# **TRAINING ON DATA MANAGEMENT** & ANALYSIS OF STD/HIV DATA

**RESOURCE BOOK** 





MINISTRY OF HEALTH SRI LANKA







# DATA MANAGEMENT AND ANALYSIS OF HIV/AIDS DATA

# **RESOURCE BOOK**

National STD/AIDS Control Programme (NSACP), Sri Lanka

&

VHS-CDC Project

The Voluntary Health Services (VHS), India

Supported by Centers for Disease Control and Prevention (CDC/DGHT-India)

T.T.T.I. Post, Rajiv Gandhi Salai,

Taramani, Chennai – 600 113, Tamil Nadu, INDIA.

Ph.: +91-44-22541965 | Email: vhs.cdcproject@gmail.com



MINISTRY OF HEALTH SRI LANKA







Book Title	:	Resource Book on Data Management and Analysis of HIV/AIDS Data.
No. of pages	:	233 pages including wrapper pages
Year of Publication	:	August 2019
Technical Guidance & Mentoring team	:	Dr Joseph D Williams, Director Projects, VHS Dr Ariyaratne Manathunge, Consultant-Venereologist, SIMU, NSACP
Team of Authors	:	Dr T Ilanchezhian, Senior Technical Advisor, VHS-CDC Project Dr Yujwal Raj, Technical Advisor (SI), VHS-CDC Project
Compiled & designed by	:	Ms T Sudha, Senior Programme Associate, VHS-CDC Project
© Copyright	:	VHS-CDC Project, Voluntary Health Services, Chennai–113.

## Сорукіснт

This document is primarily meant for internal purposes and reproducing in any form requires written permission.

This publication was supported by the Grant or Cooperative Agreement Number 6 NU2GGH001087-05-02, funded by the Centers for Disease Control and Prevention, (CDC) US and PEPFAR. Its contents are solely the responsibility of the authors and do not necessarily represent the official views of the Centers for Disease Control and Prevention, US PEPFAR or the Department of Health and Human Services.

## INDEX

S. No.	Contents				
	Acronyms	5			
	Foreword	7			
	Acknowledgement	9			
I	Data Management and Analysis of HIV/AIDS Data – An Introduction	11			
П	Presentations and exercise formats:	12			
1	Identifying & mapping data sources for HIV/AIDS analysis	13			
1.1	EXERCISE 1: Identifying Programmatic Questions & Mapping Data Sources				
	Part-A: List all the Known HIV/AIDS Data Sources & Identify the Key Information	19			
	Areas				
	Part-B: Mapping Data Sources/ Datasets with Programmatic Questions	20			
2	Principles of Database Management	21			
2.1	EXERCISE 2: Exercise Dataset 2	36			
	EXERCISE 2: Reviewing the Database	37			
2.2	Data Management Systems - Assessment Tool	38			
2.3	Data Management Checklist	40			
3	Understanding Datasets of NSACP	42			
	Understanding Programme Datasets under NSACP	61			
3.1	EXERCISE 3: Assessment of NSACP Datasets	67			
4	Variables & Indicators	68			
4.1	EXERCISE 4: Variables and Indicators				
	PART A: Review the NSACP datasets - Write the Variable values	80			
	PART B: Classify the Indicators	81			
5	Data Quality Assessment	83			
6	Data Adjustments & Validation	107			
6.1	Exercise Dataset – Excel Sheet	115			
7	Basic Programmatic Analysis – Measures & Methods	120			
8	Basic & Advanced Functions in MS Excel	124			
9	Commonly used Formulae in MS Excel	131			
10	HIV/AIDS Specific Analysis from Programme Data	137			
11	Data Triangulation	140			
12	Communication of Data Analysis Results	157			
13	Data Use for Decision Making - A Systems Approach	162			
14	How to present the Data/ Information	175			
15	Introduction to BIOSTATISTICS in Clinical Research	196			
16	Statistical Package for the Social Sciences (SPSS)	220			
17	Key Learnings & Next Steps	231			

# Acronyms

ACASI	Audio Computer Assisted Self Interview
AIDS	Acquired Immunodeficiency Syndrome
ANC	Ante-Natal Care
AR	Attributable Risk
ART	Anti-Retroviral Treatment
CAPI	Computer Assisted Personal Interviewing
CDC	Centers for Disease Control and Prevention
CON'T	Continued
CSS	Cross-Sectional Study
CSS	Case-Control Study
C&S	Care & Support
CST	Care, Support & Treatment
DD	Data Dictionary
DGHT	Division of Global HIV & TB
DM	Data Management
DQA	Data Quality Assurance
EPI Unit	Epidemiology Unit
FcFT	Facilitator cum Feedback Team
FGD	Focus Group Discussion
FSW	Female Sex Worker
GIS	Geographical Information Systems
HIV	Human Immunodeficiency Virus
HSS	HIV Sentinel Surveillance
IBBS	Integrated Biological and Behavioral Surveillance
IDI	In-Depth Interviews
IEC	Information Education Communication
ITDM	International Training on Data Management
КАР	Knowledge Attitude and Practice
КР	Key Population
Litt.	Literature
M&E	Monitoring and Evaluation
MSM	Men who have Sex with Men
NSACP	National STD/AIDS Control Programme
OR	Operational Research
OR	Odds Ratio

PAR	Population Attributable Risk
PEPFAR	President's Emergency Plan for AIDS Relief
РНІ	Public Health Inspector
PHLT	Public Health Laboratory Technician
PHNS	Public Health Nursing Sister
PLHIV	People Living with Human Immunodeficiency Virus
PM	Project Management
РМТСТ	Prevention of Mother To Child Transmission
РРТ	Power-Point Presentation
PrEP	Pre-Exposure Prophylaxis
PRT	Peer Review Team
RCT	Randomized Controlled Trial
REC	Research Ethics Committees
RR	Risk Ratio / Relative Risk
STD	Sexually Transmitted Diseases
STI	Sexually Transmitted Infections
SI	Strategic Information
SIMU	Strategic Information Management Unit
ТА	Technical Assistance
ТВ	Tuberculosis
TNA	Training Needs Assessment
VHS	Voluntary Health Services
Vs.	Versus
WHO	World Health Organization

#### Foreword



Dr Rasanjalee Hettiarachchi Director National STD/AIDS Control Programme (NSACP) Sri Lanka

am happy to write a foreword to this Resource Book on the *"Data Management and Analysis of HIV/AIDS Data"* developed/ compiled by VHS-CDC Project by engaging experts/ consultants in consultation and coordination with SIMU-NSACP team based on the training needs identified. This Resource Book has been brought out by VHS-CDC Project with the support of CDC for capacity building of the SI team including: SIMU team, Consultant-Venereologists & Medical Officers and Public Health Inspectors and Nursing Officers.

This Resource Book on "Data Management and Analysis of HIV/AIDS Data" will be of very much useful for imparting training and using as ready reckoner/ reference material for the participants.

This Resource Book is primarily introducing the basic principles and approaches of Data Management; various methods of Data Quality Assessment; validation; adjustments; basic skills in statistical analysis of program and epidemiological data; software packages used for statistical analysis; and guidelines for presentation, dissemination and use of data for programmatic purposes. This Resource Book will also contribute for enhancing knowledge and skills supported with handson training.

Training and capacity building are the key elements of VHS-CDC Project in providing Technical Assistance to NSACP on Strategic Information with the support of CDC/DGHT-India. VHS-CDC Project with the support of CDC in partnership with NSACP is planning for conducting the training on Data Management at different levels for building the capacity of SI team in the country.

This Resource Book will be used in the training on Data Management and Analysis of HIV/AIDS Data. The training programs are designed to build the data skills of participants in order to enhance the data quality, improve the data analysis and strengthen the use of HIV/AIDS data for epidemiological & programmatic decision making.

On behalf of NSACP, I wish to express my sincere thanks to Dr Joseph D Williams, Director Projects-VHS for his immense support in ensuring partnerships and continue to support in providing TA. We also appreciate the strategic support being extended by Dr T Ilanchezhian, Senior Technical Advisor, VHS-CDC Project for coordinating with NSACP and SIMU in providing technical assistance on strategic information and managing and coordinating the series of capacity building initiatives.

Thanks to VHS-CDC Project team, resource trainer/ consultant for the support extended in development of this Resource Book and contribution in conducting training program.

My gratitude should go to Dr Melissa Nyendak, Country Director, CDC/DGHT-India for the strategic leadership and guidance in providing Technical Assistance to NSACP, Ministry of Health, Nutrition & Indigenous Medicine, Govt. of Sri Lanka and CDC team for their support and guidance in these technical assistance initiatives.

Appreciate Dr Ariyaratne Manathunge, Consultant-Venereologist & Coordinator-SIMU, NSACP for his strategic leadership in coordinating the technical cooperation initiatives on TA to NSACP on SI with VHS-CDC Project, CDC team and contributions and support extended in bringing out this Resource Book and meaningful & successful conduct of the Training on Data Management.

Dr Rasanjalee Hettiarachchi Director National STD/AIDS Control Programme (NSACP), Sri Lanka

### Acknowledgement



Dr Joseph D Williams Director Projects The Voluntary Health Services (VHS) Chennai/INDIA

The Voluntary Health Services (VHS-CDC Project) with the support of Centers for Disease Control and Prevention (CDC/DGHT-India) in partnership with National STD/AIDS Control Programme (NSACP), Ministry of Health, Nutrition & Indigenous Medicine, Govt. of Sri Lanka is providing TA to NSACP on Strategic Information through a technical partnership initiative on:

- Enhance SIM Unit capacity to utilize electronic & manual program data for decision making;
- Improve capacity of SIM Unit to carryout management, analysis, documentation, and dissemination of summary program data reports;
- Improve capacity of SIM Unit to conduct and disseminate results of operational research;
- Consultation with stakeholders on monitoring & documentation of accomplishments & sustainability plans.

As part of this technical cooperation initiatives, VHS-CDC Project is providing capacity building initiatives, system strengthening, documentation and dissemination. In accordance with the capacity building initiatives, the project is organizing a series of training programs. VHS-CDC Project with the support of CDC/DGHT-India and in partnership with NSACP is conducting *"Training on Data Management and Analysis of HIV/AIDS Data"* for SIMU team, Consultant-Venereologists & Medical Officers, Public Health Inspectors, Development Officers and Nursing Officers.

To support this training, VHS-CDC Project is bringing out a Resource Book on Data Management and Analysis of HIV/AIDS Data. This resource book has been developed for the purpose of:

- Using in the training program for participants;
- Use as a reference material during the training and as and when required;
- Use this reference material for conducting/ capacity building of their team members in the respective Peripheral STD clinics.

This Resource Book has been developed with presentations, exercises, tools and other related reading materials by international professional trainers.

We thank Dr Rasanjalee Hettiarachchi, Director-NSACP for her leadership, supportive guidance in technical cooperation initiative.

We wish to acknowledge & thank Dr Ariyaratne Manathunge, Consultant - Venereologist, NSACP for his continuous support, strategic guidance and cooperation being extended in execution of this technical cooperation initiatives. Appreciate his strenuous support in bringing out this resource book and for conducting the training. Acknowledge the support extended by SIMU team, senior consultants in NSACP, SI team in peripheral STD clinics and key stakeholders.

We sincerely thank & acknowledge the technical guidance & support being extended by Dr Melissa Nyendak, Country Director, CDC/DGHT-India, Mr Lokesh Upadhyaya, Associate Director for Management & Operations, CDC/DGHT-India and CDC team. Wish to thank Ms Srilatha Sivalenka, Public Health Specialist, CDC/DGHT-India for her support in this cooperation initiatives.

We would like to thank Dr Yujwal Raj, Technical Advisor (SI), VHS-CDC Project for his support and contribution in developing resource materials.

We would like to thank Dr T Ilanchezhian, Senior Technical Advisor for his initiative, systematic support, strategic planning and contributions in bringing out this Resource Book.

We thank Ms T Sudha, Senior Programme Associate, VHS-CDC Project for her support in compiling, designing and bringing out this Resource Book.

We thank Mr B Kamalakar, Finance Controller and Mr S Sathyaraju, Associate Manager Finance, VHS-CDC Project and admin team for their support in publishing this Resource Book.

We greatly appreciate the fullest cooperation extended by NSACP & SIMU in this technical cooperation initiatives and for conducting capacity building initiatives.

Dr Joseph D Williams Director Projects The Voluntary Health Services (VHS), Chennai/INDIA

### I. DATA MANAGEMENT AND ANALYSIS OF HIV/AIDS DATA – AN INTRODUCTION:

GOAL: To build the data skills of NSACP staff in order to enhance the data quality, improve the data analysis and strengthen the use of HIV/AIDS data for epidemiological & programmatic decision making under NSACP.

#### **OBJECTIVES:**

- To build the understanding of the NSACP staff on the programmatic & epidemiological databases under NSACP;
- To introduce the basic principles and approaches of data management;
- To orient participants on methods of data quality assessment, validation & adjustments;
- To build the basic skills in statistical data analysis of program and epidemiologic data;
- To briefly introduce various software packages used for statistical analysis; and
- To improve the presentation, dissemination and use of data for programmatic purposes.

#### **METHODOLOGIES:**

- Active Learning through discussions & review of examples & case studies
- Learning by Doing
- Individual Hands-on/ Practical Exercises
- Group Exercises
- Parallel work on selected data

#### OUTCOMES:

- Identified important questions/ topics of programmatic relevance suitable for secondary data analysis.
- Exposed participants to basic principles and methods of data management.
- Enhanced knowledge and skills on analyzing the data and use of data under NSACP through hands-on practice on examples and actual program data.
- Improved skills on effective use of data to make evidence-based decision making under the program.
- Evolved a data analysis plan as a follow-up to the training and identified the next steps.

# II. PRESENTATIONS AND EXERCISE FORMATS



# Planning Program Data Analysis (1)

- 1. Identify programmatic questions/ issues for decision making in your programme
- 2. Map the data sources/ datasets (Review these questions & identify, which data can help answer these questions)
- 3. Assess the datasets for key aspects before analysis
- 4. Undertake data quality assessments, adjust & validate the data wherever necessary (Data Cleaning)
- 5. Decide the key outcomes of the analysis
- 6. Decide the type of analysis, methods & tools
- 7. Plan resources for analysis HR, Time, Money

### Planning Program Data Analysis (2)

- 8. Decide the target audience
- Develop the presentation, dissemination & publication plan (Way to present & disseminate the results of analysis, for programmatic improvements/ changes in implementation)
- Document the impact of your analysis (details of the means & outcomes of bringing in programmatic changes due to the results of your analysis)





#### **PROGRAMMATIC QUESTIONS (1)**

#### Epidemic questions

- Why are HIV cases increasing at this clinic/ in this area?
- ▶ What is the profile of the new HIV cases detected at the clinic?
- What are the key drivers of HIV epidemic in this region?
- Which areas/ regions have higher/ rising HIV epidemic?

#### Progress & Priority questions

- ▶ How is the scale up of STD services in this province?
- Why STD/ HIV service uptake is low at a clinic?
- What has been the progress against NSP/ Global targets in Sri Lanka?
- ▶ Has the targets for prevention been achieved? Gaps? Reasons?
- Which provinces should be prioritised for EMTCT?
- How to improve ART adherence rates?

### **Programmatic Questions (2)**

#### Performance questions

- Which STD clinics are the best performing in the country?
- Why service delivery performance is poor in a province?
- What are the top ten service delivery facilities that need focus under NSACP?
- How can the performance parameters be improved at ART centres?

#### Questions on evidence & information gaps

- How to generate evidence on role of international migration in HIV epidemic?
- What are the key information gaps that SI should focus in next two years?
- Which provinces suffer from lack of adequate good quality data?
- What are the barriers to strengthen data use for decision making under the programme?

#### **Programmatic Questions?**

Let us list out & discuss the programmatic questions/ Issues for decision making, that you are concerned about, in your area of work under NSACP!!!

(Use board/ display charts for the group work)

#### Fit for Secondary Analysis?

- Now, let us review these questions and identify, which data can help answer these questions?
- And is this data already collected elsewhere (Secondary Data)?
- Or does it need to be collected fresh from the field (Primary data)?

#### **Rich Evidence Base of NSACP**

- Expanded HIV Sentinel Surveillance System
- Rich Programme Data from STD/HIV Clinics & outreach
- Data from ANC clinics & PPTCT Centres
- Data from Blood Banks, TB Clinics, Hospitals & Labs
- Key Population Mapping & Size Estimations
- Integrated Biological & Behavioural Surveillance
- Intervention Data from Partner Agencies
- HIV Modeling, Estimations & Projections
- Increasing no. of Research Studies in HIV/AIDS
- DHS, Census Projections, Migration Statistics

### List out Data Sources!!!

- Let us list out all the known available data sources on HIV/AIDS
- Consider all programme data sources, surveys, other research data, other health data that may directly or indirectly give related info
- You can also list out any unpublished datasets
- For each dataset, list out the key info areas
- Use Exercise 1 Format Part A

### Map Prog Qs with Datasets

- Use Exercise 1 Format Part B
- List the programmatic questions on the left side
- If the data needs primary field data collection, mention 'primary'
- If it can be answered from secondary data, mention the s.no. of the dataset from Part A

#### **EXERCISE 1: IDENTIFYING PROGRAMMATIC QUESTIONS & MAPPING DATA SOURCES**

#### Part-A: List all the Known HIV/AIDS Data Sources & Identify the Key Information Areas

S. No.	Known HIV/AIDS Data Sources/ Datasets	Key information areas/ Broad themes

S. No.	Programmatic Question/ Issue for Decision Making	Probable Data Sources/ Datasets for Analysis (Mention S.Nos. from Part A Table)

#### Exercise 1: Part-B: Mapping Data Sources/ Datasets with Programmatic Questions



BY HEALTH SEA



NATIONAL STD/AIDS CONTROL PROGRAMME



- Data Info Knowledge
- Basics of Data Management
- Data lifecycle
- Database structure
- Key principles

# **Discussion???**

- What does good data management means to you?
- What are good things in data management system at your STD/HIV clinic?
- What are the challenges you face in data management at your STD/HIV clinic?
- What are the things that may be improved, with respect to data management at your STD/HIV clinic?

#### DATA - Information - Knowledge

DATA:	Facts concerning people, objects, events or other entities. Databases store data.

INFORMATION: Data presented in a form suitable for interpretation.

Data is converted into information by programs and queries. Data may be stored in files or in databases. Neither one stores information.

KNOWLEDGE: Insights into appropriate actions based on interpreted data.

#### DATA MANAGEMENT

- Good data management practices ensure that data are of high quality (reliable, consistent, and complete) as well as readily available to stakeholders
- Data management entails putting personnel, policies, procedures, and organizational structures in place to ensure that data are accurate, consistent, secure, and available
- Encompasses all components of data flow from the data collection tools used during service delivery to the databases used for data storage/accessibility along with all intermediate steps

#### Purposes of Data Management Systems

- Monitoring and evaluation of control programs
- Plan actions, programs, and resources
- To prioritize the allocation of health resources
- To provide the basis for epidemiological research
- Accountability

#### Purposes of Data Management Systems

- Beneficiary Management Level
  - Provision of better healthcare
- Health-Facility Management Level
  - Health facility functions
- System Management Level
  - Planning; M&E
  - Governance
  - Public health interventions

### DATA MANAGEMENT SYSTEMS

- Data Flow: the process by which data are transferred or "moved" from the primary data sources to the database(s)
- Database: the structured storage of data that will be accessed for data use
- Data use: the active use (visualization, analyses, interpretation) of data

# DATA Lifecycle

#### Creating data:

- Design research or program data collection/ acquisition
- Plan managementformats, storage etc
- Plan consent / sharing
- Identify existing data sources
- Collect data experiment, observe, measure



## DATA Lifecycle

#### Processing data:

- Enter data, digitize, transcribe, translate
- Check, validate, clean data
- Anonymise data as required
- Describe data
- Manage and store data



# DATA Lifecycle

#### Analyzing data:

- Conduct analysis
- Interpret data
- Produce various outputs
- Develop reports and publications
- Prepare for preservation



# DATA Lifecycle

#### Preserving data:

- Migrate data to appropriate format
- Migrate data to suitable medium
- Back up and storage
- Create metadata and documentation
- Archive data



# **DATA Lifecycle**

Giving access to data:

- Distribute data
- Share data
- Control access
- Establish copy right
- Promote data



# DATA Lifecycle

Re-using data:

- Follow up research
- New activities or research
- Undertake reviews
- Scrutinize findings
- Teach and learn



# Principles of Good Data Management Plan

- 1. Data / indicator definition and audit trail
- 2. Efficient and timely data flow
- 3. Effective data verification and cleaning procedures
- 4. Usable data storage
- 5. Data access policy/procedures

## Метадата

#### "Data about data"

- Description of fields
- Display and format instructions
- Structure of files and tables
- Security and access rules
- Triggers and operational rules
- Data collection details

#### Coding & DATA Dictionary

- The person who designs the coding for data entry is not necessarily the one who does data analysis.
- Hence we need a dictionary of all codes for all variables and questions, explaining what code means what, and relates to what question, what variable...
- A data dictionary is a descriptive list of names, definitions, and attributes of data elements collected in an information system or database
- Why Data dictionary
  - Standardization of terms used in the data base for common understanding of different users
  - Enhances inter operability across systems

#### Efficient and Timely Data Flow

- Document and monitor the flow of data from point of data collection to data storage including data entry and storage, mechanisms of data transfer etc.
  - how are data first collected (e.g. paper forms, excel sheets by staff or via the Ministry of Health, directly entered into a database)?
  - How does data typically flow from the site to the district?
  - How does data flow from the district to the region?
  - Who typically are responsible for transferring data between steps?
  - Document the steps till data storage (excel, Access database, shared drive, etc.) and staff responsible for transferring data

### **DATA PREPARATION**

- Data Review
- Data Quality Assessment
- Coding & Data Dictionary
- Data Cleaning: Adjustments & Imputations

### **DATA ANALYSIS**

- Computation & Recoding
  - Retain original data
  - Give appropriate variable names
- Generation of tables
- Tests and Interpretation
- Summarising results, writing & publications

#### DATA STORAGE & RETRIEVAL

- The data storage system must be designed to meet the requirements and data use needs and allow for simple and reliable access to the data
- Hard & soft forms
- Migration to suitable form & medium
- Back-up & storage
- ▶ Tables, Identifiers & Query Options
- Central database or databases linked in such a way that data can be easily and automatically combined as required

### **DATA PROTECTION & SHARING**

- Infrastructure security measures and user access permissions
  - Physical safety Lock & key
  - Restricted access Password protect all digital files/ databases; access control based and privileges based on roles
- Confidentiality (Anonymising) agreements for collectors, entry operators and supervisors
- Sharing protocol for authorized access
- Data recovery plans where there is loss of data

## Maximizing Data Usefulness

- Data Use Plan can guide the design of the data management system
- It is essential to take the time to collect our data properly the first time to avoid re-collecting data
- Review workflows to ensure that we avoid re-processing data as much as possible increases our efficiency
- Electronic data systems and automated reporting drastically reduces data manipulation

#### LONG-TERM PLANNING

- What will happen to my data after my project ends?
- How can I appraise the value of my data?
- What are my options for archiving and preserving my data?
- What are my options for publishing and sharing data?

