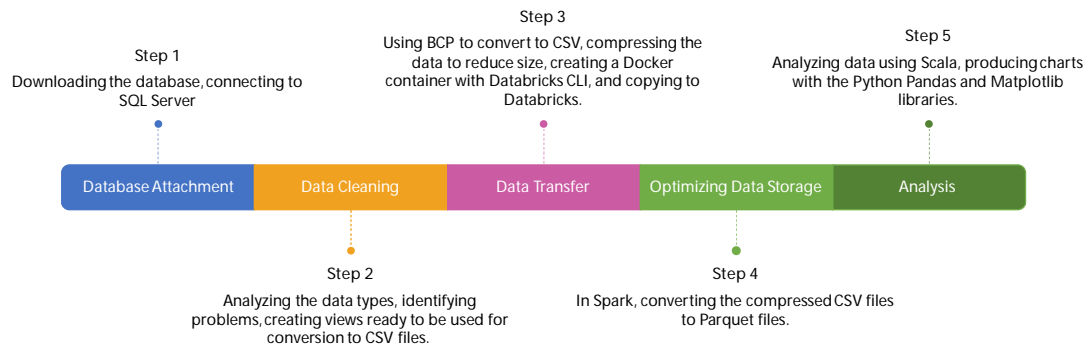# Analyzing the StackOverflow2013 Database from 2008 to 2013

Piotr Slusarczyk

The dataset under examination is the 'StackOverflow2013' database, which spans from the years 2008 to 2013. This comprehensive dataset provides a rich opportunity for various analytical explorations and can offer valuable insights into the trends, patterns, and dynamics of the Stack Overflow community during its early years.

25''

## Steps



The slide shows how the whole project looked like.

I began by downloading the database and establishing a connection to SQL Server, by attaching the downloaded database, laying the groundwork for the subsequent stages.

I then conducted data cleaning, ensuring the data types were correct and identifying any issues. This step was crucial for proper conversion to CSV.

Using BCP from command line, I converted the data into a CSV format.

I compressed the data into gzip format to make it smaller, which saves time when sending the data to Databricks.

To streamline the workflow, I set up a Docker container with Databricks CLI, simplifying the transfer of everything to Databricks.

Moving forward, I used Spark to convert the compressed CSV files into the Parquet format.

Finally, using Scala, I analyzed the data, and with Python's Pandas and Matplotlib libraries, I created charts that clearly depict the insights I've gleaned.

Each step was pivotal in transforming the data to successfully complete the project.

80''

Data Source

How to Download the Stack Overflow Database

In this part of the presentation, I'm focusing on the data source that I used for the analysis. I have chosen the Stack Overflow database for a few key reasons:

Firstly, this database is derived from the public Stack Overflow data export, which is a rich dataset used widely in the tech community for analysis and training. It's not just a simple dataset; it reflects real-world distributions of numbers, dates, and strings, making it an excellent candidate for realistic data analysis scenarios.

For my purposes, I've selected the medium option – the 50GB StackOverflow2013 database, which expands from a 10GB compressed file. This particular dataset contains data from 2008 to 2013.

The data comes in the form of .mdf files, which are the primary database file types used by SQL Server. This makes it convenient as they can be attached directly to SQL Server without any additional conversion, saving time and simplifying the setup process.

68''

# Data Model



On this slide, you can see the visual representation of the StackOverflow2013 Database's data model. This model is composed of nine tables, which together capture the essence of the Stack Overflow community's interactions.

The 'Posts' table is central to this model, linked directly to 'PostTypes', which is composed of two main types of entries: Questions and Answers.

Users are the lifeblood of Stack Overflow, and here you can see their interactions with the system. A user can own a post, meaning they've contributed content to the community, or they can answer posts, providing solutions to questions posed by others.

The 'Comments' table allows for discourse on these posts, enabling users to seek clarification, discuss content, and offer brief insights.

In the 'Votes' table, the community's feedback on posts is captured, where users can vote to signify the usefulness and accuracy of the content provided. It's connected to 'VoteTypes', which defines the nature of these votes.

