

UNIVERSITY *of* WASHINGTON

Data Science UW

Methods for Data

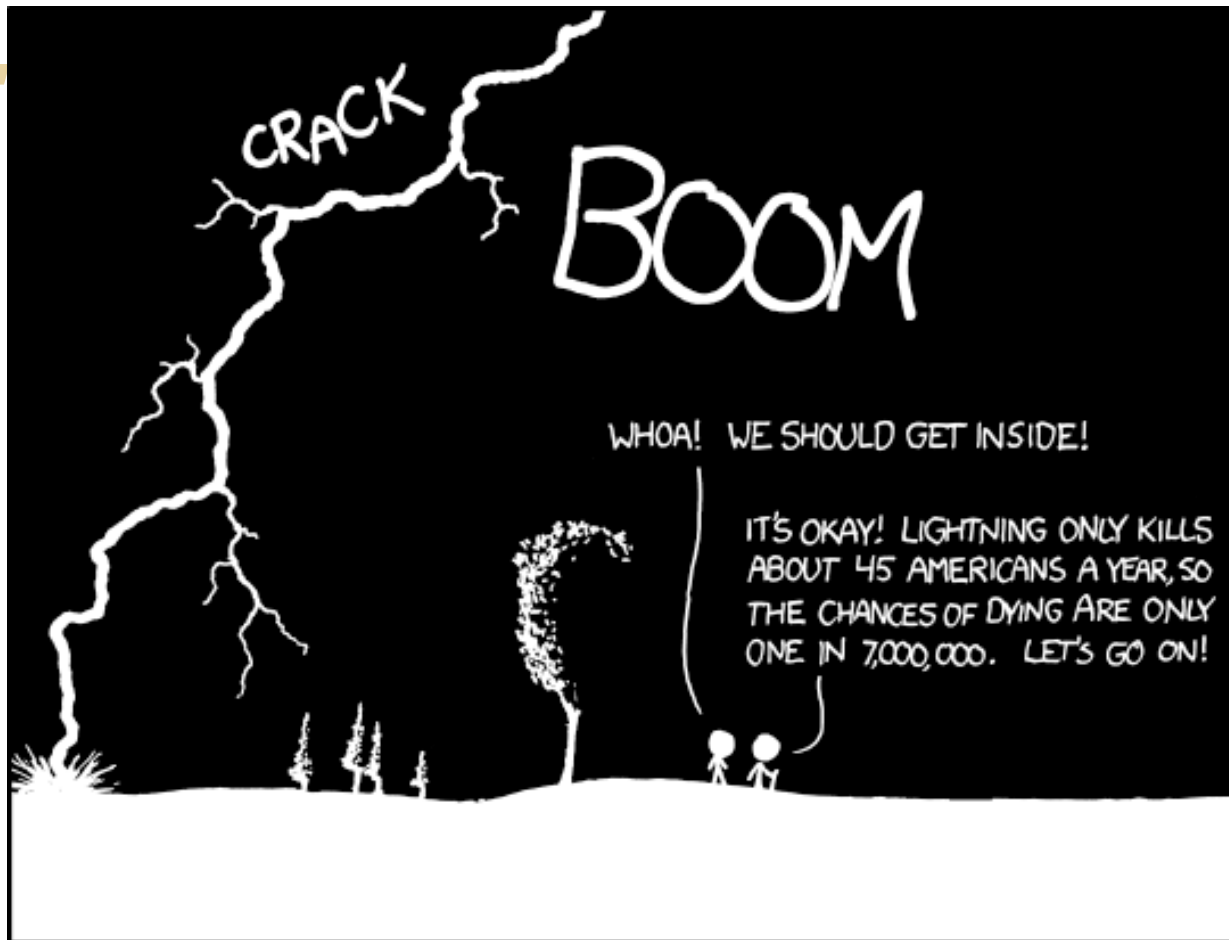
Analysis

Probability and More on Distributions

Lecture 2

Nick McClure





THE ANNUAL DEATH RATE AMONG PEOPLE
WHO KNOW THAT STATISTIC IS ONE IN SIX.

