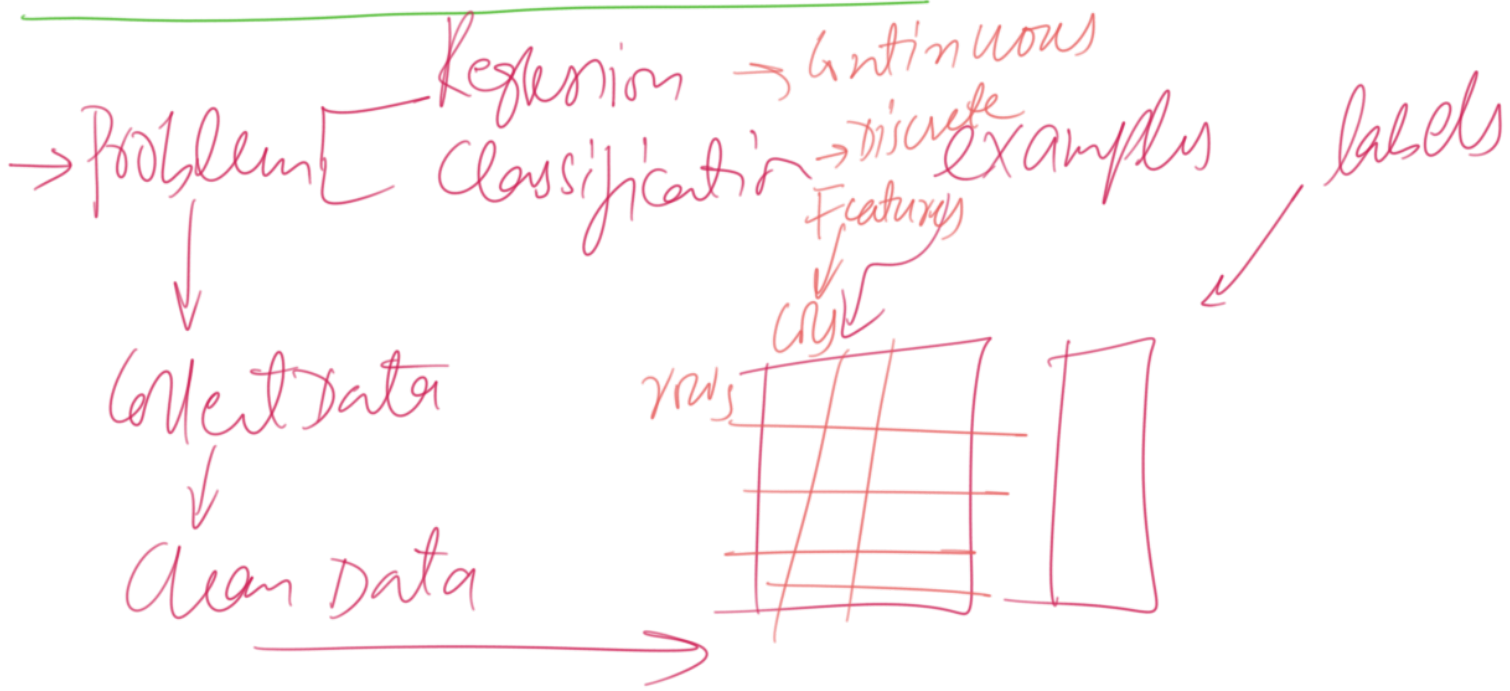
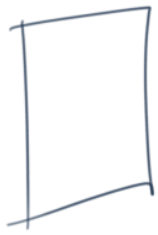


Machine Learning





$m \times n$

↓
example

→ split

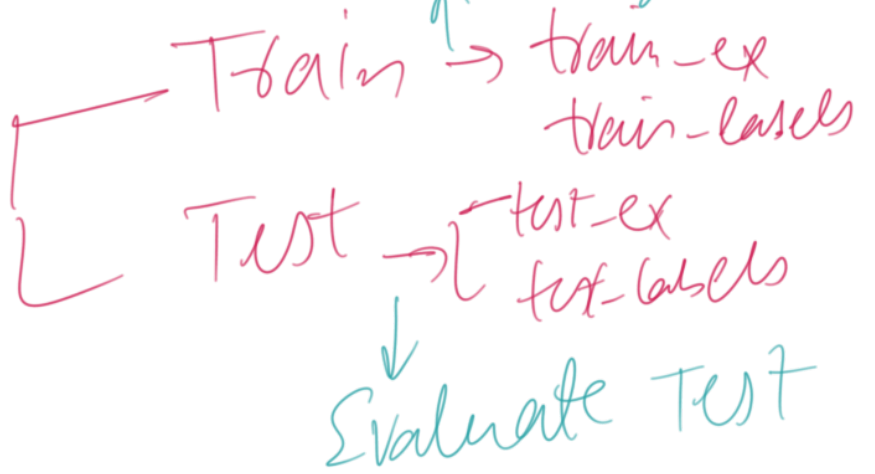
$m \rightarrow$ examples
 $n \rightarrow$ features



