

---

# Global School-Based Student Health Survey (GSHS)

Data Processing and Weighting

---

# Overview

- GSHS cleaning, weighting and analysis\* have been standardized to ensure the integrity of the global system and comparability of results.
- WHO has developed a set of R code to automate these tasks.
- These slides present how GSHS data is cleaned and weighted and the input files needed to use the R code.

---

# Data Cleaning

---

# Overview

- As part of the standard cleaning process, GSHS data are checked for quality and consistency.
- Checking occurs at both the variable level as well as at the record level.
  - Thus, it is possible for a **variable** to be dropped or a student's **entire response** to be dropped.
- The following slides in this section describe the data edits performed.

---

# Out-of-range edits

- If a student selects a response that does not correspond to one of the possible responses for a question, then the response is set to missing.
- *Example:* If “A” and “B” are the valid response options for a question and a student selects “C”, “D”, “E”, “F”, “G”, or “H,” then the response is set to missing for that student.

---

# Multi-response edits

- If a student selects more than one response for a question, then the question is set to missing.
- GSHS questions never allow for multiple responses.

---

## BMI-related edits

- BMI is calculated using the height and weight measures reported in each student's response.
  - If either height or weight are missing, BMI is set to missing.
- Height, weight and BMI are checked to see if they are outside the biologically plausible range\*.
  - If height, weight or BMI are implausible, all are set to missing.
  - If age or sex is missing, height, weight and BMI are set to missing since plausible ranges vary by age and sex and it would be impossible to determine plausibility.

---

# Logical consistency edits

- Logical consistency checks are made for questions in 6 of the core modules.
- These checks ensure that responses are *internally consistent*
  - Example: If a student responds that they did not clean their teeth in the past 30 days but in the following question responds that the toothpaste usually used to brush their teeth in the past 30 days contains fluoride – these responses would not be internally consistent.
- If a check fails, then the responses to **both** questions are set to missing *except* if one of the questions is AGE (age is never set to missing).
- Consistency checks are not exhaustive and there are no consistency checks done for core-expanded or country-specific questions.
- All 46 edits are listed on the following slide.



## Hygiene

1. HY\_CLTEETH = A AND HY\_FLUORIDE = B,C,D

## Injury

2. IN\_TIMESINJ = A AND IN\_TYPEINJ = B, C, D, E, F, G, H
3. IN\_TIMESINJ = A AND IN\_CAUSEINJ = B, C, D, E, F, G, H

## Tobacco Use

4. TO\_TRIEDCIG = B AND TO\_AGE CIG = B, C, D, E, F, G, H
5. TO\_TRIEDCIG = B AND TO\_DAYS CIG = B, C, D, E, F, G
6. DE\_AGE = A AND TO\_AGE CIG = E,F,G,H
7. DE\_AGE = B AND TO\_AGE CIG = F,G,H
8. DE\_AGE = C AND TO\_AGE CIG = F,G,H
9. DE\_AGE = D AND TO\_AGE CIG = G,H
10. DE\_AGE = E AND TO\_AGE CIG = G,H
11. DE\_AGE = F AND TO\_AGE CIG = H
12. DE\_AGE = G AND TO\_AGE CIG = H

## Alcohol Use

13. AL\_AGE = A AND AL\_DAYS = B, C, D, E, F, G
14. AL\_AGE = A AND AL\_DRINKS = B, C, D, E, F, G
15. AL\_AGE = A AND AL\_INAROW = B, C, D, E, F, G, H
16. AL\_AGE = A AND AL\_SOURCE = B, C, D, E, F
17. AL\_AGE = A AND AL\_TROUBLE = B, C, D, E, F
18. AL\_AGE = A AND AL\_DRUNK = B, C, D, E, F
19. DE\_AGE = A AND AL\_AGE = E,F,G,H
20. DE\_AGE = B AND AL\_AGE = F,G,H
21. DE\_AGE = C AND AL\_AGE = F,G,H
22. DE\_AGE = D AND AL\_AGE = G,H
23. DE\_AGE = E AND AL\_AGE = G,H
24. DE\_AGE = F AND AL\_AGE = H
25. DE\_AGE = G AND AL\_AGE = H

## Drug Use

26. DR\_AGE = A AND DR\_CANLIFE = B, C, D, E, F
27. DR\_AGE = A AND DR\_CAN30 = B, C, D, E, F
28. DR\_AGE = A AND DR\_AMPHLIFE = B, C, D, E, F
29. DE\_AGE = A AND DR\_AGE = E,F,G,H

