

Statistical estimation

- 7.1 Parametric distribution estimation
- 7.2 Nonparametric distribution estimation
- 7.3 Optimal detector design and hypothesis testing

김민주 이종현

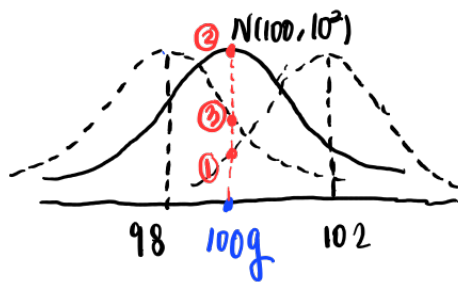
7.1 Parametric distribution estimation

01 Maximum likelihood estimation

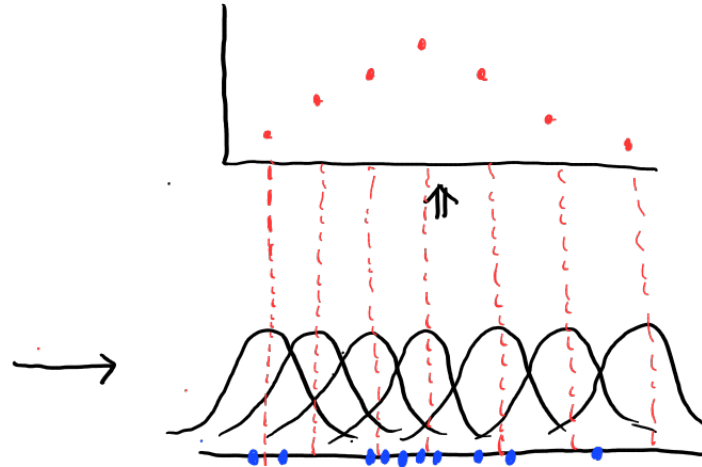
02 MLE

03 ① <확률> $P(X|D)$

② <가능도> $P(D|X)$



<likelihood>



01 Maximum likelihood estimation

02

 $x \in \mathbf{R}^n, y \in \mathbf{R}^m$ 에 대해

03

likelihood function: $p_x(y)$ log likelihood function: $l(x) = \log p_x(y)$ Maximum likelihood estimation: $\hat{x}_{ml} = \operatorname{argmax}_x p_x(x) = \operatorname{argmax}_x l(x)$ maximize $l(x) = \log p_x(y)$ subject to $x \in C$

Maximum likelihood estimation problem이 convex optimization problem이 될 조건

- $l(x)$ concave
- linear equality constraints
- convex inequality constraints

01 Linear measurements with IID noise

02
$$y_i = a_i^T x + v_i, i = 1, \dots, m$$

03 where $x \in \mathbf{R}^n$: parameter vector, $y_i \in \mathbf{R}$: observed quantitiesand v_i 's are IID

likelihood function:
$$p_x(y) = \prod_{i=1}^m p(y_i - a_i^T x)$$

log likelihood function:
$$l(x) = \log p_x(y) = \sum_{i=1}^m \log p(y_i - a_i^T x)$$

maximize
$$\sum_{i=1}^m \log p(y_i - a_i^T x)$$

 $-p_x(y)$ 가 concave면 convex optimization

: Gaussian noise

For $v_i \sim N(0, \sigma^2)$, $p(z) = (2\pi\sigma^2)^{-\frac{1}{2}} e^{-z^2/2\sigma^2}$

$$l(x) = -(m/2) \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|Ax - y\|_2^2$$

$$x_{ml} = \operatorname{argmin}_x \|Ax - y\|_2^2$$

Uniform noise

For $v_i \sim U[-a, a]$, $p(z) = 1/(2a)$, $-a < z < a$

$$x_{ml}: \|Ax - y\|_\infty \leq \underline{a} \text{를 만족하는 any } x$$

Counting problems with Poisson distribution

y : 사건 발생 횟수

$y \sim \text{Poi}(\mu)$ 를 따를 때 $\text{prob}(y = k) = \frac{e^{-\mu} \mu^k}{k!}$ 이고,

$$\mu = a^T u + b$$

(u : 설명 변수, $a \in \mathbf{R}^n$, $b \in \mathbf{R}$ 은 parameter)

