

## Initiation à R et Dataviz

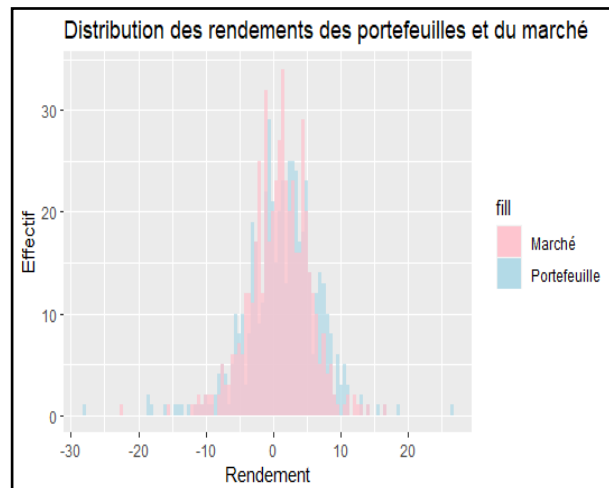
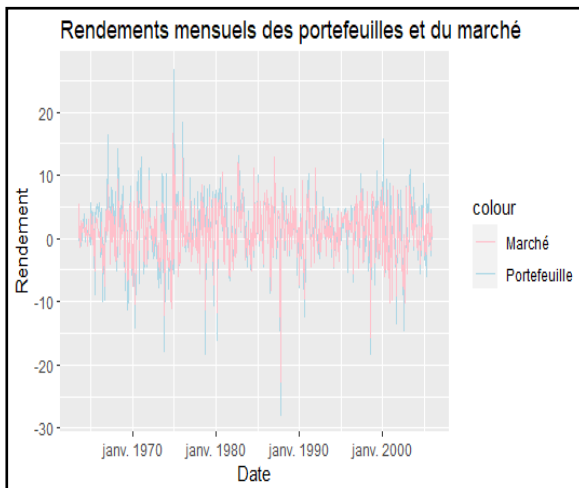
### Projet Final

#### Exercice 1 : Estimation modèle MEDAF sur plusieurs périodes

Période 01 : De juillet 1963 à décembre 2005 :

Statistiques descriptives, graphiques et estimation du modèle linéaire du MEDAF :

```
> summary(donnees_periode_1)
rend_pf_ME1_BM2   rend_pf_marche      RF
Min.   :-27.892   Min.   :-22.6400   Min.   :0.0600
1st Qu.: -1.348   1st Qu.: -1.7425   1st Qu.:0.3200
Median :  1.571   Median :  1.2150   Median :0.4300
Mean    :  1.353   Mean    :  0.9451   Mean    :0.4716
3rd Qu.:  4.702   3rd Qu.:  3.9500   3rd Qu.:0.5800
Max.    : 26.742   Max.    : 16.6100   Max.    :1.3500
> ecart_type_ME1_BM2   > ecart_type_pf_marche
[1] 4.430562             [1] 5.369225
```



```
> summary(modele_medaf_1)

Call:
lm(formula = rend_pf_ME1_BM2 ~ rend_pf_marche + RF, data = donnees_periode_1)

Residuals:
    Min       1Q   Median       3Q      Max
-11.1111  -1.5388   0.0219   1.4240  12.2451

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.63240    0.27701   2.283   0.0228 *
rend_pf_marche 1.04802    0.02693  38.921 <2e-16 ***
RF          -0.57219    0.52552  -1.089   0.2768
---

```

```

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.69 on 507 degrees of freedom
Multiple R-squared:  0.75,    Adjusted R-squared:  0.7491
F-statistic: 760.7 on 2 and 507 DF, p-value: < 2.2e-16

