

Khiops

an end-to-end solution for
Automated Machine Learning

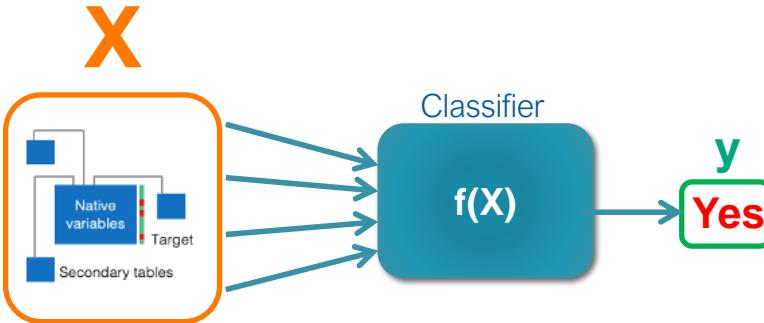
Nicolas Voisine

Marc Boullé

NormaSTIC 26 juin 2025

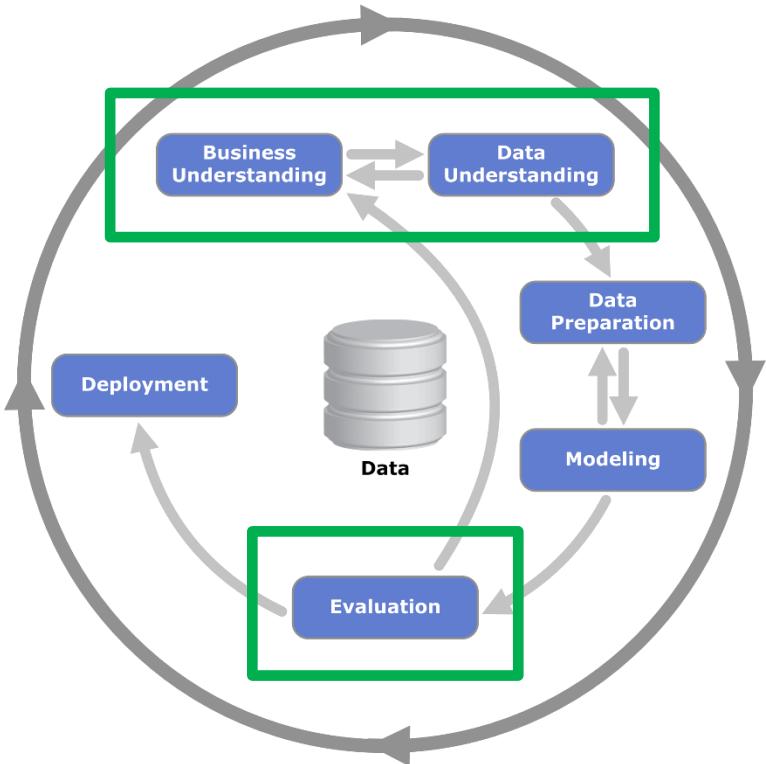


Orange needs to empower data experts (20 years ago)



1. **Business-Driven and Impact-Oriented:** Agnostic to use cases (e.g., detection of fraud, churn or anomalies), robust to poor quality data and focus on business understanding
2. **Production-Ready and Cost-Effective Scalability:** Built to process complex multi-table data structures efficiently and ensure reliable predictions in production environments
3. **Interpretability and robustness:** Ensured that models are interpretable and robust, enabling business experts to extract meaningful insights and understand the story behind the data

Focus on data and business understanding, the rest is automated



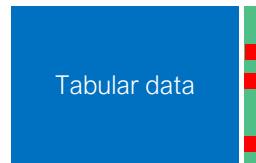
Most use cases don't justify allocating 30 guys for marginal improvements:

- Small improvements rarely influence real-world decisions and lever actions
- Over-processing just adds waste (muda) without return on investment

The green steps are for you; the rest is for the robots!

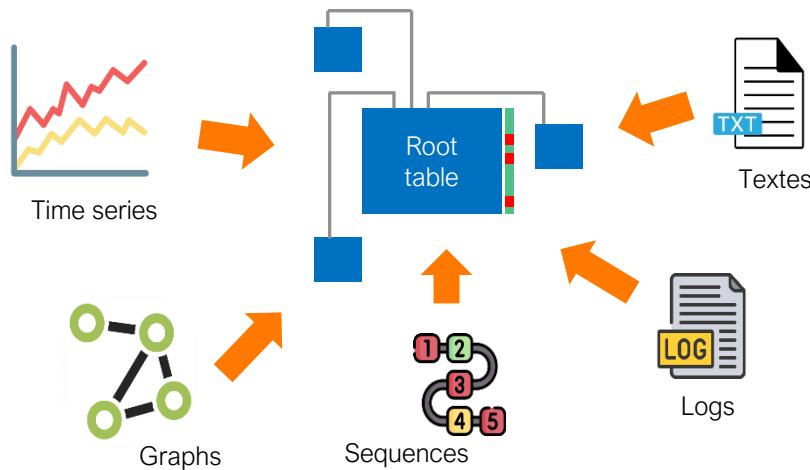
Tabular data in real-world

At school, we learn Data science with toy datasets



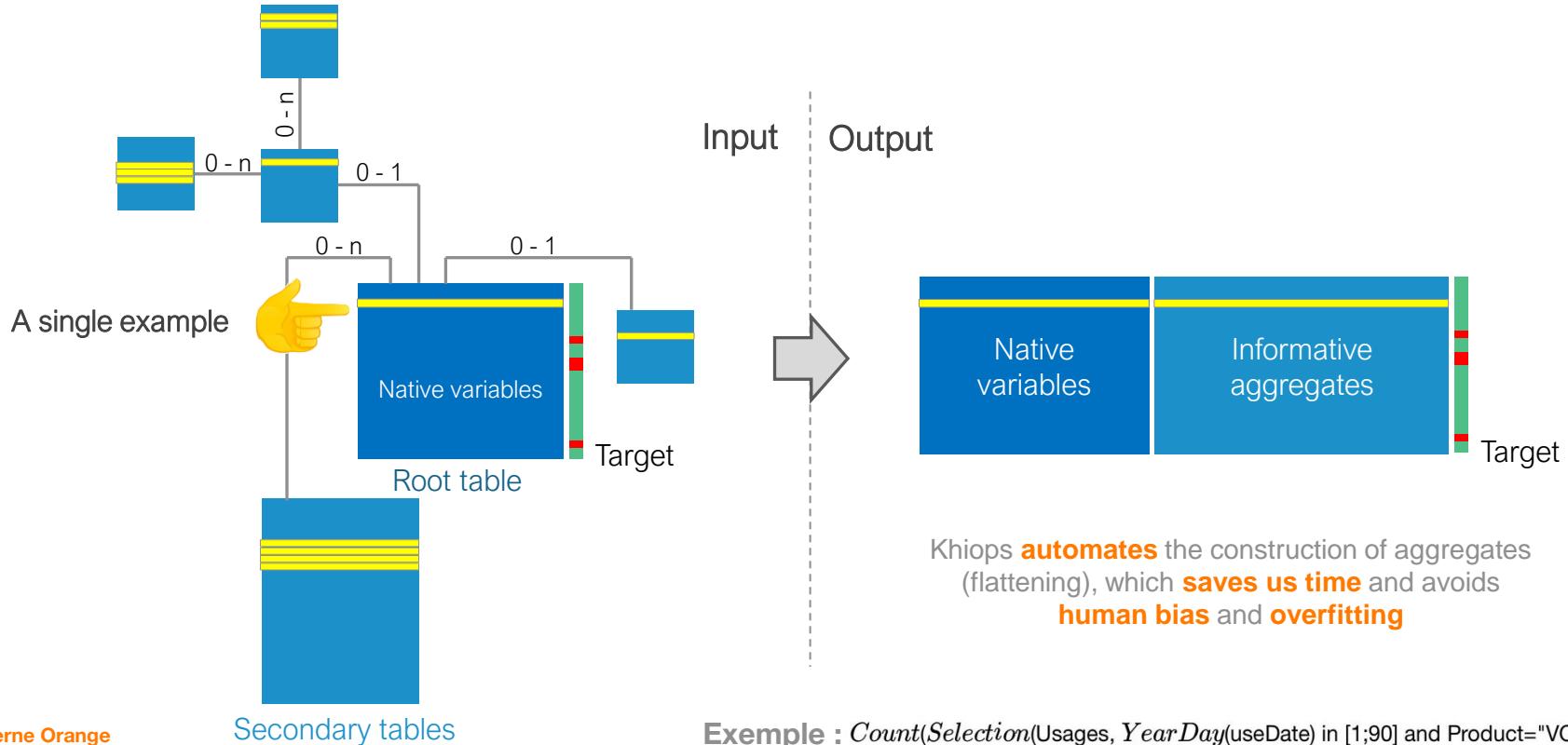
A statistical individual = A single row
Very easy !!

