# Assignment 7

## TASK 1

Program to implement wordcount using Pig is
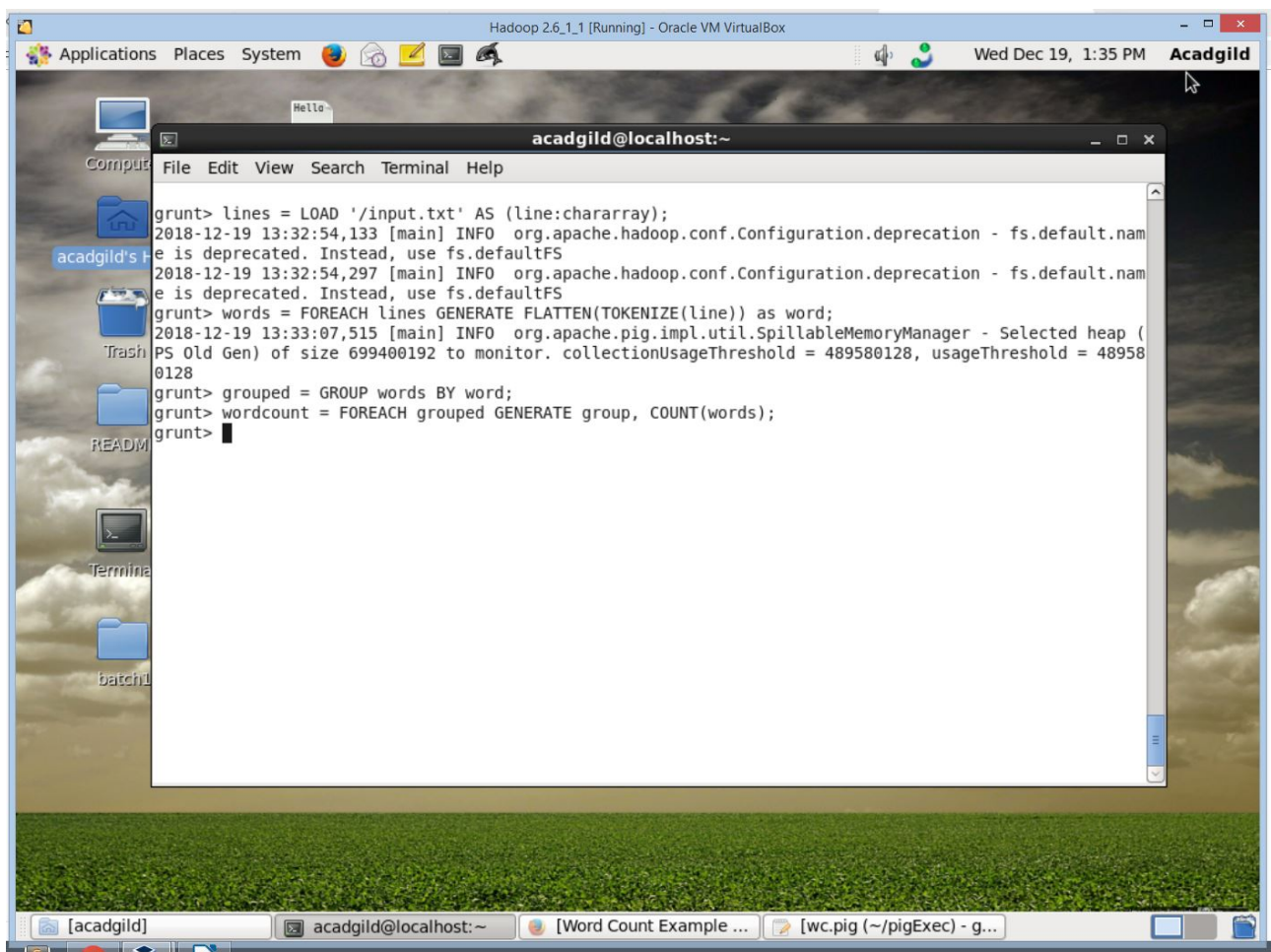
lines = LOAD '/input.txt' AS (line:chararray);

words = FOREACH lines GENERATE FLATTEN(TOKENIZE(line)) as word;

