



EA Stephen
Data scientist

Segmentez des clients d'un site e-commerce

olist

SOMMAIRE

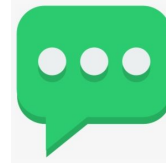
- Présentation de la problématique, du cleaning effectué, du feature engineering et de l'exploration
- Présentation des différentes pistes de modélisation effectuées et du modèle final sélectionné
- Présentation de la simulation pour définir le délai de maintenance du modèle (contrat de maintenance)



Olist (solution de vente sur les marketplaces en ligne souhaite fournir à ses équipes d'e-commerce)

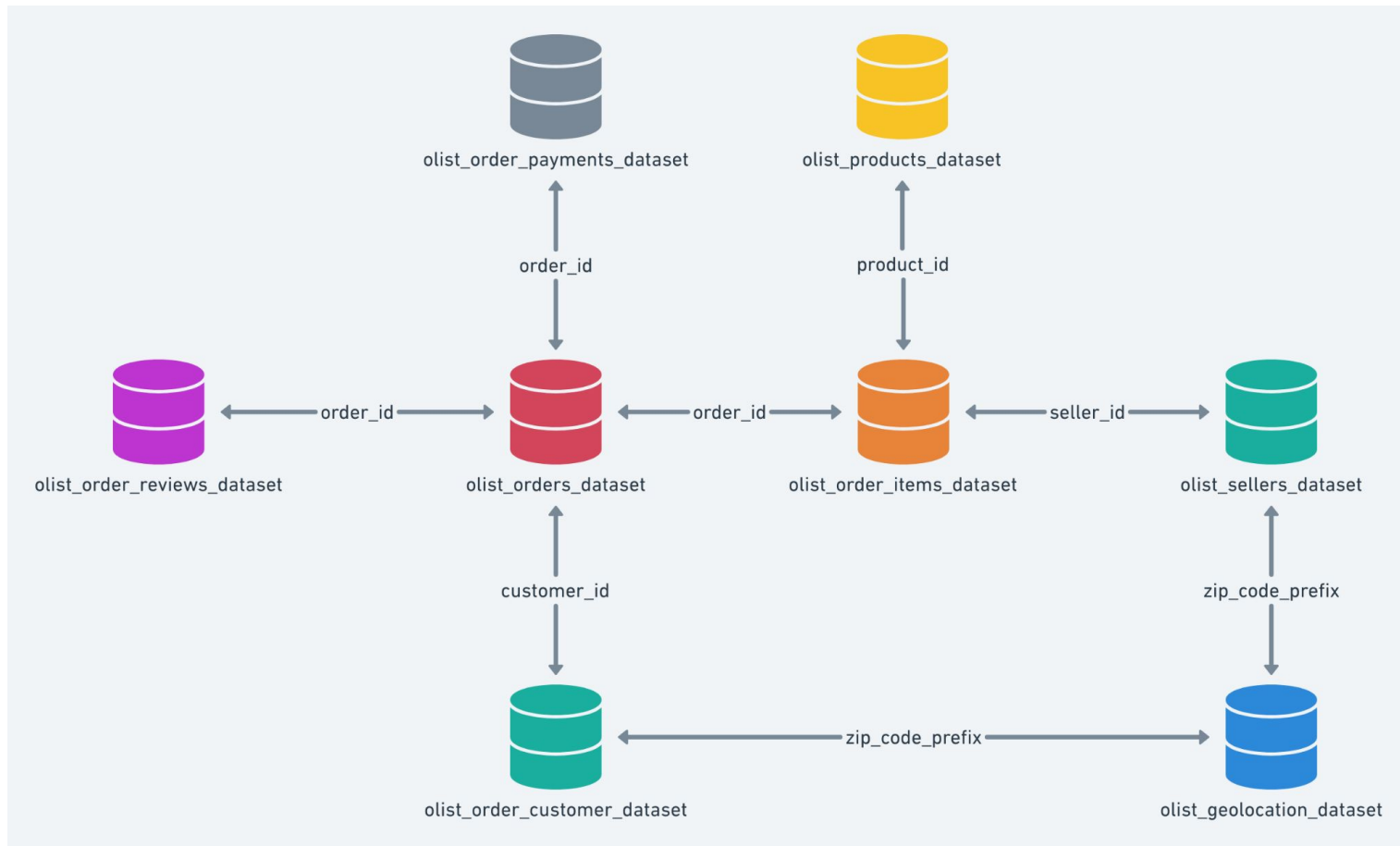


Une segmentation des clients pour leurs campagnes de communication.
La segmentation proposée doit être exploitable et facile d'utilisation pour l'équipe marketing



MAINTENANCE

Une recommandation de fréquence à laquelle la segmentation doit être mise à jour pour rester pertinente, afin de pouvoir effectuer un devis de contrat de maintenance.



Le dataset : **DF_PRODUCT_CATEGORY**

Les dimensions du df_product_category : (71, 2)

Les valeurs manquantes dans chaque colonne :

	Count_NaN	%_NaN_col	Types	Total_NaN_in_dataset	%_NaN_in_dataset
product_category_name	0	0.0	object	0	0.0
product_category_name_english	0	0.0	object	0	0.0

product_category_name product_category_name_english

0 beleza_saude health_beauty

Jointure left

product_id product_category_name product_name_length product_description_length product_photos_qty product_weight_g product_length_cm product_height_cm product_width_cm

0 1e9e8ef04d0bcf14541ed26657ea517e5 perfumaria 40.0 287.0 1.0 225.0 16.0 10.0 14.0

df_product_final

32951 rows × 10 columns

Le dataset : **DF_PRODUCTS**

Les dimensions du df_products : (32951, 9)

Les valeurs manquantes dans chaque colonne :

	Count_NaN	%_NaN_col	Types	Total_NaN_in_dataset	%_NaN_in_dataset
product_category_name	610	1.85	object	2448	0.83
product_name_length	610	1.85	float64	2448	0.83
product_description_length	610	1.85	float64	2448	0.83
product_photos_qty	610	1.85	float64	2448	0.83
product_weight_g	2	0.01	float64	2448	0.83
product_length_cm	2	0.01	float64	2448	0.83
product_height_cm	2	0.01	float64	2448	0.83
product_width_cm	2	0.01	float64	2448	0.83
product_id	0	0.0	object	2448	0.83

Fusion entre customers et orders

df_customers

Shape: (99441, 5)

df_orders

Shape: 96478,8

Orders_status: delivered

Joinure inner

	customer_id	customer_unique_id	customer_zip_code_prefix	customer_city	customer_state	order_id	order_status	order_purchase_timestamp	order_approved_at	order_delivered_carrier_date	order_delivered_customer_date	order_estimated_delivery_date
0	06b8999e2fba1a1fbc88172c00ba8bc7	1	14409	franca	SP	00e7ee1b050b8499577073aeb2a297a1	delivered	2017-05-16 15:05:35	2017-05-16 15:22:12	2017-05-23 10:47:57	2017-05-25 10:35:35	2017-06-05 00:00:00
1	18955e83d3371d6b2de6b18a428ac77	2	9790	sao bernardo do campo	SP	29150127e6685892b6eab3eac79f59c7	delivered	2018-01-12 20:48:24	2018-01-12 20:58:32	2018-01-15 17:14:59	2018-01-29 12:41:19	2018-02-06 00:00:00
2	4e7b3e00288586ebd08712fd0374a03	3	1151	sao paulo	SP	b2059ed67ce144a36e2aa97d2c59ead2	delivered	2018-05-19 16:07:45	2018-05-20 16:19:10	2018-06-11 14:31:00	2018-06-14 17:58:51	2018-06-13 00:00:00
3	b2b6027bc5c5109e5294dc65358b12c3	4	8775	mogi das cruzeiras	SP	95167092359f4fe4a63112aa7306eba	delivered	2018-03-13 16:06:38	2018-03-13 17:29:19	2018-03-27 23:22:42	2018-03-28 16:04:25	2018-04-10 00:00:00
4	4f2d8ab171c80ec83647c12e35b23ad	5	13056	campinas	SP	6b7d50bd145f6fc733ceabbd7e49dd0f	delivered	2018-07-29 09:51:30	2018-07-29 10:10:09	2018-07-30 15:16:00	2018-08-09 20:55:48	2018-08-15 00:00:00
...
96473	17ddf5dd5d51696bb3d7c6291687be6f	96092	3937	sao paulo	SP	6760e20addcf0121e9d58f21f114298	delivered	2018-04-07 15:48:17	2018-04-07 16:08:45	2018-04-11 02:08:36	2018-04-13 20:06:37	2018-04-25 00:00:00
96474	e7b71a9017aa05c9a71d292d71485e8	96093	6764	taboao da serra	SP	9ec0c8947d973db4fe8dcdf1fbfa8f1b	delivered	2018-04-04 08:20:22	2018-04-04 08:35:12	2018-04-05 18:42:35	2018-04-11 18:54:45	2018-04-20 00:00:00
96475	5e28dfe12db7fb50a4b2f691faecae5e	96094	60115	fortaleza	CE	fed4434add09aef332eaa398af0656a5c	delivered	2018-04-08 20:11:50	2018-04-08 20:30:03	2018-04-09 17:52:17	2018-05-09 19:03:15	2018-05-02 00:00:00
96476	56b18e2166679b8a959d72dd06da27f9	96095	92120	canoas	RS	e31ec91cea1ecf97797787471f98a8c2	delivered	2017-11-03 21:08:33	2017-11-03 21:31:20	2017-11-06 18:24:41	2017-11-16 19:58:39	2017-12-05 00:00:00
96477	274fa071e5e17fe303b9748641082c8	96096	6703	colia	SP	28db69209a75e59f20ccb5c36a20b90	delivered	2017-12-19 14:27:23	2017-12-19 18:50:39	2017-12-21 19:17:21	2017-12-26 18:42:36	2018-01-08 00:00:00

