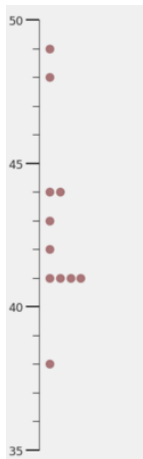


# CSCI 3022 Intro to Data Science

## Intro to Probability

Draw the box-whisker plot for the data to the left.



# Announcements and To-Dos

## Announcements:

1. HW 1 Posted, due Monday!
2. Another nb day this Friday!

## Last time we learned:

1. Drawing pictures out of our data: Histograms (which have *lots* of information!) versus Box-plots, which focus mainly on a few key measures.

## To do:

1. Start that HW! Ensure you can load the data and work with it. Practice your TeX/markdowns.

# Boxplots

A boxplot is a convenient way of graphically depicting groups of numerical data through the five number summary: minimum, first quartile, median, third quartile, and maximum.

1. The box extends from  $Q1$  to  $Q3$
2. The median line displays the median
3. The whiskers extend to farthest data point within  $1.5 \times IQR$  of each quartile
4. The fliers or outliers are any points outside of the whiskers
5. The width of the box is unimportant
6. Can be horizontally or vertically oriented

## Boxplots

Why do we use box plots?

1. They depict centrality via the median.
2. They depict dispersion through both the range and the IQR
3. Major outliers are shown
4. The median's location within the IQR suggests skewness; so too may lopsided whisker lengths or outliers

When might a box-whisker plot be misleading?

- 

When might a box-whisker plot be particularly useful?

-

## Boxplots

Why do we use box plots?

1. They depict centrality via the median.
2. They depict dispersion through both the range and the IQR
3. Major outliers are shown
4. The median's location within the IQR suggests skewness; so too may lopsided whisker lengths or outliers

When might a box-whisker plot be misleading?

- No indication of how data are dispersed (is there “no-man’s land”?)

When might a box-whisker plot be particularly useful?

- Comparing medium numbers of variables or columns quickly (say, 3-10); and much easier than histograms

## Overview: Probability

Many aspects of the world seem random and unpredictable.

1. Are we tall or short?
2. Do we have Mom's eyes or Dad's?
3. Is the hurricane going to hit Alabama?
4. Which team will win the NFC North?
5. How long until the Stampede bus shows up?
6. Which grocery store line should I get in?

One main objective of statistics/data science is to help make good decisions under conditions of uncertainty.

## Overview: Probability

Many aspects of the world seem random and unpredictable.

1. Are we tall or short?
2. Do we have Mom's eyes or Dad's?
3. Is the hurricane going to hit Alabama?
4. Which team will win the NFC North?
5. How long until the Stampede bus shows up?
6. Which grocery store line should I get in?



One main objective of statistics/data science is to help make good decisions under conditions of uncertainty.

## Overview: Definitions

**Definition:** *Set*

A *set* is a collection of objects.

**Definition:** *Probabilistic Process*

A *probabilistic process* is system/experiment whose outcomes are uncertain.

**Definition:** *Outcome*

An *outcome* is a possible result of a probabilistic process .

**Definition:** *Sample Space*

A *sample space* (denoted  $\Omega$ ) of a probabilistic process is the set of all possible outcomes of that process.



## Discrete vs. Continuous

Sets can contain many types of objects, both discrete and continuous. Our associated mathematics will shift accordingly.

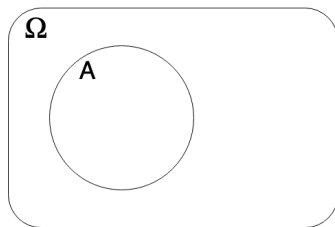
### *Discrete* (Structures)

1. Math: summation, counting, sorting
2. Sets: times, intervals

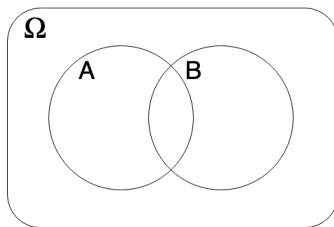
### *Continuous*

1. Math: integrals, derivatives, smooth functions
2. Sets: times, intervals

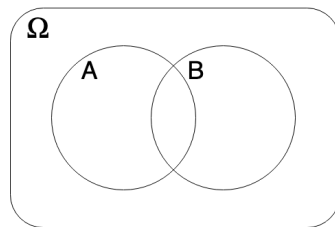
# Basic Set Operations



*Complement;*  
 $A^C$ ;  
"Not"

