

25/10/2023 | Par : Abderrahmane BOUFARES



Rapport statistique

Analyse de base de données

Introduction

Présentation de la problématique.

Problématique :

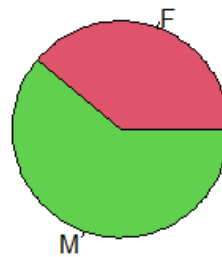
Les rectorats des universités du Sud-Ouest (UPPA, Bordeaux et Toulouse), cherche à établir des règles de décision en modélisation pour chaque licence, la probité de réussite à partir des données observées sur la population étudiante des années précédentes. Afin d'envoyer aux futurs étudiants, des avis reflétant leurs chances de réussite dans la licence choisie, grâce à un programme d'orientation active de lycéens. Notre base de donnée destinée à l'étude se compose de deux type de variables, **des variables** qualitative tel que le sexe (Féminin, Masculin), la profession des parents (Défavorable, Moyen, Favorable et Très Favorable), le retard de l'obtention du BAC (Bac réussi dès la première année, réussi en deuxième année ou réussi en troisième année), la série du BAC (ES, L, S, SMS, STI ou STT), type du BAC (General ou Technologique), la qualité de la mention obtenue au BAC (A,B,C,D ou P), l'obtention du BAC (Obtenue ou non obtenue), la mention (Obtenue ou non obtenue), la réussite en première de licence (Réussie ou non réussie) et **une variable quantitative**, celle de la note obtenue à l'épreuve de Math du BAC, entre 1 et 19. Pour faire cette étude, on va travailler avec le logiciel Rstudio et une base de donnée sous R, de nom : L1eco.rdata. Ce logiciel nous permettra de faire des tables de contingences, des tests statistiques sur cette base de données et créer des graphiques, sous forme d'histogramme, diagramme en barre, circulaire et mosaïque, boîte à moustache. L'objectif de cette étude, est d'estimer la dépendance et la relation entre les variables croisées, elle permet aussi de comparer l'intensité du lien entre les variables croisées et leurs corrélations, elle nous permet aussi de vérifier l'Indépendance entre les variables.

Développement

Analyse univariée :

- Sexe :

Repartition de la population
selon le sexe



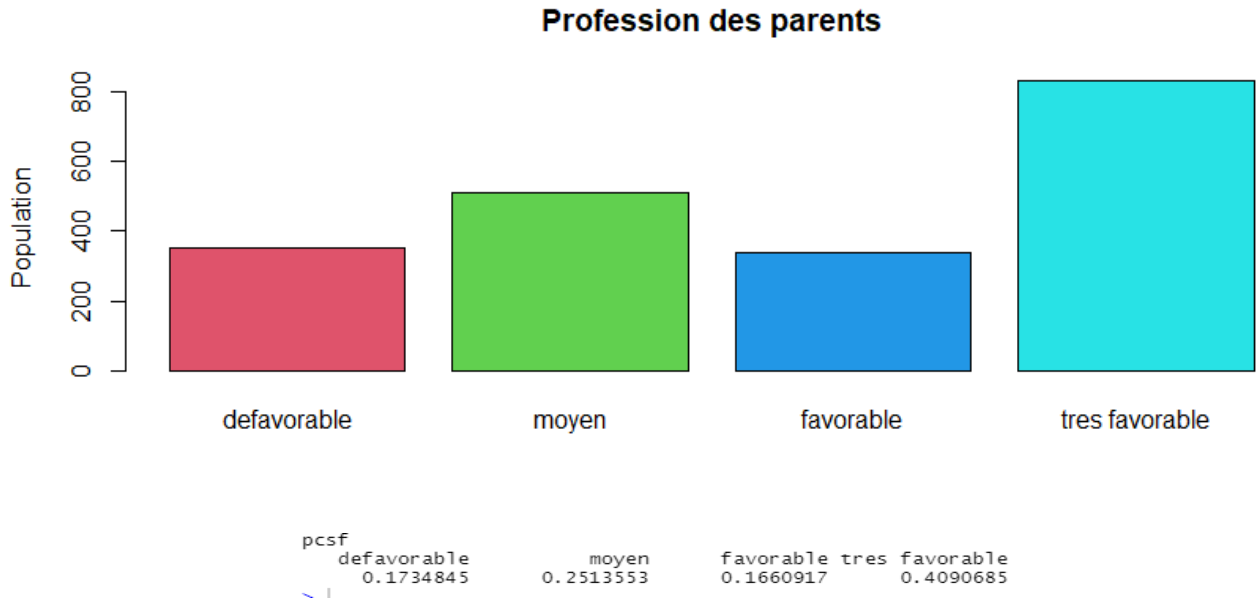
F M
0.3908329 0.6091671

D'après le diagramme circulaire et la table de contingence on distingue, que notre population est divisée en deux niveaux, le sexe Féminin qui représente 39% de l'échantillon et le sexe Masculin qui représente 60% de l'échantillon.

Développement

Analyse univarié :

- Profession des parents :

