

6 Exploration statistique de données

6.1 Pré-requis

La première étape de toute analyse de données est l'exploration. Avant de se lancer dans des tests statistiques et des procédures complexes, et à supposer que les données dont vous disposez sont déjà dans un format approprié, il est toujours très utile :

1. d'explorer visuellement les données dont on dispose en faisant des graphiques nombreux et variés, afin de comprendre, notamment quelle est la distribution des variables numériques, quelles sont les catégories les plus représentées pour les variables qualitatives (facteurs), et quelles sont les relations les plus marquantes entre variables numériques et/ou catégorielles, etc. Vous avez appris au [chapitre 2](#) comment produire toutes sortes de graphiques avec ggplot2. À présent, il faut vous poser la question du choix des graphiques à produire du point de vue de l'exploration statistique de données inconnues.
2. d'explorer les données en calculant des indices statistiques descriptives. Ces indices relèvent en général de 2 catégories : les indices de position (e.g. moyennes, médianes, quartiles...) et les indices de dispersion (e.g. variance, écart-type, intervalle inter-quartile). Nous avons déjà vu comment utiliser la fonction `summary()` pour argument calculer des moyennes ou des effectifs pour plusieurs sous-groupes de nos jeux de données. Dans ce chapitre, nous irons plus loin, et nous découvrirons d'une part comment calculer d'autres indices statistiques plus utiles, et comment utiliser d'autres fonctions encore.

Nous verrons ensuite comment installer les indices d'incertitude et de dispersion, il ne faudra pas confondre indices de dispersion et d'incertitude. Et enfin, avant de passer aux tests statistiques, nous verrons comment visualiser dispersion et incertitude. Chapitre 8

Afin d'explorer ces questions, nous aurons besoin des packages suivants :

```
library(dplyr)
library(arsy)
library(palmerpenguins)
library(nycflights13)
```

Comme vous le savez maintenant, les packages du tidyverse (Wickham 2023) permettent de manipuler facilement des tableaux de données et de réaliser des graphiques. Le tidyverse permet d'accéder, entre autres, aux packages readr (Wickham, Hester, et Bryan 2024), pour importer facilement des fichiers texte, lubridate (Wickham, Vaughan, et Glynn 2024) et et al. (2023) pour manipuler des tableaux de données en colonne, ggplot2 (Wickham et al. 2024) pour produire des graphiques, skimr (Wickham et al. 2022) permet de calculer des résumés de données très informatifs, palmerpenguins (Harris, Hill, et Gorman 2022) et nycflights13 (Wickham 2021) fournissent des jeux de données qui seront faciles à manipuler pour illustrer le chapitre (et les suivants).

Important

Si vous avez installé dplyr avant le printemps 2023, vous devez réinstaller le package (pour dplyr) à effet et mis à jour courant 2023, et nous aurons besoin de sa version v1.1.0 ou d'une version plus récente pour utiliser certaines fonctions. Chargez-le ensuite en mémoire avec library(dplyr)

Attention

Pensez à installer tous les packages listés ci-dessus
de les charger en mémoire si vous ne l'avez pas déjà fait.
Si vous ne savez plus comment faire, consultez d'urgence
la [Section](#)

Pour travailler dans de bonnes conditions, et puisqu'abordons maintenant les statistiques à proprement parler, je vous conseille de créer un nouveau script dans le même dossier. Et encore, si vous ne savez plus comment faire consulter la Section

6.2 Créer des résumés avec la fonction `reframe()`

Dans la `Section 5.8`, nous avons vu comment utiliser la fonction `summary()` éventuellement pour argument calculer des statistiques descriptives variées. N'hésitez pas à relire cette section si vous n'êtes pas sûr d'avoir tout ou tout retenu. Les calculs que nous pouvons faire grâce à la fonction `summary()` quant des fonctions statistiques qui ne renvoient qu'une valeur à la fois lorsqu'on leur passe une série de valeurs. Par exemple, si on dispose d'un vecteur numérique (les entiers compris entre 1 et 100 pour l'exemple) :

1: 1 0 0

[1]	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
[19]	19	20	21	22	23	24	25	26	27	28	29	30	31	32		
[37]	37	38	39	40	41	42	43	44	45	46	47	48	49	50		
[55]	55	56	57	58	59	60	61	62	63	64	65	66	67	68		
[73]	73	74	75	76	77	78	79	80	81	82	83	84	85	86		
[91]	91	92	93	94	95	96	97	98	99	100						

La fonction `mean` renvoie qu'une valeur, la moyenne des 100 nombres contenus dans le vecteur :

me a(1,1 0)0

```
[ 1] 50.5
```

De même pour les `sd()` (donnée à nos) toutes les autres fonctions listées à la fin de la Section

```
sd(1:10)
```

```
[ 1] 29.01149
```

```
medi(1:10)
```

```
[ 1] 50.5
```

Il existe toutefois des fonctions qui renvoient plus leur à la fois. Par exemple, la fonction `quant()` que nous avons utilisée dans un autre exemple, dont le contexte à la Section

1. la valeur minimale contenue dans le vecteur (ou quantile 0%) : c'est la valeur la plus faible contenue dans la série de données
2. le premier quartile du vecteur (Q1 ou quantile 25%) est la valeur coupant l'échantillon en deux : 25% des observations du vecteur y sont inférieures et 75% y sont supérieures
3. la médiane du vecteur (Q2 ou quantile 50%) est la valeur coupant l'échantillon en deux : 50% des observations du vecteur sont inférieures à cette valeur et 50% y sont supérieures
4. le troisième quartile du vecteur (Q3 ou quantile 75%) est la valeur coupant l'échantillon en deux : 75% des observations du vecteur y sont inférieures et 25% y sont supérieures
5. la valeur maximale contenue dans le vecteur (ou quantile 100%) : c'est la valeur la plus élevée contenue dans la série de données.

Par exemple, toujours avec le vecteur des entiers compris entre 1 et 100 :

```
quant(1:100)
```

0%	25%	50%	75%	100%
1.00	25.75	50.50	75.25	100.00

L'objet obtenu est un vecteur dont chaque élément porte un nom. Pour transformer cet objet en tibble, on utilise `enframe()`

```
enframe(name = 1:10, value = 0)
```

```
# A tibble: 5 x 2
  name value
<chr> <dbl>
1 0%      1
2 25%    25.8
3 50%    50.5
4 75%    75.2
5 100%   100
```

Il peut être très utile de calculer ces différentes valeurs pour plusieurs variables à la fois, ou pour plusieurs sous-ensembles d'un jeu de données. Le problème est que nous ne pouvons pas utiliser `summarise()` la fonction `summarise()` ne renvoie pas qu'une unique valeur. Par exemple, pour calculer les quantiles des longueurs de becs pour chaque espèce de chats, on pourrait être tenté de taper ceci :

```
penguins %>% summarise(
  quantiles = quantile(beak_length, probs = 0:4/5, na.rm = TRUE)
  ) by = species)
```

Warning: Returning more (or less) than 1 row per `summarise` group. Please use `reframe()` instead.

i Please use `reframe()` instead.
i When switching from `summarise()` to `reframe()`, always returns an ungrouped data frame and adjust accordingly.

```
# A tibble: 15 x 2
  species indices
<fct>    <dbl>
1 Adelie 32.1
2 Adelie 36.8
```

```

3 Adelie 38.8
4 Adelie 40.8
5 Adelie 46
6 Gentoo 40.9
7 Gentoo 45.3
8 Gentoo 47.3
9 Gentoo 49.6
10 Gentoo 59.6
11 Chinstrap 40.9
12 Chinstrap 46.3
13 Chinstrap 49.6
14 Chinstrap 51.1
15 Chinstrap 58

```

C'est dans ces situations que nous aurons besoin d'une fonction
Elle joue le même rôle que les fonctions `summary()` et `table()` dans les situa-
tion où les fonctions statistiques renvoient plus d'un résultat :

```

penguins
  summarise(
    ref = ref()
  )

```

```

# A tibble: 15 x 2
  species indices
  <fct>    <dbl>
1 Adelie 32.1
2 Adelie 36.8
3 Adelie 38.8
4 Adelie 40.8
5 Adelie 46
6 Gentoo 40.9
7 Gentoo 45.3
8 Gentoo 47.3
9 Gentoo 49.6
10 Gentoo 59.6
11 Chinstrap 40.9
12 Chinstrap 46.3
13 Chinstrap 49.6
14 Chinstrap 51.1
15 Chinstrap 58

```

Au contraire, `summary()` ne renvoie pas de message d'avertissement dans cette situation. Dans l'exemple, on ne sait malheureusement pas à quoi correspondent les chiffres renvoyés puisque l'information sur les quartiles a disparu (quelles valeurs correspondent aux médianes ou aux premiers quartiles par exemple). Pour y remédier, on doit transformer le vecteur renvoyé en `quantile(tibble)`. Nous avons déjà vu comment le faire grâce à `as.data.frame()` ou `as_tibble()`, puisque la fonction va maintenant renvoyer un tableau, on n'a pas besoin de lui fournir de nom de colonnes (je retire donc `indices` mon code) :

```
penguins
  summarise(quantile = quantile(black_lip_length, probs = 0:4,
                                . by = species))
```

```
# A tibble: 15 x 3
  species    name value
  <fct>      <chr> <dbl>
1 Adelie    0%      32.1
2 Adelie    25%     36.8
3 Adelie    50%     38.8
4 Adelie    75%     40.8
5 Adelie   100%     46
6 Gentoo    0%      40.9
7 Gentoo    25%     45.3
8 Gentoo    50%     47.3
9 Gentoo    75%     49.6
10 Gentoo   100%     59.6
11 Chinstrap 0%      40.9
12 Chinstrap 25%     46.3
13 Chinstrap 50%     49.6
14 Chinstrap 75%     51.1
15 Chinstrap 100%     58
```

Enfin, comme décrit à la section 4.3.2, il est possible de modifier la forme de `tbl_summary()` pour le rendre plus facilement et éventuellement l'intégrer dans un rapport ou compte-rendu :

```

penguins
  summarise(
    ref = ref(quantile(bill_length_mm),
              . by = species)
    pivot_wider(names_from = species,
                values_from = ref)
  )

```

```

# A tibble: 5 x 4
  name      Adelie Gentoo Chinstrap
  <chr>    <dbl>   <dbl>   <dbl>
1 0%      32.1    40.9    40.9
2 25%      36.8    45.6    45.6
3 50%      38.8    47.9    47.9
4 75%      40.8    49.1    49.1
5 100%     46     59.5    59.5

```

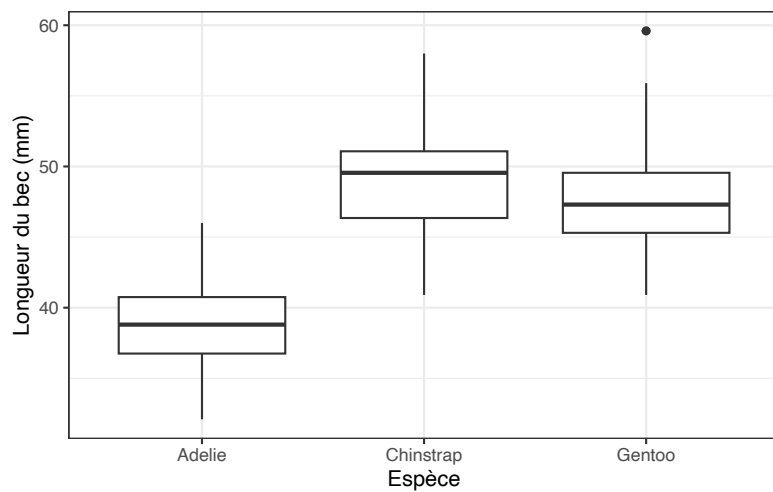
Ces statistiques nous permettent de constater que les manchots de l'espèce Adélie semblent avoir des becs plus longs que les 2 autres espèces (les 5 quantiles le confirment). Les manchots Gentoo et Chinstrap ont en revanche des becs de longueur à peu près similaires, bien que ceux des Chinstrap soient peut-être très légèrement plus longs (Q1, médiane, Q3 supérieurs à ceux des Gentoo). On peut vérifier tout cela graphiquement avec des boîtes à moustaches, puisque les valeurs de quantiles sont justement celles qui sont utilisées pour tracer les boîtes à moustaches :

```

penguins
  ggplot(aes(x = species, y = bill_length_mm))
  geom_boxplot()
  labs(x = "Espèce", y = "Longueur du bec (mm)")
  theme_bw()

```

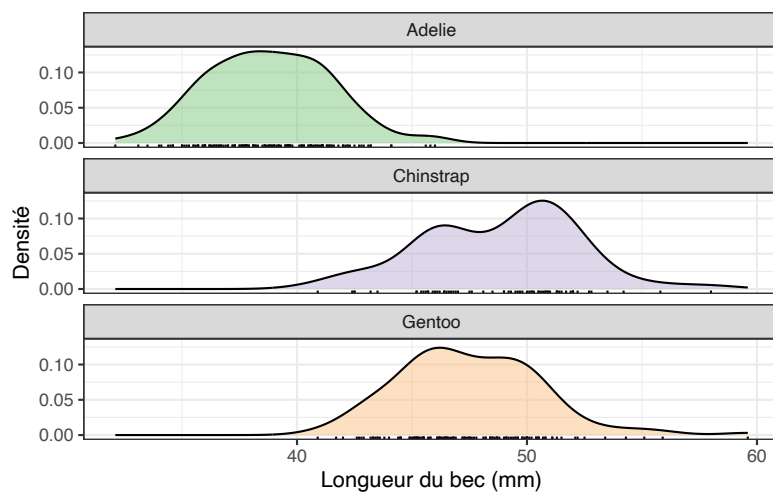
Warning: Removed 2 rows containing non-finite outliers (stat_boxplot()).



