

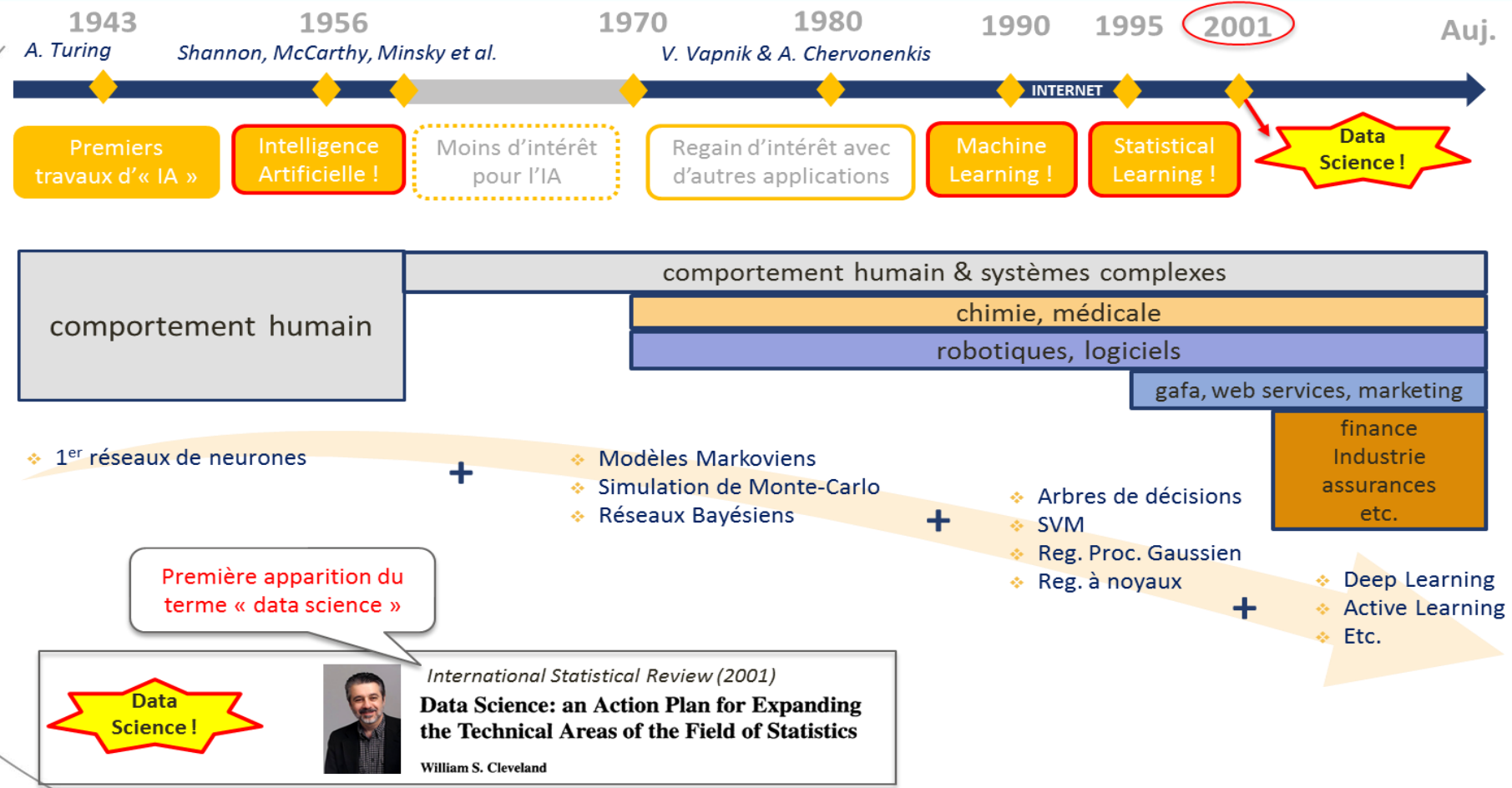
The logo consists of three white triangles of increasing size, stacked and pointing to the right, creating a sense of upward and forward movement.

**INSTITUT DES
ACTUAIRES**

La Data Science et son apport pour mieux détecter les besoins de prévention des assurés en santé

**Alexandra BARRAL, Romain GAUCHON, Nabil RACHDI
ACTUARIS**

Les origines



Quelle(s) différence(s) ?

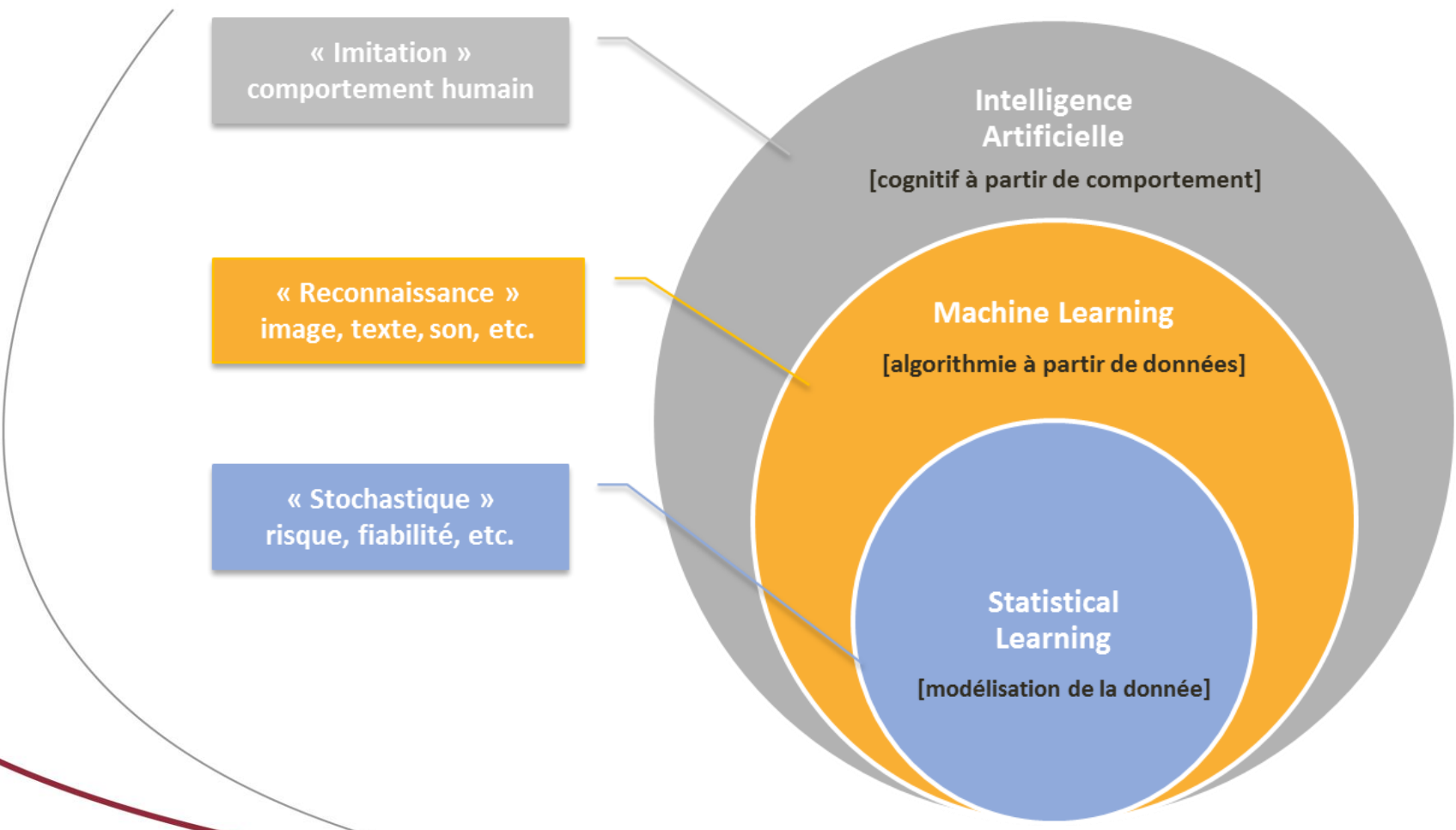
Principalement culturelle !

► **Glossaire** de R. Tibshirani (Pr. à Stanford)
co-auteur de « *The Elements of Statistical Learning* », 2008

- MACHINE Learning
= approche algorithmique
- STATISTICAL Learning
= approche par modélisation de la donnée

Glossary	
Machine learning	Statistics
network, graphs	model
weights	parameters
learning	fitting
generalization	test set performance
supervised learning	regression/classification
unsupervised learning	density estimation, clustering
large grant = \$1,000,000	large grant= \$50,000
nice place to have a meeting: Snowbird, Utah, French Alps	nice place to have a meeting: Las Vegas in August

En résumé



Exemple : Les GLM revisitées

Vue « classique »

$Y_i =$ fréquences
 $X_i =$ contrats

Vue « Statistical Learning »

1/ Modélisation des données

- $Y \sim \mathcal{P}(\lambda)$
- $\log(\mathbb{E}(Y | \mathbf{X} = \mathbf{x})) = \sum_{i=1}^d \beta_i x_i = \mathbf{x}^T \boldsymbol{\beta}$

Soit

$$(Y | \mathbf{X} = \mathbf{x}) \approx \mathcal{P}[e^{\mathbf{x}^T \boldsymbol{\beta}}]$$

2/ Maximum de Vraisemblance des données observées

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\text{Argmax}} \prod_{i=1}^n \mathcal{P}[e^{\mathbf{x}_i^T \boldsymbol{\beta}}](y_i)$$

Maximum de Vraisemblance

Exemple : Les GLM revisitées

Vue « classique »

$Y_i =$ fréquences
 $X_i =$ contrats

Vue « Statistical Learning »

1/ Modélisation des données

- $Y \sim \mathcal{P}(\lambda)$
- $\log(\mathbb{E}(Y|\mathbf{X} = \mathbf{x})) = \sum_{i=1}^d \beta_i x_i = \mathbf{x}^T \boldsymbol{\beta}$

Soit

$$(Y|\mathbf{X} = \mathbf{x}) \approx \mathcal{P}[e^{\mathbf{x}^T \boldsymbol{\beta}}]$$

2/ Maximum de Vraisemblance des données observées

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\text{Argmax}} \prod_{i=1}^n \mathcal{P}[e^{\mathbf{x}_i^T \boldsymbol{\beta}}](y_i)$$

Maximum de Vraisemblance

1/ Quantité d'intérêt cible (« target »)

loi jointe de $(\mathbf{X}, Y) \rightarrow f^*(\mathbf{x}, y)$

2/ Modélisation de la cible (« model »)

modèle : $f_{\mathbf{X}}(\mathbf{x}) \mathcal{P}[e^{\mathbf{x}^T \boldsymbol{\beta}}](y) \rightarrow f_{\boldsymbol{\beta}}(\mathbf{x}, y)$

3/ Fonction de perte associée à la cible (« loss function »)

métrique considérée \rightarrow Entropie Relative
aussi appelée divergence de Kullback-Leibler $KL(f_1, f_2) = \int \log\left(\frac{f_1}{f_2}\right) f_1$

4/ Minimisation Empirique de la fonction de perte (« ERM »)

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \underset{\boldsymbol{\beta}}{\text{Argmin}} KL(\hat{f}^*, f_{\boldsymbol{\beta}}) \\ &= \underset{\boldsymbol{\beta}}{\text{Argmax}} \sum_{i=1}^n \log\left(\mathcal{P}[e^{\mathbf{x}_i^T \boldsymbol{\beta}}](y_i)\right) \end{aligned}$$

Exemple : Les GLM revisitées

Vue « classique »

$Y_i =$ fréquences
 $\mathbf{X}_i =$ contrats

Vue « Statistical Learning »

1/ Modélisation des données

- $Y \sim \mathcal{P}(\lambda)$
- $\log(\mathbb{E}(Y|\mathbf{X} = \mathbf{x})) = \sum_{i=1}^d \beta_i x_i = \mathbf{x}^T \boldsymbol{\beta}$

