

INPF12 : projet 2023-2024

Tom Gilgenkrantz

2 mai 2024

1 Les arbres de décisions

1.1 Introduction

Un arbre de décision est un outil de prise de décision sous forme d'arbre avec des nœuds, des feuilles (nœuds terminaux) et une profondeur. Il est largement utilisé en intelligence artificielle et en apprentissage automatique pour résoudre des problèmes de classification et de régression. Chaque nœud interne de l'arbre correspond à un test sur une caractéristique, chaque branche représente le résultat du test et chaque feuille représente une prédiction de classe.

On s'intéressera ici à la construction d'un arbre de décision pour la classification de points colorés (rouges ou bleus) d'un dataset en deux classes. On commence tout d'abord par définir un type point qui est un couple de flottants choisis ici dans $\llbracket 0, 1 \rrbracket^2$.

```
1 type pos = { x:float ; y:float; } ;;
```

On définit ensuite un type couleur qui est un triplet d'entiers de $\llbracket 0, 255 \rrbracket^3$.

```
1 type color = {r:int ; g:int ; b:int} ;;
```

On peut enfin définir le type dataset qui est une liste de coordonnées colorées.

```
1 type 'a dataset = (pos*color) list;;
```

1.2 Type tree

Avant de pouvoir construire un arbre de décision, il est nécessaire de déterminer comment les tests seront effectués afin de séparer le dataset initial en sous-ensembles. L'idée la plus simple est tout d'abord de couper au milieu du dataset alternativement verticalement et horizontalement afin de déterminer 2 sous-ensembles de ce dernier. On pourrait ainsi définir le tree comme un arbre binaire avec pour nœuds des conditions sous forme de flottants de $\llbracket 0, 1 \rrbracket$, et deux sous-arbres, l'un avec une coordonnée inférieure au nœud parent, l'autre avec une coordonnée supérieure ou égale.

```

1 type tree =
2   | Leaf of color
3   | Node of float * tree * tree ;;

```

Ce n'est pourtant pas ce type que j'ai choisi de garder. En choisissant de séparer le dataset en 4 sous-ensembles, on peut définir un arbre de décision comme un arbre avec pour noeuds des conditions sous forme de couples de flottants de $[0,1]^2$, et 4 sous-arbres, chacun correspondant à un sous-ensemble du dataset initial. Cette implémentation permet de beaucoup simplifier l'implémentation puisqu'il n'est plus nécessaire d'alterner entre coupe horizontale et verticale. On peut donc définir le type tree comme suit :

```

1 type tree =
2   | Leaf of color
3   | Node of pos * tree * tree * tree * tree ;;

```

On définit enfin des couleurs de base pour les points rouges et bleus.

```

1 let red = {r=255 ; g=0 ; b=0} ;;
2 let blue = {r=0 ; g=0 ; b=255} ;;

```

2 Fonctions auxiliaires

Avant de commencer à implémenter la construction de l'arbre de décision, il est nécessaire de définir quelques fonctions auxiliaires qui seront utiles pour la suite. On commence par définir une fonction qui permet de compter le nombre de points rouges et bleus dans un dataset. Cela sera particulièrement utile pour déterminer la couleur moyenne d'un sous-ensemble de points, ou simplement vérifier qu'un dataset n'est pas vide ou d'une couleur unique.

```

1 (* ('a * color) list -> int -> int -> int * int *)
2 let rec compter_couleurs training rouge bleu =
3   match training with
4   | [] -> rouge , bleu
5   | head::tail ->
6     if (snd head) = red
7     then compter_couleurs tail (rouge+1) bleu
8     else
9       compter_couleurs tail rouge (bleu+1) ;;

```

On définit ensuite une fonction qui permet de séparer un dataset en 4 sous-ensembles en fonction d'une pos donnée. On peut alors définir la fonction suivante :

```

1 (* (pos * 'a) list ->
2   (pos * 'a) list ->

