

# Statistika

Bor Bregant

8. junij 2024



# Kazalo

|     |  |    |
|-----|--|----|
| 0.1 | Uvod . . . . .                               | 1  |
| 0.2 | Opisna statistika in vizualizacija . . . . . | 6  |
| 0.3 | Inferenčna statistika . . . . .              | 8  |
| 0.4 | Bivariantna statistika . . . . .             | 12 |



## 0.1 Uvod

Statistika je veda, ki se ukvarja z zbiranjem, analiziranjem, interpretacijo, predstavitvijo in organizacijo podatkov. Izhaja iz *statisticum* (državni), saj je prvotno označevala analizo podatkov o državi. Njene uporabe presegajo matematiko in je temeljno orodje za raziskovanje na vseh področjih znanosti. Pomembno je, da jo uporabljamo odgovorno in etično ter kritično ocenjujemo kontekst in vir statističnih informacij, da se izognemo napačnim interpretacijam in zavajanju (zgled zavajanja).

### Podatki in informacije

**Podatki** so niz vrednosti kvalitativnih spremenljivk.

**Informacije** so obdelani in interpretirani podatki.

Primer: temperatura po vsem svetu za zadnjih 100 let so podatki; analiza, ki ugotavlja, da globalna temperatura narašča, je informacija

### Enote opazovanja in Spremenljivke

Enote opazovanja, za katere se zbirajo podatki, npr. oseba, gospodinjstvo, podjetje, izdelek, stavba, dogodek, država ...

Spremenljivka je katero koli značilno število ali količina, ki jo je mogoče izmeriti ali prešteti, na primer:

- Spol, starost, izobrazba, poklic, dohodek, narodnost ... (za osebo)
- Število članov, vrsta, lastništvo stanovanja, dohodek, internetna povezava ... (za gospodinjstvo)
- Število zaposlenih, lokacija, vrsta, sektor, prihodki ... (za podjetje)
- Dimenzije, teža, starost, barva, temperatura, cena ... (za izdelek)
- Dimenzije, lokacija, starost, lastništvo, materiali, cena ... (za stavbo)
- Površina, prebivalstvo, število podjetij, politična ureditev ... (za državo)

### Tipi spremenljivk glede na vrednost in glede na mersko lestvico

#### Glede na vrednost

**Kategorične** (atributne) spremenljivke, npr. spol, izobrazba, barva, sektor, vrsta, regija

**Številske** spremenljivke:

- Zvezne (lahko imajo poljubne vrednosti), npr. točna starost, dohodek, cena, dimenzije, dolžina, širina, trajanje ...
- Diskretne (imajo le celoštevilске vrednosti), npr. letnica rojstva, velikost gospodinjstva, velikost podjetja, število udeležencev ...

#### Glede na mersko lestvico

- **Nominalne** spremenljivke: vrednosti se lahko razlikujejo le med seboj, razvrščanje ni možno, npr. spol, poklic, sektor, narodnost, regija ...
- **Ordinalne** spremenljivke: vrednosti so lahko razvrščene od najmanjše do največje, vendar razdalje med vrednostmi niso znane, npr. izobrazba, šolska ocena, stopnja strinjanja, stopnja zadovoljstva, tesnoba ...

- **Intervalne** spremenljivke: razlika med dvema vrednostma je smiselna, vendar ni dejanske ničelne vrednosti, je samo arbitrarna, npr. temperatura na lestvici Celzija, pH, koledarsko leto ...
- **Razmernostne** spremenljivke: imajo edinstveno in nearbitrarno ničelno vrednost, zato lahko izračunamo tudi razmerja, npr. temperatura po Kelvinovi lestvici, starost, dolžina, širina, višina, teža, velikost razreda, število udeležencev dogodka, dohodek ...

| Kaj lahko izračunamo     | Nominalna | Ordinalna | Intervalna | Razmernostna |
|--------------------------|-----------|-----------|------------|--------------|
| Frekvenčna porazdelitev  | ✓         | ✓         | ✓          | ✓            |
| Modus                    | ✓         | ✓         | ✓          | ✓            |
| Vrsti red vrednosti      |           | ✓         | ✓          | ✓            |
| Mediana                  |           | ✓         | ✓          | ✓            |
| Povprečje                |           |           | ✓          | ✓            |
| Razlika med vrednostmi   |           |           | ✓          | ✓            |
| Seštevanje in odštevanje |           |           | ✓          | ✓            |
| Množenje in deljenje     |           |           |            | ✓            |

