

UVOD

Namen raziskave je bil sistematično preučiti ekstremalne grafe glede na število geodetskih podpoti (angl. *geodesic path number*, okrajšano *gpn*) za različne strukturne razrede povezanih grafov. Specifično nas je zanimala identifikacija grafov, ki v danem razredu dosežejo maksimalno možno vrednost *gpn* za fiksno število vozlišč n .

Reševanja problema smo se lotili v treh glavnih korakih. V prvem koraku smo izvedli analizo števila geodetskih podpoti vseh enostavnih povezanih grafov z manj kot desetimi vozlišči. V drugem koraku smo na podlagi analize postavili hipotezo o tem, kateri grafi so optimalni za merilo dano *gpn*. V zadnjem koraku smo postavljeno hipotezo preverili z uporabo hevrističnih metod, za grafe z večjim številom vozlišč.

TEORETIČNO OZADJE

Za natančno obravnavo problema si najprej pogledjmo naslednje definicije:

Definicija Pot grafa

Naj bo G povezan, enostaven graf. Pot dolžine ℓ , je zaporedje vozlišč $(v_0, v_1, \dots, v_\ell)$ brez ponovitev, kjer je vsak par zaporednih vozlišč povezan s povezavo grafa G , tj:

$$\forall i \in \{1, 2, \dots, \ell\} : v_{i-1}v_i \in E(G)$$

Trivialna pot dolžine 0 vsebuje eno vozlišče.

Definicija Geodetska pot grafa

Naj bo G povezan, enostaven graf, in naj bosta $u, v \in V(G)$ dve vozlišči grafa. Geodetska pot med u in v je pot $(v_0, v_1, \dots, v_\ell)$, ki ima minimalno možno dolžino med vsemi potmi, ki povezujejo u in v .

Dolžino geodetske poti označimo z $d_G(u, v)$ in jo imenujemo *razdalja* med u in v .

Definicija Število geodetskih podpoti

Naj bo G povezan, enostaven graf. Število geodetskih podpoti, označeno $gpn(G)$, je definirano kot število vseh geodetskih (najkrajših) poti v grafu, vključno s trivialnimi potmi:

$$gpn(G) = \sum_{u, v \in V(G)} \text{število geodetskih poti med } u \text{ in } v$$

Torej za vsak par različnih vozlišč obstaja vsaj ena geodetska pot, ter za vsako vozlišče obstaja natanko ena trivialna pot.

V raziskavi smo se osredotočili na naslednje strukturne razrede povezanih grafov:

Definicija Dvodelni graf

Graf $G = (V, E)$ je *dvodelen*, če obstaja razcep (particija) množice vozlišč $V = A \cup B$, $A \cap B = \emptyset$, tako da za vsako povezavo $uv \in E(G)$ velja:

$$\forall u, v \in V(G) : (u \in A \wedge v \in B) \vee (u \in B \wedge v \in A)$$

Torej, vse povezave povezujejo vozlišča iz različnih particij.

Definicija Graf brez trikotnikov

Graf G je *brez trikotnikov* (angl. triangle-free graph), če v njem ne obstaja cikel dolžine 3, torej če graf G ne vsebuje podgrafa, izomorfnega ciklu C_3 .

Trditev Dvodelni grafi so brez trikotnikov

Vsak dvodelni graf je hkrati graf brez trikotnikov. Obratno ne velja.

Definicija Kubičen graf

Graf $G(V, E)$ je *kubičen* (3-regularen), če ima vsako vozlišče stopnjo 3:

$$\forall v \in V(G) : \deg(v) = 3$$

Torej vsako vozlišče je sosednje natanko trem vozliščem. Število povezav je enako:

$$|E(G)| = \frac{3|V(G)|}{2}$$

Trditev Obstoj kubičnih grafov

Kubični grafi lahko obstajajo le na sodem številu vozlišč.

Definicija Ekstremalen graf

Naj bo \mathcal{G} razred grafov in $P(G)$ neka invarianta (neodvisna lastnost grafa). Graf G^* je ekstremalen glede na P , če velja:

$$G^* = \arg \max_{G \in \mathcal{G}} P(G)$$

ko želimo maksimizirati P oziroma:

$$G^* = \arg \min_{G \in \mathcal{G}} P(G)$$

ko želimo minimizirati P .

