

Στατιστική στην Πληροφορική - Εργασία 2

Αυγεράκης Βασίλης - p3210013
Παντελίδης Ιπποκράτης - p3210150

Άσκηση 1)

a) Τα δεδομένα είναι κατάλληλα για να χρησιμοποιήσουμε τις μεθόδους συμπερασματολογίας που γνωρίζουμε αφού :

- Πρόκειται για τυχαία δειγματοληψία απο πληθυσμό (ημερήσιες αιτήσεις από βάση δεδομένων).
- Το δείγμα είναι αρκετά μεγάλο, 20 περιπτώσεις ($n \geq 15$).
- Η ατυπική τιμή 284 μπορεί να υποδεικνύει ότι η κατανομή δεν είναι κανονικά κατανεμημένη, ωστόσο λόγω του ικανοποιητικού μεγέθους του δείγματος η ακρίβεια της κατανομής t είναι αρκετά καλή και σε μη κανονικά κατανεμημένους πληθυσμούς.

b) Ο τύπος που χρησιμοποιούμε για να βρούμε το διάστημα εμπιστοσύνης για την μέση τιμή για επίπεδο εμπιστοσύνης C είναι ο :

$$C: \bar{x} \pm t_* \frac{s}{\sqrt{n}}$$

Καλώντας τις παρακάτω εντολές στην R αφού εκχωρήσουμε στην μεταβλητή `data` τα δεδομένα μας και για επίπεδο εμπιστοσύνης 95% παίρνουμε ότι το διάστημα εμπιστοσύνης είναι το [51.41365, 103.38635]. Να σημειωθεί ότι για τον υπολογισμό του t_* χρησιμοποιούμε βαθμό ελευθερίας 19 δηλαδή $n - 1$.

```
> n <- length(data)
> n
[1] 20
> xbar <- mean(data)
> xbar
[1] 77.4
> sd <- sd(data)
> sd
[1] 55.52467
> tstar <- -qt(0.025, df=n-1)
> tstar
[1] 2.093024
> xbar + c(-1,1) * tstar * sd / sqrt(n)
[1] 51.41365 103.38635
```

Άσκηση 2)

a) Είναι λάθος γιατί η τυπική απόκλιση του δειγματικού μέσου με βάση τον ορισμό είναι σ / \sqrt{n} δηλαδή $12 / \sqrt{20}$ και όχι $12 / 20$.

b) Το λάθος είναι ότι μας δίνεται σαν υπόθεση να εξετάσουμε αν ο δειγματικός μέσος είναι ίσος με μία τιμή ενώ στον έλεγχο σημαντικότητας θα έπρεπε να εξετάσουμε αν ο πληθυσμιακός μέσος ισούται με 10 ως την μηδενική υπόθεση. ($H_0 : \mu = 10$).

c) Η τιμή του δειγματικού μέσου δεν ανήκει στο διάστημα της εναλλακτικής υπόθεσης (> 54) καθώς το 45 είναι μικρότερο από το 54. Αυτό σημαίνει επίσης ότι το στατιστικό λάθος είναι αρνητικό και άρα το pvalue αρκετά μεγάλο. Το λάθος είναι ότι απορρίπτουμε την μηδενική υπόθεση.

δ) Είναι λάθος διότι για να απορριφθεί η μηδενική υπόθεση το pvalue θα πρέπει να είναι αρκετά πιο μικρό.

Άσκηση 3)

Κάνουμε έλεγχο σημαντικότητας με μηδενική υπόθεση $H_a : \mu = \mu_0$ και ξέρουμε ότι το στατιστικό ελέγχου z είναι ίσο με 1.34. `> z <- 1.34`

a) Το pvalue για την εναλλακτική υπόθεση $H_a : \mu > \mu_0$ το βρίσκω με την παρακάτω εντολή στην R :

```
> 1 - pnorm(z)
[1] 0.09012267
```

b) Το pvalue για την εναλλακτική υπόθεση $H_a : \mu < \mu_0$ το βρίσκω με την παρακάτω εντολή στην R :

```
> pnorm(z)
[1] 0.9098773
```

c) Το pvalue για την εναλλακτική υπόθεση $H_a : \mu \neq \mu_0$ το βρίσκω με την παρακάτω εντολή στην R :

```
> 2 * (1 - pnorm(abs(z)))
[1] 0.1802453
```