```

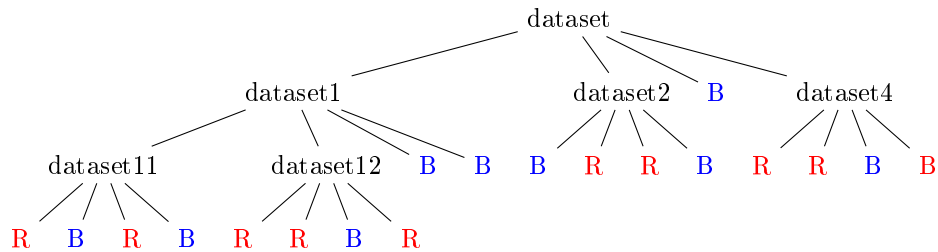
```

3   (pos * 'a) list ->
4   (pos * 'a) list ->
5   (pos * 'a) list ->
6   pos ->
7   (pos * 'a) list * (pos * 'a) list * (pos * 'a)
      list * (pos * 'a) list *)
8 let rec separer_dataset d1 d2 d3 d4 training x_y =
9   match training with
10  | [] -> d1 , d2 , d3 , d4
11  | head::tail ->
12      if (fst head).x < x_y.x && (fst head).y >= x_y.y
13      then separer_dataset (head::d1) d2 d3 d4 tail
14      x_y
15      else if (fst head).x >= x_y.x && (fst head).y >=
16      x_y.y
17      then separer_dataset d1 (head::d2) d3 d4 tail
18      x_y
19      else if (fst head).x < x_y.x && (fst head).y <
20      x_y.y
21      then separer_dataset d1 d2 (head::d3) d4 tail
22      x_y
23      else
24          separer_dataset d1 d2 d3 (head::d4) tail x_y ;;

```

3 Construction de l'arbre de décision

On peut maintenant commencer à construire l'arbre de décision. On se sert d'une fonction récursive qui permet de construire un arbre de décision à partir d'un dataset. On commence par vérifier si le dataset a une taille inférieure à 1, d'une seule couleur ou bien si on a atteint la profondeur maximale choisie, ou si tous les points sont de la même couleur. Dans ce cas, on renvoie une feuille de la couleur moyenne des points. Sinon, on sépare le dataset en 4 sous-ensembles, et on construit récursivement les sous-arbres. Voici un exemple d'arbre de décision de profondeur 4 tel que défini ci-dessus :



```

1   (* int -> 'a dataset -> tree *)
2   let build_tree prof training =

```

```

3  let rec build_node profondeur dataset separations =
4    let rouge , bleu = compter_couleurs dataset 0 0 in
5    if profondeur >= prof + 1 || rouge = 0 || bleu =
      0 || rouge + bleu < 2
6      then if rouge > bleu then Leaf red else Leaf
          blue
7    else
8      let epsilon = 1.0 /. 2.0 ** (float_of_int
          profondeur) in
9      let d1 , d2 , d3 , d4 = separer_dataset [] []
          [] [] dataset separations in
10     Node (separations,
11           build_node (profondeur+1) d1
              {x=(separations.x -. epsilon);
               y=(separations.y +. epsilon)},
12           build_node (profondeur+1) d2
              {x=(separations.x +. epsilon);
               y=(separations.y +. epsilon)},
13           build_node (profondeur+1) d3
              {x=(separations.x -. epsilon);
               y=(separations.y -. epsilon)},
14           build_node (profondeur+1) d4
              {x=(separations.x +. epsilon);
               y=(separations.y -. epsilon)}) in
15  build_node 0 training {x=0.5; y=0.5} ;;

```

Le calcul de la prochaine découpe est effectué en fonction de la profondeur de l'arbre, en prenant les séparations actuelles et en prenant les 4 nouveaux couples de flottants créés en ajoutant ou enlevant $2^{-profondeur}$ à la coordonnée x ou y, par exemple les coordonnées de la découpe de profondeur 3 seront créées à partir de la découpe de profondeur 2 (0.5, 0.5) de la manière suivante :

