

Basic Probability

Prof. dr. Siswanto Agus Wilopo, M.Sc., Sc.D.
Department of Biostatistics, Epidemiology and
Population Health
Faculty of Medicine
Universitas Gadjah Mada

Learning Objectives

In this lecture, you learn:

- **Basic probability concepts and definitions**
- **Conditional probability**
- **To use Bayes' Theorem to revise probabilities**
- **Various counting rules**

Important Terms

- **Probability** – the chance that an uncertain event will occur (always between 0 and 1)
- **Event** – Each possible outcome of a variable
- **Simple Event** – an event that can be described by a single characteristic
- **Sample Space** – the collection of all possible events

Assessing Probability

There are three approaches to assessing the probability of an uncertain event:

1. *a priori* classical probability

$$\text{probability of occurrence} = \frac{X}{T} = \frac{\text{number of ways the event can occur}}{\text{total number of elementary outcomes}}$$

2. empirical classical probability

$$\text{probability of occurrence} = \frac{\text{number of favorable outcomes observed}}{\text{total number of outcomes observed}}$$

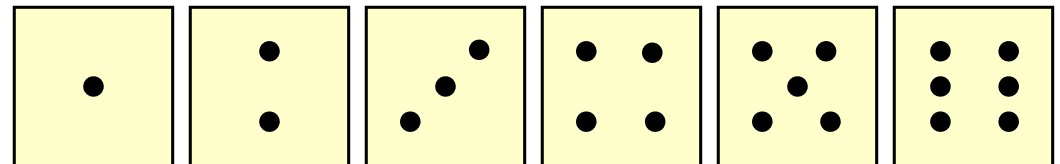
3. subjective probability

an individual judgment or opinion about the probability of occurrence

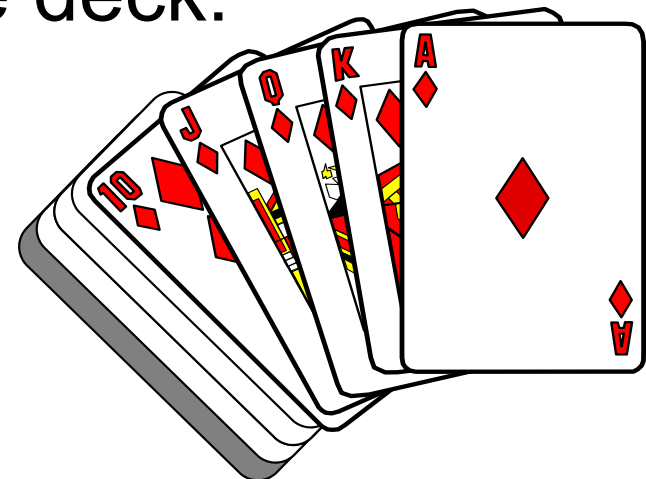
Sample Space

The **Sample Space** is the collection of all possible events

e.g. All 6 faces of a die:



e.g. All 52 cards of a bridge deck:



Events

■ Simple event

- An outcome from a sample space with one characteristic
- e.g., A red card from a deck of cards

■ Complement of an event A (denoted A')

- All outcomes that are not part of event A
- e.g., All cards that are not diamonds

■ Joint event

- Involves two or more characteristics simultaneously
- e.g., An ace that is also red from a deck of cards

