

# Selecting All Observations When Any Observation Is of Interest

Christopher Bost

## What does that mean?

- Some data have more than one obs/person
- We want to select all obs for a person
  - If at least one obs for that person meets criteria
- Examples?

## What's in it for you?

- Review how to select obs with the DATA step
- Review how to select obs with PROC SQL
- Useful techniques to have in your SAS "toolkit"



## Sample data set: COURSES

studentid	course	ap
1	BIOL101	0
2	BIOL102	1
2	CHEM102	1
3	CALC101	0
3	STAT101	0
4	PSYC201	1
4	HIST102	0
4	PHYS101	1
5	CHEM101	0
5	CALC102	1

5

## Desired data set

studentid	course	ap
2	BIOL102	1
2	CHEM102	1
4	PSYC201	1
4	HIST102	0
4	PHYS101	1
5	CHEM101	0
5	CALC102	1

6

## Select observations with DATA step

1. Subset observations of interest
2. Keep one observation per person
3. Match-merge with original data set

7

## 1. Subset observations of interest

```
data ap;  
set courses;  
where ap=1;  
run;
```

Data set AP

studentid	course	ap
2	BIOL102	1
2	CHEM102	1
4	PSYC201	1
4	PHYS101	1
5	CALC102	1

8

## 2. Keep one observation per person

```
proc sort data=ap out=ap_sort;  
by studentid;  
run;
```

```
data ap2;  
set ap_sort;  
by studentid;  
if first.studentid;  
keep studentid;  
run;
```

Data set AP2

studentid
2
4
5

9

## 3. Match-merge observations

```
proc sort data=courses out=courses_sort;  
by studentid;  
run;
```

```
data anyap;  
merge courses_sort ap2(in=inap2);  
by studentid;  
if inap2;  
run;
```

10

## Desired data set: ANYAP

studentid	course	ap
2	BIOL102	1
2	CHEM102	1
4	PSYC201	1
4	HIST102	0
4	PHYS101	1
5	CHEM101	0
5	CALC102	1

11

## Full DATA step program

```
data ap;
set courses;
where ap=1;
run;

proc sort data=ap
          out=ap_sort;
by studentid;
run;

data ap2;
set ap_sort;
by studentid;
if first.studentid;
keep studentid;
run;

proc sort data=courses
          out=courses_sort;
by studentid;
run;

data anyap;
merge courses_sort
      ap2(in=inap2);
by studentid;
if inap2;
run;
```

12

## Pros and cons

- It works
- Three DATA steps
- Two PROC SORT steps



## Select observations with PROC SQL

1. Use a subquery
2. Use GROUP BY and HAVING clauses

15

## Terminology

<b>DATA Step</b>	<b>PROC SQL</b>
Variable	Column
Observation	Row
SAS data set	Table

16



## SQL clauses: required order

<code>PROC SQL;</code>	<i>starts procedure</i>
<b>1 SELECT</b>	<b>selects variables</b>
<b>2 FROM</b>	<b>opens data sets</b>
<b>3 WHERE</b>	<b>restricts observations</b>
<b>4 GROUP BY</b>	<b>groups observations</b>
<b>5 HAVING</b>	<b>restricts groups</b>
<b>6 ORDER BY ;</b>	<b>sorts results</b>
<code>QUIT;</code>	<i>ends procedure</i>

17

## Use a subquery

```
proc sql;  
create table anyap2 as  
3 select *  
1 from courses  
2 where studentid in (select distinct studentid  
from courses  
where ap=1)  
4 order by studentid;  
quit;
```

18

## Desired data set: ANYAP2

studentid	course	ap
2	BIOL102	1
2	CHEM102	1
4	PSYC201	1
4	PHYS101	1
4	HIST102	0
5	CALC102	1
5	CHEM101	0

19

## Pros and cons

- Single step
- Two passes through data set
  - One for the subquery
  - One for the outer query

20



## SQL clauses: order of execution

<b>1 FROM</b>	<b>opens data sets</b>
<b>2 WHERE</b>	<b>restricts observations</b>
<b>3 GROUP BY</b>	<b>groups observations</b>
<b>4 HAVING</b>	<b>restricts groups</b>
<b>5 SELECT</b>	<b>selects variables</b>
<b>6 ORDER BY</b>	<b>sorts results</b>

## Use GROUP BY and HAVING clauses

```
proc sql;  
  create table anyap3 as  
4 select *  
1 from courses  
2 group by studentid  
3 having sum(ap=1) > 0  
5 order by studentid;  
quit;
```

23

## Desired data set: ANYAP3

studentid	course	ap
2	CHEM102	1
2	BIOL102	1
4	HIST102	0
4	PSYC2101	1
4	PHYS101	1
5	CALC102	1
5	CHEM101	0

24

## Pros and cons

- Single step
- **Flexibility**
  - Any condition(s) in parentheses after SUM
- **NOTE: The query requires remerging summary statistics back with the original data.**

25



## Conclusion

- Use either the DATA step or PROC SQL
- PROC SQL requires significantly less coding
- PROC SQL is not necessarily more efficient
  - Test on your own data
- Use on **any data** with groups of observations

27



## Contact information

Comments and questions are valued and encouraged.

Christopher J. Bost  
MDRC  
16 East 34<sup>th</sup> Street  
New York, NY 10016  
(212) 340-8613  
christopher.bost@mdrc.org  
chrisbost@gmail.com



29

## Alternate DATA step method

```
proc sort data=courses (where=(ap=1) keep=studentid ap)
  out=lookup (keep=studentid)
  nodupkey;
by studentid;
run;

proc sort data=courses out=courses_sort;
by studentid;
run;

data anyap;
merge courses_sort lookup(in=inlookup);
by studentid;
if inlookup;
run;
```

30

## Alternate PROC SQL method

```
proc sql;
create table anyap4 as
select courses.*
from courses inner join
    (select distinct studentid as studentid2
     from courses
     where ap=1)
on studentid=studentid2
order by studentid;
quit;
```