S&DS 242/542: Theory of Statistics Lecture 12: Maximum likelihood estimation and optimization

Maximum likelihood estimation

The joint PDF or PMF of the data, viewed as a function of the model parameters θ , is called the **likelihood function**.

For data $X_1, \ldots, X_n \stackrel{IID}{\sim} f(x \mid \theta)$ from a parametric model $f(x \mid \theta)$, this is given by

$$\mathsf{lik}(\theta) = f(X_1 \mid \theta) \times \ldots \times f(X_n \mid \theta)$$

The maximum likelihood estimator (MLE) is the value of θ in the parameter space of the model that maximizes lik(θ).

Intuitively, it is the value of θ that makes the observed data "most probable" or "most likely".

Connection to the likelihood ratio statistic

In our discussion of the Neyman-Pearson lemma, recall that for testing

$$H_0: \mathbf{X} \sim f_0$$
 vs. $H_1: \mathbf{X} \sim f_1$

the most powerful test rejects H_0 for large values of the likelihood ratio statistic

$$L(\mathbf{X}) = f_1(\mathbf{X})/f_0(\mathbf{X}).$$

In the context of a parametric model, we may consider testing $H_0: X_1, \ldots, X_n \stackrel{IID}{\sim} f(x \mid \theta_0)$ versus $H_1: X_1, \ldots, X_n \stackrel{IID}{\sim} f(x \mid \theta_1)$ for two different parameter values. Then

$$f_0(X_1,\ldots,X_n) = f(X_1 \mid \theta_0) \times \ldots \times f(X_n \mid \theta_0),$$

$$f_1(X_1,\ldots,X_n) = f(X_1 \mid \theta_1) \times \ldots \times f(X_n \mid \theta_1),$$

so the likelihood ratio statistic is exactly $L(\mathbf{X}) = \text{lik}(\theta_1)/\text{lik}(\theta_0)$.

Connection to empirical risk minimization

Maximizing lik(θ) is equivalent to maximizing its logarithm, the **log-likelihood function**. For data $X_1, \ldots, X_n \stackrel{IID}{\sim} f(x \mid \theta)$, this is

$$\ell_n(\theta) = \log(\operatorname{lik}(\theta)) = \sum_{i=1}^n \log f(X_i \mid \theta)$$

It is usually easier to work with the log-likelihood instead of the likelihood itself, because this involves a sum rather than a product.

The MLE maximizes $\ell_n(\theta)$, or equivalently, minimizes

$$\frac{1}{n}\sum_{i=1}^n -\log f(X_i \mid \theta)$$

This is an example of *empirical risk minimization*, minimizing the average of the loss function $-\log f(x \mid \theta)$ across the observed data.

Maximum likelihood in the Poisson model

Let
$$X_1, \ldots, X_n \stackrel{\text{HD}}{\sim} \text{Poisson}(\lambda)$$
.

$$f(\mathbf{x}|\lambda) = \frac{e^{-\lambda} \cdot \lambda^{\mathbf{x}}}{\mathbf{x}!} \Rightarrow \log f(\mathbf{x}|\lambda) = -\lambda + \mathbf{x} \log \lambda - \log (\mathbf{x}!)$$

$$I_n(\lambda) = \sum_{i=1}^n (-\lambda + \mathbf{x}_i \log \lambda - \log (\mathbf{x}_i!))$$

$$= -n\lambda + (\sum_{i=1}^n \mathbf{x}_i) \log \lambda - \sum_{i=1}^n \log (\mathbf{x}_i!)$$
Solve $O = I'_n(\hat{\mathbf{x}}) = -n + \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$

$$\Rightarrow \hat{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \hat{\mathbf{x}}$$
Check: $J'_n(\lambda) > 0 \quad \text{Ev} \quad \lambda < \hat{\mathbf{x}} \Rightarrow \hat{\mathbf{x}} = \hat{\mathbf{x}} \text{ is fw } MLE$

Maximum likelihood in the Pareto model

Let
$$X_1, \ldots, X_n \stackrel{HD}{\sim} \operatorname{Pareto}(\alpha, 1)$$
.

$$\begin{aligned}
& \quad f(\mathbf{x} \mid a) = \frac{\alpha}{\mathbf{x}^{\alpha+1}} \quad (f_{-r} \times > 1) \Rightarrow \log f(\mathbf{x} \mid a) = \log \alpha - (\alpha + 1) \log \mathbf{x} \\
& \quad l_n(\alpha) = \sum_{\substack{n=1 \\ n \neq n}}^n (1_{ij} \mid \alpha - (\alpha + 1) \log \mathbf{x}_i) \\
& \quad = n \log \alpha - (\alpha + 1) \sum_{\substack{n=1 \\ n \neq n}}^n \log \mathbf{x}_i \\
& \quad \text{Soluc} \quad O = \int_{-1}^n (\widehat{\alpha}) = \frac{n}{\widehat{\alpha}} - \sum_{\substack{n=1 \\ n \neq n}}^n l_{ij} \mathbf{x}_i \\
& \quad \Rightarrow \widehat{\alpha} = \frac{n}{\sum_{\substack{n=1 \\ n \neq n}}^n l_{ij} \mathbf{x}_i}
\end{aligned}$$

