

Projektna naloga pri statistiki

Gaja Jamnik

1 Prva naloga

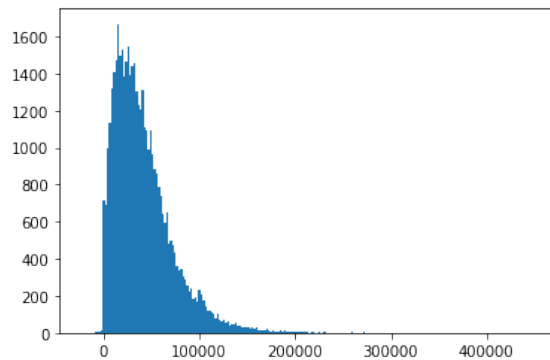
Pri prvi nalogi smo obravnavali podatke 43.886 družin iz mesta Kibergrad. Natančneje, analizirali smo dohodke družin in njihovo porazdelitev primerjali s ustrezno normalno porazdelitvijo.

1.1 A del naloge

Narisali smo histogram dohodkov vseh družin v Kibergradu. Širino in število razredov v histogramu smo določili s pomočjo formule

$$l = \frac{2(q_{3/4} - q_{1/4})}{\sqrt[3]{n}}, \quad (1)$$

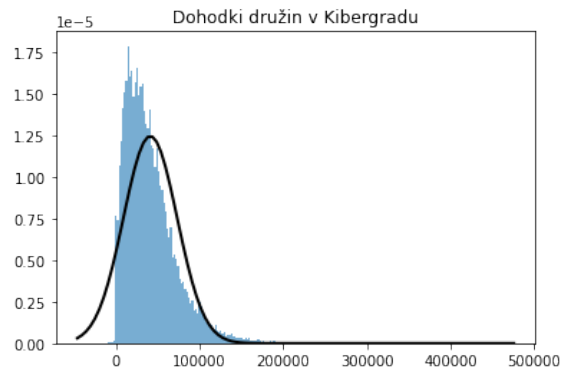
kjer sta $q_{1/4}$ in $q_{3/4}$ prvi in tretji kvartil, n pa je število enot. Tej formuli pravimo Freedman–Diaconisovo pravilo. Po formuli znaša širina posameznega razreda približno 2128. Glede na interval na katerem so razporejeni dohodki znaša to ravno 223 razredov v histogramu.



Slika 1: Histogram dohodkov družin v Kibergradu

1.2 B del naloge

V B delu naloge smo v histogram dorisali še normalno gostoto, katere pričakovana vrednost se ujema s povprečjem dohodkov družin v Kibergradu, standardni odklon pa s standardnim odklonom podatkov o dohodkih družin. Izračunano povprečje podatkov je 41335.50704096979, standardni odklon pa 32037.61941788666.

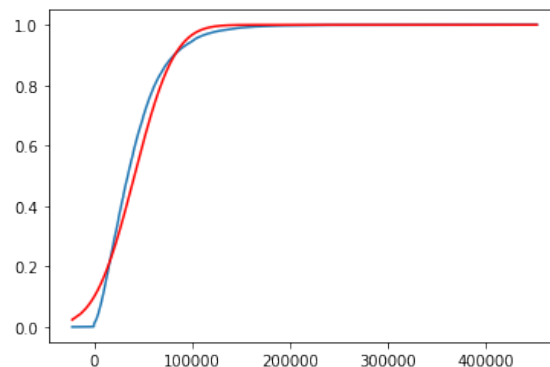


Slika 2: Histogram dohodkov družin v Kibergradu in normalna gostota

Histogram je desno asimetričen. V primerjavi z grafom normalne gostote je zamaknjen v desno, po obliki se desni del histograma lepo prilega krivulji.

1.3 C del naloge

V tem delu naloge primerjamo kumulativni porazdelitveni funkciji podatkov o dohodkih in ustrezne normalne porazdelitve, ki ima za pričakovano vrednost povprečje dohodkov družin v Kibergradu in standardni odklon, ki se ujema s standardnim odklonom dohodkov. Kumulativno porazdelitveno funkcijo porazdelitve dohodkov definiramo s pomočjo funkcije `arange`, tako, da definiramo seznam deležev. Za normalno porazdelitev pa uporabimo kar že vgrajeno funkcijo `statistics.NormalDist(mean, std).cdf(x)`.

