

9. Statistical estimation (chap 7)

- maximum likelihood estimation
- *logistic regression*
- ~~optimal detector design~~
- experiment design (*statistics, Pakelshheim '70's, ...*
multi-armed bandits (online decision making)
(G-optimal design)
lots of applications...)

Parametric distribution estimation

- distribution estimation problem: estimate probability density $p(y)$ of a random variable from observed values
- parametric distribution estimation: choose from a family of densities $p_x(y)$, indexed by a parameter x

maximum likelihood estimation

$$\left[\text{maximize (over } x) \quad \log p_x(y) \right]$$

- y is observed value
- $l(x) = \log p_x(y)$ is called log-likelihood function
- 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

Linear measurements with IID noise

linear measurement model

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

- $x \in \mathbf{R}^n$ is vector of unknown parameters

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

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

$$\begin{aligned} v_i &= y_i - a_i^T x \\ p(v) &= \prod_{i=1}^m p(v_i) \\ &= \prod_{i=1}^m p(y_i - a_i^T x) \end{aligned}$$

maximum likelihood estimate: any solution x of

$$\left[\underset{x}{\text{maximize}} \quad l(x) = \sum_{i=1}^m \log p(y_i - a_i^T x) \right]$$

(y is observed value)

examples

- Gaussian noise $\mathcal{N}(0, \sigma^2)$: $p(z) = (2\pi\sigma^2)^{-1/2} e^{-z^2/(2\sigma^2)}$,

$$\rightarrow l(x) = \underbrace{-\frac{m}{2} \log(2\pi\sigma^2)}_{\text{const.}} - \frac{1}{2\sigma^2} \underbrace{\sum_{i=1}^m (a_i^T x - y_i)^2}_{\|y - Ax\|_2^2}$$

ML estimate is LS solution

- Laplacian noise: $p(z) = (1/(2a)) e^{-|z|/a}$,

$$l(x) = -m \log(2a) - \frac{1}{a} \underbrace{\sum_{i=1}^m |a_i^T x - y_i|}_{\|y - Ax\|_1}$$

ML estimate is ℓ_1 -norm solution

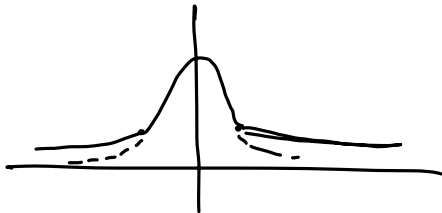
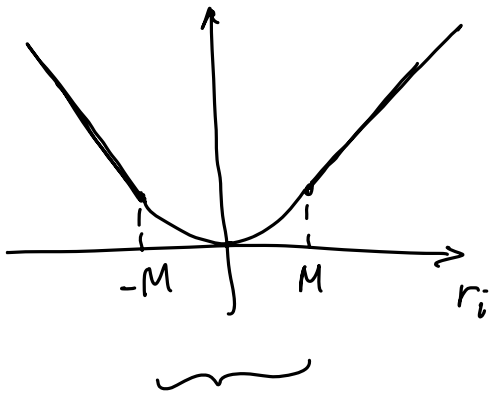
- uniform noise on $[-a, a]$:

$$l(x) = \begin{cases} -m \log(2a) & \underbrace{|a_i^T x - y_i| \leq a,}_{\text{otherwise}} \quad i = 1, \dots, m \\ -\infty & \end{cases}$$

$\|y - Ax\|_\infty$

ML estimate is any x with $|a_i^T x - y_i| \leq a$

noise dist $\sim \exp(\text{penalty})$



Huber penalty

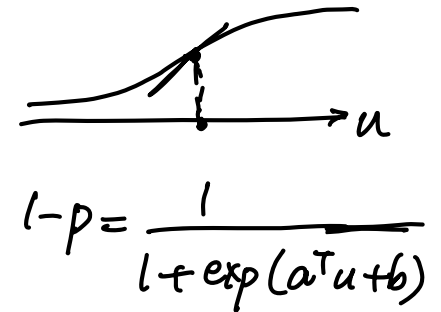
$$\left[\min_x \sum_i h(y_i - a_i^T x) \right] \equiv \begin{array}{l} \text{max} \\ \text{likelihood} \\ \text{estimate for} \\ \text{what noise?} \end{array}$$

$$r_i = y_i - a_i^T x$$

Logistic regression

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

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



- a, b are parameters; $u \in \mathbf{R}^n$ are (observable) explanatory variables
- estimation problem: estimate a, b from m observations (u_i, y_i)

log-likelihood function (for $y_1 = \dots = y_k = 1, y_{k+1} = \dots = y_m = 0$):

