



# L'application du Machine Learning dans la tarification IARD

Présentation par William James Ph.D. AIAA

Journées IARD de l'Institut des Actuaire

30 mars 2018



# L'application du Machine Learning dans la tarification IARD

## Sommaire

- **Les dimensions analytiques**
  - La dimension du problème
  - La dimension méthodologique
- **Techniques de modélisation**
  - Les plus populaires
  - Exemple: GBMs
  - Calibration
- **Comment utiliser le machine learning ?**
  - La valeur ajoutée
  - L'interprétation
  - L'utilisation du modèle
- **Faut-il passer du temps sur la méthode ou le problème?**



## Un exemple simple

- Comment prédire la probabilité de rester réveillé pendant cette présentation?
- Tout ce qu'on sait, c'est...
  - Ce que vous avez mangé pour le petit déj
  - Le nombre d'heures de sommeil hier soir
  - Si vous venez de prendre un café
  - Si vous avez votre téléphone portable allumé...
- Construisons un modèle !

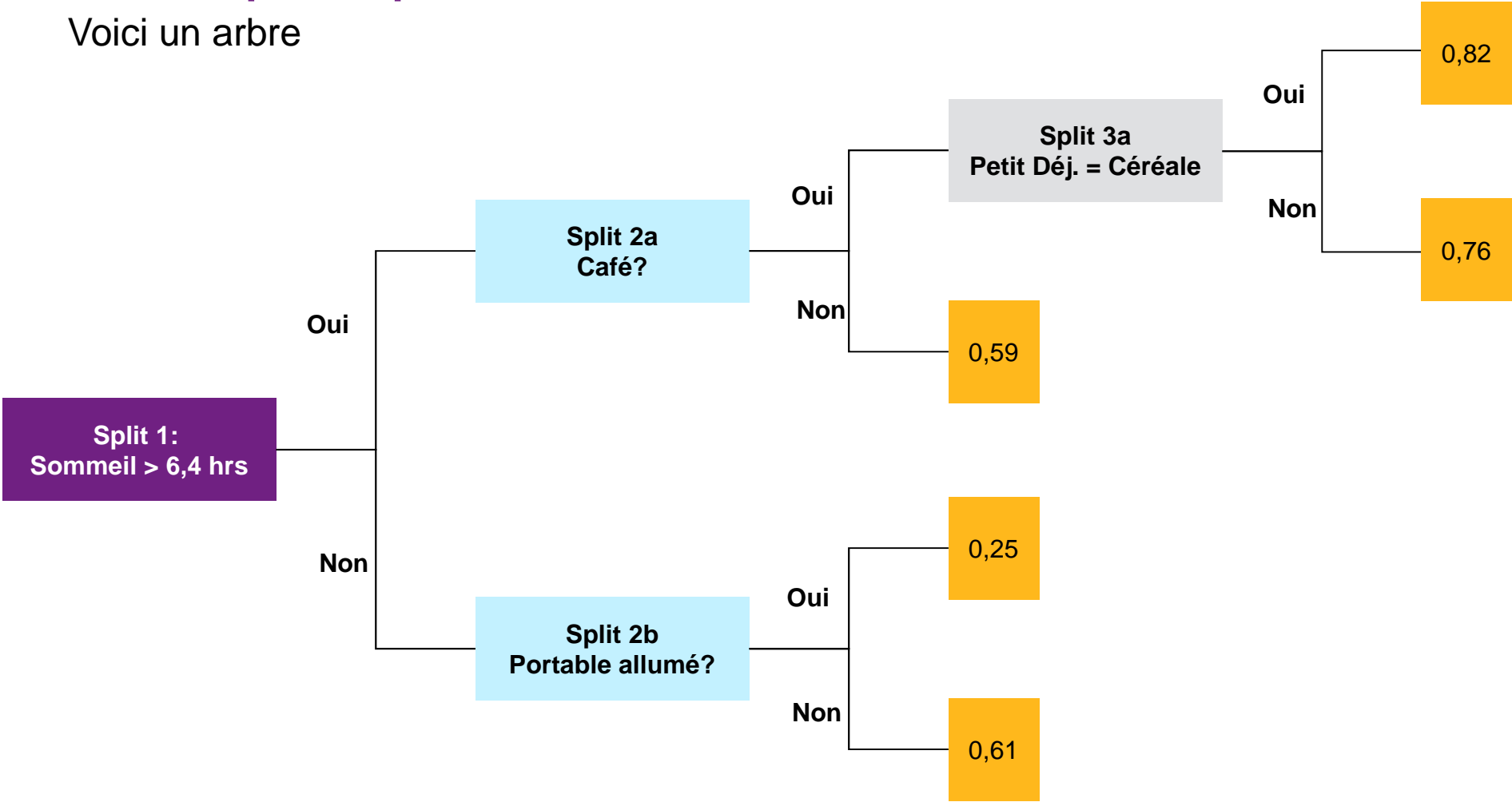
## Un exemple simple

Voici un GLM

Petit déj	Valeur	Sommeil (heures)	Valeur	Café	Valeur	Portable allumé	Value
Aucun	-0.6	0	-2	Y	0	Y	0
Tartine	0	1	-1.8	N	-1	N	1
Fruit	0	2	-1.5				
Pâtisseries	0.2	3	-1.2				
		4	-0.8				
		5	-0.4				
		6	0				
		7	0.4				
		8	0.8				
		9	1.3			Base	0

# Un exemple simple


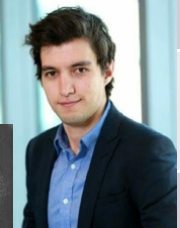





Voici un arbre



Testons avec des prédictions....

# Un exemple simple

Quelques prédictions

	Petit déj.	Sommeil	Café	Portable	GLM	Arbre
	Pâtisseries	7	Y	Y	0,65	0,76
	Céréale	5	Y	Y	0,40	0,25
	Pâtisseries	6	N	Y	0,31	0,25
	Céréale	8	N	N	0,69	0,59
	Fruit	4	N	Y	0,14	0,25
	Fruit	6	Y	N	0,73	0,61
	Céréale	8	Y	Y	0,69	0,82

## Un exemple simple

Alors... quelle méthode ?

- Quelle méthode est la meilleure ?
- Comment évalue-t-on cela ?
- Où est la valeur ?

## Un exemple simple

Mais...

- Est-ce que c'était la bonne question ?
- Les variables étaient-elles les meilleures ?
  - Étaient-elles définies correctement ?
- Avec tout l'intérêt porté aux méthodologies analytiques, souvent on oublie que la question est plus large que la sélection simple de méthodologie...



# L'application du Machine Learning dans la tarification IARD

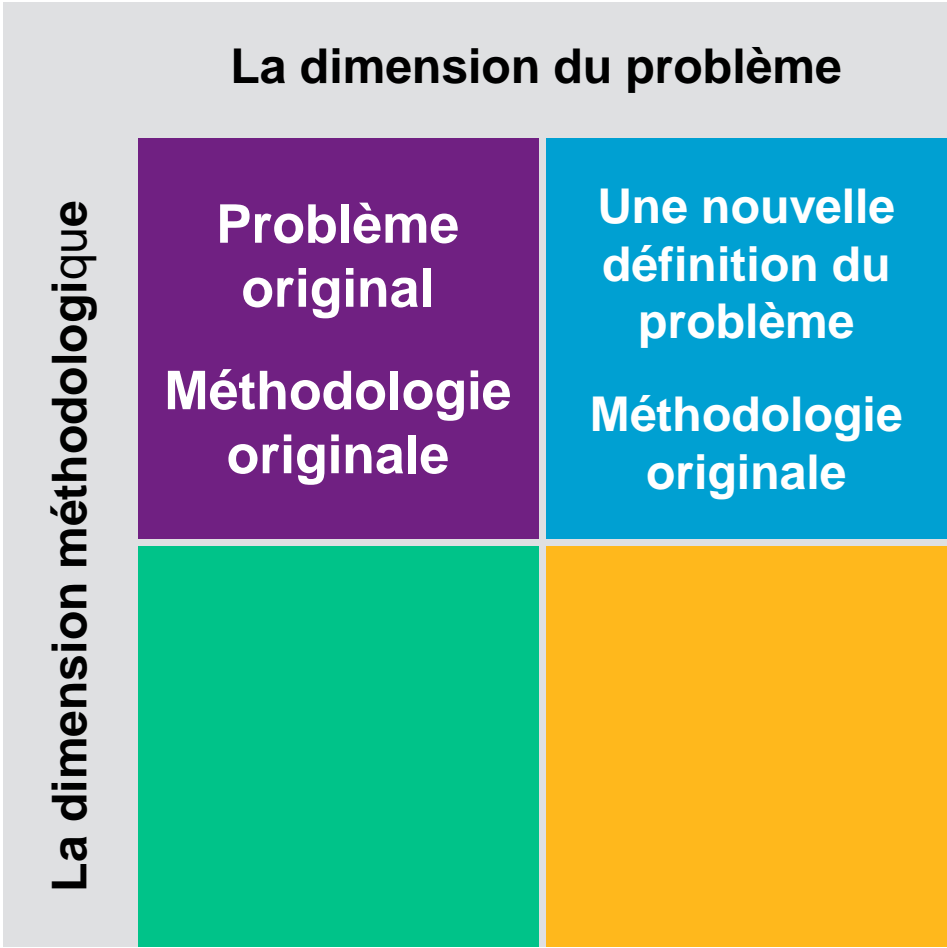
## Sommaire

- **Les dimensions analytiques**
  - La dimension du problème
  - La dimension méthodologique
- **Techniques de modélisation**
  - Les plus populaires
  - Exemple: GBMs
  - Calibration
- **Comment utiliser le machine learning ?**
  - La valeur ajoutée
  - L'interprétation
  - L'utilisation du modèle
- **Faut-il passer du temps sur la méthode ou le problème?**



