# Introduction to Cohort Studies

Beate Ritz MD PhD

Epi 200A

Fall 2012

---

**Table 1.** Validity for etiologic inference according to study designs

| Validity ranking | Types of study design |
|---|---|
| Highest | Randomized clinical trial |
| | Prospective cohort study |
| | Retrospective cohort study |
| | Nested case–control study |
| | Time–series analysis |
| | Cross-sectional study |
| | Ecologic study |
| | Cluster analysis |
| | Case study |
| Lowest | Anecdote |

---

## MacMahon and Pugh, 1970
Definition of cohort studies (in public health epidemiology)

■ The group or groups of persons to be studied are defined in terms of **characteristics manifest prior to the appearance of the disease** under investigation

■ The study group so defined are observed over a period of time to determine the **frequency of disease** among them

---

## Cohort studies

Simplistic description

• A cause 'looking' for a disease

• (*versus* case-control study: "A disease 'looking' for a cause")

---

## Cohort design:
Retrospective (historical) in terms of
  a) timing of events *or*
  b) data collection

Cohort is enumerated some time in the past and followed over *historical* time (to today)
  – time of follow-up long (20-40 years), often extends across decades
  – cohort can be large i.e. 10,000+ members

But, how do we:
° "reconstruct" the cohort - who belongs into the cohort?
• Obtain exposure and outcome information
    • Note: a historical cohort is often restricted to investigations of fatal disease (why!)

---

## Cohort design:
Prospective in terms of
  a) timing of events *or*
  b) data collected

This design is best to be used for
• short-term (common) health outcomes; e.g. for:
  – physiological changes (blood pressure and noise)
  – acute neurotoxic effects (OP pesticides)
  – pulmonary function (cotton dust)
  – skin rashes (irritants, e.g. solvents, metals)
  – injuries
  – allergic reactions, asthma attacks
• prospective medical surveillance

1

## Cohort design:
Prospective *or* retrospective in terms of
- a) timing of events *or*
- b) data collected

The major issue we want to convey is whether disease status could have influenced exposure measurement/information (such as via recall of exposure by a diseased subject)

Note that retrospective often is considered a 'less reliable' design; thus, be clear about how you use this term

### Cohort study: examples

Cohort: "Any designated group of individuals who are followed or traced over a period of time"

Historically:
- John Snow: Cholera in London (1854)
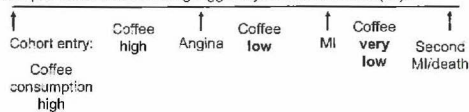- Panum: Measles on the Faroe Islands (1846)

More recent
- Framingham: cardiovascular diseases (N=5,209); bi-annual exams, medical records and deaths info
- British doctors: smoking and lung cancer among British doctors (N=34,439 male British doctors in 1951; Doll)
- Perinatal collaborative study: pregnancy and child health, cerebral palsy and allied neurological defects (N=42,000 pregnant women enrolled 1959–1966 at 12 hospitals across the United States)
- Nurses Health Study: established in 1976 from female US registered nurses ages 30-55 years who responded to a mailed questionnaire that inquired about risk factors for cancer and heart disease (N=121,700)

- HIV cohorts: 1984-2005, Multicenter AIDS Cohort Study (N=4,955 homosexual men who volunteered in Baltimore, Chicago, Los Angeles, and Pittsburgh)
- EPIC study: cancer
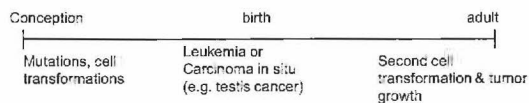- California Teachers Cohort (125,000 in 1995): Breast cancer
- And many more

## Causal Inferences in Cohort Studies

- Since the only *sine qua non* causal criteria states that **a cause precedes its effect,** it is logical to start with the exposure and follow exposed people forward in time to study the occurrence of the health endpoint of interest
- This was hardly done *prospectively* before the Framingham cohort study (baseline 1948); too expensive, too time consuming
- A cohort study is a logical design to study determinants of the changes from not having a disease to having a disease. The study can guarantee that exposure precedes the onset of clinical diagnosis (but perhaps not the real onset of pathological changes).