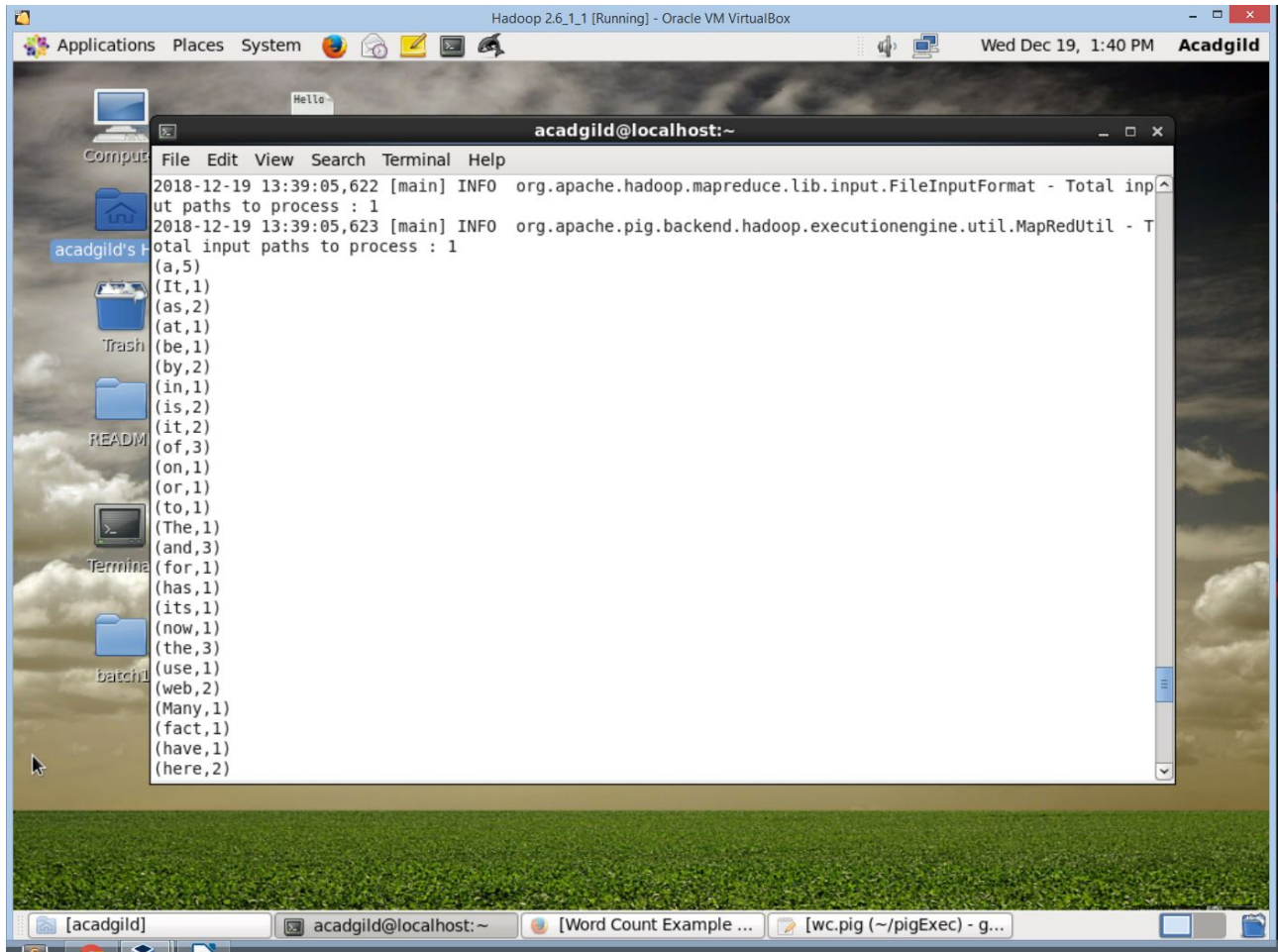grouped = GROUP words BY word;

