

Rapport Technique

Implémentation de l'Algorithmme CART

Classification and Regression Trees de Breiman (1984)

Amlal Amine - Boughnam Houda - Khald Adam

17 décembre 2025

Résumé

Ce rapport présente une analyse complète de l'algorithme CART (Classification and Regression Trees) développé par Leo Breiman et ses collaborateurs en 1984. L'algorithme CART représente une contribution majeure au domaine de l'apprentissage automatique et des statistiques. Ce document couvre les fondements théoriques, les aspects mathématiques, l'implémentation pratique ainsi que l'analyse détaillée de l'implémentation open-source disponible sur le repository GitHub [Adamkhald/breiman-cart](#). Le rapport inclut des pseudo-codes détaillés, des formulations mathématiques rigoureuses et des exemples d'application pratiques.

Table des matières

1	Introduction	4
1.1	Contexte Historique	4
1.2	Objectifs du Rapport	4
1.3	Structure du Repository GitHub	4
2	Fondements Théoriques	5
2.1	Principe Général	5
2.2	Définitions Mathématiques	5
2.2.1	Espace de Données	5
2.2.2	Structure d'Arbre	5
2.3	Règles de Division	5
2.3.1	Variables Continues	5
2.3.2	Variables Catégorielles	6
3	Critères de Division	6
3.1	Classification : Indice de Gini	6
3.2	Classification : Entropie	6
3.3	Régression : Erreur Quadratique Moyenne	7
4	Algorithmes et Pseudo-Codes	7
4.1	Algorithme Principal de Construction	7
4.2	Recherche de la Meilleure Division	8
4.3	Calcul du Gain (Classification)	9

4.4	Calcul du Gain (Régression)	9
5	Élagage (Pruning)	10
5.1	Principe de l'Élagage par Coût-Complexité	10
5.2	Séquence d'Élagage	10
5.3	Algorithme d'Élagage	11
6	Critères d'Arrêt	11
6.1	Paramètres de Contrôle	11
6.2	Algorithme de Vérification	12
7	Prédiction	13
7.1	Algorithme de Prédiction	13
8	Analyse Détailée de l'Implémentation GitHub	13
8.1	Architecture du Code	13
8.1.1	Structure des Fichiers	13
8.2	Nouvelle API Unifiée (Version 0.2.0)	14
8.2.1	BRCRegression - Régression	14
8.2.2	BRCClassification - Classification	15
8.2.3	BRCInference - Inférence avec Modèles Sauvegardés	16
8.3	Caractéristiques Principales	16
8.3.1	Gestion Native des Variables Catégorielles	16
8.3.2	Élagage par Coût-Complexité	17
8.3.3	Type Safety	18
8.4	Exemple de Régression	18
8.5	Gestion des Modèles et Persistance	19
8.5.1	Sauvegarde et Chargement avec Pickle	19
8.5.2	Export en Format JSON	19
8.5.3	Génération de Résumés	19
8.6	Logging et Débogage	20
8.7	Differences avec Scikit-Learn	20
8.8	Installation et Utilisation	20
8.8.1	Installation	20
8.8.2	Tests	21
8.9	Complexité Temporelle	21
8.9.1	Construction de l'Arbre	21
8.9.2	Élagage	21
8.9.3	Prédiction	21
8.10	Complexité Spatiale	22
8.11	Optimisations Possibles	22
9	Applications et Cas d'Usage	22
9.1	Domaines d'Application	22
9.1.1	Médecine et Santé	22
9.1.2	Finance	22
9.1.3	Marketing et E-commerce	23
10	Avantages et Limitations	23

10.1 Avantages de CART	23
10.2 Limitations de CART	23
11 Conclusion	24
11.1 Contributions de CART	24
11.2 Apports du Repository GitHub	24
11.3 Perspectives Futures	25
11.4 Recommandations Pratiques	25
11.5 Impact et Héritage	25
A Glossaire des Notations	26
B Formules Récapitulatives	27
B.1 Critères d’Impureté	27
B.2 Gain de Division	27
B.3 Coût-Complexité	27
B.4 Alpha Critique	27
B.5 Décomposition Biais-Variance	28

1 Introduction

1.1 Contexte Historique

L'algorithme CART (Classification and Regression Trees) a été introduit en 1984 par Leo Breiman, Jerome Friedman, Richard Olshen et Charles Stone dans leur ouvrage fondateur *Classification and Regression Trees*. Cette méthode révolutionnaire a transformé l'analyse de données en proposant une approche non-paramétrique pour la classification et la régression.

Leo Breiman, statisticien de renom à l'Université de Californie à Berkeley, a également apporté des contributions majeures avec le *bagging* (Bootstrap Aggregating) en 1996 et les *Random Forests* en 2001, qui sont des extensions directes de CART.

1.2 Objectifs du Rapport

Ce rapport vise à :

- Présenter les fondements théoriques de l'algorithme CART
- Développer les formulations mathématiques rigoureuses
- Fournir des pseudo-codes détaillés pour l'implémentation
- Analyser en profondeur l'implémentation du repository GitHub [Adamkhald/breiman-cart](#)
- Discuter des applications pratiques et des considérations d'implémentation
- Comparer avec d'autres implémentations notamment scikit-learn

1.3 Structure du Repository GitHub

Le repository [Adamkhald/breiman-cart](#) a été récemment restructuré (Version 0.2.0) pour offrir une API plus conviviale et une meilleure expérience utilisateur :

- **breiman_cart.py** : API unifiée principale (2000+ lignes)
 - **BRCRegression** : Classe pour tâches de régression
 - **BRCClassification** : Classe pour tâches de classification
 - **BRCIInference** : Classe pour inférence avec modèles entraînés
- **breiman_cart/** : Package interne modulaire
 - **__init__.py** : Initialisation du package
 - **node.py** : Structure de noeuds avec support catégoriel
 - **splitter.py** : Logique de division (Gini, MSE, sous-ensembles)
 - **tree.py** : Implémentation interne des arbres
 - **pruning.py** : Élagage par coût-complexité
- **tests/** : Suite de tests unitaires
- **docs/** : Documentation théorique (**theory.tex**)
- **demo.py** : Script de démonstration
- **quickstart.py** : Guide de démarrage rapide
- **test_new_api.py** : Tests pour la nouvelle API
- **fix_syntax.py** : Utilitaire de correction syntaxique
- **requirements.txt** : Dépendances (NumPy, Pandas, pickle, json)

