

Data :

	Color	Type	Origin	Stolen?
1	Red	Sports	Domestic	Yes
2	Red	Sports	Domestic	No
3	Red	Sports	Domestic	Yes
4	Yellow	Sports	Domestic	No
5	Yellow	Sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	Yes
10	Red	Sports	Imported	No

Split level 1

* Color:

→ color = Red:

+ 5 counts

$$+ P(\text{Yes} | \text{Red}) = 3/5 = 0,6$$

$$+ P(\text{No} | \text{Red}) = 2/5 = 0,4$$

$$+ \text{Gini}(\text{Red}) = 1 - (0,6)^2 - (0,4)^2 \\ = 0,48$$

— color = Yellow

+ 5 counts

$$+ P(\text{Yes} | \text{Yellow}) = 2/5 = 0,4$$

$$+ P(\text{No} | \text{Yellow}) = 3/5 = 0,6$$

$$+ \text{Gini}(\text{Yellow}) = 1 - (0,4)^2 - (0,6)^2 \\ = 0,48$$

$$\text{WeightedGini}(\text{Color}) = \frac{5}{10} \cdot 0,48 + \frac{5}{10} \cdot 0,48 \\ = 0,48$$

* Type:

- Type = Sport:

+ Counts = 6

+ $P(\text{Yes} | \text{Sport}) = 4/6 = 2/3$

+ $P(\text{No} | \text{Sport}) = 2/6 = 1/3$

+ $Gini(\text{Sport}) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2$
 $\approx 0,44$

- Type = SUV:

+ Counts = 4

+ $P(\text{Yes} | \text{SUV}) = 1/4 = 0,25$

+ $P(\text{No} | \text{SUV}) = 3/4 = 0,75$

+ $Gini(\text{SUV}) = 1 - 0,25^2 - 0,75^2$
 $= 0,375$

Weighted Gini (Type) = $\frac{6}{10} \cdot 0,44 + \frac{4}{10} \cdot 0,375$
 $= 0,414$

* Origin:

- Origin = Domestic:

$$+ \text{Counts} = 5$$

$$+ P(\text{Yes} | \text{Domestic}) = 2/5 = 0,4$$

$$+ P(\text{No} | \text{Domestic}) = 3/5 = 0,6$$

$$+ \text{Gini}(\text{Domestic}) = 1 - 0,4^2 - 0,6^2 \\ = 0,48$$

- Origin = Imported:

$$+ \text{Counts} = 5$$

$$+ P(\text{Yes} | \text{Imported}) = 3/5 = 0,6$$

$$+ P(\text{No} | \text{Imported}) = 2/5 = 0,4$$

$$+ \text{Gini}(\text{Imported}) = 1 - 0,6^2 - 0,4^2 \\ = 0,48$$

weighted Gini (Origin)

$$= \frac{5}{10} \cdot 0,48 + \frac{5}{10} \cdot 0,48$$

$$= 0,48$$

- Weighted Gini (color) = 0,48

 (type) = 0,414

 (origin) = 0,48

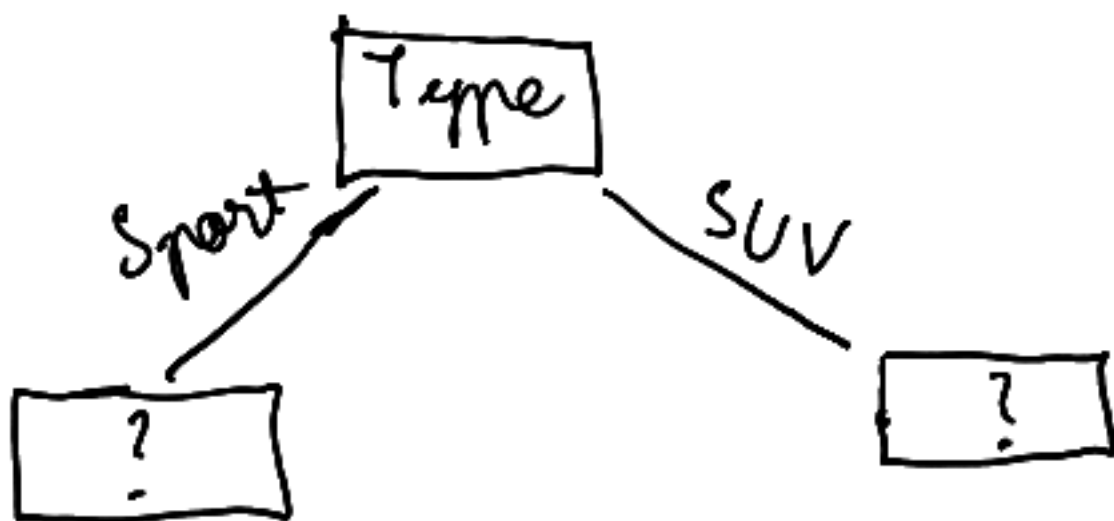
⇒ The impurities of «Type» is lowest

⇒ choose «Type» as the 1st split

Split level 2

- has 2 values: Sport & SUV

- Each of these values will be a new branch.



