

Basics of Probability

STA 102: Introduction to Biostatistics

Yue Jiang

The following material was used by Yue Jiang during a live lecture.

Without the accompanying oral comments, the text is incomplete as a record of the presentation.

What's the use of probability?

- ▶ Last time: how descriptive statistics are used to *describe* data
- ▶ Goal: Make *inferences* about a population based on a sample

To do this, we need a solid foundation of probability theory.

Probabilities come up all the time

How do we interpret the following statements?

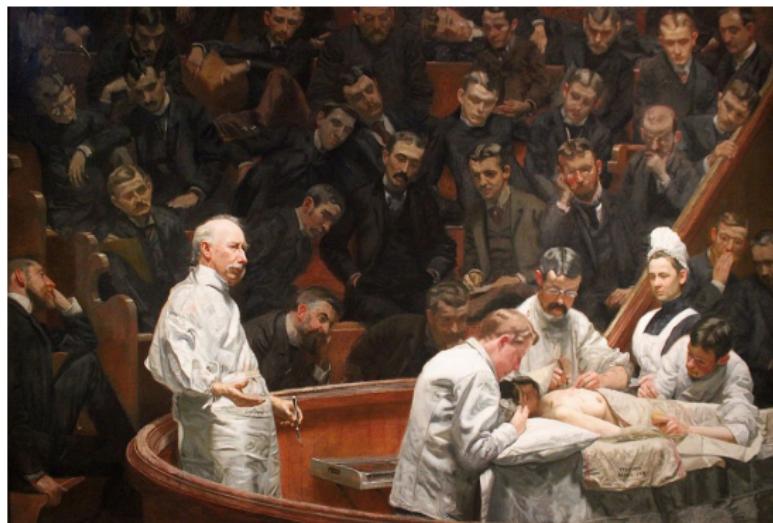
- ▶ There is a moderate chance of drought in North Carolina during the next year
- ▶ The surgery has a 50-50 probability of success
- ▶ The ten-year survival probability of invasive breast cancer among U.S women is 83%

Interpretations of probability



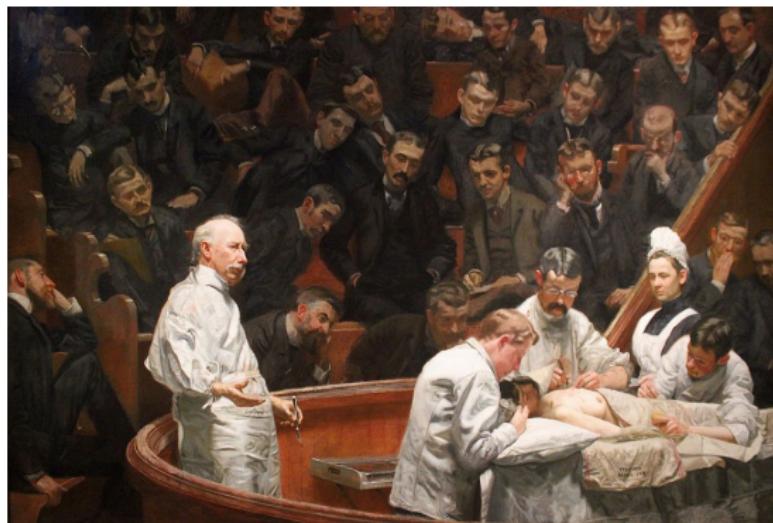
“There is a 1 in 3 chance of selecting a white ball”

Interpretations of probability



“The surgery has a 50% probability of success”

Interpretations of probability



Long-run frequencies vs. degree of belief

Probability spaces

Mathematical objects that model **random experiments**.

