

Universitatea Babeş–Bolyai, Cluj-Napoca
Facultatea de Matematică și Informatică

Învățare semisupervizată folosind kerneluri

Teză de doctorat
Rezumat

Conducător științific:
Zoltán KÁSA

Doctorand:
Zalán-Péter BODÓ

Cluj-Napoca
2009

Cuprinsul rezumatului

1	Introducere	4
1.1	Învățare semisupervizată și metode kernel	5
1.2	Structura tezei	6
2	Kerneluri de clasificare a documentelor bazate pe Wikipedia	7
2.1	Reprezentarea documentelor bazată pe Wikipedia	7
2.2	Reducerea dimensionalității kernelului Wikipedia	10
2.3	Structura pe linkuri a Wikipedia	10
2.4	Ponderarea conceptelor	11
3	Kerneluri cluster ierarhici	12
3.1	Clustering ierarhic	12
3.2	Distanțele de legătură	13
3.3	Construirea kernelului	14
3.4	Kerneluri cluster ierarhici cu distanțe pe graf	15
4	Kerneluri de reponderare	17
4.1	Kernelul bagged cluster	18
4.2	Kernelul de reponderare gaussian	19
4.3	Kerneluri de reponderare bazate pe produse scalare	19
5	Experimente	20
6	Concluzii	24

Cuprinsul tezei de doctorat

1	Introduction	1
2	Semi-supervised learning	7
2.1	Assumptions in semi-supervised learning	9
2.1.1	The smoothness assumption	9
2.1.2	The cluster assumption	9
2.1.3	The manifold assumption	10
2.2	Transduction	11
2.3	A classification of semi-supervised methods	11
2.3.1	Generative models	11
2.3.2	Low density separation	13
2.3.3	Graph-based methods	14

2.3.4	Change of representation	18
3	Kernels and kernel methods	21
3.1	A simple classification algorithm	24
3.2	Some general purpose kernels	25
3.2.1	The linear kernel	26
3.2.2	The polynomial kernel	26
3.2.3	The RBF kernel	27
3.2.4	The sigmoid kernel	28
3.3	Classification with Support Vector Machines	28
3.3.1	Hard margin SVMs	29
3.3.2	Soft margin SVMs	31
3.3.3	Kernelization	31
3.3.4	Classification with multiple classes	33
3.4	Dimensionality reduction with PCA and KPCA	36
3.4.1	Principal Component Analysis	36
3.4.2	Kernel Principal Component Analysis	39
4	Data-dependent kernels	41
4.1	The ISOMAP kernel	42
4.2	The neighborhood kernel	44
4.3	The bagged cluster kernel	44
4.4	Multi-type cluster kernel	45
4.4.1	Linear transfer function	45
4.4.2	Step transfer function	46
4.4.3	Linear step transfer function	47
4.4.4	Polynomial transfer function	47
4.5	Manifold regularization and data-dependent kernels for SSL using point cloud norms	48
5	Wikipedia-based kernels for text categorization	50
5.1	Text categorization	51
5.1.1	The bag-of-words representation	51
5.1.2	Feature selection techniques in text categorization	54
5.1.3	Machine learning in text categorization	59
5.1.4	Evaluation measures	62
5.2	String and text kernels	66
5.2.1	String kernels	66
5.2.2	The VSM kernel	68
5.2.3	The GVSM kernel	69

5.2.4	WordNet-based kernels	69
5.2.5	Latent Semantic Kernel	72
5.2.6	The von Neumann kernel	73
5.3	Wikipedia-based text kernels	74
5.3.1	Wikipedia	75
5.3.2	Wikipedia-based document representation	76
5.3.3	Dimensionality reduction for the Wikipedia kernel	79
5.3.4	Link structure of Wikipedia	80
5.3.5	Concept weighting	80
5.3.6	Experimental methodology and test results	82
5.3.7	Related methods	83
5.3.8	Discussion	84
6	Hierarchical cluster kernels	86
6.1	Motivation for a cluster kernel	86
6.2	Hierarchical clustering	88
6.2.1	Linkage distances	91
6.3	Metric multi-dimensional scaling	94
6.4	Ultrametric matrices and trees	95
6.5	The hierarchical cluster kernel	96
6.5.1	Hierarchical cluster kernel with graph distances	98
6.6	New test points	100
6.7	Experimental methodology and test results	101
6.8	Related work	105
6.9	Discussion	106
7	Variations on the bagged cluster kernel	108
7.1	The bagged cluster kernel	108
7.2	Computing the reweighting kernel	110
7.2.1	Combining kernels	110
7.2.2	Using the Hadamard product for kernel reweighting	112
7.3	Getting the clustering	116
7.4	Experimental methodology and test results	117
7.5	Discussion	120
8	Conclusions	121
A	Data sets	124
A.1	Two-moons	124
A.2	Reuters-21578	124

A.3 USPS	126
A.4 Digit1	127
A.5 COIL2	127
A.6 Text	127

Cuvinte cheie: instruire automată, învățare semisupervizată, kerneluri și metode kernel, kerneluri cluster, kerneluri semisupervizate

1 Introducere

Din anul 1956 inteligența artificială (*Artificial Intelligence – AI*) a devenit un domeniu intens studiat al informaticii, scopul ei fiind construcția mașinilor inteligente. Totuși inteligența nu este un concept bine definit și nici unul ușor de definit – iar aici ne referim la definiția informală a inteligenței. O tentativă de a demonstra inteligența automat – astfel o încercare de a o defini – a fost bine-cunoscutul test Turing: o ființă umană poartă o conversație într-o limbă naturală cu o altă ființă umană și cu o mașină, ambele încercând să pară a fi ființe umane, iar sarcina este de a determina care este ființa umană și care este mașina. Dacă acest lucru nu poate fi judecat în mod fiabil, mașina “trece examenul”. Astfel putem afirma că scopul AI este de a construi mașini care se comportă asemenea ființelor umane.

După Russel and Norvig [1995] sistemele AI pot fi organizate în patru grupuri: sisteme care gândesc într-un mod similar cu ființele umane, sisteme care pot acționa ca și ființele umane, sisteme care pot gândi rațional și sisteme care pot acționa în mod rațional. Învățarea automată (*Machine Learning – ML*) este subdomeniul AI, care dorește să modeleze cele mai importante activități ale creierului uman: clasificarea, diferențierea și predicția; mașinile instruibile aparțin primei categorii descrise mai sus. În ultimii 20 de ani științele cognitive și informatica s-au îndepărtat unul de altul și noi domenii mult mai înguste s-au format în cadrul AI. Unele tehnologii au devenit indispensabile, și în unele privințe au depășit ceea ce pot face ființele umane: să ne gândim doar la sistemele de regăsire a informației, unde milioane de documente sunt scanate pentru a prezenta utilizatorului cele care conțin informații relevante, iar în prezent acuratețea atinsă de către aceste sisteme este comparabilă cu judecata umană.

Două dintre cele mai importante subdomenii ale ML sunt învățarea supervizată și cea nesupervizată: în cazul învățării supervizate primim exemple împreună cu instrucțiuni de instruire pe care le numim etichete; învățarea nesupervizată este mai dificilă din cauza că nu sunt furnizate instrucțiuni adiționale. Sarcina obiș-

nuită în cazul învățării nesupervizate este determinarea clusterelor prin gruparea punctelor similare sau găsirea densității probabilistice a datelor.

Prezenta teză se concentrează pe învățarea semisupervizată (*Semi-Supervised Learning – SSL*): din cauza că adnotarea umană a exemplilor de instruire în general solicită expertiză în domeniu, aceasta este costisitoare și consumatoare de timp. O soluție este utilizarea numai a unei mici proporții a datelor etichetate pentru a îmbunătăți performanța algoritmului de instruire. De exemplu să presupunem că într-o problemă de clasificare a documentelor cuvântul “profesor” se dovedește a fi un bun mijloc de predicție a exemplilor pozitive pe baza datelor etichetate. Apoi, dacă datele neetichetate arată faptul că termenul “profesor” și “universitate” sunt corelate, prin utilizarea ambelor cuvinte se poate aștepta creșterea acurateții clasificatorului.

