# SESSION 5 ASSIGNMENT 1

Task_1:
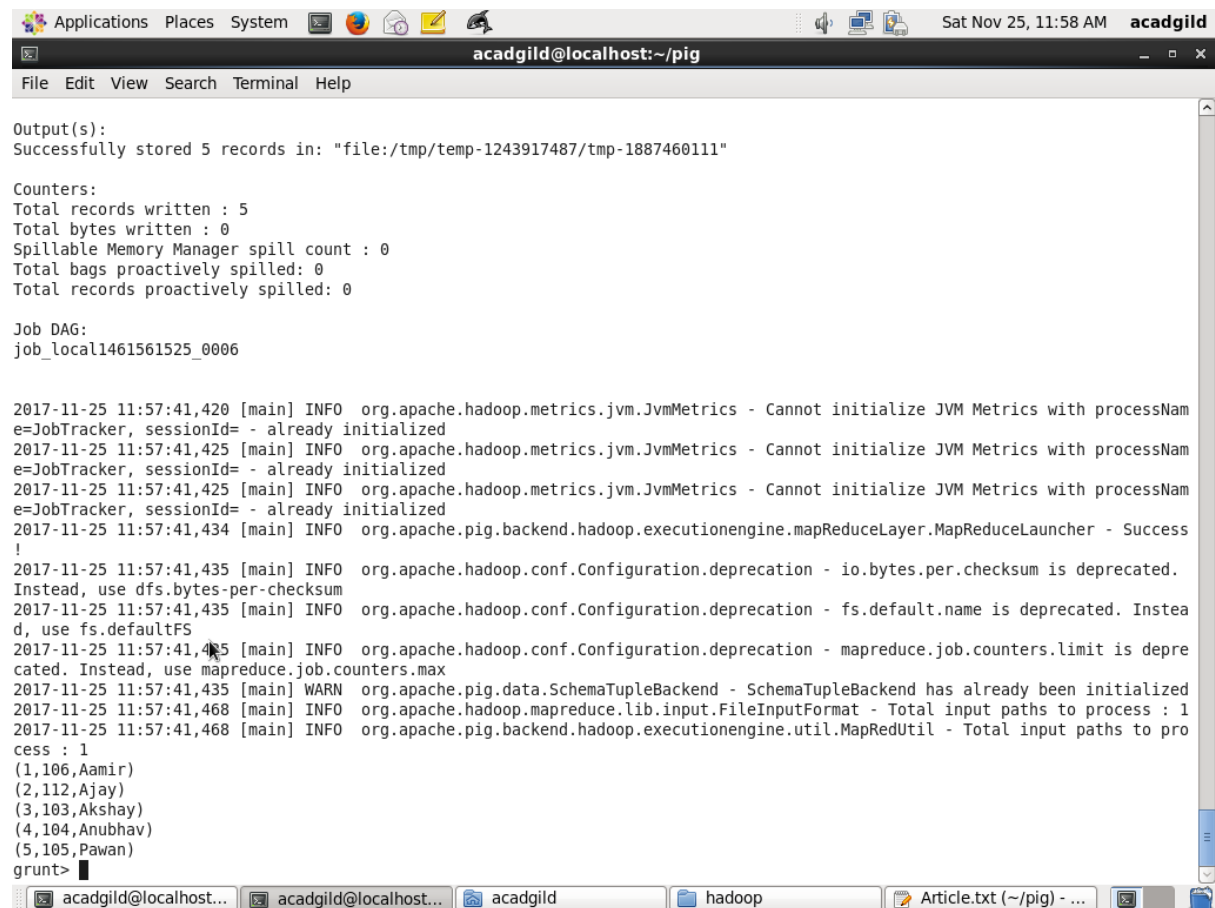
```
grunt>emp_group = GROUP emp_details by rating;

grunt> emp_5 = FOREACH emp_group {

>> da = ORDER emp_details BY rating DESC, name;

>> db = LIMIT da 1;

>> GENERATE FLATTEN(group), FLATTEN(db.id), FLATTEN(db.name);

>> }
```
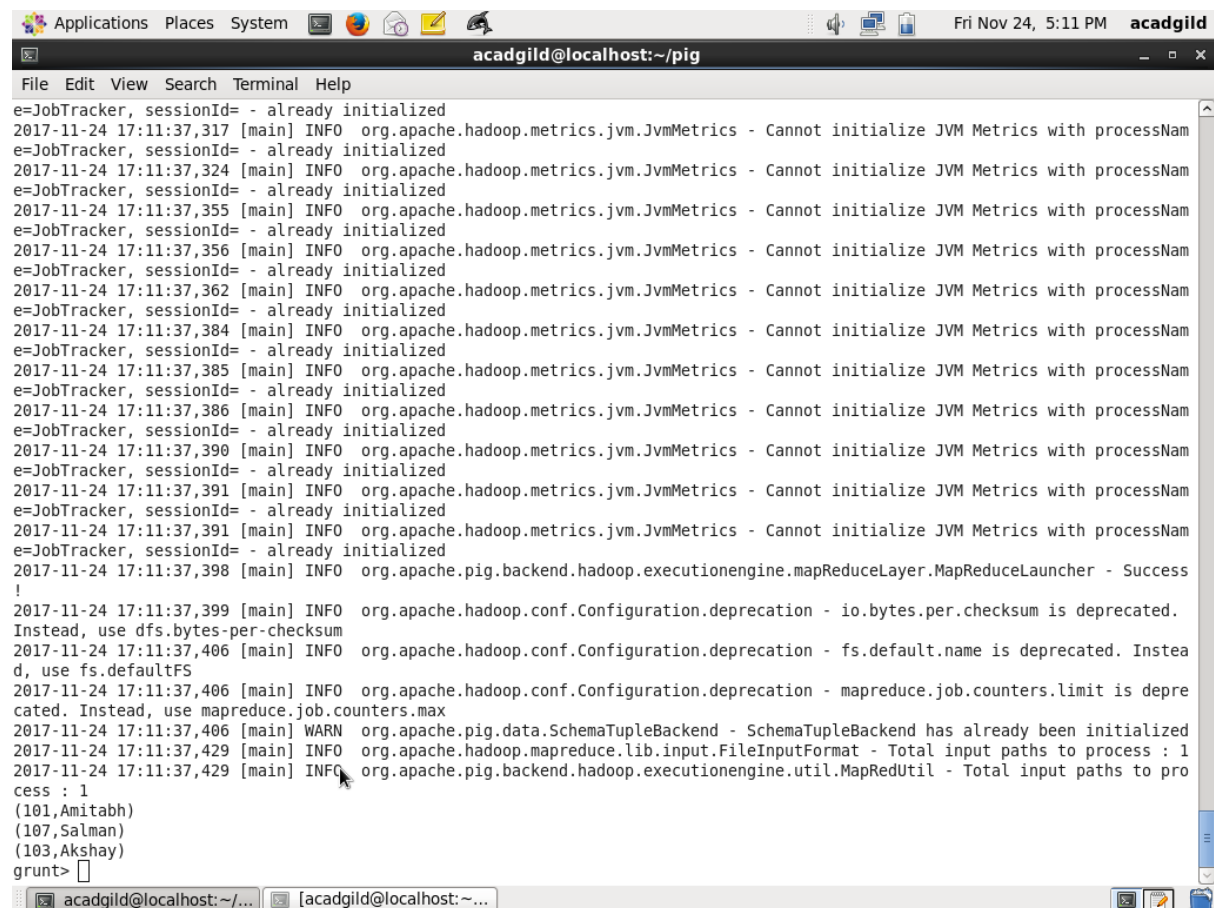
Task_2:

```
grunt> emp_odd = FILTER emp_details BY (id%2 != 0);