96478 rows x 12 columns

df_customers_orders

Le nombres de clients uniques: 93358

Le nombres de références clients: 96478

Le status des commandes: ['delivered']

Fusion entre les commandes clients et le moyen de paiement

df_customers_orders
(96478, 12)

df_order_payments
(103886, 5)

Jointure left

	customer_id	customer_unique_id	customer_zip_code_prefix	customer_city	customer_state
0	06b8999c2b1a1bcb8172c00ba8bc7	1	14409	franca	SP
1	18955e63d337b9b2d9f8b18a428ac77	2	9790	sao bernardo do campo	SP
2	4e753e0c288586eb08712600374a03	3	1151	sao paulo	SP
3	b2b6027bc5d5109e5294dc6358b12c3	4	8775	mogi das cruzeiras	SP
4	42d8ab171c80ec8384f7c12e38c23ad	5	13066	campinas	SP
...
100752	17ad55d5d51896ba3d7dc921687b9f	96092	3937	sao paulo	SP
100753	e7b71a9017aa05c9a71d292d714856e8	96093	6764	taboao da serra	SP
100754	5e28d9e12db7b50ca4b29d91faecae5e	96094	60115	fortaleza	CE
100755	56b18a2166679b8a898d72d306da27f9	96095	92120	canoas	RS
100756	274fa0771e5e176c303b9749641082c8	96096	6703	cotia	SP

100757 rows x 6 columns

order_id	order_status	order_purchase_timestamp	order_approved_at	order_delivered_carrier_date	order_delivered_customer_date	order_estimated_delivery_date	payment_sequential	payment_type	payment_installments	payment_value
00c7ee1bc50b8499577073aab2a297a1	delivered	2017-05-16 15:05:35	2017-05-16 15:22:12	2017-05-23 10:47:57	2017-05-25 10:35:35	2017-05-05 00:00:00	1.0	credit_card	2.0	148.87
29150127e6858592b6eac3eac7959c7	delivered	2018-01-12 20:48:24	2018-01-12 20:58:32	2018-01-15 17:14:59	2018-01-29 12:41:19	2018-02-06 00:00:00	1.0	credit_card	8.0	335.48
32059e067ce144a36a2aa9762c99a52	delivered	2018-05-19 16:07:45	2018-05-20 16:19:10	2018-06-11 14:31:00	2018-06-14 17:58:51	2018-06-13 00:00:00	1.0	credit_card	7.0	157.73
05187092359f4e4a63112aa7306eba	delivered	2018-03-13 16:08:38	2018-03-13 17:29:19	2018-03-27 23:22:42	2018-03-28 18:04:25	2018-04-10 00:00:00	1.0	credit_card	1.0	173.30
6b7d502d14686c733c3eabab7e49d0f	delivered	2018-07-29 09:51:30	2018-07-29 10:10:09	2018-07-30 15:16:00	2018-08-09 20:55:48	2018-08-15 00:00:00	1.0	credit_card	8.0	252.25
...
8760e20addc0121e9d580f1f14296	delivered	2018-04-07 15:48:17	2018-04-07 16:08:45	2018-04-11 02:08:36	2018-04-13 20:06:37	2018-04-25 00:00:00	1.0	credit_card	6.0	88.78
9e0c08947d973d3d44e8dc11bfa81b5	delivered	2018-04-04 08:20:22	2018-04-04 08:35:12	2018-04-05 18:42:35	2018-04-11 18:54:45	2018-04-20 00:00:00	1.0	credit_card	3.0	129.06
fed4434ad09a9d532ea398e95a5c5	delivered	2018-04-08 20:11:50	2018-04-08 20:30:03	2018-04-09 17:52:17	2018-05-09 19:03:15	2018-05-02 00:00:00	1.0	credit_card	5.0	56.04
e01ec91cea1ecd9779778747198a8c2	delivered	2017-11-03 21:08:33	2017-11-03 21:31:20	2017-11-06 18:24:41	2017-11-16 19:58:39	2017-12-05 00:00:00	1.0	credit_card	2.0	711.07
28db69209a75e5920ccbc5c36a20c90	delivered	2017-12-19 14:27:23	2017-12-19 18:50:39	2017-12-21 19:17:21	2017-12-26 18:42:36	2018-01-08 00:00:00	1.0	credit_card	1.0	21.77

df_customers_orders_payments

100757 rows x 16 columns

Le nombres de clients uniques: 93358

Le nombres de références clients: 96478

Fusion entre les commandes et la satisfaction client

df_customers_orders_payments
(100757, 16)

df_order_reviews
(99224, 7)

Jointure inner

	customer_id	customer_unique_id	customer_zip_code_prefix	customer_city	customer_state
0	00b99962b7a1fcb881720b0a8c7	1	1408	franca	SP
1	18956b3d33758c2a7b16a428c77	2	9790	sao bernardo do campo	SP
2	4670a0020886a6db87121053746c3	3	1151	sao paulo	SP
3	b2b6037b5c109a2964a6358b13c3	4	8775	moji das cruzeiras	SP
4	40285ab177d5ec39f47c12a35232ad	5	13056	campinas	SP

	order_id	order_status	order_purchase_timestamp	order_approved_at	order_delivered_carrier_date	order_delivered_customer_date	order_estimated_delivery_date	payment_sequential	payment_type	payment_installments	payment_value
0	0a7ee1050d9499577073aeb3a287a1	delivered	2017-05-16 15:05:35	2017-05-16 15:22:12	2017-05-23 10:47:57	2017-05-25 10:35:35	2017-06-05 00:00:00	1.0	credit_card	2.0	146.87
1	2510121f668b98926ba03bec7958c7	delivered	2016-01-12 20:48:24	2016-01-12 20:58:32	2016-01-15 17:34:59	2016-01-29 12:41:19	2016-02-09 00:00:00	1.0	credit_card	8.0	335.46
2	92059a6870e144436a2a07d02deba02	delivered	2018-05-19 16:07:45	2018-05-20 16:19:10	2018-05-11 14:31:00	2018-06-14 17:58:51	2018-06-13 00:00:00	1.0	credit_card	7.0	157.73
3	95167923259464a63112aa7300aba	delivered	2018-03-13 18:06:38	2018-03-13 17:29:19	2018-03-27 23:22:42	2018-04-10 00:00:00	2018-04-10 00:00:00	1.0	credit_card	1.0	173.36
4	6b785db14589c7333abab7e785c3	delivered	2018-07-29 09:51:30	2018-07-29 10:10:09	2018-07-30 15:16:00	2018-08-09 20:59:48	2018-08-19 00:00:00	1.0	credit_card	8.0	252.28

	review_id	review_score	review_comment_title	review_comment_message	review_creation_date	review_answer_timestamp
0	48b653546e0f03a8d1a2136a59630b	4	NaN	NaN	2017-05-26 00:00:00	2017-05-30 22:34:40
1	02c45a9a3e3a8f1a8e295070ee03	5	NaN	NaN	2018-01-30 00:00:00	2018-02-10 22:42:29
2	3a6695d7fee186a473a27069a487	5	NaN	NaN	2018-06-15 00:00:00	2018-06-15 12:10:59
3	03ba071a0b1f5a0b226e6f72cab79c5	5	NaN	NaN	2018-03-29 00:00:00	2018-04-02 18:36:47
4	849387958d5c477792739a2098ba	5	a melhor nota	O barbaon é excelente Amo adoro o barbaon	2018-08-10 00:00:00	2018-08-17 01:59:52

10080 rows x 22 columns

df_customers_orders_payments_reviews

100650 rows x 22 columns

Le nombres de clients uniques: 92755

Le nombres de références clients: 95832

Fusion entre les produits vendus par Olist et les catégories des produits

df_order_items

(112650, 7)

df_products

(32951, 10)