Visualizing Events

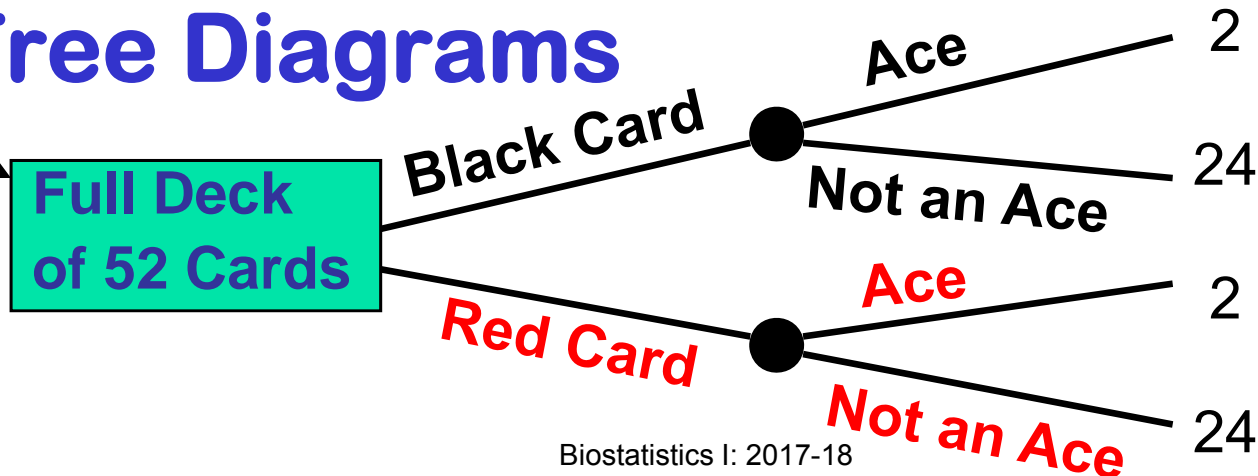
Contingency Tables

	Ace	Not Ace	Total
Black	2	24	26
Red	2	24	26
Total	4	48	52

Sample Space

Tree Diagrams

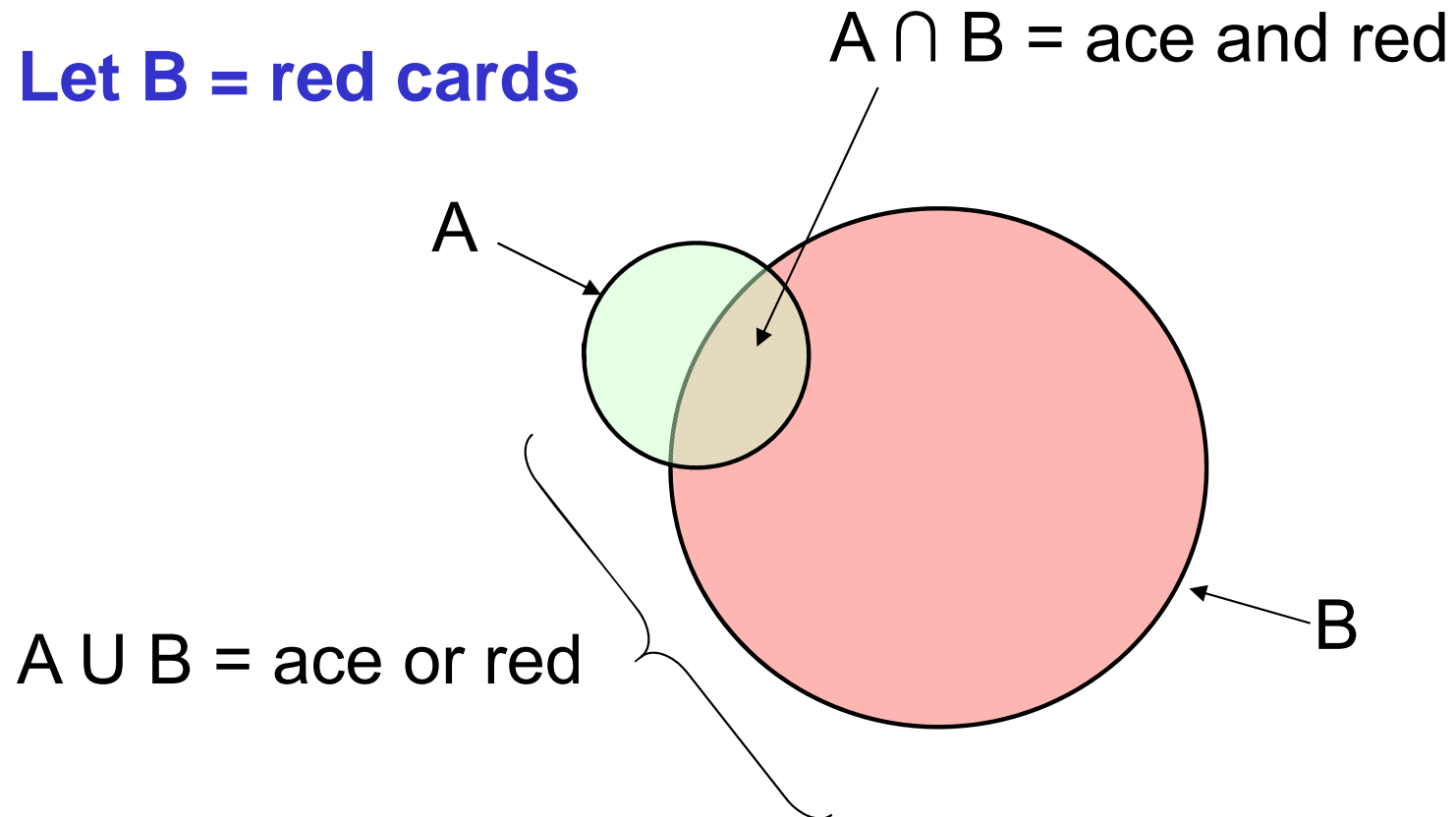
Sample Space



Visualizing Events

■ Venn Diagrams

- Let $A = \text{aces}$
- Let $B = \text{red cards}$



Mutually Exclusive Events

- **Mutually exclusive events**
 - Events that cannot occur together

example:

A = queen of diamonds; B = queen of clubs

- Events A and B are mutually exclusive

Collectively Exhaustive Events

- **Collectively exhaustive events**
 - One of the events must occur
 - The set of events covers the entire sample space

example:

**A = aces; B = black cards;
C = diamonds; D = hearts**

- Events A, B, C and D are collectively exhaustive (but not mutually exclusive – an ace may also be a heart)
- Events B, C and D are collectively exhaustive and also mutually exclusive

Probability

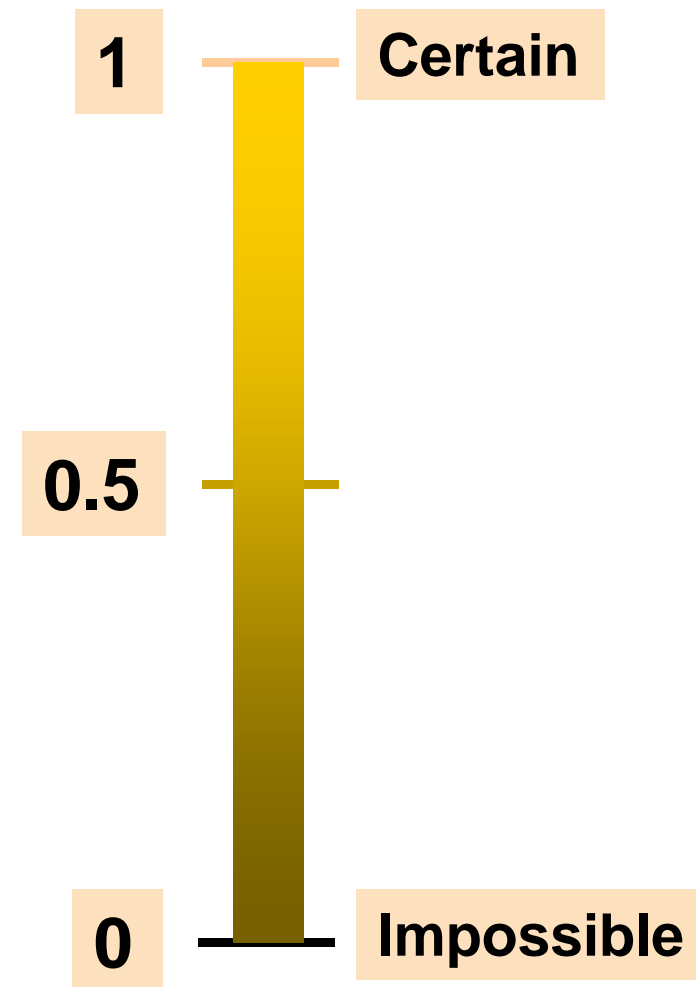
- Probability is the numerical measure of the likelihood that an event will occur
- The probability of any event must be between 0 and 1, inclusively

$$0 \leq P(A) \leq 1 \quad \text{For any event A}$$

- The sum of the probabilities of all mutually exclusive and collectively exhaustive events is 1

$$P(A) + P(B) + P(C) = 1$$

If A, B, and C are mutually exclusive and collectively exhaustive



Computing Joint and Marginal Probabilities

- The probability of a joint event, A and B:

$$P(A \text{ and } B) = \frac{\text{number of outcomes satisfying A and B}}{\text{total number of elementary outcomes}}$$

- Computing a marginal (or simple) probability:

$$P(A) = P(A \text{ and } B_1) + P(A \text{ and } B_2) + \cdots + P(A \text{ and } B_k)$$

- Where B_1, B_2, \dots, B_k are k mutually exclusive and collectively exhaustive events

Joint Probability Example

P(Red and Ace)

$$= \frac{\text{number of cards that are red and ace}}{\text{total number of cards}} = \frac{2}{52}$$

Type	Color		Total
	Red	Black	
Ace	2	2	4
Non-Ace	24	24	48
Total	26	26	52

Marginal Probability Example

P(Ace)

$$= P(\text{Ace and Red}) + P(\text{Ace and Black}) = \frac{2}{52} + \frac{2}{52} = \frac{4}{52}$$

Type	Color		Total
	Red	Black	
Ace	2	2	4
Non-Ace	24	24	48
Total	26	26	52

Joint Probabilities Using Contingency Table

Event	Event		Total
	B ₁	B ₂	
A ₁	P(A ₁ and B ₁)	P(A ₁ and B ₂)	P(A ₁)
A ₂	P(A ₂ and B ₁)	P(A ₂ and B ₂)	P(A ₂)
Total	P(B ₁)	P(B ₂)	1

Joint Probabilities

Marginal (Simple) Probabilities

General Addition Rule

General Addition Rule:

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

If A and B are mutually exclusive, then

$P(A \text{ and } B) = 0$, so the rule can be simplified:

$$P(A \text{ or } B) = P(A) + P(B)$$

For mutually exclusive events A and B

General Addition Rule Example

$$P(\text{Red or Ace}) = P(\text{Red}) + P(\text{Ace}) - P(\text{Red and Ace})$$

$$= \frac{26}{52} + \frac{4}{52} - \frac{2}{52} = \frac{28}{52}$$

Type	Color		Total
	Red	Black	
Ace	2	2	4
Non-Ace	24	24	48
Total	26	26	52

Don't count the two red aces twice!

Computing Conditional Probabilities

- A conditional probability is the probability of one event, given that another event has occurred:

$$P(A | B) = \frac{P(A \text{ and } B)}{P(B)}$$



The conditional probability of A given that B has occurred

$$P(B | A) = \frac{P(A \text{ and } B)}{P(A)}$$



The conditional probability of B given that A has occurred

Where $P(A \text{ and } B)$ = joint probability of A and B

$P(A)$ = marginal probability of A

$P(B)$ = marginal probability of B

Conditional Probability Example

- Of the cars on a used car lot, 70% have air conditioning (AC) and 40% have a CD player (CD). 20% of the cars have both.
- What is the probability that a car has a CD player, given that it has AC ?

i.e., we want to find $P(\text{CD} \mid \text{AC})$

Conditional Probability Example

(continued)

- Of the cars on a used car lot, **70%** have air conditioning (AC) and **40%** have a CD player (CD). **20%** of the cars have both.

