

# Arbres de Décision

*Octobre 2025*

Cédric Dangeard

[cedric.dangeard@orange.com](mailto:cedric.dangeard@orange.com)



# Business

# Cours 2 : Au programme

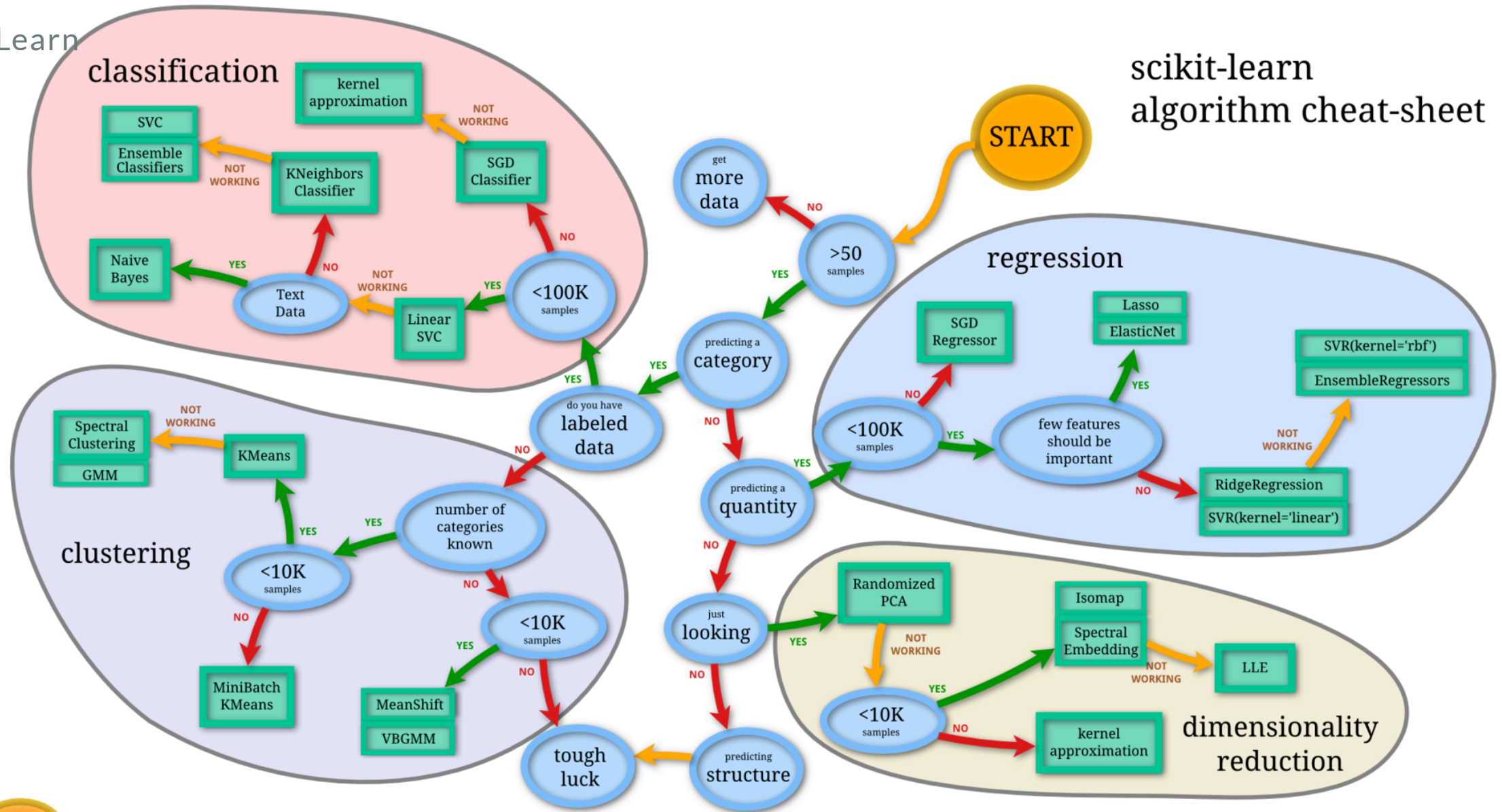
- Les bibliothèques fondamentales
- Nettoyage et Préparation des données
  - Gestion des valeurs manquantes
  - Encodage des variables catégorielles
  - Echantillonnage des données
- TP : Nettoyage des données
- Les arbres de décision
- TP : Implémentation d'un Arbre

# Les bibliothèques fondamentales

- Numpy
  - *alias* : Numerical Python
  - *objet de base* : `numpy.array`
  - *missions* : puissance et rapidité de calcul sur des vecteurs
  - *implémentation* : C
  - [code](#), [site](#)
- Scipy
  - *alias* : Scientific Python
  - *objet de base* : `numpy`
  - *missions* : Algorithmes plus haut niveau, optimisation, regression, interpolation, équations différentielles, ...
  - *implémentation* : C, Fortran, C++, Cython
  - [code](#), [site](#)

# Scikit-Learn

- Librairie de machine Learning sur Python
- Libre et Open source
- Orientée Objets
- Documentée
- Communauté active
- [code](#), [site](#)



- Estimateur

- fit : Entraîner le modèle
- transform : Transformer les données
- predict : Prédire la variable cible
- score : Évaluer les performances du modèle

```
from sklearn.base import BaseEstimator, TransformerMixin

class MyEstimator(BaseEstimator, TransformerMixin):
    def __init__(self, param : int = 1):
        self.param = 1

    def fit(self, X, y):
        self.is_fitted_ = True
        return self

    def transform(self, X):
        return X

    def predict(self, X):
        return np.full(shape=X.shape[0], fill_value=self.param)
```

# Encodage des variables catégorielles

Pourquoi encoder les variables catégorielles ?

- Les algorithmes de machine learning ne peuvent pas travailler directement avec des données non numériques.
- Les variables catégorielles doivent être converties en une représentation numérique pour être utilisées comme entrées dans les modèles.

Différentes méthodes d'encodage existent, chacune ayant ses avantages et inconvénients.