```

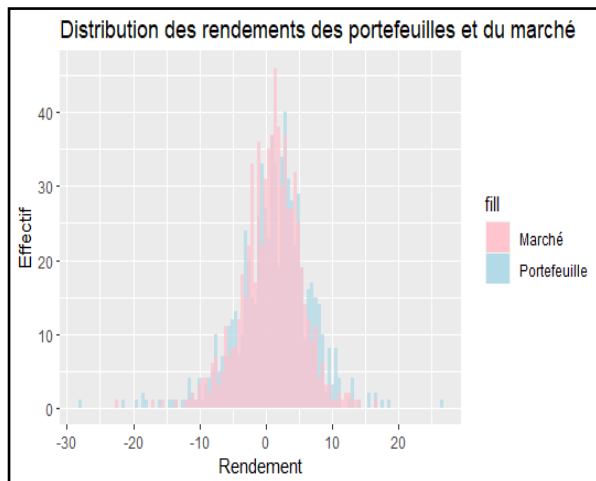
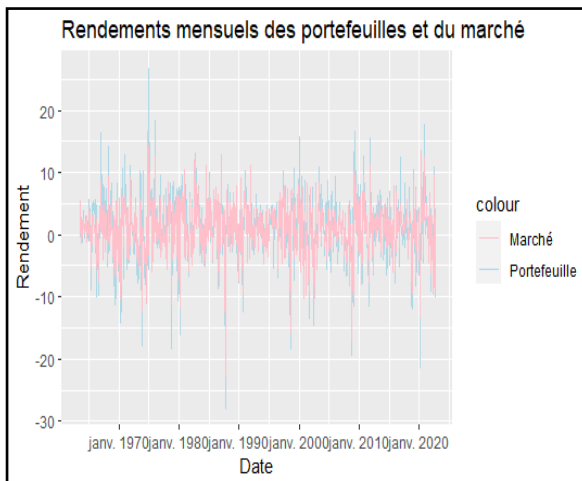
## Période 02 : De juillet 1963 à septembre 2022 :

### Statistiques descriptives, graphiques et estimation du modèle linéaire du MEDAF :

```

> summary(donnees_periode_2)
rend_pf_ME1_BM2  rend_pf_marche      RF
Min.   :-27.892   Min.   :-22.6400   Min.   :0.0000
1st Qu.: -1.696   1st Qu.: -1.6950   1st Qu.:0.1400
Median :  1.485   Median :  1.2600   Median :0.3800
Mean    :  1.214   Mean    :  0.9057   Mean    :0.3624
3rd Qu.:  4.532   3rd Qu.:  3.7550   3rd Qu.:0.5100
Max.    : 26.742   Max.    : 16.6100   Max.    :1.3500
> ecart_type_ME1_BM2  > ecart_type_pf_marche
[1] 4.474923           [1] 5.493661

```



```

> summary(modele_medaf_2)

Call:
lm(formula = rend_pf_ME1_BM2 ~ rend_pf_marche + RF, data = donnees_periode_2)

Residuals:
    Min       1Q   Median       3Q      Max
-11.1811  -1.5718  -0.0706   1.4229  12.2957

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.21711    0.16543   1.312   0.190
rend_pf_marche 1.08232    0.02179  49.678 <2e-16 ***
RF            0.04635    0.36363   0.127   0.899
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

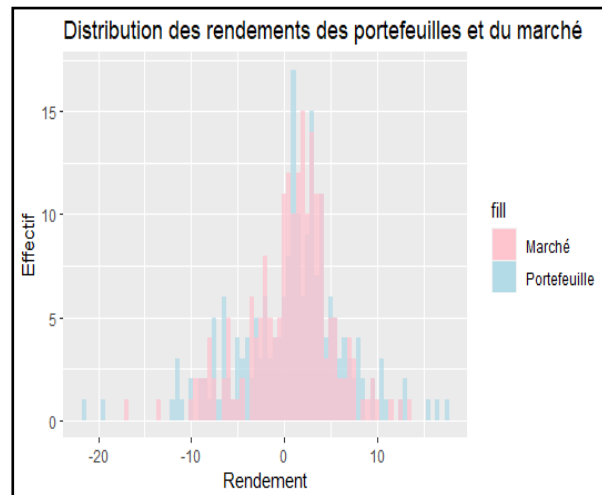
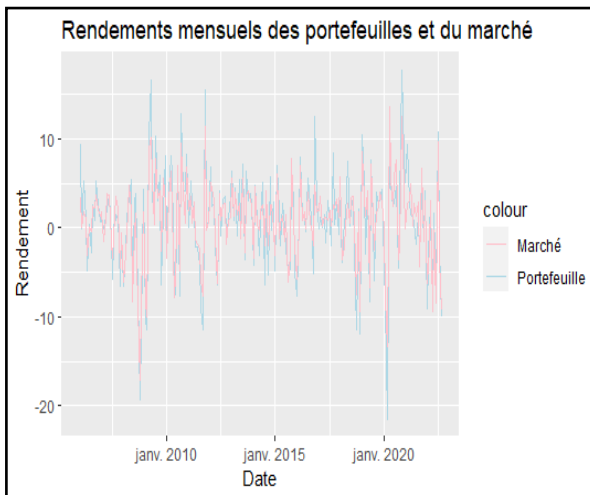
Residual standard error: 2.597 on 708 degrees of freedom
Multiple R-squared:  0.7772,    Adjusted R-squared:  0.7765
F-statistic: 1235 on 2 and 708 DF, p-value: < 2.2e-16

