Later in the course we'll refer to
this as the multinomial (probability)
distribution

(16 Apr 19)

We already
worked out that

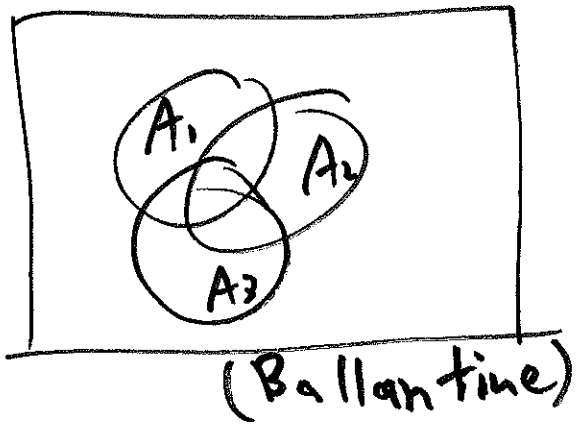How to work with $\boxed{OR}$
when you have more ↑↓
than 2 events ($A_{1:n}$, $\cup$)

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

⟨and⟩ ↓

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

we also know from kolmogorov's 3rd

Axiom that if events $A_1, ..., A_n$ 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 you can show that with 3 events,

$$P(A_1 \cup A_2 \cup A_3) = P(A_1) + P(A_2) + P(A_3)$$

$$- \left[ P(A_1 \cap A_2) + P(A_1 \cap A_3) + P(A_2 \cap A_3) \right]$$

$$+ P(A_1 \cap A_2 \cap A_3).$$

You can now probably see how this generalizes: for any events $A_1, \ldots, A_n$,

(or guess)

$$P\left( \bigcup_{i=1}^{n} A_i \right) = \sum_{i=1}^{n} P(A_i) - \sum_{i<j} P(A_i \cap A_j)$$

$$+ \sum_{i<j<k} P(A_i \cap A_j \cap A_k) + \ldots +$$

$$(-1)^{n+1} P(A_1 \cap \ldots \cap A_n).$$

(Ballantine)

**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 $P\left(\bigcup_{i=1}^{n} A_i\right)$, which can be computed with the complicated formula on the previous page.

Follow the logic detailed on Dr pp. 45-50 to obtain

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

wolfram alpha

limit (sum (-1)^(i+1)/(i!, i = 1 to n) as n → infinity

calculus result:

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

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

DS ch.4

Conditional probability

Note that kolmogorov's probability axioms defined the function $P_k(A)$, where $A$ is a set in the

collection $C$ of subsets of the sample space $S$ 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 $C$, the conditional probability of $A$ given $B$ is

$$P(A \mid B) = \begin{cases} \dfrac{P(A \cap B)}{P(B)} & \text{if } P(B) > 0 \\[2mm] \text{undefined} & \text{if } P(B) = 0. \end{cases}$$

There are other foundational theories (48)
of probability — one by the Italian
mathematician and actuary Bruno de Finetti, (deF) (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 $\underline{\underline{2}}$ inputs, not 1,
so that <u>conditional</u> probability is
the primitive concept, not
<u>unconditional</u> probability as with
Kolmogorov's $P_K(A)$. deF and CJ

were responding to the reality that

in practice, all probabilities are conditional
on background Ⓐssumptions, Ⓘnformation
and Ⓙudgments (AIJ) ⌐Example⌐ (Tay-Sachs)

we actually computed not

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

$$P\left(\begin{array}{c}\text{at least 1} \\ \text{t-s baby}\end{array} \middle| \begin{array}{c}\text{family of 5, mother} \overset{\text{and}}{\wedge} \text{and} \\ \text{father both carriers}\end{array}\right)$$

This impulse, to be explicit about your
AIJ, is <u>Bayesian</u>; Kolmogorov worked
in the <u>frequentist</u> 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 $\mathbb{C}$:

if $P(B) > 0$ then

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

and if $P(A) > 0$

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

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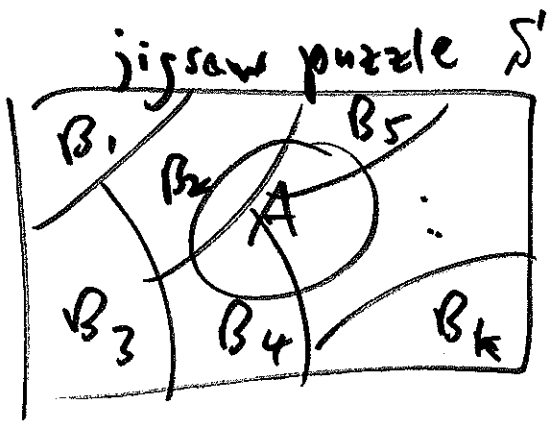are events with $P(A_1 \cap \ldots \cap A_{n-1}) > 0$;

then

$$\boxed{\text{chain rule for } \cap = \text{(and)}}$$

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

recall previous
Definition

$S$ sample space; if you
can find events $B_1, \ldots, B_k$ in $\mathbb{C}$

such that the $B_j$ are disjoint and ⑤

exhaustive $\left( \overset{n}{\underset{i=1}{\cup}} B_i = S \right)$, then you

have found a <u>partition</u> $(B_1, \ldots, B_k)$ of $S$.

jigsaw puzzle $S$



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

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

then for any event $A$ in $C$

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

this is the $\boxed{\text{law of Total probability}}$

LTP

When is the LTP useful?

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 $\{B_1, ..., B_k\}$ with known $P(B_j)$ of $S$ and ② $A$ <u>depends</u> on $B$ in such a way that the conditional probabilities $P(A \mid B_j)$ are easier to compute than $P(A)$ itself — then you can work out

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

$P(A)$ indirectly:

(Bayesian mixture modeling)

(18 Apr 19)