

Exploration des données

Arnaud Vanholderbeke, Benjamin Lorient et Natan Danous

11 Mai 2019

Contents

Description	1
Types et granularité	1
Distribution	2
Relations	11
Mesure de l'entropie et du gain d'entropie	11

Description

Ces données sont les ventes provenant d'un seul magasin lors du Black Friday. Le magasin veut mieux comprendre le comportement des consommateurs sur différents produits. Selon Kaggle, le problème principal est un problème de régression. On essaie de prédire la variable quantité d'achat à l'aide des autres variables. Ce problème peut également être vu comme un problème de régression.

Types et granularité

Table 1: Dictionnaire des données

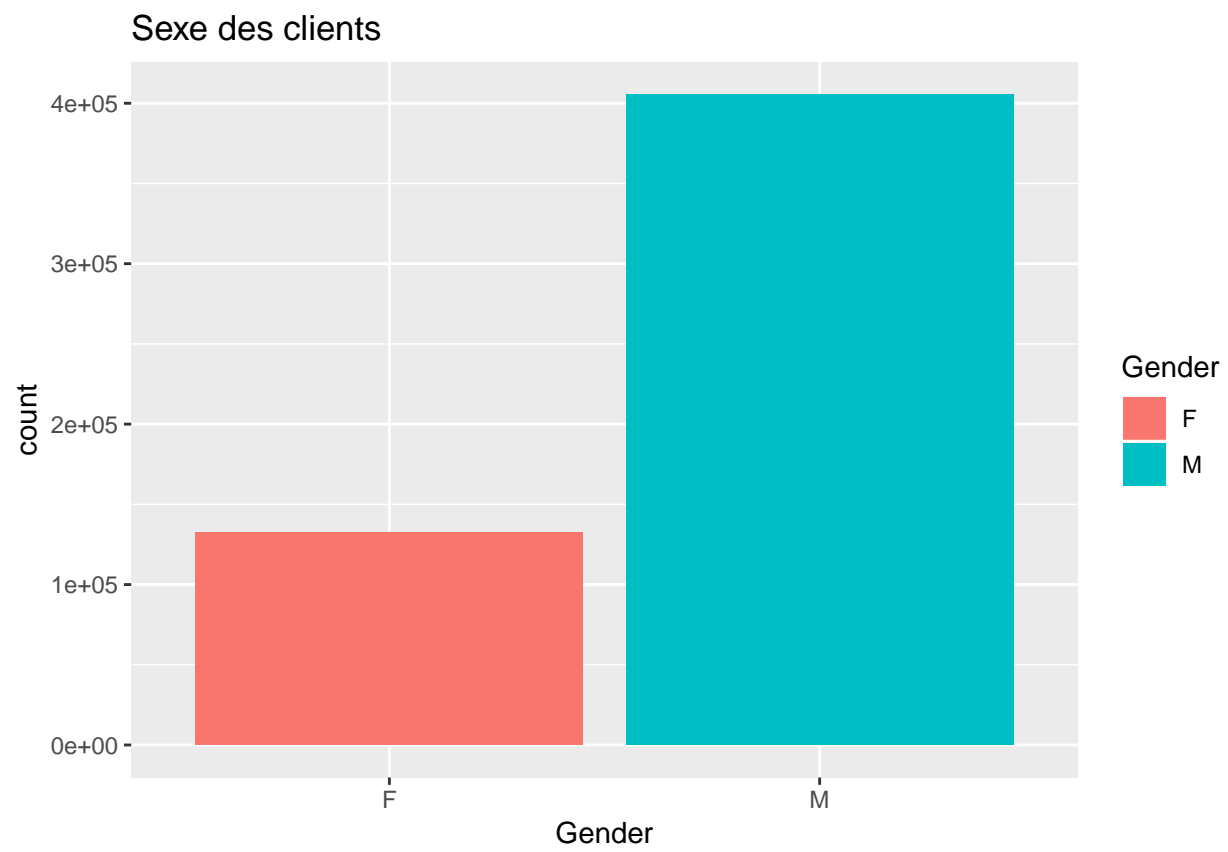
Champ	Type	Modalités
Purchase	Quantitative	
User_ID	Ordinal	
Product_ID	Nominal	
Gender	Nominal	F, M
Age	Ordinal	0-17 -> ... -> 51-55 -> 55+
Occupation	Nominal	0 -> 20
City_Category	Nominal	A, B, C
Stay_In_Current_City_Years	Ordinal	0, 1, 2, 3, 4+
Marital_Status	Nominal	0, 1
Product_Category_1	Nominal	1 -> 18
Product_Category_2	Nominal	2 -> 18 (avec NA)
Product_Category_3	Nominal	3 -> 18 (avec NA)

Remarque :

