



Epidemiology / Épidémiologie

# Biostatistics and epidemiology: measuring the risk attributable to an environmental or genetic factor

Jacques Benichou

*Unité de biostatistique, CHU de Rouen & Inserm U 657, Institut hospitalo-universitaire de recherche biomédicale, Université de Rouen, 1, rue de Germont, 76031 Rouen cedex, France*

Received 23 February 2007; accepted after revision 24 February 2007

Available online 12 April 2007

Presented by Alain-Jacques Valleron

---

## Abstract

Disease frequency is measured through estimating incidence rates or disease risk. Several measures are used for assessing exposure–disease association, with adjusted estimates based on standardization, stratification, or more flexible regression techniques. Several measures are available to assess an exposure impact in terms of disease occurrence at the population level, including the commonly used attributable risk (AR). Adjusted AR estimation relies on stratification or regression techniques. Sequential and partial ARs have been proposed to handle the situation of multiple exposures and circumvent the associated non-additivity problem. Despite remaining issues in properly interpreting AR, AR remains a useful guide to assess prevention strategies. **To cite this article: J. Benichou, C. R. Biologies 330 (2007).**

© 2007 Académie des sciences. Published by Elsevier Masson SAS. All rights reserved.

## Résumé

**Biostatistique et épidémiologie : la mesure du risque attribuable à un facteur environnemental ou génétique.** La mesure de la fréquence d'une maladie repose sur les concepts d'incidence et de risque. Plusieurs mesures d'association entre exposition et maladie existent. Leur estimation ajustée repose sur la standardisation, la stratification, ou les méthodes plus flexibles de régression. Le risque attribuable (RA) est une mesure d'impact populationnel d'une exposition sur la survenue de nouveaux cas. Son estimation ajustée repose sur la stratification ou la régression. Les RAs séquentiels et partiels permettent de prendre en compte plusieurs expositions et le problème associé de non-additivité. Malgré certaines questions d'interprétation, le RA demeure un guide utile à l'évaluation de stratégies de prévention. **Pour citer cet article : J. Benichou, C. R. Biologies 330 (2007).**

© 2007 Académie des sciences. Published by Elsevier Masson SAS. All rights reserved.

**Keywords:** Incidence; Hazard; Risk; Measure of association; Measure of impact; Attributable risk; Sequential attributable risk; Stratification; Regression; Causality

**Mots-clés :** Incidence ; Risque instantané ; Risque ; Mesure d'association ; Mesure d'impact ; Risque attribuable ; Risque attribuable séquentiel ; Stratification ; Régression ; Causalité

---

## 1. Introduction

A major aim of epidemiologic research is to measure disease occurrence in relation to various character-

---

*E-mail address:* [jacques.benichou@chu-rouen.fr](mailto:jacques.benichou@chu-rouen.fr).

istics such as exposure to environmental, occupational, or lifestyle risk factors, genetic traits or other features. The generic term exposure will be used throughout this chapter to denote these characteristics. We will start with reviewing various measures that are at the root of quantitative epidemiologic thinking. These include measures that quantify disease occurrence, associations between disease occurrence and exposures as well as their consequences in terms of disease risk (Section 2). Emphasis will be placed on measures based on occurrence of new disease cases, referred to as disease incidence. Measures based on disease prevalence, i.e., considering previously existing disease cases as well as new cases, will be mentioned only in passing.

In Section 3–7, we will focus on the main measure of impact at the population level, namely attributable risk. This measure will be introduced in some detail in Section 3. Then, we will successively review three specific problems regarding attributable risk. First, we will consider adjusted attributable risk estimation from epidemiologic study data in Section 4, an issue that has generated intensive methodological research in the last 20 years, resulting in essentially satisfactory solutions. Second, we will discuss the lack of additivity of attributable risk contributions for separate exposures and present a possible solution in Section 5. Third, we will examine conceptual issues involved in interpreting attributable risk estimates in Section 6. Final remarks will follow in Section 7.

## 2. Rates, risks and measures of association

### 2.1. Incidence and hazard rates

The incidence rate of a given disease is the number of persons who develop the disease (number of incident cases) among subjects at risk of developing the disease in the source population over a defined period of time or age. Incidence rates are not interpretable as probabilities. While they have a lower bound of zero, they have no upper bound. Units of incidence rates are reciprocal of person-time, such as reciprocals of person-years or multiples of person-years (e.g., 100 000 person-years). For instance, if five cases develop from the follow-up of 50 subjects and for a total follow-up time of two years per subject, the incidence rate is  $5/100 = 0.05$  cases per person-year (assuming an instantaneous event with immediate recovery and all 50 subjects being at risk until the end of the observation period). Usually, incidence rates are assessed over relatively short time periods compared with the time scale for disease development, e.g., intervals of five years for chronic diseases

with an extended period of susceptibility, such as many cancers.

Synonyms for incidence rate are average incidence rate, force of morbidity, person-time rate, or incidence density [1], the last term reflecting the interpretation of an incidence rate as the density of incident case occurrences in an accumulated amount of person-time [2]. Mortality rates (overall or cause-specific) can be regarded as a special case of incidence rates, the outcome considered being death rather than disease occurrence.

Incidence rates can be regarded as estimates of a limiting theoretical quantity, namely the hazard rate,  $h(t)$ , also called the incidence intensity or force of morbidity. The hazard rate at time  $t$ ,  $h(t)$ , is the instantaneous rate of developing the disease of interest in an arbitrarily short interval  $\Delta$  around time  $t$ , provided that the subject is still at risk at time  $t$  (i.e., has not fallen ill before time  $t$ ). It has the following mathematical definition:

$$h(t) = \lim_{\Delta \downarrow 0} \Delta^{-1} \Pr(t \leq T < t + \Delta \mid t \leq T) \quad (1)$$

where  $T$  is the time period for the development of the disease considered and  $\Pr$  denotes probability. Indeed, for time intervals in which the hazard rate can be assumed constant, the incidence rate as defined above represents a valid estimate of the hazard rate. Thus, this result applies when piecewise constant hazards are assumed, which can be regarded as realistic in many applications, especially when reasonably short time intervals are used, and leads to convenient estimating procedures, e.g., based on the Poisson model.

Strictly speaking, incidence and hazard rates do not coincide. Hazard rates are formally defined as theoretical functions of time, whereas incidence rates are defined directly as estimates and constitute valid estimates of hazard rates under certain assumptions (see above).

From the definitions above, it ensues that individual follow-up data are needed to obtain incidence rates or estimate hazard rates. The cohort design that incurs follow-up of subjects with various profiles of exposure is the ideal design to obtain incidence or hazard rates for various levels or profiles of exposure, i.e., exposure-specific incidence or hazard rates. In many applications, obtaining exposure-specific incidence rates is not trivial, however. Indeed, several exposures are often considered, some with several exposed levels and some continuous. Moreover, it may be necessary to account for confounders or effect-modifiers. Hence, estimation often requires modelling. Alternatively to the cohort design, in the absence of individual follow-up data, person-time at risk can be estimated as the time period width times the population size at midpoint. Such estimation makes

the assumption that individuals who disappear from being at risk, either because they succumb, or because they move in or out, do so evenly across the time interval. Thus, population data such as registry data can be used to estimate incidence rates as long as an exhaustive census of incident cases can be obtained.

Case-control data pose a more difficult problem than cohort data, because case-control data alone are not sufficient to yield incidence or hazard rates. Indeed, they provide data on the distributions of exposure respectively in diseased subjects (cases) and non-diseased subjects (controls) for the disease under study, which can be used to estimate odds ratios (see Section 2.3) but are not sufficient to estimate exposure-specific incidence rates. However, it is possible to arrive at exposure-specific incidence rates from case-control data if case-control data are complemented by either follow-up data or population data, which happens for nested or population-based case-control studies. In a nested case-control study, the cases and controls are selected from a follow-up study. In a population-based case-control study, they are selected from a specified population in which an effort is made to identify all incident cases diagnosed during a fixed time interval, usually in a grouped form (i.e., number of cases and number of subjects by age group). In both situations, full information on exposure is obtained only for cases and controls. Additionally, complementary information on composite incidence (i.e., counts of events and person-time irrespective of exposure) can be sought from the follow-up or population data. By combining this information with odds ratio estimates, exposure-specific incidence rates can be obtained as has long been recognized [1,3–7] and is a consequence of the relation [6,8]:

$$h_0 = h^* (1 - AR) \tag{2}$$

where AR is the attributable risk in the population for all exposures considered, a quantity estimable from case-control data (see Sections 3–4),  $h_0$  is the baseline incidence rate, i.e., the incidence rate for subjects at the reference (unexposed) level of all exposures considered and  $h^*$  is the composite or average incidence rate in the population that includes unexposed subjects and subjects at various levels of all exposures (i.e., with various profiles of exposure). The composite incidence rate  $h^*$  can be estimated from the complementary follow-up or population data. Eq. (2) simply states that the incidence rate for unexposed subjects is equal to the proportion of the average incidence rate in the population that is not associated with any of the exposures considered. Eq. (2) can be specialized to various subgroups or strata defined by categories of age, sex or geographic location such as

region or centre, on which incidence rates are assumed constant. From the baseline rate  $h_0$ , incidence rates for all levels or profiles of exposure can be derived using odds ratio estimates, provided odds ratio estimates are reasonable estimates of incidence rate ratios as in the case of a rare disease (see Section 2.3). Consequently, exposure-specific incidence rates can be obtained from case-control data as long as they are complemented by follow-up or population data that can be used to estimate average incidence rates.

Finally, cross-sectional designs in which a sample of the population is assessed for both exposure and disease status cannot provide any assessment of incidence rates but instead will yield estimates of disease prevalence proportions.

Exposure-specific incidence and hazard rates play a central role in quantitative epidemiology because, as will be apparent from the following sections, all measures of the disease risk, association and impact can be derived from them.

## 2.2. Measures of the disease risk

The disease risk is defined as the probability that an individual who is initially disease-free will develop a given disease over a specified time or age interval (e.g., one year, five years, or lifetime).

If the interval starting at time  $a_1$  and ending just before time  $a_2$ , i.e.,  $[a_1, a_2)$ , is considered, the disease risk can be written formally as:

$$\pi(a_1, a_2) = \int_{a_1}^{a_2} h(a) \{S(a)/S(a_1)\} da \tag{3}$$

In Eq. (3),  $h(a)$  denotes the disease hazard at time or age  $a$  (see Section 2.1). The function  $S(\bullet)$ , with  $(\bullet)$  an arbitrary argument, is the survival function, so that  $S(a)$  denotes the probability of still being disease-free at time or age  $a$ , and  $S(a)/S(a_1)$  denotes the conditional probability of staying disease-free up to time or age  $a$  for an individual who is free of disease at the beginning of the interval  $[a_1, a_2)$ . Eq. (3) integrates over the interval  $[a_1, a_2)$  the instantaneous incidence rate of developing disease at time or age  $a$  for subjects still at risk of developing the disease (i.e., still disease-free subjects). Because the survival function  $S(\bullet)$  can be written as a function of the disease hazard through:

$$S(a_2)/S(a_1) = \exp \left\{ - \int_{a_1}^{a_2} h(a) da \right\} \tag{4}$$

the disease risk is also a function of the disease hazard.