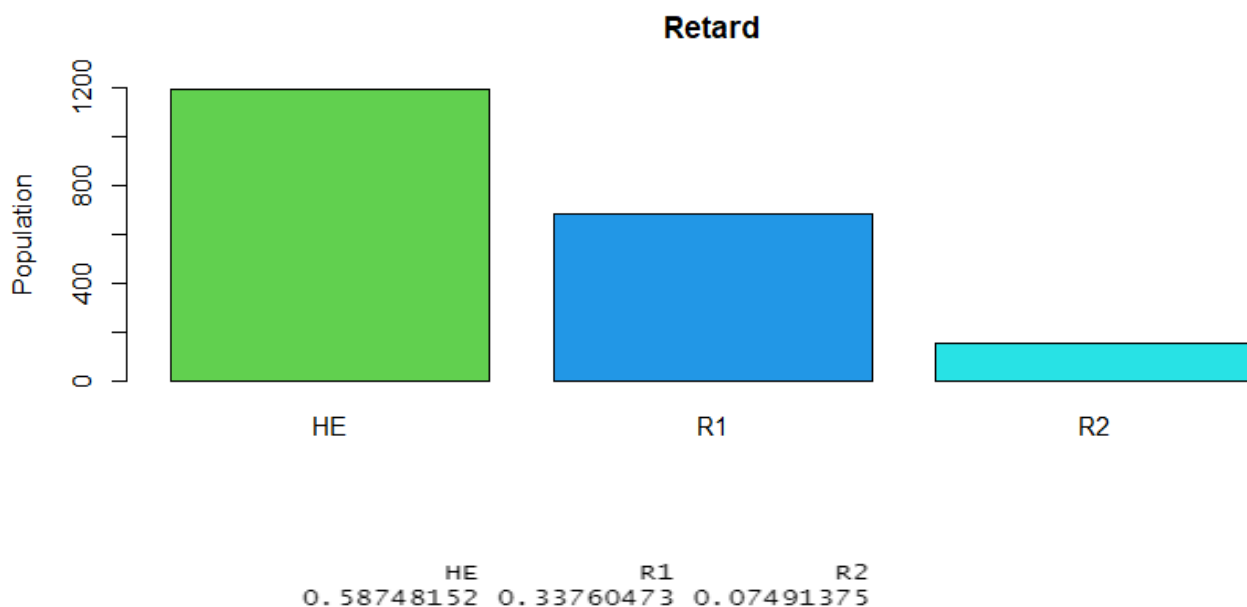


D'après le diagramme en barre et la table de contingence on distingue, que notre population est divisée en quatre niveaux, la profession défavorable qui représente 17% de l'échantillon, la profession moyen qui représente 25% de l'échantillon, , la profession favorable qui représente 16% de l'échantillon et la profession très favorable qui représente 40% de l'échantillon

Développement

Analyse univarié :

- Retard :

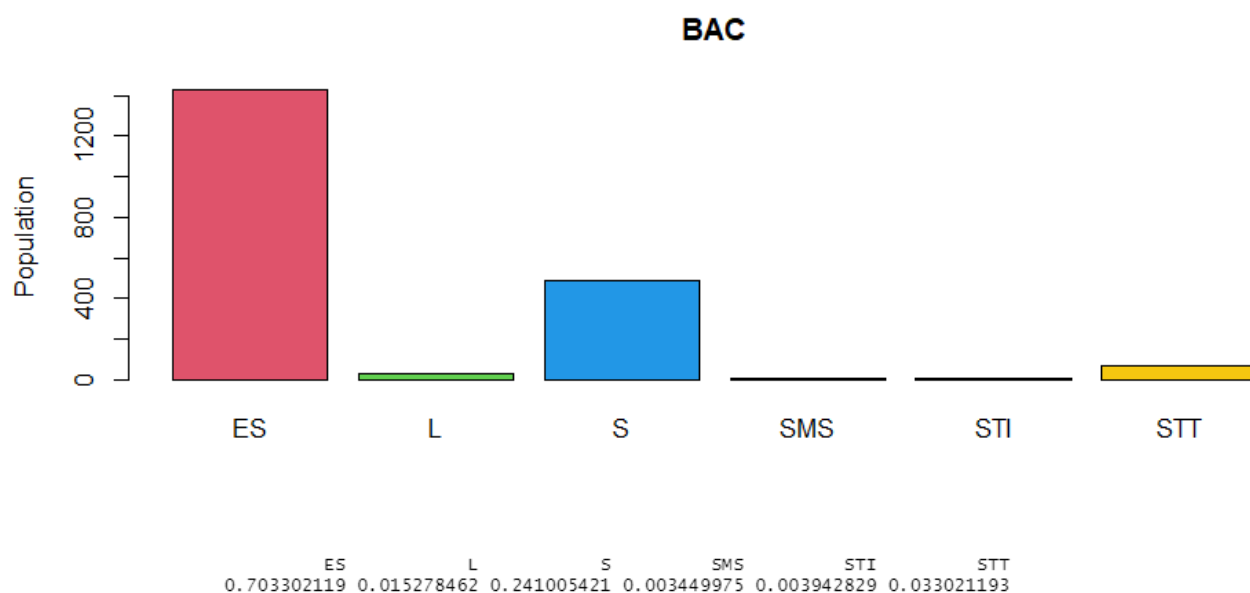


D'après le diagramme en barre et la table de contingence on distingue, que notre population est divisée en trois niveaux, les non redoublants qui représente 58% de l'échantillon, les étudiants qui ont redoublés une seule fois qui représente 33% de l'échantillon, les étudiants qui ont redoublés deux fois qui représente 16% de l'échantillon et la profession très favorable qui représente 40% de l'échantillon.

Développement

Analyse univarié :

- Série du BAC :

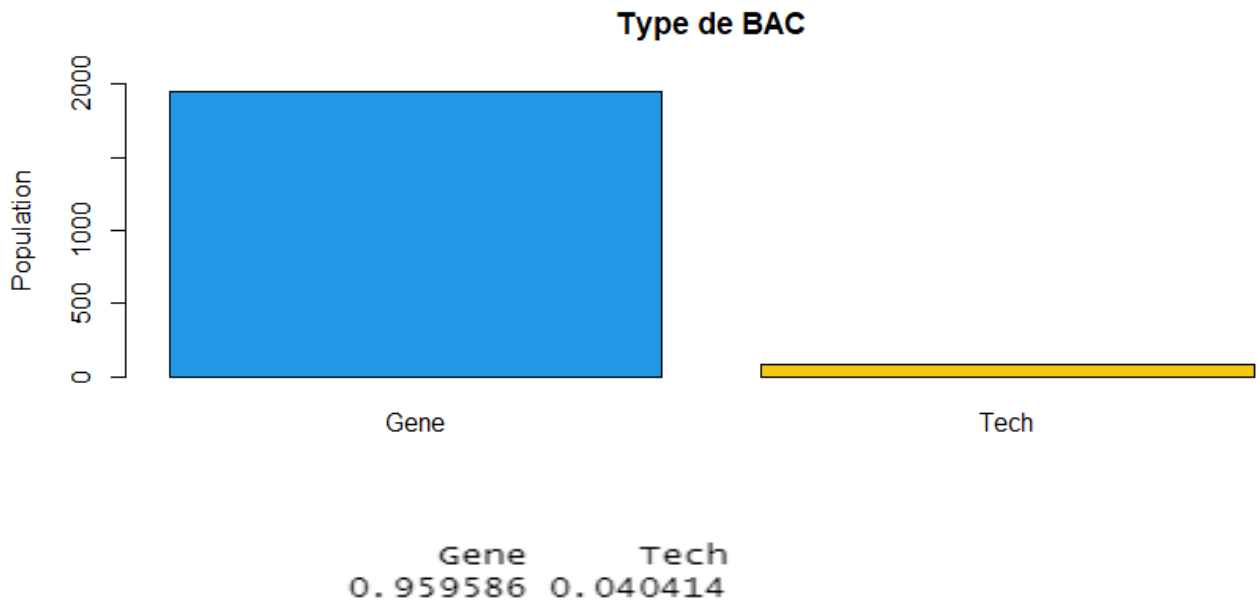


D'après le diagramme en barre et la table de contingence on distingue, que notre population est divisée en six niveaux, les séries du BAC : ES, L, S, SMS, STI et STT qui représente respectivement 70,3%, 1,5%, 24,1%, 0,34%, 0,39% et 3,3% de l'échantillon.