Ou avec un graphique de densité :

```
library(ggplot2)
penguins %>%
  ggplot(aes(x = bill_length_mm, y = density, color = species)) +
  geom_density(aes(color = species)) +
  geom_point() +
  labs(x = "Longueur du bec (mm)", y = "Densité") +
  facet_wrap(~species, color = "species") +
  scale_fill(palette = "Accent") +
  theme_bw()
```

Warning: Removed 2 rows containing non-finite outliers (``stat_density()``).



À ce stade, vous devriez être capables de créer (et d'interpréter !) ce type de graphiques. Si ce n'est le cas, relisez d'urgence les sections

À retenir

- la fonction `summary()` utilise avec des fonctions statistiques qui ne renvoient qu'une valeur (par exemple `sd()` et `var()`)
- la fonction `summary()` utilise avec des fonctions statistiques qui renvoient plusieurs valeurs (par exemple `plot()` et `plot()`)

Les fonctions `summary()` et `plot()` acceptent un argument `by()` ou la fonction `plot()` mettront donc de calculer n'importe quel indice de statistique descriptive tableau de données entier ou sur des modalités ou combinaisons de modalités de facteurs. Il existe par ailleurs nombreuses fonctions, disponibles dans les packages spécifiques, qui permettent de fournir des résumés plus ou moins automatiques de tout ou partie des variables d'un jeu de données. Nous allons maintenant en décrire quelques-unes mais il en existe beaucoup d'autres : à vous d'explorer les possibilités et d'utiliser les fonctions qui vous paraissent les plus pertinentes, les plus simples à utiliser, les plus utiles ou les plus complètes.

6.3 Créer des résumés de données avec la fonction `summary()`

La fonction `summary()` ne pas confondre avec `summary()` permet d'obtenir des résumés de données pour tous types d'objets. La fonction `summary()` change la nature des résultats obtenus. Nous verrons cela dans les chapitres suivants. Cette fonction peut être utilisée pour examiner les résultats d'analyses de variances ou de modèles de régressions linéaires. Pour l'instant, nous nous intéressons à 3 situations :

1. ce que renvoie la fonction quand on lui fournit un vecteur

2. ce que renvoie la fonction quand on lui fournit un facteur

3. ce que renvoie la fonction quand on lui fournit un tableau

6.3.1 Variable continue : vecteur numérique

Commençons par fournir un vecteur numérique à la fonction `summary()`. Nous allons pour cela extraire les données de masses corporelles des manchots du tableau

```
peng$body_mass_g
```

```
[1] 3750 3800 3250 NA 3450 3650 3625 4675 3475 4250
[16] 3700 3450 4500 3325 4200 3400 3600 3800 3950 3800
[31] 3250 3900 3300 3900 3325 4150 3950 3550 3300 4650
[46] 4600 3425 2975 3450 4150 3500 4300 3450 4050 2900
[61] 3150 4400 3600 4050 2850 3950 3350 4100 3050 4450
[76] 4250 3700 3900 3550 4000 3200 4700 3800 4200 3350
[91] 3550 4300 3400 4450 3300 4300 3700 4350 2900 4100
[106] 3550 3750 3900 3175 4775 3825 4600 3200 4275 3900
[121] 3150 3500 3450 3875 3050 4000 3275 4300 3050 4000
[136] 3900 3175 3975 3400 4250 3400 3475 3050 3725 3000
[151] 3700 4000 4500 5700 4450 5700 5400 4550 4800 5200
[166] 5850 4200 5850 4150 6300 4800 5350 5700 5000 4400
[181] 4600 5550 5250 4700 5050 6050 5150 5400 4950 5200
[196] 4750 5550 4900 4200 5400 5100 5300 4850 5300 4400
[211] 4450 5550 4200 5300 4400 5650 4700 5700 4650 5800
[226] 5200 4700 5800 4600 6000 4750 5950 4625 5450 4700
[241] 4875 5550 4950 5400 4750 5650 4850 5200 4925 4800
[256] 5500 4725 5500 4700 5500 4575 5500 5000 5950 4600
[271] 4925 NA 4850 5750 5200 5400 3500 3900 3650 3525
[286] 3700 3800 3775 3700 4050 3575 4050 3300 3700 3400
[301] 3300 4150 3400 3800 3700 4550 3200 4300 3350 4100
[316] 4500 3950 3650 3550 3500 3675 4450 3400 4300 3200
[331] 3350 3450 3250 4050 3800 3525 3950 3650 3650 4000
```

Nous avons donc 344 valeurs de masses en grammes qui correspondent aux 344 manchots du jeu de données. La fonction `summary()` renvoie le résumé suivant lorsqu'on lui fournit ces valeurs :

```
summary(biondsy_mass_g)
```

```
Min. 1st Qu. Median Mean 3rd Qu. Max. NA's
2700  3550  4050  4202  4750  6300
```

Nous obtenons ici 7 valeurs, qui correspondent aux cinq valeurs renvoyées par la fonction `summary()`, ainsi que la moyenne et le nombre de valeurs manquant. Dans l'ordre, on a donc :

- la valeur minimale observée dans le vecteur. Ici, manchot le plus léger de l'échantillon pèse donc 2700 grammes.
- la valeur du premier quartile du vecteur. Ici, 25% manchots de l'échantillon (soit 86 individus) ont une masse inférieure à 3550 grammes, et 75% des individus de l'échantillon (soit 258 individus) ont une masse supérieure à 3550 grammes.
- la valeur de médiane du vecteur. La médiane est le deuxième quartile. Ici, 50% des manchots de l'échantillon (soit 172 individus) ont une masse inférieure à 4050 grammes, et 50% des individus de l'échantillon (soit 172 individus) ont une masse supérieure à 4050 grammes.
- la moyenne du vecteur. Ici, les manchots des 3 espèces du jeu de données ont en moyenne une masse de 4202 grammes.
- la valeur du troisième quartile du vecteur. Ici, 75% manchots de l'échantillon (soit 258 individus) ont une masse inférieure à 4750 grammes, et 25% des individus de l'échantillon (soit 86 individus) ont une masse supérieure à 4750 grammes.
- la valeur maximale observée dans le vecteur. Ici, manchot le plus lourd de l'échantillon pèse donc 6300 grammes.
- le nombre de données manquantes. Ici, 2 manchots n'ont pas été pesés et présentent donc la mention `NA` (comme `NA` available) pour `body_marsis_1g`.

Utiliser la fonction `summary()` fournit donc presque les mêmes informations :

quantile englobant mais, en R, il n'y a pas de fonction pour avoir des valeurs pour chaque modalités facteur (pour chaque espèce par exemple).

0% 25% 50% 75% 100%
2700 3550 4050 4750 6300

Attention, contrairement à ce que nous avons vu plus haut, la fonction `summary()` possède pas de `by` argument. Par conséquent, il n'est pas possible de se servir de la fonction pour avoir des valeurs pour chaque modalités facteur (pour chaque espèce par exemple).

Les différents indices statistiques fournis nous renseignent la fois sur la position de la distribution et sur la dispersion des données.

- La position correspond à la tendance centrale et indique quelles sont les valeurs qui caractérisent le grand nombre d'individus. La moyenne et la médiane sont les deux indices de position les plus fréquemment utilisés. Lorsqu'une variable a une distribution parfaitement symétrique, la moyenne et la médiane sont strictement égales. Mais lorsqu'une distribution est asymétrique, la moyenne et la médiane diffèrent. En particulier, la moyenne est beaucoup plus sensible aux valeurs extrêmes que la médiane. Cela signifie que quand une distribution est très asymétrique, la médiane est souvent une meilleure indication des valeurs les plus fréquemment observées.

L'histogramme de la Figure 6.4 illustre la distribution de la taille des manchots (toutes espèces confondues). Cette distribution présente une asymétrie à droite. Cela signifie que la distribution n'est pas symétrique et que la "queue de la distribution" est plus longue à droite qu'à gauche. La plupart des individus ont une masse comprise entre 3500 et 3700 grammes, au niveau du pic principal du graphique. La médiane, en orange et qui vaut 4050 grammes est plus proche du pic que la moyenne, en rouge, qui vaut 4202 grammes. La différence entre moyenne et médiane n'est pas énorme, elle peut le devenir si la distribution est vraiment très asymétrique, par exemple, si quelques individus seulement

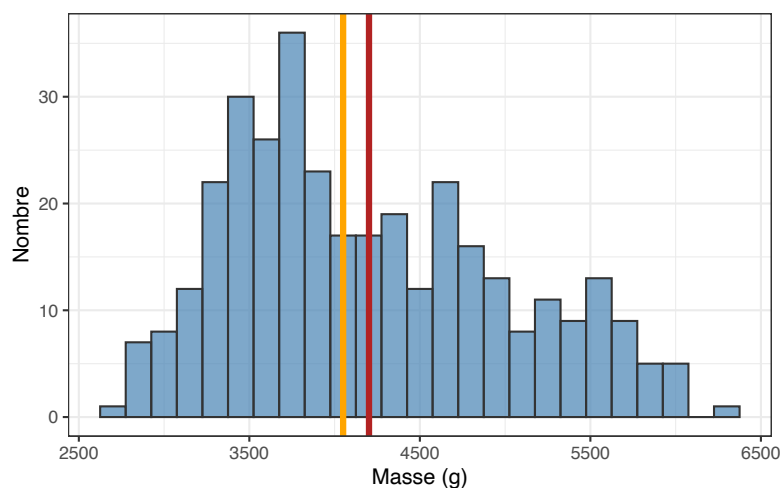


Figure 1: Distribution des masses corporelles des manchots

une masse supérieure à 7000 grammes, la moyenne serait tirée vers la droite du graphique alors que la médiane ne serait presque pas affectée. La moyenne représenterait alors moins fidèlement la tendance centrale.

Si l'on revient à un échantillon de données, on peut trouver des valeurs proches pour la moyenne et la médiane nous indiquent donc un degré de symétrie de la distribution.

