

Maximum likelihood estimation (part 1)

Lecture 13b (STAT 24400 F24)

1 / 19

The likelihood & log likelihood

Setting: $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(\cdot | \theta)$ for an unknown parameter θ

The joint density or PMF of (X_1, \dots, X_n) is:

$$\prod_{i=1}^n f(X_i | \theta) = f(X_1 | \theta) \cdot \dots \cdot f(X_n | \theta)$$

↖
called the **likelihood** of θ given the data

Sometimes it's convenient to work with the log likelihood:

$$\log \left(\prod_{i=1}^n f(X_i | \theta) \right) = \sum_{i=1}^n \log (f(X_i | \theta))$$

2 / 19

The likelihood & log likelihood (notations)

Switching notation:

- $\theta_0 \in \Theta$ is the unknown true value of the parameter
- $\theta \in \Theta$ represents any possible value of the parameter
(so that we can study the function $\text{Likelihood}(\theta)$, over all $\theta \in \Theta$, even though θ_0 is fixed)

3 / 19

The maximum likelihood estimator (MLE)

The maximum likelihood estimator (MLE) is the value of θ that maximizes the likelihood:

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} \left\{ \prod_{i=1}^n f(X_i | \theta) \right\}$$

↖
 $\max_x h(x)$ = maximum value of $h(x)$

$\operatorname{argmax}_x h(x)$ = value of x that yields maximum value of $h(x)$

Often more convenient to work with log likelihood:

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} \left\{ \sum_{i=1}^n \log (f(X_i | \theta)) \right\}$$

4 / 19

Example: MLE of Normal mean (σ^2 known)

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ for unknown $\mu \in \mathbb{R}$ (σ^2 is known)

- The density is $f(x | \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$
- Likelihood function = $\prod_{i=1}^n f(x_i | \mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_i-\mu)^2/2\sigma^2}$
- Log likelihood:

$$\sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_i-\mu)^2/2\sigma^2} \right) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_i (X_i - \mu)^2$$

- Solve for MLE:

$$0 = \frac{\partial}{\partial \mu} \left(-\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_i (X_i - \mu)^2 \right) = \frac{1}{\sigma^2} \sum_i (X_i - \mu)$$

$$\Rightarrow \hat{\mu} = \bar{X} \text{ (same as MoM)}$$

Check: it is a global maximum point (by 2nd derivative test, 1st derivative test, etc.)

5 / 19

Example: MLE of Normal μ and σ^2

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ for unknown $\mu \in \mathbb{R}$, $\sigma^2 > 0$

- The density is $f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$
- Log likelihood:

$$\sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-(X_i-\mu)^2/2\sigma^2} \right) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_i (X_i - \mu)^2$$

- Solve for MLE ($\hat{\mu}, \hat{\sigma}^2$):

$$0 = \frac{\partial}{\partial \mu} \left(-\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_i (X_i - \mu)^2 \right) = \frac{1}{\sigma^2} \sum_i (X_i - \mu)$$

$$0 = \frac{\partial}{\partial (\sigma^2)} \left(-\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_i (X_i - \mu)^2 \right) = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_i (X_i - \mu)^2$$

$$\Rightarrow \hat{\mu} = \bar{X}, \hat{\sigma}^2 = \frac{1}{n} \sum_i (X_i - \bar{X})^2 \text{ (same as MoM)}$$

6 / 19

Example: MLE for Exponential rate

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Exponential}(\lambda)$ for unknown $\lambda > 0$

- The density is (using a new, useful notation)

$$f(x | \lambda) = \lambda e^{-\lambda x} \cdot \mathbb{1}_{x>0}$$

- Log likelihood:

since $X_i > 0$ for all i
↓

$$\sum_{i=1}^n \log (\lambda e^{-\lambda X_i} \cdot \mathbb{1}_{X_i>0}) = \sum_{i=1}^n \log (\lambda e^{-\lambda X_i}) = n \log(\lambda) - \lambda \sum_i X_i$$

- Solve for MLE:

$$0 = \frac{\partial}{\partial \lambda} \left(n \log(\lambda) - \lambda \sum_i X_i \right) = \frac{n}{\lambda} - \sum_i X_i$$

$$\Rightarrow \hat{\lambda} = 1/\bar{X} \text{ (same as MoM)}$$

Check: it is a global max.