The 'Users' table holds profiles, reflecting the personal and professional backgrounds of community members, as well as their contributions and reputation within the platform. This table also connects to 'Badges', which are earned by users as a form of recognition for their contributions.

User engagement is further documented through 'PostLinks' and 'LinkTypes', which trace the relationships between different posts, such as duplicates or related content.

Together, these tables not only store data but also tell the story of how knowledge is built, shared, and valued in one of the largest technical communities online.

98''

## Data Cleaning

The main problem: HTML code in the fields

Results | Messages

| | Id | AcceptedAnswerId | AnswerCount | Body | ClosedDate | CommentCount | CommunityOwnedDate |
|---|---|---|---|---|---|---|---|
| 1 | 4 | 7 | 13 | \<p>I want to use a track-bar to change a form's opacity.\</p>  \<p>This is my code:\</p>  \<pre>\<code>decimal trans = trackBar1.Value / 5000; this.Opacity = trans; \</code>\</pre> ... | NULL | 1 | 2012-10-31 16:42:47.213 |
| 2 | 6 | 31 | 5 | \<p>I have an absolutely positioned \<code>div\</code> containing several children, one of which is a relatively positioned \<code>div\</code>. When I use a \<strong>percentage-b... | NULL | 0 | NULL |
| 3 | 7 | 0 | 0 | \<p>An explicit cast to double like this isn't necessary \</p>  \<pre>\<code>double trans = (double) trackBar1.Value / 5000.0; \</code>\</pre>  \<p>Identifying the constant as \<code... | NULL | 0 | NULL |
| 4 | 9 | 1404 | 64 | \<p>Given a \<code>DateTime\</code> representing a person's birthday, how do I calculate their age in years? \</p> | NULL | 7 | 2011-08-16 19:40:43.080 |
| 5 | 11 | 1248 | 35 | \<p>Given a specific \<code>DateTime\</code> value, how do I display relative time, like:\</p>  \<ul>  \<li>2 hours ago\</li>  \<li>3 days ago\</li>  \<li>a month ago\</li>  \</ul> | NULL | 3 | 2009-09-04 13:15:59.820 |
| 6 | 12 | 0 | 0 | \<p>Here's how I do it\</p>  \<pre class="lang-csharp prettyprint-override">\<code>var ts = new TimeSpan(DateTime.UtcNow.Ticks - dt.Ticks); double delta = Math.Abs(ts.TotalSec... | NULL | 11 | 2009-09-04 13:15:59.820 |
| 7 | 13 | 0 | 25 | \<p>Is there any standard way for a Web Server to be able to determine a user's timezone within a web page? \</p>  \<p>Perhaps from an HTTP header or part of the user-agent ... | NULL | 6 | NULL |

Solution: Creating the database views – removing carriage return (\r) and newline (\n) characters and adding quotes.

```sql
CREATE OR ALTER VIEW [dbo].[vPosts] AS
SELECT
     [Id]
    ,[AcceptedAnswerId]
    ,[AnswerCount]
    ,'"' + REPLACE(REPLACE(REPLACE(Body, CHAR(13), ''), CHAR(10), ''), '"', ' ') + '"' Body
    ,[ClosedDate]
    ,[CommentCount]
    ,[CommunityOwnedDate]
    ,[CreationDate]
    ,[FavoriteCount]
    ,[LastActivityDate]
    ,[LastEditDate]
    ,'"' + REPLACE(REPLACE(REPLACE(LastEditorDisplayName, CHAR(13), ''), CHAR(10), ''), '"', ' ') + '"' LastEditorDisplayName
    ,[LastEditorUserId]
    ,[OwnerUserId]
    ,[ParentId]
    ,[PostTypeId]
    ,[Score]
    ,'"' + REPLACE(REPLACE(REPLACE(Tags, CHAR(13), ''), CHAR(10), ''), '"', ' ') + '"' Tags
    ,'"' + REPLACE(REPLACE(REPLACE(Title, CHAR(13), ''), CHAR(10), ''), '"', ' ') + '"' Title
    ,[ViewCount]
FROM [StackOverflow2013].[dbo].[Posts]
```

On this slide, I present the crucial step in data preparation: Data Cleaning. The main challenge I faced was the presence of HTML code within the text fields. This code caused problems in export procedures.

