

nycflights13

Pauline Marie PERRIN

2022-05-06

Introduction

Ce document a pour but de rechercher et d'analyser quelques informations issues de la base de données nycflights13. Celle-ci concerne les vols d'un certain nombre de compagnies aériennes en provenance et à destination de plusieurs aéroports localisés à New York City.

Import du jeu de données

```
library(nycflights13)
flights
```

```
## # A tibble: 336,776 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     517             515           2      830           819
## 2  2013     1     1     533             529           4      850           830
## 3  2013     1     1     542             540           2      923           850
## 4  2013     1     1     544             545          -1     1004          1022
## 5  2013     1     1     554             600          -6      812           837
## 6  2013     1     1     554             558          -4      740           728
## 7  2013     1     1     555             600          -5      913           854
## 8  2013     1     1     557             600          -3      709           723
## 9  2013     1     1     557             600          -3      838           846
## 10 2013     1     1     558             600          -2      753           745
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
library(dplyr)
```

```
##
## Attachement du package : 'dplyr'

## Les objets suivants sont masqués depuis 'package:stats':
##
##   filter, lag
```

```
## Les objets suivants sont masqués depuis 'package:base':
##
## intersect, setdiff, setequal, union
```

Recherche d'informations

Dans un premier temps, on identification des variables du jeu de données nycflights13

```
head(flights)
```

```
## # A tibble: 6 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>       <dbl>   <int>         <int>
## 1  2013     1     1     517           515         2     830           819
## 2  2013     1     1     533           529         4     850           830
## 3  2013     1     1     542           540         2     923           850
## 4  2013     1     1     544           545        -1    1004          1022
## 5  2013     1     1     554           600        -6     812           837
## 6  2013     1     1     554           558        -4     740           728
## # ... with 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
## #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
## #   hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
names(flights)
```

```
## [1] "year"           "month"          "day"            "dep_time"
## [5] "sched_dep_time" "dep_delay"      "arr_time"       "sched_arr_time"
## [9] "arr_delay"      "carrier"        "flight"         "tailnum"
## [13] "origin"         "dest"           "air_time"       "distance"
## [17] "hour"           "minute"         "time_hour"
```

Ensuite, on décrit le contexte statistique

```
dim(flights)
```

```
## [1] 336776    19
```

```
ncol(flights)
```

```
## [1] 19
```

```
nrow(flights)
```

```
## [1] 336776
```

Nous pouvons ensuite mettre en évidence un certain nombre d'informations :

1. Sélection des vols à destination de Houston (IAH or HOU)

```
flights %>%
  filter(dest=="IAH" | dest=="HOU")
```

```
## # A tibble: 9,313 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     517           515           2     830           819
## 2  2013     1     1     533           529           4     850           830
## 3  2013     1     1     623           627          -4     933           932
## 4  2013     1     1     728           732          -4    1041          1038
## 5  2013     1     1     739           739           0    1104          1038
## 6  2013     1     1     908           908           0    1228          1219
## 7  2013     1     1    1028          1026           2    1350          1339
## 8  2013     1     1    1044          1045          -1    1352          1351
## 9  2013     1     1    1114           900        134    1447          1222
## 10 2013     1     1    1205          1200           5    1503          1505
## # ... with 9,303 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

2. Sélection des vols arrivés avec un retard de deux heures ou plus

```
flights %>%
  filter(flights$arr_delay >= 120)
```

```
## # A tibble: 10,200 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
## 1  2013     1     1     811           630        101    1047           830
## 2  2013     1     1     848          1835        853    1001          1950
## 3  2013     1     1     957           733        144    1056           853
## 4  2013     1     1    1114           900        134    1447          1222
## 5  2013     1     1    1505          1310        115    1638          1431
## 6  2013     1     1    1525          1340        105    1831          1626
## 7  2013     1     1    1549          1445         64    1912          1656
## 8  2013     1     1    1558          1359        119    1718          1515
## 9  2013     1     1    1732          1630         62    2028          1825
## 10 2013     1     1    1803          1620        103    2008          1750
## # ... with 10,190 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