### Python Example:

```

1      # Import necessary libraries
2      import pandas as pd
3
4      # Step 1: Read the CSV file
5      # Assume the CSV file is named 'data.csv' and located in the same
6      directory as the script
7      df = pd.read_csv('data.csv')
8
9      # Step 2: Print the first few rows of the dataframe
10     print("First few rows of the dataframe:")
11     print(df.head())
12
13     # Step 3: Check the data types of each column
14     print("\nData types of each column:")
15     print(df.dtypes)
16
17     # Step 4: Convert a specific column to float (if necessary)
18     # Let's assume we have a column named 'income' which we want to convert
19     to float
20     df['income'] = df['income'].astype(float)
21
22     # Verify the conversion
23     print("\nData types after conversion:")
24     print(df.dtypes)
25
26     # Step 5: Calculate the range of data in a numeric column
27     # Let's calculate the range for the 'income' column
28     income_range = df['income'].max() - df['income'].min()
29     print("\nRange of the 'income' column:")
30     print(income_range)
31
32     # Step 6: Label encoding for an ordinal categorical variable
33     # Let's assume we have an ordinal variable named 'education_level'
34     education_levels = {'High School': 1, 'Bachelor': 2, 'Master': 3, 'PhD':
35 : 4}
36     df['education_level'] = df['education_level'].map(education_levels)
37
38     # Verify the encoding
39     print("\nData after label encoding 'education_level':")

```

```
37     print(df.head())
38
39     # Step 7: One-hot encoding for a nominal categorical variable
40     # Let's assume we have a nominal variable named 'region'
41     df = pd.get_dummies(df, columns=['region'], prefix='region')
42
43     # Verify the one-hot encoding
44     print("\nData after one-hot encoding 'region':")
45     print(df.head())
46
47     # Additional Step: Descriptive statistics summary
48     print("\nDescriptive statistics of the dataframe:")
49     print(df.describe())
50
51     # Save the modified dataframe to a new CSV file
52     df.to_csv('modified_data.csv', index=False)
53
54     print("\nModified dataframe saved to 'modified_data.csv'.")
```

## Populacija in vzorec

**Populacija** se nanaša na skupni niz opazovanj; pomembno jo je prostorsko in časovno opredeliti, npr.

- študenti Univerze na Primorskem v študijskem letu 2023/2024
- javni vrtci v Obalno-Kraški regiji na 1. 9. 2023
- gledališke predstave v Sloveniji v tednu od 19. do 25. 2. 2024
- knjige izdane v EU v januarju 2024

**Vzorec** se nanaša na niz podatkov, izbranih iz statistične populacije po določenem postopku, npr.

- sistematični vzorec 400 študentov UP
- naključni vzorec 200 knjig

## Vrste statistične analize

Glede na namen:

- Opisna (deduktivna) statistika: analiza in opis zbranih podatkov brez težnje po posploševanju teh podatkov izven njihovega obsega
- Inferenčna (induktivna) statistika: sklepanje iz vzorca na populacijo

Glede na število sočasno analiziranih spremenljivk:

- Univariatna statistika: analiza ene spremenljivke
- Bivariatna statistika: analiza dveh spremenljivk, npr. hi-kvadrat, mere povezanosti, t-test, ANOVA, regresija, ...
- Multivariatna statistika: analiza več spremenljivk, npr. multipla regresija, analiza glavnih komponent, faktorska analiza, diskriminantna analiza ...

## Koraki statistične analize

- i Določitev vsebine in namena statistične študije, opredelitev objekta (enota in populacija) in vsebino opazovanja (spremenljivke)
- ii Statistično opazovanje (celotne populacije ali vzorca)
- iii Enostavna obdelava (urejanje, soritrnanje podatkov in izračun osnovnih karakteristik)
- iv Analitična obravnava

## Vaje

### Vaja 1

Za naslednje spremenljivke definirajte nekaj možnih vrednosti in navedite, ali so zvezne ali diskretne, ter kakšna je raven merjenja:

- Število dnevnih poslov na ljubljanski borzi
- Temperatura v Kopru v stopinjah Celzija
- Življenjska doba osebnega računalnika
- Število dni letnega dopusta za zaposlene
- Dnevno prehojena razdalja:
  - a. kilometrih
  - b. korakih
- Leto neto dohodek učitelja
- Teža solate v gramih
- Število polic v knjižni omari
- Strinjanje s trditvijo na lestvici od 1 (Sploh se ne strinjam) do 5 (Povsem se strinjam)

### Vaja 2

Za naslednje enote navedite nekaj primerov spremenljivk in določite njihovo mersko lestvico:

- Učenec
- Učitelj
- Razred
- Šola
- Učbenik

### Vaja 3

Predstavljajte si, da proučujete pojav besede “trajnostni razvoj” v slovenskih srednješolskih učbenikih izdanih v obdobju od 2010 do 2020. Med njimi naključno izberete 250 učbenikov, v katerih preštejete, kolikokrat se pojavi beseda trajnostni razvoj.



- Kaj je enota analize in kaj je spremenljivka?
- Kaj je merska lestvica spremenljivke?
- Kakšen je vzorec in kako velik je?
- Kakšna je populacija?

#### **Vaja 4**

Recimo, da je imel Pokrajinski muzej v Kopru v letu 2023 natanko 20.000 obiskovalcev. Predstavljate si, da je vsak deseti obiskovalec muzeja prejel in izpolnil kratek vprašalnik:

- Kako velika je populacija?
- Kako velik je vzorec?
- Napiši vsaj štiri vprašanja, na podlagi katerih bi lahko definirali eno nominalno, eno ordinalno, eno intervalno in eno razmernostno spremenljivko.

