

Bayesian Inference

STA 427/527, Fall 2019, Xin Wang

Probability

- An **experiment** is the process by which an observation/outcome is made.
- An **outcome** of an experiment is any possible observation of that experiment. Outcomes are often called sample points.
- The **sample space** of an experiment is the set consisting of all possible sample points. (S)
- An **event** is a set of outcomes of an experiment, or a subset of S . That is a set of sample points.

ex) flip a coin \rightarrow outcome: heads, tails

$$S = \{ \text{Heads, tails} \}$$

toss a die \rightarrow outcome: 1, 2, 3, 4, 5, 6

$$S = \{ 1, 2, 3, 4, 5, 6 \}$$

\rightarrow Event \rightarrow rolling even # $A = \{ 2, 4, 6 \}$

- **Probability:** Suppose S is a sample space, A is a subset of S . A probability measure $P(\cdot)$ is a function that maps events in S to real numbers.

\rightarrow properties of a probability function

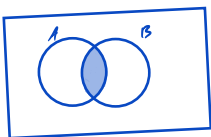
1) for any event A , $P(A) \geq 0$

2) $P(S) = 1$

3) For any countable collection A_1, A_2, \dots (finite sequence of events) of pairwise mutually exclusive events ($A_i \cap A_j = \emptyset$ if $i \neq j$)

$$P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots$$

- Consider $P(A \cap B)$



\rightarrow Different ways to write this

$$\begin{aligned} \rightarrow P(A \cap B) &= P(A|B) P(B) \\ &\downarrow \\ &= P(B|A) P(A) \end{aligned}$$

\rightarrow Conditional probability

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

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

- Bayes' theorem

$$\rightarrow P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \rightarrow P(B|A) = \dots$$

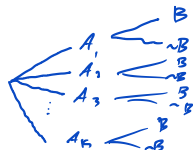
$$\downarrow = \frac{P(B|A) P(A)}{P(B)}$$

Conditional probability is very important in Bayesian statistics \Rightarrow need to know the relationships between joint, marginal, + conditional dists

\rightarrow Extension $\rightarrow B = A_1 \cup A_2 \cup \dots \cup A_k \quad A_i \cap A_j = \emptyset \quad i \neq j$

$$\rightarrow P(A_i|B) = \frac{P(B|A_i) P(A_i)}{P(B)}$$

$$\downarrow = \frac{P(B|A_i) P(A_i)}{\sum P(B|A_i) P(A_i)}$$



- **Example:** The prevalence of heart disease in a certain population is 10%. A screening test for heart disease has 99% sensitivity and 90% specificity, where sensitivity measures the probability of positive that is correctly identified as such and specificity measures the probability of negative that is correctly identified as such. Suppose everyone in the population is given the screening test. What is the probability that one individual with positive test result actually has heart disease?

\rightarrow Define events $\rightarrow A =$ has heart disease
 $B =$ positive test result

$$\rightarrow P(A|B) = ?$$

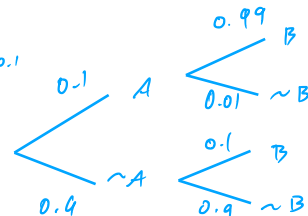
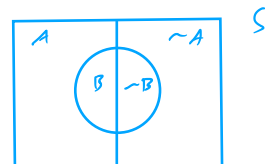
\rightarrow start with what we know

$$\rightarrow P(A) = 0.1 \Rightarrow P(\sim A) = 0.9$$

$$P(B|A) = 0.99$$

$$P(\sim B|\sim A) = 0.9$$

$$\begin{aligned} P(A|B) &= \frac{P(A \cap B)}{P(B)} \\ &= \frac{P(B|A) \cdot P(A)}{P(B|A) P(A) + P(\sim B|\sim A) \cdot P(\sim A)} \\ &\quad \hookrightarrow 1 - P(\sim B|\sim A) = 0.1 \\ &= \frac{0.99 \times 0.1}{0.99 \times 0.1 + 0.1 (0.9)} \\ &= 0.524 \end{aligned}$$



Frequentist vs Bayesian

Scenario → • Suppose θ is the unknown parameter, and $\mathbf{y} = (y_1, y_2, \dots, y_n)$ are observations.

