

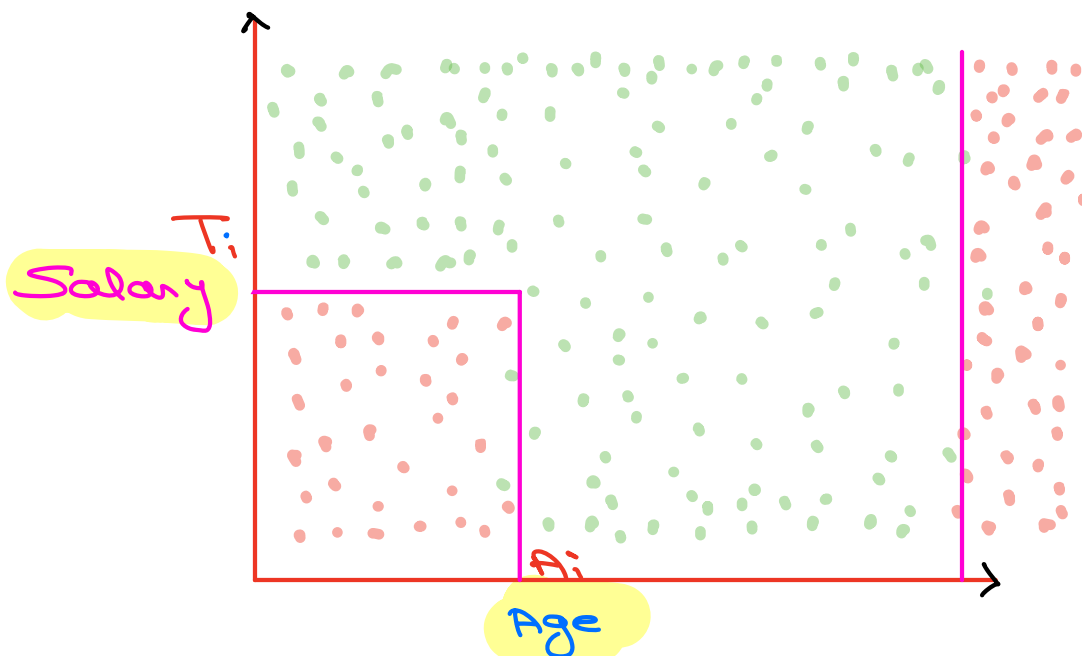
Decision Tree - 2

Agenda

- ⇒ Recap
- ⇒ Gini impurity
- ⇒ Splitting on Numerical feature
- ⇒ Imbalanced Data
- ⇒ Feature Scaling
- ⇒ Feature importance
- ⇒ DT Regression

Recap

```
if salary < T:  
    if age < A:  
        P = +1  
    else  
        P = -1  
else  
    .....
```



- ① Decision tree splits Data into Homogeneous regions using Axis parallel Hyperplanes.
- ② DT is easily interpretable

1) How do we decide which feature and Value to Split On

2) Target: Pure/Homogenous Node

3) How do we calculate purity?

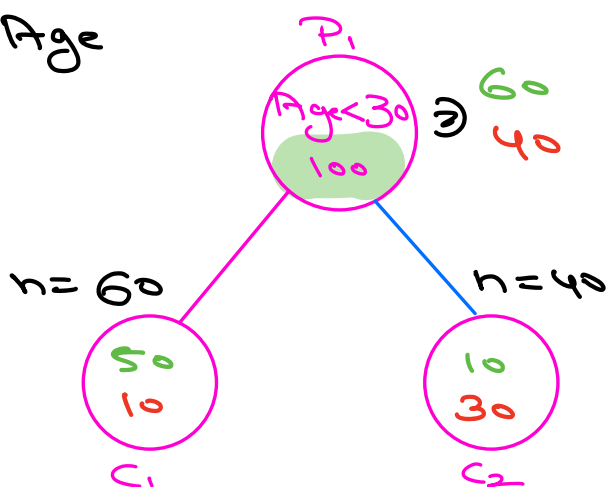
$$\text{Entropy}(H(Y)) = - \sum_{i=1}^k P(y_i) \log P(y_i)$$

