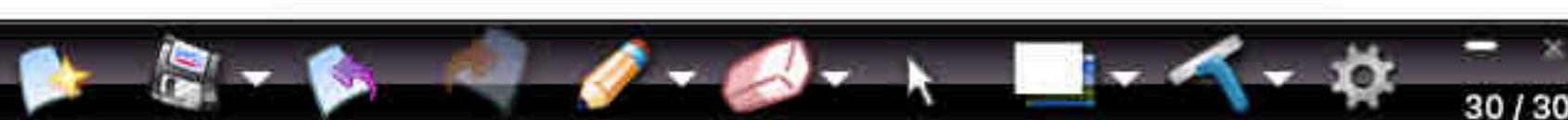
en.wikipedia.org/wiki/Gradient\_boosting#Disadvantages

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:  
$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:  
$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:  
$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
  4. Update the model:  
$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
  3. Output  $F_M(x)$ .

$$\frac{\partial L}{\partial \gamma} = 0$$



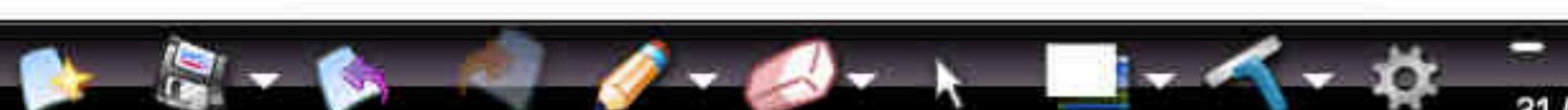
Pseudo-residuals for Gradient boosting - Wikipedia AdaBoost - Wikipedia talk.dvi 1603.02754.pdf Python API Reference Random Forests(TM) LightGBM: A Highly Efficient C++ Library for Machine Learning lightgbm.LGBMClass

en.wikipedia.org/wiki/Gradient\_boosting#Disadvantages

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:  
$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:  
$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following [one-dimensional optimization problem](#):  
$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
  4. Update the model:  
$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
3. Output  $F_M(x)$ .



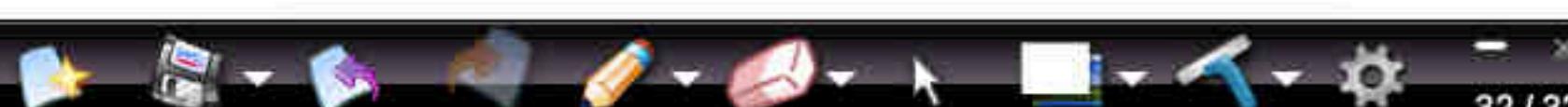
Pseudo-residuals for Gradient boosting - Wikipedia AdaBoost - Wikipedia talk.dvi 1603.02754.pdf Python API Reference Random Forests(TM) LightGBM: A Highly Efficient C++ Library for Machine Learning lightgbm.LGBMClass

en.wikipedia.org/wiki/Gradient\_boosting#Disadvantages

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:  
 $F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:  
 $r_{im} = -\left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$  for  $i = 1, \dots, n$ .
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:  
 $\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$
  4. Update the model:  
 $F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$
3. Output  $F_M(x)$ .



Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

$m: 1 \rightarrow M$

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output  $F_M(x)$ .

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output  $F_M(x)$ .

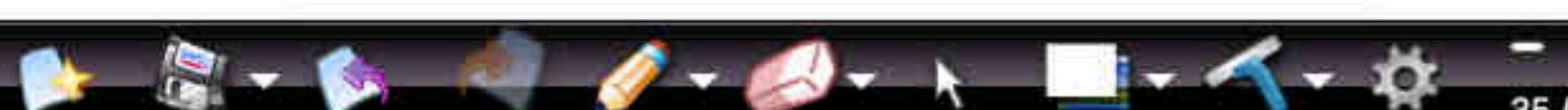
Pseudo-residuals for Gradient boosting - Wikipedia AdaBoost - Wikipedia talk.dvi 1603.02754.pdf Python API Reference Random Forests(TM) LightGBM: A Highly Efficient C++ Library for Machine Learning lightgbm.LGBMClass

en.wikipedia.org/wiki/Gradient\_boosting#Disadvantages

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:
$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:
$$\check{r}_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, \check{r}_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:
$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
  4. Update the model:
$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
  3. Output  $F_M(x)$ .



Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:  
$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:  
$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:  
$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
  4. Update the model:  
$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
  3. Output  $F_M(x)$ .

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:  

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:  

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:  

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
  4. Update the model:  

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
  3. Output  $F_M(x)$ .

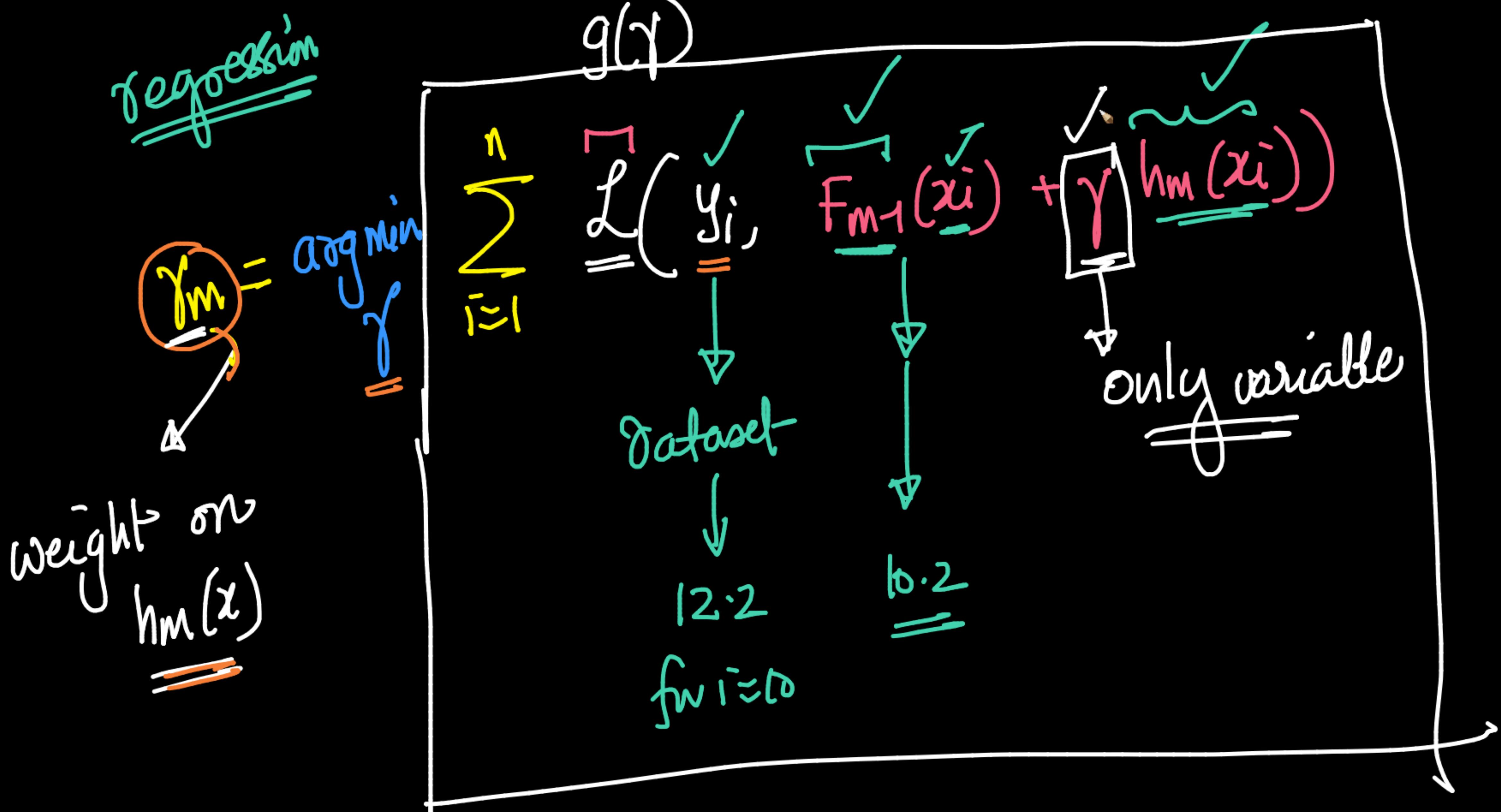
*$f_M(x) = f_0(x) + \sum_{i=1}^M \gamma_i h_i(x)$*

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

Algorithm:

1. Initialize model with a constant value:  
$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:  
$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following **one-dimensional optimization problem**:  
$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
  4. Update the model:  
$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
3. Output  $F_M(x)$ .

$$F_m = \underbrace{h_0(x)}_{\gamma_0} + \gamma_1 h_1(x) + \dots + \gamma_{M-1} h_{M-1}(x) + \boxed{\gamma_m} h_m(x)$$



en.wikipedia.org/wiki/Gradient\_boosting#Disadvantages

1. Initialize model with a constant value:
$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:
$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:
$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
4. Update the model:
$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
3. Output  $F_M(x)$ .

*scalar*

## Gradient tree boosting [edit]

Gradient boosting is typically used with decision trees (especially CART trees) of a fixed size as base learners. For this special case, Friedman

en.wikipedia.org/wiki/Gradient\_boosting#Disadvantages

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:  

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output  $F_M(x)$ .

## Gradient tree boosting [ edit ]

Gradient boosting is typically used with decision trees (especially CART trees) of a fixed size as base learners. For this special case, Friedman proposed a modification which improves the quality of fit of each base

en.wikipedia.org/wiki/Gradient\_boosting#Disadvantages

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

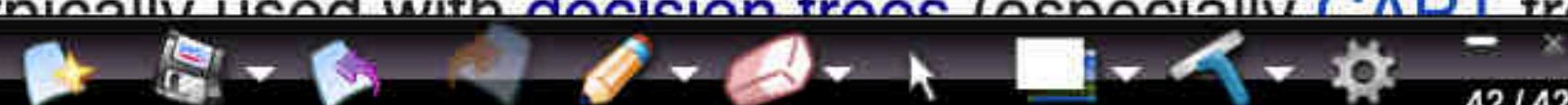
4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output  $F_M(x)$ .

## Gradient tree boosting [ edit ]

Gradient boosting is typically used with decision trees (especially CART trees) of a fixed size as base learners. For this special case, Friedman



en.wikipedia.org/wiki/Gradient\_boosting#Disadvantages

Algorithm:

1. Initialize model with a constant value:  
$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:  
$$r_{im} = -\left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:  
$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
4. Update the model:  
$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
3. Output  $F_M(x)$ .

*LOSS!*

*g(..)*

*A(γ)*

## Gradient tree boosting [edit]

Gradient boosting is a machine learning technique that builds a model by iteratively adding weak learners (base learners) of a fixed size as base learners. For this

GBDT

Algorithm:

1. Initialize model with a constant value:  

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:  

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:  

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
  4. Update the model:  

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
  3. Output  $F_M(x)$ .

Loss:  $\rightarrow$  Yes  
 $\rightarrow$  hack in some cases  
 $\uparrow$  SUM

## Gradient tree boosting [edit]

Gradient boosting is a machine learning technique for regression and classification problems, which works by combining multiple weak learners in a boosted ensemble.<sup>1</sup> It is similar to many other boosting methods, such as AdaBoost and gradient descent boosting, but typically uses regression trees as the individual learners. For this reason, it can be called "gradient tree boosting".

**Algorithm:**

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

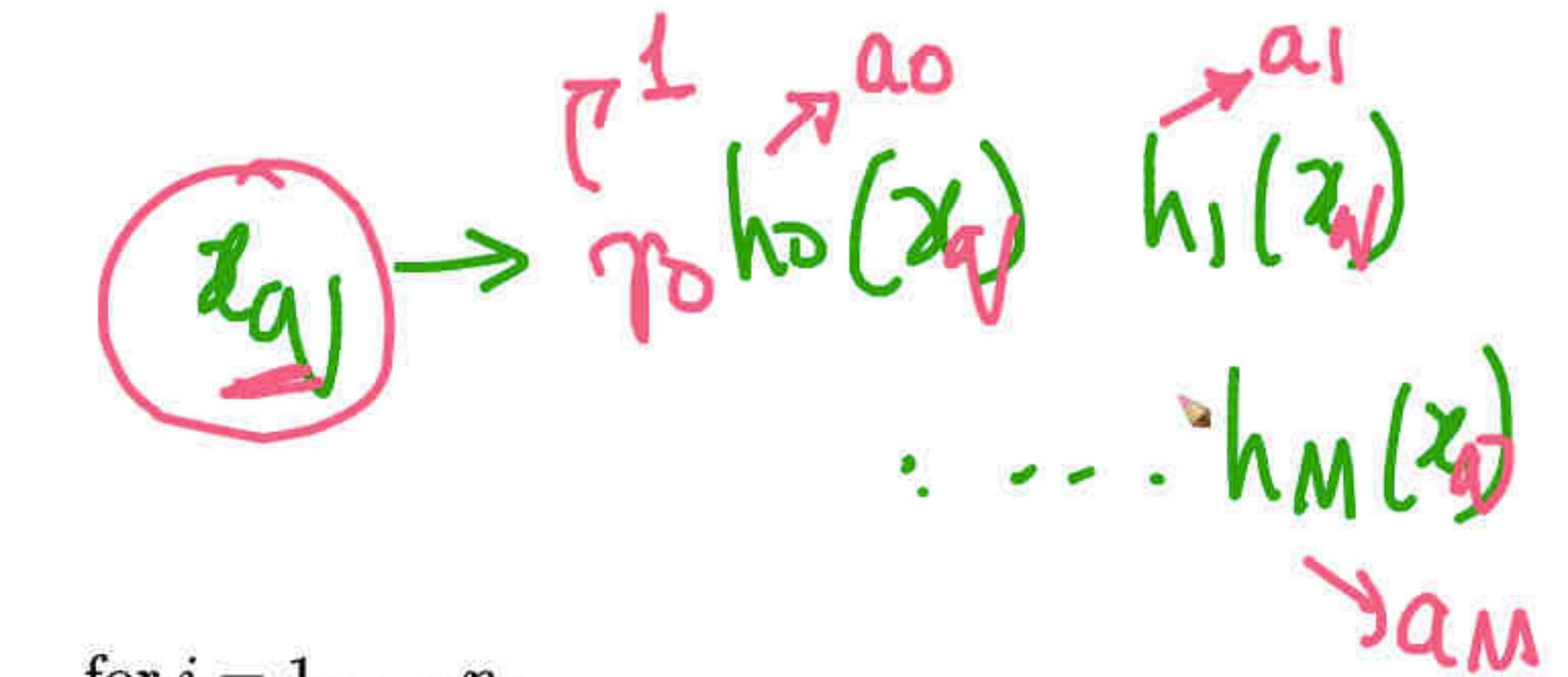
3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + \underbrace{\gamma_m}_{\text{GDT}} \underbrace{h_m(x)}_{y_m}.$$

3. Output  $F_M(x)$ .



$$y_q = \gamma_0 a_0 + \gamma_1 a_1 + \dots + \gamma_M a_M$$

**Gradient tree boosting** [ edit ]

Gradient boosting is



es) of a fixed size as base learners. For this

Algorithm:

1. Initialize model with a constant value:  

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$
2. For  $m = 1$  to  $M$ :
  1. Compute so-called *pseudo-residuals*:  

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$
  2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .
  3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:  

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$
  4. Update the model:  

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
  3. Output  $F_M(x)$ .

**GBDT** → all rows & cols  
→ all the features

$h_1(x) \ h_2(x) \dots \ h$

A handwritten diagram of a decision tree structure is shown at the bottom right, consisting of a root node branching into two child nodes, which further branch into leaf nodes.

## Gradient tree boosting [edit]

Gradient boosting is a machine learning technique for regression and classification problems, which works by combining multiple weak learners in a sequential manner. It is similar to AdaBoost, but uses a different loss function and a different way of combining the learners. The basic idea is to fit a series of simple models (base learners) sequentially, where each new model tries to correct the errors made by the previous ones. The final prediction is the weighted sum of all the base learners' predictions.