To address this, I implemented a data cleaning process through the creation of database views. These views serve a specific purpose: they remove unnecessary carriage return (\r) and newline (\n) characters from text fields. Additionally, to align with the objective of converting the data into CSV format, I needed to ensure that each text field is properly quoted.

The SQL script you see performs this transformation seamlessly. By wrapping the text fields with quotes and stripping out unwanted characters, I make the data CSV-ready. This step was not only necessary but also the most efficient route to prepare the text fields for the conversion process.

60''

5

In this slide, I'm detailing the efficient process of exporting data from SQL Server to CSV and then compressing it for optimal transfer speed.

I start by using BCP, which stands for Bulk Copy Program, a tool provided by SQL Server. It's the fastest method available to perform large-scale data exports to CSV format, ensuring speed and efficiency. I specify the view names that form the collection of data I wish to export, and then I initiate a loop to process each view individually.

Once the CSV files are generated, I use 7-Zip to compress them into gzip format. One of the key advantages of gzip is its compatibility with Databricks' file system. It can read gzip files directly, which means there's no need to decompress them on the platform, saving valuable processing time and simplifying the data pipeline.

The bottom part of the slide shows the actual script that orchestrates this entire process: from setting environment variables for BCP to executing the compression commands and cleaning up interim files. This script is the backbone of the procedure, handling everything from data extraction to final file compression, ready for sending to Databricks.

67''

Creating a Docker container with Databricks CLI, Copying Files to Databricks.

1. Defining a dockerfile

```
docker > databricks-cli-docker > docker-build > ❖ dockerfile > ...
  1   FROM python:latest
  2
  3   RUN pip install databricks-cli
  4
  5   RUN mkdir /data
  6
  7   WORKDIR /data
```

2. Creating an image

```
docker > databricks-cli-docker > ⊞ dockerbuild.cmd
  1   docker build -t databricks-cli:latest .\docker-build -f .\docker-build\dockerfile
```

3. Running the docker container

```
docker > databricks-cli-docker > ⊞ dockerrun.cmd
  1   docker run --rm -it --name databricks-cli  --volume D:\csv:/data/ databricks-cli:latest bash
```

4. Copying to Databricks

```
### the first step - configure token
databricks configure --token
here:
https://adb-6199578147525915.15.azuredatabricks.net/
**************************************************

### copy
databricks fs cp -r /data dbfs:/fall_2023_users/piotreks/csv/ --overwrite
```

5. The result

```
root@ea80dcbd59ec:/data# databricks fs ls -l dbfs:/fall_2023_users/piotreks/csv/
file      62431783  vBadges.csv.gz      1703342580000
file    1852364742  vComments.csv.gz    1703343750000
file          74    vLinkTypes.csv.gz   1703342745000
file    19150731    vPostLinks.csv.gz   1703342788000
file         147    vPostTypes.csv.gz   1703342717000
file    5647629288  vPosts.csv.gz       1703342396000
file    105608775   vUsers.csv.gz       1703342858000
file         211    vVoteTypes.csv.gz   1703342685000
file    340199412   vVotes.csv.gz       1703343114000
```

In this section of the presentation, I'm going to walk you through the process of creating a Docker container equipped with Databricks CLI for the purpose of transferring files efficiently to Databricks.

I began by defining a Dockerfile. This file is a blueprint for Docker, describing the environment needed to run Databricks CLI. It's based on a Python image, installs the Databricks CLI, and sets a working directory.

Next, I created an image from the Dockerfile. The command shown on the slide builds the Docker image, which contains all the necessary components to run the Databricks CLI.

Once the image was ready, I proceeded to run the Docker container. This step involves executing a Docker run command, which initiates a container instance where I can execute Databricks CLI commands.