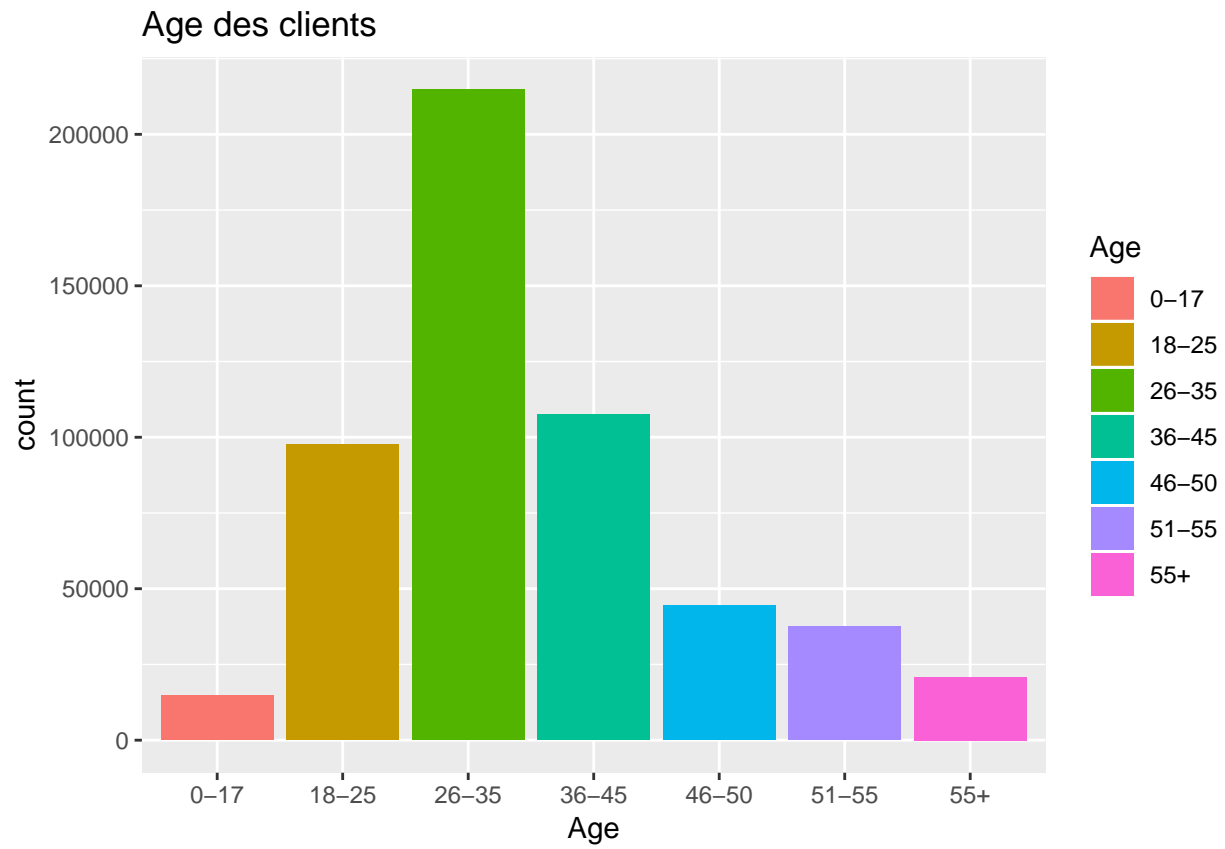
- Un produit a obligatoirement une catégorie (Product_Category_1 ne contient pas de NA). Il n'a pas obligatoirement de 2ème et 3ème catégorie.