Jointure left

order_id	order_item_id	product_id	seller_id	shipping_limit_date	price	freight_value	product_category_name	product_name_length	product_description_length	product_photos_qty	product_weight_g	product_length_cm	product_height_cm	product_width_cm	product_category_name_english
0	00010242f98c5a6d1ba2d0792db16214	1	4244733e0b67ecb4970a6e2683c13661	2017-09-19 09:45:35	58.90	13.29	cool_stuff	58.0	598.0	4.0	650.0	28.0	9.0	14.0	cool_stuff
1	0001877720320c507190d7a144bd03	1	e502d52b02189ee0588e5ca93d93a8f	2017-09-03 11:05:13	239.90	19.92	pet_shop	56.0	239.0	2.0	30000.0	50.0	30.0	40.0	pet_shop
2	000229ec39822c05f7a1c0657da4fc703e	1	c77735d618b72b67abbeef8d44d0d0d	2018-01-18 14:48:30	199.00	17.87	movele_decoracao	59.0	695.0	2.0	3050.0	33.0	13.0	33.0	furniture_decor
3	00024acbc0d0a6daa1e931b036114c75	1	7634da152a46101f595efa321f4722fc	2018-06-15 10:10:18	12.99	12.79	perfumaria	42.0	480.0	1.0	200.0	16.0	10.0	15.0	perfumery
4	00042926c09d7ce08d9fab64e50d4199	1	ad0c36230e8f30e0c30c43695e4e10089	2017-02-13 13:57:51	199.90	18.14	ferramentas_jardim	59.0	409.0	1.0	3790.0	35.0	40.0	30.0	garden_tools
...
112645	ff0c94f6ce00a00581880b5947a5ac37	1	4aa6014cecb682077f9dc4bffe0c0550	2018-05-02 04:11:01	299.99	43.41	utilidades_domesticas	43.0	1002.0	3.0	10150.0	89.0	15.0	40.0	housewares
112646	ff0c46ef2263404302a634e0577fab	1	32e070915822b0795ea448c4d074c828	2018-07-20 04:31:48	350.00	36.53	informatica_acessorios	31.0	232.0	1.0	8990.0	45.0	26.0	38.0	computers_accessories
112647	ff0c4705a9662cd70adb1304a31832d	1	72a30483855e2eaf067aee50c2900482	2017-10-30 17:14:25	99.90	16.95	esporte_lazer	43.0	869.0	1.0	967.0	21.0	24.0	19.0	sports_leisure
112648	ff0c418544fab395dfada21779c9644f	1	9c422a519119dcad7575db5af1ba540e	2017-08-21 00:04:32	55.99	8.72	informatica_acessorios	56.0	1306.0	1.0	100.0	20.0	20.0	20.0	computers_accessories
112649	ff0e41c64801cc87c081d081db308244	1	3506888f9dc1e75f9f7b3c2638365e01	2018-06-12 17:10:13	43.00	12.79	cama_mesa_banho	47.0	511.0	1.0	600.0	30.0	3.0	19.0	bed_bath_table

112650 rows x 16 columns

df_order_items_products

112650 rows x 16 columns

Fusion entre le client et les articles achetés

df_order_items_products

(112650, 16)

df_customers_orders_payments_reviews

(100650, 22)

Jointure left

order_id	order_item_id	product_id	seller_id	shipping_limit_date	price	freight_value	product_category_name	product_name_lenght	product_description_lenght	product_photos_qty	product_weight_g	product_length_cm	product_height_cm	product_width_cm	product_category_name_english
00010242f6c5a8d1ba2d4792cb16214	1	4244733a06e7ecb4970a9e2683c13661	48436dade18ac8b2bce08ec2a041202	2017-09-19 09:45:35	58.9	13.29	cool_stuff	58.0	598.0	4.0	650.0	28.0	9.0	14.0	cool_stuff

customer_id	customer_unique_id	customer_zip_code_prefix	customer_city	customer_state	order_id	order_status	order_purchase_timestamp	order_approved_at	order_delivered_carrier_date	order_delivered_customer_date	order_estimated_delivery_date	payment_sequential	payment_type	payment_installments	payment_value	review_id
00b8999e7ba1a1fbc68172d00ba8b07	1	14409	franca	SP	00b78a1b050b6499577073a02a297a1	delivered	2017-05-16 15:09:35	2017-05-16 15:22:12	2017-05-23 10:47:57	2017-05-25 10:35:35	2017-06-05 00:00:00	1.0	credit_card	2.0	146.87	88b8b02046f026a981a02136a59b30b

review_score	review_comment_title	review_comment_message	review_creation_date	review_answer_timestamp
4	NaN	NaN	2017-05-26 00:00:00	2017-05-30 22:34:40

df_customers_orders_payments_reviews_items_products

114862 rows × 37 columns

Le nombres de clients uniques: 92755

Le nombres de références clients: 95832

Fusion avec la géolocalisation

`df_customers_orders_payments_`
`reviews_items_products`
 (114862, 37)

`df_olist_geolocation`
 (738332, 5)

Jointure left

customer_id	customer_unique_id	customer_zip_code_prefix	customer_city	customer_state	order_id	order_status	order_purchase_timestamp	order_approved_at	order_delivered_carrier_date	order_delivered_customer_date	order_estimated_delivery_date	payment_sequential	payment_type	payment_installments	payment_value	review_id		
0	00b8999e2b0a1fbc88172c000aabc7	1	14409	franca	SP	00e7ee1b050b8498577073aeb2a297a1	delivered	2017-05-16 15:05:35	2017-05-16 15:22:12	2017-05-23 10:47:57	2017-05-25 10:35:35	2017-06-05 00:00:00	1.0	credit_card	2.0	146.87	888bb52646d026a8d1ad2136a59e30b	
review_score	review_comment_title	review_comment_message	review_creation_date	review_answer_timestamp	order_item_id	product_id	seller_id	shipping_limit_date	price	freight_value	product_category_name	product_name_length	product_description_length	product_photos_qty	product_weight_g	product_length_cm	product_height_cm	product_width_cm
4	NaN	NaN	2017-05-28 00:00:00	2017-05-30 22:34:40	1	a9516a079a37a6c3c3639c78b10169a8	7c67e14482009e99d59a5c6a8b010ab	2017-05-22 15:22:12	124.99	21.88	moveis_escritorio	41.0	1141.0	1.0	8883.0	54.0	64.0	31.0
product_category_name_english		customer_zip_code_prefix	geolocation_lat	geolocation_lng	geolocation_city	geolocation_state												
office_furniture		0	1037	-23.545621	-46.639292	sao paulo SP												

`df_final_fusion`

114862 rows × 41 columns

Le nombres de clients uniques: 92755

Le nombres de références clients: 95832

Les variables supprimés

customer id	0
customer zip code prefix	0
customer city	0
customer state	0
order status	0
order approved at	15
order delivered carrier date	2
order estimated delivery date	0
payment sequential	0
payment type	0
payment installments	0
review id	0
review comment title	101275
review comment message	0
review creation date	0
review answer timestamp	0
product id	0
seller id	0
shipping limit date	0
price	0
freight value	0
product category name	0
product name lenght	1626
product description lenght	1626
product photos qty	1626
product weight g	20
product length cm	20
product height cm	20
product width cm	20
product category name english	0
geolocation lat	302
geolocation lng	302
geolocation city	302
geolocation_state	302

Features sélections

customer unique id	0
order id	0
order item id	0
order purchase timestamp	0
order delivered customer date	8
payment value	3
review_score	0

Supprimer les valeurs manquantes

Feature engineering

weekday_order_purchase	0
month_order_purchase	0
hour_order_purchase	0
recence	0
Duration_delevery_order	0

Groupby

customer_unique_id

Agrégation

order id	count
order item id	mean
payment value	sum
review_score	mean
recence	mean
Duration_delevery_order	mean

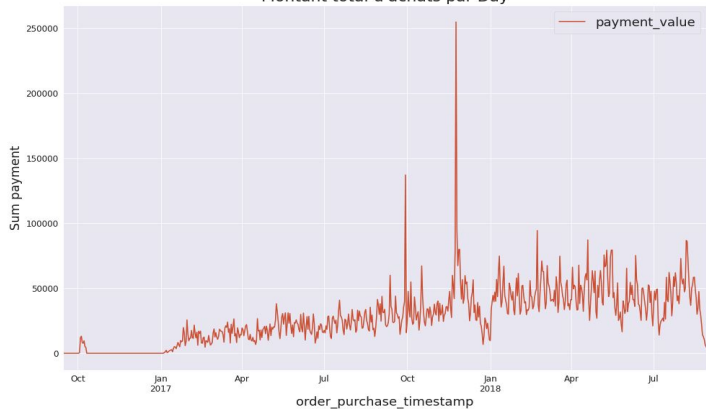
Drop outlier

Client
6186
sensibilise
le
clustering
et la
stabilité

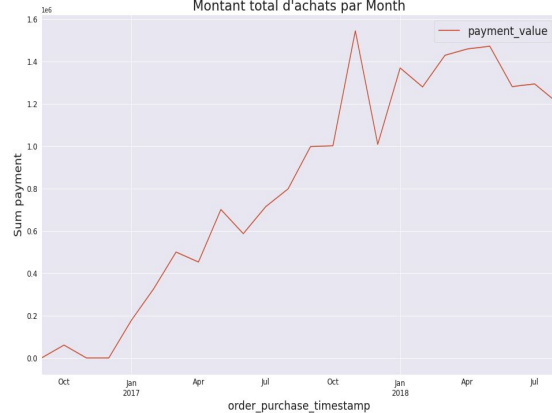
	customer_unique_id	recence	frequence	mean_order_item_id	Monnaie	review_score_mean	Duration_delevery_order_mean
0	1	470.0	1.0	1.0	146.87	4.0	9.0
1	2	229.0	1.0	1.0	335.48	5.0	17.0
2	3	102.0	1.0	1.0	157.73	5.0	26.0
3	4	169.0	1.0	1.0	173.30	5.0	15.0
4	5	31.0	1.0	1.0	252.25	5.0	11.0
...
92740	96092	144.0	1.0	1.0	88.78	4.0	6.0
92741	96093	147.0	1.0	1.0	129.06	5.0	7.0
92742	96094	143.0	1.0	1.0	56.04	1.0	31.0
92743	96095	299.0	1.0	1.0	711.07	5.0	13.0
92744	96096	253.0	1.0	1.0	21.77	5.0	7.0