1.1 Învățare semisupervizată și metode kernel

Funcțiile kernel sau kerneluri returnează similaritățile între exemple. În cazul metodelor kernel nu lucrăm cu datele în mod direct, ci cu ajutorul unei matrice numită matricea kernel, ceea ce conține similaritățile între exemple. Aceasta la rândul ei este analoagă clasificării umane, unde similaritatea între exemple joacă un rol important [Estes, 1994]. Kerneluri sunt instrumente pentru extensia neliniară a metodelor lineare: dacă un algoritm poate fi scris în termenii produselor scalare, putem în mod simplu schimba matricea de produse scalare cu o funcție/matrice arbitrară pozitiv semidefinită care să conțină produsele scalare ale datelor în așa-numitul spațiu al caracteristicilor (*feature space*), și obținem o extensie neliniară a algoritmului.

Kerneluri dependente de date dau naștere la mașini instruibile semisupervizate: funcția kernel nu depinde în continuare numai de cele două puncte în discuție, ci utilizează data completă, informațiile cuprinse în setul complet de instruire disponibil. Matematic vorbind, dacă D_1 și D_2 notează două seturi de date, $D_1 \neq D_2$, $\mathbf{x}, \mathbf{z} \in D_1 \cap D_2$, atunci

$$k(\mathbf{x}, \mathbf{z}; D_1) \approx k(\mathbf{x}, \mathbf{z}; D_2)$$

unde prin $k(\cdot, \cdot)$ notăm funcțiile kernel, “;” însemnând condiționare și “ \approx ” însemnând “nu neaparat egal”. Kernelurile dependente de date sunt folosite pentru a îmbunătăți măsura de similaritate care ia în considerare doar data etichetată. Reprezentarea în spațiul caracteristicilor este acum aleasă folosind informațiile exploatate din seturile de date etichetate și neetichetate. Kernelurile de acest fel pot fi utilizate în oricare metodă kernel în cazul în care sunt disponibile și date neetichetate.

Subiectul acestei teze este construcția acestor kerneluri dependente de date pentru învățare supervizată și semisupervizată și de a dovedi superioritatea kernelurilor dependente de date față de kerneluri convenționale.

1.2 Structura tezei

Teza este organizată în opt capitole. Capitolul 1 este capitolul introductiv, prezentând pe scurt tema și structura tezei. Capitolul 2 constă într-o prezentare a tehnicilor SSL și a conceptelor asociate. După prezentarea prezumțiilor folosite în SSL, o clasificare a metodelor SSL este dată printr-o scurtă prezentare a unei metode din fiecare categorie.

În Capitolul 3 introducem metodele kernel și prezentăm o descriere detaliată a mașinilor cu suport vectorial (*Support Vector Machines – SVMs*), a metodei *Principal Component Analysis (PCA)* și a metodei *Kernel Principal Component Analysis (KPCA)*. Metodele prezentate aici vor fi utilizate mai târziu în cadrul tezei.

Capitolul 4 prezintă tehnicile de construcție a kernelurilor dependente de date, de exemplu kernelul ISOMAP, kernelul de vecinătate, kernelul bagged cluster, kernelul cluster multi-type și un kernel dependent de date legat de regularizări de varietăți (*manifold regularization*).

Următoarele trei capitole constituie partea principală a tezei: ele conțin principalele contribuții în domeniul construcției kernelurilor dependente de date. Capitolul 5 prezintă kernelul propus de autor bazat pe Wikipedia pentru categorizarea textelor (*Text Categorization – TC*). Acest capitol oferă în același timp o introducere detaliată în domeniul categorizării textelor.

Capitolul 6 prezintă kernelul cluster ierarhic propus pentru învățarea semisupervizată. S-au implementat kernelurile dependente de date prezentate în capitolul 4 și le-am comparat în mod experimental cu kernelul ierarhic pe diferite seturi de date.

În capitolul 7 propunem trei kerneluri cluster utilizând proprietatea produsului Hadamard a matricelor pozitiv semidefinite.

Capitolul 8 rezumă și încheie teza, iar Anexa A descrie seturile de date pe care le-am utilizat în evaluarea metodelor prezentate în teză.

Contribuția autorului în domeniu poate fi rezumată în următorul fel:

- un nou kernel pentru categorizarea documentelor, bazat pe informații extrase din Wikipedia
 - propunerea includerii structurii pe linkuri a Wikipediei în kernel
 - ponderarea conceptelor în kernelul bazat pe Wikipedia utilizând algoritmul PageRank

- o nouă metodă de construire a kernelurilor cluster ierarhici pentru învățarea supervizată și semisupervizată
 - propunerea unui cadru general pentru construcția kernelurilor cluster ierarhici
 - definirea kernelurilor cluster ierarhici și celor bazate pe grafuri
- construcția kernelurilor prin utilizarea proprietății produsului Hadamard a matricelor pozitiv semidefinite
 - introducerea kernelului gaussian de reponderare
 - introducerea a doi kerneluri de reponderare prin utilizarea produselor scalare ale vectorilor de apartenență cluster

În cele de mai jos prezentăm metodele propuse în teză, incluse în capitolele 5–7.

2 Kerneluri de clasificare a documentelor bazate pe Wikipedia

Clasificarea textelor sau a documentelor constituie problema determinării categoriilor predefinite reale ale documentelor scrise în limbaje naturale pe baza unor exemple de instruire. Constituie un subdomeniu al regăsirii informației (*Information Retrieval – IR*), dar de multe ori este considerată în mod eronat a fi o problemă de prelucrare a limbajelor naturale. Reprezentarea documentelor cu ajutorul modelului de spațiu vectorial (*Vector Space Model – VSM*) cel mai performant și utilizat pe scală largă creează vectori extrem de multidimensionali și împrăștiați; din această cauză pentru ca ei să funcționeze în mod eficient în acest spațiu de reprezentare, este nevoie de aplicarea unor tehnici de reducere a dimensiunii. Metodele kernel oferă o metodă elegantă de a depăși dimensionalitatea: au nevoie doar de matricea kernel care conține similaritățile de date.

Propunem o metodă de construcție a kernelurilor dependente de date care are mari șanse să producă rezultate mai potrivite în cazul problemelor de clasificare a documentelor. Pentru a construi un kernel mai potrivit pentru similaritățile documentelor, utilizăm informațiile textuale valoroase furnizate de către Wikipedia.

2.1 Reprezentarea documentelor bazată pe Wikipedia

În VSM cu t_{idf} documentele de ponderare sunt reprezentate prin frecvența termenilor înmulțite cu factorul idf . Matricea termen \times document \mathbf{D} este definit

în modul următor:

$$\mathbf{D} = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1n_D} \\ \vdots & \vdots & \vdots & \vdots \\ w_{n_T1} & w_{n_T2} & \cdots & w_{n_Tn_D} \end{bmatrix}$$

În acest fel avem o dublă reprezentare: (i) reprezentarea documentelor (coloanele), și (ii) reprezentarea termenilor (rândurile). Fiecare document este reprezentat de termenii care apar în document, iar fiecare termen este reprezentat de documentele în care acel termen apare.

Acum trecem la o altă reprezentare a documentelor. Prima dată dăm o nouă reprezentare a termenilor, adică reprezentăm fiecare termen ca și distribuția acelui termen în câteva documente (diferite). Ceea ce privește reprezentarea unui document – a unui set de termeni – în acest spațiu al noului document formăm suma ponderată a vectorilor termenilor. Considerăm următorul exemplu: avem trei termeni de indexare iar documentul nostru arată în următorul fel: $[1 \ 0 \ 1]'$. Alegem doar două alte documente din corpus care arată în următorul fel:

$$\begin{bmatrix} 2 & 1 & 0 \\ 1 & 3 & 2 \end{bmatrix}$$

unde așezăm documentele în rândurile matricei. Acum termenii primesc următoarele reprezentări:

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \cdot [1] = \begin{bmatrix} 2 \\ 1 \end{bmatrix}; \quad \begin{bmatrix} 1 \\ 3 \end{bmatrix} \cdot [0] = \begin{bmatrix} 0 \\ 0 \end{bmatrix}; \quad \begin{bmatrix} 0 \\ 2 \end{bmatrix} \cdot [1] = \begin{bmatrix} 0 \\ 2 \end{bmatrix}$$

în cazul cărora vectorul documentului va fi $1 \cdot [2 \ 1]' + 0 \cdot [0 \ 0]' + 1 \cdot [0 \ 2]' = [2 \ 3]'$. Se poate observa ușor faptul că acesta este chiar kernelul GVSM [Wong et al., 1985], așadar transformăm un document \mathbf{d} folosind $\mathbf{D}'\mathbf{d}$, cu condiția că documentele din care extragem distribuția termenului formează același corpus din care provin documentele.

