

Τετάρτη 17/02/21

## Μη παραμετρικές μέθοδοι:

Συνήθως δεν έχουμε προγνωστική γνώση των κατανομών που περιγράφουν τα δεδομένα μας και σε τέτοιες περιπτώσεις η χρήση της παραμετρικής στατιστικής μπορεί να οδηγήσει σε λανθασμένα αποτελέσματα και συμπεράσματα.

Σε περιπτώσεις όπου δεν είναι επιβεβλητή οποιαδήποτε υποθέση για τη μορφή του γινόμενου ο ερευνητής χρησιμοποιεί το κ.ο.θ. ως να διεξάγει έλεγχο. Υπάρχουν ωστόσο προβλήματα όταν τα δεδομένα που έχουμε συλλέξει από την υπό μελέτη πληθυσμό είναι μικρά και τα δεδομένα δεν κατανοούνται σύμφωνα με την κανονική κατανομή.

~~Επιβεβλητή~~

Η βασική διαφορά της μη παραμετρικής στατιστικής συμπερασματικής είναι να εστιάζει συμπεράσματα για γνωστές ποσότητες - χαρακτηριστικά. Δεδομένα από παραπλήσια πληθυσμούς το μικρότερο είναι ομοιογενή υποθέσεων. Το κυριότερο πλεονέκτημα είναι ότι οι μη-παραμετρικές τεχνικές δεν προϋποθέτουν γνώση της κατανομής του γινόμενου.

- Απαιτούν λιγότερες υποθέσεις για τους γινόμενους από τους οποίους προέρχονται τα δεδομένα
- Επιτρέπουν στον ερευνητή να υπολογίσει σε κάποιες περιπτώσεις την ακριβή τιμή του παρατηρημένου συντελεστή σημαντικότητας  $\alpha$ .
- Πιο σίτες τεχνικές σε ετήσιμη και συχνά πιο εύκολες στην κατανόηση
- Δεν είναι ευαίσθητες στις ακραίες τιμές
- Εφαρμόζονται συχνά στις τάξεις των δεδομένων και όχι στα ίδια τα δεδομένα
- Μπορούν να χρησιμοποιηθούν και για δεδομένα που είναι γινόμενα σε κατηγορίες (κλίμακα διατάξης ή αντιστοιχική κλίμακα)
- Σε περίπτωση που είναι γνωστό ότι τα δεδομένα προέρχονται από κανονική κατανομή, οι έλεγχοι είναι πιο ~~επιβεβλημένοι~~ <sup>ισχυροί</sup>.
- Αν η κατανομή των δεδομένων δεν είναι η κανονική τότε οι μη-παραμετρικοί έλεγχοι είναι πλεονεκτήματα έναντι των παραμετρικών.

Βασικό Ερώτημα: Όταν δέν έχω / δέν γέρω πληροφορίες για σ.π.π  
Πως μπορώ να τν εκτιμώ;

Η εκτίμηση για να υλοισάνει χρειάζεται και ένα δείγμα  $X_1, \dots, X_n$  από τν  
Γνωστέον σ.π.π  $f(x)$ .

Η μέθοδος της μη παραμετρικής εκτίμησης έχει διάφορες με την παραμ. εκτίμηση

• Η παραμετρική εκτίμηση της σ.π.π

→ υιοθετεί μερική γνώση του υπό εκτίμηση μοντέλου, μυσών τουλάχιστον της ακεραιότητας  
του συντελεστή ή κατανομής

→ Δίνει πιο ακριβή αποτελέσματα από τν μη παραμετρική προσέγγιση και  
είναι το μοντέλο που υιοθετούμε είναι σωστό.

Σκοπός είναι να καταλάβουμε να χρησιμοποιήσουμε με βάση τα δεδομένα  
πότε είναι η κατάλληλη περίπτωση για να εφαρμόσουμε κάθε μια από τις δύο  
μεθόδους.

Μη παραμετρικοί εκτιμητές των σ.π.π  $f(x)$  και α.σ.κ  $F(x)$  είναι

$$\hat{f}(x) = (nh)^{-1} \sum_{i=1}^n \mathbb{1}_{\{x-h \leq X_i \leq x+h\}} \quad \rightarrow \text{δείγματα συνεισφέρον}$$

η  $\hat{f}(x)$  είναι ο ορισμός του λογιγράμματος και η παράμετρος  $h$   
είναι το είδος της μετρικής του.