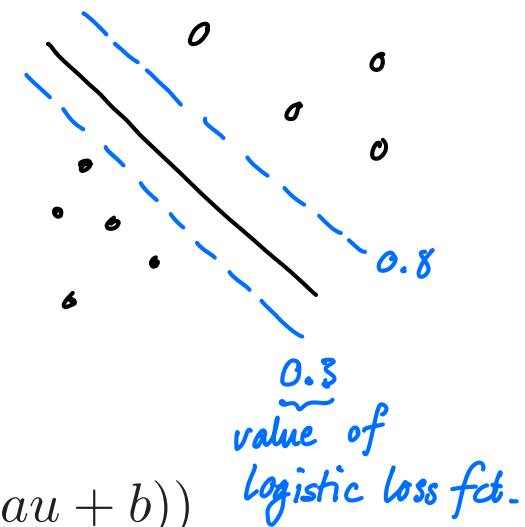
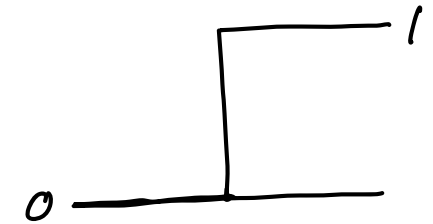
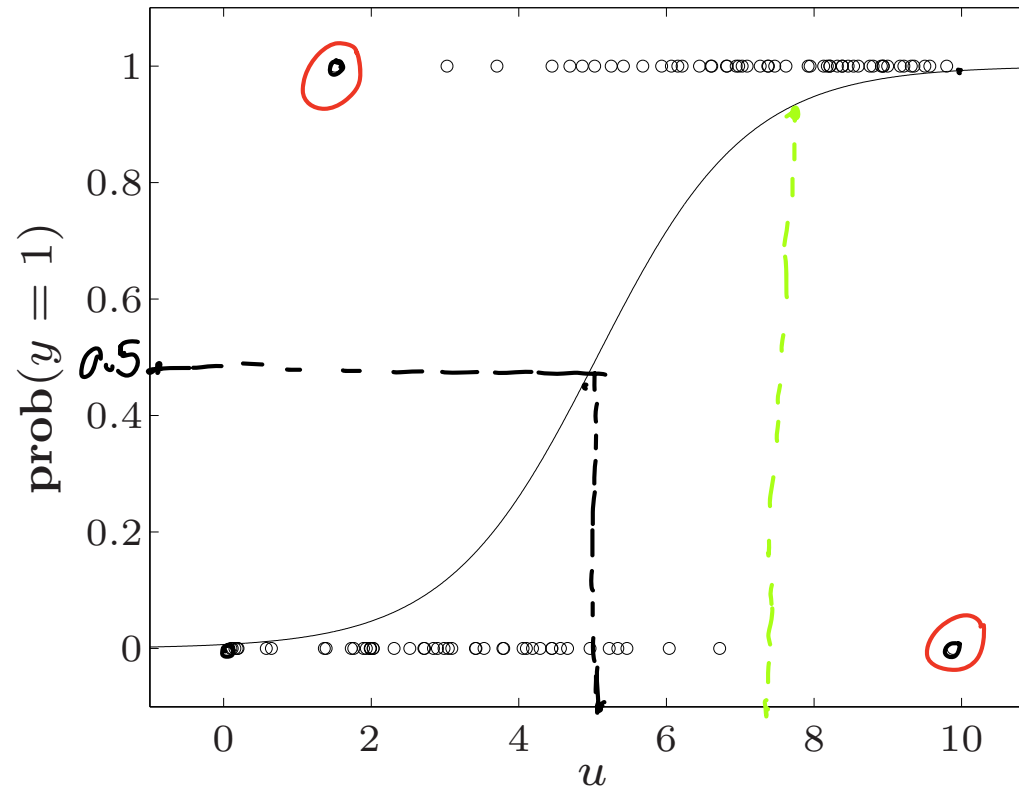
$$l(a, b) = \log \left(\underbrace{\prod_{i=1}^k \frac{\exp(a^T u_i + b)}{1 + \exp(a^T u_i + b)}}_{\text{sort the data according to labels}} \underbrace{\prod_{i=k+1}^m \frac{1}{1 + \exp(a^T u_i + b)}}_{\text{sort the data according to labels}} \right)$$

