

μη-παράμετρη Στατιστική (Φροντ.#1)-1-

Εισαγωγή στους Μη-Παράμετρικούς Ελέγχους Υποθέσεων

Αν και σε πολλές περιπτώσεις οι παράμετρες υποθέσεις μοιάζουν λογικές συχνά δεν έχουμε προηγούμενη γνώση των κατανομών που περιγράφουν τα δεδομένα μας. Σε τέτοιες περιπτώσεις η χρήση της Παράμετρικής Στατιστικής μπορεί να μας οδηγήσει σε λανθασμένα αποτελέσματα και συνεπώς υπερπέρατα.

Σε περιπτώσεις όπου δεν είναι επιτρεπτή οποιαδήποτε υπόθεση για τη μορφή του πληθυσμού ο ερευνητής χρησιμοποιεί το Κ.Ο.Θ. και στηρίζεται στην υπόθεση της κανονικότητας, ώστε να δεξιάξει τους ελέγχους που τον ενδιαφέρουν.

Τι γίνεται όμως όταν τα δείγματα που έχουμε συλλέξει είναι μικρά και τα δεδομένα δεν κατανέμονται σύμφωνα με την κανονική κατανομή;

Η βασική ιδέα της μη-Παράμετρικής Στατιστικής **Συμπερασματολογίας** είναι να επάξει υπερπέρατα για άγνωστες ποσοτικές ενδιαφέροντος, χρησιμοποιώντας δεδομένα αλλά παράλληλα θεωρώντας το μεγαλύτερο δυνατό σύνολο υποθέσεων.

Προβλήματα υπάρχουν όταν τα δείγματα που έχουμε συλλέξει από τον υπό-μέγιστο πληθυσμό είναι μικρά και τα δεδομένα δεν κατανέμονται σύμφωνα με την κανονική κατανομή. Το κυριότερο πλεονέκτημα των μη-παράμετρικών τεχνικών είναι ότι οι έλεγχοι

δεν προϋποθέτουν γνώση της κατανομής του πληθυσμού από τον οποίον προέρχεται τα υπό μελέτη στοιχεία.

Επιπλέον,

- απαιτούν λίγες υποθέσεις για τους πληθυσμούς από τους οποίους προέρχονται τα δεδομένα,
- επιτρέπουν στον ερευνητή να υπολογίσει, σε κάποιες περιπτώσεις ελέγχων την ακριβή τιμή (exact value) του παρατηρούμενου επιπέδου σημαντικότητας $\hat{\alpha}$,
- είναι (σχεδόν!) πάντα πιο απλές τεχνικές προς εφαρμογή συγκριτικά με τις αντίστοιχες παραμετρικές τεχνικές,
- συχνά είναι πιο εύκολες στην κατανόηση,
- δεν είναι τόσο ευαίσθητες στις ακραίες τιμές,
- εφαρμόζονται (συχνά) στις τάξεις μελών των δεδομένων και όχι στα ίδια τα δεδομένα,
- μπορούν να χρησιμοποιηθούν και για δεδομένα που είναι ταξινομημένα σε κατηγορίες (κατηγορικά δεδομένα σε κλίμακα διάταξης ή σε ονομαστική κλίμακα).

Όμως:

- σε περίπτωση όπου είναι γνωστό ότι τα δεδομένα προέρχονται από κανονική κατανομή τότε οι παραμετρικοί έλεγχοι είναι πιο ισχυροί.

Αλλά:

- αν η κατανομή των δεδομένων δεν είναι κανονική τότε οι μη-παραμετρικοί έλεγχοι έχουν

M_n - παραμετρικοί έλεγχοι για ένα δείγμα

Στο γενικό πρόβλημα για ένα δείγμα, τα διαθέσιμα δεδομένα συνιστούν ένα απλό σύνολο δεδομένων, οπότε ένα τυχαίο δείγμα από μια συνάρτηση κατανομής F_X .

