

III Conjugate priors and exponential family

II (1) intro

- Often, computing the posterior density is difficult because of $D = \int_{\Theta} p_{\theta}(x) \pi(\theta) d\mu(\theta)$, which may not have an analytical expression.
- Estimating the density via numerical methods approximating D is not that useful because what matters is usually rather the mean or the posterior quantiles.
- Sometimes (as in our leading example) everything goes well: the prior is parameterized by hyper parameters (α, β in our example of a Beta(α, β) distribution).

definition hyper parameters

When the prior π is chosen in a family $\mathcal{F} = \{ \pi_{\theta}, \theta \in \Gamma \}$, with $\Gamma \subset \mathbb{R}^d$, the parameter θ characterizing π_{θ} is called "hyper-parameter".

In our example, $\Gamma = \mathbb{R}_+^* \times \mathbb{R}_+^*$ and

$\mathcal{F} = \{ \text{Beta}(\alpha, \beta), (\alpha, \beta) \in \Gamma \}$ is the family of all Beta distributions. Moreover, the posterior $\pi(\theta | x)$ belongs to the same family \mathcal{F} namely $\pi(\theta | x) = \text{Beta}(\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i)$. Such a prior is called conjugate.

Definition Conjugate prior

a family \mathcal{F} of probability distributions over Θ is "conjugate" for a likelihood $p_{\theta}(x)$ if $\forall x \in \mathcal{X}, \pi(\theta | x) \in \mathcal{F}$

III (2) examples

(a) Gaussian model with known variances
 $\Theta = \mathbb{R}, p_{\theta}(x) =$

$\propto \exp \left\{ \text{quadratic function of } \theta \right\}$
 as a function of θ .

if $\pi(\theta) \propto \exp \left\{ \text{---} \right\}$

then $\pi(\theta | x) \propto p_{\theta}(x) \pi(\theta)$
 $\propto \exp \left\{ \text{quadr. fn} \right\}$

but if $\pi(\theta) \propto \exp \left\{ \text{quadr. ---} \right\}$
 then $\pi(\theta) = \mathcal{N}(\mu, \sigma^2)$ for some μ, σ^2

details $\pi(\theta) \propto \exp\left\{-\frac{1}{2\lambda^2}(\theta - \mu)^2\right\}$, $\lambda = (\lambda_1 \dots \lambda_n)$

$$\pi(\theta|z) \propto \exp\left\{-\frac{1}{2}\left\{\sum_i \frac{(x_i - \theta)^2}{\sigma^2} + \frac{(\theta - \mu)^2}{\lambda^2}\right\}\right\}$$

$$= \exp\left\{-\frac{1}{2}\left(\theta^2\left(\frac{n}{\sigma^2} + \frac{1}{\lambda^2}\right) - 2\theta\left(\frac{\sum x_i}{\sigma^2} + \frac{\mu}{\lambda^2}\right) + C\right)\right\}$$

$$\propto \exp\left\{-\frac{1}{2\left[\frac{n}{\sigma^2} + \frac{1}{\lambda^2}\right]^{-1}}\left[\theta^2 - 2\theta\left(\frac{\sum x_i}{\sigma^2} + \frac{\mu}{\lambda^2}\right) + \frac{\lambda^2 \sigma^2 / n}{\lambda^2 + \sigma^2 / n}\right]\right\}$$

$$\propto \mathcal{N}\left(\mu_n, \sigma_n^2\right) \text{ with}$$

$$\begin{cases} \mu_n = \left(\lambda^2 \frac{1}{n} \sum x_i + \frac{\sigma^2}{n} \mu\right) \times \frac{1}{\lambda^2 + \sigma^2 / n} \\ 1/\sigma_n^2 = \frac{1}{\lambda^2} + \frac{n}{\sigma^2} \end{cases}$$

rem : maximum likelihood estimate:

$$\hat{\mu}_{ML} = \frac{1}{n} \sum x_i$$

• Posterior mean, $\mu_n = \left(\lambda^2 \hat{\mu}_{ML} + \frac{\sigma^2}{n} \mu\right) \frac{1}{\lambda^2 + \sigma^2 / n}$: weighted average (regularized version)

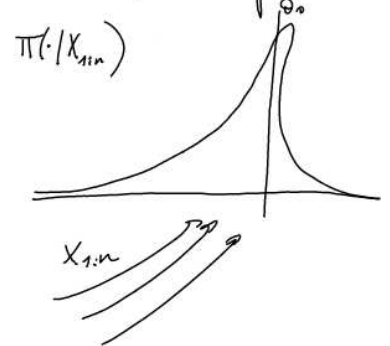
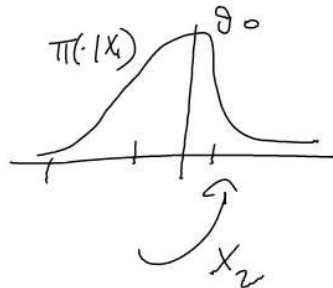
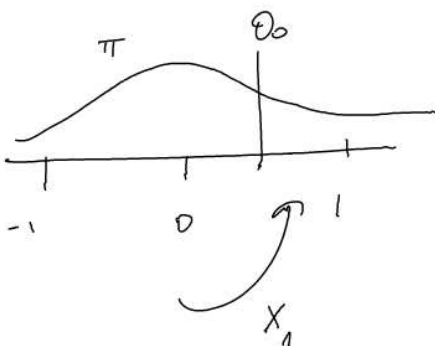
• $n = 0$: $\mu_n = \hat{\mu}_{ML}$
 $\sigma_n^2 = \lambda^2$

