# Use of censuses and surveys in record linkage studies to evaluate completeness of death registration

Chalapati Rao

# Outline

- Background

- Statistical methods

- Examples
  - Viet Nam
  - Indonesia
  - Oman

- Completeness measurement in current context of CRVS development

- <u>Generalizability</u>
  - **Coverage**
  - **Completeness**
    - Aggregated data analysis (indirect methods)
    - Record linkage and matching (direct methods)
- <u>Accuracy</u>
  - **Reliability**
  - **Validity – particularly of registered causes of death**
- <u>Policy relevance</u>
  - **Timeliness**
  - **Sub national data availability (geographical disaggregation)**

- Completeness =

$$\frac{number\ of\ registered\ events}{estimated\ number\ of\ total\ events} \times 100$$

- Comparisons of numbers/rates for same population over time for consistency/time trends

- Comparisons between populations with similar characteristics

- Comparisons between different sources for same population(e.g. census enumerations; health service records etc)

- Overall, not a satisfactory approach (both sources could be of inadequate quality)

- Using models of population growth/ change to derive expected deaths as denominator for completeness

- Models based on assumptions
  - accurate population counts;
  - no migration;
  - accurate age-reporting of population and deaths;
  - completeness invariant by age
  - In some methods – stable population (constant fertility and mortality in preceding decades)

- Vastly differing measures from different methods, with considerable uncertainty **(±25%)**

- Capture-recapture / dual record system/ matching studies

- requires two or more independent sources of information on individual members of the population

- Estimates total population size (total deaths) when a full count of the total population is unavailable or unfeasible from a single source

- Individuals 'captured' in one source and 'recaptured' when matched in 2nd source

- Matching across key variables:
  – Personal details / address variables / Event details - Date of birth/death/registration

- Linkage produces 3 sets i.e Matched records; plus sets of unique records in either source

- record linkage permits another statistical procedure (based on certain conditions) to estimate deaths not captured by either source

- Completeness estimated using denominator from reconciliation of 3 cells OR (Indian Sample Registration System)
- by including the fourth cell (estimated missed deaths) (Chinese DSP)

# Computation

**TABLE 1. Two-source model**

|  | Source Y | | |
|---|---|---|---|
|  | Yes | No | Total |
| Source Z — Yes | $a$ | $b$ | $a + b = Z_0$ |
| Source Z — No | $c$ | $x$ | |
| Total | $a + c = Y_0$ | | $N = a + b + c + x$ |

| Estimated values | | Maximum likelihood estimator (MLE) |
|---|---|---|
| Unobserved cell: | $\hat{x}$ | $bc/a$ |
| Completeness of source $Y$: | $\hat{Y}_c$ | $a/(a + b) = a/Z_0$ |
| Completeness of source $Z$: | $\hat{Z}_c$ | $a/(a + c) = a/Y_0$ |
| Total population: | $\hat{N}$ | $a + b + c + (bc/a)$ or, $(a + b)(a + c)/a$ |

Completeness of Y $= \dfrac{a+c}{a+b+c+x}$

Completeness of Z $= \dfrac{a+b}{a+b+c+x}$

- Hook, E.B. and R.R. Regal, *Capture-recapture methods in Epidemiology: Methods and limitations.* Epidemiologic Reviews, 1995. **17**(2): p. 243-64.

9

- No 'out-of-scope' events in either source

    - Correct identity/time frame/residence status/no migration

- Homogeneity of capture probability in each source

    - No selective exclusion by gender/age/ethnicity/geography/SES

- Independence of data sources (capture in one source does not influence capture in the second source)

- Accuracy of matching procedures and matching outcomes (no erroneous matches or erroneous non-matches)

