

Privacy in dynamic graphs

Guillem Garcia Dausà

Resum— La protecció de dades és un tema cada cop més rellevant en tots els àmbits. En particular, és preocupant en xarxes que generen grans volums de dades i interaccions constants, com ara les de transaccions o comunicacions. Aquestes xarxes sovint es modelen com a grafs, els quals poden evolucionar al llarg del temps amb canvis en els nodes i les relacions. Aquest component temporal comporta nous desafiaments en matèria de privacitat i dificulta la protecció de la informació confidencial. Per aquest motiu, es presenten dos mètodes de protecció de dades basats en els conceptes de *k-anonimitat* i *Privacitat diferencial*. Per avaluar els algorismes implementats, s'utilitzen mètriques de similaritat, utilitat i estructurals, adaptades específicament per comparar grafs temporals. Especialment, es vol analitzar si els grafs protegits preserven les seves característiques quan es detecten comunitats, que s'utilitza *TSCAN* per fer aquesta tasca.

Paraules clau— grafs temporals, protecció, *k-anonimitat*, *Privacitat diferencial*, similaritat, utilitat, estructurals, *TSCAN*

Abstract— Data privacy is becoming increasingly relevant across all fields. In particular, it is a growing concern in networks that generate large volumes of data and constant interactions, such as transaction or communication networks. These networks are often modeled as graphs, which can evolve over time with changes in nodes and relationships. This temporal aspect introduces new privacy challenges and makes it harder to protect sensitive information. For this reason, two data protection methods are presented, based on *k-anonymity* and *Differential privacy*. To evaluate the implemented algorithms, similarity, utility, and structural metrics are used, specially adapted to compare temporal graphs. Specifically, one of the goals is to analyze whether the protected graphs preserve their features when community detection algorithms are applied, using *TSCAN* for this task.

Keywords— temporal graphs, privacy, *k-anonymity*, *Differential privacy*, similarity, utility, structural, *TSCAN*



1 INTRODUCCIÓ I OBJECTIUS

LA protecció de dades s'ha anat convertint en un tema important en l'actualitat, on és un dilema en gairebé tot arreu, com pot ser en el sector de la salut, la tecnologia, les finances, entre altres. És preocupant aquest problema en àmbits on es generen grans volums de dades i interaccions en temps real, com en xarxes de comunicacions o de transaccions, on és fonamental garantir la seguretat i privacitat de la informació.

Aquestes dades sovint es poden modelar mitjançant grafs, on els nodes representen objectes o usuaris, i les arestes defineixen les relacions entre ells dins la xarxa. En moltes ocasions, aquests grafs no són estàtics, sinó dinàmics,

ja que els nodes i les arestes poden anar canviant al llarg del temps. Això planteja nous reptes en matèria de privacitat, pel fet que pot facilitar la re-identificació d'usuaris o l'extracció de dades personals si s'obté informació que va evolucionant. Per aquest motiu, és essencial desenvolupar tècniques adequades per protegir la privacitat tant en grafs dinàmics com en estàtics i mitigar possibles atacs dins la xarxa.

Llavors, l'objectiu principal d'aquest projecte és investigar sobre grafs que varien durant el temps, on es volen entendre els diversos processos per protegir-los, i quines són les propietats que canvien respecte als grafs originals quan s'apliquen els mètodes. Específicament, es volen assolir ordenadament els següents punts:

1. Establir les definicions bàsiques d'un graf temporal i quines propietats addicionals tenen en comparació dels grafs estàtics.
2. Aplicar diversos mètodes de privacitat a diferents conjunts de dades, des de volums de dades fàcils de tractar,

• E-mail de contacte: Guillem.GarciaD@autonoma.cat
 • Treball tutoritzat per: Guillermo Navarro Arribas (Àrea de Ciències de la Computació i Intel·ligència Artificial, UAB)
 • Curs 2024/25

fins a una gran quantitat. Pels mètodes de privacitat, es vol entendre com funcionen, i saber quines són les possibles situacions que poden succeir si una persona ataca la xarxa. Els mètodes que es volen treballar principalment són a partir dels conceptes de *Differential Privacy* i *k-anonymity*, per a l'anonimització d'arestes.

3. Fer una comparativa a escala de privacitat i utilitat dels grafs protegits i originals. La idea és utilitzar diferents mètriques i algorismes que permetin veure les diferències que es produeixen entre grafs. Els mètodes de protecció han de ser els més consistents i òptims possibles, on es vol fer un estudi de quins són les millors opcions per cada *dataset*. Es vol estudiar especialment com es detecten les comunitats en els grafs protegits, i si es preserven les mateixes característiques comparat amb els grafs originals.

Tot i que els objectius estiguin organitzats d'aquesta manera, la màxima prioritat és desenvolupar a fons el segon i el tercer. La resta tenen una importància més baixa, però complementen els conceptes en què es basen els objectius més importants.

L'article està organitzat de la següent forma: en la *Secció 2* s'explica breument com ha sigut l'organització per tal de fer el projecte, i quin és l'estat de l'art actualment. En la *Secció 3* s'introdueixen els conceptes bàsics dels grafs dinàmics, i en la *Secció 4 i 5* es fa menció dels mètodes que s'han utilitzat. Per acabar, es tenen la *Secció 6 i 7*, que es mostren els resultats en els conjunts de dades i les conclusions dels objectius esmentats. A part, per facilitar la lectura de l'article, es té la *Taula 1*, que és un glossari de tots els termes bàsics que es mencionen freqüentment en els següents apartats.

2 METODOLOGIA I ESTAT DE L'ART

Per tots els objectius, s'ha enfocat la seva resolució de la mateixa forma. S'ha començat fent recerca de l'estat de l'art. Pels mètodes de privacitat, es tenen algorismes per tal d'obtenir privacitat utilitzant *k-anonymity* [1, 2], i altres que fan servir *Differential Privacy* [3]. Per les mètriques d'avaluació dels mètodes, s'ha explorat [5, 6, 7], on aquestes mesures i algorismes es poden utilitzar per comparar els grafs, tenint en compte el factor temporal.

Seguidament, s'ha realitzat diversos dissenys per a la implementació dels mètodes, utilitzant diagrames de classes. Finalment, s'ha implementat i executat el codi mitjançant el llenguatge de programació *Python*, utilitzant essencialment les llibreries següents:

- *NetworkX* (versió 3.2.1): Per a la creació dels grafs a partir de les dades dels *datasets*.
- *Pandas* (versió 2.1.4): Per fer conversió de dades en *DataFrames*. L'ús que es dona principalment és per fer agregacions.
- *Numpy* (versió 1.26.3): Es fa ús en el moment de fer operacions matricials de forma òptima.
- *Matplotlib* (versió 3.8.2): Utilitzat per a la visualització de gràfics. Eficient per observar mètriques i comparar resultats.

