

Ensembles:

↳ group of "musicians"

ML: multiple models used together

$\{ \underline{M}_1, \underline{M}_2, \underline{M}_3, \dots, \underline{M}_K \} \rightarrow \text{base models}$



4 types:- { - Bagging (Boostrap Aggregation)

{ - Boosting
- Stacking
- Cascading

→ high-performing

→ v. powerful

→ Kaggle → most

→ very useful in real-world

~~Key aspect:~~

$M_1, M_2, M_3, \dots, M_K$

more different these models are,
the better you can combine them if — generic

{ Problem:- M_i :- expert

Bagging

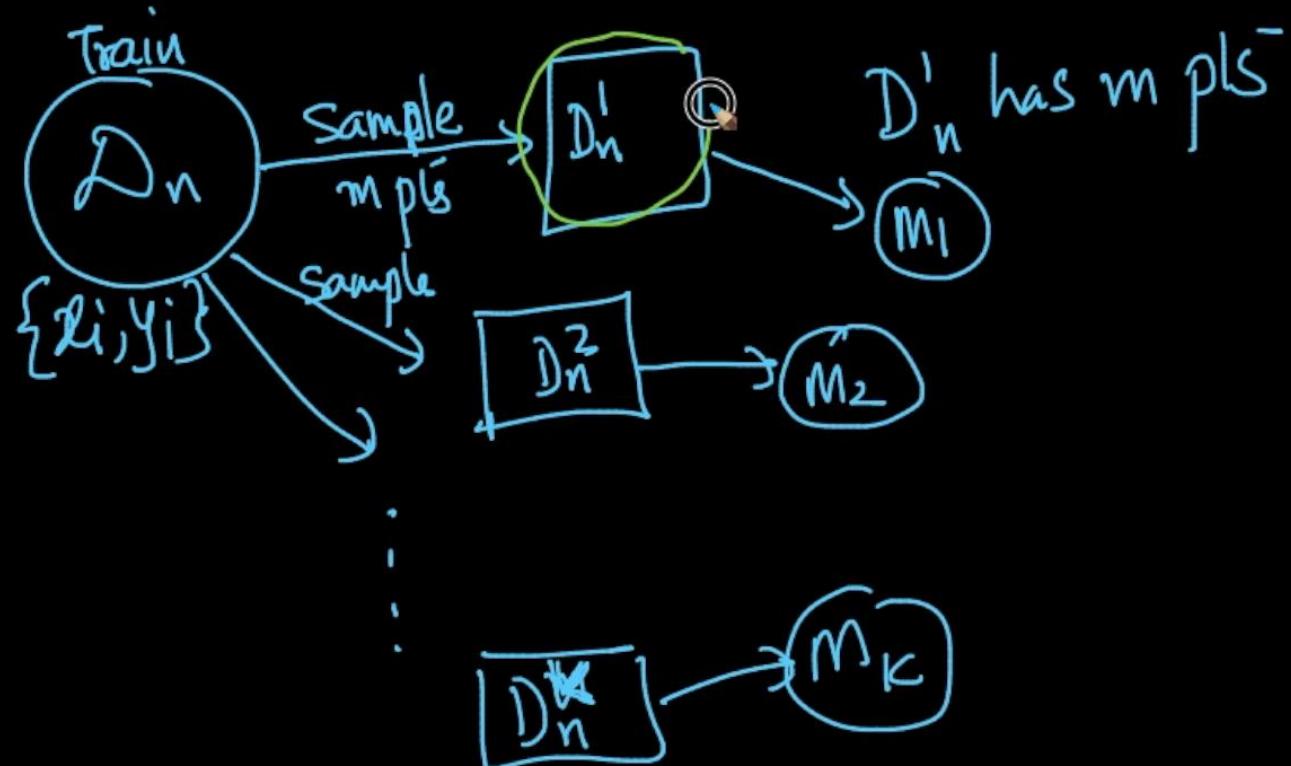
(Bootstrap Aggregation)

(Random Forest)

↳ Statistics

Intuition:

(Sampling with replacement)



Bagging

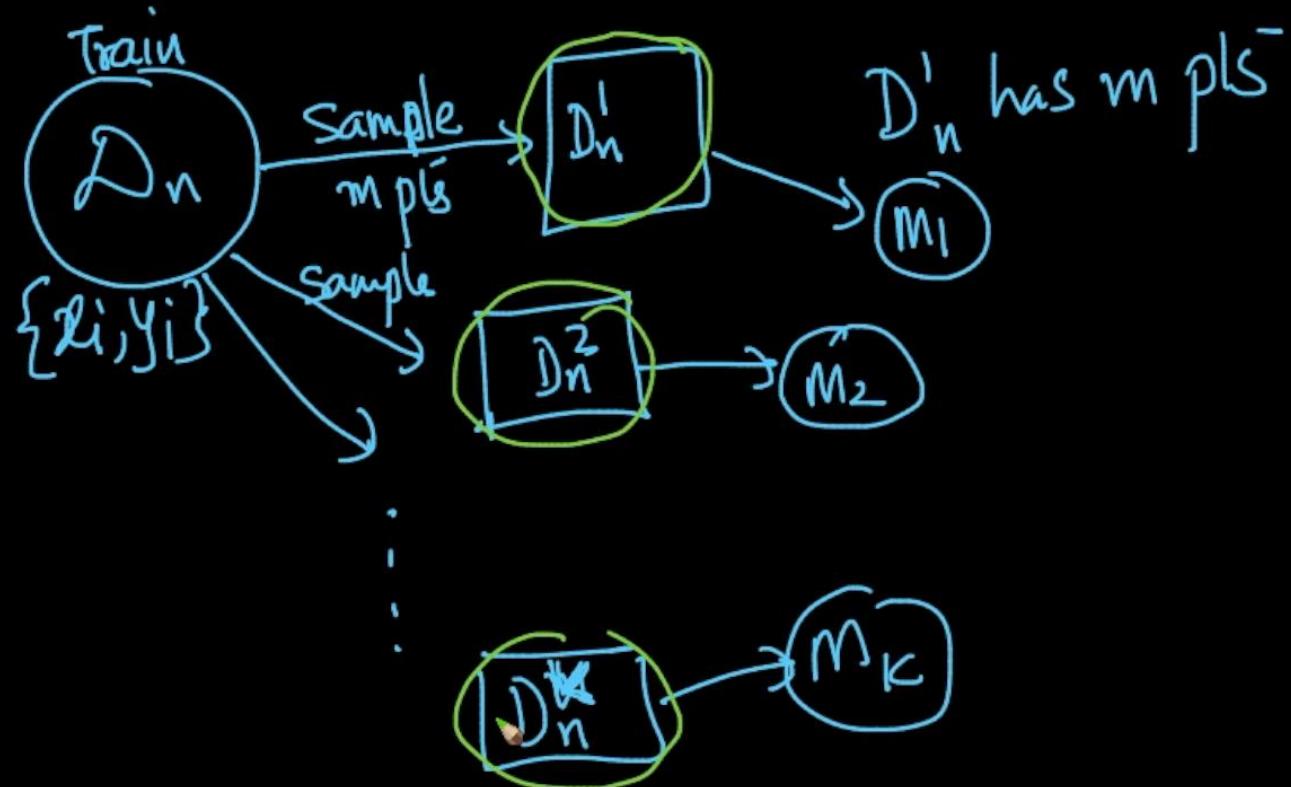
(Bootstrap Aggregation)

(Random Forest)

↳ Statistics

Intuition:

(Sampling with replacement)



Bagging

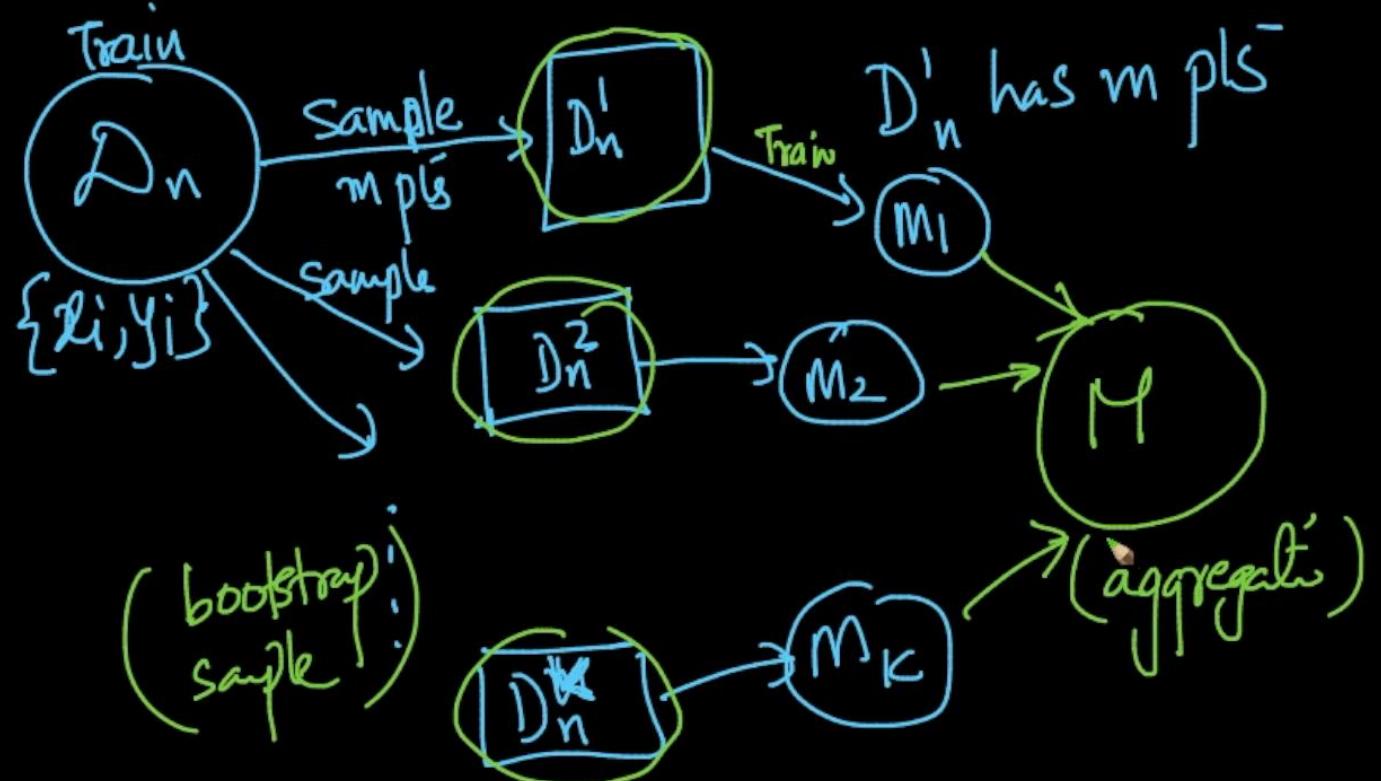
(Bootstrap Aggregation)

↳ Statistics

(Random Forest)

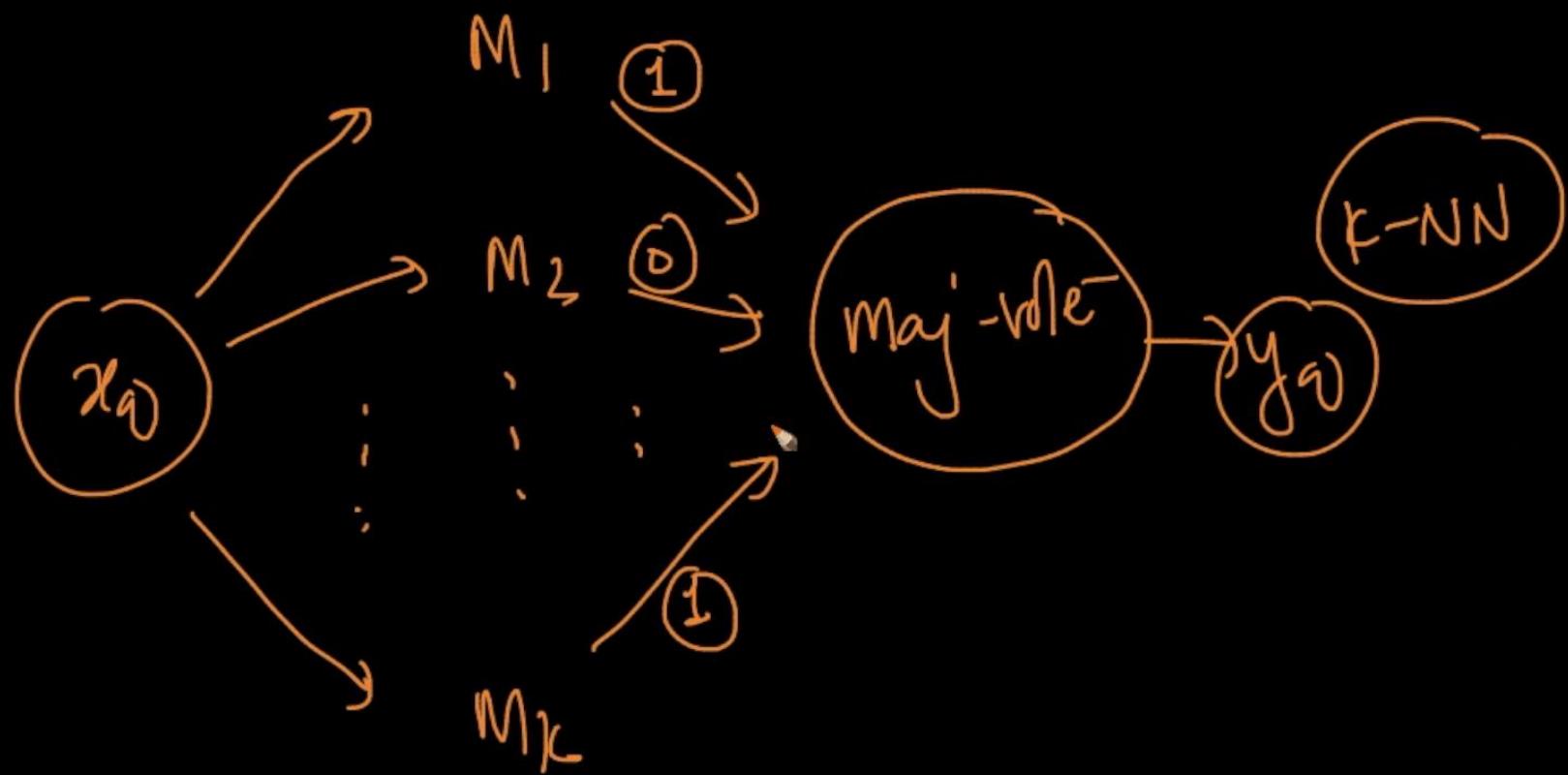
Intuition:

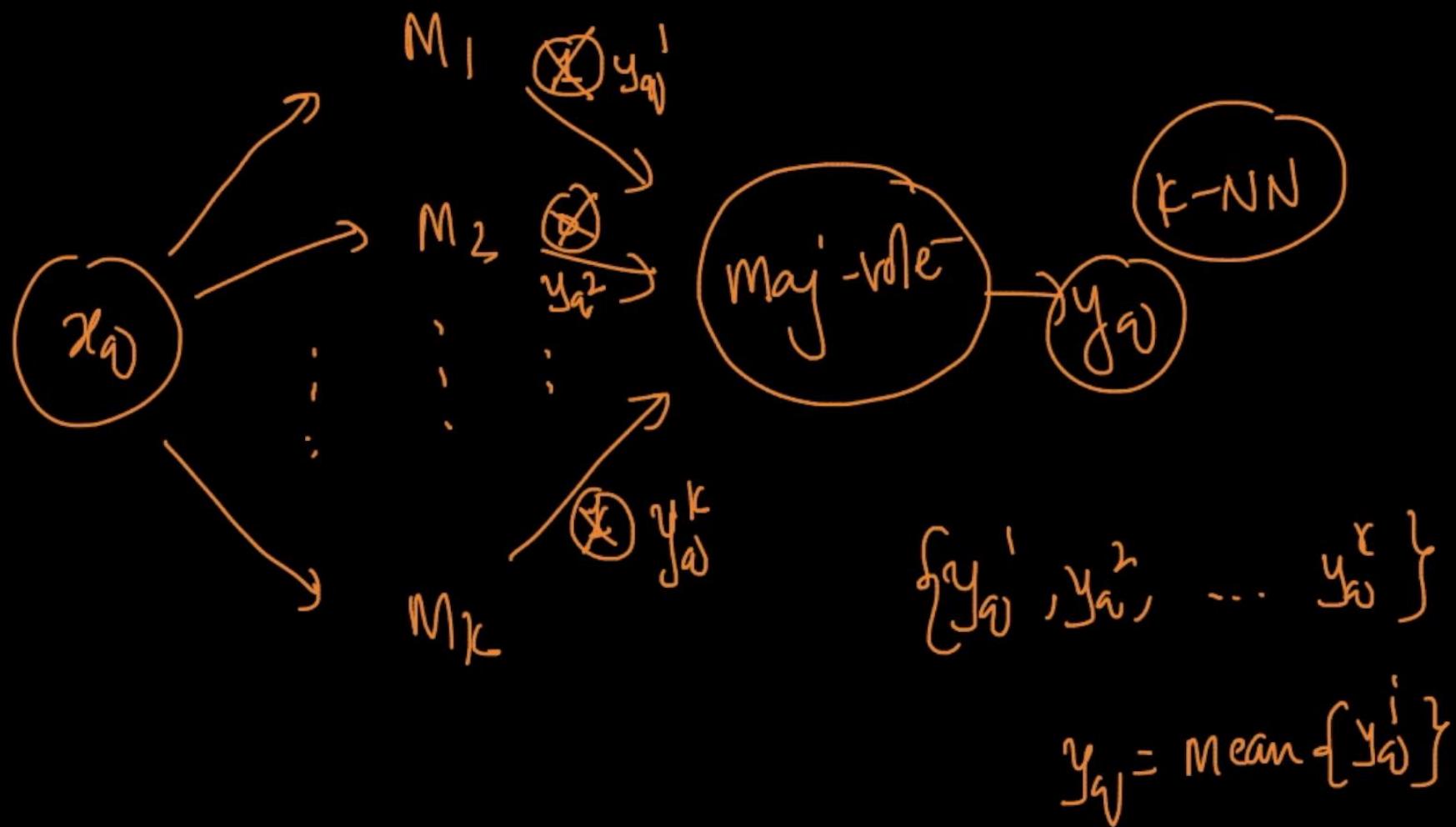
(Sampling with replacement)

