

Instrumente de clasificare pentru rețele multi-relaționale și serii de timp multidimensionale cu aplicații în biomedicină

- Rezumat -



UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică



Student Doctorand: Ramona-Carmen Stoica
Coordonator Științific: Bazil Pârv
Facultatea de Matematică și Informatică
Universitatea Babeș-Bolyai

A thesis submitted for the degree of

PhD

2017

1

Cuprinsul tezei de doctorat

Listă de figuri

Listă de tabele

0.1 Publicații / 1

1 Introducere /3

1.1 Motivația cercetării / 3

1.2 Obiectivele cercetării / 4

1.3 Contribuțiile tezei / 5

1.4 Structura tezei / 6

2 Rețele complexe / 9

2.1 Definițiile rețelelor / 9

2.2 Reprezentarea rețelelor / 11

2.2.1 Reprezentarea prin matrice / 12

2.2.2 Reprezentarea prin listă de adiacență / 12

2.2.3 Reprezentarea prin listă de incidență / 12

2.2.4 Reprezentarea prin matrice de incidență / 12

2.3 Măsurători utilizate în rețele / 12

2.3.1 Măsurători utilizate pentru noduri și muchii / 13

2.3.1.1 Măsurători pentru noduri: gradul de distribuție / 13

2.3.1.2 Măsurători pentru noduri: coeficientul de clustering/ 13

2.3.1.3 Măsurători pentru muchii și noduri: centralități / 13

2.4 Modele de rețele complexe / 16

3	Rețele biologice	/ 19
3.1	Rețele intracelulare	/ 19
3.2	Rețele metabolice	/ 20
3.3	Rețele de reglementare transcripțională	/ 24
3.4	Rețele biochimice	/ 28
3.5	Rețele de semnalizare celulară / Rețele de transducție a semnalelor și de reglare a genei	/ 31
3.6	Rețele de interacțiune genă-genă	/ 32
3.7	Rețele de interacțiune proteină-proteină (PPI)	/ 33
3.8	Rețele ale structurii proteinei	/ 35
3.9	Alte tipuri de rețele biologice	/ 36
3.9.1	Rețele de conexiuni sinaptice neuronale	/ 36
3.9.2	Rețele de conexiuni ale creierului	/ 36
3.9.3	Rețele ecologice alimentare, rețele filogenetice, rețele de corelare	/ 38
3.9.4	Rețele de asociere ale genelor	/ 39
3.9.5	Rețele de interacțiune între medicamente	/ 39
3.10	Rețele de interacțiune din lumea reală	/ 42
3.10.1	Rețele tehnologice	/ 42
3.10.2	Rețele de transport	/ 42
3.10.3	Rețele de socializare	/ 42
4	Clustering în rețele complexe	/ 45
4.1	Noțiuni de bază	/ 45
4.2	Măsuri de similaritate și disimilaritate	/ 46
4.2.1	Distanța euclidiană	/ 47
4.2.2	Distanța Manhattan	/ 47
4.2.3	Distanța maximă	/ 47
4.2.4	Distanța Minkowski	/ 48
4.2.5	Distanța medie	/ 48
4.2.6	Distanța pentru matching simplu	/ 48
4.3	Simililaritatea și disimilaritatea dintre clustere	/ 49
4.3.1	Distanța bazată pe medie	/ 49
4.3.2	Distanța celui mai apropiat vecin	/ 49
4.3.3	Distanța celui mai îndepărtat vecin	/ 49

4.4 Algoritmi pentru clustering /	51
4.4.1 Algoritmi pentru clustering ierarhic /	51
4.4.1.1 Clustering ierarhic aglomerat (clustering de jos în sus) /	51
4.4.1.2 Clustering ierarhic divizat (clustering de sus în jos) /	53
4.5 Algoritmul k-means /	54
4.6 Clustering pe serii de timp /	56
4.6.1 Clustering pe serii temporale multidimensionale /	57
4.6.2 O variantă de clustering al datelor din serii temporale multidimensionale /	60
4.6.2.1 Măsura de similaritate /	60
4.6.2.2 Setarea parametrilor: valoarea k /	61
4.6.2.3 Interfața /	61
4.6.3 Experimente numerice	61
4.6.3.1 Primul set de date /	61
4.6.3.2 Al doilea set de date /	63
4.6.3.3 Al treilea set de date /	64
4.6.4 Compararea cu metode tradiționale /	65
4.7 Rezumat /	66
5 Compararea rețelelor multi-relaționale /	67
5.1 Introducere /	67
5.2 Definiții /	67
5.3 Rețele multi-relaționale /	76
5.4 Aplicația pentru compararea rețelelor multi-relaționale /	78
5.5 Experimente și rezultate /	79
5.5.1 Seturi de date pentru compararea rețelelor multi-relaționale /	79
5.5.2 Rezultate /	79
5.6 Rezumat /	84
6 Motive în rețele multi-relaționale /	87
6.1 Introducere /	87
6.2 Defnirea motivelor /	88
6.2.1 Măsurarea semnificației motivului: Scorul Z /	92
6.2.2 Măsurarea semnificației motivului: valoarea P /	92
6.3 Motive în rețele multi-relaționale /	93
6.4 Identificarea motivelor în rețele multi-relaționale - aplicația MultiMot /	96

6.5	Experimente și rezultate /	100
6.5.1	Seturi de date pentru identificarea motivelor în rețele biologice multi-relaționale - aplicația MultiMot /	100
6.5.2	Rezultate /	100
6.6	Rezumat /	104
7	Comunități în rețele multi-relaționale /	107
7.1	Introducere /	107
7.2	Comunități în rețele multi-relaționale /	107
7.2.1	Identificarea comunităților în rețele biologice multi-relaționale - algoritmul Multi-NetCom /	110
7.2.2	Divizarea grafului în mai multe grafuri cu o singură relație /	111
7.2.3	Clustering pe grafurile cu o singură relație /	114
7.2.4	Colectarea rezultatelor individuale pentru grafurile cu o singură relație într-un singur graf /	114
7.2.5	Identificarea comunităților în graful transformat /	114
7.3	Experimente și rezultate /	121
7.3.1	Seturi de date pentru identificarea comunităților în rețele biologice multi-relaționale - aplicația MuliNetCom /	121
7.3.2	Rezultate /	122
7.4	Concluzii /	123
7.5	Rezumat /	123
8	Concluzii și direcții viitoare de cercetare /	127
8.1	Concluzii /	127
8.2	Direcții viitoare de cercetare /	129
8.3	Probleme deschise /	130

Contents

1	Cuprinsul tezei de doctorat	1
2	Introducere	7
2.1	Motivația cercetării	7
2.2	Obiectivele cercetării	8
2.3	Contribuțiile tezei	9
3	Rețele complexe	12
3.1	Definițiile rețelelor	12
3.2	Reprezentarea rețelelor	14
4	Rețele biologice	15
4.1	Rețele metabolice	15
4.2	Rețele de interacțiune genă-genă	20
5	Clustering în rețele complexe	21
5.1	Noțiuni de bază	21
5.2	Clustering pe serii de timp	22
5.3	O variantă a clustering-ului pe serii de timp multidimensionale	24
5.3.1	Măsuri de similitudine	24
5.3.2	Alegerea valorii pentru parametrul k	26
5.4	Compararea cu metodele tradiționale	27
5.5	Rezumat	27
6	Compararea rețelelor multi-relaționale	29
6.1	Introducere	29
6.2	Definiții	29

6.3	Rețele multi-relaționale	30
6.4	Experimente și rezultate	31
6.4.1	Seturi de date pentru compararea rețelelor multi-relaționale	31
6.4.2	Rezultate	31
6.5	Rezumat	31
7	Motive în rețele multi-relaționale	34
7.1	Introducere	34
7.2	Motive în rețele multi-relaționale	34
7.3	Experimente și rezultate	36
7.3.1	Seturi de date pentru detectarea motivelor în rețele multi-relaționale	36
7.3.2	Rezultate	36
7.4	Rezumat	40
8	Comunități în rețele multi-relaționale	41
8.1	Introducere	41
8.2	Comunități în rețele multi-relaționale	41
8.3	Experimente și rezultate	42
8.3.1	Seturi de date pentru detectarea comunităților în rețele biologice multi-relaționale - aplicația MultiNetCom	42
8.3.2	Rezultate	43
8.4	Rezumat	45
9	Concluzii	47
10	Cuvinte cheie	50
11	Referințele tezei	51
References		70

2

Introducere

2.1 Motivația cercetării

În zilele noastre, datele experimentale disponibile cercetătorilor și oamenilor de știință cresc exponențial. Reducerea costurilor pentru secvențializarea datelor din biologie și medicină oferă o oportunitate imensă pentru analiza, modelarea și simularea datelor. Bioinformatica, sistemele, biologia și ingineria bio-modelelor sunt considerate discipline promițătoare care vizează studiul și înțelegerea sistemelor biologice. Științele exacte, cum ar fi fizica sau chimia, precum și ingineria, utilizează modelarea matematică și computațională pentru a înțelege și analiza mai bine sistemele existente ale lumii reale, prin crearea de modele mai complexe și prin utilizarea simulărilor pentru a reproduce comportamentul lor dinamic.

Cu toate acestea, majoritatea sistemelor sunt atât de complexe încât toate ipotezele existente pot explica doar parțial comportamentul lor. Prin urmare, o întrebare firească este: cât de multe informații despre un sistem sunt necesare pentru a formula un model matematic care să îl explice pe deplin?

Motivația acestei teze provine dintr-o nevoie tot mai mare de a dezvolta modele mai bune pentru sisteme complexe, de mari dimensiuni, cu scopul de a obține noi perspective asupra comportamentului acestor sisteme. Dezvoltarea de noi metode pentru studierea interacțiunilor care apar în rețele complexe are un impact mare asupra multor discipline, de la științele sociale până la științele biologice. Cu toate acestea, această activitate se concentrează în principal pe interacțiunile bio-entitate în rețelele biomedicale. Provocările acestei teze includ dezvoltarea aplicațiilor cât mai generale astfel încât să

poată fi aplicate pe tipuri diferite de rețele biologice, precum și instrumentele practice care oferă perspective noi asupra comportamentului acestor rețele în diverse condiții (cum ar fi, de exemplu, mutațiile).

2.2 Obiectivele cercetării

Obiectivele principale ale acestei teze pot fi împărțite în două categorii de probleme din bioinformatică.

Prima categorie la care se referă această teză sunt rețelele multi-relaționale (majoritatea rețelelor biologice sunt rețele multi-relaționale). Al doilea aspect al tezei se referă la clustering-ul multidimensional al seriilor de timp (mai multe seturi de date biologice și experimente adună informații în timp).

Scopul acestei lucrări se concentrează în jurul a două tipuri de probleme și obiectivele propuse sunt dezvoltarea și analiza modelelor de rețele multi-relaționale, clustering-ul multidimensional al seriilor de timp și evidențierea unor direcții generale pentru extinderea modelelor dezvoltate în această teză, care pot evidenția o perspectivă suplimentară asupra problemelor abordate. Această teză prezintă două idei principale de abordare a problemelor biologice. Cu toate acestea, abordările noastre nu sunt limitate doar la domeniul biologic sau cel medical ci pot fi aplicate cu succes și în alte domenii.

Primul subiect se referă la rețele multi-relaționale (rețele cu mai multe legături de diferite tipuri între două noduri). Sunt abordate trei aspecte ale acestor rețele:

- Compararea rețelelor multi-relaționale,
- Identificarea de graphlets și motive în rețele multi-relaționale, precum și
- Identificarea comunităților în rețele multi-relaționale.

Provocările includ rezolvarea problemelor NP-complete (cum ar fi izomorfismul în rețea), identificarea algoritmilor eficienți pentru grafuri și motive precum și identificarea modalității corecte de reprezentare a rețelelor și a interacțiunilor dintre ele (avem mai multe rețele, unele din tipurile de muchii pot să lipsească într-una sau mai multe rețele, etc).

Al doilea obiectiv al acestei teze este abordarea datelor pentru seriile de timp multi-dimensionale precum și clustering-ul acestor entități. Cele mai multe date biologice sunt organizate ca serii de timp.

De multe ori, experimentele biologice sunt efectuate în condiții sau situații paralele, sau în medii diferite. Cu toate acestea, scopul final al experimentelor este acela de a identifica entitățile (gene, proteine, enzime etc.) care se comportă similar în același organism (sau organisme similare) în condiții diferite. Aceasta implică identificarea mai multor serii de timp care caracterizează o entitate și clustering-ul acestor entități în funcție de aceste serii de timp. Problema este interesantă și este utilă abordarea acesteia într-un mod diferit.

Un alt obiectiv al acestei teze a fost acela de a dezvolta aplicații practice care să contribuie la validarea soluțiilor propuse pentru problemele de cercetare abordate în această teză.

2.3 Contribuțiile tezei

Contribuțiile acestei teze pot fi grupate în două categorii: teoretice / conceptuale și aplicative / practice. În prima categorie menționăm algoritmul MultiNetCom (descriș în Capitolul 7), conceput pentru detectarea comunităților în rețele multi-relaționale.

Unele îmbunătățiri ale clustering-ului pentru seriile de timp multidimensionale sunt prezentate în Capitolul 4.

A doua categorie de contribuții conține patru aplicații software dezvoltate pentru analiza rețelelor multi-relaționale: (a) o aplicație care compară rețele multi-relaționale, NCTool; (b) o aplicație, MultiMot pentru identificarea motivelor; (c) o aplicație MultiNetCom pentru detectarea comunităților și (d) MDTSC, o aplicație pentru clustering-ul multidimensională a seriilor de timp.

Datele de intrare pentru aplicația NCTool constau în două sau mai multe rețele multi-relaționale. Pe aceste date de intrare se aplică câteva operații de bază, cum ar fi reuniunea dintre două rețele, intersecția dintre două rețele, subgrafurile, gradul de distribuție al nodurilor. Aceste aplicații sunt implementate pentru rețele de interacțiune genă-genă și rețele de interacțiune proteină-proteină, dar pot fi extinse și adaptate cu ușurință și altor tipuri de rețele. Utilizatorul are opțiunea de a selecta unul, două, mai multe sau toate tipurile de legături dintr-o rețea și de a rula operațiile pentru tipurile de legături selectate.

Aplicația MultiMot identifică motive de o anumită dimensiune (sunt implementate dimensiunile 3, 4 și 5 noduri), având posibilitatea de a detecta motive personalizabile.

Pot fi aplicate atât pentru grafurile directe cât și pentru cele indirecte.

Aplicația MultiNetCom implementează o extensie a algoritmului lui Fruchterman. Algoritmul propus de noi (conceput inițial pentru rețele cu o singură relație) este folosit pentru a rezolva problema detectării de comunități în rețele multi-relaționale.

Aplicația MDTSC a fost implementată pentru a fi folosită la clustering-ul entităților compuse din mai multe serii de timp. Aplicația are o interfață ușor de utilizat.

Ținând cont de domeniile de aplicabilitate, teza abordează aplicații din domeniul computațional care sunt utilizate în cea mai mare parte pentru a rezolva problemele din biologie și medicină, dar sunt destul de generale pentru a fi aplicate și în alte domenii similare.

Ținând cont de principalele domenii abordate, putem grupa contribuțiile în două categorii: *rețele* (sau *grafuri*) și *clustering*. În cadrul acestei teze ne-am focalizat pe aspectele specifice ale acestor topice și anume rețele multi-relaționale (care acceptă mai multe muchii de diferite tipuri între două noduri). De asemenea, am abordat conceptul de clustering al seriilor de timp multidimensionale (entități compuse din mai multe serii de timp).