3. Sélection des vols réalisés par les compagnies United, American, ou Delta

```
flights %>%
  filter(flights$carrier %in% c("AA", "UA", "DL"))
```

```
## # A tibble: 139,504 x 19
```

```
##      year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##      <int> <int> <int>   <int>         <int>      <dbl>   <int>         <int>
##  1  2013     1     1     517           515         2     830           819
##  2  2013     1     1     533           529         4     850           830
##  3  2013     1     1     542           540         2     923           850
##  4  2013     1     1     554           600        -6     812           837
##  5  2013     1     1     554           558        -4     740           728
##  6  2013     1     1     558           600        -2     753           745
##  7  2013     1     1     558           600        -2     924           917
##  8  2013     1     1     558           600        -2     923           937
##  9  2013     1     1     559           600        -1     941           910
## 10  2013     1     1     559           600        -1     854           902
## # ... with 139,494 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

4. Sélection des vols réalisés en été (Juillet, Août et Septembre)

```
flights %>%
  filter(flights$month %in% c(7,8,9))
```

```
## # A tibble: 86,326 x 19
##      year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##      <int> <int> <int>   <int>         <int>      <dbl>   <int>         <int>
##  1  2013     7     1     1           2029        212     236           2359
##  2  2013     7     1     2           2359         3     344           344
##  3  2013     7     1    29           2245        104     151             1
##  4  2013     7     1    43           2130        193     322            14
##  5  2013     7     1    44           2150        174     300            100
##  6  2013     7     1    46           2051        235     304           2358
##  7  2013     7     1    48           2001        287     308           2305
##  8  2013     7     1    58           2155        183     335             43
##  9  2013     7     1   100           2146        194     327             30
## 10  2013     7     1   100           2245        135     337            135
## # ... with 86,316 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

5. Sélection des vols arrivés avec plus de deux heures de retard mais qui sont partis à l'heure

```
flights %>%
  filter((flights$arr_delay > 120) & (flights$dep_delay <= 0))
```

```
## # A tibble: 29 x 19
##      year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##      <int> <int> <int>   <int>         <int>      <dbl>   <int>         <int>
##  1  2013     1    27   1419           1420        -1     1754           1550
##  2  2013    10     7   1350           1350         0     1736           1526
```

```
## 3 2013 10 7 1357 1359 -2 1858 1654
## 4 2013 10 16 657 700 -3 1258 1056
## 5 2013 11 1 658 700 -2 1329 1015
## 6 2013 3 18 1844 1847 -3 39 2219
## 7 2013 4 17 1635 1640 -5 2049 1845
## 8 2013 4 18 558 600 -2 1149 850
## 9 2013 4 18 655 700 -5 1213 950
## 10 2013 5 22 1827 1830 -3 2217 2010
## # ... with 19 more rows, and 11 more variables: arr_delay <dbl>, carrier <chr>,
## # flight <int>, tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
## # distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

6. Sélection des vols partis entre minuit et 6 heures du matin inclus

```
flights %>%
  filter(flights$hour <= 6 | flights$hour == 24)
```

```
## # A tibble: 27,905 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>      <dbl>   <int>         <int>
## 1 2013     1     1     517           515         2     830           819
## 2 2013     1     1     533           529         4     850           830
## 3 2013     1     1     542           540         2     923           850
## 4 2013     1     1     544           545        -1    1004          1022
## 5 2013     1     1     554           600        -6     812           837
## 6 2013     1     1     554           558        -4     740           728
## 7 2013     1     1     555           600        -5     913           854
## 8 2013     1     1     557           600        -3     709           723
## 9 2013     1     1     557           600        -3     838           846
## 10 2013     1     1     558           600        -2     753           745
## # ... with 27,895 more rows, and 11 more variables: arr_delay <dbl>,
## # carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## # air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

7. Sélection des vols American ayant un retard de 2h ou plus au décollage

```
flights %>%
  filter(carrier=="AA") %>%
  filter(dep_delay>120)
```

```
## # A tibble: 720 x 19
##   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>   <int>         <int>      <dbl>   <int>         <int>
## 1 2013     1     1    1856           1645       131    2212          2005
## 2 2013     1     1    2205           1720       285     46          2040
## 3 2013     1     2    1607           1030       337    2003          1355
## 4 2013     1     2    1751           1450       181    2041          1755
## 5 2013     1     3     854           630       144    1057           810
## 6 2013     1     3     909           700       129    1103           850
```

```
## 7 2013 1 3 1758 1550 128 2240 2050
## 8 2013 1 3 1821 1530 171 2131 1910
## 9 2013 1 4 1305 1030 155 1452 1210
## 10 2013 1 4 1917 1700 137 2135 1950
## # ... with 710 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>
```