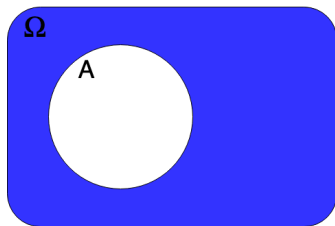


*Union;*  
 $A \cup B$ ;  
"Or"

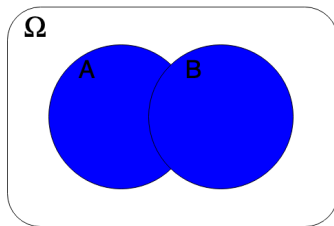


*Intersection;*  
 $A \cap B$ ;  
"And"

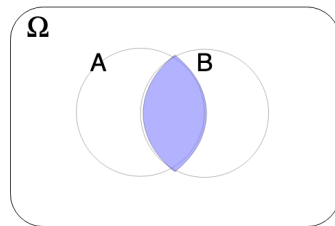
# Basic Set Operations



*Complement;*  
 $A^C$ ;  
"Not"



*Union;*  
 $A \cup B$ ;  
"Or"

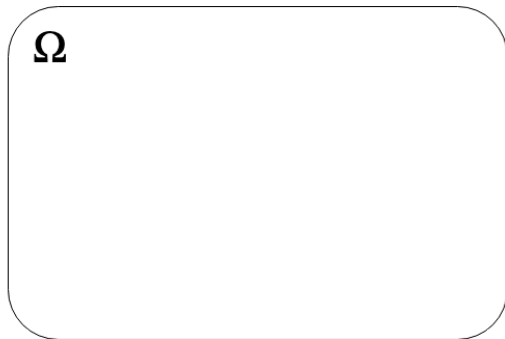


*Intersection;*  
 $A \cap B$ ;  
"And"

## Basic Set Definitions

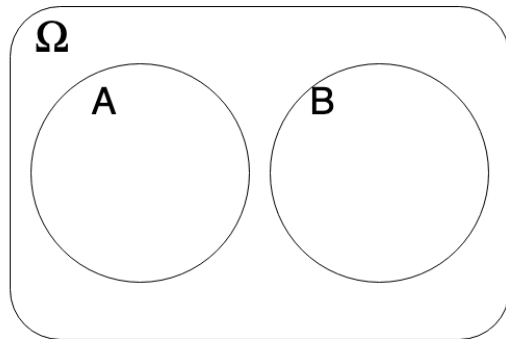


*Mutually Exclusive*

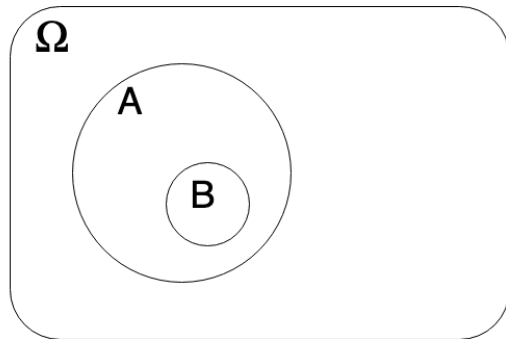


*Subset*

# Basic Set Definitions

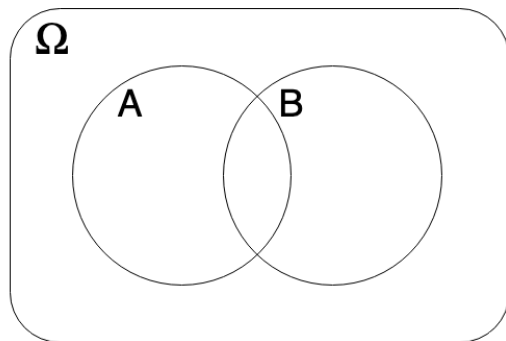


*Mutually Exclusive;*  
 $(A \cap B) = \emptyset$ ;  
 "If A, not B; If B, not A."

