
cyclistic case study

*Comment une société de vélos en libre-service
peut-elle connaître un succès rapide ?*



intro

**executive
summary**

sources

résumé

*En quoi diffèrent les usages
des abonnés et des cyclistes
occasionnels ?*

propositions

*saisonnalité
type d'usage
heure d'usage*

annexes

intro

Pour cette étude de cas qui clôt le programme de formation certifiante sur l'analyse des données proposée par Google sur le site Coursera « Google Data Analytics », **j'assume le rôle d'un Data Analyst junior chez Cyclistic**, une entreprise de vélos en libre-service basée à Chicago.

Cyclistic propose plus de 5 800 vélos dans plus de 600 stations d'accueil réparties sur Chicago. La société propose deux types de vélos classiques et électriques et les commercialise à **deux segments de clientèle : les occasionnels et les membres abonnés**. Les cyclistes occasionnels sont les clients qui achètent des laissez-passer pour un trajet simple ou une journée complète. **Les membres sont les clients qui achètent des abonnements annuels de la société Cyclistic, ces derniers sont plus rentables que les cyclistes occasionnels.**

La stratégie future de la société vise à **convertir les usagers occasionnels en membres annuels (abonnés)**. Pour ce faire, l'équipe d'analystes marketing est chargée de mieux comprendre en quoi les cyclistes membres et les cyclistes occasionnels diffèrent en analysant des données historiques sur les trajets à vélo de la société de l'année 2023 afin d'identifier les tendances et de proposer une campagne publicitaire appropriée à la stratégie de l'entreprise.

Mes interlocuteurs :

Lily Moreno : La directrice du marketing et ma N+1. Lily Moreno est responsable du développement de campagnes et d'initiatives visant à promouvoir le programme de vélos en libre-service. Il peut s'agir d'e-mails, de médias sociaux et d'autres canaux.

Équipe d'analytique marketing Cyclistic : Une équipe d'analystes de données qui sont responsables de la collecte, de l'analyse et du reporting des données qui aident à guider la stratégie marketing de Cyclistic.

Équipe de direction Cyclistic : L'équipe de direction, notoirement orientée vers les détails, décidera d'approuver ou non le programme de marketing recommandé.

executive summary

J'ai analysé le data disponible pour l'année 2023 afin d'**identifier les variations de tendances** entre l'**usage des membres annuels et des utilisateurs occasionnels**.

Dans ces deux groupes, les usages diffèrent sur plusieurs points :

- **La proportion d'utilisateurs occasionnels augmente au printemps et pendant l'été**, de juin à août.
- Toute l'année, nos membres payants sont plus nombreux.
- **L'usage varie selon le jour de la semaine :**
 - Les abonnés vont favoriser un usage tout au long de la semaine avec des durées homogènes.
 - Les occasionnels iront favoriser le week-end et une durée d'utilisation extrêmement variée.
- En regardant les horaires d'utilisation :
 - **les deux groupes favorisent la soirée avec un pic à 17h.**
 - **Les trajets matinaux sont l'apanage des abonnés.**

Les recommandations que je ferai en fin de ce document se basent sur ces trois éléments de différenciation :

- **la saisonnalité**
- **le type d'usage**
- **l'heure d'usage**

SOURCES

Pour cette analyse, j'utilise le **data mis à disposition par Cyclistic** que l'on peut retrouver [ici](#).
(Remarque: Les jeux de données ont un nom différent car Cyclistic est une entreprise fictive).

Les données ont été mises à disposition par **Motivate International Inc.** sous cette [licence](#).

J'utilise les données de **Janvier 2023 à Décembre 2023** pour baser mon analyse sur une année pleine.

Le Data a été téléchargé et des copies ont été stockées en local ainsi que sur google Drive.

