Collection of Base Models ← 

Ensemble  Models

I/P → [ $M_1$ $M_2$ $M_3$ ] → O/p -

↳ A Collection of things

Ensemble

$f_1$ $f_2$ $f_3$ ... $f_d$  $y$

1
2
3
.
.
$n = 1000$

$f_1$ $f_2$ · $f_d$
$D_n$

1. Bagging ⎫
2. Boosting ⎬ → Homogeneous
3. Stacking — Heterogeneous

1. Bagging — Reduces Variance
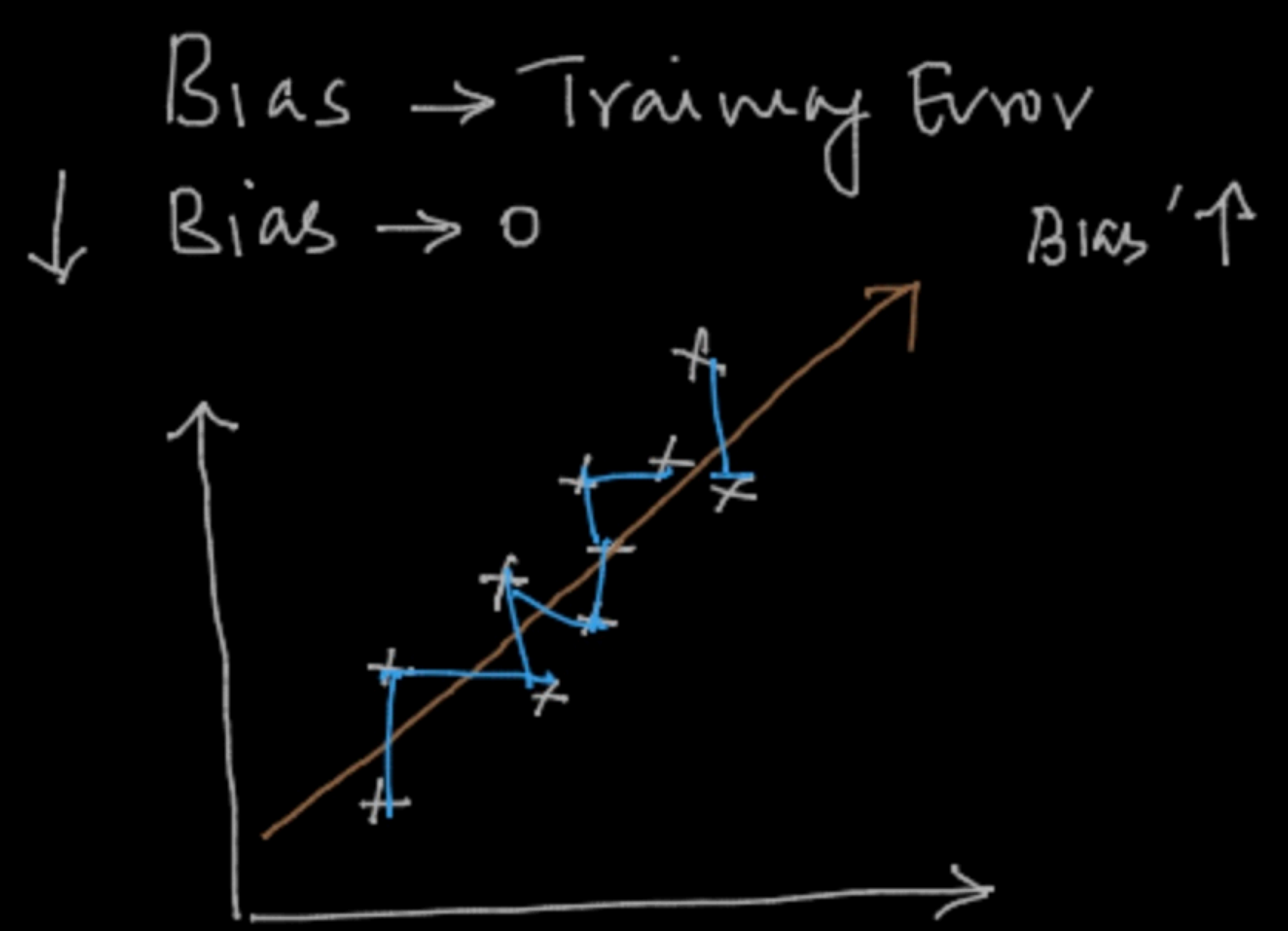
↳ Base Models
High Variance + Low Bias

Bootstrapping — Sampling with replacement

Bootstrapped Aggregation -

Bagging

Bias & Variance

$D_n$ →

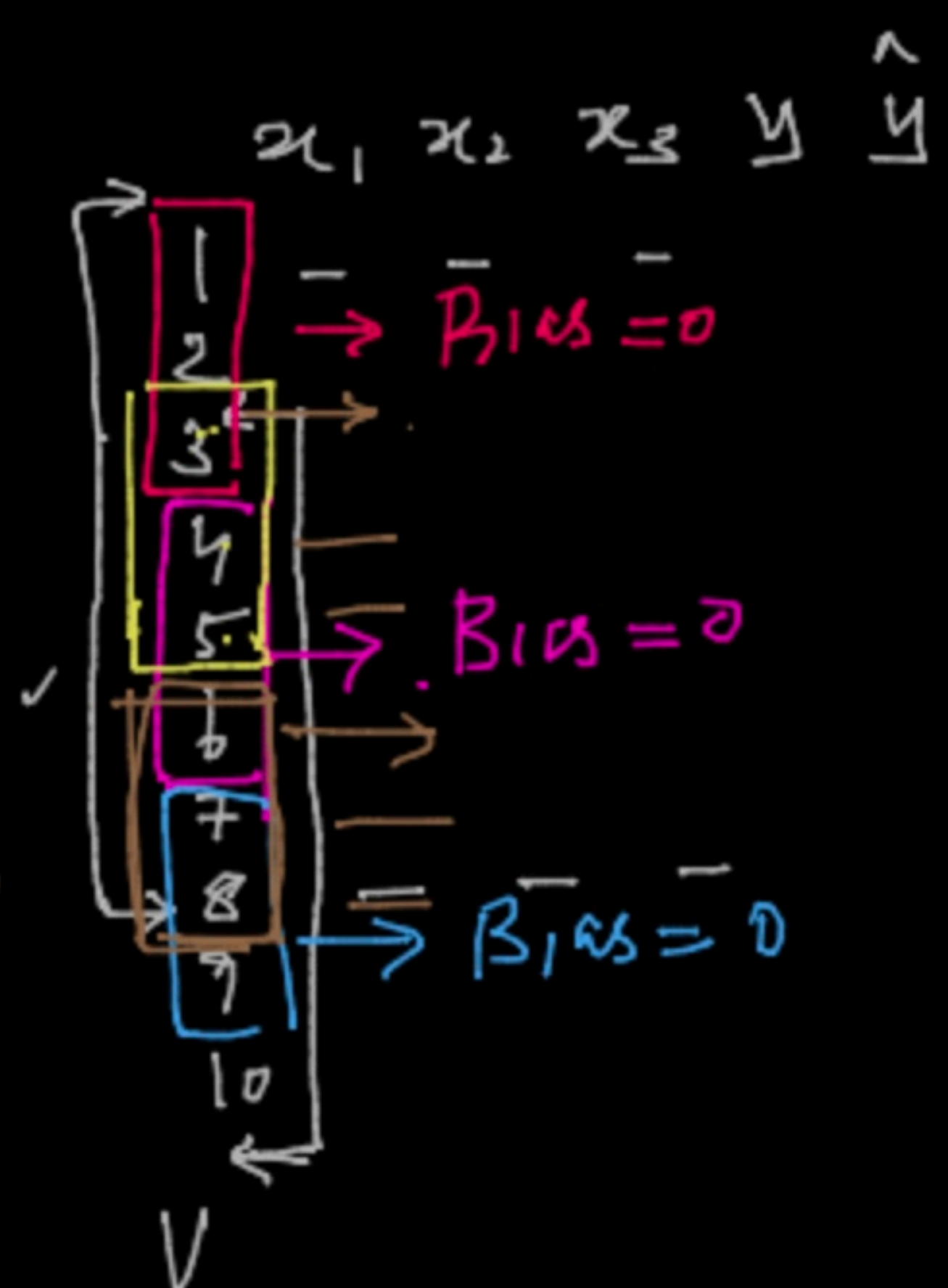$n_1 = 200$
[ $D_1$ ] — [ $M_1$ ] $\begin{smallmatrix} O_1 \\ D \end{smallmatrix}$