where k classes

4) For Binary class : $k=2$
 P or $1-P$

$$H(Y) = - (P \times \log P + (1-P) \times \log (1-P))$$

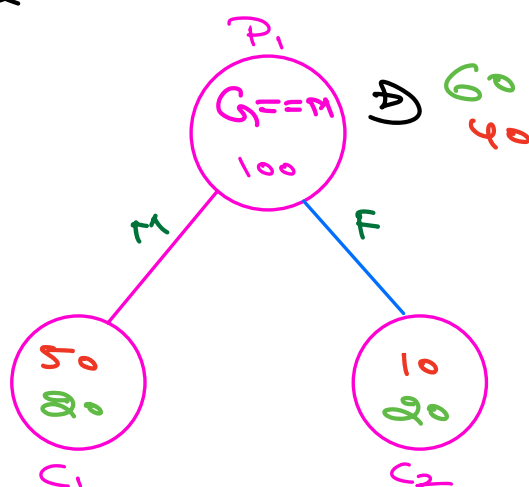
F1: Age
 Ex:



$$H(C_1) \approx 0.65$$

$$H(C_2) \approx 0.81$$

F2: Gender



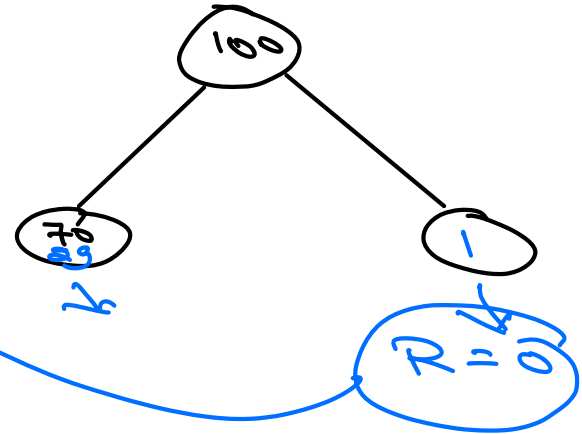
$$H(C_1) \approx 0.86$$

$$H(C_2) = 0.91$$

Total Entropy for Each Split

why not

$$\frac{C_1 + C_2}{2}$$



$$\frac{n_1}{n_{\text{Total}}} \times H(C_1) + \frac{n_2}{n_{\text{Total}}} \times H(C_2)$$

$$F_1 \Rightarrow 0.81$$

$$F_2 \Rightarrow 0.67$$

Information Gain
(Reduction in Entropy)

$$H(\text{Parent}) - H(\text{Child})$$

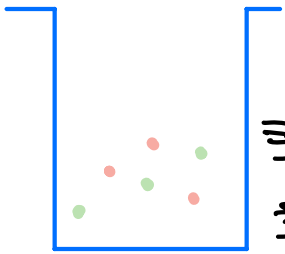
Limitation of Entropy

→ Lot of Calculations involve Log

Gini Impurity

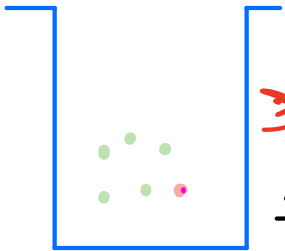
$$GI(y) = 1 - \sum_{i=1}^k (P(y_i))^2$$

$H(y)$



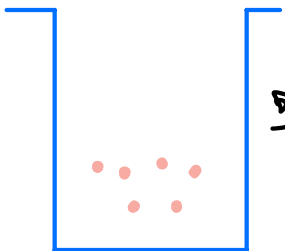
$$\Rightarrow 0.5 \times \log(0.5) + 0.5 \log(0.5)$$

→ 1



$$\Rightarrow \frac{1}{5} \times \log \frac{1}{5} + \frac{4}{5} \times \log \left(\frac{4}{5} \right)$$