Η υπερπαραστοχία αφορά κάποιο χαρακτηριστικό μέγεθος που σχετίζεται με τη συνάρτηση κατανομής F_X . Στο κλασικό πρόβλημα υπερπαραστοχίας για ένα δείγμα, τα δεδομένα χρησιμοποιούνται για να εξαχθούν

συμπέρασμα για κάποιο συγκεκριμένο χαρακτηριστικό της κατανομής του πληθυσμού. Οι μη-παραμετρικές τεχνικές είναι χρήσιμες, για παράδειγμα, στην εξαγωγή υπερπαραστοχίας ενός μέτρου θέσης για την κατανομή που μας ενδιαφέρει.

Οι υποθέσεις των ελέγχων σχετίζονται με τη διάρκεια που όπως και η μέση τιμή είναι βασικοί δείκτες της κεντρικής τάσης ενός πληθυσμού. Σε κάθε πληθυσμό η διάρκεια πάντοτε υπάρχει και είναι πιο ευσταθής επιμηκνής από τη μέση τιμή σε σχέση με τη μέση τιμή.

Έλεγχος προσήπου - Προσημικός Έλεγχος (Sign Test)

Έστω ένα τ.δ. N παρατηρήσεων X_1, X_2, \dots, X_N από ένα πληθυσμό με κατανομή F_X και άγνωστη διάρκεια

M όπου F_X θεωρείται συνεχής, τουλάχιστον στην "περιοχή" (στη γειτονιά) της διάρκειας M .

-4-

Οι N παρατηρήσεις είναι ανεξάρτητες και ταυτωτικά κατανοημένες και η διάμεσος $M = F_X^{-1}(0.5)$ μοναδική.

Η υπόθεση που μας ενδιαφέρει να ελέγχουμε αφορά την πιθανομετρική διάμεσο:

$$H_0: M = M_0$$

όπου M_0 είναι μία συγκεκριμένη τιμή έναντι της εναλλακτικής υπόθεσης είτε μονόπλευρου είτε αμφίπλευρου ελέγχου. Ισοδύναμα, μπορούμε να γράψουμε:

$$H_0: \theta = P(X > M_0) = P(X < M_0) = 0.5$$

Έστω $K := \#$ παρατηρήσεων $> M_0$. Επειδή οι παρατηρήσεις του δείγματος έχουν "δικοτομήθει" (ως προς την τιμή M_0) σιωπούν ένα σύνολο n ανεξάρτητων τ.μ. από έναν πιθανομετρικό Βερνουλλί με παράμετρο

$$\theta = P(X > M_0). \text{ Άρα, η δειγματική κατανομή της}$$

τ.μ. K είναι διωνυμική με παραμέτρους N και θ .

Όταν ισχύει η H_0 (μάτω από την H_0) τότε $\theta = 0.5$.

Επειδή ισοδύναμα μπορούμε να γράψουμε ότι:

$$K := \# \text{ " + " ανάμεσα στις } N \text{ διαφορές } X_i - M_0 \\ i = 1, 2, \dots, N$$

αυτό το μη-παραμετρικό test που βασίζεται στην ποσότητα K καλείται προσημικό test (sign test).

Συνεπώς, ελέγχουμε $H_0: M = M_0$

έναντι $H_1: M \neq M_0$ ή $H_1: M > M_0$ ή $H_1: M < M_0$

Για την εναλλακτική,

$$H_1: \mu > \mu_0 \text{ ή } \theta = P(X > \mu_0) > P(X < \mu_0)$$

περιοχή απόρριψης R (ή κρίσιμη περιοχή) της H_0
προσδιορίζεται από την ανισότητα: $K \geq K_\alpha$
όπου K_α είναι ο ελάχιστος αριθμός για τον οποίον:

$$P(K \geq K_\alpha | H_0) = \sum_{i=K_\alpha}^N \binom{N}{i} (0.5)^N \leq \alpha$$

