

Handling Categorical Data in R

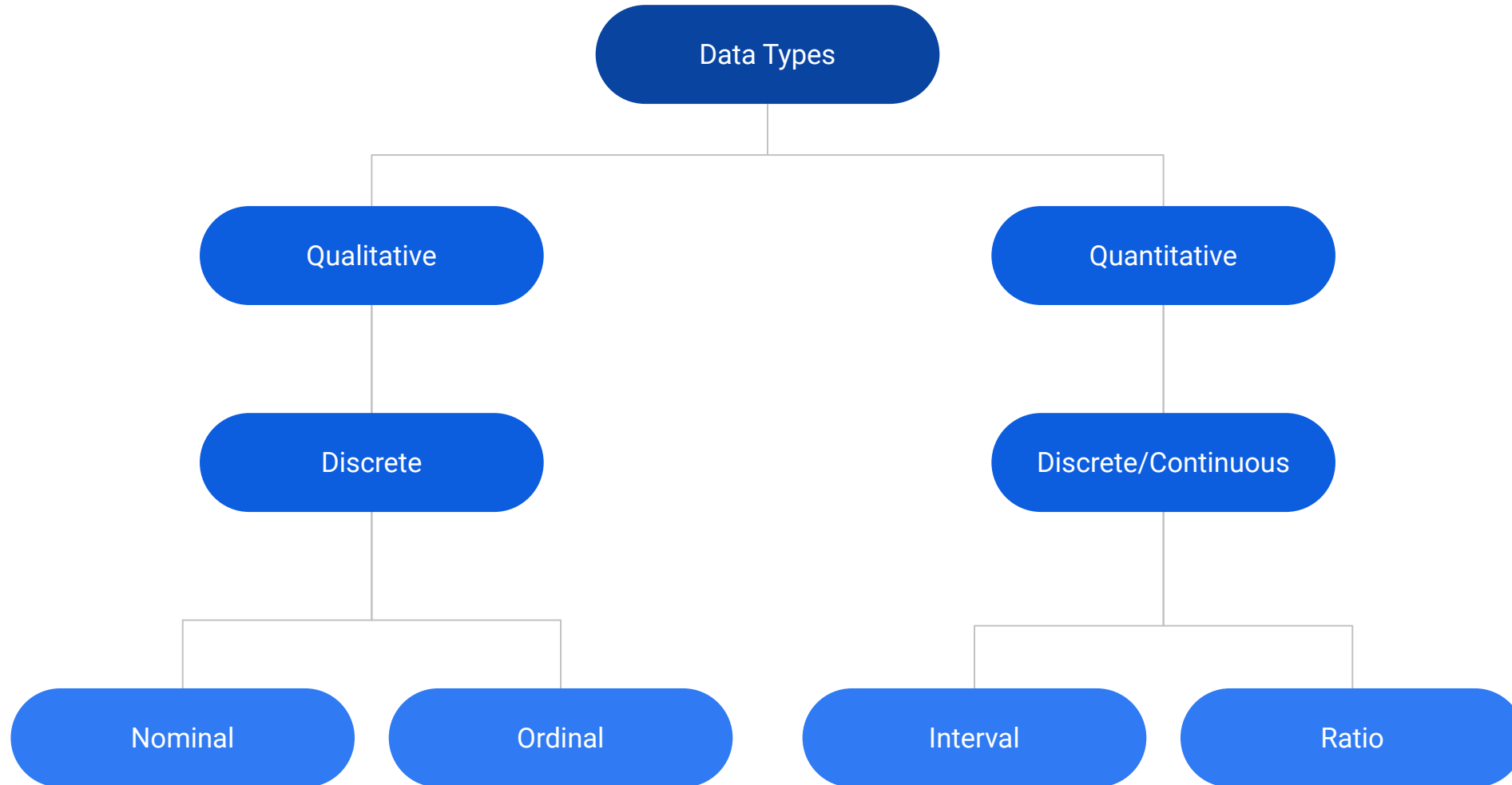


- understand categorical data
- case study intro
- factor in R
- data summarization
- data manipulation
- data visualization

- [Slides](#)
- [Code & Data](#)
- [RStudio Cloud](#)
- [Online Course](#)
- [Blog Post](#)