- La dispersion des données nous renseigne sur la dispersion des points autour des indices de position. Les quartiles et les valeurs minimales et maximales (notées par I_{min} et I_{max}) nous renseignent sur l'étalement des points. Les valeurs situées entre le premier et le troisième quartile correspondent aux 50% des valeurs de l'échantillon les plus centrales. Plus l'intervalle entre ces quartiles (notée IQR pour "intervalle interquartile") sera grande, plus la dispersion sera importante. D'ailleurs, lorsque la dispersion est très grande, les moyennes et médianes ne renseignent que très peu sur la tendance centrale. Les indices de position sont surtout pertinents lorsque la dispersion des points autour de cette tendance centrale n'est pas trop large. Par exemple, si la distribution des données ressemblait à celle-ci (Figure 2), la moyenne et la médiane seraient fort peu utiles car très éloignées de la plupart des observations :

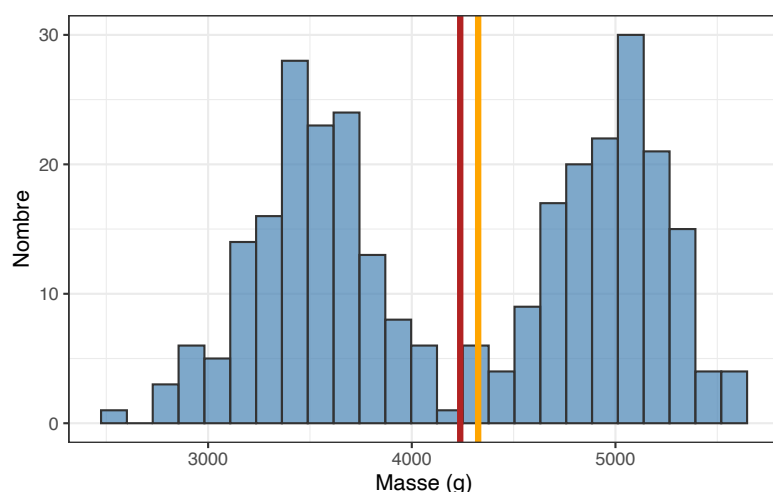


Figure 16: Distribution des masses corporelles (données fictives)

On comprend donc l'importance de considérer les indicateurs de dispersion en plus des indices de position pour caractériser et comprendre une série de données numériques. L'interquartile est toujours utile pour connaître l'étendue des données qui correspondent aux 50% des observations les plus centrales. Les autres indices de dispersion très fréquemment utilisés, mais qui ne sont pas proposés par défaut par la fonction `summary()` sont la variance et l'écart-type. Il est possible de calculer tous les indices renvoyés par la fonction `summary()` (ceux qui nous manquent grâce à la fonction `summary2()`)

```

# Calcul des statistiques descriptives pour les masses corporelles
summary2(mass)
#> # A tibble: 1 x 8
#>   min      Q1    med    moy      Q3    max    var    et
#>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1  2500  3000  3500  3500  4000  5500  10000  10000

```

```

# A tibble: 1 x 8
  min      Q1    med    moy      Q3    max    var    et
<dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>

```


6.3.2 Variable quantitative : facteur

Si l'on fournit une variable catégorielle à `summary()`, le résultat obtenu sera naturellement différent : calculs de moyennes, médianes ou quartiles n'aurait en effet pas de sens lorsque la variable fournie ne contient que des catégories.

```
summary(species)
```

```
Adelie Chinstrap   Gentoo
 152      68      124
```

Pour les facteurs, `summary()` ne compte simplement le nombre d'observations pour chaque modalité. Ici, la variable est un facteur qui compte 3 modalités. La fonction nous indique donc le nombre d'individus pour chaque modalité : notre jeu de données se compose de 152 individus de l'espèce Adélie, 68 individus de l'espèce Chinstrap et 124 individus de l'espèce Gentoo.

Comme pour les vecteurs numériques, si le facteur présente des données manquantes, `summary()` ne les ignore pas mais leur nombre :

```
summary(sex)
```

```
female  male  NA's
 165     168     11
```

Pour les facteurs, `summary()` nous fournit également une fonction équivalente à `table()` :

```
table(species)
```

```
# A tibble: 3 x 2
  species n
  <fct>    <int>
1 Adelie 152
2 Chinstrap 68
3 Gentoo 124
```

L'avantage de `count()` est qu'il est possible d'utiliser plusieurs facteurs pour compter le nombre d'observations de toutes les combinaisons de modalités (par exemple, combien d'individus de chaque sexe pour chaque espèce), ce qui n'est pas possible avec la fonction `summary()`.

6.3.3 Les tableaux à la fois

L'avantage de `summary()` par rapport à la fonction `count()` apparaît lorsque l'on souhaite obtenir des informations sur toutes les variables d'un tableau à la fois :

```
summary(penguins)
```

```

  species      island  bill_length_mm bill_depth_mm
Adelie       :152   Biscoe       :168   Min.       :32.10   Min.       :13.51
Chinstrap    :68   Dream         :124   1st Qu.    :39.23   1st Qu.    :17.80
Gentoo       :124   Torgersen    :52    Median    :44.45   Median    :19.13
              Mean    :43.92   Mean    :17.15
              3rd Qu.:48.50   3rd Qu.:18.70
              Max.    :59.60   Max.    :21.50
              NA's    :2   NA's    :2
 flipper_length_mm body_mass_g sex
Min.       :172.0   Min.       :2700   female:165   Min.       :200
1st Qu.    :190.0   1st Qu.    :3550   male :168   1st Qu.    :2700
Median    :197.0   Median    :4050   NA's  :11   Median    :2966
Mean      :200.9   Mean      :4202   Mean      :2008
3rd Qu.    :213.0   3rd Qu.    :4735   3rd Qu.    :2009
Max.      :231.0   Max.      :6300   Max.      :2009
NA's      :2     NA's      :2

```

Ici, on obtient un résumé pour chaque colonne du tableau. Les colonnes numériques sont traitées comme les vecteurs numériques (on obtient alors les minimas et maximas, les percentiles, les moyennes et médianes) et les colonnes contenant des variables catégorielles sont traitées comme des facteurs (on obtient alors le nombre d'observations pour chaque modalité).

On constate au passage que `sex` a été considérée ici comme une variable numérique, alors qu'elle devrait être

être considérée comme un facteur, ce qui nous permet de savoir combien d'individus ont été échantillonnés année :

```
penguins
mutate(year = factor(year))
summary
```

```

  species    island  bill_length_mm bill_depth_mm
Adelie     : 152    Biscoe      : 168    Min.       : 32.10    Min.       : 13
Chinstrap: 68     Dream      : 124    1st Qu.    : 39.23    1st Qu.
Gentoo     : 124    Torgersen: 52    Median     : 44.45    Median
              Mean      : 43.92    Mean      : 17.15
              3rd Qu.   : 48.50    3rd Qu.   : 18.70
              Max.      : 59.60    Max.      : 21.50
              NA's      : 2      NA's      : 2
flipper_length_mm body_mass_g
Min.      : 172.0    Min.      : 2700    female: 165    2007: 110
1st Qu.   : 190.0    1st Qu.   : 3550    male  : 168    2008: 114
Median    : 197.0    Median    : 4050    NA's   : 11    2009: 120
Mean      : 200.9    Mean      : 4202
3rd Qu.   : 213.0    3rd Qu.   : 4750
Max.      : 231.0    Max.      : 6300
NA's      : 2      NA's      : 2
```

Au final, la fonction `summary()` est très utile dans certaines situations, notamment pour avoir rapidement accès à des statistiques descriptives simples sur toutes les colonnes du tableau. Elle reste cependant limitée car d'une part, elle ne fournit pas les variances ou les écarts-types pour les variables numériques, et il est impossible d'avoir des résumés par modalité pour chaque facteur par exemple. Ici, il paraît en effet intéressant d'avoir des informations synthétiques concernant les mesures biométriques des manchots, espèce par espèce, plutôt que toutes espèces confondues. C'est là que la fonction `sketch()` intervient.

6.4 Créer des résumés de données avec la fonction `skim()`

La fonction `skim()` fait partie du package `skimr`. Avant de pouvoir l'utiliser, pensez donc à l'installer et à le mémoriser si ce n'est pas déjà fait. Comme pour la fonction `summary()`, `skim()` peut utiliser `skimr::drfp_includes()` pour sélectionner des types d'objets. Nous nous contenterons d'examiner ici le plus fréquent : celui des tableaux, groupés avec ou non.

6.4.1 Tableau non groupé

Commençons par examiner le résultat avec le tableau `penguins` non groupé :

```
skim(penguins)
```

Table Data summary

Name	penguins
Number of rows	344
Number of columns	8

Column type frequency:	
factor	3
numeric	5

Group variables	None

Variable type: factor

skim_variable					n_missing	n_complete	ratio	ordered	unique	top
species	0	1.00	FALSE	3	Ade: 152, Gen: 124, Chi: 68					
island	0	1.00	FALSE	3	Bis: 168, Dre: 124, Tor: 52					
sex	11	0.97	FALSE	2	mal: 168, fem: 165					

Variable type: numeric

```

skim_variable |> summarise(
  bill_length_mm = 43.925,
  bill_depth_mm = 17.151,
  flipper_length_mm = 199.214,
  body_mass_g = 4201.75,
  year = 2008.082
)

```

Les résultats obtenus grâce à cette fonction sont nommés `summary`. La première section nous donne des informations sur le tableau :

- son nom, son nombre de lignes et de colonnes
- la nature des variables qu'il contient (ici 3 facteurs et 3 variables numériques)
- la présence de variables utilisées pour faire des regroupements (il n'y en a pas encore à ce stade)

Ensuite, un bloc apporte des informations sur chaque facteur présent dans le tableau :

- le nom de la variable (`species`)
- le nombre de données manquantes (taux de "données manquantes")
- des informations sur les modalités (niveaux de modalités)
- le nombre d'observations pour les modalités les plus représentées

En un coup d'œil, on sait donc que 3 espèces sont présentes (et on connaît leurs effectifs), on sait que les manchots ont été échantillonnés sur 3 îles, et on sait que le sexe des individus sur 344 (soit 3%) est inconnu. Pour le reste, presque autant de femelles que de mâles.

Le dernier bloc renseigne sur les variables numériques. Pour chaque d'elle, on a :

- le nom de la variable (`bill_length_mm`)
- le nombre de données manquantes (taux de "données manquantes")
- la moyenne (et l'écart-type) (est une nouveauté par rapport à la fonction `summary`)

- les valeurs `summary(penguins$bill_length_mm)` et `summary(penguins$bill_depth_mm)` (la médiane !), de troisième quartile et la valeur maximale
- un histogramme et un `geom_histogram()` (un premier aperçu grossier de la forme de la distribution)

Là encore, en un coup d'œil, on dispose donc de toutes les informations pertinentes pour juger de la distribution et de la dispersion de chaque variable numérique de données.

6.4.2 Tableau groupé

La fonction `summary()` déjà très pratique, le devient encore plus lorsque l'on choisit de lui fournir seulement certaines variables, et qu'on fait certains regroupements. Par exemple, on peut sélectionner les variables relatives aux dimensions du bec (`bill_length_mm` et `bill_depth_mm`) et demander un résumé des données pour chaque espèce grâce à la fonction `summary_by()` que nous connaissons déjà, et demander un résumé des données pour chaque espèce grâce à la fonction `summary_by()` que nous connaissons également :

```
penguins %>% summarise(
  # Avec le tableau penguins...
  species_summary = summary(
    bill_length_mm,
    bill_depth_mm) # Je sélectionne les variables d'intérêt
  group_by(species) # Je regroupe par espèce...
) %>% summarise(
  # Et je produis un résumé des données
```

Tableau 6.4.4 summary

Name	group_by (...)
Number of rows	344
Number of columns	3

Column type frequency:	
numeric	2

Group variables	species

Variable type: numeric

skim_variables | summarise(sps1 = sd(ip1), sps2 = sd(ip2), sps3 = sd(ip3), sps4 = sd(ip4), sps5 = sd(ip5), sps6 = sd(ip6), sps7 = sd(ip7), sps8 = sd(ip8), sps9 = sd(ip9), sps10 = sd(ip10), sps11 = sd(ip11), sps12 = sd(ip12), sps13 = sd(ip13), sps14 = sd(ip14), sps15 = sd(ip15), sps16 = sd(ip16), sps17 = sd(ip17), sps18 = sd(ip18), sps19 = sd(ip19), sps20 = sd(ip20), sps21 = sd(ip21), sps22 = sd(ip22), sps23 = sd(ip23), sps24 = sd(ip24), sps25 = sd(ip25), sps26 = sd(ip26), sps27 = sd(ip27), sps28 = sd(ip28), sps29 = sd(ip29), sps30 = sd(ip30), sps31 = sd(ip31), sps32 = sd(ip32), sps33 = sd(ip33), sps34 = sd(ip34), sps35 = sd(ip35), sps36 = sd(ip36), sps37 = sd(ip37), sps38 = sd(ip38), sps39 = sd(ip39), sps40 = sd(ip40), sps41 = sd(ip41), sps42 = sd(ip42), sps43 = sd(ip43), sps44 = sd(ip44), sps45 = sd(ip45), sps46 = sd(ip46), sps47 = sd(ip47), sps48 = sd(ip48), sps49 = sd(ip49), sps50 = sd(ip50), sps51 = sd(ip51), sps52 = sd(ip52), sps53 = sd(ip53), sps54 = sd(ip54), sps55 = sd(ip55), sps56 = sd(ip56), sps57 = sd(ip57), sps58 = sd(ip58), sps59 = sd(ip59), sps60 = sd(ip60), sps61 = sd(ip61), sps62 = sd(ip62), sps63 = sd(ip63), sps64 = sd(ip64), sps65 = sd(ip65), sps66 = sd(ip66), sps67 = sd(ip67), sps68 = sd(ip68), sps69 = sd(ip69), sps70 = sd(ip70), sps71 = sd(ip71), sps72 = sd(ip72), sps73 = sd(ip73), sps74 = sd(ip74), sps75 = sd(ip75), sps76 = sd(ip76), sps77 = sd(ip77), sps78 = sd(ip78), sps79 = sd(ip79), sps80 = sd(ip80), sps81 = sd(ip81), sps82 = sd(ip82), sps83 = sd(ip83), sps84 = sd(ip84), sps85 = sd(ip85), sps86 = sd(ip86), sps87 = sd(ip87), sps88 = sd(ip88), sps89 = sd(ip89), sps90 = sd(ip90), sps91 = sd(ip91), sps92 = sd(ip92), sps93 = sd(ip93), sps94 = sd(ip94), sps95 = sd(ip95), sps96 = sd(ip96), sps97 = sd(ip97), sps98 = sd(ip98), sps99 = sd(ip99), sps100 = sd(ip100))

bi | IA_dledriget h1_mon 99 38.792.6632.136.7538.8040.7546.
bi | IC_hliensgttrha_mon 1.00 48.833.3440.946.3549.5551.0858.
bi | IG_elnetrgt h1_mon 99 47.503.0840.945.3047.3049.5559.6
bi | IA_dledpith1mm0.99 18.351.2215.517.5018.4019.0021.
bi | IC_hdmpsttr_mon 1.00 18.421.1416.417.5018.4519.4020.
bi | IG_edetp1mm0.99 14.980.9813.114.2015.0015.7017.3