Per verificar que els mètodes implementats funcionen correctament, s'han dut a terme *Unittestings* durant tot el procediment. A més, per garantir una major diversitat en els resultats, s'han emprat conjunts de dades que tenen diferents característiques, que s'explica més detalladament en la *Secció 6*.

TAULA 1: GLOSSARI DE LA SIMBOLOGIA I TERMES GENÈRICS UTILITZATS EN L'ARTICLE.

Terme	Descripció
V	Conjunt de nodes d'un graf
E_t	Conjunt d'arestes d'un graf en l'instant t
$ V $	Nombre de nodes d'un <i>dataset</i>
$ E $	Nombre d'arestes d'un <i>dataset</i>
G_t	Graf que conté la dupla (V, E_t)
d_G	Densitat d'un graf G_t
\mathcal{G}	Seqüència de grafs temporals (G_1, G_2, \dots, G_n)
ε -ELDP	Algorisme <i>Edge Local Differential Privacy</i>
k -DA	Mètode <i>k-Degree Anonymity</i>

3 INTRODUCCIÓ ALS GRAFS DINÀMICS

Partim de la definició d'un graf, que és una estructura per modelar relacions entre objectes. Es representa com:

$$G = (V, E)$$

on V és el conjunt de vèrtex, i E les arestes que connecten nodes. Llavors, seguint la notació [1, 3] definim un graf temporal com una seqüència de grafs (*snapshots*) que mostra l'evolució temporal. Més formalment:

$$\mathcal{G} = \{(G_1, G_2, \dots, G_T) : G_t = (V, E_t) \text{ for } t = 0, \dots, T\}$$

on T és el nombre de *snapshots* que es tenen. Es considera que V és equivalent per tots els grafs, on només canvia E durant el temps. A l'agafar un conjunt fix de nodes, permet que les comparacions entre grafs temporals siguin coherents, i facilita la implementació dels models de privacitat.

Resumidament, els grafs dinàmics són un conjunt de grafs estàtics, el que comporta que sigui més complex interpretar-los. Les diferències principals que comporten els grafs temporals respecte als estàtics són:

1. Per tal d'assegurar que hi hagi privacitat, s'ha de protegir cada graf individualment, però també ha d'haver-hi privacitat de forma global [1]. És a dir, que a partir de tots els grafs que es tenen, no hi hagi manera de relacionar aquests que permeti reconstruir informació confidencial.
2. Generalment, el volum de dades és major i creix exageradament ràpid, pel fet de tenir múltiples grafs que canvien constantment. Això vol dir que la complexitat dels algorismes és important en cas de voler protegir dades repetidament.

3. Tenir un conjunt de grafs permet fer agregacions entre aquests. Una característica rellevant és fer agrupacions temporals (per dies, setmanes, mesos, etc.), ja que permet compactar les dades i simplificar l'anàlisi de l'evolució estructural del sistema al llarg del temps. Tècnicament, aquestes agregacions comporten a augmentar la densitat dels grafs i reduir la granularitat temporal, facilitant la detecció de patrons o tendències globals.

4 MÈTODES DE PRIVACITAT

En aquesta secció es presenten dos models de privacitat centrats en l'anonimització de les arestes dels grafs, és a dir, en l'ocultació de les relacions entre nodes dins les xarxes. Aquests mètodes es basen en els conceptes de *Differential Privacy* i *k-anonymity*, adaptats perquè siguin funcionals en grafs dinàmics.

4.1 ϵ -Edge Local Differential Privacy

Partim de la definició de *Local Differential Privacy* [4]: Un algorisme aleatori A satisfà ϵ -Local Differential Privacy si per totes les entrades possibles i, j , i tots els outputs $k \in \text{Range}(A)$:

$$\Pr[A(i) = k] \leq e^\epsilon \Pr[A(j) = k]$$

La intenció del *Local Differential Privacy* és afegir un soroll, de forma que al obtenir una sortida k , no es pugui saber amb certesa si s'esdevé de i o de j .

Tenint aquesta informació, es planteja l'algorisme ϵ -Edge Local Differential Privacy (ϵ -ELDP) [3] com es pot observar en la Figura 1. El primer pas és obtenir una matriu de probabilitats. El càlcul de les probabilitats de G_t s'obté de la següent fórmula:

$$(p_{01}, p_{10}) = \left(\frac{1}{e^\epsilon - 1 + \frac{1}{d_G}}, 1 - \frac{e^\epsilon}{e^\epsilon - 1 + \frac{1}{d_G}} \right)$$

on d_G és la densitat del graf d'entrada. Aquestes probabilitats s'utilitzen per generar soroll, una d'aquestes per afegir arestes (p_{01}), i una altra per treure'n (p_{10}) de forma proporcionada. Per tal de fer això, es generen dos *Gilbert Graphs* (G_b) a partir d'aquestes probabilitats. La forma en què es creen els grafs és a partir d'iterar cada parell de nodes que es tenen, i es genera una aresta amb la probabilitat p_{01} per un graf, i per l'altra amb la probabilitat p_{10} . Per últim, es combinen amb el graf original per obtenir el graf protegit. Els grafs es sumen observant la diferència simètrica de les seves arestes, és a dir, fent *XOR* entre conjunts d'arestes. Formalment, els càlculs per obtenir els grafs protegits són:

$$\begin{aligned} G' &= G \oplus G_0 \oplus G_1 \\ \text{on } G_0 &= Gb(n, p_{01}) \cap \bar{G} \\ \text{on } G_1 &= Gb(n, p_{10}) \cap G \end{aligned}$$

on n és el nombre de nodes que es tenen. Aleshores, es repeteix tot aquest procediment per tots els grafs que es tenen, tal com es veu en la Figura 2. Donades les demostracions en [3], aquest mètode només és vàlid si els grafs tenen la densitat $d_G \leq \frac{1}{2}$.

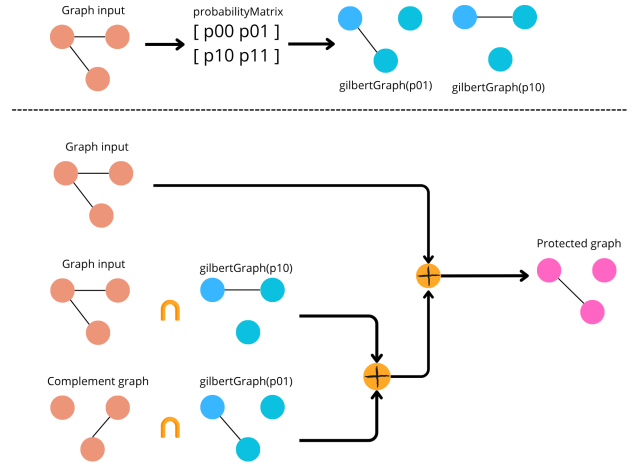


Fig. 1: Algorisme ϵ -ELDP, donat d'entrada ϵ i G_t .