| Type of data collection | Primary source[1] | Secondary source[2] | Remarks |
|---|---|---|---|
| **Continuous recording systems** | | | |
| Civil registration | Yes | | • Optimal source<br>• annual data on routine basis |
| Alternate registration | Yes | Yes | • Health system vital records e.g Vietnam, Fiji<br>• Church records in Christian societies |
| Sample registration | Yes | Can serve as a secondary source for evaluating CRVS | • Best alternative to CRVS<br>• Indian SRS (ref)<br>• Chinese DSP (ref)<br>• Bangladesh SVRS (ref) |
| Special registration | Yes | Can serve as a secondary source for evaluating CRVS or SRS | • E.g. Health and Demographic Surveillance Sites in several countries (INDEPTH Network) (ref) |
| Age based registers | | Yes | • Maternal/child health<br>• senior citizens /pensioners databases |
| Disease surveillance systems | | Yes | • tuberculosis<br>• cancers<br>• injuries<br>• stroke |
| **Periodic data collections** | | | |
| Census (total population) | Yes | Yes | • Optimal 2nd data source (national coverage) |
| National sample surveys | | Yes | • Inter censal surveys<br>• DHS program<br>• WHO NCD surveillance (STEPS) surveys<br>• UNICEF MICS surveys etc |
| Special surveys designed to assess completeness | | Yes | • Evaluation surveys for sample/special registration<br>• sporadic research based examples |

1 = data source for which completeness needs to be evaluated
2 = data source which will be used to evaluate completeness of the primary source

- <u>Scope of analysis</u> e.g national / sub national measures; by age; pop sub groups

- <u>Availability/choice</u> of primary & secondary data <u>sources</u>

- Reference time period of analysis

- Matching process
  - Manual/electronic
  - Deterministic/probabilistic/implicit rules

- Statistical procedures
  - Data reconciliation
  - Use of multiple parallel sources or partial data sources
  - DRS method ( 2source/multiple source models)
  - Hybrid models

Australian National University

- Completeness of Y $=\dfrac{a+c}{a+b+c+x}$

- Chandra-Deming proposed that if all conditions are met, then

  SE of completeness $= \sqrt{Nq1q2/p1p2}$

- Where N = total number of events estimated by the method (Table 1)

  $p1$ = the probability that an event is recorded in data source 1

  $p2$ = the probability that an event is recorded in data source 2

  $q1$ = the probability that an event is missed in data source 1

  $q2$ = the probability that an event is missed in data source 2

- RMSE of completeness estimate: RMSE = $\sqrt{variance + bias^2}$

- **Variance = sampling error (in one or both sources)**

- **Three sources of bias – out of scope/dependence/matching bias**

  - *Due to varying directions; net bias is usually less than any individual source of bias*