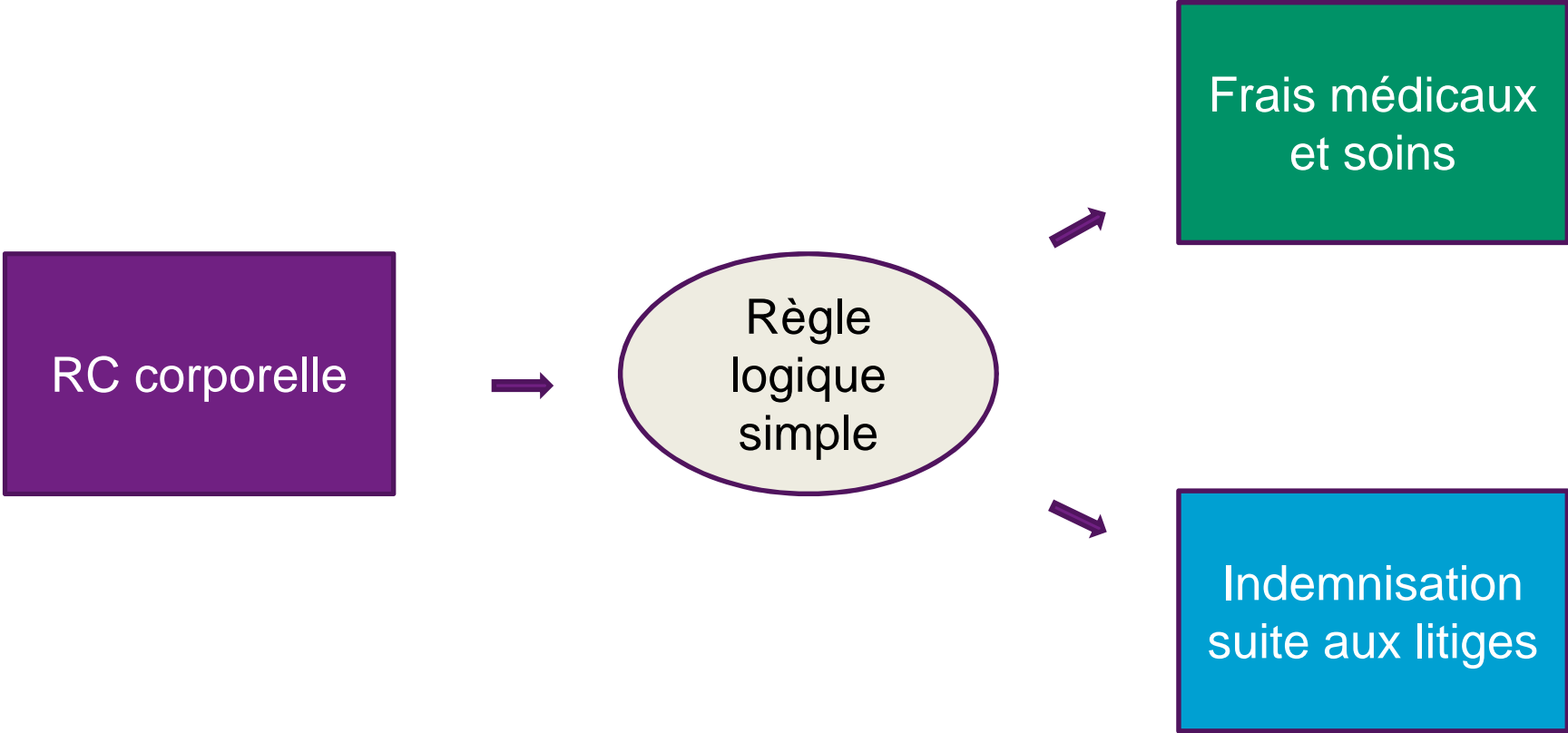
# Les dimensions analytiques

## La dimension du problème



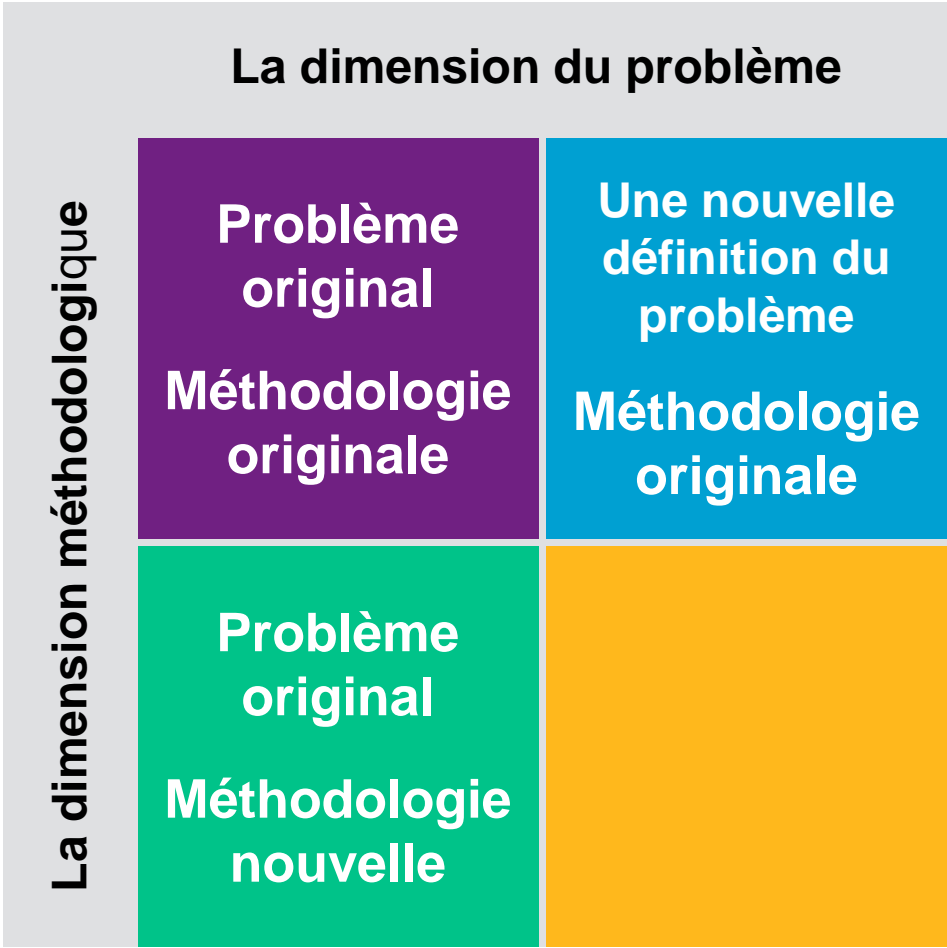
# Les dimensions analytiques

La dimension du problème



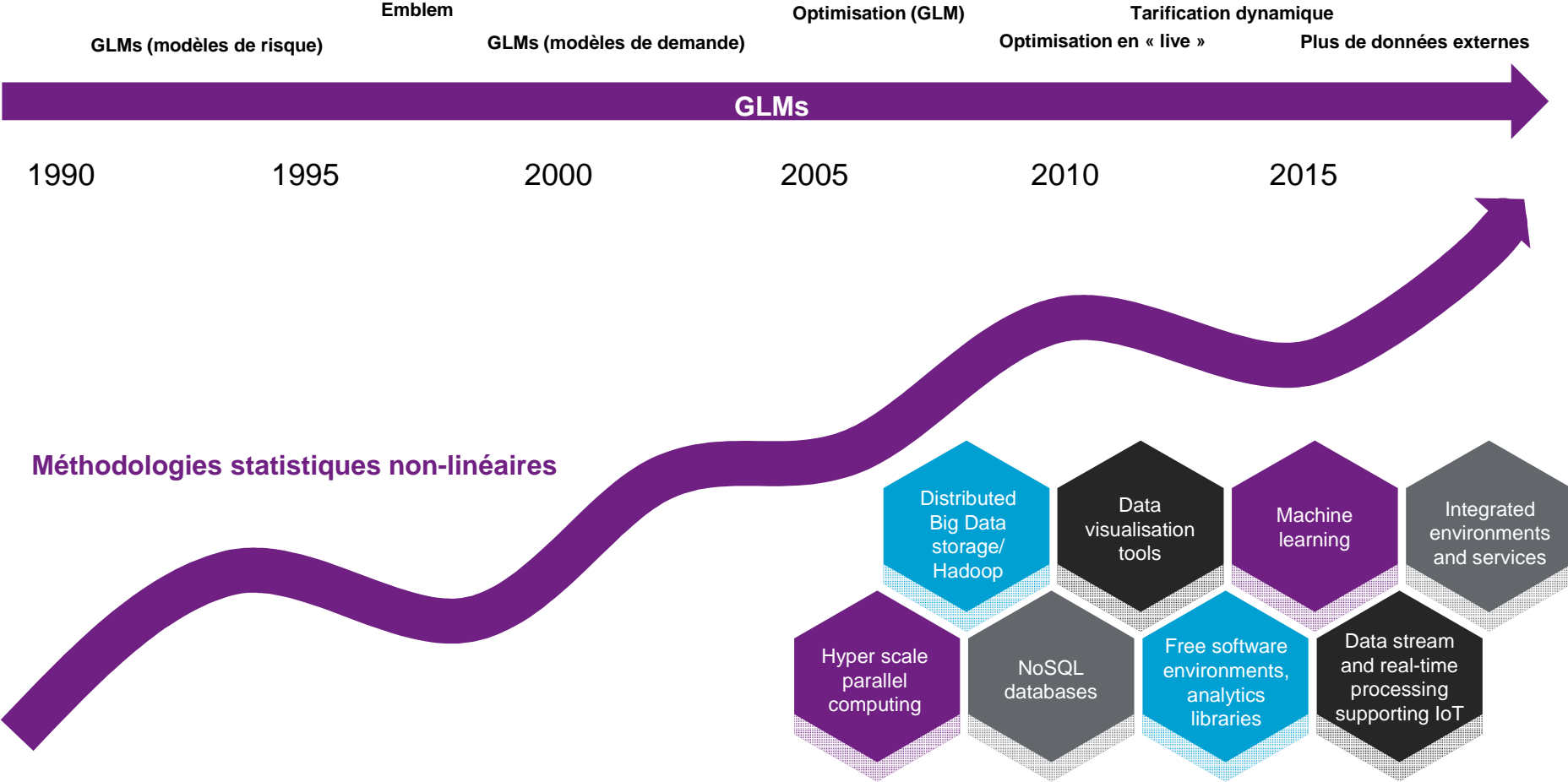
# Les dimensions analytiques

## La dimension méthodologique



# Techniques de modélisation

## L'histoire méthodologique





# Kaggle

[Host](#)
[Competitions](#)
[Datasets](#)
[Scripts](#)
[Jobs](#)
[Community](#)
[Sign up](#)
[Login](#)

Welcome to Kaggle's data science competitions.

[New to Data Science? Tutorials on the Titanic competition >](#)  
[Want to learn from other's code? Kaggle's top rated scripts >](#)

**Download**  
 Choose a competition & download the training data.

**Build**  
 Build a model using whatever methods and tools you prefer.

**Submit**  
 Upload your predictions. Kaggle scores your solution and shows your score on the leaderboard.

**Active Competitions**

All Competitions	Active Competitions
	<b>State Farm Distracted Driver Detection</b> Can computer vision spot distracted drivers? 3 months 239 teams 110 scripts \$65,000
	<b>Santander Customer Satisfaction</b> Which customers are happy customers? 18 days 3894 teams 2478 scripts \$60,000
	<b>Home Depot Product Search Relevance</b> Predict the relevance of search results on homedepot.com 11 days 1944 teams 1486 scripts \$40,000
	<b>BNP Paribas Cardif Claims Management</b> Can you accelerate BNP Paribas Cardif's claims management process? 4.4 days 2947 teams 1692 scripts \$30,000
	<b>2016 US Election</b> Explore data related to the 2016 US Election 339 scripts 699 downloads
	<b>2013 American Community Survey</b> Find insights in the 2013 American Community Survey 1077 scripts 1098 downloads
	<b>World Development Indicators</b> Explore country development indicators from around the world 147 scripts 1694 downloads

[Host](#)
[Competitions](#)
[Datasets](#)
[Scripts](#)
[Jobs](#)
[Community](#)
[Sign up](#)
[Login](#)

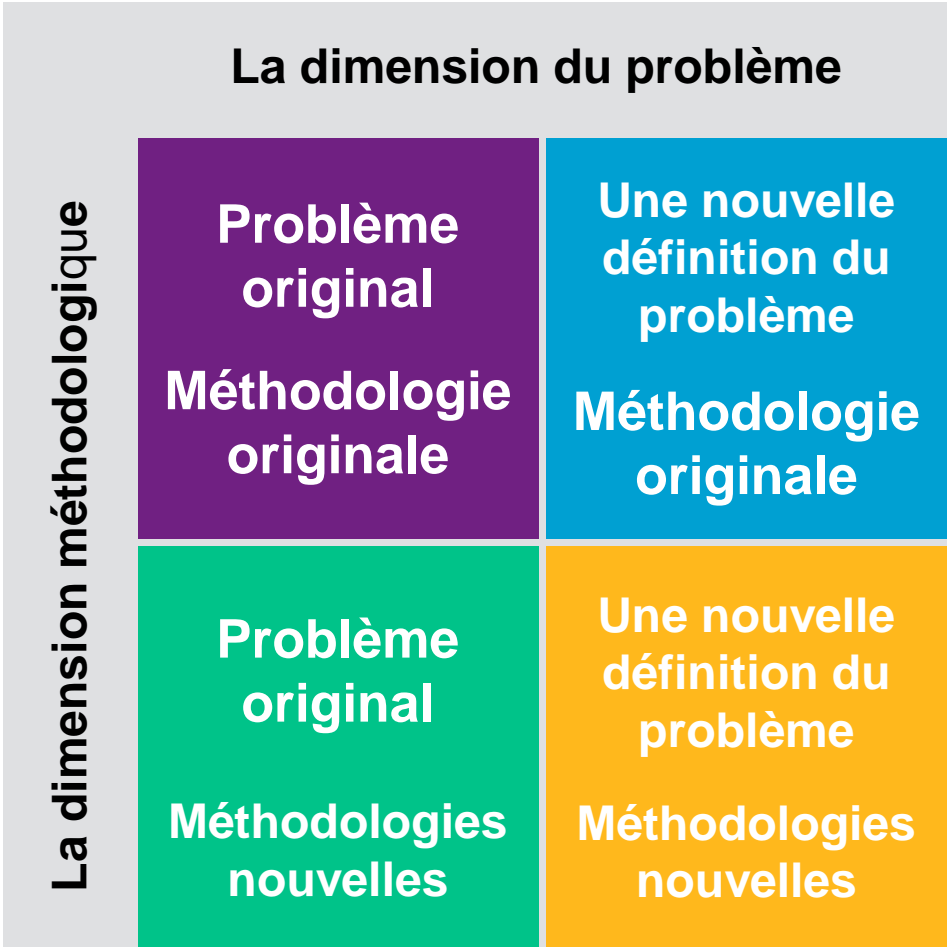
## Kaggle Rankings

Kaggle users are allocated points for their performance in competitions. This page shows the current global ranking. For more information on how we calculate points, please visit the [user ranking wiki page](#).

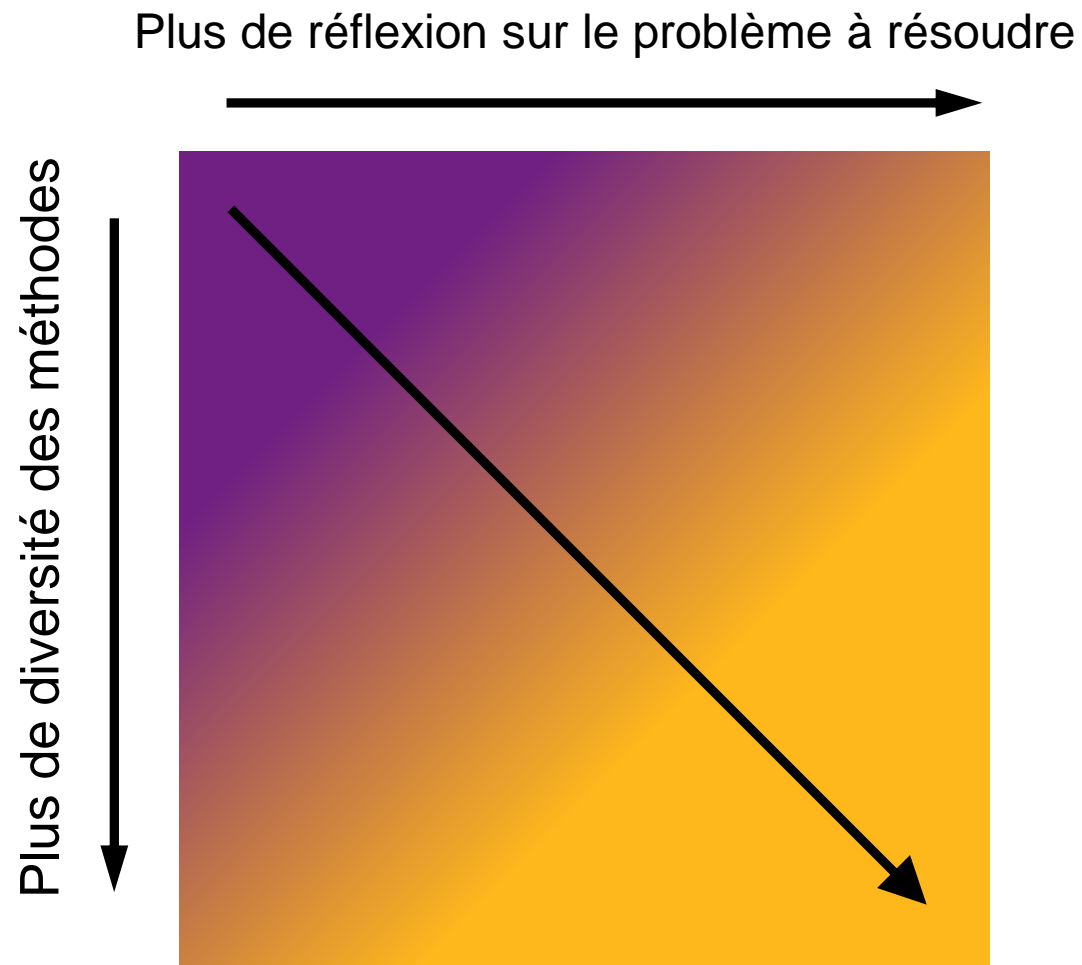
1st	191,154 pts	2nd	189,482 pts	3rd	163,407 pts	4th	144,134 pts	5th	139,658 pts
	<b>Gilberto Titericz</b> 66 competitions Curitiba Brazil		<b>Μαριος Μιχαηλιδης</b> 72 competitions Volos Greece		<b>Stanislav Semenov</b> 31 competitions Moscow Russian Federation		<b>Owen</b> 42 competitions NYC United States		<b>Kohei</b> 70 competitions Tokyo Japan
	<b>Alexander Guschin</b> 21 competitions Moscow Russia		<b>Abhishek</b> 97 competitions Berlin Germany		<b>Leustagos</b> 45 competitions Belo Horizonte Brazil		<b>Cardal</b> 4 competitions Israel		<b>Gert</b> 24 competitions Goes The Netherlands
	<b>y</b> 55 competitions South Korea		<b>Mike Kim</b> 48 competitions Washington DC United States		<b>clustifier</b> 56 competitions Israel		<b>Mario Filho</b> 17 competitions São Paulo Brazil		<b>utility</b> 15 competitions Moscow Russian Federation
	129,891 pts		122,712 pts		119,591 pts		114,004 pts		108,786 pts
	102,606 pts		100,359 pts		100,128 pts		99,000 pts		95,403 pts

# Les dimensions analytiques

La vision complète

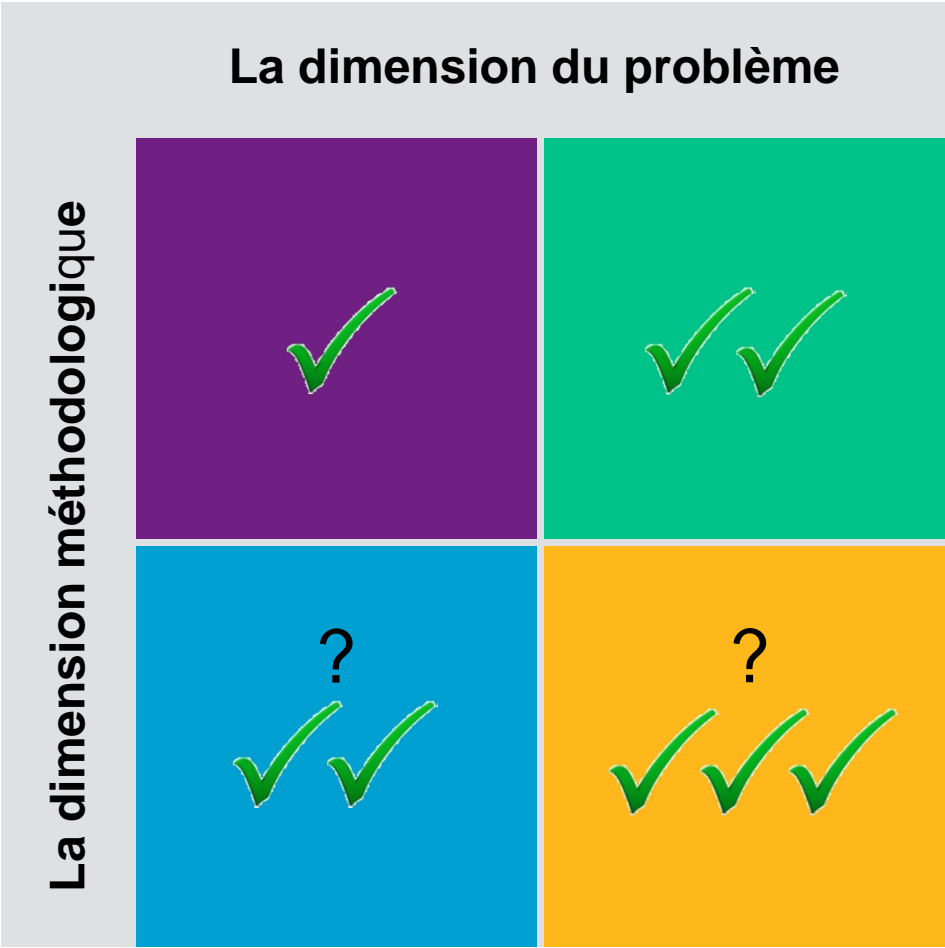


# Ce n'est pas vraiment 2 x 2, mais plutôt en continu...



# Les dimensions analytiques

Où se trouve la valeur ?