# Label Encoding

- Chaque modalité prend une valeur
- Peut être Ordinal
- Pas de nouvelles variables
- Perte d'information
- Que faire des valeurs manquantes ?

```
from sklearn.preprocessing import LabelEncoder, OrdinalEncoder

df = pd.DataFrame({
    'Name' : ['Sweet Mask', 'Bald Cape', 'Blizzard of Hell',
             'King', 'Glasses', 'Metal Bat', 'Mumen Rider'],
    'Class' : ['A', 'C', 'B', 'S', 'A', 'S', 'C']})

le = LabelEncoder()
oe = OrdinalEncoder(categories=[['S', 'A', 'B', 'C']])

df['label'] = le.fit_transform(df[['Class']])
df['ord'] = oe.fit_transform(df[['Class']])
```

	Sweet Mask	Bald Cape	Blizzard of Hell	King	Glasses	Metal Bat	Mumen Rider
Name	Sweet Mask	Bald Cape	Blizzard of Hell	King	Glasses	Metal Bat	Mumen Rider
Class	A	C	B	S	A	S	C
label	0	2	1	3	0	3	2
ord	1.0	3.0	2.0	0.0	1.0	0.0	3.0



# One Hot Encoding

- Une nouvelle variable par modalité
- Chaque variable prend la valeur 0 ou 1
- Peut ajouter beaucoup de variables
- Deux façon de faire :
  - `pandas.get_dummies`
  - `sklearn.OneHotEncoder`

```
from sklearn.preprocessing import OneHotEncoder

ohe = OneHotEncoder(categories=[['S', 'A', 'B', 'C']],
                    drop = 'first',
                    sparse_output=False)

hot_encoded = ohe.fit_transform(df[['Class']])

data_hot_encoded = pd.DataFrame(hot_encoded, index=df.index)
data_hot_encoded.columns = ohe.get_feature_names_out()

pd.concat([df, data_hot_encoded], axis=1)
```

```
pd.concat([df,
          pd.get_dummies(df[['Class']], drop_first=True)], axis=1)
```

Name	Sweet Mask	Bald Cape	Blizzard of Hell	King	Glasses	Metal Bat	Mumen Rider
Class	A	C	B	S	A	S	C
Class_A	1.0	0.0	0.0	0.0	1.0	0.0	0.0
Class_B	0.0	0.0	1.0	0.0	0.0	0.0	0.0
Class_C	0.0	1.0	0.0	0.0	0.0	0.0	1.0

# Target Encoding

- Remplacer les modalités par la valeur dérivée de la valeur à prédire (ex : la moyenne)
  - Peut créer du sur-apprentissage
  - A utiliser sur des modalités à hautes cardinalités

```
df = pd.DataFrame({
    'Name' : ['Le Bron', 'Kawamura', 'Wembanya', 'Durant',
             'Curry', 'Jokić', 'Paul', 'Davis'],
    'Taille' : [206, 173, 224, 211, 188, 211, 183, 208],
    'Club' : ['Lakers', 'Grizzlies', 'Spurs', 'Suns',
             'Warriors', 'Nuggets', 'Spurs', 'Lakers',]
})

df['encoded_Club'] = df.groupby('Club')['Taille'].transform('mean')
```

Name	Le Bron	Kawamura	Wembanya	Durant	Curry	Jokić	Paul	Davis
Taille	206	173	224	211	188	211	183	208
Club	Lakers	Grizzlies	Spurs	Suns	Warriors	Nuggets	Spurs	Lakers
encoded_Club	207.0	173.0	203.5	211.0	188.0	211.0	203.5	207.0

```
from sklearn.preprocessing import TargetEncoder

te = TargetEncoder(smooth='auto',
                  cv=2)
target_encoded = te.fit_transform(X=df[['category']], y = df[['target']])
```

# Echantillonnage des données

- Souvent les données sont déséquilibrées
- Plusieurs techniques pour rééquilibrer les classes
  - Sur-échantillonnage (Over-sampling)
  - Sous-échantillonnage (Under-sampling)
  - Synthèse de nouvelles données (SMOTE, ADASYN, ...)
- [Imbalanced-learn](#) : Librairie Python pour le rééchantillonnage

# TP : Encodage

Des Questions  
avant de s'y mettre?

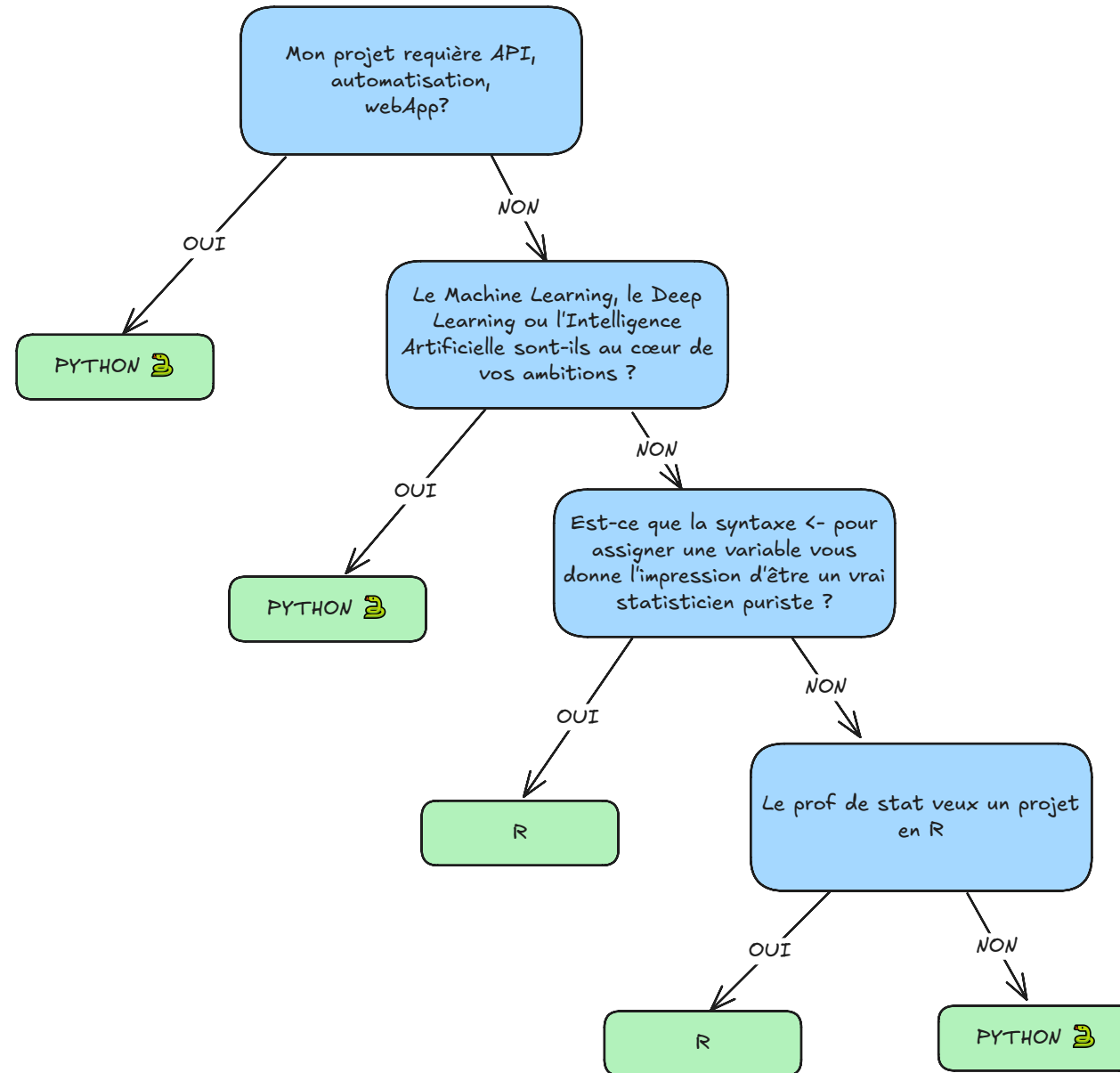


# Arbres de décisions

Qu'est ce qu'un arbre de décision ?