By specializing the meaning of functions  $h(\bullet)$  and  $S(\bullet)$ , various quantities can be obtained that measure the disease risk in different contexts. First, the time scale on which these functions as well as the disease risk are defined corresponds to two specific uses of risk. In most applications, the first relevant time scale is age, since disease incidence is usually influenced by age. Note that by considering the age interval  $[0, a_2)$ , one obtains the lifetime disease risk up to age  $a_2$ . However, in clinical epidemiology settings, risk refers to the occurrence of an event, such as relapse or death in subjects already presenting with the disease of interest. In this context, the other relevant time scale becomes time from disease diagnosis or, possibly, time from some other disease-related event, such as a surgical resection of a tumour or occurrence of a first myocardial infarction.

Second, risk definition may account or not for individual exposure profiles. If no risk factors are considered to estimate the disease hazard, the corresponding measure of disease risk defines the average or composite risk over the entire population that includes subjects with various exposure profiles. This measure, also called cumulative incidence [1], may be of value at the population level. However, the main usefulness of risk is in quantifying an individual's predicted probability of developing disease depending on the individual's exposure profile. Thus, estimates of exposure-specific disease hazard have to be available for such exposure-specific risk (also called individualized or absolute risk) to be estimated.

Third, the consideration of competing risks and the corresponding definition of the survival function  $S(\bullet)$  yields two separate definitions of risk. Indeed, although risk is defined with respect to the occurrence of a given disease, subjects can die from other causes (i.e., competing risks), which obviously precludes disease occurrence. The first option is to define  $S(a)$  as the theoretical probability of being disease-free at time or age  $a$  if other causes of death (competing risks) could be eliminated, yielding a measure of the disease risk in a setting with no competing risks. This measure may not be of much practical value. Moreover, unless unverifiable assumptions regarding incidence of the disease of interest and deaths from other causes can be made, for instance assuming that they occur independently, the function  $S(\bullet)$  will not be estimable. For these reasons, it is more feasible to define  $S(a)$  as the probability that an individual will be alive and disease-free at age  $a$  as the second option, yielding a more practical definition of disease risk as the probability of developing disease in the presence of competing causes of death [9].

From the definition of the disease risk above, it appears that it depends on the incidence rate of disease in the population considered and can also be influenced by the strength of the relationship between exposures and disease, if individual risk is considered. One consequence is that risk estimates may not be portable from one population to another, as incidence rates may vary widely among populations that are separated in time and location or even among subgroups of populations, possibly because of differing genetic patterns or differing exposure to unknown risk factors. Additionally, competing causes of death (competing risks) may also have different patterns among different populations, which might also influence the values of the disease risk.

The disease risk is a probability and therefore lies between 0 and 1, and is dimensionless. A value of 0, while theoretically possible, would correspond to very special cases such as a purely genetic disease for an individual not carrying the disease gene. A value of 1 would be even more unusual and might again correspond to a genetic disease with a penetrance of 1 for a gene carrier but, even in this case, the value should be less than 1 if competing risks are accounted for.

Beside the term 'disease risk', 'absolute risk' or 'absolute cause-specific risk' have been used by several authors [10–14]. Alternative terms include 'individualized risk' [8], 'individual risk' [15], 'crude probability' [16], 'crude incidence' [17], 'cumulative incidence' [1, 18], 'cumulative incidence risk' [6], and 'absolute incidence risk' [1]. The term 'cumulative risk' refers to the quantity  $\int_{a_1}^{a_2} h(a) da$  and approximates disease risk closely in the case where disease is rare. The term 'attack rate' defines the risk of developing a communicable disease during a local outbreak and for the duration of the epidemic or the time during which primary cases occur [19 (Chapter 5), 20 (Chapter 27)].

Upon taking individual exposure profiles into account, resulting individual disease risk estimates are useful in providing an individual measure of the probability of disease occurrence, and can therefore be useful in counselling (e.g., in breast cancer, see [8,21–23]). Individual risk is also useful in designing (i.e., for sample size calculations and definition of eligibility criteria) and interpreting trials of interventions to prevent the occurrence of a disease through a risk–benefit analysis [24]. The concept of risk is also useful in clinical epidemiology as a measure of the individualized probability of an adverse event, such as a recurrence or death in diseased subjects. In that context, it can serve as a useful tool to help define individual patient management and, for instance, the absolute risk of recurrence in the next three years might be an important element in decid-

ing whether to prescribe an aggressive and potentially toxic treatment regimen [11,17].

As is evident from its definition, the disease risk can only be estimated and interpreted in reference to a specified age or time interval. One might be interested in short time spans (e.g., five years), or long time spans (e.g., 30 years). Of course, the disease risk increases as the time span increases. Sometimes, the time span is variable such as in lifetime risk. The disease risk can be influenced strongly by the intensity of competing risks (typically competing causes of death, see above). It varies inversely as a function of death rates from other causes.

It follows from its definition that the disease risk is estimable as long as hazard rates for the disease of interest are estimable. Therefore, it is directly estimable from cohort data, but case-control data have to be complemented with follow-up or population data in order to obtain the necessary complementary information on incidence rates (see Section 2.1).

Interpretation, usefulness, and properties of the disease risk, as well as methods for its estimation from cohort data, population-based or nested case-control data have been reviewed in detail [10].

### 2.3. Measures of association

Measures of association assess the strength of associations between one or several exposures and the risk of developing a given disease. Thus, they are useful in aetiologic research to assess and quantify associations between potential risk (or protective) factors and disease risk. The question addressed is whether and to what degree a given exposure is associated with the occurrence of the disease of interest. In fact, this is the primary question that most epidemiologic studies are trying to answer.

Depending on the available data, measures of association may be based on disease rates, disease risks, or even disease odds, i.e.,  $\pi/(1 - \pi)$ , with  $\pi$  denoting the disease risk. They contrast rates, risks, or odds for subjects with various levels of exposure, e.g., risks or rates of developing breast cancer for 40-year-old women with or without a personal history of benign breast disease. They can be expressed in terms of ratios or differences of risks or rates among subjects exposed and non-exposed to given factors or among subjects with various levels of exposure.

Measures of association can be defined for categorical or continuous exposures. For categorical exposures, any two exposure level can be contrasted using the measures of association defined below. However, it is con-

venient to define a reference level to which any exposure level can be contrasted. This choice is sometimes natural (e.g., non-smokers in assessing the association of smoking with disease occurrence), but can be more problematic if the exposure considered is of continuous nature, where a range of low exposures may be considered potentially inconsequential. The choice of a reference range is important for interpreting results. It should be wide enough for estimates of measures of association to be reasonably precise. However, it should not be so wide that it compromises meaningful interpretation of the results, which depend critically on the homogeneity of the reference level. For continuous exposures, measures of association can also be expressed per unit of exposure, e.g., for each additional gram of daily alcohol consumption. The reference level may then be a precise value such as no daily alcohol consumption or a range of values such as less than 10 grams of daily alcohol consumption.

When computing a measure of association, it is usually assumed that the relationship being captured has the potential to be causal, and efforts are taken to remove the impact of confounders from the quantity. Nonetheless, except for the special case of randomized studies, most investigators retain the word ‘association’ rather than ‘effect’ when describing the relationship between exposure and outcome to emphasize the possibility that unknown confounders may still influence the relationship.

Ratio-based measures of association are particularly appropriate when the effect of the exposure is multiplicative, which means that there is a similar percent increase or decrease associated with exposure in rate, risk or odds across exposure subgroups. Effects have often been observed to be multiplicative, leading to ratios providing a simple description of the association (e.g., see [25 (Chapter 2)]). Ratio measures are dimensionless and range from 0 to infinity, with 1 designating no association of the exposure with the outcome. When the outcome is death or disease, and the ratio has the rate, risk, or odds of the outcome with the exposed group in the numerator, a value less than 1 indicates a protective effect of exposure. The exposure is then referred to as a protective factor. When the ratio in this set-up is greater than 1, there is greater disease occurrence with exposure, and the exposure is then referred to as a risk factor.

The rate ratio is the ratio between the rate of disease among those exposed and those not exposed or  $h_E/h_{\bar{E}}$ . More generally, it can also contrast rates for various levels of exposure. Conceptually, the rate ratio is identical to a hazard ratio  $HR$ . The latter term tends to be used when time dependence of the rate is emphasized, as the

hazard is a function that may depend on time. The situation of a constant rate ratio over time is referred to as proportional hazards. The proportional hazards assumption is often made in the analysis of rates. Theoretically, the hazard ratio at a given time point is the limiting value of the rate ratio as the time interval around the point becomes very short, just as the hazard is the limiting quantity for the incidence rate (see Section 2.1). The rate ratio has also been called the Incidence Density Ratio [26 (Chapter 8)].

Rate ratios refer to population dynamics, and are not as easily interpretable on the individual level. It has been argued, however, that rate ratios make more sense than risk ratios when the period subjects are at risk is longer than the observation period [26 (Chapter 8)]. Numerically, the rate ratio is further from the null than the risk ratio. When rates are low, the similarity of risks and rates leads to rate ratios being close to risk ratios, as discussed below. Further considerations of how the rate ratio relates to other ratio-based measures of association are offered by Rothman and Greenland [20 (p. 50)].

The risk ratio, relative risk or ratio of risks of disease among those exposed  $\pi_E$  and those not exposed  $\pi_{\bar{E}}$ ,  $RR = \pi_E/\pi_{\bar{E}}$ , has been viewed as the gold standard among measures of association for many years. It is interpretable on the individual level as a given-fold increase in risk of disease. Like other ratio-based measures, it tends to be more stable than the risk difference across population groups at widely different risk (see below). However, similar to rate ratios and odds ratios (see below), the risk ratio can be viewed as misleading in the public eye when the risk among both the unexposed and the exposed is very low, yet many-fold increased by exposure. The risk ratio depends on the length of the time interval considered, because risk itself refers to a specific interval (see Section 2.2). In the literature, the term relative risk is often used to denote the rate ratio as well as the risk ratio, creating some confusion.

For several reasons, the odds ratio has emerged as the most popular measure of association. The odds ratio among those exposed and not exposed is the ratio of odds,  $OR = [\pi_E/(1 - \pi_E)]/[\pi_{\bar{E}}/(1 - \pi_{\bar{E}})]$ . Historically, the odds ratio has been considered an approximation of the risk ratio obtainable from case-control studies. The reason for this is that the probabilities of being sampled into case and control groups cancel in the calculation of the odds ratio, as long as sampling is independent of exposure status. Furthermore, when  $\pi_E$  and  $\pi_{\bar{E}}$  are small, the ratio  $(1 - \pi_{\bar{E}})/(1 - \pi_E)$  has little influence on the odds ratio, making it approximately equal to the risk ratio  $\pi_E/\pi_{\bar{E}}$ . The assumption of small

Table 1

Data from the fictitious cohort study

	Exposed	Unexposed
Diseased	40	20
Non-diseased	60	80

$\pi_E$  and  $\pi_{\bar{E}}$  is referred to as the rare-disease assumption. Kleinbaum et al. [26] have pointed out that in a case-control study of a stable population with incident cases and controls being representative of non-cases, the odds ratio is the rate ratio.