$$(0.5 - 0.25, 0.5 + 0.25) = (0.25, 0.75)$$

$$(0.5 + 0.25, 0.5 + 0.25) = (0.75, 0.75)$$

$$(0.5 - 0.25, 0.5 - 0.25) = (0.25, 0.25)$$

$$(0.5 + 0.25, 0.5 - 0.25) = (0.75, 0.25)$$

4 Prédiction à l'aide de l'arbre de décision

On peut maintenant définir une fonction qui permet de prédire la couleur d'un point à l'aide de l'arbre de décision construit précédemment. On définit alors une fonction qui permet de naviguer dans l'arbre de décision en fonction des coordonnées du point à prédire. On renvoie la couleur de la feuille ainsi atteinte.

```

1 (* tree -> pos -> color *)
2 let rec predict_color tree pos =
3   match tree with
4   | Leaf color -> color
5   | Node (separations,a,b,c,d) ->
6     if pos.x < separations.x && pos.y >=
7       separations.y
8     then predict_color a pos
9     else if pos.x >= separations.x && pos.y >=
10       separations.y
11     then predict_color b pos
12     else if pos.x < separations.x && pos.y <
13       separations.y
14     then predict_color c pos
15     else predict_color d pos ;;

```

5 Résultats

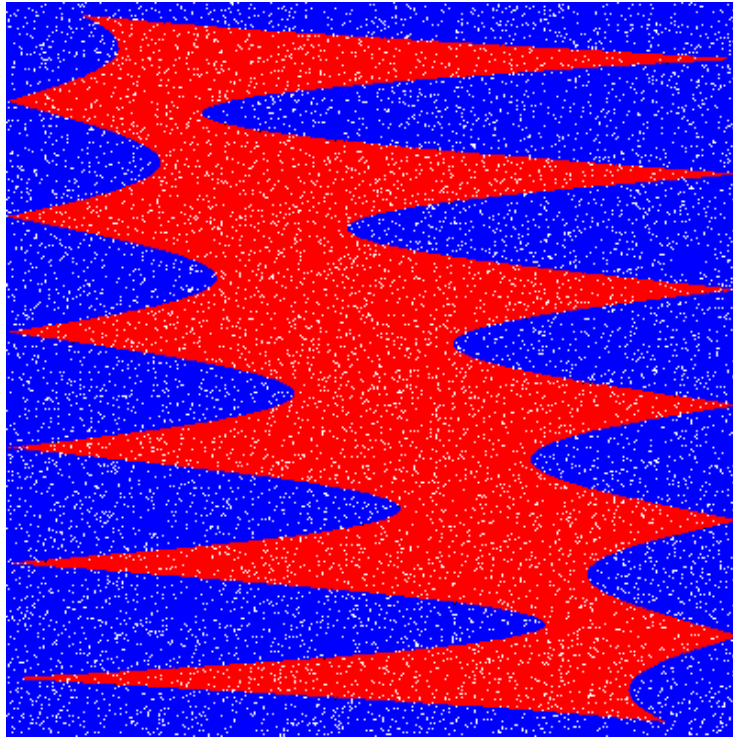
J'utiliserais pour la suite des tests trois fonctions de génération de datasets. Une fonction qui génère un cercle, une autre qui génère le graphe de la fonction graphe et une fonction qui crée un dataset complexe, voici son code.

```

1 let generate p =
2   if exp (-2.0 *. p.y) *. (abs_float (cos (20.0 *.
3     p.y))) >= p.x || 1.0 -. exp (-2.0 *. (1.0 -.
4     p.y)) *. (abs_float (cos (20.0 *. (1.0
5     -.p.y)))) < p.x
6   then blue else red ;;

