

Conjugacy, Likelihood Principle

© Hyon-Jung Kim 2024

Posterior mean and posterior variance

What will be discussed...

- More examples on
 - forming the prior distributions
 - finding 'kernels'
 - computing integral (sum) of the kernel = $1 / \text{normalizing constant}$
- What are the conjugate priors? How do we find them?
- Posterior mean and posterior variance
 - How are they denoted? How do we find them?
- Informative priors, Reference priors, improper priors

Recap – The Bayesian Framework

i) Suppose we observe an iid sample of data : $\mathbf{X} = (X_1, \dots, X_n)$

Then, our model for the distribution of the data will give us the likelihood

$$L(\theta | \mathbf{x}) = f(\mathbf{x} | \theta) = \prod_{i=1}^n f(x_i | \theta)$$

ii) We also must specify $f(\theta)$, the prior distribution for θ , based on any knowledge we have about θ **before** observing the data.

iii) By Bayes' Rule,

$$\text{Posterior distribution is } f(\theta | \mathbf{x}) = \frac{f(\mathbf{x} | \theta) f(\theta)}{f(\mathbf{x})} = \frac{f(\mathbf{x} | \theta) f(\theta)}{\int f(\mathbf{x} | \theta) f(\theta) d\theta}$$

• $f(\mathbf{x})$ is a normalizing constant that ensures $f(\theta | \mathbf{x})$ integrating to 1. We don't have to calculate it (except posterior probabilities for discrete cases).

$$f(\theta | \mathbf{x}) \propto f(\mathbf{x} | \theta) f(\theta) : \text{posterior} \propto \text{likelihood} \times \text{prior}$$

Review: Binomial-Beta example

- Derive posterior by $f(\theta|x) \propto f(x|\theta)f(\theta)$: posterior \propto Lkhd x prior
- Suppose that $X|\theta \sim \text{Binom}(n, \theta)$. Take $\theta \sim \text{Beta}(\alpha, \beta)$
- Lkhd x prior: $f(x|\theta)f(\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$

Ignoring constants (in θ),

$$\begin{aligned} f(x|\theta)f(\theta) &\propto \theta^x (1-\theta)^{n-x} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &= \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1} \quad : \text{kernel of Beta}(x+\alpha, n-x+\beta) \end{aligned}$$

Then, the posterior distribution is $\theta|x \sim \text{Beta}(x+\alpha, n-x+\beta)$

Remarks

- In a real problem, we need to specify the values of our hyperparameters α and β of our prior.
- Ideally our choices of α and β should reflect our prior beliefs about θ .
 - If we have no prior idea what θ is, we could set $\alpha = \beta = 1$, which corresponds to a Uniform(0, 1) prior for θ : completely flat, so that all values of θ are equally likely a priori.
- If we have more informative prior beliefs about the value of θ , we could choose α and β to reflect that.
 - Plots of the Beta pdf for various values of α and β can help inform the prior specification (see the plots from the previous lecture).

Forming a Beta prior distribution

- The expected value and variance of a $\text{Beta}(\alpha, \beta)$ r.v. is

$$E[\theta] = \frac{\alpha}{\alpha + \beta} \qquad \text{Var}[\theta] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

- So if our prior belief is that θ is closer to 0 than to 1, we should choose our hyperparameters α and β such that $\alpha < \beta$.
- If our prior belief is strong that θ is near a certain value, we can pick α and β so that this variance is small.
 - If our prior belief is less certain, pick α and β so that the variance is large.
- The mode (location where the pdf reaches its maximum) for $\text{Beta}(\alpha, \beta)$ pdf is $\frac{\alpha - 1}{\alpha + \beta - 2}$

Examples of forming Likelihood functions

- Suppose students in a statistics class conduct a study to estimate the proportion of cars in a Campus Parking lot that are black.
 - i) Student A observes the first 10 cars and record X , the number that are black : $NB, B, NB, NB, NB, B, NB, NB, NB, B$
 $\Rightarrow x = 3$
 - ii) Student B decides to observe cars until the third black one parks and record Y , the total number of cars that drive by until the third black one.