* Spot branch: 6 counts

Spot \rightarrow Color

- Color = Red + Yellow

- Sport. Red:

+ Counts = 4

+ $P(\text{Yes} | \text{Sport. Red}) = 3/4 = 0,75$

+ $P(\text{No} | \text{Sport. Red}) = 1/4 = 0,25$

+ Gini (Spot \rightarrow Red) = $1 - 0,75^2 - 0,25^2$
= 0,375

- Sport. Yellow:

+ Counts = 2

+ $P(\text{Yes} | \text{Sport. Yellow}) = 1/2 = 0,5$

+ $P(\text{No} | \text{Sport. Yellow}) = 1/2 = 0,5$

+ Gini (Spot \rightarrow Yellow) = $1 - 0,5^2 - 0,5^2$
= 0,5

Weighted Gini (Spot \rightarrow Color) = $\frac{4}{6} \cdot 0,375 + \frac{2}{6} \cdot 0,5$
= 0,417

$\boxed{\text{Sport} \rightarrow \text{Origin}}$

- $\text{Origin} = \text{Domestic} + \text{Imported}$

- Sport. Domestic:

+ Counts = 4

+ $P(\text{Yes} | \text{Sport. Domestic}) = 2/4 = 0,5$

+ $P(\text{No} | \text{Sport. Domestic}) = 2/4 = 0,5$

+ $\text{Gini}(\text{Sport} \rightarrow \text{Domestic}) = 1 - 0,5^2 - 0,5^2$
 $= 0,5$

- Sport. Imported

+ Counts = 2

+ $P(\text{Yes} | \text{Sport. Imported}) = 2/2 = 1$

+ $P(\text{No} | \text{Sport. Imported}) = 0/2 = 0$

+ $\text{Gini}(\text{Sport} \rightarrow \text{Domestic}) = 1 - 1^2 - 0^2$
 $= 0$

$\text{Weighted Gini}(\text{Sport} \rightarrow \text{Origin}) = \frac{4}{6} \cdot 0,5 + \frac{2}{6} \cdot 0$
 $= 0,33$

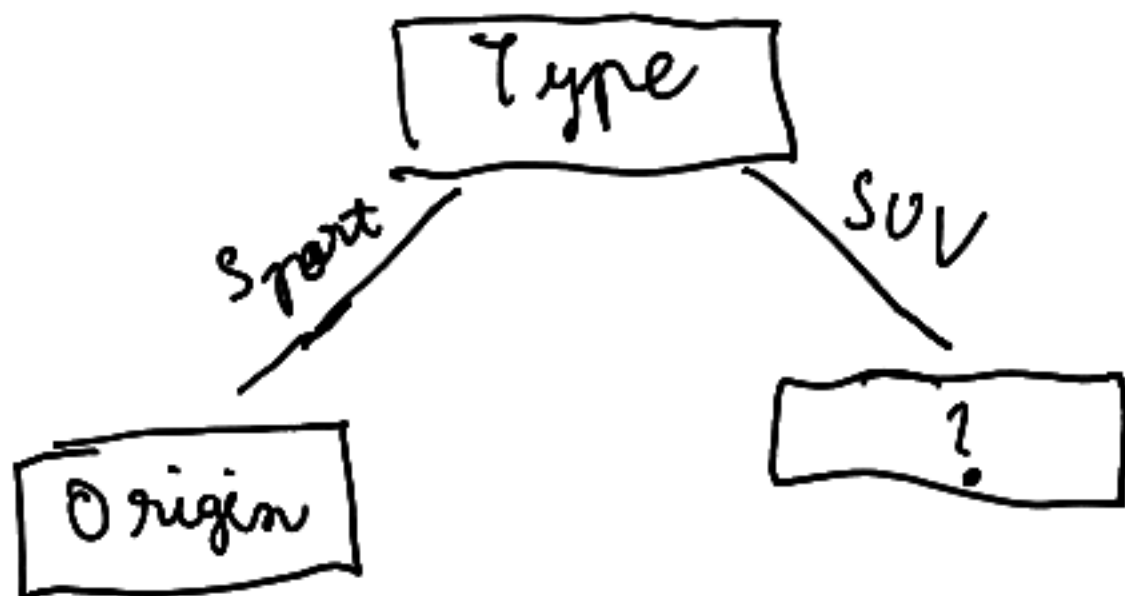
Weighted time (Sport \rightarrow dest) = 0,417

(Sport \rightarrow Origin) = 0,33

$$(Sport \rightarrow Origin) = 0,33$$

\Rightarrow "Origin" has lower impurity

\Rightarrow choose "Origin" as next node.
on "Sport" branch



* SUV branch: 4 counts.