It can be shown that numerically the odds ratio falls the furthest from the null, and the risk ratio the closest, with the rate ratio in between. For example, from Table 1, based on fictitious data from a cohort study for a disease that is not rare, we would obtain a risk ratio  $\widehat{RR} = 0.4/0.2 = 2.00$  and an odds ratio  $\widehat{OR} = [(40)(80)]/[(20)(60)] = 2.67$ . If we assume a constant hazard, so that the risk for each group is  $1 - \exp(-hT)$ , with  $T$  being the follow-up time for each subject, we have the rate ratio  $\widehat{HR} = \ln(1 - 0.4)/\ln(1 - 0.2) = 2.29$ . Hence  $1 < \widehat{RR} < \widehat{HR} < \widehat{OR}$ .

The difference in magnitude between the above ratio measures is important to keep in mind when interpreting them for diseases or outcomes that are not rare. For rare outcomes, the values of the three ratio measures tend to be close.

Difference-based measures are appropriate when effects are additive (e.g., see [25 (Chapter 2)]), which means that the exposure leads to a similar absolute increase or decrease in rate or risk across subgroups. Although additive relationships may be less common in practice, difference measures may be more understandable to the public when the outcome is rare, and relate directly to measures of impact discussed in Section 6.

The numerical ranges of difference measures depend on their component parts. The rate difference ranges from minus to plus infinity, while the risk difference is bounded between minus and plus one. The situation of no association is reflected by a difference measure of zero. When the measure is formed as the rate or risk among the exposed minus that among the non-exposed, a positive value indicates that the exposure is a risk factor, while a negative value indicates that it is a protective factor. It can be shown that the risk difference falls numerically nearer to the null than the rate difference does. For example, Table 1 yields a risk difference of  $0.40 - 0.20 = 0.20$ , while the rate difference is  $-\ln(1 - 0.40) + \ln(1 - 0.20) = 0.29$ . However, they will be close for rare outcomes.

The rate difference for exposed an unexposed subjects is defined as  $h_E - h_{\bar{E}}$ , and has been commonly

employed to compare mortality rates and other demographic rates between countries, time periods and/or regions. In such comparisons, the two rates being compared are often directly standardized to the age and sex distribution of a standard population chosen, e.g., as the population of a given country in a given census year.

For the special case of a dichotomous exposure, the rate difference, i.e., the difference between the incidence rates in the exposed and unexposed subjects has been termed ‘excess incidence’ [19,27,28], ‘excess risk’ [29], ‘Berkson’s simple difference’ [30], ‘incidence density difference’ [1], or even ‘attributable risk’ [29,31], which may have caused some confusion.

The risk difference  $\pi_E - \pi_{\bar{E}}$  is parallel to the rate difference and similar considerations apply. Due to the upper and lower limits of plus, minus one on risk, but not on rate, risk differences are more difficult to model than rate differences. The odds difference  $\pi_E/(1 - \pi_E) - \pi_{\bar{E}}/(1 - \pi_{\bar{E}})$  has virtually never been used in practice.

Because exposure-specific incidence rates and risks can be obtained from cohort data, all measures of association considered (based on ratios or differences) can be obtained as well. This is also true of case-control data complemented by follow-up or population data (see Sections 2.1 and 2.2). Case-control data alone allow estimation of odds ratios thanks to the identity between disease and exposure odds ratios, which extends to the logistic regression framework. Prentice and Pyke [32] showed that the unconditional logistic model (see also [25 (Chapter 6)]) applies to case-control data as long as the intercept is disregarded. Interestingly, time-matched case-control studies allow estimation of hazard rates (e.g., see [1,33,34]).

Measures of association have a long history as methods for estimation and statistical inference. Traditional methods adjust for confounders by direct or indirect standardization of the rates or risks involved, prior to computation of the measure of association, or by stratification, where association measures are computed separately for subgroups and then combined. For measures based on the difference of rates or risks, direct standardization and stratification can be identical, if the same weights are chosen [35]. Generally, however, direct standardization uses predetermined weights chosen for external validity, while optimal or efficient weights are chosen with stratification. Efficient weights make the standard error of the combined estimator as small as possible.

In modern epidemiology, measures of association are most often estimated from regression analysis. Regression adjustment is a form of stratification, which provides more flexibility, but most often relies on large

sample size for inference. The function applied to the rate or risk in a regression analysis is referred to as the link function in the framework of generalized linear models underlying such analyses (see [36,37] for theory and practical application). For example, linear regression would regress the risk or rate directly on exposure without any transformation, which is referred to as using the identity link. When the exposure is the only predictor in such a model, all link functions fit equally well and simply represent different ways to characterize the association. However, when several exposures or confounders are involved, or if the exposure is measured as a continuous or ordinal variable, some link functions and not others may require interaction or non-linear terms to improve the fit. Most widely used regression models are the Poisson and Cox models for rate ratio estimation from cohort data, the log linear model for risk ratio estimation from cohort data, the logistic regression model for odds ratio estimation from cohort or case-control data.

Measures of association based on prevalence parallel those for risk (for point prevalence) or incidence rates (for period prevalence). For example, one can form prevalence ratios, prevalence differences and prevalence odds ratios. They can be estimated from cross-sectional data. These measures are less useful for studying the aetiology of a disease than measures based on incidence. The reason for this is that prevalence reflects both incidence and duration of disease. For a potentially fatal or incurable disease, duration means survival and the exposures that increase incidence may reduce or increase survival, and hence the association of an exposure with prevalence may be very different from its association with incidence.

Measures of associations and related methods of inference are reviewed at length in epidemiologic textbooks (e.g., [20,25,26,29,38–43]).

### 3. Measures of impact: attributable risk

Measures of impact are used to assess the contribution of one or several exposures to the occurrence of incident cases at the population level. Thus, they are useful in public health to weigh the impact of exposure on the burden of disease occurrence and assess potential prevention programmes aimed at reducing or eliminating exposure in the population. The most commonly used measure of impact is the attributable risk.

The term ‘attributable risk’ (AR) was initially introduced by Levin in 1953 [44] as a measure to quantify the impact of smoking on lung cancer occurrence. Gradually, it has become a widely used measure to assess

the consequences of an association between an exposure factor and a disease at the population level. It is defined as the following ratio:

$$AR = \{\Pr(D) - \Pr(D|\bar{E})\} / \Pr(D) \quad (5)$$

The numerator contrasts the probability of disease,  $\Pr(D)$ , in the population, which may have some exposed,  $E$ , and some unexposed,  $\bar{E}$ , individuals, with the hypothetical probability of disease in the same population but with all exposure eliminated  $\Pr(D|\bar{E})$ . Thus, it measures the additional probability of disease in the population that is associated with the presence of exposure in the population, and AR measures the corresponding proportion. Probabilities in Eq. (5) will usually refer to the disease risk although, depending on the context, they may be replaced with incidence rates.

Unlike measures of association (see Section 2.3), AR depends both on the strength of the association between exposure and disease and the prevalence of exposure in the population,  $p_E$ . This can be seen for instance through rewriting AR from Eq. (5). Upon expressing  $\Pr(D)$  as:

$$\Pr(D|E)p_E + \Pr(D|\bar{E})p_{\bar{E}} \quad \text{with } p_{\bar{E}} = 1 - p_E$$

both in the numerator and the denominator, and noting that:

$$\Pr(D|E) = RR \times \Pr(D|\bar{E})$$

the term  $\Pr(D|\bar{E})$  cancels out and AR is obtained as [6, 45]:

$$AR = \{p_E(RR - 1)\} / \{1 + p_E(RR - 1)\} \quad (6)$$

a function of both the prevalence of exposure in the population,  $p_E$ , and the rate ratio or relative risk,  $RR$ .

An alternative formulation underscores this joint dependency in yet another manner. Upon using the same decomposition of  $\Pr(D)$  as above and again noting that:

$$\Pr(D|E) = RR \times \Pr(D|\bar{E})$$

the numerator in Eq. (5) can be rewritten as:

$$p_E \Pr(D|E) - p_E \Pr(D|\bar{E}) / RR$$

From using Bayes' theorem to express  $\Pr(D|E)$  as  $\Pr(E|D)\Pr(D)/p_E$ , it becomes equal to:

$$\Pr(D)p_{E|D}(1 - 1/RR)$$

after simple algebra. This yields [6]

$$AR = p_{E|D}(RR - 1) / RR \quad (7)$$

a function of the prevalence of exposure in diseased individuals,  $p_{E|D}$  ( $= \Pr(E|D)$ ), and the rate ratio or relative risk,  $RR$ .

A high relative risk can correspond to a low or high AR, depending on the prevalence of exposure, which leads to widely different public-health consequences. One implication is that portability is not a usual property of AR, as the prevalence of exposure may vary widely among populations that are separated in time or location. This is in contrast with measures of association such as the relative risk or rate ratio, which are more portable from one population to another, as the strength of the association between disease and exposure might vary little among populations, unless strong interactions with environmental or genetic factors are present.

When the exposure considered is a risk factor ( $RR > 1$ ), it follows from the above definition that AR lies between 0 and 1. Therefore, it is very often expressed as a percentage. AR increases both with the strength of the association between exposure and disease measured by  $RR$ , and with the prevalence of exposure in the population. A prevalence of 1 (or 100%) yields a value of AR equal to the attributable risk among the exposed individuals, i.e.,  $(RR - 1)/RR$  (see Section 6). AR approaches 1 for an infinitely high  $RR$ , provided that the exposure is present in the population (i.e., non-null prevalence of exposure).

AR takes a null value when either there is no association between exposure and disease ( $RR = 1$ ) or no subject is exposed in the population. Negative AR values are obtained for a protective exposure ( $RR < 1$ ). In this case, AR varies between 0 and  $-\infty$ , a scale on which AR lacks a meaningful interpretation. One solution is to reverse the coding of exposure (i.e., interchange exposed and unexposed categories) to go back to the situation of a positive AR, sometimes called the preventable fraction in this case [46–48]. Alternatively, one may consider a different parameter, namely the prevented fraction [6,30,46,47,49].

Some confusion in the terminology arises from the reported use of as many as 16 different terms in the literature to denote the attributable risk [50,51]. However, a literature search by Uter and Pfahlberg [52] found some consistency in terminology usage, with 'attributable risk' and 'population attributable risk' [19] being the most commonly used terms, by far followed by 'etiologic fraction' [6]. Other popular terms include 'attributable risk percentage' [45], 'fraction of aetiology' [6], and 'attributable fraction' [20 (Chapter 4), 48,53, 54].

Moreover, additional confusion may originate in the use by some authors [19,29,31] of the term 'attributable risk' to denote a measure of association, the excess incidence, that is the difference between the incidence rates in exposed and unexposed subjects (see Section 2.3).



Context will usually help the readers to detect this less common use.

