

Φροντιστήριο #3

Εισαγωγή στο γενικό πρόβλημα για δύο δείγματα

Θα ασχοληθούμε με δεδομένα από δύο ανεξάρτητα τυχαία δείγματα τα οποία μπορούν να σχηματιστούν ανεξάρτητα από καθέναν από δύο πληθυσμούς. Τα δείγματα είναι ανεξάρτητα υπό την έννοια ότι κάθε μέτρηση στο ένα δείγμα είναι ανεξάρτητη κάθε μέτρησης στο δεύτερο δείγμα. Θεωρούμε δύο πληθυσμούς X και Y με συναρτήσεις κατανομών, αντίστοιχα, F_X και F_Y . Έστω ένα δείγμα μεγέθους m από τον πληθυσμό X και ένα άλλο δείγμα μεγέθους n που σχηματίζεται ανεξάρτητα από τον πληθυσμό Y .

Έχουμε X_1, X_2, \dots, X_m

Y_1, Y_2, \dots, Y_n

Συνήθως, η υπόθεση που μας ενδιαφέρει να εξετάσουμε είναι: $H_0: F_Y(x) = F_X(x) \forall x$ δηλαδή αν τα δύο δείγματα σχηματίζονται από τον ίδιο πληθυσμό.

Αν για παράδειγμα, θεωρήσουμε κάποια παραμετρική υπόθεση όπως π.χ. ότι οι πληθυσμοί **κατανέμονται** κανονικά, τότε το t -test για τη σύγκριση των μέσων τιμών ή το F -test για την ισότητα των διακυμάνσεων

είναι οι πιο ισχυροί έλεγχχοι. Όμως αυτά και άλλα κλασικά tests μπορεί να είναι αρκετά ενάσθητα σε απουσίες από την υπόθεση της κανονικότητας ή από άλλες σημαντικές υποθέσεις. Αν υπάρχουν υποψίες για πιθανές απουσίες ή μη-επαρκής πληροφορία για να ελεγχθούν οι προϋποθέσεις της κανονικότητας μπορεί να υιοθετηθεί μια μη-παραμετρική διαδικασία έλεγχου.

Πρακτικά, χρησιμοποιούνται και υιοθετώνται κάποιες άλλες υποθέσεις για τη μορφή των υπό-εξέταση πληθυσμών. Μία πρώτη τέτοια υπόθεση καλείται μοντέλο μετατόπισης (shift model or location model). Κάτω από μία τέτοια υπόθεση, θεωρούμε ότι οι πληθυσμοί X και Y είναι ίδιοι, εκτός πιθανώς από μία μετατόπιση κατά μία ποσότητα θ , πιθανώς άγνωστη, δηλ. $F_Y(x) = P(Y \leq x) = P(X \leq x - \theta)$

$= F_X(x - \theta) \forall x$ και $\theta \neq 0$. Αυτό σημαίνει ότι $X + \theta$ και Y έχουν την ίδια κατανομή ή ότι η X μετατρέπεται όπως η $Y - \theta$. Ο πληθυσμός Y είναι ο ίδιος με τον X αν $\theta = 0$, "μετατοπίζεται" προς τα δεξιά αν $\theta > 0$ και "μετατοπίζεται" προς τα αριστερά αν $\theta < 0$. Κάτω από αυτήν την υπόθεση, οι πληθυσμοί έχουν το ίδιο σχήμα, την ίδια διασπορά και το ποσό της μετατόπισης θ θα πρέπει να είναι ίσο με τη

διαφορά των μέσων τιμών $\mu_Y - \mu_X$ (ή τη διαφορά $-3-$
των διαμέσων $M_Y - M_X$).

Άλλα υπόθεση που μπορεί να τεθεί για τη μορφή του
πληθυσμού είναι το μοντέλο κλίμακας (scale model).
Εύρφωνα με την υπόθεση αυτή, οι πληθυσμοί X και Y
είναι ίδιοι εκτός (πιθανώς) από μία ποσότητα $\theta, \theta > 0$,

$\theta \neq 1$. Μπορούμε, τότε, να γράψουμε:

$$F_Y(x) = P(Y \leq x) = P(X \leq \theta) = F_X(\theta x) \quad \forall x \text{ και}$$

$\theta > 0, \theta \neq 1$. Αυτό σημαίνει ότι οι $\frac{X}{\theta}$ και Y έχουν
την ίδια κατανομή $\forall \theta > 0$ ή ότι η X κατανέμεται ως
 θY . Επίσης $\text{var}(X) = \theta^2 \text{var}(Y)$ και $E(X) = \theta E(Y)$.