On constate ici que pour chaque variable numérique sélectionnée, des statistiques descriptives détaillées nous sont fournies pour chacune des 3 espèces. Ce premier examen semble montrer que :

- L'espèce Adélie est celle qui possède le bec le plus épais (ses valeurs de moyennes, médianes et quartiles sont plus faibles que celles des 2 autres espèces).
- L'espèce Gentoo est celle qui possède le bec le plus fin ou le moins épais (ses valeurs de moyennes, médianes et quartiles sont plus faibles que celles des 2 autres espèces)
- Il ne semble pas y avoir de fortes différences d'écart-types (donc, une dispersion comparable des points autour de leur moyenne respective) entre les 3 espèces pour chacune des 2 variables numériques, des valeurs d'écart-types comparables sont en effet observées pour les 3 espèces
- La distribution des 2 variables numériques semble approximativement suivre une distribution symétrique pour les 3 espèces, avec une forme de courbe en cloche. Les distributions devraient donc suivre à peu près une distribution normale

Note

Vous comprenez j'espère l'importance d'examiner ce genre de résumé des données avant de vous lancer dans des tests statistiques. Ils sont un complément indispensable aux explorations graphiques que vous devez également prendre l'habitude de réaliser pour mieux appréhender et comprendre la nature de vos données. Puisque chaque jeu de données est unique, vous devrez vous adapter à la situation et aux questions scientifiques

vous seront posées (ou que vous vous poserez !) : les choix qui seront pertinents pour une situation ne le seront pas nécessairement pour une autre. Mais dans tous les cas, pour savoir où vous allez et pour ne pas faire bêtise au moment des tests statistiques et de leur interprétation, vous devrez toujours explorer vos données, avec des graphiques exploratoires et des statistiques (descriptives

6.5 Exercice

En utilisant les fonctions de résumé abordées jusqu'ici, répondez aux questions suivantes :

1. Dans quel aéroport de New York les précipitations moyennes ont-elle été les plus fortes en 2013 ?
2. Dans quel aéroport de New York la vitesse du vent moyenne était-elle la plus forte en 2013 ? Quelle est cette vitesse ?
3. Dans quel aéroport de New York les rafales de vent étaient-elles les plus variables en 2013 ? Quel indicateur statistique vous donne cette information et quelle est sa valeur ?
4. Les précipitations dans les 3 aéroports de New-York ont-elles une distribution symétrique ?
5. Quelle est la température médiane observée en 2013 pour tous aéroports confondus ?
6. Tous aéroports confondus, quel est le mois de l'année où la température a été la plus variable en 2013 ? Quel indicateur statistique vous donne cette information ? Étaient les températures minimales et maximales observées ce mois-là ?

7 Dispersion et incertitude

7.1 Pré-requis

Nous avons ici besoin des packages suivants :

```
library(tidyverse)
library(ggplot2)
library(scales)
```

Pensez à les charger en mémoire si ce n'est pas déjà fait si vous venez de démarrer une nouvelle session de travail.

Le package (Wickham, Pedersen, et Seidel 2023) contient de très nombreuses fonctions particulières utiles pour améliorer l'aspect des légendes d'axes. Par exemple, pour transformer des chiffres compris entre 0 et 1 en pourcentages, ou pour ajouter le symbole d'une devise quand l'axe renseigne sur des montants en euros ou dollars (par exemple).

7.2 La notion de dispersion

Comme expliqué à la section 6.3, les indices de dispersion nous renseignent sur la variabilité des données au-delà de la valeur centrale (moyenne ou médiane) d'une population ou d'un échantillon. L'écart-type, la variance et l'interquartile sont 3 exemples d'indices de dispersion. L'exemple de l'écart-type. Un écart-type faible indique que la majorité des observations ont des valeurs proches de la moyenne. À l'inverse, un écart-type important indique que la plupart des points sont éloignés de la moyenne. L'écart-type est une caractéristique de la population que l'on estime grâce à un échantillon, au même titre que la moyenne. En travaillant sur un échantillon, on espère accéder aux grandeurs de la population. Même si ces vraies grandeurs

sont à jamais inaccessibles (on ne connaîtra jamais précisément quelle est la vraie valeur de la moyenne type de la population), on espère qu'avec un échantillonnage réalisé correctement, la moyenne de l'échantillon) et l'écart type de l'échantillon reflètent assez fidèlement les valeurs de la population générale. C'est la notion d'estimateur, intimement liée à la notion d'estimation de la moyenne de l'échantillon, que l'on connaît avec précision, est un estimateur de la moyenne de la population qui restera à jamais inconnue. C'est la raison pour laquelle la moyenne de l'échantillon n'est pas parfois notée \bar{x} ou \bar{y} . De même, l'écart type de l'échantillon sont des estimateurs de l'écart-type de la population générale. C'est la raison pour laquelle on les note s et s_y respectivement. L'accent circonflexe se prononce "chapeau". On dit donc que l'écart-type de l'échantillon) est un estimateur de l'écart-type de la population générale. Comme nous l'avons vu, les indices de dispersion doivent accompagner les valeurs de position lorsque l'on décrit des données, car présence d'une valeur de moyenne, ou de médiane seule n'a pas de sens sans savoir à quel point les données sont proches ou éloignées de la tendance centrale pour savoir si, dans la population générale, les indicateurs de position correspondent ou non aux valeurs portées par la majorité des individus.

Nous avons vu comment calculer des indices de position et de dispersion. Tout ceci devrait donc être utile pour vous à ce stade.

7.3 La notion d'incertitude

Par ailleurs, puisque on ne sait jamais avec certitude les estimations (de moyennes ou d'écart-types ou de tout autre paramètre) reflètent fidèlement ou non les vraies valeurs de la population, nous devons quantifier à quel point nos estimations s'écartent de celles de la population générale. C'est tout l'intérêt des statistiques et c'est ce que permettent les indices d'incertitude. On ne connaîtra jamais la vraie valeur de la moyenne ou d'écart-type de la population, mais on peut quantifier à quel point nos estimations (basées sur l'échantillon) sont précises ou imprécises.

Les deux indices d'incertitude les plus connus (et les plus utilisés) sont :

1. l'intervalle de confiance à 95% (ou à 99,5% ou à tout autre niveau) ; les formules sont nombreuses et il n'est pas utile de les détailler ici : nous verrons comment les calculer plus bas)
2. l'erreur standard de la moyenne, dont la formule est la suivante :

$$= \frac{s}{\sqrt{n}}$$

avec l'écart-type de l'échantillon et n la taille de l'échantillon.

Comme pour la moyenne, on peut calculer l'erreur standard d'un écart-type, d'une médiane, d'une proportion, ou de tout autre estimateur calculé sur un échantillon. Cet indice d'incertitude ne nous renseigne pas sur une grandeur de la population générale qu'on chercherait à estimer, mais bien sur l'incertitude associée à une estimation que nous faisons en travaillant sur un échantillon de taille finie et limitée. Tout processus d'échantillonnage est forcément entaché d'incertitude, causée entre autre par le hasard de l'échantillonnage (ou fluctuation d'échantillonnage). Comme nous travaillons sur des échantillons forcément imparfaits, les indices d'incertitude vont nous permettre de quantifier à quel point nos estimations s'écartent des vraies valeurs de la population. Ces "vraies valeurs", faute de pouvoir collecter tous les individus de la population, restent toujours jamais inconnues.

Autrement dit ...

Quand on étudie des populations naturelles grâce à des échantillons statistiques, on se trompe toujours un peu. Mais grâce aux indices d'incertitude, on peut savoir à quel point on se trompe et c'est déjà pas mal !

En examinant la formule de l'erreur standard de la moyenne présentée ci-dessus, on comprend intuitivement que plus la taille de l'échantillon est grande, plus l'erreur standard

(donc l'incertitude) associée à notre estimation de moyenne diminue. Autrement dit, plus les données sont abondantes sur l'échantillon, meilleure sera notre estimation de moyenne. Plus donc, moins le risque de raconter des bêtises sera grand.

L'autre indice d'incertitude très fréquemment utilisé est l'intervalle de confiance à 95% (de la moyenne, de la médiane, de la variance, ou de tout autre estimateur calculé dans l'échantillon). L'intervalle de confiance nous renseigne sur la gamme des valeurs les plus probables pour un paramètre de la population étudiée. Par exemple, si j'observe, dans un échantillon, une moyenne de 10, avec un intervalle de confiance calculé de $[7; 15]$, cela signifie que, dans la population générale, la vraie valeur de moyenne a de bonnes chances de se trouver dans l'intervalle $[7; 15]$. Dans la population générale, toutes les valeurs comprises entre 7 et 15 sont vraisemblables pour la moyenne alors que les valeurs situées en dehors de cet intervalle sont moins probables. Une façon de comprendre l'intervalle de confiance est de dire que si je récupère un grand nombre d'échantillons dans la même population, en utilisant exactement le même protocole expérimental, 95% des échantillons que je vais récupérer auront une moyenne située à l'intérieur de l'intervalle de confiance à 95%, et 5% des échantillons auront une moyenne située à l'extérieur de l'intervalle de confiance à 95%. C'est une notion qui n'est pas si évidente que ça à comprendre, donc prenez bien le temps de relire cette section si besoin et de poser des questions le cas échéant.

Concrètement, plus l'intervalle de confiance est large, plus notre confiance est grande. Si la moyenne d'un échantillon vaut 10 et que son intervalle de confiance à 95% vaut $[9,5; 11]$, la gamme des valeurs probables pour la moyenne de la population est étroite. Autrement dit, la moyenne de l'échantillon (10), a de bonnes chances d'être très proche de la vraie valeur de moyenne de la population générale (vraisemblablement comprise quelque part entre 9,5 et 11). À l'inverse, si l'intervalle de confiance à 95% de la moyenne vaut $[4; 16]$, la gamme des valeurs possibles pour la vraie moyenne de la population est grande. La moyenne de l'échantillon aura donc de grandes chances d'être assez éloignée de la vraie valeur de la population.

La notion d'intervalle de confiance à 95% est donc très proche de celle d'erreur standard. D'ailleurs, pour de nombreux

mètres, l'intervalle de confiance est calculé à partir standard.

0 À retenir !

Les paramètres sont des caractéristiques des populations que l'on étudie. On ne cherche donc pas à les estimer. Ils nous permettent en revanche de quantifier à quel point on se trompe en cherchant à estimer des paramètres de la population générale à partir d'estimateurs calculés sur un échantillon.

Autrement dit, calculer des indicateurs de dispersion ne permet pas d'apprendre des choses au sujet des populations étudiées. Calculer des indicateurs d'incertitude permet pas d'apprendre quoi que ce soit sur les populations étudiées, mais permet d'apprendre des choses sur la qualité de notre travail d'échantillonnage et d'estimation. En général, les gammes de valeurs auxquelles les vraies valeurs des paramètres des populations générales ont de bonnes chances de se trouver. Si on a bien travaillé, et qu'on dispose de beaucoup de données, ces gammes de valeurs seront peu étendues. Si à l'inverse, nous ne disposons que de trop peu de données, ces gammes de valeurs seront très étendues.