→ 0.6



→ 0

$GI(y)$

$$1 - ((0.5)^2 + (0.5)^2)$$
$$1 - (0.25 + 0.25)$$

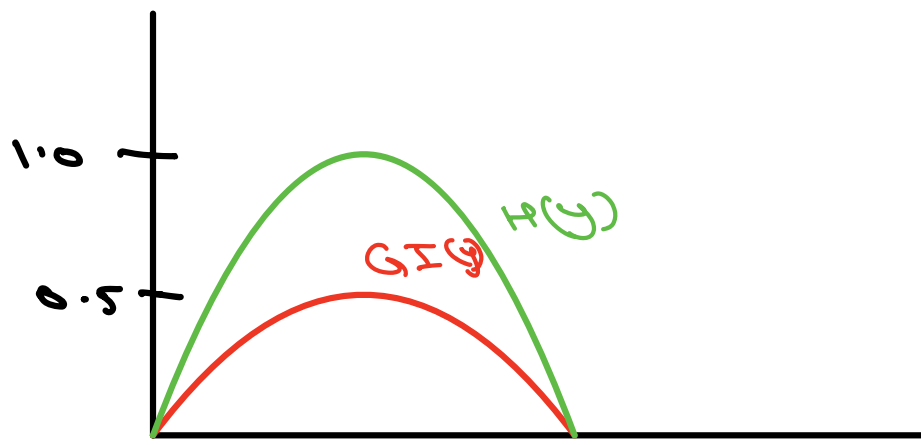
→ $1 - 0.5 = 0.5$

$$1 - \left(\left(\frac{1}{5} \right)^2 + \left(\frac{4}{5} \right)^2 \right)$$

→ 0.24

→ $1 - (1^2 + 0^2)$

→ 0



π_1

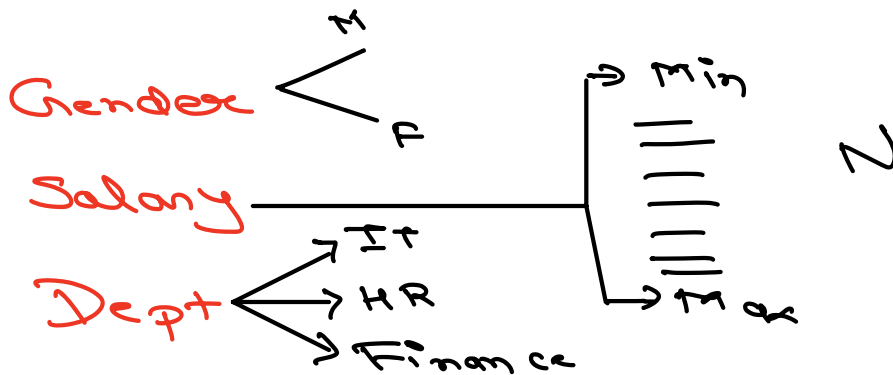
and

π_2

$$w(G(y_1)) = 0.32$$

$$w(G(y_2)) = 0.24$$

Splitting on Numerical

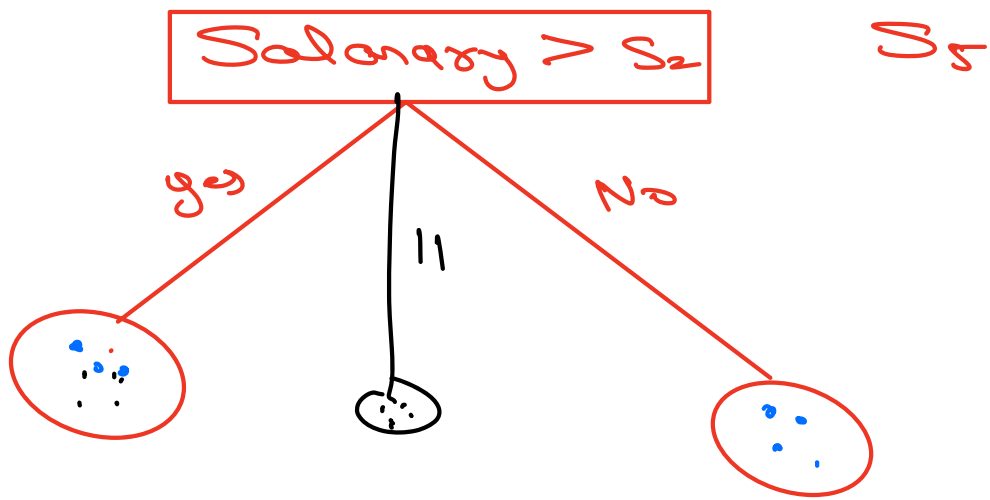


Brute force

Step 1: Get all Unique Values of Salary
 $[S_1, S_2, S_3, \dots, S_n]$

Step 2: Calculate IG for each Salary
 Value as threshold
 $[IG_1, IG_2, IG_3, \dots, IG_n]$

Step 3: Argmax of IG



