

STA 365: Applied Bayesian Statistics

Boris Babic
Assistant Professor, University of Toronto

Week 1B: Bayes' Theorem



$$\Pr(E|I) = \frac{\Pr(I|E) \Pr(E)}{\Pr(I)}$$

- Simple example: Pr of die landing on 4, conditional on it landing on an even number.
- $\Pr(4|even) = \frac{\Pr(even|4) \Pr(4)}{\Pr(even)}$
- $\Pr(4|even) = \frac{1 \times 1/6}{1/2} = 1/6 \times 1/2 = 1/3$

Total probability

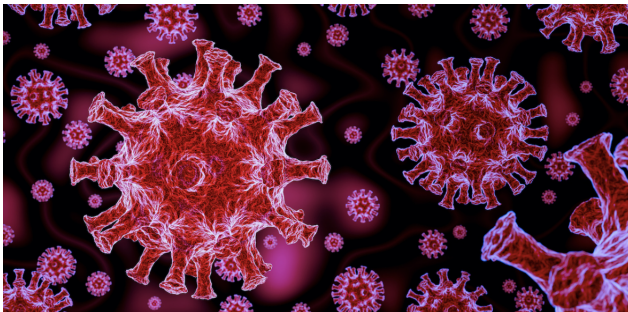
Boris Babic

Bayes'
Theorem

The
Bayesian
Approach

Law of total probability:

$$P(E) = \sum_{j=1}^n P(E|I_j)P(I_j)$$



- Suppose that the base rate (prevalence) of COVID-19 is 1% (ha!).
- Suppose PCR tests have 99% sensitivity. That is, $p(+|D) = 0.99$
- Suppose they also have specificity of 95%. That is, $p(-|\overline{D}) = 0.95$
- Find $p(D|+)$.

- This becomes:

$$p(D|+) = \frac{.99 \times .01}{.99 \times .01 + .05 \times .99} = 0.167$$

- Notice the difference between $p(D|+)$ (0.167) and $p(+|D)$ (0.99)!
- High specificity and sensitivity can still lead to extremely low posterior probability
- A lesson to remember for AI and machine learning!

Medical Screening

Boris Babic

Bayes' Theorem

The Bayesian Approach

Medical screening

$$P(D) = 0.01$$

$$P(\bar{D}) = 0.99$$

sensitivity

$$P(+|D) = 0.99$$

$$P(-|D) = 0.01$$

False -ve

specificity

$$P(-|\bar{D}) = 0.95$$

$$P(+|\bar{D}) = 0.05$$

False +ve

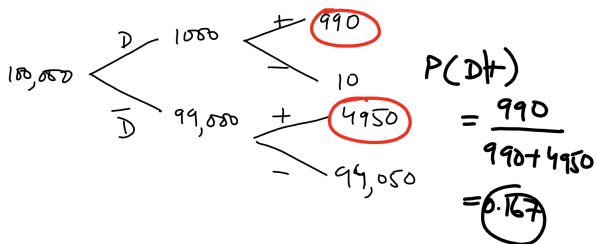
$$P(D|+) = \frac{P(D) P(+|D)}{P(+|D)P(D) + P(+|\bar{D})P(\bar{D})}$$
$$= 0.167$$

Medical Screening

Boris Babic

Bayes' Theorem

The Bayesian Approach



- Find $p(D|-)$

$$p(D|-) = \frac{p(D) \times p(-|D)}{p(D) \times p(-|D) + p(-|\bar{D}) \times p(\bar{D})}$$

- $p(D|-) = (0.01 \times 0.01) / (0.01 \times 0.01 + 0.95 \times 0.99) = 0.0001$

Marginal and Conditional Independence

Boris Babic

Bayes' Theorem

The Bayesian Approach

- E_2 is marginally independent of E_1 if learning about E_1 's occurrence doesn't affect the probability of E_2 :

$$p(E_2|E_1) = p(E_2)$$

- In this case, it would be completely redundant to learn about E_1 .
- Medical tests are not like this. Once you test positive a second time, it is still informative that you have already tested positive once.
- E_2 is conditionally independent of E_1 if learning about E_1 's occurrence doesn't affect the probability of E_1 provided that one has learned about the occurrence of a third event, E_3 :

$$p(E_2|E_1, D) = p(E_2|D)$$

- This is true of medical tests: they are independent conditional on the true disease state.

