

Lecture 9: Statistical Inference

Ziyu Shao

School of Information Science and Technology
ShanghaiTech University

May 16, 2024

Outline

- 1 History of Mathematical Statistics
- 2 Overview of Statistical Inference
- 3 Our Focus: Bayesian Statistical Inference
- 4 Beta and Gamma Distribution
- 5 Conjugate Prior: A Weapon of Bayesian
- 6 Application Case: Bayesian Ranking

Outline

- 1 History of Mathematical Statistics
- 2 Overview of Statistical Inference
- 3 Our Focus: Bayesian Statistical Inference
- 4 Beta and Gamma Distribution
- 5 Conjugate Prior: A Weapon of Bayesian
- 6 Application Case: Bayesian Ranking

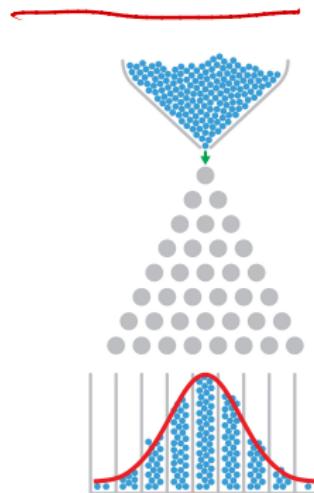
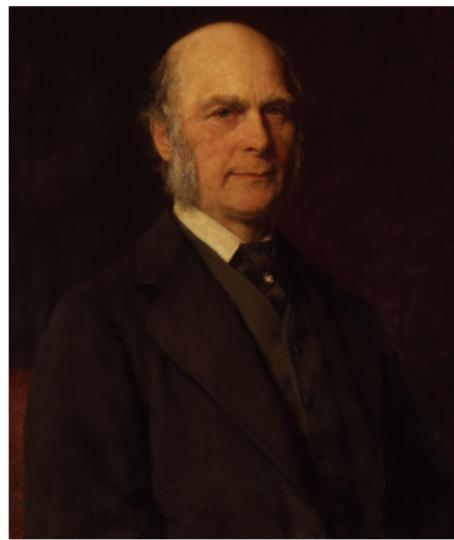
Classical Statistics

- 1800s: Linear Statistical Model and the method of least squares for estimation is often credited to Gauss (1777-1855) (1809), Adrien-Marie Legendre (1752-1833) (1805), Robert Adrain (1775-1843).
- Gauss also showed the optimality of the least-square approach (Gauss-Markov Theorem, 1823).



Classical Statistics

- 1888: Sir Francis Galton proposed the concept of correlation
- 1889: Sir Francis Galton proposed the concept of regression
- 1889: Sir Francis Galton proposed the Galton Board



Classical Statistics

- Karl Pearson (1857-1936) is credited for the establishment of the discipline of statistics. He contributed to theory of linear regression, correlation, Pearson curve, chi-square test, and the method of moments for estimation.



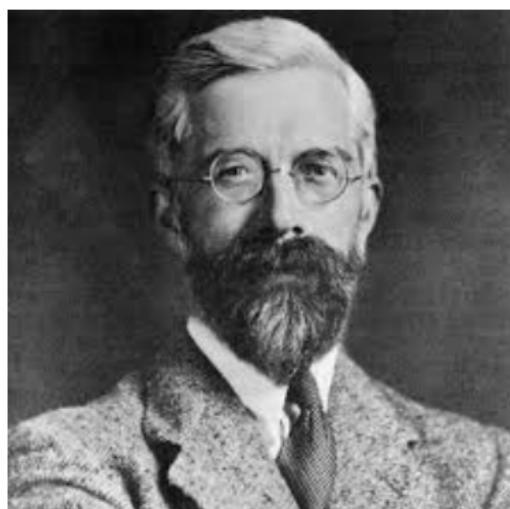
Modern Statistics: Frequency Perspective

- 1908: William Gosset (Student) (1876-1937) proposed Student t-distribution and t-test statistics
- Precursor of small-sample statistics and hypothesis testing.



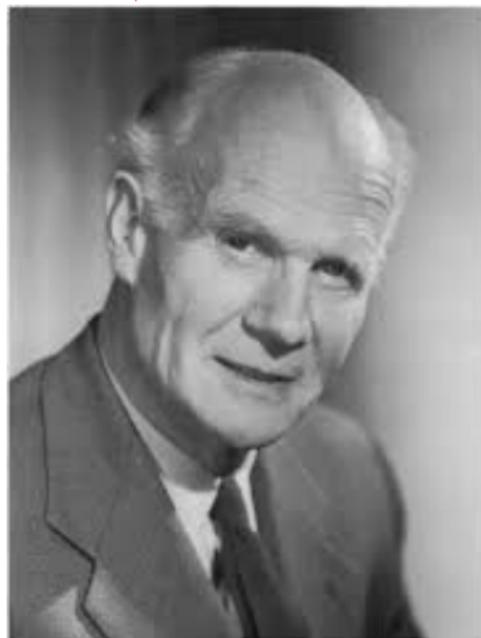
Modern Statistics: Frequency Perspective

- 1912-1922: Sir Ronald Aylmer Fisher (1890-1962) developed the notion of maximum likelihood estimator.
- He also worked on the analysis of variance (ANOVA),
F-distribution, Fisher information and design of experiment.
- Co-founder of Modern Statistics (Mathematical Statistics or Statistical Inference)



Modern Statistics: Frequency Perspective

- Egon Sharpe Pearson (1895-1980): co-founder of Neyman-Pearson Theory for hypothesis testing.



Modern Statistics: Frequency Perspective

- Jerzy Neyman (1894-1981): Co-founder of Modern Statistics (Mathematical Statistics or Statistical Inference)
- 1928-1938: Theoretical foundations of testing hypothesis, point estimation, confidence interval and survey sampling.



Modern Statistics: Frequency Perspective

- 1940s: Pao-Lu Hsu (1910-1970) obtained several exact or asymptotic distributions of important statistics in the theory of multivariate analysis.



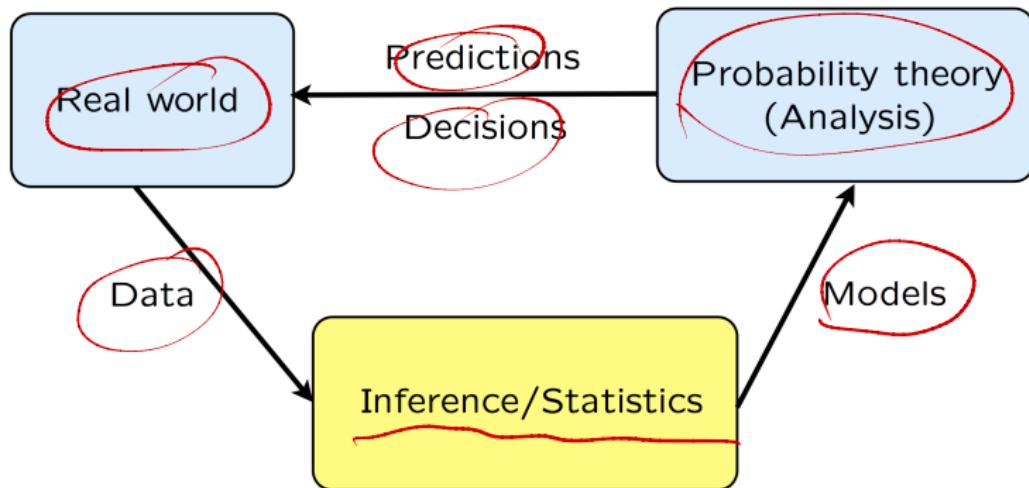
Modern Statistics: Bayesian Perspective

- 1937: Bruno de Finetti proposed a predictive inference approach to statistics, emphasizing the prediction of future observations based on past observations.
- 1939: Harold Jeffreys applied Bayesian analysis for geophysics data.
- 1941-1944: Alan Turing applied Bayesian analysis for breaking the German code (Enigma)
- 1954s: Jimmie Savage proposed Bayesian statistics systematically
- 1950s: Bayesian econometrics originated from Harvard business school prevailed in economics society
- 1950s-1988; Efficient Monte carlo methods such as Metropolis and Gibbs sampling appeared.
- 1990-present: Bayesian statistics become the focus of mathematical statistics

