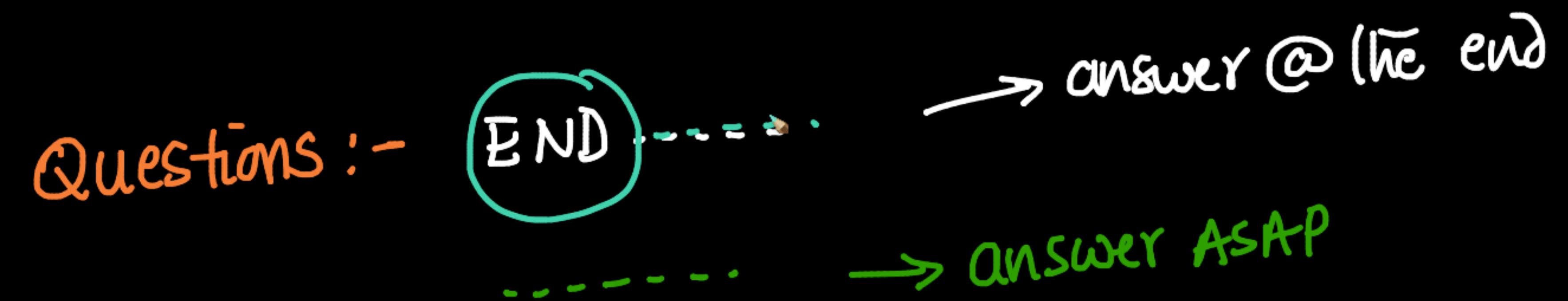


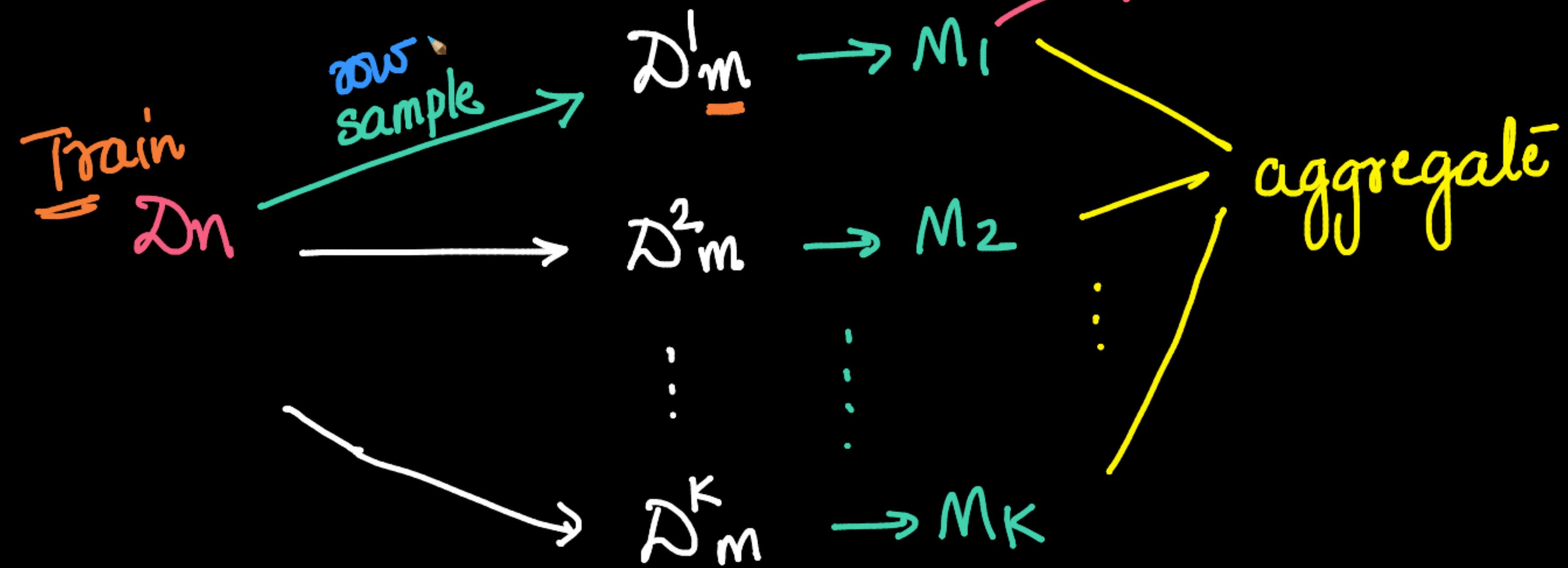
Topics:

- Random forest
- ExtraTrees
- Boosting (intuition)
- Gradient Boosting
- Code (HR-domain)



chat : interaction :-

Bagging: Bootstrapped Sampling + agg...
base-learners



let $k=100$



please
revise

$$\text{Entropy} : - \sum_{i=1}^K p(y_i) \log(p(y_i))$$

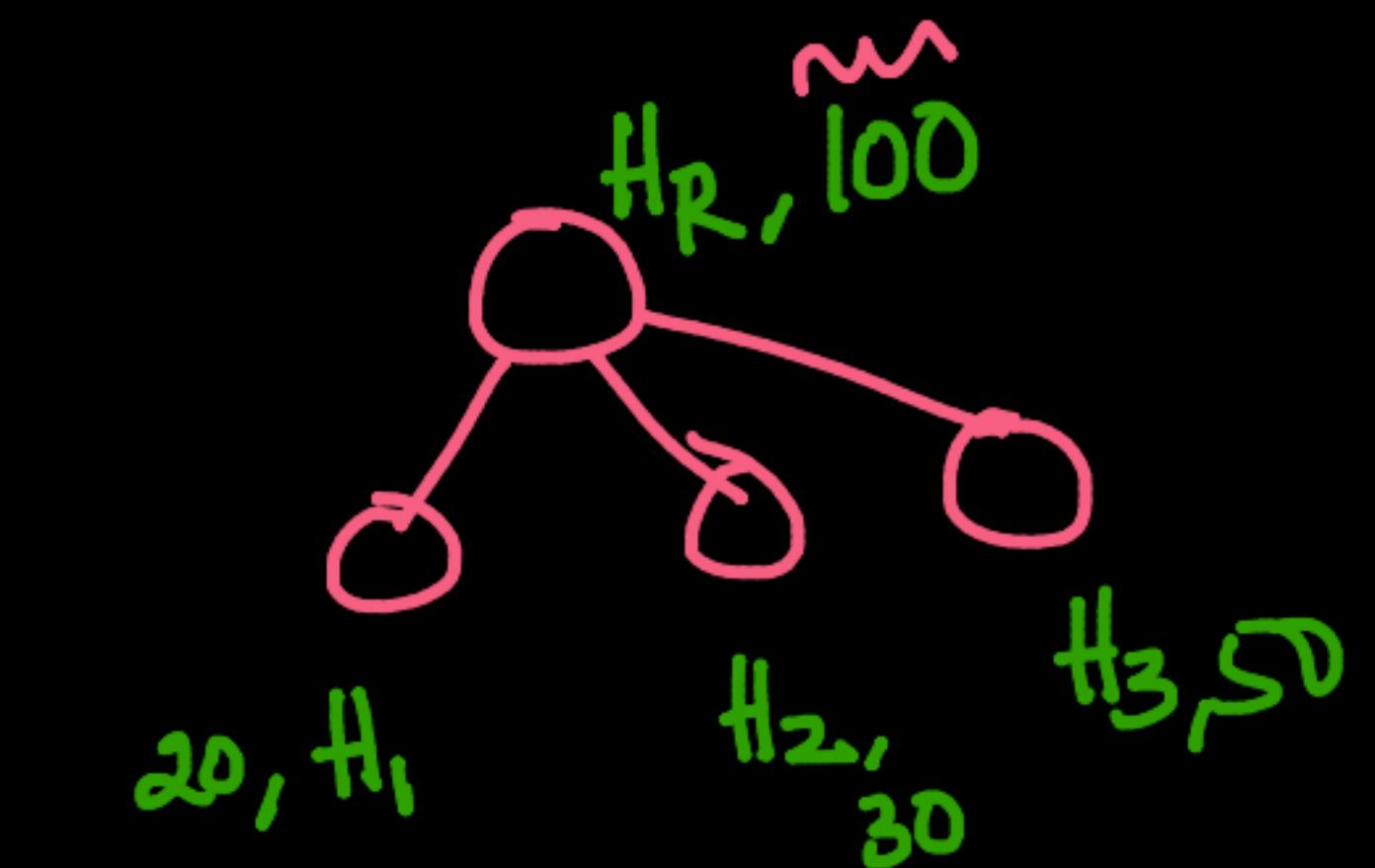
$$G.I := 1 - \sum_{i=1}^K (p(y_i))^2 \quad \checkmark$$

IG

$$\Delta \text{Entropy} \quad \text{or} \quad \Delta G.I$$

$$IG = HR$$

↑ higher



$$0.2H_1 + 0.3H_2 + 0.5H_3$$

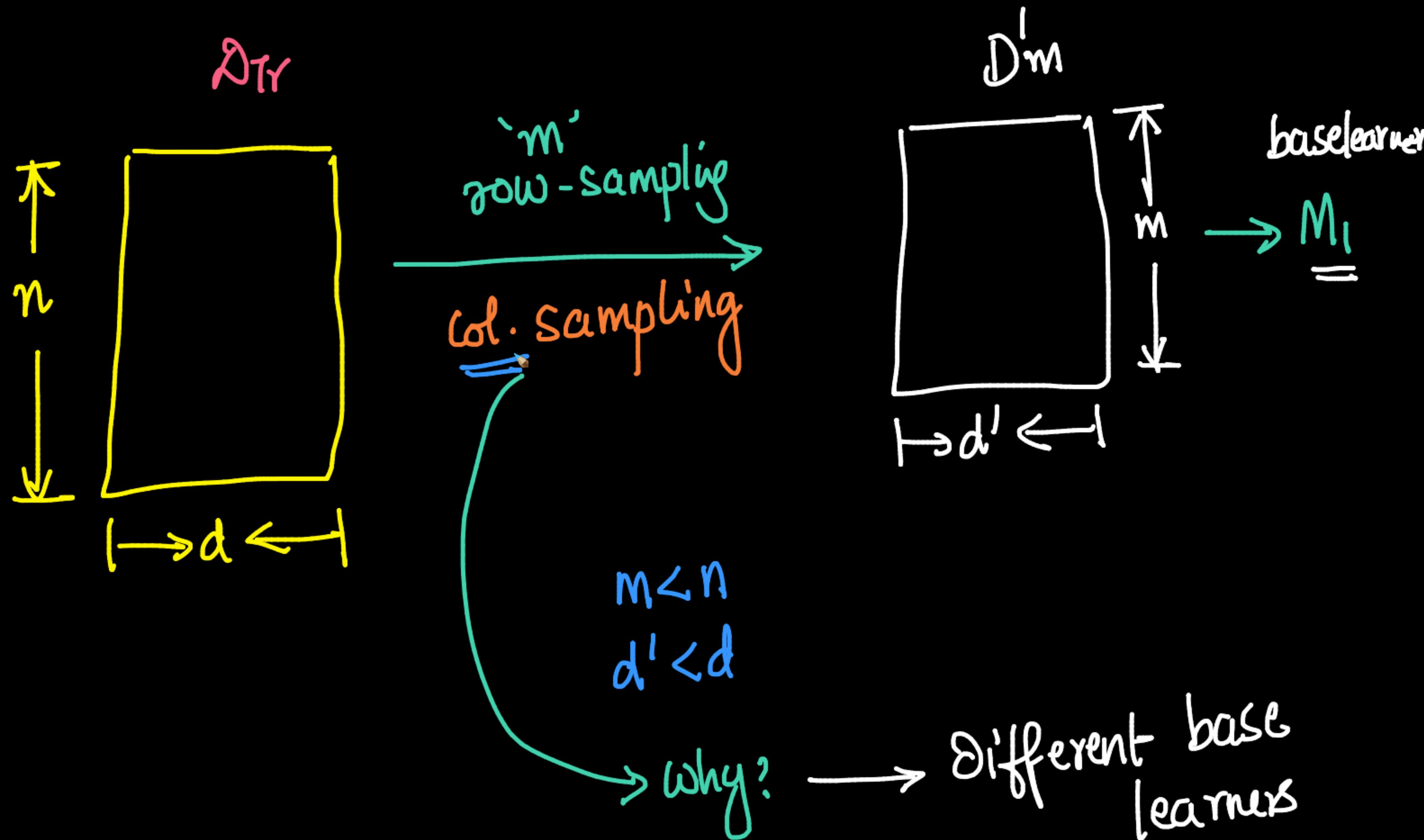
weighted H across
all child nodes

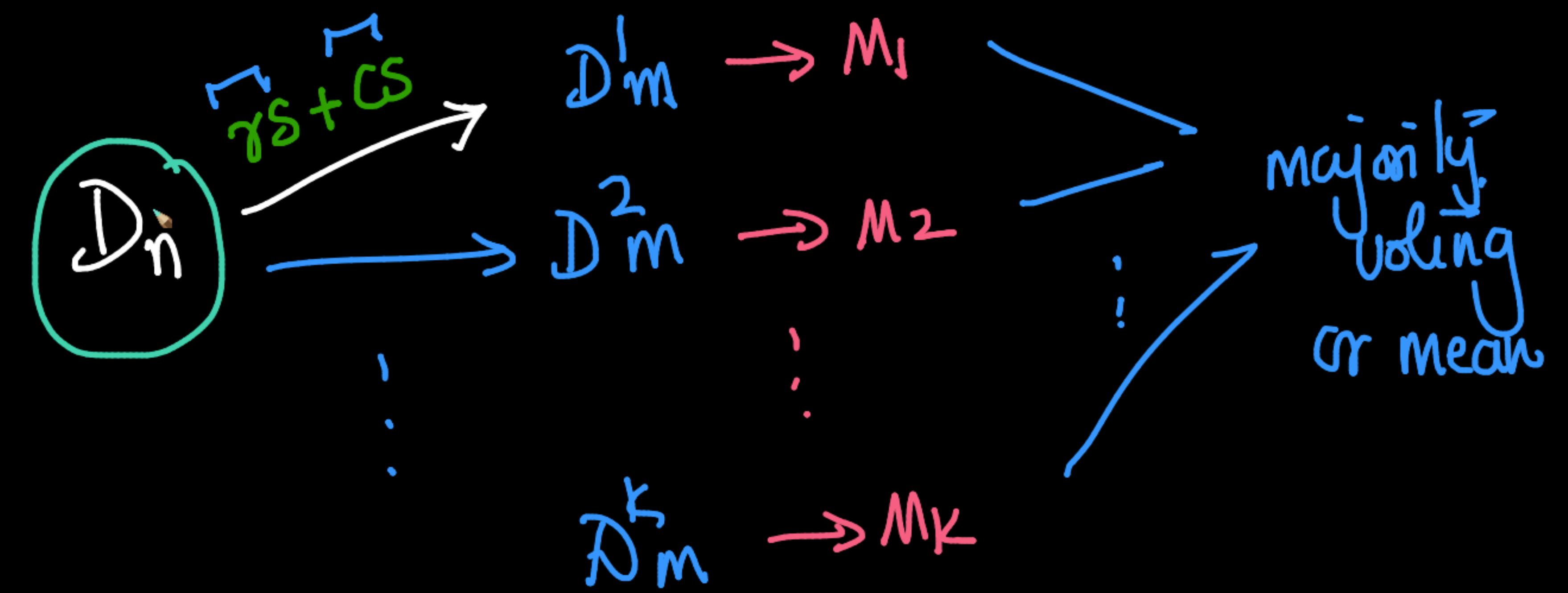


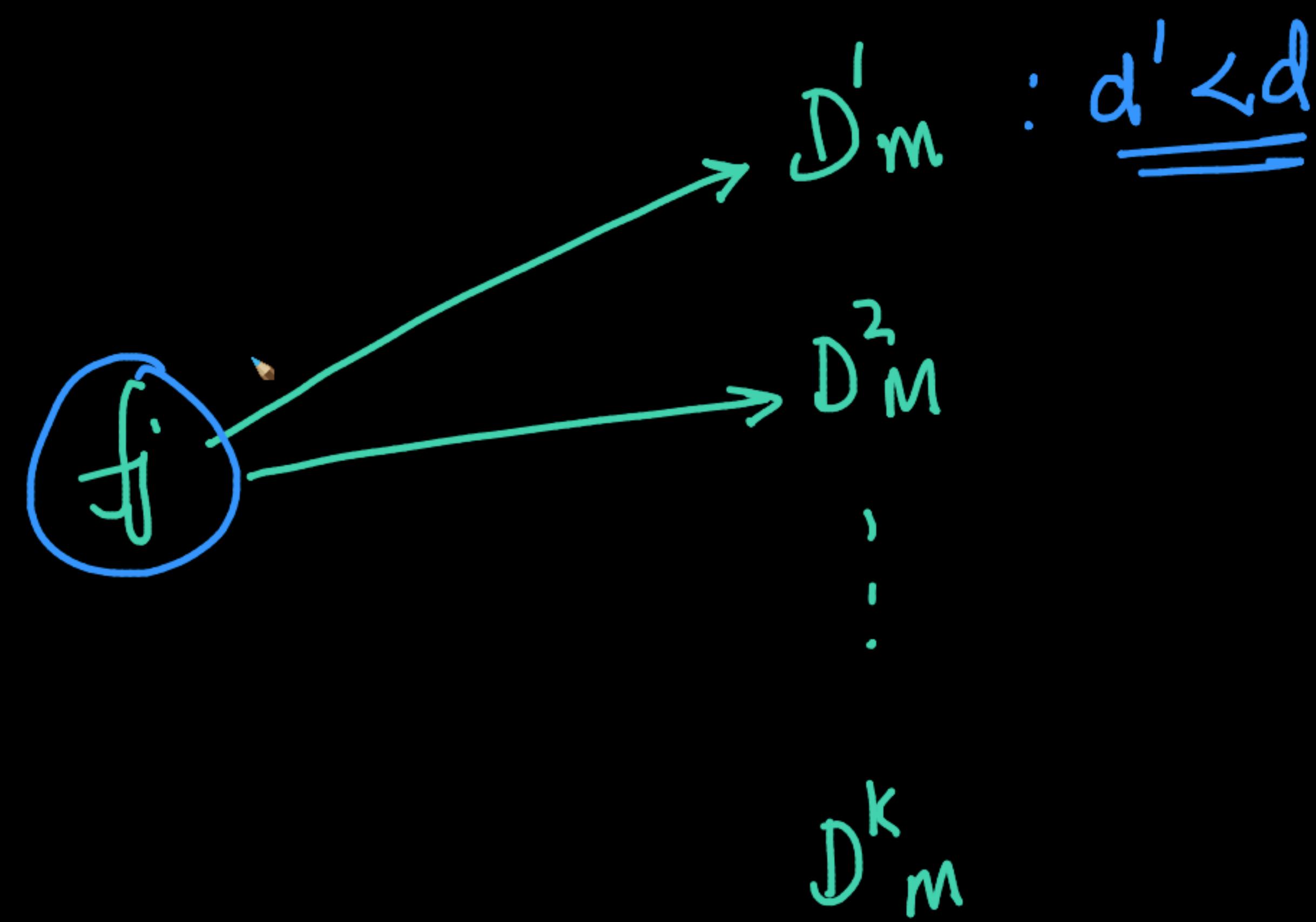
Random Forest

Bagging with DT (base-learners)