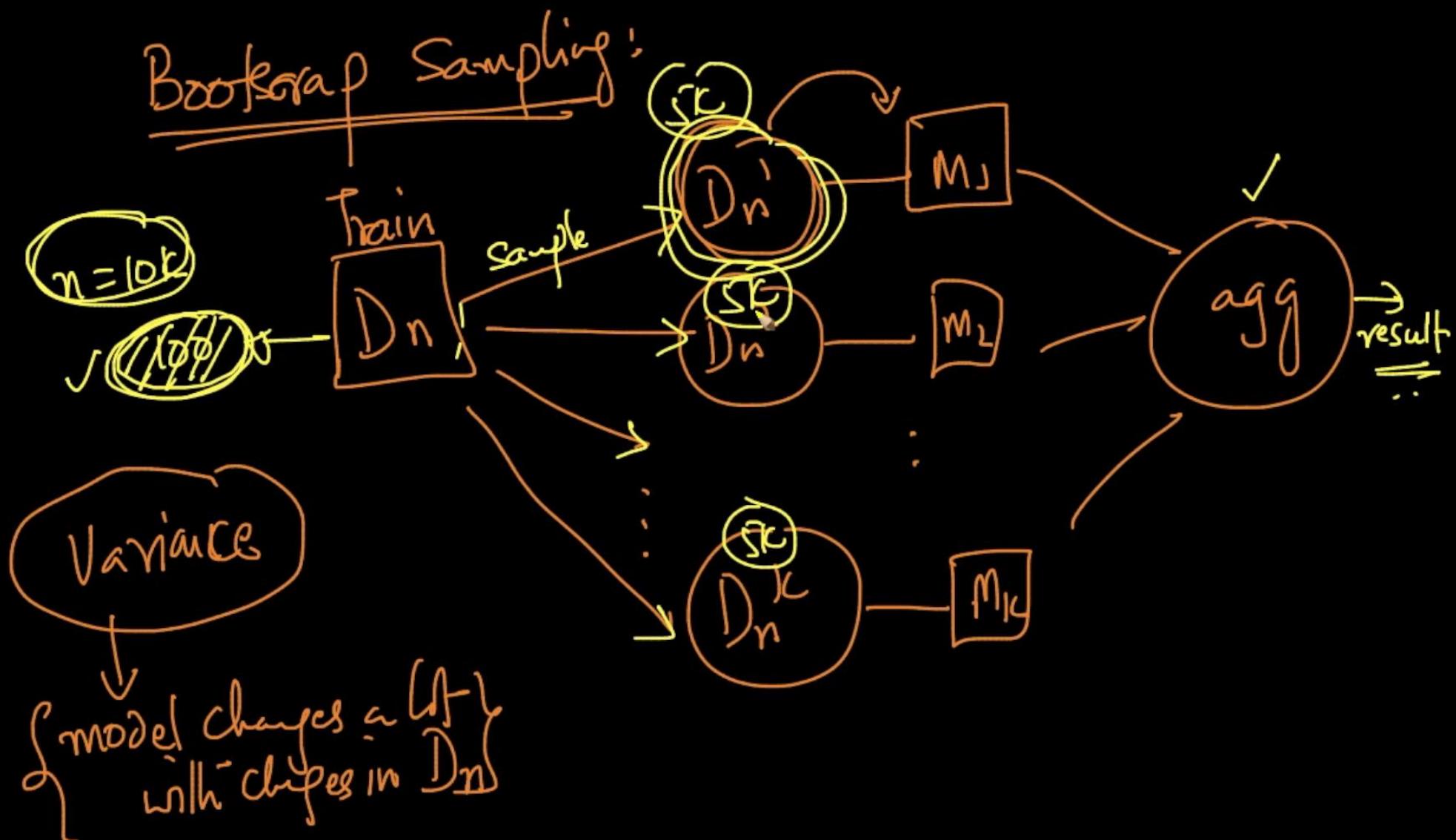


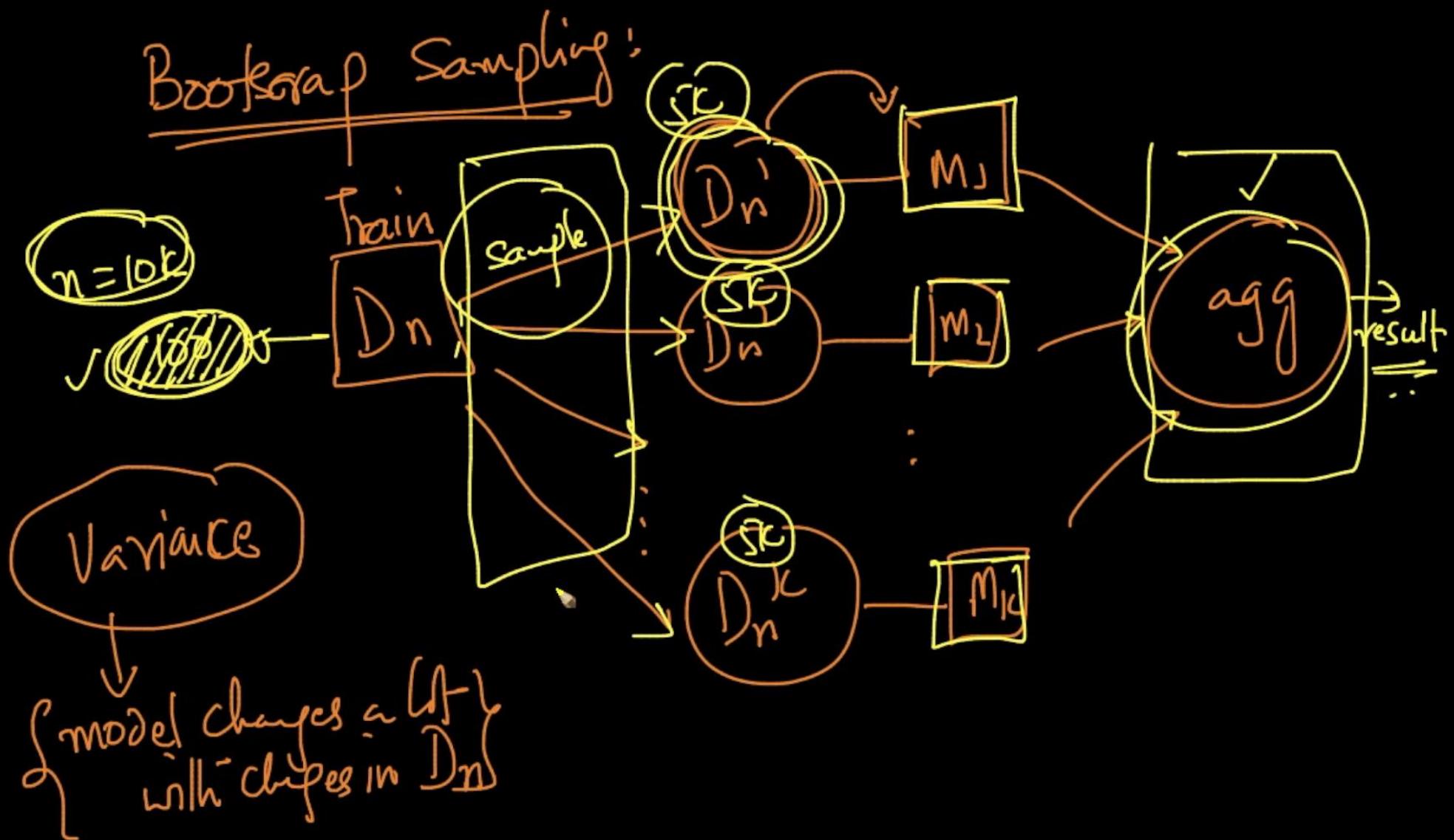
M_i is built using D_n^i of size m ($m \leq n$)
 \Rightarrow each model M_i has seen a different subset of
data

{ Aggregation :- Classification :- Majority rule
Regression :- mean / median









Bagging → can reduce variance in a model
without impacting the bias

$$\text{model error} = \underbrace{\text{Bias}^2}_{\text{low bias}} + \text{Var}$$

base-model (M_i): - low bias; $\underbrace{\text{high-var model}}$

* * * Bagging (M_i 's) → low-bias; reduced variance

Bagging:-

bunch of (low-bias, high-var) models

→ DT of depth (large)

high var
low-bias

RF

Combine them using bagging

(low-bias; reduced-var)

Bagging:-

bunch of (low-bias, high-var) models

↓
Combine them using bagging

→ DT of depth (large)

high var
low-bias

(low-bias; reduced-var)
=====

→ RF → most used & popular bagging algo

Random Forest (Bagging)

Decision Trees:-

↳ "Random" Bootstrap Sampling

RF : DT + Bagging + $\overbrace{\text{col. sampling}}$ \rightarrow feature bagging

↑
base

Random Forest (Bagging)

Decision Trees:-

"Random Bootstrap Sampling"

RF : DT

↑
base

+ Baging

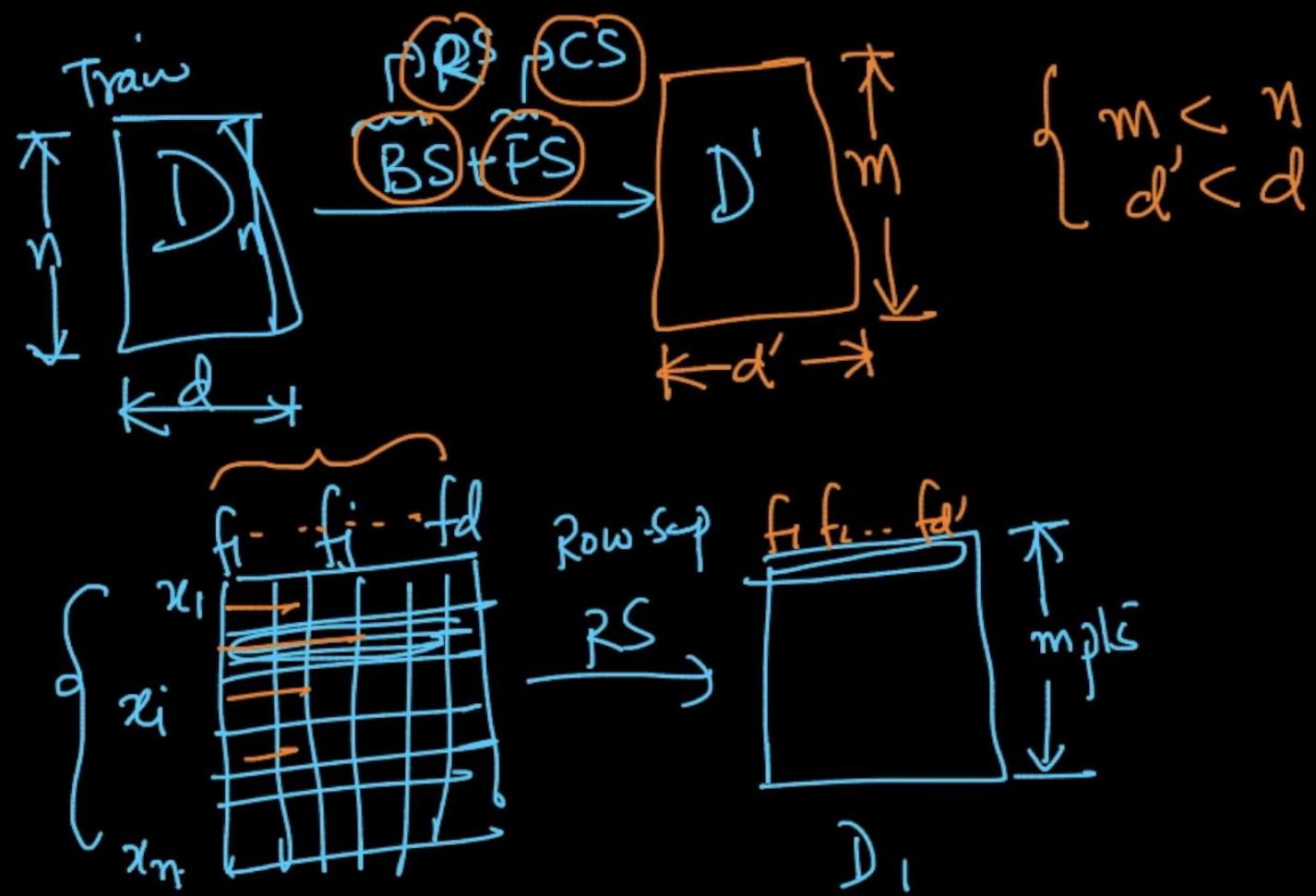
(j)