In reality, data is **much more complex**



Multi-table schema
can encode
diverse type of data

Khiops champion to handle complex data

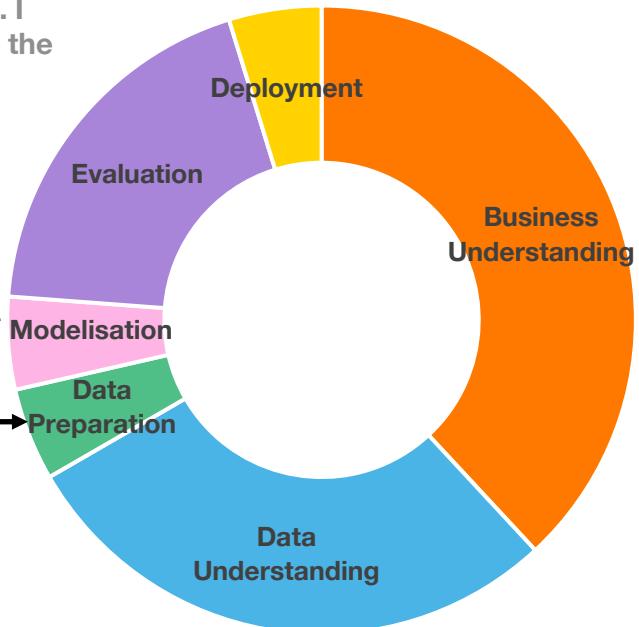
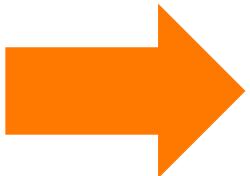


Months of savings with Khiops



Standard CRISP distribution
on multi-table data

Khiops delivers results in which I always have total confidence. I save time because I question the data rather the model.



Distribution with Khiops

Robustness,
interpretability and
efficiency facilitate
inference

Khiops natively digest
CDRs (relational data)
and manages encoding
and feature engineering

Model Interpretability for informed decision-making

Khiops results can be easily analyzed using its **built-in visualization tool**

The way each variable is linked to the target is **clear**



The final model selects only **few independant variables**

34 Variables      

name	level	weight 	importance 
0.00878912	0.929718	0.0903958	
0.00307631	0.730499	0.0474051	
0.0077074	0.683624	0.0725876	
0.00355536	0.609406	0.0465474	
0.0115712	0.532257	0.0784783	
0.00330293	0.480499	0.0398379	
0.0359826	0.448273	0.127004	
0.0588135	0.422882	0.157706	
0.00886278	0.417999	0.0608657	
0.000652961	0.343781	0.0149825	
0.025324	0.314484	0.0892412	
0.00362489	0.250031	0.0301054	
0.0303871	0.246124	0.0864812	
0.0515375	0.236359	0.110369	
0.00938844	0.222687	0.045724	
0.0234229	0.22171	0.0720632	
0.016231	0.216827	0.0593239	
0.0154212	0.203156	0.0559724	
0.0340942	0.187531	0.0799607	

A master of frugality on big data

A single end-to-end formalism, no hyperparameters and low-level implementation for outstanding performance



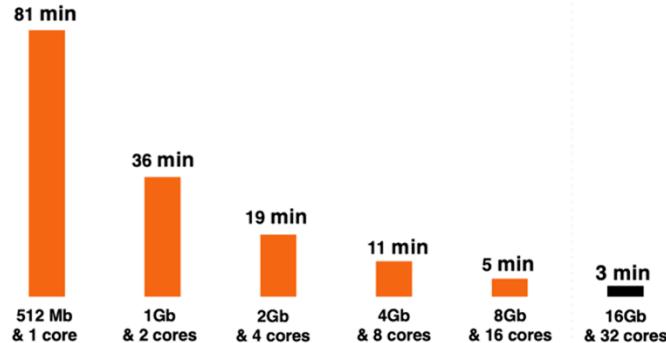
- On a use case for used battery detection:
Khiops was **50x faster** with **200x fewer resources** than Google AutoML
- On a benchmark based on 8 internal datasets, Khiops uses **50x fewer resources** than **Prevision.io** and **4x fewer** than **H2O** for comparable prediction performance

It's **green**, it's **cheaper**, and it adapts easily to available resources: **from mono-thread** with little memory (out-of-core) to **multi-machines (Kubernetes native)**



Illustration:

Khiops run on our LiveBox (very limited memory)





More than 130 scientific publications in national and internationals journals and conferences



State of the art software engineering:

- **400 000** lines of code C++ for the core part
- **10 000** lines Java code for the Desktop GUI
- **15 000** lines of code for khiops-python
- **800** sets of non-regression tests
- Evaluated with volumetric test up to:
 - Hundreds of millions of instances
 - Hundreds of thousands of variables
 - Hundreds of GB for data files

Short demo

We use the « as a service » version of Khiops
that would be move to open source shortly



DEMO



NumAcc	An	Jour	...	Mortel
201800000001	18	1	...	False
201800000002	18	2	...	False
201800000003	18	3	...	False
...	



NumAcc	Cat Route	...	Ecole
201800000001	3	...	False
201800000002	4	...	False
201800000003	5	...	False
...	

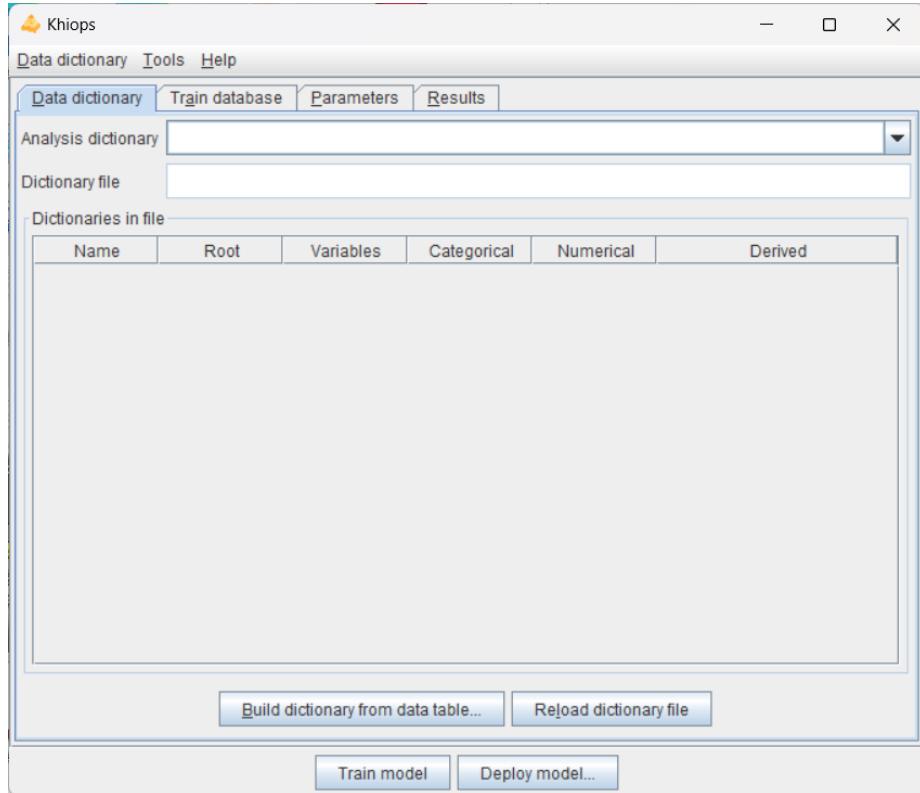


NumAcc	Num Veh	Sens circ	...	Manoeuvre
201800000001	B01	1	...	Sans chang.
201800000001	A01	1	...	Tourn. Droite
201800000002	A01	1	...	Sans chang.
...	



NumAcc	Num Veh	Place	...	Annee naissance
201800000001	B01	1	...	1928
201800000001	A01	1	...	1960
201800000002	A01	1	...	1947
...	

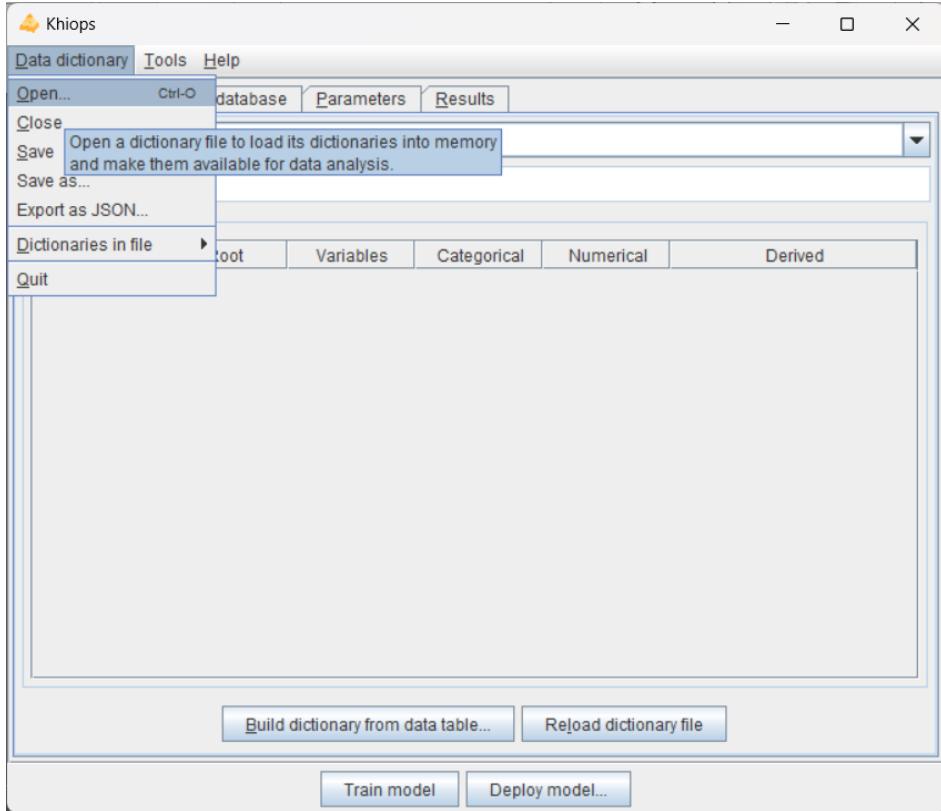
Short demo



1 Open dictionary



Short demo



1 Open dictionary



Short demo

Khiops

Data dictionary Tools Help

Data dictionary Train database Parameters Results

Database files

Data root	Path	Dictionary	Data table file
Accident		Accident	C:\Users\Public\khiops_data\samples\Accidents\Accidents.b...
Accident	Place	Place	C:\Users\Public\khiops_data\samples\Accidents\Places.txt...
Accident	Vehicles	Vehicle	C:\Users\Public\khiops_data\samples\Accidents\Vehicles.txt...
Accident	Vehicles\User:	User	C:\Users\Public\khiops_data\samples\Accidents\Users.txt...