### DATABASE & DBMS

DATABASE: A shared collection of interrelated data designed to meet the varied information needs of an organization.

DATABASE MANAGEMENT SYSTEM: A collection of programs to create and maintain a database.

Define - Construct - Manipulate

#### Advantages of Database Processing

- More information from same data Shared data
- Balancing conflicts among users Controlled redundancy Consistency

- Integrity
  Security
  Increased productivity

### **CREATING DATABASE STRUCTURE**

- Cases or Records or Rows
- Fields or Variables or Columns
- Primary key/ Unique ID/ Case ID
- Column headers
- Data Labels
- Value Labels
- Type of Data
- Data limits/ checks
- Metadata

#### Some Key Principles

- One row for one case
- No duplicate variables/ column heads
- No duplicate primary key
- No sub-totals in rows; sub-totals/ totals in columns are OK
- No merged cells; No merged headings
- No blank cells (Fill blank cells with some code or Impute)
- No two data types in one column (Text/Num/Code)
- Short & crisp variable names; Not too long
- Colour code in Excel
- Numeric data is better than text

### **Compiled DATASETS**

If dataset is a compilation of multiple datasets on the same cases:

- Ensure column alignment
- After joining the datasets, create a unique ID/ primary key
- Unique ID can be created by joining two three variables.
  E.g. Province, Facility Name & Month
- Avoid different spellings of same case
- Ensure uniformity in the data type & format under each column

## **Exercise 2: Review Database**

- Review the database shared with you
- Identify the following from the 'Exercise 2 Dataset' given to you
  - ▶ What is a case in this database?
  - ▶ No. of cases
  - ▶ No. of fields and their names
  - Primary key
  - Metadata is adequate and clear
  - > Any of the key principles violated in the dataset?

<b>EXERCISE I</b>	DATASET
-------------------	---------

				2003		2004		2005		2006		2007		2008	
State	District_name	Site_name	Site_Type	NT	NP	NT	NP	NT	NP	NT	NP	NT	NP	NT	NP
Andhra Pradesh	Visakhapatnam	Priyadarshini Service Organization, Vishakhapatnam	FSW	250	32	250	35	250	39	250	27	250	22	248	12
Andhra Pradesh	East Godavari	East Godawari	FSW	250	113	250	102	250	67	250	23	250	39	244	45
Andhra Pradesh	Prakasam	Lakshmi Development Society, Ongle, Prakasam	FSW	250	61	250	27	250	18	250	11	250	18	235	16
Andhra Pradesh	Hyderabad	Hyderabad	FSW	250	40	250	25	250	30	250	24	249	18	250	37
Andhra Pradesh	Kurnool	Parameswari, Kurnool	FSW	250	22	250	25	250	8	249	6	249	9	249	12
Andhra Pradesh	Warangal	Warangal	FSW	250	32	250	47	250	32	249	22				
Andhra Pradesh	Guntur	Needs Society, Chilakaluripet, Guntur	FSW			250	36	250	33	250	15	250	32	213	19
Andhra Pradesh	West Godavari	Action for Development, Bhimavaram (New 07)	FSW									249	41	247	35
Andhra Pradesh	Khammam	JAGRUTI (New 07)	FSW									250	39	247	66
Andhra Pradesh	Adilabad	AIRTDS,Mancherial (New 07)	FSW									217	10	249	38
Andhra Pradesh	Nalgonda	ANKITA (New 07)	FSW									250	25	227	21
Andhra Pradesh	Srikakulam	Swageti Project,Youth Club of Bejjipuram (New 07)	FSW									250	15	243	10
Andhra Pradesh	Warangal	MARI, Hnamkonda (New 07)	FSW									120	8	248	12
Andhra Pradesh	Khammam	JAGRUTI (New 07)	FSW			248	12								
Andhra Pradesh	Nalgonda	ANKITA (New 07)	FSW					250	11						

FSW: Female Sex Workers; NT: Number Tested; NP: Number Positive
#### **REVIEWING THE DATABASE**

S. No.	Observe/ Identify	Observations/ Details
1	Is it a single dataset or compiled dataset? Why do you say so?	
2	What is a case in this database?	
3	No. of cases	
4	No. of fields and their names	
5	Primary key	
6	Metadata is adequate and clear	
7	Any of the key principles violated in the dataset?	
А	One row for one case	
В	No duplicate variables/ column heads	
С	No duplicate primary key	
D	No sub-totals in rows; sub-totals/ totals in columns are OK	
E	No merged cells; No merged headings	
F	No blank cells (Fill blank cells with some code or Impute)	
G	No two data types in one column (Text/Num/Code)	
н	Short & crisp variable names; Not too long	

Review the 'Exercise 2 Dataset' given to you and fill the following table.

Data Management Systems - Assessment Tool	National	State	District	Action to be taken
Standard operating procedures have been written that define roles and responsibilities for				
data compilation, reporting, data analysis, dissemination and use.				
There is a comprehensive, singular, master list of health facilities, with unique facility				
identifier and service domain, that includes the private sector and special facilities (military,				
etc.).				
There is a formal mechanism to update and keep current the master facility list (e.g., a				
census of all facilities is conducted every 5 years).				
Data collection systems for client data (e.g. clinical episodes) are standardized across all				
implementing partners and donors.				
Personnel (clinicians and other staff) have been trained in the collection of the client data,				
and for the input of the data into the computer database (where applicable).				
Printed guidelines are available at all health facilities (and in applicable community-based				
programs) to assist with client-level data collection.				
Health data (paper or electronic) are stored appropriately and according to national policies.				
There is a schedule/plan for update, reproduction, and distribution of data collection tools.				
The data flow pattern (i.e., data flow from client encounter forms -> summary tools [e.g., a				
register or tally sheet] -> periodic aggregate reporting form) is clearly defined and				
understood by staff.				
There are printed guidelines available at all health facilities (and in applicable community-				
based programs) to assist with data compilation and reporting.				
Relevant staff at health facilities (and in applicable community-based programs) have				
received training on data compilation and reporting.				
Data disaggregations by key stratifiers (age, sex, geography) are maintained during				
compilation and transfer in order to permit equity analysis.				
Data transfer to the next level occurs in a timely way, making use of innovation and IT where				
appropriate and available.				
There is a data quality assurance plan that is shared with health programs, other government				
ministries, donors and other stakeholders to guide activities aimed at improving data quality.				
Routine health data quality assurance standards are defined and enforced, including				
completeness, timeliness, accuracy, integrity, and consistency over time.				
Roles and responsibilities for data quality are assigned at each level, including verification of				
data, summarizing data quality issues, and developing and implementing improvements				
strategies.				

Data Management Systems - Assessment Tool	National	State	District	Action to be taken
Training and capacity development for data quality assurance are provided at facility, district,	,			
and national levels using standard methods.				
Systematic and comprehensive assessments of facility data quality are conducted regularly in				
advance of health sector planning, including analysis of completeness, timeliness, accuracy,				
and consistency over time (e.g., data quality review) and which result in published reports				
describing data quality issues and plans to address them.				
Data management staff conducts regular checks of accuracy and completeness of data prior				
to submitting reports to the next level (using automated electronic checks where				
appropriate).				
Data quality assurance is linked to the health sector planning cycle in the country so that				
information on data quality is available prior to the use of data for planning.				
There is collaboration between the MOH, government agencies (e.g., national statistics				
office) and other national stakeholders (e.g., donors, universities, etc.) on data quality				
assurance so that assessments are conducted with an element of independence (i.e., no				
conflict of interest).				

#### Data Management Checklist

Databases	1	2	3	4
Name of database				
1. Type of data				
What kinds of data - survey, interview, observation,				
machine or instrument collected, physical samples,				
models, etc are you collecting?				
What formats - paper, digital, image, audio, other - will the				
 data be in?				
Will it be reproducible? What would happen if it got lost or				
 became unusable later?				
 2. Data formats and standards				
Do you have data dictionaries, code books or other				
documentation to explain terms, variable names,				
 codes and abbreviations used?				
Have you provided documentation describing how the				
 data were collected or created?				
Have you used standard collection methods, standard				
data formats, and standard file format choices (if				
these exist for your research field)?				
 3. Data access policies				
Have you removed personal or sensitive information				
from your data to ensure privacy protection?				
Have you established who owns the copyright of your				
 data?				
Do you have documentation on how institutional and				
personal credit should be acknowledged for your				
data?				
Are your data, records, and files labeled and logically				
organized?				
Have you used consistent and easy to understand file				
names?				

#### Data Management Checklist

Databases	1	2	3	4
Name of database				
4. Data use and distribution				
How will your data be made available?				
Do you plan to limit re-use or re-distribution of your				
data? If so, why and for how long?				
5. Data preservation and archiving				
Have you made arrangements for the long-term				
storage and preservation of your data (both physical				
and digital collection items)?				
Do you have data security plans in place to ensure				
that copies of your data are stored and backed up on a				
regular basis?				
Are you using data formats and software that enable				
sharing and ensure long-term validity of data, such as				
non-proprietary software and software based on open				
standards?				
When converting from one format to another, have				
you checked that no data are lost or changed in the				
process?				





# Main Data Systems of NSACP









	Case Reporting
1. HIV	
– H.:	L214 form (OLD)
	CONFIDENTIAL REQUEST FOR HIV ANTIBODY TEST/NOTIFICATION
	(To be related by the Physician PATIENT No. :
	ADDRISS 1 DATE OF REQUEST D M Y
	DATE OF REQUEST PATIENT No.1
	NPROBLED CONSENT FOR HIV TESTING GUILANDED ROM PATIENT 1, VES 2, NO NAME OF PATIENT ORIST TWO LITTERS : OF GIVEN NAME AND SURVAME ONLY)
	DISTRUCT OF RESIDENCE:
	SEX I. MALE 2. FRMALE MARITAL STATUS I. NEVER MARRIED 2. CURRENTLY MARRIED-LIVING TOGETHER 3. SEPARATEDIDIVORCED/WHO/WED
	OCCUPATION Penal Specify REASON FOR TESTING : 1. INTERSY WITH SYSAPTOMS (Confirmation on the LA content
	2. ASYMPTOMATIC: 3. VISA SCEEDING: 4. ORGAN DONOR 5. SURVEY: 6. STD CLINIC ATTENDE: 7. OTHER (Spot)

Case Reporting								
. HIV								
– H.1214 form (	modified in 2017)							
Request for Confirmatory HI of the National S	IV Testing from the Reference Laboratory							
of the Hallohar e	(VERSION: JAN 1, 2017)							
Instructions: To be completed by referring doctor/h worker at the time of requesting HIV confirmatory test reference laboratory of the National STD/AIDS Control No. 29, De Saram Place, Colombo 10, Sri Lanka.	ealthcare t from the I Programme, Date of Receipt Day Month Year							
Patient should be informed that all questions contained in this que strictly confidential and will become part of their medical record)	stionnaire are							
PART II - TESTING DETAILS AND DEMOGRAPHIC INFORMA								
PATIENT/CLIENT 1A. STD Clinic Registration Number IDENTIFICATION (For STD Clinic Clients)	1B. Sample Number (For non-STD Clinic Clients - Private Lab, TB clinic, Hospital IC							





#### Probable Mode of Transmission of New HIV Cases in 2018 (N=350)











Autional STD/AIDS Control Programme, Centrol Programme, Centrol Values Registration Number	ITEENT FORM - RECEISTRATION al STD Clinic, Colombo 19 Date of registration (differently)	STD Patient						
for samefaitale	Las same Place Place	Form						
trinanost allinsa	Please							
veferred mode of contact 1. Do not contact /contact deaths are charged during subsequent contact address:	n 2. Letter 3. Email 4. T please 5. Visit	COMPLETION OF EPISOBE OF CARE 4 to thomas 2 to 100 proteins 3 to 000 4 to thomas 2 to 100 proteins 4 to 000000 1 to 0000 proteins 4 to 0000000 1 to 000000 1 to 0000000 1 to 0000000 1 to 0000000 1 to 0000000 1 to 0000000 1 to 000000000 1 to 000000000 1 to 00000000000 1 to 00000000000000000000000000000000000						
neuer address: 	Peer Peer	1 No. 1 (201- watering 1 (201- wate						

# **STD Registers**

- 1. Main Register
- 2. Subsequent Visit Register
- 3. Outpatient Blood Testing Register
- 4. Interview and Contact Tracing Register
- 5. IEC/BCC/ Awareness Programme Register
- 6. HIV testing and counseling Register

Source: SIMU/NSACP/2014

# Cont., STD Registers

- 7. Condom Distribution Register
- 8. Commercial Sex Worker Register
- 9. Outreach Blood Survey Register
- **10. Defaulter Register**
- 11. Antenatal Syphilis Register
- **12.** Pre-employment/Visa Screening

Source: SIMU/NSACP/2014



# **Main Register**

- S.No. for episode of care
- Date DD/MM/YYYY
- Name
- Master No.
- Sex (M/F)
- Date of Birth DD/MM/YYYY
- Age (Years)
  - Address Telephone Number Email
- Marital Status<sup>1</sup>
- Reason for Attendance<sup>2</sup>
- Diagnoses<sup>3</sup> (Please include all diagnoses of an episode in a single cell)

- Age < 15 years
- Age 15-24 years
- Age 25-49 years
- Age 50 or over
- NGO escorted
- Sex worker
- MSM
- DU
- IDU
- Beach boy
- Prisoner
- HIV tested
- Received HIV results

### **QUARTERLY STD RETURN**

# **Quarterly STD Return**

QUARTERLY RETURN FROM STD CLINICS IN SRI LANKA (Revision: 29.05.2017)

Date of completion : \_\_/ \_\_/ 20\_\_\_

Table 1. Total number of new diagnoses\* by age group and sex (Source: Main Register- 2017)

				Male			Female					
SN	Name of the Disease	<15 year	15-24 year	25-49 year	50+ year	Total	<15 year	15-24 year	25-49 year	50+ year	TOTAL	
1.1	Infectious syphilis											
1.2	Late syphilis											
1.3	Early Congenital Syphilis											
1.4	Late Congenital Syphilis											
1.5	Gonorrhoea and presumptive GC											
1.6	Opthalmia neonatorum											
1.7	NGU/NGC											
1.8	Chlamydia											
1.9	Genital herpes											
1.10	Genital warts											
1.11	Pelvic inflammatory disease (PID)											
1.12	Trichomoniasis											
1.13	Candidiasis											
1.14	Bacterial vaginosis											
1.15	Other STIs											
1.16	TOTAL STI											
1.17	Non STI/Uncertain											
1.18	No illness											
1.19	GRAND TOTAL											

#### **Excel STD Database**

New ratients he	gistered of a	ii 510 Clinic	11 2017	New Fatients Reg	istered of al	i si b cimic il	12010
	New patient	s registered			New patien	ts registered	
-	Male 🝸	Female	2017 -	*	Male 🝸	Female	2018 -
Ampara	138	188	326	Ampara	135	139	274
Anuradhapura	376	330	706	Anuradhapura	397	349	746
Avissawella	26	18	44	Avissawella	207	168	375
Badulla	320	410	730	Badulla	269	401	670
Balapitiya	112	109	221	Balapitiya	131	111	242
Batticaloa	64	138	202	Batticaloa	57	131	188
Chilaw	359	496	855	Chilaw	628	557	1185
Colombo	3470	1918	5388	Colombo	4394	2229	6623
Embilipitiya	23	37	60	Embilipitiya	104	103	207
Gampaha	294	295	589	Gampaha	399	416	815
Hambanthota	364	272	636	Hambanthota	411	321	732
Jaffna	160	121	281	Jaffna	160	85	245
Kalmunai	115	93	208	Kalmunai	122	123	245
Kalubowila	814	645	1459	Kalubowila	1044	716	1760
Kalutara	560	466	1026	Kalutara	424	434	858
Kandy	515	533	1048	Kandy	636	785	1421
Kegalle	243	312	555	Kegalle	382	349	731









