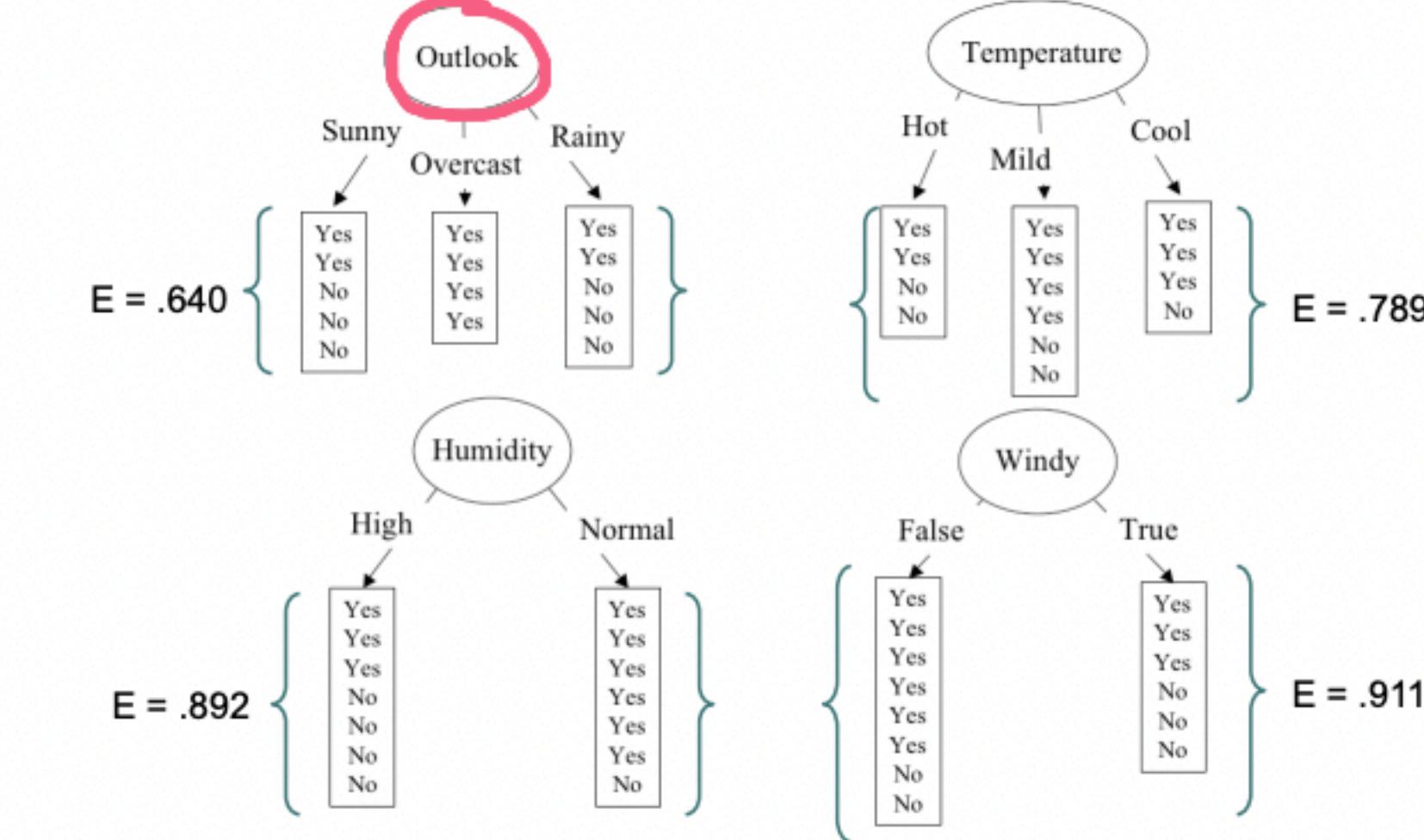


Topics:

- Questions tab
- chat for conversations
- Decision Trees (contd ..) → thoroughly ...
- Bagging → bootstrapped - aggregation
- Random forest ↗ FAANG
- Code

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Y



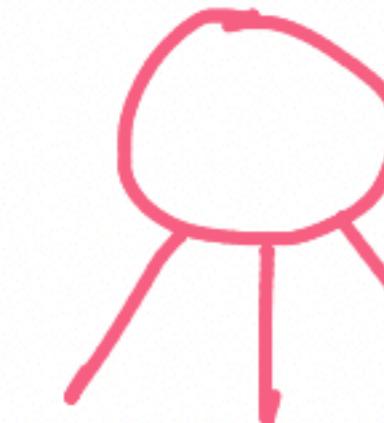
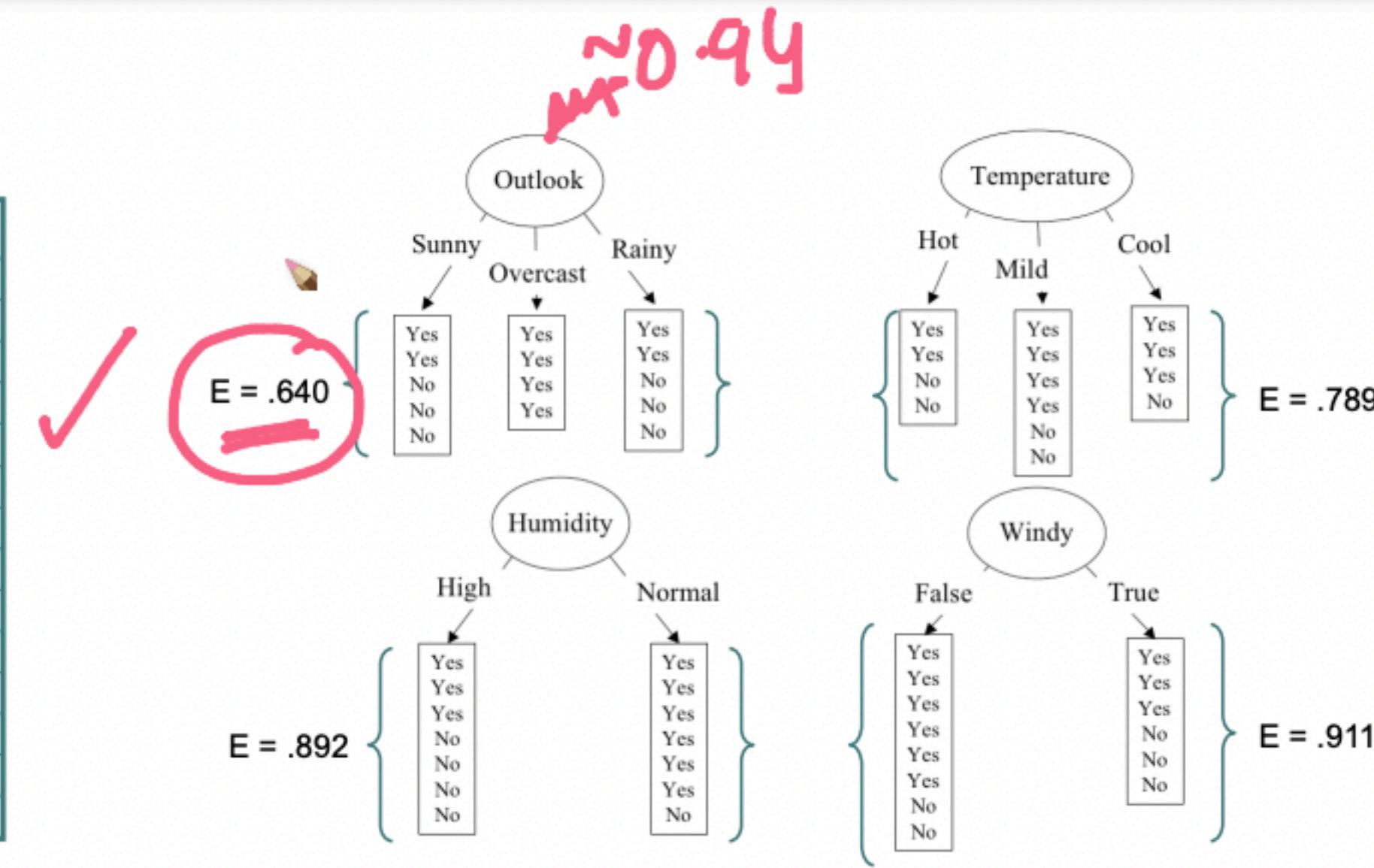
$\underline{IG(Y, \text{outlook})} = \frac{\text{weighted entropy @ child nodes}}{-\text{Entropy @ parent node}}$

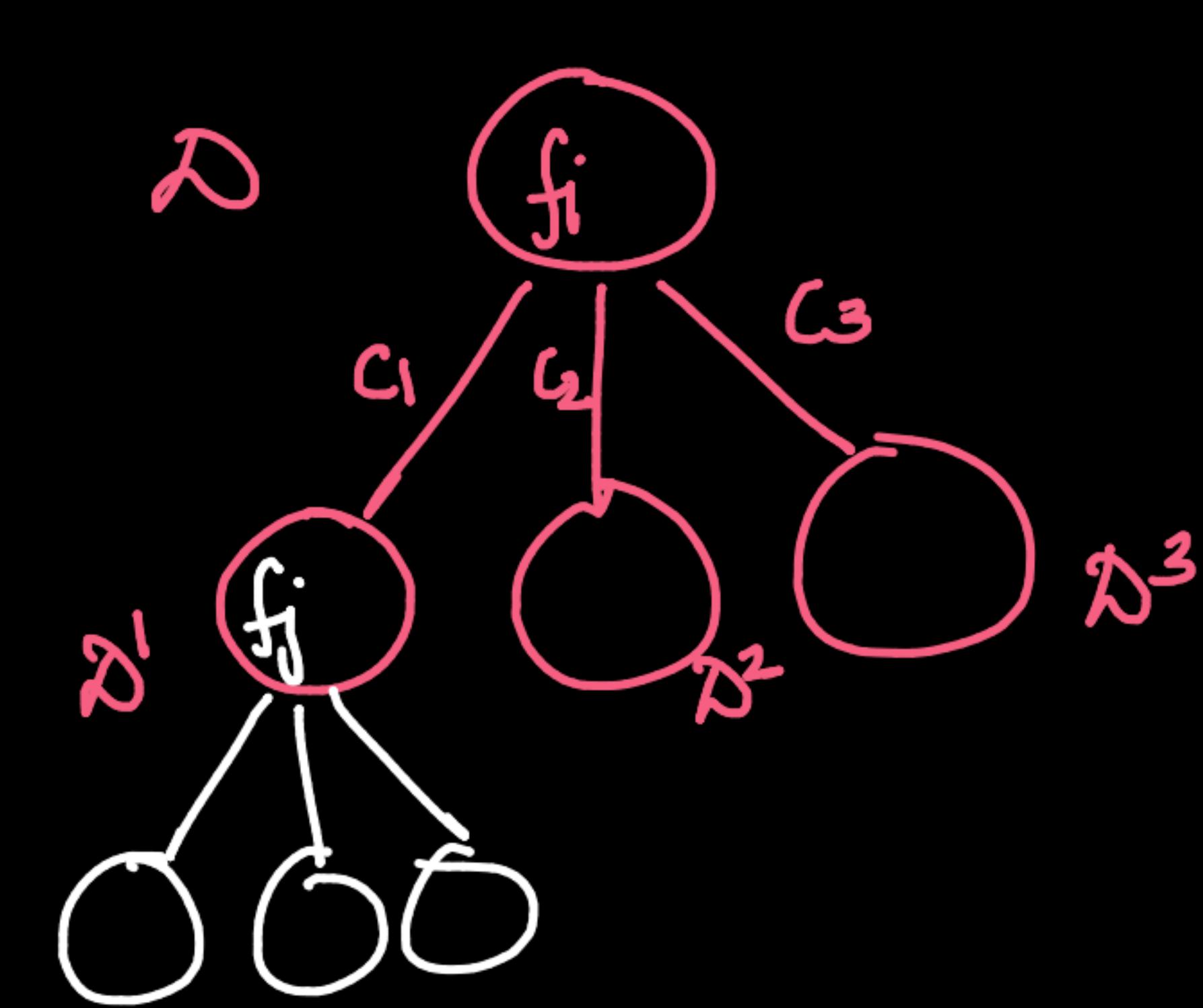


✓ ✓ ✓ ✓ yi

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

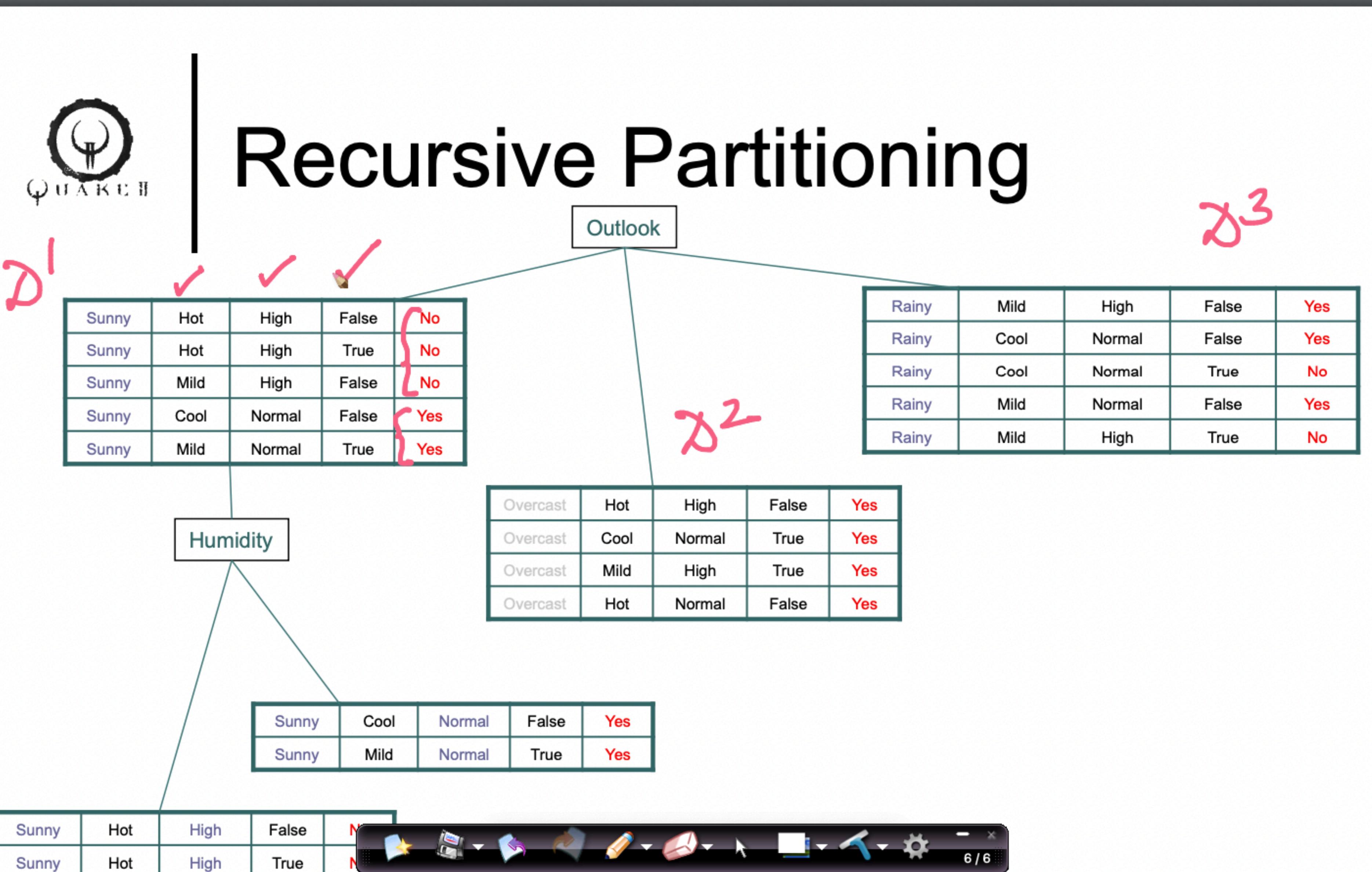
f₁ f₂ f₃ f₄

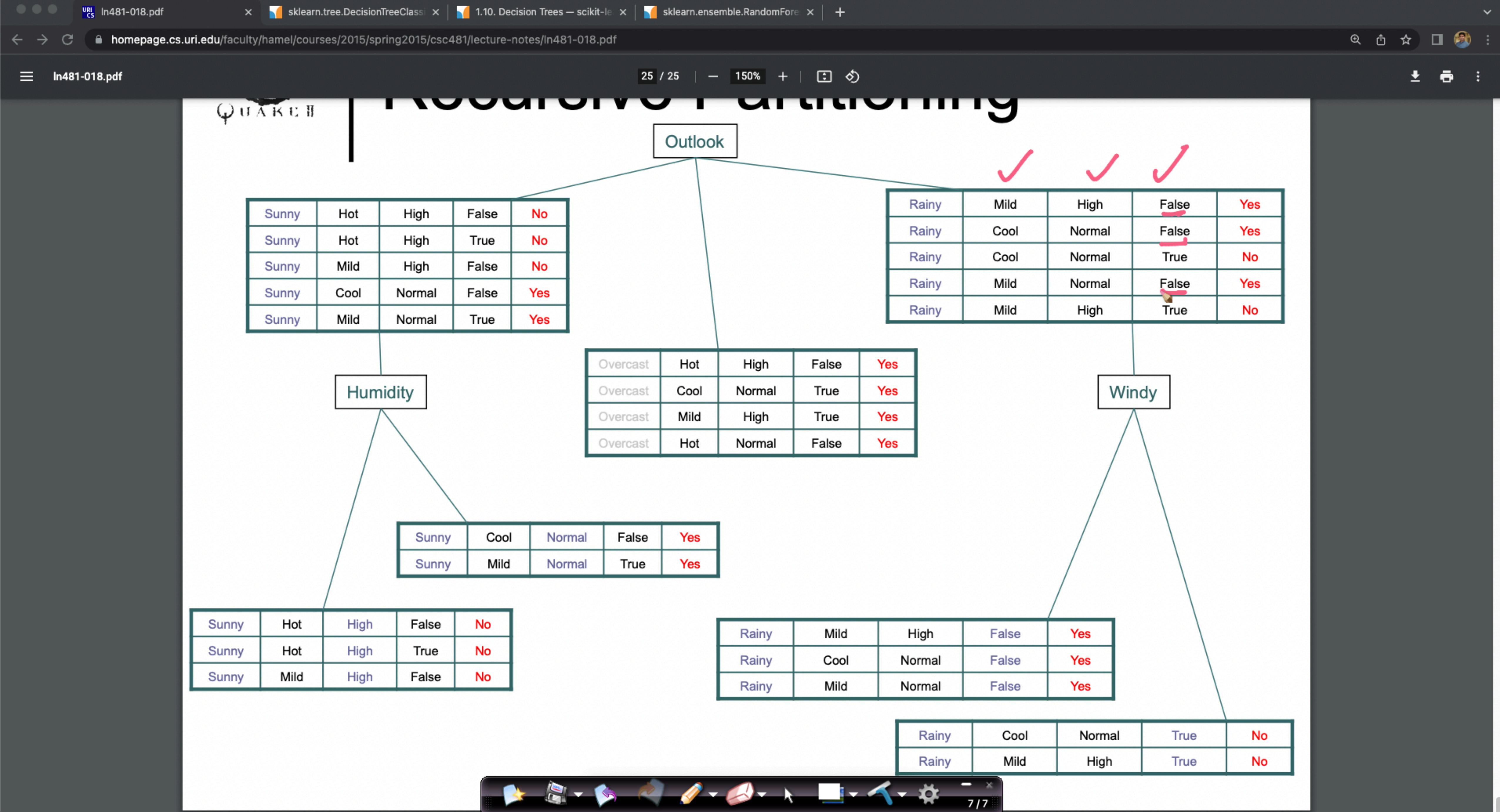


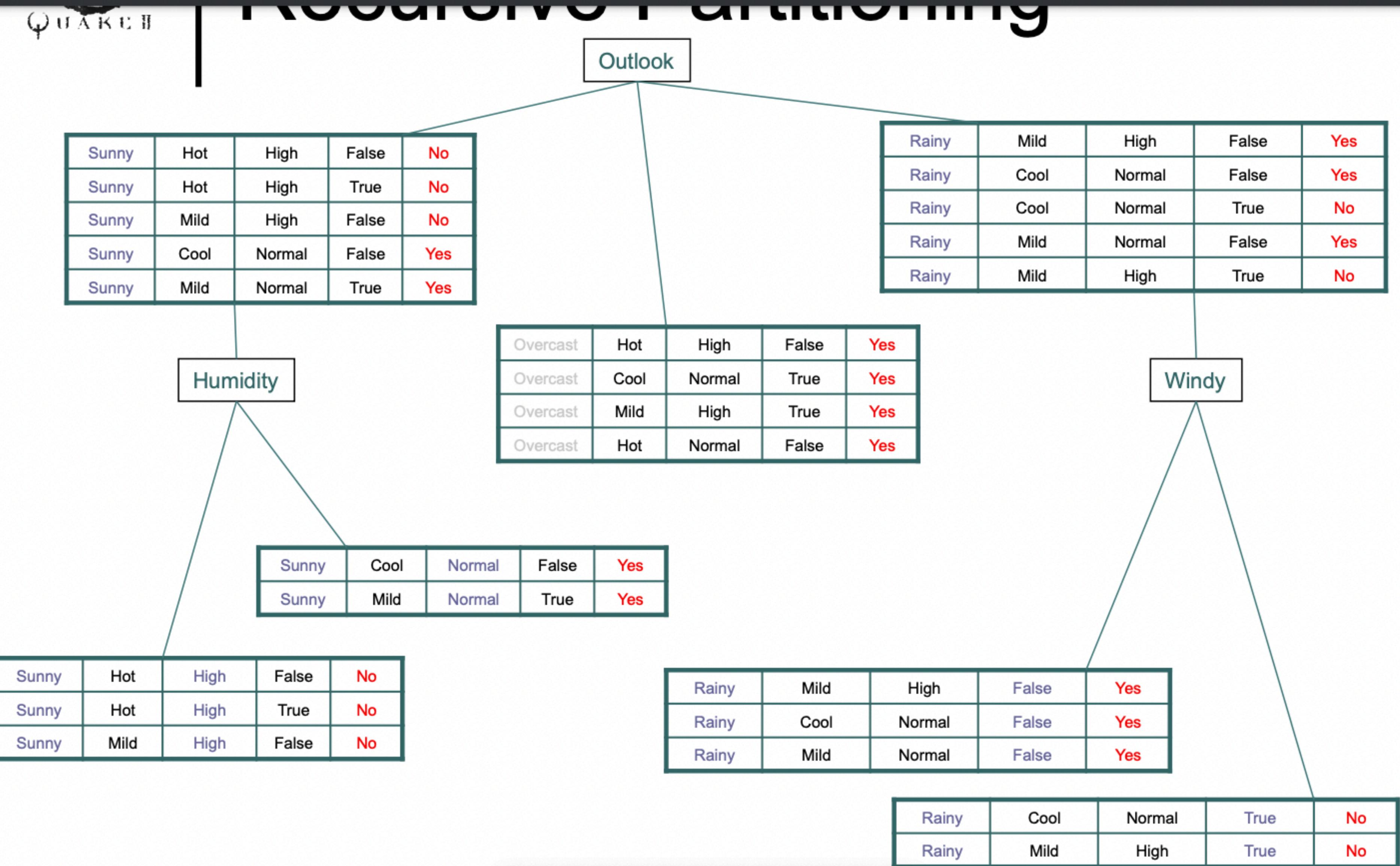


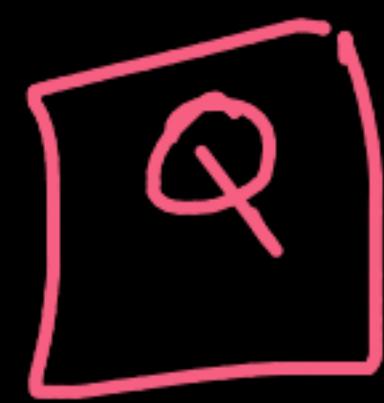
- push-nodes
=====
only $y_i = 1$
or
 $y_i = 0$
=====

* @ every node of the DT; we have to try
every feature & Split to pick the best
decision to split on









lots of features; d is high ; $\overbrace{d > 100}$

time to build the
trees \uparrow

RF
GBDT

disb · Computing
or multi-processing

$\overbrace{d = 1000}$
avoid DT

URI CS In481-018.pdf

sklearn.tree.DecisionTreeClassi x

1.10. Decision Trees – scikit-learn

sklearn.ensemble.RandomFore x

+

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

Go



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](#).

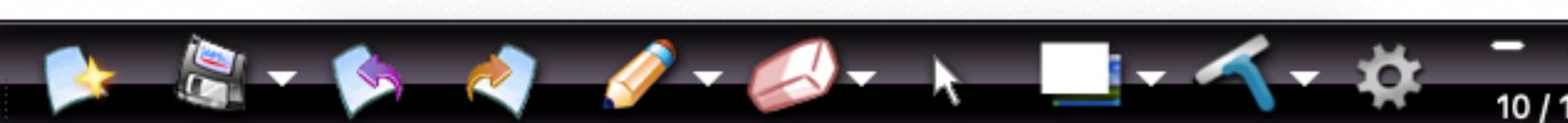
Parameters: `n_estimators : int, default=100`

The number of trees in the forest.

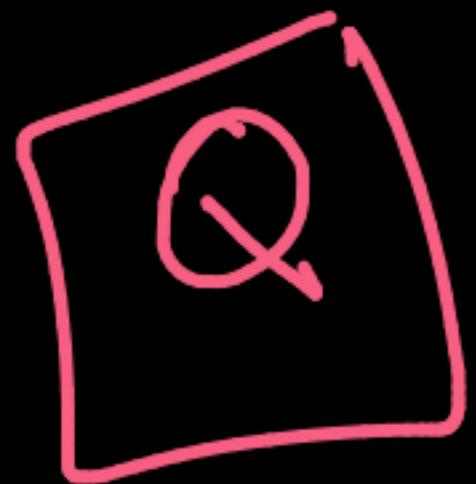
Changed in version 0.22: The default value of `n_estimators` changed from 10 to 100 in 0.22.

