

Statistika

Bor Bregant

13. avgust 2024

Kazalo

0.1	Predgovor	1
1	Uvod	3
1.1	Podatki in informacije	3
1.1.1	Enote opazovanja in Spremenljivke	3
1.1.2	Tipi spremenljivk glede na vrednost in glede na mersko lestvico	3
1.2	Populacija in vzorec	5
1.3	Vrste statistične analize	5
1.4	Koraki statistične analize	6
2	Opisna statistika in vizualizacija	9
2.1	Absolutna, relativna in grupirana frekvenčna porazdelitev	9
2.2	Grafična predstavitev frekvenčnih porazdelitev za nominalne in ordinalne spremenljivke	9
2.3	Grafična predstavitev frekvenčnih porazdelitev za intervalne in razmernostne spremenljivke	10
2.4	Normalna porazdelitev	12
2.4.1	Asimetričnost (skewness)	12
2.4.2	Sploščenost (kurtosis)	12
2.5	Rangiranje	13
2.5.1	Kvantili	13
2.6	Mere centralne tendence	14
2.7	Mere variabilnosti (dispersije)	15
2.8	Inferenčna statistika	17
3	Bivariantna analiza	21
3.1	T-test	21
3.1.1	T-test za (parne) odvisne vzorce	21
3.1.2	T-test za neodvisne vzorce	21
3.2	Analiza variance (ANOVA)	22
3.3	Neparametrične alternative	23
3.3.1	Hi-kvadrat test	23
3.3.2	Spearmanov koeficient	24
3.3.3	Pearsonov koeficient	24
3.4	Povezanost	25
3.4.1	Kovarianca	26
4	Multivariantna statistika	27
4.1	Multipla regresijska analiza (linearna)	27
4.2	Multipla regresijska analiza (logistična)	28
4.3	Razvrščanje v skupine (clustering)	29
4.4	Metode zmanjšanja dimenzionalnosti podatkov	30
4.5	Analiza zanesljivosti	31
4.6	Druge multivariantne metode	31

0.1 Predgovor

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).

Kljub temu, da je bilo gradivo pregledano, boste gotovo našli v njem napake. Vesel bom, če me boste nanje opozorili.

Poglavje 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).

1.1 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

1.1.1 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)

1.1.2 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
        directory as the script
6      df = pd.read_csv('data.csv')
7
8      # Step 2: Print the first few rows of the dataframe
9      print("First few rows of the dataframe:")
10     print(df.head())
11
12     # Step 3: Check the data types of each column
13     print("\nData types of each column:")
14     print(df.dtypes)
15
16     # Step 4: Convert a specific column to float (if necessary)
17     # Let's assume we have a column named 'income' which we want to convert
        to float
18     df['income'] = df['income'].astype(float)
19
20     # Verify the conversion
21     print("\nData types after conversion:")
22     print(df.dtypes)
23
24     # Step 5: Calculate the range of data in a numeric column
25     # Let's calculate the range for the 'income' column
26     income_range = df['income'].max() - df['income'].min()
27     print("\nRange of the 'income' column:")

```



```

28     print(income_range)
29
30     # Step 6: Label encoding for an ordinal categorical variable
31     # Let's assume we have an ordinal variable named 'education_level'
32     education_levels = {'High School': 1, 'Bachelor': 2, 'Master': 3, 'PhD'
: 4}
33     df['education_level'] = df['education_level'].map(education_levels)
34
35     # Verify the encoding
36     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'.")

```

1.2 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

1.3 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 ...

1.4 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. korakov
- 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

Poglavje 2

Opisna statistika in vizualizacija

2.1 Absolutna, relativna in grupirana frekvenčna porazdelitev

- **Frekvenčna porazdelitev** je matematična funkcija, ki prikazuje število primerov, v katerih spremenljivka zavzame vsako od svojih možnih vrednosti.
- **Frekvenca** je število pojavov podatkovne vrednosti.
- **Relativna frekvenca** je odstotek enot med vsemi enotami, ki imajo določeno vrednost spremenljivke, in se izračuna tako, da se absolutna frekvenca deli s številom vseh enot.