	CD	No CD	Total
AC	0.2	0.5	0.7
No AC	0.2	0.1	0.3
Total	0.4	0.6	1.0

$$P(\text{CD} | \text{AC}) = \frac{P(\text{CD and AC})}{P(\text{AC})} = \frac{0.2}{0.7} = 0.2857$$

Conditional Probability Example

(continued)

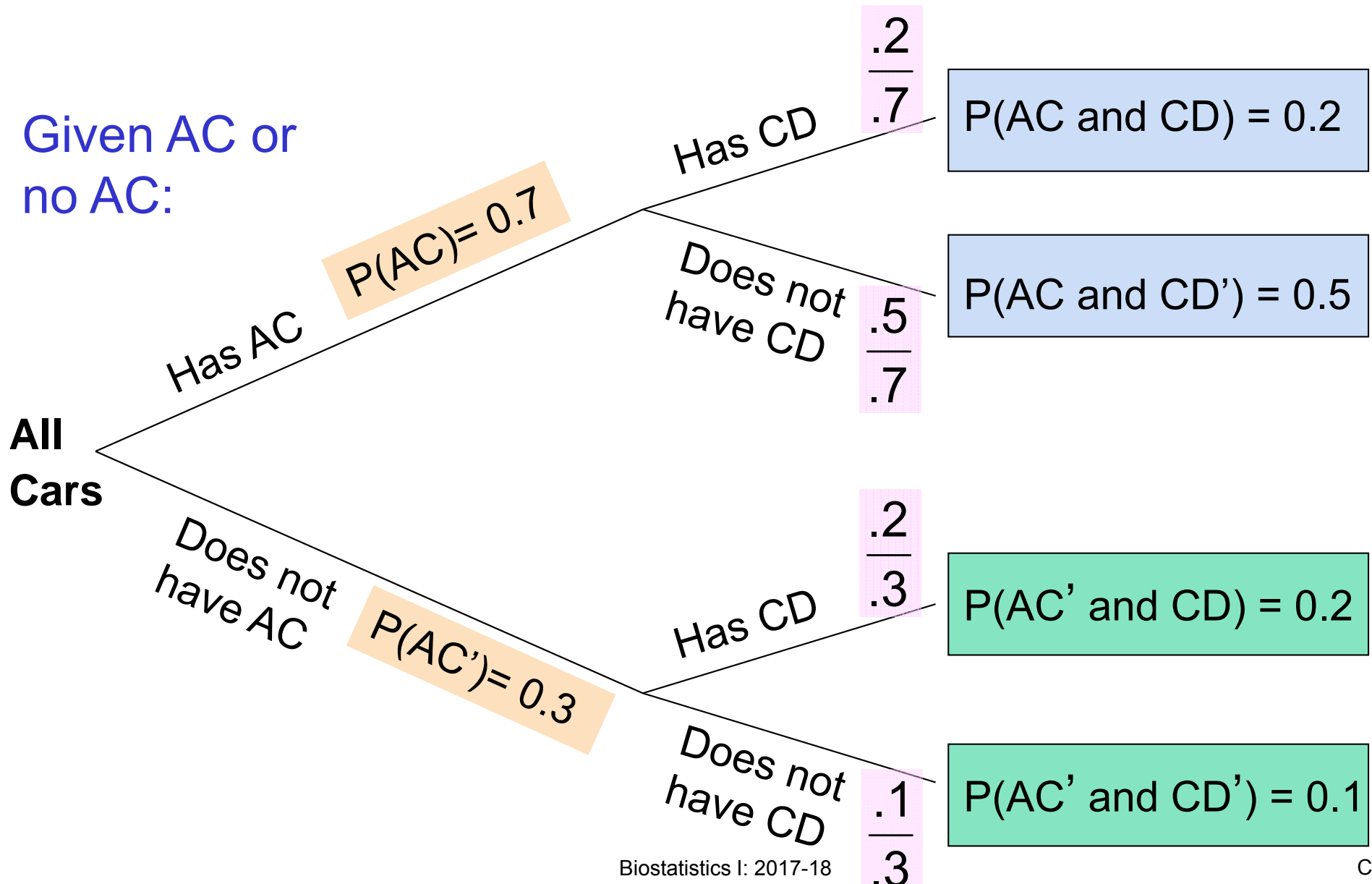
- Given AC, we only consider the top row (70% of the cars). Of these, 20% have a CD player. 20% of 70% is about 28.57%.

	CD	No CD	Total
AC	0.2	0.5	0.7
No AC	0.2	0.1	0.3
Total	0.4	0.6	1.0

$$P(\text{CD} | \text{AC}) = \frac{P(\text{CD and AC})}{P(\text{AC})} = \frac{0.2}{0.7} = 0.2857$$

Using Decision Trees

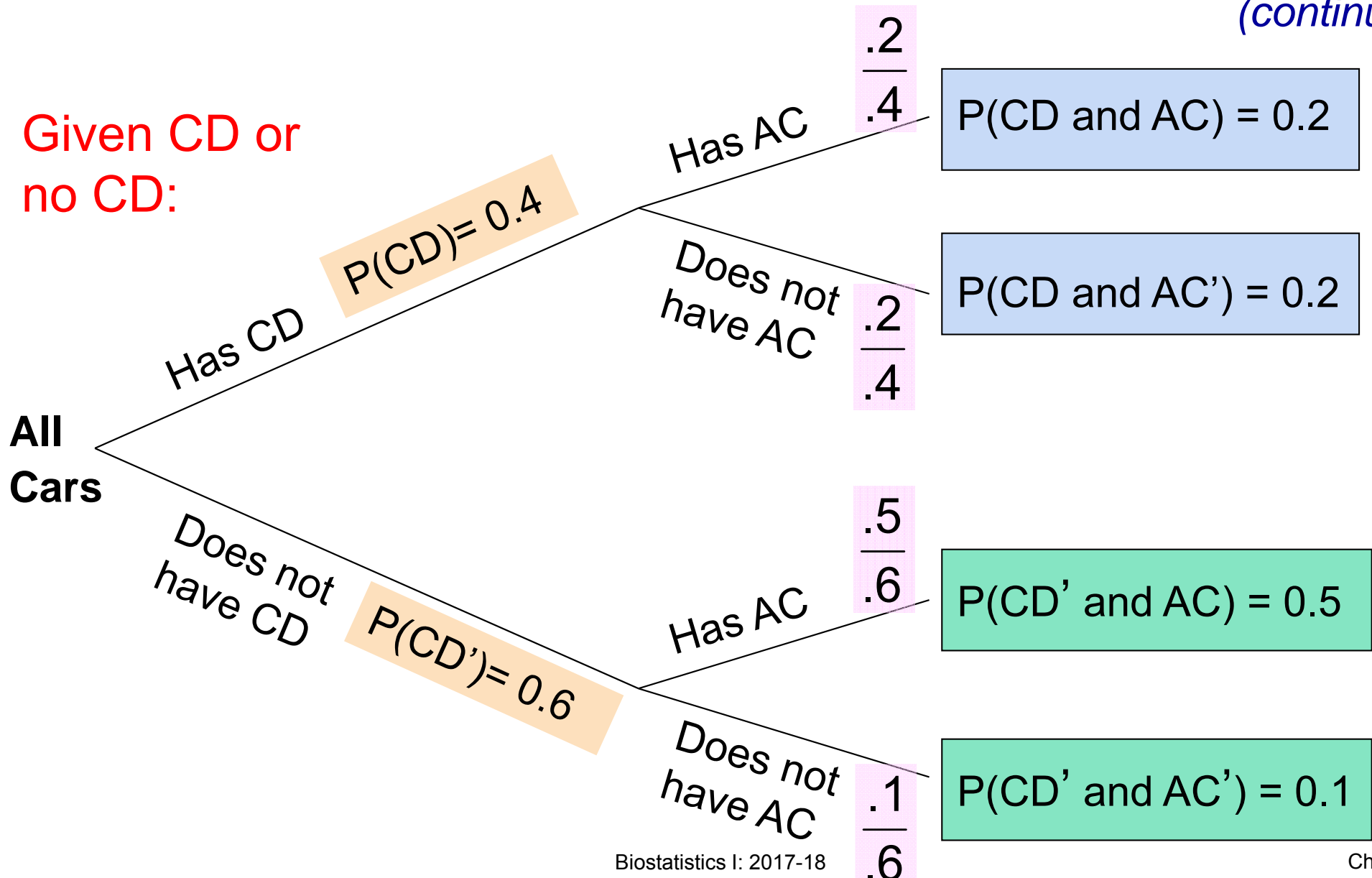
Given AC or
no AC:



Using Decision Trees

(continued)

Given CD or
no CD:



Statistical Independence

- Two events are independent if and only if:

$$P(A | B) = P(A)$$

- Events A and B are independent when the probability of one event is not affected by the other event

Multiplication Rules

- **Multiplication rule for two events A and B:**

$$P(A \text{ and } B) = P(A | B)P(B)$$

Note: If A and B are independent, then $P(A | B) = P(A)$ and the multiplication rule simplifies to

$$P(A \text{ and } B) = P(A)P(B)$$

Bayes' Theorem

Bayes' theorem is used to revise previously calculated probabilities after new information is obtained

$$P(B_i | A) = \frac{P(A|B_i)P(B_i)}{P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \dots + P(A|B_k)P(B_k)}$$

- **where:**

B_i = i^{th} event of k mutually exclusive and collectively exhaustive events

A = new event that might impact $P(B_i)$

Marginal Probability

- **Marginal probability for event A:**

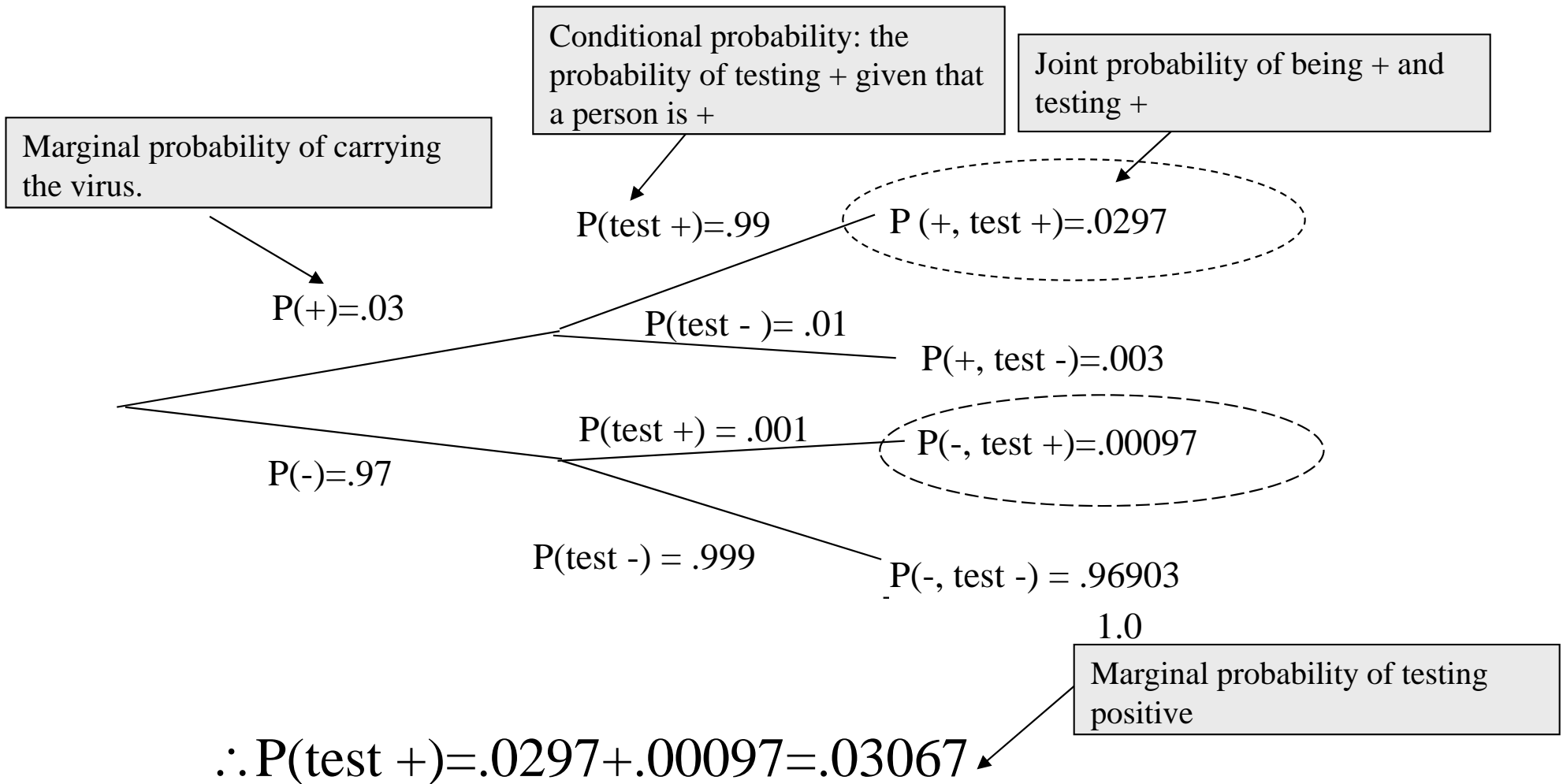
$$P(A) = P(A | B_1)P(B_1) + P(A | B_2)P(B_2) + \cdots + P(A | B_k)P(B_k)$$

- **Where B_1, B_2, \dots, B_k are k mutually exclusive and collectively exhaustive events**

Practice problem

If HIV has a prevalence of 3% in Jayapura, and a particular HIV test has a false positive rate of .001 and a false negative rate of .01, what is the probability that a random person selected off the street will test positive?

Answer



$P(+\&\text{test}+) \neq P(+)*P(\text{test}+)$
 $.0297 \neq .03*.03067 (= .00092)$
 \therefore Dependent!

Law of total probability

$$P(\text{test } +) = P(\text{test } + / \text{HIV}+)P(\text{HIV}+) + P(\text{test } + / \text{HIV}-)P(\text{HIV}-)$$

One of these has to be true (mutually exclusive, collectively exhaustive).
They sum to 1.0.

$$P(\text{test } +) = .99(.03) + .001(.97)$$

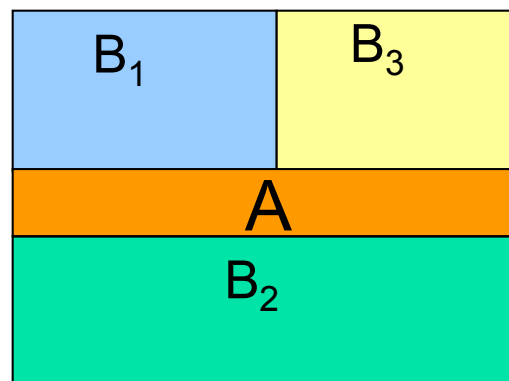
Law of total probability

- **Formal Rule: Marginal probability for event A=**

$$P(A) = P(A | B_1)P(B_1) + P(A | B_2)P(B_2) + \cdots + P(A | B_k)P(B_k)$$

$$\sum_{i=1}^k B_i = 1.0 \text{ and } P(B_i \& B_j) = 0 \text{ (mutually exclusive)}$$