Πως λειτουργεί συν γράφει το λογιγράμμο:

• Ξεκινώντας από ένα αρχικό σημείο  $x_0$ , δημιουργούμε τα διαστήματα

$$I_j = (x_0 + jh, x_0 + (j+1)h], \quad j = -L, 0, L$$

Εκεί έχουμε χυρίσει το π.ο της  $f$  σε υποδιαστήματα

• Η γνωστή  $f(x)$  ~~εκτιμάται~~ εκτιμάται μετρώντας τον αριθμό των παρατηρήσεων  
 $X_1, \dots, X_n$  που πέφτουν σε κάθε  $I_j$ .

Αυτή η εκτίμηση είναι η εμπειρική πιθανότητα να πέσει μια οποιαδήποτε  
παρατήρηση στο δ/μο  $I_j$  και προσεγγίζει τον αριθμό της σ.π.π.



- Όταν τμήμα της  $f(x)$  όλα είναι γραμμένα εκτός ίσως από το  $h$ , το οποίο δίνεται από τμήματα που εμφανίζονται μόνο από τα δεδομένα. Η τιμή του  $h$  επηρεάζει την ακρίβεια της εκτίμησης.

Παράδειγμα από τα δεδομένα geyser της R από το πακέτο MASS

$$\boxed{\text{Κανόνας Sturges}} \quad h = \lceil \log_2(299) \rceil + 1 \approx 10$$

### Ιδιότητες ιστιογράμματος:

Μέση τιμή και διακύμανση της  $\hat{f}(x)$

Αν  $n$   $f(x)$  είναι συνεχής ομακρύνη πυκνότητα πιθανότητας.

$$E(\hat{f}(x)) = f(x) (1 + o(1))$$

$$\text{Var}(\hat{f}(x)) = (nh)^{-1} f(x) + o(h^{-2})$$

### Απόδειξη:

Υποθέτουμε ότι για συνεχή  $f(x)$   $X \sim f(x)$  ισχύει:

$$E(g(x)) = \int g(x) f(x) dx$$

$X_1, \dots, X_n$  i.i.d. ενομομετρήσιμα

$$E\{\hat{f}(x)\} = E\left[(nh)^{-1} \sum_{i=1}^n \mathbb{1}_{\left\{x - \frac{h}{2} \leq X_i \leq x + \frac{h}{2}\right\}}\right] = n(nh)^{-1} E\left[\mathbb{1}_{\left\{x - \frac{h}{2} \leq X_1 \leq x + \frac{h}{2}\right\}}\right]$$

$$= h^{-1} \int_{x-h/2}^{x+h/2} f(t) dt + h^{-1} \int_{x-h/2}^{x+h/2} o(f(y)) dy =$$

$$= h^{-1} \int_{x-h/2}^{x+h/2} f(t) dt = f(x) (1 + o(1)).$$

Προκύπτει αν εφαρμόσουμε τον ισχυρισμό των μεγάλων αριθμών και τις ιδιότητες του Θεωρήματος Riemann

Βάση του στοιχείου:  $(nh)^{-1} \xrightarrow{a.s.} h^{-1} \int_{x-h/2}^{x+h/2} f(y) dy \approx f(x)$

**Τοξότης νόμος των μεγάλων αριθμών:**

Έστω  $X_1, X_2, \dots$  μια ακολουθία ανεξάρτητων και ισόκυμων τ.μ με  $E(X_i) = \mu$  και πεπεσμένη διακύμανση  $Var(X_i) = \sigma^2, i=1, 2, \dots$ . Τότε η ακολουθία συγκλίνει με πιθανότητα 1 στη θεωρητική μέση τιμή  $\mu$ :

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \text{ συγκλίνει με πιθανότητα 1 στη θεωρητική μέση τιμή } \mu$$

$$\text{δηλ. } P(\lim_{n \rightarrow \infty} \bar{X}_n = \mu) = 1$$

**Πα του Σαυβανου:**

$$Var\{g(x)\} = E g^2(x) - [E(g(x))]^2$$

$$Var\{\hat{f}(x)\} = Var\left\{(nh)^{-1} \sum_{i=1}^n \mathbb{1}_{\left\{x_i - \frac{h}{2} < X_i < x + \frac{h}{2}\right\}}\right\} = E\left\{(nh)^{-1} \sum_{i=1}^n \mathbb{1}_{\{\dots\}}\right\}^2 - [E\{\hat{f}(x)\}]^2$$

