

# Change data

## Step 1: Load packages

```
install.packages("tidyverse")
```

```
## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.0'  
## (as 'lib' is unspecified)
```

```
install.packages("skimr")
```

```
## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.0'  
## (as 'lib' is unspecified)
```

```
install.packages("janitor")
```

```
## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.0'  
## (as 'lib' is unspecified)
```

Once a package is installed, you can load it by running the `library()` function with the package name inside the parentheses:

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.2 --v ggplot2 3.3.5  
## v tibble 3.1.8      v dplyr 1.1.0  
## v tidyr 1.3.0      v stringr 1.5.0  
## v readr 2.1.3      v forcats 0.5.1-- Conflicts ----- tidyverse  
## x dplyr::filter() masks stats::filter()  
## x dplyr::lag()     masks stats::lag()
```

```
library(skimr)  
library(janitor)
```

```
##  
## Attaching package: 'janitor'  
##  
## The following objects are masked from 'package:stats':  
##  
##   chisq.test, fisher.test
```

## Step 2: Import data

```
hotel_bookings <- read_csv("hotel_bookings.csv")
```

```
## Rows: 119390 Columns: 32-- Column specification -----  
## Delimiter: ","  
## chr  (13): hotel, arrival_date_month, meal, country, market_segment, distrib...  
## dbl  (18): is_canceled, lead_time, arrival_date_year, arrival_date_week_numb...  
## date  (1): reservation_status_date  
## i Use `spec()` to retrieve the full column specification for this data.  
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

