

**100% ACTUAIRES &
100% DATA SCIENCE**

INSTITUT DES
ACTUAIRES



Vers un modèle de tarification ML interprétable : Modèle de fréquence basé sur des données telematics

Arthur MAILLART : Université Lyon 1/Forsides
Christian ROBERT : ENSAE/Chaire DAMI

29 / NOV / 2019

Hôtel Marriott Rive Gauche
Paris 14ème

1. Modèle de fréquence sinistre
2. Améliorer la connaissance de ses données grâce aux arbres
3. Extraire un ensemble parcimonieux de règles
4. Use case
5. Résultats
6. Bibliographie

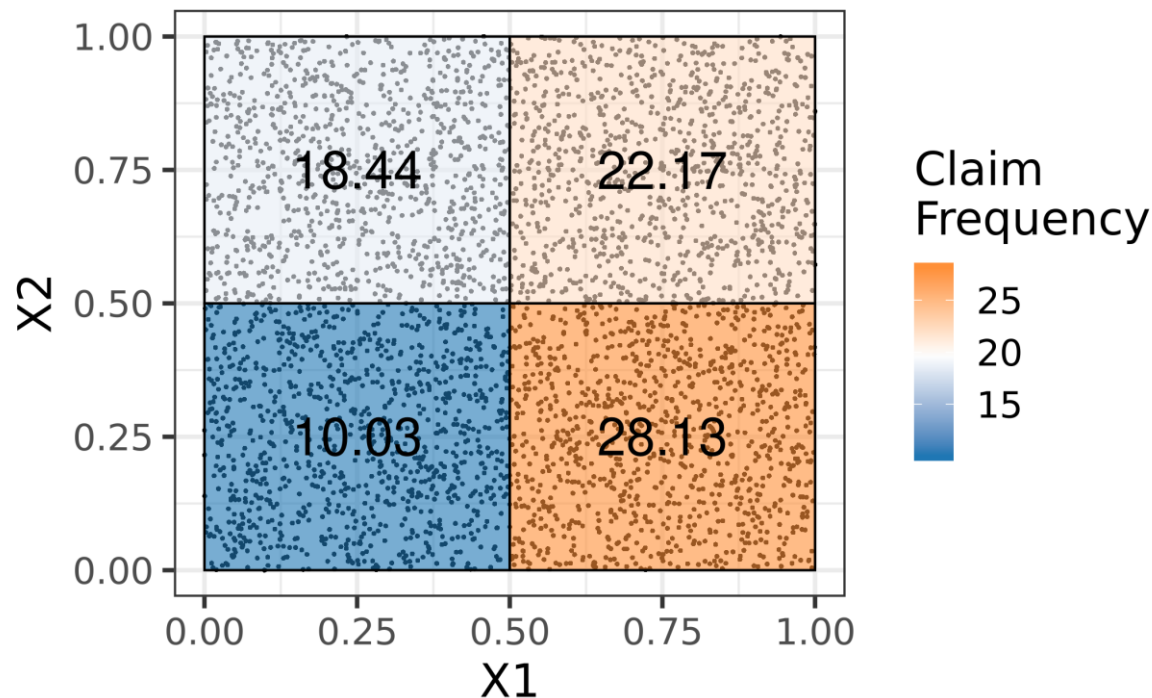
En modélisant la fréquence sinistre, on cherche à créer des groupes homogènes de risques que l'on appelle **cellules tarifaires** [Ohlsson and Johansson, 2010].

Ces dernières doivent être suffisamment larges pour qu'il y ait une **mutualisation** mais suffisamment **différenciées** pour attirer certains prospects.

Ohlsson, E. and Johansson, B. (2010).

Non-life insurance pricing with generalized linear models, volume 2.
Springer.

Pour représenter notre problème en deux dimensions et illustrer les méthodes qui vont suivre, nous créons un jeu de données simulées.



Nous avons un jeu de données *telematics* collectées dans le cadre d'une assurance PAYD et nous voulons ajuster un **modèle de fréquence**. Une méthode classique en actuariat consiste à ajuster un modèle GLM pour obtenir ce modèle de fréquence. Cependant, lorsqu'il y a des variables continues dont la relation avec la fréquence sinistre est non-linéaire, il devient nécessaire de discrétiser la variable explicative.

- Souvent, l'ajustement de ces modèles nécessite de la "connaissance d'expert" pour découper correctement ces variables.
- Il n'y a aucune garantie que le découpage est optimal.
- Cette opération peut être très longue selon le nombre de colonnes.

- Les méthodes de Machine Learning telles que Random Forest ou XGBoost constituent de bonnes **heuristiques** d'apprentissage pour des données tabulaires. Ce qui signifie qu'elles performant généralement bien sur ce genre de données.
- Elles n'ont pas besoin qu'on leur spécifie une structure a priori et prennent naturellement en compte des interactions entre les variables.
- Elles ne requièrent pas trop de paramétrage.

Seule ombre au tableau, la chaîne de décision de ces méthodes est difficile voire impossible à appréhender pour un humain. Elles sont donc difficiles à expliquer à un public non-expert.

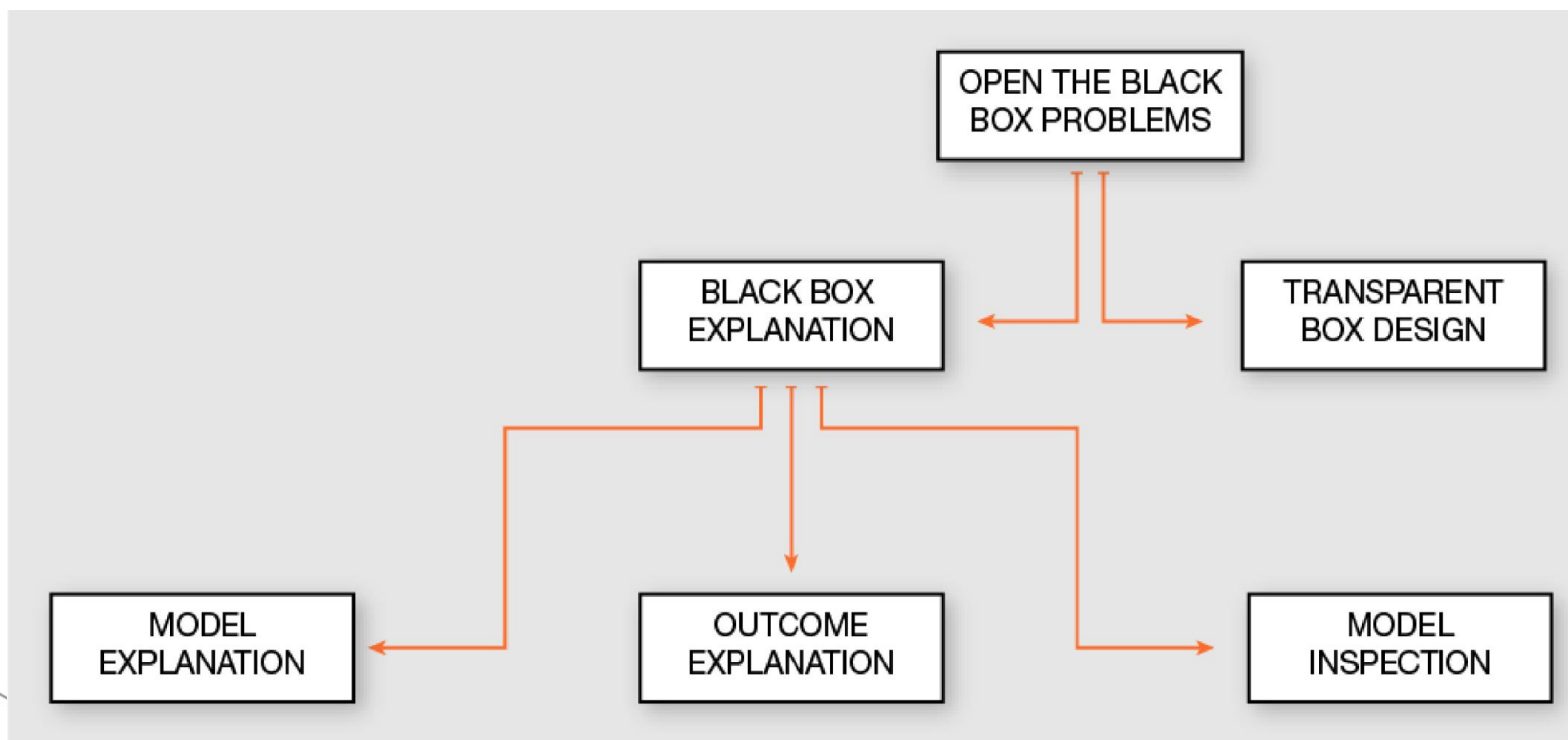
Selon [Guidotti et al., 2018], **interpréter** (en Machine Learning) signifie fournir une **explication** ou le sens d'un concept en des termes compréhensibles par l'humain.

Une **explication** est une interface entre l'humain et la machine. Cette interface doit

- être un bon proxy du modèle Black-Box
- être compréhensible pour l'humain.

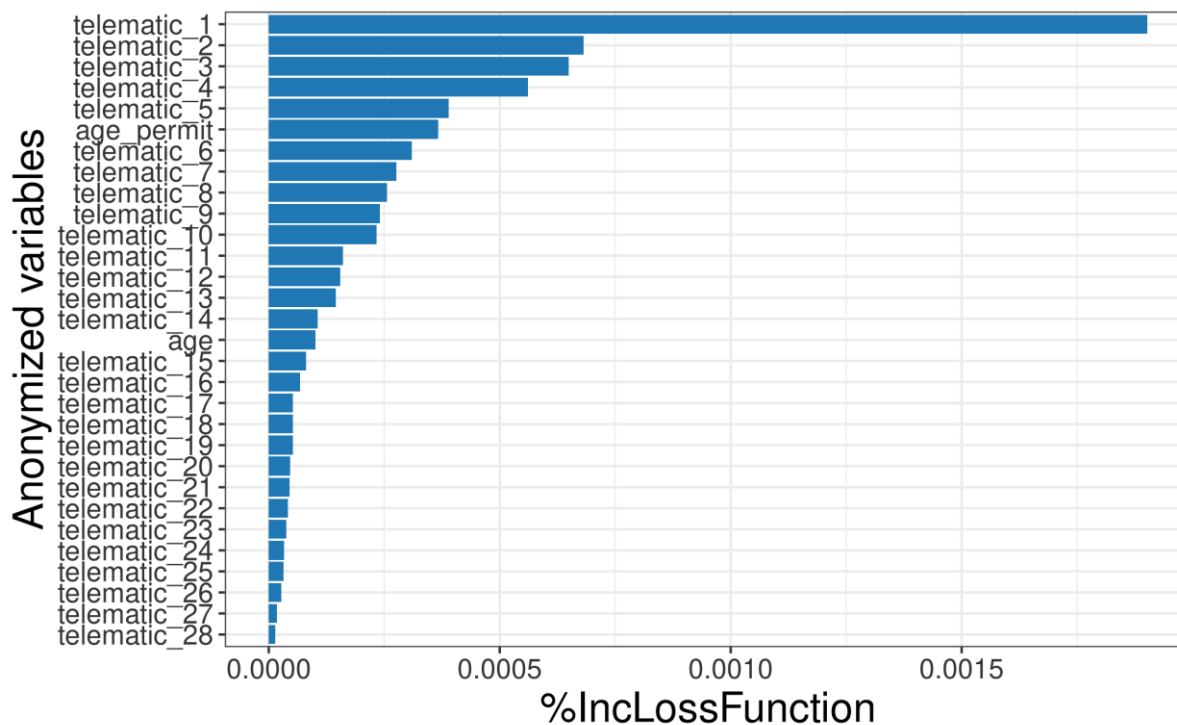
Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., and Pedreschi, D. (2018).
A survey of methods for explaining black box models.
ACM computing surveys (CSUR), 51(5):93.

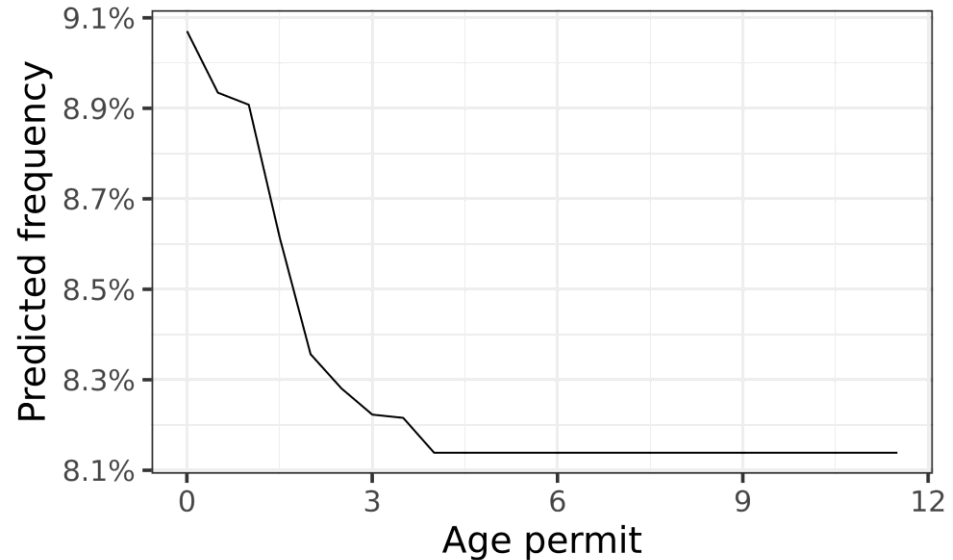
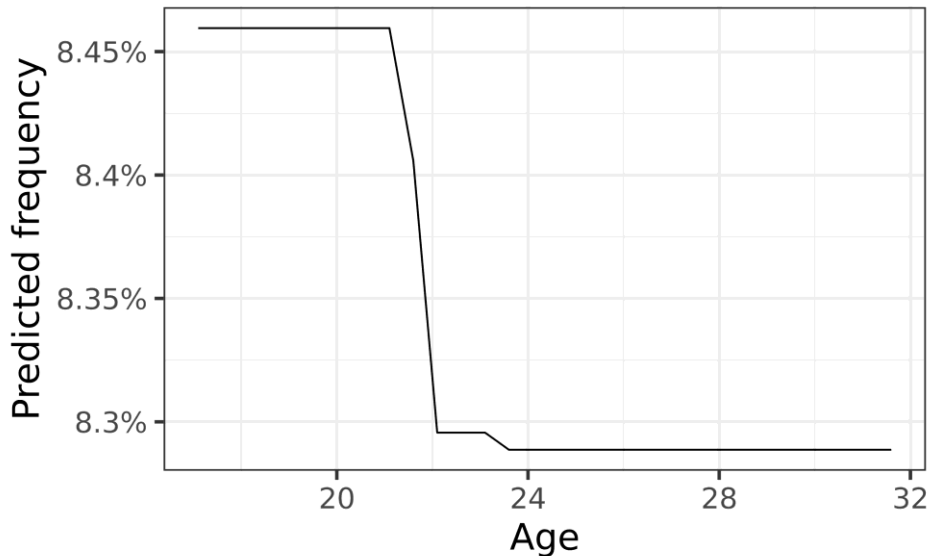
De cela découle 4 types de problèmes que l'on cherche à résoudre pour rendre les modèles interprétables.

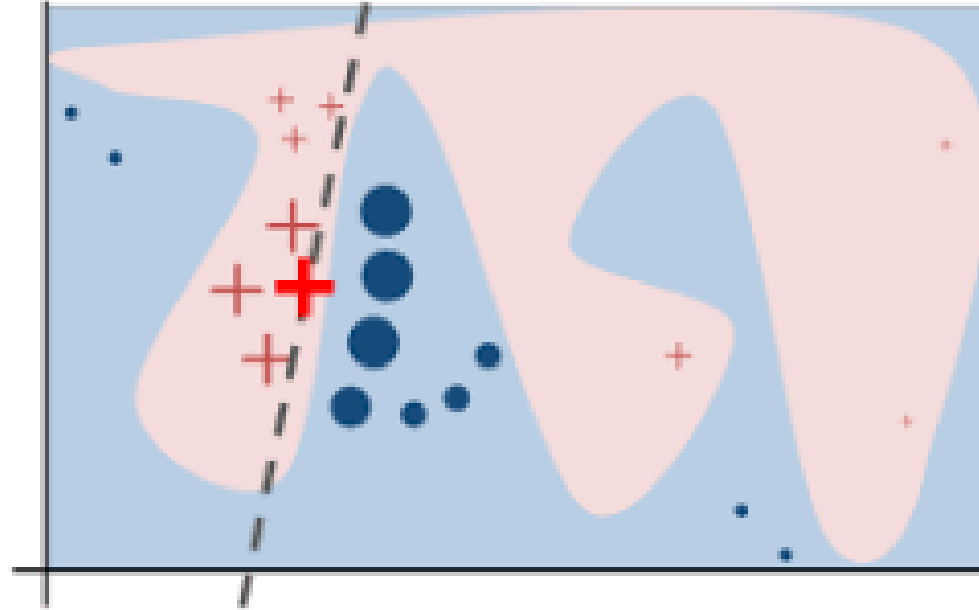


Source : [Guidotti et al., 2018]

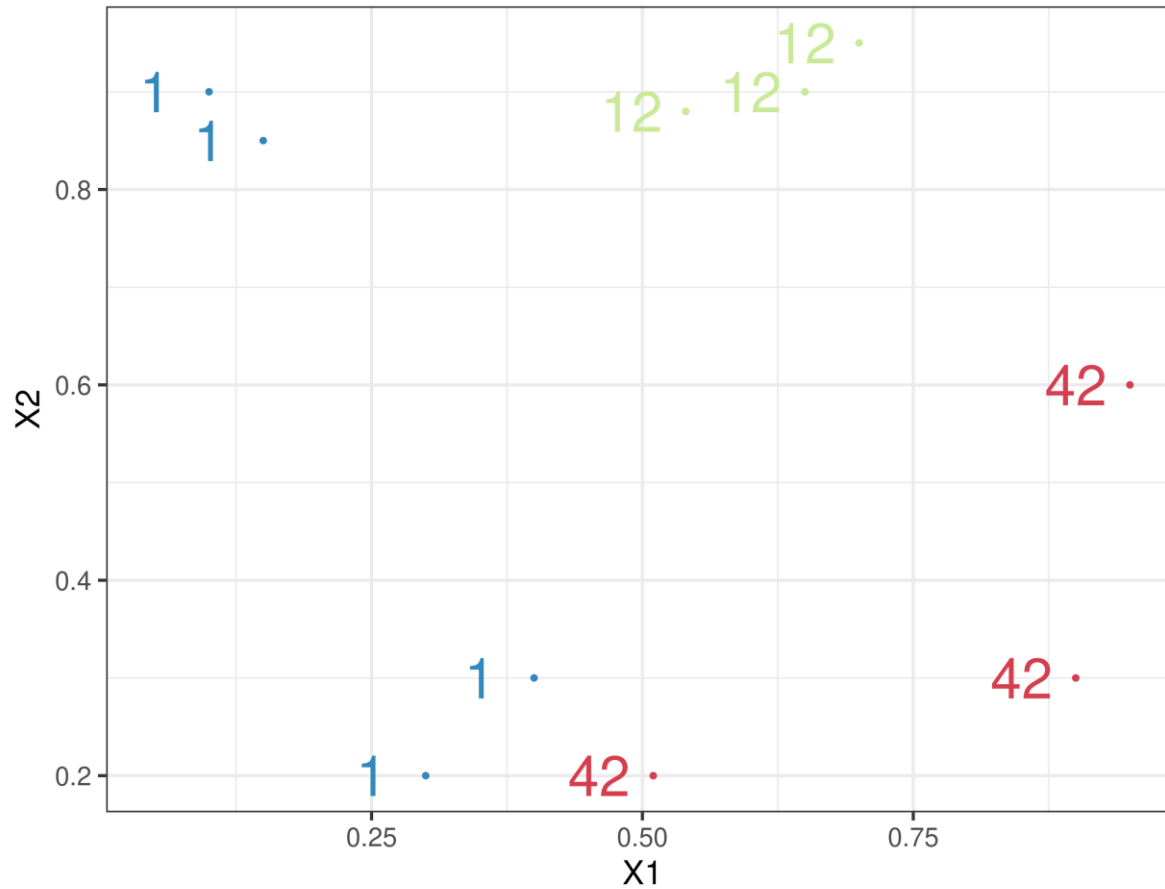
Les méthodes bien connues telles que la "Variable Importance" ou les "Partial Dependence Plots" apportent de l'information mais celle-ci est insuffisante.



