

You can play around with $P(A)$ as a function of k for fixed $n = 365$ & find that $P(A) > 0.5$ for $k \geq 23$, which many people find surprisingly low. (2 Aug 17)

Generalizing the binomial coefficients

(p.33) what if there are more than 2 possible outcomes $\binom{n}{y}$

In a generalization of the Toy-Sachs case study (T, N) ?
 $\begin{matrix} \rightarrow \\ \text{Toy baby} \end{matrix}$ $\begin{matrix} \leftarrow \\ \text{not Toy baby} \end{matrix}$ we want

n distinct elements to be divided into k different groups ($k \geq 2$) so that n_j elements fall into group j , $\sum_{j=1}^k n_j = n$

Q in how many different ways can this (39) be done?

Follow the argument in DS pp. 42-43, which generalizes the line of reasoning leading to the binomial coefficients $\binom{n}{k}$ when $k=2$:

Definition: A multinomial coefficient is of the form

$$\binom{n}{n_1, n_2, \dots, n_k} = \frac{n!}{n_1! n_2! \dots n_k!}$$

$$\frac{n!}{k!(n-k)!}$$

$n \geq 1$
 $k \geq 2$
 $1 \leq n_j \leq k$
 $\sum_{j=1}^k n_j = n$

This answers the question of how many different ways question Q above

Example: (2016 presidential election) (40)

(see IS p. 333-334) (with replacement)

Imagine randomly sampling n eligible

prospective voters from all such people
"the population" *

in the US. ~~randomly~~ (1 Aug 2016);

possible
outcome ($k=5$)

Clinton (Democrat)

Trump (Republican)

Johnson (Libertarian)

Dein (Green)

Undecided

let $X_i = \#$ people

in sample who say

they will vote for

candidate i ,

$i = 1, \dots, k=5$ as

in this table.

Suppose (unknown to us) that the

proportion of voters who favor candidate

i in the population * above is p_i ,

where $0 < p_i < 1$ and $\sum_{i=1}^k p_i = 1$. (4)

Because the people are chosen with independent identically distributed (IID) sampling (i.e., at random with replacement), each person's outcome will be independent of all the other outcomes. Thus

$P(\text{1st person favors candidate } i_1, \text{ 2nd person favors } i_2, \dots, \text{ n-th person favors } i_n) =$

$p_{i_1} p_{i_2} \dots p_{i_n}$ listed in a prespecified order

Therefore $P(\text{the sample has } x_1 \text{ people favoring candidate 1, } x_2 \text{ people favoring candidate 2, } \dots, x_k \text{ people favoring candidate } k) = p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$,

with $0 \leq x_i \leq n$ and $\sum_{i=1}^k x_i = n$. Thus (4)

$P(\text{exactly } x_1 \text{ people favor Clinton, } \dots, x_k \text{ people favor Undecided}) = \boxed{?} p_1^{x_1} \dots p_k^{x_k}$, (5)

where $\boxed{?}$ is the total # of different ways the order of the n people in the sample can be listed.

But this $\boxed{?}$ is precisely

$$\binom{n}{x_1, \dots, x_k} = \frac{n!}{x_1! x_2! \dots x_k!}$$

the multinomial coefficient defined on p. (39) above.

Thus

$$P(\overset{(7) \dots (7)}{X_1 = x_1, \dots, X_n = x_n}) = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}$$

later in the course we'll refer to (48)
this as the multinomial (probability)

distribution

We already
worked out that

How to work with OR
when you have more \updownarrow
than 2 events (union) \cup

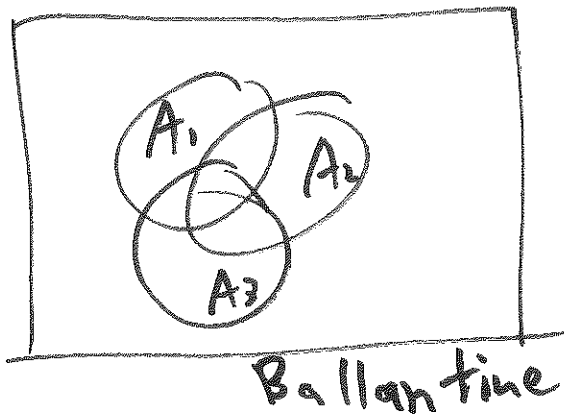
$$P(A_1 \text{ or } A_2) = P(A_1 \cup A_2)$$

\downarrow and

$$= P(A_1) + P(A_2) - P(A_1 \cap A_2).$$

we also know from Kolmogorov's 3rd
Axiom that if events A_1, \dots, A_n are
disjoint then $P\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n P(A_i).$

How do these 2 things generalize?