Contribuțiile principale ale acestei teze sunt evidențiate mai jos:

- Definierea conceptelor necesare pentru înțelegerea rețelelor multi-relaționale și extinderea de la caracteristicile specifice rețelelor cu o singură relație la caracteristicile specifice rețelelor multi-relaționale.
- Proiectarea, implementarea și testarea algoritmilor pentru compararea rețelelor biologice multi-relaționale (47).
- Proiectarea, implementarea și testarea algoritmilor pentru detectarea motivelor și graphlets în rețele multi-relaționale. Această aplicație are două părți:
 - motivele de o anumită dimensiune (date de numărul de noduri) pot fi găsite folosind un algoritm de "sampling" (dimensiunile 3, 4 și 5 sunt implementate, dar motivele sunt ușor de adaptat la dimensiuni mai mari) și
 - motive personalizabile care pot fi detectate folosind un algoritm de backtracking. Pot fi căutate simultan mai multe motive.

- Proiectarea, implementarea și testarea unui nou algoritm pentru detectarea comunităților în rețele multi-relaționale (46).

3

Rețele complexe

Sistemele complexe sunt adesea analizate ca și rețele din cauza topologiei lor comune. Acest capitol definește conceptul de rețea și noțiunile specifice într-o abordare științifică, și descrie, de asemenea, clasificarea și reprezentarea acestora. Noțiunile de bază despre rețele includ: definițiile rețelelor, clasificări, structură și reprezentare (ca matrice, listă de adiacență și listă de incidență). Sunt descrise unele măsurători utilizate în rețele (gradul de centralitate, distribuția gradului, coeficientul de clustering, centralitatea). Modelele de rețele, cum ar fi modelul Erdos-Renyi (2), modelul Watts-Strogatz (34) și modelul Barabasi-Albert (2), necesare în capitolele următoare sunt de asemenea menționate și descrise pe scurt. Ultima secțiune descrie câteva proprietăți ale rețelelor: motive în rețea și grafuri.

3.1 Definițiile rețelelor

O mare varietate de sisteme naturale și sociale poate fi descrise de rețele cu o topologie complexă. De obicei, reprezentate ca și grafuri aleatoare, este unanim recunoscut faptul că evoluția și topologia lor în lumea reală sunt conduse de reguli și principii puternice (1). Această subsecțiune introduce conceptele de bază ale rețelelor complexe și oferă o imagine de ansamblu asupra principalelor modele de rețele întâlnite în lumea reală, cum ar fi rețelele de dimensiuni reduse precum și cele aleatorii (34).

Din cauza asemănării lor cu grafurile, rețelele sunt considerate ca fiind grafuri mult mai mari și sunt definite după aceleași structuri și relații:

- Fie $G = (V, E)$ reprezentarea unui graf, unde V reprezintă setul de noduri (sau vârfuri) și E reprezintă setul de muchii (sau conexiuni, legături între noduri) (8).
- Fie $|V|$ numărul nodurilor (cardinal de V) din graf. Fie $|E|$ numărul de muchii din graf.
- Fie $Deg(v), v \in V$, numărul de legături (muchii) conectate la v .
- Un drum între nodurile $s, t \in V$ este definit de o secvență alternantă de noduri și legături, începând cu s și terminând cu t . Fiecare legătură conectează nodurile adiacente (8). Lungimea drumului este determinată de suma greutateților aflate pe muchiile sale ținând cont de funcția de greutate (8).

Rețelele pot fi clasificate în funcție de cel puțin trei criterii: zona în care se aplică, structura și modelul acestora:

1. În funcție de aplicabilitatea lor avem următoarele tipuri de rețele:

- Rețele de calculatoare
- Rețele de transport
- Rețele de conexiuni sinaptice
- Rețele funcționale ale creierului
- Rețele de boli (gene bolnave)
- Rețele ecologice care sintetizează ecosistemele
- Rețele de telecomunicații
- Rețele World Wide Web
- Rețele de socializare
- Rețele biologice, etc.

2. După structura lor avem următoarele tipuri de rețele:

- Considerăm $\omega(e), e \in E$, funcția de greutate reprezentând legăturile sale.
- O rețea neponderată și simplă are $\omega(e) = 0$ pentru toate legăturile și doar o legătură $e \in E$ între două noduri $u, v \in E$ (8).

- O rețea ponderată are aceeași structură pentru legături ca și cea neponderată, dar cu $\omega(e) \in R$ (51).
- O rețea multiplă permite mai multe legături între aceeași pereche de vârfuri și poate fi clasificată și după greutatea muchiilor: aceasta poate fi ponderată sau neponderată (51).
- O rețea direcționată este descrisă de setul de muchii care au o direcție asociată cu acestea; poate fi și ponderată sau neponderată (51).

3. După nodurile și gradul de distribuție care reprezintă modelele de rețele:

- O rețea de dimensiuni mici are o distribuție specială a vârfurilor: distanța D între două noduri crește proporțional cu logaritmul numărului de noduri $|V|$ din rețea (51).
- O rețea scale-free (majoritatea rețelelor biologice sunt scale-free) respectă o distribuție asemănătoare muchiilor în care gradul fiecărui nod respectă formula pentru putere (51).
- O rețea aleatoare nu are nici o restricție de grad: nu se respectă legea gradelor nodurilor (51).

Cu siguranță există o clasificare mai detaliată a rețelelor, cu mai multe informații, însă aceste clasificări sunt generale și acoperă cele mai populare tipuri de rețele.

3.2 Reprezentarea rețelelor

Așa cum am definit mai devreme, o rețea este definită de $G(V, E)$ având setul de vârfuri V și setul de muchii E . Rețelele pot fi reprezentate în mai multe moduri, în funcție de conexiunea existentă între noduri:

- reprezentarea ca și matrice de adiacență,
- reprezentarea ca și listă de adiacență,
- reprezentarea ca și matrice de incidență,
- reprezentarea ca și listă de adiacență.

4

Rețele biologice

Acest capitol descrie diferite tipuri de rețele biologice și prezintă câteva aspecte comune între rețelele biologice și cele reale.

Există mai multe tipuri de rețele biologice însă nu toate fac obiectul de studiu al acestei lucrări. Acest capitol insistă asupra unor tipuri de rețele care vor fi utilizate în experimente sau rețele similare care pot fi abordate ca și extensii ulterioare ale aplicațiilor implementate în această teză. Unele dintre rețelele studiate includ interacțiunile dintre gene, interacțiunile dintre proteine, interacțiunile dintre elementele biochimice (reacțiile chimice) și rețelele metabolice.

4.1 Rețele metabolice

Activitatea unei celule este stabilită în rețele complexe de reacții chimice care produc diferite funcții celulare. Procesul care ne ajută să identificăm reacțiile dintr-o rețea se numește reconstrucția rețelei. Suma proceselor fizice și metabolice care determină caracteristicile celulelor, biochimice și fiziologice, formează o rețea metabolică. Prin însăși natura sa, aceste rețele compun reacțiile chimice ale metabolismului, căile metabolice și, în același timp, interacțiunile de reglementare care conduc aceste reacții (18, 22, 25, 39, 43).

Metabolismul intermediar poate fi considerat un instrument chimic care transformă materialele de bază în energie, precum și blocurile necesare pentru a produce structuri biologice, pentru a susține celulele și pentru a menține diferitele funcții ale celulelor. Acest instrument chimic se schimbă extrem de mult, acceptă legile chimiei și ale fizicii

și, din acest motiv, este limitat de diferite constrângeri fizico-chimice. În același timp, are o structură de reglementare complicată, care îi permite să răspundă la o combinație de perturbații externe. Dezechilibrul metabolic este rădăcina principalelor boli întâlnite la oameni, cum ar fi bolile de inimă, cancerul, diabetul și obezitatea. Metabolismul este compus din două tipuri diferite de transformări chimice, cum ar fi: căile catabolice care descompun substraturile în metaboliți simpli și căile anabolice care sintetizează aminoacizii, acizii nucleici, acizii grași și alte blocuri necesare.

Pe parcursul acestor procese, există un schimb complex între diferite grupuri chimice și potențialele de reducere ale oxidării (redox) care apar într-un set de molecule. Aceste molecule de transport și proprietățile pe care le transferă astfel leagă împreună rețeaua metabolică.

Ierarhia în rețele metabolice Rețelele metabolice sunt greu de înțeles pentru mintea umană deoarece reconstrucțiile pe scară genomică ale rețelelor metabolice sunt compuse din procese mari care constau în multitudinea de metaboliți și din când în când conțin peste o mie de reacții. Prin urmare, avem nevoie de modele matematice pentru a studia proprietățile lor și a le simula funcțiile. Cu toate acestea, putem vedea proprietățile unei rețele într-un mod ierarhic pentru a simplifica abordarea (teoria) funcțiilor rețelei.

Mai jos sunt prezentate cele patru niveluri de descompunere funcțională ale metabolismului:

- Nivelul 1: întreaga celulă
- Nivelul 2: sectoarele metabolice
- Nivelul 3: căile
- Nivelul 4: reacțiile individuale

Există două metode de reconstrucție prezentate mai jos:

- Reconstrucția metabolică a genomului
- Multiple rețele genomice

În rețelele metabolice avem metaboliți și căi metabolice. Metaboliții sunt particule microscopice cum ar fi glucoza și aminoacizii, sau macromolecule precum polizaharidele, și carbohidrații.

Căile metabolice sunt secvențe de reacții biochimice consecutive pentru o anumită funcție metabolică, de exemplu glicoliza sau sinteza penicilinei, care transformă un metabolit într-altul. Enzimele sunt proteine care catalizează (acelerează) reacțiile chimice. Astfel, într-o cale metabolică, metaboliții și enzimele sunt corespondenții nodurilor, iar reacțiile metabolice corespund muchiilor. Abordările mai simple indică faptul că nodurile reprezintă metaboliți și muchiile direcționate sunt reacțiile care transformă un metabolit într-altul. Exemple ale unei părți dintr-o cale de glicoliză, o reprezentare a metabolitului, reacții, i metaboliții sunt prezentate în figurile 4.1, 4.2, și 4.3:

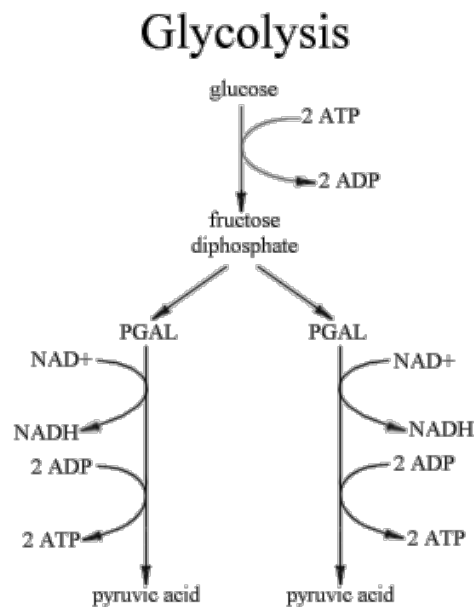


Figure 4.1: Parte a căii de glicoliză.

source: http://library.thinkquest.org/27819/ch4_4.shtml

Toate căile metabolice ale unei celule formează o rețea metabolică. O astfel de rețea este compusă din vizualizarea completă a metabolismului celular și a fluxului de material. Celulele se bazează pe această rețea pentru a digera substraturi din mediul înconjurător, pentru a genera energie și pentru a sintetiza componentele necesare din mediul înconjurător pentru creșterea și supraviețuirea lor. Aceste rețele sunt, de asemenea, utile pentru tratarea bolilor metabolice umane printr-o mai bună înțelegere a mecan-

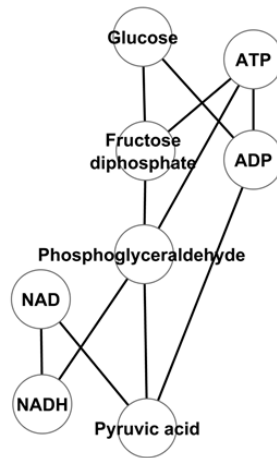


Figure 4.2: Reprezentarea unui metabolit

ismelor metabolice sau prin controlul infecțiilor patogene prin înțelegerea diferențelor metabolice dintre om și agenții patogeni (49).

Căile metabolice sunt construite parțial experimental și parțial din secvența genomului (omologie). Aceste rețele sunt generale și sunt folosite pentru multe organisme, de la bacterii până la oameni (26). Proiectul Enciclopedia de Gene și Genomul (KEGG) este o colecție vastă de baze de date online care au legătură cu genomul, căile enzimatice și substanțele chimice biologice.

KEGG este o bază de date creată cu scopul de a înțelege funcțiile și utilitățile sistemelor biologice. Ajută la învățarea sistemelor biologice cum ar fi organismul, celula și ecosistemele din punct de vedere molecular și genomic. Se compune din următoarele baze de date: informațiile chimice, informațiile rețelei și informațiile genomice (27).

Sistemul biologic este descris în figura 4.4. Reprezentarea include blocuri moleculare ale genelor și proteinelor (de exemplu, informații genomice) și substanțe chimice (informații chimice) care sunt combinate cu cunoștințele de diagrame de interacțiune, reacții și relații între rețele (informații ale sistemului) (26).

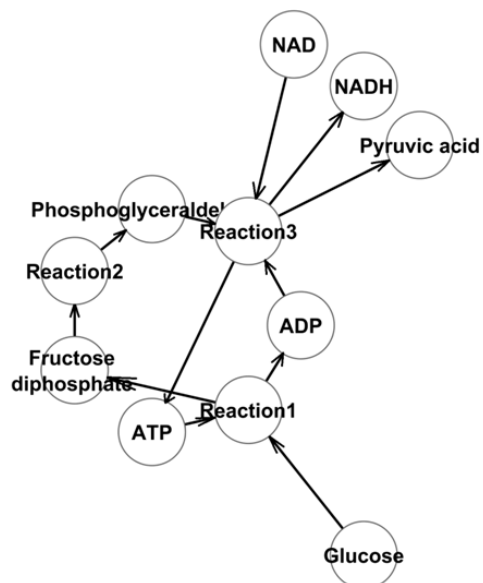


Figure 4.3: Reacții și metaboliți.

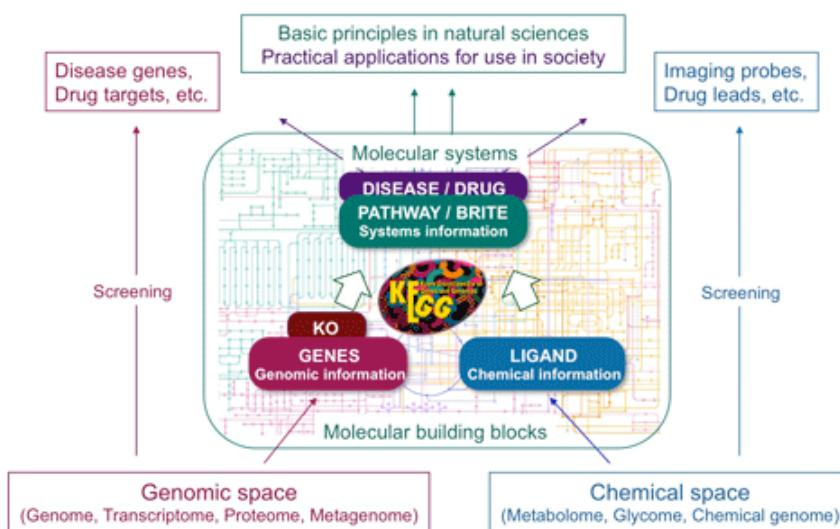


Figure 4.4: Reprezentarea digitală a sistemului biologic.
 source: <http://www.genome.jp/kegg/kegg1a.html> (26)

4.2 Rețele de interacțiune genă-genă

Rețelele de interacțiune dintre gene (GGI) sunt rețele în care nodurile reprezintă genele (de la un anumit organism) și legăturile sunt muchiile dintre ele. Aceste conexiuni sunt obținute din diferite studii, din experimente și din numeroasele baze de date existente. Există câteva instrumente GGI, dar una dintre cele mai comune și mai flexibile este Genemania (33, 37). Interacțiunile dintre genele incluse în Genemania sunt următoarele (50):

- **Co-expresie.** Două gene sunt de tipul co-expresie dacă nivelurile de expresie sunt similare în funcție de anumite condiții într-un studiu de exprimare al genei.
- **Interacțiunea fizică.** Două gene interacționează fizic dacă ele interacționează într-un studiu de interacțiune între proteine sau într-o bază de date.
- **Interacțiune genetică.** Două gene au o interacțiune genetică între ele dacă schimbările unei gene au un efect asupra modificărilor celeilalte gene.
- **Suprafețe comune de proteine.** Două gene (de exemplu proteine) au acest tip de interacțiune dacă au același domeniu proteic.
- **Co-localizare.** Două gene sunt co-localizate dacă ambele sunt exprimate în același țesut sau dacă produsele lor genetice sunt ambele identificate în aceeași locație celulară.
- **Pathway.** Două gene (de exemplu proteine) sunt conectate dacă sunt parte a aceleiași reacție într-o cale.
- **Predictibile.** Două gene au o relație funcțională predictibilă dacă există relații funcționale cunoscute de la un alt organism prin ortologie.
- **Alte interacțiuni.** Toate celelalte relații: date privind genomica chimică, corelații fenotipice de la Ensembl, informații despre boală de la OMIM.

