## **Decision Theory**



Thomas Bayes, Pierre Simon de Laplace,, Bruno de Finetti, Alan Turing, Irving Good, Leonard Jimmie Savage, Dennis Lindley, Arnold Zellner, Kathryn Chaloner, Susie Bayarri, Daniel Kahneman Who am I?

## Anders Nordgaard

Reader and Forensic specialist in statistics Swedish Police Authority – National Forensic Centre.

Former senior lecturer and director of studies at the Division of Statistics (and Machine Learning), LiU.

Nowadays, adjunct lecturer at this division (up to 20 % of full time)

Teaching this course Supervision of Master's thesis work

Easiest way of contact: andno100@gmail.com

## A course on decision making under uncertainty – Reasoning with probabilities

• <u>Course responsible and tutor:</u>

Anders Nordgaard (andno100@gmail.com, Anders.Nordgaard@liu.se)

• <u>Course web page:</u>

www.ida.liu.se/~732A66

*Note:* There is no course room in Lisam for this course (due to ignorance with the course responsible)

## • <u>Teaching:</u>

Lectures on theory Seminars with complex problems Discussion of assignments

- <u>Course literature:</u>
  - Peterson M.: An Introduction to Decision Theory 2nd ed. Cambridge University Press, 2017. ISBN 9781316606209 (paperback), 9781316585061 (digital)



## • Course literature:

 Peterson M.: An Introduction to Decision Theory 2nd ed. Cambridge University Press, 2017. ISBN 9781316606209 (paperback), 9781316585061 (digital)

#### Former course literature also works:

- Winkler R.L.: An Introduction to Bayesian Inference and Decision 2nd ed. Probabilistic Publishing, 2003 ISBN 0-9647938-4-9
- Electronic version available for purchase or lending: https://archive.org/details/introductiontoba00robe/page/n8/mode/1up
- The relevant exercises from this book will temporarily be uploaded to the course web
- Additional literature:
  - Taroni F., Bozza S., Biedermann A., Garbolino P., Aitken C. : Data analysis in forensic science A Bayesian decision perspective, Chichester: Wiley, 2010
  - Gittelson S. (2013). Evolving from Inferences to Decisions in the Interpretation of Scientific Evidence. Thèse de Doctorat, Série criminalistique LVI, Université de Lausanne. ISBN 2-940098-60-3. Available at http://www.unil.ch/esc/files/live/sites/esc/files/shared/These\_Gittelson.pdf

## • Examination:

- Assignments (compulsory to pass)
- Final oral exam (compulsory, decides the grade)

Assignments:

- There will be 4 assignments
- Co-working is permitted...
- ...but each student must submit their own solution
- Insufficient solutions will need supplementary submission

#### Oral exam:

- Normally in a group of 2 students (occasionally 1 student, never 3 or more)
- A discussion on the course contents and concepts with practical examples
- 2 hours duration (1 student: 1 hour)
- Individual feedback and grading

## Outcome of Evaluate course evaluation for study year 2021/22

- Response rate: 32%
- No questions sticking out in the multiple choice questions
- Free-text answers on question 6 and 7

6. What changes do you consider to be possible that would improve the course with respect to, for example, content, teaching principles, administration, teaching ng methods, or examination forms?

The assignments were a little bit hard to understand, it would be nice if we discuss assignments in the class too.

1. Found it sometimes difficult to follow a long some calculations that were presented on the slides. Think for these types of calculations, using the whiteboard would've been better and discuss the steps perhaps.

2. Think that it would've been great to also mention that the next session would've been a discussion seminar to better prepare

7.

# Give examples of content, teaching principles, teaching methods, examination forms, or any other aspect of the course that you consider to have been particularly successful.

Content, teaching principles, methods and examination style.

1. Found Anders way of explaining the concepts to be great. It was also great that he provided with examples on how the concepts were applied.

2. Found the problem assignments to be well planned and fun exercise to solve. The difficulty of the assignments were reasonable as well.