Soit

$$(Y|\mathbf{X} = \mathbf{x}) \approx \mathcal{P}[e^{\mathbf{x}^T \boldsymbol{\beta}}]$$

2/ Maximum de Vraisemblance des données observées

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\text{Argmax}} \prod_{i=1}^n \mathcal{P}[e^{\mathbf{x}_i^T \boldsymbol{\beta}}](y_i)$$

Maximum de Vraisemblance

1/ Quantité d'intérêt cible (« target »)

loi jointe de $(\mathbf{X}, Y) \rightarrow f^*(\mathbf{x}, y)$

2/ Modélisation de la cible (« model »)

modèle : $f_{\mathbf{X}}(\mathbf{x}) \mathcal{P}[e^{\mathbf{x}^T \boldsymbol{\beta}}](y) \rightarrow f_{\boldsymbol{\beta}}(\mathbf{x}, y)$

3/ Fonction de perte associée à la cible (« loss function »)

métrique considérée \rightarrow Entropie Relative
aussi appelée divergence de Kullback-Leibler $KL(f_1, f_2) = \int \log\left(\frac{f_1}{f_2}\right) f_1$

4/ Minimisation Empirique de la fonction de perte (« ERM »)

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \underset{\boldsymbol{\beta}}{\text{Argmin}} KL(\hat{f}^*, f_{\boldsymbol{\beta}}) \\ &= \underset{\boldsymbol{\beta}}{\text{Argmax}} \sum_{i=1}^n \log\left(\mathcal{P}[e^{\mathbf{x}_i^T \boldsymbol{\beta}}](y_i)\right) \end{aligned}$$

Minimisation de l'Entropie Relative



Exemple : Les GLM revisitées

Vue « classique »

$Y_i =$ fréquences
 $X_i =$ contrats

Vue « Statistical Learning »

1/ Modélisation des données

- $Y \sim \mathcal{P}(\lambda)$
- $\log(\mathbb{E}(Y|\mathbf{X} = \mathbf{x})) = \sum_{i=1}^d \beta_i x_i = \mathbf{x}^T \boldsymbol{\beta}$

Soit

On peut **changer** de « target »
et/ou de « model »
et/ou de « loss function » !

1/ Quantité d'intérêt cible (« target »)

2/ Modélisation de la cible (« model »)

3/ Fonction de perte associée à la cible (« loss function »)

2/ Maximum de Vraisemblance des données observées

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\text{Argmax}} \prod_{i=1}^n \mathcal{P}[e^{\mathbf{x}_i^T \boldsymbol{\beta}}](y_i)$$

4/ Minimisation Empirique de la fonction de perte (« ERM »)

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\text{Argmin}} \text{loss}(\widehat{\text{target}}, \text{target}_{\boldsymbol{\beta}})$$

Maximum de Vraisemblance

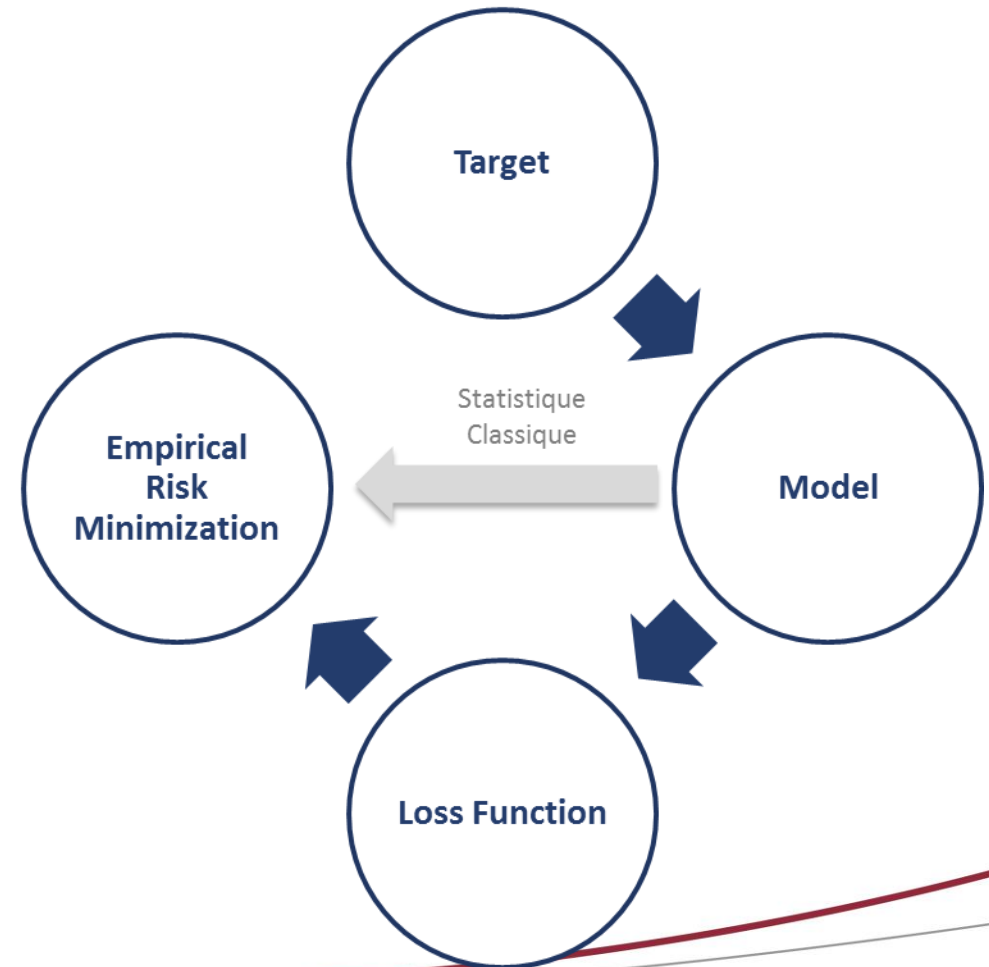
\neq

Autre(s) algorithm(e)s

En résumé

- Le Statistical Learning **étend les méthodes de prédictions statistiques classiques** à un cadre d'analyse beaucoup plus large
- donnant **plus de flexibilité** à la modélisation
- permettant l'utilisation d'une **grande variété d'algorithmes d'apprentissage**
- Il est important de savoir ce qui se cache derrière les algorithmes
(quelle target ? quelle métrique ?)

Statistical Learning process



		CONTEXTE			
		Supervisé $(X_i, Y_i)_{i=1:n}$	Non-supervisé $(X_i)_{i=1:n}$	Semi-supervisé	mix
TACHE	Régression	<ul style="list-style-type: none"> ▪ Régression linéaire (Lasso, Ridge, etc.) ▪ Régression non-linéaire (noyaux, etc.) ▪ Arbres (CART, Random Forest, etc.) ▪ kNN ▪ Réseaux de Neurones, SVR ▪ GLM ▪ Etc. 	X		<ul style="list-style-type: none"> • SVM • Graphs
	Classification	<ul style="list-style-type: none"> ▪ Régression logistique ▪ Réseaux Bayésiens ▪ Arbres ▪ Réseaux de Neurones, SVM ▪ Etc. 	X		<ul style="list-style-type: none"> • SVM • Graphs
	Clustering	X		<ul style="list-style-type: none"> • K-Means • Densités • Cartes Kohonen 	X
	Réduction Dim.	<ul style="list-style-type: none"> • ANOVA (Analyse de la variance) • Algorithmes génétiques • Etc. 	<ul style="list-style-type: none"> • Analyse des corrélations • ACP 		X

méthodes traditionnelles en assurance



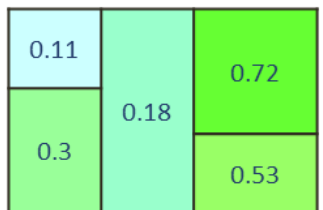
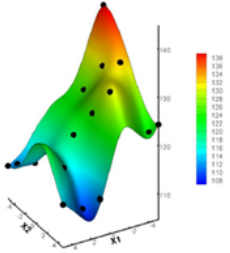
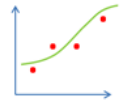

Pour un couple (Tâche, Contexte) donné, l'algorithmie dépendra de la **NATURE des données**

- ▶ **Quantitatives** (discrètes, continues)
- ▶ **Qualitatives** (ordinales, nominales)



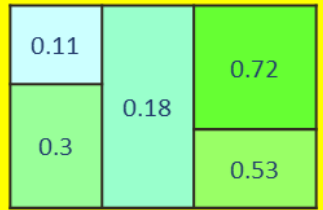
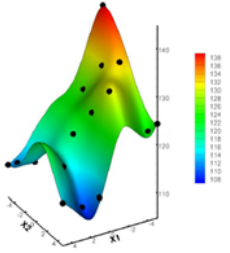
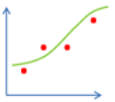

ALGO = TACHE + CONTEXTE + NATURE DONNEES

ALGO = TACHE + CONTEXTE + NATURE DONNEES

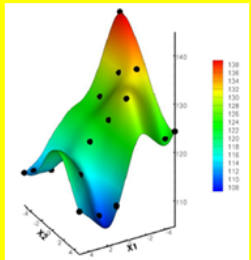
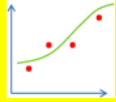
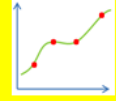
Nature des Inputs	Nature des algorithmes		
	Familles de méthodes	Méthodes	Algorithmes
<p>Catégorielles et/ou Continues</p>	<p>« Ensemblistes »</p> <p>Approximation en <i>partitionnant</i> de l'espace des entrées</p> 	<ul style="list-style-type: none"> • Arbres • Forêts 	<ul style="list-style-type: none"> • CART • AdaBoost • RandomForest • XGBoost • Gradient Boosting
<p>Continues et « peu » Catégorielles</p>	<p>« Fonctionnelles »</p> <p>Approximation par <i>représentation fonctionnelle</i></p> 	<ul style="list-style-type: none"> • Polynômes • Radial Basis Functions • Processus Gaussien 	<p>Extrapolation</p> <ul style="list-style-type: none"> • Ordinary Least Squares • Polynomial Chaos • Tensorization • MARS  <p>Interpolation</p> <ul style="list-style-type: none"> • Kriging (process. gaussien) • Sparse Kriging 

ALGO = TACHE + CONTEXTE + NATURE DONNÉES

Les données assurantielles sont très souvent un mix catégorielles/continues. D'où l'exploration croissante des méthodes ensemblistes (arbres)

Nature des Inputs	Famille méthodes		Algorithmes
	« Ensemblistes »	« Fonctionnelles »	
Catégorielles et/ou Continues	<p>Approximation en <i>partitionnant</i> de l'espace des entrées</p> 	<ul style="list-style-type: none"> • Arbres • Forêts 	<ul style="list-style-type: none"> • CART • AdaBoost • RandomForest • XGBoost • Gradient Boosting
Continues et « peu » Catégorielles	<p>Approximation par <i>représentation fonctionnelle</i></p> 	<ul style="list-style-type: none"> • Polynômes • Radial Basis Functions • Processus Gaussien 	<p>Extrapolation</p> <ul style="list-style-type: none"> • Ordinary Least Squares • Polynomial Chaos • Tensorization • MARS  <p>Interpolation</p> <ul style="list-style-type: none"> • Kriging (process. gaussien) • Sparse Kriging 

$ALGO = TACHE + CONTEXTE + NATURE\ DONNEES$

Nature des Inputs	Nature des algorithmes								
	Familles de méthodes	Méthodes	Algorithmes						
<p style="color: red; text-align: center;">Catégorielles et/ou Continues</p>	<p>« Ensemblistes »</p> <p><i>Approximation en part l'espace des en</i></p> <table border="1" style="margin: 10px auto;"> <tr> <td style="background-color: #e0f0ff;">0.11</td> <td style="background-color: #e0ffe0;">0.18</td> <td style="background-color: #e0ffe0;">0.72</td> </tr> <tr> <td style="background-color: #e0ffe0;">0.3</td> <td></td> <td style="background-color: #e0ffe0;">0.53</td> </tr> </table>	0.11	0.18	0.72	0.3		0.53	<p style="text-align: center;">L'approche « fonctionnelle » peut se révéler plus performante dans les cas avec peu de variables catégorielles</p>	<p>Gradient Boosting</p>
0.11	0.18	0.72							
0.3		0.53							
<p style="color: red; text-align: center;">Continues et « peu » Catégorielles</p>	<p>« Fonctionnelles »</p> <p><i>Approximation par représentation fonctionnelle</i></p> 	<ul style="list-style-type: none"> • Polynômes • Radial Basis Functions • Processus Gaussien 	<p>Extrapolation</p> <ul style="list-style-type: none"> • Ordinary Least Squares • Polynomial Chaos • Tensorization • MARS  <p>Interpolation</p> <ul style="list-style-type: none"> • Kriging (process. gaussien) • Sparse Kriging 						

Plus généralement : NO FREE LUNCH !