Primer

x_i	f_i	f_i (%)
Koper	60	20
Izola	150	50
Piran	90	30
Skupaj	300	100

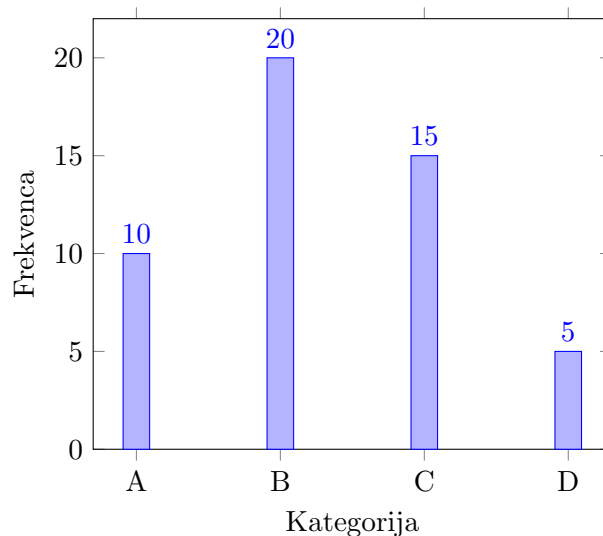
Tabela 2.1: Relativna frekvenca najljubših slovenskih obalnih mest. Enota je turist, spremenljivka pa najljubše slovensko obalno mesto.

Kumulativna frekvenca je seštevek vseh frekvenc za določene vrednosti spremenljivke do vključno določene vrednosti. Prikazuje, koliko primerov ima vrednost spremenljivke manjšo ali enako določeni vrednosti. Kumulativna frekvenca je koristna za prikazovanje porazdelitve podatkov in ugotavljanje, koliko enot dosega ali presega določeno vrednost.

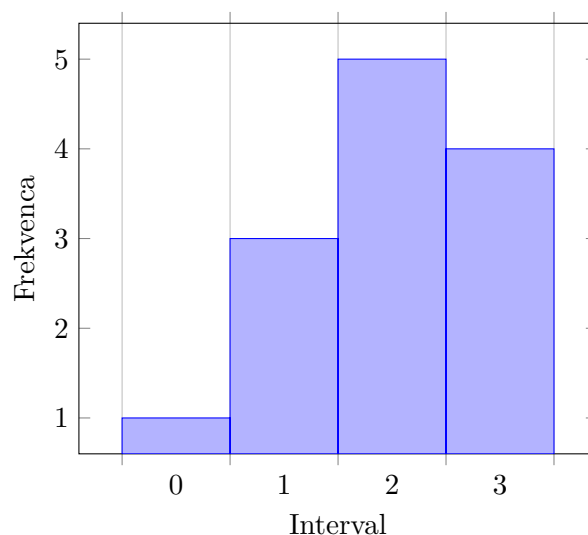
Grupirane frekvenice se uporabljajo, kadar so podatki razdeljeni v skupine ali razrede. Namesto da bi beležili frekvenice za vsako posamezno vrednost, združimo vrednosti v skupine, kar omogoča bolj pregledno analizo podatkov, še posebej pri velikih naborih podatkov. Relativne frekvenice se lahko izračunajo tudi za te skupine, da dobimo odstotek primerov v vsaki skupini.

2.2 Grafična predstavitev frekvenčnih porazdelitev za nominalne in ordinalne spremenljivke

Stolpčni grafikon prikazuje frekvenice posameznih kategorij z navpičnimi ali vodoravnimi stolpci, kjer višina oziroma dolžina stolpca predstavlja frekvenco določene kategorije. Prednost stolpčnih grafikonov je v tem, da omogočajo enostavno primerjavo med kategorijami in so primernejši za prikaz večjih množic podatkov od tortnega grafikona.



Slika 2.1: Stolpčni grafikon frekvenc za nominalne spremenljivke.



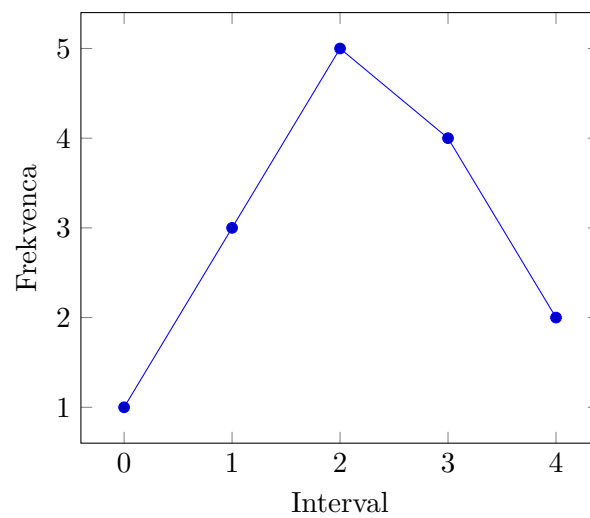
Slika 2.2: Histogram frekvenc za intervalne spremenljivke.