- **Where:**

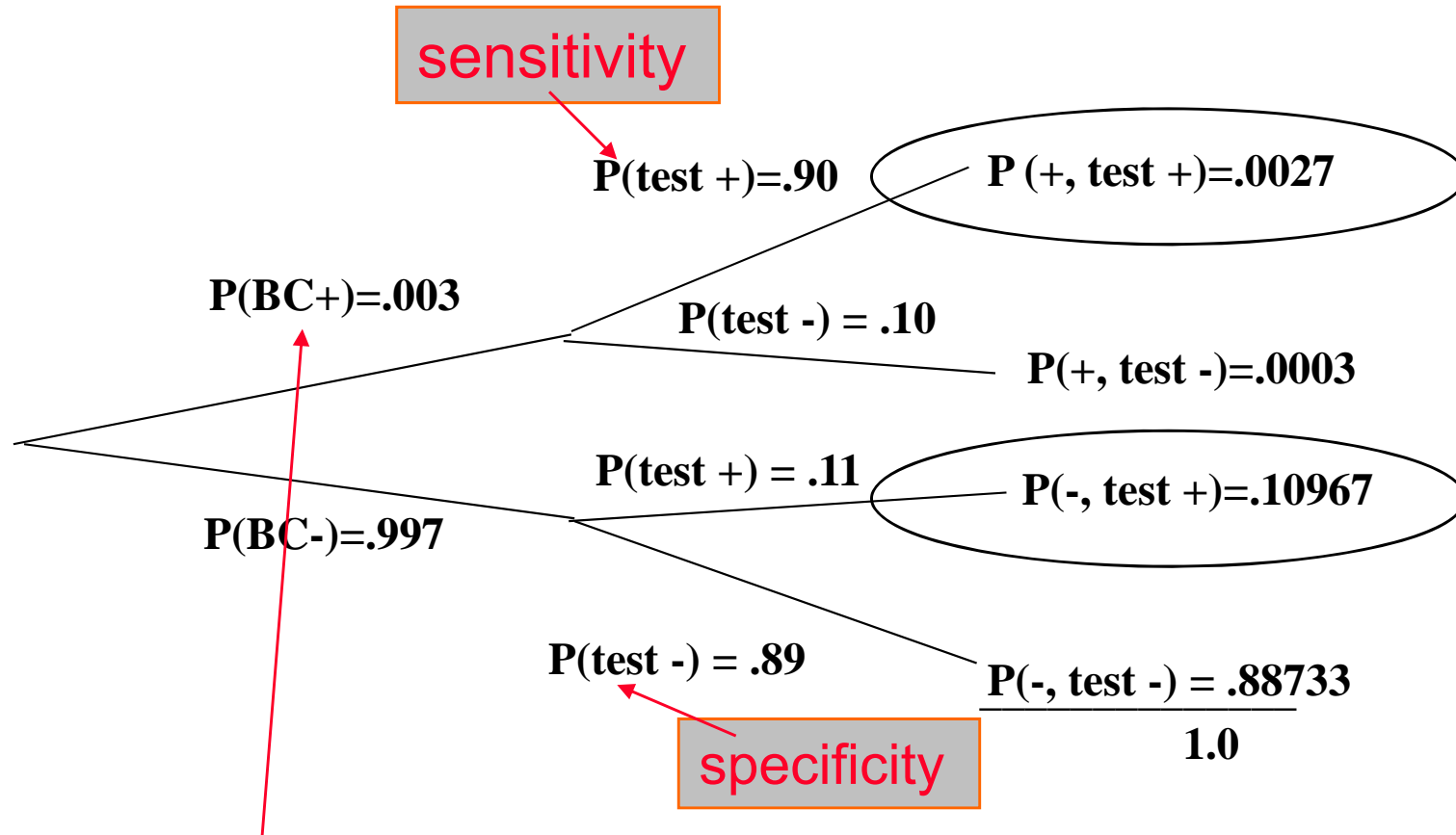


$$P(A) = (50\%)(25\%) + (0)(50\%) + \cdots + (50\%)(25\%) = 25\%$$

Example 2

- A 54-year old woman has an abnormal mammogram; what is the chance that she has breast cancer?
- Mammogram sensitivity=0.90 and specificity =0.89

Example: Mammography



Marginal probabilities of breast cancer....(prevalence among all 54-year olds)

$$P(BC/\text{test}+) = .0027 / (.0027 + .10967) = 2.4\%$$

BAYES' RULE

Bayes' Rule: derivation

■ Definition:

Let A and B be two events with $P(B) \neq 0$.
The conditional probability of A given B
is:

$$P(A / B) = \frac{P(A \& B)}{P(B)}$$

The idea: if we are given that the event B occurred, the relevant sample space is reduced to B { $P(B)=1$ because we know B is true} and conditional probability becomes a probability measure on B.

Bayes' Rule: derivation

$$P(A/B) = \frac{P(A \& B)}{P(B)}$$

can be re-arranged to:

$$P(A \& B) = P(A/B)P(B)$$

and, since also:

$$P(B/A) = \frac{P(A \& B)}{P(A)} \quad \therefore P(A \& B) = P(B/A)P(A)$$

$$P(A/B)P(B) = P(A \& B) = P(B/A)P(A)$$

$$P(A/B)P(B) = P(B/A)P(A)$$

$$\therefore P(A/B) = \frac{P(B/A)P(A)}{P(B)}$$

9/7/2017

Biostatistics I: 2017-18

36

Bayes' Rule:

$$P(A / B) = \frac{P(B / A)P(A)}{P(B)}$$

OR

$$P(A / B) = \frac{P(B / A)P(A)}{P(B / A)P(A) + P(B / \sim A)P(\sim A)}$$

From the “Law of Total Probability”

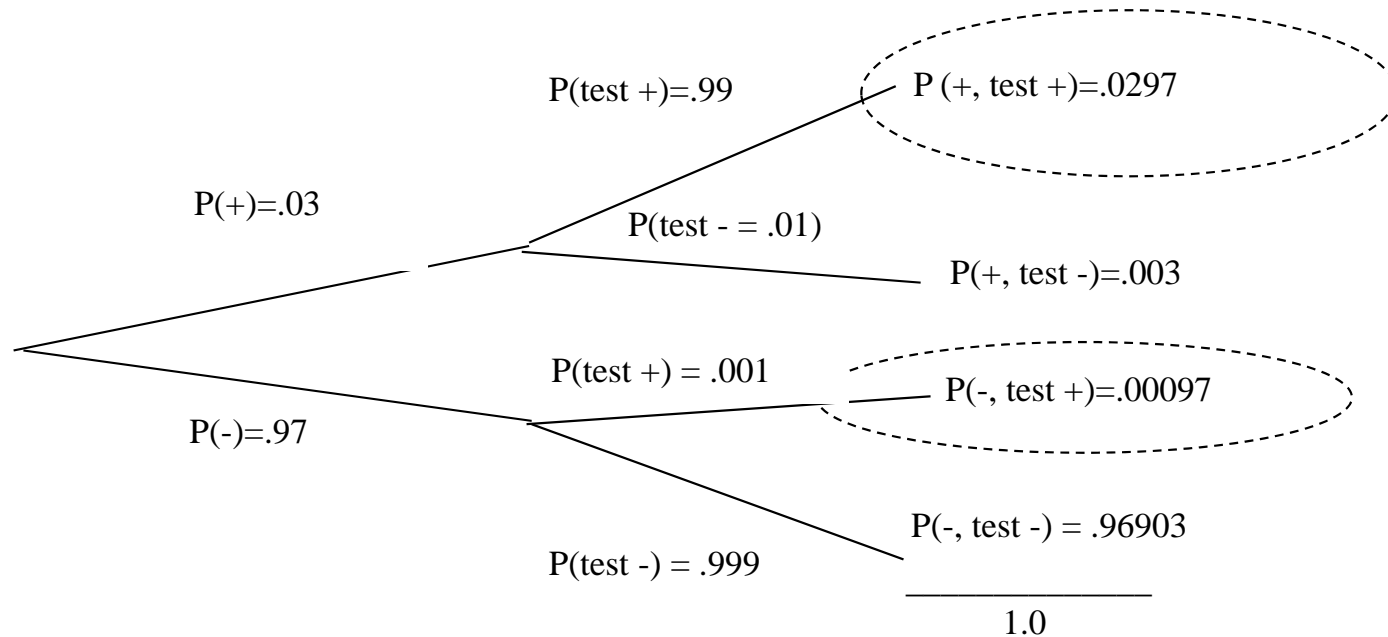
Bayes' Rule:

- **Why do we care??**
- **Why is Bayes' Rule useful??**
- **It turns out that sometimes it is very useful to be able to “flip” conditional probabilities. That is, we may know the probability of A given B, but the probability of B given A may not be obvious. An example will help...**

In-Class Exercise

- If HIV has a prevalence of 3% in Jayapura, and a particular HIV test has a false positive rate of .001 and a false negative rate of .01, what is the probability that a random person who tests positive is actually infected (also known as “positive predictive value”)?

