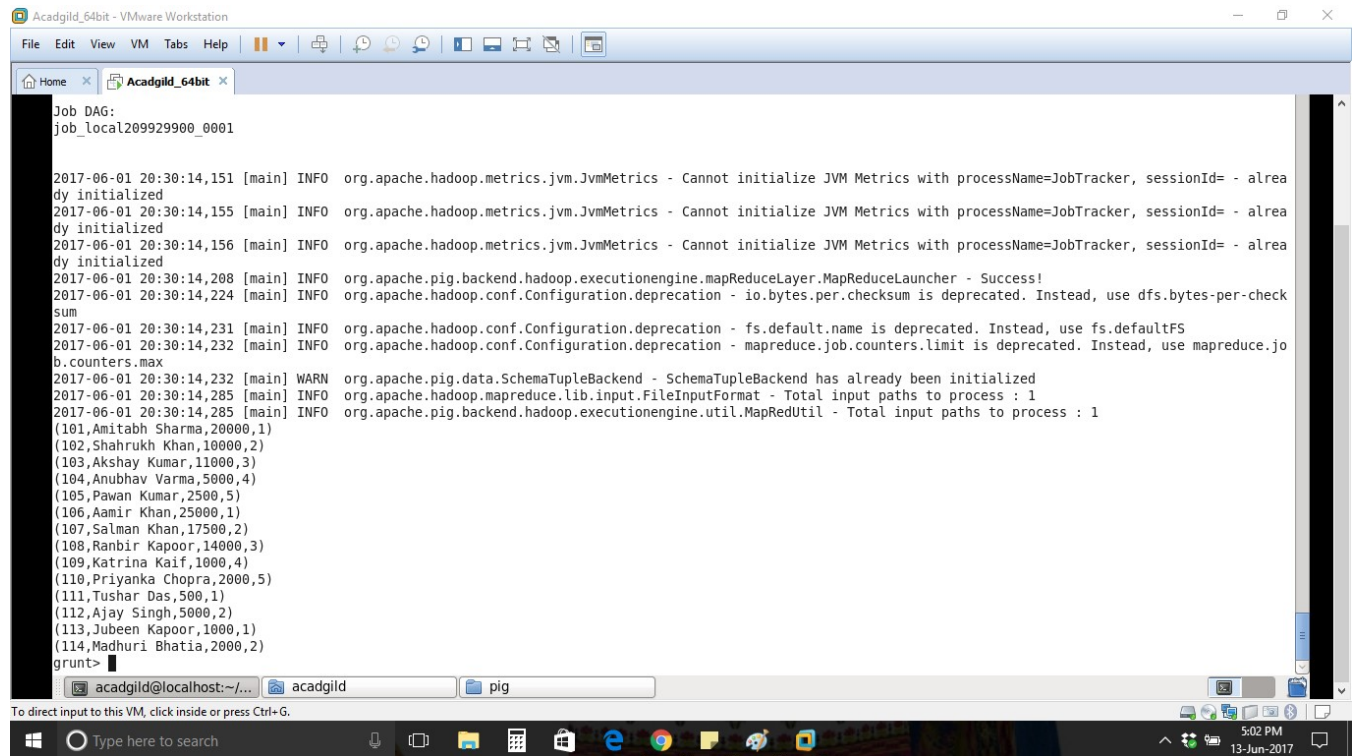# Assignment 8.2

**Employee_details file getting loaded**

emp = LOAD 'employee_details' USING PigStorage(',') AS (emp_id:int, emp_name:chararray, emp_salary:int);

dump emp;

## 8.2.1
## TOKENIZE

emp_name_tokenize = foreach emp generate TOKENIZE(emp_name);

dump emp_name_tokenize;

**8.2.2**

**CONCAT**

emp_name_detail = foreach emp generate CONCAT('Name of Employee is ' , emp_name);

dump emp_name_detail;

**8.2.3**

**SUM**

Employee details grouped by rating

emp_group_rating = group emp by emp_rating;

dump emp_group_rating;

## Sum of employee salaries grouped by rating

emp_sum_sal = foreach emp_group_rating generate group, SUM(emp.emp_salary);

dump emp_sum_sal;

