

---

Certificate Programme in Healthcare Research Analysis

## Epidemiology and Biostatistics

---

### Epidemiology:

Epidemiology is the study of the distribution and determinants of health-related states or events in specified populations and the application of this study to control health problems. It involves the study of patterns of disease and injury in human populations and the factors that influence these patterns. Epidemiologists conduct research to understand the causes of disease and other health outcomes and to develop strategies for prevention and control. Epidemiology plays a critical role in public health by helping to identify risk factors for disease and injury and by providing evidence to guide public health interventions.

### Biostatistics:

Biostatistics is the application of statistical methods to biological and health-related data. It involves the design of experiments and studies, the collection and analysis of data, and the interpretation of results. Biostatisticians work collaboratively with researchers in various fields to design studies that produce valid and reliable results. They also help researchers analyze and interpret data to draw meaningful conclusions.

### Confounding Variable:

A confounding variable is a variable that distorts the true relationship between the independent and dependent variables in a study. Confounding variables can lead to incorrect conclusions if not properly controlled for in the analysis. For example, if a study is examining the relationship between smoking and lung cancer but fails to account for age, which is a known risk factor for lung cancer, age could act as a confounding variable and distort the results.

### Descriptive Epidemiology:

Descriptive epidemiology is the branch of epidemiology that focuses on describing the distribution of health-related states or events in a population. It involves characterizing the patterns of disease and injury by time, place, and person. Descriptive epidemiology provides a basic understanding of the occurrence and distribution of health outcomes, which can inform further research and interventions.

### Incidence:

Incidence is a measure of the rate at which new cases of a disease or condition occur in a population over a specified period of time. It is typically expressed as the number of new cases per unit of population at risk. Incidence is an important measure in epidemiology because it helps researchers understand the risk of developing a particular disease or condition within a population.

### Prevalence:

Prevalence is a measure of the proportion of individuals in a population who have a specific disease or condition at a particular point in time. It is typically expressed as a percentage or a ratio. Prevalence provides information about the burden of a disease within a population and is useful for planning and evaluating public health interventions.

**Relative Risk:**

Relative risk is a measure of the strength of association between an exposure and a disease or condition. It compares the risk of developing the disease in individuals who are exposed to a particular factor to the risk in individuals who are not exposed. A relative risk of 1 indicates no association, a relative risk greater than 1 indicates an increased risk, and a relative risk less than 1 indicates a decreased risk.

**Attributable Risk:**

Attributable risk is a measure of the proportion of disease occurrence in a population that can be attributed to a specific exposure. It quantifies the additional risk of disease that is associated with the exposure compared to the risk in the absence of the exposure. Attributable risk provides information about the potential impact of removing or reducing a particular risk factor on the incidence of disease.

**Confidence Interval:**

A confidence interval is a range of values that is used to estimate the true value of a population parameter with a certain level of confidence. It provides a measure of the precision of an estimate and indicates the range within which the true value is likely to lie. For example, a 95% confidence interval indicates that there is a 95% probability that the true value of the parameter falls within the specified range.

**Hazard Ratio:**

The hazard ratio is a measure of the relative risk of an event occurring in one group compared to another group over time. It is commonly used in survival analysis to compare the risk of an event (such as death or disease recurrence) between two or more groups. A hazard ratio greater than 1 indicates an increased risk of the event in the first group compared to the second group.

**Randomized Controlled Trial (RCT):**

A randomized controlled trial is a study design in which participants are randomly assigned to receive one of two or more interventions. It is considered the gold standard for evaluating the effectiveness of a treatment or intervention because randomization helps to eliminate bias and confounding. RCTs are used to assess the causal relationship between an intervention and an outcome.

**Case-Control Study:**

A case-control study is an observational study design that compares individuals with a specific disease or condition (cases) to individuals without the disease or condition (controls) to identify factors that may be associated with the disease. Case-control studies are useful for investigating rare diseases or conditions and can provide valuable information about potential risk factors.

