

Random Forest for Regression of a Censored Variable

Yohann Le Faou

July 4, 2017

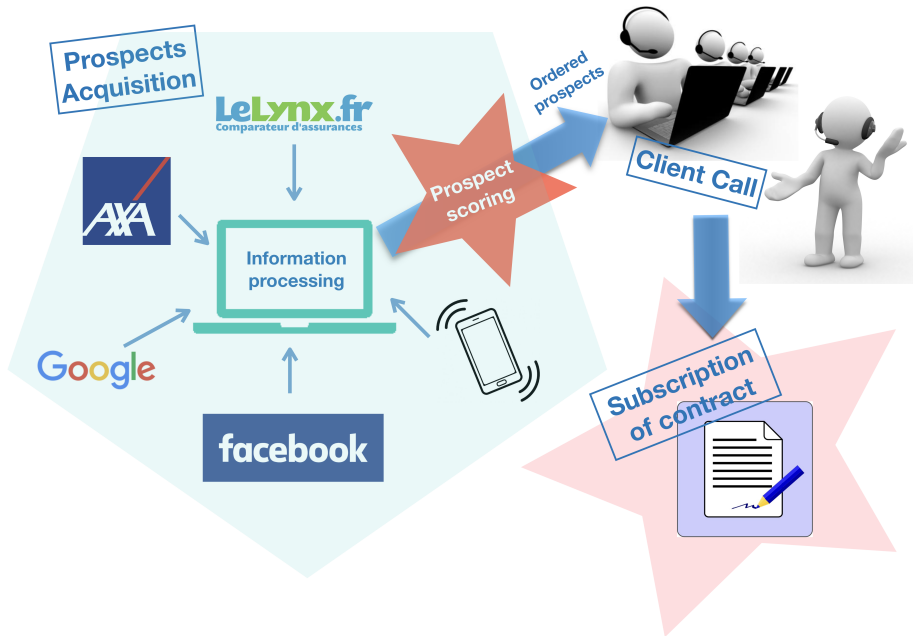
Insurance, Mathematics and Economics 2017

Table of contents

1. Introduction
2. Weighted Random Forest and IPCW method
3. Experiments

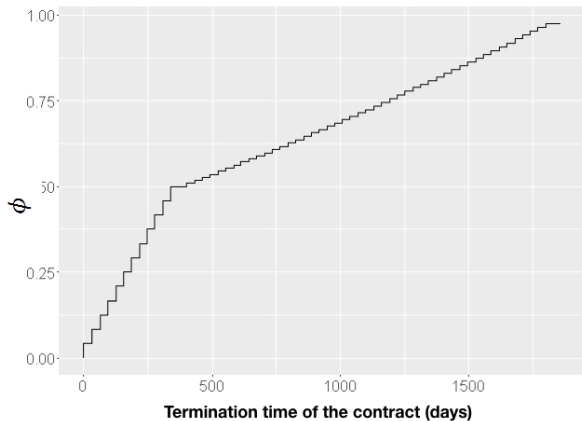
Introduction

Imagine if you were an insurance broker



Commissioning

ϕ : Commissioning function of the insurance broker (per unit of annual premium)



Mathematical Formulation

- T : Termination time of the contract (may be censored)
- C : Censoring time
- $X \in \mathbb{R}^d$: Covariates about the prospect : 6 covariates