B records $y = 10$.
 - iii) Student C decides to observe every car and records either 1 or 0 according to whether the i 'th car is black. $Z : 0\ 1\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 1$

Examples of forming Likelihood functions

- Suppose students in a statistics class conduct a study to estimate the proportion of cars in a Campus Parking lot that are black.

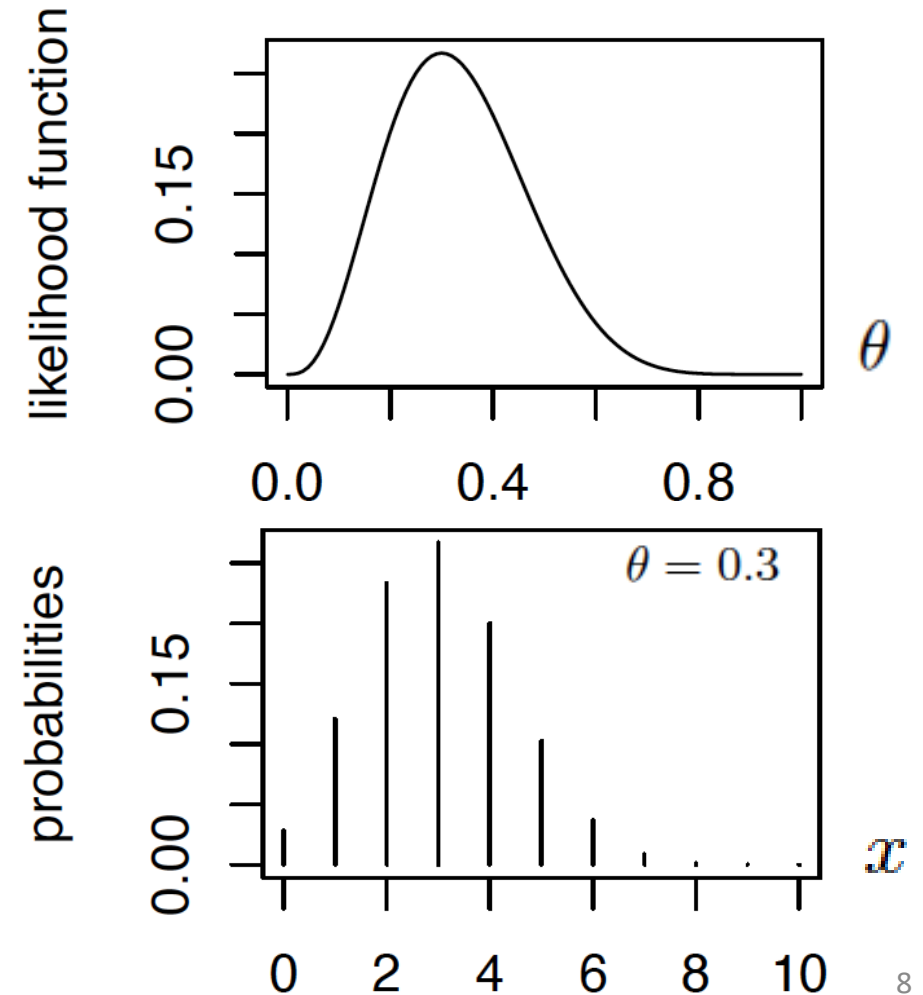
i) Student A observes the first 10 cars and record X , the number that are black :

$NB, B, NB, NB, NB, B, NB, NB, NB, B$

$$\Rightarrow x = 3$$

This is a Binomial experiment
and the statistical model is $X \sim \text{Bin}(10, \theta)$;
the likelihood function is

$$L_A(\theta | x) = \binom{10}{3} \theta^3 (1 - \theta)^7$$



ii) Student B decides to observe cars until the third black one parks and record Y , the total number of cars that drive by until the third black one.

