

TP: Le modèle linéaire

Eric Marcon

30 January 2024

Régression linéaire simple

Données du projet de dendrométrie 2020, Mont Ventoux.

```
read_csv2("data/Inv_GEEFT_Ventoux_09-2020.csv") |>
  rename(
    espece = Espèce,
    diametre = `Diamètre (cm)`,
    hauteur = `Hauteur réelle (m)`
  ) |>
  mutate(
    espece = case_match(
      espece,
      "P" ~ "Pin",
      "C" ~ "Cèdre"
    )
  ) ->
  ventoux
```

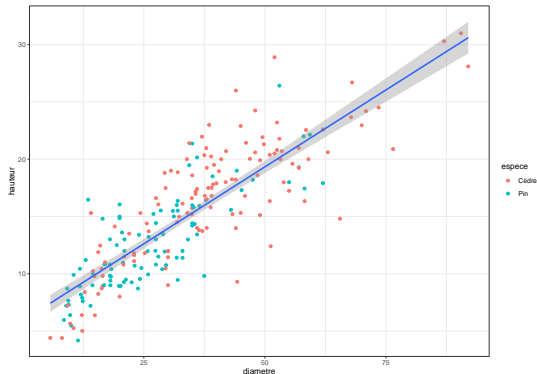
Graphique hauteur ~ diamètre

TP: Le
modèle
linéaire

Eric Marcon

Régression
linéaire simple

```
ventoux |>
  ggplot(aes(x = diametre, y = hauteur)) +
  geom_point(aes(col = espece)) +
  geom_smooth(method = "lm")
```



Modèle linéaire simple :

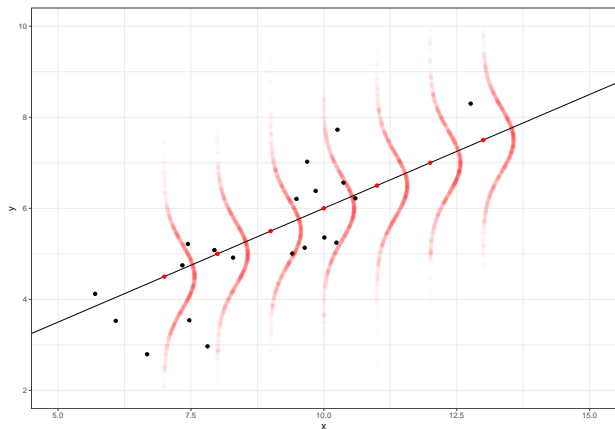
$$Y = \beta_0 + \beta_1 X + E$$

Y et X sont des vecteurs : $Y = \{y_i\}$ est l'ensemble des observations. Par abus d'écriture, Y est aussi la variable aléatoire dont les y_i sont des réalisations.

Vocabulaire : variable expliquée, exogène, coefficients, constante (intercept)...

$E = \{\epsilon_i\}$ est l'erreur du modèle. $E \sim \mathcal{N}(0, \sigma^2)$

Le modèle prédit une densité de probabilité des valeurs de Y pour toute valeur de X distribuée normalement autour de la droite de régression.



- Indépendance des erreurs : $\text{Cov}(\epsilon_i, \epsilon_j) = 0$. Assurée par le design expérimental.
- Exogénéité : X n'est pas corrélé à E .
- Homoscédasticité : la variance de l'erreur est constante sur l'étendue de X .
- Normalité des termes d'erreur : $E \sim \mathcal{N}(0, \sigma^2)$.

Générer les données du modèle.

Coefficients :

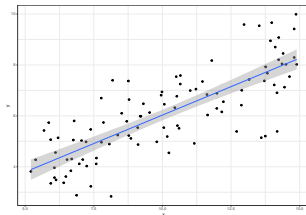
```
beta0 <- 1  
beta1 <- 0.5  
sigma <- 1
```

Tirage :

```
n <- 100  
x <- runif(n, min = 5, max = 15)  
# Jeu de points  
mod_l1 <- tibble(x, y = rnorm(n, mean = beta0 + beta1*x, sd =sigma))
```


Commencer par une figure.

```
mod_l1 |>  
  ggplot(aes(x = x, y = y)) + geom_point() + geom_smooth(method = lm)
```

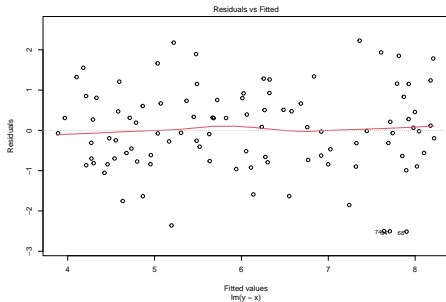


La fonction `lm()` du package *stats* estime le modèle et permet de tester les hypothèses.

```
mod_l1_lm <- lm(y ~ x, data = mod_l1)
```

Graphique $E \sim Y^*$

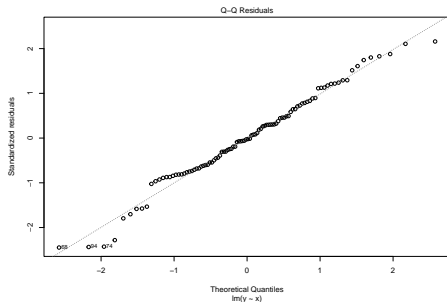
```
plot(mod_l1_lm, which = 1)
```



Les erreurs doivent être centrée sur 0 et uniformément réparties.