5

Clustering în rețele complexe

Acest capitol prezintă noțiunile de bază de clustering, metricile de similitudine și câțiva algoritmi pentru clustering. De asemenea, capitolul descrie conceptul de clustering al seriilor de timp, extensia acestor concepte de clustering de la serii de timp unidimensionale la serii de timp multidimensionale și importanța practică a acestor algoritmi. Un algoritm nou de clustering al seriilor de timp multidimensionale este propus și implementat, iar implementarea este testată pe câteva seturi de date reale.

5.1 Noțiuni de bază

Clustering-ul este procesul prin care obiectele într-o mulțime sunt divizate în submulțimi distincte, astfel încât obiectele similare aparțin la aceeași submulțime. Clustering-ul face parte din clasa algoritmilor de învățare nesupervizată, deoarece nu există nici o proprietate care ar indica cum ar trebui grupate obiectele.

Un *cluster* conține obiecte care satisfac următoarele proprietăți:

- obiectele în fiecare cluster sunt foarte similare
- obiectele în cluster diferite sunt foarte diferite

Există multe definiții de similitudine în literatura științifică (13). Aceste definiții se bazează pe:

- tipul datelor luate în considerare
- tipul similitudinii dorite

De obicei, metricele de similitudine sunt bazate pe o funcție de distanță $d(x, y)$. Ideal, funcția aleasă ar trebui să satisfacă următoarele restricții (17, 54):

- $d(x, y) \geq 0$
- $d(x, y) = 0 \iff x = y$
- $d(x, y) = d(y, x)$
- $d(x, z) \leq d(x, y) + d(y, z)$

5.2 Clustering pe serii de timp

Seriile de timp sunt secvențe de numere reale care corespund valorilor unor parametri observate la intervale egale de timp.

Seriile de timp pot fi continue (variabila e definită pentru toate punctele în timp) sau discrete. Seriile de timp folosite în clustering sunt de obicei discrete și sunt o combinație de mai multe componente (6):

1. un trend
2. fluctuații de o anumită regularitate
3. un component sezonier
4. un efect aleator

Clustering-ul seriilor de timp are multe aplicații reale în diverse domenii științifice. Seriile de timp de obicei au câteva valori extreme, mai ales în mulțimile de date foarte mari. Elementele seriilor de timp au o ordine temporală și operațiile pe astfel de date sunt des întâlnite în domeniul extragerii datelor. Algoritmii de clustering ai seriilor de timp au fost studiați extensiv și eficiența acestor algoritmi a fost testată în numeroase aplicații reale (3, 5, 11, 19, 21, 30, 41, 48).

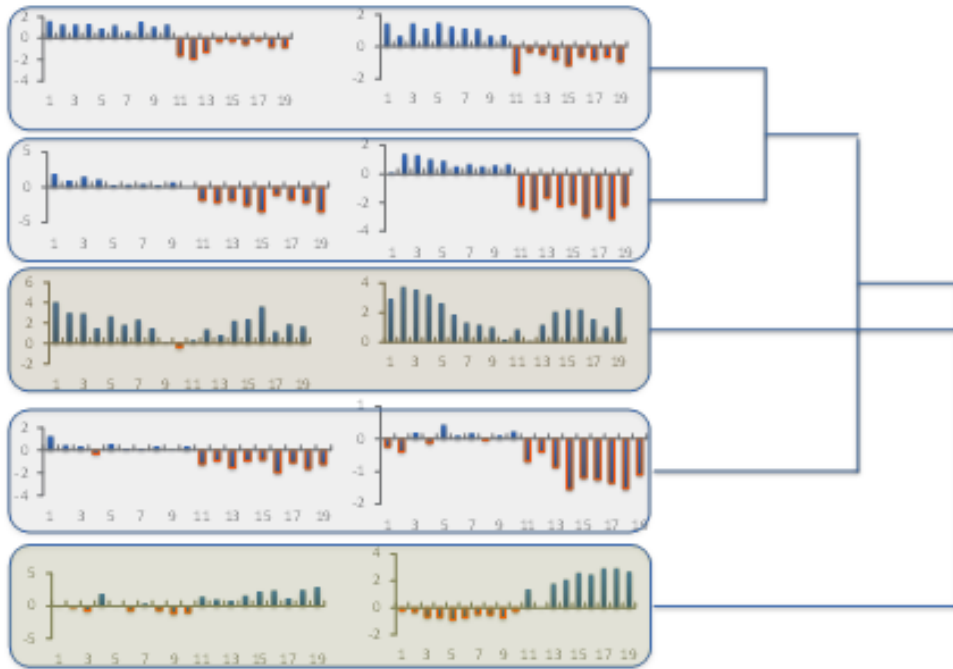


Figure 5.1: Serii de timp bi-dimensionale: exemplu de clustering ierarhică

Clustering pe serii de timp multidimensionale O serie de timp e definită printr-un șir $X = (x_1, x_2, \dots, x_n)$ de măsurări în timp a unui anumit parametru.

O *serie de timp multidimensională* e reprezentată prin:

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{pmatrix}$$

unde fiecare $X_i, 1 \leq i \leq N$, e o serie de timp separată. Seriile de timp individuale pot avea mărimi diferite.

Clustering-ul seriilor de timp multidimensionale presupune gruparea entităților de forma X . Figura 5.1 prezintă cinci entități pe care s-a făcut clustering în mod ierarhic ca serii de timp bi-dimensionale.

Clustering-ul seriilor de timp unidimensionale necesită gruparea mai multor puncte de date temporale: de exemplu, clustering-ul unor orașe în raport cu temperaturile raportate în fiecare oraș timp de doi ani. Această serie e formată din 730 puncte bi-dimensionale (2 ani câte 365 zile). În cazul seriilor de timp multidimensionale, fiecare

5.3 O variantă a clustering-ului pe serii de timp multidimensionale

punct de date e format din mai multe serii de timp. De exemplu, ar fi util de aplicat clustering-ul pe orașele pe baza vitezei vântului, presiunii, volumului precipitațiilor, etc. și nu doar pe baza temperaturilor. Figurile de mai jos prezintă exemple de entități cu două serii de timp. Unele entități sunt similare conform primei serii de timp, în timp ce altele sunt similare conform celei de-a doua serii de timp. În cazul seriilor de timp multidimensionale, clustering-ul trebuie să grupeze entitățile pe baza similitudinii ambelor serii de timp. Acest exemplu e ilustrat în Figura 5.2.

Abordarea comună în cazul multidimensional este de a concatena toate seriile de timp într-una singură și de a transforma problema într-o problemă cu o singură dimensiune a seriei de timp. Dar acest lucru poate duce la pierderea aspectelor generale ale problemei. Prin urmare, este avantajos să se trateze seriile temporale multidimensionale fără a le transforma: pe de o parte, ele oferă un punct de vedere global și prezintă unele patologii critice care decurg din discrepanțe evidente și, pe de altă parte, permit integrarea informațiilor conținute în fiecare serie de timp unidimensională a X și, prin urmare, este util atunci când fiecare matrice este redusă (15).

5.3 O variantă a clustering-ului pe serii de timp multidimensionale

O măsură de similitudine este de obicei utilizată pentru a calcula asemănarea dintre două serii de timp. În acest capitol, tratăm diferența dintre fiecare serie temporală a unei instanțe temporale multidimensionale ca o funcție obiectivă care trebuie minimizată.

5.3.1 Măsuri de similitudine

Pentru a compara similitudinea între două obiecte X și Y , unde:

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{pmatrix}, Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_N \end{pmatrix}$$

definim o funcție N -dimensională $F = (f_1, f_2, \dots, f_N)$:

5.3 O variantă a clustering-ului pe serii de timp multidimensionale

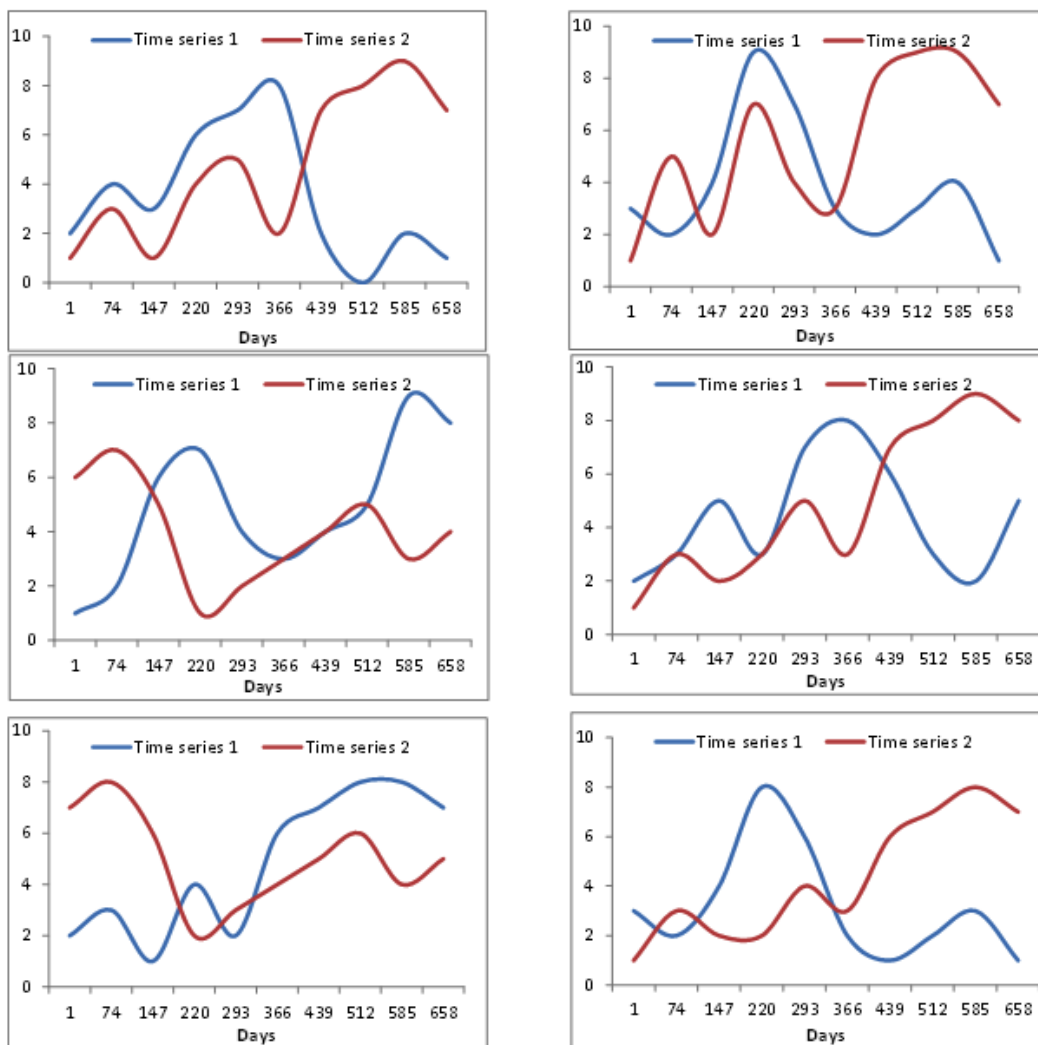


Figure 5.2: Seturile de măsurători pe parcursul celor doi ani

$$F = \begin{pmatrix} f_1 = d(X_1, Y_1) \\ f_2 = d(X_2, Y_2) \\ \vdots \\ f_N = d(X_N, Y_N) \end{pmatrix}$$

unde $d(\cdot)$ definește mărimea similitudinii.

Clasificăm seria temporală multidimensională folosind algoritmul de grupare k-means. Utilizăm funcția F pentru a determina grupul la care fiecare element trebuie atribuit. Calculăm valoarea de similitudine ca o combinație liniară a lui N și vom desemna rezultatul cu d_{sim} :

$$d_{sim} = \sum_{i=1}^N w_i f_i$$

unde w_i sunt greutatele care determină importanța fiecărei serii de timp în clustering. În experimentele noastre, toate seriile de timp au importanța egală ($w_i = 1, 1 \leq i \leq N$) și distanțele implementate sunt:

- Distanța euclidiană
- Distanța Manhattan
- Distanța maximă
- Distanța medie

5.3.2 Alegerea valorii pentru parametrul k

Utilizatorul poate selecta una din aceste patru măsuri din meniul principal, precum și procentul maxim de distanțe (setat în mod implicit la 0,6 în experimentele noastre). Valoarea inițială a k este setată la numărul total de elemente din setul de date. După fiecare repetare, semințele care nu au nici un punct de date alocat acestora sunt eliminate. Algoritmul trece prin numărul necesar de iterații pentru ca grupurile să se stabilizeze (distanțe mici între punctele din același grup și distanțe mari între punctele care aparțin la clustere diferite).

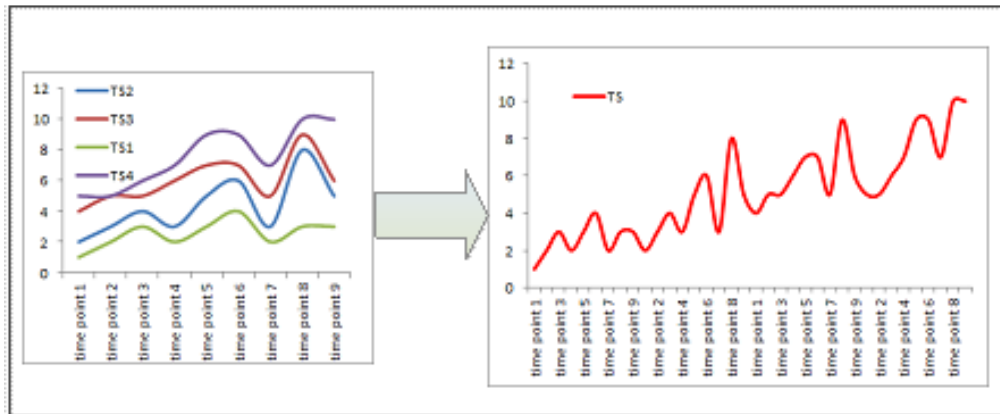


Figure 5.3: Comparația cu metodele tradiționale.

5.4 Compararea cu metodele tradiționale

Metodele clasice clasifică de obicei datele: deduc un singur parametru care caracterizează toate seriile de timp (16). Deci, în mod esențial, o serie temporală multidimensională este convertită într-o singură dimensiune, așa cum arată Figura 5.3. Am implementat și această tehnică, am comparat rezultatele cu cele obținute de abordarea noastră și am constatat că metoda tradițională tinde de obicei să împartă datele în mai multe clustere decât se aștepta. Acest lucru a fost deosebit de evident atunci când se utilizează valorile de similitudine medie și Manhattan pe primul set de date sau când se utilizează măsura medie la cel de-al doilea set de date.