Maximum likelihood in the Pareto model



True parameter $\alpha = 2$

This estimate $\hat{\alpha}_{\text{MLE}} = \frac{n}{\sum_{i=1}^{n} \log X_i}$ coincides with a generalized method of moments estimator from last lecture, and is *different* from usual the method of moments estimator $\hat{\alpha}_{\text{MoM}} = \frac{\bar{X}}{\bar{X}-1}$. The variability of $\hat{\alpha}_{\text{MLE}}$ is smaller than that of $\hat{\alpha}_{\text{MoM}}$.

Let $X_1, \ldots, X_n \stackrel{IID}{\sim} \mathcal{N}(\mu, \sigma^2).$ $f(x|_{M},\sigma^{2}) = \frac{1}{\sqrt{2\sigma^{2}}} e^{-\frac{(x-m)^{2}}{2\sigma^{2}}} \Rightarrow \log f(x|_{M},\sigma^{2}) = -\frac{(x-m)^{2}}{2\sigma^{2}} = \log e^{-\frac{1}{2\sigma^{2}}}$ $\mathcal{J}_{m}(M,\sigma^{2}) = \sum_{i=1}^{M} \left(- \frac{(x_{i}, T_{M})^{2}}{(x_{i}, T_{M})^{2}} - \frac{1}{2} \left(- \frac{2}{2} \pi \sigma^{2} \right) \right)$ $= -\frac{1}{2} \sum_{x} \sum_{x} (x, -\mu)^{2} - \frac{\mu}{2} \log \sigma^{2} - \frac{\mu}{2} \log 2\pi$ Solve $O = \frac{\partial I_n}{\partial x} (\hat{\mu}, \hat{\mu}) = \frac{1}{2\pi} \sum_{i=1}^{n} (x_i - \hat{\mu}) = \frac{1}{2\pi} \left((\hat{\Sigma}, \hat{\chi}) - n \hat{\mu} \right)$ $O = \frac{\partial l_n}{\partial x^2} \left(\hat{\mu}_n \hat{\rho}_n^2 \right) = \frac{1}{2} \sum_{i=1}^{n} \left(\hat{\mu}_n \hat{\mu}_n^2 \right)^2 - \frac{1}{2n^2}$

 $O = \frac{1}{2\hat{a}_{1}} = \hat{c}_{1} (\hat{x}_{1} + \hat{x}_{1}) = \frac{n}{2\hat{a}_{1}}$ (شرم - (² ^x ^x)) - (² ^x) $\Rightarrow \frac{n}{2\hat{\sigma}^{2}} = \frac{1}{2\hat{\sigma}^{2}} \sum_{x=1}^{2} (X_{x} - \bar{X})^{2}$ ラル= ニーズ×=× $\Rightarrow \hat{\sigma}^{2} = \int_{X} \sum_{i=1}^{2} (X_{i} - \bar{X})^{2}$

· For each fixed o'>0, in= X monimizes la(M, J) · At in= X, & maximizes ln(X, o2) => (p, 2)= (X, + 2 (X, -X)) is the MLE.

Let
$$X_1, \ldots, X_n \stackrel{HD}{\sim} \text{Gamma}(\alpha, \beta)$$
.

$$f(x \mid \alpha, \beta) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} \times^{\alpha-1} e^{-\beta x}$$

$$\Rightarrow \log f(x \mid \alpha, \beta) = \alpha \log \beta - \log \beta (\alpha) + (\alpha - 1) \log x - \beta x$$

$$f_n(\alpha, \beta) = \sum_{i=1}^{n} (\alpha \log \beta - \log \beta'(\alpha) + (\alpha - 1) \log x - \beta x_i)$$

$$= n \alpha \log \beta - n \log \beta'(\alpha) + (\alpha - 1) \sum_{i=1}^{n} \log x_i - \beta \sum_{i=1}^{n} x_i$$

$$\text{Soluc} \quad 0 = \frac{\partial e_n}{\partial \alpha} (\hat{\alpha}, \hat{\beta}) = n \log \hat{\beta} - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} + \sum_{i=1}^{n} \log x_i$$