Le data se présente sous la forme de **fichiers CSV (comma-separated values), de 13 colonnes.**

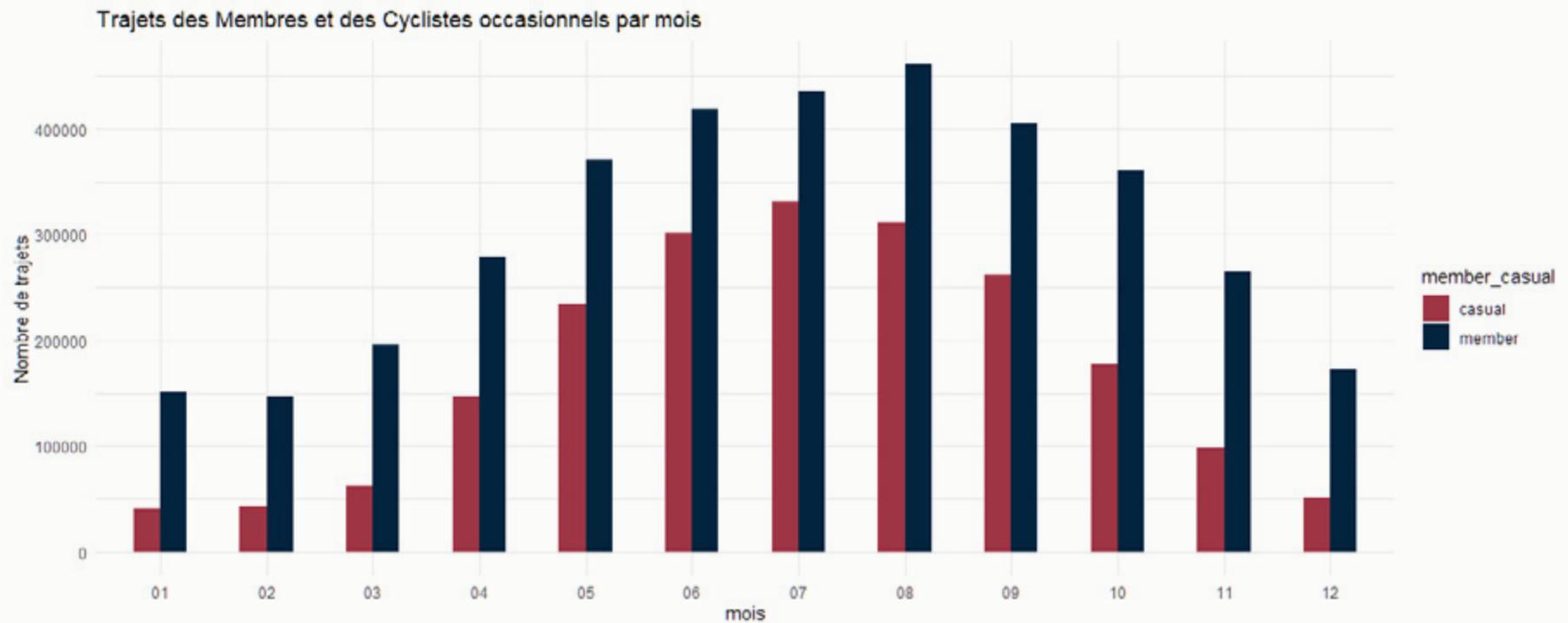
Dans l'optique du Case Study, ce dataset me permettra de répondre à la demande.
Mais des questions de confidentialité m'empêcheront d'utiliser des éléments identifiables
afin de **savoir si les utilisateurs occasionnels sont des occurrences uniques**
ou si ce sont des clients récurrents.

note: vous pouvez retrouver les liens vers le code utilisé pour l'analyse [ici](#).

*En quoi diffèrent
les usages
des abonnés
et des cyclistes
occasionnels?*

résumé de l'analyse

Augmentation de l'usage des cyclistes occasionnels pendant l'été.



