

CS 30: Discrete Math in CS (Winter 2019): Lecture 23

Date: 18th February, 2019 (Monday)

Topic: Probability: Conditional Independence and Bayes Rule

Disclaimer: These notes have not gone through scrutiny and in all probability contain errors.

Please discuss in Piazza/email errors to deeparnab@dartmouth.edu

1. Some Recap.

- $\Pr[\mathcal{A} | \mathcal{B}] = \frac{\Pr[\mathcal{A} \cap \mathcal{B}]}{\Pr[\mathcal{B}]}$.
- \mathcal{A} and \mathcal{B} are *independent* events if $\Pr[\mathcal{A} \cap \mathcal{B}] = \Pr[\mathcal{A}] \cdot \Pr[\mathcal{B}]$.
- If $\Pr[\mathcal{B}] \neq 0$, then an equivalent way of stating independence is $\Pr[\mathcal{A} | \mathcal{B}] = \Pr[\mathcal{A}]$.

2. **Mutual and Pairwise Independence.** Say $\mathcal{A}, \mathcal{B}, \mathcal{C}$ are three events. These are said to be *pairwise independent*, if any pair of them are independent events. That is,

$$\Pr[\mathcal{A} \cap \mathcal{B}] = \Pr[\mathcal{A}] \cdot \Pr[\mathcal{B}], \quad \Pr[\mathcal{B} \cap \mathcal{C}] = \Pr[\mathcal{B}] \cdot \Pr[\mathcal{C}], \quad \Pr[\mathcal{C} \cap \mathcal{A}] = \Pr[\mathcal{C}] \cdot \Pr[\mathcal{A}]$$

These three events are *mutually independent* if these are pairwise independent *and*

$$\Pr[\mathcal{A} \cap \mathcal{B} \cap \mathcal{C}] = \Pr[\mathcal{A}] \cdot \Pr[\mathcal{B}] \cdot \Pr[\mathcal{C}]$$

Exercise: Describe three events which are pairwise independent but not mutually independent.

More generally, given events A_1, A_2, \dots, A_n are mutually independent, if for any two *disjoint* subsets S and T of $\{1, 2, \dots, n\}$, the event $\bigcap_{i \in S} A_i$, that is the intersection of all events indexed by S , and the event $\bigcap_{j \in T} A_j$ are independent.

Exercise: Describe three events $\mathcal{A}, \mathcal{B}, \mathcal{C}$ such that $\Pr[\mathcal{A} \cap \mathcal{B} \cap \mathcal{C}] = \Pr[\mathcal{A}] \cdot \Pr[\mathcal{B}] \cdot \Pr[\mathcal{C}]$ but they are not pairwise independent (and therefore not mutually independent).

3. Conditional Independence.

Consider the following two events. There lies in front of you a *fair* coin. Alice tosses it. Then Bob tosses the same coin. Let \mathcal{A} be the event that Alice gets heads. Let \mathcal{B} be the event that Bob gets heads. Are these independent? Even before doing the calculation, you would say sure. Alice's toss shouldn't hinder Bob's toss. Indeed, both $\Pr[\mathcal{A}] = \Pr[\mathcal{B}] = 1/2$ and $\Pr[\mathcal{A} \cap \mathcal{B}] = 1/4$. These are independent.

Exercise: Check that \mathcal{A} and \mathcal{B} are independent even when the coin is not fair, but instead it came heads all the time, or came heads 90% of the time.

Now consider a slightly different experiment. In a box lies two coins. One is fair. The other is biased; it comes heads with probability 0.75. You pick up a coin from these two at random and place it in front of you. Alice tosses it. Bob tosses the same coin. \mathcal{A} and \mathcal{B} are same as above. Are these independent events?

To see (or at least get a hunch – in Math you should always have a hunch) that they are not before doing any calculations, take the experiment to an extreme. Suppose both the coins in the box were super un-fair; suppose one of them came tails all the time, and the other came heads all the time. Then note, if \mathcal{A} occurs, then \mathcal{B} occurs with 100% probability (if Alice sees a head, then she has for sure picked the all-heads coin, and so Bob will for sure see a heads as he is tossing the same coin). On the other hand, \mathcal{B} is not a sure-shot; if I had picked the all-tails coin, then \mathcal{B} doesn't occur. Thus, \mathcal{A} and \mathcal{B} aren't independent.

However, there is a *third* random event here. It is the event \mathcal{E} which is whether I pick the fair coin or not. I claim that \mathcal{A} and \mathcal{B} are independent *if we condition on \mathcal{E}* . That is, I claim

$$\Pr[\mathcal{A} \cap \mathcal{B} \mid \mathcal{E}] = \Pr[\mathcal{A} \mid \mathcal{E}] \cdot \Pr[\mathcal{B} \mid \mathcal{E}]$$

Indeed, if I tell you that \mathcal{E} has occurred, then the problem becomes the one asked before; given a fair coin tossed by Alice and Bob, the events that they see heads is independent.