# Overfitting - Underfitting



Regularization → avoid overfitting

①

$M = \# \text{ base learners}$

$M \uparrow \text{overfit} \uparrow$

②

depth of base-learners  $\uparrow$

overfit  $\uparrow$

$\hookrightarrow \frac{n_{\text{samples}}}{n_{\text{leaf nodes}}}$

Risk of boosting → overfit very easily  
if not careful

③ Reg/  
Sh<sup>ortage</sup>/learning-rate  
additional  
leveY

~~Shrinkage:~~  $\gamma$  : hyperparam

$$0 \leq \gamma \leq 1$$

$$f_m(x) = f_{m-1}(x) + \underbrace{\gamma \cdot g_m h_m(x)}_{\text{derivatives}} \rightarrow \text{pseudo gradients}$$

$\gamma \uparrow$   
overfitting  $\uparrow$

$$\gamma=0$$

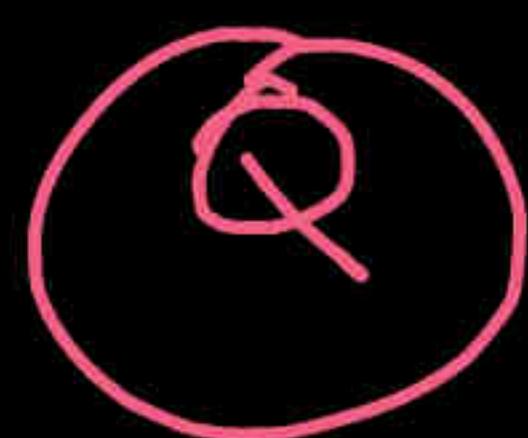
$h_0(x) \Rightarrow \text{Underfit Model}$

$$\gamma=1$$

overfitting increases

GBDT  $\rightarrow$   $\tilde{M}, \underbrace{\text{Max depth of tree} : \underline{\underline{v}}}_{\equiv}$

hyperparams



$\underline{\gamma}$



same performance

$\gamma = [\gamma_1 \ \gamma_2 \ \dots \ \gamma_M]$   $L_2\text{-reg}$  on  $\gamma_M^S$   $\rightarrow$  Yes

reg:  $L(y_i, f_M(x_i)) + \alpha \|\gamma\| + \beta (\# \text{leaf nodes})$

$M^{th\text{-model}} \rightarrow \text{control/regularize } \gamma_M$

**Algorithm:**

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

2. For  $m = 1$  to  $M$ :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

2. Fit a base learner (or weak learner, e.g. tree) closed under scaling  $h_m(x)$  to pseudo-residuals, i.e. train it using the training set  $\{(x_i, r_{im})\}_{i=1}^n$ .

3. Compute multiplier  $\gamma_m$  by solving the following one-dimensional optimization problem:

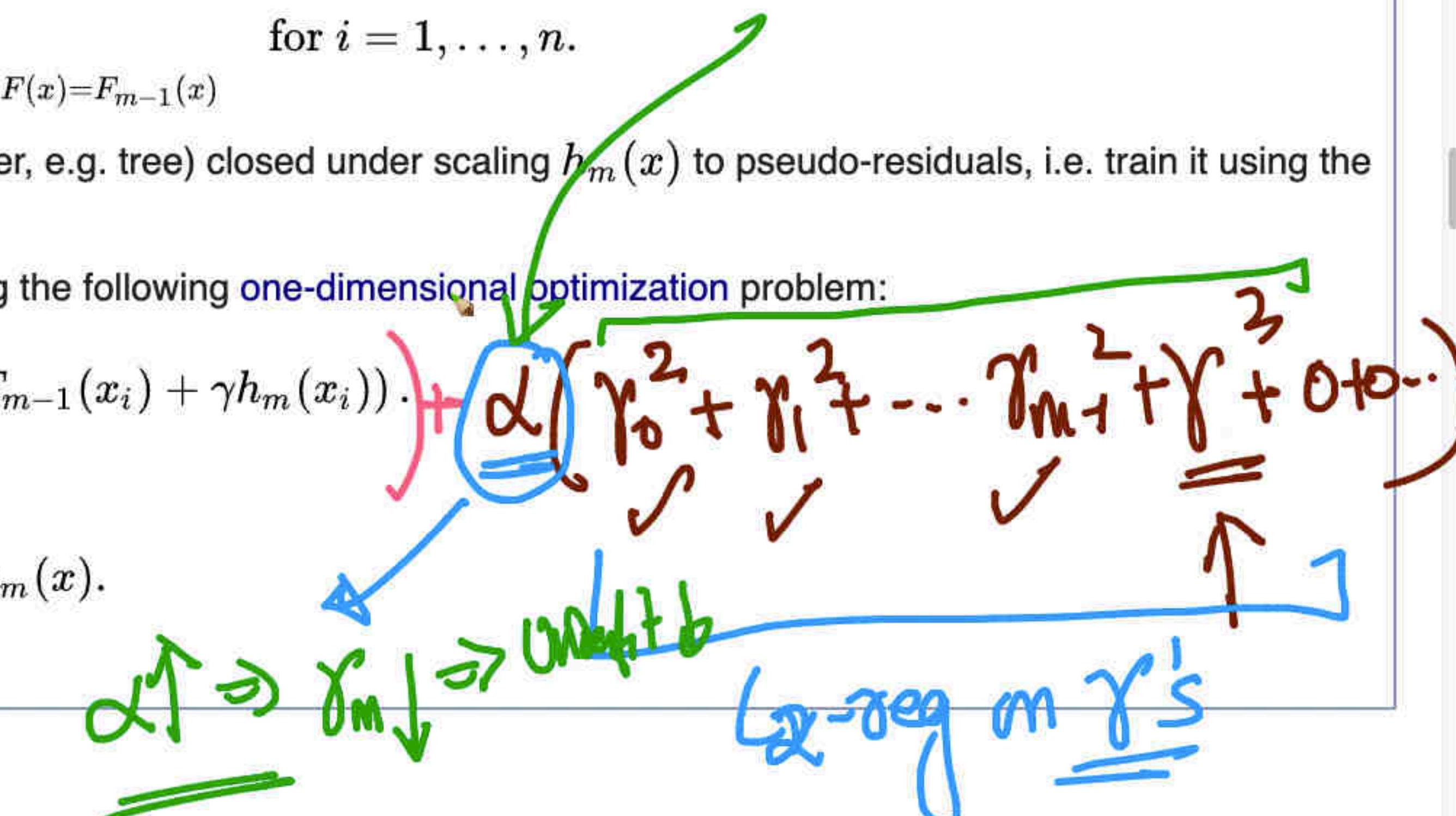
$$\boxed{\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).}$$

+  $\alpha$   $\left( \gamma_0^2 + \gamma_1^2 + \dots + \gamma_{m-1}^2 + \gamma^2 \right)$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output  $F_M(x)$ .

**Gradient tree boosting** [ edit ]

Gradient boosting is a machine learning technique for regression and classification problems, which works by combining multiple weak learners in a sequential manner. It is a special case of Friedman's proposed gradient boosting method, which improves the quality of fit of each base learner by iteratively adding new learners that focus on the residuals of the previous ones. The final prediction is the sum of all the learners' predictions.

Stochastic GB

→ RF: bagging

→ GBDT + row sampling + col-sampling

randomization as reg



Xgboost ;

Light GBM

~~sklearn~~

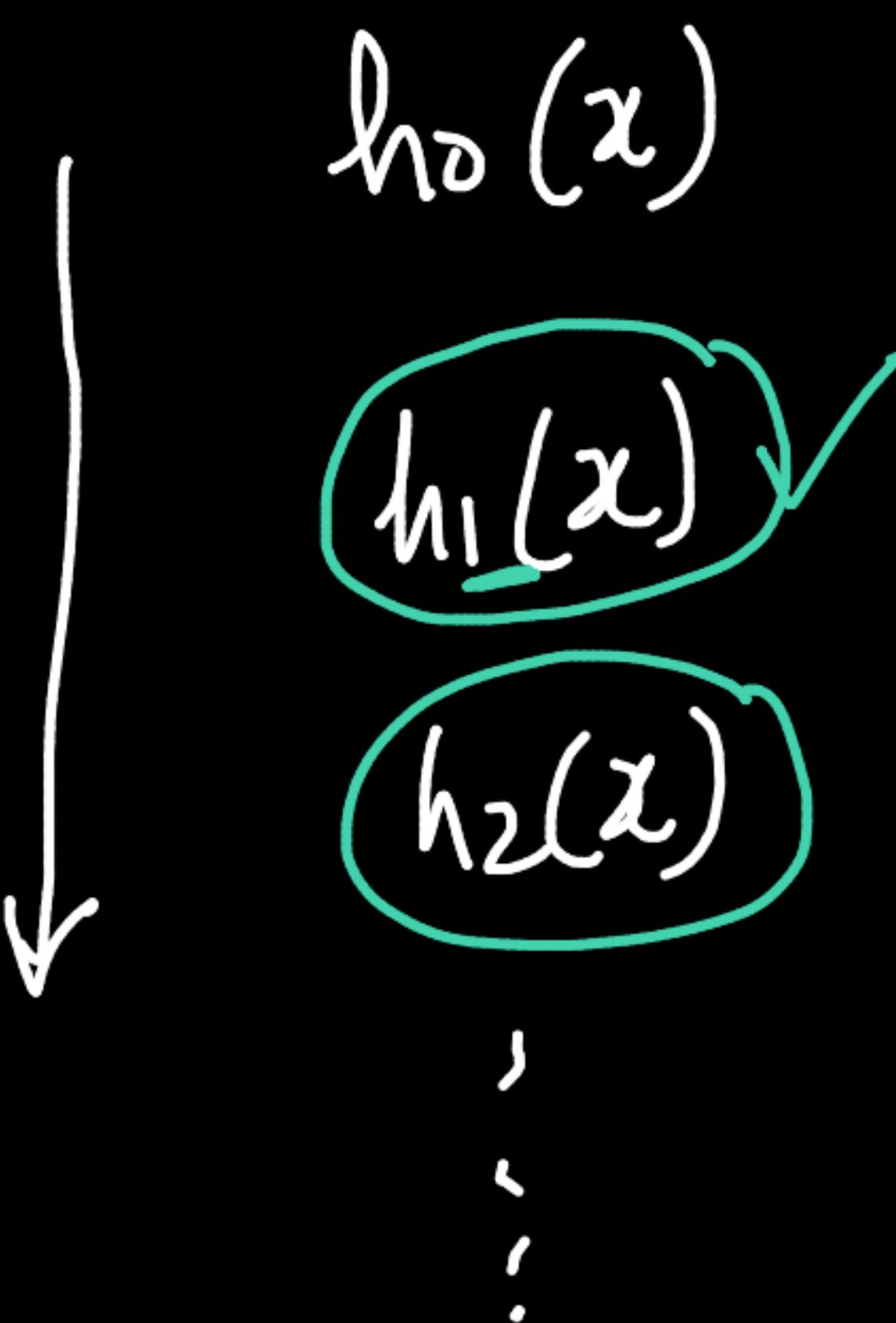
MSR

Brain time

much faster [RF :- trivially parallelize  
 $M_1, M_2, \dots, M_k \rightarrow$  different cores/  
system ✓]

GBDT  $\rightarrow$  Sequential - Model

Stochastic  
GBDT



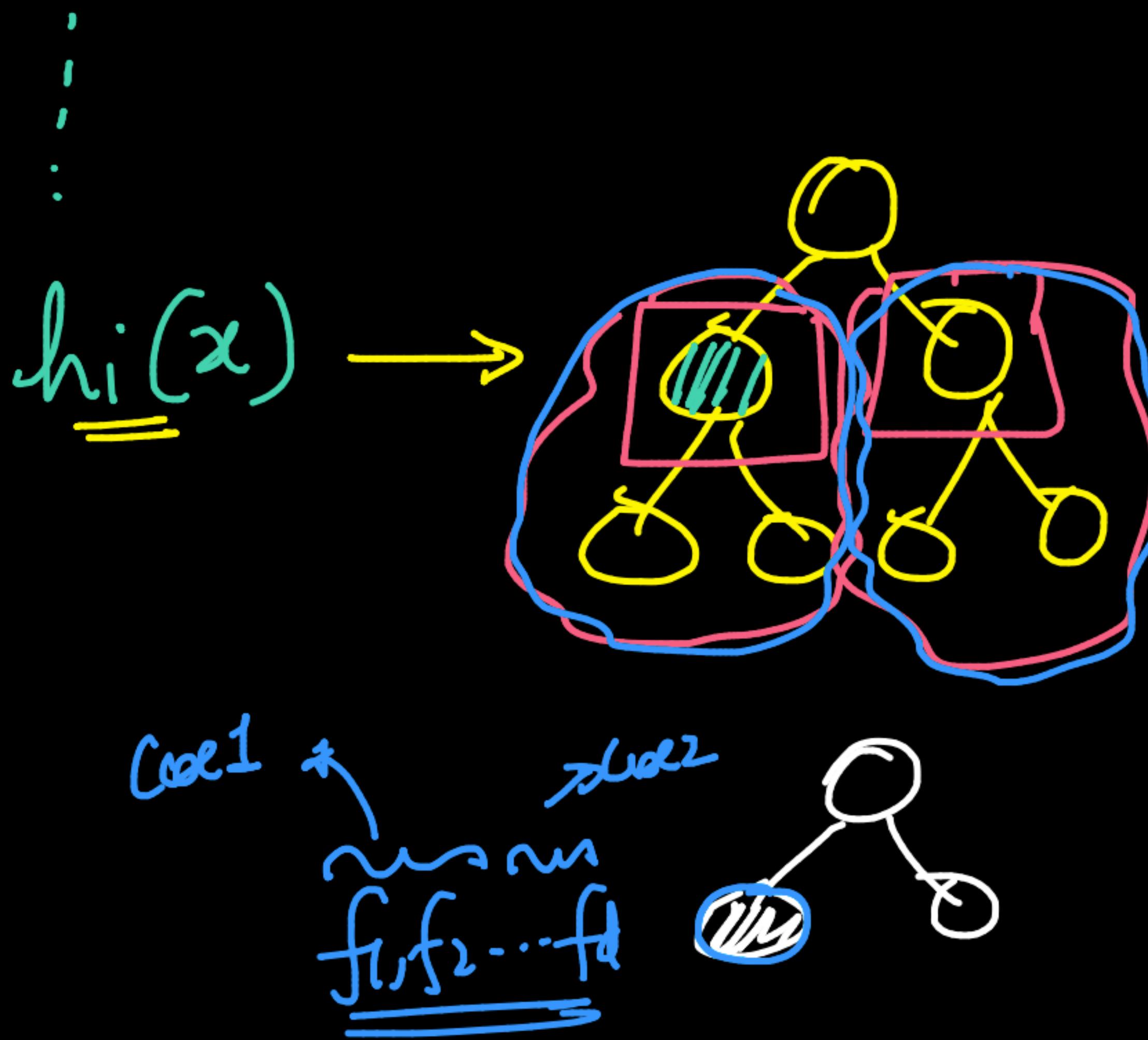
(let)  $M=100$

Slow to train



how can we Speedup  
GBDT

- row & column : small sampling : speedup

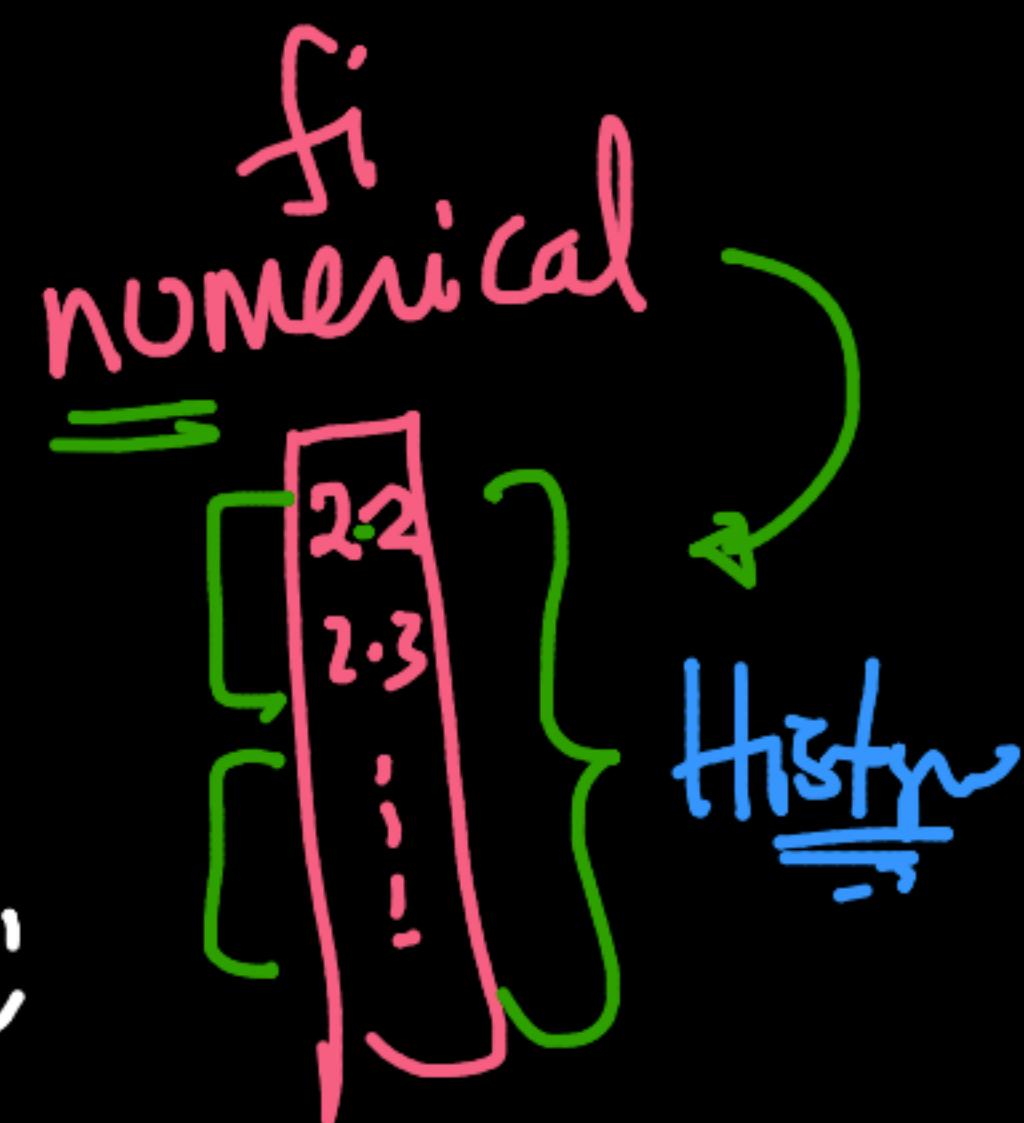


- individual tree level parallelization ✓
- IL2A for best split ✓



Xgboost: Millions  
of users

- official documentation
- top implementation → RF; GBDT ✓
- lots of code & algo optimization →
  - C++; Java; Scala;
  - Python, R . . .
- Multicore
- ~~Spark / Dask . . .~~