Example: Does coffee drinking trigger myocardial infarction (MI)?

| Cohort entry: | Coffee high | Angina | Coffee low | MI | Coffee very low | Second MI/death |
|---|---|---|---|---|---|---|

Coffee consumption high

## Experimental *vs.* Observational Studies:
Why not conduct a randomized trial?

Trials
- cannot obtain evidence for harmful agents (and sometimes for beneficial ones as well)
- deal by nature with (very) selected populations
- not practical for
  - rare outcomes (*Note:* we would expect only 50-200 lung or colon cancers and 16 Parkinson's cases per 100,000 person years of observation in most working age people)
  - long follow-up times that allow for latency
  - effects that occur late in disease progression
- focus on one (or several) specific doses only
- expensive to conduct

### Life course perspective

| Conception | birth | adult |
|---|---|---|
| Mutations, cell transformations | Leukemia or Carcinoma in situ (e.g. testis cancer) | Second cell transformation & tumor growth |

Different causal components may be operating

Note: many cohorts recruit at entry only few of those eligible (% of all eligible often not known):
- What is the impact on internal validity and external validity?

## Cohort studies: recruitment
- Recruitment to the cohort may be mandatory/ automatic
  - All in public registers = mortality, births, deaths, cancer (without informed consent)
  - Occupational cohorts using employment data from occupational plants (assess exposures retrospectively from records and outcomes from registers)
- NOTE: cohorts using "primary" data (i.e. collected during/for the investigation) are usually based upon **informed consent**

Examples:
- via General Practitioner – e.g. Danish National Birth Cohort
- Letters – e.g. to members of an organization (British doctors, CA Teachers, Nurses Health Study, Harvard Alumni)
- Advertisements – e.g. people with a given disease
- Local community: ALSPAC, Framingham
- Visitors to a website
- Participants in L.A. Marathon

## Cohort studies: follow-up

- Compliance to follow-up procedures
  - frequent contacts needed!
  - Are (health) benefit incentives given?
- Recording of endpoints
  - rely on diagnoses made by the health care system
  - repeated measurements necessary?
- Changes in other determinants/ covariates
  - questionnaires
  - interviews
  - measurements
- Participation is voluntary, participants are free to leave the cohort at any point in time
  - right to remove data from the study?

| TABLE 1. Types of outcomes for cohort |
| --- |
| Discrete events |
| Single events |
| Mortality |
| First occurrence of a disease or health-related outcome |
| Incidence (density) |
| Cumulative incidence (risk) |
| Ratios (incidence density and cumulative incidence) |
| Multiple occurrences: |
| Of disease outcome |
| Of transitions between states of health/disease |
| Of transitions between functional states |
| Level of a marker for disease or state of health |
| Change in a functional/physiologic/biochemical/anatomic marker for disease or health |
| Rate of change |
| Patterns of growth and/or decline |
| "Tracking" of markers of disease/health |
| Change in level with time (age) |

Source: Tager IB. Outcomes in cohort studies. *Epidemiologic Reviews* 1998, 20(1).

**TABLE 1. Types of outcomes in cohort morbidity studies**

| Induction period/ reversibility | Event (dichotomous) | Change in status (continuous) |
| --- | --- | --- |
| **Short (days to months)** | | |
| Reversible | Asthma attack<br>Tendonitis<br>Contact dermatitis | Cross-shift function (FEV$_1$*)<br>Temporary threshold hearing |
| Irreversible | Asthma diagnosis<br>Spontaneous abortion<br>Amputation | Annual change in FEV$_1$ |
| **Long (years)** | | |
| Reversible | Chronic bronchitis<br>Endometriosis<br>Carpal tunnel syndrome | Sperm count<br>Blood pressure |
| Irreversible | Silicosis<br>Myocardial infarction<br>Infertility | Noise-induced hearing loss<br>Atherosclerosis<br>Hepatic fibrosis |

*FEV$_1$, forced expiratory volume in 1 second.