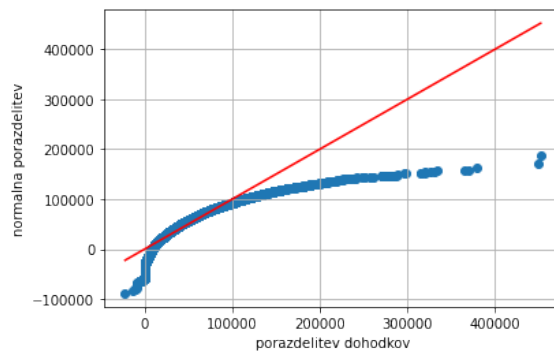


Slika 3: Primerjava kumulativnih porazdelitvenih funkcij

Grafa se od vrednosti 150000 skoraj popolnoma prilegata. Na levi strani pa je nekaj manjših odstopanj.

1.4 D del naloge

S primerjalnim kvantilnim grafikonom oz. Q-Q grafikonom smo primerjali kvantile normalne porazdelitve s kvantili porazdelitve dohodkov družin v Kibergradu.

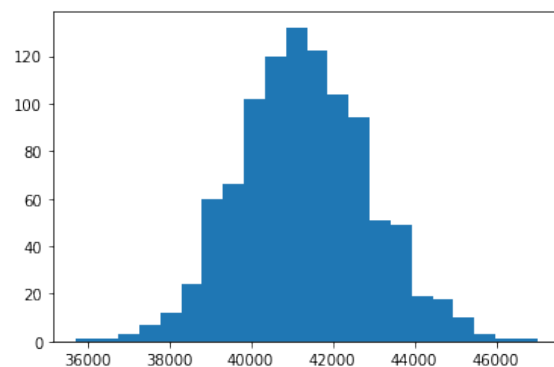


Slika 4: Primerjalni kvantilni grafikon

Če bi bili dohodki normalno porazdeljeni, bi kvantilni grafikon, narisani z modro, sovpadal s 45-stopinjsko premico, narisano z rdečo barvo. Iz grafikona razberemo, da dohodki niso porazdeljeni normalno.

1.5 E del naloge

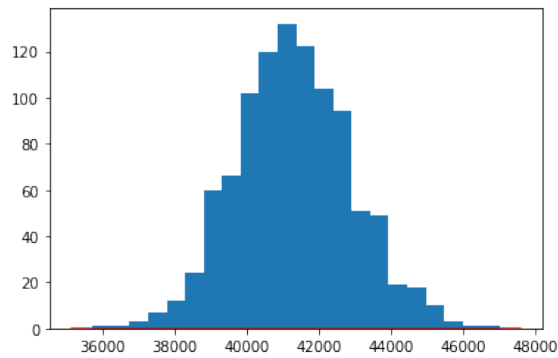
Nato smo analizirali porazdelitev dohodkov z enostavnim slučajnim vzorčenjem. Izbrali smo 1000 enostavnih slučajnih vzorcev velikosti 400. To smo naredili s pomočjo vgrajene funkcije `random.randint`. Nato smo narisali histogram vzorčnih povprečij dohodkov.



Slika 5: Histogram vzorčnih povprečij

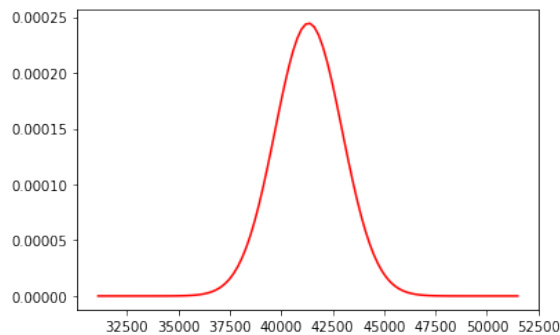
1.6 F del naloge

Na histogram smo dorisali graf normalne gostote. Za pričakovano vrednost smo vzeli povprečje dohodkov celotne populacije, za standardni odklon pa standardno napako nekega slučajnega vzorca velikosti 400. Vzorec smo ponovno generirali s pomočjo funkcije `random.randint`. Standardno napako, pa smo glede na vzorec izračunali s funkcijo `sem()`.



