



ΣΤΑΤΙΣΤΙΚΗ ΣΤΗΝ ΠΛΗΡΟΦΟΡΙΚΗ

Συμπεράσματα για Αναλογίες και Πίνακες Κατηγοριών

Εργασία 3

Μαρία Σχοινάκη
Σοφία Παπαϊωάννου

Άσκηση 1

a) Το επίπεδο εμπιστοσύνης είναι 95%, οπότε:

Συνθήκες για την ακρίβεια του z-ελέγχου:

κορώνες = 29 > 15

γράμματα = 21 > 15

Οι συνθήκες αυτές ικανοποιούνται, συνεπώς μπορούμε να χρησιμοποιήσουμε την κανονική προσέγγιση.

$$\hat{p} = \frac{29}{50} = 0.58 = 58\%$$

Χρησιμοποιούμε το τύπο για το διάστημα εμπιστοσύνης της δυαδικής πιθανότητας:

$$CI = \hat{p} \pm z \cdot \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

ή αλγοριθμικά $p_hat + c(-1, 1)*(-1)*qnorm(0.025)*sqrt(p_hat*(1 - p_hat)/n)$

Το διάστημα εμπιστοσύνης για επίπεδο εμπιστοσύνης 95% είναι [0.4431951, 0.7168049].

b) $H_0 : p = 0.5$ (Το νόμισμα είναι δίκαιο)

$H_a : p = 0.5$ (Το νόμισμα **δεν** είναι δίκαιο)

Συνθήκες για την ακρίβεια του z-ελέγχου:

1. Ο μέσος αριθμός κορώνων: $n \cdot p_0 = 50 \cdot 0.5 = 25 > 10$

2. Ο μέσος αριθμός γραμμάτων: $n \cdot (1 - p_0) = 50 \cdot 0.5 = 25 > 10$

Οι συνθήκες αυτές ικανοποιούνται, συνεπώς μπορούμε να χρησιμοποιήσουμε την κανονική προσέγγιση.

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \quad p\text{-value} = 2 \cdot P(Z > |z|)$$

ή αλγοριθμικά

$$z = (\hat{p} - 0.5) / \text{sqrt}(0.5*0.5/n) = \mathbf{1.131371}$$

$$p\text{-value} = 2*\text{pnorm}(-z) = \mathbf{0.257899} \sim \mathbf{26\%}$$

Με επίπεδο σημαντικότητας **5%**, το p-value είναι **~26% > 5%** και επομένως **δεν** απορρίπτουμε τη μηδενική υπόθεση ($p = 0.5$). Συνεπώς τα δεδομένα **δεν** είναι στατιστικά σημαντικά για την εναλλακτική υπόθεση ($p \neq 0.5$). Άρα καταλήγουμε ότι το νόμισμα μπορεί να θεωρηθεί **δίκαιο**.

c) Με χρήση του τύπου:

$$n \geq \frac{z_*^2}{4m^2}$$

Παίρνουμε:

$$n \geq \frac{1.96*1.96}{4*0.01*0.01} = \mathbf{9604}$$

Δηλαδή θα πρέπει να εκτελέσουμε **9604** ρίψεις για διάστημα εμπιστοσύνης **95%** με περιθώριο λάθους μικρότερο **1%**.

Άσκηση 2

- Επίπεδο εμπιστοσύνης: **95%**
- Περιθώριο σφάλματος: **3%** ($m=0.03$)

Ο τύπος για τον υπολογισμό του μεγέθους δείγματος είναι:

$$n \geq \frac{z^2}{4m^2}$$

Όπου:

- $Z = 1.96$: η τιμή z^* για 95% επίπεδο εμπιστοσύνης,
- m : το περιθώριο σφάλματος.

Αντικαθιστούμε τις τιμές στον τύπο:

$$n \geq \frac{1.96^2}{4 \cdot 0.03^2}$$

$$= 1067.111 \approx 1068$$

Στρογγυλοποιούμε προς τα πάνω:

$$n \approx 1100 \text{ άτομα}$$

Παρατήρηση:

Το μέγεθος δείγματος **δεν εξαρτάται από τον πληθυσμό** όταν αυτός είναι πολύ μεγάλος (π.χ., Ελλάδα: 10 εκατομμύρια, Η.Π.Α.: 300 εκατομμύρια). Το μόνο που επηρεάζει το μέγεθος δείγματος είναι το επίπεδο εμπιστοσύνης (z) και το περιθώριο σφάλματος (m). Συνεπώς, το δείγμα στις Η.Π.Α. παραμένει ίδιο με αυτό της Ελλάδας, περίπου **1100 άτομα**.

