

# Choosing a Proper Statistical Method – Part III

Yuda Chongpison, PhD, MS, MBA  
Research Affairs, Faculty of Medicine  
Chulalongkorn University

# Relationship between odds ratio and relative risk

	Disease developed	No Disease developed	Total
Exposed	a	b	a+b
Not Exposed	c	d	c+d
Total	a+c	b+d	a+b+c+d

$$RR = \frac{\frac{a}{a+b}}{\frac{c}{c+d}}$$

$$OR = \frac{\frac{a}{b}}{\frac{c}{d}}$$

*Odds ratio: a good approximation of the relative risk when the disease is rare*

# **ISSUES IN BIOSTATISTICS**

# Review: Hypotheses and Test Statistics

- Step 1: State hypotheses ( $H_0$  and  $H_a$ )
- Step 2: Summarize data (and eyeball it)
- Step 3: Think about what data might look like if  $H_0$  true ('Under  $H_0$ ')
- Step 4: Using data, gauge whether  $H_0$  true.
  - 4a: Calculate test statistic
  - 4b: p-value: evaluate probability (Test statistic under  $H_0$ )
- Step 5: Conclusion

# Type I and II errors

Research	Reality	
	$H_0$ is True (No Difference)	$H_0$ is False (Difference)
Accept $H_0$ (Conclusion: No Difference)	Correct Decision	Type II error ( $\beta$ , False Negative)
Reject $H_0$ (Conclusion: Difference)	Type I error ( $\alpha$ , False Positive)	Correct Decision

$H_0$  = no statistically significant difference between two groups

Based on the primary objective

# **SAMPLE SIZE ESTIMATIONS**

# Sample Size Determination

- Why do we determine sample size
  - To make sure that the results of the *statistical approach align with scientifically useful results* given the constraints of ethics and feasibility
- When is it appropriate to prospectively power a study?
  - Study is prospective
  - When we have information to do it
  - The analysis is sufficiently simple (e.g., bivariate analysis)

# Parameters

- Used to determine an appropriate sample size
  - Significant level
  - Desired statistical power
  - Minimal clinical difference
  - Estimated measurement variability



# Type I and II errors

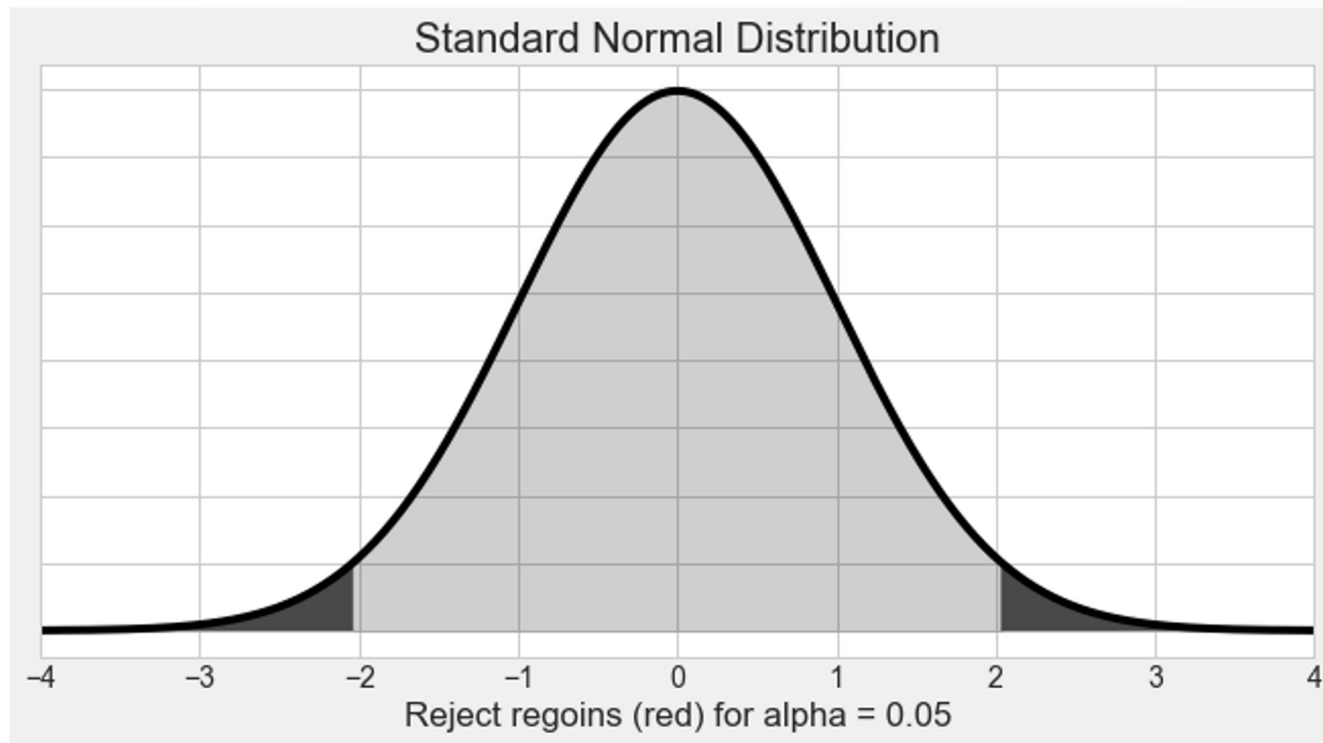
Research	Reality	
	$H_0$ is True (No Difference)	$H_0$ is False (Difference)
Accept $H_0$ (Conclusion: No Difference)	Correct Decision	Type II error ( $\beta$ , False Negative)
Reject $H_0$ (Conclusion: Difference)	Type I error ( $\alpha$ , False Positive)	Correct Decision