■ « NO FREE LUNCH » THEOREM

Aucun algorithme ne domine tous les autres pour chaque type de problème et de caractéristiques attendues

Caractéristiques du prédicteur	Algorithmes			
	GLM/Lasso	Kriging/Chaos	Deep Learning	Random Forest
Données « mixtes »	=	-	-	+
Automatique	+	+	-	+
Interprétation	+	=	-	=
Prédictivité	=	+	+	+
Facteurs d'importances	+	+	=	=
Facteurs Interactifs	-	+	+	-



Méthodologie générale

Pourquoi une méthodologie ?



Datasets
XLS

Input_1	Input_2	Input_3	Input_4	Output_1	Output_2	Output_3	IBM
0.398	0.000	0.000	0.000	1.623	0.022	-0.127	0.047
0.447	0.000	0.000	0.000	1.936	0.023	-0.143	0.054
0.477	0.000	0.000	0.000	2.123	0.024	-0.153	0.058
0.497	0.000	0.000	0.000	2.247	0.025	-0.159	0.061
0.517	0.000	0.000	0.000	2.369	0.025	-0.165	0.063
0.547	0.000	0.000	0.000	2.551	0.027	-0.173	0.068
0.597	0.000	0.000	0.000	2.863	0.029	-0.187	0.074
0.398	-1.949	0.418	0.127	1.952	0.022	0.003	0.074
0.447	-1.949	0.418	0.127	2.262	0.023	-0.013	0.049
0.477	-1.949	0.418	0.127	2.445	0.024	-0.022	0.056
0.497	-1.949	0.418	0.127	2.566	0.025	-0.028	0.060
0.517	-1.949	0.418	0.127	2.688	0.026	-0.033	0.063
0.547	-1.949	0.418	0.127	2.874	0.027	-0.042	0.065
0.596	-1.949	0.418	0.127	3.202	0.030	-0.055	0.070
0.397	-1.506	0.766	-1.038	2.347	0.023	-0.003	0.060
0.447	-1.506	0.766	-1.038	2.665	0.024	-0.018	0.076
0.477	-1.506	0.766	-1.038	2.846	0.025	-0.017	0.041
...	0.047
0.477	-1.506	0.766	-1.038	2.846	0.025	-0.027	19.025
...	0.052
0.477	-1.506	0.766	-1.038	2.846	0.025	-0.027	19.025
...	0.052
0.477	-1.506	0.766	-1.038	2.846	0.025	-0.027	19.025
...	0.052
0.477	-1.506	0.766	-1.038	2.846	0.025	-0.027	19.025
...	0.052



Algorithmes



Clouds containing the following library names:

- pandas
- statsmodel
- glmnet
- scipy
- e1071
- randomForest
- scikit-learn
- gbm
- numpy
- neuralnet



Visualisation



Clouds containing the following library names:

- ggplot2
- matplotlib
- seaborn
- googleVis
- rAmCharts
- rgl
- bokeh

Vers une approche méthodologique

Pourquoi une méthodologie ?

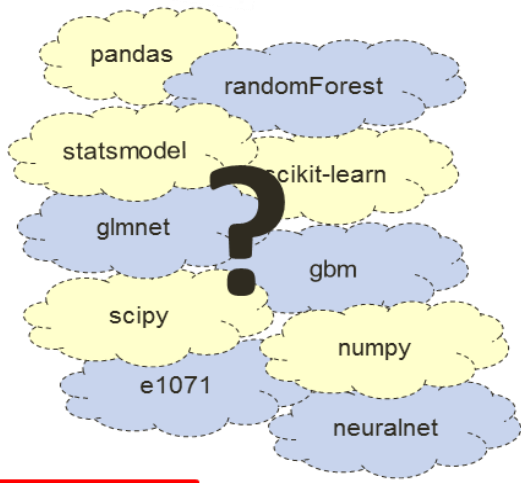


Datasets

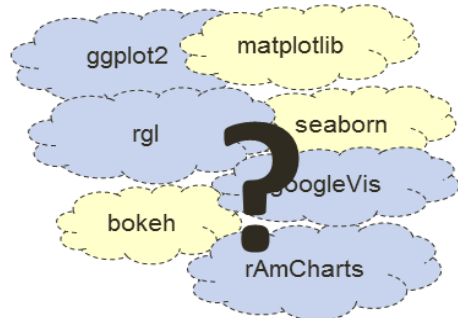


Input_1	Input_2	Input_3	Input_4	Output_1	Output_2	Output_3	BM
0.398	0.000	0.000	0.000	1.623	0.022	-0.127	0.047
0.447	0.000	0.000	0.000	1.936	0.023	-0.143	0.054
0.477	0.000	0.000	0.000	2.123	0.024	-0.153	0.058
0.497	0.000	0.000	0.000	2.247	0.025	-0.159	0.061
0.517	0.000	0.000	0.000	2.369	0.025	-0.165	0.063
0.547	0.000	0.000	0.000	2.551	0.027	-0.173	0.069
0.597	0.000	0.000	0.000	2.863	0.029	-0.187	0.074
0.398	-1.949	0.418	0.127	1.952	0.022	-0.003	0.049
0.447	-1.949	0.418	0.127	2.262	0.023	-0.013	0.049
0.477	-1.949	0.418	0.127	2.445	0.024	-0.022	0.056
0.497	-1.949	0.418	0.127	2.666	0.025	-0.028	0.060
0.517	-1.949	0.418	0.127	2.898	0.026	-0.033	0.063
0.547	-1.949	0.418	0.127	3.174	0.027	-0.042	0.065
0.596	-1.949	0.418	0.127	3.502	0.030	-0.055	0.070
0.397	-1.506	0.766	-1.038	2.347	0.023	-0.003	0.076
0.447	-1.506	0.766	-1.038	2.665	0.024	-0.018	0.041
0.477	-1.506	0.766	-1.038	2.846	0.025	-0.027	0.047

Algorithmes



Visualisation

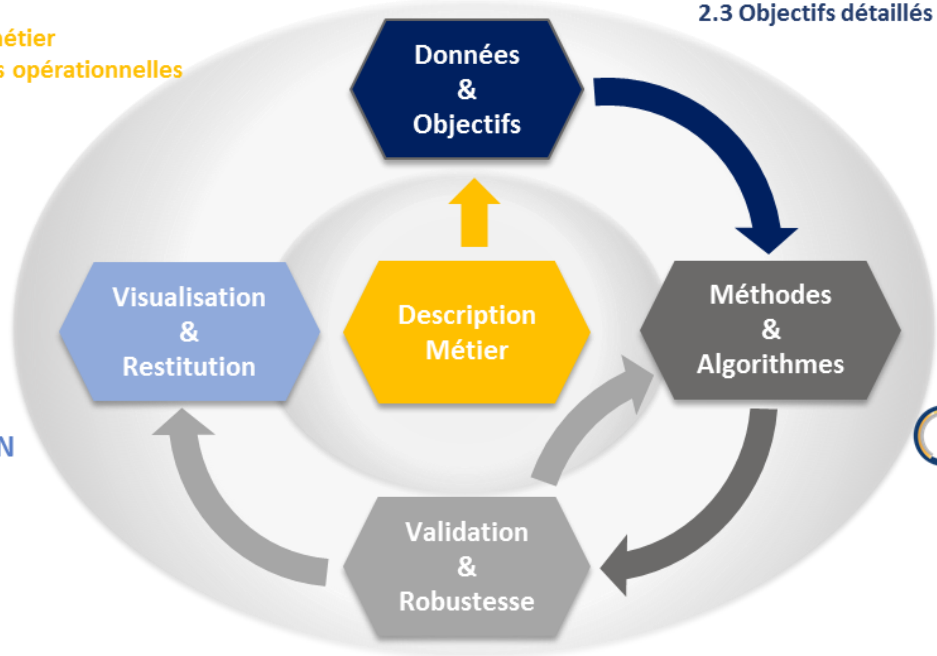


- Comment structurer la base de données ?
- pour quel(s) algorithme(s) ?
- pour quelle fin ?

Besoin d'une méthodologie structurante

- 1. DESCRIPTION METIER**
- 1.1 Contexte
 - 1.2 Objectifs métier
 - 1.3 Contraintes opérationnelles

- 2. DONNEES ET OBJECTIFS**
- 2.1 Sources des données
 - 2.2 Description de la base de données
 - 2.3 Objectifs détaillés

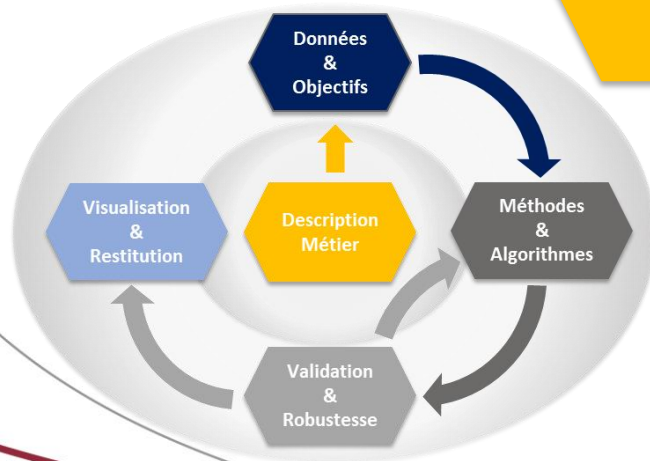


- 3. METHODES ET ALGORITHMES**
- 3.1 Méthodes utilisées
 - 3.2 Algorithmes

- 4. VALIDATION ET ROBUSTESSE**
- 4.1 Validation intrinsèque et par recoupement
 - 4.2 Robustesse des algorithmes
 - 4.3 Feedback avec les algorithmes

- 5. VISUALISATION & RESTITUTION**
- 5.1 Communication des résultats
 - 5.2 Visualisation
 - 5.3 Intégration et Déploiement

DESCRIPTION METIER



Définition de la prévention



Prévention Primaire

- Agir en amont du risque pour empêcher qu'il ne se réalise
- Définition naturelle de la prévention



Prévention Secondaire

- Détection précoce du risque
- Repérer la réalisation du risque le plus tôt possible afin de le traiter



Prévention Tertiaire

- Eviter l'aggravation du sinistre - curatif
- Prévenir les rechutes

Enjeux et problématiques



Assurés

Services personnalisés

Au-delà du paiement des soins

Meilleure connaissance des assurés



Commercial

Amélioration de l'image de l'entreprise

Créer une relation de confiance

Distinction des concurrents



Risques

Projection de l'amélioration de notre risque (moins de sinistres)