B records $y = 10$. For B, the likelihood function is

$$\begin{aligned} L_B(\theta | y) &= P[Y = 10 | \theta] = P[2 \text{ blacks among first 9 cars}] \times P[10\text{'th car is black}] \\ &= \binom{9}{2} \theta^2 (1 - \theta)^7 \times \theta = \binom{9}{2} \theta^3 (1 - \theta)^7 \end{aligned}$$

iii) Student C decides to observe every car and records either 1 or 0 according to whether the i 'th car is black. $Z: 0 \ 1 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 1$

The likelihood function is

$$L_C(\theta | z) = (1 - \theta)\theta (1 - \theta)(1 - \theta)(1 - \theta) \theta (1 - \theta)(1 - \theta)(1 - \theta)\theta = \theta^3 (1 - \theta)^7$$

- All experiments yield equal (proportional to) likelihoods and contains the same information. \Rightarrow They should produce equivalent inference about θ .

Likelihood Principle

1) Suppose that only individual binary values are observed for data:

$$X_1, \dots, X_n / \theta \sim \text{Bern}(\theta) : L(\theta: \mathbf{x}) = \prod \theta^{x_i} (1 - \theta)^{1-x_i} = \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i}$$

2) If we only observe the total number of successes, $S = \sum x_i$

$$S / \theta \sim \text{Binom}(\theta) : L(\theta: \mathbf{x}) = \binom{n}{s} \theta^s (1 - \theta)^{n-s} \propto \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i}$$

3) Suppose that we sample until we observe 's' number of successes.

(i.e. We do not fix the total number of trials = n). $S = \sum x_i$

$$N | \theta \sim \text{Negative Binom}(\theta) : L(\theta: \mathbf{x}) = \binom{n-1}{s-1} \theta^s (1 - \theta)^{n-s}$$

- In Frequentist approach, different sampling schemes yield **different probability models** (\Rightarrow different estimates for θ in general)
- But, they all give **the same likelihood function \Rightarrow same posterior**

Inference for Binary data with Beta prior

- Likelihood : $f(\mathbf{x}|\theta) \propto \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i}$ and Prior: $\theta \sim \text{Beta}(\alpha, \beta)$

$$\text{Posterior: } f(\theta|\mathbf{x}) \propto f(\mathbf{x}|\theta)f(\theta) \propto \theta^{\sum x_i + \alpha - 1} (1 - \theta)^{n - \sum x_i + \beta - 1}$$
$$\Rightarrow \theta|\mathbf{X} \sim \text{Beta}$$

- Posterior mean:

$$E[\theta | \mathbf{X}] =$$

Posterior Variance:

$$\text{Var}[\theta | \mathbf{X}] =$$

Inference for Binary data with Beta prior

- Likelihood : $f(\mathbf{x}|\theta) \propto \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i}$ and Prior: $\theta \sim \text{Beta}(\alpha, \beta)$

$$\text{Posterior: } f(\theta|\mathbf{x}) \propto f(\mathbf{x}|\theta)f(\theta) \propto \theta^{\sum x_i + \alpha - 1} (1 - \theta)^{n - \sum x_i + \beta - 1}$$

$$\Rightarrow \theta|\mathbf{X} \sim \text{Beta}(\sum x_i + \alpha, n - \sum x_i + \beta)$$

- Posterior mean:

$$E[\theta | \mathbf{X}] = \frac{\sum x_i + \alpha}{n + \alpha + \beta} = \frac{n \frac{\sum x_i}{n} + \frac{\alpha}{\alpha + \beta} (\alpha + \beta)}{n + \alpha + \beta} : \text{weighted average of } \bar{X} \text{ \& } E[\theta]$$

- If $n \rightarrow \infty$ & α, β fixed, $E[\theta | \mathbf{X}] \rightarrow \bar{X}$ (frequentist estimate).
- If $\alpha, \beta \rightarrow \infty$ & $\frac{\beta}{\alpha}$ fixed, $E[\theta | \mathbf{X}] \rightarrow \frac{\alpha}{\alpha + \beta}$ (prior mean)
- If $\alpha = \beta = 0$, prior distⁿ = $\text{Beta}(0, 0) = \text{Unif}(0, 1)$, $E[\theta | \mathbf{X}] \rightarrow \bar{X}$

EXAMPLE : Cancer trials

- A new treatment protocol is proposed for a particular form of cancer. The measure of success will be the proportion of patients that survive longer than six months after diagnosis.
- With the present treatment, this success rate is 40%. Letting θ be the success rate of the new treatment, a doctor assesses her prior beliefs about θ as follows. She judges that her expectation of θ is $E(\theta) = 0.45$, and her standard deviation is 0.07. Assume that her beliefs can be represented by a beta distribution.
- A clinical trial of the new treatment protocol is carried out. Out of 70 patients in the trial, 34 survive beyond six months from diagnosis. Is the new treatment protocol better?