W

Topics

- > Review
- > Counting
- > Axioms of Probability
- > Probability Examples
- > Conditional Probability
- > More on Distributions



Review

- > Distributions
 - Discrete: Bernoulli, Binomial, Poisson
 - Continuous: Uniform, Normal, Student's T
- > Numerical and Visual Exploration of Data
- > Transformations
- > Simpson's Paradox
- > R Code Examples
 - R Review
 - Numerical/Visual Exploration of Distributions
 - Weather_retrieval.R as a production level script



Counting

> This is one of the biggest areas of mathematics, called Combinatorics.

> Example:

- Subway has 4 different breads, 5 different meats, 4 different toppings. How many sandwich combinations?
- How many different 4-beer tasters can I have in a bar with 10 beers on tap?

> Solve these using the ‘Multiplication Principle’.

- Subway Problem:

$$\begin{array}{ccccccc} 4 & * & 5 & * & 4 & & = 80 \\ \hline (\# \text{ of breads}) & & (\# \text{ of meats}) & & (\# \text{ of toppings}) & & \end{array}$$

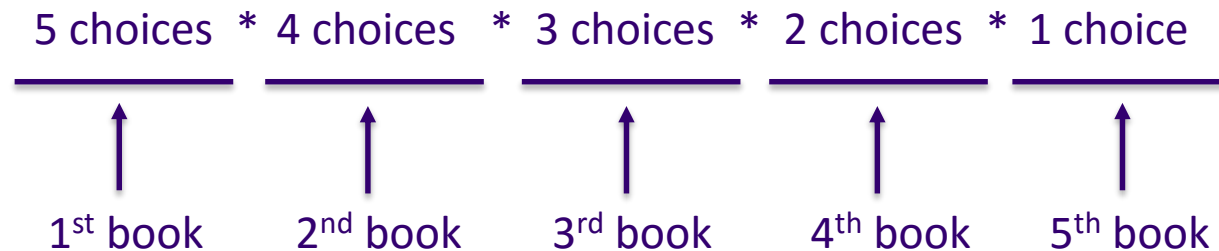
- Beer Problem:

$$\begin{array}{ccccccc} 10 & * & 9 & * & 8 & * & 7 & = 5,040 \\ \hline (\# \text{ for 1}^{\text{st}} \text{ beer}) & & (\# \text{ for 2}^{\text{nd}} \text{ beer}) & & (\# \text{ for 3}^{\text{rd}} \text{ beer}) & & (\# \text{ for 4}^{\text{th}} \text{ beer}) & \end{array}$$



Multiplication Principle

- > If there are A ways of doing task a, and B ways of doing task b, then there are $A*B$ ways of completing both tasks.
- > Example:
 - If I have 5 books, how many ways can I *order* them on the bookshelf?



$$= 5 \text{ factorial} = 5! = 120$$



Factorials

> Factorials

- Count # ways to order N things = $N!$

> Factorials get VERY large quickly.

- 21! Is larger than the biggest long-int in 64 bit.
 - > $21! = 5.1E19$
 - > Biggest long int (64 bit) = $9.2E18$
- Fun fact, every 52 card shuffle is highly likely to be the only time that shuffle has ever occurred.



Counting Subgroups

- > Revisit: 10 beers on tap, need a sample of 4 different beers.
- > Let's assume order matters, i.e., Amber-Stout-Porter-Red is different from Red-Porter-Stout-Amber.
- > Use 'Permutations' (pick):

$$10 * 9 * 8 * 7 = \frac{10!}{6!} = \frac{10!}{(10 - 4)!} = 10P4 = P(10,4)$$



Counting Subgroups

- > Now, Let's assume order doesn't matter.
- > Use 'Combinations' (choose):