Η υπόθεση του συνθετού μοντέλου (μετατόπισης-κλί-
μακας (location-scale model) ενσωματώνει ιδιότητες
και από τα δύο μοντέλα.

Ο έλεγχος Kolmogorov-Smirnov για δύο δείγματα

**Στην περίπτωση των δύο δειγμάτων, χρησι-
μοποιούμε ένα στατιστικό έλεγχο το οποίο για τα
δύο δείγματα συγκρίνει τις εμπειρικές συναρτήσεις
κατανομών τους.**

Έστω οι διατεταγμένες στατιστικές μεγέθους m
και n αντίστοιχα, από δύο πληθυσμούς με σωχικές
συναρτήσεις κατανομών F_X και F_Y που δίνονται
όπως παρακάτω: $X_{(1)}, X_{(2)}, \dots, X_{(m)}$ και

$Y_{(1)}, Y_{(2)}, \dots, Y_{(n)}$. Οι αντίστοιχες εμπειρικές συναρτήσεις -4-
 σεις κατανομών $F_m(x)$ και $F_n(x)$ ορίζονται ως εξής:

$$F_m(x) = \begin{cases} 0, & \text{αν } x < X_{(1)} \\ k/m, & \text{αν } X_{(k)} \leq x < X_{(k+1)}, \text{ για } k=1, 2, \dots, m-1 \\ 1, & \text{αν } x \geq X_{(m)} \end{cases}$$

και

$$F_n(x) = \begin{cases} 0, & \text{αν } x < Y_{(1)} \\ k/n, & \text{αν } Y_{(k)} \leq x < Y_{(k+1)}, \text{ για } k=1, 2, \dots, n-1 \\ 1, & \text{αν } x \geq Y_{(n)} \end{cases}$$

Μπορούμε να πούμε ότι $F_m(x)$ και $F_n(x)$ είναι τα ποσοστά των παρατηρήσεων των X και Y που δεν ξεπερνούν την τιμή x .

Αν ισχύει η $H_0: F_Y(x) = F_X(x) \forall x$ τότε οι πληθυσμιακές κατανομές ταυτίζονται και κατ' ουσίαν θεωρούμε ότι έχουμε δύο δείγματα από τον "ίδιο" πληθυσμό. Πολύ γρήγορα οι εμπειρικές συναρτήσεις κατανομών των X και Y μπορούν να θεωρηθούν οι κατάλληλοι εμπειρικές για τις πληθυσμιακές συναρτήσεις κατανομών των X και Y .

Χρειαζόμαστε έναν ορισμό εγγύτητας (closeness definition) για να ορίσουμε το κριτήριο Kolmogorov-Smirnov για τα δύο δείγματα. Το κριτήριο συμβολίζεται ως $D_{m,n}$ και είναι η μέγιστη απόλυτη διαφορά μεταξύ των δύο εμπειρικών συναρτήσεων κατανομών:

$$D_{m,n} = \max_x |F_m(x) - F_n(x)|$$

Η παραπάνω σ.σ. ελέγχου είναι κατάλληλη για το

Εναλλακτικό δίπλευρο έλεγχος:

-5-

$$H_A: F_Y(x) \neq F_X(x) \text{ για κάποιο } x \text{ (Έλεγχος Α)}$$

και η κρίσιμη περιοχή καθορίζεται από την ανισότητα:

$$D_{m,n} \geq c_\alpha \text{ όπου } P(D_{m,n} \geq c_\alpha | H_0) \leq \alpha$$

Η p -τιμή είναι $P(D_{m,n} \geq D_0 | H_0)$ όπου D_0 είναι η

παρατηρούμενη τιμή του $D_{m,n}$ στατιστικού που υπολογίζουμε. Η αριθμής κατανομή του στατιστικού

$D_{m,n}$ κάτω από την H_0 δίνεται από κατάλληλους πίνακες (βλέπε π.χ. Πίνακα 21, βιβλίο Ξενοδόχου).