Gabrilovich and Markovitch [2007] au modificat transformarea GVSM în următorul fel: în locul utilizării aceluiași corpus pentru detectarea corelărilor termenilor ei au folosit Wikipedia pentru extragerea distribuției termenilor. Wikipedia este compusă din aproximativ 2.5×10^7 articole¹ și astfel va oferi o mai bună, o mai bogată reprezentare. În cele de mai jos vom utiliza termenii “articol” și “concept” alternativ.

Să presupunem că Wikipedia conține n_C articole, sau am ales atâția din setul complet, și suntem interesați de distribuția în aceste articole Wikipedia ai termenilor de indexare. Transformarea documentului în spațiul conceptual al Wikipedia

¹August, 2008

devine următorul:

$$\underbrace{\begin{pmatrix} c_{11} & c_{12} & \dots & c_{1n_T} \\ c_{21} & c_{22} & \dots & c_{2n_T} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n_C1} & c_{n_C2} & \dots & c_{n_Cn_T} \end{pmatrix}}_{\mathbf{W}} \cdot \underbrace{\begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_{n_T} \end{pmatrix}}_{\mathbf{d}}$$

Denumirea acestei matrice este matricea concept \times termen Wikipedia, denumirea scurtă fiind matricea Wikipedia.

Gabrilovich and Markovitch [2007] numesc această metodă Explicit Semantic Analysis (ESA) pentru că spre deosebire de Latent Semantic Analysis (LSA) [Deerwester et al., 1990], ei aplică acești termeni/documente cu concepte explicite și nu unele latente. Putem interpreta această transformare – și în mod similar transformarea GVSM – ca transformarea vectorului documentului într-un vector al similarităților între document și conceptele Wikipedia. Astfel o comparație în această nouă reprezentare compară acești vectori în mod similar, măsurând proximitatea acestor similarități. În acest fel noul kernel de document devine

$$k(\mathbf{d}_i, \mathbf{d}_j) = \mathbf{d}_i' \mathbf{W}' \mathbf{W} \mathbf{d}_j$$

unde $\mathbf{W}'\mathbf{W}$ este matricea termen \times termen de co-apariție a termenilor din Wikipedia.

Gabrilovich and Markovitch [2007] au folosit similaritatea cosinus pentru a compara cuvinte și texte, și au testat noua reprezentare pe colecția WordSimilarity-353² și pe o colecție de 50 de documente din serviciile de știri a *Australian Broadcasting Corporation* [Lee et al., 2005]. Valorile de similaritate care au rezultat din utilizarea noii reprezentări cu similaritatea cosinus au fost comparate cu judecata umană, calculând corelațiile dintre ființele umane și calculator. Au obținut coeficienți de corelație de 0.75 și 0.72 pentru cuvinte respectiv texte, care sunt valorile cele mai ridicate produsă de un astfel de sistem automatizat. Pentru performanțele celorlalti algoritmi a se vedea studiul de referință.

Inspirat de performanța reprezentării documentelor în spațiul de concepte Wikipedia, am decis să utilizăm acest kernel pentru clasificarea textelor. Am încercat să reducem dimensionalitatea și să filtrăm zgomotul din această reprezentare.

Această metodă de transformare a documentelor poate fi considerată o tehnică semisupervizată unde un corpus larg neetichetat al articolelor Wikipedia este utilizat pentru a conferi documentelor o nouă reprezentare.

²<http://www.cs.technion.ac.il/~gabr/resources/data/wordsim353>

2.2 Reducerea dimensionalității kernelului Wikipedia

În timpul construirii kernelului Wikipedia nu am utilizat toate articolele, pentru că multe dintre ele sunt prea scurte pentru a fi folosite și sunt multe articole irelevante, de exemplu paginile colectoare sau de categorizare (e.g. `Category:Machine_learning`), pagini interne Wikipedia (e.g. `Wikipedia:Statistics`) etc. Din această cauză am selectat o fracțiune a articolelor de pe Wikipedia pe care le-am considerat utilizabile, metodologia completă fiind descrisă în teză și în Capitolul 5.

Pentru filtrarea termenilor irelevanți am decis să utilizăm LSA, adică aproximarea matricei Wikipedia cu o matrice de rang inferior,

$$\mathbf{W} \approx \widehat{\mathbf{W}} = \mathbf{U}\mathbf{S}_k\mathbf{V}'$$

unde $\mathbf{U}\mathbf{S}\mathbf{V}'$ reprezintă descompunerea valorilor singulare (*Singular Value Decomposition* – SVD) a matricei \mathbf{W} . Acesta înseamnă că noul kernel Wikipedia devine

$$k(\mathbf{d}_i, \mathbf{d}_j) = \mathbf{d}_i' \widehat{\mathbf{W}}' \widehat{\mathbf{W}} \mathbf{d}_j = \mathbf{d}_i' \mathbf{V}_k \mathbf{S}_k^2 \mathbf{V}_k' \mathbf{d}_j$$

Am experimentat cu kernelul de mai sus, dar am observat că prin folosirea \mathbf{S}_k unele dimensiuni au primit o foarte mare influență, care la rândul ei a dus la reducerea performanței. Astfel am înlocuit \mathbf{S}_k cu \mathbf{I}_k și am obținut kernelul următor:

$$k(\mathbf{d}_i, \mathbf{d}_j) = \mathbf{d}_i' \mathbf{V}_k \mathbf{V}_k' \mathbf{d}_j$$

O altă interpretare a transformării de mai sus ar putea fi următorul: cum primul component al descompunerii SVD a matricei \mathbf{D} conține vectorii proprii sau componentele principale ale spațiului feature în LSA, în acest caz \mathbf{V} îl va conține (din cauza că \mathbf{D} este o matrice termen×document, iar \mathbf{W} este o matrice concept×termen), și am asumat că articolele Wikipedia produc o covariație mai bună, de aceea am înlocuit $\mathbf{D}\mathbf{D}'$ cu $\mathbf{W}'\mathbf{W}$.

O problemă cu decompoziția de mai sus este că \mathbf{W} devine destul de mare – mai ales că n_C este mare – astfel SVD-ul matricei este foarte ineficient, și chiar dificil de efectuat în timp acceptabil. Totuși putem observa faptul că prin decompoziția pe baza valorilor proprii a matricei mult mai mici $\mathbf{W}'\mathbf{W}$ (de mărimea $n_T \times n_T$) putem obține \mathbf{V} ,

$$\mathbf{W}'\mathbf{W} = \mathbf{V}\mathbf{S}^2\mathbf{V}'$$

2.3 Structura pe linkuri a Wikipedia

Wikipedia este de asemenea structurată pe linkuri, ceea ce poate fi utilizată în mod eficient prin propagarea frecvențelor prin aceste conexiuni. Considerăm matricea de legătură concept×concept \mathbf{E} , care este definită ca $E_{ij} = 1$ în cazul în care există

un link între conceptele i și j , iar altfel ia valoarea 0. Pentru că dorim să păstrăm ponderile deja atribuite, setăm diagonala principală cu valorile 1. Astfel matricea actualizată va deveni

$$\widetilde{\mathbf{W}} = \mathbf{E}'\mathbf{W}$$

Acesta înseamnă că $\widetilde{W}_{ij} = (\mathbf{E}_{\cdot i})' \cdot \mathbf{W}_{\cdot j}$, adică adăugăm W_{ij} -ului totalul aparițiilor termenului j în conceptele $C = \{k_1, \dots, k_n\}$, unde setul C conține conceptele – sau indicile acelor concepte – din direcția cărora s-a produs un link către conceptul i .

