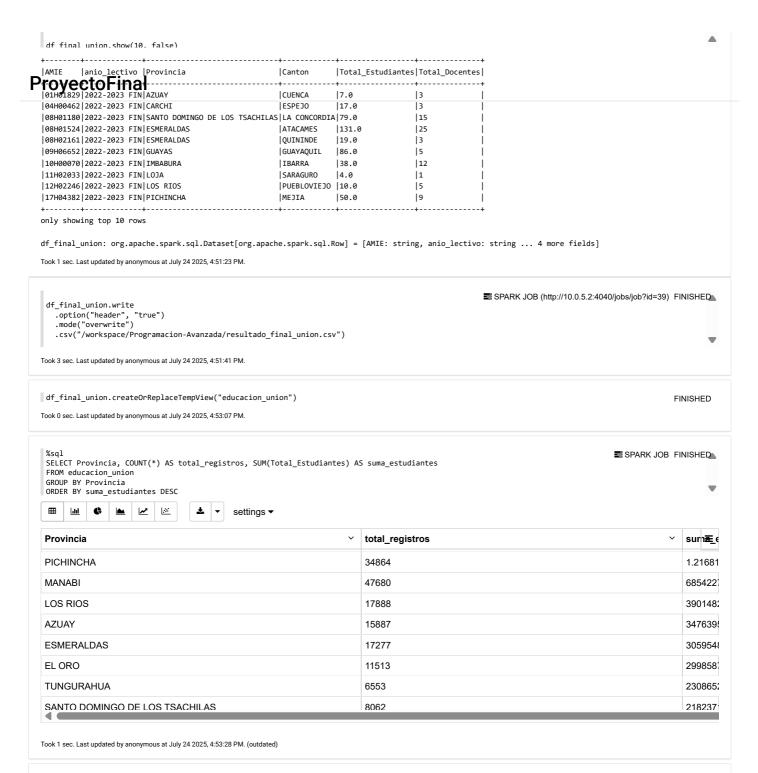
```
val df1 = spark.read
                                                                                                                                                              SPARK JOB FINISHED
Proyector ind ("true")
       .csv("/workspace/Programacion-Avanzada/2_MINEDUC_RegistrosAdministrativos_2022-2023Fin.csv")
      .option("header", "true")
.option("inferSchema", "true")
.option("delimiter", ";")
       .csv("/workspace/Programacion-Avanzada/registro-administrativo-historico_2009-2024-inicio.csv")
    val df3 = spark.read
      .option("header", "true")
.option("inferSchema", "true")
.option("delimiter", ";")
.csv("/worksnace/Programacion-
                                 ramacion-Avanzada/ineval cerectudiante2022 2023 2023diciembre ccv"\
  df1: org.apache.spark.sql.DataFrame = [A�o Lectivo: string, AMIE: string ... 248 more fields]
  df2: org.apache.spark.sql.DataFrame = [Anio_lectivo: string, Zona: string ... 25 more fields]
  df3: org.apache.spark.sql.DataFrame = [ciclo: string, grado: int ... 56 more fields]
  Took 3 sec. Last updated by anonymous at July 24 2025, 4:41:23 PM.
    val df1_clean = df1.select(
   "AMIE",
                                                                                                                                                                               FINISHED
       "A�o Ĺectivo",
       'Provincia",
       "Cant�n",
       "Total_Estudiantes",
      "Total_Docentes'
                                                                                                                                                                                        \overline{\mathbf{v}}
   ).withColumnRenamed("A♠o Lectivo". "anio lectivo")
  df1_clean: org.apache.spark.sql.DataFrame = [AMIE: string, anio_lectivo: string ... 4 more fields]
  Took 0 sec. Last updated by anonymous at July 24 2025, 4:44:12 PM.
    val df2_clean = df2.select(
                                                                                                                                                                               FINISHED
      "AMIE",
"Anio_lectivo",
       'Provincia",
       "Canton",
       "Total_Estudiantes",
      "Total_Docentes"
   ).withColumnRenamed("Anio lectivo", "anio lectivo")
  df2_clean: org.apache.spark.sql.DataFrame = [AMIE: string, anio_lectivo: string ... 4 more fields]
  Took 1 sec. Last updated by anonymous at July 24 2025, 4:44:45 PM.
    val d+3_clean = d+3.select(
                                                                                                                                                                               FINISHED
      "amie",
"ciclo",
       "grado"
        estado_eval"
      "isec", // Índice de desempeño escolar
"imat", // Matemática
      "ilyl", // Lengua y Literatura
"icn", // Ciencias Naturales
"ies" // Estudios Sociales
  "ies" // Estudios Sociales
).withColumnRenamed("amie", "AMIE")
  df3_clean: org.apache.spark.sql.DataFrame = [AMIE: string, ciclo: string ... 7 more fields]
  Took 0 sec. Last updated by anonymous at July 24 2025, 4:44:55 PM.
    // Renombra en df1_clean
val df1_ready = df1_clean
                                                                                                                                                                               FINISHED
      .withColumnRenamed("Cant�n", "Canton")
.withColumnRenamed("Provincia", "Provincia") // (Solo si hay problemas con otros nombres)
                                                                                                                                                                                        \overline{\mathbf{v}}
   val df2 ready = df2 clean
  df1_ready: org.apache.spark.sql.DataFrame = [AMIE: string, anio_lectivo: string ... 4 more fields]
  df2_ready: org.apache.spark.sql.DataFrame = [AMIE: string, anio_lectivo: string ... 4 more fields]
  Took 0 sec. Last updated by anonymous at July 24 2025, 4:46:22 PM. (outdated)
  val df12 = df1_ready.join(df2_ready, Seq("AMIE", "anio_lectivo", "Provincia", "Canton"), "outer")
                                                                                                                                                                               FINISHED
  df12: org.apache.spark.sql.DataFrame = [AMIE: string, anio_lectivo: string ... 6 more fields]
  Took 0 sec. Last updated by anonymous at July 24 2025, 4:46:49 PM.
                                                                                                                                                                               FINISHED
  val df_final = df12.join(df3_clean, Seq("AMIE"), "outer")
  df_final: org.apache.spark.sql.DataFrame = [AMIE: string, anio_lectivo: string ... 14 more fields]
  Took 0 sec. Last updated by anonymous at July 24 2025, 4:47:13 PM.
  df_final.show(10)
                                                                                                                                                              SPARK JOB FINISHED
```

```
|AMIE|anio_lectivo|Provincia|Canton|Total_Estudiantes|Total_Docentes|Total_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Cotal_Estudiantes|Total_Docentes|Total_Estudiantes|Total_Docentes|Total_Estudiantes|Total_Docentes|Total_Estudiantes|Total_Docentes|Total_Estudiantes|Total_Docentes|Total_Estudiantes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Docentes|Total_Doc
                                              null| null|
                                                                                          null|
                                                                                                                     null|
                                                                                                                                                   null|
                                                                                                                                                                              null| null| null|
                                                                                                                                                                                                                          null|null|null|null|null|
ProyectoFinahull null null null null null
                                                                                          null|
                                                                                                                     null|