The fourth step was to copy the CSV files to Databricks. To do this securely, I first configured access to Databricks using a token. This ensures that the connection to Databricks is secure and that only authorized commands are executed. Once the token was configured, I used the Databricks CLI to copy files from the local data directory to the Databricks file system, replacing any existing files with the same name.

The result, which you can see on the right, is the list of files successfully copied to Databricks. These

are the gzipped CSV files now stored in the Databricks file system, ready to be used for further conversion and data analysis.

84''

## Converting the Compressed CSV Files to Parquet files.

Preparing functions to be used during the conversion:

```scala
import org.apache.spark.sql.DataFrame
import org.apache.spark.sql.types.{StructType, StructField, StringType, IntegerType, DateType, TimestampType, LongType}

// Function for reading a csv file and applying a schema
def readCSV(filePath: String, schema: StructType): DataFrame = {
  spark.read
    .option("delimiter", ",")
    .option("header", "true")
    .option("quote", "\"")
    .schema(schema)
    .csv(filePath)
}

//Function for getting the compressed file path
def getCSVPath(entity: String): String = {
  s"/fall_2023_users/piotreks/csv/v${entity}.csv.gz"
}

//Function for getting the PARQUET files path
def getParquetPath(entity: String): String = {
  s"/fall_2023_users/piotreks/parquet/${entity}.parquet"
}
```

This slide illustrates the functions I've prepared for converting compressed CSV files to the Parquet format, optimized for processing with Apache Spark.

The first function, readCSV, is a Scala function that reads a CSV file into a Spark DataFrame, enforcing a predefined schema. It's essential for ensuring the data conforms to the expected format, which facilitates reliable analysis and processing.

To manage file paths easier, I use the getCSVPath and getParquetPath functions. They generate the storage paths for the CSV and Parquet files, respectively.

By applying these functions, I ensure that the data transformation pipeline is not only automated but also integrated, setting the stage for the subsequent analytics tasks in Spark without any manual intervention.

79''

8

## Converting the Compressed CSV Files to Parquet files.

Converting to Parquet format:

```scala
1    //Convert to parquet, use "LEGACY" because of some problems with date conversion
2    spark.conf.set("spark.sql.legacy.timeParserPolicy", "LEGACY")
3
4    val entity = "Posts"
5
6    var schema = StructType(
7      Array(
8        StructField("Id", IntegerType, nullable = false),
9        StructField("AcceptedAnswerId", IntegerType, nullable = true),
10       StructField("AnswerCount", IntegerType, nullable = true),
11       StructField("Body", StringType, nullable = false),
12       StructField("ClosedDate", TimestampType, nullable = true),
13       StructField("CommentCount", IntegerType, nullable = true),
14       StructField("CommunityOwnedDate", TimestampType, nullable = true),
15       StructField("CreationDate", TimestampType, nullable = false),
16       StructField("FavoriteCount", IntegerType, nullable = true),
17       StructField("LastActivityDate", TimestampType, nullable = false),
18       StructField("LastEditDate", TimestampType, nullable = true),
19       StructField("LastEditorDisplayName", StringType, nullable = true),
20       StructField("LastEditorUserId", IntegerType, nullable = true),
21       StructField("OwnerUserId", IntegerType, nullable = true),
22       StructField("ParentId", IntegerType, nullable = true),
23       StructField("PostTypeId", IntegerType, nullable = false),
24       StructField("Score", IntegerType, nullable = true),
25       StructField("Tags", StringType, nullable = true),
26       StructField("Title", StringType, nullable = true),
27       StructField("ViewCount", IntegerType, nullable = true)
28     )
29   )
30
31   var df = readCSV(getCSVPath(entity), schema)
32
33   df.write.mode("overwrite").parquet(getParquetPath(entity))
```

On this slide, I show the process I used to convert a compressed CSV file into the Parquet format using Spark. The schema defined here mirrors the exact structure in the data model shown earlier.

Here, you can see the code snippet for converting the 'Posts' file as an example. I've set Spark's SQL legacy time parser policy to avoid issues with date conversion, which is a common problem in data processing.