Άσκηση 3

a) Για να είναι ακριβής ο έλεγχος σημαντικότητας θα πρέπει:

- Τα δείγματα να έχουν ληφθεί με ανεξάρτητο τρόπο από τους 2 πληθυσμούς και το καθένα να είναι SRS.
- Για τις γυναίκες:
 - # «επιτυχίων» (καπνίστριες) = 14 > 5
 - # «αποτυχιών» (μη καπνίστριες) = 16 > 5
- Για τους άντρες:
 - # «επιτυχίων» (καπνιστές) = 12 > 5
 - # «αποτυχιών» (μη καπνιστές) = 18 > 5

$$H_0 : p_f = p_m \Rightarrow p_f - p_m = 0$$

$$H_a : p_f \neq p_m \Rightarrow p_f - p_m \neq 0$$

$$n_f = \text{sum}(\text{SEX} == 'F')$$

$$\hat{p}_f = \text{sum}(\text{SEX} == 'F' \ \& \ \text{SMOKER} == 'Y') / n_f = \mathbf{0.4666667}$$

$$n_m = \text{sum}(\text{SEX} == 'M')$$

$$\hat{p}_m = \text{sum}(\text{SEX} == 'M' \ \& \ \text{SMOKER} == 'Y') / n_m = \mathbf{0.4}$$

$$p = \text{sum}(\text{SMOKER} == 'Y') / (n_f + n_m) = \mathbf{0.4333333}$$

$$z = (\hat{p}_f - \hat{p}_m) / \sqrt{p(1-p)/n_f + p(1-p)/n_m} = \mathbf{0.5210501}$$

$$p\text{-value} = 2 * \text{pnorm}(-z) = \mathbf{0.6023319}$$

Το **p-value** είναι περίπου **60%**, πολύ μεγαλύτερο από το επίπεδο σημαντικότητας $\alpha=5$. Εφόσον $p\text{-value} > \alpha$, **δεν** απορρίπτουμε τη μηδενική υπόθεση ($H_0 : p_f = p_m$).

Συνεπώς, τα δεδομένα **δεν** δείχνουν στατιστικά σημαντική διαφορά στα ποσοστά καπνιστών μεταξύ ανδρών και γυναικών. Το φύλο **δεν** φαίνεται να επηρεάζει την πιθανότητα καπνίσματος.

b) Για να είναι ακριβής ο υπολογισμός του διαστήματος εμπιστοσύνης θα πρέπει:

- Τα δείγματα να έχουν ληφθεί με ανεξάρτητο τρόπο από τους 2 πληθυσμούς και το καθένα να είναι SRS.
- Για τις γυναίκες:
 - # «επιτυχίων» (καπνίστριες) = 14 > 10
 - # «αποτυχιών» (μη καπνίστριες) = 16 > 10
- Για τους άντρες:
 - # «επιτυχίων» (καπνιστές) = 12 > 10
 - # «αποτυχιών» (μη καπνιστές) = 18 > 10

$(\text{pf_hat} - \text{pm_hat}) + c(-1,1)*(-1) * \text{qnorm}(0.025) * \sqrt{\text{pf_hat}*(1-\text{pf_hat})/\text{sum}(\text{SEX} == 'F') + \text{pm_hat}*(1-\text{pm_hat})/\text{sum}(\text{SEX} == 'M')}$

Το διάστημα εμπιστοσύνης για επίπεδο εμπιστοσύνης 95% είναι:

[-0.1835364, 0.3168697]

c) Έλεγχος ανεξαρτησίας μεταξύ φύλου και καπνίσματος

Υποθέσεις

1. **Μηδενική Υπόθεση (H0):**

Το φύλο και το κάπνισμα είναι ανεξάρτητα (δεν υπάρχει σχέση μεταξύ τους).

2. **Εναλλακτική Υπόθεση (Ha):**

Το φύλο και το κάπνισμα **δεν** είναι ανεξάρτητα (υπάρχει σχέση μεταξύ τους).

```
> addmargins(table_data)
```

Πίνακας Συνάφειας:

	sex		
smoker	F	M	Sum
N	16	18	34
Y	14	12	26
Sum	30	30	60

d)

```
> chi_test <- chisq.test(table_data, correct=FALSE)
> chi_test
```

Pearson's Chi-squared test

```
data: table_data
X-squared = 0.27149, df = 1, p-value = 0.6023
```

Το **p-value** που υπολογίστηκε με τον χ^2 -έλεγχο είναι **0.6023**, το οποίο είναι ίσο με το **p-value** που προέκυψε στον **z-έλεγχο** (στο ερώτημα (α)), όπως ήταν αναμενόμενο, δεδομένου ότι ο πίνακας συνάφειας είναι διαστάσεων **2x2**, όπου οι δύο έλεγχοι είναι ισοδύναμοι. Επιπλέον, παρατηρούμε ότι το χ^2 -στατιστικό (**0.27149**) είναι ίσο με το τετράγωνο του **z-στατιστικού** ($z^2 = 0.52105^2 = 0.27149$), επαληθεύοντας τη μαθηματική σχέση μεταξύ τους. Το αποτέλεσμα αυτό επιβεβαιώνει ότι και οι δύο έλεγχοι οδηγούν στο ίδιο συμπέρασμα, δηλαδή **δεν απορρίπτεται η μηδενική υπόθεση και δεν υπάρχουν στατιστικά σημαντικές ενδείξεις ότι υπάρχει σχέση μεταξύ φύλου και καπνίσματος**.