Remark:

Conditional Independence is a tricky concept. Be wary. Here are a couple of plausible potholes.

- \mathcal{A} and \mathcal{B} are independent events. Then they are also *conditionally* independent on any event \mathcal{E} . **False.** *Example: Roll two fair dice. \mathcal{A} is the event that the first dice is odd. \mathcal{B} is the event that the second dice is odd. These are independent events. Now consider the event \mathcal{E} that the sum of the two dice is odd.. What is $\Pr[\mathcal{A} \mid \mathcal{E}]$? You can now calculate this – it is 1/2 as well. Similarly, $\Pr[\mathcal{B} \mid \mathcal{E}] = 1/2$. However, what is $\Pr[\mathcal{A} \cap \mathcal{B} \mid \mathcal{E}]$? Yep, it's zero. **Independence can be lost upon conditioning.***

- \mathcal{A} and \mathcal{B} are conditionally independent given \mathcal{E} . Then they are conditionally independent given $\neg\mathcal{E}$ as well. **False.** *In its generality this is false, although in the above example of coins, it is true. To see why it is false, we can consider again the setting of rolling two dice. However, this time \mathcal{A} occurs if the first die lands 1, and \mathcal{B} occurs if the second die lands 1. \mathcal{E} is the event that the sum is 2; $\neg\mathcal{E}$ is the event that the sum is not 2.*

*Note: $\Pr[\mathcal{A} \mid \mathcal{E}] = \Pr[\mathcal{B} \mid \mathcal{E}] = \Pr[\mathcal{A} \cap \mathcal{B} \mid \mathcal{E}] = 1$. Thus, \mathcal{A} and \mathcal{B} are conditionally independent given \mathcal{E} . On the other hand, $\Pr[\mathcal{A} \mid \neg\mathcal{E}]$ is something non-zero (figure out what it is!), and $\Pr[\mathcal{B} \mid \neg\mathcal{E}]$ is something non-zero. But, $\Pr[\mathcal{A} \cap \mathcal{B} \mid \neg\mathcal{E}]$ is certainly zero. **Conditional Independence can be lost upon the negation of the event we are complementing on.***

4. **Bayes Rule.** Put simply, Bayes' rule is the following observation: for any two events \mathcal{A} and \mathcal{B} which each occur with non-zero probability

$$\Pr[\mathcal{B} \mid \mathcal{A}] = \frac{\Pr[\mathcal{A} \mid \mathcal{B}] \cdot \Pr[\mathcal{B}]}{\Pr[\mathcal{A}]} \quad \text{(Bayes Rule)}$$

The proof is trivial after we substitute the formula of conditional probability.

We can expand it slightly more using the law of total probability to get

$$\Pr[\mathcal{B} | \mathcal{A}] = \frac{\Pr[\mathcal{A} | \mathcal{B}] \cdot \Pr[\mathcal{B}]}{\Pr[\mathcal{A} | \mathcal{B}] \cdot \Pr[\mathcal{B}] + \Pr[\mathcal{A} | \neg\mathcal{B}] \cdot (1 - \Pr[\mathcal{B}])} \quad (\text{Bayes Rule - Opened up})$$

Why is this a big deal? We will look at three examples. But in a nutshell, it states that to answer what is the probability of event \mathcal{B} given event \mathcal{A} , if we know (a) the total probability of event \mathcal{B} , and if (b) the probability of event \mathcal{A} is *easier* to figure out, then we can get our answer. The main applications come in when \mathcal{A} is an “outcome” and \mathcal{B} is a “hypothesis”; $\Pr[\mathcal{B}]$ is a “prior belief” on the hypothesis, and $\Pr[\mathcal{B} | \mathcal{A}]$ is our “posterior belief” given we see the outcome event \mathcal{A} .

Examples.

- *Arithmophobia is a quality of life debilitating condition and should be detected as early as possible. Fortunately, the pharmaceutical company HAYSTEAM have come up with a test. It's not perfect. It has a **false positive rate** of $fp = 1\%$; that is, on 1% of the healthy population, the test detects the condition, It also has a **false negative rate** of $fn = 2\%$. Again, this means that 2% of the afflicted population go undetected. It is assumed around 40% of the population may be suffering from Arithmophobia.*

You take the test and unfortunately it comes positive (the test says you have Arithmophobia). What is the probability that you actually do?

What a story! But such situations abound. Easy-peasy if you know Bayes rule and know how to set things up.

