

DE Project

Data exploration

- Explore million song dataset
- Identify Genre information

Preprocessing and data reduction

- data cleaning - remove missing value
remove duplicate
- Feature extraction - Find most relevant
- data reduction feature

Data processing solution

- Distributed data Storage (HDFS or cloud)
- Processing framework based Storage)
Apache SPARK or ? Deploy own cluster
- Parallel processing

ML for Genre classification

- Model selection (e.g. Lin, SVM or ?)
- training validation
- Parameter tuning

Computational experiment

Report

Dataset Description

Analysis group
metadata