Développement

Analyse univarié :

- Type de BAC :

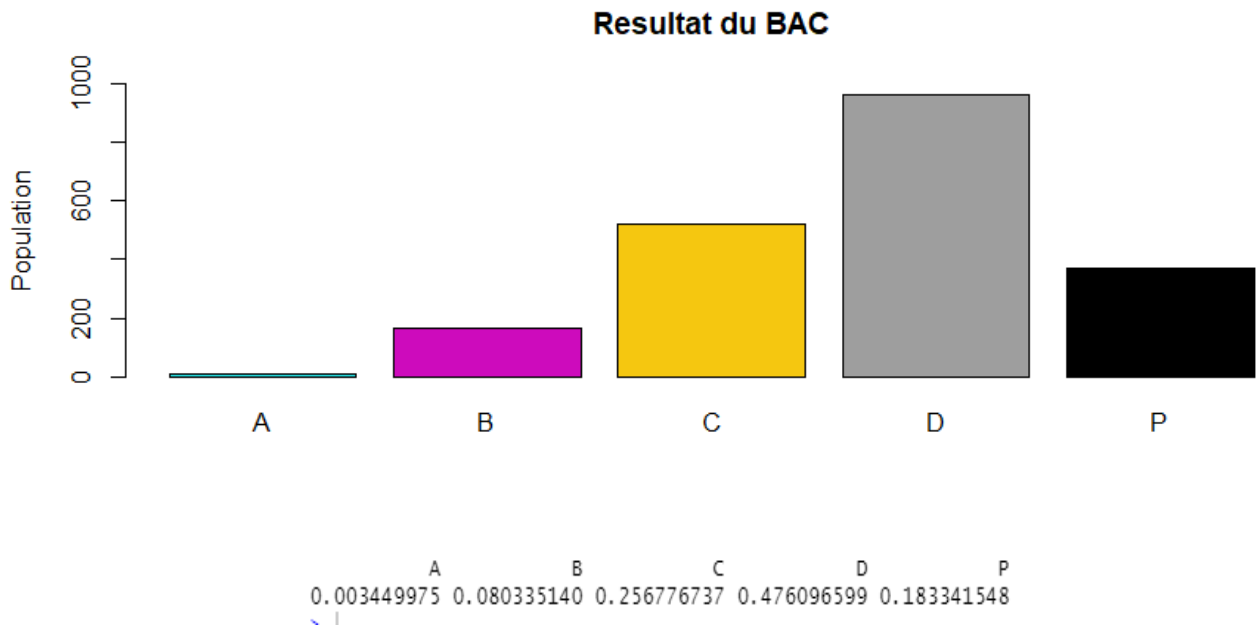


D'après le diagramme en barre et la table de contingence on distingue, que notre population est divisée en deux niveaux, les BACS généraux qui représente 96% de l'échantillon, les BACS techniques qui représente 4% de l'échantillon.

Développement

Analyse univarié :

- Qualité de mention :

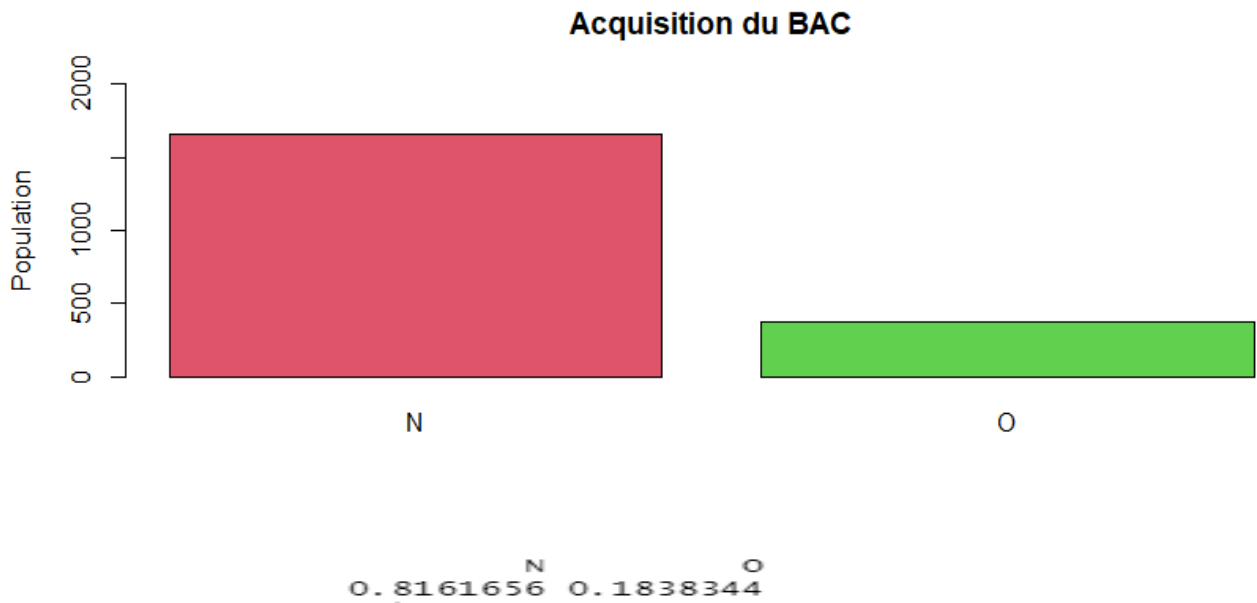


D'après le diagramme en barre et la table de contingence on distingue, que notre population est divisée en cinq niveaux, les mention BAC : A, B, C, D et P, qui représente respectivement 0,034%, 8%, 25%, 0,47%, 0,39% et 18% de l'échantillon.