# L'application du Machine Learning dans la tarification IARD

## Sommaire

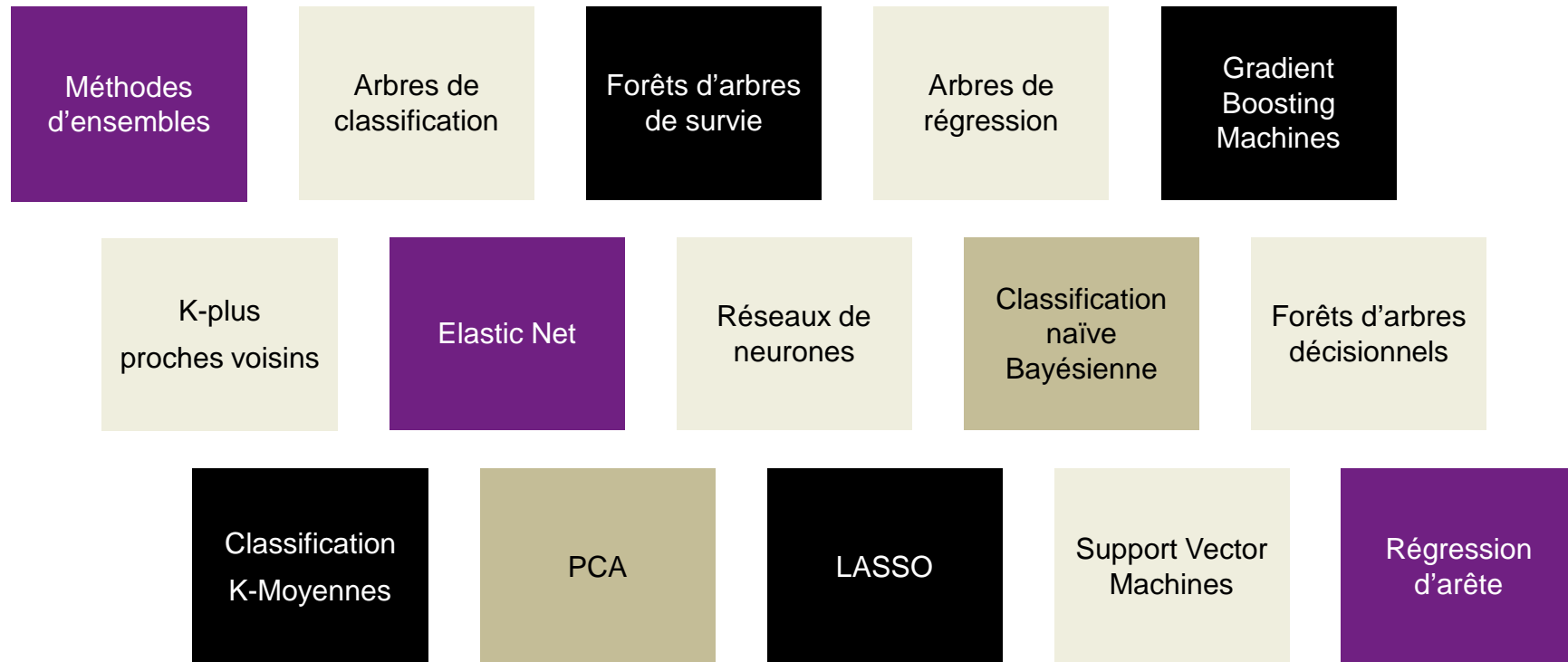
- **Les dimensions analytiques**
  - La dimension du problème
  - La dimension méthodologique
- **Techniques de modélisation**
  - Les plus populaires
  - Exemple: GBMs
  - Calibration
- **Comment utiliser le machine learning ?**
  - La valeur ajoutée
  - L'interprétation
  - L'utilisation du modèle
- **Faut-il passer du temps sur la méthode ou le problème?**





# Techniques de modélisation

Quelles sont les nouvelles méthodologies ?



## Qui a gagné Kaggle ?

- **Gradient Boosted Machines (GBM)** ont eu le plus de succès
- **Création/sélection de variables** a déterminé le succès des candidats
  - La nature de la compétition Kaggle et le partage des benchmarks amène la plupart de compétiteurs à utiliser les mêmes méthodes –ce qui différencie les résultats est la capacité de créer de variables explicatives à partir des informations données

### Nombre de fois où la méthodologie était utilisé parmi les 3 meilleurs résultats (2015)

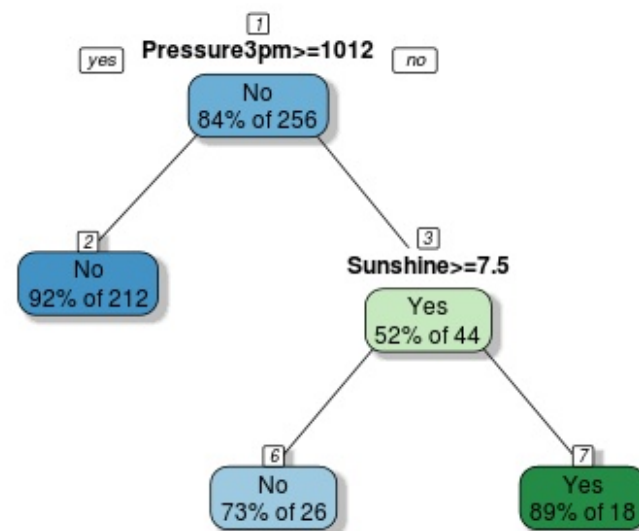
Thème de la compétition	Support Vector Machine	Gradient Boosted Machine	Neural Network	Mixed Method Ensemble	Random Forest	Total
Tout type	1	19	10	10	1	41
Assurance	-	3	-	4	-	7

# Gradient Boosted Machines

- **Boosting** : ajouter un ajustement au modèle en analysant les résidus de manière itérative
- Chaque itération fonctionne sur un échantillon pour donner un modèle robuste
- Chaque itération ajoute une fraction ( $\lambda$ ) du nouveau modèle:

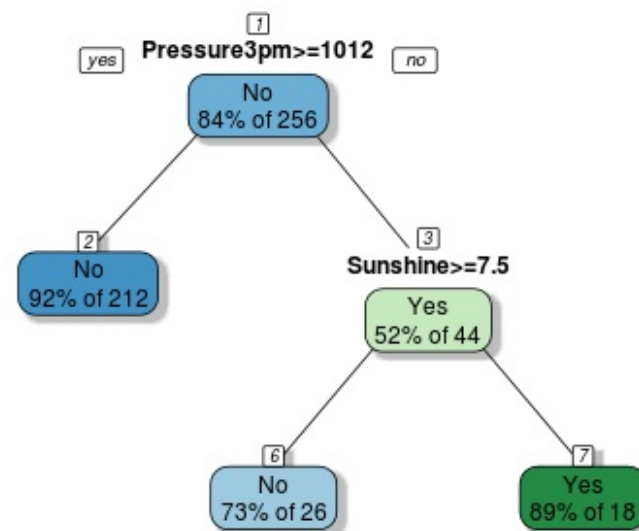
$$f(x) = \lambda \sum_{n=1}^N f_n(x)$$

- Souvent les modèles se constituent des arbres de décision, mais peut-être résultant à n'importe quel type de modèle
- Le paramétrage est calibré en utilisant la validation croisée afin de donner les résultats les plus prédictifs



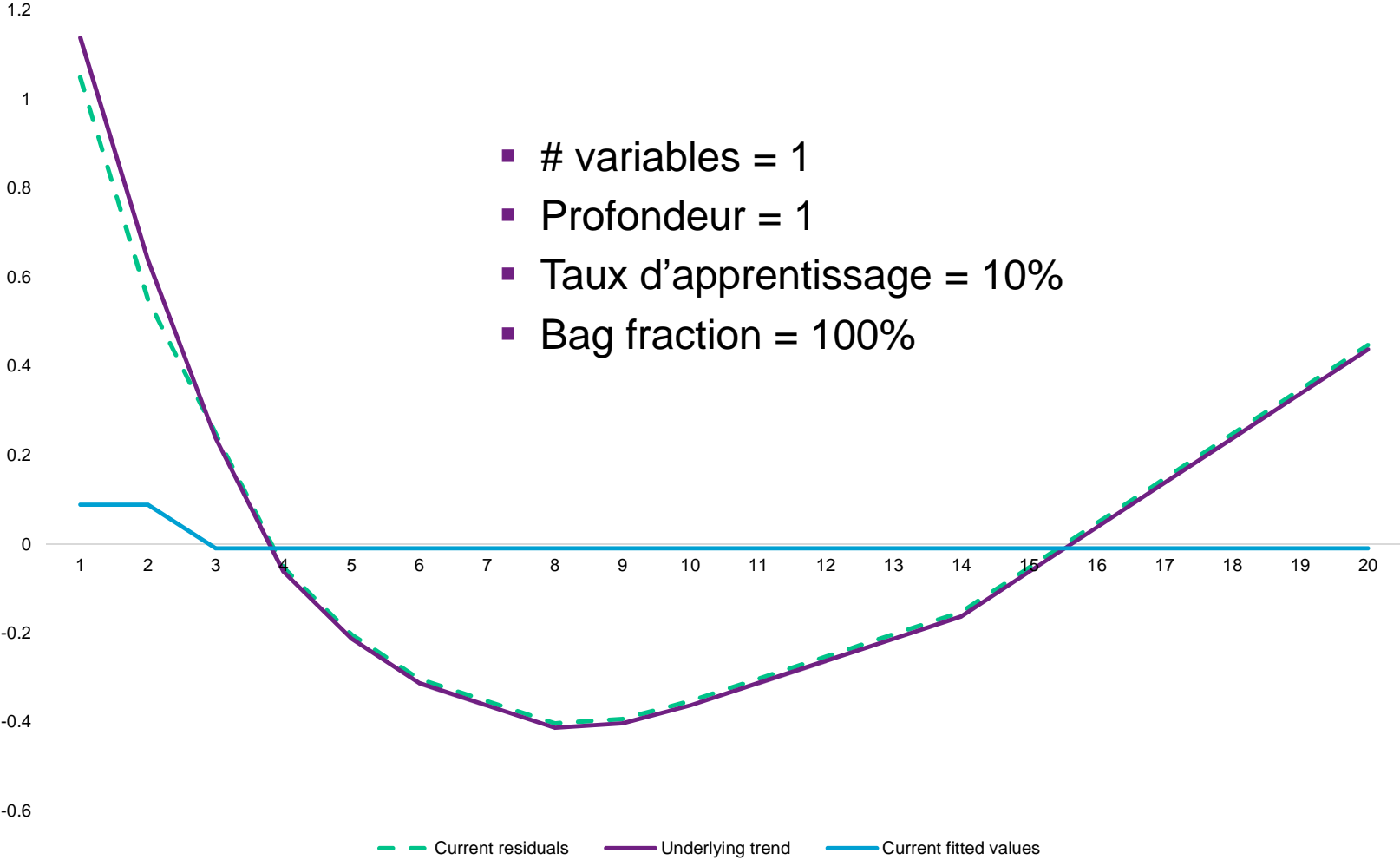
# Calibrage

- **Taux d'apprentissage ( $\lambda$ )**
  - Proportion du modèle pris en compte à chaque itération
  - $M1 = M0 + \lambda \times \text{Prédiction}$
- **Profondeur d'arbre**
  - Equivaut la taille des interactions entre variables
- **Nombre d'arbres** (itérations) maximum
- **Bag fraction**
  - Proportion des variables utilisées à chaque itération
  - On peut également travailler sur un échantillon des données



