

TP2 Statistiques 2

Juhyun Park, Angie Pineda, Nicolas Brunel

18 février 2022

L'analyse exploratoire des données: révision et extension

Air quality monitoring

[Airparif](#) exploite un système de surveillance de la qualité de l'air avec un réseau de sites dans la région de la capitale (Ile de France) sur lesquels les mesures de la qualité de l'air sont effectuées automatiquement. Ces mesures sont utilisées pour résumer les niveaux actuels de pollution atmosphérique, pour prévoir les niveaux futurs et pour fournir des données pour la recherche scientifique, contribuant à l'évaluation des risques pour la santé et des impacts environnementaux des polluants atmosphériques.

Nous examinerons *l'ozone troposphérique* (O_3). Ce polluant n'est pas émis directement dans l'atmosphère, mais est produit par des réactions chimiques entre le dioxyde d'azote (NO_2), les hydrocarbures et la lumière du soleil.

Nous nous concentrerons sur les données de deux sites de surveillance: un site urbain à Neuilly-sur-seine (**NEUIL**) et un site rural (**RUR.SE**) près de la forêt de Fontainebleau.

Les principales questions d'intérêt sont

- Comment, le cas échéant, la distribution des mesures de l'ozone varie-t-elle entre les sites urbains et ruraux?
- Comment, le cas échéant, la distribution des mesures d'ozone est-elle affectée par saison?

Les données de chaque site sont des mesures quotidiennes de la concentration moyenne horaire maximale de O_3 enregistrée en microgrammes par mètre cube ($\mu g/m^3$), de 2014 à 2019 inclusivement. Pour nous concentrer sur la question de la saison, nous comparons les données de *hiver* (novembre-février inclus) (*winter_ozone.csv*) et *été* (mai - août inclus) (*summer_ozone.csv*).

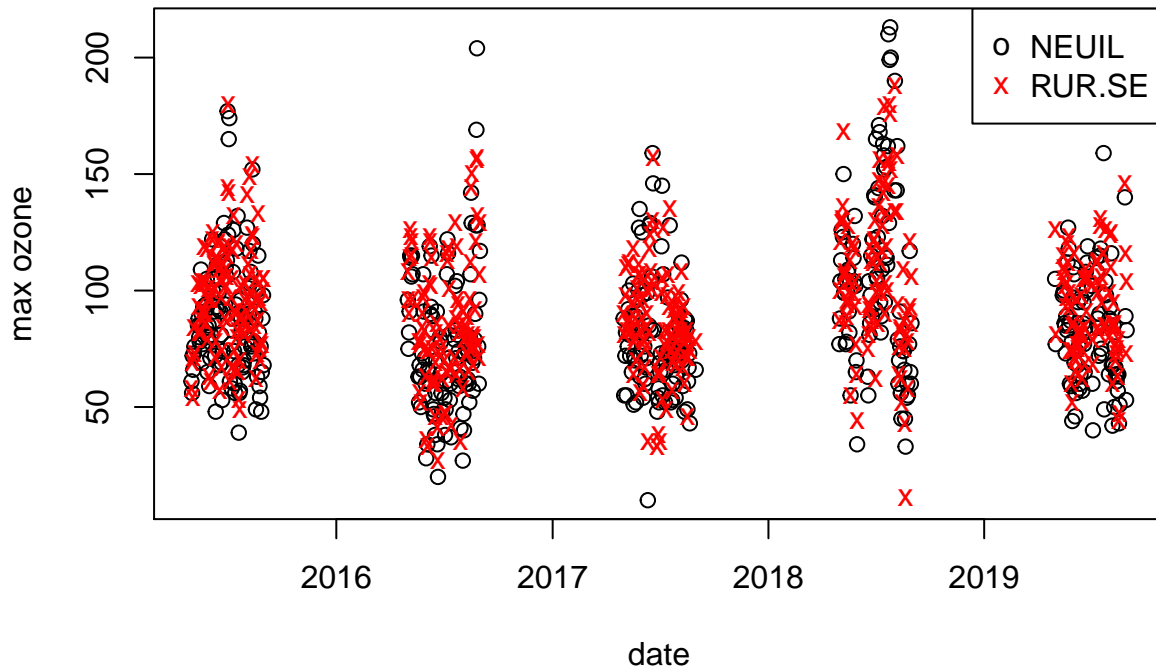
Import data from a .csv file

Nous importons le jeu de données d'été et effectuons la vérification initiale avec la commande **str**, **head** et **tail**. La fonction **names** imprime les noms des variables de colonne et **summary** donne les statistiques récapitulatives par colonne.

```
ozone.summer = read.csv("summer_ozone.csv")
str(ozone.summer) # structure
names(ozone.summer) # names of the variables
head(ozone.summer) # first few observations
tail(ozone.summer) # last few observations
ozone.summer[1:10,] # first 10 observations
summary(ozone.summer) # summary
ozone.summer$date = as.Date(ozone.summer$date2) # transform chr to date format
## plot
plot(ozone.summer$date, ozone.summer$NEUIL,
      xlab="date", ylab="max ozone", main="summer max ozone")
```

```
points(ozone.summer$date, ozone.summer$RUR.SE, col="red", pch = "x")
legend("topright", legend = c("NEUIL", "RUR.SE"),
      col=c("black","red"), pch=c("o","x"))
```