Slika 6: Histogram vzorčnih povprečij z normalno gostoto

Graf gostote normalne porazdelitve je narisana z rdečo barvo. Graf se histogramu ne prilega.. Za primerjavo sem še posebej narisala graf dane normalne porazdelitve.

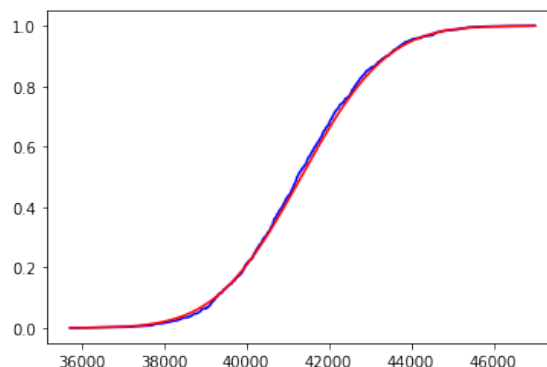


Slika 7: Gostota normalne porazdelitve

Opazimo lahko, da imata graf in histogram vrh pri približno enaki vrednosti abscise, vendar je vrh normalne gostote znatno nižji.

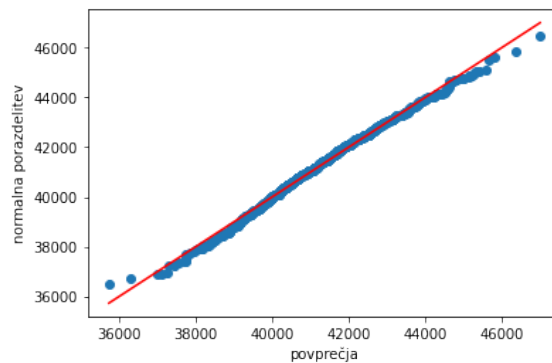
1.7 G del naloge

Tudi v tem delu naloge smo uporabili normalno porazdelitev s pričakovano vrednostjo, ki se ujema populacijskem povprečju dohodkov in standardnim odklonom, ki se ujema s standardno napako za enostavni slučajni vzorec velikosti 400. Za vzorčna povprečja smo narisali kumulativno porazdelitveno funkcijo. V isti graf smo za primerjavo z rdečo barvo narisali še kumulativno porazdelitveno funkcijo zgoraj opisane normalne porazdelitve. Iz grafa 8 opazimo, da se funkciji dobro prilegata.



Slika 8: Kumulativna porazdelitvena funkcija za vzorčna povprečja (z modro) in kumulativna porazdelitvena funkcija normalne porazdelitve (z rdečo)

Narisali smo še primerjalni kvantilni grafikon s katerim smo primerjali porazdelitev vzorčnih povprečij z normalno porazdelitvijo.



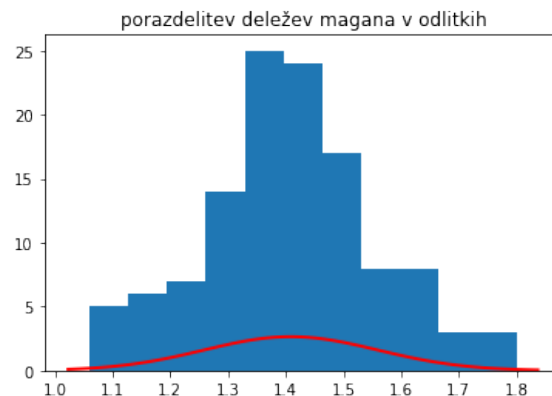
Slika 9: Kvantilni grafikon vzorčnega povprečja v primerjavi z normalno porazdelitvijo

Na grafu 9 se kvantilni grafikon, narisan z modro, dobro prilega rdeči 45-stopinjski premici. To pomeni, da ima vzorčno povprečje normalno porazdelitev.