A probability space consists of three components:

1. A **sample space**, the set of all possible **outcomes**
2. Subsets of the sample space, called **events**, which comprise any number of possible outcomes (including none of them!)
3. A function that assigns **probabilities** to events

An event **occurs** if the outcome of the random experiment is contained in that event

Sample spaces

Sample spaces depend on the random experiment in question

- ▶ Tossing a single fair coin
- ▶ Tossing two fair coins
- ▶ Sum of rolling two fair six-sided dice
- ▶ Survival (years) after cancer diagnosis

Events

Subsets of the sample space that comprise possible outcomes.
Essentially, these are all the 'plausibly reasonable' events we're interested in calculating probabilities for*:

- ▶ Tossing a single fair coin
- ▶ Tossing two fair coins
- ▶ Sum of rolling two fair six-sided dice
- ▶ Survival (years) after cancer diagnosis*

**there are some nasty mathematical details behind this seemingly simple task. Don't worry about them!*

Probabilities

A number describing the likelihood of each event's occurrence.
This maps events to a number between 0 and 1, inclusive:

- ▶ Tossing a single fair coin **A head**
- ▶ Tossing two fair coins **At least one head**
- ▶ Sum of rolling two fair six-sided dice **An odd number**
- ▶ Survival (years) after cancer diagnosis **>one year**

Probabilities

A number describing the likelihood of each event's occurrence.
This maps events to a number between 0 and 1, inclusive:

- ▶ Tossing a single fair coin **A head** **0.5**
- ▶ Tossing two fair coins **At least one head** **0.75**
- ▶ Sum of rolling two fair six-sided dice **An odd number** **0.5**
- ▶ Survival (years) after cancer diagnosis **>one year** **...harder**

Events as (sub)sets

Let's take for now the example of tossing a single fair coin and recording the outcome.

There are only two elements in the outcome space:

- ▶ A : getting a head
- ▶ B : getting a tail

We can define the simple events of A occurring or B occurring, but are there “other” events we can define?

Set operations

For two sets (or events) A and B , the most common relationships are:

- ▶ **Intersection** ($A \cap B$): A and B both occur
- ▶ **Union** ($A \cup B$): A or B occur (including when both occur)
- ▶ **Complement** (A^c): A does not occur
- ▶ **Difference** ($A \setminus B$): A occurs, but B does not occur: $A \cap B^c$

Two sets A and B are said to be **disjoint** if $A \cap B = \emptyset$

How do probabilities “work”?

Kolmogorov axioms

1. The probability of any event in the sample space is a non-negative real number
2. The probability of the entire sample space is 1
3. If A and B are **disjoint** events (**mutually exclusive**), then the probability of A or B occurring is the sum of the individual probabilities that they occur



How do probabilities “work”?

For two events A and B with probabilities $P(A)$ and $P(B)$ of occurring, the Kolmogorov axioms give us two important rules:

- ▶ **Complement Rule:** $P(A^c) = 1 - P(A)$
- ▶ **Inclusion-Exclusion:** $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

How do we extend inclusion-exclusion to more than two events?

DeMorgan's laws

- ▶ Complement of union:
$$(A \cup B)^c = A^c \cap B^c$$
- ▶ Complement of intersection:
$$(A \cap B)^c = A^c \cup B^c$$

How do we extend DeMorgan's laws to more than two events?



Conditional probability

The probability an event will occur when another event has already occurred. The **conditional probability** of event A given event B is

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

Examples come up all the time in the real-world:

- ▶ Given that a mammogram comes back positive, what is the probability that a woman has breast cancer?
- ▶ Given that a 68-year old man has suffered four previous heart attacks, what is the probability he die in the next five years?
- ▶ Given that a patient has a mutation in the *CFTR* gene, what is the probability their offspring will have cystic fibrosis?

Gunter et al. (2017) study

ORIGINAL RESEARCH

Annals of Internal Medicine

Coffee Drinking and Mortality in 10 European Countries

A Multinational Cohort Study

Coffee drinking	Died		Total
	Yes	No	
None	1039	5438	6477
Med-Low	4440	29712	34152
High	3601	24934	28535
Total	9080	60084	69164

Define the events $A = \text{died}$ and $B = \text{non-coffee drinker}$

- ▶ Marginal probability $P(A)$
- ▶ Joint probability $P(A \cap B)$
- ▶ Conditional probability $P(A|B)$

More practice

Coffee drinking	Died		Total
	Yes	No	
None	1039	5438	6477
Med-Low	4440	29712	29809
High	3601	24934	28535
Total	9080	60084	64821

- ▶ ...did not drink coffee?
- ▶ ...died during the study or did not drink coffee?
- ▶ ...did not die during the study and was a high coffee drinker?
- ▶ ...died during the study given that they do not drink coffee?