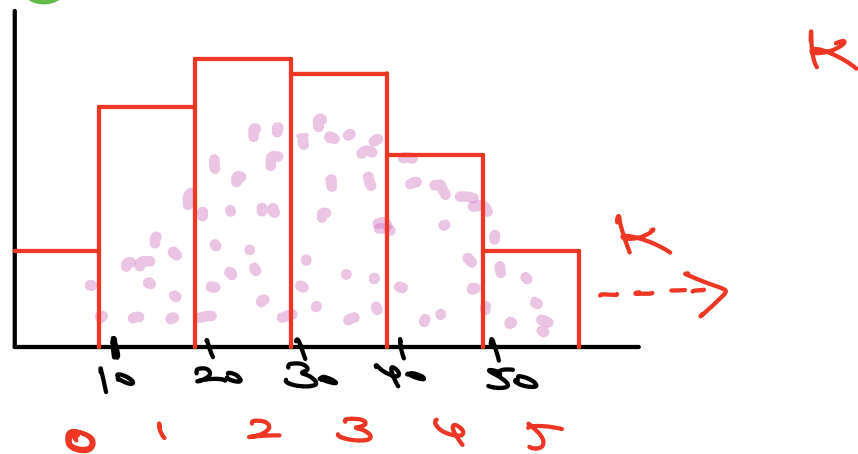
* Problem with this approach

d numerical feature
 n rows

$$O(n \times d)$$

How do we solve this

Binning



$n \leq 10000$ Unique Salary \Rightarrow n entropy
 \downarrow
 Binning
 \downarrow
 k entropy

Q - Will feature scaling have any impact on DT?

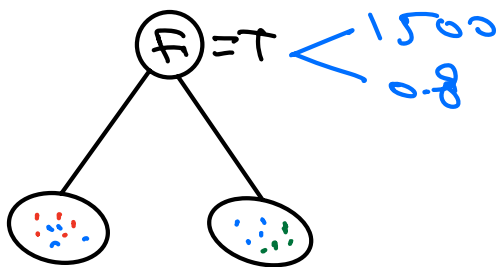
f_1 f_1 -scaled

2000 1.5

2500 1.8

1500 0.8

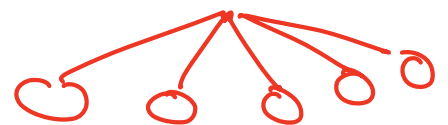
No impact



No need to Scale for DT

Q Categorical feature

⇒ 5 Unique Categorical : OHE



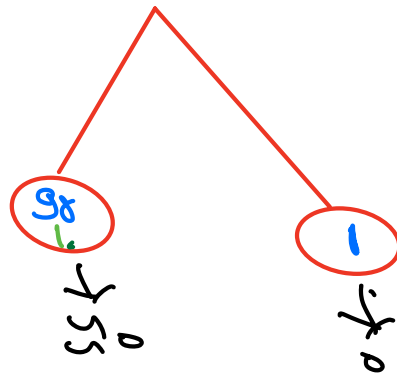
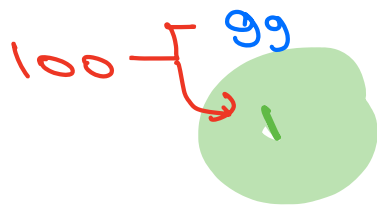
⇒ 10000 Unique Categories