→ Random + ensemble of trees + Bagging - -



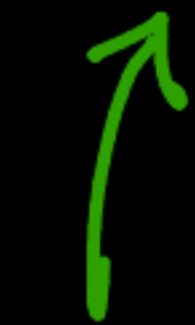




Train: $D_n \supseteq$

n - rows:

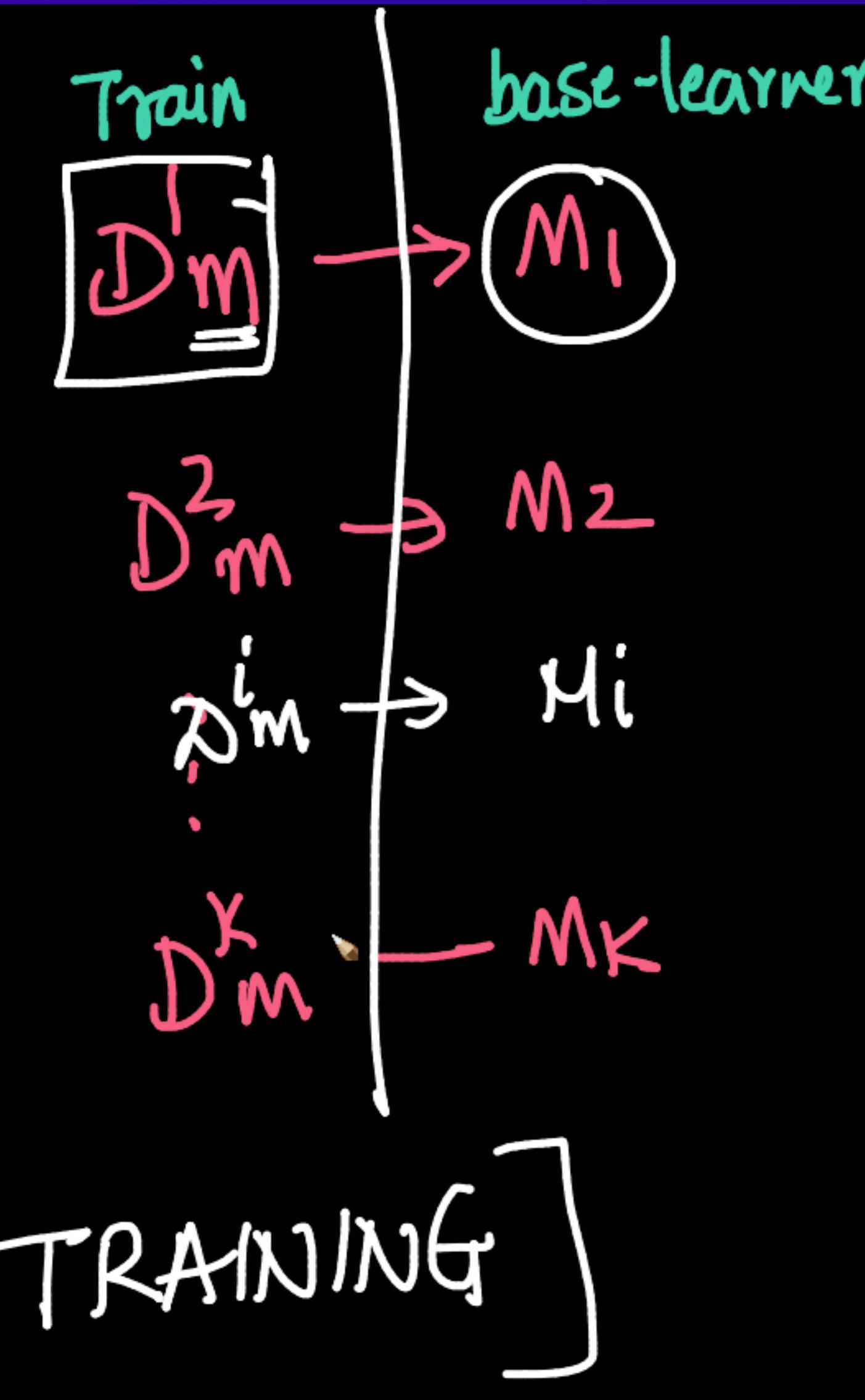
$$D_{n-m}^i = D_n - D_m^i$$



Set-difference

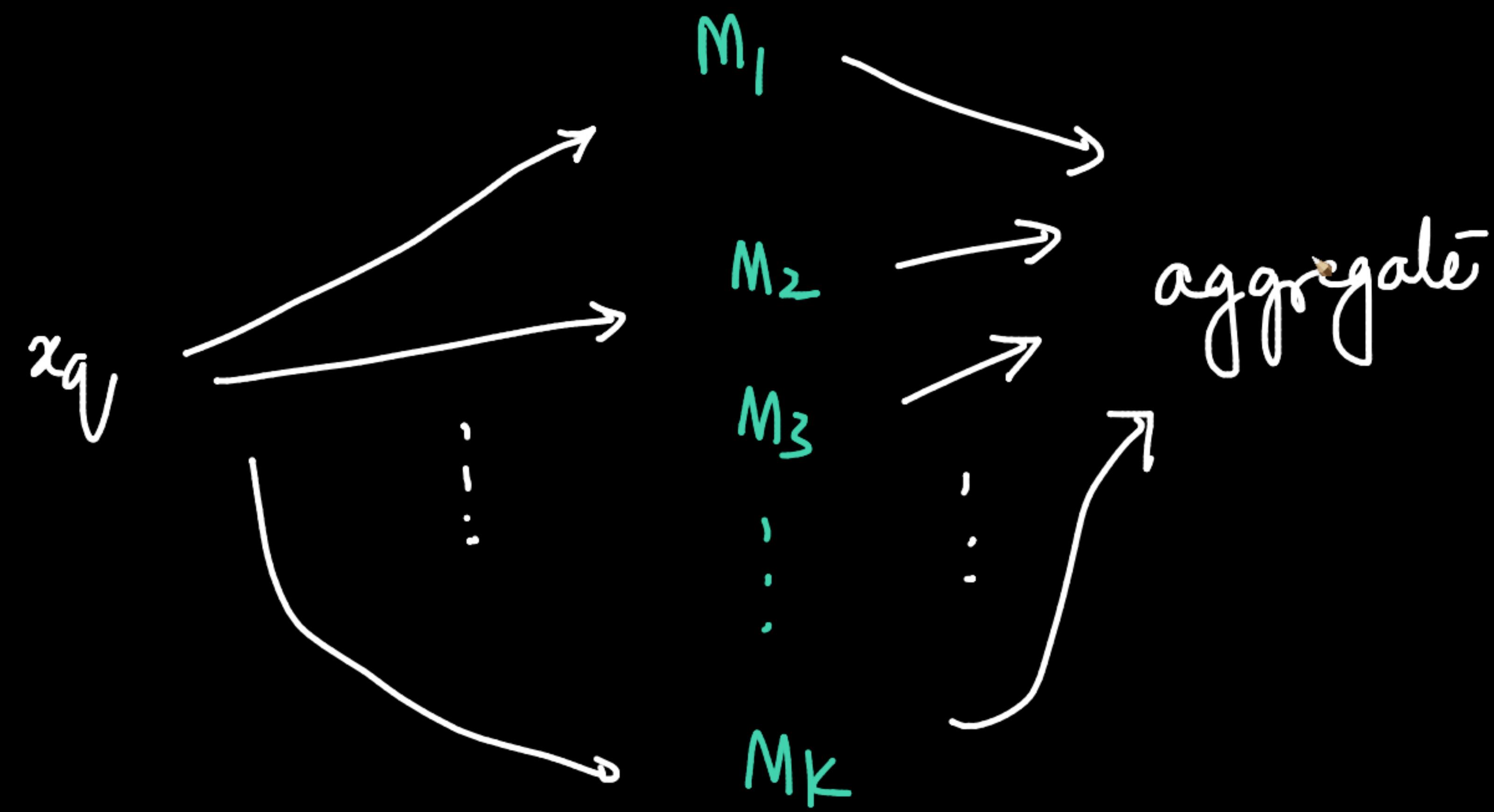
OOB CV
 $\overline{\overline{D_n}} \supseteq D_{n-m}$

$$D_{n-m}^1 \\ D_{n-m}^2 \\ \vdots \\ D_{n-m}^K$$



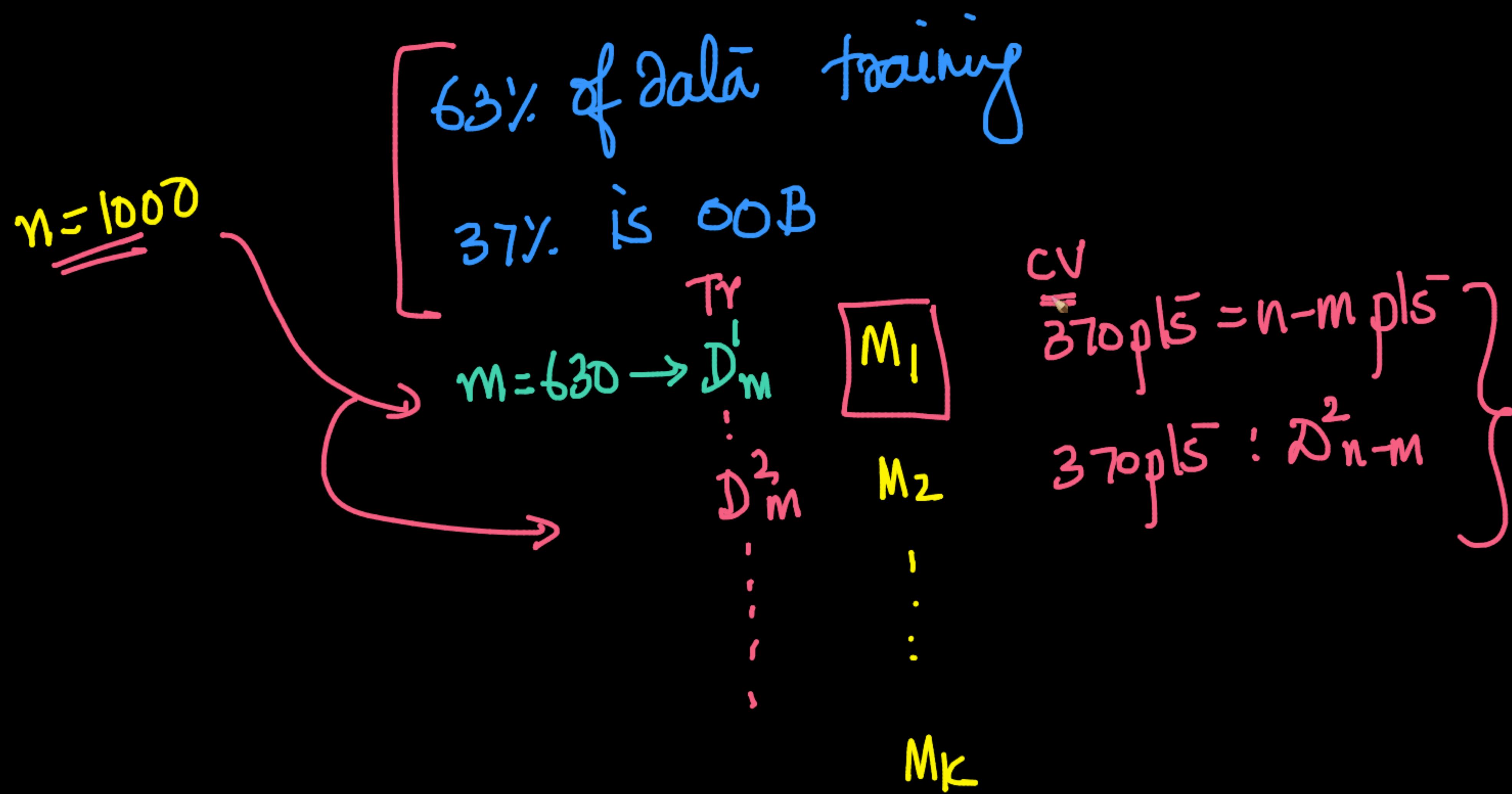
out of bag
OOB
 \downarrow
 D_{n-m}

Test-time



Feature Engineering is a must → classical ML
=

↳ skip it in DL
=

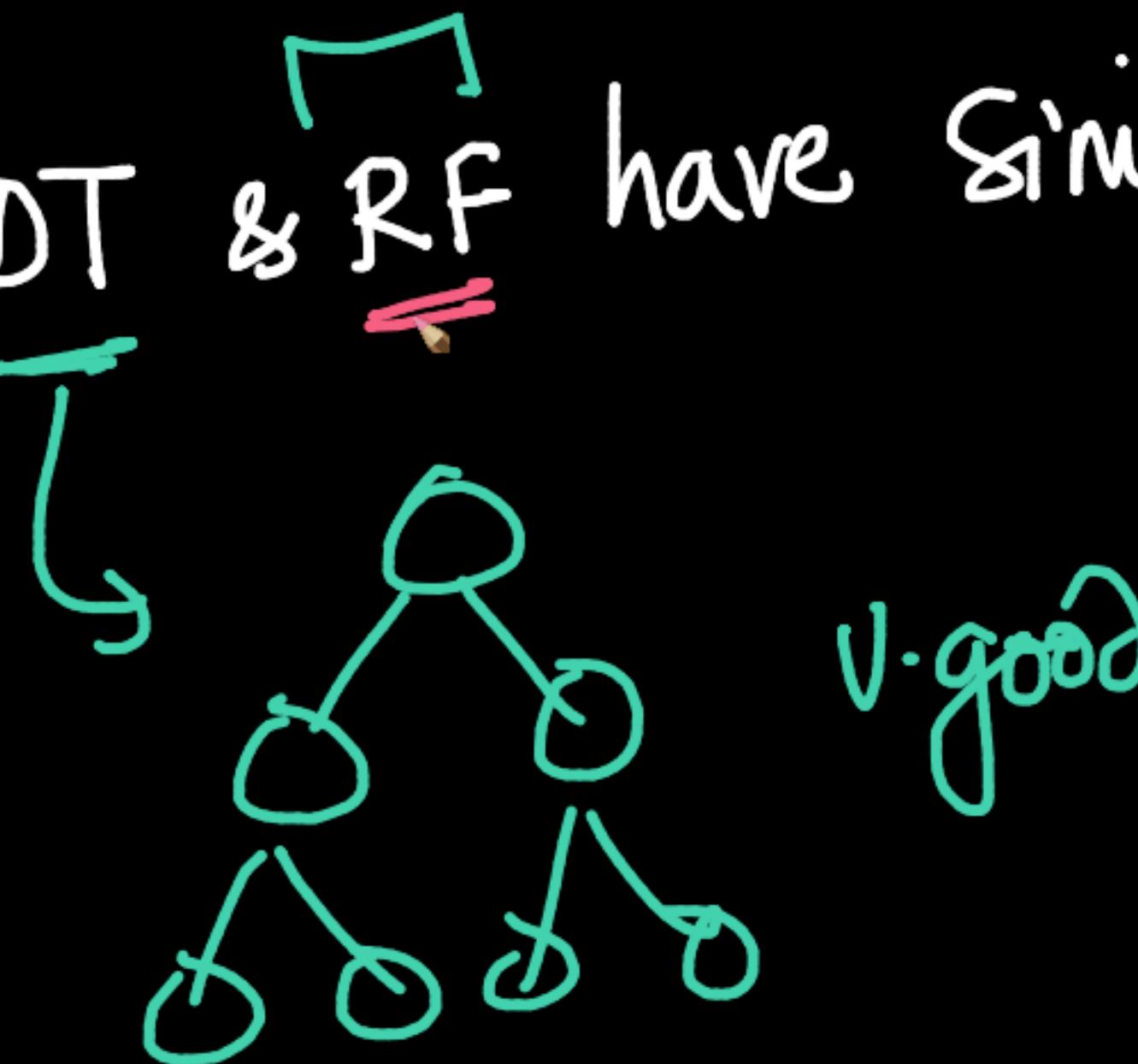


simpler-model if it works:

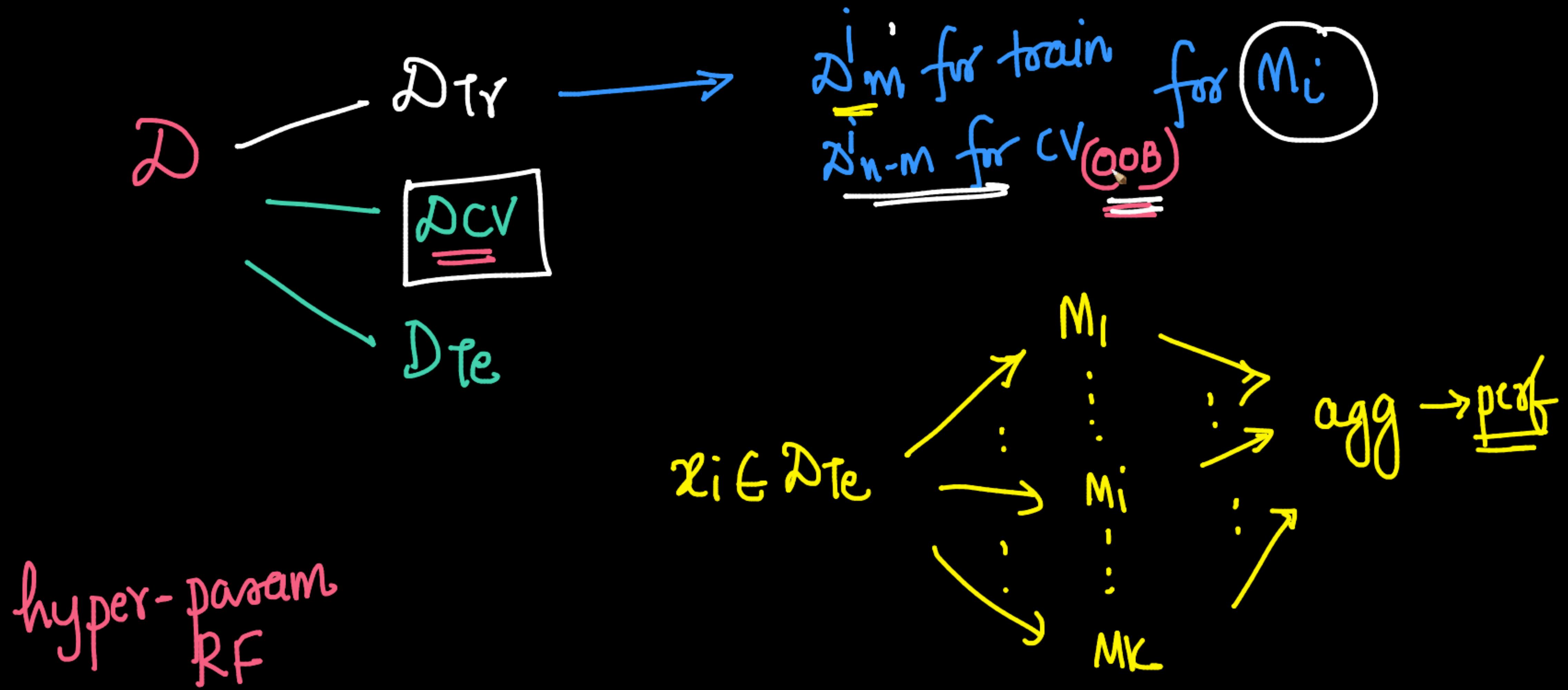
RF better than DT
(most-cases)

When do DT & RF have similar perf

Simple task:



v.good





each base-learner

M_i

is a DT

m -rows &
 d -cols

↓
deep DT (non-shallow)

→ overfit slightly

↓
Sub sample of
rows & cols

RF:  → zig-zags stat. proof
base-learners: overfit 
high-variance & low bias
on a subsample of rows & cols...


(aggregation / averaging)

Var ↓ bias →

→ reduce Varianle
without much inc in bias

Statistical ML

(Research Sessions)

$$\text{Error}_{\text{model}} = \text{Bias}^2 + \text{Var} + \text{irreducible error}$$

Typically: ~~Var~~ ↓ Bias ↑

{ Bias ↑ : underfit
Var ↑ : overfit

tradeoff

RF:-

① Bagging : RS + CS

② Base learners: overfit (high var)

③ aggregation :- Var \downarrow bias \rightarrow

↓
better overall model

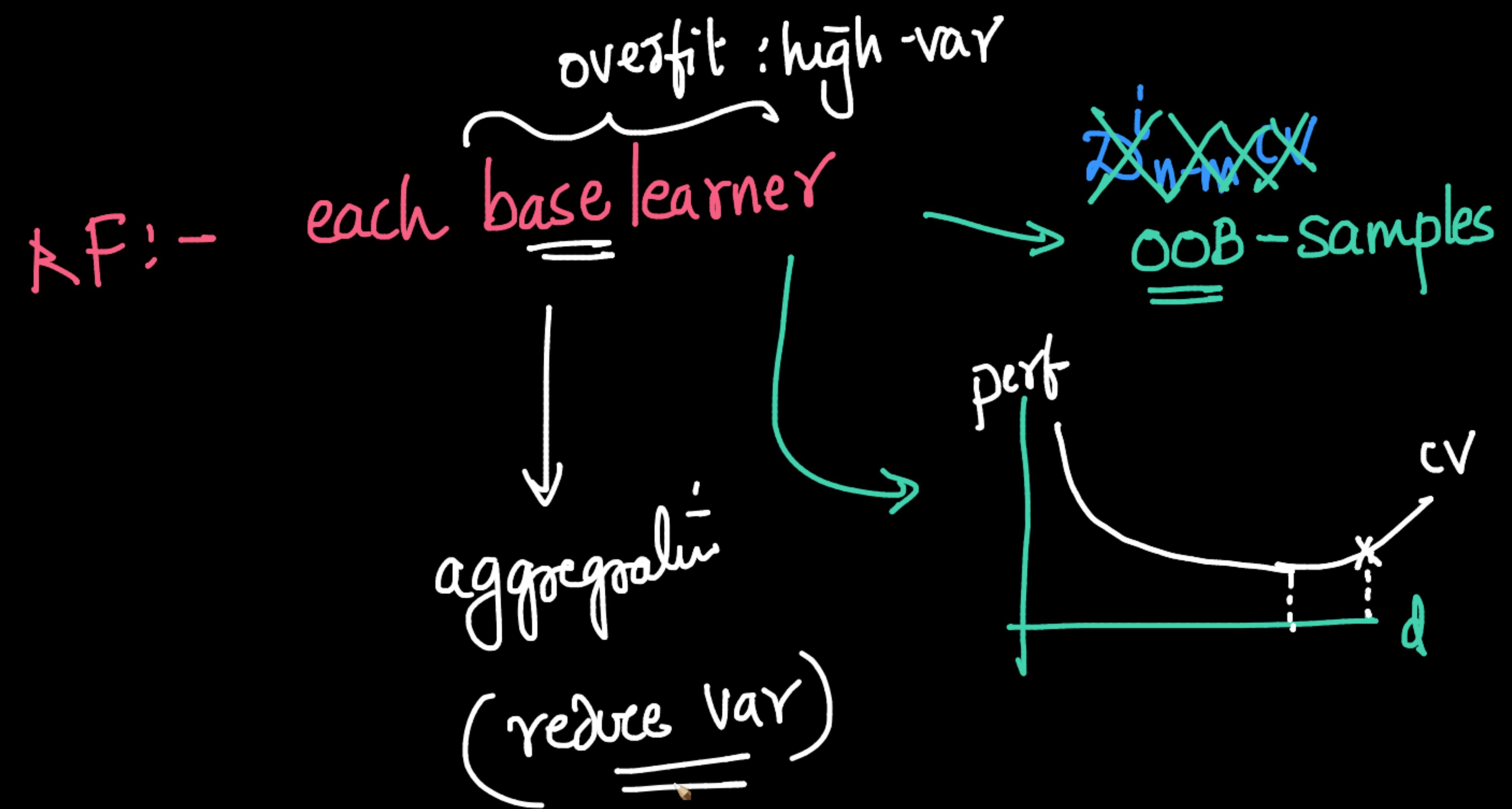
RF:

non-DT models ✓

use models with high-var (overfit)

↳ KNN with small K

↳ logistic reg with higher order
features + low λ

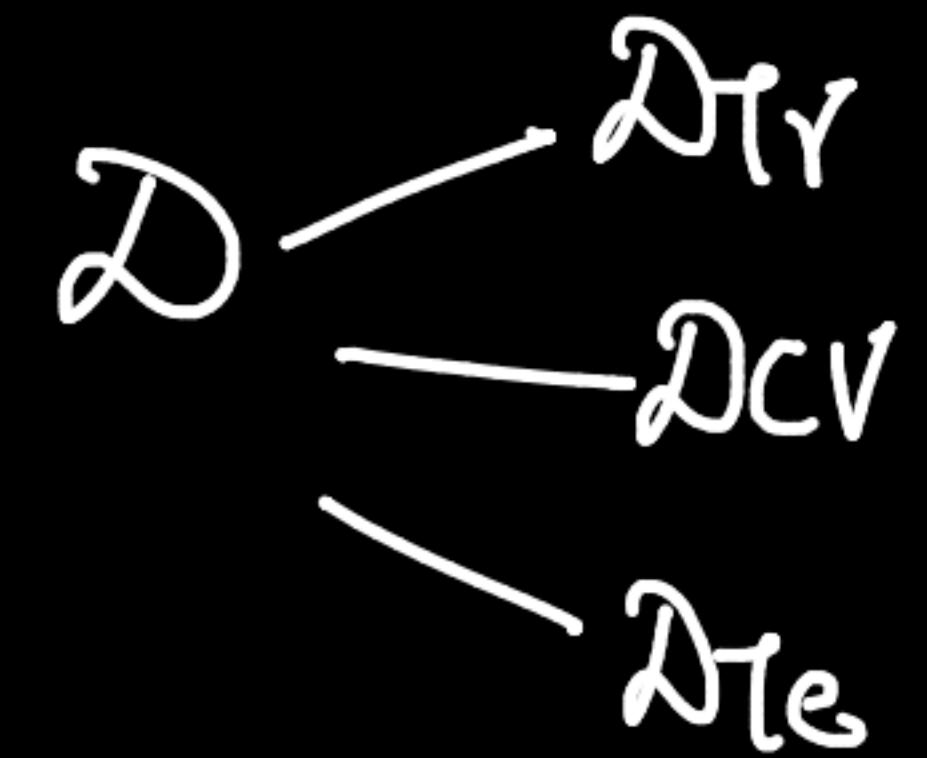


hyper-params: tuning of RF:

trees :- $K \uparrow$
(base learners)

and overfitting

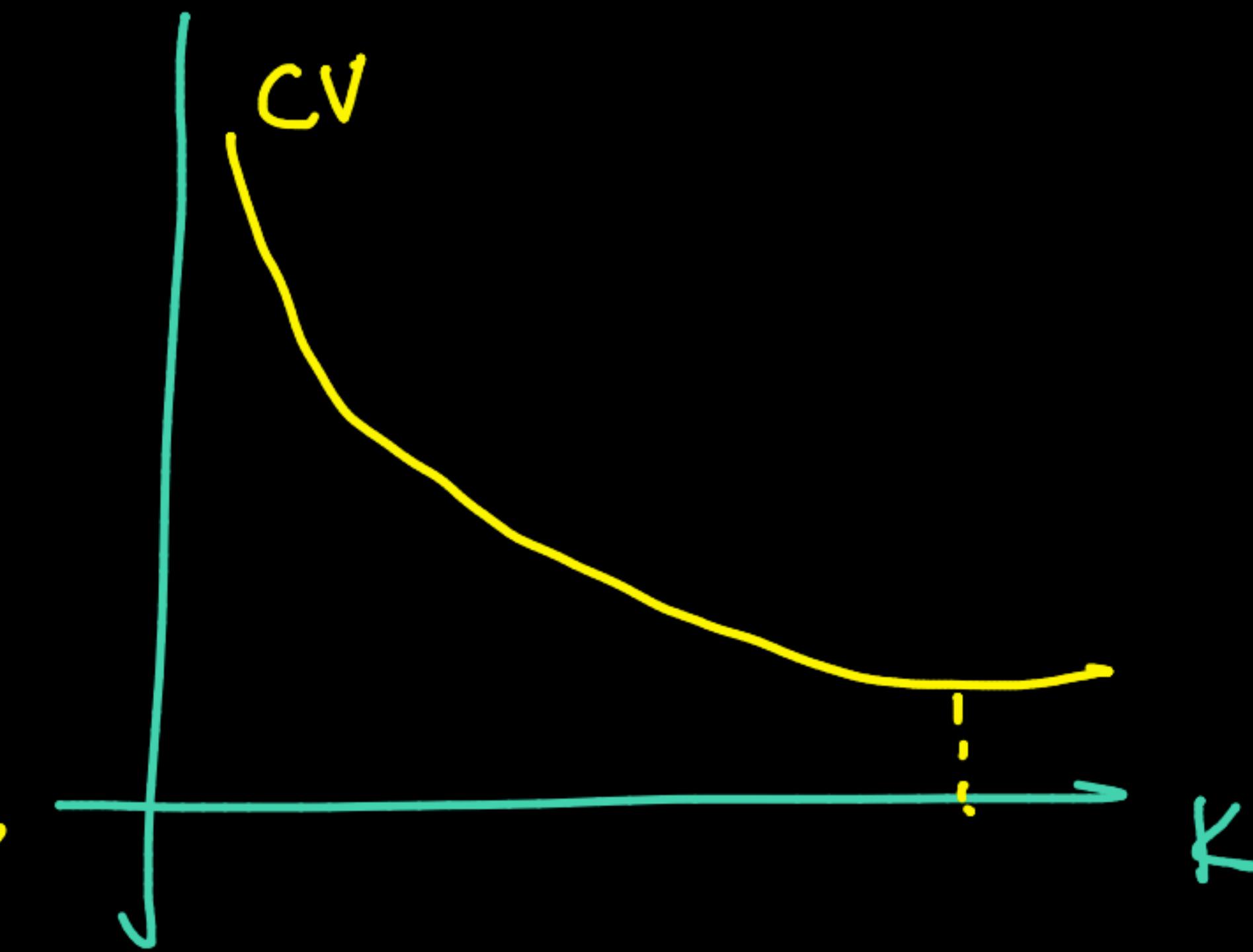
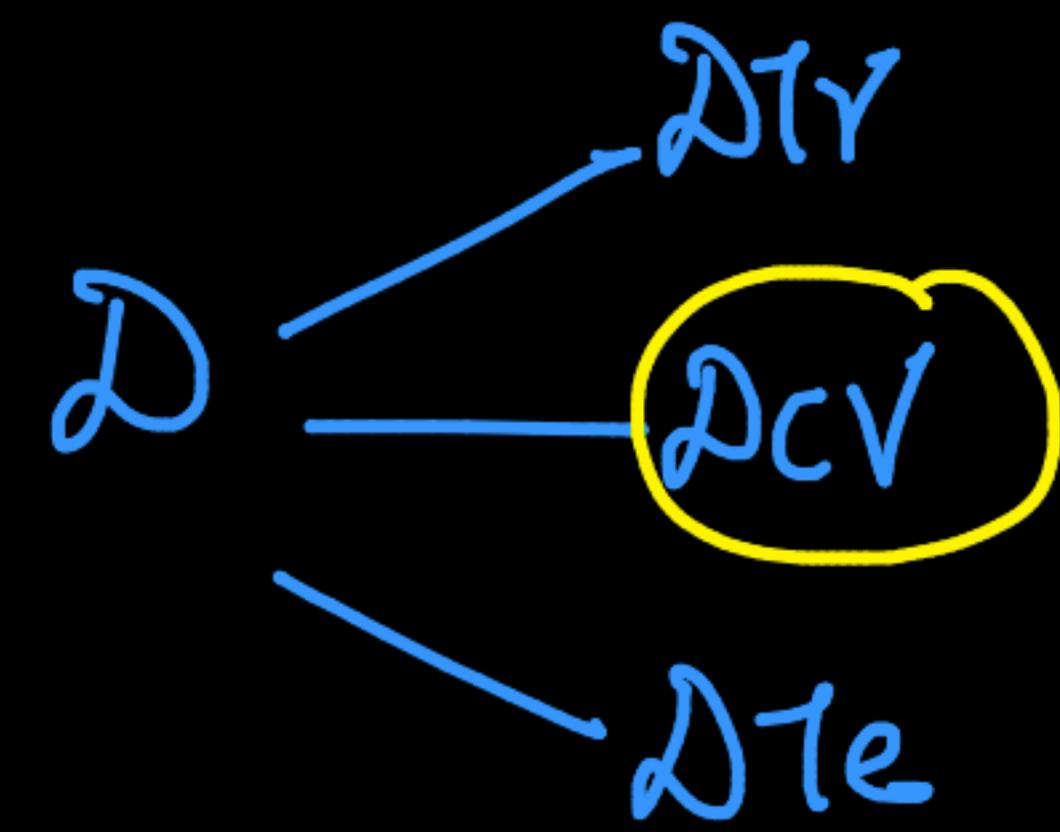
generalize better



aggregation :- $\text{Var} \downarrow$; bias \rightarrow

$k=1 \rightarrow \underline{\text{overfit}}$ slightly or massively ... (oob)