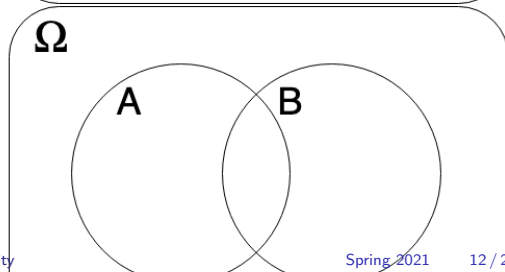
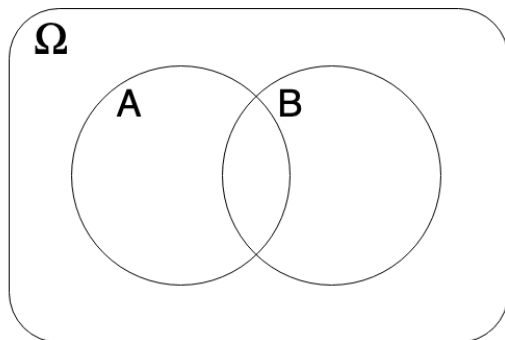


*Subset;*  
 $A \supseteq B$ ;  
 $A \supset B$ ;

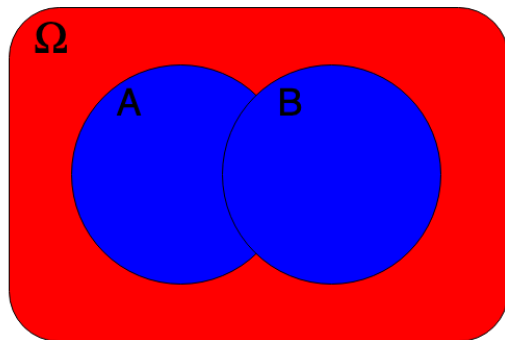
## DeMorgan's Laws



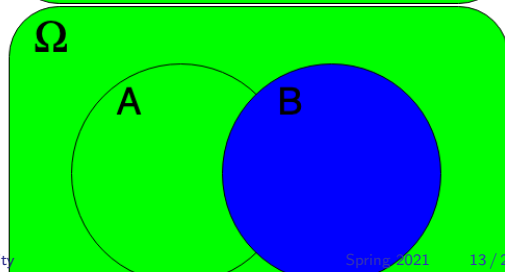
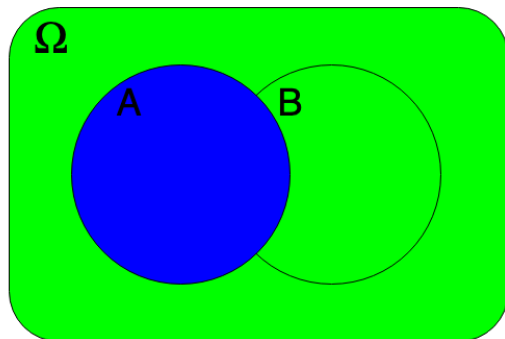
"neither A nor B"  
 compared to:  
 "not A"  
 "not B"



## DeMorgan's Laws



"neither A nor B;"  
 $(A \cup B)^C =$   
 $A^C \cap B^C;$   
 "not A AND not B"



## Some Sample Spaces

Describe sample spaces for:

1. Tossing a coin twice
2. Selecting a card from a deck
3. Measuring your commute time on a particular morning



## Some Sample Spaces

Describe sample spaces for:

1. Tossing a coin twice  
 $\{HH, HT, TH, TT\}$
2. Selecting a card from a deck  
 $\{2\clubsuit, 2\spadesuit, 2\diamondsuit, 2\heartsuit, 3\clubsuit, \dots\}$
3. Measuring your commute time on a particular morning  
 $\{t : t \in (0, T]\}$  where  $T$  is... infinity? The maximum reasonable time it *could* take?

# Event

**Definition:** *Event*

An *event* is any collection (subset) of outcomes from the sample space.

An event is simple if it consists of exactly one outcome and compound if it consists of more than one outcome.

When an experiment is performed, an event  $A$  is said to *occur* if the resulting experimental outcome is contained in  $A$ .

## Events

**Example:** Suppose that we flip a coin 3 times.

Sample space:

Some Possible Event(s):

$E_1$  : the event that we see the same flip all 3 times.