#### **Vaja 5**

Razišči bazo podatkov raziskave PISA na spletni strani OECD

## 0.2 Opisna statistika in vizualizacija

### Absolutna, relativna in grupirana frekvenčna porazdelitev

### Grafična predstavitev frekvenčnih porazdelitev za nominalne in ordinalne spremenljivke

Stolpčni grafikon > tortni grafikon

### Grafična predstavitev frekvenčnih porazdelitev za intervalne in razmernostne spremenljivke

Histogram

Poligon

Ogiva (kumulativne frekvence)

### Normalna porazdelitev

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$

Slika

Asimetričnost (skewness)

V levo ( $k < 0$ ), v desno ( $k > 0$ ), če  $k$  med -1 in 1 še vedno normalna. Slika

Sploščenost (kurtosis)

Lepto ( $k < 0$ ), mezo ( $k = 0$ ), plati ( $k > 0$ ), slika, če med -1,1 normalna

### Rangiranje

Primer

### Kvantili

Primer in imena

### Mere centralne tendence

- **Modus** - vrednost, ki se najpogosteje pojavi v nizu vrednosti podatkov
- **Mediana** - vrednost, ki ločuje zgornjo polovico obsega razpona vrednosti od spodnje polovice
- **Aritmetična sredina** - povprečje niza vrednosti
- Druge mere (geometrijska, harmonična sredina, ...)

Primerjava med modusom, mediano in aritmetično za unimodalne asimetrične - slika

### Mere variabilnosti (disperzije)

V kolikšni meri se vrednosti razlikujejo med seboj ter razlikujejo in odstopajo od povprečja. Delimo jih na:

- Absolutne mere (razpon, interkvartilni rang, absolutna deviacija aritmetične sredine/mediane, varianca in standardni odklon)

- Relativne mere (absolutne mere deljene s pripadajočo mero centralne tendence) se izračunajo samo za razmernostne spremenljivke; uporabljamo jih, ko želimo primerjati:
  - Dve porazdelitvi z zelo različno vrednostjo za nek mero centralne tendence;
  - Dve spremenljivki z različnima merskima enotama.

**Theory:**

Data visualization is the graphic representation of data. It involves producing images that communicate relationships among the represented data to viewers of the images. Visualization makes it easier to detect patterns, trends, and outliers in groups of data.

**Python Example:**

```
import seaborn as sns
import pandas as pd

# Load an example dataset
data = sns.load_dataset("iris")

# Create a pairplot
sns.pairplot(data, hue="species")
plt.show()
```

## 0.3 Inferenčna statistika

Metode sklepanja iz vzorca na populacijo. Uporabljamo teorijo verjetnosti, da ocenimo, koliko lahko zaupamo rezultatom, pridobljenim na verjetnostnem vzorcu.

### Vzorčenje

Je postopek izbire dela populacije, ki jo vključimo v raziskavo. Enote, ki so že izbrane pogosto ne vračamo v populacijo, če pa je velikost vzorca majhna, lahko to naredimo.

### Metode vzorčenja

- **Verjetnostni vzorci:** Vsaka enota v populaciji ima znano neničelno verjetnost, da bo vključena v vzorec.
  - *Enostavno slučajno vzorčenje:* Vsaka enota ima enako in znano verjetnost izbire, ki ni enaka nič, notej so vsi možni vzorci enako verjetni. Primer: Z računalnikom naključno ustvarimo vzorec 100 dijakov, vpisanih na neko šolo v šolskem leti 2025/26 na podlagi seznama dijakov.
  - *Sistematično vzorčenje:* Iz vzorčnega okvirja vzamemo vsako  $k$ -to enoto. Vsaka enota ima enako verjetnost, da je izbrana v populacijo, toda vsi vzorci niso enako verjetni (npr. ne moremo hkrati izbrati četrte in pete enote, torej vzorčenje ni enostavno). Primer: Na podlagi seznama dijakov šole, ki je urejen po abecedi izberemo vsako deseto enoto. Naključno izberemo le prvo enoto npr. 2 in nadaljujemo z 12, 22, ....
  - *Stratificirano vzorčenje:* Populacijo stratificiramo na podlagi vnaprej znanih informacij in nato izvedemo vzorčenje za vsak stratum posebej. Primer: Če ima šola 70% dijakov in 30% dijakinj, lahko vzamemo v vzorec enote proporcionalno glede na spol.
  - *Vzorčenje v skupinah:* Enote v populaciji so pogosto združene v skupinah, npr. učenci v razrede, razredi v šole, šole v države, itd. Tako lahko najprej izberemo vzorec skupin (npr. razredov) in naprej na temu vzorcu vzorčimo naprej. Primer: Na univerzi izberemo 5 od 15 programov in za vsakega od teh slučajno izberemo vzorec 100 študentov.
- **Neverjetnostni vzorci:** Verjetnosti izbir ne moremo izračunati.
  - *Priložnostni vzorci:* Primer: Državljanom pošljemo 1 milijon vprašalnikov. (slabo)
  - Ekspertna izbira
  - *Kvotno vzorčenje:* Primer: Med vzorčenjem nadziramo demografske značilnosti vzorca.