$m \times 1$

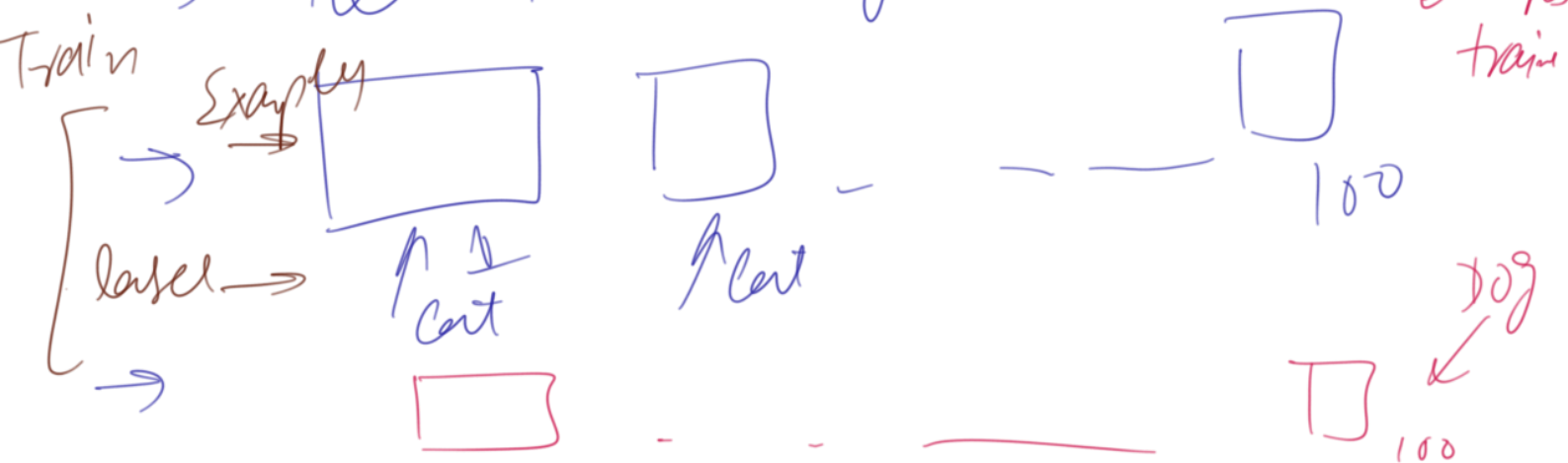
m labels
for
 m examples

Train algorithm



⇒ Goal: To teach a child to identify
Cents & dogs.

→ collect pictures of cents & dogs



Expectation: Child will learn
visual features () to
differentiate b/w cats & dogs.

This is a one time process to
teach the child to identify b/w
Cats & Dogs.

Different Pictures of Cats & Dogs



25



Cats } Test
Dogs }

25



For Evaluation to test how good the child is learning

You use this test set to measure
Child Learning Capability (via accuracy)

30% Correct Answers

test metric

Evaluation
metric

↓

↓

100%

Test set you will

stop

(Unseen set)

→ Q: What happens if I use
the same training set for evaluation?

↓
Accuracy will be 100%.

overfitting

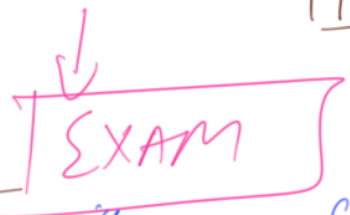
[Rote learning]

→ Class Scenario

Teach some concepts in the class
and we do some exercises on this.
Train

overfitting love them

Are we learning something? reproduce & get A



Test Same exercise

Hate Test Similar exercises

- Do Good
- Do Bad

Train set.

Underfitting

Overfitting

Test set
↓
Similar

Identical

Remembering
by heart

Generalization

Excellent

Good

Good

Bad

Bad

late learning

	Identical Test	Similar Test	
overfitting	Good	Bad	He is just memorizing without learning concepts
underfitting	Bad	Bad	
Generalization	Good	Good	not exploiting his capabilities ← Want