a, b 의 mle를 찾으려면,

$$\prod_{i=1}^m \frac{(a^T u_i + b)^{y_i} \exp(-(a^T u_i + b))}{y_i!}$$

$$l(a, b) = \sum_{i=1}^m (y_i \log(a^T u_i + b) - (a^T u_i + b) - \log(y_i!))$$

$$\text{maximize} \quad \sum_{i=1}^m (y_i \log(a^T u_i + b) - (a^T u_i + b))$$

ex) y = 교통사고 발생 횟수

$$\mu = a_1 u_1 + a_2 u_2 + b$$

\downarrow \downarrow
 교통량 강수량

Logistic regression

$$\text{prob}(y = 1) = p, \quad \text{prob}(y = 0) = 1 - p,$$

$Y: \mathcal{Y}$

ex) $P(Y=y) = \begin{cases} p & , y=1 \\ 1-p & , y=0 \end{cases}$

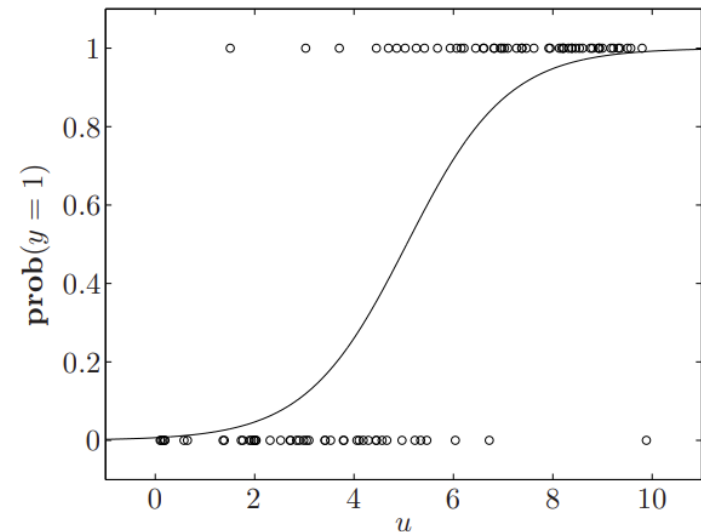
$$L(a,b) = \prod_{i=1}^q p_i \prod_{i=q+1}^m (1 - p_i)$$

$$p = \frac{\exp(a^T u + b)}{1 + \exp(a^T u + b)}$$

$$l(a,b) = \sum_{i=1}^q \log p_i + \sum_{i=q+1}^m \log(1 - p_i)$$

$$= \sum_{i=1}^q \log \frac{\exp(a^T u_i + b)}{1 + \exp(a^T u_i + b)} + \sum_{i=q+1}^m \log \frac{1}{1 + \exp(a^T u_i + b)}$$

$$= \sum_{i=1}^q (a^T u_i + b) - \sum_{i=1}^m \log(1 + \exp(a^T u_i + b)).$$



7.1.2

MAP estimation

Bayes theorem

$$P(A|B) = \frac{\overset{\text{likelihood}}{P(B|A)} \times \overset{\text{prior}}{P(A)}}{\underset{\text{posterior}}{P(B)}} \quad \leftarrow B \text{의 prior}$$

ex) X는 60% 확률로 거짓말, 거짓말 탐지기의 정확도는 90%라 가정 $\Rightarrow P(A|B)=?$
 (A) (B: 거짓말 탐지기 양성)
 (즉, $P(A)=0.6$, $P(B|A)=0.9$ 가정)

$$\Rightarrow P(\text{실제 거짓} | \text{결과 거짓}) = \frac{P(\text{결과 거짓} | \text{실제 거짓}) P(\text{실제 거짓})}{P(\text{결과 거짓})} = \frac{0.9 \times 0.6}{0.9 \times 0.6 + 0.1 \times 0.4} = \underline{0.93} \quad \leftarrow \text{posterior}$$

$\leftarrow P(\text{결과 거짓} | \text{실제 거짓}) + P(\text{결과 거짓} | \text{실제 거짓} \times)$

01

MAP estimation

02

$$p_x(x) = \int p(x, y) dy \quad p_y(y) = \int p(x, y) dx.$$

03

$$p_{y|x}(x, y) = \frac{p(x, y)}{p_x(x)}$$

$$p_{x|y}(x, y) = \frac{p(x, y)}{p_y(y)} = p_{y|x}(x, y) \frac{p_x(x)}{p_y(y)}$$

$$\begin{aligned} \hat{x}_{\text{map}} &= \operatorname{argmax}_x p_{x|y}(x, y) \\ &= \operatorname{argmax}_x p_{y|x}(x, y) p_x(x) \\ &= \operatorname{argmax}_x p(x, y). \end{aligned}$$

$$\hat{x}_{\text{map}} = \operatorname{argmax}_x (\log p_{y|x}(x, y) + \log p_x(x)).$$

01 Linear measurements with IID noise

$$y_i = a_i^T x + v_i, \quad i = 1, \dots, m,$$

$$p(x, y) = p_x(x) \prod_{i=1}^m p_v(y_i - a_i^T x)$$

$$\text{maximize} \quad \log p_x(x) + \sum_{i=1}^m \log p_v(y_i - a_i^T x)$$

ex) $v_i \sim \mathcal{U}[-a, a]$

$$\text{minimize} \quad (x - \bar{x})^T \Sigma^{-1} (x - \bar{x})$$

$$\text{subject to} \quad \|Ax - y\|_\infty \leq a,$$

7.2 Nonparametric distribution estimation

$$P(X = \alpha_k) = p_k, k=1, \dots, n$$

$$\{p \in \mathbb{R}^n \mid p \succeq 0, \mathbf{1}^T p = 1\}$$

Prior information

$$\mathbb{E} f(X) = \sum_{i=1}^n p_i f(\alpha_i)$$

$$1) \quad \text{prob}(X \in C) = c^T p, \quad c_i = \begin{cases} 1 & \alpha_i \in C \\ 0 & \alpha_i \notin C. \end{cases}$$