92745 rows x 7 columns

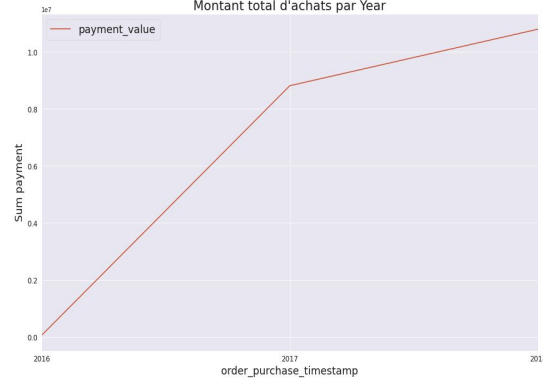
Montant total d'achats par Day



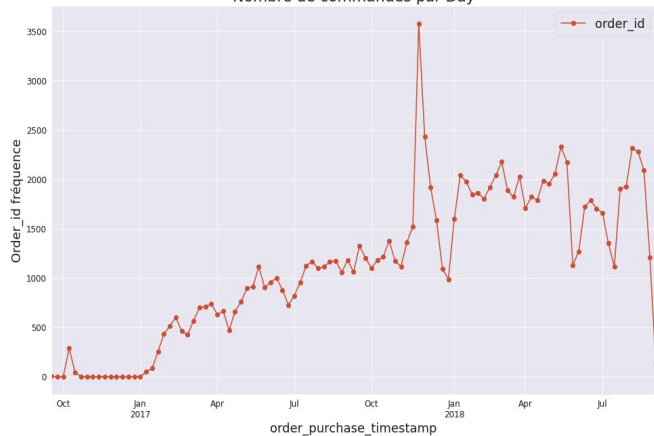
Montant total d'achats par Month



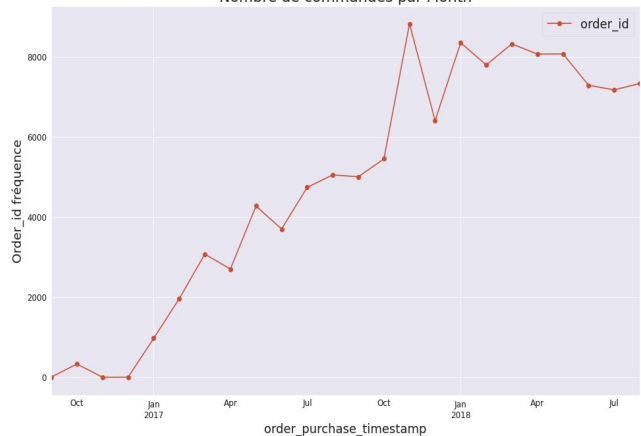
Montant total d'achats par Year



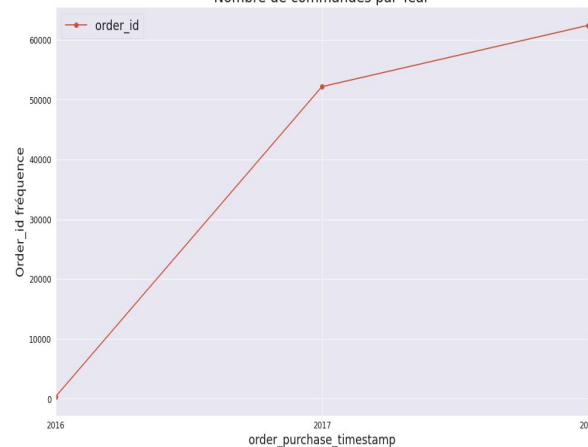
Nombre de commandes par Day



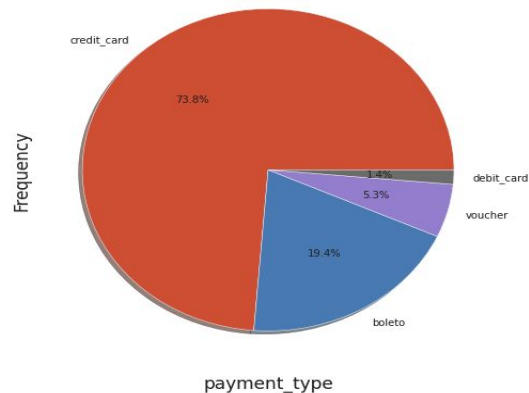
Nombre de commandes par Month



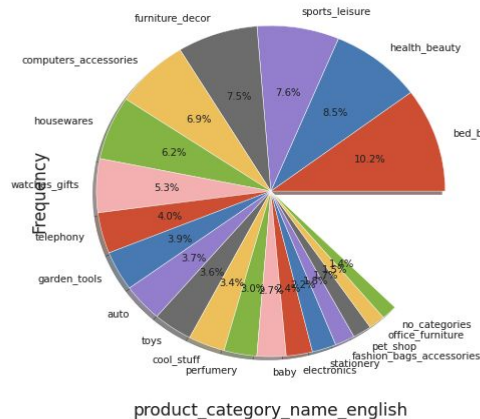
Nombre de commandes par Year



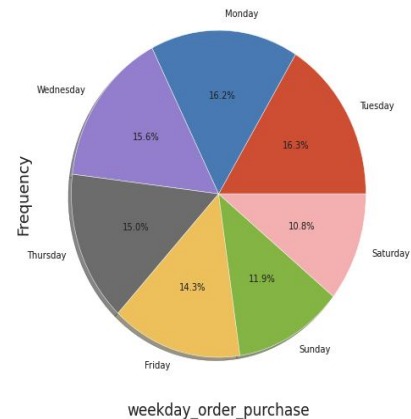
Répartition de la variable payment_type



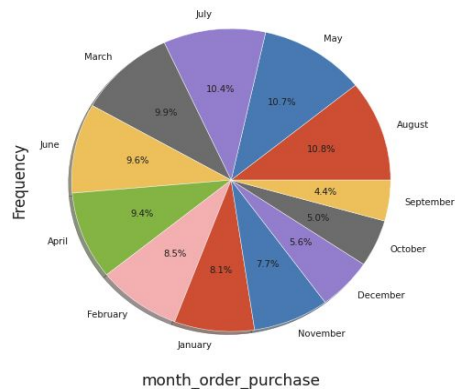
Répartition de la variable product_category_name_english



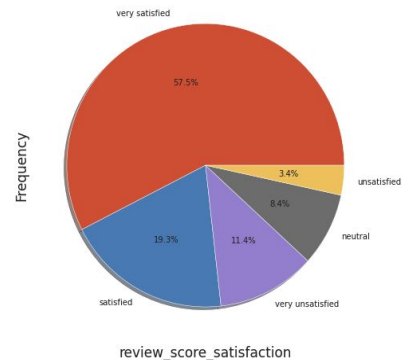
Répartition de la variable weekday_order_purchase



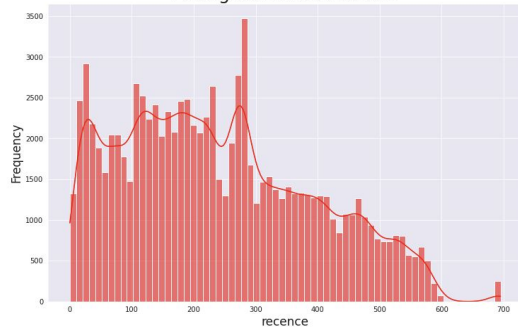
Répartition de la variable month_order_purchase



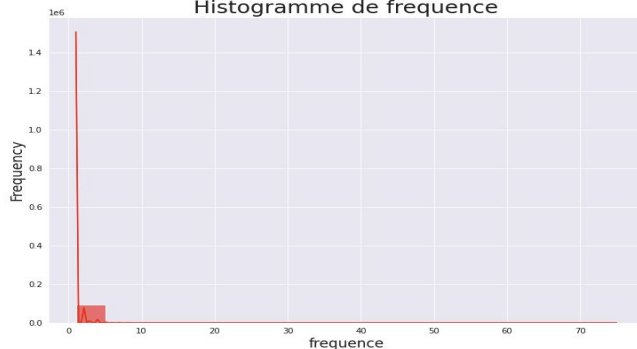
Répartition de la variable review_score_satisfaction