Class / Real life Scenario

Good	Bad
Bad	Bad
Good	Good

overfitting

Exploiting his capability

Learning Slump

& Student capacity are both low

→ Right learning Capacity
→ Right learning setup

underfitting

Right setup

Right algorithm

overfit

Simple algo
for complex

complex algo
for simple task

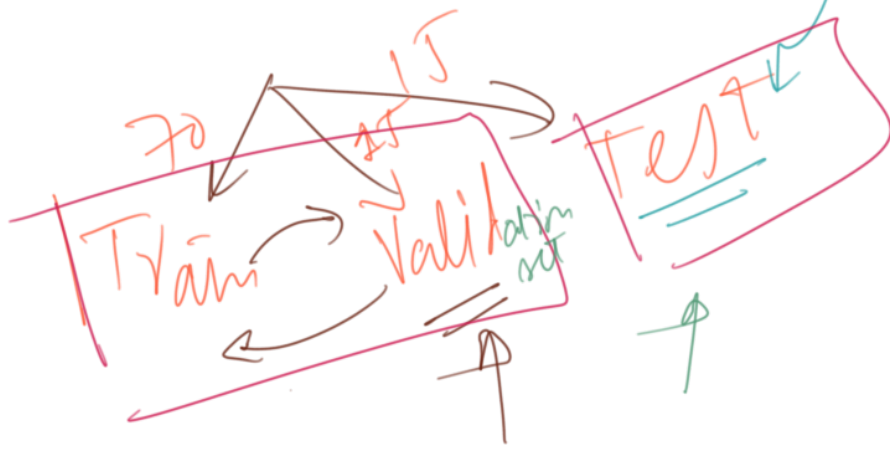
- Generalization Good on training & test
- Underfitting Bad on both
- Overfitting Good on training / Bad on test

Large Dataset

> 10k examples



We want to ~~test~~ use only once



Real life

Holy set
Reflection of our
system's performance in real life

Small Dataset $\leq 10K$
150 Examples

There is not enough data to split into train/valid/test set

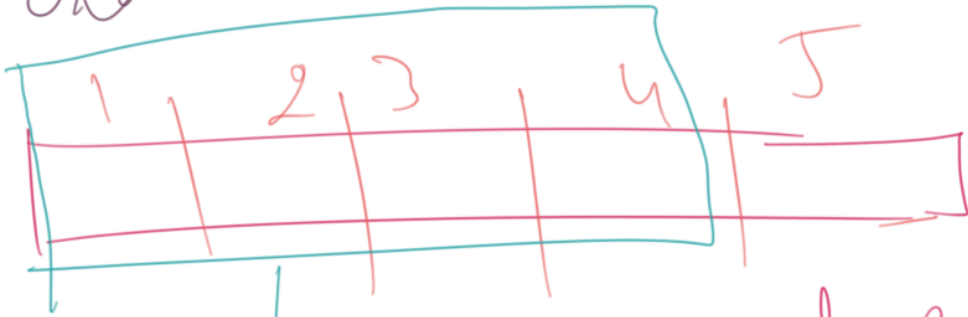
→ How to ^{Setup learning,} do train test split

→ large dataset

→ Small dataset

K-Fold Cross Validation

$\frac{4}{5}$



① Split data into 5 equal parts. ¹⁰⁰ Examples

Train $\leftarrow F[1:4]$

$F[5] \rightarrow$ Test:



- I₁: F1 — F4 → Train, F5-Test → Accuracy
- I₂: F1, F2, F3, F5 → Train, F4-Test →
- I₃: F1, F2, F4, F5 → Train; F3-Test →
- I₄: F1, F3, F4, F5 → Train, F2-Test →
- I₅: F2 — F5 → Train; F1-Test →