Για τον μονόπλευρο εναλλακτικό υπόθεση:

$$H_1: F_X(x) > F_Y(x) \text{ για τουλάχιστον ένα } x \text{ (Έλεγχος Β)}$$

Χρησιμοποιούμε το στατιστικό:

$$D_{m,n}^+ = \max_x \left[\sum_m(x) - \sum_n(x) \right]$$

ομοίως για τον μονόπλευρο εναλλακτικό υπόθεση:

$$H_1: F_X(x) < F_Y(x) \text{ για τουλάχιστον ένα } x. \text{ (Έλεγχος Γ)}$$

Χρησιμοποιούμε το στατιστικό:

$$D_{m,n}^- = \max_x \left[\sum_n(x) - \sum_m(x) \right]$$

Σχόλια Το test $K-S$ για δύο δείγματα είναι εύκολο

στην εφαρμογή, χρησιμοποιώντας την αριθμής κατανομή για κάθε τιμή των m και n στο εύρος των τιμών των πινάκων και χρησιμοποιώντας κατάλληλη ασυμπτωτική κατανομή για μεγάλα-

τερα δειγματικά μετέθνη (βλέπε πίνακα 21, βιβλίο⁻⁶⁻
Ξενοδόχου). Ο πίνακας 20 του βιβλίου της Ξενοδόχου,
χρησιμοποιείται για την περίπτωση $m=n$.

Ο έλεγχος Β είναι κατάλληλος για προβλήματα όπου έχει
έννοια να ελέγχουμε αν τα Χ είναι μετατοπισμένα προς τα
αριστερά των Υ (δηλαδή αν τα Χ τείνουν να είναι μικρό-
τερα από τα Υ). Αντίστοιχα, ο έλεγχος Γ είναι κατάλληλος
είναι κατάλληλος για τα προβλήματα στα οποία επιθυ-
μούμε να ελέγχουμε αν τα Χ είναι μετατοπισμένα προς
τα δεξιά των Υ (δηλαδή αν τα Χ τείνουν να είναι μετα-
τότερα από τα Υ).

Ακριβής κατανομή για την ελεγχοσύμβαση $D_{m,n}$

Κάτω από την H_0 , η ακριβής κατανομή της ελεγχο-
συνάρτησης $D_{m,n}$ έχει υπολογιστεί. Πιο συγκεκριμένα,
για μικρά μετέθνη δείγματα (π.χ. μικρότερα από 25
ή αν $m+n \leq 16$ ή/και $2 \leq m \leq n \leq 12$) υπάρχουν κατά-
λληλοι πίνακες τερών. Το "σώμα" αυτών των πινάκων
δίνει τις τερές $m, n, D_{m,n}$ που είναι σημαντικές σε διάφο-
ρα επίπεδα στατιστικής σημαντικότητας. Γνωρίζοντας
τις τερές των $m, n, D_{m,n}$ και αναλόγως με το αν
έχουμε μονόπλευρο ή αμφίπλευρο έλεγχο βρίσκουμε
τα κριτικά σημεία του ελέγχου. Για παράδειγμα,
για $m=6, n=8$ απορρίπτουμε την H_0 σε εσο $\alpha=1\%$
όταν $m, n, D_{m,n} \geq 38$, στον μονόπλευρο έλεγχο.
Μεγάλα δείγματα: Αν είτε m είτε n είναι μεγαλύτερα
από 25, τότε π.χ. για τον αμφίπλευρο έλεγχο, πάει

υπάρχουν κατάλληλοι πίνακες. Υπολογίζουμε από τα δεδομένα μας την τιμή του στατιστικού $D_{m,n}$ και συγκρίνουμε την τιμή αυτή με την κρίσιμη τιμή από τον πίνακα για δοθείσες τιμές των m και n και για συγκεκριμένη τιμή του επιπέδου στατιστικής σημαντικότητας α .

Αν $D_{m,n} > C_{\alpha} \sqrt{\frac{m+n}{mn}}$ τότε η H_0 απορρίπτεται σε ε.σ. α για τον δίπλευρο έλεγχο. Για παράδειγμα, για $m=55, n=60$ και για $\alpha=0.05$ για τον αμφίπλευρο έλεγχο, ο πίνακας δίνει την τιμή με την οποία η $D_{m,n}$ θα πρέπει να είναι ίση ή μεγαλύτερη ώστε να απορριφθεί η H_0 . Για αυτά τα δεδομένα, η H_0 απορρίπτεται σε ε.σ. $\alpha=5\%$ αν

$$D_{m,n} > 1.36 \sqrt{\frac{m+n}{mn}}$$

$$= 1.36 \sqrt{\frac{55+60}{(55)(60)}} = 0.254$$

Για τον μονόπλευρο έλεγχο, όταν m και n είναι μεγάλα, τότε $D_{m,n} = \max_x [S_m(x) - S_n(x)]$