LB4V 35

$n_2 = 200$
[ $D_2$ ] — [ $M_2$ ] $\begin{smallmatrix} O_2 \\ c \end{smallmatrix}$

LBHV 43

$n_3 = 300$
[ $D_3$ ] → [ $M_3$ ] $\begin{smallmatrix} O_3 \\ D \end{smallmatrix}$

LB4V 32

Bootstrapping

A g g r e g a t i o n

Regression
Avg / Median ⟶

Majority vote
classification

LBLU

[ $D_{mn}$ ] → [ $M_m$ ] $\begin{smallmatrix} O_m \\ D \end{smallmatrix}$

LBHV 41

Bias → Training Error

↓ Bias → 0　　　　　　　Bias' ↑



Variance — Testing Error.

↑ The model itself changes
　With small changes in data
Variance ↓

Bias ↓ Low Variance ↓

$x_1 \ x_2 \ x_3 \ y \ \hat{y}$

1
2　→ $\overline{Bias} = 0$
3
4
5 → $Bias = 0$
6
7
8
9 → $\overline{Bias} = 0$
10

V

Bagging Reduces Variance → Base Models –

Low Bias & High Variance

Homogeneous Ensembles

$M_1$

Bagging Regressor ( —　　)

10 ppl → Predicting House price –

1 → BTM ✓
2 → Indranagar
3 → HSR
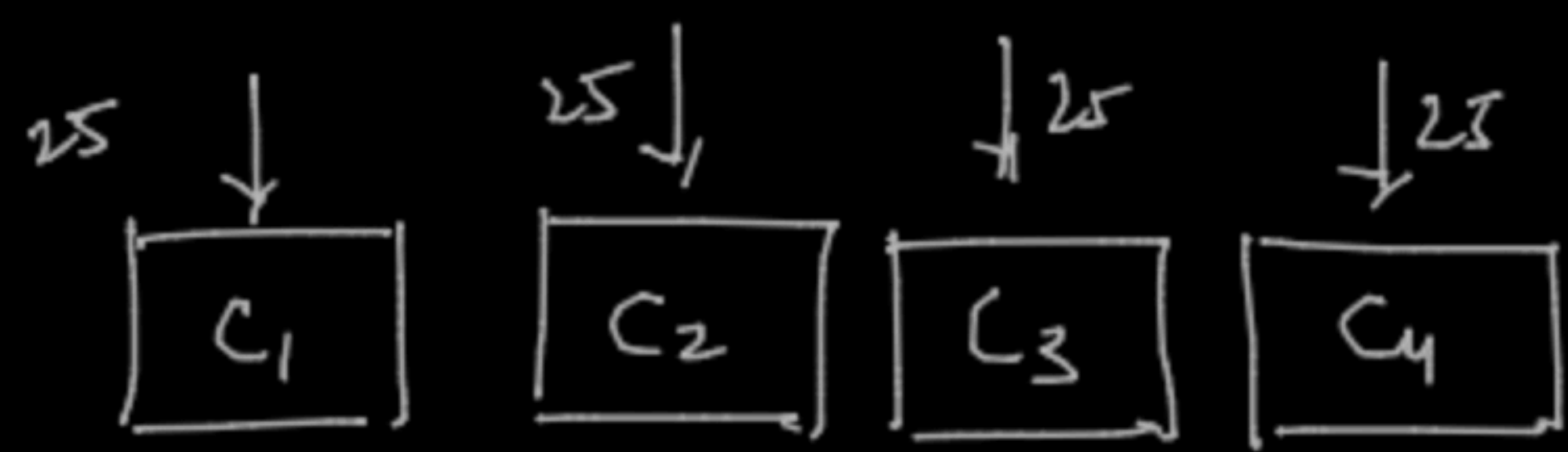4 → BTM
5 → Indranagar
·