# Methods to measure effect of dependence

TABLE 3. SELECTED ARTICLES DESCRIBING METHODS TO MEASURE BIAS AND ERROR IN COMPLETENESS ESTIMATES FROM DUAL-RECORD SYSTEMS ANALYSES
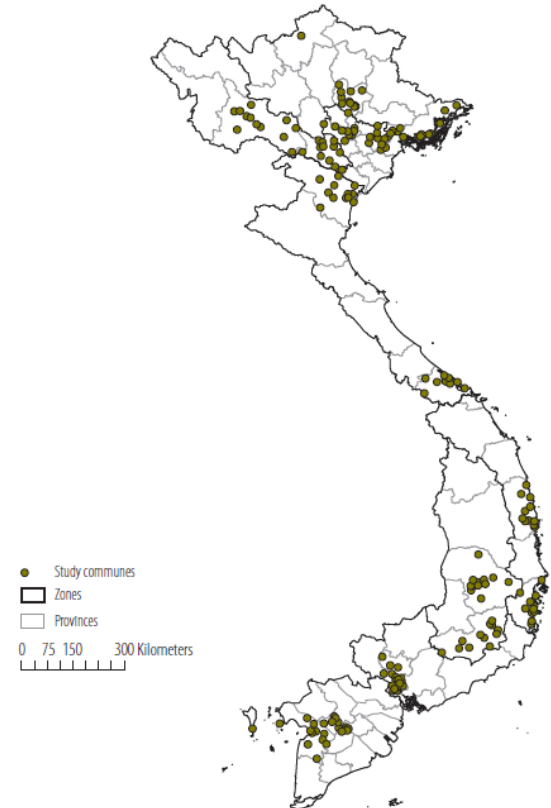
| Author/Year | Source | Title | Example | Methods/Results |
|---|---|---|---|---|
| Seltzer and Adlakha (1974) | International Program for Population Statistics, UNC, Chapel Hill, United States of America. Reprint Series 14, 1974 | On the effect of errors in the application of the Chandrasekar-Deming technique | Theoretical example | Proposes methods for estimating net relative bias from out of scope events, lack of independence, matching errors, and interactions between these three sources. No details on measurement of standard error of completeness estimate. |
| Greenfield (1976) | Journal of the Royal Statistical Society. 139 (3) 389-401 | A revised procedure for Dual-Record Systems in Estimating Vital Events | Malawi Population Change survey, 1972 | Greenfield estimate accounting for lack of independence; 10 per cent higher than standard DRS estimate; no standard error measurement. |
| Raj (1977) | Journal of the American Stats, Assoc. 72 (358) 377-81 | On Estimating the Number of Events in Demographic Surveys | Theoretical example based on 3 scenarios of completeness | Detailed methods for estimating bias from lack of independence; sampling variance; and total mean square error of completeness. |
| Nour (1982) | Journal of the Royal Statistical Society. 145 (1) 106-116 | On the estimation of the Total Number of Vital Events from Dual Systems | Malawi Population Change survey, 1972 | Nour estimate accounting for lack of independence; 4.7 per cent higher than standard DRS estimate; presents method for estimating sampling variance. |
| Chandrasekaran (1983) | Cairo Demographic Centre Working paper 6 | On two estimates of the number of events missed in a dual-record system | Theoretical examples, Indonesian Vital Registration Project | Compares Greenfield method and standard DRS, identifies that error is inversely proportional to the completeness estimate; two estimates provide a plausible range of completeness. |
| Hook and Regal (1995) | Epidemiologic Reviews. 17(2); 243-264 | Capture-recapture methods in Epidemiology | Examples of multiple source data on diseases | Implements separate models for two-source/three-source analysis of missed events, accounting for dependence in each source combination. Proposes use of range of completeness estimates from different methods, rather than any specific variance based standard error calculation of CI. |
| Ayhan (2000) | Journal of Applied Statistics. 27 (2) 157-169 | Estimators of vital events in dual-record systems | Theoretical example based on two sample sources | Method accounts for lack of independence, no details on measurement of standard error of completeness estimate, or of sampling variance. |
| El-Khorazaty (2000) | Environmetrics; 11; 435-448 | Dependent dual-record system estimation of number of Events: a capture-mark-recapture strategy | Vital registration and sample survey data for Egypt, 1974-75 | Method accounts for dependence, but assumes no geographic or matching error. Paper compares completeness estimates from data reconciliation; standard DRS; Greenfield; Nour; and El-Khorazaty. No methods for estimating standard error or variance. |
| Chatterjee and Mukherjee (2013) | arXiv:1311.3812v3[stat.ME]. https://arxiv.org/abs/1311.3812 | Approximate Bayesian solution for estimating population size from a Dual-record system | Malawi Population survey, 1972 | Models account for variations in behavioural response causing dependence between sources. Includes method for estimating Standard error and 95 per cent CI of completeness. |

Hook/Regal proposed to try as many methods as possible, and use the average of all errors

- Study population :192 communes; 2.6 million pop

- Data sources – <u>Commune health</u> (source 1) / <u>Justice system</u> (source 2)

- <u>manual matching</u> at commune level

- <u>relaxation of matching criteria</u> (age, date of death)

- Unobserved cell computed from <u>two source analysis</u>

- <u>Reconciled data used as numerator</u>

- Completeness factor used to adjust life tables etc



Fig. 1. Geographic distribution of communes included in the sample mortality surveillance system, Viet Nam, 2009

Study communes
Zones
Provinces
0  75  150    300 Kilometers

# Matching results

| | Regions | Total in reconciled list | CHC | Population Dep | Justice system | Other |
|---|---|---|---|---|---|---|
| 1 | Ha Noi | 2304 | 1723 (75%) | 1580 (69%) | 1669 (72%) | 720 (31%) |
| 2 | Thai Nguyen | 1185 | 999 (85%) | 210 (18%) | 183 (15%) | 85 (7%) |
| 3 | Hue | 2221 | 1768 (78%) | 1043 (47%) | 1311 (59%) | 777 (35%) |
| 4 | Ho Chi Minh | 2453 | 435 (18%) | 571 (23%) | 1871 (76%) | 202 (8%) |
| 5 | Can Tho | 1758 | 872 (49%) | 758 (43%) | 1081 (62%) | 535 (30%) |