2.4 Ponderarea conceptelor

În timp ce cuvintele/dimensiunile în corpusul de instruire primesc importanțe variate prin utilizarea ponderii `tfidf`, conceptele sunt considerate de către Wikipedia de a avea importanțe egale. O schemă posibilă de ponderare poate fi obținută prin clasificarea articolelor Wikipedia pe baza analizei citatelor sau a referințelor, astfel aceste importanțe pot fi extrase din structura pe linkuri discutată în capitolul anterior. Famosul algoritm de PageRank [Page et al., 1998] a Google obține acest lucru doar prin luarea în considerare a structurii de linkuri sau de citate a unui set de pagini hyperlinkate: o pagină are un rang înalt dacă deține multe back-linkuri sau dacă sunt doar câteva back-linkuri, dar acele sunt importante. Acest lucru poate fi formulat în mod recursiv prin

$$\mathbf{r}_i = \sum_{j \in N^{-1}(i)} \frac{\mathbf{r}_j}{|N^{-1}(j)|}$$

unde \mathbf{r}_i notează rangul paginii al i -lea și $N^{-1}(i)$ este setul de vecini backward a i , adică vecinii care sunt îndreptate spre i . Dacă \mathbf{P} notează matricea de adiacență având $P_{ij} = 1/|N(i)|$, unde $N(i)$ reprezintă vecinii forward ai i , vecini în direcția cărora se îndreaptă i , rangurile pot fi calculate prin rezolvarea problemei de valori proprii $\mathbf{r} = \mathbf{P}'\mathbf{r}$.

PageRank poate fi considerat un mers aleator pe un graf, unde rangurile sunt probabilitatea ca un navigator aleatoriu să fie pe pagina respectivă în pasul k . În cazul în care k este suficient de mare, probabilitățile converg într-o distribuție unică fixă. Singura problemă cu formularea de mai sus este că graful poate avea noduri fără vecini forward, și grupuri din care nu iese nici un link forward. Pentru a rezolva aceste probleme este necesar să se adauge incertitudine în procesul de navigare

$$\mathbf{r} = c\mathbf{P}'\mathbf{r} + (1 - c)\mathbf{1}, \quad c \in [0, 1]$$

unde prin $\mathbf{1}$ am notat vectorul $[1 \ 1 \ \dots \ 1]'$ de mărimea n_p având n_p pagini. În experimentele noastre am folosit valoarea $c = 0.85$ conform [Page et al., 1998].

Prin utilizarea rangurilor produse de algoritmul PageRank înlocuim \mathbf{W} cu $\text{diag}(\tilde{\mathbf{r}})\mathbf{W}$, unde $\tilde{\mathbf{r}} = \log(\mathbf{r} + \mathbf{1})$. Utilizând funcția \log am încercat să rafinăm ponderile într-o oarecare măsură, pentru că proporția conceptului de cel mai înalt și cel mai scăzut rang a fost de 8.8×10^3 . Am utilizat translația $\mathbf{r} + \mathbf{1}$ în loc de \mathbf{r} pentru a evita ponderile negative.

3 Kerneluri cluster ierarhici

În acest capitol introducem kerneluri cluster ierarhici în învățarea semisupervizată. Propunem utilizarea distanțelor induse de algoritmi de clustering în locul distanțelor euclidiene în spațiul input. Dacă datele neetichetate sunt adăugate unui set relativ mic de date etichetate, presupunem că noua distanță, obținută prin clustering și utilizarea datelor neetichetate să inducă un spațiu de reprezentare mai potrivit pentru clasificare. Pentru clustering utilizăm tehnici speciale de clustering ierarhic – acelea care rezultă în matrice de distanță ultrametrică – care are ca rezultat matrice kernel pozitiv semidefinite.

Metoda propusă se bazează pe kernelul de conectivitate [Fischer et al., 2003] și extindem asupra acestei construcții kernel prin admiterea oricărei metode de clustering ierarhic dacă acesta are ca rezultat o matrice de distanță ultrametrică. Procesul de clustering necesită setul complet de date, și doar subsetul lui etichetat este utilizat pentru construirea kernelului în problema de clasificare. Totuși, dacă datele de testare sunt disponibile în timpul instruirii, putem să le includem în setul de instruire ca date neetichetate adițional, și apoi utilizând partea corespunzătoare a matricei de distanță. În cazul seturilor de date unde prezumția de varietate (*manifold assumption*) este așteptată a ține, construim kernelul cluster ierarhic utilizând distanțele determinate pe grafurile de date kNN (*k-nearest neighbors*) și ϵ NN (*ϵ -nearest neighbors*) [von Luxburg, 2006].

3.1 Clustering ierarhic

Clustering este împărțirea unui set în grupe separate numite clustere. Cu toate că clusterelor sunt de multe ori dezmembrate, ele pot fi în același timp și suprapuse, ceea ce înseamnă că un punct poate aparține mai multor clustere. O metodă de clustering poate fi ierarhic sau partiționat. Clusteringul ierarhic construiește un arbore în pași succesivi, unde nodurile arborelui reprezintă partiționările suprapuse ale datelor. Prin partiționare înțelegem o decompoziție a setului de date în mai multe grupe. Clusteringul partiționat rezultă într-o singură partiție [Jain et al., 1999; Berkhin, 2002]. În cele de mai jos vor fi utilizate în primul rând metode de clustering ierarhic.

Avem de-a face cu un clustering suprapus în cazul în care o partiție este formată din componente fuzionate. Clusteringul ierarhic este compus dintr-o secvență de partiții cu fiecare partiție suprapusă cu următoarea partiție a secvenței. În cazul metodelor aglomerative la început fiecare punct este un cluster, iar clusterelor cele mai apropiate sunt îmbinate până la obținerea unui singur cluster mare. În cazul metodelor divizive la început tot setul de date este un singur cluster, apoi clusterelor sunt despărțite până când sunt obținute clusterelor formate de câte un punct [Jain and Dubes, 1988; Duda et al., 2001]. În cazul kernelului cluster ierarhic utilizăm o formă specială a metodei de clustering aglomerativ. Un algoritm general de clustering, care produce un dendrogram, este următorul:

Algoritm 1 Clustering aglomerativ

- 1: Se definesc clusterelor inițiale ca și punctele în sine.
 - 2: Se găsește perechea de clusterelor cele mai similare.
 - 3: Se îmbină acestea pentru a crea un nou cluster.
 - 4: Se repetă de la pasul 2 până când se obține un singur cluster.
-

Pentru a specifica algoritmul complet, este nevoie de măsurarea similarităților între clusterelor; acestea sunt denumite *distanțe de legătură*. Pe baza alegerii distanței de legătură devine posibilă designul unei mari varietăți de metode de clustering aglomerativ, prezentate în detaliu mai jos.

3.2 Distanțele de legătură

Distanțele de legătură (*linkage distances*), $D(C_1, C_2)$, măsoară distanțele între clusterelor în clusteringul aglomerativ. Aici prezentăm trei distanțe de legătură populare. Acestea sunt bazate pe $d(\mathbf{x}_1, \mathbf{x}_2)$, *distanța pe puncte* (*pointwise distance*) în spațiul input, care este de obicei distanța euclidiană, $d(\mathbf{x}_1, \mathbf{x}_2) = \|\mathbf{x}_1 - \mathbf{x}_2\|_2$.

(1) Clusteringul *single linkage* folosește următoarea funcție de distanțe cluster:

$$D(C_1, C_2) = \min \{d(\mathbf{x}_1, \mathbf{x}_2) \mid \mathbf{x}_1 \in C_1, \mathbf{x}_2 \in C_2\} \quad (1)$$

Adică alegem să fuzionăm două clusterelor unde distanța pe puncte este minimală. Se poate observa că clusteringul *single linkage* este echivalentul găririi arborelui de acoperire al grafului de date, adică dacă considerăm graful inițial ca și graful complet al punctelor de date, atunci prin alegerea repetată a muchiei cu ponderea minimală obținem arborele de acoperire minim al grafului [Duda et al., 2001].

(2) Clusteringul *complete linkage* definește distanța între clusterelor în felul următor:

$$D(C_1, C_2) = \max \{d(\mathbf{x}_1, \mathbf{x}_2) \mid \mathbf{x}_1 \in C_1, \mathbf{x}_2 \in C_2\} \quad (2)$$

și poate fi considerat ca unul care lucrează cu grafurile complete în cadrul clusterelor, adică în fiecare cluster fiecare nod este conectat cu toate celelalte noduri.

Dacă definim diametrul unui cluster fiind muchia cea mai lungă între două puncte, iar diametrul unei partiții este definit ca fiind cel mai mare diametru al clusterelor sale, în complete linkage, din formula (2), alegem să îmbinăm acele cluster care măresc cel mai puțin diametrul partiției [Duda et al., 2001].

Metodele de mai sus tind să fie sensibile punctelor outlier. De altfel, clusteringul single linkage poate lega în mod eronat cluster și poate forma cluster de o omogenitate redusă, iar clusteringul complete linkage poate rezulta în cluster care nu sunt bine separate [Jain and Dubes, 1988].

