

Last class (Sep 04)

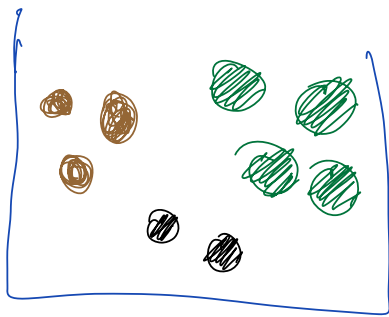
- 1) KNN with categorical features
- 2) Different distance metrics
- 3) LSH for KNN
- 4) Missing data with KNN — imputation
- 5) Employee Attrition Problem Statement
- 6) Decision Tree Intuition
- 7) How to split the nodes
- 8) Purity and Impurity of the nodes

Today's class

- 0) Quizzes ✓
- 1) Recap
- 2) Entropy
- 3) Information Gain
- 4) How to split a node — categorical & Numerical
- 5) Issue with Entropy
- 6) Gini Impurity
- 7) Feature Scaling
- 8) Overfit vs Underfit
- 9) Hyper-parameter Tuning (Time Permits)

numerical F_1	F_2 categorical
0.09	GREEN
0.1	RED
-1.5	BLUE
11.7	YELLOW

All distinct values



$$E = - \left[\frac{1}{3} \log_2 \left(\frac{1}{3} \right) + \frac{4}{9} \log_2 \left(\frac{4}{9} \right) + \frac{2}{9} \log_2 \left(\frac{2}{9} \right) \right]$$

$$P(\text{brown}) = \frac{1}{3}$$

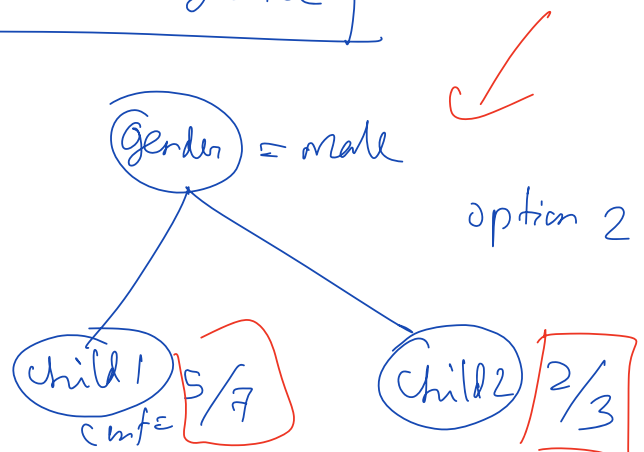
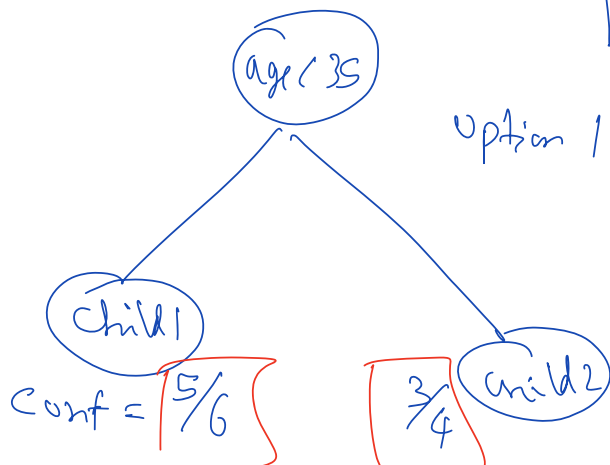
$$P(\text{green}) = \frac{4}{9}$$

$$P(\text{black}) = \frac{2}{9}$$

$$E = - \left[P_0 \log(P_0) + P_1 \log(P_1) \right]$$

\downarrow \downarrow
 $\log\text{-loss} = - \left[\gamma_0 \log(P_0) + \gamma_1 \log(P_1) \right]$
 \downarrow
 ground truths