# Vocabulaire sur les arbres

- **Nœud racine** : Le nœud supérieur de l'arbre, qui ne possède pas de parent.
- **Nœud de décision** : Un nœud à partir duquel on divise les données en fonction d'une caractéristique.
- **Feuille** (ou **nœud terminal**) : Un nœud qui ne se divise pas davantage et représente une prédiction ou une classe.
- **Branche** : Une connexion entre deux nœuds, représentant le flux de décision.



# Arbre de décision

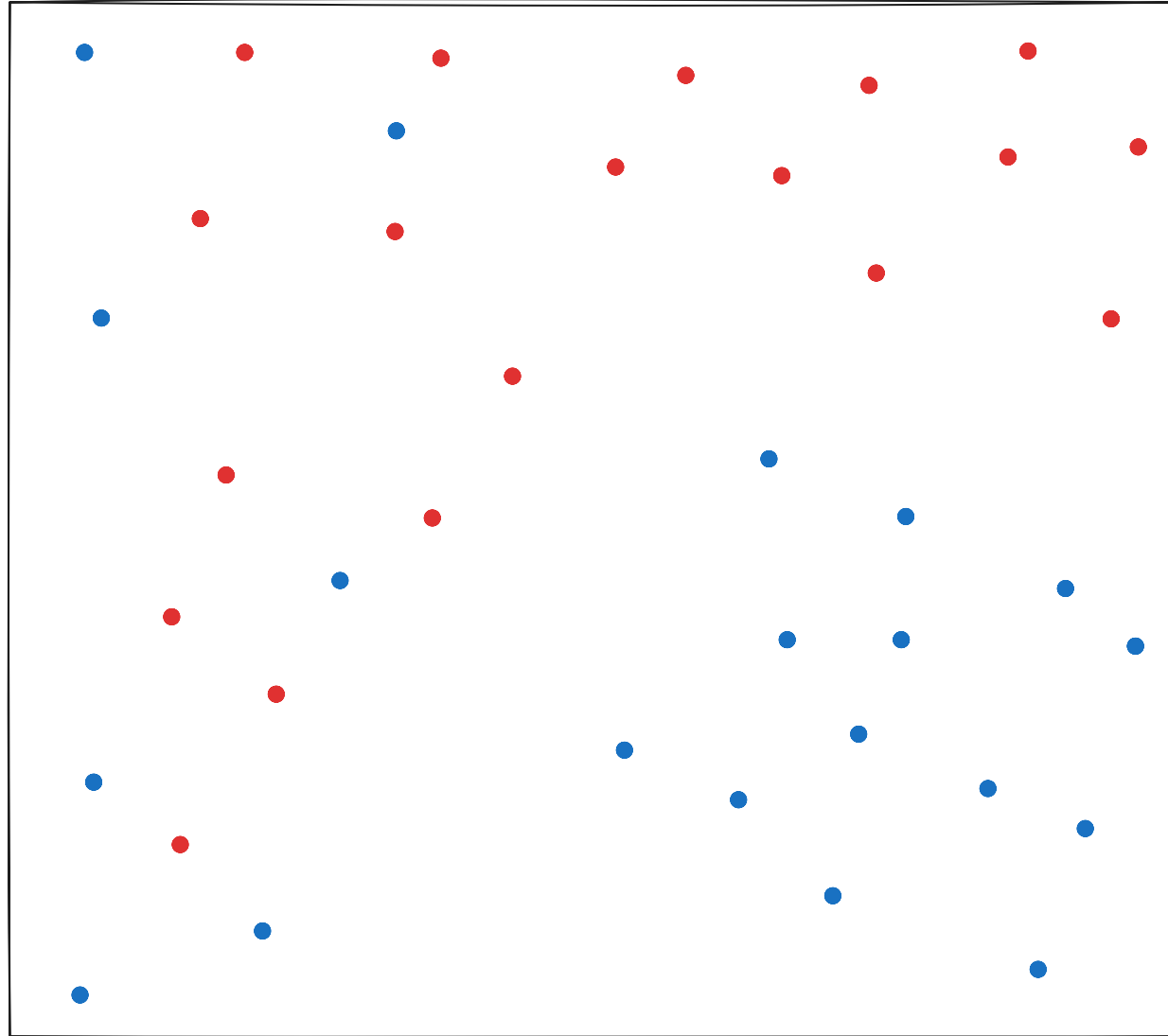
- Modèle de classification et de régression
- Structure arborescente composée de nœuds
- Chaque nœud représente une décision basée sur une caractéristique
- Chaque branche représente le résultat d'une décision
- Chaque feuille représente une prédiction ou une classe

