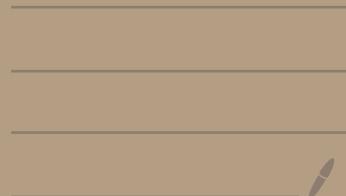


Lec II: Application ,

Statistical estimation



• Parametric distribution estimation

choose from a family of densities $P_x(y)$, indexed by a parameter x

• maximum likelihood estimation

$$\max_x P_x(y)$$

y is observed value

$l(x) = \log P_x(y)$ is called log-likelihood function

Many densities are log-concave (in x vs. in y)

can add constraints $x \in C$ explicitly, or define $P_x(y) = 0$ for $x \notin C$
a convex optimization problem if $\log P_x(y)$ is concave in x for fixed y

e.g. linear measurement with IID noise

$$y_i = a_i^T x + v_i$$

$x \in \mathbb{R}^n$ is the vector of unknown parameters

v_i is IID measurement noise, with density $p(v)$

y_i is measurement: y_i has density $p(y) = \prod_{i=1}^m p(y_i - a_i^T x)$

$$\max_x l(x) = \sum_i \log p(y_i - a_i^T x) \quad (\text{need } p \text{ to be log-concave})$$

e.g. Gaussian noise $N(0, \sigma^2)$: $p(v) = \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{v^2}{2\sigma^2}\right)$

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

$$= m \cdot \log \frac{1}{\sqrt{2\pi}\sigma} - \underbrace{\frac{1}{2\sigma^2} \cdot \sum_{i=1}^m (y_i - a_i^T x)^2}_{}$$

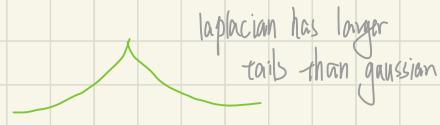
least-squares is maximum likelihood estimation with Gaussian noise

e.g. Laplacian noise $p(v) = \frac{1}{2a} \cdot \exp(-|v|/a)$

$$l(x) = \sum_{i=1}^m \log \left[\frac{1}{2a} \cdot \exp(-|y_i - a_i^T x|/a) \right]$$

$$= -m \cdot \log(2a) - \frac{1}{a} \sum_{i=1}^m |y_i - a_i^T x|$$

minimizing L₁ norm is mle with Laplacian noise



eg. Uniform noise in $[-a, a]$

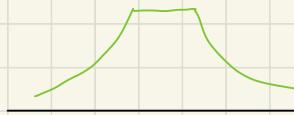
$$\ell(x) = \begin{cases} -m \log 2a & |ax-y_i| \leq a \\ -\infty & \text{otherwise} \end{cases}$$

mle estimation is any x with $|ax-y_i| \leq a$ (feasibility problem)

eg. penalty function



penalty function



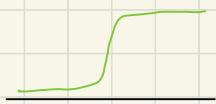
noise distribution

(uniform center and exponential tail)

eg. logistic regression

random variable $y \in \{0, 1\}$ with distribution

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



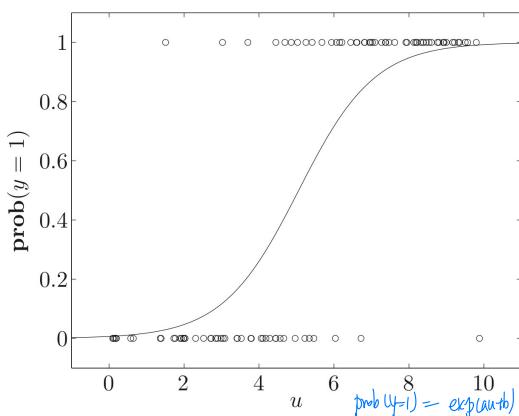
a, b are parameters, $u \in \mathbb{R}^n$ are observations

estimation problem: estimate a, b from (u_i, y_i)

$$(y_1 = \dots = y_k = 1 \quad y_{k+1} = \dots = y_m = 0)$$

$$\begin{aligned} \ell(a, b) &= \sum_{i=1}^k \left(\frac{\exp(a^T u_i + b)}{1 + \exp(a^T u_i + b)} \cdot \frac{1}{1 + \exp(a^T u_i + b)} \right) \\ &= \sum_{i=1}^k (a^T u_i + b) - \sum_{i=1}^m \log[1 + \exp(a^T u_i + b)] \end{aligned}$$