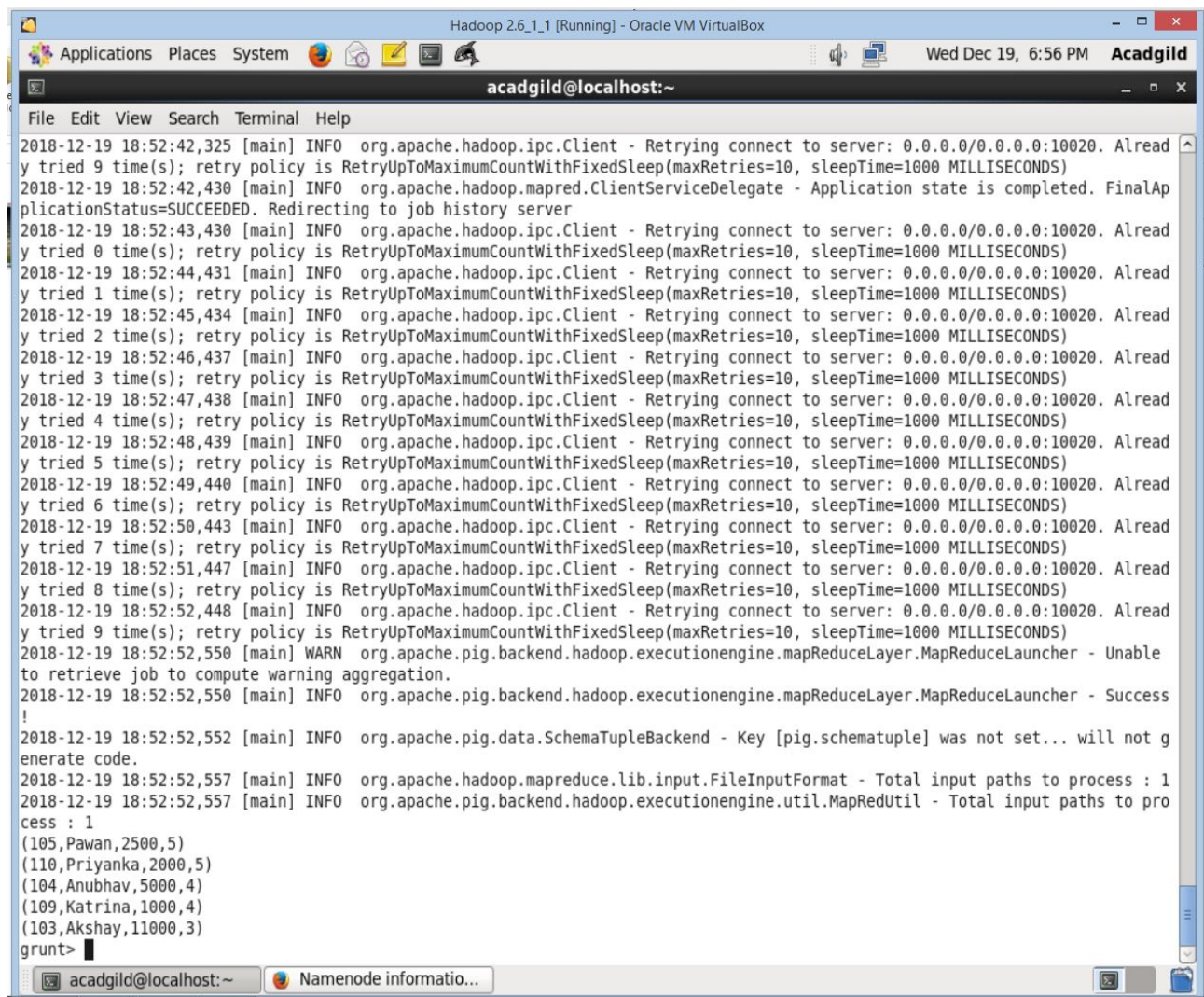wordcount = FOREACH grouped GENERATE group, COUNT(words);

DUMP wordcount;

