

YouTube 19/09/2016

Convergence TheoryPrevious review:

$$X_1, \dots, X_n \sim P$$

$$\text{CLT: } - \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$$

$$\sqrt{n}(\bar{X} - \mu) \xrightarrow{d} N(0, \sigma^2)$$

Informally $\bar{X} \approx N(\mu, \frac{\sigma^2}{n})$ (not formal but prec.)

(*) original distn arbitrary,
but limiting distn is approx normal

Theorem 17 (Multivariate CLT)

$X_1, \dots, X_n \sim P$ IID v. vectors

$$X_i = \begin{pmatrix} X_{i1} \\ \vdots \\ X_{ik} \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_k \end{pmatrix} \quad \Sigma$$

$$X_i \in \mathbb{R}^k, \mu \in \mathbb{R}^k, \Sigma \in \mathbb{R}^{k \times k}$$

(*) Intuit: Take measure multiple characteristics of a person (e.g. height, weight, ...)

(*) Elements of vector are not independent; independence is across vectors (which are observations)
(i.e. charact.) (people in this example)

$$\Sigma = \begin{pmatrix} \text{var}(X_1) & \dots & \text{cov}(X_1, X_k) \\ \text{cov}(X_2, X_1) & & \vdots \\ \text{cov}(X_k, X_1) & & \text{var}(X_k) \end{pmatrix}$$

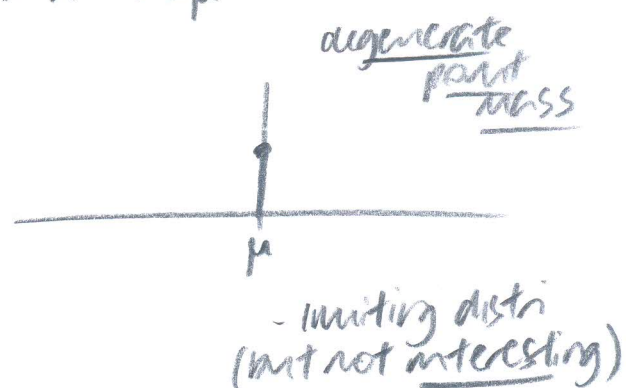
$$\bar{X} = \begin{pmatrix} \bar{X}_1 \\ \vdots \\ \bar{X}_k \end{pmatrix} \quad \bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{ji}$$

then $\sqrt{n}(\bar{X} - \mu) \xrightarrow{d} N(0, \Sigma)$

(*) Remember you want lim distn of rescaled

CLN (via convergence taxonomy)

$$\bar{X} \xrightarrow{P} \mu \Rightarrow \bar{X} \xrightarrow{d} \delta_\mu$$



5. Delta method (propagation of errors)

Q: What if I want distn of $g(\bar{X})$ for some smooth function $g(\cdot)$?

- maybe \rightarrow continuous mapping theorem?

- will not help

- we want a rescaled version of continuous mapping

(scalars)

- $g(\bar{X})$

- consider: - $\frac{\sqrt{n}(g(\bar{X}) - g(\mu))}{\sigma |g'(\mu)|} \xrightarrow{d} N(0, 1)$

- assume $g(\cdot)$ is differentiable; $g'(\mu) \neq 0$

- informally: $g(\bar{X}) \hat{=} N(g(\mu), (g'(\mu))^2 \frac{\sigma^2}{n})$ $\bar{X}_n \hat{=} N(\mu, \frac{\sigma^2}{n})$

W: When you take a function of \bar{X} , we have

- allows any limit distn of any smooth function of av.