* Learning setup

↳ large
↳ small

Check for generalization; or overfitting
& underfitting

Evaluation metrics for classification
Algorithm

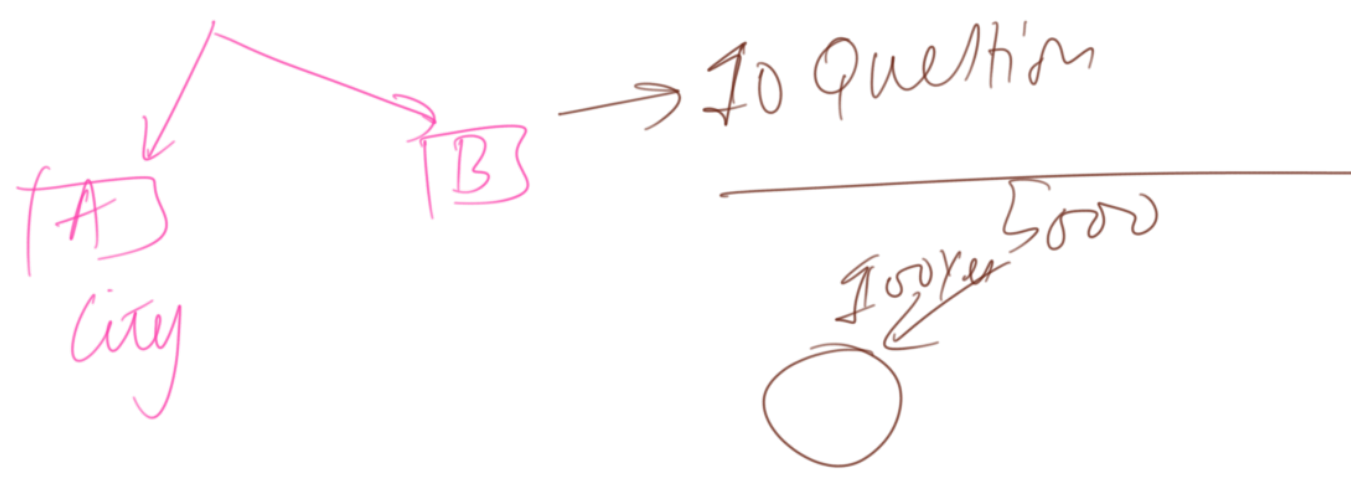
K-Nearest Neighbours

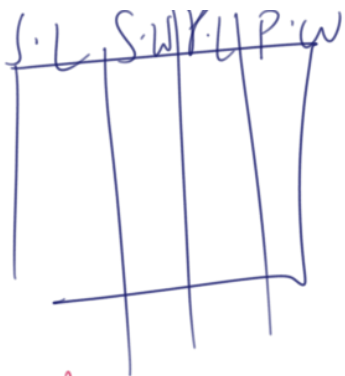
Decision Trees

Bagging | Random
Forest

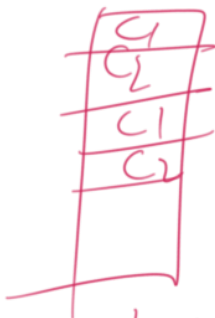
What a decision tree :-

Guess a city Game :-





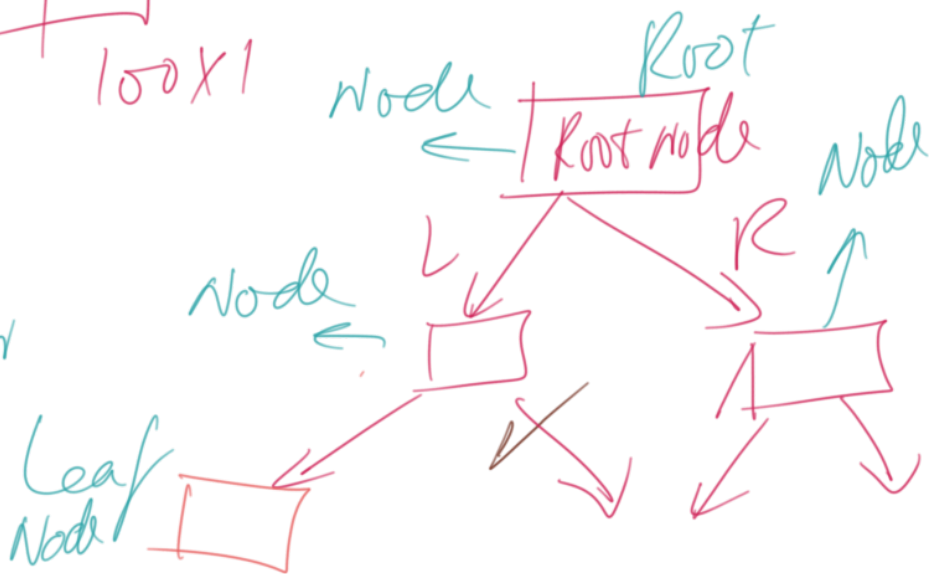
100 x 4



100 x 1

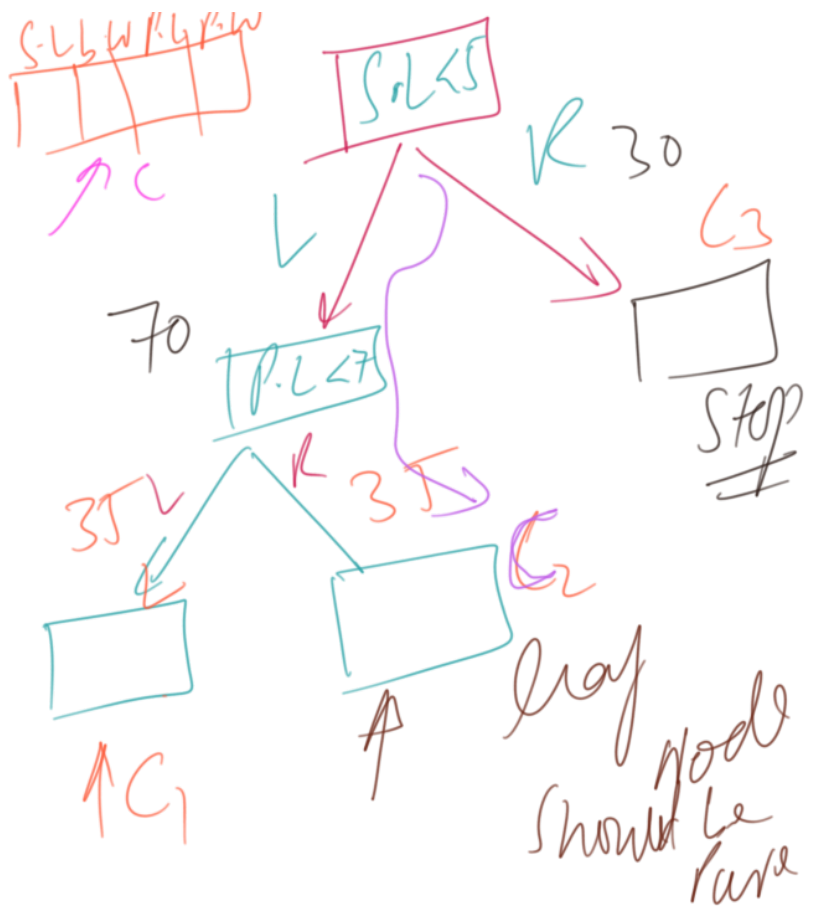