**TASK 2:**

**A)** The Pig Script is

LOAD a = 'employee_details' USING PigStorage(',') AS (e_id:int, e_name:chararray, e_salary:int, e_rating:int);

b = order a by e_rating DESC, e_name ASC;

c = LIMIT b 5;

DUMP c;

**B)** The Pig Script is

a = LOAD 'employee_details.txt' USING PigStorage(',') AS (e_id:int, e_name:chararray, e_salary:int, e_rating:int);

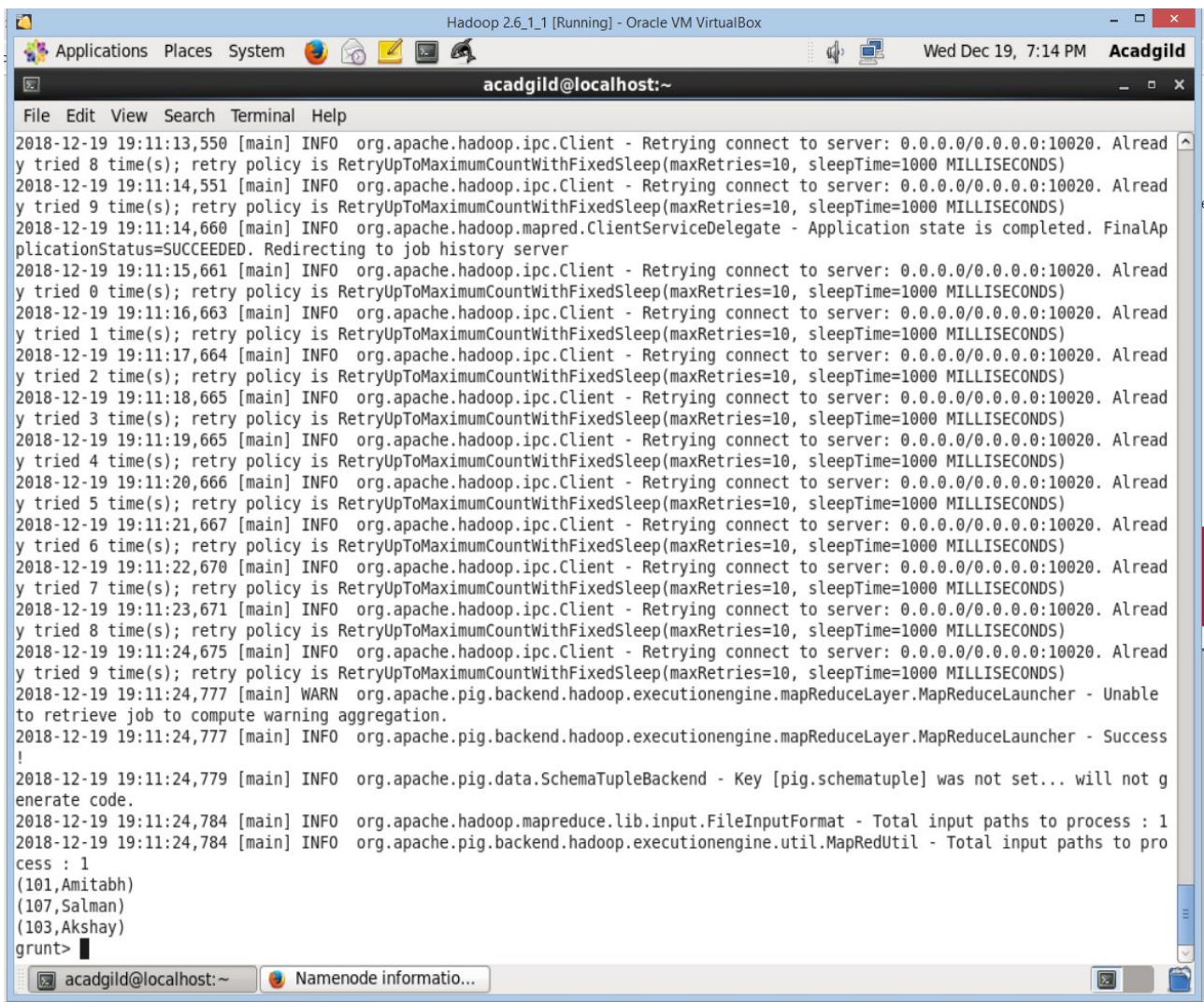b = order a by e_salary desc;

c = FILTER b by e_id%2==1;

d = FOREACH c generate  e_id,e_name;

e = LIMIT d 3;

