

# Classification et affectation des transactions aux magasins des détaillants

stage effectué chez Transaction Connect

KOMLAN JEAN-MARIE DANTODJI

Université Paris 8, LIASD  
Encadrante : Mme Rakia JAZIRI  
Tuteur : Mr Thomas MOULIN

# Plan

2/28

- 1 Introduction
- 2 Contexte
- 3 Problématique
- 4 État de l'art
- 5 Système réalisé

6 Conclusion

# Transaction Connect

3/28

- Start Up de B2B2C
- Editeur de solution numérique basé sur la donnée de paiement
- Solution déployée dans 10 pays européens et compte 60 clients

# Solutions proposées

4/28

- Transaction Connect signe des contrats avec des foncières
- Transformation de tout moyen de paiement en carte de fidélité
- Amélioration de la connaissance client aux acteurs du commerce physique
- Notification et Validation des récompenses aux clients acheteurs

# Les modes d'intégrations des Clients Business 5/28

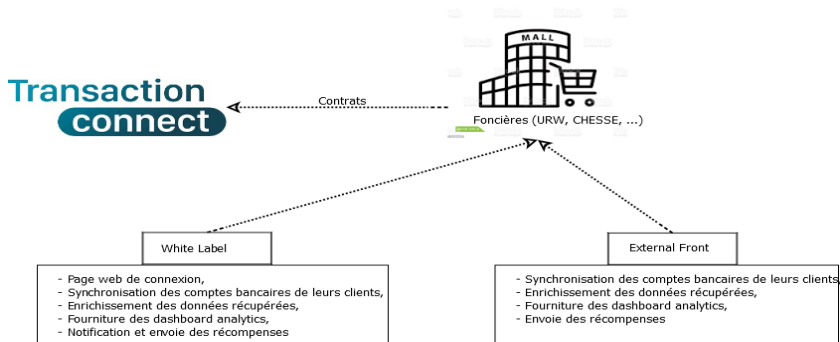


FIG. : modes d'intégration B

## 6/28



# Contexte RH

7/28

- nom du service
- composition du service
- tâches

# Contexte technique

8/28

- contexte matériel
- contexte logiciel
- contraintes



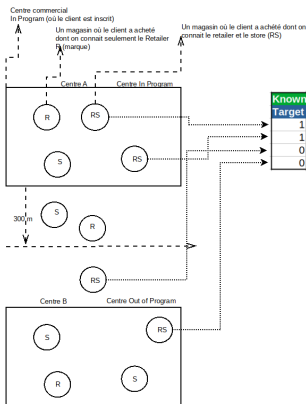
# Problème

9/28

- Classification des transactions des clients
- Affectation de transactions à un magasin

# Features Engineering

10/28



Known Stores and Retailers							Known Stores				Known Retailers		
Target	DayOfWeek	distance	distance_bin	Amount	Cannib	nb_transac	PM	SM	PNM	PC	PPM	PPNM	PPC
1	Samedi	5.1 (0,10]		10	0.166	10	1	0	1	4	1	1	3
1	Samedi	5.1 (0,10]		40	0.666	10	1	0	1	4	1	1	3
0	Samedi	5.1 (0,10]		5	0.083	10	0	0	0	4	0	0	3
0	Samedi	5.1 (0,10]		5	0.083	10	2	0	0	4	1	0	3

Pour les transactions identifiées (associées à un store et à un retailer)

Target:

Transaction In Program ou Out of Program

DayOfWeek:

jour de transaction extrait de la date de transaction

distance:

distance en (km) entre le centre commercial où le client est inscrit au store le plus proche en dehors de son centre

distance\_bin:

classe de distance qu'appartient la distance précédente (a,b]

Amount:

montant dépensé dans le store

Cannib = cannibalisation:

rapport entre le montant dépensé chez le store sur le montant total dépensé dans la journée

nb\_transac:

Nombre de transaction effectuées

PM = Purchases\_InMail

Nombre de store dans le meme centre

SM = Services\_InMail

Nombre de services utilisés (amazon locker, parking) à moins de 1km que le centre visité

