

- Bayes rule pops out of basic manipulations of probability distributions.
- Let's reach it through a very simple example.

Example

	B_1		
A_1	x	x	x
		x	
	x		x
		x	x

- Call the entire space Ω
- A_i is the i th row partition
- B_i is the i th column partition
(defined arbitrarily)

- We have points lying in this space, denoted by x 's. pick one of these points uniformly at random.
- In this case, calculating probabilities is simply a matter of counting

$$P(x \in A_1) = \frac{\#A_1}{\#\Omega}, \quad P(x \in B_1) = \frac{\#B_1}{\#\Omega}$$

- What about $P(x \in A_1 | x \in B_1)$?

This is the probability $x \in A_1$ given that I know $x \in B_1$.

This is called a conditional probability.

- Looking at the picture, and doing a clever trick...

$$P(X \in A, | X \in B_1) = \frac{\#(A_1 \cap B_1)}{\#B_1}$$

$$= \frac{\#(A_1 \cap B_1)}{\#B_1} \cdot \frac{\#\Omega}{\#\Omega} \quad \leftarrow \text{trick}$$

$$= \frac{P(X \in A_1 \& X \in B_1)}{P(X \in B_1)}$$

- We're just multiplied and divided by the same thing, but already we've made a general statement.

A more general statement

- Let A and B be two events, then

$$P(A|B) = \frac{P(A, B)}{P(B)} \quad \Rightarrow_{\text{algebra}} \quad P(A|B)P(B) = P(A, B)$$

We have some names for these:

$P(A|B)$: conditional distribution

$P(A, B)$: joint distribution

$P(B)$: marginal distribution

- This last one is tricky since it's also just "the probability of B ." However, we can also express it as follows:

Bayes rule

- So we have that $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$

- These values each have a name:

$$\text{Posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

1. Imagine we don't know A , but we have some prior belief about its value

2. We get some information about A in the form of data B .

3. Bayes rule tells us a mathematically rigorous way to incorporate this information in our belief about A in the form of the posterior distribution.

- In reference to the picture,

$$P(B) = \frac{\#B}{\#\Omega} = \frac{\sum_i \#(A_i \cap B)}{\#\Omega} = \sum_i \frac{\#(A_i \cap B)}{\#\Omega} = \sum_i P(A_i, B)$$

bring inside the summation

- Requirements: $A_i \cap A_j = \emptyset$, $\bigcup_i A_i = \Omega$

Getting to Bayes rule

- We're a few easy steps away

showed that: $P(A, B) = P(A|B)P(B)$ ← therefore the
by symmetry: $P(A, B) = P(B|A)P(A)$ ← RHS are equal

$$\begin{aligned} \text{And so: } P(A|B) &= \frac{P(B|A)P(A)}{P(B)} \\ &= \frac{P(B|A)P(A)}{\sum_i P(A_i, B)} \\ &= \frac{P(B|A)P(A)}{\sum_i P(B|A_i)P(A_i)} \end{aligned}$$

- This is called Bayes rule.

- As you can see, there are a few ways it can be written.

Example (medical test)

- This classic example shows how Bayes rule can lead to counter-intuitive results.

- We have two binary indicators.

$$A_i = \begin{cases} 1 & \text{person } i \text{ has disease} \\ 0 & \text{no disease} \end{cases}$$

← don't get to observe this

$$B_i = \begin{cases} 1 & \text{test for disease is positive} \\ 0 & \text{test is negative} \end{cases}$$

← get to observe this

- Person i tests "positive" for the disease ($B_i=1$)
What is the probability the person has it ($A_i=1$)?

- Mathematical problem: Want $P(A_i=1 | B_i=1)$.
Does Bayes rule help?

$$\begin{aligned} \text{Bayes rule: } P(A=1 | B=1) &= \frac{P(B=1 | A=1)P(A=1)}{P(B=1)} \\ &= \frac{P(B=1 | A=1)P(A=1)}{P(B=1 | A=1)P(A=1) + P(B=1 | A=0)P(A=0)} \end{aligned}$$

- We estimate from historical data that

$$P(B=1 | A=1) = 0.95, \quad P(B=1 | A=0) = 0.05$$

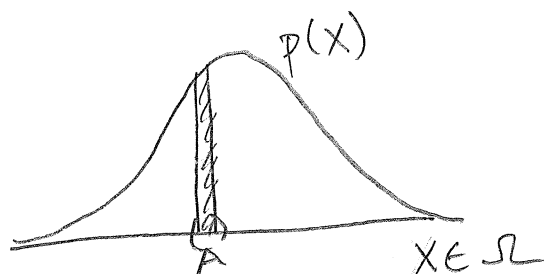
$$P(A=1) = 0.01, \quad P(A=0) = 0.99$$

- Then plugging in, $P(A=1 | B=1) = 0.16$

Continuous space

- We've been talking about discrete distributions so far. The number of values the unknowns can take is finite.
- When values are in a continuous space, we switch to continuous distributions.

