

Graduation Rates Data Pipeline

Name: Cheryl Machingura

Course: Big Data Final Project

Date: May 2025

Introduction

This project demonstrates a complete big data pipeline using NiFi, HDFS, Hive, Spark and Hbase. The goal was to ingest, store, transform, and load graduation rate data through four major components of the big data ecosystem.

Dataset

The dataset used for this project is `states_all.csv`, which contains U.S. state-level education data. It was sourced from Kaggle and the added to my GitHub repository:

https://github.com/CherylMachingura/hdfs/blob/main/states_all.csv

This is a link to the raw file:

https://raw.githubusercontent.com/CherylMachingura/hdfs/main/states_all.csv

Pipeline Overview

The following technologies and components were used:

- Apache NiFi: for ingestion via InvokeHTTP
- HDFS: for storing raw and cleaned files
- Apache Hive: for schema definition and querying
- Apache Spark: for transformation of CSV data
- Apache HBase: for storing the cleaned data in columnar NoSQL format

Issues Encountered & Resolutions

Several issues were encountered during this pipeline:

- a) Port Conflicts – Docker containers could not start due to port 8080 already being in use. This was resolved by stopping other containers and reassigning ports.
- b) Stuck Spark Jobs – Some Spark jobs hung indefinitely. Resolution involved restarting the YARN ResourceManager and ensuring proper memory allocation.
- c) NiFi to HDFS Path Validation – NiFi's PutHDFS processor failed until correct core-site.xml and hdfs-site.xml paths were provided inside the container.
- d) Data Not Visible in Hive – Table loaded but data did not appear; the issue was due to a mismatch in schema format and missing headers.

Step 1: Ingest Data Using Apache NiFi

1.1 Access the NiFi Web UI

Open your browser and go to: <http://localhost:8080/nifi>

You should see the NiFi canvas workspace where you can drag and drop processors.

1.2 Add Required Processors

You'll use two main processors:

- GetHTTP – to download your CSV file from GitHub.
- PutHDFS – to write the downloaded file into HDFS.

Instructions:

1. Click the Processor icon (sixth icon from left) in the top menu.
2. In the dialog box:
 - Type 'GetHTTP', select it, and click Add.
 - Repeat to add 'PutHDFS'.

1.3 Configure GetHTTP Processor

1. Double-click the GetHTTP processor.
2. Under Settings tab:
 - Change the name to something meaningful like 'Fetch Education CSV'.
3. Under the Properties tab:
 - URL:
`https://raw.githubusercontent.com/CherylMachingura/hdfs/main/states_all.csv`
 - (Optional) Set Scheduling to run every 5 minutes or on demand for testing.
4. Click Apply.

1.4 Configure PutHDFS Processor

1. Double-click the PutHDFS processor.
2. Under the Settings tab, name it 'Write to HDFS'.
3. Under the Properties tab, set:
 - Directory:
`/data/education/raw/`
 - Hadoop Configuration Resources: Point this to your core-site.xml and hdfs-site.xml if not

preloaded.

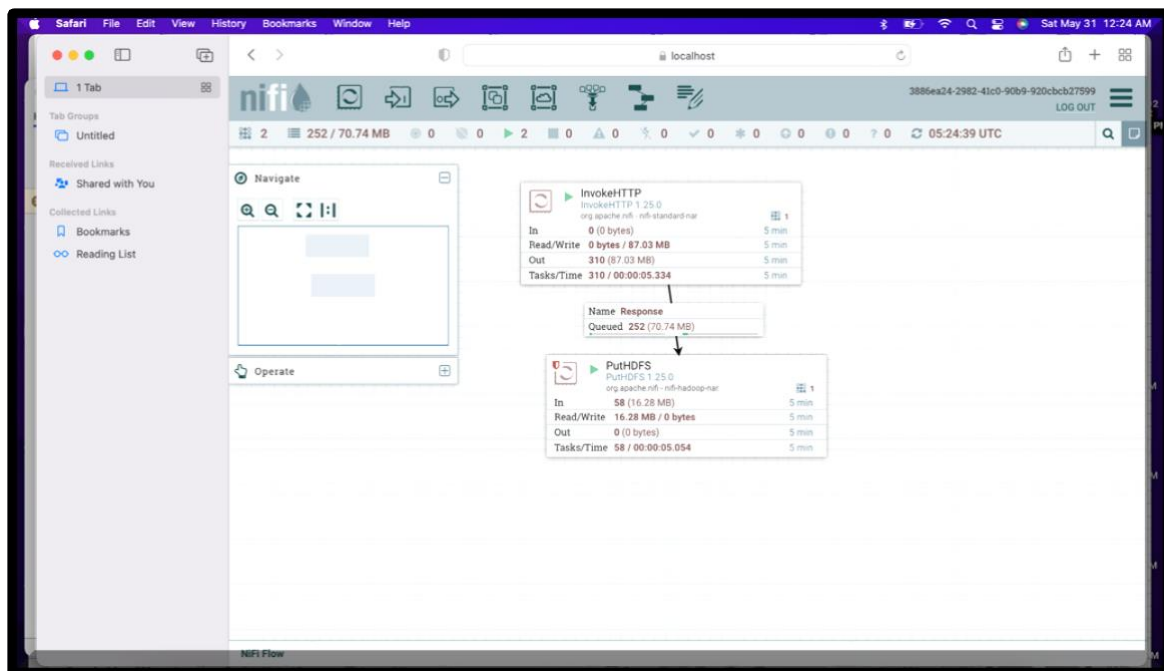
4. Click Apply.