DEFINICIJA PROBLEMA

V naši študiji smo iskali grafe, ki maksimizirajo invarianto $gpn(G)$ znotraj razreda enostavnih povezanih označenih (angl. *labeled*) grafov na fiksnem številu vozlišč. Slednji problem lahko definiramo z naslednjim celoštevilskim linearnim programom:

Celoštevilski Linearni Program za maksimizacijo geodetskega števila:

Naj bo $G \in \mathcal{G}(V, E)$ neusmerjen neutežen povezan enostaven označen graf nad n -timi vozlišči. Označimo $V(G) = \{1, \dots, n\}$ in $E(G) \subseteq \{e_{ij} \mid \forall i, j \in V(G) : e_{ij} = (i, j)\}$. Najti želimo graf z največjim geodetskim številom.

$$\text{ciljna funkcija: } \max \left(\sum_{i=1}^n \sum_{j=1}^n y_{ij} \right)$$

kjer :

1. definiramo vse matrike sosednosti za enostavne in povezane grafe

- $\forall i, j \in \{1, \dots, n\} : x_{ij} = \begin{cases} 1 & ; e_{ij} \in E(G) \\ 0 & ; \text{sicer} \end{cases} \quad (\text{sestavimo matriko})$
- $\sum_{i < j} x_{ij} \geq n - 1 \quad (\text{graf je povezan})$
- $\forall i, j \in \{1, \dots, n\} : x_{ij} = x_{ji} \quad (\text{graf je neusmerjen})$
- $\forall i \in \{1, \dots, n\} : x_{ii} = 0 \quad (\text{brez zank})$
- $\forall i, j \in \{1, \dots, n\} : x_{ij} \in \{0, 1\} \quad (\text{brez večkratnih povezav})$

2. definiramo funkcijsko odvisno spremenljivko od matrike sosednosti

- $\forall i, j \in \{1, \dots, n\} : y_{ij} = \{\text{število najkrajših poti med vozliščema } i \text{ in } j\}$
- $\forall i, j \in \{1, \dots, n\}, i < j : y_{ij} = 0 \quad (\text{vsako pot štejemo le enkrat})$
- $\forall i \in \{1, \dots, n\} : y_{ii} = 1 \quad (\text{trivialne poti})$
- $\forall i, j \in \{1, \dots, n\} : y_{ij} \in \mathbb{N} \cup \{0\}$

Če je podan CLP dopusten za izbrano število vozlišč n in je $(x_{ij}^*)_{i,j=1}^n$ optimalna resitev, ki predstavlja graf G^* , je potem to rešitev problema. Pri tem je potrebno omeniti, da lahko za dani n obstaja več neizomorfnih grafov, ki dosežejo optimum.

Eksplíciten primer reševanja programa je med prilogami. V programu smo predpostavili, da poznamo število najkrajših poti med poljubnima vozliščema, kar pa v splošnem ni trivialen problem, ki ga lahko rešimo v konstatnem času. Potreben algoritem za izračun števila najkrajših poti, ter invariante gpn je prav tako med prilogami.

METODOLOGIJA REŠEVANJA PROBLEMA

Problema smo se lotili na sledeč način:

Razvojno okolje (angl. *development environment*):

Ker je iskalni prostor (angl. *search space*) vseh enostavnih povezanih označenih grafov izjemno velik, število takih grafov za naraščajoče število vozlišč sledi zaporedju:

1, 1, 2, 6, 21, 112, 853, 11117, 261080, 11716571, ...

smo se naloge seveda lotili s pomočjo programskih jezikov in računalniške moči. Zaradi številnih sofisticiranih implementacij programske opreme za delo z grafi, smo se odločili, da uporabimo razvojno okolje **SageMath**. Slednjega smo uspostavili s sistemom Docker, saj je namestitev neposredno na določene operacijske sisteme kot je Windows brez WSL-ja (Windows Subsystem for Linux) nepraktična. Torej, s pomočjo Docker-ja smo vzpostavili virtualni vsebnik (angl. *container*), v katerem smo izvedli celotno analizo. Z virtualizacijo smo si zagotovili izolacijo okolja od gostiteljevega operacijskega sistema in enotno delovno okolje na različnih platformah. Artifakt Docker-jeve slike lahko najdete *tukaj*.

Generacija podatkov:

V prvem eksperimentalnem koraku smo morali definirati objekte iskalnega razreda. Te smo za manjše grafe generirali lokalno sami, za večje grafe pa smo se poslužili te *zbirke grafov*.

