

Bootstrapped Aggregation

Ensembles:

DT

{ Bagging + Random Forest

= class starts @ 9:03 PM

Bagging_RF.ipynb - Colaboratory +

Automatic saving failed. This file was updated repeatedly on another tab. Show diff

Saving failed since 20:32

Downloading...

[1] From: <https://drive.google.com/uc?id=12Bh2AN8LcZAlg20ehpOrEWccUDaSdsOG>
To: /content/y_test.pickle
100% 9.49k/9.49k [00:00<00:00, 14.3MB/s]

[2] import pickle
import pandas as pd
Load data (deserialize)
df = pd.read_csv("HR-Employee-Attrition.csv")

with open('preprocessed_X_sm.pickle', 'rb') as handle:
 X_sm = pickle.load(handle)

with open('y_sm.pickle', 'rb') as handle:
 y_sm = pickle.load(handle)

Using the best decision tree model as observed in the last lecture
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import KFold, cross_validate

tree_clf = DecisionTreeClassifier(random_state=7, max_depth = 4)
kfold = KFold(n_splits=10)
cv_acc_results = cross_validate(tree_clf, X_sm, y_sm, cv = kfold, scoring = 'accuracy', return_train_score = True)

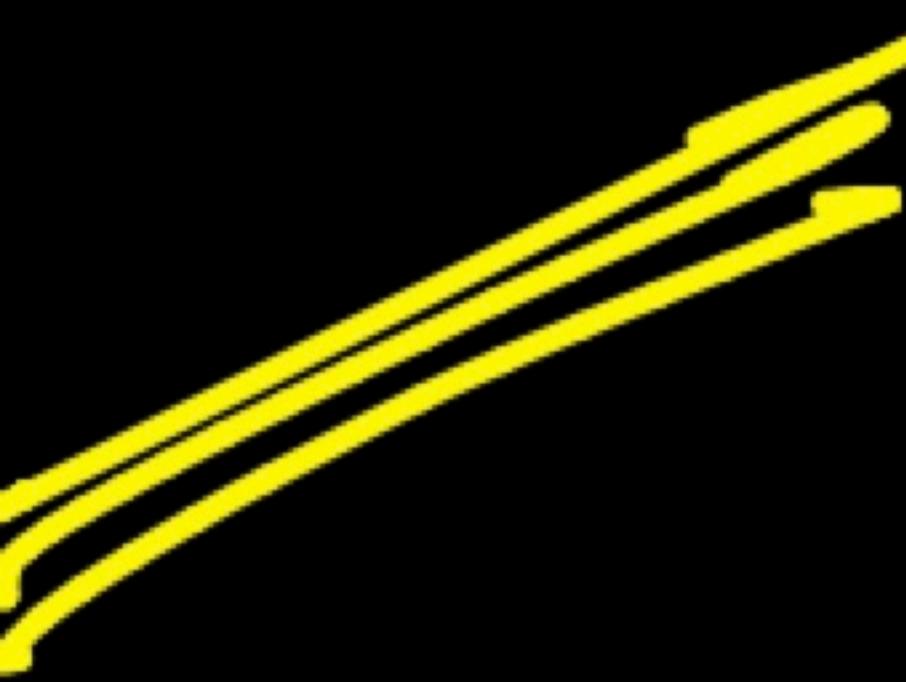
print(f"K-Fold Accuracy Mean: Train: {cv_acc_results['train_score'].mean()*100} Validation: {cv_acc_results['test_score'].mean()*100}")
print(f"K-Fold Accuracy Std: Train: {cv_acc_results['train_score'].std()*100} Validation: {cv_acc_results['test_score'].std()*100}")

K-Fold Accuracy Mean: Train: 83.50772540936214 Validation: 80.58049353701529
K-Fold Accuracy Std: Train: 0.7889219949763586 Validation: 6.337951865133366

Depth Train Acc Val Acc

{ 7 93.24 90.9 84.21

✓ 4 83.51 80.58 less

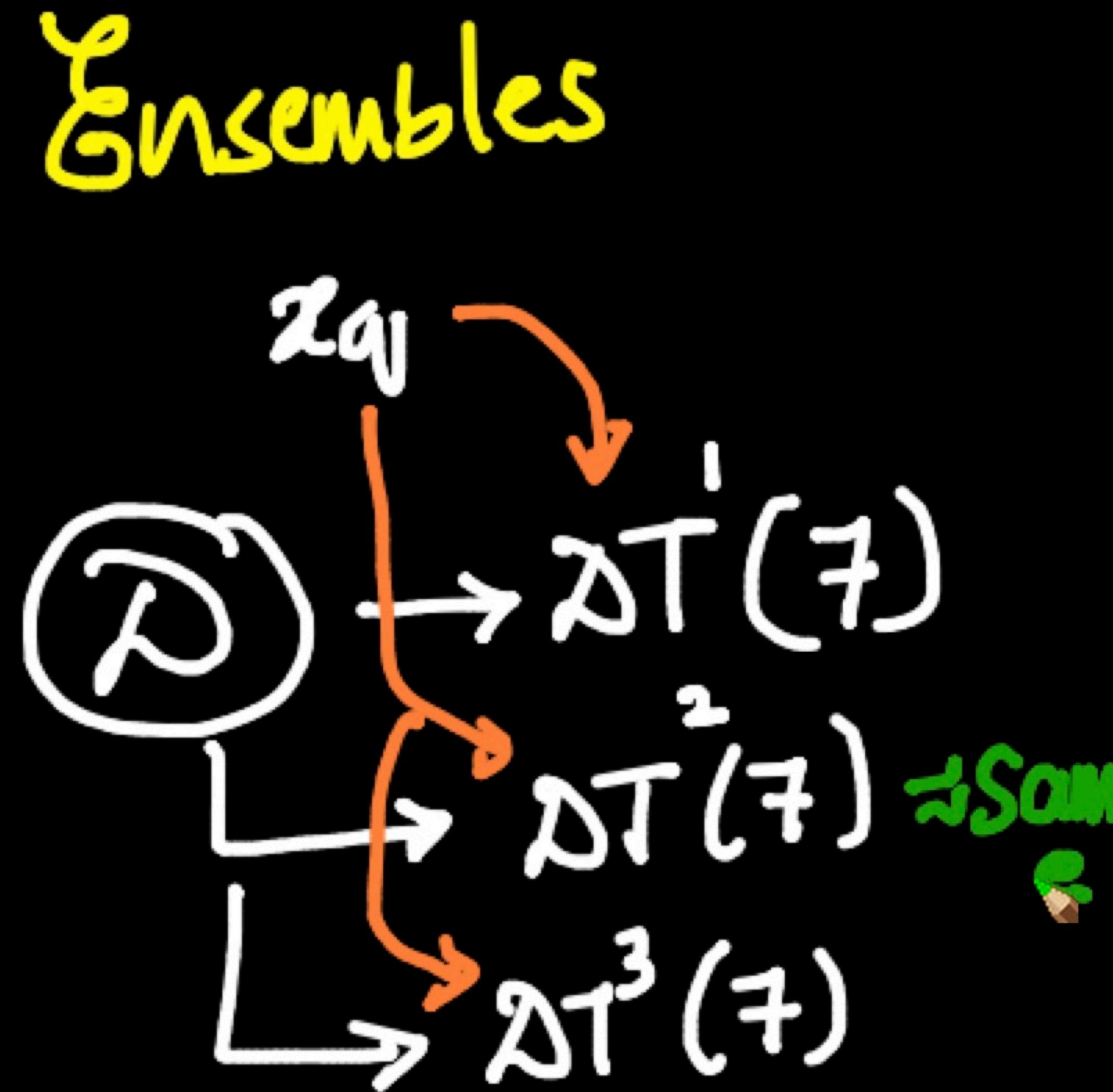


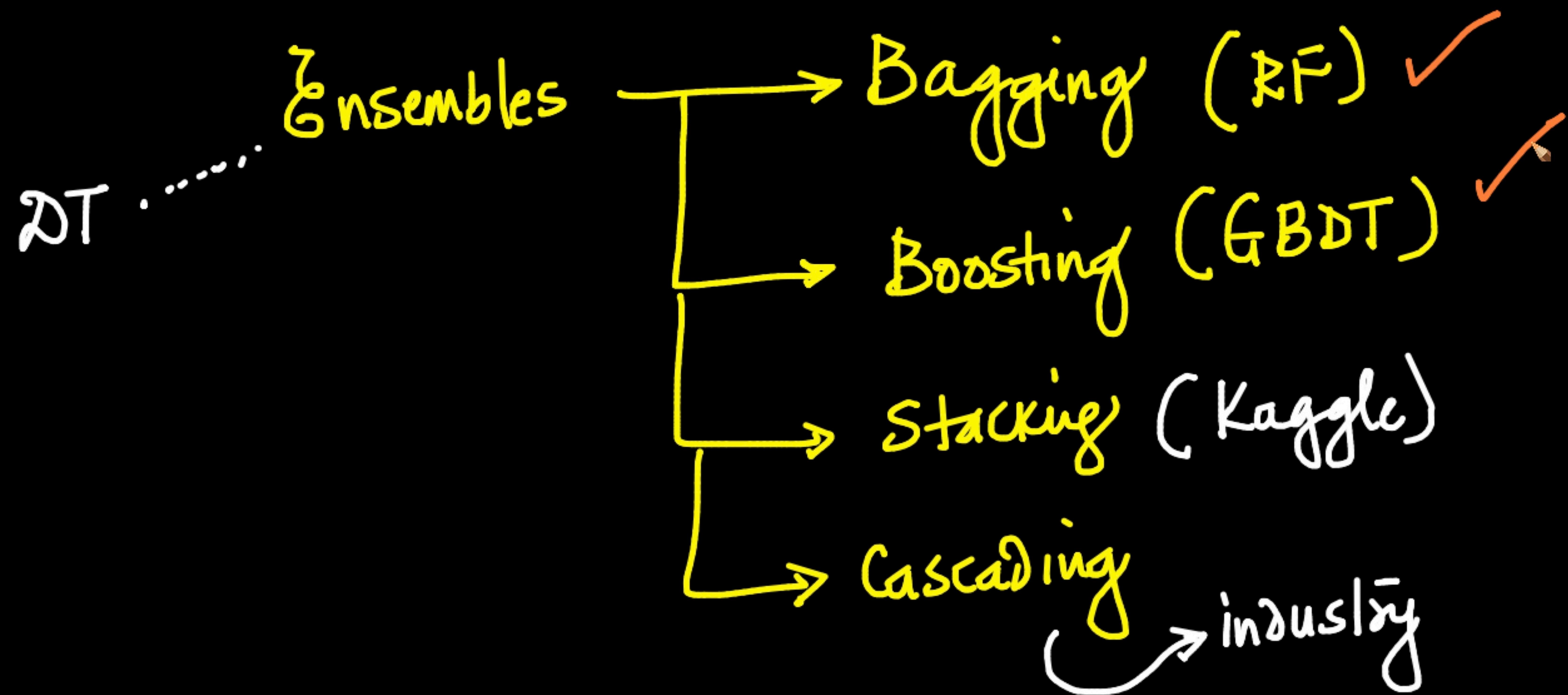
base-learner
One DT

✓ [100's of decision Trees]

Combine multiple base learner

'different'

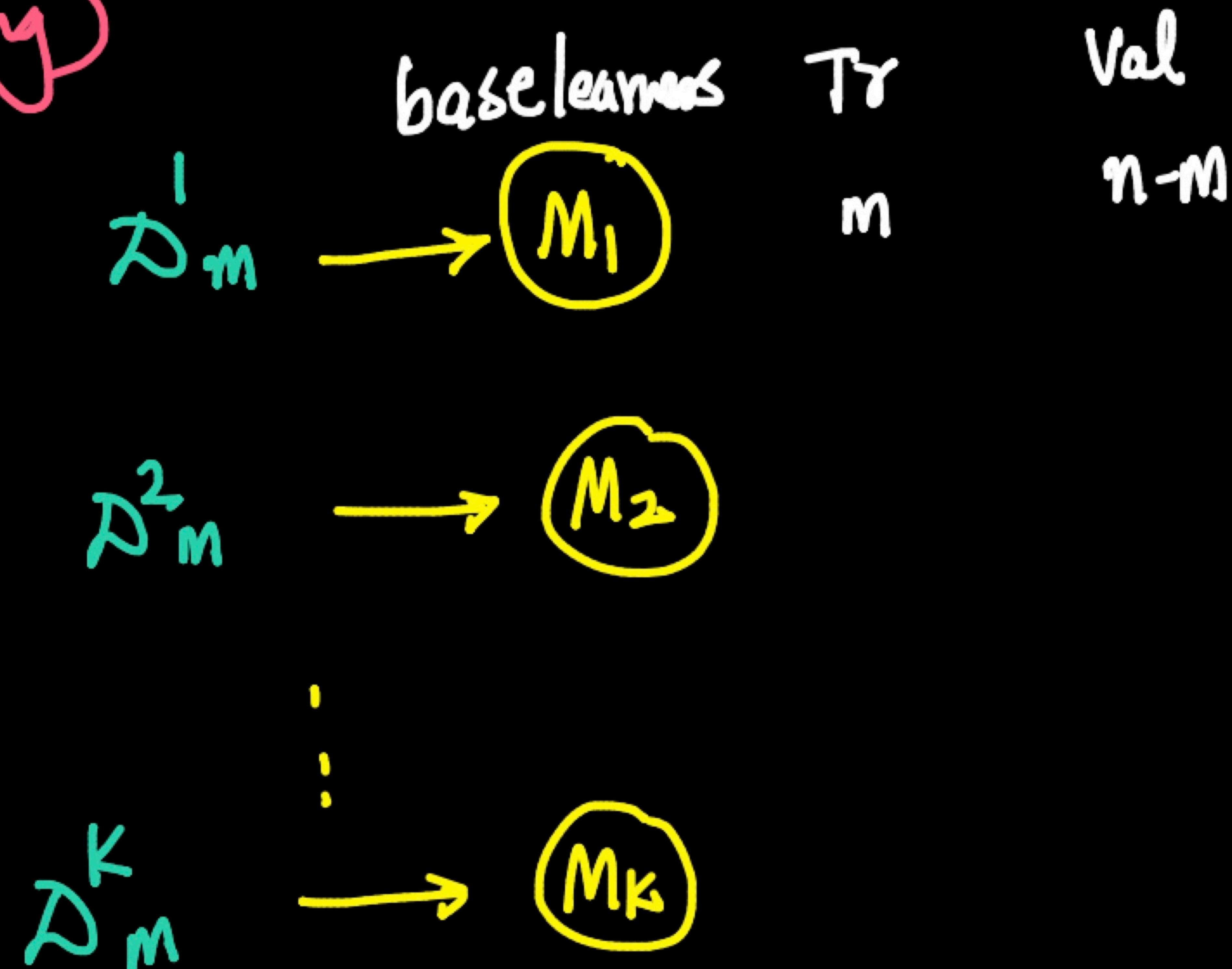


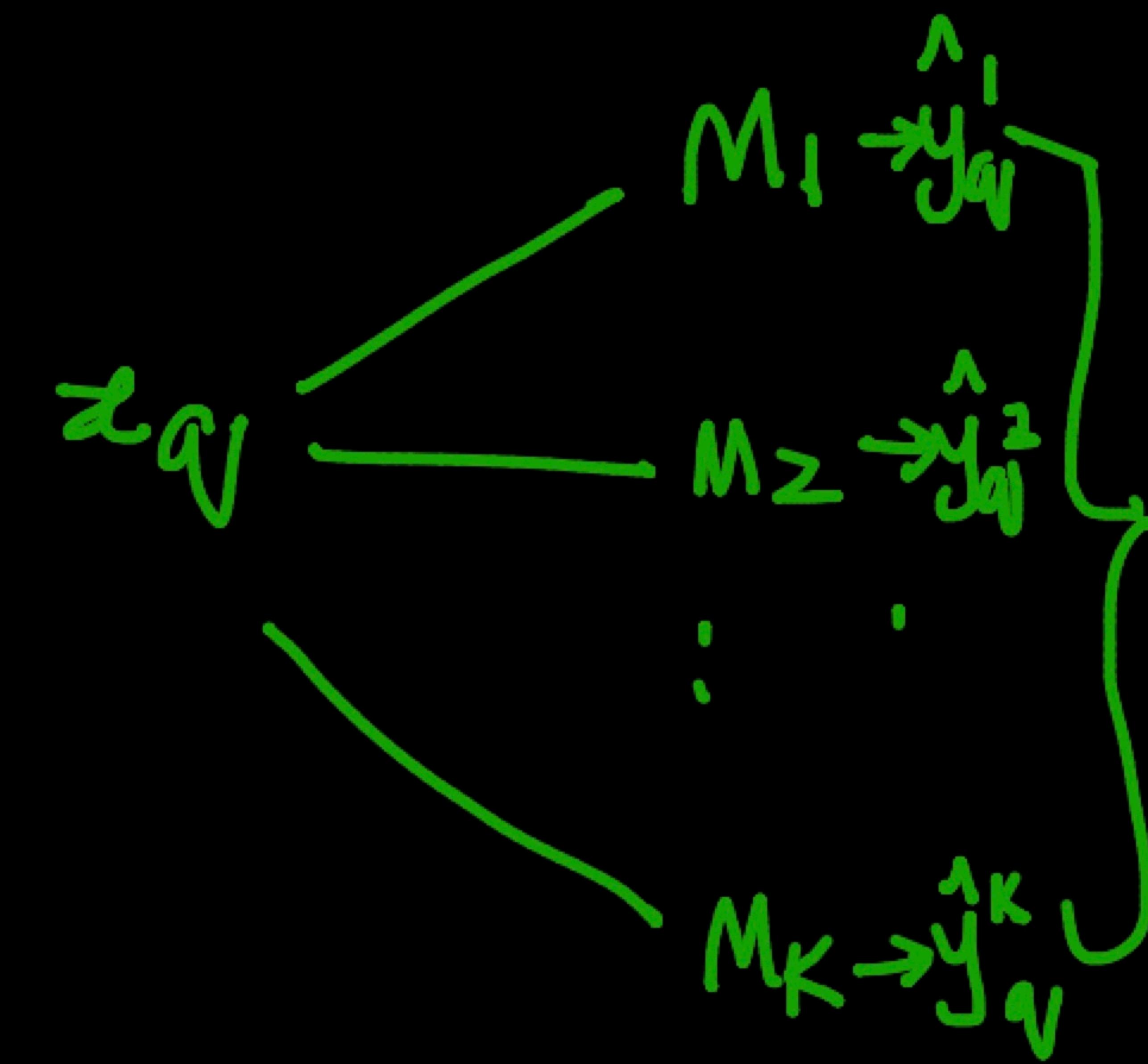


INTUITION

100K Million
 $m < n$; $K=100$

Bagging.