`criterion : {"gini", "entropy"}, default="gini"`

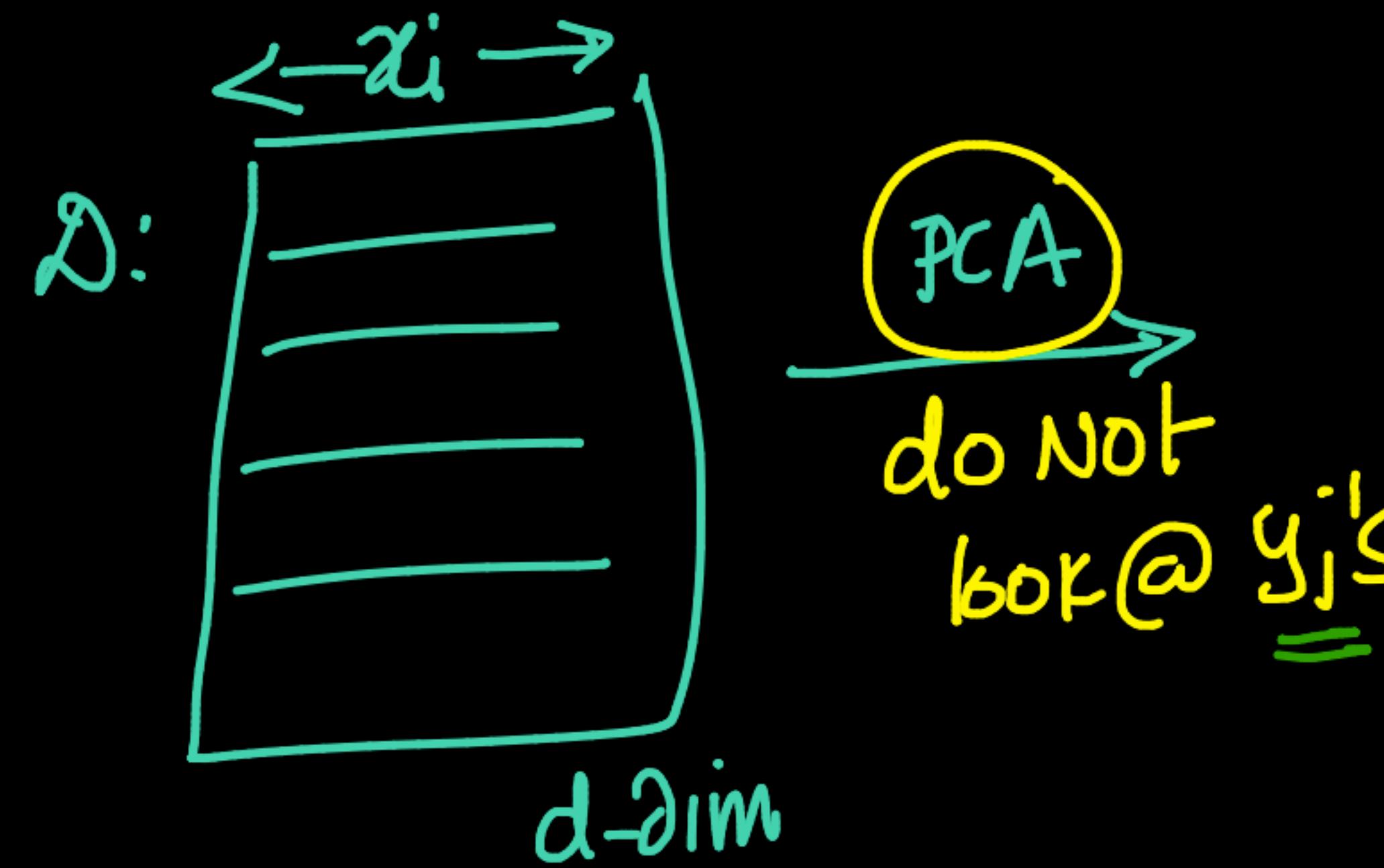
The function to measure the quality of a split. Supported criteria are “gini” for the Gini impurity and “entropy” for the information gain. Note: this parameter is tree-specific.



Toggle Menu

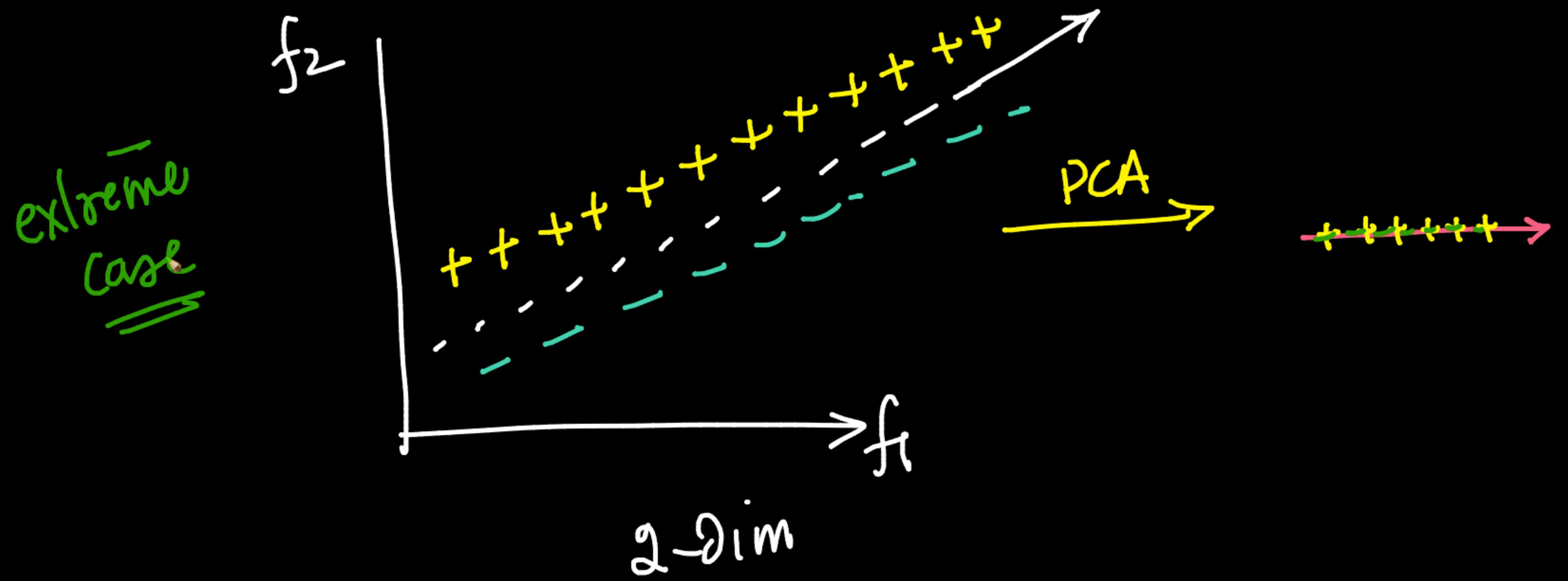


PCA : Dimensionality reduction →



$$d'\text{-dim} < d$$

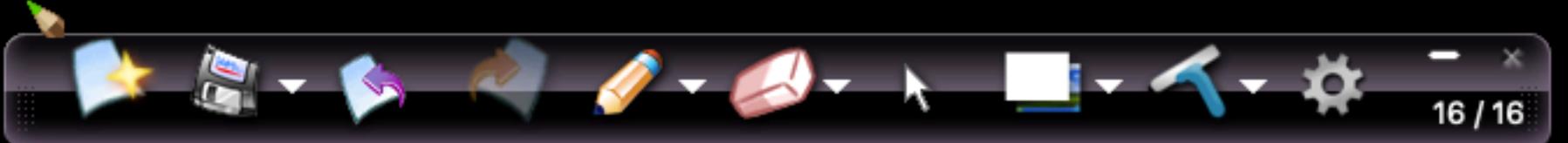
{here exist dim}
red..tech that
consider y_i 's also}

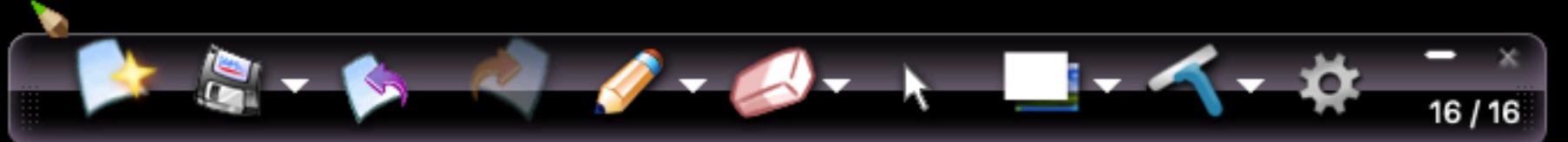


(d-dim) Try **PCA** Using top d' eigen-vcs
 $\rightarrow d' \text{ dim} \rightarrow D\bar{T}$ if it works ✓
→ if it does not work

$$x_{qj} \xrightarrow{\text{top eigen-vcs}} x_{qj}' \in \mathbb{R}^{d'}$$

Interpretability
is significantly
reduced...





~~IG~~: Info gain

$$GI(Y) = 1 - \sum_{i=1}^K (P(Y_i))^2$$

2-class: - $\underline{GI(Y)} = 1 - \left[(P(Y_i=1))^2 + (P(Y_i=0))^2 \right]$

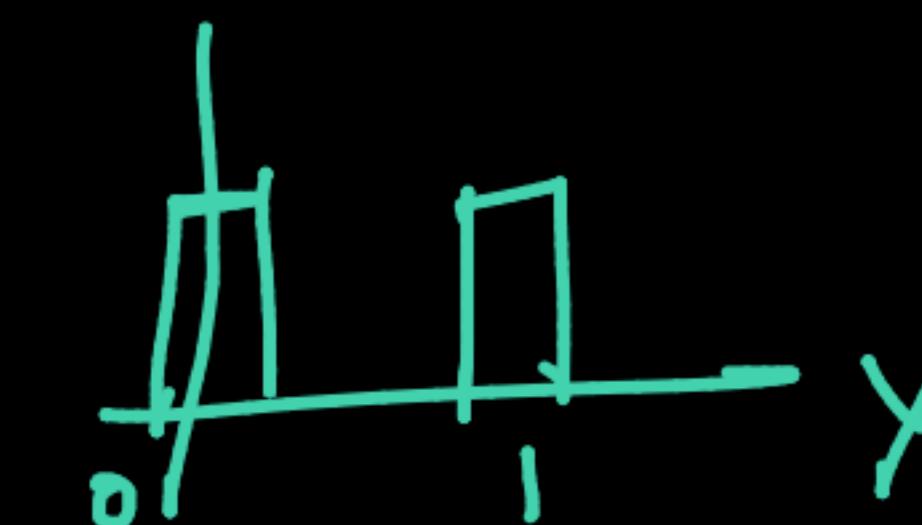
~~Gini Impurity~~
↓
Similar to
Entropy

Case I:

=

$$P(Y_+) = 0.5$$

$$P(Y-) = 0.5$$



$$\log_2 \frac{1}{\text{high}}$$

$$H(Y) = 1$$

$$GI(Y) = 1 - (0.25 + 0.25) = 0.5 \underset{\text{high}}{\underline{.5}}$$

Case 2:

$$P(Y_+) = 1$$

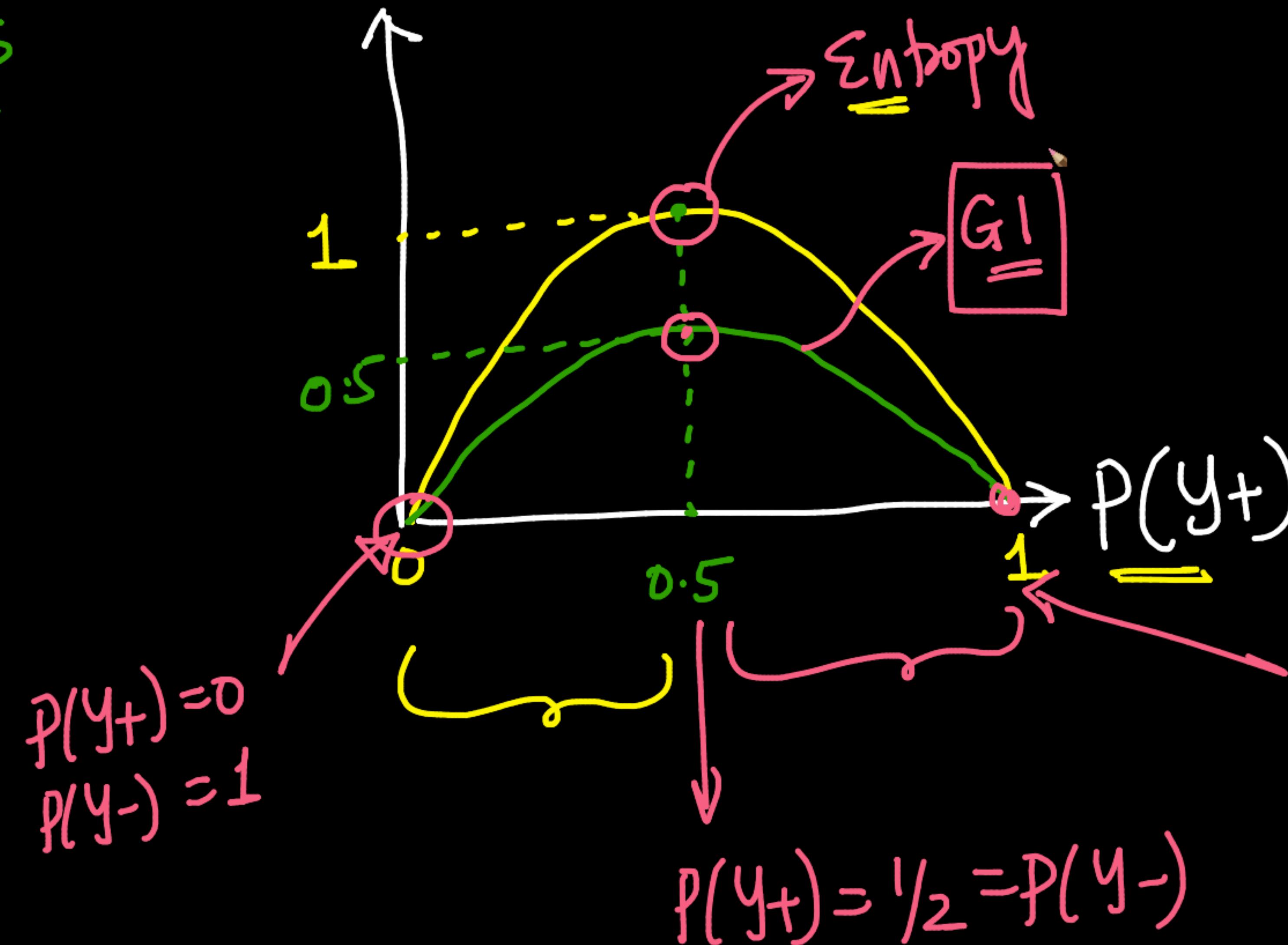
$$P(Y_-) = 0$$

$H(Y) = 0$: Min Entropy &
max purity

$$GI(Y) = 1 - (1+0) = 0$$

small

2-class

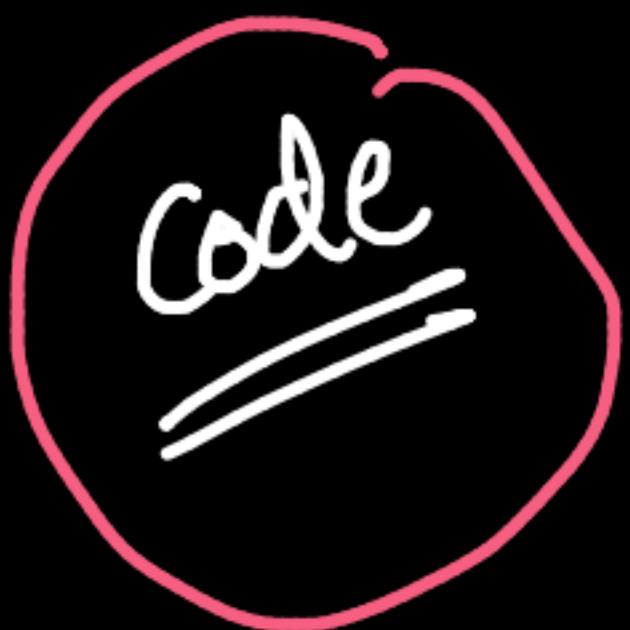


$$- \lg_2 \text{ in } H(Y)$$

$$- \hat{P}(Y_-) = 1 - \hat{P}(Y_+)$$

$$\begin{aligned} P(Y_+) &= 1 \\ P(Y_-) &= 0 \end{aligned}$$

$$P(Y_+) = 1/2 = P(Y_-)$$



$H(Y)$

VS

$G(Y)$

$k = \# \text{classes}$

$$-\sum_{i=1}^k p(y_i) \log(p(y_i))$$

$$-\sum_{i=1}^k (p(y_i))^2$$

sklearn

* log is computationally more expensive
than squaring

URI CS In481-018.pdf x sklearn.tree.DecisionTreeClassi x 1.10. Decision Trees – scikit-learn x sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html

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.tree.DecisionTreeClassifier

Examples using sklearn.tree.DecisionTreeClassifier

sklearn.tree.DecisionTreeClassifier

```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, class_weight=None, ccp_alpha=0.0)
```

[source]

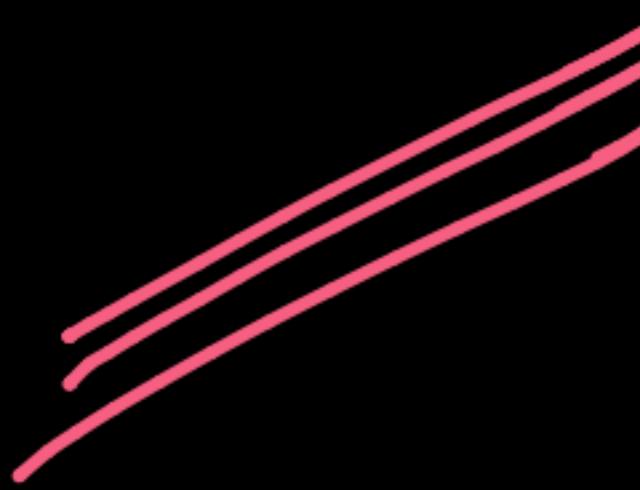
A decision tree classifier.

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

Parameters:

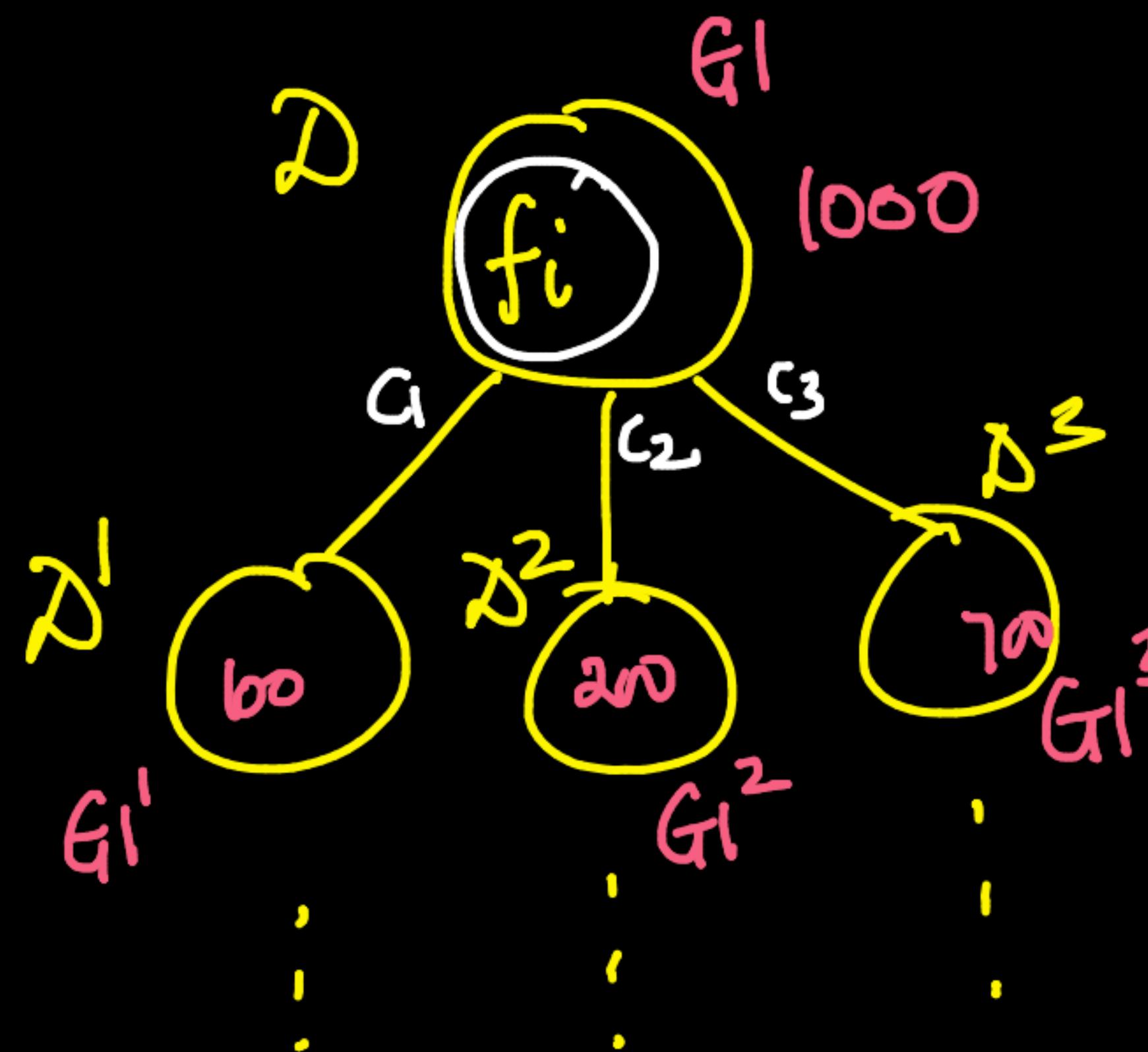
- criterion : {"gini", "entropy"}, default="gini"**
The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.
- splitter : {"best", "random"}, default="best"**
The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.
- max_depth : int, default=None**
The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.
- min_samples_split : int or float, default=2**
The minimum number of samples required to split an internal node:

Toggle Menu



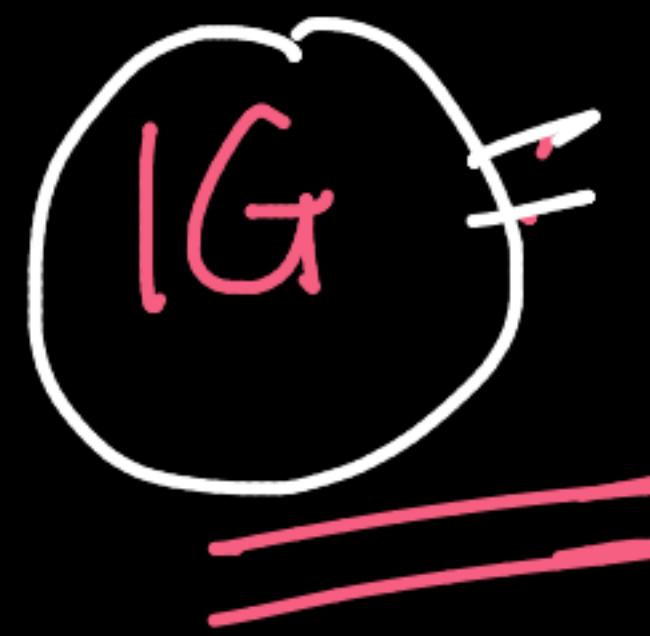
multi-class - classification:

1, 2, 3



Entropy of $G1$
g $\Delta H \approx \Delta G$

→ Split on features



Weighted $\underline{\underline{G \cdot I}}$ of the child nodes
- $G \cdot I$ of the parent node
 \hookrightarrow Entropy ✓



Prev Up Next

scikit-learn 1.0.2

[Other versions](#)Please [cite us](#) if you use the software.[sklearn.tree.DecisionTreeClassifier](#)

Examples using

[sklearn.tree.DecisionTreeClassi](#)

sklearn.tree.DecisionTreeClassifier

```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None, min_samples_split=2,  
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None,  
min_impurity_decrease=0.0, class_weight=None, ccp_alpha=0.0) [source]
```

A decision tree classifier.