i) mean becomes $g(\mu)$

ii) variance multiplies by derivative

example 19

- X_1, \dots, X_n IID with finite mean μ , finite variance σ^2

- we know $\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$

- consider

$$W_n = e^{\bar{X}_n}$$

(non linear fn)

$$g(s) = e^s \quad g'(s) = e^s$$

(W)(A1): check you can square the rest. intuitively

W: immediately know limiting distn. :-

$$\sqrt{n}(W_n - e^\mu) \xrightarrow{d} N(0, \sigma^2 (e^\mu)^2) \text{ or}$$

Heuristic derivation

- consider $g(\bar{X})$
 - use Taylor series approx, expand $g(\bar{X})$ about μ
(not including all terms \rightarrow heuristic)
- $$g(\bar{X}) \approx g(\mu) + (\bar{X} - \mu)g'(\mu)$$

$$\sqrt{n}(g(\bar{X}) - g(\mu)) \approx \sqrt{n}(\bar{X} - \mu)g'(\mu)$$

$$\sqrt{n}(\bar{X} - \mu) \xrightarrow{d} N(0, \sigma^2) \text{ via CLT}$$

- $g'(\mu)$ is constant; so multiply by constant $(g'(\mu))^2$ in variance

$$\text{Var}(aX) = a^2 \text{Var}(X)$$

Theorem 20 (Multivariate Delta Method)

- $\sqrt{n}(\bar{X} - \mu) \xrightarrow{d} N(0, \Sigma)$ (CLT) $\bar{X} \in \mathbb{R}^K$
- suppose $g: \mathbb{R}^K \rightarrow \mathbb{R}$ (scalar fn) (also one for vector fn)
- take gradient of $g(y)$ wrt y :

$$\nabla_y g(y) = \begin{pmatrix} \frac{\partial g}{\partial y_1} \\ \vdots \\ \frac{\partial g}{\partial y_K} \end{pmatrix}$$

- let $\underline{\nabla}_\mu$ denote $\nabla_y g(y)|_{y=\mu}$ (i.e. gradient of g wrt y evaluated at μ)

- Then:

$$\sqrt{n}(g(\bar{X}) - g(\mu)) \xrightarrow{d} N(0, \underline{\nabla}_\mu^T \Sigma \underline{\nabla}_\mu)$$

- example 21 $\bar{X}_i = \begin{pmatrix} X_{1i} \\ X_{2i} \end{pmatrix}$; have $X_1, X_2, X_3, \dots, X_n$

- i.e. n IID random vectors with mean (μ_1, μ_2) and variance Σ

- denote:

$$\bar{X}_1 = \frac{1}{n} \sum_{i=1}^n X_{1i} \quad \bar{X}_2 = \frac{1}{n} \sum_{i=1}^n X_{2i}$$

- then $\sqrt{n} \begin{pmatrix} \bar{X}_1 - \mu_1 \\ \bar{X}_2 - \mu_2 \end{pmatrix} \xrightarrow{d} N(0, \Sigma) \quad \sqrt{n}(\bar{X} - \mu) \xrightarrow{d} N(0, \Sigma)$

- define a function (product) $g(s_1, s_2) = s_1 s_2$

Q How do I get the limiting distri of \bar{X}_1, \bar{X}_2 ? (complic; no longer sum)

$$\nabla_{\mathbf{s}} g(\mathbf{s}) = \begin{pmatrix} s_2 \\ s_1 \end{pmatrix}$$

$$\nabla_{\mu}^T \Sigma \nabla_{\mu} = (\mu_2 \ \mu_1) \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{pmatrix} \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \mu_2^2 \sigma_{11} + 2\mu_1 \mu_2 \sigma_{12} + \mu_1^2 \sigma_{22}$$

$$\therefore \sqrt{n}(\bar{X}_1 \bar{X}_2 - \mu_1 \mu_2) \xrightarrow{d} N(0, \mu_2^2 \sigma_{11} + 2\mu_1 \mu_2 \sigma_{12} + \mu_1^2 \sigma_{22})$$

informally $g(\bar{X}) \approx N(g(\mu), \frac{1}{n} \nabla_{\mu}^T \Sigma \nabla_{\mu})$

we have probabilistic tools to do statistics/ML

- use tools to statistical inference (part II) - extracting insight from data

we: in real-world, stats is opposite to prob, in sense that:-

Probability: we have a distri

stats: - given data, figure out what generated data

YouTube 19/09/16

Statistical Inference (Lecture Notes 5)

- W: I give you data, you have to do things e.g. tell me about DGP.

(*) Difference between statistics and probability.