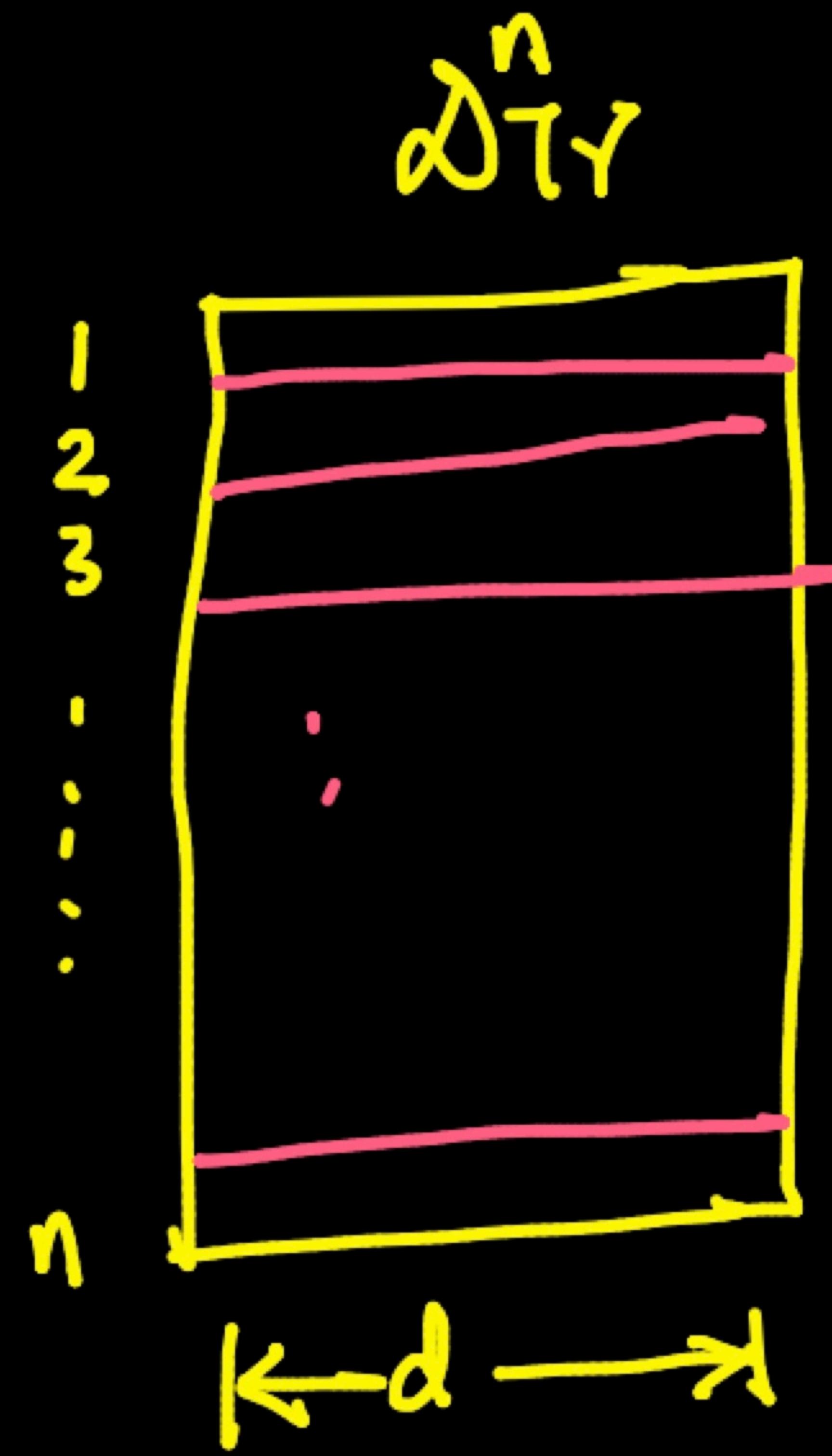




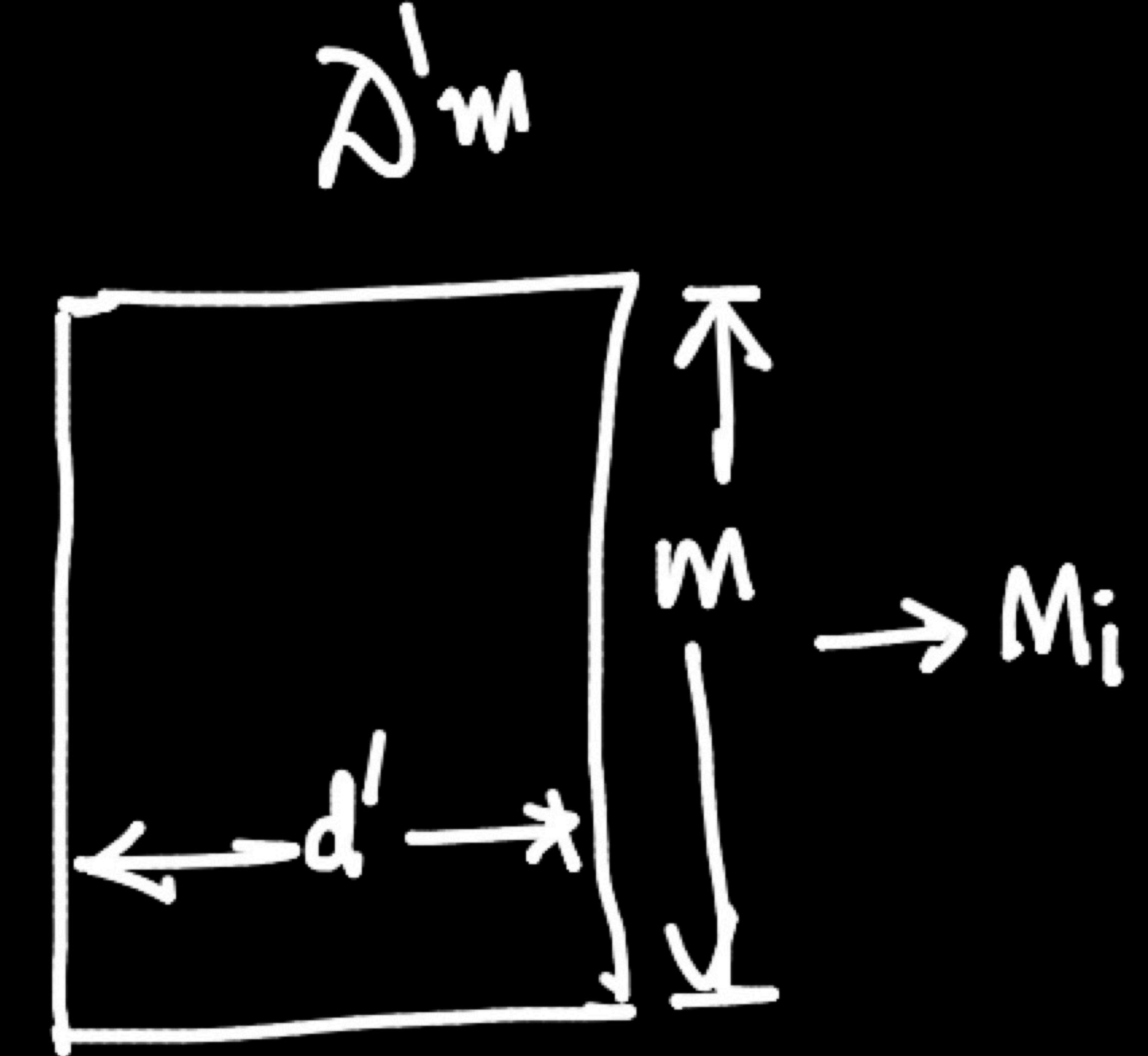
aggregation : (Classif.) Majority vote
(Reg.: Mean/Median)

if $m \ll n$ then M_i 's are more dissimilar

if $m \approx n$; $m < n$ M_i 's tend to become more similar

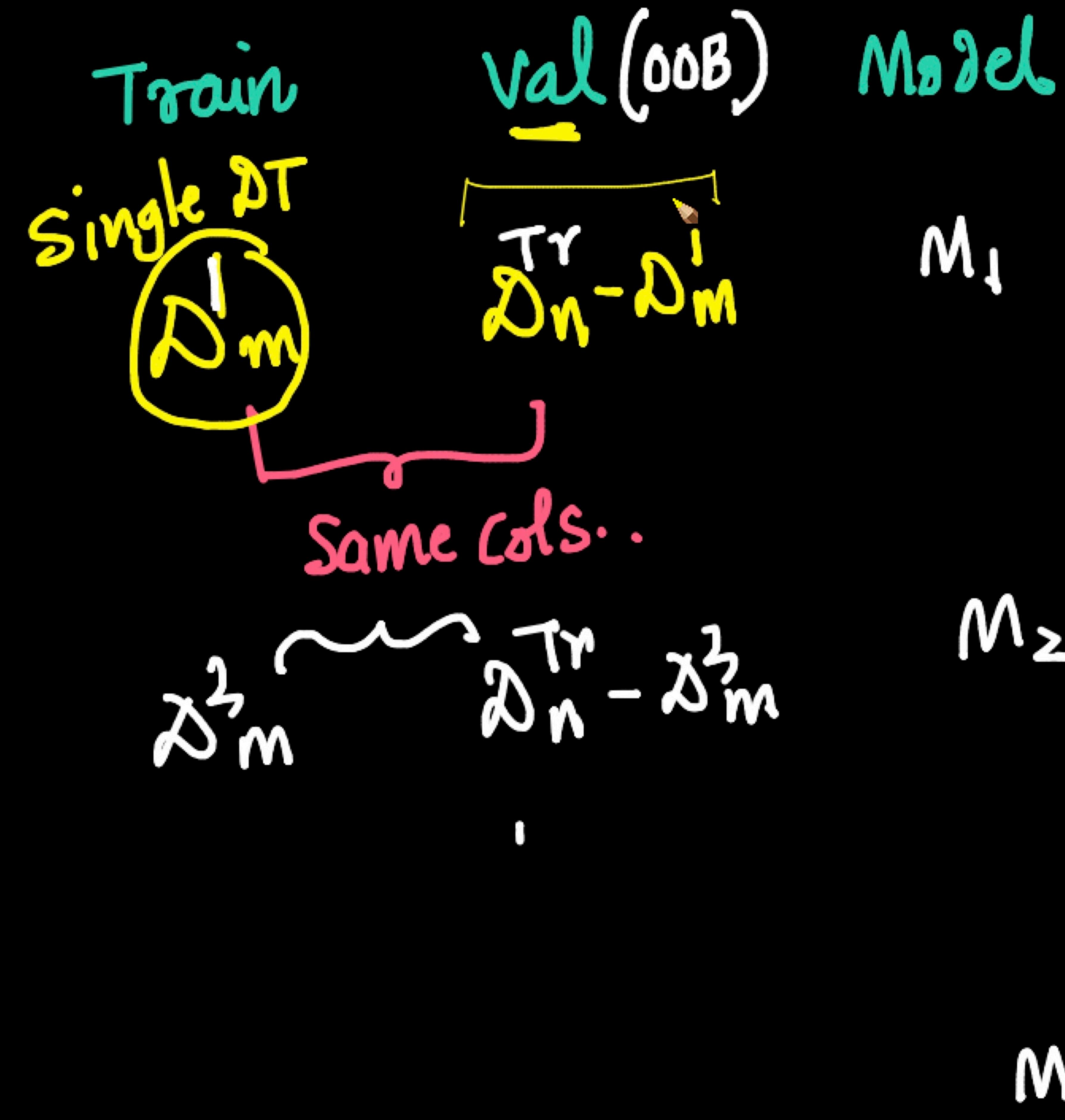
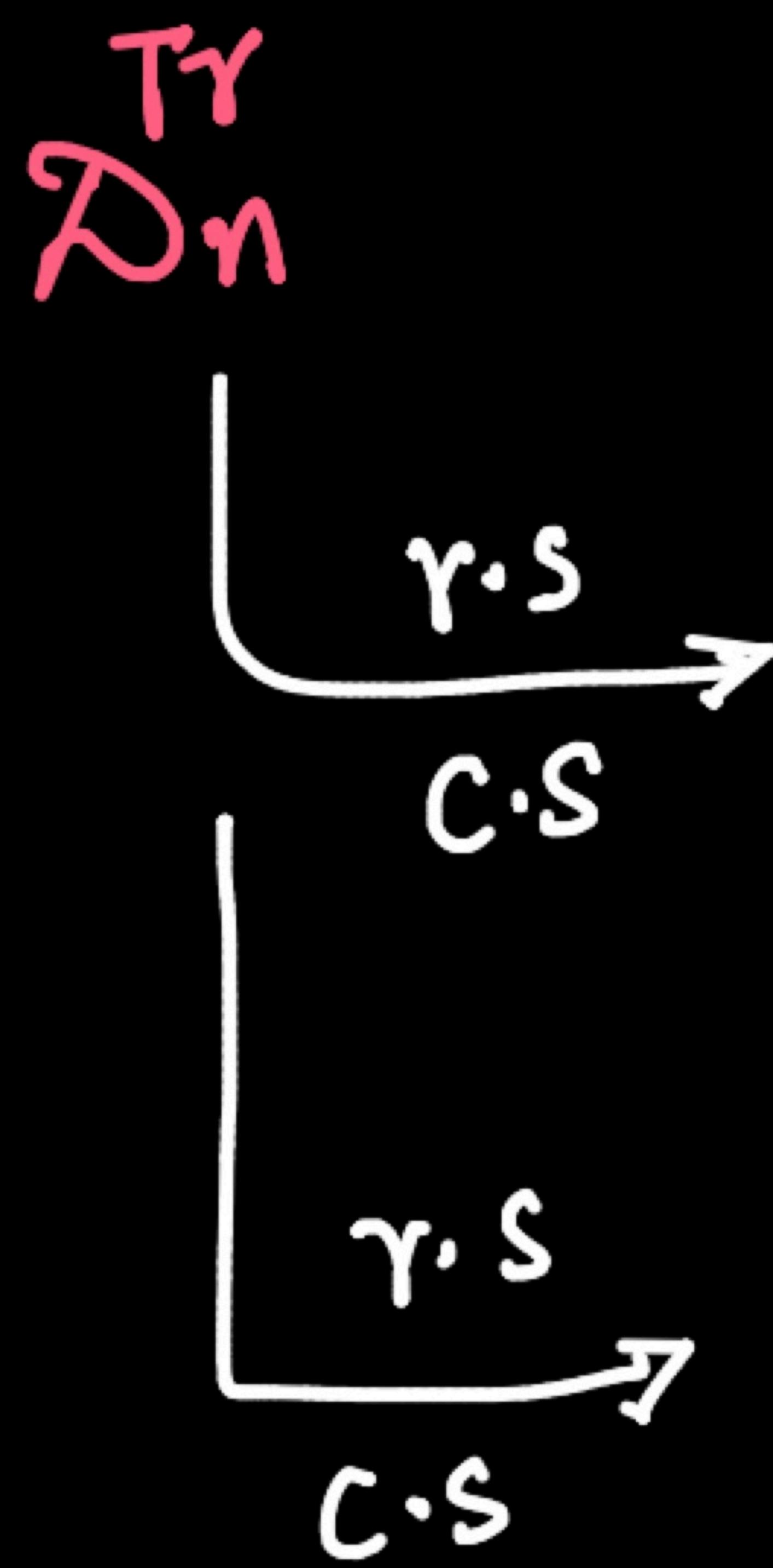