Baisse des provisions

Les prérequis à identifier

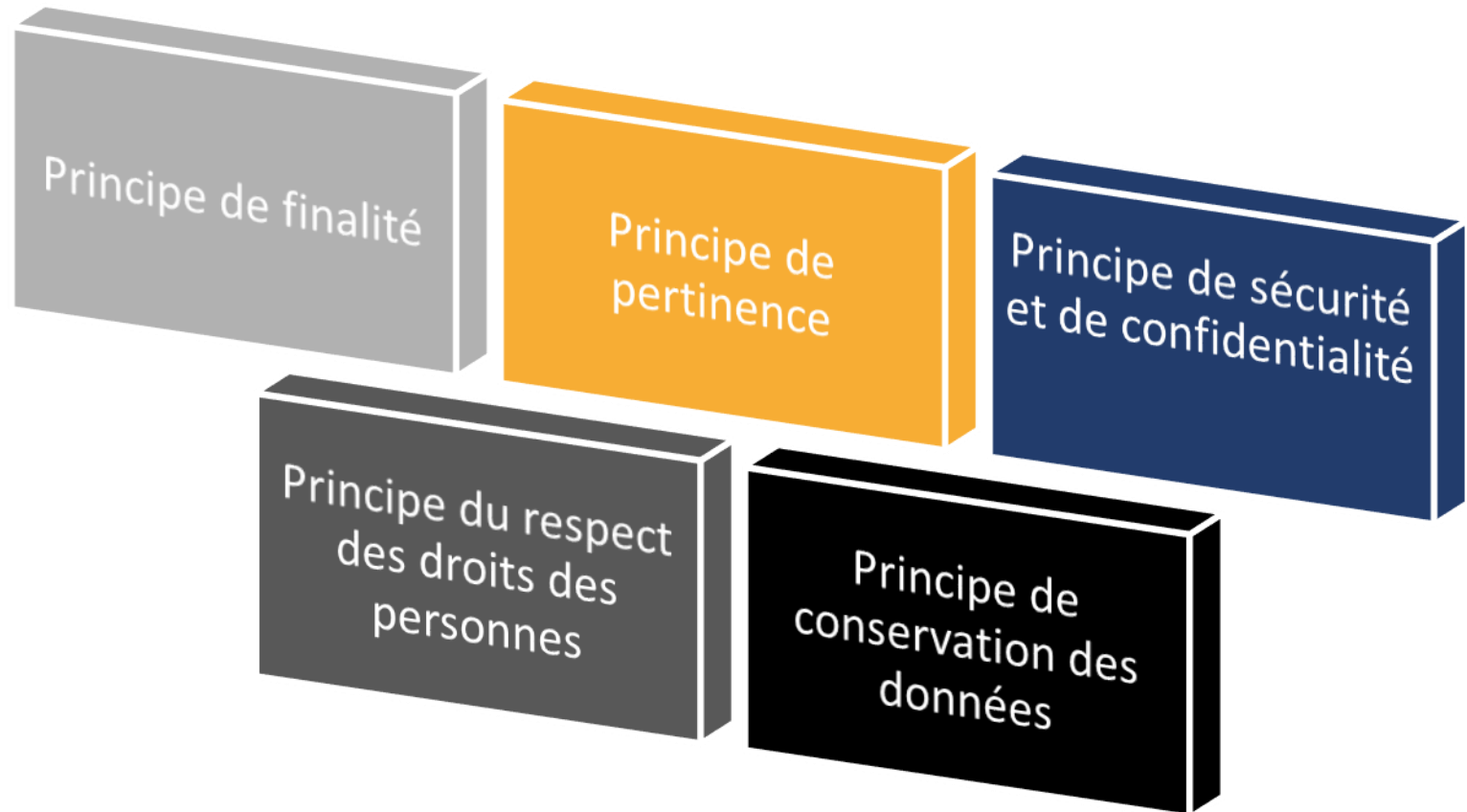


La prévention

- Quel risque veut-on diminuer ?
- Quelle action mettre en place ?
- Pour quelles personnes ?
- Les assurés vont-ils participer ?
- Comment évaluer le programme de prévention ?



Le contexte RGPD



RGPD : règlement européen pour la protection des données

Base de données *BENEFICIAIRES*

Informations relatives aux bénéficiaires :

- ▶ Identifiant du bénéficiaire
- ▶ Identifiant de l'assuré associé
- ▶ Période de couverture
- ▶ Type de bénéficiaires : Assuré, Conjoint, Enfant
- ▶ Année de naissance du bénéficiaire
- ▶ Sexe du bénéficiaire
- ▶ Identifiant de la zone géographique
- ▶ Niveau de la couverture santé
- ▶ Eventuellement : Secteur d'activité, CSP, régime spécifique, situation familiale et nombre d'enfants à charge de l'assuré...



Des variables permettant d'affiner la segmentation des bénéficiaires et d'intégrer **d'autres critères complémentaires à l'âge et au sexe**

Base de données PRESTATIONS

Informations relatives à la consommation des bénéficiaires :

- ▶ Identifiant du bénéficiaire
- ▶ Date de soins
- ▶ Acte concerné

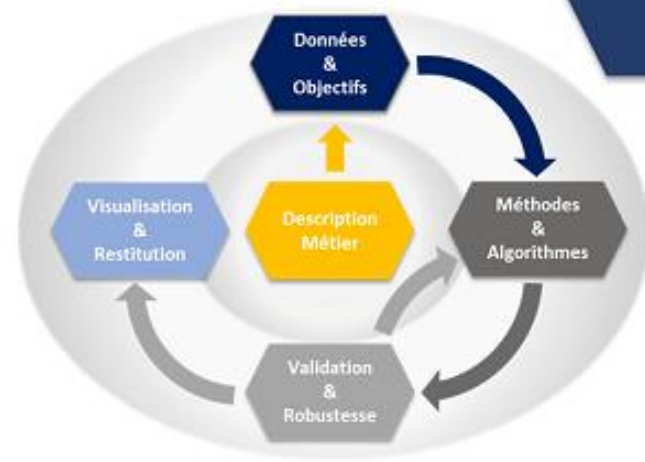
(en fonction d'une segmentation + ou – fine des actes)

- ▶ Dépense réelle
- ▶ Base de remboursement
- ▶ Remboursement obligatoire et taux de prise en charge
- ▶ Remboursement complémentaire, Reste à charge
- ▶ Eventuellement : indicateur d'une hospitalisation, d'un parcours de soins spécifique, spécialité consultée ou prescripteur...

Des informations plus ou moins agrégées, plus ou moins appauvries **dont la finesse des résultats va dépendre**



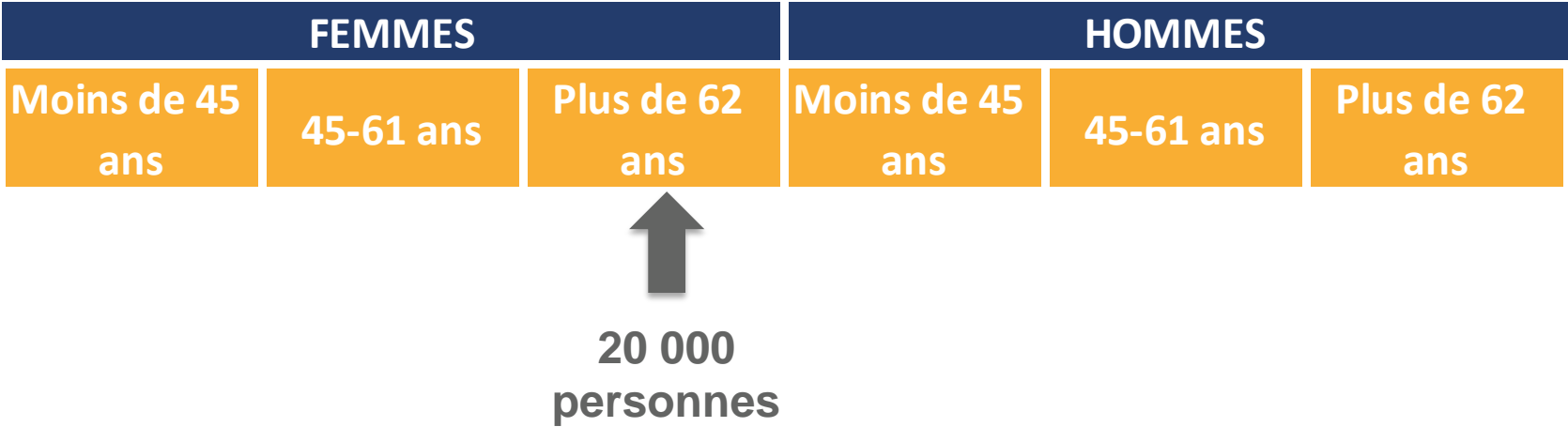
DONNEES ET OBJECTIFS



La base de données prestations

➤ Base de prestation santé utilisée pour l'étude :

- Contrat d'assurance santé à adhésion individuelle
- Un seul niveau de garantie
- 80 000 personnes observées sur une année
- 80 libellés d'actes



Le ciblage optimisé d'actions de prévention



Classifier les assurés de manière non supervisée pour proposer des actions de prévention

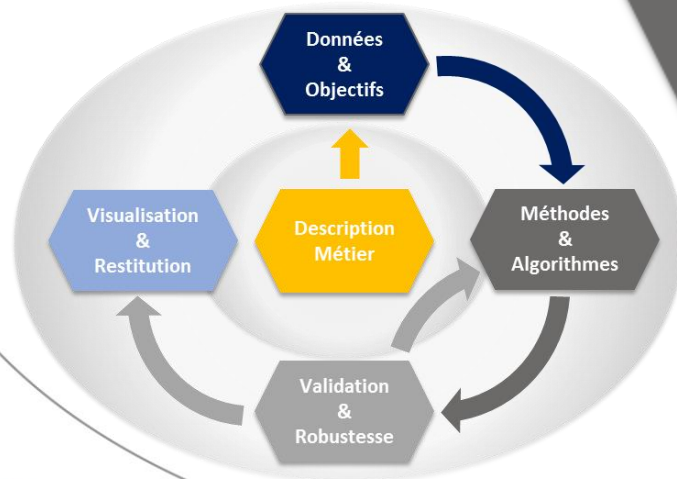


Identifier des comportements de consommation

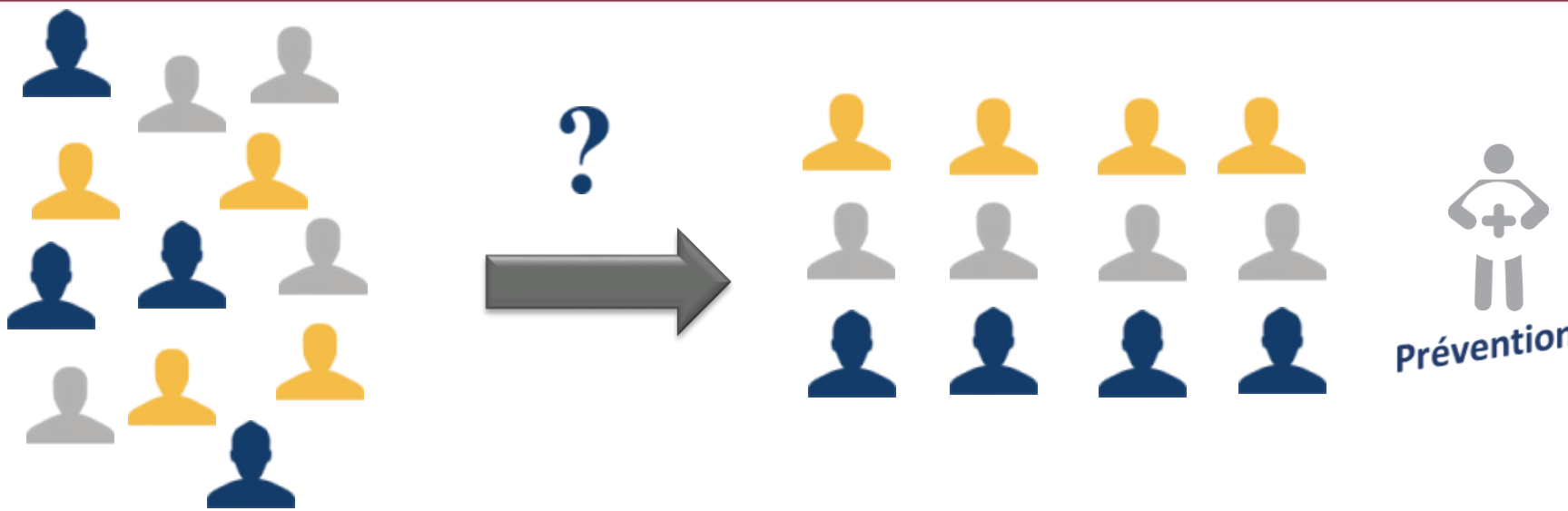


Constituer des groupes d'assurés qui ont les mêmes comportements de consommation

MODELISATION ET ALGORITHMES



Quelle approche non supervisée ?