90% →

100
180
70
80
+
÷

90%

# Random Forest → Bagging Technique

1. All Base models are Decision Trees ✓
2. Row sampling + Feature sampling
   ↓
   └→ only a subset of Features

$D_n \longrightarrow D_{n1=200}$

$f_1 f_2 f_3 \cdots f_d$
$D_n$

$(200, 3, 5, \cdots)$

**D7**
$\boxed{\begin{array}{c} f_1 f_3 f_3 \\ D_{n1} \end{array}}$ → $O_1$ $M_1$ ✓

**DT**
$\boxed{\begin{array}{c} f_1 f_2 f_4 \\ D_{n2} \end{array}}$ $O_2$ → $M_2$ ✓

**D7**
$\boxed{\begin{array}{c} f_1 f_3 f_5 \\ D_{n3} \end{array}}$ $O_3$ → $M_3$ →

**D7**
$f_1 f_2 f_5$
$\boxed{D_{n_{100}}}$ → $M_{100}$ →

Aggregation

$O_1$ ✓
$O_1$
$O/P$ →

$\to$ O/P

model = RandomForest( n-estimators = 100 )

model.fit( $X_n$ )

model.predict( )

$f_1 \ f_2 \ f_3 \ f_4 \ f_5 \cdots f_d \quad y$

**D7**

1
2
3
4
⋮
⋮

$n = 80\%$

ML Ops →

→ Data Drift

$\begin{bmatrix} \cdot \\ \end{bmatrix}$

Parallelize the model.

→ Quad Core    25%.

$\downarrow 25$      $\downarrow 25$      $\downarrow 25$      $\downarrow 25$

$\boxed{C_1}$   $\boxed{C_2}$   $\boxed{C_3}$   $\boxed{C_4}$

Only the first model predicts `y`
All subsequent models predict the errors from the previous

Boosting.

— Reduces Bias ✓
— Sequential Model.
— Additive Model

$$f_1 \ f_2 \ f_3 \quad y \quad \hat{y} \quad \varepsilon_0$$

$$x_i \leftarrow x_1 \rightarrow y_1 \quad \cdot$$
$$x_2 \leftarrow x_2 \rightarrow y_2 \quad \cdot$$
$$x_3$$
$$\vdots$$
$$x_n \leftarrow x_n \rightarrow y_n$$

Stage '0':

$$M_0$$

$$D_n \Rightarrow \{x_i, y_i\}_{i=1}^{n}$$

$$\hat{y} = f_0(x)$$

$$\varepsilon_0 = y - \hat{y}$$
$$= y - f_0(x)$$

Base Models

High Bias & Low Variance

LBLV.

Stage 1:

$$\{x_i, \varepsilon_{0i}\}_{i=1}^{n} \quad M_1$$

$$\hat{\varepsilon}_0 = f_1(x)$$

$$\varepsilon_1 = \varepsilon_0 - \hat{\varepsilon}_0$$
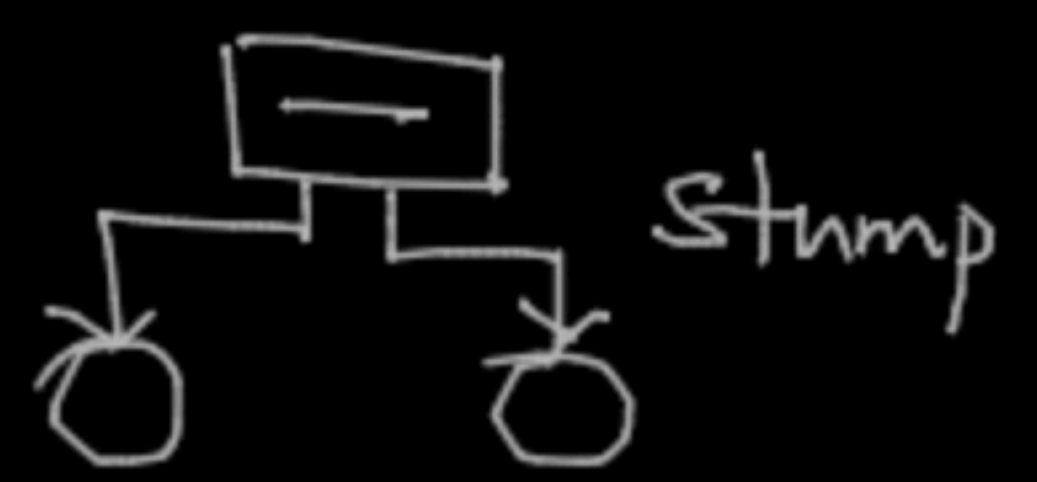$$= y - f_0(x) - f_1(x)$$

Stage 2:

$$\{x_i, \varepsilon_{1i}\}_{i=1}^{n} \quad M_2$$

$$\hat{\varepsilon}_1 = f_2(x)$$

$$\varepsilon_2 = \varepsilon_1 - \hat{\varepsilon}_1$$
$$= y - f_0(x) - f_1(x) - f_2(x)$$

$$\Downarrow$$

$$y = f_0(x) + f_1(x) + f_2(x) \longrightarrow \text{Additive}.$$

ADABOOST ✓
XhBOOST ✓

$$\begin{cases} M\_estimators = ? \\ Each \ Tree \rightarrow \end{cases}$$

Decision stump

 stump

Stage 0  $f_0(x)$

| Area | y | $\hat{y}$ | $\varepsilon_0$ |
|---|---|---|---|
| 1200 | 35 | 22 | 13 |
| 2300 | 53 | 34 | 19 |
| 1700 | 47 | 38 | 9 |
| 3000 | 96 | 81 | 15 |
| 3500 | 94 | 82 | 12 |

Stage 1  $f_1(x)$

| Area | $\varepsilon_0$ | $\hat{\varepsilon_0}$ | $\varepsilon_1$ |
|---|---|---|---|
| 1200 | 13 | 8 | 5 |
| 2300 | 19 | 10 | 9 |
| 1700 | 9 | 2 | 7 |
| 3000 | 15 | 6 | 9 |
| 3500 | 12 | 5 | 7 |

Stage 2  $f_2(x)$

| Area | $\varepsilon_1$ | $\hat{\varepsilon_1}$ | $\varepsilon_2$ |
|---|---|---|---|
| 1200 | 5 | 1 | 4 |
| 2300 | 9 | 3 | 6 |
| 1700 | 7 | 2 | 5 |
| 3000 | 9 | 3 | 6 |
| 3500 | 7 | 4 | 3 |

$-1$ $-0.8$

$$y = f_0(x) + f_1(x) + f_2(x) + \cdots$$

$$= 22 + 8 + 1 + -0.5$$

$$= 31$$

Classification

·· 36  35·2

$\rightarrow$ $\boxed{35}$  34·95

↑

$M_2 = 31·8$

↑ 0·5

$M_2 = 31$

↑ 1 ✓

$M_1 = 30$

↑ 8 ✓

Brought to 22

linear Model

''

Partition the feature space
into pure regions belonging
to each class ''

$f_2$

Feature space  $M_1$

x x
x x
x x   Decision Boundary
x x --
x-x x x

→ $f_1$

' Linearly Separable ' — linear surface

PL SL specus

PL < 4.5

Vini

No SL > 3.5 Yes

Versicolor   Setosa

SL
Plant
Ht.

Setosa

3.5

Virgin.

Versicolor

4.5

SL

fr

Decision Boundary

`Linearly Separable'

Linear Models → Logistic Regression

↓ Non Linear Decision Boundary