$$2) \quad \mathbb{E} X = \sum_{i=1}^n \alpha_i p_i = \alpha, \quad \mathbb{E} X^2 = \sum_{i=1}^n \alpha_i^2 p_i = \beta, \quad \sum_{\alpha_i \geq 0} p_i \leq 0.3,$$

$$3) \quad \text{var}(X) = \mathbb{E} X^2 - (\mathbb{E} X)^2 = \underbrace{\sum_{i=1}^n \alpha_i^2 p_i}_{\text{linear}} - \underbrace{\left(\sum_{i=1}^n \alpha_i p_i \right)^2}_{\text{quadratic}}$$

$$4) \quad \text{prob}(X \in A | X \in B) = c^T p / d^T p,$$

where

$$c_i = \begin{cases} 1 & \alpha_i \in A \cap B \\ 0 & \alpha_i \notin A \cap B \end{cases}, \quad d_i = \begin{cases} 1 & \alpha_i \in B \\ 0 & \alpha_i \notin B. \end{cases}$$

$$ld^T p \leq c^T p \leq ud^T p$$

Bounding probabilities and expected values

$$\begin{aligned} &\text{minimize} \quad \sum_{i=1}^n f(\alpha_i) p_i \\ &\text{subject to} \quad p \in \mathcal{P}. \end{aligned}$$

Maximum likelihood estimation

$$l(p) = \sum_{i=1}^n k_i \log p_i,$$

Maximum entropy

$$\begin{aligned} &\text{minimize} \quad \sum_{i=1}^n p_i \log p_i \\ &\text{subject to} \quad p \in \mathcal{P}. \end{aligned}$$

7.2 Nonparametric distribution estimation

Example 7.2 We consider a probability distribution on 100 equidistant points α_i in the interval $[-1, 1]$. We impose the following prior assumptions:

$$\begin{aligned} \mathbf{E} X &\in [-0.1, 0.1] \\ \mathbf{E} X^2 &\in [0.5, 0.6] \\ \mathbf{E}(3X^3 - 2X) &\in [-0.3, -0.2] \\ \text{prob}(X < 0) &\in [0.3, 0.4]. \end{aligned} \quad (7.8)$$

Along with the constraints $\mathbf{1}^T p = 1$, $p \succeq 0$, these constraints describe a polyhedron of probability distributions.

Figure 7.2 shows the maximum entropy distribution that satisfies these constraints. The maximum entropy distribution satisfies

$$\begin{aligned} \mathbf{E} X &= 0.056 \\ \mathbf{E} X^2 &= 0.5 \\ \mathbf{E}(3X^3 - 2X) &= -0.2 \\ \text{prob}(X < 0) &= 0.4. \end{aligned}$$

To illustrate bounding probabilities, we compute upper and lower bounds on the cumulative distribution $\text{prob}(X \leq \alpha_i)$, for $i = 1, \dots, 100$. For each value of i ,

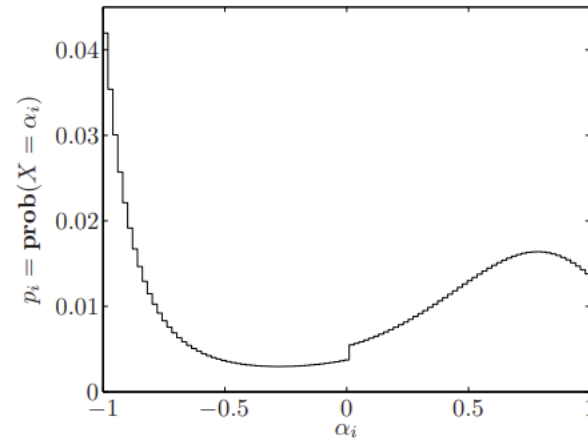


Figure 7.2 Maximum entropy distribution that satisfies the constraints (7.8).

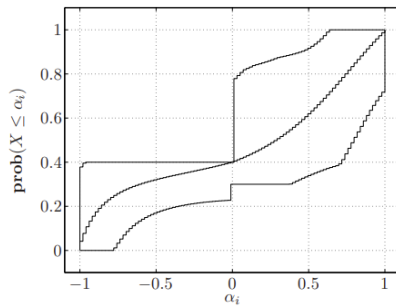


Figure 7.3 The top and bottom curves show the maximum and minimum possible values of the cumulative distribution function, $\text{prob}(X \leq \alpha_i)$, over all distributions that satisfy (7.8). The middle curve is the cumulative distribution of the maximum entropy distribution that satisfies (7.8).

we solve two LPs: one that maximizes $\text{prob}(X \leq \alpha_i)$, and one that minimizes $\text{prob}(X \leq \alpha_i)$, over all distributions consistent with the prior assumptions (7.8). The results are shown in figure 7.3. The upper and lower curves show the upper and lower bounds, respectively; the middle curve shows the cumulative distribution of the maximum entropy distribution.