100 / 1000 $k = \text{large} \rightarrow \text{avoid } \underline{\text{overfitting}}$. $\begin{cases} 1B = n \\ m = 100 \end{cases}$



K is a hyperparameter
for RF

underfit = low var & high bias
(def 1)

base-learner :- decision Stump

RF agg: - \downarrow var + bias \rightarrow

will not work as expected

hyper-params:-

① # base learners = $k \uparrow$ avoids overfitting

② $m = \text{Sample Size}$ \rightarrow all base learners are going to overfit
 m/n \uparrow or all data → overfitting \uparrow

$m=n$

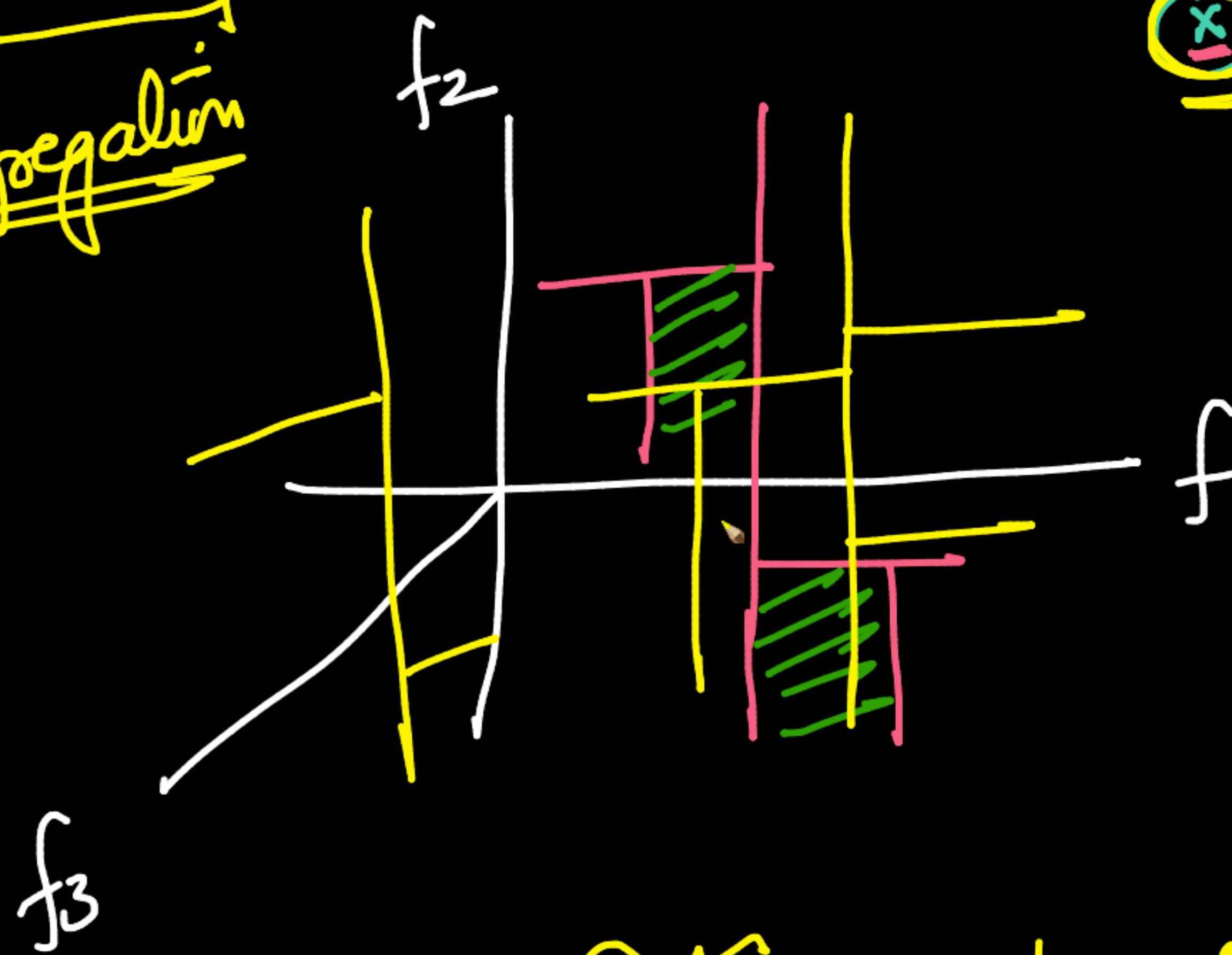
③ $d' = \# \text{cols sampled}$
 $\frac{d'}{d} \uparrow \Rightarrow \text{overfitting} \uparrow$

④ depth of base learners
depth $\uparrow \Rightarrow \text{overfit}$

$$\frac{\text{OOB Samples}}{cv} = \frac{m}{n}$$

depth $\uparrow \rightarrow \frac{m}{n} \downarrow \quad \frac{d'}{d} \downarrow \quad K \uparrow$ less chance of overfit

~~aggregation~~



≈ 1000 base-learners

~~x~~ outlier
↓
seen only some base-learners

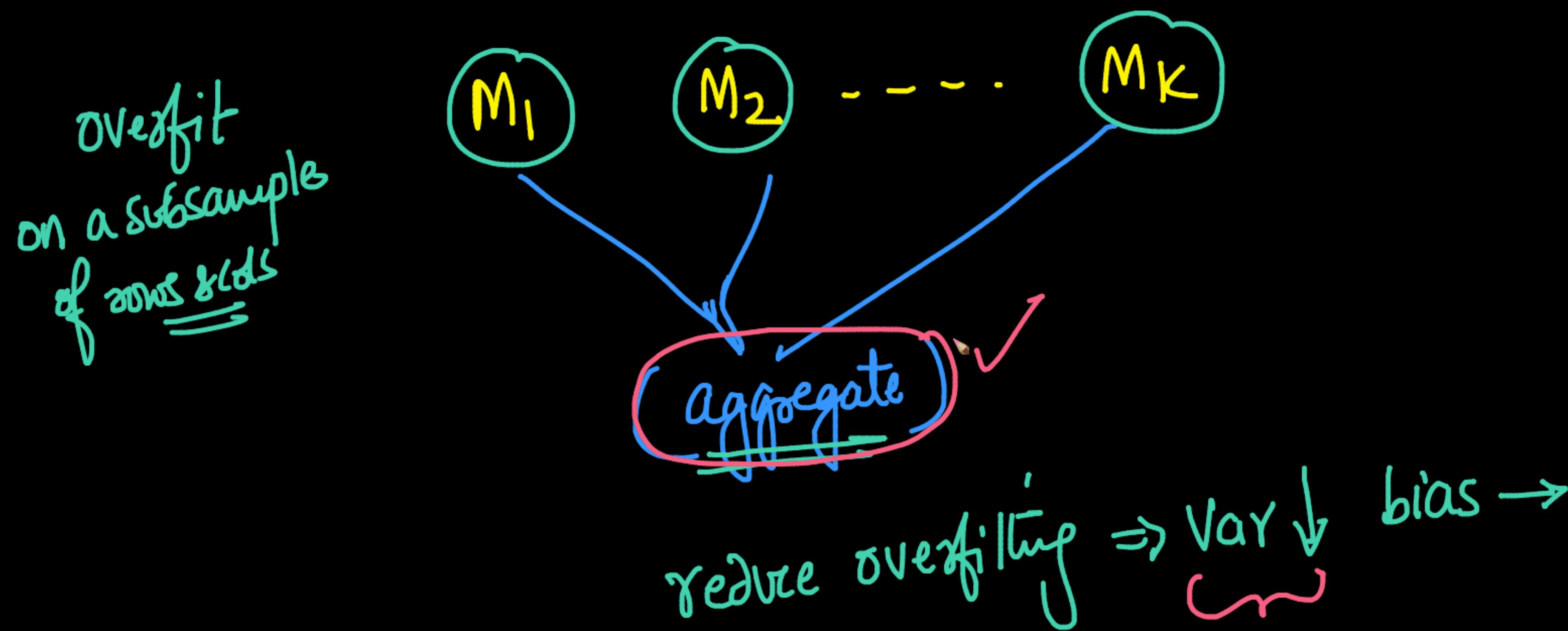
Viz a RF

is tricky

$$d=3$$

$$d^l=2$$

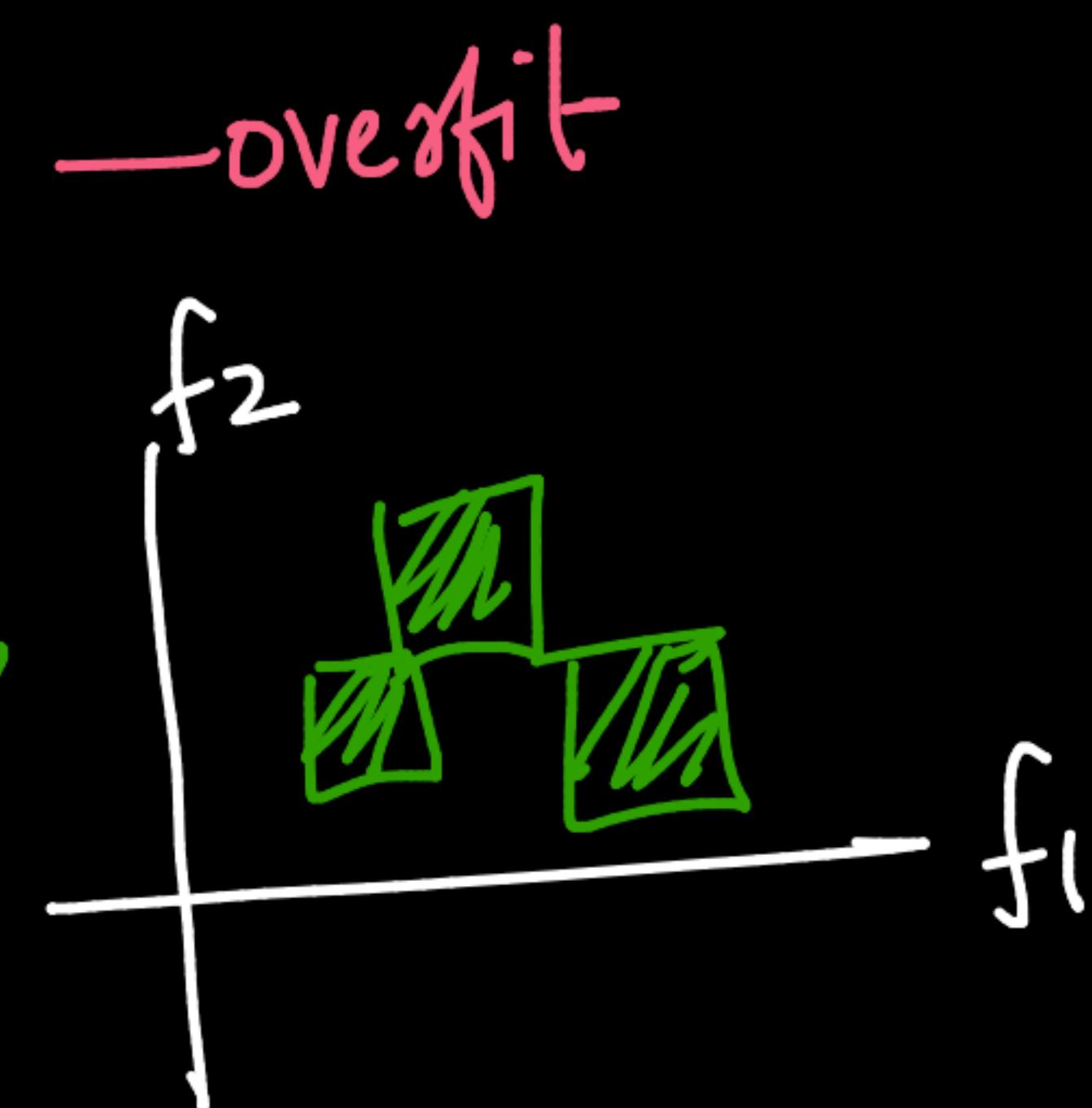
$$n=1\text{Million}$$



base learners ↑

$|C=1$ base-learner
=

DT → deep tree ⇒



ensemble of high bias & high-variance
↓
underfit

↓
overfit

Bagging → ensemble of high-var

RF:

(agg): Var ↓ bias →

GBDT: } → ensemble of high bias models --

Boosting }

DT
large depth
Linear-classifier
Underfit

high var

high bias

agg
Var ↓
bias → ↑

K : # base learners ↑ less overfitting



Typically: non-ensemble models

Var ↓ bias ↑

↓
Size of sample:

$$\boxed{m \downarrow} \quad \frac{m}{n} \downarrow$$

more base learners $\leftarrow \uparrow$
~~agg~~

↳ smaller data
{ overfit the
base learner }
Var \uparrow

Var \downarrow

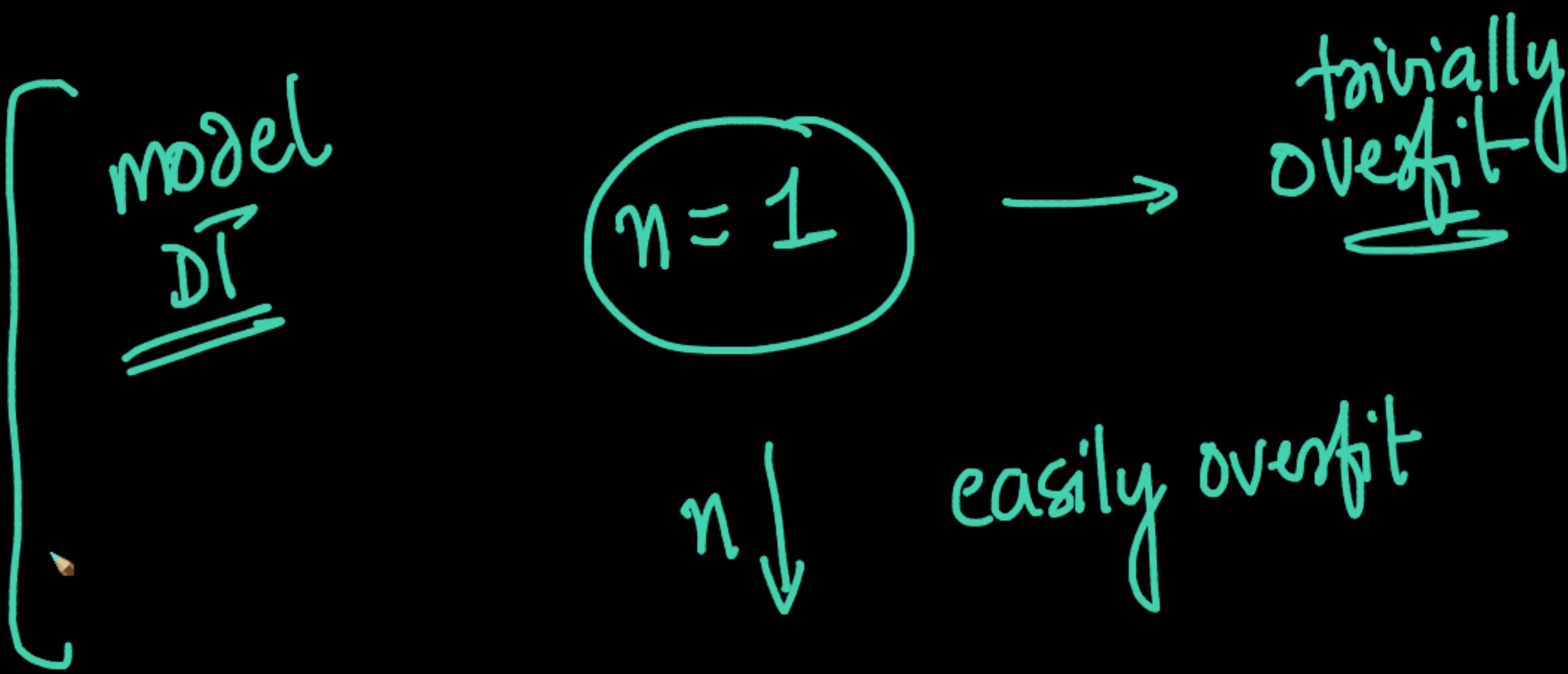
$m \approx n \Rightarrow$ base learners are similar