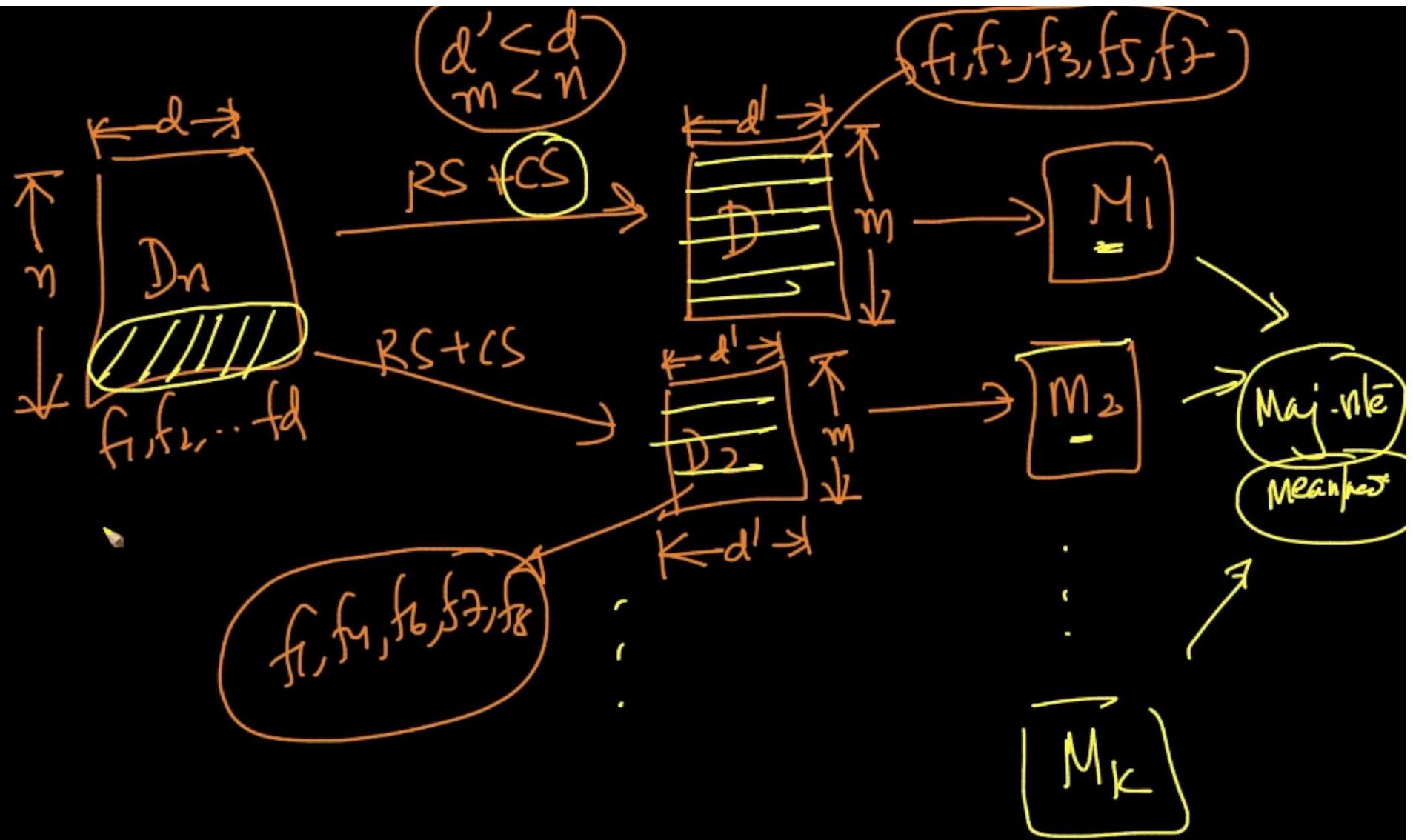
↑

Row Sampling w replacement

→ feature bagging

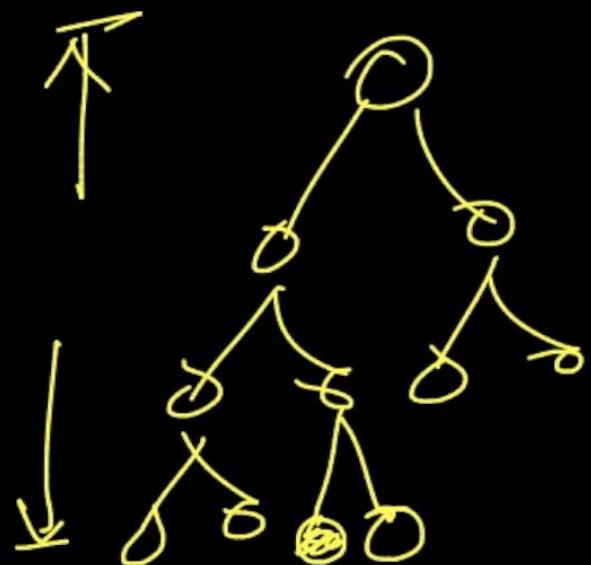
+ col sampling



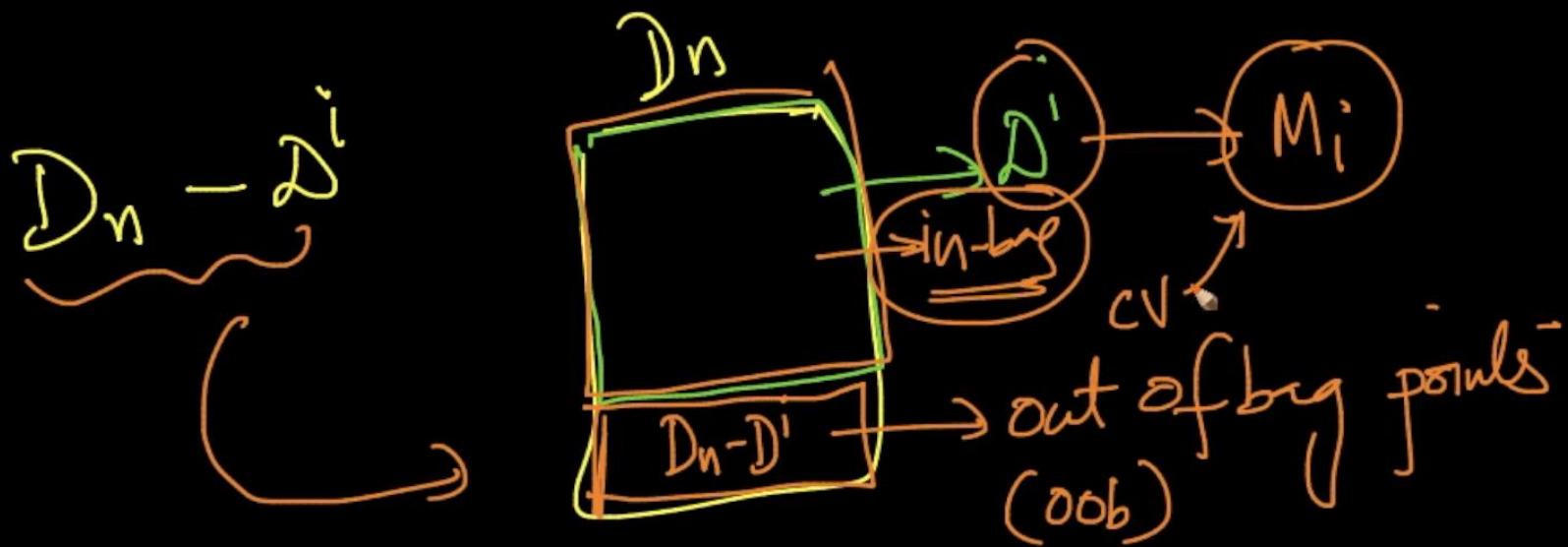


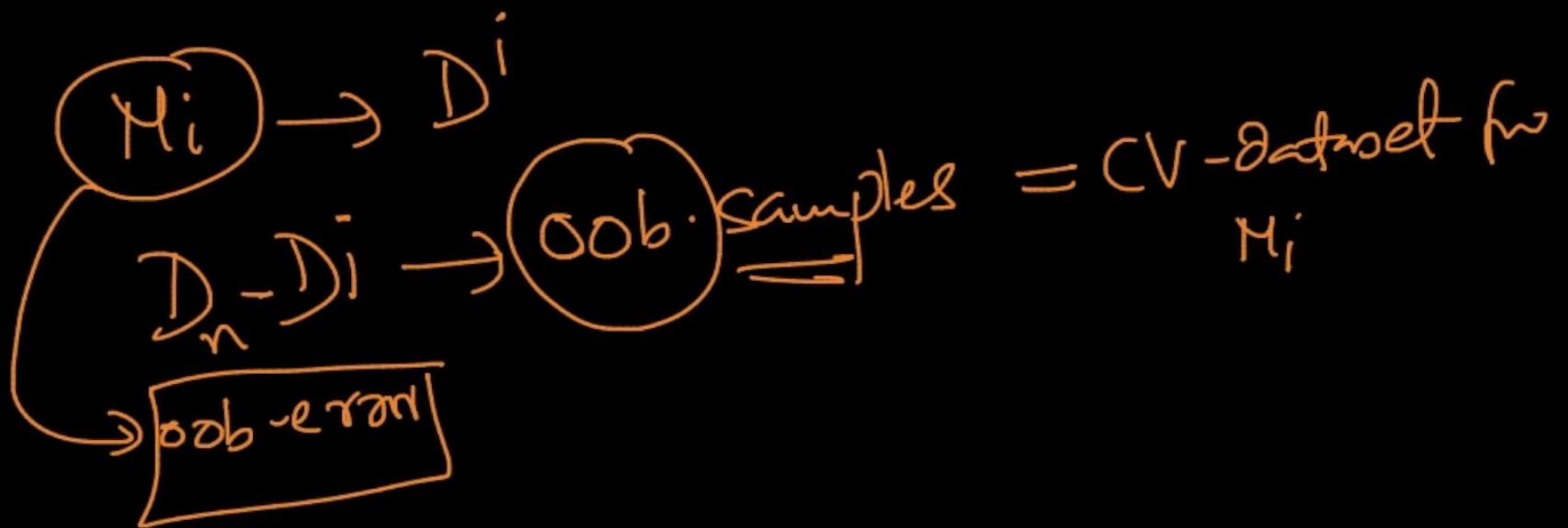
$\mathcal{H}, \mathcal{S}_j$ - DT of reasonable depth

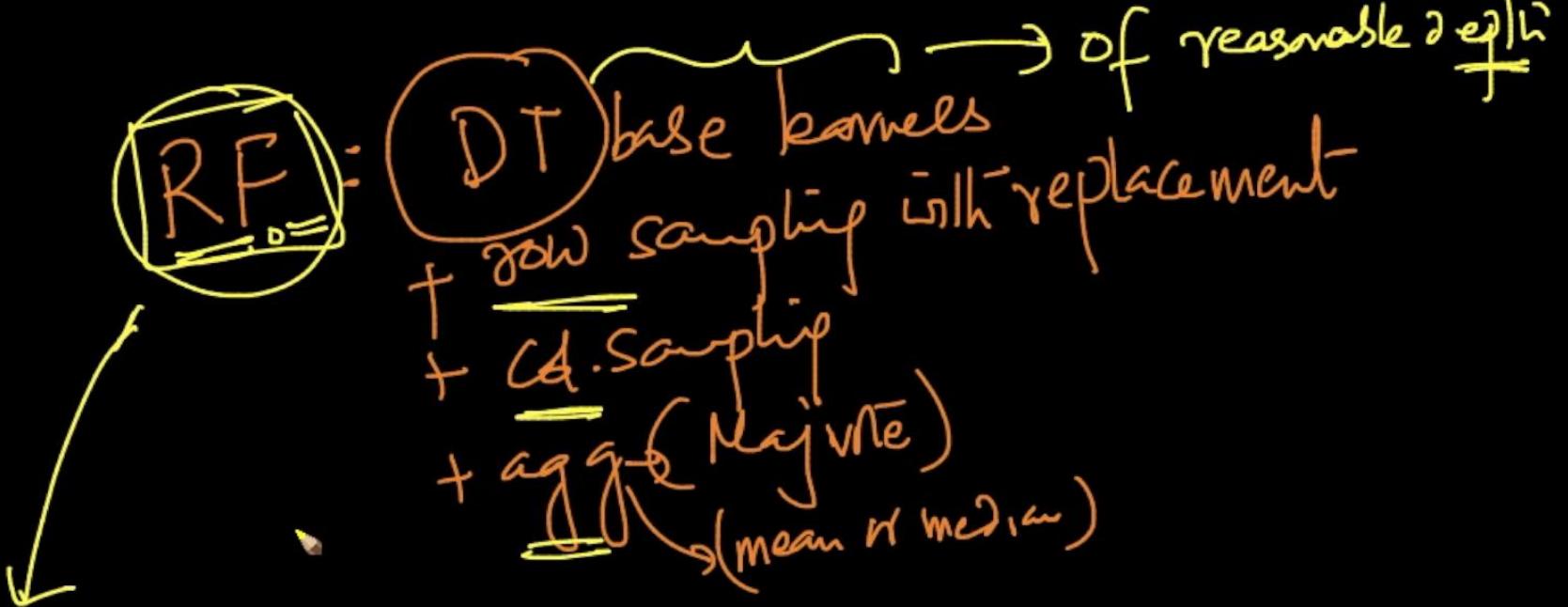
base learners: - low-bias
high-var
(reduces-variance)



$M_i \rightarrow D^i \rightarrow d' \text{ cols}$ $d' < d$
m rows.







Bias - Variance tradeoff

RF: reduce Variance



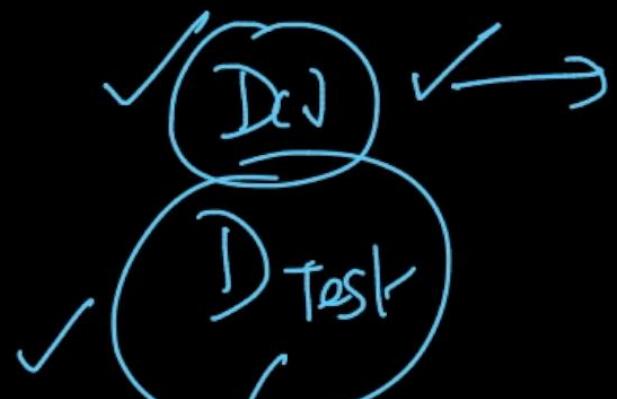
because base learners (M_i 's) are low-biased



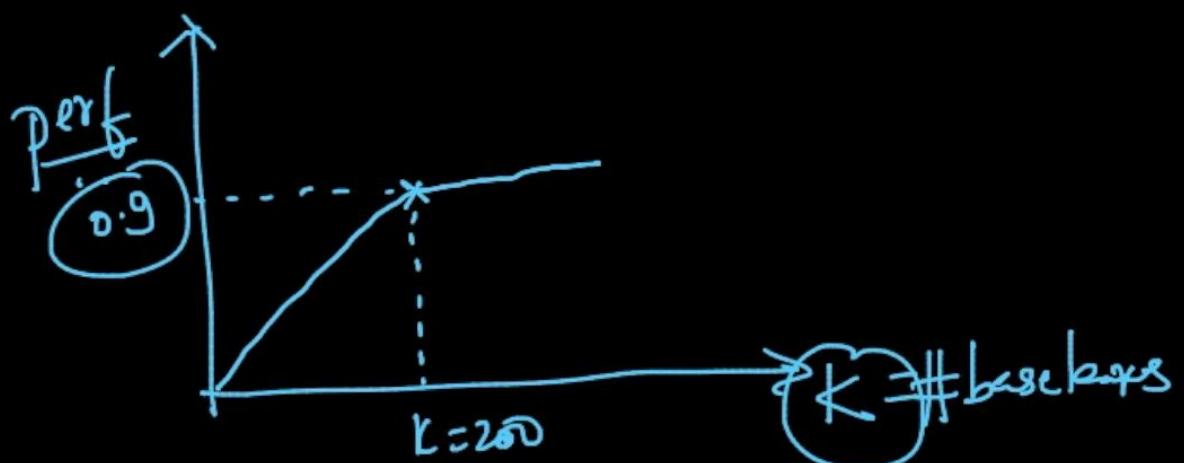
$M = \text{agg}(M_1, M_2, M_3, \dots, \underline{M}_k)$

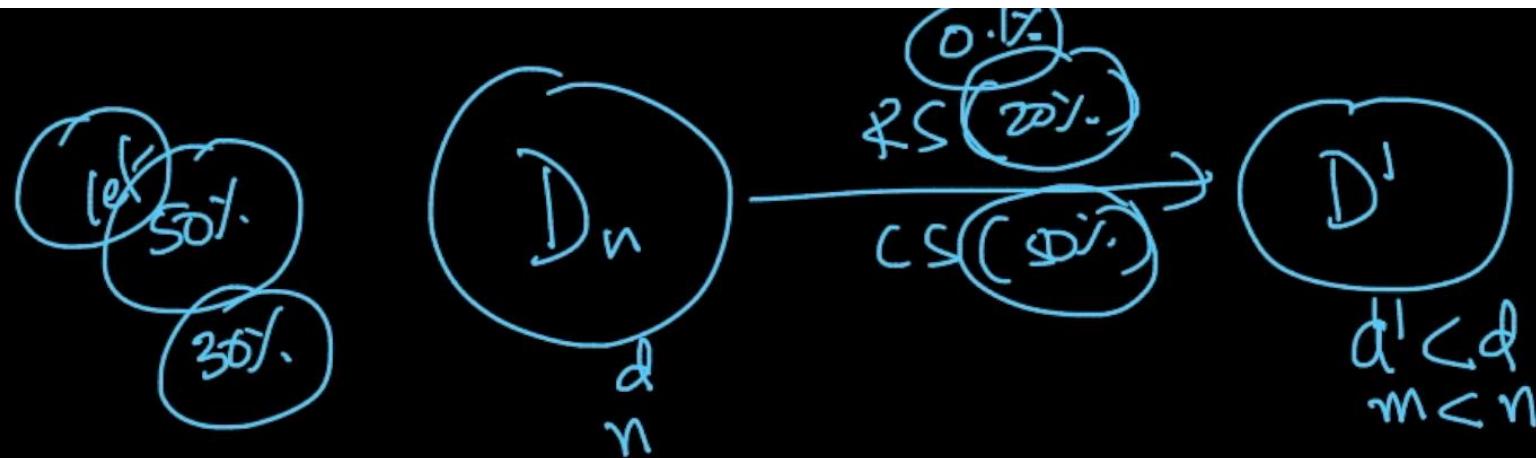
$$\text{bias}(\hat{M}) = \text{bias}(M_i)$$

$$\begin{cases} K \uparrow; \text{ Variance } \downarrow \\ K \downarrow; \text{ Variance } \uparrow \end{cases}$$



report your model





$$\left\{ \begin{array}{l} \frac{d'}{d} = \text{CS} \cdot S.R \\ \frac{m}{n} = \text{CS} \cdot S.R \end{array} \right.$$

$$\left\{ \begin{array}{l} Gf \cdot S \cdot R \downarrow \\ \text{and} \\ CS \cdot R \downarrow \end{array} \right. \quad \begin{array}{l} \xrightarrow{\text{low-variance model}} \\ \underline{\text{Var}} \downarrow \end{array}$$

$$\left\{ \begin{array}{l} k_j \rightarrow \# \text{base learners} \\ CSR \\ RSP \end{array} \right.$$

Train & Runtime Complexity

RF with K base-learners (DT)

multi-core
100-boxes

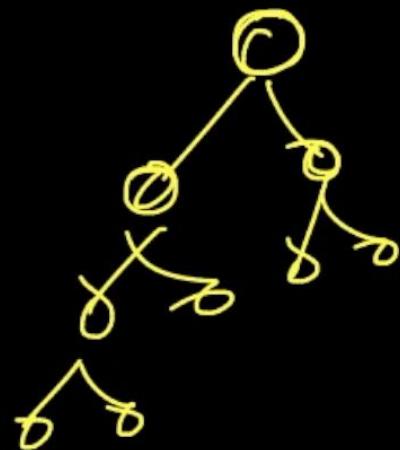
Train:

$$O(n \lg n d + k)$$

Trivially parallelizable



Runtime:-



$O(\text{depth} * k)$

↓

large (10^{20})



Space:





$O(DT * k)$

RF with 200 trees of depth = 6 → n3 MS

(if ... else)

$z_v \rightarrow y_w$

Extremely Randomized Trees: → Try out a random sample of possible values to determine C

regular DT
RF:

→ Col. Sampling; row sampling ; agg

{ Extreme Trees: →

M_i →

$f_i < C$

DT / RF

Y N

RF / DT {
 f_i : real valued values
 sort f_i
 n



(1D random values of f_i)

Try out all the possible values of f_i to determine C

Extreme Trees:- $\xrightarrow{\text{Col Sampling} + \gamma \cdot S + \text{agg.} \rightarrow RF}$
+ randomizatiⁿ when selecting C

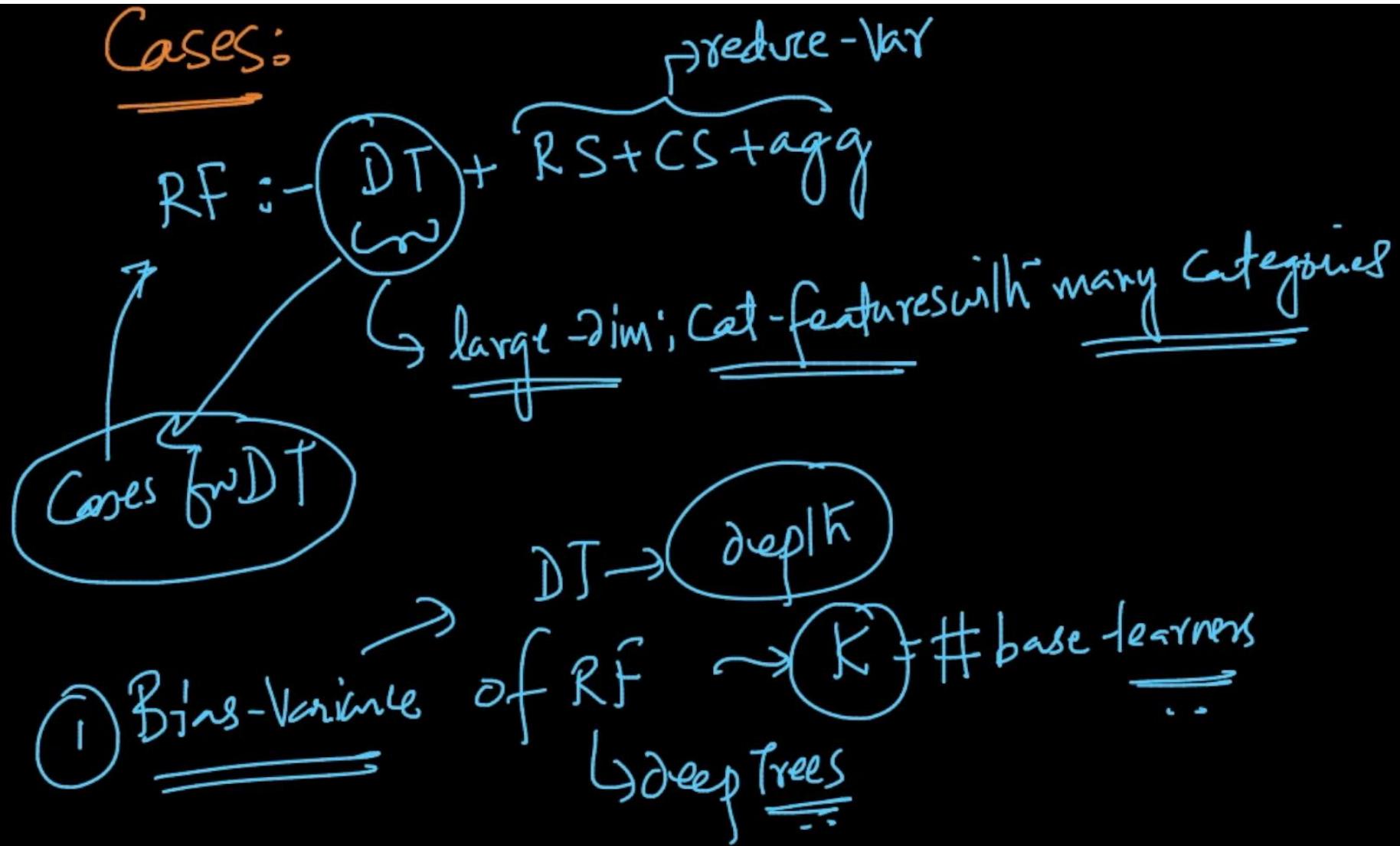
Randomization as a way to reduce Variance