Άσκηση 4)

a) Το μικρό pvalue (μάλιστα μικρότερο από την σημαντικότητα επιπέδου 5%) μας οδηγεί στην απόρριψη της μηδενικής υπόθεσης. Αυτός ο

έλεγχος σημαίνει ότι είναι σχετικά απίθανο (4%) να παρατηρήσουμε δεδομένα τέτοια ώστε να ισχύει η μηδενική υπόθεση. Από τον ορισμό του διαστήματος εμπιστοσύνης, είναι 95% πιθανό να παρατηρήσουμε (με δειγματοληψία) δεδομένα που η μέση τους τιμή να ανήκει σε αυτό το διάστημα. Είναι ευνόητο λοιπόν ότι η τιμή 30 είναι εκτός του διαστήματος αφού αν ανήκε θα ήταν πιο πιθανό από 4% να ανήκει σε αυτό το διάστημα.

b) Εξ ορισμού η μείωση του επιπέδου εμπιστοσύνης αυξάνει κατά απόλυτη τιμή την τιμή του Z^* δίνοντας μας μεγαλύτερο περιθώριο λάθους γύρω από την μέση τιμή. Όμως $p\text{-value} < (1 - 0.9) = 0.1$ οπότε με παρόμοια εξήγηση με το προηγούμενο ερώτημα το 30 δεν ανήκει ούτε σε αυτό το διάστημα αφού απορρίπτεται πάλι η H_0 .

Άσκηση 5)

Πριν απαντήσουμε τα ερωτήματα a), b) και c) παρατηρούμε ότι έχει συμβεί ένα τυπογραφικό λάθος κατά την εισαγωγή των στοιχείων αφού δεν γίνεται κάποιος ενήλικας να έχει βάρος 6 κιλά. Οπότε δεν λαμβάνουμε υπόψη αυτή την περίπτωση στους υπολογισμούς μας. Άρα $n = 24$.

a) Η κατανομή του βάρους δεν είναι και τόσο μη συμμετρική, όπως φαίνεται από το παρακάτω stemplot, και αφού έχουμε 24 περιπτώσεις ($n \geq 15$) μπορούμε να εφαρμόσουμε τις μεθόδους που βασίζονται στην κατανομή t . Περνώντας τα δεδομένα του βάρους στην μεταβλητή `weight` και εκτελώντας τις παρακάτω εντολές στην R παίρνουμε το διάστημα $[69.57826, 78.00507]$ ως το διάστημα εμπιστοσύνης 95% για το μέσο βάρος των κατοίκων της Αθήνας. Να σημειωθεί ότι για τον υπολογισμό του t χρησιμοποιούμε βαθμό ελευθερίας 23 δηλαδή `n_weight - 1`.

```
Stemplot      > n_weight <- length(weight)
               > n_weight
               [1] 24
5 | 459        > mean_weight <- mean(weight)
               > mean_weight
               [1] 73.79167
6 | 5789       > sd_weight <- sd(weight)
               > sd_weight
               [1] 9.978146
7 | 01223357   > tstar_weight <- -qt(0.025, df=n_weight-1)
               > tstar_weight
               [1] 2.068658
8 | 012336     > mean_weight + c(-1,1) * tstar_weight * sd_weight / sqrt(n_weight)
               [1] 69.57826 78.00507
9 | 12
```

b) Για να απαντήσουμε στο ερώτημα μπορούμε να θεωρήσουμε ότι οι απο τις 24 περιπτώσεις, οι περιπτώσεις που έχουν τιμή Α στην μεταβλητή “ΦΥΛΟ” είναι το ανδρικό τυχαίο δείγμα και οι υπόλοιπες το γυναικείο. Οπότε έχουμε $n_m = 13$ (male) και $n_f = 11$ (female). Παρατηρούμε ότι οι κατανομές τους δεν είναι και τόσο μη συμμετρικές, σύμφωνα με τα παρακάτω stemplots οπότε πάλι μπορούμε να χρησιμοποιήσουμε την κατανομή t.

