

Recap

Ensembles:

- ① Bagging - (Random forest) ✓
- ⇒ ② Boosting → Base learners + Aggregation
- ③ Stacking
- ④ Cascading

→ Bagging (Random forest) -

- Hyperparameters ✓
- Hyperparameter tuning.

Agenda

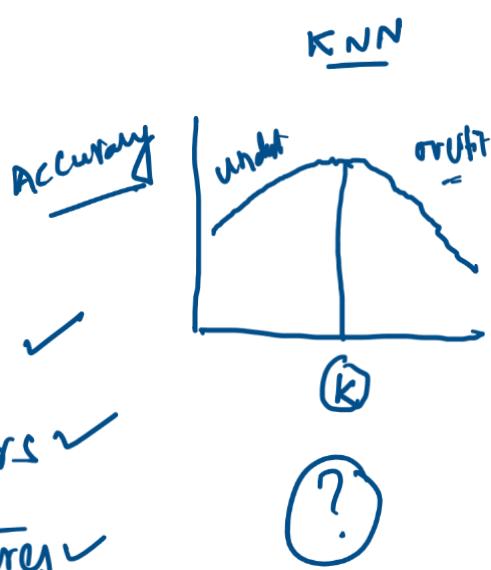
Boosting

Random forest :-

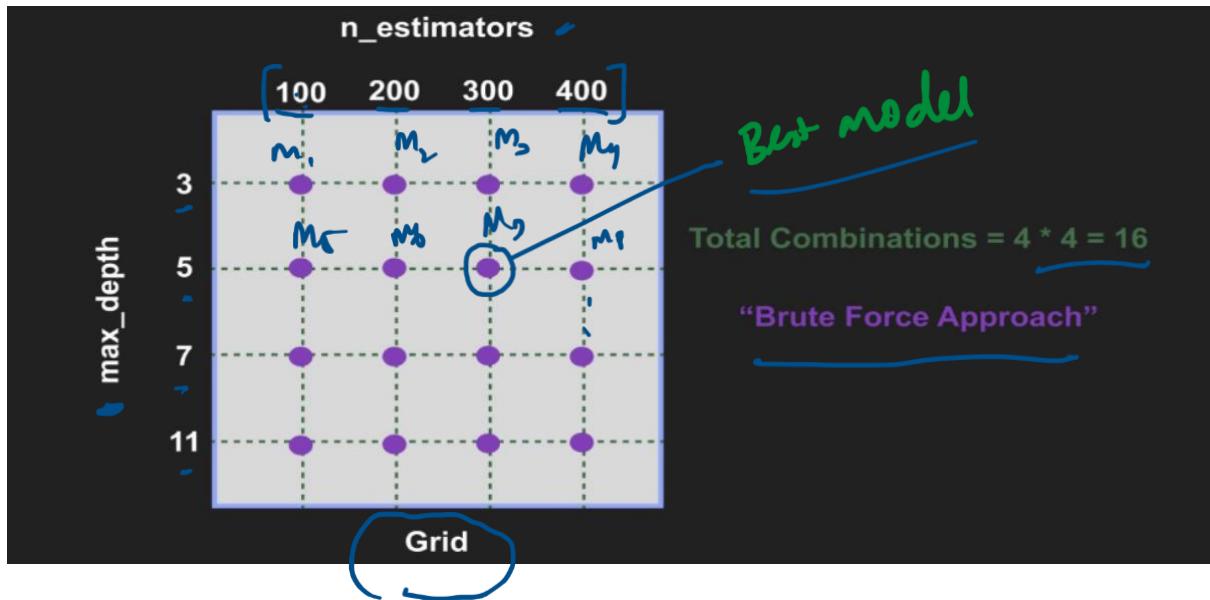
Tree parameters

- ① max no. of leafnodes
- ② class wt - - -

- max-depth ✓
- n-estimators ✓
- max-features ✓
- + tree parameters)