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

Parameters:

criterion : {"gini", "entropy"}, default="gini"

The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

splitter : {"best", "random"}, default="best"

The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.

max_depth : int, default=None

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.

min_samples_split : int or float, default=2

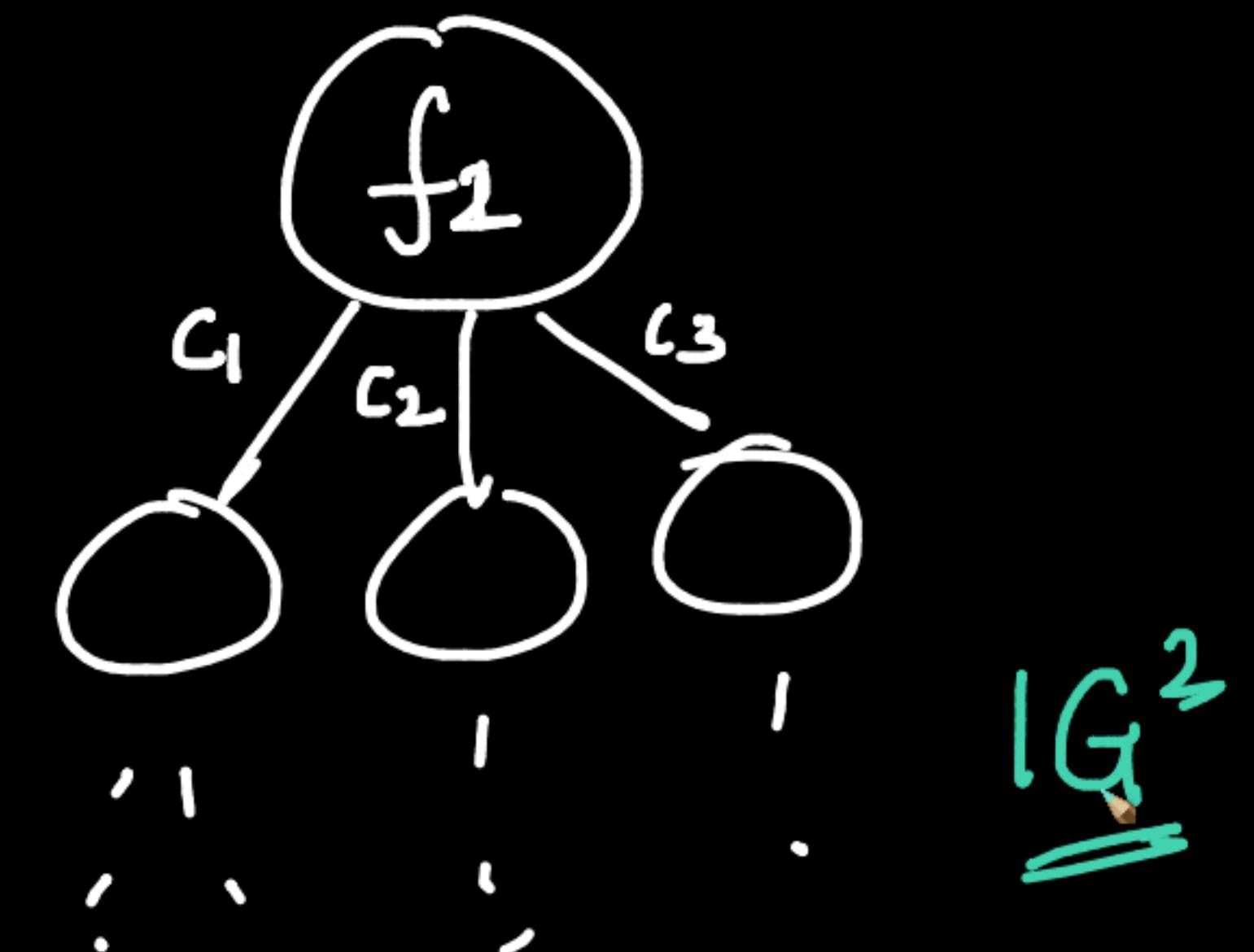
The minimum number of samples required to split an internal node:



Numerical features

f_2 : Categorical feature: c_1, c_2, c_3

↓
Splitting criterion



numerical

	f_2	$\downarrow f_1$	y_i
C_1	2.2	1	1
C_2	2.6	0	0
C_1	3.5	0	0
C_2	3.8	0	0
C_3	4.6	1	1
C_1	5.3	0	0

f_i : numerical $\rightarrow n$ -ways

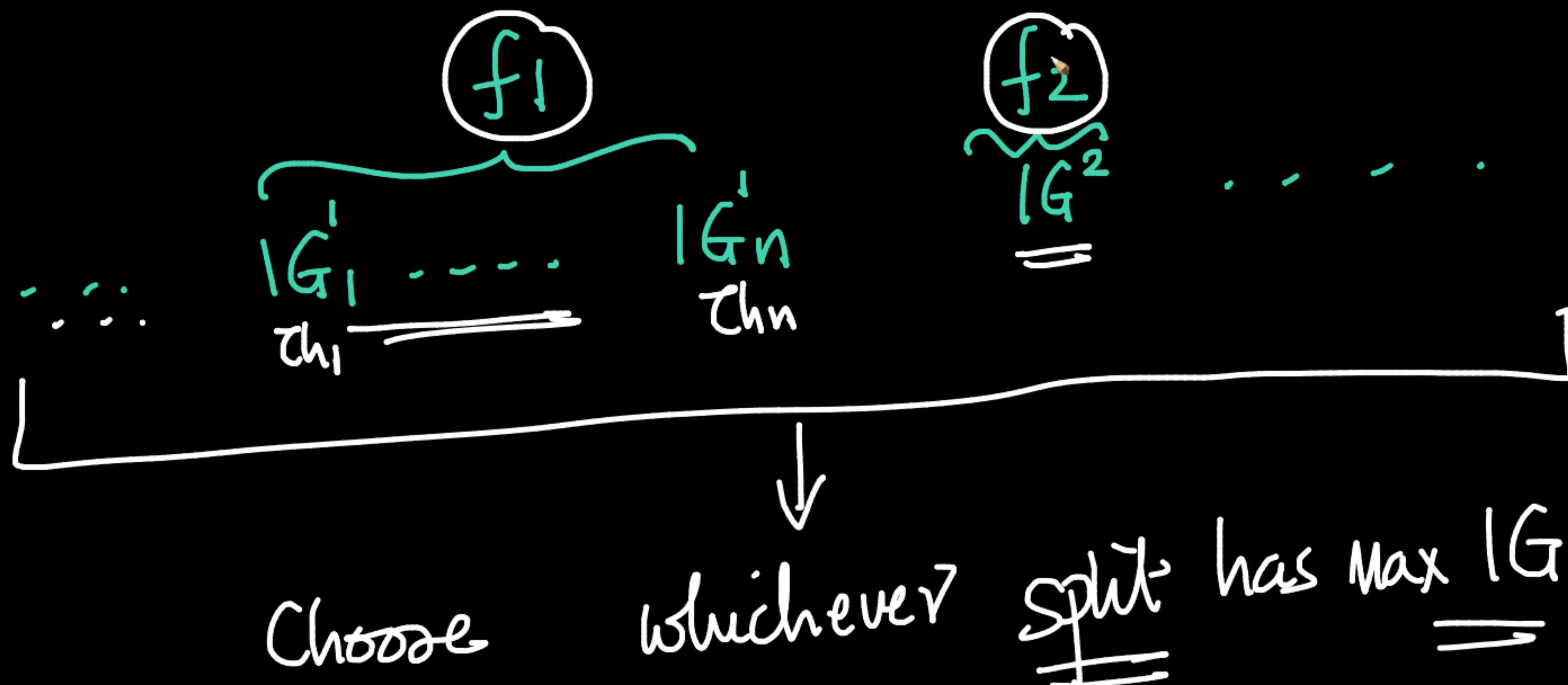
$$f_1 \leq \tau_h$$

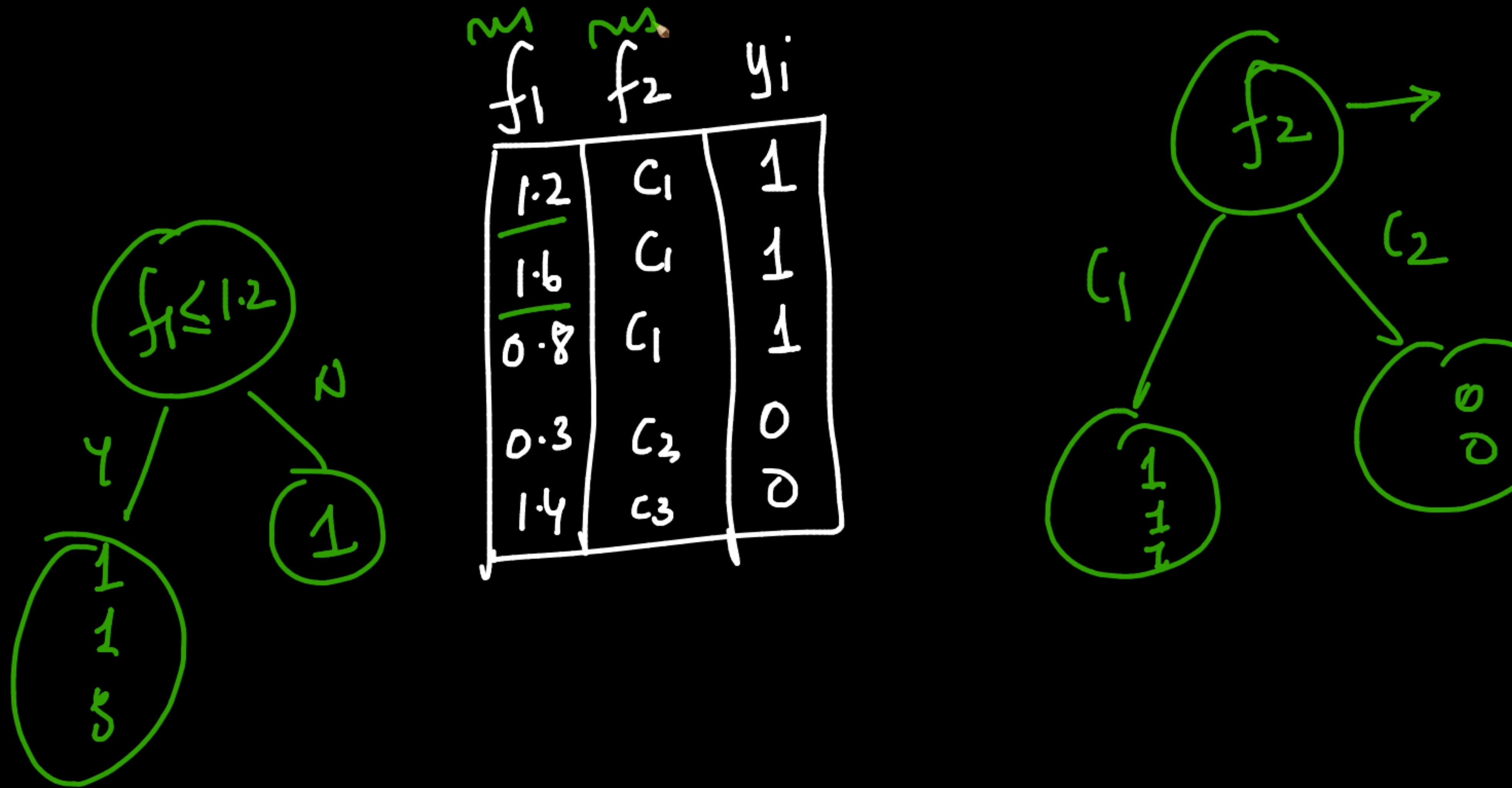
① $\tau_{h_1} = 2.2 \rightarrow |G_1|$

② $\tau_{h_2} = 2.6 \rightarrow |G_2|$

③ $\tau_{h_n} = 5.3 \rightarrow |G_n|$

⋮





$$|G|_0 = |G|^2 = \text{minimal}$$



Pick any of these equal $|G|$ splits.

f_1	y
2-3	
2-3	
:	
!	

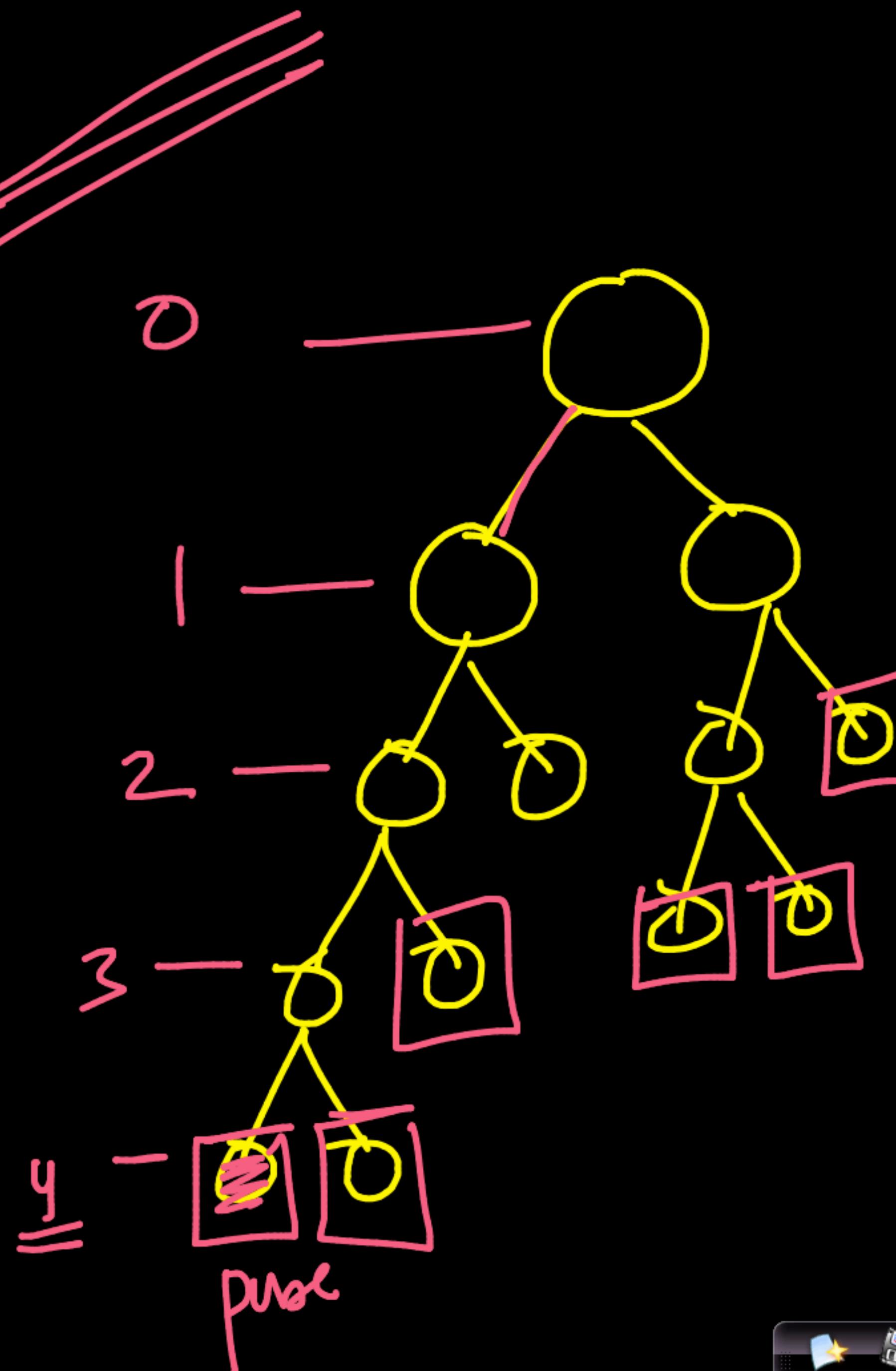
↑ rules

bin

it Carefully

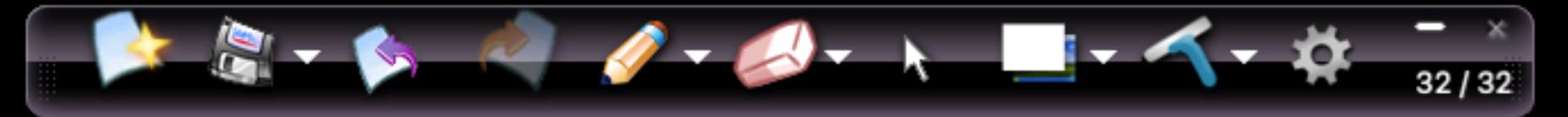
Simplest Q_1, Q_2, Q_3, Q_4

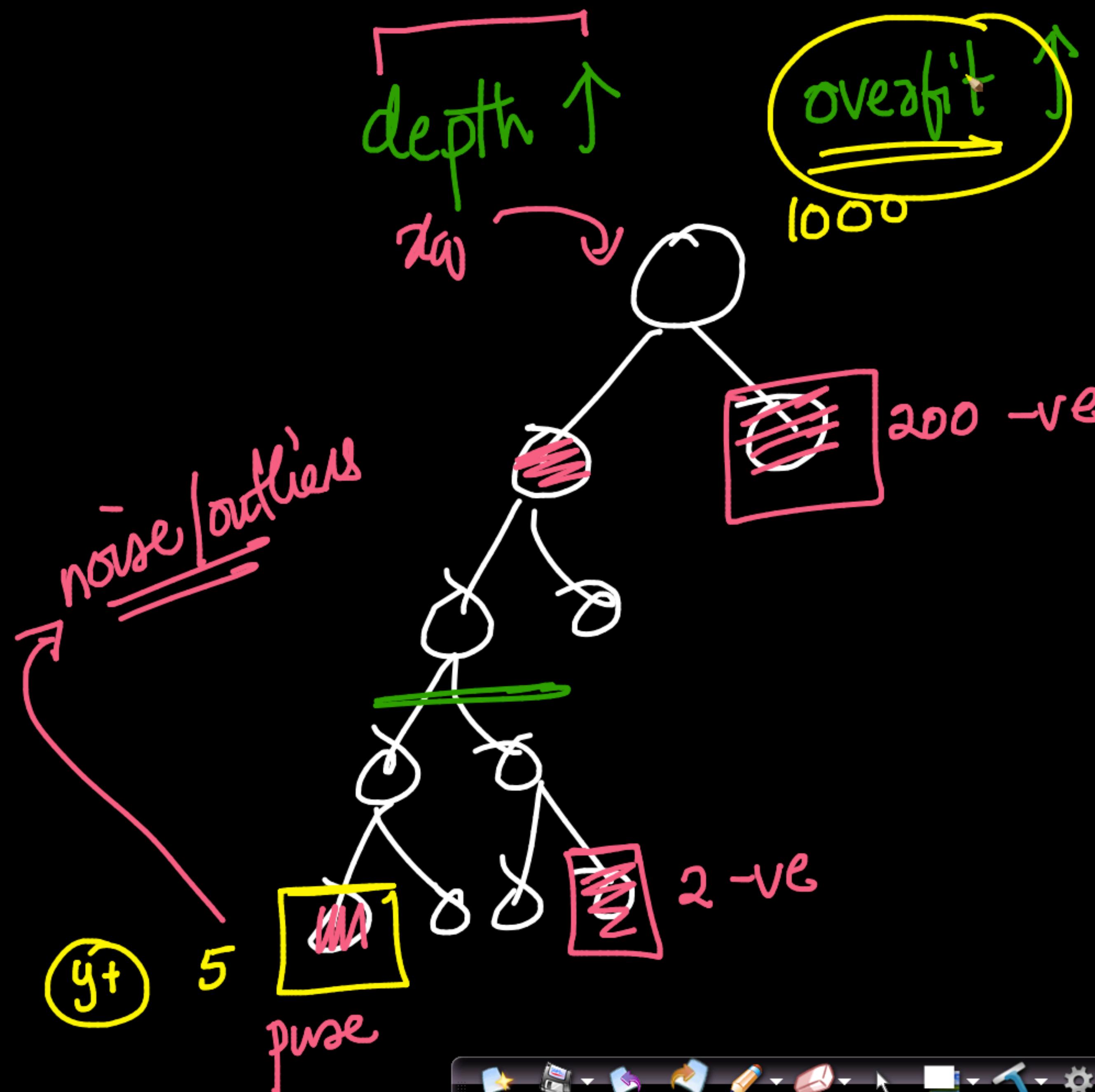
more computationally efficient



overfit vs underfit
↓
many splits

depth of a tree = 4
≠





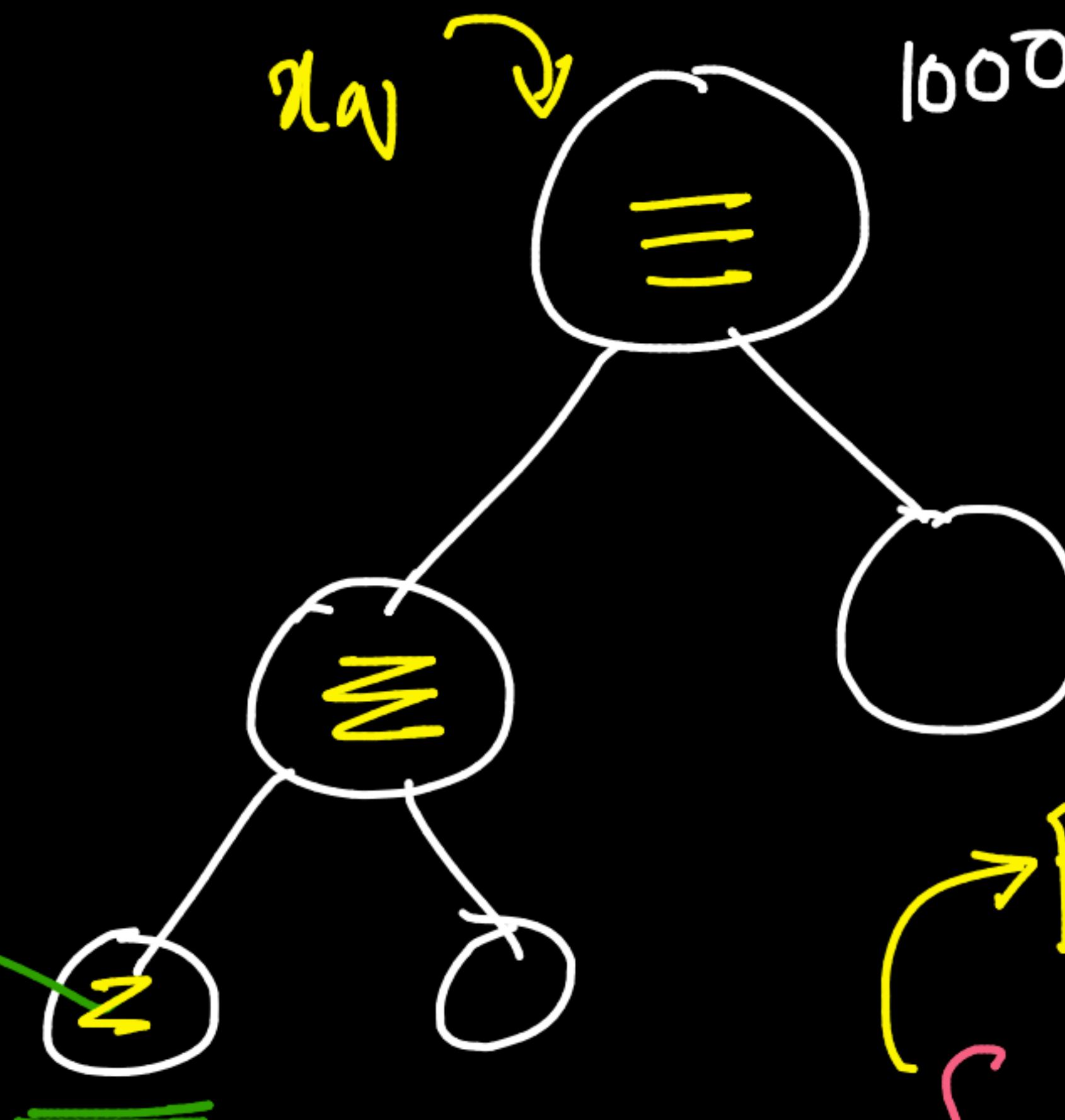
Shallow-trees:-

y_{av}  1000

depth = 2

not pure

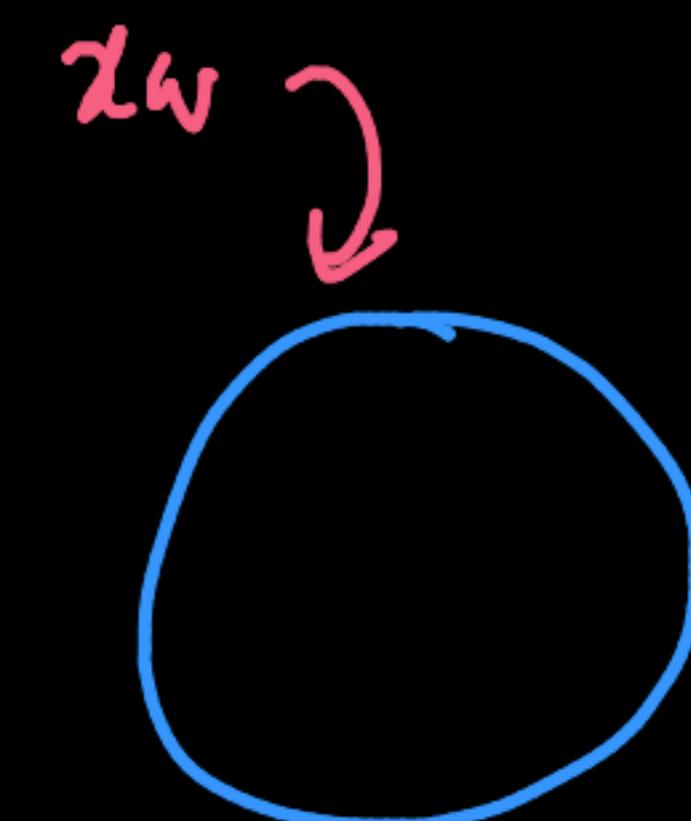
{ 500 y^+
200 y^-



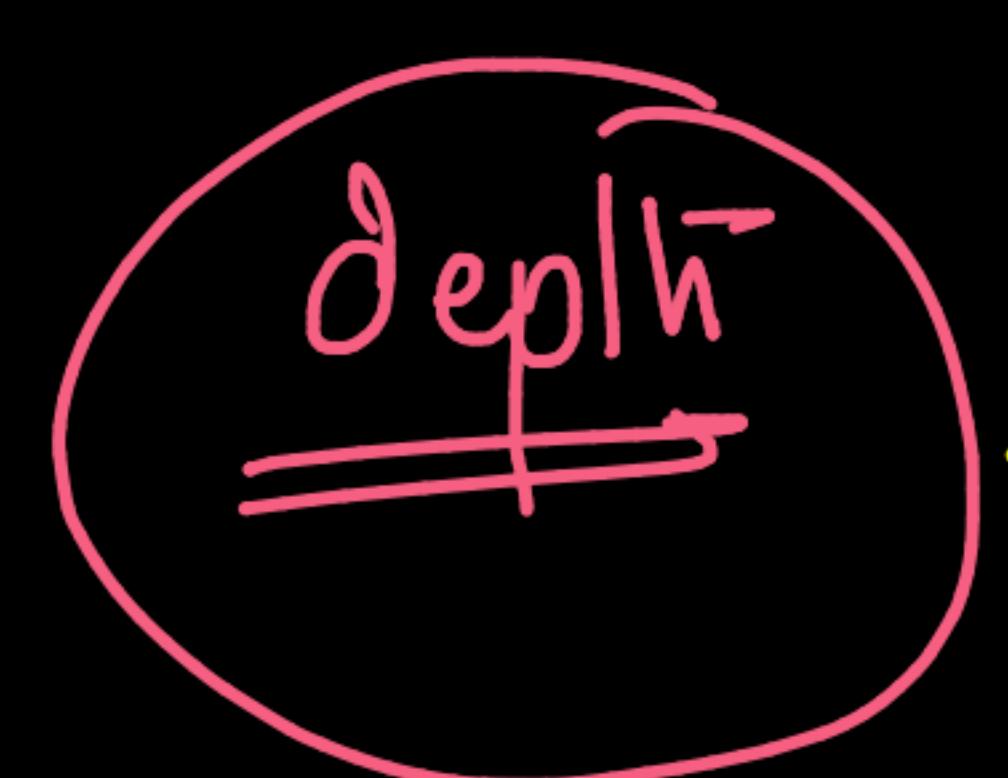
probabilistic class labels
 $y_{av} = +ve \text{ or } -ve$

$$\left\{ P(y_{av} = +ve) = \frac{500}{700} \right.$$

depth = 0



$$\omega = \frac{y_t - \underline{y}}{\bar{y} - \underline{y}} : \frac{55\%}{45\%}$$



Split atleast 20 data points
Split if IG is above some value

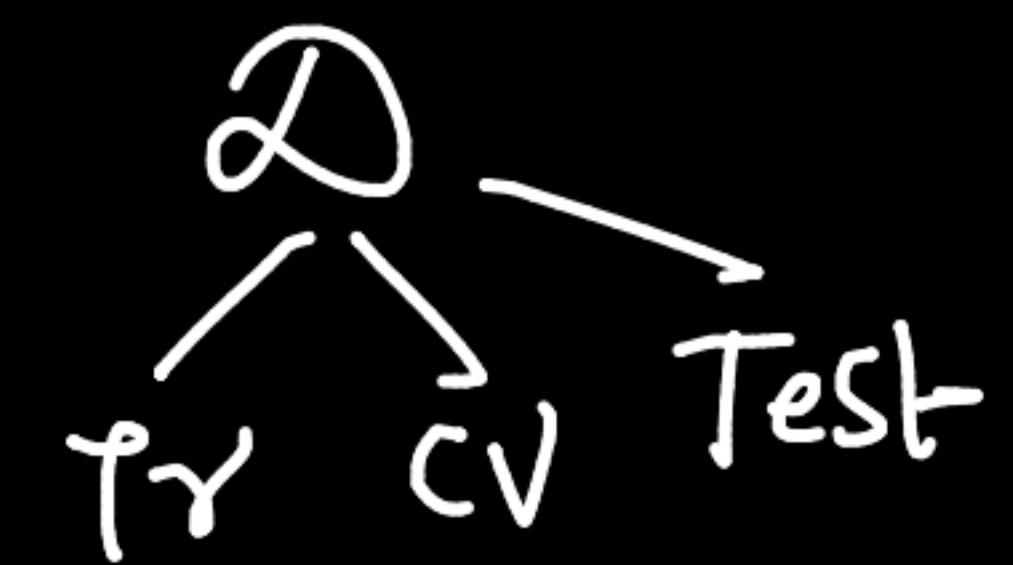
Depth : hyperparam

$G_2 \rightarrow CV\ error$

$\rightarrow 3 \rightarrow \vdots$

$\rightarrow 4 \vdots$

$\rightarrow 5 \vdots$



depth = 2

pure nodes → Great

URI CS In481-018.pdf x sklearn.tree.DecisionTreeClassi x 1.10. Decision Trees – scikit-learn x sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html

scikit learn

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

```
class sklearn.tree.DecisionTreeClassifier(*, criterion="gini", splitter="best", max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, class_weight=None, ccp_alpha=0.0)
```

[source]

A decision tree classifier.