7 / 19

Example: MLE for Binomial

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$ for unknown $p \in (0, 1) \rightsquigarrow \sum_i X_i \sim \text{Binomial}(n, p)$

- The PMF is (using a new expression)

$$f(x | p) = p^x (1-p)^{1-x}$$

- Log likelihood:

$$\sum_{i=1}^n \log (p^{X_i} (1-p)^{1-X_i}) = \sum_i X_i \cdot \log(p) + \sum_i (1-X_i) \cdot \log(1-p)$$

- Solve for MLE:

$$0 = \frac{\partial}{\partial p} \left(\sum_i X_i \cdot \log(p) + \sum_i (1-X_i) \cdot \log(1-p) \right) = \frac{\sum_i X_i}{p} - \frac{\sum_i (1-X_i)}{1-p}$$

$$\Rightarrow \hat{p} = \bar{X} \text{ (= the proportion of successes in the sample)}$$

8 / 19

Example: MLE for Uniform

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Uniform}[0, \theta]$ for unknown $\theta > 0$

- The density is

$$f(x | \theta) = \frac{1}{\theta} \cdot \mathbb{1}_{0 \leq x \leq \theta}$$

- Likelihood:

$$\prod_{i=1}^n \left(\frac{1}{\theta} \cdot \mathbb{1}_{0 \leq X_i \leq \theta} \right) \stackrel{\substack{\text{since } X_i \geq 0 \text{ for all } i \\ \downarrow}}{=} \theta^{-n} \cdot \mathbb{1}_{\theta \geq \max_i X_i}$$

- Solve for MLE:

$$\hat{\theta} = \max_i X_i = X_{(n)} \text{ (not same as MoM)}$$

9 / 19

MoM vs MLE (Example for Uniform)

Compare MoM and MLE estimators for Uniform $[0, \theta]$:

MoM estimator $\hat{\theta} = 2\bar{X}$:

- For each X_i , $\mathbb{E}(X_i) = \frac{\theta}{2}$, $\text{Var}(X_i) = \frac{\theta^2}{12}$

- Bias:

$$\mathbb{E}(\hat{\theta}) = 2\mathbb{E}(\bar{X}) = 2 \cdot \frac{\theta}{2} = \theta \Rightarrow \text{bias} = 0$$

- Variance:

$$\text{Var}(\hat{\theta}) = \text{Var}\left(\frac{2}{n} \sum_i X_i\right) = \frac{4}{n^2} \text{Var}\left(\sum_i X_i\right) = \frac{4}{n^2} \sum_i \text{Var}(X_i) = \frac{\theta^2}{3n}$$

- MSE = bias² + Var($\hat{\theta}$) = $\frac{\theta^2}{3n}$

10 / 19

MoM vs MLE (Example for Uniform)

Compare MoM and MLE estimators for Uniform $[0, \theta]$:

MLE estimator $\hat{\theta} = X_{(n)}$:

- The sampling distribution has density $f(x) = \frac{nx^{n-1}}{\theta^n} \cdot \mathbb{1}_{0 \leq x \leq \theta}$

- Bias:

$$\mathbb{E}(\hat{\theta}) = \int_{x=0}^{\theta} x \cdot \frac{nx^{n-1}}{\theta^n} dx = \frac{n\theta}{n+1} \Rightarrow \text{bias} = -\frac{\theta}{n+1}$$

- Variance:

$$\begin{aligned} \mathbb{E}(\hat{\theta}^2) &= \int_{x=0}^{\theta} x^2 \cdot \frac{nx^{n-1}}{\theta^n} dx = \frac{n\theta^2}{n+2} \\ \Rightarrow \text{Var}(\hat{\theta}) &= \frac{n\theta^2}{n+2} - \left(\frac{n\theta}{n+1}\right)^2 = \frac{n\theta^2}{(n+1)^2(n+2)} \end{aligned}$$