Target Encoding

X O O O O O O O O O

Q Impact of Outliers \Rightarrow Yes
 \hookrightarrow prune Tree

Q Data Imbalance Impact DT \Rightarrow



Entropy 50

$$1 \Rightarrow \frac{98}{99} \%$$

How do we solve this

- \hookrightarrow Class-Weights
- \hookrightarrow SMOTE
- \hookrightarrow Undersampling / oversampling

For Tmrw

Feature Importance ✓

Decision Tree for Regression ?

Q What is feature Importance

⇒ How did you determine feature Importance in Linear

⇒ Linear Model

↳ weights

(coef)

$[-10, 13]$

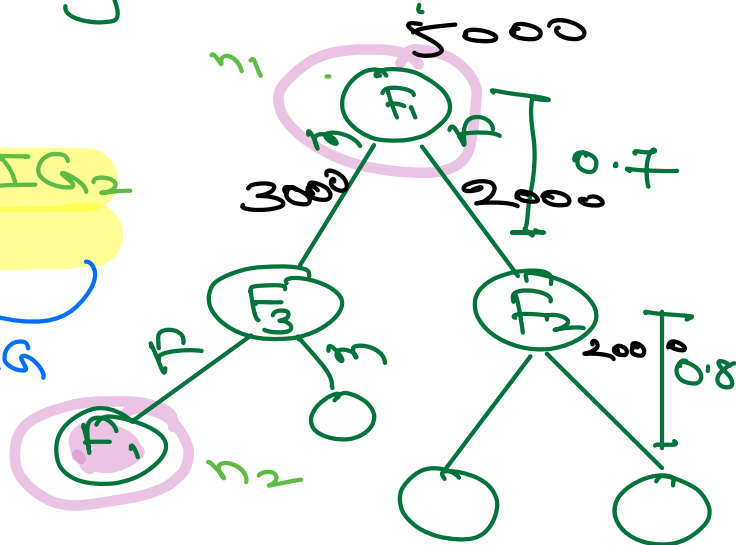
⇒ Decision Tree

No weights ✗

⇒ Higher the information Gain, more important is the feature ?