Développement

Analyse univarié :

- Obtention du BAC :

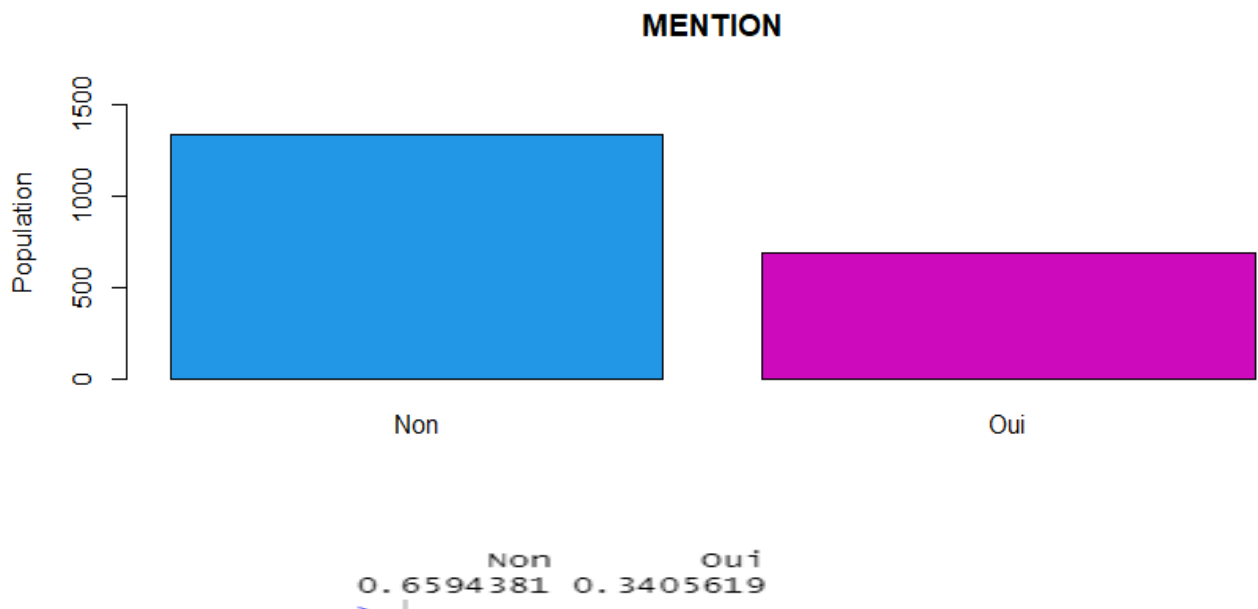


D'après le diagramme en barre et la table de contingence on distingue, que notre population est divisée en deux niveaux, 81% ont obtenu le BAC et 18% ne l'ont pas obtenu.

Développement

Analyse univarié :

- Obtention de mention :

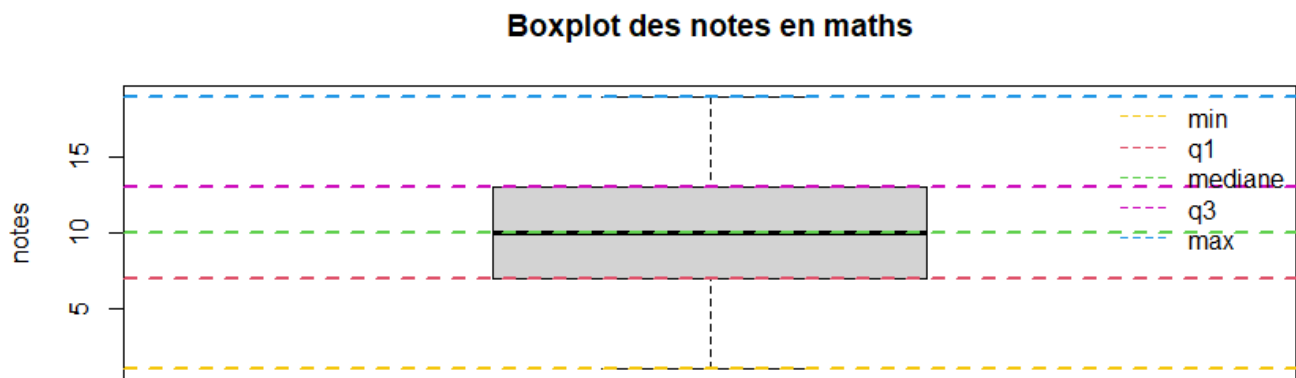
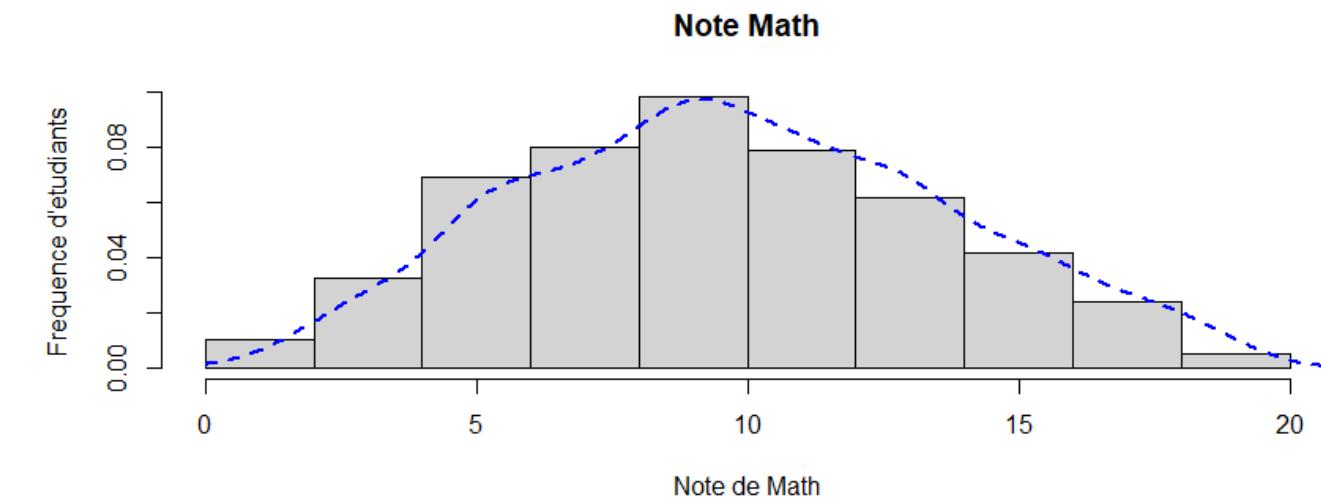


D'après le diagramme en barre et la table de contingence on distingue, que notre population est divisée en deux niveaux, 34% ont obtenu la mention au BAC et 66% ne l'ont pas obtenu.