KL Divergence



\rightarrow wt. Confidenc derived
 proportion of majority class = confidence
 of predictions

$$H[\text{parent}] = 0.97$$

$$H[C1] = 0.65$$

$$H[C2] = 0.81$$

$$C_1 = 60 \text{ samples}$$

$$P(C_1) = \frac{60}{100} = 0.6 \checkmark$$

$$C_2 = 40 \text{ samples}$$

$$P(C_2) = \frac{40}{100} = 0.4 \checkmark$$

Weighted \mathbb{E} for children = $WH[C]$

$$= P(C_1) \times H[C1] + P(C_2) \times H[C2]$$

$$= 0.6 \times 0.65 + 0.4 \times 0.81$$

$$= 0.714$$

$$\boxed{\text{Information Gain}} (IG) = \boxed{H[\text{parent}]} - WH[C]$$

$$= 0.97 - 0.714$$

$$= \boxed{0.256} \checkmark$$

Gender condition

$$H[P] = 0.97$$

$$P[C_1] = 7/10$$

$$H[C_1] = 0.86$$

$$P[C_2] = 3/10$$

$$H[C_2] = 0.91$$

$$WH[C] \equiv P[C_1] \times H[C_1]$$

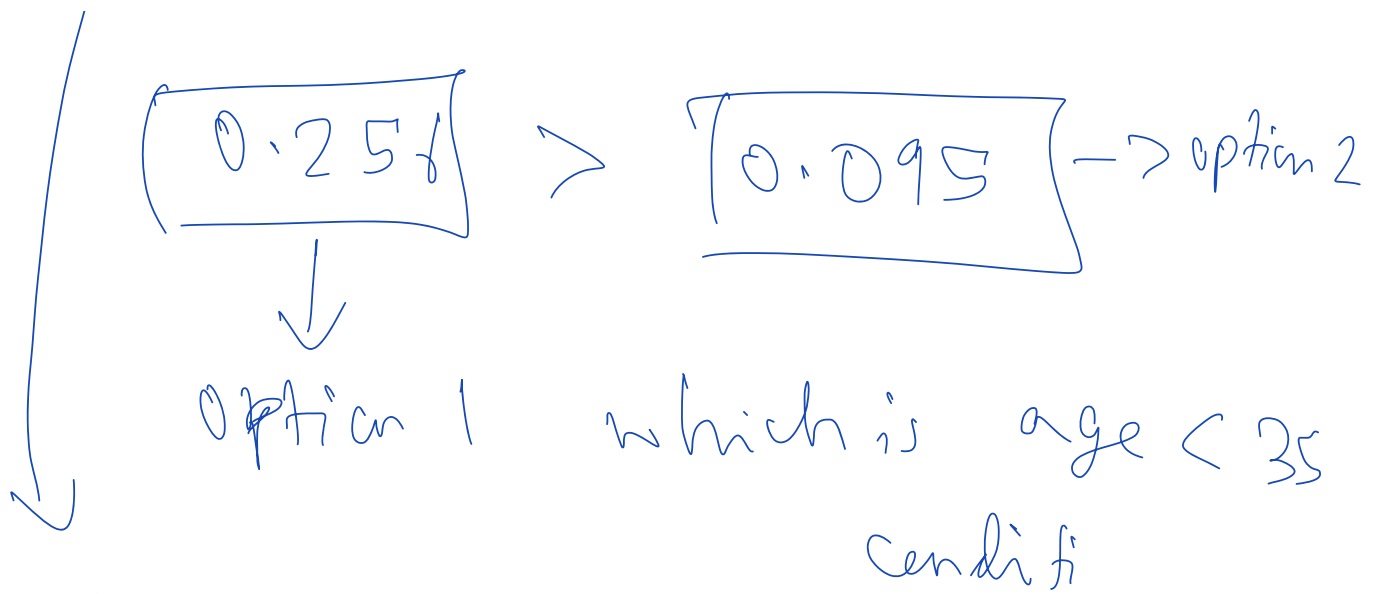
$$+ P[C_2] \times H[C_2]$$