2 Naloga 2

V drugi nalogi smo analizirali podatke o deležu mangana v železu, pridobljenem v plavžu. Podatki so vsebovali deleže mangana v odlitkih, ki so jih jemali skozi 24 dni, petkrat dnevno. Preučiti smo želeli normalnost dobljene empirične porazdelitve. Porazdelitev deležev mangana smo primerjali z normalno porazdelitvijo, katere pričakovana vrednost in standardni odklon se ujemata s povprečjem in standardnim odklonom deležev mangana iz raziskave.

Najprej smo narisali histogram vseh deležev z dorisano normalno gostoto. Širino posameznega razreda v histogramu smo kot v nalogi 1 izračunali s pomočjo Freedman-Diaconisovega pravila. Širina razreda je znašala 0.07298642394688079, kar je ravno 11 razredov.



Slika 10: Histogram deležev mangana v plavžu z dorisano normalno gostoto (z rdečo)

Graf normalne gostote se lepo prilega histogramu. Vrh histograma in grafa sta dosežena pri približno enaki vrednosti na abscisi. Histogram je tudi zvončaste oblike, kar namiguje na normalno porazdelitev.

Vendar to opažanje ni dovolj, saj normalna porazdelitev ni edina z zvončasto obliko histograma. Zato smo za primerjavo uporabili še viseči histogram razlik korenov frekvenc. Klasičen histogram ponazori frekvence posameznih vrednosti iz podatkov. Viseči histogram pa primerja korene frekvenc dane porazdelitve s frekvencami ustrezne normalne porazdelitve. To naredimo tako, da najprej poračunamo verjetnosti, da ima normalno porazdeljena spremenljivka vrednost na danih intervalih. Ocena verjetnosti na intervalu $[x_{j-1}, x_j]$ bo

$$\hat{p}_j = \Phi\left(\frac{x_j - \bar{x}}{\hat{\sigma}}\right) - \Phi\left(\frac{x_{j-1} - \bar{x}}{\hat{\sigma}}\right), \quad (2)$$

kjer je \bar{x} povprečje naših podatkov, $\hat{\sigma}$ pa ocena za standardni odklon.

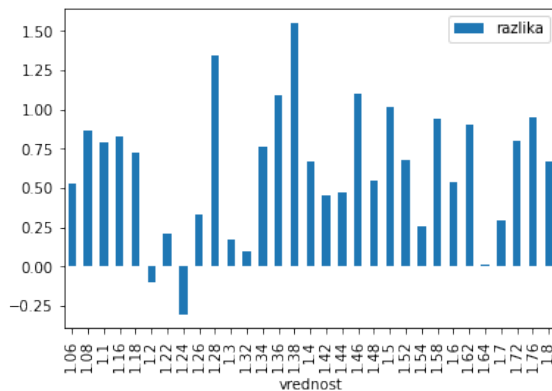
Nato določimo predvideno število pojavitev na j -tem intervalu kot

$$\hat{n}_j = np_j, \quad (3)$$

kjer je n število podatkov o deležih mangana. V zadnjem koraku poračunamo razliko korena števila ponovitev vrednosti na j -tem interval, kar označimo z n_j s korenem ocenjenih ponovitev normalne porazdelitve:

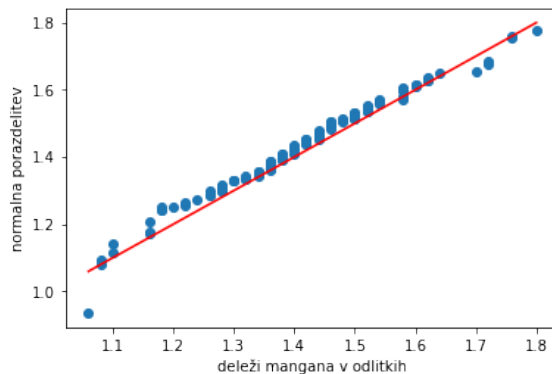
$$\sqrt{n_j} - \sqrt{\hat{n}_j}. \quad (4)$$

Poračunane vrednosti narišemo v histogram. Iz grafa 11 vidimo, da imamo vzdolž vseh vrednosti velika odstopanja v frekvencah.