$H_0$  = no statistically significant difference between two groups

# Significance Level

- The maximum p value for which a difference is considered statistically significant
- Represents by the symbol  $\alpha$  (alpha)
- Type I error (false positive rate)
- Conventionally set at 0.05

# Significance Level (cont.)

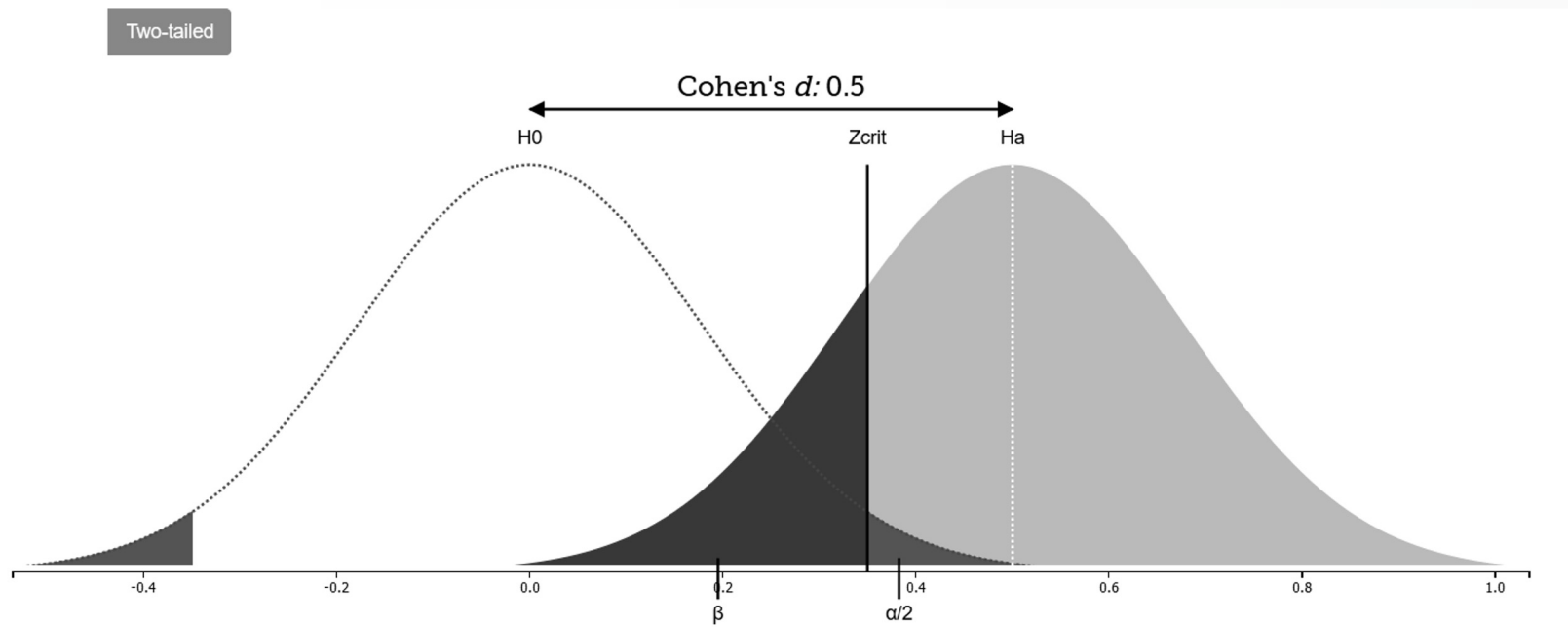


Alpha: Type I error (false positive rate)