It will build a tree

Goal: Given features
Predict type of flower





Pure means belonging to same class.



How to build Decision tree

Entropy is an agent of Entropy:-
Entropy:- measure of impurity in a set

Impure

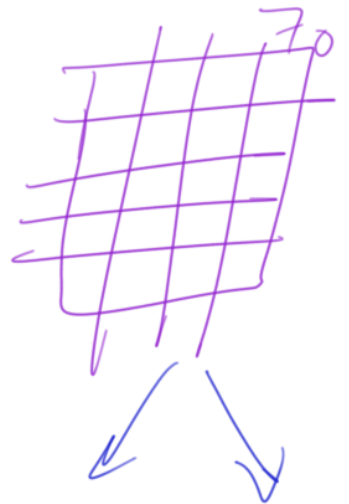
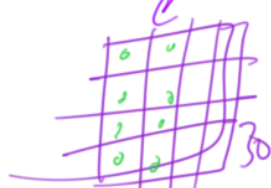
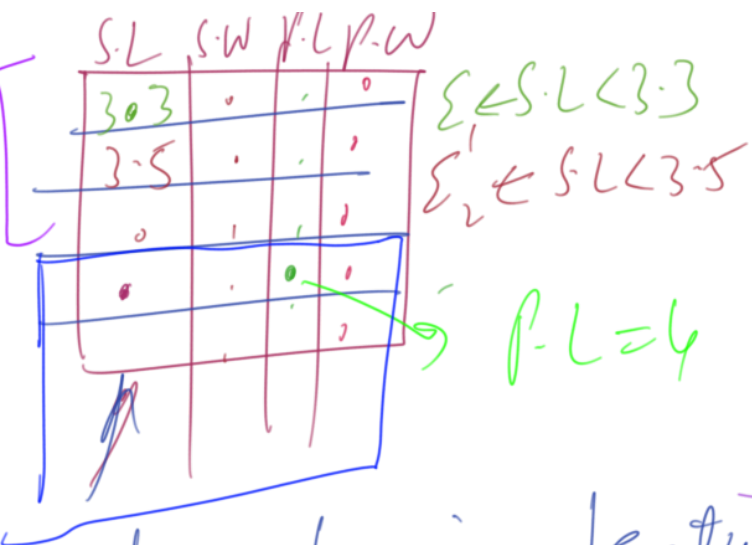


Pure

= 0

$$E = -\sum p_i \log p_i$$

$$E(S_1) > E(S_2) > E(S_3)$$



for f in features:
 for val in f:

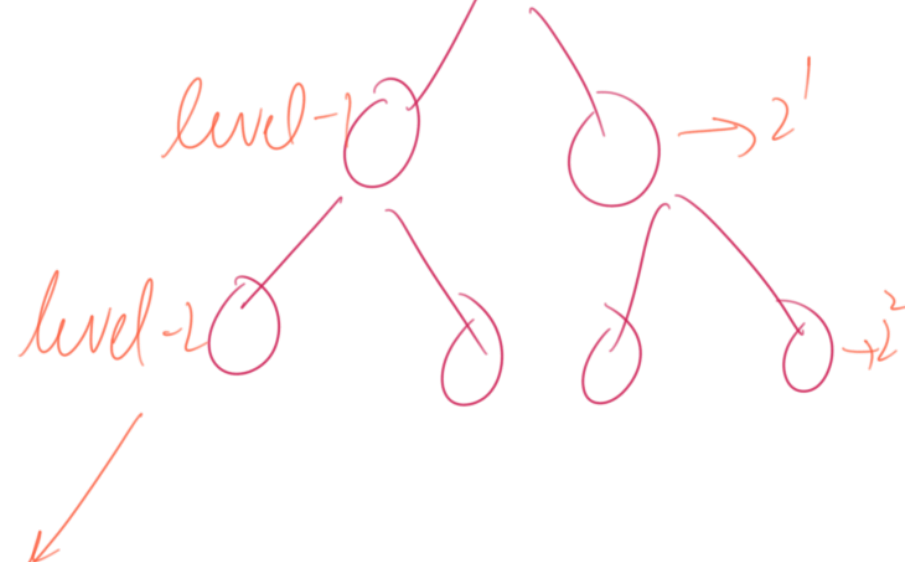
→ Depth of a tree

depth-2 has 4 nodes

depth-3 has 8 nodes

depth 5 → 32

level 0



level-3 → 2³ → 8 nodes

How to use decision tree in scikit

from sklearn.tree import DecisionTreeClassifier

- fit
- predict

height

$$h=3 \rightarrow 2^3=8$$

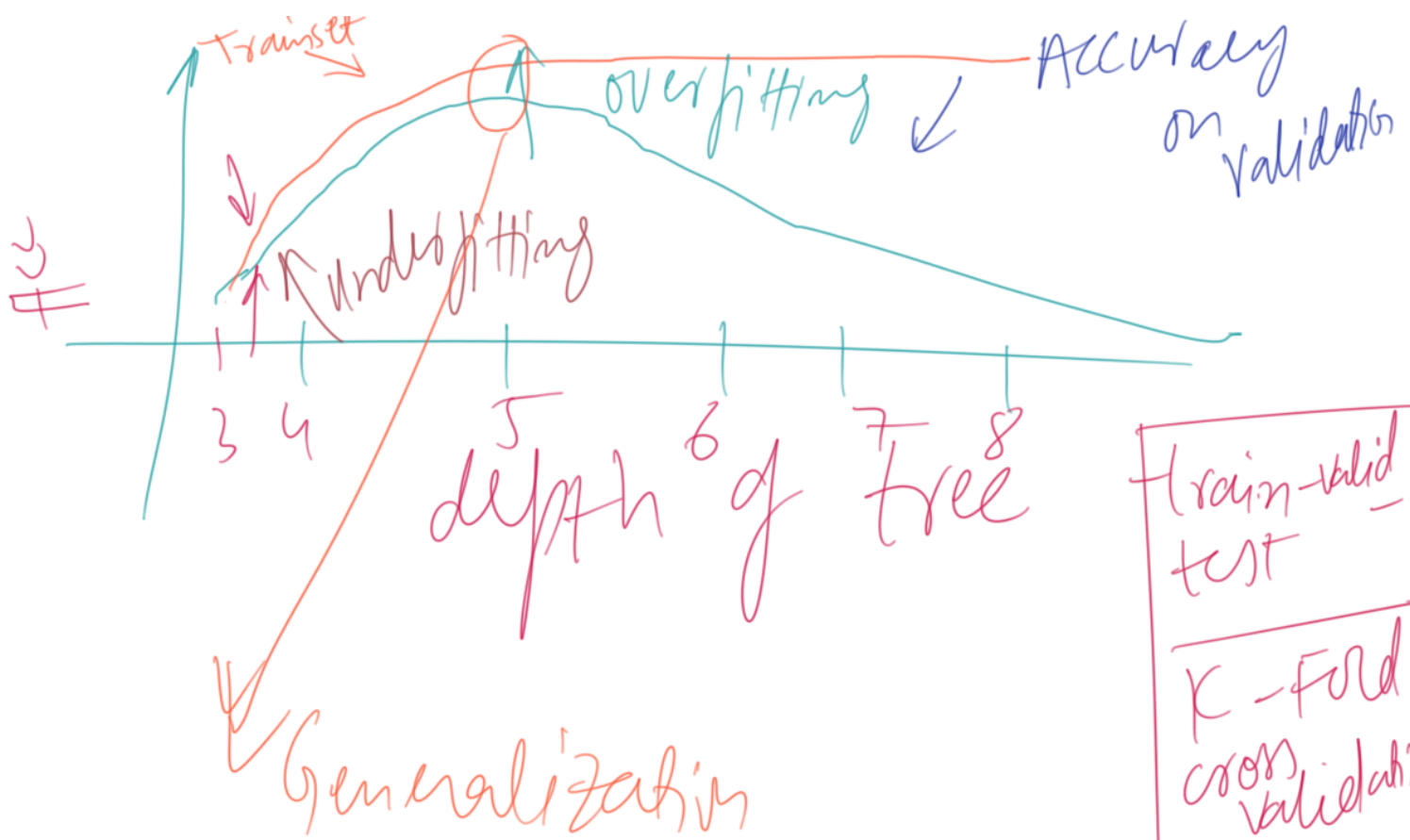
depth / Height & overfitting

→ In real life we don't build too deep trees.