7.3 Optimal detector design and Hypothesis testing

$$P = [p_{kj}] \rightarrow \theta = j \text{ 일 때 } X \text{ 의 dist.}$$

$$p_{kj} = \text{prob}(X = k \mid \theta = j).$$

7.3.1 Deterministic and randomized detectors

$$\text{Maximum likelihood detector: } \hat{\theta} = \psi_{\text{ml}}(k) = \underset{j}{\operatorname{argmax}} p_{kj}.$$

$$T: t_{ik} = \text{prob}(\hat{\theta} = i \mid X = k).$$

$$t_k \succeq 0, \quad \mathbf{1}^T t_k = 1.$$

7.3.2 Detection probability matrix

$$D_{ij} = (TP)_{ij} = \text{prob}(\hat{\theta} = i \mid \theta = j),$$

$$\text{Detection probabilities: } P_i^d = D_{ii} = \text{prob}(\hat{\theta} = i \mid \theta = i).$$

$$\text{Error probabilities: } P_i^e = 1 - D_{ii} = \text{prob}(\hat{\theta} \neq i \mid \theta = i). \quad P_i^e = \sum_{j \neq i} D_{ji}.$$

01

7.3.3 Optimal detector design

02

Limits on errors and detection probabilities

03

$$P_j^d = D_{jj} \geq L_j, \quad D_{ij} \leq U_{ij}$$

Minimax detector design

$$\begin{aligned} &\text{minimize} \quad \max_j P_j^e \\ &\text{subject to} \quad t_k \succeq 0, \quad \mathbf{1}^T t_k = 1, \quad k = 1, \dots, n, \end{aligned}$$

Bayes detector design

$$q_i = \text{prob}(\theta = i)$$

$$\begin{aligned} &\text{minimize} \quad q^T P^e \\ &\text{subject to} \quad t_k \succeq 0, \quad \mathbf{1}^T t_k = 1, \quad k = 1, \dots, n. \end{aligned}$$

Bias, mean-square error, and other quantities

$$\text{prob}(\hat{\theta} > \theta \mid \theta = i) = \sum_{j>i} D_{ji},$$

$$\text{prob}(|\hat{\theta} - \theta| > 1 \mid \theta = i) = \sum_{|j-i|>1} D_{ji},$$

$$\text{Bias:} \quad \mathbf{E}_i(\hat{\theta} - \theta) = \sum_{j=1}^m (\theta_j - \theta_i) D_{ji},$$

$$\text{Mean square error:} \quad \mathbf{E}_i(\hat{\theta} - \theta)^2 = \sum_{j=1}^m (\theta_j - \theta_i)^2 D_{ji}.$$

$$\text{Average absolute error:} \quad \mathbf{E}_i|\hat{\theta} - \theta| = \sum_{j=1}^m |\theta_j - \theta_i| D_{ji}.$$

7.3.4 Multicriterion formulation and scalarization

$$\begin{aligned} & \text{minimize} && D_{\bar{i}\bar{j}}, \quad \bar{i}, \bar{j} = 1, \dots, m, \quad \bar{i} \neq \bar{j} \\ & \text{subject to} && t_k \geq 0, \quad \mathbf{1}^T t_k = 1, \quad k = 1, \dots, n \end{aligned}$$

$$W_{ii} = 0, \quad W_{ij} \sim P^{\ell}(\hat{\theta} = i | \theta = j)$$

$$\begin{aligned} & \text{minimize} && \text{tr}(W^T D) \\ & \text{subject to} && t_k \geq 0, \quad \mathbf{1}^T t_k = 1, \quad k = 1, \dots, n, \quad \text{LP} \end{aligned}$$

$$\text{tr}(W^T D) = \text{tr}(W^T T P) = \text{tr}(P W^T T) = \sum_{k=1}^n c_k^T t_k, \quad (c_k = W P^T \text{의 열벡터})$$

$$\begin{aligned} & \text{minimize} && c_k^T t_k \\ & \text{subject to} && t_k \geq 0, \quad \mathbf{1}^T t_k = 1, \end{aligned}$$

$$\hat{\theta} = \underset{j}{\operatorname{argmin}} (W P^T)_{jk}.$$

7.3.5 Binary hypothesis testing

$$D = \begin{bmatrix} 1 - P_{fp} & P_{fn} \\ P_{fp} & 1 - P_{fn} \end{bmatrix}$$

$$\hat{\theta} = \begin{cases} 1 & W_{21}p_k > W_{12}q_k \\ 2 & W_{21}p_k \leq W_{12}q_k \end{cases} \left(\frac{p_k}{q_k} > \frac{W_{12}}{W_{21}} \right)$$

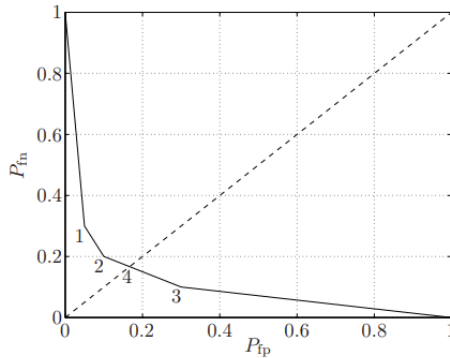


Figure 7.4 Optimal trade-off curve between probability of a false negative, and probability of a false positive test result, for the matrix P given in (7.15). The vertices of the trade-off curve, labeled 1–3, correspond to deterministic detectors; the point labeled 4, which is a randomized detector, is the minimax detector. The dashed line shows $P_{fn} = P_{fp}$, the points where the error probabilities are equal.