5.5 Rezumat

Seriile temporale multidimensionale sunt o generalizare a seriei temporale unidimensionale. Ele au mai mulți parametri care sunt utilizați pentru a descrie fiecare instanță de date și, prin urmare, gruparea acestora este mai dificilă decât gruparea seriilor temporale unidimensionale. În acest capitol, am abordat problema de grupare a seriilor de timp multidimensionale ca o problemă multi-obiectiv și am implementat mai multe măsuri de distanță geometrică în algoritmul de grupare k-means pentru a testa rezultatul măsurii de similitudine împreună cu procesul de grupare. Am validat rezultatele pe trei seturi de

date și am ajuns la concluzia că distanța cea mai frecvent utilizată pentru gruparea seriei temporale unidimensionale (distanța euclidiană) ar putea să nu fie cea mai potrivită măsură pentru gruparea seriei temporale multidimensionale. Aceste rezultate au fost publicate în cadrul conferinței Intelligent System Design and Applications (ISDA 2013) (45) și în Studia Universitatis Babeș Bolyai în 2013 (44).

6

Compararea rețelelor multi-relaționale

6.1 Introducere

În acest capitol sunt introduse definițiile și conceptele legate de rețele multi-relaționale prin comparație cu cele uni-relaționale. Unele proprietăți utile ale acestor rețele sunt prezentate și este descrisă o aplicație implementată pentru compararea rețelelor, împreună cu principalele operații asupra rețelei. Ultima secțiune analizează experimentele referitoare la compararea rețelelor de interacțiune dintre gene pentru diferite tipuri de cancer. Aceste experimente sunt efectuate pentru a valida ideile și conceptele introduse în acest capitol.

6.2 Definiții

Termenul de rețea poate fi interpretat ca un grup de persoane, organizații, locuri etc. care sunt conectate sau lucrează împreună (7). În domeniul informaticii, o rețea este descrisă ca un graf mai mare care are un set de vârfuri conectate prin muchii. Prima diferență între grafuri simple și multigrafuri este dată de numărul de interacțiuni dintre noduri: grafurile simple au cel mult o interacțiune (muchie), în timp ce în multigrafuri

pot exista interacțiuni multiple între oricare două noduri. Un multigraf are un set de vârfuri și o mulțime de perechi neordonate de vârfuri care definesc muchiile. Un multigraf are toate muchiile de același tip. În cazul rețelilor (biologice), pot exista mai multe tipuri de interacțiuni între două noduri, astfel încât o rețea biologică este o generalizare a multigrafurilor.

Combinând aceste definiții în bioinformatică, putem spune că o rețea biologică este formată dintr-un set de vârfuri (gene, proteine, metaboliți etc.) și un set de muchii reprezentând interacțiunile (de diferite tipuri) între vârfuri.

6.3 Rețele multi-relaționale

Analiza și studierea rețelilor a devenit tot mai importantă într-o varietate de domenii, cum ar fi biologia, informatica și sociologia. În acest capitol ne concentrăm asupra rețelilor biologice. Prin progresele recente în biologie, o cantitate mare de date a devenit disponibilă. Cu toate acestea, datorită dimensiunii acestor date, provocarea constă în a studia aceste date. Trebuie găsită o modalitate de a transforma aceste date în modele mai ușor de înțeles. Un astfel de model care s-a dovedit a fi foarte eficient în descrierea relațiilor complexe și a interacțiunilor în sistemele biologice este rețeaua biologică.

O rețea este descrisă ca un set de noduri (vârfuri) și un set de conexiuni (muchii) care au diferite tipuri de relații între noduri. În biologie, rețelele sunt folosite pentru a descrie atât structura, cât și dinamica unui sistem biologic, prin urmare, identificarea rețelilor biologice este esențială în biologia sistemelor.

Rețelele biologice sunt compuse din entități biologice, iar interacțiunile dintre ele, pot fi folosite pentru a descrie multe tipuri diferite de relații. Exemple comune de rețele biologice sunt: rețelele de reglementare transcripțională (GRN), rețelele metabolice și biochimice și rețelele de interacțiune dintre proteine (PPI). Aceste rețele biologice sunt, de obicei, multi-relaționale: două noduri pot avea mai multe tipuri de interacțiuni între ele, iar aceste interacțiuni sunt reprezentate prin legăturile din rețea. De exemplu, căile metabolice pot fi descrise ca rețele multi-relaționale. Moleculele sunt nodurile, iar

activitățile enzimei, transducția semnalului sau reacțiile chimice sunt muchiile (35). Spre deosebire de rețelele uni-relaționale care permit un singur tip de muchii între vârfuri, rețelele multi-relaționale sunt mult mai potrivite pentru a reprezenta interacțiunile din lumea reală, dar sunt mai greu de analizat.

6.4 Experimente și rezultate

6.4.1 Seturi de date pentru compararea rețelelor multi-relaționale

Aplicația noastră a fost testată pe trei seturi de gene care sunt reprezentative pentru cancerul endometrial, ovarian și cel mamar. Am folosit aplicația Genemania pentru a crea rețelele și a le salva în fișiere pe care le-am transformat în formatul acceptat de aplicația noastră: fiecare muchie este reprezentată sub forma (*gena1, gena2, relație tip*); unde *gena1* și *gena2* reprezintă nodurile iar *tipul de relație* definește tipul de legătură care există între cele două noduri.

6.4.2 Rezultate

Am efectuat experimente pe trei tipuri de cancer: sân, ovarian și endometrial. Datele de intrare au fost preluate din baza de date pusă la dispoziție de Institutul Sanger al Institutului de tiin Weizmann (38) și de alte câteva studii. Interacțiunile pentru cele mai semnificative gene asociate fiecărui tip de cancer au fost extrase din Genemania (4, 35) și sunt prezentate în figura 6.1 (vizualizarea folosește un aspect de cerc sortat gradat și a fost desenată folosind aplicația Cytoscape (38, 42?)).

Detaliile pentru fiecare din cele trei rețele sunt prezentate în Table 6.1.

6.5 Rezumat

În acest capitol am prezentat toate definițiile și noțiunile de bază necesare pentru extinderea de la o rețea uni-relațională la o rețea multi-relațională. Rețelele multi-relaționale pot fi aplicate pe probleme din mai multe domenii diferite, iar rețelele bio-

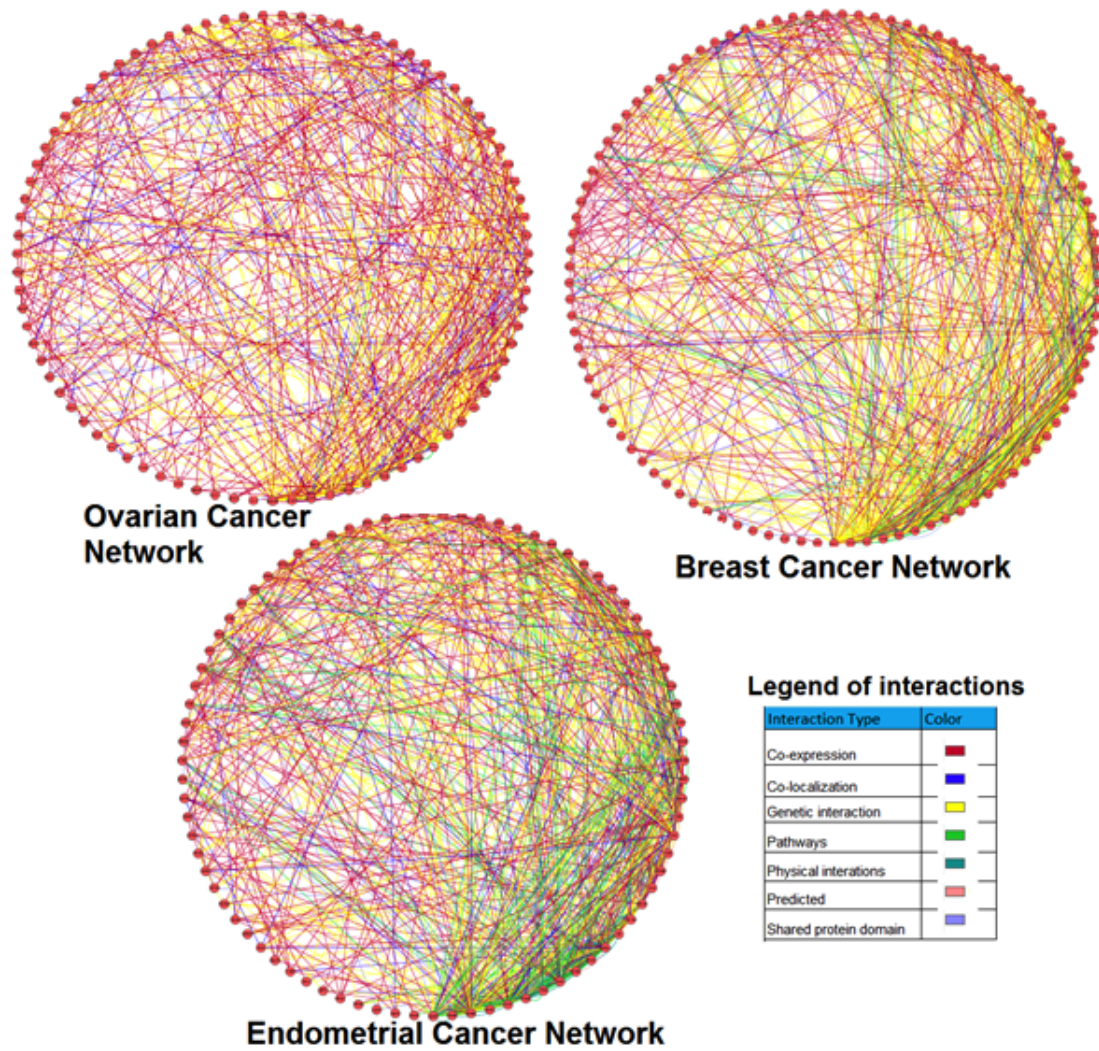


Figure 6.1: Reprezentarea rețelelor de interacțiune dintre gene pentru cele trei tipuri de cancer: ovarian, mamar și endometrial (aplicația Cytoscape (9, 38, 42) a fost folosită pentru a desena rețelele).

	Vancer ovarian	Cancer mamar	Cancer endometrial
Noduri	83	101	83
Muchii	590	817	801
Noduri cu mai multe perechi de muchii	47	74	108
Noduri izolate	0	0	0
Densitatea rețelei	0.157	0.14	0.19
Heterogenitatea rețelei	0.39	0.48	0.41

Table 6.1: Descrierea rețelelor de interacțiune dintre genele implicate în cancerul ovarian, mamar și endometrial.

logice sunt una dintre cele mai comune aplicații ale acestora. Capitolul a prezentat, de asemenea, o aplicație care este concepută pentru a compara rețele multi-relaționale. Au fost efectuate și descrise experimente privind datele reale legate de diferite tipuri de cancer feminin. Experimentele au fost efectuate pentru trei tipuri de rețele de gene implicate în trei tipuri de cancer similare (ovarian, mamar și endometrial) și rezultatele au fost prezentate la Conferința Internațională de Design și Aplicații a Sistemelor Inteligente (ISDA) (47).

7

Motive în rețele multi-relaționale

7.1 Introducere

Acest capitol prezintă motive în rețea și grafuri și este structurat în trei secțiuni. În prima secțiune sunt prezentate noțiunile de bază despre motive și grafuri. De asemenea, descriem semnificația motivelor și a grafturilor și importanța acestora în biologie. O extindere a motivelor și al modului de detectare al acestora, în cazul rețelelor multi-relaționale, este prezentată împreună cu operațiile de bază asupra motivelor din rețele. Este descrisă o aplicație care detectează motive în rețele multi-relaționale, care este o extensie a unei aplicații celebre dezvoltate în laboratorul lui Uri Alon la Institutul Weizmann, descrise în a doua parte a capitolului. Extensia nu se referă doar la rețele multi-relaționale (extinse din rețelele uni-relaționale), ci și la o aplicație de căutare a motivelor de către utilizator. Utilizatorul are posibilitatea de a-și defini propriul motiv și de a-l căuta în rețea. În ultima parte a acestui capitol sunt efectuate experimente, inclusiv date biologice reale, pentru a valida eficiența acestei aplicații.

7.2 Motive în rețele multi-relaționale

Cantitatea de informație care a devenit disponibilă în ultimii ani a permis studiul unor rețele foarte mari: World Wide Web, rețele sociale și, în special, rețele biologice,

printre care se numără rețelele de interacțiune dintre proteine (PPI) (20, 23, 24), rețelele metabolice (25) și rețelele de interacțiune dintre gene (GGI) (31, 32). Aceste rețele aduce o multitudine de date, dar extragerea de informații semnificative din aceste date are provocările sale. Trebuie dezvoltati algoritmi care sunt eficienți și în același timp robusți din punct de vedere al erorilor întâlnite în datele analizate.

Din nefericire, datorită complexității lor, unele rețele nu pot fi analizate complet. Genomul uman este un exemplu clasic al unei astfel de rețele. Aceste rețele sunt, de obicei, extrem de importante. Prin urmare, obținerea oricărei noi înțelegeri în comportamentul lor este foarte important. O modalitate de a realiza acest lucru este prin abordarea paralelă cu rețele mai simple, care sunt mai ușor de analizat. Odată ce un subsistem dintr-un organism mai simplu este bine înțeles, această experiență poate fi aplicată organismului mai complex. Această abordare permite, în esență, izolarea și analizarea în mod independent a subsistemelor mici în organismele mai complexe.

O clasă de astfel de algoritmi este *detectarea motivelor* (32). Acești algoritmi urmăresc să detecteze modelele de conectivitate care apar mai des într-o rețea și să ofere, de asemenea, o imagine asupra modularității și structurii unei rețele (20, 31, 36, 53). S-a demonstrat că aceleași motive se găsesc în multe organisme diferite, care le transformă într-un instrument esențial pentru transferul de cunoștințe de la organisme simple la subsisteme mai complexe.

Pe lângă aplicațiile lor în rețelele biologice, motivele au fost aplicate cu succes și în alte zone.

Rețelele biologice prezintă un volum mare de date, modelate ca rețele de interacțiune, căi metabolice și de semnalizare, rețele de reglementare etc. Grafurile multi-relaționale sunt cele mai utilizate pentru a arăta complexitatea și amploarea lor. Caracteristicile acestui tip de rețele sunt importante în analiza interacțiunilor dintre entitățile biologice. Ele sunt numite motive în grafuri.

În general, rețelele biologice (și motivele dintr-o rețea) sunt multi-relaționale. Acest lucru complică procesul de identificare a motivelor.

Există mai multe aplicații (29, 52) care, în grafuri simple unidirecționale și multidirecționale,

identifică motive de dimensiuni diferite dintr-o rețea. Aplicația propusă de noi detectează motivele din rețelele multi-relaționale.

Oferim două moduri de identificare a motivelor. În primul rând, *template matching*, utilizatorul poate specifica unul sau mai multe template-uri de grafuri pe care le caută.

În al doilea rând, abordarea de eșantionare a subgrafelor este o variație a algoritmului pentru grafurile multi-relaționale (28). Estimează frecvența subgrafurilor de dimensiune N folosind subgrafurile k obținute în mod aleator din rețea.

7.3 Experimente și rezultate

7.3.1 Seturi de date pentru detectarea motivelor în rețele multi-relaționale

Considerăm rețeaua definită astfel:

$node_a node_b edge_1$

$node_a node_d edge_3$

$node_e node_a edge_1$

$node_b node_e edge_2$

$node_d node_a edge_1$

$node_e node_f edge_1$

unde $node_x$ reprezintă un nod (x este eticheta sa) și $edge_i$ reprezintă o relație între două noduri (i este tipul relației).