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

\uparrow
 e^0
log-sum-exp

concave in a, b

example ($n = 1$, $m = 50$ measurements)



- circles show 50 points (u_i, y_i)
- solid curve is ML estimate of $p = \exp(au + b) / (1 + \exp(au + b))$

can add regularizers on a : $\dots + \lambda \|a\|_2^2 + \lambda \|a\|_1 \rightsquigarrow$ encourages sparse $a \Rightarrow$ feature selection⁹⁻⁶

Experiment design

m linear measurements $y_i = \underline{a_i^T} x + \underline{w_i}$, $i = 1, \dots, \underline{m}$ of unknown $\underline{x} \in \mathbf{R}^n$
deterministic

- measurement errors w_i are IID $\mathcal{N}(0, 1)$

- ML (least-squares) estimate is

*if $\sum a_i a_i^T$ is singular
 \Rightarrow a_i do not span \mathbb{R}^n*

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

- error $e = \hat{x} - x$ has zero mean and covariance

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

confidence ellipsoids are given by $\{x \mid (x - \hat{x})^T E^{-1} (x - \hat{x}) \leq \underline{\beta}\}$

experiment design: choose $\underline{a_i} \in \{\underline{v_1}, \dots, \underline{v_p}\}$ (a set of possible test vectors) to make E 'small'

$$x \in \mathbb{R}^n$$

,

$$A \in \mathbb{R}^{m \times n}$$

$$m > n$$

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} a_1^T \\ \vdots \\ a_m^T \end{bmatrix} x + \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$

$$y = Ax + w \quad \begin{matrix} \nearrow \text{r.v.} \\ \nearrow \sim N(0,1), \text{ iid} \end{matrix}$$

$$\hat{x} = \arg \min_x \|y - Ax\|_2^2$$

$$\hat{x} = (A^T A)^{-1} A^T y$$

$$e = x - \hat{x} = x - (A^T A)^{-1} A^T y$$

$$\mathbb{E} y = Ax + \mathbb{E} \hat{\omega}^0 = Ax$$

$$\mathbb{E} e = x - (A^T A)^{-1} A^T \mathbb{E} y = x - x = 0$$

$$e = x - \hat{x} = x - (A^T A)^{-1} A^T (Ax + w) = \cancel{x - (A^T A)^{-1} A^T A x} - (A^T A)^{-1} A^T w = -(A^T A)^{-1} A^T w$$

$$\begin{aligned} \mathbb{E} e e^T &= \mathbb{E} (x - \hat{x})(x - \hat{x})^T = \mathbb{E} (A^T A)^{-1} A^T \underbrace{w w^T}_I A (A^T A)^{-1} \\ &= (A^T A)^{-1} A^T \underbrace{\mathbb{E} w w^T}_I A (A^T A)^{-1} = (A^T A)^{-1} \end{aligned}$$

$$A = \begin{bmatrix} a_1^T \\ \vdots \\ a_m^T \end{bmatrix}$$

$$A^T A = \sum_{i=1}^m a_i a_i^T$$

vector optimization formulation $E = (\sum a_i a_i^T)^{-1}$ $a_i \in \{v_1, \dots, v_p\}$
 $m_1 \quad m_p$

$$\left[\begin{array}{l} \text{minimize (w.r.t. } \mathbf{S}_+^n) \\ \text{subject to} \end{array} \right. \quad \begin{array}{l} \underline{E} = \left(\sum_{k=1}^p m_k v_k v_k^T \right)^{-1} \\ m_k \geq 0, \quad m_1 + \dots + m_p = \underline{m} \\ \underline{m_k \in \mathbf{Z}} \end{array}$$