Nauk: Velikost vzorca ni vse. Pomembnejša je njegova reprezentativnost. V idealni situaciji bi uporabljali verjetnostno vzorčenje. Ko to ni možno, je kvotno vzorčenje boljša izbira kot priložnostno.

Natančnost in točnost s sliko. (precision, accuracy)

### Velikost vzorca

Standardna napaka statistike je odvisna od velikosti vzorca.

Velikost vzorca za ocenjevanje aritmetične sredine navadno zahtevamo:

$$n > \left( \frac{z_{\frac{\alpha}{2}} \cdot \sigma}{E} \right)^2,$$

kjer  $E$  označuje smiselno razliko med opazovanimi vrednostmi (navadno 1)  
nekeje nekeje statistika

## Intervali zaupanja

Parametre lahko ocenimo točkovno ali z intervalom. Intervali zaupanja kažejo na točnost ocene in podajo informacijo o njeni zanesljivosti.

S tveganjem  $\alpha$  lahko rečemo, da interval  $(a, b)$  vsebuje parameter  $\gamma$ .

Slika normalne.

Širina intervala zaupanja je odvisna od:

- *Stopnje tveganja*: Višja je stopnja tveganja  $\alpha$ , ožji je interval
- *Velikosti vzorca*: Večja je velikost vzorca  $n$ , ožji je interval.

## Testiranje hipotez

Znanstvena metoda:

1. Opazovanje pojava
2. Postavljanje raziskovalnih vprašanj
3. Oblikovanje hipotez
4. Zbiranje podatkov
5. Sprejemanje hipotez in razvoj teorij in zakonitosti

Primer: Raziskovalno vprašanje: Ali obstajajo razlike med spoloma v matematični anksioznosti?

Hipoteza: Ženske imajo več matematične anksionosti kot moški ( $H_0 : \mu_z - \mu_m > 0$ ).

Opomba: Hipoteza je vedno trditev, ki jo lahko poskušamo zavrniti. Hipoteza ni potrjena in tehnično ni pravilno reči, da je hipoteza sprejeta. Hipotezo lahko ali zavrnemo ali ne zavrnemo.

Ničelno hipotezo  $H_0$  lahko neposredno preverimo in če je zavrnjena, je alternativna hipoteza pravilna (to je običajno naši cilj)

Alternativno hipotezo  $H_1$  preverimo le posredno: če ničelne hipoteze ne moremo zavrniti, ničelne hipoteze ne sprejmemo, ampak sklenemo, da ni dovolj podatkov, da bi rekli, ali je razlika statistično pomembna

Primer: Ničelna hipoteza: Med ženskami in moškimi ni razlike v sposobnosti večopravnosti.

Alternativna hipoteza: Obstajajo razlike v sposobnosti večopravnosti med ženskami in moškimi (dvostranski test). Alternativna hipoteza: Ženske so boljše pri večopravnosti kot moški (enostransko testiranje)

Dve vrsti napake.

## Stopnja značilnosti

Izberemo največjo verjetnost  $\alpha$  stopnjo, do katere smo pripravljeni tvegati napako tipa I. Običajno se odločimo za 5% stopnjo značilnosti ( $\alpha = 0.05$ ), lahko pa je tudi 1% ( $\alpha = 0.01$ ) ali nižja. Na podlagi izbranega  $\alpha$  določimo kritično območje, kjer bo ničelna hipoteza zavrnjena. Pri tem upoštevamo, ali gre za dvostranski ali enostranski test

Slika two tests same probability level one tail, two tail.

## Studentova $t$ porazdelitev

Slika porazdelitve

Razpršenost je odvisna od tako imenovanih prostostnih stopenj (*angl. degrees of freedom, df*), ki so opredeljene kot velikost vzorca minus število parametrov populacije, ki jih je treba oceniti na podlagi vzorca.

Večji kot je vzorec, bližje je  $t$  porazdelitev normalni  $Z$  porazdelitvi

## $T$ test za en vzorec

Primerjava povprečja na vzorcu z določeno vrednostjo (npr. populacijskim parametrom ali nevtralno točko na lestvici)

$$t = \frac{\text{mean} - \text{comparison}}{\text{standarderror}}$$

Predpostavke:

- Slučajno vzorčenje, neodvisni vzorci
- Normalna porazdelitev podatkov (če je  $N < 30$ )

Primer: Ali se študenti pri preizkusu odrežejo bolje od naključja?

28 študentov je opravilo test s 100 vprašanji o besedilu, ki ga niso prebrali.

Vsako vprašanje je imelo 5 možnih odgovorov, zato bi pri povsem naključni izbiri pričakovali, da bo pravilno izbranih 20 postavk.

Vendar so v povprečju pravilno odgovorili na 46,57 vprašanj. S  $t$ -testom za en vzorec (*angl. one-sample t-test*) pokažemo, da so se študenti odrezali statistično značilno bolje od naključja ( $M = 46.57, t(27) = 20.6, p < 0.001$ ).