By (tedious) enumeration 44
 you can show that
 with 3 events,

$$P(A_1 \cup A_2 \cup A_3) = P(A_1) + P(A_2) + P(A_3) \\
 - [P(A_1 \cap A_2) + P(A_1 \cap A_3) + P(A_2 \cap A_3)] \\
 + P(A_1 \cap A_2 \cap A_3).$$

You can ^{probably} now

~~(or guess)~~
 see how this generalizes: for $n > 3$
 events A_1, \dots, A_n ,

$$P(\bigcup_{i=1}^n A_i) = \sum_{i=1}^n P(A_i) - \sum_{i < j} P(A_i \cap A_j) \\
 + \sum_{i < j < k} P(A_i \cap A_j \cap A_k) + \dots + \\
 (-1)^{n+1} P(A_1 \cap \dots \cap A_n).$$

Example Get 2 decks of ordinary playing cards; order deck 1 from (1 to 52) using any sequence you like, e.g.)

- 1 = 2♠
- ⋮
- 13 = A♠
- 14 = 2♦
- ⋮
- 26 = A♦
- 27 = 2♥
- ⋮
- 39 = A♥
- 40 = 2♣
- ⋮
- 52 = A♣

(practically speaking) Shuffle deck 2 until all 52! orderings are equally likely.

Now turn the first card of each deck over; do they match? Continue through all 52 cards; P(at least one match) = ?

let $n=52$

Let $A_i =$ (a match occurs on card i), we want $P(\bigcup_{i=1}^n A_i)$, which can be computed with the complicated formula on the previous page.

Follow the logic detailed on Dr (46)
to obtain p. 49-50

$$P(\bigcup_{i=1}^n A_i) = \frac{1}{1!} - \frac{1}{2!} + \frac{1}{3!} - \dots + (-1)^{n+1} \frac{1}{n!}$$

calculus result: $\lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{(-1)^{i+1}}{i!} = 1 - \frac{1}{e} = 0.63$

This sum approaches its limit quickly; already with $n=7$ you have the first 4 significant figures: 0.6321

JS ch. 2
Conditional probability

Note that Kolmogorov's probability axioms defined the function

$P_k(A)$, where A is a set in the

collection \mathcal{C} of subsets of the sample space \mathcal{S} in which nothing weird can occur; in other words, $P_K(A)$ is a function of a single argument A .

To include the extremely useful idea of conditional probability in his setup, Kolmogorov has to define it using P_K .

Definition Given any two events A, B in \mathcal{C} , the conditional probability of A given B is

$$P(A|B) = \begin{cases} \frac{P(A \cap B)}{P(B)} & \text{if } P(B) > 0 \\ \text{undefined} & \text{if } P(B) = 0 \end{cases}$$

There are other foundational theories (48)
of probability - one by the Italian
mathematician and actuary ^(deF) Bruno de Finetti (1906-1985),
and another by the American physicists
Richard T. Cox (1898-1991) and Edwin
T. Jaynes (1922-1998) ^(CJ) - in which the
probability function $P_{deF}(A|B)$ or
 $P_{CJ}(A|B)$ has 2 inputs, not 1,
so that conditional probability is
the primitive concept, not
unconditional probability as with
Kolmogorov's $P_K(A)$. deF and CJ

were responding to the reality that (49)

in practice, all probabilities are conditional on background Assumptions, Information

and Judgments (AIJ)

Example

(Toy-Sachs)

we actually computed not

$P(\text{at least 1 t-s baby})$ but

$P(\text{at least 1 t-s baby} \mid \text{family of 5, } \overset{\text{and}}{\text{mother and father both carriers}})$

This impulse, to be explicit about your

AIJ, is Bayesian; Kolmogorov worked

in the frequentist paradigm; in this

course, focusing on $P_K(B)$, we need to

remember that it should really be $P_K(B|AIJ)$.

Consequences
of the

conditional

probability

definition

(theorems)

① A, B events in \mathcal{C} : 50

if $P(B) > 0$ then

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

and if $P(A) > 0$

$$\text{then } P(A \cap B) = P(A) P(B|A).$$

② Direct generalization: if A_1, \dots, A_n
are events with $P(A_1 \cap \dots \cap A_{n-1}) > 0$;

then

chain rule for $n = \text{and}$

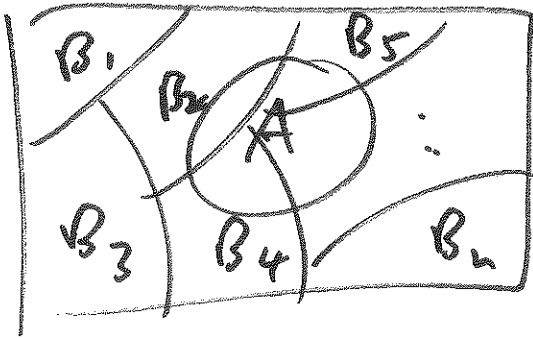
$$P(A_1 \cap \dots \cap A_n) = P(A_1) P(A_2 | A_1) P(A_3 | A_1 \cap A_2) \\ \dots P(A_n | A_1 \cap \dots \cap A_{n-1})$$

recall previous

Definition

\mathcal{S} sample space; if you
can find events B_1, \dots, B_k in \mathcal{C}

such that the B_j are disjoint and (5) exhaustive ($\bigcup_{i=1}^n B_i = S$), then you have found a partition (B_1, \dots, B_k) of S .



