

Reporte Etapa 2: Data Understanding

Abstract	3
Obtención de datos	3
Reporte inicial de recolección de datos	3
Descripción de datos	3
Reporte de descripción de los datos	3
Exploración de los datos	4
Reporte de Exploración	4
Verificación de calidad de datos	4
Reporte de Calidad de datos	4

Abstract

En esta etapa 2 “Data Understanding” de CRISP-DM, analizaremos los datos que se nos han entregado, realizando los procesos de obtención, descripción, exploración y verificación.

Obtención de datos

Reporte inicial de recolección de datos

En el proceso de obtención se nos compartió un dataset preparado para trabajar con el proyecto, sin necesidad de hacer cambios o modificaciones, es decir, previamente revisado, dichos datos se cargaron a un servidor en el que se realizará todo el procesamiento.

Legalmente no hay problemas con el tratamiento de los datos ya que han sido preprocesados con el fin de obtener el anonimato de los usuarios, con valores únicos, sin embargo fue necesario obtener información adicional a la proporcionada, de las universidades a analizar (latitudes y longitudes).

En cuanto a información relevante, el socio formador nos especificó que los datos pertenecen a alrededor del 30% de la población de Santiago de Chile, debido a que la empresa Movistar, la cual proporciona estos datos, tiene ese porcentaje de usuarios de todas las demás operadoras móviles.

Para cargar los datos utilizamos las herramientas Pandas y Spark..

Descripción de datos

Reporte de descripción de los datos

Dentro de los datos contamos con diferentes columnas con diferentes tipos de datos:

- Phone_ID: Esta columna contiene el ID del usuario en un String
- Timestamp: Contiene la fecha y hora en la que se conectó el usuario a la antena. El tipo de dato es string.
- bts_id: Contiene los ID de las antenas a la cual se hizo la conexión. El tipo de dato es String.
- Lat: contiene la latitud de la antena en cuestión. El tipo de dato es float.
- Lon: contiene la longitud de la antena en cuestión. El tipo de dato es float.

Se cuenta con 49,618,132 registros cada uno con 5 campos.

Valores únicos:

count(DISTINCT PHONE_ID)	count(DISTINCT timestamp)	count(DISTINCT bts_id)	count(DISTINCT lat)	count(DISTINCT lon)
1353435	86400	1871	1198	1264

PHONE_ID = 1,353,435
timestamp = 86,400

bts_id = 1871
lat = 1198

lon = 1264

Valores nulos:

PHONE_ID	timestamp	bts_id	lat	lon
0	0	0	0	0

PHONE_ID = 0

bts_id = 0

lon = 0

timestamp = 0

lat = 0

Exploración de los datos**Reporte de Exploración**

Para la exploración de datos decidimos omitirla, ya que nos estuvimos acercando con otros equipos para poder apoyarnos para esta etapa. Al adaptarnos a lo que realizó el grupo, de igual manera nos saltamos esta etapa y nos percatamos de esto en la etapa de preparación de los datos.

Para la exploración de datos decidimos omitirla, ya que al comunicarnos con los otros equipos obtuvimos información sobre sus hallazgos de los datos y con estos podemos tomar decisiones para las siguientes etapas.

Verificación de calidad de datos**Reporte de Calidad de datos**

De acuerdo con el reporte de descripción de los datos, podemos decir que los datos se encuentran completos y que son de calidad ya que no existen datos nulos, no hay datos faltantes y tampoco hay varios tipos de datos mezclados en una sola columna.

El formato y valores de cada columna es consistente y la cantidad de datos es inmensa para realizar un buen análisis desde nuestro punto de vista.