Spatial CART Classification Trees

Jean-Michel Poggi, Avner Bar-Hen, Servane Gey

¹ Univ. Paris & LMO, Orsay, Univ. Paris-Saclay, France
 ² CNAM, Paris, France
 ³ MAP5, UMR 8145, Univ. Paris, France

Montevideo, Feb. 2025

Introduction

- Classification And Regression Trees, Breiman et al. (1984)
- Variants and extensions of the original CART to the spatial domain
 - Ortho-CART Donoho et al. (1997), in image processing, dyadic splits + pruning using the algorithm used for the wavelet packets best basis
 - Dyadic-CART, ideas generalized in Blanchard et al. (2007)
 - Extension to spatial data with kriging type ideas see Bel et al. (2009)

• Our variant: Spatial CART

- For spatial data, extend CART for bivariate marked point processes
- New splitting criterion in Spatial CART, taking into account the spatial information, to propose a segmentation of the window into homogeneous areas for interaction between marks

Outline

Classical CART classification trees

- Binary classification
- CART Algorithm

2 Spatial CART Classification Trees

- Motivation
- Spatial CART Algorithm

CART and Spatial CART in action: Rain-forest in Paracou
 Initial resolution
 Results

Classification Trees

• Predict the unknown binary label $Y \in \{0, 1\}$ of an observation $X \in \mathbb{R}^{p}$ via a classifier

$$f: \mathbb{R}^p \to \{0; 1\}$$

• Bayes classifier: minimizer of $f \mapsto Pf := P(Y \neq f(X))$ (with $(X, Y) \sim P$)

 $f^* = \mathbb{1}_{n(x) \ge 1/2}$, with $\eta(x) = P(Y = 1 | X = x)$



A. Bar Hen | S. Gey | J-M. Poggi

CART Algorithm

Classification And Regression Trees, Breiman et al. (1984)

Growing step

- recursive partitioning by maximizing a local decreasing of heterogeneity often based on Gini index or Shannon entropy
- do not split a pure node or a node containing few data
- \Rightarrow maximal tree T_{max}
- T_{max} overfits the data

Pruning step

- Optimal tree: subtree pruned from *T_{max}*
- Reduce the number of tree candidates: minimize

$$\operatorname{crit}_{\alpha}(T) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\hat{f}_{T}(X_{i}) \neq Y_{i}} + \alpha \frac{|\widetilde{T}|}{n},$$

 $|\widetilde{T}|$ = number of leaves of T

•
$$\Rightarrow$$
 sequence $(T_k)_{1 \leq k \leq K}$

CART Algorithm (2)

Classification And Regression Trees, Breiman et al. (1984)

Theorem (Breiman et al. 84)

For all $\alpha \ge 0$, $\operatorname{argmin}_{T \preceq T_{max}} \operatorname{crit}_{\alpha}(T)$ belongs to the sequence of nested pruned subtrees $(T_k)_{1 \le k \le K}$.

Selection step (Cross-validation or Hold Out)

- Data split into a training set \mathcal{L} of size n, and a test set \mathcal{T} of size n_t
- Build $(T_k)_{1 \leq k \leq K}$ on \mathcal{L} and select

$$\hat{k} = \operatorname*{argmin}_{1 \leq k \leq K} \frac{1}{n_t} \sum_{(X_i, Y_i) \in \mathcal{T}} \mathbb{1}_{\hat{t}_{T_k}(X_i) \neq Y_i}$$

• \Rightarrow Final CART tree is given by \hat{f}_{T_k}

Outline

Classical CART classification trees

- Binary classification
- CART Algorithm

2 Spatial CART Classification Trees

- Motivation
- Spatial CART Algorithm

CART and Spatial CART in action: Rain-forest in Paracou Initial resolution Results

Spatial CART Algorithm

General idea:

- To build a tessellation of the window into homogeneous areas for interaction between marks
- To use the spatial information to build a classification tree on the observed points of the bivariate point process

Spatial CART as a variant of CART:

- *Variant in growing*: splitting criterion taking into account the spatial characterization of the data, based on the intertype function *K*_{ij}
- *Variant in pruning*: penalized criterion based on least squares criterion to estimate local mark intensities
- Variant in final selection: optimal tree selected by a variant of the slope heuristic (Massart et al.) to keep the spatial information

Bivariate spatial point process



Observation = realization of (X, M)

- Left part: blue points repulse red ones (r = 0.05)
- Right part: blue and red points are independently distributed
- Same marginal distribution (color) on left and right side

- Bivariate spatial point process: $(X, M) \in W \times \{i; j\} \sim P,$ $W \subset \mathbb{R}^2$
- Mark intensity: for ★ = i, j,
 λ_⋆ intensity of the spatial point process (X | M = ⋆)
- Intertype function: at scale $r \ge 0$,

 $K_{ij}(r) = \lambda_j^{-1} \mathbb{E} \left(N_{ij}(r) \right),$

where $N_{ij}(r)$ counts the number of type *j* points at distance at most *r* of a randomly chosen type *i* point

9/24

Interaction Examples



In black: *Estimate* of K_{ij} on observed points and In red: *theoretical* K_{ij} for independent homogeneous Poisson p.p. (with interval bounds)

• But we will use the intertype function from a different perspective

Splitting criterion

At each node t define

- node area:
- estimates of mark intensities:
- estimate of $K_{ij}(r)$:

$$\hat{(\lambda_i^t, \lambda_j^t)}$$
$$\hat{K}_{ij}^t(r) = (\hat{\lambda_i^t} \hat{\lambda_j^t} \mathbf{A}^t)^{-1} \sum_{\{i_k, j_l \in t\}} \mathbb{1}_{d_{i_k, j_l} < r}$$

 d_{i_k,j_l} Distance between individuals i_k of mark *i* and j_l of mark *j* Impurity function: for a node *t*, a splitting *s* of *t* into t_l and t_R , and r > 0

At

$$\Delta I_{ij}(s,t,r) := \hat{K}_{ij}^{t}(r) - \alpha_s \frac{\hat{\lambda}_i^{t_L} \hat{\lambda}_j^{t_L}}{\hat{\lambda}_i^t \hat{\lambda}_j^t} \hat{K}_{ij}^{t_L}(r) - (1 - \alpha_s) \frac{\hat{\lambda}_i^{t_R} \hat{\lambda}_j^{t_R}}{\hat{\lambda}_i^t \hat{\lambda}_j^t} \hat{K}_{ij}^{t_R}(r) \ge \mathbf{0}$$

with $\alpha_s = A^{t_L}/A^t$ the area proportion of $t_L \subset t$

Splitting rule of node *t* at fixed scale r > 0

$$\hat{s}(t,r) = \operatorname*{argmax}_{s} \Delta I_{ij}(s,t,r)$$

Growing maximal tree T_{max}

Input	Bivariate spatial point process,
	scale r_0 ,
Initialize	node $t = t_1$ the root of the tree, $r_t = r_0$ the scale value at node t , $\operatorname{argmax}\{\hat{\lambda}_i^t, \hat{\lambda}_j^t\}$ the label of node t .
Recursion	at node <i>t</i> Compute
	$i_0 = \operatorname*{argmax}_{\star} \hat{\lambda}^t_{\star}, j_0 = \operatorname*{argmin}_{\star} \hat{\lambda}^t_{\star},$
	$\hat{\boldsymbol{s}} = \operatorname*{argmax}_{a} \Delta \boldsymbol{I}_{i_0 j_0}(\boldsymbol{s}, \boldsymbol{t}, \boldsymbol{r}_t),$
	Set
	$t_L = \{ \text{points in } t \mid \text{answer "yes" to } \hat{s} \},$
	$t_R = \{ \text{points in } t \mid \text{answer "no" to } \hat{s} \}.$
	$r_t = \operatorname{argmax}_{r} \Delta I_{i_0 j_0}(\hat{s}, t, r),$
	left: $t = t_L$,
	right: $t = t_R$.
Output	Maximal tree T _{max} .