# Un exemple simple

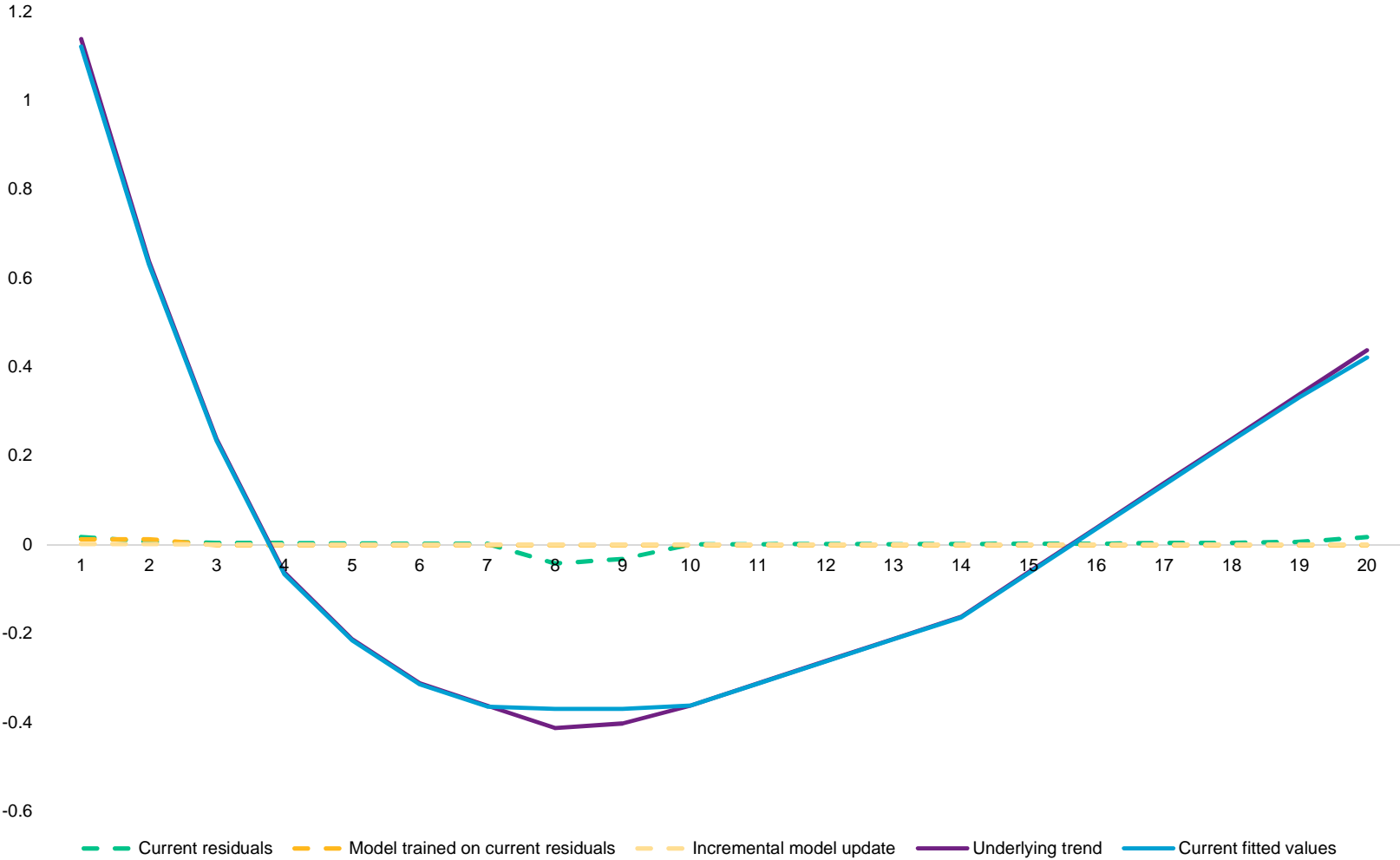
GBM results at iteration 1





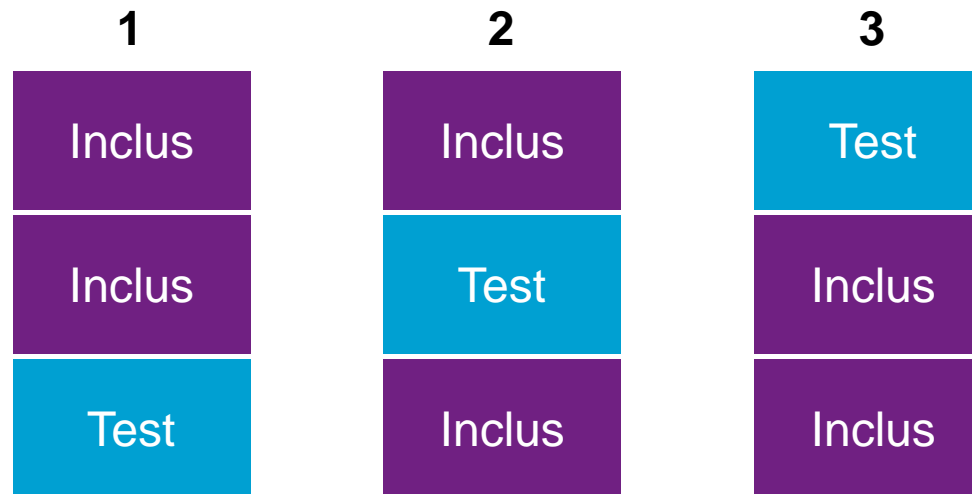
# Un exemple simple

GBM results at iteration 200



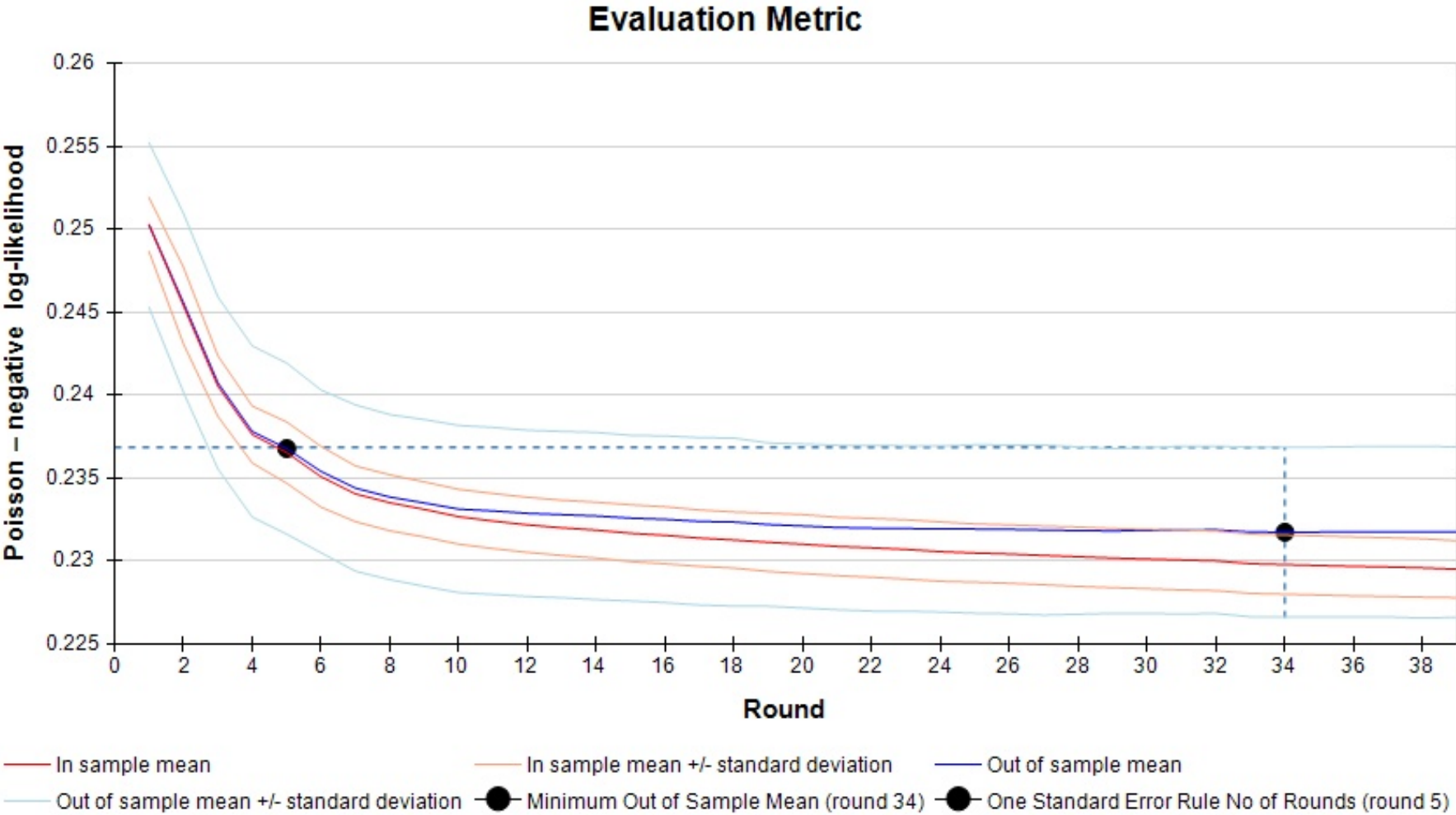
## Calibration (le “tuning”)

- La validation croisée est utilisée pour tester les valeurs possibles pour les paramètres et choisir les options qui donnent le meilleur résultat



- Les graphiques resultants peuvent être utilisés pour déterminer le choix d’hypothèse optimal.
- Notamment le nombre d’arbres à lancer.

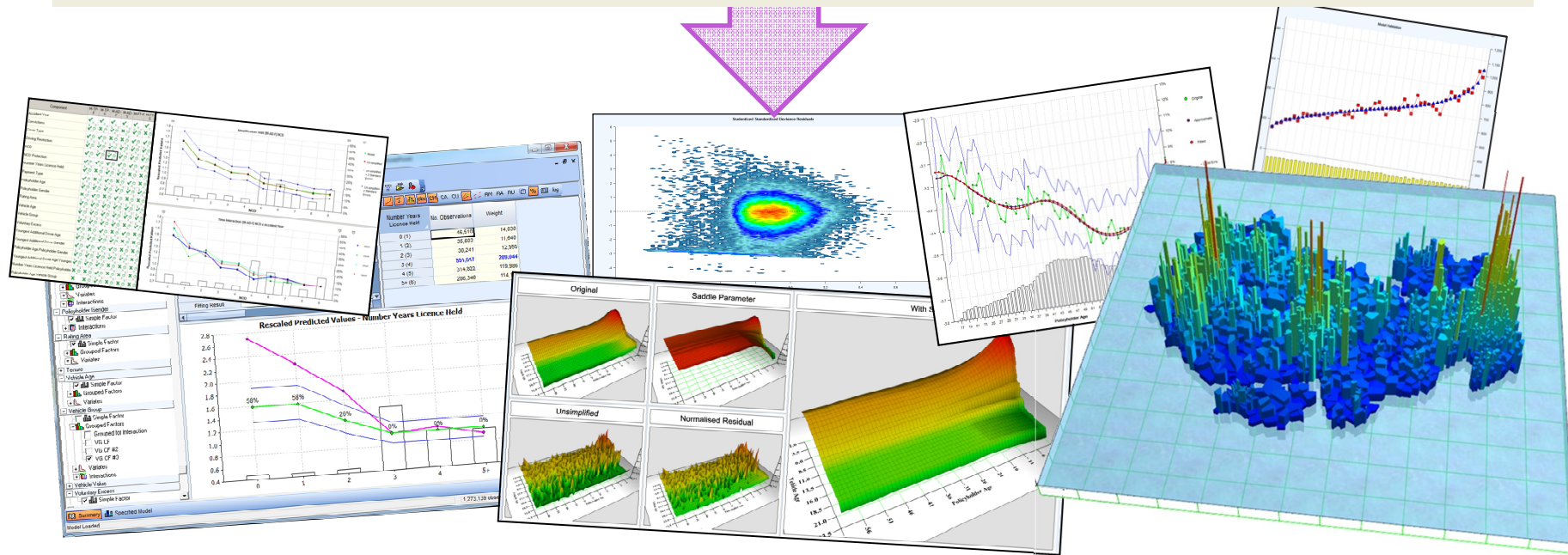
# Un exemple de la validation croisée



# Puissance prédictive versus compréhension du risque

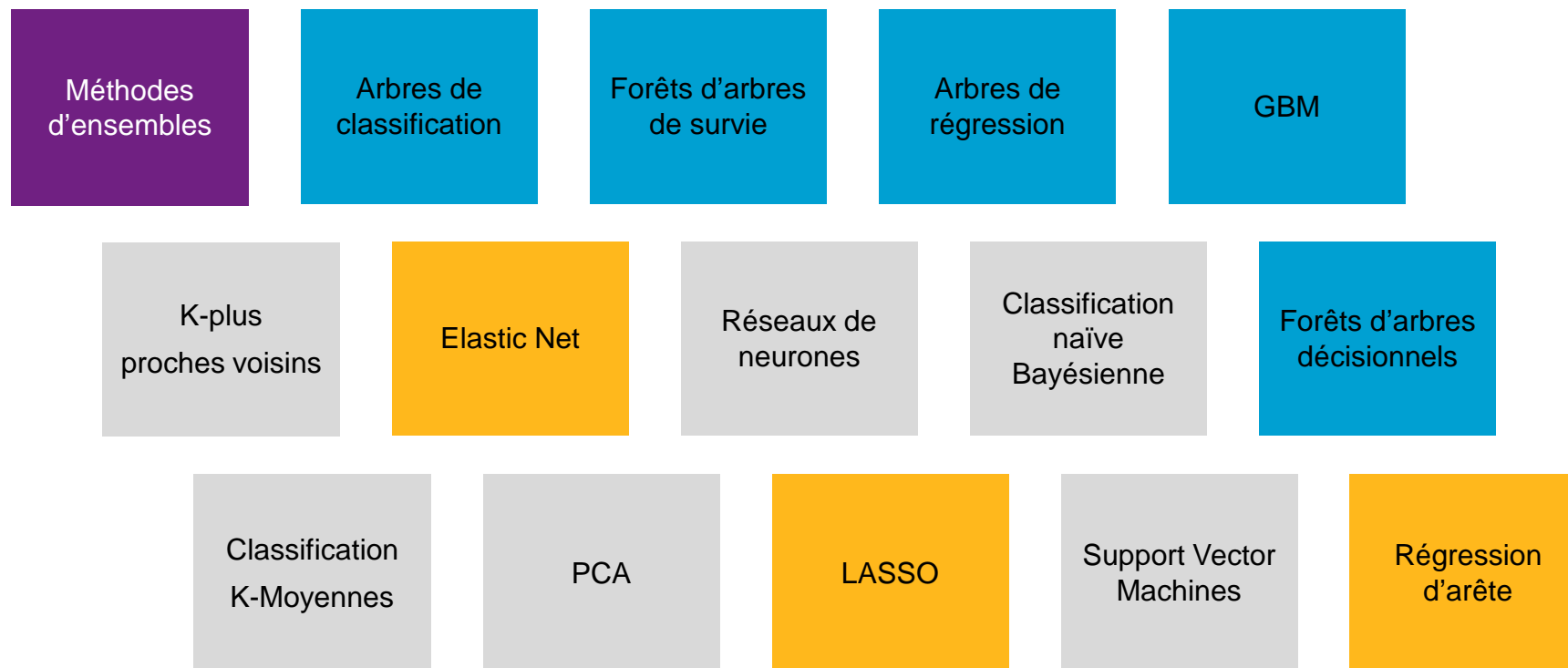


L'utilisation des méthodologies de *machine-learning* peut faciliter l'identification des facteurs explicatifs, mais au détriment de la compréhension



# Techniques de modélisation

Quelles sont les nouvelles méthodologies ?



# Régression pénalisée

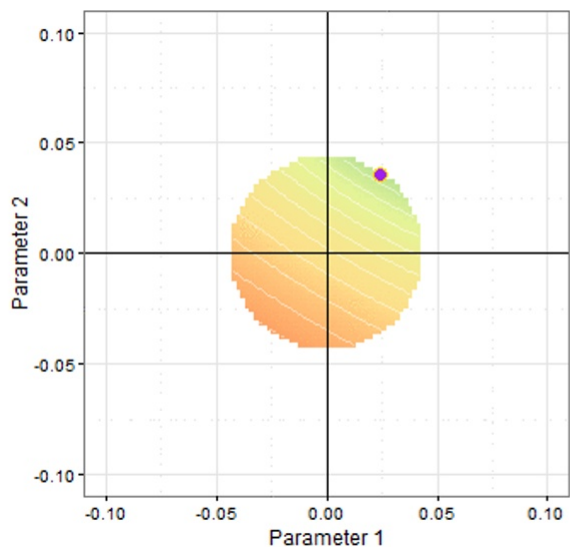
$f(\underline{x}) = g^{-1}(\mathbf{X}.\underline{\beta})$  où  $\underline{\beta}$  minimise :

GLM      Lasso      Ridge

$$L(\beta|X, y) + \lambda_1 \sum_i |\beta_i| + \lambda_2 \sum_i \beta_i^2$$

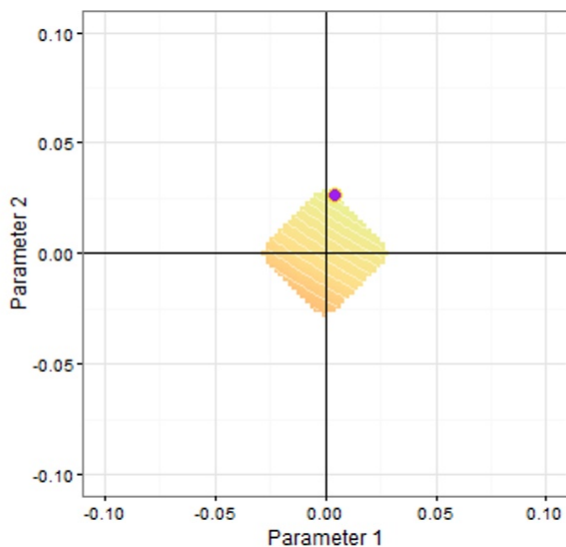
Elastic Net

Ridge  $\sum_i \beta_i^2$



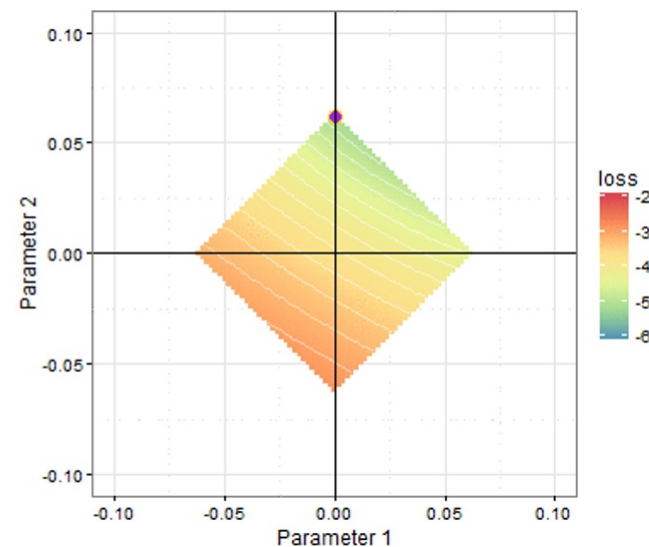
Pénaliser les grands paramètres, mais ne les réduit pas à zéro