PPM = Purchases\_NextToTheMail

Nombre de stores visités en dehors du centre et à moins de 300m

PC = Purchases\_InCity

Nombre de store visités dans la même ville que le centre visité

PPNM = Purchases\_Potential\_InMail

Nombre de store de tous les retailers inconnu dans le meme centre

PPNM = Purchases\_Potential\_NextToTheMail

Nombre de stores de tous les retailers visités en dehors du centre et à moins de 300m

PPC = Purchases\_Potential\_InCity

Nombre de store de tous les retailers visités dans la même ville que le centre visité

# Jeu de données

11/28

target	DayOfWeek	amount	distance	distance_bin	nb_transac	Purchases_InMall	Purchases_NextToTheMall	Services_InMall	Purchases_InCity	Purch
0	1	Monday	39.90	294.316896	(200.0, 500.0]	4	1	1	0	0
1	1	Monday	64.94	5.605244	(5.0, 10.0]	4	1	1	0	1
2	1	Monday	35.98	294.114634	(200.0, 500.0]	2	1	1	0	1
3	1	Monday	75.90	5.458404	(5.0, 10.0]	2	1	1	0	1
4	1	Saturday	118.96	5.627760	(5.0, 10.0]	6	4	4	0	4
5	1	Saturday	22.90	5.433361	(5.0, 10.0]	6	4	5	0	4
6	1	Saturday	10.00	294.203458	(200.0, 500.0]	6	4	4	0	3
7	1	Saturday	29.00	5.600936	(5.0, 10.0]	6	4	4	0	3
8	1	Saturday	16.90	5.458404	(5.0, 10.0]	6	4	4	0	3
9	1	Saturday	30.00	294.206739	(200.0, 500.0]	5	3	2	0	2
10	1	Saturday	13.98	473.669856	(200.0, 500.0]	5	3	2	0	3
11	1	Saturday	13.98	294.252744	(200.0, 500.0]	5	3	2	0	2
12	1	Saturday	18.00	294.247216	(200.0, 500.0]	5	3	3	0	2
13	0	Saturday	19.00	5.620421	(5.0, 10.0]	5	4	3	0	3
14	1	Thursday	39.98	294.152306	(200.0, 500.0]	3	2	2	0	2

FIG. : Transactions considérées

# Jeu de données

12/28

	Purchases_Potential_InMall	Purchases_Potential_InCity	Purchases_Potential_NextToTheMall	cannibalisation
0	0.0	0.0	0.0	0.593874
1	0.0	0.0	0.0	0.418474
1	0.0	0.0	0.0	0.484127
1	0.0	0.0	0.0	0.376787
1	0.0	0.0	0.0	0.070770
1	0.0	0.0	0.0	0.355689
3	0.0	0.0	0.0	0.569227
3	0.0	0.0	0.0	0.222634
3	0.0	0.0	0.0	0.375748
2	0.0	0.0	0.0	0.379135
3	0.0	0.0	0.0	0.661259
2	0.0	0.0	0.0	0.252171
2	0.0	0.0	0.0	0.934046
3	0.0	0.0	0.0	0.210076
2	0.0	0.0	0.0	0.528370

FIG. : Transactions considérées

# Données

13/28

- 483.725 lignes, 14 colonnes
- Données calculées grace au scoring des transactions

## Catégories d'algorithmes utilisés

14/28

- Support Vector Machine
- Decision Tree
- Random Forest
- K Nearest Neighbor
- Gradient boosting (XGBoost)
- Regression Logistique
- Naive Bayes

# Support Vector Machine

15/28

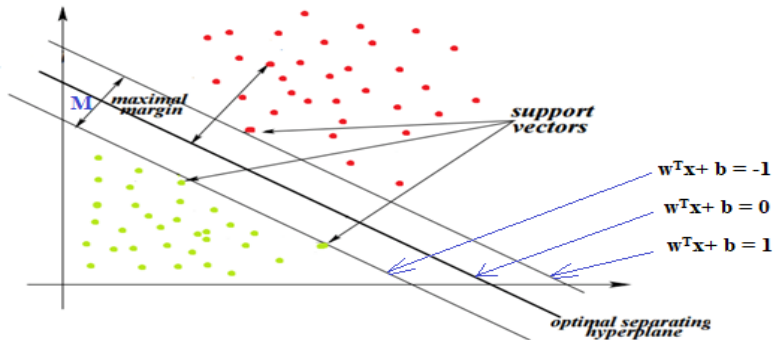


FIG. : Détermination de l'hyperplan

# SVM : Détermination d'hyperplan

16/28

$x_0$  et  $x_1$  deux vecteurs supports aux deux extrémités,  
Soit l' hyperplant