Distribution

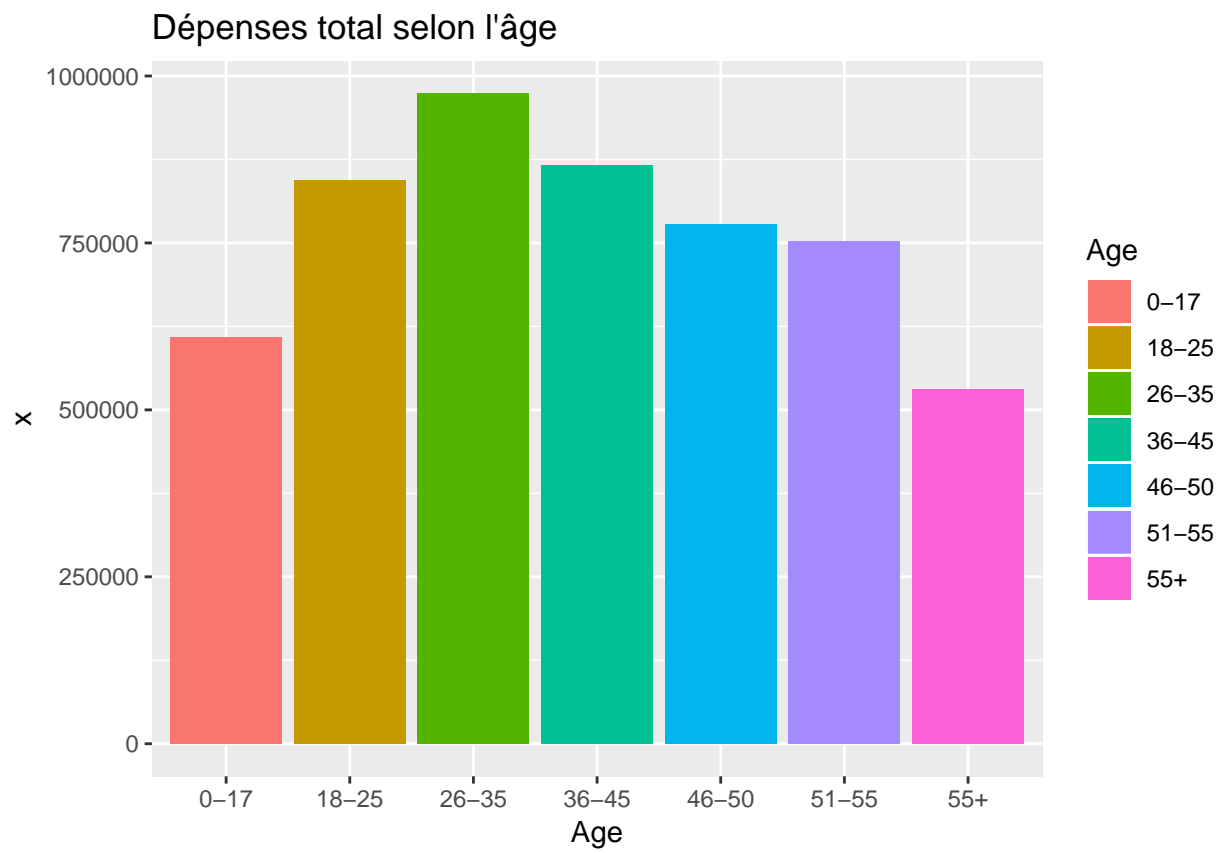
Sexe



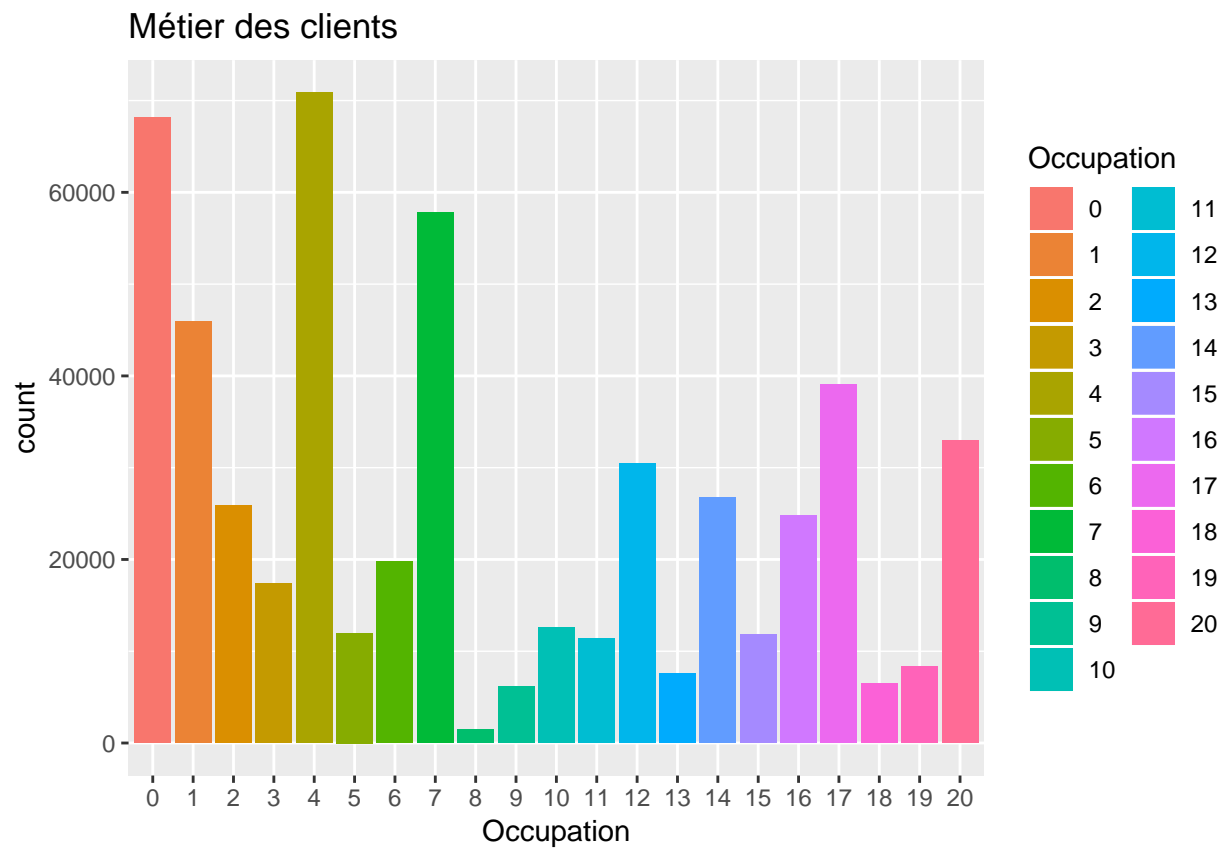
Age



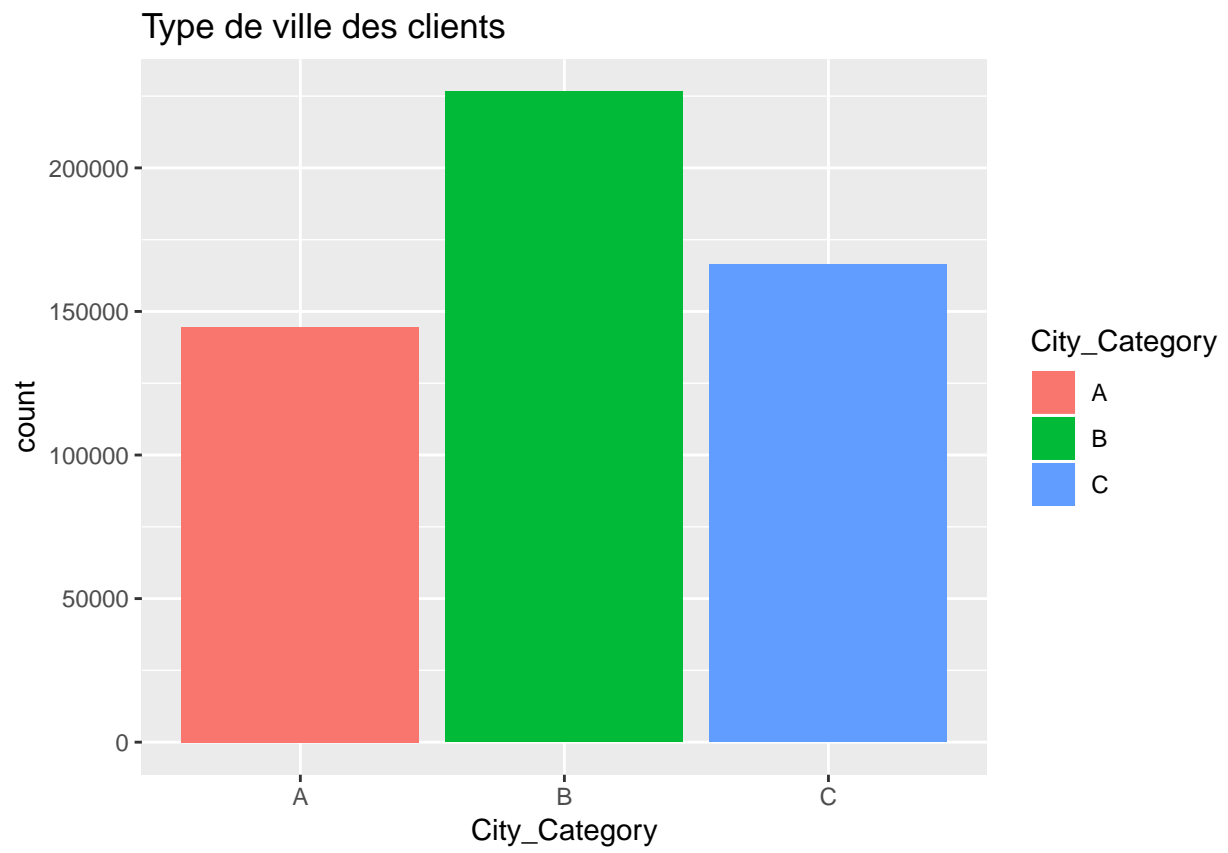
La distribution de l'âge dans une population suit habituellement une loi normale. On observe ici que ce n'est pas le cas. Il apparaît que la tranche d'âge des 26-35 effectue le plus d'achats lors du Black Friday.