2 Fondements Théoriques

2.1 Principe Général

CART est un algorithme de partitionnement récursif qui construit un arbre de décision binaire. Le principe fondamental repose sur :

1. **Partitionnement récursif** : Division successive de l'espace des variables explicatives en régions homogènes
2. **Modélisation locale** : Estimation de modèles simples (constantes) dans chaque région de la partition
3. **Structure binaire stricte** : Chaque nœud génère exactement deux branches (caractéristique distinctive de CART)
4. **Approche gloutonne** : Sélection locale optimale à chaque division sans rétroaction

2.2 Définitions Mathématiques

2.2.1 Espace de Données

Soit un ensemble de n observations :

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \quad (1)$$

où :

- $\mathbf{x}_i \in \mathbb{R}^p$: vecteur de p variables explicatives (features)
- $y_i \in \mathcal{Y}$: variable réponse (target)
- Classification : $\mathcal{Y} = \{1, 2, \dots, K\}$ (classes discrètes)
- Régression : $\mathcal{Y} = \mathbb{R}$ (valeur continue)

2.2.2 Structure d'Arbre

Un arbre T est défini par :

- $\mathcal{N}(T)$: ensemble des nœuds internes (nœuds de décision)
- \tilde{T} : ensemble des nœuds terminaux (feuilles)
- $|\tilde{T}|$: nombre de feuilles (mesure de complexité de l'arbre)
- Pour chaque nœud interne : une règle de division (j, s) où j est l'indice de la variable et s le seuil
- Pour chaque feuille : une prédiction constante

2.3 Règles de Division

2.3.1 Variables Continues

Pour une variable continue X_j , une division au point s crée deux régions :

$$R_L(j, s) = \{\mathbf{x} : x_j \leq s\} \quad (\text{région gauche}) \quad (2)$$

$$R_R(j, s) = \{\mathbf{x} : x_j > s\} \quad (\text{région droite}) \quad (3)$$

Recherche du seuil optimal : Pour une variable continue, les points de division candidats sont les points médians entre deux valeurs consécutives distinctes :

$$s_k = \frac{x_{(k)} + x_{(k+1)}}{2}, \quad k = 1, \dots, m - 1 \quad (4)$$

où $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(m)}$ sont les valeurs distinctes triées.

2.3.2 Variables Catégorielles

Pour une variable catégorielle X_j avec catégories $\mathcal{C} = \{c_1, c_2, \dots, c_m\}$, la division utilise des sous-ensembles :

$$R_L(j, S) = \{\mathbf{x} : x_j \in S\} \quad (5)$$

$$R_R(j, S) = \{\mathbf{x} : x_j \in \mathcal{C} \setminus S\} \quad (6)$$

où $S \subset \mathcal{C}$.

Note importante : Le nombre de divisions possibles pour m catégories est $2^{m-1} - 1$. Pour la classification binaire, Breiman a montré qu'on peut réduire cette complexité en ordonnant les catégories selon leur proportion de classe positive.

3 Critères de Division

3.1 Classification : Indice de Gini

L'impureté de Gini pour un nœud t est définie par :

$$I_{\text{Gini}}(t) = \sum_{k=1}^K p_k(t) \cdot (1 - p_k(t)) = 1 - \sum_{k=1}^K p_k^2(t) \quad (7)$$

où $p_k(t)$ est la proportion d'observations de classe k dans le nœud t :

$$p_k(t) = \frac{1}{n_t} \sum_{i \in t} \mathbb{1}(y_i = k) \quad (8)$$

Interprétation : L'indice de Gini mesure la probabilité de mal classer un élément choisi aléatoirement si on lui attribue une classe aléatoire selon la distribution des classes dans le nœud.

Propriétés :

- $I_{\text{Gini}}(t) = 0$ si le nœud est pur (toutes les observations de même classe)
 - $I_{\text{Gini}}(t)$ est maximale pour une distribution uniforme
 - Pour $K = 2$: $I_{\text{Gini}}(t) = 2p(1 - p)$ où p est la proportion de la classe positive
- La réduction d'impureté pour une division est :

$$\Delta I(s, t) = I(t) - \frac{n_L}{n_t} I(t_L) - \frac{n_R}{n_t} I(t_R) \quad (9)$$

3.2 Classification : Entropie

L'entropie de Shannon peut également être utilisée :

$$I_{\text{Entropy}}(t) = - \sum_{k=1}^K p_k(t) \log_2 p_k(t) \quad (10)$$

Comparaison Gini vs Entropie :

- Les deux critères produisent généralement des arbres similaires
- L'entropie est légèrement plus coûteuse à calculer (logarithmes)
- Gini favorise les divisions qui isolent la classe majoritaire
- Entropie favorise les divisions équilibrées

3.3 Régression : Erreur Quadratique Moyenne

Pour la régression, l'objectif est de minimiser la somme des carrés résiduelle (RSS) :

$$\text{RSS}(T) = \sum_{j=1}^{|T|} \sum_{i \in R_j} (y_i - \hat{y}_j)^2 \quad (11)$$

où \hat{y}_j est la prédiction dans la région R_j :

$$\hat{y}_j = \frac{1}{n_j} \sum_{i \in R_j} y_i \quad (\text{moyenne des valeurs cibles}) \quad (12)$$

La réduction de variance pour une division est :

$$\Delta \text{Var}(s, t) = \text{Var}(t) - \frac{n_L}{n_t} \text{Var}(t_L) - \frac{n_R}{n_t} \text{Var}(t_R) \quad (13)$$

avec :

$$\text{Var}(t) = \frac{1}{n_t} \sum_{i \in t} (y_i - \bar{y}_t)^2 \quad (14)$$

4 Algorithmes et Pseudo-Codes

4.1 Algorithme Principal de Construction

Algorithm 1 Construction d'Arbre CART

```

1: procedure BUILD_CART( $\mathcal{D}$ , params)
2:    $T \leftarrow \text{CreateRootNode}(\mathcal{D})$ 
3:   queue  $\leftarrow [T.\text{root}]$ 
4:   while queue  $\neq \emptyset$  do
5:      $t \leftarrow \text{queue.pop}()$ 
6:     if StoppingCriterion( $t$ , params) then
7:       MakeLeaf( $t$ )
8:       continue
9:     end if
10:     $(j^*, s^*) \leftarrow \text{FindBestSplit}(t)$ 
11:    if  $(j^*, s^*) = \text{None}$  then
12:      MakeLeaf( $t$ )
13:      continue
14:    end if
15:     $(t_L, t_R) \leftarrow \text{Split}(t, j^*, s^*)$ 
16:    queue.append( $t_L$ )
17:    queue.append( $t_R$ )
18:  end while
19:  return  $T$ 
20: end procedure