Pseudo-residuals | Gradient boosting | AdaBoost - Wikipedia | talk.dvi | 1603.02754.pdf | Python API Reference | XGBoost Document | Random Forests | LightGBM: A Highly Efficient Gradient Boosting Library | lightgbm.LGBMC | +

xgboost.readthedocs.io/en/stable/index.html

Installation Guide

Building From Source

Get Started with XGBoost

XGBoost Tutorials

Frequently Asked Questions

XGBoost User Forum

GPU Support

XGBoost Parameters

Prediction

Tree Methods

Python Package

R Package

JVM Package

Ruby Package

Swift Package

Julia Package

C Package

C++ Interface

CLI Interface

Contribute to XGBoost

## Contents

- Installation Guide
- Building From Source
- Get Started with XGBoost
- XGBoost Tutorials
  - Introduction to Boosted Trees
  - Introduction to Model IO
  - Distributed XGBoost with AWS YARN
  - Distributed XGBoost on Kubernetes
  - Distributed XGBoost with XGBoost4J-Spark
  - Distributed XGBoost with XGBoost4J-Spark-GPU
  - Distributed XGBoost with Dask
  - Distributed XGBoost with Ray
  - DART booster
  - Monotonic Constraints
  - Random Forests(TM) in XGBoost
  - Feature Interaction Constraints
  - Survival Analysis with Accelerated Failure Time
  - C API Tutorial
  - Text Input Format of DMatrix
  - Notes on Parameter Tuning
  - Custom Objective and Evaluation Metrics

Read the Docs v: stable

60 / 60

xgboost.readthedocs.io/en/stable/index.html

Installation Guide

Building From Source

Get Started with XGBoost

XGBoost Tutorials

Frequently Asked Questions

XGBoost User Forum

GPU Support

XGBoost Parameters

Prediction

Tree Methods

Python Package

R Package

VM Package

Ruby Package

Swift Package

Julia Package

C Package

C++ Interface

CLI Interface

Contribute to XGBoost

» XGBoost Documentation

Edit on GitHub

# XGBoost Documentation

XGBoost is an optimized distributed gradient boosting library designed to be highly **efficient**, **flexible** and **portable**. It implements machine learning algorithms under the [Gradient Boosting](#) framework. XGBoost provides a parallel tree boosting (also known as GBDT, GBM) that solve many data science problems in a fast and accurate way. The same code runs on major distributed environment (Hadoop, SGE, MPI) and can solve problems beyond billions of examples.

## Contents

- Installation Guide
- Building From Source
- Get Started with XGBoost
- XGBoost Tutorials
  - Introduction to Boosted Trees
  - Introduction to Model IO
  - Distributed XGBoost with AWS YARN
  - Distributed XGBoost on Kubernetes
  - Distributed XGBoost with XGBoost4J-Spark
  - Distributed XGBoost with XGBoost4J-Spark-GPU
  - Distributed XGBoost with Ray

Read the Docs v: stable

Pseudo-residuals | X Gradient boosting | W AdaBoost - Wikipedia | talk.dvi | 1603.02754.pdf | X Python API Reference | Python API Reference | Random Forests | X LightGBM: A High | lightgbm.LGBMC | X +

xgboost.readthedocs.io/en/stable/python/python\_api.html

GPU Support

XGBoost Parameters

Prediction

Tree Methods

Python Package

Python Package Introduction

Python API Reference

Callback Functions

Model

XGBoost Python Feature

Walkthrough

XGBoost Dask Feature Walkthrough

R Package

JVM Package

Ruby Package

Swift Package

Julia Package

C Package

C++ Interface

CLI Interface

Contribute to XGBoost

# Python API Reference

This page gives the Python API reference of xgboost, please also refer to Python Package Introduction for more information about the Python package.

- Global Configuration
- Core Data Structure
- Learning API
- Scikit-Learn API
- Plotting API
- Callback API
- Dask API
  - Dask extensions for distributed training
    - Optional dask configuration

sclearn



## Global Configuration

`xgboost.config_context(**new_config)`

Context manager for global XGBoost configuration.

Global configuration consists of a collection of parameters that can be applied in the global scope. See [Global Configuration](#) for the full list of parameters supported in the global config.

Read the Docs v: stable

xgboost.readthedocs.io/en/stable/parameter.html

- Installation Guide
- Building From Source
- Get Started with XGBoost
- XGBoost Tutorials
- Frequently Asked Questions
- XGBoost User Forum
- GPU Support

## XGBoost Parameters

- eta** [default=0.3, alias: `learning_rate`]
  - Step size shrinkage used in update to prevent overfitting. After each boosting step, we can directly get the weights of new features, and `eta` shrinks the feature weights to make the boosting process more conservative.
  - range: [0,1]
- gamma** [default=0, alias: `min_split_loss`]
  - Minimum loss reduction required to make a further partition on a leaf node of the tree. The larger `gamma` is, the more conservative the algorithm will be.
  - range: [0,∞]
- max\_depth** [default=6]
  - Maximum depth of a tree. Increasing this value will make the model more complex and more likely to overfit. 0 indicates no limit on depth. Beware that XGBoost aggressively consumes memory when training a deep tree. `exact` tree method requires non-zero value.
  - range: [0,∞]
- min\_child\_weight** [default=1]
  - Minimum sum of instance weight (hessian) required in a child. If the split results in a child with a sum of less than this value, the split is discarded.

Parameters for Tree Booster

Y

Read the Docs v: stable

63 / 63

xgboost.readthedocs.io/en/stable/parameter.html

- Installation Guide
- Building From Source
- Get Started with XGBoost
- XGBoost Tutorials
- Frequently Asked Questions
- XGBoost User Forum
- GPU Support

## XGBoost Parameters

- `disable_default_eval_metric` [default= `false`]
  - Flag to disable default metric. Set to 1 or `true` to disable.
- `num_feature` [set automatically by XGBoost, no need to be set by user]
  - Feature dimension used in boosting, set to maximum dimension of the feature

## Parameters for Tree Booster

- `eta` [default=0.3, alias: `learning_rate`]
  - Step size shrinkage used in update to prevent overfitting. After each boosting step, we can directly get the weights of new features, and `eta` shrinks the feature weights to make the boosting process more conservative.
  - range: [0,1]
- ~~`gamma`~~ [default=0, alias: `min_split_loss`]
  - Minimum loss reduction required to make a further partition on a leaf node of the tree. The larger `gamma` is, the more conservative the algorithm will be.
  - range: [0,∞]
- `max_depth` [default=6]
  - Maximum depth of a tree. Increasing this value will make the model more complex and more likely to overfit. 0 indicates no limit on depth. Beware that XGBoost aggressively consumes memory when training a deep tree. `exact` tree method requires non-zero value.
  - range: [0,∞]
- `min_child_weight` [default=1]
  - Minimum sum of instance weight (hessian) required in a child. If the split results in a child with a sum of less than this value, the split is discarded and no further partition is made.

depth ↓ underfit ↑ IG

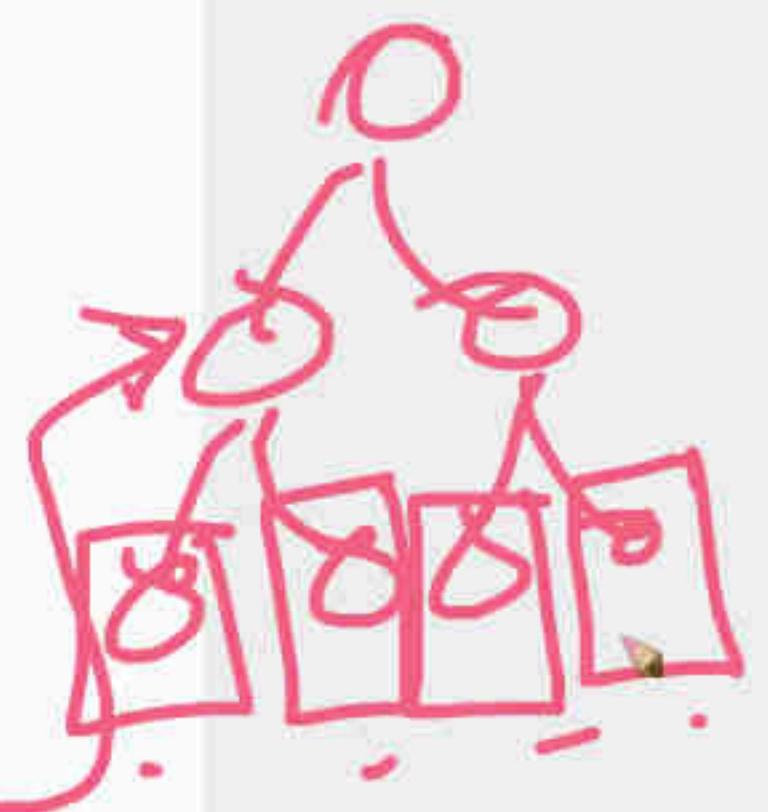
[Installation Guide](#)[Building From Source](#)[Get Started with XGBoost](#)[XGBoost Tutorials](#)[Frequently Asked Questions](#)[XGBoost User Forum](#)[GPU Support](#)

## XGBoost Parameters

[Prediction](#)[Tree Methods](#)[Python Package](#)[R Package](#)[JVM Package](#)[Ruby Package](#)[Swift Package](#)[Julia Package](#)[C Package](#)[C++ Interface](#)[CLI Interface](#)[Contribute to XGBoost](#)

# Parameters for Tree Booster

- `eta` [default=0.3, alias: `learning_rate`]
  - Step size shrinkage used in update to prevent overfitting. After each boosting step, we can directly get the weights of new features, and `eta` shrinks the feature weights to make the boosting process more conservative.
  - range: [0,1]
- `gamma` [default=0, alias: `min_split_loss`]
  - Minimum loss reduction required to make a further partition on a leaf node of the tree. The larger `gamma` is, the more conservative the algorithm will be.
  - range: [0,∞]
- `max_depth` [default=6]
  - Maximum depth of a tree. Increasing this value will make the model more complex and more likely to overfit. 0 indicates no limit on depth. Beware that XGBoost aggressively consumes memory when training a deep tree. `exact` tree method requires non-zero value.
  - range: [0,∞]
- `min_child_weight` [default=1]
  - Minimum sum of instance weight (hessian) needed in a child. If the tree partition step results in a leaf node with the sum of instance weight less than `min_child_weight`, then the building process will give up further partitioning. In linear regression task, this simply corresponds to minimum number of instances needed to be in each node. The larger `min_child_weight` is,



xgboost.readthedocs.io/en/stable/parameter.html

Search docs

- Installation Guide
- Building From Source
- Get Started with XGBoost
- XGBoost Tutorials
- Frequently Asked Questions
- XGBoost User Forum
- GPU Support
- XGBoost Parameters**
- Prediction
- Tree Methods
- Python Package
- R Package
- JVM Package
- Ruby Package
- Swift Package
- Julia Package
- C Package
- C++ Interface
- CLI Interface

• **gamma** [default=0, alias: `min_split_loss` ]

- Minimum loss reduction required to make a further partition on a leaf node of the tree. The larger `gamma` is, the more conservative the algorithm will be.
- range:  $[0, \infty]$

• **max\_depth** [default=6]

- Maximum depth of a tree. Increasing this value will make the model more complex and more likely to overfit. 0 indicates no limit on depth. Beware that XGBoost aggressively consumes memory when training a deep tree. `exact` tree method requires non-zero value.
- range:  $[0, \infty]$

• **min\_child\_weight** [default=1] 

- Minimum sum of instance weight (hessian) needed in a child. If the tree partition step results in a leaf node with the sum of instance weight less than `min_child_weight`, then the building process will give up further partitioning. In linear regression task, this simply corresponds to minimum number of instances needed to be in each node. The larger `min_child_weight` is, the more conservative the algorithm will be.
- range:  $[0, \infty]$

• **max\_delta\_step** [default=0]

- Maximum delta step we allow each leaf output to be. If the value is set to 0, it means there is no constraint. If it is set to a positive value, it can help making the update step more conservative. Usually this parameter is not needed, but it might help in logistic regression when class is extremely imbalanced. Set it to value of 1-10 might help control the update.
- range:  $[0, \infty]$



Read the Docs v: stable

xgboost.readthedocs.io/en/stable/parameter.html

Search docs

- Installation Guide
- Building From Source
- Get Started with XGBoost
- XGBoost Tutorials
- Frequently Asked Questions
- XGBoost User Forum
- GPU Support

## XGBoost Parameters

- Prediction
- Tree Methods
- Python Package
- R Package
- JVM Package
- Ruby Package
- Swift Package
- Julia Package
- C Package
- C++ Interface
- CLI Interface

- **gamma** [default=0, alias: `min_split_loss`]
  - Minimum loss reduction required to make a further partition on a leaf node of the tree. The larger `gamma` is, the more conservative the algorithm will be.
  - range:  $[0, \infty]$
- **max\_depth** [default=6]
  - Maximum depth of a tree. Increasing this value will make the model more complex and more likely to overfit. 0 indicates no limit on depth. Beware that XGBoost aggressively consumes memory when training a deep tree. `exact` tree method requires non-zero value.
  - range:  $[0, \infty]$
- **min\_child\_weight** [default=1]
  - Minimum sum of instance weight (hessian) needed in a child. If the tree partition step results in a leaf node with the sum of instance weight less than `min_child_weight`, then the building process will give up further partitioning. In linear regression task, this simply corresponds to minimum number of instances needed to be in each node. The larger `min_child_weight` is, the more conservative the algorithm will be.
  - range:  $[0, \infty]$
- **max\_delta\_step** [default=0]
  - Maximum delta step we allow each leaf output to be. If the value is set to 0, it means there is no constraint. If it is set to a positive value, it can help making the update step more conservative. Usually this parameter is not needed, but it might help in logistic regression when class is extremely imbalanced. Set it to value of 1-10 might help control the update.
  - range:  $[0, \infty]$

Read the Docs v: stable

Pseudo-residuals | Gradient boosting | AdaBoost - Wikipedia | talk.dvi | 1603.02754.pdf | Python API Reference | XGBoost Parameters | Random Forests | LightGBM: A Highly Efficient Gradient Boosting Library | lightgbm.LGBMC | +

xgboost.readthedocs.io/en/stable/parameter.html

Installation Guide

Building From Source

Get Started with XGBoost

XGBoost Tutorials

Frequently Asked Questions

XGBoost User Forum

GPU Support

**XGBoost Parameters**

Prediction

Tree Methods

Python Package

R Package

JVM Package

Ruby Package

Swift Package

Julia Package

C Package

C++ Interface

CLI Interface

Contribute to XGBoost

• `max_delta_step` [default=0]

- range:  $[0, \infty]$
- Maximum delta step we allow each leaf output to be. If the value is set to 0, it means there is no constraint. If it is set to a positive value, it can help making the update step more conservative. Usually this parameter is not needed, but it might help in logistic regression when class is extremely imbalanced. Set it to value of 1-10 might help control the update.

• `subsample` [default=1]

- Subsample ratio of the training instances. Setting it to 0.5 means that XGBoost would randomly sample half of the training data prior to growing trees. and this will prevent overfitting. Subsampling will occur once in every boosting iteration.
- range:  $(0, 1]$

• `sampling_method` [default= `uniform` ]

- The method to use to sample the training instances.
- `uniform` : each training instance has an equal probability of being selected. Typically set `subsample`  $\geq 0.5$  for good results.
- `gradient_based` : the selection probability for each training instance is proportional to the regularized absolute value of gradients (more specifically,  $\sqrt{g^2 + \lambda h^2}$ ). `subsample` may be set to as low as 0.1 without loss of model accuracy. Note that this sampling method is only supported when `tree_method` is set to `gpu_hist`; other tree methods only support `uniform` sampling.

• `cols` [ ]

Read the Docs v: stable

68 / 68

*mn -Sanjiv*

xgboost.readthedocs.io/en/stable/parameter.html

colsample\_bynode is the subsample ratio of columns for each node (split). Subsampling occurs once every time a new split is evaluated. Columns are subsampled from the set of columns chosen for the current level.

colsample\_by\* parameters work cumulatively. For instance, the combination `{'colsample_bytree':0.5, 'colsample_bylevel':0.5, 'colsample_bynode':0.5}` with 64 features will leave 8 features to choose from at each split.

Using the Python or the R package, one can set the feature\_weights for DMatrix to define the probability of each feature being selected when using column sampling. There's a similar parameter for fit method in sklearn interface.