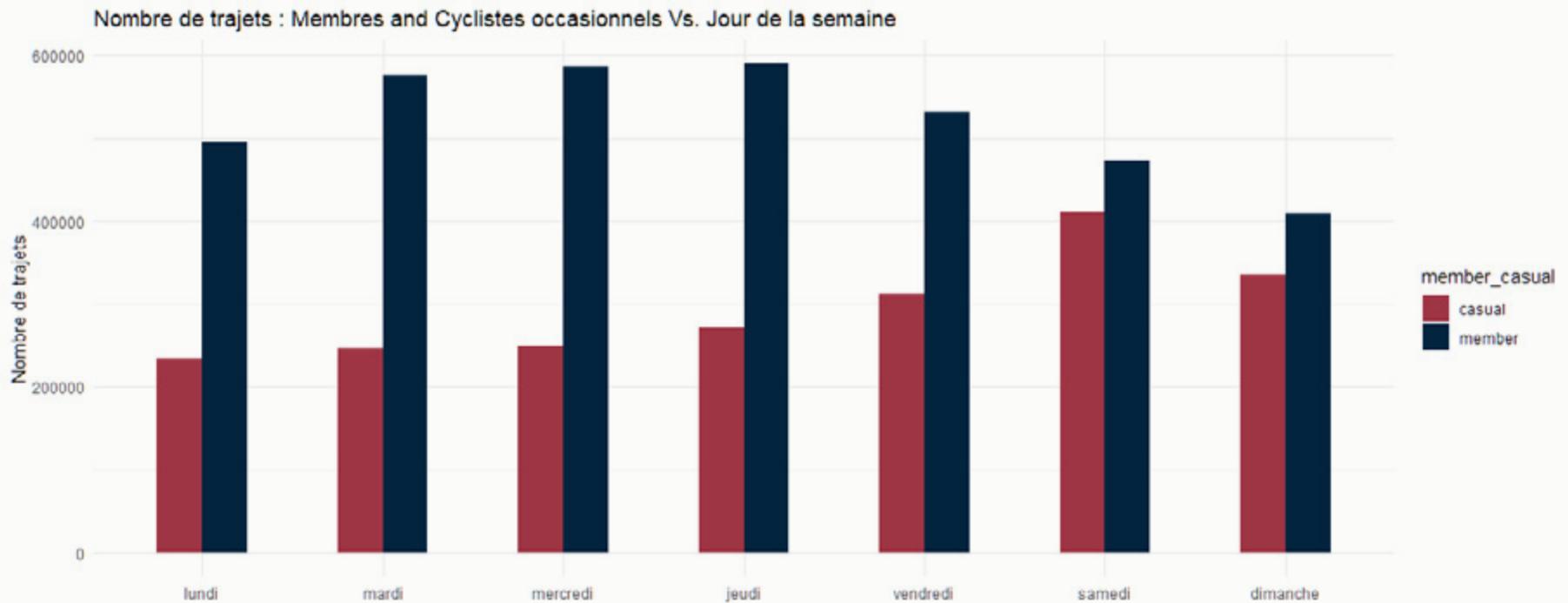
source: *cyclistic bike rides jan2023 to dec2023*

résumé de l'analyse

Jours semaine / Weekends

Les abonnés préfèrent **le milieu de la semaine**

alors que les cyclistes occasionnels voyagent majoritairement **le weekend**.