2.3 Grafična predstavitev frekvenčnih porazdelitev za intervalne in razmernostne spremenljivke

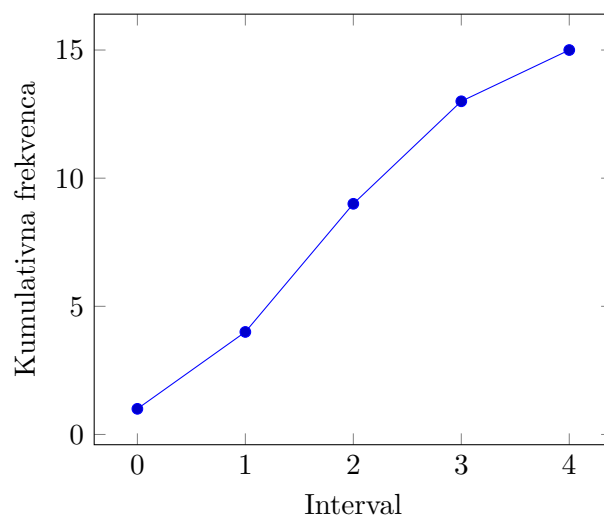
Histogram je vrsta stolpčnega grafikona, kjer stolpci predstavljajo frekvence podatkov v določenih intervalih. Histogram je uporaben za prikaz razpršenosti podatkov in za ugotavljanje oblik porazdelitev, kot so normalna, enakomerno porazdeljena ali pristranska porazdelitev.

Poligon frekvenc je črtni grafikon, ki povezuje sosednje točke, ki predstavljajo frekvence intervalov. Uporablja se za prikaz porazdelitve podatkov in omogoča enostavno primerjavo z drugimi porazdelitvami.

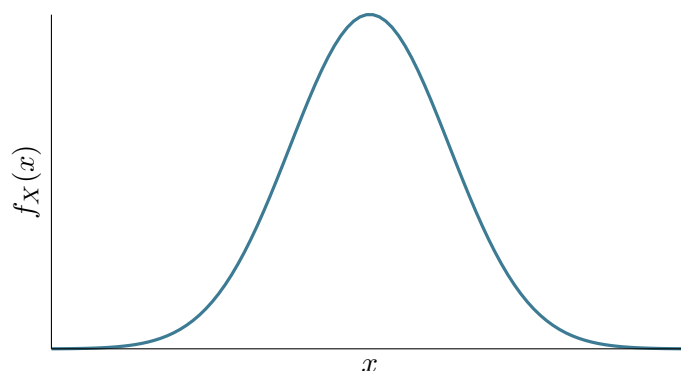
Ogiva je črtni grafikon, ki prikazuje kumulativne frekvence. To pomeni, da vsaka točka na grafu predstavlja vsoto vseh prejšnjih frekvenc do določene vrednosti. Ogiva je uporabna za določanje percentilov in za primerjavo kumulativnih porazdelitev med različnimi skupinami podatkov.



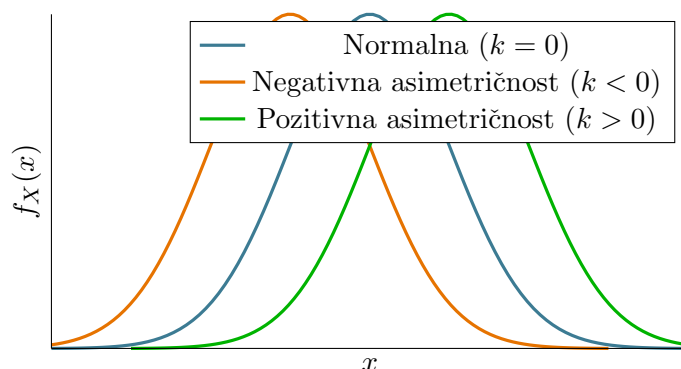
Slika 2.3: Poligon frekvenc za intervalne spremenljivke.



Slika 2.4: Ogiva za kumulativne frekvence.



Slika 2.5: Gostotna funkcija normalne porazdelitve.



Slika 2.6: Grafi normalne porazdelitve s pozitivno in negativno asimetričnostjo.

2.4 Normalna porazdelitev

Normalna porazdelitev je ena izmed najpomembnejših verjetnostnih porazdelitev v statistiki in se pogosto uporablja pri modeliranju različnih naravnih in družbenih pojavov. Gostotna funkcija normalne porazdelitve je podana z naslednjo enačbo:

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

kjer je μ povprečje (aritmetična sredina) in σ standardni odklon. Graf normalne porazdelitve je zvonaste oblike in simetričen glede na povprečje μ .

