



EXPERIAN **INNOVATION FORUM** 2019  
LA DATA AU SERVICE DE  
**L'INNOVATION**

24 octobre | Hôtel Le Casablanca, Maroc



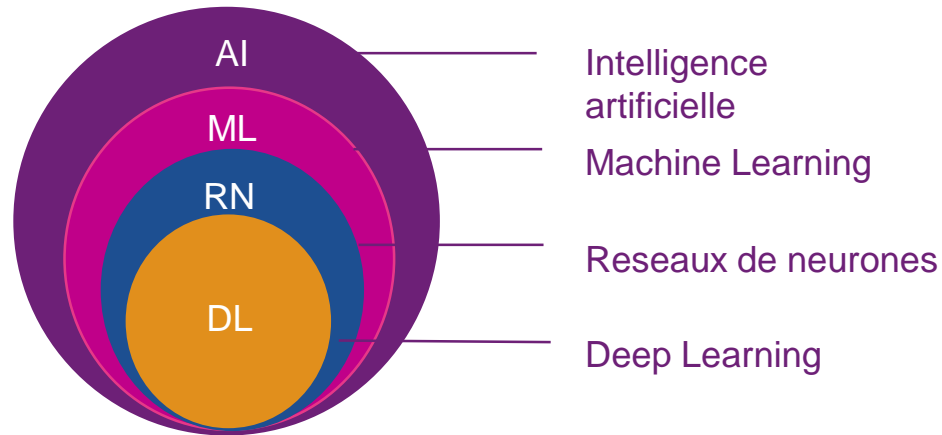


# Nouvelles générations de Machine Learning au service de l'analytics réglementaire

Frédéric Bomy  
Senior Data Modeller  
Manager, Experian



# Evolution de l'intelligence artificielle



## Quels outils pour les construire



## Où les construire

### Plateforme de ML

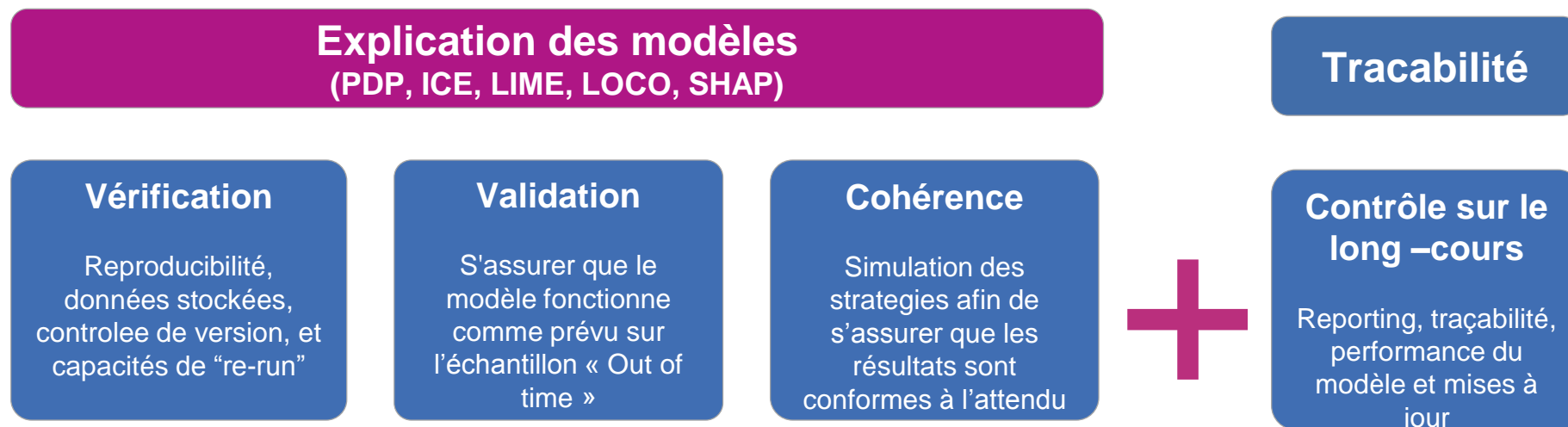
Microsoft ML  
Amazon Web Services  
Open Python  
Open R  
H2O  
Tensor Flow

1. L'intelligence artificielle est une technique qui permet aux machines d'imiter le comportement humain.
2. L'apprentissage automatique est un sous-ensemble de l'intelligence artificielle qui utilise des méthodes statistiques pour permettre aux machines de s'améliorer avec l'expérience
3. Machine learning models goal: best model for the available data and prediction problem
4. Objectif des modèles ML : Le meilleur modèle possible compte tenu de la cible et des données disponibles
5. Objectif des modèles statistiques: expliquer et formaliser les relations entre les variables

