CLASS 9 NOTES

```
data a;
  input id;
  cards;
1
1
1
1
2
3
3
;
run;
data final;
  set a;
  by id;
  if first.id=1 and last.id=1 then delete;
run;

data final;
  set sasuser.admit;
  if age gt 40 then delete;
run;
```

Output - (Untitled)    Log - (Untitled)    Editor - Untitled1 *

Edit View Tools Run Solutions Window Help

```sas
data a;
  input id;
  cards;
1
1
1
1
2
3
3
;
run;
data dups nodups;
  set a;
  by id;
  if first.id=1 and last.id=1 then outppt nodups;
  else output dups;
run;
```

Work
Final    Profile

Output - (Untitled)    Log - (Untitled)    Editor - Untitled1 *    Editor - Untitle

Address

Program :

data a;
input id marks;
cards;

| id | marks | first.id | last.id | first.marks | last.marks |
|----|-------|----------|---------|-------------|------------|
| 1 | 20 | 1 | 0 | 1 | 1 |
| 1 | 21 | 0 | 0 | 1 | 1 |
| 1 | 22 | 0 | 1 | 1 | 1 |
| 2 | 25 | 1 | 0 | 1 | 1 |
| 2 | 26 | 0 | 0 | 1 | 1 |
| 2 | 27 | 0 | 0 | 1 | 1 |
| 3 | 30 | 1 | 1 | 1 | 1 |

;
run;

Suppose if we want the non-duplicates values to be deleted and duplicate values to be printed or vice-versa. then,

data final;
set a;
by id;                    → If we want duplicates.
if first.id = 1 and last. id = 1 then delete;
run;

## Another example

```
data a;
input id;
cards;
    1
    1
    1
    2
    3
    3
;
run;

data dups nodups;
set a;
by id;
if first.id = 1 and last.id = 1 then output
                                        nodups;
else output dups;
run;
```
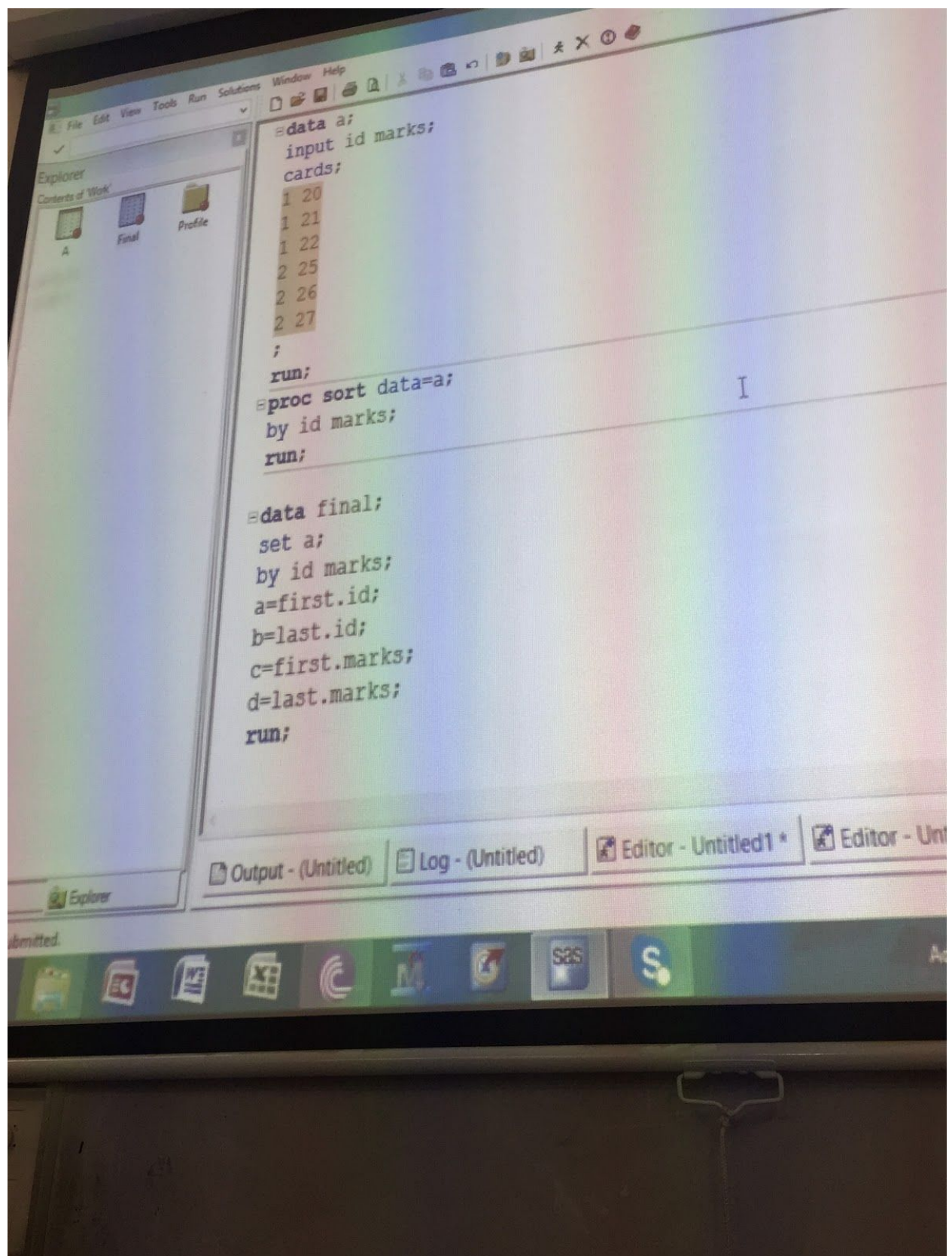
```
data a;
  input id marks;
  cards;
1 20
1 21
1 22
2 25
2 26
2 27
;
  run;
proc sort data=a;
  by id marks;
  run;

data final;
  set a;
  by id marks;
  a=first.id;
  b=last.id;
  c=first.marks;
  d=last.marks;
  run;
```