Source: Checkoway H and Eisen EA. Developments in Occupational Cohort Studies. *Epidemiologic Reviews* 1998, 20(1).
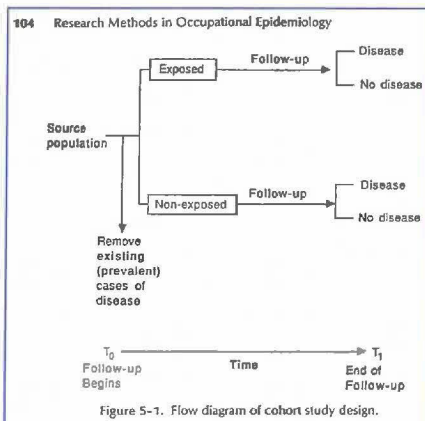
## Cohort Entry Definitions

Entry to a cohort can be defined at a fixed point in time:
- All subjects are selected at a given point (range) in time, e.g. from a registry of a type of people
  - All atomic bomb survivors in Japan on Jan 1st 1950 living in Nagasaki and Hiroshima
  - European Prospective Investigation into Cancer and Nutrition (EPIC), a multi-centre prospective cohort study in 23 study centers in ten European countries
    - E.g in Germany, recruitment was based on a random sample of subjects in targeted age range (women aged 35–65, men 40-65) from population registers between 1994 and 1998
    - participation rate was 38.5% (i.e. observed cohort is a self-selected subgroup of the underlying population)

or

  - subjects enter the cohort at different points in time; *e.g.*: all inhabitants of Framingham/MA that reach a certain age



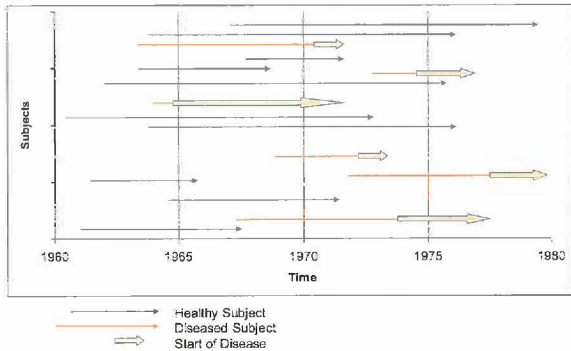Figure 5-1. Flow diagram of cohort study design.

## Cohort Exit Definitions

Subjects can be follow-up
- until a fixed point in calendar time (end of study);
  - note: some subjects are observed for a shorter time i.e. due
    - incidence of the disease under investigations,
    - death,
    - migration or
    - loss of follow-up

- or as long as they are
  - employed
  - live in the city
  - have the exposure (are "right censored" when this changes) (e.g. use of a certain type of medication)

## Study Design Overview: Identifying Diseased Subjects in a Population



Legend:
- Healthy Subject
- Diseased Subject
- Start of Disease

---

## Cohort studies: exposure assessment

- Exposure may have started at a given point in time:
  - E.g. at baseline or any other measurement point
  - and remains fixed ("ever smoker")
  - or changes over time (amount of smoking)
- Exposure can be measured as:
  - Average or cumulative exposure over time
  - exposure level at baseline
  - Note: without a prior hypothesis (or knowledge of biological mechanism) there may be numerous ways of analyzing exposure data
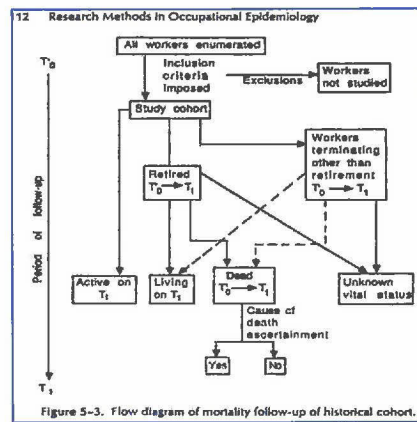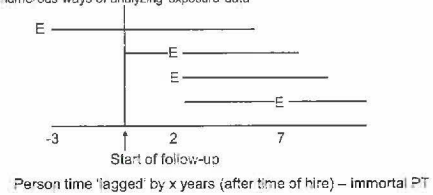