- variables are m_k ($\#$ vectors a_i equal to $\underline{v_k}$)
- difficult in general, due to integer constraint

relaxed experiment design

assume $m \gg p$, use $\lambda_k = m_k/m$ as (continuous) real variable

$$\left[\begin{array}{l} \text{minimize (w.r.t. } \mathbf{S}_+^n) \\ \text{subject to} \end{array} \right. \quad \begin{array}{l} E = (1/m) \left(\sum_{k=1}^p \lambda_k v_k v_k^T \right)^{-1} \\ \lambda \succeq 0, \quad \mathbf{1}^T \lambda = 1 \end{array}$$

- common scalarizations: minimize $\log \det E$, $\text{tr } E$, $\lambda_{\max}(E)$, \dots
- can add other convex constraints, e.g., bound experiment cost $c^T \lambda \leq \underline{B}$
 experiment k (λ_k times) has cost of c_k

D-optimal design

$$\begin{cases} \underset{\lambda}{\text{minimize}} & \log \det \left(\sum_{k=1}^p \lambda_k v_k v_k^T \right)^{-1} \\ \text{subject to} & \lambda \succeq 0, \quad \mathbf{1}^T \lambda = 1 \end{cases}$$

interpretation: minimizes volume of confidence ellipsoids

dual problem

$$\rightarrow \begin{cases} \underset{W}{\text{maximize}} & \log \det W + n \log n \\ \text{subject to} & \underline{v_k^T} W \underline{v_k} \leq 1, \quad k = 1, \dots, p \end{cases} \Leftrightarrow \text{all } v_k \text{ are in ellipsoid centered at } 0 \text{ defined by } W$$

interpretation: $\{x \mid x^T W x \leq 1\}$ is minimum volume ellipsoid centered at origin, that includes all test vectors v_k

complementary slackness: for λ , W primal and dual optimal

$$\underbrace{\lambda_k}_{\text{primal}} \underbrace{(1 - v_k^T W v_k)}_{\text{dual}} = 0, \quad k = 1, \dots, p$$

$\neq 0$

optimal experiment uses vectors v_k on boundary of ellipsoid defined by W

$$\left[\begin{array}{ll} \min_{\lambda, X} & \log \det X^{-1} \\ \text{s.t.} & X = \sum_{k=1}^p \lambda_k v_k v_k^T \quad ; \quad Z \in S^n \quad (\text{equality}) \\ & \lambda \geq 0 \quad ; \quad z \in \mathbb{R}^p \quad z \geq 0 \\ & 1^T \lambda = 1 \quad ; \quad \nu \end{array} \right.$$

$$\begin{aligned} L(X, \lambda; Z, z, \nu) &= \log \det X^{-1} + \text{Tr } Z(X - \sum \lambda_k v_k v_k^T) - z^T \lambda + \nu(1^T \lambda - 1) \\ &= \underbrace{\log \det X^{-1} + \text{Tr } Z X}_{L_1(X; Z, z, \nu)} - \underbrace{\text{Tr } Z(\sum \lambda_k v_k v_k^T) - z^T \lambda + \nu 1^T \lambda - \nu}_{L_2(\lambda; Z, z, \nu)} \end{aligned}$$

$$g(Z, z, \nu) = \inf_{X, \lambda} L(X, \lambda; Z, z, \nu)$$

$$= \inf_X (\log \det X^{-1} + \text{Tr } Z X) + \inf_{\lambda} (\nu 1^T \lambda - \text{Tr } Z(\sum \lambda_k v_k v_k^T) - z^T \lambda) - \nu$$

$$\lambda \left(\nu 1 - \begin{bmatrix} v_1^T Z v_1 \\ \vdots \\ v_p^T Z v_p \end{bmatrix} - z \right)^T \lambda$$