7.4 Calcul de l'erreur standard de la moyenne

Contrairement aux indices de position et de dispersion, n'existe pas de formule pour calculer l'erreur standard de la moyenne. Toutefois, sa formule simple nous permet de la calculer à la main quand on en a besoin grâce à la formule suivante :

Par exemple, reprenons les données de masse corporelle de 3 espèces de manchots :

que nous souhaitons étudier les différences de masses corporelles des 3 espèces, en tenant compte du sexe des individus. Pour chaque espèce et chaque sexe, nous allons calculer la masse moyenne des individus en grammes (variable `body_mass_g`). Nous prendrons soin au préalable d'éliminer les lignes pour lesquelles le sexe des individus est in-

```
penguins
  filter(!is.na(sex)) >
  summarise(moyenne = body_mass_g
    . by (species, sex) )
```

```
# A tibble: 6 x 3
  species sex    moyenne
  <fct>   <fct>   <dbl>
1 Adelie male    4043.
2 Adelie female  3369.
3 Gentoo female  4680.
4 Gentoo male    5485.
5 Chinstrap female 3527.
6 Chinstrap male   3939.
```

Pour pouvoir réutiliser ce tableau, je lui donne un nom :

```
masse-penguins
  filter(!is.na(sex)) >
  summarise(moyenne = body_mass_g
    . by (species, sex) )
```

Au final, on peut faire un graphique avec ces données. Puisqu'on dispose d'une variable numérique et 2 variables catégorielles, je fais le choix de produire un diagramme à bâtons facetés :

```
masse
  ggplot(aes(x = sex, y = moyenne, fill = species))
  geom_col(color = "grey20", width = 0.5)
  facet_wrap(~ species)
  labs(x = "", y = "Masse moyenne")
  theme_bw()
  scale_fill_manual(values = c("red", "blue", "green"))
  scale_y_continuous(labels = number_format())
```

`scale_x_discrete(labels = c("Femelle", "Mâle"))`

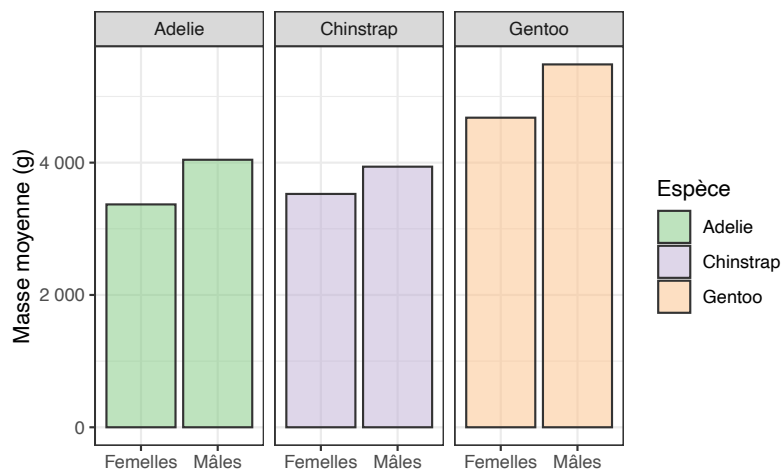


Figure 6.7 : Comparaison des tailles moyennes observées chez les mâles et les femelles de 3 espèces de manchots.

Vous remarquerez que :

1. J'utilise `geom_bar()` : j'ai déjà calculé manuellement la variable que je souhaite afficher à l'axe des ordonnées.
2. J'associe la couleur de remplissage à l'espèce, bien que ce ne soit pas indispensable puisque je fais un graphique par espèce. Cela rend toutefois le graphique plus agréable et facilite la lecture.
3. Je modifie 3 échelles :
 - avec `scale_fill_brewer()` la palette de couleur utilisée pour le remplissage des barres.
 - avec `scale_y_continuous()` l'échelle (continue) de l'axe des ordonnées. Je fais appel à la fonction `number_format()` afin d'ajouter un séparateur des milliers aux chiffres de l'axe.
 - avec `scale_x_discrete()` les termes qui apparaissent sur l'axe des abscisses (qui est un axe discret). Les catégories de la variable `sex` sont transformées en "Femelle" et "Mâle" respectivement.

Puisque chaque barre possède une valeur de moyenne, il nous faut l'indiquer également.

associée à ces calculs de moyennes, nous devons calculer l'erreur standard des moyennes :

```

penguins
  filtfilt(
    summarise(
      N_obs = n(),
      erreur_standard = mass_gram / sqrt(N_obs),
      by = c(species, sex)
    )
  )

```

```

# A tibble: 6 x 5
  species sex moyenne N_obs erreur_standard
  <fct>    <fct>    <dbl> <int>    <dbl>
1 Adelie  male    4043.    73     48.6
2 Adelie  female  3369.    37     37.35
3 Gentoo  female  4680.    35     58.0
4 Gentoo  male    5485.    40     40.1
5 Chinstrap female  3527.    34     48.34
6 Chinstrap male    3939.    62     34.1

```

Notre tableau de statistiques descriptives possède maintenant 2 colonnes supplémentaires : le nombre d'observation j'ai noté, et l'erreur standard associée à chaque moyenne, calculée grâce à la formule vue plus haut $\text{erreur_standard} = \frac{\text{moyenne}}{\sqrt{N_obs}}$ (la fonction permet de calculer la racine carrée). On constate que l'erreur standard, qui s'exprime dans la même unité que la moyenne, varie du simple au double selon les groupes. Ainsi, pour les mâles de l'espèce Chinstrap, l'incertitude est deux fois plus importante que pour les femelles de l'espèce Adélie. Cela est probablement dû à des différences de tailles d'échantillons importantes dans ces 2 catégories (73 femelles Adélie contre 34 mâles Chinstrap), mais ça n'est certainement pas la seule explication. Sinon, les femelles Chinstrap auraient elles aussi une incertitude plus grande. L'incertitude reflète aussi, de façon directe, la variabilité de la variable étudiée.

Une fois de plus, je donne un nom à ce tableau de données pour pouvoir le réutiliser plus tard :

```

masses_penguins
  filtfilt(
    summarise(
      N_obs = n(),
      erreur_standard = mass_gram / sqrt(N_obs),
      by = c(species, sex)
    )
  )

```



```

N_obs ~ dnorm(mu, sigma^2),
erreur_std ~ dnorm(mu_err, sigma_err^2),
. by (species, sex) )

```

Notez que la fonction `lapply()` permet de calculer à la fois la moyenne et erreur standard. La moyenne d'une colonne est calculée pour chaque espèce et sexe. La fonction renvoie 3 valeurs : la moyenne, l'erreur standard et la variance.

```

penguins
filter(sex != "female")
reframe(body_mass_g,
. by (species, sex))

```

```

# A tibble: 6 x 5
  species sex y_min y_max
<fct>    <fct> <dbl> <dbl> <dbl>
1 Adelie  male  4043. 4003. 4084.
2 Adelie  female 3369. 3337. 3400.
3 Gentoo  female 4680. 4643. 4717.
4 Gentoo  male  5485. 5445. 5525.
5 Chinstrap female 3527. 3478. 3576.
6 Chinstrap male  3939. 3877. 4001.

```

Les résultats obtenus ne sont pas exactement au même format :

- la colonne `y_min` contient les valeurs de moyennes
- la colonne `y_max` contient la valeur de moyenne moins une fois l'erreur standard
- la colonne `y_max` contient la valeur de moyenne plus une fois l'erreur standard

Notez également que contrairement aux fonctions `sd()` la fonction `sd()` n'a pas besoin qu'on lui précise `na.rm = TRUE` par défaut, elle ignore les valeurs manquantes.

Il ne nous restera plus qu'à ajouter des barres d'erreur sur notre graphique pour visualiser l'incertitude associée à chaque valeur de moyenne. C'est ce que nous verrons au Chapitre 8.

7.5 Calculs d'intervalles de confiance à 95 %

Comme pour les erreurs standard, il est possible de calculer des intervalles de confiance de n'importe quel est calculé à partir d'un échantillon, pour déterminer la des valeurs les plus probables pour les paramètres équ dans la population générale. Nous nous concentrerons le calcul des intervalles de confiance à 95% de la moyenne mais nous serons amenés à examiner également l'interv de confiance de la médiane [1.1.1. Intervalles de confiance de la médiane](#) et l'interv de confiance à 95% d'une différence de moyennes.

Contrairement à l'erreur standard, il n'y a pas qu'une seule bonne façon de calculer l'intervalle de confiance à 95% d'une moyenne. Plusieurs formules existent et le choix de la formule dépend en partie de la distribution des données (la distribution suit-elle une loi Normale ou non) et de la taille de l'échantillon dont on dispose (plus ou moins de 30 observations ou non?). Dans la situation idéale d'une variable quantitative à distribution Normale, les bornes inférieures et supérieures de l'intervalle de confiance à 95% sont obtenues grâce à la formule

$$-1 \underline{96} < +1 \underline{96}$$

Autrement dit, la ~~draine~~ ~~moyenne~~ ~~population~~ a de bonnes chances de se trouver dans un intervalle de plus ou moins 1.96 fois l'erreur standard de la moyenne. En première approximation, l'intervalle de confiance est donc la moitié de l'échantillon ~~tu~~ ~~si~~ ~~tu~~ ~~ou~~ ~~moins~~ 2 fois l'erreur standard (que nous avons appris à calculer à la main un peu plus tôt). On peut donc calculer à la main les bornes inférieures et supérieures de l'intervalle de confiance ainsi :

```

penguins
  filter(!is.na(sex))
  summarise(mass_mean = body_mass_mean,
            .by = (species, sex))

```

```
# A tibble: 6 x 5
  species      sex y_min y_max
```

```

  <fct>      <fct> <dbl> <dbl> <dbl>
1 Adelie     male  4043.  3964.  4123.
2 Adelie     female 3369.  3307.  3431.
3 Gentoo     female 4680.  4607.  4752.
4 Gentoo     male  5485.  5406.  5563.
5 Chinstrap female 3527.  3431.  3623.
6 Chinstrap male  3939.  3817.  4061.

```

Ici, grâce à `mutate()` et à la fonction `mean_se()`

- la colonne `ymin` contient maintenant les valeurs de moyennes moins 1.96 fois l'erreur standard
- la colonne `ymax` contient maintenant les valeurs de moyennes plus 1.96 fois l'erreur standard

Dans la pratique, puisque cette méthode reste approximative et dépend de la nature des données dont on dispose, on utilisera plutôt des fonctions spécifiques qui calculent pour nous les intervalles de confiance à 95% de nos estimateurs. C'est ce que permet en particulier la fonction `mean_cl_normal()` de `ggplot2`. Il est toutefois important de bien comprendre qu'il y a un lien étroit entre l'erreur standard (l'incertitude associée à l'estimation d'un paramètre d'une population à partir des données d'un échantillon), et l'intervalle de confiance à 95% de ce paramètre.

```

penguins
  filter(!is.na(sex))
  summarise(mean_cl = mean_cl_normal(mass_g),
            . by = c(species, sex))

```

```

# A tibble: 6 x 5
  species  sex  ymin  ymax
  <fct>    <fct> <dbl> <dbl> <dbl>
1 Adelie  male  4043.  3963.  4124.
2 Adelie  female 3369.  3306.  3432.
3 Gentoo  female 4680.  4606.  4754.
4 Gentoo  male  5485.  5405.  5565.
5 Chinstrap female 3527.  3428.  3627.
6 Chinstrap male  3939.  3813.  4065.

```

Comme dans les tableaux précédents, 3 nouvelles colonnes ont été créées :

- `y` contient toujours la moyenne des températures mensuelles pour chaque aéroport
- `ym` contient maintenant les bornes inférieures de l'intervalle à 95% des moyennes
- `ymax` contient maintenant les bornes supérieures de l'intervalle à 95% des moyennes

Pour que la suite soit plus claire, nous allons chercher des noms à ces différents tableaux en prenant soin de renommer les colonnes pour plus de clarté.