Example



observation $x \in \Omega$, $p(x)$ is its density

$$p(x) \geq 0, \int_{\Omega} p(x) dx = 1$$

$$p(x \in A) = \int_A p(x) dx$$

$$Pr(x) = \downarrow p(x) dx = 0 \text{ (theory!)}$$

- Probability theory lets us ignore probabilities and just work with the densities.
- The same rules apply as for discrete random variables.

$$\bullet p(x|\theta) = \frac{p(x, \theta)}{p(\theta)}, \quad p(\theta) = \int p(x, \theta) dx$$

$$\bullet p(x, \theta) = p(x|\theta)p(\theta) = p(\theta|x)p(x)$$

- This leads to Bayes rule for continuous variables

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} = \frac{p(x|\theta)p(\theta)}{\int p(x|\theta)p(\theta) d\theta}$$

- The difference is we're dealing with continuous functions

Bayesian modeling

- Applying Bayes rule to the unknown variables of a data modeling problem is called Bayesian modeling.

- In a simple, generic form we can write the model

$X \sim P(X|\theta)$ \leftarrow data-generating distribution
this is the model of the data

$\theta \sim P(\theta)$ \leftarrow the model prior distribution. Our belief about θ a priori.

- We want to learn θ , so we use Bayes rule

$$P(\theta|X) = \underbrace{P(X|\theta)P(\theta)}_{\text{do know this because we defined it}} / P(X) \leftarrow$$

↓
don't know this yet

do know this because we defined it

do we know this?
can we calculate
 $P(X) = \int P(X|\theta)P(\theta)d\theta$?

- This is the general form of the problems discussed in this class. However, it's non-trivial because

1. $P(X|\theta)$ can be quite complex - model definition
2. $P(\theta)$ can be even more complex
3. $P(\theta|X)$ can be intractable. We can't calculate $P(X)$, so we need an algorithm to approximate this.

Simple example: beta-Bernoulli

- We have a sequence of observations X_1, X_2, \dots, X_N where $X_i = \begin{cases} 1, & \text{"success"} \\ 0, & \text{"failure"} \end{cases}$. Think of them as coin-flips.

- We hypothesize that each X_i is generated by flipping a biased coin, $X_i \sim p(X|\theta) \Rightarrow p(X_i=1|\theta) = \theta$.

- We assume the X_i are independent and identically distributed (i.i.d.) according to $p(X|\theta) \rightarrow X_i \stackrel{iid}{\sim} p(X|\theta)$.

- This means that the X_i are conditionally independent given the bias θ .

$$p(X_1, \dots, X_N | \theta) = \prod_{i=1}^N p(X_i | \theta).$$

- Since $p(X_i | \theta) = \theta^{X_i} (1-\theta)^{1-X_i}$, we can write

$$p(X_1, \dots, X_N | \theta) = \prod_{i=1}^N \theta^{X_i} (1-\theta)^{1-X_i}$$

- We are interested in the posterior distribution of θ given X_1, \dots, X_N , i.e.,

$$\begin{aligned} p(\theta | X_1, \dots, X_N) &= p(X_1, \dots, X_N | \theta) p(\theta) / p(X_1, \dots, X_N) \\ &= \prod_{i=1}^N p(X_i | \theta) p(\theta) / p(X_1, \dots, X_N) \end{aligned}$$

- We've come across our first significant Bayesian problem: what do we set $p(\theta)$ to?

First try

- Let $P(\theta) = \text{Uniform}(0,1) \Rightarrow P(\theta) = 1(0 \leq \theta \leq 1)$

- Then by Bayes rule and i.i.d. assumption

$$\begin{aligned} P(\theta | X_1, \dots, X_N) &= \frac{\prod_{i=1}^N P(X_i | \theta) P(\theta)}{P(X_1, \dots, X_N)} = \\ &= \frac{\theta^{\sum_{i=1}^N X_i} (1-\theta)^{N - \sum_{i=1}^N X_i} 1(0 \leq \theta \leq 1)}{\int_0^1 \theta^{\sum_{i=1}^N X_i} (1-\theta)^{N - \sum_{i=1}^N X_i} 1(0 \leq \theta \leq 1) d\theta} \end{aligned}$$

- This normalizing constant is tricky, but fortunately it's been solved and we can conclude

$$P(\theta | X_1, \dots, X_N) = \frac{\Gamma(N)}{\Gamma(\sum_i X_i) \Gamma(N - \sum_i X_i)} \theta^{\sum_{i=1}^N X_i + 1 - 1} (1-\theta)^{N - \sum_{i=1}^N X_i + 1 - 1}$$

$\Gamma(\cdot)$ is a "gamma function"

- This is a very common distribution called a beta distribution

$$\text{Beta}(a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$$

- In the posterior, $a = 1 + \sum_{i=1}^N X_i$, $b = 1 + N - \sum_{i=1}^N X_i$

- Notice that when $a=b=1$, $\text{Beta}(1,1) = \text{Uniform}(0,1)$ which was our prior.