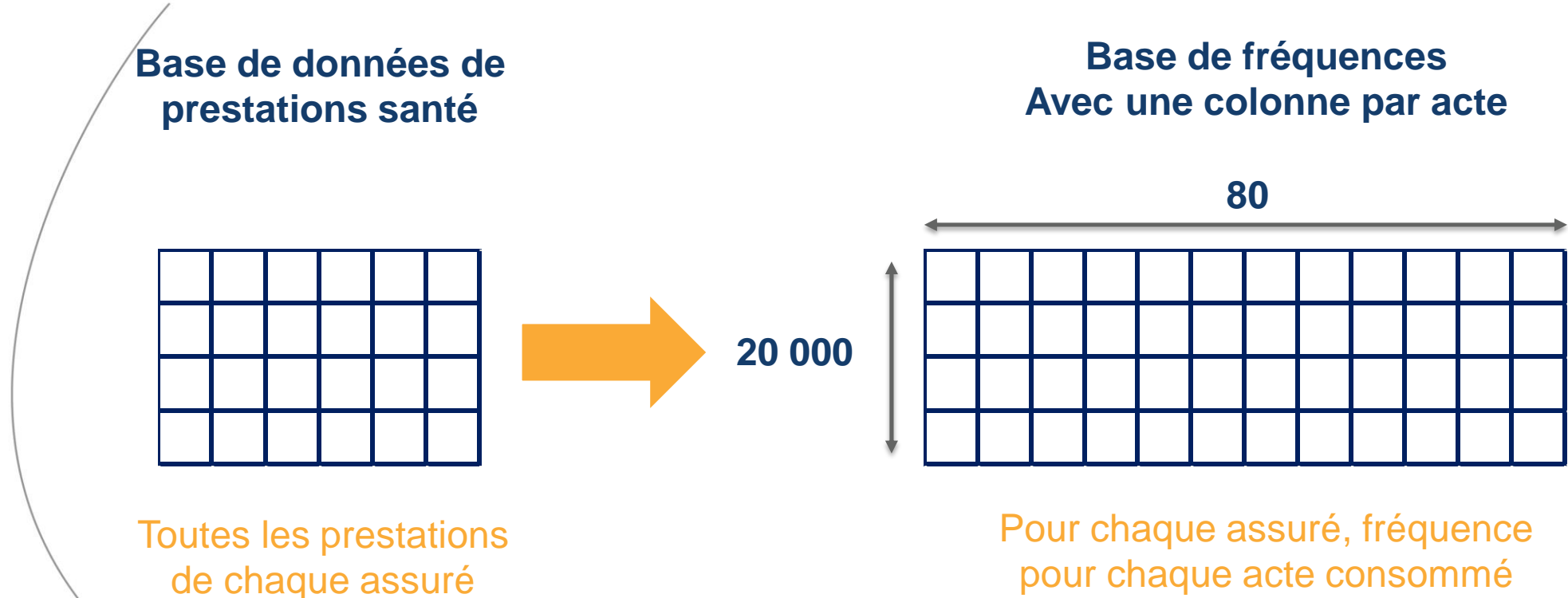
MODELISATION

**PRE-
TRAITEMENT**

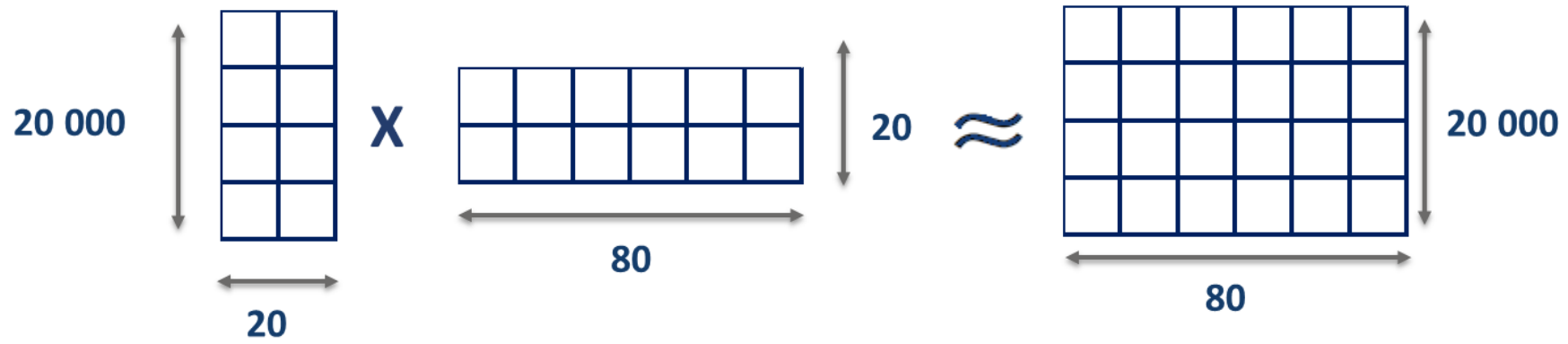
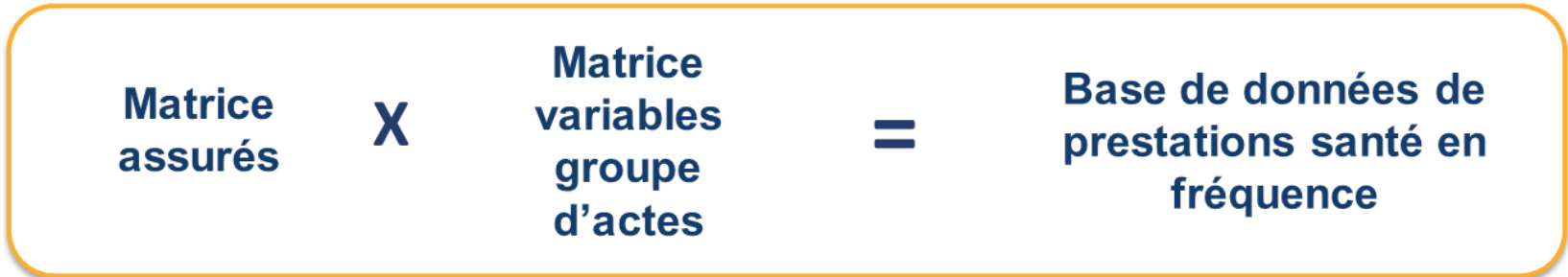
**REDUCTION
DE
DIMENSION**

CLUSTERING

Passage en fréquence



Réduction de dimension



REDUCTION DE DIMENSION

INTERPRETATION

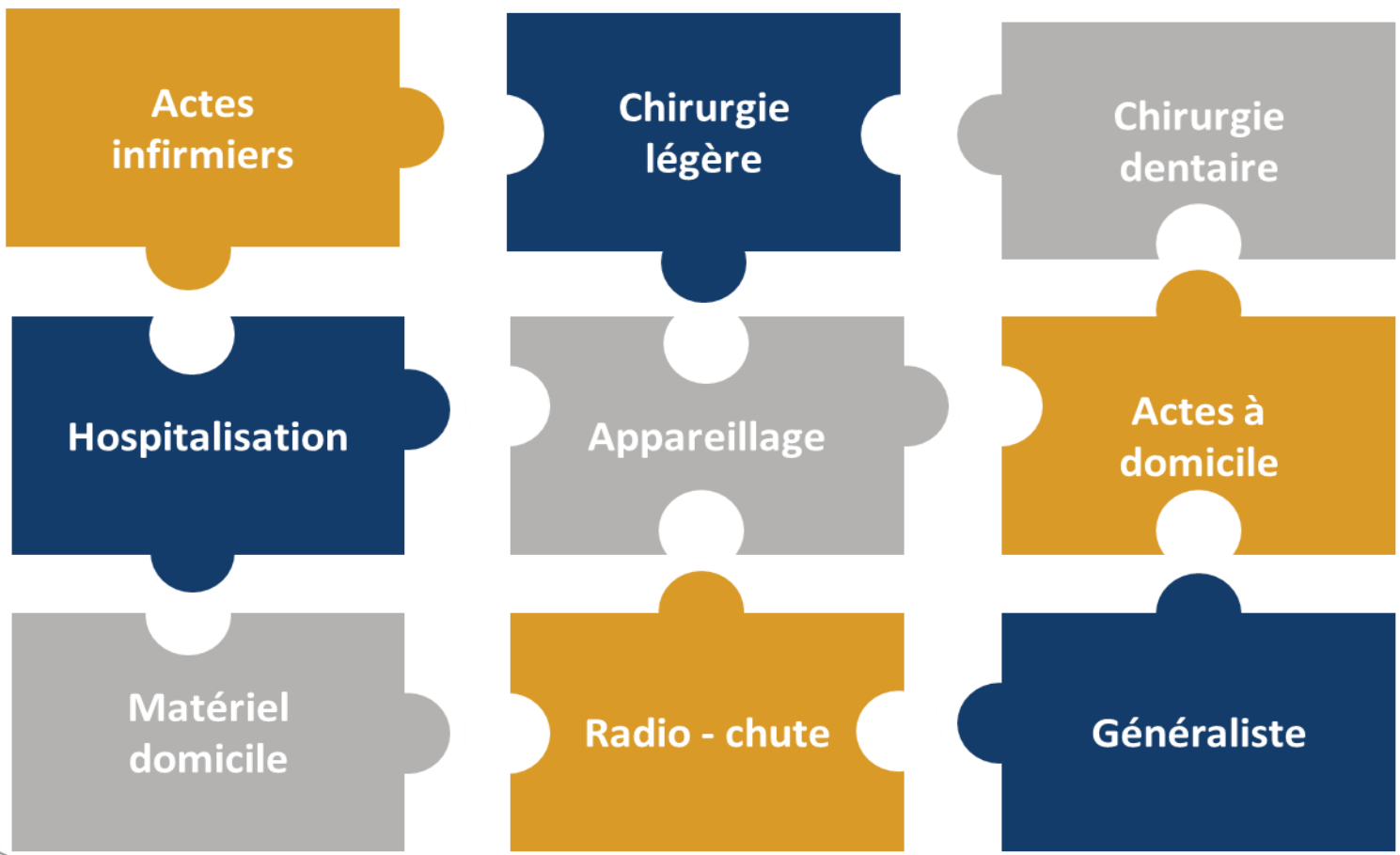
DATA

Matrice groupe d'actes – Normalisation horizontale

GROUPES D'ACTES	ACTES																																																			
	Acte à domicile	Acte d'imagerie	Acte spécialiste	Actes infirmiers	Appareil respiratoire	Appareillages divers	Auditif	Cardiologie	Chambre chirurgie	Chambre maladie	Chirurgie	Chirurgie dentaire	Consultation	Densitométrie	Dépassemement Honoraire	Dermatologue	Echographie	Forfait appareillage	Forfait cure thermique	Forfait journalier	Généraliste	Kinésithérapie	Majoration domicile	Matériel contention	Matériel domicile	Neuropsychiatrie	Optique	Orthèses	Orthoptiste	Ostéopathie	Pansements	Part Assuré Transi	Petit appareillage	Pharmacie blanche	Pharmacie bleue	Prélèvement sanguin	Prothèse dentaire	Prothèse externe	Radio	Soins Dentaires	Surveillance thermique	Taxi	Ticket modérateur	Transport	...							
1				■																																																
2																													■																							
3																																																				
4																																																				
5																												■																								
6																																																				
7																																																				
8																																																				
9																																																				
10																																																				
11																																																				
12																																																				
13																■																																				
14																																																				
15																																																				
16																																																				
17																																																				
18																																																				
19																																																				
20																																																				

