

## Previous Class - July 1

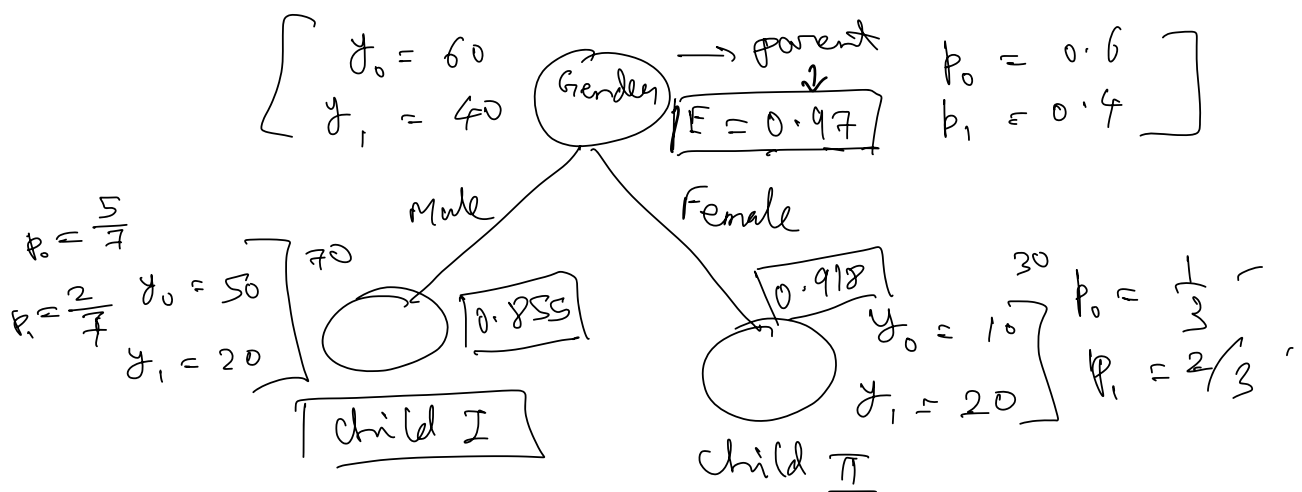
- 1) Recap- Quiz
- 2) Employee Attrition Dataset - graph
- 3) Purity of Nodes & Entropy
- 4) Plot for entropy
- 5) Weighted Entropy
- 6) Gini Impurity
- 7) Comparing Gini Impurity with Entropy
- 8) Code walkthrough (Time Permits)

## Today's agenda

- 1) Recap
- 2) Splitting on numerical features. ✓
- 3) Overfit vs Underfit. ✓
- 4) Hyperparameters ✓
- 5) Visualizing a DT. ✓
- 6) Impact of outliers → decrease the depth
- 7) Feature Scaling ✓

- 8) Encoding of categorical features
- 9) DT for high-dim data
- 10) Data Imbalance
- 11) Feature Importance ✓
- 12) DT Regression