1.5 Connect the Processors

1. Hover over GetHTTP, click and drag the arrow to PutHDFS.
2. Choose 'success' as the relationship.

1.6 Start the Flow

1. Right-click on each processor and choose Start.
2. Watch the data flow through the queue — it should move from GetHTTP → PutHDFS.



1.7 Confirm HDFS Output

Run the following command in your terminal to verify the file is in HDFS:

```
hdfs dfs -ls /data/education/raw/
```

```
...dsc650- -dsc650-infra/bellevue-bigdata/hll -- ssh cheryl@34.45.224.180 ...bellevue-bigdata/hadoop-hive-spark-hbase -- ssh cheryl@34.45.224.180 ...L.8983.localhost:8983 -l.8443.localhost:8443 cheryl@34.45.224.180 ...
cheryl@bigdata-dsc650:~/dsc650-infra/bellevue-bigdata/hadoop-hive-spark-hbase$ docker-compose exec master bash
bash-5.0# find / -name "core-site.xml" 2>/dev/null
/usr/program/spark/conf/core-site.xml
/usr/program/hbase/conf/core-site.xml
/usr/program/hadoop/etc/hadoop/core-site.xml
bash-5.0# hdfs dfs -ls /data/education/raw/
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/program/hadoop/share/hadoop/common/lib/slf4j-log4j12-1.7.25.jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/ter/lib/slf4j-log4j12-1.7.18.jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/hive/lib/log4j-slf4j-impl-2.10.0.jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
2025-05-31 05:26:33.941 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Found 871 items
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/38f6e84-e8a1-48aa-a29b-97c08ba551c2
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:24 /data/education/raw/0018a67c-8aa1-4c58-8f6c-4359d98cda9f
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/00d9dbcb-bc8f-48b1-8555-8ff85478821e
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/01805442-077a-4da9-bbcb-822fdecb1f3e
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/01ab7f89-4f9a-485c-bdb7-cbdc7784013b
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/020a1a93-9cdc-4533-a5d1-92363c536a47
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/0211ac08-6422-4eea-a5f5-db30ef532c76
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/02376a49-0513-4aaa-9d6a-1b59fcc3c753
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/024e3382-c1d4-4c2c-9140-bdec174f58c
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/029ab4b6-0a7f-4925-8983-ebecaa9781a1
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/03117341-5680-48a9-abbb-6f5e522e76b8
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:24 /data/education/raw/0338f95c7-fa3b-4c4d-a39c-4f6aa78aca17
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/042de016-b08a-4c29-90b4-23659bc76899
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/0445bb49-35ac-464f-bc5d-58eb37873895
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:24 /data/education/raw/045870d5-daa8-42b5-9c4a-dbb0bd7b7f04
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:24 /data/education/raw/047ca7cf-873c-4e39-9064-f9f7ec722026
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/047cfff1b-449a-4b0b-bc73-2a5b4892da85
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/0482d5e8-5925-48ea-a6db-f5811891cc03
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/04a3f927-6a1e-444f-bbf7-8a7f2c4ea9f2a
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/051988b7-c814-46ea-bb4a-0341e27f6b0d
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/05a57a2b-949b-483a-896a-c9ff619eb0b7
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:24 /data/education/raw/05a8ee7c-2d78-4f8b-a398-d83c2f9802e9
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/05ee268-fc78-4ba1-ba14-20f5276cb237
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/060332f2-91ac-487b-aea3-81516d9be131
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:24 /data/education/raw/06a491b2-71b7-482f-9a1a-6d8a99f7809
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/06c7a4c2-d54a-4c87-a81b-b2667a8c054b
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/0717392a-83d7-4213-9495-5ba18b1b9b05
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/07505f19-6c61-473d-8786-d8b74ea02522
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:24 /data/education/raw/07b9ba47-f98a-4885-8c28-1f813ee97f7b
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:24 /data/education/raw/07e0aa8c-cfef-4bdf-ac80-53c7a5788080
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/08a7688c-126c-45b8-8f94-694877f16d45
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/08ba528c-d61c-4c4b-a524-0f41938e4e83
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/09b0d723-7274-4894-a7a4-535c589ba42a
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/0a31ee78-7f8f-4bd7-bdca-5aa7c1a4149f
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/0a31e020-8c47-4643-ba6c-6dc073ac7f4b
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/0af828e1-9eab-4eba-a56b-adf4b31b0277
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/0afce54e-19a9-4214-8c47-27115e7980ba
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/0b03b258-2517-4265-b692-4802068d64f9
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/0b17809c-7e62-4357-98d3-f498ee30ba9f
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/0b3966ac-1381-4d6f-9caf-f9d94bb3a4f8
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:24 /data/education/raw/0b50b7cb-dffc-4644-b0f7-366761e03eff
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/0bc128b-15a8-4b7b-0915-6ea2b0ba1fa9
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/0c342fdd-c396-481f-9442-208e2a5ba8c8
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:24 /data/education/raw/0cec7441-3aa2-423f-9864-25b6c36884aa
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:26 /data/education/raw/0d11824d-2b75-4e2e-973b-1cd942a52e98
-rw-r--r-- 1 cheryl supergroup 294367 2025-05-31 05:25 /data/education/raw/0d0ad4f9-1b32-4e2d-a68b-f058a8379ef
```