Example (analysis)

- $X|\theta \sim \text{Binom}(n, \theta)$: $n = 70, x = 34$

- Prior: $\theta \sim \text{Beta}(\alpha, \beta)$ $E[\theta] = \frac{\alpha}{\alpha+\beta}$ $\text{Var}[\theta] = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$

Doctor's experience: $0.45 = \frac{\alpha}{\alpha+\beta}$ $0.07^2 = \underbrace{\frac{\alpha}{(\alpha+\beta)}}_{0.45} \underbrace{\frac{\beta}{(\alpha+\beta)}}_{0.55} \underbrace{\frac{1}{(\alpha+\beta+1)}}$

$$\Rightarrow \alpha = 0.45 \times 49.51 \approx 22.28 \quad \beta = 27.23$$

- Posterior: $\theta|X \sim \text{Beta}(x+\alpha, n-x+\beta) = \text{Beta}(56.28, 63.23)$
- $P[\theta > 0.4 | X] = 0.941$ (posterior prob.)
 $P[\theta > 0.4] = 0.758$ (prior prob.) In R, `>1-pbeta(0.4, 22.28, 27.23)`

Examples

4) Calculate the integral by identifying a kernel of a pdf

$$\int_0^{\infty} \theta^{\alpha+x-1} e^{-\theta(2\beta)} d\theta$$

- Identify the probability distribution for this kernel of density.

5) $f(x|\theta, \tau) \propto \exp(-\tau x^2 + \theta x) \quad -\infty < x < \infty, \quad \theta, \tau > 0$

6) $\exp(-2\phi^2 + 12\phi\psi)$

4)

$$\frac{1}{\Gamma(\alpha+x) \left[\frac{1}{2\beta}\right]^{\alpha+x}} \int_0^{\infty} \theta^{\alpha+x-1} e^{-\theta(2\beta)} d\theta = 1 \Rightarrow \int_0^{\infty} \theta^{\alpha+x-1} e^{-\theta(2\beta)} d\theta = \Gamma(\alpha+x) \left[\frac{1}{2\beta}\right]^{\alpha+x}$$

$$\begin{aligned} 5) \quad f(x|\theta, \tau) &\propto \exp(-\tau x^2 + \theta x) = \exp\left\{-\tau \left(x^2 - \frac{\theta x}{\tau}\right)\right\} \propto \exp\left\{-\tau \left(x^2 - \frac{\theta x}{\tau} + \left(\frac{\theta}{2\tau}\right)^2\right)\right\} \\ &= \exp\left\{-\tau \left(x - \frac{\theta}{2\tau}\right)^2\right\} = \exp\left\{-\frac{\frac{1}{\tau}}{2} \frac{1}{2}\right\} \quad : \text{kernel of } N\left(\frac{\theta}{2\tau}, \frac{1}{2\tau}\right) \end{aligned}$$

$$\text{- Note: } \left(x - \frac{\theta}{2\tau}\right)^2 = x^2 - \frac{\theta}{\tau} x + \frac{\theta^2}{4\tau^2}$$

$$\begin{aligned} 6) \quad f(\phi) &\propto \exp(-2\phi^2 + 12\phi\psi) = \exp\{-2(\phi^2 - 6\phi\psi)\} \\ &\propto \exp\{-2(\phi^2 - 6\phi\psi + (3\psi)^2)\} = \exp\{-2(\phi - 3\psi)^2\} \\ &= \exp\left\{-\frac{\frac{1}{2}}{2} \frac{1}{2}\right\} \quad : \text{kernel of Normal}(3\psi, 1/4) \end{aligned}$$

Conjugacy

- A **conjugate prior** is the one for which the **prior distribution and the posterior distribution have the same family of distributions** (same functional form), just with different (updated) parameters.
 - For example, in the Beta-binomial model, the prior is a Beta and the posterior is also a Beta, so this is a conjugate prior.
- Conjugate priors are nice because
 - we can typically derive the posterior without needing any difficult computation (of normalizing constants)
 - it is usually easy to understand the respective contributions of the prior information and the data information to the posterior.
- Examples: Poisson-Gamma, Exponential-Gamma, Normal-Normal, etc.