Développement

Analyse univarié :

- Résultat de l'épreuve de Math du BAC :



	N6	N5	N4	N3	N2	N1
	0.15075377	0.42613065	0.15829146	0.12412060	0.08341709	0.05728643
Intervalle de note	[1,5]	[5,10]	[10,12]	[12,14]	[14,16]	[16,19]

Développement

Analyse univarié :

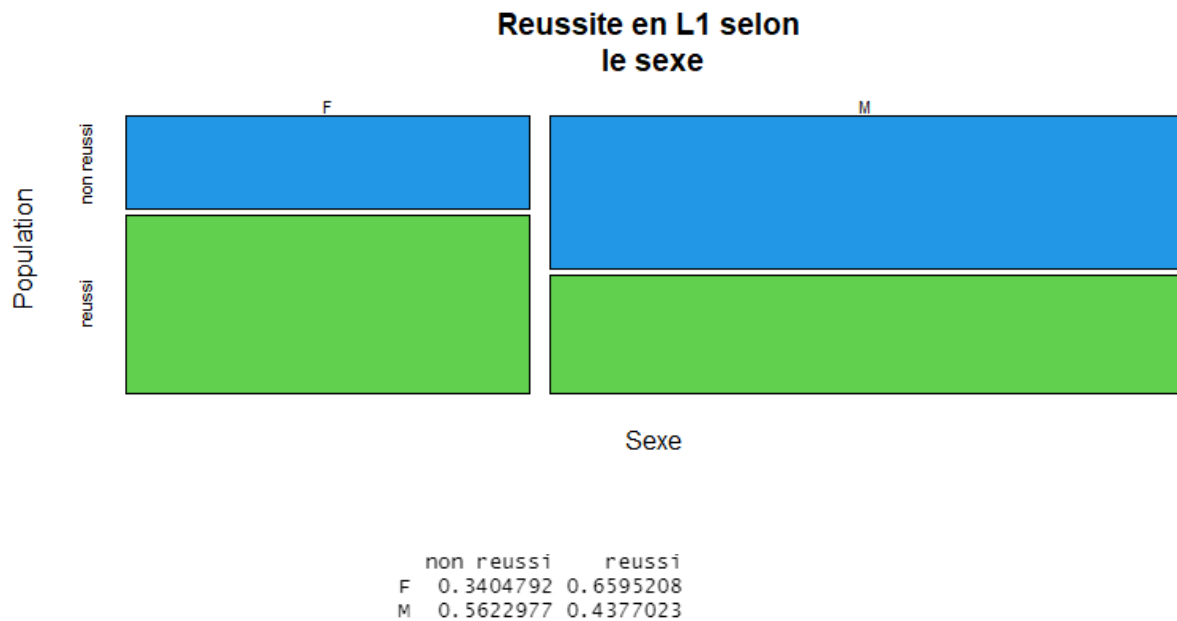
- Résultat de l'épreuve de Math du BAC :

D'après l'histogramme, la boîte à moustaches, ainsi que la table de contingence, on déduit que les notes des élèves étaient entre $[1, 19]$, avec un minimum de $1/20$ et maximum de $19/20$. Avec une moyenne de $9,83$, médiane de 10 , écart type de $3,96$ et coefficient de variance égale à $0,40$. J'ai découpé cet intervalle, en cinq intervalles, qui corresponde à la qualité de mention. Mon premier intervalle est le $[1, 5]$ qui représente $15,07\%$ de l'ensemble des notes, $[5, 10]$ qui représente $42,61\%$ de l'échantillon, $[10, 12]$ qui représente $15,82\%$, $[12, 14]$ qui représente $12,41\%$, $[14, 16]$ qui représente $8,34\%$ et $[16, 19]$ qui représente $5,72\%$. Ce qui correspond à l'histogramme sous forme de pyramide, vu que l'intervalle le plus proche rassemble plus des deux cinquièmes des notes.

Développement

Analyse entre les valeurs qualitatives :

- Réussite en L1 selon le sexe :

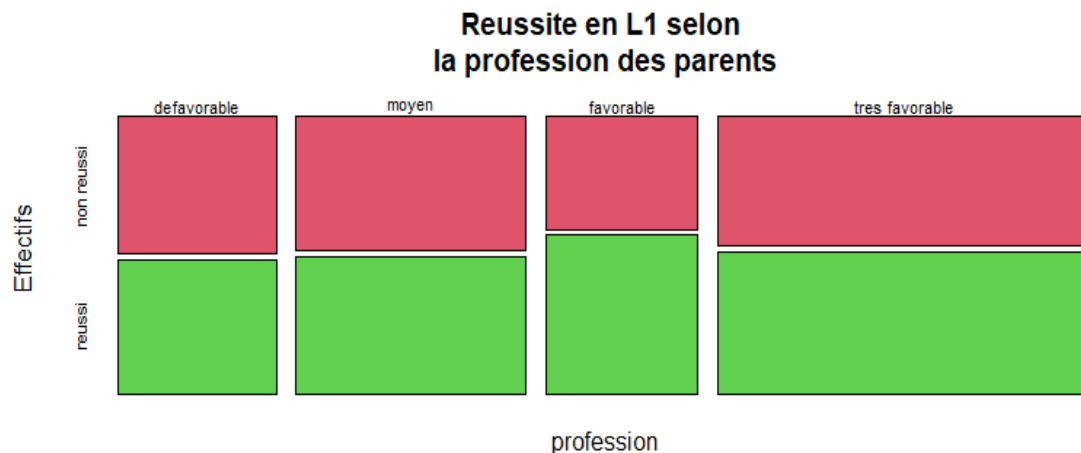


Les mosaïques et la table de contingence précédent nous montre que les femmes ont plus de chance à réussir en L1 que les hommes avec un pourcentage respectivement de 69% et 43%. On a procédé au test KHI2, afin d'être sur qu'il y a une dépendance ou indépendance entre ses deux variables et on a eu un résultat de **p-value inférieure a $2,2 \times 10^{-16}$** , ce qui est une valeur très faible par rapport au niveau de signification de 5%, ce qui implique que les deux variables **sont dépendantes**. Pour comprendre l'intensité entre ses deux valeurs, on a utilisé le test du v de Cramer, on a obtenu le résultat suivant : 0,216, ce qui explique qu'il y a vrai une dépendance entre ses variables, mais elle est à **faible intensité** car le résultat du test est inférieur a 50%.

Développement

Analyse entre les valeurs qualitatives :

- Réussite en L1 selon la profession des parents :



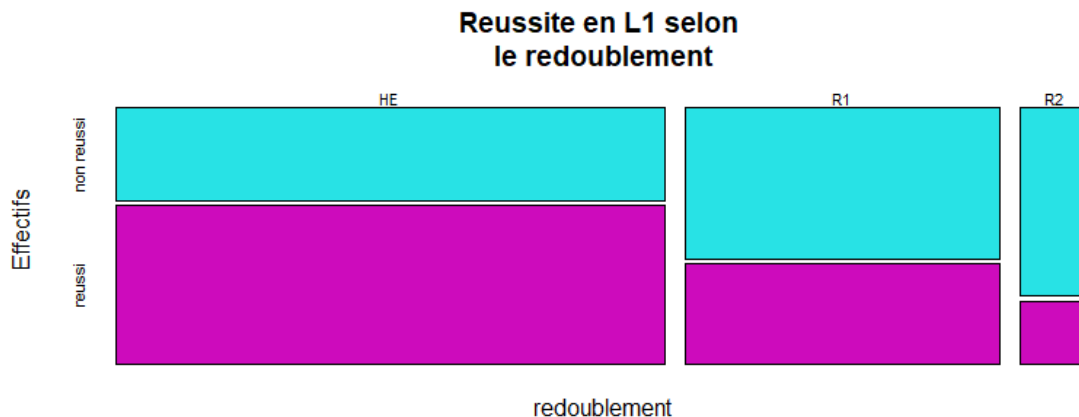
pcsf	reu	
	non reussi	reussi
defavorable	0.5056818	0.4943182
moyen	0.4941176	0.5058824
favorable	0.4154303	0.5845697
tres favorable	0.4759036	0.5240964