- *A death could be recorded in more than one system*
- ⟷ = *interdependence*

Australian National University

**Table 1.  Age- and sex-specific observed and estimated deaths[a] and completeness of mortality data, Viet Nam, 2009**

| Sex-specific age group (in years) | Sample | a[b] | b[c] | c[d] | x[e] | Other source only | Deaths | | Per cent completeness[f] (95% CI) |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Observed (a + b + c + additional) | Estimated (a + b + c + x) | |
| Males | 1 239 937 | 2138 | 1984 | 1363 | 1265 | 215 | 5700 | 6750 | 81.2 (74.1–87.1) |
| 15–59 | 873 727 | 903 | 873 | 597 | 577 | 92 | 2465 | 2950 | 80.4 (72.2–80.3) |
| 60–74 | 53 985 | 453 | 414 | 274 | 250 | 38 | 1179 | 1391 | 82.0 (74.9–87.9) |
| 75+ | 22 852 | 710 | 629 | 453 | 401 | 77 | 1869 | 2193 | 81.7 (74.7–87.4) |
| Females | 1 309 462 | 1572 | 1413 | 1026 | 922 | 181 | 4192 | 4933 | 81.3 (74.4–87.1) |
| 15–59 | 929 773 | 373 | 350 | 251 | 236 | 56 | 1030 | 1210 | 80.5 (72.5–87.1) |
| 60–74 | 72 999 | 342 | 271 | 213 | 169 | 41 | 867 | 995 | 83.0 (75.4–89.0) |
| 75+ | 37 684 | 812 | 734 | 539 | 487 | 80 | 2165 | 2572 | 81.0 (73.9–87.0) |

CI, confidence interval.

[a] Age- and sex-specific deaths deviate slightly from the totals reported in the text because 27 deaths had no age data.

[b] Number of deaths reported by the Commune Health Centre, the Commune Population and Family Planning Committee (CHC/CPFPC) and the Justice Department.

[c] Number of deaths reported by the CHC/CPFPC but not by the Justice Department.

[d] Number of deaths reported by the Justice Department but not by the CHC/CPFPC.

[e] Estimated number of deaths missing from CHC/CPFPC and Justice Department sources.

[f] Proportion of estimated deaths derived from the list obtained by reconciling the Justice Department and combined CHC/CPFPC lists. Derived with the following formula: $(a + b + c) \div (a + b + c + x) \times 100$.

• Hoa, N.P., Rao C et al., *Mortality measures from sample-based surveillance: evidence of the epidemiological transition in Viet Nam.* Bulletin of the World Health Organization, 2012. **90**(10): p. 764-772.

Table 2. **Summary sex-specific measures of mortality based on WHO, UNPD and Viet Nam census data for the 16 study provinces, Viet Nam, 2009**

| Data source | Per cent data completeness (95% CI) | Life expectancy at birth (95% CI) [e0] | Risk of death in children under 5 (deaths per 1000) [5q0] | Risk of death at ages 15–59 (deaths per 1000) [45q15] | Remaining years of life at age 60 [e60] |
|---|---|---|---|---|---|
| **Males** | | | | | |
| Surveillance sample (unadjusted) | – | 74.4 (74.0–74.8) | 7.4 | 163 | 20.9 |
| Surveillance sample (adjusted)[a] | 81.1 (74.1–87.1) | 70.4 (70.1–70.8) | 24.6[c] | 199 | 19.4 |
| Viet Nam census (unadjusted) | – | 75.2 (75.0–75.4) | 10.9 | 157 | 22.1 |
| Viet Nam census (adjusted)[b] | 65.6 (–) | 68.8 (68.6–69.0) | 16.5 | 230 | 17.9 |
| WHO (2009) | NA (modelled) | 69.8 (–) | 24.6 | 173 | 17 |
| UNPD (2005–2010) | NA (modelled) | 72.3 (–) | No data | 139 | No data |
| **Females** | | | | | |
| Surveillance sample (unadjusted) | – | 82.3 (82.0–82.7) | 5.8 | 57 | 25.1 |
| Surveillance sample (adjusted)[a] | 81.3 (74.4–87.1) | 78.7 (78.4–79.0) | 22.5[c] | 71 | 23.6 |
| Viet Nam census (unadjusted) | – | 85.2 (85.0–85.6) | 8.8 | 50 | 28.4 |
| Viet Nam census (adjusted)[b] | 57.8 (–) | 77.8 (77.5–78.0) | 15.7 | 86 | 22.4 |
| WHO (2009) | NA (modelled) | 74.5 (–) | 22.6 | 107 | 19.8 |
| UNPD (2005–2010) | NA (modelled) | 76.2 (–) | No data | 96 | No data |

CI, confidence interval; NA, not applicable; UNPD, United Nations Population Division; WHO, World Health Organization.
[a] Adjusted for data incompleteness and mortality in children under 5 years of age.
[b] Adjustment by the Preston-Coale method.
[c] WHO estimate.