While measures of association such as the rate ratio and relative risk are used to establish an association in aetiologic research, AR has a public-health interpretation as a measure of the disease burden attributable or at least related to one or several exposures. Consequently, AR is used to assess the potential impact of prevention programmes aimed at eliminating exposure from the population. It is often thought of as the fraction of disease that could be eliminated if exposure could be totally removed from the population.

However, this interpretation can be misleading because, for it to be strictly correct, the three following conditions have to be met [30]. First, estimation of AR has to be unbiased (see Section 4). Second, exposure has to be causal rather than merely associated with the disease. Third, elimination of exposure has to be without any effect on the distribution of other risk factors. Indeed, as it might be difficult to alter the level of exposure to one factor independently of other risk factors, the resulting change in disease load might be different from the AR estimate. For these reasons, various authors elect to use weaker definitions of AR, such as the proportion of disease that can be related or linked, rather than attributable, to exposure [6].

Several authors have considered an interpretation of AR in terms of aetiologic research. The argument is that if an AR estimate is available for several risk factors jointly, then its complement to 1, i.e.,  $1 - AR$ , must represent a gauge of the proportion of disease cases not explained by the risk factors used in estimating AR. Hence,  $1 - AR$  would represent the proportion of cases attributable to other (possibly unknown) risk factors. For instance, it was estimated that the AR of breast cancer was 41% for late age at first birth, nulliparity, family history of breast cancer and higher socioeconomic status, which suggested that at least 59% of cases had to be attributable to other risk factors [55]. A similar type of reasoning was used in several well-known reports of estimated percentages of cancer death or incidence attributable to various established cancer risk factors (e.g., smoking, diet, occupational exposure to carcinogens...). Some of these reports conveyed the impression that little remained unexplained by factors other than the main established preventable risk factors and that cancer was a mostly preventable illness [56–60]. Such interpretation has to be taken with great care since ARs for different risk factors may add to more than 100% because multiple exposures are usually possible (e.g., smoking and occupational exposure to asbestos). Moreover, this interpretation can be refuted based on logical arguments

regarding the fact that disease occurrence may require more than one causal factor. Furthermore, one can note that once a new risk factor is considered, the joint unexposed reference category changes from lack of exposure to all previously considered risk factors to lack of exposure to those risk factors *and* the new risk factor [61]. Because of this change in the reference category, the AR for the new risk factor may surpass the quantity  $1 - AR$  for the previously considered risk factors. Thus, while it is useful to know that only 41% of breast cancer cases can be attributed to four established risk factors in the above example, it is entirely conceivable that new risk factors of breast cancer may be elicited, which yield an AR of more than 59% by themselves in the above example.

AR can be estimated from cohort studies since all quantities in Eqs. (5)–(7) are directly estimable from cohort studies. AR estimates can differ depending on whether rate ratios, risk ratios or odds ratios are used, but will be numerically close for rare diseases. For case-control studies, exposure-specific incidence rates or risks are not available, unless data are complemented with follow-up or population-based data (see Sections 2.1 and 2.2). Thus, one has to rely on odds ratio estimates, use Eq. (6) and estimate  $p_E$  from the proportion exposed in the controls, making the rare-disease assumption also involved in estimating odds ratios rather than relative risks. Alternatively, one can use Eq. (7), in which the quantity  $p_{E|D}$  can be directly estimated from the diseased individuals (cases) and  $RR$  can be estimated from the odds ratio. It is possible to form a measure equivalent to AR based on prevalence, e.g., by substituting a prevalence odds ratio for  $RR$  in Eq. (7). This quantity can be estimated from cross-sectional data. It will be useful in weighing the impact of exposure on the overall burden of disease rather than disease occurrence and will be influenced by disease length or survival (see Section 2.3). Detailed reviews of estimability and basic estimation of AR for various epidemiologic designs can be found in Walter [30] and Benichou [62,63], who provide explicit formulae for point and standard errors estimates and consider various mathematical transformations on which to base confidence intervals.

Beside AR, other measures of impact have been proposed, notably the generalized impact fraction and the number of person-years (or potential years) of life lost. The generalized impact fraction broadens the concept of AR and is obtained by replacing the term  $\Pr(D|\bar{E})$  in Eq. (5) by  $\Pr^*(D)$ , the probability of disease under a modified distribution of exposure. Thus, it can be interpreted as the fractional reduction of disease occurrence that would result from changing the current

distribution of exposure in the population to some specific modified distribution characterized by reduced exposure, with AR corresponding to the special case of a modified distribution putting unit mass on the lowest risk configuration (i.e., with all exposure eliminated) [64–66]. For a potentially fatal or incurable disease, the number of person-years of life lost related to a given exposure is a measure defined as the difference between current life expectancy of the population and potential life expectancy with the exposure eliminated [67–69]. This quantity may be difficult to estimate. Indeed, several causes of death may have to be considered, e.g., lung cancer and pleural as well as peritoneal mesothelioma for asbestos exposure. Moreover, this measure depends on the prevalence of exposure in the population and strength of association between exposure and disease(s), but also on the age-distribution of exposure-associated diseases and their severity, i.e. case fatality.

#### 4. Adjusted attributable risk estimation

As it is the case for measures of association, unadjusted (or crude or marginal) AR estimates may be inconsistent [6,30,66,70]. The precise conditions under which adjusted AR estimates that take into account the distribution and effect of other factors will differ from unadjusted AR estimates that fail to do so were worked out by Walter [66]. If  $E$  and  $X$  are two dichotomous factors taking levels 0 and 1, and if one is interested in estimating the AR for exposure  $E$ , then the following applies. The adjusted and unadjusted AR estimates coincide (i.e., the crude AR estimate is unbiased) if and only if (a)  $E$  and  $X$  are such that  $\Pr(E = 0, X = 0) \bullet \Pr(E = 1, X = 1) = \Pr(E = 0, X = 1) \bullet \Pr(E = 1, X = 0)$ , which amounts to the independence of their distributions, or (b) exposure to  $X$  alone does not increase disease risk, namely  $\Pr(D|E = 0, X = 1) = \Pr(D|E = 0, X = 0)$ . When considering one (or several) polychotomous factor(s)  $X$  forming  $J$  levels ( $J > 2$ ), conditions (a) and (b) can be extended to a set of analogous sufficient conditions. Condition (a) translates into a set of  $J(J - 1)/2$  conditions for all pairs of levels  $j$  and  $j'$  of  $X$ , amounting to an independent distribution of  $E$  and all factors in  $X$ . Condition (b) translates into a set of  $J - 1$  conditions stating that in the absence of exposure to  $E$ , exposure to any of the other factors in  $X$ , alone or in combination, does not increase the disease risk.

The extent of bias varies according to the severity of the departure from conditions (a) and (b) above. Although no systematic numerical study of the bias of unadjusted AR estimates has been performed, Walter [66] provided a revealing example of a case-control study as-

sessing the association between alcohol, smoking, and oral cancer. In that study, severe positive bias was observed for crude AR estimates, with a very large difference between crude and adjusted AR estimates both for smoking (51.3% vs. 30.6%, a 20.7 difference in percentage points and 68% relative difference in AR estimates) and alcohol (52.2% vs. 37.0%, a 15.2% absolute difference and 48% relative difference). Thus, the prudent approach must be to adjust for factors that are suspected or known to act as confounders in a similar fashion as for estimating measures of associations.

Two simple adjusted estimation approaches discussed in the literature are inconsistent. The first approach was presented by Walter [30], and is based on a factorization of the crude risk ratio into two components, similar to those in Miettinen's earlier derivation [71]. In this approach, a crude AR estimate is obtained under the assumption of no association between exposure and disease (i.e., values of  $RR$  or the odds ratio are taken equal to 1 separately for each level of confounding). This term reflects the AR only due to confounding factors since it is obtained under the assumption that disease and exposure are not associated. By subtracting this term from the crude AR estimate that ignores confounding factors and thus reflects the impact of both exposure and confounding factors, what remains is an estimate of the AR for exposure adjusted for confounding [30]. The second approach is based on using Eq. (6) and plugging in a common adjusted  $RR$  estimate (odds ratio estimate in case-control studies), along with an estimate of  $p_E$  [45,71,72]. Both approaches, while intuitively appealing, were shown to be inconsistent [72–74] and, accordingly, very severe bias was exhibited in simulations of cross-sectional and cohort designs [51].

By contrast, two adjusted approaches based on stratification yield valid estimates. The Mantel–Haenszel approach consists in plugging-in an estimate of the common adjusted  $RR$  (odds ratio in case-control studies) and an estimate of the prevalence of exposure in diseased individuals,  $p_{E|D}$ , in Eq. (7) in order to obtain an adjusted estimate of AR [47,75–78]. In doing so, it is possible to adjust for one or more polychotomous factors forming  $J$  levels or strata. While several choices are available for a common adjusted  $RR$  or odds ratio estimator, a usual choice is to use a Mantel–Haenszel estimator of  $RR$  in cohort studies [20 (Chapters 15–16), 26 (Chapters 9 and 17), 79,80] or odds ratio in case-control studies [20 (Chapters 15–16), 25 (Chapters 4–5), 26 (Chapters 9,17), 79,81]. For this reason, the term ‘Mantel–Haenszel approach’ has been proposed to denote this approach to adjusted AR estimation [82]. When there is no interaction between exposure and fac-

tors adjusted for, Mantel–Haenszel-type estimators of *RR* or odds ratio have favourable properties, as they combine lack of (or very small) bias even for sparse data (e.g., individually matched case-control data) and good efficiency, except in extreme circumstances [25 (Chapters 4–5), 79,83–85]. Moreover, variance estimators are consistent even for sparse data (‘dually-consistent’ variance estimators) [47,86]. Simulation studies of cohort and case-control designs [47,77,78,87] showed that adjusted AR estimates are little affected by small-sample bias when there is no interaction between exposure and adjustment factors, but can be misleading if such interaction is present.

The weighted-sum approach also allows adjustment for one or more polychotomous factors forming *J* levels or strata. The adjusted AR is written as a weighted sum over all strata of stratum-specific ARs, i.e.,  $\sum_{j=1}^J w_j \text{AR}_j$  [30,88,89]. Using crude estimates of  $\text{AR}_j$  separately within each stratum *j* and setting weights  $w_j$  as proportions of diseased individuals (cases) yield an asymptotically unbiased estimator of AR, which can be seen as a maximum-likelihood estimator [88]. This choice of weights defines the ‘case-load method’. The weighted-sum approach does not require the assumption of a common relative risk or odds ratio. Instead, the relative risks or odds ratios are estimated separately for each adjustment level with no restrictions placed on them, corresponding to a fully saturated model for exposure and adjustment factors (i.e., a model with all interaction terms present). From these separate relative risk or odds ratio estimates, separate AR estimates are obtained for each level of adjustment. Thus, the weighted-sum approach not only accounts for confounding, but also for interaction. However, simulation studies of cohort and case-control designs [77,78, 87,88] show that the weighted-sum approach can be affected by small sample bias, sometimes severely. Hence, it should be avoided when analyzing sparse data, and should not be used altogether for analyzing individually matched case-control data.