# Statistical Power

- The probability that a statistical test will indicate a significant difference when, in reality, there is one.
- Analogous to the sensitivity of a diagnostic test
- The definition of sensitivity:
  - True positive rate
  - Ability of a diagnostic test to correctly diagnose the disease

# Statistical Power (cont.)



Statistical power:  $1 - \beta$

# Minimal Clinical Difference (MCD)

- Sometime called
  - Minimum expected difference
  - Effect size
- Represent
  - Scientifically meaningful difference
  - Smallest difference between comparison groups that the investigator deems clinically significant
- Based on clinical judgement and experience with the issue of interest

# Estimated Variability

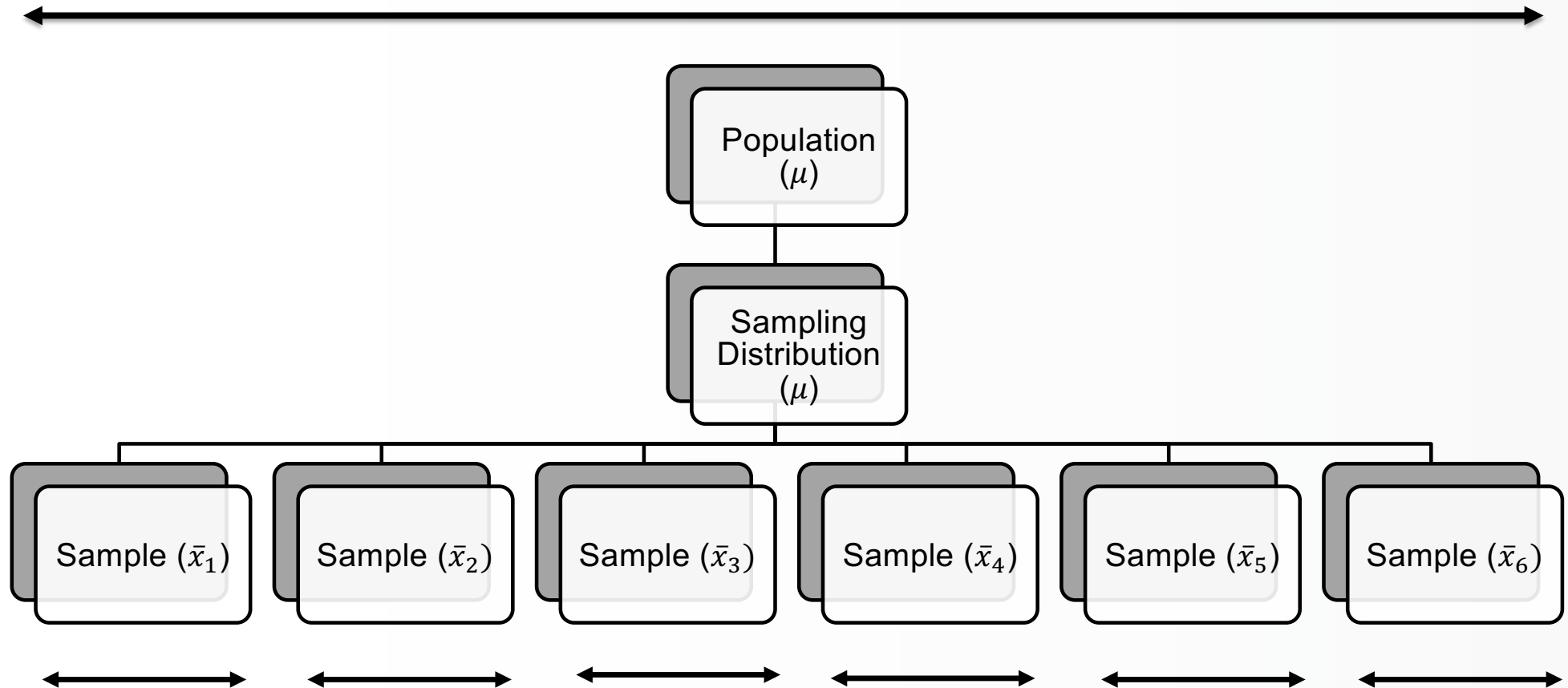
- Expected standard deviation of the outcome or measurements
- Can be determined from
  - Review literature of the same research topic in the similar study population

What's the difference?