Stemplot (Male)

```
6 | 8
7 | 223357
8 | 0136
9 | 12
```

Stemplot (Female)

```
5 | 459
6 | 579
7 | 013
8 | 23
```

Με την χρήση της R κάνουμε τους απαραίτητους υπολογισμούς και για τα δύο φύλα αφού φορτώσουμε τα δεδομένα μας στις weight_m και weight_f αντίστοιχα :

```
> n_weight_m <- length(weight_m)
> n_weight_m
[1] 13
> mean_weight_m <- mean(weight_m)
> mean_weight_m
[1] 78.69231
> sd_weight_m <- sd(weight_m)
> sd_weight_m
[1] 7.598077

> n_weight_f <- length(weight_f)
> n_weight_f
[1] 11
> mean_weight_f <- mean(weight_f)
> mean_weight_f
[1] 68
> sd_weight_f <- sd(weight_f)
> sd_weight_f
[1] 9.570789
```

Τώρα για να βρούμε ένα διάστημα εμπιστοσύνης για την διαφορά του μέσου βάρους των δύο φύλων χρησιμοποιούμε τον ακόλουθο τύπο :

$$(\bar{x}_1 - \bar{x}_2) \pm t_* \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Οπότε βρίσκουμε το [5.789155, 15.595460] ως το διάστημα εμπιστοσύνης 80% για την διαφορά της μέσης τιμής του βάρους μεταξύ ανδρών και γυναικών (ενηλίκους κατοίκους Αθηνών) εκτελώντας τις παρακάτω εντολές στην R. Να σημειωθεί ότι για τον υπολογισμό του t_* χρησιμοποιούμε βαθμό ελευθερίας 10 δηλαδή το ελάχιστο εκ των $n_{\text{weight_m}} - 1$ και $n_{\text{weight_f}} - 1$.

```
> tstar <- -qt(0.1, df = min(n_weight_m - 1, n_weight_f - 1))
> tstar
[1] 1.372184
> (mean_weight_m - mean_weight_f) + c(-1,1) * tstar *
+ sqrt(((sd_weight_m^2) / n_weight_m) + ((sd_weight_f^2) / n_weight_f))
[1] 5.789155 15.595460
```

c) Θα ελέγξουμε εάν υπάρχει σημαντική διαφορά μεταξύ του μέσου βάρους των καπνιστών έστω μ_1 από το τυχαίο δείγμα και του βάρους όσων δεν είναι καπνιστές έστω μ_2 . Άρα θα εκτελέσουμε δίπλευρο έλεγχο σημαντικότητας $H_0: \mu_1 = \mu_2$.

Τα δεδομένα μας είναι κατάλληλα για συμπερασματολογία αφού προέρχονται από απλό τυχαίο δείγμα, δεν φαίνονται να μην είναι συμμετρικά από τα παρακάτω stemplots, και οι 24 περιπτώσεις είναι αρκετές ($n \geq 15$). Χωρίζουμε τα δεδομένα σε καπνιστές και μη καπνιστές και τα αναθέτουμε στις `weight_s` (smokers) και `weight_ns` (not smokers) αντίστοιχα. Με τις παρακάτω εντολές στην R υπολογίζουμε τα απαραίτητα για το έλεγχο στοιχεία μας :

Stemplot (Smoker)

```
5 | 9
6 | 5
7 | 137
8 | 0236
9 | 2
```

Stemplot(Not Smoker)

```
5 | 45
6 | 789
7 | 022335
8 | 13
9 | 1
```

```
> n_weight_s <- length(weight_s)
> n_weight_s
[1] 10
> mean_weight_s <- mean(weight_s)
> mean_weight_s
[1] 76.8
> sd_weight_s <- sd(weight_s)
> sd_weight_s
[1] 9.975526
```

```
> n_weight_ns <- length(weight_ns)
> n_weight_ns
[1] 14
> mean_weight_ns <- mean(weight_ns)
> mean_weight_ns
[1] 71.64286
> sd_weight_ns <- sd(weight_ns)
> sd_weight_ns
[1] 9.76341
```

Έπειτα υπολογίζουμε το στατιστικό ελέγχου t σύμφωνα με τον παρακάτω τύπο και με την ακόλουθη εντολή στην R παίρνουμε $t = 1.259715$.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