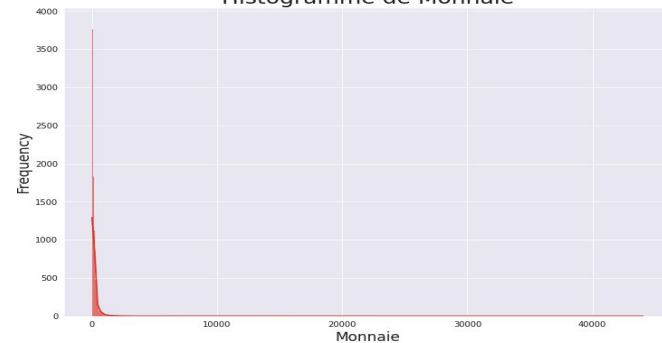
Histogramme de recence



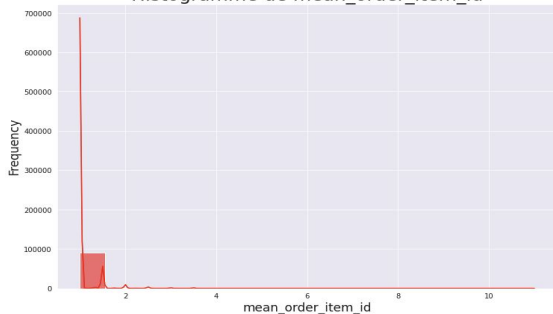
Histogramme de frequence



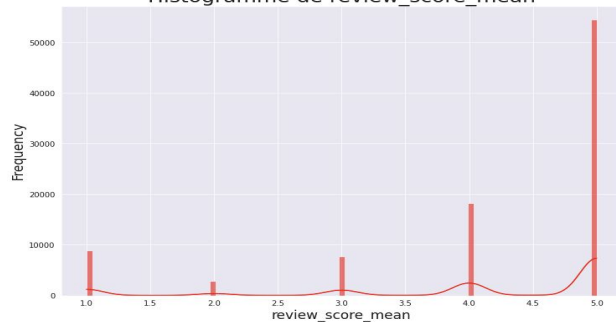
Histogramme de Monnaie



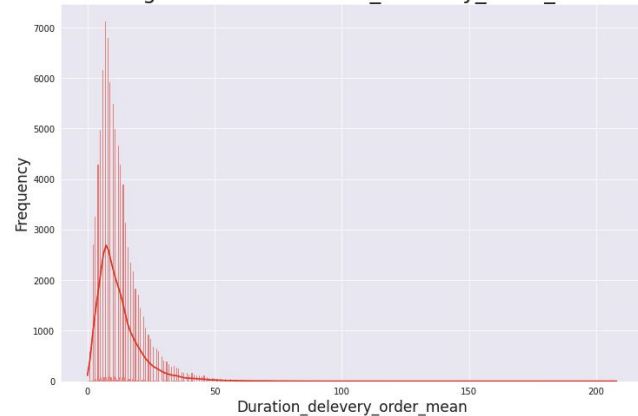
Histogramme de mean_order_item_id



Histogramme de review_score_mean



Histogramme de Duration_delevery_order_mean





Evaluation best K

n_clusters = 4

- méthode du coude , davies_bouldin_score(0.76)
- silhouette score(0.46) ,calinski_harabasz(41342)

Stabilité à l'initialisation

- Stabilité moyenne pour 10 itérations

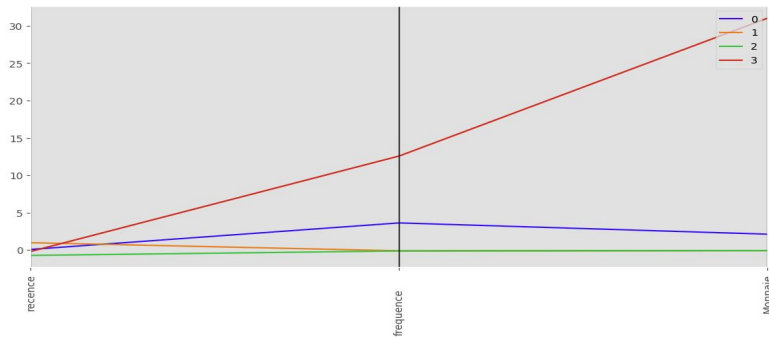
Distance

- Les clusters (1 et 2) trop similaire en terme F,M
- Cluster 3 suffisamment distant mais trop peu de clients

Effectifs des clusters

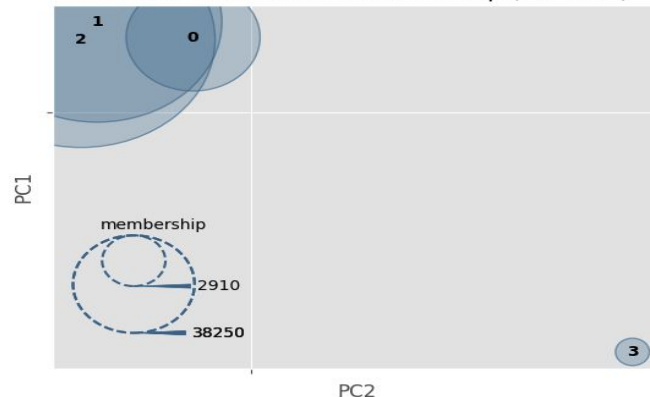
- pas homogène
- pas pertinent pour le besoin métier
- La fréquence à trop de poids dans la segmentation (97% clients ont passé 1 commande)

Parallel Coordinates plot for the Centroids



cluster 0 : 2910
 cluster 1 : 38250
 cluster 2 : 51547
 cluster 3 : 38

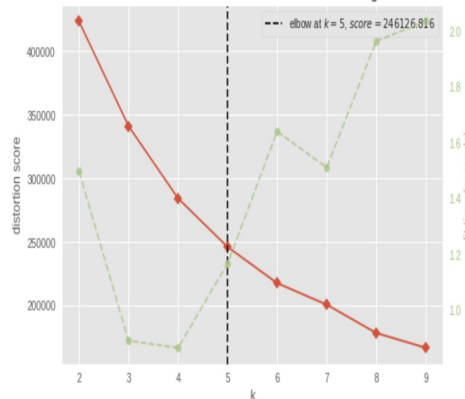
KMeans Intercluster Distance Map (via MDS)



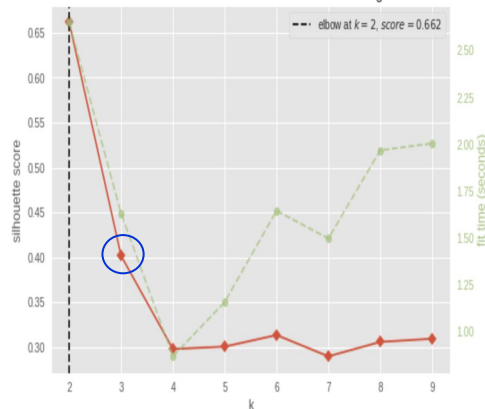
Conclusion

La segmentation RFM pas pertinente d'un point de vue métier.

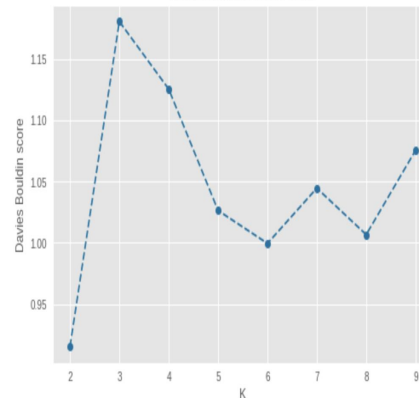
Distortion Score Elbow for KMeans Clustering



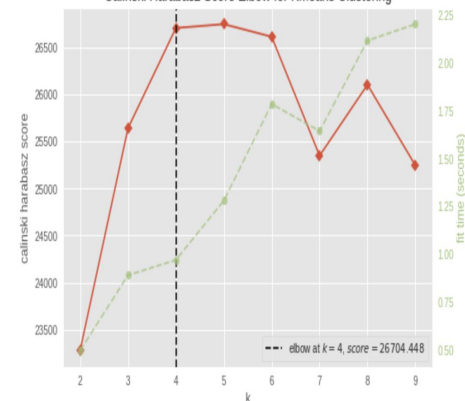
Silhouette Score Elbow for KMeans Clustering



Davies Bouldin score vs. K



Calinski Harabasz Score Elbow for KMeans Clustering



Méthode du coude (distortion)

La moyenne de la somme des carrés des distances au centroïde le proche (coude)

Silhouette score

Le coefficient de silhouette ou score de silhouette est une mesure utilisée pour calculer la qualité d'une technique de regroupement. Sa valeur est comprise entre -1 et 1.
1 : Signifie que les grappes sont bien séparées les unes des autres et clairement distinguées.
a = distance intra-cluster moyenne, c'est-à-dire la distance moyenne entre chaque point d'un cluster.

Davies Bouldin

Le score est défini comme la mesure de similarité moyenne de chaque cluster avec son cluster le plus similaire, où la similarité est le rapport des distances intra-cluster aux distances inter-cluster. Ainsi, les clusters plus éloignés et moins dispersés donneront un meilleur score. Le score minimum est de zéro, les valeurs inférieures indiquant un meilleur regroupement.

