# Assignment Report - A7

*Shabbir, Sharad*

*October 31, 2017*

# NOTE: Best viewd in HTML version!

The goal of this assignment is to cluster artists based on their genre.

# Subtask 2

## Program Design

The program is split into 3 sub-packages:

- `org.neu.pdpmr.tasks.types` : It holds datatypes responsible for reading files and parsing them to provide fields. For extensibility purposes this file will read all fields and store them internally. For use within map reduce it is recommended to extract only portion of fields from it.

- `org.neu.pdpmr.tasks.Main` : Main file that executes whole project.

- `org.neu.pdpmr.tasks.subtask2` : Holds clustering mechanism for the second sub-task.
- `KMode` : It is based on KMode clustering technique.
    1. Initially we assign $1$ as weight to all terms associated with any artist.
    2. We then take intersect count of *artist terms* to *centroid terms* as a distance measure.
    3. To calculate new centroid position we take mode (most repeating values) of, **set of artist terms**, for all artists belonging to that cluster.
    4. Finally, we calculate distance to this new centroid by using same intersection technique.
- `KMeans` : It is an experimental implementation which assumes every `term` as a dimension and does Cartesian distance between them to do clustering.
    1. Initially we assign $1$ as weight to all terms associated with any artist.
    2. Then we take a mean across the cluster for each dimension. This could be fractional depending on the actual values. We call this as score of centroid terms.
    3. Finally, we calculate distance to this new centroid by summing scores of matching centroid terms between `artist terms` and the `new centroid` (distance = taking inverse of sum).

## Assumptions and Specifications

- **Our centroids are not contrained to be a valid point in the graph.**

It could be a subset of any combination of *terms*. Thus, we cannot use the artist similarity to construct edges between artist and centroids. In our case we assume all artists are connected to each other iff they have at least $1$ common term. And distance between them is the inverse of the intersection of matching *terms*.

- We assume artist `A` is similar to artist `B` reverse is **NOT** true. As in case of $A \subset B$.

## Results

Below tag cloud shows cluster assignment of artists. This report is run on a $10000$ songs subset with approx $3800$ distinct artists.

Key information to read below figure:

1. Every iteration result is shown as bunch of tag clouds.
2. Every tag cloud is a cluster shown in row major format. And are sorted by `CLUSTER_ID`.
3. For iteration $i$ we highlight *terms* using iteration $i - 1$ centroids. So for de-referencing use previous iteration centroids.
4. Darker color represent centroid terms for that cluster.
5. The size of the word indicates the number of artists assigned to that cluster.
6. Note: Some centroids are not plotted due to shortage of screen space as decided by wordcloud library!

# Cluster assignment visualization

## Iter 1

**Iter 2**

# Iter 3

**Iter 6**

hip hop
r&b
zouk
sexy
world

funny
humorous
united states

rap
rock
funk
hardcore
reggae
gangsta
indie
soul
90s
pop

finish

electronic
world
folk
rock
80s
70s
united states
american acoustic
pop

hard rock
rock
guitar jazz
r&b
folk rock
classic rock

hardcore rap
rap
classic
funk
hardcore
gangster rap

jazz
latin
salsa
mambo cuba
pop
world
trance
remix
mashup

soul
dub roots
folk
ska
world
r&b

folk
ranchera
traditional
norteno
mexican folk

latin
rock
ballad folk
singer
90s latino

punk metal
rock
pop
alternative
indie jazz
acoustic
folk

dance
world
jazz
club
remix
indie
pop
folk
lo-fi
midlands

united states
hip hop
raga
dub
jamaica
electronic
reggae pop

jazz
latin
guitar
world fusion
electronic
pop

finish
dark

ambient
world fusion
relaxation
dreamy
new age

chanson
swing
jazz
world
rock
soul
french
60s
ballad
jazz
rock
classic
american
instrumental

requiem
world
christian
southern gospel
opera
orchestra
gospel

house
ranchera
dj
los angeles
german
emotional
latin pop
club dance

hip hop
electro
breakbeat
drum and bass

world
beats
bass

chanson francaise
french folk
los angeles
free

greek
french pop
melbourne

**Iter 10**

# Final centroids after all iterations

# Conclusion

KModes is an interesting clustering results. We observed that centroids are mostly stabilized by 6th iteration on small corpus (with some exceptions). As algorithm progresses the clusters seems to merge together to form a bigger distinction between *terms* as seen in the case of `POP` which was prominent in every cluster earlier in iteration 1. It is later aggregated to its own *term* towards the end of the algorithm.

# Task Assignment

Shabbir was responsible for sub-problem 2 including sub-problem 2 reporting and respective Makefile.

Sharad was responsible for sub-problem 1 including sub-problem 1 reporting and respective Makefile.