                                                                                                                                                     null|
                                                                                                                                                                               null| null| null|
                                                                                                                                                                                                                          null|null|null|null|null|null|
                                                                                                                                                                                                                          null|null|null|null|null|null|
                                                                                          null
                                                                                                                     null|
                                                                                                                                                     null|
                                                                                                                                                                               null| null| null|
     |null|
                             null|
                                              null| null|
                                                                                          null
                                                                                                                     null
                                                                                                                                                     null|
                                                                                                                                                                                null| null| null|
                                                                                                                                                                                                                          null|null|null|null|null|null|
    |null|
                             null|
                                             null| null|
                                                                                          null|
                                                                                                                     null|
                                                                                                                                                     null|
                                                                                                                                                                                null| null| null|
                                                                                                                                                                                                                          null|null|null|null|null|null|
    |null|
                             null|
                                             null| null|
                                                                                          null|
                                                                                                                     null|
                                                                                                                                                     null|
                                                                                                                                                                                null| null| null|
                                                                                                                                                                                                                          null|null|null|null|null|
    Inull
                            null
                                             null| null|
                                                                                          null
                                                                                                                     null
                                                                                                                                                     null
                                                                                                                                                                               null| null| null|
                                                                                                                                                                                                                          null|null|null|null|null|null|
                             null|
                                                                                                                                                     null|
                                                                                                                                                                                null| null| null|
    Inu111
                                             null| null|
                                                                                          null|
                                                                                                                     null|
                                                                                                                                                                                                                          null|null|null|null|null|
    |null|
                             null|
                                              null| null|
                                                                                          null|
                                                                                                                     null|
                                                                                                                                                     null|
                                                                                                                                                                                null| null| null|