Elastic Net



Mélange des deux

Lasso  $\sum_i |\beta_i|$



Pénalité réduit les paramètres non-significatifs à zéro – utile pour la sélection des variables

# Régression pénalisée

Cas d'exemple - la classification des véhicules



## Caractéristiques physiques

Ex., dimensions



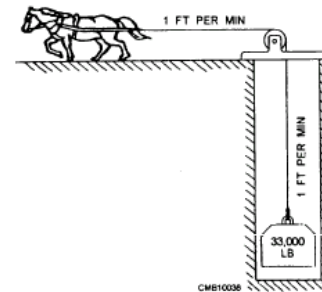
## Caractéristiques techniques

Ex., Carrosserie, modèle, sûreté



## Caractéristiques mécaniques

Ex., fuel, taille moteur



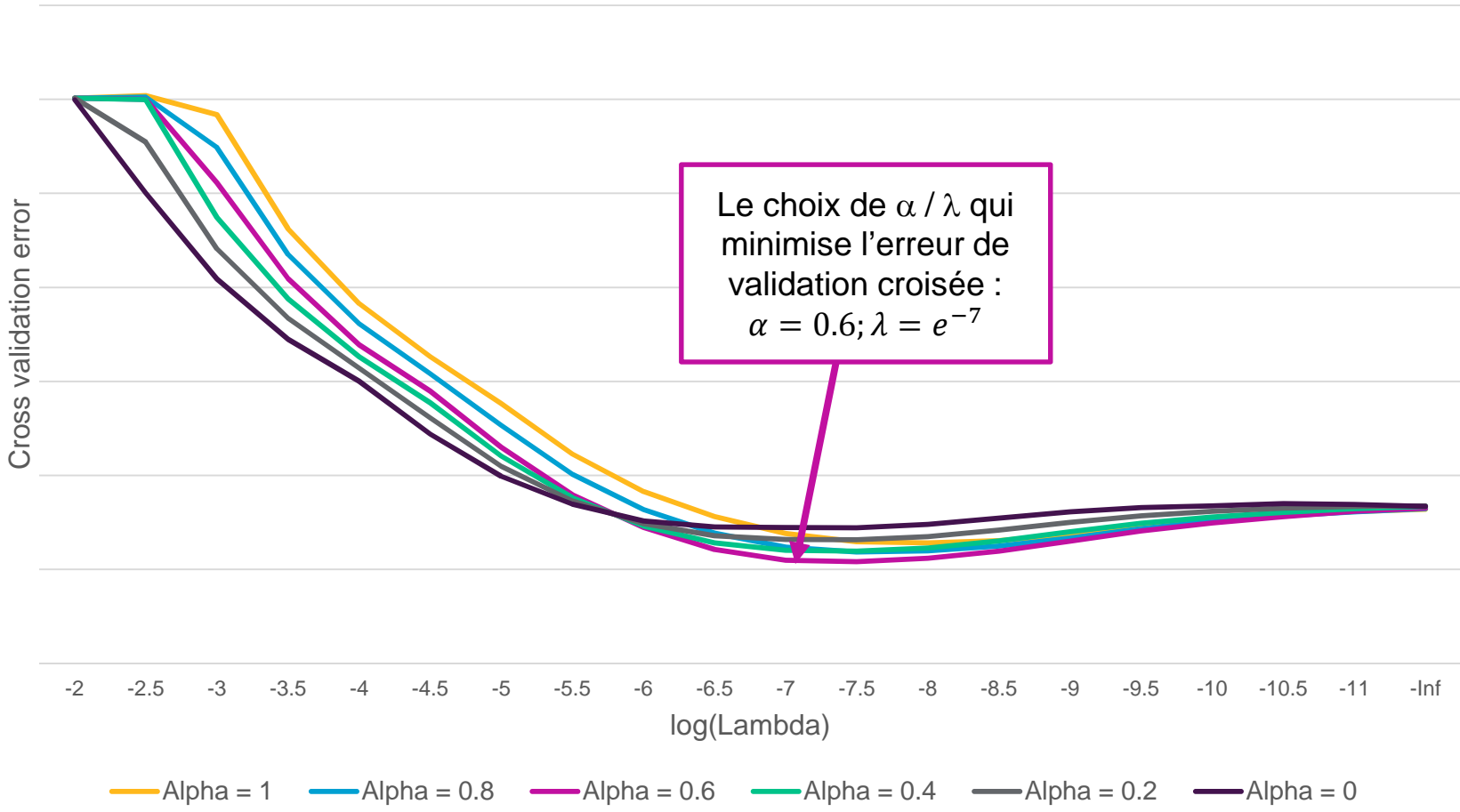
## Performance

Ex., vitesse maximum, HP/CV



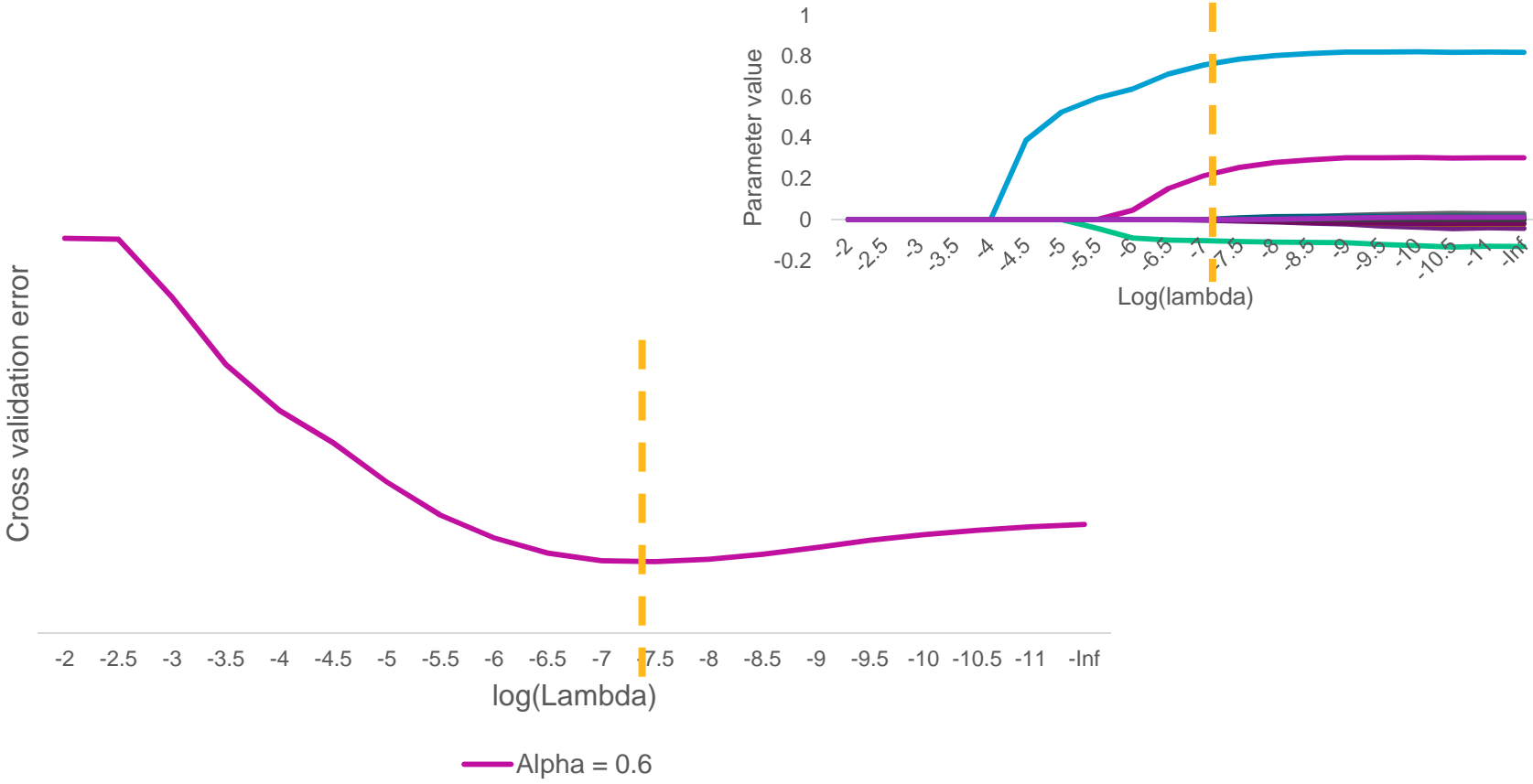
# Régression pénalisée

## Cas d'exemple - la classification des véhicules



# Régression pénalisée

## Cas d'exemple - la classification des véhicules



# Régression pénalisée

## Implémentation

Même que les GLMs!

Age	Facteur
<20	2.12
20-25	1.74
25-30	1.09
30-39	1.00
40-49	0.95
50+	0.06

Groupe	Facteur
1	0.83
2	0.91
3	0.96
4	1.00
5	1.05
6	1.17
7	1.25
8	1.42
9	1.89

Sexe	Facteur
Male	1.00
Female	0.97

# L'application du Machine Learning dans la tarification IARD

## Sommaire

- Les dimensions analytiques
  - La dimension du problème
  - La dimension méthodologique
- Techniques de modélisation
  - Les plus populaires
  - Exemple: GBMs
  - Calibration
- **Comment utiliser le machine learning ?**
  - La valeur ajoutée
  - L'interprétation
  - L'utilisation du modèle
- Faut-il passer du temps sur la méthode ou le problème?

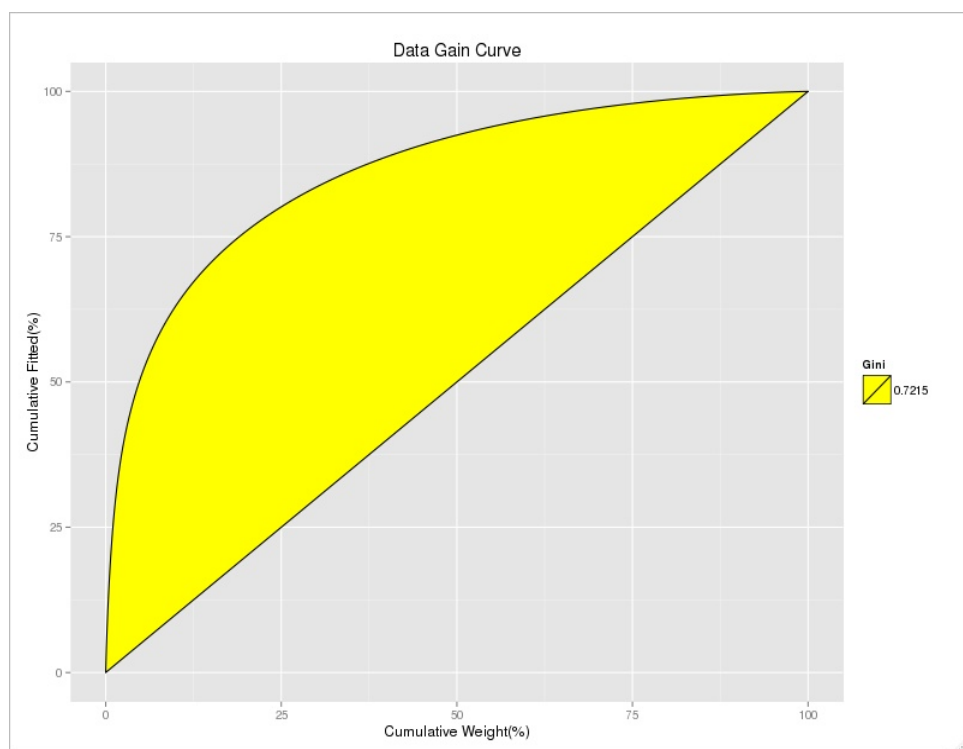


## Les questions à se poser

1. Est-ce que le modèle ajoute de la valeur ?
2. Comment interpréter le modèle ?
  - Est-ce nécessaire ?
3. Comment pouvons-nous utiliser le modèle ?

# Est-ce que le modèle ajoute de la valeur ?

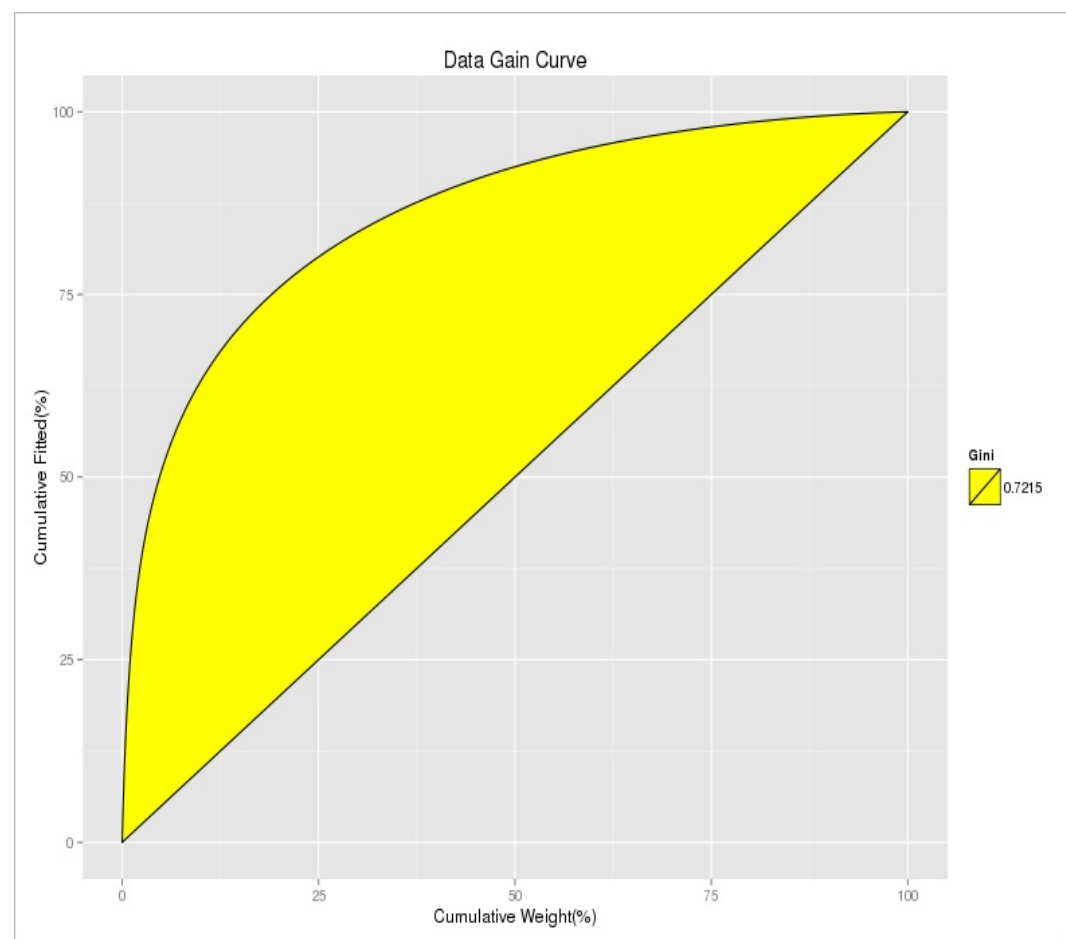
## Le coefficient Gini



- Classer les observations de validation en fonction du rang de leur **valeurs prédites**.
- Faire figurer **la réponse cumulée** versus **l'exposition cumulée** sur un graphe
- Un **meilleur modèle** expliquera **une proportion plus importante de la réponse** avec **une proportion plus petite de l'exposition**
- ...ce qui donnera un **coefficient Gini** plus grand (surface en jaune)