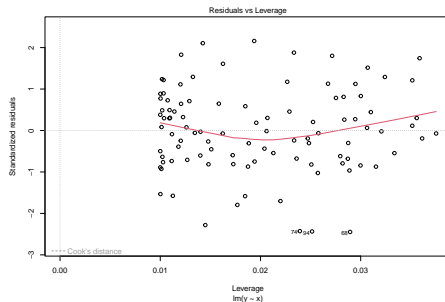
Graphique quantile - quantile (`qqplot`)

```
plot(mod_l1_lm, which = 2)
```



La non-normalité des résidus implique la non-normalité des estimateurs des coefficients.

```
plot(mod_l1_lm, which = 5)
```



Les points avec fort effet de levier forte erreur (\rightarrow grande distance de Cook) posent problème.

Affaire d'expérience.

- Éliminer les points (réellement) aberrants ;
- Transformer Y si :
 - la relation n'est pas linéaire (ex.: quadratique) ;
 - l'erreur augmente avec Y^* (\rightarrow racine carrée ou logarithme).
- Revoir les hypothèses à l'origine du modèle, le design expérimental...

Interprétation des résultats : summary

TP: Le
modèle
linéaire

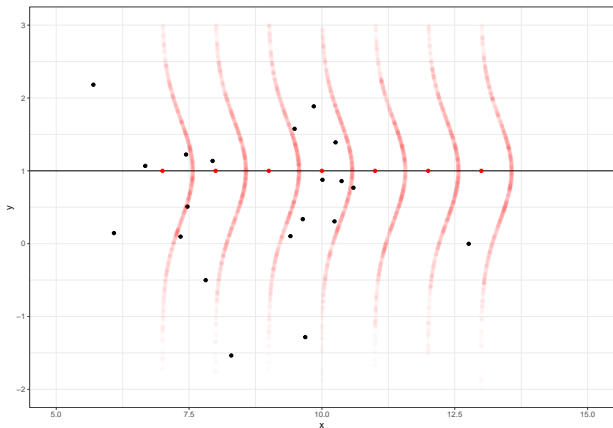
Eric Marcon

Régression
linéaire simple

```
##
## Call:
## lm(formula = y ~ x, data = mod_l1)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.51229 -0.69782 -0.02673  0.68502  2.22763
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.55838     0.37138   4.196 5.97e-05
## x            0.44729     0.03525  12.689 < 2e-16
##
## (Intercept) ***
## x            ***
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.042 on 98 degrees of freedom
## Multiple R-squared:  0.6216, Adjusted R-squared:  0.6178
## F-statistic: 161 on 1 and 98 DF, p-value: < 2.2e-16
```

La statistique F décrit la probabilité que le modèle n'explique rien.

Modèle nul: $Y = \bar{Y} = \beta_0$



R^2 mesure la proportion de la variance de Y expliquée par le modèle :

$$R^2 = \frac{\text{Var}(Y^*)}{\text{Var}(Y)} = 1 - \frac{\sigma}{\text{Var}(Y)}$$

→ Que devient R^2 en doublant σ ? Estimer rapidement puis resimuler le modèle pour vérifier.

R^2 ajusté pénalise le R^2 par le nombre de paramètres du modèle.

Les degrés de liberté sont le nombre d'observations moins le nombre de paramètres moins 1.

Les coefficients sont estimés par la méthode des moindres carrés : minimisation des écarts

$$\sum (y_i - y_i^*)^2$$

.

Résultat identique à la maximisation de la vraisemblance

$$\prod f(\epsilon_i)$$

où $f(\cdot)$ est la densité de $\mathcal{N}(0, \sigma^2)$.

L'estimateur de chaque coefficient est sa valeur la plus probable.

L'estimateur est distribué normalement (quand E est normal) :

$$\hat{\beta}_1 \sim \mathcal{N}(0.447, 0.0352^2)$$

Un test de Student donne la probabilité de se tromper en affirmant que l'estimateur n'est pas nul.

Un bon modèle a un grand R^2 et des petites p-values.

- R^2 diminue avec la variance de l'erreur ;
- L'écart-type des estimateur diminue comme \sqrt{n} .

Mais les deux dépendent du design expérimental.