V večji meri smo uporabili SageMath-ov modul *nauty*, natančneje vmesnik *nauty_geng* programa **geng**, za učinkovito generiranje neizomorfni grafov z določenimi lastnostmi. V prvotni fazi smo se osredotočili predvsem na tri strukturne razrede: kubične grafe, dvodelne grafe in grafe brez trikotnikov. Za vsakega izmed teh razredov smo v programskem jeziku Python implementirali ustrezne razrede (*OOP class*), ki omogočajo generiranje vseh možnih grafov na determinističnem številu vozlišč. Omembe vredno je, da vsak razred vključuje lastnosti grafov ter možnost vizualizacije posameznega objekta iz razreda. Za enostavnejšo implementacijo smo uporabili knjižnici *networkx* in *matplotlib*. Izvorna koda se nahaja *tukaj*.

Podrobneje, lokalno smo generirali vse neizomorfne enostavne povezane označene grafe na do vključno devetih vozliščih. Za grafe z desetimi vozlišči, smo zaradi omejenih računalniških resursov, podatke naložili. Pri zgoraj omenjenih specifičnih strukturnih razredih, nismo imeli težav z generacijo, saj so te skupine bistveno manjše (za grafe $|V(G)| \geq 10$). Izvorna koda je dostopna *tukaj*.

Vsak graf smo enkodirali z **graph6** identifikacijskim ključem, saj ta enolično določa vsak enostaven graf znotraj našega iskalnega prostora.

Izračun invariante:

Algoritem za izračun števila geodetskih podpoti je priložen med prilogami, izvorna koda pa se nahaja *tukaj*.

Podatke smo shranili v datoteke formata CSV (Comma Separated Values), ki smo jih po potrebi stisnili zaradi omejitev velikosti datotek na GitHub-u. Pri tem smo uporabili poznano knjižnico za delo s podatki *pandas* ter interaktivno računalniško okolje *Jupyter Notebook*.

Za olajšan pristop do podatkovnih zbirk so spodaj priložene povezave:

- *link*: zbirka vseh enostavnih povezanih grafov na do vključno devetih vozliščih
- *link*: zbirka vseh enostavnih povezanih grafov na desetih vozliščih
- *link*: zbirka vseh grafov iz zgoraj omenjenih specifičnih strukturnih razredov
- *link*: zbirka vseh dvodelnih grafov na 11 vozliščih
- *link*: zbirka vseh dvodelnih grafov na 12 vozliščih

Prečno iskanje (angl. *exhaustive search*):

Na podlagi prvih treh opisanih podatkovnih zbirk smo iz razredov izločili zgolj *ekstremalne grafe* glede na invarianto *gpn*. Pri obdelavi in filtriranju podatkov smo ponovno uporabili knjižnico *pandas* ter sodobnejšo knjižnico *polars*, razvito v okviru rustovega ekosistema. Glavna prednost slednje je podpora lenega nalaganja (angl. *lazy loading*), ki omogoča učinkovito filtriranje velikih CSV-datotek brez nepotrebne porabe virov.

V spodnji tabeli so zbrani ekstremalni grafi za posamezen razred glede na število vozlišč:

graph6 ID	$ V(G) $	$ E(G) $	brez trikotnikov	dvodelen	kubičen	$gpn(G)$
@	1	0	ja	ja	ne	1
A_	2	1	ja	ja	ne	3
BW	3	2	ja	ja	ne	6
Bw	3	3	ne	ne	ne	6
C]	4	4	ja	ja	ne	12
DFw	5	6	ja	ja	ne	20
EFz_	6	9	ja	ja	ja	33
F?zv_	7	11	ja	ja	ne	49
F?v_	7	12	ja	ja	ne	49
G?zVf_	8	14	ja	ja	ne	74
G?zvf_	8	15	ja	ja	ne	74
H?BvUrw	9	17	ja	ja	ne	105
H?BvVrw	9	18	ja	ja	ne	105
H?ovfbo	9	16	ja	ja	ne	105
I?BvUqw}?	10	21	ja	ja	ne	151

Prikazi zgornjih grafov so priloženi v prilogi. Na podlagi tabele lahko podamo hipotezo:

Hipoteza Ekstremalni grafi za *gpn*

Za razred enostavnih povezanih grafov nad n -timi vozlišči so ekstremalni grafi dvodelni.