$$0 = \frac{\partial e_n}{\partial p} (\hat{\alpha}, \hat{\beta}) = \frac{n \hat{\alpha}}{\beta} - \sum_{i=1}^{n} \chi_i$$

$$O = n \log \hat{\beta} - n \frac{P(\hat{\alpha})}{P(\hat{\alpha})} + \sum_{i=1}^{\infty} \int_{i=1}^{\infty} \chi_{i} = O = \frac{n \hat{\alpha}}{\hat{\beta}} - \sum_{i=1}^{\infty} \chi_{i}$$

$$\Longrightarrow O = n \log \hat{\alpha} - n \frac{P(\hat{\alpha})}{P(\hat{\alpha})} + \sum_{i=1}^{\infty} \int_{i=1}^{\infty} \chi_{i} = \hat{\alpha}$$

$$\Longrightarrow \hat{\beta} = \frac{n \hat{\alpha}}{\sum_{i=1}^{\infty} \chi_{i}} = \hat{\alpha}$$

$$\Longrightarrow O = \log \hat{\alpha} - \frac{P'(\hat{\alpha})}{P(\hat{\alpha})} - \log \chi + \int_{i=1}^{\infty} \int_{i=1}^{\infty} \int_{i=1}^{\infty} \chi_{i}$$

One may check that the function

$$f(\alpha) = \log \alpha - \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} - \log \bar{X} + \frac{1}{n} \sum_{i=1}^{n} \log X_i$$

is decreasing over $\alpha \in (0, \infty)$, and
$$\lim_{\alpha \to 0} f(\alpha) = \infty$$

$$\lim_{\alpha \to \infty} f(\alpha) = -\log \bar{X} + \frac{1}{n} \sum_{i=1}^{n} \log X_i$$

By Jensen's inequality, $\frac{1}{n} \sum_{i=1}^{n} \log X_i < \log \overline{X}$, so $f(\alpha) < 0$ for all large α . Thus $0 = f(\alpha)$ has a unique solution $\hat{\alpha}$. [There is no explicit form, and $\hat{\alpha}$ is usually computed numerically.]

The maximum likelihood estimators are $(\hat{\alpha}, \hat{\beta}) = (\hat{\alpha}, \hat{\alpha}/\bar{X})$.



True parameters $\alpha = 1$, $\beta = 1$

The maximum likelihood estimates are again different from the method of moments estimates, and have smaller variability.

Maximum likelihood in the Multinomial model

Let $(X_1, \ldots, X_k) \sim \text{Multinomial}(n, (p_1, \ldots, p_k))$. Here X_1, \ldots, X_k are not IID, but instead are integer counts summing to n.

$$\begin{split} & L_{n}(p_{1,7}p_{k}) = \log \left[\begin{pmatrix} n \\ \chi_{1},\chi_{1,7}\chi_{1k} \end{pmatrix} \cdot p_{1}^{\chi_{1}} \dots p_{k}^{\chi_{k}} \right] \\ & = \log \begin{pmatrix} n \\ \chi_{1,7}\chi_{1k} \end{pmatrix} + \chi_{1}\log p_{1}t \dots t \chi_{k}\log p_{k} \\ & Lagrange multiplier method : \\ & L_{n}(p_{1,7}p_{k},\lambda) = \mathcal{L}_{n}(p_{2,7}p_{k}) + \lambda (p_{1}t,tp_{k}-1) \\ & = \log (\chi_{1,7}\chi_{1k}) + \chi_{1}\log p_{1}t \dots t \chi_{k}\log p_{k} + \lambda (p_{1}t-p_{k}-1) \\ & = \log (\chi_{1,7}\chi_{1k}) + \chi_{1}\log p_{1}t \dots t \chi_{k}\log p_{k} + \lambda (p_{1}t-p_{k}-1) \\ & \leq \log (\chi_{1,7}\chi_{1k}) + \chi_{1}\log p_{1}t \dots t \chi_{k}\log p_{k} + \lambda (p_{1}t-p_{k}-1) \\ & \leq \log (p_{1,7}p_{k},\lambda) = \sum_{n} (p_{1,7}p_{k},\lambda) + \sum_{n} (p_{1,7}p_{n}) +$$

Maximum likelihood in the Multinomial model

 $O = \frac{\chi_1}{\beta_1} + \int \operatorname{Er} \operatorname{ench} \dot{a^2}_{j-k} = O^2 \hat{p}_{j+1} + \hat{p}_{w} - 1$

=) p:=- X: f: j=1,5k

 $= 0 = -\frac{x_1 + x_k}{r} - 1$

カデューの

=) p: = X:

I.e. (p1, pk) = (X1 Xk)

Maximum likelihood in the Multinomial model

The formal reasoning behind this Lagrange multiplier method is:

- Fixing any λ ∈ ℝ, maximizing ℓ_n(p₁,..., p_k) subject to the constraint p₁ + ... + p_k = 1 is the same as maximizing L_n(p₁,..., p_k, λ) subject to this constraint, because the additional term in L_n(p₁,..., p_k, λ) is 0.
- If we instead ignore the constraint p₁ + ... + p_k = 1, then the unconstrained maximizer of L_n(p₁,..., p_k, λ) for each λ is (p₁,..., p_k) = −(1/λ)(X₁,..., X_k) as we computed.
- ► Using the specific choice λ = −n, this unconstrained maximizer satisfies the constraint p₁ + ... + p_k = 1, so it must also be the constrained maximizer of L_n(p₁,..., p_k, λ).

Combining these three statements, $(p_1, \ldots, p_k) = (X_1, \ldots, X_k)/n$ is the constrained maximizer of $\ell_n(p_1, \ldots, p_k)$.

The Hardy-Weinberg model

In genetics, it is often assumed that the genotypes AA, Aa, and aa at a single locus satisfy *Hardy-Weinberg equilibrium* — they occur with probabilities $(1-p)^2$, 2p(1-p), and p^2 , where $p \in [0,1]$ represents the frequency of the minor allele. Thus the counts of these three genotypes in *n* samples may be modeled as

$$(X_{1}, X_{2}, X_{3}) \sim \text{Multinomial}\left(n, \left((1-p)^{2}, 2p(1-p), p^{2}\right)\right)$$

$$l_{n}(p) = \log\left[\left(\begin{pmatrix}n\\X_{i}, X_{i}, X_{3}\end{pmatrix}, \left((1-p)^{j}\right)^{X_{i}}\left(2p(1-p)\right)^{X_{2}}\left(p^{j}\right)^{X_{3}}\right]\right]$$

$$= \log\left(\begin{pmatrix}n\\X_{i}, X_{i}, X_{3}\end{pmatrix} + \left(2X_{i} + X_{i}\right) \log(1-p) + \left(2X_{3} + X_{2}\right) \log p\right)$$

$$+ X_{2} \log 2$$

$$\Rightarrow O = l_{n}(p) = \frac{2X_{i} + X_{2}}{1-p} + \frac{2K_{3} + K_{2}}{p} \Rightarrow p = \frac{2X_{3} + K_{2}}{2n}$$

Gradient ascent

Computing the MLE is an optimization problem: For $\theta \in \mathbb{R}$, we wish to maximize $\ell_n(\theta)$, or to solve the *score equation* $0 = \ell'_n(\theta)$.

When there is no explicit solution, we oftentimes use numerical optimization procedures. The simplest procedure is **gradient ascent**: Starting with an initial guess $\theta^{(0)}$, iterate

$$\theta^{(t+1)} = \theta^{(t)} + \eta \,\ell'_n(\theta^{(t)})$$

where $\eta > 0$ is a learning rate parameter.

When
$$0 = \ell'_n(\theta^{(t)})$$
, this gives $\theta^{(t+1)} = \theta^{(t)}$.

Gradient ascent



Newton's method

A second-order procedure is **Newton's method**: Given $\theta^{(t)}$, approximate the solution to $0 = \ell'_n(\theta)$ by a Taylor expansion

$$0=\ell_n'(heta)pprox\ell_n'(heta^{(t)})+\ell_n''(heta^{(t)})(heta- heta^{(t)}).$$

Solve this equation in θ and set this as $\theta^{(t+1)}$:

$$\theta^{(t+1)} = \theta^{(t)} - \frac{\ell'_n(\theta^{(t)})}{\ell''_n(\theta^{(t)})}$$

When $0 = \ell'_n(\theta^{(t)})$, again this gives $\theta^{(t+1)} = \theta^{(t)}$.

Newton's method

 $\ell_n'(\theta)$





Optimization in higher dimensions

For $\theta \in \mathbb{R}^k$, gradient ascent is the procedure

$$\theta^{(t+1)} = \theta^{(t)} + \eta \,\nabla \ell_n(\theta^{(t)})$$

where $\nabla \ell_n(\theta) \in \mathbb{R}^k$ is the gradient of $\ell_n(\theta)$, i.e. the vector of its 1st-order partial derivatives.

Newton's method is the procedure

$$\theta^{(t+1)} = \theta^{(t)} - [\nabla^2 \ell_n(\theta^{(t)})]^{-1} \nabla \ell_n(\theta^{(t)})$$

where $\nabla^2 \ell_n(\theta) \in \mathbb{R}^{k \times k}$ is the *Hessian* of $\ell_n(\theta)$, i.e. the matrix of its 2nd-order partial derivatives.