			11111	-failing theory	i		and burgers	5.	Clinical 7	and Labora	atory Inve	stigations		parenticore
🦷 1.Pa	tient Identification D	Jata Write co	naiete	TULOFELI RITOLI			BORDOBCION COMPANY	Date	WHO	Weight	Height	Perfor-	Total	CD4 co
Registration Number :	עונונים ביו	code	alinic (28	A)-code balle	m (4#)		1 1	(distimm	stage	(%0)	(cm)	A/B/C*	count	dition
Name of Treatment Unit		City:				- 1	At det utett to elimic							1
District:		Statesprown	.061			-	At ACT medical relability				ditt			
Name of patient:		. רורש		Mala 🗖	Eemala		At start of ART				shild			
Age: (di	de of bith: LUM_UL cid / mm	JUL S	×Ц	walla 🗀	remaie		At 6 months ART				dalid			
Patient's phone number						1	At 12 months ART				0.03			
Address:							At 24 months ART	1			et al l			
City/village:	District:	Sta	alprovin	108:		. 1		10000	1 10 1	inneticitenti	ial Treater	ent	170	110000
Distance from residence	to clinic (km)						19	Part of Control of Con	UDETITU	TION within	1" line St	WITCH to 2	" Ine. STOP, P	EBTART
Treatment supporter's r	ame (if applicable)						Treatment Starton	SALC: N	S	ubstitution.	Reason	Data and	Net Net	a renieren
Treatment supporter's a	ddress:						D4T30+3TC+NVP	Dat	8 5W	itch or slop	(code)	Dane rest	an nev	Viegnitori
Treatmant supporter's p	hane number:						D4T40+3TC+NVP							
Date confirmed HIV+1	est []]]/[]]/[]	<u> </u>	ace:											
	dd / mm / t	99 10.1 anna 11 11 11 11 11 11	et in 2	TR 113-0	Instant	1	CI TRUGSTONNIP							
Entry point (services o	dering the patient for H	in cardy (1) in		100 100 00 00	Soil andor	the second se	CI XELOCATO HIGH				1			
and A locations, \$10 \$ Day	diatria ID R-PMICT []	7-511 (1 B-P/W	1410 1 1 1 1	-NGO [] 10	-State server	100	TT 2TMANTCHEEV	5				And a second sec		
4-Inpatient	In a ART from writter 1	7-STE [] BAYW 13-other 13-other HV careWRT d Date transfe	inic from med in :	t the national	program		Reasons SUBSTITUT new drug available, 6 Reasons for SWITCH	TE: 1 toxi drug out i 1: 1 clinic	sity side et rf stock, 7 il treatmor	flects, 2 preg other reason st failure, 2 h	nancy, 3 ris n (apecily) nmunologic	ik of pregna	sncy, 4 newly dia virologic failure	sgnozed TI
4-Inpatient      5-Pae     114DU outreach     patient transferred in     Name previous clinic:     2. Personal Histor:     Mode	elatinc 6-PMTCT 12-CSW currentle 1 12-CSW currentle 1 1 on ART from another 1 / {Tick one choice}	AV careWRT of Date transk Montal status	nic from med in nity His	the national story (Title	program one cho Estimate	ice) id monthly id income:	Reasons SUBSTITUT new drug available, 6 Reasons for SWITCH Reasons STOP: 1 to hospitalization, 6 drug interruption, 10 differe	TE: 1 taxis drug out i 4: 1 clinics xicity side s out of str	sity side et af stock. 7 al treatmon effects, 2 ick, 7 patie	Rects, 2 prog other reason nt failure, 2 in pregnancy, ant lack of fir	nancy, 3 ria (apecily) mmunologic 3 treatment terce, 8 pai	ak of pregna ;al faiture, 3 faiture, 4 pr bern decisio	sncy, 4 newly dia virologic failure opr adherence, 0 n, 9 planned tre	sgnozed Ti 5 illnese atment
4-Inpatient      5-Pas     11-LDU outreach     patient transferred in     Name previous clinic:     2. Personal Histor     Mode     1 Comman     of HV     2 Clinic her	eliatric   6-PMTCT   12-CSW cutreach   14 n ART from another t  y (Tick one choice) at sex worker (CSW) erosesual route	7-511 [] B-44W 13-ether AtV caneWRT d Date transfe 3. Fai Montal status [] Married	nic from med in : nity Hit L) Sin ] Divore	the national story (Tick plo colseparate pericebia	program one cho Estimate househo	icae) ed monthly rid income:	Reasons SUBSTITUT new drug available, 6 Reasons for SWITCH Reasons STOP: 1 to hespitalization, 6 drug interruption, 10 others	TE: 1 taxis drug out i 3: 1 clinic xictly side ; out of str ; 7.	sity side et af stock, 7 al treatman effects, 2 vck, 7 pati Tuberce	Hects, 2 prog other reason of failure, 2 in pregnancy, ent lack of fis allosis treat	mancy, 3 ris n (apecily) mmutologic 3 treatment terren, 8 pai	ak of pregna al failure, 3 failure, 4 pr tient decision ing HIV ca	ancy, 4 newly dia virologic failure opr adherence, 5 m, 9 planned tre ine	sgnozed 11 5 illnees satment
4-Inpatient      5-Pas     11-LDU outreach     pasient transferred in     Name previous clinic:     2. Personal Histor     Mode     1 Comman     of HV     2 Clinar he     Imma     3 Mart havi     missio	eliatric 0 e-PMTCT 0 12-CSW cutreach 0 i en ART from another t y (Trick one choice) ist ses worker (CSW) erosexual racke ig acc with men (MSM) ing use (DU)	Arrow Construction	inic from med in : LJ Sin ] Divor: Not a x   Aa	story (Tick gle pplicable y HV	program one cho Estimate househo	ice) ed monlity vid income:	Reasons SUBSTITU new drug available, 6 Reasons for SWITCH Reasons STOP: 1 to hospitalization, 6 drug interruption, 10 others Disease class (lick)	TE: 1 taxis drug out of H: 1 clinic xicity side s out of str 5 7. TB 1	sity side et at stock, 7 al treatmon effects, 2 ack, 7 path Tuberca togimen (	Hects, 2 prog other reason of failure, 2 to pregnancy, ant lack of fis alloais treat lick) T	mancy, 3 ris n (specify) mmunologic 3 treatment Ismon, 8 pai tment dur B registrati	ak of pregna cal faiture, 3 faiture, 4 pr tierri decisio ing HIV ci ion	snoy, 4 newly dia i virologic failure opr adhesence, 6 yn, 9 planned tre kNe	sgnozed T 6 Illnees istment
□ 4-logation: □ 5-Pias □ 11-IDU outreach □ patient transferred in Name previous clinic: 2. Personal Histor Mode □ 1 Common of HW □ 2 clinic the minissio □ Arel new nicksio □ 2 feeding n □ 2 fielded for	eliatine [] e-PMTCT [] 12. CSW cutreach [] 13. CSW cutreach [] 4. (Trick one choice) ial sex worker (CSW) trajecutories (CSW) tra	7-STIL D BANK 13-other HV caneWART 0. Date transk Marital status Marital status Marital status Marital status Marital status Harrid Widowed Fanky menta Fanky menta	nic from med in : Divorc Not a se Age	story (Tick gle pplicable y HV t *Aktion	cree cho Estimate househo ART YN	ice) ed mostRMy pd income Pd income	ZOVASTCHEV     Reasons SUBSTITU     new drug available, 6     Reasons STOP: 1 the     resultation, 6 drug     interruption, 10 othern     Disezase class (lick)     Puinterany TB	TE: 1 basis drug out i H: 1 clinic xictly side s out of str 5 7. TB 1	pity side ei af stock. 7 al treatman effacts, 2 ack, 7 path Tuberca togimen ( lategory l	Hects, 2 prog other reason at failure, 2 in pregnancy, ent lack of fir allosis treat lick) T D	nancy, 3 ris n (specify) mmunologic 3 treatment smos, 8 pai tment dur B registrati istrict:	ak of pregna cal failure, 3 failure, 4 pr tient declais ing HIV ca ion	ancy, 4 newly dia indrologic failure oer adhesence, f an, 9 planned tre ere	agnosed Ti 5 illness atmant
Inspalent II S-Pas     Inspalent II S-Pas     Instant transferred in     pasient transferred in     Name previous disk:     Z. Personal Histor     Mode UI Common     of HV UI 2 cauch teat     in UI Sten how     Wissib DI A Inspalent	eliatric [] e-PMTCT [] 12. CSW cutreach [] 13. CSW cutreach [] 14. an ART from another to (Tick once choice) ist see worker (CSW) representations repres	7-STIL BANK 13-other Dote transk HV careWRT o Dote transk Warted Warted Widowed Family member partnerkhilder	Inic from med in : Divors	story (Tick gie colseparate ppicate s +Narison	crie cho Estimate househo ART YIN	ice) ed mostRMy sid income Places	ZOVASTCHEV     Reasons SUBSTITUT, mer drug available, 6     Reasons for SWITCH     Reasons STOP: 1 to     hespitalization, 6 drug     interruption, 10 officer      Disease class (tcs)     Constant TB     Constant TB     Smear-pacifie	TE: 1 toxi drug cut i H: 1 clinics wichy side s out of strip 7. TE 1 C	sity side et of stock. 7 al treatmen effects, 2 sck. 7 pable Tuberc4 togimen ( lategory II alegory II	Hects, 2 prog other reason at failure, 2 in pregnancy, ent lack of fai programs, 1 foolis treat fick) T D H	nancy. 3 rk n (specify) mmunologic 3 treatment ternent dur B registrati istrict: isath Centro	ak of pregna cal failure, 3 failure, 4 pr tient declaid ing HIV ca ion	ancy, 4 menty dis vérologic feiture cer admanance, 6 an, 9 planned tre ane	agnozod Ti 6 illinosis istmant
Lingatient II S-Paia     Int-IDU outreeth     paiaent transferred in     name previous disk:     Personal Histor     Mode         1 Common         Anne previous         Anne         S Destination         Anne         Annene         Annene         Annene         Annenene         Annenen	eliatric 0 +PMTCT 0 12 - CSW curreach 0 14 - en ART from another 1 4 - e	7-STIL BARW 13-other Dote transfer Dote transfer Dote transfer Marital status Marital status Marital status Marital Status Marital Status Marital Status Marital Status Marital Status Marital Status Partneridilititi Marital Status Farange Status Partneridilititi Status Partneridi Status Part	Inic from med in : Divorc Not a Not a	story (Tick gie colseparate ppicate s *Anteore	crie cho Estimate housetx	ice) ed mostitity yid income Pregist. No Pris care	ZOVASTCHEV     Reasons SUBSTITU     mev drug available, 6     Reasons for SWITCH     Reasons STOP: 1 to     keepialization, 6 drug     interruption, 10 others     Disease class (lick)     Disease class (lick)     Sincer-regalable     Sincer-regalable	TE: 1 toxic chug cut i H: 1 clinic sicily side sout of str 7. TE 3 C 0 0 0 0	city side et af stock, 7 al treatman effects, 2 ack, 7 psib Tuberco Rogimen ( lategory i altegory i 2/her spec	Hects, 2 prog other reason in failure, 2 in pregnancy, ent lack of fe illosis treat fick) T D H H ify: T	Inancy, 3 thi n (specify) mmunologic 3 treatment arrow, 8 pail treatment dur B registrati listrict anth Control B numbor B numbor	sk of pregna cal failure, 3 failure, 4 pr tierr declaid ing HIV ci ion e: satcome: C	ancy, 4 menty dis vicelogic failure cer adtesence, 6 m, 9 planned tre are	agnozed Ti 5 illness istmant
□ - Ingatient IS-Pies □ 11-12U outreach □ □ passent transferred in Name previous clinic: 2. Personal Histor Mode 0.112 □ 1 Comme 0.112 □ 1 Comme 0.112 □ 1 Comme 1.122	eliatine	7-STI    BAAW 13-ether    U canaVAT c    Date transk    Date transk    Marital status    Marital status    Midwood Fanily mental partiacchilder	Inic from med in : Inity His Divort Not a set	story (Tick gie colseparate pplication * HV * Wateroom	crie cho Estimate houseto	ice) ed mostitily sid income Placano Placano	ZDV4-STCHEV Reasons SUBSTITU Reading avoitable, 6 Reasons for SMITCH Reasons STOP: 1 to Installation, 6 drug interruption, 10 diterr Disease class (lick) Puintorany TB Smear-pail/te Encore-regil/te Encore-regil/te Encore-regil/te Reasons Rea	TE: 1 toxis disc g out of H: 1 clinic wickly side g out of str 5 7. TE 1 C C 0 0 0 0 0 0	city side et af stock, 7 al treatman effects, 2 ack, 7 pil) Tubercu Regimen ( Sategory I alegory I 2 the spec start TB	flects, 2 prog other reason nt tailure, 2 in pregnancy, ent tack of fit prognancy, ent tack of fit globia treat (Sch) T D H H T T Re: C	Inancy, 3 rbi n (apocity) mmunologic 3 treatment norce, 8 pail timent dur B registrati istrict: isath Contor B number reatment o J Rix failure	sk of pregna saf failune, 3 failune, 4 pr tiern dechte ing HIV ci ion e: ustcome: [] bied	ancy, 4 newly dia vicologic failure our adheamce, 6 ano glanned tre ano i Care C Ricci Defail C 1	agnoted Ti 5 illness istmant completed fransfer or
Lingulant IS-Plas     Int XDV outreach     passent transferred is     Name previous ditis:     Personal Histor     If XDV or Conscrete     If YDV	sianco+MPETT 12. CSW collieadi an ART from another 1 / (Tick one choice) al servenkar (CSW) collections and servenkar (CSW) and servenkar	7-STI    BAAW 13-ether    Date transk    Date transk    Marital etatus    Midriwed    Widrwed Famly mental partnerkhiller	Inic from med in : Divot a Not a	story (Tick places and the second sec	crie cho Estimate househo	ice) ed monthly isd income Plat care	ZOVASTCHEY     Reasons SUBSTITUT new drug avolable, 6     Reasons STOP: 1 the     Reasons STOP: 1 the     result of the second stop of the     result of the second stop of the     result of the second stop of the     Second	TE: 1 toxic ding cut i H: 1 clinic social yields social of stress 7. TE 1 C C Date	city side el of stack, 7 al treatman effects, 2 sck, 7 psi/ Tuberca Rogimen ( Jategory II Jategory II Jategory II Jategory II Jategory II Jategory II	Hects, 2 pro; other reason of failure, 2 in pregnaincy, ent lack of fit illoais treat (Sck) T i ity, T Rc; C C	(nancy, 3 rbi n (apecily) mmanologic 3 treatment nance, 8 pal tment dur B registrati istrict isaht comin B numbor ] Rx failure iale:	sk of pregna saf failune, 3 failune, 4 pr tiern dechte ing HIV ci ion e: 	anoy, 4 newly dic visiologic failura opr adheam.cc, f an, 9 planned tre are i Ours  Rx cc Defoult  T	agnosed Ti 6 Illness atmant completad fransfer or
© 4-inpatient = 5-Para ■ 11-10U outreach = □ patient transmerred in Name previous divit: 2. Personnal Histor Mode u   common entry = 2 common entry = 2 common entry = 2 common = 4 inpacting = 1 common = 1 commo	Liant: 0+PeFCT      Lo CSW cettered:      Lo CSW cettered:      un ART from andhet 1     y (Trick one cholos)     si sex worker (CSW)     workersex.incom     run     for the set of the se	7-STI [] B44W 13-ether HV careWAT o Date frame ] Martial status Martial status Martial status Martial status Martial status Martial status Martial status Martial status Martial status Martial status Family mental partnerchillon	Inic from med in: Divor: Not a n	the national story (Tick gls colseptratio pricesta d' HV = Addresse	crie cho Estimate houseful	ice) ed monthly sid income: Begliet.No Fit care	ZDV+3TC+E-V Reasons SUBSTITU new drug avol3bie, 6 Reasons for SWTCF Reasons STOP: The SWTCF Records Records Record Record Records Record Records Record Reco	TE: 1 toxis drug cut i H: 1 clinics wichy side s out of str s 7. TB 1 C C Date	city side el of stock, 7 al treatmen effects, 2 sck, 7 più Tuberca Rogimen ( Lategory I Lategory I Driver spec i start TB d / orm J	Hects, 2 prog other reason in tailure, 2 in pregnancy, ent lack of fin aloais treat lick) T Bick) T Rec Sys	mancy. 3 ris n (specify) mmunologic 3 treatment asmos, 8 pail treatt dar B registrati istrict: isath Canin B number reatenent o J RX failure dd J	sk of pregna cal feiture, 3 faiture, 4 pr tierri declais ing HIV ci ion e: vatcome: Died y em 1 yy	iveroign of a newly dia visiologic failure our adheatinee, f iv, 9 planned tre are	sgnozed Ti S ill nees seinent completed Transfer ou
Q 4-inpatient = 5-Para 11-100 uniterential = □ pail ant branchered in Name previous disk: 2. Personal Historia Merica 1. Comment Merica 1. Comment Merica 2. Mont two missio 1. Statuto 3. Mont two missio 1. Statuto 3. Mont two Per IDIA Substituto Per IDIA Substituto P	Stain: 0+MPECT 1 - on ART from windher 1 y (Tick one sholoo) y (Tick one sholo	7-511 [D 641W] 13-ether UIV canaVART c Date transl Mantal ethias Named Na Na Named Named Named Named Na Na Name	Inic from aned in : Inity His Divors	story (Tick gis colseparate pricessa d' HV = Addresse	crie cho Estimate househo	ice) id martify jd incenté Pis care	204-3TC+E-V Reasons SUBSTITU new drug avalabie, 8 Reasons for SMITC+ Reasons of SMITC+ Reasons STGP-114 Readination, 6 drug interruption, 10 citem Declaration of smitches Declarations of the General control (the) Ge	TE: 1 toxic drug cut i H: 1 clinics wichy side s out of str s 7. TB 1 C 0 0 0 0 0 0 0	city side el of stock, 7 al treatmas effects, 2 ack, 7 psi/e Tubercs T	Hects, 2 prog other reason in tailure, 2 in pregnancy, ent lack of fit alloais treat (sch) T Re: C Sy 8, Enel of	Inancy, 3 ris (specify) mmanologik 3 treatment, agroe, 8 pal treatt dur B registrati idrict: B number restancer o ] Rx failure late: d ] Fx failure	sk of pregna cal failure, 3 failure, 4 p tierri decisio ion e: Diod y July 4 ran 7 yy (B:	ancy, 4 newly dit velologic failure oer adheance, f ar, 0 planned tre are ) Care    Rx co    Default    1	sgnozed T 5 Illnoze adment completed Transfer or
A - Ingatient:	Static [] -PPEFCT []     on ART from and/let []     on ART from and/let []     on ART from and/let []     (Trick core choice)     all set works? (CBN)     corescut includ     grades with men (MSM)     and your (MSM)     and your (MSM)     include the set of t		Inic from aned in anily His Divors Not a a: Age a: Age Age Age Age Age Age Age Age Age Age	story (Tick ple colseparate policities of HV = ^^orecolseparate of HV = ^^orecolseparate	crie cho Estimate househo	ice) rd mariffy si incente Presist.No Fit care	Z24+3TC+E-Y      Reasons SUBSTITU nee dug untilitée     Reasons for 544     Reasons     R	TE: 1 toxic diug cut i H: 1 clinic south side sout of sto 5 7. TE 1 C C C D C D C C	city side el of stock, 7 al freatmas effacts, 2 sch, 7 psi/e Tubercs/ Tuber	Hects, 2 prog other reaso at failure, 2 h program, 2 mitosis treat fisci) the treat fack of fi mitosis treat fisci) the treat facts fisci) the treat fisci) the the the the the the the the the t	Inancy, 3 ris (specify) mmanologk 3 treatment arrow, 8 pal- tenent dur B registrati idrict: isath Centin B number restment o ] Rx failure lefe: fFollow-u	sk of pregne sal failure, 3 failure, 4 p Elent dedition ing HIV Ci ion e: ::::::::::::::::::::::::::::::::::	ancy, 4 newly GK i viologic failure protection of the second second second in the second s	S Il nose similari ompletad
Q 4-ingatient = 5-Para 11-10U outrearch =   □ pail ant branchered in Name previous cikit: 2. Personal History af KP = 1 Common af K	Stanc 0+PerfCt     to CSW callest     or an ART from another 1     (Trick one choice)     (Trick one choice)     void as worker (CSW)     rosecularization     ges.with news (MSM)     ang see (DU)     sched     void	r-san (D Seriu) 3-sthar UV canWART ( Date transl Auto transl Date transl Montal status parter biblis Family methal parter biblis Vical streatment vical streatment Series Auto Vical Streat	Inic from aned in: Divid His Divid a Not a a: Age n Not a e Chiefe Divid	story (Tick gle colseptrate pplicable e 4-Autoco	one cho Estimate househo	Kee) ed monthly til incomé fit care	D24/s1C+E-V     Reacons SLB63TUC     Reacons for SWTCH     Reacons	TE: 1 tasi ding out i H: 1 clinic south side sout of sta 5 7. TE 3 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	city side el of stock, 7 al treatman effacts, 2 ack, 7 pala Tuberest Rogimen ( 2ategory II 2ategory II	Hects, 2 prog other reason in failure, 2 is pregnancy, 2 influeis treat fack) T influeis treat fack) T influeis treat influeis	Inancy, 3 risin (specify) minumologik 3 treatment, arrow, 8 pail timent, dur B registrati listrict isath Centra B namber restment of 1 Rx failure lister dd 1 1 Ffeillow-tu	sk of pregna cal failure, 3 failure, 4 pr teerr decisie ing HIV ci ion e: plied plie	ancy, 4 menty dk viceologie failure our administrice, f. an 9 planned tre I Ourie I Ric et al.	sgnozed Ti 5 ill nees semant completed Fransfer or
A - inpatient: ■ 5-Pare     I + 10,DU outreach     I + 10,DU outreach     I + 10,DU outreach     I + 11,DU o	Later Constant	r-ssi [] Berliv Stehar [] VaniWART (: Date inansi [] Watter status [] Manter status [] Manter status [] Manter status [] Munder [] Widowski Farrity menta [] Partantohilite [] Vical typettinin arise ART	Inic from med in : ] Divorc Not a (thisto Place:	story (Tick gls relesperate policities + +Autoco	one cho Estimate househo	Kce) ad maxify yil income Pit care	Z2V+3TC+E-V     Reasons SUBSITV     read ong accitization     read ong accitization     Reasons STOF+ 1 to     Reasons Clores (lesk)     Paintersey TB     Senser-sealther     Benser-reagility     ExizyLationary     the     Death     Death     Death     Death	TE: 1 train drug cut i sticty side g out of str s 7. TB 1 C C C Date C C C C C C C C C C C C C C C C C C C	city side el of stock, 7 al treatman effacts, 2 ack, 7 pala Tuberest Regimen ( 2ategory II 2ategory II	flects, 2 prog other reason in Taikner, 2 is pregnancy, ent tack of fit if chi ficki) T B H Sign C Sign C S	Inancy, 3 risi ropecify) mmunologik 3 treatment arres, 8 pail timent dur B registratilization isath Cardin B rastratilization B rastration B rastratio B rastrati	sk of pregne cal failure, 3 tailure, 4 pr teent decisie ing HIV ci ion e: biod vice vice vice vice pied vice	ancy, 4 menty dki i virologie failura ore admenice, t u, o planned tre are	agnozed Ti 6 Il noze adment completed Fransfer or
A - inpatient: ■ 5-Pare     I + 10,DU untersch     I + 10,DU u	Stain: [] -0+PefCT []     I - conv cellent []     I - conv cellent []     Y (Tick core choice)     y (Tick core choice)     y (Tick core choice)     y (Tick core choice)     y (fick core choi	r-ssi [] BeHW T-sther 	Inic from med in : Divor: Not a Chiester Place:	story (Tick gle classifier and the second of the second se	cree cho Estimate houset ART O YIN	Kog) id manify jú inconé Pri cree	Z2V+3TC+E-V      Reasons SUBSTITU med rug available, 8      Reasons STOF 1     Descise class (int)     Putmenty TB     Smere-public     Smere-publ	TE: 1 trais ethog carl of H: 1 clinic south of strip TE 1 clinic p cost of strip TE 2 C C C C C C C C C C C C C C C C C C C	city side el of stock, 7 al treatman effects, 2 ack, 7 psi/ Tubercs Soglemen ( ategory 11 ategory 11 ategory 11 ategory 11 ategory 11 ategory 12 of error 1 Date s) Date Date	Hects, 2 prog other reason in tailure, 2 i pregramey, ent tack of fin lack of fin tack of	Inancy, 3 Hi (specify) mmanologic 3 treatment, 8 pail imment dar B registrati latrict restauent o 1 Ro failure latrict dd 1 Ffollow-u 00 ( m	sk of pregna cal failure, 3 calure, 4 po territ decision ing HIV calure ing HIV c	ancy, 4 newly 5k i vitologic failure our adressince, 1, a planned tre i Oure R or Dofran 1 Dofran 1 New clinic:	agnoted T 6 illnoce admant completed fransfer o









# **Cross-sectional ART database**

- Example of an Excel file
- CD4 count at baseline
- Viral load at 12 months of ART

	С	ross-s	sec	cti	ona	al	AR	T da	atak	base	5
AutoSave 💿 O	0 <b>.</b> 5 · 0	· L · C		Final 2	017 ART Cross-	sectional datab	ase 5.4.2018.x	lsx - Excel	exic.	K.A.M. ariyaratno	· · · ·
ile Home	Insert Dra	aw Page Layout Form	ulas Data	Review	View Help	,⊘ Tell m	e what you war	nt to do			
te Clipboard	Painter 5	$\begin{array}{c c} \bullet & 11 & \bullet & A^* & A^* \\ \bullet & 11 & \bullet & A^* & A^* \\ \bullet & 11 & \bullet & A^* & \bullet \\ \bullet & \bullet & A^* & \bullet \\ \hline Font & & 5 \\ \hline f_x & Any other clinic r \\ \hline \end{array}$	= ≫ •   = • • • Align 10.	🐉 Wrap Text	Center • \$ •	ral % ∮	Conditiona Formatting	al Format as Cell I * Table * Styles * Styles	nsert Delete Format Cells	AutoSum - A      Z     Fill - Z     Fill     Clear - Filte     Editing	& Find & r * Select *
Р	Q	т	U	v	w	x	Y	Z	AA	AB	AC
Age at ART initiation	Viral load at start of ART (+/_ 3 months)	Last Viral load in 2017	Last Viral load in 2017 <1000, 1000	Last Viral load in 2017 <1000, 1000 (excludi ng baseline	CD 4 at start of ART (+/_ 3 months)	CD 4 at start of ART (+/_3 months) <200, 200+	CD 4 at start of ART (+/_ 3 months) <350, 350+	Outcome as of end of4th Quarter 2017 (OT 1st ,OT 2nd,S, D, LFU)	Transfer In/ Transfer out/ same clinic	The clinic of currently followed up (by 4th Quarter 2017)	Current ART regime by4th Quarter 2017
33.0	× NA	46124	1000+	After Bl	202	200+	<350	OT1	same clinic	Polonnaruwa	TDE+ETC+EEV
43	NA	<34	<1000	After BL	212	200+	<350	OT1	same clinic	Colombo	TDF+FTC+EFV
39	NA	<34	<1000	After BL	305	200+	<350	OT1	same clinic	Colombo	TDF+FTC+EFV
22	43,011	466000	1000+	After BL	287	200+	<350	OT1(SUB)	same clinic	Colombo	TDF+FTC+ATV/r
44	NA	ND	<1000	After BL	217	200+	<350	OT1	same clinic	Colombo	TDF+FTC+EFV
42	NA	ND	<1000	After BL	233	200+	<350	OT1	same clinic	Colombo	TDF+FTC+EFV
30	NA	ND	<1000	After BL	317	200+	<350	OT1	same clinic	Colombo	TDF+FTC+EFV
32	NA	ND	<1000	After BL	263	200+	<350	OT1	same clinic	Colombo	TDF+FTC+EFV
→ To	2017 4th C	Net categorised Pivot	1 Data tab	olesGAM I	Datatables AR	(+)		: 4			