7.3.2 Rezultate

Pentru a ilustra modul în care funcționează aplicația MultiMot, am selectat un exemplu de rețea multi-relațională de interacțiune dintre gene. Am considerat primele 50 de gene neregulate ale unui țesut cardiac de la două specii de șoareci. Datele sunt preluate din baza de date Gene Expression Omnibus (12). Rețelele de interacțiune ale genelor (descrise în figura 7.1) au fost generate utilizând Genemania.

Rețelele de interacțiune dintre gene ale țesutului cardiac exprimate diferit sunt reprezentate în figura 7.1. Tipurile de legături reprezentate cu culori diferite sunt prezen-

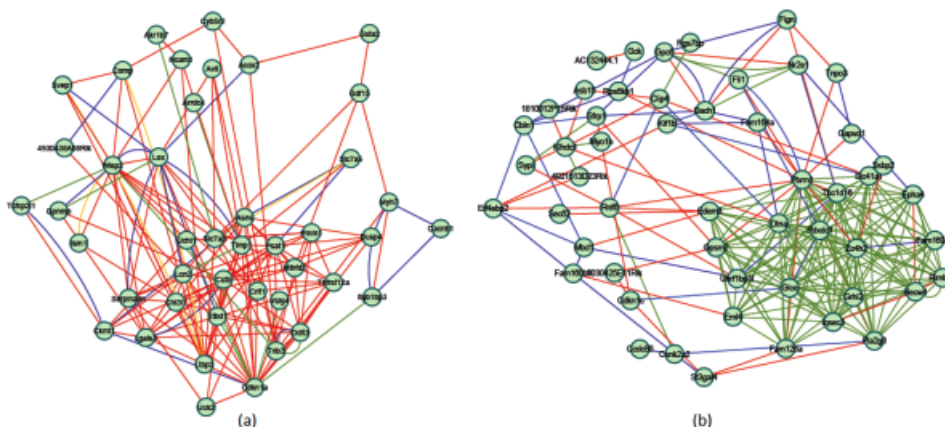


Figure 7.1: Reprezentarea rețelelor de interacțiune dintre gene ale a unui țesut cardiac exprimat diferențial.



Figure 7.2: Reprezentarea legendei rețelelor de interacțiune dintre gene ale unui țesut cardiac exprimată diferențial.

tate în figura 7.2.

Dacă metoda trebuie să găsească toate motivele de dimensiune 3 în fiecare dintre rețele, atunci rezultatul este cel din Figura 7.3 respectiv Figura 7.4.

Legenda tipurilor de legături sunt reprezentate în Figura 7.2.

Motivele de dimensiune 3 găsite de MultiMot în rețele, Figura 7.2 (a) și Figura 7.2 (b) sunt reprezentate în Figura 7.3 și Figura 7.4.


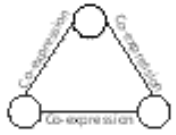

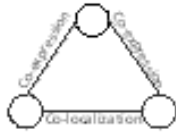



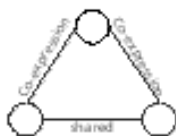


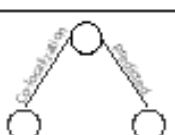





Motif	Network (a)	Network (b)	Motif	Network (a)	Network (b)
	223	22		63	5
	48	69		6	1
	60	19		2	4
	17			3	
	1	7		1	1
	4	19		4	4
	1			1	
	1			1	1

Figure 7.3: Frecvența motivelor de dimensiunea 3 găsite de MultiMot în rețelele din Figura 7.1 (a) și Figura 7.1 (b).


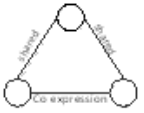
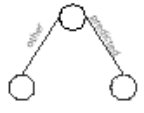
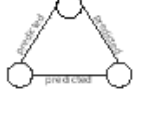
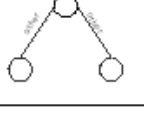
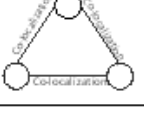
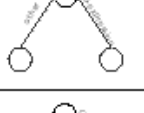
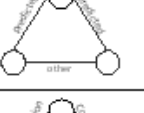
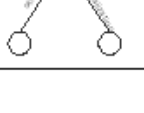
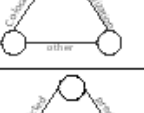
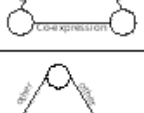
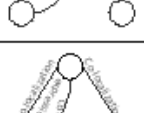
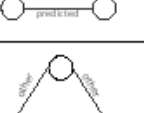

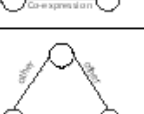
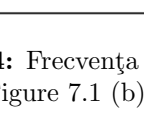
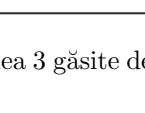
	2	219		1	
	31			0	20
	4			0	6
	9			0	1
	6			0	8
				0	10
	0	1		0	1
	3	0		0	1
	0	1		0	1

Figure 7.4: Frecvența motivelor de dimensiunea 3 găsite de MultiMot în rețelele din Figura 7.1 (a) și Figure 7.1 (b) (continuare).

7.4 Rezumat

Aceast capitol prezintă un algoritm nou de detectare al motivelor în rețele multi-relaționale.

Identificarea unor modele similare și frecvent întâlnite (motive în rețea) în rețele oferă informații importante pentru o mai bună înțelegere a bolilor și are un imens impact în studiile biologice. Această lucrare prezintă un algoritm nou de identificare a motivelor, care este o extindere a algoritmului folosit pentru detectarea motivelor în rețele uni-relaționale și se ocupă de grafuri eterogene multi-relaionale. Nu numai că găsește motive standard de o anumită dimensiune, ci caută și motive customizate (personalizate) (care pot avea un grad mai mare de complexitate). Acesta consideră că atât în grafurile direcționate, cât și în cele nedirecționate, pot fi detectate motive pe un subset de muchii și pot fi folosite metacaractere pentru etichetarea muchiilor, dacă este necesar, ceea ce la creșterea generalității.

8

Comunități în rețele multi-relaționale

8.1 Introducere

Acest capitol oferă o prezentare generală a comunităților din rețele multi-relaționale și rețele uni-relaționale și mai mulți algoritmi pentru identificarea comunităților în rețele uni-relaționale, respectiv rețele multi-relaționale. De asemenea, propune un algoritm nou pentru detectarea comunităților în rețele multi-relaționale. În continuare, în acest capitol prezentăm o aplicație implementată de noi care are la bază algoritmul pe care l-am propus, precum și rezultatele obținute prin rularea algoritmului pe mai multe seturi de date.

8.2 Comunități în rețele multi-relaționale

Rețelele complexe sunt modele matematice care sunt folosite pentru a reprezenta într-o formă citibilă de mașini multe fenomene de interacțiune care au loc în lumea reală. Următoarele rețele sunt exemple de rețele mari:

- World Wide Web: nodurile reprezintă paginile web, iar legăturile reprezintă hyperlink-uri.

- Rețele metabolice: nodurile reprezintă enzime și legăturile reprezintă reacțiile între ele.
- Facebook: nodurile reprezintă utilizatorii și legăturile reprezintă relațiile de prietenie.

Comunitățile pot fi definite ca grupuri discontinue de entități într-un graf, astfel încât fiecare entitate să fie "mai aproape" de toate celelalte entități din același grup decât de entitățile din afara acesteia. Detectarea comunităților într-o rețea reprezintă un pas important în identificarea modelelor din acea rețea (10).

Detectarea comunităților este foarte importantă în domeniul informaticii și al biologiei, unde informațiile sunt reprezentate ca rețele. Cele mai multe rețele în natură sunt multi-relaționale. Cu toate acestea, detectarea comunităților în rețele multi-relaționale este dificilă, deoarece complexitatea crește odată cu dimensiunea rețelei și cu numărul de tipuri de conexiuni. Prin urmare, atunci când căutăm comunități, este adesea mai ușor să reducem rețeaua multi-relațională la o rețea uni-relațională. Acest lucru duce în mod inevitabil la pierderea informațiilor, iar rezultatele nu sunt întotdeauna exacte.

8.3 Experimente și rezultate

8.3.1 Seturi de date pentru detectarea comunităților în rețele biologice multi-relaționale - aplicația MultiNetCom

Considerăm 11 seturi de date luate din (14) și din (40). Acestea conțin gene implicate în 11 tipuri de cancer cum ar fi: leucemia limfoblastică acută (ALL), leucemia mieloidă acută (AML), carcinomul mamar (BC), adenocarcinomul colorectal (CA), limfomul non-Hodgkin (NHL) Carcinomul spastic (ST), cancerul gastric (GC), cancerul hepatic (LIC), cancerul pulmonar (LUC), cancerul sistemului nervos (NSC). Apoi, utilizăm Genemania (50), o platformă software care generează rețele de interacțiune gene-gene din milioane de publicații disponibile în diferite baze de date medicale și biologice. Obținem o rețea multi-relațională pentru fiecare tip de cancer. Genele reprezintă nodurile rețelei

și interacțiunile dintre gene reprezintă legăturile sau muchiile. În rețelele pe care le-am obținut, s-au generat următoarele tipuri de interacțiuni: interacțiuni prezise, co-localizare, coexpresie, domenii de proteine partajate, interacțiuni fizice, interacțiuni genetice, interacțiuni pe căi.

Tabelul 8.1 conține detaliile pentru fiecare rețea: numărul nodurilor, numărul muchiilor și numărul tipurilor de muchii.

Tip de cancer	Datele de cancer		
	Nr noduri	Nr muchii	Tipul muchiilor
ALL	47	260	6
AML	92	822	7
BC	26	163	7
CA	39	411	7
NHL	19	49	6
NSCLC	35	166	6
ST	21	52	5
GC	511	28148	7
LIC	322	9694	7
LUC	74	687	7
NSC	47	20276	7

Table 8.1: Descrierea seturilor de date folosite in experimente.

8.3.2 Rezultate

Realizăm două teste pentru fiecare rețea, separat: pentru prima, reducem rețeaua multi-relațională la o rețea uni-relațională și apoi aplicăm algoritmul inițial al lui Fruchterman pentru a obține comunitățile. Pentru al doilea test, considerăm rețeaua multi-relațională, și apoi aplicăm versiunea nouă a algoritmului lui Fruchterman propusă de noi în această lucrare, pentru a obține comunitățile. Figura 8.1 arată numărul de comunități și numărul de gene din fiecare comunitate obținute de fiecare algoritm pentru cele 11 seturi de date.

Din rezultatele prezentate aici, este evident faptul că folosind cele două abordări obținem comunități diferite. Numărul comunităților obținute pentru rețeaua uni-relațională este aproape întotdeauna diferit de numărul comunităților obținute de rețeaua multi-

8.3 Experimente și rezultate

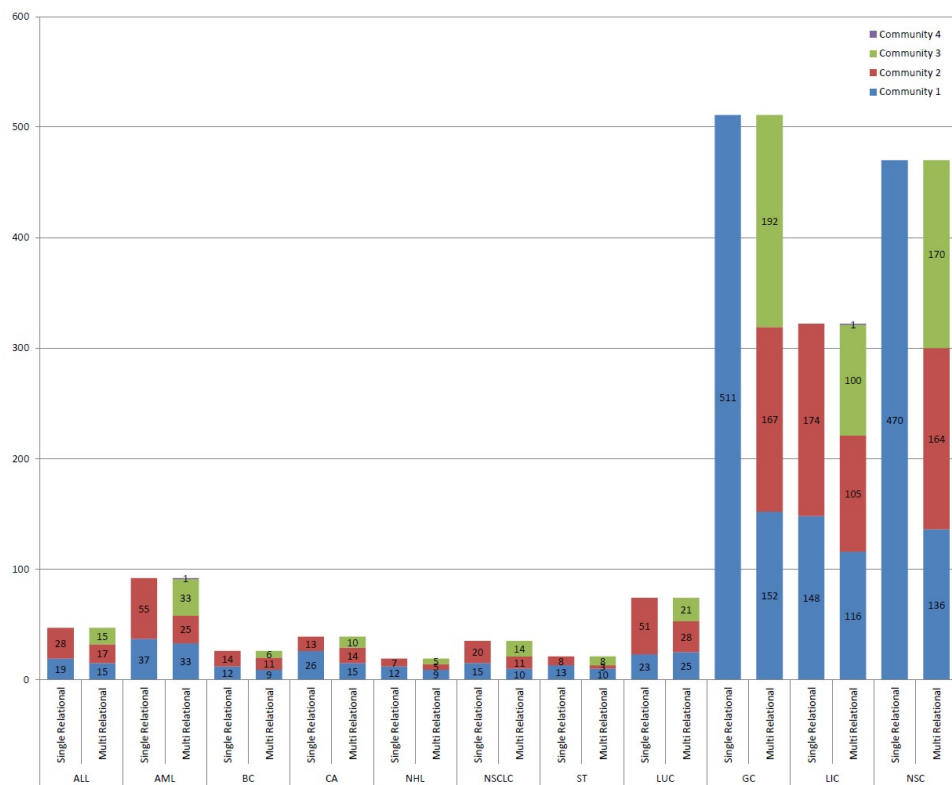


Figure 8.1: Numărul comunităților și genelor în fiecare comunitate obținută prin una din două metode: rețea uni-relațională transformată și rețeaua multi-relațională originală.

relațională, așa cum se poate vedea în figura 8.1. Mai mult decât atât, raportul de suprapunere între oricare două comunități este foarte scăzut pentru anumite rețele, fără gene comune (ca în cazul rețelelor ALL, AML, NHL, ST și LIC) sau doar o genă în comun (ca în cazul Din rețelele AML, BC, CA și LIC). În unele cazuri, unele comunități obținute în rețeaua multi-relațională sunt incluse în comunitățile obținute în rețeaua uni-relațională. Merită menționat faptul că testele multi-relaționale generează întotdeauna mai multe comunități decât cele unice.

8.4 Rezumat

Lucrarea prezentată aici are scopul de a demonstra că rețelele multi-relaționale și operațiile computaționale aplicate acestora trebuie să păstreze aspectul multi-relațional și nu pot fi reduse la rețele uni-relaționale. Transformarea unei rețele multi-relaționale într-o rețea uni-relațională poate reduce în mod semnificativ costul calculului, dar duce, de obicei, la pierderea informațiilor. În această lucrare propunem un algoritm pentru detectarea comunităților în rețele multi-relaționale, dezvoltat inițial pentru rețele uni-relaționale. Algoritmul a fost extins la rețele multi-relaționale, iar comparațiile dintre cei doi algoritmi au fost efectuate pe 11 seturi de date luate din date implicate în 11 tipuri de cancer care reprezintă rețele genetice de interacțiune (interacțiune genă-gena). Rezultatele arată diferențe semnificative între cele două abordări care ar putea fi luate în considerare de către cercetătorii care lucrează în domeniul biomedical. Algoritmul pe care l-am propus este unul general, astfel încât orice algoritm (nu doar algoritmul lui Fruchterman prezentat aici) poate fi extins cu ușurință de la rețelele uni-relaționale la cele multi-relaționale.

În acest capitol se propune un nou algoritm de detectare a comunităților pentru rețele multi-relaționale (46). Am aplicat acest algoritm în rețele biologice multi-relaționale și am comparat rezultatele cu cele produse de un algoritm de detectare a comunităților similare pentru rețelele uni-relaționale. Am arătat că algoritmul nostru este capabil să detecteze mai multe comunități, în comparație cu un algoritmul implementat pentru

rețele uni-relaționale, deoarece informațiile se pierd atunci când tratăm rețeaua inițială multi-relațională ca fiind o rețea uni-relațională. Rezultatele prezentate în acest capitol sunt incluse într-un articol trimis la revista PlosOne.

9

Concluzii

Ultimul capitol analizează succint eforturile de dezvoltare și analiză a modelelor de rețele multi-relaționale și a clustering-ului pe serii de timp multidimensionale și identifică câteva direcții generale pentru extinderea modelelor dezvoltate în această teză, care pot da perspective suplimentare. În continuare sunt descrise mai multe probleme deschise în analiza rețelelor multi-relaționale și gruparea multidimensională a seriilor de timp.