$$10 * 9 * 8 * 7 = \frac{10!}{6!} = \frac{10!}{(10 - 4)!} = 10P4 = P(10,4)$$

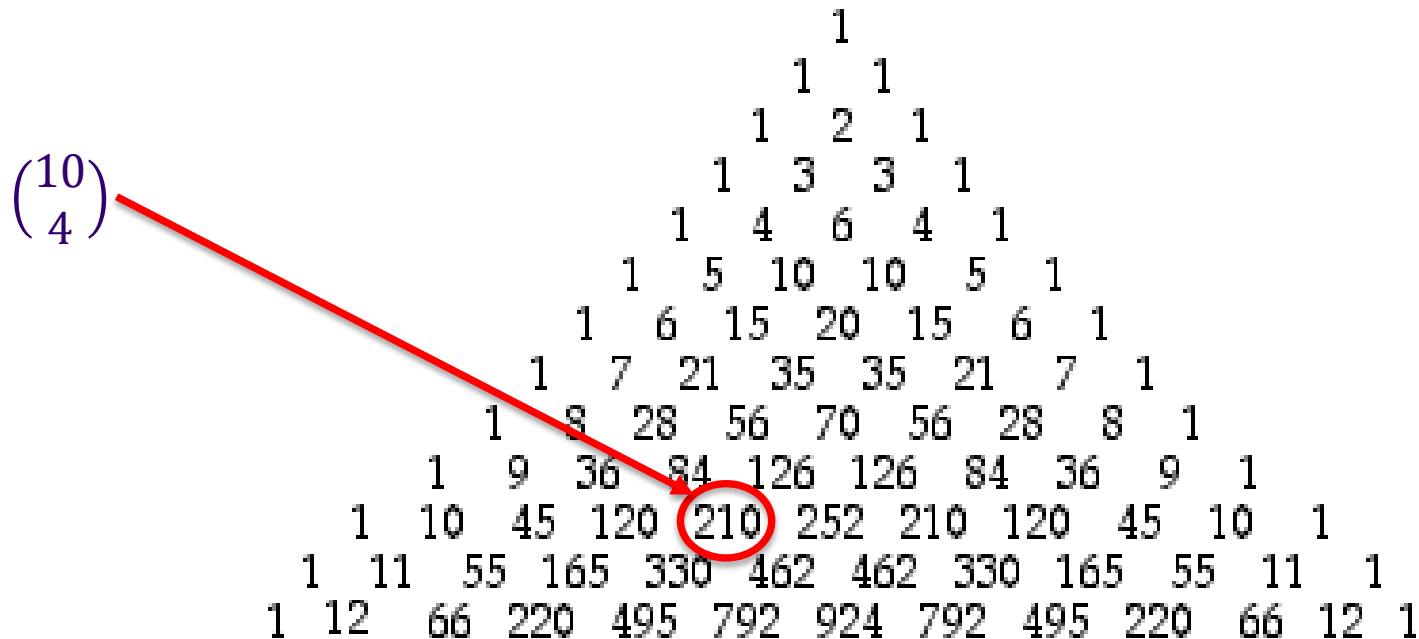
(# of orderings of 4 beers) = 4!

$$= \frac{10!}{4! (10 - 4)!} = 10C4 = C(10,4) = \binom{10}{4}$$



More on Combinations

- > Combinations appear on the Pascal's Triangle!
- > $C(N,x)$ appears on the Nth row, xth number (starting at 0)



Counting Examples

- > There are 10 Light beers on tap, and 10 Dark beers on tap, how many ways can Rick get a 4-beer sampler that contains exactly 1 light beer? (ordering doesn't matter)

$$\frac{(\# \text{ of ways for light beer}) \cdot (\# \text{ of ways for dark beer})}{(\# \text{ of ways to order 1L and 3D})}$$