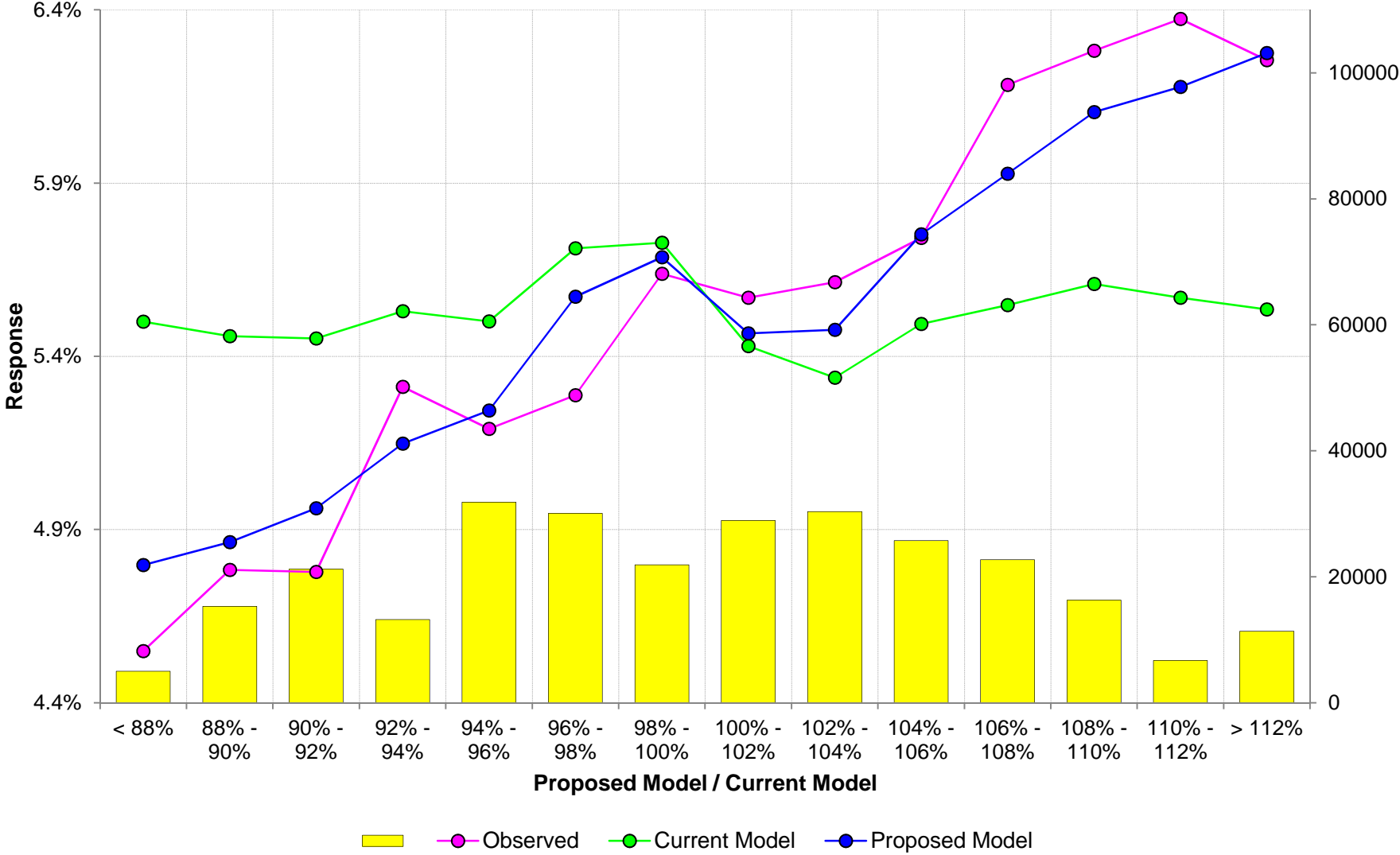
## Mais...

- Prenez votre modèle...
- Multipliez par 123
- Prendre le carré
- Ajouter 74 milliards
  
- ...et le coefficient Gini reste le même!



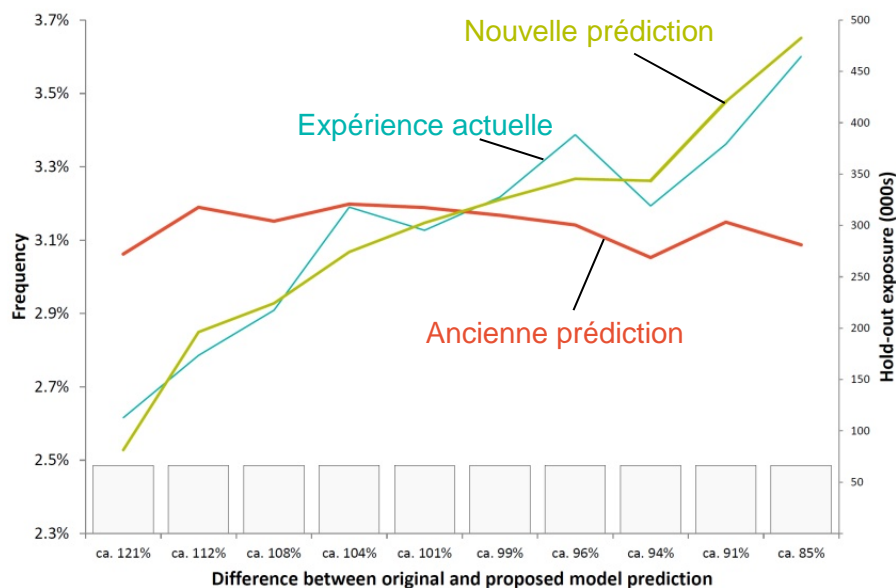


# Comparaison des modèles



# L'impact financier

- Les erreurs de tarification ne sont pas symétriques
- Les tests de scénarios peuvent mesurer l'impact
- Modèle simplifié peut faire une approximation afin de donner un test heuristique de la valeur créée, avec
  - Une hypothèse de l'élasticité
  - Une hypothèse de changement min/max



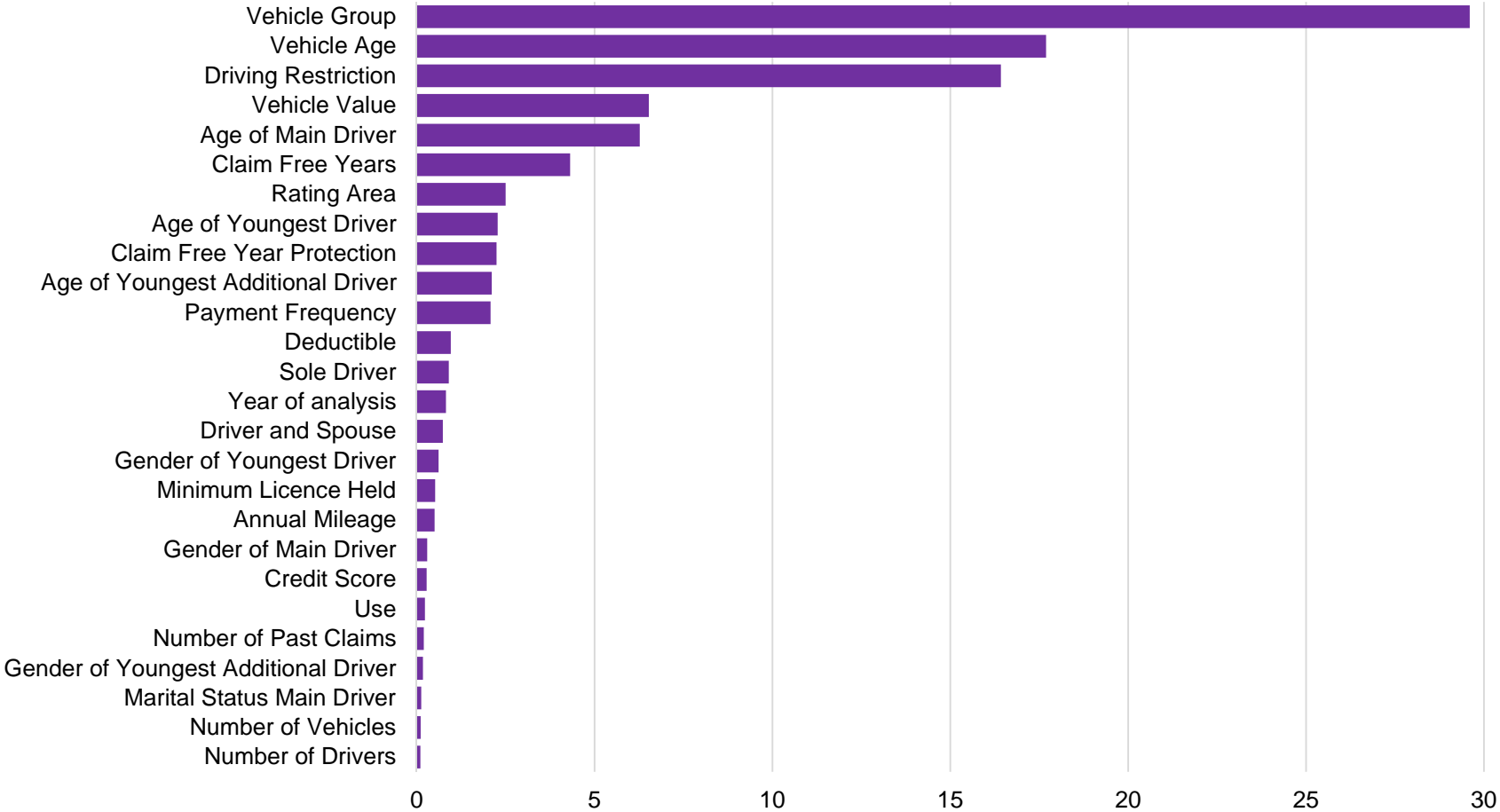
Ratio	Nouvelle prime	Volume attendue	Sinistres	$\Delta$ profit
121%	$P_1$	$V_1$	$C_1$	$X_1$
...	$P_2$	$V_2$	$C_2$	$X_2$
...	...	...	...	...
...	$P_{99}$	$V_{99}$	$C_{99}$	$X_{99}$
85%	$P_{100}$	$V_{100}$	$C_{100}$	$X_{100}$
<b>Valeur créée</b>				<b>€X</b>

## Les questions à se poser

1. Est-ce que le modèle ajoute de la valeur ?
- 2. Comment interpréter le modèle ?**
  - Est-ce nécessaire ?
3. Comment pouvons-nous utiliser le modèle ?

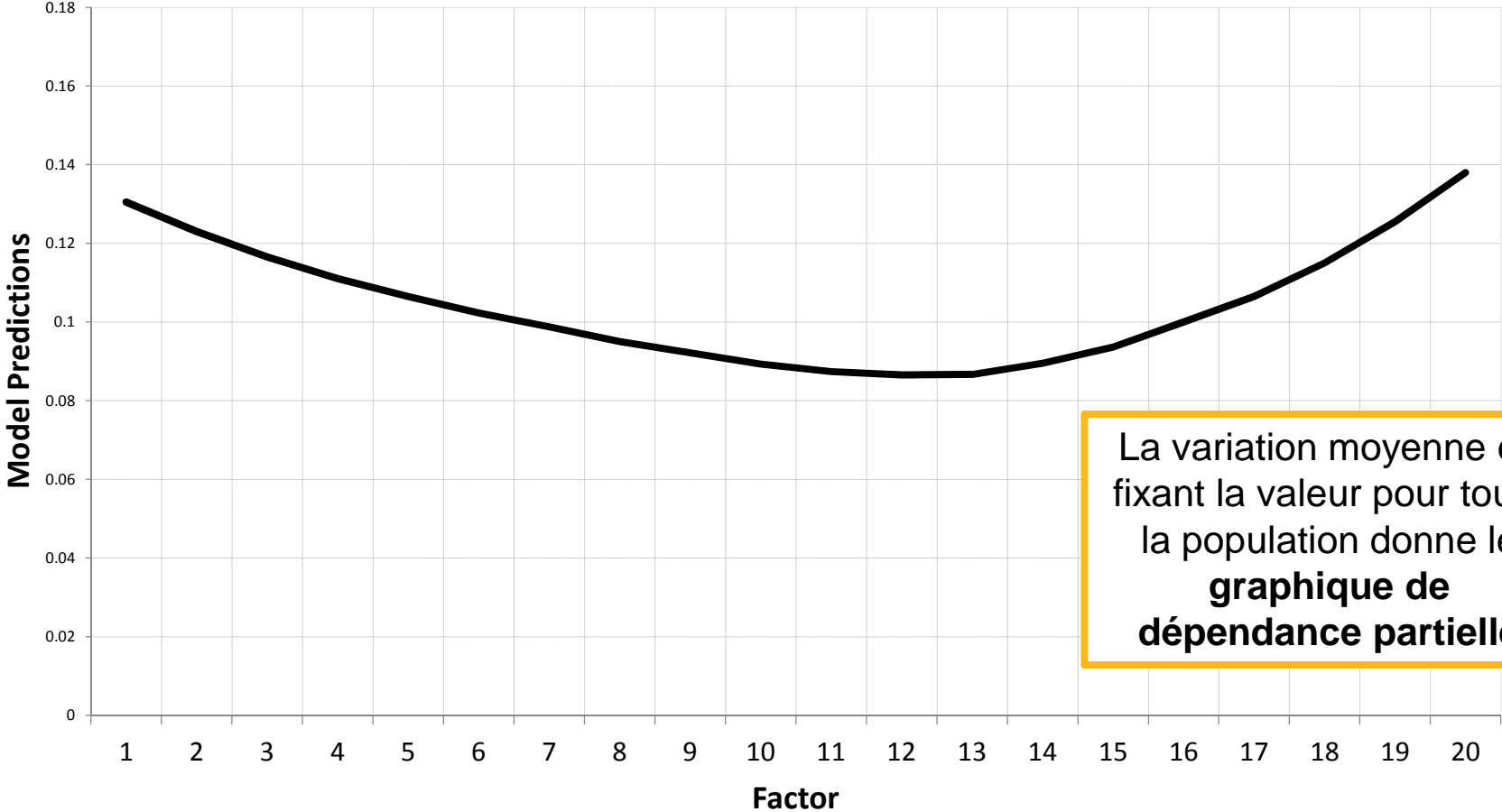
# Importance des facteurs

L'influence relative des variables peut être mesurée par la réduction totale de l'erreur attribuée à des règles de partition définie par chaque variable, à travers tous les arbres du GBM.