30. DE\_AGE = B AND DR\_AGE = F,G,H
31. DE\_AGE = C AND DR\_AGE = F,G,H
32. DE\_AGE = D AND DR\_AGE = G,H
33. DE\_AGE = E AND DR\_AGE = G,H
34. DE\_AGE = F AND DR\_AGE = H
35. DE\_AGE = G AND DR\_AGE = H

## Sexual Behaviors

36. DE\_AGE = A AND SX\_AGE = E,F,G,H
37. DE\_AGE = B AND SX\_AGE = F,G,H
38. DE\_AGE = C AND SX\_AGE = F,G,H
39. DE\_AGE = D AND SX\_AGE = G,H
40. DE\_AGE = E AND SX\_AGE = G,H
41. DE\_AGE = F AND SX\_AGE = H
42. DE\_AGE = G AND SX\_AGE = H
43. SX\_EVERSEX = B AND SX\_AGE = B, C, D, E, F, G, H
44. SX\_EVERSEX = B AND SX\_NUMBER = B,C,D,E,F,G
45. SX\_EVERSEX = B AND SX\_CONDOM = B,C
46. SX\_EVERSEX = B AND SX\_BC = B,C,D,E,F,G,H

---

## Variable-level edits

- After all other checks have been implemented, each variable is checked to ensure at least 60% of students have responded.
- If the response rate for a variable is less than 60%, the variable is set to missing for all students.

---

## Record-level edits

- After all other checks have been implemented, each record is checked to ensure there are at least 20 valid responses
- If a student's response has fewer than 20 valid responses, the response is deleted.
- If a record has 15 or more identical responses in a row, other than "A", the entire record is deleted.

---

# Weighting

---

# Overview

- Once data have been cleaned, the weighting process can begin.
- Weighting accounts for:
  - the probability of selection of schools and classes
  - non-responding schools, classes and students
  - the distribution of the target population (i.e. students in the targeted grades) by grade and sex
- In addition to analysis weights, PSU and Stratum will also be generated which inform the statistical software about the design of your sample.

---

# Requirements

- All of the following conditions must be met in order to weight GSHS data:
  - the sample was scientifically selected from an up-to-date and complete sampling frame
  - all school-level and class-level forms were accurately completed
  - a high (>60%) overall response rate was obtained

---

# Weight calculation

- The formula used to calculate analysis weights for most GSHS data sets is

$$\text{weight} = w1 * w2 * f1 * f2 * f3$$

where:

**Base weight**  $\left\{ \begin{array}{l} w1 = \text{the inverse probability of selecting each school} \\ w2 = \text{the inverse probability of selecting each class} \end{array} \right.$

**Non-response adjustment**  $\left\{ \begin{array}{l} f1 = \text{a school-level non-response adjustment factor} \\ f2 = \text{a student-level non-response adjustment factor (calculated per class)} \end{array} \right.$

**Post-stratification adjustment**  $f3 = \text{a post-stratification adjustment factor (calculated by sex within each grade)}$

---

# PSU and Stratum

- PSU and Stratum describe the complex sample design of the survey
- These numbers are generated as follows:
  - **Schools selected with certainty\*** : assign a unique stratum to each school and a unique PSU to each class in each school
  - **All other schools** : sort schools by school weight\*\* and group schools into pairs (if there is an odd number, make one group of three), assign a unique stratum to each pair of schools (or group of three) and a unique PSU to all classes within a given school



# PSU and Stratum - example

School Weight	School	Classes	Stratum	PSU	
1.0	A	1	1	1	}
		3	1	2	
		6	1	3	
1.0	B	2	2	4	
		4	2	5	
		6	2	6	
1.0	C	1	3	7	}
		3	3	8	
		4	3	9	
		6	3	10	
1.27	E	1	4	11	
		2	4	11	
		3	4	11	}
1.38	F	1	4	12	
		3	4	12	
		5	4	12	
1.79	G	2	5	13	
		4	5	13	
		6	5	13	}
		8	5	13	
1.83	H	1	5	14	
		3	5	14	
		5	5	14	
1.90	I	3	5	15	
		6	5	15	}
		9	5	15	

**Schools selected with certainty**

**All other schools (i.e. smaller schools)**