- lambda [default=1, alias: reg\_lambda]  $\alpha (\sum w_i^2)$   
L2 regularization term on weights. Increasing this value will make model more conservative.
- alpha [default=0, alias: reg\_alpha]
  - L1 regularization term on weights. Increasing this value will make model more conservative.
- tree\_method string [default= auto]
  - The tree construction algorithm used in XGBoost. See description in the [reference paper](#) and [Tree Methods](#).
  - XGBoost supports approx, hist and gpu\_hist for distributed training. Experimental support for external memory is available for approx and gpu\_hist.
  - Choices: auto, exact, approx, hist, gpu\_hist, this is a combination of commonly used updaters. For other updaters like refresh, set the parameter updater directly.
  - auto · Use heuristic to choose the fastest method

Read the Docs v: stable

Pseudo-residuals | Gradient boosting | AdaBoost - Wikipedia | talk.dvi | 1603.02754.pdf | Python API Reference | XGBoost Parameters | Random Forests | LightGBM: A Highly Efficient Gradient Boosting Library | lightgbm.LGBMC | +

xgboost.readthedocs.io/en/stable/parameter.html

parameter for `fit` method in sklearn interface.

• `lambda` [default=1, alias: `reg_lambda`]  
◦ L2 regularization term on weights. Increasing this value will make model more conservative.

• `alpha` [default=0, alias: `reg_alpha`]  
◦ L1 regularization term on weights. Increasing this value will make model more conservative.

• `tree_method` string [default= `auto`]  
◦ The tree construction algorithm used in XGBoost. See description in the [reference paper](#) and [Tree Methods](#).  
◦ XGBoost supports `approx`, `hist` and `gpu_hist` for distributed training. Experimental support for external memory is available for `approx` and `gpu_hist`.  
◦ Choices: `auto`, `exact`, `approx`, `hist`, `gpu_hist`, this is a combination of commonly used updaters. For other updaters like `refresh`, set the parameter `updater` directly.  
▪ `auto`: Use heuristic to choose the fastest method.  
▪ For small dataset, exact greedy (`exact`) will be used.  
▪ For larger dataset, approximate algorithm (`approx`) will be chosen. It's recommended to try `hist` and `gpu_hist` for higher performance with large dataset. (`gpu_hist`) has support for `external memory`.  
▪ Because old behavior is always use exact greedy in single machine, user will get a message when approximate algorithm is chosen to notify this choice.  
✓ `exact`: Exact greedy algorithm. Enumerates all split candidates.  
▪ `approx`: Approximate greedy algorithm using quantile sketch and gradient histogram.  
▪ `hist`: Faster histogram optimized approximate greedy algorithm.

dmlc  
**XGBoost**

stable

Search docs

Installation Guide

Building From Source

Get Started with XGBoost

XGBoost Tutorials

Frequently Asked Questions

XGBoost User Forum

GPU Support

XGBoost Parameters

Prediction

Tree Methods

Python Package

R Package

JVM Package

Read the Docs v: stable



stable

Search docs

Installation Guide

Building From Source

Get Started with XGBoost

XGBoost Tutorials

Frequently Asked Questions

XGBoost User Forum

GPU Support

XGBoost Parameters

Prediction

Tree Methods

Python Package

R Package

JVM Package

Ruby Package

Swift Package

Julia Package

- `alpha` [default=0, alias: `reg_alpha`]
  - L1 regularization term on weights. Increasing this value will make model more conservative.
- `tree_method` string [default= `auto`]
  - The tree construction algorithm used in XGBoost. See description in the [reference paper](#) and [Tree Methods](#).
  - XGBoost supports `approx`, `hist` and `gpu_hist` for distributed training. Experimental support for external memory is available for `approx` and `gpu_hist`.
  - Choices: `auto`, `exact`, `approx`, `hist`, `gpu_hist`, this is a combination of commonly used updaters. For other updaters like `refresh`, set the parameter `updater` directly.
    - `auto` : Use heuristic to choose the fastest method.
      - For small dataset, exact greedy (`exact`) will be used.
      - For larger dataset, approximate algorithm (`approx`) will be chosen. It's recommended to try `hist` and `gpu_hist` for higher performance with large dataset. (`gpu_hist`) has support for `external memory`.
      - Because old behavior is always use exact greedy in single machine, user will get a message when approximate algorithm is chosen to notify this choice.
    - `exact` : Exact greedy algorithm. Enumerates all split candidates.
    - `approx` : Approximate greedy algorithm using quantile sketch and gradient histogram.
    - `hist` : Faster histogram optimized approximate greedy algorithm.
    - `gpu_hist` : GPU implementation of `hist` algorithm.
- `sketch_eps` [default=0.03]
  - Only used for `updater=grow_local_histmaker`
  - T compared to directly select



stable

Search docs

Installation Guide

Building From Source

Get Started with XGBoost

XGBoost Tutorials

Frequently Asked Questions

XGBoost User Forum

GPU Support

XGBoost Parameters

Prediction

Tree Methods

Python Package

R Package

JVM Package

Ruby Package

Swift Package

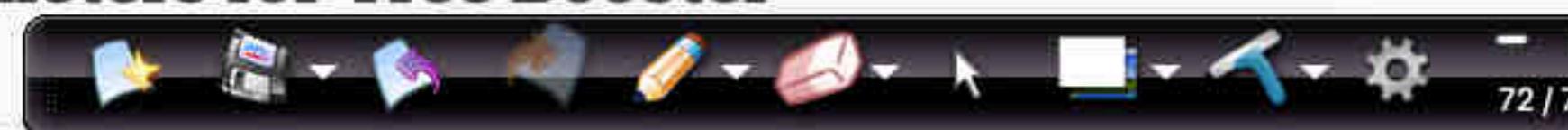
Java Package

## General Parameters

- **booster** [default= `gbtree`] 

Which booster to use. Can be `gbtree`, `gblinear` or `dart`; `gbtree` and `dart` use tree based models while `gblinear` uses linear functions.
- **verbosity** [default=1]
  - Verbosity of printing messages. Valid values are 0 (silent), 1 (warning), 2 (info), 3 (debug). Sometimes XGBoost tries to change configurations based on heuristics, which is displayed as warning message. If there's unexpected behaviour, please try to increase value of verbosity.
- **validate\_parameters** [default to `false`, except for Python, R and CLI interface]
  - When set to True, XGBoost will perform validation of input parameters to check whether a parameter is used or not. The feature is still experimental. It's expected to have some false positives.
- **nthread** [default to maximum number of threads available if not set]
  - Number of parallel threads used to run XGBoost. When choosing it, please keep thread contention and hyperthreading in mind.
- **disable\_default\_eval\_metric** [default= `false`]
  - Flag to disable default metric. Set to 1 or `true` to disable.
- **num\_feature** [set automatically by XGBoost, no need to be set by user]
  - Feature dimension used in boosting, set to maximum dimension of the feature

## Parameters for Tree Booster





stable

Search docs

Installation Guide

Building From Source

Get Started with XGBoost

XGBoost Tutorials

Frequently Asked Questions

XGBoost User Forum

GPU Support

XGBoost Parameters

Prediction

Tree Methods

Python Package

R Package

JVM Package

Ruby Package

Swift Package

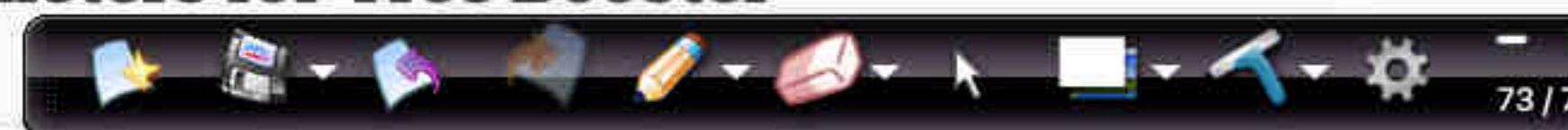
## General Parameters

- **booster** [default= `gbtree` ]
  - Which booster to use. Can be `gbtree`, `gblinear` or `dart`; `gbtree` and `dart` use tree based models while `gblinear` uses linear functions.
- **verbosity** [default=1]
  - Verbosity of printing messages. Valid values are 0 (silent), 1 (warning), 2 (info), 3 (debug). Sometimes XGBoost tries to change configurations based on heuristics, which is displayed as warning message. If there's unexpected behaviour, please try to increase value of verbosity.
- **validate\_parameters** [default to `false`, except for Python, R and CLI interface]
  - When set to True, XGBoost will perform validation of input parameters to check whether a parameter is used or not. The feature is still experimental. It's expected to have some false positives.



- **nthread** [default to maximum number of threads available if not set]
  - Number of parallel threads used to run XGBoost. When choosing it, please keep thread contention and hyperthreading in mind.
- **disable\_default\_eval\_metric** [default= `false`]
  - Flag to disable default metric. Set to 1 or `true` to disable.
- **num\_feature** [set automatically by XGBoost, no need to be set by user]
  - Feature dimension used in boosting, set to maximum dimension of the feature

## Parameters for Tree Booster



# hyper-params Xgboost

- ✓ ↗ M
  - ↗ max-depth
  - ↗  $\gamma$ : learning rate
  - ↗ col. Sampling / row sampling
- =====

LightGBM  
        
(MS)

→ faster GBDT

⇒



Search docs

Installation Guide

Building From Source

Get Started with XGBoost

XGBoost Tutorials

Frequently Asked Questions

XGBoost User Forum

GPU Support

XGBoost Parameters

Prediction

Tree Methods

## Python Package

Python Package Introduction

Python API Reference

## Scikit-Learn API

Scikit-Learn Wrapper interface for XGBoost.

`class xgboost.XGBRegressor(*, objective='reg:squarederror', **kwargs)`

Bases: `xgboost.sklearn.XGBModel`, `sklearn.base.RegressorMixin`

Implementation of the scikit-learn API for XGBoost regression.

Parameters:

- `n_estimators` (`int`) – Number of gradient boosted trees. Equivalent to number of boosting rounds.
- `max_depth` (`Optional[int]`) – Maximum tree depth for base learners.
- `max_leaves` – Maximum number of leaves; 0 indicates no limit.
- `max_bin` – If using histogram-based algorithm, maximum number of bins per feature
- `grow_policy` – Tree growing policy. 0: favor splitting at nodes closest to the node, i.e. grow depth-wise. 1: favor splitting at nodes with highest loss change.
- `learning_rate` (`Optional[float]`) – Boosting learning rate (xgb's "eta")
- `verbosity` (`Optional[int]`) – The degree of verbosity. Valid values are 0 (silent) - 3 (debug).
- `objective` (`Union[str, Callable[[numpy.ndarray, numpy.ndarray], Tuple[numpy.ndarray, numpy.ndarray]], NoneType]`) – Specify the learning task and the corresponding learning objective or a custom objective function to be used (see note below).

xgboost.readthedocs.io/en/stable/tutorials/rf.html

## Standalone Random Forest With XGBoost API

The following parameters must be set to enable random forest training.

- `booster` should be set to `gbtree` as we are training forests. Note that as this is the default, this parameter needn't be set explicitly.
- `subsample` must be set to a value less than 1 to enable random selection of training cases (rows).
- One of `colsample_by*` parameters must be set to a value less than 1 to enable random selection of columns. Normally, `colsample_bynode` would be set to a value less than 1 to randomly sample columns at each tree split.
- `num_parallel_tree` should be set to the size of the forest being trained.
- `num_boost_round` should be set to 1 to prevent XGBoost from boosting multiple random forests. Note that this is a keyword argument to `train()`, and is not part of the parameter dictionary.
- `eta` (alias: `learning_rate`) must be set to 1 when training random forest regression.
- `random_state` can be used to seed the random number generator.

Other parameters should be set in a similar way they are set for gradient boosting. For instance, `objective` will typically be `reg:squarederror` for regression and `binary:logistic` for classification, `lambda` should be set according to a desired regularization weight, etc.

If both `num_parallel_tree` and `num_boost_round` are greater than 1, training will use a combination of random forest and gradient boosting strategy. It will perform `num_boost_round` rounds, boosting a random forest of `num_parallel_tree` trees in each round. If `subsample` is set to 1 and sampling is not enabled, the final model will be a random forest.

Max-depth = 10

Distributed XGBoost with XGBoost4J-Spark-GPU

Distributed XGBoost with Dask

Distributed XGBoost with Ray

DART booster

Monotonic Constraints

Random Forests(TM) in XGBoost

Feature Interaction Constraints

Survival Analysis with Accelerated Failure Time

C API Tutorial

Text Input Format of DMatrix

Notes on Parameter Tuning

Using XGBoost External Memory

Version

Custom Objective and Evaluation Metric

Categorical Data

Multiple Outputs

Frequently Asked Questions

XGBoost User Forum

Read the Docs v: stable

xgboost.readthedocs.io/en/stable/tutorials/rf.html

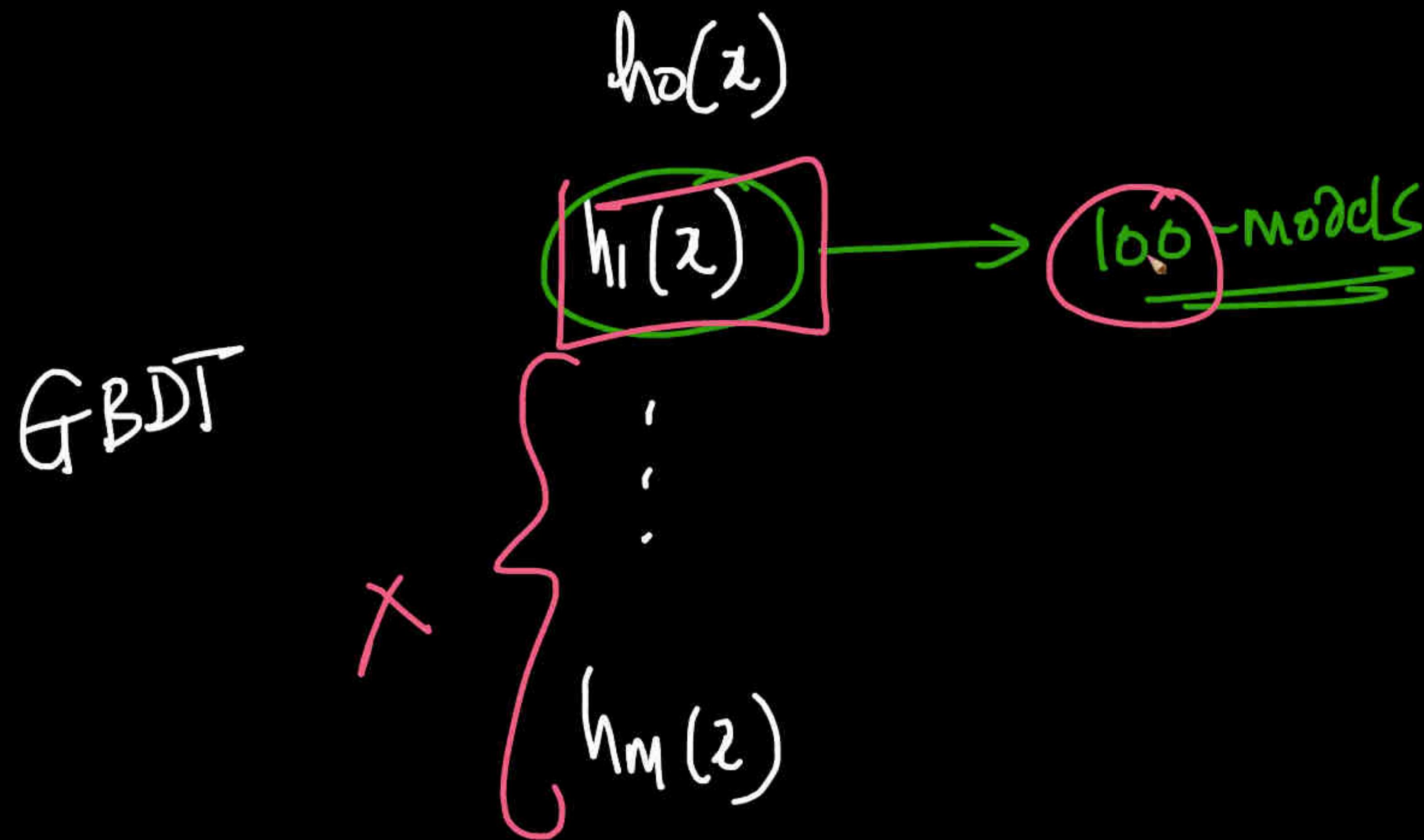
- One or `colsample_by*` parameters must be set to a value less than 1 to enable random selection of columns. Normally, `colsample_bynode` would be set to a value less than 1 to randomly sample columns at each tree split.
- `num_parallel_tree` should be set to the size of the forest being trained.
- `num_boost_round` should be set to 1 to prevent XGBoost from boosting multiple random forests. Note that this is a keyword argument to `train()`, and is not part of the parameter dictionary.
- `eta` (alias: `learning_rate`) must be set to 1 when training random forest regression.
- `random_state` can be used to seed the random number generator.

Other parameters should be set in a similar way they are set for gradient boosting. For instance, `objective` will typically be `reg:squarederror` for regression and `binary:logistic` for classification, `lambda` should be set according to a desired regularization weight, etc.

If both `num_parallel_tree` and `num_boost_round` are greater than 1, training will use a combination of random forest and gradient boosting strategy. It will perform `num_boost_round` rounds, boosting a random forest of `num_parallel_tree` trees at each round. If early stopping is not enabled, the final model will consist of `num_parallel_tree * num_boost_round` trees.