Această teză a prezentat două idei principale care sunt utilizate aici pentru abordarea problemelor biologice. Cu toate acestea, ele nu sunt construite special pentru biologie sau medicină și pot fi aplicate și în alte domenii.

Primul subiect se referă la rețele multi-relaționale (rețele în care există mai multe linkuri și mai mult de un tip de muchii între două noduri). Ea abordează trei aspecte ale acestor rețele:

- Comparația rețelelor multi-dimensionale,
- Căutarea de graphlets și motive în rețele multi-dimensionale,
- Detectarea comunităților în rețele multi-relaționale.

Provocările întâlnite în această parte a tezei au fost:

- Abordarea unor problem NP-complete (cum ar fi izomorfismul rețelelor),

-
- Depistarea algoritmilor eficienți pentru graphlets și motive,
 - Găsirea modului corect de a reprezenta rețelele și interacțiunile acestora,
 - Implementarea măsurilor potrivite pentru compararea rețelelor,
 - Detectarea comunităților în rețele multi-relaționale.

Totuși, încă mai trebuie depuse multe eforturi pentru o comparație cuprinzătoare a rețelelor și pentru implementarea unor algoritmi și mai eficienți pentru găsirea motivelor și detectarea comunităților. Am realizat până în prezent, următoarele:

- Proiectarea și implementarea unei aplicații de comparație a rețelelor:
 - Este nevoie de două sau mai multe rețele
 - Este conceput pentru rețele de interacțiune genetică și rețele de interacțiune proteină-proteină, dar este ușor de extensibil și adaptabil
 - Implementează câteva operații de bază, cum ar fi uniunea, intersecția, sub-grafurile comune, distribuția gradelor
 - Permite selectarea a unul, două, mai multe sau toate tipurile de link-uri existente și efectuează operațiuni pentru tipurile de linkuri selectate
- Proiectarea și implementarea unei aplicații pentru detectarea de graphlets și a motive:
 - Proiectat pentru rețele multi-relaționale
 - Găsește motive de o anumită dimensiune dată (sunt implementate mărimi de 3, 4 și 5 noduri)
 - Găsește motive customized
 - Lucrează cu grafuri direcționate și nedirecționate (rețele)
- Experimente și aplicații pentru problemele biologice:

-
- Pentru rețelele de gene implicate în trei tipuri similare de cancer feminin (ovarian, endometrial și sân)
 - Compararea rețelelor de interacțiune genetică bazate pe motivele conținute (adică blocurile de bază)
 - Proiectarea și implementarea unei aplicații pentru clustering-ul multidimensional:
 - Entitățile formate din mai multe serii de timp sunt grupate
 - Este efectuată o comparație a unei similitudini bazate pe distanță geometrică

10

Cuvinte cheie

- bioinformatică
- rețele complexe
- rețele cu o singură relație
- rețele multi-relaționale
- rețele biologice multi-relaționale
- clusters
- clustering multidimensional al seriilor de timp
- motive în graf
- motive în rețele multi-relaționale
- detectarea comunităților
- comunități în rețele multi-relaționale
- rețele de interacțiune genetică
- analiza rețelelor
- aplicații soft

Referințele tezei

- [1] AARTS/KORST. Simulated annealing and boltzmann machines. A stochastic approach to combinatorial optimization and neural computing. John Wiley., 1990.
- [2] Gaurav Agarwal and David Kempe. Modularitymaximizing graph communities via mathematical programming. *The European Physical Journal B*, 66(3):409-418, 2008.
- [3] Edoardo M Airoldi, David M Blei, Stephen E Fienberg, and Eric P Xing. Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9(Sep):1981-2014, 2008.
- [4] Reka Albert and Albert-Laszlo Barabasi. Statistical mechanics of complex networks. *Reviews of modern physics*, 74(1):47, 2002.
- [5] Reka Albert, Hawoong Jeong, and Albert-Laszlo Barabasi. Error and attack tolerance of complex networks. *arXiv preprint cond-mat/0008064*, 2000.
- [6] Charles J Alpert and So-Zen Yao. Spectral partitioning: the more eigenvectors, the better. In *Proceedings of the 32nd annual ACM/IEEE Design Automation Conference*, pages 195-200. ACM, 1995.
- [7] Khaled Alsabti, Sanjay Ranka, and Vineet Singh. An efficient k-means clustering

algorithm. 1997.

[8] Michael R Anderberg. Cluster analysis for applications. monographs and textbooks on probability and mathematical statistics, 1973.

[9] E Michael Azoff. Neural network time series forecasting of financial markets. John Wiley and Sons, Inc., 1994.

[10] Gary D Bader, Doron Betel, and Christopher WV Hogue. Bind: the biomolecular interaction network database. *Nucleic acids research*, 31(1):248-250, 2003.

[11] Albert-Laszlo Barabasi and Zoltan N Oltvai. Network biology: understanding the cell's functional organization. *Nature Reviews Genetics*, 5(2):101-113, 2004.

[12] Alain Barrat, Marc Barthelemy, Romualdo Pastor-Satorras, and Alessandro Vespignani. The architecture of complex weighted networks. *Proceedings of the National Academy of Sciences of the United States of America*, 101(11):3747-3752, 2004.

[13] Federico Battiston, Jacopo Iacovacci, Vincenzo Nicosia, Ginestra Bianconi, and Vito Latora. Emergence of multiplex communities in collaboration networks. *PloS one*, 11(1):e0147451, 2016.

[14] Anais Baudot, Gonzalo Gomez-Lopez, and Alfonso Valencia. Translational disease interpretation with molecular networks. *Genome biology*, 10(6):221, 2009.

[15] Laura Bennett, Aristotelis Kittas, Gareth Muirhead, Lazaros G Papageorgiou, and Sophia Tsoka. Detection of composite communities in multiplex biological networks. *Scientific reports*, 5, 2015.

[16] Laura Bennett, Songsong Liu, Lazaros G Papageorgiou, and Sophia Tsoka. Detection of disjoint and overlapping modules in weighted complex networks. *Advances in Complex Systems*, 15(05):1150023, 2012.

[17] Michele Berlingerio, Michele Coscia, Fosca Giannotti, Anna Monreale, and Dino

Pedreschi. Foundations of multidimensional network analysis. In Advances in Social Networks Analysis and Mining (ASONAM), 2011 International Conference on, pages 485-489. IEEE, 2011.

[18] Bela Bollobas, Gautam Das, Dimitrios Gunopulos, and Heikki Mannila. Time-series similarity problems and well-separated geometric sets. In Proceedings of the thirteenth annual symposium on Computational geometry, pages 454-456. ACM, 1997.

[19] David E Booth. Time series, 1992.

[20] Stefan Bornholdt and Heinz Georg Schuster. Handbook of graphs and networks: from the genome to the internet. John Wiley and Sons, 2006.

[21] Laurie A Boyer, Tong Ihn Lee, Megan F Cole, Sarah E Johnstone, Stuart S Levine, Jacob P Zucker, Matthew G Guenther, Roshan M Kumar, Heather L Murray, Richard G Jenner, et al. Core transcriptional regulatory circuitry in human embryonic stem cells. *cell*, 122(6):947-956, 2005.

[22] Ulrik Brandes. A faster algorithm for betweenness centrality*. *Journal of Mathematical Sociology*, 25(2):163- 177, 2001.

[23] Rainer Breitling, David Gilbert, Monika Heiner, and Richard Orton. A structured approach for the engineering of biochemical network models, illustrated for signalling pathways. *Briefings in Bioinformatics*, 9(5):404-421, 2008.

[24] Karin Breuer, Amir K Foroushani, Matthew R Laird, Carol Chen, Anastasia Sribnaia, Raymond Lo, Geoffrey L Winsor, Robert EW Hancock, Fiona SL Brinkman, and David J Lynn. Innatedb: systems biology of innate immunity and beyond recent updates and continuing curation. *Nucleic acids research*, 41(D1):D1228-D1233, 2013.

[25] Piotr Brodka, Tomasz Filipowski, and Przemyslaw Kazienko. An introduction to community detection in multi-layered social network. In World Summit on Knowledge Society, pages 185-190. Springer, 2011.

-
- [26] Deng Cai, Zheng Shao, Xiaofei He, Xifeng Yan, and Jiawei Han. Community mining from multi-relational networks. In European Conference on Principles of Data Mining and Knowledge Discovery, pages 445-452. Springer, 2005.
- [27] Muffy Calder, Stephen Gilmore, and Jane Hillston. Modelling the influence of rkip on the erk signalling pathway using the stochastic process algebra pepa. In Transactions on computational systems biology VII, pages 445-452. Springer, 2006.
- [28] Melissa S Cline, Michael Smoot, Ethan Cerami, Allan Kuchinsky, Neri Landys, Chris Workman, Rowan Christmas, Iliana Avila-Campilo, Michael Creech, Benjamin Gross, et al. Integration of biological networks and gene expression data using cytoscape. *Nature protocols*, 2(10):2366-2382, 2007.
- [29] Michele Coscia, Fosca Giannotti, and Dino Pedreschi. A classification for community discovery methods in complex networks. *Statistical Analysis and Data Mining*, 4(5):512-546, 2011.
- [30] Nello Cristianini, John Shawe-Taylor, and Jaz S Kandola. Spectral kernel methods for clustering. In *Advances in neural information processing systems*, pages 649-655, 2002.
- [31] Gautam Das, Dimitrios Gunopulos, and Heikki Mannila. Finding similar time series. In *European Symposium on Principles of Data Mining and Knowledge Discovery*, pages 88-100. Springer, 1997.
- [32] Darcy A Davis and Nitesh V Chawla. Exploring and exploiting disease interactions from multi-relational gene and phenotype networks. *PloS one*, 6(7):e22670, 2011.
- [33] Kara Dolinski, Andrew Chatr-aryamontri, and Mike Tyers. Systematic curation of protein and genetic interaction data for computable biology. *BMC biology*, 11(1):43, 2013.
- [34] Benjamin S Duran and Patrick L Odell. *Cluster analysis: a survey*. Springer-Verlag

New York, 1974.

[35] Ron Edgar, Michael Domrachev, and Alex E Lash. Gene expression omnibus: Ncbi gene expression and hybridization array data repository. *Nucleic acids research*, 30(1):207-210, 2002.

[36] William H Elliott, Daphne C Elliott, John R Jefferson, and John Wheldrake. *Biochemistry and molecular biology*. Oxford University Press Oxford, 1997.

[37] Graphical Enumeration. F. harary and em palmer, 1973.

[38] Paul Erdos and Alfred Renyi. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci*, 5(1):17- 60, 1960.

[39] Brian S Everitt. Unresolved problems in cluster analysis. *Biometrics*, pages 169-181, 1979.

[40] <http://cancer.sanger.ac.uk/census/>.

[41] Santo Fortunato. Community detection in graphs. *Physics reports*, 486(3):75-174, 2010.

[42] Andrea Franceschini, Damian Szklarczyk, Sune Frankild, Michael Kuhn, Milan Simonovic, Alexander Roth, Jianyi Lin, Pablo Minguéz, Peer Bork, Christian von Mering, et al. String v9. 1: protein-protein interaction networks, with increased coverage and integration. *Nucleic acids research*, 41(D1):D808-D815, 2013.

[43] Marco Franciosi and Giulia Menconi. Multi-dimensional sparse time series: feature extraction. *arXiv preprint arXiv:0803.0405*, 2008.

[44] Thomas MJ Fruchterman and Edward M Reingold. Graph drawing by force-directed placement. *Software: Practice and experience*, 21(11):1129-1164, 1991.

[45] Sylvia Fruhwirth-Schnatter, Christoph Paminger, Rudolf Winter-Ebmer, and Andrea Weber. Model-based clustering of categorical time series with multinomial logit

classification. In AIP Conference Proceedings, volume 1281, pages 1897-1900. AIP, 2010.

[46] Guojun Gan, Chaoqun Ma, and Jianhong Wu. Data clustering: theory, algorithms, and applications, volume 20. Siam, 2007.

[47] Michael R Garey and David S Johnson. Crossing number is np-complete. SIAM Journal on Algebraic Discrete Methods, 4(3):312-316, 1983.

[48] Claude Gerard, Didier Gonze, and Albert Goldbeter. Effect of positive feedback loops on the robustness of oscillations in the network of cyclin-dependent kinases driving the mammalian cell cycle. FEBS Journal, 279(18):3411-3431, 2012.

[49] Michelle Girvan and Mark EJ Newman. Community structure in social and biological networks. Proceedings of the national academy of sciences, 99(12):7821-7826, 2002.

[50] Roger Guimera and Luis A Nunes Amaral. Functional cartography of complex metabolic networks. Nature, 433(7028):895-900, 2005.

[51] Roger Guimera, Marta Sales-Pardo, and Luis A Nunes Amaral. A network-based method for target selection in metabolic networks. Bioinformatics, 23(13):1616-1622, 2007.

[52] Dimitrios Gunopulos and Gautam Das. Time series similarity measures (tutorial pm-2). In Tutorial notes of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining, pages 243-307. ACM, 2000.

[53] Jing-Dong J Han, Nicolas Bertin, Hao Tong, Debra S Goldberg, et al. Evidence for dynamically organized modularity in the yeast protein-protein interaction network. Nature, 430(6995):88, 2004.

[54] Frank Harary and Edgar M Palmer. Graphical enumeration. Technical report,

DTIC Document, 1973.

[55] John A Hartigan. Clustering algorithms. John Wiley and Sons, Inc., 1975.

[56] John A Hartigan and Manchek A Wong. Algorithm as 136: A k-means clustering algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1):100-108, 1979.

[57] Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. *Social networks*, 5(2):109-137, 1983.

[58] William M Holmes. Time series: Sir maurice kendall and j. keith ord, (edward arnold, great britain, 1990) pp. 296, 1992.

[59] Thomas C Ings, Jose M Montoya, Jordi Bascompte, Nico Bluthgen, Lee Brown, Carsten F Dormann, Francois Edwards, David Figueroa, Ute Jacob, J Iwan Jones, et al. Review: Ecological networks-beyond food webs. *Journal of Animal Ecology*, 78(1):253-269, 2009.

[60] Kuhn Ip, Caroline Colijn, and Desmond S Lun. Analysis of complex metabolic behavior through pathway decomposition. *BMC systems biology*, 5(1):91, 2011.

[61] Negin Iranfar, Danny Fuller, and William F Loomis. Transcriptional regulation of post-aggregation genes in dictyostelium by a feed-forward loop involving gbf and lagc. *Developmental biology*, 290(2):460-469, 2006.

[62] Shalev Itzkovitz, Reuven Levitt, Nadav Kashtan, Ron Milo, Michael Itzkovitz, and Uri Alon. Coarse-graining and self-dissimilarity of complex networks. *Physical Review E*, 71(1):016127, 2005.

[63] Ariel Jaimovich, Gal Elidan, Hanah Margalit, and Nir Friedman. Towards an integrated protein-protein interaction network: A relational markov network approach. *Journal of Computational Biology*, 13(2):145-164, 2006.