$E_2$  : the event that flip # 2 is heads.

What outcomes or elements(s) of  $\Omega$  are in  $E_1 \cap E_2$ ?

## Events

**Example:** Suppose that we flip a coin 3 times.

Sample space:

$\{HHH, HHT, HTH, THH, HTT, THT, TTH, TTT\}$

Some Possible Event(s):

$E_1$  : the event that we see the same flip all 3 times.

$E_2$  : the event that flip # 2 is heads.

What outcomes or elements(s) of  $\Omega$  are in  $E_1 \cap E_2$ ? : just  $\{HHH\}$

# Probability Axioms

**Definition:** *Probability*

*Probability* is a function that takes in sets (and later, we'll see, random variables) and outputs numbers according to the following rules:

1. *Non-negativity:*
2. *Unity:*
3.  *$\sigma$ -additivity:*

# Probability Axioms

**Definition:** *Probability*

*Probability* is a function that takes in sets (and later, we'll see, random variables) and outputs numbers according to the following rules:

1. *Non-negativity:* For every  $A \in \Omega$ ,  $P(A) \geq 0$ .
2. *Unity:* Given a sample space  $\Omega$ ,  $P(\Omega) = 1$ .
3.  *$\sigma$ -additivity:* If  $A$  and  $B$  are disjoint (mutually exclusive) sets,  $P(A \cup B) = P(A) + P(B)$ .

## Probability Theorems

The axioms of probability give us a couple of important results.

1. *Complementation:*
2. *Inclusion/Exclusion:* What is  $P(A \cup B)$ ?

## Probability Theorems

The axioms of probability give us a couple of important results.

1. *Complementation*:  $P(A^C) = 1 - P(A)$ .

Proof:

2. *Inclusion/Exclusion*: What is  $P(A \cup B)$ ?  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

Proof:



## Probability Theorems

The axioms of probability give us a couple of important results.

1. *Complementation*:  $P(A^C) = 1 - P(A)$ .

Proof: From unity,  $P(\Omega) = 1$ , and  $\Omega = A \cup A^C$ , which are disjoint sets. So  $P(\Omega) = P(A \cup A^C) = P(A) + P(A^C) = 1$ .

2. *Inclusion/Exclusion*: What is  $P(A \cup B)$ ?  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

Proof:  $A \cup B$  is "A or B," which can happen 3 disjoint ways:

0.1 "A not B;" or  $A \cap B^C$ , with probability  $P(A) - P(A \cap B)$ ;

0.2 "B not A;" or  $B \cap A^C$ , with probability  $P(B) - P(A \cap B)$ ;

0.3 "both;" with probability  $P(A \cap B)$ ;

Summing these 3 probabilities gives the desired result.

## Probability Theorems

The axioms of probability give us a couple of important results.

1. *Complementation:*
2. *Inclusion/Exclusion:* What is  $P(A \cup B)$ ?  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .  
Proof:

This idea works for more than 2 sets!

## Probabilities on Random Variables

Let  $X = \#$  of heads in three tosses of a fair coin.  $X$  is a *random variable*: it maps events (a count of heads) into real numbers through probabilities.

What is the underlying probabilistic process?

What is the sample space?

What are the possible values for  $X$ ?

What is the probability that  $X$  is equal to 1:  $P(X = 1)$ ?

## Probabilities on Random Variables

Let  $X = \#$  of heads in three tosses of a fair coin.  $X$  is a *random variable*: it maps events (a count of heads) into real numbers through probabilities.

What is the underlying probabilistic process?

Flipping a fair coin.

What is the sample space?

The same 8 flip-outcomes as before.

What are the possible values for  $X$ ?

$X \in \{0, 1, 2, 3\}$

What is the probability that  $X$  is equal to 1:  $P(X = 1)$ ?

If all outcomes are equally likely in a set, we can arrive at this by counting the elements of  $\Omega$  in  $X$  compared to all of  $\Omega$ : or  $\frac{|X|}{|\Omega|} = 3/8$

## Probabilities on Random Variables

Why stop at fair coins? What if our coin is *unfair*, and comes up heads  $p$  proportion of the time, so  $P(\{H, T\}) = \{p, q\}$ ?