Header line used

Field separator

Sample percentage 70.0

Sampling mode Include sample

Selection variable

Selection value

Test database Complementary

Detect file format

Inspect test database settings...

Train model Deploy model...

2 set data set files



Short demo

Khiops

Data dictionary Tools Help

Data dictionary Train database Parameters Results

Target variable Gravity

Main target value Lethal

Predictors Recoders Preprocessing System parameters

Selective Naïve Bayes predictor

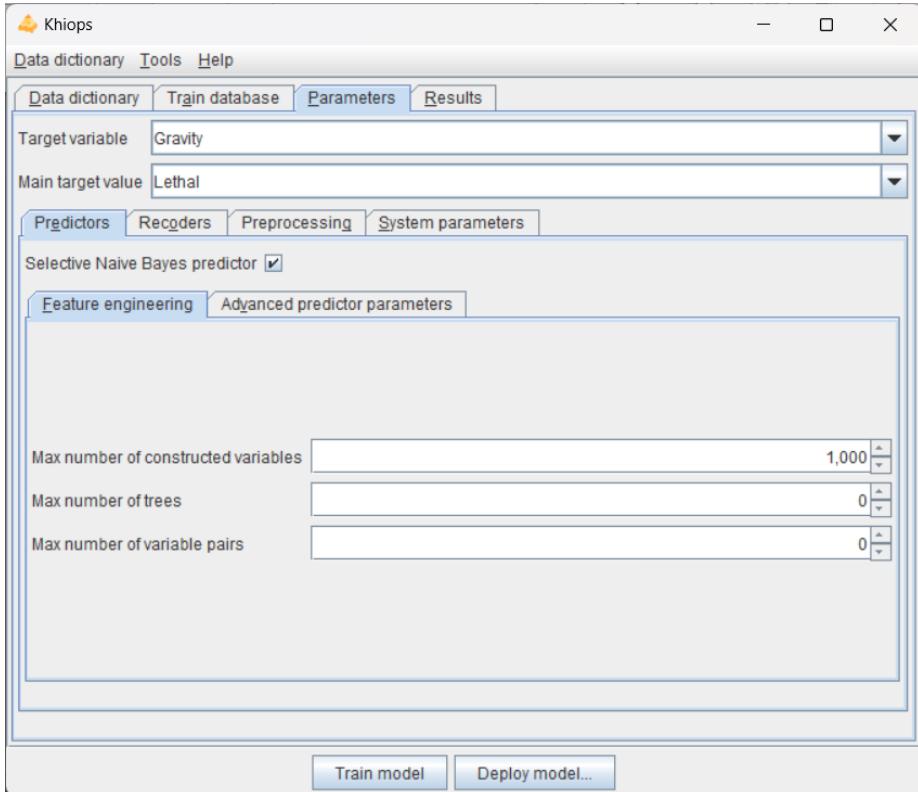
Feature engineering Advanced predictor parameters

Max number of constructed variables 1,000

Max number of trees 0

Max number of variable pairs 0

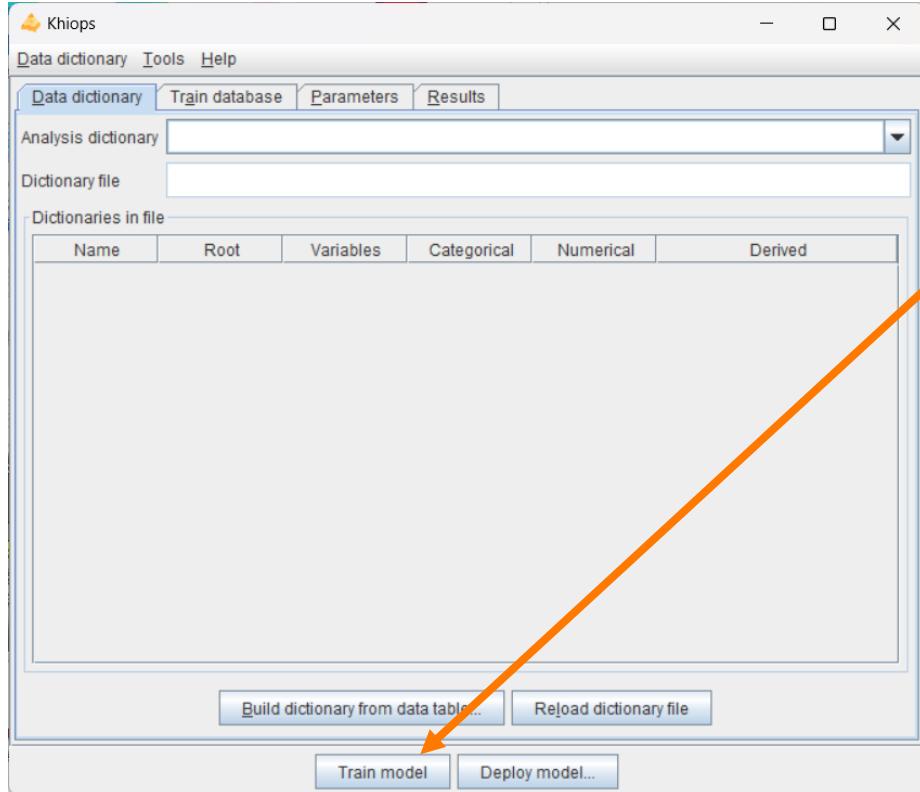
Train model Deploy model...



3 analysis configuration



Short demo



4 Run





Summary

Dictionary

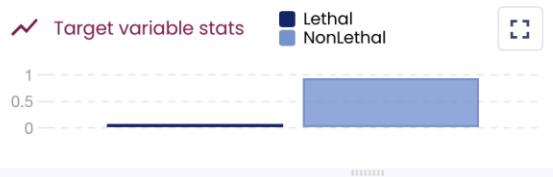
Accidents
s3://br-diod-ocp-fr01-khiops-as-a-service-prod-s3bucketclaim01/project/edc48638-babe-40e8-9f75-b61486799858/Accidents.txt_sorted

Database

Target

Select trained predictor

Selective Naive Bayes



Trained predictors

Selective Naive Bayes 54 Variables

54 Variables

Level distribution

Name	Level	Weight	Importance
Max(Vehicleless.Min(Userss.BirthYear)) where Categ	0.0119991	0.181641	0.0466855
Light	0.023844	0.162598	0.0622655
Hour.DecimalTime	0.00597559	0.640625	0.0618717
Mode(Vehicleless.Maneuver) where Category	0.0259667	0.128906	0.0578556
Mean(Vehicleless.Mean(Userss.BirthYear))	0.0164264	0.19043	0.0559292
Mode(Vehicleless.Category) where Maneuve	0.0166102	0.146484	0.0493268
Mode(Vehicleless.Category) where Direction	0.0109246	0.204956	0.0473188
Max(Vehicleless.Min(Userss.BirthYear))	0.0119991	0.181641	0.0466855
Mode(Vehicleless.Maneuver) where Maneuvre	0.0307446	0.0644531	0.044515
Mode(Vehicleless.ImpactPoint) where Categ	0.013172	0.138184	0.0426633
CountDistinct(Vehicleless.Direction)	0.00363414	0.498047	0.0425437
Mode(Vehicleless.Maneuver) where ImpactPo	0.0108922	0.158203	0.0415112
Count(Vehicleless) where Category > Categ	0.00456561	0.252516	0.0401749

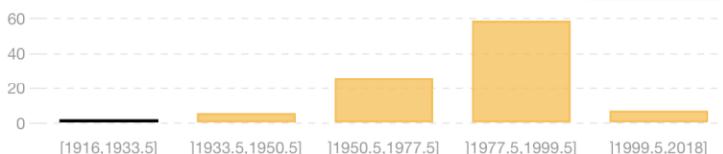
Max(Vehicleless.Min(Userss.BirthYear))