Slednjo lahko potrdimo že z diagramoma v prilogi, ki sta bila ustvarjena v okolju programskega jezika R z uporabo paketov iz skupine *tidyverse* in paketa *ggplot2* za vizualizacije. Izvorna koda se nahaja *tukaj*.

Na podlagi hipoteze, smo prečno preiskali (angl. *brute-force search*) vse dvodelne grafe na enajstih in dvanajstih vozliščih in dobili naslednje rezultate o ekstremalnih grafih:

graph6 ID	$ V(G) $	$ E(G) $	brez trikotnikov	dvodelen	kubičen	$gpn(G)$
$J??E@w\{Fo?$	11	20	ja	ja	ne	209
$K??E@w\{Fo^?$	12	25	ja	ja	ne	303

Z dodatnim opazovanjem strukture optimalnih grafov, si lahko še dodatno skrčimo iskalni prostor, saj opazimo da so grafi sistematično dvodelni:

Hipoteza Ekstremalni grafi za gpn

Največje število geodetskih podpoti za razred enostavnih povezanih grafov nad n -timi vozlišči dosežejo **uravnoteženi dvodelni** in precej regularni grafi, kar pomeni da je particija vozlišč čim bolj enakomerna:

- n je sod \implies optimalna particija: (k, k)
- n je lih \implies optimalna particija: $(k, k + 1)$

Hevristične metode iskanja:

Na podlagi zgornje hipoteze smo implementirali razred, ki generira naključen uravnotežen dvodelen graf, pri katerem lahko stohastično uravnavamo pričakovano gostoto povezav. Izvorna koda se nahaja *tukaj*.

Ker se je iskalni prostor znatno zmanjšal, smo za iskanje ekstremalnih večjih grafov od 11 do vključno 30 vozlišči uporabili metodo naključnega iskanja (angl. *random search*). Izvorna koda je dostopna *tukaj*. Hevristični rezultati se nahajajo *tukaj*.

Za nadaljnjo optimizacijo smo uporabili knjižnico *simanneal*, ki poenostavi implementacijo simuliranega žarjenja (angl. *simulated annealing, SA*). Začetno stanje smo inicializirali z optimalnim grafom iz prejšnjega koraka. Vsako naslednje stanje smo dobili z majhno perturbacijo grafa, kjer smo z enakomerno porazdeljeno verjetnostjo odstranili, dodali ali spremenili povezavo. Izboljšani rezultati so dostopni *tukaj*.

Testiranje hipoteze:

Hipotezo smo preverili s Studentovim t-testom za dva neodvisna vzorca (angl. *Student's two-sample t-test*) s pomočjo knjižnice *scipy* z uporabo modula *stats*. Za prvi vzorec smo izbrali grafe iz populacije vseh enostavnih povezanih grafov, za drugi pa le grafe, ki ustrezajo razredu, opisanemu v hipotezi. Ničelna hipoteza je predpostavljala, da je povprečje gpn vrednosti uravnoteženih dvodelnih grafov večje od povprečja gpn vrednosti grafov iz splošne populacije. Z uporabo 5-odstotne stopnje napake smo hipotezo potrdili. Izvorna koda je na voljo *tukaj*.

Priloga s primeroma reševanja celoštevilskega linearnega programa:

Za $|V| = 3$ velja:

Obstajata le dve dopustni matriki sosednosti:

$$X_{P_3} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad \text{in} \quad X_{K_3} = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

potem so izpeljane matrike najkrajših poti enake:

$$Y_{P_3} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \quad \text{in} \quad Y_{K_3} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

in so nato vrednosti objektivne funkcije enake:

$$gpn(P_3) = gpn(K_3) = \sum_{i=1}^3 \sum_{j=1}^3 y_{ij} = 6$$

Za $|V| = 4$ velja:

Obstaja 6 dopustih matrik sosednosti za vse neizomorfne grafe:

$$X_{P_4} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}, \quad X_{S_4} = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}, \quad X_{C_4} = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{bmatrix}$$

$$X_{\text{diamant}} = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}, \quad X_{K_4^-} = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix}, \quad X_{K_4} = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

potem so izpeljane matrike najkrajših poti enake:

$$Y_{P_4} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \quad Y_{S_4} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \quad Y_{C_4} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 2 & 1 & 1 & 0 \\ 1 & 2 & 1 & 1 \end{bmatrix}$$