Example 7.4 We consider a binary hypothesis testing example with $n = 4$, and

$$P = \begin{bmatrix} 0.70 & 0.10 \\ 0.20 & 0.10 \\ 0.05 & 0.70 \\ 0.05 & 0.10 \end{bmatrix}. \quad (7.15)$$

The optimal trade-off curve between P_{fn} and P_{fp} , *i.e.*, the receiver operating curve, is shown in figure 7.4. The left endpoint corresponds to the detector which is always negative, independent of the observed value of X ; the right endpoint corresponds to the detector that is always positive. The vertices labeled 1, 2, and 3 correspond to the deterministic detectors

$$T^{(1)} = \begin{bmatrix} 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix},$$

$$T^{(2)} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix},$$

$$T^{(3)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix},$$

respectively. The point labeled 4 corresponds to the nondeterministic detector

$$T^{(4)} = \begin{bmatrix} 1 & 2/3 & 0 & 0 \\ 0 & 1/3 & 1 & 1 \end{bmatrix},$$

which is the minimax detector. This minimax detector yields equal probability of a false positive and false negative, which in this case is $1/6$. Every deterministic detector has either a false positive or false negative probability that exceeds $1/6$, so this is an example where a randomized detector outperforms every deterministic detector.

ESC 2023 spring WEEK1

About bounds and Experiment Design

학술부: 김민주,이종현

March 9,2023

Three Bounds

- $P(|X| \geq c) \leq E(|X|^r)/c^r$:Markov Bound
- $P(|X - \mu| \geq 1) \leq \sigma^2$:Chebyshev Bound
- $P(X \geq u) \leq \inf_{\lambda \geq 0} E(e^{\lambda(X-u)})$: Chernoff Bound

can be derived by convex optimization!

Chebyshev Bound: Objective

X : RV on $S \subseteq \mathbb{R}^m$, $C \subseteq S$

Find the bound of $P(X \in C) = E(1_C(x))$

Chebyshev Bound: Terms

prior knowledge: $E(f_i(X)) = a_i, E(f_0(X)) = 1$

linear combination and expectation: $f(z) = \sum_{i=0}^n x_i f_i(z) \rightarrow E(f(X)) = a^T x$

key idea: $f(z) \geq 1_C(z)$ for all $z \in C$

$\rightarrow E(f(X)) = a^T x \geq E(1_C(Z)) = P(X)$

Chebyshev Bound Problem

$$\underset{x}{\text{minimize}} \quad x_0 + a_1 x_1 + \dots + a_n x_n$$

$$\text{subject to} \quad f(z) = \sum_{i=0}^n x_i f_i(z) \geq 1, \text{ for } z \in C$$

$$f(z) = \sum_{i=0}^n x_i f_i(z) \geq 0, \text{ for } z \in S \setminus C$$

- $1 - \inf_z f(z) \leq 0$ is convex
- $-\inf_{z \in S \setminus C} f(z) \leq 0$ is convex

⇒ Above problem is convex optimization

Chebyshev Bound Problem: Markov

let $S \subseteq \mathbb{R}_+$, $C = [1, \infty)$, $f_0(z) = 1$, $f_1(z) = z$, $E(f_1(X)) = \mu \leq 1$

$$\begin{aligned} & \underset{x}{\text{minimize}} && x_0 + \mu x_1 \\ & \text{subject to} && x_0 \geq 0, x_1 \geq 0 \\ & && x_0 + x_1 \geq 1 \end{aligned}$$

→ This problem becomes simple LP, where it yields:
 $P(X \geq 1) \leq \mu$: Markov inequality

Chebyshev Bound Problem:with First and Second moment

let $E(X) = a \in R^m, E(XX^T) = \Sigma \in S^m$
 $f(z) = z^T Pz + 2q^T z + r$

$$\begin{aligned} E(f(X)) &= E(X^T P X + 2q^T X + r) \\ &= E(\text{tr}(PXX^T) + 2E(q^T X) + r) \\ &= \text{tr}(\Sigma P) + 2q^T a + r \end{aligned}$$

Chebyshev Bound Problem:with First and Second moment

$f(z) \geq 0$ for all z :

$$\rightarrow \begin{bmatrix} P & q \\ q^T & r \end{bmatrix} \succeq 0, \quad P \succeq 0$$

C we are looking for:

$$C = R^m \setminus \mathcal{P}, \quad \mathcal{P} = \{z | a_i^T z < b_i, \quad i = 1, \dots, k\}$$

Chebyshev Bound Problem:with First and Second moment

$f(z) \geq 1$ for all $z \in C$:

$$a_i^T z \geq b_i \Rightarrow z^T P z + 2q^T z + r \geq 1$$

There exist $\tau_1, \dots, \tau_k \geq 0$ such that

$$\begin{bmatrix} P & q \\ q^T & r-1 \end{bmatrix} \succeq \tau_i \begin{bmatrix} 0 & a_i/2 \\ a_i^T/2 & -b_i \end{bmatrix}, \quad i = 1, \dots, k$$