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

Parameters: **criterion : {"gini", "entropy"}, default="gini"**

The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

splitter : {"best", "random"}, default="best"

The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.

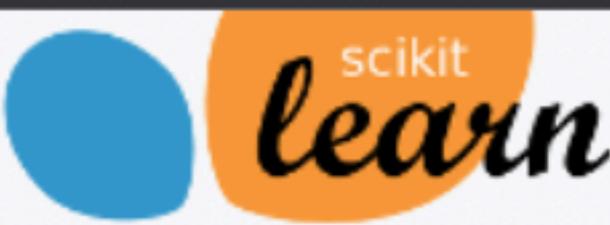
max_depth : int, default=None

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.

min_samples_split : int or float, default=2

The minimum number of samples required to split an internal node:

Extremely Random Forest



Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.tree.DecisionTreeClassifier

Examples using

sklearn.tree.DecisionTreeClassi



```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None,  
min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None,  
random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, class_weight=None,  
ccp_alpha=0.0)
```

[\[source\]](#)

A decision tree classifier.

Read more in the [User Guide](#).**Parameters:** **criterion : {"gini", "entropy"}, default="gini"**

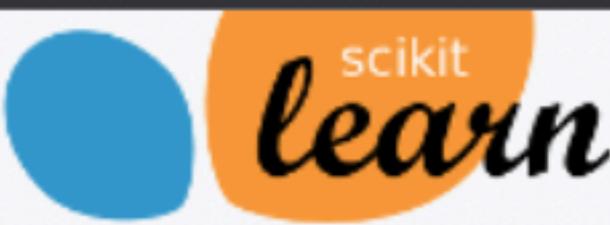
The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

splitter : {"best", "random"}, default="best"

The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.

max_depth : int, default=None

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.



Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.tree.DecisionTreeClassifier

Examples using

sklearn.tree.DecisionTreeClassi



```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None,  
min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None,  
random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, class_weight=None,  
ccp_alpha=0.0)
```

[\[source\]](#)

A decision tree classifier.

Read more in the [User Guide](#).**Parameters:** **criterion : {"gini", "entropy"}, default="gini"**

The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

splitter : {"best", "random"}, default="best"

The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.

max_depth : int, default=None

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.



Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.tree.DecisionTreeClassifier

Examples using

sklearn.tree.DecisionTreeClassi

```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None,  
min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None,  
random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, class_weight=None,  
ccp_alpha=0.0)
```

[\[source\]](#)

A decision tree classifier.

Read more in the [User Guide](#).**Parameters:** **criterion : {"gini", "entropy"}, default="gini"**

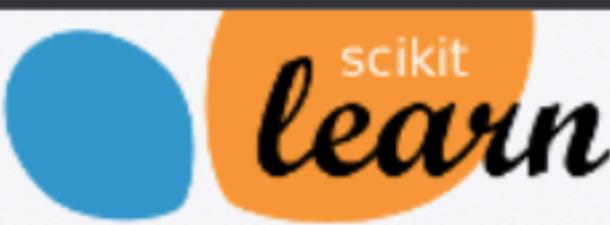
The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

splitter : {"best", "random"}, default="best"

The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.

max_depth : int, default=None

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.



Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.tree.DecisionTreeClassifier

Examples using

sklearn.tree.DecisionTreeClassifi



Toggle Menu

```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None,  
min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None,  
random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, class_weight=None,  
ccp_alpha=0.0)
```

[\[source\]](#)

A decision tree classifier.

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

Parameters: **criterion : {"gini", "entropy"}, default="gini"**

The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

splitter : {"best", "random"}, default="best"

The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.

max_depth : int, default=None

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.

URI CS In481-018.pdf x sklearn.tree.DecisionTreeClassi x 1.10. Decision Trees – scikit-learn x sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html

scikit learn Install User Guide API Examples Community More ▾

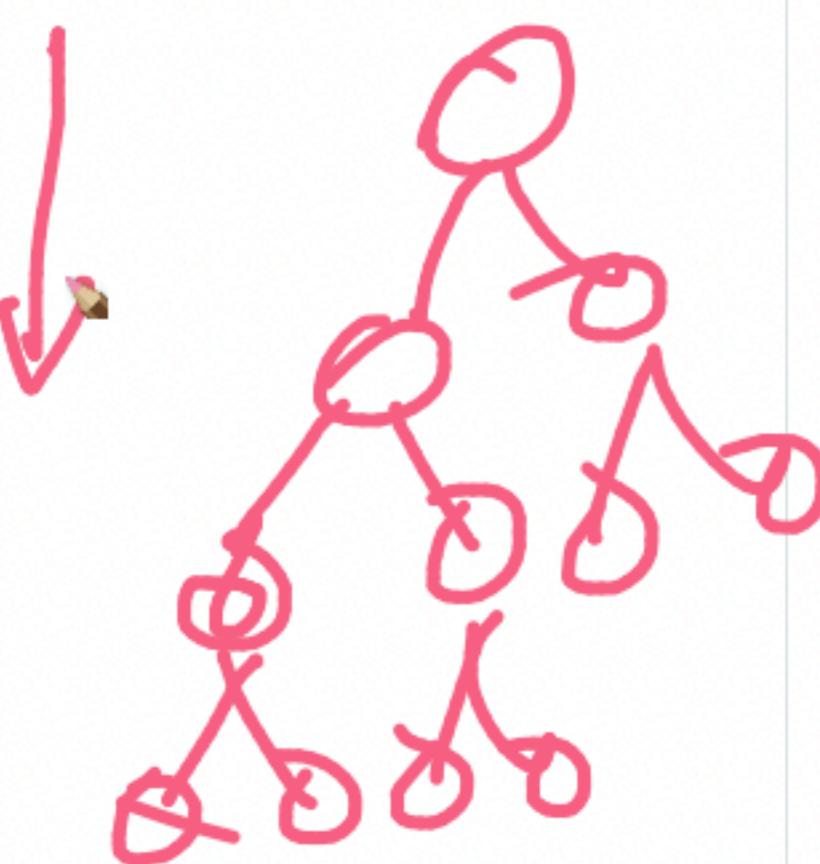
Prev Up Next

scikit-learn 1.0.2 Other versions

Please cite us if you use the software.

sklearn.tree.DecisionTreeClassifier

Examples using sklearn.tree.DecisionTreeClassifier



`class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, class_weight=None, ccp_alpha=0.0)`

[source]

A decision tree classifier.

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

Parameters:

- criterion : {"gini", "entropy"}, default="gini"**
The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.
- splitter : {"best", "random"}, default="best"**
The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.
- max_depth : int, default=None**
The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.

URI CS In481-018.pdf x sklearn.tree.DecisionTreeClassi x 1.10. Decision Trees – scikit-learn x sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html

scikit learn

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.tree.DecisionTreeClassifier

Examples using sklearn.tree.DecisionTreeClassifier

f1 f2 ... fJ

n_samples) are the minimum number of samples for each node.

Changed in version 0.18: Added float values for fractions.

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 : int, float or {"auto", "sqrt", "log2"}, default=None

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 int(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).
- 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.

45 / 45

URI CS In481-018.pdf x sklearn.tree.DecisionTreeClassi x 1.10. Decision Trees – scikit-learn x sklearn.ensemble.RandomFore x +

scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.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.tree.DecisionTreeClassifier

Examples using sklearn.tree.DecisionTreeClassifier

sklearn.tree.DecisionTreeClassifier

```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, class_weight=None, ccp_alpha=0.0)
```

[source]

A decision tree classifier.

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

Parameters:

criterion : {"gini", "entropy"}, default="gini"
The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

splitter : {"best", "random"}, default="best"
The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.

max_depth : int, default=None

then nodes are expanded until all leaves are

URI CS In481-018.pdf x | sklearn.tree.DecisionTreeClassi x | 1.10. Decision Trees – scikit-learn 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`

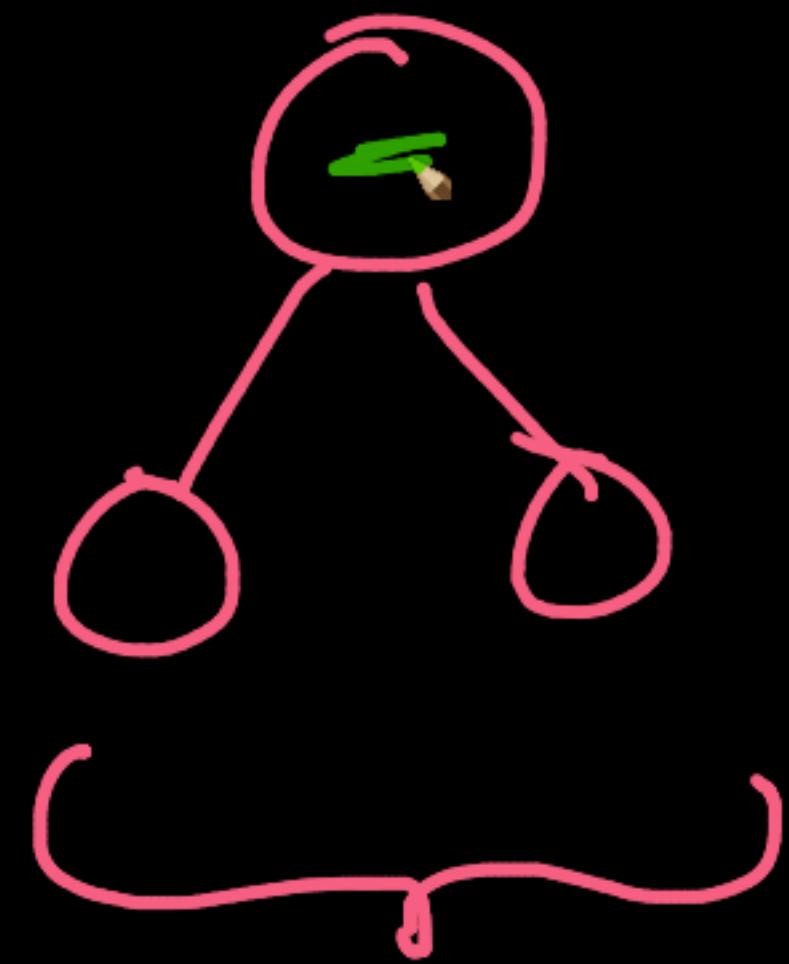
The number of trees in the forest.

Toggle Menu

decision

stump:

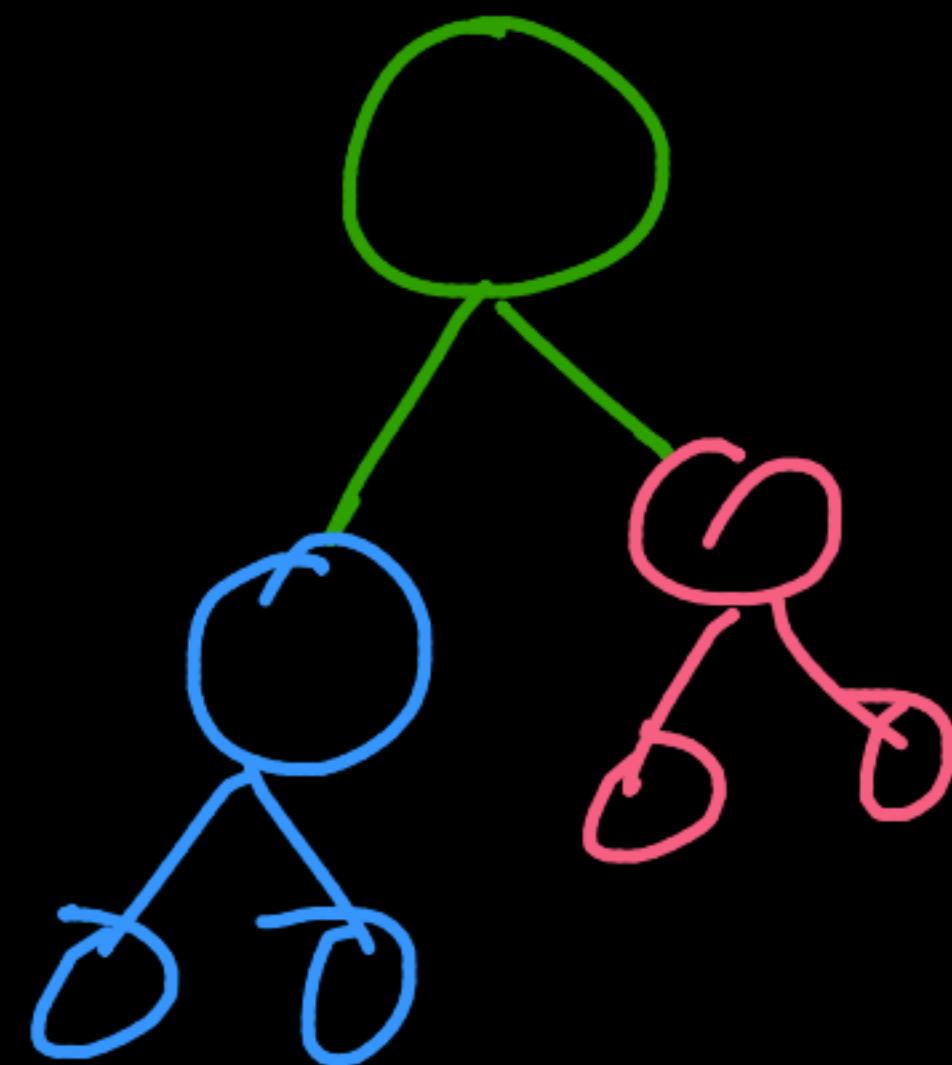
↳ DT of depth = 1



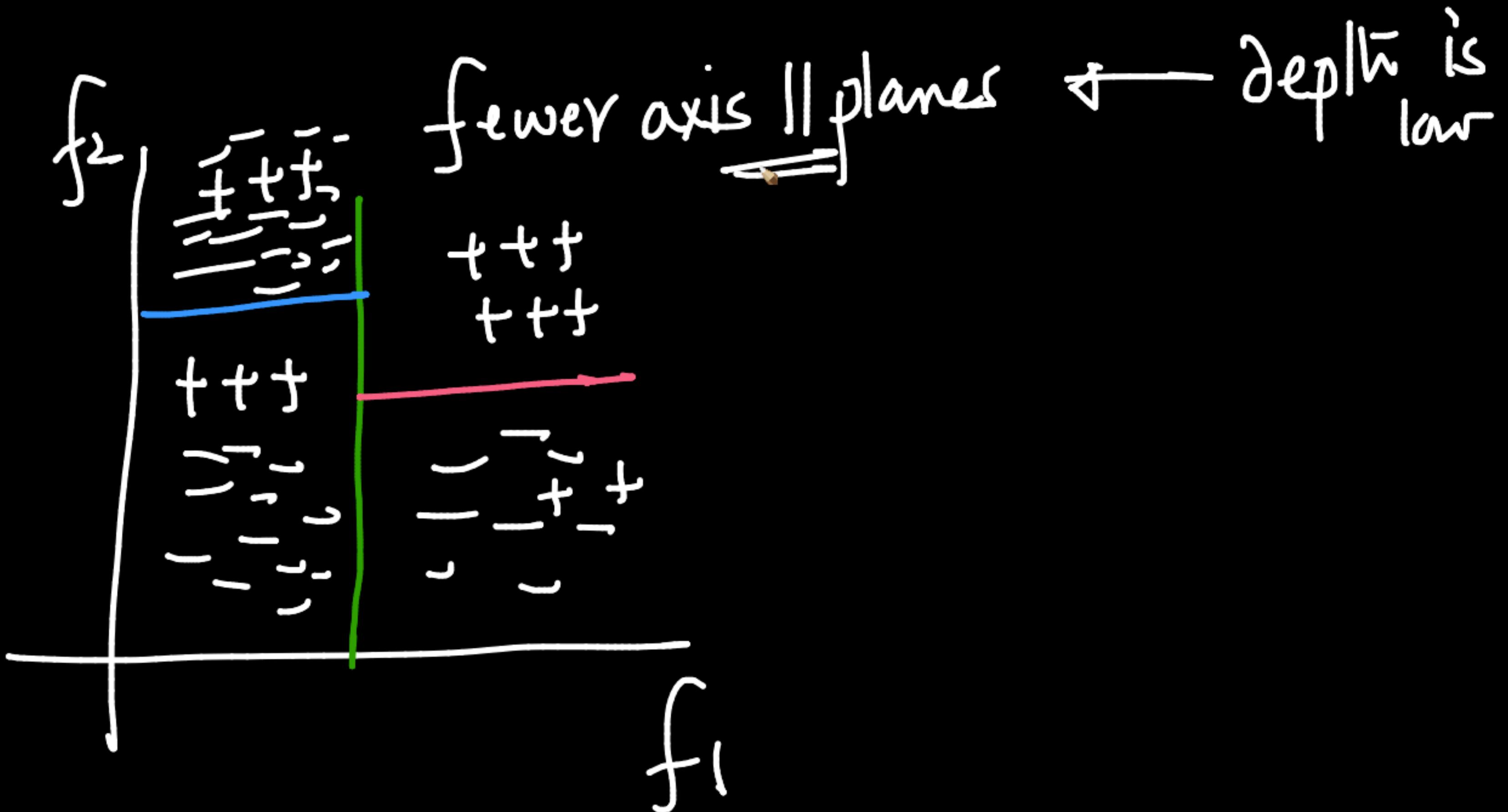
Shallow tree: depth is small

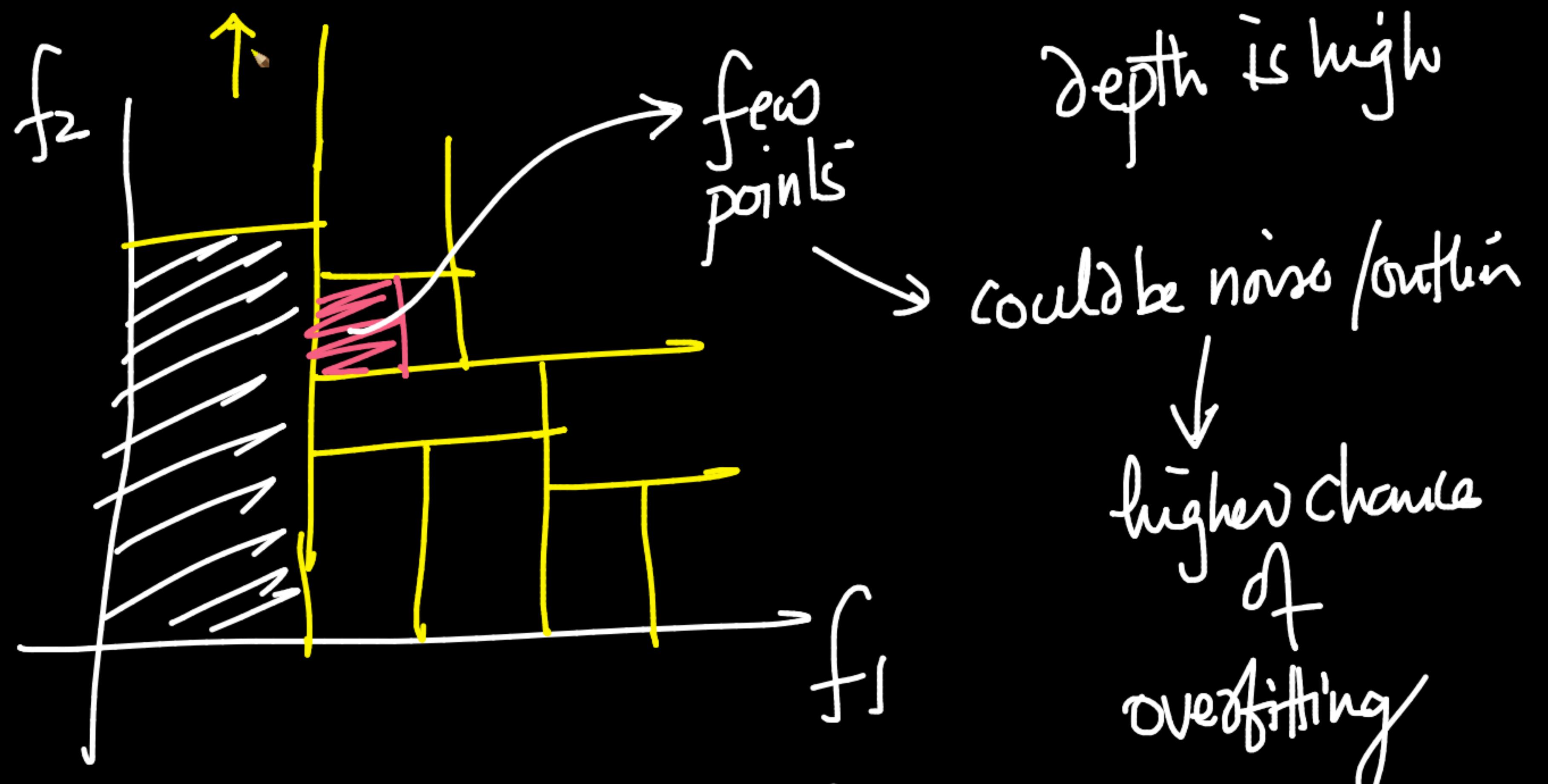
Deep tree: depth is large..