- Probability \rightarrow start with a distri P

- Statistics \rightarrow start with a model (statistical)

(*) - i.e. what are the possible distributions (aka distri family)
that would have generated the data

1. Statistical Models

- A statistical model \mathcal{P} is a set of probability distributions (or densities)

- W: All we know is that the 'true' distribution is within that set of distributions

- W: Introduce artificial distinction between nonparametric and parametric models

Parametric model

W: Set of distri that can be indexed by a finite no. params

$$\mathcal{P} = \left\{ p(x; \theta) : \theta \in \Theta \right\} \quad \Theta \subset \mathbb{R}^d \quad \Theta - \text{param space}$$

- set of normal densities: $\left\{ p(x; \theta) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \right\}$

$$\theta = (\mu, \sigma^2) \quad \Theta = \mathbb{R} \times \mathbb{R}^+$$

(param space)

W: Highly restrictive; in general shouldn't use

W: learn off parametric models later (but only starting point)

Nonparametric models

- A model (collection of prob. distri) that cannot be specified by a finite dimensional param.

W: should be infinitely-parametric

2 senses of non-parametric

example $\mathcal{P} = \{\text{all distri on } \mathbb{R}^d\}$

→ non-param

example 2 $\mathcal{P} = \left\{ p: \int (p''(x))^2 dx < \infty \right\}$ "set of all distri with densities (which have 2nd derivatives and are square integrable (?))
- smoothness constraint on \mathcal{P}

(*) start with parametric models

W: given data from some distri; and that (assumed) distri comes from some distri within parametric model

- Have to find parameters → can we estimate them from a finite sample?

W: data reduction, sufficiency

- Ronald Fisher

- no computers then

- Q: can I reduce data to a few nos → how much info loss (lots of emphasis on sufficiency in pre-computer age)

- sufficiency became less important with advent of computers

- But BIG DATA → sufficiency

Q: what is sufficiency (in context of a statistic)

W: A statistic is a function of the data (itself a random variable as a function of random variables)

(*) $\bar{X} - \mu$ is not a statistic; not a function of the data, but param μ

(*) key distinction between i.v.; parameters \rightarrow these index the distri

(*) Assume parameters unknown.

Q: Which statistics capture all info in the data \rightarrow sufficient statistic

(*) formal definition; intuition

- sufficient statistics

- "parametric model"
; θ - means family of densities (freq.)
indexed by the parameters (fixed)

- Suppose $X_1, \dots, X_n \sim p(x; \theta) = p_\theta(x)$

- $T(X_1, \dots, X_n)$ is sufficient for θ if the condit. distri: - (*) has realisations t
(*) condition on $T=t$

$$p(x_1, \dots, x_n | T=t, \theta) = p(x_1, \dots, x_n | T=t)$$

(*) intuition: - distri depends on θ

e.g. Normal: - $p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2}}$
 $x \sim N(0,1)$

$$p(x_1, x_2, \dots, x_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_i-\mu)^2}{2}}$$

- suppose I give you a statistic; and get you to estimate μ

- \therefore ms left

(*) amounts to condit. distri; conditioning on statistic $T=t$; then
the cond. distri no longer depends on parameters

(*) density no longer contains parameters conditioning on $T=t$

(*) possible to take statistic as the data: -

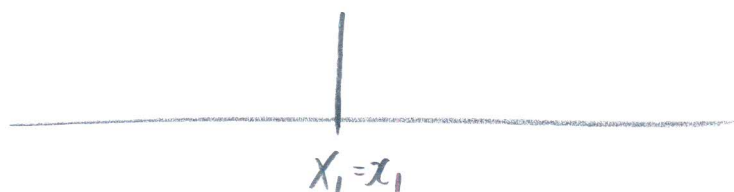
$$T(X_1, \dots, X_n) = (X_1, \dots, X_n)$$

- what is distri of single

i.v. given I tell you what it is?
(one obs.)

- point mass at x

(no free param.)



(*) The statistic equal to whole dataset is always sufficient
(data are sufficient)

(*) Tautological; but tells us that there are good and bad statistics

(*) We want to build to minimal sufficient statistics (*)

- \bar{x}_n better at data reduction than (x_1, \dots, x_n)

(*) Not enough to be sufficient. (*) sufficiency is only
really useful for parametric models

(*) next lecture