SUV \rightarrow Color

- Color = Red + Yellow

- SUV, Red

+ Counts = 1

+ $P(\text{Yes} | \text{SUV, Red}) = 0/1 = 0$

+ $P(\text{No} | \text{SUV, Red}) = 1/1 = 1$

+ $\text{Gini}(\text{SUV} \rightarrow \text{Red}) = 1 - 0^2 - 1^2$
 $= 0$

- SUV, Yellow

+ Counts = 3

+ $P(\text{Yes} | \text{SUV, Yellow}) = 1/3$

+ $P(\text{No} | \text{SUV, Yellow}) = 2/3$

+ $\text{Gini}(\text{SUV} \rightarrow \text{Yellow}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2$
 $= 0,44$

Weighted $\text{Gini}(\text{SUV} \rightarrow \text{Color}) = \frac{1}{4} \cdot 0 + \frac{3}{4} \cdot 0,44$
 $= 0,33$

[SUV \rightarrow Origin]

- Origin = Domestic + Imported

- SUV, Domestic:

$$+ \text{Counts} = 1$$

$$+ P(\text{Yes} | \text{SUV, Domestic}) = 0/1 = 0$$

$$+ P(\text{No} | \text{SUV, Domestic}) = 1/1 = 1$$

$$+ \text{Gini}(\text{SUV} \rightarrow \text{Domestic}) = 1 - 0^2 - 1^2 = 0$$

- SUV, Imported:

$$+ \text{Counts} = 3$$

$$+ P(\text{Yes} | \text{SUV, Imported}) = 1/3$$

$$+ P(\text{No} | \text{SUV, Imported}) = 2/3$$

$$+ \text{Gini}(\text{SUV} \rightarrow \text{Imported}) = 1 - \frac{1}{9} - \frac{4}{9} = 0,44$$

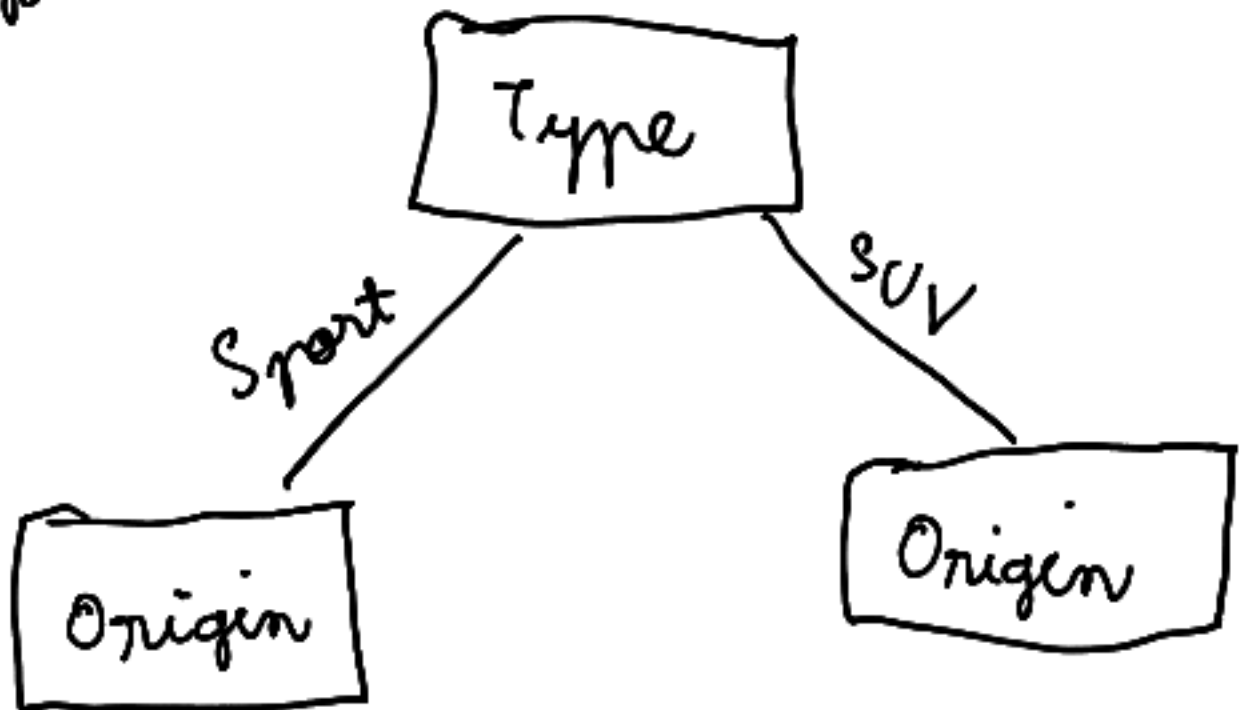
Weighted Gini (SUV \rightarrow Origin)