## Vaja 1

Izbrali smo vzorec šestnajstih otrok in jih stehali. Določi interval zaupanja za pravo vrednost aritmetične sredine s 5% tveganja.

X: 35 37 29 26 31 32 28 40 27 33 33 34 31 30 29 38

## Vaja 2

Imamo podatke za vzorec učencev o tem, koliko ur tedensko porabijo za učenje doma. 78 jih je odgovorilo, da se pripravljajo na pouk 2 do 3 ure tedensko, 125 s jih uči 3 do 4 ure na teden, 103 se učijo vsak teden več kot 4 ure. Določi odstotek tistih učencev, ki tedensko posvečajo učenju najmanj časa in oceni ta odstotek v osnovni množici (z 1% tveganjem).

## Vaja 3

Učitelja matematike zanima, ali je povprečno število točk njegovih učencev na preizkusu znanja višje od nacionalnega povprečja, ki je 75. Izračunajte vrednost  $t$ -statistike pri stopnji tveganja 5%, določite kritično območje in sklepajte o rezultati za naslednje podatke:

Za sodelovanje v raziskavi je bil izbran vzorec 25 učencev.

Povprečna ocena iz matematičnega izpita v vzorcu je 80.

Standardni odklon v vzorcu je 6.

## Vaja 4

Želimo določiti minimalno velikost vzorca za pilotno študijo novega izobraževalnega programa. Standardni odklon učinka programa na bralne sposobnosti je 3 točke. Želite imeti 95% stopnjo zaupanja, da bo vaša ocena učinka natančna. Izračunajte minimalno velikost vzorca.

## 0.4 Bivariantna statistika

### T-test

#### T-test za (parne) odvisne vzorce

Primerjava povprečij dveh pogojev, v katerih so sodelovale iste enote.

Primer: 20 študentov je dobilo test v izpolnjevanje pred študijem določenega predmeta in nato ponovno po zaključku tega predmeta.

Testna statistika  $t = \frac{\bar{d}}{SE(\bar{d})}$ .

Izračun:

1. Postavimo ničelno in alternativno hipotezo:
  - $H_0$ : Ni razlik v znanju študentov pred in po študiju tega modula.
  - $H_1$ : Obstajajo razlike v znanju študentov pred in po študiju tega modula.
2. Izračunamo razlike med pari opazovanj:  $d_i = y_i - x_i$  (za vajo lahko to naredimo v SPSS).
3. Izračunamo povprečje razlik:  $\bar{d}$ .
4. Standardni odklon razlik:  $s_d$ .
5. Standardna napaka povprečne razlike:  $SE(\bar{d}) = \frac{s_d}{\sqrt{n}}$ .
6. T-statistika:  $t = \frac{\bar{d}}{SE(\bar{d})}$  (empirična vrednost);  $df = n - 1$ .
7. V tabeli poiščemo kritično vrednost pri  $\alpha = 5\%$ .
8. Empirična vrednost  $t$  pade v kritično območje, ki ga določa teoretična vrednost  $t$  pri dani stopnji zaupanja, zato lahko zavrnilo ničelno hipotezo in sprejmemo alternativno.....  
Napiši kaj če ne pade v kritično pade!!!!
9. Interval zaupanja za resnično vrednost razlike povprečij je:

$$\bar{d} \pm (t \cdot SE(\bar{d}))$$

#### T-test za neodvisne vzorce

Primerjava preizkusa domneve o srednjih vrednostih dveh skupin enot.

Primer: Primerjava kalorične vsebnosti dveh vrst štrudlja.

Testna statistika  $t = \frac{\bar{x}_1 - \bar{x}_2}{SE(\bar{x}_1 - \bar{x}_2)}$ .

Opomba: Test predpostavi enakost varianc neodvisnih vzorcev. Če to ni zagotovljeno, uporabimo Welchov test  $t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$ .

Izračun:

1. Postavimo ničelno in alternativno hipotezo:
  - $H_0$ : Ni razlik v kalorični vsebnosti med dvema vrstama hotdoga.
  - $H_1$ : So razlike v kalorični vsebnosti med dvema vrstama hotdoga.
2. Razlika med povprečnima vrednostima:  $\bar{x}_1 - \bar{x}_2$ .
3. Skupni standardni odklon (pod predpostavko enakih varianc):

$$s_p = \sqrt{\frac{(n_1 - 1) \cdot s_1^2 + (n_2 - 1) \cdot s_2^2}{n_1 + n_2 - 2}}$$



4. Standardna napaka:

$$SE(\bar{x}_1 - \bar{x}_2) = s_p \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

5. Empirična t-vrednost:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{SE(\bar{x}_1 - \bar{x}_2)}$$

6. Prostostne stopnje:  $df = n_1 + n_2 - 2$ .

7. V tabeli poiščemo kritično vrednost pri  $\alpha = 5\%$ .

8. Interval zaupanja za resnično vrednost razlike povprečij:

$$\bar{x}_1 - \bar{x}_2 \pm (t \cdot SE(\bar{x}_1 - \bar{x}_2))$$

### Analiza variance (ANOVA)

Primerjava povprečij večih skupin (če dve skupini je ANOVA enaka T testu).

