Youtube Lecture - 17/10/16

Lecture notes 9: Asympt. Norm (cont.)

(47)

### Review

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N\left(0, \frac{1}{I(\theta)}\right)$$

Theorem 12  when $\hat{\theta}_n$ is MLE

$$\hat{\theta}_n = \theta + O_p\left(\frac{1}{\sqrt{n}}\right)$$

· regularity conditions: (for asymp. Norm.)
  - density three times diff. wrt $\theta$
  - key exception is uniform; where range depends on $\theta$
- Good proof :- (from 10/10/16)
- Not: $\ell(\hat{\theta}) = \ell(\hat{\theta}_n) = \ell(\hat{\theta}_{mle})$

$$\ell'(\hat{\theta}) = \ell'(\theta) + (\hat{\theta} - \theta)\ell''(\theta) + \dots = 0$$

$$\sqrt{n}(\hat{\theta} - \theta) \approx \frac{\frac{1}{\sqrt{n}}\ell'(\theta)}{-\frac{1}{n}\ell''(\theta)} \quad = \frac{A}{B}$$

consider A, B :-

$$A = \frac{1}{\sqrt{n}}\ell'(\theta) = \sqrt{n} \cdot \frac{1}{n}\sum_{i=1}^{n} S(\theta, X_i) = \sqrt{n} \cdot \frac{1}{n}\sum_{i=1}^{n}\left(S(\theta, X_i) - 0\right)$$

$$= \sqrt{n}(\bar{S} - 0) \qquad (*)$$

· $S(\theta, X_i)$ - score fn based on single $X_i$

· $S(\theta, X_i) = \frac{\partial}{\partial\theta}\log p(X_i, \theta)$

- (*) Also, derivative of log-like over all data $(X_1, \dots, X_n)$ is sum of derivatives of log-likelihood over a single instance wrt param. (due to IID obs)

$$\ell'(\theta) = \sum_{i=1}^{N}\frac{\partial\log p(X_i, \theta)}{\partial\theta} = \sum_{i=1}^{N}S(X_i, \theta)$$

- $\bar{S}_n = \frac{1}{n}\sum_{i=1}^{n}S(\theta, X_i)$

- Sove know form

(*): sample mean / average of IID r.v.s centred and rescaled

   - Think CLT

**via CLT:-**

$$A = \sqrt{n}(\bar{S} - 0) \xrightarrow{d} N(0, I(\theta))$$

$\bar{S}$

(.) $\mathbb{E}[S(\theta, x_i)] = 0$

$\text{var}[S(\theta, x_i)] = I(\theta)$

(*) Key notational trick you've seen; but is explained explicitly:-

$$Y \sim N(0, \sigma^2)$$

      - In general for mean $0$, $\sigma^2$ variance Normal r.v.

$$\Rightarrow \quad Y = \sigma Z \quad \text{where} \quad Z \sim N(0, 1)$$

(*) can always write a general Normal r.v. as a standard deviation times standard Normal.

**Hence:**
$$A = \sqrt{n}(\bar{S} - 0) \xrightarrow{d} N(0, I(\theta))$$
$$= \sqrt{I(\theta)} \; Z$$

(A) - Tidy up notation for consistency eg indexing
   - espere.

**consider B:**

$$B = \frac{1}{n} \ell''(\theta) \xrightarrow{p} -\mathbb{E}\left[\frac{\partial^2 \log p(X, \theta)}{\partial \theta^2}\right] = I(\theta)$$

· $\ell'' = \frac{\partial^2}{\partial \theta^2} \ell(\theta)$
      sum over $n$ observations

(*) classic trick - exchange $\frac{\partial^2}{\partial \theta^2}$, $\Sigma$ ; sove have $\ell''(\theta)$ is sum over $n$
   2nd derivatives

(*) Hence $\frac{1}{n}\ell''(\theta)$ is an average ('sample mean')