Calinski

Le score est défini comme le rapport entre la dispersion intra-cluster et la dispersion inter-cluster.

```

-----
Iter_1(n_cluster 3):
-----
Fit_time: 2.384185791015625e-07
Le score ARI 0.9999521516389113
Le score HOMO 0.9997709474649992
Le score AMI 0.9997857656100236
-----
Iter_2(n_cluster 3):
-----
Fit_time: 7.152557373046875e-07
Le score ARI 0.9959405702565177
Le score HOMO 0.9901246797240195
Le score AMI 0.9888804399395965
-----
Iter_3(n_cluster 3):
-----
Fit_time: 9.5367431640625e-07
Le score ARI 0.9999521516389113
Le score HOMO 0.9997709474649992
Le score AMI 0.9997857656100236
-----
Iter_4(n_cluster 3):
-----
Fit_time: 4.76837158203125e-07
Le score ARI 0.999904303273828
Le score HOMO 0.9995678771782378
Le score AMI 0.9995975095952087
-----
Iter_5(n_cluster 3):
-----
Fit_time: 2.384185791015625e-07
Le score ARI 0.9964650441254504
Le score HOMO 0.991210947206712
Le score AMI 0.9901260596176928
-----

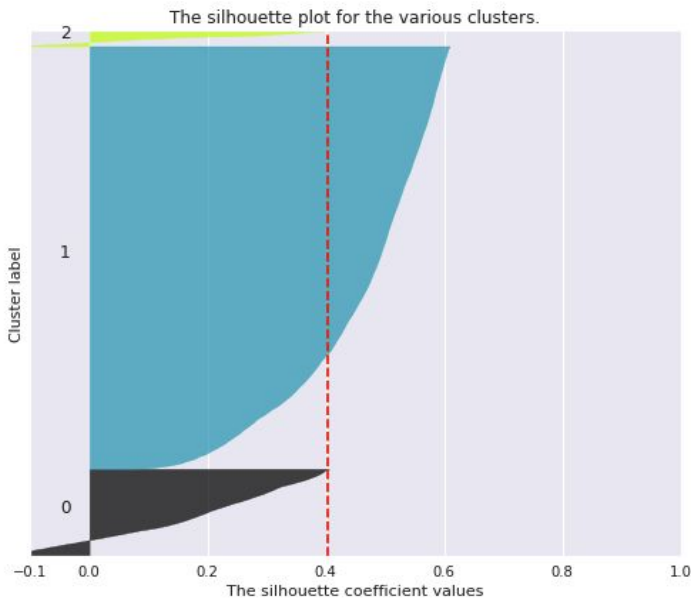
```

```

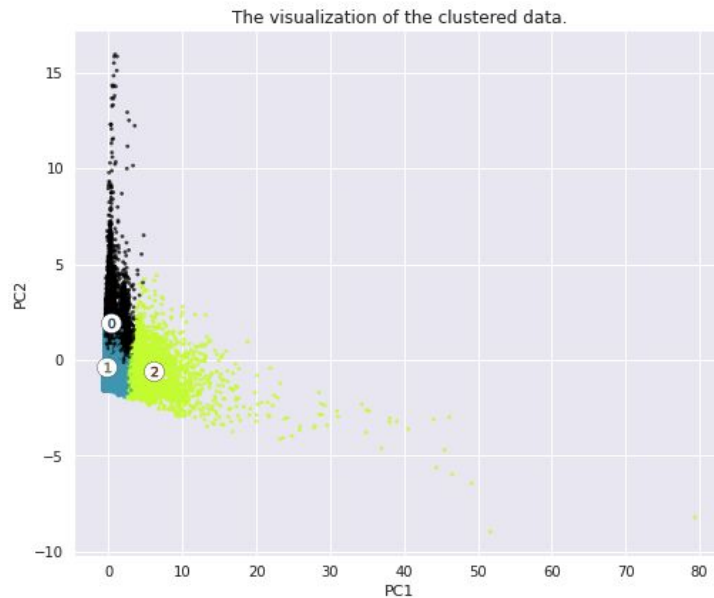
-----
Iter_6(n_cluster 3):
-----
Fit_time: 7.152557373046875e-07
Le score ARI 0.9999521516389113
Le score HOMO 0.9997709474649992
Le score AMI 0.9997857656100236
-----
Iter_7(n_cluster 3):
-----
Fit_time: 1.1920928955078125e-06
Le score ARI 0.9998086065314611
Le score HOMO 0.9991877183510793
Le score AMI 0.9992469672824071
-----
Iter_8(n_cluster 3):
-----
Fit_time: 2.384185791015625e-07
Le score ARI 1.0
Le score HOMO 1.0
Le score AMI 1.0
-----
Iter_9(n_cluster 3):
-----
Fit_time: 9.5367431640625e-07
Le score ARI 0.999904303273828
Le score HOMO 0.9995678771782378
Le score AMI 0.9995975095952087
-----
Iter_10(n_cluster 3):
-----
Fit_time: 2.384185791015625e-07
Le score ARI 0.999904303273828
Le score HOMO 0.9995678771782378
Le score AMI 0.9995975095952087
-----

```

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 3$

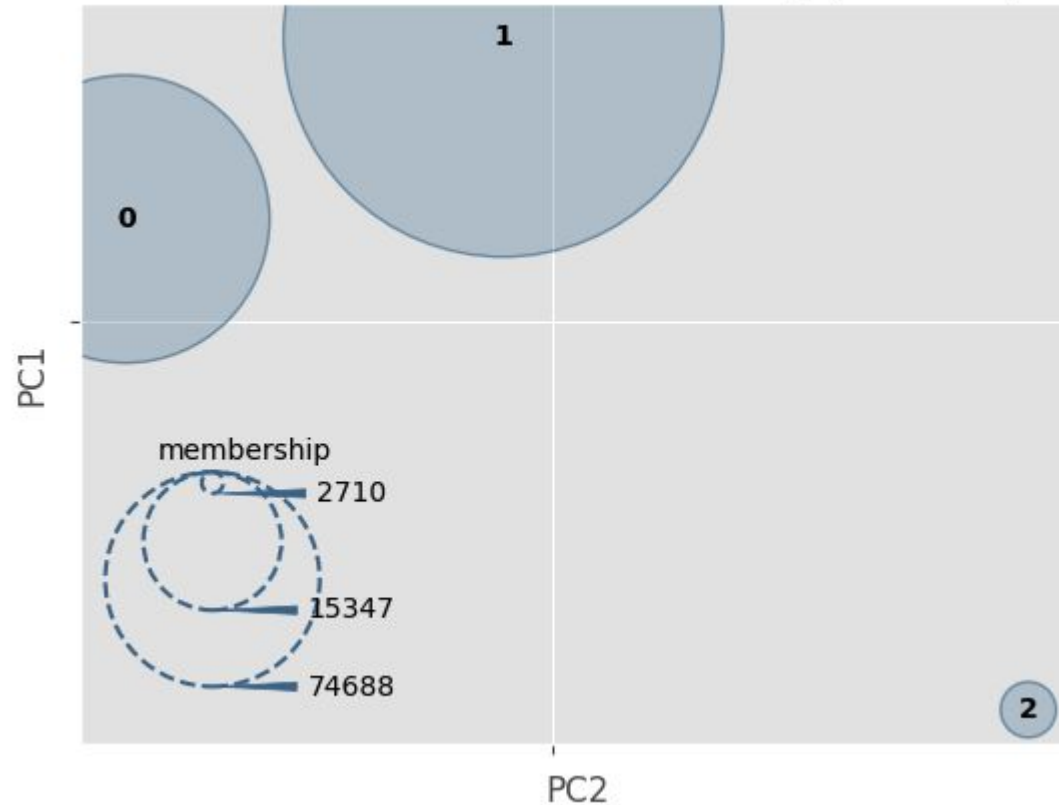


Cluster 0: 15347 clients
Cluster 1: 74688 clients
Cluster 2: 2710 clients



- Silhouette Score: 0.40
- calinski_harabasz Score: 25600
- davies_bouldin Score: 1.18

KMeans Intercluster Distance Map (via MDS)



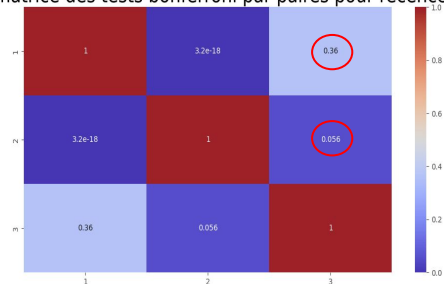
La distance intercluster permet de voir une dissimilarité entre les clusters

Kruskal-wallis Test 5%

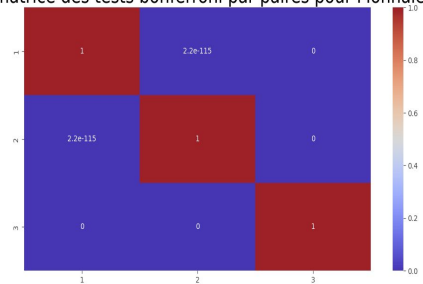
- `stat=60105.071, p=0.000`
`recence moyenne`
- `stat=19657.082, p=0.000`
`fréquence`
- `stat=6645.377, p=0.000`
`Monnaie`
- `stat=41808.725, p=0.000`
`review_score_mean`
- `stat=12074.576, p=0.000`
`Duration_delevery_order_mean`
- `stat=25552.940, p=0.000`
`mean_order_item_id`

Déterminer si les clusters par paires sont différents

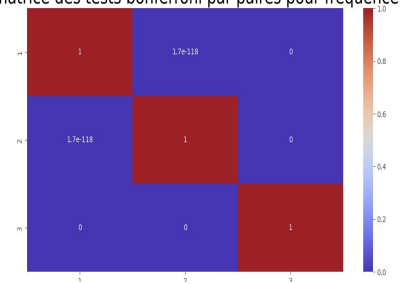
La matrice des tests bonferroni par paires pour recence



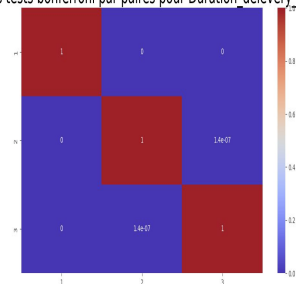
La matrice des tests bonferroni par paires pour Monnaie



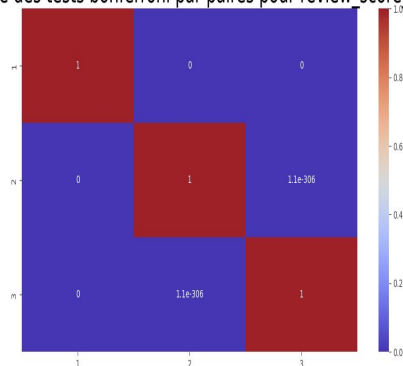
La matrice des tests bonferroni par paires pour frequence