Module 1 **Introduction**







one



two



three



- it is always discrete
- it may be divided into groups
- consists of names or labels
- takes on limited & fixed number of possible values
- arises in situation when counting is involved
- analysis generally involves the use of data tables



yes



no



yes



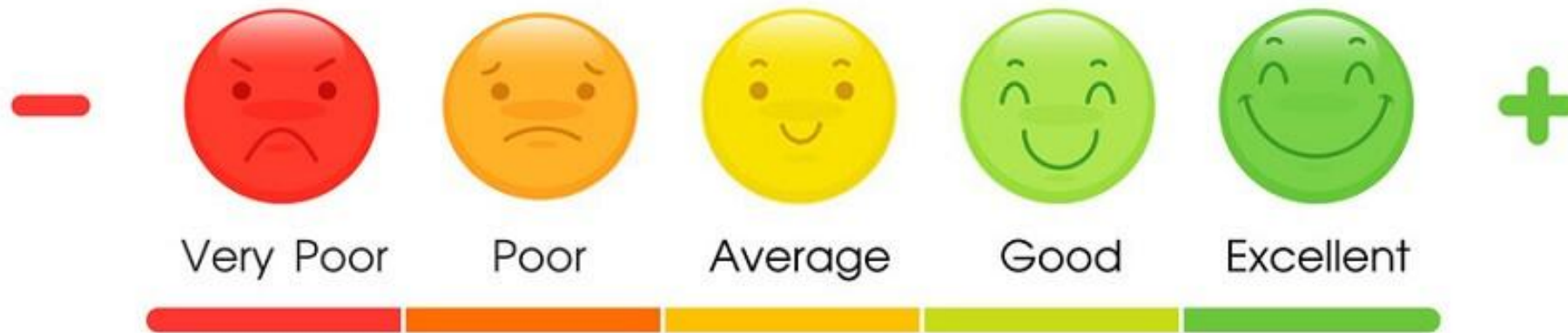
maybe



no



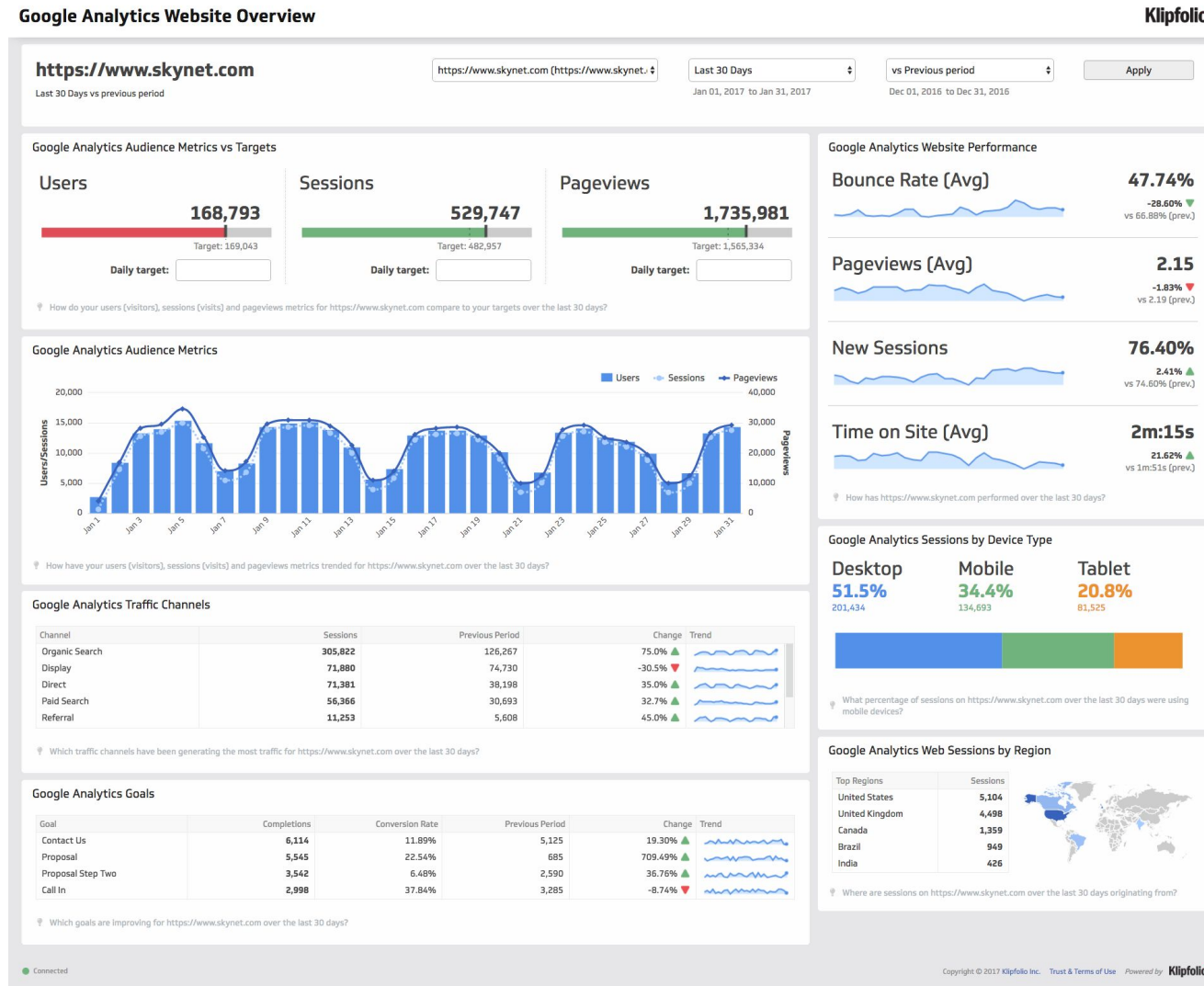
CUSTOMER SATISFACTION



- data can be quantitative or qualitative
- qualitative data is always discrete
- dichotomous data consists of only 2 groups/levels
- polychotomous data consists of more than 2 groups/levels
- nominal data do not have an intrinsic order
- in ordinal data
 - categories can be ordered or ranked
 - and difference between the categories cannot be determined