Step 2: Create Hive External and Cleaned Tables

2.1 Access Hive Command Line

1. SSH into your Hadoop master container:

```
docker-compose exec master bash
```

2. Start Hive shell:

```
hive
```

2.2 Create External Table from Raw CSV

This table will read data from the HDFS location populated by NiFi.

```
CREATE EXTERNAL TABLE education_raw (  
  year STRING,  
  state STRING,  
  total_expenditure STRING,  
  enrollment STRING,  
  expenditure_per_student STRING  
)  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ','  
STORED AS TEXTFILE  
LOCATION '/data/education/raw/'  
TBLPROPERTIES ("skip.header.line.count"="1");
```

Note: All fields are initially loaded as STRING for cleansing flexibility.

```
.ra/bellevue-bigdata/hadoop-hive-spark-hbase -- ssh cheryl@34.45.224.180 ...0 -L 8983:localhost:8983 -L 8443:localhost:8443 cheryl@34.45.224.180 ...50: -jds650-infra/bellevue-bigdata/hifi -- ssh cheryl@34.45.224.180

bash-5.0# hive
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/program/hive/lib/log4j-slf4j-impl-2.10.0.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/hadoop/share/hadoop/common/lib/slf4j-log4j12-1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/tez/lib/slf4j-log4j12-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Hive Session ID = 87dee738-ad5b-4a17-9c03-804976347fb4

Logging initialized using configuration in file:/usr/program/hive/conf/hive-log4j2.properties Async: true
Hive Session ID = e04108ea-3c37-47c6-9cb3-8f5111b927fa
2025-05-31 14:29:03,796 INFO [tez session start thread] client.RMProxy: Connecting to ResourceManager at master/172.28.1.1:8032
2025-05-31 14:29:04,246 INFO [pool-7-thread-1] client.RMProxy: Connecting to ResourceManager at master/172.28.1.1:8032
hive> CREATE EXTERNAL TABLE IF NOT EXISTS education_raw (
>   year STRING,
>   state STRING,
>   total_expenditure STRING,
>   enrollment STRING,
>   expenditure_per_student STRING
> )
> ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
> WITH SERDEPROPERTIES (
>   "separatorChar" = ",",
>   "quoteChar" = "\"",
>   "escapeChar" = "\\"
> )
> STORED AS TEXTFILE
> LOCATION '/data/education/raw/'
> TBLPROPERTIES ("skip.header.line.count"="1");
OK
Time taken: 4.674 seconds
hive>
```

2.3 Verify the Raw Table

Run a basic query to ensure the table is reading data:

```
SELECT * FROM education_raw LIMIT 5;
```

