

# SQL В АНАЛИЗЕ ДАННЫХ 101

РЕШАЕМ ПРИКЛАДНУЮ ЗАДАЧУ СЕГМЕНТАЦИИ

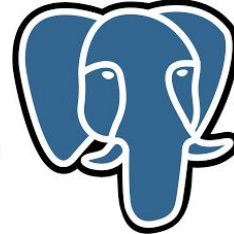
# ПЛАН

- Немного теории про Spark и RFM
- Разбираем код для решения на Spark SQL в Databricks
- Строим дашбордик и публикуем в [PUBLIC.TABLEAU.COM](https://public.tableau.com)
- Обсуждаем вопросы
- Happy End



# ПОЧЕМУ SQL?

- ТРАДИЦИОННЫЕ БД (MS SQL, POSTGRES, MYSQL, ORACLE)
- NOSQL (MONGODB, CASSANDRA, NEO4J, HADOOP)
- MPP (VERTICA, EXASOL, BQ, REDSHIFT, GREENPLUM)
- SQL ENGINES (PRESTO, SPARK)
- ETL TOOLS (DBT, LUIGI, AIRFLOW)
- BI TOOLS
- OTHER (CLICKHOUSE, KAFKA)



Amazon Redshift



Google BigQuery



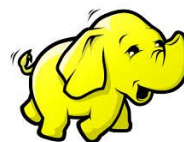
EXASOL



cassandra



mongoDB



# SPARK

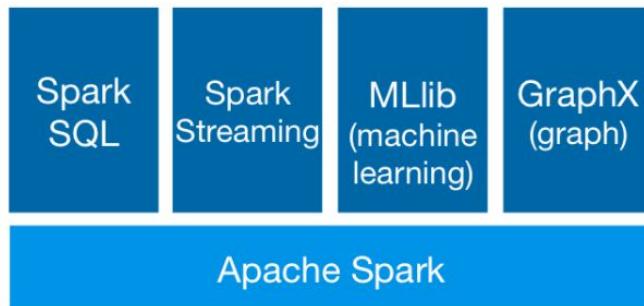
- SCALA, JAVA, PYTHON, R AND SQL (KOTLIN, CLOJURE)

- [HTTPS://SPARK.APACHE.ORG/](https://spark.apache.org/)

- DATAFRAME API

- SQL (HIVE SQL OR PG SQL)

- [HTTPS://WWW.WAITINGFORCODE.COM/APACHE-SPARK-SQL/WHAT-NEW-APACHE-SPARK-3-POSTGRESQL-FEATURE-PARITY/READ](https://www.waitingforcode.com/apache-spark-sql/what-new-apache-spark-3-postgresql-feature-parity/read)



## Runs Everywhere

Spark runs on Hadoop, Apache Mesos, Kubernetes, standalone, or in the cloud. It can access diverse data sources.

You can run Spark using its [standalone cluster mode](#), on [EC2](#), on [Hadoop YARN](#), on [Mesos](#), or on [Kubernetes](#). Access data in [HDFS](#), [Alluxio](#), [Apache Cassandra](#), [Apache HBase](#), [Apache Hive](#), and hundreds of other data sources.