Les mosaïques et la table de contingence précédent, nous montre que la profession des parents n'a pas d'influence sur la réussite en L1. On a procédé au test KHI2, afin d'être sûr qu'il y a une dépendance ou indépendance entre ses deux variables et on a eu un résultat de **p-value égale a 0,07**, ce qui est une valeur supérieure au niveau de signification de 5%, ce qui implique que les deux variables **sont indépendantes**. Pour comprendre l'intensité entre ses deux valeurs, on a utilisé le test du v de Cramer, on a obtenu le résultat suivant : 0,058.

Développement

Analyse entre les valeurs qualitatives :

- Réussite en L1 selon le redoublement au BAC :



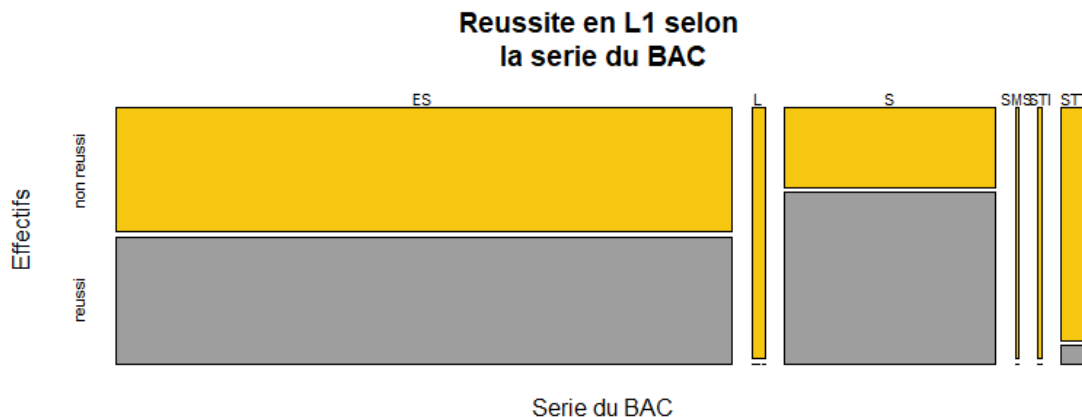
	reu	
	non reussi	reussi
HE	0.3691275	0.6308725
R1	0.6000000	0.4000000
R2	0.7500000	0.2500000

Les mosaïques et la table de contingence précédent nous montre que les étudiants qui redouble, on moins de chance à réussir par rapport aux étudiants qui ont obtenus leurs BAC dès la première année. On a procédé au test KHI2, afin d'être sûr qu'il y a une dépendance ou indépendance entre ses deux variables et on a eu un résultat de **p-value inférieure a $2,2 \times 10^{-16}$** , ce qui est une valeur très faible par rapport au niveau de signification de 5%, ce qui implique que les deux variables **sont dépendantes**. Pour comprendre l'intensité entre ses deux valeurs, on a utilisé le test du v de Cramer, on a obtenu le résultat suivant : 0,265, ce qui explique qu'il y a une dépendance entre ses variables, mais **son intensité est faible**, car le résultat du test est inférieur a 50%.

Développement

Analyse entre les valeurs qualitatives :

- Réussite en L1 selon la série du BAC :



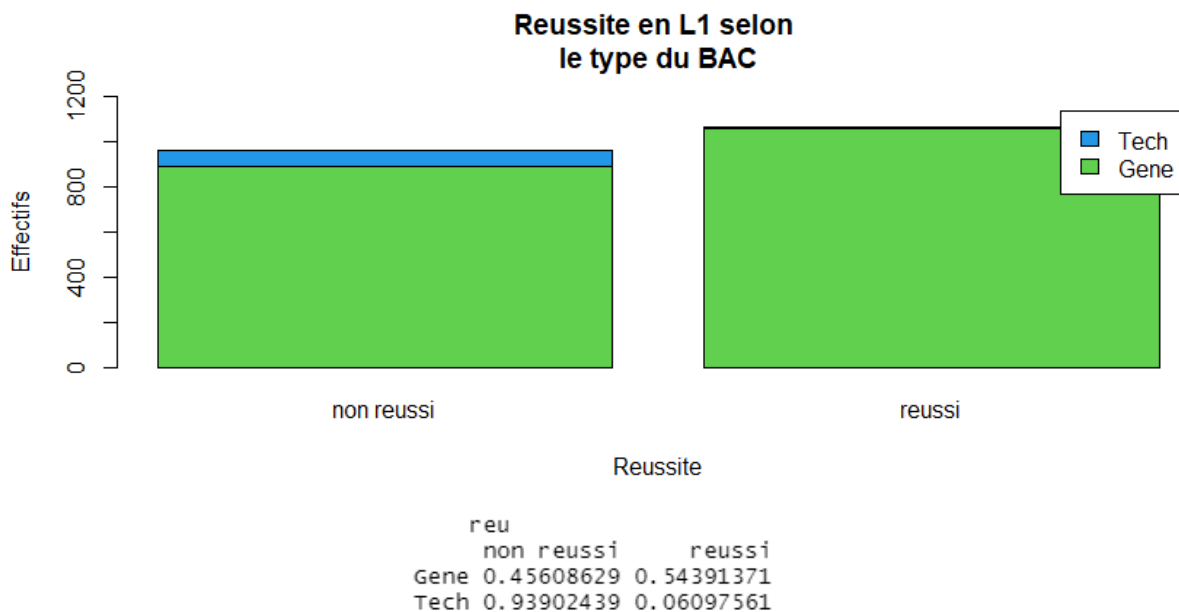
	reu	
	non reussi	reussi
ES	0.49264191	0.50735809
L	1.00000000	0.00000000
S	0.31492843	0.68507157
SMS	1.00000000	0.00000000
STI	1.00000000	0.00000000
STT	0.92537313	0.07462687

Les mosaïques et la table de contingence précédent nous montre que les étudiants de la série du BAC (S et STT), on distingue aussi que ceux de la série (L, SMS et STI) n'ont aucune chance de réussir. On a procédé au test KHI2, afin d'être sûr qu'il y a une dépendance ou indépendance entre ses deux variables et on a eu un résultat de **p-value inferieure a $2,2 \times 10^{-16}$** , ce qui est une valeur très faible par rapport au niveau de signification de 5%, ce qui implique que les deux variables **sont dépendantes**. Pour comprendre l'intensité entre ses deux valeurs, on a utilisé le test du v de Cramer, on a obtenu le résultat suivant : 0,28, ce qui explique qu'il y a vrai une dépendance entre ses variables, mais elle est à **faible intensité** car le résultat du test est inferieur a 50%.