Answer: using probability tree



A positive test places one on either of the two “test +” branches. But only the top branch also fulfills the event “true infection.” Therefore, the probability of being infected is the probability of being on the top branch given that you are on one of the two circled branches above.

$$P(+ / \text{test}+) = \frac{P(\text{test } + \& \text{true}+)}{P(\text{test}+)} = \frac{.0297}{.0297 + .00097} = 96.8\%$$

Answer: using Bayes' rule

$$P(\text{true} + / \text{test} +) = \frac{P(\text{test} + / \text{true} +)P(\text{true} +)}{P(\text{test} + / \text{true} +)P(\text{true} +) + P(\text{test} + / \text{true} -)P(\text{true} -)} =$$
$$\frac{.99(.03)}{.99(.03) + .001(.97)} = 96.8\%$$

In-class exercise

An insurance company believes that drivers can be divided into two classes—those that are of high risk and those that are of low risk. Their statistics show that a high-risk driver will have an accident at some time within a year with probability .4, but this probability is only .1 for low risk drivers.

- a) Assuming that 20% of the drivers are high-risk, what is the probability that a new policy holder will have an accident within a year of purchasing a policy?
- b) If a new policy holder has an accident within a year of purchasing a policy, what is the probability that he is a high-risk type driver?

Answer to (a)

Assuming that 20% of the drivers are of high-risk, what is the probability that a new policy holder will have an accident within a year of purchasing a policy?

Use law of total probability:

$P(\text{accident}) =$

$P(\text{accident/high risk}) * P(\text{high risk}) +$

$P(\text{accident/low risk}) * P(\text{low risk}) =$

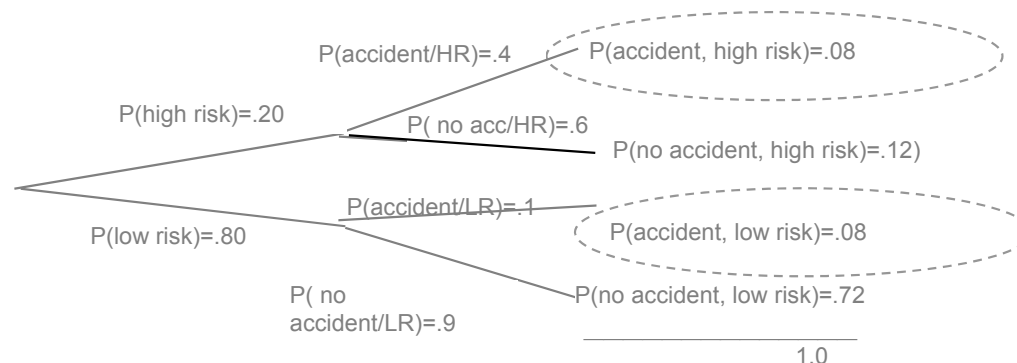
$.40(.20) + .10(.80) = .08 + .08 = .16$

Answer to (b)

If a new policy holder has an accident within a year of purchasing a policy, what is the probability that he is a high-risk type driver?

$$\begin{aligned} P(\text{high-risk/accident}) &= \\ P(\text{accident/high risk}) * P(\text{high risk}) / P(\text{accident}) \\ &= .40(.20) / .16 = 50\% \end{aligned}$$

Or use tree:



$$P(\text{high risk/accident}) = .08 / .16 = 50\%$$

Conditional Probability for Epidemiology:

The odds ratio and risk ratio as
conditional probability

The Risk Ratio and the Odds Ratio as conditional probability

In epidemiology, the association between a risk factor or protective factor (exposure) and a disease may be evaluated by the “risk ratio” (RR) or the “odds ratio” (OR).

Both are measures of “relative risk”—the general concept of comparing disease risks in exposed vs. unexposed individuals.

Odds and Risk (probability)

Definitions:

Risk = $P(A)$ = cumulative probability (you specify the time period!)

For example, what's the probability that a person with a high sugar intake develops diabetes in 1 year, 5 years, or over a lifetime?

Odds = $P(A)/P(\sim A)$

For example, “the odds are 3 to 1 against a horse” means that the horse has a 25% probability of winning.

Note: An odds is always higher than its corresponding probability, unless the probability is 100%.

Odds vs. Risk=probability

If the risk is...	Then the odds are...
$\frac{1}{2}$ (50%)	1:1
$\frac{3}{4}$ (75%)	3:1
$\frac{1}{10}$ (10%)	1:9
$\frac{1}{100}$ (1%)	1:99

Note: An odds is always higher than its corresponding probability, unless the probability is 100%.

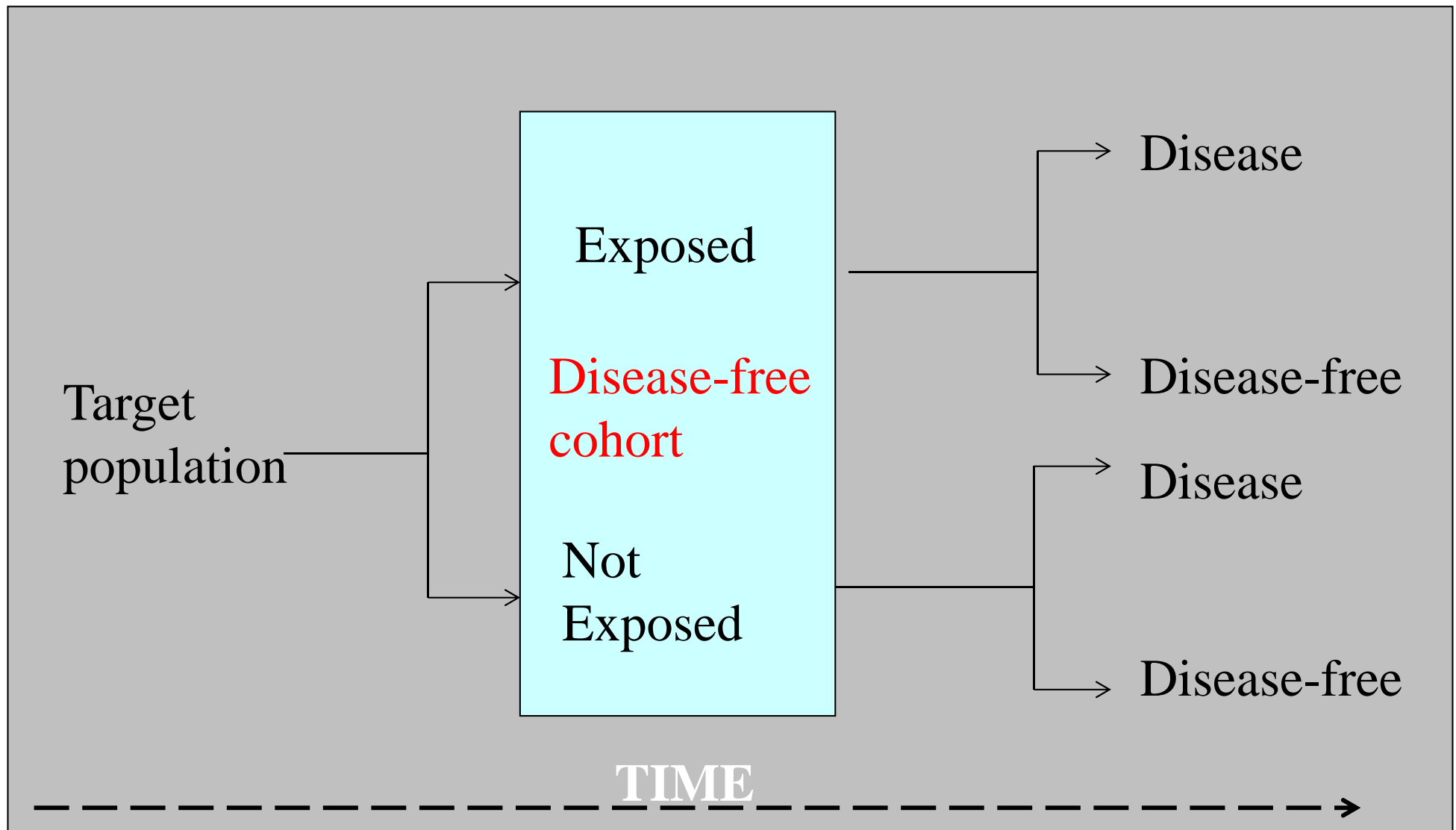
Introduction to the 2x2 Table

	Exposure (E)	No Exposure ($\sim E$)	
Disease (D)	a	b	$a+b = P(D)$
No Disease ($\sim D$)	c	d	$c+d = P(\sim D)$
	$a+c = P(E)$	$b+d = P(\sim E)$	