Person time 'lagged' by x years (after time of hire) – immortal PT

---

## Cohort studies: exposure assessment

- Exposures can be lagged (i.e. exclude exposure during time irrelevant for the disease)
  - E.g. exposure too close to disease onset
- Exposure contrast
  - Generally we like to examine as large an exposure contrast as possible - thus, we want to establish a cohort with different exposure levels (e.g. workers in a copper-smelter compared to the general population)

- Select the non-exposed subjects as close to the **counterfactual** ideal as possible
  - Non-exposed subjects should have the same disease risk as the exposed had they not been exposed

---



Figure 5-3. Flow diagram of mortality follow-up of historical cohort.

---

## Start of follow-up in a cohort study

- hire date or fixed time/date after hiring

- first monitoring date (e.g. radiation monitoring, blood lead monitoring)

- fixed date (such as Jan 1970)



Figure 5-2. Cohort membership enrolment options.

---

## End of follow-up in a cohort study

- end of follow-up for the cohort reached
- death or incidence from outcome of interest
- death from competing causes
- last known date alive (after that we call them 'lost to follow up')

Or

- should we assume a worker is alive if no information is found that indicates that the subject died (and thus continues to add person-time)?



Fig. 1. Follow-up over a 34 year-period of a hypothetical cohort.

Figure 5–3. Flow diagram of mortality follow-up of historical cohort.

# Summary: Cohort Studies

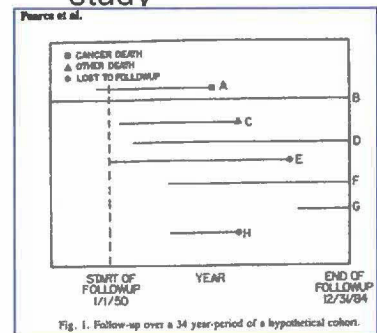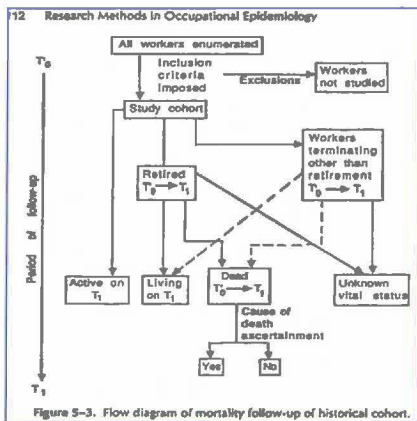- Generally most accepted in scientific community

- Include the entire available study population

- Most similar to standard experimental strategies
  - determine (rather than apply) a toxin or preventative agent among subjects disease-free at baseline
  - follow-up subjects over time
  - observe adverse or positive health effects in exposed and non-exposed subjects

The goal is to estimate the risk of (various or one) disease/s among the exposed subjects relative to the background risk experienced by "comparable" unexposed persons:
- comparable refers to the "exchangeability assumption" or "counterfactual"
  - *what would have happened to this group of exposed subjects if they had NOT been exposed?*

# Summary: Cohort Studies

- Select non-exposed as close to the counterfactual ideal as possible:
  - Non-exposed should have the same disease risk as the exposed had they not been exposed
- Recruitment to the cohort
  - based upon informed consent if primary data are collected
  - Without informed consent if all are followed in public registers = mortality, births, deaths
- Historical cohorts: e.g. use existing data but need not be 'retrospective'