**Cohort Study:**

A cohort study is an observational study design that follows a group of individuals over time to evaluate the association between exposures and outcomes. Participants in a cohort study are classified based on their exposure status at the beginning of the study and are followed to assess the development of the outcome. Cohort studies are useful for investigating the natural history of diseases and for identifying risk factors.

**Meta-Analysis:**

A meta-analysis is a statistical technique that combines the results of multiple studies on a particular topic

to derive a single summary estimate of effect. It allows researchers to synthesize evidence from different studies and obtain a more precise estimate of the true effect size. Meta-analyses are commonly used in epidemiology and biostatistics to summarize the findings of multiple studies and to identify patterns and trends across studies.

#### Population Health:

Population health is a field of study that focuses on the health outcomes of groups of individuals within a population. It examines the distribution of health outcomes across populations and the factors that influence these outcomes. Population health considers social, economic, and environmental determinants of health and seeks to improve the health of entire populations through interventions at the community level.

#### Validity:

Validity refers to the degree to which a study accurately measures what it is intended to measure. It is an essential concept in epidemiology and biostatistics because valid study results are essential for drawing accurate conclusions and making informed decisions. Validity can be assessed in terms of internal validity (the extent to which the study design minimizes bias) and external validity (the generalizability of study results to other populations or settings).

#### Reliability:

Reliability refers to the consistency and repeatability of study results. A reliable study produces consistent results when repeated under similar conditions. Reliability is important in epidemiology and biostatistics because it ensures that study findings are trustworthy and can be replicated by other researchers. Reliability can be assessed through measures such as test-retest reliability and inter-rater reliability.

#### Propensity Score Matching:

Propensity score matching is a statistical technique used to reduce bias in observational studies by matching individuals with similar characteristics. It involves estimating the probability of receiving a particular treatment based on a set of covariates and then matching individuals who have similar propensity scores. Propensity score matching helps to balance the distribution of confounding variables between treatment groups and improve the validity of study results.

#### Survival Analysis:

Survival analysis is a statistical method used to analyze the time until an event of interest occurs, such as death, disease recurrence, or recovery. It is commonly used in epidemiology and biostatistics to estimate survival probabilities over time and to compare survival curves between groups. Survival analysis accounts for censored data (individuals who have not experienced the event by the end of the study) and provides valuable information about the timing of events.

#### Power Analysis:

Power analysis is a statistical technique used to determine the sample size required to detect a significant effect in a study. It involves calculating the statistical power of a study, which is the probability of correctly rejecting a null hypothesis when the alternative hypothesis is true. Power analysis is important in epidemiology and biostatistics because studies with insufficient power may fail to detect true effects, leading to inconclusive results.

**Confidence Level:**

The confidence level is the probability that a confidence interval will contain the true value of a population parameter. It is commonly expressed as a percentage, such as 95% or 99%. A higher confidence level indicates greater certainty that the true value falls within the specified range. For example, a 95% confidence level means that there is a 95% probability that the confidence interval contains the true value of the parameter.

**Cluster Randomized Trial:**

A cluster randomized trial is a study design in which groups of individuals, rather than individual participants, are randomly assigned to receive different interventions. It is commonly used in public health research when it is not feasible to randomize individuals. Cluster randomized trials are used to evaluate the effectiveness of interventions at the community or population level and account for clustering effects within groups.

**Selection Bias:**

Selection bias occurs when the selection of study participants is not random and is related to both the exposure and the outcome of interest. It can lead to an overestimation or underestimation of the true association between variables in a study. Selection bias is a common source of error in epidemiological research and can distort study results if not properly controlled for.

**Publication Bias:**

Publication bias occurs when research findings are systematically distorted or suppressed due to the selective publication of studies with positive results. It can lead to an overestimation of the true effect size of an intervention or exposure. Publication bias is a significant challenge in epidemiology and biostatistics because it can skew the evidence base and lead to misleading conclusions in systematic reviews and meta-analyses.