Interprétation des groupes d'actes – Comportement de consommation

Réduction de la dimension : 80 actes → 20 groupes d'actes

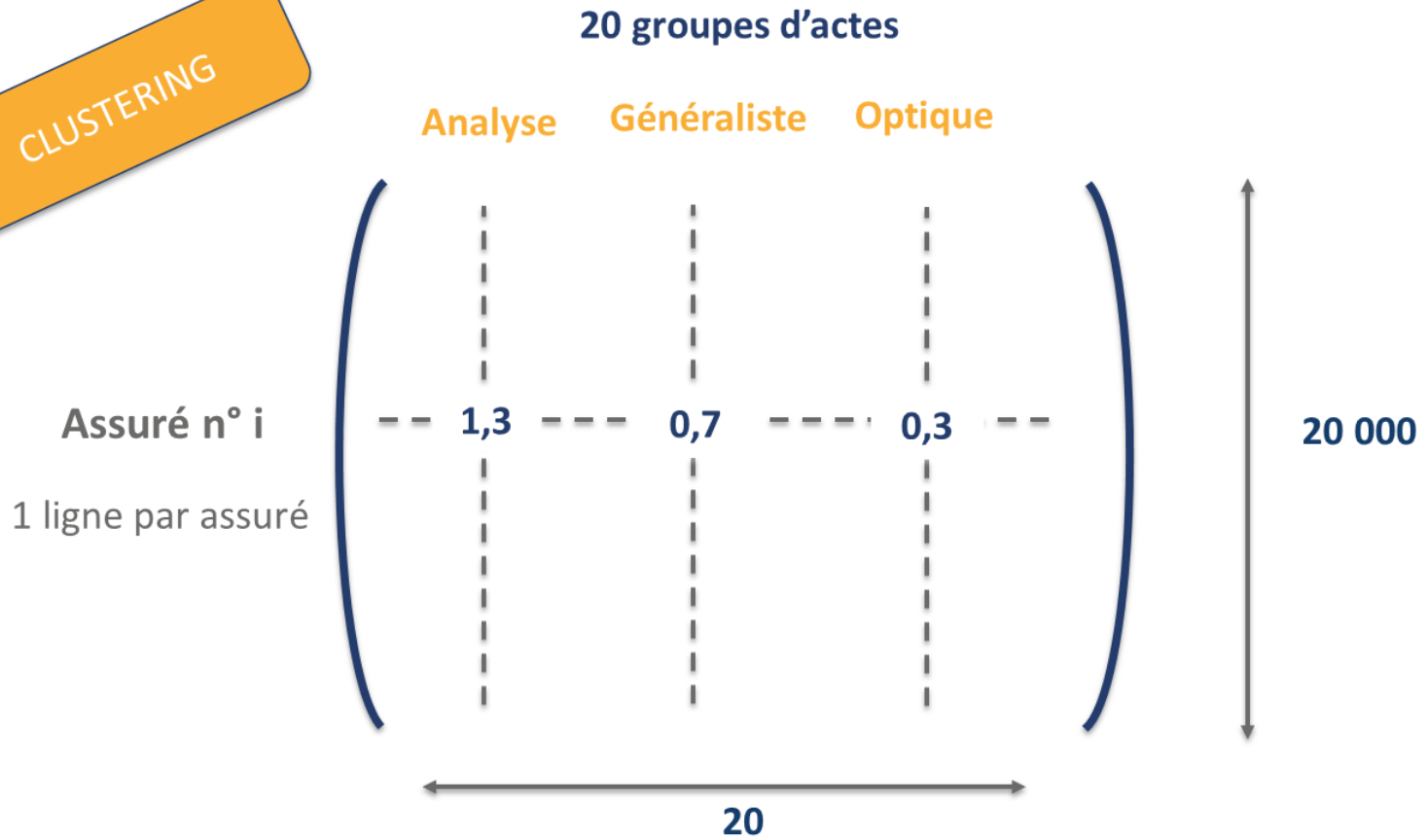


Interprétation des groupes d'actes selon les tranches d'âges

	FEMMES	
EXEMPLES	45-61 ans	Plus de 62 ans
Groupes d'actes communs	Pharmacie Optique Dentaire Analyse	
Groupes d'actes spécifiques	Orthoptie Psychiatrie Ostéopathie	Soins à domicile Matériel domicile Hospitalisation

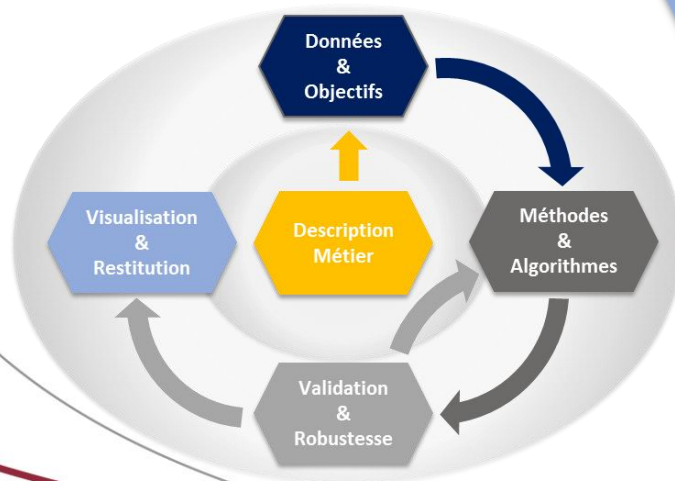
Matrice assurés consommant ces groupes d'actes

CLUSTERING

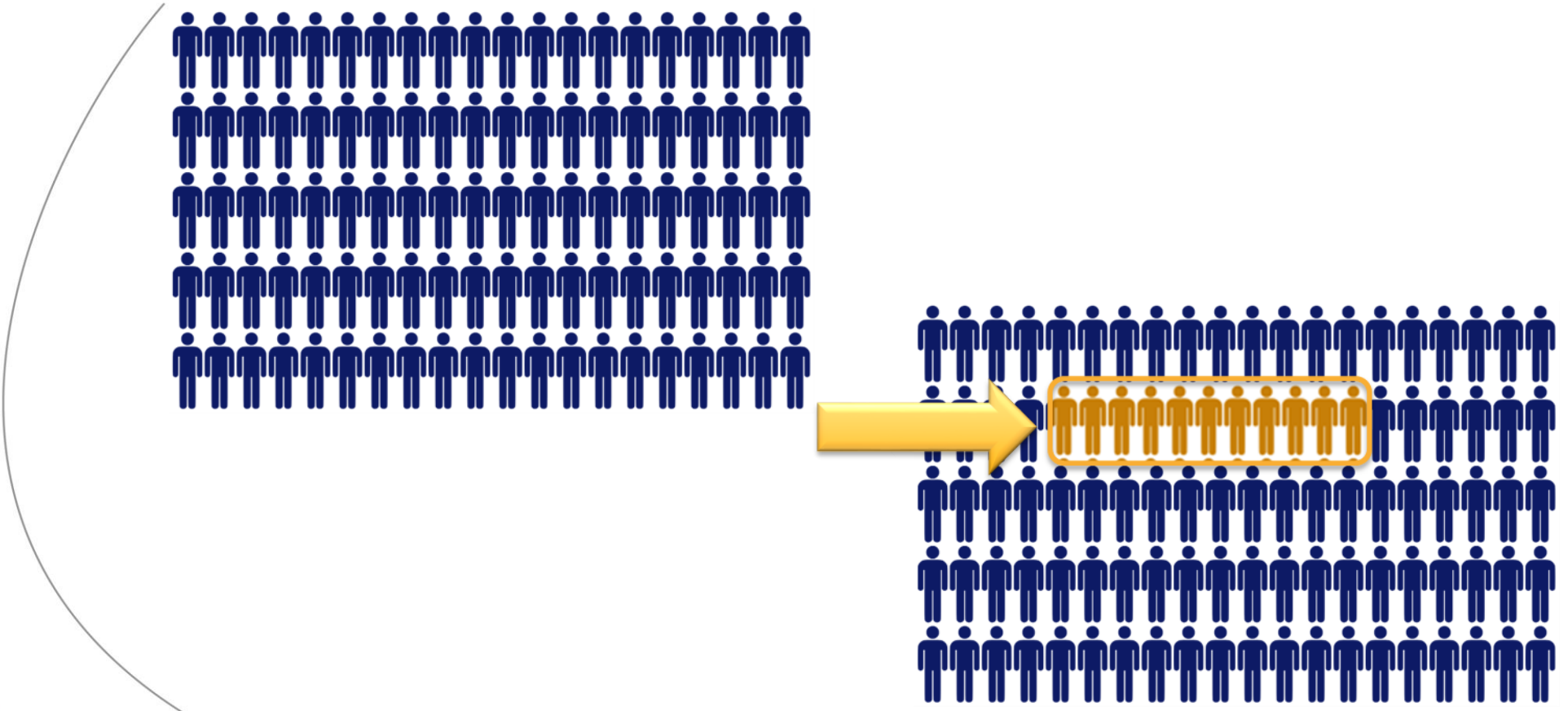


- La dimension réduite permet de classifier nos assurés

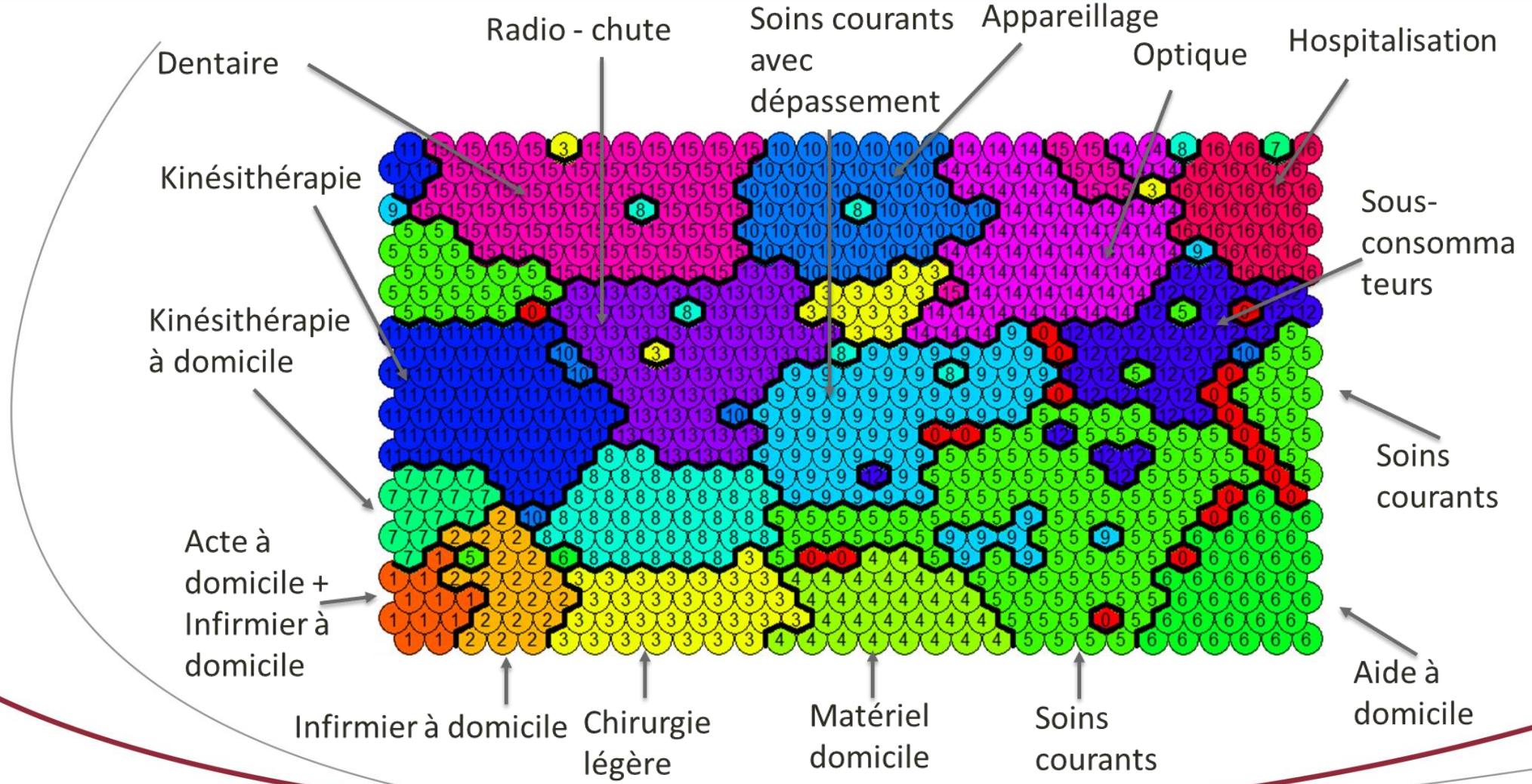
VISUALISATION ET RESTITUTIONS



Le ciblage de segments d'assurés



Quels segments ?



