

Bayes Estimation

Outline

- 1) Bayes risk, Bayes estimator
- 2) Examples
- 3) Conjugate priors

Frequentist Motivation

Model $\mathcal{P} = \{P_\theta : \theta \in \Sigma\}$ for data X

Notation change: for now, reserving Θ to denote random variable

Loss $L(\theta, d)$, Risk $R(\theta; \delta)$

The Bayes risk is the average-case risk, integrated wrt some measure Λ , called prior

For now, assume $\Lambda(\Sigma) = 1$ (prob. meas.)

Later we will allow to be improper ($\Lambda(\Sigma) = \infty$)

(Note Λ and $c\Lambda$ for $c > 0$ functionally equiv.)

$$R_{\text{Bayes}}(\Lambda, \delta) = \int_{\Sigma} R(\theta, \delta) d\Lambda(\theta)$$

$$= \mathbb{E} R(\Theta; \delta) \quad \text{where } \Theta \sim \Lambda$$

$$= \mathbb{E} L(\Theta, \delta(x)) \quad \begin{array}{l} \Theta \sim \Lambda \\ X | \Theta = \theta \sim P_\theta \end{array}$$

An estimator δ minimizing $R_{\text{Bayes}}(\Lambda; \cdot)$ is called Bayes (a Bayes estimator).

Depends on \mathcal{P} , Λ , L

Note avg-case loss makes sense even if we don't "believe" the parameter is random.

Bayes Estimator

Thm

Suppose $\Theta \sim \Lambda$

$$X | \Theta = \theta \sim P_\theta$$

$$L(\theta, d) \geq 0 \quad \forall \theta, d$$

$$R_{\text{Bayes}}(\Lambda; \delta_0) < \infty \quad \text{for some } \delta_0(x)$$

Then $\delta_\Lambda(x) \in \operatorname{argmin}_d \mathbb{E}[L(\Theta, d) | X=x]$ a.e. x

$\Leftrightarrow \delta_\Lambda(x)$ is Bayes with $R_{\text{Bayes}}(\Lambda; \delta_\Lambda) < \infty$

Interpretation: we can find the Bayes estimator by minimizing $\mathbb{E}[L(\Theta, d) | X=x]$ "one x at a time."

Proof (\Rightarrow) Let δ be any other estimator

$$\begin{aligned} R_{\text{Bayes}}(\Lambda; \delta) &= \mathbb{E} L(\Theta, \delta(X)) \\ &= \mathbb{E} \left[\mathbb{E} [L(\Theta, \delta(X)) \mid X=x] \right] \\ &\geq \mathbb{E} \left[\mathbb{E} [L(\Theta, \delta_{\Lambda}(x)) \mid X=x] \right] \\ &= R_{\text{Bayes}}(\Lambda, \delta_{\Lambda}) \\ &< \infty \quad \text{if } \delta = \delta_0 \end{aligned}$$

(\Leftarrow) Let $E_x(d) = \mathbb{E} [L(\Theta, d) \mid X=x]$

Define

$$\delta^*(x) = \begin{cases} \delta_{\Lambda}(x) & \text{if } \delta_{\Lambda}(x) \in \arg \min E_x \\ \delta_0(x) & \text{if } E_x(\delta_0(x)) < E_x(\delta_{\Lambda}(x)) \\ d^*(x) & \text{otherwise} \end{cases}$$

where $E_x(d^*) < E_x(\delta_{\Lambda}(x))$

Then $E_x(\delta^*(x)) \stackrel{\text{a.s.}}{\leq} E_x(\delta_0(x))$

and $E_x(\delta^*(x)) \stackrel{\text{a.s.}}{\leq} E_x(\delta_{\Lambda}(x))$

with ineq. strict on a set of measure > 0 .

□

Prior, Posterior

Usual interp. of Δ is prior belief about \textcircled{H} before seeing the data

Epistemic uncertainty: "I think there is a 50% chance that..."

[Mathematically, it makes no more or less sense to take Θ as fixed or random, but a matter of philosophy whether this is scientifically appropriate]
More on this later...

Conditional dist. of \textcircled{H} given X , which we will write $\mathcal{Q}(\textcircled{H} | X)$, ^(law) is called the posterior distribution, aka our beliefs after seeing the data.