## 8.2.4

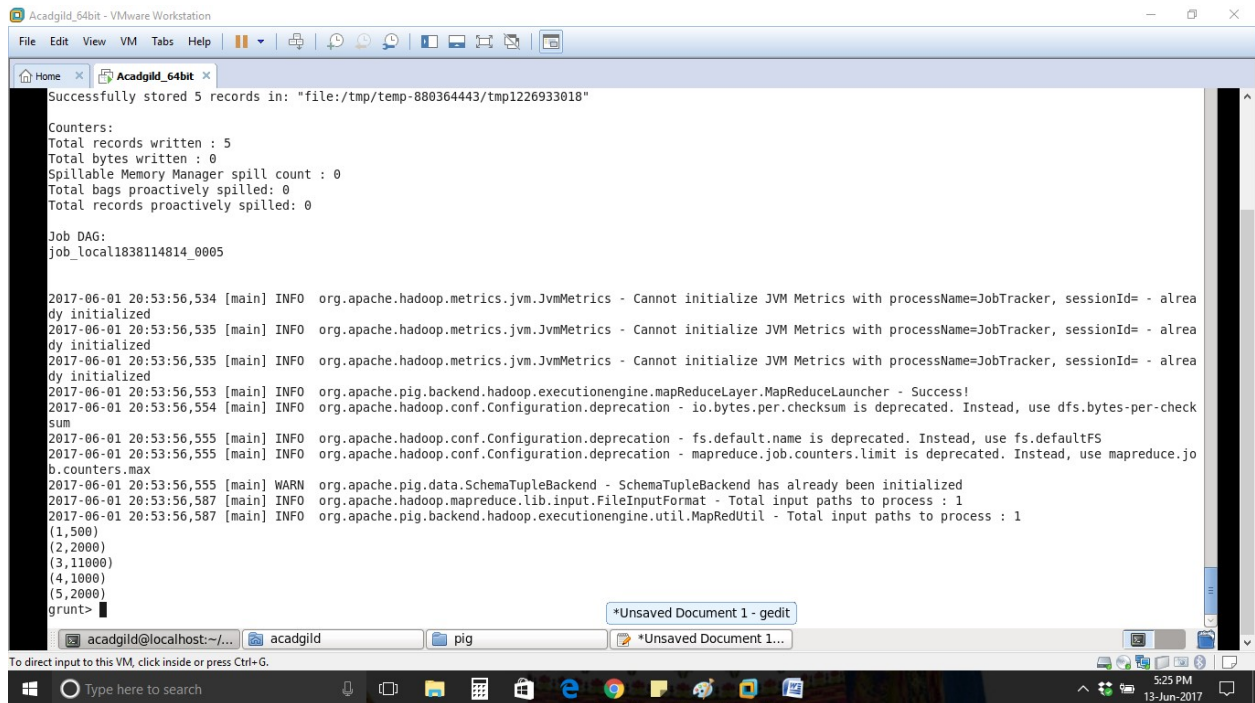## MIN

emp_min_sal = foreach emp_group_rating generate group, MIN(emp.emp_salary);

dump emp_min_sal;

## 8.2.5

## MAX

emp_max_sal = foreach emp_group_rating generate group, MAX(emp.emp_salary);

dump emp_max_sal;

## 8.2.6

## LIMIT

first_five_emp = limit emp 5;

dump first_five_emp;

## 8.2.7

## STORE

store first_five_emp INTO './pig_output/' USING PigStorage(',');