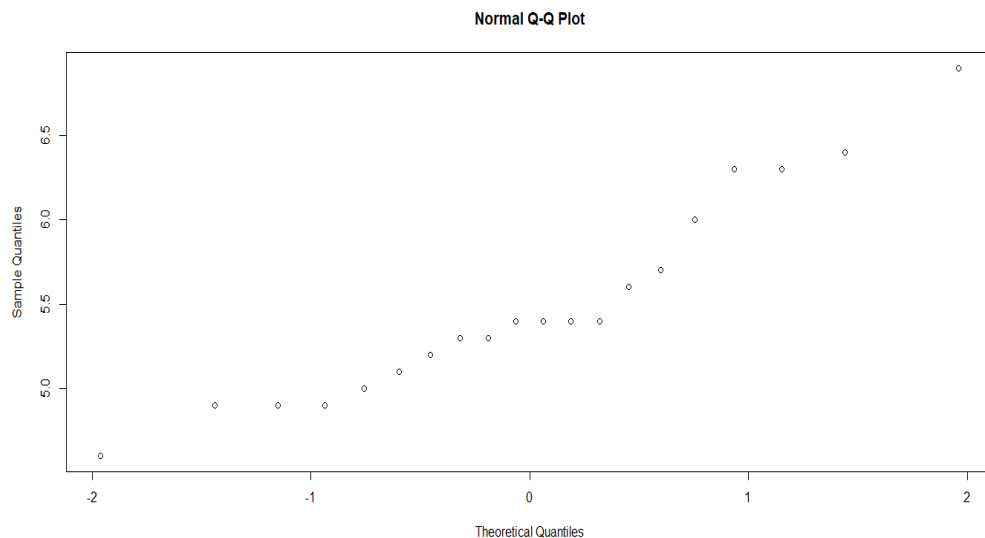
```
> t <- (mean_weight_s - mean_weight_ns) /
+ sqrt(((sd_weight_s^2) / n_weight_s) + ((sd_weight_ns^2) / n_weight_ns))
> t
[1] 1.259715
```

Οπότε προκύπτει με το παρακάτω υπολογισμό $pvalue = 0.2394573$, τιμή για την οποία δεν απορρίπτουμε την μηδενική υπόθεση. Να σημειωθεί ότι για τον υπολογισμό του $pvalue$ χρησιμοποιούμε βαθμό ελευθερίας 9 δηλαδή το ελάχιστο εκ των $n_weight_s - 1$ και $n_weight_ns - 1$.

```
> (1 - pt(t, df = min(n_weight_s - 1, n_weight_ns - 1))) * 2
[1] 0.2394574
```

Άσκηση 6)

a) Τα δεδομένα που μας δόθηκαν είναι κατάλληλα για συμπερασματολογία αφού αρχικά αποτελούν τυχαίο δείγμα (SRS). Εξετάζοντας το παρακάτω normal quantile plot παρατηρούμε ότι η κατανομή των δεδομένων δεν είναι και τόσο κανονική ωστόσο το πλήθος των περιπτώσεων είναι επαρκές $20 \geq 15$ οπότε τα δεδομένα μας να είναι κατάλληλα.



b) Αφού φορτώσουμε τα δεδομένα μας στην μεταβλητή `fuel_cons` κάνουμε τους ακόλουθους υπολογισμούς στην R :

```
> mean_fuel_cons <- mean(fuel_cons)
> mean_fuel_cons
[1] 5.5
> sd_fuel_cons <- sd(fuel_cons)
> sd_fuel_cons
[1] 0.6008766
```

c) Εκτελώντας τις παρακάτω εντολές στην R παίρνουμε το διάστημα $[5.218781, 5.781219]$ ως το διάστημα εμπιστοσύνης 95% για την μέση τιμή της απόδοσης του αυτοκινήτου. Να σημειωθεί ότι για τον υπολογισμό του t χρησιμοποιούμε βαθμό ελευθερίας 19 δηλαδή $n_{\text{fuel_cons}} - 1$.

```
> n_fuel_cons <- length(fuel_cons)
> n_fuel_cons
[1] 20
> tstar <- -qt(0.025, df= n_fuel_cons - 1)
> tstar
[1] 2.093024
> mean_fuel_cons + c(-1,1) * tstar * sd_fuel_cons / sqrt(n_fuel_cons)
[1] 5.218781 5.781219
```

Η ακρίβεια του ποσοστού αυτού είναι και ζήτημα της κατανομής απο την οποία προήλθαν τα δεδομένα γιατί κάνουμε υπόθεση “κανονικής” ενώ το δείγμα δεν πληρεί ικανοποιητικά όλα τα κριτήρια όπως αναδεικνύει και το ερώτημα α).