$$\frac{m}{n} \uparrow$$

RF will
overfit



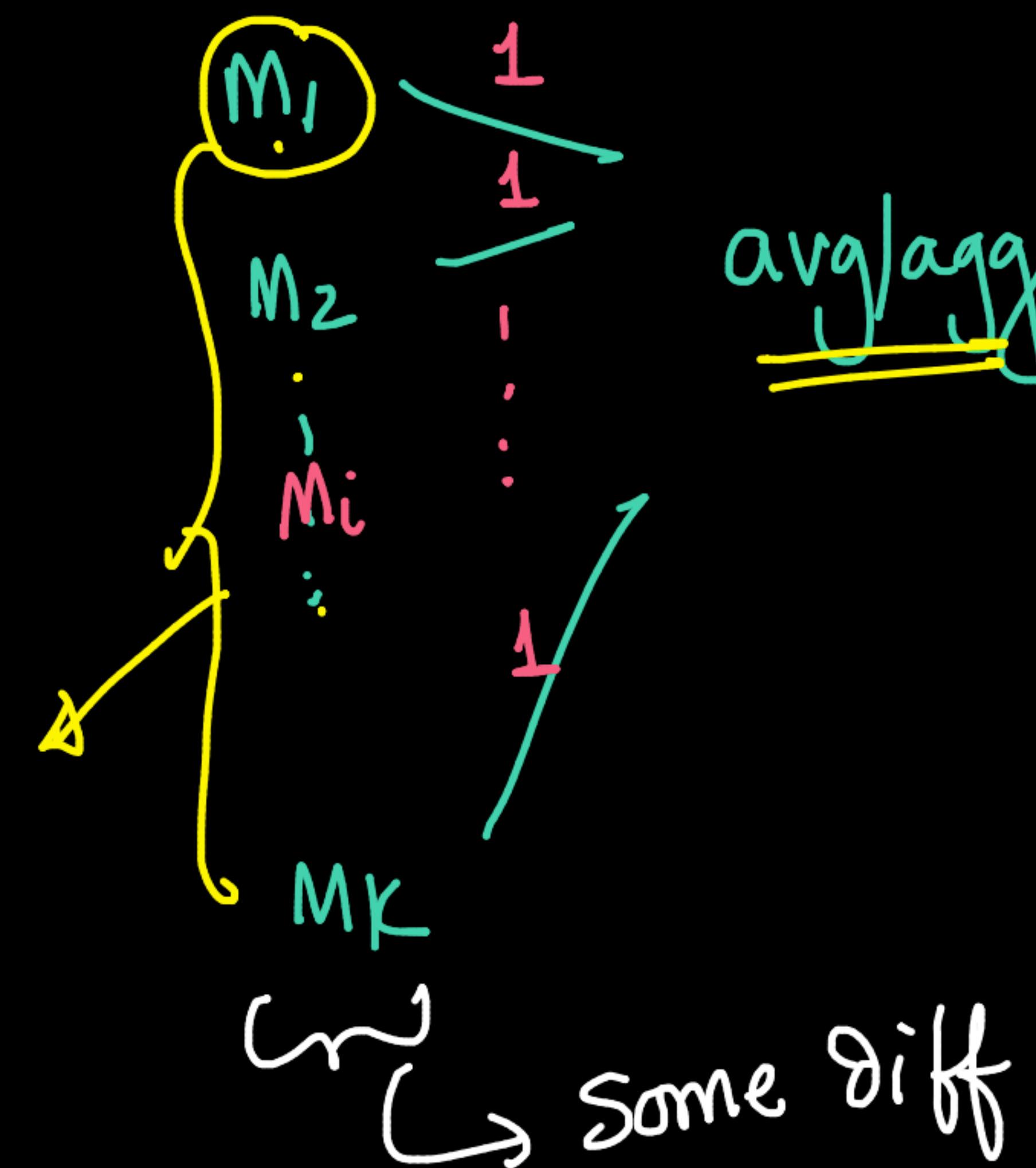
(agg) doesn't work well



$w_i : \text{(errors) on CV}$

-1^{-1} oob

$\sqrt{\frac{\text{Var}}{\text{bias}}}$



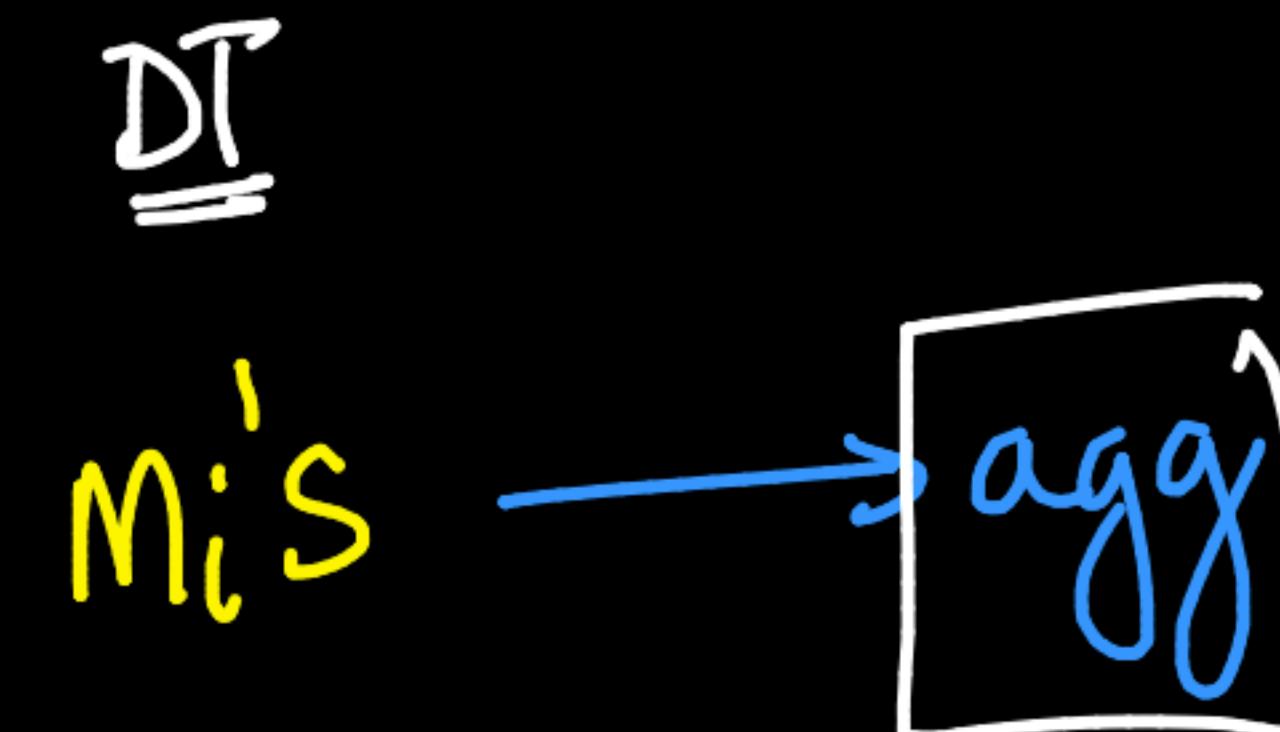
Some diff in error datarw }
oob - CV

BETTER

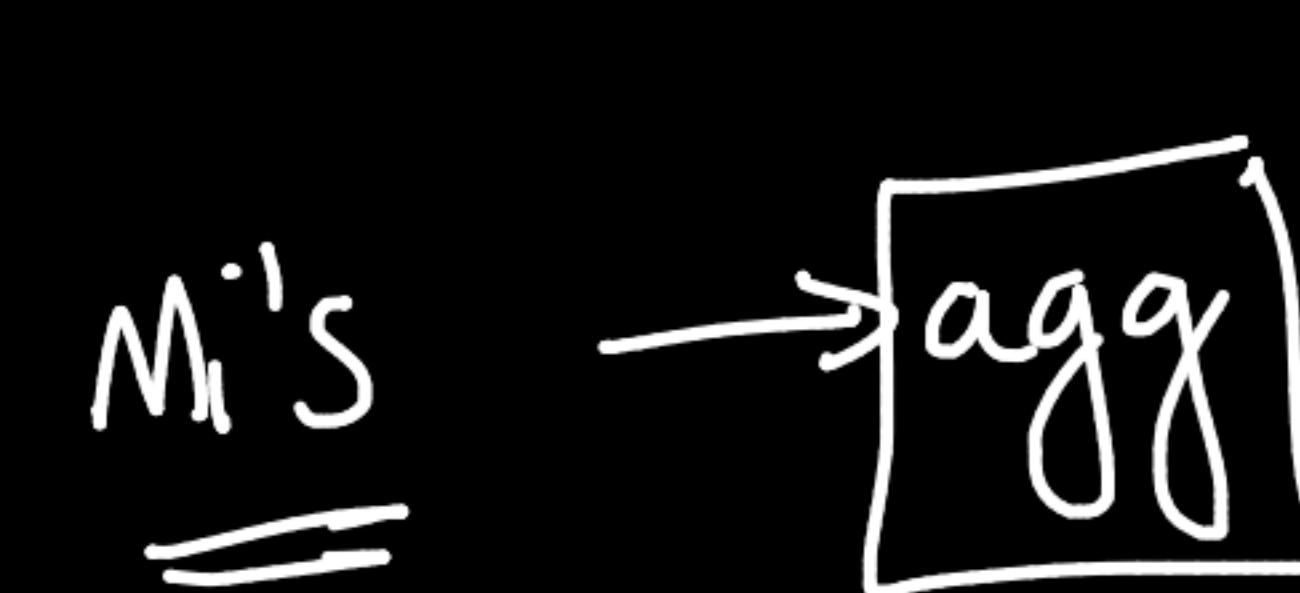
high-var
low bias

WORSE

med-var
med-bias



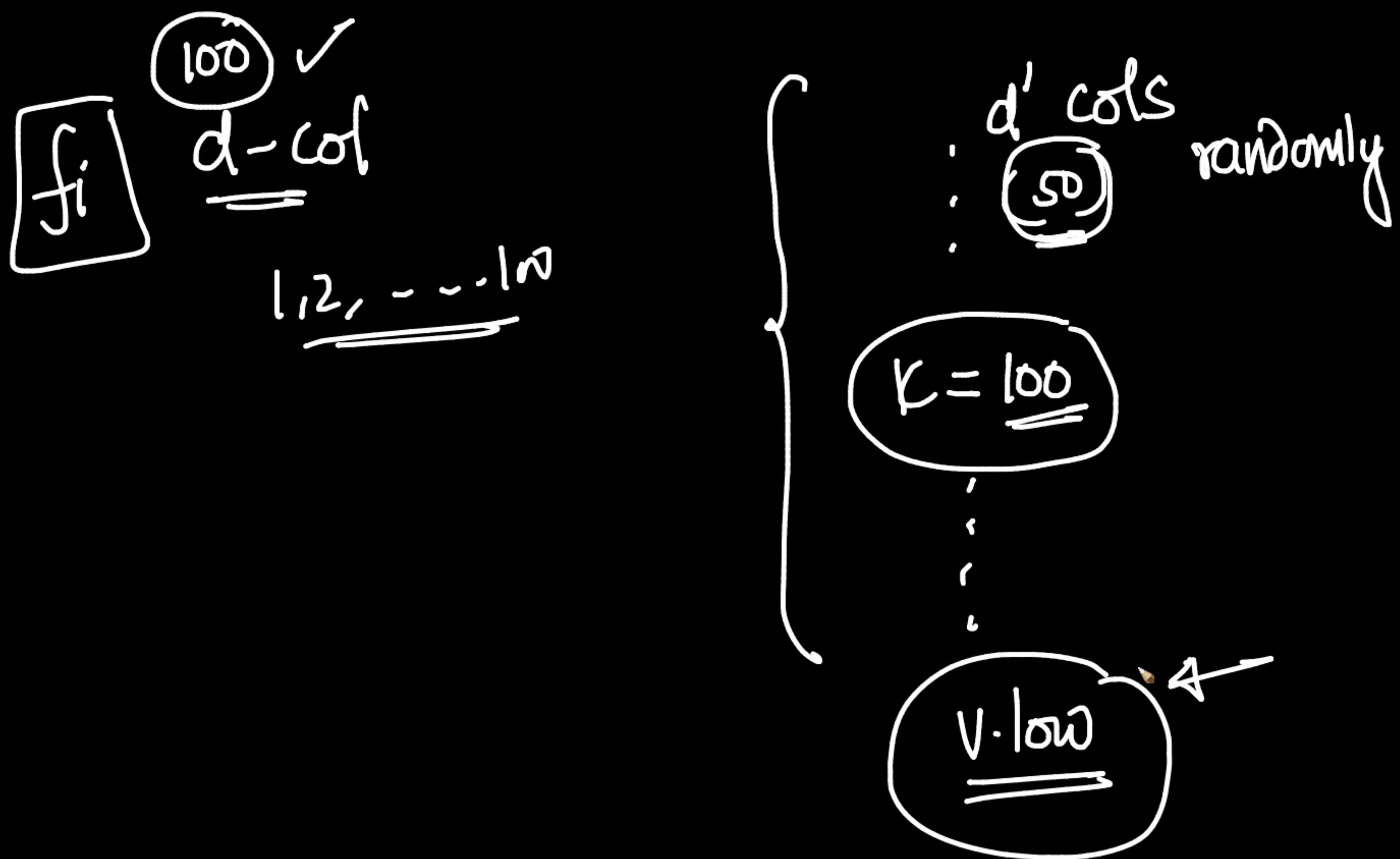
Var \downarrow bias $\xrightarrow{\text{low}}$