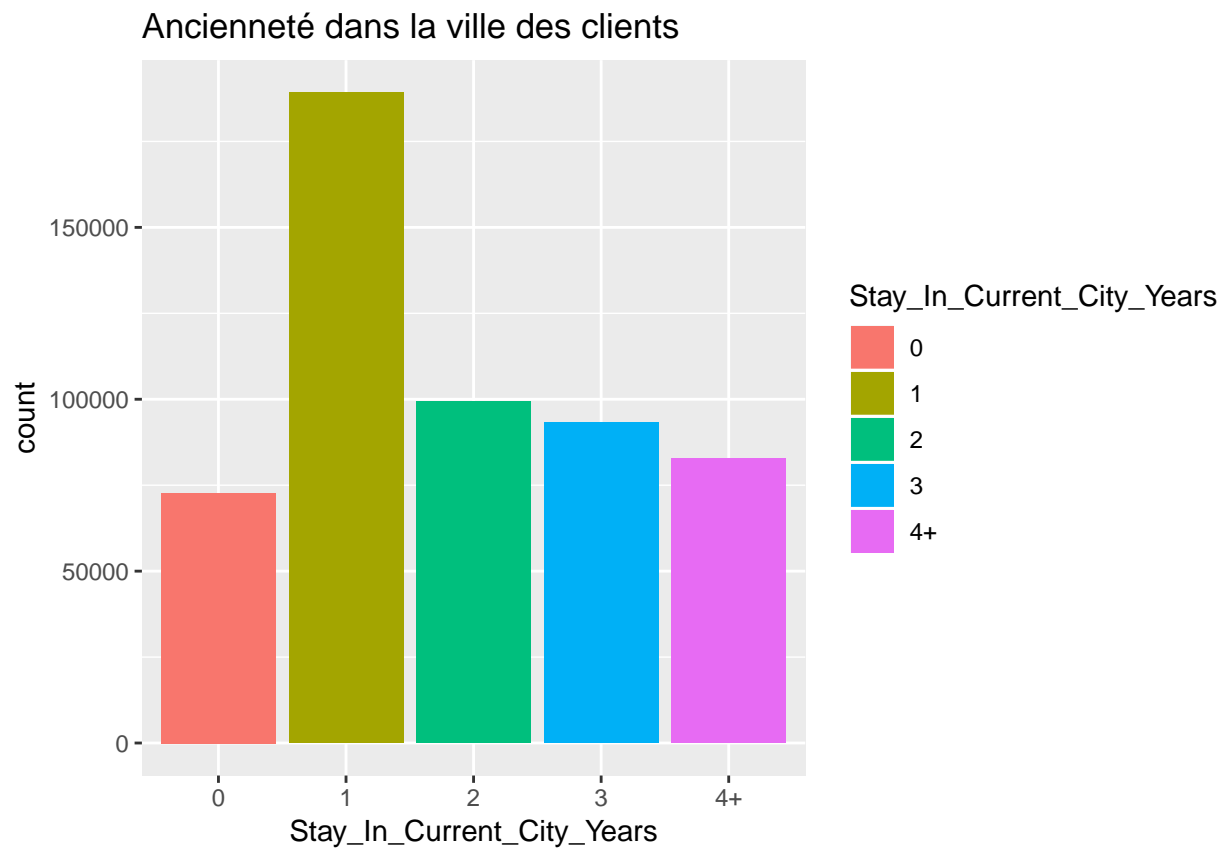
Métier



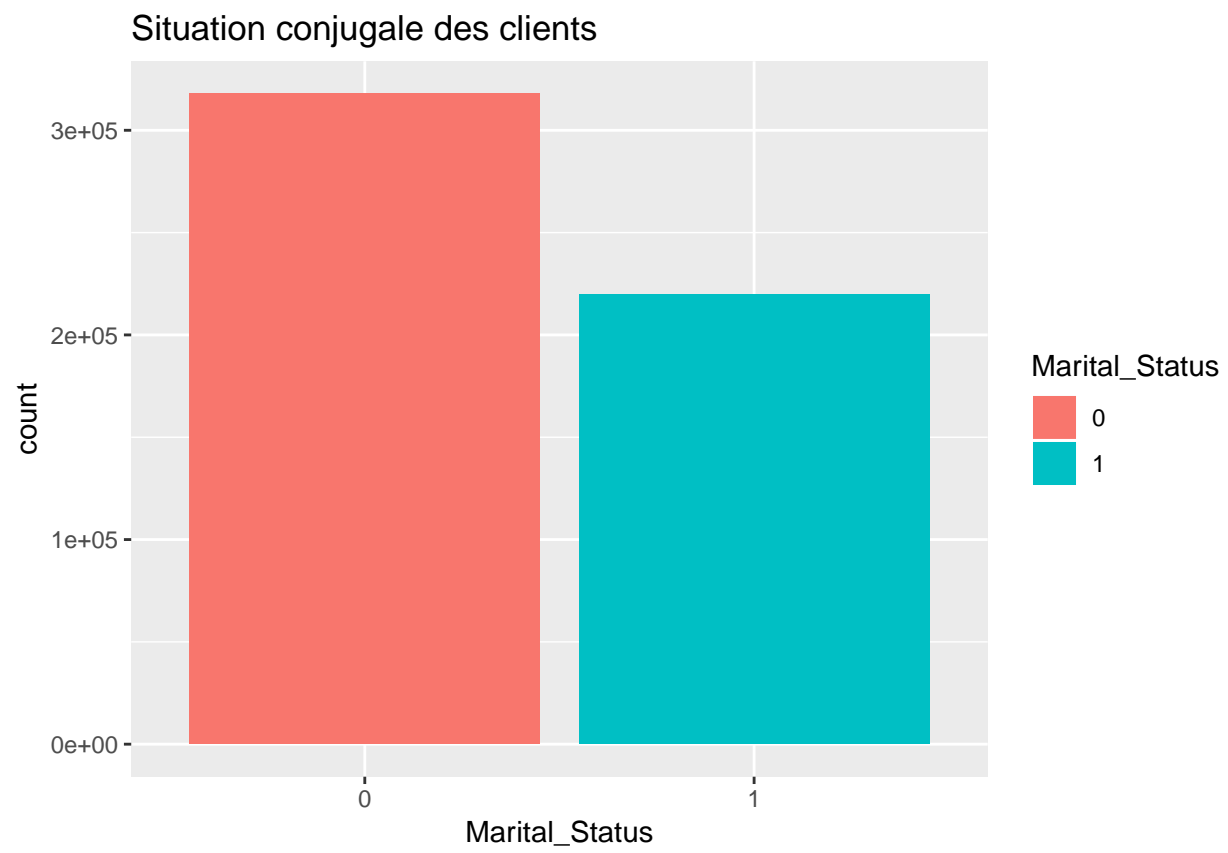