Hyper Parameters



$$n_{\text{estimators}} - [100 \quad 200 \quad 300 \quad 400] \quad (4)$$

$$\text{max_depth} - [3 \quad 5 \quad 7 \quad 9] \quad (4)$$

(4 × 16)

	max depth	max dep	mat dep
n estimators	100, 3	100, 5	100, 7 -

$$n_{\text{estimators}} - [100 \quad 200 \quad 300 \quad 400 \quad 500] \quad 5$$

$$\text{max_depth} - [5 \quad 6 \quad 7] \quad 3$$

$$\text{Max_feature} - [4 \quad 5 \quad 6 \quad 7] \quad 4$$

60

60 different models

	n-est	max-depth	max-features	Accuracy
M_1	- 100,	4	, 5	0.81
M_2	- 100	4	6	0.82
1

3 hyperparameters $a \times b \times c$ Grid.

Grid Search + CV

GridSearchCV

hyperparameter tuning

+ cross validation

→ Randomised Search

RandomSearchCV

$n - (5)$ count

Max - (5.)

max-f (5.)

1200*

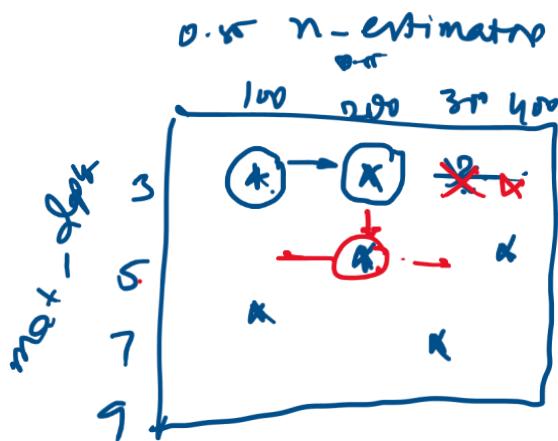
cv - 10

n-estimators

100 200 300 400

CCP-alpha -

0.1 — 0.5 increment by 0.05



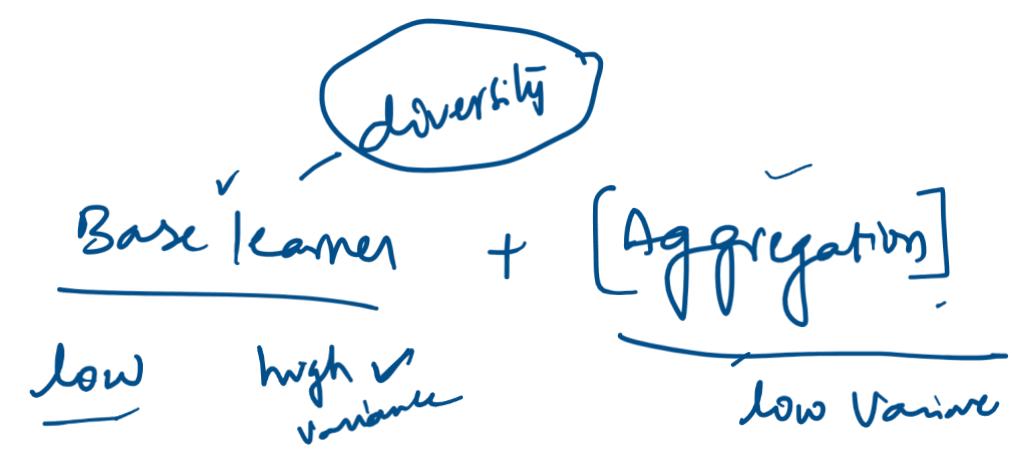
Grid.

