

MEM-262 Παραμετρική Στατιστική

Τμήμα Μαθηματικών και Εφ. Μαθηματικών, Πανεπιστήμιο Κρήτη

Διδάσκων : Κώστας Σμαραγδάκης

Διάλεξη 5 : 16-10-2020

Θεώρημα Mann-Wald

Έστω ακολουθία τυχαίων μεταβλητών X_1, X_2, \dots στον ίδιο χώρο πιθανότητας και g συνεχής συνάρτηση τότε:

$$X_n \xrightarrow{a.s.} X \Rightarrow g(X_n) \xrightarrow{a.s.} g(X)$$

$$X_n \xrightarrow{p} X \Rightarrow g(X_n) \xrightarrow{p} g(X)$$

$$X_n \xrightarrow{d} X \Rightarrow g(X_n) \xrightarrow{d} g(X)$$

Σύγκλιση της δειγματικής διασποράς

$$\sigma^2 = \mathbb{E}[(X - \mathbb{E}(X))^2]$$

$$X_1, X_2, \dots \text{ iid} \quad \mathbb{E} \sum X_k \zeta = \mu$$

Θα δείξουμε ότι

$$S_n^2 = \frac{\sum_{k=1}^n (X_k - \bar{X}_n)^2}{n-1} \xrightarrow{P} \sigma^2$$

$$\text{Var} \sum X_k \zeta = \sigma^2$$

$$\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k \quad \text{Διαστηματικός τύπος}$$

Δειγματική διασπορά

$$S_n \xrightarrow{P} \sigma$$

$$S_n^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2$$

$$S_n^2 = \frac{n}{n-1} \left[\frac{1}{n} \sum_{k=1}^n (X_k - \bar{X}_n)^2 \right]$$

$$\frac{1}{n} \sum_{k=1}^n X_k = \bar{X}_n \quad - \frac{1}{n} \mu$$

$$\frac{1}{n} \sum_{k=1}^n (X_k - \bar{X}_n)^2 = \frac{1}{n} \sum_{k=1}^n [(X_k - \mu) - (\bar{X}_n - \mu)]^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \mu)^2 - \underbrace{2(\bar{X}_n - \mu) \frac{1}{n} \sum_{k=1}^n (X_k - \mu)}_{2(\bar{X}_n - \mu)^2} + \frac{1}{n} \sum_{k=1}^n (\bar{X}_n - \mu)^2 =$$

$$= \frac{1}{n} \sum_{k=1}^n (X_k - \mu)^2 - (\bar{X}_n - \mu)^2$$

NMA: $\bar{X}_n \xrightarrow{P} \mu$ επιλέγω $g(x) = (x - \mu)^2$

$$g(\bar{X}_n) = (\bar{X}_n - \mu)^2 \quad g(\mu) = 0$$

$$* \underbrace{(X_1 - \mu)^2}_{Y_1}, \underbrace{(X_2 - \mu)^2}_{Y_2}, \dots, \underbrace{(X_n - \mu)^2}_{Y_n} \text{ iid}$$

NMA: $\frac{1}{n} \sum_{k=1}^n Y_k \xrightarrow{P} \mathbb{E}(Y_1) = \mathbb{E}[(X_1 - \mu)^2] = \text{Var}(X_1) = \sigma^2$

$$\frac{1}{n} \sum_{k=1}^n (X_k - \mu)^2 \xrightarrow{P} \sigma^2$$

$$Z_n = \frac{n}{n-1} \xrightarrow{P} 1 \quad \mathbb{P}(|Z_n - 1| > \varepsilon) \xrightarrow{n \rightarrow \infty} 0$$

Σύγκλιση της δειγματικής διασποράς

$$\left. \begin{array}{l} X_n \xrightarrow{d} X \\ Y_n \xrightarrow{p} c \end{array} \right\} X_n Y_n \xrightarrow{d} cX$$

$$Z_n = \frac{1}{n} \sum_{k=1}^n (X_k - t)^2 \xrightarrow{d} \sigma^2$$

$$S_n^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x}_n)^2 \quad \text{Var}(X) = \mathbb{E}((X - \mathbb{E}(X))^2)$$

$$S_n^2 \xrightarrow{p} \sigma^2$$

$$\text{δηλαδή } \mathbb{E}(S_n^2) = \sigma^2$$

Αν $X_n \xrightarrow{d} c$ σταθερά τότε $X_n \xrightarrow{p} c$

$$\begin{aligned} \mathbb{E}\left(\sum_k (X_k - \bar{X}_n)^2\right) &= \mathbb{E}\left(\sum_k X_k^2 - 2\bar{X}_n \sum_{k=1}^n X_k + n\bar{X}_n^2\right) = \sum_{k=1}^n \mathbb{E}(X_k^2) - \mathbb{E}(n\bar{X}_n^2) \\ &= \underbrace{\sum_{k=1}^n \mathbb{E}(X_k^2)}_{n\sigma^2} - \underbrace{n \mathbb{E}(\bar{X}_n^2)}_{n(\frac{\sigma^2}{n} + t^2)} = n\sigma^2 - \sigma^2 = \underline{(n-1)\sigma^2} \end{aligned}$$