# Advantages of the cohort method

- In principle, can provide a **complete description of experience of cohort members** subsequent to exposure, including rates of progression to and staging of disease, and natural history of disease
- Allows study of **multiple potential effects** of a given exposure, thereby obtaining information on potential benefits as well as risks
- Allows for the **calculation of rates** of disease in exposed and unexposed individuals and time to event
- Permits **flexibility in choosing variables** to be systematically recorded
- Allows for **thorough quality control in measurement** of study variables (not in historical cohort studies though)

# Disadvantages of the cohort method

- Large numbers of subjects required (thus, low feasibility to study rare diseases)
- Relatively expensive to conduct
- Potentially long duration for follow-up necessary
- Exposures may change, making findings irrelevant unless the exposure assessment is adapted
- Maintaining follow-up may be difficult
- The cohort is generally not representative of the general population

# Example: The Agricultural Health Study Cohort (AHS)

- Collaborative effort to study the effects of pesticide exposures among farmers
  - National Cancer Society (NCI)
  - National Institute of Environmental Health Sciences (NIEHS)
  - U.S. Environmental Protection Agency (EPA)

http://aghealth.nci.nih.gov/

5

## The AHS Cohort study:
### Retro- and prospective data collection

- Phase I, initial cohort recruitment, 1994-1997:
- 89,658
  - private pesticide applicators, and
  - spouses of private applicators, and
  - commercial pesticide applicators
- Recruited at Iowa and North Carolina state pesticide applicator licensing facilities
- Each pesticide applicator asked to complete a 21-page enrollment questionnaire
  - a. Demographic data
  - b. Pesticides used (50 pesticides), other pesticide-related questions
  - c. Lifestyle (i.e., smoking, alcohol, vegetable, and fruit consumption)
  - d. Brief medical history
  - e. Family history of cancer, kidney failure, diabetes, and heart disease
  - f. Farm exposures other than pesticides (not in commercial pesticide applicator version)
  - g. Personal identifiers, spouse identifiers, children identifiers

Farmer applicators completing the enrollment questionnaire are given three take-home questionnaires (scanable) for
- the applicator (licensing exam taker)
- spouse, and
- female and family health questionnaires

## The AHS Cohort

Take Home Questionnaires:
Farmer Applicator/Commercial Applicator

a. Farm exposures (comprehensive)
b. Pesticide use information (i.e., methods of application, additional pesticides used)
c. Work practices used currently versus those used 10 years ago
d. Other occupational exposures
e. Leisure and work physical activity, physical attributes (e.g., height, weight, eye color, skin pigmentation category)
f. Dietary and cooking practices
g. Medical history (comprehensive)
f. Personal identifiers

## The AHS Cohort

- Cancer and non-cancer outcomes
  - Linkage with
    - » cancer registries
    - » vital statistics
    - » United States Renal Data System (USRDS)
  - Exposure data collection
    - » Baseline questionnaire at licensing exam
    - At follow-up
    - » telephone interviews (CATI)
    - » food frequency questionnaire and
    - » cheek cell collection

- Phase II: follow-up in 1999-2003
- Phase III: follow-up in 2004-2008

## The AHS Cohort

1. Cohort studies
   - All cause and cancer mortality
   - cancer incidence
2. Cross-sectional studies:
   - Using questionnaire data, functional measures, biomarkers, and GIS
   - E.g. cross sectional immunology study of atrazine applicators/corn farmers in Iowa
3. Nested case-control studies
   - High pesticide exposure events
   - Parkinson's disease study
4. Exposure assessment and validation studies

## The AHS Cohort

### Table 1. Composition of Cohort and Data Collection Progress

| | Phase I (Complete) | Phase II (In Progress) [2] | | |
|---|---|---|---|---|
| | Contacts Completed | Main Qx Admin | Buccal Cell Collection | Dietary Health Qx Admin |
| Private Applicators | 52,395 | 26,575 | 14,577 | 14,882 |
| Spouses | 32,347 | 20,856 | 12,030 | 13,224 |
| Commercial Applicators [1] | 4,916 | 0 | 0 | 0 |
| Total | 89,658 | 47,431 | 26,607 | 28,106 |