- Reduces Sparseness ✓
- ↪ faster computation

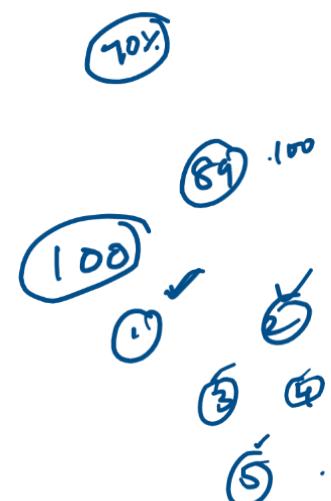
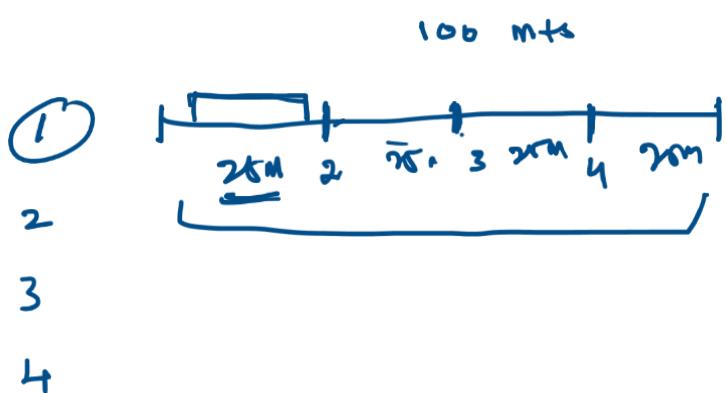
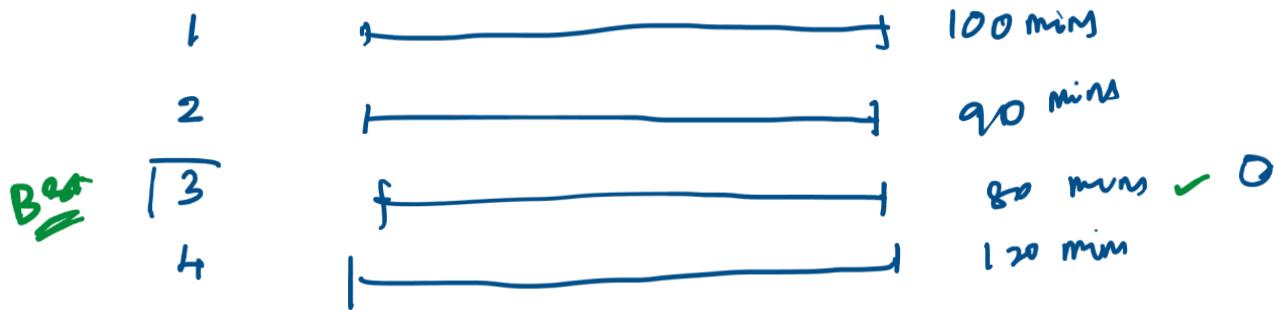
150
n estimators [100 200 300]
(100 125 150 175 200 250)

Boosting

Bagging -



- (Row sampling + col sampling) *



Boosting

Base learners (weak learner)

High bias, low variance.

Boosting

Base learners + additive
(weak learners)
high bias

x_1 Height	x_2 Gender	y Weight	\bar{y}	$\frac{60}{=}$
$\frac{1.6}{1.5}$	M	$\frac{82}{55}$	$\frac{82-60}{55-60}$	$\frac{y^{(i)}}{x^{(i)}}$
F	!	!	$\left[\begin{matrix} x^{(i)} \\ y^{(i)} \end{matrix} \right]$	$y^{(i)}$
.	.	.	$[n \times d]$	$y^{(i)}$

Regression on $[x^{(i)}] \times \in \mathbb{R}^d$ $[y^{(i)}] \in \mathbb{R}^1$
 std.