```

### Période 03 : De janvier 2006 à septembre 2022 :

#### Statistique descriptive, graphiques et estimation du modèle linéaire du MEDAF :

```
> summary(donnees_periode_3)
rend_pf_ME1_BM2  rend_pf_marche      RF
Min.   :-21.4683  Min.   :-17.1500  Min.   :0.00000
1st Qu.: -2.1295  1st Qu.: -1.5600  1st Qu.:0.00000
Median :  1.3336  Median :  1.3400  Median :0.01000
Mean   :  0.8617  Mean   :  0.8056  Mean   :0.08527
3rd Qu.:  4.0535  3rd Qu.:  3.4000  3rd Qu.:0.14000
Max.   : 17.7379  Max.   : 13.6500  Max.   :0.44000
> ecart_type_ME1_BM2      > ecart_type_pf_marche
[1] 4.595261                [1] 5.79632
```



```
> summary(modele_medaf_3)
Call:
lm(formula = rend_pf_ME1_BM2 ~ rend_pf_marche + RF, data = donnees_periode_3)

Residuals:
    Min       1Q   Median       3Q      Max
-6.0212 -1.6343 -0.0713  1.2101  6.9084

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  -0.03087    0.19830  -0.156   0.876
rend_pf_marche  1.15798    0.03547  32.648 <2e-16 ***
RF            -0.47280    1.27122  -0.372   0.710
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.298 on 198 degrees of freedom
Multiple R-squared:  0.8444,    Adjusted R-squared:  0.8428
F-statistic: 537.2 on 2 and 198 DF, p-value: < 2.2e-16
```

### Comparaison entre les trois périodes :

	Rendement de portefeuille ME1 BM2	Rendement du portefeuille de marché
<b>Période 01</b>	Mean : 1,353 Median : 1,571 Max : 26,742 Sd : 4,430562	Mean : 0,9451 Median : 1,215 Max : 16,61 Sd : 5,369225
<b>Période 02</b>	Mean : 1,214 Median : 1,485 Max : 26,742 Sd : 4,474923	Mean : 0,9057 Median : 1,26 Max : 16,61 Sd : 5,493661
<b>Période 03</b>	Mean : 0,8617 Median : 1,3336 Max : 17,7379 Sd : 4,595261	Mean : 0,8056 Median : 1,34 Max : 13 ,65 Sd : 5,79632

On remarque que le rendement moyen du portefeuille ME1 BM2 est relativement supérieur sur la période de juillet 1963 à décembre 2005 par rapport aux autres périodes. De plus, le rendement moyen du portefeuille ME1 BM2 est relativement supérieur sur la période de juillet 1963 à septembre 2022 par rapport à la période qui s'étend de janvier 2006 à septembre 2022 => On peut conclure que le rendement moyen du portefeuille ME1 BM2 a significativement baissé à partir de janvier 2006.

On remarque aussi que l'écart-type sur la période de janvier 2006 à septembre 2022 est légèrement plus élevé que celui de la période de juillet 1963 à décembre 2005, qui peut indiquer une certaine augmentation de volatilité des investissements qui peut être vu comme un risque pour un certain type d'investisseur. Concernant le rendement du portefeuille de marché, on peut voir qu'il suit presque la même tendance que le rendement du portefeuille ME1 BM2 (chose qu'on peut remarquer sur les graphiques ci-dessus), avec une baisse significative du rendement moyen et une légère augmentation de volatilité à partir de janvier 2006.