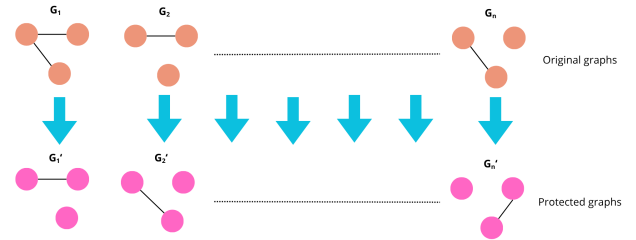


Fig. 2: Procediment de com s'aplica l'algorisme ϵ -ELDP per tots els grafs d'un dataset.

4.2 k -Degree Anonymity

Un conjunt de dades és k anònim si per cada registre es té almenys $k-1$ registres que tenen els mateixos valors. La motivació d'utilitzar k -anonymity és afegir incertesa per tal de dificultar que un atacant obtingui informació de la xarxa. Per definició, com més gran és la k , s'està afegint més incertesa. Dintre dels grafs, considerem això pels graus dels nodes, el que vol dir que per tots els graus que es tenen, almenys ha d'haver k nodes amb el mateix grau.

Ara bé, per tal d'assegurar que hi hagi k -Degree Anonymity (k -DA) en grafs dinàmics [1] s'ha de complir dues restriccions:

1. La seqüència de graus de cada graf ha de ser k anònim.
2. De forma global, ha de complir k anonimitat, de manera que per cada seqüència de graus que aparegui, s'ha de repetir almenys $k-1$ cops més durant tota l'evolució temporal. El motiu d'això és perquè si solament es té en compte les seqüències de graus per separat, es pot identificar de forma única els grafs que es tenen, el que pot facilitar obtenir informació.

LLavors, per implementar l'algorisme k -DA s'ha de seguir les fases que s'observen en la Figura 3, com es realitza en [1]. El procés que es segueix és:

1. Primer s'anonimitzen les seqüències de graus per tots els grafs. Específicament, s'ha decidit seguir el procediment de la Figura 4. Dels grafs que tenim, obtenim una matriu de totes les seqüències de graus i es calcula a partir d'això la matriu de medianes. S'obté a partir

de permutar cada seqüència, separar la llista en grups de mida $\frac{n}{k}$ (on n és el nombre de nodes), i calcular la mediana de cada grup. La seqüència anonimitzada de graus s'assignen a partir de les medianes obtingudes, de k en k .

2. El segon pas és comprovar que les seqüències siguin realitzables. És a dir, que sigui possible representar-ho en un graf simple. Per tal d'assolir això, es comprova el teorema de *Erdős Gallai* [10]. Llavors, en cas de no poder-se representar, o bé que la seqüència no aparegui k cops en tota la matriu, s'aproxima a la seqüència de graus més propera a aquesta, utilitzant la distància de *Manhattan* [11] com a mesura de proximitat.
3. Finalment, es grafiquen els grafs a partir de les seqüències de graus realitzables. La forma en què s'ha fet és amb *Havel-Hakimi* [12], que reconstrueix grafs a partir d'una llista de graus.

En el cas de ser un graf dirigit, el mètode s'ha de canviar lleugerament, al tenir graus d'entrada i de sortida. S'ha optat fer-ho de forma simple, on s'aplica l'algorisme per un dels dos graus, i seguidament s'obtenen els altres a partir d'una permutació dels resultats obtinguts. Això és perquè s'ha de tenir el mateix nombre de graus d'entrada i de sortida quan s'ha de reconstruir el graf.



Fig. 3: Passos per tal d'obtenir k -DA en un conjunt de dades.

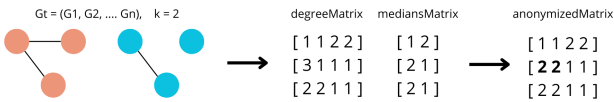


Fig. 4: Exemple del procediment per anonimitzar les seqüències de graus, donat com entrada k i G_t

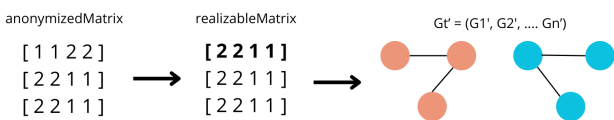


Fig. 5: Exemple del procediment per realitzar les seqüències de graus i graficar els grafs protegits.

5 MÈTRQUES PER AVALUAR ELS MÈTODES DE PRIVACITAT

Després d'implementar els mètodes de privacitat, s'han establert les mètriques per observar com canvien a nivell d'utilitat, a nivell estructural entre grafs, i quina és la pèrdua d'informació. Les mètriques que s'han escollit són les següents:

- **Índex de Jaccard:** Per una banda ens serveix per observar la similitud entre connexions dels grafs protegits

i originals. Per obtenir el percentatge, es fa la intersecció dividit entre la unió d'arestes entre el graf original $G_t = (V, E1)$ i el graf protegit $G'_t = (V, E2)$:

$$Jaccard(G_t, G'_t) = \frac{|E1 \cap E2|}{|E1 \cup E2|}$$

Per l'altra banda, s'ha utilitzat com a mesura de similitut entre els nodes més centrals dels algorismes de *Betweenness*, *Closeness* i *Degree Centrality*. En aquest cas, s'agafa un percentatge dels nodes més centrals de cada graf, i es computa l'índex de *Jaccard* amb aquests conjunts.

- **DeltaCon** [6]: Mesura de similitut estructural entre dos grafs, que es basa en l'afinitat entre nodes. A diferència de l'índex de *Jaccard*, aquest no només considera els enllaços directes, sinó també els camins indirectes. És a dir, té en consideració la influència dels nodes amb la resta de la xarxa.
- **Densitats dels grafs i graus dels nodes:** Informació d'utilitat per comparar si els grafs que es generen comparteixen el nombre de connexions respecte els grafs originals.

Seguidament, es va decidir prioritzar l'anàlisi de generació de comunitats i comparar-los de forma dinàmica. Per tal de dur a terme aquesta tasca, s'utilitza [7], que implementa diversos algorismes inspirats en el mètode *TSCAN*, els quals tenen com a objectiu identificar nodes que actuïn com a *StableCores* dins les xarxes temporals. Per formar *StableCores*, els nodes han de mantenir una similitut estructural (ϵ) amb un cert nombre de veïns (μ) en múltiples *snapshots* temporals (τ) consecutius. Es tenen els algorismes:

- **TSCAN-B:** Aquesta és la versió bàsica del algorisme, on no s'utilitza cap eina de *pruning* per fer *clustering* de *StableCores*.
- **TSCAN-A:** A diferència de *TSCAN-B*, s'utilitzen tècniques de *pruning* per descartar en primer lloc els nodes que no siguin candidats a *StableCores*.
- **TSCAN-S:** És una variant que utilitza directament *StrongCores*, una relaxació dels *StableCores*, com a nucli per a formar les comunitats, obtenint una major eficiència a canvi d'una lleugera pèrdua de precisió.