Densities: prior $\lambda(\theta)$, likelihood $p_{\theta}(x)$

$\Rightarrow q(x) = \int_{\Omega} \lambda(\theta) p_{\theta}(x) d\theta$ marginal density of x

$\lambda(\theta | x) = \frac{\lambda(\theta) p_{\theta}(x)}{q(x)}$ posterior density

Bayes est: $\delta_{\lambda}(x) = \operatorname{argmin}_d \int_{\Omega} L(\theta, d) \lambda(\theta | x) d\theta$

Posterior Mean

If $L(\theta, d) = (g(\theta) - d)^2$ then the Bayes estimator is the posterior mean:

$$\begin{aligned}\mathbb{E}[(g(\Theta) - d)^2 | X] \\&= \mathbb{E}[(g(\Theta) - \mathbb{E}[g(\Theta) | X] + \mathbb{E}[g(\Theta) | X] - d)^2 | X] \\&= \text{Var}(g(\Theta) | X) + (\mathbb{E}[g(\Theta) | X] - d)^2\end{aligned}$$

(why is the cross-term 0?)

$$\Rightarrow \delta_{\Delta}(x) = \mathbb{E}[g(\Theta) | X=x]$$

Weighted sq. error:

$$L(\theta, d) = w(\theta) (g(\theta) - d)^2$$

e.g. $(\frac{\theta-d}{\theta})^2$
sq. rel. error

$$\begin{aligned}\mathbb{E}[(d - g(\Theta))^2 w(\Theta) | X] \\&= d^2 \mathbb{E}[w(\Theta) | X] - 2d \mathbb{E}[w(\Theta)g(\Theta) | X] \\&\quad + \mathbb{E}[\cancel{w(\Theta)g(\Theta)^2} | X]\end{aligned}$$

no dep. on d

$$\text{min at } d = \frac{\mathbb{E}[w(\Theta)g(\Theta) | X]}{\mathbb{E}[w(\Theta) | X]} \quad (= \delta_{\Delta}(x))$$

Example: Beta-Binomial

$$X | \Theta = \theta \sim \text{Binom}(n, \theta) = \theta^x (1-\theta)^{n-x} \binom{n}{x}$$

$$\Theta \sim \text{Beta}(\alpha, \beta) = \theta^{\alpha-1} (1-\theta)^{\beta-1} \underbrace{\frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha+\beta)}}_{\text{normalizing const.}}$$

θ is r.v. here

Marginal dist. of X called Beta-Binomial

Posterior:

$$\lambda(\theta | x) = \lambda(\theta) p_{\theta}(x) / q(x)$$

we can drop factors that don't depend on θ \rightarrow

$$\propto_{\theta} \theta^{\alpha-1} (1-\theta)^{\beta-1} \theta^x (1-\theta)^{n-x}$$
$$= \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}$$

$$\Rightarrow \Theta | X=x \sim \text{Beta}(x+\alpha, n-x+\beta)$$

$$\mathbb{E}[\Theta | X] = \frac{x+\alpha}{n+\alpha+\beta}$$

convex combo of $\frac{x}{n}$, $\frac{\alpha}{\alpha+\beta}$ \rightarrow

$$\overset{\text{UMVU}}{=} \frac{x}{n} \cdot \frac{n}{n+\alpha+\beta} + \underbrace{\frac{\alpha}{\alpha+\beta}}_{\text{Prior Expectation}} \cdot \underbrace{\frac{\alpha+\beta}{n+\alpha+\beta}}_{\left(1 - \frac{n}{n+\alpha+\beta}\right)}$$

Interp.: $k = \alpha + \beta$ "pseudo-trials," α successes

(Recall $\frac{x+3}{n+6}$ from Lec. 2)

Example: Normal mean

$$X | \Theta = \theta \sim N(\theta, \sigma^2) \propto_{\theta} e^{-(x-\theta)^2/2\sigma^2}$$

$$\Theta \sim N(\mu, \tau^2) \propto_{\theta} e^{-(\theta-\mu)^2/2\tau^2}$$

$$\lambda(\theta | x) \propto_{\theta} \exp \left\{ -\frac{(x-\theta)^2}{2\sigma^2} - \frac{(\theta-\mu)^2}{2\tau^2} \right\}$$

$$\propto_{\theta} \exp \left\{ \frac{x\theta}{\sigma^2} - \frac{\theta^2}{2\sigma^2} - \frac{\theta^2}{2\tau^2} + \frac{\theta\mu}{\tau^2} \right\}$$

$$= \exp \left\{ \underbrace{\theta \left(\frac{x}{\sigma^2} + \frac{\mu}{\tau^2} \right)}_b - \underbrace{\theta^2 \left(\frac{\sigma^{-2} + \tau^{-2}}{2} \right)}_{a^2} \right\}$$