### Step 3: Check data

```
head(hotel_bookings)
```

```
## # A tibble: 6 x 32
##   hotel    is_ca~1 lead_~2 arriv~3 arriv~4 arriv~5 arriv~6 stays~7 stays~8 adults
##   <chr>      <dbl>   <dbl>   <dbl> <chr>      <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Resort~      0     342    2015 July      27       1       0       0       2
## 2 Resort~      0     737    2015 July      27       1       0       0       2
## 3 Resort~      0       7    2015 July      27       1       0       1       1
## 4 Resort~      0      13    2015 July      27       1       0       1       1
## 5 Resort~      0      14    2015 July      27       1       0       2       2
## 6 Resort~      0      14    2015 July      27       1       0       2       2
## # ... with 22 more variables: children <dbl>, babies <dbl>, meal <chr>,
## #   country <chr>, market_segment <chr>, distribution_channel <chr>,
## #   is_repeated_guest <dbl>, previous_cancellations <dbl>,
## #   previous_bookings_not_canceled <dbl>, reserved_room_type <chr>,
## #   assigned_room_type <chr>, booking_changes <dbl>, deposit_type <chr>,
## #   agent <chr>, company <chr>, days_in_waiting_list <dbl>,
## #   customer_type <chr>, adr <dbl>, required_car_parking_spaces <dbl>, ...
```

```
str(hotel_bookings)
```

```
## spc_tbl_ [119,390 x 32] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ hotel                : chr [1:119390] "Resort Hotel" "Resort Hotel" "Resort Hotel" "Resort Hotel" ...
## $ is_canceled           : num [1:119390] 0 0 0 0 0 0 0 0 0 1 1 ...
## $ lead_time             : num [1:119390] 342 737 7 13 14 14 0 9 85 75 ...
## $ arrival_date_year     : num [1:119390] 2015 2015 2015 2015 2015 ...
## $ arrival_date_month    : chr [1:119390] "July" "July" "July" "July" ...
## $ arrival_date_week_number : num [1:119390] 27 27 27 27 27 27 27 27 27 27 ...
## $ arrival_date_day_of_month : num [1:119390] 1 1 1 1 1 1 1 1 1 1 ...
## $ stays_in_weekend_nights : num [1:119390] 0 0 0 0 0 0 0 0 0 0 ...
## $ stays_in_week_nights   : num [1:119390] 0 0 1 1 2 2 2 2 3 3 ...
## $ adults                : num [1:119390] 2 2 1 1 2 2 2 2 2 2 ...
## $ children              : num [1:119390] 0 0 0 0 0 0 0 0 0 0 ...
## $ babies                : num [1:119390] 0 0 0 0 0 0 0 0 0 0 ...
## $ meal                  : chr [1:119390] "BB" "BB" "BB" "BB" ...
## $ country                : chr [1:119390] "PRT" "PRT" "GBR" "GBR" ...
## $ market_segment        : chr [1:119390] "Direct" "Direct" "Direct" "Corporate" ...
## $ distribution_channel   : chr [1:119390] "Direct" "Direct" "Direct" "Corporate" ...
## $ is_repeated_guest      : num [1:119390] 0 0 0 0 0 0 0 0 0 0 ...
## $ previous_cancellations : num [1:119390] 0 0 0 0 0 0 0 0 0 0 ...
## $ previous_bookings_not_canceled: num [1:119390] 0 0 0 0 0 0 0 0 0 0 ...
## $ reserved_room_type     : chr [1:119390] "C" "C" "A" "A" ...
## $ assigned_room_type     : chr [1:119390] "C" "C" "C" "A" ...
## $ booking_changes        : num [1:119390] 3 4 0 0 0 0 0 0 0 0 ...
## $ deposit_type           : chr [1:119390] "No Deposit" "No Deposit" "No Deposit" "No Deposit" ...
## $ agent                  : chr [1:119390] "NULL" "NULL" "NULL" "304" ...
## $ company                : chr [1:119390] "NULL" "NULL" "NULL" "NULL" ...
## $ days_in_waiting_list   : num [1:119390] 0 0 0 0 0 0 0 0 0 0 ...
## $ customer_type          : chr [1:119390] "Transient" "Transient" "Transient" "Transient" ...
## $ adr                    : num [1:119390] 0 0 75 75 98 ...
## $ required_car_parking_spaces : num [1:119390] 0 0 0 0 0 0 0 0 0 0 ...
## $ total_of_special_requests : num [1:119390] 0 0 0 0 1 1 0 1 1 0 ...
## $ reservation_status     : chr [1:119390] "Check-Out" "Check-Out" "Check-Out" "Check-Out" ...
```

```