```

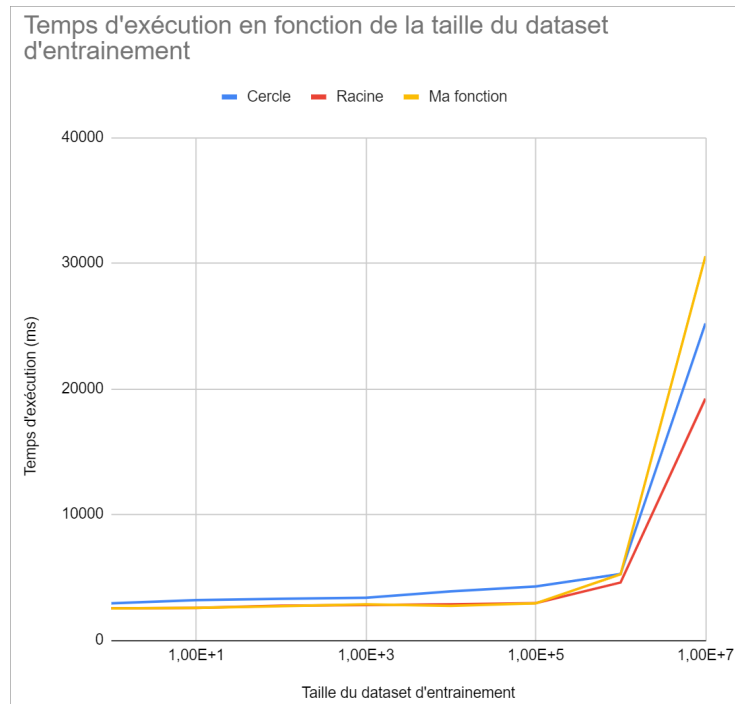
Cette fonction me semble interessante puisqu'elle genère un dataset avec une frontière de décision complexe.



On commence par mesurer complexité temporelle de notre programme en fonction de la taille du dataset. On génère un dataset de taille 10^n avec n allant de 0 à 6, et on mesure la durée d'exécution à l'aide de la commande Windows suivante :

1 `Measure-Command {ocamlrun test}`

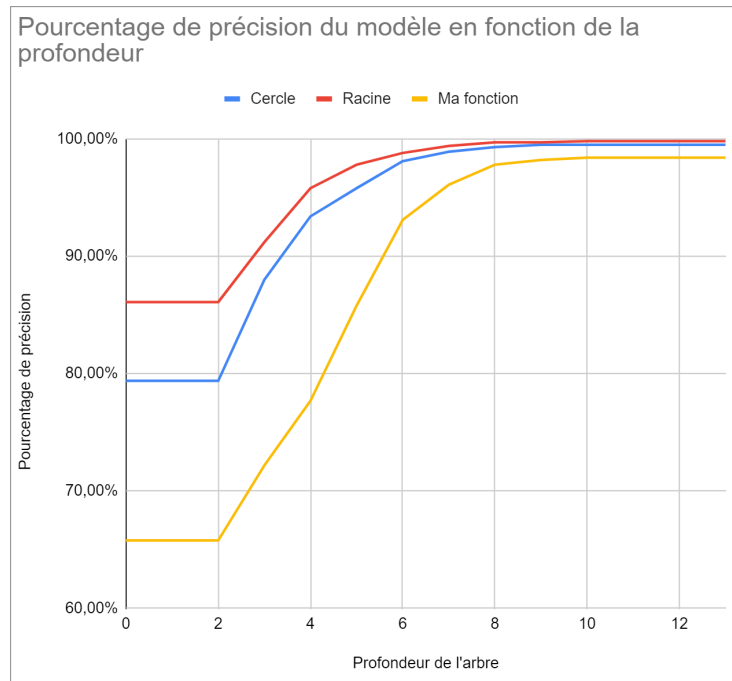
Ce qui nous donne le graphique suivant :



On constate que la complexité temporelle de notre programme est exponentielle en fonction de la taille du dataset. Cependant la durée d'exécution reste acceptable (quelques secondes) pour des datasets d'une taille inférieure à 10^7 données.

On peut ensuite tester notre programme sur un dataset de taille fixe 100000, et faire varier la profondeur maximale de l'arbre de décision pour mesurer la précision de celui-ci. On fait cela à l'aide de la commande Windows suivante :

```
1 ocamlrn test -md profondeur
```



Ce graphique permet de montrer que notre arbre de décision est pertinent puisque la précision de celui-ci augmente en fonction de la profondeur de l'arbre pour tendre vers 100% de précision. On peut aussi voir que la précision stagne pour des profondeurs supérieures à 10, ce qui montre que l'arbre de décision est optimal pour une profondeur de 10 (d'où le choix de la profondeur maximale par défaut de 10).