Output - (Untitled)    Log - (Untitled)    Editor - Untitled1 *    Editor - Unt

Explorer

bmitted.

```
data a;
input id marks;
cards;
1    20
1    21
1    22
2    25
2    26
2    27
;
run
_____

proc sort data = a;
by id marks;
run;
_____

data final;
set a;
by id marks;
a = first.id;
b = last.id;
c = first.marks;
d = last.marks;
run;
```

first.id    last.id

| first.id | last.id |
|----------|---------|
| 1 | 0 |
| 0 | 1 |
| 0 | 1 |
| 0 | 0 |
| 1 | 0 |
| 0 | 1 |

first.marks    last.marks

| first.marks | last.marks |
|-------------|------------|
| 1 | 1 |
| 1 | 1 |
| 1 | 1 |
| 1 | 1 |
| 1 | 1 |
| 1 | 1 |

→ whatever variable we put in by, their first dot and last dot are made.

Here we have two groups – primary(id) and secondary (marks)

'Id' will work according to

the primary group and marks being the secondary
group will work according to id.

## Program

```
data a;
input id marks;
cards;
 1    20
 1    21
 1    29
 2    25
 2    26
 2    27
 ;
run;

proc sort data = a;
by id marks;          ⟶ This will run.
run;

data final;
set a;                ⟶ This will show error because
by marks;                the dataset is not sorted by
run;                     marks, it is sorted by id.
```

File Edit View Tools Run Solutions Window Help

Explorer
Contents of 'Work'

A    Final    Profile

```
data a;
  input id marks;
  cards;
1 20
1 21
1 29
2 25
2 26
2 27
;
run;
proc sort data=a;
  by id marks;
run;

data final;
  set a;
  by marks;

run;
```

```
data a;
  input id marks;
  cards;
1 20
1 21
1 29
2 25
2 26
2 27
;
run;
proc sort data=a;
by descending id;  ;
run;


data final;
set a;
by id;
run;
```

Output - (Untitled)    Log - (Untitled)    Editor - Untitled1 *    Editor -

## Program

```
data a;
input id marks;
cards;
1    20
1    21
1    29
2    25
2    26
2    27
;
run;

proc sort data=a;
by descending id;
run;

data final;
set a;
by id;
run;
```
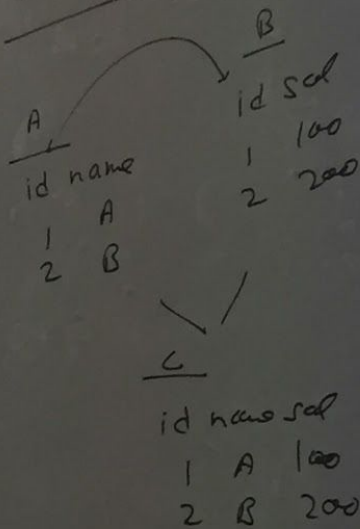
Here 'id' is sorted in descending order.

This has to match

Here in ascending order. so, will show error.

Merging in SAS / VLOOKUP / JOINS

what?