(3) If (B_1, \dots, B_k) is a partition of S

with $P(B_j) > 0$ for all $j = 1, \dots, k$,

then for any event A in C

$$P(A) = \sum_{j=1}^k P(B_j) P(A|B_j) -$$

this is the

Law of Total Probability

LTP

When is the LTP useful!

(51.1)

You're trying to compute $P(A)$ and you find it

hard to compute directly. If you can find some aspect B of the

world satisfying 2 properties -

- ① B defines a partition $\{B_1, \dots, B_k\}$ of \mathcal{S} with known $P(B_j)$
- ② A depends on B in such

a way that the conditional probabilities

$P(A|B_j)$ are easier to compute than

$P(A)$ itself - then you can work out

$$P(A) \text{ indirectly: } P(A) = \sum_{j=1}^k P(B_j) P(A|B_j)$$

(Bayesian mixture modeling)

Extension of
the LTP (4)

Assuming all conditional probabilities are defined (52)

in what follows, if C is in \mathcal{C} then

$$P(A|C) = \sum_{j=1}^k P(B_j|C) P(A|B_j \cap C).$$

Definition

Events A, B are
independent iff

$$P(A \cap B) = P(A) \cdot P(B) \quad \leftarrow \text{(freq)}$$

which (as long as $P(A) > 0, P(B) > 0$)

is equivalent to

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

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

$\left. \vphantom{\begin{matrix} P(A|B) = P(A) \\ P(B|A) = P(B) \end{matrix}} \right\} \leftarrow \text{(Bayesian)}$

Consequences
of the
definition
of independence

① If A and B are 53
independent, then so are
 A and B^c , A^c and B ,
and A^c and B^c .

② Extension of the definition to
more than 2 events:

Definition:

Given events A_1, \dots, A_k , they are
(mutually) independent if, for

every subset A_{i_1}, \dots, A_{i_j} of (A_1, \dots, A_k)
($j = 2, \dots, k$),

$$P(A_{i_1} \cap \dots \cap A_{i_j}) = P(A_{i_1}) \cdot \dots \cdot P(A_{i_j})$$

(Bayesian)
Interpretation
of independence

A, B independent \iff 54
information about A

doesn't change the chances associated with B , and vice versa.

Definition Another ^{useful} extension of independence

Events $\{A_1, \dots, A_k\}$ are conditionally independent given event B if for every subset $\{A_{i_1}, \dots, A_{i_j}\}$ of $\{A_1, \dots, A_k\}$ ($j = 2, \dots, k$)

$$P(A_{i_1} \cap \dots \cap A_{i_j} | B) = \prod_{l=1}^j P(A_{i_l} | B)$$

← product

Statistical Example | Suppose that
There is a machine that can ⁽⁵⁵⁾ take an ordinary coin and produce IID

tosses of the coin with $P(H) = \theta$,
and θ can be set to any value in $[0, 1]$

with a dial on the machine's control panel.

Someone sets the dial to a θ value that's unknown to you and

starts producing coin tosses I_1, I_2, \dots

Suppose the first 10 tosses come out

1 0 1 1 1 0 0 1 1 1
HTHTHTTTHHH

← "bits" (binary digits)
(7 H, 3 T).

Q: Is there information in these first 10 tosses that helps you to predict I_{11} ?

A: Yes, definitely: it looks like (56)

θ is around $\frac{7}{10}$, so you would predict

$I_{11} = H$. Thus I_{11} depends on I_1, \dots, I_{10} probabilistically.

Now, suppose instead that you watched the guy with the machine

set the dial to $\theta = 0.81$, so that

θ is now known to you. The next 10

tosses come out H H H T H T H H H H

(8 H, 2 T). **Q:** Is there information

in these 10 tosses that helps you

to predict the next toss?

A: No; you know that $\theta = 0.81$, so there's no information in any of the I_v

that helps you to predict any of (57)

the other I_j .

given θ , the I_i are indep.

Thus the I_i are

unconditionally dependent but

conditionally independent given θ .

(4 Aug 17)

Bayes's
Theorem
for events

(a finite partition)

Suppose that the events B_1, \dots, B_k partition the

sample space in such a way that

$P(B_j) > 0$ for all $j = 1, \dots, k$. If A

is an event with $P(A) > 0$, then for

each

$i = 1, \dots, k$

$$P(B_i | A) = \frac{P(B_i) P(A | B_i)}{P(A)}$$