Développement

Analyse entre les valeurs qualitatives :

- Réussite en L1 selon le type du BAC :

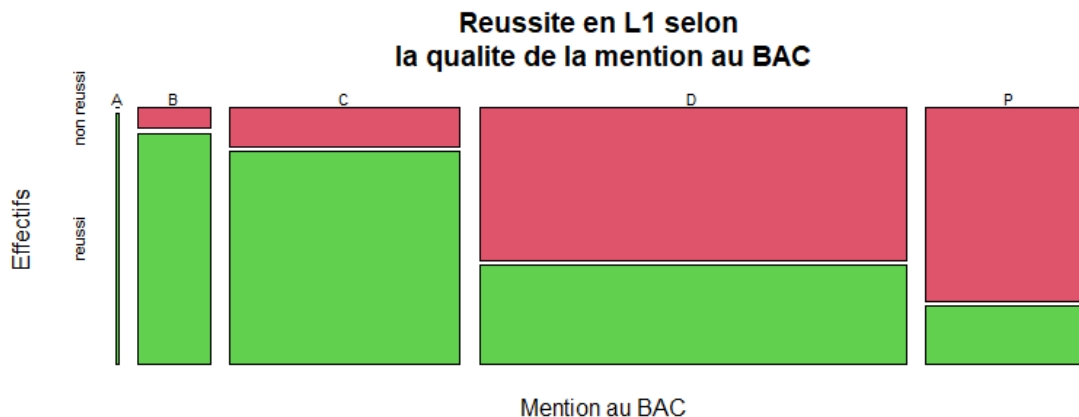


Le digramme en barre et la table de contingence précédent nous montre que les étudiants du type BAC général ont plus de chance à réussir en L1 que les autres avec un pourcentage respectivement de 54% et 6%. On a procédé au test KHI2, afin d'être sûr qu'il y a une dépendance ou indépendance entre ses deux variables et on a eu un résultat de **p-value inférieure à $2,2 \times 10^{-16}$** , ce qui est une valeur très faible par rapport au niveau de signification de 5%, ce qui implique que les deux variables **sont dépendantes**. Pour comprendre l'intensité entre ses deux valeurs, on a utilisé le test du v de Cramer, on a obtenu le résultat suivant : 0,19, ce qui explique qu'il y a une dépendance entre ses variables, mais elle est à **faible intensité**, car le résultat du test est inférieur à 50%.

Développement

Analyse entre les valeurs qualitatives :

- Réussite en L1 selon la qualité de la mention du BAC :



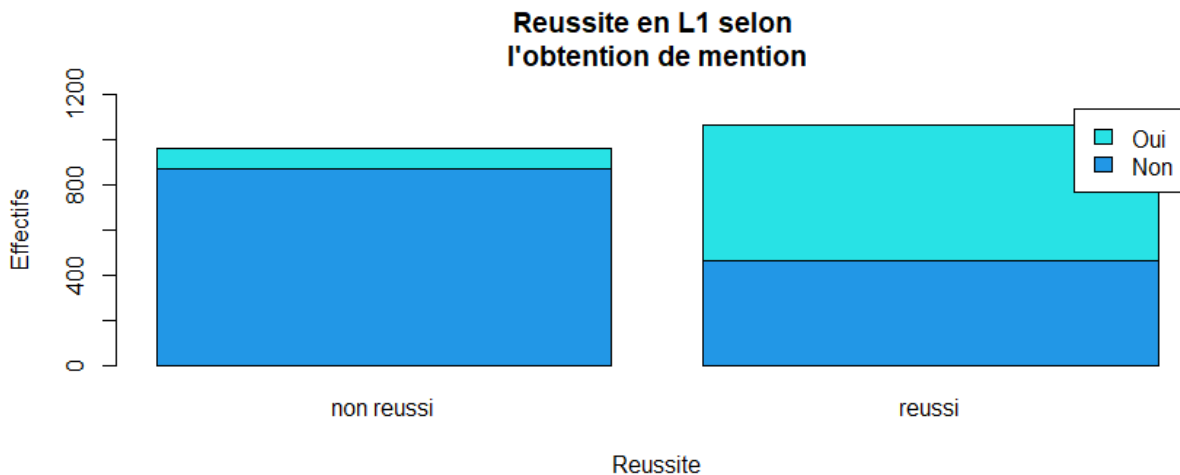
	reu	
	non reussi	reussi
A	0.0000000	1.0000000
B	0.0797546	0.9202454
C	0.1516315	0.8483685
D	0.6076605	0.3923395
P	0.7688172	0.2311828

Les mosaïques et la table de contingence précédent nous montre que plus la qualité de mention au BAC est meilleure, plus l'étudiant a de chance pour réussir son L1. On a procédé au test KHI2, afin d'être sûr qu'il y a une dépendance ou indépendance entre ses deux variables et on a eu un résultat de **p-value inférieure a $2,2 \times 10^{-16}$** , ce qui est une valeur très faible par rapport au niveau de signification de 5%, ce qui implique que les deux variables **sont dépendantes**. Pour comprendre l'intensité entre ses deux valeurs, on a utilisé le test du v de Cramer, on a obtenu le résultat suivant : 0,5, ce qui explique qu'il y a vrai une dépendance entre ses variables, avec une **moyen intensité** entre eux, car le résultat du test est égal a 50%.