• $n \rightarrow \infty$: $\mu_n \sim \hat{\mu}_{ML} \rightarrow \mu_0$ (true mean)
 $\sigma_n^2 \sim \frac{\sigma^2}{n} \rightarrow 0$: posterior becomes peaked around $\hat{\mu}_{ML}$

• $\lambda^2 \rightarrow \infty$: $\sigma_n^2 \sim \frac{\sigma^2}{n}$

Illustration

: concentration of the posterior distribution around the true parameter



(b) Gaussian model, unknown variance, known mean
 Conjugate prior? Working with $d = \frac{1}{\sigma^2}$.

$$p(d) \propto d^{n/2} \exp\left\{-\frac{d}{2} \sum_{i=1}^n (x_i - \mu)^2\right\}$$

$$\propto d^{n/2} \exp\left\{-\frac{d}{2} S^2\right\} \quad S^2 = \sum_{i=1}^n (x_i - \mu)^2$$

remind the Gamma distribution: ($a, b > 0$)

$$\mathcal{G}_{a,b}(d) = \frac{1}{\Gamma(a)} d^{a-1} \exp\{-bd\} : \quad \mathbb{E}_{ab}(d) = \frac{a}{b}, \quad \text{Var}_{ab}(d) = \frac{a}{b^2}$$

thus $p(d|x) \propto \mathcal{G}\left(\frac{n}{2} + 1, \frac{S^2}{2}\right)$

take $\pi(d) = \mathcal{G}_{a,b}(d)$

$$\text{then } \pi(d|x) \propto d^{a-1 + \frac{n}{2}} \exp\left\{-d\left(\frac{S^2}{2} + b\right)\right\}$$

$$= \mathcal{G}\left(a + \frac{n}{2}, b + \frac{S^2}{2}\right)(d)$$

$$\Rightarrow \mathbb{E}(d|x) = \frac{a + n/2}{b + S^2/2} \quad \rightarrow \quad \hat{d} = \frac{a + n/2}{b + S^2/2}$$

$$\text{Var}(d|x) = \frac{a + n/2}{(b + S^2/2)^2}$$

Remark: $\sigma_{PL}^2 = \frac{S^2}{n}$ if we take

$$\hat{\sigma} = \frac{1}{\hat{\mu}_{\text{mean}}}, \text{ we get } \hat{\sigma} = \frac{b + S^2/2}{a + n/2} \rightarrow \text{regularized version of } \hat{\sigma}_{PL}$$

(c) Gaussian model, mean and variance unknown.

$$\theta = (\mu, d) \quad d = \frac{1}{\sigma^2}$$

$$P(\mu, d) (x) = \prod_{i=1}^n \left(\frac{d}{2\pi} \right)^{\frac{1}{2}} \exp \left\{ -\frac{d}{2} (x_i - \mu)^2 \right\}$$

$$= \left[d^{\frac{1}{2}} \exp \left\{ -\frac{d\mu^2}{2} \right\} \right]^n \exp \left\{ d\mu \sum_{i=1}^n x_i - \frac{d}{2} \sum_{i=1}^n x_i^2 \right\}$$

We need a prior $\pi(\mu, d)$ with a similar form

$$\pi(\mu, d) \propto \left[d^{\frac{1}{2}} \exp \left\{ -\frac{d\mu^2}{2} \right\} \right]^\beta \exp \left\{ c d \mu - d d \right\}$$

$$= d^{\beta/2} \exp \left\{ -\frac{d\beta}{2} \left(\mu - \frac{c}{\beta} \right)^2 \right\} \exp \left\{ -d \left(d - \frac{c^2}{2\beta} \right) \right\}$$

writes $\pi(\mu, d) = \pi_1(\mu|d) \pi_2(d)$

with $\pi_1(\mu|d)$ gaussian $\left(\mu_0 = \frac{c}{\beta} \right)$
 $1/\sigma_0^2 = \beta d$

$$\pi_2(d) = \text{Gamma} \left(\underbrace{\frac{\beta}{2} + 1}_a, \underbrace{d - \frac{c^2}{2\beta}}_b \right)$$

"Normal-Gamma model":

$$\left\{ \begin{array}{l} d \sim \text{Gamma}(a, b) \\ \mu|d \sim \mathcal{N} \left(\mu_0, \sigma_0^2 = \frac{1}{d\beta} \right) \end{array} \right.$$

→ "Normal-Gamma" (μ_0, β, a, b)

Posterior: $a_N = a + \frac{N}{2}$; $b_N = b + \frac{1}{2} \left\{ S^2 + \beta \frac{n(\bar{x} - \mu_0)^2}{\beta + n} \right\}$; $\mu_N = \frac{\beta \mu_0 + n \bar{x}}{\beta + n}$; $\beta_N = \beta + n$.

$$S^2 = \sum (x_i - \bar{x})^2; \quad \bar{x} = \frac{1}{n} \sum x_i$$

(d) Multivariate case: conjugate priors $X \sim \mathcal{N}(\mu, \Sigma) \in \mathbb{R}^D$

(i) unknown mean: → $\Pi(\mu)$: multivariate gaussian is conjugate

(ii) unknown covariance:

$$L = \Sigma^{-1} \quad W(L) = B |L|^{(D-D-1)/2} \exp \left\{ -\frac{1}{2} \text{Tr}(W^{-1}L) \right\}$$

"Wishart distribution with

ν degrees of freedom $\in \mathbb{R}^+$

$W \in \mathbb{R}^{D \times D}$; $B = \text{normalization const}$