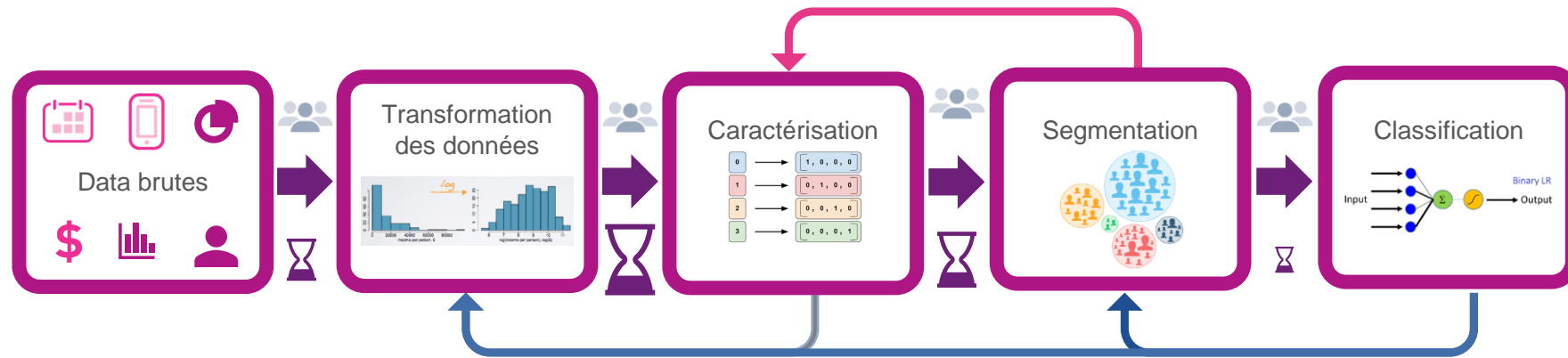
# ML pour les modèles réglementaires

Les régulateurs n'acceptent toujours pas les modèles basés sur du Machine Learning, car ils les considèrent comme des boîtes noires et non cadrés quant à la validation et à la gouvernance.

La **transparence**, l'**explicabilité**, la **réplicabilité** et la **validation** des modèles d'apprentissage automatique deviennent de plus en plus importants pour qu'ils soient acceptés par les régulateurs



# Notre approche pour la modélisation et le déploiement



## Complément

Utilise des modèles standards.

Cherche à améliorer la performance des modèles traditionnels en utilisant l'approche ML comme complément : segmentation, sélection des variable, génération de nouvelles variables ...

## Benchmark

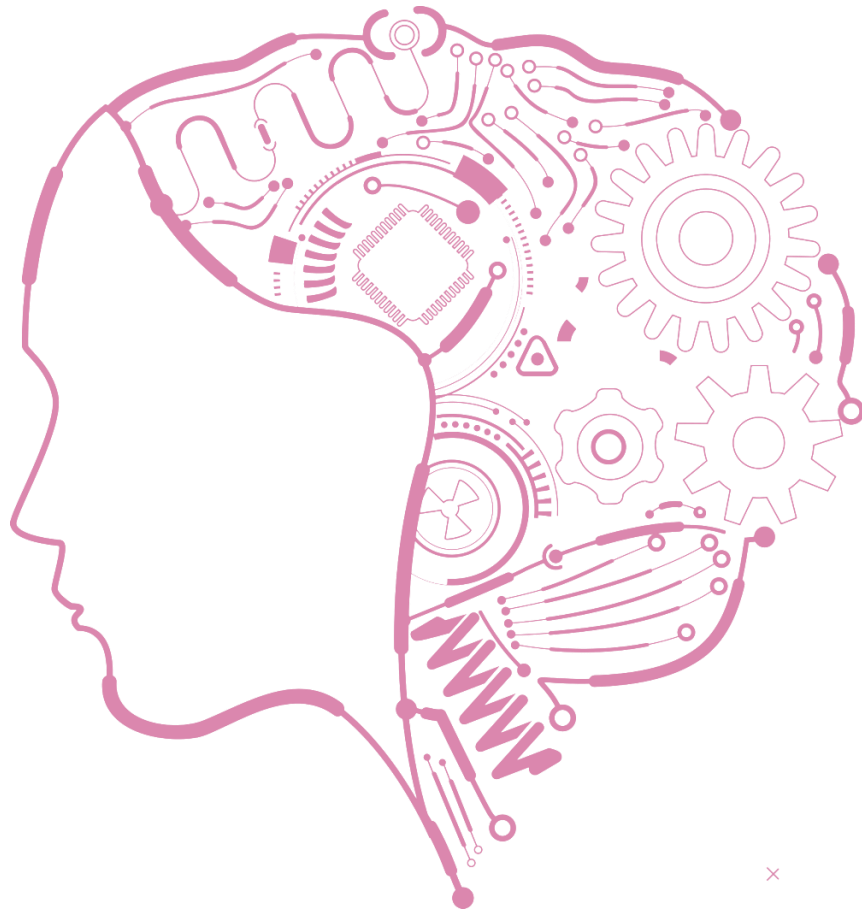
Approche Champion / Challenger.

Fait tourner parallèlement les modèles standards. Le focus est mis sur la performance et la stabilité pour mesurer l'apport potentiel du bénéfice. S'assure de la transparence et de la traçabilité.

## Full ML

Remplacement des modèles traditionnelles par du machine learning

Le focus est mis sur la performance et la mise à jour régulière des modèles



## Exemples de domaines d'application du Machine Learning pour le réglementaire

# Trois domaines d'application du ML

## Données alternatives

- Données Web, transactionnelles, fraude
- Data-driven
- Complètement transparente
- Facilité d'implémentation
- Facile à incorporer dans les modèles traditionnelles

## Segmentation

- Data-driven
- Automatisable
- Personnalisable
- Traitement flexible de différents types de formats de données

## Modèles

- Précis
- Rapide à mettre en œuvre
- Scalable
- Explicable

# Données alternatives

## Web data - performance



90% de “hit rate”

**Données fraîches et plus larges**

90% des entreprises présentes sur le web



+20% Gini pour les scores “traditionnels”

**Les données Web permettent de “booster” le Gini des scores traditionnels**



Gini 43% avec uniquement les données Web

**Réduction du risque sur le portfolio des entreprise**



Réponse <1sec

**Solution temps réel**  
Moins d'une seconde pour le traitement des données et le calcul du score



# Données alternatives

## Web data- méthodologie

Portefeuille de clients



Portefeuille des entreprises

Web Crawler



Informations requises:

- Raison sociale
- Nom d'usage
- Adresse / Ville
- Date de demande
- Comportement observé («G/B Flag»)

En output



Profil Web

Recolte et analyse de l'information web

Modélisation ML



Machine Learning



Score intégrable

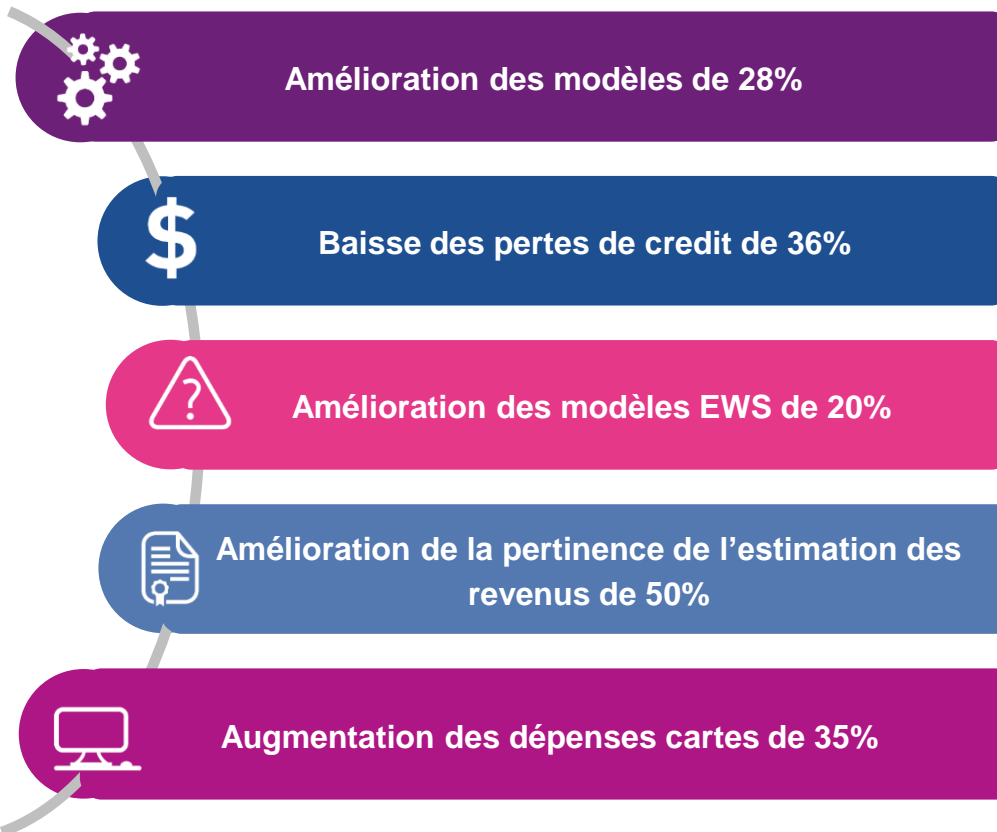
Résultat



Web Data Score

# Données transactionnelles

## Quelques résultats



# Comme utiliser les données transactionnelles

## Usage des données transactionnelles

### 1. Stand alone/ données calculées

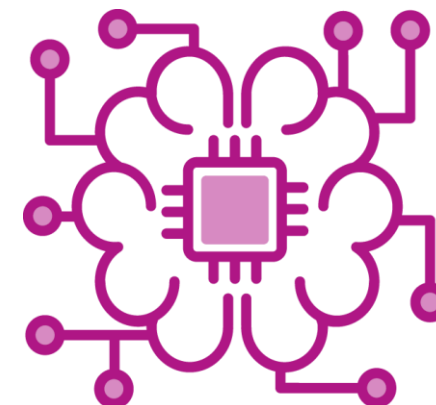
- Données au niveau transaction
- Données segmentées par catégories de transaction
- Agrégées au niveau compte/client
- Combinées avec d'autres données

### 2. Problématiques

- Profil / Risque / fraude / Salaire / réseau

### 3. Algorithmes

- Machine learning
- Modèles traditionnels
- Modèles "hybrides"



# Modèles transactionnels

## Catégorisation des transactions

	<u>Signe</u>
1. Transactions de type Bonus / Remise et Litige / Perte.	-
2. Transactions de type paiement.	
3. Transactions de type intérêt	+
4. Virement (telecom)	+
5. Virement (pas de code marchand)	+
6. Retrait d'espèces "domestiques".	+
7. Retraits d'espèces non domestiques et plusieurs types d'opérations d'achat.	+
8. Transferts de solde et plusieurs types de transactions d'achat.	-
9. Divers types de transactions d'achats	-
10. Divers types de transactions d'achats	+

	<u>Signe</u>
11. Divers types de transactions d'achats.	-
12. Divers types de transactions d'achats.	+
13. Divers types de transactions d'achats.	-
14. Divers types de transactions d'achats.	-
15. Divers types de transactions d'achats.	-
16. Divers types de transactions d'achats.	+
17. Éliminé par le Neural Learner.	
18. Éliminé par le Neural Learner.	
19. Éliminé par le Neural Learner.	
20. Éliminé par le Neural Learner.	