Tout d'abord, nous donnons un nom au tableau qui contient, les masses moyennes des mâles et des femelles de 3 espèces, et les erreurs standard de ces moyennes :

```
masses_se
```

```
# A tibble: 6 x 5
  species sex moyenne N_obs erreur_standard
  <fct>   <fct>   <dbl> <int> <dbl>
1 Adelie male    4043.    40  6
2 Adelie female  3369.    37 35
3 Gentoo female  4680.    37 80
4 Gentoo male    5485.    40  1
5 Chinstrap female 3527.48 394
6 Chinstrap male   3939.    62 34
```

Ensuite, nous avons produit un tableau presque équivalent que nous allons nommer `masses_bornes` pour lequel nous allons modifier les noms des colonnes

```
masses_se %>%
  filter(!is.na(sex)) %>%
  summarise(moyenne =
    mean(body_mass_g),
    ym = by(species, sex) %>%
      summarise(moyenne =
        mean(body_mass_g),
        ymax =
        mean(body_mass_g))
  )
masses_se_bornes
```

```
# A tibble: 6 x 5
  species sex      moyenne moyenne_moins_se moyenne_pl
  <fct>    <fct>    <dbl> <dbl> <dbl>
1 Adelie   male      4043 4003.    4084.
2 Adelie   female    3369 3337.    3400.
3 Gentoo   female    4684 4643.    4717.
4 Gentoo   male      5485 5445.    5525.
5 Chinstrap female    3524 3477.8.   3576.
6 Chinstrap male      3939 3877.    4001.
```

Nous avons ensuite calculé manuellement des intervalles de confiance approximatifs en utilisant la fonction `mean_cl_normal()` de la bibliothèque `ggplot2`. Encore, nous allons stocker cet objet dans un tableau sous le nom `masses_ci` puis nous allons modifier le `type` des données.

```
masses_ci <- append(masses,
  file = "masses_ci.csv",
  ref = mean_cl_normal(masses$body_masse,
    by = c(species, sex))
  )
# A tibble: 6 x 5
  species sex      moyenne ci_borne_inf ci_borne_sup
  <fct>    <fct>    <dbl> <dbl> <dbl>
1 Adelie   male      4043 4003.    4084.
2 Adelie   female    3369 3337.    3400.
3 Gentoo   female    4684 4643.    4717.
4 Gentoo   male      5485 5445.    5525.
5 Chinstrap female    3524 3477.8.   3576.
6 Chinstrap male      3939 3877.    4001.
```

```
# A tibble: 6 x 5
  species sex      moyenne ci_borne_inf ci_borne_sup
  <fct>    <fct>    <dbl> <dbl> <dbl>
1 Adelie   male      4043 4003.    4084.
2 Adelie   female    3369 3337.    3400.
3 Gentoo   female    4684 4643.    4717.
4 Gentoo   male      5485 5445.    5525.
5 Chinstrap female    3524 3477.8.   3576.
6 Chinstrap male      3939 3877.    4001.
```

Enfin, nous avons calculé les intervalles de confiance en utilisant une fonction spécialement dédiée à cette tâche : la fonction `mean_cl_normal()` de la bibliothèque `ggplot2`. Nous allons stocker cet objet dans un tableau sous le nom `masses_ci` puis nous allons modifier le nom des colonnes `ci_borne_inf` et `ci_borne_sup` en `ci_min` et `ci_max`.

```

masses_penguins
  fill(!is.na(sex)) >
  refr(amc_n_ci_(bordr_ymass_g),
    . by c(species, sex))
rend(moyenne =
  ci_borne_y_min, =
  ci_borne_y_max) =

masses_ci

```

```

# A tibble: 6 x 5
  species sex moyenne ci_borne_inf ci_borne_sup
<fct> <fct> <dbl> <dbl> <dbl>
1 Adelie male 4034.96 3. 4124.
2 Adelie female 3336.96 3. 3432.
3 Gentoo female 4468.06 4. 4754.
4 Gentoo male 5458.50 5. 5565.
5 Chinstrap female 3352.7 3. 3627.
6 Chinstrap male 3939. 4. 4065.

```

Maintenant, si l'on compare les 2 tableaux contenant l'culs d'intervalles de confiance de la moyenne, on const les résultats sont très proches :

```

masses_ci_approx
masses_ci

```

```

# A tibble: 6 x 5 # A tibble: 6 x 5
  species sex smoyenne ci_borne_inf ci_borne_sup
<fct> <fct> <dbl> <dbl> <dbl>
1 Adelie male 4034.96 3. 4124.
2 Adelie female 3336.96 3. 3432.
3 Gentoo female 4468.06 4. 4754.
4 Gentoo male 5458.50 5. 5565.
5 Chinstrap female 3352.7 3. 3627.
6 Chinstrap male 3939. 4. 4065.

```

Les bornes inférieures et supérieures des intervalles de confiance à 95% des moyennes ne sont pas égales quand on les calcule manuellement de façon approchée et quand on calcule de façon exacte, mais les différences sont min

8 Visualiser l'incertitude de dispersion

8.1 Pré-requis

Nous avons ici besoin des packages suivants :

```
library(triDyverse)
library(ape)
library(merpenguins)
library(arycflights13)
library(ascyales)
library(ascyolorspace)
```

Le package `ascyolorspace` (plahaek et al., 2024) permet d'utiliser de très nombreuses palettes de couleurs. Ces palettes ont de nombreux avantages (dont la propriété de permettre aux toniens de distinguer correctement les différentes couleurs des palettes catégorielles). [Le site de ce package](#) est très complet et présente de nombreux exemples.

Nous aurons aussi besoin des tableaux créés au Chapitre 3 (masses, masses_se, masses_ci, masses_ci_atp, prox masses)_ci

Donc pensez bien à charger en mémoire les packages et à lancer les commandes de vos scripts qui vous ont permis de créer ces tableaux s'ils ne sont plus en mémoire dans votre environnement de travail. Il existe plusieurs façons de présenter visuellement les positions et les dispersions incertitudes. Concernant les positions et les dispersions, d'abord, nous avons déjà vu plusieurs façons de faire à la Chapitre 3 en particulier dans les parties consacrées aux histogrammes, aux stripcharts et aux boxplots. Nous reprenons ici brièvement chacun de ces 3 types de graphique afin de remettre en contexte avec ce que nous avons appris ici.

Dans un dernier temps, nous verrons comment visualiser l'incertitude liée à des calculs de moyennes ou de variances grâce aux barres d'erreurs ou aux encoches boîtes à moustaches.

8.2 Position et dispersion : les histogrammes

Je vous renvoie à la section 3.1 pour vous rafraîchir la mémoire. Vous pouvez aussi jeter un coup d'œil sur les histogrammes facettés, section 3.9.

Les histogrammes permettent de déterminer à la fois où se trouvent les valeurs les plus fréquemment observées (la position du pic principal correspond à la tendance centrale) et la dispersion (ou variabilité) des valeurs autour de la tendance centrale. Par exemple, l'opérateur `geom_histogram()` permet de faire des histogrammes des températures pour chaque port de New York et chaque mois de l'année 2013 :

```
library(ggplot2)
penguins <- read_csv("penguins.csv")
ggplot(penguins, aes(x = body_mass_g, y = density)) +
  geom_histogram(bins = 20, color = "grey", fill = "white",
                 facet_wrap(~ species, scales = "free_y")) +
  labs(x = "Masse corporelle (g)", y = "Densité") +
  scale_x_continuous(labels = number_format()) +
  theme_bw()
```

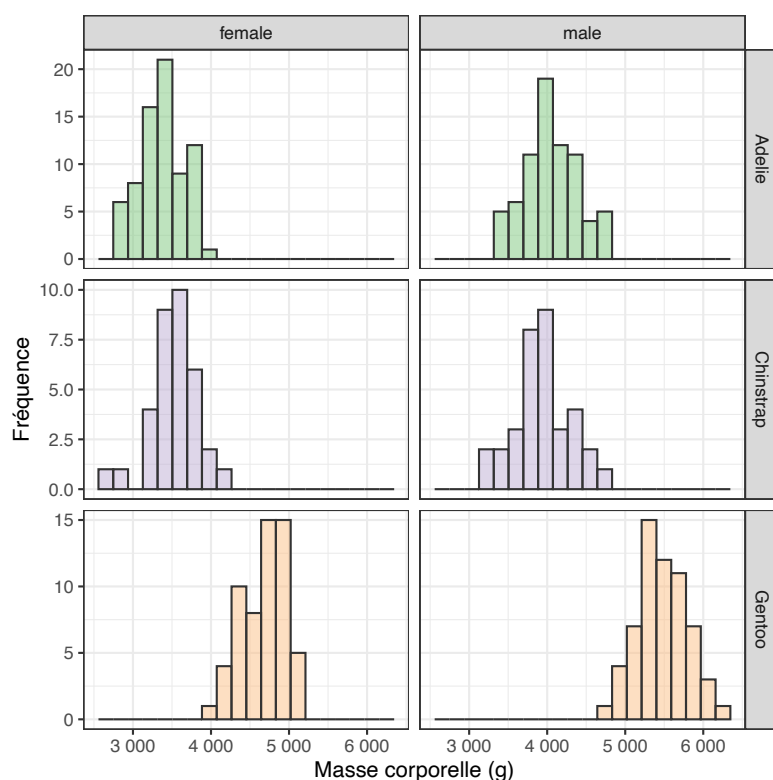



Figure 1: Distribution des masses corporelles chez les mâles et les femelles de 3 espèces de manchots

Ici, 6 histogrammes sont produits. Ils permettent de conclure :

- les masses sont à peu près toutes distribuées selon une courbe en cloche
- les masses moyennes sont plus élevées chez les Gentoo que chez les deux autres espèces. C'est bien la position des pics sur l'axe des abscisses qui nous renseigne dessus.
- pour chaque espèce, les masses moyennes sont globalement plus élevées chez les mâles que chez les femelles. Par exemple, pour l'espèce Adélie, le pic se situe à 3500 grammes pour les femelles, et autour de 4000 grammes pour les mâles.
- la variabilité des masses est comparable pour les 3 espèces et les 2 sexes. Cette fois, c'est l'étalement des histogrammes qui nous renseigne sur la dispersion. L'étalement est toujours d'environ 1500 grammes, s

peut-être pour les femelles Gentoo dont l'étalage
d'environ 1000 grammes.

8.3 Position et dispersion : les stripcharts

Une autre façon de visualiser à la fois les tendances et les dispersions consiste à produire un nuage de points "stripchart". Là encore, je vous renvoie à la partie sur les charts, [3 secondes](#) nous avez besoin de vous rafraîchir la mémoire.

```
penguins
  filter(!is.na(sex))
  ggplot(aes(x = sex, y = body_mass_g, size = length))
  geom_jitter(aes(color = row.leader),
             width = 1, height = 5,
             alpha = 0.5)
  facet_wrap(~species, scales = "free_y")
  labs(x = "Sex", y = "Masse corporelle (g)")
  scale_fill_manual(values = c("black", "red", "blue"))
  scale_y_continuous(labels = number_format())
  theme_bw()
```

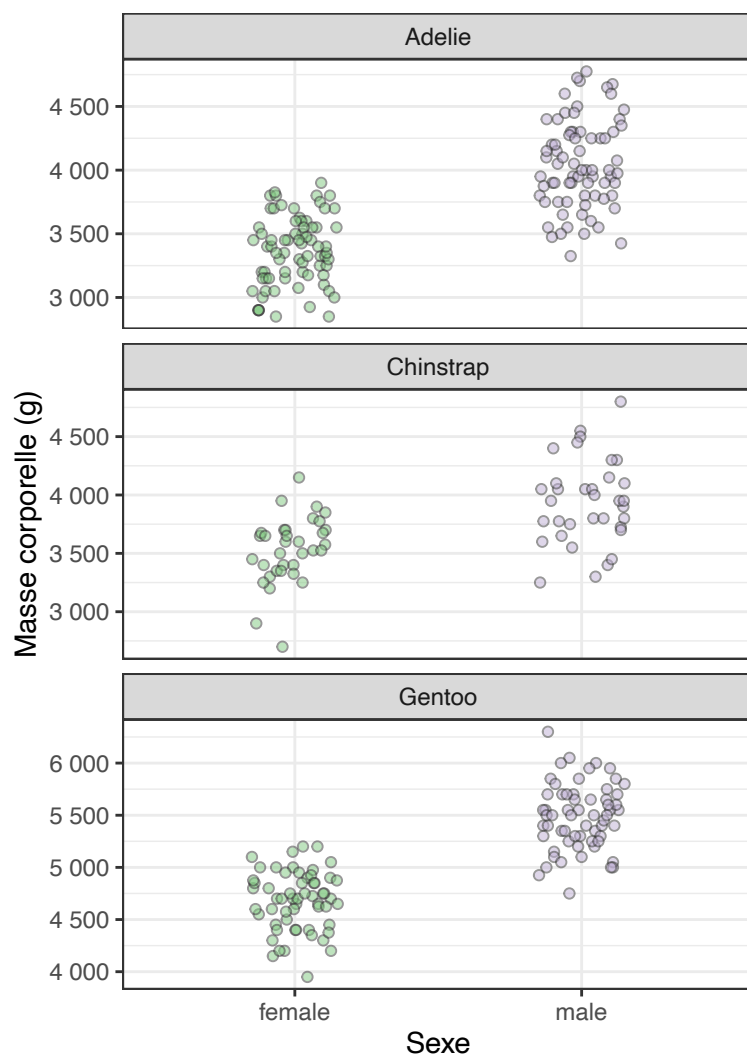


Figure 12: Distribution des masses corporelles chez les mâles et les femelles de 3 espèces de manchots

Cette fois, nous visualisons la totalité des données d'ici et non les données regroupées dans des classes plus ou moins arbitraires. Mais là encore, on peut facilement comparer la composition de chaque série de données : pour chaque espèce, les mâles ont des masses corporelles plus importantes que les femelles, et globalement, les Gentoo ont des masses corporelles plus élevées que les autres espèces. La dispersion des données est aussi facile à comparer entre les groupes : ici l'étendue du nuage de points sur l'axe des ordonnées nous permet de le faire.

