

Lecture 2

Cezar Ionescu

Probability theory

- The first question of probability theory is “probability of what?”.
 - Richard von Mises (1883-1953): “The term probability is meaningful for us only with regard to a clearly defined collective¹ (or population.” (*Mathematical Theory of Probability and Statistics*, page 17).
 - Sir Harold Jeffreys (1891-1989): “there is a valid primitive idea expressing the degree of confidence that we may reasonably have in a proposition” (*Theory of Probability*, 3rd Ed., page 15).
- The currently accepted mathematical model was developed by Kolmogorov in 1933 and covers both interpretations.
 - Ingredients:
 - * Ω , the set of all possible outcomes
 - * **Event**, the set of all events $\text{Event} \subseteq \mathbb{P}(\Omega)^2$, which are what we actually want to assign probabilities to. The set of events is assumed to satisfy certain elementary properties (forming a σ -algebra of Ω).
 - * $p : \text{Event} \rightarrow [0, 1]$, the probability measure, satisfying
 1. $p(\emptyset) = 0$, $p(\Omega) = 1$
 2. for all events X_1, \dots, X_n, \dots pairwise disjoint, $\sum p(X_i) = p(\bigcup X_i)$

Examples

1. Drawing a card from a standard deck. Any of the 52 cards is a possible result, so $\Omega = \{2\clubsuit, 2\spadesuit, 2\heartsuit, 2\diamonds, \dots\}$. An event is what we want to assign a probability to, for example “the card drawn is a spade”, or “a red card between 7 and 10”. Any event can be represented by a subset of Ω ; when Ω is finite, we take *all* subsets, so that $\text{Event} = \mathbb{P}(\Omega)$.
2. Choosing a point of the unit disc. This can be seen as an idealised model for a game of darts. $\Omega = \{(x, y) \mid x, y \in \mathbb{R}, x^2 + y^2 \leq 1\}$. Events are of the form “the point will be chosen from this or that subset of the unit disc”, so in principle we would

¹Intuitively, a *collective* is an infinite sequence of results of a repeated experiment.

² $\mathbb{P}(X)$ denotes the set of all subsets of X , including \emptyset and X itself.

like again to have $\text{Event} = \mathbb{P}(\Omega)$. It is an annoying but inescapable mathematical fact that we cannot, since that would make it impossible to define a probability measure satisfying the conditions above. Therefore, we limit the number of subsets that we assign probabilities to. The standard mathematical choice is that of a *Borel σ -algebra*, which is “almost as big” as $\mathbb{P}(\Omega)$. We can be even less ambitious and choose Event to be the set of all *computable* subsets of Ω (i.e., all those whose characteristic function $X : \Omega \rightarrow \{0, 1\}$ can be implemented as a program).

The classical model of probability

Ω is finite, $\text{Event} = \mathbb{P}(\Omega)$ and the probability measure is defined by

$$p(X) = \text{card}(X) / \text{card}(\Omega)$$

where $\text{card}(X)$ is the number of elements in the finite set X .

Interpretations:

- Frequentist: if we were to repeat the experiment many times, every elementary outcome $\omega \in \Omega$ would turn up with approximately the same frequency as any other.
- Bayesian: in the absence of any information, we have no reason to prefer one elementary outcome over any other, so we must have $p(\omega_i) = p(\omega_j)$ for all i, j .

Example: rolling a die

- $\Omega = \{1, 2, 3, 4, 5, 6\}$, $\text{card}(\Omega) = 6$
- rolling an even number: $\text{even} = \{2, 4, 6\}$, $\text{card}(\text{even}) = 3$
- probability of rolling an even number: $p(\text{even}) = \text{card}(\text{even}) / \text{card}(\Omega) = 3/6 = 0.5$

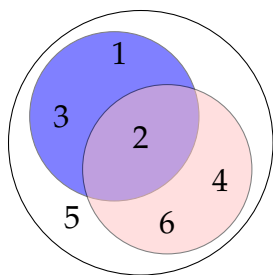
Operations with events:

Since events and predicates are modelled as sets, we can perform on them set-theoretical operations.