A natural alternative to generalize these approaches is to use adjustment procedures based on regression models, in order to take advantage of their flexible and unified approach to efficient parameter estimation and hypothesis testing. Regression models allow one to take into account adjustment factors as well as interaction of exposures with some or all adjustment factors. This approach was first used by Walter [30], Sturmans et al. [90] and Fleiss [91] followed by Deubner et al. [92] and Greenland [47]. The full generality and flexibility of the regression approach was exploited by Bruzzi et al. [93], who developed a general AR estimate based on

rewriting AR as:

$$1 - \sum_{j=1}^J \sum_{i=0}^I \rho_{ij} RR_{i|j}^{-1}$$

Quantities  $\rho_{ij}$  represent the proportion of diseased individuals with level *i* of exposure (*i* = 0 at the reference level, *i* = 1, ..., *I* for exposed levels) and *j* of confounding and can be estimated from cohort or case-control data (or cross-sectional survey data) using the observed proportions. The quantity  $RR_{i|j}^{-1}$  represents the inverse of the rate ratio, risk ratio or odds ratio, depending on the context, for level *i* of exposure at level *j* of confounding. It can be estimated from regression models, both for cohort and case-control data (as well as cross-sectional data), which allows confounding and interactions to be accounted for. Hence, this regression-based approach to AR estimation allows control for confounding and interaction, and can be used for the main epidemiologic designs. Depending on the design, conditional or unconditional logistic, log-linear or Poisson models can be used. Variance estimators were developed based on an extension of the delta-method to implicitly related random variables in order to take into account the variability in estimates of terms  $\rho_{ij}$  and  $RR_{i|j}^{-1}$  as well as their correlations [94–96]. This regression approach includes the crude and two stratification approaches as special cases and offers additional options [82]. The unadjusted approach corresponds to models for  $RR_{i|j}^{-1}$  with exposure only. The Mantel–Haenszel approach corresponds to models with exposure and confounding factors, but no interaction terms between exposure and adjustment factors. The weighted-sum approach corresponds to fully saturated models with all interaction terms between exposure and confounding factors. Intermediate models are possible, for instance models allowing for interaction between exposure and only one confounder, or models in which the main effects of some confounders are not modelled in a saturated way.

A modification of the approach by Bruzzi et al. was developed by Greenland and Drescher [97] in order to obtain full maximum likelihood estimates of AR. The modification consists in estimating the quantities  $\rho_{ij}$  from the regression model rather than simply relying on the observed proportions of cases. The two model-based approaches seem to differ very little numerically [97]. Greenland and Drescher’s approach might be more efficient in small samples, although no difference was observed in simulations of the case-control design, even for samples of 100 cases and 100 controls [97]. It might be less robust to model misspecification, however, as it

Table 2

Illustration of the phenomenon of non-additivity of attributable risks for two exposures  $E_1$  and  $E_2$  and multiplicative risks

Exposure to factor $E_1$	Exposure to factor $E_2$	Prevalence	Relative risk	Risk	Risk in the absence of factor $E_1$	Risk in the absence of factor $E_2$
Yes	Yes	0.25	81	0.81	0.09	0.09
Yes	No	0.25	9	0.09	0.01	0.09
No	Yes	0.25	9	0.09	0.09	0.01
No	No	0.25	1	0.01	0.01	0.01

relies more heavily on the  $RR$  or the odds ratio regression model used. Finally, it does not apply to the conditional logistic model, and if that model is to be used (notably, in case-control studies with individual matching), the original approach of Bruzzi et al. is the only possible choice.

Detailed reviews of adjusted AR estimation [63,82,87,98] are available. Alternative methods to obtain estimates of variance and confidence intervals for AR have been developed either based on resampling techniques [52,87,99–102] or on quadratic equations [103–105].

### 5. Non-additivity of attributable risks for separate exposures

AR is frequently estimated in multifactorial situations when trying to evaluate the joint and individual impact of multiple exposures. In this context, separate ARs can be estimated for each exposure as well as the overall AR for all or several exposures jointly. This raises a problem since individual contributions of exposures to attributable risk are usually non-additive.

Indeed, Walter [70] showed that the sum of separate ARs for each exposure is not equal to the joint AR unless at least one of the two following specific conditions is fulfilled: there is no joint exposure to the different exposures in the population and the effects of the exposures on disease risk are additive. For two exposures, the latter condition means that the relative risk for exposure to the two factors,  $RR_{12}$ , is linked to the relative risks for exposures 1 and 2 separately,  $RR_1$  and  $RR_2$  respectively, by the formula  $(RR_{12} - 1) = (RR_1 - 1) + (RR_2 - 1)$ . If none of the two conditions above is verified, then the sum of the separate for each exposure differs from the joint AR and the difference can be very substantial.

Table 2 taken from Begg [61] illustrates this problem. It considers two dichotomous exposures  $E_1$  and  $E_2$  with prevalence in the population of 0.25 for each of the four joint categories. Each of these exposures multiplies the disease risk by 9 with a joint multiplicative effect such that the risk is multiplied by 81 in the case of joint exposure to  $E_1$  and  $E_2$ . By using either Eq. (5) or Eq. (6), one obtains an 80% AR for both factors  $E_1$  and  $E_2$  sep-

arately. Indeed, with Eq. (6), for example, the AR for factor  $E_1$  is:

$$AR_1 = 0.50 \times (9 - 1) / \{1 + 0.50 \times (9 - 1)\} = 0.80$$

i.e.,  $AR_1 = 80\%$ . The same applies to  $E_2$ , since the problem is perfectly symmetrical in this particular case. Thus, the sum of the separate ARs for exposures  $E_1$  and  $E_2$ , i.e.,  $AR_1 + AR_2$ , cannot be equal to the joint attributable risk for factors  $E_1$  and  $E_2$ , since this sum is greater than 100%! The joint AR for exposures  $E_1$  and  $E_2$ ,  $AR_{12}$ , can be obtained by using Eq. (5), which is equivalent to:

$$AR_{12} = 1 - \Pr(D|\bar{E}) / \Pr(D)$$

where  $\Pr(D|\bar{E})$  is the risk of developing the disease in subjects that are neither exposed to  $E_1$  nor  $E_2$ , i.e. 0.01. The probability  $\Pr(D)$  represents the risk of developing the disease in the population. Upon taking into account the joint distribution of exposures  $E_1$  and  $E_2$  it is given by:

$$\Pr(D) = 0.25 \times (0.81 + 0.09 + 0.09 + 0.01) = 0.25$$

and the joint AR for exposures  $E_1$  and  $E_2$  is:

$$AR_{12} = 1 - 0.01/0.25 = 0.96$$

i.e.,  $AR_{12} = 96\%$ , which is clearly lower than the sum  $AR_1 + AR_2$ .

The non-additivity problem comes from the fact that by forming the sum  $AR_1 + AR_2$ , one is not considering the same reference levels as when considering the joint AR, namely  $AR_{12}$ . For the latter, the reference level is the category that corresponds to an absence of exposure to  $E_1$  and  $E_2$ . In the case of the AR for exposure  $E_1$ , i.e.,  $AR_1$ , the reference level corresponds to an absence of exposure to  $E_1$  only and therefore includes subjects both exposed and unexposed to  $E_2$  (in equal proportion in this example). Similarly, for  $AR_2$ , the reference level corresponds to an absence of exposure to  $E_2$  only, and therefore includes subjects both exposed and unexposed to  $E_1$  (again, in equal proportion in this example). This means that the contribution of the category of subjects exposed to both  $E_1$  and  $E_2$  is taken into account more

than once in the sum  $AR_1 + AR_2$ , which explains the inadequateness of calculating  $AR_1 + AR_2$ , except in the specific cases described by Walter [70].

Because non-additivity is somewhat counter-intuitive and generates misinterpretations, three alternative approaches have been suggested, one based on considering variance decomposition methods [106] rather than estimating AR, one based on estimating assigned share or probability of causation of a given exposure with relevance in litigation procedures for individuals with multiple exposures [107–113], and one based on an extension of the concept of AR [114,115]. This last approach relies on partitioning techniques [116,117] and keeps with the framework of AR estimation by introducing the sequential AR that generalizes the concept of AR. The principle is to define an order among the exposures considered. Then, the contribution of each exposure is assessed sequentially according to that order. The contribution of the first exposure considered is calculated as the standard AR for that exposure separately. The contribution of the second exposure is obtained as the difference between the joint AR estimate for the first two exposures and the separate AR estimate for the first exposure, the contribution of the third exposure is obtained as the difference between the joint AR estimates for the first three and first two exposures, etc. Thus, a multidimensional vector consisting of contributions of each separate exposure is obtained.

These contributions are meaningful in terms of potential prevention programmes that consider successive rather than simultaneous elimination of exposures from the population. Indeed, each step yields the additional contribution of the elimination of a given exposure once higher-ranked exposures are eliminated. At some point, additional contributions may become very small, indicating that there is not much point in considering extra steps. By construction, these contributions sum to the overall AR for all exposures jointly, which constitutes an appealing property. Of course, separate vectors of contributions are obtained for different orders. Meaningful orders depend on practical possibilities in implementing potential prevention programmes in a given population. Average contributions can be calculated for each given step (i.e., the first step, second step, etc.) by calculating the mean of contributions corresponding to that step over all possible orders. These average contributions have been termed partial ARs [114], and they represent another potentially useful measure.

For the data in Table 2, sequential ARs would be equal to 80% and  $96 - 80 = 16\%$  for exposures ranked first and second respectively, and partial ARs would be equal to  $(80 + 16)/2 = 48\%$  for each factor.

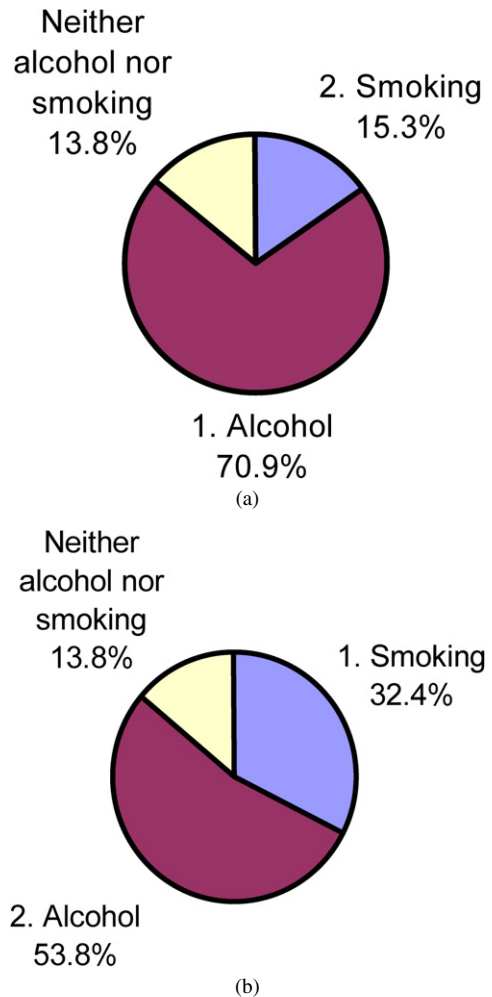


Fig. 1. Sequential attributable risk estimates for elevated alcohol consumption ( $40 + \text{g/day}$ ) and heavy smoking ( $10 + \text{g/day}$ ) for two different orders of removal (top panel (a): alcohol, then smoking; bottom panel (b): smoking, then alcohol) – Case-control data on oesophageal cancer [119].