Enfin, les stripcharts facettés sont particulièrement utiles lorsque le nombre de séries de données est grand. Par exemple, dans `nyairports`, qui contient des données météo enregistrées toutes les heures de l'année 2013. Si l'on souhaite comparer l'évolution des températures mensuelles dans chacun des 3 aéroports de New York, voilà ce qu'on peut faire :

```
library(ggplot2)
weather %>%
  mutate(temp_cel = (temp_f - 32) / 1.8) %>%
  ggplot(aes(x = factor(month), y = temp_cel, fill = airport)) +
  geom_jitter(aes(color = airport), width = 0.1, height = 0.5, show.legend = FALSE) +
  facet_wrap(~airport, ncol = 1) +
  scale_fill_manual(values = c("Blue", "Red", "Green")) +
  labs(x = "Mois", y = "Températures (°C)") +
  theme_bw()
```

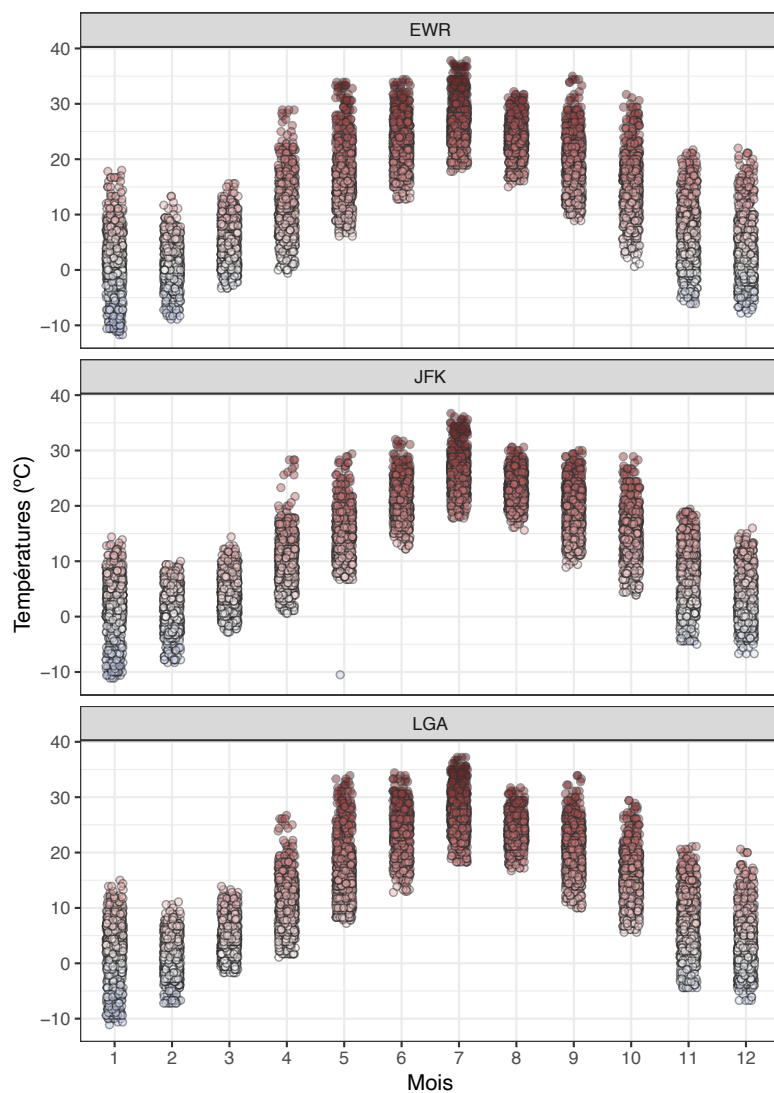


Figure 8.3 : Distribution des températures mensuelles dans les 3 aéroports de New York en 20123

8.4 Position et dispersion : les boxplots

La dernière façon classique de visualiser à la fois les centrales et les dispersions consiste à produire un graphique à boîtes à moustaches, ou "boxplot". Là encore, je vous renvoie à la partie sur les [boxplots](#) avec [besoin](#) de vous rafraîchir la mémoire. Les boîtes à moustaches sont également très pratiques pour comparer de nombreux groupes.

les uns avec les autres. Avec les données de température voilà à quoi ça ressemble :

```
weather %>%
  mutate(temp_cel = (temp - 32) / 1.8)
ggplot(aes(x = factor(month), y = temp_cel)) +
  geom_boxplot(aes(fill = factor(month))) +
  facet_wrap(~, ncol = 1) +
  labs(x = "Mois", y = "Températures (°C)") +
  theme_bw()
```

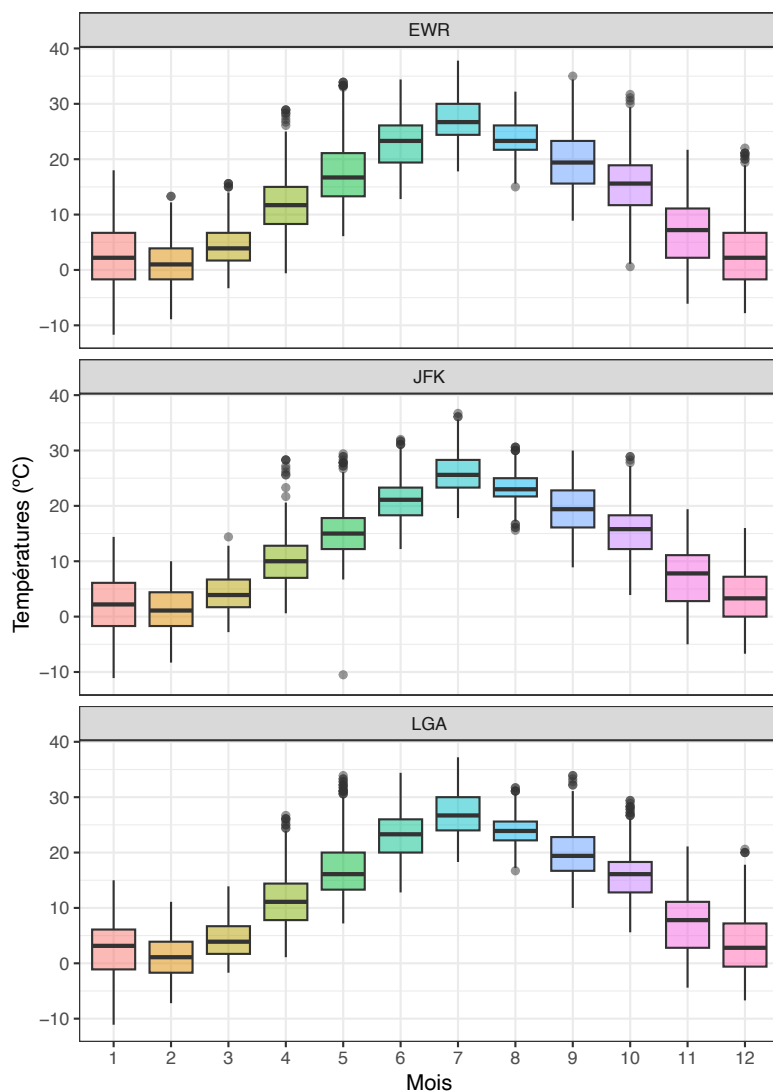


Figure 18: Distribution des températures mensuelles dans les 3 aéroports de New York en 2012

Vous voyez que le code est très proche pour produire un `boxplot` ou un `boxplot`. Comme il est dit au Chapitre 6, les différents éléments de chaque boîte nous renseignent sur la position et sur la dispersion des données pour chaque mois de l'année pour chaque aéroport :

- La limite inférieure de la boîte correspond au premier quartile : 25% des données de l'échantillon sont situées au-dessous de cette valeur.

- La limite supérieure de la boîte correspond au troisième quartile : 75% des données de l'échantillon sont situées au-dessous de cette valeur.
- Le segment épais à l'intérieur de la boîte correspond au second quartile : c'est la médiane de l'échantillon. Cela nous renseigne sur la position de la distribution. Les données de l'échantillon sont situées au-dessus de cette valeur, et 50% au-dessous.
- La hauteur de la boîte correspond à l'étendue (ou intervalle) interquartile ou Inter Quartile Range (IQR) en anglais. On trouve dans cette boîte 50% des observations de l'échantillon. C'est une mesure de la dispersion des 50% des données les plus centrales. Une boîte plus allongée indique donc une plus grande dispersion.
- Les moustaches correspondent à des valeurs qui sont en dessous du premier quartile (pour la moustache du bas) et au-dessus du troisième quartile (pour la moustache du haut). La règle statistique des moustaches s'étendent jusqu'aux valeurs minimales et maximales de l'échantillon, mais elles ne peuvent en aucun cas s'étendre au-delà de 1,5 fois la hauteur de la boîte ($1,5 \text{ fois l'IQR}$) vers le haut et le bas. Si des points apparaissent au-delà des moustaches (vers le haut ou le bas), ces points sont appelés "outliers". On peut observer ici pour plusieurs mois et pour les 3 aéroports (par exemple, en avril dans les 3 aéroports). Ce sont des points qui s'éloignent du centre de la distribution. C'est un fait important puisqu'ils sont au-delà de 1,5 fois l'IQR de part et d'autre du premier ou du troisième quartile. Ils peuvent s'agir d'anomalies de mesures, d'anomalies de saisie des données, ou tout simplement, d'enregistrements tout à fait valides mais atypiques ou extrêmes ; ils ne s'agit donc pas toujours de points aberrants. J'attire votre attention sur le fait que la définition de ces outliers est relativement arbitraire. Nous pourrions choisir d'étendre les moustaches jusqu'à 1,8 fois l'IQR (ou 2, ou 2,5). Nous observerions alors beaucoup moins d'outliers. D'une façon générale, la longueur des moustaches renseigne sur la variabilité des données en dehors de la zone centrale. Plus elles sont longues, plus la variabilité est importante. Très souvent, l'examen des outliers est utile car il nous permet d'en apprendre plus sur le comportement extrême de certaines observations.

tions.

Lorsque les boîtes ont une forme à peu près symétrique de part et d'autre de la médiane (c'est le cas pour l'exemple dans la plupart des catégories), cela signifie qu'un histogramme des mêmes données serait symétrique également.

Les stripcharts et les boxplots sont donc un bon moyen de comparer rapidement la position et la dispersion d'un nombre de séries de données : ici, en quelques lignes de code nous en comparons 12 pour chacun des 3 aéroports de New York.

Les histogrammes sont plus utiles lorsqu'il y a moins de catégories à comparer, comme le boxplot permet de le voir en outre de mieux visualiser les distributions non symétriques ou qui présentent plusieurs pics (distribution bimodales).