Quadrupler l'effort d'échantillonnage divise par deux l'intervalle de confiance

```
mod_l1x4 <- tibble(  
  x = rnorm(n * 4, mean = 10, sd = 2), # x est calculé avant y  
  y = rnorm(n * 4, mean = beta0 + beta1 * x, sd = sigma) # y utilise x  
)  
mod_l1x4_lm <- lm(y ~ x, data = mod_l1x4)  
summary(mod_l1x4_lm)$coefficients
```

```
##              Estimate Std. Error   t value  
## (Intercept) 1.8435413 0.25344512  7.273927  
## x           0.4226806 0.02443802 17.296027  
##              Pr(>|t|)  
## (Intercept) 1.870813e-12  
## x           2.180247e-50
```

Choix économique.

Retirer les valeurs intermédiaires de X augmente le R^2 (*design factoriel*) alors que σ ne change pas.

```
mod_l1x4 |>
  filter(x < 6 | x > 14) %>% # pas /> pour "data = ."
  lm(y ~ x, data = .) |>
  summary() |>
  pluck("r.squared")
```

```
## [1] 0.8153628
```

contre 0.4291063 avec toutes les données.

Le R^2 d'un modèle avec des données individuelles est plus faible qu'avec des données agrégées.

- Estimer le modèle hauteur \sim diamètre des données Ventoux.
- Regrouper les données par espèce. → Estimer le modèle à nouveau.

Considérer R^2 et p-values en fonction du modèle :

- beaucoup de données individuelles \rightarrow faible R^2 mais petites p-values pour montrer l'influence d'un facteur ;
- possibilité d'un très grand R^2 sans aucun coefficient significatif si peu de points ;
- un grand R^2 et des petites p-values permettent de faire des prédictions.

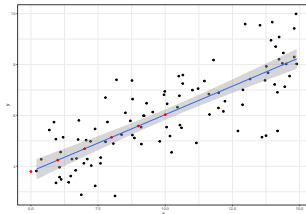
`predict()` permet d'extrapoler le modèle.

```
mod_l1_lm |> predict(newdata = data.frame(x = 5:10))
```

```
##           1           2           3           4           5
## 3.794805 4.242090 4.689375 5.136661 5.583946
##           6
## 6.031231
```


Ajout des points sur la figure :

```
# Estimation du modèle
mod_l1 |>
  ggplot(aes(x, y)) + geom_point() + geom_smooth(method = lm) ->
  mod_l1_ggplot
# Choix des x pour lesquels y est à prédire
mod_l1_predict <- data.frame(x = 5:10)
# Ajout des prédictions
mod_l1_predict$y <- predict(mod_l1_lm, newdata = mod_l1_predict)
# Ajout des points à la figure précédente
mod_l1_ggplot +
  geom_point(data = mod_l1_predict, aes(x = x, y = y), col = "red")
```



Intervalles de confiance et de prédiction

La zone grisée de `geom_smooth` est l'intervalle de confiance de l'espérance de $Y|X$, c'est-à-dire de la moyenne des prédictions.

Il est bien plus étroit que l'intervalle de prédiction, qui correspond à 95% des prédictions :

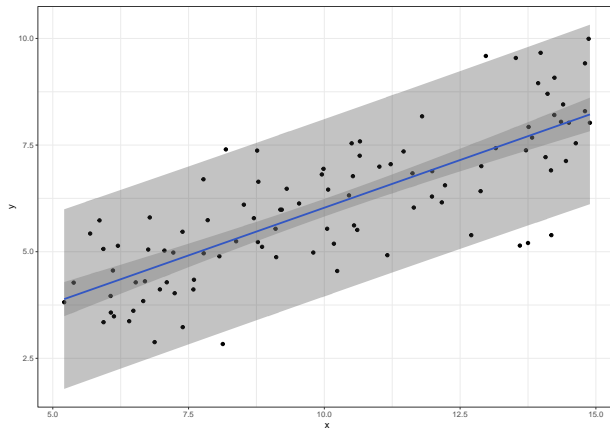
```
mod_l1_predict <- data.frame(
  x = seq(from = min(mod_l1$x), to = max(mod_l1$x), length.out = 50)
)
mod_l1_predict <- cbind(
  mod_l1_predict,
  predict(
    mod_l1_lm,
    newdata = mod_l1_predict,
    interval = "prediction"
  )
)
mod_l1_ggplot +
  geom_ribbon(
    data = mod_l1_predict,
    aes(y = fit, ymin = lwr, ymax = upr),
    alpha = 0.3
  ) -> mod_l1_ggplot_predict
```

Intervalles de confiance et de prédiction

TP: Le
modèle
linéaire

Eric Marcon

Régression
linéaire simple



TP: Le
modèle
linéaire

Eric Marcon

Régression
linéaire simple