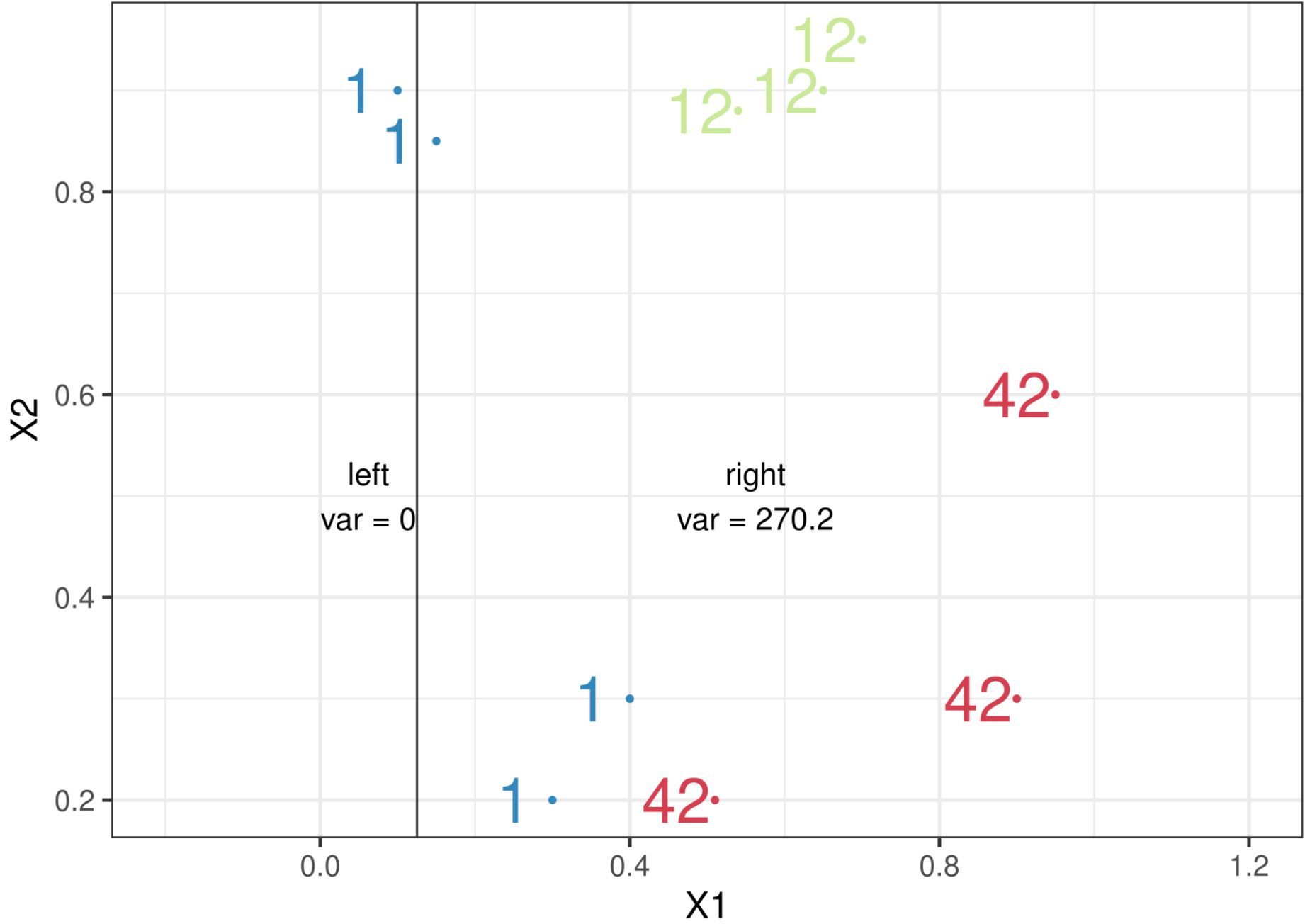




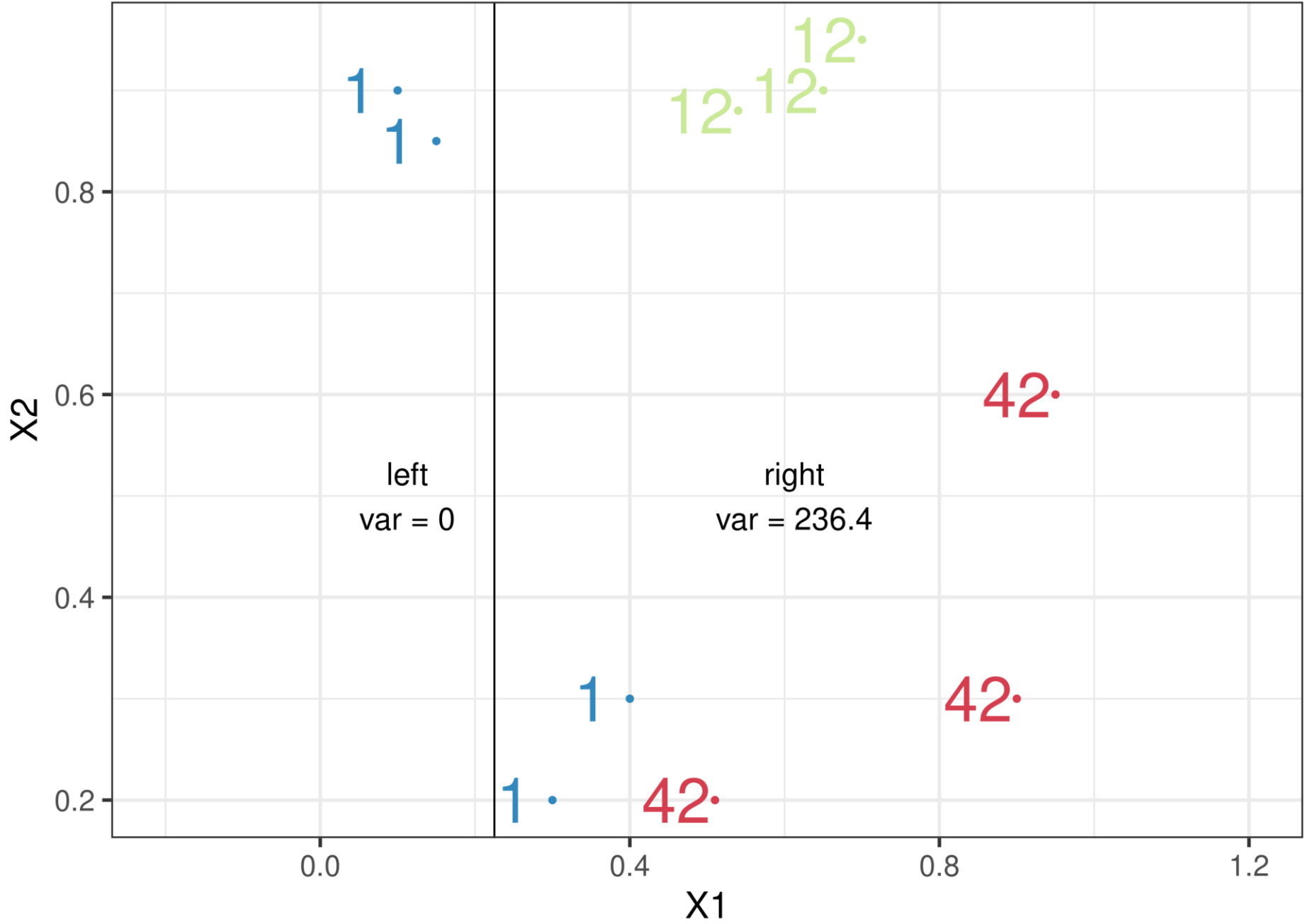
Source : <https://arxiv.org/pdf/1602.04938.pdf>



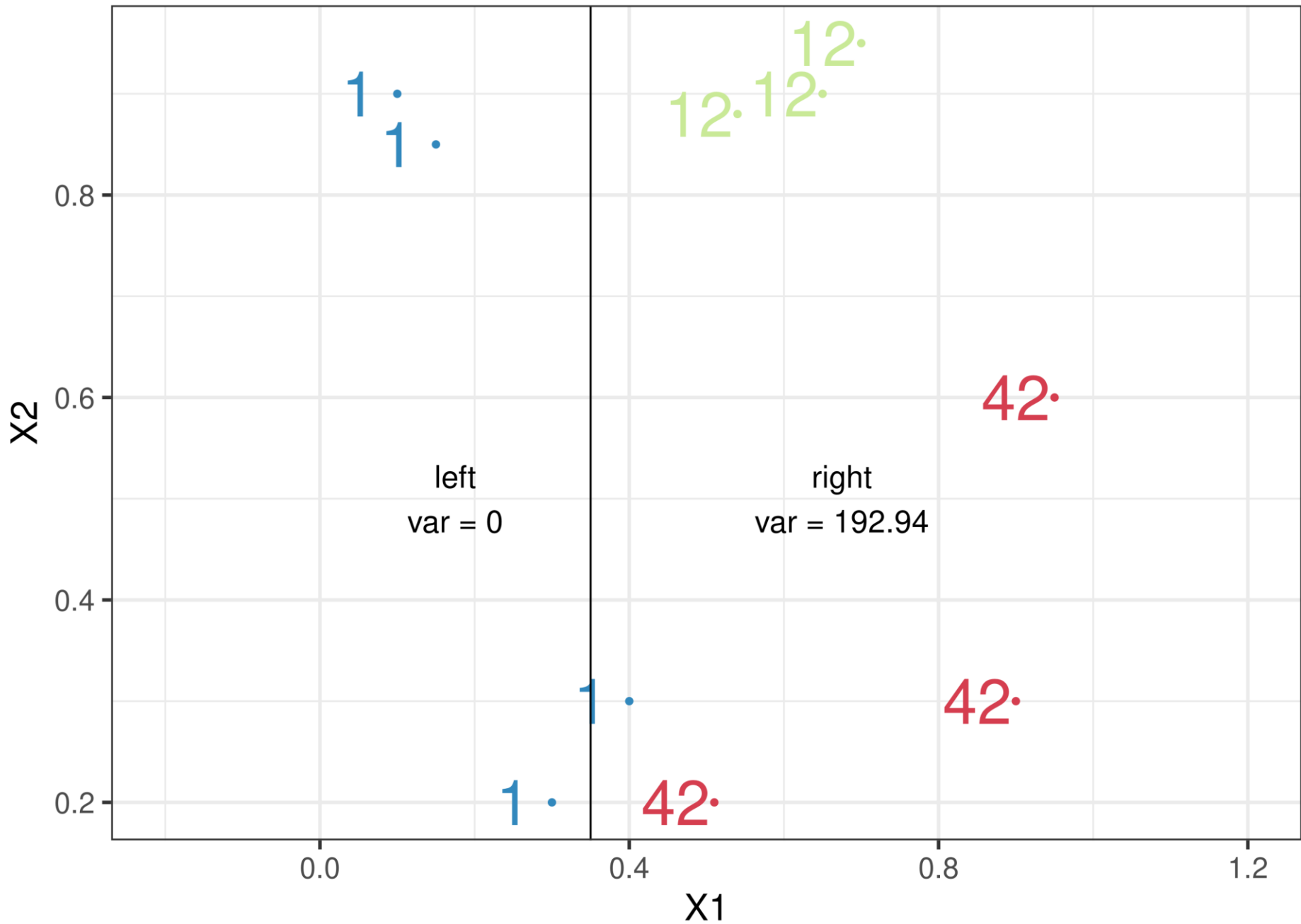
critierion = 297.24 - 0 - 270.2 = 27.04



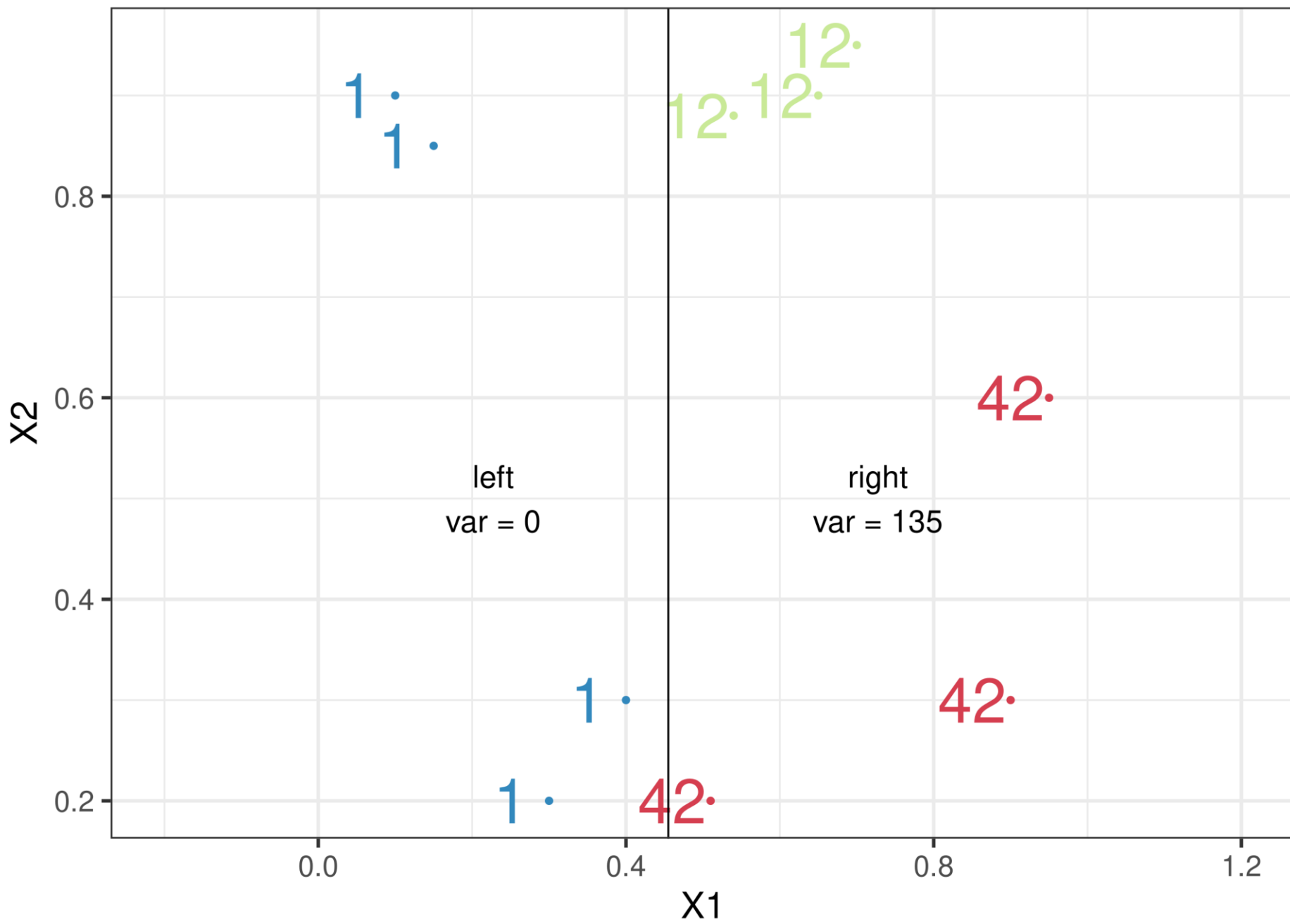
critierion = 297.24 - 0 - 236.4 = 60.84



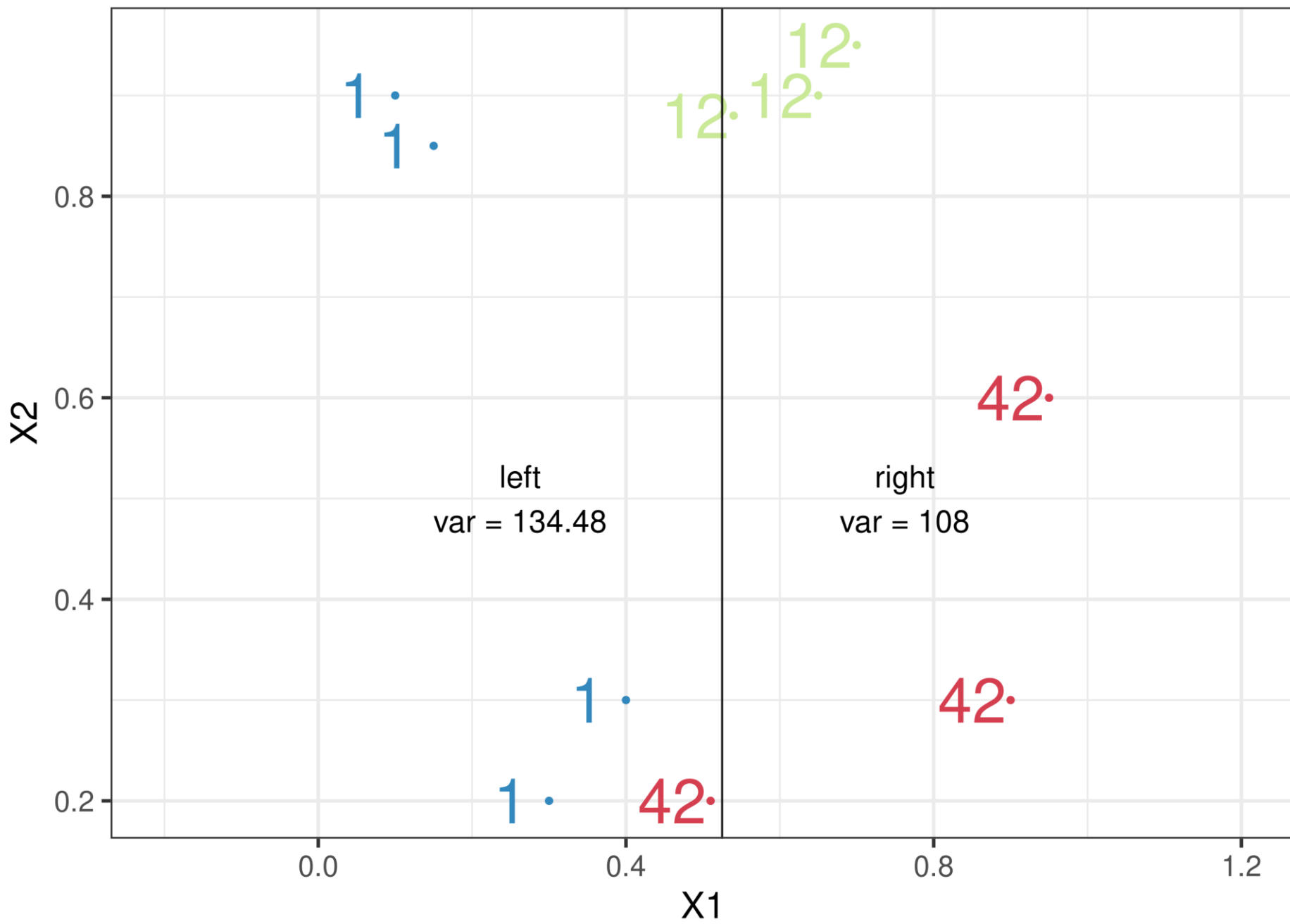
critierion = 297.24 - 0 - 192.94 = 104.3