NOW, IT'S YOUR TURN





Column Name	Description
device	Device used to visit the website
os	Operating system of the device
browser	Browser used to visit the website
user_type	New or returning visitor
channel	Source of traffic
gender	Gender of the visitor
frequency	Number of visits to the website including the current one
recency	Number of days since the last visit to the website
page_depth	Number of pages browsed on the website
hour_of_day	Hour of day

Column Name	Description
age	Age of the visitor
duration	Time spent on the website (in seconds)
landing_page	Page on which the visitor landed (first page)
exit_page	Page from which the visitor exited the website (last page)
country	Country of the visitor
city	City of the visitor
quantity	Number of units purchased
revenue	Revenue from the visitor
purchase_flag	Whether the visitor made a transaction or not
user_rating	Rating given by customer

Import Text Data

File/URL:
J:/R/ebooks/forcats/analytics_raw.csv Browse...

Data Preview:

X1 (double)	device (character)	os (character)	browser (character)	user_type (character)	channel (character)	gender (character)	frequency (double)	recency (double)	page_depth (double)	hour_of_day (character)
1	Desktop	Windows	Chrome	New Visitor	Organic Search	female	1	0	1	02
2	Mobile	iOS	Safari	Returning Visitor	Organic Search	NA	3	1	1	20
3	Desktop	Chrome OS	Chrome	New Visitor	Direct	NA	1	0	5	05
4	Desktop	Macintosh	Chrome	Returning Visitor	Organic Search	NA	2	0	1	17
5	Desktop	Macintosh	Chrome	Returning Visitor	Referral	NA	5	8	1	04
6	Mobile	Android	Chrome	New Visitor	Organic Search	NA	1	0	5	00
7	Desktop	Windows	Chrome	New Visitor	Organic Search	NA	1	0	4	03

Previewing first 50 entries.

Import Options:

Name: analytics_raw ☒ First Row as Names Delimiter: Comma Escape: None

Skip: 0 ☒ Trim Spaces Quotes: Default Comment: Default

☒ Open Data Viewer Locale: Configure... NA: Default

Code Preview:

```
library(readr)
analytics_raw <- read_csv("analytics_raw.csv")
View(analytics_raw)
```

? Reading rectangular data using readr Import Cancel

```
library(readr)
read_csv("analytics_raw.csv", col_types =
  cols_only(
    device = col_factor(levels = c("Desktop", "Tablet", "Mobile")),
    gender = col_factor(levels = c("female", "male", "NA")),
    user_rating = col_factor(levels = c("1", "2", "3", "4", "5"),
                              ordered = TRUE)))
```

Module 2 **Introduction to Factors**



- introduction to factor
- how to recognize factor variables
- how to coerce other data types to factor
- handle missing values
- handle ordinal data
- specify order of levels/categories

Function	Description
<code>is.factor()</code>	Identify
<code>is.ordered()</code>	
<code>as.factor()</code>	Convert
<code>as_factor()</code>	
<code>as.ordered()</code>	
<code>factor()</code>	Create
<code>ordered()</code>	

- R uses **factor** to handle categorical data
- Use **as.factor()** or **as_factor()** to coerce other data types to factor
- Use **is.factor()** and **is.ordered()** to identify factor and ordered factor respectively
- Use **factor()** to
 - Specify labels
 - Modify labels
 - Handle missing data
 - Create ordered factors
 - Specify order of levels
- Use **ordered()** to create ordered factors

NOW, IT'S YOUR TURN



Module 3
Summarize Factors



- <https://forcats.tidyverse.org/>
- <https://r4ds.had.co.nz/factors.html>
- <https://recipes.tidymodels.org/reference/discretize.html>
- <https://ggplot2.tidyverse.org/>
- <https://haleyjeppson.github.io/ggmosaic/>
- <https://rpkgs.datanovia.com/ggpubr/reference/ggdonutchart.html>

- [Website](#)
- [Free Online R Courses](#)
- [R Packages](#)
- [Shiny Apps](#)
- [Blog](#)
- [GitHub](#)
- [YouTube](#)
- [Twitter](#)
- [Linkedin](#)



Thank You

For more information please visit our website
www.rsquaredacademy.com