$F$  porazdelitev ima dve prostorski skopnji

Slika variance between and within.

Primer: Tri skupine desetih slučajno izbranih študentov so postavljene v tri različne učilnice. A ima konstantno glasbo v ozadju, B variabilno glasbo v ozadju, C brez glasbe. Po enem mesecu nas zanima, ali glasba pomaga pri učenju.

Izračun:

1. Postavimo ničelno in alternativno hipotezo:

- $H_0$ : Med skupinami ni razlik v vsrkavanju informacij.
- $H_1$ : Med skupinami so razlike v vsrkavanju informacij.

2. Izračunamo povprečja. Skupno povprečje je  $\bar{x} = 5,1$ , povprečja posameznih skupin pa so  $\bar{x}_1 = 7$ ,  $\bar{x}_2 = 4$ , in  $\bar{x}_3 = 4,3$ .

3. Vsota kvadratov:

- Med skupinami:  $SS_{between} = 54,6$
- Znotraj skupin:  $SS_{within} = 90,1$

4. Prostostne stopnje:

- Med skupinami:  $df_{between} = 2$
- Znotraj skupin:  $df_{within} = 27$

5. Povprečni kvadrat:

- Med skupinami:  $MS_{between} = \frac{SS_{between}}{df_{between}} = \frac{54,6}{2} = 27,3$
- Znotraj skupin:  $MS_{within} = \frac{SS_{within}}{df_{within}} = \frac{90,1}{27} = 3,34$

6. Empirična F-vrednost:

$$F = \frac{MS_{between}}{MS_{within}} = \frac{27,3}{3,34} = 8,18$$

7. V tabeli poiščemo kritično vrednost pri  $\alpha = 5\%$ :  $F_{critical} = 2,03$ .

8. Ker je empirična F-vrednost (8,18) večja od kritične vrednosti (2,03), zavrnemo ničelno hipotezo  $H_0$ .

9. Izračunamo eta kvadrat ( $\eta^2$ ), ki je merilo učinka:

$$\eta^2 = \frac{SS_{between}}{SS_{between} + SS_{within}} = \frac{54,6}{54,6 + 90,1} = 0,38$$

$$\eta = \sqrt{\eta^2} = \sqrt{0,38} \approx 0,62$$

### Neparametrične alternative

Če podatki niso normalno porazdeljeni, moramo parametrični test zamenjati z ustreznim neparametričnim testom.

| Cilj   | Parametrični test   | Neparametrični test      |
|--|---------------------|--------------------------|
| Testiranje razlike med dvema odvisnima nizoma enot         | Odvisni $t$ -test   | Wilcoxon test predznačen |
| Testiranje razlike med dvema neodvisnima nizoma enot       | Neodvisni $t$ -test | Mann-Whitney $U$ test    |
| Testiranje razlike med tremi ali več neodvisnimi nizi enot | ANOVA               | Kruskal-Wallis $H$ test  |

Tabela 1: Pregled parametričnih in neparametričnih testov

### Hi-kvadrat test

Uporablja se za ugotavljanje, ali obstaja statistično značilna razlika med pričakovanimi in opazovanimi frekvencami v eni ali več kategorijah (torej če imamo nominalne spremenljivke)

Ničelna hipoteza  $H_0$ : V populaciji ni povezanosti med spremenljivkama.

Slika hi kvadrat porazdelitve

**Izračunavanje stopnje prostosti ( $df$ ):** Stopnja prostosti ( $df$ ) za hi-kvadrat test se izračuna po formuli:

$$df = (\text{vrstice} - 1) \cdot (\text{stolpci} - 1)$$

### Postopek izračuna:

1. Zberemo podatke in uredimo frekvence v kontingenčni tabeli.
2. Izračunamo pričakovane frekvence za vsako celico tabele.
3. Uporabimo formulo za hi-kvadrat vrednost:

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

kjer je  $O_i$  opazovana frekvenca,  $E_i$  pa pričakovana frekvenca.

**Uporaba kontingenčnih koeficientov:** Ker vrednosti hi-kvadrat same po sebi niso primerljive med različnimi tabelami, pogosto uporabimo kontingenčne koeficiente, kot je Cramérjev  $V$ , ki ga izračunamo po formuli:

$$V = \sqrt{\frac{\chi^2}{N \cdot (k - 1)}},$$

kjer je  $\chi^2$  hi-kvadrat vrednost,  $N$  skupno število opazovanj,  $k$  pa število kategorij v najbolj številni spremenljivki.

**Spearmanov koeficient**

Spearmanov koeficient korelacije meri moč in smer monotone povezave med dvema spremenljivkama. Uporablja se za ordinalne spremenljivke ali za kvantitativne spremenljivke, ki niso nujno normalno porazdeljene.