$$= \frac{1}{4} \cdot 0 + \frac{3}{4} \cdot 0,44 = 0,11$$

Weighted Cini (SVV \rightarrow color) = 0,33

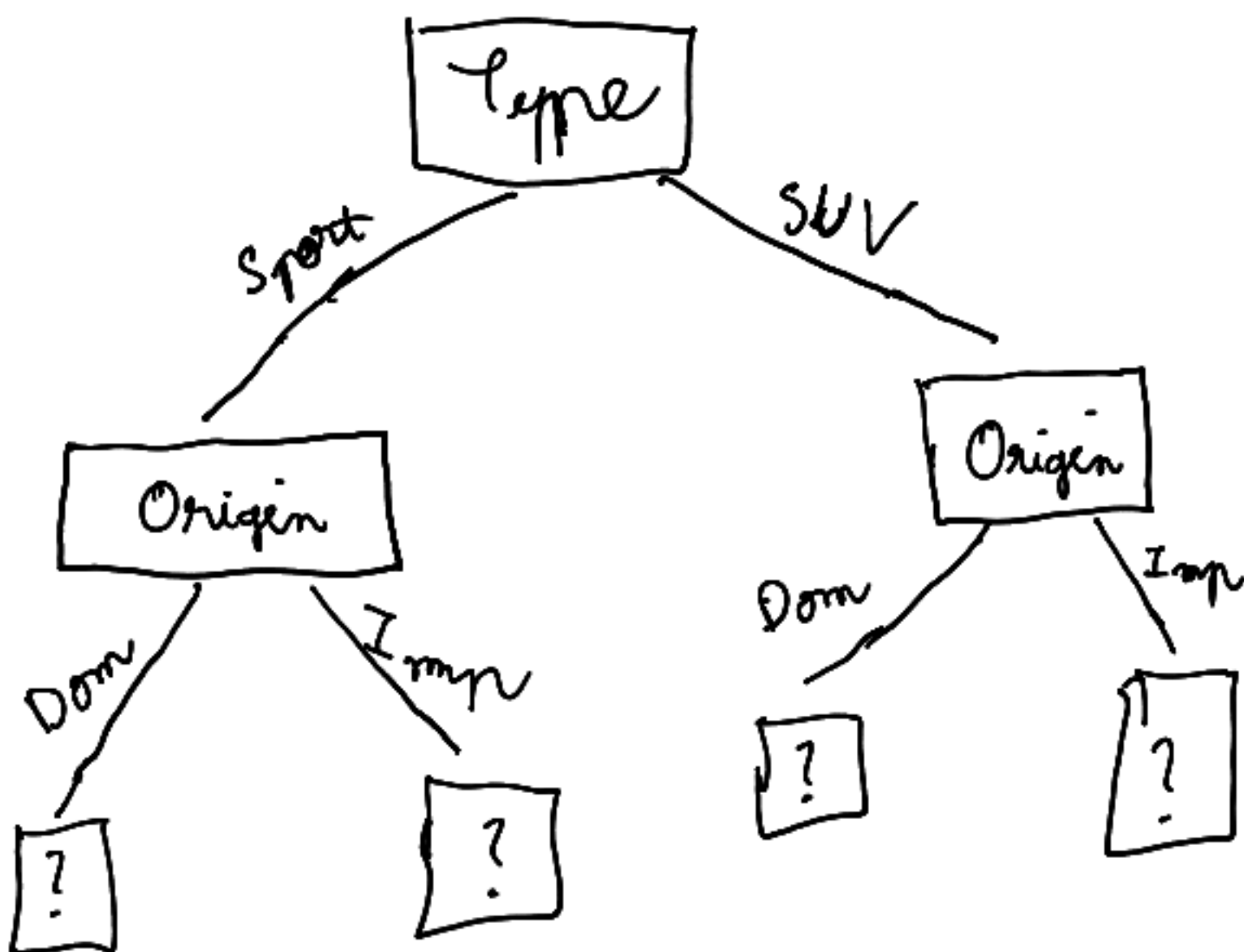
———— (SUV → Origin) = 0,11

\Rightarrow choose "origin" as next node
for "SU_V" branch



Split level 3

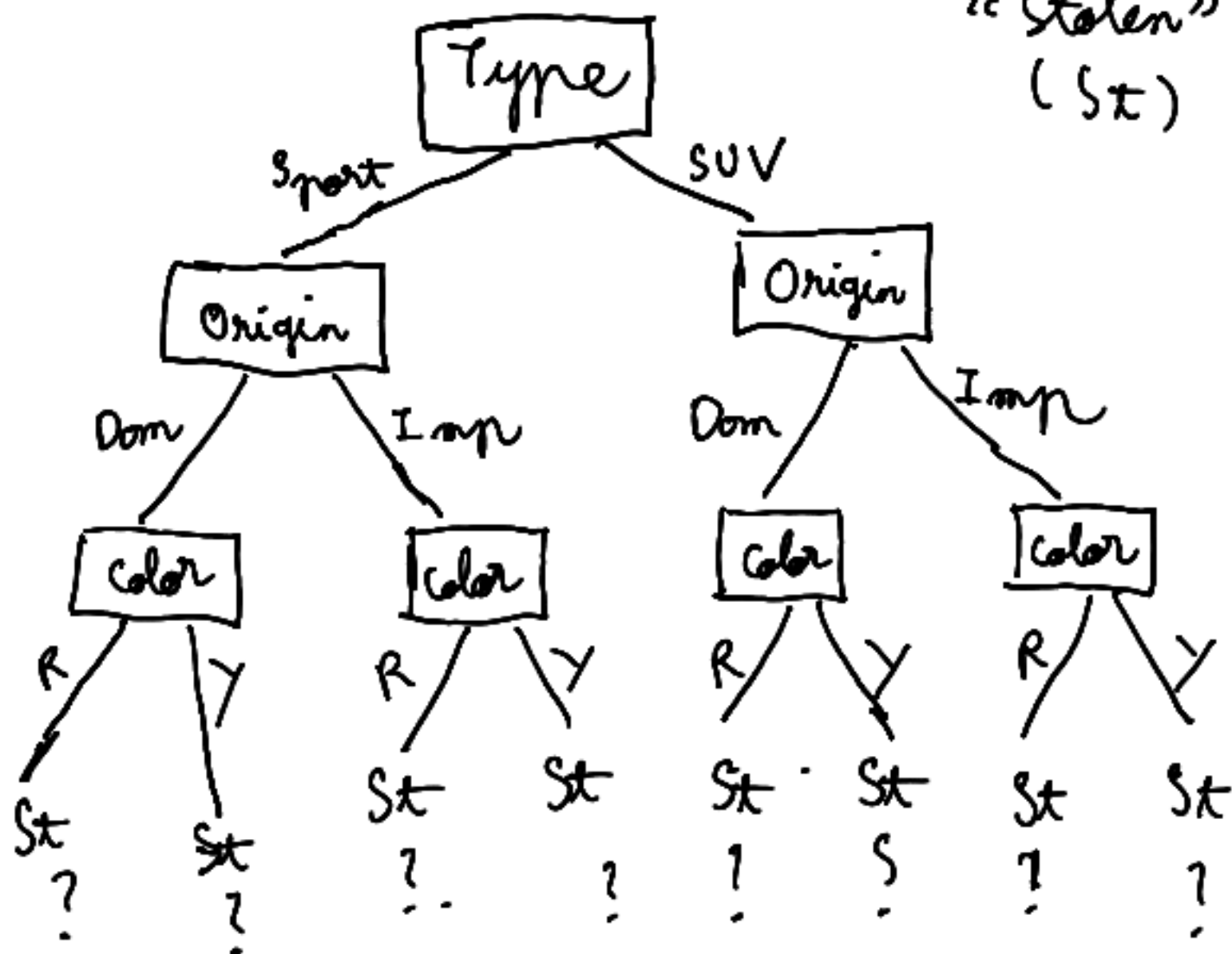
- Origin has 2 values : Domestic and Imported
⇒ 2 branches



- "Color" is the only feature left,
so it's the final node (leaf node)

- Color \longrightarrow 2 branches: Red + Yellow

- These 2 branches lead to outcome
"Stolen" (St)



$$- P(\text{Yes} \mid \text{Sport. Dom. Red}) = 2/3$$

$$P(\text{No} \mid \text{Sport. Dom. Red}) = 1/3$$

$$\Rightarrow \text{"Stolen"} = \text{Yes}$$

$$- P(\text{Yes} \mid \text{Sport. Dom. Yellow}) = 0/1 = 0$$

$$P(\text{No} \mid \text{Sport. Dom. Yellow}) = 1/1 = 1$$

$$\Rightarrow \text{"Stolen"} = \text{No}$$

$$= P(\text{Yes} \mid \text{Sport. Imp. Red}) = 1/1 = 1$$

$$P(\text{No} \mid \text{Sport. Imp. Red}) = 0/1 = 0$$

$$\Rightarrow \text{"Stolen"} = \text{Yes}$$

$$- P(\text{Yes} \mid \text{Sport. Imp. Yellow}) = 1/1 = 1$$

$$P(\text{No} \mid \text{Sport. Imp. Yellow}) = 0/1 = 0$$

$$\Rightarrow \text{"Stolen"} = \text{Yes}$$

$$\begin{aligned}
 - P(\text{Yes} | \text{SUV, Dom. Red}) &= X \\
 P(\text{No} | \text{SUV, Dom. Red}) &= X
 \end{aligned}
 \left. \vphantom{\begin{aligned} P(\text{Yes} | \text{SUV, Dom. Red}) \\ P(\text{No} | \text{SUV, Dom. Red}) \end{aligned}} \right\} \begin{array}{l} \text{No data} \\ \text{No predict} \end{array}$$