Slika 11: Viseči histogram razlik korenov frekvenc

Za nadaljno primerjavo uporabimo še kvantilni (Q-Q) grafikon.



Slika 12: Primerjalni kvantilni grafikon

Iz grafikona 12 opazimo, da se kvantili prilegajo 45-stopinjski premici, na robovih pa pride do rahlih odstopanj. Kvantilni grafikon kaže, da so deleži mangana porazdeljeni približno normalno.

3 Naloga 3

V tretji nalogi smo analizirali podatke o dolžini zob morskih prašičkov, ki so jim dodajali vitamin C v različnih količinah na dva različna načina: bodisi neposredno bodisi s pomarančnim sokom.

3.1 A del naloge

Najprej nas je zanimalo ali dodajanje vitamina C sploh vpliva na rast zob. Za začetek smo izračunali povprečno dolžino zobu, glede na količino dodatka vitamina C. Za količino 0.5 je povprečna dolžina zoba 10.605, za količino 1.0 je 19.735, za količino 2.0 pa 26.100. To nam da misliti,

da dolžina zobu raste z količino dodatka vitamina C. Preizkusimo to še s preizkusom hipoteze. Definirajmo ničelno in alternativno hipotezo:

$$H0 : \mu_{0.5} \geq \mu_{1.0} \quad (5)$$

$$H1 : \mu_{0.5} < \mu_{1.0}. \quad (6)$$

Hipotezo preizkusimo s Studentovim t-testom. Izberimo stopnjo tveganja $\alpha = 0.05$. Naj slučajna spremenljivka X označuje porazdelitev velikosti zob pri dodajanju količine 0.5 vitamina C, spremenljivka Y pa pri dodajanju količine 1.0.

Na vajah smo pokazali, da za spremenljivko T velja

$$T := \frac{\bar{X} - \bar{Y}}{S} \sqrt{\frac{mn}{m+n}} \sim Student(n+m-2),$$

kjer je

$$S := \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^m (Y_i - \bar{Y})^2}{m+n-2}}.$$

V našem primeru je $n = m = 30$. Iz formule

$$P(|\bar{M} - \bar{Z}| \leq C_\alpha) \leq \alpha$$

nato poračunamo C_α . Iz vzorcev izračunamo povprečji \bar{X} in \bar{Y} , ter nato na podlagi teh podatkov in C_α izračunamo zgornjo verjetnost p . Če je $p < \alpha$ ovržemo ničelno hipotezo, sicer jo obdržimo. Verjetnost smo izračunali kar z že vgrajeno funkcijo `stats.ttest_ind()`. Dobimo vrednost $p = 6.331484806608257e - 08$, kar je manjše kot α , zato zavržemo ničelno hipotezo.

Podobno storimo še s primerjavo pričakovane vrednosti dolžin zob z dodajanjem količine 1.0 in 2.0. Z enakim postopkom dobimo $p = 9.054142680908653e - 06$, kar je manj kot α , zato bo veljalo, da je $\mu_{1.0} < \mu_{2.0}$. Od tod sledi, da res dodajanje vitamina C pozitivno vpliva na rast zob.

3.2 B del naloge

Dalje nas je zanimalo, kateri način dodajanja vitamina C je učinkovitejši. Ponovno naredimo preizkus hipoteze s Studentovim t-testom. Naj P označuje slučajno spremenljivko za dolžino zob pri dodajanju s pomarančnim sokom, N pa pri neposrednem dodajanju. Zapišimo hipotezi:

$$H0 : \mu_P \leq \mu_N \quad (7)$$

$$H1 : \mu_P > \mu_N \quad (8)$$

Ponovno uporabimo vgrajeno funkcijo in dobimo $p = 0.03019668561206424$, kar je manjše kot $\alpha = 0.05$. Torej ničelno hipotezo ponovno zavržemo. Sledi, dodajanje s pomarančnim sokom je učinkovitejše. Vendar bi za stopnjo tveganja $\alpha = 0.01$ ničelno hipotezo držali, in trdili bi nasprotno.