Geom:



Shallow - few



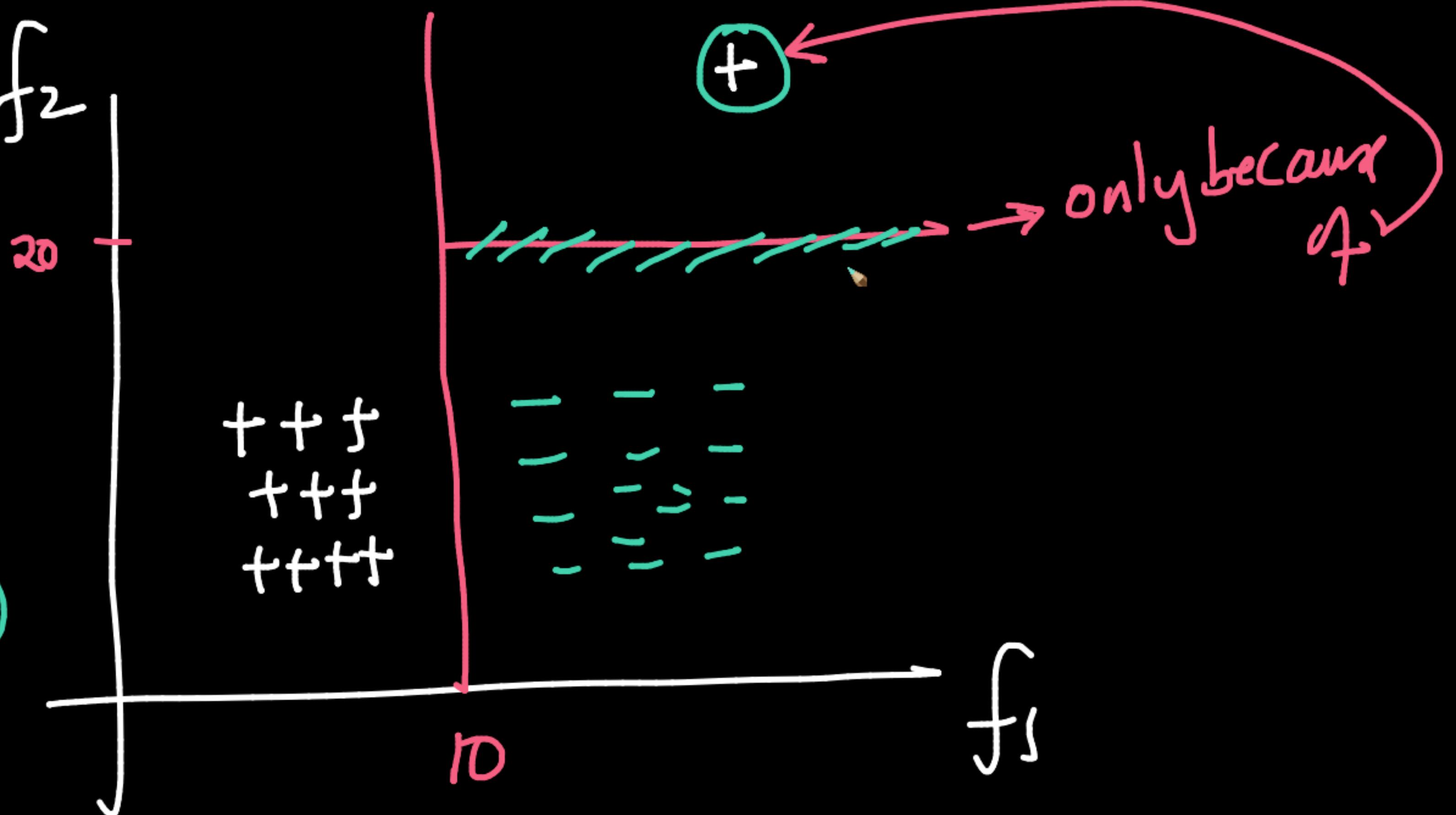
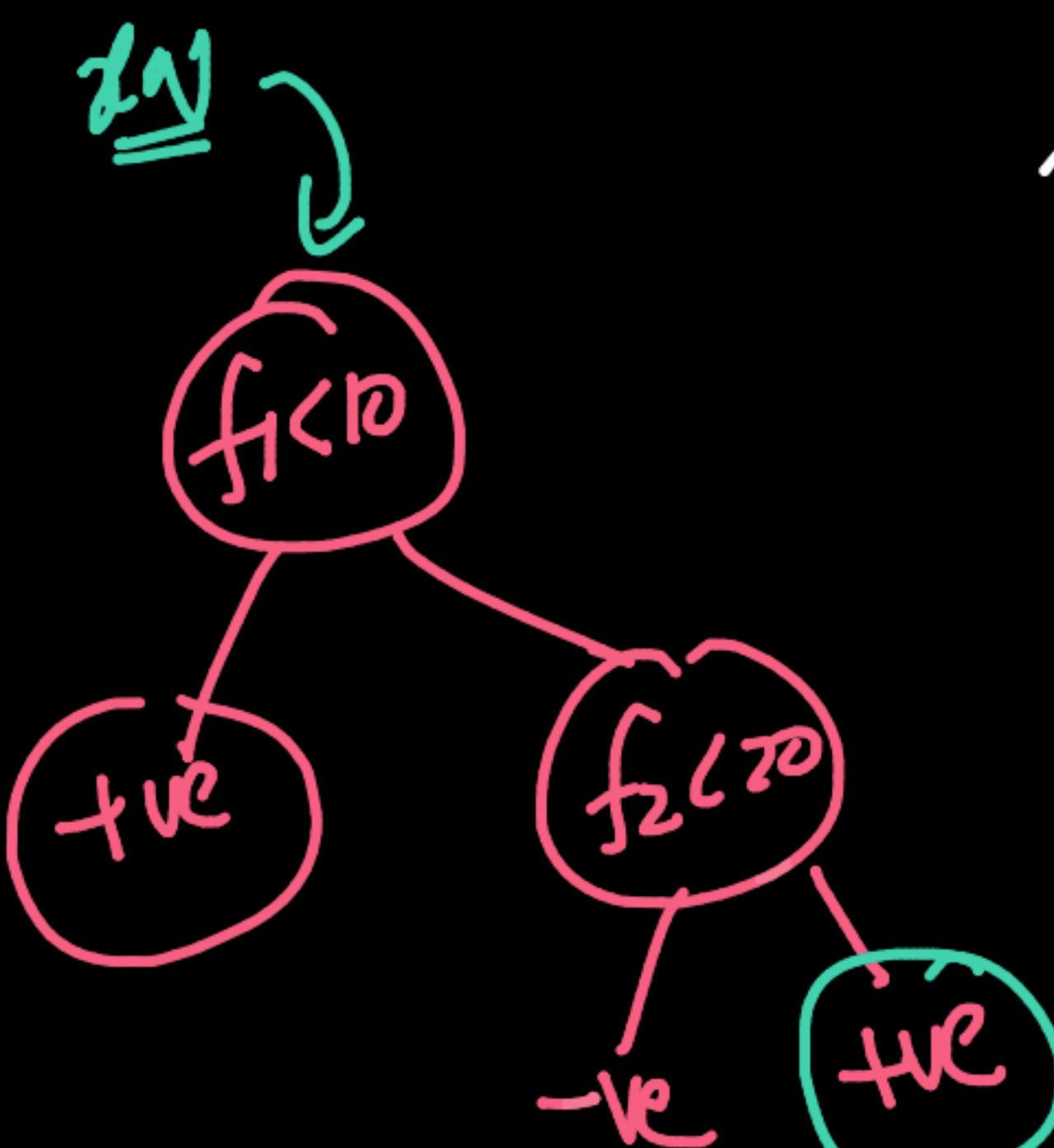


2D-rectangle

3D-cuboid

d-dim: hyper
cuboid

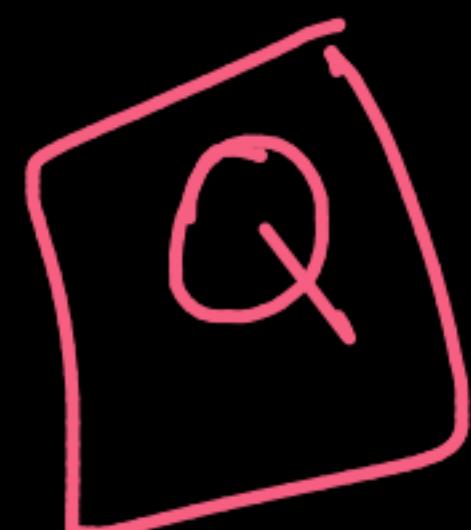
outliers → impact when depth is large



{ → non-linear data
→ 'd' is not too large

Tabular-data

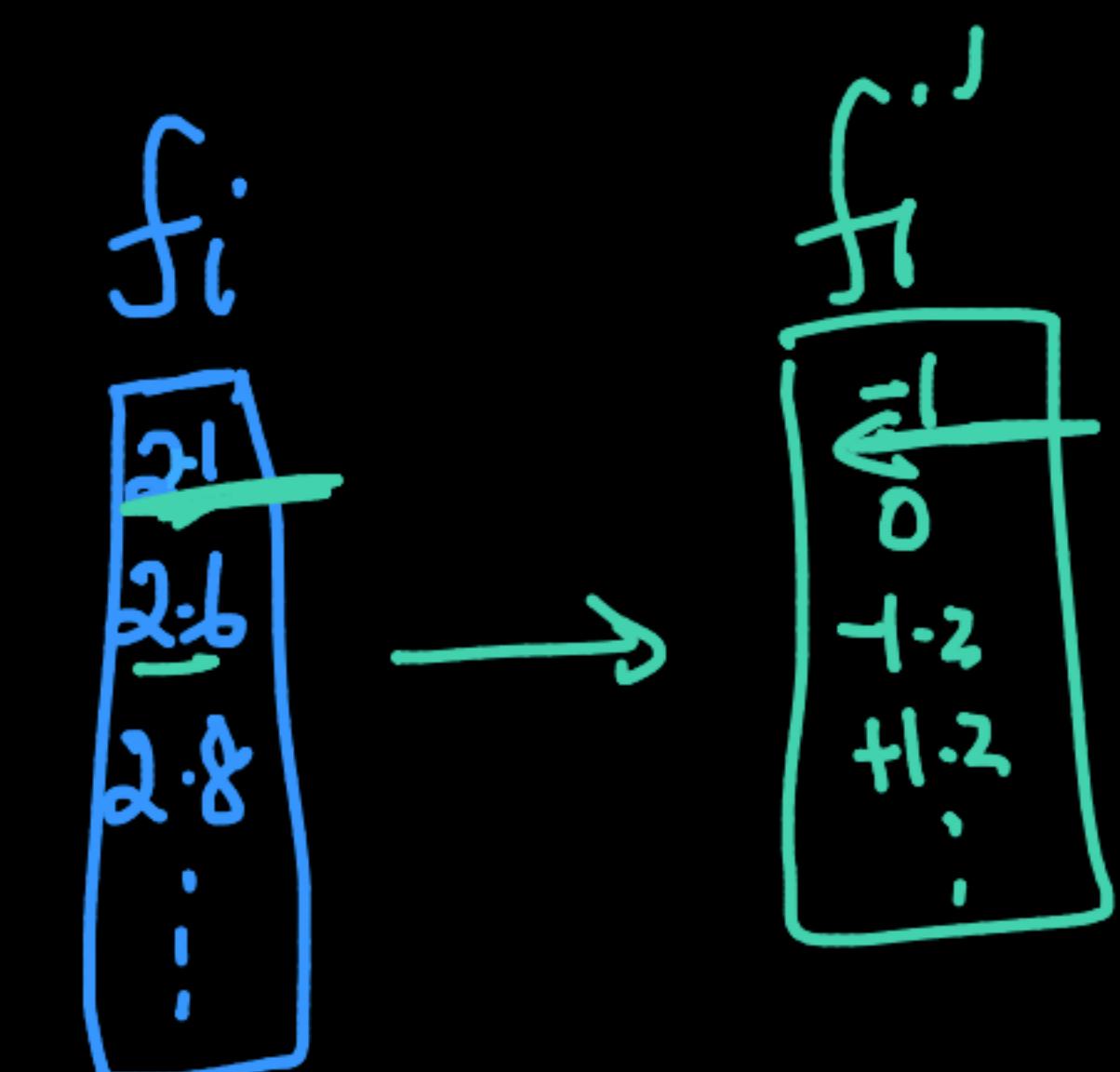
DT → RF & GBDT

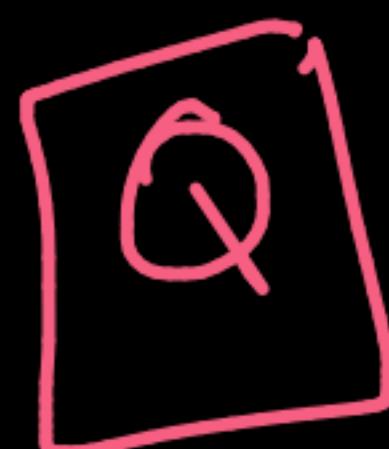


Linear & logistic & PCA → Standardize the data

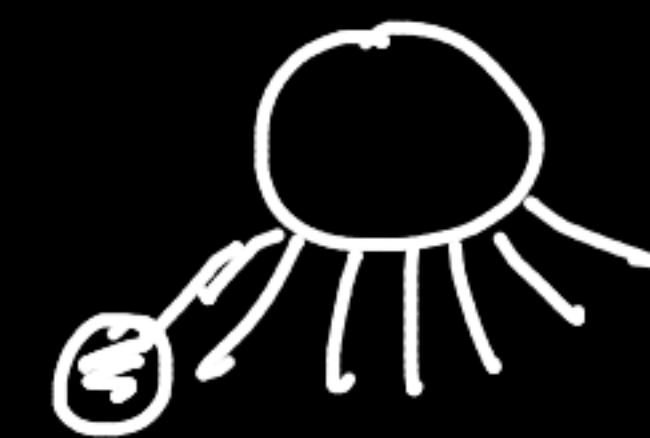
Is standardization of data needed for DT

does not impact IG





Categorical = features:



[Zipcode \rightarrow City]

Zipcode
~~10,000~~



encoding for DT

- no encoding \times \rightarrow (entropy
Measuring)

- Target = encoding \checkmark

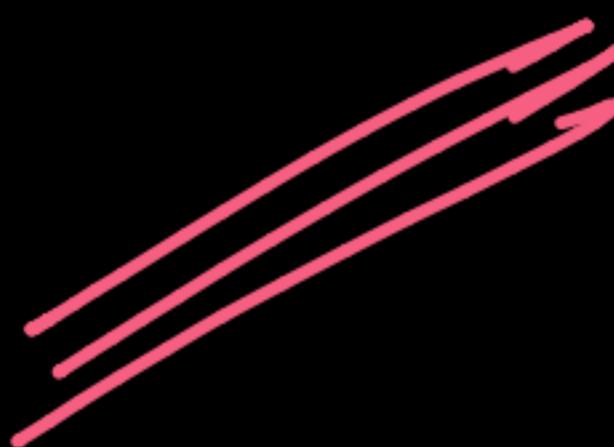
- One hot = encoding \times

$d \uparrow$

$n \uparrow$

$=$

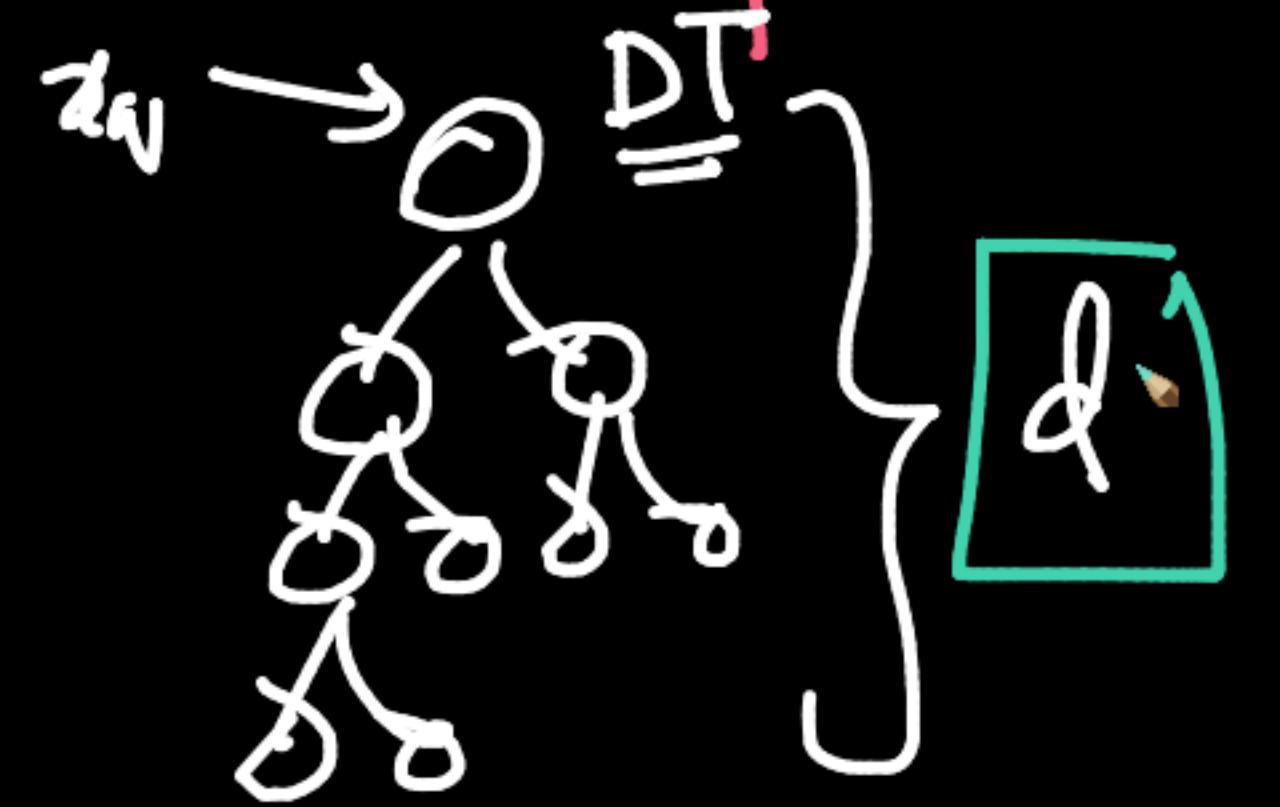
* appropriate feature
encoding also depends on model using



Time : $\mathcal{O}(\text{depth}) \rightsquigarrow (\lg m)$

Space: $\mathcal{O}(m) =$

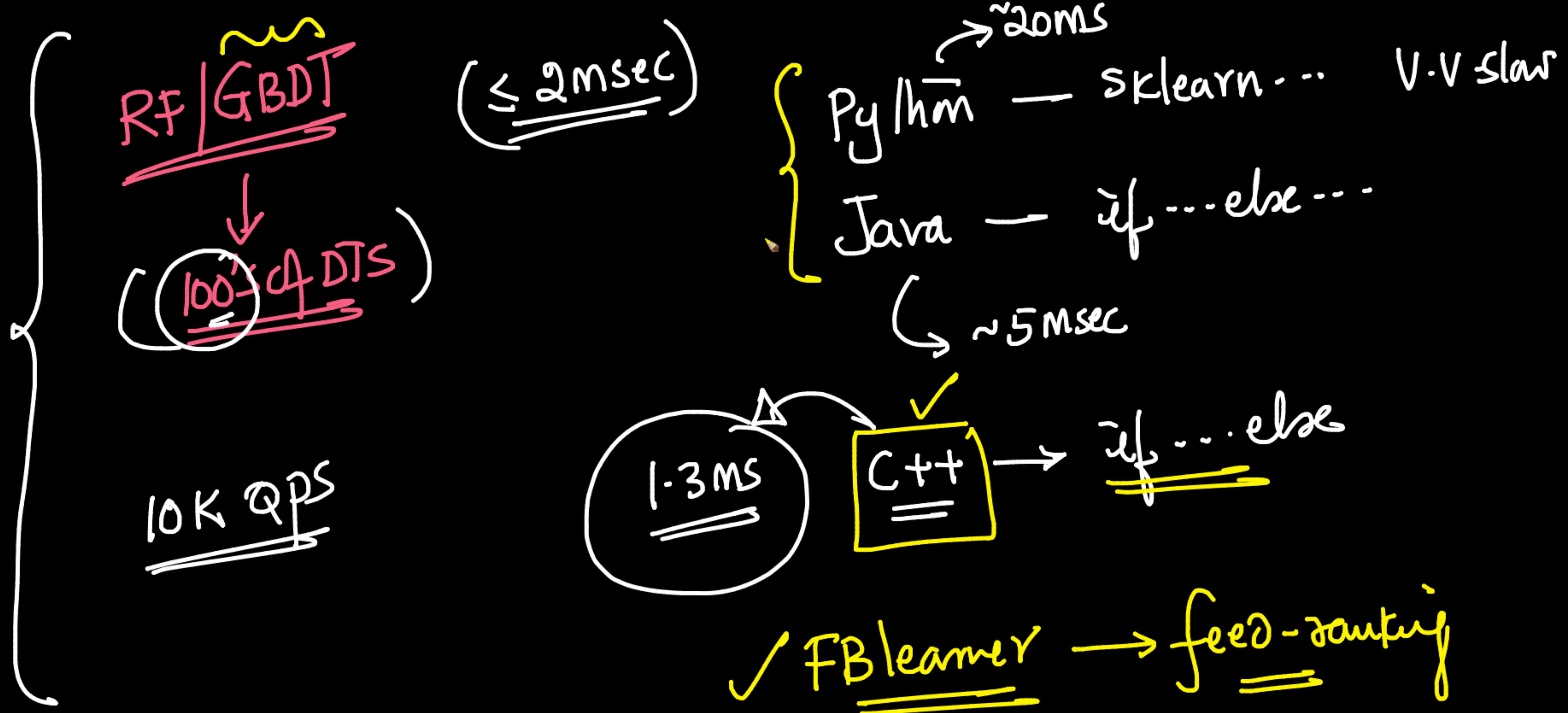
Runtime - complex

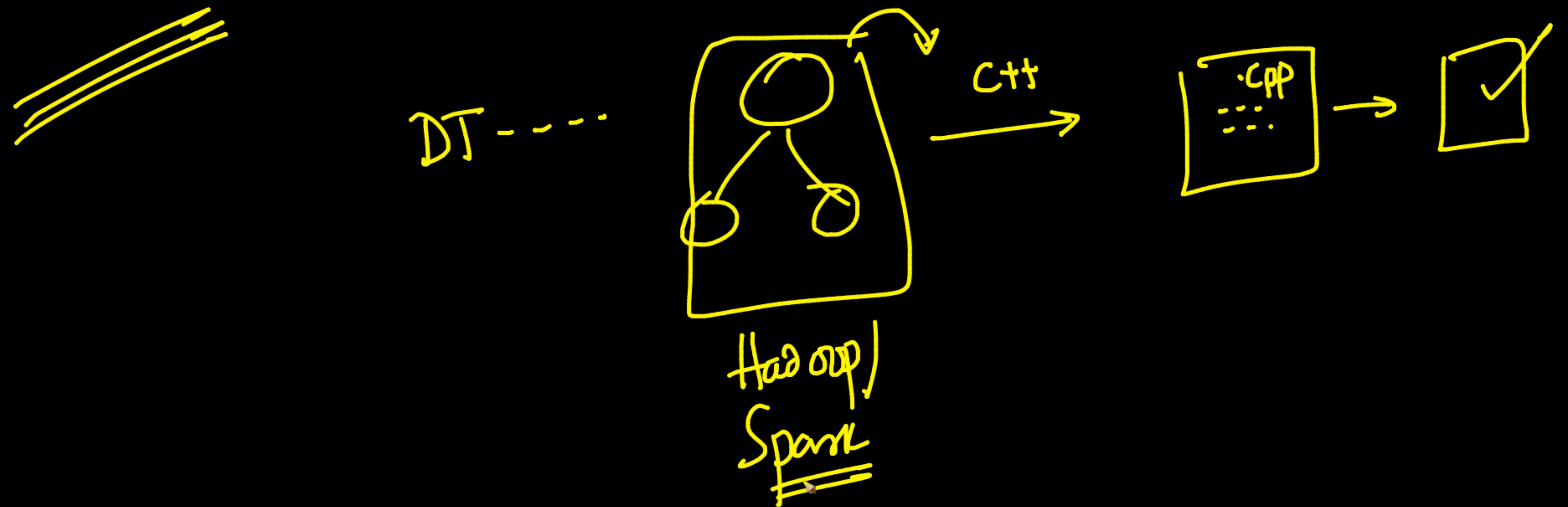


$m = \# \text{ nodes}$

$n = \# \text{ datapts}$

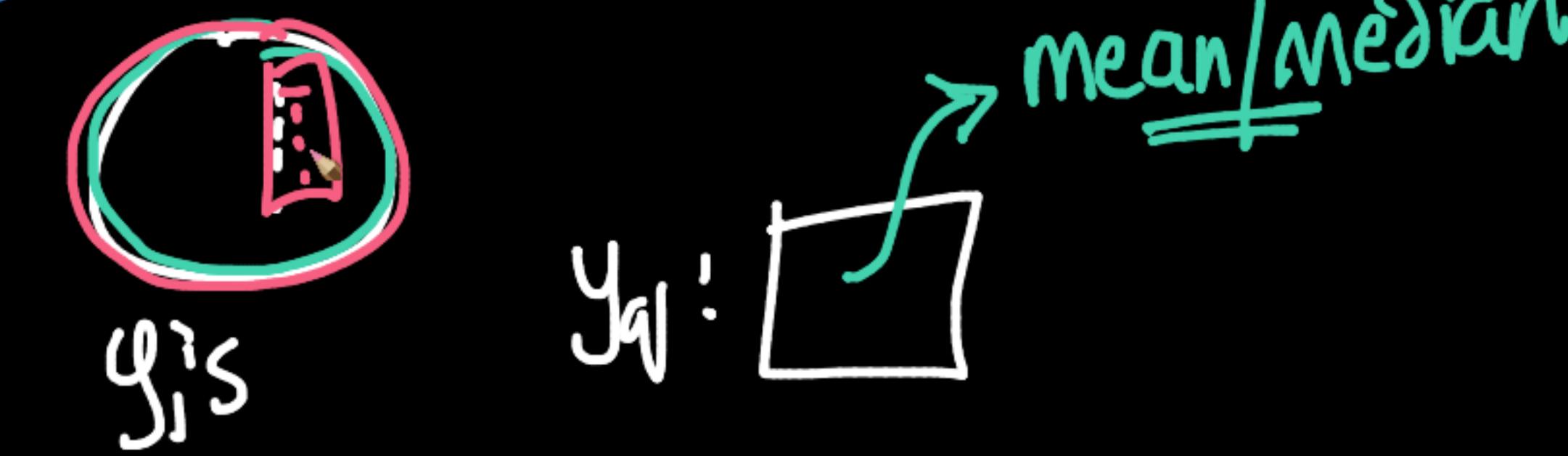
depth best using CV \rightarrow DT is very efficient
@ runtime