```

4.2 Recherche de la Meilleure Division

Algorithm 2 Recherche de la Meilleure Division

```

1: procedure FINDBESTSPLIT( $t$ )
2:   best_gain  $\leftarrow -\infty$ 
3:   best_split  $\leftarrow \text{None}$ 
4:   for  $j = 1$  to  $p$  do                                 $\triangleright$  Pour chaque variable
5:     if  $X_j$  est continue then
6:       candidates  $\leftarrow \text{GetUniqueSortedValues}(X_j, t)$ 
7:       for  $s \in \text{candidates}$  do
8:         gain  $\leftarrow \text{ComputeGain}(t, j, s)$ 
9:         if gain > best_gain then
10:          best_gain  $\leftarrow \text{gain}$ 
11:          best_split  $\leftarrow (j, s)$ 
12:        end if
13:      end for
14:    else                                          $\triangleright$  Variable catégorielle
15:      subsets  $\leftarrow \text{GenerateSubsets}(X_j)$ 
16:      for  $S \in \text{subsets}$  do
17:        gain  $\leftarrow \text{ComputeGain}(t, j, S)$ 
18:        if gain > best_gain then
19:          best_gain  $\leftarrow \text{gain}$ 
20:          best_split  $\leftarrow (j, S)$ 
21:        end if
22:      end for
23:    end if
24:  end for
25:  return best_split
26: end procedure
  
```

4.3 Calcul du Gain (Classification)

Algorithm 3 Calcul du Gain d'Information (Gini)

```

1: procedure COMPUTEGAIN( $t, j, s$ )
2:    $(t_L, t_R) \leftarrow \text{PartitionNode}(t, j, s)$ 
3:    $n_t \leftarrow \text{—samples in } t\text{—}$ 
4:    $n_L \leftarrow \text{—samples in } t_L\text{—}$ 
5:    $n_R \leftarrow \text{—samples in } t_R\text{—}$ 
6:    $I_t \leftarrow \text{GiniImpurity}(t)$ 
7:    $I_L \leftarrow \text{GiniImpurity}(t_L)$ 
8:    $I_R \leftarrow \text{GiniImpurity}(t_R)$ 
9:    $\text{gain} \leftarrow I_t - \frac{n_L}{n_t}I_L - \frac{n_R}{n_t}I_R$ 
10:  return gain
11: end procedure
12: procedure GINIIMPURITY( $t$ )
13:   gini  $\leftarrow 1$ 
14:   for  $k = 1$  to  $K$  do
15:      $p_k \leftarrow \frac{1}{n_t} \sum_{i \in t} \mathbb{1}(y_i = k)$ 
16:     gini  $\leftarrow \text{gini} - p_k^2$ 
17:   end for
18:   return gini
19: end procedure

```

4.4 Calcul du Gain (Régression)

Algorithm 4 Calcul du Gain pour Régression

```

1: procedure COMPUTEREGRESSIONGAIN( $t, j, s$ )
2:    $(t_L, t_R) \leftarrow \text{PartitionNode}(t, j, s)$ 
3:    $n_t \leftarrow \text{—samples in } t\text{—}$ 
4:    $n_L \leftarrow \text{—samples in } t_L\text{—}$ 
5:    $n_R \leftarrow \text{—samples in } t_R\text{—}$ 
6:    $\text{Var}_t \leftarrow \text{Variance}(t)$ 
7:    $\text{Var}_L \leftarrow \text{Variance}(t_L)$ 
8:    $\text{Var}_R \leftarrow \text{Variance}(t_R)$ 
9:    $\text{gain} \leftarrow \text{Var}_t - \frac{n_L}{n_t} \text{Var}_L - \frac{n_R}{n_t} \text{Var}_R$ 
10:  return gain
11: end procedure
12: procedure VARIANCE( $t$ )
13:    $\bar{y} \leftarrow \frac{1}{n_t} \sum_{i \in t} y_i$ 
14:   var  $\leftarrow \frac{1}{n_t} \sum_{i \in t} (y_i - \bar{y})^2$ 
15:   return var
16: end procedure

```

5 Élagage (Pruning)

5.1 Principe de l'Élagage par Coût-Complexité

L'élagage vise à réduire le sur-apprentissage en simplifiant l'arbre. Le critère de coût-complexité est défini par :

$$R_\alpha(T) = R(T) + \alpha|\tilde{T}| \quad (15)$$

où :

- $R(T)$: erreur de l'arbre sur l'ensemble d'apprentissage
- $\alpha \geq 0$: paramètre de complexité (hyperparamètre de régularisation)
- $|\tilde{T}|$: nombre de feuilles (pénalisation de la complexité)

Intuition :

- $\alpha = 0$: arbre complet (sur-apprentissage possible)
- $\alpha \rightarrow \infty$: arbre réduit à la racine (sous-apprentissage)
- α optimal : compromis biais-variance optimal

5.2 Séquence d'Élagage

Pour chaque valeur de α , on trouve l'arbre optimal :

$$\hat{T}_\alpha = \arg \min_{T \subseteq T_{\max}} R_\alpha(T) \quad (16)$$

Le théorème de Breiman garantit l'existence d'une séquence finie d'arbres imbriqués :

$$T_{\max} \succ T_1 \succ T_2 \succ \dots \succ T_M = \{\text{root}\} \quad (17)$$