Άσκηση 4

α) Για να είναι ακριβής ο υπολογισμός θα πρέπει:

- Το δείγμα να έχει ληφθεί με SRS από τον πληθυσμό (*πληθυσμός τα κόκκινα και μπλε smarties που παρασκευάζονται και δείγμα τα 34 συνολικά smarties της σακούλας που είναι είτε μπλε είτε κόκκινα*)
- Μέσο πλήθος επιτυχιών $n(p_0) = 34 \cdot 0.5 = 17 > 10$
- Μέσο πλήθος αποτυχιών $n(1-p_0) = 34 \cdot 0.5 = 17 > 10$

H_0 : $\text{pred} \leq 0.5$ “Δεν παρασκευάζονται περισσότερα κόκκινα από μπλε”

H_a : $\text{pred} > 0.5$ “Παρασκευάζονται περισσότερα κόκκινα από μπλε”

$$\hat{p}_{\text{red}} = 19 / (19+15) = 19 / 34 = 0.5588235$$

$$z = (\hat{p}_{\text{red}} - 0.5) / \sqrt{0.5(1 - 0.5)/34} = 0.68599434$$

$$p - \text{value} = 1 - \text{pnorm}(z) = 0.2463583$$

Το **p-value** είναι περίπου **25%**, αρκετά μεγαλύτερο από το επίπεδο σημαντικότητας (**$\alpha=5\%$**) γεγονός που σημαίνει ότι **δεν απορρίπτουμε τη μηδενική υπόθεση ($p_{\text{red}} \leq 0.5$)**. Με βάση αυτό, τα δεδομένα δεν υποστηρίζουν την εναλλακτική υπόθεση ότι παρασκευάζονται περισσότερα κόκκινα smarties από μπλε. Συνεπώς, **δεν** μπορούμε να συμπεράνουμε ότι τα κόκκινα smarties υπερτερούν αριθμητικά των μπλε.

b) H_0 : Η κατανομή του πληθυσμού συμφωνεί με την:
 $p = <0.198, 0.178, 0.176, 0.196, 0.252>$

H_a : Ο πληθυσμός έχει διαφορετική κατανομή

```
> d <- c(22, 19, 16, 15, 8)
> p <- c(0.198, 0.178, 0.176, 0.196, 0.252)
>
> chisq.test(d, p = p)
```

Chi-squared test for given probabilities

```
data: d
X-squared = 11.613, df = 4, p-value = 0.02048
```

Το **p-value** είναι περίπου **2%**, το οποίο είναι **μικρότερο** από το επίπεδο σημαντικότητας ($\alpha=5\%$) Αυτό σημαίνει ότι **απορρίπτουμε** τη μηδενική υπόθεση. Συνεπώς, τα δεδομένα παρέχουν στατιστικά σημαντικές ενδείξεις ότι η κατανομή των χρωμάτων έχει **αλλάξει** σε σχέση με το **2009**.

c) Έλεγχος ομοιογένειας

H_0 : Η αναλογία χρωμάτων στα smarties είναι ίδια με αυτή στα M&Ms

H_a : Οι δύο αναλογίες είναι διαφορετικές

```
> t <- matrix(c(22, 19, 16, 15, 8, 10, 12, 20, 9, 5), 5, 2)
> colnames(t) <- c('Smarties', 'M&Ms')
> rownames(t) <- c('Brown', 'Red', 'Yellow', 'Blue', 'Green')
> t <- as.table(t)
>
> t
```

	Smarties	M&Ms
Brown	22	10
Red	19	12
Yellow	16	20
Blue	15	9
Green	8	5

```
> chisq.test(t, correct = FALSE)
```

Pearson's Chi-squared test

```
data: t
X-squared = 4.6262, df = 4, p-value = 0.3278
```

Το **p-value** είναι **~33%** και άρα αρκετά μεγάλο ώστε **δεν** απορρίπτουμε τη μηδενική υπόθεση. Συνεπώς τα δεδομένα **δεν** είναι στατιστικά σημαντικά για την εναλλακτική υπόθεση. Οπότε συμπεραίνουμε ότι η αναλογία χρωμάτων πρέπει να είναι η **ίδια**.