참고:Appendix B.2

Chebyshev Bound Problem: with First and second moment

Putting it all together, the Chebyshev bound problem (7.17) can be expressed as

$$\begin{aligned} & \text{minimize} && \text{tr}(\Sigma P) + 2q^T a + r \\ & \text{subject to} && \begin{bmatrix} P & q \\ q^T & r - 1 \end{bmatrix} \succeq \tau_i \begin{bmatrix} 0 & a_i/2 \\ a_i^T/2 & -b_i \end{bmatrix}, \quad i = 1, \dots, k \\ & && \tau_i \geq 0, \quad i = 1, \dots, k \\ & && \begin{bmatrix} P & q \\ q^T & r \end{bmatrix} \succeq 0, \end{aligned} \tag{7.19}$$

optimal value α : *upper bound for* $P(X \in C)$

$1 - \alpha$: *lower bound of* $P(X \in \mathcal{P})$

Chernoff Bound

$$P(X \geq u) \leq \inf_{\lambda \geq 0} E(e^{\lambda(X-u)})$$

$$\log P(X \geq u) \leq \inf_{\lambda \geq 0} \left\{ -\lambda\mu + \log E(e^{\lambda X}) \right\}$$

Terms:

- $\log E(e^{\lambda X})$: cumulant generating function (log-mgf)

-ex) $\log E(e^{\lambda X}) = \lambda^2/2$, when $X \sim N(0, 1)$

$\rightarrow P(X \geq u) \leq e^{-u^2/2}$

Chernoff Bound

- $\lambda \in R^m, \mu \in R, f: R^m \rightarrow R$
- $f(z) = e^{\lambda^T z + \mu}$
- $f(z) \geq 1_C(z)$

$$\rightarrow P(X \in C) = E(1_C(X)) \leq E(f(X))$$

$$f(z) \geq 0 \text{ and } f(z) \geq 1 \text{ iff } \lambda^T z + \mu \geq 0$$

Chernoff Bound

when $-\lambda^T z \leq \mu$:

$$P(X \in C) \leq E(\exp(\lambda^T X + \mu))$$
$$\log P(X \in C) \leq \mu + \log E(\exp(\lambda^T X))$$

Chernoff Bound

From this we obtain a general form of Chernoff's bound:

$$\begin{aligned}\log \mathbf{prob}(X \in C) &\leq \inf \{ \mu + \log \mathbf{E} \exp(\lambda^T X) \mid -\lambda^T z \leq \mu \text{ for all } z \in C \} \\ &= \inf_{\lambda} \left(\sup_{z \in C} (-\lambda^T z) + \log \mathbf{E} \exp(\lambda^T X) \right) \\ &= \inf (S_C(-\lambda) + \log \mathbf{E} \exp(\lambda^T X)),\end{aligned}$$

where S_C is the support function of C . Note that the second term, $\log \mathbf{E} \exp(\lambda^T X)$, is the cumulant generating function of the distribution, and is always convex (see example 3.41, page 106). Evaluating this bound is, in general, a convex optimization problem.

Chernoff Bound: Gaussian Polyhedron

$$X \sim N(0, I), \log E(\exp(\lambda^T X)) = \lambda^T \lambda / 2, C = \{x | Ax \preceq b\}$$

$$\begin{aligned} S_c(y) &= \sup \{y^T x | Ax \preceq b\} \\ &= -\inf \{-y^T x | Ax \preceq b\} \\ &= -\sup \{-b^T u | A^T u = y, u \succeq 0\} \\ &= \inf \{b^T u | A^T u = y, u \succeq 0\} \end{aligned}$$

$$\begin{aligned} \log P(X \in C) &\leq \inf_{\lambda} (S_C(-\lambda) + \log E \exp(\lambda^T X)) \\ &= \inf_{\lambda} \inf_u \{b^T u + \lambda^T \lambda / 2 | u \succeq 0, A^T u + \lambda = 0\} \end{aligned}$$

Chernoff Bound: Gaussian Polyhedron

This problem has an interesting geometric interpretation. It is equivalent to

$$\begin{array}{ll}\text{minimize} & b^T u + (1/2)\|A^T u\|_2^2 \\ \text{subject to} & u \succeq 0,\end{array}$$

which is the dual of

$$\begin{array}{ll}\text{maximize} & -(1/2)\|x\|_2^2 \\ \text{subject to} & Ax \preceq b.\end{array}$$

In other words, the Chernoff bound is

$$\mathbf{prob}(X \in C) \leq \exp(-\mathbf{dist}(0, C)^2/2), \quad (7.22)$$

where $\mathbf{dist}(0, C)$ is the Euclidean distance of the origin to C .