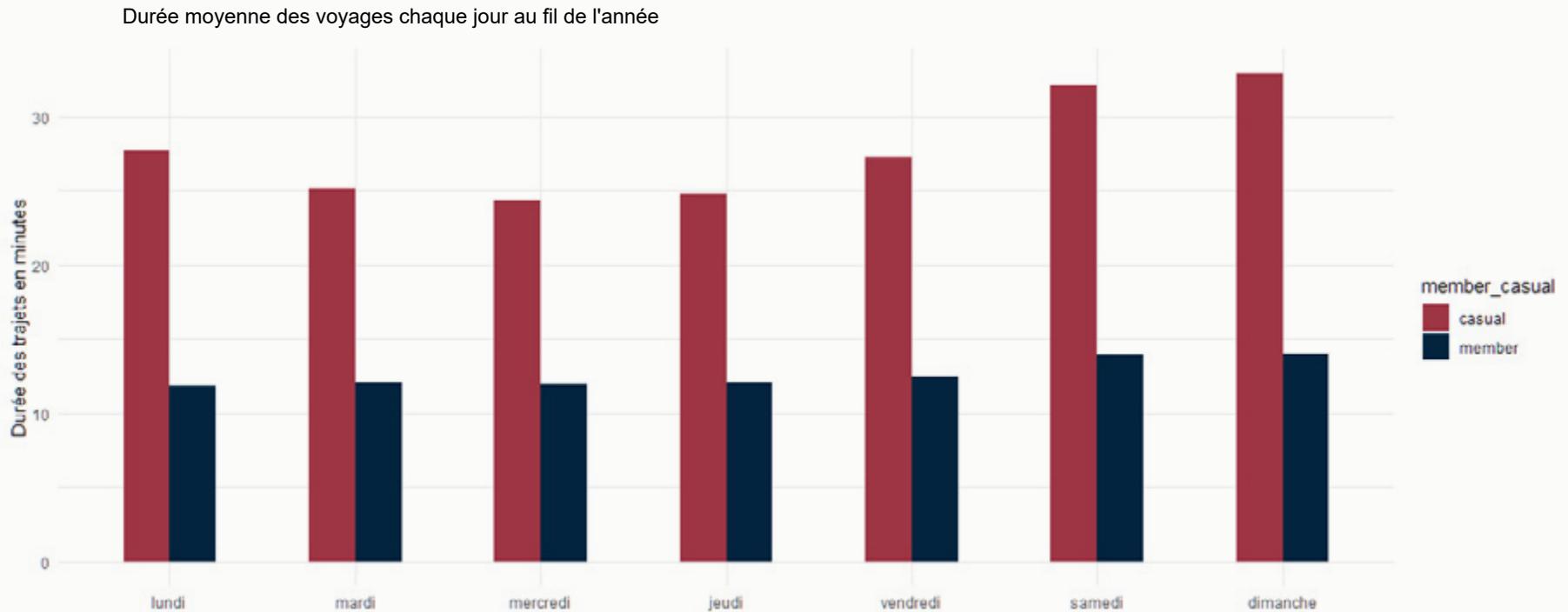


source: cyclistic bike rides jan2023 to dec2023

résumé de l'analyse

Des trajets stables pour les abonnés.

Les abonnés font des **trajets courts** et avec **globalement la même durée** tout au long de la semaine.

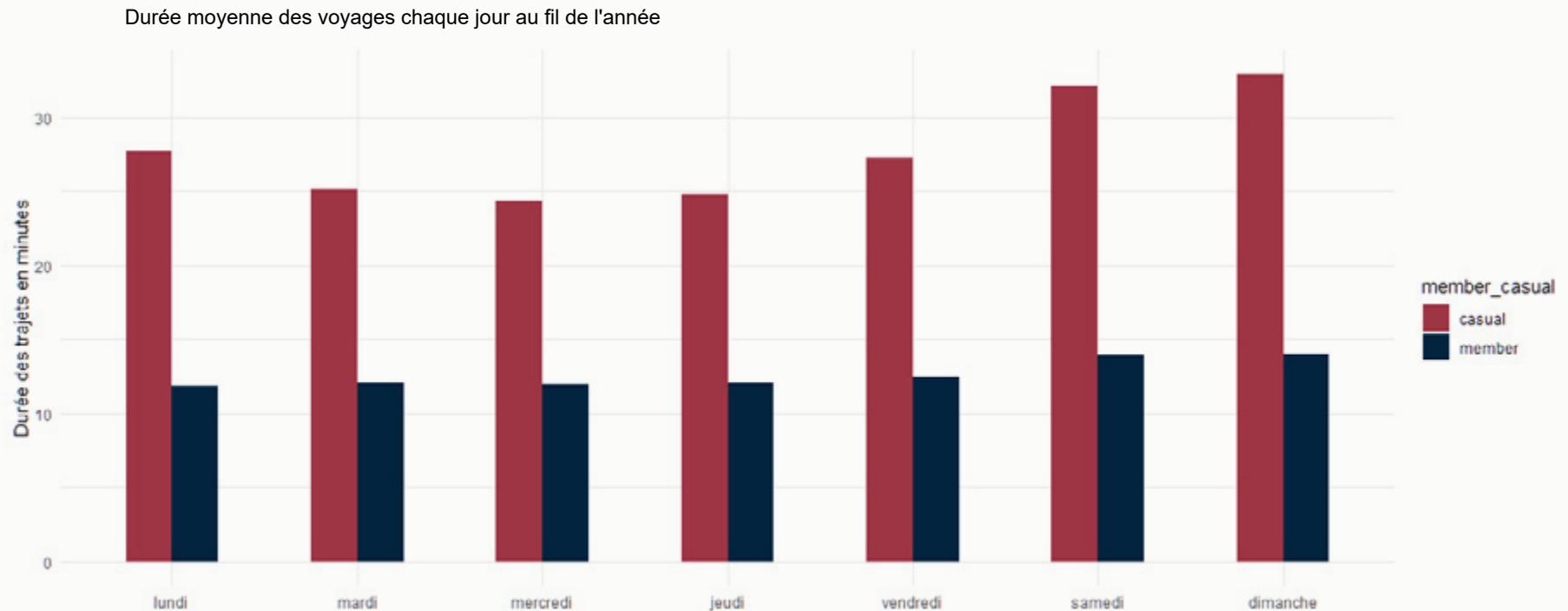


source: cyclistic bike rides jan2023 to dec2023

résumé de l'analyse

Une utilisation "week-end" pour les cyclistes occasionnels.