**Ničelna hipoteza  $H_0$ :** Med dvema spremenljivkama ni monotone povezave.

**Formula:** Spearmanov koeficient ( $\rho$ ) se izračuna po formuli:

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

kjer je  $d_i$  razlika med vrstnimi številkami vsakega opazovanja in  $n$  število opazovanj.

**Postopek izračuna:**

1. Razvrstimo podatke v naraščajočem vrstnem redu za vsako spremenljivko.
2. Izračunamo vrstne številke za vsako spremenljivko.
3. Izračunamo razlike med vrstnimi številkami ( $d_i$ ) za vsako opazovanje.
4. Uporabimo zgornjo formulo za izračun Spearmanovega koeficienta.

**Interpretacija:** Spearmanov koeficient se giblje med -1 in 1. Vrednost blizu 1 kaže na močno pozitivno monotono povezavo, vrednost blizu -1 kaže na močno negativno monotono povezavo, vrednost blizu 0 pa kaže na odsotnost monotone povezave.

**Zaključek:** Spearmanov koeficient je uporaben za analizo povezav med ordinalnimi ali kvantitativnimi spremenljivkami, ki niso normalno porazdeljene. Pomaga nam razumeti moč in smer monotone povezave med spremenljivkami.

**Pearsonov koeficient**

Pearsonov koeficient korelacije meri linearno povezanost med dvema kvantitativnima spremenljivkama. Uporablja se za intervalne ali razmerne spremenljivke, ki so normalno porazdeljene.

**Ničelna hipoteza  $H_0$ :** Med dvema spremenljivkama ni linearne povezave.

**Formula:** Pearsonov koeficient ( $r$ ) se izračuna po formuli:

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}},$$

kjer je  $x_i$  in  $y_i$  vsaka opazovanja za spremenljivki  $x$  in  $y$ ,  $\bar{x}$  in  $\bar{y}$  pa povprečji za spremenljivki  $x$  in  $y$ .

**Postopek izračuna:**

1. Izračunamo povprečje za vsako spremenljivko.
2. Izračunamo odstopanja vsakega opazovanja od povprečja.
3. Pomnožimo odstopanja za ustrezna opazovanja in izračunamo vsoto.
4. Izračunamo kvadrate odstopanj za vsako spremenljivko in vsoto kvadratov.
5. Uporabimo zgornjo formulo za izračun Pearsonovega koeficienta.

**Interpretacija:** Pearsonov koeficient se giblje med -1 in 1. Vrednost blizu 1 kaže na močno pozitivno linearno povezavo, vrednost blizu -1 kaže na močno negativno linearno povezavo, vrednost blizu 0 pa kaže na odsotnost linearne povezave.

### Povezanost

Povezanost med dvema spremenljivkama pomeni, da spremembe v eni spremenljivki sovpadajo s spremembami v drugi spremenljivki. Povezanost lahko merimo na različne načine, odvisno od narave spremenljivk in vrste povezave.

### Funkcionalna povezanost

Funkcionalna povezanost pomeni, da obstaja točno določena matematična funkcija, ki povezuje dve spremenljivki. Na primer, če  $y = f(x)$ , potem je  $y$  popolnoma določeno s  $x$ . To je najmočnejša oblika povezanosti, saj vsaka vrednost ene spremenljivke določa točno eno vrednost druge spremenljivke.

### Korelacijska povezanost

Korelacijska povezanost meri stopnjo in smer linearne povezave med dvema kvantitativnima spremenljivkama. Najpogosteje uporabljamo Pearsonov koeficient korelacije ( $r$ ), ki meri linearno povezanost, in Spearmanov koeficient korelacije ( $\rho$ ), ki meri monotono povezanost.

### Močna in šibka povezanost

Moč povezanosti se nanaša na velikost koeficienta korelacije.

- **Močna povezanost:** Koeficient korelacije blizu 1 ali -1 kaže na močno povezanost.
- **Šibka povezanost:** Koeficient korelacije blizu 0 kaže na šibko povezanost.

### Linearna in nelinearna povezanost

- **Linearna povezanost:** Povezanost, kjer lahko povezavo med spremenljivkama najbolj opišemo z ravno črto. Pearsonov koeficient korelacije ( $r$ ) je primeren za merjenje linearne povezanosti.
- **Nelinearna povezanost:** Povezanost, kjer je povezava med spremenljivkama bolj zapletena in je ni mogoče opisati z ravno črto. Spearmanov koeficient korelacije ( $\rho$ ) je primeren za merjenje monotone (nelinearne) povezanosti.

### Pozitivna in negativna povezanost

- **Pozitivna povezanost:** Ko se ena spremenljivka povečuje, se druga spremenljivka tudi povečuje. Koeficient korelacije je pozitiven.
- **Negativna povezanost:** Ko se ena spremenljivka povečuje, se druga spremenljivka zmanjšuje. Koeficient korelacije je negativen.

### Kovarianca

Kovarianca je mera, ki opisuje, kako se dve spremenljivki sočasno spreminjata. Pozitivna kovarianca pomeni, da se večje vrednosti ene spremenljivke povezujejo z večjimi vrednostmi druge spremenljivke, medtem ko negativna kovarianca pomeni, da se večje vrednosti ene spremenljivke povezujejo z manjšimi vrednostmi druge spremenljivke.

