# kmeans-spark

This project shows example of training K-means model using spark. ClickHouse was used as a data source. Predictions were stored there either.

## Dataset

OpenFoodFacts dataset consists of the descriptions of different food products. More info could be found <u>here</u> <u>(https://world.openfoodfacts.org/data)</u>

## Data preparation

Data was preprocessed with removing of unimportant features and null columns filling.

## Project structure

1. <u>Research notebook (notebooks/)</u>
2. <u>Preprocessor (src/preprocessing/)</u>
3. <u>Model trainer (src/model.py)</u>
4. <u>Training with ClickHouse using docker-compose (docker-compose.yml)</u>

Consider put *clickhouse-jdbc-0.4.6-all.jar* in <u>jars (jars)</u> folder (used for clickhouse connection).