

# We are IntechOpen, the world's leading publisher of Open Access books Built by scientists, for scientists

6,900

Open access books available

185,000

International authors and editors

200M

Downloads

Our authors are among the

154

Countries delivered to

TOP 1%

most cited scientists

12.2%

Contributors from top 500 universities



WEB OF SCIENCE™

Selection of our books indexed in the Book Citation Index  
in Web of Science™ Core Collection (BKCI)

Interested in publishing with us?  
Contact [book.department@intechopen.com](mailto:book.department@intechopen.com)

Numbers displayed above are based on latest data collected.  
For more information visit [www.intechopen.com](http://www.intechopen.com)



---

# Limitations and Biases in Cohort Studies

---

Muriel Ramirez-Santana

Additional information is available at the end of the chapter

<http://dx.doi.org/10.5772/intechopen.74324>

---

## Abstract

Good practice in research involves considering diverse sources of biases when designing a study for later validation of results. If they are recognized beforehand, it is possible to minimize or avoid them. Selection biases may originate at the time of enrolling the subjects of study, making it necessary to clearly state the selection criteria of the exposed and nonexposed individuals. If people get lost from the original sample, bias may be introduced by the consequences of reducing the sample. Biases of information could originate in loss of evidence at the moment of recording the data. The definition of follow-up protocols may also help to keep registers of all variables, so information will not be missed from the individuals under study or from the observers who conduct the follow-up. It is necessary to apply the same protocols and instruments for measuring and evaluating the health outcomes in exposed and nonexposed individuals in order to avoid biases of misclassification. Confusion biases can be avoided at the time of designing the study, with the inclusion of confounding variables from the onset. Matching by age and gender is strongly recommended, and finally, adjustment techniques are used at the time of the data analysis.

**Keywords:** systematic error, selection bias, information bias, confusion, interaction, cohort studies

---

## 1. Introduction

The external validity of the results of an analytical study (including cohort studies) is determined by the possibility that the results can be extrapolated to larger populations, making the representativeness and randomness of the sample(s) important. However, there is controversy about the real need of representativeness when other situations are more relevant in the study, for example, some practical reasons, restrictions in the selection criteria or focus in certain population groups [1].

Internal validity, however, is determined by a series of factors that can lead to systematic errors or biases [2]. Bias can originate both in the design stage of the study, such as sample selection, data collection or analysis, but can be minimized with good planning of the study protocol or using statistical analysis techniques in this phase of the study [3]. The sample size will determine the validity in terms of the statistical power necessary to reject or approve the working hypothesis. An adequate sample size will make it easier to avoid random errors in the results of the study.

Although cohort studies have a lower risk of presenting biases than other types of epidemiological studies (ecological, cross-sectional or prevalence studies, cases and controls), they are not free of them. This chapter highlights the types of biases, their origin, their effects on the validity of the study and ways to avoid or minimize them. The chapter also gives examples that allow better understanding of the concepts as well as practical advice when carrying out a cohort study.

### 1.1. Feasibility considerations

Study protocols should always adhere to the evaluation of duly accredited Scientific Ethics Committees. Ethical principles indicate that all participants must adhere to informed consent before beginning to participate in the study, being able to understand all the implications of participating and to decline his/her participation at any moment. Authorizations of the managers in charge of the administration of any institution (healthcare centers, schools, municipalities, hospitals or others) are usually required to access the registered data or to collect the health information of the users. In studies of occupational health, authorization of the workplaces is required to perform the evaluations of jobs and workers exposed to occupational hazards. In studies about infants or children, framed in the educational sector, the assent of the minors is required, in addition to informed consent of the parents/guardians/proxies, and authorization of the executives of participating educational facilities. Collaboration agreements, purchase of laboratory services, transport, locations, surveyors, data analysis, computer support and other technical and logistical requirements that involve carrying out a follow-up study of people, usually for several years, must also be managed. When a large research team is involved, protocols must be in place for recruitment, evaluations, transporting and storage of samples and materials, laboratory procedures, recording data, backing up information and so on.

## 2. Bias in cohort studies

Certainly, among analytical epidemiological research, cohort studies are less prone to have bias than the case-control ones, specifically regarding memory bias. But as any other epidemiological study, several biases could be present in cohort studies. In this sense, researchers must be aware of those biases in advance and take them into account at the moment of selecting participants, designing the study (collection tools/instruments), when registering the data during field work (data base design) and, later on, at the moment of analyzing and interpreting the data (statistical analysis).

We understand bias as systematic errors that can lead to mistaken results or interpretation regarding the association under study, when the purpose of a study is assessing the association of certain factors toward supporting the causality of a health event or outcome [4].

There are several ways of classification of biases. For academic proposes, we will use the following classification [2, 4]:

Selection bias: originated from the way the participants of the study are selected or followed and can affect the apparent association between the exposure and outcome.

Information biases: could originate in the observed individuals, in the observers or in the instruments used to assess the outcomes.

Confusion bias: their origin is in the relationship that other variables that are not the exposition are related to the outcome, and can modulate the effect(s) of the exposition, contributing to a spurious association.

We will now review each kind of bias in detail and with some examples.

## 2.1. Selection bias

In cohort studies, the researcher must select exposed and nonexposed individuals. In the first place, it should be understood that both groups are representative of the general population from where they are taken, in order to facilitate the external validity of the study (basic condition to generalize the results in order to support causality). This condition, however, would not necessarily affect the internal validity. In other words, the internal validity is due to systematic errors sourced in stubborn participation of individuals.