$RF \rightarrow C \cdot S \& \gamma \cdot S$

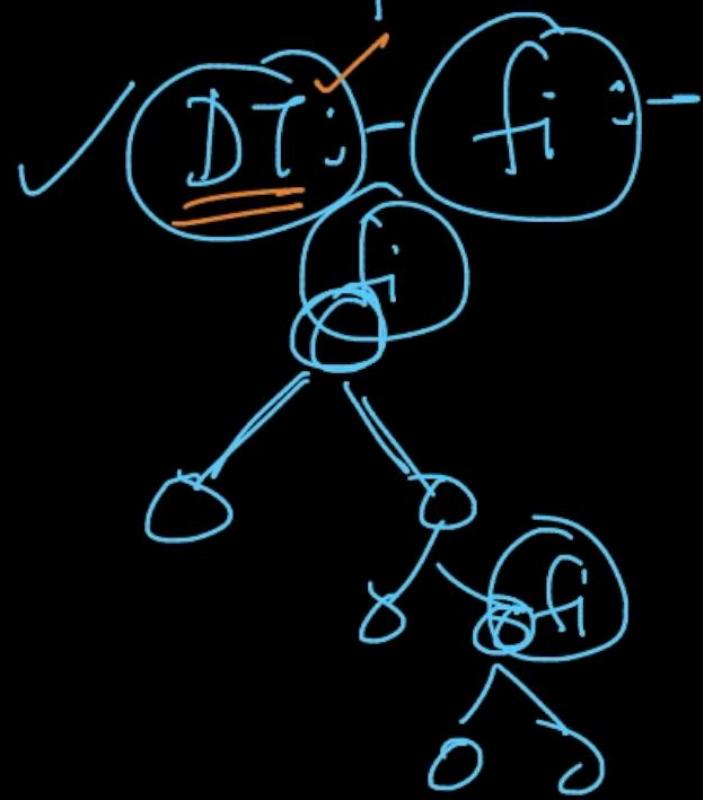
$ET \rightarrow C \cdot S + \gamma \cdot S + \underbrace{\text{random } C}_{\text{reduce Variance}}$

(reduce Variance better than RF)

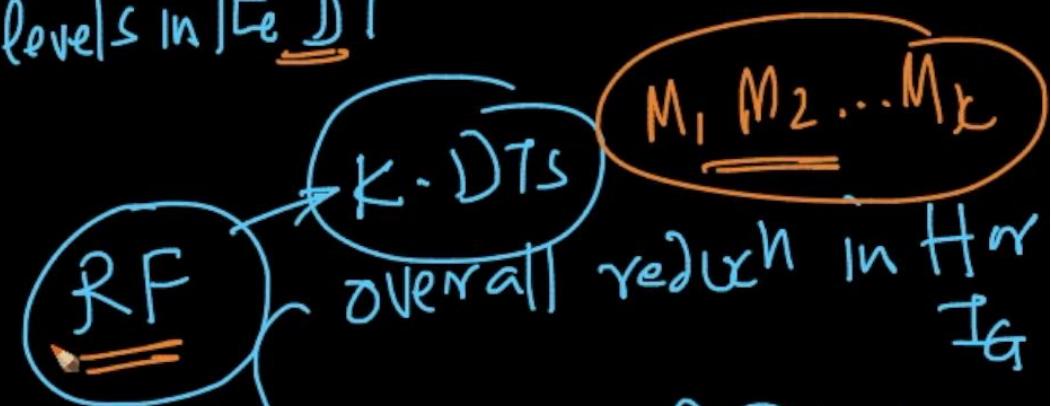
Cases:



② Feature Imp:



overall reduction in Entropy or I_G
because of this feature @ various
levels in the DT



because of f_i @ various
levels of each of M_i 's

✓ Boosting: intuition

Bagging: - high Var, low bias
base-models

+ reduce - Var,
randomization + aggregation
(CS, RS)

Boosting:
↓
low-var high bias
Error + bias + Var + E

+ additively combine
reduce bias
while keep our var low ..

GB

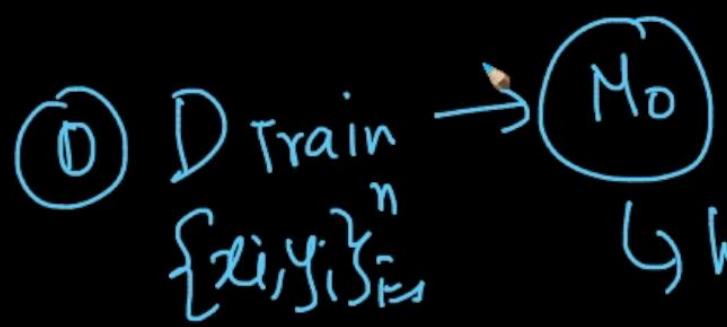
Core-Idea

: - reduces bias

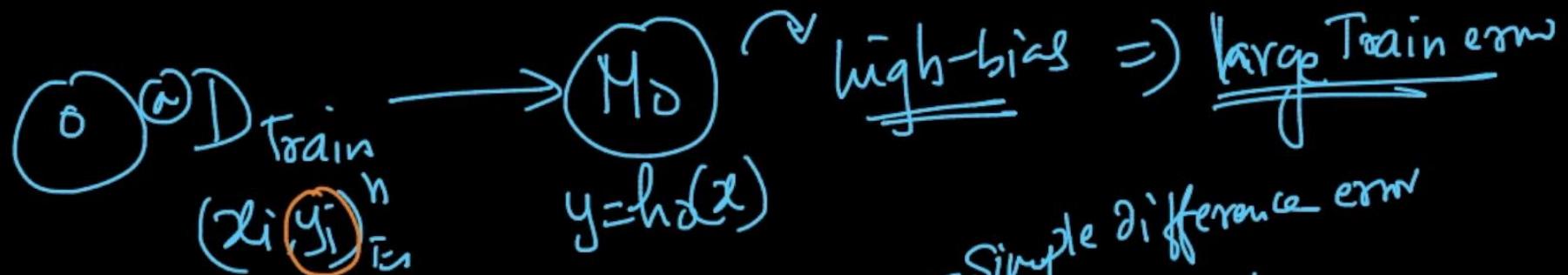
{ classifn ✓
regsn. ✓

$$\mathcal{D}_{\text{train}} = \{x_i, y_i\}_{i=1}^n$$

high-bias \Rightarrow high Train-err ~~✓~~

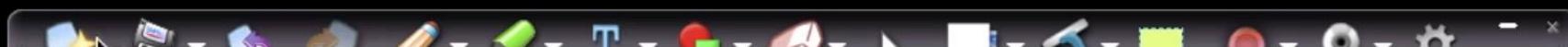
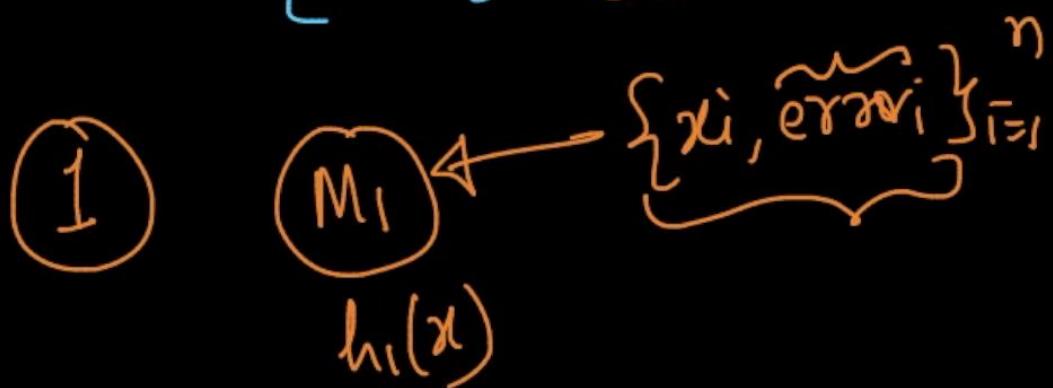


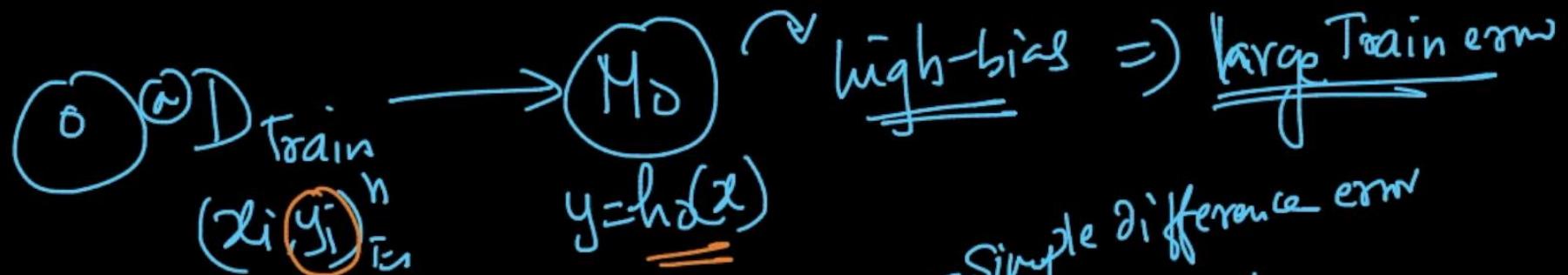
high-bias { e.g: DT which is shallow }



$y_i - \underline{h_0(x_i)} = \underline{\text{error}_i}$ $\left\{ \begin{array}{l} \text{Simple difference err} \\ \text{Sq.-error} \\ \text{log err} \\ \text{true err} \end{array} \right.$

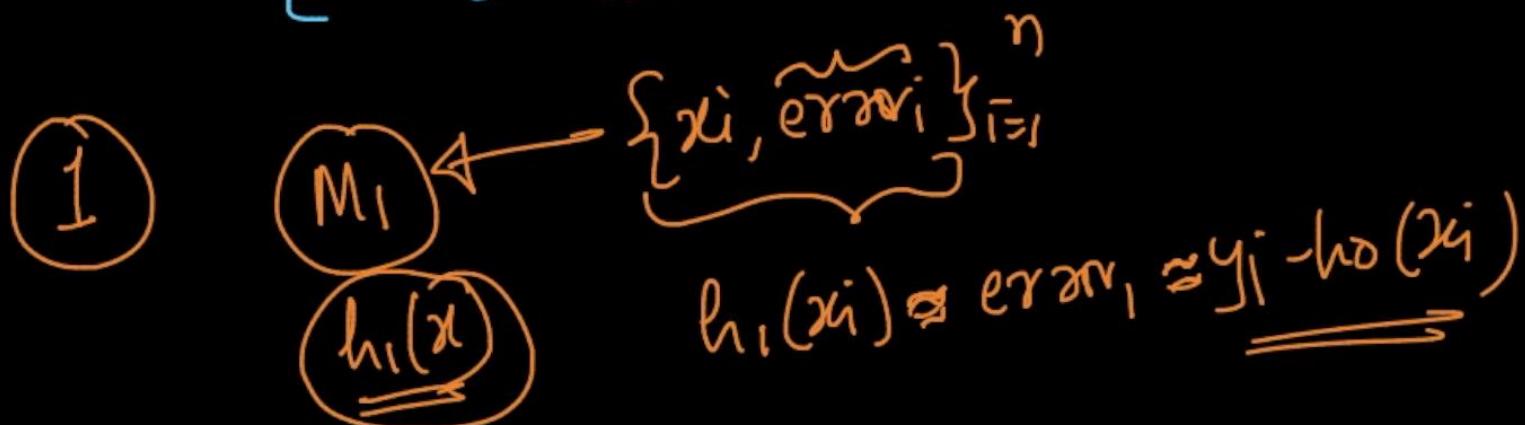
$\{x_i, y_i, \text{error}_i\}_{i=1}^n$





$$(b) \underline{\underline{y_i - h_D(x_i)}} = \underline{\underline{\text{error}_i}} \leftarrow \begin{cases} \text{Simple difference err} \\ \text{Sq.-error} \\ \text{log err} \\ \text{true err.} \end{cases}$$

$\{x_i, y_i, \text{error}_i\}_{i=1}^n$



$F_1(x)$ = model at end of Stage 1

$$F_1(x) = \underbrace{\alpha_0 h_0(x)}_{\text{bias term}} + \underbrace{\alpha_1 h_1(x)}_{\text{base model}} : - \text{weighted sum of 2-base models}$$

② $\{x_i, \text{err}_i\} \rightarrow M_2$ $h_2(x)$
 $y_i - F_1(x_i)$

$$F_2(x) = \underbrace{\alpha_0 h_0(x)}_{\text{bias term}} + \underbrace{\alpha_2 h_2(x)}_{\text{new base model}} + \underbrace{\alpha_1 h_1(x)}_{\text{old base model}}$$

end of stage k :-

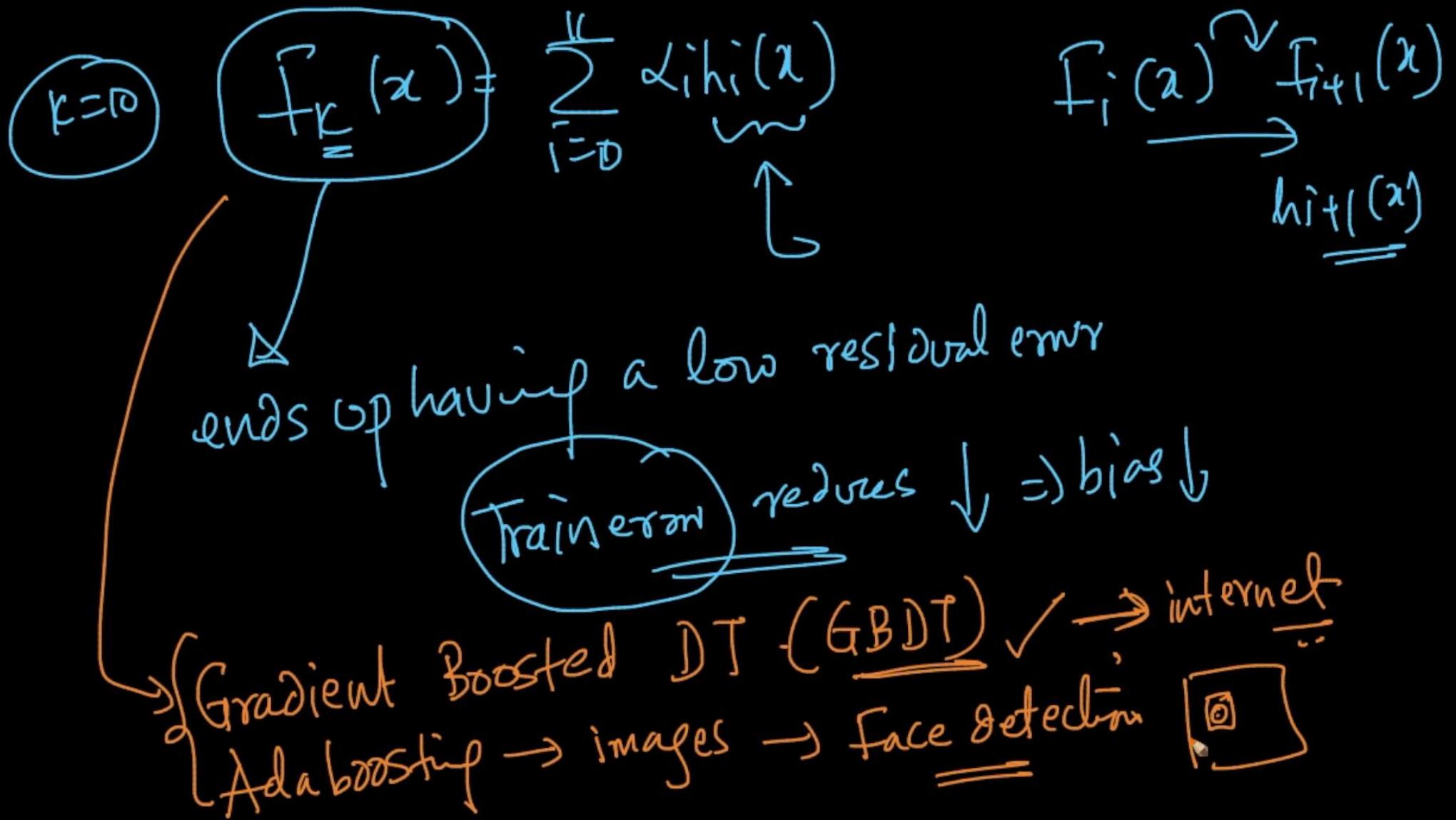
$$F_k(x) = \sum_{i=0}^k \alpha_i h_i(x)$$

trained to fit the
residual error @ end of
the prev stage

additive weighted model

$$h_i(x) \leftarrow \{x_i, \text{err}_i\}$$

$y_i - F_{i-1}(x)$
↑ residual error @ end of
stage $(i-1)$



Residuals, Loss-fns & Gradients → Gradient
Boosting

$$f_k(x) = \sum_{i=0}^k \alpha_i h_i(x)$$

residual } → err_i = y_i - f_k(x)
@ end of
stage k ↑ Residual

$$M_{k+1} \leftarrow \{x_i, err_i\}$$



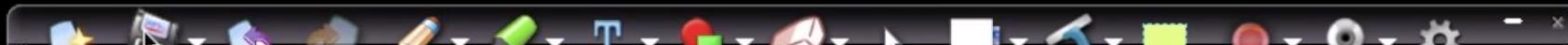
{ Loss-minimization: - ✓ Logistic-loss ← classfn
 ✓ Lr. regr ← Sq. loss) ← easier
 SVM ← hinge-loss

regression:

$$L(y_i, f_k(x_i)) = \underbrace{(y_i - f_k(x_i))^2}_{z_i} \quad \begin{cases} \text{let} \\ f_k(x_i) = z_i \end{cases}$$