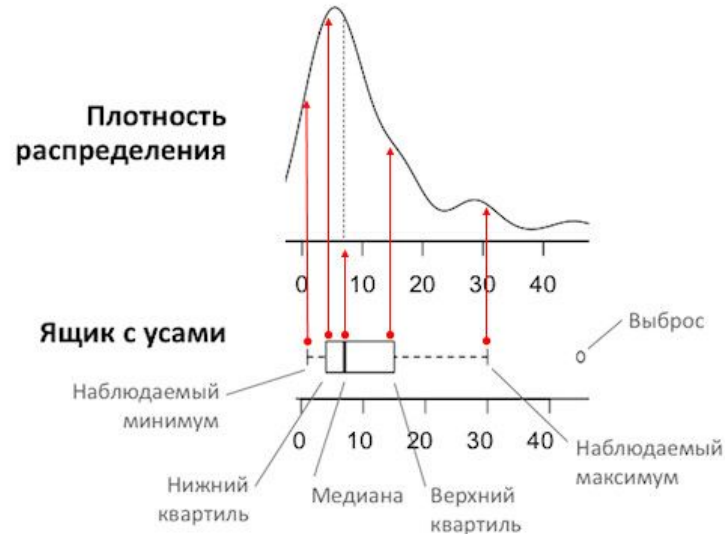
# RFM

|   | InvoiceNo | StockCode | Description                         | Quantity | InvoiceDate        | UnitPrice | CustomerID | Country        |
|---|-----------|-----------|-------------------------------------|----------|--------------------|-----------|------------|----------------|
| 1 | 536365    | 85123A    | WHITE HANGING HEART T-LIGHT HOLDER  | 6        | 01.12.2010 8:26:00 | 2.55      | 17850      | United Kingdom |
| 2 | 536365    | 71053     | WHITE METAL LANTERN                 | 6        | 01.12.2010 8:26:00 | 3.39      | 17850      | United Kingdom |
| 3 | 536365    | 84406B    | CREAM CUPID HEARTS COAT HANGER      | 8        | 01.12.2010 8:26:00 | 2.75      | 17850      | United Kingdom |
| 4 | 536365    | 84029G    | KNITTED UNION FLAG HOT WATER BOTTLE | 6        | 01.12.2010 8:26:00 | 3.39      | 17850      | United Kingdom |
| 5 | 536365    | 84029E    | RED WOOLLY HOTTIE WHITE HEART.      | 6        | 01.12.2010 8:26:00 | 3.39      | 17850      | United Kingdom |
| 6 | 536365    | 22752     | SET 7 BABUSHKA NESTING BOXES        | 2        | 01.12.2010 8:26:00 | 7.65      | 17850      | United Kingdom |
| 7 | 536365    | 21730     | GLASS STAR FROSTED T-LIGHT HOLDER   | 6        | 01.12.2010 8:26:00 | 4.25      | 17850      | United Kingdom |

Showing the first 1000 rows

- RECENCY (R) — ДАВНОСТЬ ПОСЛЕДНЕЙ ПОКУПКИ
- FREQUENCY (F) — СУММАРНАЯ ЧАСТОТА ПОКУПОК
- MONETARY (M) — ОБЪЁМ ПОКУПОК

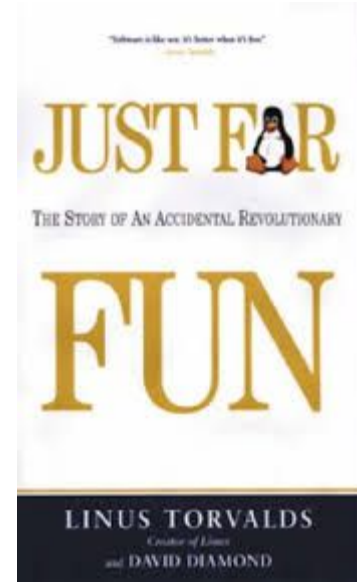
- <https://esputnik.com/blog/rfm-segmentaciya>
- <https://community.tableau.com/s/question/0D54T00000G54JKSAB/rfm-segmentation>
- <https://medium.com/analytics-vidhya/marketing-analytics-rfm-modeling-855ebec18014>
- <https://habr.com/ru/company/ptarum/blog/431520/>
- <https://towardsdatascience.com/find-your-best-customers-with-customer-segmentation-1>
- <https://towardsdatascience.com/dynamic-customer-analytics-19135b2a8854b>
- <https://guillaume-martin.github.io/rfm-segmentation-with-python.html>



# ЧТО НАДО СДЕЛАТЬ ЧТОБ НАЧАТЬ?

- [HTTPS://DATABRICKS.COM/TRY-DATABRICKS](https://databricks.com/try-databricks) OR [HTTPS://COMMUNITY.CLOUD.DATABRICKS.COM/LOGIN.HTML](https://community.cloud.databricks.com/login.html)
- [HTTPS://PUBLIC.TABLEAU.COM/S/](https://public.tableau.com/s/)
- [HTTPS://PUBLIC.TABLEAU.COM/EN-US/S/DOWNLOAD](https://public.tableau.com/en-us/s/download)
-