The appropriated assessment of the exposure is the first crucial step. Auto-selection is one of the circumstances that could lead to inaccurate selection. As an example, a study conducted among pregnant women in Norway intended to evaluate auto-selection bias by comparing two cohorts; one group was taken from the Medical Birth Registry (2000–2006) as a population-based cohort, and the second group was from women who agreed to participate in the Norwegian Mother and Child Cohort Study. The results suggested that the prevalence estimates of exposures and the outcomes were biased due to self-selection in the Norwegian Mother and Child Cohort Study. Nevertheless, the estimates of exposure-outcome associations were not biased [5]. But in other cases, the associations could also be flawed.

Another example for selection bias might occur when the compared cohorts are part of a population who receive public health interventions, so the exposure can be misled by this influence. That is the example given by researchers who studied the association of bad water quality (measured by *E. coli* burden) and development of diarrhea in Bangladesh. Interventions to purify water (use of chlorine) may interfere by reducing the pathogens and misclassify the exposure [6].

Selection criteria must be clearly defined from the beginning of the study in a way that ensures that biases are avoided. For example, a research was conducted with the objective of assessing the association between exposure to pesticides and neurocognitive impairment, including fine

motor coordination [7]. The researchers used Purdue Pegboard and MOART reaction time tests to measure the outcome. If any right- or left-handed people were selected, bias may be introduced when evaluating the outcomes due to the way that the tests are performed. Both tests have separated evaluation of left and right hands, giving certain scoring to the performance. Then, one important inclusion criterion to consider was right-handed people only; so the responses were standardized under the same criteria, and bias was eluded.

Another example of selection bias can happen in large multicenter cohort studies evaluating the association between diet and cancer. In this case, systematic errors may originate in the measurement of the exposure by dietary questionnaires that are not easy to standardize for all locations. Researchers suggest to use the calibration approach for such cases [8].

Another type of selection bias is known as the nonresponders or no-participation bias, which are less frequent in prospective cohort studies due to the need for strict following-up of the participants, strengthening the evaluations during the follow-up visits, encouraging participants and evaluators (observers) to always respond and/or register the records properly. Nevertheless, missing data could be present in retrospective cohort studies, where previously registered data are used. This will be explained in detail later, related to the information biases (Section 2.2).

In prospective cohort studies, loss of follow-up may occur, giving rise to selection bias. Loss of follow-up bias is caused by the loss of individuals from one or more exposure groups. Because cohort studies take normally several months or years of following the participants, it is expected that life situations will vary from time to time, causing some of the participants to get lost during the development of the study. Individuals can be lost homogeneously in the groups to be compared, causing bias of poor global miss-classification, which generally leads the estimate toward the null value [9]. Or individuals from a single group can be lost, causing bias of poor differential miss-classification. In the first case, the estimated risk would not be severely affected, because the incidence rates would keep similar in both groups, but the power of the results may be lost. In the latter case, the results to be obtained may be underestimating or overestimating the association. For example, if the people who are exposed and develop the outcome (disease) are lost, the incidence rate may be lower among the exposed individuals and the relative risk (RR) would be underestimated. On the other hand, if people who are not exposed and do not get the disease after time of follow-up get lost, then the incidence rate among the nonexposed will be higher and the RR would be overestimated.

Here is a hypothetical example showing the four possibilities of losing individuals:

**Original data:**

Size of the exposed cohort = 1000.

Size of the nonexposed cohort = 1000.

Number of individuals with the outcome among the exposed = 100.

Number of individuals with the outcome among the nonexposed = 10.

**Correct results:**

Incidence rate in exposed =  $100/1000 = 0.1$ .

Incidence rate in nonexposed =  $10/1000 = 0.01$ .

Relative risk = 10.

**Loss of 50 individuals during follow-up with the disease among the exposed:**

Incidence rate in exposed =  $50/1000 = 0.05$ .

Incidence rate in nonexposed =  $10/1000 = 0.01$ .

Relative risk = 5.

**Loss of five individuals during follow-up with the disease among the nonexposed:**

Incidence rate in exposed =  $100/1000 = 0.1$ .

Incidence rate in nonexposed =  $5/1000 = 0.005$ .

Relative risk = 20.

**Loss of 100 individuals during follow-up without the disease among the exposed:**

Incidence rate in exposed =  $100/800 = 0.125$ .

Incidence rate in nonexposed =  $10/1000 = 0.01$ .

Relative risk = 12.5.

**Loss of 200 individuals during follow-up without the disease among the nonexposed:**

Incidence rate in exposed =  $100/1000 = 0.1$ .

Incidence rate in nonexposed =  $10/800 = 0.0125$ .

Relative risk = 8.

As you can see, the estimated association variation is given by the number of people who completed the follow-up schedule. The general recommendation is that 60–80% of the individuals complete the timeframe defined originally, but a study that simulated a cohort of 500 observations with 1000 replications in computer found utterly biased estimates of the risks with low ranks of loss to follow-up [10]. On the other hand, as was already said, the results of the diminution in the number of subjects can also affect the statistical power of the results. Then, in the design of the study, at least 10% sample loss must be considered, so this proportion must be added to the minimum calculated sample size for the study. During the field work phase, measures need to be taken in advance in order to avoid losing individuals. To ensure the permanence of the individuals during the follow-up time, it is suggested to include incentives for the participants. These incentives do not necessarily have to be monetary, and a food and transportation voucher can be offered for those who attend scheduled evaluations.