2.4.1 Asimetričnost (skewness)

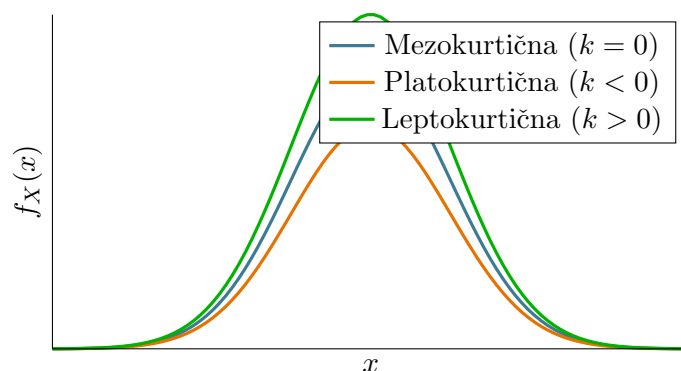
Asimetričnost meri, kako je porazdelitev nagnjena v levo ali desno. Normalna porazdelitev je popolnoma simetrična in ima koeficient asimetričnosti $k = 0$. Vendar pa v realnih podatkih pogosto srečamo porazdelitve, ki niso povsem simetrične:

- Negativna asimetričnost (v levo, $k < 0$): Rep porazdelitve je daljši na levi strani.
- Pozitivna asimetričnost (v desno, $k > 0$): Rep porazdelitve je daljši na desni strani.

Če je k med -1 in 1, se porazdelitev še vedno lahko šteje za normalno.

2.4.2 Sploščenost (kurtosis)

Sploščenost meri, kako visoka in ostra je konica porazdelitve v primerjavi z normalno porazdelitvijo:



Slika 2.7: Grafi normalne porazdelitve s pozitivno in negativno sploščenostjo.

- Leptokurtična ($k > 0$): Porazdelitev ima ožjo in višjo konico ter debelejšše repe.
- Mezokurtična ($k = 0$): Porazdelitev je normalno zvonaste oblike.
- Platokurtična ($k < 0$): Porazdelitev ima ploščato konico in tanjše repe.

Če je k med -1 in 1, se porazdelitev še vedno lahko šteje za normalno.

2.5 Rangiranje

Rangiranje vključuje razvrščanje podatkovnih točk v naraščajočem ali padajočem vrstnem redu. Rang določene vrednosti v podatkovnem naboru je njeno mesto v tem vrstnem redu. Rangiranje je uporabno za prepoznavanje relativnega položaja podatkovne točke znotraj nabora podatkov.

Primer:

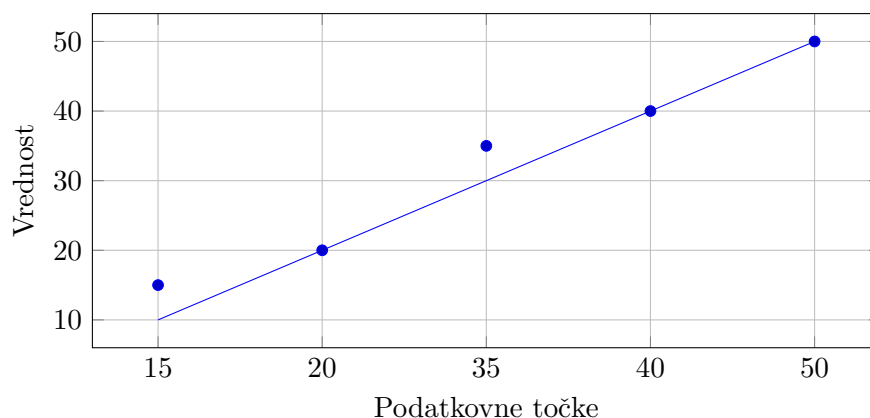
- Podatkovni niz: 45, 32, 67, 23, 89, 56
- Urejeni podatkovni niz: 23, 32, 45, 56, 67, 89
- Rangiranje:
 - 23 ima rang 1
 - 32 ima rang 2
 - 45 ima rang 3
 - 56 ima rang 4
 - 67 ima rang 5
 - 89 ima rang 6

2.5.1 Kvantili

Kvantili so vrednosti, ki delijo urejen niz podatkov na enake dele. Najpogosteje uporabljeni kvantili so kvartili, decili in percentili.

Kvartili: Kvartili delijo podatkovni niz na štiri enake dele:

- Prvi kvartil (Q_1): 25. percentil
- Drugi kvartil (Q_2): 50. percentil (median)
- Tretji kvartil (Q_3): 75. percentil



Slika 2.8: Prikaz podatkovnega niza in kvantilov.

Decili: Decili delijo podatkovni niz na deset enakih delov:

- Prvi decil (D_1): 10. percentil
- Drugi decil (D_2): 20. percentil
- Tretji decil (D_3): 30. percentil
- itd.

Percentili: Percentili delijo podatkovni niz na sto enakih delov:

- 10. percentil: vrednost pod katero leži 10% podatkov
- 50. percentil: vrednost pod katero leži 50% podatkov (median)
- 90. percentil: vrednost pod katero leži 90% podatkov

Primer: Podatkovni niz: 15, 20, 35, 40, 50

- Q_1 (25. percentil): $20 + 0.25 \times (35 - 20) = 23.75$
- Q_2 (50. percentil): Median = 35
- Q_3 (75. percentil): $35 + 0.75 \times (50 - 35) = 46.25$

2.6 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

2.7 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.