Methods for visualizing sequential and partial ARs are provided by Eide and Heuch [118]. An illustration is given by Fig. 1 based on data from the case-control study of oesophageal cancer conducted in the Ille-et-Vilaine district of France. This study included 200 cases and 775 controls selected by simple random sampling from electoral lists [119]. The assessment of associations between alcohol consumption and smoking with oesophageal cancer has been the focus of detailed illustration by Breslow and Day [25], who presented various approaches to odds ratio estimation. Upon considering 0–39 g/day as the reference category for alcohol consumption, 29 cases and 386 controls were in the reference category, while 171 cases and 389 controls were in the exposed (i.e.,  $40 + \text{g/day}$ ) category.

The corresponding crude odds ratio was estimated as  $(171 \times 386)/(29 \times 389) = 5.9$  and the crude AR as 70.9% from Eq. (6). Upon considering 0–9 g/day as the reference category for smoking, 78 cases and 447 controls were in the reference category, while 122 cases and 328 controls were in the exposed (i.e., 9+ g/day) category. The corresponding crude odds ratio was estimated as  $(122 \times 447)/(78 \times 328) = 2.1$  and the crude AR as 32.4% from Eq. (6). Moreover, there were nine cases and 252 controls in the joint reference level of alcohol consumption and smoking (i.e., 0–39 g/day of alcohol and 0–9 g/day of tobacco), which yielded a crude joint odds ratio estimate of 10.2 and a crude joint AR estimate for drinking at least 40 g/day of alcohol or smoking at least 10 g/day of tobacco of 86.2%.

Hence, considering the first order of risk factor removal (i.e., eliminating alcohol consumption above 39 g/day followed by eliminating smoking above 9 g/day) yields sequential AR estimates of 70.9% for elevated daily alcohol consumption and  $(86.2 - 70.9)\% = 15.3\%$  for heavy smoking, so that, once elevated alcohol consumption is eliminated, the additional impact of eliminating heavy smoking appears rather limited (Fig. 1a). Considering the second order (i.e., eliminating heavy smoking first) yields sequential AR estimates of 32.4% for heavy smoking and  $(86.2 - 32.4)\% = 53.8\%$  for elevated alcohol consumption so that, once heavy smoking is eliminated, the additional impact of eliminating elevated alcohol consumption remains major (Fig. 1b). A summary of these results is provided by partial ARs for elevated alcohol consumption and heavy smoking, with estimated values of 62.4% and 23.9%, respectively, again reflecting the higher impact of elevated alcohol consumption on oesophageal cancer.

A detailed review of properties, interpretation, and variants of sequential and partial ARs was provided by Land et al. [115].

## 6. Conceptual issues in estimating attributable risk

As first pointed out by Greenland and Robins [53], there are three distinct measures within the concept of AR. It is easier to make this point using the formulation of AR in Eq. (7) and considering the attributable risk in the exposed individuals. In Eq. (7), AR is expressed as the product of two terms, the prevalence of exposure in diseased individuals,  $p_{E|D}$ , and the quantity  $(RR - 1)/RR$ . The latter term is equivalent to the attributable risk in the exposed ( $AR_E$ ) or attributable fraction in the exposed individuals, which is defined as the following ratio [6,19,44,45]:

$$AR_E = \{\Pr(D|E) - \Pr(D|\bar{E})\} / \Pr(D|E) \quad (8)$$

where  $\Pr(D|E)$  is the probability of disease in the exposed individuals and  $\Pr(D|\bar{E})$  is the hypothetical probability of disease in the same subjects, but with every exposure eliminated.  $AR_E$  can be regarded as a relative rate or risk difference (see Section 2.3). Depending on the context, probabilities in Eq. (8) will refer to the disease risk or the incidence rates (see Sections 2.1–2.2), as will be discussed further below. Hence, the logic behind Eq. (7) is that the impact of exposure at the population level depends on the impact among exposed subjects and the prevalence of exposure among disease cases.

Greenland and Robins [53] proposed to distinguish three separate measures of impact that are conceptually different. They made these distinctions for  $AR_E$ , but they apply equivalently to AR. The first measure, denoted ‘excess fraction’, corresponds to the definition of  $AR_E$ , with probabilities in Eq. (8) representing risks. The corresponding AR definition is that considered throughout this paper with probabilities in Eqs. (5)–(7) representing risks and the measure of association  $RR$  representing the risk ratio (i.e., the relative risk). The second measure denoted ‘excess rate’, ‘rate fraction’ [20 (Chapter 4)], or ‘assigned share’ in the risk-assessment literature [110,120], corresponds to the definition of  $AR_E$  with probabilities in Eq. (8) representing incidence rates or hazards. The corresponding AR definition is again that considered throughout this paper, but with probabilities in Eqs. (5)–(7) representing incidence rates or hazards, and the measure of association  $RR$  representing the rate ratio or hazard rate. Both AR (and  $AR_E$ ) measures are estimable from cohort data, but estimability from case-control data requires that case-control data be complemented with follow-up or population-based data (see Sections 2.1, 2.2, 3, and 4). Usually, these two AR measures will be substantially different numerically. For instance, from the data in Table 1 and for a 10% prevalence of the exposed category, the AR counterpart of the excess fraction will be estimated as 9.1% and the AR counterpart of the excess rate as 11.4%, a 25.7% relative difference. The two measures will be close only in special circumstances, for instance for a closed cohort and a rare disease [53]. Generally speaking, the AR counterpart of the rate fraction will be more useful than the AR counterpart of the excess fraction in situations where the question of when the disease occurred matters as much or more than the question of whether it occurred or not in a given delay. Both measures aim at assessing the impact of exposure by quantifying the increase in disease risk or rate caused by exposure, which amounts to considering only extra cases that would not have become cases in the absence of exposure. Although it is usually impossible to dis-

tinguish individual exposed cases for which exposure played an aetiologic role from those where exposure was irrelevant, such risk or rate increases are estimable from epidemiologic data.

The third AR measure is defined as the proportion of disease cases in which exposure played an aetiologic role either by contributing to disease occurrence through making the case's incidence time earlier than it would have been in the absence of exposure (i.e., disease would have occurred in the absence of exposure, only later) or by causing disease occurrence (i.e., disease would not have occurred in the absence of exposure). It is the AR counterpart to the 'aetiologic fraction' definition of  $AR_E$  proposed by Greenland and Robins [53]. Strong and usually unverifiable assumptions are needed to estimate it. Greenland and Robins [53] gave a theoretical example in which a 20-year leukaemia risk was estimated in a cohort of former military personnel who had all been exposed to radiation during a nuclear-weapon test. In this example, AR and  $AR_E$  were equivalent, because all subjects had been exposed. Overall, 24 cases developed including six cases in the last five years (i.e., 15–20 years after radiation exposure). Under one biologic assumption, the effect of radiation could be to accelerate leukaemia incidence by five years, in which case the six subjects who developed leukaemia in 15–20 years would have remained leukaemia-free at 20 years in the absence of exposure. Thus, the excess fraction (and its AR counterpart) would be  $6/24 = 0.25$  (25%), but the aetiologic fraction (and its AR counterpart) would be  $24/24 = 1.00$  (100%) because radiation exposure was a contributory cause in every one of the 24 cases. Under an alternative biologic assumption, the effect of radiation could be carcinogenic in 25% (i.e., six) of the exposed cases, but without any effect in the remaining 75% (i.e., 18) 'spontaneous' cases, which would have become cases irrespective of exposure, this differential effect being possibly explained by a difference in genetic make-ups of subjects (i.e., exposure-gene interaction). Thus, the excess fraction would be  $6/24 = 0.25$  (25%) as before, but the aetiologic fraction would be 25% as well in this case, because radiation exposure did not play a contributory causal role in any one of the remaining 18 cases. This considerable change in aetiologic fraction estimates results from the consideration of different biologic mechanisms, although the same epidemiologic data are used. This example helps to grasp the distinction between the excess and aetiologic fractions and their AR counterparts and stresses that, depending on the underlying biologic assumptions, the aetiologic fraction estimates can vary enormously. Other similar examples were given by Beyea and Green-

land [121]. Thus, the use of the aetiologic fraction and of its AR counterpart is problematic, because it requires usually unverifiable biologic assumptions to be estimated. Moreover, AR counterparts to excess fractions and rates seem better suited to quantify the impact of exposure on disease occurrence at the population level in view of their definitions. However, the aetiologic fraction may be useful from a biological standpoint, and it has also become the key measure used in legal thinking for awarding compensation for harmful exposure, at least in the United States. In this context, the term 'probability of causation' has been proposed to denote the aetiologic fraction and has been extended to account for several exposed levels or multiple exposures [107–113]. Whereas aetiologic and excess fractions will usually differ numerically, sometimes very substantially, it has been shown that they are equal in certain biologic models [86,120].

## 7. Final remarks

Disease frequency is measured through the computation of incidence rates or estimation of disease risk. Both measures are directly accessible from cohort data. They can be obtained from case-control data only if they are complemented by follow-up or population data. Using regression techniques, methods are available to derive incidence rates or risk estimates specific to a given exposure profile. Exposure-specific risk estimates are useful in individual prediction. A wide variety of options and techniques are available for measuring association. Adjustment for confounding is a key point in all analyses of observational studies, and can be pursued by standardization, stratification, and by regression techniques. The flexibility of the latter, especially in the generalized linear model framework, and the availability of computer software, has made it widely applied in the last several years.

Several measures are available to assess the impact of an exposure in terms of the occurrence of new disease cases at the population level, among which AR is the most commonly used. Several approaches have been developed to derive adjusted AR estimates from case-control as well as cohort data, either based on stratification or on more flexible regression techniques. Sequential and partial ARs have been proposed to handle the situation of multiple exposures and circumvent the associated non-additivity problem. Although there remain issues in properly interpreting the concept of AR, AR remains a useful measure to assess the potential impact of exposure at the population level and can serve as

a suitable guide in practice to assess and compare various prevention strategies.

General problems of AR definition, interpretation and usefulness as well as properties have been reviewed in detail [6,30,51,62,122,123]. Special issues were reviewed by Benichou [62,63]. They include estimation of AR for risk factors with multiple levels of exposure or with a continuous form, multiple risk factors, recurrent disease events, and disease classification with more than two categories. They also include assessing the consequences of exposure misclassification on AR estimates. Specific software for attributable risk estimation [100,124,125] as well as a simplified approach to confidence interval estimation [126] has been developed to facilitate implementation of methods for attributable risk estimation. Finally, much remains to be done to promote proper use and interpretation of AR, as illustrated in a recent literature review [127].