DUMP e;

**C)** The Pig Script is

a = LOAD 'employee_details.txt' USING PigStorage(',') AS (e_id:int, e_name:chararray, e_salary:int, e_rating:int);

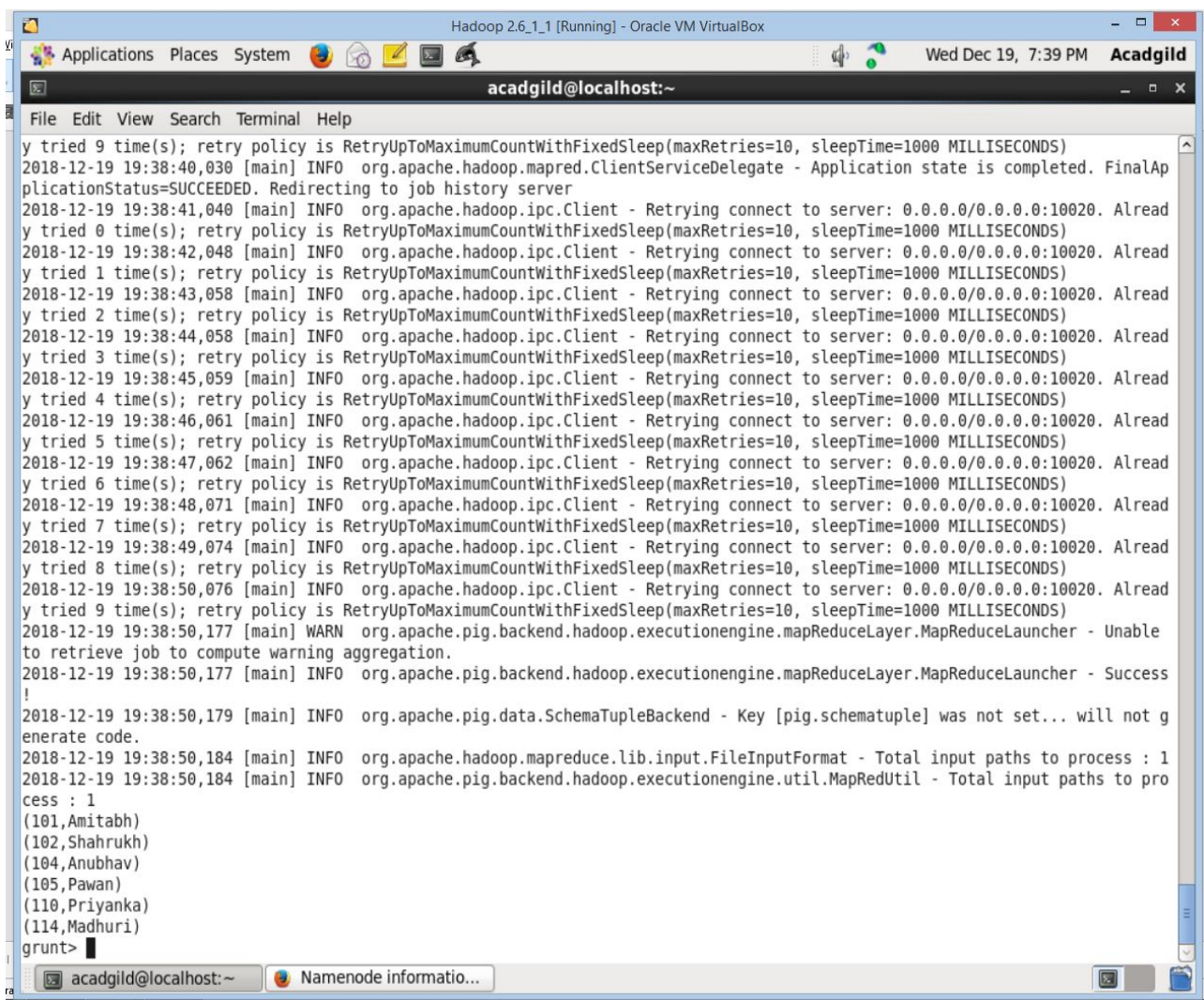b = LOAD '/employee_expenses.txt' AS (e_id:int, e_exp:int);

c = JOIN a BY e_id, b BY e_id;

d = FOREACH c GENERATE a::e_id, a::e_name;

e = DISTINCT d;

dump e;