(3) Clusteringul *average linkage* – reprezentând un compromis – ia în considerare media distanțelor pe puncte dintre toate perechile de elemente din cele două cluster:

$$D(C_1, C_2) = \frac{1}{|C_1||C_2|} \sum_{\mathbf{x}_{1i} \in C_1} \sum_{\mathbf{x}_{2j} \in C_2} d(\mathbf{x}_{1i}, \mathbf{x}_{2j}) \quad (3)$$

În această teză experimentăm doar cu cele trei distanțe de legătură prezentate mai sus. Alte tehnici populare, bazate pe variația clusteringului *average linkage* (denumit *UPGMA – unweighted pair group method using arithmetic mean*) include *weighted average linkage clustering* (WPGMA), *average group linkage clustering* (UPGMC), *weighted average group linkage clustering* (WPGMC) și metoda lui Ward [Jain et al., 1999; Berkhin, 2002].

Toate cele trei metode au următoarele calități utilizate pentru construcția kernelurilor pozitiv semidefinite pentru clusteringul aglomerativ: să presupunem că alegem să îmbinăm trei cluster, C_1 , C_2 and C_3 în următoarea ordine: prima dată îmbinăm C_1 cu C_2 ceea ce rezultă în C_{12} , apoi vom îmbina acesta din urmă cu C_3 . Dacă

$$D(C_1, C_2) \leq D(C_1, C_3) \quad \text{și} \quad D(C_1, C_2) \leq D(C_2, C_3)$$

atunci

$$D(C_1, C_2) \leq D(C_{12}, C_3)$$

Această calitate este denumită *proprietatea ultrametrică* sau ultrametricitate. Pe baza unei metode de clustering aglomerativ care utilizează distanțe de legătură ultrametrice, putem defini o matrice de distanțe ultrametrice, și pe rând o funcție kernel care poate fi utilizată pentru o reprezentare mai potrivită.

3.3 Construirea kernelului

Așa cum am afirmat mai sus, clusteringul ierarhic rezultă într-un dendrogram, nodurile căruia sunt etichetate cu distanța între clusterelor care sunt îmbinate la nodul respectiv. Construim o matrice de distanțe prin considerarea etichetei atașate strămoșului comun cel mai apropiat al punctelor de pe arbore. Pentru a transforma

distanțele în produse scalare utilizăm o metodă similară scalării multidimensionale (*Multi-Dimensional Scaling – MDS*) [Borg and Groenen, 2005]:

$$\mathbf{K} = -\frac{1}{2}\mathbf{J}\mathbf{M}\mathbf{J} \quad \text{cu} \quad \mathbf{J} = \mathbf{I} - \frac{1}{N}\mathbf{1}\mathbf{1}'$$

unde \mathbf{M} conține distanțele bazate pe dendrogramă la pătrat, și \mathbf{J} este matricea de centrare construită din matricea de identitate și produsul tensorial al vectorului cu toate elementele egale cu 1. Matricea care rezultă din aceasta conține produsele scalare între un set de vectori $\{\mathbf{z}_i\}_{i=1}^N$ cu distanțele euclidiene la pătrat $\|\mathbf{z}_i - \mathbf{z}_j\|_2^2 = M_{ij}$. Utilizăm următoarea teoremă care declară că prin această transformare matricele de distanță ultrametrică rezultă întotdeauna în matrice Gram.

Theorem 1 (Fischer et al. [2003]). *Dat fiind un \mathbf{M} ultrametric care conține distanțele bazate pe dendrogramă, matricea $\mathbf{K} = -\frac{1}{2}\mathbf{J}\mathbf{M}\mathbf{J}$ este o matrice Gram pozitiv semidefinită.*

În cele de mai jos construim kernelul cluster cu utilizarea distanțelor de legătură din clusteringul ierarhic. Astfel putem aplica punctele într-un spațiu feature unde distanțele pe puncte sunt egale cu distanțele cluster în spațiul input. Pașii acestuia sunt:

Algoritm 2 Kernelul cluster ierarhic

- 1: Efectuăm un clustering aglomerativ pe datele etichetate și neetichetate – de exemplu prin utilizarea funcțiilor de legături single, complete sau average.
 - 2: Definim matricea \mathbf{M} cu elemente $M_{ij} =$ distanța de legături în arborele ultrametric la subsumarea cel mai puțin generală a i și j ; $M_{ii} = 0, \forall i$.
 - 3: Definim matricea kernel ca și $\mathbf{K} = -\frac{1}{2}\mathbf{J}\mathbf{M}\mathbf{J}$.
-

Kernelul care rezultă în acest fel este în mod evident unul dependent de date. Folosim datele neetichetate în pasul de clustering pentru a determina distanțe pe puncte “mai bune” care duc la kernel. Așteptarea noastră este obținerea unei mai bune similarități prin includerea părții neetichetate. Este important faptul că pentru a calcula funcția kernel pentru setul de testare, le includem în setul neetichetat. Acesta înseamnă că dacă punctul de testare nu este disponibil la momentul instruirii, întregul proces de clustering ar trebui repetat prin încetinirea clasificării. În această teză prezentăm o metodă care poate fi utilizată pentru a evita recalcularea kernelului.

3.4 Kerneluri cluster ierarhici cu distanțe pe graf

În construirea kernelului cluster ierarhic am folosit doar prezumția cluster. Aici extindem kernelul de mai sus pentru a exploata prezumția de varietate utilizând

un kernel cluster ierarhic bazat pe grafuri. Aproximăm distanțele prin utilizarea drumurilor cele mai scurte care sunt bazate pe grafurile kNN sau ϵ NN, similare cu ISOMAP [Tenenbaum et al., 2000]. În acest proces substituim distanțele pe puncte $d(\cdot, \cdot)$ cu distanțele pe graf.

Rezultatul este că algoritmul de clustering ierarhic este precedat de următorii trei pași:

Algoritmul 3 Kernel cluster ierarhic bazat pe graf

- 2: Determinăm cei mai apropiați k vecini sau ϵ -vecinătatea fiecărui punct și a se considera distanța celorlalte puncte a fi egală cu zero.
 - 1: Calculăm drumurile cele mai scurte pentru fiecare pereche de puncte – utilizând de exemplu algoritmul lui Dijkstra.
 - 0: Utilizăm aceste distanțe în clusteringul pentru distanțele pe puncte $d(\cdot, \cdot)$.
-

Am început numărarea în mod deliberat de la -2 pentru a accentua faptul că acești pași *preced* algoritmul din subcapitolul anterior. Subliniem faptul că acești pași sunt opționali: ei trebuie folosiți în cazul în care prezumția de varietate se ține în setul de date.

Utilizăm aici cele mai scurte drumuri calculate folosind graful de cei mai apropiați k vecini sau de ϵ -vecinătate a datelor, astfel – dacă datele se află pe o varietate de dimensiuni scăzute – aproximând distanțele pe această varietate [vezi Bernstein et al., 2000, pentru condiții].

Graful construit folosind cei mai apropiați k vecini sau ϵ -vecinătatea punctelor poate conține mai multe componente deconectate, de exemplu dacă k sau ϵ este prea mic. În [Tenenbaum et al., 2000] se afirmă că clusterelor mici deconectate de componenta “gigantică” conțin de obicei outlieri, și din această cauză pot fi neglijate. Aici lăsăm detectarea de outlieri ca și un pas de preprocesare, însemnând faptul că nu neglijăm intrări și dorim să obținem un singur component conectat. Urmăm [Yong and Jie, 2005] pentru a obține un component: prin \mathbf{M} notăm matricea de distanțe pe puncte după rărirea matricei luând cei mai apropiați k vecini sau ϵ -vecinătatea a punctelor. În mod similar notăm prin \mathbf{G} matricea de drumuri cele mai scurte bazată pe \mathbf{M} rărit. Dacă vom găsi valori $G_{ij} = \infty$, înseamnă că graful nu este conectat în totalitate. Astfel alegem perechea de \mathbf{x}_i și \mathbf{x}_j neconectat având distanța minimă euclidiană, și le conectăm prin utilizarea unei distanțe relativ mari,

$$\hat{G}_{ij} = g_{\max} + \frac{d_{\min}}{d_{\max}}$$

unde g_{\max} este drumul maxim în \mathbf{G} , d_{\min} și d_{\max} sunt distanțele minime și maxime – fie euclidiene, fie bazate pe graf – între punctele de date. După schimbarea \hat{G}_{ij} , actualizăm fiecare element în \mathbf{G} pentru care $G_{kl} = \infty$. Pentru orice astfel de

pereche k și ℓ actualizarea este:

$$G_{k\ell} = \min\{G_{ik} + \widehat{G}_{ij} + G_{j\ell}, G_{kj} + \widehat{G}_{ij} + G_{i\ell}\}$$

Dacă totuși rămân componente neconectate, procedura de mai sus trebuie repetată până la obținerea unui singur component conectat.