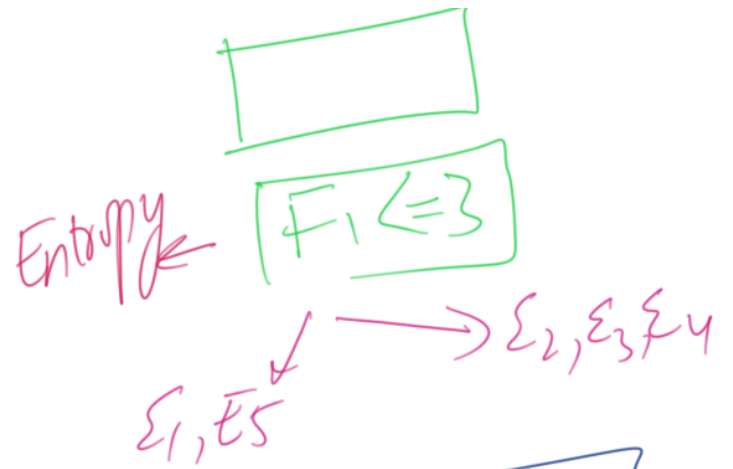
How to choose tree depth?

Job d in $[4, 5, 6, 7, 8, 9, 10]$:

~~train-pred = tree~~
~~acc-train =~~
accuracy-score (train labels, train-pred)
tree = DecisionTreeClassifier (max depth = d)
 \neq tree-fit (trainex, train-labels)
pred = tree.predict(testex)
acc = accuracy-score (testex, test-labels)



	F_1	F_2	
ϵ_1	3	-1	+
ϵ_2	4	-2	+
ϵ_3	5	-3	-
ϵ_4	4	1	-
ϵ_5	3	2	



[" F_1 " : (3, 7), " F_1 " : (4, 5), " F_1 " : (4, 5)]

→ Entropy



↑ 2 labels

$$P(+)=3/5$$

$$P(-)=2/5$$

$$E(S) = \sum_{\text{labels}} P_{\text{label}} \log P_{\text{label}}$$



↑ 2 labels



↑ 2 labels

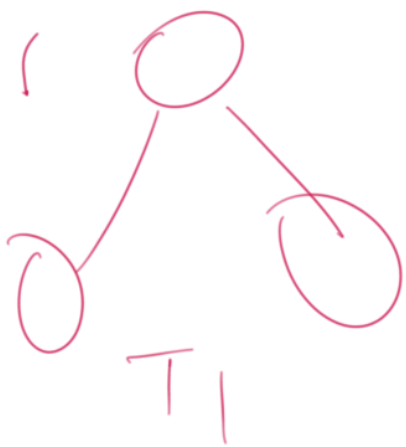
$$\rightarrow = 3/5 \log 3/5 + 2/5 \log 2/5$$



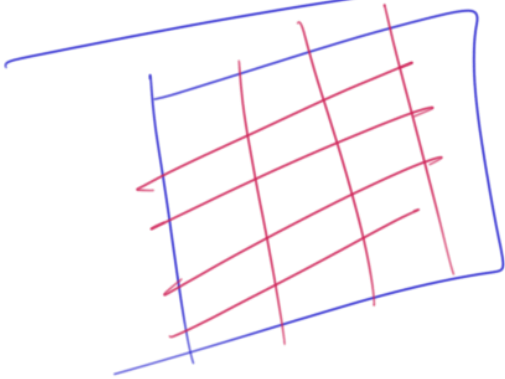
~~Random Forest~~ ← many trees
← Ensemble

"A committee of trees is
always better than a single tree".

Instead of a single tree
we build many trees that take
the majority vote.