[1] Phase II data collection on Commercial Applicators not yet begun
[2] Progress through October 12, 2001

## The AHS Cohort

### Table 2a: Post-enrollment (Incident only) Malignant Cancer Cases by Site and Phase II Data Collection progress [1,2,3]

| Cancer Site | Total with Cancer | Post-enrollment Cases Only | | |
|---|---|---|---|---|
| | | Completed Phase II Qx | Returned Buccal Sample | Returned Dietary History Qx |
| Breast | 268 | 181 | 131 | 142 |
| Prostate | 572 | 337 | 215 | 210 |
| Colon | 224 | 106 | 64 | 73 |
| Lung | 180 | 41 | 21 | 23 |
| NHL | 79 | 29 | 23 | 25 |
| Other [4] | 789 | 320 | 217 | 216 |
| Total | 2112 | 1014 | 671 | 689 |

Table 2b: Pre- and Post-enrollment (Prevalent and Incident) Malignant Cancer Cases by Site and Phase II Data Collection progress [1,2,3]

## Ag-Health study topics

Cancer mortality and incidence in Applicators and Spouses
Pesticide Exposure Assessment, Applicators, Spouses and Children – questionnaires
Pesticide Exposure Assessment - Field Studies – Acute exposures
Biologic and Functional Effects of Chronic Pesticide Exposure
Biomarkers and Molecular Genetics
Injury
Lifestyle and Diet
Non-pesticide Exposures, Exposure to Animals
Respiratory Disease and Function
Neurological Disease and Function
Reproductive Health, Child and Adolescent Health
Autoimmune Disease and Immune Function
Other Non-cancer Chronic Disease

## Pooling of cohorts

Advantages:
- Can study rare outcomes
- Conduct subgroup analyses for effect measure modifiers (e.g. sex, race etc)
- Wide geographic distribution allows spread of exposures
- Availability of prospective data; stored serum blood samples can be analyzed by same lab

Disadvantages
- Usually no common data elements, i.e. diverse data collection methods need to be reconciled
- Some variables may not have been collected at all; how to handle missing data?

## Vid D and type 2 diabetes: meta-analysis



Fig. 1. Association between vitamin D and Type 2 diabetes, measured as circulating blood and dietary exposures.



Fig. 2. Associations of vitamin D with type 2 diabetes and other metabolic outcomes.

## Person time

Incidence Proportion: A/N     A= case number N=initial population size
Person-time instead of persons:
A/T observed rate [A= observed cases and T= person-time units in study group]

Poisson model
$$Pr(A=a) = exp(-I^*T)(I^*T)^a/a!$$

I = the rate parameter (average rate we would observe if we repeated the study over and over under the same conditions with the same amount of person-time T observed each time(i.e. end the follow-up when we reach T)
Note: Under the Poisson model     A/T is the MLE estimator of I

Immortal person time
The study has a criterion for a minimum of time before a subject is eligible to be in the study:
E.g. in occupational cohort studies when workers are required to have worked for a minimum of x-years. All workers who did not work for this length of time are automatically not enrolled in this cohort and all of those who are could not be censored prior to 2 years i.e. could not have died if included in the cohort.