**Formula**

Kovarianca dveh spremenljivk  $X$  in  $Y$  se izračuna po formuli:

$$\text{Cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}),$$

kjer:

- $X_i$  in  $Y_i$  sta posamezni opazovanji spremenljivk  $X$  in  $Y$ ,
- $\bar{X}$  in  $\bar{Y}$  sta povprečji spremenljivk  $X$  in  $Y$ ,
- $n$  je število opazovanj.

**Interpretacija**

- **Pozitivna kovarianca:** Ko se vrednosti ene spremenljivke povečujejo, se povečujejo tudi vrednosti druge spremenljivke.
- **Negativna kovarianca:** Ko se vrednosti ene spremenljivke povečujejo, se vrednosti druge spremenljivke zmanjšujejo.
- **Kovarianca blizu nič:** Ni jasno izražene povezave med spremenljivkama.

**Razlika med kovarianco in korelacijo**

Kovarianca meri le smer sočasnih sprememb dveh spremenljivk, vendar ni standardizirana, zato njena vrednost ni omejena. Po drugi strani pa je korelacija standardizirana mera, ki pove, kako močna in v katero smer je linearna povezava med dvema spremenljivkama, njena vrednost pa je vedno med -1 in 1.

**Vaja 1**

...

## 0.5 Multivariantna statistika

Povezanost med spremenljivkama ne pomeni nujno, da med njima obstaja vzročna povezava. Spremenljivke so lahko povezane tudi navidezno in jih pojasni uvedba tretje spremenljivke (npr. poletni čas pojasni povezavo med napadi morskih psov in prodajo sladoleda).

### Multipla regresijska analiza (linearna)

Multipla regresijska analiza je metoda za preučevanje razmerja med odvisno spremenljivko in dvema ali več neodvisnimi spremenljivkami. Namen je napovedovanje vrednosti odvisne spremenljivke z nizom neodvisnih spremenljivk.

Regresijski model je podan kot:

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

Delež variabilnosti, ki ga model pojasni, je izražen z  $R^2$ . Velikost vzorca naj bo vsaj 10-krat večja od števila neodvisnih spremenljivk.

Predpostavke multiple linearne regresije:

- **Neodvisnost:** Opazovanja so neodvisna.
- **Normalnost:** Napake v regresijskem modelu so normalno porazdeljene.
- **Homoskedastičnost:** Variabilnost napak je konstantna pri vseh nivojih neodvisnih spremenljivk.
- **Linearnost:** Obstaja linearna povezava med odvisno in neodvisnimi spremenljivkami.

### Multipla regresijska analiza (logistična)

Logistična regresija se uporablja, kadar je odvisna spremenljivka nominalna. Namesto linearne funkcije se uporablja logistična funkcija, ki omogoča napovedovanje kategorije odvisne spremenljivke.

Logistični model je podan kot:

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k,$$

kjer je  $p$  verjetnost, da se odvisna spremenljivka pojavi v določeni kategoriji.

### Razvrščanje v skupine (clustering)

Razvrščanje v skupine ali clustering je metoda za razdelitev podatkov v več homogenih skupin na podlagi podobnosti med posameznimi podatkovnimi točkami. Najpogostejše metode vključujejo k-means, hierarhično razvrščanje in DBSCAN.

### Metode zmanjšanja dimenzionalnosti podatkov

**Analiza glavnih komponent (PCA):** PCA je tehnika za zmanjšanje dimenzionalnosti podatkov, ki pretvori več povezanih spremenljivk v manj nepovezanih komponent. To omogoča lažjo vizualizacijo in analizo podatkov.

**Faktorska analiza:** Faktorska analiza identificira latentne spremenljivke ali faktorje, ki pojasnjujejo vzorce korelacij med opazovanimi spremenljivkami. Faktorji so teoretične konstrukte, ki so lahko osnovni vzroki za opazovane korelacije.

## Analiza zanesljivosti

**Cronbachov  $\alpha$  koeficient:** Cronbachov  $\alpha$  koeficient je mera za notranjo konsistenco (zanesljivost) skale ali vprašalnika. Visoka vrednost  $\alpha$  (bližje 1) nakazuje na visoko zanesljivost skale.

## Druge multivariantne metode

**Kanonična korelacijska analiza:** Kanonična korelacijska analiza je metoda za preučevanje povezave med dvema sklopoma spremenljivk. Omogoča določitev linearnih kombinacij spremenljivk iz obeh sklopov, ki so medsebojno najbolj povezane.

**Diskriminantna analiza:** Diskriminantna analiza je tehnika za razvrščanje opazovanj v predhodno določene skupine na podlagi neodvisnih spremenljivk. Uporablja se za napovedovanje kategorijske odvisne spremenljivke.

**Strukturni modeli:** Strukturni modeli so kompleksni statistični modeli, ki omogočajo preučevanje vzročnih odnosov med več spremenljivkami. Združujejo elemente regresijske analize, faktorske analize in poti.

## Vaja 1

...