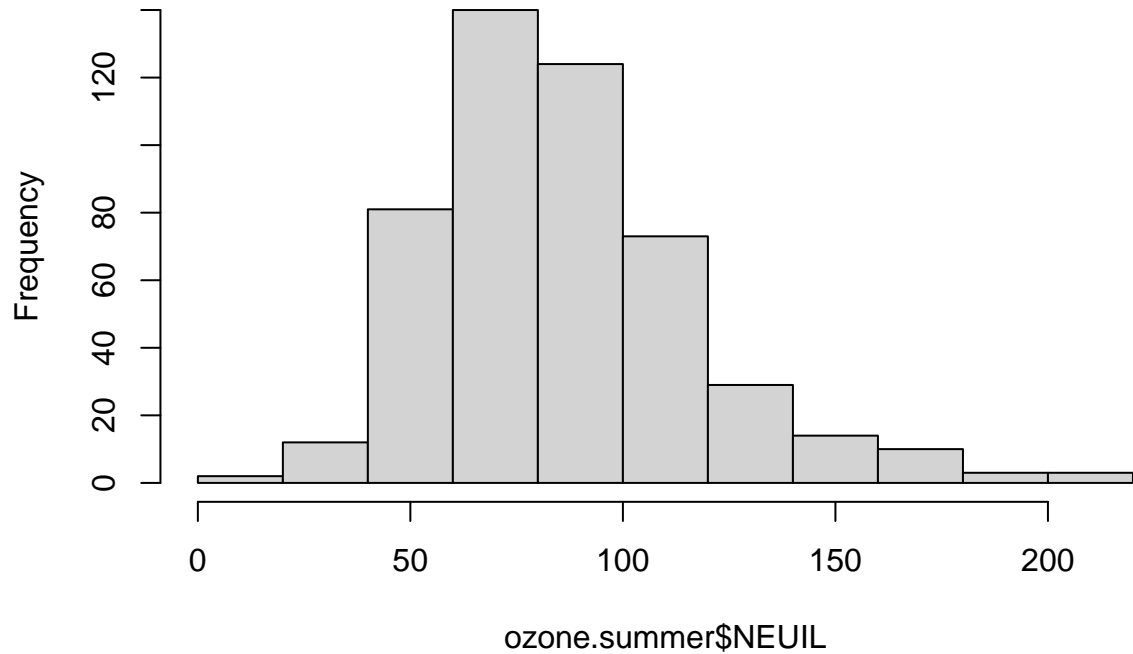
summer max ozone



1. Quelle est la taille des données? Y a-t-il une différence entre les sites urbains et ruraux? Y a-t-il une variation annuelle?
2. Faites les histogrammes des données d'été sur l'ozone pour les deux sites et comparez-les. Pour les rendre comparables sur l'échelle de densité, vous pouvez utiliser l'option *prob = TRUE*.

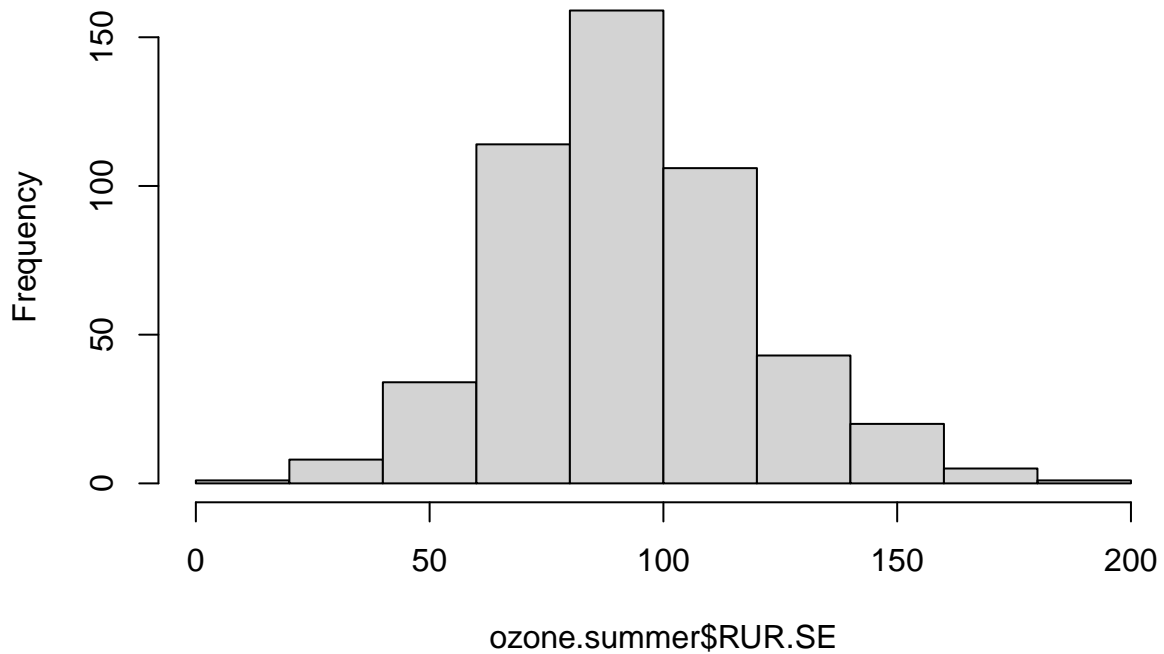
```
## histogram
hist(ozone.summer$NEUIL)
```