- Acknowledgement: This study was a PhD thesis by Dr Salah al Muzahmi passed by the University of Queensland, Australia in 2016

- Study covering entire population of Omani nationals (excl expats)

- Data sources – Health system routine data 2010 (Source 1)
  Census 2010 one year recall (Source 2)

- Three rounds of matching – electronic plus manual

- Analysis – capture-recapture adjustment of completeness of death notification data

Table 1  Variables by source

| Variable | BDNS database | Census 2010 database |
|---|---|---|
| Notification number | √ | |
| Reported institution | √ | |
| Name of deceased | √ | |
| Name/tribe name of applicant* | √ | √ |
| Governorate/region | √ | √ |
| *Wilayat* (district) | √ | √ |
| Town/village | √ | √ |
| Locality or compound | | √ |
| Sex | √ | √ |
| Date of death | √ | √ |
| Age at death | √ | √ |
| Date of birth | √ | |

* The applicant for death registration, as well as the census respondent, is assumed to be from the same household and tribe as the deceased. Hence the tribe name of the deceased would be the same as the tribe name of the BDNS applicant as well as the census respondent. Hence, the tribe names were used in the matching process.

Table 1  Missing/duplication of the primary variables.

| Items | Birth and death notification system database | Census |
|---|---|---|
| Total records | 6,039 | 5,400 |
| Missing date of death | 0 | 0^ |
| Duplicates | 3 | 19 |
| Missing age | 652 | 0 |
| Missing sex | 18 | 0 |
| Missing governorate | 457 | 0 |
| Missing *Wilayat* | 535 | 0 |
| Missing nationality | 18 | 0 |
| Missing *Wilayat* and governorate | 457 | 0 |
| Records used in matching | 6,036 | 5,381 |

^ Date of death in the census dataset is divided into three variables (year, month and day); there are 153 records with unknown day and month

## FIRST ROUND

Table 14 Summary findings of the first phase of the matching process

| | Records |
|---|---|
| Matched records in the first round | 568 (9.5%) |
| Not matched from Death notification | 5468 |
| Missing age | 500 |
| Missing governorate | 435 |
| Missing *wilayat* | 502 |
| Missing village/locality | 1022 |

**Reasons for mismatch**
- Variations in
- Spellings
- age
- address
- date of death

## SECOND ROUND

Table 15 Summary findings of the phase two of matching process

| | Records |
|---|---|
| Matched according to age | 2,983 |
| Matched according to date of death | 3,078 |
| Matched according to gender | 3,252 |
| Matched according to *wilayat*/village | 3,284 |
| Total matched records on all variables | 2,983 (49.5%) |

**Correction strategy**
- Corrected spellings, address variables,
- 5 year margin for age, if matched on other variables
- One month margin for date, if matched on other variables

## THIRD ROUND

Table 17   Summary findings of the third round of matching process

| | Records |
|---|---|
| Matched records after third corrections | 4.819 (79%) |
| Not matched | 1,217 |
| Reasons for un-matched records* | |
| Missing age | 192 |
| Missing governorate | 168 |
| Missing *wilayat*/village | 179 |
| Under-recorded events in census | 650 |

\* Some records remained unmatched due to > 1 missing variable

**Correction strategy**
- Field verification of variables for unmatched cases from health records
- 10 year margin for age for deaths above 65 years, if matched on other variables
- Two month margin for date, if matched on other variables