$$(P) : w^T x + b = 0$$

$$\begin{aligned} M = d(x_0, P) + d(x_1, P) &= \frac{|w^T x_0 + b|}{\sqrt{w^T w}} + \frac{|w^T x_1 + b|}{\sqrt{w^T w}} \\ &= \frac{|1|}{\sqrt{w^T w}} + \frac{|-1|}{\sqrt{w^T w}} = \frac{2}{\sqrt{w^T w}} \end{aligned}$$

Maximiser M revient à minimiser

$$\frac{\sqrt{w^T w}}{2} = \frac{\|w\|}{2}$$



# Arbre de décision

17/28

Time	Rain	Walk
30	1	No
15	1	No
5	1	No
10	0	No
5	0	No
15	0	Yes
20	0	Yes
25	0	Yes
30	0	Yes
30	0	Yes

Best feature: **Time**  
Threshold: **[5,10,15,20,25,30]**  
Best Split: **Time > 10**

Rain = 1 ?

Best feature: **Rain**  
Threshold: **[0, 1]**  
Best Split: **Rain = 1**

Time > 10 ?

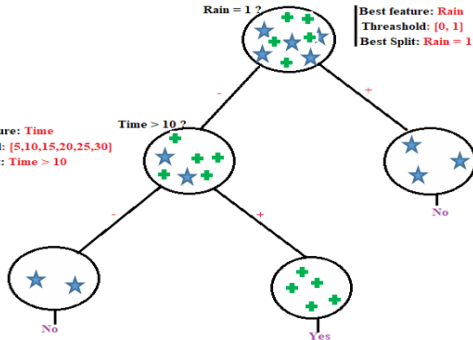


FIG. : Arbre de décision

# Arbre de décision

18/28

Soit  $X_i(\text{label}) \in [\text{"Yes"}, \text{"No"}]$

$$\text{Posons } P(X_i) = \frac{\text{nb\_label\_i\_in\_node}}{\text{total\_population}}$$

$$\text{Pour Entropie : } E = - \sum_{i=0}^{\text{nb\_labels}} P(X_i) * \log_2(P(X_i))$$

$$\text{Pour Gini : } G = 1 - \sum_{i=0}^{\text{nb\_labels}} P(X_i)^2$$

# Arbre de décision

19/28

Déterminer la meilleure variable et coupure qui correspond au  $\text{Max}(IG)$  :

$$IG = E(\text{parent}) - \sum_{i=0}^{nb\_childs} \frac{\text{total\_population\_in\_node}}{\text{total\_population}} E(\text{child\_}i)$$

$$IG = G(\text{parent}) - \sum_{i=0}^{nb\_childs} \frac{\text{total\_population\_in\_node}}{\text{total\_population}} G(\text{child\_}i)$$