• Frequentist

- ① parameter θ is an unknown fixed value
- ② use \mathbf{y} to estimate θ , denoted $\hat{\theta}$
- ③ $\hat{\theta}$ is a random variable, $\hat{\theta}(\mathbf{y})$
→ different $\hat{\theta}$'s for different sets of \mathbf{y}
- ④ The inference is based on the assumption that the data are "repeatable"

→ Example → $Y_i \sim N(\mu, \sigma^2)$, $i=1, \dots, n$, σ^2 is unknown

observations y_1, y_2, \dots, y_n

Construct a confidence interval for μ 100(1- α)% ($\alpha=0.05$)

$$\bar{y} \pm t_{n-1, \alpha/2} \frac{s}{\sqrt{n}} \Rightarrow \left[\bar{y} - t_{n-1, \alpha/2} \frac{s}{\sqrt{n}}, \bar{y} + t_{n-1, \alpha/2} \frac{s}{\sqrt{n}} \right]$$

95% CI deals w/ foundation of inference
(i.e. how we collected data + estimate $\hat{\theta}$)

• Bayesian

- ① θ is a random variable (prior)

→ This is the biggest difference in assumptions to the frequentist approach

- ② How can we make inference about θ ?

→ In order to make inference on θ , we need to have a posterior distribution

→ First we have a general idea of the shape of the distribution of θ (before collecting data)

→ This is called a prior

→ Then we get observations $\mathbf{y} = (y_1, \dots, y_n)^T$

→ Then we update the prior w/ the new data by finding the distribution of θ conditional on the given dataset \mathbf{y}
(i.e. the posterior)

– Steps in Bayesian

1. Data distribution (likelihood)

→ specify a model $p(\mathbf{y}|\theta)$

→ Likelihood example → $Y_i \sim N(\mu, \sigma^2)$, σ^2 is known → (y_1, \dots, y_n)

$$\rightarrow L(\theta|\mathbf{y}) = \prod_{i=1}^n p(y_i|\theta)$$

↓
= $p(\mathbf{y}|\theta)$ → data distribution (ie our observations) are based on an unknown parameter

2. Prior: → prior knowledge of θ → $\pi(\theta)$

3. Posterior $\star \rightarrow p(\theta|\mathbf{y}) = \frac{p(\mathbf{y}|\theta) \pi(\theta)}{p(\mathbf{y})}$

$$\rightarrow \text{marginal of data} \rightarrow p(\mathbf{y}) = \int_{\theta} p(\mathbf{y}, \theta) d\theta = \int_{\theta} p(\mathbf{y}|\theta) \pi(\theta) d\theta$$

"omega" = parameter space

$$p(\mathbf{A} \cap \mathbf{B}) = p(\mathbf{B}|\mathbf{A}) p(\mathbf{A})$$

$$\rightarrow \text{in practice} \rightarrow \boxed{p(\theta|\mathbf{y}) \propto p(\mathbf{y}|\theta) \pi(\theta)}$$

↓
proportional to

Bayes theorem

$$p(\mathbf{A}|\mathbf{B}) = \frac{p(\mathbf{A}|\mathbf{B}) p(\mathbf{A})}{p(\mathbf{B})}$$

↓ ↓
 $\theta \mathbf{y}$

$$p(\theta|\mathbf{y}) = \frac{p(\mathbf{y}|\theta) \pi(\theta)}{p(\mathbf{y})} \rightarrow \text{prior}$$

(normalizing constant)

- **Example:** A novelty coin company produces coins with varying levels of bias. For an individual coin, the probability of spinning "heads" is θ , where θ is drawn from a Uniform(0, 1) distribution. Consider the following thought experiment. Suppose it were possible to spin each coin in the population 100 times. Let Y represent the number of "heads" resulting from a coin's 100 spins.

1. Find the distribution of Y . That is, what is the distribution of the numbers of "heads" observed in 100 spins?

→ setup → flip a coin 100 times

$Y = \# \text{ heads}$

$p(Y) = ?$ → This is often difficult to find & computationally expensive if using numerical methods

$$\rightarrow p(Y) = \int_0^1 p(Y|\theta) \pi(\theta) d\theta$$

→ [data dist] → $Y|\theta \sim \text{Binomial}(100, \theta)$ → choose p.d.f.s based on situation or shape of data after collection

$$p(Y|\theta) = \binom{100}{Y} \theta^Y (1-\theta)^{100-Y}, \quad Y = 0, 1, \dots, 100$$