$$y^{(i)} \rightarrow \bar{y}$$

Base model

$$\hat{y} = \bar{y} \quad (\text{Base model})$$

Mean model

high Bias



$$M^0 - \hat{y} = \text{mean } (\bar{y})$$

$$M \rightarrow f(x) \quad \underline{M^0 - \text{mean } (\bar{y})} \checkmark$$

$$\hat{y} = \bar{y} \quad (\text{mean model})$$

$$M^1 \rightarrow \text{err}^0 = \hat{y}^{(i)} - \bar{y}$$

$$\hat{y}^{(i)} = M^0(\bar{y}) + \text{err}^0$$

$$82 = 60 + 22 \cdot ? \downarrow$$

hypotherm
(x)

$$M^0 - \underline{h_0(x)} - \bar{y}$$

$$\underline{M'} \quad \left\{ \underline{x_i^{(r)}}, \underline{\text{err}^0} \right\}$$

weights

$$\underline{82} = \underline{\text{avg}}_{60} + \underline{22} \cdot ?$$

$$M' - \underline{h_1(x)}$$

additive combining.

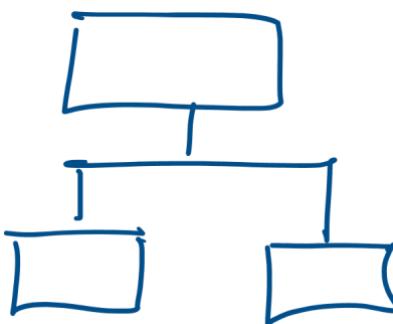
Mean model M_1

$$M = h_0(x) + \underline{h_1(x)} + \text{err}$$

weak model

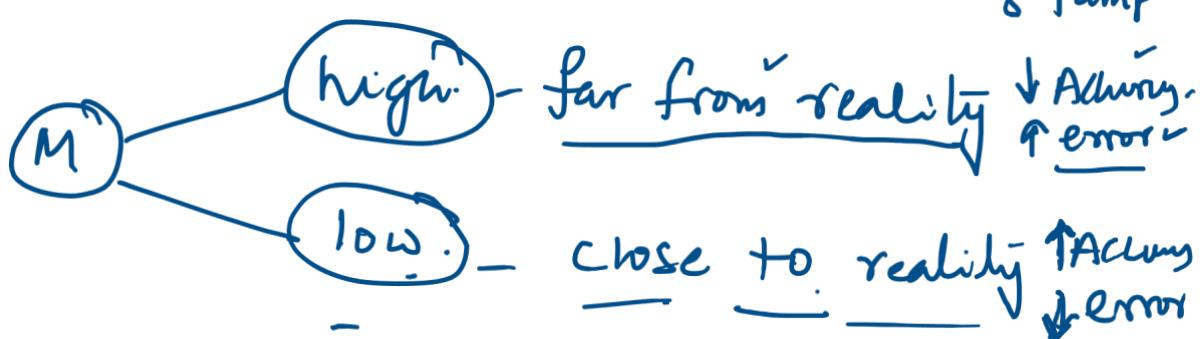
Tree Shallow depth - max-depth = 2

max-depth = 1 ✓



One Split tree

Decision stump