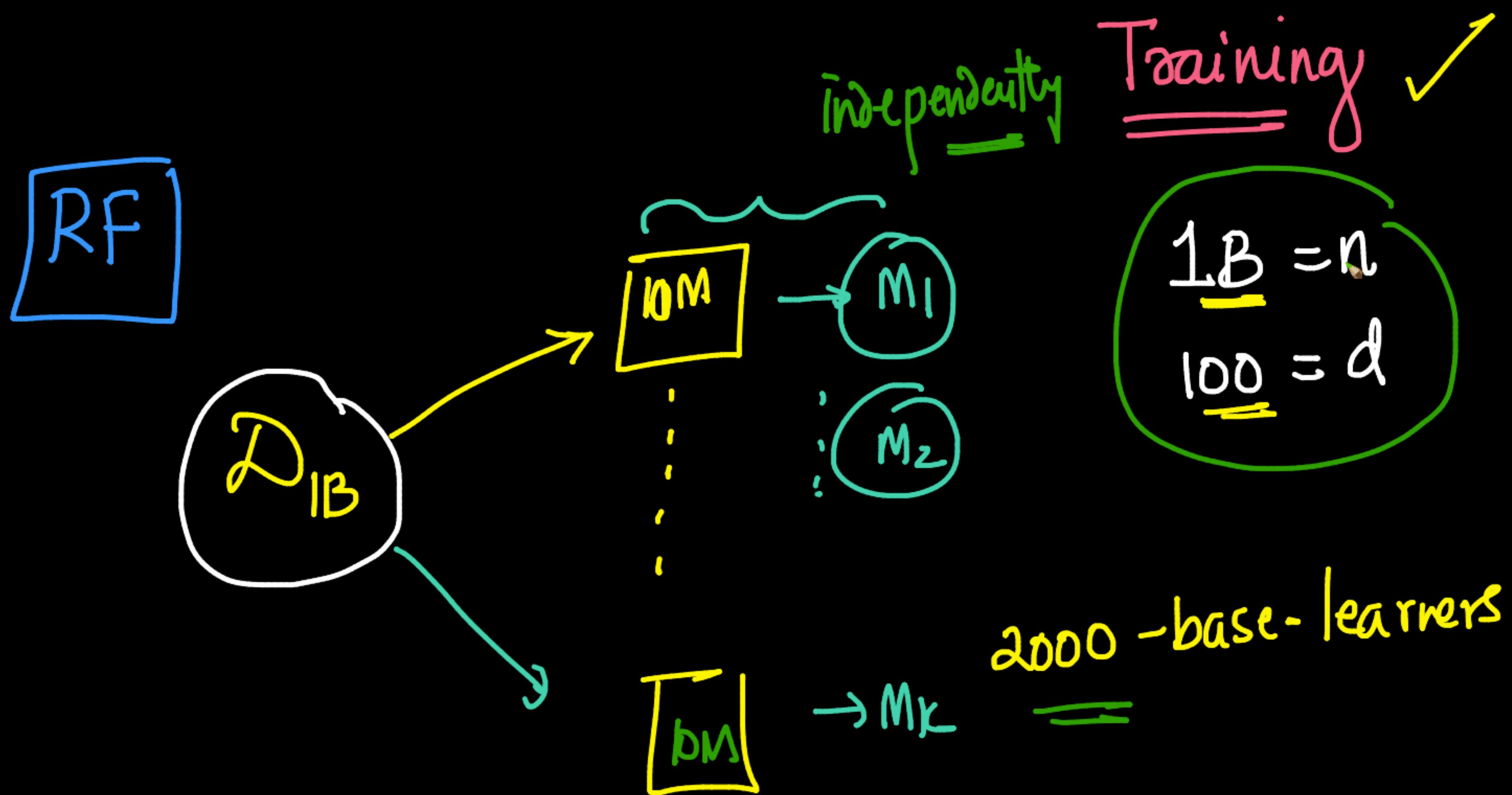


Var \downarrow bias: $\xrightarrow{\text{med}}$

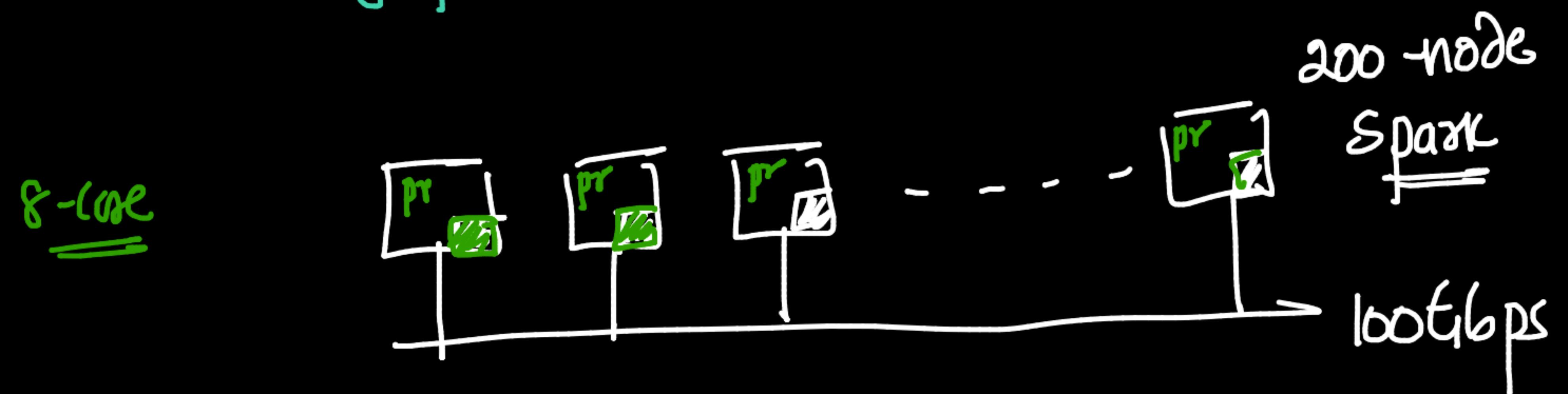
$$E = \underline{\text{bias}}^2 + \underline{\text{Var}} + \text{irr. com}$$

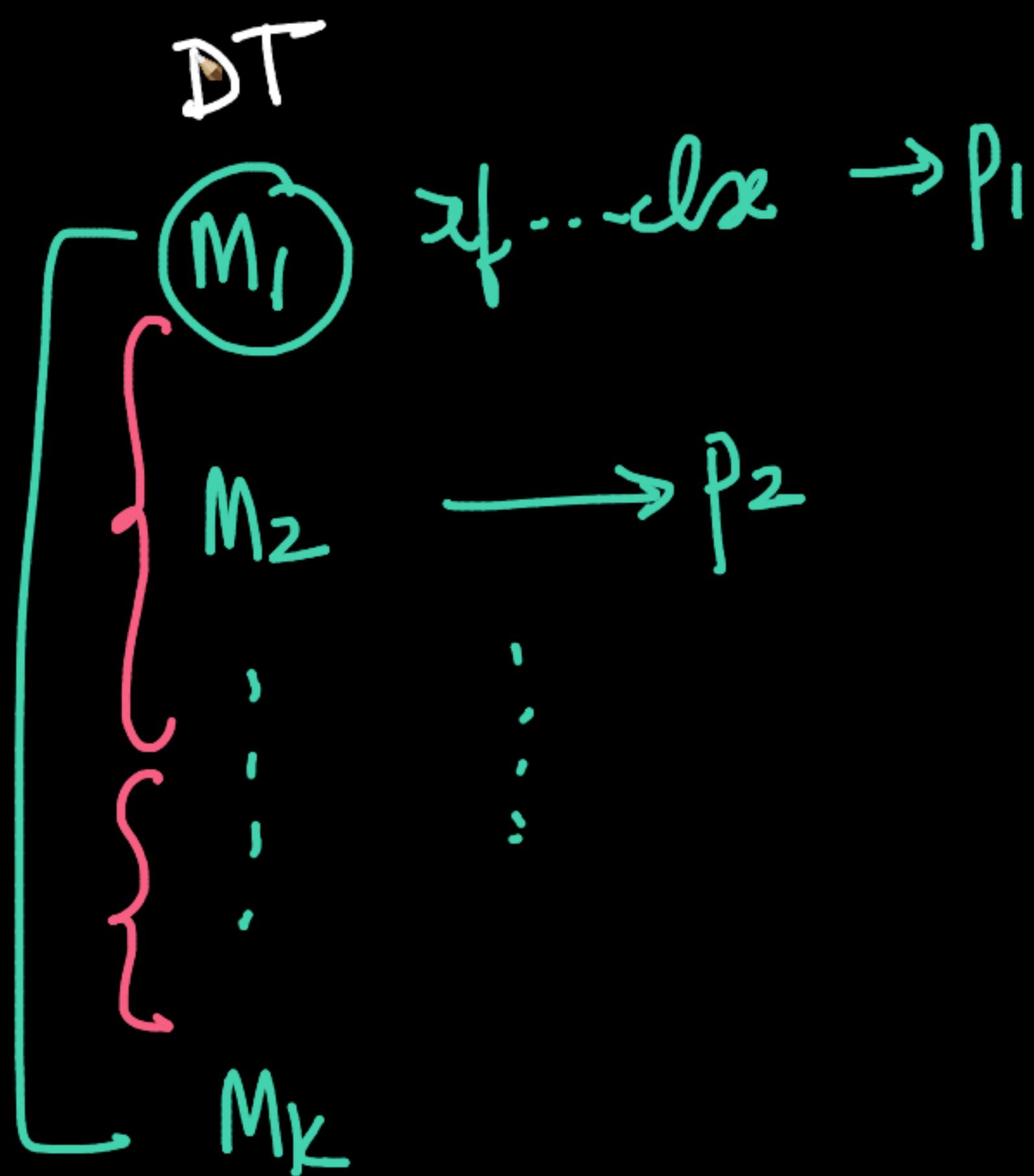
RF → ~~GB for CV~~
max depth = 15
base learner
min depth = 5



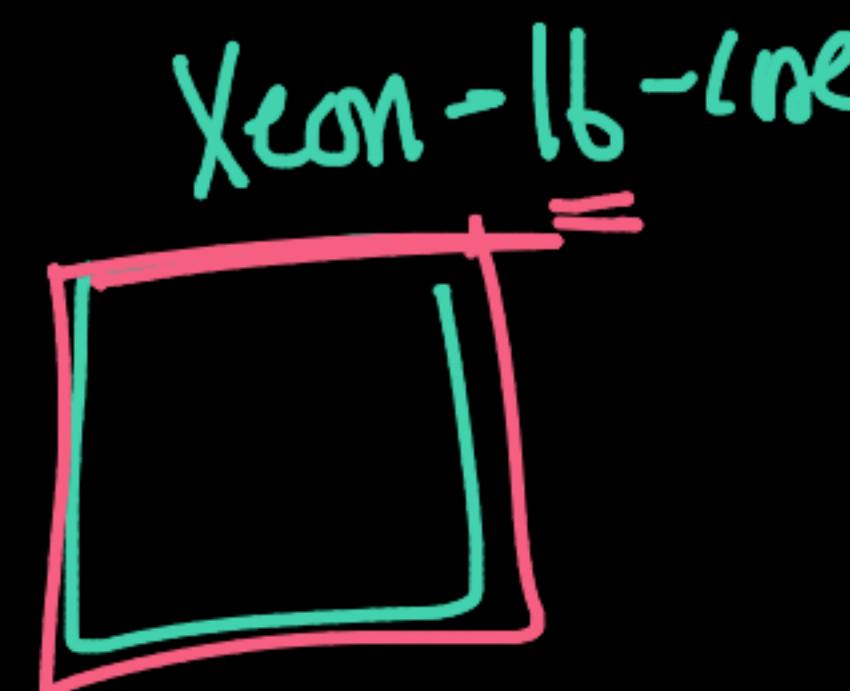


trivially parallelizable & very fast to train





$$\mathcal{O}(K \times \underbrace{\text{depth-max}}_{d})$$



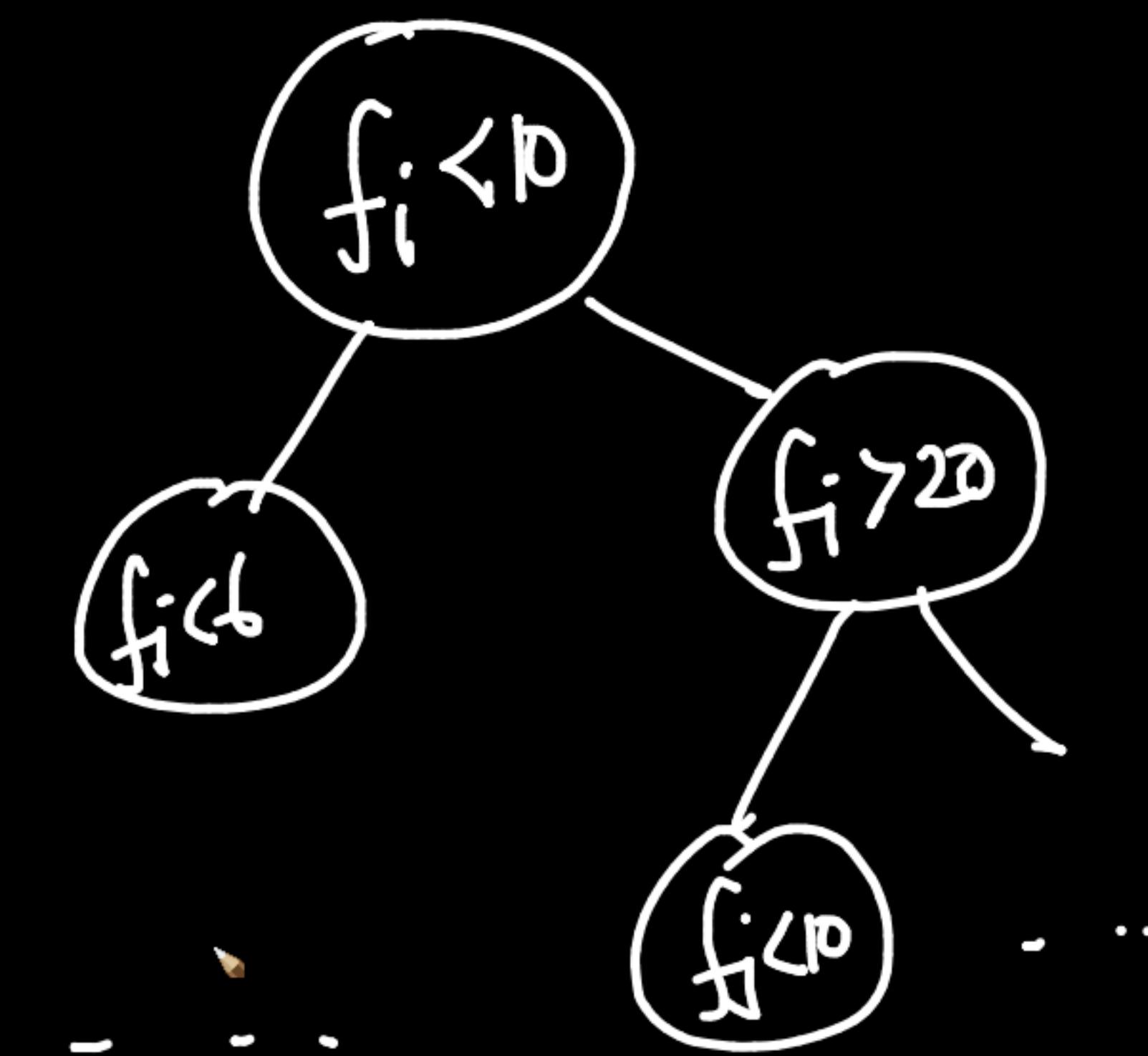
Xeon-16-line

parallelize @ runtime

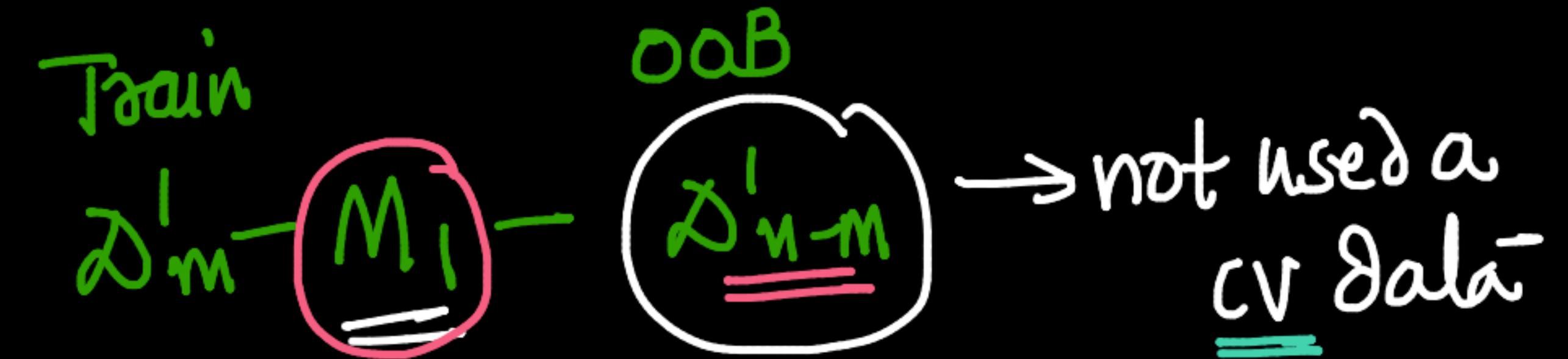
multi-core

max-depth $\neq d$

Run-time



OOB-Score:



Train-data perf. \neq \approx OOB.

overfit

Vari↑ low bias

perf. of M_i , m

D_{n-m} : $\boxed{\text{OOB}_i}$

OOB-Scores

min error

max perf

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

not the best impl

scikit learn Install User Guide API Examples Community More ▾ Go

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Xg boost ← GBDT

Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`

The number of trees in the forest.

Toggle Menu

47 / 49 value of `n_estimators` changed from 10 to

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Go

scikit learn Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`

The number of trees in the forest.

Toggle Menu

48 / 50 value of `n_estimators` changed from 10 to

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

X OOB as CV

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`

The number of trees in the forest.

Toggle Menu

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores
tClassifier

Examples using
sklearn.ensemble.RandomForestC

sklearn.ensemble.RandomForestClassifier

`class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)`

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

base-learner - overfit

K ↑ for my model to work

Parameters: `n_estimators : int, default=100`
The number of trees in the forest.

value of `n_estimators` changed from 10 to

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Go

scikit learn Install User Guide API Examples Community More ▾

BOBX-XKGV

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`

The number of trees in the forest.

Toggle Menu

51 / 53 value of `n_estimators` changed from 10 to

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Go

scikit learn Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC

2.

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`

The number of trees in the forest.

Toggle Menu

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores
tClassifier

Examples using
sklearn.ensemble.RandomForestC

base learners
underfit

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None,  
min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto',  
max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None,  
random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0,  
max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`
The number of trees in the forest.

Toggle Menu

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC

1

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

Shallow base-learners

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`

The number of trees in the forest.

Toggle Menu

54 / 56 value of `n_estimators` changed from 10 to

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Go

scikit learn Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`

The number of trees in the forest.

Toggle Menu

value of `n_estimators` changed from 10 to

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

scikit learn

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC1

min_weight_fraction_leaf : float, default=0.0

The minimum weighted fraction of the sum total of weights (of all the input samples) required to be at a leaf node. Samples have equal weight when sample_weight is not provided.

max_features : {"auto", "sqrt", "log2"}, int or float, default="auto"

The number of features to consider when looking for the best split:

- If int, then consider max_features features at each split.
- If float, then max_features is a fraction and round(max_features * n_features) features are considered at each split.
- If "auto", then max_features=sqrt(n_features).
- If "sqrt", then max_features=sqrt(n_features) (same as "auto").
- If "log2", then max_features=log2(n_features).
- If None, then max_features=n_features.

Note: the search for a split does not stop until at least one valid partition of the node samples is found, even if it requires to effectively inspect more than max_features features.

max_leaf_nodes : int, default=None

Grow trees with max_leaf_nodes in best-first fashion. Best nodes are defined as relative reduction in impurity. If None, then unlimited number of leaf nodes.

Toggle Menu

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`

The number of trees in the forest.

Toggle Menu

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`

The number of trees in the forest.