concave in a, b



$$\text{prob}(y=1) = \exp(b)/[1 + \exp(a^T u + b)]$$

b controls the natural point
 a controls the stretch

eg. Covariance estimation for Gaussian variables

Suppose $\mathbf{y} \in \mathbb{R}^n$ is a gaussian random variable

$$E[\mathbf{y}] = \mathbf{0}, R = \text{cov}(\mathbf{y}) = E[\mathbf{y}\mathbf{y}^T]$$

$$p_R(\mathbf{y}) = \frac{1}{(2\pi)^{n/2} |\mathbf{R}|^{1/2}} \exp\left(-\frac{\mathbf{y}^T \mathbf{R}^{-1} \mathbf{y}}{2}\right) \quad \mathbf{R} \in \mathbb{S}_{++}^n$$

$$\ell(R) = \sum_i^m \log p(\mathbf{y}_i)$$

$$= -m \cdot \left[\frac{n}{2} \cdot \log 2\pi + \frac{1}{2} \log \det(R) \right] - \frac{1}{2} \sum_{i=1}^m \mathbf{y}_i^T \mathbf{R}^{-1} \mathbf{y}_i$$

$$= -m \left[\frac{n}{2} \cdot \log 2\pi + \frac{1}{2} \log \det(R) \right] - \frac{m}{2} \text{tr}(R^{-1} \mathbf{Y})$$

let $S = R^{-1}$ (the information matrix)

$$= -\frac{mn}{2} \cdot \log 2\pi + m/2 \log \det S - \frac{m}{2} \text{tr}(SY)$$

is concave in S

$$\max. \log \det(S) - \text{tr}(SY)$$

1° no constraints on S

$$\max \log \det(S) - \text{tr}(SY)$$

$$\text{st } S \in \mathbb{S}_{++}^n \quad (S = R^{-1})$$

$$\nabla_S (\log \det(S) - \text{tr}(SY)) = S^{-1} - Y \Rightarrow$$

$$R = S^{-1} = Y \text{ if } Y \in \mathbb{S}_{++}^n$$

the problem is unbounded below if $Y \notin \mathbb{S}_{++}^n$

if there's no prior assumptions on R ,

then mle estimation of R is sample covariance

Sample covariance

$Y = 1/m \cdot \sum_i \mathbf{y}_i \mathbf{y}_i^T$ is the

• maximum a posteriori probability estimation.
(a bayesian version of MLE)

assume that x (to be estimated) and y (the observation) are random variables with a joint probability p_{xy}

the prior density is given by $p_x(x) = \int p_{xy} dy$ $p_y(y) = \int p_{xy} dy$

the conditional density is given by $p_{yx}(y|x) = \frac{p_{xy}}{p_x(x)}$

$$\hat{x}_{map} = \arg \max_x p_{xy}(x|y)$$

$$= \arg \max_x p_{xy}(x|y) \cdot p_y(y)$$

$$= \arg \max_x p_{xy}(x|y)$$

$$= \arg \max_x p_{yx}(y|x) \cdot p_x(x) \quad \text{taking prior knowledge of } x \text{ into account compared to MLE}$$

e.g. linear measurement with IID noise

$$y_i = a_i^T x + v_i$$

v_i are iid noise with density p_v on \mathbb{R} , x has prior density p_x on \mathbb{R}^n

$$p_{xy} = p_x(x) \cdot \prod_{i=1}^m p_v(y_i - a_i^T x)$$

↓

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

when v_i is uniform on $[c, d]$, and the prior distribution of x is gaussian

$$\max_x \log \frac{1}{(2\pi)^{n/2} |\Sigma|^1/2} \cdot \exp(-\frac{1}{2} (x - \bar{x})^T \Sigma^{-1} (x - \bar{x}))$$

$$\text{s.t. } \|Ax - b\|_w \leq a$$

↓

$$\min_x (x - \bar{x})^T \Sigma^{-1} (x - \bar{x})$$

$$\text{s.t. } \|Ax - b\|_w \leq a$$

e.g. MAP with perfect measurements

$x \in \mathbb{R}^n$ is a vector of parameters to be estimated, p_x have m perfect, noise-free measurements $y = Ax$

$$\text{max. } \log p_x(x)$$

$$\text{s.t. } Ax = y$$

• (Binary) hypothesis testing

1. detection (hypothesis testing) problem

given observation of a random variable $X \in \{1, \dots, n\}$ chose between:

hypothesis 1: X was generated by distribution $p = (p_1, \dots, p_n)$

$$1^T p = 1 \quad p \geq 0$$

hypothesis 2: X was generated by distribution $q = (q_1, \dots, q_n)$

$$1^T q = 1 \quad q \geq 0$$

2. randomized detector

a nonnegative matrix $T \in \mathbb{R}^{n \times n}$, with $1^T T = 1$

$$\begin{bmatrix} .9 & .8 & .6 \\ .1 & .2 & .4 \end{bmatrix}$$

if we observe $X=k$, we choose hypol with prob T_{1k} , hypo+ with prob T_{2k}

if T has all elements 0 or 1, its a deterministic detector

3. detection probability matrix

$$D = T \cdot [p \quad q] = [Tp \quad Tq] = \begin{bmatrix} 1 - P_{fp} & P_{fn} \\ P_{fp} & 1 - P_{fn} \end{bmatrix}$$

P_{fp} is probability of choosing hypothesis 2 when X generated by distribution 1

e.g. multicriterion formulation of detector design.

$$\text{min. (W.r.t } R^c) \quad (P_{fp}, P_{fn}) = (Tp)_2, (Tq)_1$$

s.t.

$$t_{1k} + t_{2k} = 1$$

$$t_{ik} \geq 0$$

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\text{min. } \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} (Tp)_2 & (Tq)_1 \end{bmatrix}$$

$$\text{s.t. } t_{1k} + t_{2k} = 1$$

$$t_{ik} \geq 0$$

$$\text{a LP with a simple analytic solution } (t_{1k}, t_{2k}) = \begin{cases} (1, 0) & P_k \geq \lambda \cdot q_k \\ (0, 1) & P_k < \lambda \cdot q_k \end{cases}$$

- a deterministic detector, given by a likelihood ratio test $(P_k/q_k \stackrel{?}{\geq} \lambda)$
- when $\lambda=1$, it's minimizing $P(\text{being wrong}) \Leftrightarrow \max_i P(\text{being right})$ (MLE)

e.g. minimax detector

$$\text{min. } \max \{P_{fn}, P_{fp}\} = \max \{Tp)_2, (Tq)_1\}$$

$$\text{s.t. } t_{1k} + t_{2k} = 1 \quad t_{ik} \geq 0$$

a LP; solution is usually not deterministic

Experiment design

m linear measurements $y_i = a_i^T x + w_i$ if unknown $x \in \mathbb{R}^n$
 $w_i \sim N(0, 1)$ iid

$$A = \begin{bmatrix} -\alpha_1 \\ \vdots \\ -\alpha_m \end{bmatrix}$$

least-square estimate is: $\hat{x} = (\sum a_i a_i^T)^{-1} \sum y_i a_i$ $(A^T A)^{-1} A^T y$

error $e = \hat{x} - x$ has zero mean and covariance: $E = E[ee^T] = (\sum a_i a_i^T)^{-1}$

$$\begin{aligned}\hat{x} &= (A^T A)^{-1} A^T y \\ &= (A^T A)^{-1} A^T (Ax + w) \\ &= (A^T A)^{-1} A^T A x + (A^T A)^{-1} A^T w \\ &= x + (A^T A)^{-1} A^T w \\ E[\hat{x}] &= x \quad (w \sim N(0, I)) \\ E[(\hat{x} - x)(\hat{x} - x)^T] &= E[((A^T A)^{-1} A^T w) \cdot ((A^T A)^{-1} A^T w)^T | A] \\ &= (A^T A)^{-1} A^T E[w w^T] \cdot A (A^T A)^{-1} \\ E &= (A^T A)^{-1} \quad (\text{if } A \text{ large, then it has high signal-to-noise ratio, so covariance small})\end{aligned}$$

confidence ellipsoid is given by: $\{x \mid (x - \bar{x})^T E^{-1} (x - \bar{x}) \leq \beta\}$

experiment design: choose a set of possible tests to make E small
 $\{ \text{choose } a_i \text{ for } m_i \text{ times} \mid \sum m_i = m \}$

eg sensor a_1 has the highest signal-to-noise ratio.
 if choose a_1 for two times and no other a_i ,
 then $\sum a_i a_i^T$ is not invertible (rank 1)

↓

$$\min_{\text{w.r.t. } S^1} E = \left(\sum_{k=1}^p m_k \cdot V_k V_k^T \right)^{-1}$$

$$\text{s.t. } m_k \geq 0 \quad m_1 + \dots + m_p = m$$

$$m_k \in \mathbb{Z}$$

(hard to solve due to the integer constraint)

