

statistical estimation problems

Maximum Likelihood estimation (MLE)

$$\max_{\underline{\theta} \in \mathbb{R}^r} \underbrace{L(\underline{\theta})}_{\begin{array}{l} \text{Parameter} \\ \text{vector} \end{array}} = \max_{\underline{\theta} \in \mathbb{R}^r} \log f(\underline{x}, \underline{\theta})$$

↗ Log-Likelihood
 ↗ probability density function

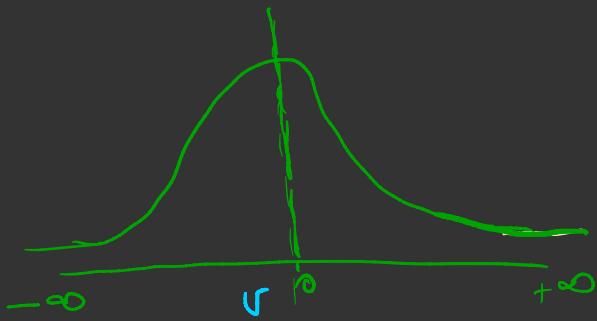
\underline{y} : realization of some random vector

If ρ is a function of $\underline{\Theta}$ (parameter vector)
is Log-concave, then the MLE problem is a convex
optimization problem.

Example : linear measurement with i.i.d noise.
independent,
and identically,
distributed,

$$y_i = \alpha_i^T \underline{\Theta} + \underline{U}_i \quad , i=1, \dots, m$$

measurement
noise, random variable



We assume $v \in \mathbb{R}$
measured noise (i.i.d)
has PDF on \mathbb{R}

$$y_i = \underline{\alpha}_i^T \underline{\theta} + v_i \Rightarrow \underbrace{v_i = y_i - \underline{\alpha}_i^T \underline{\theta}}_{\text{independent}}$$

$$\Rightarrow \overline{f}(y, \underline{\theta}) = \prod_{i=1}^m f(y_i - \underline{\alpha}_i^T \underline{\theta})$$

Joint
PDF

$$\Rightarrow \log \tilde{f} = \sum_{i=1}^m \log f(y_i - \underline{\alpha}_i^\top \underline{\theta})$$

MLE problem

$$\max_{\underline{\theta} \in \mathbb{R}^n} \sum_{i=1}^m \log f(y_i - \underline{\alpha}_i^\top \underline{\theta})$$

convex optimization problem, if f is log-concave
in $\underline{\theta}$ (parameter vector)

Special case: Gaussian noise: $v_i \sim \mathcal{N}(0, \sigma^2)$

$$\Rightarrow p(z) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{z^2}{2\sigma^2}\right)$$

$$L(\underline{\theta}) = -\left(\frac{m}{2}\right) \log(2\pi\sigma^2) - \underbrace{\frac{1}{2\sigma^2} \|A\underline{\theta} - \underline{y}\|_2^2}$$

A has rows $\underline{\alpha}_1^T, \dots, \underline{\alpha}_m^T$

$$\underline{\theta}_{\text{ML}}^* = \arg \min \|A\underline{\theta} - \underline{y}\|_2^2$$

ordinary least square (OLS)

MLE
is robust
against Gaussian
noise

Logistic regression :

Bernoulli a.k.a two class Logistic regression :

n experiment : $\{\underline{x}_i, y_i\}$, $y_i \in \{0, 1\}$

covid test : x_i : age, blood pressure, heart rate

n patients : y_i : $y_i = 0$ No, $y_i = 1$ Yes

spam : x_i : capital letter + attachments

y_i : $y_i = 0$ No, $y_i = 1$ Yes

$\mathbb{Y} \in \{0, 1\}$ with probability mass function: $P(\mathbb{Y}=1) + P(\mathbb{Y}=0) = 1$

Bernoulli(p) , $p \in (0, 1)$, $P(\mathbb{Y}=1) = p$, $P(\mathbb{Y}=0) = 1-p$

Logistic function : $f_{\text{logistic}}(z) = \frac{\exp(z)}{1 + \exp(z)}$ $\in (0, 1)$, $z \in \mathbb{R}$

we want to model p as the composition of f_{logistic} :

$$P := f_{\text{Logistic}}(\underline{\theta}^T \underline{x}) = \frac{\exp(\underline{\theta}^T \underline{x})}{1 + \exp(\underline{\theta}^T \underline{x})} \quad \underline{\theta} \in \mathbb{R}^d$$

to be determined.