# Transparent neural learner (TNL) pour données transactionnelles

## Classification automatisée des transactions

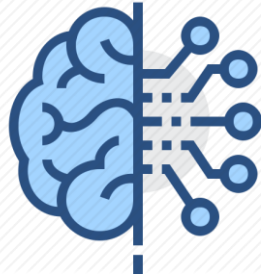
1



- *Modèles de classification concurrents*
- *Explicabilité complète des classes*
- *Possibilité d'ajouter des classes métiers*

## Calcul des "caractéristiques"

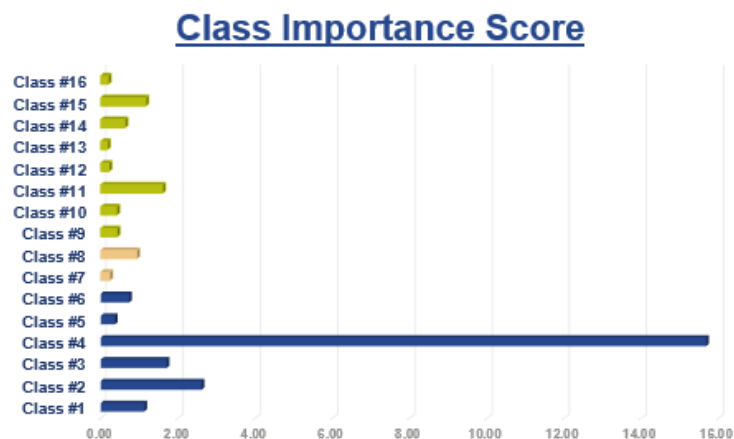
2



- *Prend en charge pour l'ensemble des classes le maximum, la moyenne, le nombre de transactions sur les 3/6/12 périodes*
- *Même signe pour toutes les "caractéristiques" d'une classe*
- *Relation linéaire entre les features et la sortie.*

# Modèles transactionnels

## Interprétabilité



	Max. Amnt. 3M	Total Amnt. 3M	Avg. Amnt. 3M	Max # Trans. 3M	Total # Trans. 3M	Avg. # Trans. 3M	Max. Amnt. 12M	Total Amnt. 12M	Avg. Amnt. 12M	Max # Trans. 12M	Total # Trans. 12M	Avg. # Trans. 12M
Class #16	0.001	-	0.001	-	-	0.001	0.021	-	0.024	0.007	0.005	0.002
Class #15	(0.001)	-	(0.022)	(0.101)	(0.045)	(0.039)	(0.004)	(0.002)	(0.028)	(0.009)	(0.003)	(0.005)
Class #14	(0.030)	(0.005)	(0.015)	-	(0.027)	-	(0.034)	(0.001)	(0.044)	(0.005)	-	(0.003)
Class #13	-	-	(0.001)	(0.016)	(0.004)	(0.002)	-	(0.001)	(0.002)	(0.001)	(0.004)	(0.002)
Class #12	0.001	0.001	0.004	0.001	0.077	0.018	0.009	0.001	0.016	0.002	0.005	0.003
Class #11	(0.061)	(0.001)	(0.001)	(0.319)	(0.008)	(0.002)	(0.035)	(0.117)	(0.175)	(0.045)	(0.003)	(0.226)
Class #10	0.041	0.007	0.018	0.023	0.016	-	0.029	0.014	0.007	-	-	0.008
Class #9	(0.002)	-	-	(0.026)	(0.001)	-	(0.001)	-	(0.014)	(0.001)	(0.004)	-
Class #8	(0.010)	(0.004)	-	-	(0.011)	(0.010)	-	(0.012)	(0.002)	(0.001)	(0.022)	(0.023)
Class #7	0.002	-	0.001	-	0.001	-	-	0.018	0.010	0.023	-	-
Class #6	0.026	0.001	0.004	0.031	0.003	0.006	0.001	0.060	0.001	-	0.001	-
Class #5	0.141	0.109	0.023	-	0.001	-	0.063	0.036	0.180	0.106	0.001	0.047
Class #4	0.722	0.085	0.003	0.004	0.085	0.002	0.140	0.218	0.361	0.001	0.408	0.163
Class #3	0.074	0.004	0.221	0.222	0.003	0.092	0.051	0.254	0.161	-	0.031	0.001
Class #2	(0.216)	-	(0.002)	(0.012)	(0.002)	(0.036)	(0.002)	(0.001)	(0.078)	(0.053)	(0.001)	(0.014)
Class #1	(0.008)	(0.014)	-	(0.029)	-	-	(0.013)	(0.056)	(0.155)	(0.073)	(0.001)	(0.082)

Le score d'importance de classe est défini comme le rapport entre la proportion de score contribué à la population de modélisation par les caractéristiques de classe et la proportion de population de modélisation affectée par la classe.

Rapport complet avec les coefficients des « features » agrégés dérivés en multipliant toutes les transformations linéaires au sein du TNL