(*) So $B \xrightarrow{P} I(\theta)$

Putting this together:-

$$A \xrightarrow{d} \sqrt{I(\theta)} \, Z \qquad B \xrightarrow{P} I(\theta) = c \text{ (constant)} \qquad \left(\text{and note } \xrightarrow{P} \Rightarrow \xrightarrow{d}\right)$$

Hence by Slutsky:

$$\frac{A}{B} \xrightarrow{d} \frac{\sqrt{I(\theta)} \, Z}{I(\theta)} = \frac{Z}{I(\theta)} \sim N\left(0, \frac{1}{I(\theta)}\right) \quad \blacksquare$$

(*) definitions of score, fisher info arise naturally

Here:- $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N\left(0, \frac{1}{I(\theta)}\right)$

informally: $\hat{\theta} \approx N\left(0, \frac{1}{nI(\theta)}\right) = N\left(0, \frac{1}{I_n(\theta)}\right)$

(*) v. important:

- If you want variance of MLE, approximate as above, it is approx $\frac{1}{I_n(\theta)}$

- standard deviation of error i.e. standard error $\overset{approx}{B}$:-

$$se(\hat{\theta}) \approx \sqrt{\frac{1}{I_n(\theta)}} = \sqrt{\frac{1}{nI(\theta)}} \quad \text{use Slutsky}$$

(*) we do not know $\theta$ in practice, so we use estimates, to get estimated standard error

in practice $\hat{se} = se(\hat{\theta}) = \sqrt{\frac{1}{I_n(\hat{\theta})}}$

/ odd ratio is going to 1 in prob.

(*) via Slutsky, we insert an estimate of 'this' - which?

- use fisher information evaluated at MLE $\hat{\theta}$ - $I_n(\hat{\theta})$

- As long as fisher info is smooth function (contin function, then by continuous mapping theorem $\hat{\theta} \xrightarrow{P} \theta$ will imply

(A2) - some steps missing
- clarify
- Do we need cont. mapping here to extend Slutsky to division?

that $\hat{se}$ is a consistent estimator of se.

(*) we can also write :-

$$\hat{se} = \sqrt{\frac{1}{I_n(\hat{\theta})}} = \sqrt{\frac{1}{n I(\hat{\theta})}}$$

(due to property of fisic info relating individual - samp) obs.

(ii) : This is the most common method in science for computing standard errors

---

## Theorem 14

- let $\tau$ be a smooth function of $\theta$ (via delta method)

- then:- $\sqrt{n}(\tau(\hat{\theta}_n) - \tau(\theta)) \xrightarrow{d} N\left(0, \frac{(\tau'(\theta))^2}{I(\theta)}\right)$

(*) Say interested in function of $\theta$ , $\tau(\theta)$

- The standard error of $\hat{\tau} = \tau(\hat{\theta})$ is :-

$$se = \sqrt{\frac{|\tau'(\theta)|^2}{n I(\theta)}} = \sqrt{\frac{|\tau'(\theta)|^2}{I_n(\theta)}}$$

- The estimated standard error :-

$$\hat{se}(\hat{\tau}) = \sqrt{\frac{|\tau'(\hat{\theta})|^2}{I_n(\hat{\theta})}}$$

- LW: can get MLE of parameter, function of that param
   standard error for MLE of param; function of param

---

## Example 15

- $X_1, \ldots, X_n \sim Exp(\theta)$
- $p(x; \theta) = \theta e^{-\theta x}$     $x > 0$

(*) exponential distri family (not exp. family)
comes up a lot in lifetimes of components etc.

$$L(\theta) = \prod_{i=1}^{n} p(x_i, \theta) = \prod_{i=1}^{n} \theta e^{-\theta x_i} = \theta^n e^{-\theta \sum_{i=1}^{N} x_i} = \theta^n e^{-n\theta \bar{x}_n}$$

Hence $\ell(\theta) = -n\theta \bar{x}_n + n \log \theta$

$$S(\theta) = \ell'(\theta) = \frac{n}{\theta} - n\bar{x}_n$$

- MLE $\hat{\theta} = \frac{1}{\bar{x}_n}$

· $\ell''(\theta) = \frac{-n}{\theta^2}$ so $I_n(\theta) = \mathbb{E}\left[-\ell''(\theta)\right] = \frac{n}{\theta^2}$ (?)     ✓ ⑬ Review/ clarify