8. Sélection des variables textuelles dans le tableau

```
flights %>%
  select(where(is.character))
```

```
## # A tibble: 336,776 x 4
##   carrier tailnum origin dest
##   <chr>    <chr>    <chr> <chr>
## 1 UA      N14228 EWR   IAH
## 2 UA      N24211 LGA   IAH
## 3 AA      N619AA  JFK   MIA
## 4 B6      N804JB  JFK   BQN
## 5 DL      N668DN  LGA   ATL
## 6 UA      N39463  EWR   ORD
## 7 B6      N516JB  EWR   FLL
## 8 EV      N829AS  LGA   IAD
## 9 B6      N593JB  JFK   MCO
## 10 AA     N3ALAA  LGA   ORD
## # ... with 336,766 more rows
```

9. Création d'une nouvelle colonne avec la durée de vol en heures à partir d'une variable existante

On a pu remarquer que les durées de vol étaient données en minutes. Pour mettre en forme nos données, il peut être intéressant de créer une nouvelle colonne à partir d'une variable existante avec ces durées de vol en heures :

```
flights %>%
  mutate(duree = air_time/60) %>%
  select(flight, duree, air_time) %>%
  arrange(air_time)
```

```
## # A tibble: 336,776 x 3
##   flight duree air_time
##   <int> <dbl>    <dbl>
## 1 4368 0.333      20
## 2 4631 0.333      20
## 3 4276 0.35      21
## 4 4619 0.35      21
## 5 4368 0.35      21
## 6 4619 0.35      21
## 7 2132 0.35      21
```

```
## 8 3650 0.35 21
## 9 4118 0.35 21
## 10 4276 0.35 21
## # ... with 336,766 more rows
```

Création de graphiques

Pour créer nos graphiques, nous avons recours au package `ggplot`. Il s'agit d'une extension de `tidyverse` qui permet de concevoir des graphiques plus attractifs et plus complexes.