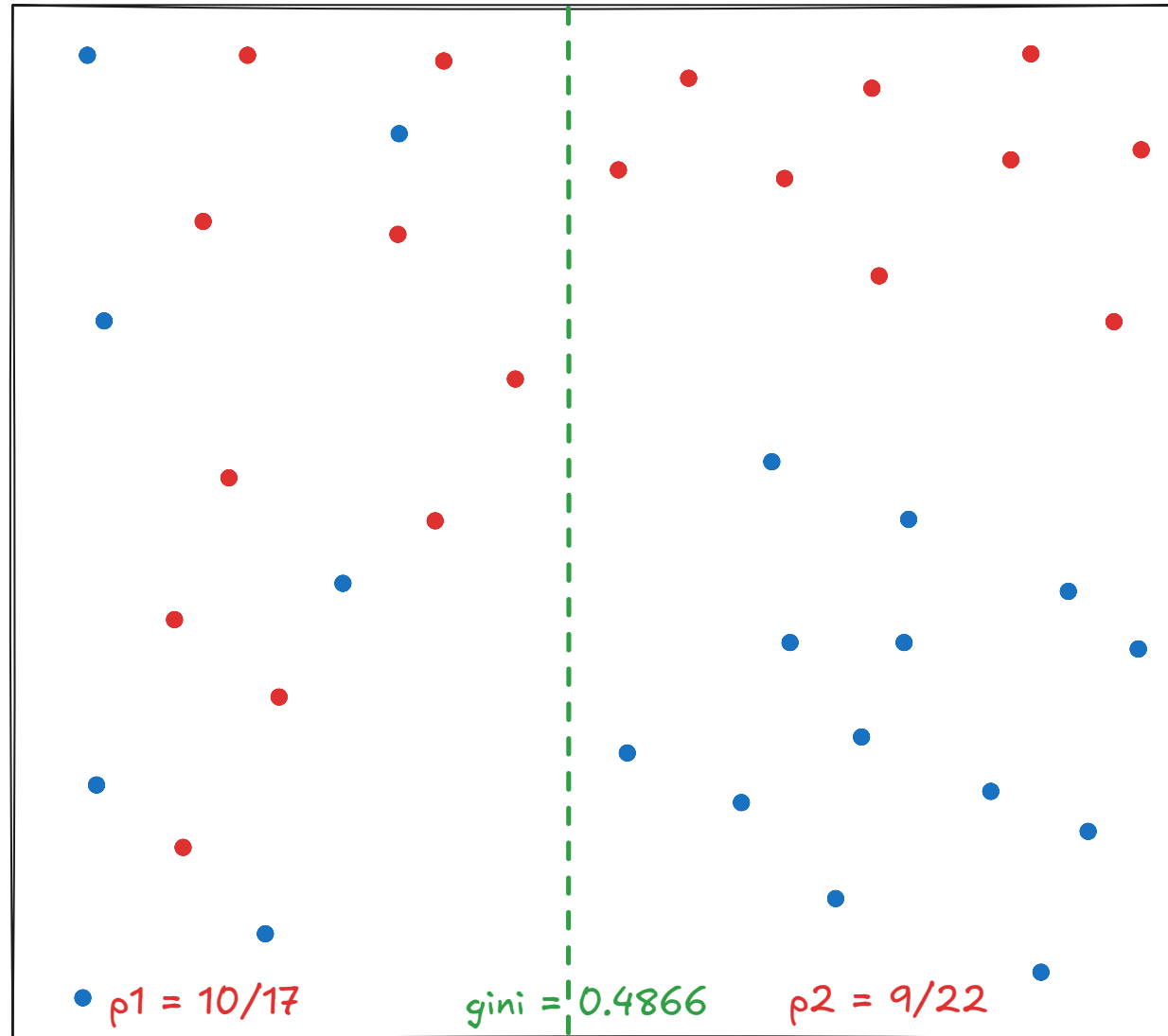


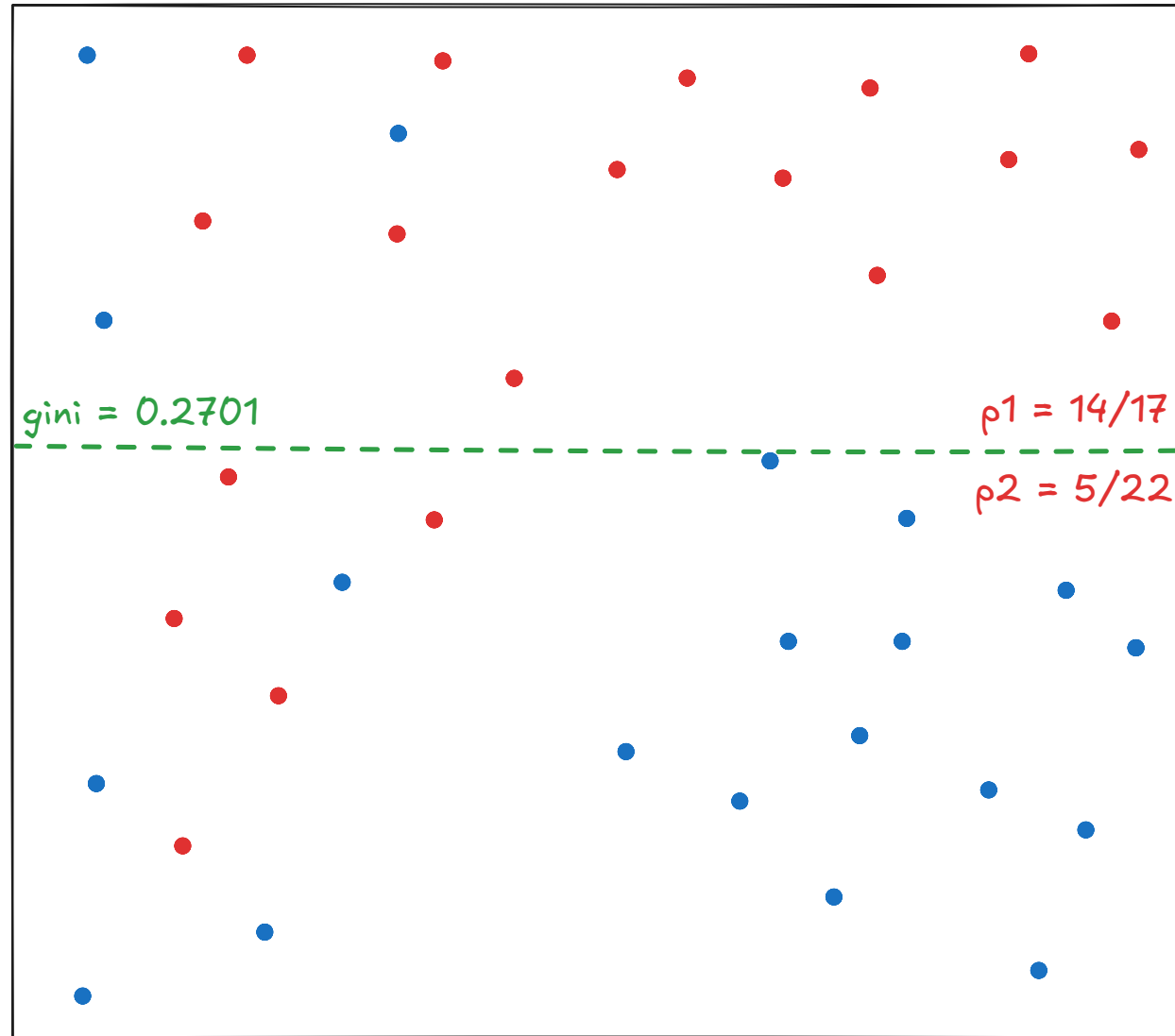
# Construction d'un arbre de décision