5.3 Algorithme d'Élagage

Algorithm 5 Élagage par Coût-Complexité

```

1: procedure PRUNETREE( $T$ ,  $\mathcal{D}_{\text{val}}$ ,  $y_{\text{val}}$ )
2:    $\alpha_{\text{sequence}} \leftarrow \text{ComputeAlphaSequence}(T)$ 
3:   trees  $\leftarrow []$ 
4:   for  $\alpha \in \alpha_{\text{sequence}}$  do
5:      $T_\alpha \leftarrow \text{PruneAtAlpha}(T, \alpha)$ 
6:     score  $\leftarrow \text{Evaluate}(T_\alpha, \mathcal{D}_{\text{val}}, y_{\text{val}})$ 
7:     trees.append(( $T_\alpha$ , score))
8:   end for
9:    $T^* \leftarrow \text{SelectBestTree}(\text{trees})$ 
10:  return  $T^*$ 
11: end procedure
12: procedure COMPUTEALPHASEQUENCE( $T$ )
13:    $\alpha_{\text{seq}} \leftarrow []$ 
14:    $T_{\text{current}} \leftarrow T$ 
15:   while  $T_{\text{current}}$  has more than 1 leaf do
16:      $\alpha_{\min} \leftarrow \infty$ 
17:     for  $t \in$  internal nodes of  $T_{\text{current}}$  do
18:        $R_t \leftarrow \text{Error}(t \text{ as leaf})$ 
19:        $R_{T_t} \leftarrow \text{Error}(\text{subtree at } t)$ 
20:        $|\tilde{T}_t| \leftarrow \text{Number of leaves in subtree at } t$ 
21:        $\alpha_t \leftarrow \frac{R_t - R_{T_t}}{|\tilde{T}_t| - 1}$ 
22:       if  $\alpha_t < \alpha_{\min}$  then
23:          $\alpha_{\min} \leftarrow \alpha_t$ 
24:          $t_{\text{prune}} \leftarrow t$ 
25:       end if
26:     end for
27:      $\alpha_{\text{seq}}.\text{append}(\alpha_{\min})$ 
28:     Prune subtree at  $t_{\text{prune}}$ 
29:   end while
30:   return  $\alpha_{\text{seq}}$ 
31: end procedure
  
```

6 Critères d'Arrêt

6.1 Paramètres de Contrôle

Les critères d'arrêt principaux incluent :

1. **Profondeur maximale** : $\text{depth}(t) \geq \text{max_depth}$
2. **Nombre minimal d'échantillons pour diviser** : $n_t < \text{min_samples_split}$
3. **Amélioration minimale** : $\Delta I < \text{min_impurity_decrease}$
4. **Nombre minimal d'échantillons par feuille** : $\min(n_L, n_R) < \text{min_samples_leaf}$
5. **Pureté du nœud** : $I(t) = 0$ (classification)

6. **Homogénéité des features** : Toutes les observations ont les mêmes valeurs de features

6.2 Algorithme de Vérification

Algorithm 6 Vérification des Critères d'Arrêt

```

1: procedure STOPPINGCRITERION( $t$ , params)
2:   if  $\text{depth}(t) \geq \text{params.max\_depth}$  then
3:     return true
4:   end if
5:   if  $n_t < \text{params.min\_samples\_split}$  then
6:     return true
7:   end if
8:   if  $I(t) = 0$  then                                 $\triangleright$  Nœud pur
9:     return true
10:  end if
11:  if AllSamplesSameFeatureValues( $t$ ) then
12:    return true
13:  end if
14:  return false
15: end procedure

```

7 Prédition

7.1 Algorithme de Prédition

Algorithm 7 Prédiction sur Nouvelles Données

```

1: procedure PREDICT( $T, \mathbf{x}$ )
2:    $t \leftarrow T.\text{root}$ 
3:   while  $t$  is not a leaf do
4:      $(j, s) \leftarrow \text{GetSplitInfo}(t)$ 
5:     if  $x_j$  est continue then
6:       if  $x_j \leq s$  then
7:          $t \leftarrow t.\text{left}$ 
8:       else
9:          $t \leftarrow t.\text{right}$ 
10:      end if
11:    else                                      $\triangleright$  Catégorielle
12:      if  $x_j \in s$  then                   $\triangleright s$  est un ensemble de catégories
13:         $t \leftarrow t.\text{left}$ 
14:      else
15:         $t \leftarrow t.\text{right}$ 
16:      end if
17:    end if
18:  end while
19:  return  $t.\text{prediction}$ 
20: end procedure
21: procedure BATCHPREDICT( $T, \mathbf{X}$ )
22:   predictions  $\leftarrow []$ 
23:   for  $\mathbf{x} \in \mathbf{X}$  do
24:      $\hat{y} \leftarrow \text{Predict}(T, \mathbf{x})$ 
25:     predictions.append( $\hat{y}$ )
26:   end for
27:   return predictions
28: end procedure
  
```

8 Analyse Détailée de l'Implémentation GitHub

8.1 Architecture du Code

Le repository [Adamkhald/breiman-cart](#) présente une implémentation modulaire de qualité production qui respecte fidèlement la méthodologie originale de 1984.

8.1.1 Structure des Fichiers

Fichier/Dossier	Description
-----------------	-------------

breiman_cart.py	API unifiée principale (v0.2.0, 2000+ lignes). Expose BRCRegression, BRCClassification et BRCInference pour usage simplifié.
breiman_cart/_init_.py	Initialisation du package. Gère les imports et expose l'API publique.
breiman_cart/node.py	Structure de données pour les noeuds avec support natif des catégories. Implémente la classe Node avec attributs pour divisions, prédictions et métadonnées.
breiman_cart/splitter.py	Logique de division incluant calcul de Gini, MSE et génération de sous-ensembles pour variables catégorielles.
breiman_cart/tree.py	Implémentation interne des arbres CART. Classes de base pour construction et traversée.
breiman_cart/pruning.py	Implémentation complète de l'élagage par coût-complexité avec calcul de la séquence alpha.
tests/	Suite de tests unitaires pour valider chaque composant et l'API.
docs/theory.tex	Documentation mathématique détaillée de la théorie CART.
demo.py	Script de démonstration avec exemples complets d'utilisation.
quickstart.py	Guide de démarrage rapide montrant l'utilisation de base.
test_new_api.py	Tests spécifiques pour la nouvelle API (v0.2.0).
fix_syntax.py	Utilitaire pour validation et correction syntaxique.
requirements.txt	Dépendances : numpy, pandas, pickle, json, logging, warnings, itertools.
pyproject.toml	Configuration du package Python (build, métadonnées).
MANIFEST.in	Fichiers à inclure dans la distribution.
LICENSE	Licence MIT.
ReadMe.md	Documentation principale du projet.