\mathcal{A} be the event that you have the affliction. Now, you don't know whether you do or not (that's why, presumably, you take the test). Before taking the test, you just look at the statistics and believe that you are as likely as anyone else to have this condition. Since 40% of the population have it, you conclude (reasonably)

$$\Pr[\mathcal{A}] =: p_A = 0.4$$

\mathcal{P} be the event that the test comes out positive on you. We are really interested in figuring out $\Pr[\mathcal{A} | \mathcal{P}]$. We will do so by applying Bayes rule.

First, is $\Pr[\mathcal{P} | \mathcal{A}]$ easy? That is, if you did have the affliction, what is the probability that the test would catch it? The answer is $(1 - fn)$; you would be tested positive unless we got a *false negative*. Thus,

$$\Pr[\mathcal{P} | \mathcal{A}] = (1 - fn) = 0.98$$

How about $\Pr[\mathcal{P} | \neg\mathcal{A}]$? This is precisely the *false positive* rate. So,

$$\Pr[\mathcal{P} | \neg\mathcal{A}] = fp = 0.01$$

Now to apply Bayes rule,

$$\Pr[\mathcal{A} | \mathcal{P}] = \frac{\Pr[\mathcal{P} | \mathcal{A}] \cdot \Pr[\mathcal{A}]}{\Pr[\mathcal{P} | \mathcal{A}] \cdot \Pr[\mathcal{A}] + \Pr[\mathcal{P} | \neg\mathcal{A}] \cdot (1 - \Pr[\mathcal{A}])}$$

which simplifies to

$$\Pr[\mathcal{A} | \mathcal{P}] = \frac{(1 - \text{fn})p_A}{(1 - \text{fn})p_A + \text{fp}(1 - p_A)} = 0.985$$

Thus, if the test comes positive, then the chances you have the affliction gets close to 98.5%.

It is instructive to repeat the above calculations when $\Pr[\mathcal{A}]$ is small, say $\Pr[\mathcal{A}] = 0.1$, that is, only 10% of the population have Arithmophobia. In that case, we get $\Pr[\mathcal{A} | \mathcal{P}] = \frac{(0.98)(0.1)}{(0.98)(0.1) + (0.01)(0.9)} = 0.915$. After the positive test, my belief that I have the affliction goes from 10% to more than 90%.

If $\Pr[\mathcal{A}] = 1\%$, however, then if you repeat the calculation you get that $\Pr[\mathcal{A} | \mathcal{P}] = 0.497$. Thus, if the affliction is so rare that less than 1% of the people have the disease, then a positive test (with the given rates) shouldn't take your belief to more than a random (fair) coin toss.

More generally, to have "high belief probabilities" the false-negative and false-positive scores should be *substantially* smaller than the (prior) probability of the affliction (hypothesis) itself.

- *Spam Filters. We are trying to train a (Bayesian) Spam Filter. We start with a corpus with 2000 spam messages and 1000 real messages. We observe that the word "Congratulations" appears in 100 spam messages, and 10 real messages. We also observe that the word "Account" appears in 160 spam messages and 20 real messages. Assume you believe that any incoming email is possible spam with probability 40%. What is the probability an incoming message is spam given it contains the word "Congratulations"? What is the probability an incoming message is spam given it contains the word "account"? What is the probability that the incoming message is spam, given it contains **both** words "account" and "congratulations"? If we set a threshold of 90% to mark spam or not, in which of these cases would we mark spam.*

Remark: *I couldn't cover this example in class. It is an interesting and informative example. In fact, before moving ahead, please try it yourself.*

Consider an incoming email. Let \mathcal{S} be the event that it is spam. The assumption we are making is that $\Pr[\mathcal{S}] = 0.4$.

Let \mathcal{A} be the event that the word “account” appears in the email. Let \mathcal{C} be the event that the word “congratulations” appears in the email. From the data, we *deduce* that in a random spam message, the chances of seeing “congratulations” is $\frac{100}{2000} = 0.05$. Thus, we conclude

$$\Pr[\mathcal{C} | \mathcal{S}] = 0.05$$

Similarly, we conclude,

$$\Pr[\mathcal{C} | \neg\mathcal{S}] = \frac{10}{1000} = 0.01$$

since $\neg\mathcal{S}$ implies a ‘real’ message. Also, we conclude

$$\Pr[\mathcal{A} | \mathcal{S}] = \frac{160}{2000} = 0.08$$

and

$$\Pr[\mathcal{A} | \neg\mathcal{S}] = \frac{20}{1000} = 0.02$$

Now, we can apply Bayes rule to get

$$\Pr[\mathcal{S} | \mathcal{A}] = \frac{\Pr[\mathcal{A} | \mathcal{S}] \cdot \Pr[\mathcal{S}]}{\Pr[\mathcal{A} | \mathcal{S}] \cdot \Pr[\mathcal{S}] + \Pr[\mathcal{A} | \neg\mathcal{S}] \cdot \Pr[\neg\mathcal{S}]} = \frac{(0.08) \cdot (0.4)}{(0.08)(0.4) + (0.02)(0.6)}$$

which computes to 0.727. That is, if we see the word “account” in an incoming mail, we would believe the probability it is spam is around 72.7%. Thus, our spam-filter won’t mark it spam.