$$= \underline{7/10} \times 0.86 + \underline{3/10} \times 0.91$$

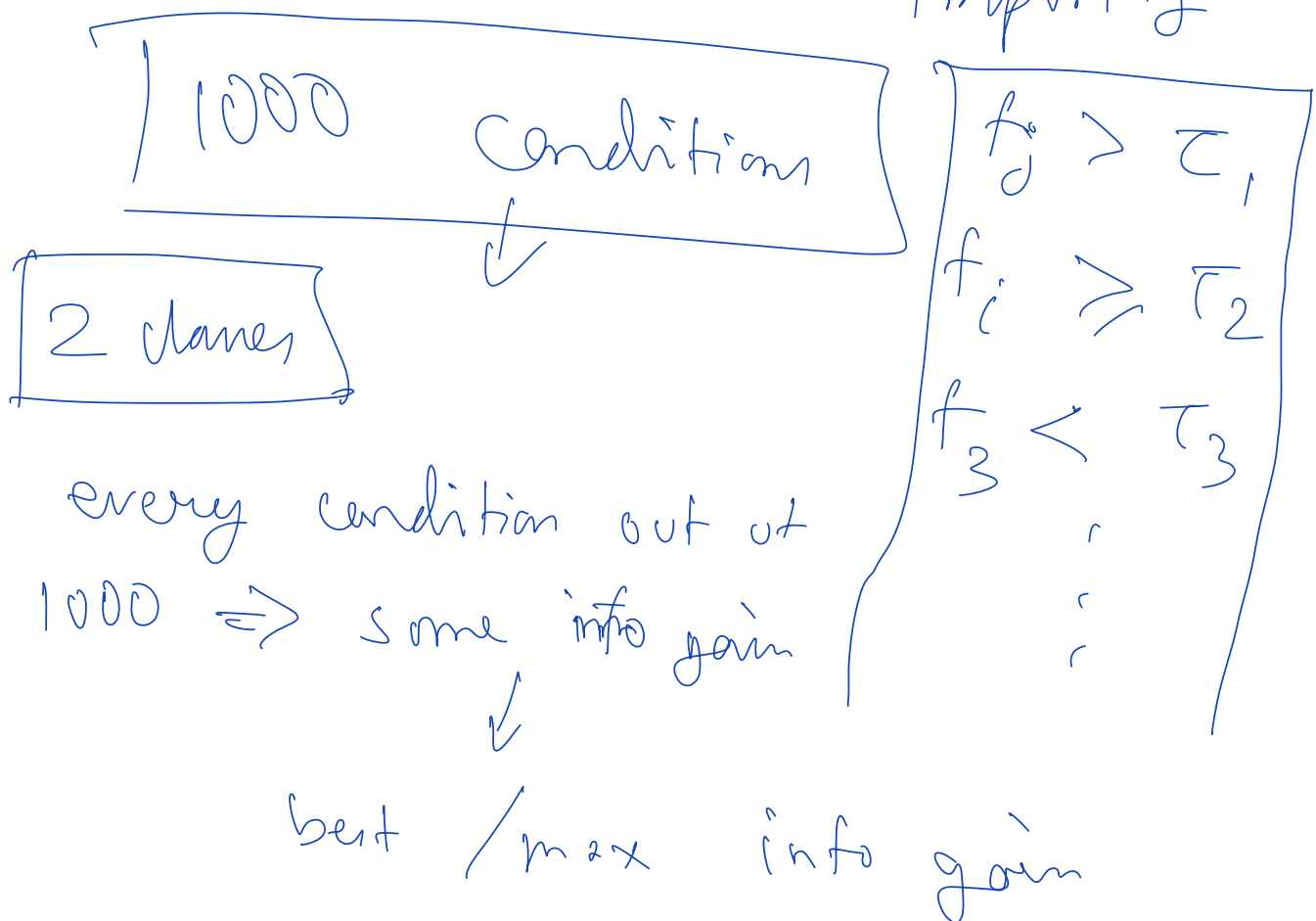
$$= 0.875$$

$$IG = 0.97 - 0.875$$

$$= \boxed{0.095} \quad \checkmark$$



Information Gain \propto Reduction
in impurity



$$O(1000) \rightarrow O(n)$$

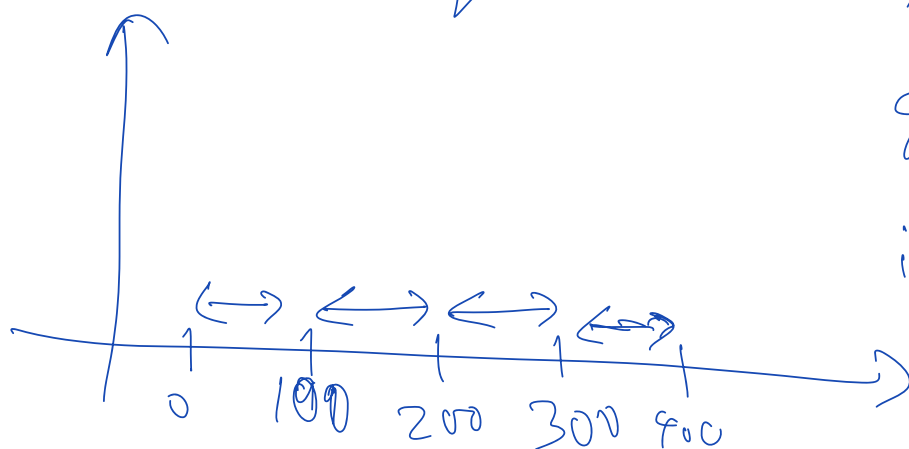
number of conditions

age = [1 to 1000]

range

99, 999, 70, ...

grouped
values
into bins



1000 conditions

10 conditions

1, 2, 3, 4, ...

