

# Introduction to Bayesian Statistics (I)

Xiaoquan (William) Wen

July 3, 2019

# AN OVERVIEW

- ▶ Named after Thomas Bayes (1701 - 1761)
- ▶ What is Bayesian statistics
  - ▶ a mathematical procedure that applies probabilities to statistical problems
  - ▶ provides the tools to update people's beliefs in the evidence of new data.
- ▶ Bayesian approach is trending in big data era

# A BRIEF HISTORY OF BAYESIAN STATISTICS

- ▶ 1700s, Bayes' Theorem
- ▶ 1800s, Pierre-Simon Laplace formalized and popularized Bayesian inference
- ▶ 1940s, Alan Turing's Bayesian system decoded German Enigma Machine, but in general Bayesianism is considered in decline
- ▶ 1960s, revival of the Bayes' theorem: theory and computation work
- ▶ Current day practice:
  - ▶ Election prediction (FiveThirtyEight.com)
  - ▶ Aviation incidents investigations
  - ▶ Broadly used in medicine, economy and all branches of sciences

# CONDITIONAL PROBABILITY

$$\Pr(A \mid B)$$

- ▶ Probability of event A given event B has occurred
- ▶  $\Pr(A \mid B) = \frac{\Pr(A \cap B)}{\Pr(B)}$
- ▶ Fundamental in probability theory and statistics

# BAYES' THEOREM

$$\begin{aligned}\Pr(\theta \mid \text{data}) &= \frac{\Pr(X) \Pr(\text{data} \mid X)}{\Pr(\text{data})} \\ &= \frac{\Pr(X) \Pr(\text{data} \mid X)}{\Pr(\text{data})} \\ &= \frac{\Pr(X) \Pr(\text{data} \mid X)}{\Pr(X) \Pr(\text{data} \mid X) + \Pr(X^c) \Pr(\text{data} \mid X^c)}\end{aligned}$$

- ▶  $\Pr(X)$ : prior probability/distribution
- ▶  $\Pr(\text{data} \mid X)$ : likelihood
- ▶  $\Pr(X \mid \text{data})$ : posterior probability/distribution

## APPLICATION OF BAYES THEOREM

If 1% of a population have a specific form of cancer, for a screening test with 80% sensitivity and 95% specificity, **What is the chance that a patient has the cancer if he tests positive?**

## APPLICATION OF BAYES THEOREM

If 1% of a population have a specific form of cancer, for a screening test with 80% sensitivity and 95% specificity, **What is the chance that a patient has the cancer if he tests positive?**

- ▶ Sensitivity:  $\Pr(\text{test} + \mid \text{cancer}) = 80\%$
- ▶ Specificity:  $\Pr(\text{test}- \mid \text{no cancer}) = 95\%$

## APPLICATION OF BAYES THEOREM

If 1% of a population have a specific form of cancer, for a screening test with 80% sensitivity and 95% specificity, **What is the chance that a patient has the cancer if he tests positive?**

- ▶ Sensitivity:  $\Pr(\text{test} + \mid \text{cancer}) = 80\%$
- ▶ Specificity:  $\Pr(\text{test}- \mid \text{no cancer}) = 95\%$

$$\begin{aligned}\Pr(\text{cancer} \mid \text{test} +) &= \frac{\Pr(\text{cancer}) \Pr(\text{test} + \mid \text{cancer})}{\Pr(\text{test} +)} \\ &= \frac{0.01 \times 0.80}{0.01 \times 0.80 + 0.99 \times 0.05} \\ &\approx 13.9\%\end{aligned}$$



## APPLICATION OF BAYES THEOREM (CONT'D)

- ▶ Most positive tests ( $\approx 86\%$ ) are actually false alarms
  
- ▶ But is the prior  $\Pr(\text{cancer}) = 0.01$  reasonable to use here?