També es tenen mètriques per mesurar la qualitat de les comunitats generades, que són les següents:

- **Separabilitat (AS):** Mesura fins a quin punt una comunitat està ben separada de la resta de la xarxa. Es calcula com la proporció entre el nombre d'arestes temporals internes (entre nodes de la comunitat) i el nombre d'arestes temporals externes (entre nodes de dins i fora de la comunitat). Com més alt el valor, indica una comunitat més ben delimitada.
- **Densitat (AD):** Indica la connectivitat interna de les comunitats. Es defineix com el nombre mitjà d'arestes temporals per node dins de cada comunitat. Com més alt sigui aquest valor, més interaccions hi ha entre els membres de la comunitat.

- **Cohesió (AC):** Reflecteix la dificultat de dividir una comunitat en subcomunitats. Com més alt és el valor les comunitats són més difícil de separar.

6 RESULTATS

Després de decidir amb quins paràmetres comparar els grafs, s'han calculat i visualitzat les corresponents mètriques. Atès que ambdós mètodes de privacitat incorporen processos aleatoris, s'ha decidit executar-los cinc vegades per analitzar el seu comportament de manera més robusta.

Per tal que els mètodes de privacitat funcionin en tota mena de situacions, s'ha optat per recol·lectar diversos conjunts de dades de diferents característiques [8, 9]. En la *Taula 2* es mostra la llista de *datasets* que es contenen, i paràmetres que han estat importants en el moment d'escollir-los.

TAULA 2: DATASETS SELECCIONATS, AMB LES CARACTERÍSTIQUES PRINCIPALS DELS SEUS GRAFS TEMPORALS.

Nom	V	E	Direcció	#Snapshots
Aves-sparrow	52	516	No dirigit	2
Mammalia-voles	1480	4569	No dirigit	61
Insecta-ant	152	194K	No dirigit	41
Enron-employees	151	50.5K	Dirigit	16067
CollegeMsg	1899	59.8K	Dirigit	58911

TAULA 3: NOMBRE DE TIMESTAMPS PER A CADA AGRUPACIÓ TEMPORAL DELS DATASETS AMB UNIX TIMESTAMPS.

Nom	#DAYS	#WEEKS	#MONTHS
Enron-employees	867	161	38
CollegeMsg	193	29	7

A més, pels conjunts de dades que tenen un gran volum de *timesteps*, s'ha calculat les seves mètriques fent agrupacions temporals. Els *datasets* que s'han aplicat diferents agrupacions són els de la *Taula 3*, que aquests particularment tenen els temps representats en segons reals (*Unix Timestamps*). El motiu per ajuntar *snapshots* és per executar els algorismes més ràpidament, a part de proporcionar comparacions sobre com canvia el mateix *dataset* quan es compacten les dades.

En primer lloc, s'ha visualitzat les mitjanes de les mètriques de similaritat, corresponents a les *Figures 6, 7, 8 i 9*. Cal recalcar que pel càlcul de les mètriques de centralitat, s'ha escollit fer l'índex de *Jaccard* pel 5% de nodes més centrals. Es poden fer vàries observacions en aquests gràfics:

1. En l'algorisme ε -ELDP es pot controlar la pèrdua d'informació i d'utilitat segons el paràmetre ε . Quan es tria una major ε , s'està afegint menys soroll, el que implica que siguin més similars els grafs protegits en

comparació als originals. Es pot acabar de comprovar aquest anàlisi amb l'exemple de la *Figura 10*, on es veu detalladament com augmenten les similaritats quan s'escull una major ε per cada *timestamp*.

2. L'algorisme k -DA no implica el que passa en ε -ELDP, i es veu que els resultats són similars per totes les k provades. Sobretot, són sorprenents els valors obtinguts. El motiu de ser baixos per tots els conjunts de dades, es deu a la reconstrucció dels grafs amb *Havel-Hakimi*. Utilitzant *Havel-Hakimi* no t'assegura que els nodes dels grafs tinguin els veïns que es tenien en l'original, el que afecta directament a aquestes mesures.
3. De les mesures de similaritat, *DeltaCon* és més sensible al soroll (s'obtenen menors percentatges generalment) a comparació de l'índex de *Jaccard*. Això per definició és lògic, ja que addicionalment considera la connectivitat dels nodes amb la resta de la xarxa, cosa que no calcula l'índex de *Jaccard*.

També s'ha notat diferències entre els algorismes en les gràfiques de densitats i de graus, com es poden veure en les *Figures 11, 12 i 13*. En el cas de ε -ELDP, la densitat es preserva quan s'aplica la protecció als grafs en tots els casos. En k -DA, es preserva la densitat segons la k , i segons quina densitat tenen els grafs originals. Per exemple, si la densitat i els graus són molt baixos per defecte, com en *Mammalia-voles*, els grafs que generen són buits, pel simple fet que les medianes que es calculen en l'algorisme i s'assignen és igual a 0 en la majoria de seqüències.

El mateix succeeix en les *Figures 15 i 16*, que es compara les densitats del *dataset CollegeMsg* utilitzant k -DA en diferents agrupacions temporals. En aquests exemples es pot observar que les densitats van creixent i cada cop es tenen menys *timesteps*. Però encara això, menys per $k=2$, s'obtenen resultats similars en la resta de paràmetres.

Seguidament, s'han calculat les mètriques de detecció de comunitats amb *TSCAN*. Com a començament, s'ha establert un valor (ϵ - μ - τ) per defecte, i després s'han realitzat experiments variant els paràmetres. Els paràmetres utilitzats són els de la *Taula 4*.

TAULA 4: PARÀMETRES UTILITZATS PER A LA DETECCIÓ DE COMUNITATS EN TSCAN.

	ϵ	μ	τ
Valor per defecte	0.2	3.0	3.0
Valors utilitzats	[0.2, 0.5, 0.8]	[3.0, 5.0, 10.0]	[3.0, 5.0, 10.0]

Dels resultats més destacables que hem obtingut en les comparacions anteriors, hem vist com es generen les seves comunitats. En les *Figures 17 i 18* s'observen les mètriques de detecció de comunitats pel *dataset Mammalia-voles*. Es pot notar que els resultats obtinguts amb ε -ELDP són molt similars comparant amb els seus grafs originals. Com a observació, s'obtenen diferents resultats en algunes mètriques depenent del model *TSCAN* utilitzat. És a dir, encara que hi hagi models que permeten realitzar l'algorisme *TSCAN* més eficientment, no assegura tenir la mateixa qualitat de les comunitats generades. Per les gràfiques obtingudes, és possible que en alguns casos amb certs paràmetres no es detecten comunitats, i per tant no apareixen en les imatges.