## $ reservation_status_date      : Date[1:119390], format: "2015-07-01" "2015-07-01" ...
## - attr(*, "spec")=
## .. cols(
## ..   hotel = col_character(),
## ..   is_canceled = col_double(),
## ..   lead_time = col_double(),
## ..   arrival_date_year = col_double(),
## ..   arrival_date_month = col_character(),
## ..   arrival_date_week_number = col_double(),
## ..   arrival_date_day_of_month = col_double(),
## ..   stays_in_weekend_nights = col_double(),
## ..   stays_in_week_nights = col_double(),
## ..   adults = col_double(),
## ..   children = col_double(),
## ..   babies = col_double(),
## ..   meal = col_character(),
## ..   country = col_character(),
## ..   market_segment = col_character(),
## ..   distribution_channel = col_character(),
## ..   is_repeated_guest = col_double(),
## ..   previous_cancellations = col_double(),
## ..   previous_bookings_not_canceled = col_double(),
## ..   reserved_room_type = col_character(),
## ..   assigned_room_type = col_character(),
## ..   booking_changes = col_double(),
## ..   deposit_type = col_character(),
## ..   agent = col_character(),
## ..   company = col_character(),
## ..   days_in_waiting_list = col_double(),
## ..   customer_type = col_character(),
## ..   adr = col_double(),
## ..   required_car_parking_spaces = col_double(),
## ..   total_of_special_requests = col_double(),
## ..   reservation_status = col_character(),
## ..   reservation_status_date = col_date(format = "")
## .. )
## - attr(*, "problems")=<externalptr>
```

```
glimpse(hotel_bookings)
```

```
## Rows: 119,390
## Columns: 32
## $ hotel      <chr> "Resort Hotel", "Resort Hotel", "Resort~
## $ is_canceled <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, ~
## $ lead_time  <dbl> 342, 737, 7, 13, 14, 14, 0, 9, 85, 75, ~
## $ arrival_date_year <dbl> 2015, 2015, 2015, 2015, 2015, 2015, 201~
## $ arrival_date_month <chr> "July", "July", "July", "July", "July",~
## $ arrival_date_week_number <dbl> 27, 27, 27, 27, 27, 27, 27, 27, 27, 27,~
## $ arrival_date_day_of_month <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ stays_in_weekend_nights <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ stays_in_week_nights <dbl> 0, 0, 1, 1, 2, 2, 2, 2, 3, 3, 4, 4, ~
## $ adults     <dbl> 2, 2, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, ~
## $ children   <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ babies     <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ meal       <chr> "BB", "BB", "BB", "BB", "BB", "BB", "BB~
```

```
## $ country          <chr> "PRT", "PRT", "GBR", "GBR", "GBR", "GBR~
## $ market_segment  <chr> "Direct", "Direct", "Direct", "Corporat~
## $ distribution_channel <chr> "Direct", "Direct", "Direct", "Corporat~
## $ is_repeated_guest <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ previous_cancellations <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ previous_bookings_not_canceled <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ reserved_room_type <chr> "C", "C", "A", "A", "A", "A", "C", "C", ~
## $ assigned_room_type <chr> "C", "C", "C", "A", "A", "A", "C", "C", ~
## $ booking_changes <dbl> 3, 4, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ deposit_type <chr> "No Deposit", "No Deposit", "No Deposit~
## $ agent <chr> "NULL", "NULL", "NULL", "304", "240", "~
## $ company <chr> "NULL", "NULL", "NULL", "NULL", "NULL", ~
## $ days_in_waiting_list <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ customer_type <chr> "Transient", "Transient", "Transient", ~
## $ adr <dbl> 0.00, 0.00, 75.00, 75.00, 98.00, 98.00, ~
## $ required_car_parking_spaces <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ total_of_special_requests <dbl> 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 3, ~
## $ reservation_status <chr> "Check-Out", "Check-Out", "Check-Out", ~
## $ reservation_status_date <date> 2015-07-01, 2015-07-01, 2015-07-02, 20~
```

```
colnames(hotel_bookings)
```

```
## [1] "hotel" "is_canceled"
## [3] "lead_time" "arrival_date_year"
## [5] "arrival_date_month" "arrival_date_week_number"
## [7] "arrival_date_day_of_month" "stays_in_weekend_nights"
## [9] "stays_in_week_nights" "adults"
## [11] "children" "babies"
## [13] "meal" "country"
## [15] "market_segment" "distribution_channel"
## [17] "is_repeated_guest" "previous_cancellations"
## [19] "previous_bookings_not_canceled" "reserved_room_type"
## [21] "assigned_room_type" "booking_changes"
## [23] "deposit_type" "agent"
## [25] "company" "days_in_waiting_list"
## [27] "customer_type" "adr"
## [29] "required_car_parking_spaces" "total_of_special_requests"
## [31] "reservation_status" "reservation_status_date"
```