Type de ville



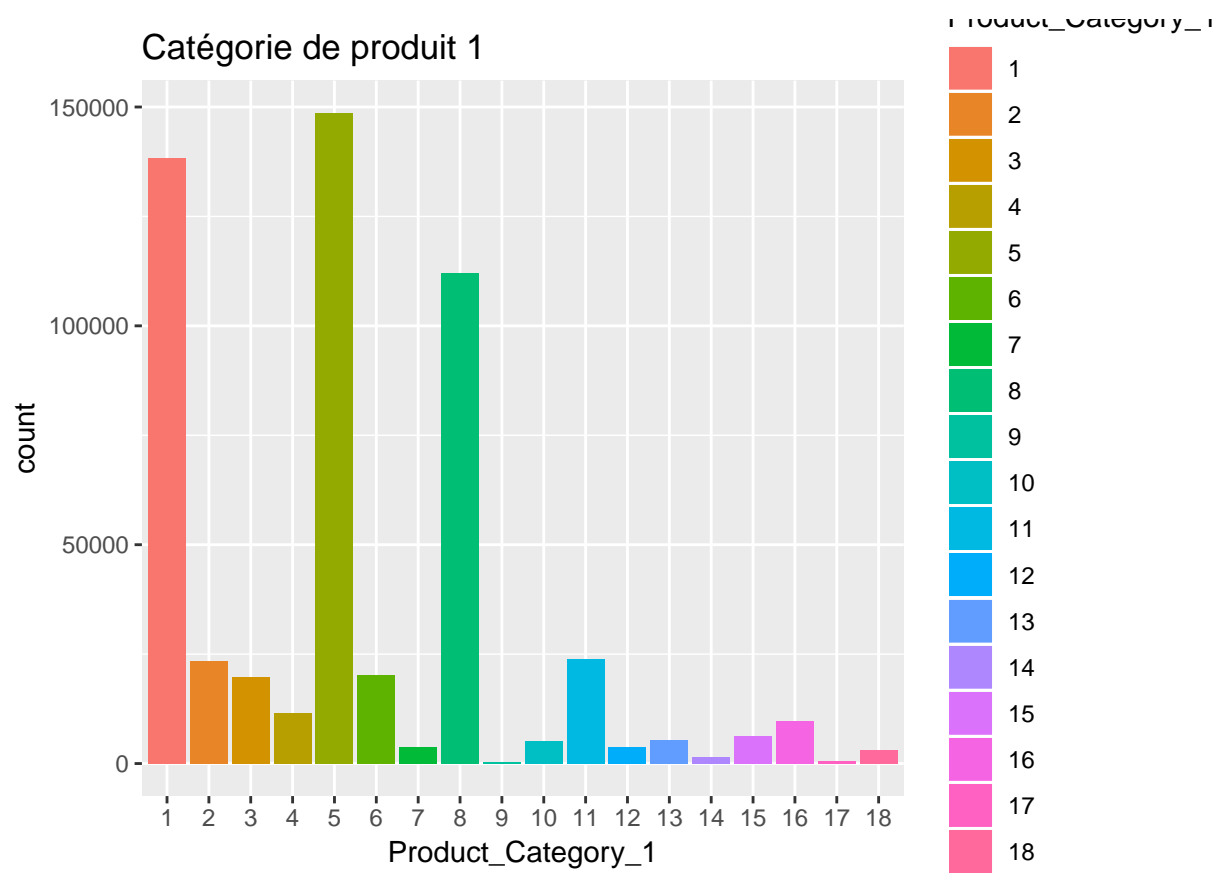
Ancienneté dans la ville