De manière générale, les occasionnels font des trajets à la durée plus longue et avec une augmentation de la durée pendant le week-end.

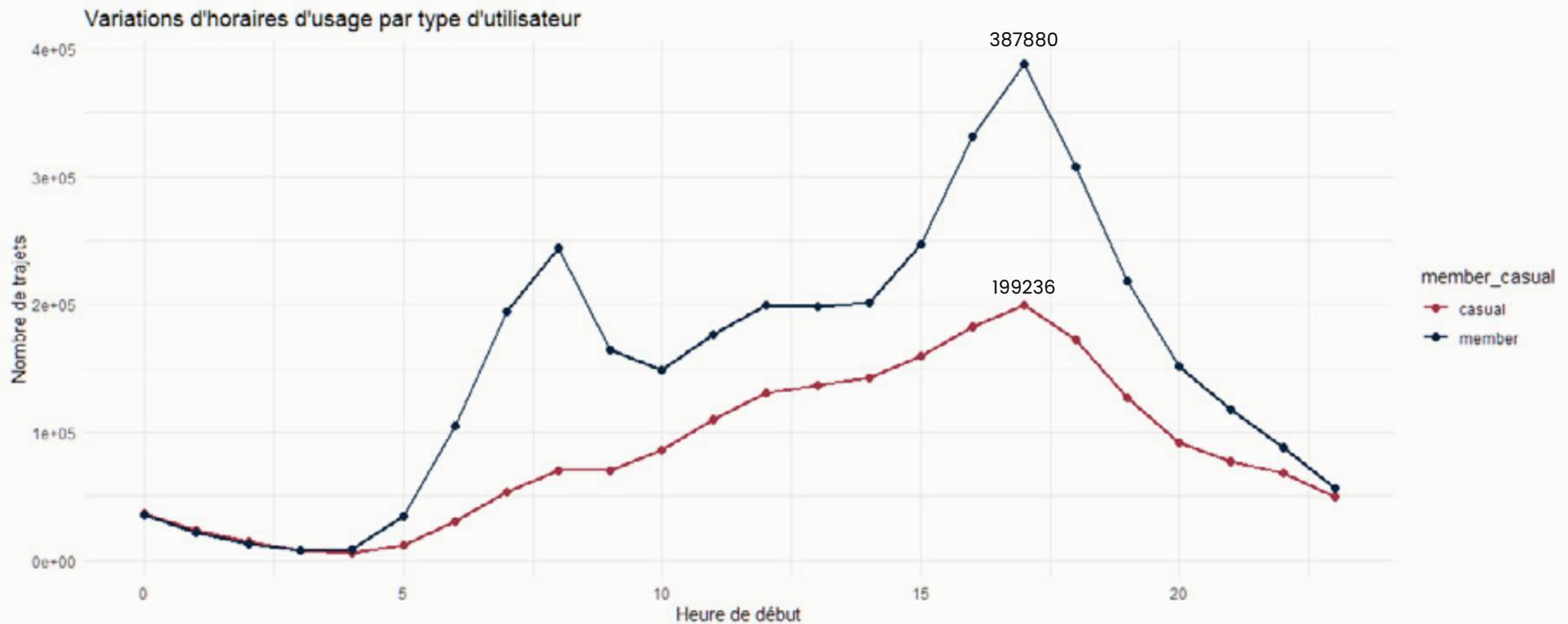


source: cyclistic bike rides jan2023 to dec2023

résumé de l'analyse

Des courbes semblables sauf le matin.

La fréquence d'utilisation est semblable pour les abonnés comme pour les cyclistes occasionnels avec la période du matin (8h) comme différence majeure.



source: cyclistic bike rides jan2023 to dec2023

*Comment ces insights
vont transformer
ces cyclistes occasionnels
en abonnés annuels?*

recommandations

saisonnalité : l'augmentation du printemps et de l'été

L'augmentation de la demande de nos vélos pendant le printemps et l'été nous offre une fenêtre de tir toute trouvée pour des campagnes de promotions et de publicité.

proposition :

Créer une campagne de mailing ciblant spécifiquement les cyclistes occasionnels pendant la fin du printemps et tout au long de l'été.