-
- [64] Anil K Jain and Richard C Dubes. Algorithms for clustering data. Prentice-Hall, Inc., 1988.
- [65] E Ferrell James Jr. Feedback loops and reciprocal regulation: recurring motifs in the systems biology of the cell cycle. *Current Opinion in Cell Biology*, 25(6):676- 686, 2013.
- [66] Hawoong Jeong, Sean P Mason, Albert-Laszlo Barabasi, and Zoltan N Oltvai. Lethality and centrality in protein networks. arXiv preprint cond-mat/0105306, 2001.
- [67] Hawoong Jeong, Balint Tombor, Reka Albert, Zoltan N Oltvai, and A-L Barabasi. The large-scale organization of metabolic networks. *Nature*, 407(6804):651-654, 2000.
- [68] Stephen C Johnson. Hierarchical clustering schemes. *Psychometrika*, 32(3):241-254, 1967.
- [69] Bjorn H Junker and Falk Schreiber. Analysis of biological networks, volume 2. John Wiley and Sons, 2008.
- [70] Minoru Kanehisa. A database for post-genome analysis. *Trends in genetics: TIG*, 13(9):375, 1997.
- [71] Minoru Kanehisa. The kegg database. In *Novartis Found Symp*, volume 247, pages 91-101, 2002.
- [72] Minoru Kanehisa, Susumu Goto, Shuichi Kawashima, Yasushi Okuno, and Masahiro Hattori. The kegg resource for deciphering the genome. *Nucleic acids re- search*, 32(suppl 1):D277-D280, 2004. 23
- [73] Ravi Kannan, Santosh Vempala, and Adrian Vetta. On clusterings: Good, bad and spectral. *Journal of the ACM (JACM)*, 51(3):497-515, 2004.
- [74] Tapas Kanungo, David M Mount, Nathan S Netanyahu, Christine D Piatko, Ruth Silverman, and Angela Y Wu. An efficient k-means clustering algorithm: Analysis

and implementation. *IEEE transactions on pattern analysis and machine intelligence*, 24(7):881-892, 2002.

[75] Nadav Kashtan, Shalev Itzkovitz, Ron Milo, and Uri Alon. Mfinder tool guide. Department of Molecular Cell Biology and Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot Israel, Tech. Rep, 2002.

[76] Nadav Kashtan, Shalev Itzkovitz, Ron Milo, and Uri Alon. Efficient sampling algorithm for estimating subgraph concentrations and detecting network motifs. *Bioinformatics*, 20(11):1746-1758, 2004.

[77] Leonard Kaufman and Peter J Rousseeuw. Finding groups in data: an introduction to cluster analysis, volume 344. John Wiley and Sons, 2009.

[78] Jinki Kim and Gwan-Su Yi. Rmod: A tool for regulatory motif detection in signaling network. *PloS one*, 8(7):e68407, 2013.

[79] Mikko Kivela, Alex Arenas, Marc Barthelemy, James P Gleeson, Yamir Moreno, and Mason A Porter. Multilayer networks. *Journal of complex networks*, 2(3):203-271, 2014.

[80] Mark Kozdoba and Shie Mannor. Community detection via measure space embedding. In *Advances in Neural Information Processing Systems*, pages 2890-2898, 2015.

[81] Cho Kwang-Hyun, Shin Sung-Young, Kim Hyun-Woo, Olaf Wolkenhauer, Brian McFerran, and Walter Kolch. Mathematical modeling of the influence of rkip on the erk signaling pathway. In *Computational methods in systems biology*, pages 127-141. Springer, 2003.

[82] Godfrey N Lance and William Thomas Williams. A general theory of classificatory sorting strategies 1. hierarchical systems. *The computer journal*, 9(4):373-380, 1967.

[83] Tong Ihn Lee, Nicola J Rinaldi, Francois Robert, Duncan T Odom, Ziv Bar-Joseph,

Georg K Gerber, Nancy M Hannett, Christopher T Harbison, Craig M Thompson, Itamar Simon, et al. Transcriptional regulatory networks in *saccharomyces cerevisiae*. *science*, 298(5594):799-804, 2002.

[84] Xutao Li, Michael K Ng, and Yunming Ye. Multicomm: Finding community structure in multi-dimensional networks. *IEEE Transactions on Knowledge and Data Engineering*, 26(4):929-941, 2014.

[85] T Warren Liao. Clustering of time series dataa survey. *Pattern recognition*, 38(11):1857-1874, 2005.

[86] Luana Licata, Leonardo Briganti, Daniele Peluso, Livia Perfetto, Marta Iannuccelli, Eugenia Galeota, Francesca Sacco, Anita Palma, Aurelio Pio Nardoza, Elena Santonico, et al. Mint, the molecular interaction database: 2012 update. *Nucleic acids research*, 40(D1):D857-D861, 2012.

[87] M. Lichman. UCI machine learning repository, 2013.

[88] Aristidis Likas, Nikos Vlassis, and Jakob J Verbeek. The global k-means clustering algorithm. *Pattern recognition*, 36(2):451-461, 2003.

[89] Azi Lipshtat, Sudarshan P Purushothaman, Ravi Iyengar, and Avi Maayan. Functions of bifans in context of multiple regulatory motifs in signaling networks. *Biophysical journal*, 94(7):2566-2579, 2008.

[90] James MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281-297. Oakland, CA, USA., 1967.

[91] Shmoolik Mangan and Uri Alon. Structure and function of the feed-forward loop network motif. *Proceedings of the National Academy of Sciences*, 100(21):11980-11985, 2003.

-
- [92] Sergei Maslov, Kim Sneppen, and Uri Alon. Correlation profiles and motifs in complex networks. *Handbook of Graphs and Networks: From the Genome to the Internet*, pages 168-198, 2003.
- [93] Louis L McQuitty. Elementary linkage analysis for isolating orthogonal and oblique types and typical relevancies. *Educational and Psychological Measurement*, 1957.
- [94] Marina Meila and Jianbo Shi. Learning segmentation by random walks. In *NIPS*, volume 14, 2000.
- [95] Manuel Middendorf, Etay Ziv, and Chris H Wiggins. Inferring network mechanisms: the drosophila melanogaster protein interaction network. *Proceedings of the National Academy of Sciences of the United States of America*, 102(9):3192-3197, 2005.
- [96] Ron Milo, Shalev Itzkovitz, Nadav Kashtan, Reuven Levitt, Shai Shen-Orr, Inbal Ayzenshtat, Michal Sheffer, and Uri Alon. Superfamilies of evolved and designed networks. *Science*, 303(5663):1538-1542, 2004.
- [97] Ron Milo, Shai Shen-Orr, Shalev Itzkovitz, Nadav Kashtan, Dmitri Chklovskii, and Uri Alon. Network motifs: simple building blocks of complex networks. *Science*, 298(5594):824-827, 2002.
- [98] Alexander Y Mitrophanov and Eduardo A Groisman. Positive feedback in cellular control systems. *Bioessays*, 30(6):542-555, 2008.
- [99] Sara Mostafavi, Debajyoti Ray, David Warde-Farley, Chris Grouios, Quaid Morris, et al. Genemania: a realtime multiple association network integration algorithm for predicting gene function. *Genome Biol*, 9(Suppl 1):S4, 2008.
- [100] Panagiotis Moulos, Julie Klein, Simon Jupp, Robert Stevens, Jean-Loup Bascands, and Joost P Schanstra. The kupnetviz: a biological network viewer for multipleomics datasets in kidney diseases. *BMC bioinformatics*, 14(1):235, 2013.
- [101] KA Abdul Nazeer and MP Sebastian. Improving the accuracy and efficiency of the

k-means clustering algorithm. In Proceedings of the World Congress on Engineering, volume 1, pages 1-3, 2009. 109 [102] Mark EJ Newman. Models of the small world. Journal of Statistical Physics, 101(3):819-841, 2000.

[103] Mark EJ Newman. Clustering and preferential attachment in growing networks. Physical review E, 64(2):025102, 2001.

[104] Mark EJ Newman. The structure of scientific collaboration networks. Proceedings of the National Academy of Sciences, 98(2):404-409, 2001.

[105] Mark EJ Newman. The structure and function of complex networks. SIAM review, 45(2):167-256, 2003.

[106] Mark EJ Newman and Michelle Girvan. Finding and evaluating community structure in networks. Physical review E, 69(2):026113, 2004.

[107] Andrew Y Ng, Michael I Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In Advances in neural information processing systems, pages 849-856, 2002 [108] Michael Kwok-Po Ng, Xutao Li, and Yunming Ye. Multirank: co-ranking for objects and relations in multirelational data. In Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 1217-1225. ACM, 2011.

[109] BO Palsson. Properties of reconstructed networks. Cambridge: Systems Biology, 2006.

[110] Tony Pawson and Rune Linding. Network medicine. FEBS letters, 582(8):1266-1270, 2008.

[111] Matteo Pellegrini, Edward M Marcotte, Michael J Thompson, David Eisenberg, and Todd O Yeates. Assigning protein functions by comparative genome analysis: protein phylogenetic profiles. Proceedings of the National Academy of Sciences, 96(8):4285-4288, 1999.

-
- [112] Steven J Phillips. Acceleration of k-means and related clustering algorithms. In *Algorithm Engineering and Experiments*, pages 166-177. Springer, 2002.
- [113] V Pihur, Somnath Datta, and Susmita Datta. Finding common genes in multiple cancer types through meta- analysis of microarray experiments: A rank aggregation approach. *Genomics*, 92(6):400-403, 2008.
- [114] Ivan Plavec, Oksana Sirenko, Sylvie Privat, Yuke Wang, Maya Dajee, Jennifer Melrose, Brian Nakao, Evangelos Hytopoulos, Ellen L Berg, and Eugene C Butcher. Method for analyzing signaling networks in complex cellular systems. *Proceedings of the National Academy of Sciences of the United States of America*, 101(5):1223-1228, 2004.
- [115] Teresa Przytycka. An important connection between network motifs and parsimony models. In *Research in Computational Molecular Biology*, pages 321-335. Springer, 2006.
- [116] Marcos G Quiles, Elbert EN Macau, and Nicolas Rubido. Dynamical detection of network communities. *Scientific reports*, 6, 2016.
- [117] Filippo Radicchi, Claudio Castellano, Federico Cecconi, Vittorio Loreto, and Domenico Parisi. Defining and identifying communities in networks. *Proceedings of the National Academy of Sciences of the United States of America*, 101(9):2658-2663, 2004.
- [118] Silvia Rausanu, Crina Grosan, Zujian Wu, Ovidiu Parvu, Ramona Stoica, and David Gilbert. Computational models for inferring biochemical networks. *Neural Computing and Applications*, 26(2):299-311, 2015.
- [119] Erzsebet Ravasz, Anna Lisa Somera, Dale A Mongru, Zoltan N Oltvai, and A-L Barabasi. Hierarchical organization of modularity in metabolic networks. *science*, 297(5586):1551-1555, 2002.
- [120] F James Rohlf. 12 single-link clustering algorithms. *Handbook of Statistics*, 2:267-284, 1982.

-
- [121] Michal Ronen, Revital Rosenberg, Boris I Shraiman, and Uri Alon. Assigning numbers to the arrows: parameterizing a gene regulation network by using accurate expression kinetics. *Proceedings of the national academy of sciences*, 99(16):10555-10560, 2002.
- [122] Jill A Rosenfeld, Dina Amrom, Eva Andermann, Frederick Andermann, Martin Veilleux, Cynthia Curry, Jamie Fisher, Stephen Deputy, Arthur S Aylsworth, Cynthia M Powell, et al. Genotype-phenotype correlation in interstitial 6q deletions: a report of 12 new cases. *neurogenetics*, 13(1):31-47, 2012.
- [123] Martin Rosvall and Carl T Bergstrom. Maps of random walks on complex networks reveal community structure. *Proceedings of the National Academy of Sciences*, 105(4):1118-1123, 2008.
- [124] Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53-65, 1987.
- [125] Louis A Saddic, Barbel Huvermann, Staver Bezhani, Yanhui Su, Cara M Winter, Chang Seob Kwon, Richard P Collum, and Doris Wagner. The leafy target *lmi1* is a meristem identity regulator and acts together with *leafy* to regulate expression of cauliflower. *Development*, 133(9):1673-1682, 2006.
- [126] Rintaro Saito, Michael E Smoot, Keiichiro Ono, Johannes Ruscheinski, Peng-Liang Wang, Samad Lotia, Alexander R Pico, Gary D Bader, and Trey Ideker. A travel guide to cytoscape plugins. *Nature methods*, 9(11):1069-1076, 2012.
- [127] Regina Samaga and Steffen Klamt. Modeling approaches for qualitative and semi-quantitative analysis of cellular signaling networks. *Cell communication and signaling*, 11(1):43, 2013.
- [128] S Schuster, DA Fell, and T Dandekar. A general definition of metabolic pathways

useful for systematic organization and analysis of complex metabolic networks. *Nature Biotechnology*, 18(3):326-332, 2000.

[129] <http://ccgd-starrlab.oit.umn.edu/search.php>.

[130] Jimmy Shadbolt and John Gerald Taylor. *Neural Networks and the Financial Markets: Predicting, Combining, and Portfolio Optimisation*. Springer, 2002.

[131] Paul Shannon, Andrew Markiel, Owen Ozier, Nitin S Baliga, Jonathan T Wang, Daniel Ramage, Nada Amin, Benno Schwikowski, and Trey Ideker. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome research*, 13(11):2498-2504, 2003.

[132] Junming Shao, Zhichao Han, and Qinli Yang. Community detection via local dynamic interaction. *arXiv preprint arXiv:1409.7978*, 2014.

[133] Peter HA Sneath. The application of computers to taxonomy. *Journal of general microbiology*, 17(1):201-226, 1957.

[134] Daniel A Spielmat and Shang-Hua Teng. Spectral partitioning works: Planar graphs and finite element meshes. In *Foundations of Computer Science, 1996. Proceedings., 37th Annual Symposium on*, pages 96-105. IEEE, 1996.

[135] Jorg Stelling, Steffen Klamt, Katja Bettenbrock, Stefan Schuster, and Ernst Dieter Gilles. Metabolic network structure determines key aspects of functionality and regulation. *Nature*, 420(6912):190-193, 2002.

[136] RAMONA STOICA. Multiobjective approach of multidimensional time series clustering. *Studia Universitatis Babes-Bolyai, Informatica*, 59(1), 2014.

[137] Ramona Stoica, Mihaela Ola, and Mihai Paraschivescu. Multidimensional temporal clustering: geometrical similarity measures analysis in k-means. *Intelligent Systems*

Design and Applications, 2013.

[138] Ramona Stoica, Bazil Parv, and Crina Grosan. Communities detection in multi-relational networks. *PlosOne*, 26(2):299-311, 2017.

[139] Ramona Stoica and Liviu Stirb. A tool for comparing multirelational networks from biology. In *Intelligent Systems Design and Applications (ISDA), 2013 13th International Conference on*, pages 242-246. IEEE, 2013.

[140] Gemma Swiers, Roger Patient, and Matthew Loose. Genetic regulatory networks programming hematopoietic stem cells and erythroid lineage specification. *Developmental biology*, 294(2):525-540, 2006.

[141] Damian Szklarczyk, Andrea Franceschini, Michael Kuhn, Milan Simonovic, Alexander Roth, Pablo Minguez, Tobias Doerks, Manuel Stark, Jean Muller, Peer Bork, et al. The string database in 2011: functional interaction networks of proteins, globally integrated and scored. *Nucleic acids research*, 39(suppl 1):D561-D568, 2011.

[142] Bosiljka Tadic, Miroslav Andjelkovic, Biljana Mileva Boshkoska, and Zoran Levnajic. Algebraic topology of multi-brain connectivity networks reveals dissimilarity in functional patterns during spoken communications. *PloS one*, 11(11):e0166787, 2016.

[143] Bing Tian Dai, Freddy Chong Tat Chua, and Ee-Peng Lim. Structural analysis in multi-relational social networks. In *Proceedings of the 2012 SIAM International Conference on Data Mining*, pages 451-462. SIAM, 2012.

[144] Hanghang Tong, Christos Faloutsos, and Jia-Yu Pan. Random walk with restart: fast solutions and applications. *Knowledge and Information Systems*, 14(3):327-346, 2008.

[145] Ruey S Tsay. Analysis of financial time series. *Financial econometrics*, a wiley-interscience publication, 2002.

[146] Joshua R Tyler, Dennis M Wilkinson, and Bernardo A Huberman. E-mail as

spectroscopy: Automated discovery of community structure within organizations. *The Information Society*, 21(2):143-153, 2005.

