Big Data HW1 README:

1. Submitted by: dxd132630 Name: deepti deshpande
2. The submission file contains :
    1. dxd132630HW1.jar file
    2. README.txt

Contents of JAR file and the execution details:
1. The JAR file contains all the necessary BUILD libraries
2. Three class files for each part of the question respectively:
    1. AgeBasedList.java : This Class file is the implementation of the Mapreduce function for HW1 part 1.
       The class file has only the implementation of the Map function with NULLwritable value.

       Command to execute the class file :
        hadoop jar <path_to_Jar_dir>/dxd132630HW1.jar AgeBasedList input output
        Where input and output are dfs locations for input and output directories.

       Expected output : The list of Male UserIds Whose age is less then or equal to 7

    2. AgeGrouping.java :This class file is the implementation of the Mapreduce function for HW1 part 2
        The class file has both map and reduce method implementation. The mapper class emits the Age and
gender as key and 1 as value
        The reduce method counts the number of male or female based on the age group
        Assumptions :
        51 is emited when age is >=45 and <=54
        56 is emited when age is >=55 and <=61

        Command to execute the class file:

        hadoop jar <path_to_Jar_dir>/dxd132630HW1.jar AgeGrouping input ouput
        Where input and output are dfs locations for input and output directories.

        Expected output : Count of Male and female users based on the age group:
        56 F 350
        56 M 146
        62 F 278
        62 M 102
        7 F 144
        7 M 78

    3. GenreBasedList.java : This class file is the implementation of the Mapreduce function for HW1 part
3.The class file has only Map menthod implementation. It reads the input files and lists the movies that
match the Genre mentioned by the user in the command line. ***The genre passed in the command line is
case sensitive****

        Command to execute the class file:

        hadoop jar <path_to_Jar_dir>/dxd132630HW1.jar GenreBasedList input ouput <genre>
        example : hadoop jar dxd132630HW1.jar GenreBasedList input ouput Drama
        Where input and output are dfs locations for input and output directories.

        Expected output : The list of the movies which belong to the mentioned genre in the command line