Παρατηρούμε ότι  $\left(\sum_{i=1}^n X_i\right)^2 = \sum_{i=1}^n X_i^2 + \sum_{i \neq j} X_i X_j$

$$\text{Έτσι έχουμε } E\{\hat{f}(x)\}^2 = E\left\{(nh)^{-1} \sum_{i=1}^n \mathbb{1}_{\{\dots\}}\right\}^2 =$$

$$= (nh)^{-2} \left\{ \sum_{i=1}^n E \mathbb{1}_{\{\dots\}}^2 + \sum_{i \neq j} E \mathbb{1}_{\{\dots\}} \mathbb{1}_{\{\dots\}} \right\} \text{ (1)}$$

$$= (nh)^{-2} \left\{ n E \mathbb{1}_{\{\dots\}}^2 + n(n-1) E \mathbb{1}_{\{\dots\}} \mathbb{1}_{\{\dots\}} \right\} \text{ (1)}$$

$$h^{-1} E \mathbb{1}_{\{\dots\}}^2 = f(x) (1 + o(1)) \text{ (2)}$$

$$h^{-2} E \mathbb{1}_{\{\dots\}} \mathbb{1}_{\{\dots\}} = h^{-1} E \mathbb{1}_{\{\dots\}} \cdot h^{-1} E \mathbb{1}_{\{\dots\}} =$$

$$= \left( \int_{x-h/2}^{x+h/2} f(t) dt \right) \left( \int_{x-h/2}^{x+h/2} f(z) dz \right) = f(x)^2 (1 + o(1)) \text{ (3)}$$



Επιπλέον  $E\{\hat{f}(x)\}^2 = nh^{-1}f(x) + \frac{n(n-1)}{n^2} f(x)^2$

Τελικά  $\text{Var}\{\hat{f}(x)\} = nh^{-1}f(x) - \frac{f(x)^2}{n}$

Κεντρικό πρόβλημα παύει να προσέτα  $x_i \in I_j$ .

Η κατανομή αυτής της ποσότητας συστήνεται με τον ορισμό της διωνυμικής κατανομής. όπου η πιθανότητα να πέσει η παρατήρηση στο δίκτυο  $I_j$  είναι  $p_j = \int_{I_j} f(x) dx$ , όπου  $f(x)$  η πυκνότητα.

αλλά έχουμε κατανομή

Κατανομή της  $f(x_i)$

Για  $x_j \in I_j = [x_j - h/2, x_j + h/2]$  τότε  $f(x_j) \sim \text{binomial}(n, p_j)$

Επιπλέον της α.ο.κ.  $F(x)$

$$\hat{F}(x) = n^{-1} \sum_{i=1}^n \mathbb{1}\{x_i \leq x\}$$

Θεωρούμε ότι ο  $\hat{F}(x)$  προκύπτει από την  $\hat{F}(x)$ .

$\hat{F}(x) = n^{-1} \sum_{i=1}^n \mathbb{1}\{x_i < x\}$  ένας διηρημένος εκτιμητής της  $F(x) = P(X \leq x)$ .

$$F(x+h) - F(x-h) \approx n^{-1} \sum_{i=1}^n \mathbb{1}\{x_i < x+h\} - n^{-1} \sum_{i=1}^n \mathbb{1}\{x_i - h \leq x_i\} =$$

$$= n^{-1} \sum_{i=1}^n \mathbb{1}\{x-h < x_i < x+h\} \quad \text{και προσαρμόζουμε}$$

$$\hat{F}(x) = \frac{1}{nh} \sum_{i=1}^n \mathbb{1}\{x-x_i < h\} = \frac{1}{nh} \sum_{i=1}^n \mathbb{1}\left\{\frac{x-x_i}{h} < 1\right\}$$

Παρατηρούμε ότι  $\sum_{i=1}^n \mathbb{1}\{x_i \leq x\} = n \hat{F}(x) \sim B(n, F(x))$ .

