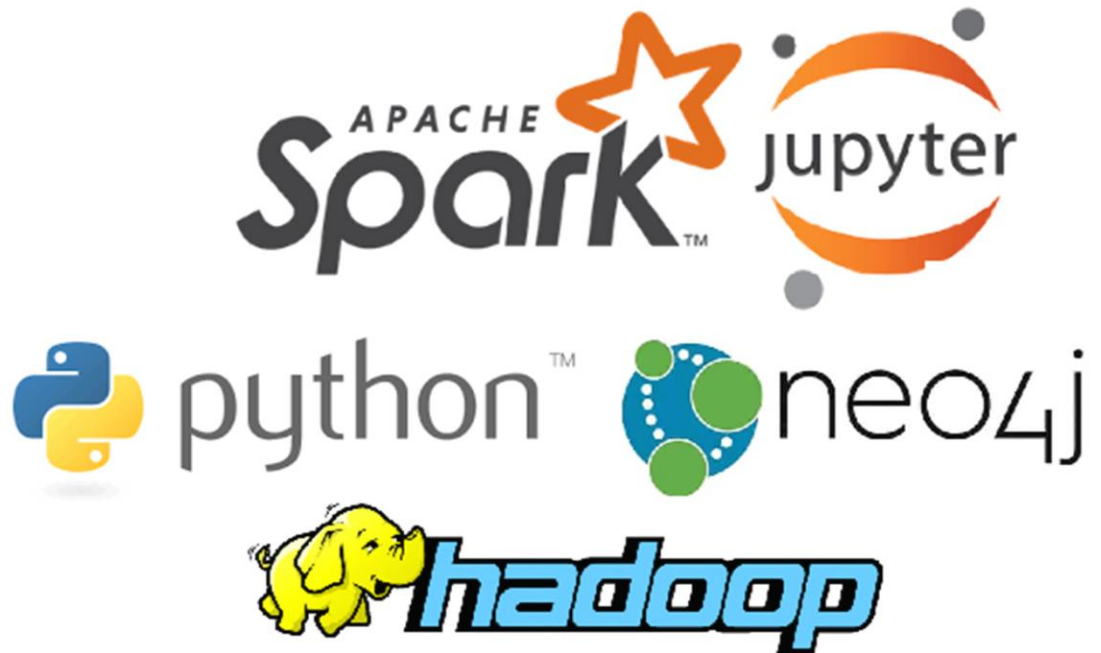




Pontificia Universidad  
**JAVERIANA**  
Cali

2021

## Big Data Processing - Project



JUAN SEBASTIAN REYES

ALEXANDER CASTRO CONTRERA

JHAN CARLOS DEL RIO

23-9-2021

# Contenido

---

Introducción.....	2
Pre-procesamiento de la Data .....	4

# Introducción

El conjunto de datos base para la elaboración del proyecto se toma directamente de reportes generados por la Alcaldía de Medellín. Es un conjunto de datos muy completo, contiene 1'725.243 registros de las instituciones médicas de Medellín donde se trataron distintas enfermedades a pacientes en el periodo 2009 – 2017.

<b>Características:</b>	Multivaluado
<b>Número de Registros:</b>	1'725.243
<b>Área:</b>	Atenciones, Morbilidad, Urgencias, Salud, Servicio, prestación
<b>Tipo de Datos:</b>	String, Char, Integer, Float
<b>Numero de Atributos:</b>	18
<b>Idioma:</b>	Español
<b>Autor:</b>	Secretaría de Salud
<b>Nivel de Acceso Público:</b>	Público
<b>Fecha de modificación:</b>	20/04/2028
<b>Fecha de publicación:</b>	19/09/2024
<b>Licencia:</b>	<a href="https://opendefinition.org/licenses/cc-by-sa/">https://opendefinition.org/licenses/cc-by-sa/</a>

## Información de los Atributos

- **Identificadas:**
  1. **consecutivo:** Número único (**PK**) de cada registro.
  2. **año:** Año en que se realizó el registro.
  3. **cod\_eas:** Identificador EPS (**FK** referencia '**nombre\_eas**').
  4. **nombre\_eas:** Nombre EPS.
  5. **edad:** Edad del conjunto de personas atendidas asociadas al '**consecutivo**'.
  6. **sexo:** Sexo común del conjunto de personas atendidas asociadas al '**consecutivo**'.
  7. **cod\_departamento:** Identificador Departamento. En este caso todos los registros están asociados al **ID:5** = '**Antioquia**'.
  8. **cod\_municipio:** Identificador Municipio. En este caso todos los registros están asociados al **ID:1** = '**Medellín**'.
  9. **zona:** Identifica el tipo de zona del conjunto de personas atendidas asociadas al '**consecutivo**'.
  10. **cod\_ips:** Identificador Institución (**FK** referencia '**nombre\_institucion**').
  11. **nombre\_institucion:** Nombre Institución.
  12. **cod\_dx\_salida:** Identificador Enfermedad (**FK** referencia '**nombre\_dx**').
  13. **nombre\_dx:** Nombre Enfermedad

**14. servicio:** Identificador del tipo de servicio. En este caso todos los registros están asociados a '**URGENCIAS**'.

**15. total\_atenciones:** Representa el numero de pacientes atendidos en la misma institución, por la misma enfermedad, del mismo sexo asociadas al '**consecutivo**'.

- **Desconocidas:**
  1. **tipo\_usuario**
  2. **tipo\_edad**
  3. **causa\_externa.**

## Clasificación de los Atributos

Evidenciamos que el conjunto de datos maneja datos con variables cualitativas y cuantitativas. A continuación, mostramos la clasificación:

Cualitativa (Nominal)	Cualitativa (Ordinal)	Cuantitativa (Intervalo)	Cuantitativa (Razón)
cod_eas nombre_eas sexo zona cod_ips nombre_institucion cod_dx_salida nombre_dx servicio	consecutivo año tipo_usuario tipo_edad cod_departamento cod_municipio causa_externa	edad	total_atenciones

## Valores Únicos

Como podemos observar se tiene que los atributos como servicio, **cod\_municipio** y **cod\_departamento** solo presentan 1 elemento dentro de su rango de opciones. Por otro lado, se tiene que el atributo **zona**, presenta un problema en la distinción entre mayúsculas y minúsculas, específicamente con la letra **U/u**. Además, se puede evidenciar que los datos de edad presentan algunos valores anormales.

# Pre-Procesamiento de la Data

---

Este es el proceso inicial más importante, la idea es que al finalizar este proceso las siguientes preguntas claves estén resueltas: **¿Tienen sentido los datos? ¿Los datos siguen las reglas apropiadas para su campo? ¿Los datos tienen problemas de calidad?**

Para ello, seguimos una serie de recomendaciones estrictas que permitieron un mayor desempeño del conjunto de datos a la hora de encaminarnos hacia la variable objetivo.

A continuación, detallamos el proceso:

**1. Eliminamos observaciones duplicadas o irrelevantes:**

Las columnas de servicio-cod\_municipio-cod\_departamento.

**2. Analizamos y corregimos errores estructurales:**

- Columna zona, caso de mayúsculas-minúsculas.
- Columna Edad en algunos registros no cumplía con el formato.

**3. Filtramos valores atípicos y evaluamos la importancia en el conjunto de datos:**

El hecho de que exista un valor atípico no significa que sea incorrecto. Este paso es necesario para determinar la validez de ese dato.

**4. Etapa de QA:**

En esta etapa de prueba lo que buscamos fue identificación de patrones de errores para reducir estos mismos en el futuro. Después de esta etapa hicimos una iteración mas por estos mismos 4 pasos.