Example: Poisson-Gamma

- The Poisson distribution is a common model for count data, characterized by a rate λ (expected occurrences per unit time, volume, or other unit), and the number of occurrences x (taking on integers).
- If our data consists of a random sample on n such counts, then the likelihood function is the joint density function $f(\mathbf{x}|\lambda) = f(x_1|\lambda)f(x_2|\lambda)\cdots f(x_n|\lambda)$, since X_1, \dots, X_n are independent.
- **Likelihood:** $X_1, \dots, X_n / \lambda \sim \text{ind. Poisson}(\lambda)$ i.e. $f(x_i|\lambda) = \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$

$$L(\lambda: \mathbf{x}) = \prod \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} = \frac{\lambda^{x_1} \cdots \lambda^{x_n} e^{-\lambda} \cdots e^{-\lambda}}{x_1! \cdots x_n!} = \frac{\lambda^{\sum_{i=1}^n x_i} e^{-n\lambda}}{\prod x_i!}$$

- The Gamma distribution is a good choice for the prior, since its support is $(0, \infty)$ as $\lambda > 0$.

- **Prior:** $\lambda \sim \text{Gamma}(\alpha, \beta), \quad \lambda > 0$

$$f(\lambda) = \frac{1}{\Gamma(\alpha)\beta^\alpha} \lambda^{\alpha-1} \exp\left(-\frac{\lambda}{\beta}\right) \quad E[\lambda] = \alpha\beta \quad \text{Var}[\lambda] = \alpha\beta^2$$

$$\text{(Then, } \int_0^\infty \lambda^{\alpha-1} \exp\left(-\frac{\lambda}{\beta}\right) d\lambda = \Gamma(\alpha)\beta^\alpha : \text{1/N.C.)}$$

- The parameterization of the Gamma distribution that we will use in this course is a Gamma pdf with a **shape** parameter α and a **scale** parameter β .
- Note that the scale parameter is the **reciprocal of the rate** parameter used in the other parameterization.

Posterior distribution with conjugacy

- **Posterior:** $f(\lambda|\mathbf{x}) = \frac{f(\mathbf{x}|\lambda) f(\lambda)}{\int f(\mathbf{x}|\lambda) f(\lambda) d\lambda} \propto f(\mathbf{x}|\lambda) f(\lambda)$

$$f(\lambda|\mathbf{x}) = \frac{\lambda^{\sum x_i} e^{-n\lambda}}{\prod x_i!} \frac{1}{\Gamma(\alpha)\beta^\alpha} \lambda^{\alpha-1} \exp\left(\frac{-\lambda}{\beta}\right)$$

$$\propto \lambda^{A-1} \exp\left(-\frac{\lambda}{B}\right)$$

: kernel of ?

Thus, $\lambda|\mathbf{x} \sim ?$

Posterior distribution with conjugacy

- **Posterior:** $f(\lambda|\mathbf{x}) = \frac{f(\mathbf{x}|\lambda) f(\lambda)}{\int f(\mathbf{x}|\lambda) f(\lambda) d\lambda} \propto f(\mathbf{x}|\lambda) f(\lambda)$

$$f(\lambda|\mathbf{x}) = \frac{\lambda^{\sum x_i} e^{-n\lambda}}{\prod x_i!} \frac{1}{\Gamma(\alpha)\beta^\alpha} \lambda^{\alpha-1} \exp\left(\frac{-\lambda}{\beta}\right)$$

$$\propto \lambda^{\sum x_i + \alpha - 1} \exp\left(-n\lambda - \lambda \frac{1}{\beta}\right)$$

$$= \lambda^{\sum x_i + \alpha - 1} \exp(-\lambda(n + 1/\beta)) : \text{kernel of Gamma}(\sum x_i + \alpha, \frac{1}{n + \frac{1}{\beta}})$$