Razpon

Absolutni razpon

$$R_{abs} = x_{max} - x_{min}$$

Relativni razpon

$$R_{rel} = \frac{2(x_{max} - x_{min})}{x_{max} + x_{min}}$$

Interkvartilni rang

$$IQR = \frac{Q_3 - Q_1}{2}$$

Absolutna deviacija mediane/aritmetične sredine

$$AD_{Me, \mu} = \frac{1}{N} \sum_{i=1}^N |x_i - Me, \mu|$$

Varianca in standardni odklon

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

$$\sigma = \sqrt{\sigma^2}$$

Variabilnost normalne porazdelitve

Slika naslednjega

- 68,3% enot je znotraj enega standardnega odklona od povprečja $[\mu - \sigma, \mu + \sigma]$
- 95,4% enot je znotraj dveh standardnih odklonov od povprečja $[\mu - 2\sigma, \mu + 2\sigma]$
- 99,7% enot je znotraj treh standardnih odklonov od povprečja $[\mu - 3\sigma, \mu + 3\sigma]$

Standardizacija

$$z_i = \frac{x_i - \mu_X}{\sigma_X}$$

Rezultat je standardizirana spremenljivka Z , kjer standardizirane vrednosti z_i predstavljajo relativna odstopanja od aritmetične sredine. Standardizacija nam omogoča primerjavo vrednosti različnih spremenljivk, ki praviloma niso primerljive.

Primer primerjave teže in višine dojenčkov, kjer standardiziramo glede na spol.

Dojenček	Spol	Teža (kg)	μ_X (kg)	σ_X (kg)	z_i
1	Moški	3.5	3.4	0.5	$z_1 = \frac{3.5-3.4}{0.5} = 0.2$
2	Ženski	3.0	3.2	0.4	$z_2 = \frac{3.0-3.2}{0.4} = -0.5$
3	Moški	4.0	3.4	0.5	$z_3 = \frac{4.0-3.4}{0.5} = 1.2$
4	Ženski	2.8	3.2	0.4	$z_4 = \frac{2.8-3.2}{0.4} = -1.0$

Tabela 2.2: Standardizacija teže dojenčkov glede na spol.

Vaja 1

Učence smo testirali v znanju matematike in dosegli so naslednje rezultate:

20 17 25 20 12 25 18 16 25
 17 21 22 22 16 17 22 20 19 16
 24 20 15 14 23 15 19 21 26
 18 29 30 20 16 12 24 30 27 25
 15 25 24 23 26 21

1. Poišči rezultat z rangom 8 (od najslabšega).
2. Koliko % učencev je doseglo več kot 22 točk?
3. Kateri rezultat (ali rezultati) se najpogosteje pojavlja?
4. Rezultate testiranja razvrsti v frekvenčno porazdelitev.
5. Podatke prikaži grafično.

Vaja 2

Skupino desetih učencev smo testirali s psihodiagnostičnim testom in dobili naslednje rezultate:

17 26 22 15 18 27 16 20 18 24

1. Izračunaj kvartile.
2. Določi kvantilni rang vrednosti $x = 22.50$.
3. Določi kvantilni rang rezultata $x = 22.00$.

Vaja 3

Na eni izmed avtošol smo se pozanimali, Koliko ur vožnje so kandidati potrebovali, preden so uspešno opravili vozniški izpit. Izračunaj aritmetično sredino, varianco in standardni odklon.

2.8 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:* Vzorec, ki naj bi bil reprezentativen izbire strokovnjak področja raziskave.
 - *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)
nekej nekej statistika

Intervali zaupanja

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

S tveganjem α lahko rečemo, da interval (a, b) vsebuje parameter γ .

Slika normalne.

Širina intervala zaupanja je odvisna od:

- *Stopnje tveganja*: Višja je stopnja tveganja α , 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 naš 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)

V statistiki ločimo med dvema vrstama napak, ki ju lahko naredimo pri testiranju hipotez:

- **Napaka prve vrste (α -napaka)**: To je napaka, ki jo naredimo, ko zavrnemo ničelno hipotezo (H_0), čeprav je v resnici resnična. Gre torej za "lažno pozitivni" rezultat. Verjetnost, da se zgodi napaka prve vrste, imenujemo raven značilnosti α , ki je običajno nastavljena na 0,05 ali 5 %. To pomeni, da obstaja 5-odstotna možnost, da napačno zavrnemo resnično ničelno hipotezo.