Histogram of ozone.summer\$NEUIL



```
hist(ozone.summer$RUR.SE)
```

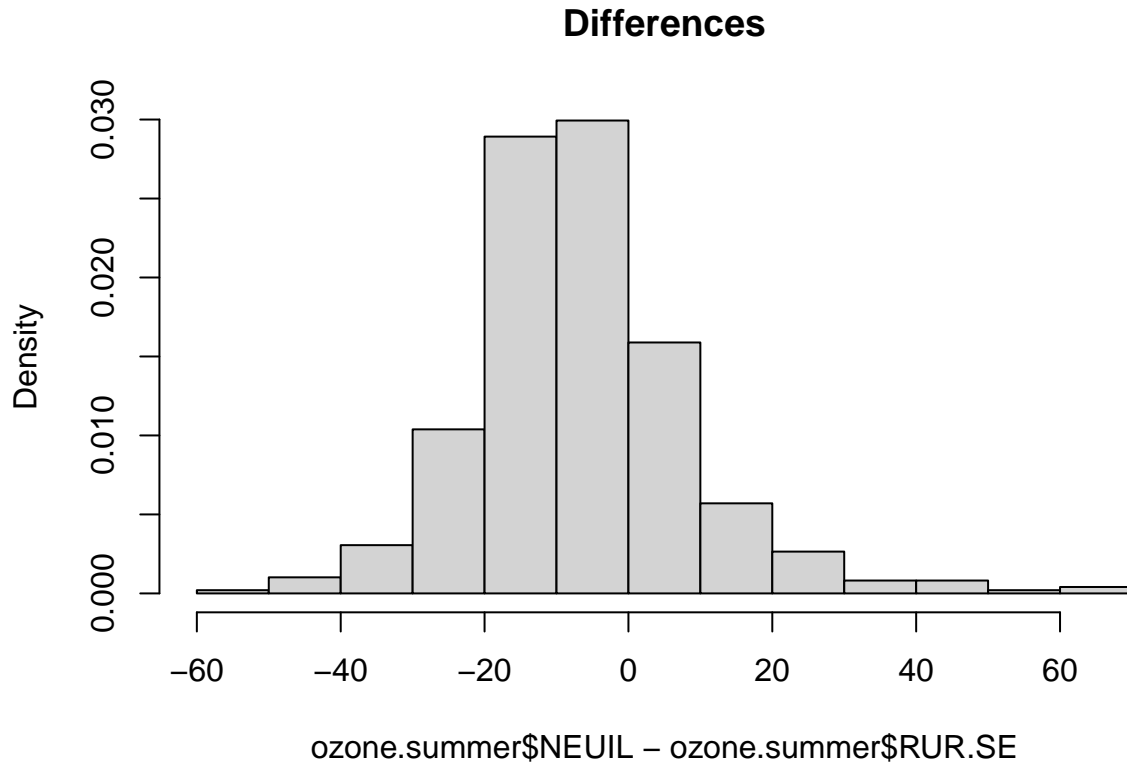
Histogram of ozone.summer\$RUR.SE



Une partie de la variabilité de chacune des distributions est due aux conditions climatiques, qui seront similaires puisqu'elles sont relativement proches les unes des autres. Comme nous nous intéressons à la différence des distributions, nous n'avons pas nécessairement à regarder la distribution séparée elle-même.

Nous désignons les données sur l'ozone du site urbain par x_1, \dots, x_n et le site rural par y_1, \dots, y_n , l'indice indiquant les n jours différents pour lesquels nous avons des mesures. L'histogramme ci-dessous montre la différence $d_i = x_i - y_i$ pour $i = 1, \dots, n$ pour les jours d'été.

```
hist(ozone.summer$NEUIL - ozone.summer$RUR.SE, prob=TRUE, main = "Differences")
```

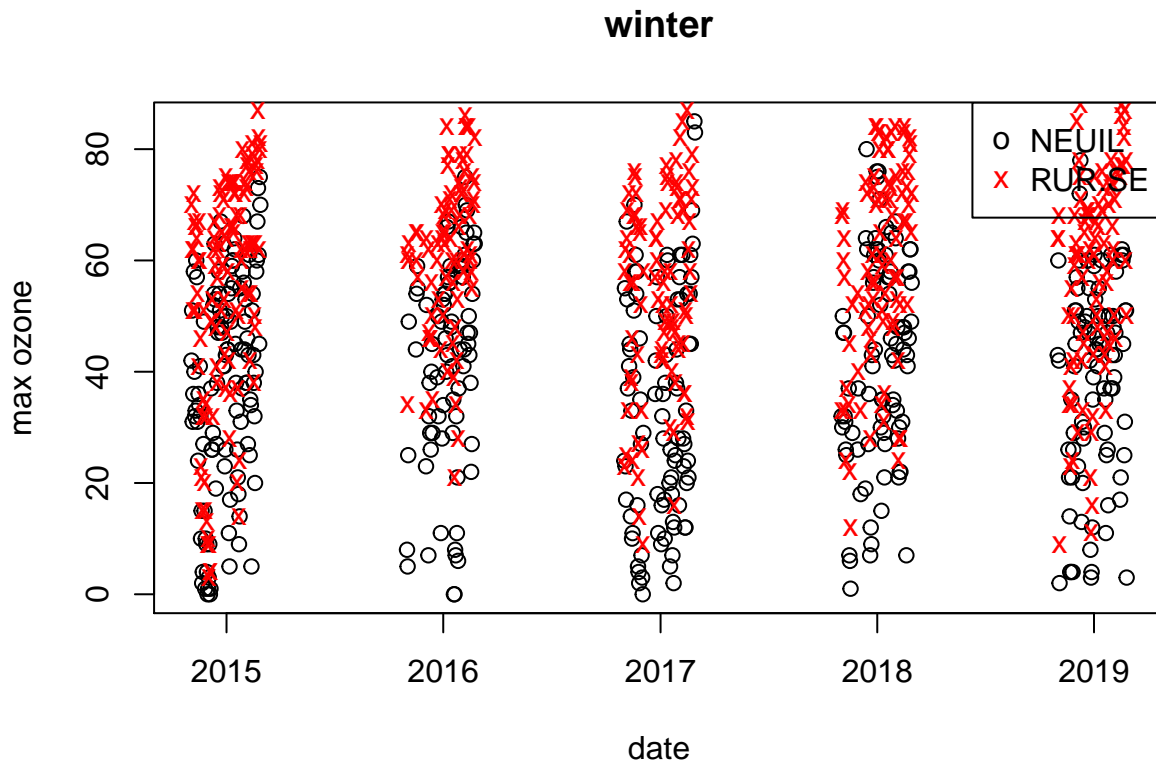


La variabilité de ces données différenciées est considérablement moindre que la variabilité des mesures effectuées sur les sites séparés. Cela indique que des facteurs communs affectant les deux sites influencent la variation des valeurs d'ozone. La plupart des différences sont négatives, ce qui suggère qu'en général les mesures rurales sont plus grandes que les mesures urbaines, ce qui coïncide avec les attentes scientifiques.