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

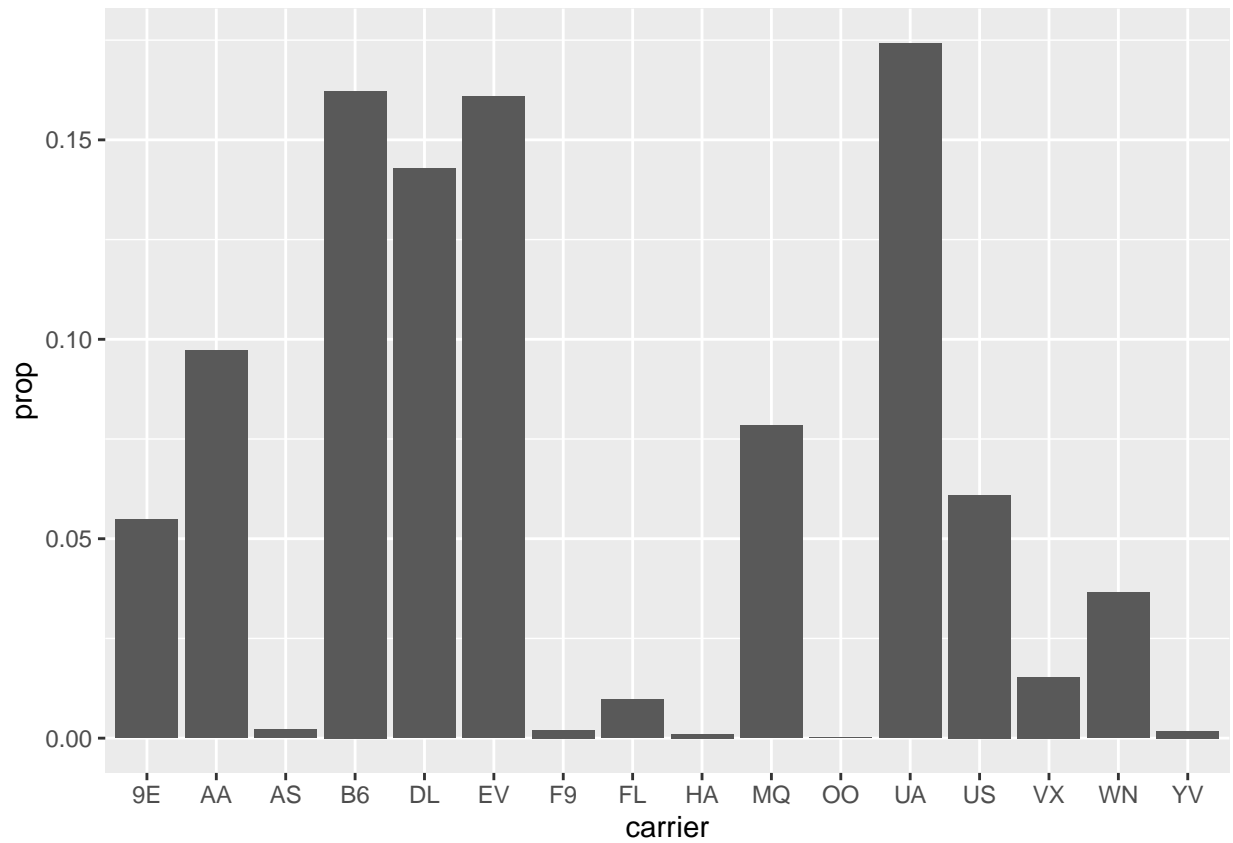
```
## v ggplot2 3.3.6    v purrr 0.3.4
## v tibble 3.1.7     v stringr 1.4.0
## v tidyr 1.2.0      v forcats 0.5.1
## v readr 2.1.2
```

```
## -- Conflicts ----- tidyverse_conflicts() --
```

```
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

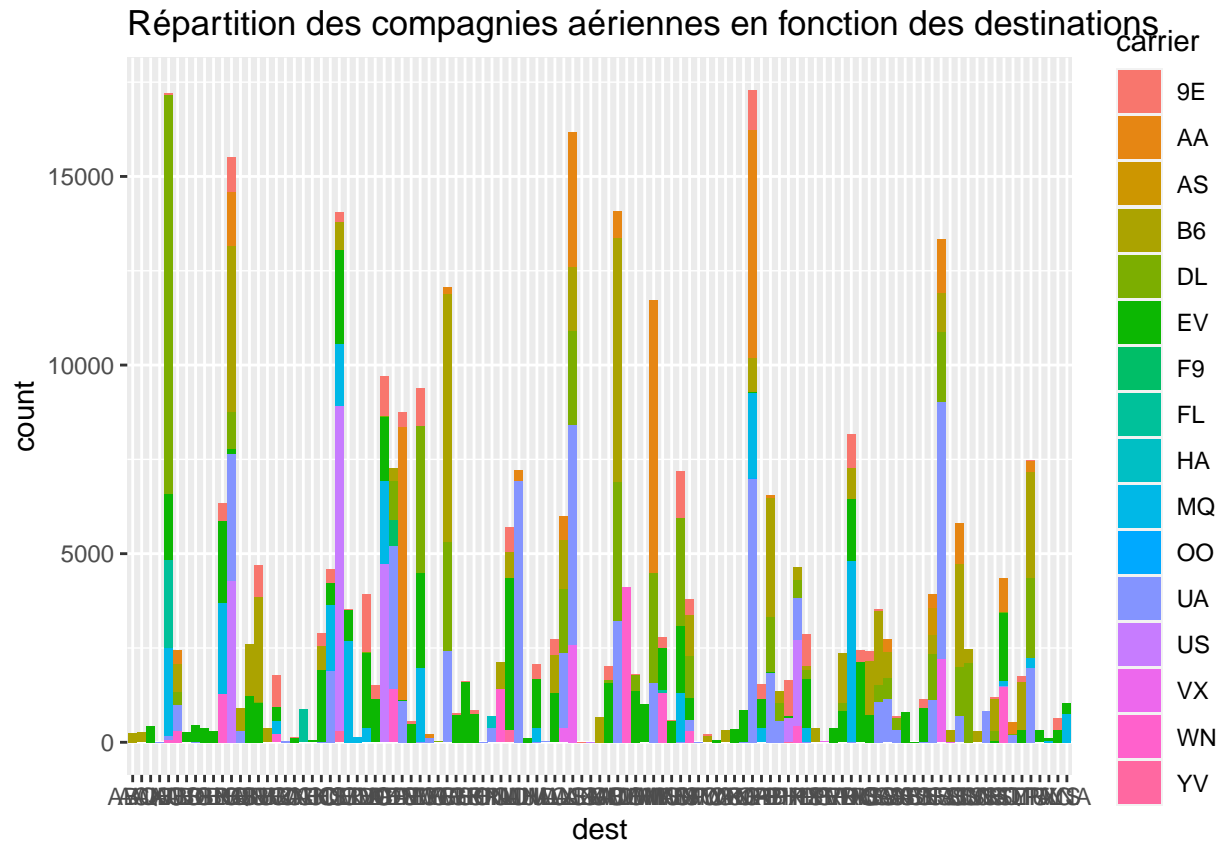
1. Graphique des proportions de vols réalisés par chacune des compagnies aériennes étudiées