#### **PROGRAMME DATA – COMMON ISSUES**

- Collected for monitoring purposes; Using data for other purposes should be attempted with caution
- Mostly aggregate data; not individual level data
- Limited scope for slice & dice
- Varying quality from different centres at different times
- Changing formats and reporting mechanisms
- Changing programme strategies affect the data reported
- Scale up of centres
- Interest & ability; absence/ change of personnel affects data
- Commodity supplies & stock outs affect data
- Largely paper-based; needs computerisation

### **HIV TESTING DATA**

- Mostly no. of tests; not individuals tested
- Limited information on risk profiling
- Beneficiary segmentation
- Routes of transmission data
- Couple testing and sero-discordance
- Linkages to ART

### **PPTCT DATA**

- Mostly no. of tests; not individual pregnant women
- Documentation of trimester and gravida
- Background profile data
- Documentation of multiple testing is challenging
- Public sector Private sector overlaps
- Linkages with ART
- Positive pregnancy follow-up data, baby testing & outcomes
- Longitudinal data often not maintained

## **STD Clinic DATA**

- Footfalls vs STD Episodes vs Individuals
- Background profiles of STD Clinic Attendees
- Syndromic data vs Etiologic data
- Linkage to HIV testing
- Re-occurrence of STDs; Incidence vs Prevalence
- Correct denominators for case rates, incidence, prevalence etc.
- Partner testing data

## ART DATA

- Aggregate vs Individual data
- Cross-sectional vs longitudinal data
- Demographic, Epidemiological, clinical data
- PLHIV profiles
- Sero-discordance
- Viral suppression data
- Key population treatment data
- Building testing & treatment cascades at sub-national level

## **Assessment of Data Sources**

#### Utility & Usability:

S.No	Criteria	Score: 3	Score: 2	Score: 1
1	Explains Epidemic	Explains very well	Explains moderately	Doesn't explain
2	Reflects Programme Performance	Reflects very well	Reflects moderately	Doesn't Reflect
3	Availability of Data at the Desired Level	Easily Available	Not Easily Available	Not Available
4	Feasibility of Extraction & Use	Easy	Difficult	Very Difficult

## ISSUES AT VARIOUS STEPS (1)

- Collecting data/ Gathering Info from the beneficiaries or patients
  - Skipping questions
  - Way of asking questions & eliciting information
  - Interviewer bias
  - Judgmental approach
  - Workload, timings & staff cooperation

#### Documenting in Registers

- Using the correct registers & formats
- Inefficient documentation Info split into multiple registers
- Duplication
- Transcribing errors
- ▶ Filling the register without enquiring
- Blank cells filled later or left unfilled
- Illegible writing/ improper noting down

#### ISSUES AT VARIOUS STEPS (2)

#### Compilation/ Aggregation

- Lack of clarity of definitions/ Absence of data dictionary
- Errors in manual counting
- Wrong method of including/ excluding criteria
- Computation Errors/ Errors in data that needs to be computed from register

#### Data entry

- Typo errors in data entry
- Mismatch b/w register & data entered
- Incomplete data entry/ gaps/ blanks
- Entry in wrong fields/ wrong format

#### ISSUES AT VARIOUS STEPS (3)

#### Reporting

- Non-reporting/ Irregular reporting
- Lack of timeliness in reporting
- Reporting in wrong/ outdated formats

#### Compilation & Analysis

- Compilation errors
- Dataset merging issues
- Missed/ excluded in compilation
- Aggregation errors
- Errors/ Wrong approaches in analysis

### **Exercise 7: Assess NSACP Datasets**

- Review the NSACP Dataset that you have brought for the following. Each group to take up one prog dataset.
  - Time period
  - Geographic scope
  - Data lifecycle for the dataset
  - Assessment of Utility & Usability
  - Probable issues with the dataset at various steps
    - Collecting data
    - Documenting in Registers
    - Compilation/ Aggregation
    - Data entry
    - Reporting

#### **EXERCISE 3: ASSESSMENT OF NSACP DATASETS**

S.No.	Attributes of Dataset	<b>Observations/ Details</b>		
1	Name of the dataset			
2	Reference Period of dataset			
3	Geographic Scope	National/ Provincial/ District/ Clinic Details:		
4	Data Lifecycle	Who performs this function?	When/ At what frequency?	
А	Collecting data		• •	
В	Creating dataset			
С	Processing data			
D	Analysing data			
E	Preserving data			
F	Have access to data			
G	Publish results			
Н	Using & reusing data			
5	Assessment of Utility & Usability	Score the following on a scale of 3, where 3 is good, 2 is moderate, 1 is low		
А	Explains epidemic			
В	Reflects programme performance			
С	Availability of data at the desired level			
D	Feasibility of extraction & use			
6	Probable issues with the dataset at	Discuss & write the key issues affecting		
	various steps	the dataset at each step below.		
A	Collecting data			
В	Documenting in Registers			
С	Counting & Aggregation			
D	Data entry			
E	Reporting			
F	Compilation & Analysis			

Review the NSACP Dataset brought by you in your group and fill the following table.



# Overview

- What & Why
- Levels of Measurement
- Types of Variables
- Indicators
- Types of Indicators
- Computing Indicators

### DEFINITION

**Definition of an Element in Chemistry** 

- ... substance consisting of atoms
- … chemically the simplest substances
- ... cannot be broken down using chemical methods.
- ... can only be changed into other elements using nuclear methods.

**Definition of Variable – The Element of Research** 

- … consisting information
- … in simplest form
- … cannot be broken down into more basic info
- ... can be changed into other variables by some technique

#### Variables... Diff. Viewpoints

 ... that changes from person to person, from observation to observation (think in terms of height, weight, Hb status)

 ... either a result of some force or is itself the force that causes a change in another variable (think in terms of Obesity & Hypertension)

... characteristics of interest in a study that has different values for different subjects or objects (think in terms of 'proportion of ANC having Hb<7gm/dl)</p>

### Why to focus on key variables?

- > ... Provide focus when writing the Introduction section.
- ... Major terms to use when searching for research articles for the Literature Review.
- ... The key variables are the terms to be operationally defined.
- > ... Provide focus to the Methods section.
- ... The Instrument will measure the key variables.

Variables are the "things" we're measuring, or collecting data on, or forming groups on, in order to conduct our research study & answer the questions

#### Why Levels of Measurement?

- Indicate the relationship among the values assigned to different attributes of a variable
- E.g. Hypertension High, Medium, Low; >180/100 mmHg; 180/100 – 120/80 mmHg; <120/80 mmHg</p>
- helps to decide how to interpret the data from that variable..
- helps to decide what statistical analysis is appropriate on the values that were assigned.

#### LET'S CHECK THE VISUAL ACUITY OF A PATIENT

- Vision normal / Vision not normal
- Normal vision/ Visual impairment / Severe visual impairment/ Blind
- VA-6/6 to 6/18 , VA-6/18 to 6/60 , VA-6/60 to 3/60, VA-3/60 to NPL

Same characteristic measured in different ways



# Qualitative vs Quantitative



## Levels of Measurement

#### Nominal

- Qualities can not be graded
- No order preserving transformation
- E.g. Site of malignancy, constipation present, vomiting absent, Ram, Black

#### Ordinal

- Qualities can be graded
- Order preserving transformation
- E.g. Stage I, stage II, stage III CaCx
- Mild, Moderate and Severe Hypertension
- For convenience
- When quantitative scale not available or difficult

### Before we go further... Let's visit Sikkim


#### Levels of Measurement... Metric Scale

#### ▶ Interval

- Measured and have constant, equal distances between values
- Zero point arbitrary; No Absolute Zero
- E.g. IQ test difference between a 100 and a 110 equal to the scoring distance between 120 and 130, but no true zero on this test and an IQ of 140 is not twice as high as an IQ of 70.
- Distance between two cities
- Time interval between two doses
- Cannot say one is double or triple or so many times the other; (Becomes a ratio)

# LEVELS OF MEASUREMENT... METRIC SCALE

#### ▶ <u>Ratio</u>

- Agreed absolute zero; Same reference point
- Zero is meaningful
- E.g. Parity 4 and Parity 0; Height from the sea level;
- Laboratory tests;
- Weight of a person (A person weighs 100 kilos twice as heavy as a person who weighs 50 kilos (measure of zero kilos meaningful)
- Pulse rate, Respiratory rate, etc.



# LET'S CLASSIFY THE VARIABLES

Characteristic	Cases	Controls
Number of women (n)	832	846
Mean (SD) age (years)	55.3 (8.60)	55.7 (8.58)
Education		
No formal education or just elementary school	204 (24.5)	234 (27.7)
Middle school	503 (60.5)	513 (60.6)
College and above	125 (15.0)	99 (11.7)
Marital status		
Unmarried	14 (1.7)	10 (1.2)
Married or cohabiting	724 (87.0)	742 (87.7)
Separated, divorced, or widowed	94 (11.3)	94 (11.1)
Per capita income in previous year (yuan	)	
≤4166.7	230 (27.7)	244 (28.9)
4166.8-6250.3	243 (29.2)	242 (28.6)
6250.4-8333.3	57 (6.9)	50 (5.9)
≥8333.3	301 (36.2)	309 (36.6)
No of pregnancies		
None	62 (7.5)	35 (4.1)
1	137 (16.5)	109 (12.9)
2	199 (23.9)	208 (24.6)
3	194 (23.3)	207 (24.5)
4	141 (17.0)	157 (18.6)
25	99 (11.9)	130 (15.4)
cancer among first degree relatives	289 (34.7)	228 (27.0)
Oral contraceptive use	147 (17.7)	207 (24.5)
Regular exercises	253 (30.4)	287 (33.9)
Age at menarche*	14 (13 to 16)	15 (13 to 16)
Age at menopause (among postmenopausal women)*	50.1 (48.6 to 52.5)	49.4 (47.1 to 51.
Rody mass index*	25.1 (22.7 to 27.9)	23.7 (21.4 to 26.

# What do you say...

- A story of 100 meter race: Three runners are participating from three different districts of Madhya Pradesh . Each runner is assigned a color-coded T-shirt (BLUE/RED/GREEN) to differentiate from each other (Nominal scale). The winner is declared along with the declaration of first runner up and second runner up based on the criteria that who reaches the destination first, second and last. The rank order of runners such as "second runner up as 3", "first runner up as 2" and the "winner as 1" (Ordinal scale). During the tournament, 3 consequent races are arranged at a time interval of 10 day (Interval scale). The time spent by each runner in completing the race (Ratio scale) was also announced.
- Distance between two cities in kilometers
- Divorced, married and widowed
- Mild, moderate and Severe Anaemia
- Number of umbrellas sold in April, July, October
- GE cases among fish eaters and non- fish eaters

# Types of Variable Values

- Text variables
  - Coded
  - Un-coded
- Numeric variables
  - Continuous Measures
  - Codes
- Alphanumeric codes

### **Other Issues**

- Identify the correct Source Registers & Sections/ Columns in the register from where the variable is extracted
- > Explore the Dis-aggregations available for the variable
  - By gender
  - By age or age category (child/youth/adult/elderly)
  - By occupation
  - By marital status
  - By type of infection (viral/bacterial)
  - By source of referral
  - **b** By place of residence/ location of clinic

### Indicators

- Measures computed from one or more variables to reflect a process or an outcome
  - Program element that needs tracking
  - Measures an aspect of a program's performance
- Could be a number, percentage, rate or ratio
- E.g. No. of HIV tests done in a month; % of pregnant women receiving ANC services; STD cure rate; HIV transmission rate; Sex ratio;

# **Types of Indicators**

- Indicators of need/ Epidemic indicators
  - Size of KP, HIV prevalence
- Input Indicators
  - Finance, HR, Commodities
- Process Indicators
  - No. of outreach camps; No. of trainings; Counseling duration; Waiting time at ART clinic
- Output Indicators
  - No. of facilities opened; No. of tests done; No. of KP reached; Cascade indicators
- Outcome Indicators
  - > Condom uptake; No. of sexual partners; Health seeking behaviour
- Impact Indicators
  - ▶ New HIV infections/ Incidence; AIDS deaths/ Mortality; Survival

# **Computing Indicators**

- Indicator Title
- Indicator Type
- Indicator Definition
- Geographic Scope
- ▶ Time Reference
- Numerator & Denominator
- Inflation/ Deflation factors
- Units
- Use/ Utility/ Importance/ Application

### **Indicator Eq**

- Prevalence of Syphilis
- Epidemic Indicator/ Outcome Indicator
- Definition: Percentage of STD patients found reactive for Syphilis in last one year in a fixed geographic area/ at a STD clinic
- > Num: No. of STD patients found reactive for Syphilis in last 1 yr
- Den: No. of STD patients tested for Syphilis in last 1 yr
- Inflation/ Deflation factor: Proportion of STD patients tested for Syphilis out of all those who visited STD clinic
- ▶ Units: %
- Monitor STD epidemics; Reflect the success of STD control programme

#### Is it a Variable or Indicator?

- As of December 1, 2016, more than 1 million PLHIV were on antiretroviral treatment in India
- 45.3% of FSW respondents from Western Province in the IBBS (who reported having regular clients) reported having anal sex with regular clients
- > 7000 KP were tested for HIV during the FY 2015-16
- Herpetic to Non-herpetic STD Ratio at Colombo clinic is 1.2
- Condoms were distributed to 5,500 youth during the outreach camps

# **Exercise** 4

- Review the NSACP datasets that you have brought
- ▶ List out at least 20 variables from the dataset.
- Write the variable values.
- Classify them based on the level of measurement
- Mention the source register & the dis-aggregations available for the variable.
- Evolve 12 indicators 2 of each type based on the variables listed above. Classify them and indicate their num, den, units & importance

#### **EXERCISE 4: VARIABLES & INDICATORS**

PART A: Review the NSACP datasets that you have brought. List out at least 20 variables from the dataset. Write the variable values. Classify them based on the level of measurement. Mention the source register & dis-aggregations available for the

Variable Values S. Variable Level of Source Disaggregations Available No Measurement Register

variable.

Indicator 1	
Туре	
Num	
Den	
Units	
Importance	
Indicator 2	
Туре	
Num	
Den	
Units	
Importance	
Indicator 3	
Туре	
Num	
Den	
Units	
Importance	
Indicator 4	
Туре	
Num	
Den	
Units	
Importance	
Indicator 5	
Туре	
Num	
Den	
Units	
Importance	
Indicator 6	
Туре	
Num	
Den	
Units	
Importance	
Indicator 7	
Туре	
Num	
Den	
Units	
Importance	

#### PART B: Evolve 12 indicators – 2 of each type – based on the variables listed above. Classify them and indicate their num, den, units & importance

Indicator 8	
Туре	
Num	
Den	
Units	
Importance	
Indicator 9	
Туре	
Num	
Den	
Units	
Importance	
Indicator 10	
Туре	
Num	
Den	
Units	
Importance	
Indicator 11	
Туре	
Num	
Den	
Units	
Importance	
Indicator 12	
Туре	
Num	
Den	
Units	
Importance	



# OUTLINE

- What is Data Quality?
- Why DQA? Importance of DQ
- Reasons for poor DQ
- Conclusions from DQA
- Attributes of DQ



# DATA QUALITY

#### Services: the REAL world

In the *real world*, project activities are implemented in the field. These activities are designed to produce results that are quantifiable.

# Data: the INFORMATION SYSTEM

An information system represents these activities by collecting the results that were produced and mapping them to a recording system.







### Why DQA?

- An essential step before actual analysis
  - It is careful consideration of data collected, before embarking on analysis
- To help preparation of the working dataset/ database
  - To be used for further analysis
- It may open up further questions
  - E.g. HIV positivity rates from ICTC & HSS at same site, different from each other needs exploration of both datasets
- May require further more analysis
  - Eg. Settling an outlier identified in data
- Empowers the investigator to either adjust, correct, or totally reject the data for consideration in final dataset
  - Ability to make an informed decision

### Why Is IT IMPORTANT TO HAVE Quality Data?

- Accountability for funding and results reported increasingly important
- Quality data needed at program level for decision making
- Quality Data needed for any sort of operations research
- ▶ Quality data  $\rightarrow$  Robust evidence  $\rightarrow$  Appropriate policy



# **REASONS FOR POOR DATA QUALITY?**

- Not so well structured or well designed data collection tools
- Ill-framed questions
- Untrained personnel
- Biases
- Documentation errors
- Data entry errors
- Compilation errors
- Fraudulent practices

#### COMMON ERRORS dURING dATA COLLECTION

- Collecting incorrect information (e.g. information collected not as per the indicator definition)
- Collecting information from incorrect or unauthentic document
- Collecting information for incorrect reporting period
- Documenting the information improperly (e.g. making inadequate/rough notes which become difficult to comprehend later)
- Error in consolidating the data before entering into SIMS

### COMMON ERRORS dURING dATA ENTRY

- Entering data incorrectly (e.g. typing 10 instead of 01)
- Entering data against the wrong indicator
- Incomplete data entry (e.g. cells missed in the data entry format)
- Improper storing of offline data entry formats (e.g. not maintaining files properly in a folder, no or irregular data backup, etc.)
- Not doing version control- so that outdated files are uploaded (i.e. not uploading the latest version after revision)

# Common Errors during data validation

- Inbuilt validation check signals are not addressed and corrected
- Additional validation checks like checking outliers/out of range and missing values for indicators, consistency checks across different indicators
- Errors are made in cleaning the data even as validation checks are conducted

# DATA QUALITY - AT EVERY STEP of the Way

- Designing questionnaire/ Registers/ Tools
- Training of data collectors/ facility staff
- Testing the tools
- Recurrent training
- Developing an operational manual for questionnaire to:
  - Standardize the way of asking
  - Standardize the definitions

# **Conclusions from DQA**

- Data quality issues Present/ Absent
- Level of data quality Good/ Average/ Poor
- Identified specific quality issues
- Localise the quality issue (indicator; year/month; site/facility)
- Decisions about the quality issues
  - Accept as they are (Decide on acceptance limits)
  - Verify & Correct from original source (Decide on feasibility)
  - Impute, Adjust or Replace data (Decide the methods)
  - Reject the data/ Exclude from analysis (Decide on rejection limits)

# **Reject/ Exclude data, when...**

- The quality issue cannot be corrected
- It is a true reflection of the reality, but adversely affects the analysis
- The quality issue cannot be adjusted

## Impute missing data, when...

- Leaving blank affects the power of analysis
- A large segment of data is missing
- A small segment of data is missing & you don't want to have blank cells
- The variable or the case/ facility/ unit is important for analysis that you don't want to exclude, to avoid loss of other important attributes

# KEY ATTRIBUTES OF DATA QUALITY

- Availability/ Reporting Status
- Completeness
- Correctness/ Accuracy/ Validity
- Consistency
  - Consistency over time
  - Consistency Internal/ Reliability
  - Consistency External/ Validity/ Representativeness
- Precision
- Others
  - Timeliness
  - Integrity
  - Confidentiality

# DQA needs Most Granular Data

- DQA should be performed on the most granular data; not on aggregated data – Facility-wise, month/quarter-wise data; Individual-level data; Not on annual, province level data
- Granular data allows you to localise the error/ drill down to the exact issue
- Granular data allows you to fix the exact issue, without altering the other data which is good
- Recommended to do DQA for programme data at facility or district level; NOT ABOVE



# Availability/ Reporting Status

#### Availability of relevant indicators in reporting format

- Facility level (All RUs in the State):
  - ▶ No. of months in a year the RU has reported
  - Expressed as % out of 12
  - Categorised as 0-24%, 25-49%, 50-74% & 75-100%
- District & State levels:
  - Total no. of RU-Months reported
  - Expressed as % out of total no. of RU-Months (Tot no. of RUs in district or state x 12)

# Summarise

S.No.	Facility Name	Facility Type	No. of months reported: 2015-16	No. of months reported: 2016-17	Reporting %: 2015-16	Reporting %: 2016-17
1	Facility 1	STD Clinic	10	12	83	100
2	Facility 2	STD Clinic	6	8	50	66
3	Facility 3					

# How to assess & conclude?

- Need compiled dataset of one or more facilities
- Calculate Reporting Status
- Identify the non-reported months
- Try to obtain the non-reported data
- Decide acceptance limits:
  - > 75% reporting status in a year: Include in analysis
  - < 75% reporting status in a year: Exclude the facility from analysis</p>
- Impute the missing data if within acceptance limits
  - With average of other months/ neighbouring months

# Completeness

Completeness of data means that all variables for all reporting units are being collected and reported; Reporting formats are completely filled, without any blanks



# Completeness

No. of months the report is completely filled for the specified indicators

- Facility level (All RUs in the State):
  - No. of months in a year the RU has reported completely for the specified indicators
  - Expressed as % out of the no. of months in a year the RU has reported
- District & State level:
  - No. of RU-Months that have reported completely for the relevant indicators
  - Expressed as % out of total no. of RU-Months reported

### HOW TO ASSESS & CONCLUDE?

- Identify the variables needed from the dataset for the analysis.
- Review completeness for the identified variables for each facility
- Identify the facilities & months where the data is missing
- Try to obtain the missing data from the registers
- Decide acceptance limits:
  - > 75% completeness in a year: Include in analysis
  - < 75% completeness in a year: Exclude the facility from analysis
- Impute the missing data if within acceptance limits

# Summarise

S.No.	Facility Name	Facility Type	No. of months reported	No. of months reported completely	Comple teness %
1	Facility 1		10	10	100
2	Facility 2		14	7	50
3	Facility 3				

# **CORRECTNESS/ ACCURACY**

- Also sometimes known as validity (Valid data vs Invalid data)
- Accurate data are considered correct: the data measures what it is intended to measure
- Range checking/matching is conducted to check for accuracy
- **Example:** 
  - Recording Male STI syndrome (scrotal swelling) in a female STI syndrome data cell
  - Mentioning age as 234 years (could be 23 or 34)
  - Typing 39 instead of 93 (Typo errors)
  - > Ticking the wrong code on the questionnaire