- union: $A \cup B$
- intersection: $A \cap B$
- complement: $\neg A = \Omega - A$

Examples:

- Example: probability that the result is even **and** ≤ 3
- Example: probability that the result is even **or** ≤ 3
- Example: probability that the result is **not** ≤ 3

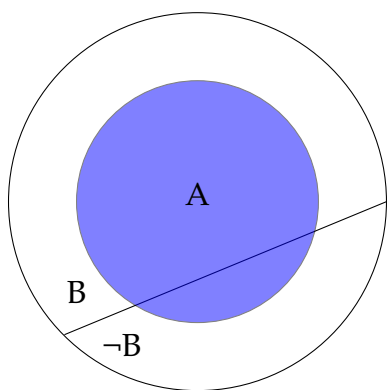


Remark: Operations on sets correspond to logical operations.

Exercise:

1. Derive the formula for $p(\neg A)$.
2. Derive the formula for $p(A \cup B)$
3. Assume that B_1, \dots, B_n are pairwise disjoint. Show that

$$p(A) = p(A \cap B_1) + \dots + p(A \cap B_n)$$



Conditional probability

The following definition is perhaps the most misunderstood aspect of elementary probability theory.

Definition: Let $(\Omega, \text{Event}, p)$ as above, and let $X \in \text{Event}$ such that $p(X) \neq 0$. Assume that event X has been realised (the elementary result of an experiment, ω , is an element of X). The probability that another event, $Y \in \text{Event}$, has *also* been realised is denoted by $p(Y \mid X)$ and defined as

$$p(Y \mid X) = p(Y \cap X) / p(X)$$

The notation $p(Y \mid X)$ is **very bad**. It suggests that \mid is a set-theoretical operation, just like \cup or \cap . After all, p is defined on events, and it seems to be applied to $Y \mid X$, which therefore must be an event, hence a subset of Ω . The notation also suggests a symmetry between X and Y which does not exist: the two events have completely different roles. The notation used by Kolmogorov in his original publication of 1933 was $p_X(Y)$. This makes it clear that we are dealing with **another probability measure**, namely, one that reflects the

assumption that X has been realised, and there is no hint of symmetry. Unfortunately, the inferior notation has been universally adopted, and we have no choice but to follow suit.

Exercise (the law of total probability): Assume that B_1, \dots, B_n are pairwise disjoint, such that $p(B_i) \neq 0$ for all i . Show that

$$p(A) = p(A \mid B_1) * p(B_1) + \dots + p(A \mid B_n) * p(B_n)$$

Bayesian learning

Recall the setting from lecture 1. We have a set of hypotheses, H , and are given some data, $d \in D$. We want to find the “best” hypothesis that fits the data. The question is: what does “best” mean? In the context of **Find-S**, “best” meant the most specific hypothesis consistent with the data. In the much wider (and more reasonable) context of Bayesian learning, “best” means *the most probable* hypothesis, given the data.

In other words, the aim of Bayesian learning is to determine the **maximum a posteriori (MAP) hypothesis**, i.e., the hypothesis that satisfies

$$h_{\text{map}} = \operatorname{argmax} p(_ \mid d)$$

Remarks:

1. The partial function argmax satisfies

$$\operatorname{argmax} : (X \rightarrow \mathbb{R}) \rightarrow X$$

$$\operatorname{argmax} f = x \text{ iff for all } x' \in X \ f(x) \geq f(x')$$

2. The notation $p(h \mid d)$ is an abbreviation for $p(\{h\} \mid \{d\})$ (remember that events are *sets*).

3. The expression $p(_ \mid d)$ denotes a function of argument h :

$$p(_ \mid d) : H \rightarrow [0, 1]$$

$$p(_ \mid d)(h) = p(h \mid d)$$

4. Alternative notations for $p(_ \mid d)$ are $h \mapsto p(h \mid d)$, so we could have written

$$h_{\text{map}} = \operatorname{argmax} (h \mapsto p(h \mid d))$$

or

$$h_{\text{map}} = \operatorname{argmax} (h \in H \mapsto p(h \mid d))$$

to make the domain of the function explicit. In the context of argmax , the most frequent notation is

$$h_{\text{map}} = \operatorname{argmax}_h p(h \mid d)$$

Exercise:

- What specification does the partial function \max satisfy?
- Why are \max and argmax partial?