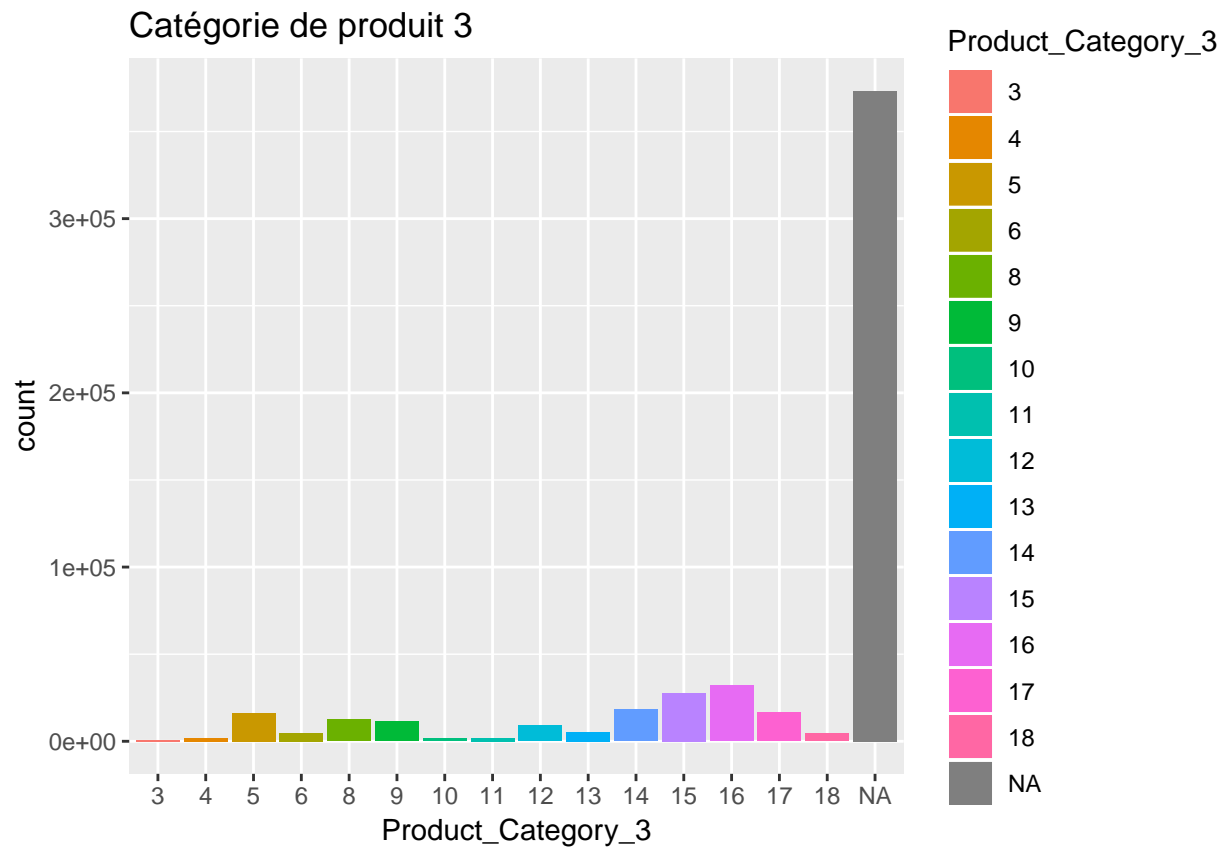
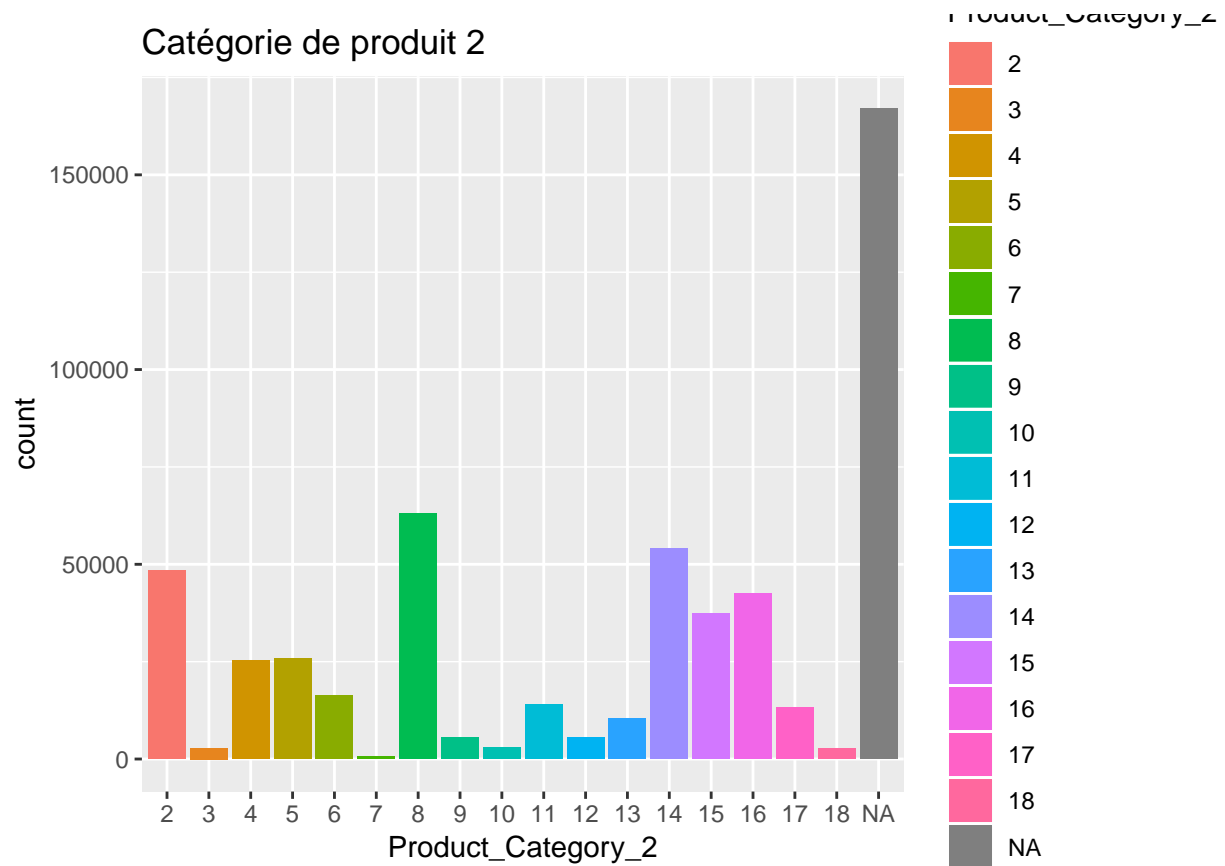