## Manipulating data

```
arrange(hotel_bookings, lead_time)
```

```
## # A tibble: 119,390 x 32
##   hotel is_ca~1 lead_~2 arriv~3 arriv~4 arriv~5 arriv~6 stays~7 stays~8 adults
##   <chr>   <dbl>   <dbl>   <dbl> <chr>       <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Resor~     0     0   2015 July       27     1     0     2     2
## 2 Resor~     0     0   2015 July       27     1     0     1     2
## 3 Resor~     0     0   2015 July       27     2     0     1     2
## 4 Resor~     0     0   2015 July       27     2     0     1     2
## 5 Resor~     0     0   2015 July       27     2     0     1     2
## 6 Resor~     0     0   2015 July       28     5     1     0     2
## 7 Resor~     0     0   2015 July       28     6     0     0     1
## 8 Resor~     0     0   2015 July       28     7     0     1     1
## 9 Resor~     0     0   2015 July       28     7     0     1     3
```

```
## 10 Resor~      0      0    2015 July      28      7      0      1      1
## # ... with 119,380 more rows, 22 more variables: children <dbl>, babies <dbl>,
## #   meal <chr>, country <chr>, market_segment <chr>,
## #   distribution_channel <chr>, is_repeated_guest <dbl>,
## #   previous_cancellations <dbl>, previous_bookings_not_canceled <dbl>,
## #   reserved_room_type <chr>, assigned_room_type <chr>, booking_changes <dbl>,
## #   deposit_type <chr>, agent <chr>, company <chr>, days_in_waiting_list <dbl>,
## #   customer_type <chr>, adr <dbl>, required_car_parking_spaces <dbl>, ...
arrange(hotel_bookings, desc(lead_time))
```

```
## # A tibble: 119,390 x 32
##   hotel is_ca~1 lead_~2 arriv~3 arriv~4 arriv~5 arriv~6 stays~7 stays~8 adults
##   <chr>   <dbl>   <dbl>   <dbl> <chr>       <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Resor~      0     737    2015 July      27       1       0       0       2
## 2 Resor~      0     709    2016 Februa~      9     25       8     20       2
## 3 City ~      1     629    2017 March      13     30       0       1       1
## 4 City ~      1     629    2017 March      13     30       0       1       1
## 5 City ~      1     629    2017 March      13     30       0       2       2
## 6 City ~      1     629    2017 March      13     30       0       2       2
## 7 City ~      1     629    2017 March      13     30       0       2       2
## 8 City ~      1     629    2017 March      13     30       0       2       2
## 9 City ~      1     629    2017 March      13     30       0       2       2
## 10 City ~     1     629    2017 March      13     30       0       2       2
## # ... with 119,380 more rows, 22 more variables: children <dbl>, babies <dbl>,
## #   meal <chr>, country <chr>, market_segment <chr>,
## #   distribution_channel <chr>, is_repeated_guest <dbl>,
## #   previous_cancellations <dbl>, previous_bookings_not_canceled <dbl>,
## #   reserved_room_type <chr>, assigned_room_type <chr>, booking_changes <dbl>,
## #   deposit_type <chr>, agent <chr>, company <chr>, days_in_waiting_list <dbl>,
## #   customer_type <chr>, adr <dbl>, required_car_parking_spaces <dbl>, ...
```

```
head(hotel_bookings)
```

```
## # A tibble: 6 x 32
##   hotel is_ca~1 lead_~2 arriv~3 arriv~4 arriv~5 arriv~6 stays~7 stays~8 adults
##   <chr>   <dbl>   <dbl>   <dbl> <chr>       <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 Resort~      0     342    2015 July      27       1       0       0       2
## 2 Resort~      0     737    2015 July      27       1       0       0       2
## 3 Resort~      0       7    2015 July      27       1       0       1       1
## 4 Resort~      0     13    2015 July      27       1       0       1       1
## 5 Resort~      0     14    2015 July      27       1       0       2       2
## 6 Resort~      0     14    2015 July      27       1       0       2       2
## # ... with 22 more variables: children <dbl>, babies <dbl>, meal <chr>,
## #   country <chr>, market_segment <chr>, distribution_channel <chr>,
## #   is_repeated_guest <dbl>, previous_cancellations <dbl>,
## #   previous_bookings_not_canceled <dbl>, reserved_room_type <chr>,
## #   assigned_room_type <chr>, booking_changes <dbl>, deposit_type <chr>,
## #   agent <chr>, company <chr>, days_in_waiting_list <dbl>,
## #   customer_type <chr>, adr <dbl>, required_car_parking_spaces <dbl>, ...
```