Marginal probability
of exposure

Marginal probability of
disease

Cohort Studies



The Risk Ratio, or Relative Risk (RR)

	Exposure (E)	No Exposure (~E)
Disease (D)	a	b
No Disease (~D)	c	d
	a+c	b+d

risk to the exposed

$$RR = \frac{P(D/E)}{P(D/\sim E)} = \frac{a/(a+c)}{b/(b+d)}$$

risk to the unexposed

Hypothetical Data

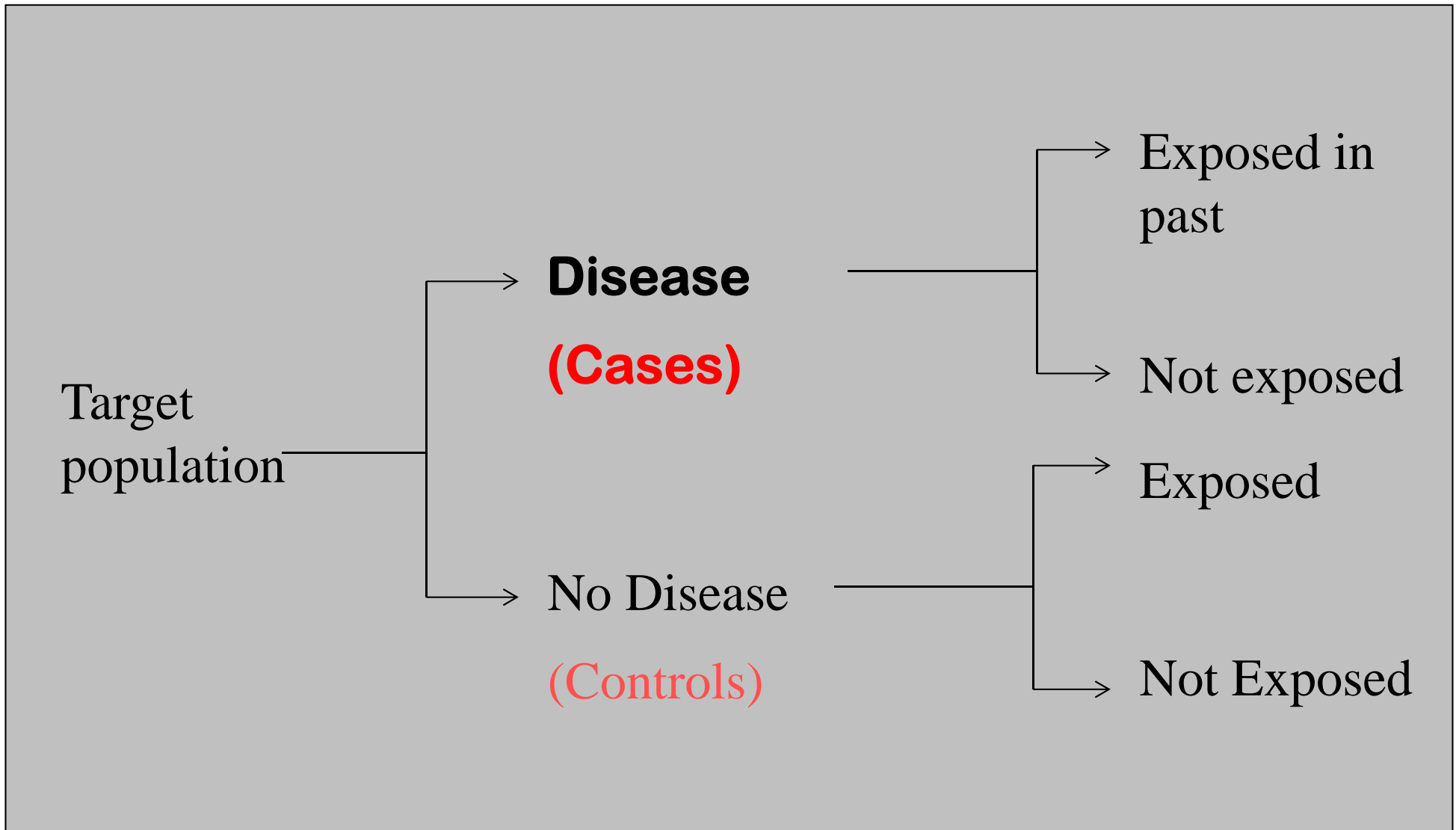
	High Systolic BP	Normal BP
Congestive Heart Failure	400	400
No CHF	1100	2600
	1500	3000

$$RR = \frac{400 / 1500}{400 / 3000} = 2.0$$

The odds ratio...

- This risk ratio seems like the perfect measure of relative risk. Why not stop here? Why introduce the more complicated odds ratio??
- We cannot calculate a risk ratio from a case-control study. Case-control studies are a popular study design in epidemiology, because they are useful for studying rare diseases.
- In a case-control study, we sample *conditional on disease status*, so we cannot calculate risk of disease.

Case-Control Studies



The Odds Ratio (OR)

	Hep C +	Hep C -	
Cases: Liver cancer	90	10	100
Controls	30	70	100

What are $P(D/E)$ and $P(D/\sim E)$ here?

We can't tell, because, by design, we have fixed the proportion of liver cancer cases in this sample at 50% simply by selecting half controls and half cases.

All these data give us is: $P(E/D)$ and $P(E/\sim D)$.

The Odds Ratio (OR)

	Hep C +	Hep C -	
Cases: Liver cancer	90	10	100
Controls	30	70	100

by Bayes' Rule...

Luckily, $P(E/D)$ [=the quantity you have] $\rightarrow P(E/D) = \frac{P(D/E)P(E)}{P(D)}$

Unfortunately, our sampling scheme precludes calculation of the marginals: $P(E)$ and $P(D)$, but turns out we don't need these if we use an odds ratio because the marginals cancel out!

$$\begin{array}{l}
 \left. \frac{P(E/D)}{P(\sim E/D)} \right\} \leftarrow \text{Odds of exposure in the cases} \\
 \frac{\frac{P(E/D)}{P(\sim E/D)}}{\frac{P(E/\sim D)}{P(\sim E/\sim D)}} \left\{ \leftarrow \text{Odds of exposure in the controls} \right. \\
 \\
 \text{Bayes' Rule} \rightarrow \frac{\frac{P(D/E)P(E)}{P(D)}}{\frac{P(D/\sim E)P(\sim E)}{P(D)}} \\
 \\
 \frac{\frac{P(\sim D/E)P(E)}{P(\sim D)}}{\frac{P(\sim D/\sim E)P(\sim E)}{P(\sim D)}} =
 \end{array}$$

$$\begin{array}{l}
 \left. \frac{P(D/E)}{P(\sim D/E)} \right\} \leftarrow \text{Odds of disease in the exposed} \\
 \frac{\frac{P(D/E)}{P(\sim D/E)}}{\frac{P(D/\sim E)}{P(\sim D/\sim E)}} \left\{ \leftarrow \text{Odds of disease in the unexposed} \right.
 \end{array}$$

What we want!

The Odds Ratio (OR)

	Exposure (E)	No Exposure (~E)
Disease (D)	a	b
No Disease (~D)	c	d

Odds of exposure for the cases.

