

Recap: decision theory.

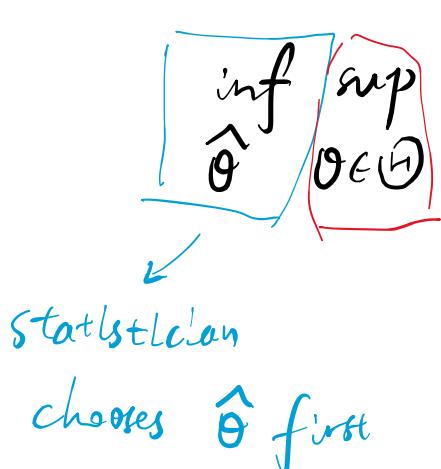
- Bayes criteria.

$$r_\pi(\hat{\theta}) = \int_{\Theta} R(\theta; \hat{\theta}) \pi(\theta) d\theta.$$

- Minimax criteria.

$$\bar{R}(\hat{\theta}) = \sup_{\theta \in \Theta} R(\theta; \hat{\theta})$$

"minimax optimal":



$R(\theta, \hat{\theta})$
nature (adversary)
choose θ after
observing $\hat{\theta}$.

Optimal $\hat{\theta}$
is from
equilibrium
of
this game

Analogy to rock-scissors-paper game:

If we randomize, adversary cannot make advantage out of it.

Basic results from game theory:

Alice choose $a \in A$

Bob choose $b \in B$.

$$\sup_{b \in B} \inf_{a \in A} R(a, b) \leq$$

$$\inf_{a \in A} \sup_{b \in B} R(a, b)$$

Alice chooses after Bob

Alice chooses before Bob.

. von Neumann's minimax theorem

$$\sup_{\pi_b} \inf_{\pi_a} \mathbb{E}_{\substack{a \sim \pi_a \\ b \sim \pi_b}} [R(a, b)]$$

(under some

regularity

conditions).

$$\inf_{\pi_a} \sup_{\pi_b} \mathbb{E}_{\substack{a \sim \pi_a \\ b \sim \pi_b}} [R(a, b)]$$

For minimax estimation.

$$\inf_{\hat{\theta}} \sup_{\pi} r_{\pi}(\hat{\theta}) = \sup_{\pi} \left[\inf_{\hat{\theta}} r_{\pi}(\hat{\theta}) \right]$$

optimal Bayes risk
under prior π .

Implication: under certain regularity conditions,
 minimax estimator is the Bayes estimator
 under "least favorable prior".

Simple special cases.

Constant-risk Bayes estimator $\hat{\theta}$ is minimax optimal.

Proof: Suppose $\tilde{\theta}$ is another estimator

$$\sup_{\theta \in \Theta} R(\theta; \tilde{\theta}) \geq \int_{\Theta} R(\theta; \tilde{\theta}) \pi(\theta) d\theta$$

$$\geq \int_{\Theta} R(\theta; \hat{\theta}) \pi(\theta) d\theta$$

(since $\hat{\theta}$ is Bayes optimal under π).

$$= \sup_{\theta \in \Theta} R(\theta; \hat{\theta}).$$

e.g. $X \sim \text{Binom}(n, p)$, $R(p; \hat{p}) = \mathbb{E}_p[(p - \hat{p})^2]$.

Natural choice: $\hat{p} = \frac{X}{n}$. $X = \sum_{i=1}^n Z_i$

$$R(p; \hat{p}) = \frac{\text{var}_p(Z_i)}{n} = \frac{p(1-p)}{n}$$

$$\bar{R}(\hat{p}) = \frac{1}{4n}, \text{ achieved at } p = \frac{1}{2}.$$

To construct prior distribution, we use conjugate families.

Conjugate family:

Suppose we work with $(P_\theta : \theta \in \Theta)$. $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} P_\theta$

Let \mathcal{F} to be a class of prior distributions

We call it conjugate if $\forall \pi \in \mathcal{F}$

$$\pi(\cdot | x_1, \dots, x_n) \in \mathcal{F}.$$

(useful for computational convenience).

Beta-Binomial conjugacy.

If $\pi = \text{Beta}(\alpha, \beta)$

then $\pi(\cdot | x) = \text{Beta}(\alpha + x, \beta + n - x)$.

(interpreted as "pseudo-counts").