Ο Goodman (1954) απέδειξε ότι η σ.σ.

$$\chi^2 = 4 D_{m,n}^2 \frac{mn}{m+n}$$

την κατανομή χ^2 με $df=2$ αν τα δειγματικά μεγέθη αυξάνουν ($n, m \rightarrow \infty$). Σε πλα τέτοια περίπτωση χρησιμοποιούμε κατάλληλους πίνακες της κατανομής χ^2 . Η προσέγγιση με την χ^2 μπορεί να χρησιμοποιηθεί και με μικρά δείγματα αλλά αυτό συνήθως οδηγεί σε έναν συντηρητικό έλεγχο. Δηλαδή

το λάθος στη χρήση της χ^2 -προσέγγισης με ριζικά -8- δείγματα οδηγεί σχεδόν πάντα τον έλεγχο προς τη σωστή κατεύθυνση (π.χ. σωστή απόρριψη της H_0). Σε κάθε περίπτωση προτιμάται η αυριβής μεταφορά για το στατιστικό $D_{m,n}$ σε ριζικά δείγματα.

Ας δούμε ένα παράδειγμα.

Παράδειγμα Ερευνητής θέλει να ελέγξει αν δύο τάξεις παιδιών συμπεριφέρονται στατιστικά ίδια ως προς μία ψυχολογική μέθοδο που τους εφαρμόζεται. Ο ερευνητής ελέγχει την υπόθεση:

H_0 : Δεν υπάρχει στατιστικά σημαντική διαφορά στα ποσοστά λαθών που κάνουν οι δύο τάξεις παιδιών, έναντι της

H_1 : τα παιδιά της μεγαλύτερης τάξης έχουν ριζότερα (λιγότερα) ποσοστά λαθών από τα παιδιά της ριζότερης τάξης.

Δίνεται ο ακόλουθος πίνακας με τα ποσοστά των σωστικών λαθών στις δύο τάξεις.

Τάξη με παιδιά 11 ετών	Τάξη με παιδιά 7 ετών
35.2	39.1
39.2	41.2
40.9	45.2
38.1	46.2
34.4	48.4
29.1	48.7
41.8	55.0
24.3	40.6
32.4	52.1
	47.2

Επιλέγουμε $\alpha = 0.01, m = 9, n = 10$

Εφαρμόσουμε το one-sided K-S two sample test.

Τα δεδομένα μας καταχωρούνται σε κλάσεις διατεταγμένα ώστε να υπολογιστούν οι δύο εμπειρικές συνάρτησεις

κατανορών $F_m(x), F_n(x)$.

	$F_m(x)$	$F_n(x)$	$F_m(x) - F_n(x)$
24-27	1/9	0/10	0.111
28-31	2/9	0/10	0.222
32-35	5/9	0/10	0.556
36-39	7/9	4/10	0.678
40-43	9/9	3/10	0.7
44-47	9/9	6/10	0.4
48-51	9/9	8/10	0.2
52-55	9/9	10/10	0
	$F_m(x)$	$F_n(x)$	

Παρατηρούμε ότι η μέγιστη διαφορά παρατηρείται για τιμή $D_{m,n} = 0.7$. Άρα $m \cdot n \cdot D_{m,n} = 9 \cdot 10 \cdot 0.7 = 63$.

Για $\alpha = 0.01$ η κριτική τιμή από τον κατάλληλο πίνακα είναι: $C_{0.01} = 61$. Επειδή $m \cdot n \cdot D_{m,n} > C_{0.01}$

υπάρχουν ισχυρές ενδείξεις απόρριψης της H_0 .

Συμπεραίνουμε ότι τα παιδιά 11 ετών φαίνεται ότι μάθω λιγότερα γάθη από τα παιδιά των 7 ετών όταν και στις δύο ομάδες παιδιών εφαρμοστεί η ψυχολογική μέθοδος.

Για την επιλογή του ελέγχου Kolmogorov-Smirnov για τα

δύο δείγματα χρησιμοποιούμε τη συνάρτηση ks-test, στην R.

Χρήσιμα ορίσματα της συνάρτησης είναι τα ακόλουθα:

alternative = $\begin{cases} "g" \\ "l" \end{cases}$ για τον καθορισμό του μονόπλευρου ελέγχου

και exact = $\begin{cases} TRUE \\ FALSE \end{cases}$ για τον καθορισμό της μεθόδου υπολογισμού της p-τιμής (p-value).