Scale chart

Distribution

Coverage

Coverage

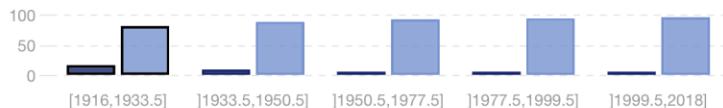


Target distribution

Lethal NonLethal

Values

Probabilities



Current interval

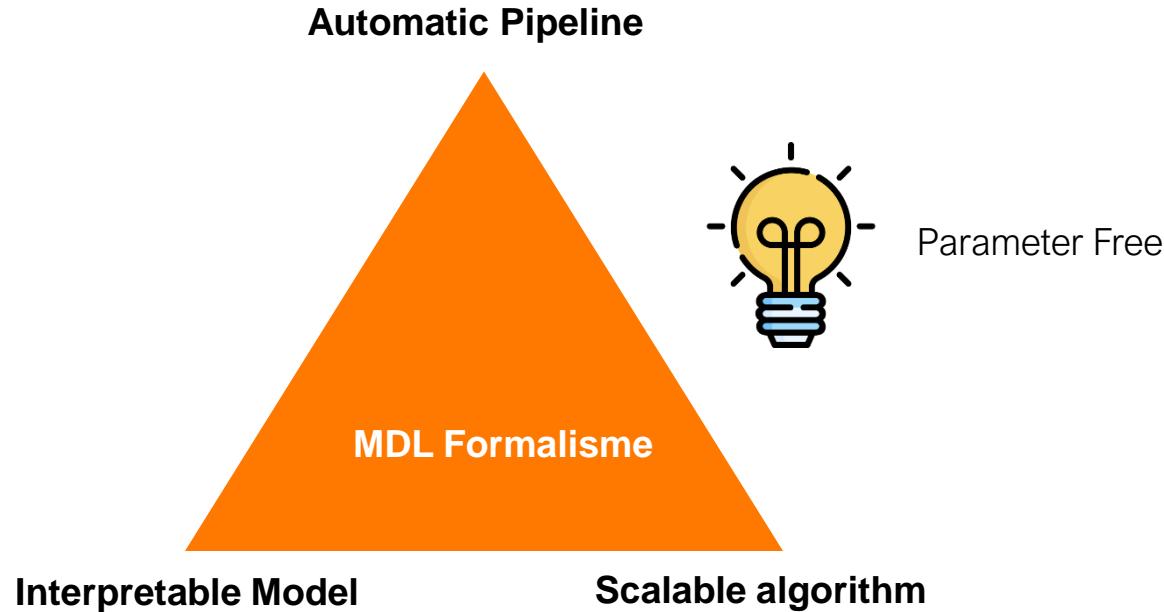
Conclusion

« Khiops est un outil *productif*,
qui *s'adapte* à beaucoup de cas d'usage,
et fait *gagner* énormément de *temps* dans les projets de data-science.»



Nous pouvons intervenir chez vous pour des ateliers d'accompagnement.
Contactez-nous à **khiops.team@orange.com**

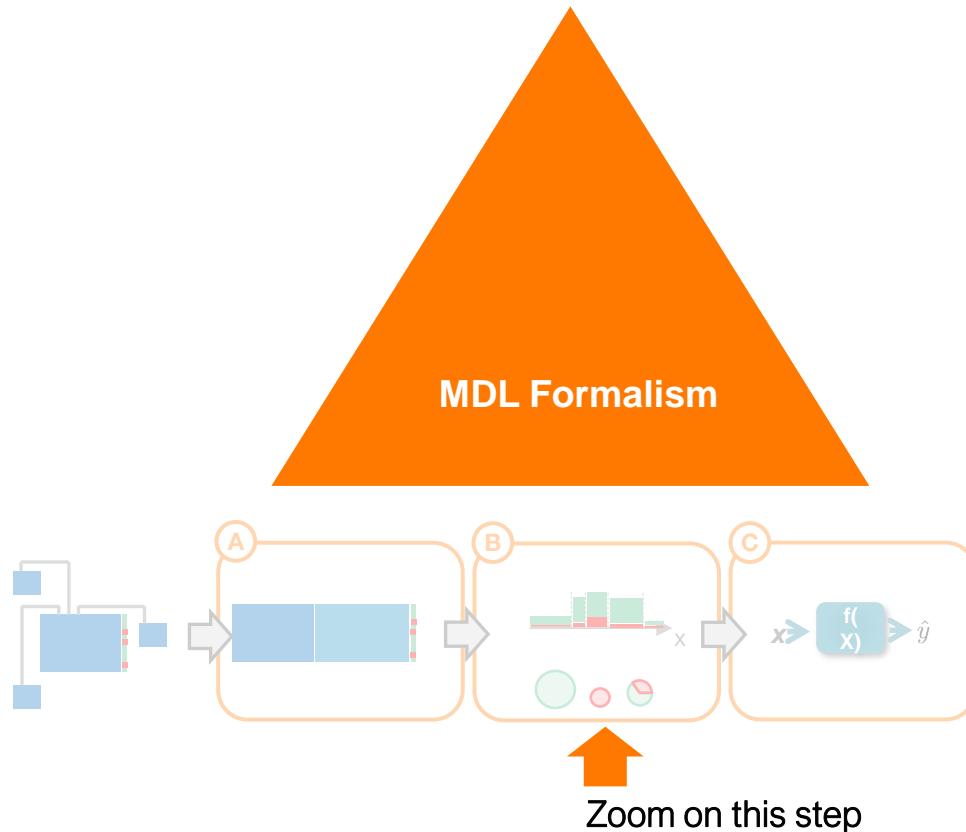
Khiops take advantage of unique formalism base on MDL*



* Minimum Description Length

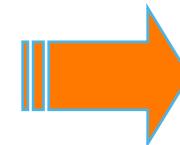
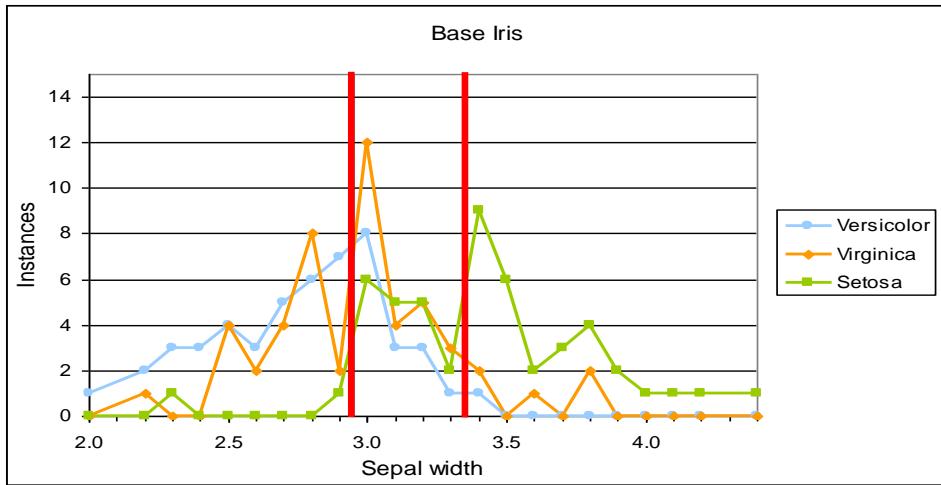
MDL is used end to end !

Optimal Encoding

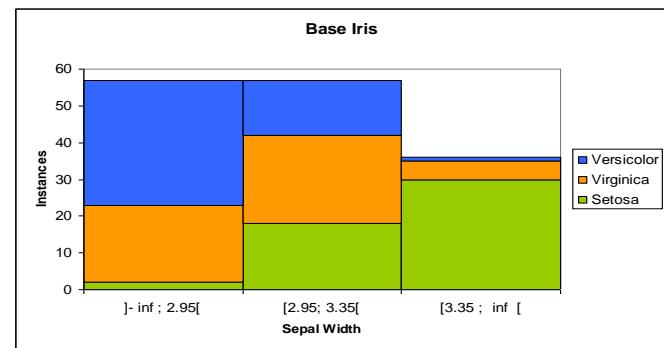


Formalisme MDL appliqué à la discréétisation

Étape B du pipeline : prétraitement univarié



Versicolor	34	15	1
Virginica	21	24	5
Setosa	2	18	30
Sepal width]- inf ; 2.95[[2.95; 3.35[[3.35 ; inf [



Quel est le meilleur modèle ?

Formalisme MDL appliqué à la discréétisation

Une approche bayésienne de sélection de modèle

Sélection du modèle le plus probable connaissant les données:

$$\text{ArgMax } P(M|D) = \frac{\text{Prior } P(M) \text{ Likelihood } P(D|M)}{P(D)}$$



Formalisme MDL appliqué à la discréétisation

Un critère d'optimisation régularisé

$$\mathcal{C}(M) = \log N + \log \binom{N+I-1}{I-1} + \sum_{i=1} \log \binom{N_i+J-1}{J-1} + \sum_{i=1} \log \frac{N_i!}{N_{i1}!N_{i2}!\dots N_{iJ}!}$$

Choix du nombre d'intervalles

Choix de la position des bornes

Description de la distribution des classes

Vraisemblance des données connaissant le modèle



Régularisation intrinsèque

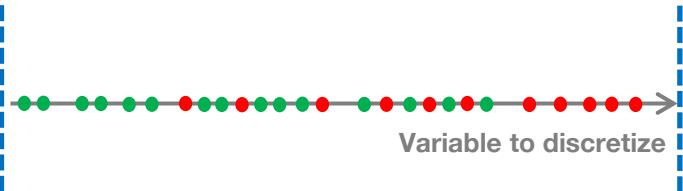
Formalisme MDL appliqué à la discréétisation