$$- P(\text{Yes} | \text{SUV, Dom. Yellow}) = 0/1 = 0$$

$$P(\text{No} | \text{SUV, Dom. Yellow}) = 1/1 = 1$$

$$\Rightarrow \text{"Stolen"} = \text{No}$$

$$- P(\text{Yes} | \text{SUV, Imp. Red}) = 0/1 = 0$$

$$P(\text{No} | \text{SUV, Imp. Red}) = 1/1 = 1$$

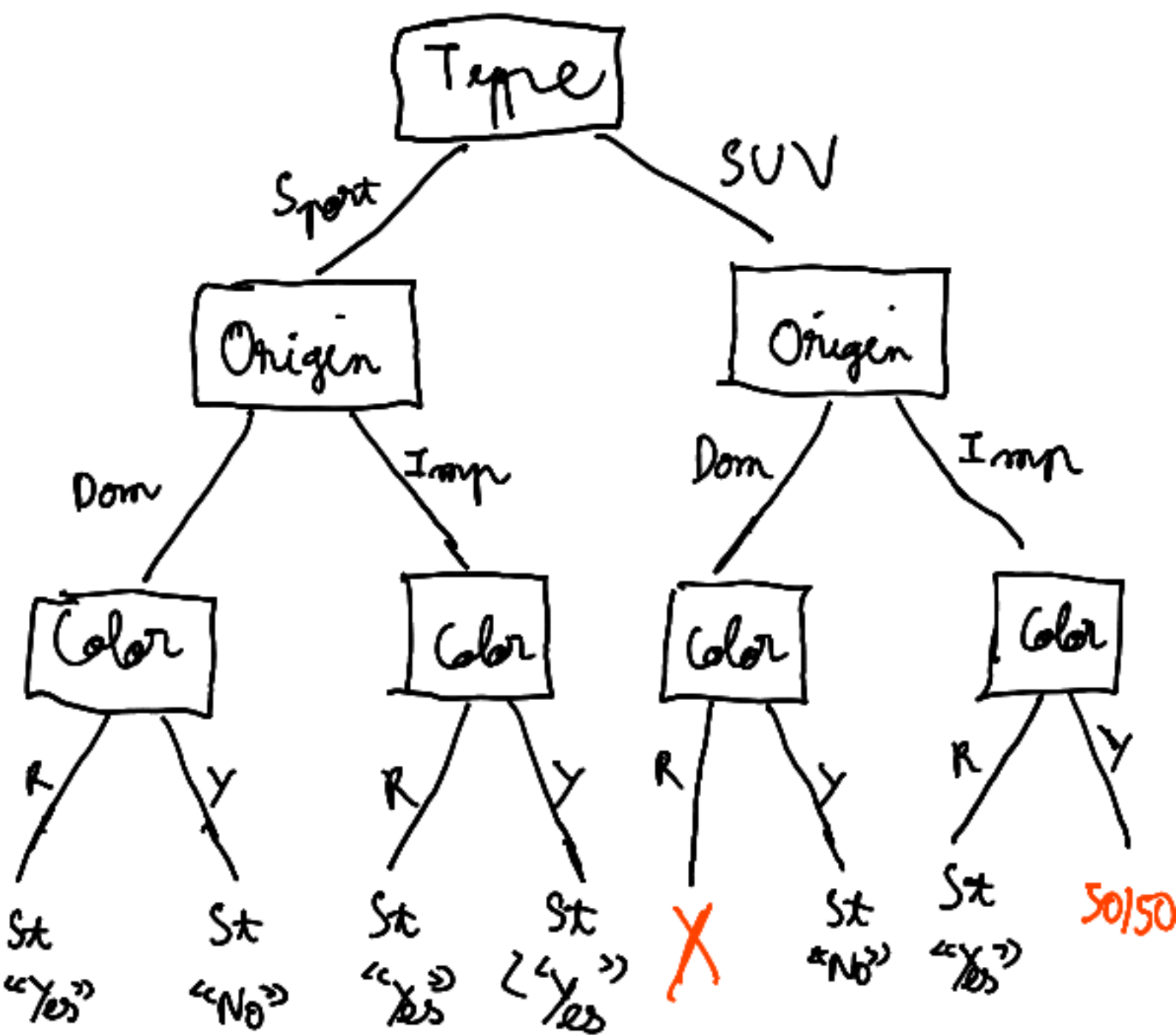
$$\Rightarrow \text{"Stolen"} = \text{Yes}$$

$$- P(\text{Yes} | \text{SUV, Imp. Yellow}) = 1/2$$

$$P(\text{No} | \text{SUV, Imp. Yellow}) = 1/2$$

$$\Rightarrow \text{"Stolen"} = \dots \text{unsure} \dots$$

Final tree:



- New data:

Color	Type	Origin
Yellow	Sport	Imported

Type = Sport
 ↳ origin = Imp
 ↳ color = Yellow → Stolen: "Yes"

Simplify the tree

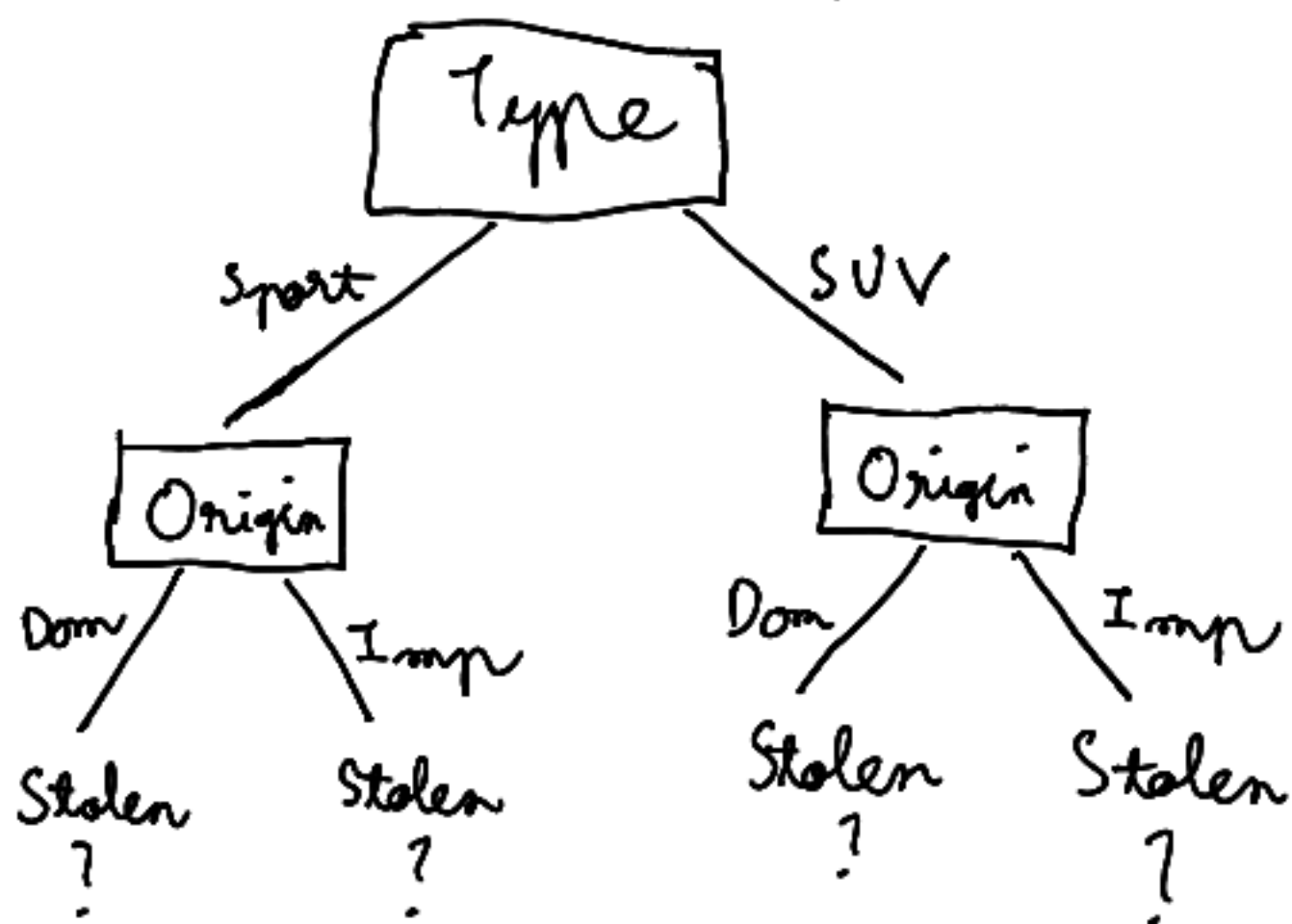
- If the tree is too detailed
 - ↳ can cause "Overfitting"

(as we can see in the example above, the tree cannot predict for "SUV, Domestic, Red" case due to lack of data,

and cannot predict "SUV, Imported, Yellow" due to tied probabilities)
$$P(\text{Yes}) = P(\text{No}) = \frac{1}{2}$$

⇒ Must simplify the tree, by stop the splitting at a certain level.