## Directory pig_output created by STORE function
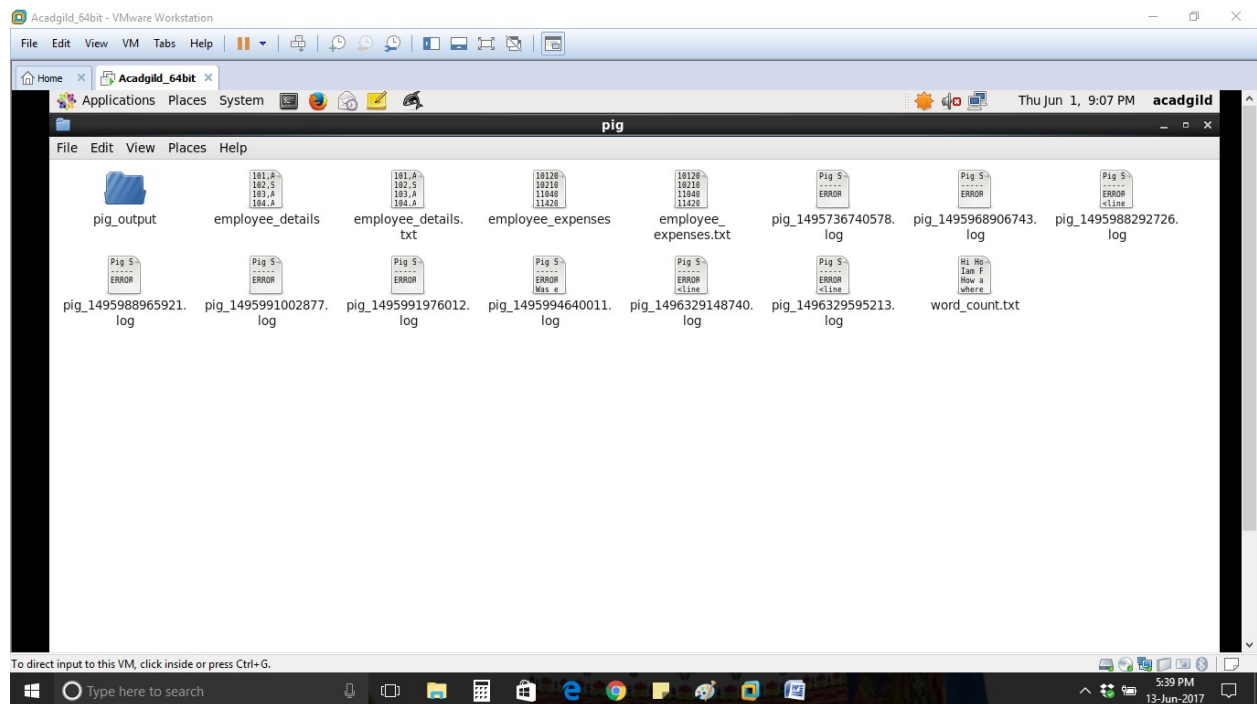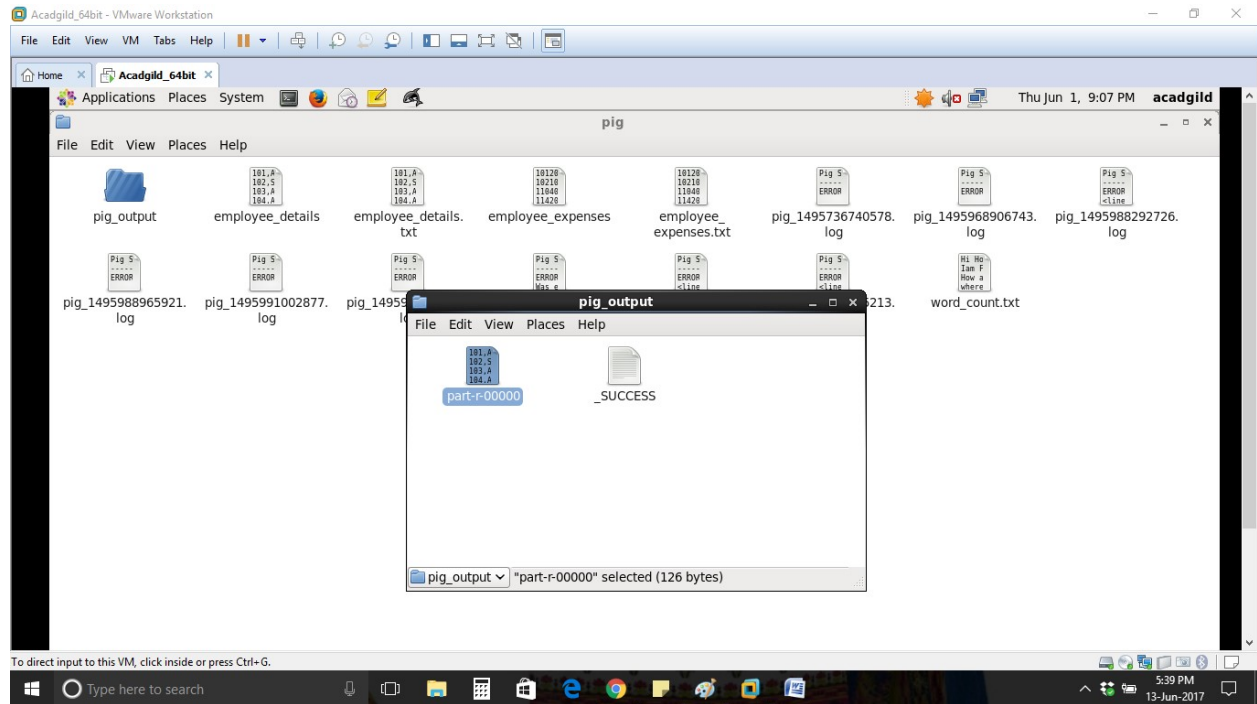