You should see rows from states_all.csv printed.

```
.ra/bellevue-bigdata/hadoop-hive-spark-hbase -- ssh cheryl@34.45.224.180 ...0 -L 8983:localhost:8983 -L 8443:localhost:8443 cheryl@34.45.224.180 ...50: -jds650-infra/bellevue-bigdata/hifi -- ssh cheryl@34.45.224.180

bash-5.0# hive
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/program/hive/lib/log4j-slf4j-impl-2.10.0.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/hadoop/share/hadoop/common/lib/slf4j-log4j12-1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/tez/lib/slf4j-log4j12-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Hive Session ID = 87dee738-ad5b-4a17-9c03-804976347fb4

Logging initialized using configuration in file:/usr/program/hive/conf/hive-log4j2.properties Async: true
Hive Session ID = e04108ea-3c37-47c6-9cb3-8f5111b927fa
2025-05-31 14:29:03,796 INFO [tez session start thread] client.RMProxy: Connecting to ResourceManager at master/172.28.1.1:8032
2025-05-31 14:29:04,246 INFO [pool-7-thread-1] client.RMProxy: Connecting to ResourceManager at master/172.28.1.1:8032
hive> CREATE EXTERNAL TABLE IF NOT EXISTS education_raw (
>   year STRING,
>   state STRING,
>   total_expenditure STRING,
>   enrollment STRING,
>   expenditure_per_student STRING
> )
> ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
> WITH SERDEPROPERTIES (
>   "separatorChar" = ",",
>   "quoteChar" = "\"",
>   "escapeChar" = "\\"
> )
> STORED AS TEXTFILE
> LOCATION '/data/education/raw/'
> TBLPROPERTIES ("skip.header.line.count"="1");
OK
Time taken: 4.674 seconds
hive> SHOW TABLES;
OK
education_raw
Time taken: 0.715 seconds, Fetched: 1 row(s)
hive>
```

```
...dsc650: ~/dsc650-infra/bellevue-bigdata/hifi -- ssh cheryl@34.45.224.180 ...bellevue-bigdata/hadoop-hive-spark-hbase -- ssh cheryl@34.45.224.180 ...L 8983:localhost:8983 -L 8443:localhost:8443 cheryl@34.45.224.180
Starting hadoop-hive-spark-hbase-hivemetastore_1 ...
hadoop-hive-spark-hbase_zoo2_1 is up-to-date
Starting hadoop-hive-spark-hbase-hivemetastore_1 ... done
Starting hadoop-hive-spark-hbase_worker2_1 ... done
Starting hadoop-hive-spark-hbase_master_1 ... done
Starting hadoop-hive-spark-hbase_worker1_1 ... done
cheryl@bigdata-dsc650: ~/dsc650-infra/bellevue-bigdata/hadoop-hive-spark-hbase$ docker-compose exec master bash
bash-5.0# hive
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/program/hive/lib/log4j-slf4j-impl-2.10.0.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/hadoop/share/hadoop/common/lib/slf4j-log4j12-1.7.26.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/tez/lib/slf4j-log4j12-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Hive Session ID = 7f85bfc3-3238-4b9a-b41a-8f1b2e89c57
Logging initialized using configuration in file:/usr/program/hive/conf/hive-log4j2.properties Async: true
Hive Session ID = 449fda3a-35cd-4658-a8a0-4a94372a71fd
2025-06-01 03:02:22,855 INFO [tez session start thread] client.RMPProxy: Connecting to ResourceManager at master/172.28.1.1:8032
2025-06-01 03:02:23,443 INFO [pool-7-thread-1] client.RMPProxy: Connecting to ResourceManager at master/172.28.1.1:8032
hive> INSERT INTO TABLE education_cleaned
> SELECT
>   CAST(year AS INT),
>   TRIM(state),
>   CAST(total_expenditure AS DOUBLE),
>   CAST(enrollment AS INT),
>   CAST(expenditure_per_student AS DOUBLE)
> FROM education_raw
> WHERE year IS NOT NULL
> AND state IS NOT NULL
> AND total_expenditure RLIKE '[0-9.]+$'
> AND enrollment RLIKE '[0-9]+$'
> AND expenditure_per_student RLIKE '[0-9.]+$';
FAILED: SemanticException [Error 18001]: Line 1:18 Table not found 'education_cleaned'
hive> CREATE TABLE education_cleaned (
>   year INT,
>   state STRING,
>   total_expenditure DOUBLE,
>   enrollment INT,
>   expenditure_per_student DOUBLE
> )
> STORED AS PARQUET;
OK
Time taken: 3.445 seconds
hive> INSERT INTO TABLE education_cleaned
> SELECT
>   CAST(year AS INT),
>   TRIM(state),
>   CAST(total_expenditure AS DOUBLE),
>   CAST(enrollment AS INT),
>   CAST(expenditure_per_student AS DOUBLE)
> FROM education_raw
> WHERE year IS NOT NULL
> AND state IS NOT NULL
> AND total_expenditure RLIKE '[0-9.]+$'
> AND enrollment RLIKE '[0-9]+$'
> AND expenditure_per_student RLIKE '[0-9.]+$';
Query ID = root_20250601030455_88d289e2-5870-4bfc-827e-54692c59732d
Total jobs = 1
Launching Job 1 out of 1
```