	Estimation du modèle du MEDAF
<b>Période 01</b>	<pre> (Intercept)    0.63240 rend_pf_marche 1.04802 RF              -0.57219 Multiple R-squared:  0.75 F-statistic: 760.7 p-value: &lt; 2.2e-16 </pre>

<b>Période 02</b>	<pre>(Intercept)    0.21711 rend_pf_marche 1.08232 RF             0.04635 Multiple R-squared:  0.7772 F-statistic: 1235 p-value: &lt; 2.2e-16</pre>
<b>Période 03</b>	<pre>(Intercept)    -0.03087 rend_pf_marche  1.15798 RF             -0.47280 Multiple R-squared:  0.8444 F-statistic: 537.2 p-value: &lt; 2.2e-16</pre>

Sur les trois périodes, on remarque que le coefficient de détermination est suffisamment élevé, qui veut dire que nos variables explicatives expliquent une partie importante de notre modèle.

On remarque aussi que l'estimateur du rendement du portefeuille de marché est positif, c'est-à-dire qu'il y a une relation positive entre le rendement du portefeuille du marché et le rendement de portefeuille ME1 BM2 et aussi l'estimateur est statistiquement significatif qui veut dire qu'il y a une forte corrélation entre eux.

On peut voir qu'il y a une relation négative entre le rendement sans risque et le rendement de portefeuille ME1 BM2.

La p-value du modèle est inférieur à 5% pour les trois périodes => c'est-à-dire que l'estimation du modèle du MEDAF est statistiquement significatif.

En conclusion, on peut dire que les résultats de notre estimation viennent confirmer notre premier raisonnement qui est que le rendement du portefeuille ME1 BM2 et le rendement du portefeuille du marché suivent la même tendance et presque la même distribution et la preuve est qu'il y a une relation positive entre les deux (c-à-d que les rendements augmentent et baissent ensemble, chose qu'on peut remarquer sur les graphiques) et aussi la forte corrélation qui existe entre les deux rendements.

## **Exercice 2 : Estimation du modèle Fama-French à 3 facteurs**

Voir script R.

## Exercice 3 : Tabagisme et âge

### Question 1 :

Nous avons à notre disposition un ensemble de données médicales sur les maladies pulmonaires et le tabac. Le fichier comprend 7 variables qui nous donnent des informations sur les individus étudiés (chacun a un numéro ID affecté). Les informations données sont l'âge, le genre, la situation conjugale, le niveau de consommation de tabac, l'exposition au tabagisme passif et enfin la présence de problème pulmonaires. La variable ID n'est pas pertinente à étudier car elle donne simplement un numéro à chaque individu.

Parmi ces variables 4 sont qualitatives :

\_Le sexe qui renvoi « homme » ou « femme »

\_La situation conjugale de l'individu qui peut être « marie », « en couple », « célibataire » ou « veuf »

\_Le tabagisme passif qui nous dit si l'individu y est exposé ou non avec « TRUE » ou « FALSE »

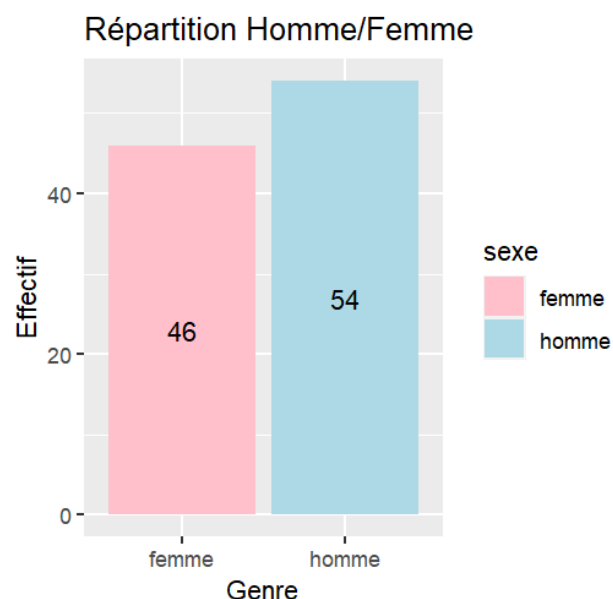
\_La présence éventuelle de problème pulmonaires avec « TRUE » ou « FALSE »