# Modèles transactionnels

## Résultats

**+13%**

Meilleure discrimination  
/ benchmark modèle  
comportemental

*(Augmentation de Gini de 70.2% à  
79.2%)  
Gros éch. – 499 511 obs.*

**+15%**

Meilleure discrimination  
/ benchmark modèle  
Cross selling

*(Augmentation de Gini de 44% à  
50.6%)  
Petit éch. – 3 044 obs.*

**-80%**

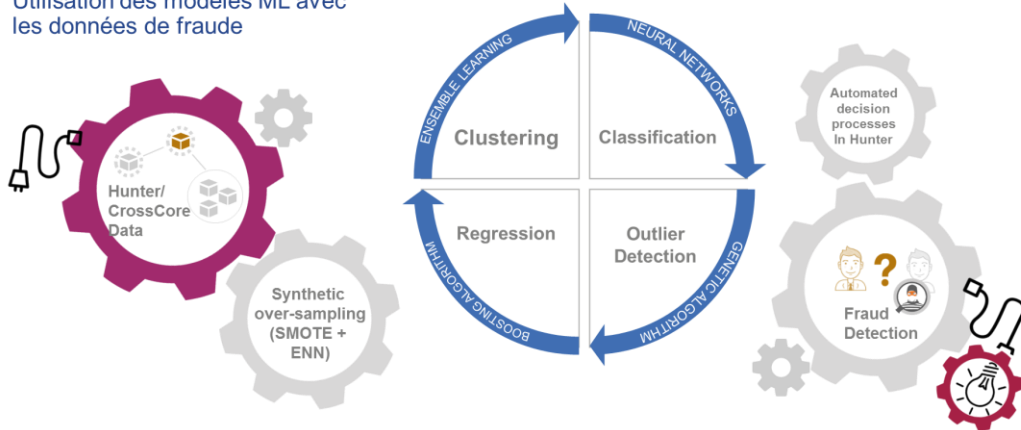
Sur le temps de  
préparation nécessaire  
à un modèle  
transactionnel

*(Baisse de 2.5 semaines à 3 jours)*

# Données alternatives

## Modèle fraude

Utilisation des modèles ML avec les données de fraude



### Bénéfices

#### Financiers

#### Telecom Portfolio

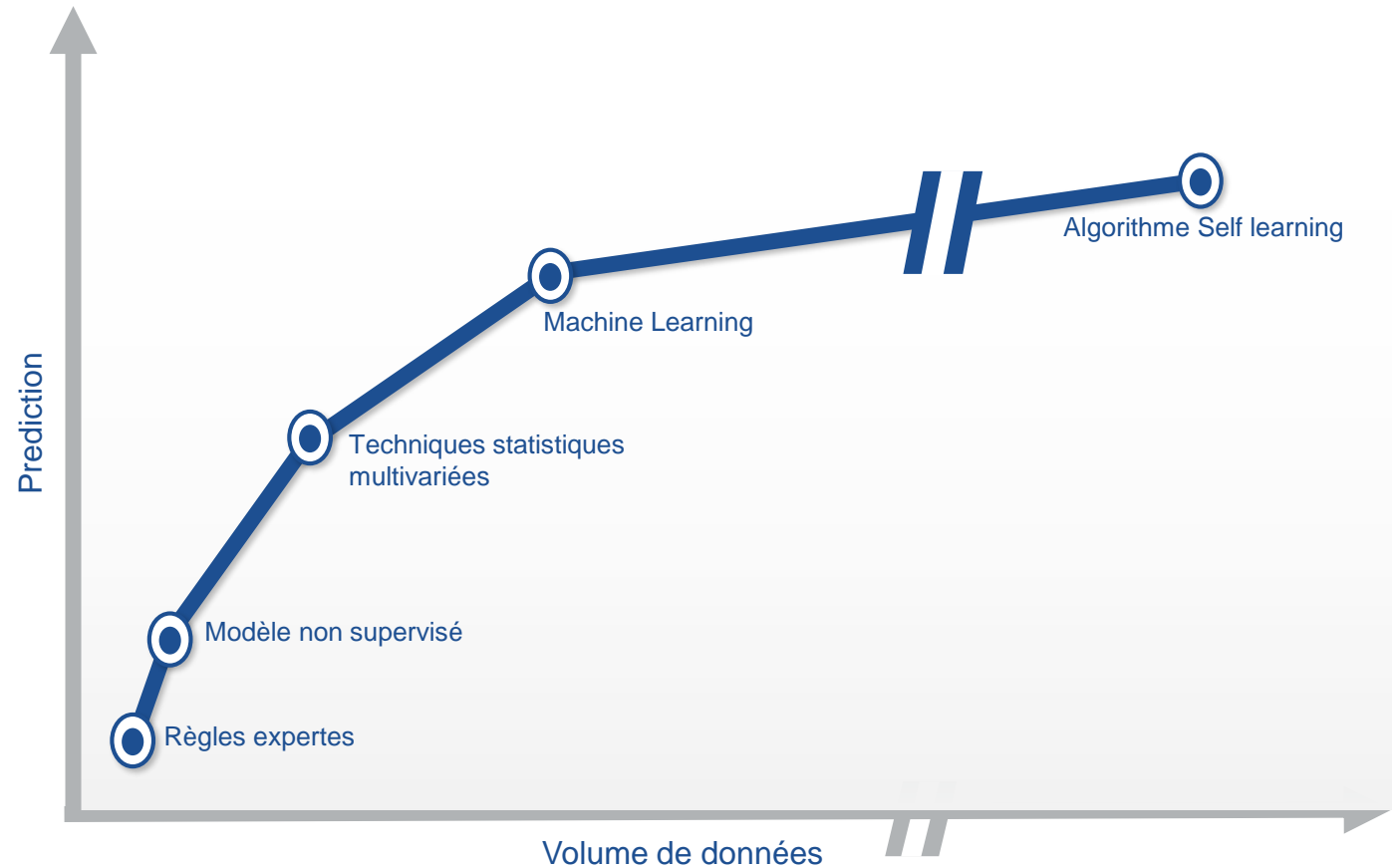
- + **12%** amélioration de la fraude identifiée
- + **39MM** de revenu additionnel

#### Financial Services

- + **19%** amélioration de la fraude identifiée
- + **2MM** de revenu additionnel

#### Opérationnels

- Process automatique de décision
- Plus de précision dans la stratégie anti fraude