Here is a sample parameter dictionary for training a random forest on a GPU using xgboost:

```
params = {
    'colsample_bynode': 0.8,
    'learning_rate': 1,
    'max_depth': 5,
    'num_parallel_tree': 100,
    'objective': 'binary:logistic',
    'sub'
    'tre'
```



Version
Custom Objective and Evaluation Metric
Categorical Data
Multiple Outputs
Frequently Asked Questions
XGBoost User Forum
GPU Support
XGBoost Parameters
Prediction
Tree Methods
Python Package
R Package
JVM Package
Ruby Package
Swift Package
Julia Package
C Package
C++ Interface
CLI Interface
Contribute to XGBoost

lambda should be set according to a desired regularization weight, etc.

If both num\_parallel\_tree and num\_boost\_round are greater than 1, training will use a combination of random forest and gradient boosting strategy. It will perform num\_boost\_round rounds, boosting a random forest of num\_parallel\_tree trees at each round. If early stopping is not enabled, the final model will consist of num\_parallel\_tree \* num\_boost\_round trees.

Here is a sample parameter dictionary for training a random forest on a GPU using xgboost:

```
params = {  
    'colsample_bynode': 0.8,  
    'learning_rate': 1,  
    'max_depth': 5,  
    'num_parallel_tree': 100,  
    'objective': 'binary:logistic',  
    'subsample': 0.8,  
    'tree_method': 'gpu_hist'  
}
```

RF with 100 learn

A random forest model can then be trained as follows:

```
bst = train(params, dmatrix, num_boost_round=1)
```

## Standalone Random Forest With Scikit-Learn-Like API

Version
Custom Objective and Evaluation Metric
Categorical Data
Multiple Outputs
Frequently Asked Questions
XGBoost User Forum
GPU Support
XGBoost Parameters
Prediction
Tree Methods
Python Package
R Package
JVM Package
Ruby Package
Swift Package
Julia Package
C Package
C++ Interface
CLI Interface
Contribute to XGBoost

lambda should be set according to a desired regularization weight, etc.

If both num\_parallel\_tree and num\_boost\_round are greater than 1, training will use a combination of random forest and gradient boosting strategy. It will perform num\_boost\_round rounds, boosting a random forest of num\_parallel\_tree trees at each round. If early stopping is not enabled, the final model will consist of num\_parallel\_tree \* num\_boost\_round trees.

Here is a sample parameter dictionary for training a random forest on a GPU using xgboost:

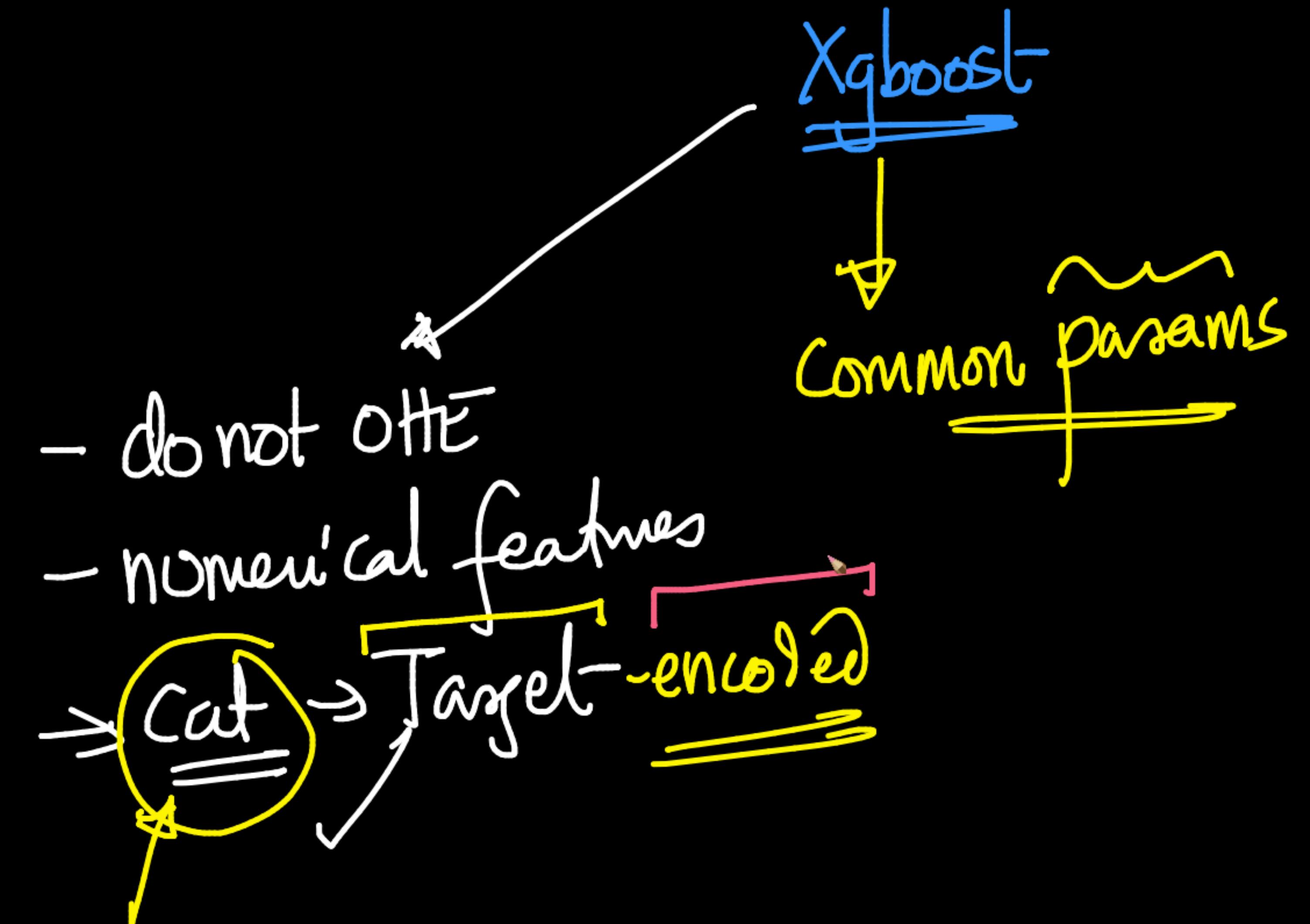
```
{  
    'colsample_bynode': 0.8,  
    'learning_rate': 1,  
    'max_depth': 5,  
    'num_parallel_tree': 100,  
    'objective': 'binary:logistic',  
    'subsample': 0.8,  
    'tree_method': 'gpu_hist'  
}
```

A random forest model can then be trained as follows:

```
bst = train(params, dmatrix, num_boost_round=1)
```

## Standalone Random Forest With Scikit-Learn-Like API

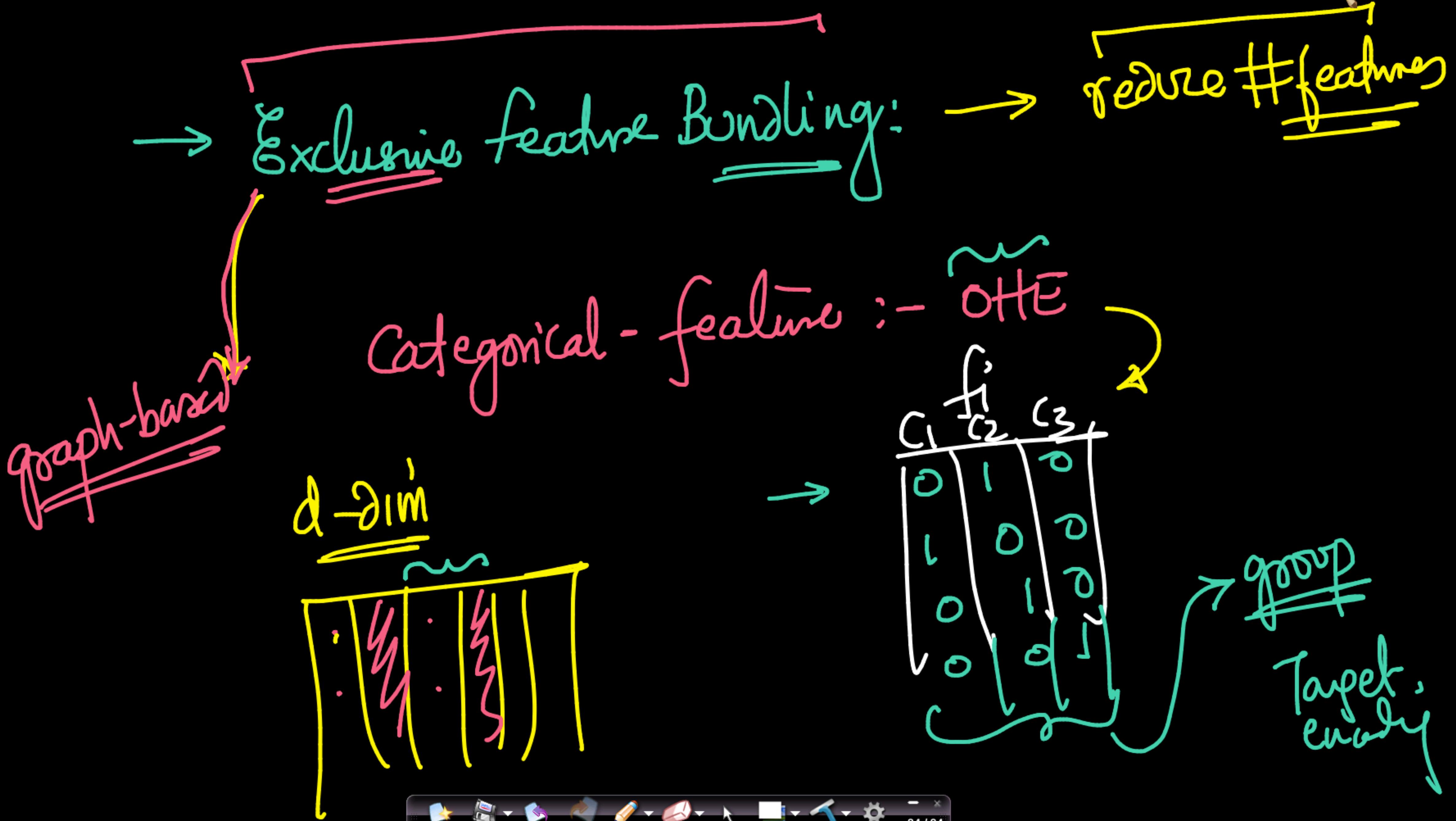
XGBoost provides a scikit-learn like API for using the XGBoost random forest functionality. They are basically versions of `XGBRFCL`, `XGBRF`, `XGBRFM` and `XGBRFSS` that train random forest instead of

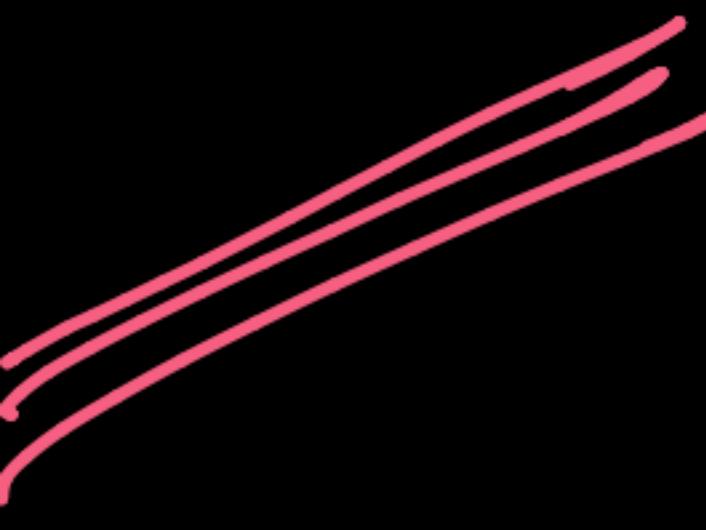


→ faster than XgBoost

→ GOSS: Gradient based  
m<sup>th</sup> model:  $\{(x_i, r_i^m)\}_{i=1}^n$   
# datapoints ↓

lightGBM  
(2017)  
one side Sampling  
 $\hookrightarrow$  drop where it is small  
 $\hookrightarrow$  Sampling ✓





GBT >

flexibility  $\downarrow$   
choosing  $L$

AdaBoost  
 $\downarrow$   
adaptive

Pseudo-residuals X Gradient boosting X AdaBoost - Wikipedia X talk.dvi X 1603.02754.pdf X Python API Reference X XGBoost Parameters X Random Forests X LightGBM: A Highly Efficient Gradient Boosting Library X lightgbm.LGBMC X +

csie.ntu.edu.tw/~mhyang/course/u0030/papers/schapire.pdf

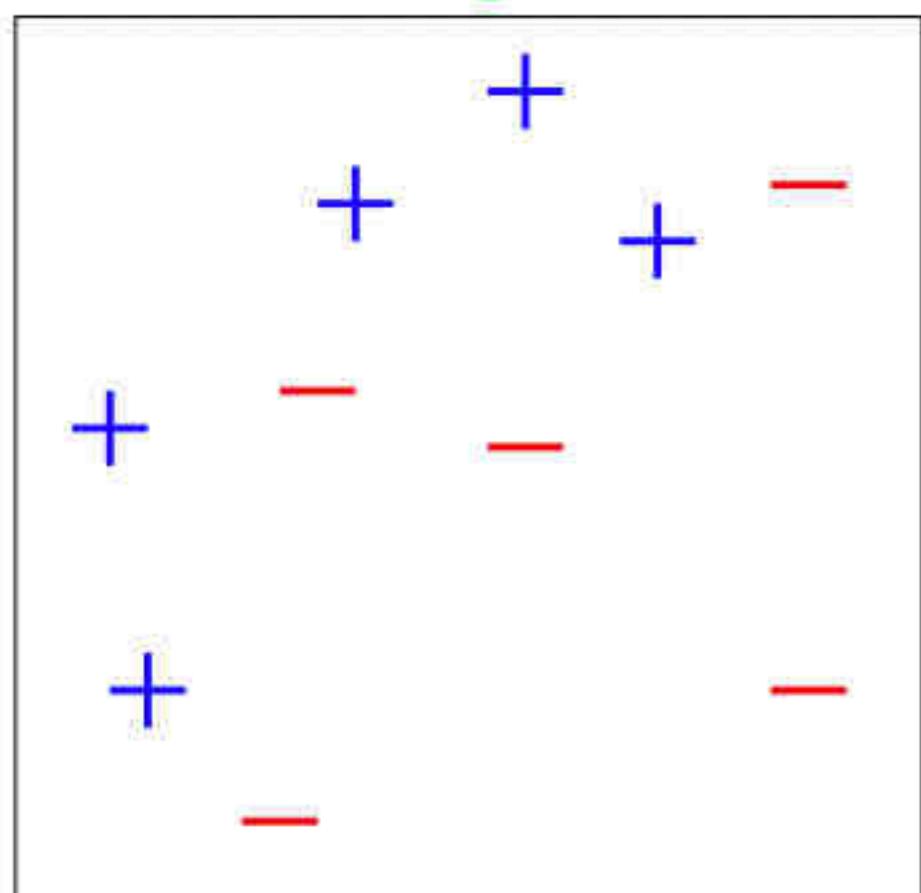
14 / 74 | - 125% + | ☰ ⚡

talk.dvi

Toy Example

"AdaBoost"