Sq. loss

$$\begin{aligned}
 \frac{\partial L}{\partial f_k(x_i)} &= \frac{\partial L}{\partial z_i} = \frac{\partial}{\partial z_i} (y_i - z_i)^2 \\
 &= (-1) * 2 * (y_i - z_i) \\
 \frac{\partial L}{\partial z_i} &= -2(y_i - z_i)
 \end{aligned}$$

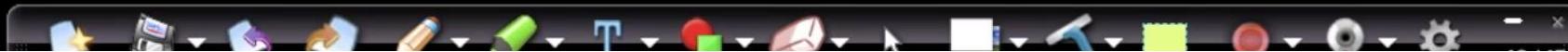


$$-\frac{\partial \tilde{L}}{\partial F_K(x_i)} = \boxed{2} \underbrace{(y_i - \hat{F}_K(x_i))}_{\text{residual}} \overset{\text{err}_i}{\sim}$$

neg derivative

* neg. gradient \approx residual

\uparrow pseudo-residual



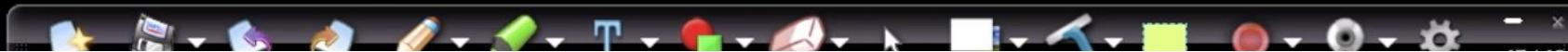
$$h_i(x) \leftarrow x_i, e_{\text{rr}_i} = y_i - f_{i-1}(x)$$

$\checkmark e_{\text{rr}_i} \rightarrow \text{residual}_i$

$e_{\text{rr}_i} \leftarrow \text{pseudo-residual}$

$$-\frac{\partial L}{\partial F_{i-1}(x)}$$

Let us have
any loss fn which
is differentiable



Stage: i

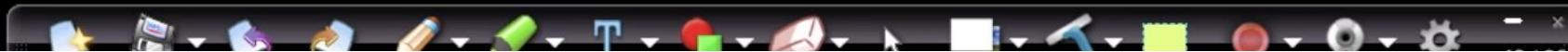
$f_{i+1}(x)$

$(x_i, \text{err}_i) \xrightarrow{\sim} M_i$

$h_i(x)$

pseudo-residual

$$-\frac{\partial \mathcal{L}}{\partial f_{i+1}(x)}$$



✓ RF → min hinge loss

✓ GB → min any loss
differentiable

↳ Super-powerful

Ada Boost

base-models → combine pseudo-residuals
high bias ↳ $\frac{\partial L}{\partial f_{i-1}(x)}$



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 (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

3. Compute multiplier γ_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 $F_M(x)$.

↑ logistic-loss

↑ SQ-loss

↑ hinge-loss

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

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 (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

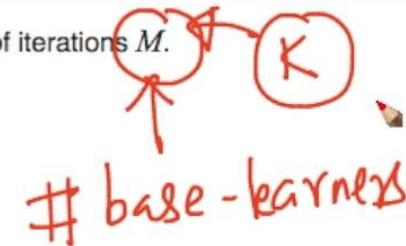
3. Compute multiplier γ_m by solving the following one-dimensional optimization problem:

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

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

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:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma). \quad \text{constant} \checkmark \rightarrow \text{(SO-loss)}$$

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 (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

3. Compute multiplier γ_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 commonly used with decision trees (especially CART trees) of a fixed size as base learners. For this special case, we can write the update rule as:

APPLIED COURSE

Secure | https://en.wikipedia.org/wiki/Gradient_boosting



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:

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

Zeroth model
 $y_i = \text{mean}(y_i)$

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 (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

3. Compute multiplier γ_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 $F_M(x)$.

$$\arg \min_{\gamma} \sum_{i=1}^n (y_i - \gamma)^2 \quad \left\{ \begin{array}{l} \gamma = \text{mean}(y_i) \\ \hline \end{array} \right.$$

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

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:

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

$h_1(x), h_2(x) \dots h_M(x)$

$h_m(x)$
←
 m^{th} stage

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 (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

3. Compute multiplier γ_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

Secure | https://en.wikipedia.org/wiki/Gradient_boosting



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*:

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

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

3. Compute multiplier γ_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:

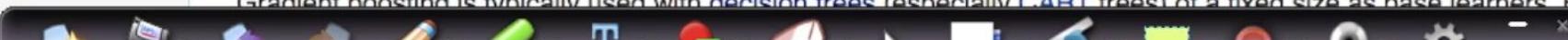
$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Output $F_M(x)$.

err^m
err_i

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



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 (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

3. Compute multiplier γ_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 $F_M(x)$.

$$F_M(x) = \underbrace{h_0(x)}_{\gamma_0 h_0(x)} + \underbrace{\gamma_1 h_1(x)}_{\gamma_1 h_1(x)} + \underbrace{\gamma_2 h_2(x)}_{\gamma_2 h_2(x)} + \dots + \underbrace{\gamma_m h_m(x)}_{\gamma_m h_m(x)}$$

mth iterⁿ

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

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:

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

$$\underline{\underline{F_M(x)}}$$

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 (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

3. Compute multiplier γ_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 $F_M(x)$.

$$\left\{ \begin{array}{l} F_M(x) = F_{m-1}(x) + \gamma_m h_m(x) \\ \underline{\underline{F_M(x)}} \end{array} \right.$$

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

Secure | https://en.wikipedia.org/wiki/Gradient_boosting



*high bias
low var*

GB

GBDT

Shallow DT

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 (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

3. Compute multiplier γ_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 $F_M(x)$.

1,2,3,4,5

Decision Stump

$h_m(x)$

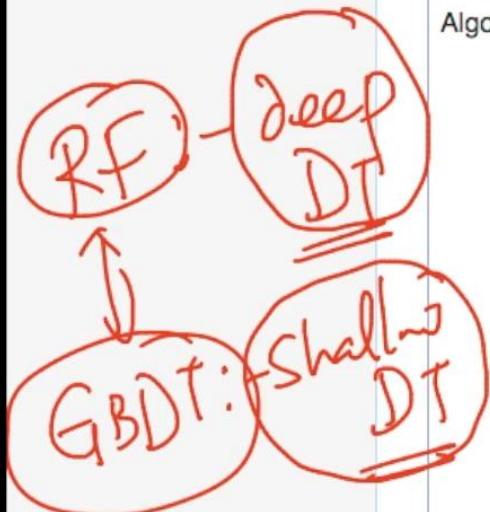
decision tree

Gradient tree boosting [edit]

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

APPLIED COURSE

Secure | https://en.wikipedia.org/wiki/Gradient_boosting



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 (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

3. Compute multiplier γ_m by solving the following one-dimensional optimization problem:

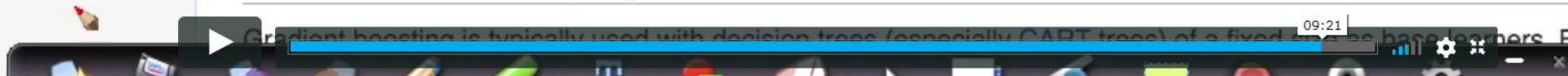
$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

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]

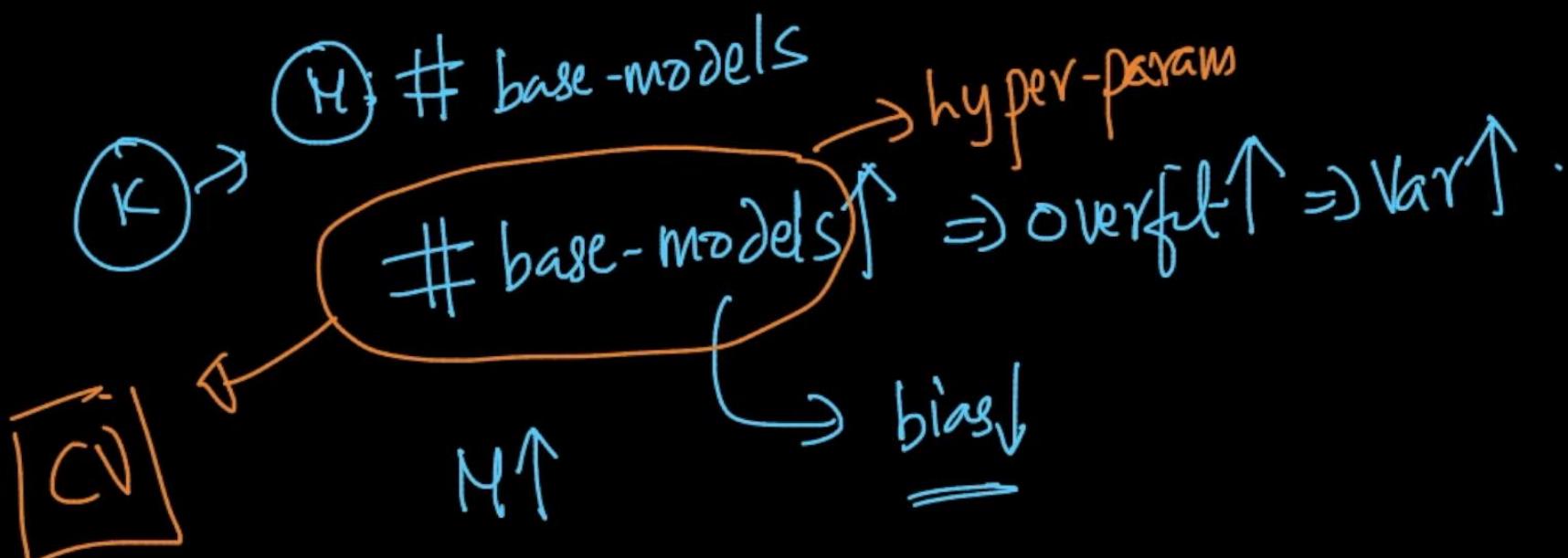


GB

Regularization & Shrinkage:

$$F_M(x) = h_0(x) + \sum_{m=1}^M \gamma_m h_m(x)$$

bias - Var
— — —



[Secure | https://en.wikipedia.org/wiki/Gradient_boosting](#)

Regularization [edit]

Fitting the training set too closely can lead to degradation of the model's generalization ability. Several so-called regularization techniques reduce this [overfitting](#) effect by constraining the fitting procedure.

One natural regularization parameter is the number of gradient boosting iterations M (i.e. the number of trees in the model when the base learner is a decision tree). Increasing M reduces the error on training set, but setting it too high may lead to overfitting. An optimal value of M is often selected by monitoring prediction error on a separate validation data set. Besides controlling M , several other regularization techniques are used.

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

Shrinkage [edit]

An important part of gradient boosting method is regularization by shrinkage which consists in modifying the update rule as follows:

$$F_m(x) = F_{m-1}(x) + \nu \cdot \gamma_m h_m(x), \quad 0 < \nu \leq 1,$$

where parameter ν is called the "learning rate".

$\nu \approx 0.1$
Hyperparam \rightarrow CV

ν is small

0.0001

Empirically it has been found that using small learning rates (such as $\nu < 0.1$) yields dramatic improvements in model's generalization ability over gradient boosting without shrinking ($\nu = 1$).^[7] However, it comes at the price of increasing computational time both during training and querying: lower learning rate requires more iterations.

Stochastic gradient boosting [edit]

Soon after the introduction of gradient boosting Friedman proposed a minor modification to the algorithm, motivated by Breiman's bagging method.^[3] Specifically, he proposed that at each iteration of the algorithm, a base learner should be fit on a

Regularization [edit]

Fitting the training set too closely can lead to degradation of the model's generalization ability. Several so-called regularization techniques reduce this [overfitting](#) effect by constraining the fitting procedure.

One natural regularization parameter is the number of gradient boosting iterations M (i.e. the number of trees in the model when the base learner is a decision tree). Increasing M reduces the error on training set, but setting it too high may lead to overfitting. An optimal value of M is often selected by monitoring prediction error on a separate validation data set. Besides controlling M , several other regularization techniques are used.

↓
Overfitting
↓
Var↓

Shrinkage [edit]

An important part of gradient boosting method is regularization by [shrinkage](#) which consists in modifying the update rule as follows:

$$F_m(x) = F_{m-1}(x) + \nu \cdot \gamma_m h_m(x), \quad 0 < \nu \leq 1,$$

where parameter ν is called the "learning rate".

$M = \# \text{base-models}$
 $\nu \uparrow \Rightarrow \text{overfitting} \uparrow$

Empirically it has been found that using small learning rates (such as $\nu < 0.1$) yields dramatic improvements in model's generalization ability over gradient boosting without shrinking ($\nu = 1$).^[7] However, it comes at the price of increasing computational time both during training and querying: lower learning rate requires more iterations.

Stochastic gradient boosting [edit]

Soon after the introduction of gradient boosting Friedman proposed a minor modification to the algorithm, motivated by Breiman's bagging method.^[3] Specifically, he proposed that at each iteration of the algorithm, a base learner should be fit on a



GB

or GBDT :- overfit

"M" is large

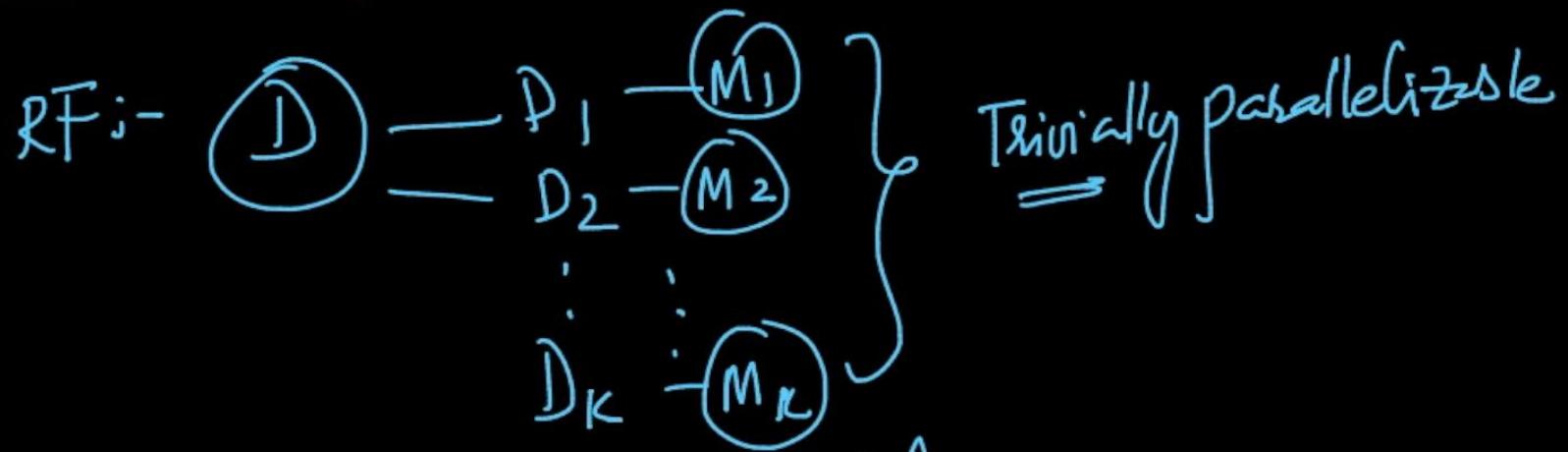
overfit

$M = 2,000$

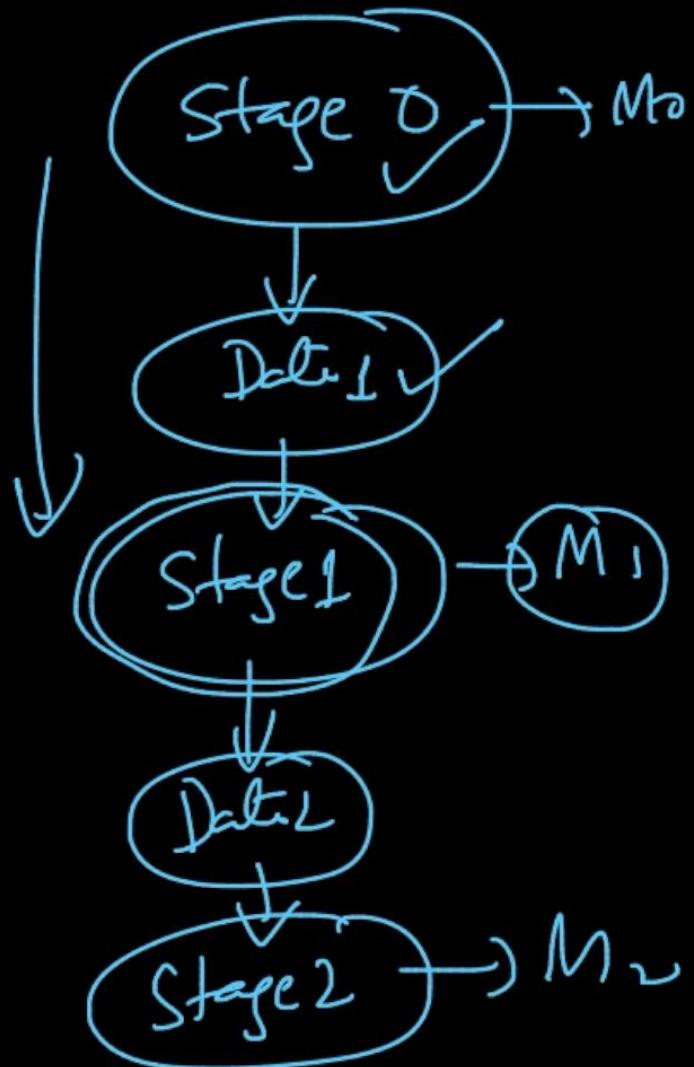
Train & Run time complexity



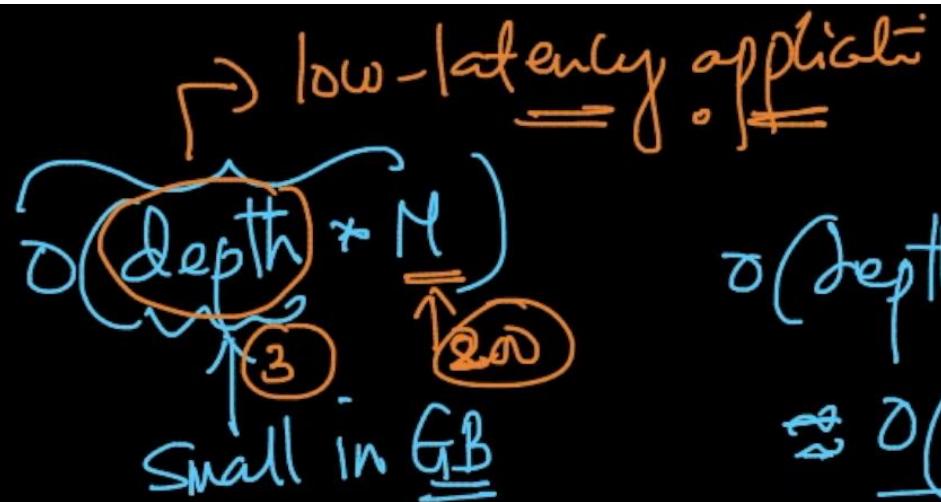
Train: $O(n \lg n d * M)$ M : #base-learners



