

Proiect Probabilitate și Statistică

Vîlculescu Mihai-Bogdan - Grupa 232

Martac Dana-Maria - Grupa 232

Drăguțescu Mihai-Valentin - Grupa 232

Set de date = randu

0. RANDU

Set-ul de date utilizat de echipa noastră a fost randu. Acest set de date este alcătuit din 400 de triplete de numere aleatoare succesive. Acestea sunt memorate ca un data frame cu 400 de observații pe 3 variabile: x, y și z.

1. Operații de statistică descriptivă

Pentru primul exercițiu, se cerea să efectuăm diverse operații de statistică descriptivă, precum: medie, varianță, quantile, boxplot.

- **Media:** este valoarea medie a unui set de date. În R, am calculat media pentru fiecare din cele 3 variabile (x, y, z). Acest lucru se realizează utilizând funcția "**mean**". Aceasta primește ca parametru un set de date, și calculează media acelor date.
- **Mediana:** este valoarea din mijloc a unui set de date. Am procedat la fel ca la medie, și anume am calculat mediana pentru x, y și z. Funcția utilizată se numește "**median**" și se comportă la fel ca funcția pentru calcul a mediei.
- **Deviația standard:** reprezintă deviația medie a valorilor dintr-un set de date față de valoarea medie. Aceasta se calculează în limbajul R, folosind o funcție numită "**sd**".
- **Varianța / Dispersia:** reprezintă pătratul deviației standard. Funcția folosită se numește "**var**".
- **Quartile:** reprezintă o modalitate de a împărți un sample de date în mai multe subgrupuri de dimensiuni egale și adiacente. Funcția folosită se cheamă intuitiv "**quantile**". Aceasta nu întoarce nimic, ci afișează în consolă distribuția setului de date în aceste quartile.

- **Supply**: o funcție din R care permite afișarea în consolă a valorilor unui set de date în momentul aplicării unei funcții asupra sa. Spre exemplu, se poate apela astfel: "**supply(my_data, mean)**", iar acest apel calculează media pentru toate variabilele set-ului my_data. Astfel, putem să vedem valorile mediei fără să apelăm manual funcția pentru fiecare variabilă.
- **Summary**: o funcție care afișează detaliile statistice cele mai importante ale unei anumite variabile dintr-un set de date furnizat ca parametru. (quartilele, media, mediana, min, max).
- **Boxplot**: este un tip de grafic folosit pentru a afișa modele de date cantitative. Ea oferă informații privind tendința centrală și forma distribuției studiate. Pe această diagramă, sunt observabile următoarele elemente statistice importante: minimul, prima quartilă, mediana, a treia quartilă și maximul.
- **Interpretare**: după apelul tuturor funcțiilor, putem încerca interpretarea valorilor obținute.
 - **Mean** = 0.526 pentru x, 0.486 pentru y și 0.481 pentru z. Deci putem afirma că mediile se apropie de 0.5 pentru toate cele 3 variabile.
 - **Median** = 0.540 (x), 0.483 (y), 0.463 (z). Astfel, medianele sunt foarte apropiate de valorile medii, la nu mai mult de 0.02 distanță.
 - **Standard Deviation** = 0.285 (x), 0.293 (y), 0.279 (z). Deviația medie a valorilor față de medie este destul de mică, dar ținând cont că valorile din setul de date randu sunt cuprinse între 0 și 1, această deviație demonstrează că valorile sunt destul de răspândite.
 - **Variance** = 0.081 (x), 0.086 (y), 0.077 (z). Varianța este mică dacă privim obiectiv, dar relativ la setul nostru de date, această varianță confirmă ce am spus la deviația standard (cum era și firesc), și anume că valorile nu sunt tocmai apropiate una de cealaltă.
 - **Quantiles** = quartilele pentru cele 3 variabile, împart valorile în subgrupuri. Astfel, pentru x, 25% din valori sunt mai mici decât 0.3003. Pentru y, 75% din valori sunt mai mici decât 0.7399195. Quartilele ne ajută să observăm anumite distribuții ale valorilor.
 - **Boxplot** = diagrama indică într-un mod grafic ceea ce putem citi din valorile regăsite în quartile. Mai exact, putem observa cum sunt distribuite valorile față de mediană și față de minim și maxim. Aici,

valorile lui x sunt mai mari decât y și z, inclusiv mediana lui x are o valoare mai mare. În schimb, y și z sunt mai apropiate ca minim, maxim și mediană, dar valorile din z sunt mai puțin răspândite, lucru observabil și în rezultatul deviației standard și al varianței.