4 Kerneluri de reponderare

În continuare propunem trei tehnici ușor diferiți de reponderare a unui kernel de bază utilizând prezumția cluster a învățării semisupervizate. Înaintea acestuia enumerăm combinațiile kernel care au ca rezultat kerneluri pozitiv semidefinite și sunt utilizate în cele de mai jos [Lütkepohl, 1996; Abadir and Magnus, 2005]:

- (i) Dacă \mathbf{K}_1 și \mathbf{K}_2 sunt matrice pozitiv semidefinite, atunci $\mathbf{K}_1 + \mathbf{K}_2$ este de asemenea pozitiv semidefinită.
- (ii) Dacă \mathbf{K} este o matrice pozitiv semidefinită, iar $\alpha > 0$, atunci $\alpha \mathbf{K}$ este de asemenea pozitiv semidefinită.
- (iii) Dacă \mathbf{K}_1 și \mathbf{K}_2 sunt matrice pozitiv semidefinite, atunci $\mathbf{K}_1 \odot \mathbf{K}_2$ este de asemenea pozitiv semidefinită, unde \odot notează produsul Hadamard.

În acest capitol creăm kerneluri dependente de date prin reponderare: utilizăm prezumția cluster și mărirea similarității pentru punctele din același cluster și scăderea acesteia când clusterelor punctelor sunt diferite [Weston et al., 2006].

Dezvoltăm tehnici care reponderesc matricea kernel prin exploatarea structurii cluster a datelor de instruire. Astfel în cazul în care două puncte sunt în același cluster, similaritatea lor primește o pondere ridicată, ≈ 1 sau > 1 , iar dacă sunt în diferite cluster, similaritatea lor primește o pondere mai scăzută, adică de < 1 sau $\ll 1$, ceea ce noi numim *kernel de reponderare*, sau $k_{rw}(\mathbf{x}_1, \mathbf{x}_2)$. Ponderile de similaritate sunt combinate cu valorile kernelului original sau de bază $k_b(\mathbf{x}_1, \mathbf{x}_2)$, și vor forma o altă matrice kernel pentru a garanta natura semidefinită a kernelului final. Rezumând, noul kernel cluster este

$$k(\mathbf{x}_1, \mathbf{x}_2) = k_{rw}(\mathbf{x}_1, \mathbf{x}_2) k_b(\mathbf{x}_1, \mathbf{x}_2)$$

unde $k_{rw}(\cdot, \cdot)$ este reponderarea și $k_b(\cdot, \cdot)$ este kernelul de bază. Folosind notația matriceală acesta poate fi scris astfel

$$\mathbf{K} = \mathbf{K}_{rw} \odot \mathbf{K}_b$$

Ne confruntăm cu două probleme în construcția kernelului cluster de mai sus: (i) kernelul de reponderare trebuie să fie pozitiv semidefinit, (ii) matricea kernel de

bază trebuie să fie pozitiv semidefinit și *pozitiv*. Prima cerință este evidentă: este nevoie de ea pentru a garanta natura pozitiv semidefinită a kernelului construit.

A doua condiție este crucială mai ales pentru că în cazul valorilor negative în matricea kernel de bază este nevoie de efectuarea unei reponderări diferite. Pentru a evita complicațiile datorate negativității, solicităm un kernel de bază pozitiv, $k_b(\mathbf{x}_1, \mathbf{x}_2) \geq 0$. Utilizând produse scalare aceasta poate fi îndeplinită prin deplasarea fiecărei direcții a datelor în semispațiul pozitiv, sau pot fi folosite doar kerneluri pozitive ca și cel gaussian sau kernelurile polinomiale de gradul doi [Schölkopf and Smola, 2002].

În cele de mai jos prezentăm kernelul bagged cluster, apoi propunem trei metode pentru construirea kernelurilor de reponderare pozitiv semidefinite.

4.1 Kernelul bagged cluster

Kernelul bagged cluster, propus în [Weston et al., 2006], reponderește valorile kernelului de bază ținând cont de probabilitatea că punctele aparțin aceluiași cluster. Pentru a calcula această probabilitate, kernelul bagged cluster utilizează clusteringul k-means [Jain and Dubes, 1988], împreună cu proprietatea că alegerea centrelor clusterului inițial afectează ieșirea algoritmului. Presupunând că avem N itemi de date și K cluster, kernelul este construit așa cum apare în Algoritmul 4 [Weston et al., 2006].

Algoritmul 4 Kernelul bagged cluster

- 1: Lansăm k-means de t ori, ceea ce rezultă în atribuire cluster $c_j(\mathbf{x}_i)$, $j = 1, \dots, t$, $i = 1, \dots, N$, $c_j(\cdot) \in \{1, \dots, K\}$.
- 2: Construim kernelul bagged în următorul fel:

$$k_{\text{bag}}(\mathbf{x}, \mathbf{z}) = \frac{\sum_{j=1}^t [c_j(\mathbf{x}) = c_j(\mathbf{z})]}{t}$$

Următorul pas este înmulțirea kernelului de bază cu kernelul bagged obținut în pasul precedent, adică folosind notația matriceală:

$$\mathbf{K}_{\text{bag}} = \frac{1}{T} \sum_{t=1}^T \mathbf{u}^{(t)'} \mathbf{u}^{(t)} \quad \text{și} \quad \mathbf{K} = \mathbf{K}_{\text{bag}} \odot \mathbf{K}_b$$

cu matricea de apartenență cluster \mathbf{U} de dimensiune $K \times N$, unde prin $\mathbf{u}^{(t)}$ am notat matricea obținută în pasul t .

4.2 Kernelul de reponderare gaussian

Să supunem că sunt date ieșirile unui algoritm de clustering, matricea de apartenență cluster \mathbf{U} de tipul $K \times N$. Distingem cele două familii de metode clusterig: metodele clustering partiționale și fuzzy. Metodele partiționale produc o matrice de apartenență cluster în care fiecare coloană conține o singură valoare de 1, toate celelalte fiind nule; un algoritm fuzzy de clustering produce o matrice, unde fiecare coloană este o distribuție probabilistică asupra clusterilor, adică U_{ij} conține probabilitatea că punctul al j -ilea aparține clusterului i . În ambele cazuri presupunem că cele două puncte aparțin aceluiași cluster(e) dacă vectorii de apartenență cluster sunt similare sau apropiate unul de celălalt. Definim similaritatea prin kernelul gaussian în felul următor:

$$k_{rw}(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\|\mathbf{u}_{\mathbf{x}_1} - \mathbf{u}_{\mathbf{x}_2}\|^2}{2\sigma^2}\right) \quad (4)$$

unde $\mathbf{u}_{\mathbf{x}}$ notează vectorul de apartenență cluster al punctului \mathbf{x} , adică coloana lui \mathbf{x} în \mathbf{U} . Știm că matricea construită este pozitiv semidefinită [Schölkopf and Smola, 2002] cu fiecare element fiind între 0 și 1. În acest caz parametrul σ definește cantitatea separării între punctele similare și disimilare: dacă σ este mare, gaura dintre valorile exprimând similaritate și disimilaritate devine mai mic, iar pentru un σ mai mic, aceste valori se deplasează mai departe una de cealaltă.

4.3 Kerneluri de reponderare bazate pe produse scalare

O altă posibilitate de utilizare a vectorilor de apartenență cluster este definirea următorului kernel de reponderare:

$$\mathbf{K}_{rw} = \mathbf{U}'\mathbf{U} + \alpha \mathbf{1}\mathbf{1}' \quad (5)$$

unde \mathbf{U} notează matricea de apartenență cluster și $\alpha \in [0, 1)$. (Pentru o formă mai generală s-ar putea folosi kernelul polinomial $k(\mathbf{x}_1, \mathbf{x}_2) = (\alpha + \mathbf{x}_1'\mathbf{x}_2)^p$ cu α pozitiv și p pozitiv natural.) Primul termen $\mathbf{U}'\mathbf{U}$ punctează similaritatea pe baza apartenenței cluster, iar al doilea termen este utilizat pentru a evita similaritățile zero: dacă doi vectori de apartenență sunt ortogonale, conducând la valoarea de zero în matricea produselor scalare. Dacă presupunem că clusteringul primit nu este foarte “de încredere”, putem folosi o valoare mică α , astfel apartenența la un cluster este diminuat cu ajutorul termenului $\alpha \mathbf{1}\mathbf{1}'$.