## References

- [1] O.S. Miettinen, Estimability and estimation in case-referent studies, *Am. J. Epidemiol.* 103 (1976) 226–235.
- [2] H. Morgenstern, D. Kleinbaum, L.L. Kupper, Measures of disease incidence used in epidemiologic research, *Int. J. Epidemiol.* 9 (1980) 97–104.
- [3] J. Cornfield, A method for estimating comparative rates from clinical data: Applications to cancer of the lung, breast and cervix, *J. Natl Cancer Inst.* 11 (1951) 1269–1275.
- [4] J. Cornfield, A statistical problem arising from retrospective studies, in: J. Neyman (Ed.), *Proc. Third Berkeley Symposium*, vol. IV, University of California Press, Monterey, CA, USA, 1956, pp. 133–148.
- [5] B. MacMahon, Prenatal X-ray exposure and childhood cancer, *J. Natl Cancer Inst.* 28 (1962) 1173–1191.
- [6] O.S. Miettinen, Proportion of disease caused or prevented by a given exposure, trait or intervention, *Am. J. Epidemiol.* 99 (1974) 325–332.
- [7] R.R. Neutra, M.E. Drolette, Estimating exposure-specific disease rates from case-control studies using Bayes' theorem, *Am. J. Epidemiol.* 108 (1978) 214–222.
- [8] M.H. Gail, L.A. Brinton, D.P. Byar, D.K. Corle, S.B. Green, C. Schairer, J.J. Mulvihill, Projecting individualized probabilities of developing breast cancer for white females who are being examined annually, *J. Natl Cancer Inst.* 81 (1989) 1879–1886.
- [9] R.L. Prentice, J.D. Kalbfleisch, A.V. Peterson, N. Flourmoy, V.T. Farewell, N.E. Breslow, The analysis of failure times in the presence of competing risks, *Biometrics* 34 (1978) 541–554.
- [10] J. Benichou, Absolute risk, in: M.H. Gail, J. Benichou (Eds.), *Encyclopedia of Epidemiologic Methods*, Wiley, Chichester, UK, 2000, pp. 1–17.
- [11] J. Benichou, M.H. Gail, Estimates of absolute cause-specific risk in cohort studies, *Biometrics* 46 (1990) 813–826.
- [12] J. Benichou, M.H. Gail, Methods of inference for estimates of absolute risk derived from population-based case-control studies, *Biometrics* 51 (1995) 182–194.
- [13] D.W. Dupont, Converting relative risks to absolute risks: A graphical approach, *Stat. Med.* 8 (1989) 641–651.
- [14] B. Langholz, O. Borgan, Estimation of absolute risk from nested case-control data, *Biometrics* 53 (1997) 767–774.
- [15] D. Spiegelman, G.A. Colditz, D. Hunter, E. Hertzmark, Validation of the Gail et al. model for predicting individual breast cancer risk, *J. Natl Cancer Inst.* 86 (1994) 600–607.
- [16] C.L. Chiang, *Introduction to Stochastic Processes in Biostatistics*, Wiley, New York, 1968.
- [17] E.L. Korn, F.J. Dorey, Applications of crude incidence curves, *Stat. Med.* 11 (1992) 813–829.
- [18] R.J. Gray, A class of  $k$ -sample tests for comparing the cumulative incidence of a competing risk, *Ann. Stat.* 16 (1988) 1141–1151.
- [19] B. MacMahon, T.F. Pugh, *Epidemiology: Principles and Methods*, Little, Brown and Co, Boston, MA, 1970.
- [20] K.J. Rothman, S. Greenland, *Modern Epidemiology*, Lippincott-Raven, Philadelphia, PA, USA, 1998.
- [21] J. Benichou, M.H. Gail, J.J. Mulvihill, Graphs to estimate an individualized risk of breast cancer, *J. Clin. Oncol.* 14 (1996) 103–110.
- [22] M.H. Gail, J. Benichou, Validation studies on a model for breast cancer risk (editorial), *J. Natl Cancer Inst.* 86 (1994) 573–575.
- [23] K.F. Hoskins, J.E. Stopfer, K. Calzone, S.D. Merajver, T.R. Rebbeck, J.E. Garber, B.L. Weber, Assessment and counseling for women with a family history of breast cancer. A guide for clinicians, *J. Am. Med. Assoc.* 273 (1995) 577–585.
- [24] M.H. Gail, J.P. Costantino, J. Bruant, R. Croyle, L. Freedman, K. Helzlsouer, V. Vogel, Weighing the risks and benefits of Tamoxifen treatment for preventing breast cancer, *J. Natl Cancer Inst.* 91 (1999) 1829–1846.
- [25] N.E. Breslow, N.E. Day, *Statistical Methods in Cancer Research*, vol. 1: The Analysis of Case-Control Studies, International Agency for Research on Cancer Scientific Publications, No. 32, Lyons, France, 1980.
- [26] D.G. Kleinbaum, L.L. Kupper, H. Morgenstern, *Epidemiologic Research: Principles and Quantitative Methods*, Lifetime Learning Publications, Belmont, CA, USA, 1982.
- [27] J. Berkson, Smoking and lung cancer. Some observations on two recent reports, *J. Am. Stat. Assoc.* 53 (1958) 28–38.
- [28] J.S. Mausner, A.K. Bahn, *Epidemiology: An Introductory Text*, W.B. Saunders, Philadelphia, PA, USA, 1974.
- [29] J.J. Schlesselman, *Case-control Studies. Design, Conduct and Analysis*, Oxford University Press, New York, 1982.
- [30] S.D. Walter, The estimation and interpretation of attributable risk in health research, *Biometrics* 32 (1976) 829–849.
- [31] R.E. Markush, Levin's attributable risk statistic for analytic studies and vital statistics, *Am. J. Epidemiol.* 105 (1977) 401–406.
- [32] R.L. Prentice, R. Pyke, Logistic disease incidence models and case-control studies, *Biometrika* 66 (1979) 403–411.
- [33] S. Greenland, D.C. Thomas, On the need for the rare disease assumption, *Am. J. Epidemiol.* 116 (1982) 547–553.
- [34] R.L. Prentice, N.E. Breslow, Retrospective studies and failure time models, *Biometrika* 65 (1978) 153–158.
- [35] H.A. Kahn, C.T. Sempos, *Statistical Methods in Epidemiology*, Monographs in Epidemiology and Biostatistics, vol. 12, Oxford University Press, Oxford, New York, 1989.
- [36] P. McCullagh, J.A. Nelder, *Generalized Linear Models*, second ed., CRC Press, Boca Raton, FL, USA, 1989.
- [37] M. Palta, *Quantitative Methods in Population Health: Extensions of Ordinary Regression*, John Wiley & Sons, Hoboken, NJ, 2003.
- [38] N.E. Breslow, N.E. Day, *Statistical Methods in Cancer Research*, vol. II: The Design and Analysis of Cohort Studies,



- International Agency for Research on Cancer Scientific Publications, No. 82, Lyons, France, 1987.
- [39] O.S. Miettinen, *Theoretical Epidemiology. Principles of Occurrence Research in Medicine*, Delmar Publishers Inc., Albany, NY, USA, 1985.
- [40] D.E. Lilienfeld, P.D. Stolley, *Foundations of Epidemiology*, third ed., Oxford University Press, New York, 1994.
- [41] B. MacMahon, D. Trichopoulos, *Epidemiology: Principles and Methods*, second ed., Little, Brown and Co, Boston, 1996.
- [42] W. Ahrens, I. Pigeot (Eds.), *Handbook of Epidemiology*, Springer-Verlag, Berlin, 2005.
- [43] M.H. Gail, J. Benichou (Eds.), *Encyclopedia of Epidemiologic Methods*, John Wiley and Sons, Chichester, UK, 2000.
- [44] M.L. Levin, The occurrence of lung cancer in man, *Acta Unio Internationalis contra Cancrum* 9 (1953) 531–541.
- [45] P. Cole, B. MacMahon, Attributable risk percent in case-control studies, *Br. J. Prev. Soc. Med.* 25 (1971) 242–244.
- [46] J. Benichou, Preventable fraction, in: M.H. Gail, J. Benichou (Eds.), *Encyclopedia of Epidemiologic Methods*, Wiley, Chichester, 2000, pp. 736–737.
- [47] S. Greenland, Variance estimators for attributable fraction estimates, consistent in both large strata and sparse data, *Stat. Med.* 6 (1987) 701–708.
- [48] J.M. Last, *A dictionary of Epidemiology*, Oxford University Press, New York, 1983.
- [49] P.M. Gargiullo, R. Rothenberg, H.G. Wilson, Confidence intervals, hypothesis tests, and sample sizes for the prevented fraction in cross-sectional studies, *Stat. Med.* 14 (1995) 51–72.
- [50] O. Gefeller, Theory and application of attributable risk estimation in cross-sectional studies, *Stat. Appl.* 2 (1990) 323–331.
- [51] O. Gefeller, Definitions of attributable risk-revisited, *Public Health Rev.* 23 (1995) 343–355.
- [52] W. Uter, A. Pfahlberg, The concept of attributable risk in epidemiological practice, *Biom. J.* 41 (1999) 985–999.
- [53] S. Greenland, J.M. Robins, Conceptual problems in the definition and interpretation of attributable fractions, *Am. J. Epidemiol.* 128 (1988) 1185–1197.
- [54] B.L. Ouellet, J.M. Romeder, J.M. Lance, Premature mortality attributable to smoking and hazardous drinking in Canada, *Am. J. Epidemiol.* 109 (1979) 451–463.
- [55] M.P. Madigan, R.G. Ziegler, J. Benichou, C. Byrne, R.N. Hoover, Proportion of breast cancer cases in the United States explained by well-established risk factors, *J. Natl Cancer Inst.* 87 (1995) 1681–1685.
- [56] B.N. Ames, L.S. Gold, W.C. Willett, The causes and prevention of cancer, *Proc. Natl Acad. Sci. USA* 254 (1995) 1131–1138.
- [57] G. Colditz, W. DeJong, D. Hunter, D. Trichopoulos, W. Willett (Eds.), *Harvard Report on Cancer Prevention*, vol. 1, *Cancer Causes Control* 7 (suppl.) (1996) S3–S59.
- [58] G. Colditz, W. DeJong, D. Hunter, D. Trichopoulos, W. Willett (Eds.), *Harvard Report on Cancer Prevention*, vol. 2, *Cancer Causes Control* 8 (suppl.) (1997) S1–S50.
- [59] R. Doll, R. Peto, *The Causes of Cancer*, Oxford University Press, New York, 1981.
- [60] B.E. Henderson, R.K. Ross, M.C. Pike, Toward the primary prevention of cancer, *Science* 254 (1991) 1131–1138.
- [61] C.B. Begg, The search for cancer risk factors: When can we stop looking? *Am. J. Public Health* 91 (2001) 360–364.
- [62] J. Benichou, Attributable risk, in: M.H. Gail, J. Benichou (Eds.), *Encyclopedia of Epidemiologic Methods*, Wiley, Chichester, UK, 2000, pp. 50–63.
- [63] J. Benichou, A review of adjusted estimators of the attributable risk, *Stat. Methods Med. Res.* 10 (2001) 195–216.
- [64] K. Drescher, H. Becher, Estimating the generalized attributable fraction from case-control data, *Biometrics* 53 (1997) 1170–1176.
- [65] H. Morgenstern, E.S. Bursic, A method for using epidemiologic data to estimate the potential impact of an intervention on the health status of a target population, *J. Community Health* 7 (1982) 292–309.
- [66] S.D. Walter, Prevention for multifactorial diseases, *Am. J. Epidemiol.* 112 (1980) 409–416.
- [67] P. Morfeld, Years of life lost due to exposure: causal concepts and empirical shortcomings, *Epidemiol. Perspect. Innovations* 1 (2004) 5.
- [68] L. Smith, Person-years of life lost, in: P. Armitage, T. Colton (Eds.), *Encyclopedia of Biostatistics*, Wiley, Chichester, UK, 1998, pp. 3324–3325.
- [69] J.M. Robins, S. Greenland, Estimability and estimation of expected years of life lost due to a hazardous exposure, *Stat. Med.* 10 (1991) 79–93.
- [70] S.D. Walter, Effects of interaction, confounding and observational error on attributable risk estimation, *Am. J. Epidemiol.* 117 (1983) 598–604.
- [71] O.S. Miettinen, Components of the crude risk ratio, *Am. J. Epidemiol.* 96 (1972) 168–172.
- [72] H. Morgenstern, Uses of ecologic analysis in epidemiological research, *Am. J. Public Health* 72 (1982) 1336–1344.
- [73] A. Ejigou, Estimation of attributable risk in the presence of confounding, *Biom. J.* 21 (1979) 155–165.
- [74] S. Greenland, H. Morgenstern, Morgenstern corrects a conceptual error [letter], *Am. J. Public Health* 73 (1983) 703–704.
- [75] S. Greenland, Bias in methods for deriving standardized mortality ratio and attributable fraction estimates, *Stat. Med.* 3 (1984) 131–141.
- [76] S.J. Kuritz, J.R. Landis, Attributable risk estimation from matched-pairs case-control data, *Am. J. Epidemiol.* 125 (1987) 324–328.
- [77] S.J. Kuritz, J.R. Landis, Summary attributable risk estimation from unmatched case-control data, *Stat. Med.* 7 (1988) 507–517.
- [78] S.J. Kuritz, J.R. Landis, Attributable risk estimation from matched case-control data, *Biometrics* 44 (1988) 355–367.
- [79] J.R. Landis, T.J. Sharp, S.J. Kuritz, G. Koch, Mantel–Haenszel methods, in: M.H. Gail, J. Benichou (Eds.), *Encyclopedia of Epidemiologic Methods*, Wiley, Chichester, UK, 2000, pp. 499–512.
- [80] R.E. Tarone, On summary estimators of relative risk, *J. Chron. Dis.* 34 (1981) 463–468.
- [81] N. Mantel, W. Haenszel, Statistical aspects of the analysis of data from retrospective studies of disease, *J. Natl Cancer Inst.* 22 (1959) 719–748.
- [82] J. Benichou, Methods of adjustment for estimating the attributable risk in case-control studies: A review, *Stat. Med.* 10 (1991) 1753–1773.
- [83] M.W. Birch, The detection of partial associations, I: The  $2 \times 2$  case, *J. R. Stat. Soc., Ser. B* 27 (1964) 313–324.
- [84] N.E. Breslow, Odds ratio estimators when the data are sparse, *Biometrika* 68 (1981) 73–84.
- [85] J.R. Landis, E.R. Heyman, G.G. Koch, Average partial association in three-way contingency tables: A review and discussion of alternative tests, *Int. Stat. Rev.* 46 (1978) 237–254.
- [86] J.M. Robins, S. Greenland, Estimability and estimation of excess and etiologic fractions, *Stat. Med.* 8 (1989) 845–859.