Άσκηση 7)

Για να ελέγξουμε την μηδενική υπόθεση, ότι δεν υπάρχει υπερεκτίμηση των ζημιών από το συνεργείο θα ελέγξουμε αν η μέση εκτίμηση της ζημιάς από το συνεργείο είναι μεγαλύτερη απο αυτή του εμπειρογνώμονα. Μπορούμε να παρατηρήσουμε ότι οι δύο αυτές εκτιμήσεις ζημιάς σχετίζονται. Καθώς όταν το συνεργείο κάνει μία υψηλή εκτίμηση, για παράδειγμα για κάποια μεγάλη ζημιά, είναι πολύ πιθανό και ο εμπειρογνώμονας να κάνει εξίσου υψηλή εκτίμηση. Αρα οι εκτιμήσεις δεν είναι ανεξάρτητες. Οπότε μπορούμε να χρησιμοποιήσουμε την διαφορά των δύο εκτιμήσεων και άρα τα δεδομένα μας είναι τα εξής :

```
[1] 100 50 -50 0 -50 200 250 200 150 300
```

Το δείγμα είναι τυχαίο (μετρήσεις για 10 αυτοκίνητα) και σύμφωνα με το παρακάτω stemplot θα μπορούσε η κατανομή του να ανήκει στην κανονική. Ωστόσο το μέγεθος του δείγματος είναι μικρό οπότε συνεχίζουμε τις μετρήσεις με ενδιαασμούς για την ακρίβεια των υπολογισμών.

Οπότε κάνουμε τον έλεγχο σημαντικότητας $H_0: \mu = 0$ και με εναλλακτική υπόθεση $H_a: \mu > 0$, όπου μ είναι η μέση τιμή των διαφορών. Η εναλλακτική υπόθεση λαμβάνει υπόψη μόνο το την περίπτωση ($>$) καθώς θέλουμε να ελέγξουμε αν υπερεκτιμάται η ζημιά από το συνεργείο δηλαδή αν η διαφορά γίνεται θετική. Με τις παρακάτω εντολές στην R κάνουμε τους απαραίτητους υπολογισμούς και υπολογίζουμε το στατιστικό ελέγχου t .

Stemplot

-0 | 55

0 | 05

1 | 05

2 | 005

3 | 0

```
> n <- length(data)
> n
[1] 10
> mean <- mean(data)
> mean
[1] 115
> sd <- sd(data)
> sd
[1] 124.8332
> t <- mean / (sd / sqrt(n))
> t
[1] 2.913182
```


Τέλος, με την ακόλουθη εντολή βρίσκουμε ότι το `rvalue` είναι ίσο με την τιμή 0.008610911, που είναι τόσο μικρή που αποκλείουμε την μηδενική υπόθεση. Η τιμή είναι τόσο μικρή που ακόμα και η κατανομή του πληθυσμού να μην είναι κανονικά κατανεμημένη αλλά κοντά σε αυτή, πάλι θα αποκλείαμε την μηδενική υπόθεση. Άρα τελικά, το συνεργείο υπερεκτιμά τις ζημιές. Να σημειωθεί ότι για τον υπολογισμό του `rvalue` χρησιμοποιούμε βαθμό ελευθερίας 9 δηλαδή το $n - 1$.

```
> 1 - pt(t, df= n - 1)
[1] 0.008610911
```

Ασκηση 8)

Πριν ξεκινήσουμε φορτώνουμε στην R τα δεδομένα απο το ερωτηματολόγιο του 2023 καθώς θα τα χρησιμοποιήσουμε στα ακόλουθα ερωτήματα.

a) Για να βρούμε ένα διάστημα εμπιστοσύνης για την διαφορά ύψους μεταξύ των ανδρών και των γυναικών του τμήματος πληροφορικής θα πρέπει από τα συνολικά δεδομένα να ξεχωρίσουμε τις μεταβλητές `male_height` και `female_height` που είναι το ύψος των ανδρών και των γυναικών αντίστοιχα. Αυτό το πετυχαίνουμε με τις παρακάτω εντολές στην R

```
> male_height <- data[gender == "M",]$height
> female_height <- data[gender == "F",]$height
```