Outline

- 1 History of Mathematical Statistics
- 2 Overview of Statistical Inference
- 3 Our Focus: Bayesian Statistical Inference
- 4 Beta and Gamma Distribution
- 5 Conjugate Prior: A Weapon of Bayesian
- 6 Application Case: Bayesian Ranking

Probability & Statistics



Statistical Inference

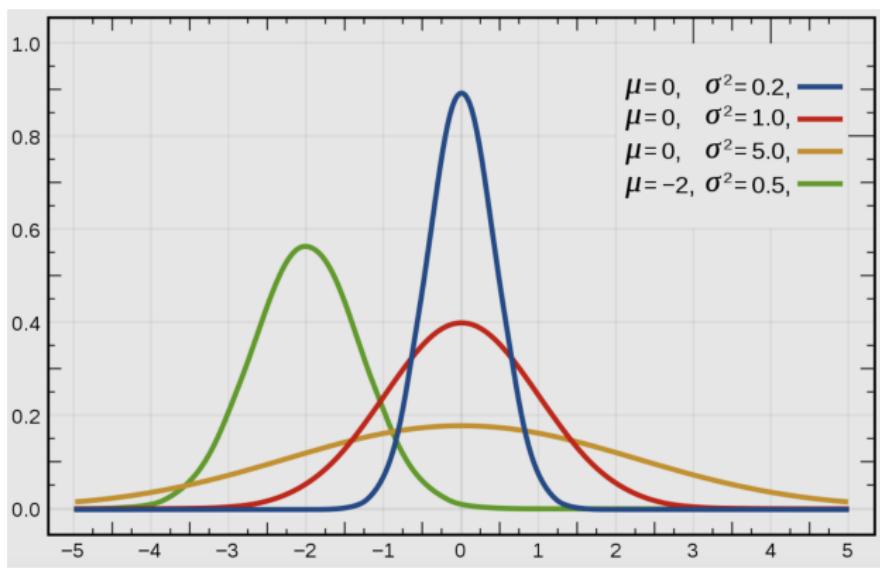
- The process of extracting information from available data
- Called "Learning" in CS
- Called "Signal Processing" in EE
- A typical question is: given n samples $X_1, \dots, X_n \sim F$ how do we infer F or some features of F (e.g. mean of F)?

Statistical Model

- A statistical model is a set of distributions (PMFs or PDFs)
- Our focus is the parametric model: a statistical model that can be parameterized by a finite number of parameters
- Example of $\text{Bern}(p)$: Bernoulli distribution with parameter p

Example: Parameterized Normal Distribution

$$\left\{ f(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right), \mu \in \mathbb{R}, \sigma > 0 \right\}.$$



Core Parts of Statistical Inference

- Point Estimation
- Interval Estimation (Confidence Interval)
- Hypothesis Testing

Our Focus: Parameterized Statistical Inference

- Given a statistical model $f(x; \theta) : \theta \in \Theta$
- θ is an unknown parameter in the parameter space Θ
- Now given n samples from such model: X_1, \dots, X_n
- How to make parameterized statistical inference?

Parameterized Statistical Inference: Bayesian versus Frequentist

- Difference relates to the nature of the unknown parameter θ
- Treated as an random variable with prior (known) distribution:
Bayesian approach with statistical model $f(x|\theta)$
- Treated as an unknown constant: **frequentist approach** with statistical model $f(x, \theta)$

Parameterized Statistical Inference: Bayesian versus Frequentist

- Bayesian:
 - ▶ View the world probabilistically, rather than as a set of fixed phenomena that are either known or unknown.
 - ▶ Prior information abounds and it is important and helpful to use it.
- Frequentist:
 - ▶ The parameters of interest are fixed and unchanging under all realistic circumstances
 - ▶ No information prior to the model specification.

Parameterized Statistical Inference: Bayesian versus Frequentist

- Bayesian:
 - ▶ Data are observed from the realized sample
 - ▶ Parameters are unknown and described probabilistically
 - ▶ Data are fixed
- Frequentist:
 - ▶ Data are a repeatable random sample: there is a frequency
 - ▶ Underlying parameters remain constant during this repeatable process.
 - ▶ Parameters are fixed

Point Estimation: Frequentist Perspective

- Point Estimation refers to providing a single “best guess” of parameter θ based on n random samples
- Estimator $\underline{g(X_1, \dots, X_n)}$: a function of $\underline{X_1, \dots, X_n}$
- Estimate $\underline{g(x_1, \dots, x_n)}$: when $\underline{X_1 = x_1, \dots, X_n = x_n}$

Interval Estimation: Frequentist Perspective

- Usually called “Confidence Interval”
- A $1 - \alpha$ confidence interval for a parameter θ is an interval $C_n = (a, b)$ where $a = a(X_1, \dots, X_n)$ and $b = b(X_1, \dots, X_n)$ are functions of the sample data such that

$$P_{\theta}(\theta \in C_n) \geq 1 - \alpha, \forall \theta \in \Theta.$$

$\alpha = 0.05$

- In words, (a, b) traps θ with probability $1 - \alpha$. We call $1 - \alpha$ the coverage of the confidence interval.

Hypothesis Testing: Frequentist Perspective

- A hypothesis is a statement about the data
- Start with a finite number of competing hypotheses and use the available sample data to decide which of them is true.
- Example 1: given a noisy picture, decide whether there is a person in the picture or not
- Example 2: given a noisy received signal, decide whether symbol 1 or 0 was sent by the transmitter
- Typically a hypothesis test is specified in terms of a test statistic $W(X_1, \dots, X_n)$, a function of the sample data

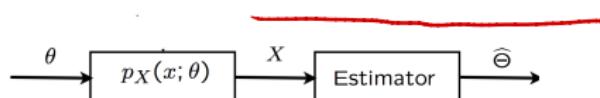
Example: Testing if a coin is Fair

- Samples $X_1, \dots, X_n \sim \text{Bern}(p)$ are results of n independent coin flips.
- H_0 : the coin is fair
- H_1 : the coin is not fair
- Parameterized hypothesis: $H_0 : p = \frac{1}{2}$ versus $H_1 : p \neq \frac{1}{2}$
- It seems reasonable to reject H_0 if test statistic W is large:

$$W(X_1, \dots, X_n) = \left| \frac{X_1 + \dots + X_n}{n} - \frac{1}{2} \right|.$$

Statistical Inference: Frequentist Perspective

- Classical statistics: unknown constant θ



- Hypothesis testing: $H_0 : \theta = 1/2$ versus $H_1 : \theta = 3/4$
- Composite hypotheses: $H_0 : \theta = 1/2$ versus $H_1 : \theta \neq 1/2$
- Estimation: design an **estimator** $\widehat{\theta}$, to “keep estimation **error** $\widehat{\theta} - \theta$ small”