Odds of disease for the exposed

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

Odds of exposure for the controls

Odds of disease for the unexposed

The rare disease assumption

$$OR = \frac{\frac{P(D/E)}{\cancel{P(\sim D/E)}_1}}{\frac{P(D/\sim E)}{\cancel{P(\sim D/\sim E)}_1}} \approx \frac{P(D/E)}{P(D/\sim E)} = RR$$

When a disease is rare:
 $P(\sim D) = 1 - P(D) \cong 1$

Example

	Hep C +	Hep C -
Cases: Liver cancer	90	10
Controls	30	70

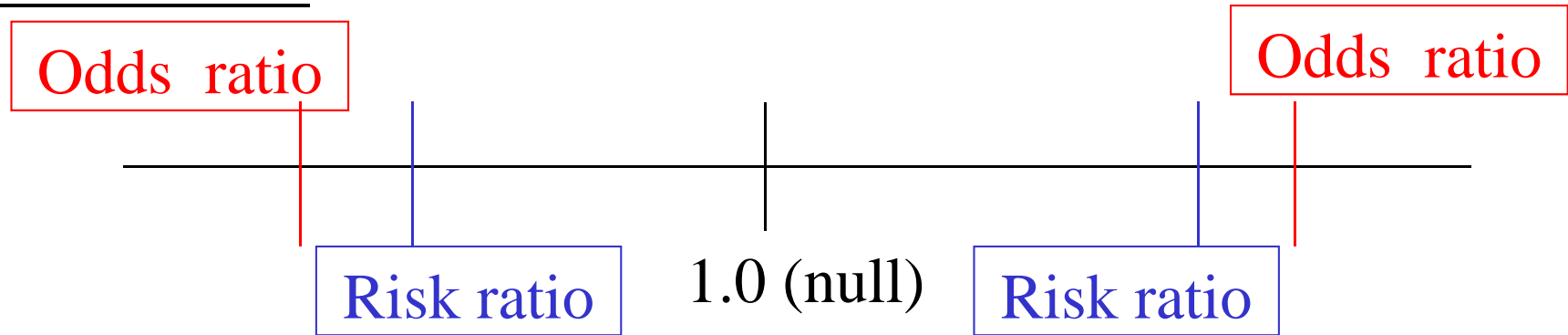
$$OR = 90 \cdot 70 / 10 \cdot 30 = 21.0$$

Note: This indicates that those with Hep C infection have a 21-fold increase in their *odds* of developing liver cancer (not in their risk!).

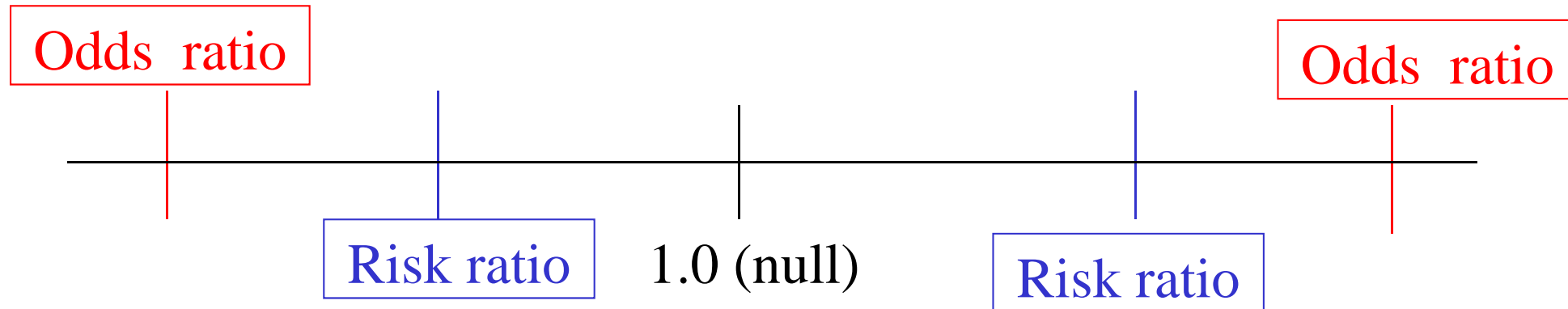
- The odds ratio will always be bigger than the corresponding risk ratio if $RR > 1$ and smaller if $RR < 1$ (the harmful or protective effect always appears larger)
- The magnitude of the inflation depends on the prevalence of the disease.

The odds ratio vs. the risk ratio

Rare Outcome



Common Outcome



In-Class Exercise:

1. Suppose the following data were collected on a random sample of subjects (the researchers did not sample on exposure or disease status).

	Neck pain	No Neck Pain
Own a cell phone	143	209
Don't own a cell phone	22	69

- Calculate the odds ratio and risk ratio for the association between cell phone usage and neck pain.

Answer

	Neck pain	No Neck Pain
Own a cell phone	143	209
Don't own a cell phone	22	69

- **$OR = (69 \cdot 143) / (22 \cdot 209) = 2.15$**
- **$RR = (143/352) / (22/91) = 1.68$**

In-Class Exercise:

- 2. Suppose the following data were collected on a random sample of subjects (the researchers did not sample on exposure or disease status).

	Brain tumor	No brain tumor
Own a cell phone	5	347
Don't own a cell phone	3	88

Calculate the odds ratio and risk ratio for the association between cell phone usage and brain tumor.

Answer

	Brain tumor	No brain tumor
Own a cell phone	5	347
Don't own a cell phone	3	88

- $OR = (5 \cdot 88) / (3 \cdot 347) = .42267$
- $RR = (5/352) / (3/91) = .43087$

Counting Rules

- Rules for counting the number of possible outcomes
- Counting Rule 1:
 - If any one of k different mutually exclusive and collectively exhaustive events can occur on each of n trials, the number of possible outcomes is equal to

$$k^n$$

Counting Rules

(continued)

■ Counting Rule 2:

- If there are k_1 events on the first trial, k_2 events on the second trial, ... and k_n events on the n^{th} trial, the number of possible outcomes is

$$(k_1)(k_2)\dots(k_n)$$

■ Example:

- You want to go to a park, eat at a restaurant, and see a movie. There are 3 parks, 4 restaurants, and 6 movie choices. How many different possible combinations are there?
- Answer: $(3)(4)(6) = 72$ different possibilities

Counting Rules

(continued)

■ Counting Rule 3:

- The number of ways that n items can be arranged in order is

$$n! = (n)(n - 1)\dots(1)$$

■ Example:

- Your restaurant has five menu choices for lunch. How many ways can you order them on your menu?
- Answer: $5! = (5)(4)(3)(2)(1) = 120$ different possibilities

Counting Rules

(continued)

■ Counting Rule 4:

- **Permutations:** The number of ways of arranging X objects selected from n objects in order is

$${}_n P_x = \frac{n!}{(n-X)!}$$

■ Example:

- Your restaurant has five menu choices, and three are selected for daily specials. How many different ways can the specials menu be ordered?

- Answer: ${}_n P_x = \frac{n!}{(n-X)!} = \frac{5!}{(5-3)!} = \frac{120}{2} = 60$ different possibilities

Counting Rules

(continued)

■ Counting Rule 5:

- **Combinations:** The number of ways of selecting X objects from n objects, irrespective of order, is

$${}_n C_x = \frac{n!}{X!(n-X)!}$$

■ Example:

- Your restaurant has five menu choices, and three are selected for daily specials. How many different special combinations are there, ignoring the order in which they are selected?

- Answer: ${}_n C_x = \frac{n!}{X!(n-X)!} = \frac{5!}{3!(5-3)!} = \frac{120}{(6)(2)} = 10$ different possibilities

Chapter Summary

- **Discussed basic probability concepts**
 - **Sample spaces and events, contingency tables, simple probability, and joint probability**
- **Examined basic probability rules**
 - **General addition rule, addition rule for mutually exclusive events, rule for collectively exhaustive events**

Chapter Summary

- **Defined conditional probability**
 - **Statistical independence, marginal probability, decision trees, and the multiplication rule**
- **Discussed Bayes' theorem**
- **Examined counting rules**