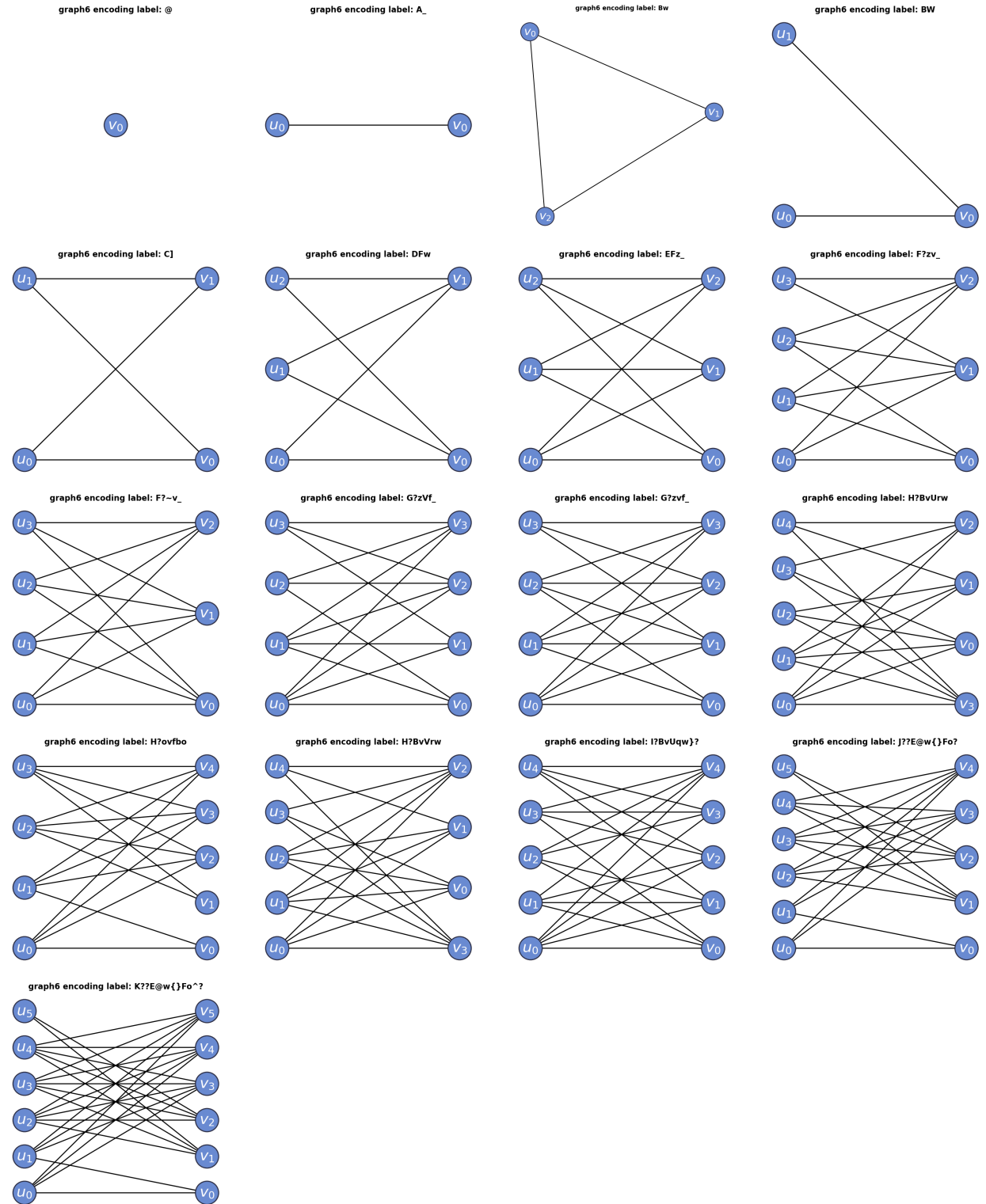
$$Y_{\text{diamant}} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \quad Y_{K_4^-} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \quad Y_{K_4} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

in so nato vrednosti objektivne funkcije enake:

$$gpn(P_4) = gpn(S_4) = gpn(K_{\text{diamant}}) = gpn(K_4) = 10, \quad gpn(K_4^-) = 11, \quad gpn(K_4) = 12$$