Toggle Menu

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

scikit learn

Prev Up Next

scikit-learn 1.0.2 Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier Examples using sklearn.ensemble.RandomForestC

oob_score : bool, default=False

Whether to use out-of-bag samples to estimate the generalization score. Only available if bootstrap=True.

{ ④

n_jobs : int, default=None

The number of jobs to run in parallel. `fit`, `predict`, `decision_path` and `apply` are all parallelized over the trees. `None` means 1 unless in a `joblib.parallel_backend` context. `-1` means using all processors. See [Glossary](#) for more details.

random_state : int, RandomState instance or None, default=None

Controls both the randomness of the bootstrapping of the samples used when building trees (if `bootstrap=True`) and the sampling of the features to consider when looking for the best split at each node (if `max_features < n_features`). See [Glossary](#) for details.

verbose : int, default=0

Controls the verbosity when fitting and predicting.

warm_start : bool, default=False

When set to `True`, reuse the solution of the previous call to fit and add more estimators to the ensemble, otherwise, just fit a whole new forest. See [the Glossary](#).

class_weight : {"balanced", "balanced_subsample"}, dict or list of dicts, default=None

form `{class_label: weight}`. If not given, all

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

max_samples 3/6 ^ v x

scikit learn

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores
tClassifier

Examples using
sklearn.ensemble.RandomForestCl

m

For multi-output, the weights of each column of y will be multiplied.

Note that these weights will be multiplied with sample_weight (passed through the fit method) if sample_weight is specified.

ccp_alpha : non-negative float, default=0.0

Complexity parameter used for Minimal Cost-Complexity Pruning. The subtree with the largest cost complexity that is smaller than `ccp_alpha` will be chosen. By default, no pruning is performed. See [Minimal Cost-Complexity Pruning](#) for details.

New in version 0.22.

max_samples : int or float, default=None

If `bootstrap` is True, the number of samples to draw from X to train each base estimator.

- If None (default), then draw `X.shape[0]` samples.
- If int, then draw `max_samples` samples.
- If float, then draw `max_samples * X.shape[0]` samples. Thus, `max_samples` should be in the interval `(0.0, 1.0)`.

New in version 0.22.

Attributes: `base_estimator : DecisionTreeClassifier`

create the collection of fitted sub-estimators.

scikit-learn

Gradient boosting - Wikipedia | ESLII.pdf

scikit-learn.org/stable/modules/tree.html#minimal-cost-complexity-pruning

1.10. Decision Trees – scikit-learn

Prev Up Next

scikit-learn 1.0.2 Other versions

Please cite us if you use the software.

1.10.8. Minimal Cost-Complexity Pruning

CCP

Minimal cost-complexity pruning is an algorithm used to prune a tree to avoid over-fitting, described in Chapter 3 of [BRE]. This algorithm is parameterized by $\alpha \geq 0$ known as the complexity parameter. The complexity parameter is used to define the cost-complexity measure, $R_\alpha(T)$ of a given tree T :

$$R_\alpha(T) = R(T) + \alpha|\tilde{T}|$$

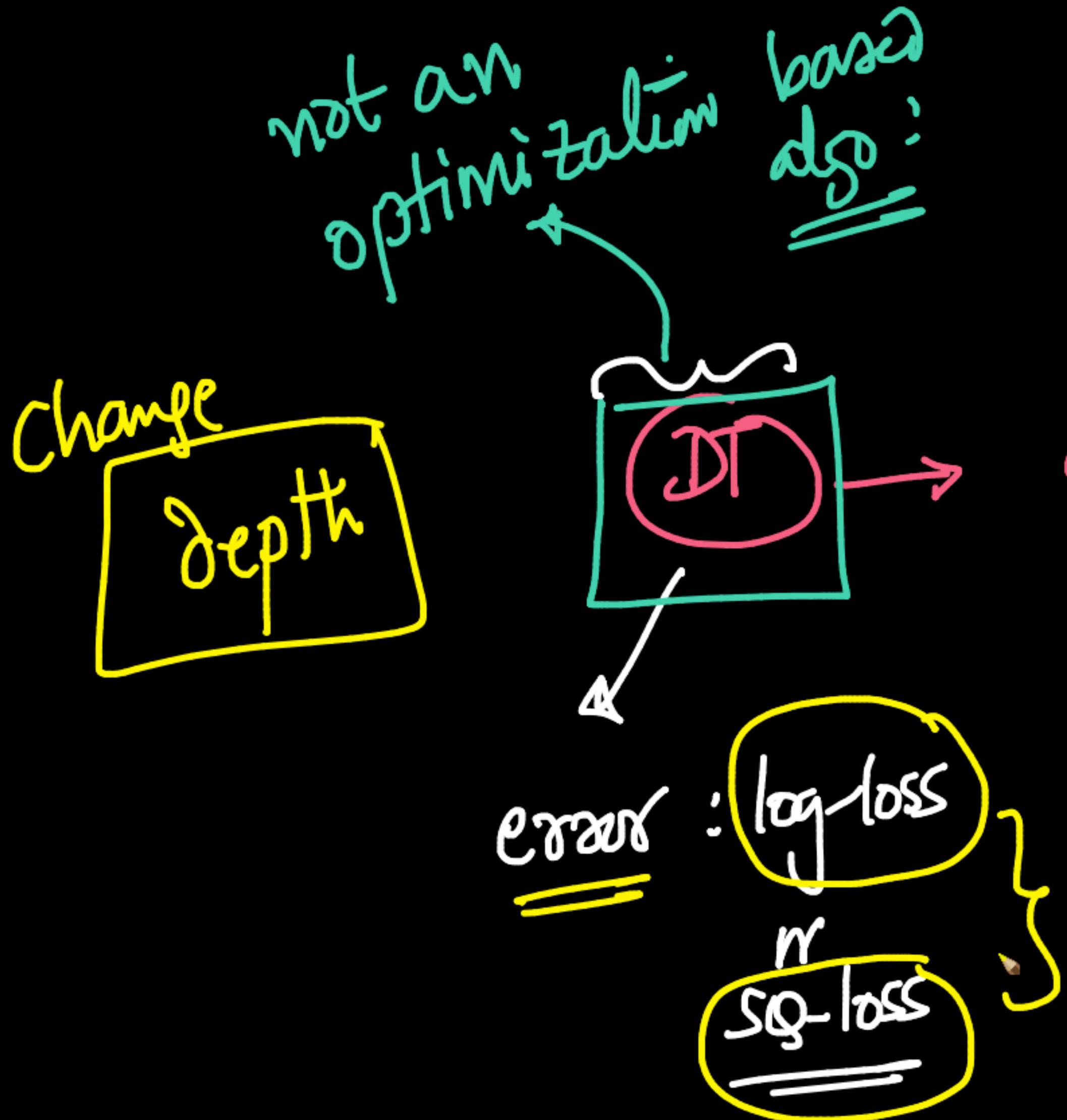
where $|\tilde{T}|$ is the number of terminal nodes in T and $R(T)$ is traditionally defined as the total misclassification rate of the terminal nodes. Alternatively, scikit-learn uses the total sample weighted impurity of the terminal nodes for $R(T)$. As shown above, the impurity of a node depends on the criterion. Minimal cost-complexity pruning finds the subtree of T that minimizes $R_\alpha(T)$.

The cost complexity measure of a single node is $R_\alpha(t) = R(t) + \alpha$. The branch, T_t , is defined to be a tree where node t is its root. In general, the impurity of a node is greater than the sum of impurities of its terminal nodes, $R(T_t) < R(t)$. However, the cost complexity measure of a node, t , and its branch, T_t , can be equal depending on α . We define the effective α of a node to be the value where they are equal, $R_\alpha(T_t) = R_\alpha(t)$ or $\alpha_{eff}(t) = \frac{R(t) - R(T_t)}{|T|-1}$. A non-terminal node with the smallest value of α_{eff} is the weakest link and will be pruned. This process stops when the pruned tree's minimal α_{eff} is greater than the `ccp_alpha` parameter.

Examples:

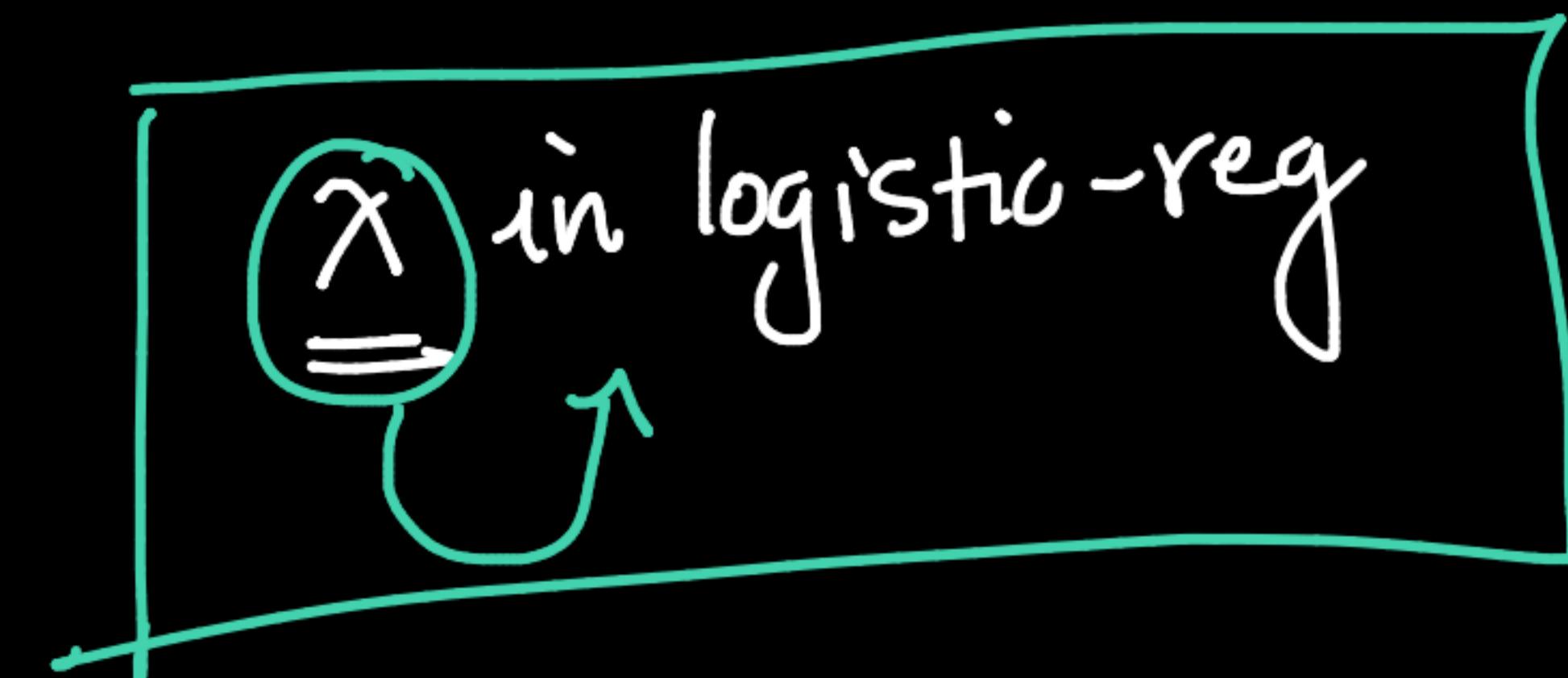
- Post-pruning decision trees with cost-complexity pruning

Toggle Menu



$$\boxed{\min \left(\text{loss} + \alpha \# \text{leaves in the tree} \right)}$$

overfit vs underfit



CP

depth



$$\min \left(\frac{\text{train loss}}{T} + \alpha \# \text{leaf-nodes} \right)$$

$\alpha \uparrow$:
do not overfit
lower depth

regularization
as it controls
the tree from
overfitting

Gradient boosting - Wikipedia x | ESLII.pdf x | sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html

Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomFores tClassifier

Examples using sklearn.ensemble.RandomForestC

sklearn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source] ↑ *base-learners Shallow*

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

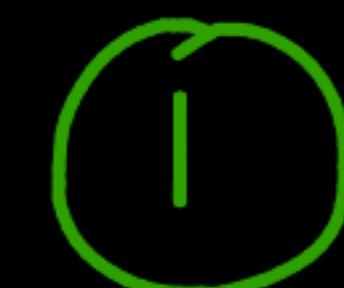
Read more in the [User Guide](#).

Parameters: `n_estimators : int, default=100`
The number of trees in the forest.

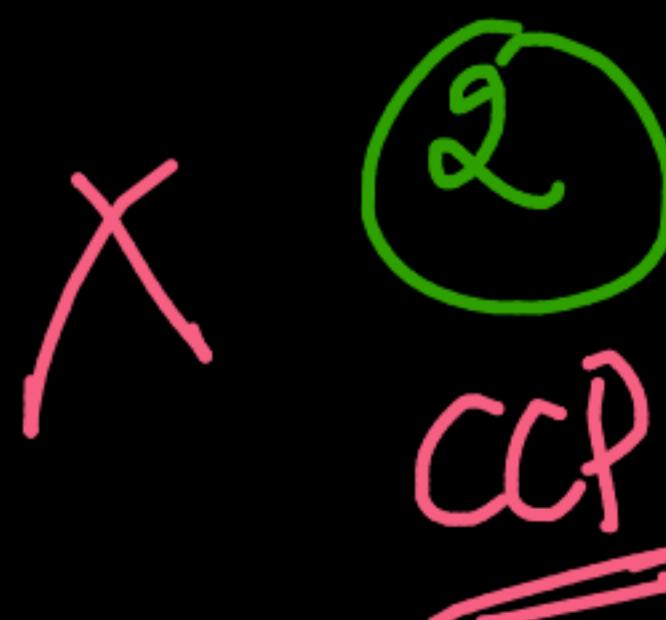
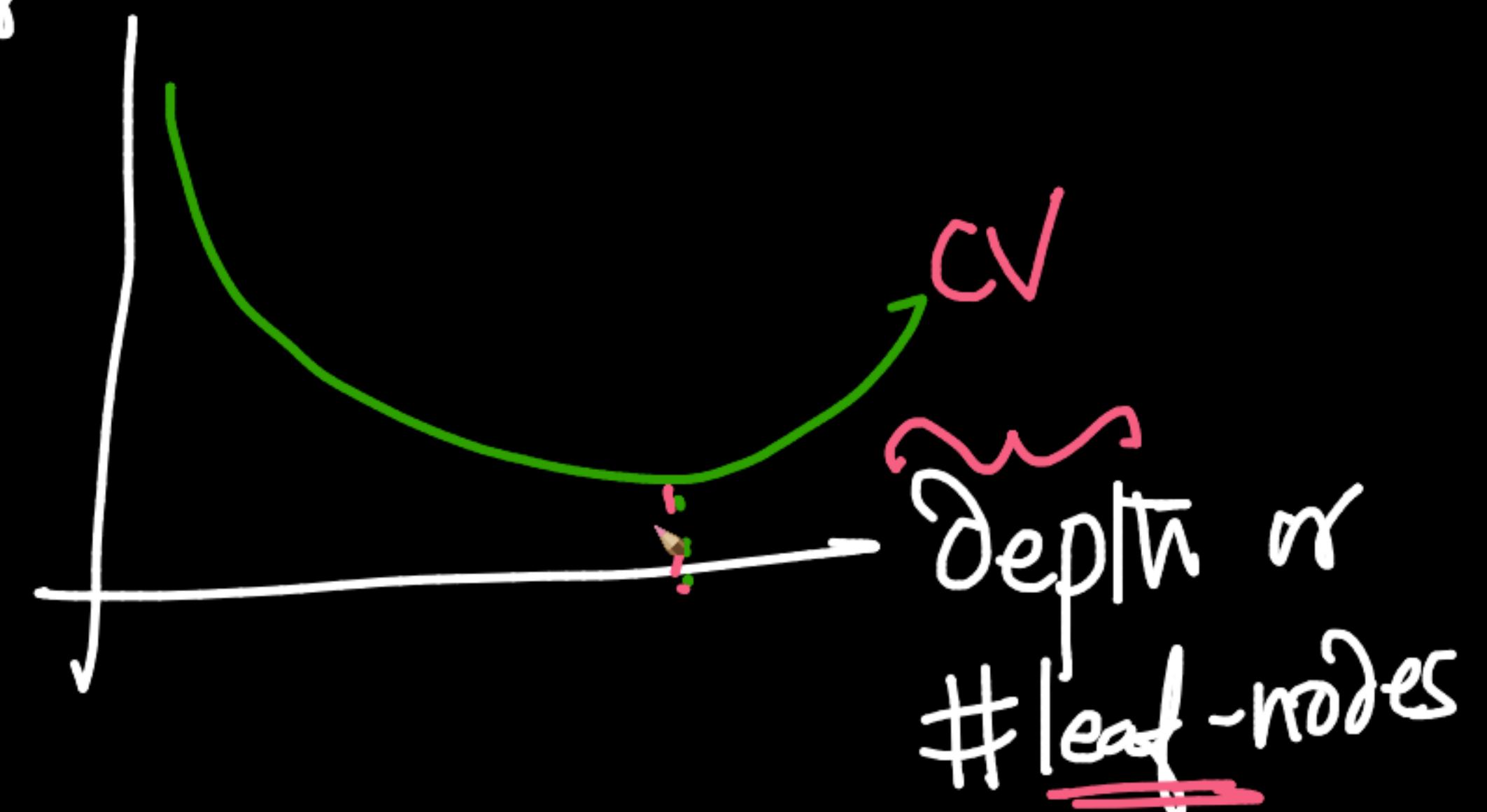
Toggle Menu