- **Napaka druge vrste (β -napaka):** To je napaka, ki jo naredimo, ko ne zavrnilimo ničelne hipoteze (H_0), čeprav je v resnici napačna. Gre za "lažno negativni" rezultat. Verjetnost, da se zgodi napaka druge vrste, označujemo z β . Komplement te verjetnosti ($1 - \beta$) predstavlja moč testa, ki je verjetnost, da bomo pravilno zavrnili napačno ničelno hipotezo.

V praksi se moramo zavedati, da je zmanjšanje verjetnosti ene vrste napake običajno povezano s povečanjem verjetnosti druge vrste napake, zato je pomembno, da testiranje hipotez izvedemo premišljeno in uravnotežimo tveganja za obe vrsti napak.

Stopnja značilnosti

Izberemo največjo verjetnost α 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 α 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.

Poglavje 3

Bivariantna analiza

3.1 T-test

3.1.1 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 območje!
9. Interval zaupanja za resnično vrednost razlike povprečij je:

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

3.1.2 T-test za neodvisne vzorce

Primerjava preizkusa domneve o srednjih vrednosti dveh skupin enot.

Primer: Primerjava kalorične vrednosti 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:

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

3.2 Analiza variance (ANOVA)

Primerjava povprečij večih skupin (če sta samo 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), zavrnamo ničelno hipotezo H_0 .

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

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

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

3.3 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 3.1: Pregled parametričnih in neparametričnih testov

3.3.1 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 χ^2 hi-kvadrat vrednost, N skupno število opazovanj, k pa število kategorij v najbolj številni spremenljivki.

3.3.2 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 (ρ) 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.

3.3.3 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.

3.4 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 (ρ), 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 najbolje 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 (ρ) 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.

3.4.1 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

V datoteki data_cleaned.xlsx določi naslednje:

- Ali obstaja povezava med matematično anksioznostjo in motivacijo za matematiko? Če je povezava znatna, določi še enačbo linearne regresije, kjer si sam izbereš odvisno in neodvisno spremenljivko.
- Ali obstaja povezava med spolom in anksioznostjo?
- Ali obstaja povezava med profesorjem in razredom?

Imej v mislih, da so kategorične spremenljivke v dokumentu kodirane označevalno.

Poglavje 4

Multivariantna statistika

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

4.1 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:

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

kjer je:

- Y odvisna spremenljivka,
- X_1, \dots, X_k neodvisne spremenljivke,
- β_0 presečišče (intercept),
- β_1, \dots, β_k regresijski koeficienti, ki kažejo na spremembo v Y pri enotni spremembi X_i ,
- ϵ napaka modela.

Izračun koeficientov β_i Regresijski koeficienti β_i se izračunajo z metodo najmanjših kvadratov, ki minimizira vsoto kvadratov odstopanj predvidenih vrednosti od dejanskih vrednosti:

$$\hat{\beta} = (X'X)^{-1}X'Y$$

kjer je X matrica neodvisnih spremenljivk, Y pa vektor odvisne spremenljivke.

Koeficient determinacije R^2 Delež variabilnosti, ki ga model pojasni, je izražen z R^2 , ki ga izračunamo kot:

$$R^2 = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2},$$

kjer:

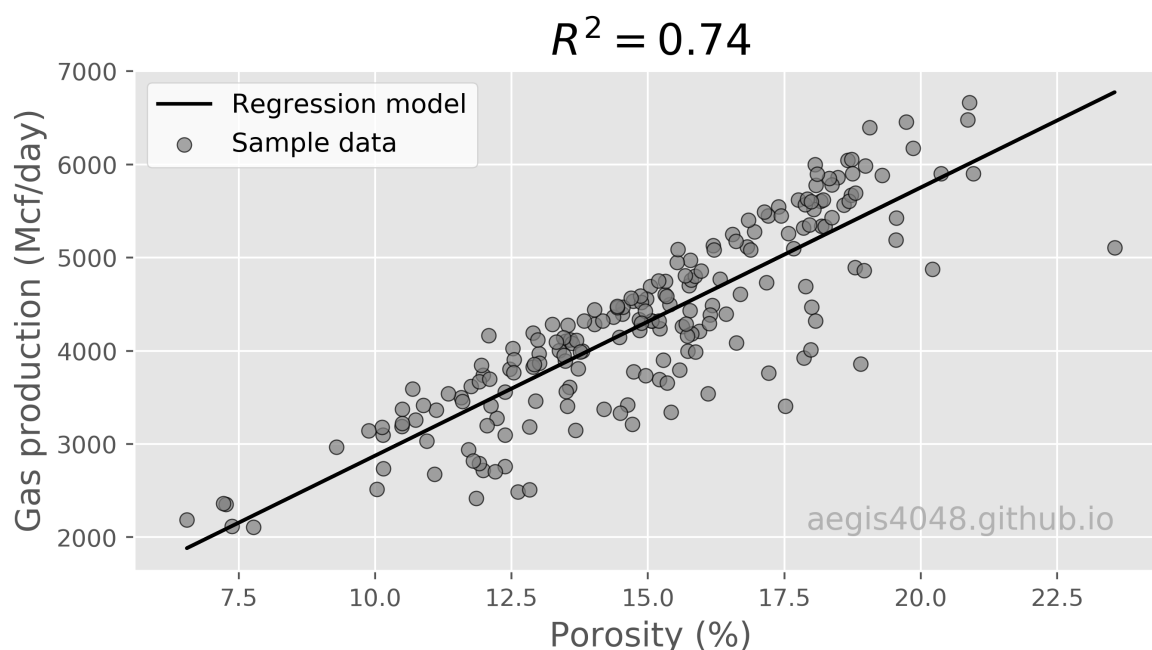
- \hat{Y}_i so napovedane vrednosti odvisne spremenljivke,
- Y_i so dejanske vrednosti odvisne spremenljivke,
- \bar{Y} je povprečje odvisne spremenljivke.