Situation conjugale

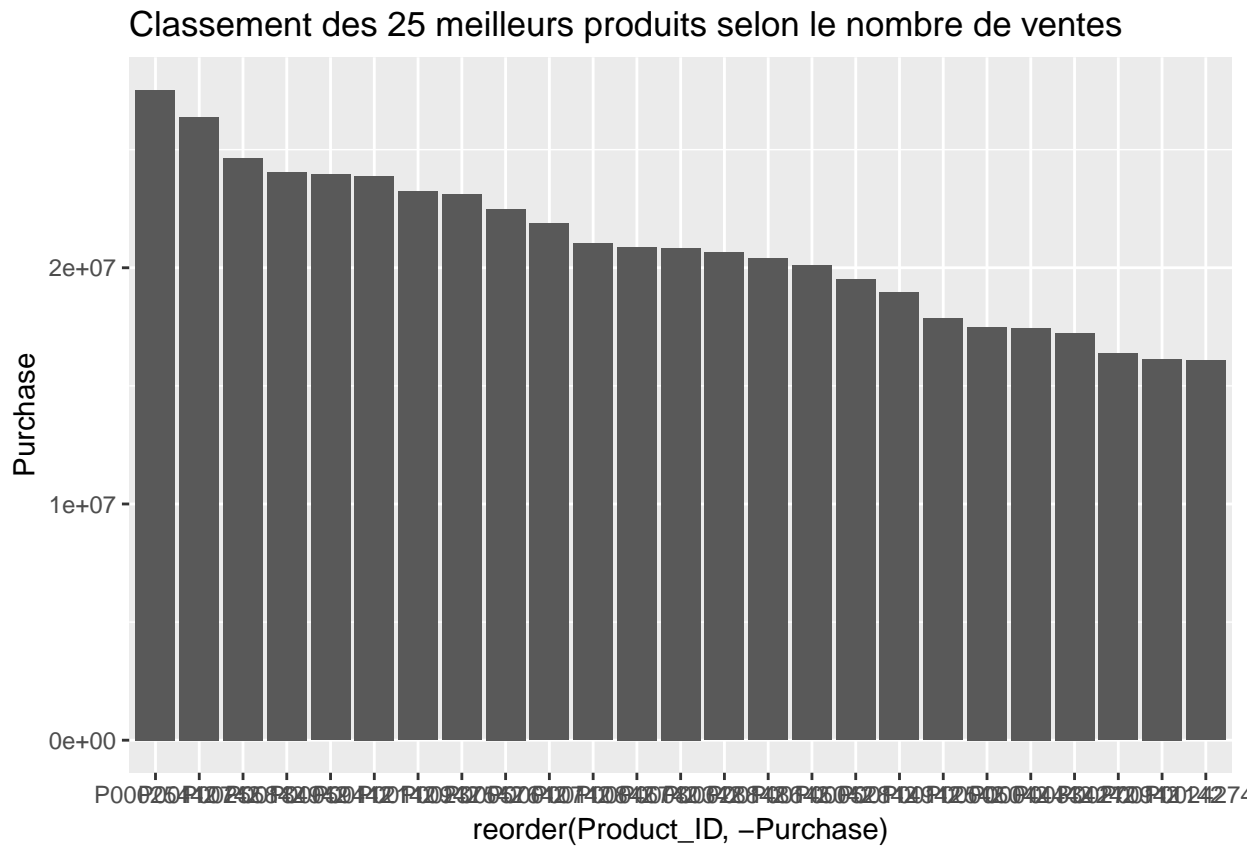


Catégories des produits





Performance des produits



Pas de produit que se détache du reste (du moins dans les 25 premiers, il faut garder à l'esprit qu'il y en a +3000).

Relations

Sociologie

- Age
- Métier
- Ancienneté dans la ville
- Marié, pas marié
- Sexe

=> Montant achat en fonction de l'âge, métier, etc...

Commercial

- Performance des produit (via ID)
- Catégorie de produit qui se vendent le mieux

Mesure de l'entropie et du gain d'entropie

TODO: A expliquer et justifier et synthétiser.

##	Age	Gender
##	0.069374015	0.022694835

## Stay_In_Current_City_Years	Occupation
## 0.013115811	0.091392940
## Marital_Status	City_Category
## 0.007119339	0.019069753

Deux catégories importantes : Age et Occupation.