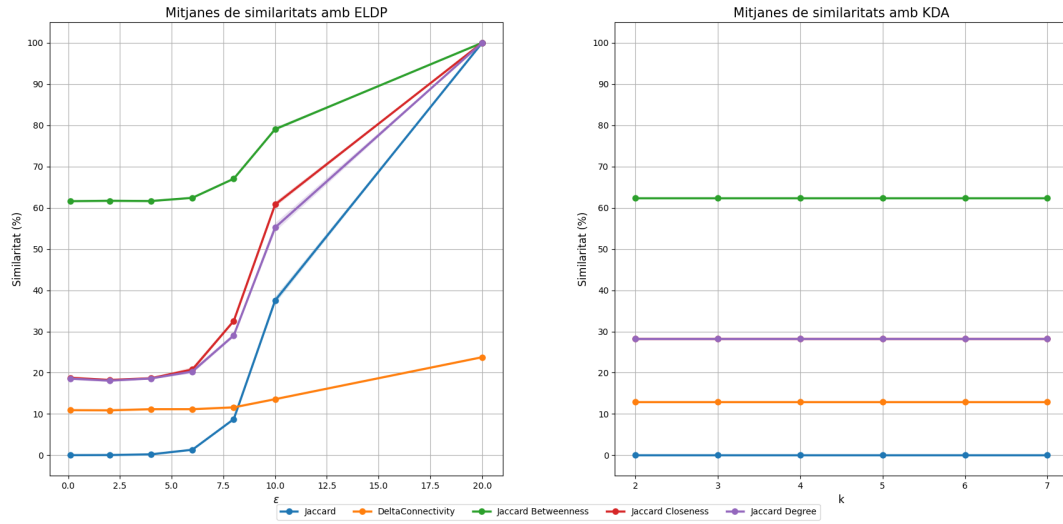


Fig. 6: Mitjana de mètriques de similaritat en *Mammalia-voles*. Els gràfics representen el valor (en percentatge) de les mètriques per cada paràmetre dels algorismes de privacitat utilitzats.

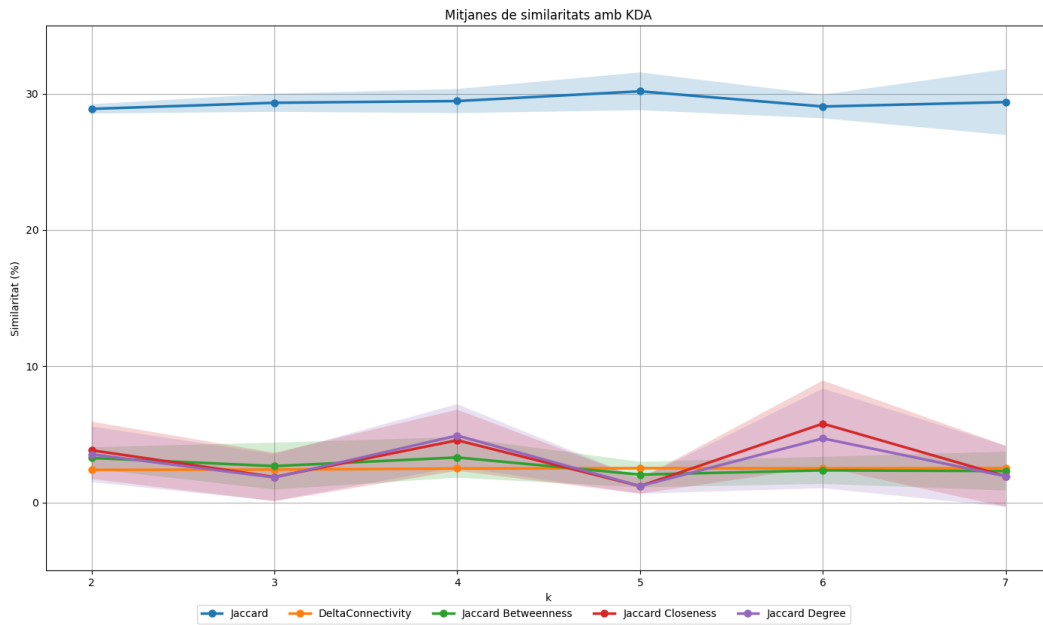


Fig. 7: Mitjana de mètriques de similaritat en *Insecta-ant*. En aquest cas només s'ha aplicat k -DA, perquè les densitats dels grafs són majors a $\frac{1}{2}$, el que vol dir que no es pot realitzar ϵ -ELDP per definició.

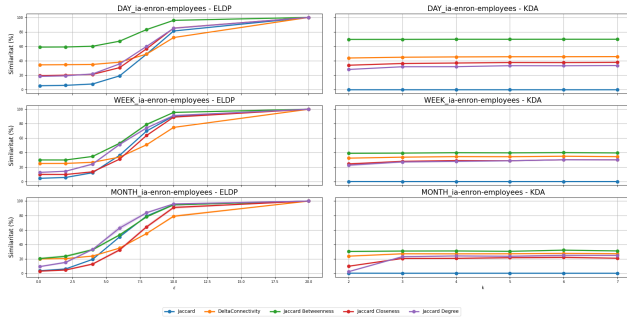


Fig. 8: Mitjana de mètriques de similaritat en *Enron-employees*. Els gràfics representen el valor de les mètriques per cada paràmetre dels algorismes de privacitat utilitzats per tots els agrupaments temporals.

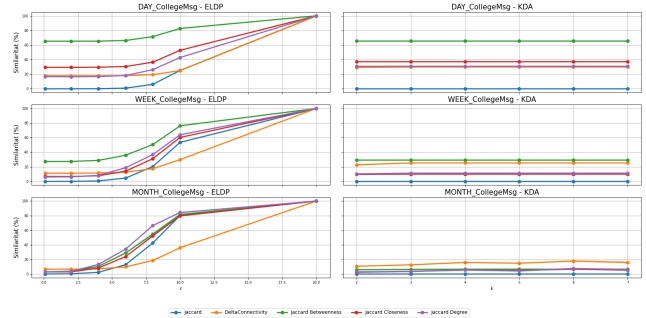


Fig. 9: Mitjana de mètriques de similaritat en *CollegeMsg*. Els gràfics representen el valor de les mètriques per cada paràmetre dels algorismes de privacitat utilitzats per tots els agrupaments temporals.



Fig. 10: Mapes de calor de les similituds en cada *timestamp* en *Mammalia-voles*. Es pot observar més detalladament com va evolucionant els valors de les similituds segons les ϵ triades en ϵ -ELDP, i com es menté segons la k escollida en k -DA.