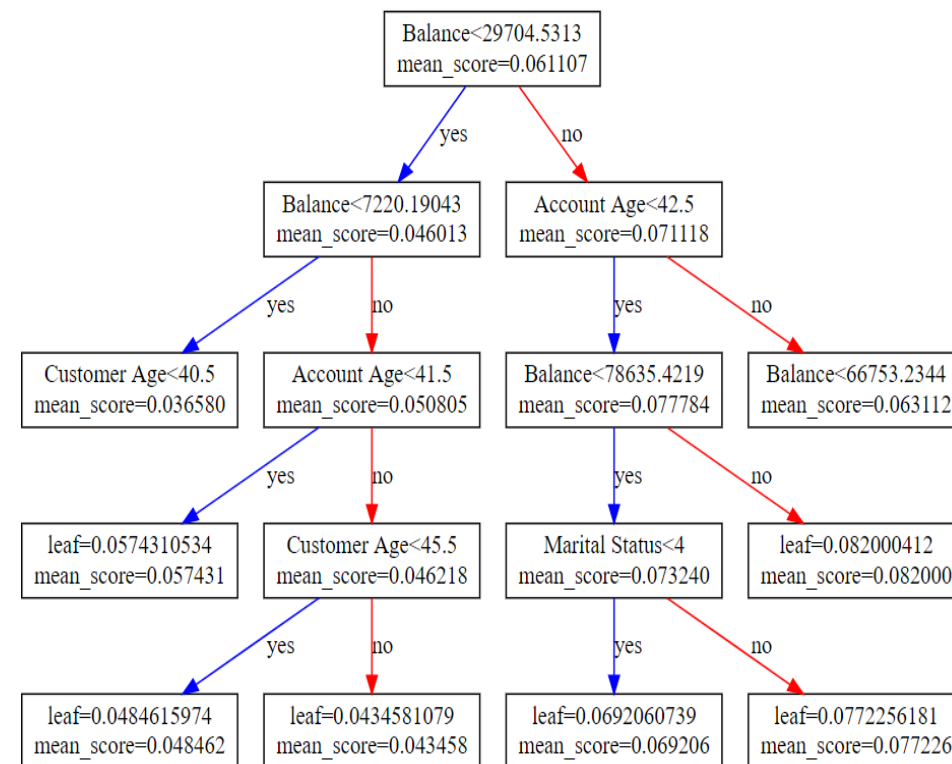


# Modèle de ML pour la segmentation

## IFRS9

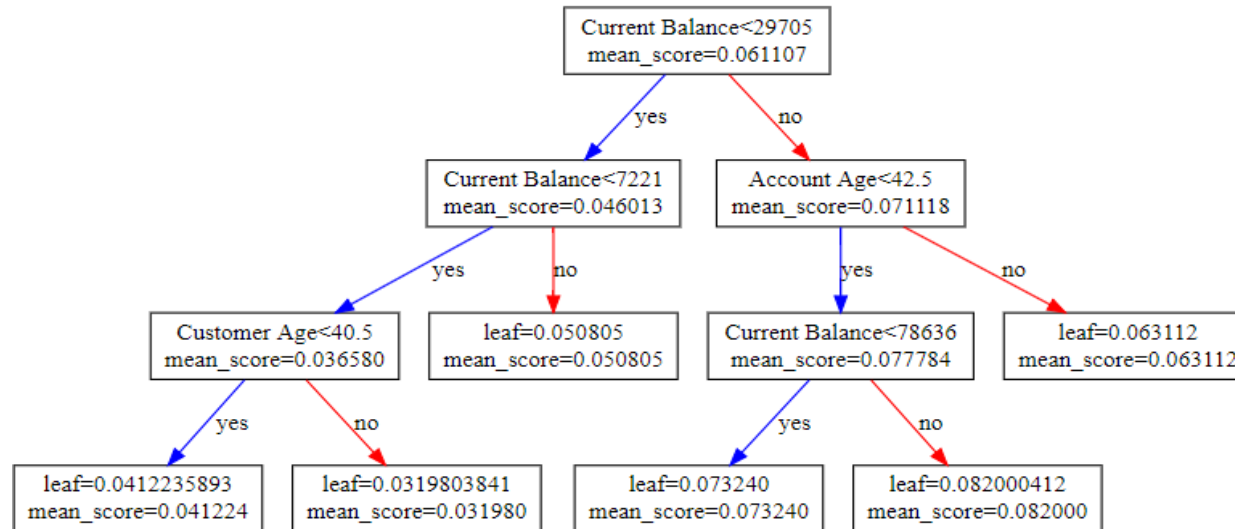
Les segments produits doivent répondre à certaines exigences pour être considérés comme «bons»:

- ❑ **Considérations de stabilité** - Les segments doivent être suffisamment grands (taille préférable entre 5% et 30% de la population)
- ❑ **Considérations statistiques** - Les segments doivent différer suffisamment les uns des autres en termes de LGD moyenne de la population qu'ils couvrent.
- ❑ **Considérations Business** – Nous devons nous assurer que les découpages ne contredisent pas les attendus métier.