	S Age	Itatistics				
1	N Vali	id 1	129			
	Mis	sing	0			
			Age			
÷		a – a	1	N	Cumulative	18
		Frequency	Percent	Valid Percent	Percent	
Valid	46	11	8.5	8.5	8.5	
	47	8	6.2	6.2	14.7	
	48	22	17.1	17.1	31.8	
	49	3	2.3	2.3	34.1	
	50	21	16.3	16.3	50.4	
	51	4	3.1	3.1	53.5	
	52	7	5.4	5.4	58.9	
	53	3	2.3	2.3	61.2	
	54	1	.8	.8	62.0	
	55	6	4.7	4.7	66.7	
	58	3	2.3	2.3	69.0	
	59	1	.8	.8	69.8	
	60	5	3.9	3.9	73.6	
	62	5	3.9	3.9	77.5	
	63	1	.8	.8	78.3	
	64	1	.8	.8	79.1	
	88	1	.8	8.	79.8	
	102	4	3.1	3.1	82.9	· · · · · · · · · · · · · · · · · · ·
2	245	4	3.1	3.1	80.0	
	543	4	3.1	3.1	69.1	
	875	4	2.2	0.1	52.2 04 B	
	997	3	2.5	2.3	97.0	
	999	7	2.2	2.2	100.0	
1	Total	120	100.0	100.0	100.0	

# Correctness (1)

Consistency b/w registers & reports and b/w two registers for each relevant indicator in selected months

- ► Facility level:
  - Verification Ratio b/w registers & reports (Max value 1)
    - Total no. in the source register in the selected month/ Total no. reported in the monthly report of the selected month
  - Cross-check consistency b/w two registers (Max value 100%)
    - No. of cases taken from Register 1 also found in Register 2/ Total no. of cases taken from Register 1
  - Reverse cross-check consistency b/w two registers
    - No. of cases taken from Register 2 also found in Register 1/ Total no. of cases taken from Register 2

### **CORRECTNESS** (2)

Consistency b/w registers & reports and b/w two registers for each relevant indicator in selected months

- District & State levels:
  - No. of RUs that have Verification Ratio of 1 in the selected month for the selected indicator
    - Expressed as % out of total no. of RUs where Trace & Verify is carried out in the district/state
  - No. of RUs that have Cross-check consistency of ≥ 90% in the selected month for the selected indicator
    - Expressed as % out of total no. of RUs where cross-check is carried out in the district/state
  - No. of RUs that have Reverse Cross-check consistency of ≥ 90% in the selected month for the selected indicator
    - Expressed as % out of total no. of RUs where reverse crosscheck is carried out in the district/state

## How to assess and conclude?

- Conduct range checks for all numeric variables & identify out of range values
- Errors within the range & text data errors can be identified by
  - Verification & Cross-check with registers Gold standard; Random checks of a few data points, if not all
  - Identifying & verifying implausible data data that does not seem to be correct/ not likely to be correct
  - Identifying implausible combinations of variables e.g. agespecific variables & values (age & marital status, age & occupation, age & education), gender-specific variables (female & house-wife, female & received ANC services), etc.
  - Identifying data points that stand out from the rest/ looks like an outlier

### How to assess and conclude?

- Highlight the incorrect data points
- Delete confirmed errors & exclude from analysis
- Impute, adjust or replace errors by an appropriate method

#### **Summarise**

Correctnes	Correctness Trace & Verify Cros		oss Check				
Facility Name	Facility Type	No. in Source Register	No. in Monthly Report	Verificati on Ratio	No. of cases taken from Register 1	No. of cases found in Register 2	Cross- check Consist ency %
Facility 1	STD Clinic	10	10	1.00	10	10	100
Facility 2	STD Clinic	14	15	0.93	10	8	80
Facility 3	STD Clinic	9	5	1.80	10	5	50

# **CONSISTENCY - INTERNAL**

- Also known as Reliability
- Refers to the degree of similarity of information obtained when the data are measured and collected in the same method every time
- Can be checked by comparing reported data with different response checks or other primary sources like random respondent check.
- Example:
  - Number of individuals tested: "received post-test counseling" or "collected tests within 7 days"
  - While collecting data on male sexual behavior the number of regular sexual partners is reported to be more than total number of sexual partners in last three months
  - Respondent says s/he is not aware of HIV but states later that s/he has availed HIV testing services.
  - The birth date and age of the respondent are inconsistent

# **CONSISTENCY** - INTERNAL

#### Example:

- Number of individuals tested & received pre-test counseling; received post-test counseling & collected test results within 7 days
- While collecting data on male sexual behavior the number of regular sexual partners is reported to be more than total number of sexual partners
- Respondent says s/he is not aware of HIV but states later that s/he has availed HIV testing services.
- > The birth date and age of the respondent are inconsistent

# Outliers

- Outliers are disturbing...it is best to recognize them and make an informed decision to include or exclude the value..
- <u>Outlier</u>: specific values that differ very markedly from the other 'usual' values
- These specific values have an unduly large effect on statistics and may need special handling
- Their presence makes it difficult to draw valid conclusion
- Ideally, all analysis should be performed with and without including the outlier value
  - There might, in fact, be a true reason for the apparently strange data that should not be ignored

# **C**ONSISTENCY OVER TIME

Absence of outliers for each relevant indicator

- Facility level
  - No. of months in a year where there are no outliers in monthly data
  - Expressed as % out of the no. of months in a year the RU has reported
- District & State levels
  - Tot no. of RU-Months where there are no outliers
  - Expressed as % out of total no. of RU-Months reported



Compare numbers tested (Target ANC=400, others=250) HIV positives among ANC attendees, by site and year

District X	HIV Sentinel Surveillance (ANC)				
	Number	Number			
Year	Tested	Positive	% positive		
2002	400	1	0.25		
2003	400	13	3.25		
2004	400	4	1		
2005	400	3	0.75		
2006	400	0	0		
2007	400	2	0.5		

- ▶ Number tested is 400 across years (if <75%, disregard)
- The site is a consistent site for last six years

### Another example

- While examining trend of HIV positivity among STD clinic attendees, you identify an outlier in one district in one year
- On drilling down, you identify that the outlier is coming from one clinic in the district
- On further drilling down, you identify that the outlier is coming from one quarter data reported by that clinic in that year
- Fixing this one data point will fix the trend

#### HOW TO ASSESS & CONCLUDE

- Identify the outlier
  - Plot a trend line & simple eyeballing of a trend line
  - Build confidence intervals Identify non-overlapping confidence intervals
  - Figures falling outside 3 standard deviations
- Investigate the outlier
  - Look for data entry errors
  - Examine if the sample size is very low leading to large data fluctuations
  - Look for clustering by day or place or site
  - Check if the profile of the beneficiaries is different from the rest
  - For testing data, look for any consecutive positives due to probable lab contaminations
- Adjust the outlier Mean imputation; 3-yr moving average





#### **CONSISTENCY – EXTERNAL VALIDITY**

- External consistency or Validity refers to the representativeness of the data to the population from which it is collected
- This is more applicable to survey data
- For clinic based data, this is not directly applicable as the segment of population who attend the clinic are not exactly same as the general population; But, they can be used as proxy indicators for specific sub-populations
- If coverage levels are high, near 90%, usually, external validity is assumed to be present
- Data with very low sample sizes/ denominators is usually excluded, due to poor external validity



#### Precision

This means that the data have sufficient detail. Precise data collection plans/systems ensure that all data disaggregations that are required, are collected. Inadequate data affects the quality of the entire analysis

#### Example:

- Totals for positives put on treatment, but no mention of gender or age disaggregations, or type of key population that it serves
- In a study exploring causes of gender disparity in availing HIV counseling & testing in a particular district, a data collection plan/system lacks precision if it is not designed to record the attitudes of men & women separately

#### Timeliness

- Data are timely when they are up-to-date (current), and when the information is available on time.
- Example:
  - Service delivery sites submit the previous month's report at different times in a month instead of by a certain date in the month as per reporting timelines.
  - Data collection on "treatment seeking behavior among migrant men for STI/STDs" starts in 2008 when there is no STD clinic in the migration habitation area and ends in 2012 when there are two STD clinics set up by the new TIs in the catchment area.
  - Data is collected for a period which is not within the study period, or from a period wherein no intervention was being carried out.

#### Integrity

Integrity is when data generated/collected are protected from deliberate bias or manipulation for political or personal reasons. Manipulation could be done to benefit the project or program.

#### Example:

- Data would lack integrity if an ART centre inflates the number of PLHIV on ART in order to reach the targets defined by NACO
- Interviewer deliberately records "no sex happened in last one month" to avoid the section on current sexual behavior
- During training data collectors are encouraged by the implementers to get positive feedback from beneficiaries on an intervention which is being evaluated

#### Confidentiality

Confidentiality means that respondents are assured that their information will be maintained according to national and/or international standards for data. This means that personal data are not disclosed inappropriately, and that data in hard copy and electronic form are treated with appropriate levels of security.

#### Example:

- Care should be taken to ensure that during different steps of research , confidentiality of data is not breached
- Basic ethical principles of the Belmont Report (1979) should be adhered to- Respect for Persons, Beneficence, Justice
- Registers containing personally identifiable information (PII) should be kept in locked cabinets
- Confidentiality agreements for data collectors

SUMMING UP						
Data Quality Attribute	Indicator to Assess	Key thing to look for				
Availability/ Reporting Status	Reporting %	Non-reporting facilities & quarters				
Completeness	Completeness %	Missing Data				
Correctness/ Accuracy	Verification Ratio & Cross-check Consistency %	Incorrect entries; Invalid data;				
Internal Consistency/ Reliability		Cross-checks within data				
Consistency over time	% of facilities without outliers over a year	Outliers				
External Consistency/ Validity/ Representativeness		Sample size & Coverage %				
Precision		Disaggregation available				



# STEPS TO IMPROVE DATA QUALITY

- Good understanding of registers, columns, variables & definitions
- Good idea of normal or usual range or pattern of data
- > Attention & clarity during data recording in the registers
- Correctness & completeness in registers
- Good understanding of mapping of variables/columns in registers to fields/ cells in the reporting formats
- Accuracy in counting, keeping in mind the defining variables & disaggregations
- Focus during data entry to avoid data entry errors
- Self-check/ Self-review of the entered data for DQ before submitting

# **Adjustment & Validation**

- Adjustment refers to taking action to fix the data quality issues and improve the quality of data
- Adjustment also refers to the actions to make the dataset more suitable and appropriate for analysis
- Validation refers to comparing the quality-controlled data with other data/ other source to ensure that we are dealing with plausible data

# Summing Up...DQA

Data Quality Attribute	Indicator to Assess	Key thing to look for
Availability/ Reporting Status	Reporting %	Non-reporting facilities & quarters
Completeness	Completeness %	Missing Data
Correctness/ Accuracy	Verification Ratio & Cross-check Consistency %	Incorrect entries; Invalid data;
Internal Consistency/ Reliability		Cross-checks within data
Consistency over time	% of facilities without outliers over a year	Outliers
External Consistency/ Validity/ Representativeness		Sample size & Coverage %
Precision		Disaggregation available
#### **Conclusions from DQA**

- Data quality issues Present/ Absent
- Level of data quality Good/ Average/ Poor
- Identify specific quality issues
- Localise the quality issue (variable; year/month; site/facility)
- Decisions about the quality issues
  - Accept as they are (Decide on acceptance limits)
  - Verify & Correct from original source (Decide on feasibility)
  - Impute, Adjust or Replace data (Decide the methods)

#### Actions to fix quality issues

- Leave it as it is, if it is not of great importance or if it does not affect your analysis significantly
- Delete only the concerned data point or Delete the entire case (facility or month or individual); Known as List Deletion
- Verify and correct from the source data/ registers, if possible and feasible
- Impute the missing values or go with missing data/ blank cells if it is acceptable
- Recode the missing values, if required



# Adjustments

Data Quality Attribute	Key thing to look for	Adjustment
Availability/ Reporting Status	Non-reporting facilities & quarters; Unreported data	Leave, Correct or Impute
Completeness	Missing Data	Leave, Correct or Impute
Correctness/ Accuracy	Incorrect entries; Invalid data;	Leave, Correct or Impute
Internal Consistency/ Reliability	Cross-checks within data	Leave or exclude; Can't do anything
Consistency over time	Outliers	Leave, Correct, Smoothen, adjust or replace
External Consistency/ Validity/ Representativeness	Sample size & Coverage %	Leave or exclude; Can't do anything
Precision	Disaggregation available	Leave; Can't do anything

## Analysis of Missing Patterns

- Missing At Random (MAR)
- Missing Completely at Random (MCAR)
- Missing Not At Random (MNAR)



### IMPUTE MISSING dATA, WHEN...

- Leaving blank affects the power of analysis
- A large segment of data is missing
- A small segment of data is missing & you don't want to have blank cells
- The variable or the case/ facility/ unit is important for analysis that you don't want to exclude, to avoid loss of other important attributes

#### **Methods of Imputation**

- Mean Imputation Replace with average of all months/ all facilities; average of before & after; moving averages;
- Mode Imputation Replace with the most common value
- Nearest Neighbourhood Imputation Replace with the value from the nearest similar case/ facility/ district
- Univariate/ Single Imputation Methods
  - Linear imputation
  - Logistic regression imputation
- Multiple Imputation by Chained Equations (MICE)
  - Using Linear Regression/ Predictive Mean Matching method for scale variables
  - Using Logistic Regression for categorical variables
- More advanced methods also available

### REJECT/ Exclude data, when...

- The quality issue cannot be corrected
- It is a true reflection of the reality, but adversely affects the analysis
- The quality issue cannot be adjusted

#### Don't Overkill

- > Too much manipulation of original data to be avoided
- Too much of imputation of missing data not advised
- Leave the originality of data as much as possible
- Adjust and impute only where it is of significant value addition for the analysis

# Other Data Adjustments

- Inflation or deflation factor
- Adjust for specific demographic variables. E.g. Age adjustment
- Standardisation Adjust to a standard population
- Correction factors
- Calibration

# Validation of QC DATA

- Compare with external gold standard Values/ Distribution/ Patterns
- Compare with averages from similar complete data for other years, states, countries, etc
- Validate against known epidemiological and programmatic understanding of the issue and expert consensus

				Exercise	Dataset									
Sno	Province	District	Month	Name of the Facility	Total tested during the month	Total clients received post- test and results	Total ANC positives during the month	Total ANC Positives linked to ART during the month	Total Non ANC positives during the month	Total Non ANC clients linked to ART during the month	Total Positives (ANC+Non ANC)	Positivity Rate (%)	Total Positive clients linked to ART(ANC+Non ANC)	Linkage loss during the month
1	WP	Dist1	JAN 18	ICTC AMALAPURAM	702		0	0	14	12				
2	WP	Dist1	JAN 18	ICTC RAMACHANDRAPURM	483		1	1	12	11				
3	WP	Dist1	JAN 18	ICTC TUNI	652		0	0	10	9				
4	WP	Dist1	JAN 18	ICTC PEDDAPURAM	437		2	2	16	22				
5	WP	Dist1	JAN 18	ICTC DH – RAJAMEHINDRAVARAM	742		0	0	76	69				
6	WP	Dist1	JAN 18	ICTC RMC KAKINADA	722		0	0	46	44				
7	WP	Dist1	JAN 18	ICTC TBIDH GGH KAKINADA	583		0		23	18				
8	WP	Dist3	JAN 18	SA AH - BAPATLA (P&V)	815		0		154	15				
9	WP	Dist3	JAN 18	SA AH - NARSARAOPET (P&V)	1243		2		34	32				
10	WP	Dist3	JAN 18	SA CHC - CHILAKALURIPETA (V)	590		2	2	7	7				
11	WP	Dist3	JAN 18	SA CHC - MACHERLA (V)	350			0	5	7				
12	WP	Dist3	JAN 18	SA CHC - VINUKONDA (I)	315		1	1	9	8				
13	WP	Dist3	JAN 18	SA DH - TENALI (V)	673			0	39	35				
14	WP	Dist3	JAN 18	SA GGH - GUNTUR (V)	1480			0	54	46				
15	WP	Dist3	JAN 18	SA TBH - GUNTUR (I)	422			0	15	13				
16	WP	Dist2	JAN 18	SA AH - GUDIVADA (P&V)	1089		0	0	10	8				
17	WP	Dist2	JAN 18	SA AH - NUZIVEEDU (P&V)	874		0	0	14	13				
18	WP	Dist2	JAN 18	SA CHC - NANDIGAMA-OLD (V)	401		1	3	20	20				
19	WP	Dist2	JAN 18	SA CHC - RAIIVNAGAR-II (V)	602		0	0	17	15				
20	WP	Dist2	JAN 18	SA DH - MACHILIPATNAM (V)	508		0	0	12	15				
21	WP	Dist2	JAN 18	SA GGH - SIDDHARTHA MEDICAL COLLEGE (V)	1670		0	0	78	57				
22	WP	Dist2	JAN 18	SA GGH - VIJAVAWADA (P)	1022		1	1	3	3				
22	ND	Dist2	JAN 10		1521	1455	0	0	16	9				
23	ND	Dist4	JAN 10	Shatabdi Goyandi	1196	1455	0	0	10	12				
24	ND	Dist4	JAIN 10		2061	2559	0	0	15	20				
25	NP	Dist4	JAIN 10 JAIN 19	ITMCH	2922	2404	0	0	49	39				
20	ND	Dist4	JAN 10		2020	2434	2	3	18	20				
27	ND	Dist4	JAN 10	Nair Hos	2046	1079	1	1	40	12				
28	ND	Dist4	JAN 10	G T Hospital	2040	1112	1		51	45				
29	NP	Dist4	JAN 10	G. T. Hospital	307	1045	0	0	24	0				
30	NP	Dist4	JAN 18	n. Bridgawati nospital	2015	1845	0	0	12	8				
31	NP	Dist4	JAN 18	Sidufiaf Lina Nagar Hospital	1641	704	1	0	13	0				
32	NP	DISt4	JAN 18		1041	1494	1	1	14	8				
33	NP	Dist5	JAN 18	A F M C, PUNE VCTC (SDLIC2725000220983 )	1386	1369	0	0	22	13				
34	NP	DIST5	JAN 18	BHARATHTHOSPITAL & MEDICAL COLLEGE (SDLIC2725000201912)	405	405	0	0	5					
35	NP	Dist5	JAN 18	BHOSARI HOSPITAL (SDLIC2725000202045)	/12	683	0	0	9	2				
36	NP	Dist5	JAN 18	BJMC PUNE VCTC (SDLIC2/250002020049 )	0	0	0	0	0	0				
37	NP	Dist5	JAN 18	CHEST & DISTRICT HOSPITAL AUNDH (SDLIC2/25000201/97)	702	692	0	0	14	14				
38	NP	Dist5	JAN 18	DR.D Y PATIL MEDICAL COLLEGE & HOSPITAL (SDLIC2/25000202047)	327	327	0	0	5	5				
39	NP	Dist5	JAN 18	MAHARASHTRA AROGYA MANDAL (SDLIC2725000201911)	436	436	0		0	0				
40	NP	Dist5	JAN 18	MIMER MEDICAL COLLEGE (SDLIC2/25000202046 )	/14	/14	0		4	4				
41	NP	Dist5	JAN 18	RAJIV GANDHI HOSPITAL (KHANSAHEB PESTANJI MATERNITY HOME) (SDLIC2725000216075 )	0	0	0		0	0				
42	NP	Dist5	JAN 18	RH MANCHAR (SDLIC2725000202042 )	630	630	0	0	6	6				
43	NP	Dist5	JAN 18	RH SHIRUR (SDLIC2725000220634 )	244	244	0	0	3	0				
44	NP	Dist5	JAN 18	SDH INDAPUR (SDLIC2725000220633 )	287	284	1	1	5	1				
45	NP	Dist5	JAN 18	SILVER JUBILEE GOVT RURAL HOSPITAL,BARAMATI (SDLIC2725000216081 )	239	239	0	0	16	16				
46	NP	Dist5	JAN 18	SMT KASHIBAI NAVALE HOSPITAL PPP (SDLIC2725000222763 )	1659	1542	0	0	15	10				
47	NP	Dist5	JAN 18	SONAWANE HOSPITAL PMC (SDLIC2725000201798 )	398	400	0	0	8	0				
48	NP	Dist5	JAN 18	YASHWANTRAO CHAVAN MEMORIAL HOSPITAL (SDLIC2725000203522 )	854	847	2	2	35	29				
49	NP	Dist6	JAN 18	V S GENERAL HOSPITAL I (SDLIC2721001414877 )	366	366	0	0	6	6				
50	NP	Dist6	JAN 18	MEENATAI THAKARE THANE(SHIVAJI NAGAR) (SDLIC2721001420585 )	198	197	1	1	2	2				
51	NP	Dist6	JAN 18	DR A G JOSHI HOSPITAL (SDLIC2721001421951 )	152	152	0	0		2				
52	NP	Dist6	JAN 18	RAJIV GANDHI MEDICAL COLLEGE, KALVA (SDLIC2721001414883 )	990	1157		1	15	15				
53	NP	Dist6	JAN 18	SUTIKAGRUH HOSPITAL (SDLIC2721001421955 )	413	407		1	2	2				
54	NP	Dist6	JAN 18	RUKMINIBAI HOSPITAL (SDLIC2721001416388 )	710	707	0	0	14	14				