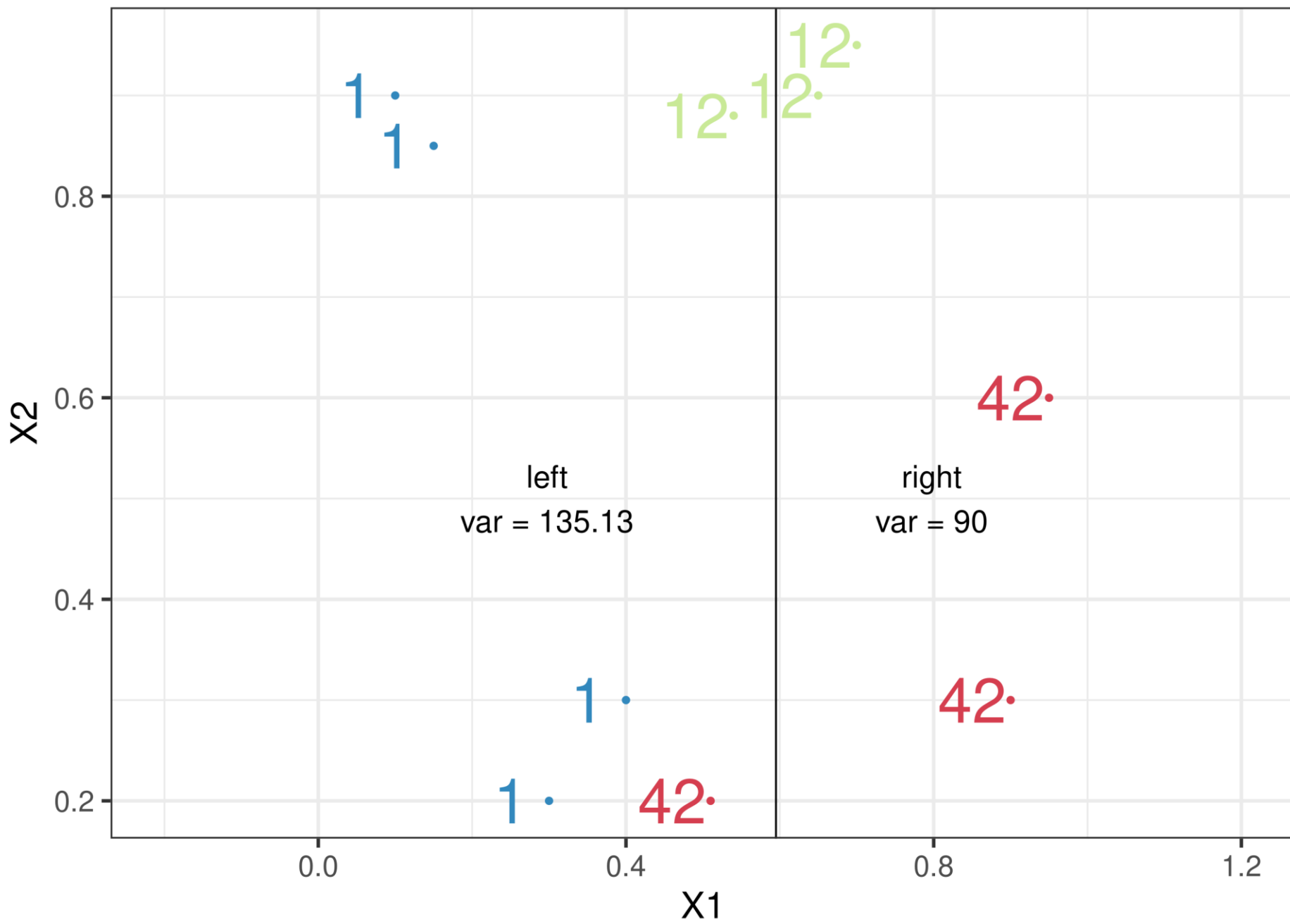
critierion = 297.24 - 0 - 135 = 162.24



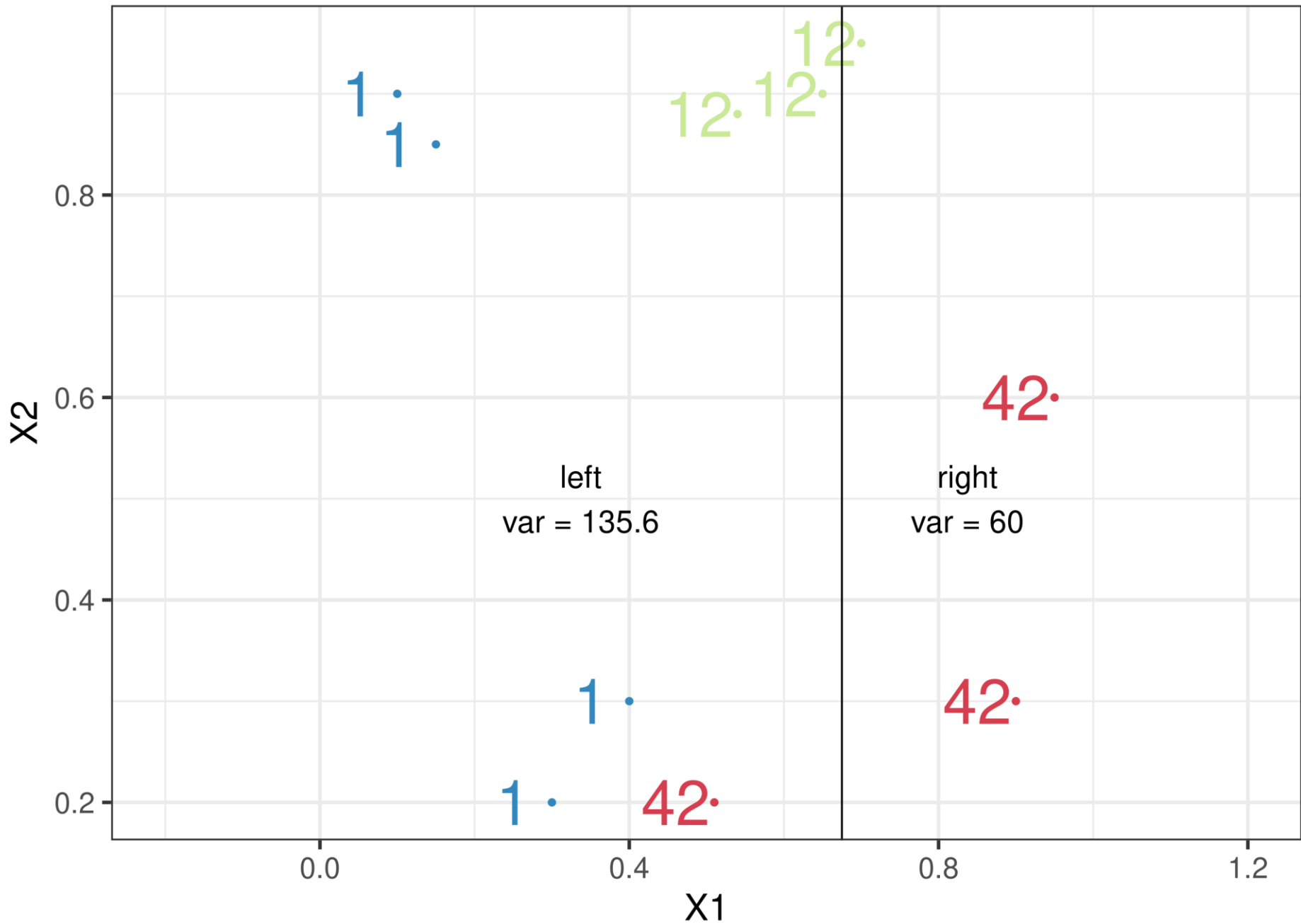
$$\text{criterion} = 297.24 - 134.48 - 108 = 54.76$$



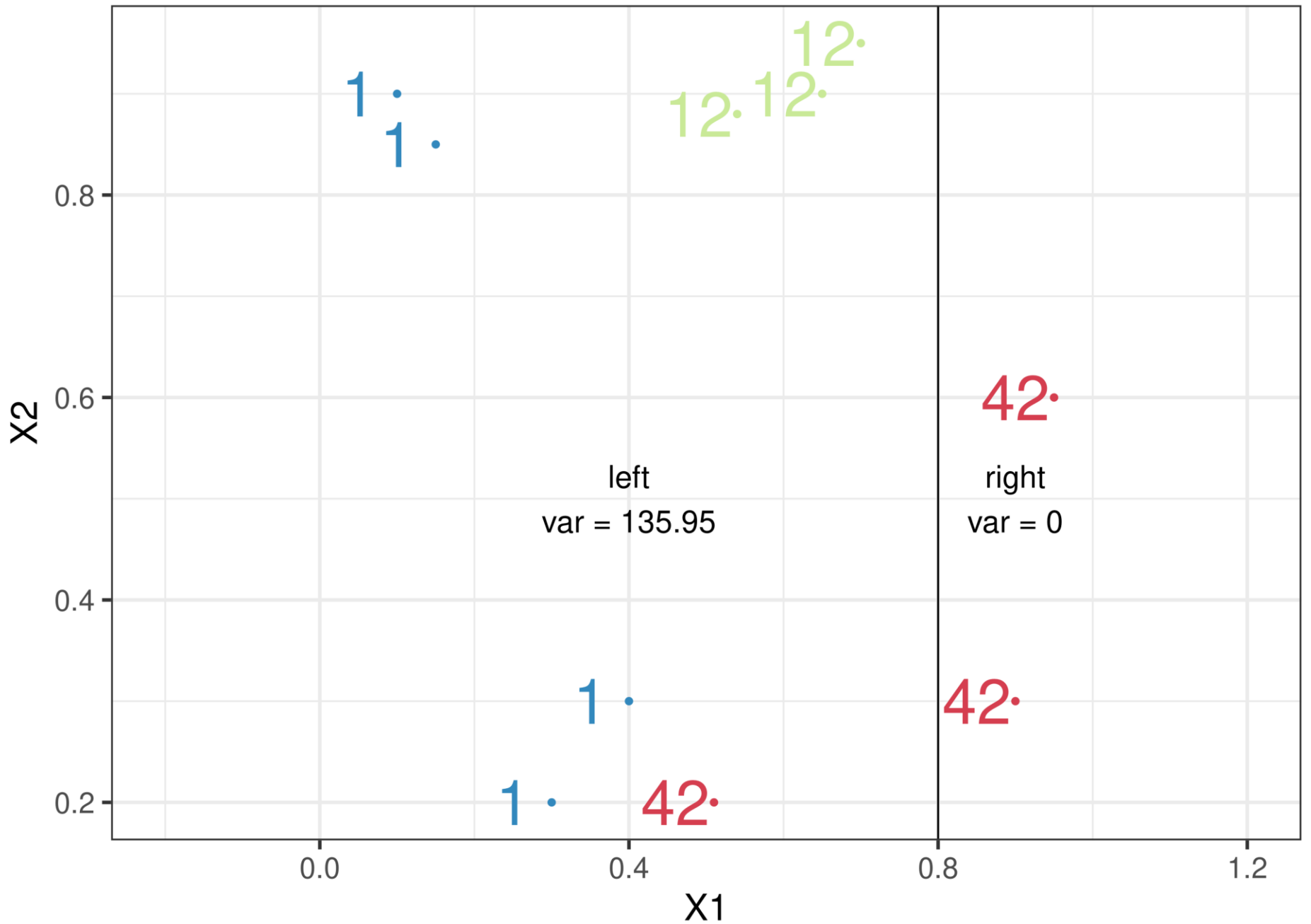
$$\text{criterion} = 297.24 - 135.13 - 90 = 72.11$$



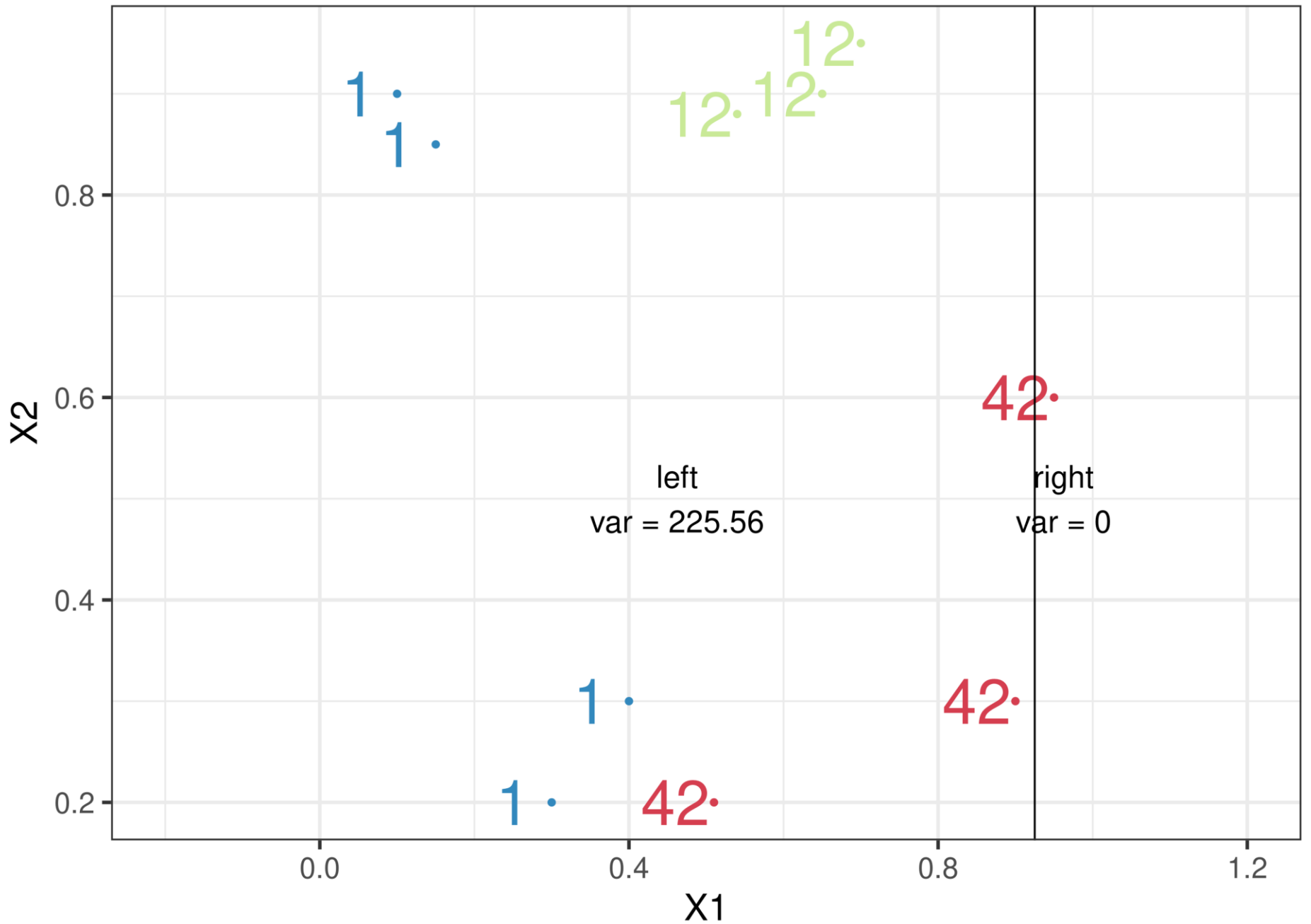
critierion = $297.24 - 135.6 - 60 = 101.64$



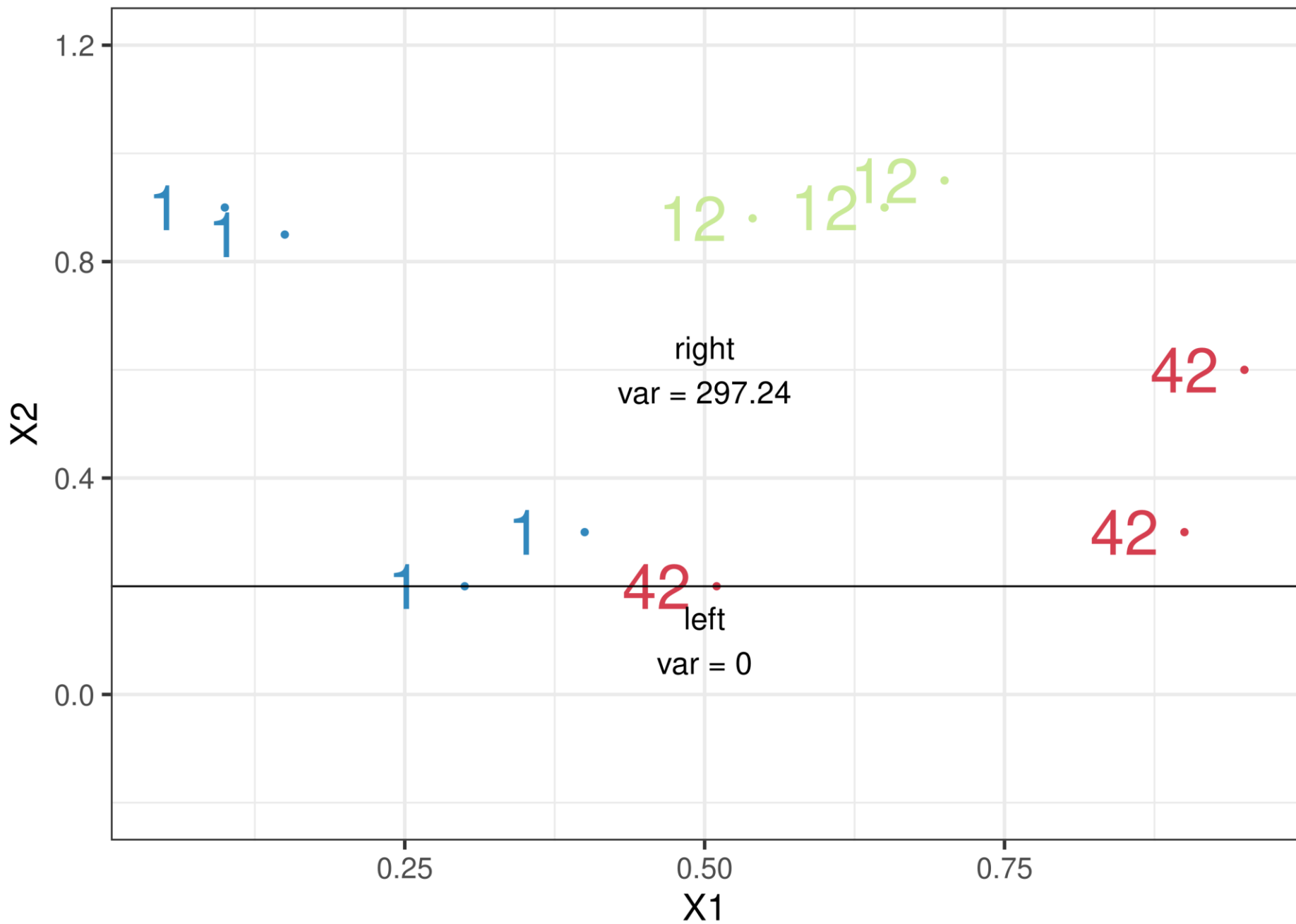
$$\text{criterion} = 297.24 - 135.95 - 0 = 161.29$$



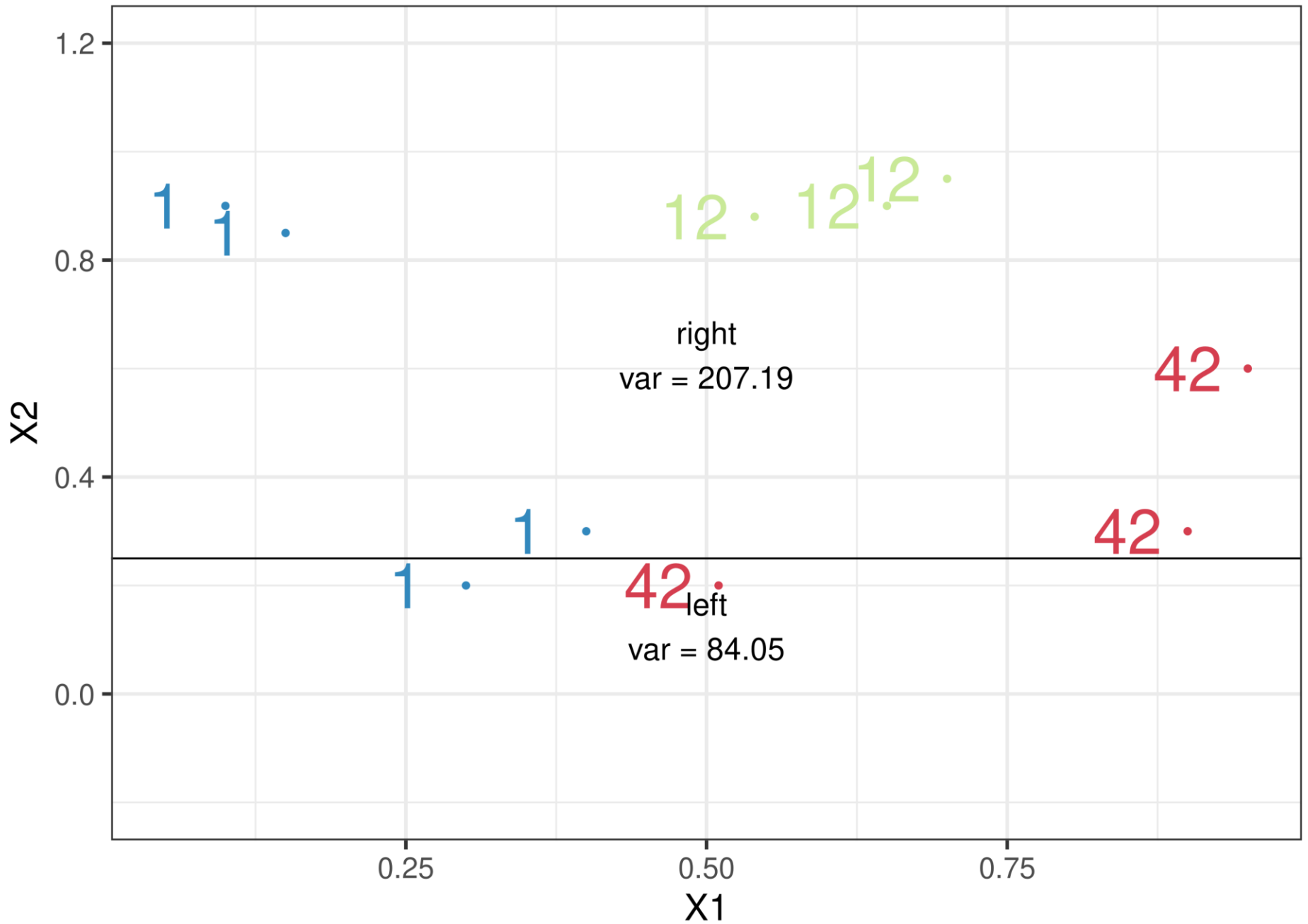
critierion = 297.24 - 225.56 - 0 = 71.68