$M < n$
(M-rows)
row - sampling
column - samp
 $d' < d$



↳ different base-models

Random forest
Set of DTs
bootstrapping & col. sampling



Random Forest

Training

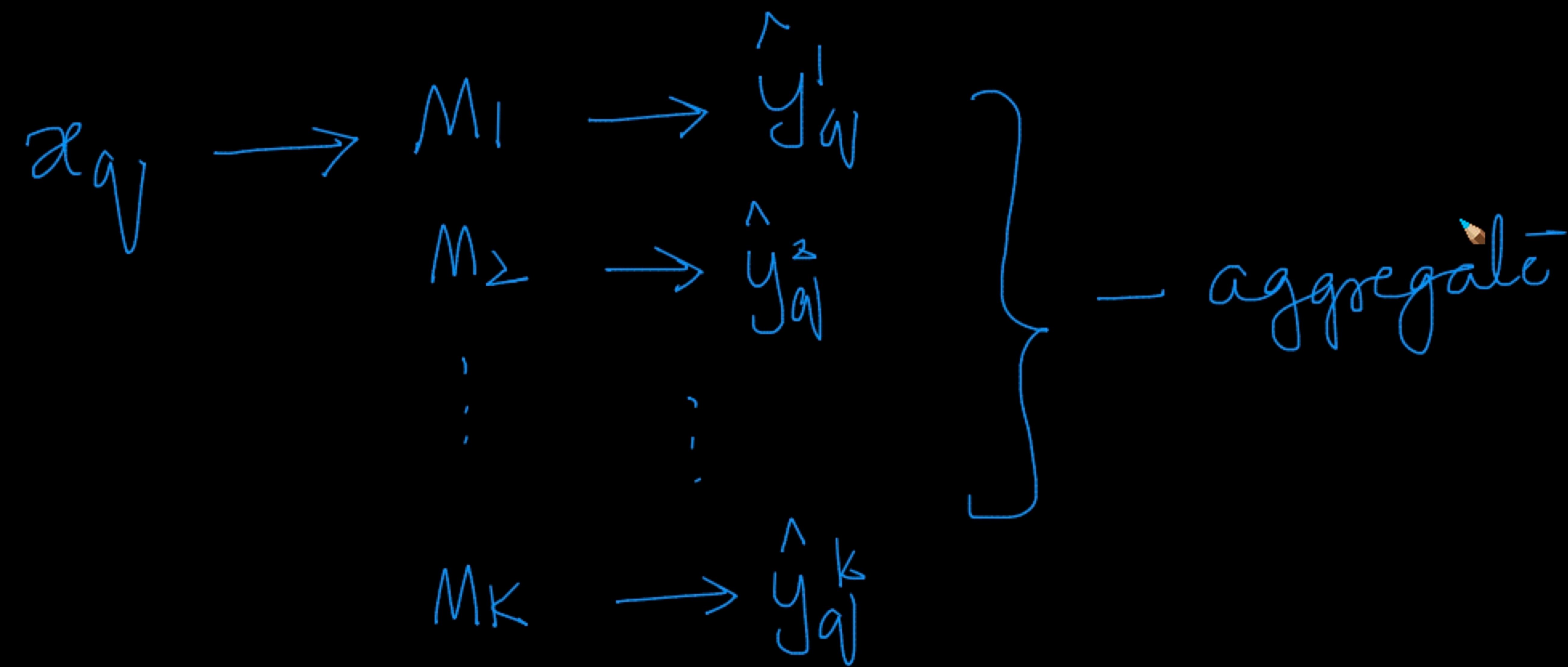
(Q)

Can't we create dissimilar base learners (DT)

using diff hyper-params

↳ Can do it
~~==>~~

Run-time



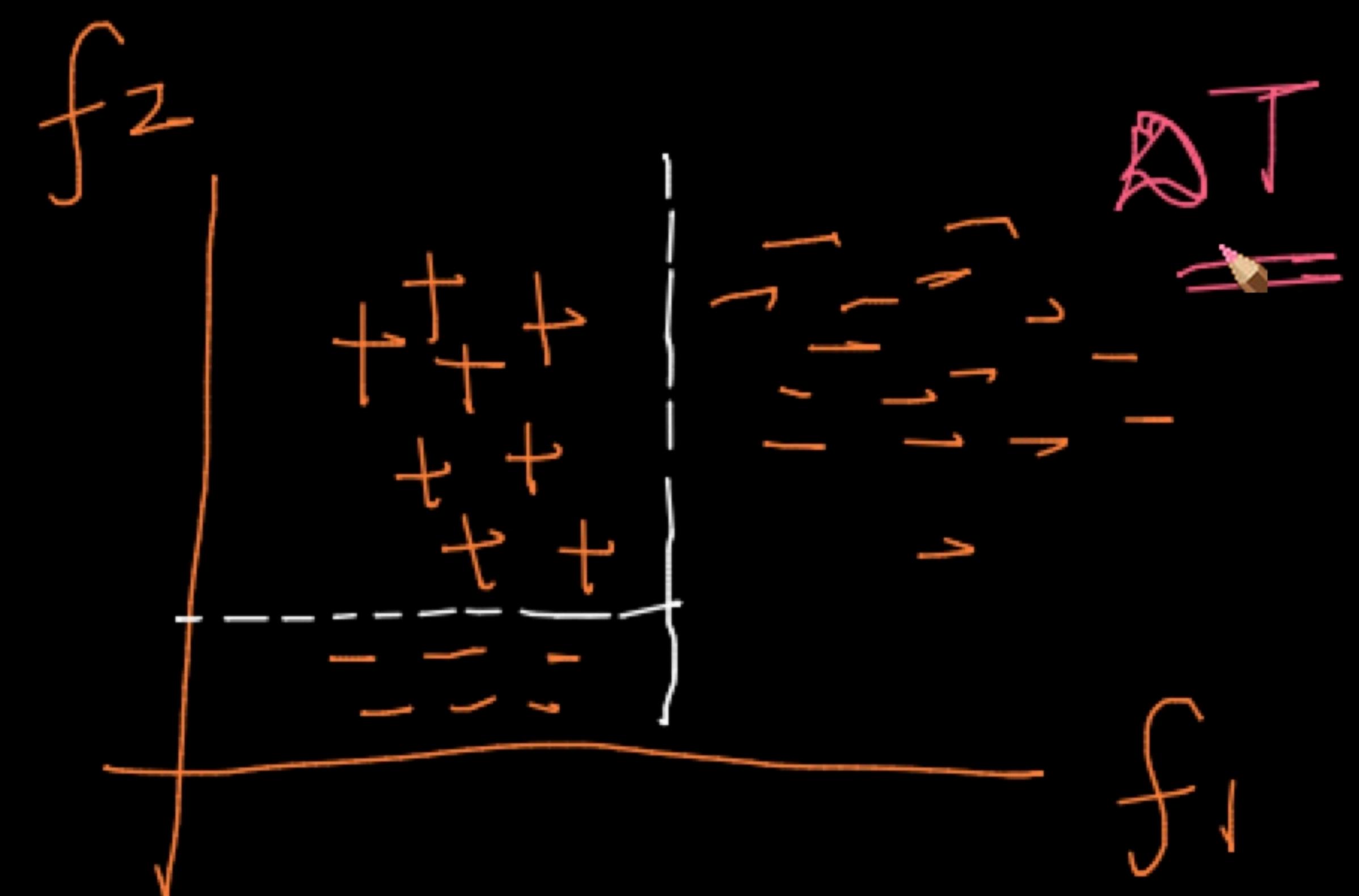
Practical

RF works better than DT (most cases)

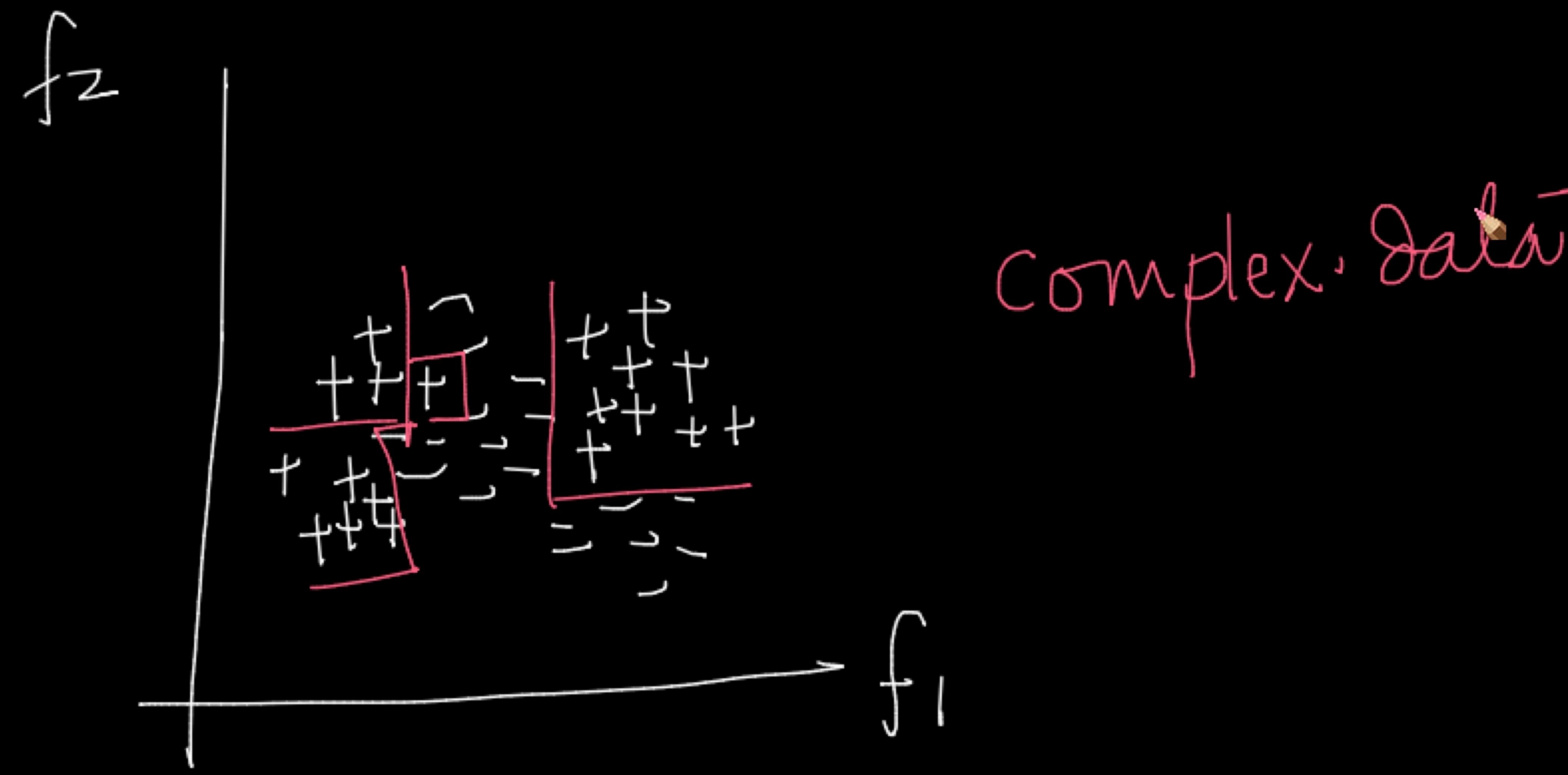
(Q) When would RF & DT have same performance?

$\rightarrow n = m \times \text{Col. Sampling}$

\rightarrow small datasets



Simple task



Hyperparams

RF

✓ $K = \# \text{ base-learners/trees}$

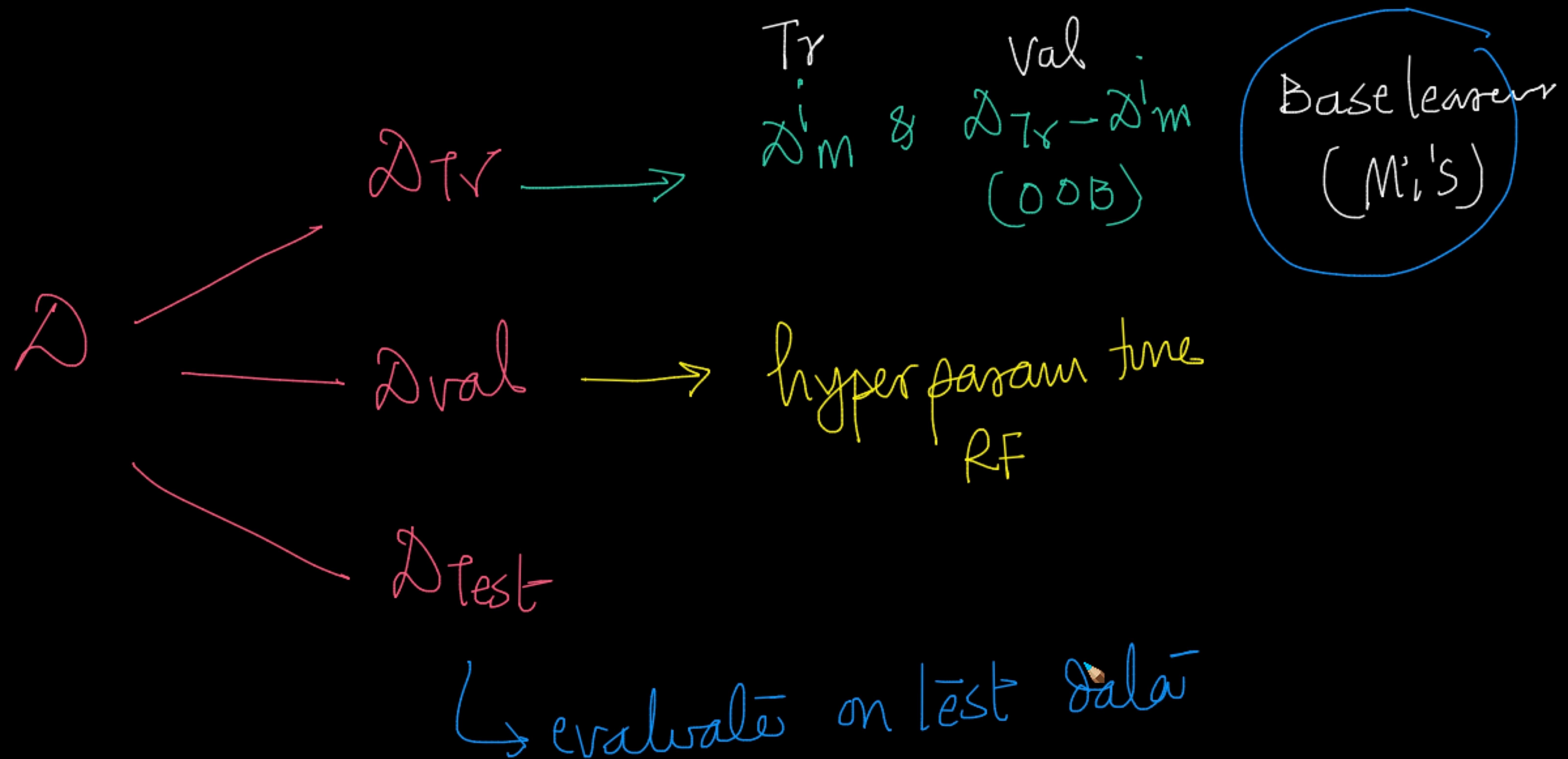
✓ Col. Sampling ratio = d'/d

✓ Row. Sampling ratio = m/n

bias-variance

tradeoff

[depth of DT \rightarrow ODB-val data]