In addition, it may happen that nonexposed individuals enter into the exposed group or vice versa. An example of this could occur when studying the association of tobacco consumption and a certain outcome. Then, during the study, people who smoke can leave the consumption and/or people who do not smoke can start smoking. In those cases, it is suggested to use the *incidence density* indicator instead of the cumulative incidence. The incidence density is interpreted as exposure measured in units of *person-time*, for example, person-weeks or person-days. Person-time is the sum of the time periods of observation of each person who has been observed for all or part of the entire time period [4].

The incidence density is calculated as follows:

$$\frac{\text{Number of new cases of a disease occurring in a population during a specified period of time}}{\text{Total person – time (the sum of the time periods of observation of each person who was observed for all or part of the entire time period)}} \times 1.000. \quad (1)$$

It is important to mention that the person-time unit is not in all occasions equivalent to the person-time of all individuals. For example, one person-year could represent one person being followed for 1 year or two people being followed for 6 months. But in any case, this is a way of measuring incidence that is very useful in cohort studies because it avoids the issue of subjects shifting from one exposure group to the other.

Finally, we have the selective survival bias. This bias is known in occupational health as the *healthy worker effect* and occurs when workers who have the health effect (disease or outcome) abandon the work, so a greater proportion of healthy exposed workers finally lead to underestimation of the health effect or outcome. This situation may happen when the exposed individuals have the condition already for certain time (prevalent cohort), so the probability to express the outcome is greater than individuals who were recently exposed (incident cohort). That effect is known as *left truncation bias* or *time related (immortal time bias, time lag bias)* [11, 12]. This influence has been described in several studies: occupational settings, development of AIDS among HIV patients, cancer survival, obstetric research, use of acetylsalicylic acid and myocardial infarction [13–15]. The last is a good example, showing how the use of a cohort recently diagnosed with myocardial infarction has differences in baseline characteristics and prognosis compared to the group that has had the disease for some time (prevalent cohort), even though they were taken from the same population. Then, the researcher suggested studying incident cohorts when estimating survival of a defined outcome [14]. Another example of occupational health has been published utilizing simulation with the Monte Carlo technique. Results showed that prevalent jobs contribute to descendant bias in an occupational cohort. This arises because individuals who are less susceptible to the exposure's effect continue to be exposed, thus undervaluing the association [13].

## 2.2. Information biases

Loss during follow-up may cause information bias that was already explained in detail in Section 2.1 [10].

Usually in prospective cohorts, information bias is easy to elude, because measures may be taken during the design by including all variables in the registration forms (instruments), in order to not miss variables of interest. On the other hand, in retrospective cohorts, already existing records may be used. In that case, there could be missing data due to poor registration quality or due to variables that were not considered to be registered in advance. In both cases, the origin of missing information can lead to information bias. To minimize this effect on large population-based cohorts, it is possible to exclude individuals who have missing data from the analyses. But, this is a decision that researchers can take when the size of the remaining cohort still allows for sufficient statistical power to validate the results. That was the case presented in a large study conducted among the Danish population assessing the association between lifestyle and colorectal cancer [16]. From a total of 160,725 potential participants, several hundreds were not included due to nonresponse, cancer diagnosis and missing data (N = 997). Finally, a cohort of over 55,000 people was included in the investigation.

One important source of bias in cohort studies can occur when diagnosing the health event or outcome. It is necessary to apply the same protocol for measuring or evaluating the health outcomes in exposed and nonexposed individuals in order to avoid the biases of misclassification [9]. Similarly to what happen in the previously explained bias caused by loss of follow-up, the final effect of misclassification will depend on whether the inaccuracy in the evaluated outcome influences both exposure groups (global misclassification bias) or only affects one of them (differential misclassification bias).

Let us have a look at a hypothetical example in a study that evaluates the risk of having myocardial infarction due to exposure to a high-fat diet. **Table 1** shows the correct classification.

If the evaluation of the exposure is misled in both groups due to the mistakes in the daily food register, this results in a non-differential misclassification. Imagining that 20% of the exposed people go to the nonexposed group and 20% of nonexposed goes to the exposed group, we could have the following situation (**Table 2**).

In that case, the relative risk is diminished due to a higher incidence among the nonexposed group.

Now, suppose that the evaluators applied two diagnostic tests to the exposed that resulted in an increased diagnosis of myocardial infarction among the exposed group. This will result in a differential misclassification due to the mistaken diagnosis in the outcomes (**Table 3**).

High-fat diet	Myocardial infarction	
	Disease	No disease
Exposed	250	450
Nonexposed	100	900

$RR = (250/700)/(100/1000) = 0.357/0.1 = 3.57$ .

**Table 1.** High-fat diet and acute myocardial infarction, *correct classification*.



High-fat diet	Myocardial infarction	
	Disease	No disease
Exposed	290	410
Nonexposed	260	740

$RR = (290/700)/(260/1000) = 0.414/0.26 = 1.59.$

**Table 2.** High-fat diet and acute myocardial infarction, non-differential misclassification.

High-fat diet	Myocardial infarction	
	Disease	No disease
Exposed	295	405
Nonexposed	100	900

$RR = (295/700)/(100/1000) = 0.421/0.1 = 4.21.$

**Table 3.** High-fat diet and acute myocardial infarction, differential misclassification.

In this last case, when 10% of the exposed people without myocardial infarction moved to the disease group, the result is a higher relative risk due to a higher incidence among the exposed group.