**Bayesian Statistics:**

Bayesian statistics is a method of statistical inference that uses Bayes' theorem to update the probability of a hypothesis as new evidence becomes available. It incorporates prior knowledge or beliefs about the likelihood of an event occurring into the analysis. Bayesian statistics is used in epidemiology and biostatistics to estimate parameters, make predictions, and perform sensitivity analyses based on available data and expert judgment.

**Confounding Bias:**

Confounding bias occurs when a third variable influences the relationship between an exposure and an outcome in a study. It can lead to a spurious association between variables if not properly controlled for in the analysis. Confounding bias is a common source of error in epidemiological research and can result in incorrect conclusions about the causal relationship between variables.

**Regression Analysis:**

Regression analysis is a statistical technique used to model the relationship between one or more independent variables and a dependent variable. It helps researchers understand how changes in the independent variables are associated with changes in the dependent variable. Regression analysis is

commonly used in epidemiology and biostatistics to quantify the effect of exposures on outcomes and to control for potential confounding variables.

#### Sensitivity Analysis:

Sensitivity analysis is a statistical method used to assess the robustness of study results to changes in assumptions or methods. It involves varying key parameters or assumptions in the analysis to determine their impact on the results. Sensitivity analysis helps researchers evaluate the reliability and validity of study findings and provides insights into the potential sources of bias or uncertainty.

#### Specificity:

Specificity is a measure of the proportion of true negative results in a diagnostic test. It indicates the ability of a test to correctly identify individuals who do not have a particular disease or condition. Specificity is an important measure in epidemiology and biostatistics because it helps evaluate the accuracy of a diagnostic test and the likelihood of false-negative results.

#### Receiver Operating Characteristic (ROC) Curve:

A receiver operating characteristic curve is a graphical plot that illustrates the performance of a binary classifier system as its discrimination threshold is varied. It is commonly used in epidemiology and biostatistics to evaluate the diagnostic accuracy of a test. The ROC curve plots the true positive rate against the false positive rate at various threshold settings and provides information about the trade-off between sensitivity and specificity.

#### Bayesian Network:

A Bayesian network is a graphical model that represents probabilistic relationships between variables in a directed acyclic graph. It uses Bayesian statistics to model the conditional dependencies between variables and to make inferences about the joint probability distribution of the variables. Bayesian networks are used in epidemiology and biostatistics to model complex systems and to assess the impact of interventions on outcomes.

#### Negative Predictive Value:

Negative predictive value is a measure of the proportion of true negative results in a diagnostic test among individuals who test negative. It indicates the probability that a negative test result is correct and that an individual does not have the disease or condition. Negative predictive value is important in epidemiology and biostatistics because it helps assess the reliability of a diagnostic test in ruling out disease.

#### Positive Predictive Value:

Positive predictive value is a measure of the proportion of true positive results in a diagnostic test among individuals who test positive. It indicates the probability that a positive test result is correct and that an individual has the disease or condition. Positive predictive value is important in epidemiology and biostatistics because it helps assess the reliability of a diagnostic test in confirming disease.

#### Kaplan-Meier Estimator:

The Kaplan-Meier estimator is a non-parametric method used to estimate the survival function in survival analysis. It provides an estimate of the probability of survival over time based on observed data, even when

some individuals are censored. The Kaplan-Meier estimator is commonly used in epidemiology and biostatistics to analyze time-to-event data and to compare survival curves between groups.

#### Cox Proportional Hazards Model:

The Cox proportional hazards model is a regression model used in survival analysis to assess the relationship between covariates and the hazard rate of an event occurring. It allows researchers to estimate the effect of multiple variables on the time to an event while accounting for censoring. The Cox proportional hazards model is widely used in epidemiology and biostatistics to analyze survival data and to identify risk factors for disease outcomes.

#### Case-Cohort Study:

A case-cohort study is a hybrid study design that combines elements of a case-control study and a cohort study. In a case-cohort study, a random sample of the cohort is selected as the comparison group, allowing for the estimation of incidence rates and risk ratios. Case-cohort studies are useful for investigating rare diseases or conditions and for examining multiple outcomes within the same study population.