Australian National University

Table 18 Overall completeness of reporting of deaths

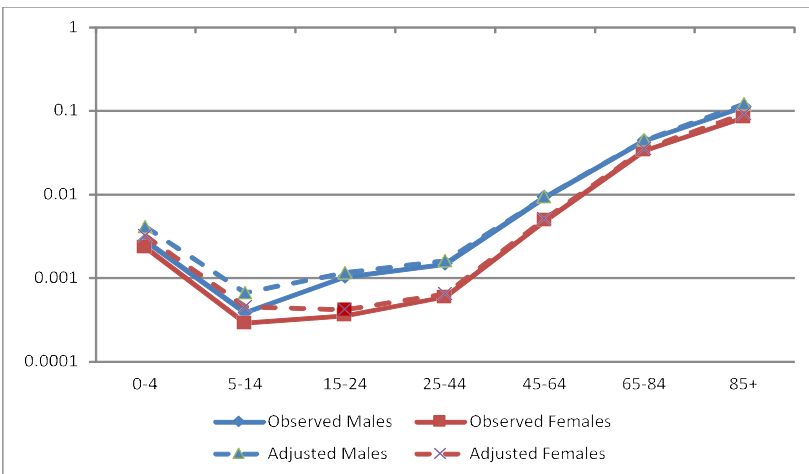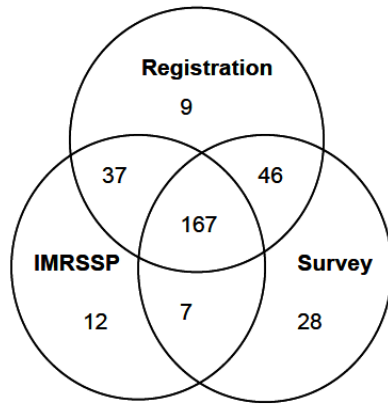|  | | census 2010 | | |
| --- | --- | --- | --- | --- |
|  | | Yes | No | Total |
| Death notification system | Yes | 4,819 | 1,217 | 6,036 |
|  | No | 562 | 142 | 644 |
|  | Total | 5,381 | 1,359 | 6,740 |



**Figure 1: Log plot of Age and sex specific death rate (Observed vs Adjusted), Oman 2010**

| Sex | Governorate | Completeness rate (95% CI) | Adjusted LE in years (95% CI) |
| --- | --- | --- | --- |
| Males | Ad Dhakhliyah | 92 (90 - 95) | 73.7 (72.4 - 74.9) |
|  | Ad Dhahira | 86 (83 - 91) | 72.1 (70.4 - 73.9) |
|  | Al Buraymi^ | 81 (71 - 91) | 81.0 (77.1 - 84.7) |
|  | Dhofar | 87 (82 - 91) | 75.3 (73.5 - 77.0) |
|  | Musandam^ | 83 (72 - 95) | 83.6 (77.5 - 89.6) |
|  | Muscat | 87 (84 - 90) | 74.2 (73.1 - 75.2) |
|  | N Al Batinah | 93 (91 - 95) | 73.0 (71.9 - 74.0) |
|  | N. Sharqiyah | 98 (86 - 93) | 74.1 (72.3 - 75.8) |
|  | S Al Batinah | 91 (89 - 94) | 73.3 (72.0 - 74.6) |
|  | S. Sharqiyah | 92 (89 - 95) | 77.0 (75.0 - 79.0) |
|  | **Total** | 90 (89 - 91) | 73.7 (73.3 - 74.2) |
| Females | Ad Dhakhliyah | 91 (88 - 94) | 78.8 (77.4 - 80.2) |
|  | Ad Dhahira | 87 (82 - 93) | 82.1 (79.7 - 84.3) |
|  | Al Buraymi^ | 84 (72 - 97) | 83.4 (79.2 - 87.6) |
|  | Dhofar | 88 (83 - 93) | 80.2 (78.4 - 82.0) |
|  | Musandam^ | 67 (46 - 87) | 79.6 (76.1 - 83.2) |
|  | Muscat | 82 (78 - 86) | 80.3 (79.1 - 81.5) |
|  | N Al Batinah | 97 (95 - 98) | 80.6 (79.3 - 82.0) |
|  | N. Sharqiyah | 90 (86 - 95) | 79.3 (77.5 - 81.3) |
|  | S Al Batinah | 90 (86 - 93) | 81.5 (80.0 - 83.2) |
|  | S. Sharqiyah | 89 (84 - 94) | 86.3 (84.1 - 88.4) |
|  | **Total** | 89 (88 - 90) | 80.0 (79.5 - 80.4) |

# Indonesia : Three independent sources

- Central Java – record linkage/matching across three sources (<u>health system</u>, <u>vital registration</u>, <u>independent survey</u>)

- Independent survey and record linkage/matching conducted only in a sample of villages from the overall study population