LET'S CODE



# ЧТО ЧИТАТЬ И СМОТРЕТЬ?

- [HTTPS://MODE.COM/SQL-TUTORIAL/](https://mode.com/sql-tutorial/)
- [HTTPS://MODE.COM/BLOG/FINDING-USER-SESSIONS-SQL/](https://mode.com/blog/finding-user-sessions-sql/)
- [HTTPS://EXACADEMY.EXASOL.COM/COURSES/COURSE-V1:EXASOL+ESSENTIALS+X/COURSE/](https://exacademy.exasol.com/courses/course-v1:EXASOL+ESSENTIALS+X/course/)
- [HTTPS://EXACADEMY.EXASOL.COM/COURSES/COURSE-V1:EXASOL+PERF+X/COURSE/](https://exacademy.exasol.com/courses/course-v1:EXASOL+PERF+X/course/)
- [HTTPS://WWW.KAGGLE.COM/YAMAERENAY/SPOTIFY-DATASET-19212020-160K-TRACKS](https://www.kaggle.com/yamaerenay/spotify-dataset-19212020-160k-tracks)

## Announcements

### Free Exasol Certifications

Starting November 23rd, we have decided to offer free vouchers for Exasol Certifications permanently! To be eligible, you must score 70% on the Hands-on Exercises.

[More Details](#)

Use this opportunity to check out our Course Offering and get up to date on your Exasol certifications !

Exasol

Exasol Certificate

This certifies that  
Eugene Kudashev  
is recognized by Exasol as an  
Exasol Certified Performance Expert

*Uwe Hesse*

Uwe Hesse, Training Coordinator

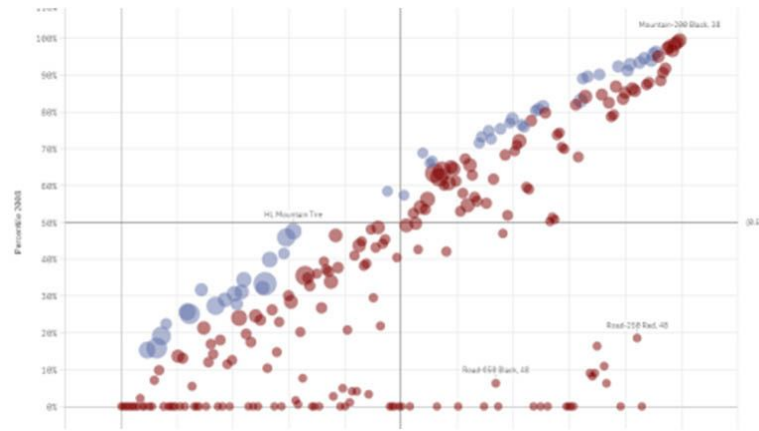
29th Jun 2020

Certificate ID: CLVPRVFRYF-SKVQZTLQ-XHMMWHXWRT



## ЧТО ЕЩЕ ПОДЕЛАТЬ?

- ДОБАВИТЬ ВОЗМОЖНОСТЬ РАСЧЕТА ПО СТРАНАМ (ROLLUP + TABLEAU FILTERS)
- ПОПРОБОВАТЬ POSTGRES SQL ДЛЯ SPARK
- ПЕРЕНЕСТИ ВСЮ ЛОГИКУ В TABLEAU + ДИНАМИЧЕСКИЕ ПОРОГИ ДЛЯ RFM
- [HTTP://BLOG.ATKCG.RU/PERCENTIL-ALTERNATIVNYJ-VZGLYAD-NA-DANNYE/](http://blog.atkcg.ru/percentil-alternativnyj-vzglyad-na-dannye/)
- 



СПАСИБО ЧТО ДОТЕРПЕЛИ ДО HAPPY END-A

ВОПРОСЫ?