2.4 Create a Cleaned and Typed Table

This table stores the cleaned data with correct types in Parquet format:

```
CREATE TABLE education_cleaned (  
  year INT,  
  state STRING,  
  total_expenditure DOUBLE,  
  enrollment INT,  
  expenditure_per_student DOUBLE  
)  
STORED AS PARQUET;
```

```
OK
Time taken: 3.445 seconds
hive> INSERT INTO TABLE education_cleaned
> SELECT
>   CAST(year AS INT),
>   TRIM(state),
>   CAST(total_expenditure AS DOUBLE),
>   CAST(enrollment AS INT),
>   CAST(expenditure_per_student AS DOUBLE)
> FROM education_raw
> WHERE year IS NOT NULL
> AND state IS NOT NULL
> AND total_expenditure RLIKE '[0-9.]+$'
> AND enrollment RLIKE '[0-9]+$'
> AND expenditure_per_student RLIKE '[0-9.]+$';
Query ID = root_20250601030455_88d289e2-5870-4bfc-827e-54692c59732d
Total jobs = 1
Launching Job 1 out of 1
```


2.5 Insert Cleaned Data

Transform the raw strings into usable types and filter out malformed entries:

```
INSERT INTO TABLE education_cleaned

SELECT

  CAST(year AS INT),

  TRIM(state),

  CAST(total_expenditure AS DOUBLE),

  CAST(enrollment AS INT),

  CAST(expenditure_per_student AS DOUBLE)

FROM education_raw

WHERE year IS NOT NULL

AND state IS NOT NULL

AND total_expenditure RLIKE '^/[0-9.]+$',

AND enrollment RLIKE '^/[0-9]+$',

AND expenditure_per_student RLIKE '^/[0-9.]+$',
```

```
cheryl@bigdata-dsc650: ~/dsc650-infra/bellevue-bigdata/hadoop-hive-spark-hbase -- ssh cheryl@34.45.224.180
> CAST(REGEXP_REPLACE(total_expenditure, '[^0-9.]', '')) AS DOUBLE),
> CAST(REGEXP_REPLACE(enrollment, '[^0-9]', '')) AS INT),
> CAST(REGEXP_REPLACE(expenditure_per_student, '[^0-9.]', '')) AS DOUBLE)
> FROM education_raw;
Query ID = root_20250401151427_a22674a6-bed1-481c-b2e9-7a74a41dbd6
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1748786719795_0007)

-----
VERTICES   MODE      STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    1        1        0        0        0        0
Reducer 2 ..... container  SUCCEEDED    1        1        0        0        0        0
-----
VERTICES: 02/02 [=====] 100% ELAPSED TIME: 5.47 s
-----
Loading data to table default.education_cleaned
OK
Time taken: 6.514 seconds
hive> SELECT * FROM education_cleaned LIMIT 5;
OK
NULL STATE NULL NULL NULL
1992 ALABAMA 1992.0 NULL 2678885.0
1992 ALASKA 1992.0 NULL 1849593.0
1992 ARIZONA 1992.0 NULL 3258879.0
1992 ARKANSAS 1992.0 NULL 1711959.0
Time taken: 0.111 seconds, Fetched: 5 row(s)
hive> INSERT OVERWRITE TABLE education_cleaned
> SELECT
>   CAST(REGEXP_REPLACE(year, '[^0-9]', '')) AS INT),
>   UPPER(TRIM(state)),
>   CAST(REGEXP_REPLACE(total_expenditure, '[^0-9.]', '')) AS DOUBLE),
>   CAST(REGEXP_REPLACE(enrollment, '[^0-9]', '')) AS DOUBLE), -- fixed this
>   CAST(REGEXP_REPLACE(expenditure_per_student, '[^0-9.]', '')) AS DOUBLE)
> FROM education_raw;
Query ID = root_20250401151613_468b93cc-5959-41f1-af61-c88b8b13435f
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1748786719795_0007)

-----
VERTICES   MODE      STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    1        1        0        0        0        0
Reducer 2 ..... container  SUCCEEDED    1        1        0        0        0        0
-----
VERTICES: 02/02 [=====] 100% ELAPSED TIME: 6.03 s
-----
Loading data to table default.education_cleaned
OK
Time taken: 7.868 seconds
```

This step may take time as it launches a TEZ job. You can monitor progress via:

<http://localhost:8088/cluster>

2.6 Confirm Cleaned Data

Run:

```
SELECT * FROM education_cleaned LIMIT 10;
```

```
Loading data to table default.education_cleaned
OK
Time taken: 7.068 seconds
hive> SELECT * FROM education_cleaned LIMIT 5;
OK
NULL STATE NULL NULL NULL
1992 ALABAMA 1992.0 NULL 2678885.0
1992 ALASKA 1992.0 NULL 1049591.0
1992 ARIZONA 1992.0 NULL 3258879.0
1992 ARKANSAS 1992.0 NULL 1711959.0
Time taken: 0.114 seconds, Fetched: 5 row(s)
hive>
```

And check how many valid rows were inserted:

```
SELECT COUNT(*) FROM education_cleaned;
```