GBDT :- not easy to parallelize :-

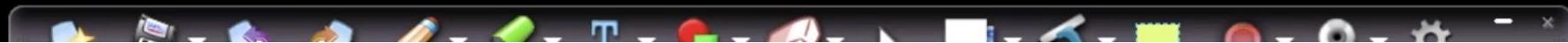
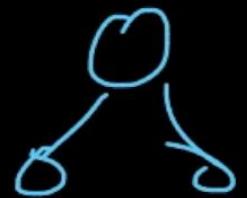
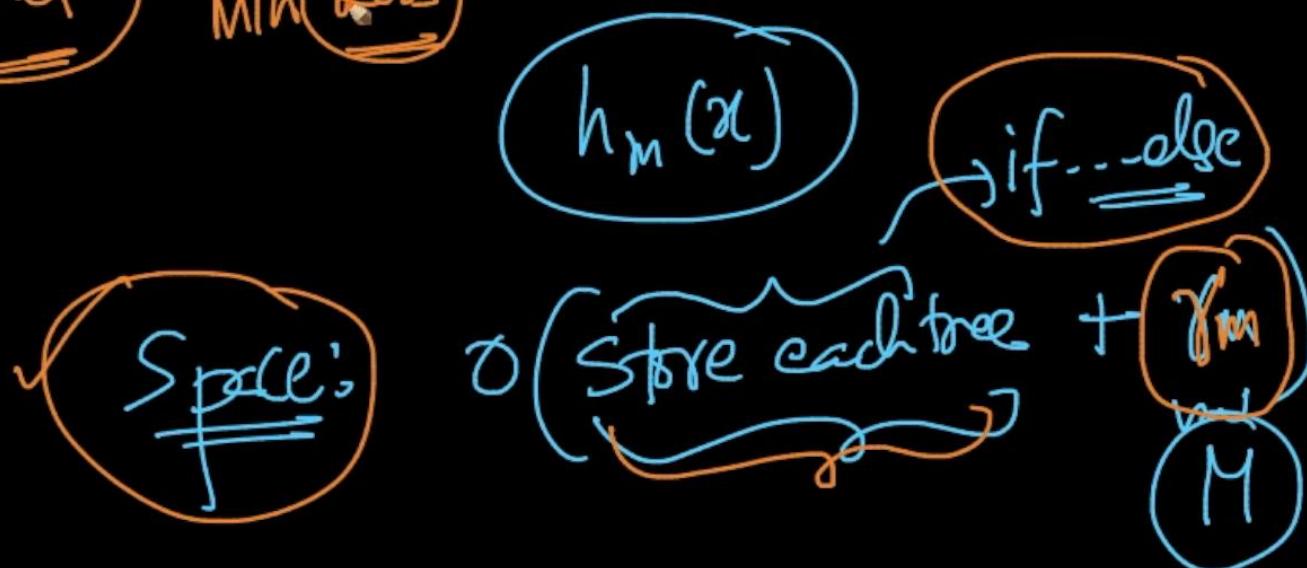


$\left\{ \begin{array}{l} \text{GBDT take more} \\ \text{time to train than} \\ \text{RF} \end{array} \right.$



$$\mathcal{O}(\text{depth} \times M + M)$$

$$\approx \mathcal{O}(\text{depth} + M)$$



XGBoost

[Get Started](#)
[Tutorials](#)
[How To](#)
[Packages](#)
[Knobs](#)

Scikit-Learn API

Scikit-Learn Wrapper interface for XGBoost.

```
class xgboost.XGBRegressor(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True, objective='reg:linear',
booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1, max_delta_step=0, subsample=1,
colsample_bytree=1, colsample_bylevel=1, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, base_score=0.5, random_state=0,
seed=None, missing=None, **kwargs)
```

Bases: `xgboost.sklearn.XGBModel`, `object`

Implementation of the scikit-learn API for XGBoost regression.

Parameters

`max_depth : int`

Maximum tree depth for base learners.

`learning_rate : float`

Boosting learning rate (xgb's "eta")

`n_estimators : int`

Number

✓ best GBDT

GBDT: - p.r + additive

SK :- GBDT + r.s

XG BOOST: GBDT + $\frac{r.s}{c-s}$

RF

Same-add

SKlearn

XG Boost



A toy example from Schapire's tutorial

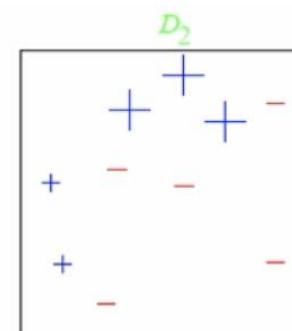
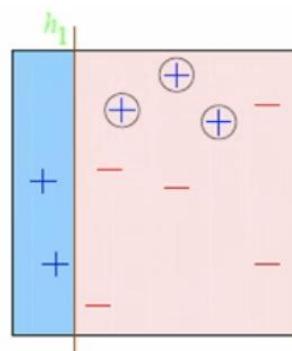
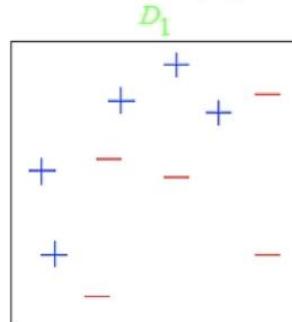
Consider this toy dataset in 2D and assume that our weak classifiers are decision stumps (vertical or horizontal half-planes):

GBDT

Boosting! - Ada-boost

The first round:

Image
face detection

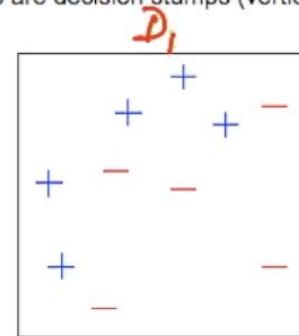


The second round:



A toy example from Schapire's tutorial

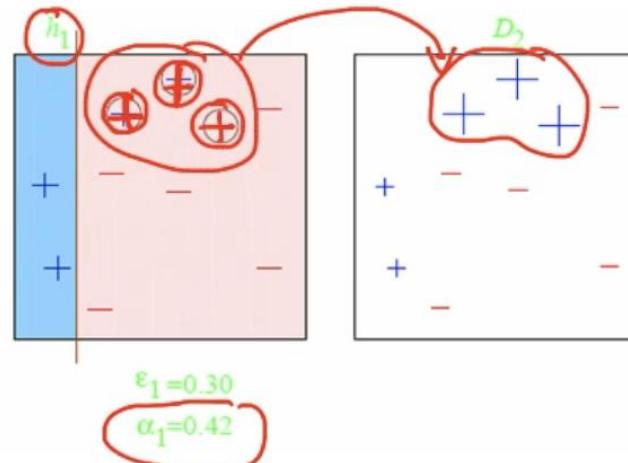
Consider this toy dataset in 2D and assume that our weak classifiers are decision stumps (vertical or horizontal half-planes):



DT. with depth=1
↓
Decision Stump

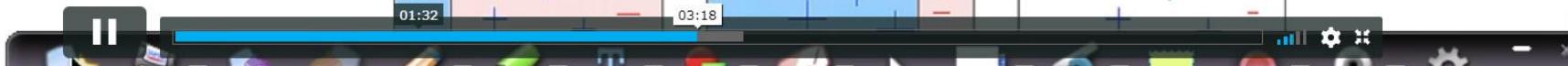
The first round:

$d_1 h_1(x)$



weighted models
↓
upsample

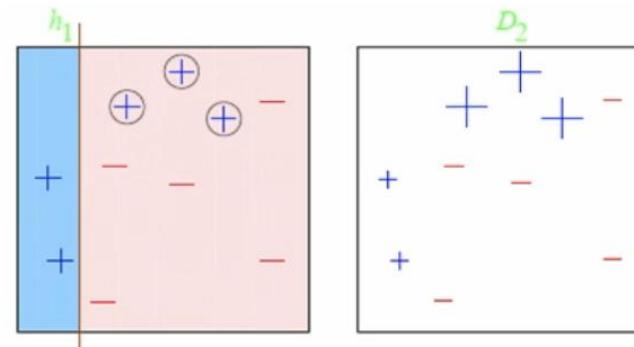
The second round:



Secure | https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=lectures.boosting

The first round:

$\alpha_1 h_1$



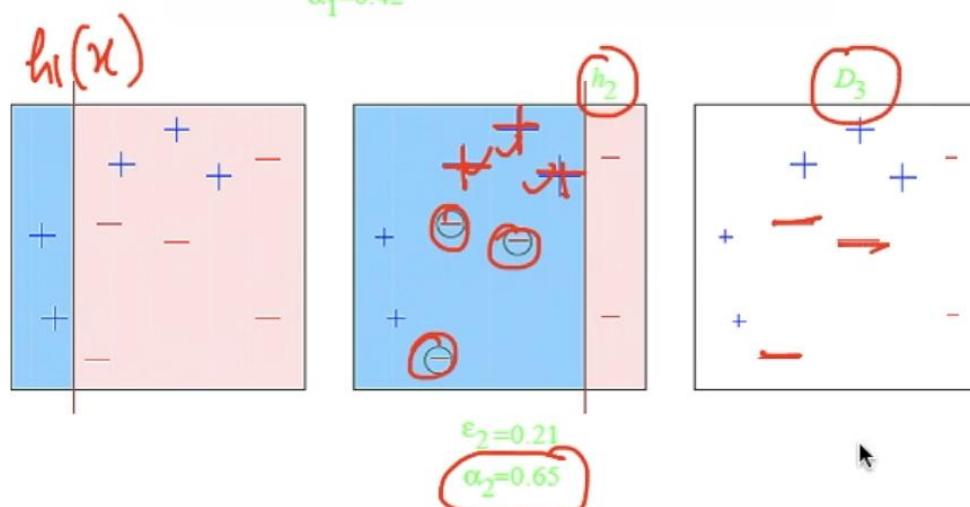
$$\epsilon_1 = 0.30$$

$$\alpha_1 = 0.42$$

~~exponentially~~

The second round:

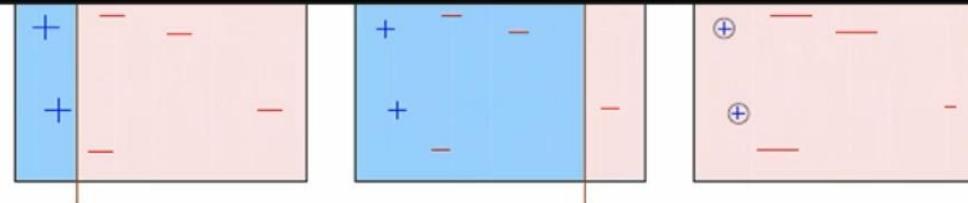
$\alpha_2 h_2$



The third round:



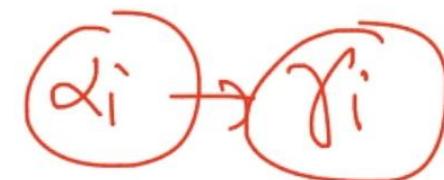
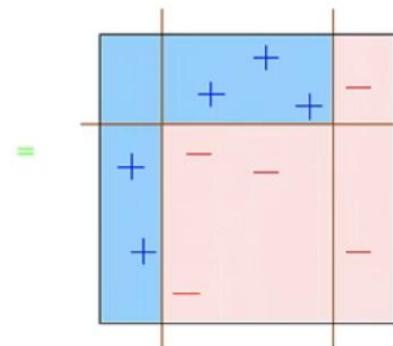
Secure | https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=lectures.boosting



The final classifier

$$F_3(x) = \alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3$$

$$H_{\text{final}} = \text{sign} \left(\alpha_1 \begin{array}{|c|c|} \hline \text{blue} & \text{red} \\ \hline 0.42 & \\ \hline \end{array} + \alpha_2 \begin{array}{|c|c|} \hline \text{blue} & \text{red} \\ \hline & 0.65 \\ \hline \end{array} + \alpha_3 \begin{array}{|c|c|} \hline \text{blue} & \text{red} \\ \hline & 0.92 \\ \hline \end{array} \right)$$



Contents - Google Docs x StackingClassifier - mlxtend x Chekuri Srikan...

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

Overview

$h_1(x)$ $h_2(x)$

Stacking is an ensemble learning technique to combine multiple classification models via a meta-classifier. The individual classification models are trained based on the complete training set; then, the meta-classifier is fitted based on the outputs -- meta-features -- of the individual classification models in the ensemble. The meta-classifier can either be trained on the predicted class labels or probabilities from the ensemble.

```

graph TD
    TS[Training set] --> C1((G1))
    TS --> C2((G2))
    TS --> Cm((Gm))
    C1 --> P1((P1))
    C2 --> P2((P2))
    Cm --> Pm((Pm))
    P1 --> MC{Meta-Classifier}
    P2 --> MC
    Pm --> MC
    MC --> Pf((Pf))
    ND[New data] --> C1
    
```

Training set

Classification models

Predictions

Meta-Classifier

Final prediction

New data

m-models

ideally parallel & independent

Contents - Google Docs StackingClassifier - mlxtend

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

Overview

Stacking is an ensemble learning technique to combine multiple classification models via a meta-classifier. The individual classification models are trained based on the complete training set; then, the meta-classifier is fitted based on the outputs -- meta-features -- of the individual classification models in the ensemble. The meta-classifier can either be trained on the predicted class labels or probabilities from the ensemble.

```

graph TD
    TS[Training set] --> C1((C1))
    TS --> C2((C2))
    TS --> ...((...))
    TS --> Cm((Cm))
    C1 -- prediction --> P1((P1))
    C2 -- prediction --> P2((P2))
    ...((...)) -- prediction --> ...((...))
    Cm -- prediction --> Pm((Pm))
    P1 --> MC[Meta-Classifier]
    P2 --> MC
    ...((...)) --> MC
    Pm --> MC
    MC --> Pf((Pf))
    ND[New data] --> C1
    ND --> C2
    ND --> Cm
  
```

The diagram illustrates the Stacking process. It starts with a 'Training set' which feeds into multiple 'Classification models' (C_1, C_2, \dots, C_m). Each model produces a 'Prediction' (P_1, P_2, \dots, P_m). These predictions are then fed into a 'Meta-Classifier', which finally outputs the 'Final prediction' (P_f). A 'New data' input feeds into all the classification models.

Handwritten notes:

- C_1 : Ex-SUM (circled in purple)
- C_2 : RBF-SUM
- C_3 : GBDT(20)
- C_4 : NB
- C_5 : K-NN ($K=5$)
- C_6 : K-NN ($K=3$)

Contents - Google Docs X StackingClassifier - mlxtend Chekuri Srikan...

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

StackingClassifier

- Overview
- References
- Example 1 - Simple Stacked Classification
- Example 2 - Using Probabilities as Meta-Features
- Example 3 - Stacked Classification and GridSearch
- Example 4 - Stacking of Classifiers that Operate on Different Feature Subsets
- API
- Methods

Overview

Stacking is an ensemble learning technique to combine multiple classification models via a meta-classifier. The individual classification models are trained based on the complete training set; then, the meta-classifier is fitted based on the outputs -- meta-features -- of the individual classification models in the ensemble. The meta-classifier can either be trained on the predicted class labels or probabilities from the ensemble.

Training set

Classification models

Predictions

Meta-Classifier

New data

Final prediction

Handwritten notes:

- A circled 'RF' is followed by handwritten text: 'C → high var' and 'agg → high var'.
- A circled 'C' is followed by handwritten text: 'C → low bias'.
- A circled 'y' is followed by handwritten text: 'y → high var'.
- A circled 'Model' is followed by a circled 'Meta-classifier'.

Contents - Google Docs x StackingClassifier - mlxtend x Chekuri Srikan...

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

StackingClassifier

- Overview
- References
- Example 1 - Simple Stacked Classification
- Example 2 - Using Probabilities as Meta-Features
- Example 3 - Stacked Classification and GridSearch
- Example 4 - Stacking of Classifiers that Operate on Different Feature Subsets
- API
- Methods

Overview

Stacking is an ensemble learning technique to combine multiple classification models via a meta-classifier. The individual classification models are trained based on the complete training set; then, the meta-classifier is fitted based on the outputs -- meta-features -- of the individual classification models in the ensemble. The meta-classifier can either be trained on the predicted class labels or probabilities from the ensemble.

Training set

Classification models

Predictions

Meta-Classifier

Final prediction

New data

RF
C → high var
low bias
agg → high var

Handwritten notes:

RF

$C_1 \rightarrow$ high var
 $C_2 \rightarrow$ low bias
 $\text{agg} \rightarrow$ high var

Hand-drawn diagram:

A hand-drawn diagram illustrates the stacking process. It shows a cluster of circles labeled C_1, C_2, \dots, C_m representing individual classifiers. Arrows point from each C_i to a corresponding circle labeled $y_{i1}, y_{i2}, \dots, y_{in}$, representing predictions for different data points. A large bracket groups all these circles under the label "Model". An arrow points from this group to a box labeled "Meta-classifier". Another arrow points from the "Meta-classifier" box to a final circle labeled y_f , representing the final prediction. The handwritten text "L2 DT" is written near the bottom left of the diagram.