A good example of this kind of misclassification bias could be given regarding the use of mortality records, which are frequently used in epidemiological studies. The registered codes of the diagnoses may be mistaken and lead to misclassification of the outcomes. That was studied recently by Deckert, who reported the results of a simulation study based on real data of cardiovascular disease mortality [17]. He reported that non-differential bias can lead to a null hypothesis, whereas differential misclassification leads the observed Standardized Mortality Ratios to be incorrect, in either direction or magnitude. Differences were from 10 to 30%, depending on the sensitivity and specificity characteristics of the diagnosis of cardiovascular disease [17]. Statistic techniques like quantitative bias analysis (QBA) or bootstrapping disease status imputation could be used to correct misclassification bias due to correct diagnostic codes [18]. Although statistical adjustments are possible to do in cases where standard information is available, these techniques are not always enough to overcome the bias. An example is reported by Candice Johnson et al., related to the misclassification of self-reported obesity and diabetes, adjusted by the National Health and Nutrition Examination Survey [19].

Regarding the accuracy in gathering information during the follow-up visits, one could have the temptation to assess more strictly the exposed individuals than the nonexposed or to evaluate the exposed persons more frequently than the nonexposed. The advice is to apply the same protocol and instruments to both groups of people, and in that way, bias introduced by the observer or the instruments is avoided. We understand as instruments the questionnaires, weighting scale, sphygmomanometer, altimeter, laboratory tests/techniques and others. Additionally, if the person(s) who observe and diagnose the outcome are aware of the exposure status,

a preconception may lead to overdiagnose the exposed people and/or underdiagnose the nonexposed people. Alike in randomized trials, the best way to avoid this bias is blinding the observers.

There is also a possibility that bias may originate from the observed individuals. It can happen that people, who know they are under observation, change their behavior. This has been called the Hawthorne effect and is due to the effects that the research can produce in the participants (observers and/or studied individuals). This was first described in a factory near Chicago between 1924 and 1936, in which a group of workers who knew they were under strict supervision significantly improved their productivity, compared to workers who were not aware of being observed [20]. There is still some controversy about the real predisposition effect of the participant's observation and the amount of bias that could cause. Some studies have found such an effect, but others have not [20, 21]. For example, a study conducted in Tanzania regarding malaria treatment did find a modest suggestion that the health professionals maintained better practice during the study [22].

Finally, all types of epidemiological studies may be affected by partiality in the phase of the analysis. The way to avoid this analytical bias is by *masking or blinding* the statistician. That means the statistician, or the person performing the analysis, does not know the exposure condition at the time of the analysis.

### 2.3. Confusion bias and interaction

We understand a confounder as a variable that is associated with the exposure as well as to the health event or outcome, but not being necessarily a cause of the event. For example, an inaccurate causal inference can be made between drinking coffee and pancreatic cancer, when drinking coffee has been associated with a smoking habit [4]. This is known as *spurious association*.

The most common confusion variables to be considered during the design of any epidemiological study are gender and age. As cohort studies are observational, people are not randomly assigned to the exposure and nonexposure group; it is not always possible to match both groups by certain variables such as sex, age, or other confounders. Depending on the exposure or events being studied, other variables could work as confounders; therefore, before designing any study, it is important for researchers to read previous studies and develop the design with all evidence that highlight confounders.