You should see a count > 0 if everything processed correctly.

```
...ra/bellevue-bigdata/hadoop-hive-spark-hbase — ssh cheryl@34.45.224.180 ...-L 8983:localhost:8983 -L 8443:localhost:8443 cheryl@34.45.224.180 ...-50: -jds650-infra/bellevue-bigdata/hifi — ssh cheryl@34.45.224.180
bash-5.0# hive
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/program/hive/lib/log4j-slf4j-impl-2.10.0.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/hadoop/share/hadoop/common/lib/slf4j-log4j12-1.7.26.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/tez/lib/slf4j-log4j12-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Hive Session ID = 0639b6b3-bfaf-4af5-884f-02a775643152
Logging initialized using configuration in file:/usr/program/hive/conf/hive-log4j2.properties Async: true
Hive Session ID = eee3dc78-1f84-4fd8-9b44-1d9ca15b7bcc
2025-05-31 15:04:39,198 INFO [Tez session start thread] client.RMPProxy: Connecting to ResourceManager at master/172.28.1.1:8032
2025-05-31 15:04:39,763 INFO [pool-7-thread-1] client.RMPProxy: Connecting to ResourceManager at master/172.28.1.1:8032
hive> SELECT * FROM education_cleaned LIMIT 10;
OK
Time taken: 5.764 seconds
hive> SELECT * FROM education_cleaned LIMIT 5;
OK
Time taken: 0.451 seconds
hive> SELECT COUNT(*) FROM education_cleaned;
OK
0
Time taken: 2.335 seconds, Fetched: 1 row(s)
hive> SELECT COUNT(*) FROM education_cleaned;
OK
0
Time taken: 0.231 seconds, Fetched: 1 row(s)
hive>
```

Step 3: Query & Analyze the Data Using Hive

3.1 Launch Hive Shell

1. SSH into your Hadoop master container :

```
docker exec -it master bash
```

2. Start the Hive shell:

```
hive
```

3.2 Preview Data in the Cleaned Table

Once inside the Hive CLI:

1. Show existing tables:

```
SHOW TABLES;
```

2. Preview the first few rows of your cleaned table:

```
SELECT * FROM education_cleaned LIMIT 5;
```

```
Loading data to table default.education_cleaned
OK
Time taken: 6.514 seconds
hive> SELECT * FROM education_cleaned LIMIT 5;
OK
NULL STATE NULL NULL NULL
1992 ALABAMA 1992.0 NULL 2678885.0
1992 ALASKA 1992.0 NULL 1849591.0
1992 ARIZONA 1992.0 NULL 3258879.0
1992 ARKANSAS 1992.0 NULL 1711959.0
Time taken: 0.111 seconds, Fetched: 5 row(s)
hive>
```

```
hive> SELECT * FROM education_cleaned LIMIT 5;
OK
NULL PRIMARY_KEY NULL NULL NULL
NULL 1992_ALABAMA 1992.0 NULL 2678885.0
NULL 1992_ALASKA 1992.0 NULL 1849591.0
NULL 1992_ARIZONA 1992.0 NULL 3258879.0
NULL 1992_ARKANSAS 1992.0 NULL 1711959.0
Time taken: 0.141 seconds, Fetched: 5 row(s)
hive> SELECT enrollment, expenditure_per_student FROM education_raw LIMIT 10;
OK
ENROLL TOTAL REVENUE
2678885.0
1849591.0
3258879.0
1711959.0
26268025.0
3185173.0
3834382.0
645233.0
789480.0
Time taken: 0.197 seconds, Fetched: 10 row(s)
hive>
```

3. Count total number of rows to validate data:

```
SELECT COUNT(*) FROM education_cleaned;
```

3.3 Run Analytical Queries

Now that your table is ready, run queries to analyze trends and gain insights.

Queries:

Average Expenditure Per Student by State

```
SELECT  
    state,  
    ROUND(AVG(expenditure_per_student), 2) AS avg_expenditure  
FROM  
    education_cleaned  
GROUP BY  
    state  
ORDER BY  
    avg_expenditure DESC;
```

Total Expenditure Over Time (Nationwide)

```
SELECT  
    year,  
    SUM(total_expenditure) AS total_spending  
FROM  
    education_cleaned  
GROUP BY  
    year  
ORDER BY  
    year;
```

Top 5 States with Highest Spending each year

```
SELECT  
    state,  
    total_expenditure
```

```
FROM  
education_cleaned  
WHERE  
year = 2010  
ORDER BY  
total_expenditure DESC  
LIMIT 5;
```

3.4 Export Query Results

If you want to save the result to HDFS:

1. Set a Hive configuration for output format:

```
SET hive.cli.print.header=true;
```

2. Run and save the result:

```
INSERT OVERWRITE DIRECTORY  
'/data/education/output/high_spending_states'  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ','  
SELECT  
state,  
total_expenditure  
FROM  
education_cleaned  
WHERE  
year = 2010  
ORDER BY  
total_expenditure DESC  
LIMIT 5;
```