**D)** The Pig Script is

a = LOAD 'employee_details.txt' USING PigStorage(',') AS (e_id:int, e_name:chararray, e_salary:int, e_rating:int);
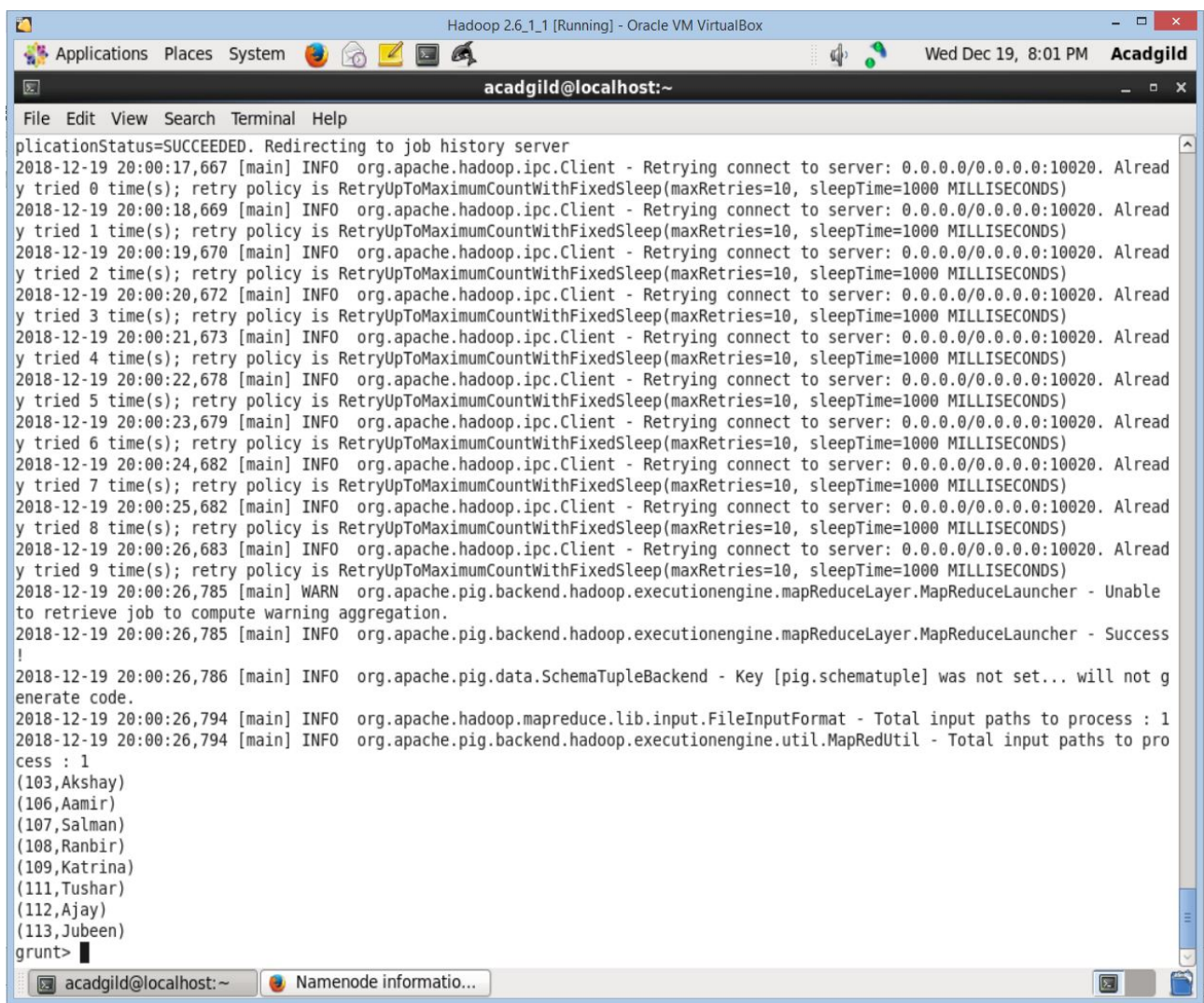
b = LOAD 'employee_expenses.txt' AS (e_id:int, e_exp:int);