Experiment Design

Consider estimating a vector x in $y_i = a_i^T x + w_i$, $i = 1, \dots, m$, $w_i \sim N(0, 1)$

$$\hat{x} = \left(\sum_{i=1}^m a_i a_i^T \right)^{-1} \sum_{i=1}^m y_i a_i$$

$$E = E(ee^T) = \left(\sum_{i=1}^m a_i a_i^T \right)^{-1}$$

Note: Think about Regression Analysis:

- $\hat{x} = (A^T A)^{-1} A^T y$, $A = (a_1^T, \dots, a_m^T)(\text{rows})$
- $\text{cov}(\hat{\beta}) = (X^T X)^{-1} \sigma^2$

Experiment Design

why care about E ?:

→ *Because of* α confidence level ellipsoid

$$\zeta = \{z | (z - \hat{x})^T E^{-1} (z - \hat{x}) \leq \beta\}$$

Experiment Design

Possible choice of a_i : v_1, \dots, v_p

p 개의 v 중 m 개를 뽑은게 a 가 된다.(중복 허용하는듯)

m_j : v_j 가 몇번 뽑혔는가

$\rightarrow m_1 + \dots + m_p = m$

Objective: E 를 가장 작게하는 choice들을 찾아라.

$$E = \left(\sum_{i=1}^m a_i a_i^T \right)^{-1} = \left(\sum_{i=1}^p m_i v_i v_i^T \right)^{-1}$$

ex) v_1 이 a_1, a_2 에서 2개 뽑혔다

$\rightarrow m_j = 2$

Experiment Design

$$\begin{aligned} & \underset{w.r.t. S_+^n}{\text{minimize}} && E = \left(\sum_{i=1}^p m_i v_i v_i^T \right)^{-1} \\ & \text{subject to} && m_i \geq 0, m_1 + \dots + m_p \\ & && m_i \in \mathbf{Z} \end{aligned}$$

interpretation:

$E \preceq \tilde{E} \rightarrow$ the first ellipsoid contained in second
first experiment design estimate the variance of $q^T \hat{x}$ better

Relaxed Experiment Design

when m is sufficiently large compared to n ,

→ use relaxed experiment design

$$\lambda_i = m_i/m$$

$$E = \frac{1}{m} \left(\sum_{i=1}^p \lambda_i v_i v_i^T \right)^{-1}$$

Relaxed Experiment:

$$\begin{aligned} & \underset{w.r.t. S_+^n}{\text{minimize}} && E = \frac{1}{m} \left(\sum_{i=1}^p \lambda_i v_i v_i^T \right)^{-1} \\ & \text{subject to} && \lambda \succeq 0, \mathbf{1}^T \lambda = 1 \end{aligned}$$

This is a convex optimization since E is an S_+^n convex optimization of λ

Usefulness of Relaxed Experiment Design

앞에 꺼보다 constraint가 하나 적으므로, RED는 ED의 lower bound제공 λ_i 가 integer m_i 에 비례($1/m$)한다는 걸 빼고본다면,

$$m_i = \mathbf{round}(m\lambda_i), \quad i = 1, \dots, p$$
$$\tilde{\lambda}_i = (1/m)\mathbf{round}(m\lambda_i)$$

- λ :비례 뺀, $\tilde{\lambda}$:비례 함
- λ_i 와 $\tilde{\lambda}_i$ 의 차이는 $1/m$ 에 비례

→ m 이 크면, 비례제약이 없는 λ 를 구하고 $\tilde{\lambda}$ 를 구해도 상관이 없다.
→ 제약 하나를 없앨 수 있어서 relaxed experiment인것.

Scalarization:D-optimal design

D-optimal design

The most widely used scalarization is called *D-optimal design*, in which we minimize the determinant of the error covariance matrix E . This corresponds to designing the experiment to minimize the volume of the resulting confidence ellipsoid (for a fixed confidence level). Ignoring the constant factor $1/m$ in E , and taking the logarithm of the objective, we can pose this problem as

$$\begin{aligned} & \text{minimize} && \log \det \left(\sum_{i=1}^p \lambda_i v_i v_i^T \right)^{-1} \\ & \text{subject to} && \lambda \succeq 0, \quad \mathbf{1}^T \lambda = 1, \end{aligned} \tag{7.26}$$

which is a convex optimization problem.

About Ellipsoid

A related family of convex sets is the *ellipsoids*, which have the form

$$\mathcal{E} = \{x \mid (x - x_c)^T P^{-1} (x - x_c) \leq 1\}, \quad (2.3)$$

where $P = P^T \succ 0$, i.e., P is symmetric and positive definite. The vector $x_c \in \mathbf{R}^n$ is the *center* of the ellipsoid. The matrix P determines how far the ellipsoid extends in every direction from x_c ; the lengths of the semi-axes of \mathcal{E} are given by $\sqrt{\lambda_i}$, where λ_i are the eigenvalues of P . A ball is an ellipsoid with $P = r^2 I$. Figure 2.9 shows an ellipsoid in \mathbf{R}^2 .

Scalarization:E-optimal design

E-optimal design

In *E-optimal design*, we minimize the norm of the error covariance matrix, *i.e.*, the maximum eigenvalue of E . Since the diameter (twice the longest semi-axis) of the confidence ellipsoid \mathcal{E} is proportional to $\|E\|_2^{1/2}$, minimizing $\|E\|_2$ can be interpreted geometrically as minimizing the diameter of the confidence ellipsoid. *E*-optimal design can also be interpreted as minimizing the maximum variance of $q^T e$, over all q with $\|q\|_2 = 1$.