The readCSV function is used to read the CSV file into a DataFrame, applying the schema directly as it's defined. Then, with the write method, the DataFrame is saved in the Parquet format. While this slide shows the conversion for one file, I want to note that the same approach was used for the rest of the files. Each CSV file was converted in a similar manner, ensuring a uniform and efficient batch processing.

46''

9

## Data Analysis



Load All Parquet Files to A Corresponding Dataframe

```
1    // Load the DataFrames
2    val badgesDF = spark.read.parquet("dbfs:/fall_2023_users/piotreks/parquet/Badges.parquet")
3
4    val commentsDF = spark.read.parquet("dbfs:/fall_2023_users/piotreks/parquet/Comments.parquet")
5
6    val postsDF = spark.read.parquet("dbfs:/fall_2023_users/piotreks/parquet/Posts.parquet")
7
8    val linkTypesDF = spark.read.parquet("dbfs:/fall_2023_users/piotreks/parquet/LinkTypes.parquet")
9
10   val postTypesDF = spark.read.parquet("dbfs:/fall_2023_users/piotreks/parquet/PostTypes.parquet")
11
12   val postLinksDF = spark.read.parquet("dbfs:/fall_2023_users/piotreks/parquet/PostLinks.parquet")
13
14   val usersDF = spark.read.parquet("dbfs:/fall_2023_users/piotreks/parquet/Users.parquet")
15
16   val votesDF = spark.read.parquet("dbfs:/fall_2023_users/piotreks/parquet/Votes.parquet")
17
18   val voteTypesDF = spark.read.parquet("dbfs:/fall_2023_users/piotreks/parquet/VoteTypes.parquet")
```

▶ (18) Spark Jobs

▶ 🖽 badgesDF:  org.apache.spark.sql.DataFrame = [Id: integer, Name: string ... 2 more fields]
▶ 🖽 commentsDF:  org.apache.spark.sql.DataFrame = [Id: integer, CreationDate: timestamp ... 4 more fields]
▶ 🖽 postsDF:  org.apache.spark.sql.DataFrame = [Id: integer, AcceptedAnswerId: integer ... 18 more fields]
▶ 🖽 linkTypesDF:  org.apache.spark.sql.DataFrame = [Id: integer, Type: string]
▶ 🖽 postTypesDF:  org.apache.spark.sql.DataFrame = [Id: integer, Type: string]
▶ 🖽 postLinksDF:  org.apache.spark.sql.DataFrame = [Id: integer, CreationDate: timestamp ... 3 more fields]
▶ 🖽 usersDF:  org.apache.spark.sql.DataFrame = [Id: integer, AboutMe: string ... 12 more fields]
▶ 🖽 votesDF:  org.apache.spark.sql.DataFrame = [Id: integer, PostId: integer ... 4 more fields]

This slide shows the first step of my data analysis, where I have loaded all Parquet files into Spark DataFrames. Each DataFrame corresponds to a specific component of the Stack Overflow data captured in Parquet files, which are highly efficient for this kind of processing.

By loading these datasets into DataFrames, I've laid the foundation for further data exploration.

18''

Data Analysis

This slide presents an analysis aimed at identifying the top 5 active users on Stack Overflow by post count, along with their most engaged tags. I've used Spark to filter and process posts, and then to count and rank user activity. By exploding the tags column, I ensured each tag associated with a post was considered individually, which allowed me to determine the most common tags per user.

After aggregating the total post counts and the most active tags for each user, I joined this data with the users' details to provide a more comprehensive view. The result of this join operation gave me a table showcasing users' display names, their total post counts, and their primary tags.

Finally, I visualized this data using Python's matplotlib and pandas libraries to create the pie chart you see here. This chart highlights the proportion of contributions from each of the top users.

53''

This slide demonstrates an approach to counting annual posts in a forum dataset, focusing on content related to the Scala programming language. Initially, I used Spark's DataFrame API to filter the dataset for posts tagged with 'Scala'. Next, I extracted the year from each post's creation date and aggregated the count of posts per year. This process involved sorting and grouping the data to facilitate a clearer analysis.