Sno	Province	District	Month	Name of the Facility	Total tested during the month	Total clients received post- test and results	Total ANC positives during the month	Total ANC Positives linked to ART during the month	Total Non ANC positives during the month	Total Non ANC clients linked to ART during the month	Total Positives (ANC+Non ANC)	Positivity Rate (%)	Total Positive clients linked to ART(ANC+Non ANC)	Linkage loss during the month
55	NP	Dist6	JAN 18	SHASHTRINAGAR GENERAL HOSPITAL (SDLIC2721001416392 )	458	458	0	0	6	6				
56	NP	Dist6	JAN 18	CENTRAL HOSPITAL ULHASNAGAR (SDLIC2721001401795 )	765	696	0	0	135	13				
57	NP	Dist6	JAN 18	BALKUM (CHHAYA HOSPITAL, ABERNATH ) (SDLIC2721001416961 )	689	469	0	0	7	7				
58	NP	Dist6	JAN 18	HEALTH CENTRE BHAYANDAR-W (SDLIC2721001420583 )	593	592	0	0	3					
59	NP	Dist6	JAN 18	HEALTH CENTRE MIRA ROAD (SDLIC2721001421947 )	586	558	2	2	13	13				
60	NP	Dist6	JAN 18	I G M HOSPITAL, BHIWANDI (SDLIC2721001416389 )	404	411	3	3	28	28				
61	NP	Dist6	JAN 18	NAVI MUMBAI (GEN.HOSPITAL),VASHI (SDLIC2721001401796 )	1193	1185	0	0	30	30				
62	NP	Dist6	JAN 18	RAJMATA JIJAU MATA VAL RUGNALAYA,AIROLI (SDLIC2721001416959 )	1010	997	0	0	7	5				
63	NP	Dist6	JAN 18	PAD. DR D Y PATIL MEDICAL COLLEGE & HOSPITAL, NERUL NAVI MUMBAI (SDLIC2721001414879)	832	831	0		2	2				
64	NP	Dist6	FEB 18	V S GENERAL HOSPITAL I (SDLIC2721001414877 )	330	330	0		10	10				
65	NP	Dist6	FEB 18	MEENATAI THAKARE THANE(SHIVAJI NAGAR) (SDLIC2721001420585 )	112	112	0		0	0				
66	NP	Dist6	FEB 18	DR A G JOSHI HOSPITAL (SDI IC2721001421951 )	177	176	0	0	2	2				
67	NP	Dist6	FEB 18	RAIV GANDHI MEDICAL COLLEGE, KALVA (SDLIC2721001414883.)	990	960	1	1	15	13				
68	NP	Dist6	FFB 18	SUTIKAGRUH HOSPITAL (SDUC2721001421955.)	378	375	0	0	4	4				
60	NP	Dist6	FEB 18	RUKMINIBALHOSPITAL (SDUC2721001/15388.)	633	629	0	0	10	9				
70	NP	Dist6	FEB 18	SHASHTRINAGAR GENERAL HOSPITAL (SDLIC2721001416392)	444	441	0	1	7	7				
70	ND	Disto	EED 10	CENTRAL HOSPITAL III HASNAGAR (SDLIC2721001410352)	765	726	0	0	24	22				
71	ND	Disto	EED 10		691	207	0	0	24					
72	NP	Disto	FED 10	BALKOWI (CHHATA HOSPITAL, ABERNATH ) (SDEIC2/21001410501 )	600	307	0	0	0	0				
73	NP	Disto	FED 10	HEALTH CENTRE MIRA ROAD (SDUC2721001420363 )	500 512	498	0	0	4	4				
74	NP ND	Disto	FED 10	HEALTH CENTRE MIRA ROAD (SDEIC2721001421947)	313	300	1	0	10	10				
75	NP	Dist6	FEB 18	TG M HOSPITAL, BHIWANDI (SDLIC2/21001416389 )	293	292	1	2	30	30				
/6	NP	Dist6	FEB 18	NAVI MUMBAI (GEN.HOSPITAL), VASHI (SDLIC2/21001401/96)	1014	1006	0	0	2	23				
//	NP	DISTO	FEB 18	RAJMATA JIJAU MATA VAL RUGNALAYA,AIRULI (SDLIC2/21001416959 )	883	862	0	0	3	3				
78	NP	Dist6	FEB 18	PAD. DR D Y PATIL MEDICAL COLLEGE & HOSPITAL, NERUL, NAVI MUMBAI (SDLIC2/210014148/9)	699	693	0	0	4	4				
79	NP	Dist5	FEB 18	A F M C, PUNE VCTC (SDLIC2725000220983)	1082	1122	0	0	9	/				
80	NP	Dist5	FEB 18	BHARATHI HOSPITAL & MEDICAL COLLEGE (SDLIC2725000201912 )	190	190	1	1	4	4				
81	NP	Dist5	FEB 18	BHOSARI HOSPITAL (SDLIC2725000202045 )	282	280	3	3	3	3				
82	NP	Dist5	FEB 18	BJMC PUNE VCTC (SDLIC2725000202049 )	951	923	0	0	77	44				
83	NP	Dist5	FEB 18	CHEST & DISTRICT HOSPITAL AUNDH (SDLIC2725000201797 )	348	348	8	8	14	14				
84	NP	Dist5	FEB 18	DR.D Y PATIL MEDICAL COLLEGE & HOSPITAL (SDLIC2725000202047 )	236	236	1	1	10	10				
85	NP	Dist5	FEB 18	MAHARASHTRA AROGYA MANDAL (SDLIC2725000201911 )	225	225	7	7	10	10				
86	NP	Dist5	FEB 18	MIMER MEDICAL COLLEGE (SDLIC2725000202046 )	364	364	4	4	0	0				
87	NP	Dist5	FEB 18	RAJIV GANDHI HOSPITAL (KHANSAHEB PESTANJI MATERNITY HOME) (SDLIC2725000216075 )	0	0	1	1	0	0				
88	NP	Dist5	FEB 18	RH MANCHAR (SDLIC2725000202042 )	301	301	1	1	5					
89	NP	Dist5	FEB 18	SDH INDAPUR (SDLIC2725000220633 )	152	152	2	2	5	2				
90	NP	Dist5	FEB 18	SILVER JUBILEE GOVT RURAL HOSPITAL,BARAMATI (SDLIC2725000216081 )	119	119	2	2	14	14				
91	NP	Dist5	FEB 18	SMT KASHIBAI NAVALE HOSPITAL PPP (SDLIC2725000222763 )	1136	1068	15	15	23	11				
92	NP	Dist5	FEB 18	SONAWANE HOSPITAL PMC (SDLIC2725000201798 )	240	234	2	2	66	5				
93	NP	Dist5	FEB 18	YASHWANTRAO CHAVAN MEMORIAL HOSPITAL (SDLIC2725000203522 )	356	350	14	3	34	32				
94	NP	Dist4	FEB 18	Dr.R.N.Cooper	1352	1320	0	0	16	5				
95	NP	Dist4	FEB 18	Shatabdi Govandi	814	772	1	1	17	23				
96	NP	Dist4	FEB 18	KEM Hos.	3686	3876	33	1	35	21				
97	NP	Dist4	FEB 18	LTMGH	2521	2146	1	1	53	32				
98	NP	Dist4	FEB 18	J.J. Hos.	3112	2951	2	2	55	45				
99	NP	Dist4	FEB 18	Nair Hos.	1772	1727	0	0	35	31				
100	NP	Dist4	FEB 18	G.T. Hospital	822	806	0	0	14	1				
101	NP	Dist4	FEB 18	H. Bhagawati Hospital	1771	1682	1	1		6				
102	NP	Dist4	FEB 18	Siddhartha Nagar Hospital	596	571	-	ō	19	11				
103	NP	Dist4	FEB 18	Rajawadi Hospital	1395	1309	2	1	9	6				
104	WP	Dist1	FEB 18	ICTC AMALAPUBAM	793		0	0	10	10				
104	W/P	Dist1	FEB 18		544		1	1	20	20				
105	W/D	Dist1	FFP 19		6592		0	0	14	17				
107	\\/D	Dist1	EEB 19		522		2	2	10	15				
100	\\/D	Dist1	EEB 19		977		2	2 0	10/	10				
100	\\/D	Dist1	EEB 19		807		0	0	104	50				
103	VVP	DISLI	LED TQ	ICTC NIVIC KAKINADA	007		U	U	44	44				

Sno	Province	District	Month	Name of the Facility	Total tested during the month	Total clients received post- test and results	Total ANC positives during the month	Total ANC Positives linked to ART during the month	Total Non ANC positives during the month	Total Non ANC clients linked to ART during the month	Total Positives (ANC+Non ANC)	Positivity Rate (%)	Total Positive clients linked to ART(ANC+Non ANC)	Linkage loss during the month
110	WP	Dist1	FEB 18	ICTC TBIDH GGH KAKINADA	739		0	0	27	24				
111	WP	Dist3	FEB 18	SA AH - BAPATLA (P&V)	951		0	0	14	14				
112	WP	Dist3	FEB 18	SA AH - NARSARAOPET (P&V)	1352		1	1	31					
113	WP	Dist3	FEB 18	SA CHC - CHILAKALURIPETA (V)	618		1	1	13	13				
114	WP	Dist3	FEB 18	SA CHC - MACHERLA (V)	35		0	0	7	6				
115	WP	Dist3	FEB 18	SA CHC - VINUKONDA (I)	332		0	0	10	8				
116	WP	Dist3	FFB 18	SA DH - TENALI (V)	625		0	0	45	39				
117	WP	Dist3	FFB 18	SA GGH - GUNTUB (V)	1521		0	0	66	76				
118	WP	Dist3	FFB 18	SA TBH - GUNTUR (I)	358		0	0	14	13				
110	W/P	Dist2	FEB 18	SA AH - GUDIVADA (P&V)	1068		2	2	124	12				
120	W/P	Dist2	FEB 18	SA AH - NUZIVEEDU (P&V)	890		4	2	18	12				
120	W/P	Dist2	FEB 18	SA CHC - NANDIGAMA-OLD (V)	574		0	0	15	15				
121	W/D	Dist2	EED 10	SA CHC - NANDIGAMA OLD (V)	602		0	0	15	12				
122	W/P	Dist2	EED 10		668		12		25	24				
123	W/D	Dist2	FED 10		10842		12	0	2.5	24				
124	WP	Dist2	FEB 18		19842		1	0	91	/1				
125	WP	DISLZ	FEB 18	SA GGE - VIJATAWADA (P)	1089		1	1	14	14				
126	WP	Dist2	MAR 18	SA AH - GUDIVADA (P&V)	1316		0	0	15	15				
127	WP	Dist2	IVIAR 18	SA AH - NUZIVEEDU (P&V)	1146		2	2	10	37				
128	WP	Dist2	MAR 18	SA CHC - NANDIGAMA-OLD (V)	532		1	1	18	9				
129	WP	Dist2	MAR 18	SA CHC - KAJIVNAGAK-U (V)	384		0	0	15	15				
130	WP	Dist2	MAR 18	SA DH - MACHILIPATNAM (V)	835			0	23	23				
131	WP	Dist2	MAR 18	SA GGH - SIDDHARTHA MEDICAL COLLEGE (V)	2297			0	104					
132	WP	Dist2	MAR 18	SA GGH - VIJAYAWADA (P)	1421		4	5	11	11				
133	WP	Dist3	MAR 18	SA AH - BAPATLA (P&V)	808		1	1	9	9				
134	WP	Dist3	MAR 18	SA AH - NARSARAOPET (P&V)	1334		1	1	26	21				
135	WP	Dist3	MAR 18	SA CHC - CHILAKALURIPETA (V)	531		0	0	13	13				
136	WP	Dist3	MAR 18	SA CHC - MACHERLA (V)	352		0	0	5	3				
137	WP	Dist3	MAR 18	SA CHC - VINUKONDA (I)	513		0	0	9	6				
138	WP	Dist3	MAR 18	SA DH - TENALI (V)	679		0	0	27	24				
139	WP	Dist3	MAR 18	SA GGH - GUNTUR (V)	1598		0	0	78	71				
140	WP	Dist3	MAR 18	SA TBH - GUNTUR (I)	390		0	0	20	18				
141	WP	Dist1	MAR 18	ICTC AMALAPURAM	823		1	1	17	17				
142	WP	Dist1	MAR 18	ICTC RAMACHANDRAPURM	605		1	1	12	12				
143	WP	Dist1	MAR 18	ICTC TUNI	677		0	0	11	11				
144	WP	Dist1	MAR 18	ICTC PEDDAPURAM	477		1	1	12	12				
145	WP	Dist1	MAR 18	PPTCT GGH – KAKINADA	926		2	2	5	5				
146	WP	Dist1	MAR 18	ICTC RMC KAKINADA	767		0	0	36	36				
147	WP	Dist1	MAR 18	ICTC TBIDH GGH KAKINADA	570		0	0	18	17				
148	NP	Dist4	MAR 18	Dr.R.N.Cooper	1518	1259	0	0	17	7				
149	NP	Dist4	MAR 18	Shatabdi Govandi	553	518	2	2	16					
150	NP	Dist4	MAR 18	KEM Hos.	3322	3544	1	1	42	27				
151	NP	Dist4	MAR 18	LTMGH	2009	1782	1	1		23				
152	NP	Dist4	MAR 18	J.J. Hos.	3042	2934	12	0	65	38				
153	NP	Dist4	MAR 18	Nair Hos.	1700	1652	0	0	32	25				
154	NP	Dist4	MAR 18	G.T. Hospital	742	740		0	10	1				
155	NP	Dist4	MAR 18	H. Bhagawati Hospital	1908	2233	1	1	23	3				
156	NP	Dist4	MAR 18	Siddhartha Nagar Hospital	591	555	0	0	6	8				
157	NP	Dist4	MAR 18	Rajawadi Hospital	1463	1332	1	3	8	5				
158	NP	Dist5	MAR 18	A F M C, PUNE VCTC (SDLIC2725000220983 )	1059	1048	0	0	18	8				
159	NP	Dist5	MAR 18	BHARATHI HOSPITAL & MEDICAL COLLEGE (SDLIC2725000201912 )	403	403	0	0	5	3				
160	NP	Dist5	MAR 18	BHOSARI HOSPITAL (SDLIC2725000202045)	521	513	0	0	6	6				
161	NP	Dist5	MAR 18	RIMC PLINE VCTC	1030	1030	0	0	64	50				
167	NP	Dist5	MAR 18	CHEST & DISTRICT HOSPITAL AUNDH (SDLIC2725000201797.)	729	729	1	1	18	17				
163	NP	Dist5	MAR 18	DR.D.Y.PATIL MEDICAL COLLEGE & HOSPITAL (SDLIC27250002047.)	328	328	0	0	3	3				
164	NP	Dist5	MAR 18	MAHARASHTRA AROGYA MANDAL (SDLIC2725000201911)	441	441	0 0	0	9	6				
						· · -	-	-	-	-				

Sno	Province	District	Month	Name of the Facility	Total tested during the month	Total clients received post- test and results	Total ANC positives during the month	Total ANC Positives linked to ART during the month	Total Non ANC positives during the month	Total Non ANC clients linked to ART during the month	Total Positives (ANC+Non ANC)	Positivity Rate (%)	Total Positive clients linked to ART(ANC+Non ANC)	Linkage loss during the month
165	NP	Dist5	MAR 18	MIMER MEDICAL COLLEGE (SDLIC2725000202046 )	787	442	0	0	1	1	ANO		ANO)	
166	NP	Dist5	MAR 18	RAJIV GANDHI HOSPITAL (KHANSAHEB PESTANJI MATERNITY HOME) (SDLIC2725000216075 )	53	53	0	0	0	0				
167	NP	Dist5	MAR 18	RH MANCHAR (SDI IC2725000202042 )	498	508	0	0	3	3				
168	NP	Dist5	MAR 18	SDH INDAPUB (SDI IC2725000220633.)	185	185	34	0	2	2				
169	NP	Dist5	MAR 18	SILVER ILIBILEE GOVT RUBAL HOSPITAL BARAMATI (SDLIC2725000216081.)	271	271	0	0	12	12				
170	NP	Dist5	MAR 18	SMT KASHIBAI NAVALE HOSPITAL PPP (SDLIC2725000222763.)	1498	1551	1	1	185	14				
171	NP	Dist5	MAR 18	SONAWANE HOSPITAL PMC (SDLIC2725000201798.)	463	450	1	1	6	3				
171	ND	Dist5	MAD 19		990	994	0	0	30	29				
172	ND	Dist	MAD 10		330	271	0	0	55	50				-
173	NP	Disto	MAD 10	V 3 GENERAL HUSPHALT (SDEICZ/210014146/7)	146	146	0	0	5	5				
174	NP	Disto	MAD 10		140	140	0	0	5	3				
175	NP ND	Disto	IVIAN 10		140	140	0	0	25	3				
1/6	NP ND	Disto	IVIAR 18	RAJIV GANDHI MEDICAL COLLEGE, KALVA (SDLICZ/Z1001414883)	1045	1019	0	0	25	24				
177	NP	Dist6	MAR 18	SUTIKAGRUH HUSPITAL (SULICZ/21001421955)	462	452	1	1	4	4				
178	NP	Dist6	MAR 18	RUKMINIBAI HOSPITAL (SDLIC2/21001416388 )	667	663	0	0	18	18				
179	NP	Dist6	MAR 18	SHASHTRINAGAR GENERAL HOSPITAL (SDLIC2721001416392 )	397	398		0	7	9				
180	NP	Dist6	MAR 18	CENTRAL HOSPITAL ULHASNAGAR (SDLIC2721001401795 )	762	807	1	1	35	34				
181	NP	Dist6	MAR 18	BALKUM (CHHAYA HOSPITAL, ABERNATH ) (SDLIC2721001416961 )	705	705	0	0	6	6				
182	NP	Dist6	MAR 18	HEALTH CENTRE BHAYANDAR-W (SDLIC2721001420583 )	540	538	0	0	10	10				
183	NP	Dist6	MAR 18	HEALTH CENTRE MIRA ROAD (SDLIC2721001421947 )	576	544	0	0	9	9				
184	NP	Dist6	MAR 18	I G M HOSPITAL, BHIWANDI (SDLIC2721001416389 )	622	619	1	1	30	30				
185	NP	Dist6	MAR 18	NAVI MUMBAI (GEN.HOSPITAL),VASHI (SDLIC2721001401796 )	1239	1332	0	0	29	28				
186	NP	Dist6	MAR 18	RAJMATA JIJAU MATA VAL RUGNALAYA,AIROLI (SDLIC2721001416959 )	0	0	2	2	4	4				
187	NP	Dist6	MAR 18	PAD. DR D Y PATIL MEDICAL COLLEGE & HOSPITAL, NERUL, NAVI MUMBAI (SDLIC2721001414879 )	589	586	3	2	4	4				
188	WP	Dist1	APR 18	ICTC AMALAPURAM	708		0	0	14	13				
189	WP	Dist1	APR 18	ICTC RAMACHANDRAPURM	525		1	1	11	8				
190	WP	Dist1	APR 18	ICTC TUNI	658		0	0	8	8				
191	WP	Dist1	APR 18	ICTC PEDDAPURAM	447		0	0	16	15				
192	WP	Dist1	APR 18	ICTC DH - RAJAMEHINDRAVARAM	936		0	0	803	62				
193	WP	Dist1	APR 18	ICTC RMC KAKINADA	690		0	0	46	42				
194	WP	Dist1	APR 18	ICTC TRIDH GGH KAKINADA	602		Ű	0	37	37				
105	W/P	Dist3	APR 18	SA AH - BAPATIA (P&V)	813		1	1	13	13				
195	W/P	Dist3	APR 18	SA AH - NARSARAOPET (P&V)	1368		1	1	15	29				
196		Dist3	APR 10		6712		1	1	10	10				
197	VVP	Dist3	APR 10		0715		1	1	10	10				
198	WP	Dist3	APR 18		379		0	0	12	6				
199	WP	DISL3	APR 18	SA CHC - VINUKUNDA (I)	280		2	2	13	10				
200	WP	DISL3	APR 18	SA DH - TENALI (V)	528		0	0	39	35				
201	WP	Dist3	APR 18	SA GGH - GUNTUR (V)	1582		0	0	86	53				
202	WP	Dist3	APR 18	SA IBH - GUNIUR (I)	300		0	0	18	15				
203	WP	Dist2	APR 18	SA AH - GUDIVADA (P&V)	1147		0	0	18					
204	WP	Dist2	APR 18	SA AH - NUZIVEEDU (P&V)	909		1	1	12	12				
205	WP	Dist2	APR 18	SA CHC - AVANIGADDA (V)	275		0	0	5	4				
206	WP	Dist2	APR 18	SA CHC - GANNAVARAM (V)	344		0	2	4	4				
207	WP	Dist2	APR 18	SA CHC - JAGGAIAHPET-OLD (V)	335		2	2	10	10				
208	WP	Dist2	APR 18	SA CHC - KAIKALUR (V)	232		0	0	3	3				
209	WP	Dist2	APR 18	SA CHC - MYLAVARAM-OLD (V)	283		1	1	12	11				
210	WP	Dist2	APR 18	SA CHC - NANDIGAMA-OLD (V)	441		0	0	16	16				
211	WP	Dist2	APR 18	SA CHC - RAJIVNAGAR-U (V)	56		0	0	21	19				
212	WP	Dist2	APR 18	SA CHC - TIRUVURU-OLD (V)	201		0	0	7	7				
213	WP	Dist2	APR 18	SA CHC - VUYYURU (V)	213		12	0	2	2				
214	WP	Dist2	APR 18	SA DH - MACHILIPATNAM (P)	453		0	0	0	0				
215	WP	Dist2	APR 18	SA DH - MACHILIPATNAM (V)	776	İ	0	0	20	19				
216	WP	Dist2	APR 18	SA GGH - SIDDHARTHA MEDICAL COLLEGE (V)	1744	1	0	0	85	78				
217	WP	Dist2	APR 18	SA GGH - VUAYAWADA (P)	119	1	2	2	7	7				
218	WP	Dist2	APR 18	SA PHC - KANKIPADU (Previously NP)	526		0	0		3				
210	W/P	Dist2	APR 19	SA PVT-MC - PINNAMANENI INSTITUTE OF MEDICAL SCIENCES (I)	472		0	0	5	3				
219	***	DISCZ	7111 10	on the second se	1 7/2	1			5	5				

Sno	Province	District	Month	Name of the Facility	Total tested during the month	Total clients received post- test and results	Total ANC positives during the month	Total ANC Positives linked to ART during the month	Total Non ANC positives during the month	Total Non ANC clients linked to ART during the month	Total Positives (ANC+Non ANC)	Positivity Rate (%)	Total Positive clients linked to ART(ANC+Non ANC)	Linkage loss during the month
220	WP	Dist2	APR 18	SA UPHC - RUDRAPAKA (I)	216		0	0	24	0				
221	WP	Dist1	MAY 18	ICTC AMALAPURAM	843		0	0	12	11				
222	WP	Dist1	MAY 18	ICTC RAMACHANDRAPURM	582		0	0	11	9				
223	WP	Dist1	MAY 18	ICTC TUNI	695		4	4						
224	WP	Dist1	MAY 18	ICTC PEDDAPURAM	573			0	16	11				
225	WP	Dist1	MAY 18	ICTC DH – RAJAMEHINDRAVARAM	1717		2	2	90	66				
226	WP	Dist1	MAY 18	ICTC RMC KAKINADA	733				43	39				
227	WP	Dist1	MAY 18	ICTC TBIDH GGH KAKINADA	473				29	24				
228	WP	Dist2	MAY 18	SA AH - GUDIVADA (P&V)	1159		0	0	11	10				
229	WP	Dist2	MAY 18	SA AH - NUZIVEEDU (P&V)	871		2	2	25	23				
230	WP	Dist2	MAY 18	SA CHC - NANDIGAMA-OLD (V)	420		0	0	12	12				
231	WP	Dist2	MAY 18	SA CHC - RAJIVNAGAR-U (V)	375		0	0	8	8				
232	WP	Dist2	MAY 18	SA DH - MACHILIPATNAM (V)	808		0	0	36	34				
233	WP	Dist2	MAY 18	SA GGH - SIDDHARTHA MEDICAL COLLEGE (V)	1821		0	0	89	75				
234	WP	Dist2	MAY 18	SA GGH - VIJAYAWADA (P)	1288		3	3	12	12				
235	WP	Dist3	MAY 18	SA AH - BAPATLA (P&V)	828		2	2	4	4				
236	WP	Dist3	MAY 18	SA AH - NARSARAOPET (P&V)	1559		0	0	32	32				
237	WP	Dist3	MAY 18	SA CHC - CHILAKALURIPETA (V)	770		0	3	9	9				
238	WP	Dist3	MAY 18	SA CHC - MACHERLA (V)	359		0	0	44	4				
239	WP	Dist3	MAY 18	SA CHC - VINUKONDA (I)	289		23	0	1	1				
240	WP	Dist3	MAY 18	SA DH - TENALI (V)	6410		0	0	50	53				
241	WP	Dist3	MAY 18	SA GGH - GUNTUR (V)	1925		0	0	77	50				
242	WP	Dist3	MAY 18	SA TBH - GUNTUR (I)	386		0	0	23	21				