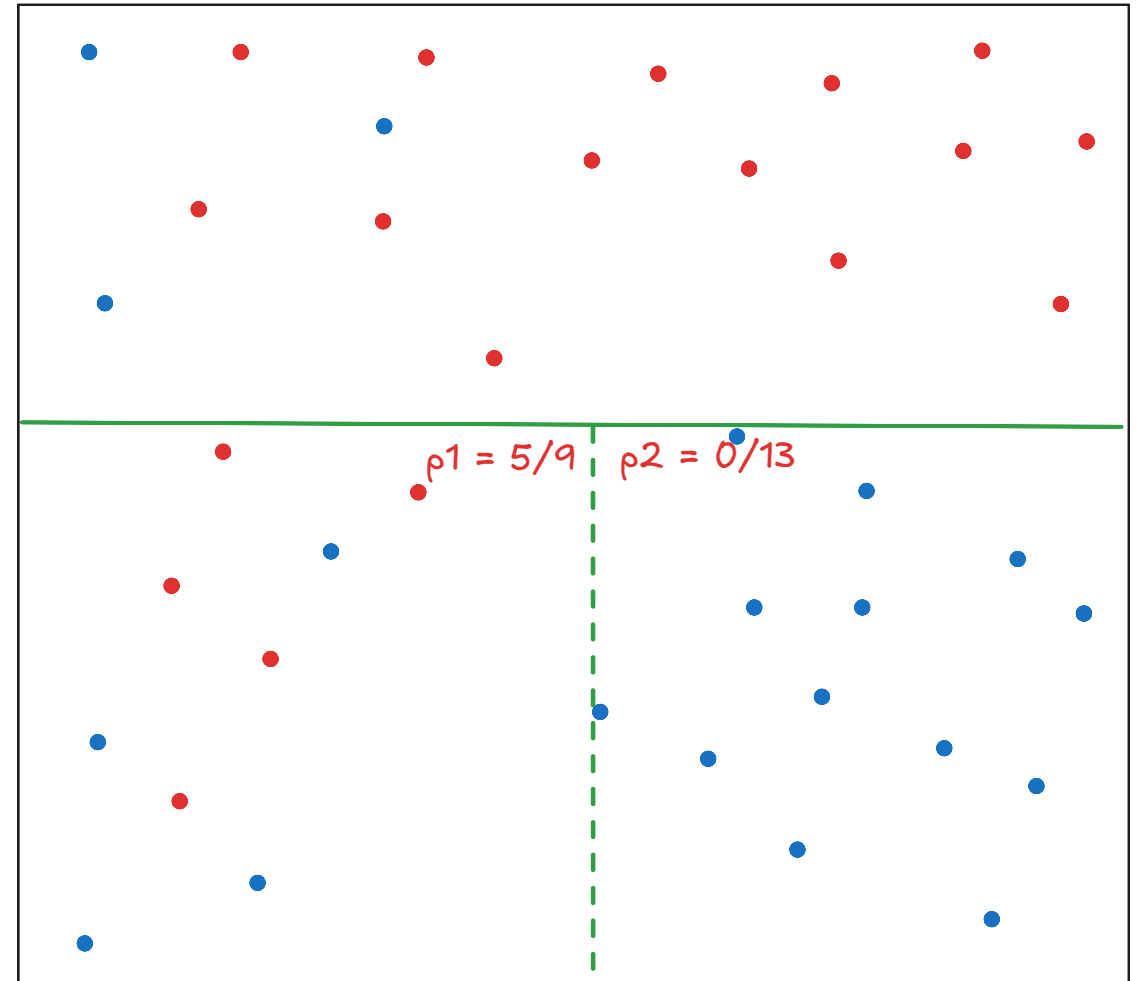
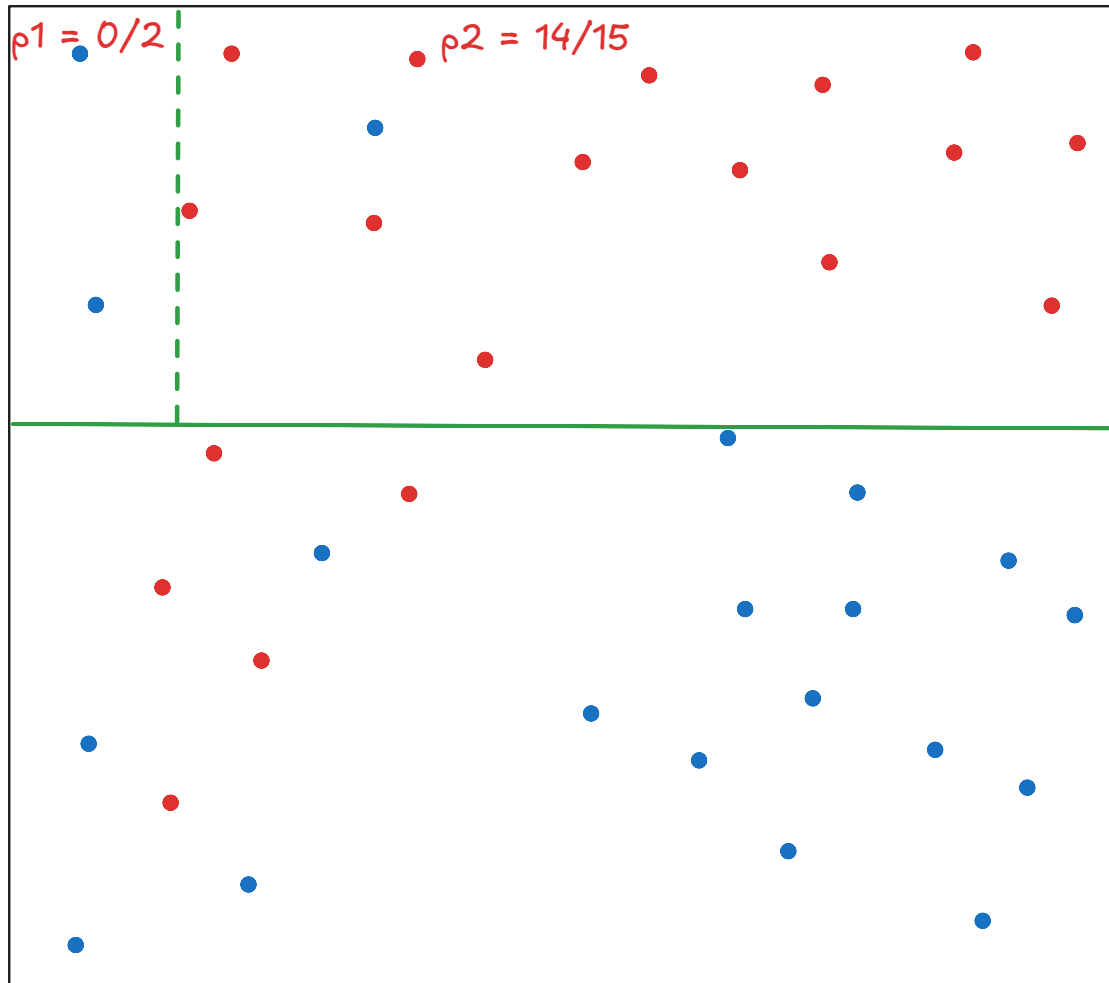
La construction se fait par *partitionnement récursif* des données.