Par conséquent 2 variables quantitatives :

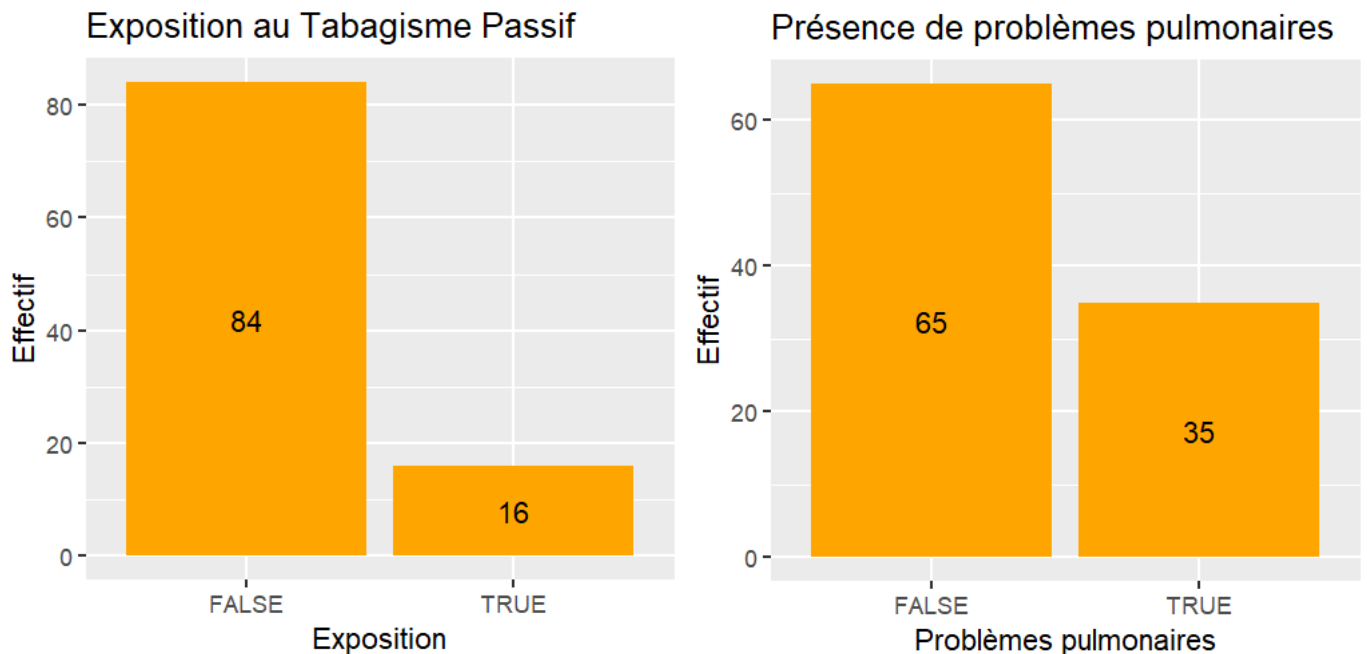
\_L'âge de l'individu qui renvoi un endroit un nombre entier

\_Le niveau de consommation de tabac de chaque individu traduit par un nombre entier allant de 0 à 14