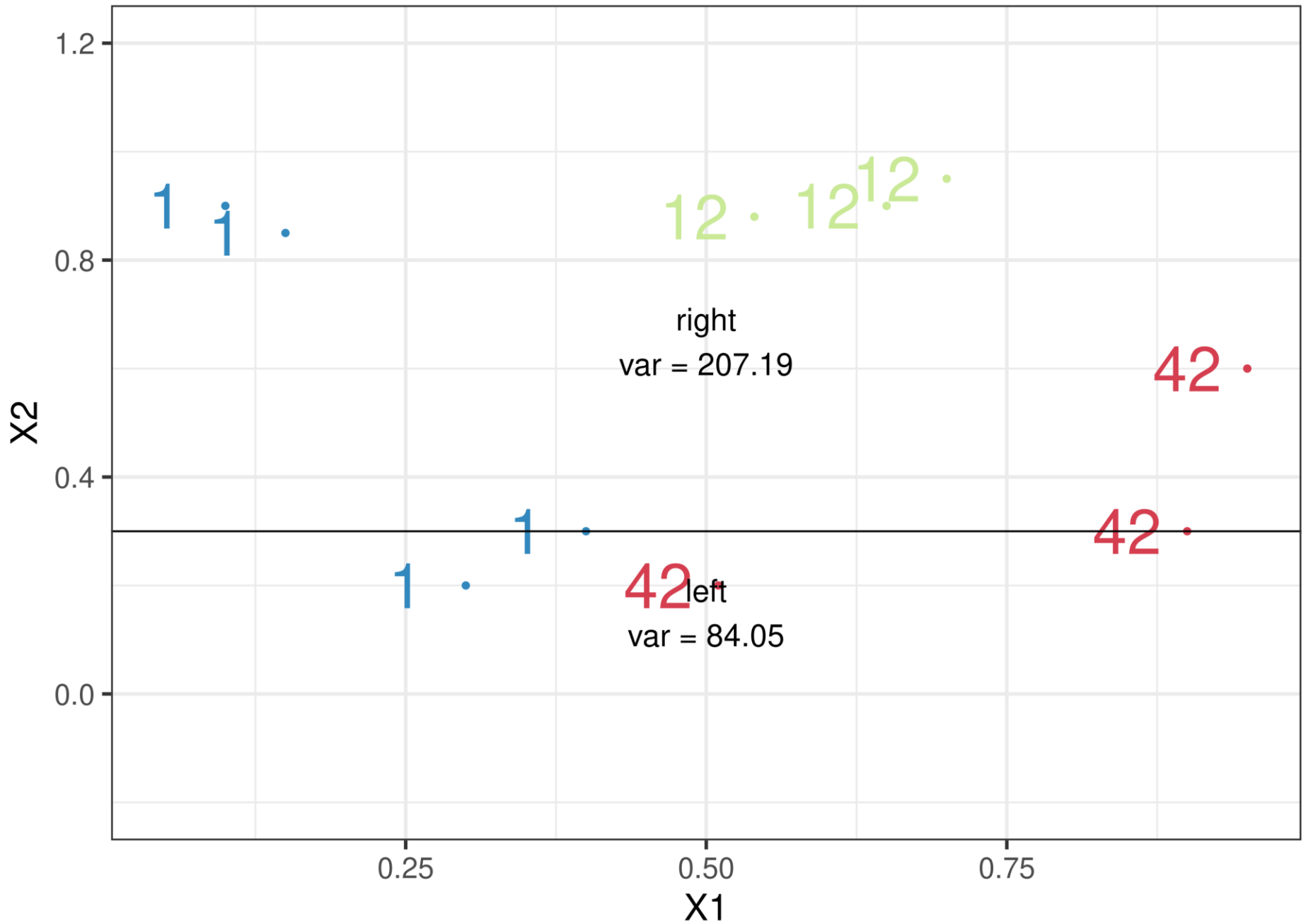
critierion = 297.24 - 0 - 297.24 = 0



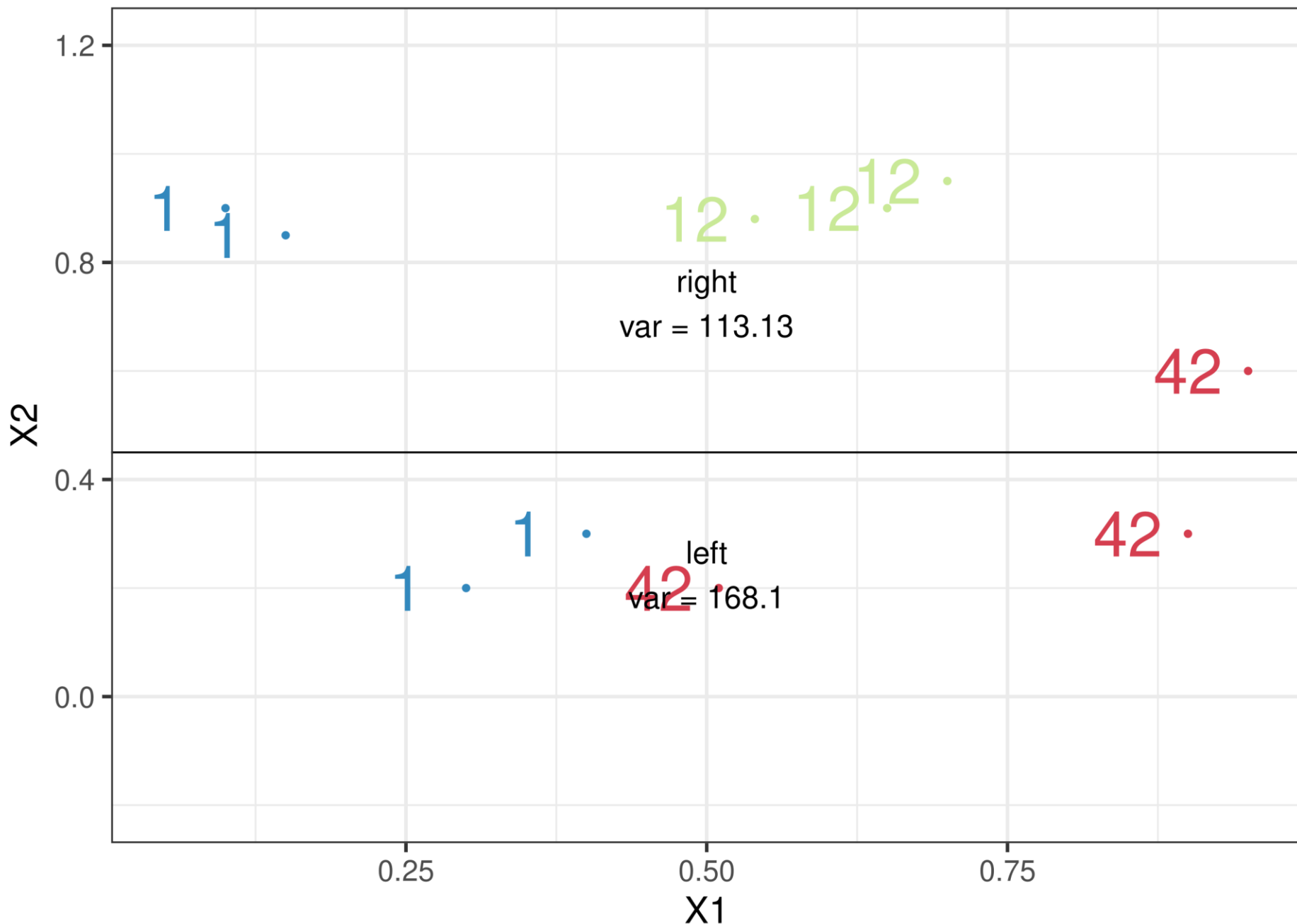
$$\text{criterion} = 297.24 - 84.05 - 207.19 = 6$$



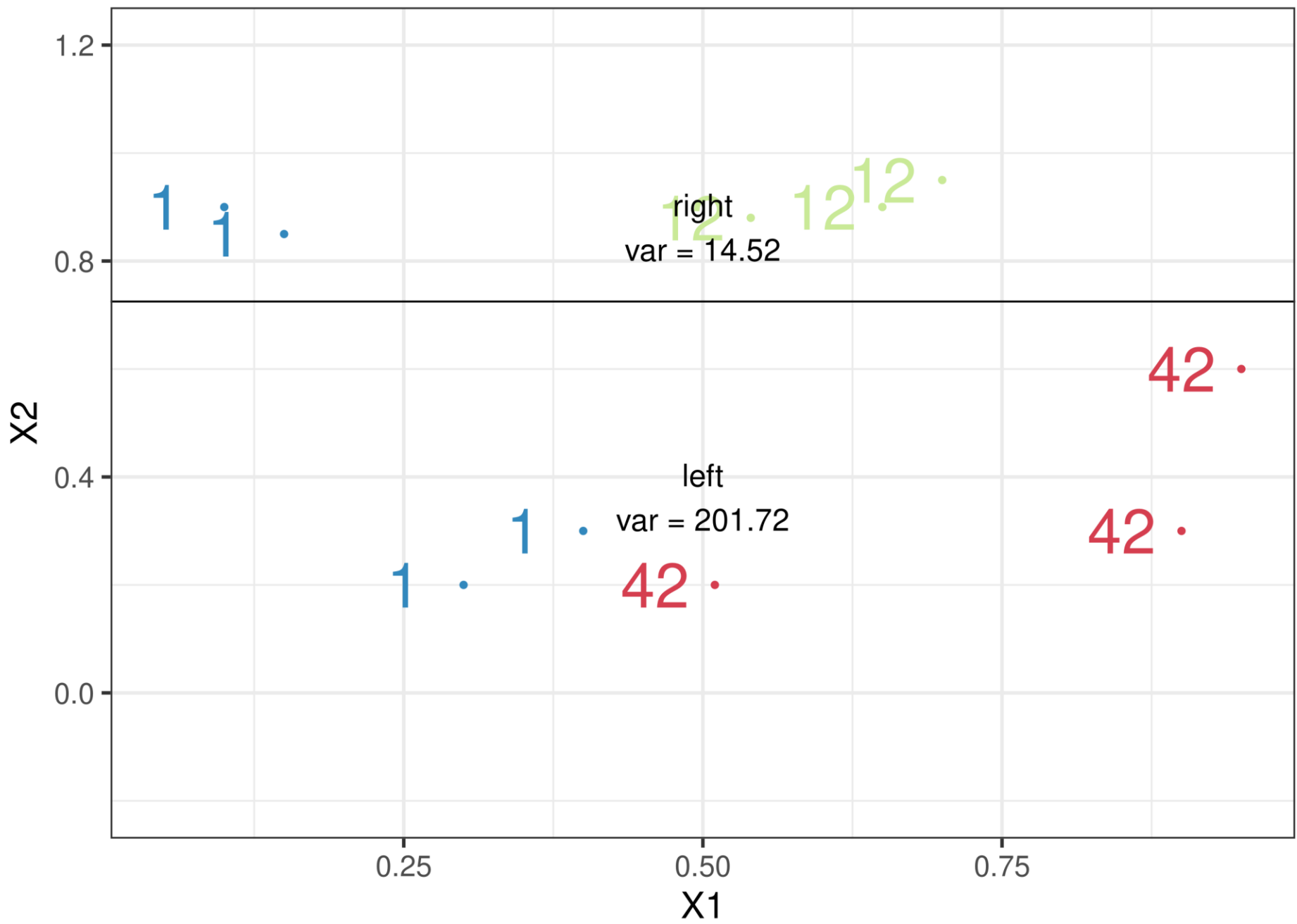
$$\text{criterion} = 297.24 - 84.05 - 207.19 = 6$$



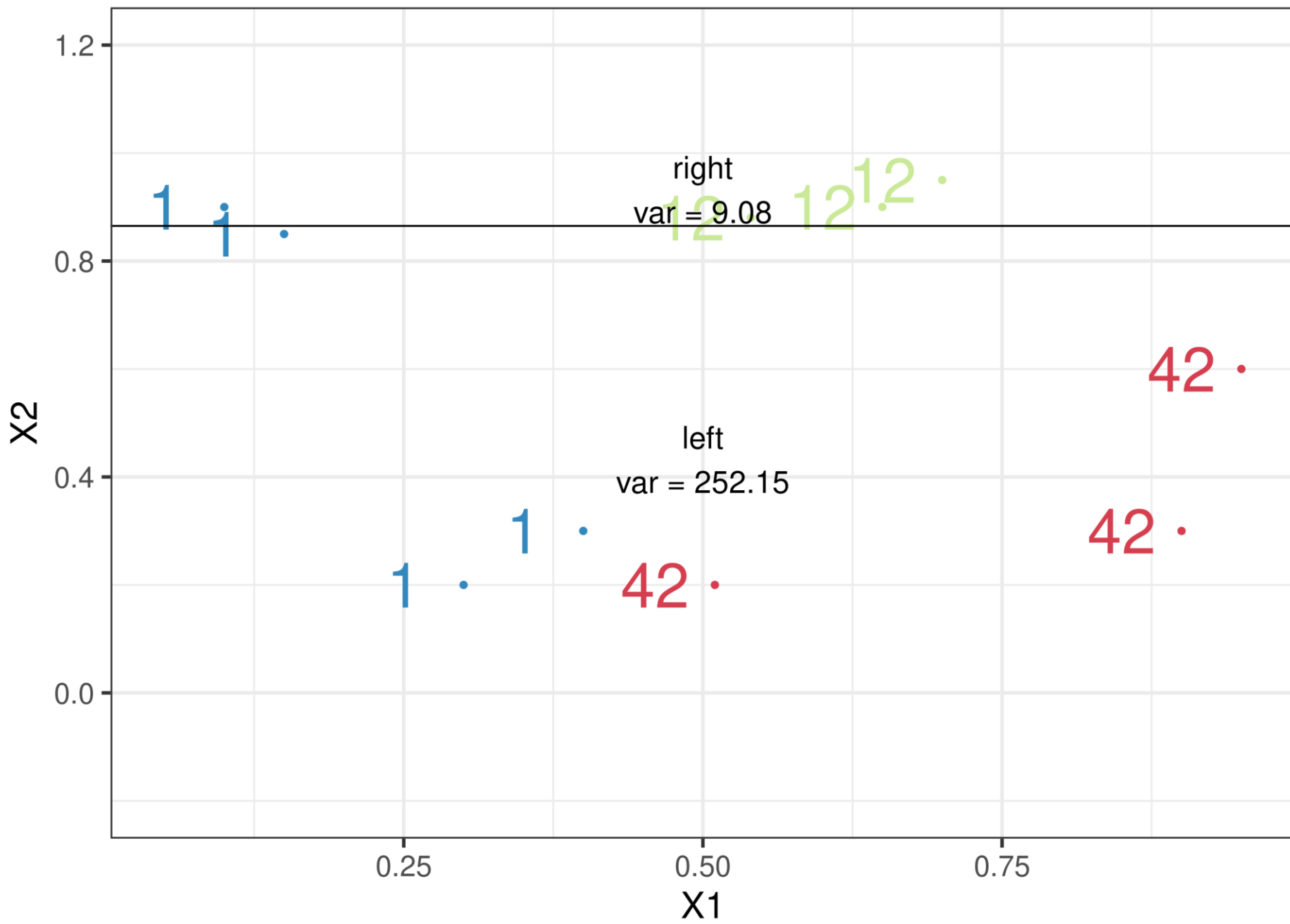
critierion = $297.24 - 168.1 - 113.13 = 16.01$



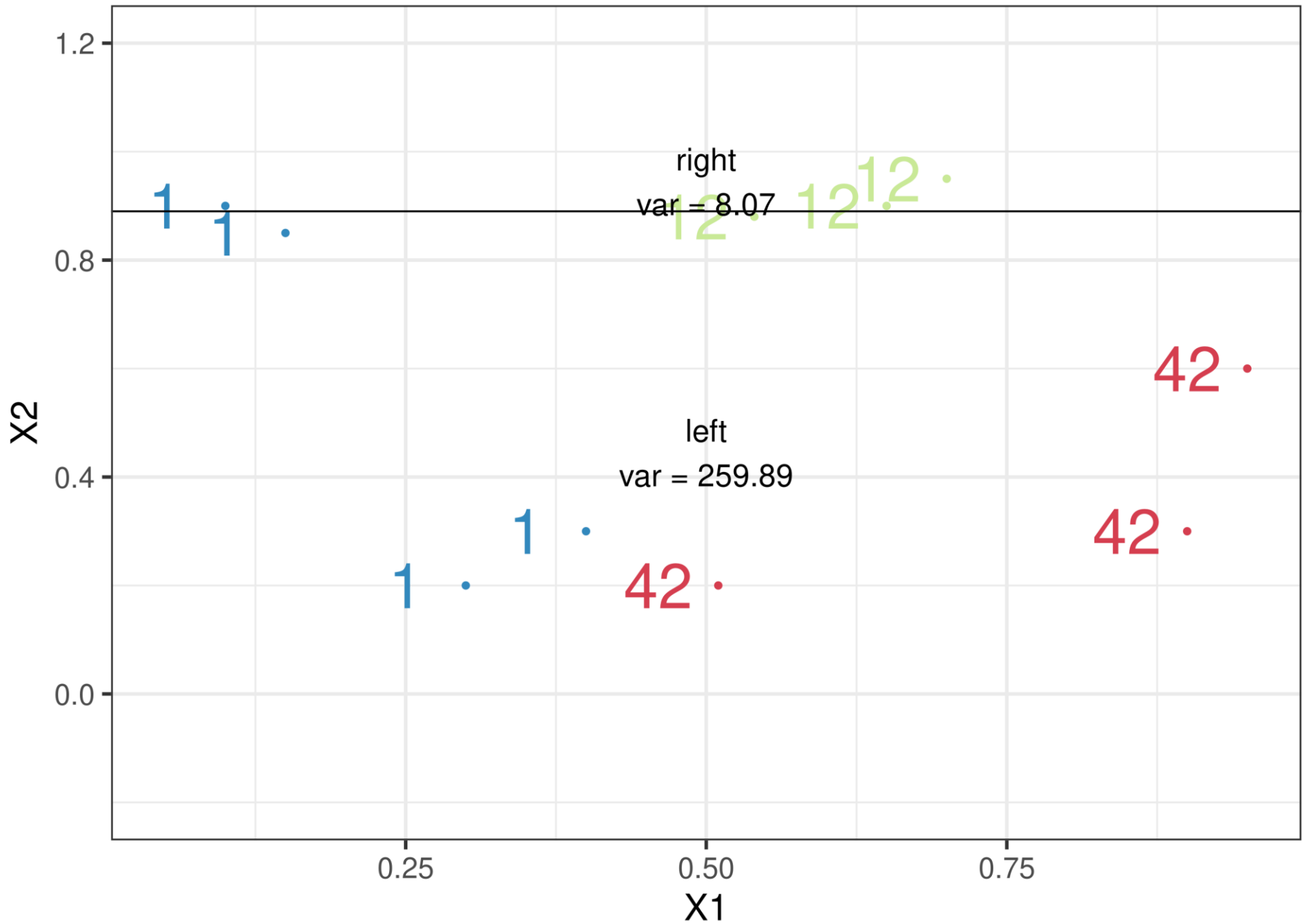
critierion = $297.24 - 201.72 - 14.52 = 81$