Este clar faptul că kernelul de reponderare din formula (5) este pozitiv semidefinit: primul și al doilea termen este un produs tensorial, astfel pozitiv semidefinit, și în baza proprietăților definite în Capitolul 4 putem deduce faptul că kernelul este într-adevăr unul pozitiv semidefinit.

	#eigv	mP	mR	mBEP	mF1	MP	MR	MBEP	MF1
χ^2_{5209}	–	88.46	84.59	86.52	86.48	71.61	61.25	66.43	66.02
χ^2_{4601}	–	87.88	83.09	85.49	85.42	64.21	54.62	59.41	59.03
covarianță Wikipedia	–	48.95	35.42	42.18	41.10	8.79	5.35	7.07	6.65
LSA	4500	88.05	83.65	85.85	85.80	62.48	54.38	58.43	58.15
LSA+linkuri	4500	87.89	83.71	85.80	85.75	65.11	56.05	60.58	60.24
LSA+PageRank	4000	87.44	83.68	85.56	85.52	59.75	53.86	56.80	56.65

Tabela 1: Performanța metodei pe corpusul Reuters exprimat în procente. Notății: mP=micro-precision, mR=micro-recall, mBEP=micro-breakeven, mF1=micro-F₁, MP=macro-precision, MR=macro-recall, MBEP=macro-breakeven, MF1=macro-F₁

O altă versiune a kernelului din (5) este

$$\mathbf{K}_{rw} = \beta \mathbf{U}'\mathbf{U} + \mathbf{1}\mathbf{1}' \quad (6)$$

unde $\beta \in (0, \infty)$. Aici valorile kernel pentru care matricea de produs scalar a vectorilor de apartenență cluster corespund cu zero, prin $\mathbf{1}\mathbf{1}'$ rămân aceleași, totuși dacă punctele se află în același cluster, $\beta\mathbf{U}'\mathbf{U}$ conferă o pondere mai mare decât zero, astfel această valoare kernel va fi mărită.

Pentru clusteringul fuzzy, unde o coloană din \mathbf{U} conține probabilități de apartenență cluster, prima dată normalizăm coloanele pentru a obține valoarea 1 în cazul vectorilor identici de apartenență cluster.

Pentru a construi kernelul de reponderare punctele trebuie clusterate, adică trebuie implicat un algoritm de clustering. În caz optim dacă setăm numărul clusterelor la numărul claselor în problema de clasificare, vom obține o partiție care corespunde separării adecvate de clasă a datelor, dar din nefericire acest lucru se întâmplă foarte rar; astfel se poate experimenta cu utilizarea diferitelor numere de cluster.

5 Experimente

În această parte descriem pe scurt experimentele și rezultatele obținute.

În experimentele noastre de categorizare am utilizat depozitul Wikipedia (*dump*) în limba engleză din luna noiembrie 2006³ care conținea aproximativ 1.6 milioane de articole, care este egal cu aproximativ 8 gigabyte de date textuale – în afara imaginilor și altor documente adiționale, dar incluzând tagurile XML. În mod asemănător cu [Gabrilovich and Markovitch, 2007], pentru a filtra articolele irelevante (paginile de categorizare, pagini interne, cioturi Wikipedia [vezi `Wikipedia:Stub`], etc.), am eliminat următoarele articole:

- articole care conțin mai puțin de 500 de cuvinte

³<http://download.wikimedia.org/enwiki>

- articole care conțin mai puțin de 5 linkuri forward către alte articole Wikipedia

După această reducere am construit un index inversat pentru cuvintele care apar în Wikipedia, în afara cuvintelor stop și a celor mai frecvente 300 de cuvinte. Am exclus de asemenea cuvinte care apar în mai puțin de 10 articole. Excluderea acestor cuvinte a fost efectuată pentru a ajuta filtrarea termenilor irelevanți de indexare, fiindcă am presupus că un cuvânt care nu este important în Wikipedia este de asemenea irelevant în indexarea corpusului de clasificare a documentelor. Astfel au rămas 327 653 articole.

Pentru testarea kernelului Wikipedia asupra categorizării de documente am utilizat corpusul Reuters-21578, împărțirea ModApté cu 90 + 1 categorii, în cadrul căruia categoria “unknown” (*necunoscut*) nu a fost utilizată.

În același timp am utilizat o metodă de selectare a termenilor de filtrare pentru a selecta cei mai relevanți termeni din corpusul Reuters. În acest scop am utilizat selectarea de termeni χ^2 și am selectat 5209 termeni (am obținut un bun rezultat cu acest număr de trăsături în [Minier et al., 2006]). Pentru aceste cuvinte am construit vectori de cuvinte reprezentând distribuția acestor cuvinte în articolele Wikipedia. Toate cuvintele extrase din articolele Wikipedia și din corpusul Reuters au fost prelucrate utilizând algoritmul lui Porter [Baeza-Yates and Ribeiro-Neto, 1999].

După construirea vectorilor de documente am folosit SVM-uri pentru a învăța datele de instruire; în acest scop am utilizat implementarea LIBSVM [Chang and Lin, 2001].

Rezultatele sunt prezentate în Tabela 1. Pentru evaluarea sistemului am utilizat măsura *precision*, *recall*, punctul de echilibru (*breakeven*) *precision–recall* și F_1 . χ^2_{5209} și χ^2_{4601} prezintă rezultatele obținute utilizând metoda de selectare a termenilor χ^2 . Rândurile covarianței Wikipedia, LSA, LSA+linkuri și LSA+PageRank arată rezultatele obținute în urma utilizării kernelului Wikipedia, adică kernelul Wikipedia fără reducere, kernelul Wikipedia cu LSA, kernelul Wikipedia cu LSA și cu utilizarea matricei link, și la sfârșit kernelul Wikipedia cu LSA și ponderarea conceptelor cu ajutorul algoritmului PageRank. Coloana #eigv arată numărul vectorilor proprii alese.

Cele mai bune rezultate au fost obținute prin utilizarea selecției χ^2 a termenilor, iar pe locul doi se află utilizarea kernelului Wikipedia cu LSA.

Tabela 2 prezintă rezultatele obținute în urma utilizării kernelurilor dependente de date pe seturile de date USPS, Digit1, COIL2 și Text [Chapelle et al., 2006]. Am utilizat precizia (*accuracy*) ca și măsura de evaluare, iar rezultatele sunt date în procente. În cazul fiecărui set de date am indicat cele mai bune rezultate obținute. Acest lucru nu a fost posibil în tabela anterioară fiindcă acolo am avut două

	USPS		Digit1		COIL2		Text	
	10	100	10	100	10	100	10	100
linear	72.82	86.43	81.07	90.86	60.74	80.43	58.26	67.86
Gaussian	80.07	89.71	56.11	93.86	57.38	82.50	59.06	56.43
ISOMAP	85.10	86.71	94.43	97.43	62.62	80.64	59.80	72.43
vecinătate	76.31	94.14	87.11	94.21	64.43	84.43	51.68	62.79
bagged	87.38	92.79	93.29	96.93	71.28	85.57	63.29	66.14
multi-type, pas	80.07	92.86	91.01	91.29	55.77	84.86	53.56	74.79
multi-type, pas linear	80.07	92.86	91.01	91.36	55.77	84.86	53.02	75.29
multi-type, polinomial	80.07	80.29	48.86	65.07	54.23	82.29	50.60	56.71
HCK, single	80.07	81.79	48.86	70.21	67.85	96.00	66.78	73.14
HCK, complete	82.01	89.50	60.67	89.71	55.64	86.36	50.27	49.57
HCK, average	81.48	92.86	71.75	93.79	68.05	91.71	64.63	50.14
gHCK, single	80.07	81.79	48.86	70.21	60.60	93.86	66.78	73.14
gHCK, complete	88.26	95.64	75.50	93.71	68.52	88.79	56.17	67.71
gHCK, average	89.26	95.64	94.70	95.21	60.54	90.64	47.32	66.86
RCK1, k-means	84.45	92.98	83.14	94.28	58.55	83.76	–	–
RCK1, ierarhic	86.17	95.29	89.06	94.94	62.08	85.93	62.35	68.07
RCK1, spectral	81.43	90.87	88.32	95.20	58.22	83.83	63.26	66.93
RCK2, k-means	83.86	92.45	84.58	94.08	58.95	83.76	–	–
RCK2, ierarhic	86.11	95.50	89.06	95.29	62.08	85.64	61.28	71.14
RCK2, spectral	81.63	91.39	88.32	94.64	58.03	83.37	61.50	70.07
RCK3, k-means	83.66	92.59	84.13	92.96	58.60	83.28	–	–
RCK3, ierarhic	84.97	95.29	89.06	94.57	62.95	86.07	59.13	71.21
RCK3, spectral	81.16	91.56	88.32	94.73	55.83	83.20	59.26	71.00

Tabela 2: Rezultatele de precizie (*accuracy*) după utilizarea diferitelor kerneluri. Rezultatele sunt prezentate în procente. Pentru fiecare set de date cele mai bune rezultate sunt înramate.

linii de referință.