# Modèles Segmentation

## Amélioration de la pertinence



	Erreur 12M	Erreur 24M	Erreur 36M	Erreur 48M
XGBoost	-4.23%	-5.52%	-3.95%	-1.74%
Chaid	-4.43%	-6.03%	-5.08%	-3.63%

# ML pour la prediction du risque

## Evaluation de modèles alternatifs

Sample parameters	Techniques (single model)
<ul style="list-style-type: none"><li>• Generic bureau data samples<ul style="list-style-type: none"><li>• Auto</li><li>• Bankcard</li></ul></li><li>• 90+ DPD performance flag</li><li>• 24-month outcome period</li></ul>	<ul style="list-style-type: none"><li>• Logistic regression (LR)</li><li>• Neural network (NN)</li><li>• Random forest (RF)</li><li>• Support vector machines (SVM)</li><li>• Extreme gradient boosting (XGB)</li></ul>

Approx. 5% de lift sur OOT

Validation Gini					
	LR	NN	RF	SVM	XGB
Auto	71.34	73.87	73.21	73.98	74.80
Bankcard	69.30	72.11	72.31	72.22	73.18

# ML pour la prediction du risque

## Modèles PD

### Projet 1: Suède Portfolio: Carte de credit

#### Chiffre-clés du jeu de données:

- 112k enregistrements, 100 variables

#### Indice de Gini:

- Scorecard: 0.752
- Modèles XGBoost:
  - Toutes variables, sans contraintes de monotonie: 0.883
  - Toutes variables, contraintes de monotonie: 0.873
  - Variables du score uniquement, contraintes de monotonie: 0.866

#### Rapidité de construction:

- Score classique: 7 jours
- Modèle XGBoost (toutes variables): 1-2 jours

### Project 2: Norvège Portfolio: Carte de credit

#### Chiffre-clés du jeu de données:

- 265k enregistrements, 375 variables

#### Indice de Gini:

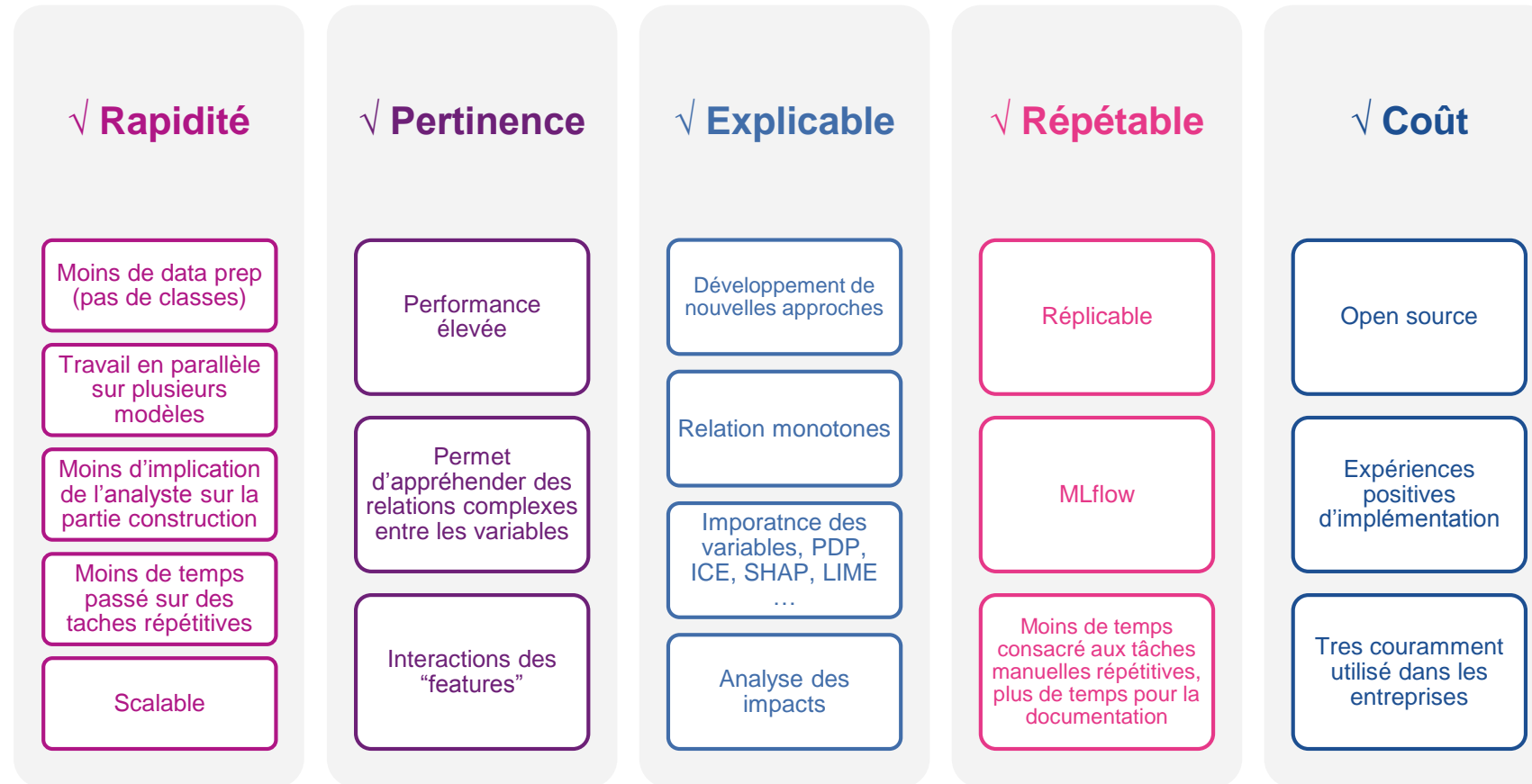
- Scorecard: 0.786
- Modèles XGBoost:
  - Toutes variables, sans contraintes de monotonie: 0.891
  - Toutes variables, contraintes de monotonie: 0.911
  - Variables du score uniquement, contraintes de monotonie : 0.889