                                                                                                                                                                                                                           null|null|null|null|null|
    |null|
                             null|
                                             null| null|
                                                                                          null|
                                                                                                                     null|
                                                                                                                                                     null|
                                                                                                                                                                                \verb|null|| \verb|null|| \verb|null|| \verb|null|| \verb|null|| \verb|null|| |
     +----+
                                                                                                                    ----+--
                                                                                                                                                    ----+-
    only showing top 10 rows
    Took 4 sec. Last updated by anonymous at July 24 2025, 4:47:26 PM.
     df1_ready.select("AMIE", "anio_lectivo", "Provincia", "Canton").show(10, false)
df2_ready.select("AMIE", "anio_lectivo", "Provincia", "Canton").show(10, false)
df3_clean.select("AMIE").show(10, false)
                                                                                                                                                                                                                                           SPARK JOB FINISHED
                                                                                                                                                                                                                                                                                 AMIE | anio lectivo | Provincia
                                                                                                   | Canton
    |01H01829|2022-2023 Fin|AZUAY
                                                                                                    CUENCA
    |04H00462|2022-2023 Fin|CARCHI
                                                                                                    LESPE TO
    |08H01180|2022-2023 Fin|SANTO DOMINGO DE LOS TSACHILAS|LA CONCORDIA|
    |08H01524|2022-2023 Fin|ESMERALDAS
                                                                                                   ATACAMES
    |08H02161|2022-2023 Fin|ESMERALDAS
                                                                                                    |QUININDE
    |09H06652|2022-2023 Fin|GUAYAS
                                                                                                    |GUAYAQUIL
    |10H00070|2022-2023 Fin|IMBABURA
                                                                                                    IIBARRA
    |11H02033|2022-2023 Fin|LOJA
                                                                                                    SARAGURO
    |12H02246|2022-2023 Fin|LOS RIOS
                                                                                                    |PUEBLOVIEJO
    |17H04382|2022-2023 Fin|PICHINCHA
    only showing top 10 rows
    LAMTE lania lastina | DecuincialConton|
    Took 0 sec. Last updated by anonymous at July 24 2025, 4:48:49 PM.
                                                                                                                                                                                                                                                                   FINISHED
      import org.apache.spark.sql.functions.
      val df1_norm = df1_ready
          .withColumn("AMIE")))
.withColumn("AMIE")))
.withColumn("anio_lectivo", trim(upper(col("anio_lectivo"))))
.withColumn("Provincia", trim(upper(col("Provincia"))))
           .withColumn("Canton", trim(upper(col("Canton"))))
     val df2_norm = df2_ready
.withColumn("AMIE", trim(upper(col("AMIE"))))
.withColumn("anio_lectivo", trim(upper(col("anio_lectivo"))))
.withColumn("Provincia", trim(upper(col("Provincia"))))
.withColumn("Provincia", trim(upper(col("Canton"))))
      val df3 norm = df3 clean
         .withColumn("AMIE", trim(upper(col("AMIE"))))
    import org.apache.spark.sql.functions.
    df1_norm: org.apache.spark.sql.DataFrame = [AMIE: string, anio_lectivo: string ... 4 more fields]
    df2_norm: org.apache.spark.sql.DataFrame = [AMIE: string, anio_lectivo: string ... 4 more fields]
    df3 norm: org.apache.spark.sql.DataFrame = [AMIE: string, ciclo: string ... 7 more fields]
    Took 0 sec. Last updated by anonymous at July 24 2025, 4:50:05 PM
      import org.apache.spark.sql.functions.
                                                                                                                                                                                                                                                                   FINISHED
      val cols_finales = Seq("AMIE", "anio_lectivo", "Provincia", "Canton", "Total_Estudiantes", "Total_Docentes")
      val df3 union = df3 norm
         al dr3_union = dr3_norm
.withColumn("anio_lectivo", lit(null).cast("string"))
.withColumn("Provincia", lit(null).cast("string"))
.withColumn("Canton", lit(null).cast("string"))
.withColumn("Total_Estudiantes", lit(null).cast("int"))
.withColumn("Total_Docentes", lit(null).cast("int"))
           select(cols finales.head. cols finales.tail:
    import org.apache.spark.sql.functions._
    cols_finales: Seq[String] = List(AMIE, anio_lectivo, Provincia, Canton, Total_Estudiantes, Total_Docentes)
    df3_union: org.apache.spark.sql.DataFrame = [AMIE: string, anio_lectivo: string ... 4 more fields]
    Took 0 sec. Last updated by anonymous at July 24 2025, 4:50:51 PM.
      val df final union = df1 norm
                                                                                                                                                                                      SPARK JOB (http://10.0.5.2:4040/jobs/job?id=38) FINISHED
         .select(cols_finales.head, cols_finales.tail: _*)
          .unionByName(
             df2_norm.select(cols_finales.head, cols_finales.tail: *)
          .unionByName(
             df3_union
```