Hence :- $\hat{\theta} \approx N\left(\theta, \frac{\theta^2}{n}\right)$     $\hat{se} = \frac{\hat{\theta}}{\sqrt{n}}$

---

## Example 16 · Bernoulli

- $X_1, \ldots, X_n \sim Ber(p)$

MLE: $\hat{p} = \bar{X}$

Fisher info: $I(p) = \frac{1}{p(1-p)}$     so $\sqrt{n}(\hat{p}-p) \xrightarrow{d} N(0, p(1-p))$
(n=1)

HW: This result can be attained via application of CLT; where the above machinery (asympt. normality) gets useful is when the MLE is a complicated non-linear function

⑭ - this deduction (simple manip. of normal)
- asymptotic variance is $\frac{p(1-p)}{n}$

(*) informally; $\hat{p} \approx N\left(p, \frac{p(1-p)}{n}\right)$

- asymptotic variance $\frac{p(1-p)}{n}$, estimated via $\frac{\hat{p}(1-\hat{p})}{n}$

- estimated standard error of MLE : $\hat{se} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

- Suppose we want to estimate

$$\tau = \frac{p}{(1-p)}$$

- we know $\hat{\tau} = \frac{\hat{p}}{1-\hat{p}}$    so    $\frac{\partial}{\partial p}\frac{p}{1-p} = \frac{1}{(1-p)^2}$
  MLE

- The estimated standard error :-

$$\hat{se}(\hat{\tau}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \times \frac{1}{(1-\hat{p})^2} = \sqrt{\frac{\hat{p}}{n(1-\hat{p})^3}}$$

- $\frac{W9\,typo}{p9}$ :-    $\hat{se} = \sqrt{\frac{|\tau'(\hat{\theta})|^2}{I_n(\hat{\theta})}}$

- LW: optimality of MLE (claimed by Fisher)
- There exist better estimators for more complex distns, but even in one para
  settings. For finite dim models with reg. conditions
- LW: Any well behaved estimator $\hat{\theta}$ satisfies

$$\tilde{\theta} = \theta + \frac{1}{n}\sum_{i=1}^{n} \psi(X_i) + o_p(n^{-\frac{1}{2}}) \quad or \quad \sqrt{n}(\hat{\theta}-\theta) \xrightarrow{d} N(0, V(\theta))$$

But also:    $V(\theta) \geq \frac{1}{I(\theta)}$     - this is the
sense in which
MLE is 'optimal'

- LW: Proven by LeCam in 70s/80s
  - see van de Vaart, asymp. stats
       optimal/ec ot
- we really mean MLE is efficient i.e. smallest possible variance
- Asymptotic variance is small compared to other estimators.
- A lot of technical apparatus to describe some of insights on
  regularity, efficiency.
- LW: Asymptotic theory gives use useful info about approx. distn of an
  estimator, and optimality. But also for comparing estimators

# 8. Relative efficiency

(*) Another way to compare well-behaved estimators (i.e. have asymptotic normal distributions).

(*) Done through asymptotic relative efficiency (ARE).

- for estimators $W_n, V_n$, if

$$\sqrt{n}(W_n - \tau(\theta)) \xrightarrow{d} N(0, \sigma_W^2)$$

$$\sqrt{n}(V_n - \tau(\theta)) \xrightarrow{d} N(0, \sigma_V^2)$$

then asymptotic rel. efficiency is:-

$$ARE(V_n, W_n) = \frac{\sigma_W^2}{\sigma_V^2} \qquad \text{(ratio of asymptotic variances)}$$

## Example 17

$X_1, \ldots, X_n \sim \text{Poisson}(\lambda)$

- MLE of $\lambda$ is $\bar{X}$

- Suppose you want to estimate $\tau = e^{-\lambda}$ (e.g. situations where you require $P(X_i = 0)$), then $\hat{\tau} = e^{-\hat{\lambda}}$ (MLE)   (e.g. no. of crashes/week)

- A few methods for estimating $P(X=0)$.