→ [prior dist] → $\pi(\theta) = 1, \quad 0 < \theta < 1$

→ uniform prior

→ marginal of data → $p(Y) = \int_0^1 \binom{100}{Y} \theta^Y (1-\theta)^{100-Y} d\theta$

$$= \binom{100}{Y} \int_0^1 \theta^Y (1-\theta)^{100-Y} d\theta$$

$$= \downarrow \int_0^1 \theta^{(Y+1)-1} (1-\theta)^{(101-Y)-1} d\theta$$

$$= \binom{100}{Y} \left[\frac{\Gamma(Y+1) \Gamma(101-Y)}{\Gamma(Y+1 + 101-Y)} \right]$$

$$= \downarrow \left[\frac{\Gamma(Y+1) \Gamma(101-Y)}{\Gamma(102)} \right]$$

$$= \frac{100!}{(102-1)!} \left[\frac{Y! (102-Y)!}{101!} \right]$$

$$= \frac{1}{101}, \quad Y = 0, 1, 2, \dots, 100$$

~ Discrete uniform (0, 100)

if $X \sim \text{Beta}(a, b)$ $0 < x < 1$

$$\rightarrow f(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$$

$$\rightarrow 1 = \int_0^1 f(x) dx$$

$$= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1} dx$$

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 x^{a-1} (1-x)^{b-1} dx$$

$$\frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} = \int_0^1 x^{a-1} (1-x)^{b-1} dx$$

$$\Gamma(x) = (x-1)! \quad \text{if } x \in \mathbb{Z}^+ \text{ positive integer}$$

$$\Gamma(0) = 1$$

2. Now think about the case that exactly 50 "heads" are observed in the coins we observed. Among these coins, determine the distribution of θ , the actual probabilities of "heads".

$$\begin{aligned}
 \rightarrow [\text{posterior dist}] \rightarrow p(\theta|Y) &= \frac{p(Y|\theta) \pi(\theta)}{p(Y)} \\
 &= \frac{\left[\binom{100}{Y} \theta^Y (1-\theta)^{100-Y} \right] \cdot 1}{1/101} \\
 &= 101 \left(\frac{100!}{(100-Y)! Y!} \right) \theta^Y (1-\theta)^{100-Y} \\
 &= \frac{(101)!}{(100-Y)! Y!} \theta^Y (1-\theta)^{100-Y} \\
 &= \frac{\Gamma(102)}{\Gamma(101-Y) \Gamma(Y+1)} \theta^{(Y+1)-1} (1-\theta)^{(101-Y)-1} \\
 &\sim \text{Beta}(Y+1, 101-Y)
 \end{aligned}$$

showing exact calculation of the posterior
(i.e. not using proportionality)
→ much easier in practice to use proportionality
+ just look for kernel

$$\rightarrow \text{if } Y=50 \rightarrow \theta|Y \sim \text{Beta}(51, 51)$$

$$\begin{aligned}
 \rightarrow [\text{posterior dist}] &\rightarrow p(\theta|Y) \propto p(Y|\theta) \pi(\theta) \\
 &\propto \left[\binom{100}{Y} \theta^Y (1-\theta)^{100-Y} \right] \cdot 1 \\
 &\propto \theta^Y (1-\theta)^{100-Y} \\
 &\propto \theta^{(Y+1)-1} (1-\theta)^{(101-Y)-1} \\
 &\sim \text{Beta}(Y+1, 101-Y)
 \end{aligned}$$

Just a constant
(w/ respect to θ)
⇒ drop it

→ Kernel

$$X \sim \text{Beta}(a, b)$$

$$0 < x < 1, a, b > 0$$

$$\rightarrow f(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$$

Kernel of Beta w/ shape parameters a + b

$$\rightarrow E(x) = \frac{a}{a+b}$$

$$\begin{aligned}
 \rightarrow [\text{posterior mean}] &\rightarrow E(\theta|Y) = \frac{Y+1}{(Y+1)+(101-Y)} \\
 &= \frac{Y+1}{102} \\
 &\downarrow \\
 &\text{Bayesian point estimate}
 \end{aligned}$$

frequentist approach

$$\rightarrow [\text{data dist}] \rightarrow p(Y|\theta) = \binom{100}{Y} \theta^Y (1-\theta)^{100-Y}$$

$$\begin{aligned}
 \rightarrow [\text{log-likelihood}] &\rightarrow \ell(\theta|Y) = \log \left[\binom{100}{Y} \theta^Y (1-\theta)^{100-Y} \right] \\
 &\downarrow \\
 &= \log \left(\binom{100}{Y} \right) + Y \log(\theta) + (100-Y) \log(1-\theta)
 \end{aligned}$$

→ ... → ... derivative + set to zero ... →

$$\rightarrow \text{point estimate} \rightarrow \hat{\theta}_{MLE} = \frac{Y}{100}$$