$$\frac{(10) \cdot \binom{10}{3}}{4} = \frac{10 * 120}{4} = 300$$



Counting Examples

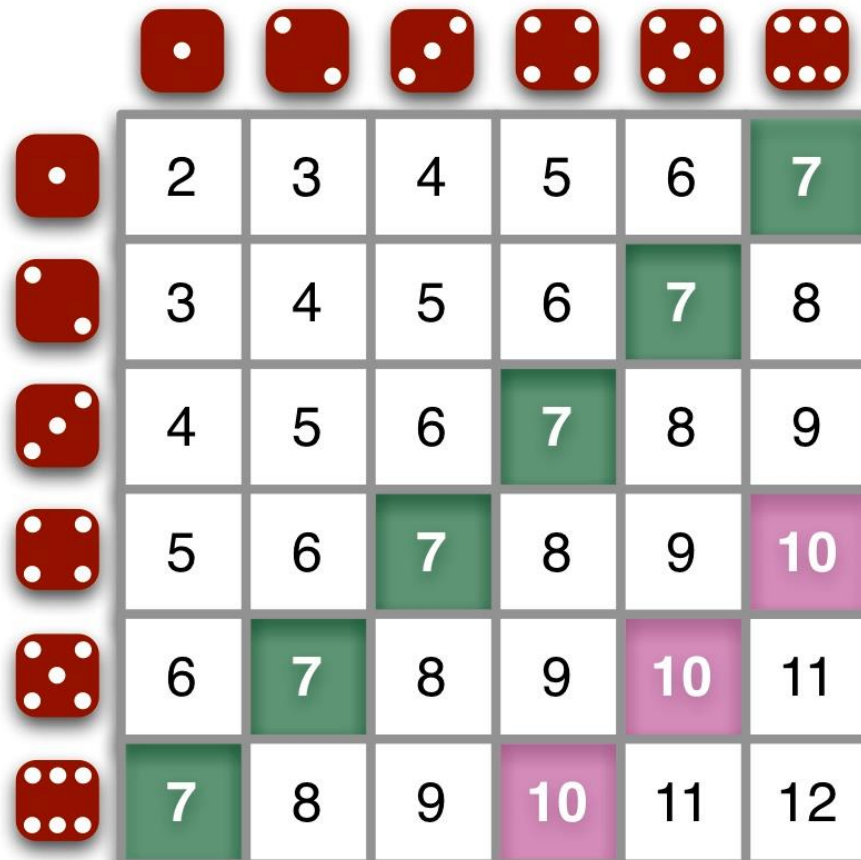
- > 6:5 Blackjack is dealt with a 6 shoe deck ($52 \times 6 = 312$ cards). How many ways can someone get dealt two rank 10 cards?

$$\binom{6decks * 4ranks * 4suits}{2} = \binom{96}{2} = \frac{96!}{2! (94!)} = \frac{96 * 95}{2} = 4560$$



Counting Examples

> How many ways can two dice be rolled to get a sum of 10?



	1	2	3	4	5	6
1	2	3	4	5	6	7
2	3	4	5	6	7	8
3	4	5	6	7	8	9
4	5	6	7	8	9	10
5	6	7	8	9	10	11
6	7	8	9	10	11	12

W

Counting in R

- > `expand.grid()` – function that creates a data frame from all combinations of vectors supplied.
- > R-demo



Probability

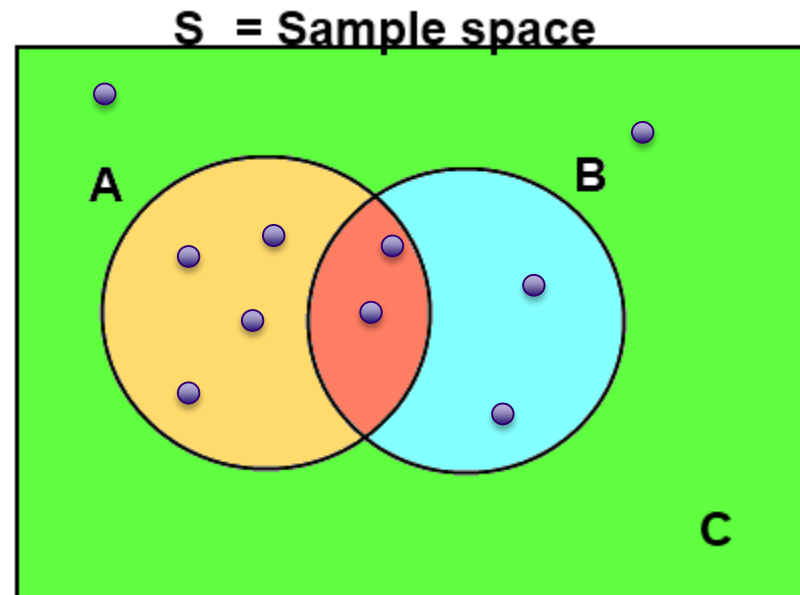
- > The Probability of an event, A, is the number of ways A can occur, divided by the number of total possible outcomes in our Sample Space, S.