- [87] O. Gefeller, The bootstrap method for standard errors and confidence intervals of the adjusted attributable risk [letter], *Epidemiology* 3 (1992) 271–272.
- [88] A.S. Whittemore, Statistical methods for estimating attributable risk from retrospective data, *Stat. Med.* 1 (1982) 229–243.
- [89] A.S. Whittemore, Estimating attributable risk from case-control studies, *Am. J. Epidemiol.* 117 (1983) 76–85.
- [90] F. Sturmans, P.G.H. Mulder, H.A. Walkenburg, Estimation of the possible effect of interventive measures in the area of ischemic heart diseases by the attributable risk percentage, *Am. J. Epidemiol.* 105 (1977) 281–289.
- [91] J.L. Fleiss, Inference about population attributable risk from cross-sectional studies, *Am. J. Epidemiol.* 110 (1979) 103–104.
- [92] D.C. Deubner, W.E. Wilkinson, M.J. Helms, H.A. Tyroler, C.G. Hames, Logistic model estimation of death attributable to risk factors for cardiovascular disease in Evans County, Georgia, *Am. J. Epidemiol.* 112 (1980) 135–143.
- [93] P. Bruzzi, S.B. Green, D.P. Byar, L.A. Brinton, C. Schairer, Estimating the population attributable risk for multiple risk factors using case-control data, *Am. J. Epidemiol.* 122 (1985) 904–914.
- [94] S. Basu, J.R. Landis, Model-based estimation of population attributable risk under cross-sectional sampling, *Am. J. Epidemiol.* 142 (1995) 1338–1343.
- [95] J. Benichou, M.H. Gail, A delta-method for implicitly defined random variables, *Am. Stat.* 43 (1989) 41–44.
- [96] J. Benichou, M.H. Gail, Variance calculations and confidence intervals for estimates of the attributable risk based on logistic models, *Biometrics* 46 (1990) 991–1003.
- [97] S. Greenland, K. Drescher, Maximum-likelihood estimation of the attributable fraction from logistic models, *Biometrics* 49 (1993) 865–872.
- [98] S.S. Coughlin, J. Benichou, D.L. Weed, Attributable risk estimation in case-control studies, *Epidemiol. Rev.* 16 (1994) 51–64.
- [99] S. Greenland, The bootstrap method for standard errors and confidence intervals of the adjusted attributable risk [letter], *Epidemiology* 3 (1992) 271.
- [100] M.J. Kahn, W.M. O’Fallon, J.D. Sicks, Generalized Population Attributable Risk Estimation, Technical Report #54, Mayo Foundation, Rochester, MN, USA, 1998.
- [101] C. Kooperberg, D.B. Petitti, Using logistic regression to estimate the adjusted attributable risk of low birthweight in an unmatched case-control study, *Epidemiology* 2 (1991) 363–366.
- [102] J. Llorca, M. Delgado-Rodríguez, A comparison of several procedures to estimate the confidence interval for attributable risk in case-control studies, *Stat. Med.* 19 (2000) 1089–1099.
- [103] K.J. Lui, Interval estimation of the attributable risk in case-control studies with matched pairs, *J. Epidemiol. Community Health* 55 (2001) 885–890.
- [104] K.J. Lui, Notes on interval estimation of the attributable risk in cross-sectional sampling, *Stat. Med.* 20 (2001) 1797–1809.
- [105] K.J. Lui, Interval estimation of the attributable risk for multiple exposure levels in case-control studies with confounders, *Stat. Med.* 22 (2003) 2443–2557.
- [106] C.B. Begg, J.M. Satagopan, M. Berwick, A new strategy for evaluating the impact of epidemiologic risk factors for cancer with applications to melanoma, *J. Am. Stat. Assoc.* 93 (1998) 415–426.
- [107] J. Benichou, Re: “Methods of adjustment for estimating the attributable risk in case-control studies: A review” (letter), *Stat. Med.* 12 (1993) 94–96.
- [108] L.A. Cox, Probability of causation and the attributable proportion of risk, *Risk Anal.* 4 (1984) 221–230.
- [109] L.A. Cox, A new measure of attributable risk for public health applications, *Manage. Sci.* 7 (1985) 800–813.
- [110] S.W. Lagakos, F. Mosteller, Assigned shares in compensation for radiation-related cancers (with discussion), *Risk Anal.* 6 (1986) 345–380.
- [111] P. McElduff, J. Attia, B. Ewald, J. Cockburn, R. Heller, Estimating the contribution of individual risk factors to disease in a person with more than one risk factor, *J. Clin. Epidemiol.* 55 (2002) 588–592.
- [112] F.A. Seiler, Attributable risk, probability of causation, assigned shares, and uncertainty, *Environ. Int.* 12 (1986) 635–641.
- [113] F.A. Seiler, B.R. Scott, Mixture of toxic agents and attributable risk calculations, *Risk Anal.* 7 (1987) 81–90.
- [114] G.E. Eide, O. Gefeller, Sequential and average attributable fractions as aids in the selection of preventive strategies, *J. Clin. Epidemiol.* 48 (1995) 645–655.
- [115] M. Land, C. Vogel, O. Gefeller, Partitioning methods for multifactorial risk attribution, *Stat. Methods Med. Res.* 10 (2001) 217–230.
- [116] O. Gefeller, M. Land, G.E. Eide, Averaging attributable fractions in the multifactorial situation: Assumptions and interpretation, *J. Clin. Epidemiol.* 51 (1998) 437–451.
- [117] M. Land, O. Gefeller, A game-theoretic approach to partitioning attributable risks in epidemiology, *Biom. J.* 39 (1997) 777–792.
- [118] G.E. Eide, I. Heuch, Attributable fractions: fundamental concepts and their visualization, *Stat. Methods Med. Res.* 10 (2001) 159–193.
- [119] A.J. Tuyns, G. Pequignot, O.M. Jensen, Le cancer de l’œsophage en Ille-et-Vilaine en fonction des niveaux de consommation d’alcool et de tabac, *Bull. Cancer* 64 (1977) 45–60.
- [120] L.A. Cox, Statistical issues in the estimation of assigned shares for carcinogenesis liability, *Risk Anal.* 7 (1987) 71–80.
- [121] J. Beyea, S. Greenland, The importance of specifying the underlying biologic model in estimating the probability of causation, *Health Phys.* 76 (1999) 269–274.
- [122] B. Rockhill, B. Newman, C. Weinberg, Use and misuse of population attributable fractions, *Am. J. Public Health* 88 (1998) 15–21.
- [123] B. Rockhill, C. Weinberg, B. Newman, Population attributable fraction estimation for established breast cancer risk factors: considering the issues of high prevalence and unmodifyability, *Am. J. Epidemiol.* 147 (1998) 826–833.
- [124] M. Mezzetti, M. Ferraroni, A. Decarli, C. La Vecchia, J. Benichou, Software for attributable risk and confidence interval estimation in case-control studies, *Comput. Biomed. Res.* 29 (1996) 63–75.
- [125] US National Cancer Institute’s Division of Epidemiology and Genetics Interactive Risk Attributable Program, <http://dceg.cancer.gov/IRAP/>, 2002, accessed on 15 November 2006.
- [126] L.E. Daly, Confidence limits made easy: interval estimation using a substitution method, *Am. J. Epidemiol.* 147 (1998) 783–790.
- [127] W. Uter, A. Pfahlberg, The application of methods to quantify attributable risk in medical practice, *Stat. Methods Med. Res.* 10 (2001) 231–237.