

Lecture 4 - scribbles

Friday, September 15, 2017
13:23

today: - continue Bayesian approach
- MLE

(continuation from last class)

as a Bayesian, posterior $p(\theta | x=z)$ contains all info we need
observation

eg. here is a question: what is the probability that the next coin flip is = 1

$$P(\text{next flip} = 1 | \theta) = \theta$$

\uparrow \uparrow
 F $\theta = \theta$

as a Bayesian: $P(F=1 | X=x) = \int_{\Theta} P(F=1, \theta | X=x) d\theta$

(always true)
 $= \int_{\Theta} P(F=1 | \theta, X=x) P(\theta | X=x) d\theta$
|| (by our model)

$\underbrace{P(F=1 | \theta)}_{\theta}$

conditional expectation
 \downarrow

$= \int_{\Theta} \theta P(\theta | X=x) d\theta = \mathbb{E}[\theta | X=x]$

"posterior mean" of θ

a meaningful Bayesian estimator of θ

is $\hat{\theta}_{\text{Bayes}}(x) \triangleq \mathbb{E}[\theta | X=x]$ (posterior mean)

... $\hat{\theta}$... statistical

relation: $\hat{\theta} : \text{obs} \rightarrow \mathbb{R}^D$ statistical "estimator"

"frequentist statistics"
(traditional statistics)

consider multiple possible estimators

- MLE
- moment matching
- Bayesian posterior mean
- MAP

and then analyze their ^{statistical} properties:

- biased?
- variance?
- consistent?

coming back to coin example: $p(\theta | X=x) = \text{Beta}(\theta; \alpha=x+1, \beta=n-x+1)$ parameters

mean of a Beta P.V. = $\frac{\alpha}{\alpha+\beta}$

thus $E[\theta | X=x] = \frac{x+1}{n+2} = \hat{\theta}_{\text{Bayes}}(x)$

here is a binomial P.V.

compare with $\hat{\theta}_{\text{MLE}}(x) = \frac{x}{n}$

here (biased, but asymptotically unbiased)

(unbiased)

ie. $E_{X|\theta}[\hat{\theta}_{\text{MLE}}(X)] = \theta$

Maximum Likelihood principle

setup: given a parametric family $p(x; \theta)$ for $\theta \in \Theta$

we want to estimate θ

$$\hat{\theta}_{ML}(x) \triangleq \underset{\theta \in \Theta}{\operatorname{argmax}} p(x; \theta) \quad \text{i.e. } \hat{\theta}_{ML}(x) \text{ maximizes } p(x; \cdot)$$

$\hookrightarrow \triangleq L(\theta)$ "likelihood function (θ)"

example: n coin flips $\Omega_x = 0:n$

$$X \sim \text{Bin}(n, \theta) \quad p(x; \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

trick: to maximize $\log L(\theta)$ instead of $L(\theta)$

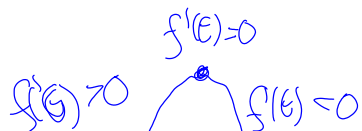
justification: $\log(\cdot)$ is strictly increasing

$$\text{i.e. } a < b \iff \log a < \log b$$

$$\implies \underset{\theta \in \Theta}{\operatorname{argmax}} \log(p(x; \theta)) = \underset{\theta \in \Theta}{\operatorname{argmax}} p(x; \theta)$$

log-likelihood $\log p(x; \theta)$

$$= \underbrace{\log \binom{n}{x}}_{\text{constant}} + x \log \theta + (n-x) \log(1-\theta) = \ell(\theta)$$

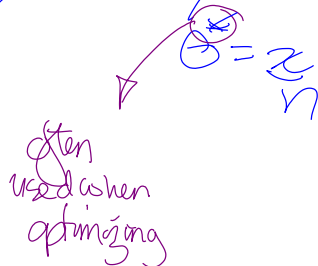




look for $\frac{dL(\theta)}{d\theta} = 0$

i.e. $\frac{x}{\theta} - \frac{n-x}{1-\theta} = 0$

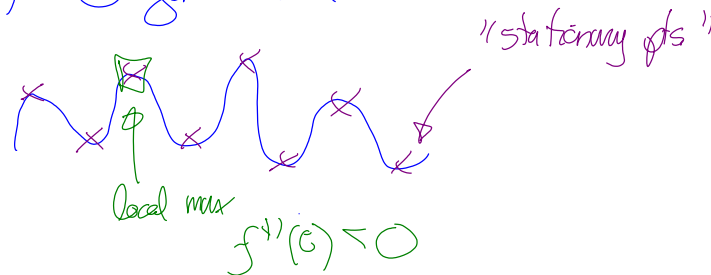
$x - \theta x - (n-x)\theta = 0$



here $\hat{\theta}_{MLE}(x) = \frac{x}{n}$
i.e. relative frequency

some optimization comments:

$(\nabla f(\theta) = 0)$ • $f'(\theta) = 0$ is necessary condition for local max when θ in interior(Θ)
→ also need to check $f''(\theta) < 0$ for a local max



defn matrix

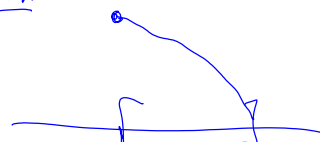
→ only local result in general

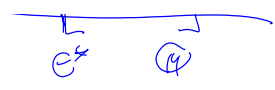
$(\text{Hessian}(\theta) \preceq 0)$
negative
semi-definite

but if $f'(\theta) \leq 0 \forall \theta \in \Theta$, function is "concave"

in this case, $f'(\theta) = 0$ is sufficient for global max
condition

• be careful with boundary cases i.e. $\theta^* \in \text{boundary}(\Theta)$





Some notes about MLE

- does not always exist [e.g. $\hat{\theta} \in \text{bd}(\Theta)$ but Θ is open]

e.g. $\Theta =]0, 1[$

- is not necessarily unique (i.e. multiple maxima)