$D_1$



: base-learner

weak classifier

86 / 86

Pseudo-residuals | Gradient boosting | AdaBoost - Wikipedia | talk.dvi | 1603.02754.pdf | Python API Reference | XGBoost Parameters | Random Forests | LightGBM: A Highly Efficient Gradient Boosting Library | lightgbm.LGBMC | +

csie.ntu.edu.tw/~mhyang/course/u0030/papers/schapire.pdf

15 / 74 | - 125% + | ☰ ⚡

talk.dvi

Round 1

$h_0(x) + \underline{h_1(x)}$

$h_1$

$\oplus$

$\oplus$

$\oplus$

$\ominus$

$\ominus$

$\ominus$

$\oplus$

$\ominus$

$\oplus$

$\ominus$

$\ominus$

$\ominus$

$D_2$

$\oplus$

$\ominus$

$\oplus$

$\ominus$

$\oplus$

$\ominus$

$\oplus$

$\ominus$

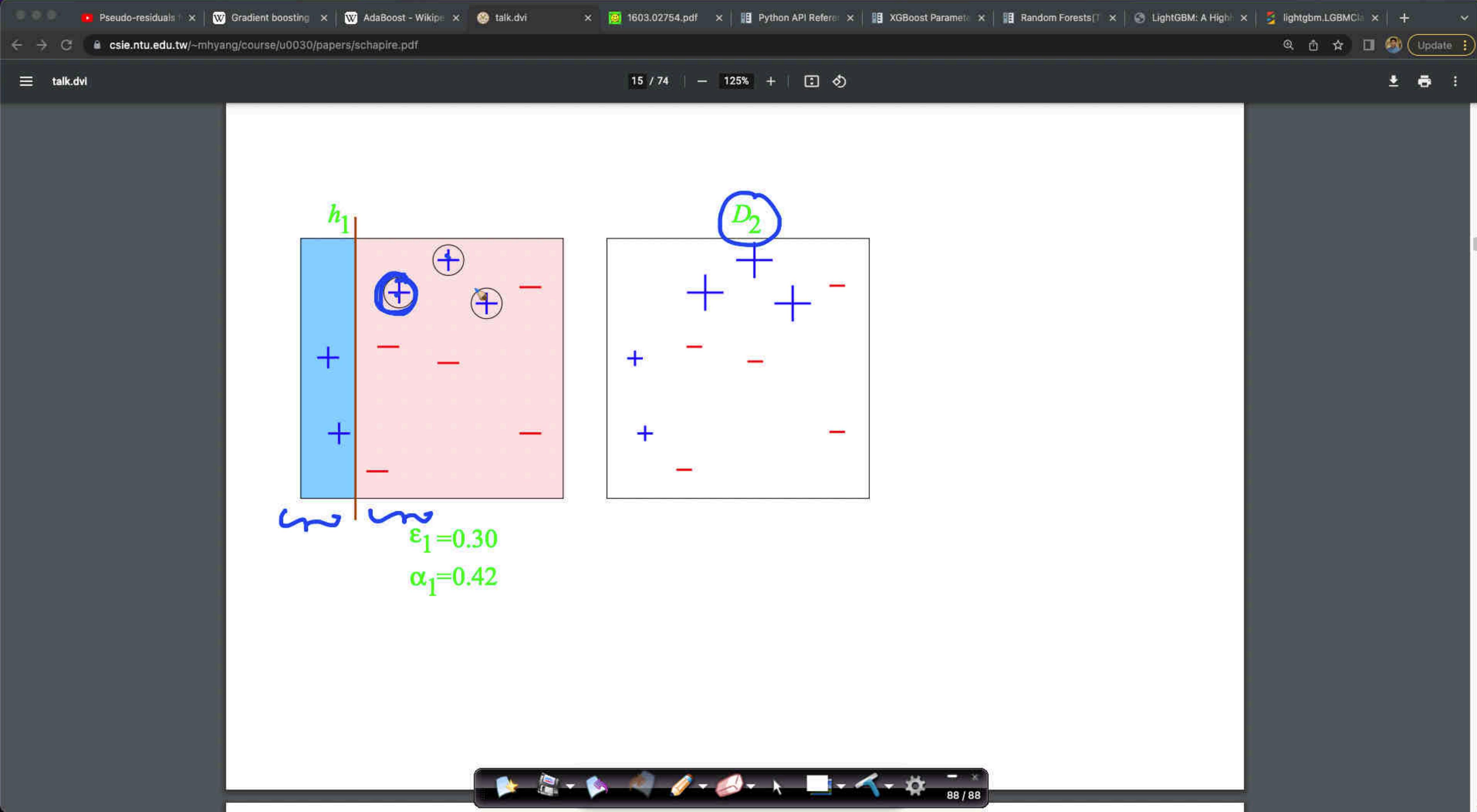
$\ominus$

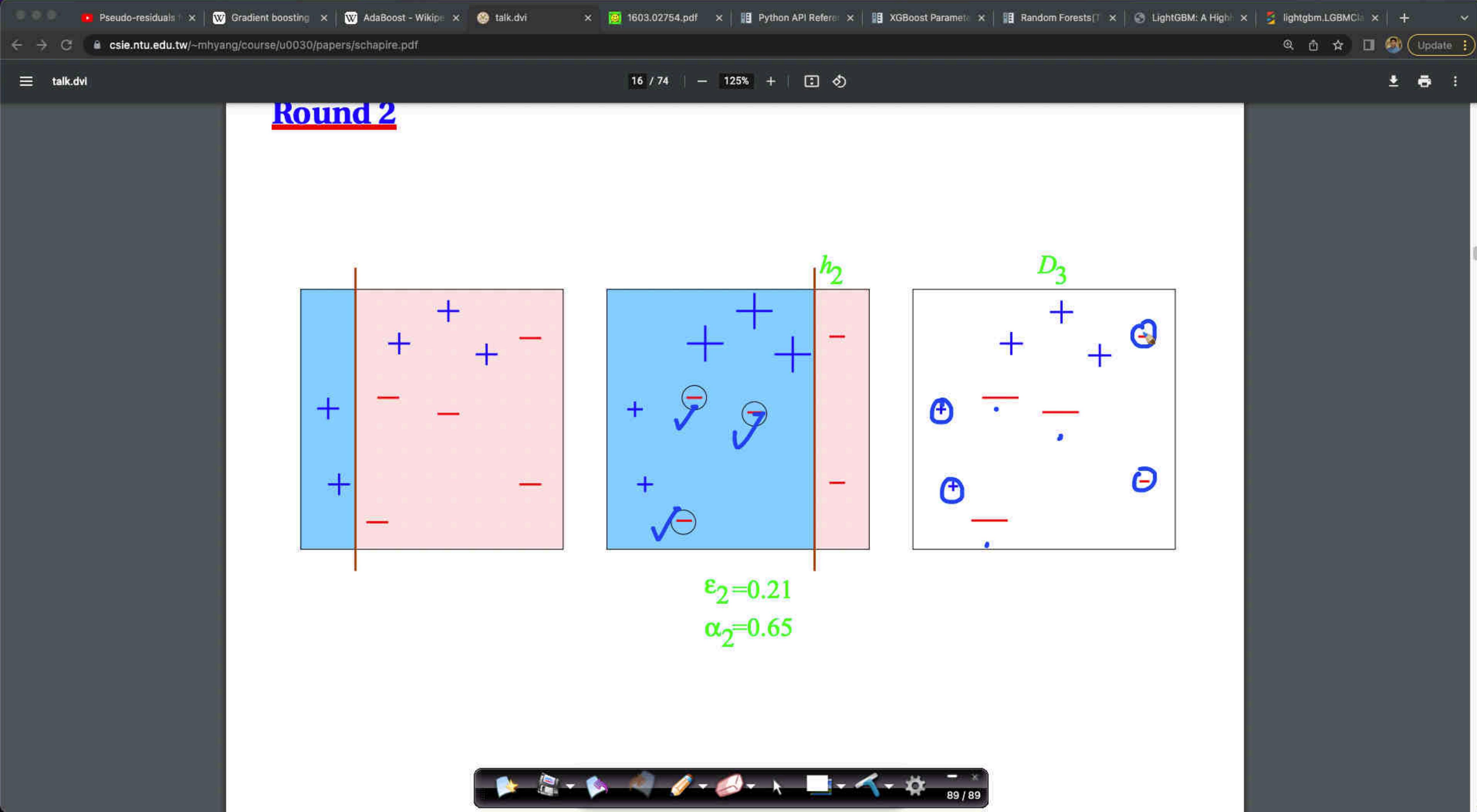
$\ominus$

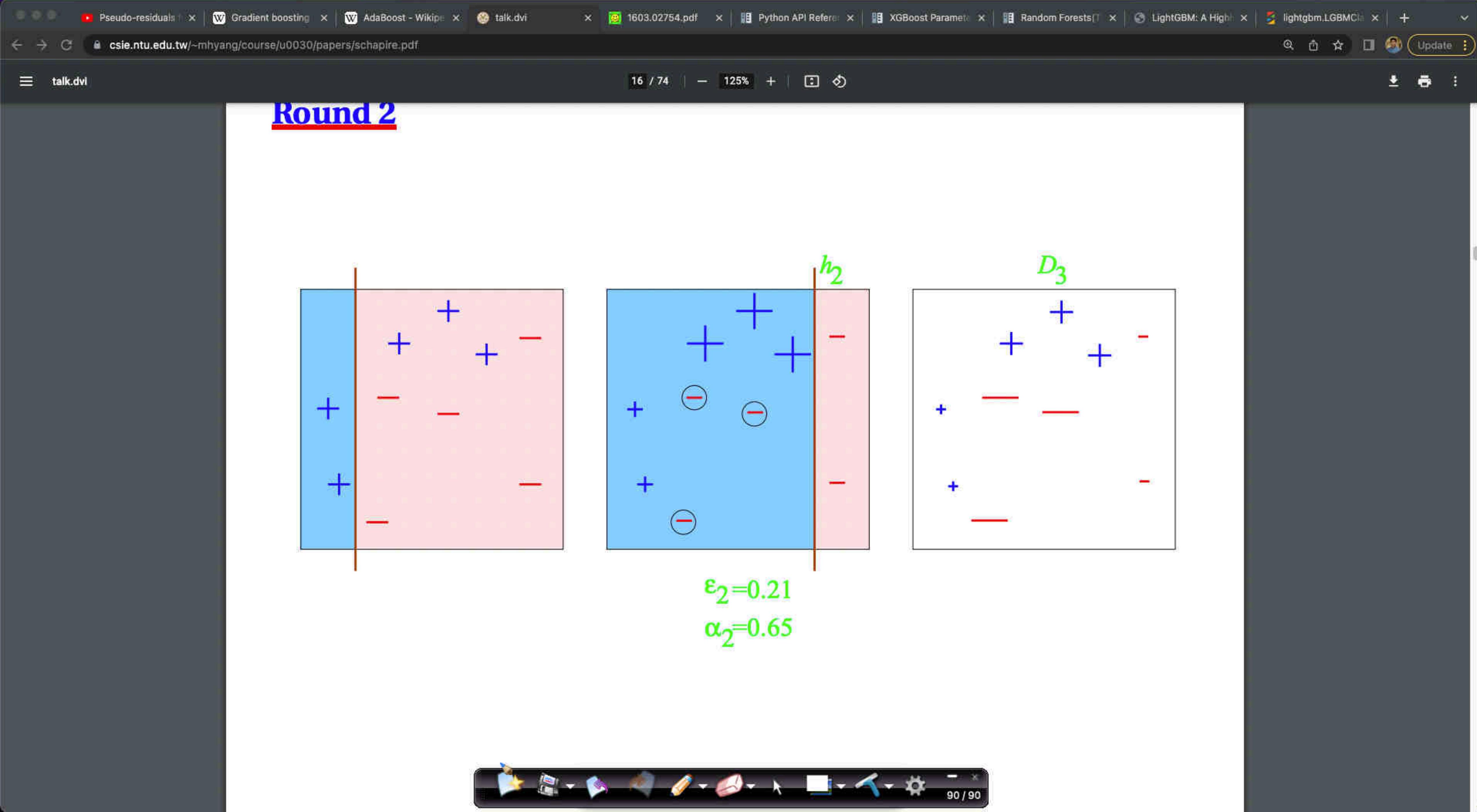
$\epsilon_1 = 0.30$

$\alpha_1 = 0.42$

87 / 87







Pseudo-residuals x Gradient boosting x AdaBoost - Wikipedia x talk.dvi x 1603.02754.pdf x Python API Refer... x XGBoost Param... x Random Forests x LightGBM: A High... x lightgbm.LGBMC x +

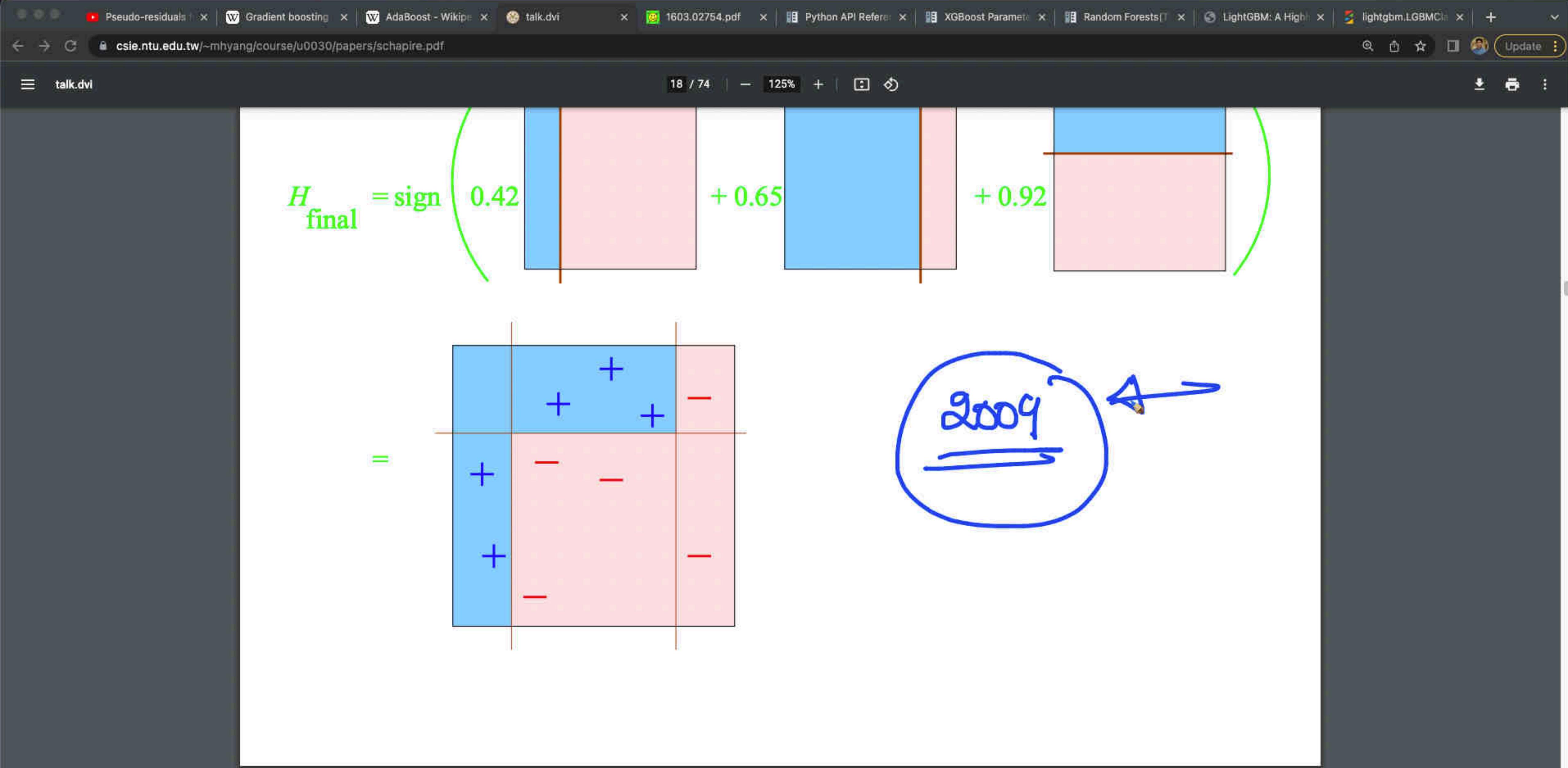
csie.ntu.edu.tw/~mhyang/course/u0030/papers/schapire.pdf