## Opinions taken up at oral exams:

- Less technical slides better slide structures
- Case examples
- More problem discussions
- Tutorials
- More on Bayesian networks
- Relations to machine learning reinforcement learning
- Half-time study rate instead of quarter-time
- Use Lisam
- Pre-scheduled time-points for exams

## Amendments due to last year's course evaluation

- Problem seminars announced in timetable (on course web pages)
- Some slides further reworked with respect to technicalities
- Better follow-ups of assignments in class
- Pres-scheduling time-points for exam

## Lecture 1: Repeat and extend...

Probability and likelihood



## The concept of probability



· th		

Category	Frequency	Probability
		?
	9	0.6
<b>X</b>	3	0.2
No. Contraction of the second	3	0.2

The *probability* of an event is...

- the degree of belief in the event (that the event has happened)
- a measure of the size of the event relative to the size of the universe



The universe, all events in it and the probabilities assigned to each event constitute the *probability space*. Probability of event= *P*(*Event*)

- $0 \le P(Event) \le 1$
- P(Universe) = 1
- If two events, *A* and *B* are mutually exclusive then

P(A or B) = P(A) + P(B)

**"Kolmogorov axioms"** (finite additivity variant)

This does not mean that...

"probabilities and stable relative frequencies are equal" (*Frequentist definition of probability*)

merely...

If any event is assigned a probability, that probability must satisfy the axioms.

#### Example: Coin tossing

Suppose you toss a coin. One possible event is "heads", another is "tails"

If you assign a probability p to "heads" and a probability q to "tails they both must be between 0 and 1.

As "heads" cannot occur simultaneously with "tails", the probability of "heads or tails" is p + q.

If <u>no other event is possible</u> then "heads or tails" = Universe  $\rightarrow$ p + q = 1



## Relevance, Conditional probabilities

An event *B* is said to be *relevant* for another event *A* if the probability (degree of belief) that *A* is true depends on the state of *B*.

The *conditional* probability of A <u>given</u> that B is true is

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



If *B* is true then its *complement*  $\overline{B}(B^C, \neg B)$  is *irrelevant* to consider.

If *A* is to be true under these conditions, only the part of *A* inside *B* should be considered.

This part coincides with (A,B)

The measure of the size of this event must be relative to the size of B

#### Example:

Assume you believe that approx. 1% of all human beings carry both a gene for developing disease *A* and a gene for developing disease *B*.

Further you believe that 10% of all human beings carry the gene for developing disease *B*.

Then as a consequence your degree of belief that a person who has developed disease *B* also carries the gene for developing disease *A* should be 10% (0.01/0.10)

Carrying the gene for *B* is relevant for carrying the gene for *A*.





Reversing the definition of conditional probability:

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

but also... 
$$P(A,B) = P(B|A) \cdot P(A)$$

$$\Rightarrow P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \text{ and } P(B|A) = \frac{P(A|B) \cdot P(B)}{P(A)}$$

 $\rightarrow$  For sorting out conditional probabilities it is not necessary to assign the probabilities of intersections

## "All probabilities are conditional..."

How a probability is assigned <u>depends on</u> background knowledge.

E.g. if you assign the probability 0.5 for the event "heads" in a coin toss, you have assumed that

- the coin is fair
- the coin cannot land endways



...but it may be the case that you cannot assign any probability to the background knowledge

Let I denote all background knowledge relevant for A

 $\Rightarrow P(A) = P(A|I)$ 

**Extensions:** 

$$P(A, B|I) = P(A|B, I) \cdot P(B|I)$$
  

$$P(A_1, A_2, ..., A_n|I) =$$
  

$$= P(A_1|I) \cdot P(A_2|A_1, I) \cdot \cdots \cdot P(A_n|A_1, A_2, ..., A_{n-1}, I)$$

*Example*: Suppose you randomly pick 3 cards from a well-shuffled deck of cards. What is the probability you will <u>in order</u> get a spade, a hearts and a spade?

I = The deck of cards is well-shuffled  $\Rightarrow$  It does not matter how you pick your cards.

Let  $A_1$  = First card is a spade;  $A_2$  = Second card is a hearts;  $A_3$  = Third card is a spade

$$\Rightarrow P(A_1, A_2, A_3 | I) = P(A_1 | I) \cdot P(A_2 | A_1, I) \cdot P(A_3 | A_1, A_2, I) = = \frac{13}{52} \cdot \frac{13}{51} \cdot \frac{12}{50} \approx 0.015$$

## Relevance and (conditional) independence

If *B* is relevant for *A* then  $P(A|B,I) \neq P(A|I)$ 

If *B* is *irrelevant* for *A* then P(A|B,I) = P(A|I)which in turn gives  $P(A,B|I) = P(A|I) \cdot P(B|I)$ 

In this case *A* and *B* are said to be <u>conditionally independent</u> events. (In common statistical literature only *independent* is used as term.)