Identification d'un parcours de fragilité



69-70 ans

Optique
Dentaire
Sous-consommateur
Soins courants avec ou
sans dépassement
Radio – chute



73-74 ans

Appareillage
Chirurgie légère
Infirmier à domicile



77-82 ans

Hospitalisation
Acte à domicile
Acte à domicile + infirmier à
domicile
Kinésithérapie à domicile

Quel programme de prévention ciblée ?



70 ans

- ✓ **Activités communes** au senior pour éviter l'isolement
- ✓ **Visite d'étudiant en psychologie** (signe de dépression)
- ✓ ...



74 ans

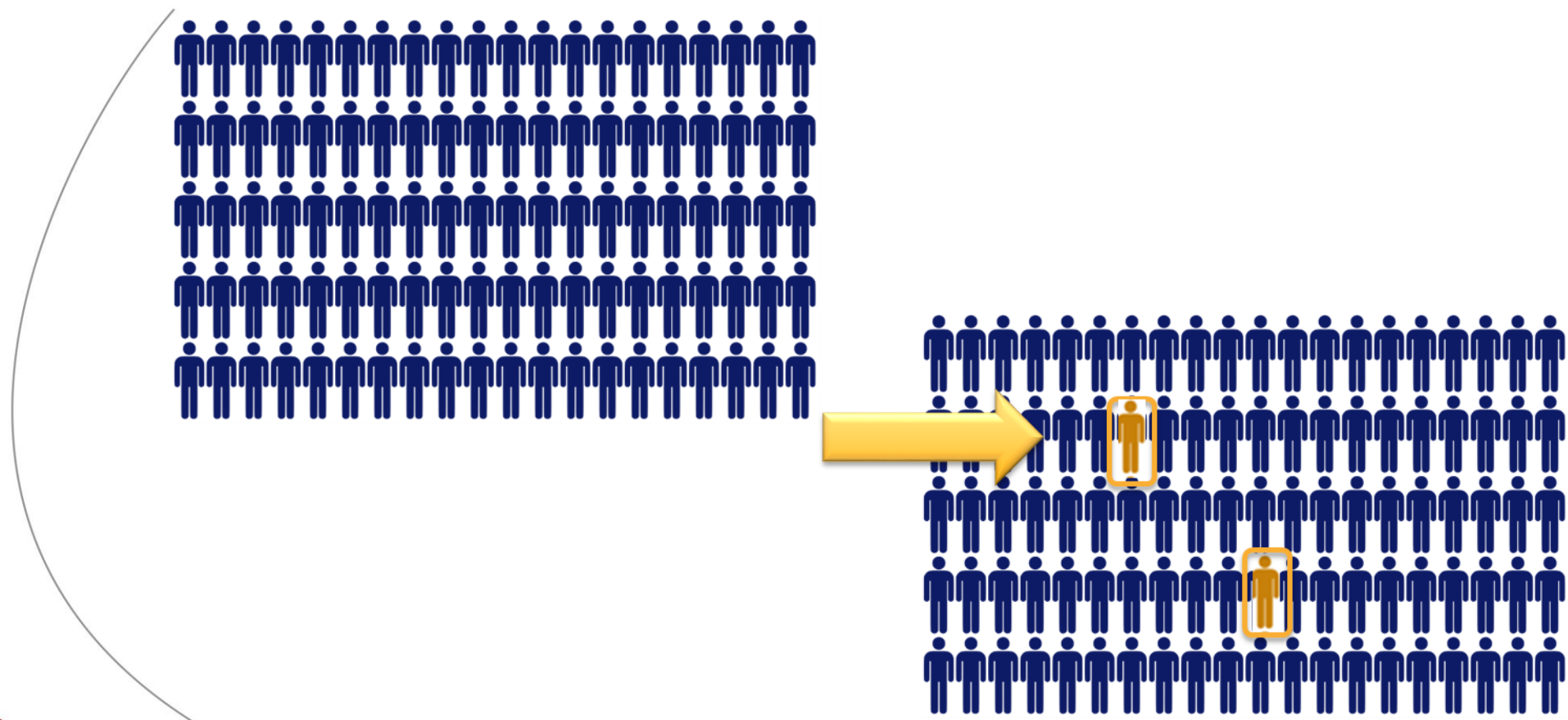
- ✓ **Sport adapté** pour éviter l'entrée en dépendance lourde
- ✓ **Formation de la famille** à la détection des signes de dépendance
- ✓ ...



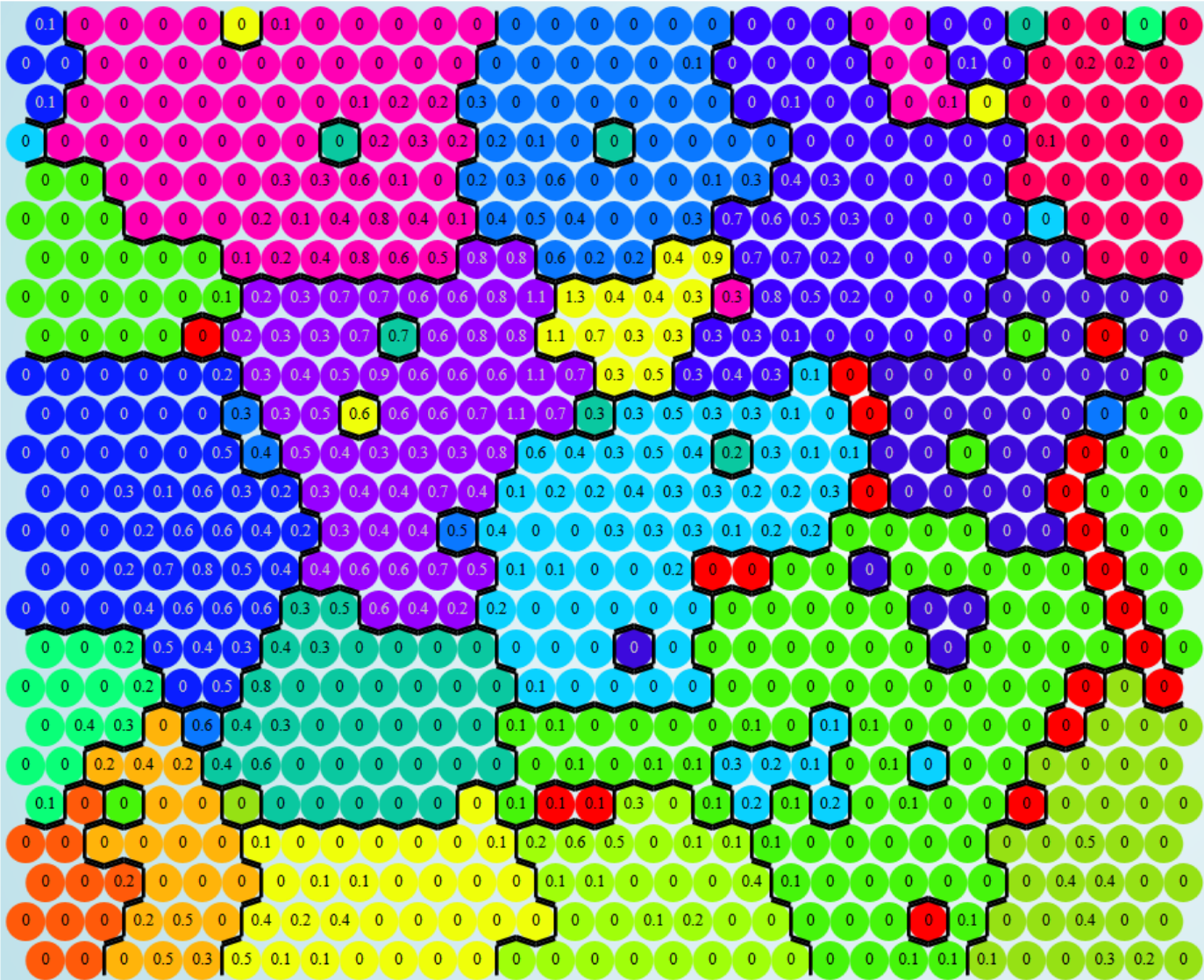
80 ans

- ✓ **Télémédecine** pour favoriser le suivi à domicile
- ✓ **Objets connectés** pour favoriser le suivi quotidien des personnes dépendantes et alerter si accident ou anomalie mesurée
- ✓ ...

Le ciblage d'assurés



Quel besoin de prévention pour un assuré ?



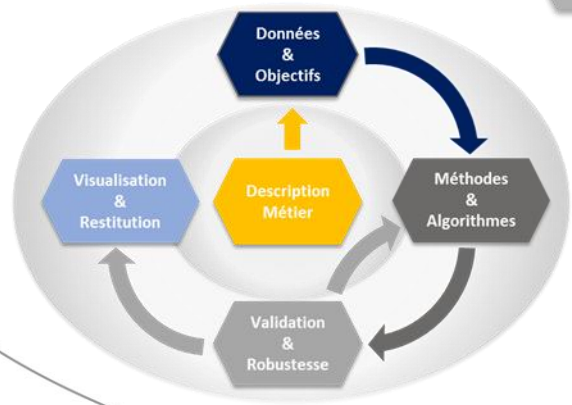
Probabilité d'appartenance à une classe de risque

L'assuré n°96 appartient à la classe :

- « Radio-chute » avec 30% de chances
- à la classe « kinésithérapeute » avec 11% de chances
- à la classe « Chirurgie légère » avec 10% de chance

Les autres classes n'étant pas significatives

VALIDATION ET ROBUSTESSE



Comment valider notre méthode de classification non supervisée ?

Robustesse



**Classe fictive
d'assurés**



**Fort
consommateur**



**Autre base de
données
prestation**



Text mining



**Comparaison
méthodes**

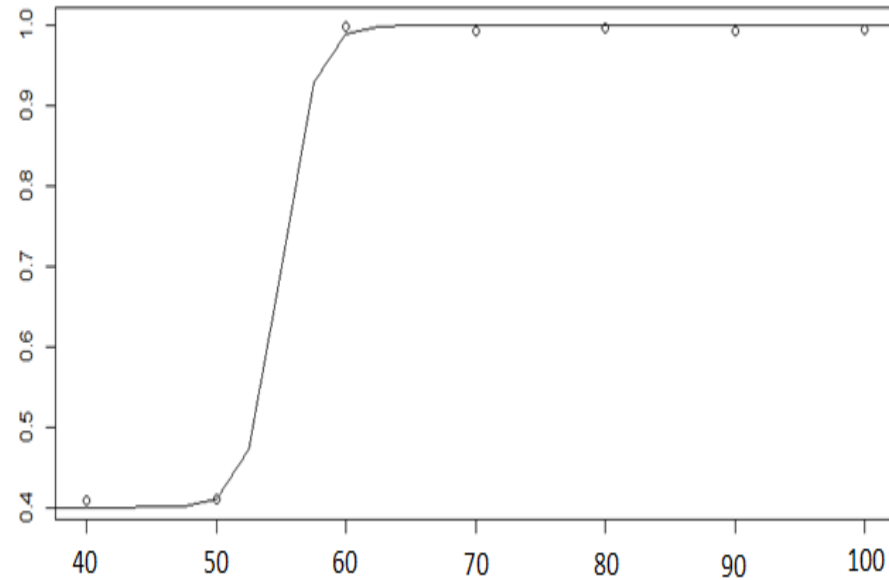


**Classification
seule**

Validation

Classe fictive d'assurés

- **Tester la robustesse**
 - ✓ Création d'une classe fictive d'assurés
 - ✓ A partir de combien d'assurés fictifs cette classe est-elle détectée ?
- Exemple :
 - ✓ Création de corrélation entre libellés



Quels segments de prévention pour les forts consommateurs ?