- In the above example,
let's stop at "Origin"



$$- P(\text{Yes} | \text{Sport. Domestic}) = 2/4 = 1/2$$

$$P(\text{No} | \text{Sport. Domestic}) = 2/4 = 1/2$$

$$\Rightarrow \dots \text{ Unsure.}$$

$$- P(\text{Yes} | \text{Sport. Imported}) = 2/2 = 1$$

$$P(\text{No} | \text{Sport. Imported}) = 0/2 = 0$$

$$\Rightarrow \text{“Stolen”} = \text{Yes}$$

$$- P(\text{Yes} | \text{SUV. Domestic}) = 0/1 = 0$$

$$P(\text{No} | \text{SUV. Domestic}) = 1/1 = 1$$

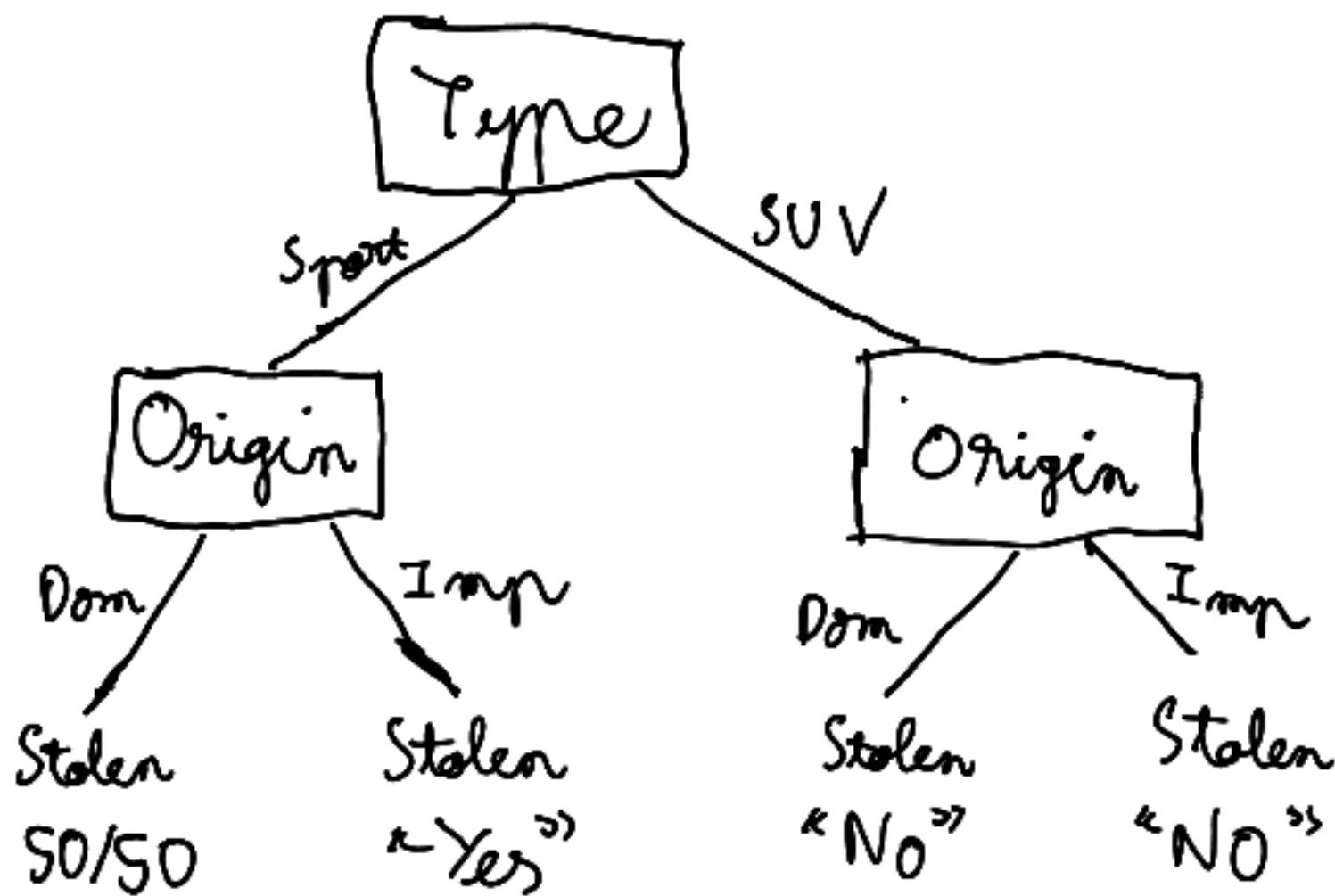
$$\Rightarrow \text{“Stolen”} = \text{No}$$

$$- P(\text{Yes} | \text{SUV. Imported}) = 1/3$$

$$P(\text{No} | \text{SUV. Imported}) = 2/3$$

$$\Rightarrow \text{“Stolen”} = \text{No}$$

Final Simplified Tree



- New data:

Color	Type	Origin
Yellow	Sport	Imported

check Type = Sport

↳ check Origin = Imp

↳ Stolen
"Yes"