Παρόμοια, για την εναλλακτική: $H_1: \mu < \mu_0$

$$\text{ή } \theta = P(X > \mu_0) < P(X < \mu_0)$$

η κρίσιμη περιοχή επιπέδου α
από την ανισότητα: $K \leq K'_\alpha$ όπου K'_α είναι ο μέγιστος
αριθμός για τον οποίον:

$$\sum_{i=0}^{K'_\alpha} \binom{N}{i} (0.5)^N \leq \alpha$$

Για την εναλλακτική, $H_1: \mu \neq \mu_0$ ή $\theta = P(X > \mu_0) \neq P(X < \mu_0)$

η κρίσιμη περιοχή επιπέδου α προσδιορίζεται από
την ανισότητα $K \geq K_{\alpha/2}$ ή $K \leq K'_{\alpha/2}$ όπου $K_{\alpha/2}$ και $K'_{\alpha/2}$
είναι αντίστοιχα, ο ελάχιστος και ο μέγιστος αριθμός
τέτοιοι ώστε:

$$\sum_{i=K_{\alpha/2}}^N \binom{N}{i} (0.5)^N \leq \frac{\alpha}{2} \text{ και } \sum_{i=0}^{K'_{\alpha/2}} \binom{N}{i} (0.5)^N \leq \frac{\alpha}{2}$$

Προφανώς, $k_{\alpha/2} = N - k'_{\alpha/2}$.

Υπολογισμός της p -τιμής: Έστω ότι $H_1: \mu > \mu_0$.

Έστω k_0 ότι είναι η παρατηρούμενη τιμή του προσημ-
κούς ελέγχου. Τότε η p -τιμή δίνεται από την άνω ουρά
της διωνυμικής πιθανότητας:

$$\sum_{i=k_0}^N \binom{N}{i} (0.5)^N$$

Ανάλογα, προκύπτουν οι p -τιμές και για τις υπόλοιπες
μορφές των εναλλακτικών υποθέσεων.

Προσέγγιση με την Κανονική Κατανομή

Για οποιοδήποτε μέγεθος δείγματος N , p μπορούμε να
εφαρμόσουμε τον προσημικό έλεγχο χρησιμοποιώντας
την διωνυμική κατανομή που είναι η ακριβής κατανομή
για τη διενέργεια του ελέγχου. Σημειώνουμε ότι υπάρχουν
κάτάλληλοι πίνακες που για διάφορες τιμές των N και
 α μπορούν να μας δώσουν τις κριτικές περιοχές. Από
τους πίνακες αυτές μπορούμε επίσης να προσδιορίσουμε
και τις p -τιμές.

Οποιοδήποτε, γνωρίζουμε ότι η κανονική προσέγγιση
στη διωνυμική είναι "καλή" όταν $\theta = 0.5$. Έτσι, για
μέτριες τιμές του N (τουλάχιστον 12) η κανονική
προσέγγιση στη διωνυμική μπορεί να χρησιμοποιηθεί
για τον καθορισμό των κριτικών περιοχών. Επειδή
έχουμε μία συνεχής προσέγγιση σε μία διακριτή κατα-

υορή μπορούμε να συρπεριλάβουμε στους υπολογισμούς -7- και τη διόρθωση συνέχειας.

Όπως είδαμε η σ.δ. ελέγχου δίνεται από τον στωχαστικό αριθμό των X_1, X_2, \dots, X_N που είναι μεγαλύτεροι από την τιμή M_0 δηλ. $S'_N = \sum_{i=1}^N I(X_i > M_0)$ όπου

I είναι η δείκτη συνάρτηση. Έτσι, μερικές φορές της S'_N οδηγούν σε απόρριψη της H_0 .

$$S'_N \stackrel{H_0}{\sim} \text{Bin}(N, 1/2) \text{ δηλ. } P(S'_N = k) = \binom{N}{k} \left(\frac{1}{2}\right)^N$$

$$E_{H_0}(S'_N) = \sum_{i=1}^N \left(\frac{1}{2}\right) = \frac{N}{2}, \text{ Var}_{H_0}(S'_N) = \sum_{i=1}^N \left(\frac{1}{4}\right) = \frac{N}{4}$$