3. Répétez l'analyse pour les données d'hiver sur l'ozone. Résumez vos découvertes.

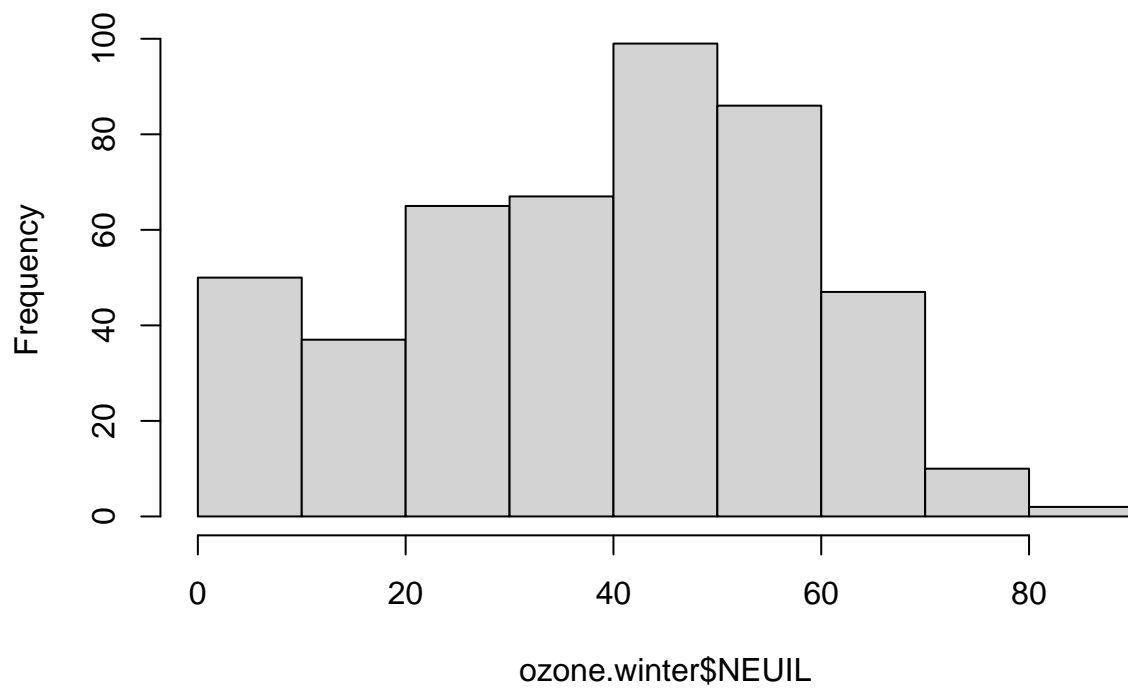
```
ozone.winter = read.csv("winter_ozone.csv")
str(ozone.winter)
ozone.winter$date = as.Date(ozone.winter$date2)

plot(ozone.winter$date, ozone.winter$NEUIL, xlab="date", ylab="max ozone", main="winter")
points(ozone.winter$date, ozone.winter$RUR.SE, col="red", pch = "x")
legend("topright", legend = c("NEUIL", "RUR.SE"), col=c("black","red"), pch=c("o","x"))
```



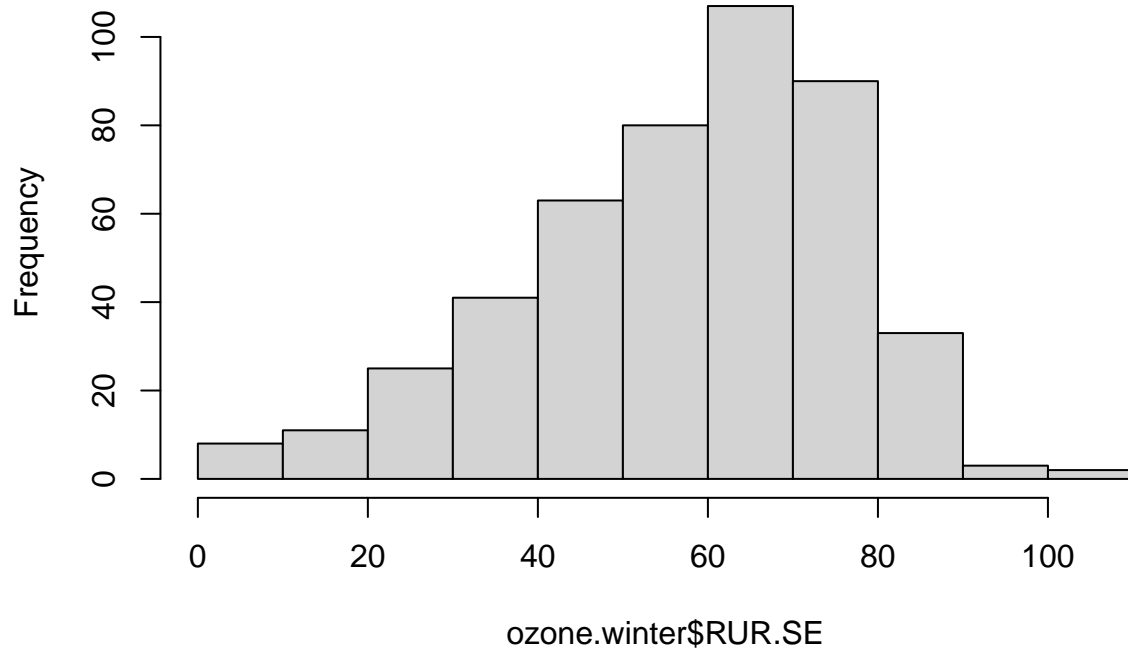
```
## histogram
hist(ozone.winter$NEUIL)
```