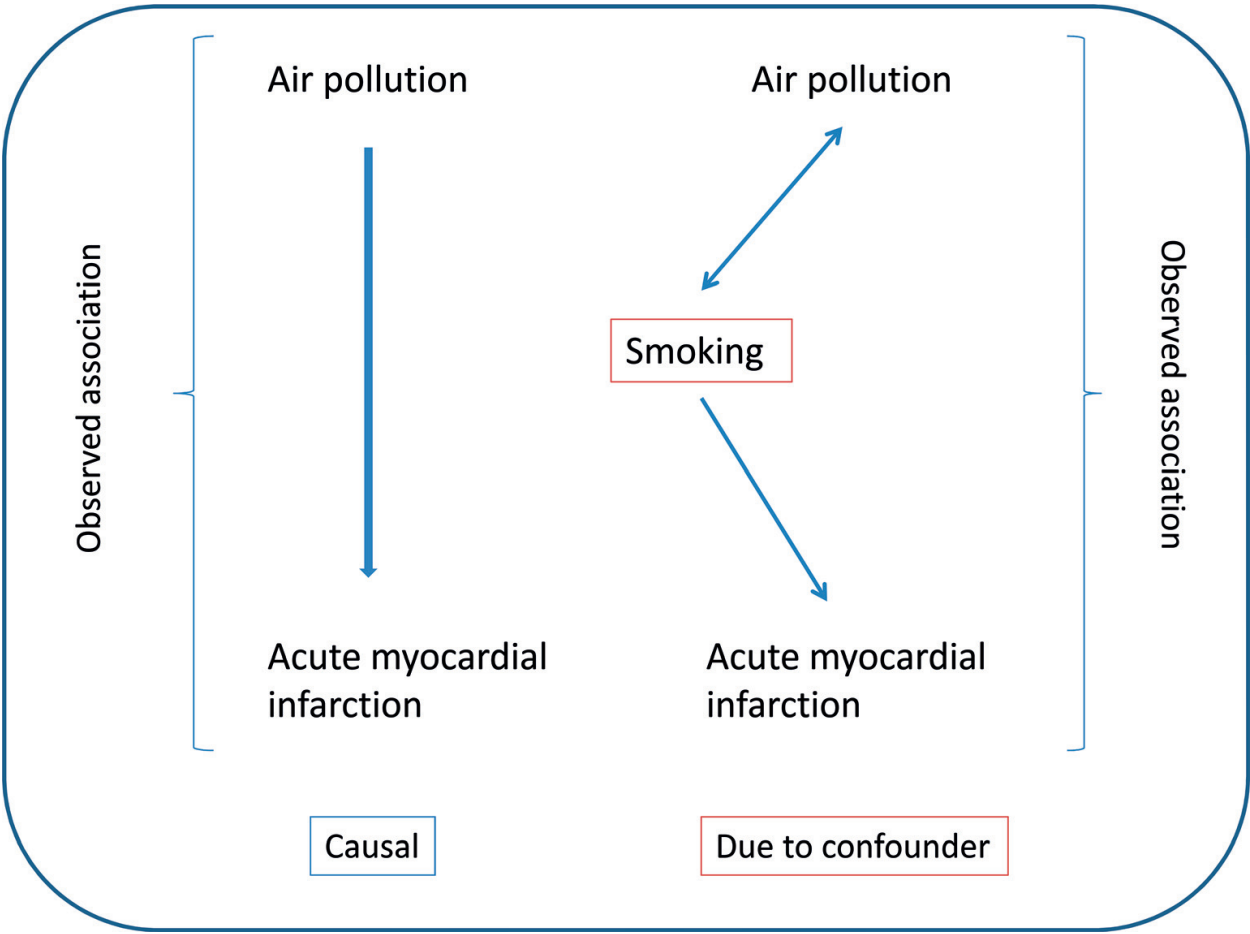
Several examples can be given in this matter: (1) In the aforementioned study about acetylsalicylic acid exposure and major bleeding, confounders considered were age, sex, previous hospitalization for alcoholism, non-bleeding ulcer disease, other non-bleeding conditions, and comorbidities [15]; then the researchers could adjust the risk ratios according to those variables. (2) In the study relating endometriosis and infertility, the considered confounders were menstrual cycle pattern, hirsutism, participant's birthweight, race, household income, husband's education, BMI at age 18 years, alcohol consumption, oral contraception use, any analgesic use, health screening behavior, personal history of cardiovascular disease, and personal history of diabetes [23]. In that case, researchers could evaluate if any of those

variables truly acted as confounders or not. (3) In a study relating air pollution and mortality risks, past exposure to pollution and habit of tobacco consumption were not considered; because they may act as confounders for the results of the association and causal inference could be misled [24]. For a better understanding, please refer to **Figure 1**.

There are some cases in which the dose of exposure may introduce confusion. Such is the case of age, smoking, or drinking alcohol. Age involves the quantity of years (or months) itself, and for the analysis, it will be possible to use as a continuous variable or create ranges of age for stratification analysis. On the other hand, for consumption, it is recommended to consider registering the quantities being consumed by the individuals, so strata can be made during the analysis phase. Tobacco may be registered by number of cigarettes per day. The case of alcohol is rather difficult. The suggestion is to ask for quantity (number of glasses) and types of drinks consumed and then transform it into grams of pure alcohol consumed daily or weekly.

Confounding variables can be controlled in several ways: restriction, matching, stratification and more sophisticated multivariate techniques [2].

Restriction is a simple way of avoiding the introduction of already known confounders, by excluding people who present that factor from the beginning. The problem is that this could



**Figure 1.** Scheme of confounding (smoking) in relation to the exposure (air pollution) and the outcome (acute myocardial infarction).

limit recruitment and the representativeness of certain population groups. So, while increasing the internal validity, this may reduce external validity [2].

As it was said before, matching is rather difficult to do in prospective studies because the first enrolment criterion is the exposure. Matching normally is used in case control studies, but researchers could emphasize that the proportion of women and men would be 50% each or that a ratio of young/old people was similar in the exposure groups.

Stratification is a simple statistic technique that could be used during analysis, but that requires forethought concerning the possible confusion variables and registering them.

The technique consists of separating the analysis of association, according to strata of the confusing factor, for example, perform separate analyses of men and women when it is suspected that gender may be a confounder. Then, when the difference between the calculated raw risk and the risk calculated by strata is over 15%, we could say that confusion is present.

Let's see an example. In a rural area, a study proposed to evaluate the relationship between indoor exposure to smoke—from the combustion of wood stoves—and the occurrence of tuberculosis (TB). The results obtained are presented in **Table 4**.

Therefore, the factor *indoor exposure to wood smoke* for food cooking turned out to be positively associated with the disease. In other words, the incidence of the disease among the exposed group was significantly more than two times greater than that among the nonexposure group.

Given that there is a suspicion that cigarette smoking could modify the effect of indoor contamination on the risk of acquiring tuberculosis, smoking habits were considered. Then, it was possible to assess if this condition acted as a confounder in the association between tuberculosis and indoor smoke exposure using stratification.

The stratification is shown below, where the smoking habit was coded as “never” or “past or present” (**Tables 5 and 6**).

As a conclusion of the stratification results, the factor *indoor exposure to smoke* was found to be positively and significantly associated with the disease, both in nonsmokers and in past or present smokers. However, in past or present smokers, the risk of suffering from tuberculosis is 44% higher than in nonsmokers (from 2.57 to 2.12), when the indoor pollution was present, confirming that smoking habit acts like a confounder in the association between indoor smoke and tuberculosis incidence.