8.2 Nouvelle API Unifiée (Version 0.2.0)

La version 0.2.0 introduit une refonte majeure de l'API avec trois classes principales pour une expérience utilisateur améliorée :

8.2.1 BRCRegression - Régression

Classe dédiée aux tâches de régression avec interface simplifiée :

Listing 1 – Utilisation de BRCRegression

```

1 from breiman_cart import BRCRegression
2 import pandas as pd
3 import numpy as np
4
5 # Données de régression
6 X = pd.DataFrame({

```

```

7     'sqft': [1200, 1800, 1500, 2200, 1900],
8     'neighborhood': ['A', 'B', 'A', 'C', 'B'],
9     'age_years': [10, 5, 15, 2, 7]
10)
11y = np.array([300000, 450000, 350000, 550000, 420000])
12
13# Cr ation et entra nement du mod le
14model = BRCRegression(
15    max_depth=10,
16    min_samples_split=20,
17    min_samples_leaf=10
18)
19model.fit(X, y, categorical_features=['neighborhood'])
20
21# Pr dictions
22predictions = model.predict(X)
23
24# Sauvegarde du mod le
25model.save('regression_model.pkl')

```

8.2.2 BRCClassification - Classification

Classe optimisée pour les tâches de classification :

Listing 2 – Utilisation de BRCClassification

```

1 from breiman_cart import BRCClassification
2
3 # Donn es de classification
4 X = pd.DataFrame({
5     'age': [25, 45, 35, 50, 23, 40, 55, 30],
6     'city': ['NYC', 'LA', 'NYC', 'SF', 'LA', 'SF', 'NYC', 'LA'],
7     'income': [50000, 80000, 60000, 95000, 45000, 75000, 100000,
8                 55000]
9 })
10y = np.array([0, 1, 0, 1, 0, 1, 1, 0])
11
12# Cr ation du classificateur
13model = BRCClassification(
14    max_depth=5,
15    min_samples_split=2,
16    criterion='gini' # ou 'entropy'
17)
18model.fit(X, y, categorical_features=['city'])
19
20# Pr dictions
21predictions = model.predict(X)
22probabilities = model.predict_proba(X)
23
24# valuation
25accuracy = model.score(X, y)
26print(f"Accuracy:{accuracy:.3f}")

```

8.2.3 BRCInference - Inférence avec Modèles Sauvegardés

Classe pour charger et utiliser des modèles pré-entraînés :

Listing 3 – Utilisation de BRCInference

```

1 from breiman_cart import BRCInference
2
3 # Chargement d'un modèle sauvegardé
4 model = BRCInference.load('regression_model.pkl')
5
6 # Prédictions sur nouvelles données
7 new_data = pd.DataFrame({
8     'sqft': [1600, 2000],
9     'neighborhood': ['A', 'B'],
10    'age_years': [8, 3]
11})
12
13 predictions = model.predict(new_data)
14 print(f"Prix prédits : {predictions}")
15
16 # Export du modèle en différents formats
17 model.export_json('model.json')
18 model.export_summary('model_summary.txt')

```

8.3 Caractéristiques Principales

8.3.1 Gestion Native des Variables Catégorielles

Innovation majeure : L'implémentation traite véritablement les variables catégorielles par division de sous-ensembles (subset splitting), contrairement à scikit-learn qui utilise l'encodage one-hot.

Avantages :

- Fidélité à l'algorithme original de Breiman
- Réduction de la dimensionnalité (pas d'explosion du nombre de variables)
- Meilleure interprétabilité des divisions
- Respect de la sémantique des variables catégorielles

Listing 4 – Exemple d'utilisation avec la nouvelle API v0.2.0

```

1 from breiman_cart import BRCClassification
2 import pandas as pd
3 import numpy as np
4
5 # Données avec variable catégorielle 'city'
6 X = pd.DataFrame({
7     'age': [25, 45, 35, 50, 23],
8     'city': ['NYC', 'LA', 'NYC', 'SF', 'LA']
9 })
10 y = np.array([0, 1, 0, 1, 0])
11
12 # Construction du classificateur avec nouvelle API
13 model = BRCClassification()

```

```

14     max_depth=5,
15     min_samples_split=2,
16     criterion='gini'
17 )
18
19 # Entrancement (spécifier catgorielles dans fit())
20 model.fit(X, y, categorical_features=['city'])
21
22 # Prédictions
23 predictions = model.predict(X)
24 print("Prédictions:", predictions)
25
26 # Probabilités
27 probabilities = model.predict_proba(X)
28 print("Probabilités:", probabilities)
29
30 # Sauvegarde du modèle
31 model.save('classifier.pkl')

```

Avantages de la nouvelle API :

- Interface plus intuitive et cohérente
- Noms de classes explicites (BRCClassification vs CARTClassifier)
- Gestion simplifiée des modèles sauvegardés avec BRCInference
- Support natif de sauvegarde/chargement (pickle et JSON)
- Meilleure séparation des responsabilités
- Documentation inline améliorée

8.3.2 Élagage par Coût-Complexité

L'implémentation complète de l'élagage suit rigoureusement la méthodologie originale de Breiman avec :

- Calcul de la séquence complète des valeurs α
- Génération de tous les sous-arbres optimaux
- Validation sur ensemble de validation séparé
- Sélection du meilleur arbre selon la performance de validation

Listing 5 – Élagage de l'arbre

```

1 # Division des données
2 X_train, X_test, y_train, y_test = train_test_split(
3     X, y, test_size=0.2, random_state=42
4 )
5 X_train, X_val, y_train, y_val = train_test_split(
6     X_train, y_train, test_size=0.2, random_state=42
7 )
8
9 # Entrancement sur ensemble d'apprentissage
10 tree.fit(X_train, y_train)
11
12 # Élagage avec ensemble de validation
13 tree.prune(X_val, y_val)
14

```

```

15 # valuation sur ensemble de test
16 accuracy = tree.score(X_test, y_test)

```

8.3.3 Type Safety

Le code utilise des annotations de type Python (type hints) pour :

- Améliorer la lisibilité
- Faciliter la maintenance
- Permettre la détection d'erreurs statique
- Documenter implicitement les interfaces