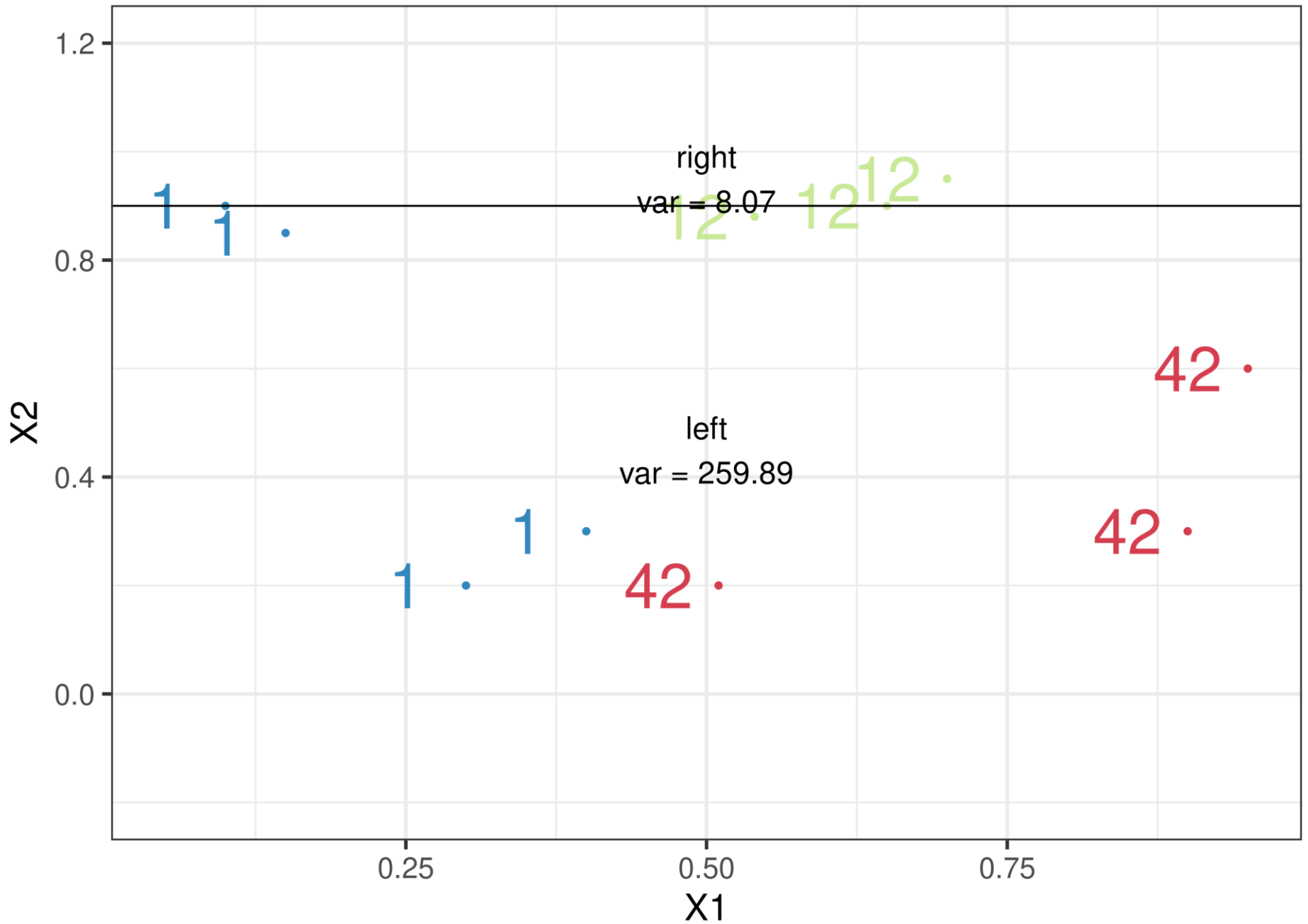
$$\text{criterion} = 297.24 - 252.15 - 9.08 = 36.01$$



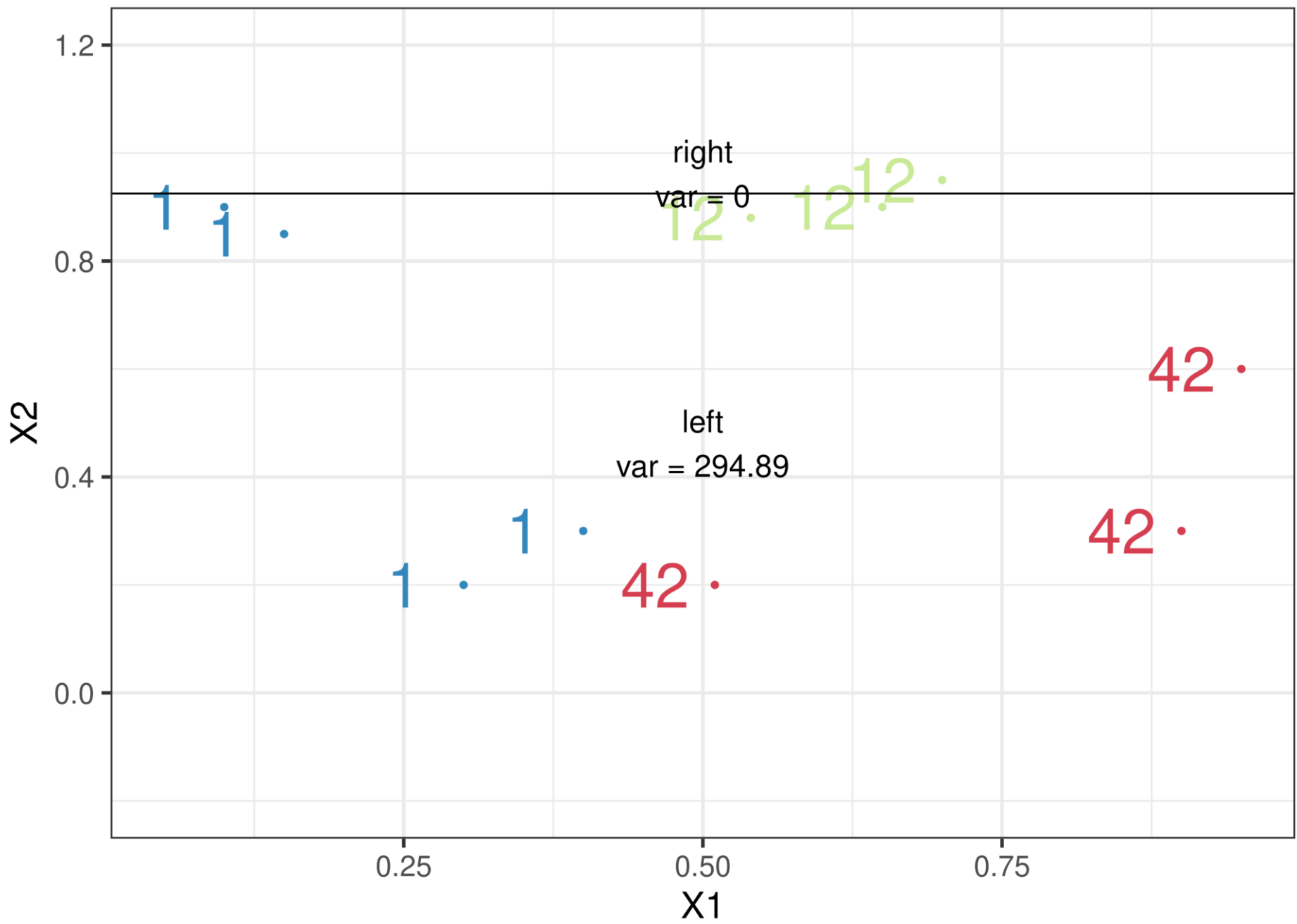
$$\text{criterion} = 297.24 - 259.89 - 8.07 = 29.29$$

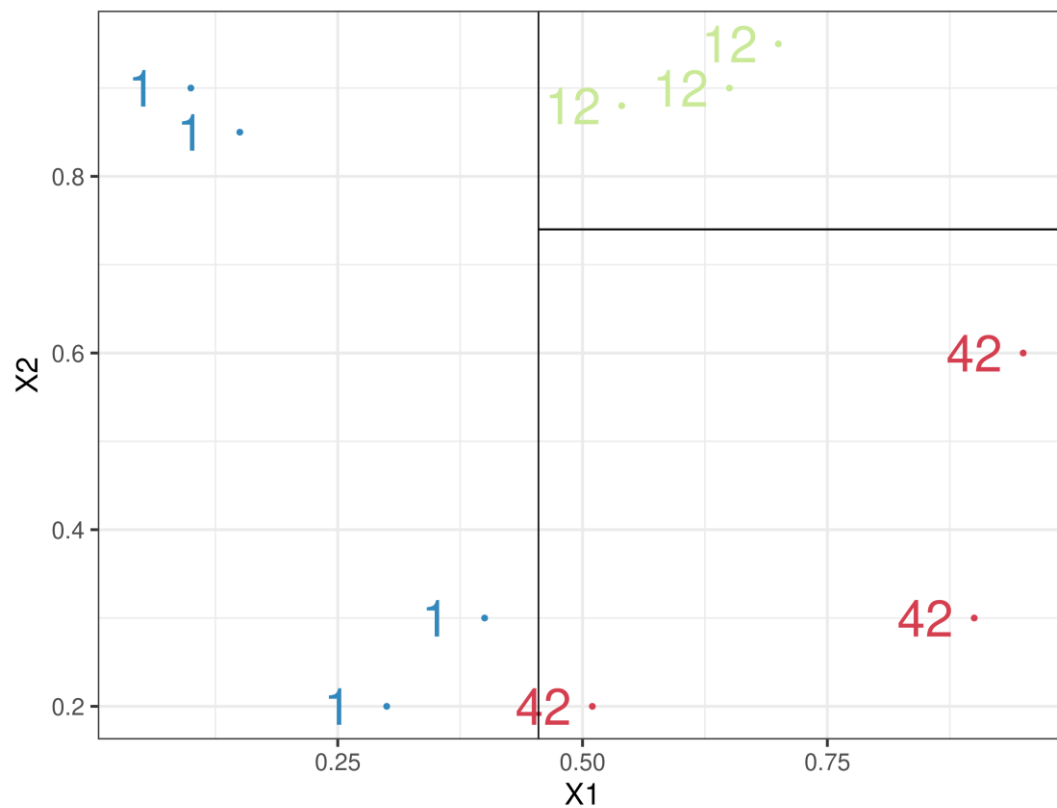
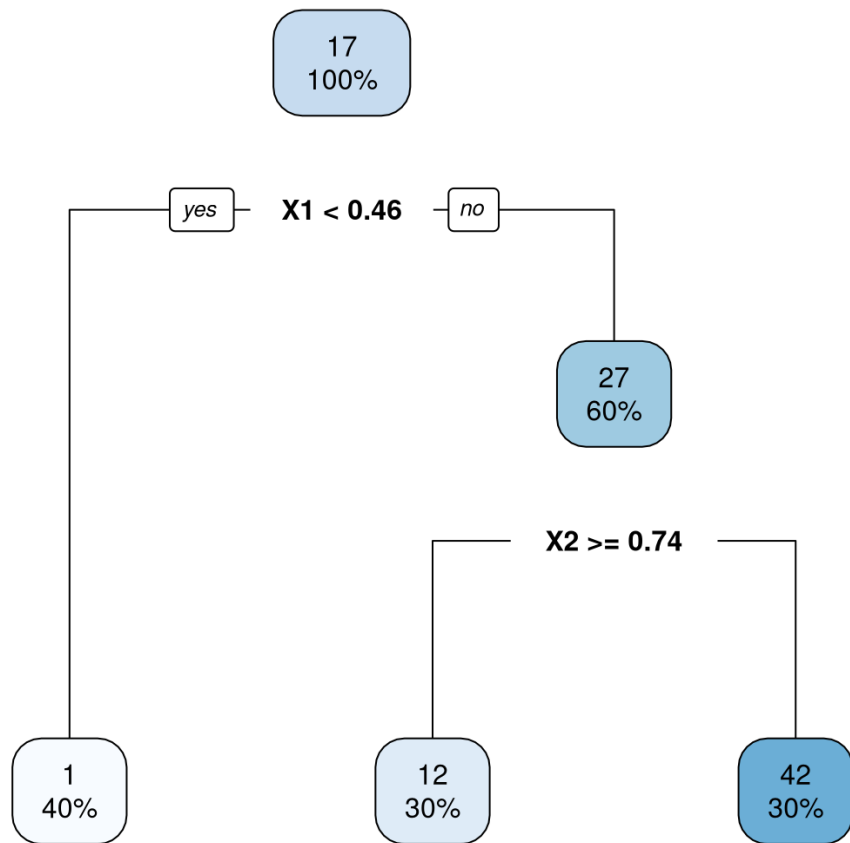


$$\text{criterion} = 297.24 - 259.89 - 8.07 = 29.29$$



critierion = 297.24 - 294.89 - 0 = 2.35





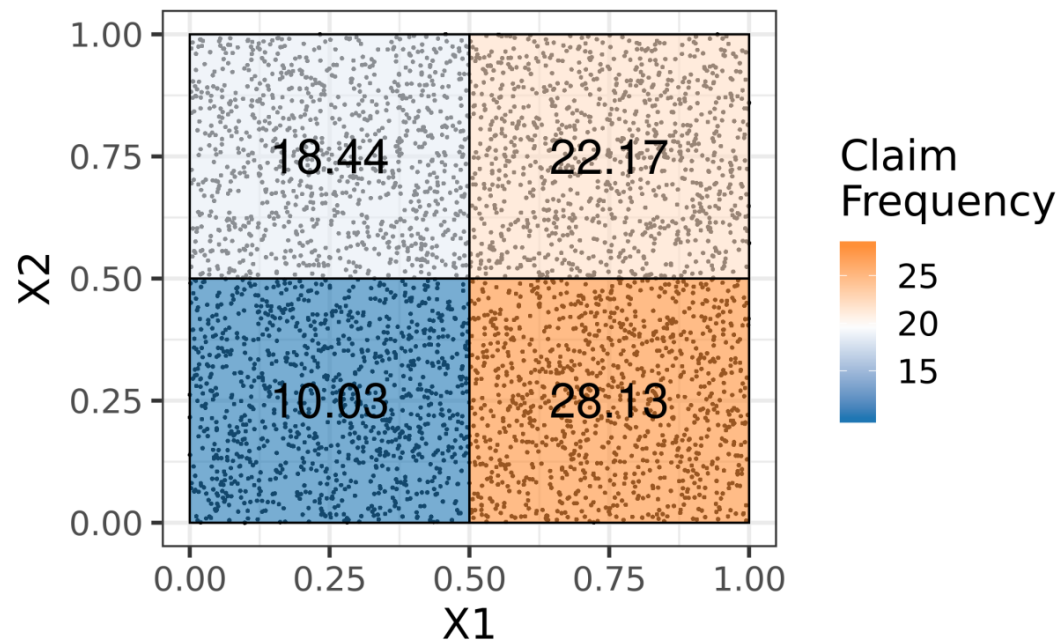
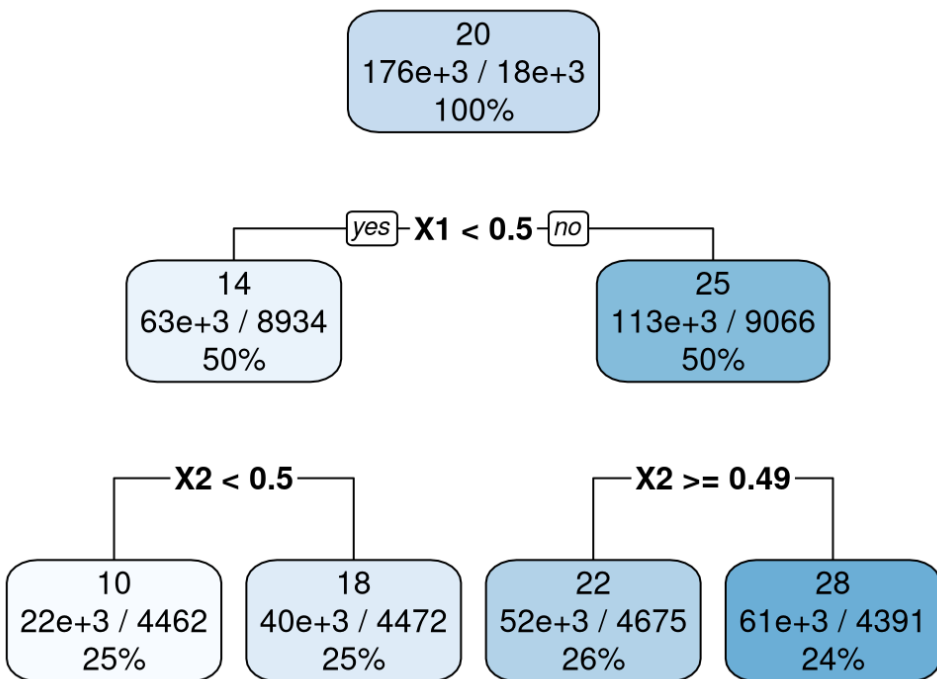
Pour chaque *split* potentiel, la somme des déviiances des deux noeuds/régions filles :

$$D_{\tilde{R}_{l,0}}^*(N, \tilde{\lambda}_{l,0}) + D_{\tilde{R}_{l,1}}^*(N, \tilde{\lambda}_{l,1})$$

Où la déviance s'exprime pour chaque région fille $\tilde{R}_{l,\tau}$, $\tau = \{0, 1\}$

$$D_{\tilde{R}_{l,\tau}}^*(N, \tilde{\lambda}_{l,\tau}) = \sum_{k=1}^n 2N_k \left[\frac{\tilde{\lambda}_{l,\tau}\omega_k}{N_k} - 1 - \log \left(\frac{\tilde{\lambda}_{l,\tau}\omega_k}{N_k} \right) \right] \mathbb{1}_{x_k \in \tilde{R}_{l,\tau}}$$

avec N le nombre total de sinistres observés, N_k le nombre de sinistres pour l'assuré k , $\tilde{\lambda}_{l,\tau}$ la fréquence sinistre prédite au sein d'une région fille $\tilde{R}_{l,\tau}$, $\tau = \{0, 1\}$ et ω_k l'exposition pour l'assuré k .



Comme on peut s'y attendre, l'arbre obtenu est cohérent avec les données initiales.

Un Additive Tree Model (ATM) [Cui et al., 2015] est un ensemble d'arbres de régression ou de classification qui s'écrit sous la forme :

$$F(x_i) = \sum_{t=1}^T w_t f_t(x_i)$$

où $f_t(x_i)$ est l'output de l'arbre t pour l'observation x_i et $w_t \in \mathcal{W}$ est le poids de l'arbre t parmi l'ensemble des T arbres.

Cui, Z., Chen, W., He, Y., and Chen, Y. (2015).

Optimal action extraction for random forests and boosted trees.

In Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining, pages 179–188. ACM

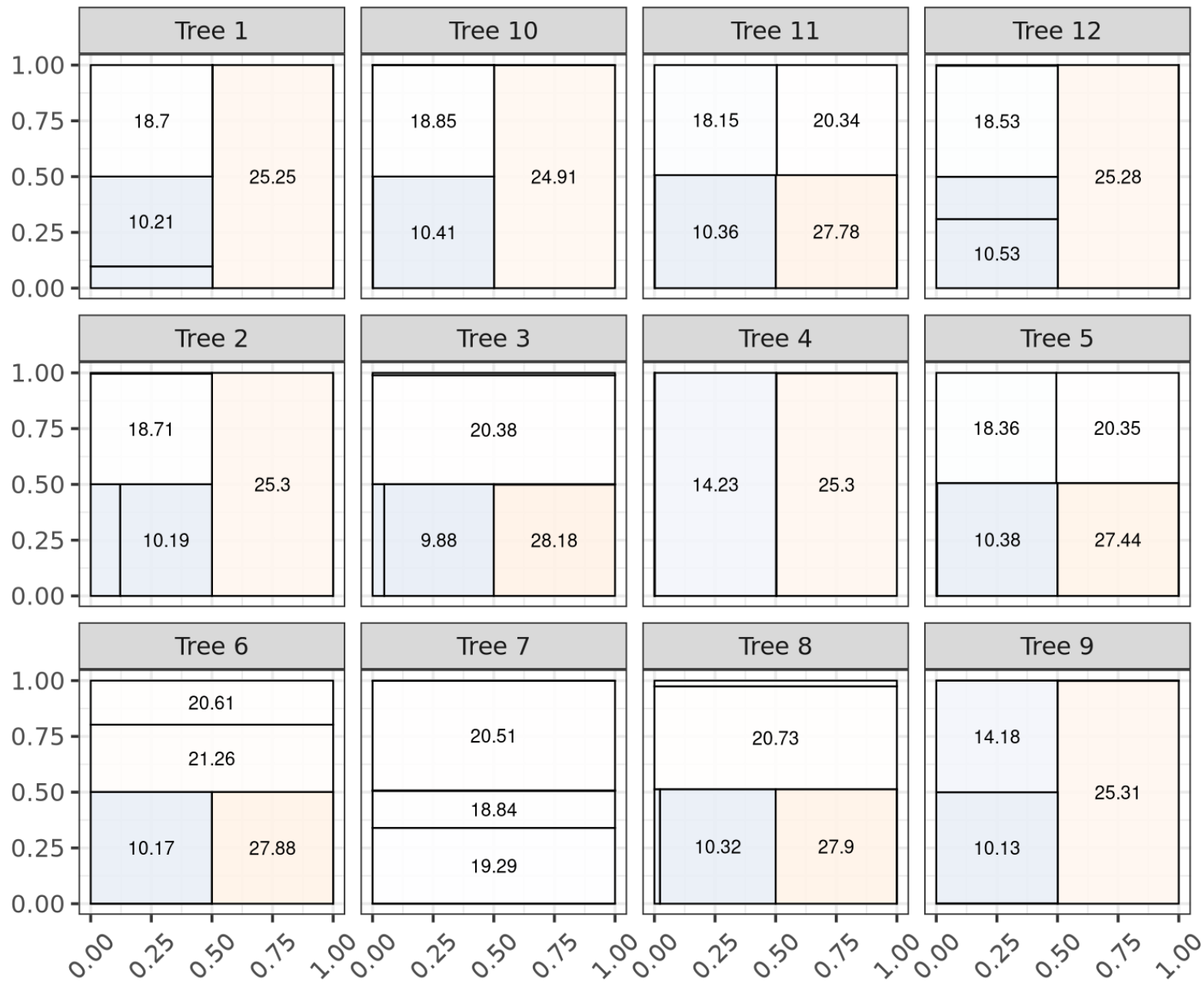
Le **Poisson Random Forest** est une combinaison de T arbres de régression de Poisson. En conséquence, il est possible d'écrire ses prédictions :

$$\forall i \in \llbracket 1, n \rrbracket, F(x_i) = \sum_{t=1}^T \frac{1}{T} f_t(x_i; \tilde{\mathcal{R}}_t, \tilde{\Lambda}_t, I_t)$$

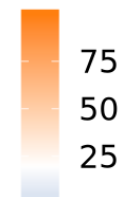
où $\tilde{\mathcal{R}}_t$ est la partition formée par l'arbre t , $\tilde{\Lambda}_t$ est l'ensemble des fréquences sinistres prédites par l'arbre et I_t le nombre de noeuds.