Thus, $\lambda|\mathbf{x} \sim \text{Gamma}(\sum x_i + \alpha, \frac{1}{n + \frac{1}{\beta}})$

Deriving Posterior with all constants – not needed in practice (Aside)

$$f(\mathbf{x}|\lambda)f(\lambda) = \frac{\lambda^{\sum x_i} e^{-n\lambda}}{\prod x_i!} \frac{1}{\Gamma(\alpha)\beta^\alpha} \lambda^{\alpha-1} \exp\left(-\frac{\lambda}{\beta}\right)$$

$$\begin{aligned} f(\mathbf{x}) &= \int_0^\infty f(\mathbf{x}|\lambda)f(\lambda)d\lambda = \frac{1}{\prod x_i! \Gamma(\alpha)\beta^\alpha} \int_0^\infty \lambda^{\sum x_i + \alpha - 1} e^{-\lambda(n+1/\beta)} d\lambda \\ &= \frac{1}{\prod x_i! \Gamma(\alpha)\beta^\alpha} \Gamma(\sum x_i + \alpha) \left[\frac{1}{n + \frac{1}{\beta}} \right]^{\sum x_i + \alpha} \end{aligned}$$

$$f(\lambda|\mathbf{x}) = \frac{f(\mathbf{x}|\lambda)f(\lambda)}{\int f(\mathbf{x}|\lambda)f(\lambda)d\lambda} = \frac{\frac{\lambda^{\sum x_i} e^{-n\lambda}}{\prod x_i!} \frac{1}{\Gamma(\alpha)\beta^\alpha} \lambda^{\alpha-1} \exp(-\lambda/\beta)}{\frac{1}{\prod x_i! \Gamma(\alpha)\beta^\alpha} \Gamma(\sum x_i + \alpha) \left[\frac{1}{n + \frac{1}{\beta}} \right]^{\sum x_i + \alpha}} = \frac{\lambda^{\sum x_i + \alpha - 1} \exp\left(-\lambda \left(n + \frac{1}{\beta}\right)\right)}{\Gamma(\sum x_i + \alpha) \left[\frac{1}{n + \frac{1}{\beta}} \right]^{\sum x_i + \alpha}}$$

Inference with Poisson-Gamma model

- **Posterior distribution:** $\lambda|\mathbf{x} \sim \text{Gamma}(\sum x_i + \alpha, \frac{1}{n + \frac{1}{\beta}})$
 - Posterior mean: $E[\lambda|\mathbf{x}] =$

: Posterior mean compared with \bar{X} and $E[\lambda]$?
 - Posterior variance: $\text{Var}[\lambda|\mathbf{x}] =$

Inference with Poisson-Gamma model

- **Posterior distribution:** $\lambda|\mathbf{x} \sim \text{Gamma}(\sum x_i + \alpha, \frac{1}{n + \frac{1}{\beta}})$

$$\text{Posterior mean: } E[\lambda|\mathbf{x}] = \frac{\sum x_i + \alpha}{n + \frac{1}{\beta}} = \frac{n \frac{\sum x_i}{n} + \frac{\alpha\beta}{\beta}}{n + \frac{1}{\beta}} : \text{weighted average of } \bar{X} \text{ \& } E[\lambda]$$

- As $n \rightarrow \infty$, $E[\lambda|\mathbf{x}] \rightarrow \bar{X}$ (frequentist estimate).

As $\frac{1}{\beta} \rightarrow \infty$, $E[\lambda|\mathbf{x}] \rightarrow \alpha\beta$ (prior mean)

- $\text{Var}[\bar{X}] = \frac{\lambda}{n} \Rightarrow \text{precision of } \bar{X} \propto \frac{n}{\hat{\lambda}}$ Also, precision of $\lambda \propto \frac{1}{\beta}$

- Posterior variance: $\text{Var}[\lambda|\mathbf{x}] = (\sum x_i + \alpha) \left[\frac{1}{n + \frac{1}{\beta}} \right]^2 = \frac{\sum x_i + \alpha}{\left[n + \frac{1}{\beta} \right]^2}$