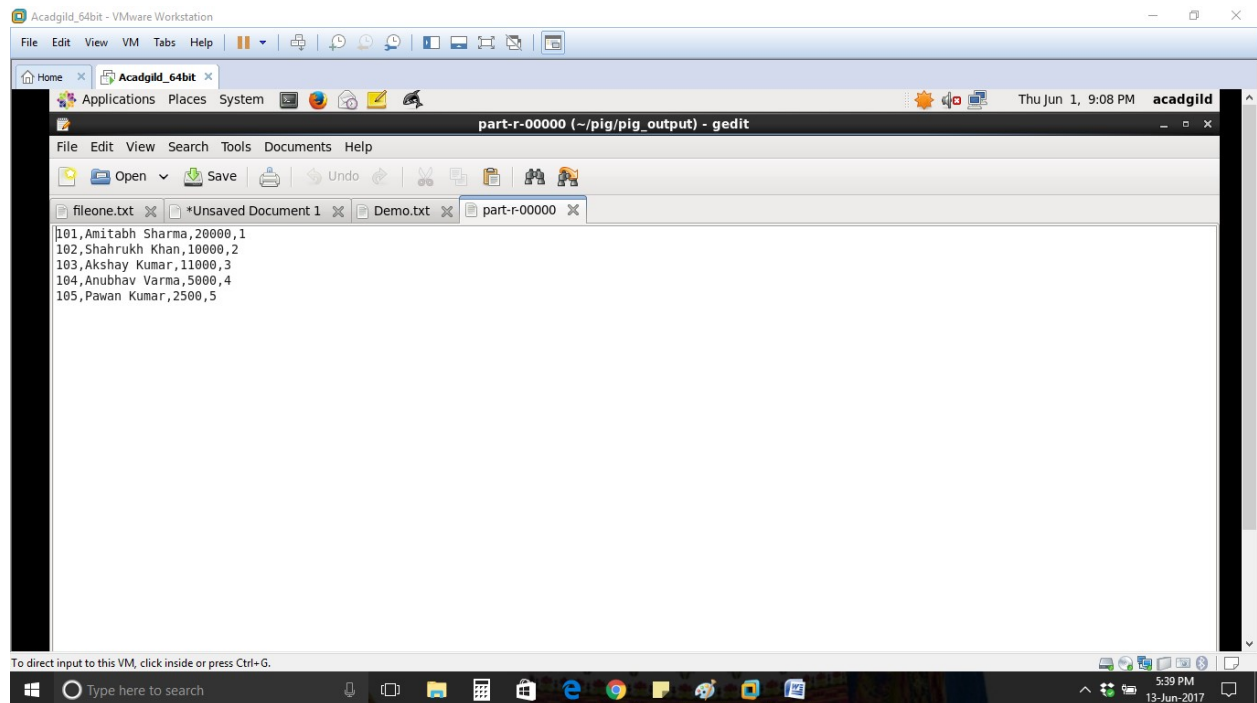
**Two files created inside the pig_output folder**
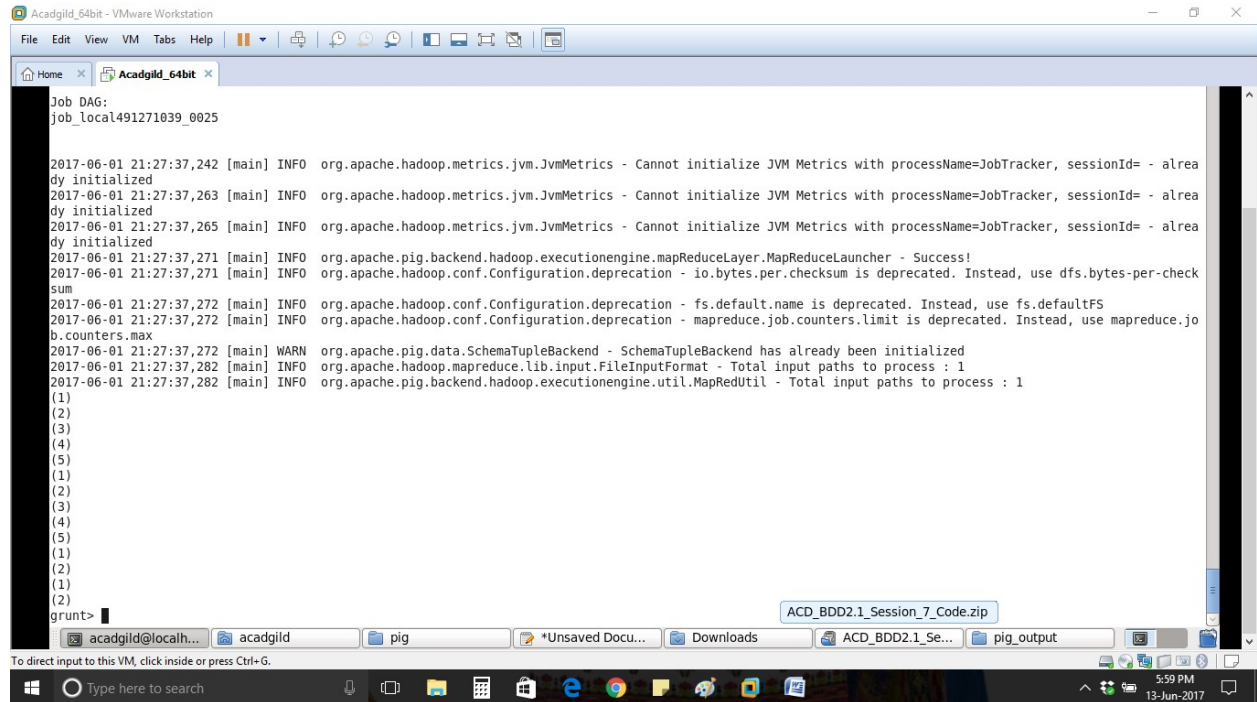