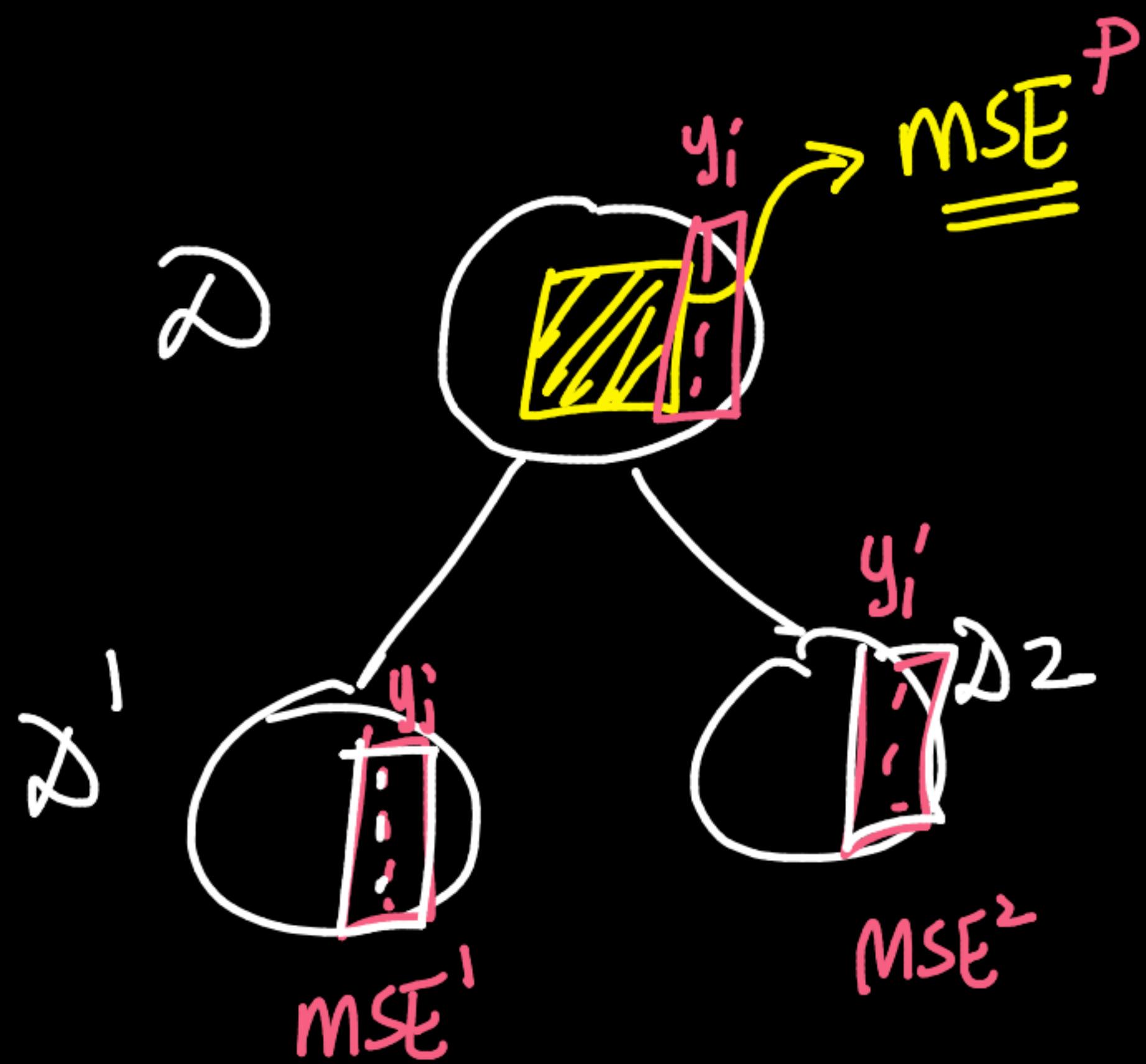
Regression on DT

- Entropy & GJ → alternative → MSE or MAE
- leaf-node



y_l : mean/median

Classfn:- $y_+ = 500$; $y_- = 200 \rightarrow$ we r $\frac{5}{7}$



$y_i \in \mathbb{R}$

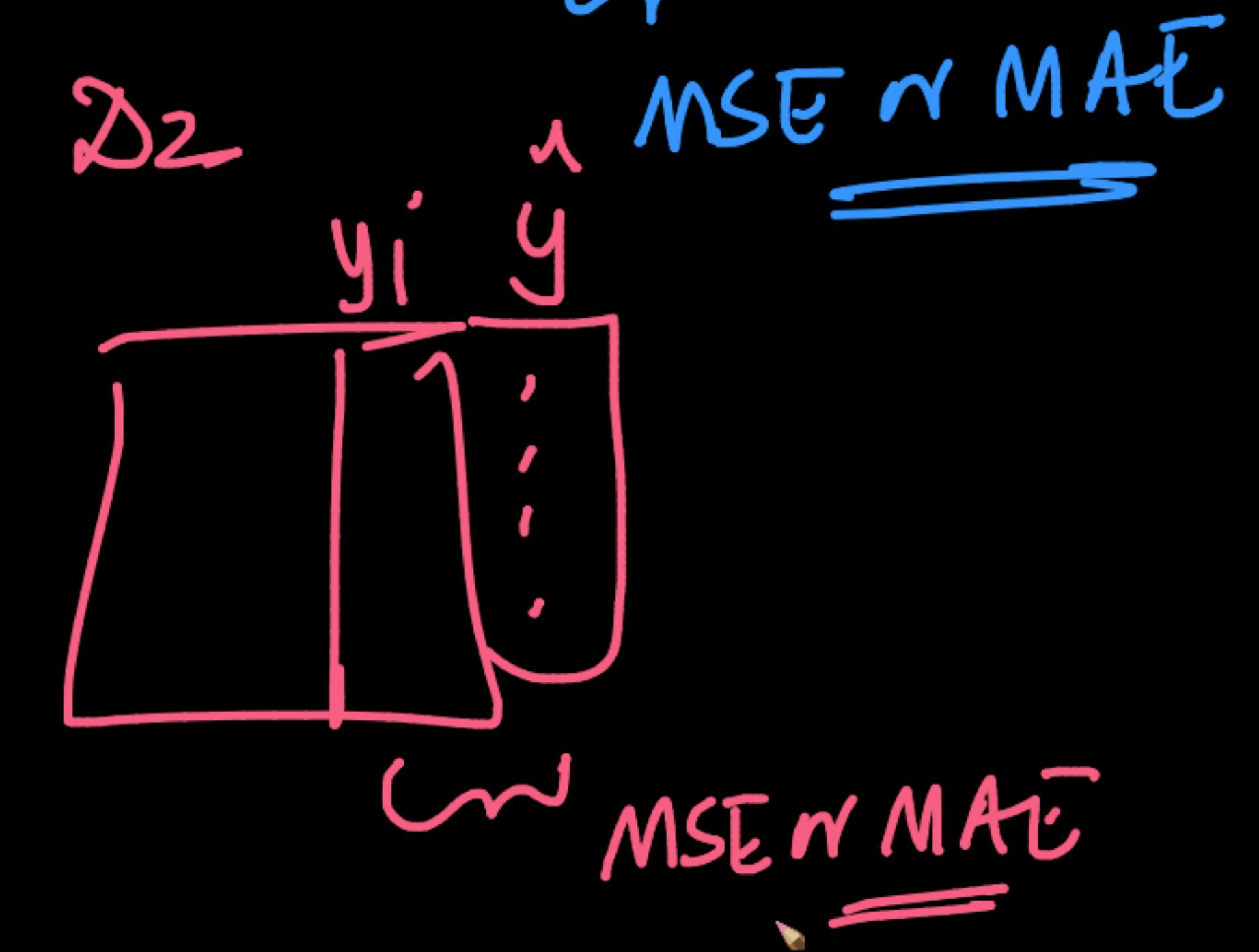
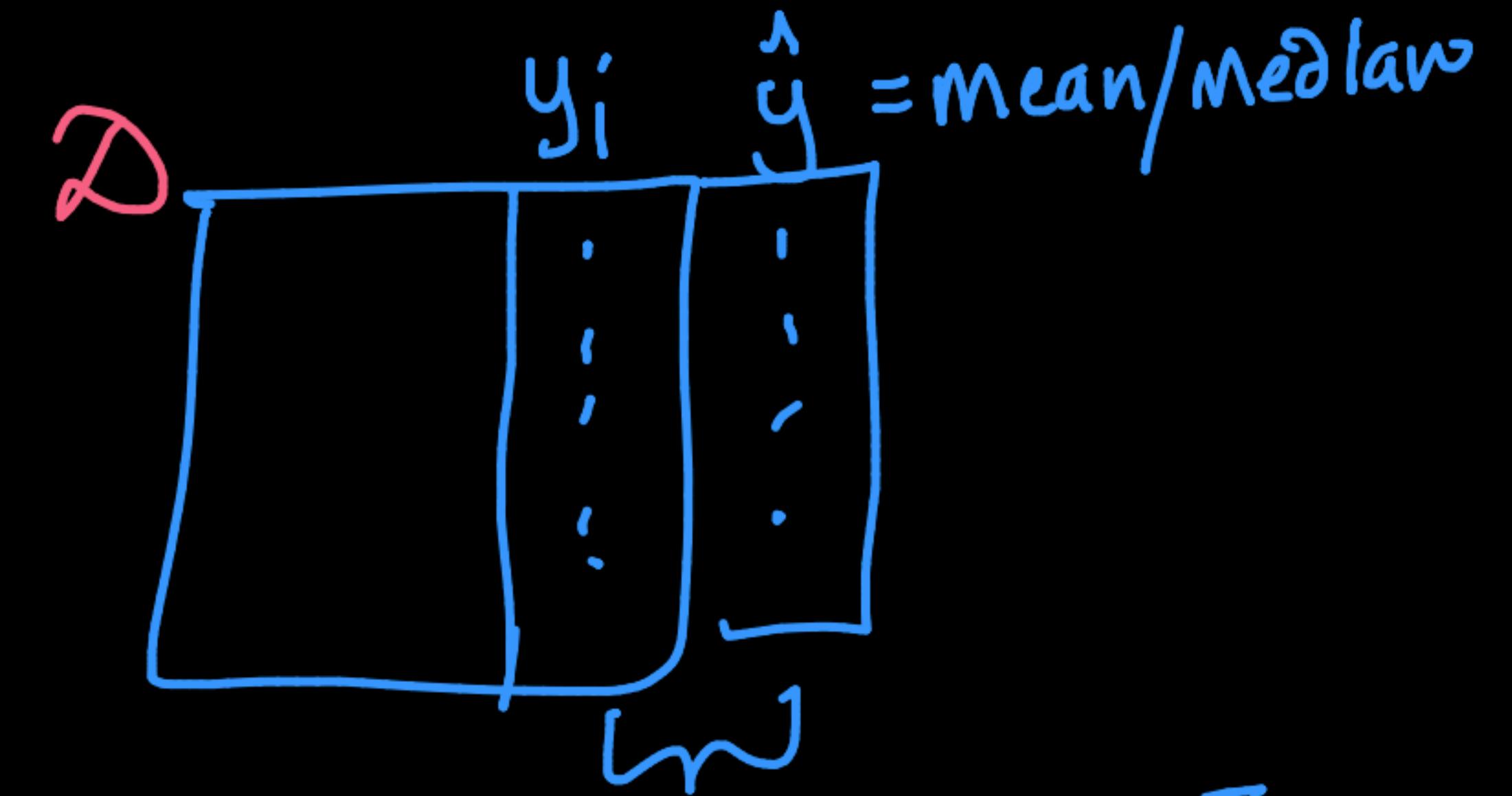
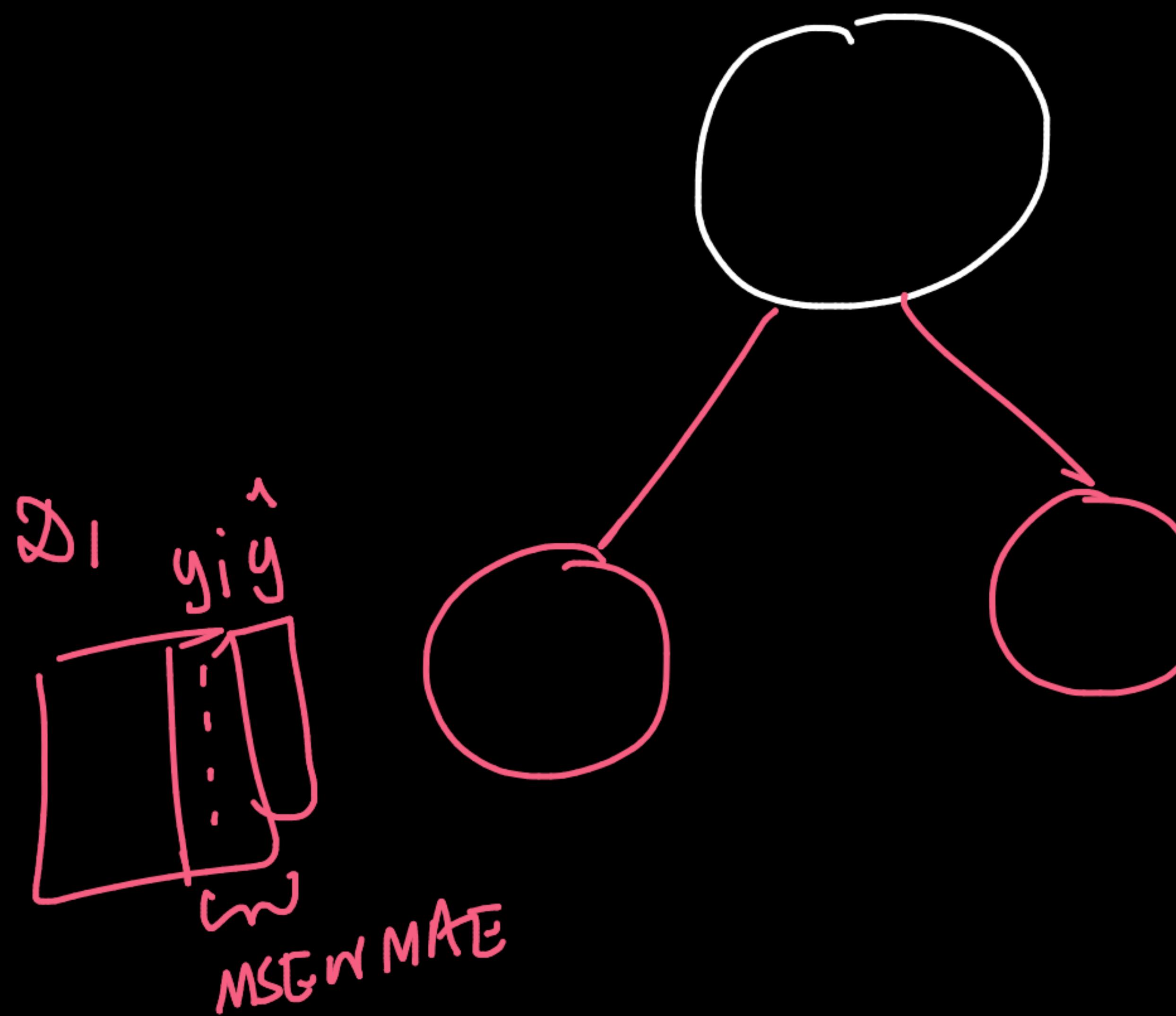
all y_i 's are same @
anode

\downarrow
 MSE or $MAE = 0$

{ weighted MSE of child nodes
- MSE of parent nodes

MSE : lowest \rightarrow all y_i 's are same

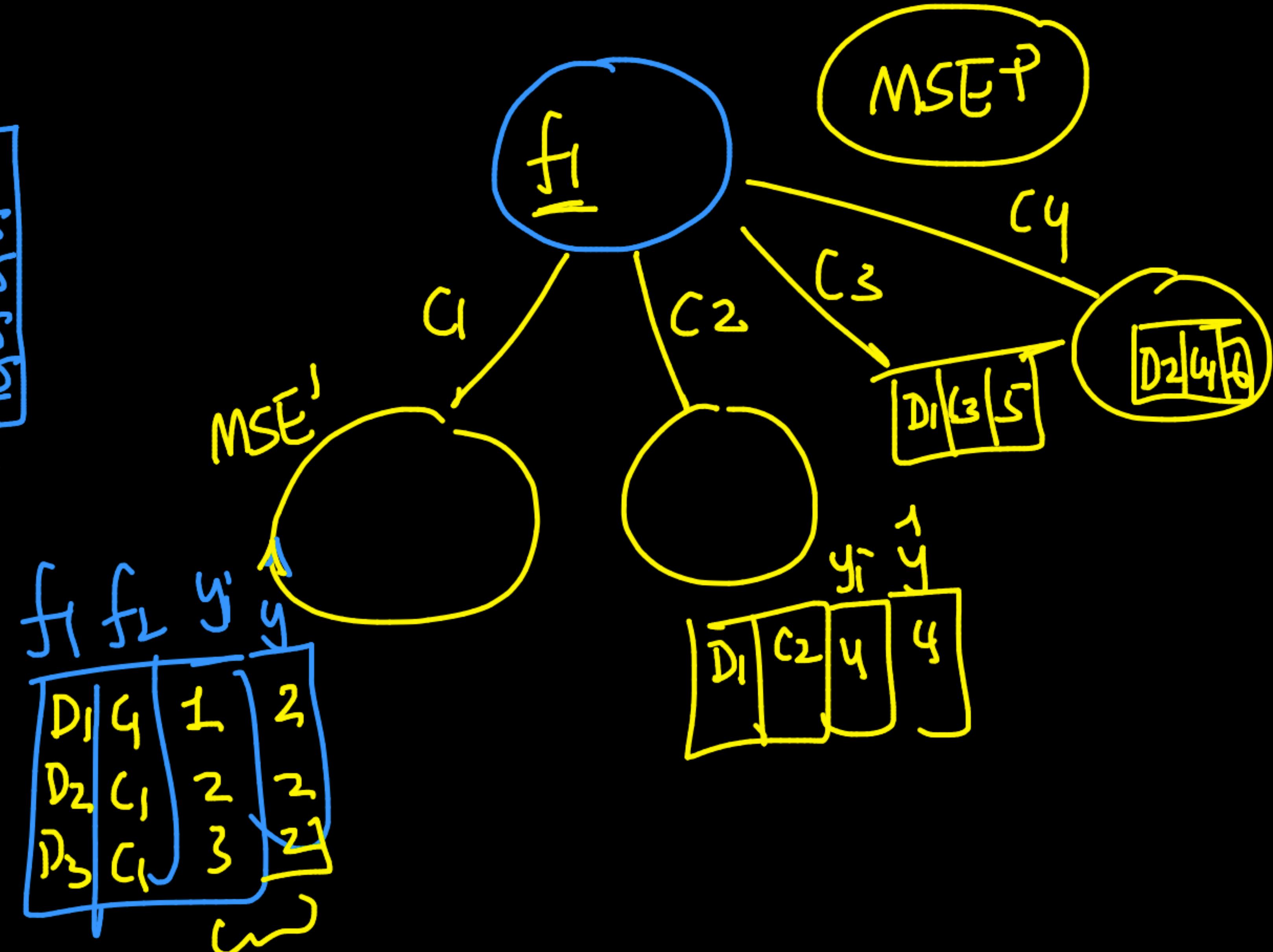
MSE : Max \rightarrow y_i 's are very divergent

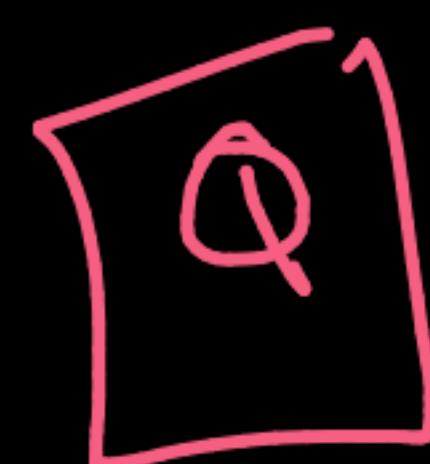


$n=6$

f_2	f_1	y_i	\hat{y}
D_1	C_1	1	3.5
D_2	C_1	2	3.5
D_3	C_1	3	3.5
D_1	C_2	4	3.5
D_1	C_3	Σ	3.5
D_2	C_4		

$\Delta \bar{MSE} \rightarrow f_1$
 $\Delta \bar{MSE} \rightarrow f_2$





Entropy or G · | ?

→
$$-\sum_{i=1}^k \log(\tilde{P}(y_i)) \tilde{P}(y_i)$$



Q

$\overbrace{\text{MSE}}$ or $\overbrace{\text{MAE}}$?