Prévoir des promotions et du matériel promotionnel en lien : Tarifs réduits sur l'abonnement annuel sur les 3 mois de l'été, des propositions de balades été et leurs contreparties hiver pour promouvoir l'usage all year long, faire des ventes flash de l'abonnement à prix cassés sur un an, etc.

recommandations

type d'usage: le même vélo, des habitudes différentes

Les abonnés utilisent nos vélos pour aller bosser: pendant la semaine, des pics le matin et le soir, avec des durées stables.

Les cyclistes occasionnels ont un usage récréatif: le week-end, pour aller se balader en soirée l'été.

proposition:

Créer des variations "sur mesure" de l'abonnement annuel.

Offre spéciale «Riders du week-end» pour pousser l'usage en week-ends.

Offre spéciale «Long rider» qui va chercher les aficionados des grandes balades.

Offre spéciale «Half Week» avec un tarif réduit le week-end selon l'usage en pleine semaine.

Bref, des offres qui pivotent autour de l'usage du week-end ou des variations de distance.

Organiser des événements spécifiquement pour les abonnés pendant les weekends.

Jouer sur le FOMO des occasionnels pour les amener à s'abonner.

recommandations

l'heure d'usage: l'opportunité de la tranche matinale

La seule différence notable dans les horaires d'utilisation des vélos se situe dans la tranche matinale: entre 05:00 et 10:00.

proposition:

Communiquer sur les bienfaits sanitaires de Cyclistic pour les trajets quotidiens.

Appuyer sur les bienfaits du vélo pour la santé avec des campagnes marketing ciblées. Les cyclistes abonnés et occasionnels ne diffèrent pas tant que ça en horaire d'usage, l'occasion présentée par le vide de la tranche matinale en terme de croissance est non négligeable.

packages Rstudio

```
tidyverse  
lubridate  
ggplot2  
dplyr  
readr  
tibble
```

code pour traiter (langage: R)

```
jan_23_data <- read_csv(«202301-divvy-tripdata.csv»)  
fev_23_data <- read_csv(«202302-divvy-tripdata.csv»)  
mar_23_data <- read_csv(«202303-divvy-tripdata.csv»)  
avr_23_data <- read_csv(«202304-divvy-tripdata.csv»)  
mai_23_data <- read_csv(«202305-divvy-tripdata.csv»)  
jun_23_data <- read_csv(«202306-divvy-tripdata.csv»)  
jul_23_data <- read_csv(«202307-divvy-tripdata.csv»)  
aug_23_data <- read_csv(«202308-divvy-tripdata.csv»)  
sept_23_data <- read_csv(«202309-divvy-tripdata.csv»)  
oct_23_data <- read_csv(«202310-divvy-tripdata.csv»)  
nov_23_data <- read_csv(«202311-divvy-tripdata.csv»)  
dec_23_data <- read_csv(«202312-divvy-tripdata.csv»)  
  
# maintenant on va vérifier les noms de colonnes  
# avant de pouvoir tout réunir pour avoir du data sur l'année  
  
colnames(jan_23_data)  
colnames(fev_23_data)  
colnames(mar_23_data)  
colnames(avr_23_data)  
colnames(mai_23_data)  
colnames(jun_23_data)  
colnames(jul_23_data)  
colnames(aug_23_data)  
colnames(sept_23_data)  
colnames(oct_23_data)
```

annexes

```
colnames(nov_23_data)
colnames(dec_23_data)

# c'est bon, il y a bien 13 colonnes partout et toutes ont les mêmes intitulés
# on réunit tout ça en un seul jeu maintenant

trip_data <- bind_rows(jan_23_data, fev_23_data, mar_23_data, avr_23_data, mai_23_data, jun_23_
data, jul_23_data, aug_23_data, sept_23_data, oct_23_data, nov_23_data, dec_23_data)

# on découpe les dates en jours / mois / années

trip_data$date <- as.Date(trip_data$started_at)
trip_data$month <- format(as.Date(trip_data$date), «%m»)
trip_data$day <- format(as.Date(trip_data$date), «%d»)
trip_data$year <- format(as.Date(trip_data$date), «%Y»)
trip_data$day_of_week <- format(as.Date(trip_data$date), «%A»)
colnames(trip_data) #to get the names of all the columns

# on rajoute une colonne de 'durée de voyage'

trip_data$ride_length <- difftime(trip_data$ended_at, trip_data$started_at)
str(trip_data) #on inspecte l'ajout des colonnes

trip_data$ride_length <- as.numeric(as.character(trip_data$ride_length))#on passe en numérique pour faciliter
les calculs
is.numeric(trip_data$ride_length) #le [1] TRUE nous le confirme
str(trip_data)#on le revérifie parce que

#maintenant on rajoute une colonne distance parcourue en km
trip_data$ride_distance <- distGeo(matrix(c(trip_data$start_lng, trip_data$start_lat), ncol=2), matrix
(c(trip_data$end_lng, trip_data$end_lat), ncol=2))
trip_data$ride_distance <- trip_data$ride_distance/1000 #distance en km

# petit nettoyage des entrées où le ride_length est négatif ou égal à 0
trip_data_clean <- trip_data[!(trip_data$ride_length <= 0),]
glimpse(trip_data_clean) #création nouvelle table avec les données clean

View(trip_data_clean) #on retourne vérifier la version nettoyée avant de passer à la phase d'analyse
```