$$P(A) = \frac{N(A)}{N(S)}$$

- > If \bullet is an event, then

$$P(A) = \frac{6}{10} = \frac{3}{5}$$

$$P(B) = \frac{4}{10} = \frac{2}{5}$$



W

Probability

> If \bullet is an event, then

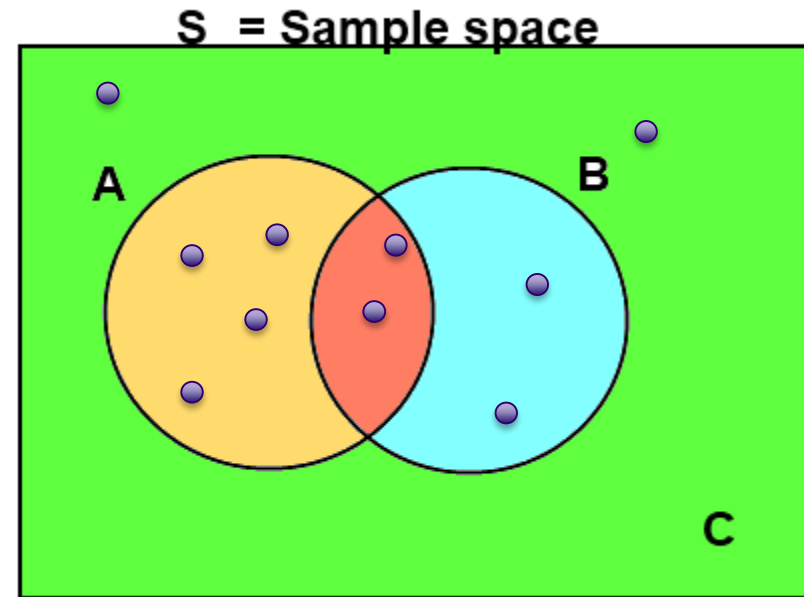
– Intersection: $P(A \cap B) = \frac{2}{10} = \frac{1}{5}$

– Union: $P(A \cup B) = \frac{8}{10} = \frac{4}{5}$

– Negation: $P(A') = \frac{6}{10} = \frac{3}{5}$

$$P((A \cup B)') = P(C) = \frac{2}{10} = \frac{1}{5}$$

$$P(A' \cap B') = P(C) = \frac{2}{10} = \frac{1}{5}$$



Axioms of Probability

- > Probability is bounded between 0 and 1.

$$0 \leq P(A) \leq 1$$

Note: “Percent” literally means per one hundred

- > Probability of the Sample Space = 1.

$$P(S) = 1$$

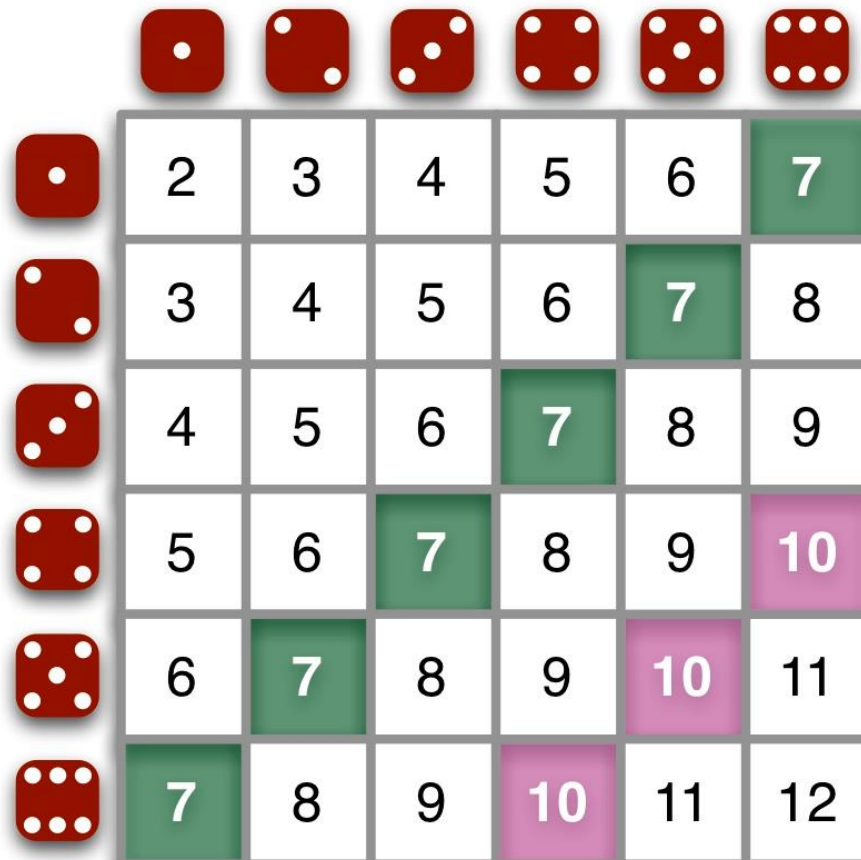
- > The probability of finite *mutually exclusive* unions is the sum of their probabilities.