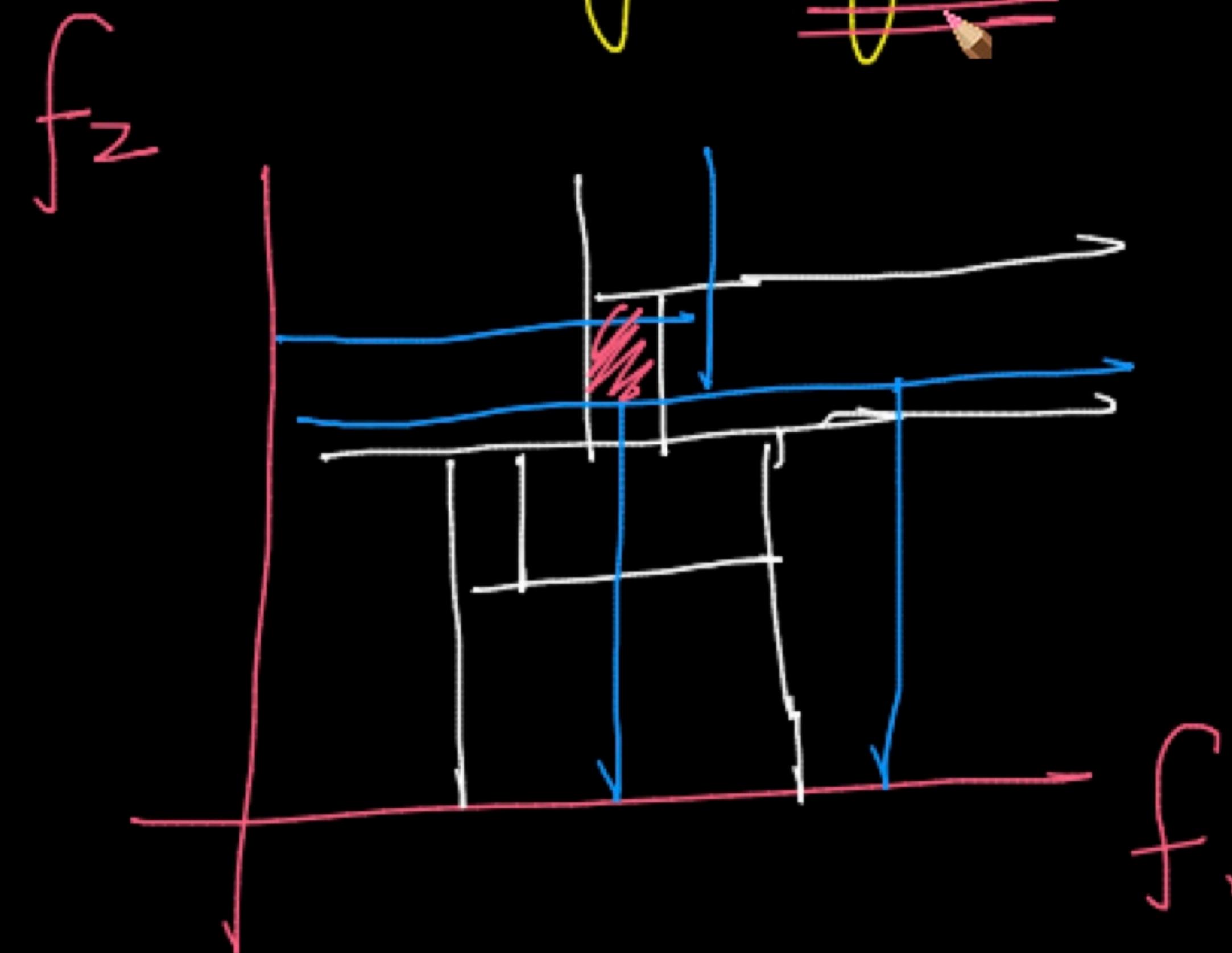


RF

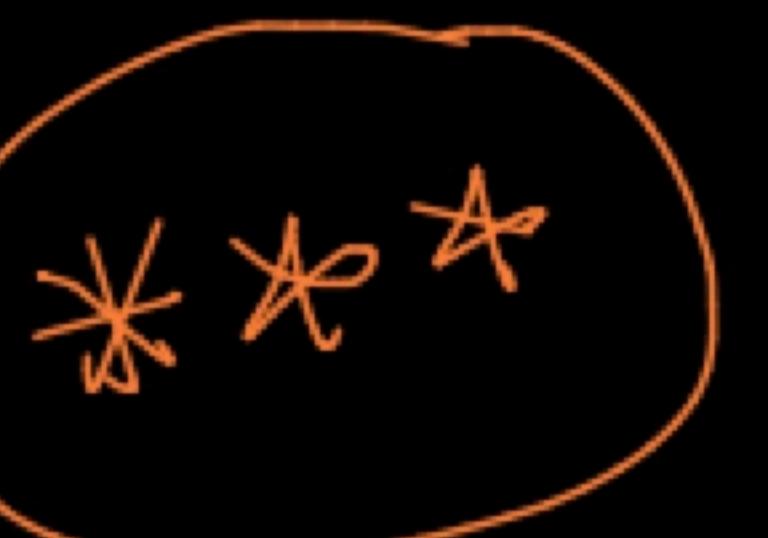
each base learner is slightly deep (e.g. $g=7$)

trained on a
subset of rows
& cols

↓
slightly overfit



Bagging



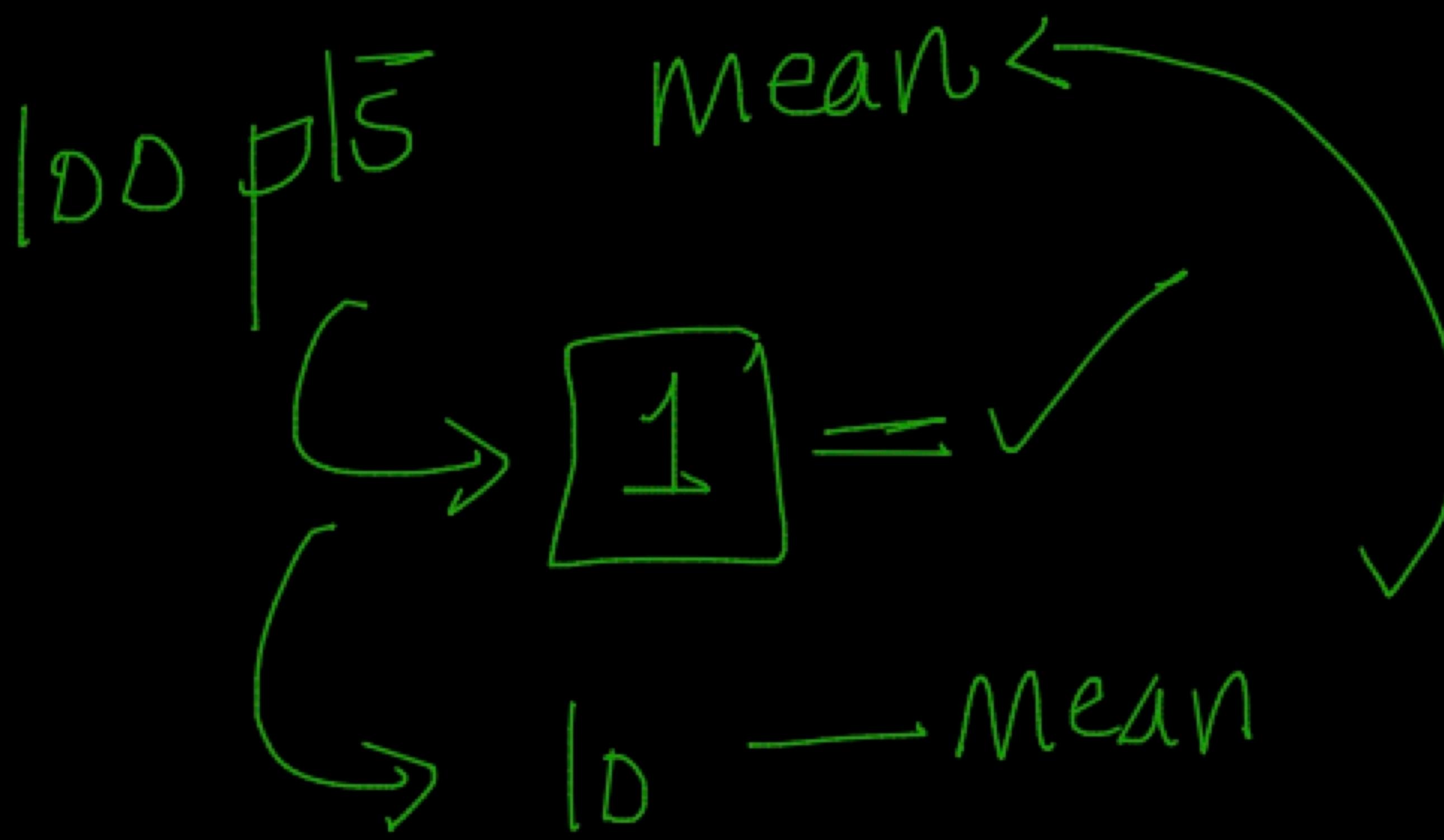
slightly overfit
base models

different from one another
(randomization)

→ aggregation

(majority vote
mean / median)

e.g.



results in
reducing of
overfitting

(Q)

$$\text{depth} = \begin{cases} 1 & \rightarrow M_1 \\ 2 & \rightarrow M_2 \\ 2 & \rightarrow M_3 \end{cases}$$

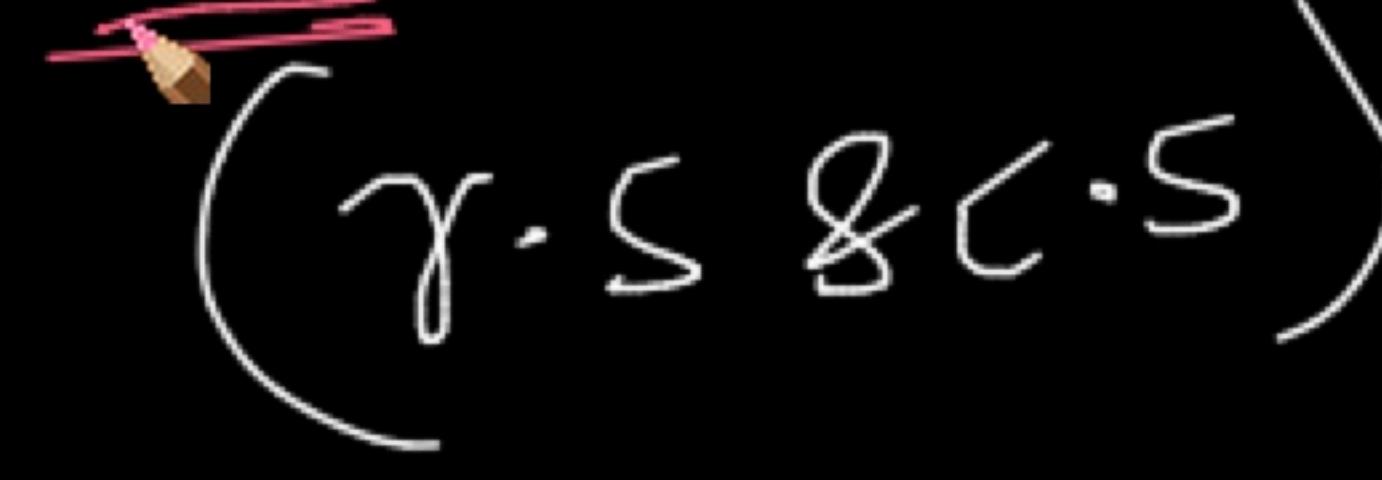
underfitting

$$3 \rightarrow M_4$$

$$3 \rightarrow M_5$$

;

overfitting

 VanDongen zali^m + aggregatum
 for

regularizali^m

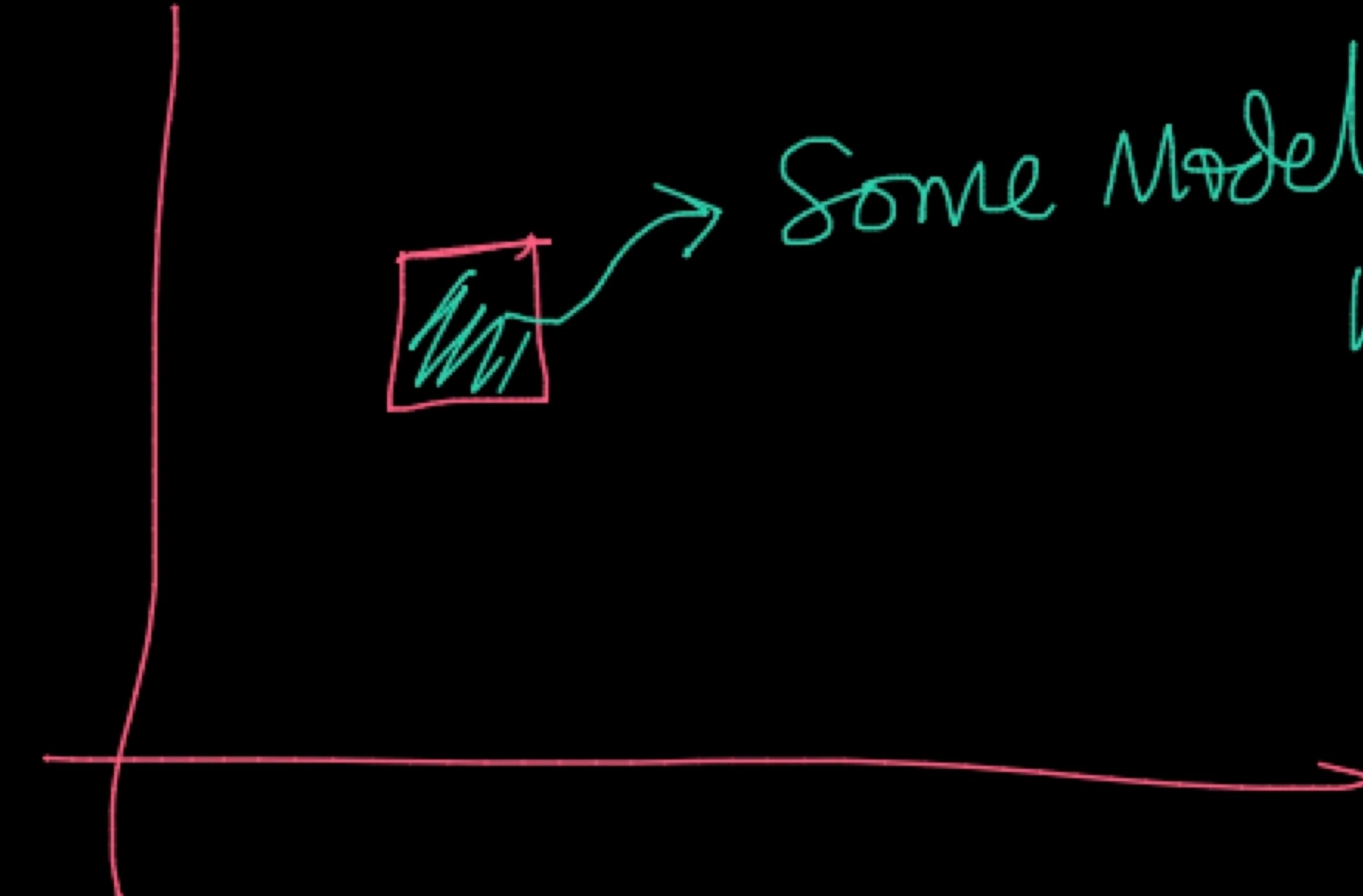
(Δ)