)



Proceso completo de Integración, Limpieza y Carga de Datos Educativos a MySQL usando Apache Zeppelin, Spark y Bash

FINISHED

1. Carga de archivos en Spark

Se cargaron tres archivos CSV desde distintas fuentes oficiales:

- 2_MINEDUC_RegistrosAdministrativos_2022-2023Fin.csv
- registro-administrativo-historico_2009-2024-inicio.csv
- ineval_serestudiante2022_2023_2023diciembre.csv

El objetivo era unificar la información relevante sobre instituciones educativas a nivel nacional.

2. Inspección y Estandarización de Columnas

Proyectorinal

- AMIE (código de institución)
- anio lectivo
- Provincia
- Canton
- Total_Estudiantes
- Total Docentes

Se estandarizaron nombres y se agregaron columnas nulas cuando fue necesario para compatibilidad.

3. Unión y Homologación de DataFrames

Se procesaron los tres DataFrames para garantizar que tuvieran la misma estructura. Se usaron funciones como withColumnRenamed y withColumn en Spark/Scala.

Luego, se unificaron usando unionByName para obtener un solo DataFrame consolidado.

4. Exportación y Limpieza Final en Bash

El DataFrame final se exportó a un CSV único.

Se aplicaron scripts bash para:

- Reemplazar valores vacíos ("" , ,) por NULL .
- Corregir líneas con problemas de formato.
- Eliminar filas basura o vacías.

Ejemplo de comandos utilizados:

```
sed -i 's/,"",/,NULL,/g; s/,"$/NULL/g; s/,,,,NULL,/g; s/,$/,NULL/' resultado_final_union_todo.csv
\verb"awk'!/^""/' resultado\_final\_union\_todo.csv > resultado\_final\_union\_todo\_limpio.csv > resultado\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_union\_todo\_final\_u
sed -i 's/,""$/\,NULL/' resultado_final_union_todo_limpio.csv
El archivo limpio fue movido a /var/lib/mysql-files/ y cargado en una tabla SQL con:
sal
Copiar
Editar
LOAD DATA INFILE '/var/lib/mysql-files/resultado_final_union_todo_limpio.csv'
INTO TABLE union educacion
FIELDS TERMINATED BY ',
OPTIONALLY ENCLOSED BY '"'
LINES TERMINATED BY '\n'
IGNORE 1 LINES
(AMIE, anio_lectivo, Provincia, Canton, Total_Estudiantes, Total_Docentes);
6. Validación y Consultas Iniciales
Se comprobó la carga correcta y la calidad de los datos con consultas SQL, por ejemplo:
SELECT COUNT(*) FROM union_educacion;
SELECT Provincia, COUNT(DISTINCT AMIE) as Instituciones FROM union_educacion GROUP BY Provincia;
SELECT * FROM union_educacion WHERE Total_Estudiantes IS NULL OR Total_Docentes IS NULL LIMIT 10;
```

Took 3 sec. Last updated by anonymous at July 24 2025, 5:18:36 PM.

%md READY