$$P(A \cup B) = P(A) + P(B) \quad \text{If A and B are M.E.}$$



Probability Examples

> Probability of rolling a sum of 10?



A 6x6 grid of dice rolls. The columns are labeled with dice showing 1 to 6 dots at the top. The rows are labeled with dice showing 1 to 6 dots on the left. The grid contains the following values:

	1	2	3	4	5	6
1	2	3	4	5	6	7
2	3	4	5	6	7	8
3	4	5	6	7	8	9
4	5	6	7	8	9	10
5	6	7	8	9	10	11
6	7	8	9	10	11	12

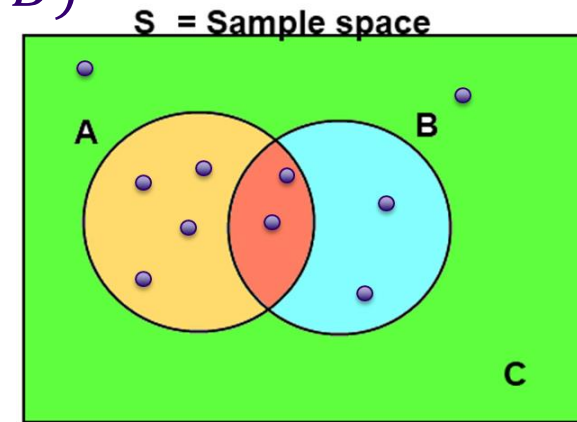
The value 10 is highlighted in pink in the cells (4,6), (5,5), and (6,4). The values 7 are highlighted in green in the cells (1,6), (2,5), (3,4), and (4,3).

W

Mutually Exclusive Events

- > In all cases, the probability of the union of A and B takes the form:

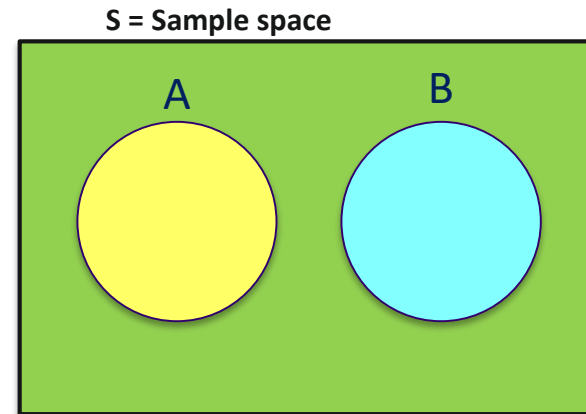
$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