Histogram of ozone.winter\$NEUIL



```
hist(ozone.winter$RUR.SE)
```

Histogram of ozone.winter\$RUR.SE



Empirical distribution function

La fonction de distribution empirique est donnée par

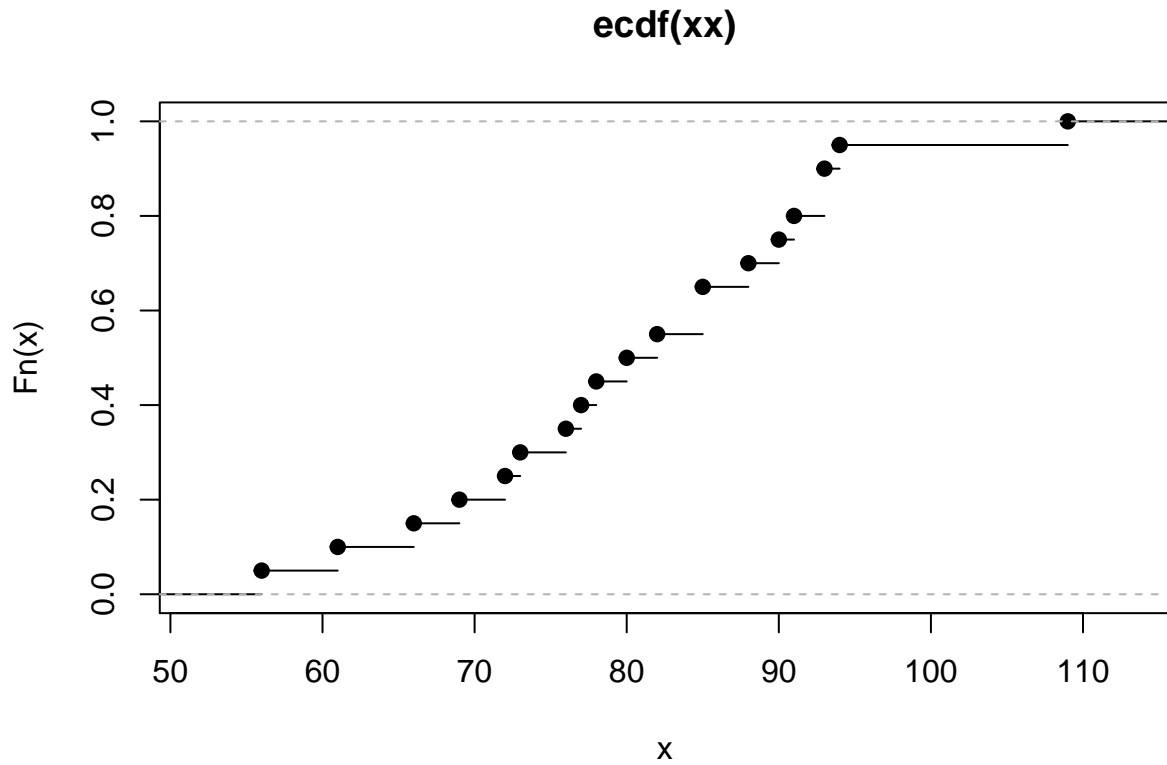
$$F_n(x_{(i)}) = \frac{i}{n}$$

aux points de données ordonnés $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$. Pour les valeurs de x entre les points de données, nous avons

$$F_n(x) = \frac{i}{n}, \text{ où } x_{(i)} \leq x < x_{(i+1)}$$

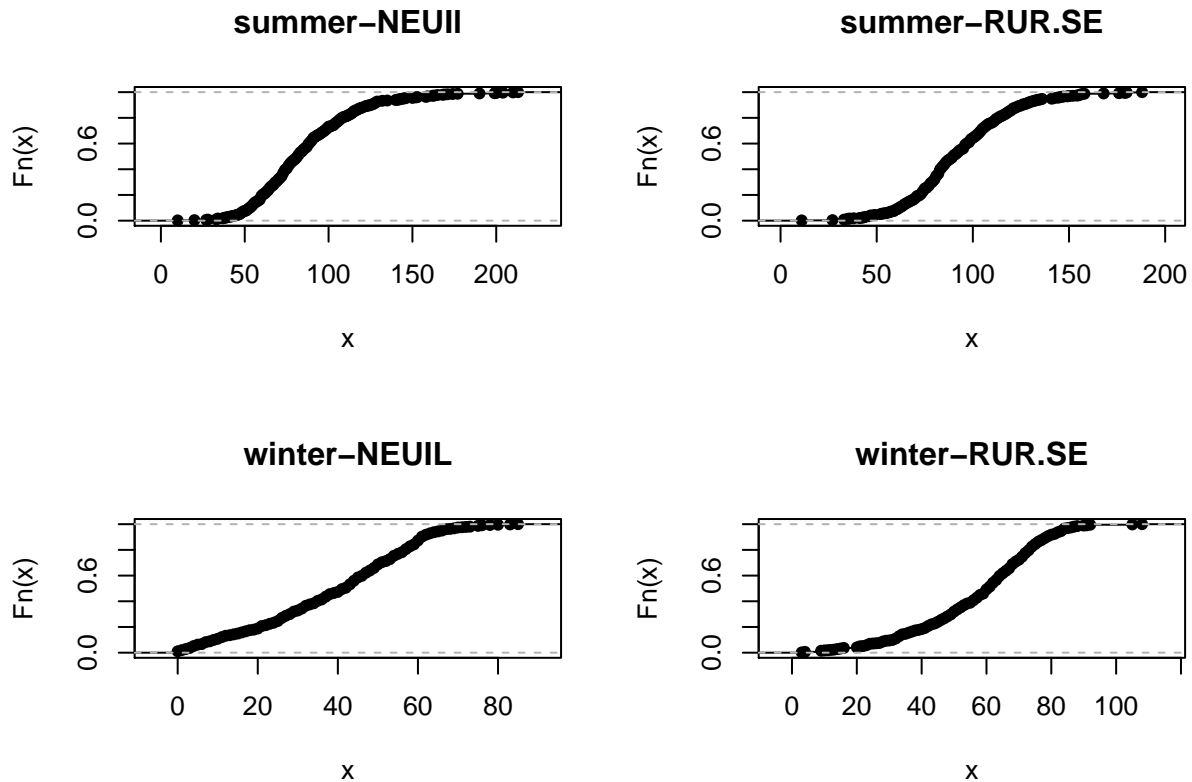
Par exemple, les 20 premières observations des mesures d'ozone en milieu urbain sont

```
xx = ozone.summer$NEUIL[1:20]
sort(xx)
plot(ecdf(xx))
```