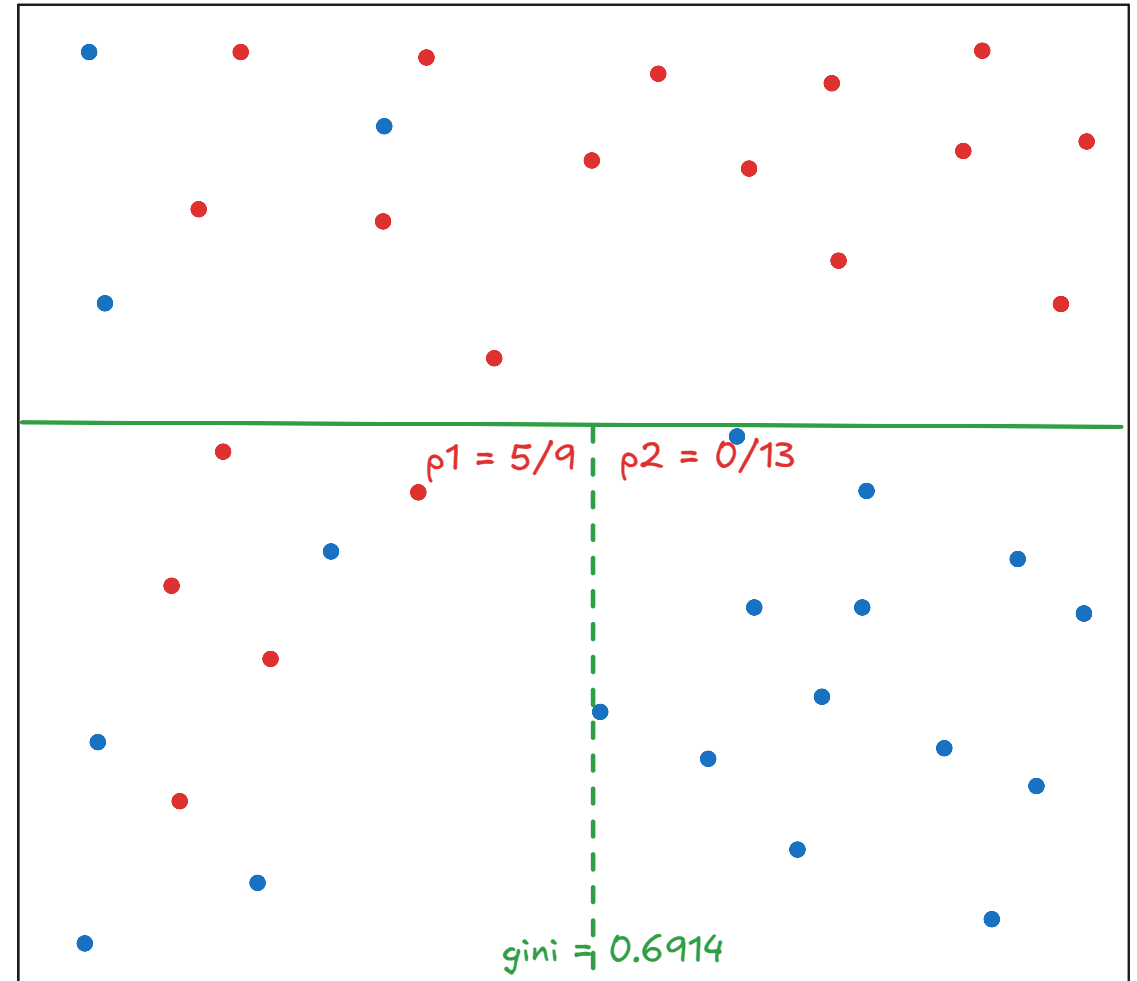
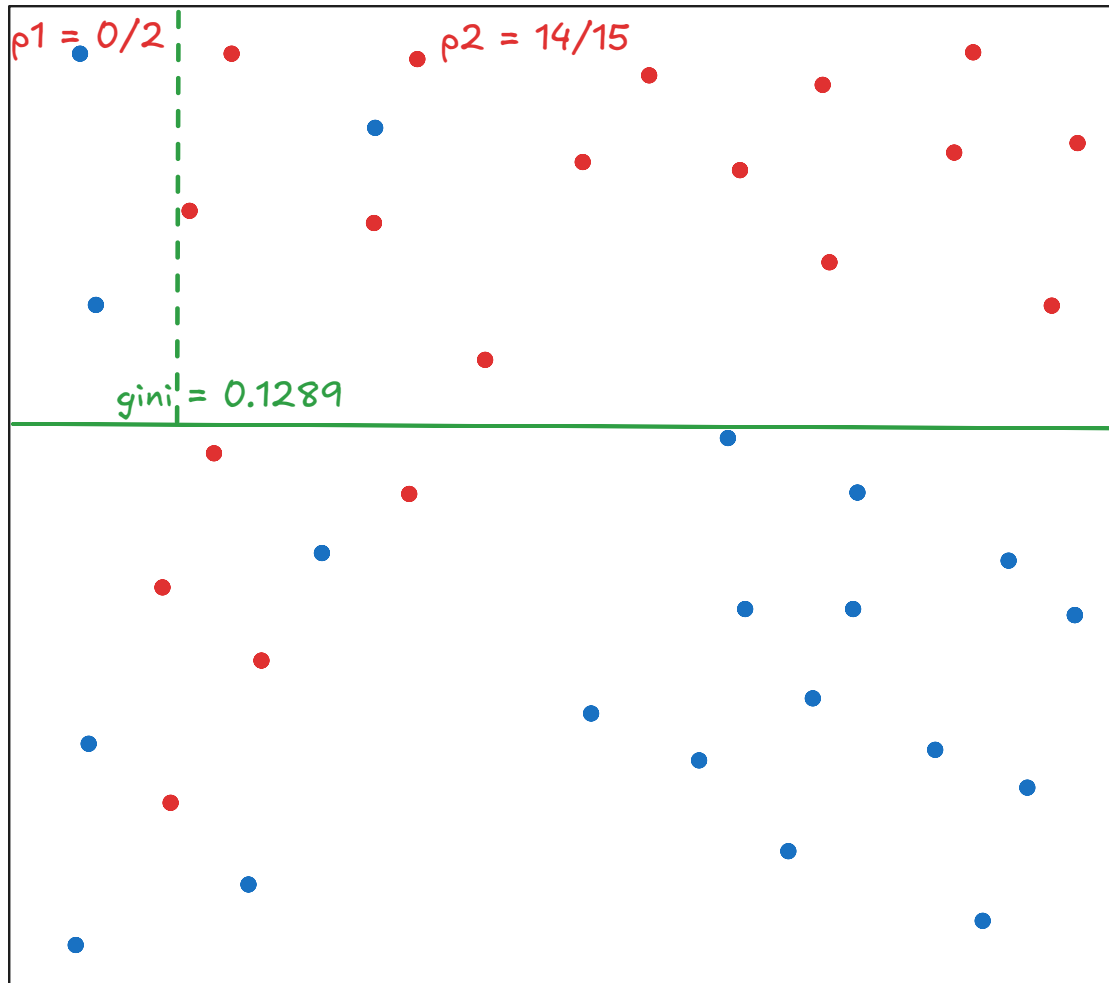
1. Choisir la meilleure division pour diviser les données.
2. Diviser les données en sous-ensembles basés sur cette caractéristique.
3. Répéter le processus pour chaque sous-ensemble jusqu'à ce qu'un critère d'arrêt soit atteint