2. Regresia liniară

Pentru al doilea exercițiu, trebuie să creăm 2 modele de regresie liniară pentru setul nostru de date (unul de regresie liniară și unul de regresie multiplă). Pașii efectuați pentru îndeplinirea cerinței sunt:

A. Verificare

- a. Pentru a putea crea modelele de regresie, trebuie, în primul rând, să ne asigurăm că setul nostru de date este unul care se pretează asupra regresiei liniare.
- b. Apelăm funcția **cor** având ca parametri variabilele x și y, pe care intenționăm să le folosim ca predicator și răspuns. Această funcție testează gradul de corelație dintre aceste 2 variabile. Mai exact, în funcție de valoarea întoarsă de cor, ne putem da seama dacă cât de dependente sunt cele două variabile una față de cealaltă.
- c. Funcția în cazul nostru întoarce **-0.048**. Această valoare este mai mică decât 0, fapt care îmi sugerează că variația variabilei răspuns nu poate fi susținută de variabila predictor. Astfel, ar trebui să alegem variabile mai dependente una de cealaltă. Dar, dacă apelăm funcția pentru toate combinațiile celor 3 variabile, obținem mereu o valoare scăzută și apropiată de 0, ceea ce implică faptul că nu vom putea crea un model de regresie liniară care să funcționeze corespunzător cu aceste variabile.
- d. Am construit modelul liniar pentru întregul set de date apelând funcția **lm**. Apoi, am aplicat **summary** pentru modelul creat. Apoi am extras din acel summary, coeficienții pentru a putea calcula beta estimat și eroarea standard. Cu aceste valori, am calculat valorile p și t. Cunoaștem că, pentru o valoare mai mare a lui t, setul de date este mai semnificativ din punct de vedere statistic, dar noi am obținut valoarea **-0.968**. În continuare, valoarea lui p este **0.3335**, iar noi ar fi trebuit să obținem un p mai apropiat de 0.05. Astfel, rezultatele de la acest punct confirmă concluzia de la punctul c.

B. Construirea modelului

- a. Pentru ambele tipuri de regresie efectuăm în mare parte aceleași operații. Prima dată, alegem seturile de date de training și de test. Acestea vor fi alese random din datele inițiale folosind funcția **sample**. Am ales să împărțim datele inițiale în raport 80:20 (training:test).
- b. Creăm modelul liniar simplu, respectiv multiplu. Pentru cel simplu, specificăm întâi variabila răspuns apoi predictorul (**response ~ predictor**). Pentru cel multiplu, adăugăm separați prin +, mai multe variabile de tip predictor (**response ~ predictor1 + predictor2**). Apoi, facem predicțiile apelând funcția **predict**, care primește ca parametri modelul realizat în raport cu setul de date training și al doilea parametru este setul de date test.
- c. Următoarele instrucțiuni reprezintă verificări ale predicțiilor făcute, dar după cum am precizat anterior, acestea vor fi inexacte.

3. Repartiția Beta

Pentru al treilea exercițiu, am selectat repartiția beta, pentru care trebuia să afișăm în 2 reprezentări alăturate funcția de densitate și cea de repartiție. Astfel, am ales un șir de numere de la -5 la 5, cu pasul de 0.001 (-5, -4.999, ...). Pentru acel șir am apelat funcția **plot** cu șirul ca prim parametru și funcția corespunzătoare ca al doilea parametru. (**dbeta, pbeta**). Acesta a fost procesul de construire a reprezentărilor.

Interpretări

- ☐ Un lucru care se observă imediat la cele 2 grafice este acela că funcția de densitate nu este continuă, în timp ce cea de repartiție este. Funcția de densitate ia valori foarte mici până într-un anumit punct, unde crește brusc și începe din nou să scadă.
- ☐ Funcția de repartiție crește mai puțin brusc și după nu mai scade, rămânând constantă.

Aplicații: Rule of Succession, Deducția Bayesiană, Order Statistics, Logica subiectivă.