R^2 meri delež skupne variabilnosti v odvisni spremenljivki, ki jo je mogoče pojasniti z regresijskim modelom. Vrednost R^2 se giblje med 0 in 1, pri čemer 1 pomeni popolno prileganje modela podatkom.

Pomembno je, da velikost vzorca zadostuje za stabilnost ocenjenih koeficientov. Priporočeno je, da je velikost vzorca 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.



Slika 4.1: 2D regresijski model, ki ga bomo seveda spremenili v avtorsko delo.

4.2 Multipla regresijska analiza (logistična)

Logistična regresija se uporablja, kadar je odvisna spremenljivka nominalna (ali binarna). 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.

Izračun koeficientov β_i Regresijski koeficienti β_i se izračunajo z metodo največjega verjetja (Maximum Likelihood Estimation, MLE), ki išče koeficiente, ki maksimizirajo verjetnost opazovanih podatkov glede na model.

Interpretacija koeficientov β_i Koeficient β_i predstavlja spremembo v $\text{logit}(p)$ (logaritmu verjetnostnega razmerja) za enoto spremembe v X_i . Ekspontentiranje koeficientov daje t.i. *odds ratio* (razmerje verjetnosti):

$$\text{OR} = e^{\beta_i}$$

Odnos razmerja verjetnosti razlaga, kako se verjetnost pojavitve dogodka spremeni z enoto spremembe v X_i . Če je $\text{OR} > 1$, se verjetnost dogodka poveča, če je $\text{OR} < 1$, se verjetnost zmanjša.

Pseudo- R^2 Za oceno kakovosti prileganja modela se v logistični regresiji pogosto uporablja t.i. pseudo- R^2 . Ena izmed pogosto uporabljenih formul je McFaddenov pseudo- R^2 :

$$R_{\text{McFadden}}^2 = 1 - \frac{\log L_{\text{model}}}{\log L_{\text{null}}}$$

kjer je:

- $\log L_{\text{model}}$ logaritmična verjetnost ustreznosti modela,
- $\log L_{\text{null}}$ logaritmična verjetnost modela brez neodvisnih spremenljivk (samo z intercep-
tom).

McFaddenov pseudo- R^2 se giblje med 0 in 1, pri čemer višje vrednosti nakazujejo boljše prileganje modela.

Predpostavke logistične regresije

- **Neodvisnost opazovanj:** Vsako opazovanje mora biti neodvisno od drugih.
- **Neprekinjenost neodvisnih spremenljivk:** Neodvisne spremenljivke morajo biti kontinuirane ali binarne.
- **Linearnost logit transformacije:** Obstajati mora linearna povezava med $\text{logit}(p)$ in neodvisnimi spremenljivkami.

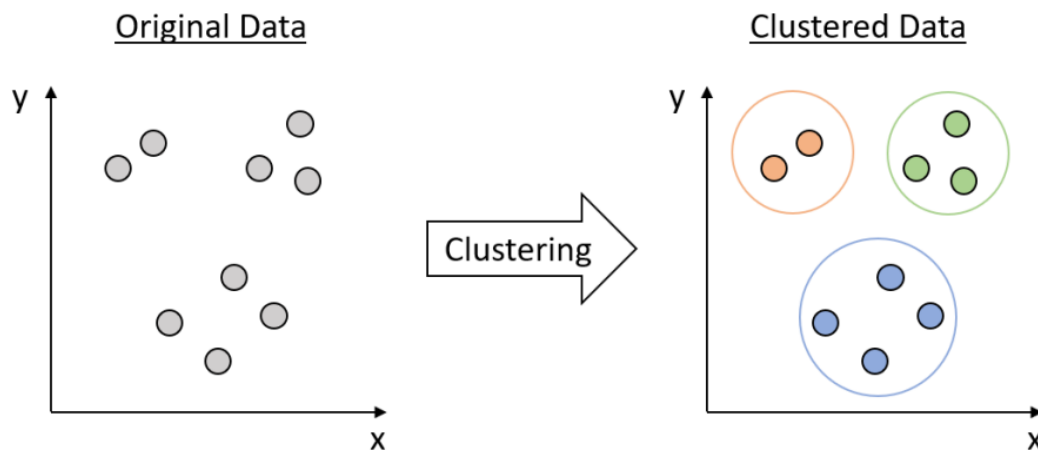
4.3 Razvrščanje v skupine (clustering)

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*.