Παρατήρηση: Στον έλεγχο K-S για δύο δείγματα υποθέ-

τούμε ότι οι πληθυσμοί X και Y είναι ίδιοι εντός (ως από μία πιθανή διαφοροποίηση στις θέσεις τους που εκφράζεται από μία μετατόπιση θ (δηλ. ότι X + θ και Y έχουν την ίδια κατανομή ή ότι X και Y - θ έχουν την ίδια κατανομή). Άρα, εναλλακτικά οι έλεγχοι υποθέσεων:

$H_0: F_X(x) = F_Y(x) \quad \forall x$ έναντι:

$H_1: F_X(x) \neq F_Y(x)$
για ένα τουλάχιστον x

$H_1: F_X(x) \geq F_Y(x)^*$ $H_1: F_X(x) \leq F_Y(x)$
ισχύει η γνήσια ανισότητα για ένα τουλάχιστον x

* Αν $F_X(x) \geq F_Y(x)$ με αυστηρή ανισότητα για κάποιο x τότε λέμε ότι η τ.ρ. Y είναι στοχαστικά μεγαλύτερη από την τ.ρ. X. Συμβολισμός: $Y \stackrel{st}{\geq} X$.

Μπορούν να γραφούν ως εξής: (για το μοντέλο θέσης)

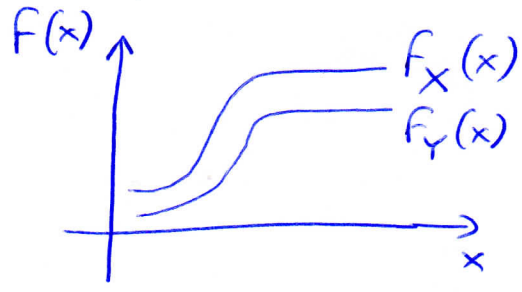
$H_0: \theta = 0$ έναντι:

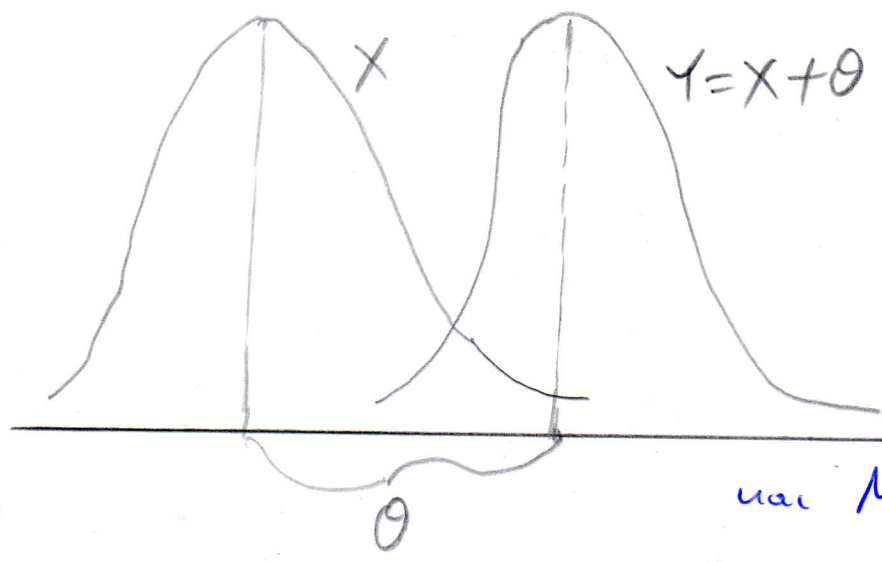
$H_1: \theta \neq 0$

$H_1: \theta > 0$

$H_1: \theta < 0$

Σχήμα για (*)





Οι X και Y έχουν ίδια μορφή
 ίδια διακύμανση
 ενώ η θ συνήθως είναι:
 $\theta = \mu_Y - \mu_X$
 ή $\theta = M_Y - M_X$
 όπου μ_X, μ_Y : μέσες τιμές
 και M_X, M_Y διάρροσι X, Y .