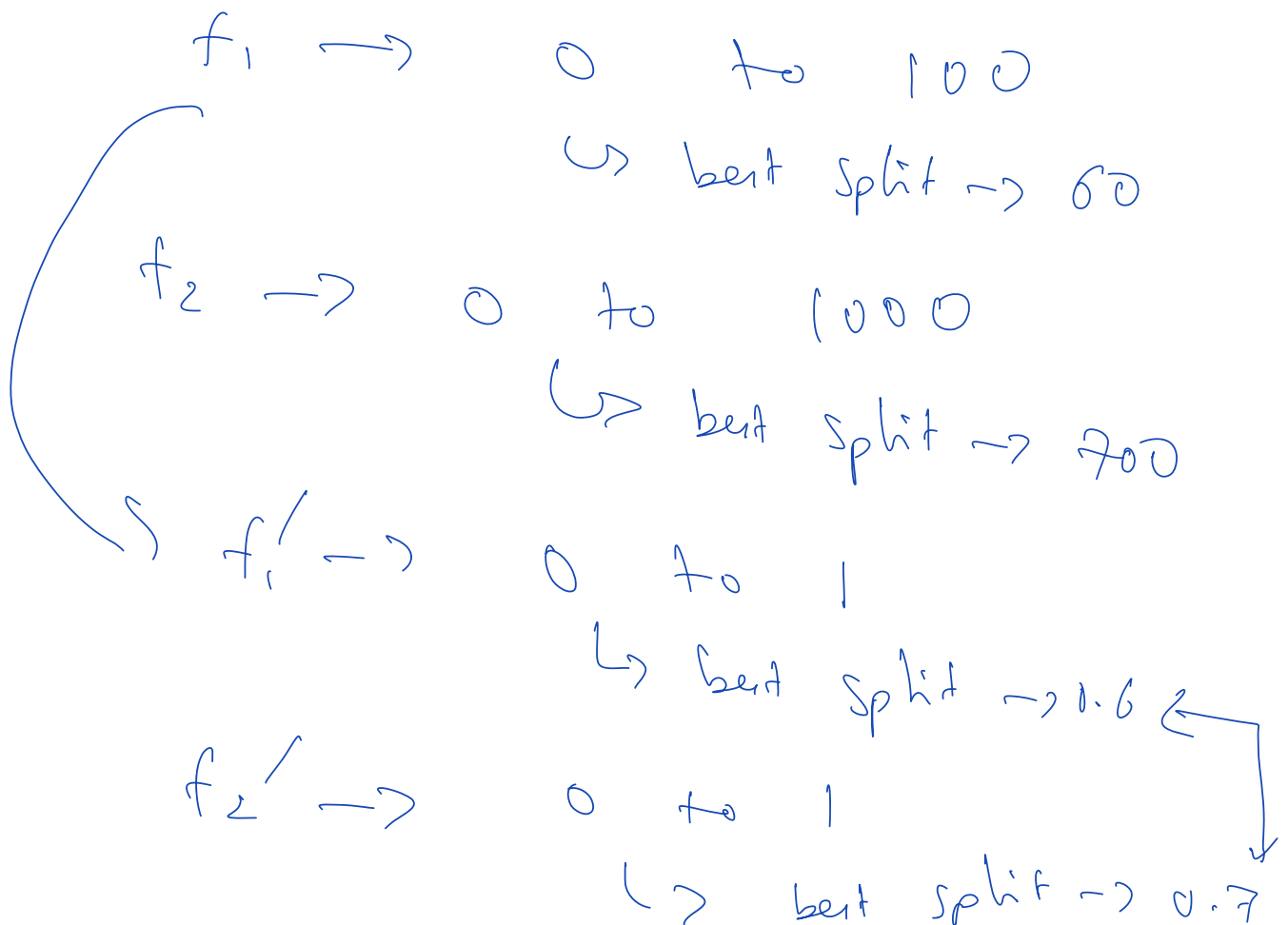
0 - 100, 100 - 200, 200 - 300, ...