Observations

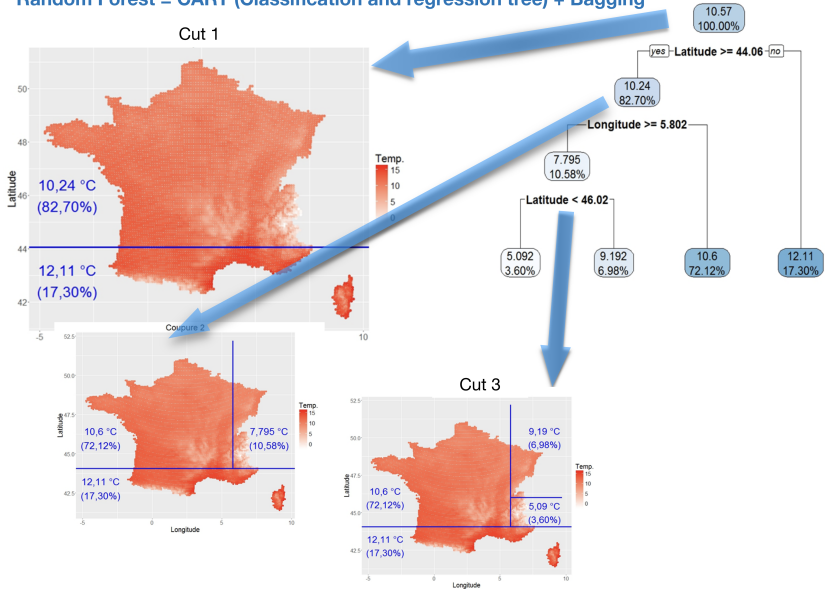
We observe $(Y_i, \delta_i, X_i)_{1 \leq i \leq n}$ i.i.d. with :

- $Y = \min(T, C)$
 - $\delta = \mathbb{1}_{T \leq C}$
-
- Goal : Build a model for $f(x) = E[\phi(T)|X = x]$

Weighted Random Forest and IPCW method

Random Forest

Random Forest = CART (Classification and regression tree) + Bagging



To solve our problem

- We want to estimate $f(x) = E[\phi(T)|X = x]$
- We know :

$$f = \underset{g}{\operatorname{argmin}} E [(\phi(T) - g(X))^2] \quad (1)$$

⇒ Need an estimate of $E [(\phi(T) - g(X))^2]$ with T censored

⇒ More generally, for any bounded ψ , we can estimate $E [\psi(T, X)]$ with T censored using IPCW principle

- IPCW : Inverse Probability of Censoring Weighting

IPCW principle]

Let $p(t, x) = P(\delta = 1 | T = t, X = x)$

Then for any bounded function ψ ,

$$E [W \cdot \psi(Y, X)] = E [\psi(T, X)] \text{ with } W = \frac{\delta}{p(Y, X)}$$

Reminder

- $Y = \min(T, C)$
- $\delta = \mathbb{1}_{T \leq C} = \mathbb{1}_{Y=T}$

Proof

$$\begin{aligned} E \left[\frac{\delta}{p(Y, X)} \cdot \psi(Y, X) \right] &= E \left[\frac{\delta}{p(T, X)} \cdot \psi(T, X) \right] \\ &= E \left[\frac{\psi(T, X)}{p(T, X)} \cdot \underbrace{E[\delta | T, X]}_{p(T, X)} \right] \\ &= E [\psi(T, X)] \end{aligned}$$

Hypothesis

H1 : $P(T \leq C|X, T) = S_C(T)$ (true if $C \perp\!\!\!\perp (T, X)$)

H2 : $P(T \leq C|X, T) = S_C(T|X)$ (true if $C \perp\!\!\!\perp T$ conditionally on X)

- Under **H1** :

$$p(t, x) = P(t \leq C|T = t, X = x) = S_C(t)$$

- Under **H2** :

$$p(t, x) = P(t \leq C|T = t, X = x) = S_C(t|X = x)$$

Weighted Random Forest

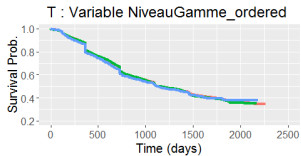
- Let $\hat{W}_i = \frac{\delta_i}{\hat{\mathcal{S}}_c(Y_i)}$ or $\frac{\delta_i}{\hat{\mathcal{S}}_c(Y_i|X_i)}$
- We estimate $E[(\phi(T) - g(X))^2]$ by

$$\frac{1}{n} \sum_{i=1}^n \hat{W}_i \cdot (\phi(Y_i) - g(X_i))^2$$

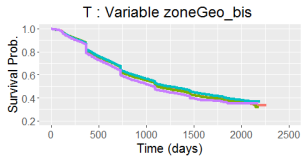
- **Weights are taken into account in the bootstrap of the Random Forest**