This time should not be used to calculate person-time for those included in the cohort

Figure 3-4: Example of a small closed population with end of follow-up at 19 years
see ME3 p.42



|  | Start | Outcome event times (tk) | | | | End |
|---|---|---|---|---|---|---|
|  | 0 | 2 | 4 | 8 | 14 | 19 |
| Index (k) | 0 | 1 | 2 | 3 | 4 | 5 |
| No. of outcome events (Ak) | 0 | 1 | 2 | 1 | 1 | 0 |
| No. at risk (Nk) | 9 | 9 | 8 | 6 | 5 | 4 |
| Prop. Surviving (Sk) |  | 8/9 | 6/8 | 5/6 | 4/5 | 4/4 |
| Length of interval (Δtk) |  | 2 | 2 | 4 | 6 | 5 |
| Person time (NkΔtk) |  | 18 | 16 | 24 | 30 | 20 |
| Incidence rate (Ik) |  | 1/18 | 2/16 | 1/24 | 1/30 | 0/20 |

108 PY total

**EXAMPLE:**
Incidence rate ratios (IRR) for epilepsy among children exposed to pre-eclampsia or eclampsia

| Pre-eclampsia or Eclampsia | Entire Birth Cohort | | | | | Cohort of children without cerebral palsy or a low Apgar score† | | | |
|---|---|---|---|---|---|---|---|---|---|
|  | Person years | No. of epilepsy cases | IR | Crude IRR (95%CI) | Adjusted* IRR (95%CI) | Person years | No. of epilepsy cases | IR | Adjusted* IRR (95%CI) |
| Non-exposed | 17,850,197 | 19,441 | 108.9 | 1.00 | 1.00 (Ref) | 16,651,803 | 15,734 | 94.5 | 1.00 (Ref) |
| Pre-eclampsia |  |  |  |  |  |  |  |  |  |
| Mild | 458,558 | 620 | 135.2 | 1.27 | 1.20 (1.11-1.30) | 418,784 | 485 | 115.8 | 1.20 (1.10-1.32) |
| Severe | 78,386 | 135 | 172.2 | 1.54 | 1.14 (0.96-1.36) | 68,957 | 94 | 136.3 | 1.22 (0.99-1.49) |
| Eclampsia | 7,672 | 15 | 195.5 | 1.78 | 1.35 (0.81-2.24) | 6,604 | 10 | 151.4 | 1.35 (0.73-2.52) |
| Unspec. | 43,328 | 49 | 113.1 | 1.04 | 0.95 (0.72-1.26) | 40,002 | 42 | 105.0 | 1.05 (0.77-1.42) |

IR: incidence rate /100,000 person-years

BRESLOW AND DAY

Fig. 2.1  Schema showing the follow-up of one person in a cohort study



BRESLOW AND DAY

Fig. 3.1  Schematic diagram illustrating proper and improper methods of allocation of person-years. ×, death from cause of interest; O, withdrawal



# Person-time calculations

Table 2.1  Calculation of exact and approximate age- and year-specific person-years at risk

| Point[a] | Coordinates (year, age) | Quinquennium | | Person-years | |
|---|---|---|---|---|---|
|  |  | Year | Age | Exact | Approximate |
| A | (1956.03, 43.71) | 1955–1959 | 40–44 | 1.29 | 1.50 |
| B | (1957.32, 45.00) | 1955–1959 | 45–49 | 2.68 | 2.00 |
| C | (1960.00, 47.68) | 1960–1964 | 45–49 | 2.32 | 3.00 |
| D | (1962.32, 50.00) | 1960–1964 | 50–54 | 2.68 | 2.00 |
| E | (1965.00, 52.68) | 1965–1969 | 50–54 | 2.68 | 2.00 |
| F | (1967.15, 54.83) | 1965–1969 | 50–54 | 2.15 | 2.50 |
| Total |  |  |  | 11.12 | 11.00 |

[a] See Figure 2.1

## Incorrect *vs.* correct person-time calculations

Table 3.1  Reanalysis of data by Duck *et al.* showing original *versus* revised numbers of expected deaths and SMRs by duration of exposure and cause of death[a]

| Cause of death | Duration of exposure (years) | No. of observed deaths | No. of expected deaths | | SMR | |
|---|---|---|---|---|---|---|
|  |  |  | Original | Revised | Original | Revised |
| All causes | 0–14 | 111 | 100.92 | 118.97 | 110 | 94 |
|  | 15+ | 25 | 41.30 | 24.15 | 61 | 104 |
| Total cancers | 0–14 | 27 | 25.55 | 29.93 | 106 | 90 |
|  | 15+ | 8 | 10.89 | 6.51 | 73 | 123 |
| Digestive system cancers | 0–14 | 7 | 7.77 | 9.10 | 90 | 77 |
|  | 15+ | 4 | 3.31 | 1.98 | 121 | 202 |
| Lung cancer | 0–14 | 13 | 10.73 | 12.57 | 121 | 103 |
|  | 15+ | 3 | 4.80 | 2.96 | 62 | 101 |