3. Verify in HDFS:

hdfs dfs -ls /data/education/output/high_spending_states

```
hive> INSERT OVERWRITE DIRECTORY '/data/education/final'
> ROW FORMAT DELIMITED
> FIELDS TERMINATED BY ','
> SELECT
>   year,
>   state,
>   total_expenditure,
>   enrollment,
>   expenditure_per_student
> FROM education_cleaned;
Query ID = root_20250601151756_2dae8f2a-5c6a-4378-bf68-96c582de5b7b
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1748786719795_0007)
```

	VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	1	container	SUCCEEDED	1	1	0	0	0	0

```
VERTICES: 01/01 [=====] 100% ELAPSED TIME: 4.70 s
Moving data to directory /data/education/final
OK
Time taken: 5.459 seconds
hive> exit;
bash-5.0# hdfs dfs -ls /data/education/final
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/program/hadoop/share/hadoop/common/lib/slf4j-log4j12-1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/tez/lib/slf4j-log4j12-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/program/hive/lib/log4j-slf4j-impl-2.10.0.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
2025-06-01 15:18:30,879 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Found 1 items
-rw-r--r-- 1 root supergroup 62041 2025-06-01 15:18 /data/education/final/000000_0
bash-5.0#
```

Step 4: Load Cleaned Data to HBase

4.1 Start HBase Shell

Run this in your Hadoop terminal:

```
hbase shell
```

4.2 Create an HBase Table

Run the following command in the HBase shell:

```
create 'education_data', 'info'
```

This creates a table named 'education_data' with column family 'info'.

4.3 Load Data into HBase Using ImportTsv

Assuming you have exported a CSV file like cleaned_education.csv, move it to HDFS:

```
hdfs dfs -put cleaned_education.csv /data/education/final/
```