Complete square:

$$a^2 \theta^2 - b\theta = \left(\theta a - \frac{b}{2a} \right)^2 - c(a, b)$$

$$= \left(\theta - \frac{b}{2a^2} \right)^2 \cdot a^2 - c$$

$$\rightarrow \propto_{\theta} \exp \left\{ - \left(\theta - \frac{x\sigma^{-2} + \mu\tau^{-2}}{\sigma^{-2} + \tau^{-2}} \right)^2 / 2(\sigma^{-2} + \tau^{-2})^{-1} \right\}$$

$$\propto_{\theta} N \left(\underbrace{\frac{x\sigma^{-2} + \mu\tau^{-2}}{\sigma^{-2} + \tau^{-2}}}_{\text{precision-weighted average of } x, \mu}, \underbrace{\frac{1}{\sigma^{-2} + \tau^{-2}}}_{\text{harmonic mean of } \sigma^2, \tau^2} \right)$$

precision-weighted
average of x, μ

harmonic mean
of σ^2, τ^2

$$\mathbb{E}[\Theta | x] = x \cdot \frac{\sigma^{-2}}{\sigma^{-2} + \tau^{-2}} + \mu \cdot \frac{\tau^{-2}}{\sigma^{-2} + \tau^{-2}}$$

Gaussian iid sample

$$\theta \sim N(\mu, \tau^2) \quad (\text{dropping } \textcircled{+} \text{ for computations})$$

$$X_i | \theta \stackrel{\text{iid}}{\sim} N(\theta, \sigma^2), \quad i = 1, \dots, n$$

$$\bar{X} | \theta \sim N(\theta, \frac{\sigma^2}{n})$$

$$\Rightarrow \mathbb{E}[\theta | X] = X \cdot \frac{n\sigma^{-2}}{n\sigma^{-2} + \tau^{-2}} + \mu \cdot \frac{\tau^{-2}}{n\sigma^{-2} + \tau^{-2}}$$

$$= X \cdot \frac{n}{n + \sigma^2/\tau^2} + \mu \cdot \frac{\sigma^2/\tau^2}{n + \sigma^2/\tau^2}$$

Interp: $k = \sigma^2/\tau^2$ pseudo-observations, mean μ

If $n \gg k$, "data swamps prior"

If $n \ll k$, "prior swamps data"

Note in both examples:

- Prior & Likelihood have similar fcn. form
- Posterior comes from same exp. fcn. as prior
- Prior can be interp. as "pseudo-obs."

This is because we chose conjugate priors
for the parameter; topic for next lecture.

Bias and Bayes

[Bayes estimators are almost always biased, especially when the parameter value is extreme relative to the prior]

Theorem The posterior mean is biased unless

$$\delta_{\Delta}(x) \stackrel{\text{a.s.}}{=} g(\Theta)$$

Proof Suppose δ_{Δ} is unbiased. Then

$$g(\Theta) = \mathbb{E}[\delta(x) | \Theta]$$

$$\delta_{\Delta}(x) = \mathbb{E}[g(\Theta) | x] \quad (\text{by def.})$$

Condition on x :

$$\begin{aligned} \mathbb{E}[\delta_{\Delta}(x) g(\Theta) | x] &= \delta_{\Delta}(x) \mathbb{E}[g(\Theta) | x] \\ &= \delta_{\Delta}(x)^2 \end{aligned}$$

Cond. on Θ :

$$\mathbb{E}[\delta_{\Delta}(x) g(\Theta) | \Theta] = g(\Theta)^2$$

$$\Rightarrow \mathbb{E}[\delta_{\Delta} g(\Theta)] = \mathbb{E}[\delta_{\Delta}^2] = \mathbb{E}[g(\Theta)^2]$$

$$\begin{aligned} \mathbb{E}[(\delta_{\Delta} - g(\Theta))^2] &= \mathbb{E} \delta_{\Delta}^2 + \mathbb{E} g(\Theta)^2 - 2 \mathbb{E}[\delta_{\Delta} g(\Theta)] \\ &= 0 \quad \square \end{aligned}$$