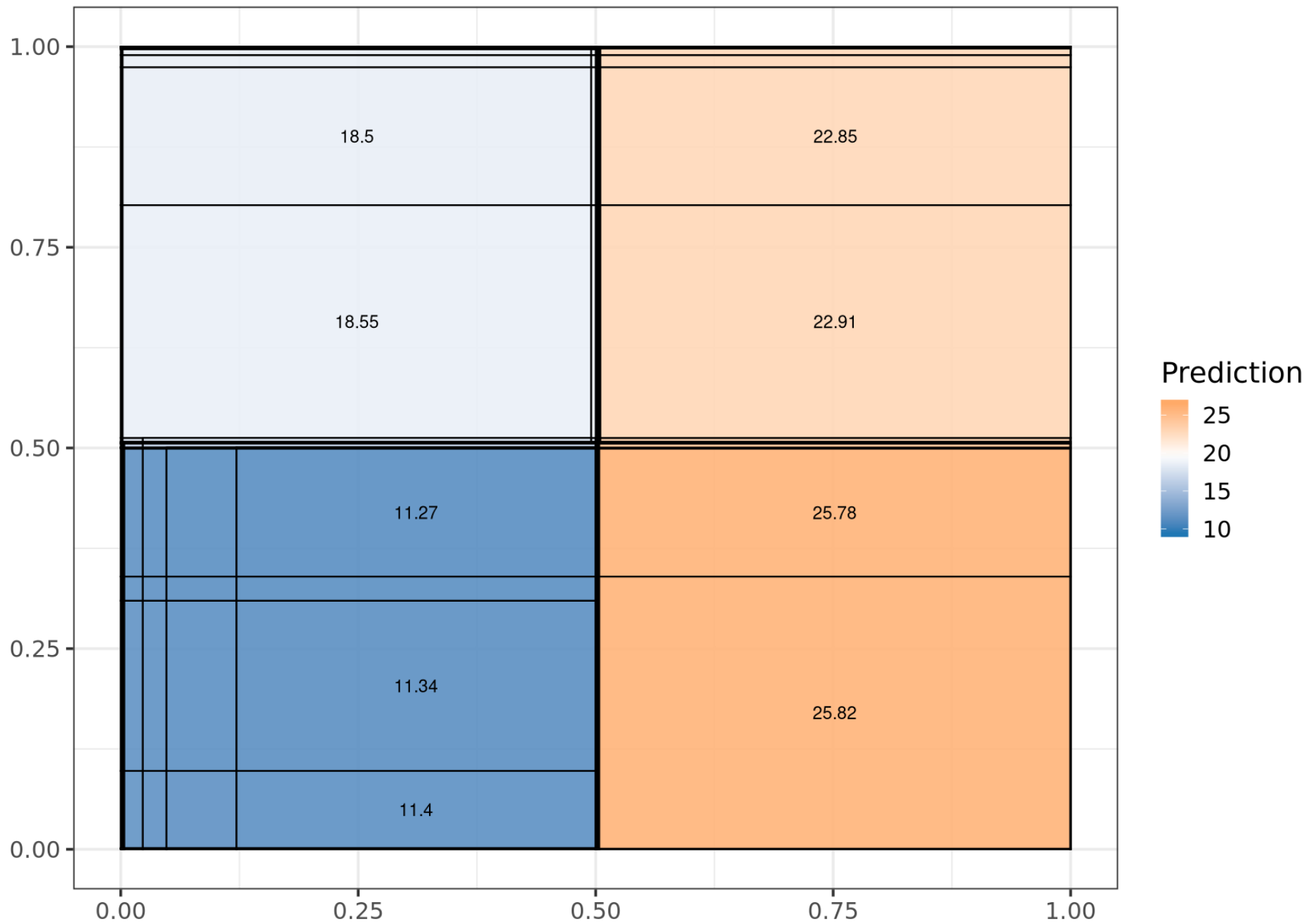
Dès lors, en posant $w_t = 1/T$ il devient aisé de voir que le **Poisson Random Forest** est un ATM.

Par définition, un modèle **Gradient Boosting** appartient à la classe des ATM.



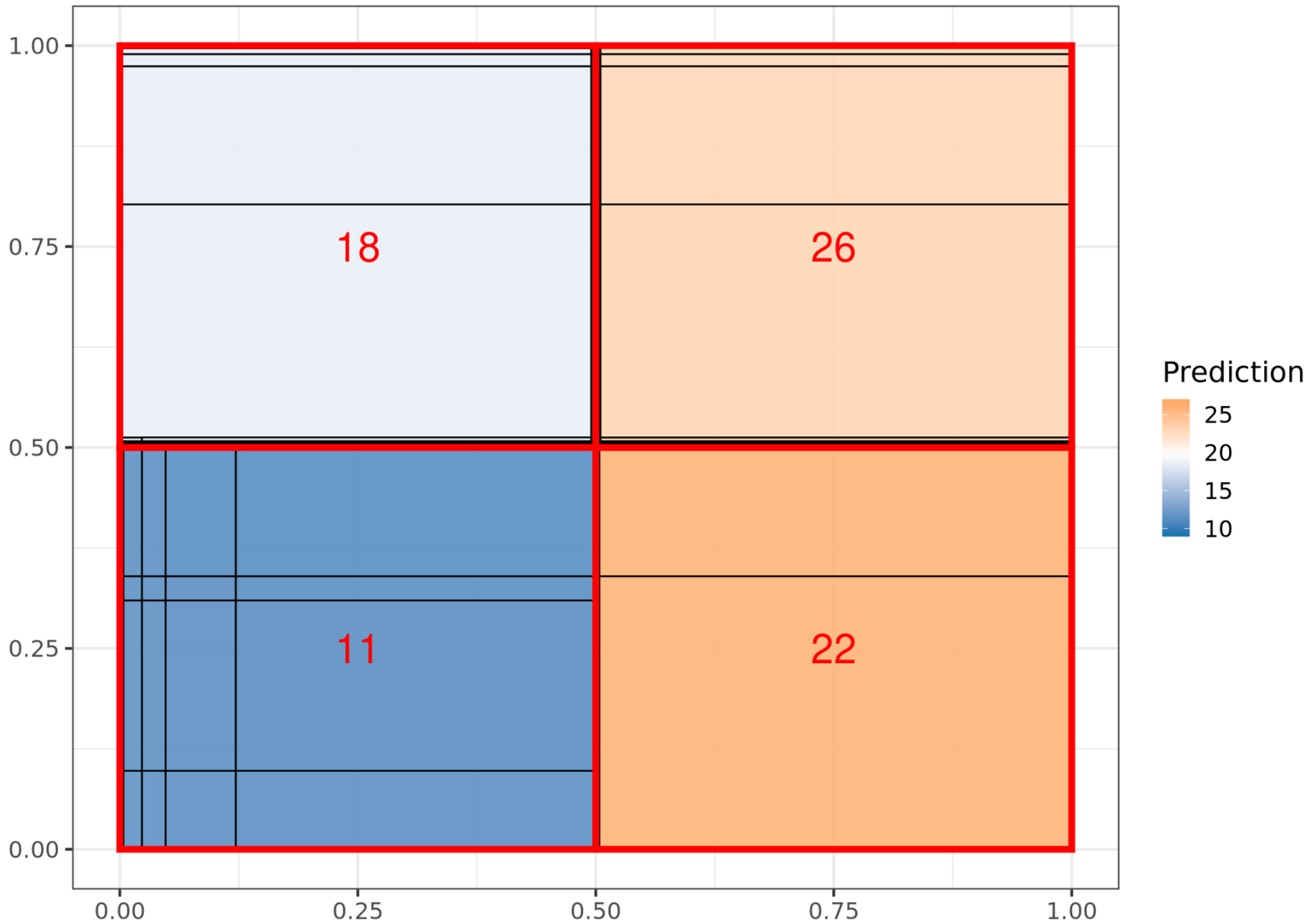
Prediction

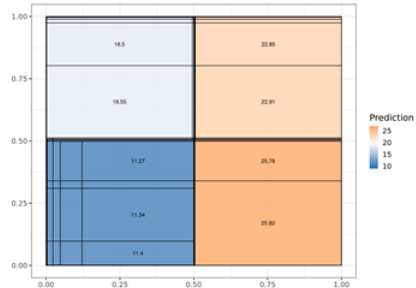
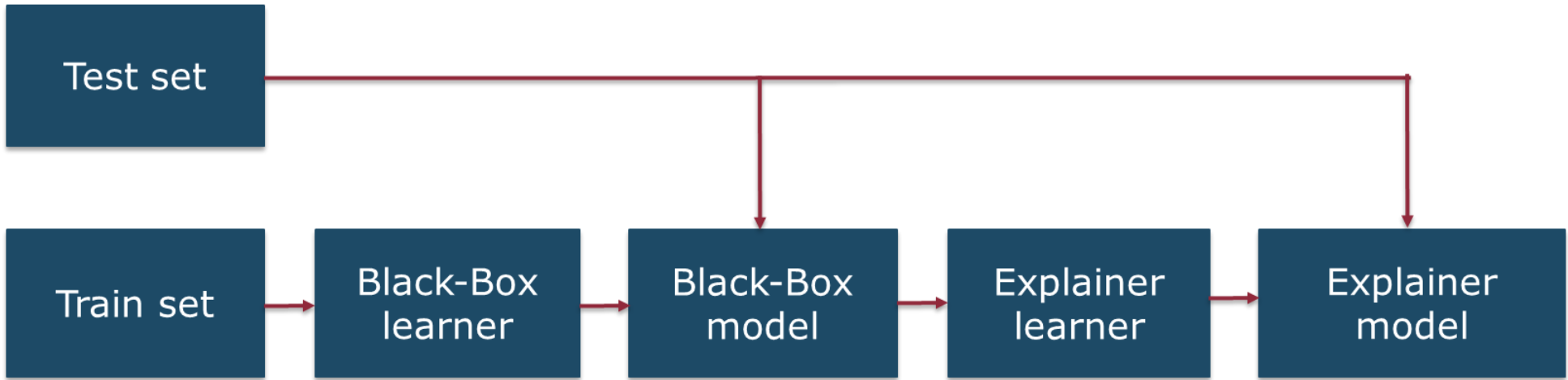




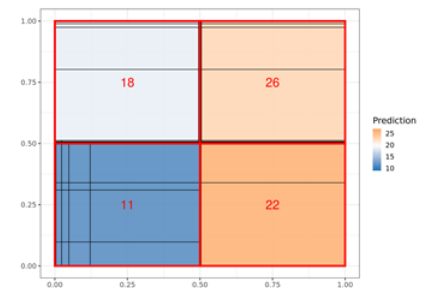
L'exemple précédent montre que sur un ATM bien connu (Random Forest), le problème d'intelligibilité peut être formulé comme un problème de **parcimonie**.

En effet, la prédiction du Random Forest a été rendue plus fine que celle de l'arbre grâce aux multiples régions qui ont été créées. Cependant, cette **fragmentation** rend le modèle peu parcimonieux et donc très compliqué à comprendre.





Prediction



Explanation

Idée générale : agréger des fragments d'espace pour simplifier la partition de la Black-Box et en extraire un modèle surrogat.

Inputs : Les observations du dataset, la prédiction du modèle Black-Box (Random Forest, GBM, XGBoost) ainsi que leurs splits.

Output : Un arbre de décision à K feuilles.

L'ATM génère $G \in \mathbb{N}$ régions notées $\mathcal{R} = \{\mathcal{R}_g\}_{g=1}^G$ dont les valeurs prédites sont $\Lambda = \{\lambda_g\}_{g=1}^G$. Nous voulons trouver $K \in \mathbb{N}$ régions avec $K \ll G$, $\Lambda' = \{\lambda'_k\}_{k=1}^K$ et $\mathcal{R}' = \{\mathcal{R}'_k\}_{k=1}^K$ telles que

$$R^2 = 1 - \frac{\sum_{i=1}^n (f_{bb}(x_i, \omega_i; \Lambda, \mathcal{R}, G) - f_{surrogate}(x_i, \omega_i; \Lambda', \mathcal{R}', K))^2}{\sum_{i=1}^n (f_{bb}(x_i, \omega_i; \Lambda, \mathcal{R}, G) - \bar{y}_{bb})^2} > C$$

où \bar{y}_{bb} est la moyenne des prédictions de l'ATM Black-Box et C est une constante définie par l'utilisateur.

Comme le proposent [Hara and Hayashi, 2016], nous voulons transformer le problème en un problème d'optimisation de vraisemblance. Pour cela, il nous faut une expression generative probabiliste de l'ATM $p(y, \mathbf{x}; \Pi)$. Le problème de maximisation devient,

$$\max_{\Pi} p(y, \mathbf{x}; \Pi)$$

où Π est l'ensemble des paramètres de l'ATM, \mathbf{x} les observations et y les prédictions du modèle Black-Box.

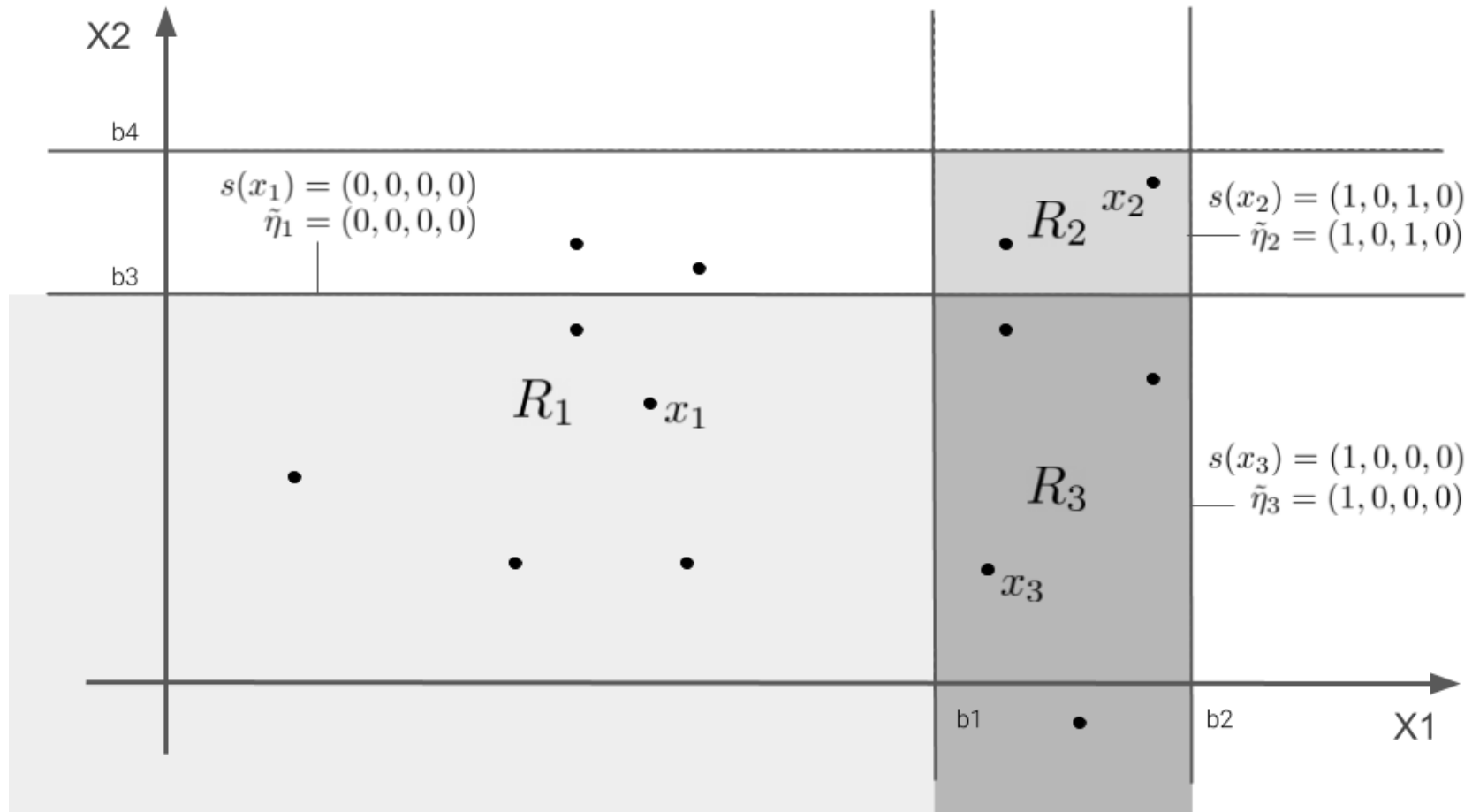
Malheureusement, il n'y a pas de formule fermée pour obtenir les paramètres. Néanmoins, nous pouvons utiliser l'**algorithme Expectation Maximisation** qui est courante pour des problèmes complexes d'optimisation de log-vraisemblance. Un cas d'école est le mélange gaussien.

Hara, S. and Hayashi, K. (2016).

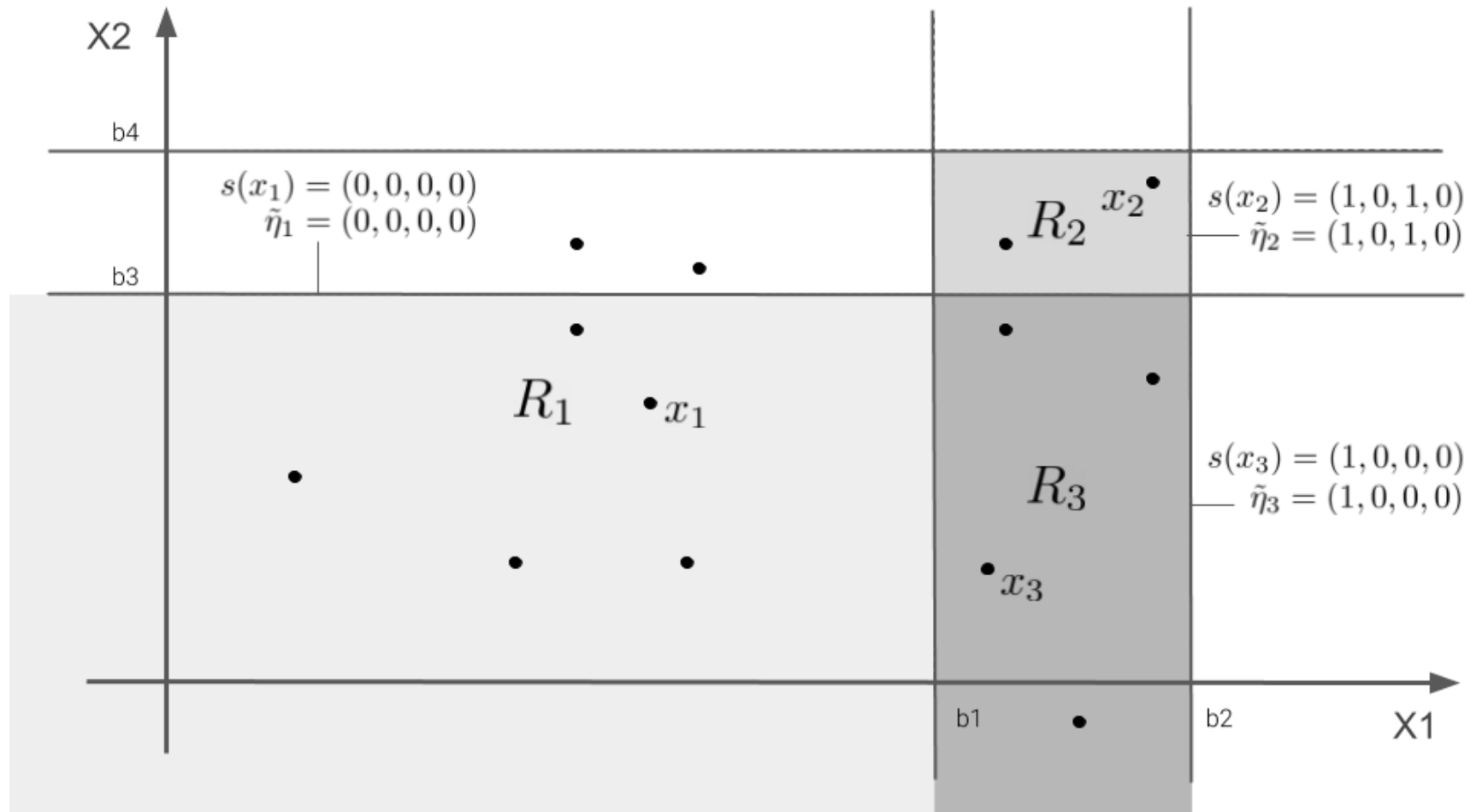
Making Tree Ensembles Interpretable: A Bayesian Model Selection Approach.