Posterior mean (i.e. Bayes optimal estimator)

$$\hat{p}_{n, \pi} = \frac{\alpha + x}{\alpha + \beta + n}$$

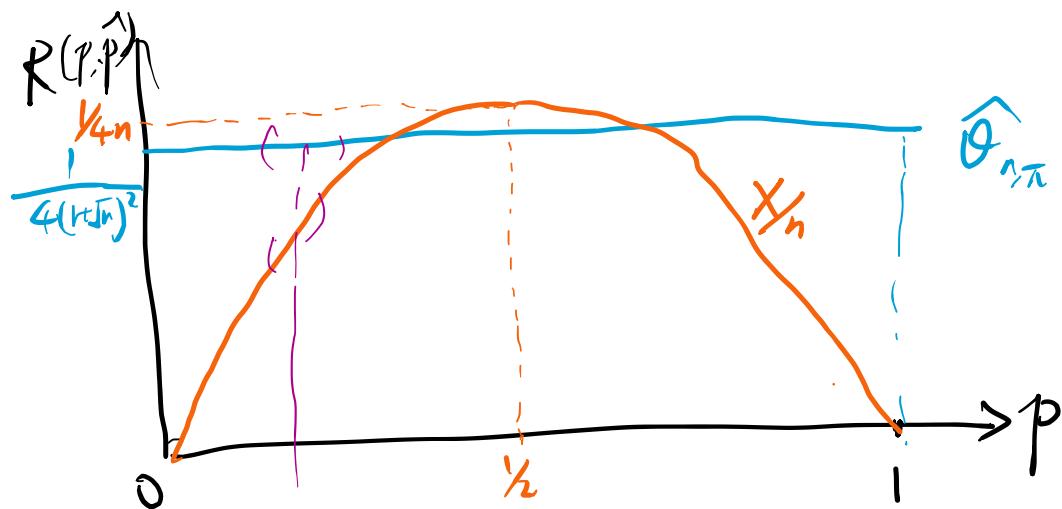
$$R(p, \hat{p}_{n,\pi}) = \frac{np(1-p) + (\alpha(1-p) - \beta p)^2}{(\alpha + \beta + n)^2}$$

Want it to be indep of p .

$$\text{Solve for } (\alpha, \beta) \Rightarrow \begin{cases} \alpha = \sqrt{n}/2 \\ \beta = \sqrt{n}/2 \end{cases}$$

$$\bar{R}(\hat{\theta}_n) = r_n(\hat{\theta}_n) = \frac{1}{4(1+\sqrt{n})^2}$$

Does it mean that we should use $\frac{\sqrt{n}}{2}$ pseudo-count?



This shows limitation of minimax criteria.

(In literature, people considered more refined optimality notion, e.g. local minimax)

Extension of this "simple case"

Fact: Suppose that we have a sequence of priors $(\pi_j)_{j=1}^{+\infty}$

r_j = optimal Bayes risk under π_j

If \exists an estimator $\hat{\theta}$

$$\text{st. } \lim_{j \rightarrow +\infty} r_j \geq \sup_{\theta} R(\theta; \hat{\theta})$$

then $\hat{\theta}$ is minimax optimal.

Proof: consider any other estimator $\tilde{\theta}$

$$\begin{aligned} \sup_{\theta \in \Theta} R(\theta; \tilde{\theta}) &\geq \lim_{j \rightarrow +\infty} \int_{\Theta} R(\theta; \tilde{\theta}) \pi_j(\theta) d\theta \\ &\geq \lim_{j \rightarrow +\infty} r_j \geq \sup_{\theta} R(\theta, \hat{\theta}). \end{aligned}$$

e.g. $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\theta, 1)$

Conjugate prior $\pi_j = \mathcal{N}(\theta, j^2)$

$$\pi(\cdot | x_1, \dots, x_n) = \mathcal{N}\left(\frac{j^2 n \cdot \bar{x}_n}{j^2 n + 1}, \frac{j^2}{j^2 n + 1}\right).$$

$$r_j = \frac{j^2}{j^2 n + 1} \xrightarrow{j \rightarrow +\infty} \frac{1}{n}.$$

We also know $\sup_{\theta \in \mathbb{R}} \mathbb{E}_{\theta} [|\bar{x}_n - \theta|^2] = \frac{1}{n}$

So \bar{X}_n is minimax optimal.

Admissibility.

Let $\hat{\theta}$ be an estimator, if $\exists \tilde{\theta}$

$$\text{s.t. } R(\theta; \hat{\theta}) \leq R(\theta; \tilde{\theta}) \quad (\forall \theta \in \Theta)$$

and furthermore, $\exists \theta_0 \in \Theta$, s.t.

$$R(\theta_0; \hat{\theta}) < R(\theta_0; \tilde{\theta})$$

then we call $\hat{\theta}$ inadmissible.

We call $\hat{\theta}$ admissible when not inadmissible.