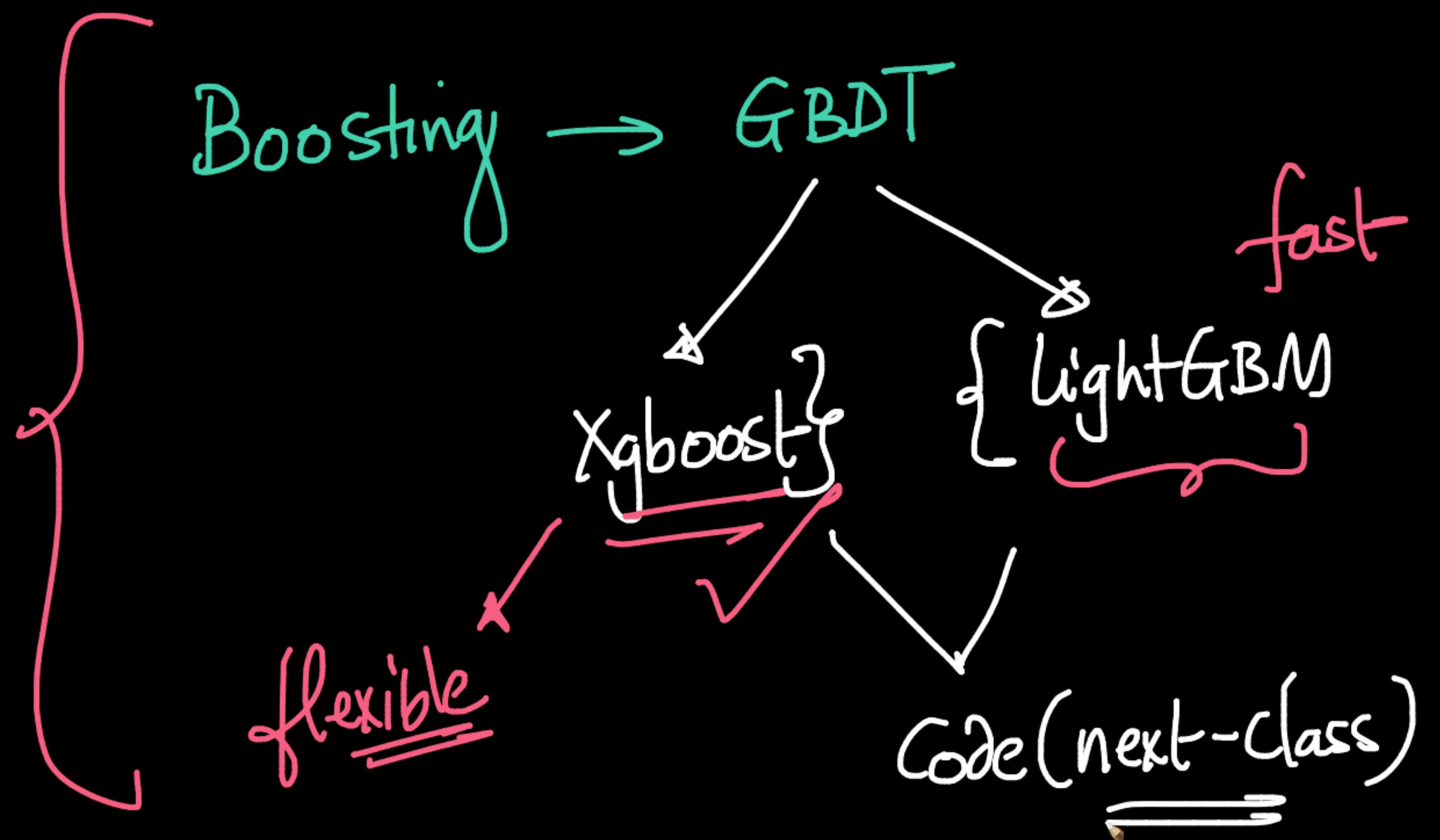
18 / 74 | - 125% + | ☰ ⚡

## Final Classifier

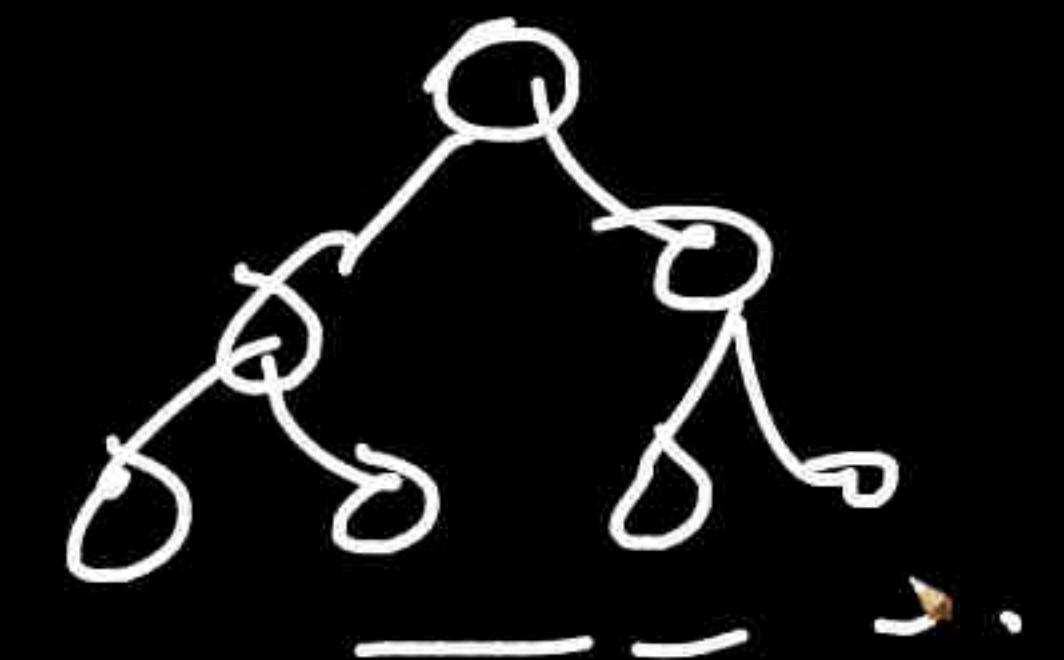
$H_{\text{final}} = \text{sign} (0.42 h_1 + 0.65 h_2 + 0.92 h_3)$