100's or 1000's

Slightly overfit
models

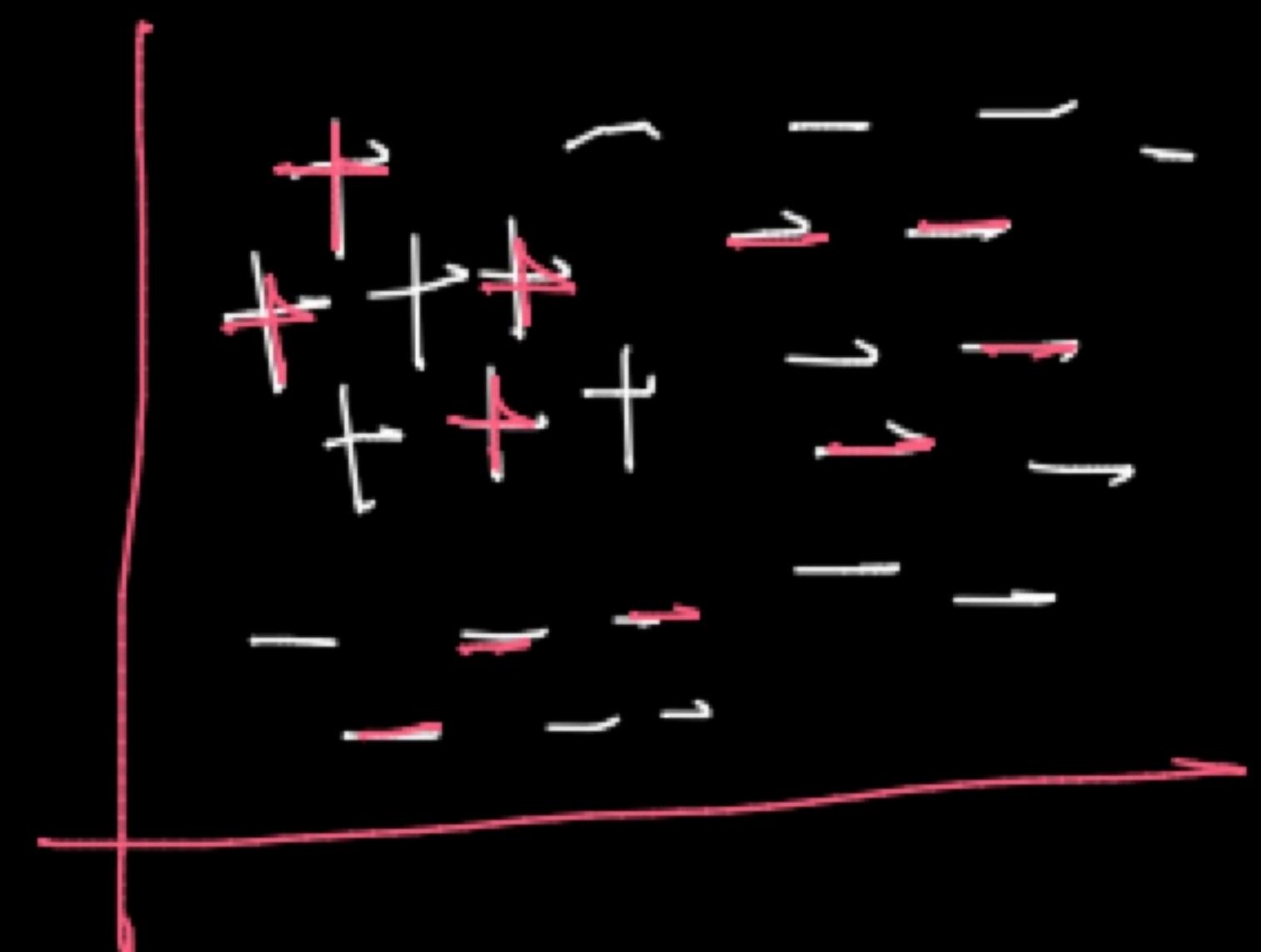
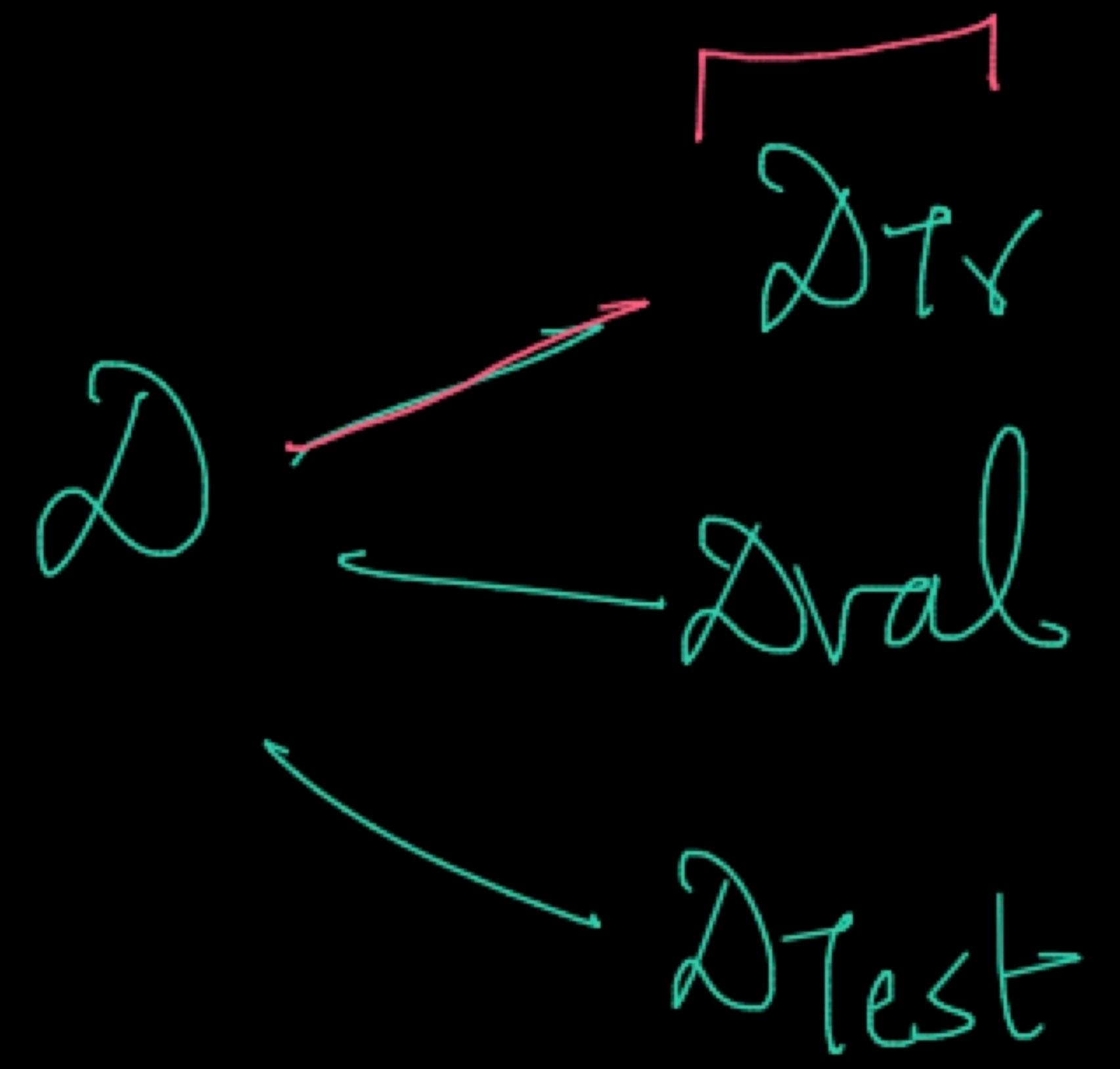
→ aggregating
↓

Some Models not all
will overfit here



(Q)

No optimization

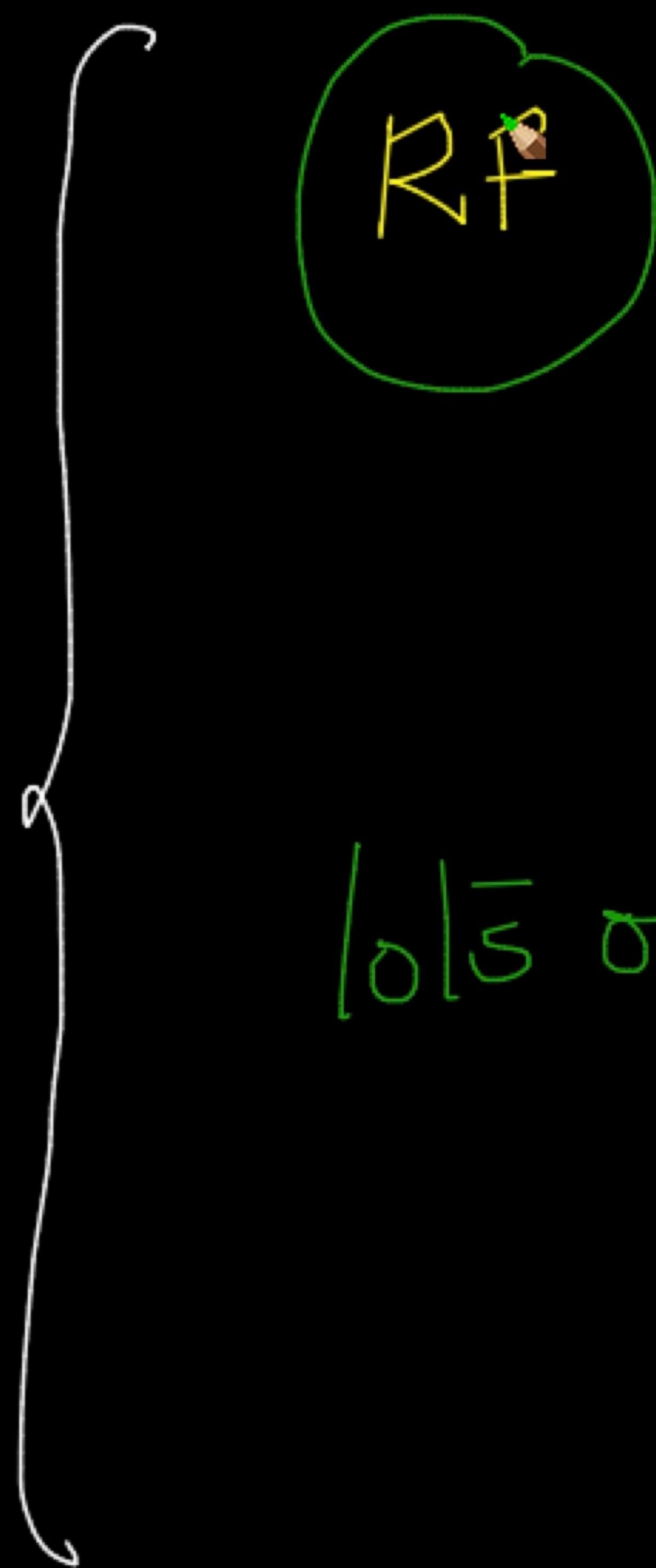


Lr. Regression

needs less data
to train as

Compared to DT / RFs





→ rigorous Math proof (Stanford Stats)
early 2000's

lots of slightly overfit base learners
(high-variance)

↓
aggregation : reduction in variance \Rightarrow less overfit

aggregatim +
randomizatim

100 heights

Obj:
Mean-height

1 random person = 155cm → high variance
(Overfitting)

10 random people → avg → lower variance

Model → picks 1 person
random ; O/P height

↓
155cm }
160cm }
175cm }

(Q)

Linear / logistic reg model as base learner

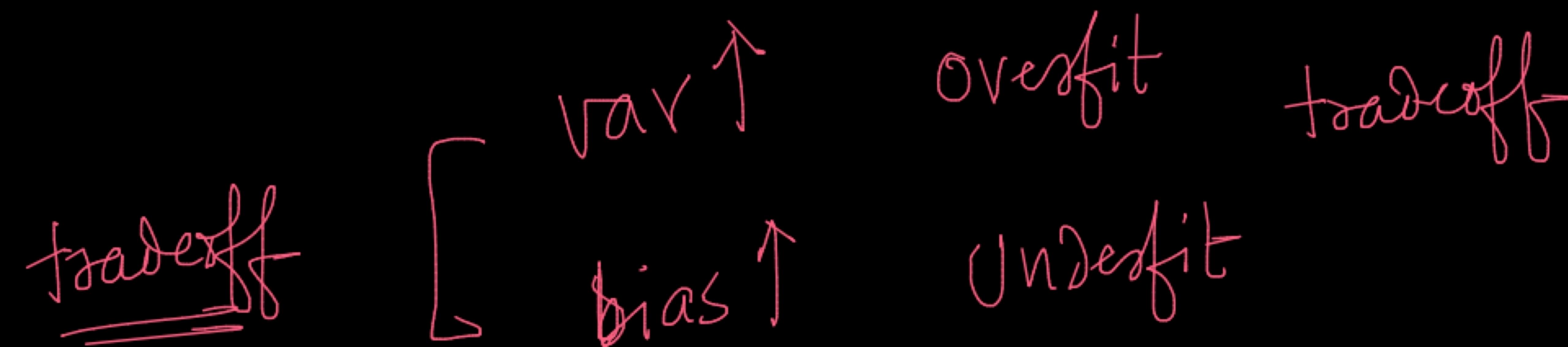
↳ Polynomial versions → slightly overfit
of Linear Model

✓ DT → model interpretable

Statistical ML

(Optimal Research Sessions)

$$\text{Error}_{\text{Model}} = \underbrace{\text{Bias}^2}_{\rightarrow} + \underbrace{\text{Var}}_{\downarrow} + \underbrace{\text{irreducible error}}_{\rightarrow}$$



INTUTION



{ variance \rightarrow \approx



how much variability would exist
in my model if I chose
different subsets of \mathcal{D} to train

on

{ bias \rightarrow \approx

how simple the model is?

Linear Model

lowest bias

\longrightarrow Non linear Mod

deep - DT

high bias \approx

RF:

Bagging : RS + CS

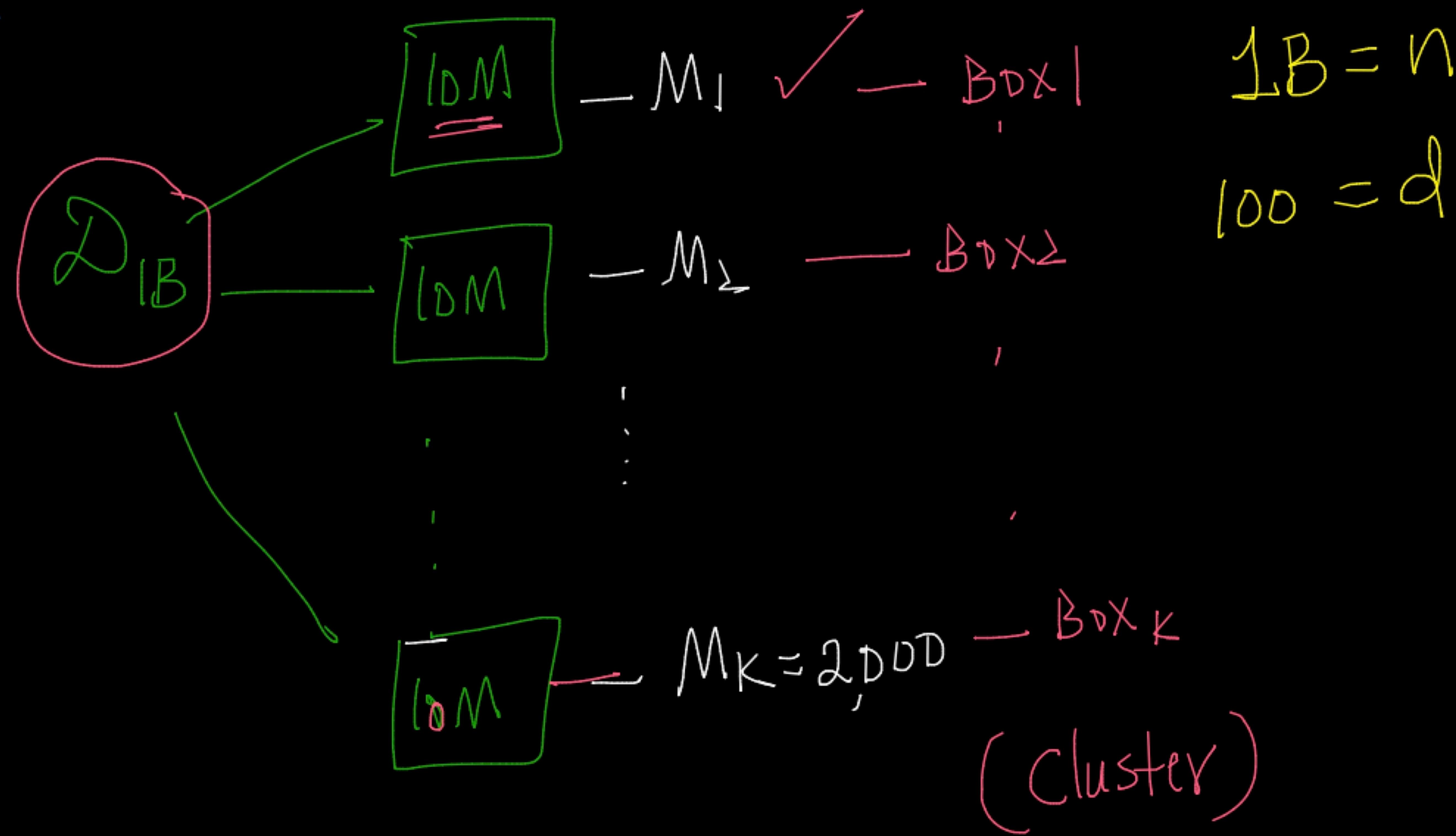
slightly overfit base-learners

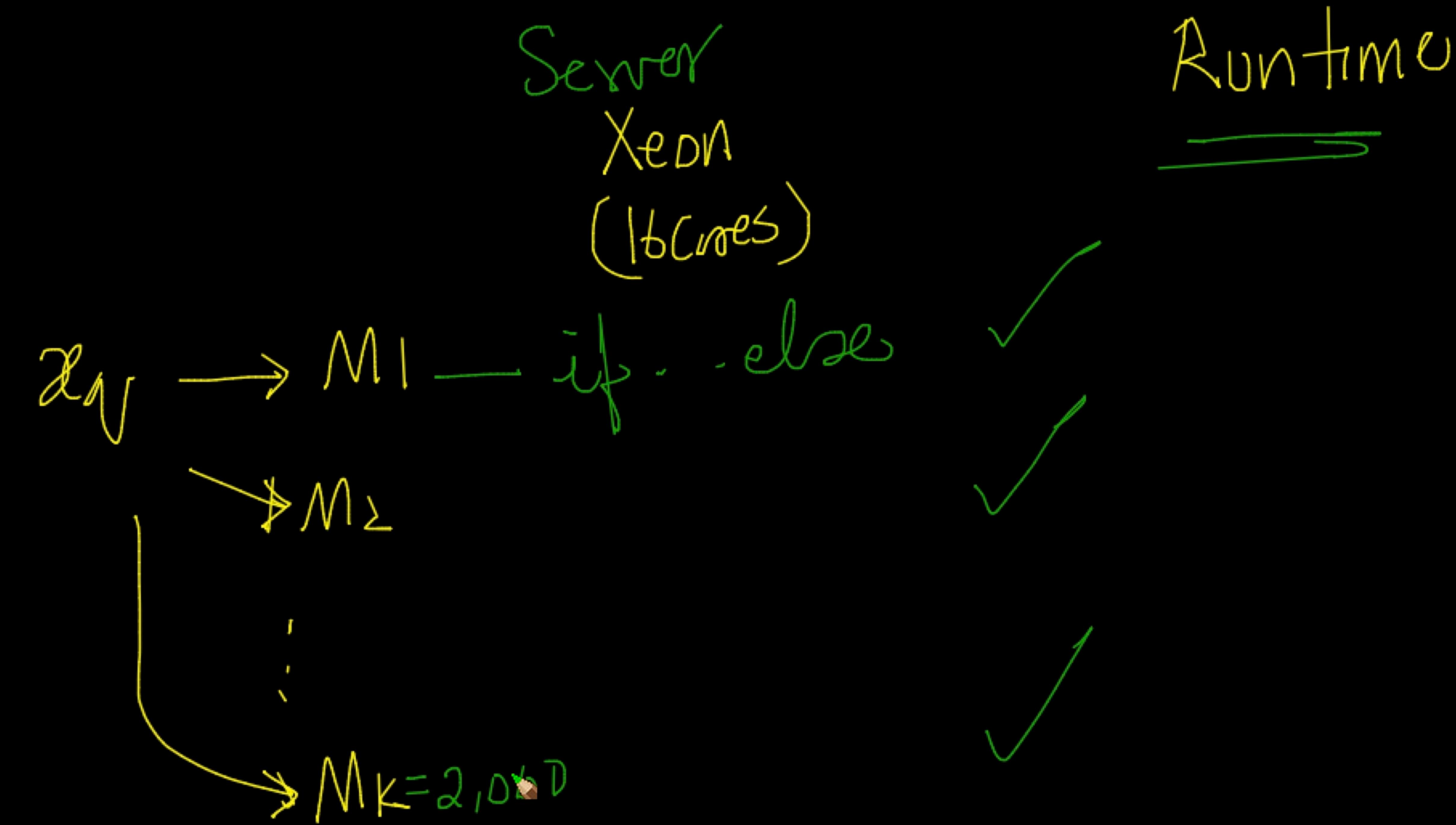
aggregating to several learners while keeping
the bias the
same as base
learners

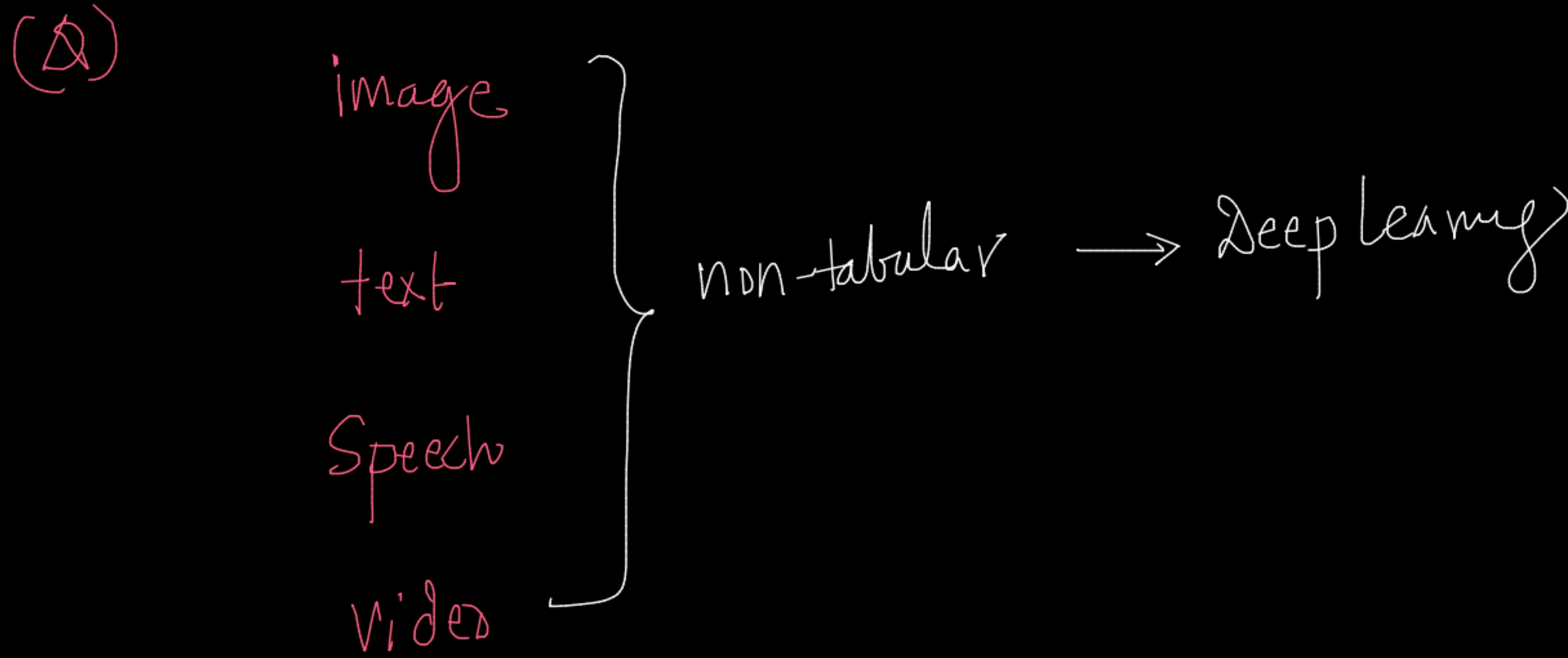
MWPS
(SparkML)

trivially parallelizable

Training







Tabular → Ensembles (RF / GBDT)

(Q) when not to use a RF?

→ Decision Trees?

→ too many features

Linearly
separable

Bagging_RF.ipynb - Colaboratory | Random forest - Wikipedia | sklearn.ensemble.RandomForestClassifier

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

Go



Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.2.1

Other versions

Please [cite us](#) if you use the software.

[sklearn.ensemble.RandomForestClassifier](#)

[RandomForestClassifier](#)

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='sqrt', 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]

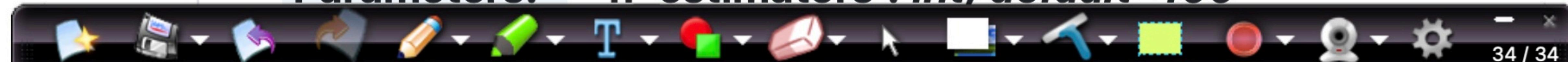
A random forest classifier.

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`

Toggle Menu



Bagging_RF.ipynb - Colaboratory | Random forest - Wikipedia | sklearn.ensemble.RandomForestClassifier

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

Go



Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.2.1

Other versions

Please [cite us](#) if you use the software.

[sklearn.ensemble.RandomForestClassifier](#)

[RandomForestClassifier](#)

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='sqrt', 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]

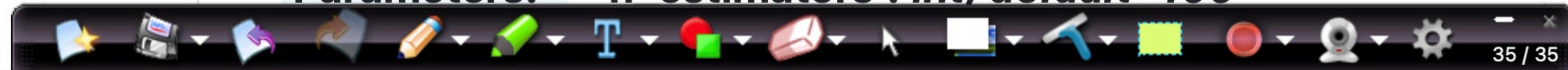
A random forest classifier.

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`

Toggle Menu



Bagging_RF.ipynb - Colaboratory | Random forest - Wikipedia | sklearn.ensemble.RandomForestClassifier

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

Prev Up Next

scikit-learn 1.2.1 Other versions

Please cite us if you use the software.

sklearn.ensemble.RandomForestClassifier

RandomForestClassifier

Examples using

sklearn.ensemble.RandomForestCl

(m)

0.7

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:

estimator_ : *DecisionTreeClassifier*
Estimator used to grow the ensemble.

base_estimator_ : *DecisionTreeClassifier*

Toggle Menu

scikit-learn 1.2.1

Other versions

Please **cite us** if you use the
software.

sklearn.ensemble.RandomForestClassifier

RandomForestClassifier

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='sqrt', 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)
```

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

Bagging_RF.ipynb - Colaboratory Random forest - Wikipedia sklearn.ensemble.RandomForestClassifier

Automatic saving failed. This file was updated remotely or in another tab. [Show diff](#)

+ Code + Text Saving failed since 20:32 RAM Disk

1s ✓ print(f"K-Fold Accuracy Mean: Train: {cv_acc_results['train_score'].mean()*100} Validation: {cv_acc_results['validation_score'].mean()*100}")
print(f"K-Fold Accuracy Std: Train: {cv_acc_results['train_score'].std()*100} Validation: {cv_acc_results['validation_score'].std()*100}")

{x} K-Fold Accuracy Mean: Train: 93.24203753064435 Validation: 84.20945945945945
K-Fold Accuracy Std: Train: 0.8887486809601636 Validation: 8.453236055736403

5s [9] from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import KFold, cross_validate

tree_clf = RandomForestClassifier(random_state=7, max_depth=7, n_estimators=100)
kfold = KFold(n_splits=10)
cv_acc_results = cross_validate(tree_clf, X_sm, y_sm, cv = kfold, scoring = 'accuracy', return_train_score=True)

print(f"K-Fold Accuracy Mean: Train: {cv_acc_results['train_score'].mean()*100} Validation: {cv_acc_results['validation_score'].mean()*100}")
print(f"K-Fold Accuracy Std: Train: {cv_acc_results['train_score'].std()*100} Validation: {cv_acc_results['validation_score'].std()*100}")

K-Fold Accuracy Mean: Train: 96.09794552823907 Validation: 89.78025851938895
K-Fold Accuracy Std: Train: 0.6293464780221507 Validation: 6.684410725948877

<>

Automatic saving failed. This file was updated remotely or in another tab. [Show diff](#)

38/38

+ Code + Text Saving failed since 20:32

✓ RAM Disk

```
[10] print(f"K-Fold Accuracy Mean: Train: {cv_acc_results['train_score'].mean()*100} Validation: {cv_acc_results['validation_score'].mean()*100}")
print(f"K-Fold Accuracy Std: Train: {cv_acc_results['train_score'].std()*100} Validation: {cv_acc_results['validation_score'].std()*100}")

K-Fold Accuracy Mean: Train: 93.24203753064435 Validation: 84.20945945945945
K-Fold Accuracy Std: Train: 0.8887486809601636 Validation: 8.453236055736403
```

41s

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import KFold, cross_validate

tree_clf = RandomForestClassifier(random_state=7, max_depth=7, n_estimators=1000)
kfold = KFold(n_splits=10)
cv_acc_results = cross_validate(tree_clf, X_sm, y_sm, cv = kfold, scoring = 'accuracy', return_train_error=True)

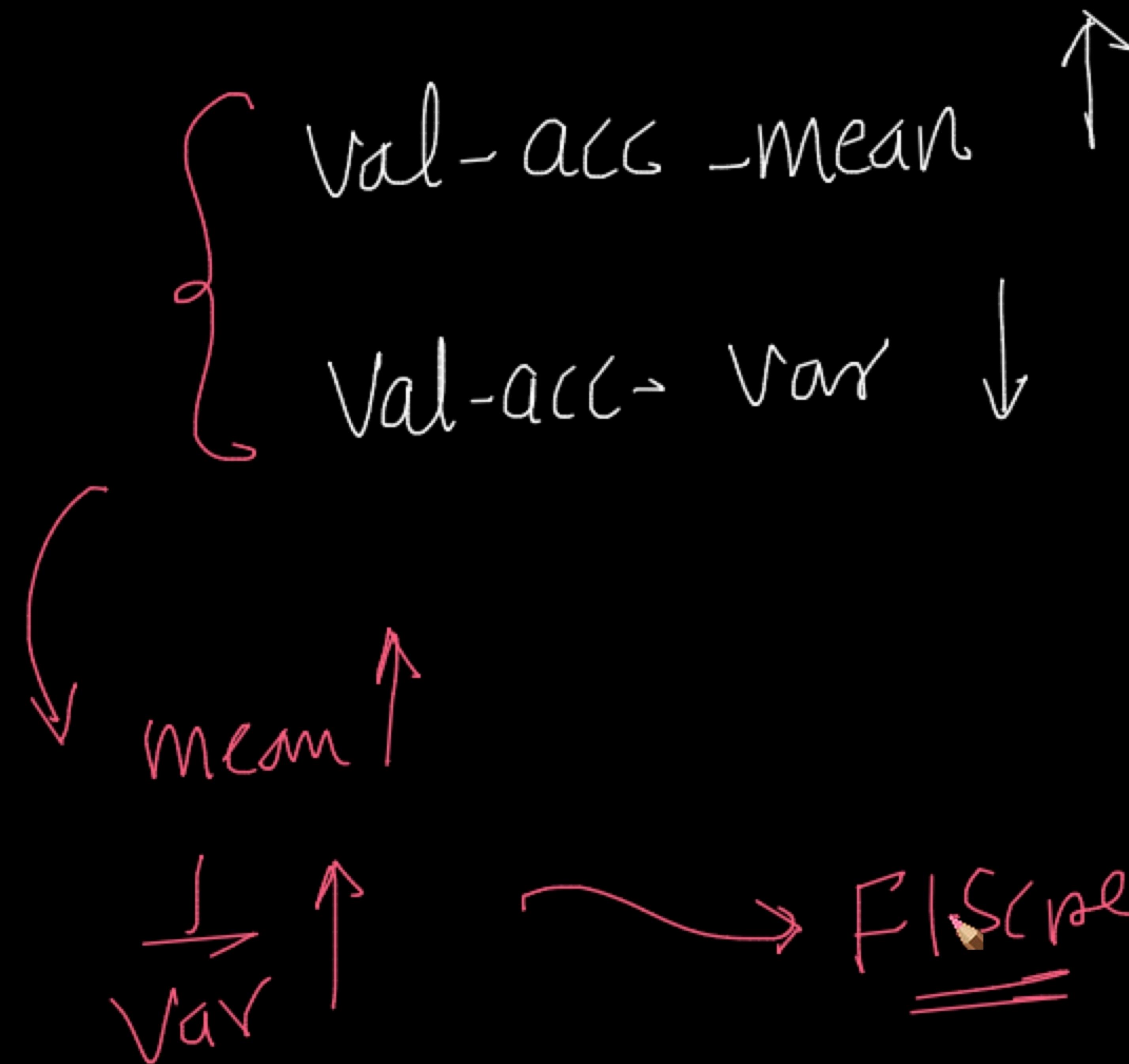
print(f"K-Fold Accuracy Mean: Train: {cv_acc_results['train_score'].mean()*100} Validation: {cv_acc_results['validation_score'].mean()*100}")
print(f"K-Fold Accuracy Std: Train: {cv_acc_results['train_score'].std()*100} Validation: {cv_acc_results['validation_score'].std()*100}")

K-Fold Accuracy Mean: Train: 96.23023656852769 Validation: 90.37573443008226
K-Fold Accuracy Std: Train: 0.5861570220746954 Validation: 6.764979471421996
```

Automatic saving failed. This file was updated remotely or in another tab. [Show diff](#)

Metric

K-fold CV

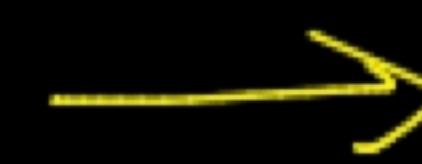


1 $\{ \text{mean} - 2\sqrt{\text{var}}$
= $\text{mean} + 2\sqrt{\text{var}} \}$

2 F1 Score: $P_T \uparrow R_T \uparrow$

Col. Samp = 0.1

f_i



M_1

M_2

take models where f_i
was considered

~~M_k~~

F_i^1 F_i^2

signals

{ → mean ✓
→ weighted mean:
→ medians

Weight = OOB performance

Bagging_RF.ipynb - Colaboratory Random forest - Wikipedia sklearn.ensemble.RandomForestClassifier