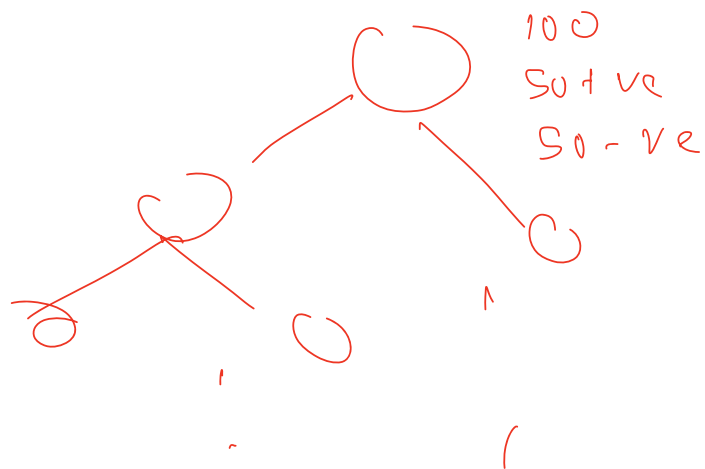
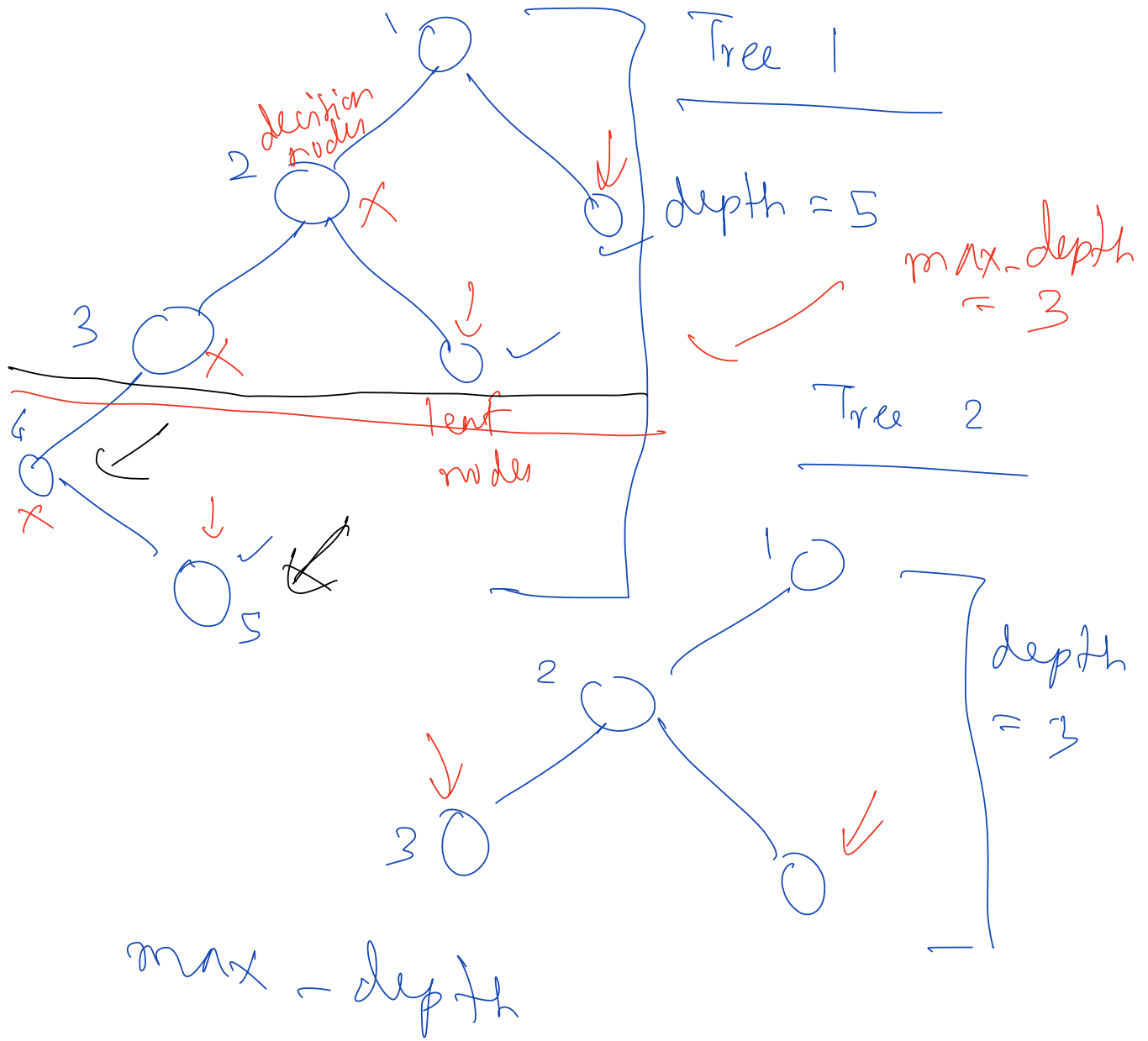
$$1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right]$$