### Recap

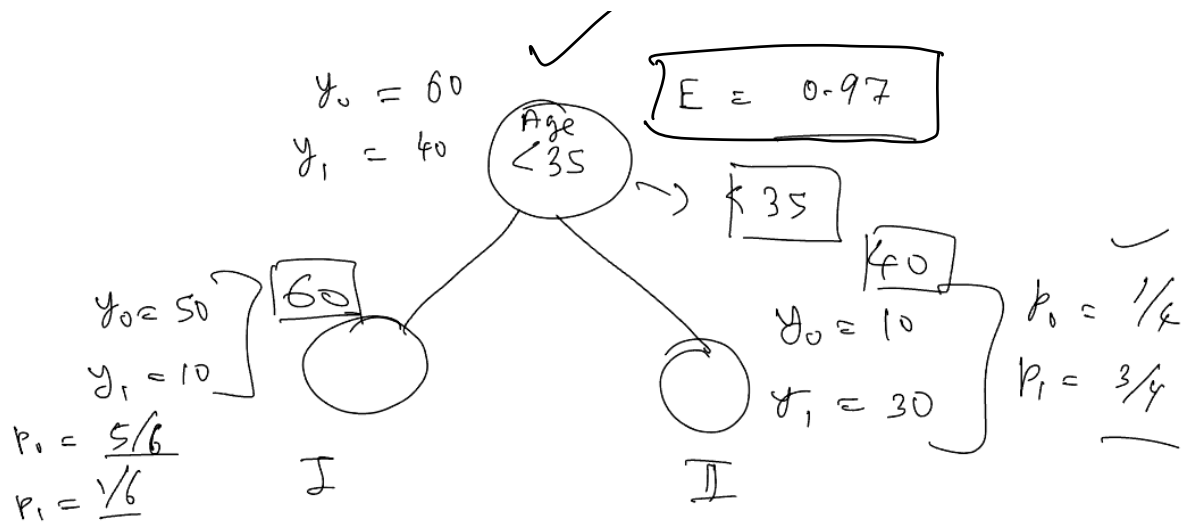


$$WE = \underline{0.874}$$

$$IG = 0.97 - 0.874 = \underline{0.096}$$



$$0.2566 > 0.096$$



$$WE = 0.7144$$

$$IG = 0.97 - 0.7144 = 0.2566$$

$$f_1 = ['A', 'B', 'C', 'D']$$

$$f_2 = ['E', 'F', 'G']$$

$$f_1 = 'A'$$

Yes No

$$f_1 = 'B'$$

Yes No

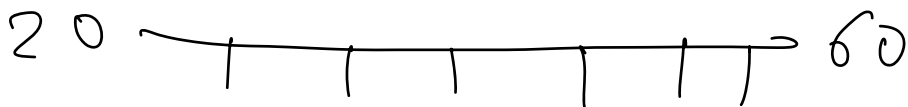
$$f_2 = 'E'$$

Yes No

$$4 + 3 = 7$$

$f_1$	$f_2$	$y$
0.1	10.9	
0.3	11.2	
1.6	-5.6	
1.7	17.1	

↳ unique values of every feature



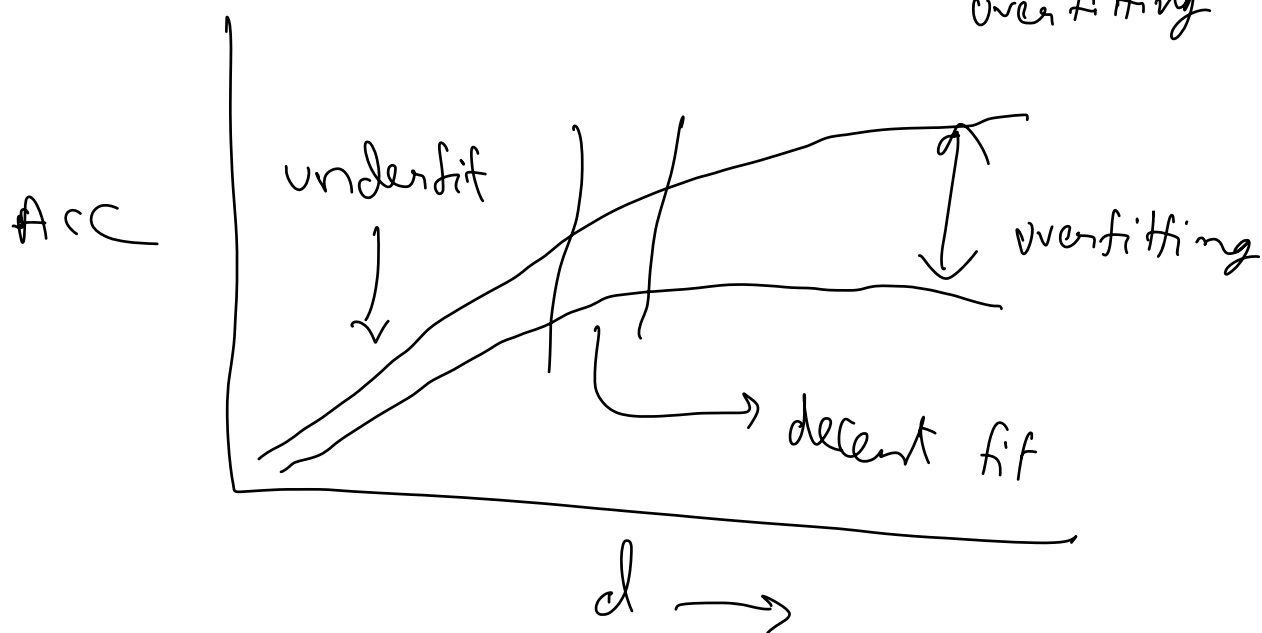
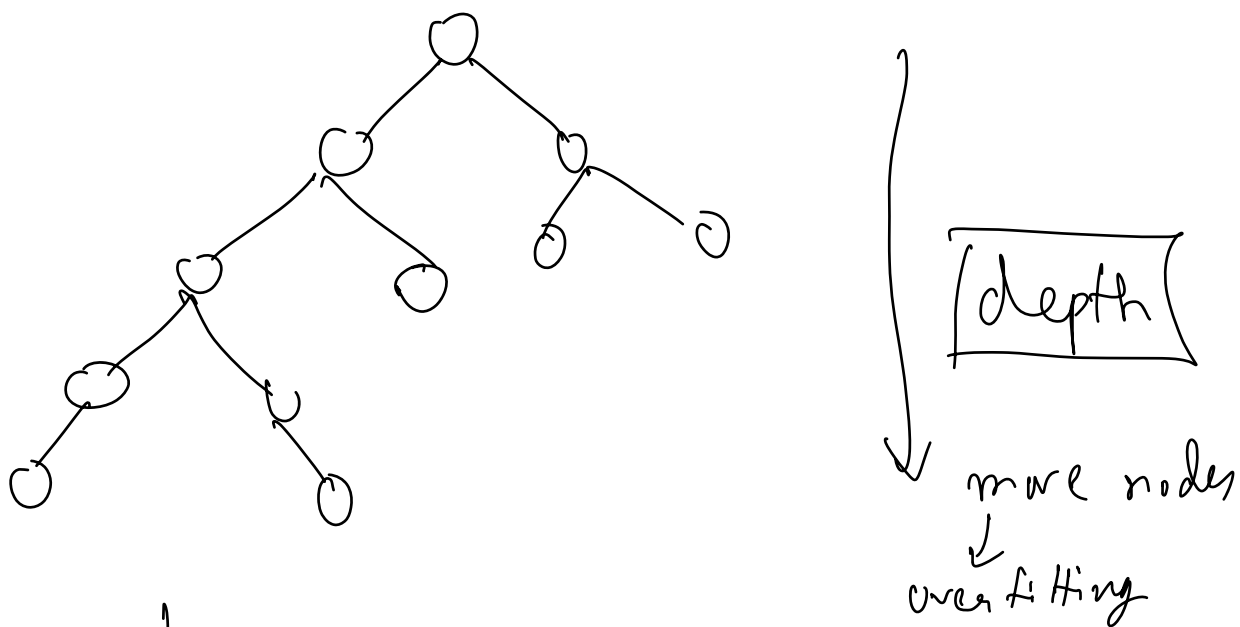
$$\frac{40}{10} = 4$$