Équilibre entre *prior* et *vraisemblance* pour lutter contre le sur-apprentissage

$$\log N + \log \binom{N+I-1}{I-1} + \sum_{i=1} \log \binom{N_i+J-1}{J-1}$$

Optimisation du « Prior » seul

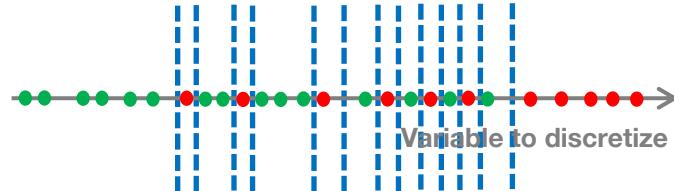
A single interval



$$\sum_{i=1} \log \frac{N_i!}{N_{i1}!N_{i2}!\dots N_{iJ}!}$$

Optimisation de la « vraisemblance » seule

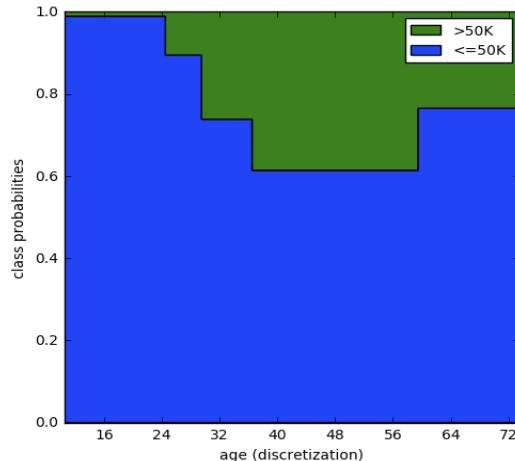
On interval for each « pure » zone



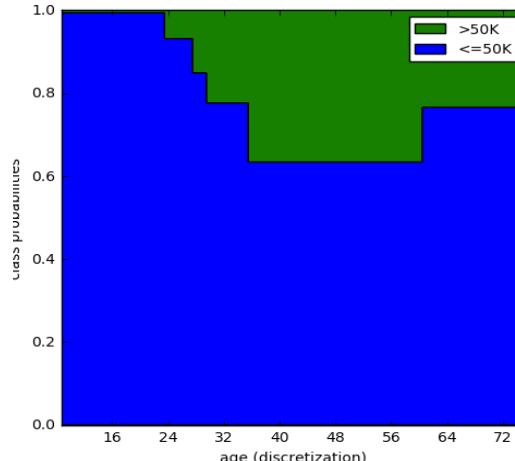
Formalisme MDL appliqué à la discréétisation

Le point d'équilibre dépend de la taille du jeu de données

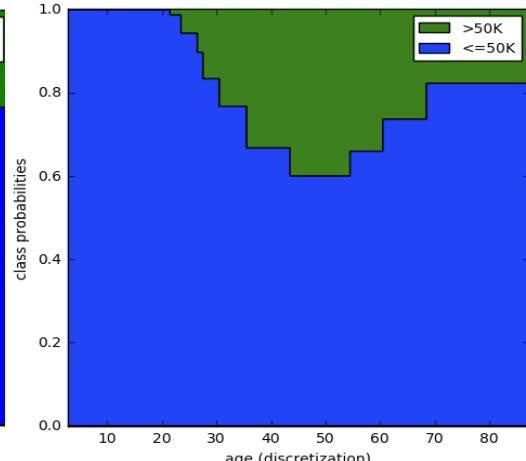
L'exemple de la variable « Age » du dataset « Adult »



$N / 2$



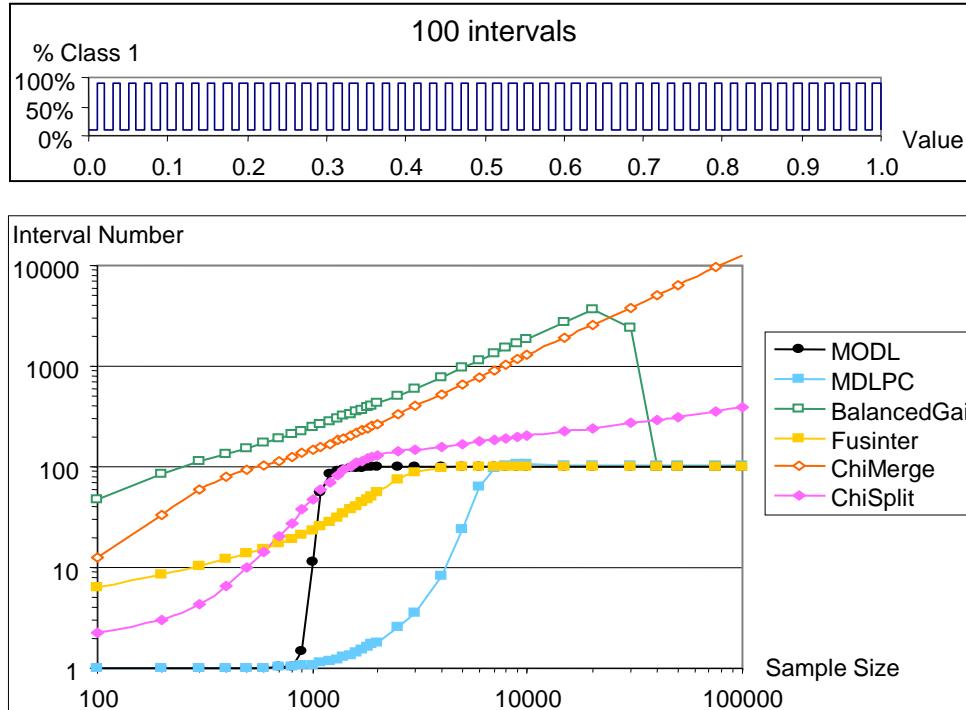
N



$N \times 2$

Univariate preprocessing : optimization criterion for the discretization

The example of the crenel pattern



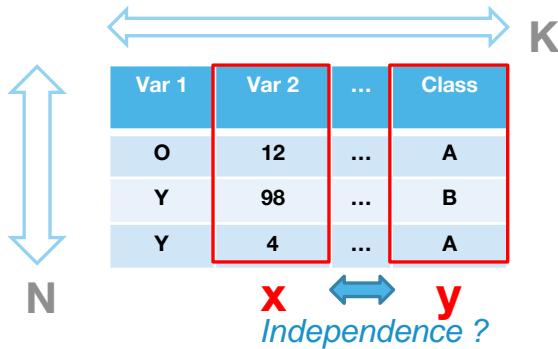
MODL correctly identifies the relevant information with a minimal number of instances

Variable Selection : compression gain based filtering

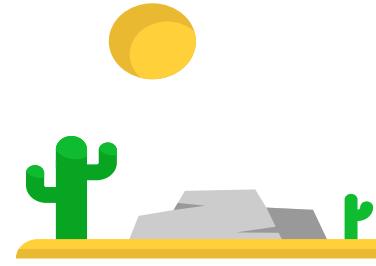


Why we want to select variables ?

1 – Scaling of the learning algorithms



2 – Accuracy of the learned models



Curse of dimensionality

In the literature, there are two families of variable selection approaches :

- **envelope** : The learned classifier is used to select the useful variable (ex: threshold on features importance)
- **filter** : Variables are selected before learning the classifier

Variable Selection : compression gain based filtering



What is compression gain ? (or level in the Khiops interface)

$$CG = 1 - \frac{-\log(P(M | D))}{-\log(P(M_0 | D))}$$

Model to be evaluated

Null model that includes a single interval / group

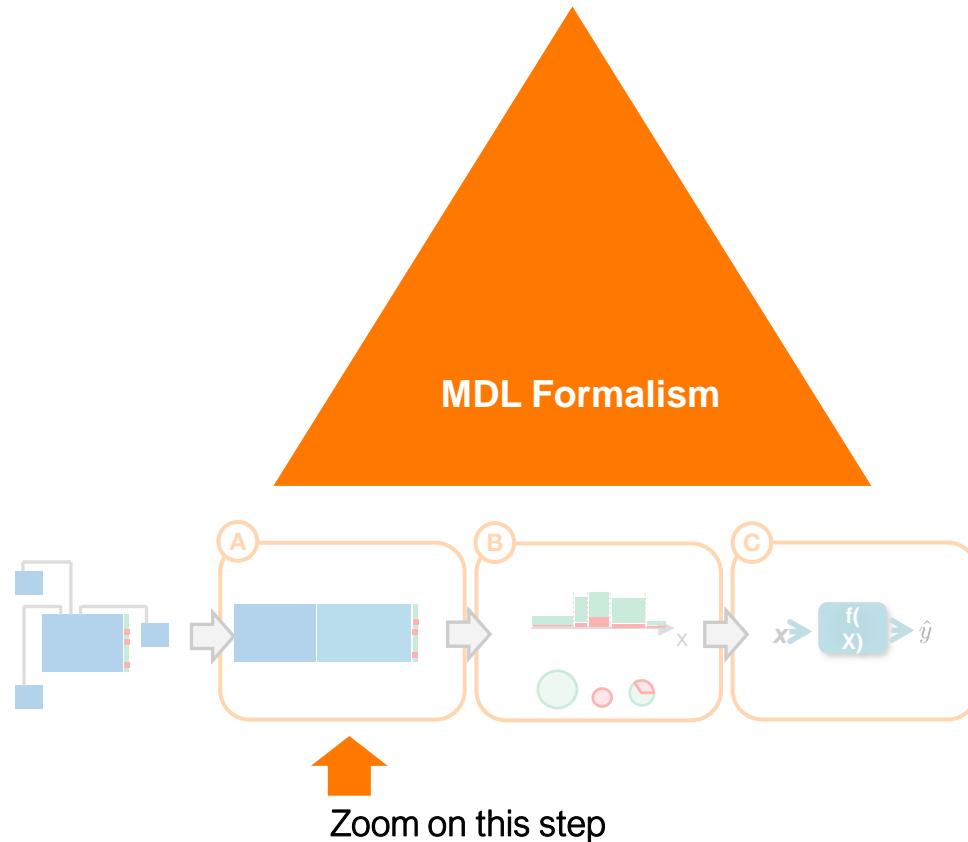
Ratio of coding length

- **Filtering method** : We keep only the variables with a positive compression gain
- **Intuition** : We remove variables for which the optimal model is less probable than the null model