Inference Rule of Frequentist: Maximum Likelihood Estimation (MLE)

- Given parameterized statistical model PMF $P_X(x; \theta)$ (or PDF $f_X(x; \theta)$)
- θ : unknown parameter.
- n samples: X_1, \dots, X_n
- We observe particular sample values x_1, \dots, x_n , then a **maximum likelihood estimate** (MLE) is a value of the parameter θ that maximizes the numerical function $P_X(x_1, \dots, x_n; \theta)$ (or $f_X(x_1, \dots, x_n; \theta)$) over all θ :

$$\hat{\theta}_n = \arg \max_{\theta} P_X(x_1, \dots, x_n; \theta)$$
$$\hat{\theta}_n = \arg \max_{\theta} f_X(x_1, \dots, x_n; \theta)$$

MLE under Independent Case

- When X_i are independent, we have

$$\log[P_X(x_1, \dots, x_n; \theta)] = \log \prod_{i=1}^n P_{X_i}(x_i; \theta) = \sum_{i=1}^n \log[P_{X_i}(x_i; \theta)]$$

$$\log[f_X(x_1, \dots, x_n; \theta)] = \log \prod_{i=1}^n f_{X_i}(x_i; \theta) = \sum_{i=1}^n \log[f_{X_i}(x_i; \theta)]$$

MLE under Independent Case

- Thus a **maximum likelihood estimate** (MLE) under independent case is shown as follows:

$$\hat{\theta}_n = \arg \max_{\theta} \sum_{i=1}^n \log[P_{X_i}(x_i; \theta)]$$

$$\hat{\theta}_n = \arg \max_{\theta} \underbrace{\sum_{i=1}^n \log[f_{X_i}(x_i; \theta)]}_{\text{red wavy line}}$$

Example: Biased Coin Problem

Let $\text{Bin}(n, p)$
 p : unknown parameter.

1^o. n coin tosses.

p is a constant.

n independent samples $X_1, \dots, X_n \sim \text{Bin}(n, p)$.

$$\begin{aligned} i=1, 2, \dots, n & \Pr(X_i=1) = p \\ & \Pr(X_i=0) = 1-p \end{aligned} \quad \left. \Pr(X_i=x) = p^x (1-p)^{n-x} \right\} \quad x=0 \text{ or } 1$$

2^o. $X_i = x_i, x_i = 0 \text{ or } 1$

$$\text{Likelihood} \quad \underline{P_X(x|p)} = \prod_{i=1}^n P_{X_i}(x_i; p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{\sum_{i=1}^n x_i} \cdot (1-p)^{n - \sum_{i=1}^n x_i}$$

$$S_n = \underline{x_1 + \dots + x_n} : \# \text{ of heads in } n \text{ coin tosses.} ; \quad = p^{S_n} (1-p)^{n-S_n}$$

$$3^o. \log(P_X(x|p)) = \underline{S_n(\log p + (n-S_n)\log(1-p))} = f(p)$$

$$\begin{bmatrix} f'(p)=0 \\ f''(p) \leq 0 \end{bmatrix}$$

$$\hat{p}_{MLE} \notin \underset{p}{\operatorname{argmax}} f(p)$$

$$\Rightarrow \hat{p}_{MLE} = \frac{1}{n} S_n = \underline{\frac{1}{n} (x_1 + \dots + x_n)}$$

Example: Biased Coin Problem

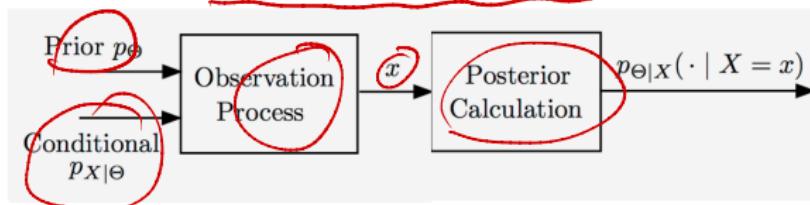
Outline

- 1 History of Mathematical Statistics
- 2 Overview of Statistical Inference
- 3 Our Focus: Bayesian Statistical Inference
- 4 Beta and Gamma Distribution
- 5 Conjugate Prior: A Weapon of Bayesian
- 6 Application Case: Bayesian Ranking

Statistical Inference: The Bayesian Perspective

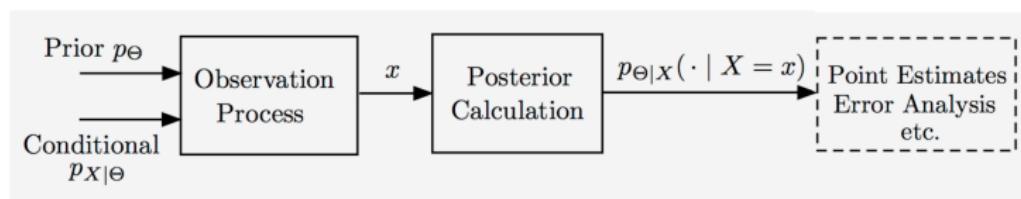
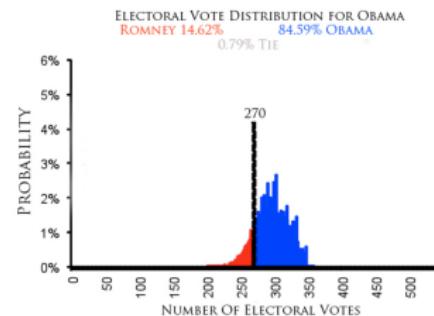
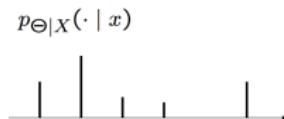
- Unknown Θ
 - treated as a random variable
 - prior distribution p_Θ or f_Θ
- Observation X Samples
 - observation model $p_{X|\Theta}$ or $f_{X|\Theta}$
- Use appropriate version of the Bayes rule to find $p_{\Theta|X}(\cdot | X = x)$ or $f_{\Theta|X}(\cdot | X = x)$

- Where does the prior come from?
 - symmetry
 - known range
 - earlier studies
 - subjective or arbitrary



The Output of Bayesian Statistical Inference

The complete answer is a posterior distribution:
PMF $p_{\Theta|X}(\cdot | x)$ or PDF $f_{\Theta|X}(\cdot | x)$



Recall: General LOTP

	Y discrete	Y continuous
X discrete	$P(X = x) = \sum_y P(X = x Y = y)P(Y = y)$	$P(X = x) = \int_{-\infty}^{\infty} P(X = x Y = y)f_Y(y)dy$
X continuous	$f_X(x) = \sum_y f_X(x Y = y)P(Y = y)$	$f_X(x) = \int_{-\infty}^{\infty} f_{X Y}(x y)f_Y(y)dy$

Recall: General Bayes' Rule

	Y discrete	Y continuous
X discrete	$P(Y = y X = x) = \frac{P(X=x Y=y)P(Y=y)}{P(X=x)}$	$f_Y(y X = x) = \frac{P(X=x Y=y)f_Y(y)}{P(X=x)}$
X continuous	$P(Y = y X = x) = \frac{f_X(x Y=y)P(Y=y)}{f_X(x)}$	$f_{Y X}(y x) = \frac{f_{X Y}(x y)f_Y(y)}{f_X(x)}$