#### Rapidité de construction:

- Score classique: 10 jours
- Modèle XGBoost (toutes variables): 1-2 jours

# Modèles de Gradient boosting

## Avantages

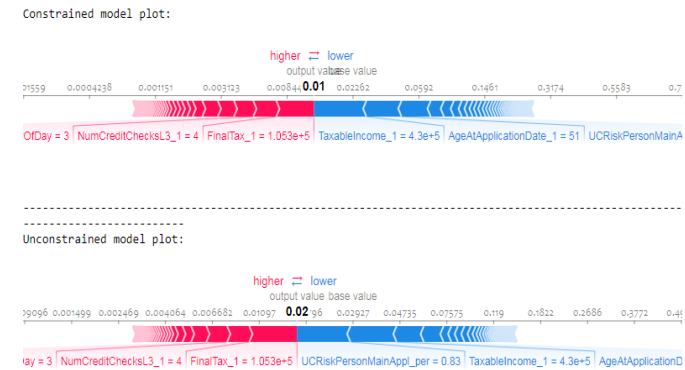
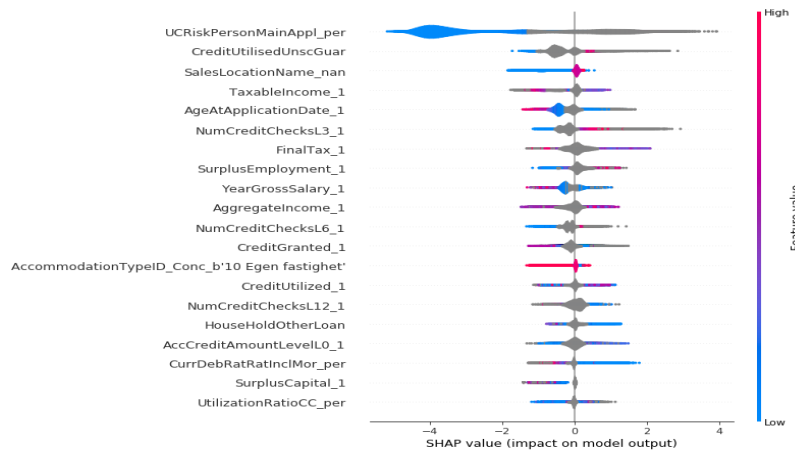
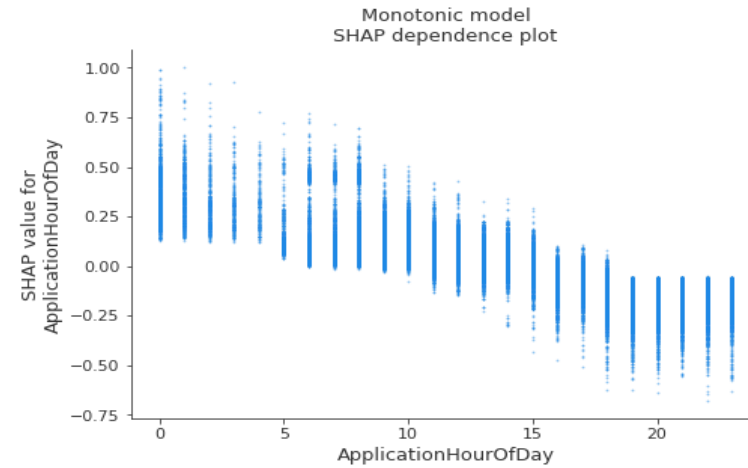


# Modèles de Gradient boosting

## SHAP – Outil pour l'interprétabilité

Constrained model table:  
Contributions of the Top-25 variables to prediction: 1.1421680441753017

	value	SHAP	abs(SHAP)
TaxableIncome_1	430000.00	-0.796813	0.796813
FinalTax_1	105300.00	0.571930	0.571930
AgeAtApplicationDate_1	51.00	-0.533548	0.533548
UCRiskPersonMainAppl_per	0.83	-0.443313	0.443313
NumCreditChecksL3_1	4.00	0.349832	0.349832
ApplicationHourOfDay	3.00	0.293296	0.293296
CreditUtilisedUnscGuar	359166.00	0.245020	0.245020
YearGrossSalary_1	264000.00	0.234493	0.234493
CreditGranted_1	523166.00	-0.200734	0.200734
CreditUtilisedInstalment_1	20000.00	-0.151029	0.151029
NumCreditChecksL12_1	5.00	0.146806	0.146806
NumOfCredits_1	9.00	-0.119769	0.119769
ApplicationDayOfWeek	6.00	-0.088433	0.088433
ULBehaviourPD	NaN	0.083023	0.083023
SalesLocationName_nan	1.00	0.081720	0.081720
CreditUtilized_1	389796.00	-0.075891	0.075891
LimitCreditCards_1	0.00	0.069056	0.069056
AccommodationTypeID_Conc_b'7 Hyresr'xe4tt'	1.00	-0.067559	0.067559
NumCreditCards_1	0.00	0.067519	0.067519
AvgUtilRatioCCs_1_per	NaN	0.066347	0.066347
HouseHoldOtherLoan	389796.00	-0.065782	0.065782
EmploymentTypeID_1_Conc_b'13 Vikarie/Projektanstxe4llid'	1.00	-0.064664	0.064664
CCAppL6M	NaN	-0.052865	0.052865
AccCreditAmountLevelL0_1	10630.00	-0.050864	0.050864
ULAppl12M	NaN	0.049089	0.049089



# Une méthode pour ouvrir la Boite Noire