[a] From Duck *et al.* (1975); Duck & Carter (1976)

8

Table 2.2 Exact and approximate* person-years of observation in the Montana cohort, by age and calendar year

BRESLOW AND DAY

RATES AND RATE STANDARDIZATION

Table 2.3 Number of deaths and death rates (per 1000 person-years)* from all causes in the Montana cohort, by age and calendar year

## Role of Statistical Modeling

Construction of a probability model that explicitly recognizes
- the role of chance mechanism in producing some variation in the rates;
- i.e. observed rates are regarded as just one of the many possible realizations of an underlying random process.

Parameters in the model describe systematic effects of
- exposure of interest
- confounding variables such as age, period, length of follow-up etc.

Estimates of these parameters, obtained during the process of fitting the model, serve as summary statistics analogous to SMR or MH estimates of relative risk.

## Role of Statistical Modeling

Advantage of model fitting over standardization:
- facilitates simultaneous consideration of several different exposure variables at risk
- estimates of relative risk obtained by model fitting generally have greater numerical stability than those computed from standardized rates.

Disadvantage of model fitting:
- parametric specification of the model due to statistical rather than biological criteria. Note: epidemiologic data are rarely extensive enough to allow to discriminate between closely related models (according to model fit criteria).

## Risk set approach in a cohort study

- each subject that enters the cohort at some *entry time* is *at risk*
- each subject exits the study either as a *failure* i.e. contracting or dying of the disease of interest or is *censored*, i.e. is alive at the end of study, is lost to follow-up or does not contract the disease
- associated with each subject is a covariate history – fixed or time-dependent –, including factors that are known or believed to be related to the rate of the disease of interest
- At each failure a *risk set* is formed of the size *m* that included the *case* (failure at that failure time) and all *controls*, i.e. any other cohort member who is at *risk* at the failure time.

Note: The approach that organizes the cohort data by risk sets leads to data which looks just like a matched case-control study and hence we can use the conditional logistic likelihood for the analysis

*also note*: the risk sets are not independent, i.e. subjects can be sampled as controls in multiple risk sets and failures can serve as controls in risk sets prior to their failure times.

## Risk set approach in a cohort study

Confounder control can be achieved by either
- Modeling the effect of the confounder
- Restricting each risk set to those who have similar (or the same) confounder values (=matching).

*Note*: if the matching factors are categorical this approach corresponds to stratification in the Cox model

# Sampling from Risk Sets

- Risk set sampling designs are intrinsically related to semiparametric estimation methods for parameters in the Cox proportional hazards model used in the analysis of full cohort data.

- A sampled risk set of size $m$ is a subset of the risk set that contains
  - the case and $m-1$ sampled controls
  - e.g. 1:1 simple nested case-control sampling: each risk set consists of the case and one control randomly sampled from all the controls in the risk set.
  - *note*: one can use the $(m-1)/m$ relative efficiency rule for control sampling versus full cohort analysis for testing associations between single exposures and diseases (Breslow and Patton, 1979)
  - Thus, we have for 1 case and 4 controls (or 4/5=0.8 or 80% efficiency but then for one case and 5 controls 5/6=0.83 or 83% power, and for 9/10=0.90 or 90% power, thus, we need to add 4 controls to gain10% efficiency, i.e. double your efforts to increase efficiency only slightly; it gets worse after that add another 10 controls and you get 19/20=0.95 only 5% efficiency added