$$1 - [0.5^2 + 0.5^2]$$


$$1 - [0.25 + 0.25]$$

$$1 - [0.5] = 0.5$$



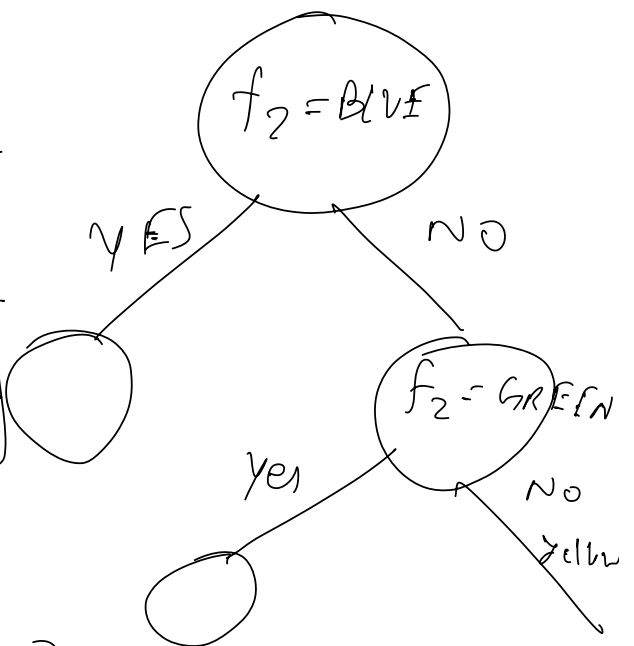


$n=1$
+ve 

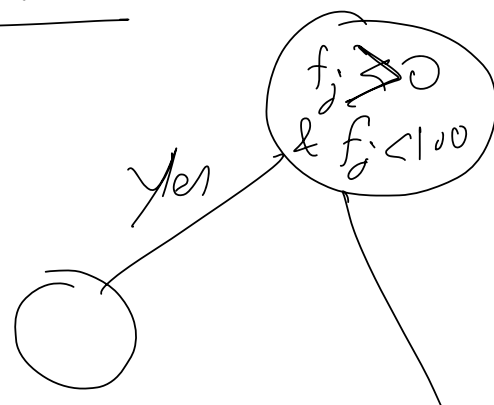
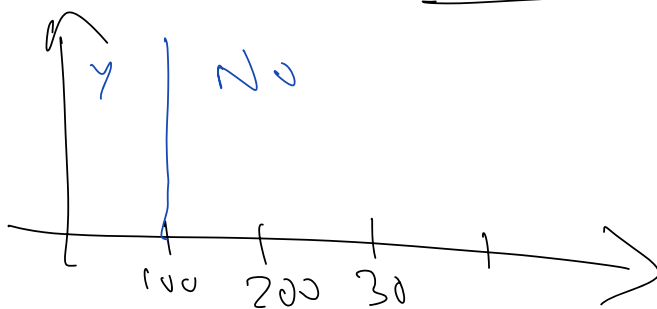
$n=2$
-ve 

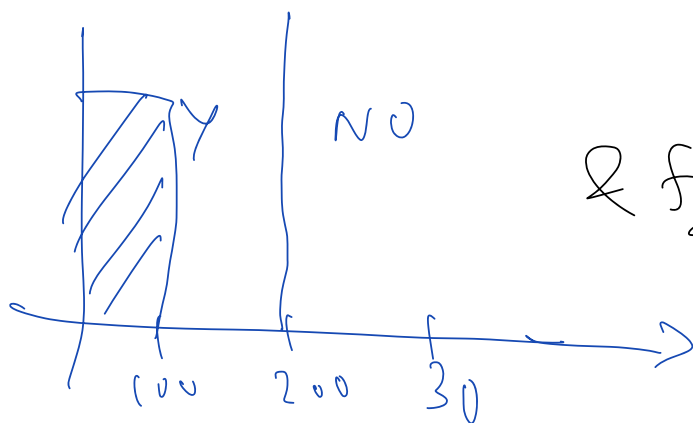
increase max-depth \rightarrow overfit
decrease " \rightarrow underfit

f_1	f_2
	BLUE
	YELLOW
	GREEN



numerical:
1 to 1000





$$f_j > 100$$

$$\& f_x < 200$$

Yes

$$200 - 1000$$