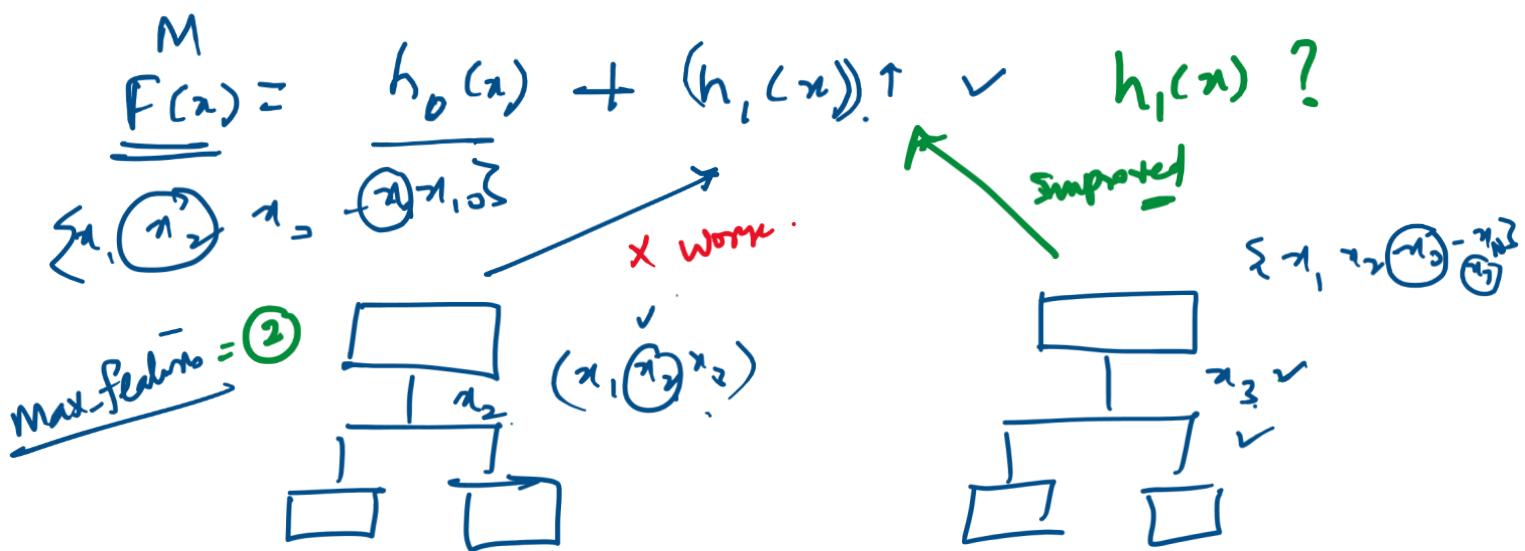


$$\underline{M}' \quad \underline{h_1(x)}$$

$$\left\{ x^{(i)}, e^{(i)} \right\}$$

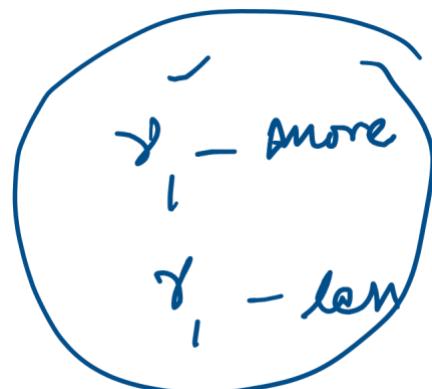
$h_1(x)$ high bias -
 by
 Baseline

Decision Stumps
 Shallow trees



Regularization addition

$$h_0(x) + \gamma_1 h_1(x)$$



γ - hyper Parameters - Learn Parameters

$$M f(x) = h_0(x) + \gamma_1 h_1(x) + \gamma_2 h_2(x) + \gamma_3 h_3(x) + \dots + \gamma_m h_m(x)$$

learn combine $\gamma_1 \quad \gamma_2 \quad \gamma_3 \dots \gamma_m$

$$\text{err}^h = F_{h-1}(x) - \underline{\text{err}}^{h-1}$$

$h-1$ - models. $F_{h-1} = h_0(x) + \gamma_1 h_1(x) + \gamma_2 h_2(x)$
 $+ \dots \gamma_{h-1} h_{h-1}(x)$

$F_K(x)$ = $h_0(x) + \gamma_1 h_1(x) + \gamma_2 h_2(x) + \dots + \gamma_K h_K(x)$

$$h_{K+1}(x) \rightarrow \left\{ \underline{x}^{(i)}, \underline{\text{err}}^K \right\}$$

$$\boxed{\text{err}_i^k = \underline{y}^{(i)} - F_K(\underline{x})} \quad \checkmark$$

$$\boxed{F_K(x)}$$

Gradient descent?