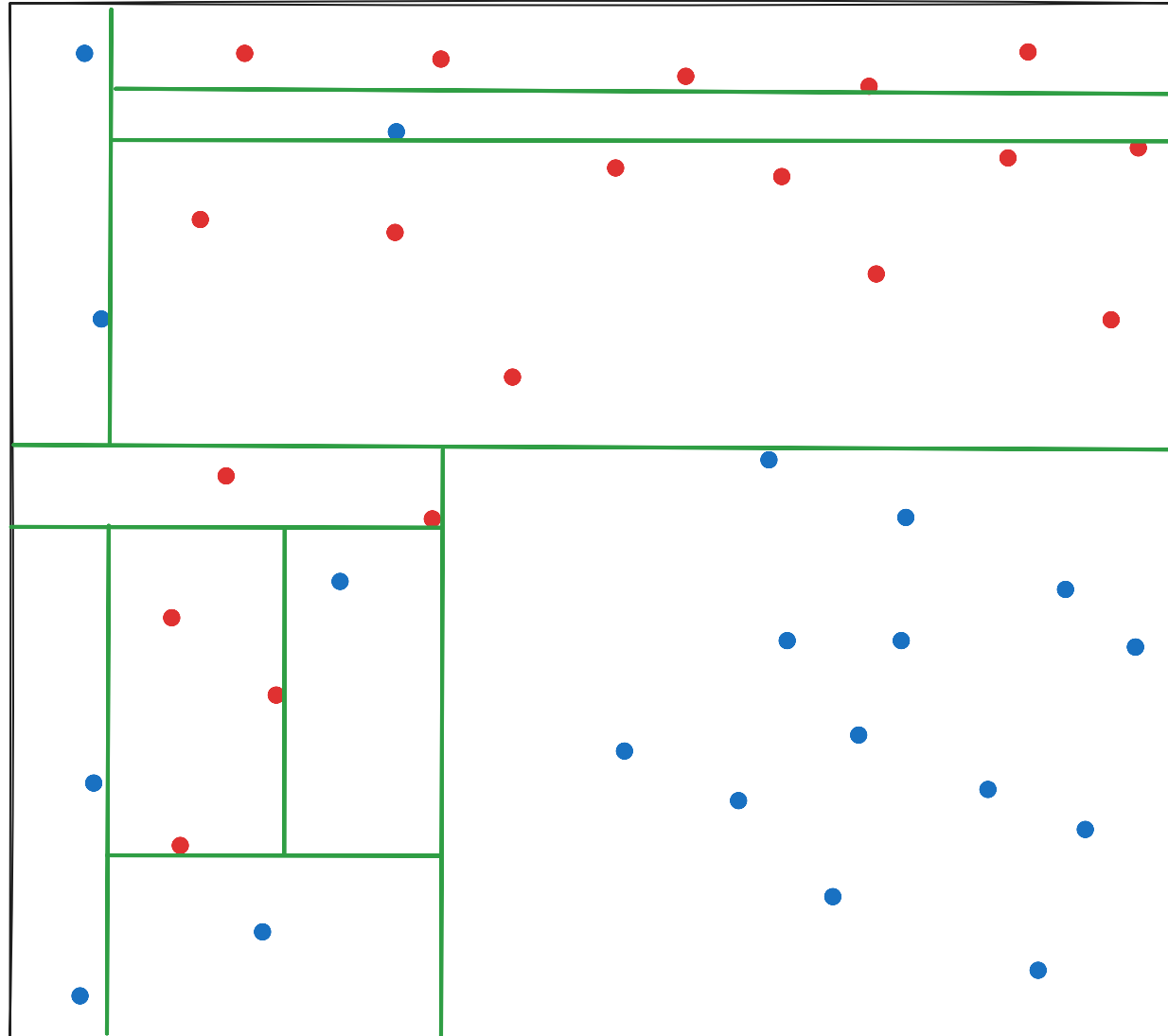




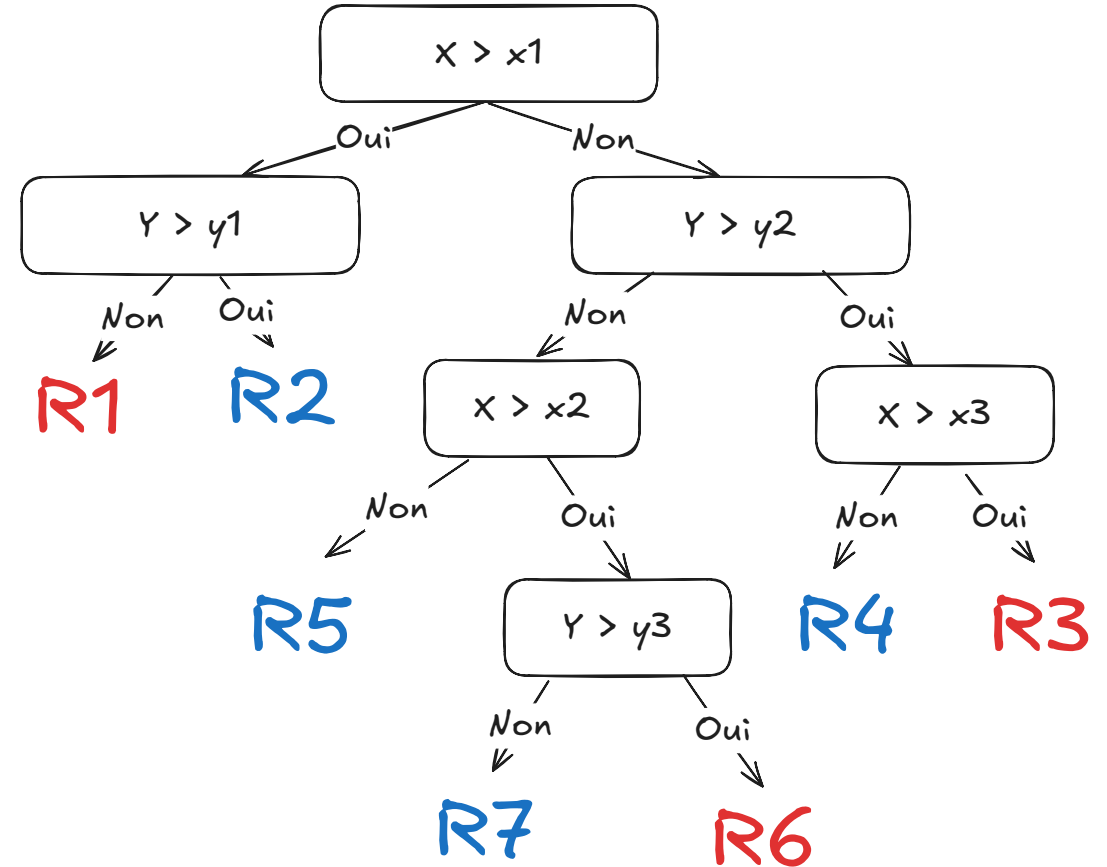
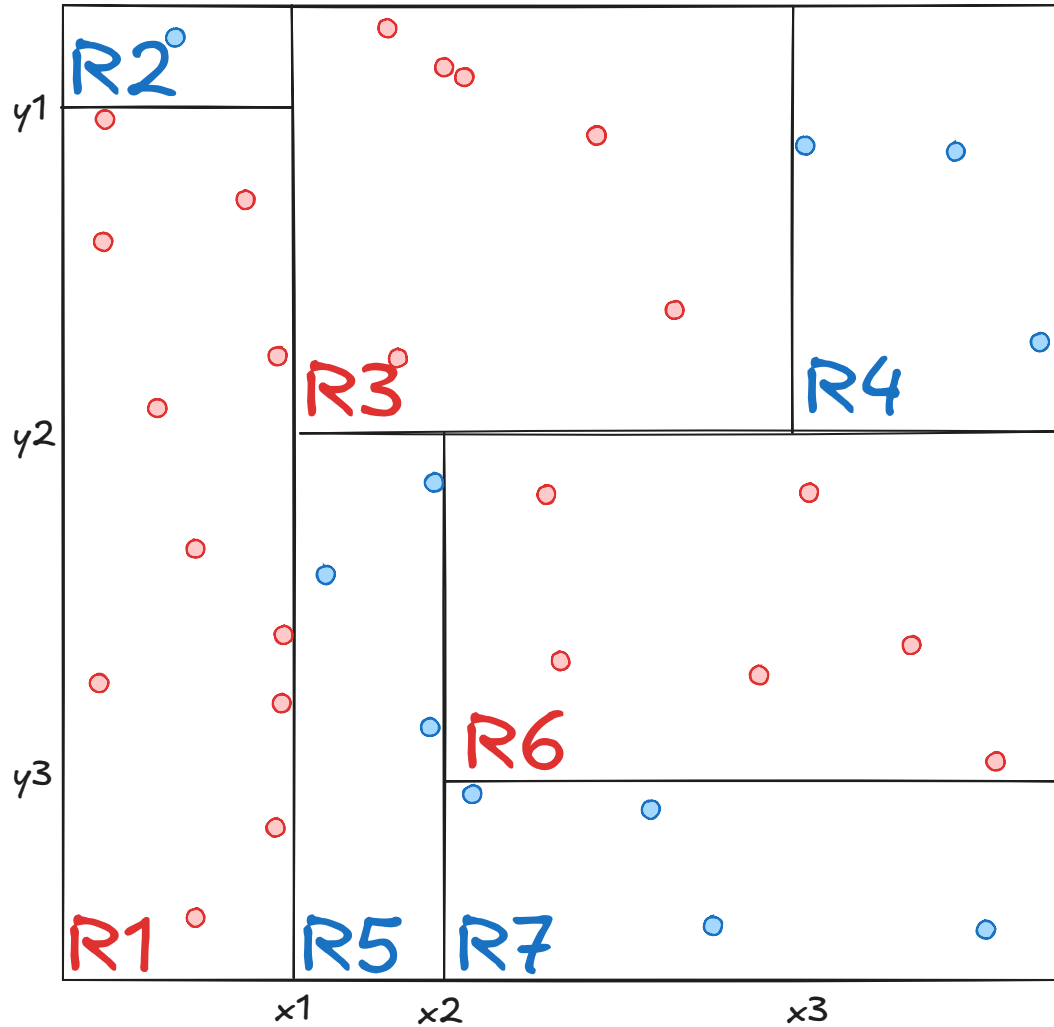








# Arbres de décisions





# Comment choisir la meilleure division ?

- Quelle variable choisir ?
- Quel seuil de division ?

# CART : L'algorithme pour un nœud

- Définir un critère d'homogénéité
  - CART : Indice de gini :  $I_G = 1 - \sum_{i=1}^m f_i^2$   
où  $f_i$  est la proportion de classe  $i$  dans le jeu de données.
  - Autres : [Entropie de Shanon](#), Gain d'information, ...
- Calculer ce critère pour un ensemble de segmentations des données.
- Choisir la segmentation qui minimise le critère.

# CART : Où s'arreter ?

Un individu par feuille ?

Jamais le même des feuilles mêmes classes ?

Toutes des feuilles, quelques des feuilles ?

Profondeur de l'arbre ?

# **CART : Où s'arreter ? Que cherche t-on à éviter?**

- Un individu par feuille ?
- Une classe par feuille ?
- Taille des feuilles, nombre de feuilles ?
- Profondeur de l'arbre ?

Le choix de ce critère d'arrêt, va constituer l'un des hyperparamètres à optimiser de notre modèle.

# CART : Prunning

Plus un arbre est profond, plus la variance est élevée, et le biais est faible.