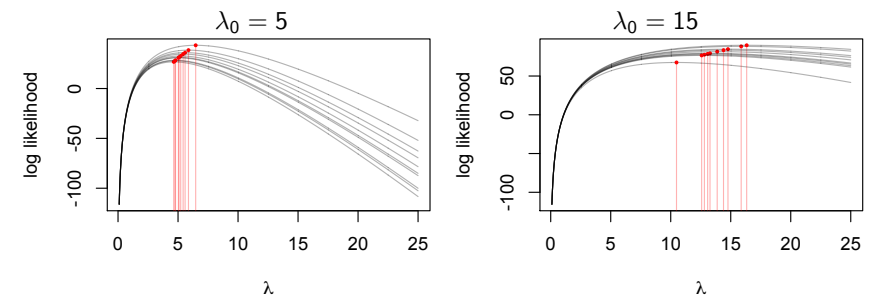
- MSE = bias² + Var($\hat{\theta}$) = $\frac{2\theta^2}{(n+1)(n+2)}$ (comparison: By MoM, MSE = $\frac{\theta^2}{3n}$)

11 / 19

Accuracy of the MLE

Example: suppose $X_1, \dots, X_{50} \stackrel{\text{iid}}{\sim} \text{Exponential}(\lambda)$

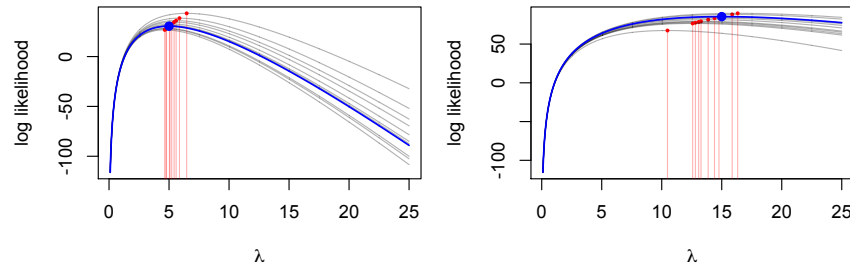
Here is a plot of the log likelihood function, and the MLE, over 10 trials:



12 / 19

Accuracy of the MLE

Add in the expected log likelihood curve (with λ_0 highlighted):



\Rightarrow higher curvature leads to a more accurate estimate

13 / 19

Fisher information (Example: Exponential rate)

(Now we are going to work with one-dimensional θ .)

The **Fisher information** is defined as:

$$\mathcal{I}(\theta) = \mathbb{E} \left(\left(\frac{\partial}{\partial \theta} \log(f(X | \theta)) \right)^2 \right) = \mathbb{E} \left(-\frac{\partial^2}{\partial \theta^2} \log(f(X | \theta)) \right)$$

\uparrow
 with some regularity conditions
 (smoothness of $\log(f)$ as a function of θ)

Example:

- Exponential(λ):

$$-\frac{\partial^2}{\partial \lambda^2} \log(f(X | \lambda)) = -\frac{\partial^2}{\partial \lambda^2} \log(\lambda e^{-\lambda X}) = -\frac{\partial^2}{\partial \lambda^2} [\log(\lambda) - \lambda X] = \frac{1}{\lambda^2}$$

$$\mathcal{I}(\lambda) = \mathbb{E} \left(-\frac{\partial^2}{\partial \lambda^2} \log(f(X | \lambda)) \right) = \mathbb{E} \left(\frac{1}{\lambda^2} \right) = \frac{1}{\lambda^2}$$

14 / 19

Fisher information example (Normal mean)

Example:

- $N(\mu, \sigma^2)$ with σ^2 known:

$$\begin{aligned} -\frac{\partial^2}{\partial \mu^2} \log(f(X | \mu)) &= -\frac{\partial^2}{\partial \mu^2} \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-(X-\mu)^2/2\sigma^2} \right) \\ &= -\frac{\partial^2}{\partial \mu^2} \left[-\frac{1}{2} \log(2\pi\sigma^2) - \frac{(X-\mu)^2}{2\sigma^2} \right] = \frac{1}{\sigma^2} \end{aligned}$$

$$\mathcal{I}(\mu) = \mathbb{E} \left(\frac{1}{\sigma^2} \right) = \frac{1}{\sigma^2}$$