Automatic saving failed. This file was updated remotely or in another tab Show diff

+ Code + Text Saving failed since 20:32

[12] 41s

```
print("K-Fold Accuracy Mean: Train: {} Validation: {}".format(np.mean(cv_acc_results['train']), np.mean(cv_acc_results['validation'])))  
print("K-Fold Accuracy Std: Train: {} Validation: {}".format(np.std(cv_acc_results['train']), np.std(cv_acc_results['validation'])))
```

Feature importance

```
import numpy as np  
import matplotlib.pyplot as plt  
  
clf = RandomForestClassifier(random_state=7, max_depth=4, n_estimators=100)  
clf.fit(X_sm, y_sm)  
importances = clf.feature_importances_
```

→

```
indices = np.argsort(importances)[::-1] # Sort feature importances in descending order  
names = [X_sm.columns[i] for i in indices] # Rearrange feature names so they match the sorted feature importances  
plt.figure(figsize=(15, 7)) # Create plot  
plt.title("Feature Importance") # Create plot title  
plt.bar(range(X_sm.shape[1]), importances[indices]) # Add bars  
plt.xticks(range(X_sm.shape[1]), names, rotation=90) # Add feature names as x-axis labels  
plt.show() # Show plot
```

RAM Disk

hyper-params:

① $\# \text{ base learners} = K \uparrow$

Overfitting \downarrow

② $m \uparrow$

now samples \uparrow

Overfitting \uparrow

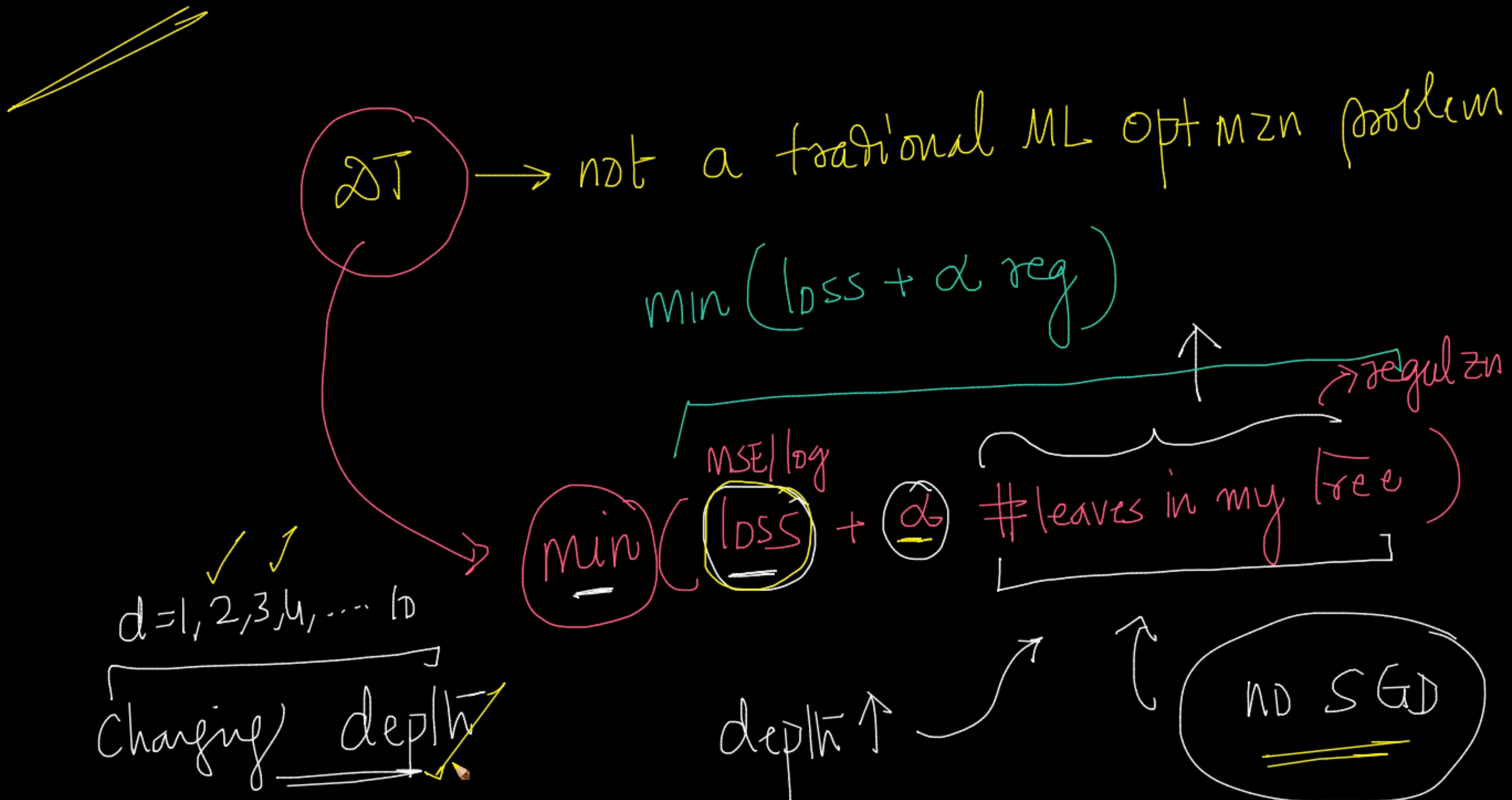
$m \uparrow \Rightarrow$ base models are
more similar

5

$d' = \# \text{cols sampled} \uparrow \Rightarrow \text{overfit} \uparrow$

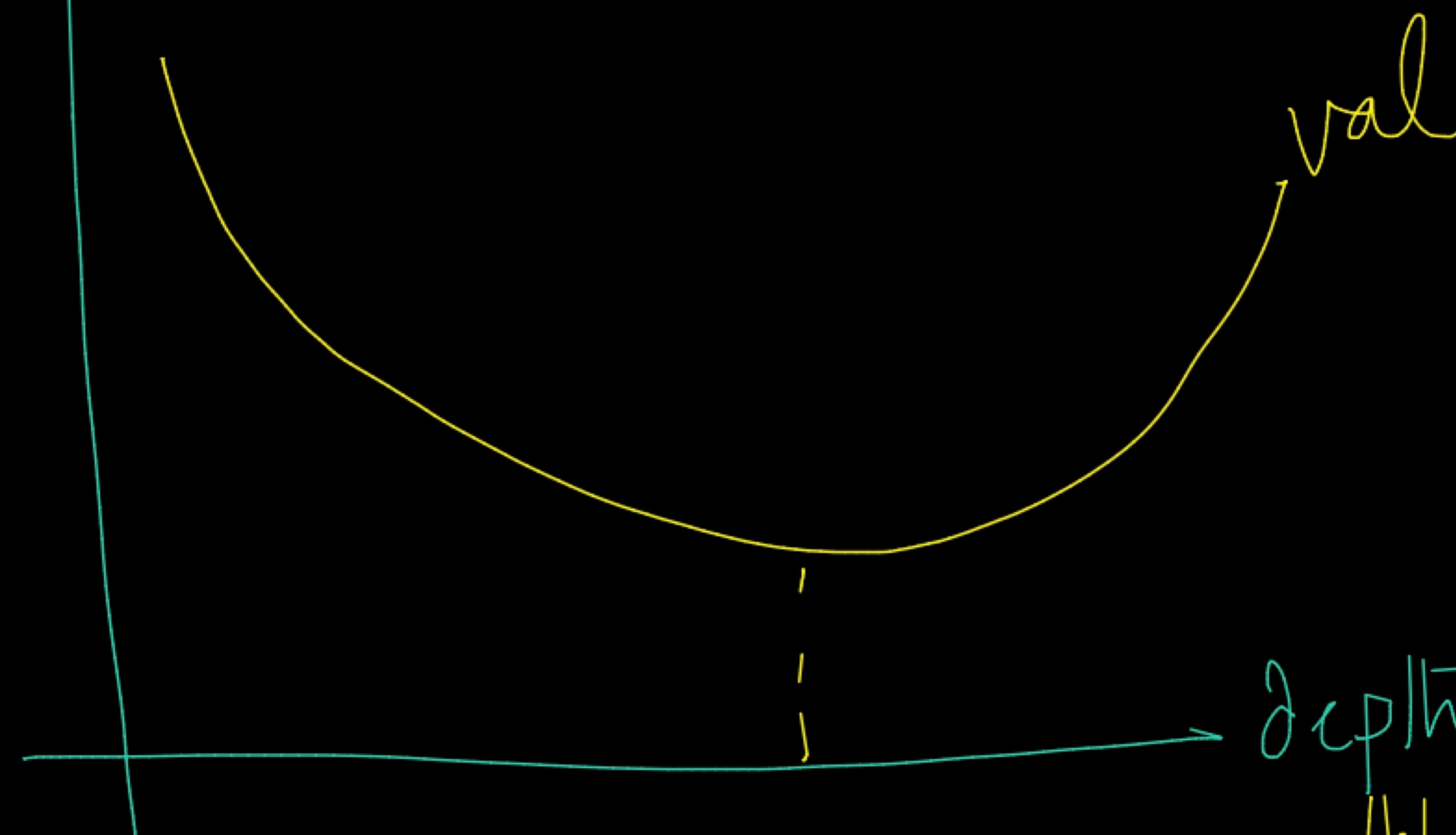
4

depth of base learners $\uparrow \Rightarrow \text{overfit} \uparrow$
(not using CV)

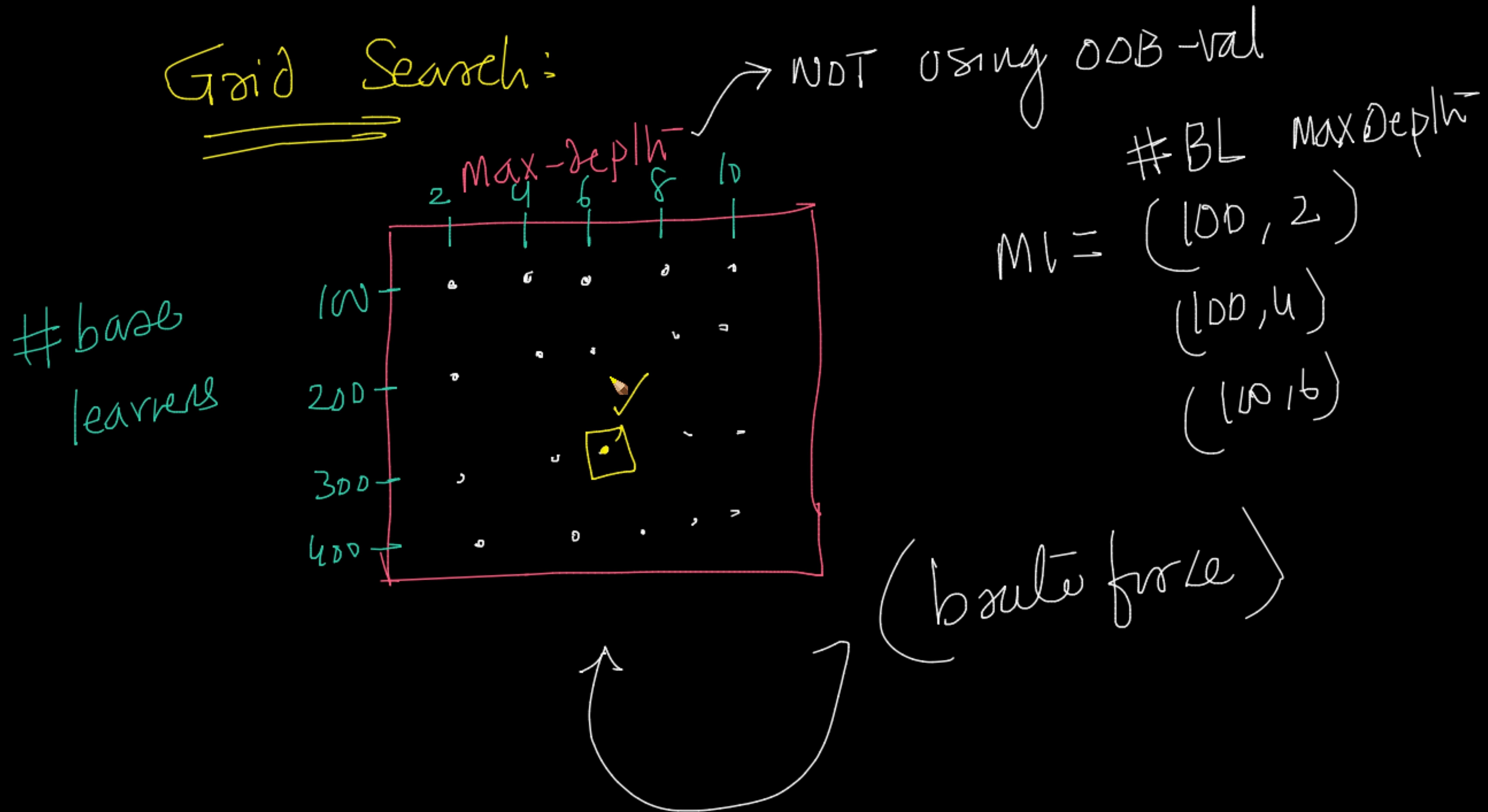


depth ↑ ⇒ #leaves ↑
⇒ (loss + d #leaves) ↑

7SS+



Depth of
leaf nodes



Automatic saving failed. This file was updated remotely or in another tab. Show diff

Bagging_RF.ipynb - Colaboratory Random forest - Wikipedia sklearn.ensemble.RandomForestClassifier

Saving failed since 20:32

[14] 0s

MaritalStatus_Single Job StockOptionLvl Department_Research & Development EducationF MaritalStatus_Married JobSatisfac JobL JobInvolvem MonthlyInc Department_Sci EnvironmentSatisfac BusinessTravel_Travel_Freque RelationshipSatisfac MaritalStatus_Divorced WorkLifeBalanc TotalWorkingY YearsWithCurrManag YearsInCurrent YearsAtComp DistanceFromHome OverTime TrainingTimesLast BusinessTravel_Travel_Rate MonthlyBusinessTravel_NonTrav HourlyBusinessTravel DailyBusinessTravel Ger NumCompaniesWorked PercentSalary YearsSinceLastPromotion Educa Department_HumanResources PerformanceRate RAM Disk

Defining Parameters

```
params = {
    'n_estimators' : [100,200,300,400],
    'max_depth' : [3,5,10],
    'criterion' : ['gini', 'entropy'],
    'bootstrap' : [True, False],
    'max_features' : [8,9,10]
}
```

Handwritten annotations on the code:

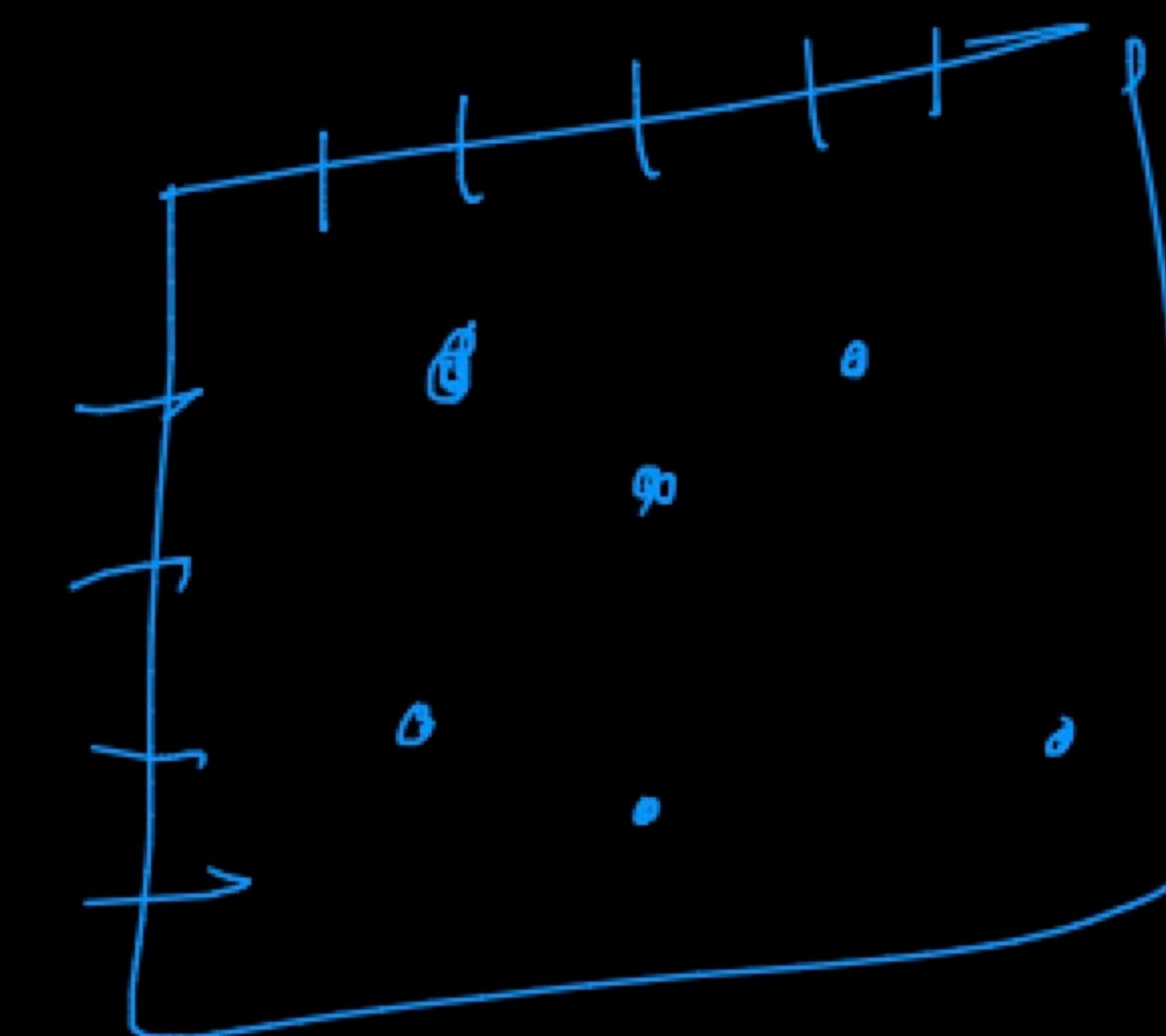
- An arrow points from the first '3' in the 'max_depth' list to a circled '3'.
- An arrow points from the second '3' in the 'max_depth' list to a circled '2'.
- An arrow points from the first '2' in the 'max_features' list to a circled '2'.
- An arrow points from the second '2' in the 'max_features' list to a circled '3'.
- A large bracket on the right side groups the numbers 144, with a checkmark and a pencil icon above it.