- Completeness of health system data calculated as a proportion of total deaths obtained from the reconciled list of unique deaths
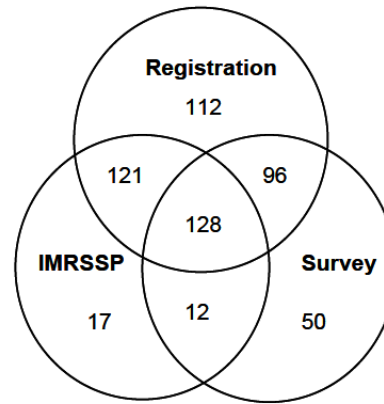
**PEKALONGAN**

Registration
9

37          46

167

IMRSSP          Survey

12      7      28

Total deaths = 306

**SURAKARTA**

Registration
112

121          96

128

IMRSSP          Survey

17      12      50

Total deaths = 536

Completeness = 73%          Completeness = 55%          25

- Conditions for using capture-recapture methods are 'data driven' as compared to the demographic assumptions of underlying fertility /mortality/population growth/migration patterns in the study population

- The data collection procedures allow direct assessment of bias and error

- Independent survey findings can identify systemic weaknesses in registration

- Involvement of local staff in matching builds awareness for improving registration

- Age specific measures of completeness

- Data reconciliation from additional sources helps fill data gaps in cause of death information

- Availability of <u>computerised data sources</u> from registration and census/surveys

- Electronic linkage vastly <u>reduces logistical challenges</u> of manual matching

- Wider use and recording of <u>Unique Identifiers</u> which are invaluable for linkage

- Potential to <u>improve data quality</u> of recorded variables used in matching (name spellings; address variables, age, date of death etc)

- Explicit rules and <u>probabilistic approach</u> using computerised datasets can be applied to test a range of scenarios and judge cut points for specific criteria

- <u>Routine application</u> of these methods <u>in  India and China</u> serve as robust examples of their general acceptability

Mortality estimates by age, sex and cause are universally recognised as essential data for population health assessment

To the extent that the dictum since 1990 has been

𝕎𝕳𝕰𝕽𝕰 𝕿𝕳𝕰𝕽𝕰 𝕴𝕾 ℕ𝕺 𝔻𝔸𝕿𝔸, 𝕸𝕺𝔻𝕰𝕷 𝕴𝕿

Currently, modelling is guided by national mortality data availability score

'**Percent well certified**' = *completeness (%) * (1 - % 'ill-defined causes')*

Lower the score, higher the extent of statistical modelling for estimation **(GBD)**

Negligible = 0-34%;    Partial = 35 – 84%;    Adequate = 85%

**Table 2** Distribution of countries by mortality data quality rating* according to geography and population size

| WHO region† | Population | Data quality rating | | | |
| | | Negligible | Partial | Adequate | Total |
|---|---|---|---|---|---|
| **Africa** | | | | | |
| | <10m | 16 | 2 | 1 | 19 |
| | 10–50m | 22 | 0 | 0 | 22 |
| | >50m | 4 | 2 | 0 | 6 |
| **Americas** | | | | | |
| | <10m | 1 | 16 | 5 | 22 |
| | 10–50m | 2 | 5 | 5 | 12 |
| | >50m | 0 | 1 | 3 | 4 |
| **Eastern Mediterranean** | | | | | |
| | <10m | 6 | 2 | 0 | 8 |
| | 10–50m | 8 | 3 | 0 | 11 |
| | >50m | 1 | 2 | 0 | 3 |
| **Europe** | | | | | |
| | <10m | 3 | 13 | 16 | 32 |
| | 10–50m | 1 | 8 | 5 | 14 |
| | >50m | 0 | 3 | 3 | 6 |
| **South east Asia** | | | | | |
| | <10m | 2 | 1 | 0 | 3 |
| | 10–50m | 2 | 1 | 0 | 3 |
| | >50m | 1 | 4 | 0 | 5 |
| **Western Pacific** | | | | | |
| | <10m | 10 | 4 | 2 | 16 |
| | 10–50m | 2 | 1 | 1 | 4 |
| | >50m | 1 | 4 | 0 | 5 |
| **World** | | | | | |
| | <10m | 38 | 38 | 24 | 100 |
| | 10–50m | 37 | 18 | 11 | 66 |
| | >50m | 7 | 16 | 6 | 29 |
| **Total** | | 82 | 72 | 41 | 195 |

- Completeness estimation -a combination of science and art

- Existing and future 'market' for completeness estimation over next 3 decades, as CRVS systems develop globally

- Current market monopoly at global level

- Need for new players at country level, along with simple methods for error measurement