# Forêt aléatoire

20/28

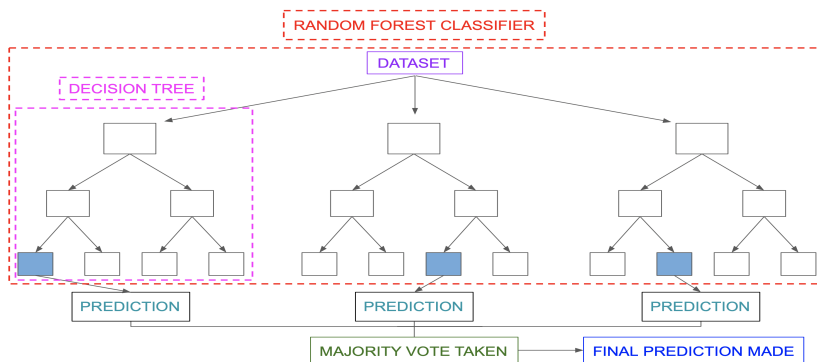


FIG. : Foret aléatoire

# Fonctionnement du K-NN

21/28

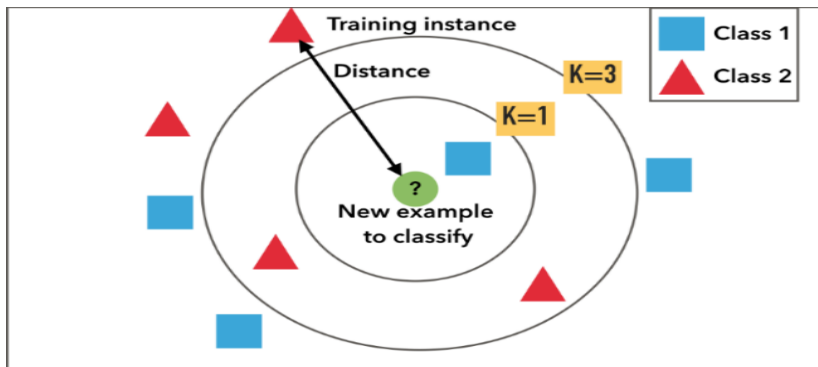


FIG. : K-Nearest Neighbor

# Les types de distances

22/28

- Distance euclidienne

$$d(A, X) = \sqrt{\sum_{i=1}^n (a_i - x_i)^2}$$

- Distance de Manhattan

$$d(A, X) = \sum_{i=1}^n |a_i - x_i|$$

- Distance de Minkowski

$$d(A, X) = \sqrt[p]{\sum_{i=1}^n |a_i - x_i|^p}$$

# Choix du paramètre K

23/28

- Utilisation de K

$$K = \sqrt{\text{nombre} - \text{de} - \text{donnees}}$$

- Choisir K suivant celui qui donne une meilleure prédiction

# Informations données

24/28

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 483725 entries, 0 to 483724
Data columns (total 14 columns):
target                                483725 non-null int64
DayOfWeek                            483725 non-null object
amount                               483725 non-null float64
distance                             483725 non-null float64
distance_bin                          483725 non-null object
nb_transac                           483725 non-null int64
Purchases_InMall                     483725 non-null int64
Purchases_NextToTheMall              483725 non-null int64
Services_InMall                      483725 non-null int64
Purchases_InCity                     483725 non-null int64
Purchases_Potential_InMall           483725 non-null float64
Purchases_Potential_InCity           483725 non-null float64
Purchases_Potential_NextToTheMall    483725 non-null float64
cannibalisation                      483725 non-null float64
dtypes: float64(6), int64(6), object(2)
memory usage: 51.7+ MB
```

FIG. : Les types de features



Ici une conclusion qui met en valeur votre travail et indique ce qui reste à faire

## Références

26/28

- ▶ Yingjie Tian, Yong Shi, Xiaohui Liu. RECENT ADVANCES ON SUPPORT VECTOR MACHINES RESEARCH. in TECHNOLOGICAL AND ECONOMIC DEVELOPMENT OF ECONOM, 2012 Volume 18(1) : 5–33
- ▶ Jehad Ali, Rehanullah Khan, Nasir Ahmad, Imran Maqsood. Random Forest and Decision Tree. In IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 5, No 3, September 2012 ISSN (Online) : 1694-0814

## Références

27/28

- ▶ Gongde Guo, Hui Wang, David Bell, Yaxin Bi, and Kieran Greer. KNN Model-Based Approach in Classification. In School of Computing and Mathematics, University of Ulster Newtownabbey, BT37 0QB, Northern Ireland, UK
- ▶ Ramraj S, Nishant Uzir, Sunil R and Shatadeep Banerjee. Experimenting XGBoost Algorithm for Prediction and Classification of Different Datasets. In International Journal of Control Theory and Applications ISSN : 0974–5572 International Science Press Volume 9 ■ Number 40 , 2016

## Références

28/28

- C. Mitchell Dayton. LOGISTIC REGRESSION ANALYSIS. Department of Measurement, Statistics and Evaluation. In Room 1230D Benjamin Building University of Maryland September 1992

Merci pour votre attention