outlier,
impact is
high

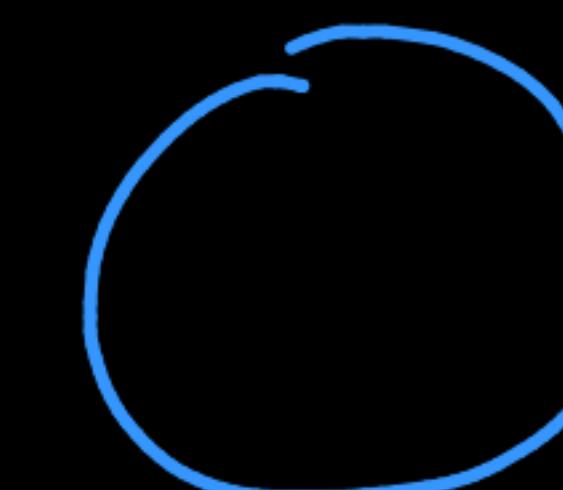
↳ outlier resistant

Var vs MAD

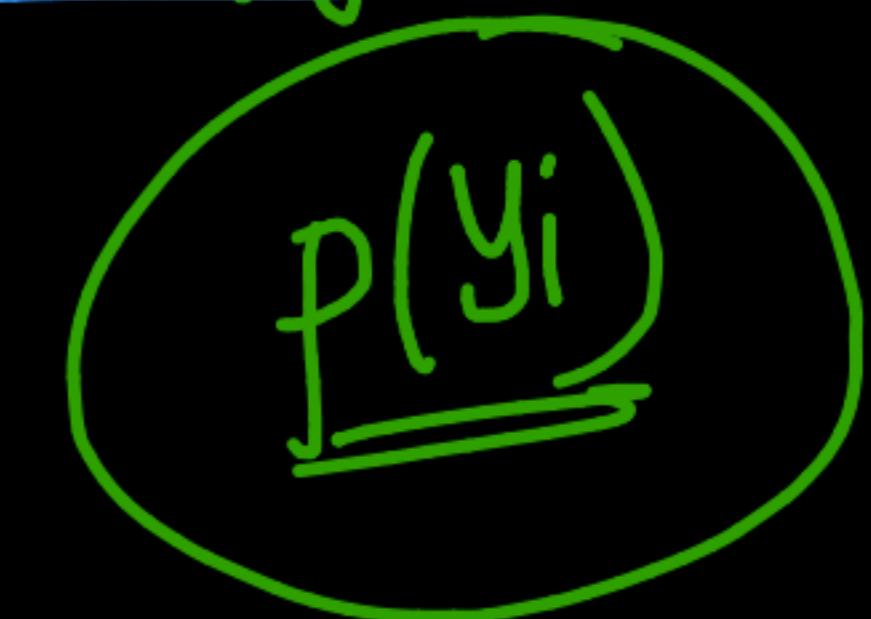


Imbalanced data in DT \rightarrow class weights; oversampling
SMOTE

\hookrightarrow rebalance it

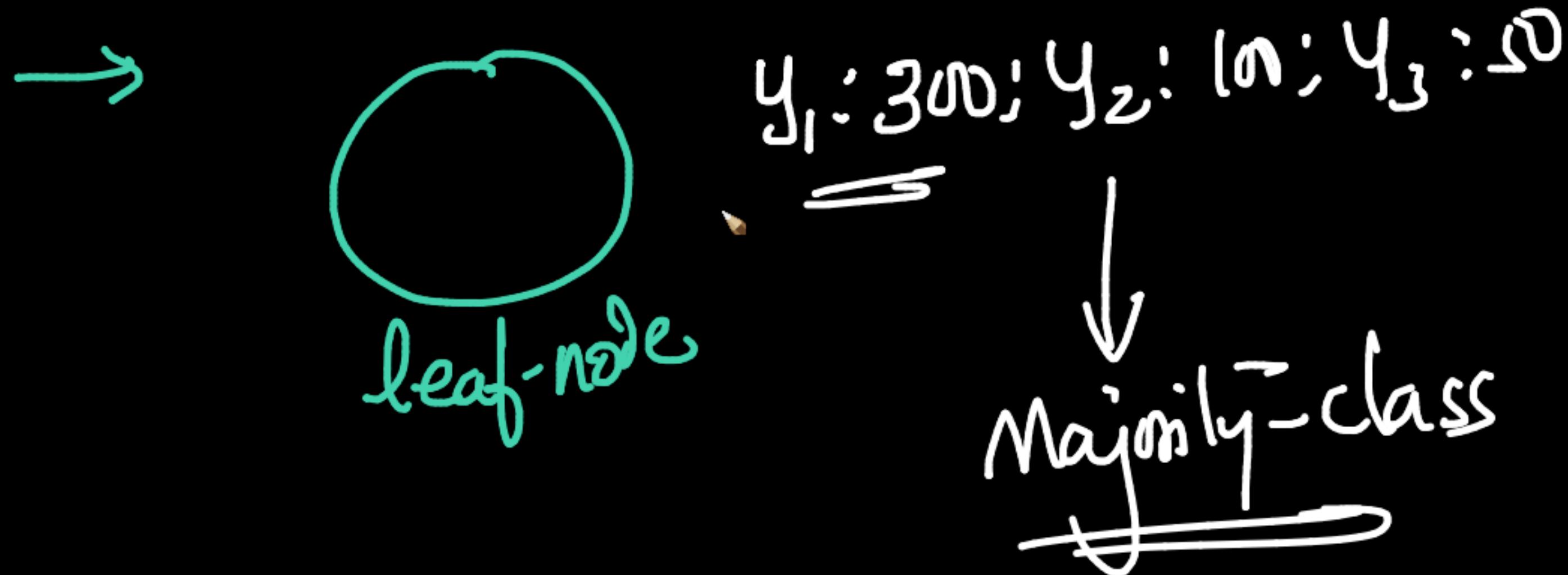


$y_+ : 99\%$
 $y_- : 1\%$



Multi-class classfn:-

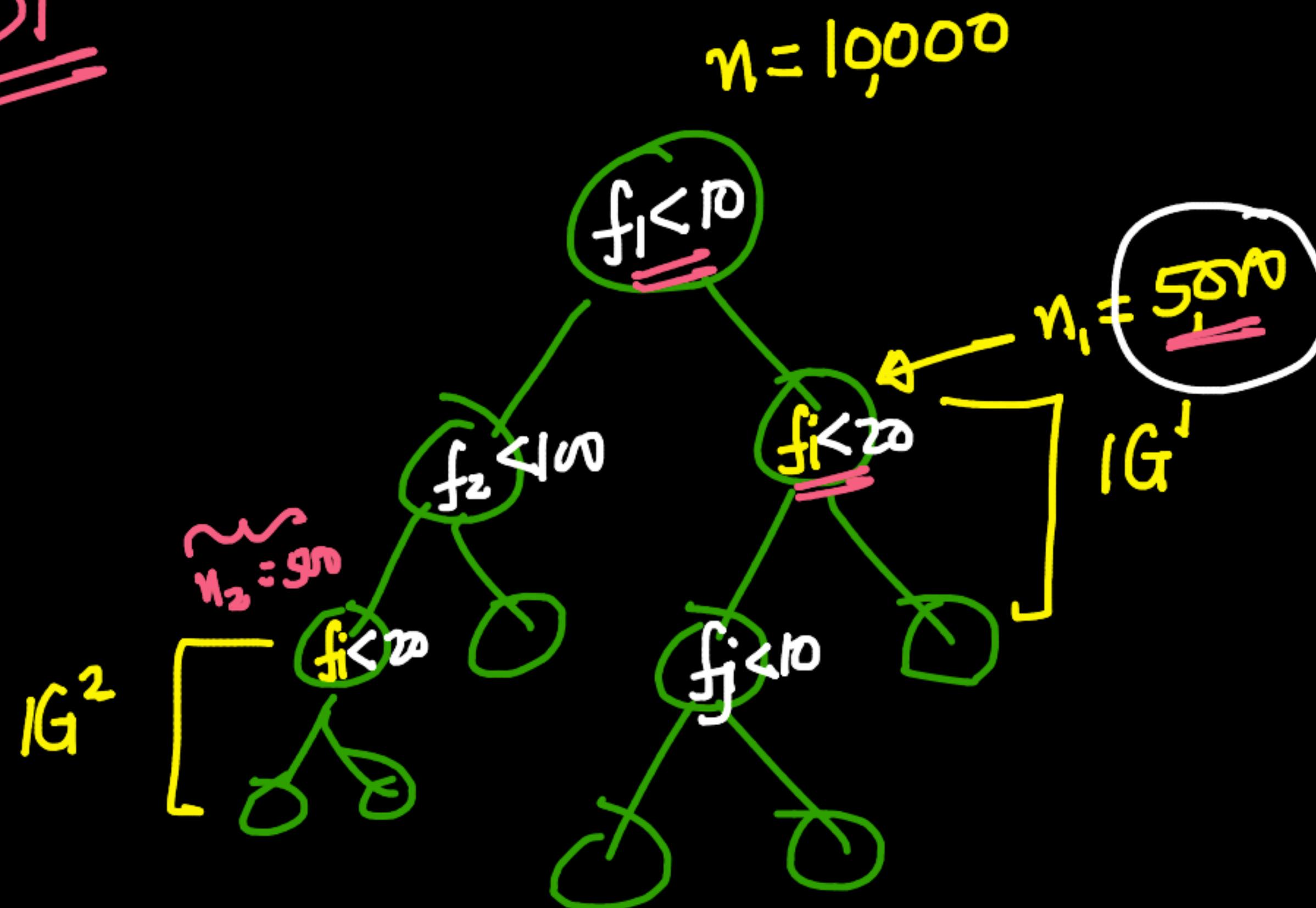
→ entropy & G.I ✓
⋮



Interpretable:-

if ... else ... } → non-expert

DT

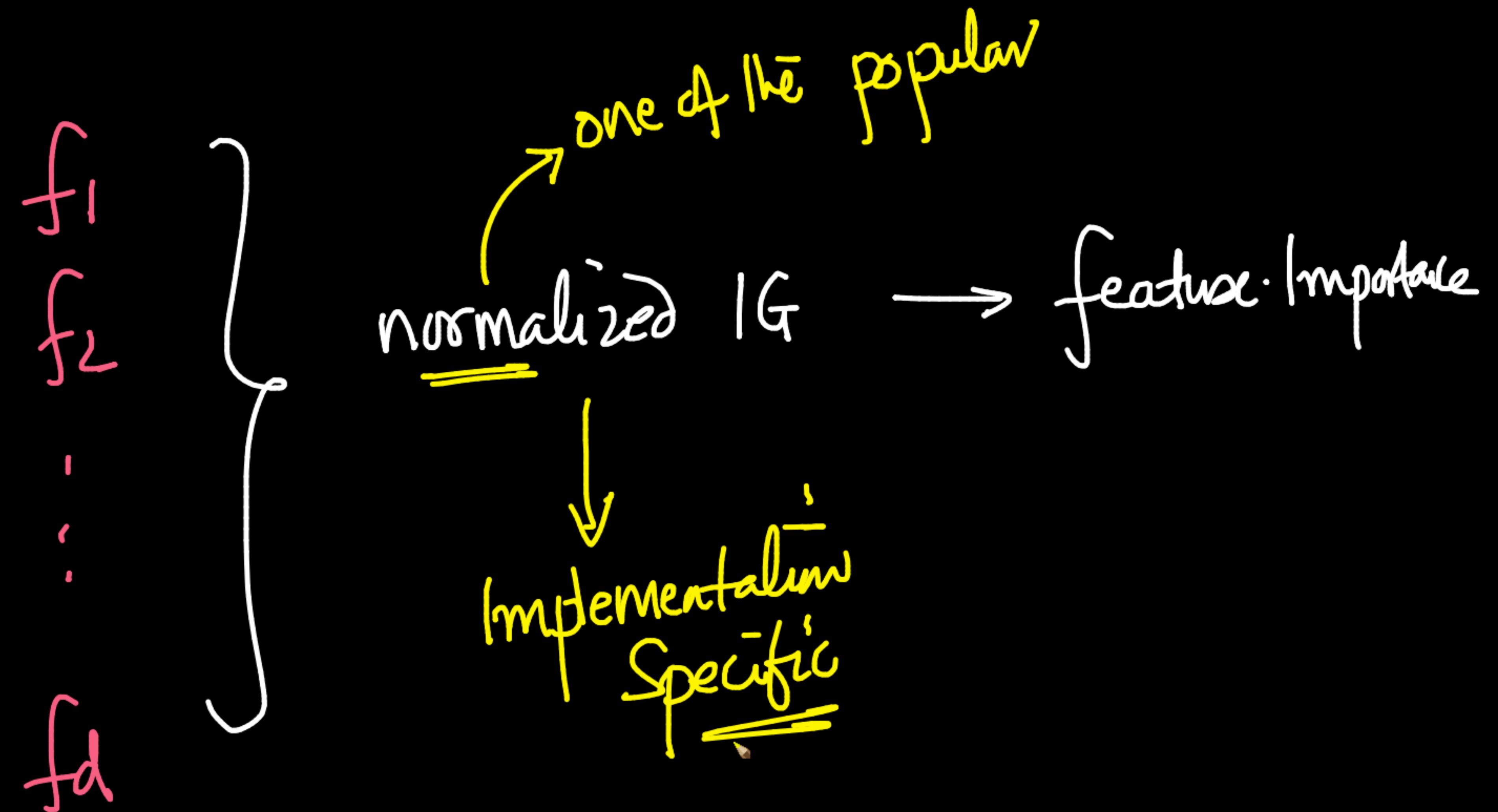


$$f_i := \underbrace{IG^1}_{} + \underbrace{\frac{5000}{10000}}_{IG^2} + \checkmark$$
$$IG^2 + \frac{500}{10000}$$

Feature Importance

→ IG ✓
→ order of nodes
normalized info-gain ...

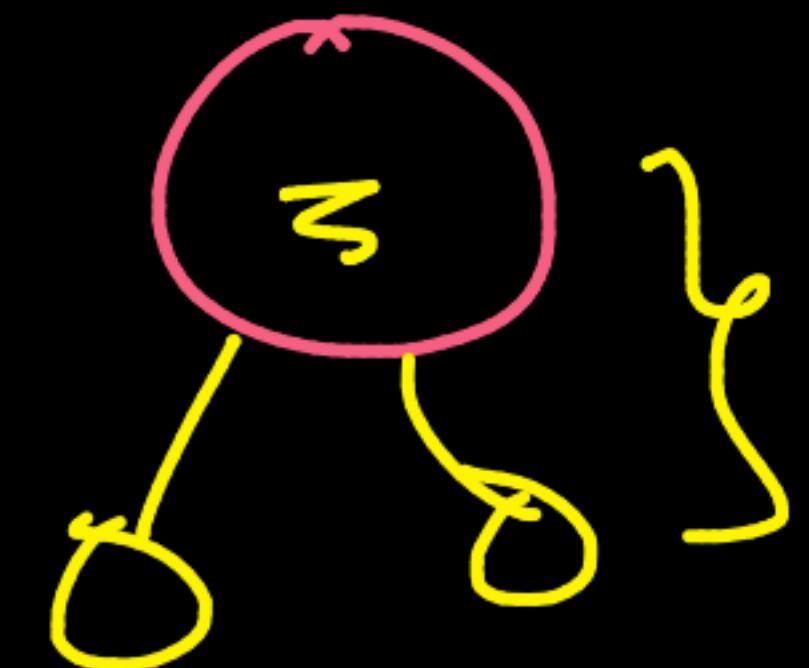




FYI

Optimal DT \rightarrow NP-complete (exp. han
time compl.)

greedy approx' malem

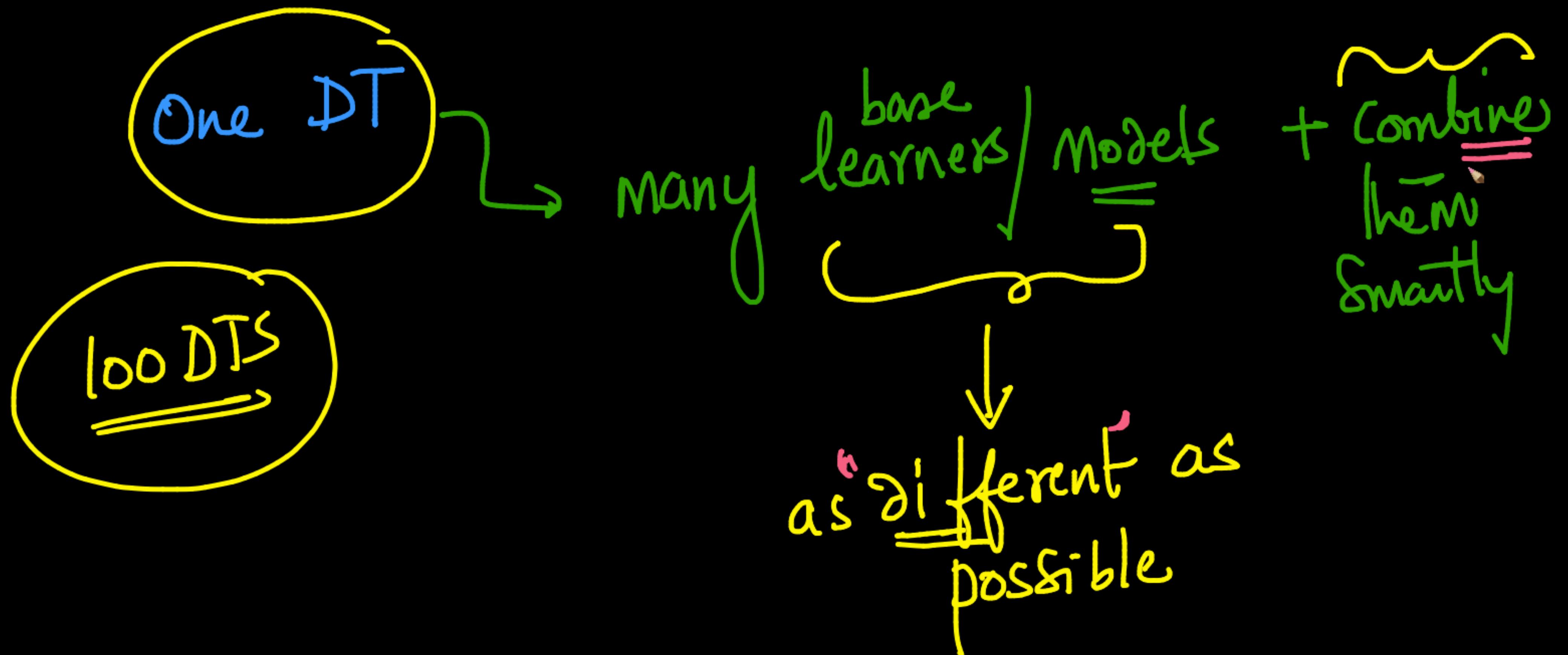


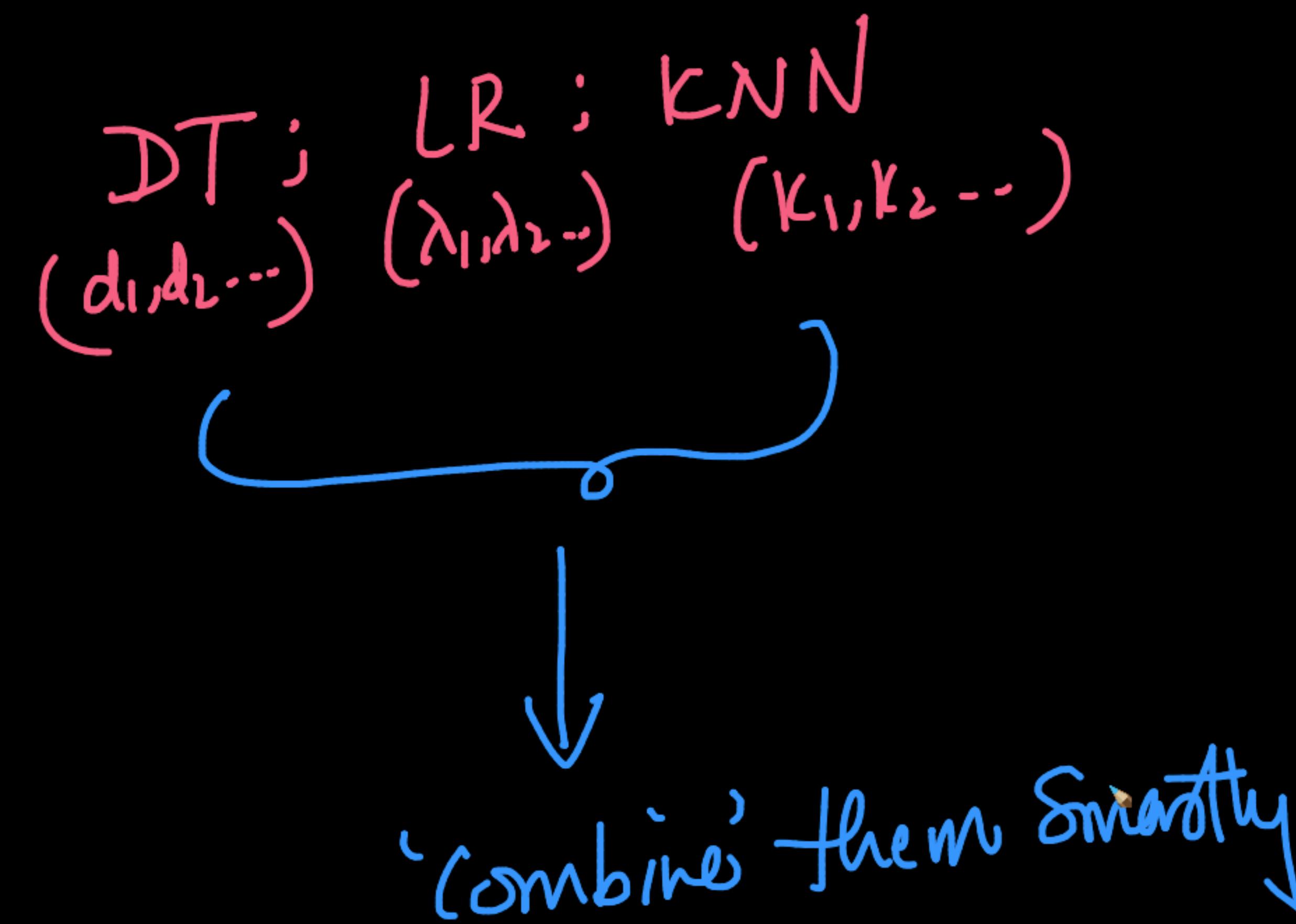
DT

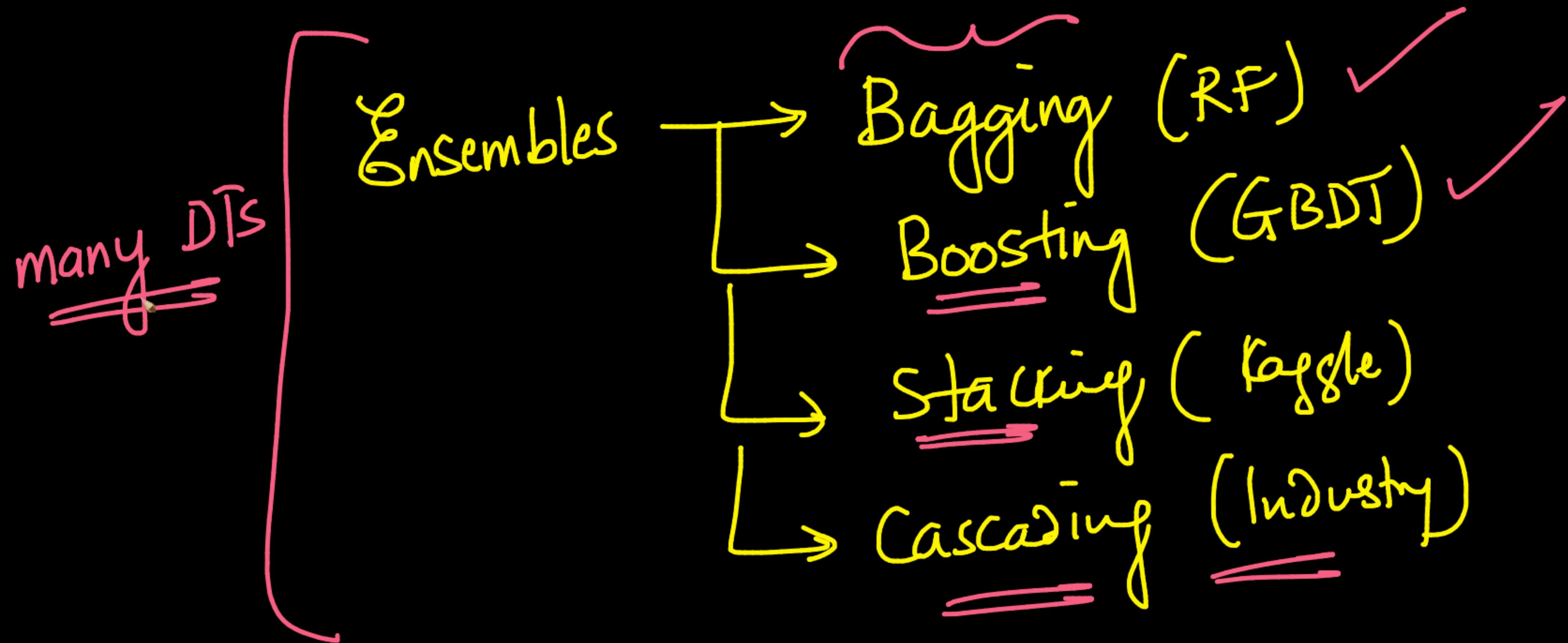
- d is not too large
- low runtime complexity
- no standardization needed
- depth is the key hyperparam

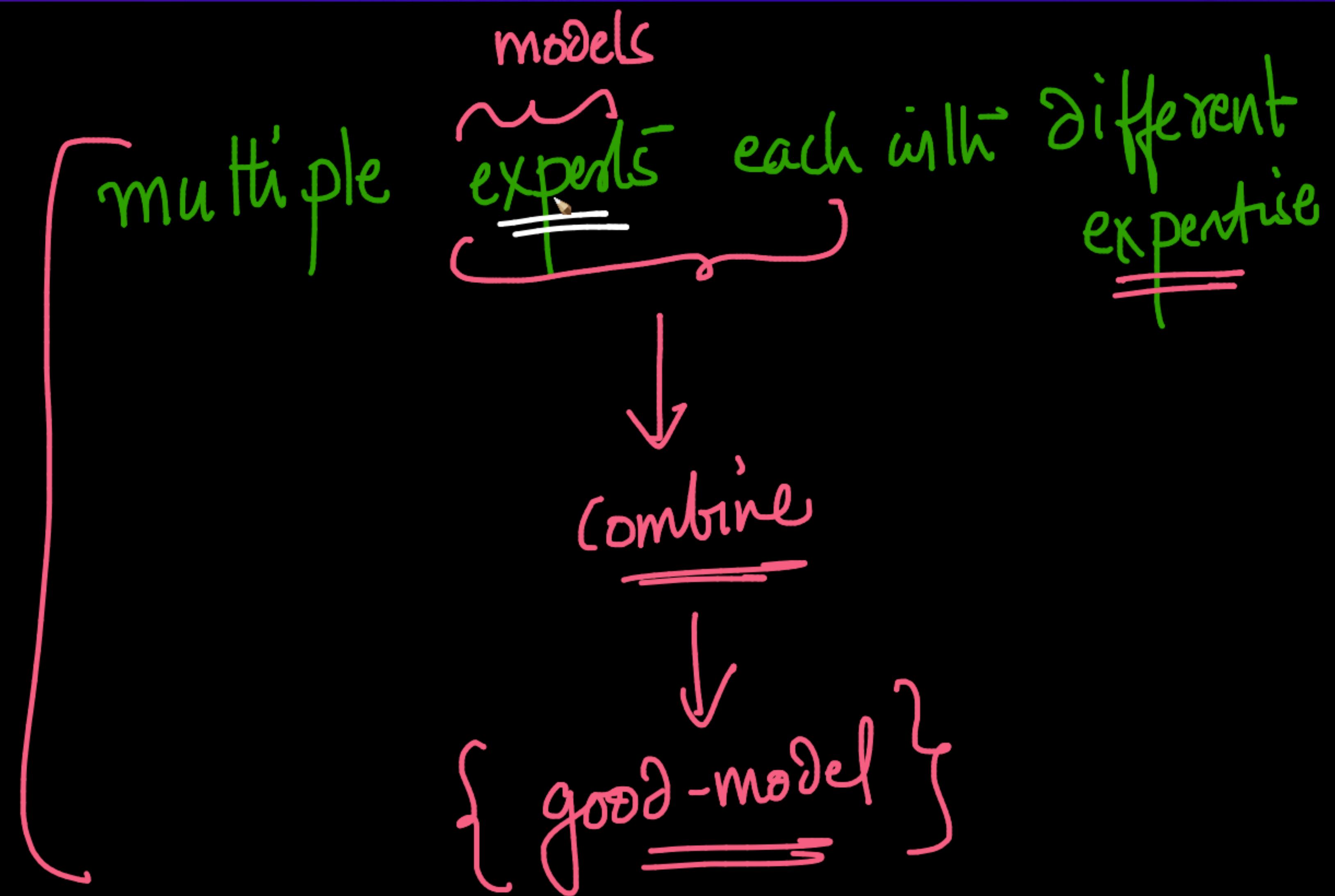


Ensembles:









INTUITION:

$$\text{look } m < n$$

~~Bagging~~:

Bootstrapped Aggregation

Train

$$\mathcal{D}_n$$

sample
m with
replacement

$$\mathcal{D}_m^1$$

$$M_1$$

& CV on remain
 $n-m$ pts

$$\mathcal{D}_m^2$$

$$M_2$$

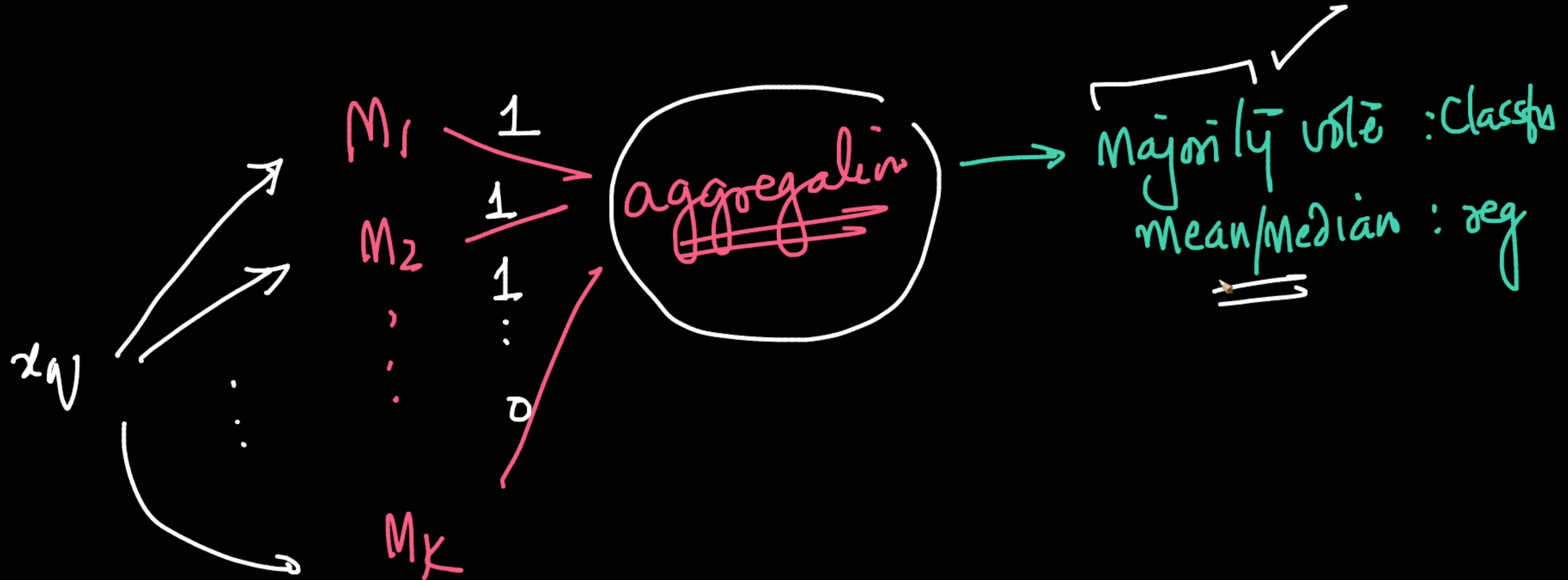
" "

$$\mathcal{D}_m^K$$

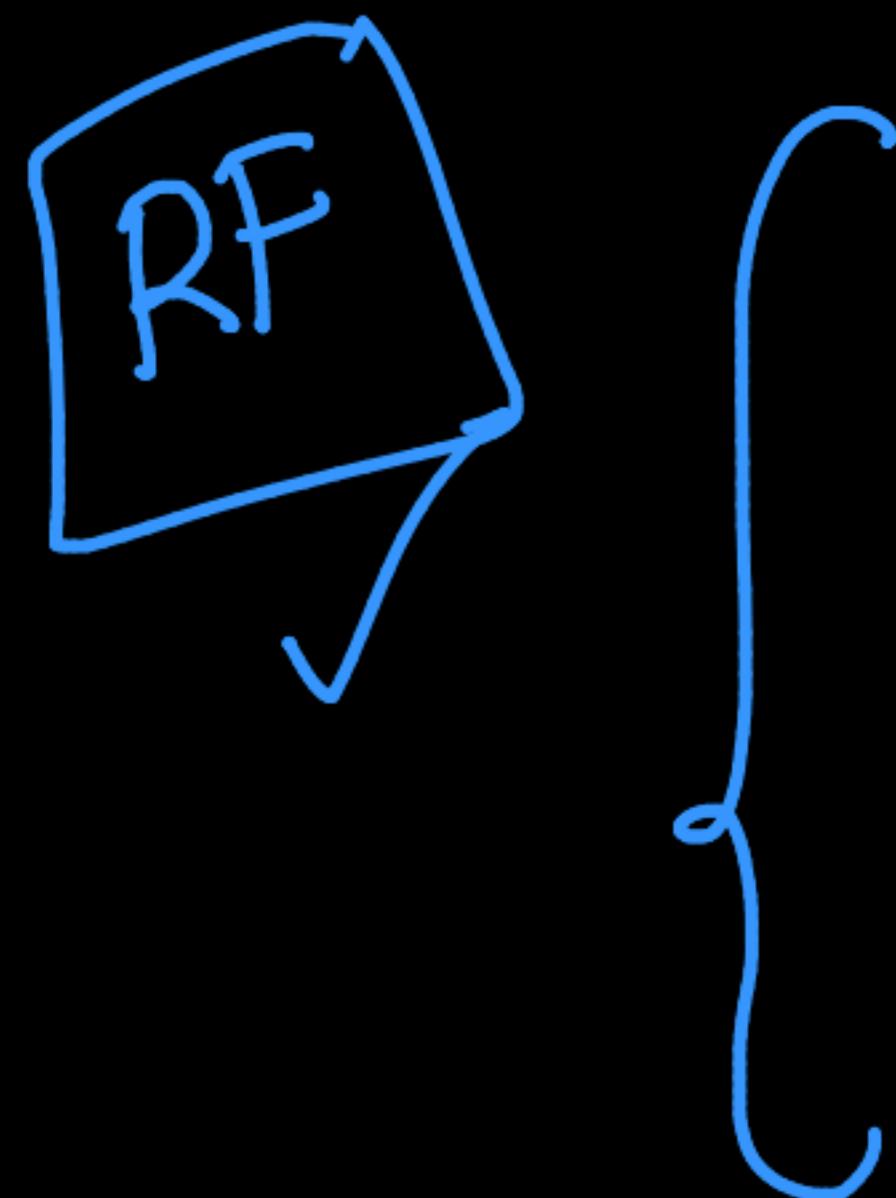
$$M_K$$

" "
base-learners

$$K=100$$



Challenge!



already happening
in m sampling

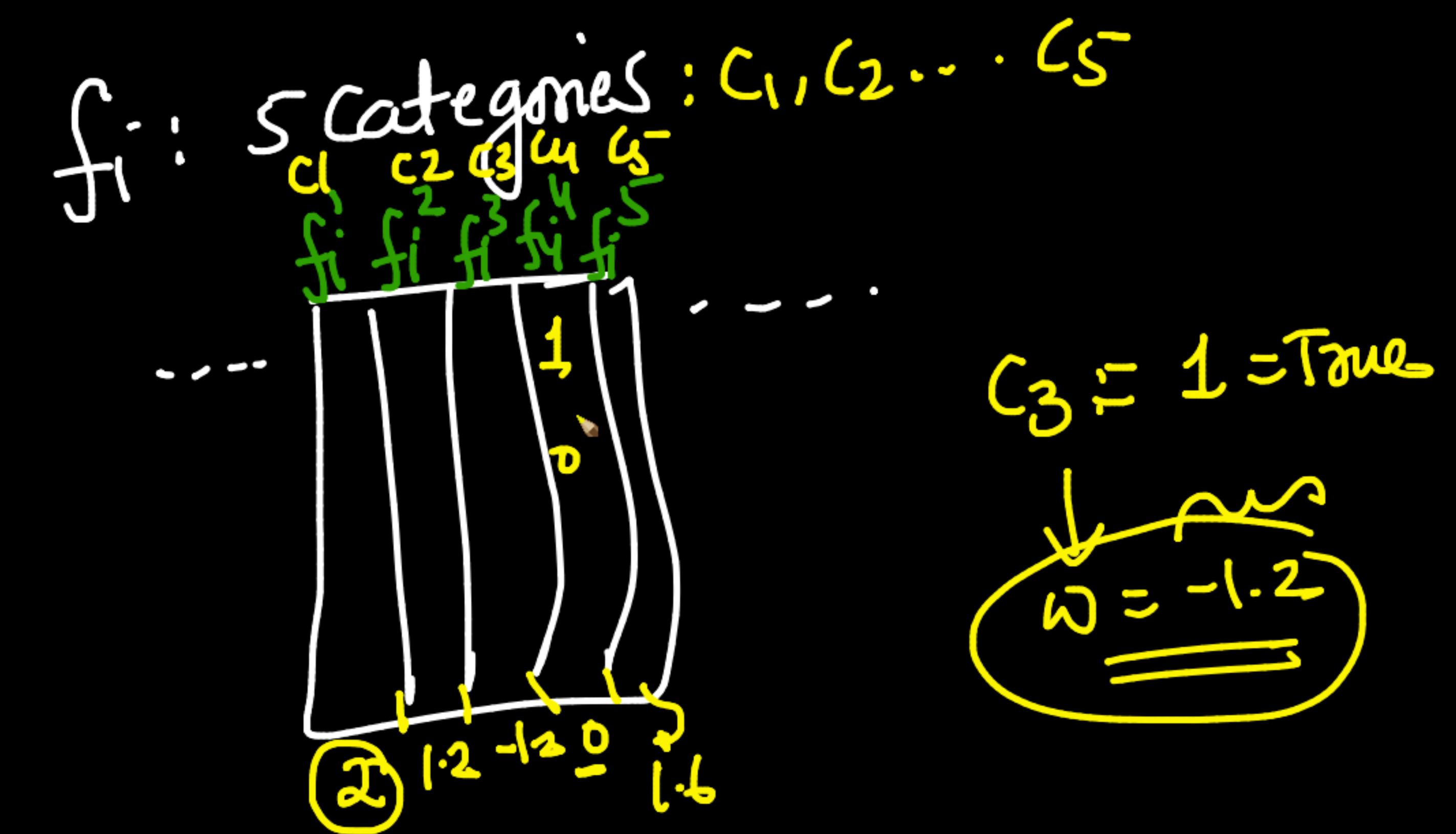
ideas:

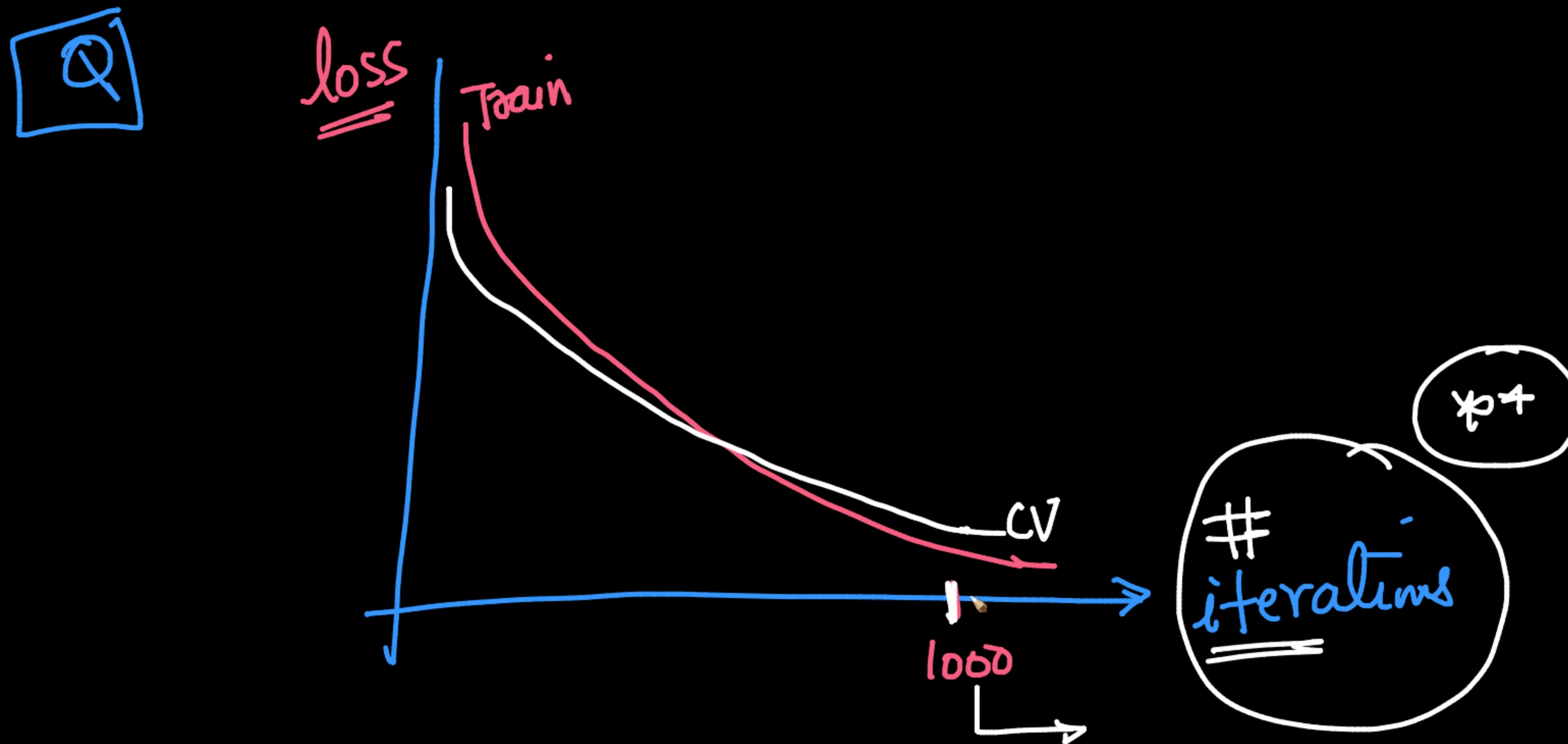
- different sets of features for each tree : feature sampling
- tune the depth
- $M_i \& M_j$ would be more diff if m is smaller

DP
 M_1, M_2, \dots, M_K are different
from one another?
 $k \approx \underline{100's}$
 $d' \leq d$

Logistic Reg → OHE
5-15 categories

$c_1 = 1 = \text{True}$
 $w = 2 \rightarrow$





Regression Using DT



MSE or MAE

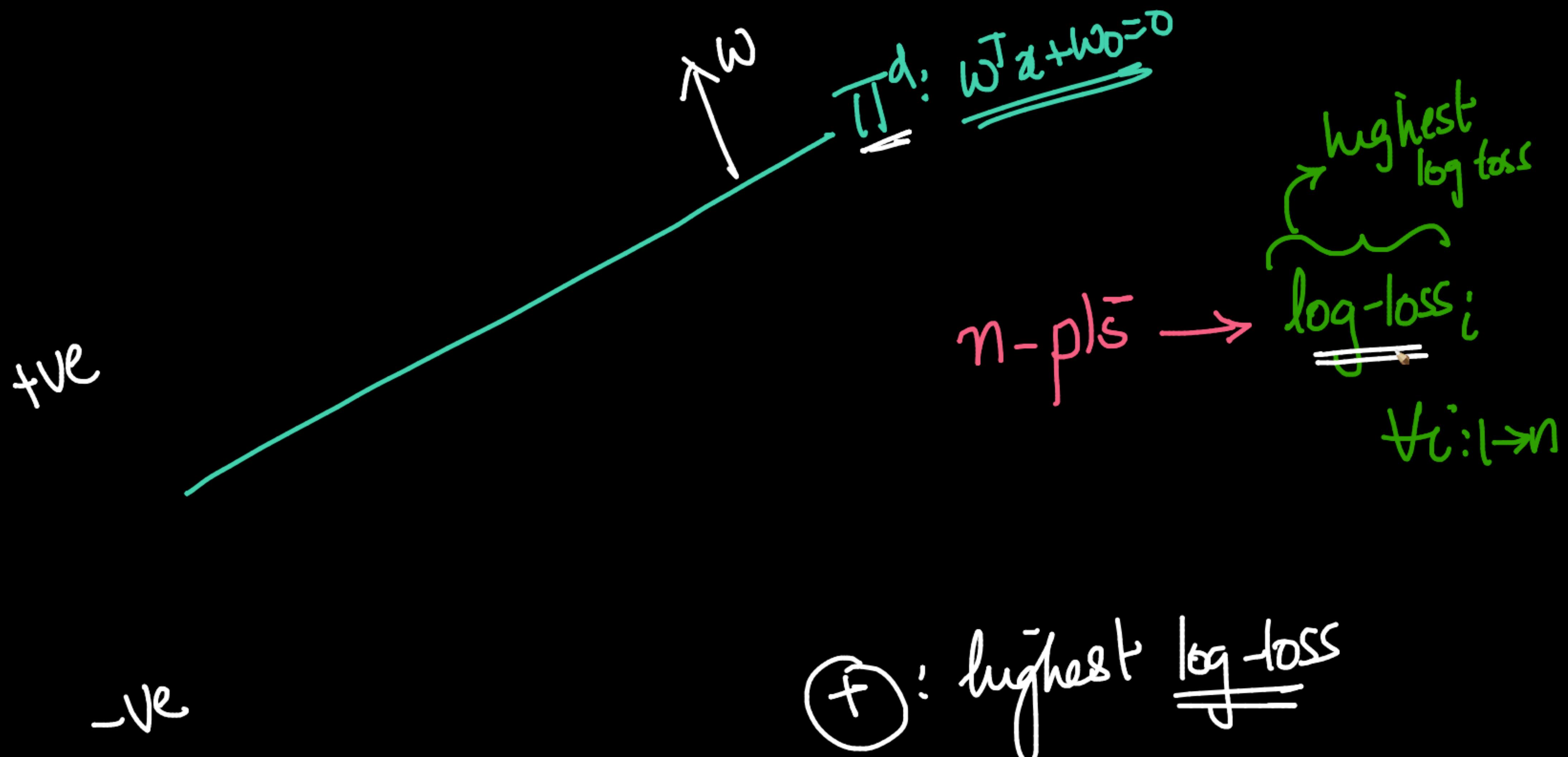


Variable



MAD

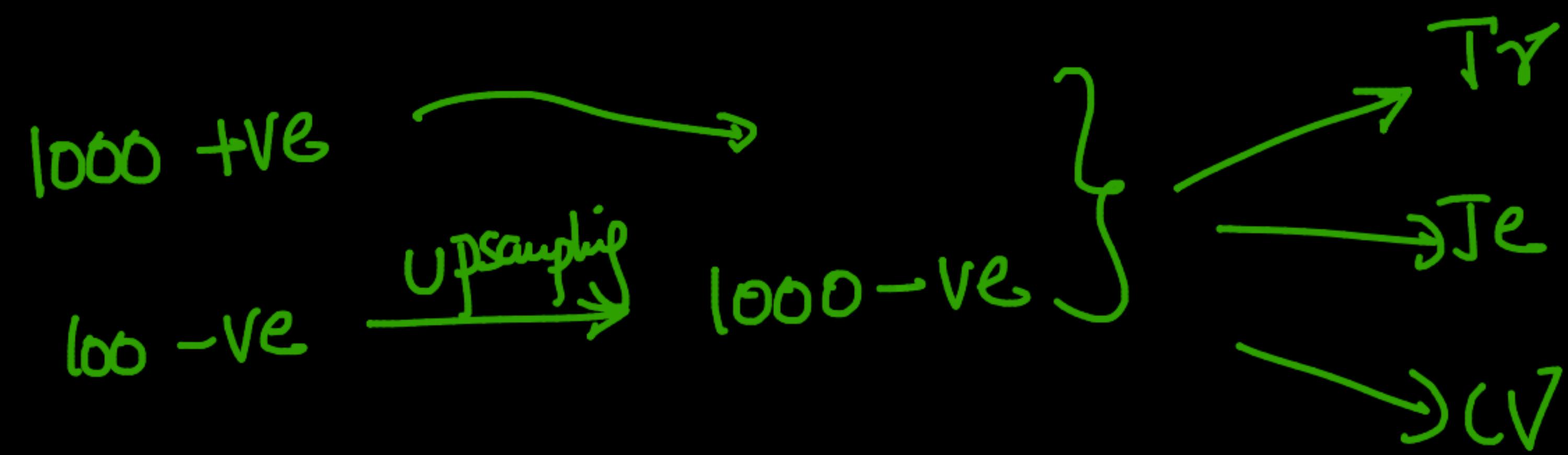




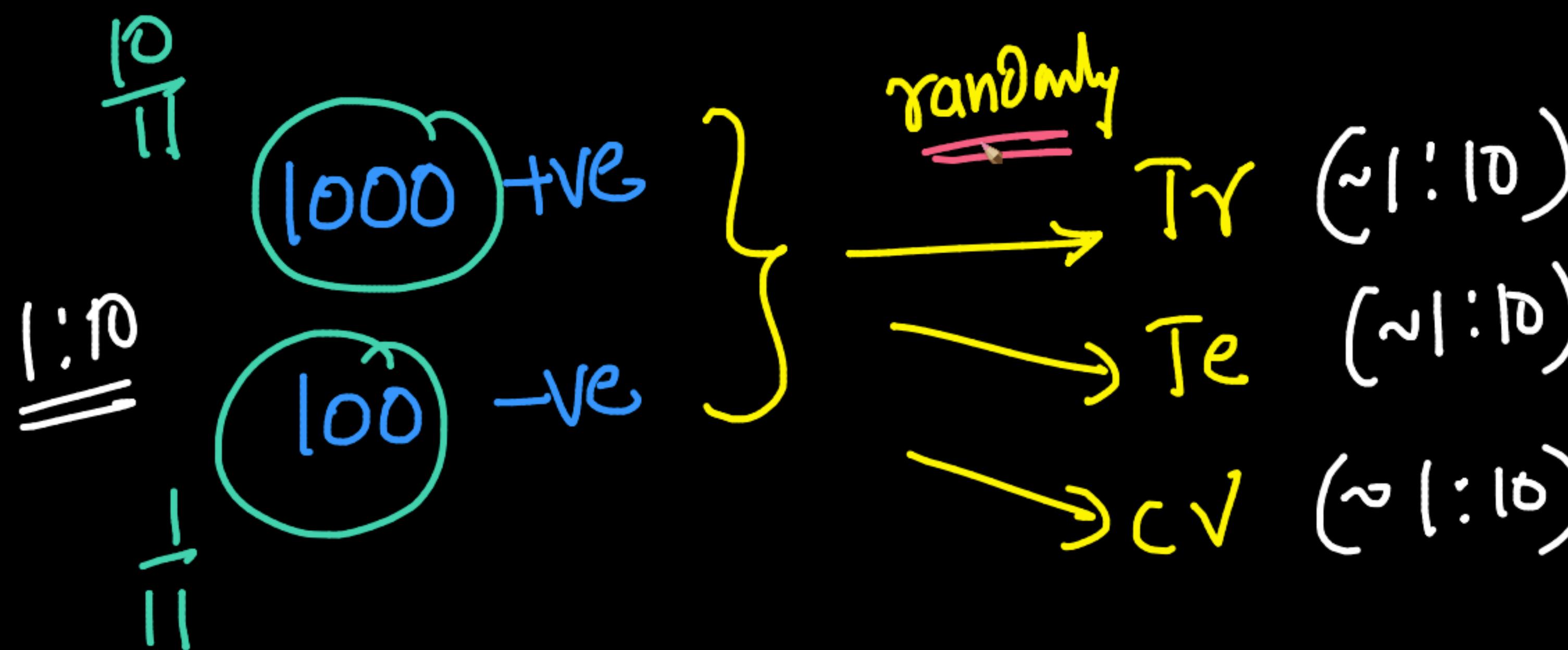
1000 +ve
100 -ve

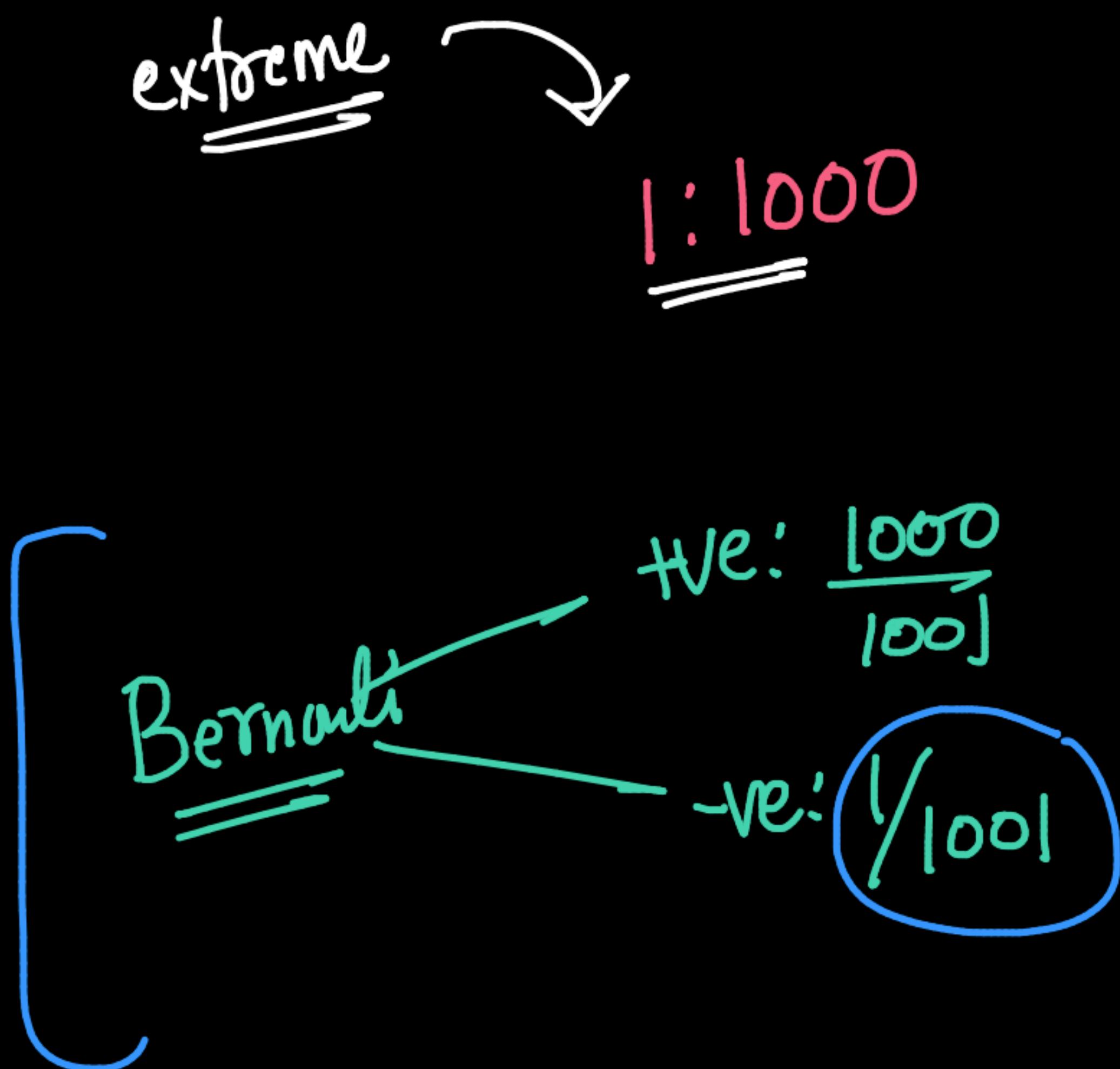
100 +ve (?)
100 -ve

The diagram illustrates a comparison or transformation between two groups of values. The left group consists of 1000 and 100, both labeled '+ve'. The right group consists of 100 and 100, where the top 100 is labeled '+ve' followed by a question mark, and the bottom 100 is labeled '-ve'. Blue arrows point from the first group to the second group, indicating a relationship or mapping between them.



class-weights





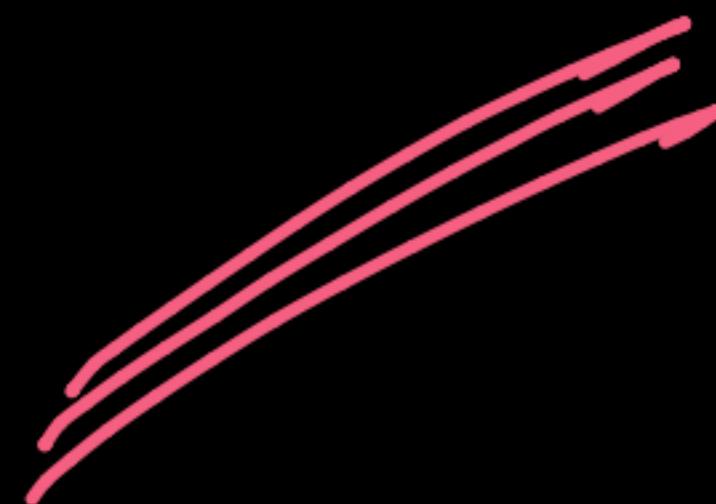
+ve: 10000

-ve: 10

v-few

$$X \sim \text{Binomial}\left(n=500; p = \frac{1}{100}\right)$$

$$P(X=1) =$$



R^2 : adjusted R^2

cover it later



Recursive Partitioning