# **STANDARD ERROR VS STANDARD DEVIATION**



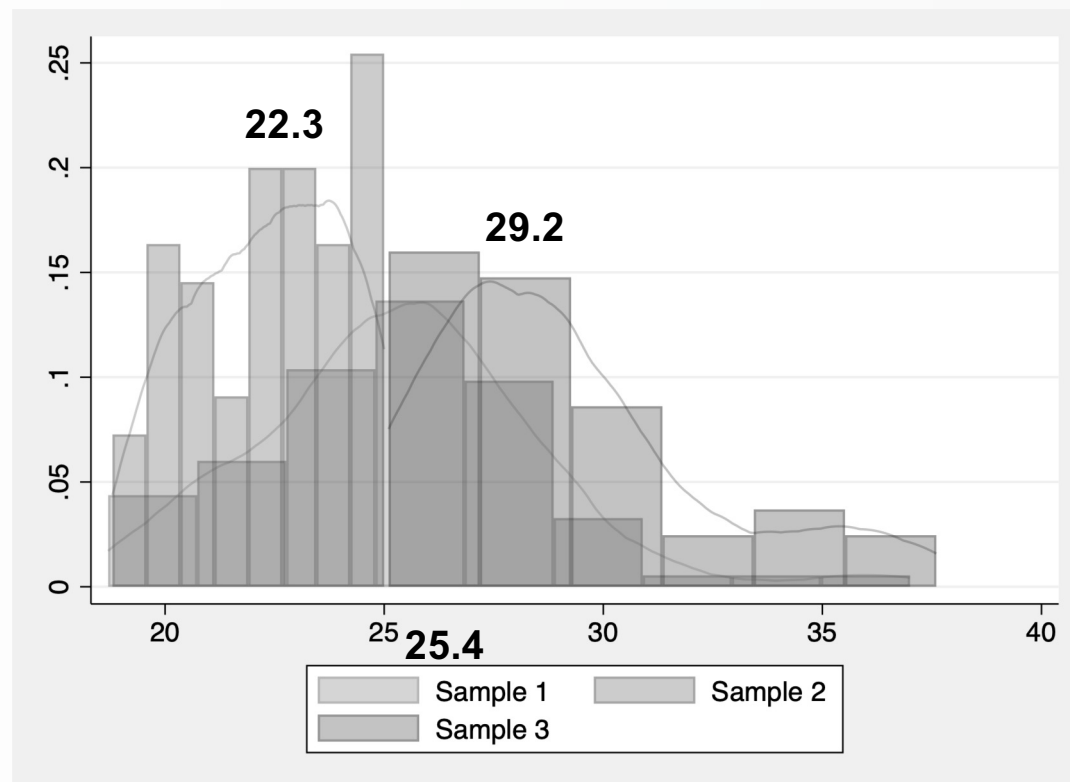
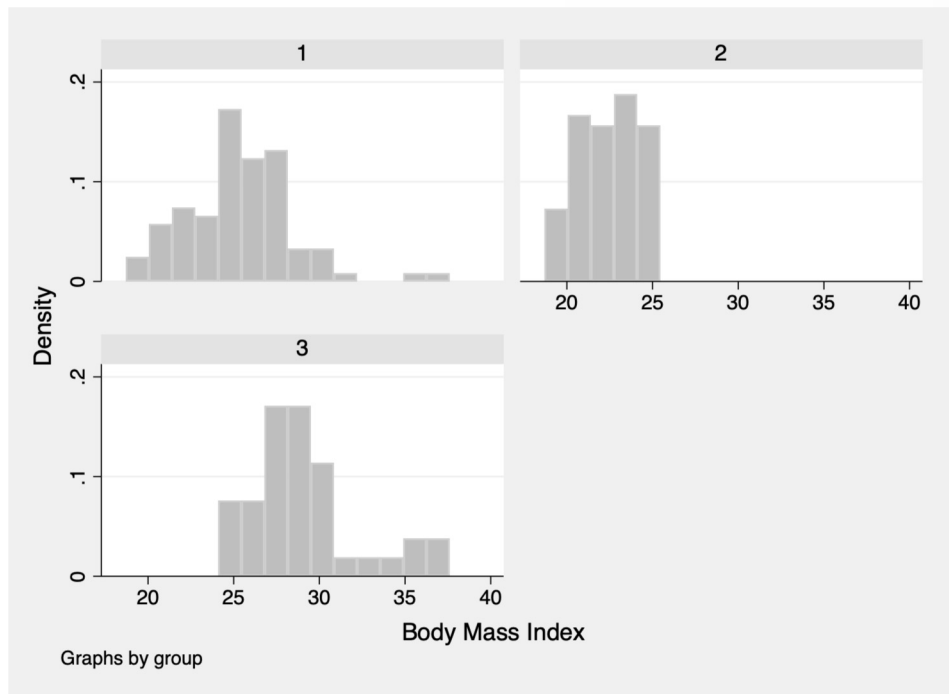
Standard Error: variability across samples