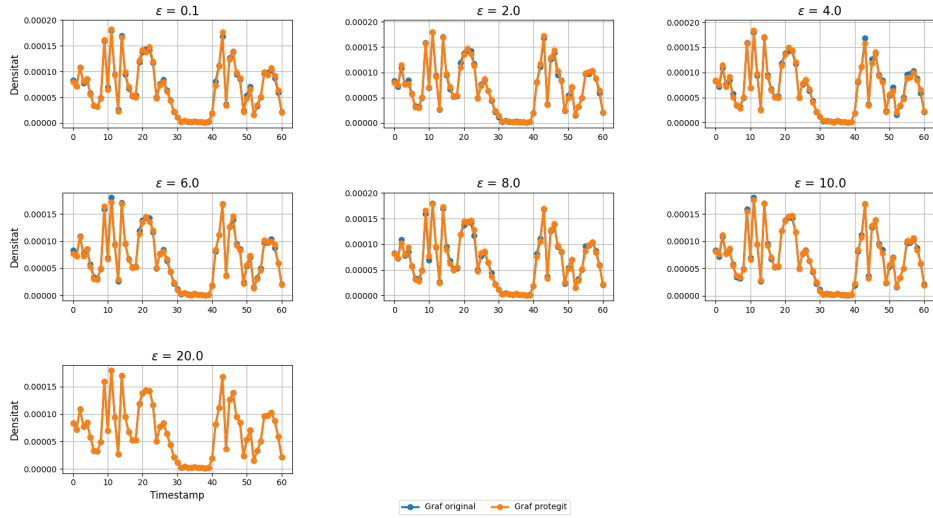


Fig. 11: Densitats pel mètode ϵ -ELDP del dataset *Mammalia-voles*. Es pot comprovar en aquesta imatge com l'algorisme preserva la densitat en els grafs protegits, independentment de la ϵ triada.

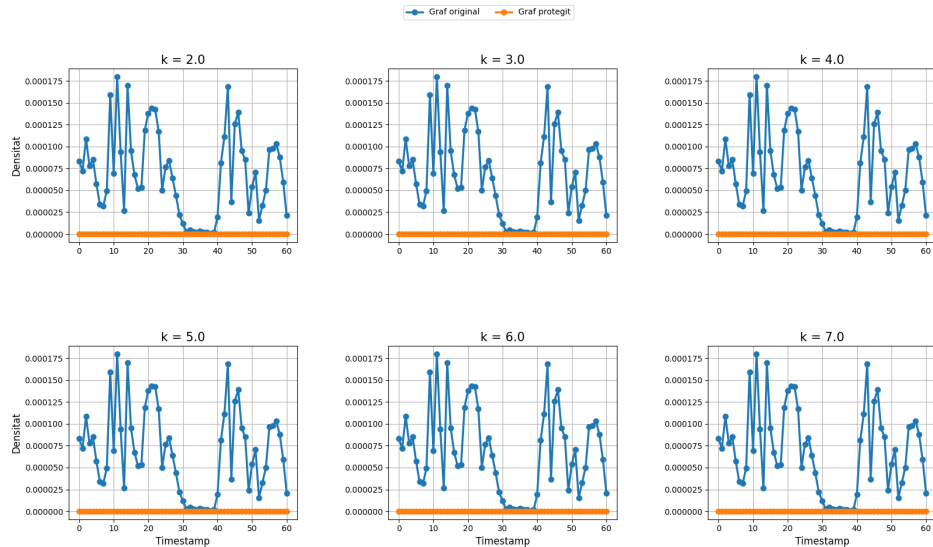


Fig. 12: Densitats pel mètode k -DA en *Mammalia-voles*. Aquí les densitats dels grafs protegits és 0, el que vol dir que generen grafs buits. Això succeeix perquè les densitats en els grafs originals són molt baixes.

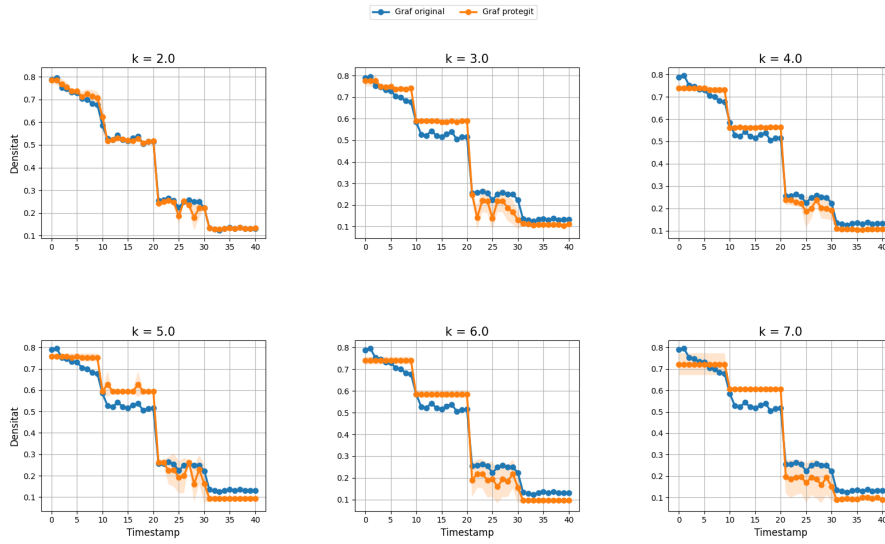


Fig. 13: Densitats pel mètode k -DA del *dataset Insecta-ant*. Es pot veure que més o menys en aquest cas es preserva la densitat. Quan menor és el paràmetre k , les densitats s'assemblen més als grafs originals. El mètode funciona perquè els grafs d'aquest *dataset* són densos.

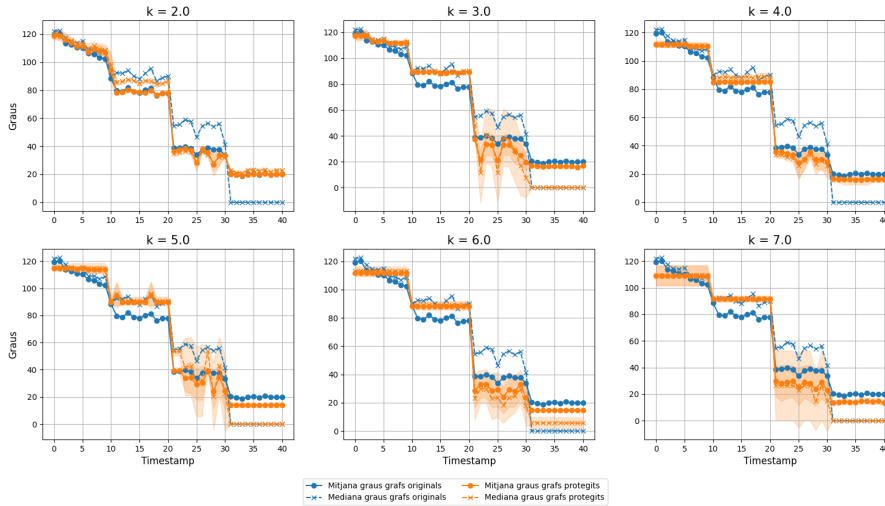


Fig. 14: Mitjana de graus i medianes pel *dataset Insecta-ant*. Es pot veure de forma més precisa com canvien les mitjanes i medianes de graus durant l'evolució temporal dels grafs.