Indoor exposure to smoke	Tuberculosis		
	Disease	No disease	Total
Exposed	50	21	71
Nonexposed	238	524	762
Total	288	545	833

The RR calculation is presented as:  $RR = (50/71)/(238/762) = 0.704/0.312 = 2.25$ .

**Table 4.** Tuberculosis and indoor exposure to smoke from wood burning.

Indoor exposure to smoke	Tuberculosis		
	Disease	No disease	Total
Exposed	33	17	50
Nonexposed	186	411	597

Risk ratio calculation among *never smokers*:  $RR = (33/50)/(186/597) = 0.66/0.311 = 2.12$ .

**Table 5.** TB and indoor exposure to smoke from wood burning of *never smokers*.

Indoor exposure to smoke	Tuberculosis		
	Disease	No disease	Total
Exposed	17	4	21
Nonexposed	52	113	165

Risk ratio calculation among *smokers past or present*:  $RR = (17/21)/(52/165) = 0.81/0.315 = 2.57$ .

**Table 6.** TB and indoor exposure to smoke from wood burning of *past or present smokers*.

Smoking past or present	Tuberculosis		
	Disease	No disease	Total
Yes	52	113	165
No smoking	186	411	597

Risk calculation of smoking among nonexposed to smoke from wood burning:  $RR = (52/165)/(186/597) = 0.315/0.311 = 1.01$ .

**Table 7.** TB and smoking habits (without indoor exposure to smoke from wood burning).

In addition to confusion, we have the concept of *interaction* that refers to the effect that two of more factors have by increasing or reducing the incidence of a disease when they are together. Then, the incidence resulting when the factors are together differs from the incidence when the factors are isolated.

Let us try to find interaction in the same example.

In order to assess interaction, it will be necessary to calculate the association between smoking and tuberculosis alone (without the indoor exposure to wood burning smoke) **Table 7**.

The result shows that the relative risk of developing TB due exclusively to the habit of smoking is almost nil. But, to know if there is interaction, we should estimate if the presence of both exposures together differs or not from the expected effects if the two exposures were simply the sum of both.

From the previous tables and calculations, we have that the incidences are the following:



- Incidence rate of TB without any smoke exposure = 31.1%
- Incidence rate of TB with smoking only = 31.5%
- Incidence rate of TB with indoor pollution only = 66%
- Incidence rate of TB with smoking and indoor pollution = 81%

In order to know whether interaction is present, we should clear the incidences from the underground risk of developing TB (baseline incidence). Then, we should start by calculating the attributable risks (ARs), as follows:

AR to smoking = (TB incidence due to smoking – baseline TB incidence) =  $31.5 - 31.1 = 0.4$ .

AR to indoor pollution = (TB incidence if indoor pollution – baseline TB incidence) =  $66 - 31.1 = 34.9$ .

The expected attributable risk to both factors would be the addition of the TB incidence of (smoking + indoor contamination) =  $34.9 + 0.4 = 35.3\%$ . Then, the expected incidence will be  $(31.1 + 35.3) = 66.4\%$ .

But the real TB incidence with both exposure factors was 81%. The difference between 81 and 66.4 would be attributable to the interaction, which is 14.6%.

In other way, the incidence when both factors are together is higher than the addition of incidences when the factors are alone, taking into consideration that we have to clear the underground risk (incidence of TB in population free of exposures).

Effectively, we have shown that interaction is present, because the incidences of both exposures together differ from the expected effects if the two exposures were simply the sum of both. As a conclusion, the indoor pollution is a risk factor to develop TB in that setting, but this risk increases substantially more if people smoke indoors. For a better understanding, please refer to **Figure 2**.

Coming back to the control of confusion bias, adjustment techniques using statistical models require computer training and have the advantage of working with two or more possible confounding variables; opposite to stratification that permits working on one factor only. When using modeling multivariate techniques, logistic regression or proportional hazard regression might be used, but researchers must be aware of how to interpret the results properly [2].

An important comment about the confusion is that finding a confounder is not always an issue to be worried about. It could also be useful. For example, in the mentioned study about pesticide exposure in agricultural workers and cognitive impairment, gender turned out to be a confounder [7]. This resulted from the type of work performed differing between men and women. Men used to perform tasks like mixing, blending and applying pesticides; while women pruned to collecting fruits, so men were directly exposed to the toxins. Then, knowing that men were more exposed and, consequently, more susceptible to the health damage, the preventive measures may be oriented by strengthening them toward men, but still keeping care on women.

Finally, confounders are not a mistake in the research, but a phenomenon that is present must be understood by the investigators in order to finally consider them when interpreting the results of the study [4].



