

Comprendre le principe des arbres de décision et performances

① a- Le meilleur test à choisir à la racine est t1, car il va maximiser le gain en information au regard des deux autres tests.

En effet, par toutes les branches issues de ce test, on arrive à des nœuds dans lesquels les exemples se répartissent très majoritairement dans une seule classe, ce qui mène à un indice de Gini (ou d'entropie) le plus faible.

Le test t2 n'aide en rien la répartition des exemples dans des classes uniques; le test t3 ne mène que partiellement à une telle discrimination en classes.

b- Choisissons t1.

$$A(t1=1) = \frac{9}{16} = (+)$$

$$A(t1=2) = \frac{14}{15} = (-)$$

$$A(t1=3) = \frac{5}{5} = (+)$$

} On choisit la classe majoritaire des exemples arrivant à chaque feuille

Seuls 2 exemples sont à déplorer: l'exemple négatif classé en (+) car sm t1=1, et l'exemple positif classé en (-) car sm t1=2.

$$\text{erreur} = \frac{2}{30} = 0.07 \text{ sur ens. d'apprentissage}$$

② a+b

$$t1(x_1) = +$$

$$y_1 = +$$

OK

$$t1(x_2) = +$$

$$y_2 = +$$

OK

$$t1(x_3) = -$$

$$y_3 = +$$

erreur

$$t1(x_4) = +$$

$$y_4 = +$$

OK

$$t1(x_5) = +$$

$$y_5 = -$$

erreur

$$t1(x_6) = -$$

$$y_6 = -$$

OK

$$t1(x_7) = -$$

$$y_7 = -$$

OK

$$t1(x_8) = +$$

$$y_8 = +$$

OK

Nous comptons 2 erreurs sur 8 exemples

$$\text{donc } \text{erreur} = 0.25 \text{ sur ens. test.}$$

La matrice de confusion est:

<div> <div>classe</div> <div>pred</div> </div> <div> <div>vraie classe</div> <div></div> </div>	+	-
+	4 TP	1 FN
-	1 FP	2 TN

les erreurs sont hors diagonales, nous venons

$$e = \frac{2}{8} = 0.25$$

$$\underline{c} \quad \underline{\text{rappel}(+)} = \frac{TP}{TP + FN} = \frac{4}{5} = \underline{0.8}$$

Construction d'un arbre de décision

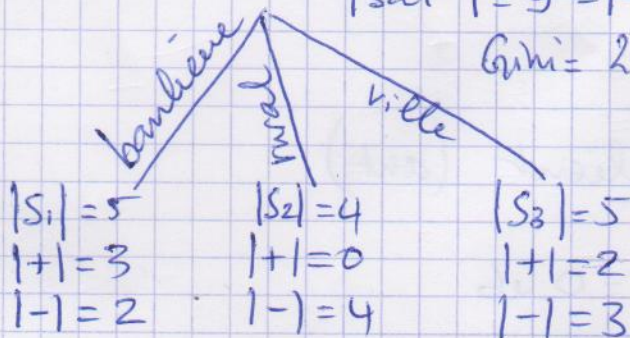
Construisons un stump: il s'agit de choisir le meilleur test à la racine, ie celui qui maximise le gain en information.

Calculons ce gain pour chaque test basique possible.

1. Test sur emplacement (epl)

insat = +
sat = -

empl $|S| = 14$
 $|insat| = 5 = |+|$
 $|sat| = 9 = |-|$



$$Gini = \frac{2 \times 5 \times 9}{14^2} = 0.46 = Gini(E)$$

$$Gini = \frac{2 \times 3 \times 2}{5^2} = 0.48$$

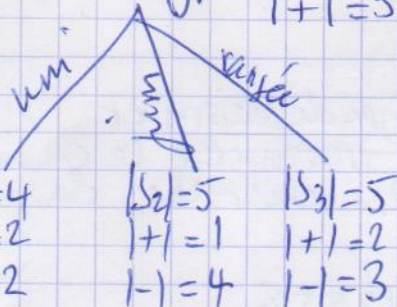
$$Gini = \frac{2 \times 0 \times 4}{4^2} = 0$$

$$Gini = \frac{2 \times 2 \times 3}{5^2} = 0.48$$

$$\begin{aligned} \text{Gain}(E, \text{empl}) &= Gini(E) - \sum_i P_i Gini(S_i) \\ &= 0.46 - \left(\frac{5}{14} \cdot 0.48 + \frac{4}{14} \cdot 0 + \frac{5}{14} \cdot 0.48 \right) \\ &= 0.117 \end{aligned}$$