Bayesian Posterior Calculation

The Four Versions of Bayes' Rule

- Θ discrete, X discrete:

$$p_{\Theta|X}(\theta | x) = \frac{p_\Theta(\theta)p_{X|\Theta}(x | \theta)}{\sum_{\theta'} p_\Theta(\theta')p_{X|\Theta}(x | \theta')}.$$

- Θ discrete, X continuous:

$$p_{\Theta|X}(\theta | x) = \frac{p_\Theta(\theta)f_{X|\Theta}(x | \theta)}{\sum_{\theta'} p_\Theta(\theta')f_{X|\Theta}(x | \theta')}.$$

- Θ continuous, X discrete:

$$f_{\Theta|X}(\theta | x) = \frac{f_\Theta(\theta)p_{X|\Theta}(x | \theta)}{\int f_\Theta(\theta')p_{X|\Theta}(x | \theta') d\theta'}.$$

- Θ continuous, X continuous:

$$f_{\Theta|X}(\theta | x) = \frac{f_\Theta(\theta)f_{X|\Theta}(x | \theta)}{\int f_\Theta(\theta')f_{X|\Theta}(x | \theta') d\theta'}.$$

Example of Inference Rule: The Maximum A Posteriori Probability (MAP)

- Given the observation value x , the MAP rule selects a value $\hat{\theta}$ that maximizes over θ the posterior distribution $p_{\Theta|x}(\theta|x)$ (if Θ is discrete) or $f_{\Theta|x}(\theta|x)$ (if Θ is continuous).
- Equivalently, it selects $\hat{\theta}$ that maximizes over θ :
 - $p_{\Theta}(\theta)p_{X|\Theta}(x|\theta)$: if both Θ and X are discrete
 - $p_{\Theta}(\theta)f_{X|\Theta}(x|\theta)$: if Θ is discrete and X is continuous
 - $f_{\Theta}(\theta)p_{X|\Theta}(x|\theta)$: if Θ is continuous and X is discrete
 - $f_{\Theta}(\theta)f_{X|\Theta}(x|\theta)$: if both Θ and X are continuous

Outline

- 1 History of Mathematical Statistics
- 2 Overview of Statistical Inference
- 3 Our Focus: Bayesian Statistical Inference
- 4 Beta and Gamma Distribution
- 5 Conjugate Prior: A Weapon of Bayesian
- 6 Application Case: Bayesian Ranking

Beta Distribution

$$a = b = 1 \quad f(x) \text{ constant} \quad 0 < x < 1.$$

$$\int_0^1 f(x) dx = 1 \quad \Rightarrow \quad f(x) = 1$$

Unif. f(x)

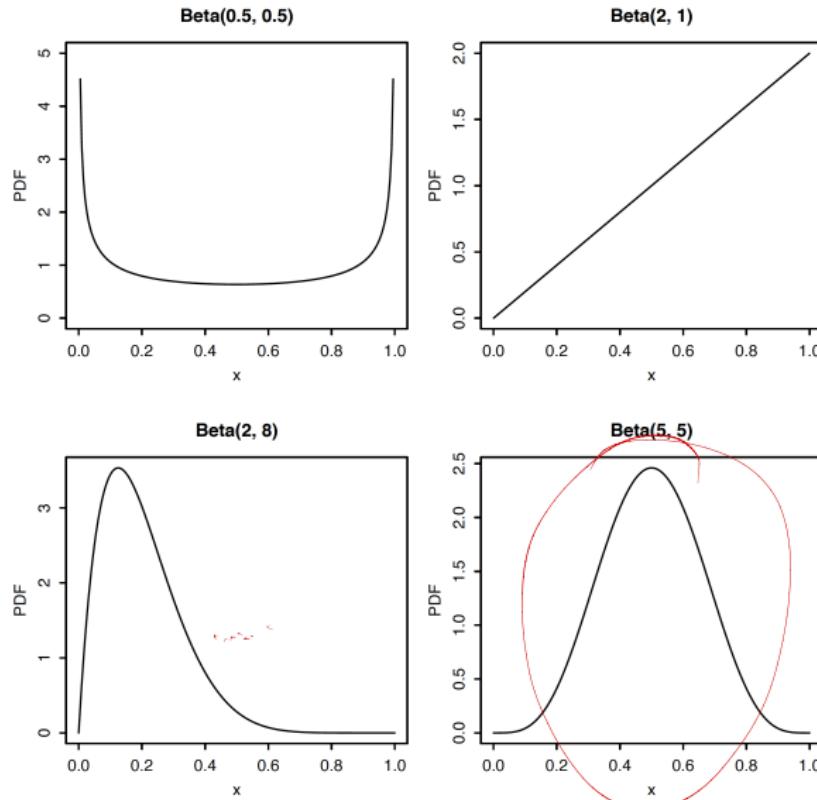
Definition

An r.v. X is said to have the *Beta distribution* with parameters a and b , $a > 0$ and $b > 0$, if its PDF is

$$f(x) = \frac{1}{\beta(a, b)} x^{a-1} (1-x)^{b-1}, \quad 0 < x < 1,$$

where the constant $\beta(a, b)$ is chosen to make the PDF integrate to 1. We write this as $X \sim \text{Beta}(a, b)$. Beta distribution is a generalization of uniform distribution.

PDF of Beta Distribution



Beta Integral

$$\beta(a, b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx.$$

Gamma Function

$$1 = \int_0^\infty \frac{1}{\Gamma(a)} \cdot x^{a-1} e^{-x} dx$$

$(0, \infty)$ PDF.

$$= \int_0^\infty f(x) dx$$

Gamma ($a, 1$)

Definition

The gamma function Γ is defined by

$$\Gamma(a) = \int_0^\infty x^a e^{-x} \frac{dx}{x},$$

for real numbers $a > 0$.

Property of Gamma Function

$$\Gamma(1) = \int_0^\infty t e^{-t} dt = 1$$

$$\Gamma(1+1) = 1 \cdot \Gamma(1) = 1$$

- $\Gamma(a+1) = a\Gamma(a)$ for all $a > 0$.
- $\Gamma(n) = (n-1)!$ if n is a positive integer.
 $\Gamma(3) = \Gamma(2+1)$

$$= 2\Gamma(1)$$

$$= \textcircled{2}$$

Gamma Distribution

$$a=1; f(y) = \frac{1}{\Gamma(1)} (\lambda y)^1 e^{-\lambda y} \cdot \frac{1}{y}$$