grunt> emp_odd_order = ORDER emp_odd BY salary desc, name;

grunt> emp = LIMIT emp_odd_order 3;

grunt> emp_result = FOREACH emp GENERATE id,name;
```
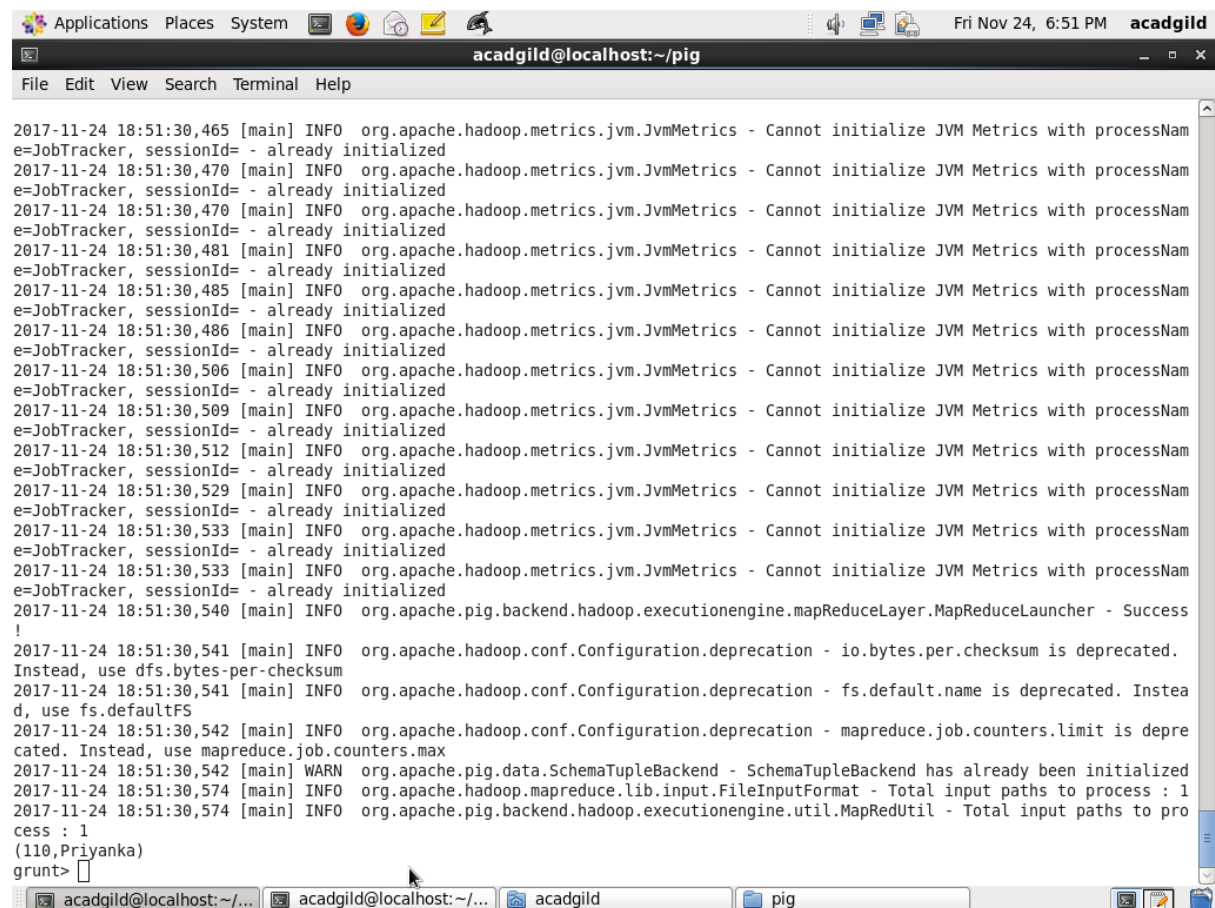
Task_3:

```
grunt> A = JOIN emp_details BY id, emp_expenses BY id;

grunt> B = ORDER A BY  expense DESC,name;

grunt> C = LIMIT B 1;

grunt> D= FOREACH C GENERATE emp_details::id,name;
```

Task_4:

```
grunt> joined = JOIN emp_details BY id, emp_expenses BY id;

grunt> emp = FOREACH joined GENERATE emp_details::id, name;

grunt> emp_data = DISTINCT emp;
```

Task_5:

```
grunt> joined = GROUP emp_details BY id, emp_expenses BY id;

grunt> a= FILTER joined BY IsEmpty(emp_expenses);

grunt> b = FOREACH a GENERATE FLATTEN(emp_details.id), FLATTEN(emp_details.name);
```