Contents - Google Docs StackingClassifier - mlxtend Chekuri Srikan...

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

StackingClassifier

Overview References Example 1 - Simple Stacked Classification Example 2 - Using Probabilities as Meta-Features Example 3 - Stacked Classification and GridSearch Example 4 - Stacking of Classifiers that Operate on Different Feature Subsets API Methods

Final prediction P_f

The algorithm can be summarized as follows (source: [1]):

Algorithm 19.7 Stacking

Input: Training data $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^m$ ($\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathcal{Y}$)
Output: An ensemble classifier H

```
1: Step 1: Learn first-level classifiers
2: for  $t \leftarrow 1$  to  $T$  do
3:   Learn a base classifier  $h_t$  based on  $\mathcal{D}$ 
4: end for
5: Step 2: Construct new data sets from  $\mathcal{D}$ 
6: for  $i \leftarrow 1$  to  $m$  do
7:   Construct a new data set that contains  $\{\mathbf{x}'_i, y_i\}$ , where  $\mathbf{x}'_i = \{h_1(\mathbf{x}_i), h_2(\mathbf{x}_i), \dots, h_T(\mathbf{x}_i)\}$ 
8: end for
9: Step 3: Learn a second-level classifier
10: Learn a new classifier  $h'$  based on the newly constructed data set
11: return  $H(\mathbf{x}) = h'(h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_T(\mathbf{x}))$ 
```

$\mathcal{D} = \circlearrowleft$

$h_1, h_2, h_3 \dots h_T$

bias
var

References

- [1] Tang, J., S. Alelyani, and H. Liu. "Data Classification: Algorithms and Applications." Data Mining and Knowledge Discovery Series, CRC Press (2015): pp. 498-500.
- [2] Wolpert, David H. "Stacked generalization." Neural networks 5.2 (1992): 241-259.

Contents - Google Docs StackingClassifier - mlxtend Chekuri Srikan...

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

StackingClassifier

Overview References Example 1 - Simple Stacked Classification Example 2 - Using Probabilities as Meta-Features Example 3 - Stacked Classification and GridSearch Example 4 - Stacking of Classifiers that Operate on Different Feature Subsets API Methods

D'

Final prediction

The algorithm can be summarized as follows (source: [1]):

Algorithm 19.7 Stacking

Input: Training data $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^m$ ($\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathcal{Y}$)

Output: An ensemble classifier H

1: Step 1: Learn first-level classifiers

2: **for** $t \leftarrow 1$ to T **do**

3: Learn a base classifier h_t based on \mathcal{D}

4: **end for**

5: Step 2: Construct new data sets from \mathcal{D}

6: **for** $i \leftarrow 1$ to m **do**

7: Construct a new data set that contains $\{\mathbf{x}'_i, y_i\}$, where $\mathbf{x}'_i = \{h_1(\mathbf{x}_i), h_2(\mathbf{x}_i), \dots, h_T(\mathbf{x}_i)\}$

8: **end for**

9: Step 3: Learn a second-level classifier

10: Learn a new classifier h' based on the newly constructed data set

11: **return** $H(\mathbf{x}) = h'(h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_T(\mathbf{x}))$

3 → built T-base models
 $\mathbf{x}'_i \in \mathbb{R}^T$

References

- [1] Tang, J., S. Alelyani, and H. Liu. "Data Classification: Algorithms and Applications." Data Mining and Knowledge Discovery Series. CRC Press (2015). pp. 402-500.

49 / 49

Contents - Google Docs StackingClassifier - mlxtend Chekuri Srikan...

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

StackingClassifier

- Overview
- References
- Example 1 - Simple Stacked Classification
- Example 2 - Using Probabilities as Meta-Features
- Example 3 - Stacked Classification and GridSearch
- Example 4 - Stacking of Classifiers that Operate on Different Feature Subsets
- API
- Methods

Final prediction

$y_{\hat{g}} = h' \left(\underline{h_1(x_0)}, \underline{h_2(x_0)}, \dots, \underline{h_T(x_0)} \right)$

The algorithm can be summarized as follows (source: [1]):

Algorithm 19.7 Stacking

Input: Training data $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^m$ ($\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathcal{Y}$)
Output: An ensemble classifier H

- 1: Step 1: Learn first-level classifiers
- 2: **for** $t \leftarrow 1$ to T **do**
- 3: Learn a base classifier h_t based on \mathcal{D}
- 4: **end for**
- 5: Step 2: Construct new data sets from \mathcal{D}
- 6: **for** $i \leftarrow 1$ to m **do**
- 7: Construct a new data set that contains $\{\mathbf{x}'_i, y_i\}$, where $\mathbf{x}'_i = \{h_1(\mathbf{x}_i), h_2(\mathbf{x}_i), \dots, h_T(\mathbf{x}_i)\}$
- 8: **end for**
- 9: Step 3: Learn a second-level classifier
- 10: Learn a new classifier h' based on the newly constructed data set
- 11: **return** $H(\mathbf{x}) = h'(h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_T(\mathbf{x}))$

References

- [1] Tang, J., S. Alelyani, and H. Liu. "Data Classification: Algorithms and Applications." Data Mining and Knowledge Discovery Series. CRC Press (2015). pp. 402-500.

Contents - Google Docs StackingClassifier - mlxtend Installation - mlxtend Chekuri Srikan...

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

StackingClassifier

Overview References Example 1 - Simple Stacked Classification Example 2 - Using Probabilities as Meta-Features Example 3 - Stacked Classification and GridSearch Example 4 - Stacking of Classifiers that Operate on Different Feature Subsets API Methods

```
iris = datasets.load_iris()
X, y = iris.data[:, 1:3], iris.target
```

```
from sklearn import model_selection
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.ensemble import RandomForestClassifier
from mlxtend.classifier import StackingClassifier
import numpy as np

clf1 = KNeighborsClassifier(n_neighbors=1)
clf2 = RandomForestClassifier(random_state=1)
clf3 = GaussianNB()
lr = LogisticRegression()
scf = StackingClassifier(classifiers=[clf1, clf2, clf3],
                         meta_classifier=lr)

print('3-fold cross validation:\n')

for clf, label in zip([clf1, clf2, clf3, scf],
                      ['KNN',
                       'Random Forest',
                       'Naive Bayes',
                       'StackingClassifier']):

    scores = model_selection.cross_val_score(clf, X, y,
                                              cv=3, scoring='accuracy')
    print("Accuracy: %0.2f (+/- %0.2f) [%s]"
          % (scores.mean(), scores.std(), label))
```

3-fold cross validation:

$h_1(x) = \text{clf1} \rightarrow \text{KNN}$

$h_2(x) = \text{clf2} \rightarrow \text{RF}$

$h_3(x) = \text{clf3} \rightarrow \text{GNB}$

$h^*(x) = \text{LR} \Rightarrow \text{LR}$

Contents - Google Docs StackingClassifier - mlxtend Installation - mlxtend Chekuri Srikan...

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

StackingClassifier

An ensemble-learning meta-classifier for stacking.

from mlxtend.classifier import StackingClassifier

Overview

Stacking is an ensemble learning technique to combine multiple classification models via a meta-classifier. The individual classification models are trained based on the complete training set; then, the meta-classifier is fitted based on the outputs -- meta-features -- of the individual classification models in the ensemble. The meta-classifier can either be trained on the predicted class labels or probabilities from the ensemble.

Team!

① $h_1(x), h_2(x), h_3(x)$
more different
the base models,
the better it is

② Kaggle: Stacking is must

③ real-world!

The diagram illustrates the Stacking Classifier process. It starts with a 'Training set' which is split into 'training' and 'New data'. The 'training' part is used to train 'Classification models' labeled C_1, C_2, \dots, C_m . These models produce 'Predictions' labeled P_1, P_2, \dots, P_m . A red circle labeled 'Team!' points to the 'Classification models' and 'Predictions' stages. A red arrow points from the 'StackingClassifier' heading to the 'Classification models' stage. Three red circles numbered 1, 2, and 3 point to the text annotations on the right.

Contents - Google Docs StackingClassifier - mlxtend Installation - mlxtend Chekuri Srikan...

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

```
from mlxtend.classifier import StackingClassifier
```

StackingClassifier

- Overview
- References
- Example 1 - Simple Stacked Classification
- Example 2 - Using Probabilities as Meta-Features
- Example 3 - Stacked Classification and GridSearch
- Example 4 - Stacking of Classifiers that Operate on Different Feature Subsets
- API
- Methods

Overview

Stacking is an ensemble learning technique to combine multiple classification models via a meta-classifier. The individual classification models are trained based on the complete training set; then, the meta-classifier is fitted based on the outputs -- meta-features -- of the individual classification models in the ensemble. The meta-classifier can either be trained on the predicted class labels or probabilities from the ensemble.

Diagram of Stacking Process:

```

graph TD
    TS[Training set] --> CM[Classification models]
    CM --> P[Predictions]
    P --> MC[Meta-Classifier]
    MC --> FP[Final prediction]
    ND[New data] --> CM
  
```

Handwritten Annotations:

- A large oval on the left contains handwritten text: "1230", "0.1%", and "log loss". Below this is a bracketed note: "1st: 0.1230" and "2nd: 0.1231".
- A large oval on the right contains handwritten text: "Stacking being used less often than boosting (GB) & bagging (RF)". Above this is the handwritten note: "competition → Kaggler".
- A small oval at the bottom right contains handwritten text: "0.1%".

Contents - Google Docs StackingClassifier - mlxtend Installation - mlxtend Chekuri Srikan...

Secure | https://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

mlxtend Home User Guide API Installation About Search Previous Next GitHub

```
from mlxtend.classifier import StackingClassifier
```

StackingClassifier

- Overview
- References
- Example 1 - Simple Stacked Classification
- Example 2 - Using Probabilities as Meta-Features
- Example 3 - Stacked Classification and GridSearch
- Example 4 - Stacking of Classifiers that Operate on Different Feature Subsets
- API
- Methods

Overview

Stacking is an ensemble learning technique to combine multiple classification models via a meta-classifier. The individual classification models are trained based on the complete training set; then, the meta-classifier is fitted based on the outputs -- meta-features -- of the individual classification models in the ensemble. The meta-classifier can either be trained on the predicted class labels or probabilities from the ensemble.

Diagram of Stacking Process:

The diagram shows a flow from a 'Training set' through 'Classification models' (C_1, C_2, \dots, C_m) to 'Predictions' (P_1, P_2, \dots, P_m). These predictions are then processed by a 'Meta-Classifier' to produce the 'Final prediction' (P_f). A 'New data' box is shown pointing to the 'Classification models'.

Handwritten Annotations:

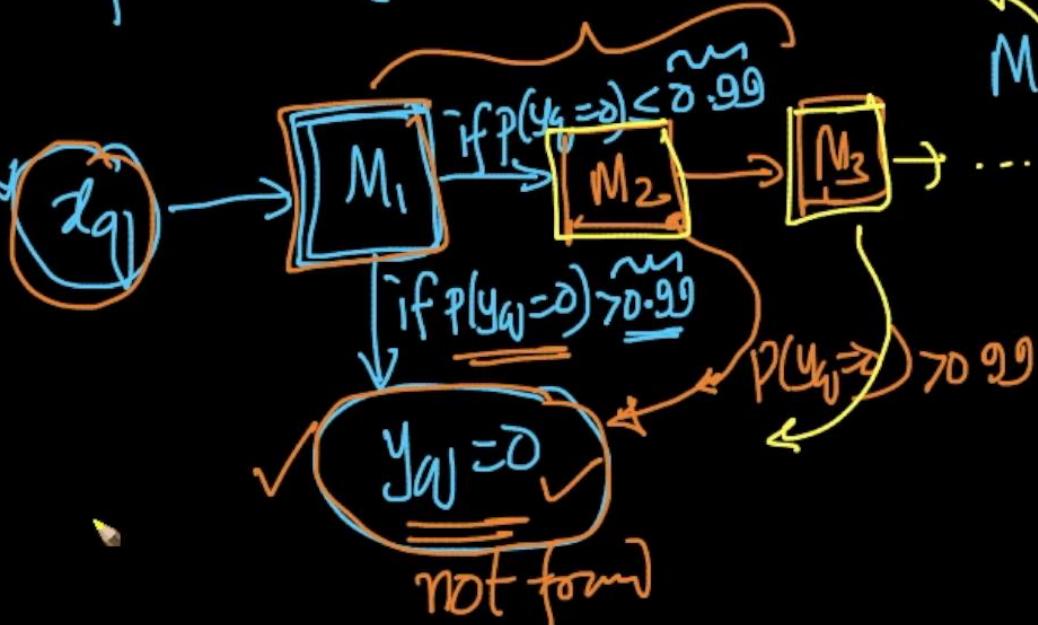
- A large oval on the left contains handwritten numbers: 1, 2, 3, 0 and 0, 1, %.
- A large oval on the right contains handwritten text: 'Stacking being used less often than boosting (GB) & bagging (RF)'.
- Other handwritten notes include: 'log loss' crossed out, '1st: 0.1230', '2nd: 0.1231', and '0.1%'.
- A blue arrow points from the text 'competition → Kaggle' to the handwritten note on the right.

"Cascading" models:

→ Bagging, boosting, Stacking

→ e.g.: - predicts {credit card transaction} is fraudulent or not

Transaction



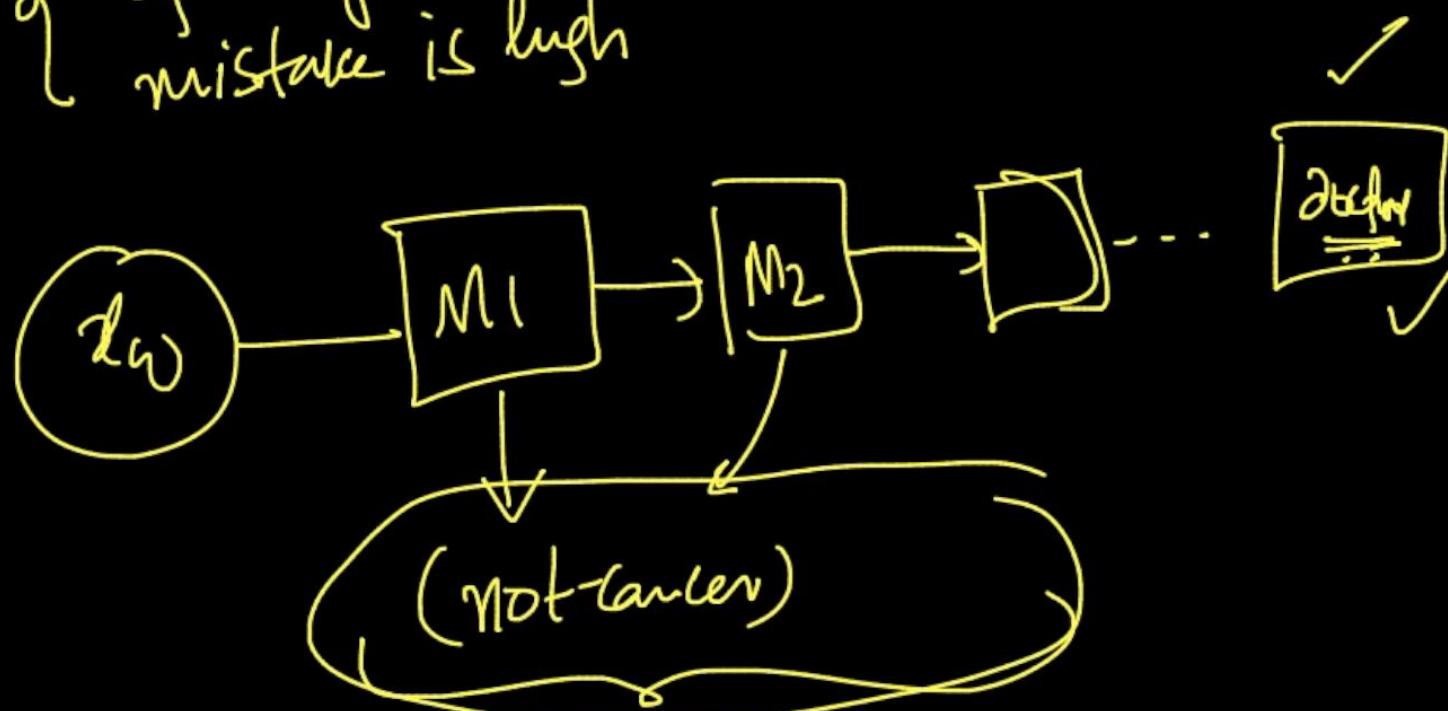
$$\begin{aligned} M_1: & p(y_{qj}=0) \\ & p(y_{qj}=1) \end{aligned}$$