Στην συνέχεια πάλι με την R κάνουμε τους απαραίτητους υπολογισμούς ώστε να βρούμε το διάστημα εμπιστοσύνης.

```
> n_male_height <- length(male_height) > n_female_height <- length(female_height)
> n_male_height > n_female_height
[1] 86 [1] 43
> mean_male_height <- mean(male_height) > mean_female_height <- mean(female_height)
> mean_male_height > mean_female_height
[1] 1.787093 [1] 1.666047
> sd_male_height <- sd(male_height) > sd_female_height <- sd(female_height)
> sd_male_height > sd_female_height
[1] 0.05978129 [1] 0.05113658
```

```
> tstar <- -qt(0.025, df = min(n_male_height - 1, n_female_height - 1))
> tstar
[1] 2.018082
```

Οπότε με τον τύπο που χρησιμοποιήσαμε και προηγουμένως παίρνουμε το διάστημα [0.1006281, 0.1414649] ως το 95% διάστημα εμπιστοσύνης

για την διαφορά του μέσου ύψους ύψους μεταξύ ανδρών και γυναικών φοιτητών πληροφορικής του ΟΠΑ.

```
> (mean_male_height - mean_female_height) + c(-1,1) * tstar *  
+ sqrt(((sd_male_height^2) / n_male_height) + ((sd_female_height^2) / n_female_height))  
[1] 0.1006281 0.1414649
```

Από την άλλη με την χρήση της εντολής t.test βγάζουμε το ακόλουθο διάστημα :

```
> t.test(male_height, female_height)  
  
Welch Two Sample t-test  
  
data: male_height and female_height  
t = 11.964, df = 96.701, p-value < 2.2e-16  
alternative hypothesis: true difference in means is not equal to 0  
95 percent confidence interval:  
 0.1009648 0.1411282  
sample estimates:  
mean of x mean of y  
 1.787093  1.666047
```

b) Για να εξετάσουμε αν ο μέσος βαθμών των ανδρών φοιτητών πληροφορικής, που έχουν πάρει το μάθημα “Στατιστική στην Πληροφορική”, στο μάθημα “Πιθανότητες” είναι μεγαλύτερος από τον μέσο βαθμό των γυναικών θα εφαρμόσουμε τον μονόπλευρο έλεγχο σημαντικότητας $H_0: \mu_1 = \mu_2$ με εναλλακτική υπόθεση $H_a: \mu_1 > \mu_2$ όπου μ_1 και μ_2 είναι ο μέσος βαθμός των ανδρών και των γυναικών στο μάθημα “Πιθανότητες”. Με τις παρακάτω εντολές στην R φορτώνουμε τα δεδομένα μας στις αντίστοιχες μεταβλητές και στην συνέχεια κάνουμε τους απαραίτητους υπολογισμούς ώστε να κάνουμε τον έλεγχο σημαντικότητας.

```
> male_prob <- data[gender == "M",]$prob  
> male_prob <- male_prob[!is.na(male_prob)]  
> female_prob <- data[gender == "F",]$prob  
> female_prob <- female_prob[!is.na(female_prob)]
```

```
> n_male_prob <- length(male_prob) > n_female_prob <- length(female_prob)  
> n_male_prob > n_female_prob  
[1] 81 [1] 30  
> mean_male_prob <- mean(male_prob) > mean_female_prob <- mean(female_prob)  
> mean_male_prob > mean_female_prob  
[1] 6.104938 [1] 5.966667  
> sd_male_prob <- sd(male_prob) > sd_female_prob <- sd(female_prob)  
> sd_male_prob > sd_female_prob  
[1] 2.838041 [1] 2.466791
```