c = JOIN a BY e_id LEFT OUTER, b BY e_id;

d = FILTER c BY b::e_id is null;
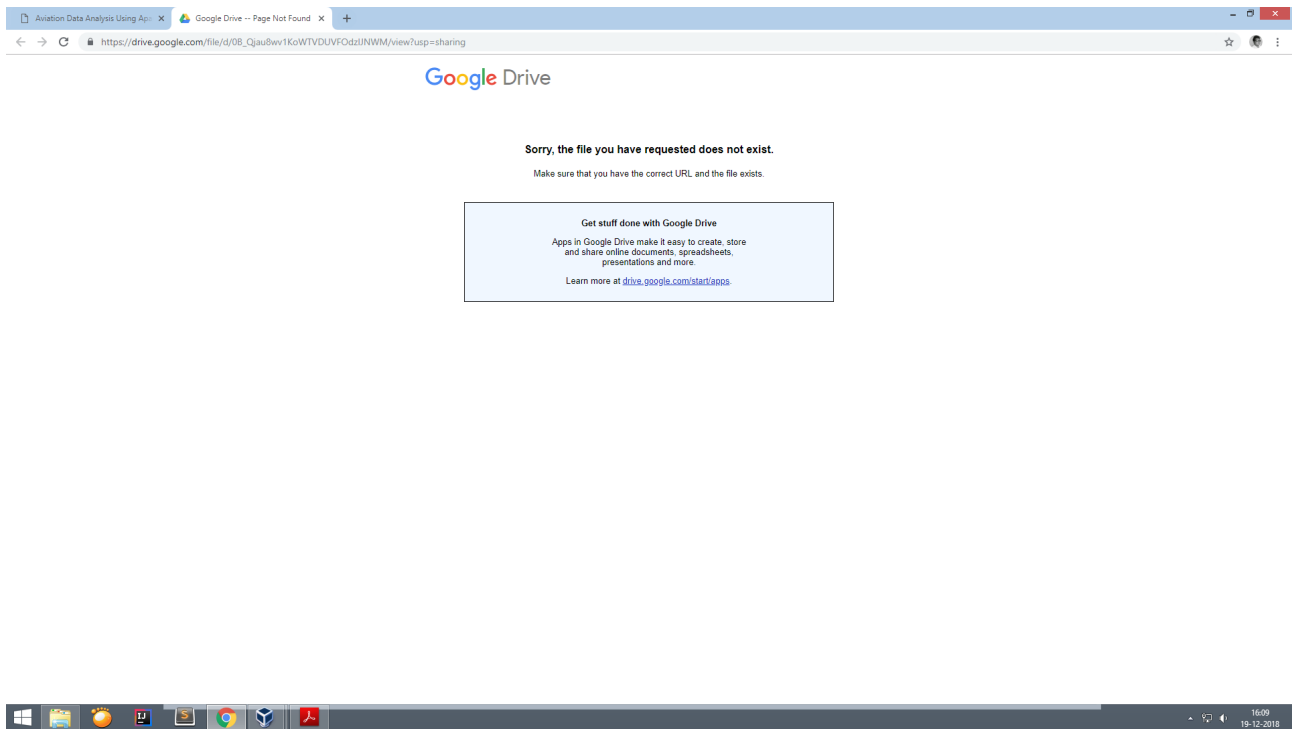
e = FOREACH d GENERATE a::e_id, a::e_name;

DUMP e;

## TASK 3:



The data sets for task 3 is not available in the provided link.