$$\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2 \Rightarrow \mathbb{E}(X^2) = \text{Var}(X) + (\mathbb{E}(X))^2 = \sigma^2 + t^2$$

$$\mathbb{E}(\bar{X}_n^2) = \text{Var}(\bar{X}_n) + (\mathbb{E}(\bar{X}_n))^2$$

$$\mathbb{E}(S_n^2) = \sigma^2.$$

Κεντρικό οριακό θεώρημα

Κεντρικό οριακό θεώρημα

Έστω X_1, X_2, \dots i.i.d και $\mu = \mathbb{E}(X_1)$, $\sigma^2 = \text{Var}(X_1) < \infty$

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} Z \sim \mathcal{N}(0, 1)$$

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

$$(\bar{X}_n - \mu) \xrightarrow{d} \left(\frac{\sigma}{\sqrt{n}} Z \right) \sim \mathcal{N}\left(0, \frac{\sigma^2}{n}\right)$$

$$\bar{X}_n \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

όπου Z τυχαία μεταβλητή η οποία ακολουθεί την τυπική κανονική κατανομή $\mathcal{N}(0, 1)$.

Αποτέλεσμα για άγνωστη διασπορά σ^2

$$g(x) = \frac{\sigma}{x}$$

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} \xrightarrow{d} Z \sim \mathcal{N}(0, 1) \quad ?$$

$$S_n \xrightarrow{p} \sigma \quad \frac{\sigma}{S_n} \xrightarrow{p} 1$$

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} = \underbrace{\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma}}_Z \cdot \underbrace{\left(\frac{\sigma}{S_n}\right)}_{\xrightarrow{p} 1} \xrightarrow{d} Z \sim \mathcal{N}(0, 1)$$

- ▶ Στη στατιστική συμπερασματολογία παρατηρήσιμα δεδομένα μοντελοποιούνται ως πραγματοποιήσεις τυχαίων μεταβλητών.
- ▶ Αξιοποιώντας αυτά τα δεδομένα μελετούμε το μηχανισμό (στατιστική κατανομή) για τη δημιουργία τους.

Βασικό κομμάτι λοιπόν της στατιστικής συμπερασματολογίας αποτελεί η επεξεργασία και η ανάλυση παρατηρήσεων.

Τα παρατηρήσιμα δεδομένα θα αναπαριστώνται ως \mathbf{x} με

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T$$

- ▶ Κάθε x_i περιγράφει μια παρατήρηση που όπως είπαμε θεωρείται πραγματοποίηση μιας τυχαίας μεταβλητής την οποία θα συμβολίζουμε ως X_i .
- ▶ Τα x_i μπορούν να είναι πέρα από αριθμούς και διανύσματα. Σε αυτή τη περίπτωση το \mathbf{x} θα είναι πίνακας.

Σύμφωνα με τα παραπάνω μπορούμε να δούμε το \mathbf{x} σαν μια πραγματοποίηση μιας διανυσματικής τυχαίας μεταβλητής \mathbf{X} όπου:

$$\mathbf{X} = (X_1, X_2, \dots, X_n)^T$$

Η \mathbf{X} εν γένει θα λαμβάνει τιμές στο \mathbb{R}^n (ή θα περιγράφεται από πραγματικούς πίνακες εάν x_i είναι διανύσματα). **Γενικά τον χώρο στον οποίο λαμβάνει τιμές η \mathbf{X} τον συμβολίζουμε με \mathcal{X} και θα τον καλούμε δειγματικό χώρο της \mathbf{X} .**

Παράδειγμα (ref: Liero-Zwanzig)

Ένας κατασκευαστής στυλό πραγματοποιεί τυχαία δειγματοληψία 400 μονάδων ανά ημέρα από τα 40000 στυλό της ημερήσιας παραγωγής. Οι παρατηρήσεις είναι ο αριθμός των ελαττωματικών προϊόντων ανά ημέρα. Για παράδειγμα:

Mon	Tue	Wed	Thu	Fri	Sat	Sun
8	5	9	4	6	8	10
$0-400$	$0-400$	$0-400$	- - - -	- - - -	- - - -	- - - -

Θα περιγράψουμε τα δεδομένα, την τυχαία μεταβλητή που θεωρείται υπεύθυνη για την παραγωγή των δεδομένων και τον δειγματικό της χώρο.

$$\mathbf{x} = (8, 5, 9, 4, 6, 8, 10)^T \quad \mathbf{X} = (X_1, X_2, X_3, X_4, X_5, X_6, X_7)$$

$$\mathcal{X} = \{0, 1, 2, \dots, 400\}^7$$

Παράδειγμα (ref: Liero-Zwanzig)

Ο παρακάτω πίνακας παρουσιάζει την ημερήσια παραγωγή γαλακτός (kg/day) και τη αντίστοιχη παραγωγή πρωτεΐνη γαλακτός (kg/day) για 5 (Friesian) αγελάδες.