Prior of parameter
 $(0, +\infty)$

$$= \lambda e^{-\lambda y}, y > 0.$$

Definition

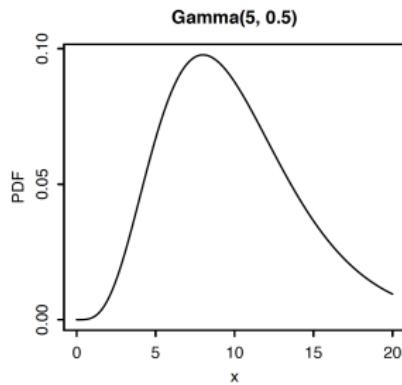
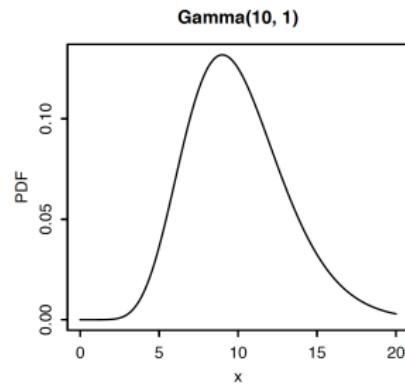
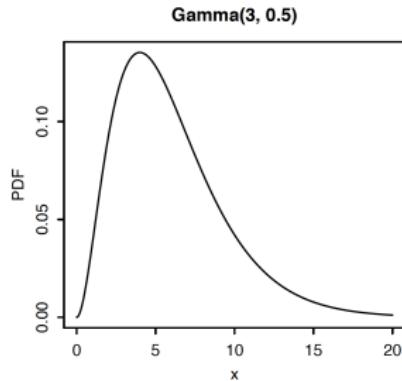
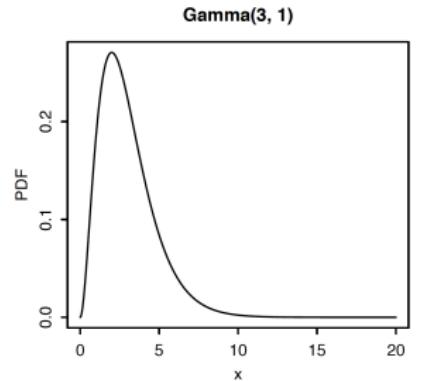
Expo(λ)

An r.v. Y is said to have the *Gamma distribution* with parameters a and λ , $a > 0$ and $\lambda > 0$, if its PDF is

$$f(y) = \frac{1}{\Gamma(a)} (\lambda y)^a e^{-\lambda y} \frac{1}{y}, \quad y > 0.$$

We write $Y \sim \underline{\text{Gamma}(a, \lambda)}$. Gamma distribution is a generalization of the exponential distribution.

PDF of Gamma Distribution



Moments of Gamma Distribution

$$Y \sim \text{Gamma}(\alpha, \lambda)$$

$$E(Y) = \frac{\alpha}{\lambda}$$

$$\text{Var}(Y) = \frac{\alpha}{\lambda^2}$$

Gamma: Convolution of Exponential

Theorem

Let X_1, \dots, X_n be i.i.d. $\text{Expo}(\lambda)$. Then

$$\underline{X_1 + \dots + X_n} \sim \text{Gamma}(n, \lambda)$$

MGF

Proof

Beta-Gamma Connection

When we add independent Gamma r.v.s X and Y with the same rate λ , the total $X + Y$ has a Gamma distribution, the fraction $\frac{X}{X+Y}$ has a Beta distribution, and the total is independent of the fraction.

Story: Bank–post Office

While running errands, you need to go to the bank, then to the post office. Let $X \sim \text{Gamma}(a, \lambda)$ be your waiting time in line at the bank, and let $Y \sim \text{Gamma}(b, \lambda)$ be your waiting time in line at the post office (with the same λ for both). Assume X and Y are independent. What is the joint distribution of $T = X + Y$ (your total wait at the bank and post office) and $W = \frac{X}{X+Y}$ (the fraction of your waiting time spent at the bank)?

Story: Bank-post Office

$$T = X + Y \quad ; \quad W = \frac{X}{X+Y} \quad ;$$

$$\textcircled{1} \quad \begin{array}{l} t > 0 \\ w > 0 \end{array} \quad \begin{array}{l} t = x+y \\ w = \frac{x}{x+y} \end{array} \quad \Rightarrow \quad \begin{array}{l} x = tw \\ y = t - w \end{array} \quad \Rightarrow \quad \frac{\partial(x,y)}{\partial(t,w)} = \begin{pmatrix} w & t \\ tw & -t \end{pmatrix}$$

$$\det\left(\frac{\partial(x,y)}{\partial(t,w)}\right) = -t < 0$$

$X \sim \text{Gamma}(a, \lambda)$
 $Y \sim \text{Gamma}(b, \lambda)$

X and Y are independent

$$\textcircled{2} \quad f_{T,W}(t,w) = f_{X,Y}(x,y) \cdot \left| \det\left(\frac{\partial(x,y)}{\partial(t,w)}\right) \right| = f_X(x) \cdot f_Y(y) \cdot t$$

$$= \frac{1}{\Gamma(a)} (\lambda x)^a e^{-\lambda x} \cdot \frac{1}{\Gamma(b)} (\lambda y)^b e^{-\lambda y} \cdot \frac{1}{t} \cdot t \quad (x = tw, y = t - w)$$

$$= \frac{1}{\Gamma(a)} \cdot \frac{1}{\Gamma(b)} \cdot w^{a-1} (t-w)^{b-1} \cdot (at)^{a+b} e^{-at} \cdot \frac{1}{t}$$

$$= \frac{\Gamma(a+b)}{\Gamma(a) \cdot \Gamma(b)} w^{a-1} (t-w)^{b-1} \cdot$$

$$\frac{g(t)}{\Gamma(a+b) \cdot (at)^{a+b} e^{-at} \cdot \frac{1}{t}}$$

Story: Bank–post Office

③ T and W are independent.

$$T \sim \text{Gamma}(a+b, \lambda)$$

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

$$\frac{1}{\beta(a,b)} w^{a-1} (1-w)^{b-1}$$

$$\frac{\pi(a+b)}{\pi(a) \cdot \pi(b)} w^{a-1} (1-w)^{b-1}$$

④ $\beta(a, b) = \frac{\pi(a) \cdot \pi(b)}{\pi(a+b)}$

Story: Bank–post Office

$$\textcircled{5} \quad W \sim \text{Beta}(a, b)$$

$$a > 0, b > 0.$$

$$E(W) ?$$

$$T = X + Y, \quad W = \frac{X}{X+Y};$$

T and W are independent.

$$E[T \cdot W] = E[T] \cdot E[W]$$

$$\begin{cases} X \sim \text{Gamma}(a, \lambda) \\ E(X) = \frac{a}{\lambda} \end{cases}$$

$$\begin{cases} Y \sim \text{Gamma}(b, \lambda) \\ E(Y) = \frac{b}{\lambda} \end{cases}$$

$$\Rightarrow E[W] = \frac{E[T \cdot W]}{E[T]} = \frac{E[X]}{E[X+Y]} = \frac{E[X]}{E[X]+E[Y]}$$

$$= \frac{\frac{a}{\lambda}}{\frac{a}{\lambda} + \frac{b}{\lambda}} = \frac{a}{a+b}, \quad a > 0, b > 0.$$

Outline

- 1 History of Mathematical Statistics
- 2 Overview of Statistical Inference
- 3 Our Focus: Bayesian Statistical Inference
- 4 Beta and Gamma Distribution
- 5 Conjugate Prior: A Weapon of Bayesian
- 6 Application Case: Bayesian Ranking

Conjugate Prior

- Before Monte Carlo, posterior calculation is hard
- Conjugate Prior: reduce the computing complexity of posterior distribution
- Loosely speaking, a prior distribution is conjugate to the likelihood model if both the prior and posterior distribution stay in the same distribution family.

Story: Beta-Binomial Conjugacy

① p : unknown, model it as a r.v. $f[0,1]$

Prior distribution. $p \sim \text{Beta}(a, b)$.

② Data model : n tosses of coin; X : # of heads.
Likelihood model $\times | p = P \sim \text{Bin}(n, p) \rightarrow \text{real number}$

We have a coin that lands Heads with probability p , but we don't know what p is. Our goal is to infer the value of p after observing the outcomes of n tosses of the coin. The larger that n is, the more accurately we should be able to estimate p .

③ $f(p)$: prior PDF of p .