20, 24, 28, 32, 36, ..., 60

Underfit vs overfit

More nodes  $\rightarrow$  Overfitting

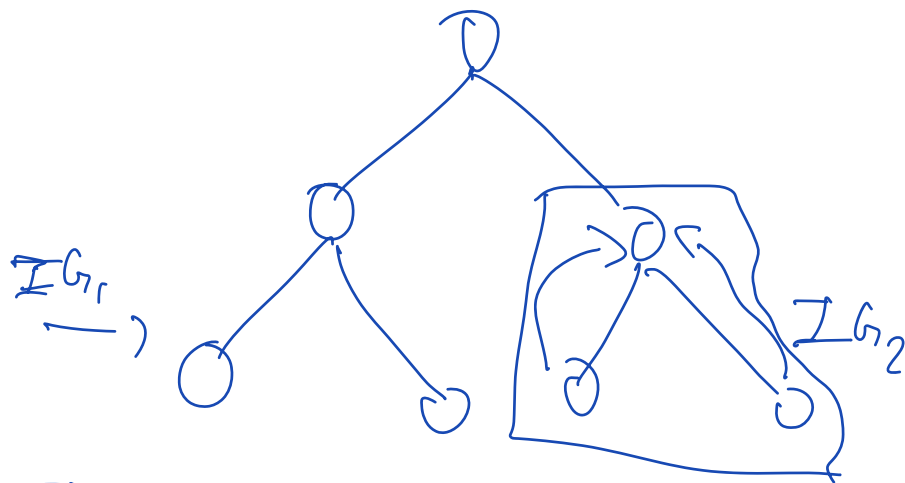
less nodes  $\rightarrow$



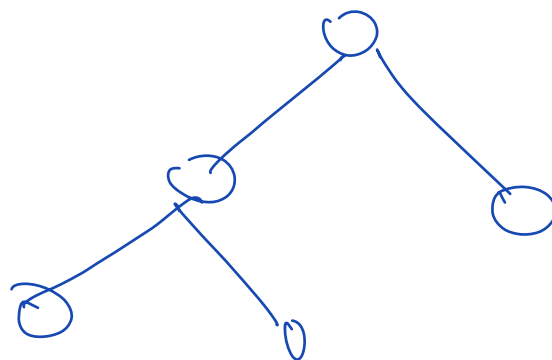
$$- [ p_0 \log_2(p_0) + p_1 \log_2(p_1) ]$$