## **Basic Measures**

- Sum/Total
- Sub-totals
- Percentage
- Distribution
- Min-Max
- Cut-off based
- Top & bottom
- Quartiles
- Average/ Mean
- Median

# Computing new variables

- Recoding text to numeric variables
- Converting continuous to categorical variables
- Computing new variables age groups, percentages

## Indicator Estimation

- Number
- Percentage
- Rate
- Ratio

#### **Useful Excel Functions**

- Sort & Filter
- Replace
- Go to
- Pivot tables
- V-look Up
- Conditional Formatting
- Remove Duplicates
- Sum, Average, Median, Quartiles

#### Levels of Measurement

#### Nominal Variable

- Categorical/ Discrete data
- Distinct groups
- No ordering/ Order is not meaningful
- Gender, Occupation, Marital status, Colour, Yes/No,

#### Ordinal Variable

- Categorical/ Discrete data
- Distinct groups
- Ordering meaningful
- Distance/gap between two level in the order is not meaningful
- Stages of disease, Severity of disease, Mild-Moderate-Severe, Low-Medium-High, Education, Ranks

#### Interval Variable

- Numeric/ Continuous data
- Distance/gap between two levels can be measured & is meaningful; Subtraction is possible between values
- No absolute zero/ No Common Reference Point; Zero is not meaningful
- Temperature, Age groups, Distance b/w two points

#### Ratio Variable

- Numeric/ Continuous data
- Absolute zero is meaningful
- Levels can be expressed as ratio or number of times of one another
- Height, Weight, Prevalence rate

# SUMMARY MEASURES

Level of Measurement	Туре	Measures of Central Tendency	Measures of Dispersion	Other Tests
Nominal	Categorical/ Discrete	Count/ Frequency, Mode	Distribution	Chi-square, T- test
Ordinal	Categorical/ Discrete	+ Median	+ Percentiles, Quartiles, Inter-quartile range	Chi-square, T- test
Interval	Continuous/ Numeric	+ Mean	+ Standard Deviation	One-way ANOVA
Ratio	Continuous/ Numeric	+ Ratio	+ Coefficient of Variation	Regression

# Exercise 6

- Estimate the basic measures from the NSACP dataset you have cleaned in DQA exercise
- Convert continuous to categorical variable
- Estimate four indicators



## Overview

- Understanding Excel Layout
- Sheet Visualisation Functions
- Dealing with Rows & Columns
- Cell Formatting
- Cell Navigation
- Paste Special
- Dataset Functions
- Analytic Functions
- Basic Formulae
- Advanced Formulae
- Advanced Operations

# Understanding Excel Layout

- Three views Normal, Print, Pagebreak
- Sheets
- Ribbon & Tabs

# **SHEET VISUALISATION FUNCTIONS**

- Adding headers
- Colour cells
- Freeze panes
- Naming & Renaming sheets

# Dealing with Rows & Columns

- Inserting Rows & Columns
- Deleting Rows & Columns
- Copying & Pasting Rows & Columns
- Adjusting Column Widths & Row Heights

# **Cell Formatting**

- Understanding Cell formats General, Number, Date
- Wrap Text
- Adding borders

# **Cell Navigation**

- Using arrow buttons to navigate
- Using Ctrl & Shift buttons with arrow buttons to select cells
- Go To Dialogue Box & Keyboard Shortcut (Ctrl+G)

## **PASTE SPECIAL**

- Paste values
- Paste formulas
- Paste column widths
- Paste formats
- Transpose
- Link Cells

### **DATASET FUNCTIONS**

- Sort
- ▶ Filter
- Conditional Formatting Blanks, Duplicates, Text/ Numeric Condition
- Remove Duplicates
- Find & Replace
- Increasing & Decreasing Decimal values

## **Analytic Functions**

- Writing Formulas
- Using Brackets & Commas
- Locking Formulas
- Relative Reference

## **Basic Formulae**

- Count, CountA, Countblank, Countif, Countifs
- Sum, Sumif, Sumifs
- Other Fundamental operations Subtract, Multiply, Divide
- Percentage & Distribution
- Min-Max
- Cut-off based
- Top & bottom
- Quartiles
- Average/ Mean
- Median

# **Advanced Formulae**

- ▶ If function
- Nested If functions And, Or
- Left, Right, Mid functions
- Join function

# **Advanced Operations**

- Data migration using Vlookup
- Analysis using Pivot Tables
- Creating Graphs
- Data Validation



Overview General Rules	
Number Functions	Logical Functions
▶ Count	▶ And
▶ Sum	► Or
► Min-Max	Special Functions
Quartiles	▶ If
Mean/Average	Vlookup
Median	
Text Functions	
► Join	

Left, Right & Mid

### **GENERAL RULES OF WRITING FORMULAE**

- Formula starts with '=' followed by the formula name followed by opening bracket
- Formula can be typed in capital or small letters. Excel automatically converts it into capitals.
- Every formula has a defined syntax/ structure/ components, that should be followed
- Parts of syntax are separated by commas
- Formula always ends with bracket
- No. of opening brackets and closing brackets should be equal
- Numbers in formulae can be entered directly
- Text in formulae is entered within double quotation marks "OK"
- Cell reference is entered in Alphanumeric code A3, B5
- Cell range is entered with a colon between two cell references A3:B5
- There should not be any blank spaces in formulae

#### **COUNT FUNCTIONS**

#### 5 types of Count Functions

- Count: Counts only numbers
  - ▶ Formula: =Count(A2:A30)
- CountA: Counts all non-blank cells numbers & text
  - ▶ Formula: =CountA(A2:A30)
- Countblank: Counts blank cells
  - Formula: =Countblank(A2:A30)
- Countif: Counts the cells that satisfy a condition
  - ► Formula: =Countif(Range,Criteria). Eg: =Countif(A2:A30,">5")
- Countifs: Counts the cases that satisfy more than one condition in the same or different cells
  - ► Formula: =Countifs(Range1,Criteria1,Range2,Criteria2.....).
  - ▶ Eg: =Countifs(A2:A30,">5", A2:A30,"<10")

#### **SUM FUNCTIONS**

#### 3 types of Sum Functions

- Count: Gives total of the numbers in the range
  - Formula: =Sum(A2:A30)
- Sumif: Adds the cells that satisfy a condition
  - ▶ Formula: =Sumif(Range,Criteria,Sum\_Range).
  - ▶ Eg: =Sumif(A2:A30,">5",B2:B30)
- Sumifs: Adds the cases that satisfy more than one condition in the same or different cells
  - ► Formula:
    - =Sumifs(Sum\_range,Criteria\_Range1,Criteria1,Criteria\_Range2,Criteria2.....)
  - Eg: =Sumifs(A2:A30,A2:A30,">5", A2:A30,"<10")</p>

#### Min, Max, Quartile & Percentile Functions

- Min: Gives the minimum number in the range
  - ► Formula: =Min(A2:A30)
- Max: Gives the maximum number in the range
  - Formula: = Max(A2:A30)
- Quartile: Gives the number at the defined quartile in the range
  - Formula: =Quartile(Range,Quart); Eg: =Quartile(A2:A30,2)
  - ▶ Quart values: 0-Min, 1-1<sup>st</sup> quartile, 2 median, 3 3<sup>rd</sup> quartile, 4 Max)
- Percentile: Gives the number at the defined percentile in the range. Percentile should be given in decimal value (E.g. 50<sup>th</sup> percentile – 0.5, 90<sup>th</sup> percentile – 0.9)
  - Formula: =Percentile(Range,k); Eg: =Percentile(A2:A30,0.5)

#### MEAN/ AVERAGE & MEDIAN

- Average: Gives the average/ mean of the in the range
  - Formula: =Mean(A2:A30)
- Median: Gives the middle number/ 50<sup>th</sup> percentile in the range
  - ▶ Formula: = Median(A2:A30)

#### **TEXT FUNCTIONS**

- Join Function: Physically joins the characters and creates a text string
  - ► Formula: =A2&" is OK"
- Left: Extracts the defined number of characters from the left of a text string
  - Formula: =Left(Cell Reference/ Text,No. of characters)
  - E.g. =Left(A2,2); =Left("Sri Lanka",3)
- Right: Extracts the defined number of characters from the right of a text string
  - Formula: =Right(Cell Reference/ Text,No. of characters)
  - E.g. = Right(A2,2); =Right("Sri Lanka",3)
- Mid: Extracts the defined number of characters from the Middle of a text string
  - Formula: =Mid(Cell Reference/ Text,Start\_Num,No. of characters)
  - E.g. = Mid(A2,2,3); =Left("Sri Lanka",4,5)

#### **IF** Function

- One of the most powerful & most productive function in MS Excel
- If Function: Returns a predefined result if a condition is fulfilled and another predefined result if the condition is not fulfilled
- Multiple If conditions, upto 32, can be nested inside one another, to evaluate a series of conditions and give results
- If function has three parts Condition, Value if True, Value if False.
- Condition can be specified using any other function in Excel
- Value if True and False can be a number, text (enclosed in double quotation marks), blank (double quotation marks without anything in between), a cell reference, a formula or another IF function.
- Used to convert one data into the other (Text to Numbers, Numbers to Categories, coding, etc.)
- ► Formula: =IF(Condition, Value if True, Value if False)
- E.g. =IF(A3>0,"OK","Error")

#### **Examples of IF Function**

- =IF(A3>0, "OK", "ERROR")
- =IF(A3=0,"",B3/A3)
- =IF(OR(A3="NORTH",A3="SOUTH"),"NS","EW")
- IF(A3<15,"<15",IF(AND(A3>=15,A3<=24),"15-24","25 & ABOVE")))</pre>

#### **VLOOKUP** Function

- Used to bring data of same cases from a different dataset or different part of the same dataset.
- Useful in combining two different datasets
- Need at least one common field Lookup Value
- Lookup value should be in the first column of the lookup range
- Column number is relative to the first column of the lookup range; Not the column number in the worksheet
- Match types: True Approximate Match; False Exact Match
- =VLOOKUP(LOOKUP VALUE,LOOKUP RANGE,COLUMN NO.,MATCH TYPE)
- =VLOOKUP(A3,SHEET1!A2:P100,4,FALSE)



# OUTLINE

- Epidemic Analysis
- Progress Analysis
- Performance Analysis
- Cascade & Cohort Analysis

### **Epidemic Analysis**

- HIV & STD Positivity rates
  - Levels
  - Trends
  - Differentials by demographic & risk characteristics
- Size of beneficiaries
- Profile of beneficiaries
  - STD Clinic attendees
  - PLHIV
- Vulnerabilities/ Risk Behaviours among KP
  - > Typology, Partner volume, condom use, STI uptake

#### **Progress Analysis**

- Progress against targets
  - Levels
  - Trends
  - Differentials by province/ district
  - Differentials by other characteristics
- Gaps in Achievement
- Different Denominators

## **Performance Analysis**

- Performance Indicators
  - Levels
  - Trends
  - Differentials by province/ district
  - Differentials by other characteristics
- Best & worst performing units
- Performance Quartiles
- Performance Quadrants

## Cascade & Cohort Analysis

- Testing and treatment cascade
- Prevention cascade
- PLHIV cohort analysis
  - Incidence of Opportunistic infections
  - Mortality rates
  - Survival analysis
- Positive Pregnant Women & Exposed Baby cohort analysis



### **DATA TRIANGULATION**

- Concept & Principles
- What it does?
- When to use?
- Advantages
- Questions for Data Triangulation
- Data Sources
- Data Quality Assessments
- Methods of analysis
- Hypothesis-building/ Epidemiological Framework

#### DATA TRIANGULATION

Data Triangulation is an Analytical Approach that synthesizes data from multiple sources, to improve the understanding of a public health issue and guide programmatic decisionmaking to address the issue.

- By putting different bits of information from different sources into a meaningful framework, it explains and improves the understanding of HIV/AIDS scenario in the district.
- By providing answers to vital programme questions, it helps in taking effective decisions for planning and implementation of HIV prevention & control efforts

#### **BASIC PREMISE**

- Works on already available evidence; No primary data collection involved
- Exploratory in nature (arrive at plausible hypothesis); May not always give definitive conclusions with statistical rigour
- Better suited for broader issues at population/ community/ region level; May not work well for micro-issues
- Always guided by a specific question to be answered
- Need a basic understanding of epidemiological considerations
- Element of judgment is critical

#### Integrates diverse information

- Info. on diverse subjects/ thematic areas
- Data of different types (quantitative, qualitative, anecdotal, referential, experiential, etc.)
- Data generated through different methodologies (survey, evaluation, rapid assessment, OR etc.)
- Data generated for different purposes (Surveillance, monitoring, academic research etc.)
- Data generated by different sources (health, non-health)
- Data applicable to different populations/ regions
- Data with different time reference

#### Basic principle of data triangulation

To analyse and interpret a dataset in the light of information emerging from other datasets so that the synthesis offers a better

understanding of the issues than what will be inferred from a single

dataset





#### When to use Triangulation Approaches?

- When data is plentiful but dissimilar
- When data is scanty
- When the quality of data is not optimal
- When quick answers are needed
- ▶ When there are no resources for rigorous research
- When broad policy / strategy level decisions are to be made
- Before planning any major data collection activities

Research analysis	Triangulation analysis
<ul> <li>Focus on statistical analysis</li> <li>Designed to provide data that can be generalized</li> <li>Variables from a single dataset</li> <li>Focus on internal validity: "Did A cause B to change among group C?"</li> <li>Emphasis on generating the highest scientific rigor of data for interpretation</li> <li>Long delay between data collection and presentation of results</li> </ul>	<ul> <li>May or may not use statistics.Use of statistical analysis will depend on available data</li> <li>Variables from multiple datasets</li> <li>Focus on external validity: "Can observed effects in group C be attributed the larger population as well?"         Emphasis on the "best possible" interpretation of existing data for policy and programme decision making         Quick turnaround between secondary data capture and presentation of results       </li> </ul>
### META-ANALYSIS AND TRIANGULATION

- Meta-analysis: combines rigorous scientific data of similar quality and design to conduct statistical analyses.
- Triangulation: seeks to make use of data from diverse sources and study designs, and incorporates judgments, findings and interpretations on each data source's limitations.

#### Advantages of Triangulation

- Can be easily employed in programme settings
- Can be applied to virtually every programme question, but more appropriate while seeking answers to broad questions on programme and policy
- Helps in identifying more data sources
- Helps in bringing out information gaps
- Provides new insights and helps in generating new hypotheses
- Less resource intensive (relatively)
- Makes best possible use of available evidence; a good alternative to new data generation

#### **QUESTIONS FOR TRIANGULATION**

- Importance: The question should address a current and pertinent issue. The answer to that question should significantly influence your understanding of and response to the HIV/AIDS situation.
- <u>Actionability:</u> The results of the data-triangulation process should be useful to make improvements in HIV/AIDS programme.
- Answerability (Data availability & Feasibility): The data required for answering the question should be available in usable format. At least, it should be possible to extract it from existing reports or documents. Primary data collection of any form is not a part of data triangulation.
- Appropriateness: Finally, it is important to ascertain whether triangulation is the appropriate method to use to answer the question. Could the question be better answered by any other research methods of analysis, an expert panel, or another type of study?



Start with a broad set of questions. Revisit and Refine the questions At every step

#### **Steps in Data Triangulation**

- Step 1: Understanding Thematic Areas & Questions for Triangulation
- Step 2: Review of Data Sources and Assessment of Data Availability in the District
- Step 3: Decision on Questions to be answered for the district
- Step 4: Compilation of Secondary Data
- Step 5: Data Quality Assessments
- Step 6: Data Validation, Adjustments & Filling Data Gaps
- Step 7: Preparation of Data Tables with clean data for analysis
- Step 8: Data Analysis, Interpretation and Inferences; Describe Thematic Areas
- Step 9: Data Triangulation (Hypotheses Building; Answer Triangulation Questions)
- Step 10: Preparation of Draft Reports
- Step 11: Discussion & Consultation with stakeholders & Local experts on draft reports
- Step 12: Finalisation of conclusions and recommendations
- REVISIT QUESTIONS AT EVERY STEP

# **Epidemiological Considerations**



#### DATA QUALITY ASSESSMENT

- 1. Review of methodology, design, data collection processes, possible biases
- 2. Validity and Reliability
- 3. Completeness
- 4. Correctness
- 5. Consistency

Followed by cleaning, adjustments and validation

# FRAMEWORK FOR DATA TRIANGULATION

Components of Data Triangulation	What it Does?	Guiding Elements	Action To Do	Output
Descriptive Analysis	Describes (What? Who? When? Where?)	Themes	Analyse Data & Describe the Themes	Descriptive Section of Report
Triangulation	Explains (How? Why?)	Questions	Triangulate Data & Answer the Questions	Synthesis Section of Report





### Triangulation of info on same data element from different sources

- HIV Prevalence among GP: Compare HSS ANC data with PPTCT data.
- HIV Prevalence among HRG: Compare HSS HRG data with ICTC General Clients data and TI NGO data (wherever available)
- STI Prevalence and Trends: Compare data on STIs from STI Clinics with that collected in behavioural surveys, STI treatment data from TI projects and VDRL positivity rates from ANC HSS.
- Size of HRG: Compare mapping data of HRG with TI Needs Assessment and coverage data of HRG, both refer to the size of HRG.
- Risk Behaviours: Compare behavioural survey data from BSS/ IBBA with the behavioural data from TI Projects, wherever available.
- Risk Behaviours: Compare routes of transmission reported from ICTC with that from ART registration data
- Size of PLHA: PLHA density from ART registration data with that from Positive People Networks & ICTC detections

### Triangulation of information on different data elements

- Triangulation in time plane
- Triangulation in geographical plane
  - It establishes the relationships between different data elements.
  - It helps in identifying the blocks/ talukas where the HIV epidemic or its vulnerabilities are high and hence require priority attention under the programme.
  - Thirdly, it shows the gaps in the programme response in terms of geographic coverage i.e. talukas/ blocks where there is need but there are no interventions.

## Triangulating time data

- Timing of HIV trends among different groups
- HIV trends among different groups and see if the information available on HRG size at different points in time has any relation with HIV trends.
- HIV trends among different groups and see if the trends in STIs correlate with HIV trends.
- HIV trends among different groups and see if the data on risk behaviours and changes in risk behaviours correlate with HIV trends.
- HIV trends among different groups and mark the timing of start of interventions in the district (preventive interventions, counseling & testing and ART) and see if any message emerges out.
- If any information is available on any other data element such as migration, establishing industries or factories, starting of Transshipment locations or truck halt points, etc., they should be captured in the time trend analysis and examined for any meaningful insights.

### Triangulating geographic patterns

- HIV Positivity among pregnant women from PPTCT data
- ▶ HIV Positivity among general clients from ICTC data.
- High Risk Group mapping data
- ART Registration data
- Presence of programme interventions (TIs, ICTCs, PPTCT centres, ART centres, LACs, CCCs, PLHA networks)
- Sources of out-migration & Destinations of in-migration

#### Building an Epidemiological Framework...

- Starting with contextual factors and background vulnerabilities
- Understanding the presence of risk groups, their distribution and profiles
- Understanding the risk behaviour patterns among HRG and GP
- Understanding the important vulnerabilities such as STIs and Migration
- Describing how all these interact and cause HIV epidemic
- Understanding the HIV/AIDS scenario and patterns in different risk groups and general population
- Appreciation of possible future vulnerabilities and direction of epidemic

Making a story out of it...