8.4 Exemple de Régression

Listing 6 – Utilisation pour régression avec nouvelle API

```

1 from breiman_cart import BRCRegression
2
3 # Données pour prédiction de prix immobilier
4 X = pd.DataFrame({
5     'sqft': [1200, 1800, 1500, 2200, 1900],
6     'neighborhood': ['A', 'B', 'A', 'C', 'B'],
7     'age_years': [10, 5, 15, 2, 7]
8 })
9 y = np.array([300000, 450000, 350000, 550000, 420000])
10
11 # Construction du régisseur
12 model = BRCRegression(
13     max_depth=4,
14     min_samples_split=2,
15     min_impurity_decrease=0.01
16 )
17 model.fit(X, y, categorical_features=['neighborhood'])
18
19 # Prédictions
20 predictions = model.predict(X)
21 print("Prix prédictifs:", predictions)
22
23 # Métriques
24 mse = np.mean((y - predictions)**2)
25 rmse = np.sqrt(mse)
26 r2 = model.score(X, y)
27 print(f"RMSE: {rmse:.2f}")
28 print(f"R² : {r2:.3f}")
29
30 # Export et sauvegarde
31 model.save('house_price_model.pkl')
32 model.export_json('model_config.json')

```

8.5 Gestion des Modèles et Persistance

La version 0.2.0 introduit des fonctionnalités robustes de sauvegarde et chargement des modèles :

8.5.1 Sauvegarde et Chargement avec Pickle

Listing 7 – Persistance des modèles

```

1 from breiman_cart import BRCClassification, BRCInference
2
3 # Entrainement et sauvegarde
4 model = BRCClassification(max_depth=5)
5 model.fit(X_train, y_train, categorical_features=['city'])
6 model.save('my_model.pkl')
7
8 # Chargement et utilisation
9 loaded_model = BRCInference.load('my_model.pkl')
10 predictions = loaded_model.predict(X_test)

```

8.5.2 Export en Format JSON

Listing 8 – Export JSON pour interopérabilité

```

1 # Export de la configuration complète
2 model.export_json('model_config.json')
3
4 # Le fichier JSON contient:
5 # - Structure de l'arbre
6 # - Hyperparamètres
7 # - Statistiques d'entraînement
8 # - Variables catégorielles
9 # - M étadonnées du modèle

```

8.5.3 Génération de Résumés

Listing 9 – Génération de résumé lisible

```

1 # Génération d'un résumé textuel
2 model.export_summary('model_summary.txt')
3
4 # Contient:
5 # - Nombre de nœuds et feuilles
6 # - Profondeur de l'arbre
7 # - Variables importantes
8 # - Performance du modèle
9 # - Règles de décision principales

```

8.6 Logging et Débogage

Le package intègre un système de logging complet pour faciliter le débogage et la surveillance :

Listing 10 – Configuration du logging

```

1 import logging
2
3 # Configuration du niveau de log
4 logging.basicConfig(
5     level=logging.INFO,
6     format='%(asctime)s - %(name)s - %(levelname)s - %(message)s'
7 )
8
9 # Le mode log automatiquement:
10 # - tapes de construction de l'arbre
11 # - Divisions choisies chaque nœud
12 # - Métriques d'impureté
13 # - Warnings pour données problématiques
14 # - Erreurs et exceptions
15
16 model = BRCClassification(max_depth=5, verbose=True)
17 model.fit(X, y) # Affiche des logs de tailles

```

8.7 Différences avec Scikit-Learn

Aspect	breiman-cart (v0.2.0)	scikit-learn
API Principale	BRCClassification, BRCRegression, BRCInference	DecisionTreeClassifier, DecisionTreeRegressor
Variables catégorielles	Division par sous-ensembles native	Encodage one-hot requis
Élagage	Séquence complète de coût-complexité	Élagage simplifié via ccp_alpha
Sauvegarde	Méthodes save()/load() intégrées + JSON	Nécessite joblib/pickle externe
Logging	Système de logging intégré	Logging minimal
Transparence	Code Python pur, éducatif	Code optimisé en Cython
Dépendances	NumPy/Pandas + stdlib	Scipy, joblib, etc.
Performance	Plus lent mais pédagogique	Très rapide
Objectif	Éducation et recherche	Production
Fidélité Breiman	Implémentation fidèle 1984	Optimisations modernes
Export	JSON, pickle, résumés	pickle/joblib

8.8 Installation et Utilisation

8.8.1 Installation

Listing 11 – Installation du package

```

1 # Installation des dépendances
2 pip install numpy pandas
3
4 # Clonage du repository
5 git clone https://github.com/Adamkhald/breiman-cart.git
6 cd breiman-cart
7
8 # Installation en mode développement (optionnel)
9 pip install -e .

```

8.8.2 Tests

Listing 12 – Exécution des tests

```

1 # Installation de pytest
2 pip install pytest
3
4 # Execution de tous les tests
5 python -m pytest tests/
6
7 # Tests avec couverture
8 pip install pytest-cov
9 python -m pytest tests/ --cov=src

```

8.9 Complexité Temporelle

8.9.1 Construction de l'Arbre

- Meilleur cas : $O(np \log n)$ (arbre équilibré)
- Pire cas : $O(n^2 p)$ (arbre dégénéré linéaire)
- Cas moyen : $O(np \log n)$

Détail par niveau :

- Nombre de niveaux : $O(\log n)$ (équilibré) à $O(n)$ (dégénéré)
- Par niveau : $O(np)$ pour évaluer toutes les divisions possibles
- Tri initial des variables continues : $O(np \log n)$

Variables catégorielles :

- Cas général : $O(2^m)$ divisions à évaluer
- Classification binaire optimisée : $O(m \log m)$

8.9.2 Élagage

- Calcul de la séquence α : $O(|\text{nodes}|^2)$
- Validation de chaque arbre dans la séquence : $O(|\alpha_{\text{seq}}| \cdot n_{\text{val}} \cdot \text{depth})$
- Total : $O(|\text{nodes}|^2 + |\text{nodes}| \cdot n_{\text{val}} \cdot \log n)$

8.9.3 Prédiction

- Une prédiction : $O(\text{depth})$ où $\text{depth} \leq \log n$ (équilibré)
- Batch de m prédictions : $O(m \cdot \text{depth})$