15 / 19

Fisher information example (Bernoulli)

Example:

- Bernoulli(p):

$$\begin{aligned} -\frac{\partial^2}{\partial p^2} \log(f(X | p)) &= -\frac{\partial^2}{\partial p^2} \log(p^X (1-p)^{1-X}) \\ &= -\frac{\partial^2}{\partial p^2} [X \log(p) + (1-X) \log(1-p)] = \frac{X}{p^2} + \frac{1-X}{(1-p)^2} \end{aligned}$$

$$\mathcal{I}(p) = \mathbb{E} \left(-\frac{\partial^2}{\partial p^2} \log(f(X | p)) \right) = \frac{\mathbb{E}(X)}{p^2} + \frac{1 - \mathbb{E}(X)}{(1-p)^2} = \frac{1}{p(1-p)}$$

16 / 19

Asymptotic distribution of the MLE — Fisher's Theorem

The Fisher information determines the (approximate) variance of the MLE.

Informally: if $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(\cdot | \theta_0)$ and $\hat{\theta}$ is the MLE,

the distribution of $\hat{\theta}$ is $\approx N(\theta_0, \frac{1}{n\mathcal{I}(\theta_0)})$

More formally: under some regularity conditions,

$$\sqrt{n\mathcal{I}(\theta_0)} \cdot (\hat{\theta} - \theta_0) \rightarrow N(0, 1)$$

“converges in distribution”

This means that the CDF converges — i.e., for all fixed $x \in \mathbb{R}$,

$$\mathbb{P}(\sqrt{n\mathcal{I}(\theta_0)} \cdot (\hat{\theta} - \theta_0) \leq x) \rightarrow \Phi(x)$$

The same holds with $\mathcal{I}(\hat{\theta})$ in place of $\mathcal{I}(\theta_0)$: (very useful in practice)

$$\sqrt{n\mathcal{I}(\hat{\theta})} \cdot (\hat{\theta} - \theta_0) \rightarrow N(0, 1)$$

17 / 19

Asymptotic distribution of the MLE (examples)

Examples:

- Exponential(λ): $\hat{\lambda} = 1/\bar{X}$ and $\mathcal{I}(\lambda) = 1/\lambda^2$, so:

$$\hat{\lambda} \approx N(\lambda_0, \frac{\lambda_0^2}{n}) \text{ or } \approx N(\lambda_0, \frac{\hat{\lambda}^2}{n})$$

- $N(\mu, \sigma^2)$ with σ^2 known: $\hat{\mu} = \bar{X}$ and $\mathcal{I}(\mu) = 1/\sigma^2$ so:

$$\hat{\mu} \approx N(\mu_0, \frac{\sigma^2}{n})$$

(In fact, in this case we know this is the *exact* distribution!)

- Bernoulli(p): $\hat{p} = \bar{X}$ and $\mathcal{I}(p) = \frac{1}{p(1-p)}$, so:

$$\hat{p} \approx N(p_0, \frac{p_0(1-p_0)}{n}) \text{ or } \approx N(p_0, \frac{\hat{p}(1-\hat{p})}{n})$$

18 / 19

Asymptotic distribution of the MLE (counterexamples)

Examples:

- Uniform $[0, \theta]$: in this case the regularity conditions do not hold.

We need $\log(f(X | \theta))$ to be smooth as a function of θ ,
but $\log(f(X | \theta)) = \log(0) = -\infty$ if $X > \theta$.

To confirm the theorem doesn't hold:

We've calculated $\text{Var}(\hat{\theta}) = \frac{n\theta^2}{(n+1)^2(n+2)} = \mathcal{O}(\frac{1}{n^2})$,

while asymptotic normality of the MLE would yield $\text{Var}(\hat{\theta}) = \mathcal{O}(\frac{1}{n})$

In fact, no approximation is needed here, since we actually know the *exact* distribution of the MLE in this case (via order statistics)

19 / 19