Note that it is the background knowledge *I* that determines whether this holds or not.

Note also that if P(A|B,I) = P(A|I) then P(B|A,I) = P(B|I)

Irrelevance is reversible!

Below are four rectangles. Each rectangle represents the universe, so its area is equal to one (1=100%)

Assume that the sets A (green) and B (yellowish) are drawn according to scale (the sizes of the sets are proportional to the probabilities of the events).

In which of the cases below are *A* and *B* <u>definitely</u> conditionally <u>dependent</u> (given *I*)?



## Further conditioning...



## $P(A, B|I) \neq P(A|I) \cdot P(B|I)$



 $P(A, B|C, I) = P(A|C, I) \cdot P(B|C, I)$ 

Two events that are conditionally dependent under one set of assumptions may be conditionally *independent* under another set of assumptions



The law of total probability:

 $P(A|I) = P(A, B|I) + P(A, \overline{B}|I) =$ =  $P(A|B, I) \cdot P(B|I) + P(A|\overline{B}, I) \cdot P(\overline{B}|I)$ 

 $\Rightarrow$  Bayes' theorem:

 $P(A|B,I) = \frac{P(B|A,I) \cdot P(A|I)}{P(B|A,I) \cdot P(A|I) + P(B|\overline{A},I) \cdot P(\overline{A}|I)}$ 

## Example:

Assume a method for detecting a certain kind of dye on banknotes is such that

• it gives a positive result (detection) in 99 % of the cases when the dye is present, i.e. the proportion of false negatives is 1%

• it gives a negative result in 98 % of the cases when the dye is absent, i.e. the proportion of false positives is 2%

The presence of dye is rare: prevalence is about 0.1 %

Assume the method has given positive result for a particular banknote.

What is the conditional probability that the dye is present?



#### Solution:

Let *A* = "Dye is present" and *B* = "Method gives positive result" What about *I* ?

• We must assume that the particular banknote is as equally likely to be exposed to dye detection as any banknote in the population of banknotes.

• Is that a realistic assumption?

Now, 
$$P(A) = 0.001; P(B|A) = 0.99; P(B|\overline{A}) = 0.02$$

#### Applying Bayes' theorem gives

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B|A) \cdot P(A) + P\left(B|\overline{A}\right) \cdot P(\overline{A})} =$$
$$= \frac{0.99 \cdot 0.001}{0.99 \cdot 0.001 + 0.02 \cdot 0.999} =$$

which makes ...?

## Odds and Bayes' theorem on odds form

The *odds* for an event A "is" a quantity equal to the probability:

$$Odds(A) = \frac{P(A)}{P(\overline{A})} = \frac{P(A)}{1 - P(A)} \implies P(A) = \frac{Odds(A)}{Odds(A) + 1}$$

Why two quantities for the same thing?

### Example: An "epidemiological" model

Assume we are trying to model the probability p of an event (i.e. the prevalence of some disease).

The *logit link* between p and a set of k explanatory variables  $x_1, x_2, \ldots, x_k$  is

$$\operatorname{logit}(p) = \ln \frac{p}{1-p} = \beta_0 + \beta_1 \cdot x_1 + \dots + \beta_k \cdot x_k$$

This link function is common in *logistic regression analysis*.



Note that we are modelling the natural logarithm of the odds instead of modelling p.

As the odds can take any value between 0 and  $\infty$  the logarithm of the odds can take any value between  $-\infty$  and  $\infty \rightarrow$  Makes the model practical.

## Conditional odds

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

express the updated belief that A holds when we take into account that B holds

Like probabilities, all odds are conditional if we include background knowledge *I* as our basis for the calculations. P(A|I) P(A|B,I)

$$Odds(A|I) = \frac{P(A|I)}{P(\overline{A}|I)}; \quad Odds(A|B,I) = \frac{P(A|B,I)}{P(\overline{A}|B,I)}$$

The odds ratio:

$$OR = \frac{Odds(A|B,I)}{Odds(A|I)} = \frac{\frac{P(A|B,I)}{P(\overline{A}|B,I)}}{\frac{P(A|I)}{P(\overline{A}|I)}}$$

expresses how the belief that A holds updates when we take into account that B holds.