Standard Deviation: variability within a sample

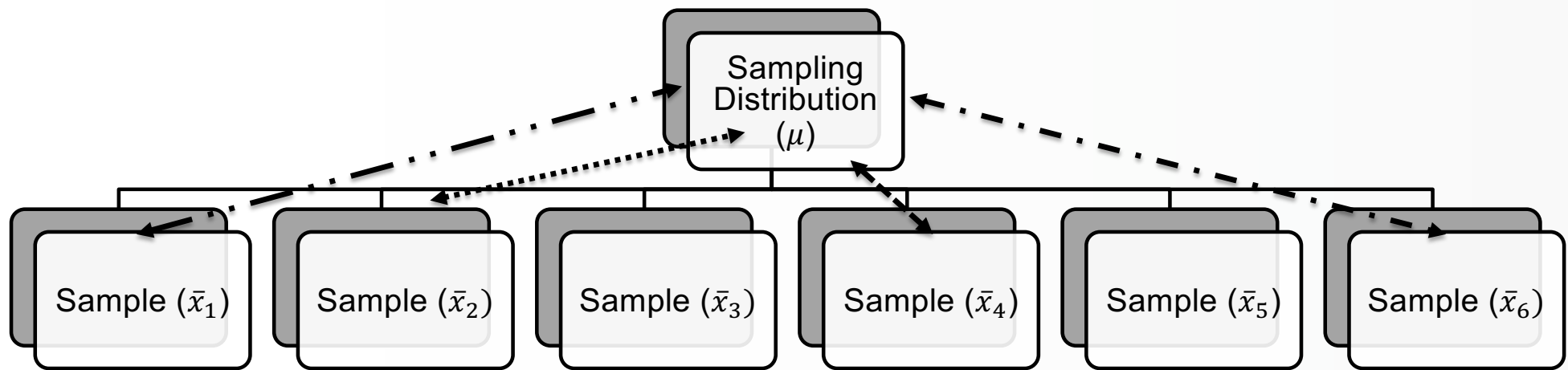
# Standard Deviation

- Standard deviation for each sample
- $s_1 = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x}_1)^2}{n-1}}$
- Measure of the spread of values within a set of data



# Standar Error

- $SE_{\mu_x} = \frac{s}{\sqrt{n}}$
- The standard deviation of its sampling distribution (a sample distribution)
- A sample mean deviates from the actual mean of a population; this deviation is the standard error of the mean