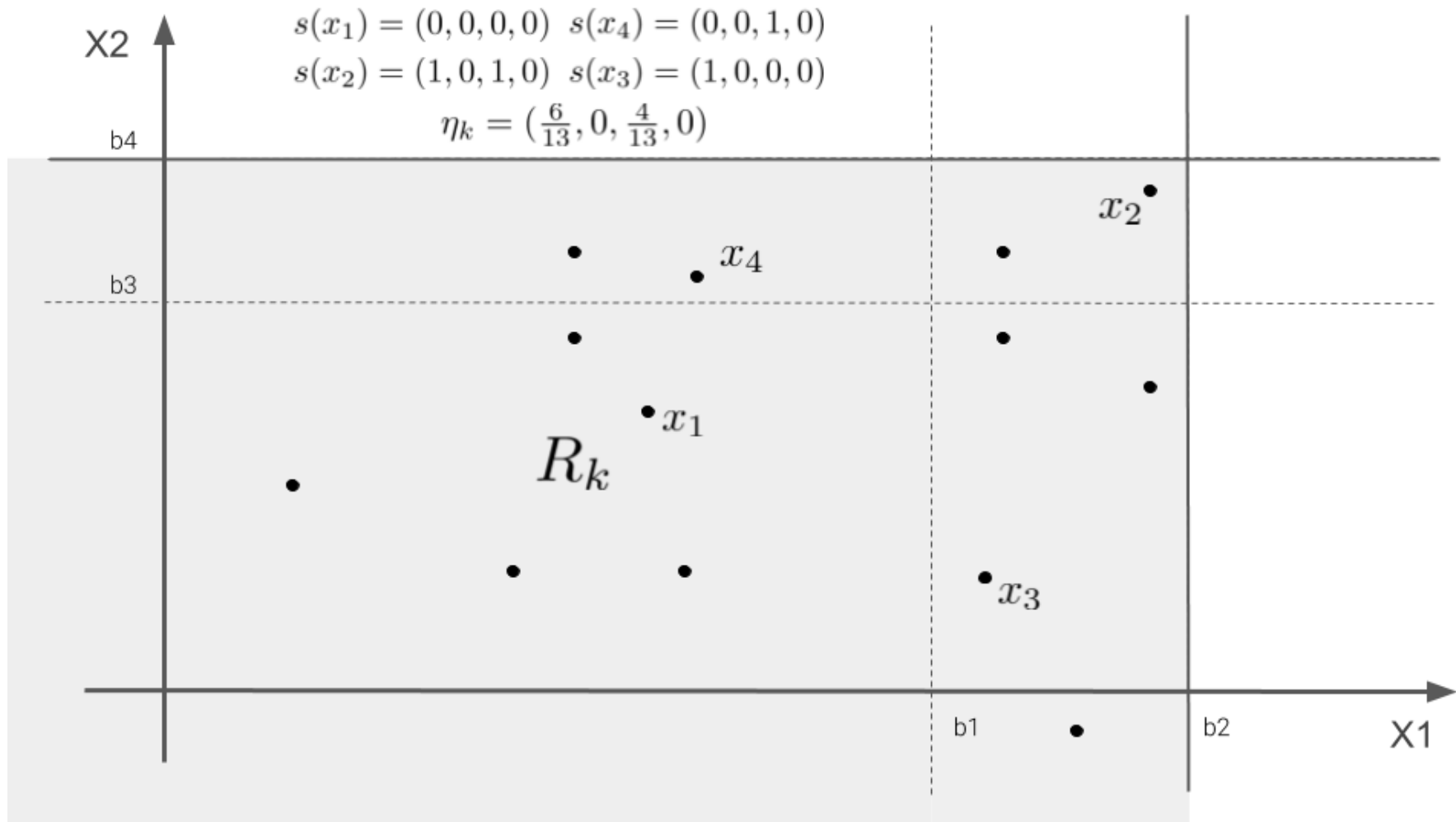
arXiv:1606.09066 [stat]

La forêt générée au total L splits. Nous pouvons exprimer chaque région R_g de la partition comme un vecteur binaire $\tilde{\eta}_g \in \{0, 1\}^L$ dont les composantes $\tilde{\eta}_{gl}$ valent 1 si la région satisfait $\forall \mathbf{x} \in R_g, x_{.,d_l} > b_l$, où d_l est l'index de la variable splittée au noeud l .



Chaque \mathbf{x} appartient à une unique région R_g de la partition formée par l'output de l'ATM. Ce qui signifie que la représentation binaire de x notée $\mathbf{s}(\mathbf{x})$ est égale à la représentation binaire $\tilde{\boldsymbol{\eta}}_g$ de la région R_g si et seulement si $\mathbf{x} \in R_g$.





Considérons maintenant une potentielle région R_k créée à partir de la concaténation des régions R_1 , R_2 et R_3 . Désormais, certains $\mathbf{x} \in R_k$ vérifient $x_{.,d_2} > b_2$ alors que d'autres non. Mais tous les $\mathbf{x} \in R_k$ vérifient $x_{.,d_1} > b_1$ et $x_{.,d_3} \leq b_3$ par exemple. Il est maintenant possible d'exprimer la probabilité d'observer une représentation binaire pour une région k par :

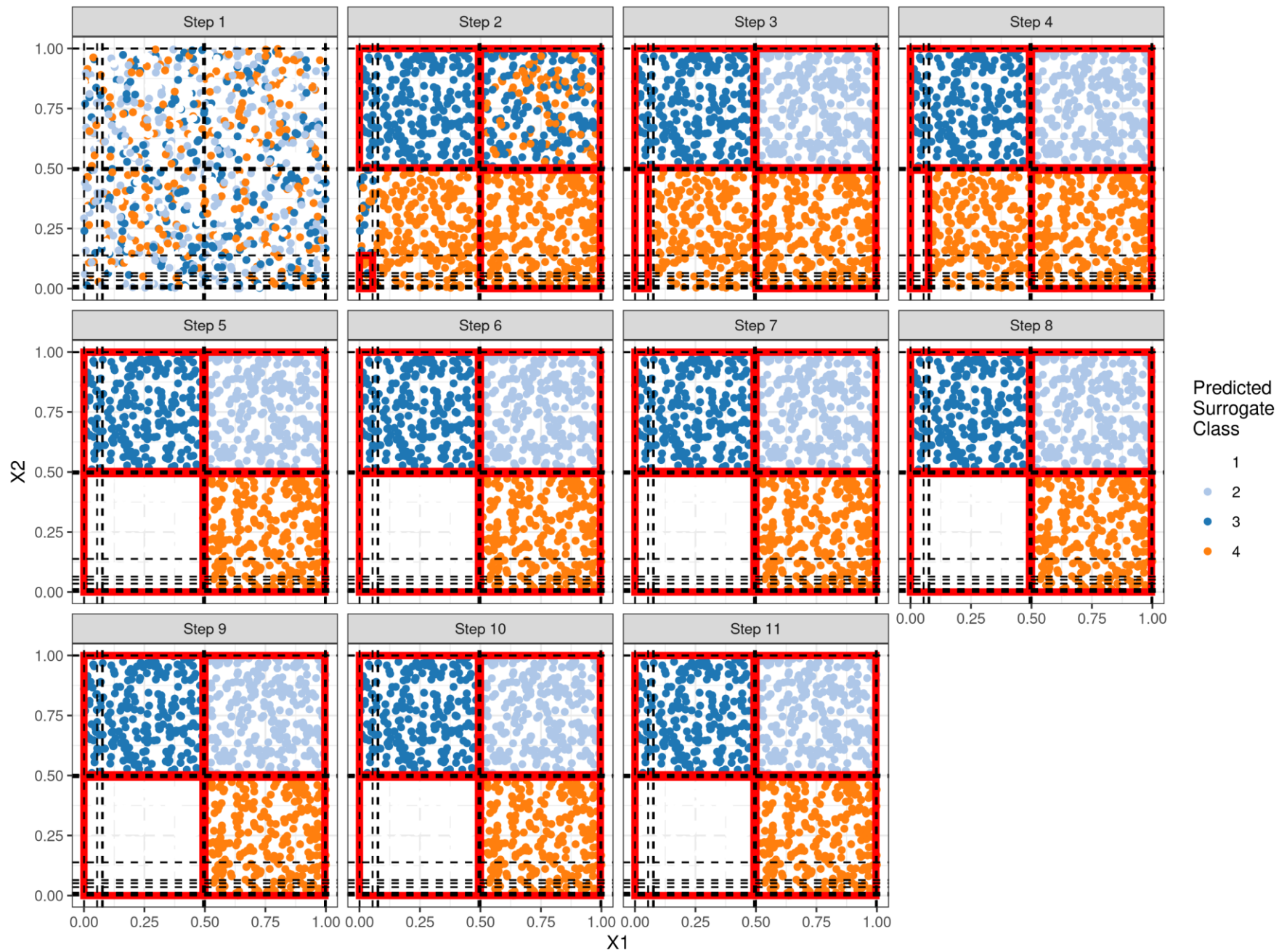
$$p(\mathbf{s}(\mathbf{x}); \boldsymbol{\eta}, k) = \prod_{l=1}^L \eta_{kl}^{s_l} (1 - \eta_{kl})^{(1-s_l)}$$

Pour rentrer dans le cadre de l'optimisation par l'algorithme EM, nous avons besoin d'une variable latente \mathbf{u} . Ce vecteur $\mathbf{u} \in \{0, 1\}^K$ décrit à quelle région \mathbf{x} appartient ($u_k = \mathbb{1}_{\mathbf{x} \in R_k}$). Nous pouvons ainsi définir la probabilité d'observer le couple $(y, \mathbf{s}(\mathbf{x}))$ sachant \mathbf{u} :

$$p(y, \mathbf{s}(\mathbf{x}) | \mathbf{u}) = \prod_{k=1}^K (p(y; \boldsymbol{\lambda}, \boldsymbol{\omega}, k) \times p(\mathbf{s}(\mathbf{x}); \boldsymbol{\eta}, k))^{u_k}$$

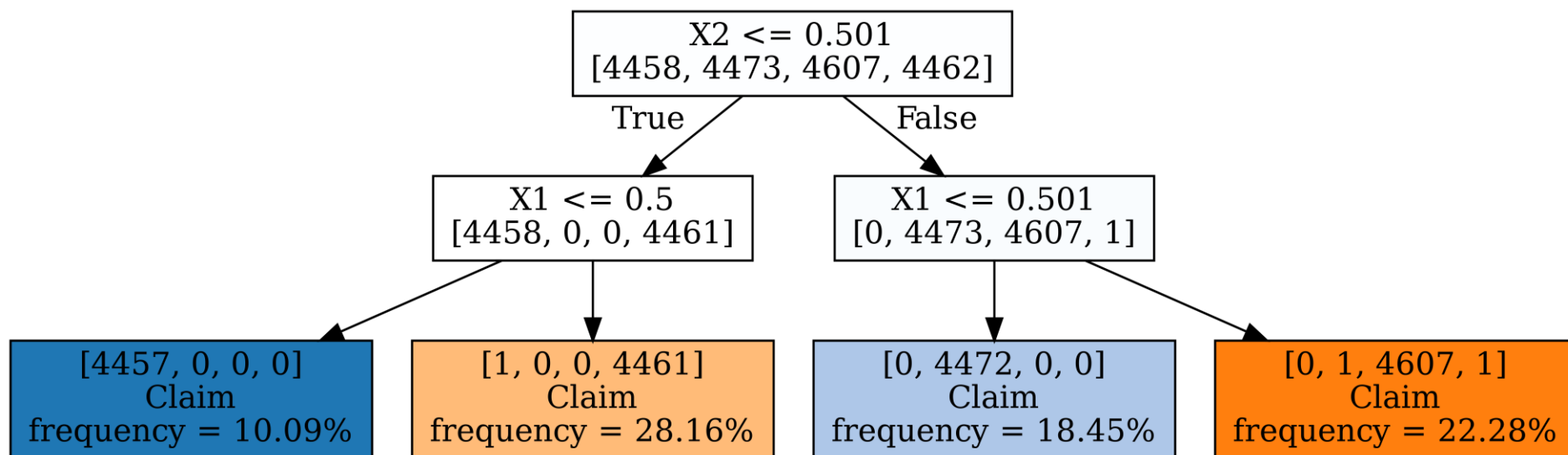
Nous choisissons une distribution de Poisson $p(y; \boldsymbol{\lambda}, \boldsymbol{\omega}, k) \sim \mathcal{P}(\lambda_k \omega_k)$ pour être cohérent avec notre ATM Black-Box.

$$p(y, \mathbf{s}(\mathbf{x}), \mathbf{u}; \boldsymbol{\lambda}, \boldsymbol{\omega}, \boldsymbol{\eta}, \boldsymbol{\alpha}, K) = \prod_{k=1}^K (p(y; \boldsymbol{\lambda}, \boldsymbol{\omega}, k) \times p(\mathbf{s}(\mathbf{x}); \boldsymbol{\eta}, k) \times \alpha_k)^{u_k}$$



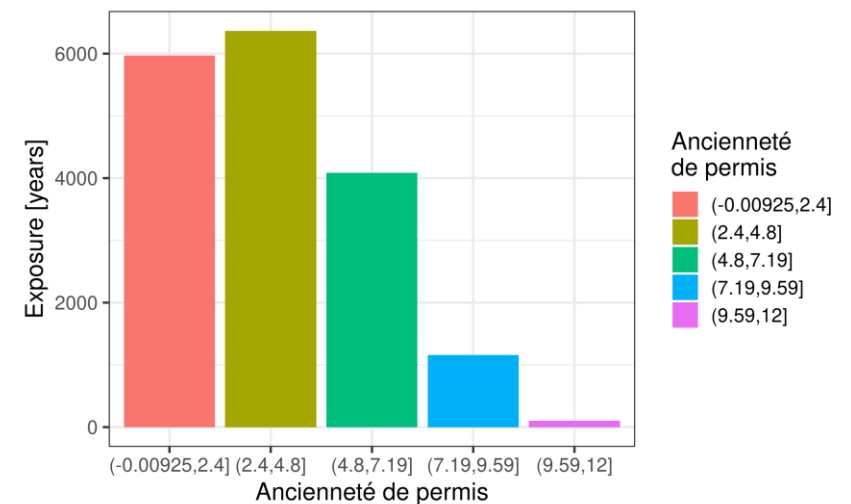
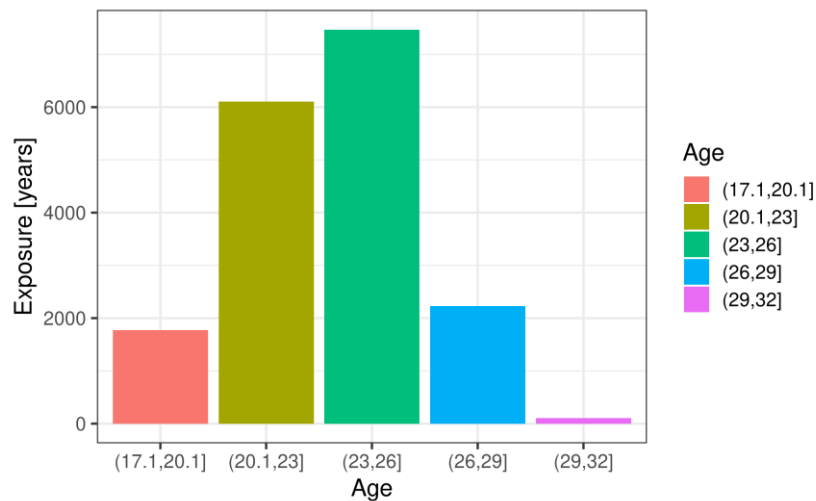
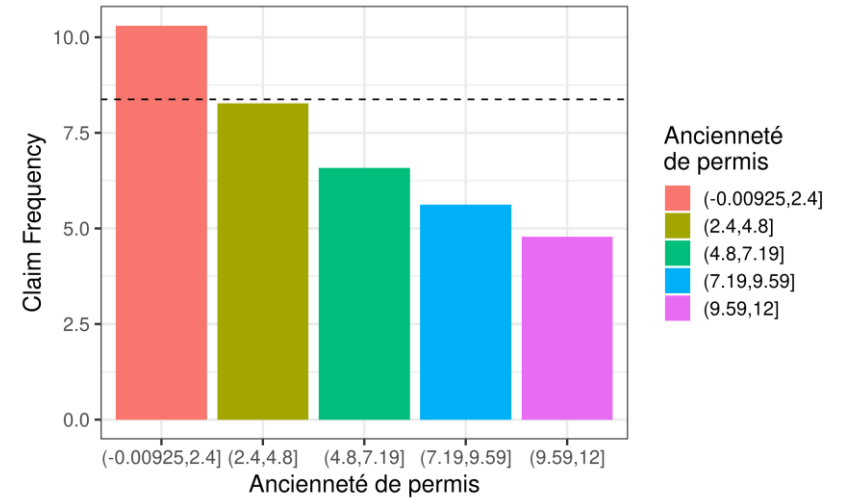
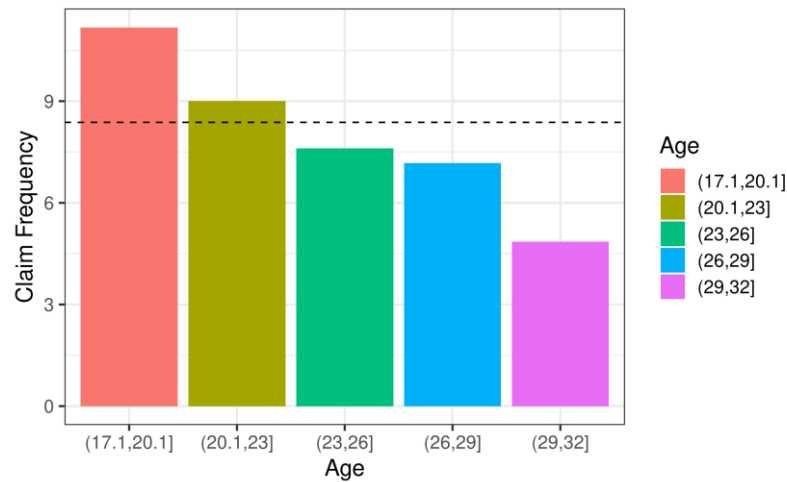
L'algorithme converge vers les paramètres optimaux : η^* , λ^* , ω^* , α^* . Ces derniers forment une "soft partition" de l'espace. En conséquence, une observation de l'espace peut appartenir à plusieurs des K classes du modèle surrogate.

Pour remédier à ce problème nous proposons d'utiliser une version modifiée de l'arbre de classification pour reformer une partition de l'espace. Cet arbre apprend sur les observations et a pour cible la classe affectée par le surrogate. Sa recherche est restreinte aux splits initiaux pour assurer la cohérence de l'explication.