- is not "admissible" in general [see next class]

Example #2: Multinomial dist.

suppose X_i is discrete R.V. on K choices \rightarrow Multinoulli

(we could choose $\Omega_{X_i} = \{1, \dots, K\}$)

but instead, convenient to encode with the unit basis in \mathbb{R}^K

i.e. $\Omega_{X_i} = \{e_1, \dots, e_K\}$ where $e_i \in \mathbb{R}^K$ "one hot encoding"

$e_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$ i th coordinate

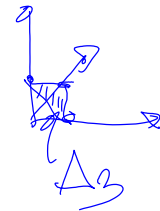
parameter for discrete R.V. $\pi \in \Delta_K$

probability simplex on K choices

$\Delta_K \triangleq \left\{ \pi \in \mathbb{R}^K : \pi_j \geq 0 \forall j, \sum_{j=1}^K \pi_j = 1 \right\}$

$\Theta = \Delta_K$

we will write $X_i \sim \text{Mult}(\pi)$ "multinoulli"
 \uparrow
 parameter



* consider $X_i \stackrel{iid}{\sim} \text{Mult}(\pi)$

then $X = \sum_{i=1}^n X_i \sim \text{Mult}(n, \pi)$ (analog of binomial for k classes)
 \downarrow
 $\in \mathbb{N}^k$ "multinomial"

$$\Omega_X = \left\{ (n_1, \dots, n_k) : \sum_{j=1}^k n_j = n \right\}$$

consider MLE for R.V. $X = \sum_{i=1}^n X_i$
 \uparrow
 φ
 (vector)

$$p(x|\pi)$$

$$p(x_i|\pi) = \text{Mult}(x_i|\pi) = \prod_{j=1}^k \pi_j^{x_{ij}}$$

$\leftarrow j^{\text{th}}$ component of vector π_i

$$p(x|\pi) = \prod_{i=1}^n p(x_i|\pi) \stackrel{\text{by indep.}}{=} \prod_{i=1}^n \left(\prod_{j=1}^k \pi_j^{x_{ij}} \right)$$

$$= \prod_{j=1}^k \left(\prod_{i=1}^n \pi_j^{x_{ij}} \right)$$

$$= \prod_{j=1}^k \pi_j^{\sum_{i=1}^n x_{ij}}$$

log-likelihood: $l(\pi) = \log p(x|\pi) = \sum_{j=1}^k n_j \log \pi_j$
 \hookrightarrow note that $n_j(x)$

we want $\max_{\pi} l(\pi)$
 s.t. $\pi \in \Delta_k$ } constraints

two options a) could reparameterize with $\pi_1, \dots, \pi_{k-1} \in [0, 1]$
 with constraint $\sum_{j=1}^{k-1} \pi_j \leq 1$
 and use $\pi_k = 1 - \sum_{j=1}^{k-1} \pi_j$



and then do
 unconstrained optimization over π_1, \dots, π_{k-1}

b) use Lagrange multiplier approach for equality constraint

look for stationary points (0 gradient) of

$$J(\pi, \lambda) \triangleq \underbrace{l(\pi)}_{g(\pi)} + \lambda \left(1 - \sum_{j=1}^k \pi_j\right)$$

Lagrange multiplier

$$\sum_{j=1}^k \pi_j = 1$$

$$1 - \sum_{j=1}^k \pi_j = 0$$

$$g(\pi) = 0$$

$$\max f(\pi)$$

$$\text{s.t. } g(\pi) = 0$$

ie. want $\nabla_{\pi} J(\pi, \lambda) = 0$
 and $\nabla_{\lambda} J(\pi, \lambda) = 0 \rightarrow$ equivalent $g(\pi) = 0$

$$l(\pi) = \sum_{j=1}^k n_j \log \pi_j \quad \frac{\partial J}{\partial \pi_j} = 0 \Rightarrow \frac{n_j}{\pi_j} - \lambda = 0$$

$$\Rightarrow \pi_j^* = \frac{n_j}{\lambda} \text{ [scaling constant]}$$

$$\text{want } \sum_j \pi_j^* = 1 \Rightarrow \lambda = \sum_j n_j = n$$

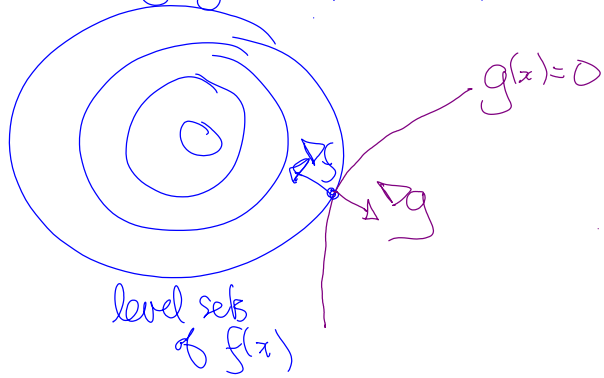
$$\Rightarrow \boxed{\pi_j^* = \frac{n_j}{n}}$$

MLE for multinomial

Side note: $\pi = n^{-1} \mathbf{n}$

Sidenote: $\pi_i = \frac{h_i}{h} \geq 0$

picture behind Lagrange multiplier technique



$$J(x, \lambda) = f(x) + \lambda g(x)$$
$$\nabla_x J(x, \lambda) = 0 \Rightarrow \nabla f(x) = -\lambda \nabla g(x)$$