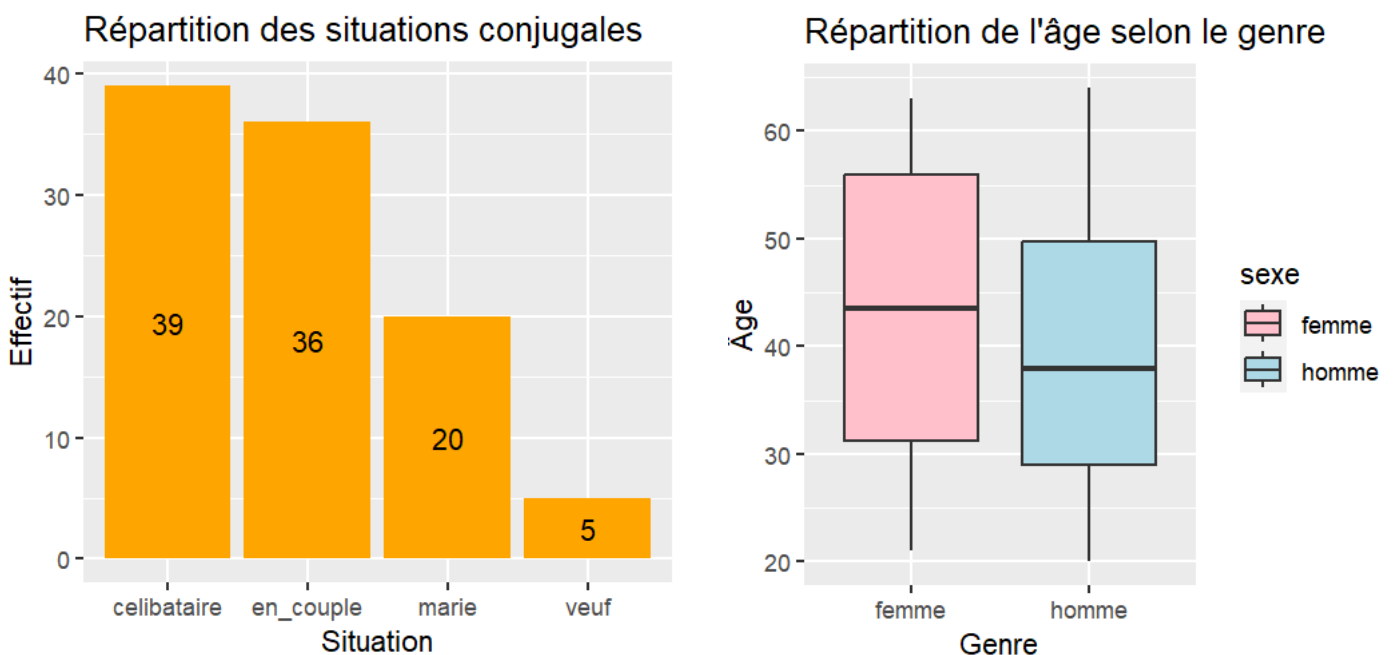
Nous pouvons désormais effectuer une analyse statistique des données :



On remarque qu'il y a 100 individus dans notre étude, ainsi chaque valeur peut être interprétée comme un pourcentage. Il y a 46 femmes et 54 hommes, la parité est relativement respectée.