$f(p|X=k)$: posterior PDF of p .

Story: Beta-Binomial Conjugacy $\Pr \sim \text{Beta}(a, b)$

$$f(p|X=k) = \frac{\Pr(X=k|p)}{\Pr(X \geq k)} \cdot f(p) = \frac{\binom{n}{k} p^k (1-p)^{n-k}}{\Pr(X \geq k)} \cdot \frac{1}{\Pr(X \geq k)} \frac{\Pr(X=k|p)}{\Pr(X \geq k)} \cdot \frac{\Pr(X=k|p)}{\Pr(X \geq k)} \cdot \frac{\Pr(X=k|p)}{\Pr(X \geq k)} \cdot \frac{\Pr(X=k|p)}{\Pr(X \geq k)} \cdot \frac{\Pr(X=k|p)}{\Pr(X \geq k)}$$

$$\Pr(X=k) \stackrel{\text{Def}}{=} \int_0^1 \Pr(X=k|p) f(p) dp = \int_0^1 \binom{n}{k} p^k (1-p)^{n-k} \frac{1}{\Pr(X \geq k)} p^{a-1} (1-p)^{b-1} dp$$

④ $f(p|X=k)$: a function of p . [every item does not depend on p can be regard as a constant]

$$f(p|X=k) \propto \underbrace{p^k a^{-1} \cdot (1-p)^{n-k} b^{-1}}_{C \cdot p^k a^{-1} (1-p)^{n-k} b^{-1}} \sim \text{Beta}(a, b)$$

$$\Rightarrow p | X=k \sim \text{Beta}(a+k, b+n-k)$$

Story: Beta-Binomial Conjugacy

Likelihood model

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

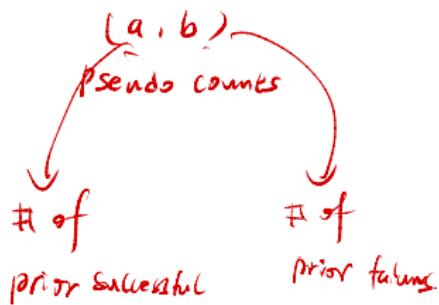
↑
Prior

+

$$X|P=p \sim \text{Bin}(n, p)$$

$X=k$ out of n tosses

$$P|X=k \sim \text{Beta}(a+k, b+n-k)$$



$$\xrightarrow{X=k} (a+k, b+n-k)$$

k success
 $n-k$ failures

↓
 $X'=m$
out of
 n tosses.

$$(a+k+m, b+n-k+m)$$

Story: Beta-Binomial Conjugacy

Story: Beta-Binomial Conjugacy

- Furthermore, notice the very simple formula for updating the distribution of p .
- We just add the number of observed successes, k , to the first parameter of the Beta distribution.
- We also add the number of observed failures, $n - k$, to the second parameter of the Beta distribution.
- So a and b have a concrete interpretation in this context:
 - ▶ a as the number of prior successes in earlier experiments
 - ▶ b as the number of prior failures in earlier experiments
 - ▶ a, b : pseudo counts

Mean vs. Bayesian Average

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

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

- Infer the value of p (probability of coin lands heads)
- Observed k heads out of n tosses of the coin
- Mean: $\frac{k}{n}$ MLE
- Bayesian Average: $E(p|X=k) = \frac{a+k}{a+b+n}$ $\sim \text{Beta}(a+k, b+n-k)$
- Suppose the prior distribution is Unif(0,1): $a=1, b=1$
- Bayesian Average: $\frac{k+1}{n+2}$
- When $k=n$, we have: 1 (mean) vs. $\frac{n+1}{n+2}$ (Bayesian average)
 $\xrightarrow{n \rightarrow \infty} \frac{n+1}{n+2} \rightarrow 1$

Story: Beta-Binomial Conjugacy

If we have a Beta prior distribution on p and data that are conditionally Binomial given p , then when going from prior to posterior, we don't leave the family of Beta distributions. We say that **the Beta is the conjugate prior of the Binomial.**

Example: Inference of A Biased Coin

$$1^{\circ}. \quad \Theta \sim \text{unif}(0,1) = \text{Beta}(1,1) \quad \# \text{ of heads } X | \Theta = \theta$$

By Beta-Binomial Conjugacy - $\Theta | X=k \sim \text{Beta}(1+k, n-k) \sim \text{Bin}(n, \theta)$

$$\hat{\theta} = E[\theta | X=k] = \frac{k+1}{n+2}$$

We wish to estimate the probability of landing heads, denoted by θ , of a biased coin. We model θ as the value of a random variable Θ with a known prior PDF $f_\Theta \sim \text{Unif}(0,1)$. We consider n independent tosses and let X be the number of heads observed. Find the MAP estimator of Θ .

2^o. MAP estimation

$$f_{\Theta|X=k}(\theta) \propto \frac{\theta^k (1-\theta)^{n-k}}{f_{\Theta}(0,1)}$$

$$\hat{\theta}_{\text{MAP}} = \underset{\Theta}{\operatorname{argmax}} \frac{f_{\Theta|X=k}(\theta)}{= \frac{k}{n}} = \hat{\theta}_{\text{MLE}}$$

$\equiv \operatorname{argmax} \log f_{\Theta|X=k}(\theta)$

Solution

Solution

Example: Revisit Biased Coin Problem

1^o. MMSE : $E[\Theta|X=k]$: $\Theta \sim \text{Unif}(0,1) = \text{Beta}(1,1)$
 $E[\Theta|X]$: # of heads $X|\Theta=\theta \sim \text{Binom}(\theta)$

By Beta-Binomial conjugacy, $\Theta|X=k \sim \text{Beta}(1+k, n-k)$.

We wish to estimate the probability of landing heads, denoted by θ , of a biased coin. We model θ as the value of a random variable Θ with a known prior PDF $f_\Theta \sim \text{Unif}(0,1)$. We consider n independent tosses and let X be the number of heads observed. Find the MMSE $E(\Theta|X)$ and LLSE $L(\Theta|X)$

$$\Rightarrow E[\Theta|X=k] = \frac{k+1}{n+2}$$

$$\Rightarrow E[\Theta|X] = \left(\frac{X+1}{n+2} \right)$$

Solution 2° LLSE : $L[\theta | X] = \underline{E[\theta]} + \frac{\text{cov}(\theta, x)}{\text{var}(x)} (x - \underline{E(x)})$

$$\theta \sim \text{unif}(0,1) \Rightarrow E(\theta) = \frac{1}{2}, \text{var}(\theta) = \frac{1}{12}, E(\theta^2) = \frac{1}{3}$$

$$X|\theta = \theta \sim \text{Bin}(n, \theta) \Rightarrow E[X|\theta = \theta] = n\theta \Rightarrow E[X|\theta] = \underline{n\theta}$$

$$\text{Var}[X|\theta = \theta] = n\theta(1-\theta) \Rightarrow \text{Var}[X|\theta] = n\theta(1-\theta)$$

$$\Rightarrow E[X] = E[E[X|\theta]] = E[n\theta] = n E[\theta] = \underline{\frac{n}{2}},$$

$$\text{Var}[X] = E[\underline{\text{Var}[X|\theta]}] + \text{Var}[E[X|\theta]]$$

$$= E[n\theta(1-\theta)] + \text{Var}[\underline{n\theta}]$$

$$= n(E[\theta] - E[\theta^2]) + n^2 \underline{\text{Var}[\theta]}$$

$$= n\left(\frac{1}{2} - \frac{1}{3}\right) + n^2 \cdot \frac{1}{12} = \frac{n}{12}(n+2)$$