$$\text{Conditional Independence : } P(t_2 | D, t_1) = P(t_2 | D)$$

Total Redundancy

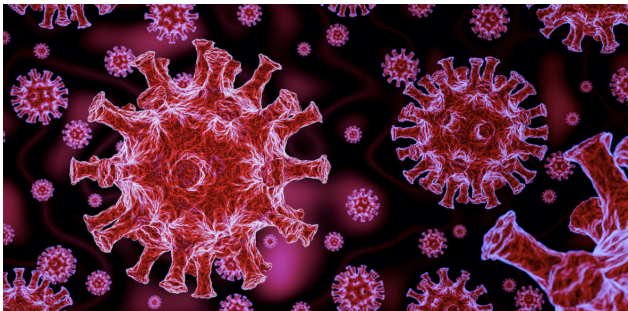
$$\text{Marginal independence : } P(t_2 | t_1) = P(t_2)$$

with conditional independence :

$$P(D) \xrightarrow{1^{\text{st}} \text{ test}} P(D | t_1)$$

$$P(D | t_1) \xrightarrow{2^{\text{nd}} \text{ test}} P(D | t_1, t_2)$$

$$P(D | t_1, t_2) = \frac{P(D) P(t_1, t_2 | D)}{P(t_1, t_2 | D) P(D) + P(t_1, t_2 | \bar{D}) P(\bar{D})}$$



- Suppose again that the base rate (prevalence) of COVID-19 is 1%.
- As before PCR tests have 99% sensitivity. That is, $p(+|D) = 0.99$
- And they also have specificity of 95%. That is, $p(-|\overline{D}) = 0.95$
- You have tested positive once and found that $p(D|+) = 0.167$
- Find $p(D|+_1, +_2)$.

$$P(D|t_1, t_2) = ??$$

$$P(D|t_1, t_2) = \frac{P(D|t_1) P(t_2|D, t_1)}{P(t_2|D, t_1) P(D|t_1) + P(t_2|\bar{D}, t_1) P(\bar{D}|t_1)}$$

$$\left. \begin{aligned} P(t_2|D, t_1) &= P(t_2|D) \\ P(t_2|\bar{D}, t_1) &= P(t_2|\bar{D}) \end{aligned} \right\} \leadsto \text{conditional independence!}$$

with conditional independence,

$$P(D|t_1, t_2) = \frac{0.167(0.99)}{0.167(0.99) + 0.833(0.05)} = 0.80$$

Generalizing Further

Boris Babic

Bayes' Theorem

The Bayesian Approach

For more events

E_1, E_2, \dots, E_K form a partition on S

$$P(E_i | I) = \frac{P(E_i)P(I|E_i)}{\sum_{i=1}^K P(I|E_i)P(E_i)} = P(I)$$

K prior probs: $P(E_1), \dots, P(E_K)$

K likelihoods: $P(I|E_1), \dots, P(I|E_K)$

K posterior probs: $P(E_1|I), \dots, P(E_K|I)$

If you test positive for a certain disease...

You may want to ask

- Do I have the disease or not?
- What is the chance that I have the disease?

Possible answers

- Frequentist: I do not know. You're asking the wrong question. Whether you have the disease or not is not a random variable. It is a fixed value. Therefore, the question does not make sense.
- Bayesian: The chance that you have the disease is ... % (How?)

Differences from frequentist statistics

- On the Bayesian approach, the parameter θ is considered as a random quantity.
- We describe our uncertainty about θ by a probability distribution, referred to as the prior distribution.
- A sample is taken from a population indexed by θ , and the prior is then updated, using Bayes' Rule, to get a posterior distribution for θ given the sample.
- Inferences are then made from the posterior distribution.

The Three Steps of the Bayesian Approach

Boris Babic

Bayes'
Theorem

The
Bayesian
Approach

- Set up a full probability model
 - A joint probability distribution for all observable and unobservable quantities in a problem.
 - The model should be consistent with knowledge about the underlying scientific problem and the data collection process
- Condition on the observed data
 - Calculate and interpret the appropriate posterior distribution
 - We are interested in the conditional probability distribution of the unobserved quantities of interest given the observed data
 - We are often interested in the marginal distribution of a subset of unobserved quantities
- Evaluate the model
 - Does the model fit the data?
 - Are substantive conclusions reasonable?
 - How sensitive are results to model assumptions?

Random Variables and the probability density function

Boris Babic

Bayes' Theorem

The Bayesian Approach

- In Bayesian inference a random variable is defined as an unknown numerical quantity about which we make probability statements.
- Quantitative outcomes of experiments are random variables.
- And fixed but unknown parameters are also random variables.
- Recall that $F(y) = \Pr(Y \leq y)$ is called the cumulative distribution function..
- If F is continuous we say that Y is a continuous random variable.
- A theorem from mathematics says that for every continuous cdf F there exists a positive function $\pi(y)$ such that

$$F(a) = \int_{-\infty}^a \pi(y) d(y)$$

- $\pi(y)$ is the probability density function of Y
- Its essential properties are that it is positive, and it integrates to 1.
- Due to the linearity of integration it is also additive.
- Integration for continuous distributions behaves similarly to summation for discrete distributions.
- In fact, integration can be thought of as a generalization of summation for situations in which the sample space is not countable.
- $p(y)$ is not the probability that $Y = y$.