Experiments

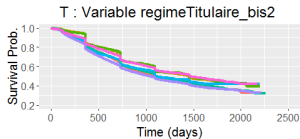
Survival curves by subgroup of individuals



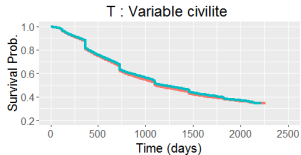
Levels : lev_1 lev_2 lev_3



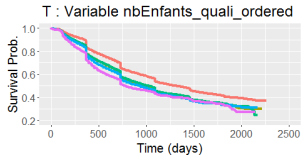
Levels : lev_1 lev_2 lev_3 lev_4



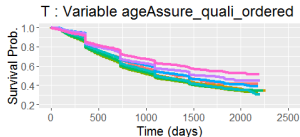
Levels : lev_1 lev_2 lev_3 lev_4 lev_5 lev_6 lev_7



Levels : lev_1 lev_2



Levels : lev_1 lev_2 lev_3 lev_4 lev_5



Levels : lev_1 lev_2 lev_3 lev_4 lev_5 lev_6 lev_7 lev_8

Setting of the Experiments

Weighted RF

Train data	ϕ	$\widehat{W} = \delta / \widehat{S}_c(\cdot X)$
.....		

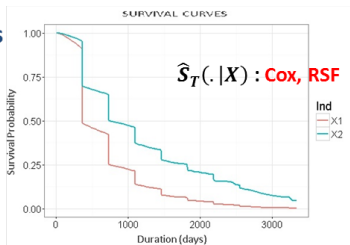
$$\widehat{S}_c(\cdot | X) \begin{cases} KM \\ RSF \\ Cox \end{cases}$$

$\widehat{\phi} = \widehat{f}(x)$: prediction of the random forest

$$\widehat{f} = \underset{g}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \widehat{W}_i (\phi(Y_i) - g(X_i))^2$$

Alternative methods

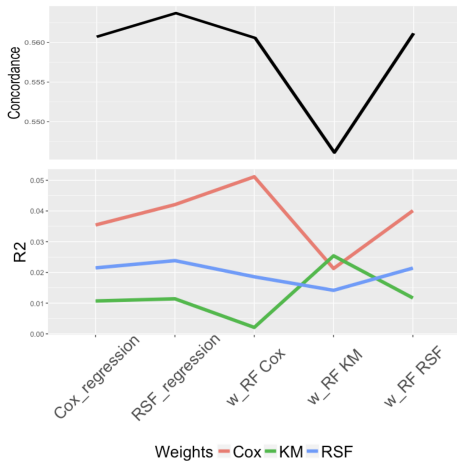
Train data
.....



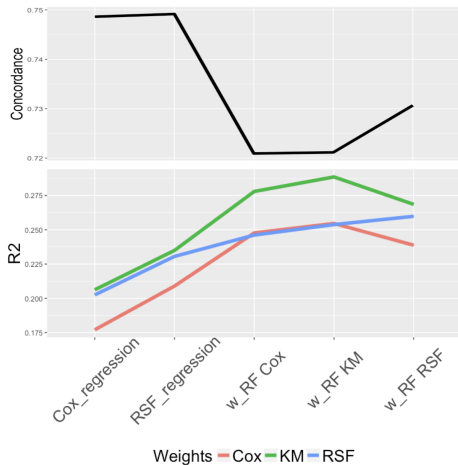
$$\widehat{\phi} = - \int \phi d\widehat{S}_T(\cdot | X)$$

Results

Insurance broker data



Transplant data



Summary

- We can adapt the Random Forest algorithm to the case where the target Y is censored using IPCW principle.
- Weighted Random Forest is competitive with other standard methods

Outlook

- Implementation of the method in a R package
- Theoretical study of the consistency of the method

Acknowledgments

Thanks to Olivier Lopez (UPMC, Paris, France), Guillaume Gerber (Forsides, Paris, France) and Michael Trupin (Groupe Santiane, Paris, France)

Thank you for listening

mail : yohannlefa@gmail.com