A
---
id name
  A
1
2  B

B
---
id Sal
1  100
2  200

C
---
id name Sal
1  A  100
2  B  200

why? (To Consolidate data), to make it information Rich.

Prerequisite:-

1 :- Common Primary key. (id)

2 :- Sorted by the Primary key B/D

3 :- The type of Primary key to be same.

MERGING in SAS / vlookup (Excel) / Join (Sql)

what is merging :

Merging is combining two or more data sets "head on"
to make a child data set that has attributes
(variables) from all the parent data set.

| A | |
|---|---|
| id | name |
| 1 | A |
| 2 | B |

| B | |
|---|---|
| id | sal |
| 1 | 100 |
| 2 | 200 |

Merge

| C | | |
|---|---|---|
| id | name | sal |
| 1 | A | 100 |
| 2 | B | 200 |

* In general, Appending increase rows and
merging increases columns.

Here, C is the merged product of A and B or
C is the child of A and B.

## Why do we merge?

In order to consolidate data or to make it information rich, we do merging.

## What are the pre-requisite (requirements) of merging:

1. Common key (eg: 'id')

2. Sorted by the primary key (either ascending or descending).

3. The type of primary key should be same, (either character or numeric) eg-if it numeric on one side then it should be numeric on other side also.

## How to do merging:

```
data a;
input id name $;
cards;
 1    A
 2    B
;
run;
```

```
data b;
input id sal;
cards;
 1    100
 2    200
;
run;
```

# Merge

data c;     → child of 'a' and 'b'.

| merge a b; |    → Merge statement

by id;     → sorted by 'id'

run;

## Output

| id | name | sal |
|----|------|-----|
| 1 | A | 100 |
| 2 | B | 200 |

## Types of merging

① 1 to 1 merging

Eg:

| id | name |
|----|------|
| (1 | A |
| (2 | B |

↓
two
bygroups

| id | sal |
|----|-----|
| (1 | 100 |
| (2 | 200 |

| id | name | sal |
|----|------|-----|
| 1 | A | 100 |
| 2 | B | 200 |

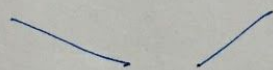Here, by groups are two and each by group has single row in it. So, this type of merging is called __1-1 merging__.

Note : In 1-1 merging, no. of __rows__ in the by group should be __one__.

Another eg:

① 

| id | name |
|----|------|
| 1 | A |
| 2 | B |
| 3 | C |

| id | sal |
|----|-----|
| 1 | 100 |
| 2 | 200 |

| id | name | sal |
|----|------|-----|
| 1 | A | 100 |
| 2 | B | 200 |
| 3 | C | . → numeric field |

---

②

| id | name |
|----|------|
| 1 | A |
| 2 | B |

| id | sal |
|----|-----|
| 1 | 100 |
| 2 | 200 |
| 7 | 300 |

→

| id | name | sal |
|----|------|-----|
| 1 | A | 100 |
| 2 | B | 200 |
| 7 | □ | 300 |

↓
character field

② 1 - many merging

```
Data c;
merge a b;
by id;
run;
```

| A | |
|---|---|
| id | name |
| 1 | A |
| 2 | B |
| 3 | C |

one to many

| B | |
|---|---|
| id | sub |
| 1 | H |
| 1 | M |
| 2 | E |
| 3 | M |
| 3 | S |

Output

C

| id | name | sub |
|---|---|---|
| 1 | A | H |
| 1 | A | M |
| 2 | B | E |
| 3 | C | M |
| 3 | C | S |

3 bygroups

No. of bygroups are three, but in each bygroup no., of rows are one or more than one, so one to many merging.

③ Many to 1 merging

Data c;
merge b a;
by id;
run;

B

| id | sub |
|----|-----|
| 1 | H |
| 1 | M |
| 2 | E |
| 3 | M |
| 3 | S |

many to 1 →

A

| id | name |
|----|------|
| 1 | A |
| | B |
| 2 | B |
| 3 | C |

C

| id | sub | name |
|----|-----|------|
| 1 | H | A |
| 1 | M | A |
| 2 | E | B |
| 3 | M | C |
| 3 | S | C |

## 1-1 merging

**I.**

A
| id | name |
|----|------|
| 1  | A    |
| 2  | B    |
| 3  | C    |

B
| id | Sal |
|----|-----|
| 1  | 100 |
| 2  | 200 |

→

| id | name | Sal |
|----|------|-----|
| 1  | A    | 100 |
| 1  | B    | 200 |
| 2  | C    | .   |
| 3  |      |     |