[147] Ikuo Uchiyama. Mbgd: microbial genome database for comparative analysis. *Nucleic acids research*, 31(1):58- 62, 2003.

[148] Olga Vechtomova. Introduction to information retrieval christopher d. manning, prabhakar raghavan, and hinrich schutze (stanford university, yahoo! research, and university of stuttgart) cambridge: Cambridge university press, 2008, xxi+ 482 pp; hard-bound, isbn 978-0- 521-86571-5, 2009.

[149] David Warde-Farley, Sylva L Donaldson, Ovi Comes, Khalid Zuberi, Rashad Badrawi, Pauline Chao, Max Franz, Chris Grouios, Farzana Kazi, Christian Tannus Lopes, et al. The genemania prediction server: biological network integration for gene prioritization and predicting gene function. *Nucleic acids research*, 38(suppl 2):W214-W220, 2010.

[150] Duncan J Watts and Steven H Strogatz. Collective dynamics of small-worldnetworks. *nature*, 393(6684):440- 442, 1998.

[151] Yen-Chuen Wei and Chung-Kuan Cheng. Towards efficient hierarchical designs by ratio cut partitioning. In *Computer-Aided Design, 1989. ICCAD-89. Digest of Technical Papers.*, 1989 IEEE International Conference on, pages 298-301. IEEE, 1989.

[152] Sebastian Wernicke. Efficient detection of network motifs. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 3(4), 2006.

[153] Dennis M Wilkinson and Bernardo A Huberman. A method for Finding communities of related genes. *pro- ceedings of the national Academy of sciences*, 101(suppl 1):5241-5248, 2004.

[154] WT Williams and JM t Lambert. Multivariate methods in plant ecology: V. similarity analyses and information-analysis. *The Journal of Ecology*, pages 427-445, 1966.

-
- [155] Elisabeth Wong, Brittany Baur, Saad Quader, and Chun-Hsi Huang. Biological network motif detection: principles and practice. *Briefings in bioinformatics*, 13(2):202-215, 2011.
- [156] Elisabeth A Wong and Brittany Baur. On network tools for network motif finding: a survey study, 2010.
- [157] Zhiang Wu, Zhan Bu, Jie Cao, and Yi Zhuang. Discovering communities in multi-relational networks. In *User Community Discovery*, pages 75-95. Springer, 2015.
- [158] Stefan Wuchty. Scale-free behavior in protein domain networks. *Molecular biology and evolution*, 18(9):1694-1702, 2001.
- [159] Ioannis Xenarios, Danny W Rice, Lukasz Salwinski, Marisa K Baron, Edward M Marcotte, and David Eisenberg. Dip: the database of interacting proteins. *Nucleic acids research*, 28(1):289-291, 2000.
- [160] Gang Xu, Laura Bennett, Lazaros G Papageorgiou, and Sophia Tsoka. Module detection in complex networks using integer optimisation. *Algorithms for Molecular Biology*, 5(1):36, 2010.
- [161] Yoshihiro Yamanishi, Michihiro Araki, Alex Gutteridge, Wataru Honda, and Minoru Kanehisa. Prediction of drug-target interaction networks from the integration of chemical and genomic spaces. *Bioinformatics*, 24(13):i232-i240, 2008.
- [162] Muhammed A Yildirim, Kwang-Il Goh, Michael E Cusick, Albert-Laszlo Barabasi, and Marc Vidal. Drug- target network. *Nature biotechnology*, 25(10):1119, 2007.
- [163] Chang Hun You, Lawrence B Holder, and Diane J Cook. Application of graph-based data mining to metabolic pathways. In *Data Mining Workshops, 2006. ICDM Workshops 2006. Sixth IEEE International Conference on*, pages 169-173. IEEE, 2006.
- [164] Alon Zaslaver, Avi E Mayo, Revital Rosenberg, Pnina Bashkin, Hila Sberro,

Miri Tsalyuk, Michael G Surette, and Uri Alon. Just-in-time transcription program in metabolic pathways. *Nature genetics*, 36(5):486, 2004.

[165] Bin Zhang and Sargur N Srihari. Properties of binary vector dissimilarity measures. In *Proc. JCIS Intl Conf. Computer Vision, Pattern Recognition, and Image Processing*, volume 1, 2003.

[166] Bin Zhang and Sargur N Srihari. Fast k-nearest neighbor classification using cluster-based trees. *IEEE Transactions on Pattern analysis and machine intelligence*, 26(4):525-528, 2004.

[167] Yan Zhang, Zhifeng Yang, and Xiangyi Yu. Ecological network and emergy analysis of urban metabolic systems: model development, and a case study of four chinese cities. *Ecological Modelling*, 220(11):1431-1442, 2009.

[168] Chunxiang Zheng, Chad R Weisbrod, Juan D Chavez, Jimmy K Eng, Vagisha Sharma, Xia Wu, and James E Bruce. Xlink-db: Database and software tools for storing and visualizing protein interaction topology data. *Journal of proteome research*, 12(4):1989-1995, 2013.

[169] Xiaowei Zhu, Mark Gerstein, and Michael Snyder. Getting connected: analysis and principles of biological networks. *Genes and development*, 21(9):1010-1024, 2007.

References

- [1] Réka Albert and Albert-László Barabási. Statistical mechanics of complex networks. *Reviews of modern physics*, 74(1):47, 2002. 12
- [2] Réka Albert, Hawoong Jeong, and Albert-László Barabási. Error and attack tolerance of complex networks. *arXiv preprint cond-mat/0008064*, 2000. 12
- [3] E Michael Azoff. *Neural network time series forecasting of financial markets*. John Wiley & Sons, Inc., 1994. 22
- [4] Albert-László Barabási and Zoltan N Oltvai. Network biology: understanding the cell's functional organization. *Nature Reviews Genetics*, 5(2):101–113, 2004. 31
- [5] Béla Bollobás, Gautam Das, Dimitrios Gunopulos, and Heikki Mannila. Time-series similarity problems and well-separated geometric sets. In *Proceedings of the thirteenth annual symposium on Computational geometry*, pages 454–456. ACM, 1997. 22
- [6] David E Booth. Time series, 1992. 22
- [7] Stefan Bornholdt and Heinz Georg Schuster. *Handbook of graphs and networks: from the genome to the internet*. John Wiley & Sons, 2006. 29
- [8] Ulrik Brandes. A faster algorithm for betweenness centrality*. *Journal of Mathematical Sociology*, 25(2):163–177, 2001. 13
- [9] Melissa S Cline, Michael Smoot, Ethan Cerami, Allan Kuchinsky, Neri Landys, Chris Workman, Rowan Christmas, Iliana Avila-Campilo, Michael Creech, Benjamin Gross, et al. Integration of biological networks and gene expression data using cytoscape. *Nature protocols*, 2(10):2366–2382, 2007. 32
- [10] Michele Coscia, Fosca Giannotti, and Dino Pedreschi. A classification for community discovery methods in complex networks. *Statistical Analysis and Data Mining*, 4(5):512–546, 2011. 42
- [11] Gautam Das, Dimitrios Gunopulos, and Heikki Mannila. Finding similar time series. In *European Symposium on Principles of Data Mining and Knowledge Discovery*, pages 88–100. Springer, 1997. 22
- [12] Ron Edgar, Michael Domrachev, and Alex E Lash. Gene expression omnibus: Ncbi gene expression and hybridization array data repository. *Nucleic acids research*, 30(1):207–210, 2002. 36
- [13] Brian S Everitt. Unresolved problems in cluster analysis. *Biometrics*, pages 169–181, 1979. 21
- [14] <http://cancer.sanger.ac.uk/census/>. 42
- [15] Marco Franciosi and Giulia Menconi. Multi-dimensional sparse time series: feature extraction. *arXiv preprint arXiv:0803.0405*, 2008. 24
- [16] Sylvia Frühwirth-Schnatter, Christoph Pamminger, Rudolf Winter-Ebmer, and Andrea Weber. Model-based clustering of categorical time series with multinomial logit classification. In *AIP Conference Proceedings*, volume 1281, pages 1897–1900. AIP, 2010. 27
- [17] Guojun Gan, Chaoqun Ma, and Jianhong Wu. *Data clustering: theory, algorithms, and applications*, volume 20. Siam, 2007. 22
- [18] Roger Guimera and Luis A Nunes Amaral. Functional cartography of complex metabolic networks. *Nature*, 433(7028):895–900, 2005. 15
- [19] Dimitrios Gunopulos and Gautam Das. Time series similarity measures (tutorial pm-2). In *Tutorial notes of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 243–307. ACM, 2000. 22
- [20] Jing-Dong J Han, Nicolas Bertin, Hao Tong, Debra S Goldberg, et al. Evidence for dynamically organized modularity in the yeast protein-protein interaction network. *Nature*, 430(6995):88, 2004. 35
- [21] William M Holmes. Time series: Sir maurice kendall and j. keith ord, (edward arnold, great britain, 1990) pp. 296, 1992. 22
- [22] Kuhn Ip, Caroline Colijn, and Desmond S Lun. Analysis of complex metabolic behavior through pathway decomposition. *BMC systems biology*, 5(1):91, 2011. 15
- [23] Ariel Jaimovich, Gal Elidan, Hanah Margalit, and Nir Friedman. Towards an integrated protein-protein interaction network: A relational markov network approach. *Journal of Computational Biology*, 13(2):145–164, 2006. 35

REFERENCES

- [24] Hawoong Jeong, Sean P Mason, Albert-Laszlo Barabasi, and Zoltan N Oltvai. Lethality and centrality in protein networks. *arXiv preprint cond-mat/0105306*, 2001. 35
- [25] Hawoong Jeong, Bálint Tombor, Réka Albert, Zoltan N Oltvai, and A-L Barabási. The large-scale organization of metabolic networks. *Nature*, 407(6804):651–654, 2000. 15, 35
- [26] Minoru Kanehisa. The kegg database. In *Novartis Found Symp*, volume 247, pages 91–101, 2002. 18, 19
- [27] Minoru Kanehisa, Susumu Goto, Shuichi Kawashima, Yasushi Okuno, and Masahiro Hattori. The kegg resource for deciphering the genome. *Nucleic acids research*, 32(suppl 1):D277–D280, 2004. 18
- [28] Nadav Kashtan, Shalev Itzkovitz, Ron Milo, and Uri Alon. Mfinder tool guide. *Department of Molecular Cell Biology and Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot Israel, Tech. Rep*, 2002. 36
- [29] Jinki Kim and Gwan-Su Yi. Rmod: A tool for regulatory motif detection in signaling network. *PLoS one*, 8(7):e68407, 2013. 35
- [30] T Warren Liao. Clustering of time series data survey. *Pattern recognition*, 38(11):1857–1874, 2005. 22
- [31] Shmoolik Mangan and Uri Alon. Structure and function of the feed-forward loop network motif. *Proceedings of the National Academy of Sciences*, 100(21):11980–11985, 2003. 35
- [32] Ron Milo, Shai Shen-Orr, Shalev Itzkovitz, Nadav Kashtan, Dmitri Chklovskii, and Uri Alon. Network motifs: simple building blocks of complex networks. *Science*, 298(5594):824–827, 2002. 35
- [33] Sara Mostafavi, Debajyoti Ray, David Warde-Farley, Chris Grouios, Quaid Morris, et al. Genemania: a real-time multiple association network integration algorithm for predicting gene function. *Genome Biol*, 9(Suppl 1):S4, 2008. 20
- [34] Mark EJ Newman. Models of the small world. *Journal of Statistical Physics*, 101(3):819–841, 2000. 12
- [35] V Pihur, Somnath Datta, and Susmita Datta. Finding common genes in multiple cancer types through meta-analysis of microarray experiments: A rank aggregation approach. *Genomics*, 92(6):400–403, 2008. 31
- [36] Michal Ronen, Revital Rosenberg, Boris I Shraiman, and Uri Alon. Assigning numbers to the arrows: parameterizing a gene regulation network by using accurate expression kinetics. *Proceedings of the national academy of sciences*, 99(16):10555–10560, 2002. 35
- [37] Jill A Rosenfeld, Dina Amrom, Eva Andermann, Frederick Andermann, Martin Veilleux, Cynthia Curry, Jamie Fisher, Stephen Deputy, Arthur S Aylsworth, Cynthia M Powell, et al. Genotype–phenotype correlation in interstitial 6q deletions: a report of 12 new cases. *neurogenetics*, 13(1):31–47, 2012. 20
- [38] Rintaro Saito, Michael E Smoot, Keiichiro Ono, Johannes Ruscheinski, Peng-Liang Wang, Samad Lotia, Alexander R Pico, Gary D Bader, and Trey Ideker. A travel guide to cytoscape plugins. *Nature methods*, 9(11):1069–1076, 2012. 31, 32
- [39] S Schuster, DA Fell, and T Dandekar. A general definition of metabolic pathways useful for systematic organization and analysis of complex metabolic networks. *Nature Biotechnology*, 18(3):326–332, 2000. 15
- [40] <http://ccgd-starrlab.oit.umn.edu/search.php>. 42
- [41] Jimmy Shadbolt and John Gerald Taylor. *Neural Networks and the Financial Markets: Predicting, Combining, and Portfolio Optimisation*. Springer, 2002. 22
- [42] Paul Shannon, Andrew Markiel, Owen Ozier, Nitin S Baliga, Jonathan T Wang, Daniel Ramage, Nada Amin, Benno Schwikowski, and Trey Ideker. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome research*, 13(11):2498–2504, 2003. 31, 32
- [43] Jörg Stelling, Steffen Klamt, Katja Bettenbrock, Stefan Schuster, and Ernst Dieter Gilles. Metabolic network structure determines key aspects of functionality and regulation. *Nature*, 420(6912):190–193, 2002. 15
- [44] RAMONA STOICA. Multiobjective approach of multidimensional time series clustering. *Studia Universitatis Babeş-Bolyai, Informatica*, 59(1), 2014. 28
- [45] Ramona Stoica, Mihaela Ola, and Mihai Paraschivescu. Multidimensional temporal clustering: geometrical similarity measures analysis in k-means. *Intelligent Systems Design and Applications*, 2013. 28
- [46] Ramona Stoica, Bazil Parv, and Crina Grosan. Communities detection in multi-relational networks. *PLoSOne*, 26(2):299–311, 2017. 11, 45
- [47] Ramona Stoica and Liviu Stirb. A tool for comparing multirelational networks from biology. In *Intelligent Systems Design and Applications (ISDA), 2013 13th International Conference on*, pages 242–246. IEEE, 2013. 10, 33
- [48] Ruey S Tsay. Analysis of financial time series. financial econometrics, a wiley-intercience publication, 2002. 22

REFERENCES

- [49] Ikuo Uchiyama. Mbgd: microbial genome database for comparative analysis. *Nucleic acids research*, 31(1):58–62, 2003. 18
- [50] David Warde-Farley, Sylva L Donaldson, Ovi Comes, Khalid Zuberi, Rashad Badrawi, Pauline Chao, Max Franz, Chris Grouios, Farzana Kazi, Christian Tannus Lopes, et al. The genemania prediction server: biological network integration for gene prioritization and predicting gene function. *Nucleic acids research*, 38(suppl.2):W214–W220, 2010. 20, 42
- [51] Duncan J Watts and Steven H Strogatz. Collective dynamics of small-worldnetworks. *nature*, 393(6684):440–442, 1998. 14
- [52] Elisabeth A Wong and Brittany Baur. On network tools for network motif finding: a survey study, 2010. 35
- [53] Alon Zaslaver, Avi E Mayo, Revital Rosenberg, Pnina Bashkin, Hila Sberro, Miri Tsalyuk, Michael G Surette, and Uri Alon. Just-in-time transcription program in metabolic pathways. *Nature genetics*, 36(5):486, 2004. 35
- [54] Bin Zhang and Sargur N Srihari. Properties of binary vector dissimilarity measures. In *Proc. JCIS Intl Conf. Computer Vision, Pattern Recognition, and Image Processing*, volume 1, 2003. 22