code pour analyser (langage: R)

```
str(trip_data_clean) #on reprend (depuis la veille) avec la structure du data
summary(trip_data_clean) #on revérifie les détails

#on va regrouper les cyclistes occasionnels et les membres et comparer leurs min, max, average et median ride_
length (durée)
trip_data_clean %>%
  group_by(member_casual) %>%
  summarise(average_ride_length = mean(ride_length), median_length = median(ride_length),
            max_ride_length = max(ride_length), min_ride_length = min(ride_length))

#on va regrouper les cyclistes occasionnels et les membres et comparer le nombre de rides totaux
trip_data_clean %>%
  group_by(member_casual) %>%
  summarise(ride_count = length(ride_id))

#on aimerait bien avoir ce tibble en ggplot, donc je le mets dans une nouvelle table
ride_count <- trip_data_clean %>%
  group_by(member_casual) %>%
  summarise(ride_count = length(ride_id))

# je visualise ma nouvelle table qui devrait être comme le tibble
View(ride_count)

ggplot(ride_count, aes(x = member_casual, y = ride_count, fill = member_casual)) +
  geom_bar(stat = «identity») +
  geom_text(aes(label = ride_count), vjust = -0.5) +
  labs(title = «Nombre de trajets : membres Vs cyclistes occasionnels»,
        x = «Type d'utilisateur»,
        y = «Nombre de trajets») +
  theme_minimal() +
  scale_fill_manual(values = c(«member» = «#042540», «casual» = «#A13646»))
# on remarque la large différence entre le nombre de voyages des utilisateurs abonnés et des occasionnels

# maintenant on va découper ces résultats par jour

# organisation des jours
trip_data_clean$day_of_week <- ordered(trip_data_clean$day_of_week, levels=c(«lundi», «mardi», «mercredi»,
«jeudi», «vendredi», «samedi», «dimanche»))
```