# THE PROCESS OF BAYESIAN INFERENCE

## The Bayesian Machinery

1. Define a parametric model (prior, likelihood)
2. Apply Bayes Theorem and compute the posterior for the parameters of interest
3. Posterior distributions contain full information of inference result

# THE BAYESIAN PHILOSOPHY

- ▶ The Bayesian inference process is a byproduct of multiple statistical principles
- ▶ They start from different perspectives and all conclude that statistical inference results should be summarized in form of posterior distributions
- ▶ This also leads to different interpretations of probabilities
  - ▶ Bayesian: probability is simply a quantification of uncertainty
  - ▶ Frequentist: probability reflects a long-run frequency

# ARGUMENT 1: COHERENCE OF DECISION MAKING

- ▶ Need principled approach to make decision accounting for uncertainty
- ▶ Consider make a prediction,  $\delta(x)$ , with respect to an unknown parameter  $\theta$  based on observed data  $x$
- ▶ Coherent decision should be informed by the posterior distribution  $p(\theta | x)$
- ▶ Inevitably, it requires a prior distribution and apply Bayes theorem

## ARGUMENT 2: THE LIKELIHOOD PRINCIPLE

- ▶ Sufficiency Principle (S): irrelevance of observations independent of a sufficient statistic
- ▶ Conditionality Principle (C): irrelevance of (component) experiments not actually performed
  - ▶ The voltmeter story: [https://en.wikipedia.org/wiki/Likelihood\\_principle](https://en.wikipedia.org/wiki/Likelihood_principle)
- ▶ Likelihood Principle (L): irrelevance of outcomes not actually observed

## THE LIKELIHOOD PRINCIPLE (CONT'D)

- ▶ (almost) All statisticians accept S and C
- ▶ It has been shown (Birnbaum, 1962) that

$$S + C \rightarrow L$$

i.e., all data scientists should accept L

- ▶ Bayesian inference process follows the likelihood principle
- ▶ Some commonly used frequentist procedures, e.g., p-values, confidence intervals, violate the likelihood principle

## ARGUMENT 3: EXCHANGEABILITY

Bayesian inference provides more **flexible** and **realistic** modeling options

Consider tossing a coin with a sequence of outcomes :  
 $X_1, X_2, \dots$

- ▶ The random sequence is often modeled as *independent identically distributed (i.i.d)*
- ▶ Are the sequence of outcomes really independent?  
Note that, independence implies

$$\Pr(X_1, X_2, \dots, X_p) = \prod_{i=1}^p \Pr(X_i)$$

$$\Pr(X_p \mid X_1, X_2, \dots, X_{p-1}) = \Pr(X_p)$$

## ARGUMENT 3: EXCHANGEABILITY

Bayesian inference provides more **flexible** and **realistic** modeling options

Consider tossing a coin with a sequence of outcomes :  
 $X_1, X_2, \dots$

- ▶ The random sequence is often modeled as *independent identically distributed (i.i.d)*
- ▶ Are the sequence of outcomes really independent?  
Note that, independence implies

$$\Pr(X_1, X_2, \dots, X_p) = \prod_{i=1}^p \Pr(X_i)$$

$$\Pr(X_p \mid X_1, X_2, \dots, X_{p-1}) = \Pr(X_p)$$

- ▶ But the sequence of outcomes share the information on the biasness of the coin!



## EXCHANGEABILITY (CONT'D)

- ▶ A more realistic modeling assumption is to treat the sequence *exchangeable*
- ▶ Mathematically, it means  $\Pr(X_1, \dots, X_p)$  is invariant to the permutations of indexes  $(1, \dots, p)$ .
- ▶ An independent sequence is obviously exchangeable, but an exchangeable sequence does not need to be independent!

## DE FINETTI THEOREM

The de Finetti theorem indicates

$$\Pr(X_1, \dots, X_p) = \int \left[ \prod_{i=1}^p \Pr(X_i | \theta) \right] p(\theta) d\theta,$$

for any exchangeable sequence.

- ▶  $\theta$  represents the biasness of the coin
- ▶ Conditional on  $\theta$ , the sequence is i.i.d

$$\Pr(X_1, \dots, X_p | \theta) = \prod_{i=1}^p \Pr(X_i | \theta)$$

- ▶ Because  $\theta$  is unknown, the probability of the sequence has to be averaged over the uncertainty of  $\theta$  (a prior distribution!)

# MODEL EXCHANGEABILITY

- ▶ Requires a prior distribution
- ▶ Give rise to a *hierarchical model*
- ▶ Hierarchical model as a probabilistic generative model

# SUMMARY

- ▶ What is Bayesian statistics
- ▶ The machinery of Bayesian inference
- ▶ The foundations of Bayesianism

Next time

- ▶ Apply Bayesian principle to build statistical models for data analysis problems