Η p -value του ελέγχου είναι ίση με: (για τον αφίπτερο έλεγχο)

$P(D_{n,m} > d_{n,m})$ όπου $d_{n,m}$ είναι η παρατηρούμενη τιμή της $D_{n,m}$

Είδαμε ότι η ακριβής κατανομή της $D_{n,m}$ είναι δύσκολο να υπολογιστεί και έχουν κατασκευαστεί ειδικοί πίνακες με τα ποσοστιαία σφάλματα της $D_{n,m}$. Μπορεί ωστόσο, να αποδειχθεί ότι:

$$\lim_{n,m \rightarrow \infty} P\left(\sqrt{\frac{nm}{n+m}} D_{n,m} \leq d\right) = 1 - 2 \sum_{i=1}^{\infty} (-1)^{i-1} e^{-2id^2}$$

άρα για προσεγγιστική τιμή της p -value του ελέγχου για μεγάλα δείγματα είναι η:

$$p\text{-value} = P(D_{n,m} > d_{n,m}) = 1 - P(D_{n,m} \leq d_{n,m})$$

$$\approx 2 \sum_{i=1}^{\infty} (-1)^{i-1} e^{-2 \frac{nm}{n+m} i^2 d_{n,m}^2}$$

Στο παράδειγμά μας έχουμε:

$$H_0: F_X(x) = F_Y(x) \quad \forall x$$

vs

$$H_1: F_X(x) \geq F_Y(x) \text{ για ένα τουλάχιστον } x.$$

$$D_{n,m}^+ = \max_x [F_n(x) - F_m(x)]$$

παιδιά μεγαλύτερης τάξης (11 ετών)
παιδιά μικρότερης τάξης (7 ετών)

$$\rho \text{ ε } D_{n,m}^+ = 0.7$$

$$\text{και } p\text{-τιμή} = P(D_{n,m}^+ \geq 0.7) = 0.009637$$

Για τα συνήθη επίπεδα σ.σ. η H_0 απορρίπτεται άρα υπάρχουν ισχυρές ενδείξεις υπέρ της H_1 . Η υλοποίηση του ελέγχου στην R δίνεται στην επόμενη σελίδα μαζί με την υλοποίηση των γραφημάτων για τις ερπειδικές σ.κ. για τα δύο δείγματα, των παιδιών μεγαλύτερης τάξης 11 ετών και των παιδιών μικρότερης τάξης 7 ετών.



Για την κατασκευή των γραφημάτων των δύο
εμπειριικών σ.κ. για τα δύο δείγματα χρησιμοποιούμε
την library (ggplot2).

Θέτουμε $cdf1$ και $cdf2$ τις δύο ε.σ.κ. για τα δύο
δείγματα.

Ο τίτλος του γραφήματος δίνεται ως όρισμα της
εντολής plot, π.χ. $main = "ECDF"$

Με την εντολή lines συνδέουμε τα δύο γραφήματα.

Τα υπόλοιπα γραφήματα καθορίζουν την μορφή και
το χρώμα των γραμμών που θα έχουν τα δύο γραφήματα.
(col.points, lty).

Με την εντολή legend βάζουμε στη θέση (30, 0.8)
την επεξηγηματική "ετικέτα" με τα ονόματα
"sample 1" και "sample 2".



-14-

```
library(ggplot2)
sample1=c(35.2,39.2,40.9,38.1,34.4,29.1,41.8,24.3,32.4)
sample2=c(39.1,41.2,45.2,46.2,48.4,48.7,55,40.6,52.1,47.2)
# create ECDF of data
cdf1 = ecdf(sample1)
cdf2 = ecdf(sample2)
plot(cdf1, verticals=TRUE, do.points=FALSE, col.points=2,
col.hor=2, col.vert=2, lty=1, main="ECDFS")
lines(cdf2, verticals=TRUE, do.points=FALSE, col.points=3,
col.hor=3, col.vert=3, lty=2)
leg.names=c("sample1", "sample2")
legend(30, 0.8, leg.names, lty=1:2, col=2:3)
#KS_test_two_samples
ks.test(sample1, sample2, exact=TRUE)
ks.test(sample1, sample2, alternative="g")
```

Two-sample Kolmogorov-Smirnov test

data: sample1 and sample2
D = 0.7, p-value = 0.007036
alternative hypothesis: two-sided

```
> ks.test(sample1, sample2, alternative="g")
```

Two-sample Kolmogorov-Smirnov test

data: sample1 and sample2
D⁺ = 0.7, p-value = 0.009637
alternative hypothesis: the CDF of x lies above that of y

-25-

ECDFS