A. Bar Hen | S. Gey | J-M. Poggi

CART (left) and Spatial CART (right) maximal trees



Spatial CART: initial scale $r_0 < r_{repuls} = 0.05$

Penalized criterion Class Probability Trees *Breiman et al. 84*

- If X is locally stationary, estimating local mark intensities amounts to estimating local mark rates
- Use penalized criterion derived from Gini index to prune T_{max}

$$\operatorname{crit}_{\alpha}^{G}(T) = \frac{1}{n} \sum_{t \in \widetilde{T}} n_{t} \left(1 - \sum_{\star = i, j} \hat{p}(\star | t)^{2} \right) + \alpha \frac{|\widetilde{T}|}{n},$$

where *n* = number of observed points; n_t = number of points falling in node *t*; $\hat{p}(\star|t)$ proportion of points of type \star in node *t*

• \Rightarrow sequence of nested pruned subtrees $(T_k)_{1 \leq k \leq K}$

Final tree selection





Complexity

- Identifying the "largest complexity plateau" (red circle) or the modified "largest dimension jump" (blue triangle), more agressive over-penalizing
- Related to the slope heuristic proposed by Birge, Massart in the 2000s (see Baudry et al. (2012) for a recent survey)

A. Bar Hen | S. Gey | J-M. Poggi

Spatial CART

15/24

Outline

1) Classical CART classification trees

- Binary classification
- CART Algorithm

2 Spatial CART Classification Trees

- Motivation
- Spatial CART Algorithm

CART and Spatial CART in action: Rain-forest in Paracou Initial resolution

Results

Data description (Gourlet-Fleury et al. 2004, Traissac 2003) DISPOSITIF EXPERIMENTAL DE PARACOU



Réalisation CIRAD-Forêt, Janvier 1998

A. Bar Hen | S. Gey | J-M. Poggi

Spatial CART

Rain-forest in Paracou: focus on two species



- Two tree species: Vouacapoua americana and Oxandra asbeckii
- Elevation is the environmental factor that drives their spatial distribution and this creates a strong interaction between both repartitions
- Competition is *high* for the hill at the bottom of the plot and very *low* at the top left of the plot

A. Bar Hen | S. Gey | J-M. Poggi

Spatial CART

Choice of initial scale r_0



Figure: Difference between estimated and theoretical K_{ij} ; blue: r = 6, red: r = 24.

- r = 6: species begin to interact
- r = 24: the interaction between species increases rapidly
- Initial median scale value r = 15 for SpatCART is sufficiently large to capture interaction, and not too large to avoid deeper maximal trees

Spatial CART partition



• SpatCART (with r = 15) recovers the spatial structure and the interaction-based (on the $K_{ij}(r_t)$ for all the nodes t) colormap is meaningful

CART partition (the largest plateau variant)



- CART results are not informative from the spatial viewpoint: it highlights the regions according to the specie distribution, not w.r.t. the interaction
- CART cannot catch the mixed structure of species

Perspectives

- Extension to spatial Bagging or spatial Random Forests to cope with instability issue
- Use the sensitivity of CART with respect to rotation to generate several tessellations
- Extension to multi-marked point processes by combining one-versus-rest classifiers and then obtain several tessellations and select the partition maximizing some global measure of heterogeneity between cells
- Incorporate covariables:
 - in the example, elevation could be introduced as a third spatial coordinate,
 - more generally, we could imagine to first perform a classical CART using additional covariables and then, in each leaf, to perform a SpatCART and finally select the best one

Paper and Reproducible research

- A. Bar-Hen, S. Gey and J-M. Poggi, "Spatial CART Classification Trees" *Computational Statistics*, 36, 2591-2613, 2021
- An R package spatcart, and the R codes to reproduce experiments of the paper, are available on https://github.com/ Servane-Gey/Spatial-classification-trees.
- Package spatcart may also be directly installed with R package devtools from the github repository Servane-Gey/spatcart. Package spatcart requires the following R packages to implement the results:
 - spatstat to deal with point processes, and in particular to compute ΔI_{ii} in the construction of the maximal tree,
 - tree to deal with tree structures.

Thank you!