annexes

```
trip_data_clean %>%
  group_by(member_casual, day_of_week) %>% #groupage par member_casual
  summarise(number_of_rides = n() #calcul du total de voyages et de leur durée moyenne
,average_ride_length = mean(ride_length),.groups=>drop) %>% # calcul de la durée moyenne
  arrange(member_casual, day_of_week) #sort

#on met ça dans une nouvelle table qu'on visionne
ride_count_daily <- trip_data_clean %>%
  group_by(member_casual, day_of_week) %>% #groupage par member_casual
  summarise(number_of_rides = n() #calcul du total de voyages et de leur durée moyenne
,average_ride_length = mean(ride_length/60),.groups=>drop) %>% # calcul de la durée moyenne en mi-
nutes
  arrange(member_casual, day_of_week) #sort
View(ride_count_daily)

# on ggplot le tout
ggplot(ride_count_daily, aes(x = day_of_week, y = number_of_rides, fill = member_casual)) +
  labs(title =»Nombre de trajets : Membres and Cyclistes occasionnels Vs. Jour de la semaine«,
        x = «»,
        y = «Nombre de trajets») +
  geom_col(width=0.5, position = position_dodge(width=0.5)) +
  scale_y_continuous(labels = fonction(x) format(x, scientific = FALSE))+
  theme_minimal() +
  scale_fill_manual(values = c(«member» = «#042540», «casual» = «#A13646»))

# tant qu'à faire, avec cette table on profite pour regarder la durée moyenne
ggplot(ride_count_daily, aes(x = day_of_week, y = average_ride_length, fill = member_casual)) +
  labs(title =»Durée moyenne des trajets : Membres and Cyclistes occasionnels Vs. Jour de la semaine«,
        x = «»,
        y = «Durée des trajets en minutes») +
  geom_col(width=0.5, position = position_dodge(width=0.5)) +
  scale_y_continuous(labels = fonction(x) format(x, scientific = FALSE))+
  theme_minimal() +
  scale_fill_manual(values = c(«member» = «#042540», «casual» = «#A13646»))

# maintenant on va vérifier l'utilisation mensuelle
trip_data_clean %>%
  group_by(member_casual, month) %>%
  summarise(number_of_rides = n(),.groups=>drop) %>%
  arrange(member_casual, month)
```

```
number_of_rides <- trip_data_clean %>%
  group_by(member_casual, month) %>%
  summarise(number_of_rides = n(), .groups => drop) %>%
  arrange(member_casual, month)

ggplot(number_of_rides, aes(x = month, y = number_of_rides, fill = member_casual)) +
  labs(title = »Trajets des Membres et des Cyclistes occasionnels par mois«,
        x = «mois»,
        y = «Nombre de trajets») +
  theme(axis.text.x = element_text(angle = 45)) +
  geom_col(width=0.5, position = position_dodge(width=0.5)) +
  scale_y_continuous(labels = function(x) format(x, scientific = FALSE)) +
  theme_minimal() +
  scale_fill_manual(values = c(«member» = «#042540», «casual» = «#A13646»))

# comparons la distance moyenne parcourue entre membres et cyclistes occasionnels

trip_data_clean %>%
  group_by(member_casual) %>% drop_na() %>%
  summarise(average_ride_distance = mean(ride_distance))

average_ride_distance <- trip_data_clean %>%
  group_by(member_casual) %>% drop_na() %>%
  summarise(average_ride_distance = mean(ride_distance))
View(average_ride_distance)

ggplot(average_ride_distance, aes(x = member_casual, y = average_ride_distance, fill = member_casual)) +
  geom_bar(stat = «identity») +
  geom_text(aes(label = round(average_ride_distance, 2)), vjust = -0.5) +
  labs(title = «Distance moyenne des trajets par type d'utilisateur»,
        x = «»,
        y = «Distance moyenne (km)») +
  theme_minimal() +
  scale_fill_manual(values = c(«member» = «#042540», «casual» = «#A13646»))

# testons voir les heures de trajet

# Ajouter une colonne pour l'heure de début des trajets
trip_data_clean <- trip_data_clean %>%
  mutate(start_hour = lubridate::hour(started_at))

# Calculer le nombre de trajets par heure pour les membres et les utilisateurs occasionnels
```

```
hourly_usage <- trip_data_clean %>%
  group_by(member_casual, start_hour) %>%
  summarise(ride_count = n()) %>%
  ungroup()

# Visualiser les variations d'horaires d'usage
#Calculer les valeurs maximales pour chaque groupe
max_values <- hourly_usage %>%
  group_by(member_casual) %>%
  summarize(max_ride_count = max(ride_count),
            max_hour = start_hour[which.max(ride_count)])

ggplot(hourly_usage, aes(x = start_hour, y = ride_count, color = member_casual, group = member_casual)) +
  geom_line(size = 1) +
  geom_point(size = 2) +
  labs(title = «Variations d'horaires d'usage par type d'utilisateur»,
       x = «Heure de début»,
       y = «Nombre de trajets») +
  theme_minimal() +
  scale_color_manual(values = c(«member» = «#042540», «casual» = «#A13646»))

daily_avg_duration <- ride_count_daily %>%
  group_by(day_of_week, member_casual) %>%
  summarize(avg_duration = mean(duration))

view(daily_avg_duration)

ggplot(number_of_rides, aes(x = day_of_week, y = avg_duration, fill = member_casual)) +
  labs(title = «Durée moyenne des voyages chaque jour au fil de l'année»,
       x = «mois»,
       y = «Nombre de trajets») +
  theme(axis.text.x = element_text(angle = 45)) +
  geom_col(width=0.5, position = position_dodge(width=0.5)) +
  scale_y_continuous(labels = function(x) format(x, scientific = FALSE))+
  theme_minimal() +
  scale_fill_manual(values = c(«member» = «#042540», «casual» = «#A13646»))
```

annexes

sites webs et documentation

datanovia.com

palletton.com

documentation RStudio

Github

Reddit/R

kaggle

retrouvez moi sur www.charlesgrillet.com

merci