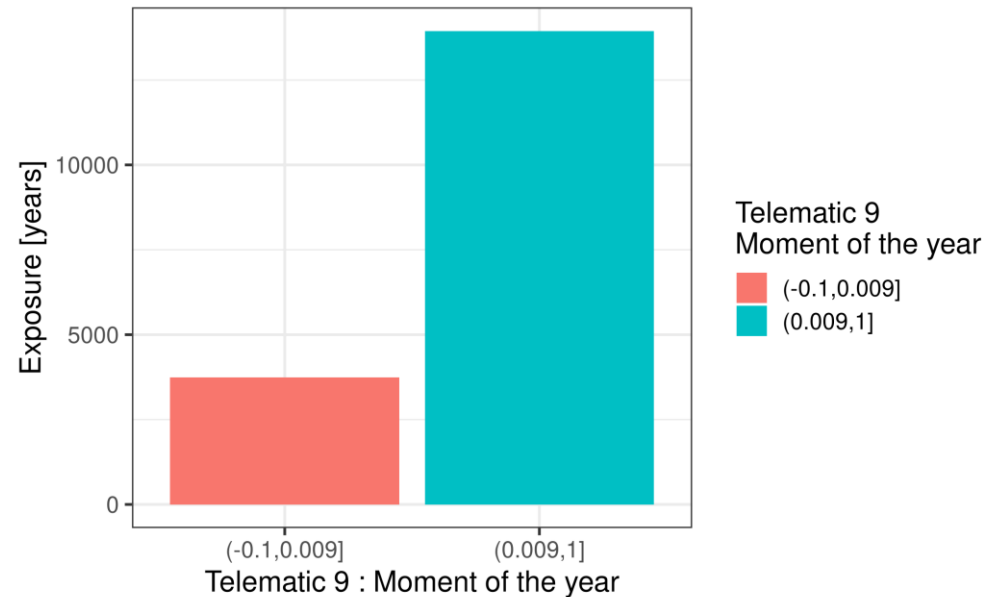
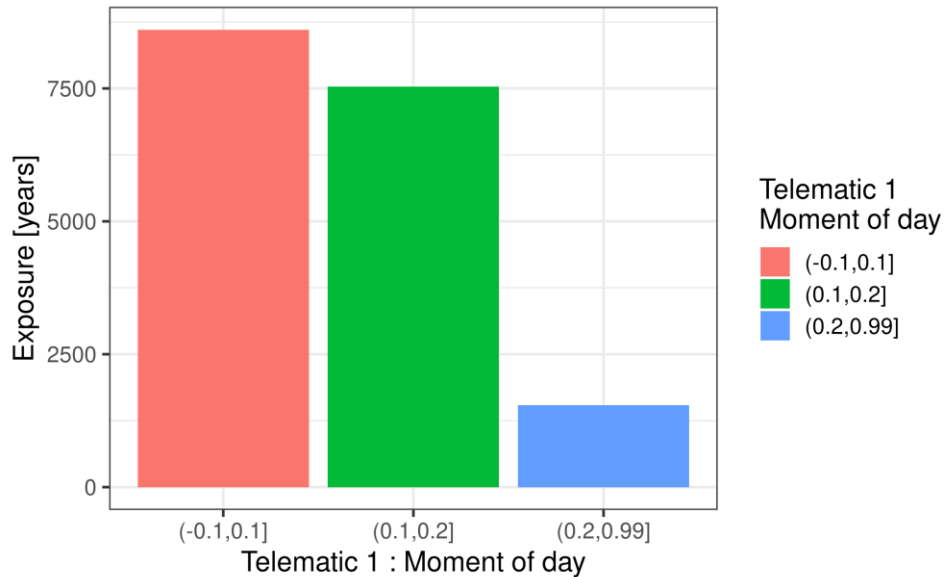
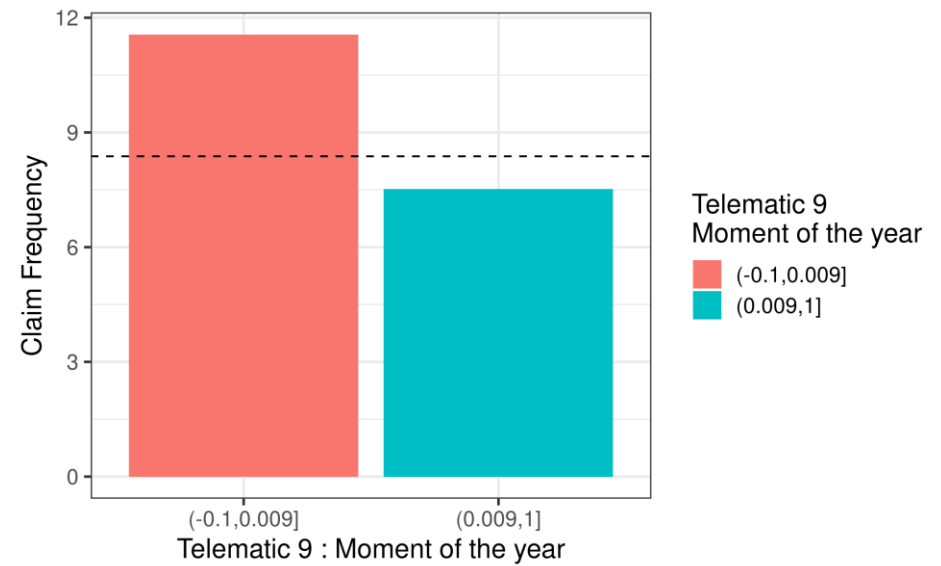
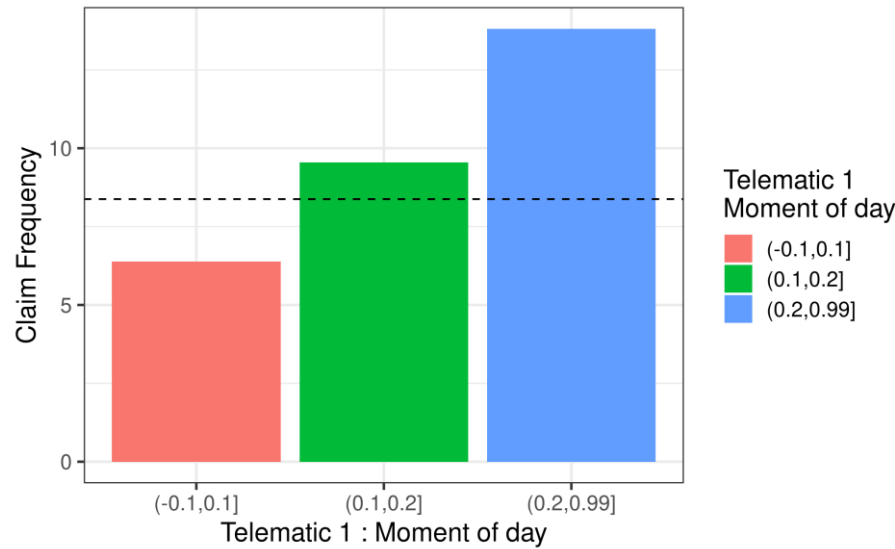


- Le jeu de données contient 33 259 observations
- Le portefeuille est majoritairement composé de jeunes conducteurs.
- Il y a une fréquence sinistre moyenne de 8.38% sur le portefeuille global si l'on considère le temps comme exposition.

Dans le jeu de données, nous avons des variables classiques de la police telles que l'âge, l'ancienneté de permis et le sexe par exemple. Ces variables semblent être discriminantes pour la prédiction de la fréquence sinistre.

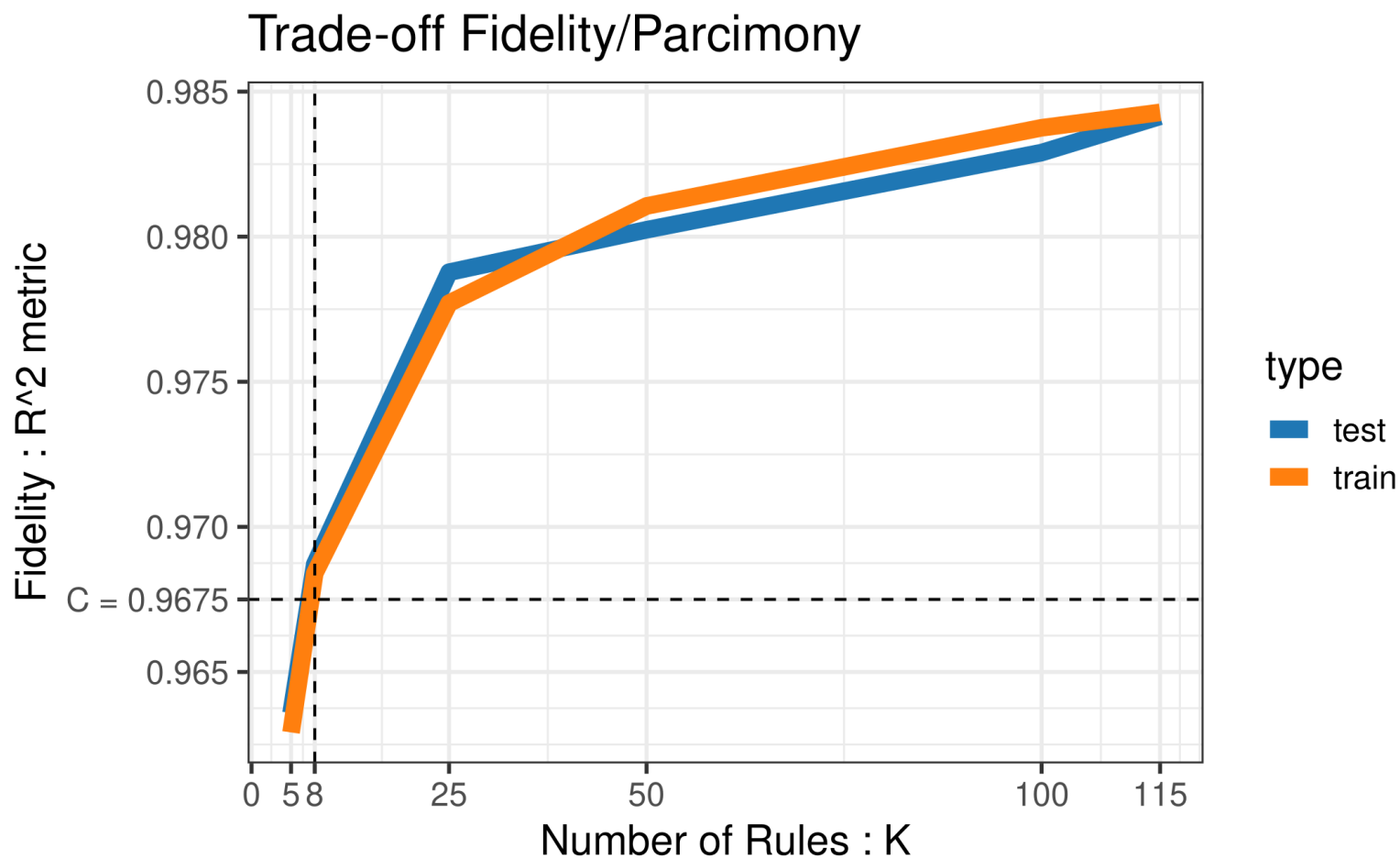


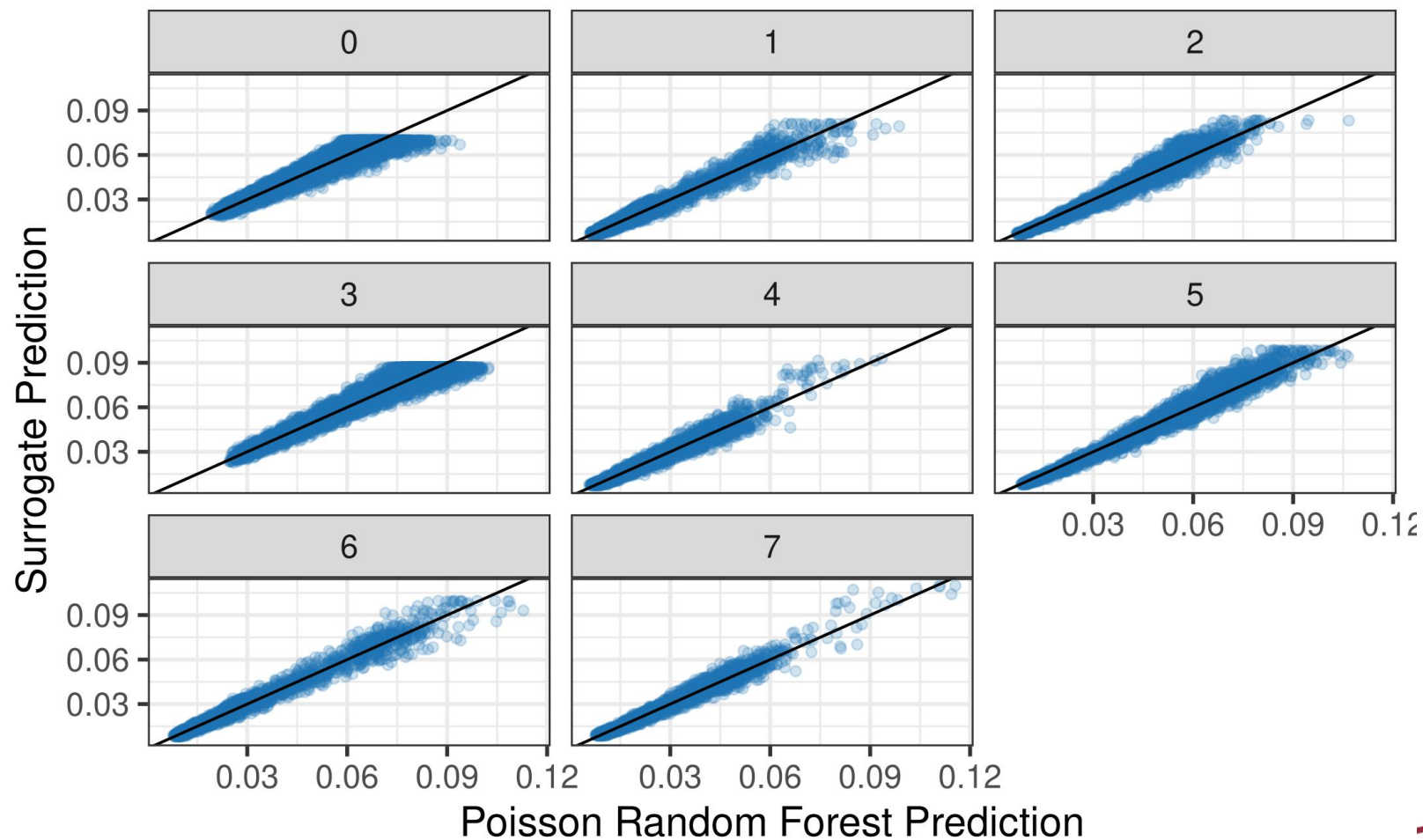
Le jeu de données contient aussi des variables originales collectées par des boîtiers embarqués dans les voitures pendant les trajets.

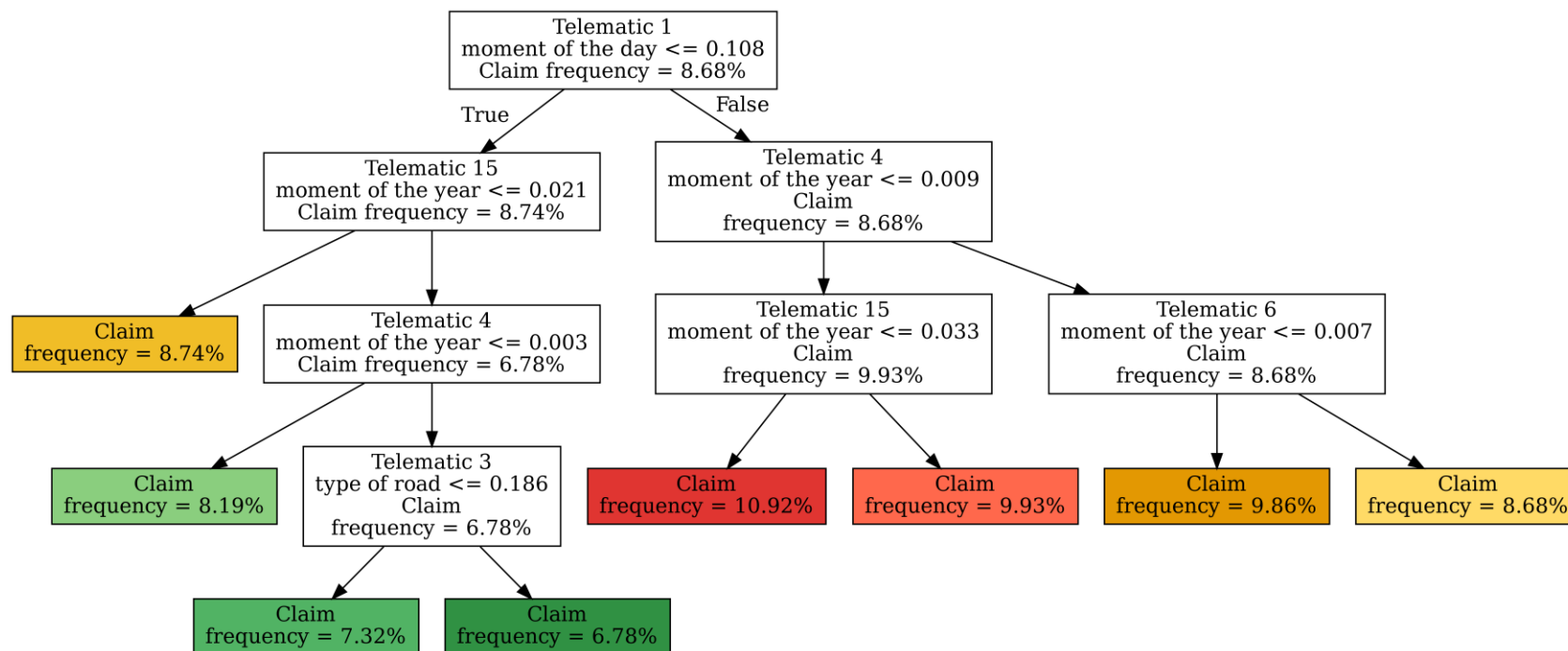


Pour modéliser la fréquence sinistre avec des ATM poissonniens, nous avons choisi le Poisson Random Forest implémenté dans le package R **RfCountData** ainsi que le Gradient Boosting implémenté dans le package **XGBoost**

Tous les deux permettent de gérer la variable d'offset aussi appelée **exposure** chez les actuaires.







- DefragTrees permet de donner un bon surrogate en termes de fidélité
- Le framework est assez général pour s'adapter à la majorité des distributions utilisées en actuariat
- L'algorithme ne fournit pas une suite d'arbres emboîtés

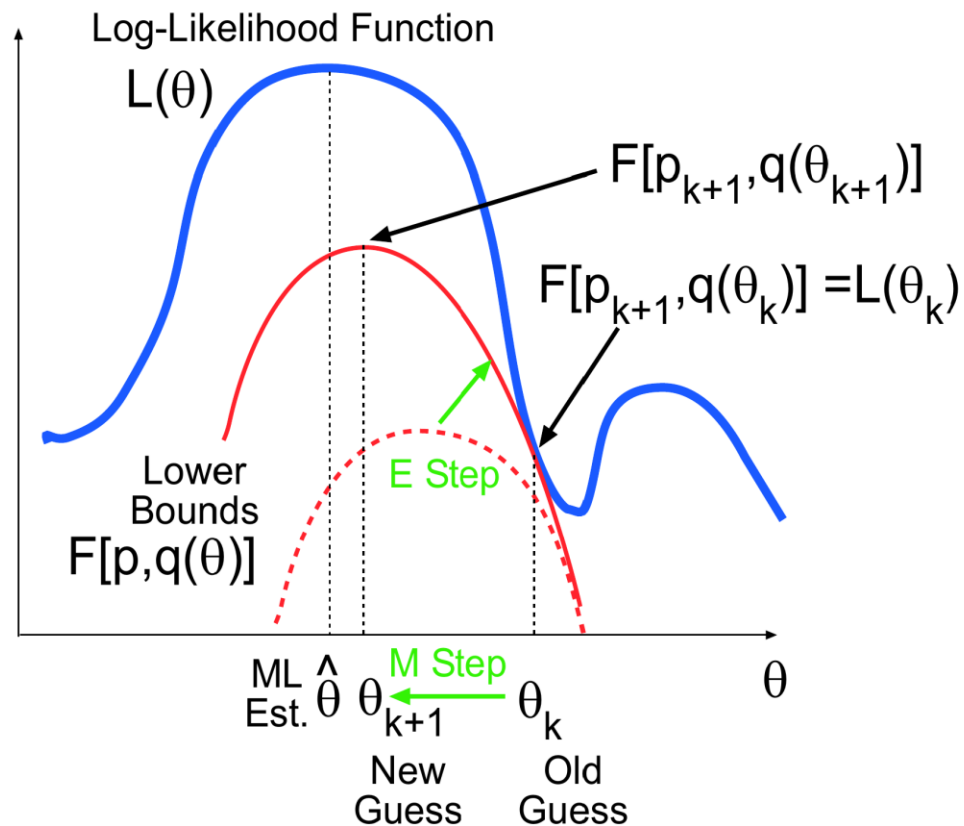


Figure 1: The EM algorithm alternates between finding a greatest lower bound (“E step”), and maximizing this bound (“M step”).

Grâce à l'expression précédente, il est possible d'obtenir une expression pour les étapes **Expectation** (E) et **Maximization** (M) de l'algorithme EM.

E step :

Pour les paramètres courants $\Pi = \{\lambda, \omega, \eta, \alpha, k\}$

$$q(u_k^{(n)} = 1) = \frac{p(y^{(n)} | \lambda, \omega, k) p(s^{(n)} | \eta, k) p(k; \alpha)}{\sum q(u_k^{(n)}) p(y^{(n)}, s^{(n)}, u_k^{(n)} | \Pi, K)}$$

M step :

En posant $q(u_k^{(n)} = 1) = \beta_k^{(n)}$

$$\eta_{k,l} = \frac{\sum_{n=1}^N \beta_k^{(n)} s_l^{(n)}}{\sum_{n=1}^N \beta_k^{(n)}}, \quad \lambda_k = \frac{\sum_{n=1}^N \beta_k^{(n)} y^{(n)} \omega_k^{(n)}}{\sum_{n=1}^N \beta_k^{(n)} \omega_k^{(n)}}, \quad \alpha_k = \frac{1}{N} \sum_{n=1}^N \beta_k^{(n)}$$

Pour implémenter ceci, nous avons modifié le package python **DefragTrees**.