```
hotel_bookings_v2 <-
  arrange(hotel_bookings, desc(lead_time))
```

```
head(hotel_bookings_v2)
```

```
## # A tibble: 6 x 32
##   hotel    is_ca~1 lead_~2 arriv~3 arriv~4 arriv~5 arriv~6 stays~7 stays~8 adults
##   <chr>      <dbl>  <dbl>  <dbl> <chr>      <dbl>  <dbl>  <dbl>  <dbl>  <dbl>
## 1 Resort~    0    737   2015 July      27      1      0      0      2
## 2 Resort~    0    709   2016 Februa~  9     25      8     20      2
## 3 City H~    1    629   2017 March    13     30      0      1      1
## 4 City H~    1    629   2017 March    13     30      0      1      1
## 5 City H~    1    629   2017 March    13     30      0      2      2
## 6 City H~    1    629   2017 March    13     30      0      2      2
## # ... with 22 more variables: children <dbl>, babies <dbl>, meal <chr>,
## #   country <chr>, market_segment <chr>, distribution_channel <chr>,
## #   is_repeated_guest <dbl>, previous_cancellations <dbl>,
## #   previous_bookings_not_canceled <dbl>, reserved_room_type <chr>,
## #   assigned_room_type <chr>, booking_changes <dbl>, deposit_type <chr>,
## #   agent <chr>, company <chr>, days_in_waiting_list <dbl>,
## #   customer_type <chr>, adr <dbl>, required_car_parking_spaces <dbl>, ...
```

```
max(hotel_bookings$lead_time)
```

```
## [1] 737
```

```
min(hotel_bookings$lead_time)
```

```
## [1] 0
```

```
mean(hotel_bookings$lead_time)
```

```
## [1] 104.0114
```

```
mean(hotel_bookings_v2$lead_time)
```

```
## [1] 104.0114
```

```
hotel_bookings_city <-
  filter(hotel_bookings, hotel_bookings$hotel=="City Hotel")
```

```
head(hotel_bookings_city)
```

```
## # A tibble: 6 x 32
##   hotel    is_ca~1 lead_~2 arriv~3 arriv~4 arriv~5 arriv~6 stays~7 stays~8 adults
##   <chr>      <dbl>  <dbl>  <dbl> <chr>      <dbl>  <dbl>  <dbl>  <dbl>  <dbl>
## 1 City H~    0      6   2015 July      27      1      0      2      1
## 2 City H~    1     88   2015 July      27      1      0      4      2
## 3 City H~    1     65   2015 July      27      1      0      4      1
## 4 City H~    1     92   2015 July      27      1      2      4      2
## 5 City H~    1    100   2015 July      27      2      0      2      2
## 6 City H~    1     79   2015 July      27      2      0      3      2
## # ... with 22 more variables: children <dbl>, babies <dbl>, meal <chr>,
## #   country <chr>, market_segment <chr>, distribution_channel <chr>,
## #   is_repeated_guest <dbl>, previous_cancellations <dbl>,
## #   previous_bookings_not_canceled <dbl>, reserved_room_type <chr>,
## #   assigned_room_type <chr>, booking_changes <dbl>, deposit_type <chr>,
## #   agent <chr>, company <chr>, days_in_waiting_list <dbl>,
## #   customer_type <chr>, adr <dbl>, required_car_parking_spaces <dbl>, ...
```

```
mean(hotel_bookings_city$lead_time)
```

```
## [1] 109.7357
```

```
hotel_summary <-  
  hotel_bookings %>%  
  group_by(hotel) %>%  
  summarise(average_lead_time=mean(lead_time),  
            min_lead_time=min(lead_time),  
            max_lead_time=max(lead_time))
```

```
head(hotel_summary)
```

```
## # A tibble: 2 x 4
```

```
##   hotel          average_lead_time min_lead_time max_lead_time  
##   <chr>                <dbl>          <dbl>          <dbl>  
## 1 City Hotel           110.              0             629  
## 2 Resort Hotel         92.7              0             737
```