What is the probability that I flip a biased coin twice and both flips come up heads?

Sample space for one flip:

Sample space for both flips (a product of sample spaces!):

Should the probability of the second flip change based on the result of the first?

## Probabilities on Random Variables

Why stop at fair coins? What if our coin is *unfair*, and comes up heads  $p$  proportion of the time, so  $P(\{H, T\}) = \{p, q\}$ ? Note:  $q = 1 - p$ .

What is the probability that I flip a biased coin twice and both flips come up heads?

Sample space for one flip:  $\{H, T\}$

Sample space for both flips (a product of sample spaces!):  $\{HH, HT, TH, TT\}$

Should the probability of the second flip change based on the result of the first?

Not usually: we call these *independent*... Not everything is independent!

## Probabilities on Random Variables

Our coin is *unfair*, and comes up heads  $p$  proportion of the time. What is the probability that I flip a biased coin twice and both flips come up heads?

Sample space for one flip:

Sample space for both flips (a product of sample spaces!):

Should the probability of the second flip change based on the result of the first?

## Probabilities on Random Variables

Our coin is *unfair*, and comes up heads  $p$  proportion of the time. What is the probability that I flip a biased coin twice and both flips come up heads?

Sample space for one flip:  $\{H, T\}$

Sample space for both flips (a product of sample spaces!):  $\{HH, HT, TH, TT\}$

Should the probability of the second flip change based on the result of the first?

Not usually: we call these *independent*... Not everything is independent! **Idea:** two or more trials are *independent* if they don't affect each other



## Independence and Probabilities

Our coin is *unfair*, and comes up heads  $p$  proportion of the time. What is the probability that I flip a biased coin twice and both flips come up heads?

If two outcomes are independent, probabilities on their intersection ("and") becomes a product.

Result: What are  $P(\{HH\})$  and  $P(\{TT\})$ ?

If two outcomes are disjoint, probabilities on their union ("or") becomes a sum.

Result: What is  $P(\{HT\} \text{ OR } \{HT\})$ ?

Sanity check: did we just add up to 1?

## Independence and Probabilities

Our coin is *unfair*, and comes up heads  $p$  proportion of the time. What is the probability that I flip a biased coin twice and both flips come up heads?

If two outcomes are independent, probabilities on their intersection ("and") becomes a product.

Result: What are  $P(\{HH\})$  and  $P(\{TT\})$ ?

A:  $p \cdot p$  and  $q \cdot q = (1 - p)^2$

If two outcomes are disjoint, probabilities on their union ("or") becomes a sum.

Result: What is  $P(\{HT\} \text{ OR } \{TH\})$ ?

A:  $P(\{HT\}) = pq$  PLUS  $P(\{TH\}) = qp$

Sanity check: did we just add up to 1?

## Counting outcomes

Finally, what is the probability of I flip our biased coin five times and get *exactly* one heads?

## Counting outcomes

Finally, what is the probability of I flip our biased coin five times and get *exactly* one heads?

This is the set of events  $\{HTTTT, THTTT, TTHTT, TTTHT, TTTTH\}$ .

## Counting outcomes

Finally, what is the probability of I flip our biased coin five times and get *exactly* one heads?

This is the set of events  $\{HTTTT, THTTT, TTHTT, TTTHT, TTTTH\}$ .

Each is composed of 5 independent flips, so the probability of any one of these events is the product  $pq^4$

## Counting outcomes

Finally, what is the probability of I flip our biased coin five times and get *exactly* one heads?

This is the set of events  $\{HTTTT, THTTT, TTHTT, TTTHT, TTTTH\}$ .

Each is composed of 5 independent flips, so the probability of any one of these events is the product  $pq^4$

Each outcome is disjoint/exclusive, so the full cumulative probability is the sum of 5 of these:  
 $5pq^4$

## Moving Forward

Suppose we have a coin and we don't know if it's biased... what could we do? (nb04, lecture next week to come!)

# Daily Recap

Today we learned

1. A review of probability
2. Think about when we can use “all outcomes equally likely” and then just *count* those outcomes. This is a big part of *independence*.

Moving forward:

- Friday: making some histograms, boxplots, and playing around with data frames: scrubbing data!

Next time in lecture:

- We probably talk even more about probability!