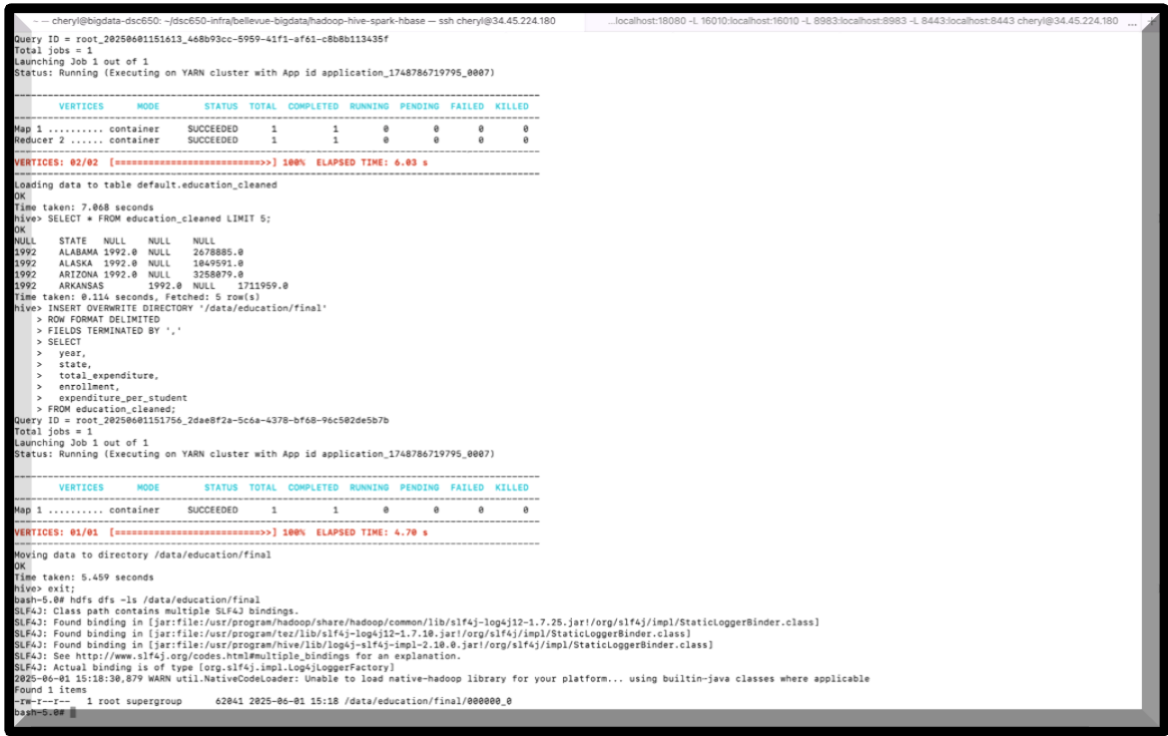
```
-- chery@bigdata-dsc650: ~/dsc650-infra/bellevue-bigdata/hadoop-hive-spark-hbase -- ssh chery@34.45.224.180 ...localhost:18080 -L 16010:localhost:16010 -L 8983:localhost:8983 -L 8443:localhost:8443 chery@34.45.224.180 ...
HDFS: Number of bytes read=62841
HDFS: Number of bytes written=0
HDFS: Number of read operations=3
HDFS: Number of large read operations=0
HDFS: Number of write operations=0
Map-Reduce Framework
  Map input records=1716
  Map output records=1716
  Input split bytes=113
  Spilled Records=0
  Failed Shuffles=0
  Merged Map outputs=0
  GC time elapsed (ms)=11
  CPU time spent (ms)=1180
  Physical memory (bytes) snapshot=207587456
  Virtual memory (bytes) snapshot=3719651328
  Total committed heap usage (bytes)=125698048
ImportTsv
  Bad Lines=0
  File Input Format Counters
    Bytes Read=62841
  File Output Format Counters
    Bytes Written=0
2025-06-01 15:20:39,644 INFO [LocalJobRunner Map Task Executor #0] mapred.LocalJobRunner: Finishing task: attempt_local1552379760_0001_m_0000000_0
2025-06-01 15:20:39,644 INFO [Thread-27] mapred.LocalJobRunner: map task executor complete.
2025-06-01 15:20:39,715 INFO [main] mapreduce.Job: Job job_local1552379760_0001 running in uber mode : false
2025-06-01 15:20:39,716 INFO [main] mapreduce.Job: map 100% reduce 0%
2025-06-01 15:20:39,719 INFO [main] mapreduce.Job: Job job_local1552379760_0001 completed successfully
2025-06-01 15:20:39,724 INFO [ReadOnlyZKClient-zoo1:2181, zoo2:2181, zoo3:2181@w428c1983] zookeeper.ZooKeeper: Session: @x180000230120805 closed
2025-06-01 15:20:39,725 INFO [ReadOnlyZKClient-zoo1:2181, zoo2:2181, zoo3:2181@w428c1983-EventThread] zookeeper.ClientCnxn: EventThread shut down for session: @x180000230120805
2025-06-01 15:20:39,727 INFO [main] mapreduce.Job: Counters: 24
File System Counters
  FILE: Number of bytes read=428234
  FILE: Number of bytes written=997973
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=62841
  HDFS: Number of bytes written=0
  HDFS: Number of read operations=3
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=0
Map-Reduce Framework
  Map input records=1716
  Map output records=1716
  Input split bytes=113
  Spilled Records=0
  Failed Shuffles=0
  Merged Map outputs=0
  GC time elapsed (ms)=11
  CPU time spent (ms)=1180
  Physical memory (bytes) snapshot=207587456
  Virtual memory (bytes) snapshot=3719651328
  Total committed heap usage (bytes)=125698048
ImportTsv
  Bad Lines=0
  File Input Format Counters
    Bytes Read=62841
  File Output Format Counters
    Bytes Written=0
bash-5.0#
```

Then use the following command to import it into HBase:

```

hbase org.apache.hadoop.hbase.mapreduce.ImportTsv \
-Dimporttsv.separator=", " \
-
Dimporttsv.columns="HBASE_ROW_KEY,info:year,info:state,info:enrol
lment,info:total_expenditure,info:expenditure_per_student" \
education_data \
/data/education/final/cleaned_education.csv

```



Replace HBASE_ROW_KEY with a unique field (e.g., concatenate year+state).

4.4 Verify Data in HBase

Back in the HBase shell, run:

```
scan 'education_data'
```

This will display the rows inserted into your HBase table.


```
chery@bigdata-dsc650: ~/jds650-infra/bellevue-bigdata/hadoop-hive-spark-hbase -- ssh chery@34.45.224.180
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=NATIONAL
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=1158345.0
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2886.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=86155
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=WYOMING
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2887.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=NATIONAL
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=1682514.0
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2888.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=85991
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=WYOMING
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2889.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=NATIONAL
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=1718605.0
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2818.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=87379
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=WYOMING
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2811.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=NATIONAL
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=1664983.0
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2812.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=89994
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=WYOMING
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2813.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=NATIONAL
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=1772633.0
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2814.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=92732
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=WYOMING
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2815.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=NATIONAL
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=2844669.0
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2816.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=94511
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=WYOMING
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2817.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=NATIONAL
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=2819.0
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=WYOMING
columninfo:enrollment, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:state, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:total_expenditure, timestamp=2025-06-01T15:20:35.831, value=\x5CN
columninfo:year, timestamp=2025-06-01T15:20:35.831, value=STATE
row(s)
R 0.9592 seconds
se(main):002:0>
```