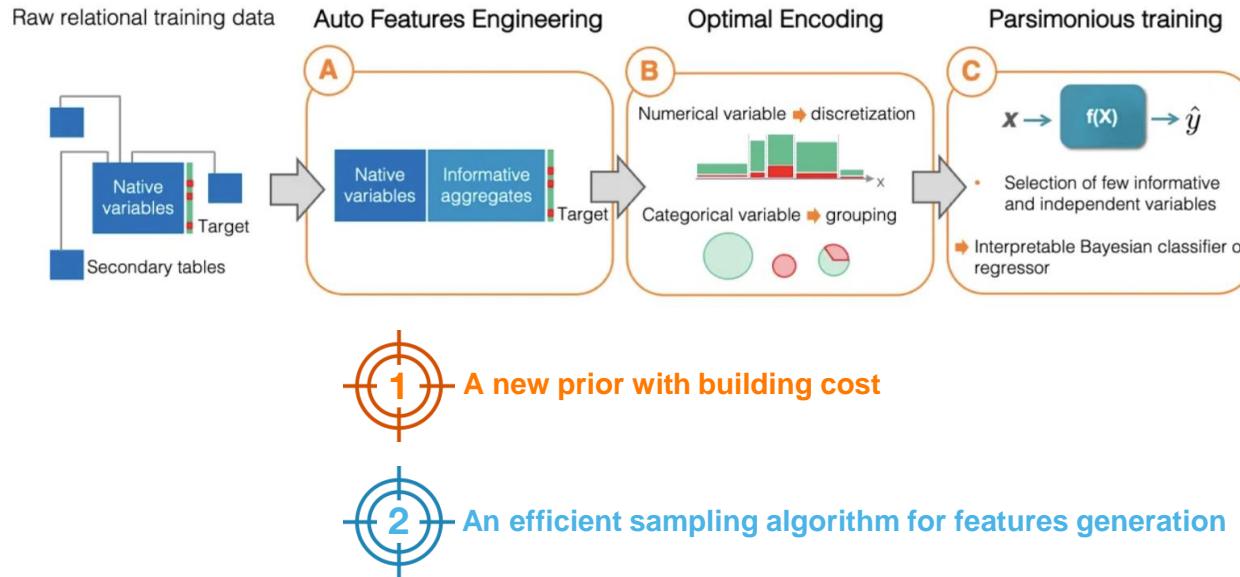


MDL is used end to end !

Auto Features Engineering

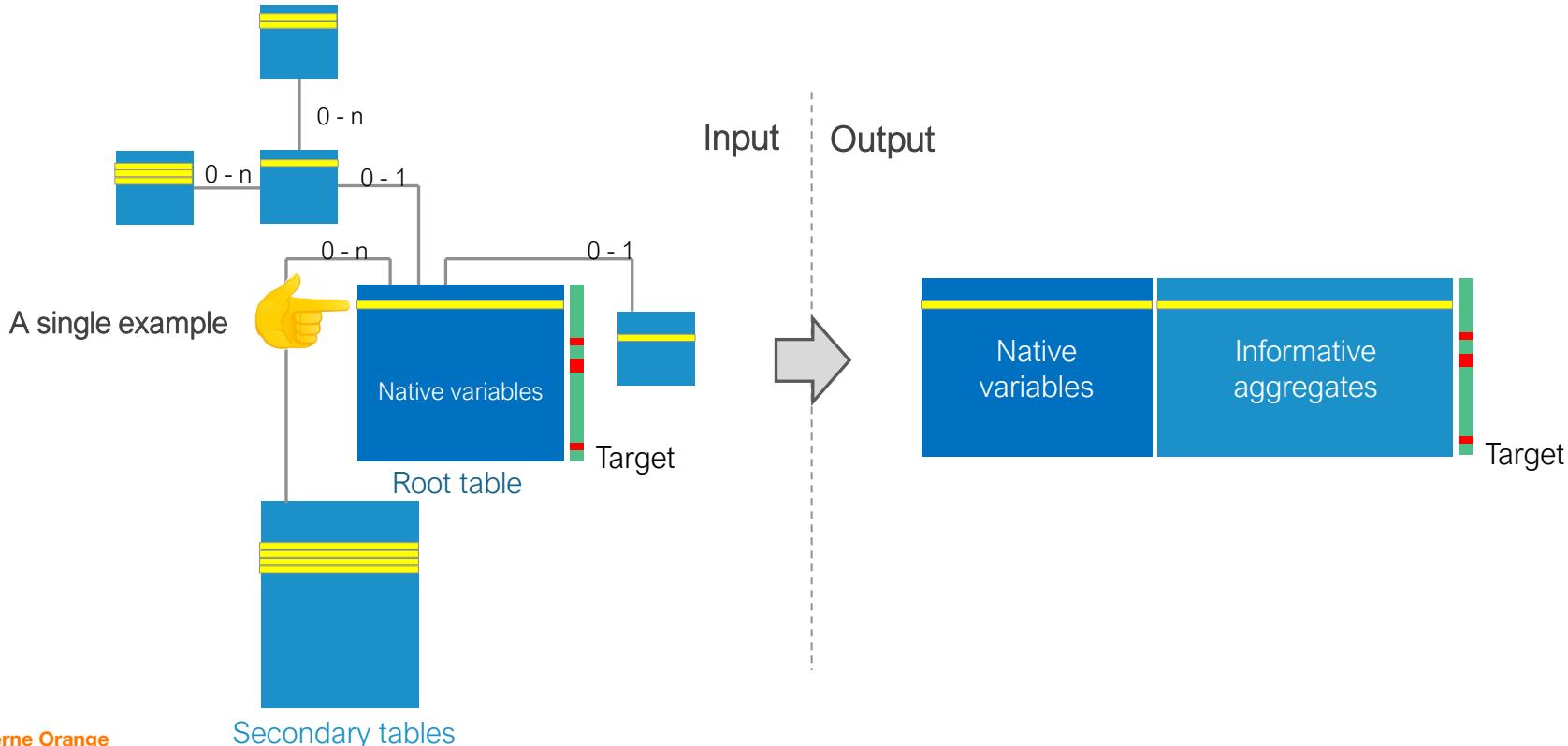


MODL and Auto Features Engineering



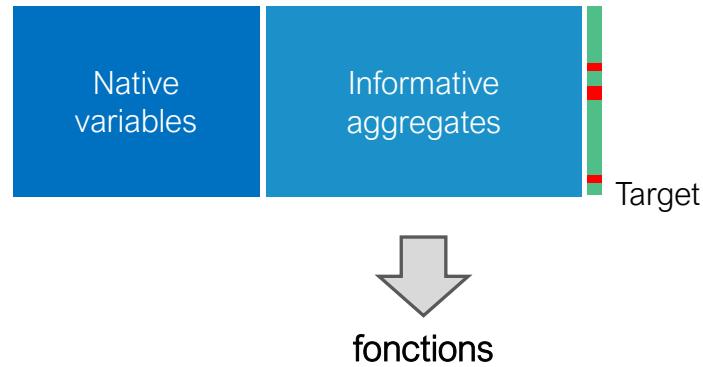
Auto Feature Engineering save time

Avoids time consuming of data-scientists



Auto Feature Engineering generate important productivity gain

Aggregates generation



Regularization:
Aggregates simpler are
more probable



Aggregates are interpretable, here some examples:

1. $\text{Mode}(\text{Usages}, \text{Product})$
2. $\text{Max}(\text{Usages}, \text{YearDay}(\text{useDate}))$
3. $\text{Count}(\text{Selection}(\text{Usages}, \text{YearDay}(\text{useDate}) \text{ in } [1;90] \text{ and Product} = \text{"VOD"})$

1 - A new prior

- Generic expression for discretization and grouping criteria:

$$c(X) = L(M_P(X)) + L(D_Y | M_P(X), D_X)$$

Coding length: describe the model encode the target by using the model

- The multi-table approach implies an additional term within the prior:


$$c(X) = L(M_C(X)) + L(M_P(X)) + L(D_Y | M_P(X), D_X)$$

construction cost
which penalizes complex aggregate features

1 - A new prior

- The new prior has a **recursive** form due to function composition ($f \circ g$)

$$L(M_C(X)) = \log(K + 1) + \log R + \sum_{op \in \mathcal{R}} L(M_C(X_{op}))$$

choice of constructing a new feature

choice of the construction rule

for each operand

Choice of constructing a new feature

- Case of original features

$$L(M_C(X)) = \log K$$

1 - A new prior

- What's appended with variable filtering ?