- > If A and B are mutually exclusive that means that

$$P(A \cap B) = 0$$

$$P(A \cup B) = P(A) + P(B)$$



W

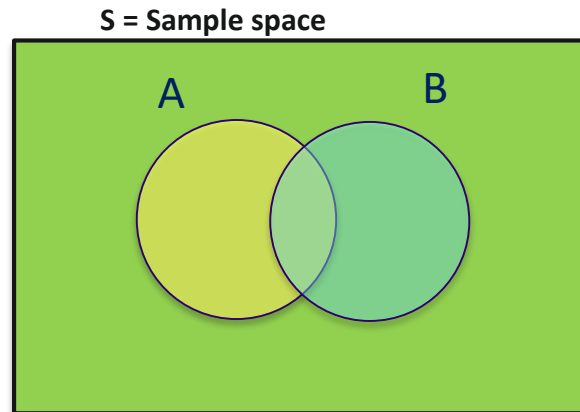
Conditional Probability

> The probability of A *given* B is written:

$$P(A|B)$$

> And is equal to:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad , \text{ compare to: } P(E) = \frac{P(E)}{P(S)}$$



W

Independent Events

- > Events A is independent of B if and only if:

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

- > A being independent of B does NOT imply B is independent of A.

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

 \Rightarrow

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

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

E.g. The event that my boss takes vacation has an impact on when I take vacation, but when I take vacation has no impact on when my boss takes vacation. (i.e., his vacation is independent of mine, but not vice versa)

W

Independence vs. Mutually Exclusive

> These are not related AT ALL and in fact, are nearly opposite ideas.

> If A is M.E. of B then: $P(A|B) = 0$



B occurring has a HUGE impact on $P(A)$

> If A is independent of B then: $P(A|B) = P(A)$

Example: The probability the sidewalk is wet given it is raining is very high,
But the probability that it is raining given the sidewalk is wet is lower (if I run
my sprinklers often).



Odds

- > Odds are expressed as (Count in event favor):(Count not in event favor)
 - Make sure you reduce the fraction first

$$P(A) = \frac{n}{m} = \frac{n}{n + (m - n)}$$

↑ ↑
Count in Count not in
favor of A favor of A

- Implies the odds are:

$$n : (m - n)$$

Examples:

If $P(A)=5/6$, then the odds are 5:1. 'Five to one'.

If the odds are 3:20, then $P(A)=3/23$

A straight up sports bet in Vegas has odds 1:1 (50%), but pays 0.95Xbet.

Monty Hall Problem

- > Famous conditional probability problem that divided statisticians when it came out.
 - Start with 3 doors. One prize behind unknown door. Pick a door. Host reveals a separate door with no prize. Then contestant can switch. Should they?



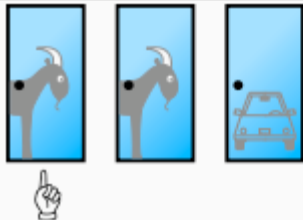
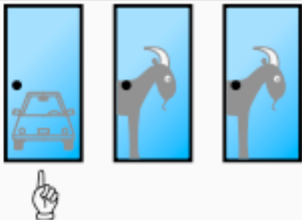
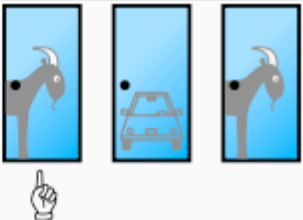
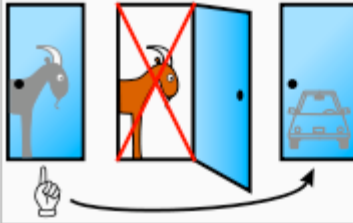
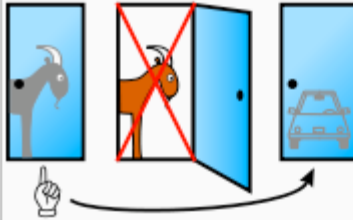


Monty Hall Problem

- > Start with 3 doors. One prize behind unknown door. Pick a door. Host reveals a separate door with no prize. Then contestant can switch. Should they?

Car hidden behind Door 3	Car hidden behind Door 1	Car hidden behind Door 2
Player initially picks Door 1		
		
Host must open Door 2	Host randomly opens either goat door	Host must open Door 3

Monty Hall Problem

- > Start with 3 doors. One prize behind unknown door. Pick a door. Host reveals a separate door with no prize. Then contestant can switch. Should they?