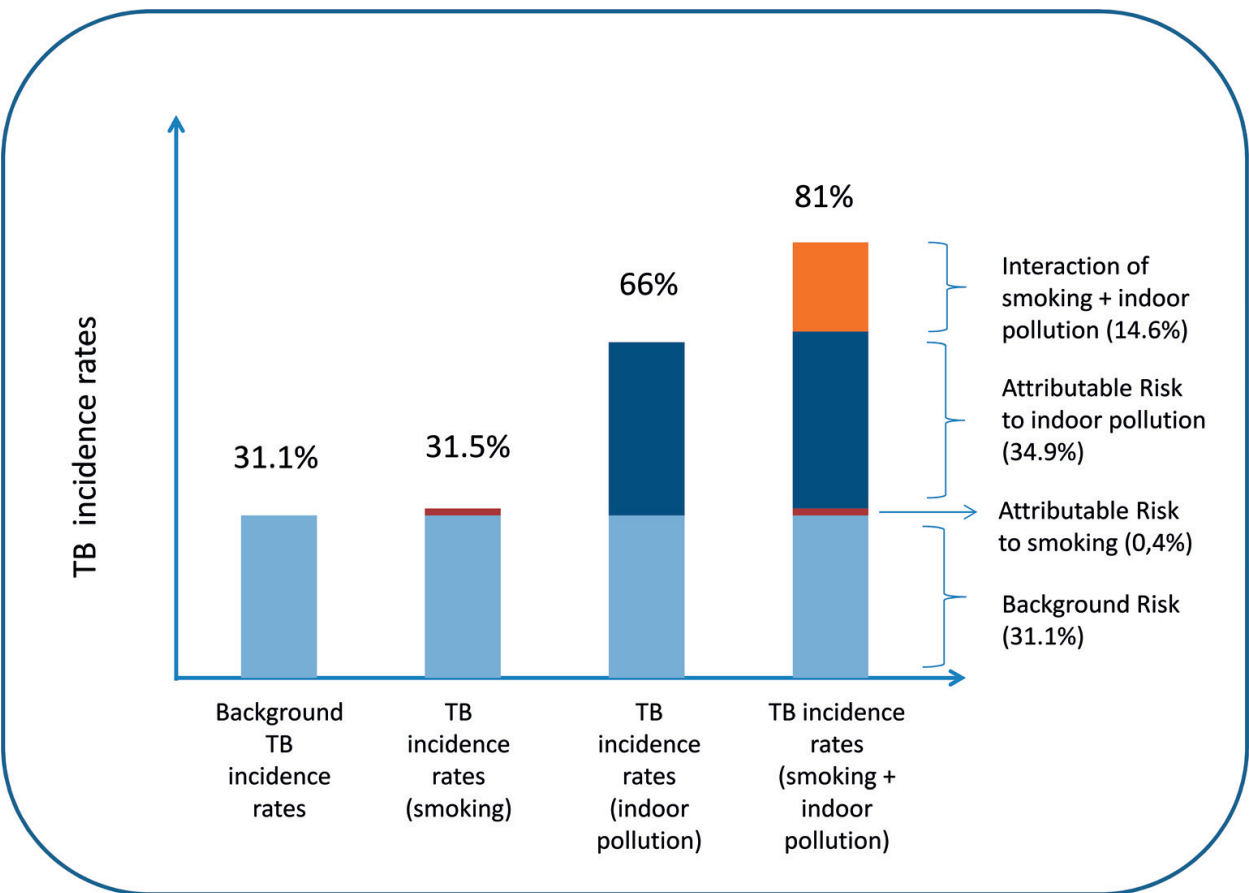


Figure 2. Incidence rates and attributable risk to factors related to TB incidence and their interaction.

### 3. Conclusions

As shown, biases can be present in any study, originating from multiple steps of the investigation. Their presence should not be grounds for rejection of the results due to the poor quality of the study, but careful attention is required when interpreting the results. To the extent that the researcher is able to recognize the biases, he/she can be proactive in mitigating them, either by way of improving the design or applying statistical techniques (stratification or multivariate adjustment) when analyzing the results. Therefore, when clinicians or researchers look for good quality of articles to read and use as references, they must recognize them when interpreting the results and acknowledge the limitations that the studies may have.

It should be noted that biases are more frequent among retrospective cohort, given by missing information when using existing records (information bias) or by selection bias, because individuals are selected after the outcome has occurred, so both conditions (exposure and outcome) are present at the moment of enrollment. In that case, it is easier that exposed or unexposed subjects would be related to the result of interest, causing selection bias. On the other hand, a prospective cohort design could be affected by the loss of follow-up. Both types of cohort studies may be influenced by information bias, confusion or interaction.

Interesting tools for weighting quality and predisposition to unfairness in observational studies have been gathered and reported by Sanderson et al. [25]. Those included items for selection methods, measurement of study variables, design-specific sources of bias, control of confounding variables and use of statistics.

Finally, it is considered that cohort studies are used normally as a source of information of systematic reviews and meta-analysis. In those cases, publication bias and outcome reporting bias must be taken into consideration. This is because the journals are prone to publish positive results rather than negative ones, a situation that has been shown [26].

## Conflict of interest

The author of this chapter declares no conflict of interest.

## Author details

Muriel Ramirez-Santana

Address all correspondence to: [mramirezs@ucn.cl](mailto:mramirezs@ucn.cl)

Public Health Department, Faculty of Medicine, Universidad Catolica del Norte, Coquimbo, Chile

## References

- [1] Richiardi L, Pizzi C, Pearce N. Commentary: Representativeness is usually not necessary and often should be avoided. *International Journal of Epidemiology*. 2017;**42**:1018-1022
- [2] Grimes DA, Schulz KF. Bias and causal associations in observational research. *Lancet*. 2002;**359**:248-252
- [3] Mantel N. Avoidance of bias in cohort studies. *National Cancer Institute Monograph*. May 1985;**67**:169-172
- [4] Gordis L. *Epidemiology*. 4th ed. Philadelphia: Elsevier; 2009
- [5] Roy MN, Vollset SE, Gjessing HK, Skjærven R, Melve KK, Schreuder P, Alsaker ER, Haug K, Daltveit AK, Per Magnus. Self selection bias in a large prospective pregnancy cohort in Norway. *Paediatric and Perinatal Epidemiology* 2009;**23**:507-608
- [6] Ercumen A, Arnold BF, Naser AM, Unicomb L, Colford JM, Luby SP. Potential sources of bias in the use of *Escherichia coli* to measure waterborne diarrhoea risk in low-income settings. *Tropical Medicine & International Health*. 2017;**22**:2-11