```
ggplot(data=flights) +
  geom_bar(mapping = aes(x = carrier, y = stat(prop), group = 1))
```



2. Répartition des compagnies aériennes en fonction des destinations

```
ggplot(data=flights) +  
  geom_bar(mapping = aes (x = dest, fill = carrier)) +  
  ggtitle("Répartition des compagnies aériennes en fonction des destinations")
```

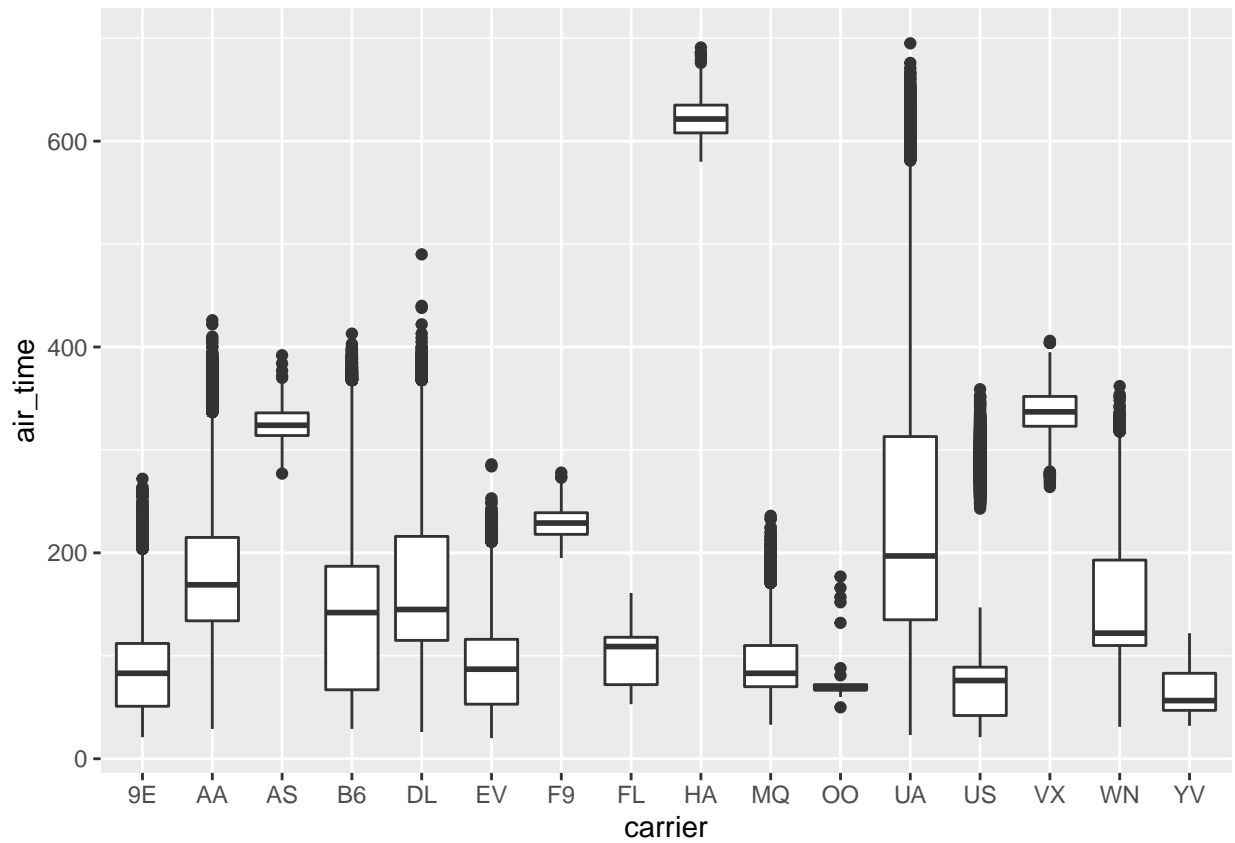



3. Répartition des retards au départ en fonction des compagnies aériennes

Il est également possible de réaliser des barplots :

```
flights%>%
  select(carrier, air_time) %>%
  arrange(carrier) %>%
  ggplot() +
  geom_boxplot(aes(x=carrier, y=air_time))
```

Warning: Removed 9430 rows containing non-finite values (stat_boxplot).



Par cet exemple graphique, il est possible d'observer la répartition des retards au décollage en fonction des compagnies aériennes.