Car hidden behind Door 3	Car hidden behind Door 1		Car hidden behind Door 2
Player initially picks Door 1			
			
Host must open Door 2	Host randomly opens either goat door		Host must open Door 3
			
Probability 1/3	Probability 1/6	Probability 1/6	Probability 1/3
Switching wins	Switching loses	Switching loses	Switching wins
If the host has opened Door 2, switching wins twice as often as staying		If the host has opened Door 3, switching wins twice as often as staying	

W

Monty Hall Problem

- <http://www.stayorswitch.com/>

W

Simulations in R

- > Simulations are used to verify probabilities.
- > With these, we can also estimate variation in probabilities.
- > Use `system.time()` from base or `microbenchmark()` from `microbenchmark` package.
- > Clean up after yourself:
 - `gc()` or `invisible(gc())`
- > R demo



Dealing with Missing Data

> Reasons for missing data

- Recording failure (mechanical/software failures)
- Reporting failure (human decisions)
- Translation failure (data transferring/parsing errors)

> Many shapes and types

- Shapes: block, regular, random, sparse
- Types:

- > Missing At Random (MAR): a particular variable has randomly omitted data.
- > Missing Completely At Random (MCAR) : every piece of data has equal chance of being omitted.
- > Missing Not At Random (MNAR): The value of data is related to chance of being omitted.

> Outliers may also be treated as missing data.



Dealing with Missing Data

Type	Benefits	Disadvantages	Notes
Drop Missing	-Speed	-Data Loss	
Mean/Median/Mode Fill	-No Data Loss	-Variance Reduction	
$X \sim F(\text{independents})$	-More Accurate -No Data Loss	-Slower	-Needs most columns to be filled out -Harder on ind. data
knn	-More Accurate -No Data Loss	-Slower -Dependent on distance function	
$X \sim F(y, \text{independents})$	-Very accurate -No Data Loss	-Slower -Need y	-Only on training set!



Dealing with Missing Data: Variance and Multiple Imputation

- > Dealing with imputation, it is important to try and keep the intrinsic variance in the data set.
- > To achieve this, multiple different predictions are made for each missing data point. (Using previous methods)
- > These data sets are kept and future hypothesis testing and predictions are made on all imputed sets to gauge the variance in the outcomes.
- > R package 'Amelia' does this and creates a nested list of data frames.
- > Amelia R demo



Dealing with Missing Data: Using Outside or New Data Sources

- > Don't forget to explore outside or new data sources to help fill-in missing data.
- > With the advent of free public data and bigger data sources, this is gaining popularity as a tool for imputation.
- > Unstructured text is a major source of data.
- > Ex:
 - Caesar's uses public reviews on websites to mine for customer sentiment about hotel rooms.
 - Zillow uses text descriptions of properties to fill in missing data about # bedrooms, # bathrooms, sq. footage, and various amenities.
 - Subject to human stupidity.

Yelp Rating for Circus-Circus: 2/5

Text Description: "My son and I stayed here. The service was great, the room was great, but it turns out my son is deathly afraid of clowns."



Getting Data

> Files

- Csv: read.csv
- Txt: read.table

> Web/HTML

- readLines
- XML, xpath
- http://gastonsanchez.com/work/webdata/getting_web_data_r4_parsing_xml_html.pdf

> API

- Twitter Example
- Get consumer/access keys here:
 - > <https://dev.twitter.com/apps>



Storing Data

- > .csv – write.csv()

- > .txt – write.txt()

- > .Rdata – save()

 - Workspaces are very compressed compared to csv

- > Databases

 - Sqlite: sqldf, RSQLite packages

 - > Sqlite example

 - MongoDB: rmongodb package

 - Postgresql: RPostgreSQL package



Assignment

> Complete Homework 2:

- Write an R-script to verify the Monty Hall Probabilities with simulations (get probabilities AND variances for switching and not switching).
- You should submit:
 - > **ONE Production level R-script** that outputs the probabilities and variances.
 - > Submission should include a text document/log file of your results.
- Read Intro to Data Science Chapter 7 and 10.
- Read Statistical Thinking for Programmers Ch. 4.