Bayes' theorem

Consider events X, Y , such that $p(X) * p(Y) \neq 0$. From the definition of conditional probability, we have

$$p(X | Y) = p(X \cap Y) / p(Y), \text{ therefore } p(X \cap Y) = p(X | Y) * p(Y)$$

$$p(Y | X) = p(Y \cap X) / p(X), \text{ therefore } p(Y \cap X) = p(Y | X) * p(X)$$

But $X \cap Y = Y \cap X$, therefore $p(X \cap Y) = p(Y \cap X)$, therefore

$$p(X | Y) = p(Y | X) * p(X) / p(Y)$$

This relatively trivial result is sometimes called "Bayes' theorem".

Example: (Mitchell, pages 157-158)

- $H = \{\text{healthy}, \text{ill}\}$
- There is a test with possible results \ominus, \oplus ("negative", i.e., the disease is not present, and "positive", i.e., the disease is present).
- We are told that only 0.008 of the population have the disease.
- We are also told that the test is not perfectly reliable. If the patient has the disease, the test will be positive only in 98% of the cases. If the patient does not have the disease, then the test will still be positive in 3% of the cases.
- Finally, we are given the data: a patient's test has come back positive. The question is: what is the MAP hypothesis? Is it likelier that the patient is ill, or healthy?

We need express the problem in the language of probability theory. That is, we must find Ω , **Event**, p . This is, at least in simple situations, the most difficult (and important) step. For example:

- Ω = the set of people in the population
- since Ω is finite, **Event** = $\mathbb{P}(\Omega)$
- the experiment consists of drawing a random person from the population, with equal probability, and applying the test. Therefore, the classical model applies: $p(X) = \text{card}(X) / \text{card}(\Omega)$
- the events of interest are: **Healthy**, the subset of the population who are healthy, **Ill**, the subset of the population whose members have the disease, **Pos**, the subset of the population for which the test result was positive, **Neg**, for which the test was negative.
- the given percentages are translated as follows: $p(\text{healthy}) = 0.992$, $p(\text{ill}) = 0.008$, $p(\oplus | \text{ill}) = 0.98$, $p(\oplus | \text{healthy}) = 0.03$
- the question we are asked is to compute $p(\text{ill} | \oplus)$.

We can now try to answer the question by applying Bayes' theorem:

$$\begin{aligned} & p(\text{ill} | \oplus) \\ = & p(\oplus | \text{ill}) * p(\text{ill}) / p(\oplus) \\ = & 0.98 * 0.008 / p(\oplus) \\ = & 0.00784 / p(\oplus) \end{aligned}$$

$$\begin{aligned}
& p(\text{healthy} \mid \oplus) \\
= & p(\oplus \mid \text{healthy}) * p(\text{healthy}) / p(\oplus) \\
= & 0.03 * 0.992 / p(\oplus) \\
= & 0.02976 / p(\oplus)
\end{aligned}$$

Therefore, the MAP hypothesis is healthy.