# Le graphique de dépendance partielle

## Exemple



## Les questions à se poser

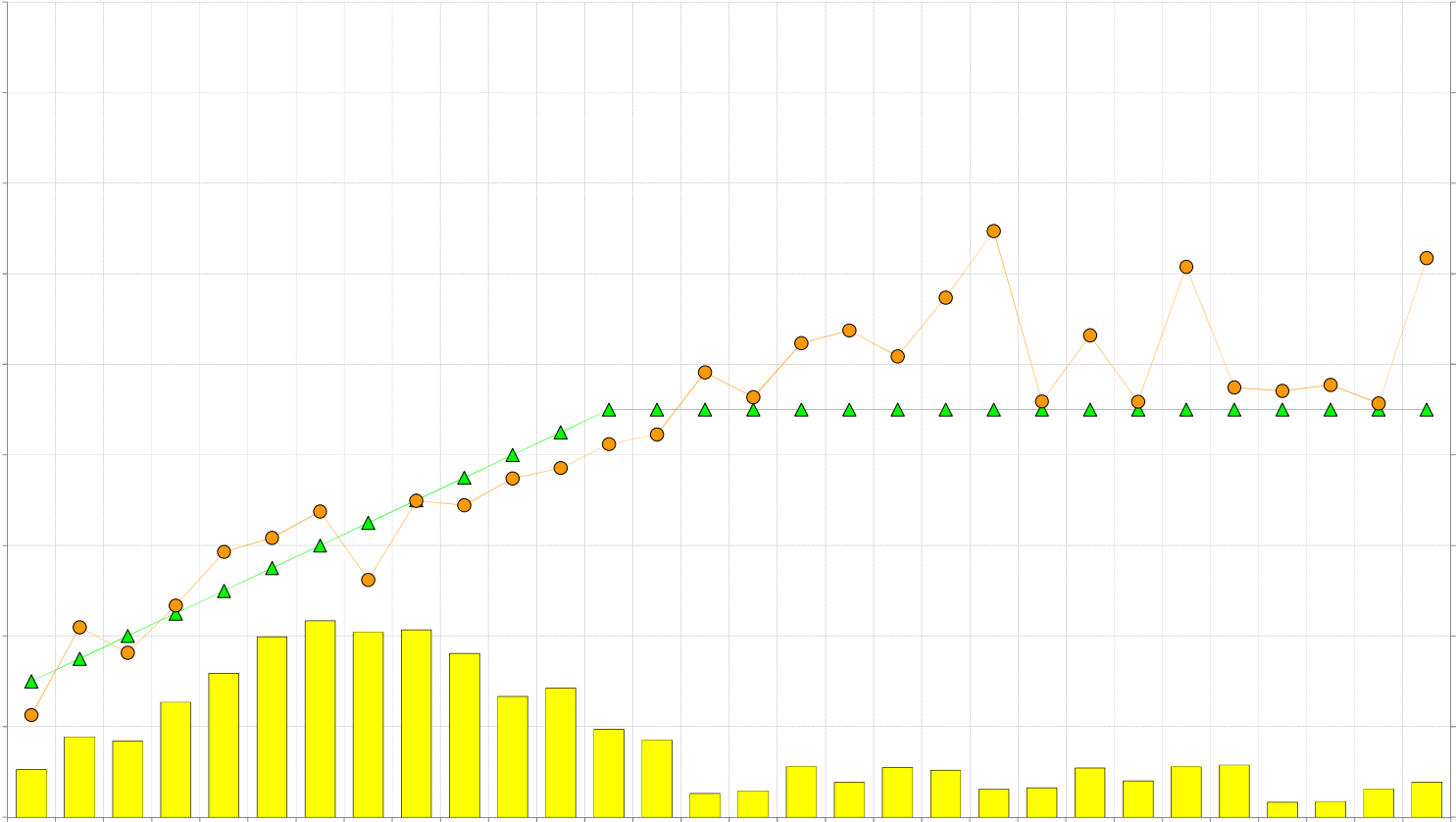
1. Est-ce que le modèle ajoute de la valeur ?
2. Comment interpréter le modèle ?
  - **Est-ce nécessaire ?**
3. Comment pouvons-nous utiliser le modèle ?

# Comment pouvons-nous utiliser le modèle ?

Décrire le modèle par un GLM

Choisir les variables du GLM

Utiliser directement



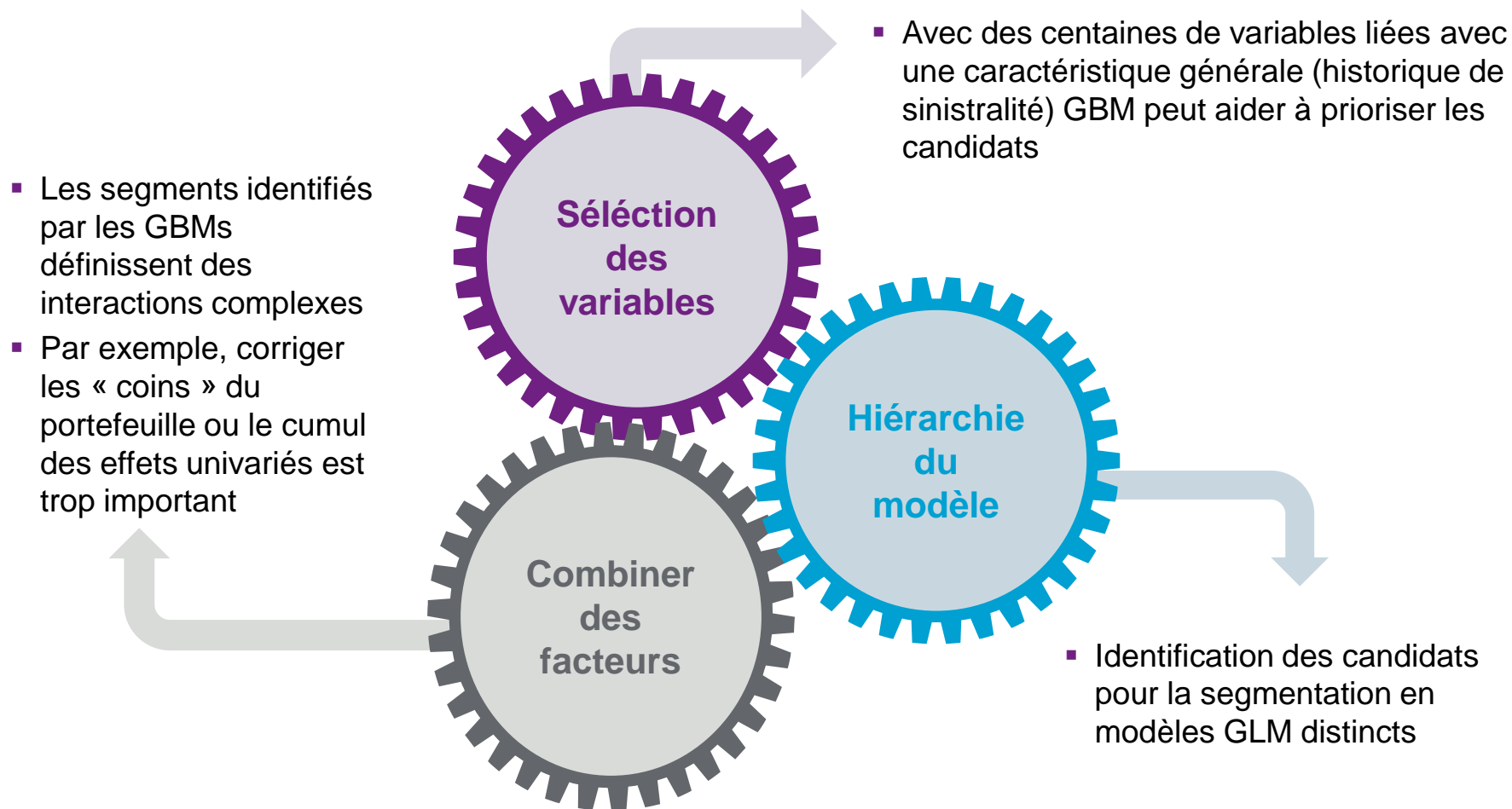


## Comment pouvons-nous utiliser le modèle ?

Décrire le modèle par un GLM

Choisir les variables du GLM

Utiliser directement

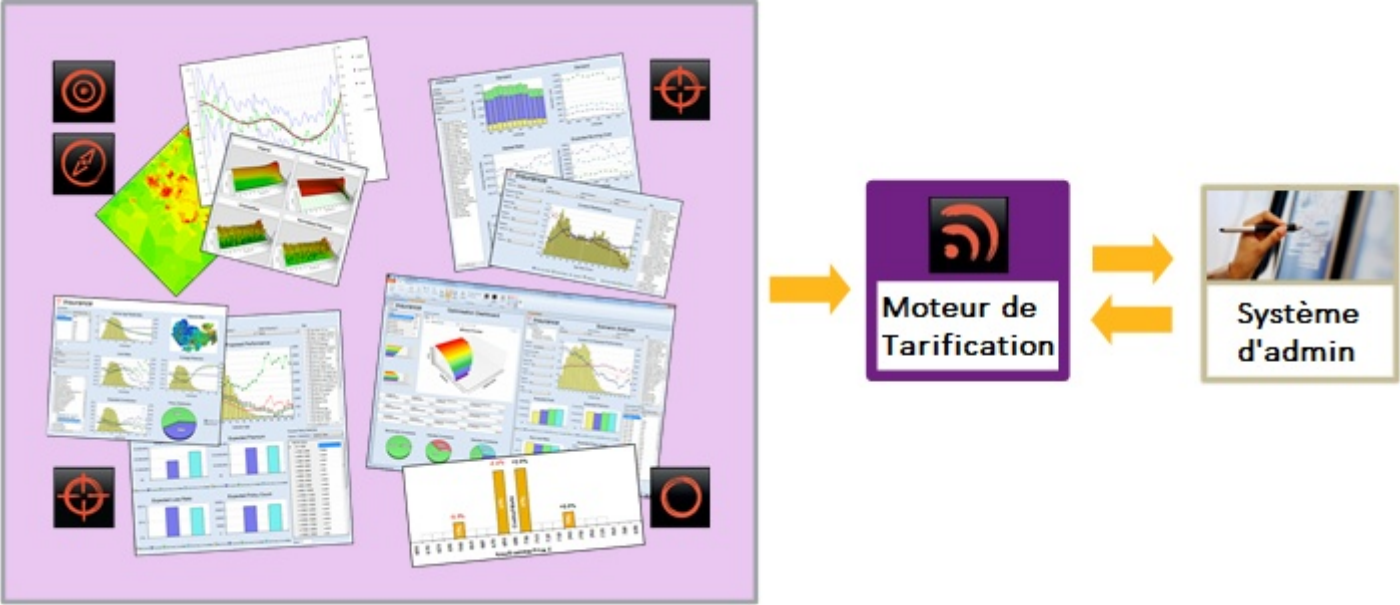


# Comment pouvons-nous utiliser le modèle ?

Décrire le modèle par un GLM

Choisir les variables du GLM

Utiliser directement



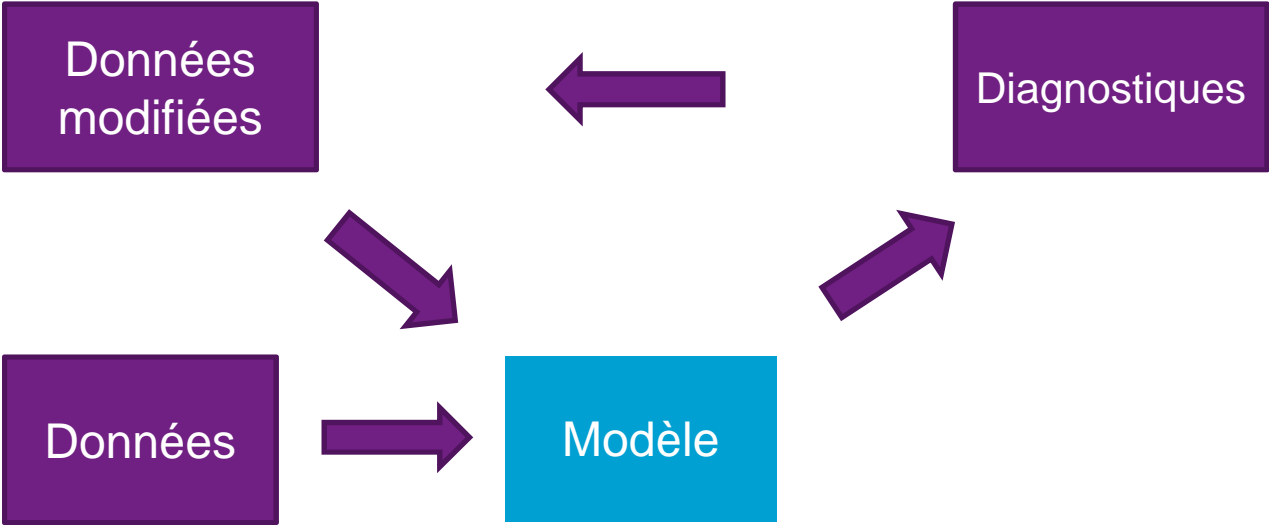
De quelle manière est-on prêt à utiliser un modèle construit automatiquement ?

# Comment pouvons-nous utiliser le modèle ?

Décrire le modèle par un GLM

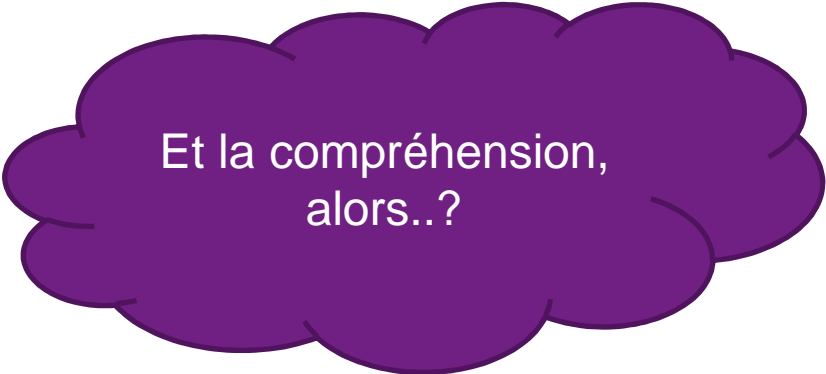
Choisir les variables du GLM

Utiliser directement



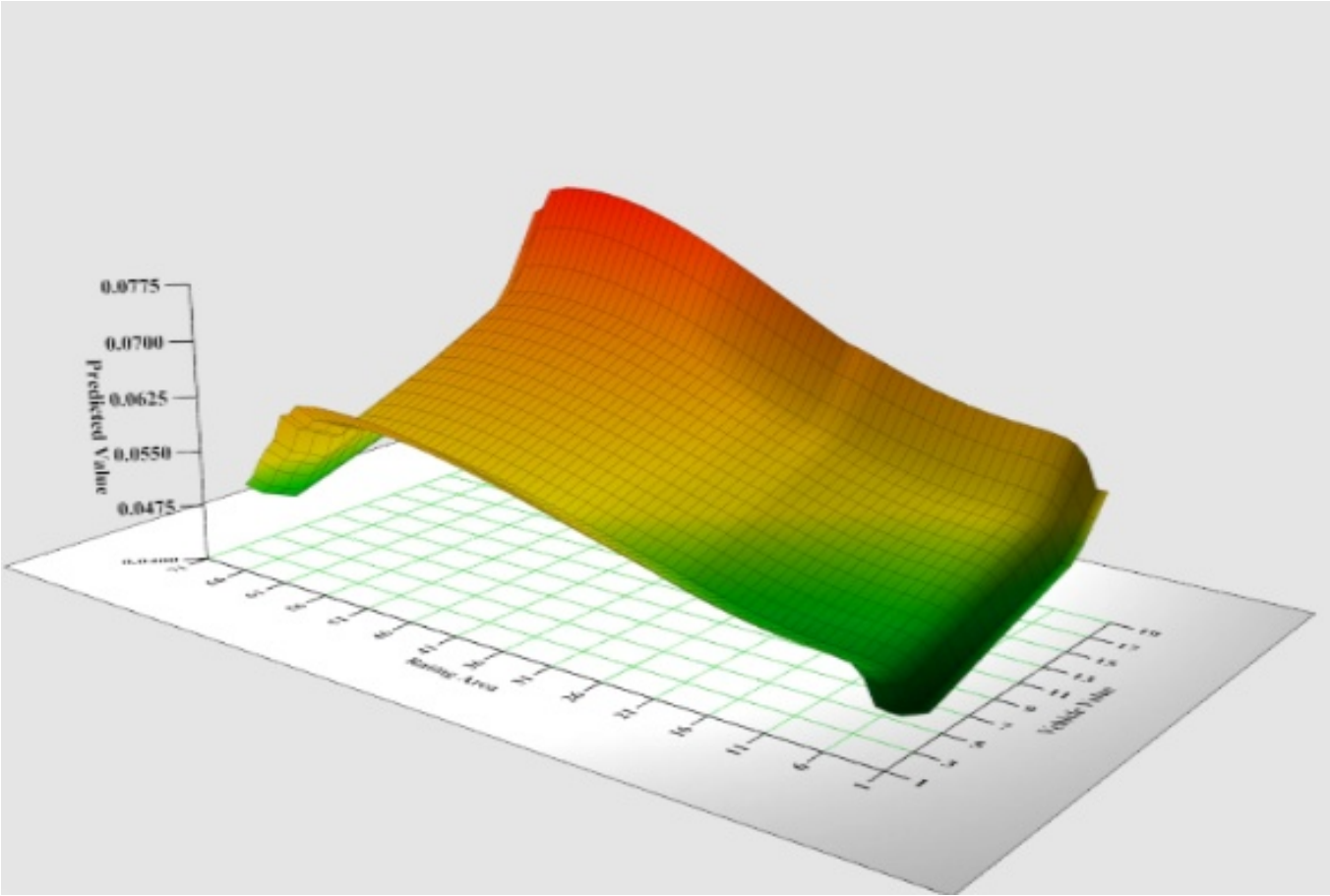
## Les questions à se poser

1. Est-ce que le modèle ajoute de la valeur ?
2. Comment interpréter le modèle ?
  - Est-ce nécessaire ?
3. **Comment pouvons-nous utiliser le modèle ?**



Et la compréhension,  
alors..?