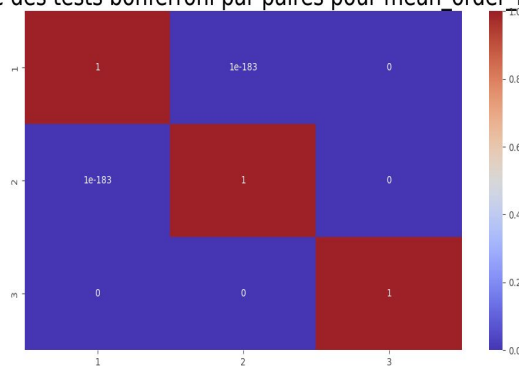
La matrice des tests bonferroni par paires pour Duration_delevery_order_mean



La matrice des tests bonferroni par paires pour review_score_mean



La matrice des tests bonferroni par paires pour mean_order_item_id



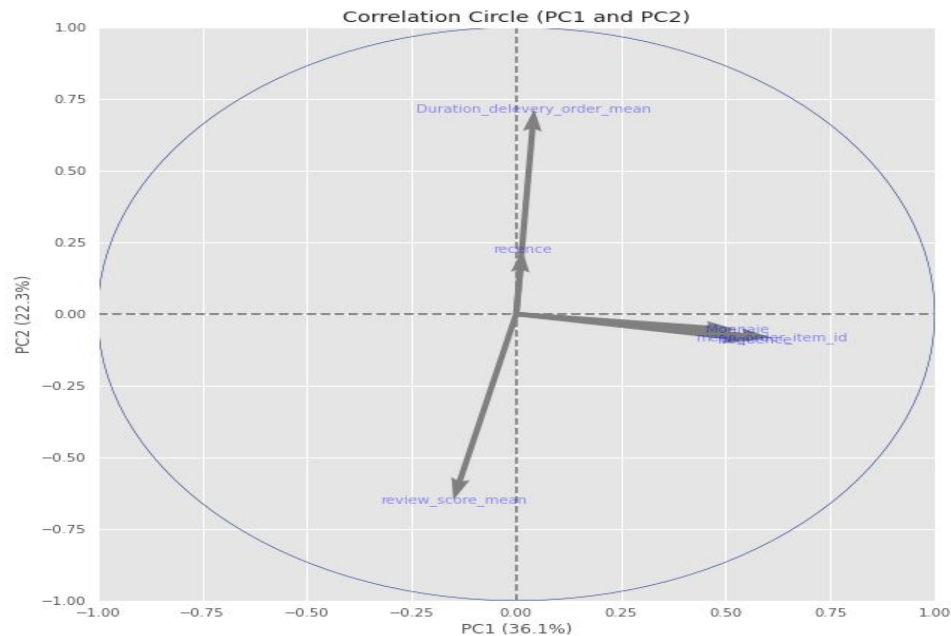
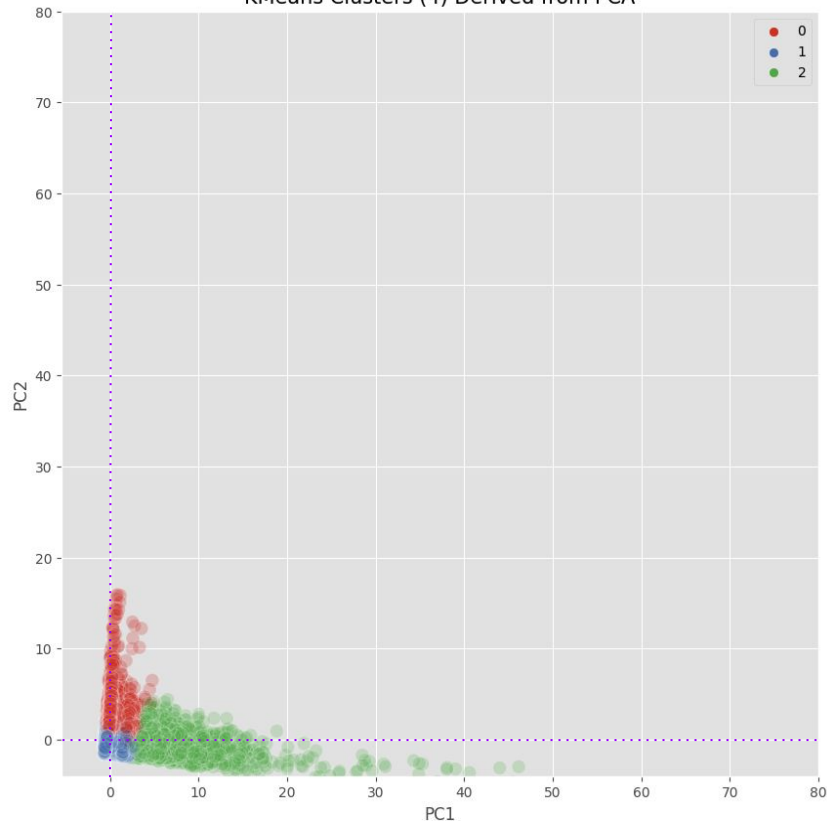
Test Dunn 5%

Hypothèses des test paramétriques non respectées (Variance et normalité)

Les clusters non différents sont (1 et 3) (2 et (2 et 3) pour la recence

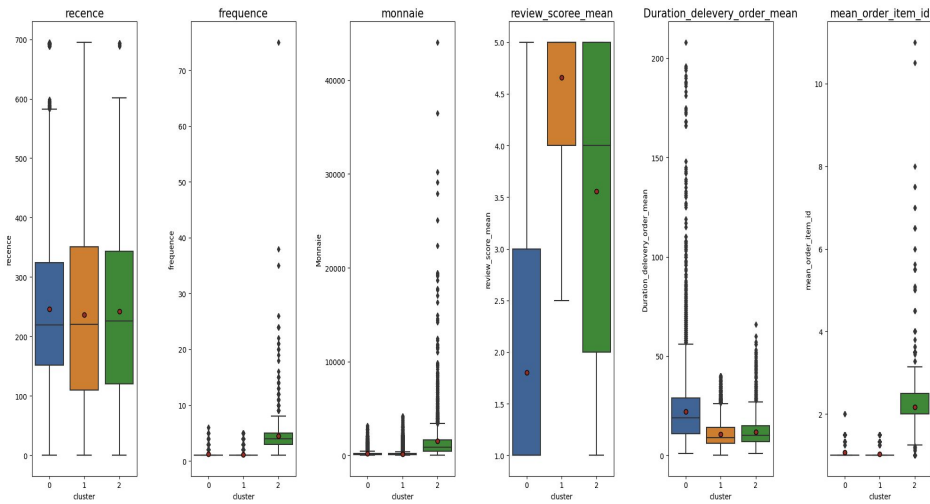
Les autres paires sont différentes pour les autres features

KMeans Clusters (4) Derived from PCA



PC1 : Monnaie, mean_order_item_d, fréquence

PC2 : review_score_mean,
Duration_delevery_order_mean



Parallel Coordinates plot for the Centroids



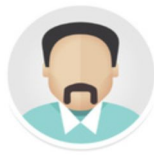
cluster 0



**Client
très
insatisfait
peu
dépensier
16.55 %**

- fréquence : 1
- Moyenne articles : 1
- Dépense moyenne : 203
- Satisfaction moyenne : 1.80
- Durée livraison moyenne : 22 jours

cluster 1



**Client
très
satisfait
peu
dépensier
80.53 %**

- fréquence : 1
- Moyenne articles : 1
- Dépense moyenne : 165.26
- Satisfaction moyenne : 4.65
- Durée livraison moyenne : 10 jours

Cluster 2



**Best client
satisfaction
moyenne
très
dépensier
2.92 %**

- fréquence : 4
- Moyenne articles : 2
- Dépense moyenne : 1509.46
- Satisfaction moyenne : 3.55
- Durée livraison moyenne : 11 jours