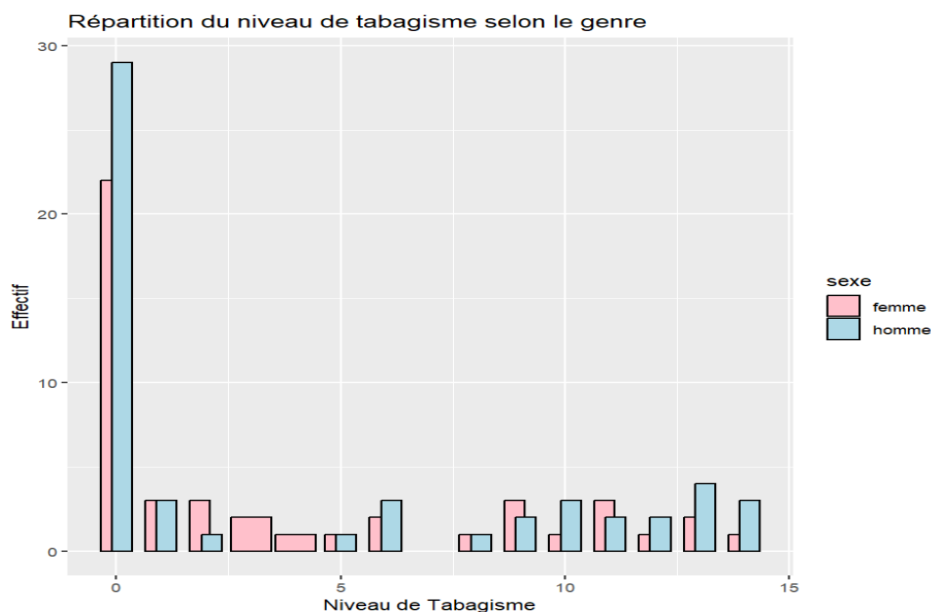


Ci-dessus, deux graphiques nous donnant des informations sur l'exposition au tabagisme passif et sur la présence de problèmes pulmonaires chez les individus de l'étude. On remarque que 84 personnes ne sont pas exposées au tabagisme passif et par conséquent 16 y sont exposés. Ensuite, 65 personnes n'ont pas de problèmes pulmonaires et 35 en ont.



Ci-dessus à gauche, une représentation de la répartition des différentes situations conjugales des individus, 39 sont célibataires, 36 sont en couple, 20 mariés et enfin 5 sont veufs.

Enfin, en haut à droite un graphique « en boîtes à moustache » nous donnant des informations sur l'âge des individus selon leur genre. On remarque que l'étendu est plus élevée chez les hommes mais la médiane et de manière générale les femmes sont plus âgées que les hommes dans notre étude.



Ici, un histogramme en barre représentant la répartition du niveau de tabagisme de 0 à 14 par genre. On remarque que les hommes fument moins que les femmes dans notre étude. Ensuite, les effectifs sont relativement similaires selon le genre.

age	sexe	situation	tabac	tabagisme_passif
Min. :20.00	femme:46	celibataire:39	Min. : 0.0	FALSE:84
1st Qu.:29.75	homme:54	en_couple :36	1st Qu.: 0.0	TRUE :16
Median :41.00		marie :20	Median : 0.0	
Mean :41.38		veuf :5	Mean :3.9	
3rd Qu.:52.25			3rd Qu.: 9.0	
Max. :64.00			Max. :14.0	
probleme_pulmonaire				
FALSE:65				
TRUE :35				



## Question 2 :

Le modèle linéaire choisi pour étudier le lien entre l'âge de la personne et le niveau de tabagisme est avec les variables portant sur l'âge, le niveau de tabagisme, le genre et la présence de problèmes pulmonaires. L'objectif étant de déterminer l'impact réel de l'âge et du genre de l'individu sur le niveau de tabagisme et de constater si les problèmes pulmonaires ont un lien significatif avec le modèle linéaire. Ci-dessous le résumé du modèle étudié :

```
Call:
lm(formula = tabac ~ age + sexe + probleme_pulmonaire, data = data)

Residuals:
    Min       1Q   Median       3Q      Max
-9.4857 -1.4225 -0.1925  1.3640  6.4332

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    4.41445    1.12181   3.935 0.000157 ***
age           -0.08344    0.02416  -3.454 0.000823 ***
sexehomme     -0.17884    0.62100  -0.288 0.773983
probleme_pulmonaireTRUE  8.67115    0.65495  13.239 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.049 on 96 degrees of freedom
Multiple R-squared:  0.6477,    Adjusted R-squared:  0.6367
F-statistic: 58.83 on 3 and 96 DF,  p-value: < 2.2e-16
```

Pour commencer, on remarque que pour toutes les variables égales à 0 le niveau de tabagisme estimé est de 4.414. Selon le modèle, une année supplémentaire dans l'âge de l'individu traduit une diminution de -0.834 unité sur le niveau de tabagisme, cela est significatif de part la faible pvalue. De plus, être un homme traduit une diminution de -0.178 unité sur le niveau de tabagisme mais cela n'est pas significatif car la pvalue est élevée (0.774). Quant aux problèmes pulmonaires, la pvalue est très faible donc la variable est significative et montre que la présence de problème augmente de 8.67 unité le niveau de tabagisme.

Pour conclure, dans notre modèle le genre n'est pas significatif. Cependant, l'âge a une relation significative et négative avec le niveau de tabagisme. De plus, la présence de problèmes pulmonaires est fortement lié à la consommation de tabac. Enfin, le niveau de tabagisme est expliqué par environ 65% des variables du modèle linéaire qui est lui-même significatif de part sa pvalue extrêmement faible.