4. Faites le c.d.f empirique des jeux de données complets pour chaque site. Expliquez comment utiliser les graphiques pour estimer la médiane.

```
par(mfrow=c(2,2))
plot(ecdf(ozone.summer$NEUIL), main="summer-NEUIL")
plot(ecdf(ozone.summer$RUR.SE), main="summer-RUR.SE")
plot(ecdf(ozone.winter$NEUIL), main="winter-NEUIL")
plot(ecdf(ozone.winter$RUR.SE), main="winter-RUR.SE")
```



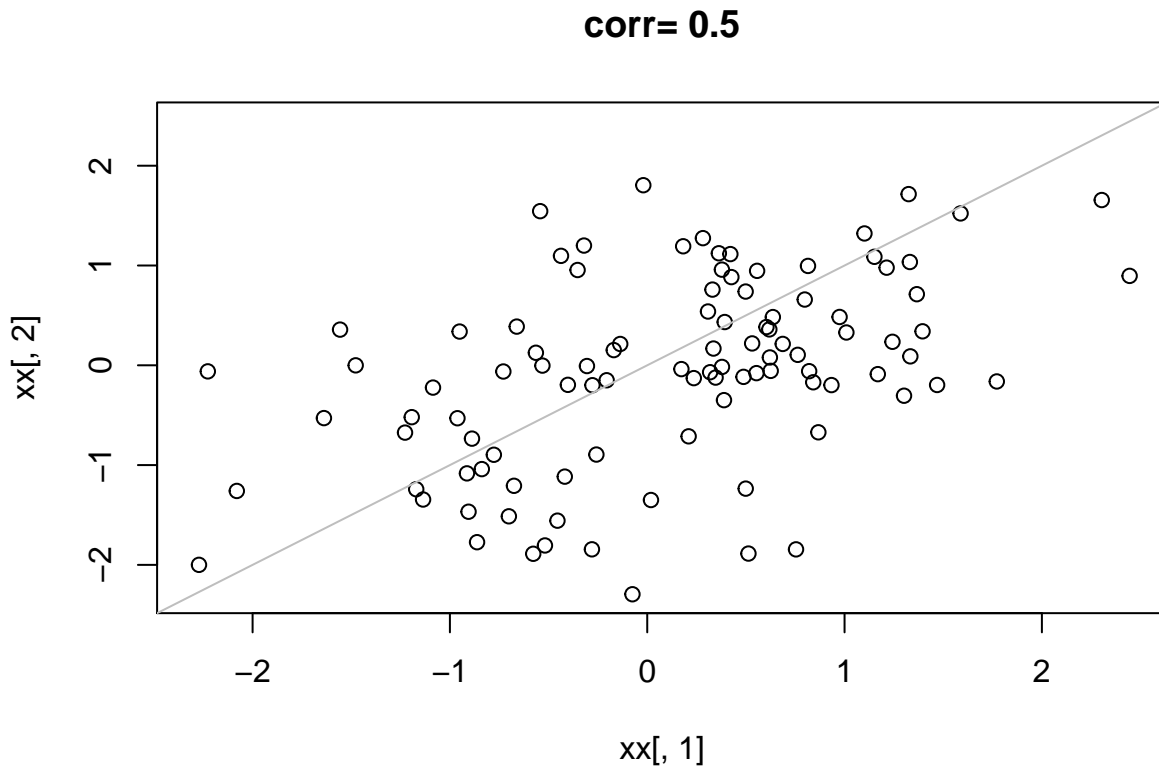
Sample covariance and sample correlation

La **corrélation** entre deux variables aléatoires X et Y est

$$\rho = \text{Corr}(X, Y) = \frac{E[(X - E[X])(Y - E[Y])]}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

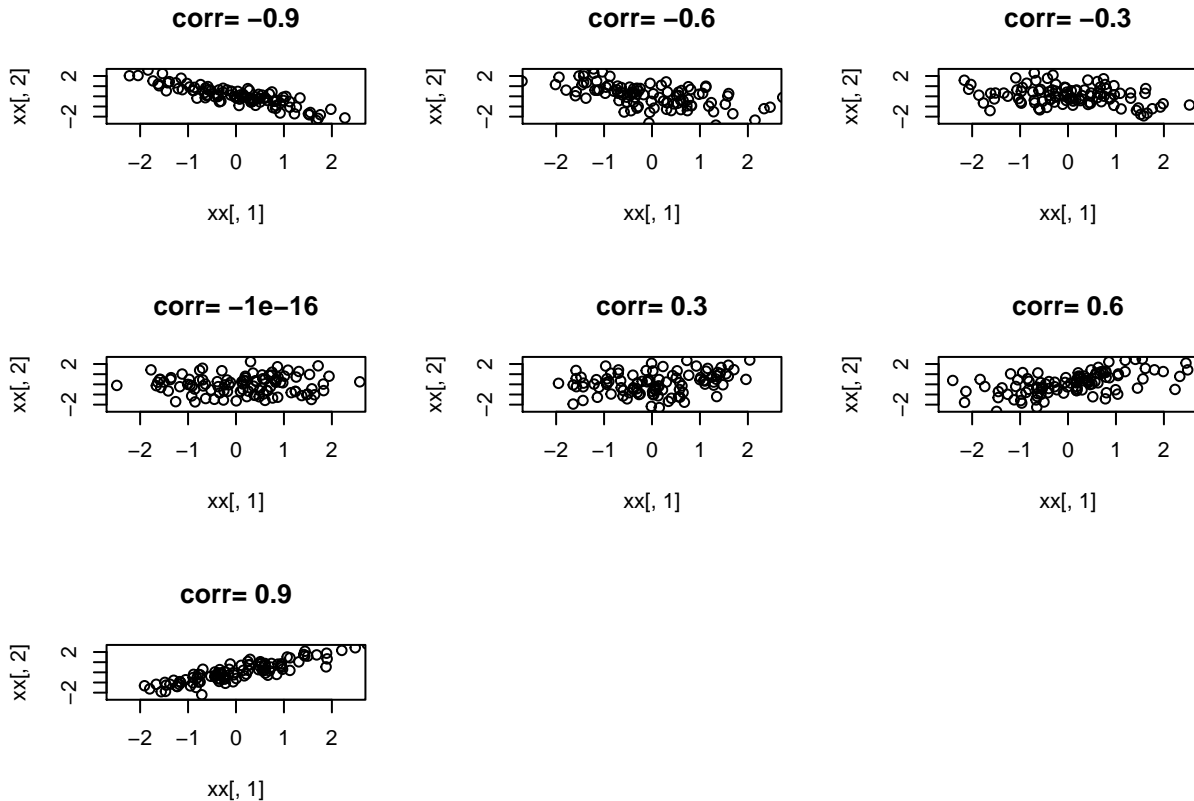
Le code suivant simule à partir d'une distribution normale bivariée avec une moyenne nulle, une variance unitaire et une corrélation ρ , $-1 \leq \rho \leq 1$:

```
#install.packages("MASS") # if not installed already
library(MASS)
?mvrnorm # help
rho = 0.5
xx = mvrnorm(n=100, mu=c(0,0), Sigma = matrix(c(1,rho, rho, 1), ncol=2))
xx[1:10,]
xxlim = c(min(xx), max(xx))
plot(xx[,1], xx[,2], main = paste("corr=", rho), xlim = xxlim, ylim=xxlim)
abline(a=0, b=1, col='gray') # add line y=x
```