The *E*-optimal experiment design problem is

$$\begin{aligned} & \text{minimize} && \left\| \left(\sum_{i=1}^p \lambda_i v_i v_i^T \right)^{-1} \right\|_2 \\ & \text{subject to} && \lambda \succeq 0, \quad \mathbf{1}^T \lambda = 1. \end{aligned}$$

The objective is a convex function of λ , so this is a convex problem.

The *E*-optimal experiment design problem can be cast as an SDP

$$\begin{aligned} & \text{maximize} && t \\ & \text{subject to} && \sum_{i=1}^p \lambda_i v_i v_i^T \succeq tI \\ & && \lambda \succeq 0, \quad \mathbf{1}^T \lambda = 1, \end{aligned} \tag{7.27}$$

with variables $\lambda \in \mathbf{R}^p$ and $t \in \mathbf{R}$.

Proof to SDP cast

$$\sum_{i=1}^p \lambda_i v_i v_i^T = A \text{ and this is sym}$$

use

$$\|A\|_2 \leq t \text{ iff } A^T A \leq t^2 I. \text{ (eq 1709)}$$

$$\|A^{-1}\|_2 \leq t \text{ iff } (A^{-1})^T A^{-1} \leq t^2 I.$$

$$(A^{-1})^T = (A^T)^{-1} = A^{-1}$$

$$\therefore A^{-2} \leq t^2 I \Rightarrow \sum_{i=1}^p \lambda_i v_i v_i^T \geq t I.$$

$$\text{So the prob becomes } \max t \text{ for } \sum_{i=1}^p \lambda_i v_i v_i^T \geq t I.$$

Scalarization:A-optimal

A-optimal design

In *A-optimal experiment design*, we minimize $\text{tr } E$, the trace of the covariance matrix. This objective is simply the mean of the norm of the error squared:

$$\mathbf{E} \|e\|_2^2 = \mathbf{E} \text{tr}(ee^T) = \text{tr } E.$$

The *A-optimal* experiment design problem is

$$\begin{aligned} &\text{minimize} && \text{tr} \left(\sum_{i=1}^p \lambda_i v_i v_i^T \right)^{-1} \\ &\text{subject to} && \lambda \succeq 0, \quad \mathbf{1}^T \lambda = 1. \end{aligned} \tag{7.28}$$

This, too, is a convex problem. Like the *E-optimal* experiment design problem, it can be cast as an SDP:

$$\begin{aligned} &\text{minimize} && \mathbf{1}^T u \\ &\text{subject to} && \begin{bmatrix} \sum_{i=1}^p \lambda_i v_i v_i^T & e_k \\ e_k^T & u_k \end{bmatrix} \succeq 0, \quad k = 1, \dots, n \\ &&& \lambda \succeq 0, \quad \mathbf{1}^T \lambda = 1, \end{aligned}$$

u_k is the diagonal element of the objective

Proof to SDP

see <https://www.tandfonline.com/doi/full/10.1080/03610918.2015.1030414>

D-optimal's Geometric Meaning

Optimal experiment design and duality

The Lagrange duals of the three scalarizations have an interesting geometric meaning.

The dual of the D -optimal experiment design problem (7.26) can be expressed as

$$\begin{aligned} & \text{maximize} && \log \det W + n \log n \\ & \text{subject to} && v_i^T W v_i \leq 1, \quad i = 1, \dots, p, \end{aligned}$$

with variable $W \in \mathbf{S}^n$ and domain \mathbf{S}_{++}^n (see exercise 5.10). This dual problem has a simple interpretation: The optimal solution W^* determines the minimum volume ellipsoid, centered at the origin, given by $\{x \mid x^T W^* x \leq 1\}$, that contains the points v_1, \dots, v_p . (See also the discussion of problem (5.14) on page 222.) By complementary slackness,

$$\lambda_i^* (1 - v_i^T W^* v_i) = 0, \quad i = 1, \dots, p, \quad (7.29)$$

i.e., the optimal experiment design only uses the experiments v_i which lie on the surface of the minimum volume ellipsoid.

D-optimal's Geometric Meaning

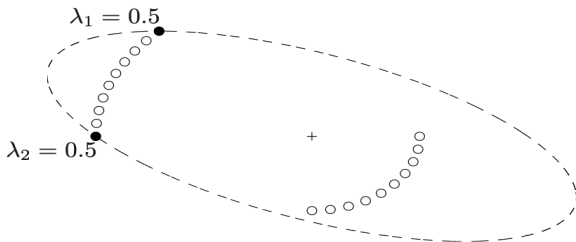


Figure 7.9 Experiment design example. The 20 candidate measurement vectors are indicated with circles. The D -optimal design uses the two measurement vectors indicated with solid circles, and puts an equal weight $\lambda_i = 0.5$ on each of them. The ellipsoid is the minimum volume ellipsoid centered at the origin, that contains the points v_i .

Coding

https://jump.dev/JuMP.jl/stable/tutorials/conic/experiment_design/#A-optimal-design