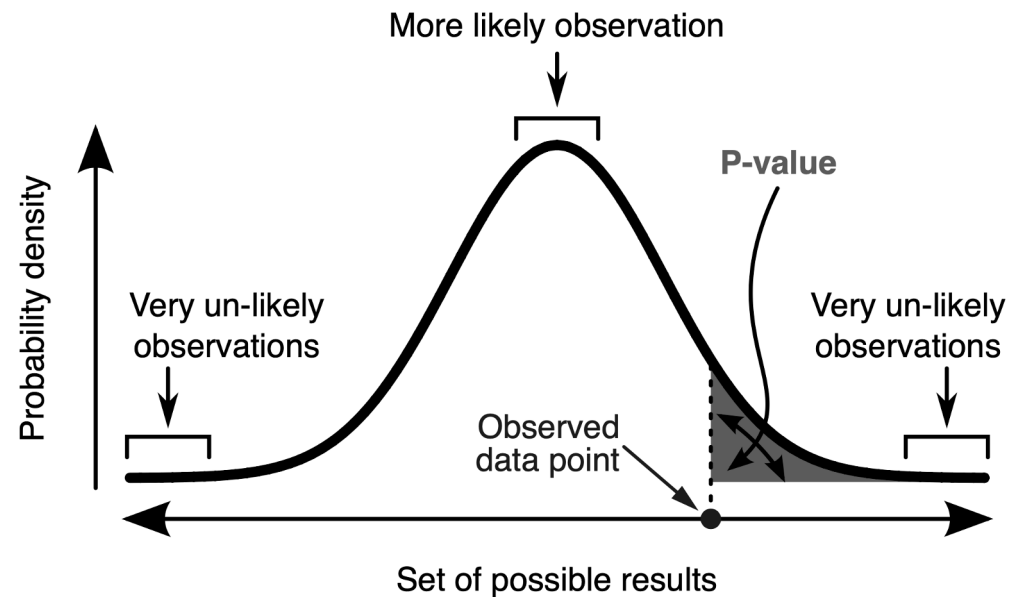


*A sample mean deviates from the actual mean of a population; this deviation is the standard error of the mean*

# **P-VALUE AND CONFIDENCE INTERVAL**

# Definition of P-value

- The P-value represents *the chance of seeing a statistic* (e.g., a mean difference or an odds ratio) *as extreme, by random chance alone, given the null hypothesis is true*
- By convention, if this chance is less than 1/20 ( $p < 0.05$ ), we throw the null hypothesis away (conclude treatment works)
  - i.e., reject the null hypothesis and accept the alternative hypothesis



# Confidence Intervals

- How precise is our estimate in terms of the population? To gauge this, we can use ***confidence intervals***
- A common misinterpretation of the 95% confidence interval (for example) is an interval that has a 95% chance of containing the true population value
- However, operationally, we can view the confidence interval as *an estimate of the range of plausible values for the population parameter*. It is this way that we can use confidence intervals for hypothesis testing



# Statistical Significance $\neq$ Clinical Significance

- Statistical significance: the probability or likelihood that there is a difference
- Clinical significance or clinical importance: decided by judgment whether the observed difference is important enough for you to act (e.g., to start treating patients with a treatment)

One dependent variable and at least one independent variable(s)

# **REGRESSION**

# Regression models

- Most used methods in Biostatistics
  - General linear models: linear regression and analysis of variance
  - Generalized linear models: logistic regression, Poisson regression
  - Survival analysis method: Cox proportional hazard regression
  - Methods for longitudinal (panel) data: linear mixed models, generalized estimating equations, generalized linear mixed model
- Choosing models depends on the dependent variable (types of measurement)

# Regression model

- A mathematical representation of the relation between an outcome variable and a set of explanatory variables.

$$y = f(x_1, x_2, x_3, x_4, \dots, x_n) + e$$

- *Cure from infection =  $f(\text{ABX used, type of organism, immune system}) + \text{residual}$*
- *Cholesterol level =  $f(\text{diet, body weight, genetic}) + \text{residual}$*

# Steps in regression analysis

- Fit the regression model, and obtain estimates of  $\beta_s$ ,  
(Using Statistical Software)
- Assess model adequacy (Question: is it the 'right' model?)
  - Model significance: Overall model significance, Xs significance
  - Explanatory power: How good is the model?
  - Model validity: Does it work properly? Are any assumptions violated?

## Steps in regression analysis (cont.)

- Prediction: Sometimes the model is 'good' enough to predict the response variable from values of the explanatory variable (rarely case in an 'observational' setting)

# Regression Usage

- RCT: baseline adjustments
- Non-Randomized Studies
  - Cohort study
  - Case-control study
  - Cross-sectional study

# Choice of Model

Type of Outcome Variable	Regression Model
Continuous	Linear regression
Binary	Logistic regression
Categorical	Ordinal, nominal logistic regression
Count	Poisson regression
Time-to-event	Cox-proportional hazard model

One outcome, measured once 32



# Linear Regression and ANOVA

	ANOVA	Linear Regression
<b>Outcome variable</b>	Continuous	Continuous
<b>Independent variables</b>	At least one <b>categorical</b> variables	One or more <b>continuous</b> variables
<b>Analysis</b>	Present differences in means of outcome between categories	Estimate regression coefficients
<b>Note</b>	Is a special case of regression	

Both ANOVA and linear regression can be used for one continuous outcome, two independent variables (one continuous and one categorical variables)

# Indicator (Dummy) variables

- An indicator variable or proxy variable takes value of 0 or 1 to indicate the absence or presence of some categorical effect
- Number of dummy variables = number of categories - 1
- Categorical variable: severity status (mild, moderate, severe), 2 dummy variables, mild is a reference group

	Dummy 1 (sev1)	Dummy 2 (sev2)
Mild	0	0
Moderate	1	0
Severe	0	1