Solution

$$\begin{aligned}\Rightarrow \text{cov}(x, \theta) &= \underline{E[\theta x]} - E[\theta] \cdot E[x] \\&= E[\underline{E[\theta x | \theta]}] - E[\theta] \cdot E[x] \\&= E[\theta \underline{E[x | \theta]}] - E[\theta] \cdot E[x] \\&= E[\theta \cdot n\theta] - E[\theta] \cdot E[x] \\&= n E[\theta^2] - E[\theta] \cdot E[x] \\&= n \cdot \frac{1}{2} - \frac{1}{2} \cdot \frac{n}{2} = \frac{1}{2}n\end{aligned}$$

$$\begin{aligned}\Rightarrow \text{LLSE } L[\theta | x] &= E[\theta] + \frac{\text{cov}(\theta, x)}{\text{var}(x)} [x - E(x)] \\&= \frac{1}{2} + \frac{\frac{1}{2}n}{\frac{n}{2}(\text{ave}_n)} \left(x - \frac{n}{2}\right) = \frac{x+1}{n+2} = \text{MMSE}\end{aligned}$$

Solution

Recall: Story of Multinomial Distribution

Beta - Binomial conjugacy.
↓ ↓
D Multinomial.

Each of n objects is independently placed into one of k categories. An object is placed into category j with probability p_j , where the p_j are nonnegative and $\sum_{j=1}^k p_j = 1$. Let X_1 be the number of objects in category 1, X_2 the number of objects in category 2, etc., so that $X_1 + \dots + X_k = n$. Then $X = (X_1, \dots, X_k)$ is said to have the Multinomial distribution with parameters n and $\mathbf{p} = (p_1, \dots, p_k)$. We write this as $\mathbf{X} \sim \text{Mult}_k(n, \mathbf{p})$.

Recall: Multinomial Joint PMF

Theorem

If $\mathbf{X} \sim \text{Mult}_k(n, \mathbf{p})$, then the joint PMF of \mathbf{X} is

$$P(X_1 = n_1, \dots, X_k = n_k) = \frac{n!}{n_1! n_2! \dots n_k!} \cdot p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}$$

for n_1, \dots, n_k satisfying $n_1 + \dots + n_k = n$.

$$\cancel{\frac{n!}{n_1! n_2! \dots n_k!}} \underbrace{p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}}$$

Dirichlet Distribution

Beta distribution
PDF (a, b)

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

$$\propto p^{a-1} (1-p)^{b-1}$$

$p_1 = p, p_2 = 1-p; p_1 + p_2 = 1$

The Dirichlet distribution is parameterized by a vector α of positive real numbers.

- The PDF is:

$$f(p_1, p_2, \dots, p_k; \alpha_1, \alpha_2, \dots, \alpha_K) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{i=1}^k p_i^{\alpha_i - 1}$$

where $p_1 + \dots + p_k = 1$ and $0 < p_i < 1$.

- The marginal mean of P_j is:

$$E(P_j) = \frac{\alpha_j}{\sum_{i=1}^K \alpha_i}$$

Story: Dirichlet-Multinomial Conjugacy

If we have a Dirichlet prior distribution on p and data that are conditionally Multinomial given p , then when going from prior to posterior, we don't leave the family of Dirichlet distributions. We say that the Dirichlet is the conjugate prior of the Multinomial.

Likelihood Model: Discrete

Sample Space	Sampling Dist.	Conjugate Prior	Posterior
$\mathcal{X} = \{0, 1\}$	<u>Bernoulli(θ)</u>	<u>Beta(α, β)</u>	<u>Beta($\alpha + n\bar{X}, \beta + n(1 - \bar{X})$)</u>
$\mathcal{X} = \mathbb{Z}_+$	<u>Poisson(λ)</u>	<u>Gamma(α, β)</u>	<u>Gamma($\alpha + n\bar{X}, \beta + n$)</u>
$\mathcal{X} = \mathbb{Z}_{++}$	<u>Geometric(θ)</u>	<u>Gamma(α, β)</u>	<u>Gamma($\alpha + n, \beta + n\bar{X}$)</u>
$\mathcal{X} = \mathbb{H}_K$	<u>Multinomial(θ)</u>	<u>Dirichlet(α)</u>	<u>Dirichlet($\alpha + n\bar{X}$)</u>

Likelihood Model: Continuous

Sampling Dist.	Conjugate Prior	Posterior
Uniform(θ)	Pareto(ν_0, k)	Pareto $(\max\{\nu_0, X_{(n)}\}, n + k)$
<u>Exponential(θ)</u>	<u>Gamma(α, β)</u>	<u>Gamma($\alpha + n, \beta + n\bar{X}$)</u>
<u>$N(\mu, \sigma^2)$, known σ^2</u>	<u>$N(\mu_0, \sigma_0^2)$</u>	<u>$N\left(\left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)^{-1} \left(\frac{\mu_0}{\sigma_0^2} + \frac{n\bar{X}}{\sigma^2}\right), \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)^{-1}\right)$</u>
$N(\mu, \sigma^2)$, known μ	InvGamma(α, β)	InvGamma $\left(\alpha + \frac{n}{2}, \beta + \frac{n}{2}(X - \mu)^2\right)$
$N(\mu, \sigma^2)$, known μ	ScaledInv- $\chi^2(\nu_0, \sigma_0^2)$	ScaledInv- $\chi^2\left(\nu_0 + n, \frac{\nu_0 \sigma_0^2}{\nu_0 + n} + \frac{n(X - \mu)^2}{\nu_0 + n}\right)$
$N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, known $\boldsymbol{\Sigma}$	$N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$	$N\left(\mathbf{K} \left(\boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + n \boldsymbol{\Sigma}^{-1} \bar{\mathbf{X}}\right), \mathbf{K}\right)$, $\mathbf{K} = (\boldsymbol{\Sigma}_0^{-1} + n \boldsymbol{\Sigma}^{-1})^{-1}$
$N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, known $\boldsymbol{\mu}$	InvWishart(ν_0, \mathbf{S}_0)	InvWishart($\nu_0 + n, \mathbf{S}_0 + n\bar{\mathbf{S}}$), $\bar{\mathbf{S}}$ sample covariance

Outline

- 1 History of Mathematical Statistics
- 2 Overview of Statistical Inference
- 3 Our Focus: Bayesian Statistical Inference
- 4 Beta and Gamma Distribution
- 5 Conjugate Prior: A Weapon of Bayesian
- 6 Application Case: Bayesian Ranking

Rating System

- Consumers rely on the collective intelligence of other consumers: rating
- A common metric: 5 star rating
- Requirement: many ratings are needed to make this system work
- Quality of rating system depends on
 - ▶ average number of stars
 - ▶ average number of reviews

Which One to Choose?

- 1. Presto Coffee Pot - average rating of 5 (1 review).
- 2. Cuisinart Brew Central - average rating of 4.1 (78 reviews).

Example: Movie Ranking

- Data Set : <http://grouplens.org/datasets/movielens/>
- Top Ten Movies

Top 10 Movies chosen by Mean

title	count	mean
Aiqing wansui (1994)	1	5
They Made Me a Criminal (1939)	1	5
Great Day in Harlem, A (1994)	1	5
Saint of Fort Washington, The (1993)	2	5
Entertaining Angels: The Dorothy Day Story (1996)	1	5
Someone Else's America (1995)	1	5
Star Kid (1997)	3	5
Santa with Muscles (1996)	2	5
Prefontaine (1997)	3	5
Marlene Dietrich: Shadow and Light (1996)	1	5

Tool: Bayesian Estimation

- Mean of star reviews with a limited number of observations
- Useful for recommender services and other predictive algorithms that use preference space measures like star reviews.

