

NoSQL Data Modeling

<https://www.github.com/bryangoodrich/python-exercises>

Data Collection

This is far from a real-world example, but it's sufficient to demonstrate how you can rethink your SQL modeling strategies when working on NoSQL systems, like Spark or BigQuery.

Using the UCI Wine Quality dataset, my aim to use the categorical "quality" field to group collections of values of another field, "density", into a vector.

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, size, collect_list, collect_set
import pandas as pd

spark = SparkSession.builder.getOrCreate()
uci = "https://archive.ics.uci.edu/ml/machine-learning-databases"
url = f"{uci}/wine-quality/winequality-red.csv"
src = pd.read_csv(url, sep=";")
src.columns = src.columns.str.replace(" ", "_") # snake_case
src.info()
# <class 'pandas.core.frame.DataFrame'>
# RangeIndex: 1599 entries, 0 to 1598
# Data columns (total 12 columns):
# #   Column                Non-Null Count  Dtype
# ---  ---
# 0   fixed_acidity          1599 non-null   float64
# 1   volatile_acidity       1599 non-null   float64
# 2   citric_acid            1599 non-null   float64
# 3   residual_sugar         1599 non-null   float64
# 4   chlorides              1599 non-null   float64
# 5   free_sulfur_dioxide    1599 non-null   float64
# 6   total_sulfur_dioxide   1599 non-null   float64
# 7   density                1599 non-null   float64
# 8   pH                     1599 non-null   float64
# 9   sulphates              1599 non-null   float64
# 10  alcohol                1599 non-null   float64
# 11  quality                1599 non-null   int64
# dtypes: float64(11), int64(1)
```

Data Modeling

Using the `collect_list` and `collect_set` aggregate functions, we create both a list vector and a distinct set of density values (sets remove duplicates). The results show we can store the entire table of density values as collections, not rows of atomic values, including jagged arrays—i.e., not equal number of “rows” per group.

```
density = collect_list("density").alias("density")
dense_set = collect_set("density").alias("set")

df = (spark
      .createDataFrame(src)
      .groupBy("quality")
      .agg(density, dense_set)
)

cnt = size("density").alias("counts")
dist = size("set").alias("unique")
projection = ["quality", "density", "set", cnt, dist]

df.select(projection).sort("quality").show()
# +-----+-----+-----+-----+
# |quality|      density|      set|counts|unique|
# +-----+-----+-----+-----+
# |      3|[1.0008, 0.9994, ...|[0.99476, 0.99705...|    10|    10|
# |      4|[0.9974, 0.994, 0...|[0.99672, 0.99517...|    53|    48|
# |      5|[0.9978, 0.9968, ...|[0.99784, 0.99334...|   681|   239|
# |      6|[0.998, 0.9969, 0...|[0.9921, 0.9997, ...|   638|   266|
# |      7|[0.9946, 0.9968, ...|[0.9997, 0.99834,...|   199|   122|
# |      8|[0.9973, 0.9976, ...|[0.9917, 0.99472,...|    18|    17|
# +-----+-----+-----+-----+

spark.stop()
```

 bryangoodrich.xyz

💡 If you like this post, follow me for daily #python tips, and hit that like button so the algorithms help others see it, too. For full code and data on this and earlier exercises, visit <https://www.github.com/bryangoodrich/python-exercises>