2016-10-03 Dataset 2018-08-29

2016-10 1 an 2017-10

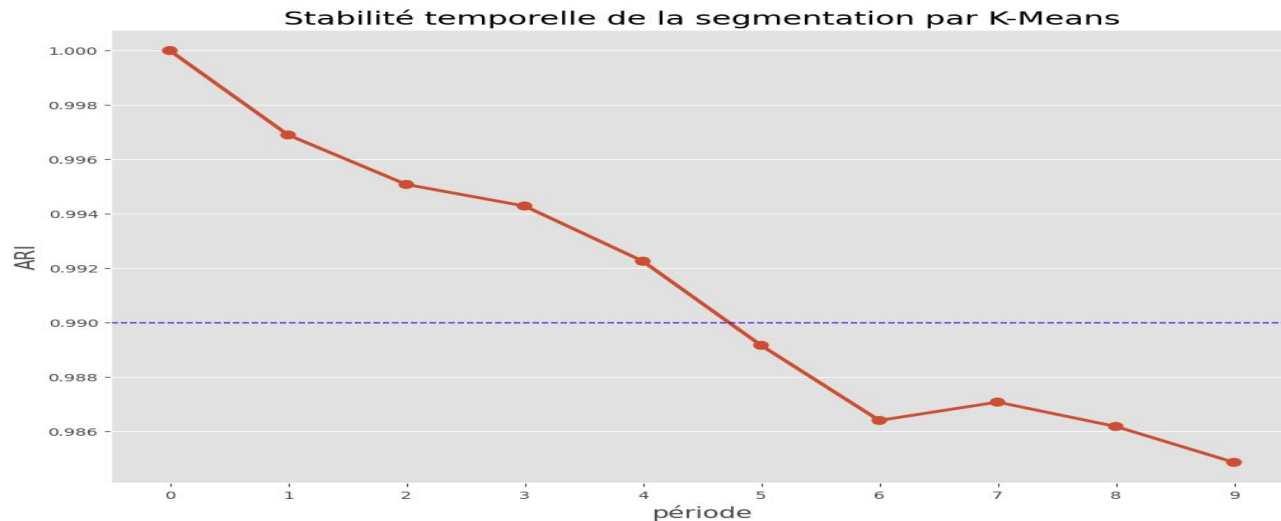
Labels True

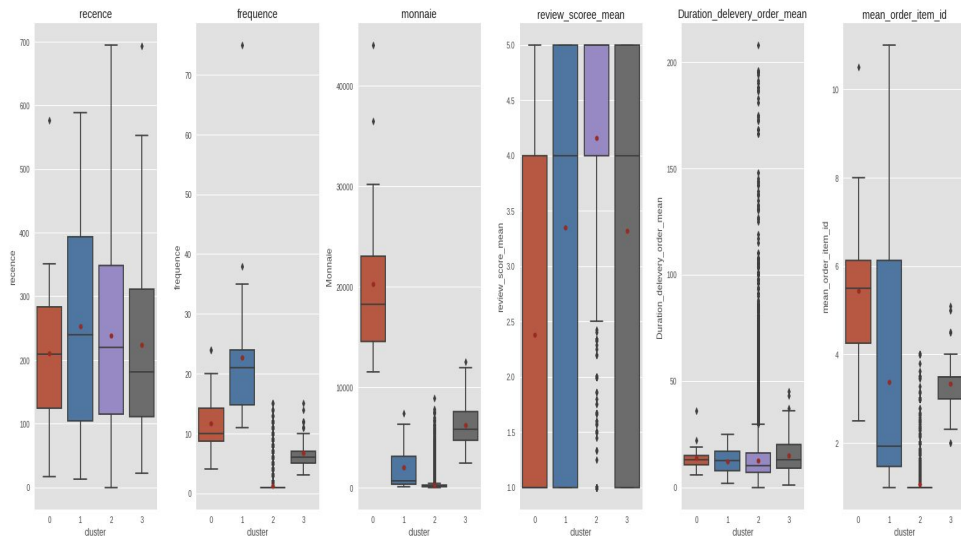
ARI

Labels pred

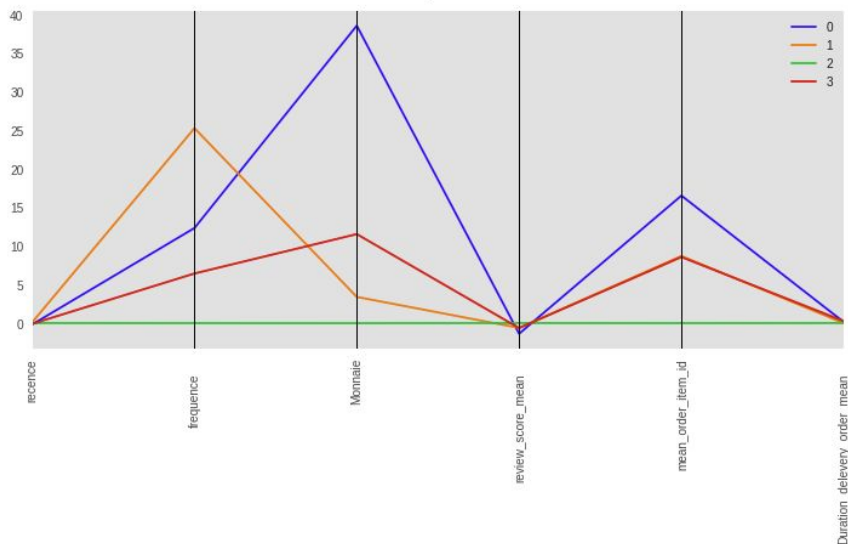
- Clients communs pour chaque périodes par rapport à la période de 1 an
- Agg des variables
- PCA
- Kmeans(n_clusters=3,init='k-means++', random_state=39)

periode		ARI
0	0	1.000000
1	1	0.996899
2	2	0.995078
3	3	0.994286
4	4	0.992258
5	5	0.989166
6	6	0.986405
7	7	0.987079
8	8	0.986183
9	9	0.984864





Parallel Coordinates plot for the Centroids



Effectifs

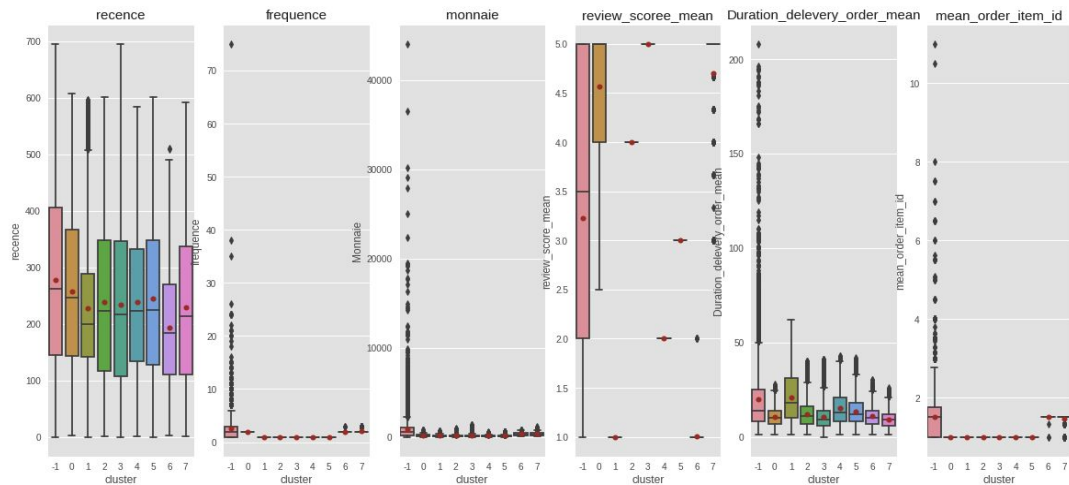
cluster 0 : 24
 cluster 1 : 24
 cluster 2 : 92589
 cluster 3 : 108

Evaluation

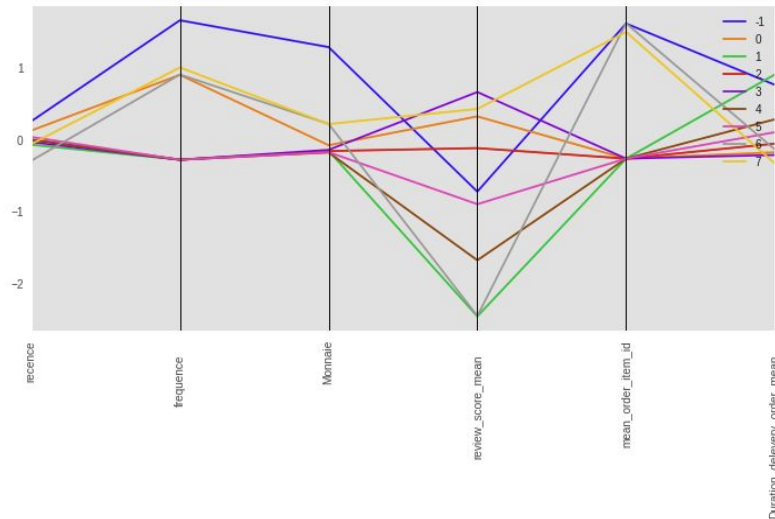
- Silhouette Score: 0.84
- calinski_harabasz Score: 6164.79
- davies_bouldin Score: 0.68

Conclusion

Le clustering ne permet pas de répondre à la problématique métier (Les effectifs des clusters sont trop hétérogène pour effectuer des stratégies marketing)



Parallel Coordinates plot for the Centroids



Effectifs

cluster -1 (bruit): 8350 cluster 3: 15506
 cluster 0: 3003 cluster 4: 1897
 cluster 1: 6092 cluster 5: 6055
 cluster 2: 46826 cluster 6: 867
 cluster 7: 4149

Evaluation

- Silhouette Score: 0.09
 - calinski_harabasz Score: 5985.50
 - davies_bouldin Score: 2.62

Conclusion

- Le clustering ne permet pas de répondre à la problématique métier (Les effectifs des clusters sont trop hétérogène pour effectuer des stratégies marketing)
- Les métriques indiquent que la segmentation n'est pas pertinente

La segmentation finale porte sur 6 variables :

- Recence
- Fréquence
- Monnaie
- Review_score_mean
- Mean_order_item_id
- Duration_delevery_order_mean

Le modèle final choisis est le Kmeans:

- 3 clusters
- Stable à l'initialisation avec 10 itérations
- Le clustering permet d'identifier des profils clients afin d'adapter une stratégie marketing pour chaque clusters
- Le clustering permet de segmenter les bons et moins bons clients en termes de commandes et de satisfaction.

contrat de maintenance

- Stabilité fortement dépendante des variables choisies
- La fréquence à laquelle la segmentation doit être mise à jour est 5 mois pour qu'elle reste pertinente