8.10 Complexité Spatiale

- Stockage de l'arbre : $O(|\text{nodes}|) = O(2^{\text{depth}} - 1)$
- Données triées (si pré-calcul) : $O(np)$
- Indices pour chaque noeud : $O(n \cdot \text{depth})$
- Total : $O(np + n \cdot \text{depth} + 2^{\text{depth}})$

8.11 Optimisations Possibles

1. **Tri pré-calculé** : Trier chaque variable une seule fois au début
2. **Statistiques incrémentales** : Mise à jour efficace des moyennes et variances
3. **Parallélisation** : Évaluation parallèle des variables
4. **Early stopping** : Arrêt anticipé si gain insuffisant

9 Applications et Cas d'Usage

9.1 Domaines d'Application

9.1.1 Médecine et Santé

Exemple : Prédiction du risque cardiovasculaire

- **Variables** : âge, pression artérielle, cholestérol, IMC, antécédents familiaux, habitudes (tabac, exercice)
- **Avantage CART** : Règles de décision interprétables pour les cliniciens
- **Exemple de règle** : "Si âge ≥ 55 ET pression ≥ 140 ET cholestérol ≥ 240 ALORS risque élevé"

Autres applications médicales :

- Diagnostic de maladies
- Stratification des patients
- Prédiction de réponse aux traitements
- Identification de facteurs de risque

9.1.2 Finance

Exemple : Évaluation du risque de crédit

- **Variables** : revenu, historique de crédit, ratio d'endettement, durée emploi, historique bancaire
- **Avantage CART** : Règles transparentes pour conformité réglementaire (explicabilité requise)
- **Résultat** : Score de crédit binaire ou probabilité de défaut

Autres applications financières :

- Détection de fraude
- Prédiction de churn client
- Évaluation d'investissements
- Pricing d'assurances

9.1.3 Marketing et E-commerce

Exemple : Segmentation de clientèle

- **Variables** : démographiques (âge, revenu), comportementales (fréquence achat, panier moyen), géographiques
- **Avantage CART** : Segments actionnables et interprétables
- **Utilisation** : Campagnes marketing ciblées, recommandations personnalisées

10 Avantages et Limitations

10.1 Avantages de CART

1. **Interprétabilité exceptionnelle**
 - Structure visuelle claire (arbre)
 - Règles de décision compréhensibles par des non-experts
 - Facilite l'adhésion des parties prenantes
 - Important pour domaines réglementés (finance, santé)
2. **Pas d'hypothèses distributionnelles**
 - Méthode non-paramétrique
 - Pas d'hypothèse sur la distribution des données
 - Robuste aux violations d'hypothèses
3. **Gestion native des types de données**
 - Variables continues sans transformation
 - Variables catégorielles sans encodage
 - Variables ordinaires respectées
 - Mix de types dans le même modèle
4. **Détection automatique des interactions**
 - Captures automatiques des interactions complexes
 - Pas besoin de spécification manuelle
 - Hiérarchie naturelle des interactions
5. **Robustesse aux valeurs aberrantes**
 - Divisions basées sur les rangs/ordres
 - Peu sensible aux valeurs extrêmes
 - Pas de normalisation nécessaire
6. **Sélection automatique de variables**
 - Seules les variables informatives utilisées
 - Gestion naturelle de la haute dimensionnalité
 - Feature importance intrinsèque
7. **Gestion des valeurs manquantes**
 - Méthodes de surrogate splits
 - Pas d'imputation nécessaire
 - Utilisation optimale des données disponibles

10.2 Limitations de CART

1. **Instabilité élevée**
 - Petites variations ⇒ arbres très différents

- Problème de généralisation
 - Solution : ensembles d'arbres (Random Forests, Boosting)
2. **Sur-apprentissage facile**
 - Arbres profonds mémorisent les données
 - Forte variance du modèle
 - Solution : élagage rigoureux, validation croisée
 3. **Biais de sélection de variables**
 - Favorise les variables avec plus de points de division
 - Variables continues avantagées vs catégorielles
 - Impact sur l'importance des variables
 4. **Frontières de décision rectangulaires**
 - Limitation pour relations obliques
 - Approximation grossière de frontières lisses
 - Nécessite beaucoup de divisions pour frontières diagonales
 5. **Difficulté d'extrapolation**
 - Prédictions constantes dans chaque région
 - Pas d'extrapolation au-delà des données d'entraînement
 - Problématique pour séries temporelles
 6. **Sensibilité au déséquilibre des classes**
 - Performance dégradée sur classes minoritaires
 - Solution : pondération des classes, stratification
- ```
((y_ttest == 0)(y_predit_ttest == 1))[0]print(f"de faux positifs : len(fp_indices)") if len(fp_indices) > 0 : print("Exemples de faux positifs :") print(X_ttest.iloc[fp_indices[:5]])
Faux négatifs fn_indices = np.where((y_ttest == 1)(y_predit_ttest == 0))[0]print(f"de faux négatifs : len(fn_indices)") if len(fn_indices) > 0 : print("Exemples de faux négatifs :") print(X_ttest.iloc[fn_indices[:5]])
```

## 11 Conclusion

### 11.1 Contributions de CART

L'algorithme CART de Breiman et al. (1984) a apporté des contributions majeures et durables à l'apprentissage automatique :

1. **Méthodologie unifiée** : Cadre cohérent pour classification et régression
2. **Fondements théoriques solides** : Base mathématique rigoureuse avec preuves de consistance
3. **Praticité exceptionnelle** : Implémentation efficace et interprétable
4. **Influence durable** : Fondation pour Random Forests, Gradient Boosting, XG-Boost
5. **Adoption industrielle** : Utilisation généralisée dans de nombreux domaines

### 11.2 Apports du Repository GitHub

L'implémentation Adamkhald/breiman-cart se distingue par :

- **Fidélité historique** : Respect scrupuleux de la méthodologie de 1984

- **Gestion native des catégories** : Implémentation correcte du subset splitting
- **Élagage complet** : Séquence alpha complète selon Breiman
- **Code pédagogique** : Clarté et lisibilité pour l'apprentissage
- **Minimalisme** : Dépendances réduites (NumPy/Pandas uniquement)

### 11.3 Perspectives Futures

Les directions de recherche et développement actuelles incluent :