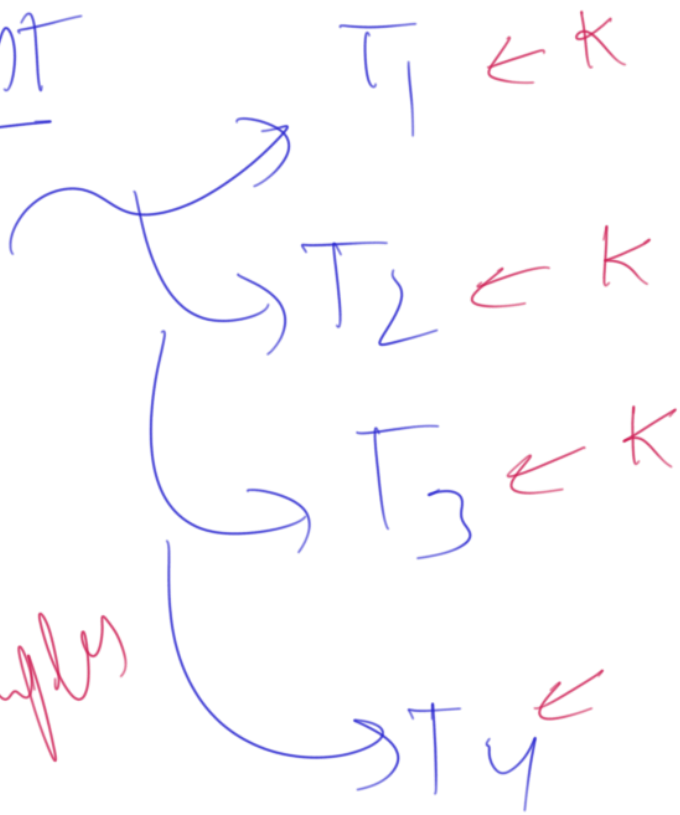


Random Forest



$X \rightarrow m \times n$

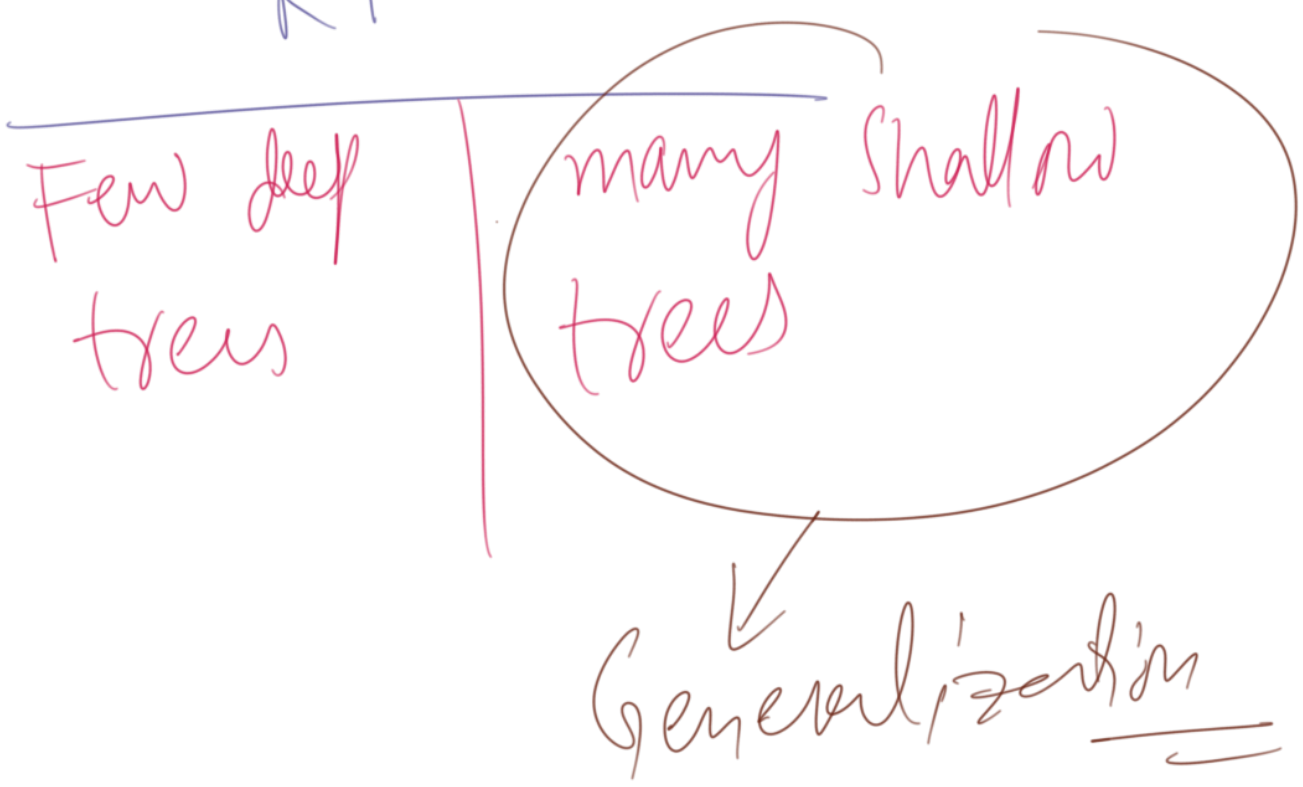
m examples



→ BIAS → Underfitting

→ VARIANCE → overfitting

RF



Few deep trees

many shallow trees

Generalization

RF