2. Test sur type de maison (typ)

typ $|S| = 14$
 $|+| = 5$ $|-| = 9$



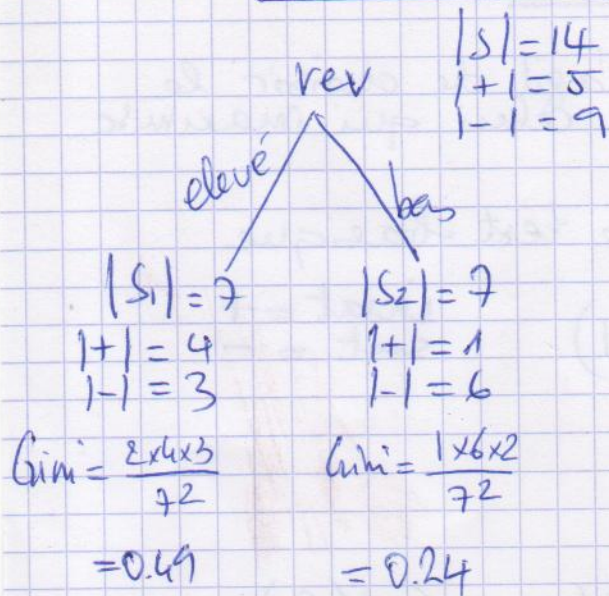
$$Gini = \frac{2 \times 2 \times 2}{4^2} = 0.5$$

$$Gini = \frac{1 \times 4 \times 2}{5^2} = 0.32$$

$$Gini = \frac{2 \times 3 \times 2}{5^2} = 0.48$$

$$\begin{aligned} \text{Gain}(E, \text{typ}) &= Gini(E) - \sum_i P_i Gini(S_i) \\ &= 0.46 - \left(\frac{4}{14} \cdot 0.5 + \frac{5}{14} \cdot 0.32 + \frac{5}{14} \cdot 0.48 \right) \\ &= 0.031 \end{aligned}$$

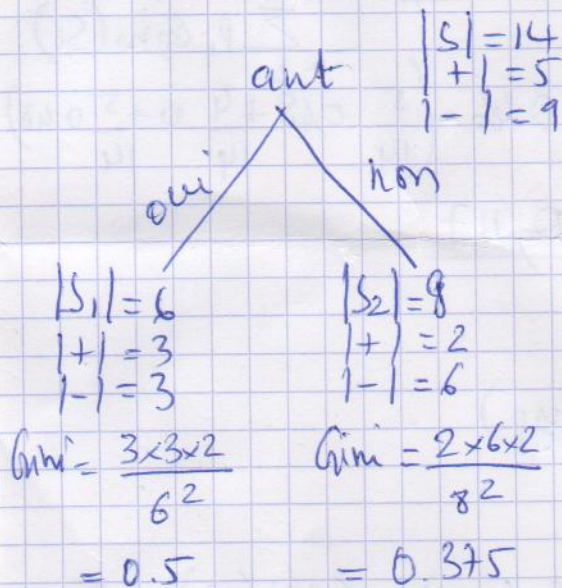
3. Test sur revenu (rev)



$$Gini(\mathcal{E}) = 0.46$$

$$Gain(\mathcal{E}, rev) = 0.46 - \left(\frac{7}{14} \cdot 0.49 + \frac{7}{14} \cdot 0.24 \right) = 0.095$$

4. Test sur antériorité du client (ant)



$$Gini(\mathcal{E}) = 0.46$$

$$Gain(\mathcal{E}, ant) = 0.46 - \left(\frac{6}{14} \cdot 0.5 + \frac{8}{14} \cdot 0.375 \right) = 0.031$$

Conclusion

Le gain en information est maximal pour le test sur l'emplacement, que nous choisissons donc à la racine.

Nous attribuons la classe majoritaire pour la décision en chaque feuille.