- **Deep Learning et arbres** :
  - Neural Decision Trees
  - Soft Decision Trees différentiables
  - Tree-structured neural networks
- **Arbres causaux** :
  - Inférence causale via structures arborescentes
  - Causal CART pour hétérogénéité d'effets
  - Policy learning trees
- **Équité et explicabilité** :
  - Garanties formelles d'équité
  - Fair CART avec contraintes
  - Explications contrefactuelles
- **Scalabilité** :
  - Algorithmes distribués (MapReduce, Spark)
  - Approximations pour très grandes données
  - Online learning pour streaming

### 11.4 Recommandations Pratiques

Pour une utilisation efficace de CART en production :

1. **Toujours élaguer** : Utiliser un ensemble de validation pour l'élagage
2. **Validation rigoureuse** : Validation croisée systématique
3. **Ensembles pour production** : Random Forests ou Boosting pour robustesse
4. **Monitoring de stabilité** : Vérifier sur échantillons bootstrap
5. **Interprétation prudente** : Distinguer corrélation et causalité
6. **Feature engineering** : Malgré l'automatisme, les transformations aident
7. **Déséquilibre des classes** : Utiliser class\_weight ou stratification
8. **Valeurs manquantes** : Implementer surrogate splits si nécessaire

### 11.5 Impact et Héritage

CART reste aujourd'hui :

- Un des algorithmes les plus utilisés en pratique
- La base de nombreux algorithmes modernes performants
- Un outil d'exploration de données incontournable
- Un standard d'interprétabilité en ML
- Une méthode d'enseignement fondamentale

## Remerciements

Ce rapport a été élaboré en analysant l'implémentation open-source disponible sur le repository GitHub [Adamkhald/breiman-cart](#) et les fondements théoriques de l'article original de Breiman et al. (1984). Nous remercions les contributeurs du repository pour leur travail rigoureux et pédagogique.

## Références

- [1] Breiman, L., Friedman, J., Olshen, R., & Stone, C. (1984). *Classification and Regression Trees*. Wadsworth International Group, Belmont, CA.
- [2] Breiman, L. (1996). Bagging predictors. *Machine Learning*, 24(2), 123–140.
- [3] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
- [4] Breiman, L. (2001). Statistical modeling : The two cultures. *Statistical Science*, 16(3), 199–215.
- [5] Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning : Data Mining, Inference, and Prediction* (2nd ed.). Springer.
- [6] James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical Learning with Applications in R*. Springer.
- [7] Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, 1(1), 81–106.
- [8] Quinlan, J. R. (1993). *C4.5 : Programs for Machine Learning*. Morgan Kaufmann.
- [9] Therneau, T. M., & Atkinson, E. J. (2000). *An Introduction to Recursive Partitioning Using the RPART Routines*. Technical Report 61, Mayo Clinic.
- [10] Loh, W. Y. (2011). Classification and regression trees. *Wiley Interdisciplinary Reviews : Data Mining and Knowledge Discovery*, 1(1), 14–23.
- [11] Strobl, C., Boulesteix, A. L., Zeileis, A., & Hothorn, T. (2007). Bias in random forest variable importance measures : Illustrations, sources and a solution. *BMC Bioinformatics*, 8(1), 25.
- [12] Athey, S., & Imbens, G. (2016). Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27), 7353–7360.
- [13] Repository GitHub : <https://github.com/Adamkhald/breiman-cart>

## A Glossaire des Notations

| Notation         | Description                                 |
|------------------|---------------------------------------------|
| $\mathcal{D}$    | Ensemble de données d'apprentissage         |
| $n$              | Nombre total d'observations                 |
| $p$              | Nombre de variables explicatives (features) |
| $\mathbf{x}_i$   | Vecteur de variables pour l'observation $i$ |
| $y_i$            | Variable réponse pour l'observation $i$     |
| $T$              | Arbre de décision complet                   |
| $\mathcal{N}(T)$ | Ensemble des nœuds internes de l'arbre      |
| $\tilde{T}$      | Ensemble des nœuds terminaux (feuilles)     |

|               |                                                  |
|---------------|--------------------------------------------------|
| $ \tilde{T} $ | Nombre de feuilles (complexité de l'arbre)       |
| $R_j$         | Région correspondant à la feuille $j$            |
| $n_t$         | Nombre d'observations dans le nœud $t$           |
| $I(t)$        | Impureté du nœud $t$                             |
| $\alpha$      | Paramètre de coût-complexité pour l'élagage      |
| $R_\alpha(T)$ | Critère de coût-complexité                       |
| $R_L, R_R$    | Régions gauche et droite après division          |
| $n_L, n_R$    | Nombre d'observations dans régions gauche/droite |
| $p_k(t)$      | Proportion de classe $k$ dans le nœud $t$        |
| $\hat{y}_j$   | Prédiction pour la région $j$                    |
| $K$           | Nombre de classes (classification)               |
| $m$           | Nombre de catégories (variable catégorielle)     |

## B Formules Récapitulatives

### B.1 Critères d'Impureté

Gini (Classification) :

$$I_{\text{Gini}}(t) = 1 - \sum_{k=1}^K p_k^2(t) \quad (18)$$

Entropie (Classification) :

$$I_{\text{Entropy}}(t) = - \sum_{k=1}^K p_k(t) \log_2 p_k(t) \quad (19)$$

Variance (Régression) :

$$I_{\text{Var}}(t) = \frac{1}{n_t} \sum_{i \in t} (y_i - \bar{y}_t)^2 \quad (20)$$

### B.2 Gain de Division

$$\Delta I(s, t) = I(t) - \frac{n_L}{n_t} I(t_L) - \frac{n_R}{n_t} I(t_R) \quad (21)$$

### B.3 Coût-Complexité

$$R_\alpha(T) = R(T) + \alpha |\tilde{T}| \quad (22)$$

### B.4 Alpha Critique

$$\alpha_t = \frac{R_t - R_{T_t}}{|\tilde{T}_t| - 1} \quad (23)$$

## B.5 Décomposition Biais-Variance

$$\mathbb{E}[(y - \hat{f}(\mathbf{x}))^2] = \text{Var}(\hat{f}(\mathbf{x})) + [\text{Bias}(\hat{f}(\mathbf{x}))]^2 + \sigma^2 \quad (24)$$