$$\nabla_X L = 0 \Rightarrow -X^{-1} + Z = 0$$

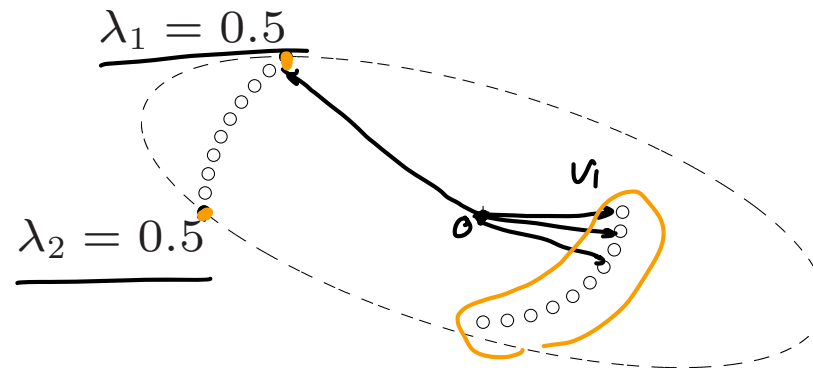
$$\begin{aligned} \nu 1 - \underbrace{\text{Tr } Z v_k v_k^T}_{= \text{Tr}(v_k^T Z v_k)} - z \\ = \nu 1 - \text{Tr}(v_k^T Z v_k) - z \\ = \nu 1 - \nu 1 - z \end{aligned}$$

$$\text{Tr}(\sum \lambda_k Z v_k v_k^T) = \sum \lambda_k \underbrace{\text{Tr}(Z v_k v_k^T)}_{= v_k^T Z v_k} = \sum \lambda_k v_k^T Z v_k$$

$$\inf_{\lambda} L_2(\lambda; Z, z, \nu) = \begin{cases} 0 & \nu - v_k^T Z v_k - z_k = 0 \quad k=1, \dots, p \\ -\infty & \text{else} \end{cases}$$

$$v_1, \dots, v_{20} \in \mathbb{R}^2$$

example ($p = 20$)



design uses two vectors, on boundary of ellipse defined by optimal W

$$y_i = \underbrace{a_i^T}_{\uparrow} x + w_i \quad w_i \sim N(0, 1) \quad \|a_i\|_2 \equiv \text{SNR} \quad (\text{large is good})$$

- pick exp. with good SNR
- pick exp. that are most informative (as orthogonal as possible)
less redundancy / overlap

derivation of dual of page 9-9

first reformulate primal problem with new variable X :

$$\begin{cases} \text{minimize} & \log \det X^{-1} \\ \text{subject to} & X = \sum_{k=1}^p \lambda_k v_k v_k^T, \quad \lambda \succeq 0, \quad \mathbf{1}^T \lambda = 1 \end{cases}$$

$$\underline{L(X, \lambda, Z, z, \nu) = \log \det X^{-1} + \text{tr} \left(Z \left(X - \sum_{k=1}^p \lambda_k v_k v_k^T \right) \right) - z^T \lambda + \nu (\mathbf{1}^T \lambda - 1)}$$

- minimize over X by setting gradient to zero: $\underline{-X^{-1} + Z = 0} \rightarrow Z = X^{-1}$
- minimum over $\underline{\lambda_k}$ is $-\infty$ unless $\underline{-v_k^T Z v_k - z_k + \nu = 0} \rightarrow Z_k = \nu - v_k^T Z v_k \geq 0$

dual problem

$$\begin{cases} \text{maximize}_{Z, \nu} & \cancel{\log \det X} + \log \det Z - \nu \\ \text{subject to} & v_k^T Z v_k \leq \nu, \quad k = 1, \dots, p \end{cases}$$

change variable $W = Z/\nu$, and optimize over ν to get dual of page 9-9

$$\begin{aligned} & \cancel{\log \det X} \\ & \nu^T W \nu_k \leq 1 \end{aligned}$$

$$\left[\begin{array}{l} \max_{\nu, W} \log \det(\nu W) - \nu = \log(\nu^n \det W) - \nu = \underline{\underline{n \log \nu}} + \log \det W - \underline{\underline{\nu}} \\ \text{s.t.} \quad \nu_k^T W \nu_k \leq 1 \end{array} \right.$$

$$\min_{\nu} \nu - n \log \nu$$

$$1 - \frac{n}{\nu} = 0 \Rightarrow \nu = n$$

$$\left[\begin{array}{l} \max_W \log \det W \\ \text{s.t.} \quad \nu_k^T W \nu_k \leq 1 \end{array} \right.$$