Joint Distribution

- To use Bayesian estimation to compute the posterior probability for star ratings, we must use a joint distribution.
- We are not estimating the distribution of some scalar value X but, rather, the joint distributions of the probability estimate of whether or not the reviewer will give the movie a 1, 2, 3, 4, or 5 star rating (not just a simple thumbs up or down).
- In this case, the random variable is a categorical distribution because it can take some value within 1,2,3,4,5 with probabilities as follows:

$$p_1 + p_2 + p_3 + p_4 + p_5 = 1$$

Multinomial Distribution

- We can compute our posterior probability with N observations for five categories with corresponding numbers K_1, K_2, K_3, K_4, K_5 as follows:

$$Pr(O|p_1, p_2, p_3, p_4, p_5) \propto p_1^{K_1} p_2^{K_2} p_3^{K_3} p_4^{K_4} p_5^{K_5}$$

where $K_1 + \dots + K_5 = N$.

$$0 < p_i < 1$$

$$i=1, \dots, 5$$

- This is a multinomial distribution.

Dirichlet Distribution: Prior $\alpha^0 = (\alpha_1^0, \alpha_2^0, \dots, \alpha_r^0)$

- If we include our prior as a distribution of the exact same form in the proportionality equation (e.g. a Dirichlet distribution with parameter α^0), then

$$Pr(p_1, p_2, p_3, p_4, p_5 | O) \propto \prod_{j=1}^5 p_j^{K_j + \alpha_j^0 - 1}$$

- This is another Dirichlet distribution with another parameter α^1 :

$$\alpha_j^1 = K_j + \alpha_j^0, \forall j$$

Expected Average

- What is the expected value of the average rating given a posterior in the shape of our Dirichlet distribution?
- The expected value of the average rating based on the posterior is then computed for our star ratings as follows:

$$E(\underline{p_1} + \underline{2p_2} + \underline{3p_3} + \underline{4p_4} + \underline{5p_5}|O) = \sum_{i=1}^5 iE(p_i|O)$$

- Using our Dirichlet distribution we can compute the probability of a star value given our observations as the ratio of the Dirichlet parameter for that star to the sum of the Dirichlet parameters:

$$E(p_i|O) = \frac{\alpha_i^1}{\sum_{j=1}^5 \alpha_j^1}$$

Intra-Item: Bayesian Average Rating

$$\sum_{i=1}^5 i \cdot \frac{\alpha_i^0}{\sum_{j=1}^5 \alpha_j^0} = \frac{\sum_{i=1}^5 i \cdot \alpha_i^0}{\sum_{j=1}^5 \alpha_j^0} = \frac{\sum_{i=1}^5 i (\alpha_i^0 + K_i)}{\sum_{j=1}^5 (\alpha_j^0 + K_j)}$$

$$\text{Bayes Average Rating} = \frac{\sum_{i=1}^5 i \alpha_i^0 + \sum_{i=1}^5 i K_i}{N + \sum_{i=1}^5 \alpha_i^0}$$

- N : the number of reviews
- $\sum_{i=1}^5 i K_i$: sum of all review scores
- $\sum_{i=1}^5 \alpha_i^0$: prior(given) number of reviews
- $\sum_{i=1}^5 i \alpha_i^0$: prior sum of all review scores

$$= \frac{\sum_{i=1}^5 i \alpha_i^0 + \sum_{i=1}^5 i K_i}{\sum_{j=1}^5 \alpha_j^0 + \sum_{j=1}^5 K_j}$$

(1)

Intra-Item: Bayesian Average Rating

$$1^o \quad C=0 \quad \Rightarrow \quad \frac{\sum(\text{ratings})}{N}$$

$$2^o \quad N=0 \quad \Rightarrow \quad m$$

$$\text{Bayes Average Rating} = \frac{C \cdot m + \sum(\text{ratings})}{C + N}$$

- N : the number of reviews
 - m : a prior for the average of review scores
 - C : a prior for the number of reviews
-

Example: Movie Ranking

- Data Set : <http://grouplens.org/datasets/movielens/>
- Top Ten Movies

Case 1: $m = 3.25$ & $C = 50$

title	bayes	count	mean
One Flew Over the Cuckoo's Nest (1975)	4.125796	264	4.291667
Raiders of the Lost Ark (1981)	4.145745	420	4.252381
Rear Window (1954)	4.167954	209	4.387560
The Silence of the Lambs (1991)	4.171591	390	4.289744
The Godfather (1972)	4.171706	413	4.283293
The Usual Suspects (1995)	4.206625	267	4.385768
Casablanca (1942)	4.250853	243	4.456790
The Shawshank Redemption (1994)	4.265766	283	4.445230
Star Wars (1977)	4.270932	583	4.358491
Schindler's List (1993)	4.291667	298	4.466443

Case 2: $m = 2$ & $C = 6$

title	count	bayes	mean
One Flew Over the Cuckoo's Nest (1975)	264	4.244526	4.291667
The Godfather (1972)	413	4.252955	4.283293
The Silence of the Lambs (1991)	390	4.257500	4.289744
Star Wars (1977)	583	4.335582	4.358491
The Usual Suspects (1995)	267	4.335740	4.385768
<u>The Wrong Trousers (1993)</u>	118	4.351562	4.466102
<u>A Close Shave (1995)</u>	112	4.368852	4.491071
The Shawshank Redemption (1994)	283	4.395904	4.445230
Casablanca (1942)	243	4.399209	4.456790
Schindler's List (1993)	298	4.418831	4.466443

Inter-Items: Pseudo Bayesian Average Rating

$$\bar{m}_i = \frac{C_i \cdot m_i + \sum(\text{ratings})}{C_i + N}$$

- \bar{m}_i : bayesian average rating for item i
- N : the number of reviews for all items
- m_i : average of review scores for item i
- C_i : the number of reviews for item i

Example: Bayesian Changes Order

$$N = 10 + 15 + 228 + 150 + 124 = 527$$

$$\sum(\text{rat} \cdot s) = 10 \times 4.920 + 15 \times 4.667 + 228 \times 4.535 \\ + 150 \times 4.310 + 124 \times 4.298 = 2332.637$$

MacBook	No. Ratings	Ave. Rating	Rank	Bayesian Rating	Bayesian Rank
MB991LL	10	4.920	1	4.436	2
MB403LL	15	4.667	2	4.433	3
MB402LL	228	4.535	3	4.459	1
MC204LL	150	4.310	4	4.401	5
MB061LL	124	4.298	5	4.402	4

$$\bar{m}_1 = \frac{\sum m_i + \sum(\text{rat} \cdot s)}{C_1 + N} = \frac{10 \times 4.920 + 2332.637}{10 + 527} = 4.436$$

Reverse Engineering Amazon

- Bayesian adjustment
- Recency of view
- Reputation score

Key Factors

- Bayesian ranking
- Too few or too outdated reviews penalized
- Very high quality reviews help a lot

Summary

- Average ratings scalarize a vector and ranks
- Number of ratings should matter, Bayesian ranking does that
- Other statistical methods help too

References

- Chapters 9 of **BH**
- Chapters 4 & 6 & 8 of **BT**