- [7] Ramírez-Santana M, Zúñiga L, Corral S, Sandoval R, Scheepers PT, Van Der Velden K, Roeleveld N, Pancetti F. Assessing biomarkers and neuropsychological outcomes in rural populations exposed to organophosphate pesticides in Chile—Study design and protocol Environmental and occupational health. *BMC Public Health*. 2015;**15**:116. DOI: 10.1186/s12889-015-1463-5
- [8] Kaaks R, Plummer M, Riboli E, Estève J, Van Staveren W. Adjustment for bias due to errors in exposure assessments in multicenter cohort studies on diet and cancer: A calibration approach. *The American Journal of Clinical Nutrition*. 1994;**59**:2455-2505
- [9] Copeland KT, Checkoway H, McMichael AJ, Holbrook RH. Bias due to misclassification in the estimation of relative risk. *American Journal of Epidemiology*. 1997;**105**(5):188-495. DOI: <https://doi.org/10.1093/oxfordjournals.aje.a112408>
- [10] Kristman V, Manno M, Côte P. Loss to follow-up in cohort studies: how much is too much? *European Journal of Epidemiology*. 2004;**19**:751-760
- [11] Lévesque Linda E, Hanley James A, Kezouh Abbas SS. Problem of immortal time bias in cohort studies: Example using statins for preventing progression of diabetes. *British Medical Journal*. 2010;**340**:b5087. DOI: <https://doi.org/10.1136/bmj.b5087>
- [12] Suissa S. Lower risk of death with SGLT2 inhibitors in observational studies: Real or bias? *Diabetes Care*. Jan. 2018;**41**:6-10
- [13] Applebaum KM, Malloy EJ, Eisen EA. Left truncation, susceptibility, and bias in occupational cohort studies. *NIH Public Access*. 2014;**22**:599-606
- [14] Buckley BS, Simpson CR, McLernon DJ, Hannaford PC, Murphy AW. Considerable differences exist between prevalent and incident myocardial infarction cohorts derived from the same population. *Journal of Clinical Epidemiology*. 2010;**63**:1351-1357
- [15] Pedersen L, Stürmer T. Conditioning on future exposure to define study cohorts can induce bias: The case of low-dose acetylsalicylic acid and risk of major bleeding. *Clinical Epidemiology*. 2017;**9**:611-626
- [16] Sedgwick P. Cohort studies: Source of bias. *British Medical Journal*. 2011;**343**:d7839. DOI: <https://doi.org/10.1136/bmj.d7839>
- [17] Deckert A. The existence of standard-biased mortality ratios due to death certificate misclassification - A simulation study based on a true story. *BMC Medical Research Methodology*. 2016;**16**:1-9
- [18] Walraven C. A comparison of methods to correct for misclassification bias from administrative database diagnostic codes. *International Journal of Epidemiology*. 2017;**0**:1-12
- [19] Johnson CY, Flanders WD, Strickland MJ, Honein MA, Howards PP. Potential sensitivity of bias analysis results to incorrect assumptions of nondifferential or differential binary exposures misclassification. *Epidemiology*. 2014;**15**:902-909

- [20] McCambridge J, Witton J, Elbourne DR. Systematic review of the Hawthorne effect: New concepts are needed to study research participation effects. *Journal of Clinical Epidemiology*. 2014;**67**:267-277
- [21] Henry SG, Jerant A, Iosif A-M, Feldman MD, Cipri C, Kravitz RL. Analysis of threats to research validity introduced by audio recording clinic visits: Selection bias, Hawthorne effect, both, or neither? *Diagnostic Microbiology and Infectious Disease*. 2016;**28**:1304-1314
- [22] Leurent B, Reyburn H, Muro F, Mbakilwa H, Schellenberg D. Monitoring patient care through health facility exit interviews: An assessment of the Hawthorne effect in a trial of adherence to malaria treatment guidelines in Tanzania. *BMC Infectious Diseases*. 2016;**16**(59):1-9
- [23] Prescott J, Farland LV, Tobias DK, Gaskins AJ, Spiegelman D, Chavarro JE, Rich-Edwards JW, Barbieri RL, Missmer SA. A prospective cohort study of endometriosis and subsequent risk of infertility. *Human Reproduction*. 2016;**31**:1475-1482
- [24] Hansell A, Ghosh RE, Blangiardo M, Perkins C, Vienneau D, Goffe K, Briggs D, Gulliver J. Historic air pollution exposure and Long-term mortality risks in England and Wales: Prospective longitudinal cohort study. *Thorax*. 2016;**71**:330-338
- [25] Sanderson S, Tatt ID, Higgins JPT. Tools for assessing quality and susceptibility to bias in observational studies in epidemiology : A systematic review and annotated bibliography. *International Journal of Epidemiology*. 2007;**36**:666-676
- [26] Dwan K, Altman DG, Arnaiz JA, et al. Systematic Review of the empirical evidence of study publication bias and outcome reporting bias. *PLoS One*. 2008;**3**:1-30