**II.**

A
| id | name |
|----|------|
| 1  | A    |
| 2  | B    |

B
| id | Sal |
|----|-----|
| 1  | 100 |
| 2  | 200 |
| 7  | 300 |

(by id)

→

| id | name | Sal |
|----|------|-----|
| 1  | A    | 100 |
| 2  | B    | 200 |
| 7  | □    | 300 |

## 1-many merging

I,

A
id name
( 1  A
( 2  B
( 3  C

B
id Sub
( 1  H
( 1  M
( 2  E
( 3  M
( 3  S

→

| id | name | Sub |
|----|------|-----|
| 1  | A    | H   |
| 1  | A    | M   |
| 2  | B    | E   |
| 3  | C    | M   |
| 3  | C    | S   |

```
Data C;
merge A B;
by id;
run,
```

## merge merging

I;

A
| id | name |
|----|------|
| 1  | A    |
| 2  | B    |
| 3  | C    |

B
| id | sub |
|----|-----|
| 1  | H   |
| 1  | M   |
| 2  | E   |
| 3  | M   |
| 3  | S   |

→

| id | name | sub |
|----|------|-----|
| 1  | A    | H   |
| 1  | A    | M   |
| 2  | B    | E   |
| 3  | C    | M   |
| 3  | C    | S   |

```
Data C;
merge B A
by id;
run;
```

```
data a;
  input id name $;
  cards;
. A
2 B
;
run;

data b;
  input id sal;
  cards;
1 100
2 200
;
run;

data c;
  merge a b;
  by id;
run;
```

Profile

Output - (Untit...  Log - (Untitled)  Editor - Untitle...  Edit

File   Edit   View   Tools   Run   Solutions   Window   Help

```
data a;
  input id name $;
  cards;
1 A
2 B
;
run;

data b;
  input id sal;
  cards;
1 100
2 200
;
run;

data c;
  merge a b;
  by id;
run;
```

Output - (Untitle...   |   Log - (Untitled)   |   Editor - Untitled1 *   |   Edito

Submitted.

# Interview Question (CASE STUDY - 1)

Merge behaves like _append_ when you _don't_ 'a

have _common_ values.

Eg:

```
data a;                      data b;
input id name$;              input id sal;
cards;                       cards;
1  A                         11   100
2  B                         22   200
;                            ;
run;                         run;
```

```
data c;
merge a b;
by id;
run;
```

Output

| id | name | sal |
|----|------|-----|
| 1  | A    |     |
| 2  | B    |     |
| 11 |      | 100 |
| 22 |      | 200 |

# CASE STUDY - 2

```
data a;
input id name $ age;
cards;
1    A    12
2    B    15
;
run;
```

```
data b;
input id sal age;
cards;
1    100    45
2    200    55
;
run;
```

```
data c;
merge a b;
by id;
run;
```

Here, we have four variables — id, name, age, sal.

But, the issue is whether the age from dataset 'a' will come or dataset 'b'.

Acc., to concept of merging overlepping is done.

Here, the age from dataset 'b' will overwrite dataset 'a' age's value.

But here also the condition is that dataset 'b' should have the updated values of age, only then it will overwrite the age value of 'a'.

Output

| id | name | sal | age |
|----|------|-----|-----|
| 1  | A    | 100 | 45  |
| 2  | B    | 200 | 55  |

# CASE STUDY - 3

Suppose, from the previous example, if dataset 'a' have the updated value of age and we don't want dataset 'b' age's value to overwrite 'a'. In that case, drop the variable of <u>age</u> from dataset 'b', so that it should not overwrite.

```
data c;
merge a b (drop = age);    →  By this, output
by id;                         will not be
run;                           overlapped.
```

Output

| id | name | sal | age |
|----|------|-----|-----|
| 1  | A    | 100 | 12  |
| 2  | B    | 200 | 15  |

'age' from dataset 'a'.

Program

data ctc;
input empid sal;
cards;

1     100
2     300
3     900
;
run;

---

data new;
input empid sal;
cards;

1     100
1     110
1     120
1     135
2     395
3     900
3     950
;
run;

Here, we want to update the data of 'ctc' by 'new' by emp id (which is sorted) means with every empid, the updated salary should come.

1st updated salary = 135

2s    "        "    = 395

3s    "        "    = 950

↓

Program → data ctc;
          update ctc new;
          by emp id;
          run;

Output

|   | empid | sal |
|---|-------|-----|
| 1 | 1     | 135 |
| 2 | 2     | 395 |
| 3 | 3     | 950 |