# Paramétrage automatique d'un GLM



# L'application du Machine Learning dans la tarification IARD

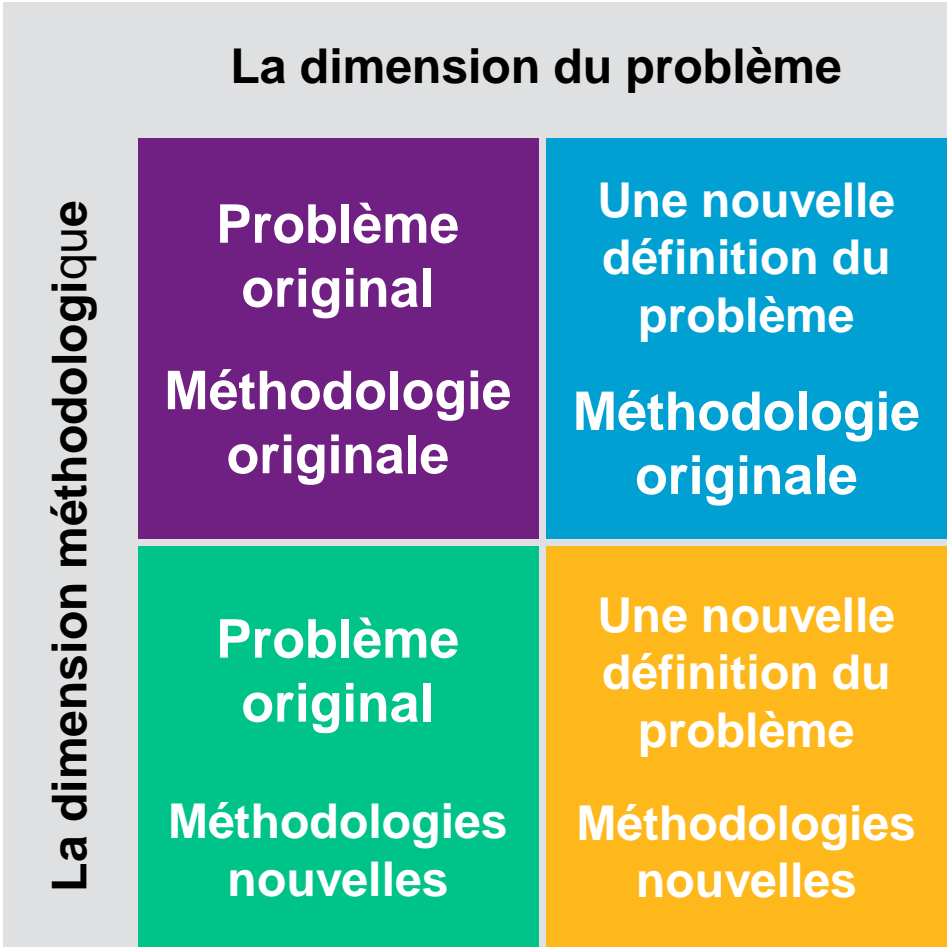
## Sommaire

- Les dimensions analytiques
  - La dimension du problème
  - La dimension méthodologique
- Techniques de modélisation
  - Les plus populaires
  - Exemple: GBMs
  - Calibration
- Comment utiliser le machine learning ?
  - La valeur ajoutée
  - L'interprétation
  - L'utilisation du modèle
- **Faut-il passer du temps sur la méthode ou le problème?**



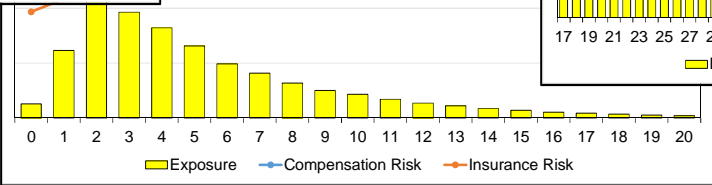
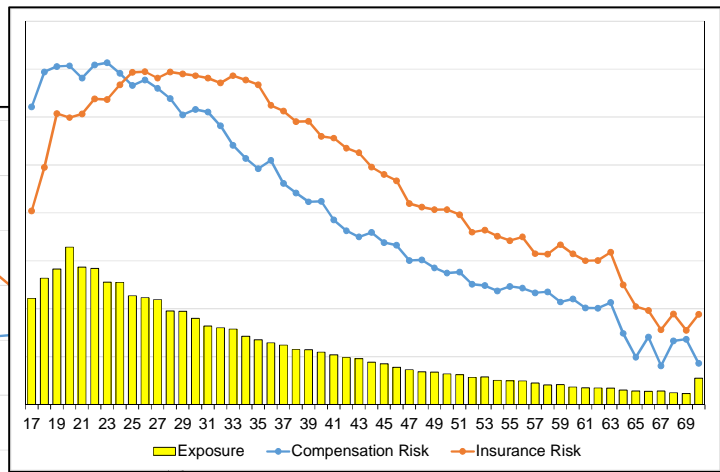
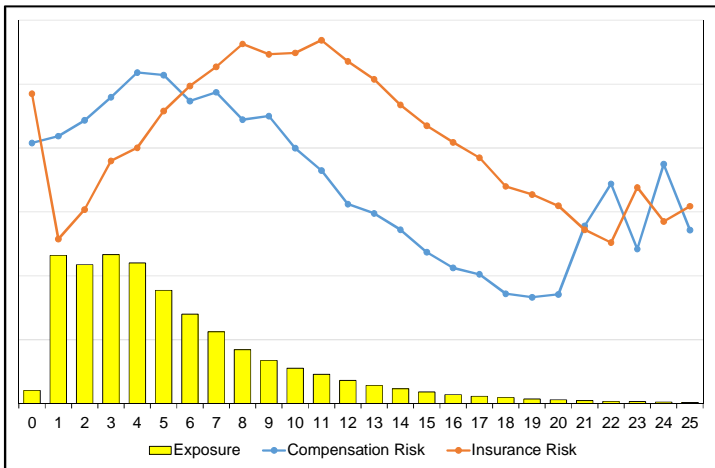
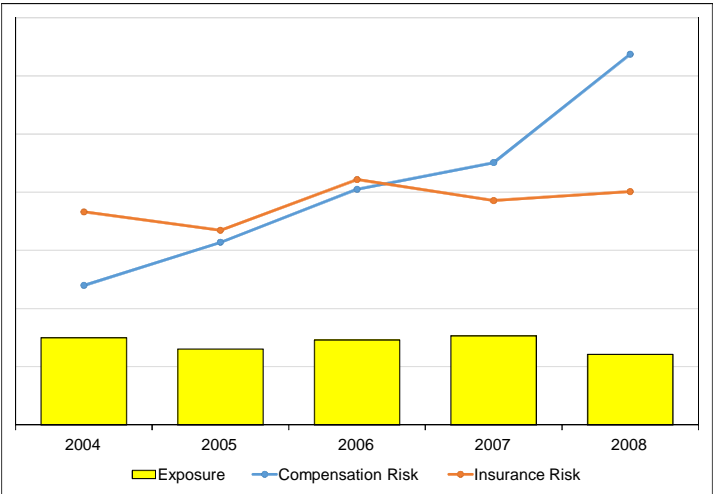
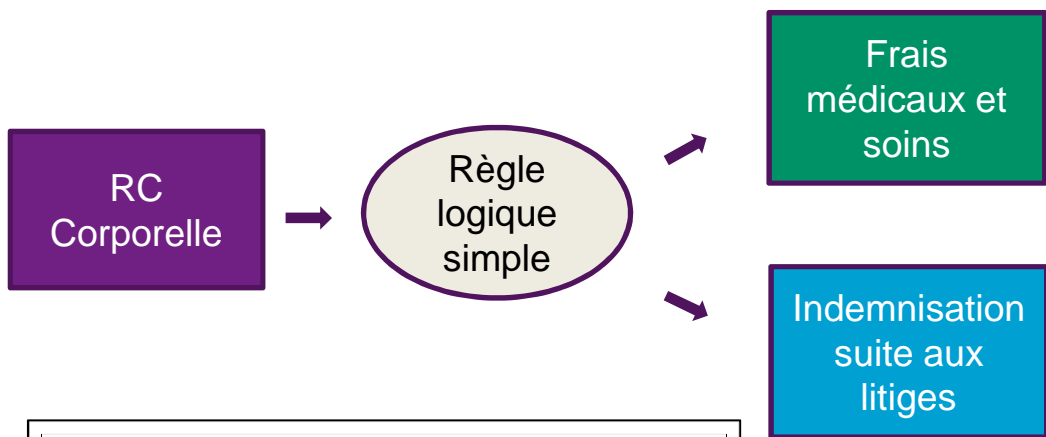
# Les dimensions analytiques

La vision complète

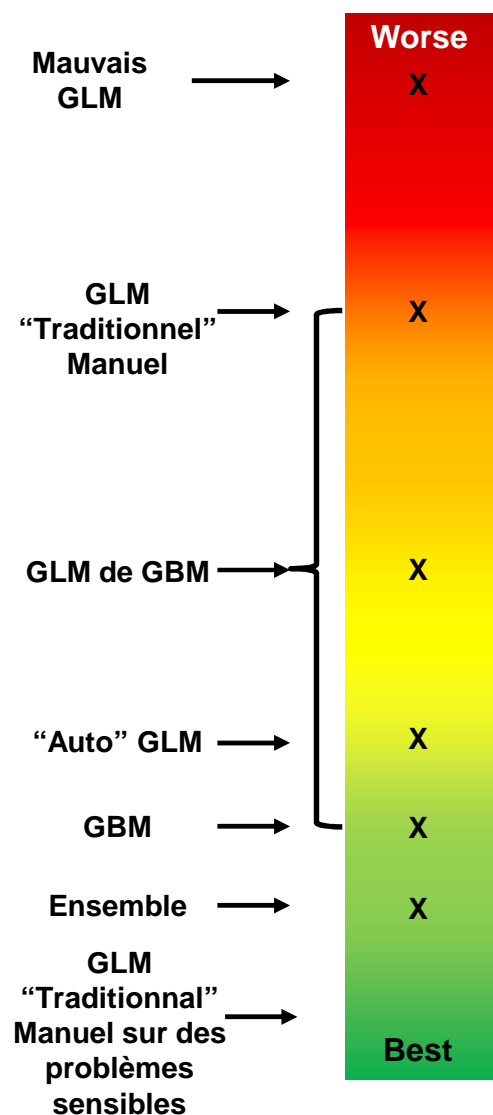




# Sélection des variables et de la réponse



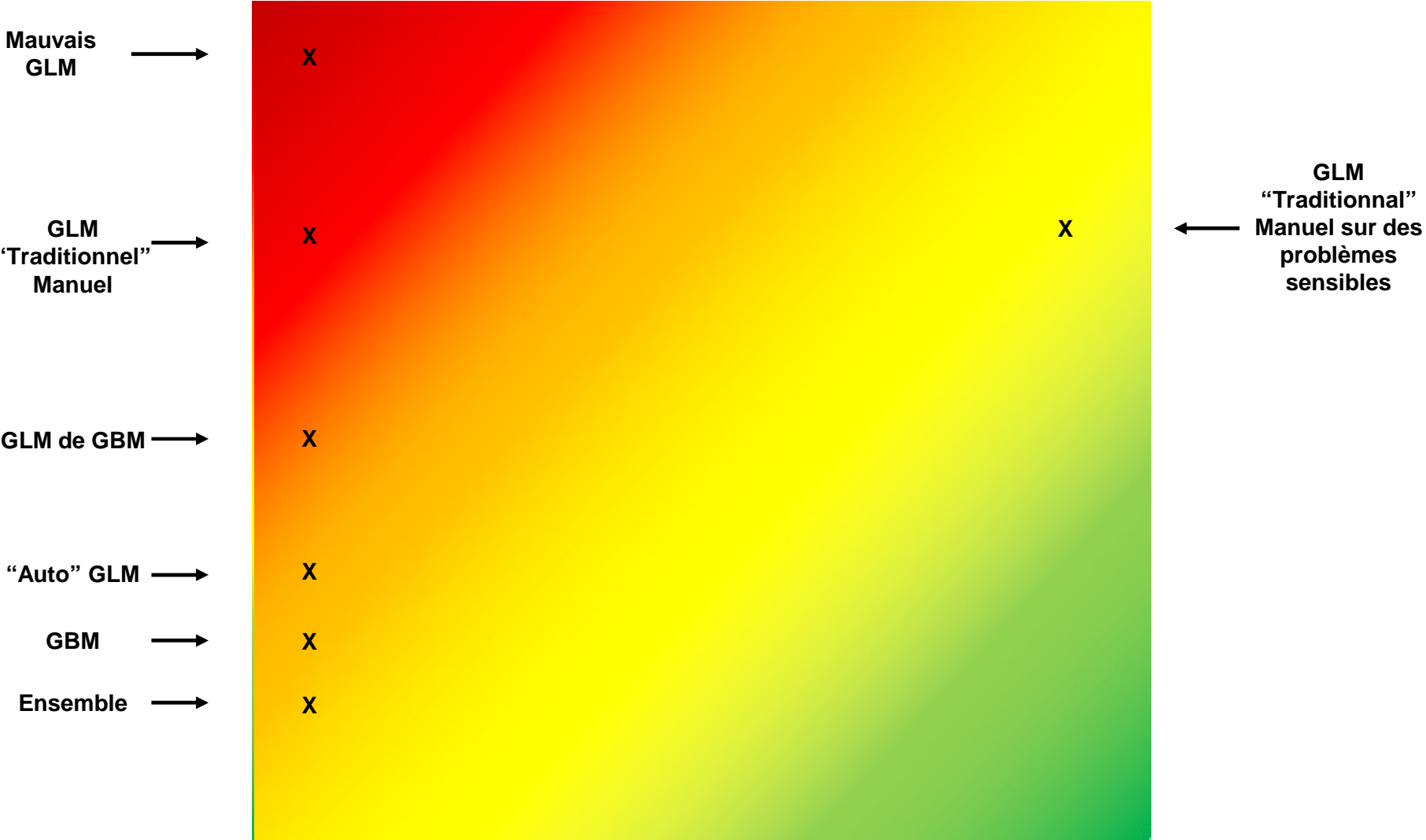
## Conclusions



Si vous pouvez..

- Accepter le fait de ne pas voir le modèle, seulement des diagnostics
  - Accepter des erreurs "dans les coins"
  - Et votre régulateur aussi...
  - Et vous avez un moteur de tarification qui l'implémente
  - Et que vous pouvez traiter de grandes volumes de données ..alors il y a un bénéfice de l'utilisation directe du machine learning
  - Dans d'autres domaines opérationnels (ex : marketing, gestion de sinistres) l'utilisation est moins problématique
  - Si non, les modèles offrent de nouvelles manières de trouver une meilleure compréhension, implémentée dans les GLMs
    - Si vous acceptez des modèles difficiles à interpréter, les GLMs peuvent être calibrés automatiquement également ...
- Ne perdez pas de vue l'importance de la réflexion sur la question en elle-même et la valeur de l'expérience métier

# Conclusions



# Discussion



**Willis Towers Watson** 