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$, $\hat{\theta}_n = 0$ admissible.

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, I_d)$, $\theta \in \mathbb{R}^d$

\bar{X}_n is the natural estimator

(unbiased, minimax, MLE, ...).

Stein's phenomenon: \bar{X}_n inadmissible for $d \geq 3$.

James-Stein estimator

$$\hat{\theta}_{JS} = \left(1 - \frac{(d-2)}{n\|\bar{X}_n\|_2^2}\right) \bar{X}_n$$

It adapts the shrinkage towards 0.

$$\mathbb{E}_{\theta} \left[\| \hat{\theta}_{JS} - \theta \|_2^2 \right] = \frac{d}{n} - \mathbb{E}_{\theta} \left[\left(\frac{d}{n \| \bar{x}_n \|_2} \right)^2 \right].$$

$$< \frac{d}{n} = \mathbb{E}_{\theta} \left[\| \bar{x}_n - \theta \|_2^2 \right].$$

Fact. Bayes estimators are admissible. (for discrete case).

Fact. Under some regularity conditions, admissible estimators are Bayes.

Bayesian stats (ct'd).

Interval estimation in Bayesian setting:

Given π , $\theta | x_1, \dots, x_n \sim \pi(\cdot | x_1, \dots, x_n)$.

"Posterior interval / credible interval".

Find an interval C_n that depends on data

$$\text{s.t. } \pi(\theta \in C_n | x_1, \dots, x_n) = 1 - \alpha.$$

Why: this is different from confidence interval!
Conditionally on data, C_n is deterministic, θ is random.

Asymptotic properties for Bayes estimators.

Given prior π , suppose $X_1 \dots X_n \stackrel{iid}{\sim} P_{\theta^*}$
(We assume a frequentist model, Bayes estimator
is just a choice of method).

θ^* is determinate but unknown.

Under two conditions:

- $\pi(\theta \text{ at a local neighborhood of } \theta^*) > 0$.

[e.g. when π 's density is positive around θ^*]

- For the testing problem:

$$H_0: \theta = \theta^* \quad \text{v.s.} \quad H_1: \|\theta - \theta^*\| > \epsilon$$

\exists a test ϕ_n s.t. $n \rightarrow \infty$

type-I, type-II err $\rightarrow 0$

Thm (Schwartz). Under above two conditions,

posterior dist \xrightarrow{P} point mass at θ^* .

$$\text{i.e. } H \Sigma > 0, \quad \pi(\|\theta - \theta^*\| > \epsilon \mid X_1 \dots X_n) \xrightarrow{P} 0.$$

Asymptotic shape of posterior:

Thm (Bernstein-von-Mises).

$X_1 - X_n \stackrel{iid}{\sim} P_{\theta^*}$, under regularity conditions

(same conditions as required by

MLE asymptotic normality)

(In particular, $I(\theta^*)$ invertible).

Posterior $\approx N(\hat{\theta}_n, (nI(\theta^*))^{-1})$.

i.e. $d_{TV}\left(\pi(\cdot | X_1, \dots, X_n), N(\hat{\theta}_n, (nI(\theta^*))^{-1})\right) \xrightarrow{P} 0$.

Implication. (in 1-D).

$\pi(\cdot | X_1, \dots, X_n) \approx N(\hat{\theta}_n, (nI(\theta^*))^{-1})$.

Credible interval $C_n \approx \left[\hat{\theta}_n \pm \frac{z_{1-\alpha/2}}{\sqrt{nI(\theta^*)}}\right]$. \approx CI for θ using frequentist MLE

Sometimes we use "improper prior".

Given a density function π (may not be a pdf)
(we may have $\int \pi(x)dx = \infty$)

$$\Pi(\theta | x_1 \dots x_n) = \frac{\pi(\theta) \cdot p_\theta(x_1) \dots p_\theta(x_n)}{\int \pi(\theta') p_{\theta'}(x_1) \dots p_{\theta'}(x_n) d\theta'}$$

Can be still well-defined.

Computational methods.

Usually $\Pi(\cdot | x_1 \dots x_n)$ is difficult to compute.

Mainstream solutions:

• MCMC. Construct a Markov chain $(Z_t)_{t \geq 0}$

st. $Z_t \xrightarrow{d} \Pi(\cdot | x_1 \dots x_n) \quad (t \rightarrow \infty)$

We can simulate posterior by

similarity $(Z_t)_{t \geq 0}$.

• VI. Given a class \mathcal{Q} of tractable distributions

minimize $D_{KL}(Q || \Pi(\cdot | x_1 \dots x_n))$.

$Q \in \mathcal{Q}$