Είναι άρα άμεσος Bernoulli τ.κ. με π.θ.  $p = F(x)$ .

Εστιάμε ότι αν  $X \sim B(n, p)$  τότε  $EX = n \cdot p$ ,  $Var(X) = np(1-p)$

Αρα  $n \hat{F}(x) = n^{-1} \sum_{i=1}^n \mathbb{1}\{X_i \leq x\}$  έχει Διωνυμική κατανομή

με

$$E\hat{F}(x) = n^{-1} n F(x) = F(x)$$

$$Var(\hat{F}(x)) = n^{-2} n F(x)(1-F(x)) = \frac{1}{n} F(x)(1-F(x))$$

Με βάση τις ιδιότητες και με εφαρμογή του Κ.Θ προκύπτει

**Ασπυρωτική κατανομή της  $\hat{F}(x)$ .**

Κάθως το  $n \rightarrow \infty$ .

$$\hat{F}(x) \sim N(F(x), n^{-1} F(x)(1-F(x)))$$

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq x\}, X_1, \dots, X_n \text{ στατιστικώς ανεξάρτητα δείγματα από την}$$

μικρότερη από μεγαλύτερη ~~κατανομή~~ <sup>παρατηρηθών</sup>

Βλέπουμε ότι  $\hat{F}(x)$  βασίζεται στην περιγραφή με διαδοχικά ύψη ανά για  $\frac{1}{n}$  σε κάθε σημείο  $X_i$  δηλ. δίνει κλάση μεθυσμένων  $\frac{1}{n}$  σε κάθε

$X_i$ , άρα είναι αόξερα και δέξια ανεξάρτητα.

# ΜΗ ΠΑΡΑΜΕΤΡΙΚΗ ΣΤΑΤΙΣΤΙΚΗ ΜΑΘΗΜΑ 1 ΕΚΤΙΜΗΣΗ ΣΠΠ ΜΕ ΙΣΤΟΓΡΑΜΜΑ ΣΤΗΝ R

π.χ. 1

```

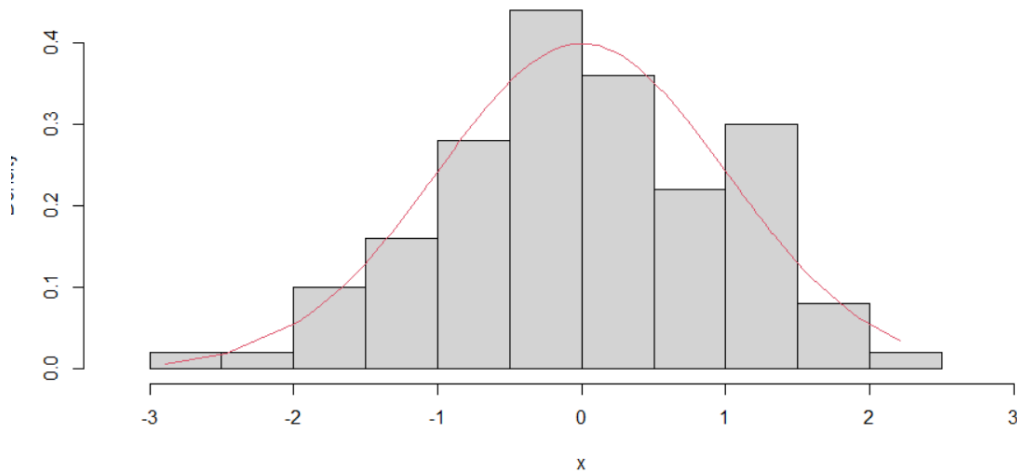
1 #μη παραμετρική στατιστική εκτίμηση της σ.π.π με
2 #ιστόγραμμα
3 set.seed(20)
4 x<-rnorm(100)
5 x<-sort(x)
6 hist(x,probability=TRUE,ylim=c(0,0.45),xlim=c(-3.2,3.2))
7 lines(x,dnorm(x),col=2)
8
9 set.seed(20)
10 x<-rnorm(200)
11 x<-sort(x)
12 hist(x,probability=TRUE,ylim=c(0,0.45),xlim=c(-3.2,3.2))
13 lines(x,dnorm(x),col=2)
14 #με την αύξηση των παρατηρήσεων μειώθηκε η διακύμανση όπως μπορούμε να δούμε και από το θεώρημα στις σημειώσεις
15 #γενικά όσο πιο πολλές παρατηρήσεις έχουμε τόσο πιο καλή θα είναι η εκτίμηση
16
17
18 #δοκιμάζουμε για n=10000
19 set.seed(20)
20 x<-rnorm(10000)
21 x<-sort(x)
22 hist(x,probability=TRUE,ylim=c(0,0.45),xlim=c(-3.2,3.2))
23 lines(x,dnorm(x),col=2)
24 #παρατηρούμε στο ιστόγραμμα ότι έχουμε αισθητά καλύτερη προσέγγιση
    
```