To visualize the trends, I employed Python's matplotlib and pandas libraries again, creating a bar chart.

25''

## Data Analysis

3. Who is the Scala Guru? - Identifying the User with the Max Score in Scala Posts and Comments

```scala
1   import org.apache.spark.sql.functions._
2   import spark.implicits._
3
4   // Filter Scala posts
5   val scalaPostsDF = postsDF.filter($"Tags".contains("<scala>"))
6
7   // Prepare comments related to Scala posts
8   val scalaCommentsDF = commentsDF
9     .join(scalaPostsDF, commentsDF("PostId") === scalaPostsDF("Id"))
10    .select(commentsDF("UserId").alias("CommentUserId"), commentsDF("Score").alias("CommentScore"))
11
12  // Aggregate scores by user from posts
13  val postScores = scalaPostsDF
14    .groupBy("OwnerUserId")
15    .agg(sum("Score").alias("PostScore"))
16
17  // Aggregate scores by user from comments
18  val commentScores = scalaCommentsDF
19    .groupBy("CommentUserId")
20    .agg(sum("CommentScore").alias("CommentScore"))
21
22  // Combine scores from posts and comments
23  val combinedScores = postScores
24    .join(commentScores, postScores("OwnerUserId") === commentScores("CommentUserId"), "outer")
25    .na.fill(0) // Replace null values with 0
26    .withColumn("TotalScore", $"PostScore" + $"CommentScore")
27    .select($"OwnerUserId", $"TotalScore")
28
29  // Find the top user
30  val scalaGuruDF = combinedScores
31    .join(usersDF, $"OwnerUserId" === usersDF("Id"))
32    .orderBy(desc("TotalScore"))
33    .limit(1)
34    .select("DisplayName", "Location", "TotalScore")
35
36  scalaGuruDF.collect().foreach { row =>
37    val name = row.getAs[String]("DisplayName")
38    val location = row.getAs[String]("Location")
39    val score = row.getAs[Long]("TotalScore")
40    println(s"The Scala guru was $name from $location with a total score of $score.")
41  }
```

The Scala guru was oxbow_lakes from London, England United Kingdom with a total score of 3,223.

This slide focuses on finding the top participant regarding Scala topics. I began by filtering the posts for those tagged with 'Scala' and then prepared a dataset of comments related to these posts. By joining the two datasets and selecting relevant fields, I aggregated scores from both posts and comments for each user.

Next, I combined these scores to get a comprehensive view of user engagement. The final step was to identify the top contributor, which I achieved by ordering the combined scores and limiting the results to the single highest score.

The outcome of this analysis highlighted a user, showcased on the right of the slide, who emerged as the Scala guru with an impressive total score, reflecting his expertise and active participation in the Scala community.

42''

This slide presents a targeted analysis to pinpoint potential Polish users within a dataset. I began by establishing regex patterns to match Polish names and locations. Then, I applied filters on the user data to identify names or contain typical Polish characters, and also to find users from locations in Poland.

Subsequently, I filtered posts and comments for matches with the same patterns. By combining these filtered datasets, I could create a comprehensive list of users likely to be Polish. The final step was counting these users, which revealed a total of 12,454 potential Polish users.

38''

# Conclusions

1. Data cleaning required the most substantial time investment.

2. Spark worked really well with Parquet files when the schema was set up in advance.

3. Regular expression operations significantly slowed down the process

4. A more efficient strategy may involve utilizing Azure SQL Database with a direct Databricks integration.

This slide summarizes the key findings from the analysis. The most significant observation is that the bulk of the time was spent on cleaning the data, which underscores the importance of having clean data for accurate analysis. I noted that Spark was highly efficient in handling Parquet files, though it's crucial to provide the correct schema upfront to ensure smooth processing.

Another point of interest is the observation that using regular expressions proved to be quite slow, which suggests the need for more efficient text processing methods in future analyses. Lastly, considering the overall data pipeline performance, deploying a database on Azure SQL Database and connecting it directly to Databricks could offer a more streamlined and faster approach.

42''