error



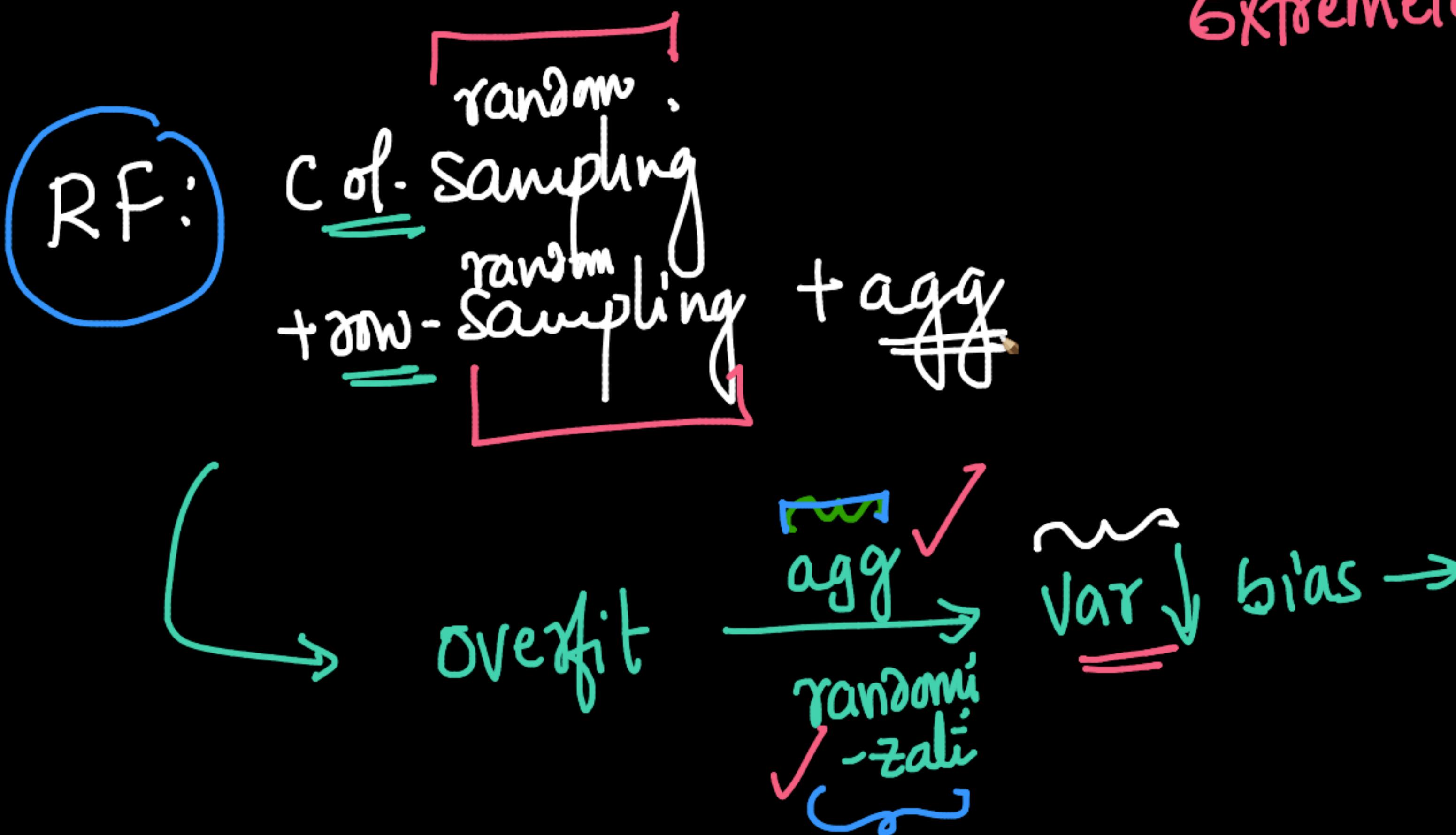
=



$$\min \left(\text{↑loss} + \frac{\alpha}{\text{↑train}} \# \text{leaf-nodes} \right)$$

CCP

Extra Trees or Extremely Randomized Trees



Randomization is a useful strategy for regularization → 2048

(*) Var ↓

{ RF; Xgboost
ExtraTrees; DL } ↗ dropouts

Explain trees:

- Col. sampling
- Row. sampling
- Aggregation
- randomly picks the 'c' to split
for numerical features

~~recall:~~

f_i	y_i
2.1	1
2.2	1
2.3	6
2.4	1
2.8	1
2.9	0

DT

$f_i < t$

$$\tau_1 = 2 \cdot 1 \rightarrow \Delta H : IG_2$$

$$\tau_2 = 2 \cdot 2 \rightarrow IG_2$$

⋮
⋮
⋮

$$\tau_6 = 2 \cdot 9 \rightarrow \underline{IG_L}$$

Extra trees!

more
randomization

	f_i	y
2.1	1	1
2.2	1	0
2.3	0	0
2.4	1	0
2.6	0	0

$$\begin{cases} \tau_1 = 2 \cdot 3 \rightarrow 1G_{2 \cdot 3} \\ \tau_2 = 2 \cdot 6 \rightarrow 1G_{2 \cdot 6} \\ \vdots \\ \tau_{10} \end{cases}$$

mislabel

(2006)

Extra trees: - one more

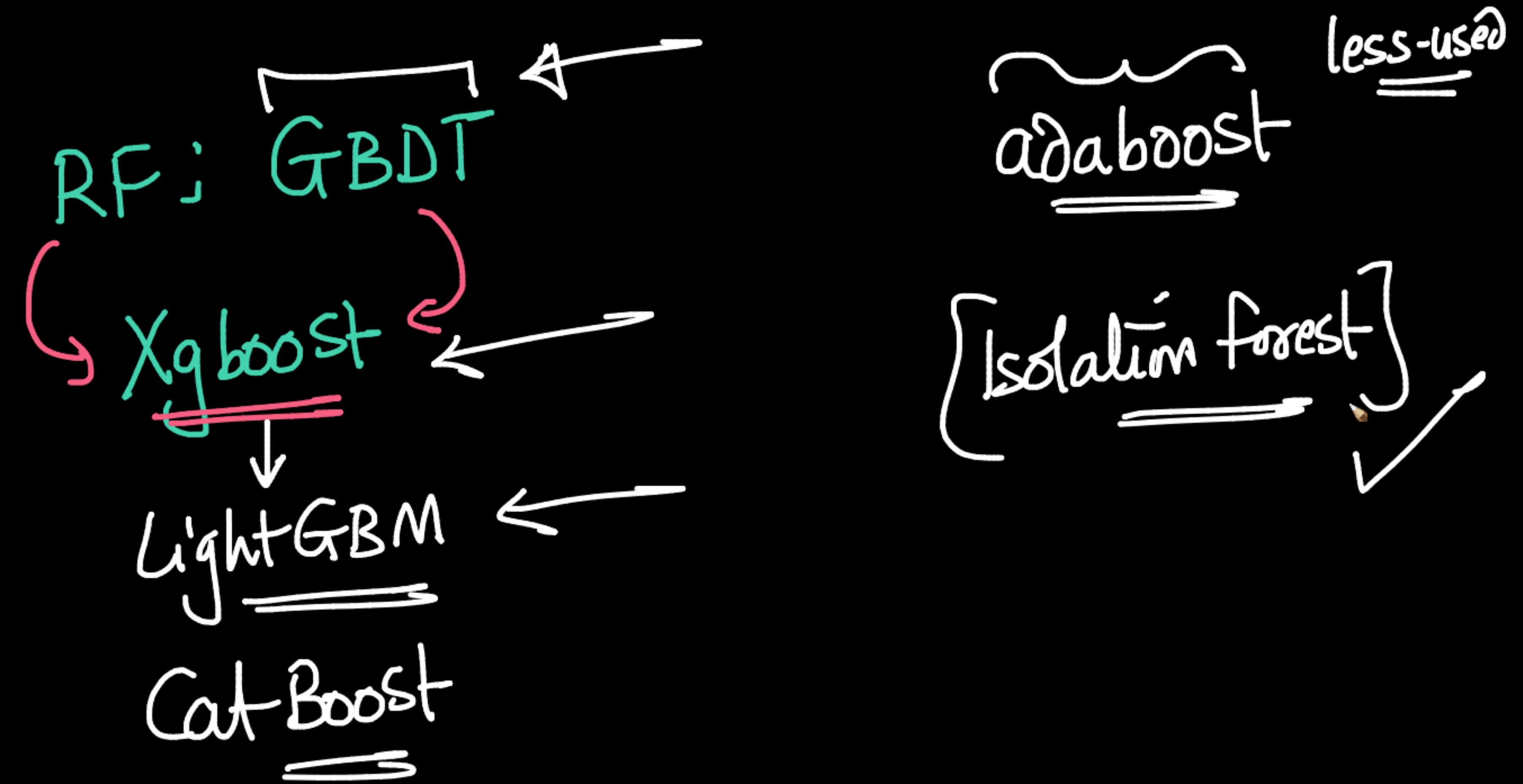
randomization 

↳ m is large (useful) → speed up

typical base-learners to perform well
most

not widely used

RF; GBDT

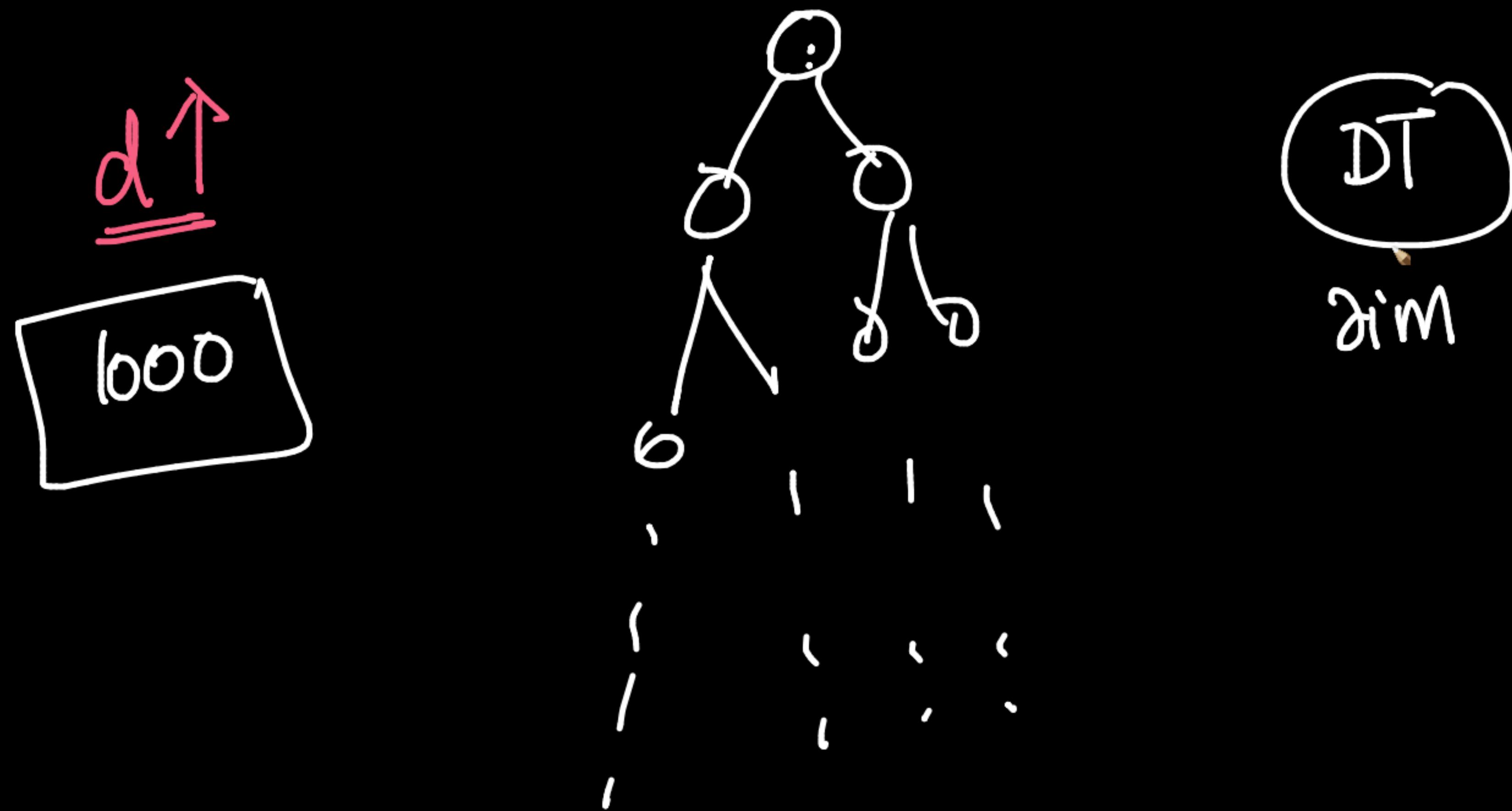


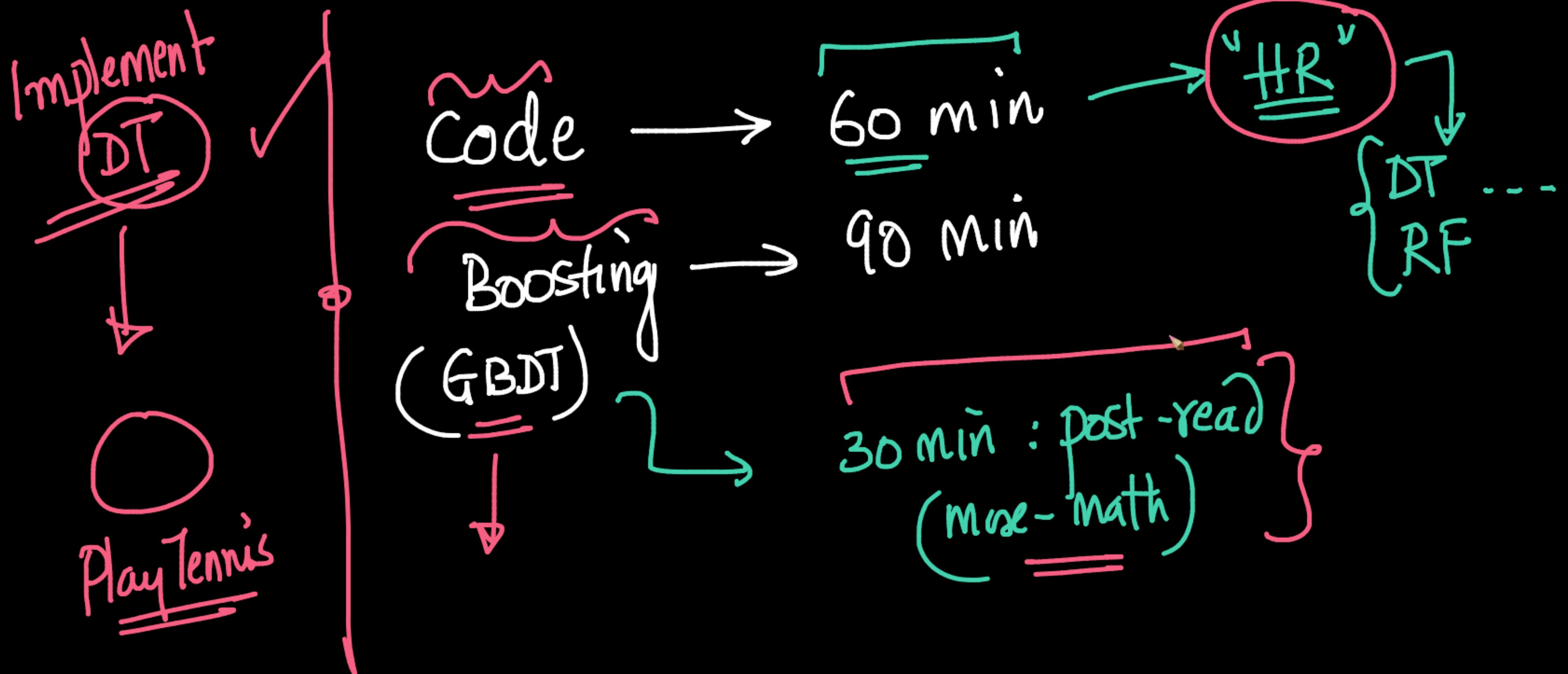
RF:

→ \dim is not too large

→ trivially parallelizable: n is large...

→ $K; \frac{m}{n}; \frac{d'}{d}$





Gradient boosting - Wikipedia

ESLII.pdf

https://stefvanbuuren.name/fin

mar wikipedia missing values -

stats361.pdf

web.stanford.edu/~swager/stats361.pdf



stats361.pdf

1 / 127 | - 175% + | ☰ ⚡



STATS 361: Causal Inference

Stefan Wager
Stanford University

Spring 2020