Similarly, for “congratulations”, we get

$$\Pr[\mathcal{S} | \mathcal{C}] = \frac{\Pr[\mathcal{C} | \mathcal{S}] \cdot \Pr[\mathcal{S}]}{\Pr[\mathcal{C} | \mathcal{S}] \cdot \Pr[\mathcal{S}] + \Pr[\mathcal{C} | \neg\mathcal{S}] \cdot \Pr[\neg\mathcal{S}]} = \frac{(0.05) \cdot (0.4)}{(0.05)(0.4) + (0.01)(0.6)}$$

which computes to around 0.769. That is, if we see the word “congratulations” in an incoming mail, we would believe the probability it is spam is around 77%. The spam-filter won’t mark this spam.

How do we solve the next question – when we see both “congratulations” and “account”. Well, we need to find

$$\Pr[\mathcal{S} | \mathcal{A} \cap \mathcal{C}] = \frac{\Pr[\mathcal{A} \cap \mathcal{C} | \mathcal{S}] \cdot \Pr[\mathcal{S}]}{\Pr[\mathcal{A} \cap \mathcal{C}]} \quad (1)$$

We **don’t know** how to calculate $\Pr[\mathcal{A} \cap \mathcal{C} | \mathcal{S}]$. This is where (another) assumption, called the **Naive Bayes Assumption** is made. In the setting of Spam Filters, it states that the events \mathcal{A} and \mathcal{S} are *conditionally independent* on *both* spam (that is \mathcal{S}) and *real* messages. What it says that it does recognize that the distribution of these words (“congratulations”, “account”) may not behave independently on the whole email corpus; but if we focus our attention to the classes at hand, then it does. Again, this is an *assumption*, which is actually made out there many times in the real world.

$$\Pr[\mathcal{A} \cap \mathcal{C} | \mathcal{S}] = \Pr[\mathcal{A} | \mathcal{S}] \cdot \Pr[\mathcal{C} | \mathcal{S}], \quad \Pr[\mathcal{A} \cap \mathcal{C} | \neg \mathcal{S}] = \Pr[\mathcal{A} | \neg \mathcal{S}] \cdot \Pr[\mathcal{C} | \neg \mathcal{S}]$$

(Naive Bayes)

Once we make it, then our calculations can start again. We get:

$$\Pr[\mathcal{A} \cap \mathcal{C}] = \Pr[\mathcal{S}] \cdot \Pr[\mathcal{A} \cap \mathcal{C} | \mathcal{S}] + \Pr[\neg \mathcal{S}] \cdot \Pr[\mathcal{A} \cap \mathcal{C} | \neg \mathcal{S}]$$

and the RHS, with the Naive Bayes assumption, becomes

$$\Pr[\mathcal{A} \cap \mathcal{C}] = \Pr[\mathcal{S}] \cdot \Pr[\mathcal{A} | \mathcal{S}] \cdot \Pr[\mathcal{C} | \mathcal{S}] + \Pr[\neg \mathcal{S}] \cdot \Pr[\mathcal{A} | \neg \mathcal{S}] \cdot \Pr[\mathcal{C} | \neg \mathcal{S}]$$

Substituting in the Bayes rule formula (1), we get

$$\Pr[\mathcal{S} | \mathcal{A} \cap \mathcal{C}] = \frac{\Pr[\mathcal{A} | \mathcal{S}] \cdot \Pr[\mathcal{C} | \mathcal{S}] \cdot \Pr[\mathcal{S}]}{\Pr[\mathcal{S}] \cdot \Pr[\mathcal{A} | \mathcal{S}] \cdot \Pr[\mathcal{C} | \mathcal{S}] + \Pr[\neg \mathcal{S}] \cdot \Pr[\mathcal{A} | \neg \mathcal{S}] \cdot \Pr[\mathcal{C} | \neg \mathcal{S}]}$$

which evaluates to

$$\Pr[\mathcal{S} | \mathcal{A} \cap \mathcal{C}] = \frac{(0.05)(0.08)(0.4)}{(0.05)(0.08)(0.4) + (0.02)(0.01)(0.6)} = 0.9302$$