5. Expérimentez avec une plage de valeurs de ρ et comparez les nuages de points. Décrivez l'effet de ρ .

```
vrho = seq(-0.9, 0.9, by=0.3)
k = length(vrho)
xxlim = c(-2.5, 2.5)
par(mfrow=c(3,3))
for (ik in 1:k){
  rho = vrho[ik]
  xx = mvrnorm(n=100, mu=c(0,0), Sigma = matrix(c(1,rho, rho, 1), ncol=2))
  plot(xx[,1], xx[,2], main = paste("corr=", signif(rho,1)), xlim = xxlim, ylim=xxlim)
}
```



Le **coefficient de corrélation d'échantillon** de n paires d'observations $(x_1, y_1), \dots, (x_n, y_n)$ est noté r et est donné par

$$r = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

où \bar{x} et \bar{y} sont les moyennes de l'échantillon et s_x et s_y sont les écarts types de l'échantillon.

La corrélation mesure une dépendance linéaire d'une manière indépendante de l'échelle. La **covariance** entre deux variables aléatoires X et Y est définie de manière similaire:

La **covariance** entre les variables aléatoires X et Y est

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

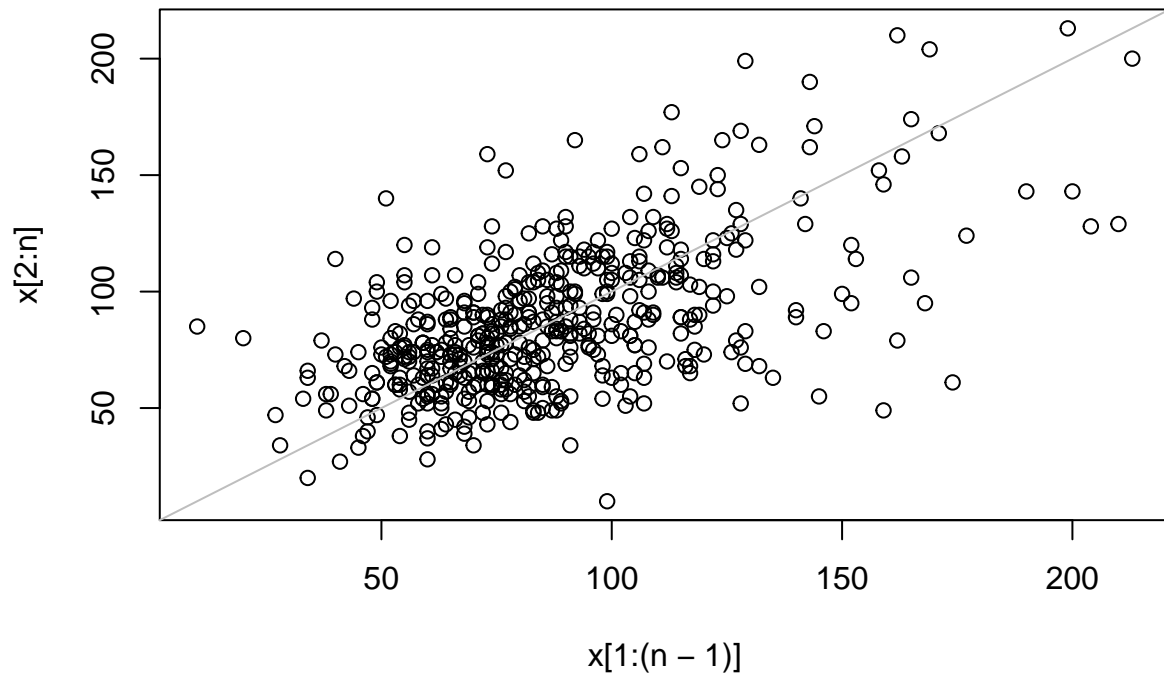
La **covariance d'échantillon** de n observations appariées $(x_1, y_1), \dots, (x_n, y_n)$ est donnée par

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = r s_x s_y$$

6. La production d'ozone peut persister pendant plusieurs jours. Le code suivant crée le nuage de points de x_t contre x_{t-1} pour toutes les valeurs de t pour l'ozone urbain en été.

```
x = ozone.summer$NEUIL
n = length(x)
plot(x[1:(n-1)], x[2:n], main="summer-NEUIL")
abline(a=0, b=1, col="gray")
```