$$FI(f_i) = \frac{n_1 IG_1 + n_2 IG_2}{n}$$

normalised IG



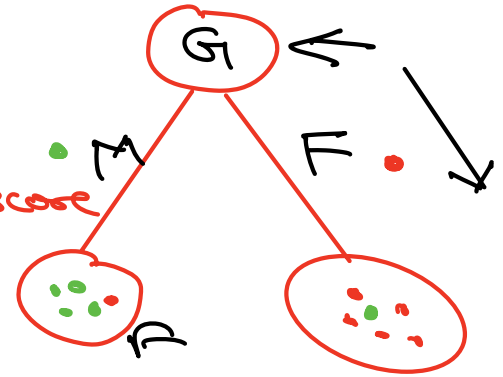
How Can we Use DT for Regression

DataPoint \Rightarrow π

What will be label and prob-score

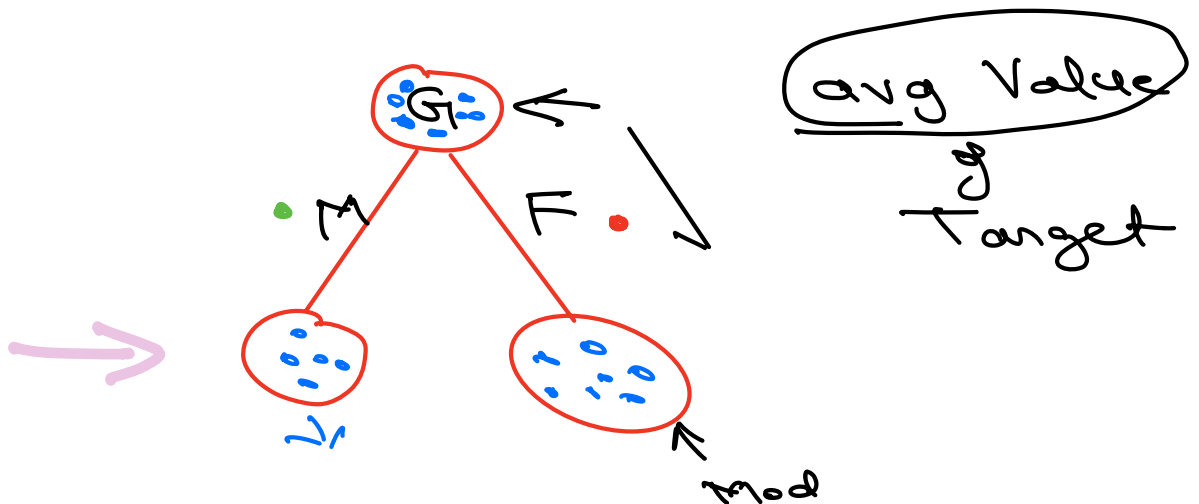
label \Rightarrow Red

Prob-score \Rightarrow $5/6$



Inference $\hookrightarrow O(1) \Rightarrow O(d)$
* DT is very fast for Inference

Regression



Inference \Rightarrow Mean or Average ①

Split \Rightarrow MSE or RMSE ②

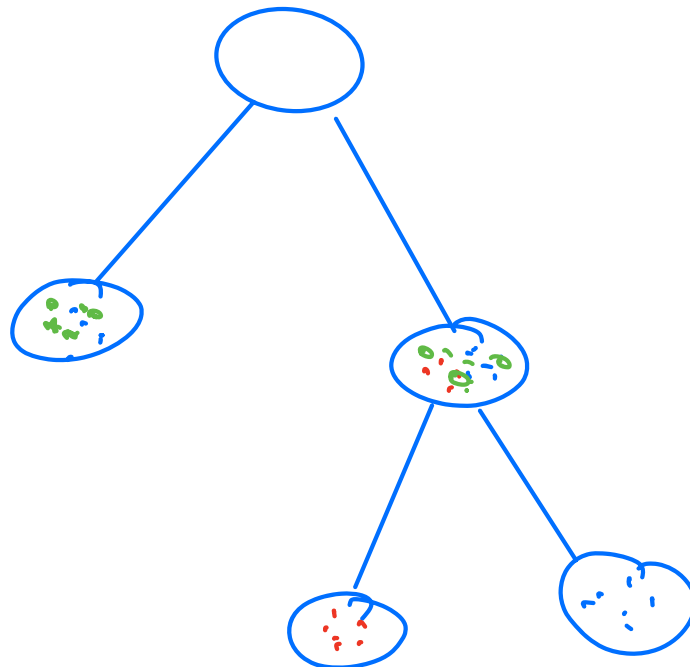
IG \rightarrow

Smote D To create Artificial D.P



k-points

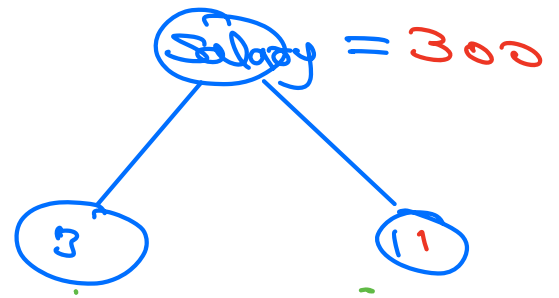
$$P' = \underbrace{\lambda_1 P_1 + \lambda_2 P_2}_{\text{...}} + \underbrace{\lambda_3 P_3}_{\text{...}}$$



min-leaf

	1,1	1,2	1,3
	2,1	2,2	2,3
Max-dep	10,1	10,2	10,3

Salary	Churn
5000	1
3000	1
2000	0
5000	1
1000	1



Salary	IG
5000	x_1
3000	x_2
2000	x_3
1000	1

$x_3 \rightarrow$ IG is Max