#### Propensity Score:

A propensity score is a conditional probability that estimates the likelihood of an individual receiving a particular treatment based on a set of covariates. Propensity scores are commonly used in observational studies to balance the distribution of confounding variables between treatment groups. Matching on propensity scores helps to reduce bias and improve the validity of study results.

#### Survival Function:

The survival function is a mathematical representation of the probability that an individual survives beyond a certain time point. It is commonly estimated in survival analysis using methods such as the Kaplan-Meier estimator or the Cox proportional hazards model. The survival function provides valuable information about the probability of surviving over time and allows researchers to compare survival curves between groups.

#### Exposure:

An exposure is a factor or variable that is hypothesized to influence the risk of developing a disease or condition. In epidemiology and biostatistics, exposures can include lifestyle factors, environmental exposures, genetic factors, or interventions. Understanding the relationship between exposures and outcomes is essential for identifying risk factors for disease and developing effective prevention strategies.

#### Outcome:

An outcome is a health-related event or condition that is of interest in a study. In epidemiology and biostatistics, outcomes can include diseases, injuries, symptoms, or other health-related states. Researchers investigate the relationship between exposures and outcomes to understand the factors that influence health outcomes and to develop interventions to improve health.

#### Confidence Interval:

A confidence interval is a range of values that is used to estimate the true value of a population parameter with a certain level of confidence. It provides a measure of the precision of an estimate and indicates the range within which the true value is likely to lie. For example, a 95% confidence interval indicates that there

---

is a 95% probability that the true value of the parameter falls within the specified range.

#### Hazard Ratio:

The hazard ratio is a measure of the relative risk of an event occurring in one group compared to another group over time. It is commonly used in survival analysis to compare the risk of an event (such as death or disease recurrence) between two or more groups. A hazard ratio greater than 1 indicates an increased risk of the event in the first group compared to the second group.

#### Randomized Controlled Trial (RCT):

A randomized controlled trial is a study design in which participants are randomly assigned to receive one of two or more interventions. It is considered the gold standard for evaluating the effectiveness of a treatment or intervention because randomization helps to eliminate bias and confounding. RCTs are used to assess the causal relationship between an intervention and an outcome.

#### Case-Control Study:

A case-control study is an observational study design that compares individuals with a specific disease or condition (cases) to individuals without the disease or condition (controls) to identify factors that may be associated with the disease. Case-control studies are useful for investigating rare diseases or conditions and can provide valuable information about potential risk factors.

#### Cohort Study:

A cohort study is an observational study design that follows a group of individuals over time to evaluate the association between exposures and outcomes. Participants in a cohort study are classified based on their exposure status at the beginning of the study and are followed to assess the development of the outcome. Cohort studies are useful for investigating the natural history of diseases and for identifying risk factors.

#### Meta-Analysis:

A meta-analysis is a statistical technique that combines the results of multiple studies on a particular topic to derive a single summary estimate of effect. It allows researchers to synthesize evidence from different studies and obtain a more precise estimate of the true effect size. Meta-analyses are commonly used in epidemiology and biostatistics to summarize the findings of multiple studies and to identify patterns and trends across studies.

#### Population Health:

Population health is a field of study that focuses on the health outcomes of groups of individuals within a population. It examines the distribution of health outcomes across populations and the factors that influence these outcomes. Population health considers social, economic, and environmental determinants of health and seeks to improve the health of entire populations through interventions at the community level.

#### Validity:

Validity refers to the degree to which a study accurately measures what it is intended to measure. It is an essential concept in epidemiology and biostatistics because valid study results are essential for drawing accurate conclusions and making informed decisions. Validity can be assessed in terms of internal validity (the extent to which the study design minimizes bias) and external validity (the generalizability of study

results to other populations or settings).

**Reliability:**

Reliability refers to the consistency and repeatability of study results. A reliable study produces consistent results when repeated under similar conditions. Reliability is important in epidemiology and biostatistics because