Milk production	Milk protein
42.7	1.20
40.2	1.16
38.2	1.07
37.6	1.13
32.2	0.96

$$X = (X_1, \dots, X_5)$$

$$x_1 = (42.7, 1.2)$$

$$x_2 = (40.2, 1.16)$$

⋮

$$X = \mathbb{R}_+^{2,5}$$

$$\begin{bmatrix} 42.7 & 1.2 \\ 40.2 & 1.16 \\ \vdots & \vdots \\ 32.2 & 0.96 \end{bmatrix}$$

Θα περιγράψουμε τα δεδομένα, την τυχαία μεταβλητή που θεωρείται υπεύθυνη για την παραγωγή των δεδομένων και τον δειγματικό της χώρο.

Παραμετρικό στατιστικό μοντέλο

- ▶ Δειγματικός χώρος \mathcal{X} της τυχαίας μεταβλητής \mathbf{X} που περιγράφει τις παρατηρήσεις.
- ▶ Παραμετρικός χώρος Θ επί τη βάση του οποίου ορίζεται η οικογένεια των πιθανών μέτρων πιθανότητας $\mathcal{P} = \{P_\theta : \theta \in \Theta \subseteq \mathbb{R}^d\}$

i.i.d Συνήθως θεωρούμε επαναλαμβανόμενα πειράματα που διεξάγονται υπό τις ίδιες συνθήκες και στα οποία η κάθε πραγματοποίηση δεν επηρεάζει τις υπόλοιπες.

- ▶ Τυχαίο δείγμα πειράματος:

X_1, X_2, \dots, X_n i.i.d

- ▶ Μέτρο πιθανότητας:

$$P_\theta = P_{1,\theta} \otimes P_{1,\theta} \otimes \dots \otimes P_{1,\theta}$$

- ▶ Συνάρτηση μάζας (διακριτές τμ) ή συνάρτηση πυκνότητας (συνεχείς τμ) πιθανότητας

$$p(\mathbf{x}; \theta) \quad \text{ή} \quad f(\mathbf{x}; \theta)$$

$$\mathcal{N}(\mu, \sigma^2) \quad \theta = (\mu, \sigma^2) \quad f(x; \theta) \quad \text{ή} \quad f(x; \mu, \sigma^2)$$

Θα ασχοληθούμε με παραμετρικά μοντέλα για τα οποία ισχύει:

Για $\theta_1, \theta_2 \in \Theta$ με $\theta_1 \neq \theta_2$ οι κατανομές P_{θ_1} και P_{θ_2} είναι διαφορετικές.

Έστω ένα ενδεχόμενο A του \mathcal{X} X

- ▶ Για \mathbf{X} διακριτή τυχαία μεταβλητή έχουμε:

$$P_{\theta}(A) = \sum_{\mathbf{x} \in A} \cancel{P_{\theta}(\mathbf{x})} p(\mathbf{x}; \theta)$$

- ▶ Για \mathbf{X} συνεχή τυχαία μεταβλητή έχουμε:

$$P_{\theta}(A) = \int_A f(\mathbf{x}; \theta) dx$$

Τρόποι στατιστικής συμπερασματολογίας

$$\begin{array}{l} \mathbf{x} = (\dots) \\ \mathbf{x} = \left[\begin{array}{c} \\ \end{array} \right] \end{array} \longrightarrow \begin{array}{l} N(\mu, \sigma^2) \\ \hat{\mu} \quad \hat{\sigma}^2 \\ \mu \in [1, 2] \\ \sigma^2 \in [10, 15] \end{array} \quad \begin{array}{l} W = [1, 2] \times [10, 15] \\ \mathcal{P}(\theta \in W^c) < \alpha. \end{array}$$

1. Σημειακή εκτίμηση : $\mathbf{x} \rightarrow \hat{\theta}(\mathbf{x}) \approx \theta$ $\mathcal{P} = 0.05$
2. Περιοχές εμπιστοσύνης : $\mathbf{x} \rightarrow W \subset \Theta$ τω $P(\theta \notin W) < \rho$ για ρ κοντά στο 0.
3. Έλεγχοι υποθέσεων : Λέμε ότι ισχύει μια υπόθεση H_0 (μηδενική) για το θ εκτός και εάν μπορεί να καταρριφθεί από το \mathbf{x} για χάρη της εναλλακτικής υπόθεσης H_1 .

$$H_0 : \mu = 2$$

$$H_1 : \mu \neq 2$$