$$\boxed{\gamma_1 \quad \gamma_2}$$

$$F_2(x) = h_0(x) + \gamma_1 h_1(x) + \gamma_2 h_2(x)$$

$$\underline{\text{err}} = \underline{y}^{(i)} - F_2(x) \quad \text{small } 1$$

Regression loss fn - "MSE"

$$F_m(x) = h_0(x) + \sum_{i=1}^K \gamma_i h_i(x)$$

K -hyperparameters γ_i - Parameter

→ $h_2(x)$ → without $h_1(x)$

$$h_2(x) \quad \{x^{(i)}, \text{err}^{(i)}\}$$

→ Boosting, additive., Sequential learning

Combining, Learn weights that optimizes the error.

$$\underbrace{F_K(x)}_{=} = h_0(x) + \gamma_1 h_1(x) + \gamma_2 h_2(x) \\ \dots + \gamma_K h_K(x)$$

$$\text{err}^k = (y^{(i)} - p_k(x)) \quad \downarrow$$

GB DT [Gradient Boosting Decision Trees]

$$F_K(x) = h_0(x) + \bar{\delta}_1 h_1(x) + \dots + \bar{\delta}_K h_K(x)$$

$\boxed{\bar{\delta}_1 - \dots - \bar{\delta}_K}$

Loss function MSE

$$L(\hat{y}^{(i)}, \hat{y}^{(i)}) = \sum_{i=1}^n (\hat{y}^{(i)} - \hat{y}^{(i)})^2$$

$$\boxed{\frac{\partial L}{\partial \hat{y}} = -2 (\hat{y}^{(i)} - \hat{y}^{(i)})}$$

$$\frac{-\partial L}{\partial \hat{y}} = 2 (\hat{y}^{(i)} - \hat{y}^{(i)})$$

error
Residual

-ve gradient \downarrow \sim Residual \downarrow

$$\gamma' = \gamma_0 - \eta \cdot \frac{\partial L}{\partial \gamma}$$

-ve gradient \downarrow \sim error \downarrow

gradient $\rightarrow 0$ ✓ gradients

pseudo residual \rightarrow

-ve gradient ✓

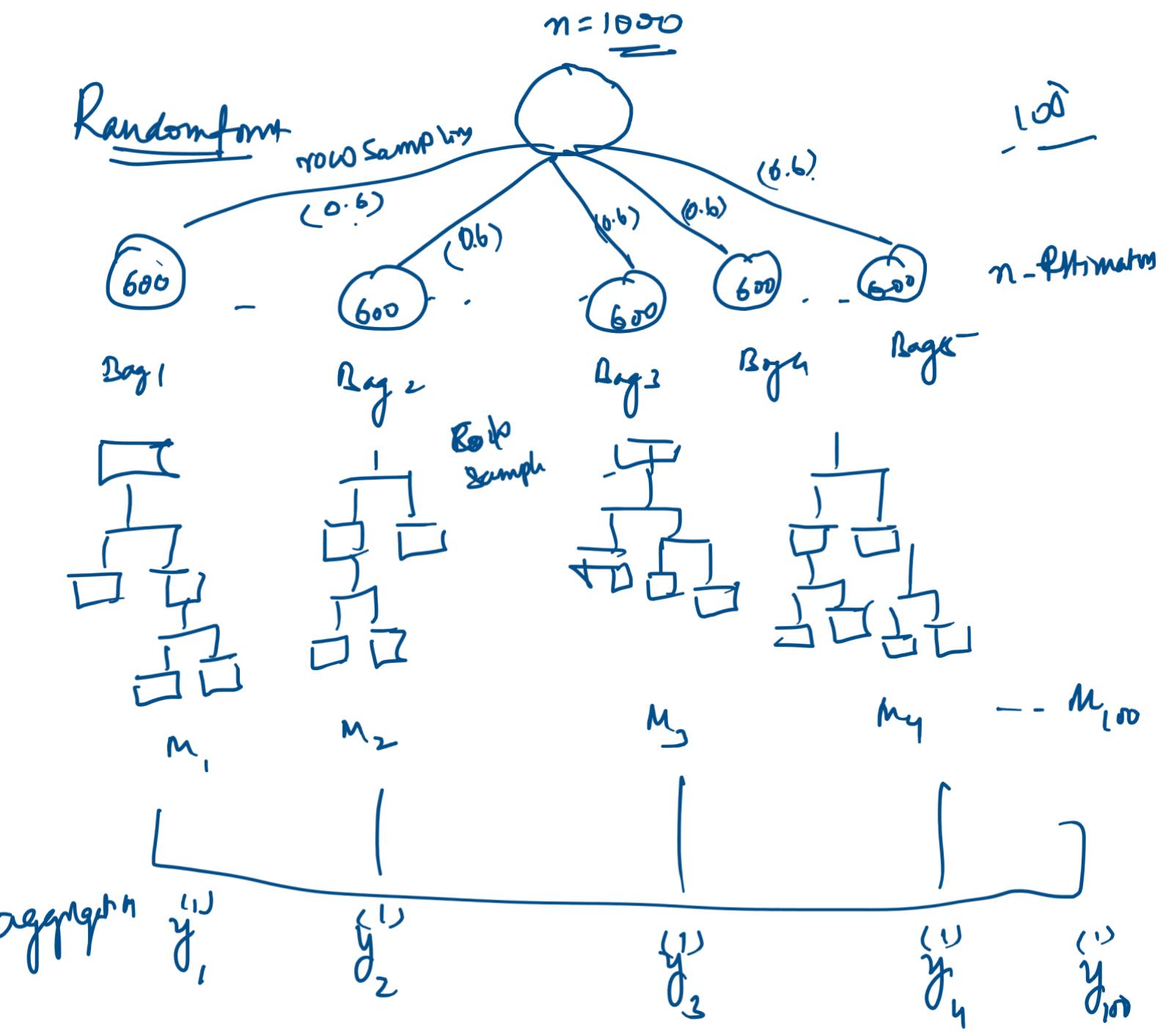
loss function - differentiable. -ve gradient