Développement

Analyse entre les valeurs qualitatives :

- Réussite en L1 selon l'obtention de la mention au BAC :



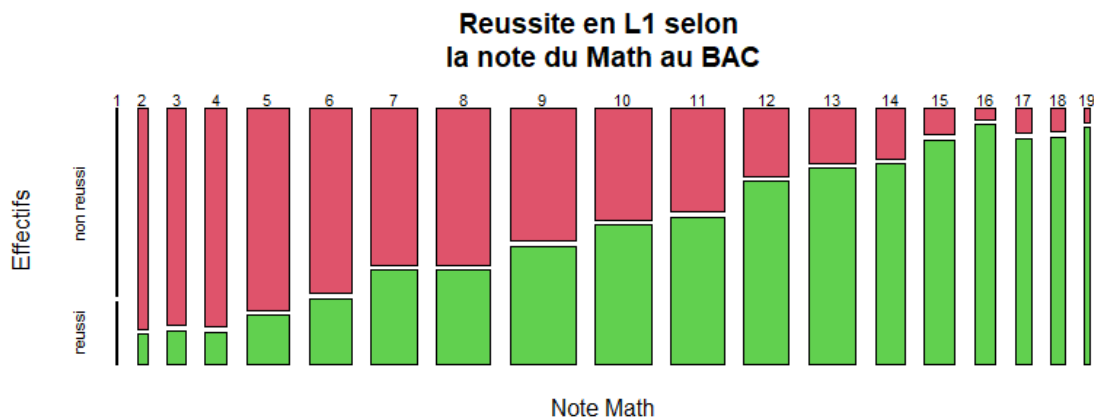
	reu	
	non reussi	reussi
Non	0.6524664	0.3475336
Oui	0.1331404	0.8668596

Le diagramme en barre et la table de contingence précédent nous montre que les étudiants qui ont obtenu leur BAC avec mention sont les plus favorisé à réussir leur L1. On a procédé au test KHI2, afin d'être sûr qu'il y a une dépendance ou indépendance entre ses deux variables et on a eu un résultat de **p-value inférieure a $2,2 \times 10^{-16}$** , ce qui est une valeur très faible par rapport au niveau de signification de 5%, ce qui implique que les deux variables **sont dépendantes**. Pour comprendre l'intensité entre ses deux valeurs, on a utilisé le test du v de Cramer, on a obtenu le résultat suivant : 0,49, ce qui explique qu'il y a vrai une dépendance entre ses variables, avec une **moyen intensité** entre eux, car le résultat du test est égal à 50%.

Développement

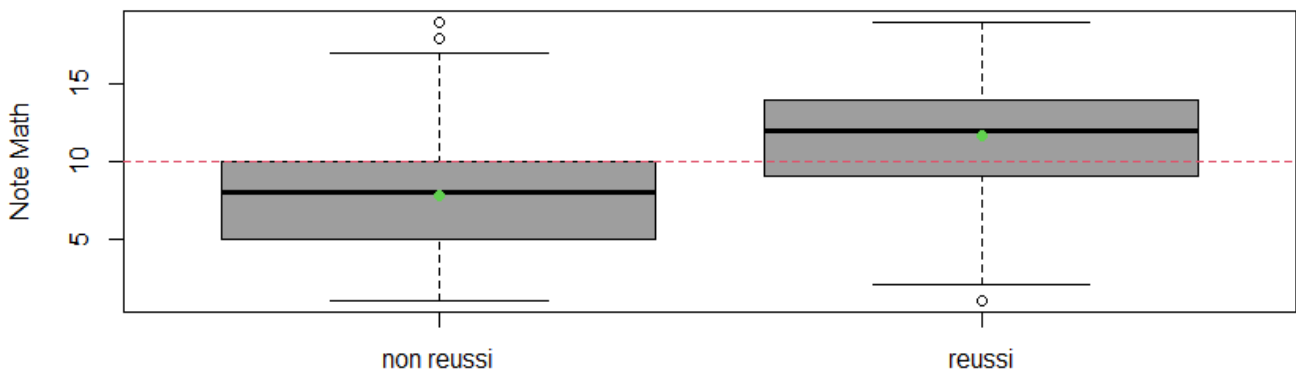
Analyse entre une variable qualitative et une variable quantitative :

- Réussite en L1 selon l'obtention de la note de l'épreuve de Math au BAC :



	reussi	non reussi
1	0.250000000	0.750000000
2	0.121212121	0.878787879
3	0.135593222	0.864406778
4	0.128571429	0.871428571
5	0.195652174	0.804347826
6	0.262773372	0.737226628
7	0.376712333	0.623287667
8	0.375722544	0.624277456
9	0.471428571	0.528571429
10	0.554945055	0.445054945
11	0.588235294	0.411764706
12	0.731034483	0.268965517
13	0.782894737	0.217105263
14	0.800000000	0.200000000
15	0.896907222	0.103092778
16	0.956521739	0.043478261
17	0.903846154	0.096153846
18	0.906976744	0.093023256
19	0.947368421	0.052631579

Distribution des effectifs selon les notes de math



Développement

- Analyse entre une variable qualitative et une variable quantitative :
- Réussite en L1 selon l'obtention de la note de l'épreuve de Math au BAC :

Les mosaïques, la boîte à moustache et la table de contingence précédent nous montre que les étudiants ayant les bonnes notes à l'épreuve des maths au BAC, sont les étudiants qui réussissent mieux en L1, selon la table de contingence précédente, dès que la note de l'épreuve augmente, la fréquence liée à la réussite augmente et vice versa. Via la boîte à moustache, on observe que les trois quarts des étudiants qui ont eu moins ou égale à 10 à l'épreuve des math au BAC, n'ont pas pu réussir leur première année. On a procédé au test KHI2, afin d'être sûr qu'il y a une dépendance ou indépendance entre ses deux variables et on a eu un résultat de **p-value inférieure à $2,2 \times 10^{-16}$** , ce qui est une valeur très faible par rapport au niveau de signification de 5%, ce qui implique que les deux variables **sont dépendantes**. Pour comprendre l'intensité entre ses deux valeurs, on a utilisé le test du χ^2 de Cramer, on a obtenu le résultat suivant : 0,5, ce qui explique qu'il y a vraiment une dépendance entre ses variables, avec une **moyenne intensité** entre eux, car le résultat du test est égal à 50%.

Conclusion

- Conclusion :

D'après les études sur la base de données, on distingue que la réussite des étudiants en L1 est dépendante de toutes les variables sauf la profession que son influence est presque nul. Parmi ses variables dépendantes y ont à qui sont plus intense que les autres, comme je le montre avec le tableau suivant :

Variable/Réussite	P-value	Etat de la relation	V-cramer	Intensité
Type de BAC	<2,2 ^e -16	Dépendance	0,19	Faible
Sexe	<2,2 ^e -16	Dépendance	0,216	Faible
Redoublement	<2,2 ^e -16	Dépendance	0,265	Faible
Série du BAC	<2,2 ^e -16	Dépendance	0,28	Faible
Obtention de mention	<2,2 ^e -16	Dépendance	0,49	Moyenne
Qualité de mention	<2,2 ^e -16	Dépendance	0,5	Moyenne
Note a l'épreuve de Math au BAC	<2,2 ^e -16	Dépendance	0,5	Moyenne
Profession des parents	0,07	Indépendance	0,05	Faible