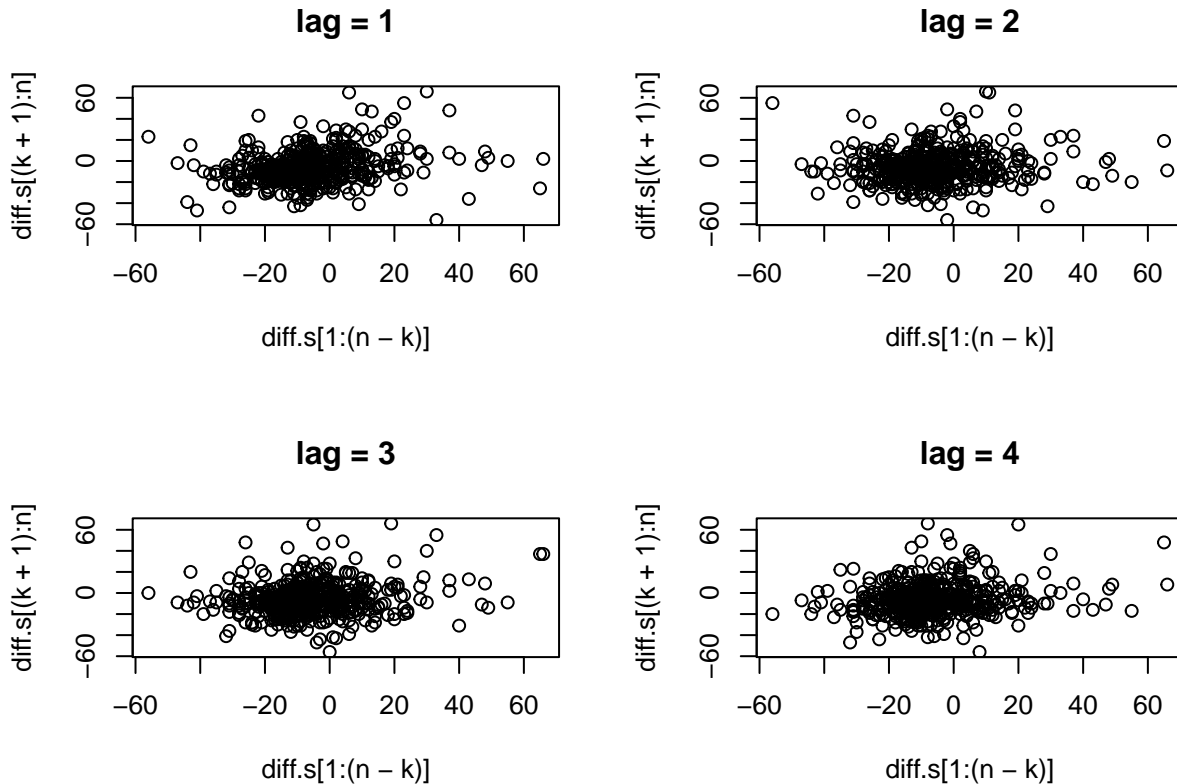
summer-NEUIL



Qu'observez-vous? Estimez le coefficient de corrélation.

7. Explorez la dépendance pour la série différenciée $d_i, 1 \leq i \leq n$. Qu'observez-vous?

```
diff.s = ozone.summer$NEUIL - ozone.summer$RUR.SE
n = length(diff.s)
par(mfrow=c(2,2))
for (k in 1:4){
  plot(diff.s[1:(n-k)], diff.s[(k+1):n], main = paste("lag =", k))
}
```



no strong dependence

Moyenne et phénomène de concentration

Nous allons montrer que la moyenne d'une variable aléatoire est un résumé déterministe d'une v.a., dont la qualité est contrôlée par la variance.

8. Rappeler l'inégalité de Bienaymé-Chebychef dans les cas Gaussien et Poisson.
9. Estimer par Monte Carlo les probabilités de déviation d'une variable aléatoire de sa moyenne.
 - (a) Exprimer $P(|X - \mu| \geq \delta)$ comme l'espérance d'une certaine variable aléatoire Z .
 - (b) Simuler un échantillon de taille N Z_1, Z_2, \dots, Z_N de même loi que Z (dans le cas Gaussien, Pareto et Poisson) - on prendra N grand. Déterminer une estimation de $P(|X - \mu| \geq \delta)$. Pouvez vous déterminer la précision de cette estimation?
 - (c) Comparer avec les bornes obtenues par Bienaymé Chebychev pour plusieurs δ . Faites varier σ .
 - (d) Dans le cas Gaussien et Poisson, comparer les estimations Monte-Carlo de $P(X - \mu \geq \delta)$ avec les bornes données par les inégalités Chernoff pour plusieurs δ et σ (cf. cours).
10. Simuler un échantillon de taille $n = 20$ pour les lois de Gauss et de Poisson (*choisir σ, λ approprié*)
 - (a) Calculer les bornes de Chernoff dans le cas échantillon pour \bar{X}_n . Faites varier $n = 20, 100, 1000$.
 - (b) En déduire un estimateur de μ et λ respectivement.
11. Simuler un échantillon de taille $n = 20$ d'une loi de Cauchy $\mathcal{C}(\theta)$ de densité? $f(x, \theta) = \frac{1}{\pi} \frac{1}{1+(x-\theta)^2}$.
 - (a) Calculer la moyenne empirique \bar{X}_n . Faites varier la taille de l'échantillon $n = 20, 100, 1000$ et 10000. Qu'en déduire ?

- (b) Expliquer ce comportement. On se rappellera notamment que la fonction caractéristique s'écrit $\phi_{\theta}(t) = \exp(i\theta t - |t|)$.
- (c) Quelle est la médiane d'une loi de Cauchy $\mathcal{C}(\theta)$? En déduire un estimateur de θ pour $n = 20, 100, 1000$.