Custom loss function - differentiable ✓

minimize - pseudo-error, -ve gradients

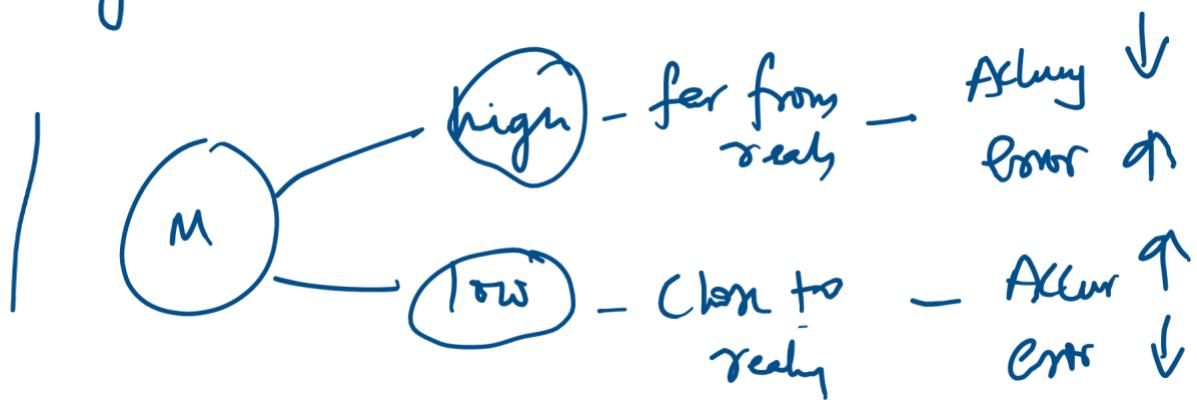
Regression - MSE -

Classification - log loss -



Regm angle $\hat{y}_i^{(1)}$

Claim majority $(\hat{y}_i^{(1)})$ — ✓



$$\underline{n_jobs = 1}$$



CPU 4 core, 8 threads

$$n_Jobs = 4$$