A "Conjugate" prior

- The beta distribution $\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$ looks a lot like the likelihood term, $\theta^{\sum x_i} (1-\theta)^{N-\sum x_i}$.
- Also, because $\text{Unif}(0,1)$ is a special case, a beta prior would give us more options to express bel. of.
- Beta prior:
"proportional to"
 \downarrow
$$P(\theta | x_1, \dots, x_N) \propto P(x_1, \dots, x_N | \theta) P(\theta)$$
$$\propto \left[\theta^{\sum x_i} (1-\theta)^{N-\sum x_i} \right] \left[\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1} \right]$$
$$\propto \theta^{a+\sum_{i=1}^N x_i - 1} (1-\theta)^{b+N-\sum_{i=1}^N x_i - 1}$$

Comments

1. Where did $\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$ go? We are writing the posterior as a proportionality. The functions

$$g(\theta) \propto f(\theta) \text{ if } g(\theta) = \frac{1}{Z} f(\theta) \text{ for some}$$

constant Z (i.e., Z is not a function of θ) but may be function of other things

2. In general, if $g(\theta) = \frac{f(\theta)}{\int f(\theta) d\theta}$, then $g(\theta) = \frac{\frac{1}{Z} f(\theta)}{\int \frac{1}{Z} f(\theta) d\theta}$

In both cases we can write $g(\theta) \propto f(\theta) \propto \frac{1}{Z} f(\theta)$

3. In general, we can multiply the joint likelihood by any constant to make it look nicer. In this case we chose $\frac{T'(a)T'(b)}{T'(a+b)}$

- So we have $p(\theta | x_1, \dots, x_N) \propto \theta^{a + \sum_{i=1}^N x_i - 1} (1-\theta)^{b + N - \sum_{i=1}^N x_i - 1}$

Solution:
$$p(\theta | x_1, \dots, x_N) = \frac{\theta^{a + \sum_{i=1}^N x_i - 1} (1-\theta)^{b + N - \sum_{i=1}^N x_i - 1}}{\int_0^1 \theta^{a + \sum_{i=1}^N x_i - 1} (1-\theta)^{b + N - \sum_{i=1}^N x_i - 1} d\theta}$$

Trick: Notice $\theta^{a + \sum_{i=1}^N x_i - 1} (1-\theta)^{b + N - \sum_{i=1}^N x_i - 1} \propto \text{Beta}(a', b')$
 where $a' = a + \sum_{i=1}^N x_i$, $b' = b + N - \sum_{i=1}^N x_i$

- The posterior distribution is in the same family as the prior (beta). We just update its parameters.

- Conjugate priors: Let $X \sim p(X|\theta)$ and $\theta \sim p(\theta)$. If the posterior $p(\theta|X)$ is in the same family as the prior, $p(\theta)$, then $p(\theta)$ is conjugate to the likelihood $p(X|\theta)$.

~~Example: gamma prior~~

What do we gain by being Bayesian?

- Consider the expectation and variance under the posterior

$$E[\theta] = \int_0^1 \theta p(\theta | \vec{X}) d\theta = \frac{a + \sum_{i=1}^N x_i}{a + b + N}$$

$$\text{Var}(\theta) = \int_0^1 (\theta - E[\theta])^2 p(\theta | \vec{X}) d\theta = \frac{(a + \sum_{i=1}^N x_i)(b + N - \sum_{i=1}^N x_i)}{(a + b + N)^2 (a + b + N - 1)}$$

- As N increases,

1. $E(\theta) \rightarrow$ empirical success rate

2. $\text{Var}(\theta) \rightarrow 0$ at the rate $1/N$

- Compare with maximum likelihood

$$\theta_{ML} = \arg \max_{\theta} P(x_1, \dots, x_N | \theta) = \frac{1}{N} \sum_{i=1}^N x_i$$

• The Bayesian approach is capturing our uncertainty about the quantity we are interested in.

• Maximum Likelihood doesn't do this

• There are more compelling reasons in the case of data modeling for machine learning applications that we will discuss throughout the semester.