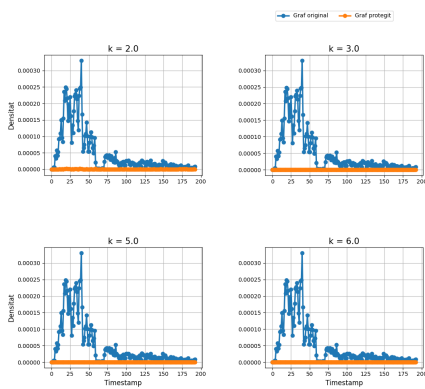


Fig. 15: Densitats pel mètode k -DA del *dataset CollegeMsg* agrupat per dies. En aquest cas es pot observar que cada gràfic de línies conté gairebé 200 *timestamps*, encara que les densitats són molt baixes en tots, i els grafs protegits que es generen no tenen connexions.

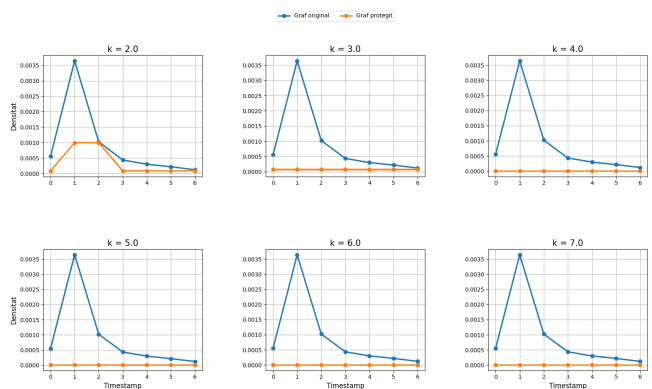


Fig. 16: Densitats pel mètode k -DA del *dataset CollegeMsg* agrupat per mesos. Comparat amb l'agrupació per dies, es tenen menys *snapshots*, però aquestes tenen una major densitat. Menys per $k=2$, els grafs que es generen en tots els paràmetres són buits.

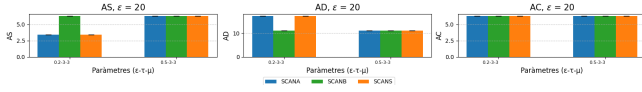


Fig. 17: Mètriques de detecció de comunitats en el mètode ε -ELDP en *Mammalia-voles*. S'obtenen resultats diferents segons el mètode TSCAN utilitzat.

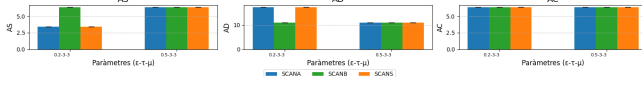


Fig. 18: Mètriques de detecció de comunitats pels grafs originals en *Mammalia-voles*. S'obtenen resultats molt similars comparant amb el mètode ε -ELDP.

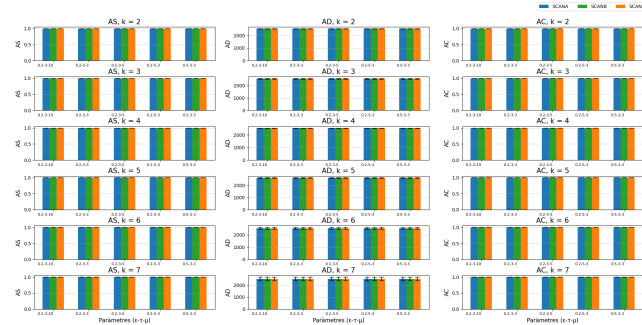


Fig. 19: Mètriques de detecció de comunitats en el dataset *Insecta-ant* utilitzant l'algorisme k -DA. En aquest cas es tenen la mateixa qualitat de comunitats per tots els mètodes i paràmetres.

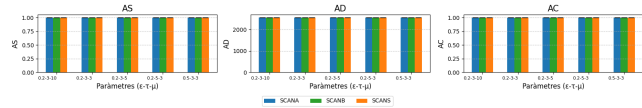


Fig. 20: Mètriques de detecció de comunitats en el dataset *Insecta-ant* en els seus grafs originals. Es tenen resultats molt similars comparat amb el mètode k -DA.

Això vol dir que no s'ha format cap *StableCore*, que és a causa d'agafar paràmetres (ε - μ - τ) que no compleixen els requisits. Com més grans són aquests paràmetres, més restrictiva és la forma de trobar comunitats.

Ara bé, en les Figures 19 i 20 es tenen els resultats pel dataset *Insecta-ant*. En aquest cas es té més variabilitat de *StableCores*, on els resultats pels tres mètodes TSCAN són molt similars tant utilitzant k -DA com en els grafs originals.

Per acabar, s'ha realitzat la detecció de comunitats amb els datasets amb agrupacions temporals, com es poden veure d'exemple les Figures 21, 22 i 23. Com a observació general, hi ha molta més variabilitat de resultats, segons els paràmetres que s'escullen per la detecció de *StableCores*, com també segons les variables dels mètodes de privacitat. Si es comparen amb els grafs originals, el dos datasets ha funcionat millor ε -ELDP, pel fet de ser poc densos els grafs.

Com a reflexió final de tots els resultats, el mètode ε -ELDP ha funcionat per tots els conjunts de dades que compleixen les condicions de l'algorisme. A part, aquest mètode permet controlar eficientment quanta informació es vol perdre a canvi d'obtenir privacitat. En canvi, si es té la intenció de només protegir els grafs, és útil k -DA per fer

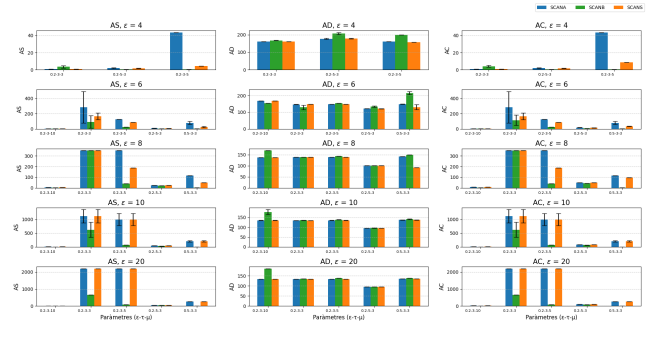


Fig. 21: Mètriques de detecció de comunitats en el dataset *Enron-employees*, agrupat per setmanes, aplicant el mètode ε -ELDP. Es pot observar que hi ha alta varietat de resultats.