- define:-

$$Y_i = \mathbb{I}(X_i = 0) \qquad W_n = \frac{1}{n}\sum_{i=1}^{n} Y_i \qquad \mathbb{E}[W_n] = \tau$$

(*) How to find limiting distn of $\hat{\tau}$?                     (A$_5$) review.

- from $\lambda$ - variance is Fisic info.
- Take a function of $\lambda$ - use Delta method (test 9.)

(*)  $\sqrt{n}(\hat{\tau} - \tau) \xrightarrow{d} N(0, \lambda e^{-2\lambda})$      (via Delta method)

(*) can also use CLT.

$$\sqrt{n}(W_n - \tau) \xrightarrow{d} N(0, e^{-\lambda}(1 - e^{-\lambda}))$$

(*)

$$ARE = \frac{\lambda}{e^{\lambda}-1} < 1$$

⌣

- as mle has smallest variance

- can trial different values of $\lambda$ to search for other efficient est.

(*) If model is wrong, mle is also not very good
(i.e. not Poisson)

LW: There is a tradeoff between robustness and efficiency; how much
you trust the model; tempers claim of mle optimality.
must include caveat that it is optimal IF the model is 'correct'

## 9. Multivariate case

- $\theta = (\theta_1, ..., \theta_R)$     $|\theta|$ is fixed

- in this case:-

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, I^{-1}(\theta))$$

$I^{-1}(\theta)$ - Fisher inverse
matrix
info.

- Approximate
standard error $\hat{se}(\hat{\theta}_j) = \sqrt{\dfrac{I^{-1}_{j,j}(\hat{\theta})}{n}}$

②

- If $\tau = g(\theta)$ with $g: \mathbb{R}^K \to \mathbb{R}$, then by delta method:-

$$\sqrt{n}(\hat{\tau} - \tau) \xrightarrow{d} N(0, (g')^T I^{-1} g')$$     where $g'$ is gradient of
$g$ eval at $\theta$

LW: Eg. what's MLE of $\dfrac{\mu}{\sigma}$?

- MLE $= \dfrac{\hat{\mu}}{\hat{\sigma}}$

- standard error: $g(\mu, \sigma) = \dfrac{\mu}{\sigma}$

- compute gradient i.e. $\dfrac{\partial g}{\partial \mu}, \dfrac{\partial g}{\partial \sigma}$; multiply by inverse Fisher info;
get limiting distr. of funct. of
params.

LW: Many of parametric models covered seem to be 'nice', as they are in exponential family. Not to be confused with exponential distr. (although this is an exp. family)

## 10. Exp. families

A density of the form:

$$p(x; \theta) = c(\theta) h(x) e^{\theta^T t(x)}$$

- MLE is obtained by solving:-

$$\mathbb{E}_\theta [t(X)] = \frac{1}{n} \sum_{i=1}^{n} t(X_i)$$

- In general $t(x)$ is a vector, with same dim as $\theta$

- Note: $\frac{1}{n} \sum_{i=1}^{n} t(X_i)$ is a minimal sufficient statistic

- Fisher info: $I(\theta) = a''(\theta)$ ; $a(\theta) = -\log c(\theta)$

LW: Favors non-parametric stats
   - Everything said so far depends on whether parametric model is correct

LW: A parametric model is 'never correct'
   - Deal with this by finding estimators that are robust (i.e. allow deviations from correctness of models)
   - (other than nonparametric methods)

LW: Will not go through 11. Robustness

   - LW calculates ARE of (median, mle) = 0.64

- But if data is not Normal, there are outliers; median is better → more robust

(57): The whole subfield of robust statistics/estimators concerns tradeoff between with efficiency and robustness

LW: We deal with this issue of model correctness through non-parametric methods (soon)

point estim ⟶ hypoth test ⟶ confidence int.

- We've seen much on finding good point estimators, next hypothesis testing