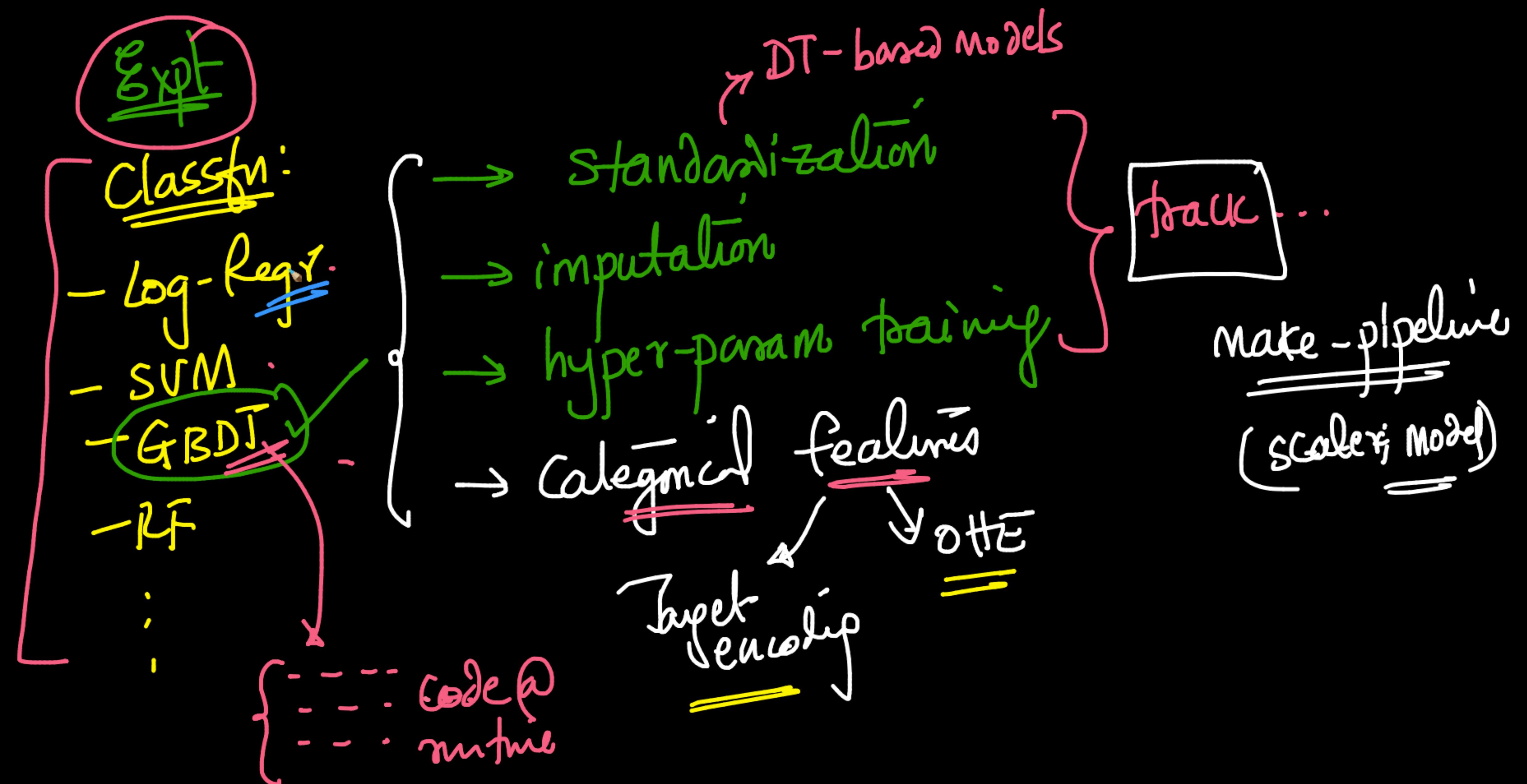
=

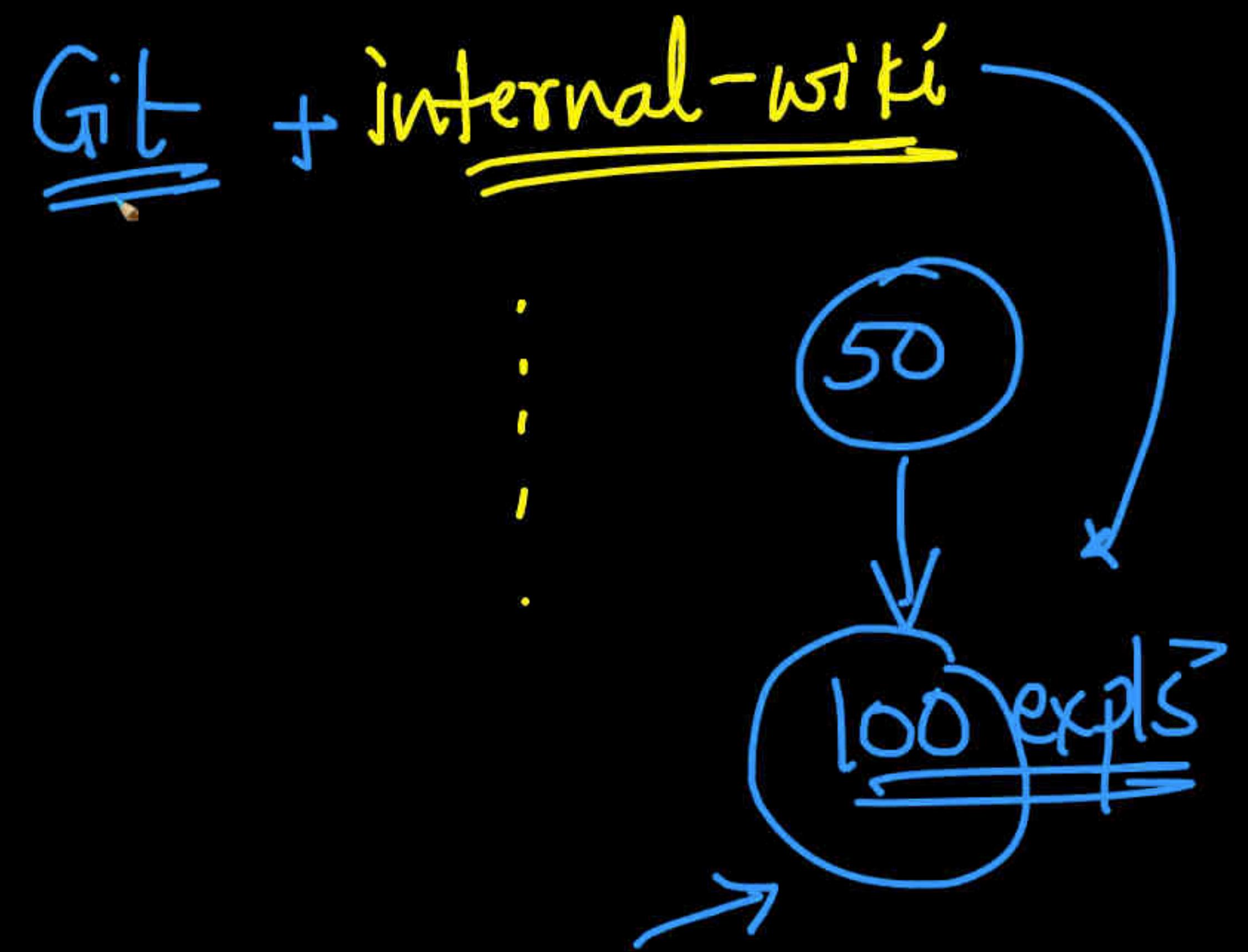




See the trees

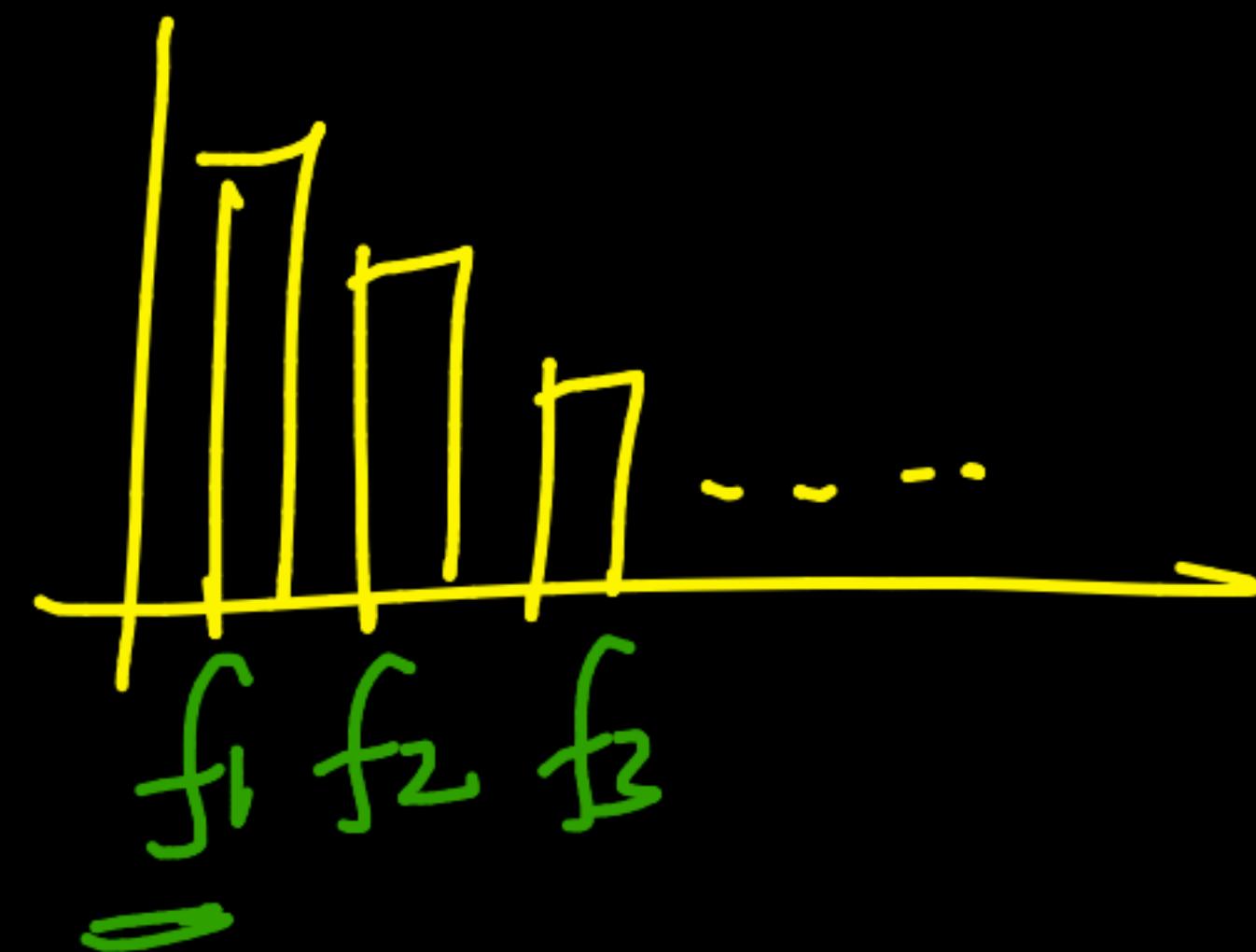






$f_1$

+ve-values



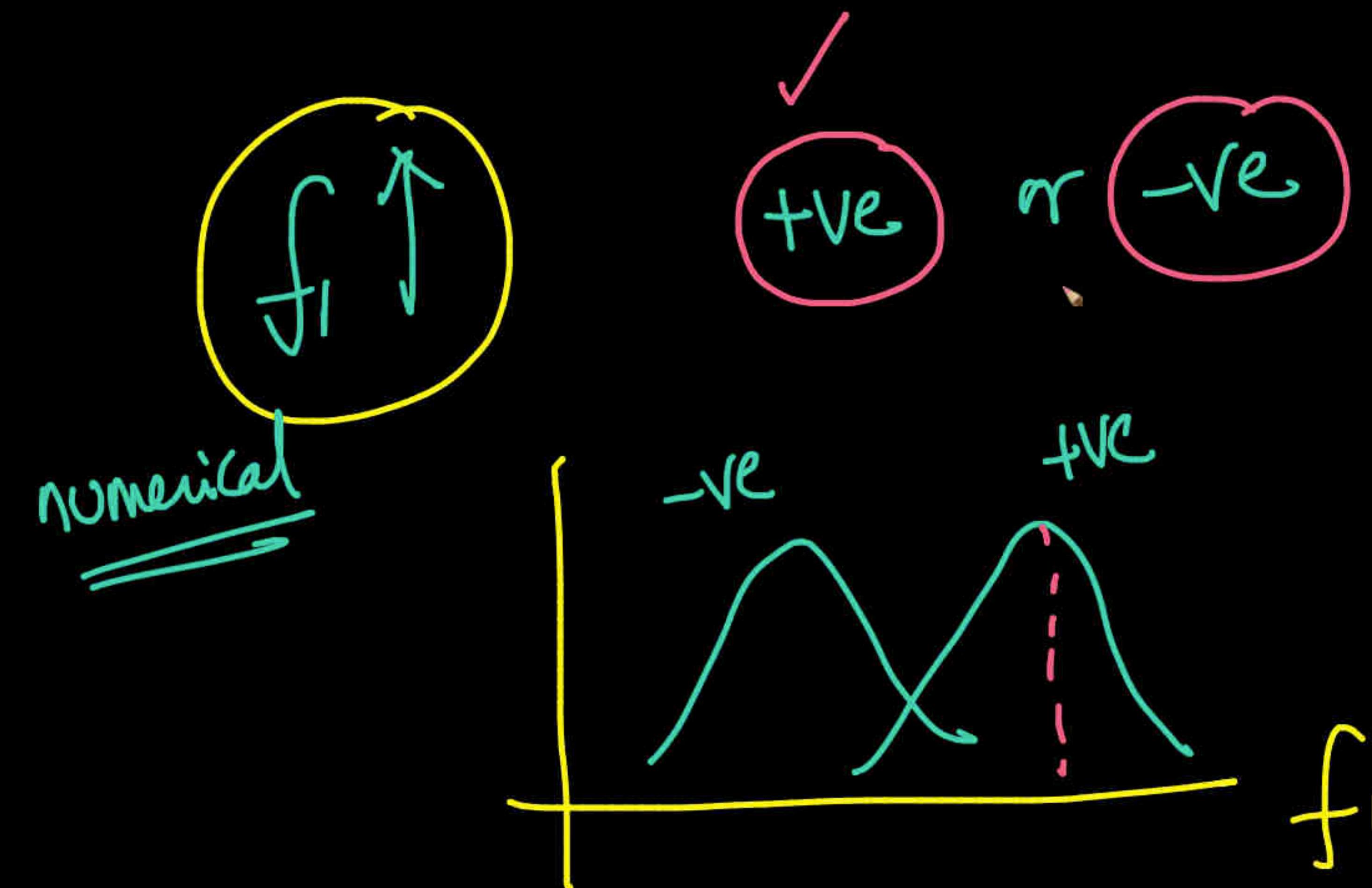
+ve  $\rightarrow$  1 Cancer

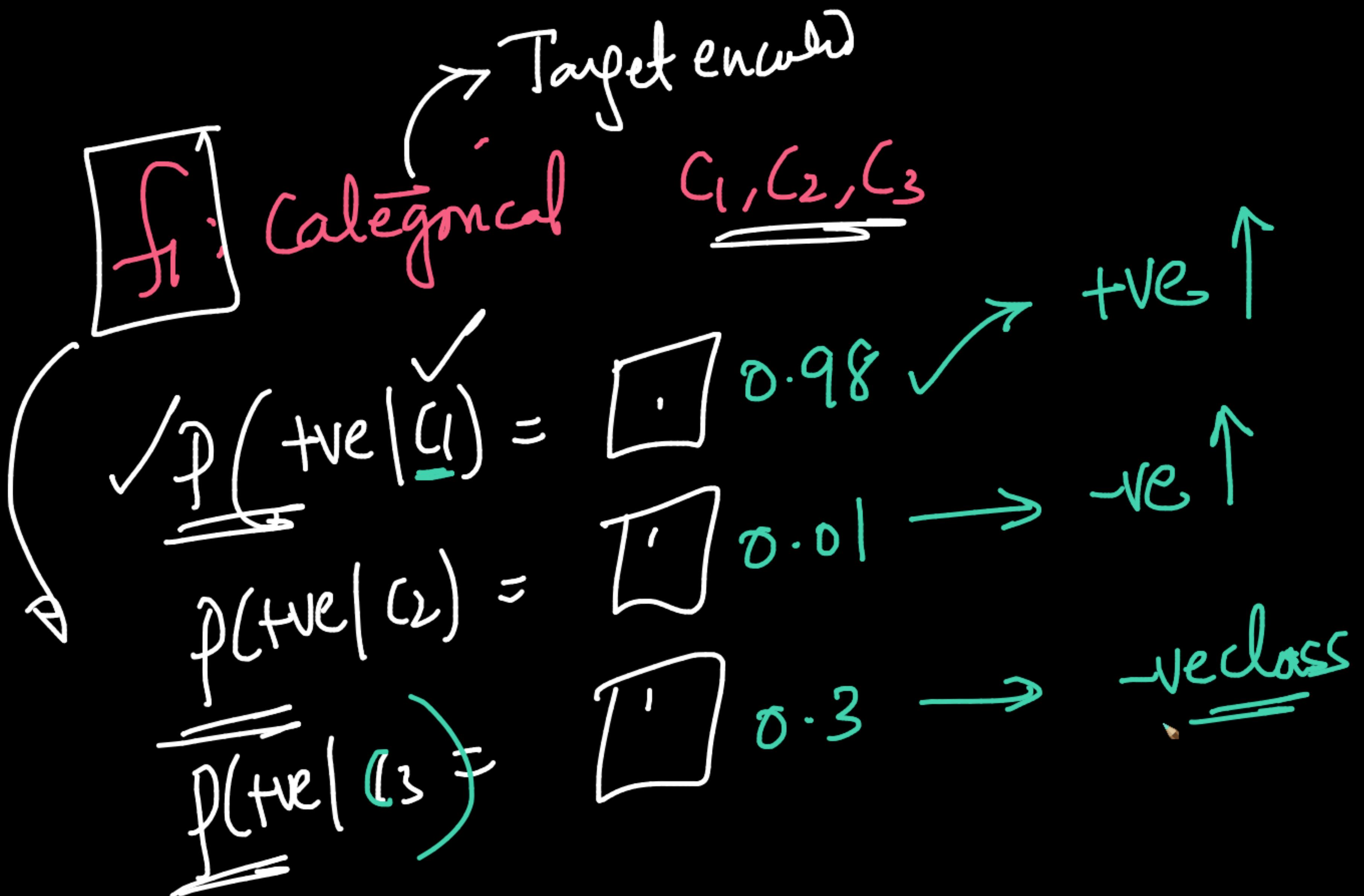
-ve  $\rightarrow$  0 not-cancer

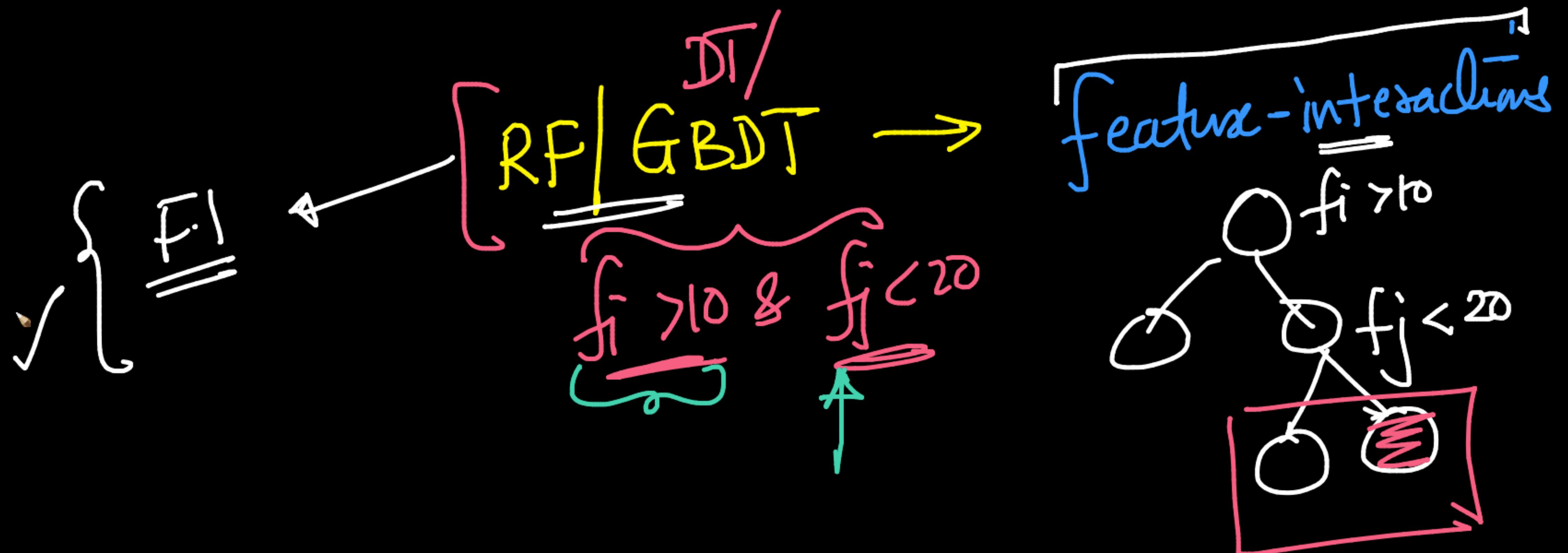
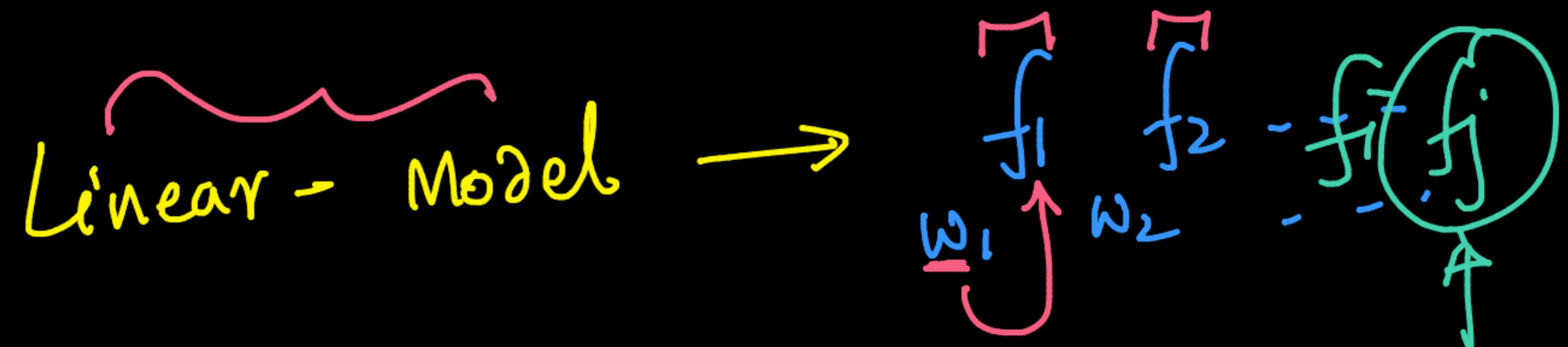
$f_i$   $\rightarrow$  +ve or -ve

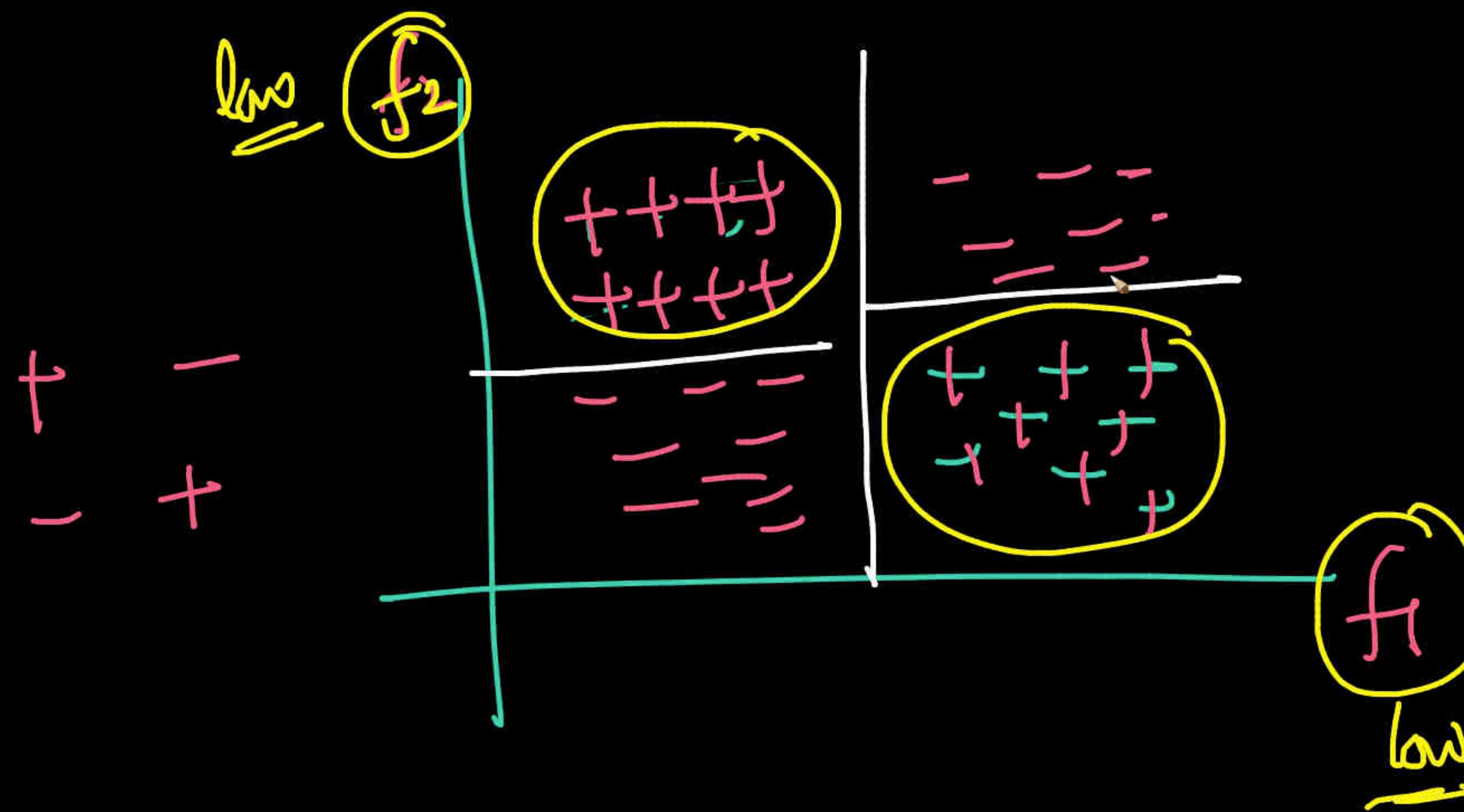
{IsSmoking}

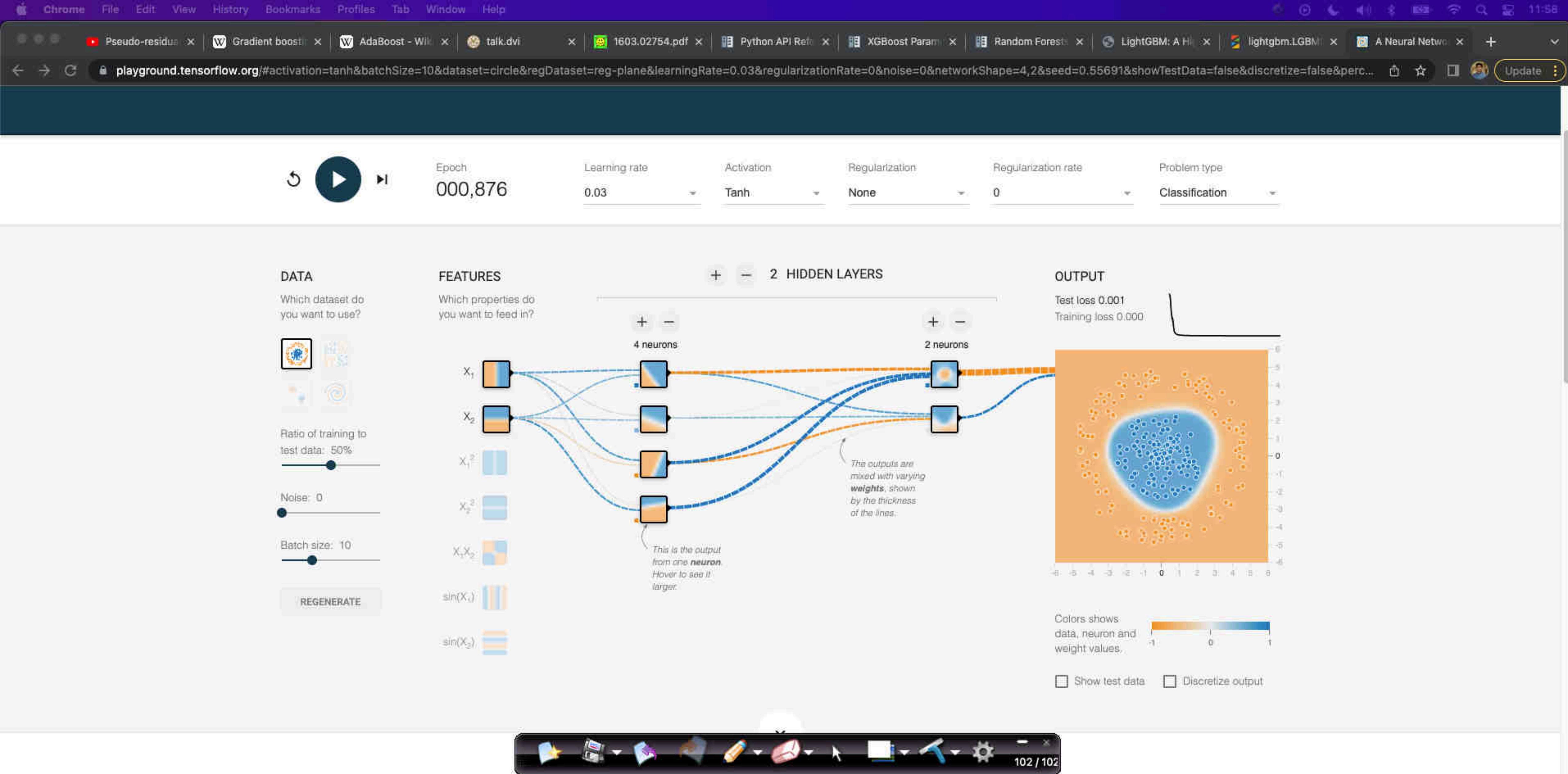
+ve weight











It's a technique for building a computer program that learns from data. It is based very loosely on how we think the human brain works. First, a