8.5 Visualiser l'incertitude : les barres d'erreur

Comme évoqué plus haut, il est important de ne pas confondre les données et les moyennes calculées à partir des données d'un échantillon. Lorsque l'on visualise des moyennes calculées à partir des données d'un échantillon, il est important de faire apparaître des barres d'erreur qui correspondent en général :

- soit à l'erreur standard de la moyenne
- soit à l'intervalle de confiance à 95% de la moyenne

Puisque deux choix sont possibles, il sera important d'indiquer systématiquement dans la légende du graphique, le type de barres représentées. Revenons aux données de masses corporelles des manchots, et commençons par visualiser les masses moyennes avec les erreurs standards. Pour cela nous allons reprendre le code que nous avons écrit précédemment :

```
masses %>% summarise(
  ggplot(aes(x = sex, y = moyenne, lty = "specie")) +
  geom_col(aes(fill = "specie", color = "specie")) +
  geom_errorbar(aes(ymin = moyenne - 1.96 * se, ymax = moyenne + 1.96 * se))
)
```

```

    y_max = y_moyenne + erreur_standard),
    width = 10) 5
facet_wrap(~species, ncol = 3)
labs(x = "Sex", y = "Masse moyenne (g)", title = "Masse moyenne par sexe et espèce")
theme_bw()
scale_fill_manual(values = c("Adelie" = "#4daf4a", "Chinstrap" = "#377eb8", "Gentoo" = "#f4a460"))
scale_y_continuous(labels = function(x) format(x, big.mark = ",", scientific = FALSE))
scale_x_discrete(labels = c("Femelles", "Mâles"))

```

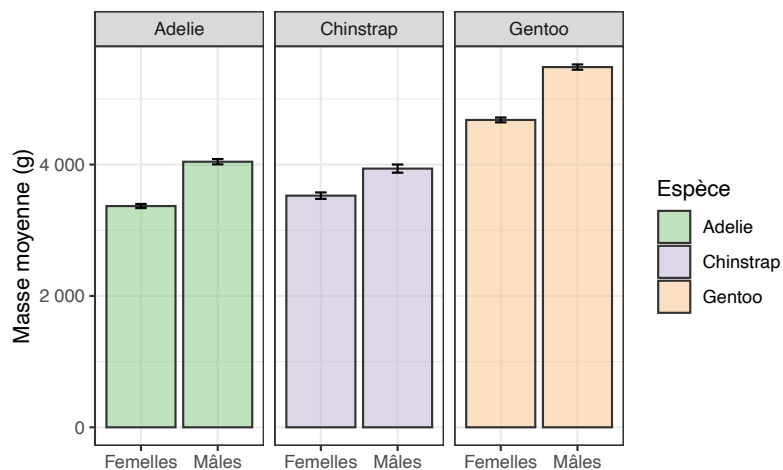


Figure 6.8 : Comparaison des masses moyennes observées chez les mâles et les femelles de 3 espèces de manchots. Les barres d'erreur sont les erreurs standard.

Vous remarquerez que :

1. la fonction `geom_errorbar()` de nouvelles caractéristiques esthétiques qu'il nous faut obtenir : les extrémités inférieures et supérieures des barres d'erreur. Il nous faut donc associer 2 variables à ces caractéristiques esthétiques. Ici, nous nous intéressons à la borne inférieure des barres d'erreur, et la borne supérieure des barres d'erreur. Les variables sont `se_min` et `se_max`. Elles sont dans la partie du tableau `se` du dataset `se` (vous pouvez le vérifier).
2. l'argument `width` de la fonction `geom_errorbar()` permet d'indiquer la longueur des segments horizontaux qui apparaissent à chaque extrémité des barres d'erreur.

Nous pouvons arriver au même résultats en utilisant `masses_selbornes` qui contient des variables différentes :

```
masses_selbornes
ggplot(aes(x = sex, y = moyenne, lsp = species))
geom_col(color = grey20, phi = 5)
geom_errorbar(aes(ymin = moyenne_moins_se,
                  ymax = moyenne_plus_se),
              width = 1)
facet_wrap(species, scales = "y")
labs(x = "Sex", y = "Masse moyenne (g)", title = "Gentoo")
theme_bw()
scale_fill(palette = "dodger")
scale_y_continuous(labels = "Masse (g)")
scale_x_discrete(labels = "Sex")
```

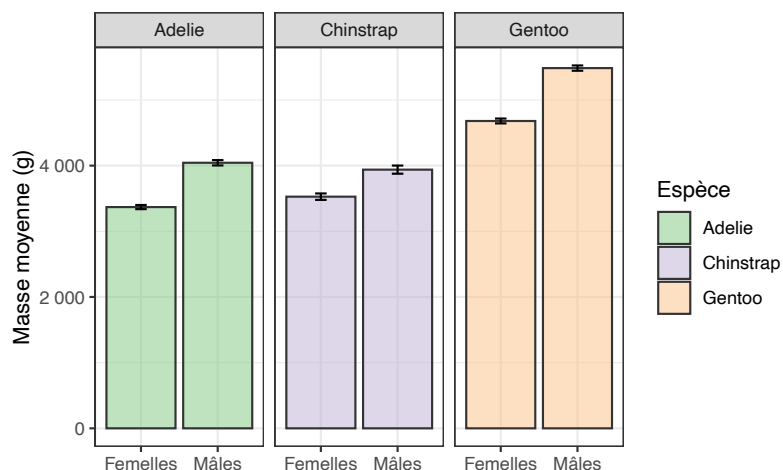


Figure 6.8 : Comparaison des masses moyennes observées chez les mâles et les femelles de 3 espèces de manchots. Les barres d'erreur sont les erreurs standard.

Seule la spécification de `geom_errorbar()` a changé puisque le tableau contient des variables différentes du tableau.

De la même façon, nous pouvons parfaitement faire apparaître, au lieu des erreurs standards, les intervalles de confiance à 95% de chaque masse moyenne. Il nous suffit pour cela d'utiliser `geom_errorbar()` avec les

valeurs de moyennes et des bornes supérieures et inférieures de ces intervalles :

```
masses <- ggplot(aes(x = sexe, y = masse_moyenne, lwr = lower, upr = upper)) +
  geom_bar(aes(fill = sexe), width = 0.5) +
  geom_linerange(aes(lwr = lower, upr = upper), width = 0.5) +
  facet_wrap(~ species) +
  labs(x = "sexe", y = "Masse moyenne (g)", title = "Masse moyenne par espèce") +
  theme_bw() +
  scale_fill_manual(values = c("Femelles" = "#F08080", "Mâles" = "#4682B4")) +
  scale_x_discrete(labels = c("Femelles", "Mâles"))
```

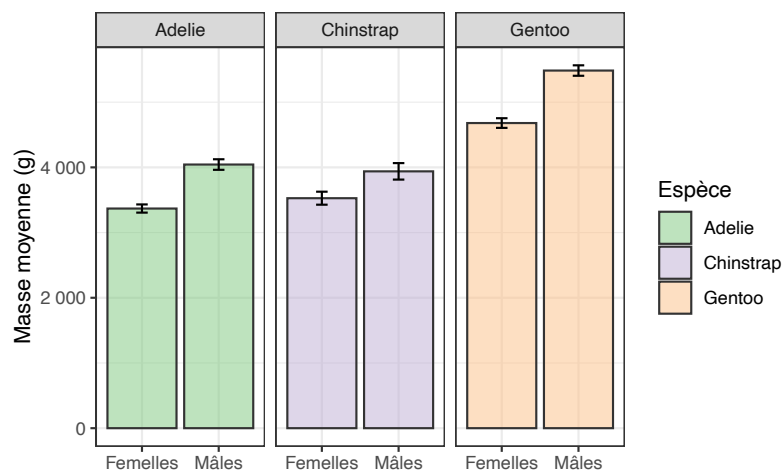


Figure 6.8 : Comparaison des masses moyennes observées chez les mâles et les femelles de 3 espèces de manchots. Les barres d'erreur sont les intervalles de confiance à 95 % des masses moyennes.

Comme vous voyez, les barres d'erreurs sont maintenant longues que sur la Figure 6.7. Cela est normal car rappelez-vous que les intervalles de confiance sont à peu près équivalents à 2 fois les erreurs standards. L'intérêt de représenter les intervalles de confiance est qu'ils sont directement liés aux tests statistiques que nous aborderons dans les chapitres suivants. Globalement, quand 2 séries de données ont des intervalles de confiance qui se chevauchent largement (c'est-à-dire quand les intervalles de confiance se chevauchent largement), on peut conclure qu'il n'y a pas de différence significative entre les deux séries.

par exemple pour les mâles Adélie et Chinstrap), alors le test d'hypothèses conclurait presque toujours à l'absence de différence significative entre les 2 groupes. À l'inverse, si les 2 séries de données ont des intervalles de confiance qui se chevauchent pas du tout (comme les mâles et les femelles Adélie par exemple), alors, un test d'hypothèses conclurait presque toujours à l'existence d'une différence significative entre les 2 groupes. Lorsque les intervalles de confiance des 2 catégories se chevauchent faiblement ou partiellement (comme entre les femelles Adélie et Chinstrap), la situation est moins tranchée, et nous devons nous en remettre aux résultats du test pour savoir si la différence observée peut être considérée comme significative ou non.

8.6 Visualiser l'incertitude : les boîtes à moustaches

Outre les informations de position et de dispersion, les boîtes à moustaches permettent également de visualiser l'incertitude associée aux médianes. Il suffit pour cela d'ajouter l'argument `notch = TRUE` à la fonction `geom_boxplot()` :

```
library(ggplot2)
penguins <- read_csv("data/penguins.csv")
ggplot(penguins, aes(x = sex, y = body_mass_g, fill = species)) +
  geom_boxplot(aes(fill = grey20), notch = TRUE) +
  facet_wrap(~ species) +
  labs(x = "Sex", y = "Masse corporelle (g)", title = "Penguins")
theme_bw()
scale_fill_manual(values = c("black", "white", "red"))
scale_y_continuous(labels = function(x) format(x, unit = "g"))
scale_x_discrete(labels = c("F", "M"))
```

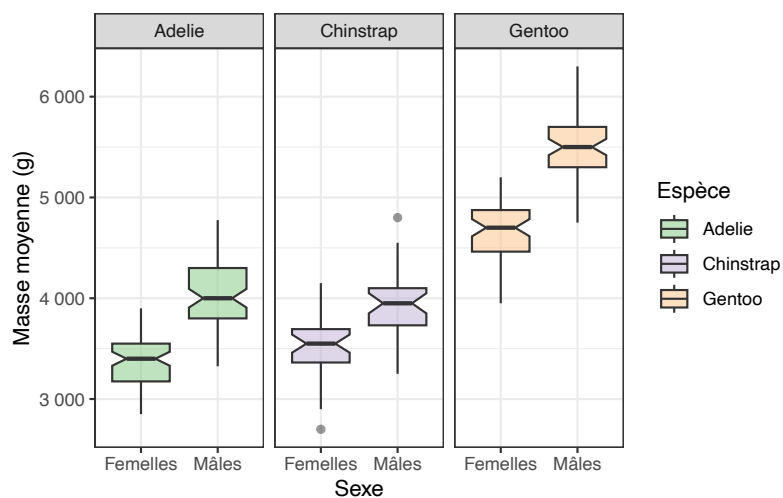


Figure 8.6 : Comparaison des masses corporelles des mâles et femelles de 3 espèces de manchots.

Des encoches ont été ajoutées autour de la médiane de chaque boîte à moustache. Ces encoches sont des encoches d'incertitudes. Les limites inférieures et supérieures des encoches correspondent aux bornes inférieures et supérieures de l'intervalle de confiance à 95% des médianes. Comme pour les moyennes, le chevauchement ou l'absence de chevauchement entre les encoches de 2 séries de données nous renseignent sur l'issue probable des futurs tests statistiques que nous n'aurons pas amenés à réaliser. Notez que tout ce que nous avons dit jusqu'à présent sur le chevauchement des intervalles de confiance des moyennes se retrouve ici pour les intervalles de confiance des médianes (pas de chevauchement entre les encoches des mâles Adélie et Chinstrap, absence de chevauchement entre femelles et mâles Adélie, faible chevauchement entre femelles Adélie et Chinstrap). Il sera donc important d'examiner ces encoches en amont des tests statistiques pour éviter de faire/dire des bêtises...

8.7 Exercice

1. Avec le tableau ci-dessous, calculez les grandeurs suivantes pour chaque espèce de manchot et chaque sexe :
 - la moyenne de la longueur des nageoires

- la variance de la longueur des nageoires
- l'écart-type de la longueur des nageoires
- l'erreur standard de la longueur moyenne des nageoires
- la moyenne de l'épaisseur du bec
- la variance de l'épaisseur du bec
- l'écart-type de l'épaisseur du bec
- l'erreur standard de l'épaisseur du bec

Attention : pensez à retirer les individus dont le sexe est inconnu.

2. Vérifiez avec `sklearn.metrics.mean_squared_error` que les moyennes et écart-types calculés ci-dessus sont corrects.

3. Avec ces données synthétiques faites le graphique suivant :

Moyennes (et erreurs standard) des longueurs de nageoires chez les mâles et les femelles de trois espèces de manchots