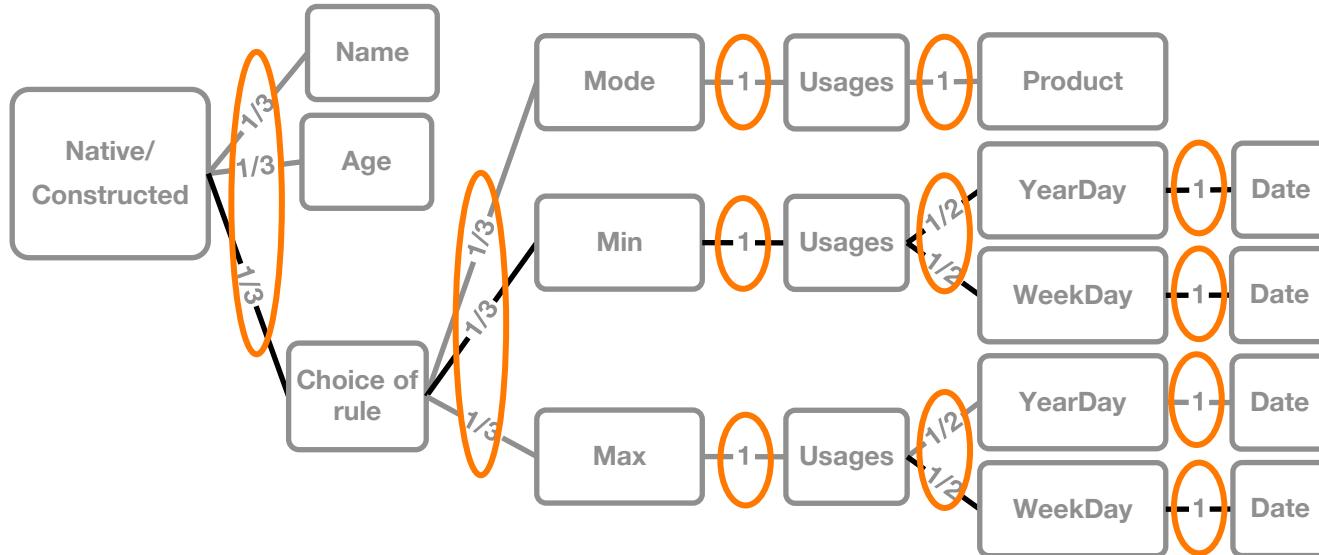
- Variables are filtered after the Auto Feature Engineering step (*univariate discretization and grouping*)
- Compression gains take into account the new prior for constructed variables



The filtering step favors **simple and informative** constructed variables ☺

2 – An efficient sampling algorithm

- **Objective :** draw K constructed / native variables from this prior distribution
- This prior can be represented by a tree structure
- And it consists of a **Hierarchy of Multinomial Distributions with potentially Infinite Depth (HMDID)**



■ Major version

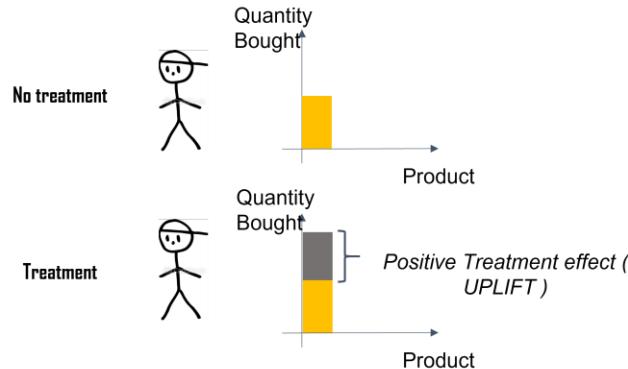
- Text data
- SNB classifier for sparse data
- Shapley
- Random forests for regression
- Khiops interpretation
- Histograms
- Coclustering instances x variables
- New visualization tools
- Simplified ergonomy

UMODL an MODL approach for Uplift

Mina Rafla, Nicolas Voisine et Bruno Cremilleux



Uplift



- Uplift modeling¹ aims at estimating the change in an individual's behavior caused by a treatment.
- Y : output variable, T : treatment and X : features
Real Uplift for each individual = $Y(T = 1) - Y(T = 0)$
- Both values cannot be known simultaneously.
- So, we consider two groups: treatment and control groups.

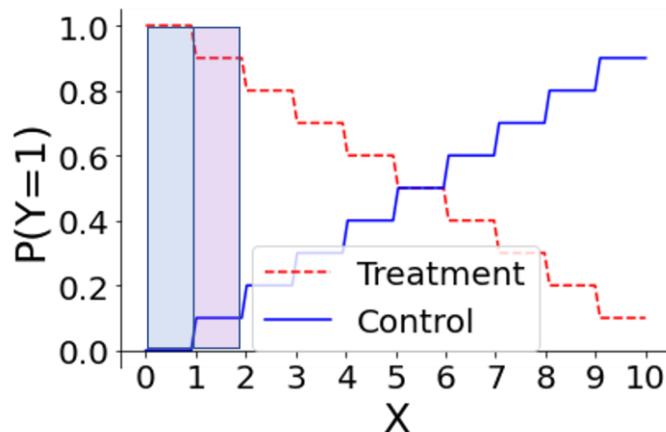
$$CATE : \hat{r}_n = E[Y_n(T = 1)|X_n] - E[Y_n(T = 0)|X_n] \quad (1)$$

Motivation

Motivation: detect areas with homogeneous treatment effect.

- To enable us to estimate uplift for each individual
- Know feature importance

Idea: creating a discretization method that separates areas with different outcome densities in treatment and control groups.



UMODL Discretization

Proposed method: UMODL a method based on the MODL approach² and designed for uplift.

M : an uplift discretization model. D : data. From a Bayes point of view:

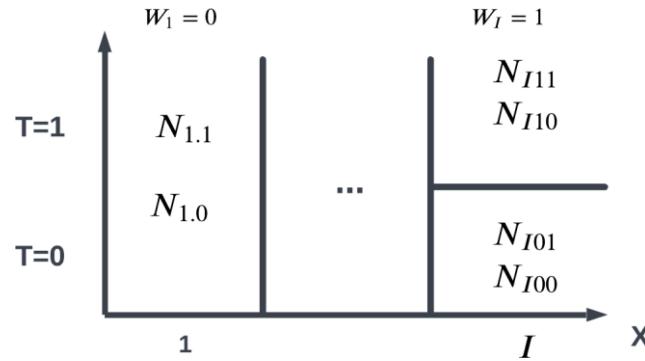
$$P(M | D) = \frac{P(M) P(D | M)}{P(D)}$$

- $P(M)$: prior
- $P(D | M)$: likelihood
- $P(D)$: constant

Objective: Maximize $P(M | D)$, hence maximize $P(M) P(D | M)$ → Tradeoff between prior and likelihood

UMODL Discretization

We define an uplift discretization model M by the hierarchy of parameters:



- I : number of intervals
- N_i : number of instances in the interval i
- W_i : boolean term indicating if the treatment has an effect in interval i ($W_i=1$) or not ($W_i=0$)
- $\{N_{i,j}\}_{W_i=0}$: number of instances in the interval i of class j
- $\{N_{ij}\}_{W_i=1}$: number of instances in the interval i of class j and the treatment t

UMODL Discretization

By exploiting the hierarchy of parameters and assuming the independence of the local distributions for each interval:

$$P(M)P(D|M) = P(I) \times P(\{N_i\}|I) \times \\ \prod_i P(W_i|I) [(1 - W_i) \times P(\{N_{ij}\}|I, \{N_j\}) + W_i \times \prod_t P(\{N_{it}\}|I, \{N_{it}\})] \times \\ \prod_i P(D_i|M)$$

A discretization model is Bayes optimal if its cost is minimal. Model's cost $C(M)$:

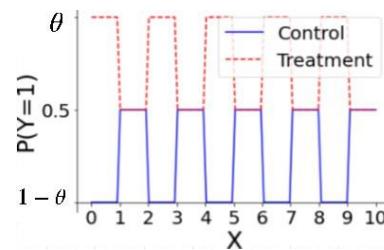
$$C(M) = -\log(P(M)P(D|M)) = \log N + \log \left(\frac{N+I-1}{I-1} \right) + I \times \log 2 \\ + \sum_{i=1}^I (1 - W_i) \log \left(\frac{N_i + J - 1}{J - 1} \right) + \sum_{i=1}^I (1 - W_i) \underbrace{\log \frac{N_i!}{N_{i1}! \dots N_{iJ}!}}_{\text{Likelihood}} \\ + \sum_{i=1}^I W_i \sum_t \log \left(\frac{N_{it} + J - 1}{J - 1} \right) + \sum_{i=1}^I W_i \underbrace{\sum_t \log \frac{N_{it}!}{N_{it1}! \dots N_{itJ}!}}_{\text{Likelihood}}$$

A greedy search algorithm is implemented to find the parameters that minimize $C(M)$.

UMODL Evaluation

Goal: Evaluate whether UMODL is a good uplift estimator

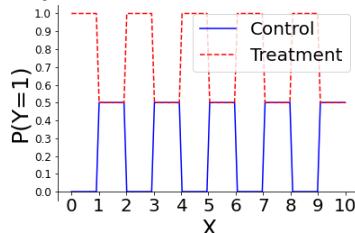
- ① Define **synthetic** data patterns, for example:



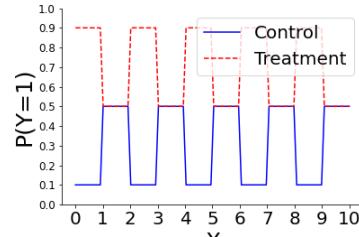
UMODL Evaluation

Goal: Evaluate whether UMODL is a good uplift estimator

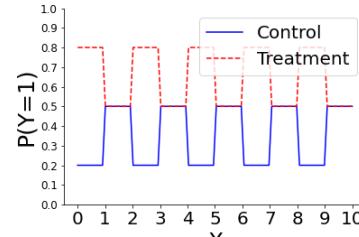
- ① Define *synthetic* data patterns, for example:



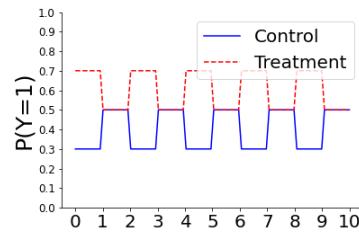
$$\theta = 1$$



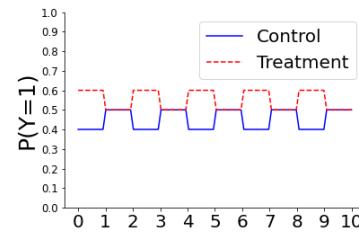
$$\theta = 0.9$$



$$\theta = 0.8$$



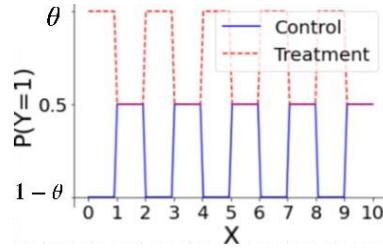
$$\theta = 0.7$$



$$\theta = 0.6$$

UMODL Evaluation : protocol

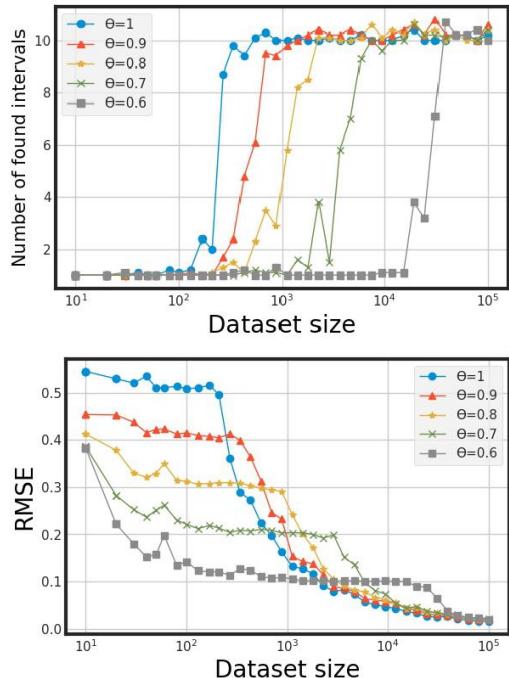
- ① Define **synthetic** data patterns, for example:



- ② Generate several datasets:
 - different sizes from 10 to 10,000 instances.
 - data is uniformly distributed on $[0, 10]$ for $T = 1$ and $T = 0$.
- ③ Generate a test set of 10K examples for the same uplift pattern.
- ④ Apply UMODL to find the Bayes optimal model.
- ⑤ Test the model by comparing the RMSE for each data point between:
 - the CATE estimation in the found interval and
 - the real CATE value.
- ⑥ Observe the number of found intervals and the RMSE, how do they behave with the data size.

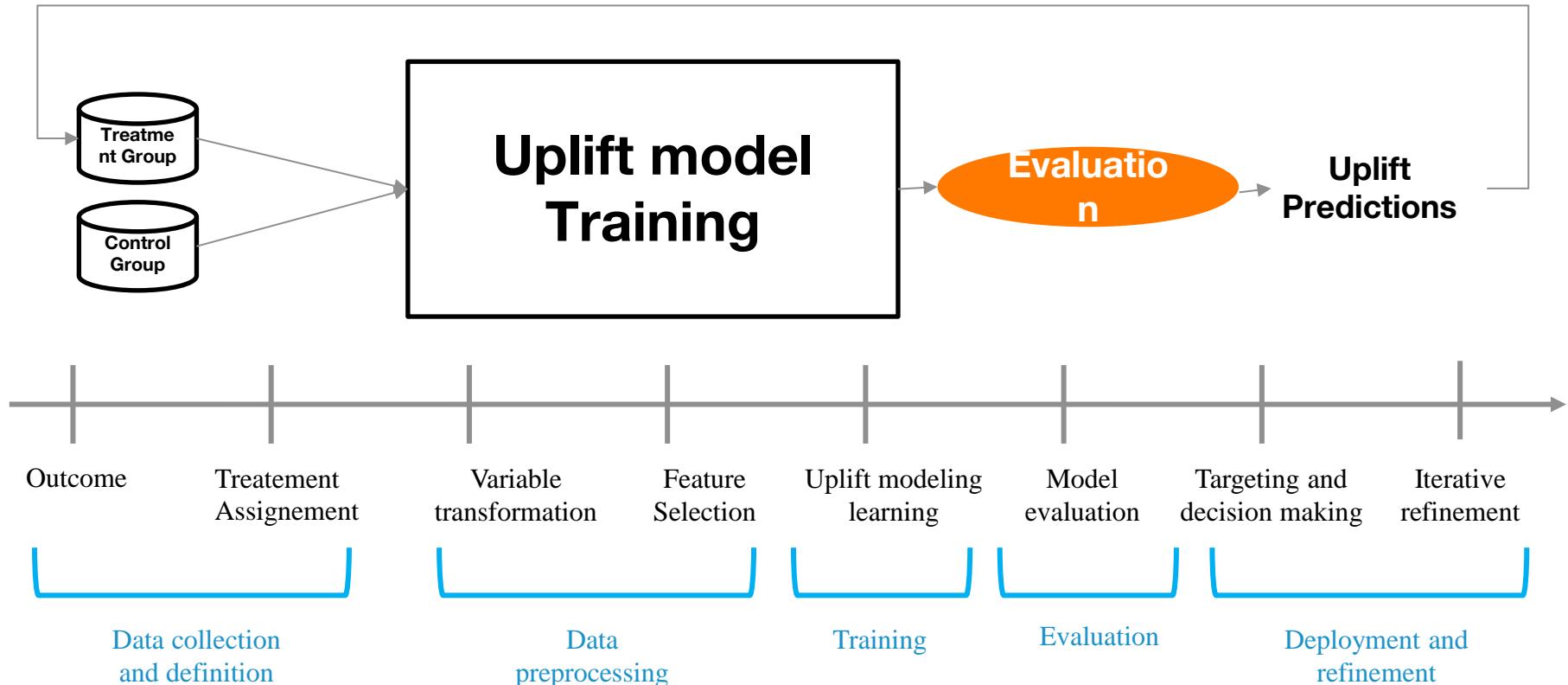
Other patterns were tested as well!

UMODL Evaluation : results



- **UMODL is a good estimator of uplift:** uplift data patterns are found (i.e correct number of intervals and a RMSE close to 0.)
- **Number of instances:** when the difference of densities between intervals gets smaller, UMODL needs more data to consider a separation to two intervals.
- **UMODL does not overfit:** when the data size increases significantly, UMODL does not consider extra intervals.

Uplift Modeling Life Cycle



1. Feature selection using UMODL-FS:

```
Import kuplift as kp  
fs = kp.FeatureSelection()  
important_vars = fs.filter(df[column_names], df["treatment"], df["outcome"])
```

2. Variable transformation using UMODL:

```
Import kuplift as kp  
ue = kp.UnivariateEncoding()  
encoded_data = ue.fit_transform(df[column_names], df["treatment"], df["outcome"])
```

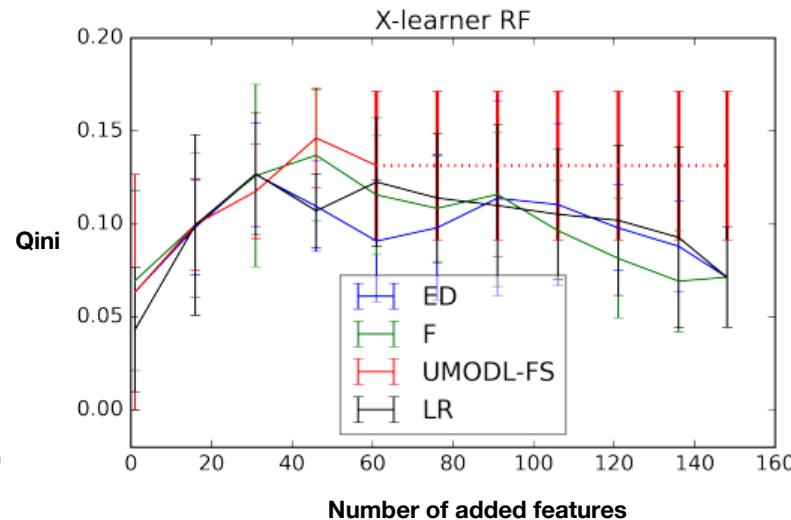
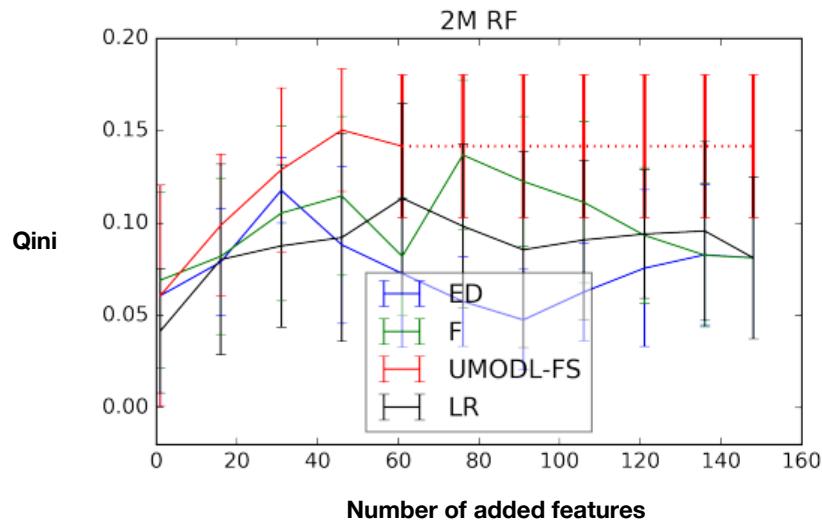
1. Decision trees using UB-DT:

```
tree = kp.BayesianDecisionTree()  
tree.fit(df[column_names], df["treatment"], df["outcome"])  
preds = tree.predict(df[column_names])
```

2. Random forests using UB-RF:

```
forest = kp.BayesianRandomForest(n_trees=4)  
forest.fit(df[column_names], df["treatment"], df["outcome"])  
preds = forest.predict(df[column_names])
```

KUPLIFT Package Python



Merci