# Mandal Wise VCTC Positivity, Guntur Dt., AP





				PPTCT	Size/TI	PLHA	
VCTC - 2008	Men	Women	Total	2008	FSW	MSM	registered
Hanumakonda	Sam	ple very s	mall	1.70	270		657
Warangal	24.8	29.6	26.7	0.68	117	1655	1234
Wardhannapet	3.7	2.5	3.3	0.66	262	268	101
Janagaon	9.0	9.0	9.0	0.59	752	542	164
Mahabubabad	4.4	6.7	5.3	0.52	850	1147	155
Narsampet	3.6	5.6	4.2	0.41	260	211	
Parkal	6.6	4.5	5.6	0.39	522	906	135
Mullugu	2.5	6.0	3.5	0.36			- OF
Cheryal	4.3	4.1	4.2	0.18			orogic
Station-Ghanpur	5.6	3.0	4.4	0.00		oti	on P 109
Eturunagaram	1.4	3.2	1.9	0.00		orever	
Gudur	2.2	1.7	1.9	0.00	itiOt	eP	
Chityal	1.1	1.3	1.2	0.00	is to inin		
Warangal	6.6	8.0	7.0	Nee	ds		
Bhopalpally					745	453	
Thorur					271	273	111
Keesamudram	No infe				276	256	
Nekonda		mation	On Lui		254	266	
Hasanbharathy				prevalor	332		180
Dharmasagar				dien	<b>Ce</b> 51		147
Atmakur							121



# OUTLINE

- Scientific Abstracts
- Scientific Papers/ Articles
- Policy Briefs
- Handouts
- Posters
- Reports
- Presentation

# **Scientific Abstracts**

- Short and Summarised; 250-300 words
- For conference and journal publications
- Features
  - Objectives
  - Methods
  - Results
  - Conclusion
  - Key words

# Scientific Paper/ Article

- Published in Peer Reviewed Journals
- Detailed version of the data sources, methods and results
- Sections
  - Introduction
  - Methods
  - Results
  - Discussion The most important section; Author's contribution
  - References
  - Acknowledgements
  - Conflict of Interest statement
  - Tables & Figures

# **Policy Briefs**

- Targeted at Policy makers, senior level programme managers & experts
- Key take home messages to be very crisply worded and adequately highlighted; Should be at the top or on the cover page
- Extract of only the top-line, high priority findings
- Presented in a visually appealing manner for quick consumption; Graphs, Infographics, Pictures; Not in tables; Very limited text
- Very limited info on methods, if required
- Not more than 2 pages
- Sometimes done as wall chart
- Colourful and elegantly designed

# HANDOUTS/ FOLDERS

- Larger than policy briefs; Smaller than Papers; Much smaller than reports
- 4-8 pages; designed as a folder/ folder chart/ booklet
- Targeted at key stakeholders prog managers, communities, field staff, academicians, experts
- Easy to mail and circulate
- Presents the topline findings and other key findings
- Graphs & Tables
- Brief paragraphs or text on the intro/background/ methods may be included
- Conclusions/ Key take home messages highlighted
- Contact details

## Posters

- Large sized visual presentation of the methods and results
- Mostly use graphs, maps, infographics; Limited text
- > Targeted towards academic dissemination & conferences
- May also be used as standees during official meetings/ programme events
- Sections:
  - Introduction
  - Methods
  - Results Major section
  - Conclusions

# Presentations

- Using MS Power-point or Prezi
- Based on target audience
  - Detailed or Brief
  - Text or Graphics or both
  - Academic or Policy
  - High impact presentations vs Discussion-oriented
  - Interesting or boring
- Clean, less crowded slides (not more than 8 lines)
- Can embed videos and animation

# Reports

- More detailed and elaborate
- Complete details of the project, process, outcomes and discussion
- Sections and chapters
- > Detailed presentation of results in text, tables & graphics forms
- Foreword/ Preface, List of abbreviations, List of tables, List of figures
- Detailed references and acknowledgements
- Targeted at universal audience
- Final and complete consolidation of all the process and the results



# SUCCESSFUL CONTROL OF HIV Epidemic NEEds..

- Effectiveness in delivery of services at every level....targeting the right thing
- Understanding local scenario and factors responsible
- District specific action plans with identified focus areas
- Customised choice of programme interventions
- Contextualised emphasis on programme components
- EFFECTIVE & TIMELY USE OF DATA AT ALL LEVELS

## DATA USE - AN approach, a mindset

- Developing the approach of using data for decision-making and programme planning at all levels
- Inculcating, among programme managers, a habit of looking at data regularly
- Encouraging simple analytical methods that anyone can employ
- Emphasizing the importance of local knowledge and contextual understanding
- Capacity-building of state & district level institutes/personnel for sustainability

## A long term process...

- Trigger the interest in data analysis and use
- Make programme managers work on data of their own state/district & reflect upon the insights
- Expose them to real time examples that demonstrate the use of data for decision making & programme planning
- Develop guidelines & tools to assist them in data use
- Develop HR plans that sustain the interest & facilitate data use as an on-going process
- SYSTEM STRENGTHENING FOR DATA USE

















# Knowledge Management



# Four components of Knowledge Management

- Creation of Knowledge Data Analysis (Converting data into information & knowledge)
- Collection and storage of Knowledge from other sources Data Archiving
- Sharing of knowledge Dissemination and Communication of Knowledge Products to stakeholders
- Translation of Knowledge Data Use for programmatic action



## **Knowledge Collection and Archiving**

- While the former refers to creating knowledge from data generated under the programme, this component refers to
  - Collection and documentation of knowledge generated by other sources (communities, partners, academia, etc.)
  - Collection and synthesis of tacit or experiential knowledge from beneficiaries, service providers and other stakeholders
- Data Archives and Repositories
- Experience sharing mechanisms

# **Knowledge Sharing and Dissemination**

- Passive dissemination through website, knowledge hub, datacentre, e-libraries, newsletters, reports and publications, etc.
- Active dissemination through seminars/ workshops/ conferences, discussion forums, e-learning tools etc.
- Addressing needs of wider audience lay man, programme manages, policy makers, academia & researchers

# Knowledge Translation (Ensuring Effective Data Use)

- Effective use of knowledge for programmatic decision-making and policy or strategy formulation
- Starts with understanding programme requirements for evidence and analysis and identifying key questions to be answered
- Generating/compiling/synthesising body of knowledge to answer the questions and communicating them to the programme
- Taking programmatic decisions based on the evidence

# Need for System Strengthening for Effective Knowledge Management

> Not as a one-time activity, but as a systematic approach of

- Shaping perceptions and promoting 'data use for decisionmaking' as a regular practice
- Standardisation of methods & tools
- Capacity building of prog. Managers & data teams
- Devising mechanisms for collection of knowledge from other sources
- Developing strategy for timely communication of analytical outputs
- Building institutional resource pools
- Promote scientific writing within the programme

#### TRANSLATING KNOWLEDGE INTO PRACTICE – A SIMPLE MODEL FOR PUBLIC HEALTH















# Ways to Present Data

- Text
- Tables
- Figures and illustrative graphs



# Text

Use text whenever there are small amounts of data to be summarized.

"During 2018, there is 60% increase in clinic attendance compare to last year"



#### Tables are arrangements of numbers or words in columns and rows that display data or relationships

#### Tables

#### STIs reported from STD clinics during 2018 by sex

Diserseis	Male		Fer	nale	Total	
Diagnosis	No.	%	No.	%	No.	%
Genital herpes	1268	26%	1813	32%	3081	29%
Non-gono. infections	886	18%	1956	35%	2842	27%
Genital warts	1416	29%	1156	21%	2572	25%
Syphilis*	552	11%	263	5%	815	8%
Gonorrhoea	206	4%	76	1%	282	3%
Trichomoniasis	11	0%	42	1%	53	1%
Other STIs	528	11%	322	6%	850	8%
Total STIs	4867	100%	5628	100%	10495	100%

\* All types of syphilis





# An example for a scientific paper

#### Table X

Proportion of Errors in Younger and Older Groups

Level of difficulty		Younger	r	Older			
	n	M (SD)	95% CI	n	M (SD)	95% CI	
Low	12	.05 (.08)	[.02, .11]	18	.14 (.15)	[.08, .22]	
Moderate	15	.05 (.07)	[.02, .10]	12	.17 (.15)	[.08, .28]	
High	16	.11 (.10)	[.07, .17]	14	.26 (.21)	[.15, .39]	

Note. CI = confidence interval.
































































# **Role of Biostatistician**

- Protocol development
- Study Implementation
  - Data Analysis
  - Data Management
- Report/Manuscript writing

## What is Biostatistics?

It is an application of Statistics to the analysis of biological and medical data.

It is the field of study concerned with the collection, organization, summarization and analysis and interpretation of data.

## Sources of Data

**Routinely Kept Records** 

Surveys

**Experiments** 

**External sources** 

# Population

It is the complete collection of data to be studied and it contains all the subjects and variables of Interest.

# SAMPLE

A sample is a subset of the population.



### Variable

A Variable is a charecteristic which differs from a person, place or things etc

### **QUANTITATIVE VARIABLE**

Measurements made on quantitative variables convey information regarding amount. Ex – height. Weight, age

#### **Qualitative Variable**

Measurements made on qualitative variables convey information regarding attribute. Ex – color of the eye

### MEASUREMENT AND MEASUREMENT SCALES

It is the assignment of numbers to objects or events according to a set of rules.

**Nominal Scale** 

**Ordinal Scale** 

**Interval Scale** 

**Ratio Scale** 

#### SAMPLING

It is the Process of collecting information from the sample.

## Types of Sampling

- Simple Random sampling
- Stratified Sampling
- Systematic Sampling
- Multistage Sampling
- Cluster sampling

## Types of statistics

#### **Descriptive Statistics**

It describes or summarize the parameters of the population.

#### **Inferential Statistics**

It is the procedure by which we reach a conclusion about a population on the basis of information from the sample.

## **Descriptive Statistics**

Tabulations

Diagrams

Graphs

Measures of central Tendency

Measures of dispersion

**Measures of Association** 

TA	Ы	LES

	Column heading	Column heading
Row identifier		
Row identifier		
Footnote: Details body of the table	of the abbreviati	ions given in the

The table should be Complete and Accurate.

No repetition with texts or figures.





# Other Diagrams and Graphs

Line graphs

**Box Plots** 

**Dot Plots** 

Histogram

**Bubble plots** 

**Pyramids** 

## **Basic Measures of Central Tendency**

- 1. Mean
- 2. Median
- 3. Mode
- 4. Geometric Mean
- 5. Harmonic Mean
- 6. Weighted Mean
- 7. Truncated mean

## **Basic Measures of Dispersion**

Range

Inter Quartile Range

Standard deviation

## **MEASURES OF ASSOCIATION**

Correlation (refers to a linear relationship between two variables or quantities)

Cross Tabulations (chi square tests)

# **Inferential Statistics**

Types of statistical Inference

Estimation

**Hypothesis Testing** 

## **Estimation Types**

**Point Estimation** 

Interval Estimation: Confidence Intervals

## Estimation

Objectives ——> Non – Comparative

Example,

Objective: To determine the prevalence of Rheumatic heart disease in children age 5 – 15, India.

Summary Measure: % of children diagnosed with RHD based on an echocardiogram screening.

#### Example

Population: 18, 000 children between 5 – 15 years of age Sample: 2, 400 children (Multistage random sample)

No. of children diagnosed with RHD: 72 Point Estimate: 72/2400 = 3%

Interval Estimation: 95% confidence Interval (2.2 – 3.8%)

Interpretation: We have 95% confidence that the prevalence of RHD is between 2.2% and 3.8%

## Hypothesis Testing

Objectives ——> Comparative

Example,

Objective: To determine the effectiveness of HPV vaccine in the prevention of the transmission.

Study Design: Placebo controlled RCT

Hypothesis: The new HPV vaccine is effective for the prevention of transmission.

Summary measure: % of people infected with HPV in a two year period.

## COMPONENTS AND STEPS

**Statistical Hypotheses** 

Level of significance

**Statistical Test** 

**Decision Rule** 

**Data Analysis Decision** 

Interpretation

## **STATISTICAL Hypotheses**

Null Hypothesis (H0): Hypothesis of no difference.

Alternate Hypothesis (H1): An alternate statement about the null hypothesis (one sided or two sided).

H0: μ1 = μ2

H1:  $\mu$ 1  $\neq$   $\mu$ 2 or H1:  $\mu$ 1  $\geq$   $\mu$ 2 or H1:  $\mu$ 1  $\leq$   $\mu$ 2

### Example

H0: % of people infected with HPV is the same in both the groups. (New vaccine is ineffective).

H1: % of people infected with HPV is lower among the people vaccinated with new HPV vaccine than with the placebo (New vaccine is effective)

# Level of Significance

The Level of significance (a) is the probability of rejecting the null hypothesis.

Possible action	True	False	
Fail to reject H0	Correct action (1-a)	Type II error (β)	
Reject H0	Type I error (a)	Correct action (1-β)	
(1-a) – Confidence Interval			
(1-β) – Power of	f the test		

## Decision based on the data

P value: Probability of obtaining the results as observed with the collected data under the assumption that H0 is true.

The P value depends on the statistical tests such as

Z –test t – test Chi square test Model based test etc.

## **Decision Rule**

If p-value < a, we do not reject H0

If p-value > a, we reject H0

### Example

H0: % of people infected with HPV is the same in both the groups. (New vaccine is ineffective).

H1: % of people infected with HPV is lower among the people vaccinated with new HPV vaccine than with the placebo (New vaccine is effective)

Level of significance (a) = 0.05

## Decision

% of people infected with HPV in a new vaccine group is 3%

% of People infected with HPV in a placebo group is 10%

Statistical test used Z –test (Test for proportions)

P-value - 0.012

Decision 0.012<0.05, we reject H0. The vaccine is effective.

## STATISTICAL TOOLS USED FOR DATA ANALYSIS

- Analysis of Variance
- Simple Linear Regression
- Multiple linear Regression
- ✤ Logistic Regression
- Chi square tests
- Non Parametric Tests
- Vital Statistics

#### Simple and Multiple Linear Regression

Regression analysis is widely used for prediction.

It is used to obtain the linear relationship between a dependent variable and one or more independent variables.

#### **Logistic Regression**

When the dependent variable is a dichotomous variable, we can make use of logistic regression.

#### Chi square Test

Types of Chi – Square Tests

- > Tests of goodness of fit
- > Tests of Independence
- Tests of homogeneity
- Fisher's Exact Test
- Relative Risk and Odds Ratio
- Survival Analysis

#### **Relative Risk**

It is the ratio of risk of developing a disease among subjects with the risk factor to the risk of developing the disease among subjects without the risk factor.

#### Example

Subjects with and without the risk factor who became cases of preterm labor.

Risk factor	Cases of preterm labor	Non cases of preterm labor	Total
Extreme exercising	22	216	238
Not exercising	18	199	217
Total	44	415	455

#### Relative Risk = (22/238)/(18/217) = 1.1

This data indicates that the risk of experiencing preterm labor when a woman exercises heavily is 1.1 times as great as it is among women who do not exercise at all.

#### **Odds Ratio**

The odds for success are the ratio of probability of success to the probability of failure.

#### Example

Subject classified according to obesity status and mother's smoking status during pregnancy

Smoking status	Cases (obese)	Non cases (non obese)	Total
Smoked	64	342	406
Not smoked	68	3496	3564
Total	132	3838	3970

Odds ratio = (64/342)/(68/3496) = 9.6

The obese children are 9.6 times as likely as non obese children to have had a mother who smoked throughout the pregnancy.

## SURVIVAL ANALYSIS

It is a branch of statistics which studies the amount of time that it takes before a particular event occurs.

The event may be death or time to failure or time to occurrence of an event.

# **Analytic Techniques**

- \* Life Table Analysis
- ✤ Distribution plots
- \* Kaplan Meier product limit estimator
- \* Regression Models

# Vital Statistics

Branch of statistics concerning the important events In human life such as births, deaths, marriages, and migrations.

**Death Rates and Ratios** 

**Measures of Fertility** 

**Measures of Morbidity**
## DEATH RATES

- \* Annual crude death rate
- Annual specific death rate
- Standardized death rate
- \* Maternal mortality rate
- ✤ Infant mortality rate
- Neonatal mortality rate
- ✤ Fetal death rate
- ✤ Perinatal mortality rate

### **MEASURES OF FERTILITY**

- Crude birth rate
- General fertility rate
- Age specific fertility rate
- Standardized fertility rate

#### **MEASURES OF MORDIDITY**

- Incidence Rate
- Prevalence Rate
- Case fatality Ratio
- Immaturity Ratio
- Secondary attack rate

#### **Incident Rate**

It is the ratio of total number of new cases of a specific disease during a year to the total population.

#### **PREVALENCE RATE**

It is the ratio of total number of cases, new or old at a point in time to the total population at that point in time.

This rate is especially useful in the study of chronic diseases, but it may also be computed for acute diseases.

#### **Case fatality Ratio**

This ratio is useful in determining how well the treatment Programme for a certain disease is succeeding.

It is the ratio of total number of deaths due to a disease to the total number of cases due to the disease.

## DATA MANAGEMENT

- ✤ Data Capture
- ✤ Data Transcription
- ✤ Data Transfer
- Data Entry
- Data Cleaning
- Storage of hard copies of CRFs
- ✤ Storage of electronic data
- Data coding
- Dataset creation
- Data backup and recovery



### INTRODUCTION - Why Use SPSS

- SPSS has been around since the late 1960s.
- SPSS is the statistical package most widely used by social scientists
- Of the major packages it seems to be the easiest with less coding
- One can use it with either a Windows point-and-click approach or through syntax (i.e., writing out of SPSS commands.). Each has its own advantages, and the user can switch between the approaches.

#### Syntax

- Originally, SPSS was written like a programming language. Users wrote SPSS syntax (often on a mainframe computer and even with key-punch cards) that performed the tasks they wanted.
- In SPSS-Windows, users can still use syntax by using the syntax editor.
- They would open the syntax window by clicking on File, dragging down to New, and choosing Syntax;
- Type the SPSS syntax that they want to run;
- Click on Run and drag down to All. (Alternatively, if they want to run only a few commands, they would highlight those commands, click on Run, and drag down to Selection.)



#### Choosing Appropriate Scales & measures

- There are many different ways of collecting data
- We need to measure output or performance on some objective criteria
- In choosing appropriate scales, need to aware of reliability and validity
- The reliability of a scale indicates how free it is from random error

### **Reliability and Validity**

- Reliability means consistency
- Nunnally (1978) recommends a minimum level of .7 Cronbach Alpha value
- Validity refers to the degree to which it measures what it is supposed to measure.
- ▶ Research design → Objectives and goals of study → Questionnaire construction → code book preparation → data management and mining → selecting appropriate statistical tools and techniques to explore data towards achieving research goals and objectives.

### Scales of Measurement

- The variables are of two types: quantitative and qualitative, and the measurement of variables or levels of measurement are of four types. They are nominal, ordinal, interval and ratio (N,O,I,R)
- Levels of measurement are very important and serve as a basis for which statistical tests are permitted on a given set of data and type of research.

### **Preparing a codebook**

- Preparing the codebook involves deciding (and documenting) how you will go about
- Defining and labeling each of the variable
- Assigning numbers to each of the possible responses
- > You should have a unique ID for each respondent or case
- Keep minimum on open ended questions

#### Rules for naming of variables

- Must begin with a letter (not a number)
- Cannot include period, blanks or other characteristics like (!, ?, \*)
- Cannot include words used as commands by SPSS (all, ne, eq, to, by, or, gt, and, not, with etc)
- Cannot exceed 64 characters (SPSS V12) or 8 characters for earlier versions of SPSS

## Opening an existing data file

- Like opening a word document, double click on SPSS .sav data file
- You will get an untitled blank spreadsheet like Microsoft Excel data sheet
- You can open an existing data file by using Open menu and then pick Data
- > The data file will open which has data view and variable view
- You can click on save to save your file

t var	SPSS 12.0 for Windows	1W	197 197	191
2 3 4 4 5 6 7 7 6 3 10 10 11 12 13 14 15 14 15 14 15 16 17 10 10 10 10 10 10 10 10 10 10	What would you like to do?   What would you like to do?			
23 27 28 38	DK. Cavel	-		

#### CREATING A DATA FILE AND ENTERING DATA

- Open an untitled data file
- Click on the 'Variable View" tab at the bottom of the Data Editor
- Variable view has 'Name', Type, width etc appear at the top of each column. Each row corresponds to a variable.
- Let us say you type "Age" under variable name, the other values appear as they are
- set by default. You can change all the values as per your requirements.

Name	Type	Width	Decimals	Label	Values	Missing	Calumns	Align	Measure	
Age	Numeric	0	2		None	None	0	Right	Scole	
						1.1.1.1.1				
1										
1		_			_		_	_		
-		-		-			_	_		
		-						_		
						_	-			
		-					-			
		_								
		_		-	_	_				
		-		-			_	_		
		-		-				_		
		-					-	_		
1				1						
		-	-				-			
1		_								
		-								
		-	-	-				_		
		-		-				_		
		-				_	_			
		-			_	_				
-		-			-	_	-			
4		_					_			
		-					-	_		
		-		1	1		-			
		-				-	-	_		
		-								
		-				-	-	_		
		-	-		-		-			
1										

## VARIABLE NAMES

- First you need to decide on names for each of your variables
- ▶ Follow the rules involved in creating variable names
- Our tip: Avoid using symbols in variable names.....that way you won't need to remember which ones are okay and which one are not
- Specify correct variable type, column width, number of decimal points.

#### Defining variables and value labels

- Variable LABELS and Value Labels can be used to make output easier to understand.
- A variable label is used to descriptively label the variables
- Its use makes the output easier to read and can be very useful if the output is used over a long time period. Each label cannot exceed 60 characters.

### Value Labels

- A Value label is used to descriptively label the values of a variable.
- **Example of a value label:**
- Lets do a value label for a variable named community
  - A value of 1 refer to Asian, 2 refer to African American, 3 means Hispanic and so one
- For Age Category let us say 1 = 0-10 years; 2 = 11 -20; 3=20-45; 4=45-65 etc

### Missing data

- The presence of missing data is very common in any kind of research. A research will come across dead mice, sick children, non-compliant household head, unfilled forms, lost samples and so on.
- We can't really do anything about it other than planning for it while creating the data file
- You can enter a special code like '8 for refusal' or '9=don't know'; or you can leave the item blank
- A good thing about SPSS is that any blank or period (".") is considered missing data unlike some other software programs which consider a blank as zero.

### Changing the SPSS option

- It is always a good idea to check the SPSS options. Especially when you are working in a computer lab, students will change these options which can dramatically influence how the program appears. You can reset the option
- To reset the option, click on edit menu and select Option. You will see this option window

Data	Currency	Scripts
General Viewer Draft View	ver Output Labels Charl	s Interactive Pivot Tables
Initial Output State	Title Font	
Contents are initially: Contents are initiall	Text Output Page Size Width: Standard (80 characters) Wide (132 characters) Custom: 80 Text Output Font	Length: C Standard (59 lines) Infinite C Custom: 59
Display commands in the log	Courier New 🚽 10	B I U     ■

### DATA ENTRY USING EXCEL

- You can create your data file in Excel and import it to SPSS
- Warning: Remember, Excel can deal with 256 columns of data and if the data is likely to be larger, it is probably easier to set it up in SPSS
- Once you import the data in SPSS, you have to do variable and value labeling. The variable name must conform to the SPSS rules for naming variables.
- Make sure you save the file with .sav extension

# Screening and cleaning the data

- 1. Checking for errors: Check each of your variable by doing a descriptive statistics (frequencies)
- 2. Check minimum and maximum values and declare missing values appropriately



### Key Learnings

Hands-on Training on

- Understanding datasets, components, structure & database management principles
- Variables & Indicators Types and how to manage
- Data Quality Assessment & Adjustments using Excel
- Data Management using SPSS
- Exposure to Cohort Database using MS Access
- Data Triangulation
- Communication of Data Analysis Results

## **NEXT STEPS...**

- Work more on NSACP datasets
- Explore Excel functions and options
- Use SPSS with survey data
- Identify topics for analysis and commission
- Practice.... Practice.... & Practice....



National STD/AIDS Control Programme (NSACP), Sri Lanka & The Voluntary Health Services (VHS), India Supported by Centers for Disease Control and Prevention (CDC/DGHT-India) (VHS-CDC Project)