48 ans

Psychiatrie et
dépassement
d'honoraires



54 ans

Soins dentaires
Optique
Echographie



64 ans

Matériel domicile
Chirurgie légère et
spécialiste
Pharmacie



71 ans

Hospitalisation

Autre base de prestations santé

	FEMMES	
EXEMPLES	Base 1	Base 2
Groupes d'actes communs	Acte à domicile Chirurgie Analyse	
Différences dues à la construction de la base	Ticket modérateur Radio / Chute	Actes biologie Kiné orthopédique Radio
Autres différences	Appareillage	Neuropsychiatrie

Text Mining

- **Tester sur un jeu de données déjà classifié**
 - ✓ Ancien jeu de données classifié
 - ✓ Jeu de données issu du text mining
 - ✓ Permet de comparer les techniques entre elles

Mail 1	La data science, c'est quoi ?
Mail 2	C'est la science des données. Elle regroupe la collecte et l'exploitation de la donnée afin mieux comprendre la source dont elles proviennent.
Mail 3	Donc si j'ai bien compris la data-science, c'est apprendre à connaître le monde grâce à l'informatique !



	Data-science	Comprendre	Donnée	la
Mail 1	1	0	0	1
Mail 2	0	1	2	4
Mail 3	1	1	0	1

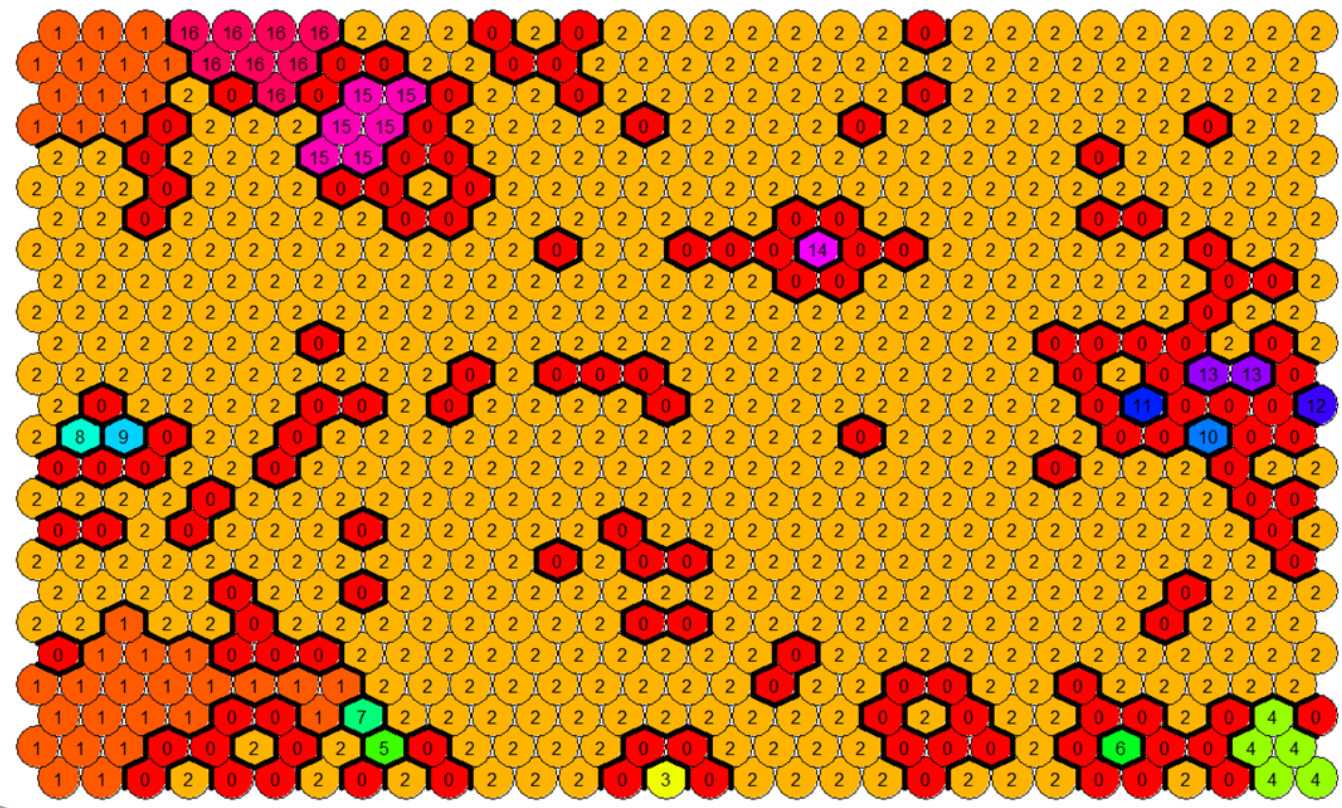
Comparaison avec ACP

- **Comparer différentes méthodes entre elles**
 - ✓ Avec un jeu de données externe
 - ✓ Avec notre jeu de données
 - ✓ Quelle méthode donne les classes les plus cohérentes ?
- **Exemple : Comparaison de deux méthodes**
 - ✓ Quelle méthode regroupe le mieux les actes d'optique ?

Acte de chirurgie	Acte de spécialité	Actes infirmiers	Actes radiologie	Analyses biologie	Appareillage divers	Cardiologie	Chambre Maladie	Forfait journalier	Généraliste	Kinésithérapie	Lentilles	Lentilles acceptées	Monture de lunettes	Neuropsychiatrie	Orthophoniste	Orthoptiste	Pansements	Réparation verres	Soins dentaires	Soins infirmiers	Spécialiste	Verres de lunettes	Vignettes blanches	Vignettes bleues	

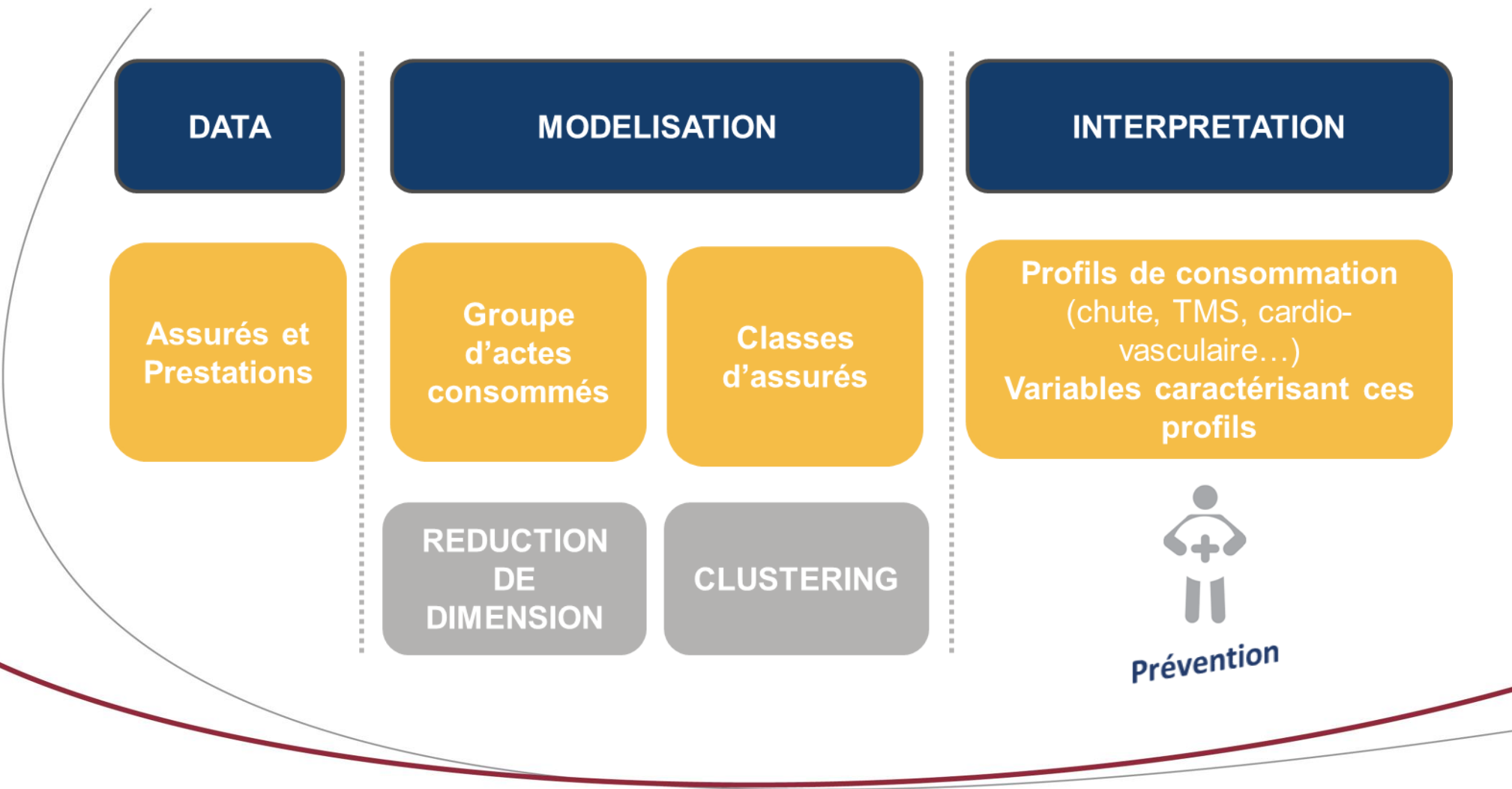
Suppression de la réduction de dimension

- Comparer différents process entre eux
- Pourquoi réduire la dimension ?
 - ✓ Fléau de la dimension



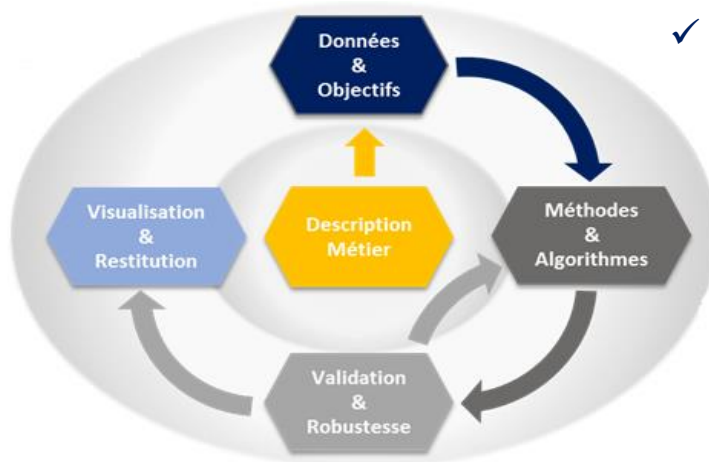
CONCLUSION

Le ciblage optimisé d'actions de prévention en 3 étapes



Une méthodologie structurée pour designer les méthodes Data Science

- ✓ La méthodologie permet de **structurer** les projets Data Science



- ✓ Selon un **même process**, quel que soit le domaine d'application

- ✓ Permet une **traçabilité** des hypothèses et des méthodes

- ✓ Capitalise et **mutualise** les développements

- ✓ **Facilite les échanges** avec les experts métier

- ✓ **Favorise la réussite** des projets de Data Science