Now  

$$\frac{Odds(A|B,I)}{P(\overline{A}|B,I)} = \frac{P(A|B,I)}{P(\overline{A}|B,I)} = \frac{\frac{P(B|A,I) \cdot P(A|I)}{P(B|I)}}{\frac{P(B|\overline{A},I) \cdot P(\overline{A}|I)}{P(B|I)}} = \frac{P(B|A,I)}{P(B|\overline{A},I)} \cdot \frac{P(A|I)}{P(\overline{A}|I)} = \frac{P(B|A,I)}{P(B|\overline{A},I)} \cdot Odds(A|I)$$

"Bayes' theorem on odds form"



is a special case of what is called a *likelihood ratio* (the concept of "likelihood" will follow)

 $LR = \frac{P(B|A, I)}{P(B|C, I)}$ 

where we have substituted C for  $\overline{A}$  and we no longer require A and C to be complementary events (not even mutually exclusive ).

 $\frac{P(A|B,I)}{P(C|B,I)} = \frac{P(B|A,I)}{P(B|C,I)} \cdot \frac{P(A|I)}{P(C|I)}$ 

always holds, but the ratios involved are not always odds

"The updating of probability ratios when a new event is observed goes through the likelihood ratio based on that event."

## Probability and Likelihood – Synonyms?

An event can be *likely* or *probable*, which for most people would be the same. Yet, the definitions of probability and likelihood are different.

## In a simplified form:

- The probability of an event measures the degree of belief that this event is true and is used for reasoning about not yet observed events
- The likelihood of an event is a measure of how likely that event is in light of another *observed* event
- Both are objected to probability calculus

## More formally...

Consider the *unobserved* event A and the *observed* event B.

There are probabilities for both representing the degrees of belief for these

events in general: P(A|I), P(B|I)

However, as *B* is observed we might be interested in

P(A|B,I)

which measures the *updated* degree of belief that A is true once we know that B holds. Still a probability, though.



 $P(B \mid A, I)$  might look meaningless to consider as we have actually observed *B*. However, it says something about *A*.

We have observed *B* and if *A* is relevant for *B* we may compare P(B | A, I) with  $P(B | \overline{A}, I)$ .

Now, even if we have not observed A or  $\overline{A}$ , one of them must be true (as a consequence of A and B being relevant for each other).

If  $P(B | A, I) > P(B | \overline{A}, I)$  we may conclude that A is more *likely* to have occurred than is  $\overline{A}$ , or better phrased:

"A is a better *explanation* to why B has occurred than is  $\overline{A}$ ".

P(B | A, I) is called the *likelihood* of A given the observed B (and  $P(B | \overline{A}, I)$ ) is the likelihood of  $\overline{A}$ ).

*Note!* This is different from the conditional probability of A given B: P(A | B, I).

#### Potential danger in mixing things up:

When we say that an event is the more likely one in light of data we do not say that this event has the <u>highest probability</u>.

Using the likelihood as a measure of how likely is an event is a matter of *inference to the best explanation*.

Logics: Implication:

 $A \rightarrow B$ 

- If *A* is true then *B* is true, i.e.  $P(B | A, I) \equiv 1$
- If *B* is false then *A* is false, i.e.  $P(A | \overline{B}, I) \equiv 0$

• If *B* is true we cannot say anything about whether *A* is true or not (implication is different from equivalence)

"Probabilistic implication":

 $A \xrightarrow{P} B$ 

- If A is true then B may be true, i.e. P(B|A, I) > 0
- If *B* is false the *A* may still be true, i.e.  $P(A|\overline{B}, I) > 0$
- If *B* is true then we may decide which of *A* and  $\overline{A}$  that is the best explanation

Inference to the best explanation:

- *B* is observed
- $A_1, A_2, \ldots, A_m$  are potential alternative explanations to B

• If for each  $j \neq k$   $P(B \mid A_k, I) > P(B \mid A_j, I)$  then  $A_k$  is considered the best explanation for *B* and is provisionally accepted

#### LINKÖPING UNIVERSITY