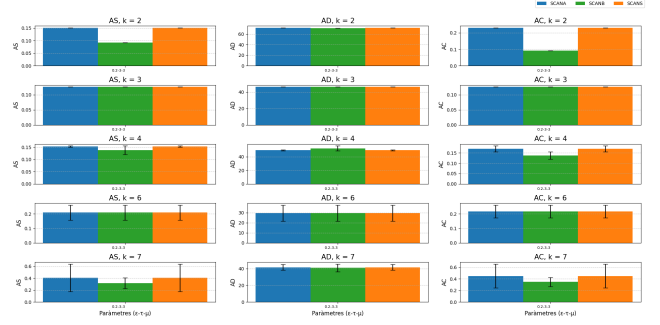


Fig. 22: Mètriques de detecció de comunitats en el dataset *Enron-employees*, agrupat per setmanes, aplicant el mètode k -DA. Aquí no hi ha variabilitat de resultats, però la qualitat de les comunitats és pitjor que amb ε -ELDP.

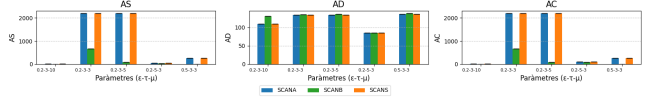


Fig. 23: Mètriques de detecció de comunitats en el dataset *Enron-employees*, agrupat per setmanes, pels grafs originals. El mètode que s'aproxima més a aquests resultats és ε -ELDP.

aquesta tasca de forma eficient, encara que pels resultats obtinguts és recomanable que els grafs siguin densos. No obstant, cal remarcar també que k -DA ha estat desenvolupat com a proposta pròpia amb l'objectiu de disposar d'un punt de comparació directe amb ε -ELDP. Per aquest motiu, seria interessant continuar desenvolupant-lo de manera més exhaustiva per explorar tot el seu potencial. Per a la detecció de comunitats, s'ha trobat exemples amb els dos mètodes on les mesures de qualitat són similars comparades amb els grafs originals, encara que les comunitats que es puguin generar siguin diferents.

7 CONCLUSIONS

En aquest projecte s'ha investigat, dissenyat, implementat i analitzat dues tècniques de privacitat per a grafs: ε -Edge-Local Differential Privacy i k -Degree Anonymity. Les metodologies proposades s'han aplicat amb èxit sobre diversos conjunts de dades, assolint els objectius més prioritaris i obtenint els resultats esperats per cada mètode.

Tanmateix, un dels objectius inicials (la integració dels

grafs anonimitzats en una xarxa neuronal *Graph Neural Network*) es va haver de descartar per limitacions de temps. Aquest aspecte representa una línia clara de continuació del treball, amb l'objectiu de validar l'impacte de les tècniques de privacitat en tasques d'aprenentatge automàtic.

A més, el projecte es podria ampliar explorant nous enfocaments de protecció de privacitat no tractats en aquest treball, com ara la incorporació de pesos a les arestes dels grafs, una característica especialment rellevant en escenaris com les xarxes de transaccions.

AGRAÏMENTS

Vull expressar el meu sincer agraïment al meu tutor, Guillermo Navarro Arribas, per la seva extraordinària tutorització al llarg d'aquest treball, així com per haver proposat la idea original del projecte. També vull apreciar l'ajuda de Julián Salas Piñón, que ha estat present durant tot el desenvolupament del projecte.

Finalment, voldria reconèixer la col·laboració amb el departament d'Enginyeria de la Informació i les Comunicacions (*dEIC*), que ha facilitat els recursos tècnics necessaris per a l'execució del codi.

REFERÈNCIES

- [1] L. Rossi, M. Musolesi i A. Torsello, "On the k-Anonymization of Time-Varying and Multi-Layer Social Graphs", *Proceedings of the International AAAI Conference on Web and Social Media*, 9(1), 377-386, 2021. Disponible en: <https://ojs.aaai.org/index.php/ICWSM/article/view/14605> [Darrer accés: 13-jun-2025].
- [2] J. Herrera, J. Casas i V. Torra, "An algorithm for k-degree anonymity on large networks", *Proceedings of the International Conference on Modeling Decisions for Artificial Intelligence (MDAI)*, 12-23, 2010. Disponible en: <https://www.researchgate.net/profile/Jordi-Herrera-Joancomarti/publication/262398542> [Darrer accés: 13-jun-2025].
- [3] S. Paul, J. Salas i V. Torra, "Edge Local Differential Privacy for Dynamic Graphs", In *International Symposium on Security and Privacy in Social Networks and Big Data* (pp. 224-238). Singapore: Springer Nature Singapore, 2023, Agost. Disponible en: https://link.springer.com/content/pdf/10.1007/978-981-99-5177-2_13.pdf [Darrer accés: 11-jun-2025].
- [4] B. Bebensee, "Local Differential Privacy: a tutorial", *arXiv preprint arXiv:1907.11908*, Juliol 2019. Disponible en: <https://arxiv.org/abs/1907.11908> [Darrer accés: 11-jun-2025].
- [5] E. Castrillo, E. León, i J. Gómez, "Dynamic Structural Similarity on Graphs", *arXiv preprint arXiv:1805.01419*, 2018. Disponible en: <https://arxiv.org/pdf/1805.01419> [Darrer accés: 02-abr-2025].
- [6] D. Koutra, J. T. Vogelstein i C. Faloutsos, "DeltaCon: A Principled Massive-Graph Similarity Function with Attribution", *SIAM International Conference on Data Mining (SDM)*, 2013. Disponible en: https://web.eecs.umich.edu/~dkoutra/papers/DeltaCon_KoutraVF_withAppendix.pdf [Darrer accés: 19-abril-2025].
- [7] H. Qin, R.-H. Li, G. Wang, X. Huang, Y. Yuan i J. X. Yu, "Mining Stable Communities in Temporal Networks by Density-Based Clustering", *IEEE Transactions on Big Data*, vol. 8, núm. 3, pp. 671-684, 1 juny 2022. Disponible en: <https://doi.org/10.1109/TBDATA.2020.2974849> [Darrer accés: 04-maig-2025].
- [8] J. Leskovec, Stanford Network Analysis Project (SNAP). Disponible en: <https://snap.stanford.edu/index.html> [Darrer accés: 10-mar-2025].
- [9] Ryan A. Rossi i Nesreen K. Ahmed, *The Network Data Repository with Interactive Graph Analytics and Visualization*, 2015. Disponible en: <https://networkrepository.com/dynamic.php> [Darrer accés: 10-mar-2025].
- [10] "Erdős-Gallai theorem", Wikipedia, l'enciclopèdia lliure. Disponible en: https://en.wikipedia.org/wiki/ErdosGallai_theorem [Darrer accés: 14-abr-2025].
- [11] "Manhattan distance", Wikipedia, l'enciclopèdia lliure. Disponible en: https://simple.wikipedia.org/wiki/Manhattan_distance [Darrer accés: 14-abr-2025].
- [12] "Havel-Hakimi algorithm", Wikipedia, l'enciclopèdia lliure. Disponible en: https://en.wikipedia.org/wiki/HavelHakimi_algorithm [Darrer accés: 13-abr-2025].