K-means clustering *K-means* je ena najpogostejše uporabljenih metod za razvrščanje v skupine. Cilj algoritma je razdeliti podatke na k skupin tako, da se minimizira vsota kvadratov odstopanj podatkovnih točk od njihovega pripadajočega centroida (povprečne vrednosti v skupini). Algoritem iterativno prilagaja centriodi do konvergence:

$$\min \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2,$$

kjer je C_i skupina in μ_i centroid skupine i .



Slika 4.2: clustering, ki ga bomo seveda spremenili v avtorsko delo.

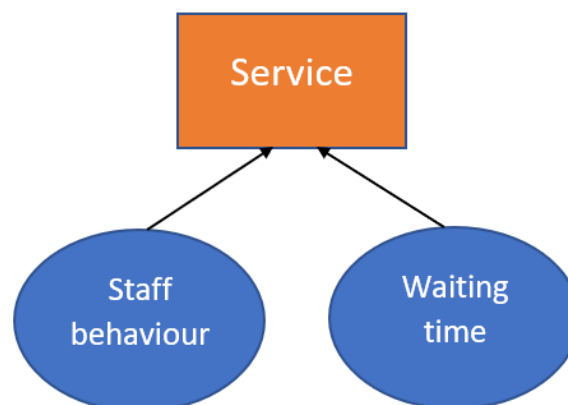
4.4 Metode zmanjšanja dimenzionalnosti podatkov

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. Model faktorske analize predpostavlja, da vsako opazovano spremenljivko X_i lahko izrazimo kot linearno kombinacijo faktorjev F_j in napake ϵ_i :

$$X_i = \lambda_{i1}F_1 + \lambda_{i2}F_2 + \dots + \lambda_{im}F_m + \epsilon_i,$$

kjer so λ_{ij} faktorski uteži, ki kažejo na vpliv faktorja F_j na spremenljivko X_i . Faktorska analiza se uporablja za zmanjšanje števila spremenljivk, hkrati pa ohranja interpretabilne latentne strukture.



Slika 4.3: faktorska analiza.

4.5 Analiza zanesljivosti

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

Formula za izračun Cronbachovega α koeficienta

Cronbachov α se izračuna s pomočjo naslednje formule:

$$\alpha = \frac{k}{k-1} \left(1 - \frac{\sum_{i=1}^k \sigma_i^2}{\sigma_{\text{skupaj}}^2} \right),$$

kjer:

- k je število postavk (vprašanj) v lestvici,
- σ_i^2 je varianca posamezne postavke,
- σ_{skupaj}^2 je varianca celotne lestvice (skupne ocene).

Interpretacija Cronbachovega α

- $\alpha \geq 0.9$: Odlična (zelo visoka) zanesljivost
- $0.8 \leq \alpha < 0.9$: Dobra zanesljivost
- $0.7 \leq \alpha < 0.8$: Sprejemljiva zanesljivost
- $0.6 \leq \alpha < 0.7$: Vprašljiva zanesljivost
- $0.5 \leq \alpha < 0.6$: Slaba zanesljivost
- $\alpha < 0.5$: Nesprejemljiva zanesljivost

Pomembne točke

- **Uporaba Cronbachovega α :** Najpogosteje se uporablja v psihometriji, pri preučevanju zanesljivosti vprašalnikov in testov. Pomembno je, da je lestvica homogena, kar pomeni, da vse postavke merijo en sam koncept.
- **Omejitve Cronbachovega α :** Čeprav je α priljubljena mera, ni popolna. Visoka vrednost α ne pomeni nujno, da lestvica meri tisto, kar namerava meriti (validnost). Poleg tega lahko visoka vrednost α nastane tudi zaradi velikega števila postavk, ne pa nujno zaradi dejanske notranje konsistence.
- **Alternativne mere:** V nekaterih primerih se uporabljajo tudi druge mere zanesljivosti, kot je McDonaldov ω , ki lahko poda natančnejšo sliko zanesljivosti v primerih, ko so postavke lestvice med seboj različno povezane.

4.6 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

V datoteki ... združi stolpce ... v eno spremenljivko ... in določi notranjo konsistentnost vprašalnika.