Feature vector $\underline{x} \in \mathbb{R}^d$

regressor parameter

N experiment
Reorder the data

$$\underline{x}_1, \dots, \underline{x}_{N_1} \Rightarrow y=1 \quad , \quad \underline{x}_{N_1+1}, \dots, \underline{x}_N \Rightarrow y=0$$

$$f(\underline{\theta}) = \left(\prod_{c=1}^{N_1} P_c \right) \left(\prod_{c=N_1+1}^N (1-P_c) \right)$$

where $P_c = f_{\text{Logistic}}(\underline{\theta}, \underline{x}_c)$

the likelihood function

$$\log p(\underline{\theta}) = \sum_{i=1}^{N_1} \log \frac{\exp(\underline{\theta}, \underline{x}_i)}{1 + \exp(\underline{\theta}, \underline{x}_i)} + \sum_{i=N_1+1}^{N_2} \log \frac{1}{1 + \exp(\underline{\theta}, \underline{x}_i)}$$

$\arg \max_{\underline{\theta} \in \mathbb{R}^d} \log L(\underline{\theta}) \Rightarrow$ This is an unconstrained
convex optimization problem
in variable $\underline{\theta} \in \mathbb{R}^d$
(convexity by Hessian of obj)

The development above can be generalized for the multi-trial case.

categorical a.k.a. a multiclass logistic regression.

$$\{(x_c, \mathbb{1}_c)\}, \quad \mathbb{1} \in \underbrace{\{0, 1, \dots, k\}}_{n+1} \quad \text{the recorded labels}$$

are the realization
of a categorical
random variable

$$\text{PMF cat}(p_1, \dots, p_k) \Rightarrow p_0 = 1 - \sum_{n=1}^k p_n$$

$$P(\mathbb{1}=1) = p_1, \dots, P(\mathbb{1}=k) = p_k \Rightarrow P(\mathbb{1}=0) = 1 - (p_1 + \dots + p_k)$$

Regressor parameter $\underline{\Theta} := \{\underline{\Theta}_1, \dots, \underline{\Theta}_K\}^T \in \mathbb{R}^{d \times K}$

The regression model

$$P_1 = \frac{\exp(\underline{\Theta}_1, \underline{x})}{1 + \sum_{r=1}^K \exp(\underline{\Theta}_r, \underline{x})}, \dots, P(K) = 1 - (P_1 + \dots + P_K)$$
$$= \frac{1}{1 + \sum_{r=1}^K \exp(\underline{\Theta}_r, \underline{x})}$$

$$f(\theta) = \left(\prod_{i=1}^{N_1} p_{1,i} \right) \cdots \left(\prod_{i=N_{k-1}+1}^{N_k} p_{k,i} \right) \left(\prod_{i=N_k+1}^N (1 - (p_{1,i} + \cdots + p_{k,i})) \right)$$

$$p_{ri} = \frac{\exp(\underline{\theta}_r, \underline{x}_i)}{1 + \sum_{r=1}^k \exp(\underline{\theta}_r, \underline{x}_i)}$$

MLE

$$\arg \max_{\underline{\theta} \in \mathbb{R}^d} \log f(\theta) = \sum_{i=1}^{N_1} \underline{\theta}_1^T \underline{x}_i + \sum_{i=N_1+1}^{N_2} \underline{\theta}_2^T \underline{x}_i + \cdots + \sum_{i=1}^N \log \left(1 + \sum_{r=1}^k \exp(\underline{\theta}_r^T \underline{x}_i) \right)$$

A different way to do statistical inference

is MAP (Maximum a Posteriori probability
after the fact estimation)

\underline{x} ← to be estimated } Joint PDF
 \underline{z} ← observed variable } $g(\underline{x}, \underline{z})$

In MLE, \underline{x} was deterministic (non-random)

But in MAP, both \underline{x} and \underline{z} are random

Prior of \underline{x}

$$f_x(\underline{x}) = \int p(\underline{x}, \underline{y}) d\underline{y}$$

before we observe \underline{y}

$$f_{\underline{x}}(\underline{y}) = \int p(\underline{x}, \underline{y}) d\underline{x}$$

conditional:

$$f_{\underline{y}|\underline{x}}(\underline{x}, \underline{y}) = \frac{f(\underline{x}, \underline{y})}{f_x(\underline{x})} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Bayes' rule}$$

$$p(\underline{x}, \underline{y}) = f_{\underline{y}|\underline{x}}(\underline{x}, \underline{y}) f_x(\underline{x})$$

$$\underline{f_x(\underline{y})} = \frac{f(\underline{x}, \underline{y})}{f_{\underline{y}}(\underline{y})} = \frac{f_{\underline{y}|\underline{x}}(\underline{x}, \underline{y}) f_x(\underline{x})}{f_{\underline{y}}(\underline{y})}$$