π.χ. 2

π.χ. 3

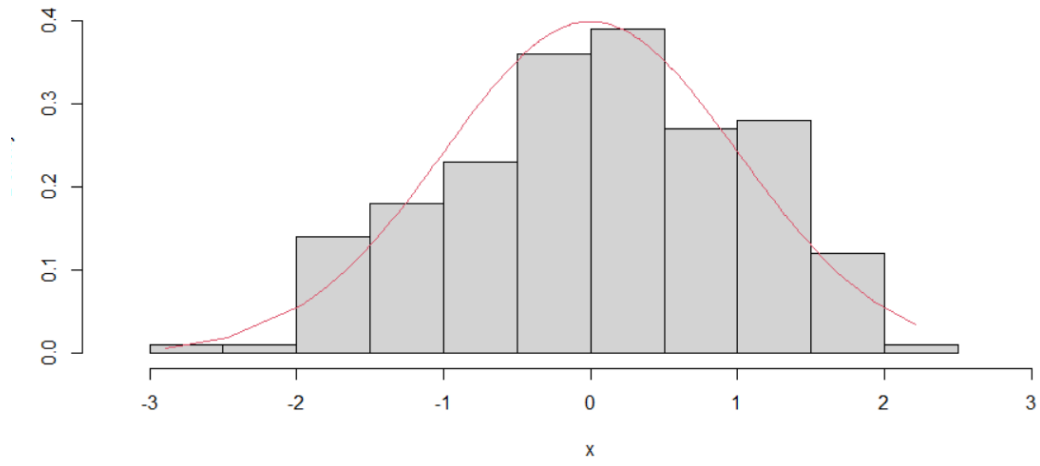
→ 100 τυχαίες τιμές από κανονική κατανομή  
 → αντίστοιχα το ιστόγραμμα

Histogram of x



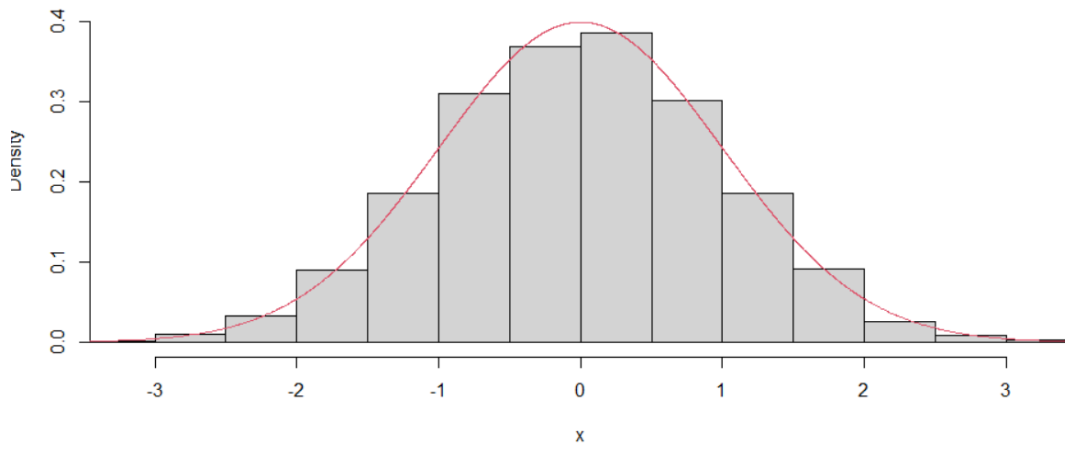
π.χ. 1

Histogram of x



n.x.2

Histogram of x



n.x.3