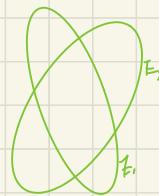
↓

Relaxed experiment design

$$\min_{\lambda} (\text{w.r.t } S^*) \quad E = \frac{1}{m} \left(\sum_{k=1}^P \lambda_k \cdot V_k \cdot V_k^T \right)^{-1}$$

$$\text{s.t.} \quad \lambda \geq 0$$

$$1^T \lambda = 1$$



Common scalarization: $\min \log\det(E), \text{tr}(E), \lambda_{\max}(E)$

e.g. D-optimal design

$$\min \log\det \left(\sum_{k=1}^P \lambda_k V_k V_k^T \right)^{-1}$$

$$\text{s.t.} \quad \lambda \geq 0 \quad 1^T \lambda = 1$$

(minimizing the volume of confidence ellipsoid)

↓

$$\min \log\det(X^T)$$

$$\text{s.t.} \quad X = \sum_{k=1}^P \lambda_k V_k V_k^T \quad (2) \quad (V_k \in \mathbb{R}^n)$$

$$\lambda \geq 0 \quad 1^T \lambda = 1$$

(2)

$$\begin{aligned} L(X, \lambda, z, \beta) &= \log\det(X^T) + \text{tr} \left\{ z \left(X - \sum_{k=1}^P \lambda_k V_k V_k^T \right) \right\} - d^T \lambda + \beta (1^T \lambda - 1) \\ &= \log\det(X^T) + \text{tr}(zX) - \text{tr} \left\{ z \left(\sum_{k=1}^P \lambda_k V_k V_k^T \right) \right\} - d^T \lambda + \beta (1^T \lambda - 1) \\ &= \log\det(X^T) + \text{tr}(zX) - \sum_{k=1}^P \lambda_k \cdot V_k^T z V_k - \sum_{k=1}^P d_k \lambda_k + \sum_{k=1}^P \beta_k \lambda_k - \beta \\ &= \log\det(X^T) + \text{tr}(zX) + \sum_{k=1}^P \lambda_k (-V_k^T z V_k - d_k + \beta) - \beta \end{aligned}$$

$$g(\lambda, z, \beta) = \inf_{X \in S^*} \left\{ L(X, \lambda, z, \beta) \right\}$$

$$= \begin{cases} \log\det(X^T) + \text{tr}(zX) - \beta - V_i^T z V_i - d_i + \beta = 0 \\ -\infty \text{ otherwise} \end{cases}$$

$$\forall \lambda \quad \log\det(X^T) + \text{tr}(zX) = 0$$

$$-X^T z = 0 \quad z = X^T$$

$$\therefore g(\lambda, z, \beta) = \begin{cases} \log\det z + n - \beta - V_i^T z V_i - d_i + \beta = 0 \\ -\infty \text{ otherwise} \end{cases}$$

↓

$$\max_z \log\det z + n - \beta$$

$$\text{s.t.} \quad \beta - V_i^T z V_i \geq 0$$

$$\Downarrow W = (1/\beta) \cdot z$$

$$\begin{aligned} \text{max. } & \log \det W + n \log \beta + n - \beta \\ \text{s.t. } & v_i^T W v_i \leq 1 \\ & \text{if } \max_i n \log \beta - \beta \text{ over } \beta : \beta = n \end{aligned}$$

dual of D-optimal

$$\begin{aligned} \text{max. } & \log \det W + n \log n \\ \text{s.t. } & v_i^T W v_i \leq 1 \end{aligned}$$

$\{x | x^T W x \leq 1\}$ is the minimum volume ellipsoid at origin that includes all test vector v_i

complementary slackness:

for x, w that's a primal-optimal - dual-optimal pair
 $d_i (1 - v_i^T W v_i) = 0$

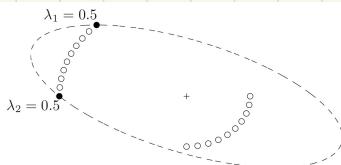


Figure 7.9 Experiment design example. The 20 candidate measurement vectors are indicated with circles. The D -optimal design uses the two measurement vectors indicated with solid circles, and puts an equal weight $\lambda_i = 0.5$ on each of them. The ellipsoid is the minimum volume ellipsoid centered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