In general, we have

$$\begin{aligned}
h_{\text{map}} &= \operatorname{argmax}_h p(h \mid d) \\
&= \operatorname{argmax}_h (p(d \mid h) * p(h) / p(d)) \\
&= \operatorname{argmax}_h (p(d \mid h) * p(h)) \text{ -- since } p(d) \text{ is constant w.r.t. } h
\end{aligned}$$

If, additionally, we assume that all hypotheses are equally likely, i.e., $p(h)$ is also constant w.r.t. h ($p(h) = 1/\text{card}(H)$ when H is finite), then we can go one step further and obtain

$$\begin{aligned}
h_{\text{map}} &= \operatorname{argmax}_h p(h \mid d) \\
&= \operatorname{argmax}_h p(d \mid h)
\end{aligned}$$

$p(d \mid h)$, the probability of encountering the data d if we assume that hypothesis h is correct, is called the **likelihood** of d given h . A hypothesis that maximises the likelihood is known, reasonably enough, as a *maximum likelihood hypothesis (ML)*. Therefore

$$h_{\text{ml}} = \operatorname{argmax}_h p(d \mid h)$$

Please note that, in general, $h_{\text{map}} \neq h_{\text{ml}}$! The equality only holds when we assume that all hypotheses are equally likely! In particular, in the example above, this is not the case.

Exercise: what is h_{ml} in the example above?

Homework

Apply Bayes' theorem to answer the following question (Elmer Mode 1966, page 53):

A class in advanced mathematics contains 10 juniors, 30 seniors, and 10 graduate students. Three of the juniors, 10 of the seniors, and 5 of the graduate students received an A in the course. If a student is chosen at random from this class and is found to have earned an A, what is the probability that he is a graduate student?

Bayesian concept learning

Recall that a concept is a subset of, or a boolean function defined on, a given set:

$$c : X \rightarrow \{0, 1\}$$

Each hypothesis $h \in H$ is such a subset or function ($h : X \rightarrow \{0, 1\}$), and we assume that $c \in H$.

The data we are given in the concept learning task consists of pairs

$$d = \{(x_1, c(x_1)), \dots, (x_m, c(x_m))\}$$

The “brute-force” Bayesian concept learning algorithms then returns a (or “the”) MAP hypothesis:

```
given  $H$ ,  $p(h)$ 
for every  $h \in H$ :
    compute  $p(h \mid d) = p(d \mid h) * p(h) / p(d)$ 
    return  $h_{\text{map}} = \text{argmax}_h p(h \mid d)$ 
```

The problem is computing $p(d \mid h)$ and $p(d)$. As we have seen above, the latter is not necessary, but it is very instructive.

To compute the likelihood, we make the following assumption: the data is completely correct, and hypotheses are not “noisy” functions. Therefore, for any $h \in H$, if we assume that $h = c$, then the data d can only arise if h is consistent with it. Therefore:

```
 $p(d \mid h) = \text{if for all } (x_i, b_i) \in d \text{ } h(x_i) = b_i$ 
               then 1
               else 0
```

Now that we have computed $p(d \mid h)$ for every $h \in H$, we can compute the a-priori probability that the data d is realised by applying the law of total probability:

$$p(d) = \sum_h p(d \mid h) * p(h)$$

This computation shows that the a-priori estimate of the probability of data must be consistent with the a-priori estimate of the probability of hypotheses. Normally, we consider data to be more *objective* than the hypotheses, which are more *subjective*, so this can cause some confusion.

Bayesian concept learning and Find-S

If we assume that all hypotheses are equally likely, then the brute-force Bayesian concept learning algorithm will return an ML hypothesis:

$$h_{\text{map}} = h_{\text{ml}} = \text{argmax}_h p(d \mid h)$$

Since $p(d \mid h) = 1$ if h is consistent with the data, and $p(d \mid h) = 0$ otherwise, the algorithm will return any one of the hypotheses in H consistent with d .

Since **Find-S** returns a consistent hypothesis, it therefore follows that it returns a MAP (and an ML) hypothesis too!

Even if p is not constant w.r.t. hypotheses, **Find-S** can still return a MAP hypothesis, if

$$h_i \subseteq h_j \Rightarrow p(h_i) \geq p(h_j)$$

(i.e., p favours the more specific hypothesis), and if there is a unique most specific hypothesis consistent with d in H .