Στην συνέχεια υπολογίζουμε το στατιστικό ελέγχου t με τον τύπο που χρησιμοποιήσαμε και προηγουμένως και εν τέλει βρίσκουμε p -value ίσο με 0.4016021, το οποίο είναι μεγαλύτερο από το επίπεδο εμπιστοσύνης 5%, οπότε δεν απορρίπτουμε την μηδενική υπόθεση. Επομένως οι άνδρες της πληροφορικής δεν επιτυγχάνουν μεγαλύτερες βαθμολογίες στο μάθημα “Πιθανότητες” από τις γυναίκες του τμήματος.

```
> t <- (mean_male_prob - mean_female_prob) /  
+ sqrt(((sd_male_prob^2) / n_male_prob) + ((sd_female_prob^2) / n_female_prob))  
> t  
[1] 0.2514961  
> 1 - pt(t, df = min(n_male_prob - 1, n_female_prob - 1))  
[1] 0.4016021
```

Από την άλλη με την χρήση της εντολής `t.test` βγάζουμε το ακόλουθο p -value, και καταλήγουμε στο ίδιο συμπέρασμα :

```
> t.test(male_prob, female_prob, alternative = "greater")  
  
Welch Two Sample t-test  
  
data: male_prob and female_prob  
t = 0.2515, df = 59.242, p-value = 0.4012  
alternative hypothesis: true difference in means is greater than 0  
95 percent confidence interval:  
-0.780425      Inf  
sample estimates:  
mean of x mean of y  
6.104938  5.966667
```

c) Πριν κάνουμε τον έλεγχο σημαντικότητας μπορούμε να σκεφτούμε ότι ο βαθμός ενός φοιτητή στο μάθημα “Μαθηματικά 1” δεν είναι τελείως ανεξάρτητος από τον βαθμό του στο μάθημα “Πιθανότητες” αφού πρόκειται για το ίδιο άτομο. Οπότε για να απαντήσουμε στο ερώτημα θα κρατήσουμε τις περιπτώσεις όπου έχουμε τον βαθμό του φοιτητή και στα δύο μαθήματα. Αυτό επιτυγχάνεται με την ακόλουθη εντολή στην R :

```
> filtered_data <- na.omit(data[c('math', 'prob')])  
> dim(filtered_data)  
[1] 107  2  
> prob <- filtered_data$prob  
> math <- filtered_data$math
```

Οπότε τώρα το δείγμα αποτελείται από 107 περιπτώσεις αντί 129 που εξακολουθεί να είναι ένας πολύ ικανοποιητικός αριθμός για συμπερασματολογία. Τώρα θα κάνουμε τον έλεγχο $H_0: \mu_1 - \mu_2 = 0$ με εναλλακτική u -

πόθεση $H_a: \mu_1 - \mu_2 \neq 0$ όπου μ_1 και μ_2 είναι μέσος βαθμός στο μάθημα των Μαθηματικών και των Πιθανοτήτων αντίστοιχα.

Με τις ακόλουθες εντολές στην R κάνουμε τους απαραίτητους υπολογισμούς :

```
> grade_differences <- math - prob
> n_gd <- length(grade_differences)
> n_gd
[1] 107
> mean_gd <- mean(grade_differences)
> mean_gd
[1] 0.2056075
> sd_gd <- sd(grade_differences)
> sd_gd
[1] 2.228623
```

Στη συνέχεια υπολογίζουμε το στατιστικό ελέγχου με τον ακόλουθο τύπο και τέλος βρίσκουμε το pvalue ίσο με 0.3420928, τιμή αρκετά μεγάλη ώστε να μην απορρίψουμε την μηδενική υπόθεση. Άρα, οι βαθμοί στα δύο μαθήματα δεν διαφέρουν και πολύ και η διαφορά δεν είναι στατιστικά σημαντική. Να σημειωθεί ότι για τον υπολογισμό του pvalue χρησιμοποιούμε βαθμό ελευθερίας 106 δηλαδή $n_{gd} - 1$.

$$t = \frac{\bar{x} - \bar{y}}{s / \sqrt{n}}$$

```
> t <- mean_gd / (sd_gd / sqrt(n_gd))
> (1 - pt(t, df = n_gd - 1)) * 2
[1] 0.3420928
```

Από την άλλη με την χρήση της εντολής t.test βγάζουμε το ακόλουθο pvalue, και καταλήγουμε στο ίδιο συμπέρασμα :

```
> t.test(grade_differences)

One Sample t-test

data:  grade_differences
t = 0.95432, df = 106, p-value = 0.3421
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 -0.2215414  0.6327563
sample estimates:
mean of x
0.2056075
```