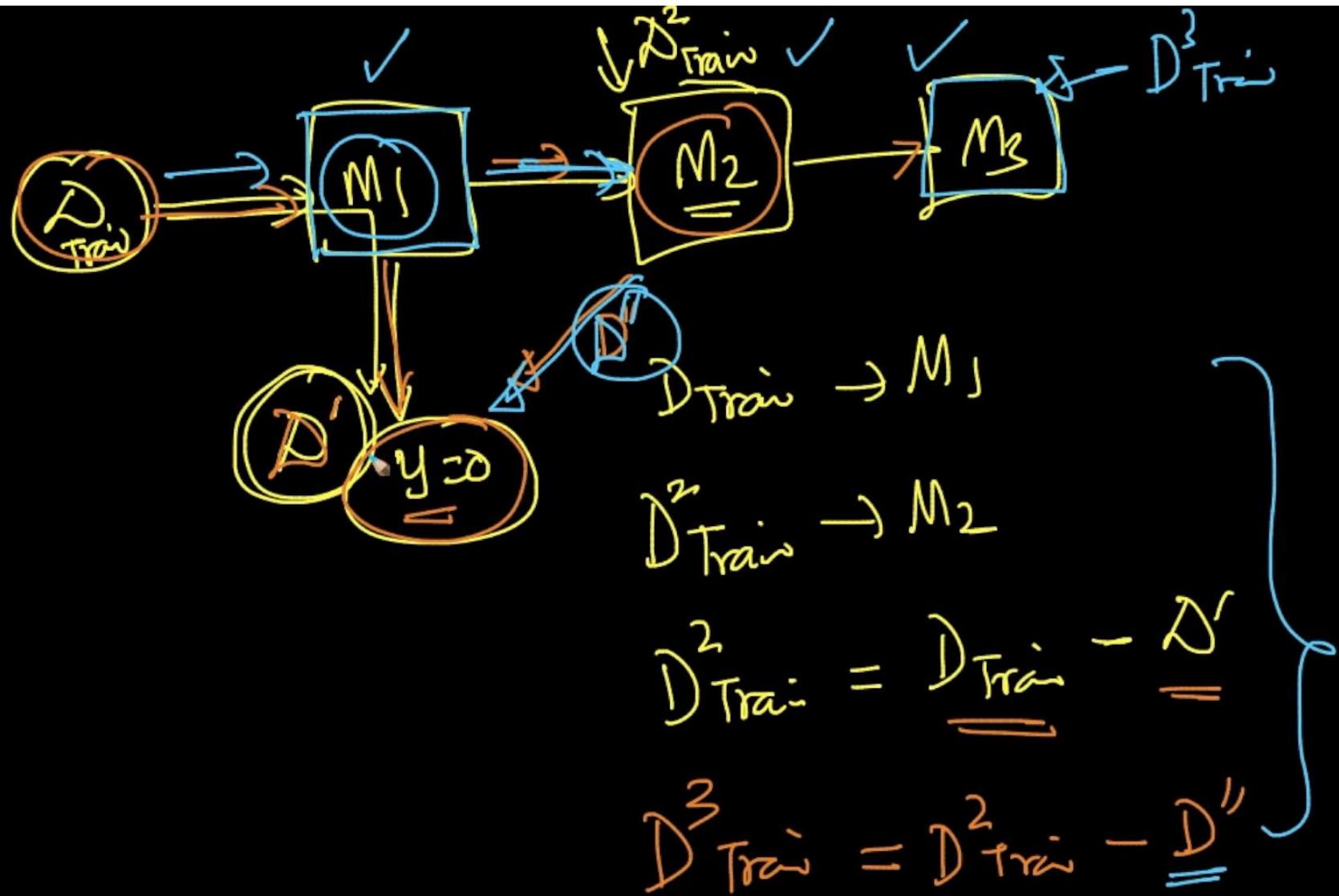
10-models

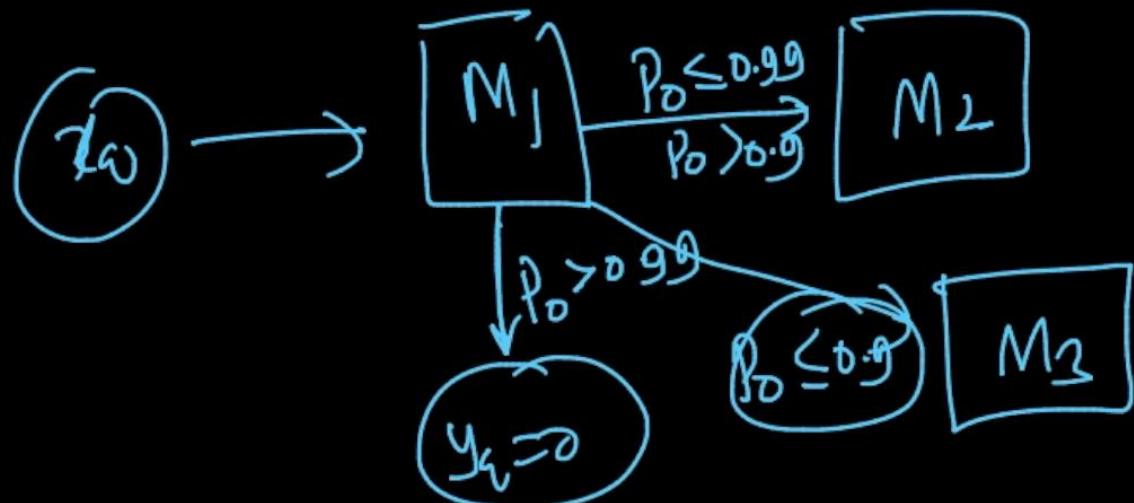
{
0 → not-fraud
1 → fraud
}

(Cascade-model)

{ Typically used when the cost of making a mistake is high







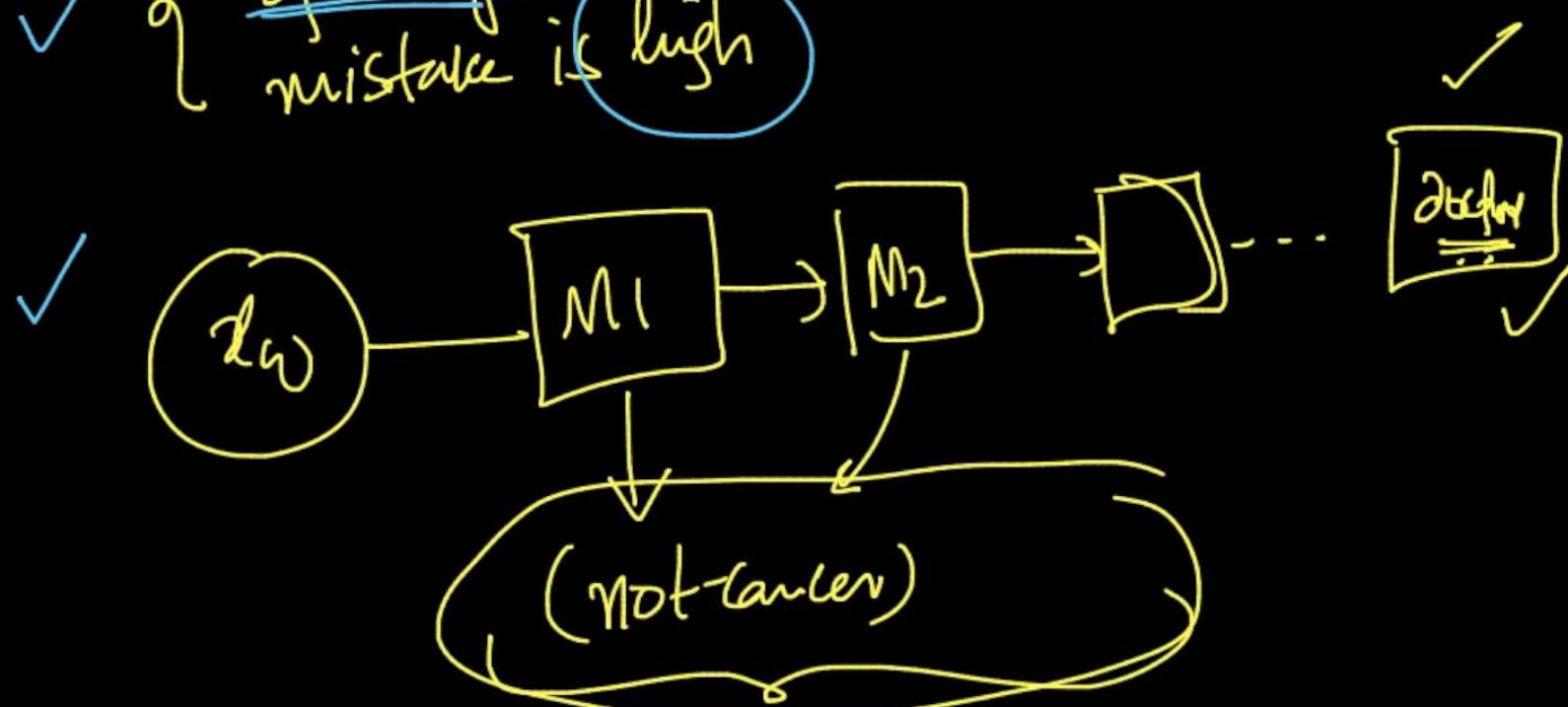
$$P_D = \beta(y_n = 0)$$

$$\beta_1 = \beta(y_n = 1)$$

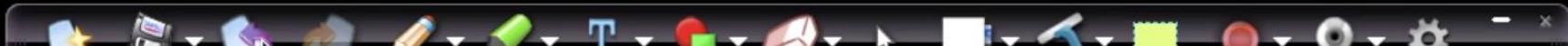
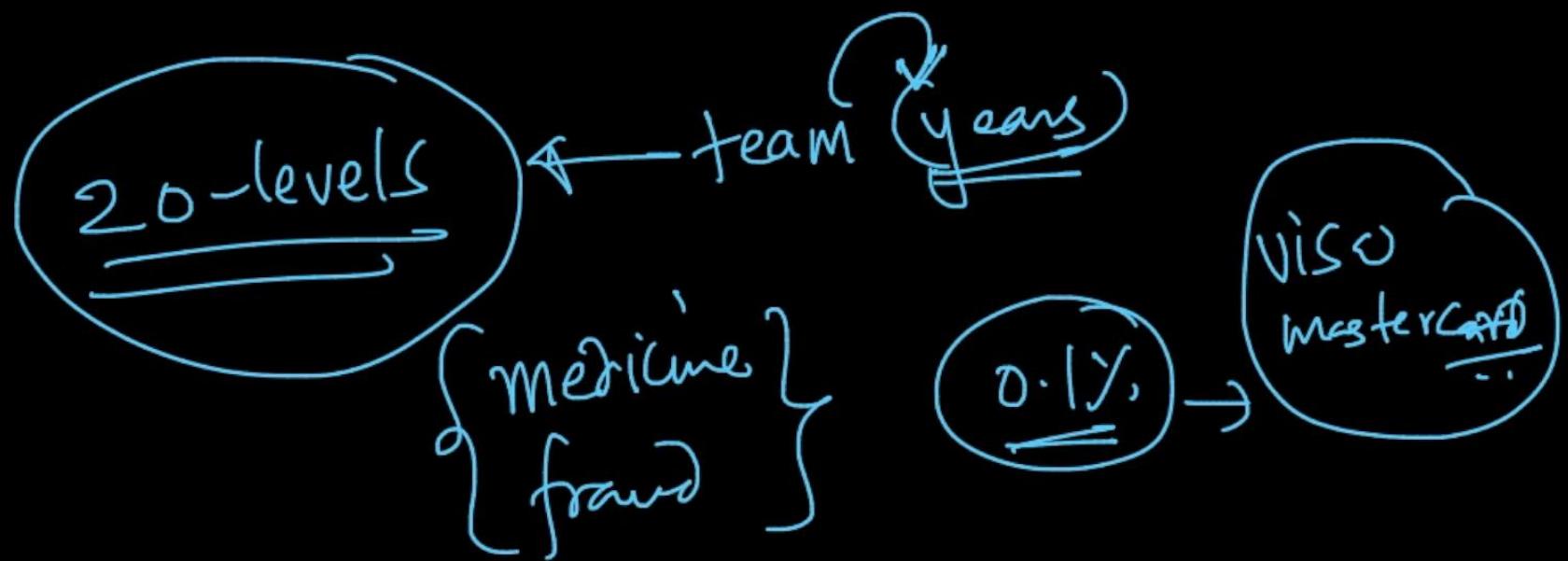
if $P_D > 0.99$
 $y_n = 0$
 else
 if $P_D > 0.9$
 use M_2
 else
 use M_3

(Cascade-model)

✓ { Typically used when the cost of making a mistake is high



Fraud-detection



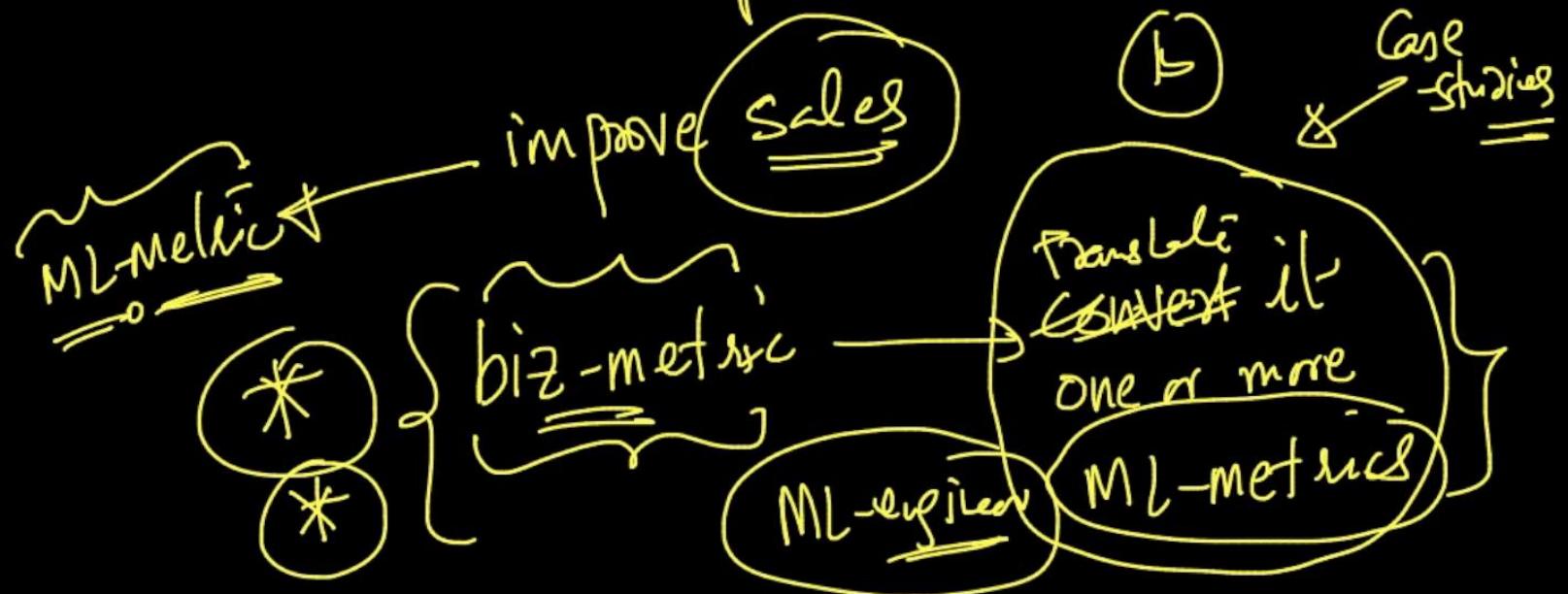
Kaggle vs Real-world

ensembles:- Kaggle, competition

- ① Kaggle:- "one" metric (log-loss, f1c, F-score)
 - ↳ competition
 - ↳ objective
 - ↳ 1st
 - ↳ 2nd
 - ↳ <0.1% =
 - ↳ Leaderboard

real-world :- multiple Metric (a)

→ primary metric
→ secondary metric



②

Very complex ensembles → improve one metric

2000

impractical

low-latency
train-time
interpretability

participate
read winners solution

③

Kaggle

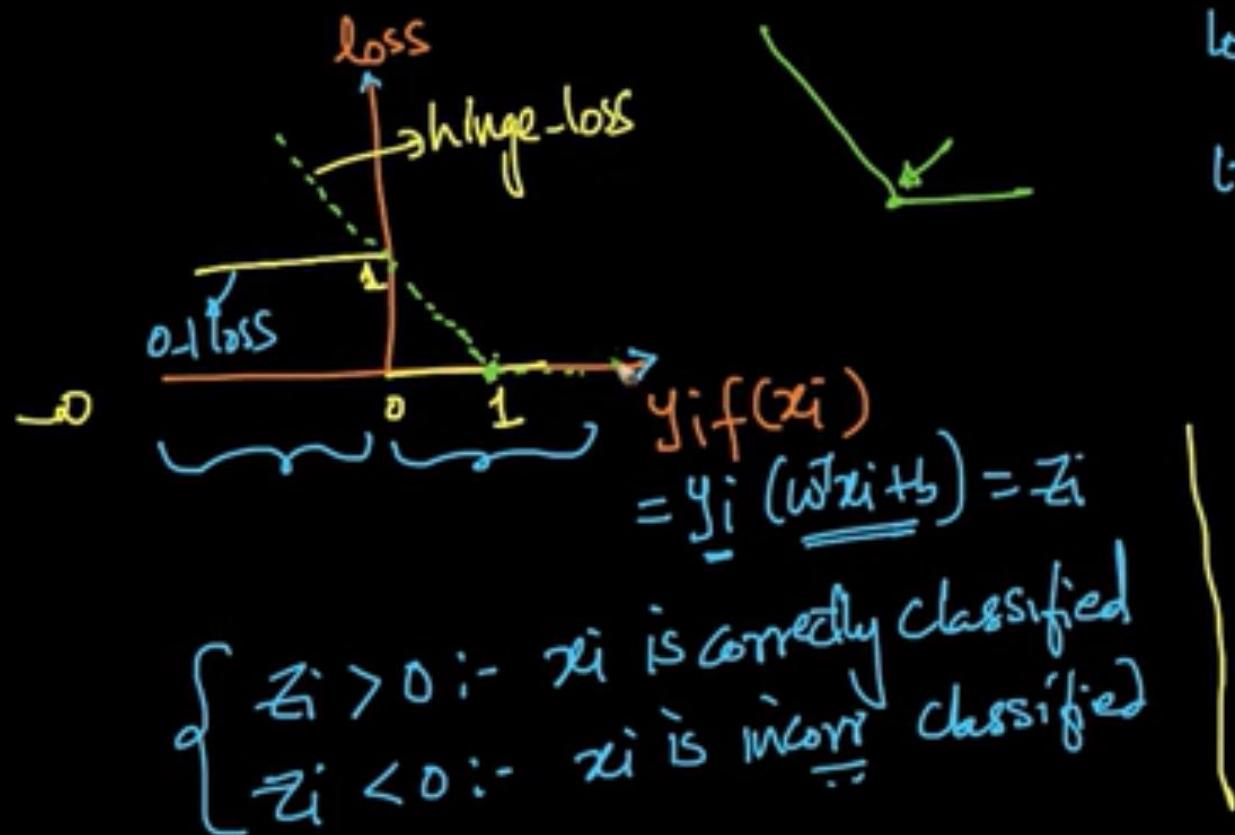
competitive

single-metric

date cleaning, pre-processing
feature-engineering, aspects of modeling



Loss-minimization = Hinge-loss



log-reg:- logistic loss + reg

lr-reg:- lr. loss + reg

SUM:- hinge-loss + reg