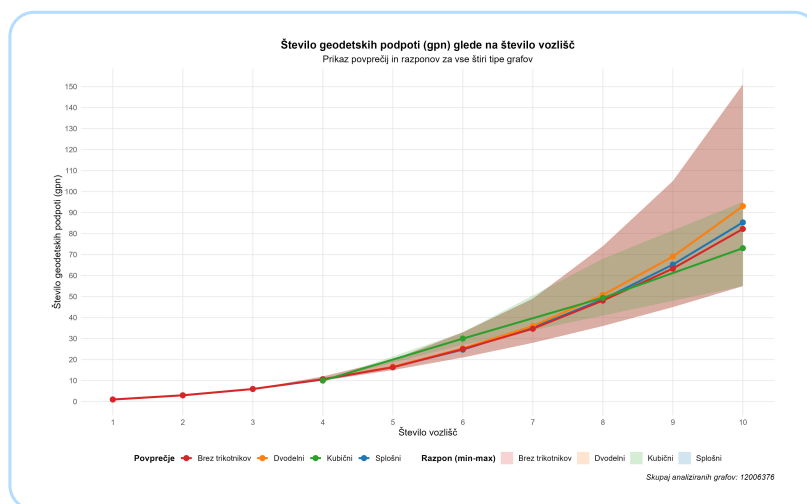
Priloga z algoritmom gpn:

Priloga z ekstremalnimi grafi:

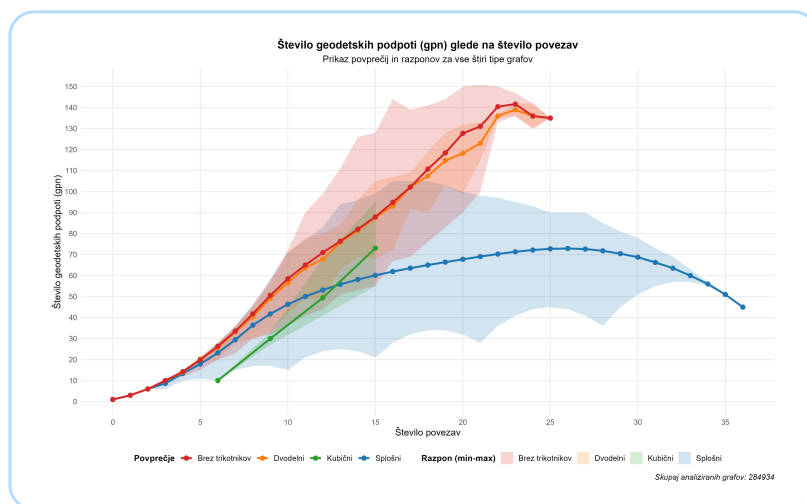


Priloga z diagrami statistične analize:

Za vse obravnavane razrede grafov smo opazili monotono naraščanje števila geodetskih podpoti s številom vozlišč, kar je dobro razvidno na črtnem diagramu:



Na podlagi analize povprečnih vrednosti so se dvodelni grafi izkazali kot optimalni razred, sledili so jim splošni grafi, medtem ko so kubični grafi dosegali najnižje vrednosti *gpn* pri enakem številu vozlišč. Ker je invarianta *gpn* precej odvisna od števila povezav že iz same definicije, smo v naslednjem diagramu prikazali glede na število povezav. Poudariti je treba metodološki vidik analize, saj predstavljeni diagram prikazuje povprečne vrednosti *gpn* glede na število povezav za grafe, ki so bili generirani odvisnosti od števila vozlišč. Torej tukaj nismo obravnavali vseh možnih grafov na *m*-tih povezavah, temveč zgolj grafe z danim številom vozlišč, razvrščenih glede na število povezav.



Iz diagrama je razvidno, da se povprečna vrednost *gpn* z večanjem števila povezav najprej povečuje, doseže maksimum, nato pa začne upadati za splošen razred. Na začetku dodatne povezave ustvarjajo alternativne geodetske poti, kar vrednost *gpn* povečuje. Vendar ko graf doseže preveliko gostoto povezav in se približa polnemu grafu, večina razdalj med vozlišči postane dolžine ena, kar posledično zmanjša število različnih najkrajših poti. Polni graf K_n , podobno kot drevesa, dosega minimalno vrednost *gpn*, tj. $gpn(K_n) = \binom{n}{2} + n$.