Grid Search → brute force (try all combinations)

✗ Binary Search → sorted

Randomized Search

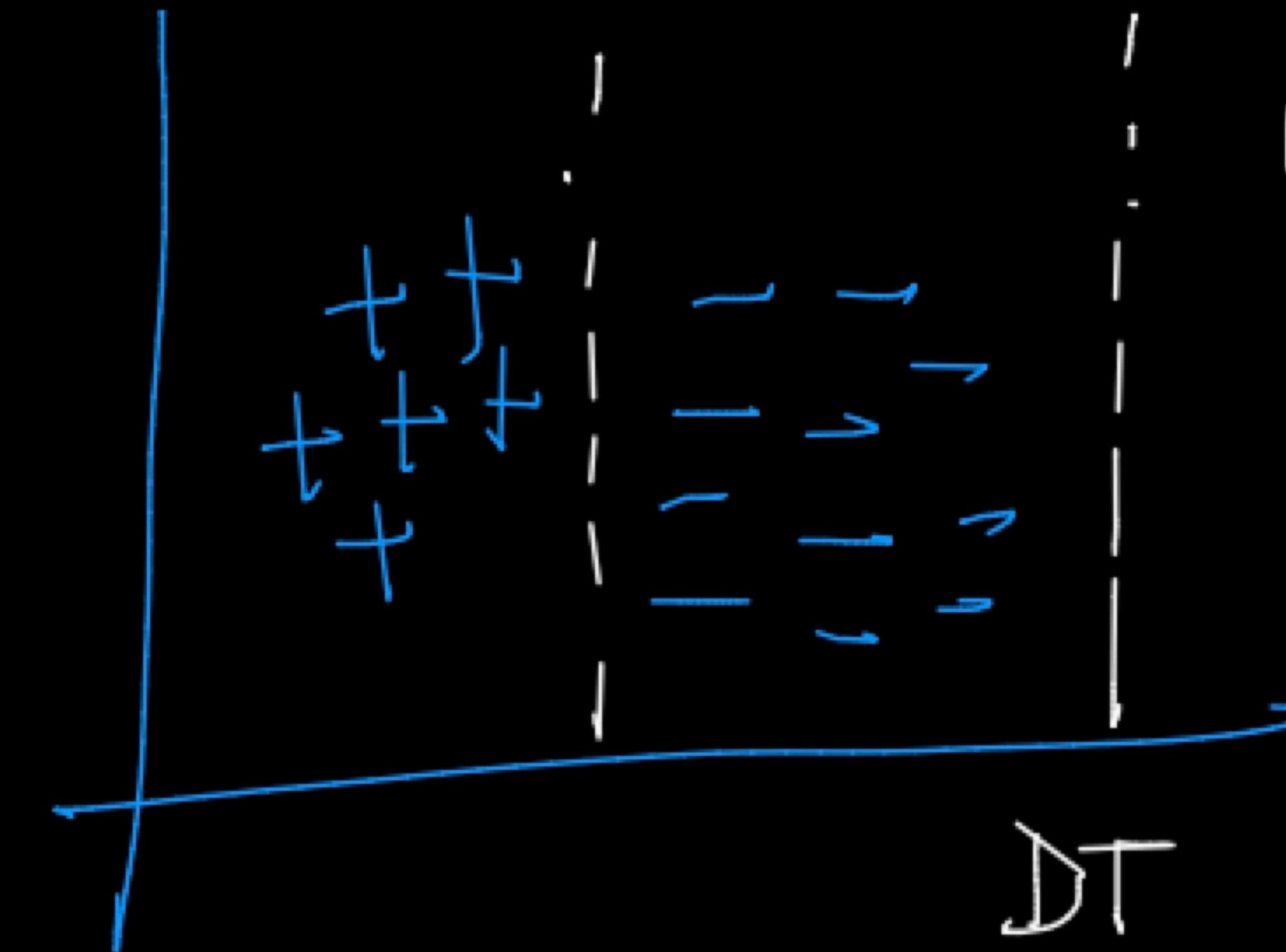
random combinations



INTUITION

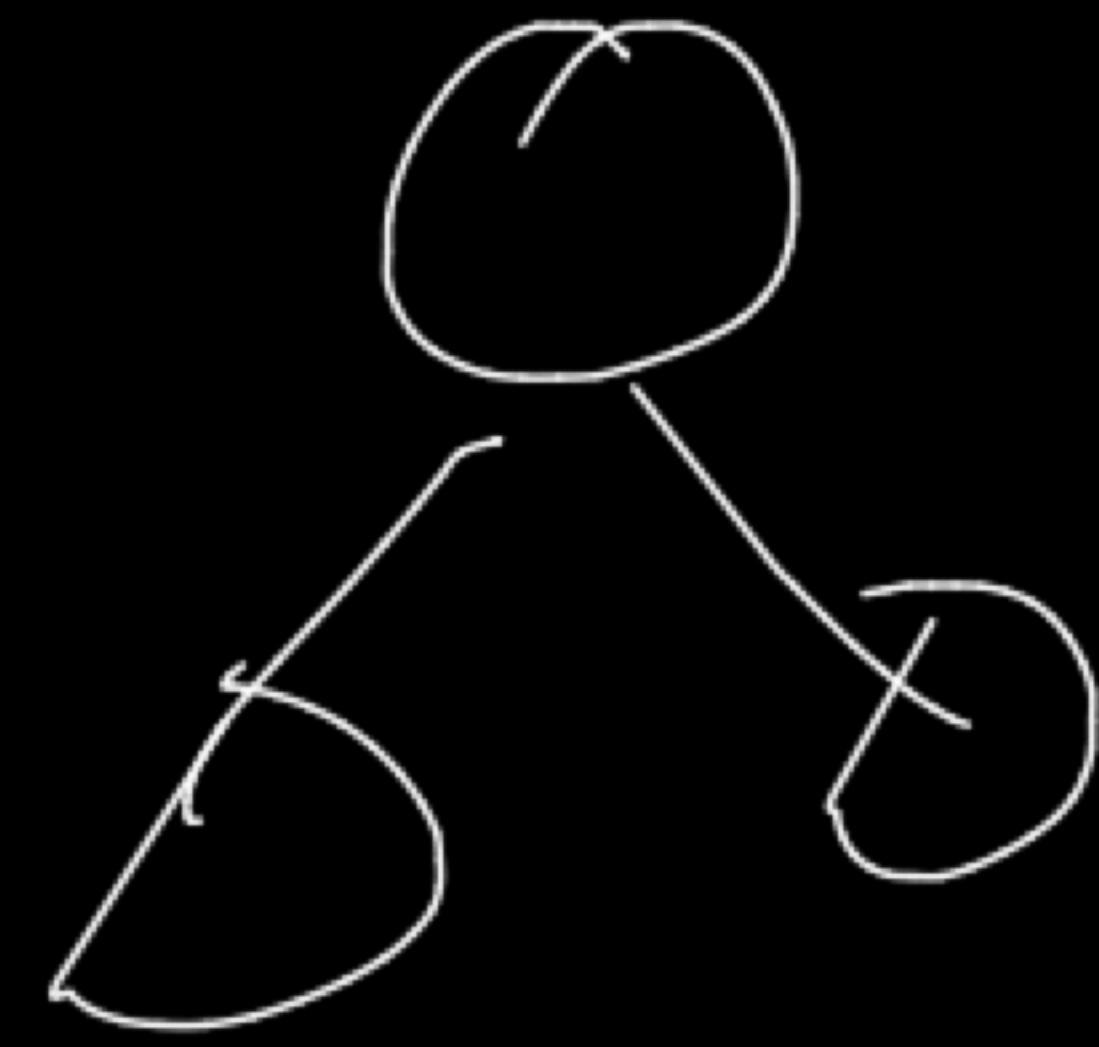
Model input
changes slightly

→ Model changes
a lot
Overfitting



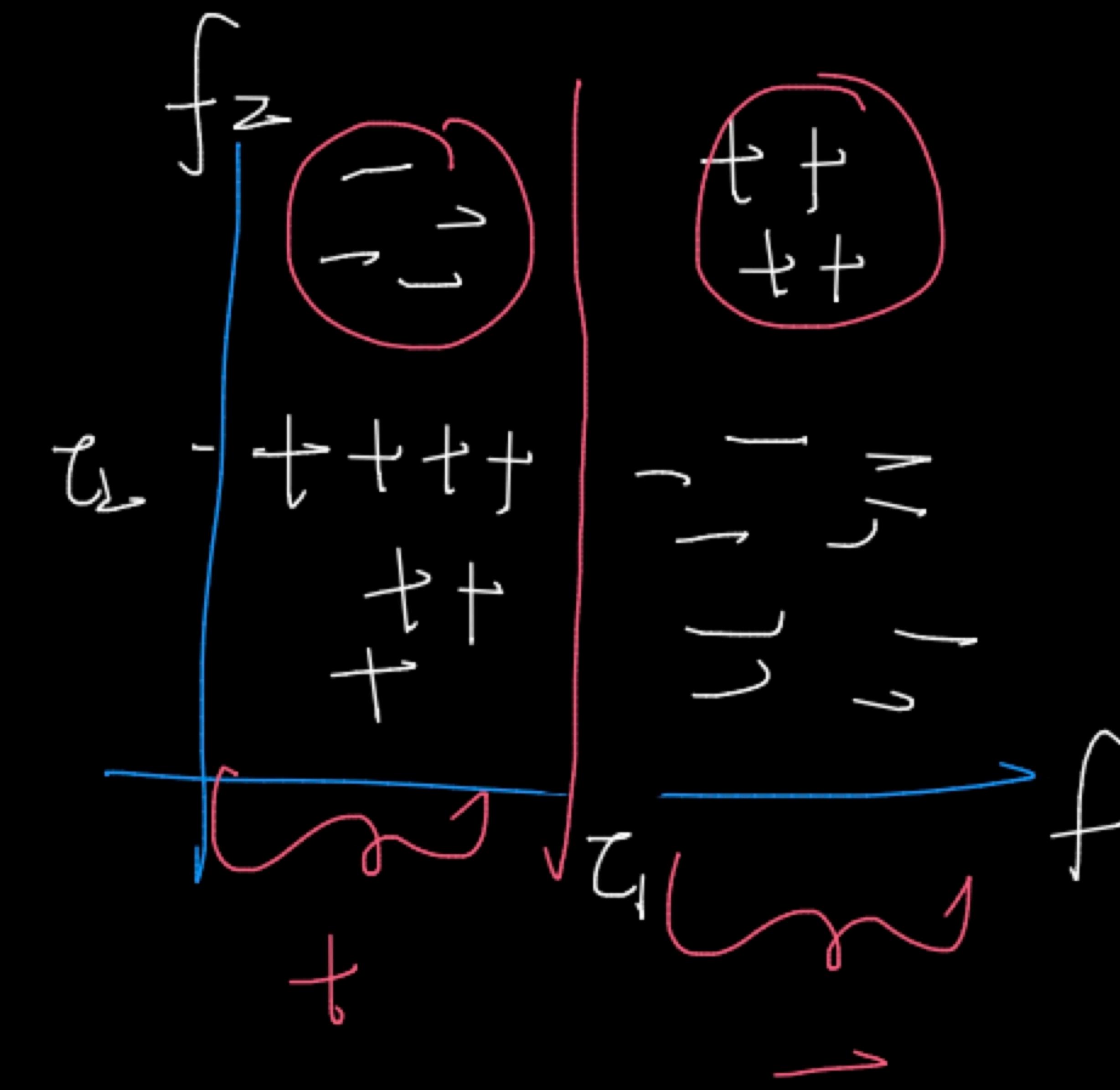
Underfitting

$$\text{Dep} | \bar{w} = 1$$



Simple & limited

Train - error is high



INTUITION



Simplex model



Complex Model



low bias

high bias \rightarrow Underfit
model



In statistics, a value of the parameter being estimated. An estimator of decision rule with zero bias is called unbiased. In statistics,

• "Is" is an *objective* property of an estimator. Bias is a distinct concept from **consistency**: consistent estimators converge in probability to the true value of the parameter, but may be biased or unbiased; see **bias versus consistency** for more.

All else being equal, an unbiased estimator is preferable to a biased estimator, although in practice, biased estimators (with generally small bias) are frequently used. When a biased estimator is used, bounds of the bias are calculated. A biased estimator may be used for various reasons: because an unbiased estimator does not exist without further assumptions about a population; because an estimator is difficult to compute (as in [unbiased estimation of standard deviation](#)); because a biased estimator may be unbiased with respect to different measures of [central tendency](#); because a biased estimator gives a lower value of some [loss function](#) (particularly [mean squared error](#)) compared with unbiased estimators (notably in [shrinkage estimators](#)); or because in some cases being unbiased is too strong a condition, and the only unbiased estimators are not useful.

Bias can also be measured with respect to the [median](#), rather than the mean (expected value), in which case one distinguishes *median-unbiased* from the usual *mean-unbiasedness* property. Mean-unbiasedness is not preserved under non-linear [transformations](#), though median-unbiasedness is (see § [Effect of transformations](#)); for example, the [sample variance](#) is a biased estimator for the population variance. These are all illustrated below.

Definition [edit]

Suppose we have a **statistical model**, parameterized by a real number θ , giving rise to a probability distribution for observed data, $P_\theta(x) = P(x | \theta)$. and a statistic $\hat{\theta}$ which serves as an **estimator** of θ based on any observed data x . That is, we assume that our data is part of this

$$\mathcal{O}\left(K \cdot \underbrace{\text{depth}}_{\uparrow}^{\max}\right) \rightarrow \text{runtime Compl}$$

$$\mathcal{O}\left(K * \underbrace{\text{time compl}}_{\text{for me DT}}\right) \rightarrow \text{Tauntive Compl}$$

an squared error

$$e^{-4\lambda} - 2e^{\lambda(1/e^2-3)} + e^{\lambda(1/e^4-1)}$$

is smaller; compare the unbiased estimator's MSE of

$$1 - e^{-4\lambda}.$$

The MSEs are functions of the true value λ . The bias of the maximum-likelihood estimator is:

$$e^{-2\lambda} - e^{\lambda(1/e^2-1)}.$$

Maximum of a discrete uniform distribution [edit]

Main article: Maximum of a discrete uniform distribution

The bias of maximum-likelihood estimators can be substantial. Consider a case where n tickets numbered from 1 through to n are placed in a box and one is selected at random, giving a value X . If n is unknown, then the maximum-likelihood estimator of n is X , even though the expectation of X given n is only $(n + 1)/2$; we can be certain only that n is at least X and is probably more. In this case, the natural unbiased estimator is $2X - 1$.

Median-unbiased estimators [edit]

The theory of median-unbiased estimators was revived by George W. Brown in 1947:^[6]