- $X = x, Y = z, Z = z$ for observed data/ observable random variables.
- θ unobservable (vector) of quantities/population parameters. Eg: $\theta' = (\mu, \sigma^2)'$.
- \tilde{y} unknown, potentially observable data (e.g., the outcome for the next patient enrolled in a trial)
- $\pi(\cdot)$ is a probability density/mass function. Sometimes we use $\Pr(\cdot)$.
- $\pi(\cdot|\cdot)$ denotes a conditional density/mass function.
- $f(x|\theta)$ is the sampling distribution of X . Equivalently: $\pi(x|\theta)$
- $\theta \sim f(\varphi)$ means that the distribution of θ is given by $f(\varphi)$ (and hence φ is a (vector) of hyperparameter(s)).
- Parameter space: Θ , sample space: \mathcal{Y} .

- Prior distribution for θ :

$$\theta \sim \pi(\theta)$$

- Sample distribution (or likelihood) of \mathbf{X} given θ :

$$\mathbf{X}|\theta \sim f(\mathbf{x}|\theta) = \pi(\mathbf{x}|\theta)$$

- Joint distribution of \mathbf{X} and θ (this is our full model):

$$f(\mathbf{x}, \theta) = f(\mathbf{x}|\theta)\pi(\theta)$$

- Recall, chain rule of probability: $p(E_1 \cap E_2) = p(E_2|E_1)p(E_1)$

- Marginal distribution of \mathbf{X} :

$$m(\mathbf{x}) = \int_{\theta \in \Omega} f(\mathbf{x}, \theta) d\theta = \int_{\theta \in \Omega} f(\mathbf{x}|\theta)\pi(\theta) d\theta$$

- Posterior distribution of θ (conditional distribution of θ given \mathbf{X}):

$$\pi(\theta|\mathbf{x}) = \frac{f(\mathbf{x}, \theta)}{m(\mathbf{x})} = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{m(\mathbf{x})} \propto f(\mathbf{x}|\theta)\pi(\theta) \quad (\text{Bayes' Rule})$$

- At the beginning of class 1A, we assumed the data generating process is *iid*. This is a typical assumption in classical inference.
- In Bayesian inference, we will assume something weaker, *exchangeability*.
- For example, a sequence of coin tosses is exchangeable if $\pi(x_1, \dots, x_n) = \pi(x_{\sigma(1)}, \dots, x_{\sigma(n)})$ for every permutation σ of the order.
- Position and order is irrelevant, for any length of the sequence.
- Exchangeability is an assumption about the underlying symmetry in the inference problem.
- For example: compare the probability of observing (H, H, T) with the probability of observing (H, T, H) .
- Independent (Bernoulli) trials with x successes and $n - x$ failures are exchangeable: for any length, the probability is proportional to $\theta^x (1 - \theta)^{n-x}$
- Can you think of an exchangeable sequence that is not iid?

- Let $Y_i \in \mathcal{Y}$ for all $i \in \{1, 2, \dots\}$. Suppose that, for any n , our belief model for Y_1, \dots, Y_n is exchangeable:

$$\pi(y_1, \dots, y_n) = \pi(y_{p1} \dots y_{pn})$$

for all permutations p of $\{1, \dots, n\}$.

- Then our model can be written as:

$$\pi(y_1, \dots, y_n) = \int \left[\prod_{i=1}^n \pi(y_i | \theta) \right] \pi(\theta) d\theta$$

- This θ , whose existence is guaranteed for exchangeable sequences, can be interpreted as the Bayesian prior. But notice that it is not imposed into the problem! Its existence follows from a symmetry assumption weaker than *iid*.

Sketch proof of exchangeability for the Bernoulli process

Boris Babic

Bayes' Theorem

The Bayesian Approach

- Let $p_{k,n}$ denote $P(X_1 = 1, \dots, X_k = 1, X_{k+1} = 0, \dots, X_n = 0)$ where X_1, \dots, X_n is an exchangeable sequence of Bernoulli random variables.

- Let $q_r = P(\sum_{i=1}^m X_i = r)$

- Then,

$$p_{k,n} = \sum_{r=0}^m \frac{(r)_k (m-r)_{n-k}}{(m)_n}$$

where $(x)_k = \prod_{j=0}^{k-1} (x-j)$

- From exchangeability, it follows that given r ones, the distribution of X_1, \dots, X_m is the same as that obtained by drawing from an urn containing r ones and $m-r$ zeros.

- Thus, the r th term of the series is

$$P\left[X_1 = 1, \dots, X_r = 1, X_{r+1} = 0, \dots, X_m = 0 \mid \sum_{j=1}^m X_j = r\right] \times P\left[\sum_{j=1}^m X_j = r\right]$$

- So we can rewrite the first eq. as

$$p_{k,n} = \int_0^1 \frac{(\theta m)_k ((1-\theta)m)_{n-k}}{(m)_n} F_m(d\theta)$$

- where F_m is the distribution function concentrated on $\{r/m : 0 \leq r \leq m\}$ whose jump at r/m is q_r .