**The file contains first five employee details**

**8.2.8**

**DISTINCT**

emp_by_rating = foreach emp generate emp_rating as rating;

dump emp_by_rating;

**Employee Ratings listed**

emp_distinct_rating = distinct emp_by_rating;

dump emp_distinct_rating
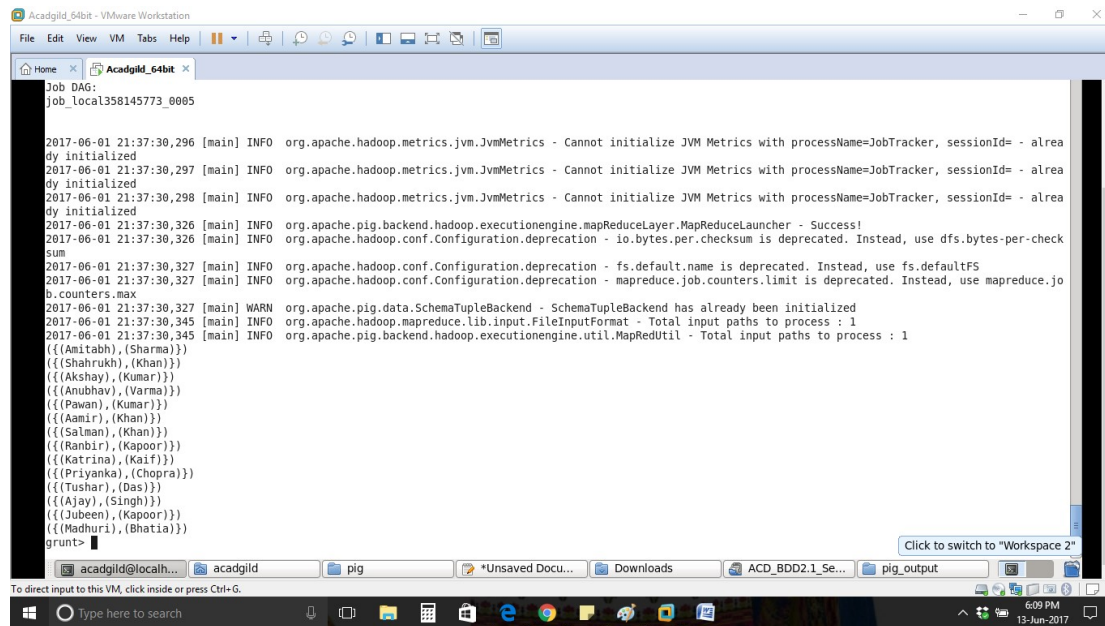
**Only Distinct values in Rating Loaded avoiding duplicates**

**8.2.9**

**FLATTEN**

emp_name_tokenize = foreach emp generate TOKENIZE(emp_name);

dump emp_name_tokenize;

**employee name details tokenized**

B = foreach emp_name_tokenize generate flatten($0);


Dump B;

**FLATTEN function on Tokenized name**



**8.2.10**

**IsEmpty()**

```
emp_no_rating = filter emp by IsEmpty(emp_rating);
```


**Giving error iterator could not be started.**