Η ασυμπτωτική κανονικότητα της τυποποιημένης μορφής:

$$S_N^* = \frac{S'_N - E_{H_0}(S'_N)}{\sqrt{\text{Var}_{H_0}(S'_N)}} = \frac{S'_N - \frac{N}{2}}{\left(\frac{N}{4}\right)^{1/2}} \text{ προκύπτει}$$

από το Κ.Θ.

Για παράδειγμα, για την εναλλακτική $H_1: \mu > M_0$ η H_0 απορρίπτεται αν $k \geq k_\alpha$ όπου

$$k_\alpha = \frac{N}{2} + z_\alpha \frac{\sqrt{N}}{2} \text{ όπου } z_\alpha \text{ είναι το } \alpha\text{-ποσοστιαίο}$$

σημείο της $N(0,1)$.

Επίσης, η p -τιμή δίνεται από την έκφραση

$$1 - \Phi\left(\frac{k_0 - \frac{N}{2}}{\sqrt{0.25N}}\right), \text{ όπου } \Phi \text{ η σ.κ. της } N(0,1).$$

Σχόλια: Οι υποθέσεις που απαιτούνται για τη διενέργεια του sign test είναι οι ελάχιστες δυνατές. -8-

- ① Ανεξαρτησία παρατηρήσεων
- ② Πληθυσμός που έχει μία συνεχή σ.κ. F_X παντού.
- ③ Πειραματικά υπάρχει ένα πρόβλημα για τις διαφορές ενδύς όπως $X_i - M_0 = 0$. Θα δοθεί πώς αντιμετωπίζεται.
- ④ Είναι εφαρμόσιμο όταν ποσοτικά δεδομένα δεν είναι εφικτό να βρεθούν
- ⑤ Μπορεί να χρησιμοποιηθεί σε διχοτομικά δεδομένα με παρατηρήσεις της μορφής ΝΑΙ-ΟΧΙ (yes-no)
- ⑥ Είναι εύκολο στη χρήση και η μαθηματική προσέγγιση είναι αρκετά ακριβής ακόμα και για μέτρες υπέρ του N .

Αντιμετώπιση του ③. Θεωρητικά, οι διαφορές αυτές δεν λειτουργούν πρόβλημα διότι υποθέτουμε συνέχεια της F_X στην γειτονιά της διαφέρει. Μπορούμε να έχουμε μηδενικές διαφορές είτε λόγω μη-αρκετών μετρήσεων είτε λόγω λανθασμένης υποτιθέμενης συνέχειας της F_X .

Η πιο κοινή αντιμετώπιση των μηδενικών διαφορών είναι να **αγνοήσουμε** αυτές τις διαφορές και να μετράσουμε αναλόγως το N . Εναλλακτικά, μπορούμε να θεωρήσουμε:

πισές από τις διαφορές αυτές ως "+" και

πισές από τις διαφορές αυτές ως "-".

```
# DATA
x=c(-4.7,3.7,22.4,23.5,14.4,13.6,8.7,9.1,20.2,6.5,-7.8,10.8,15.6,10.1,-6.9)
library("BSDA")
SIGN.test(x, md=0, alternative = "two.sided")
##SIGN.test(x, md=0, alternative= "greater")
##SIGN.test(x, md=0, alternative= "less")
```

-11-

One-sample Sign-Test

data: x
s = 12, p-value = 0.03516
alternative hypothesis: true median is not equal to 0
95 percent confidence interval:
4.198872 15.386198
sample estimates:
median of x
10.1

	Conf.Level	L.E.pt	U.E.pt
Lower Achieved CI	0.8815	6.5000	14.4000
Interpolated CI	0.9500	4.1989	15.3862
Upper Achieved CI	0.9648	3.7000	15.6000