- Élagage (*Prunning*)
  - Pré-élagage :
    - Instaurer des règles d'arrêt pendant l'apprentissage
  - Post-élagage :
    - Partir d'un arbre profond et réduire la variance en supprimant des feuilles.

# Arbres de décisions : Conclusion

- Facilement interprétable
- Variance élevé
- Sur-apprentissage pour les arbres trop profonds
- Peu performant en général

## Solutions :

- Comment améliorer les performances des arbres ?
- Comment réduire la variance ?

# CART sur SkLearn : DecisionTreeClassifier

- Paramètres :
  - **criterion**
  - **max\_depth**
  - **min\_samples\_split**
  - **min\_samples\_leaf**
  - **max\_features**

# CART sur SkLearn : DecisionTreeClassifier

- Paramètres :
  - **criterion** [*Gini, entropy, logloss*]: Critères d'homogénéité
  - **max\_depth** *int*: Profondeur de l'arbre
  - **min\_samples\_split** *int*: Individus minimum dans la feuille pour procéder à un split.
  - **min\_samples\_leaf** *int*: Nombre d'individus minimum dans les feuilles filles pour accepter le split
  - **max\_features** [*int, float, 'sqrt', 'log2', None*]: Nombre maximum de variables à tester pour créer un noeud



# TP : Arbres de décision