$$- [ p \log p + (1-p) \log (1-p) ]$$

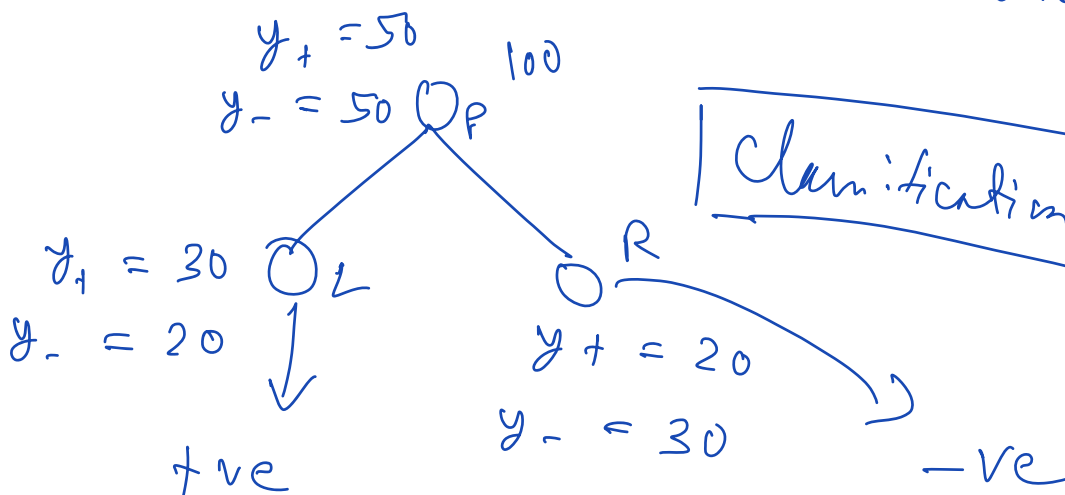
$$- [ w_1 p \log p + w_2 (1-p) \log (1-p) ]$$



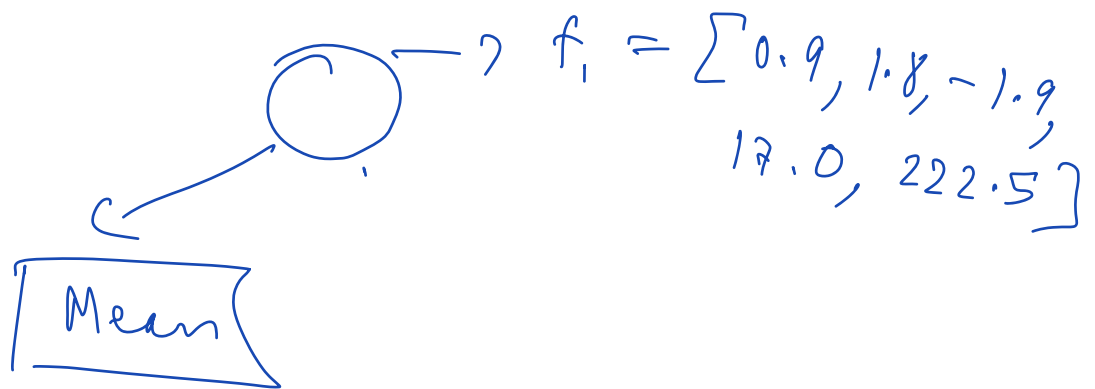
$IG_2$  is very small



Pruning  
↓  
removing  
of  
unnecessary  
branches



## Regression



$$- [p_1 \log(p_1) + p_2 \log(p_2)]$$

$\hookrightarrow 2$  classes

MSE

## Feature imp

$f_1, f_2, f_3, \dots, f_{10}$

$f_1 \rightarrow 5$  times  
 $f_2 \rightarrow 2$  times  
 $f_3 \rightarrow 1$  time

$$\begin{array}{l|l} f_1 \propto 5 & f_{f \dots 10} \propto 0 \\ f_2 \propto 2 & \\ f_3 \propto 1 & \end{array}$$

## Categorical features

- Label encoding ( $0, 1$ )  $|f| = 2$
- One hot encoding,  $|f| = 3$  to  $6$
- Target encoding,  $|f| > 6$

$N$  unique values  $\hookrightarrow$  probability  
 $\hookrightarrow (N-1)$

$$f \rightarrow p(y=1 | f_i = 4)$$

$f_i$	$y$	$f_j$
4	0	$p_4$
5	0	$p_5$
6	1	

$p(y=1 | f_i = 5)$