Primele două rânduri conțin rezultatele liniilor de referință obținute prin utilizarea kernelurilor lineare și gaussiene. În principal am dorit să îmbunătățim aceste rezultate. Hiperparametrul kernelului gaussian a fost fixat prin aplicarea a două metode – vezi [Chapelle et al., 2006] și [Zhu and Ghahramani, 2002] – pentru estimarea valorii de bază σ_0 , iar apoi am utilizat o validare încrucișată de 10 ori pe setul mai mare etichetat conținând 100 de puncte etichetate utilizând valorile $[\sigma_0/8 \ \sigma_0/4 \ \sigma_0/2 \ \sigma_0 \ 2\sigma_0 \ 4\sigma_0 \ 8\sigma_0]$. Pentru parametrii optimali găsiți de cele două metode s-a stabilit o medie de $[4.0039 \ 1.2555 \ 265.3732 \ 0.9041]$ pentru seturile de date USPS, Digit1, COIL2, respectiv Text.

Următoarele 6 rânduri arată rezultatele obținute în urma aplicării kernelului ISOMAP [Tenenbaum et al., 2000], kernelului de vecinătate [Weston et al., 2006], kernelului bagged cluster [Weston et al., 2006] și kernelul cluster multi-type cu

diferite funcții de transfer [Chapelle et al., 2002].

ISOMAP este o tehnică nelineară de reducere a dimensiunii care combină metodele PCA, MDS și aproximarea de varietate (*manifold approximation*). Pentru kernelul ISOMAP am utilizat un graf kNN cu k setat la 7.

Kernelul de vecinătate este bazat pe prezumția cluster și reprezintă fiecare punct ca și media vecinilor. Aici am încercat mai multe configurații ai parametrilor prin schimbarea numărului vecinilor de la 2 la 7; în tabelă apar cele mai bune rezultate (am omis cei mai buni parametri din cauza lipsei de spațiu). În cazul kernelului bagged cluster am setat parametrul $t = 20$, și în mod similar în cazul kernelului precedent am experimentat prin utilizarea diferitelor numere de clustere: K a fost ales din mulțimea $\{2, 3, \dots, 7\}$.

Kernelul cluster multi-type combină cu succes mai multe tehnici ca și clusteringul spectral, kernel PCA și mersurile aleatorii; funcțiile de transfer determină modul în care valorile proprii sunt transformate. În cazul acestui kernel am utilizat parametrii fixați:

- funcția transfer de pas (*step transfer function*): au fost utilizate cele mai mari valori proprii $\ell + 10$, unde ℓ este numărul exemplorilor etichetate;
- funcția transfer de pas linear: au fost utilizate cele mai mari valori proprii $\ell + 10$;
- funcția transfer polinomial: $t = 5$ (folosind notarea din [Chapelle et al., 2002]).

În cazul kernelurilor de vecinătate și kernelurilor bagged cluster am utilizat kernelul linear pentru kernelul de bază.

Următorii 15 rânduri – sub linia a doua orizontală – prezintă rezultatele obținute în urma utilizării kernelurilor noștri. HCK și gHCK notează clusterul ierarhic respectiv clusterul ierarhic bazat pe grafuri. În cazul gHCK, ca și în cazul kernelului ISOMAP, am utilizat un graf kNN cu $k = 7$. Cuvintele *single*, *complete* și *average* indică distanțele de legătură utilizate. RCK1, RCK2 și RCK3 notează kerneluri cluster de reponderare definite în formulele (4), (5), respectiv (6). Aici am experimentat cu trei tehnici de clustering: clusteringul k-means, ierarhic și spectral. Din cauza că ieșirea metodei k-means depinde puternic de centrul inițial al clusterului, am repetat clusteringul k-means de 10 ori și am determinat valorile de precizie ca media acestor 10 încercări. Pentru setul de date Text algoritmul k-means s-a dovedit a fi inadecvat din cauza dimensionalității sale mari (în acest caz inadecvat înseamnă timp îndelungat de acțiune). În cazul clusteringului ierarhic am utilizat următoarele metode: distanțe de legături single, complete, average, weighted average, centroid (legături medii pe grupă), median (legături medii ponderate pe grupă) și ward (metoda lui Ward). În tabelă am prezentat doar cele

mai bune metode. În încheiere am utilizat clusterigul spectral pentru obținerea kernelului: am utilizat metoda spectrală de clustering normalizat – varianta lui Shi și Malik [von Luxburg, 2006] cu clustering k-means. Astfel, ca și în cazul k-meansului convențional, am determinat valoarea de precizie cu determinarea mediei rezultatelor pe baza a 10 încercări. În cazul fiecărui algoritm de clustering numărul K a clusterelor a fost ales din mulțimea $\{2, 3, \dots, 7\}$ (ceea ce nu apare în tabelă). Ca și în cazul kernelurilor de vecinătate și de bagged cluster, am utilizat aici kernelul linear pentru kernelul de bază. Parametriile de mai jos ai kernelurilor de reponderare au fost fixate în decursul experimentelor:

- kernel gaussian de reponderare definit în (4): $\sigma = 1/\sqrt{2}$
- kernel de reponderare definit în (5): $\alpha = 0.3$
- kernel de reponderare definit în (6): $\beta = 0.5$

Prezentarea detaliată a rezultatelor poate fi găsită în capitolele corespunzătoare ale tezei.

6 Concluzii

Kernelurile dependente de date combină metodele kernel și de învățare semisupervizată prin construirea kernelurilor – astfel în mod implicit dând o nouă reprezentare a noilor exemple – cu utilizarea seturilor de date etichetate și neetichetate. Numim aceste kerneluri dependente de date pentru că funcția kernel depinde pe setul complet al punctelor etichetate și neetichetate disponibile în faza de instruire. Matematic vorbind, dacă $D_1 \neq D_2$, $\mathbf{x}, \mathbf{z} \in D_1 \cap D_2$,

$$k(\mathbf{x}, \mathbf{z}; D_1) \not\approx k(\mathbf{x}, \mathbf{z}; D_2)$$

unde “ $\not\approx$ ” se citește “nu neaparat egal”. În această teză am propus kerneluri dependente de date și metode de construcție a kernelurilor în condiții de învățare supervizată și semisupervizată. Am introdus următoarele kerneluri dependente de date:

- kerneluri de categorizare a documentelor bazate pe Wikipedia (Capitolul 5)
- kerneluri cluster ierarhici (Capitolul 6)
- kerneluri cluster de reponderare (Capitolul 7)

Categorizarea textelor este o sarcină de utilizare a unor exemple de instruire pentru construirea unei entități care învață supervizat, care este capabil să atribuie

una sau mai multe categorii documentelor în limbaje naturale. Kernelul bazat pe Wikipedia este construit prin aproximarea matricei Wikipedia cu o matrice de rang inferior pentru a filtra termenii irelevanți, și care în același timp asigură o reprezentare mai bogată a documentelor. Metoda propusă poate fi interpretată ca și o simplă înlocuire a matricei de covarianță a termenilor, clădită prin utilizarea documentelor de instruire, prin matricea de covarianță indus de Wikipedia. Din cauză că unii termeni de indexare aleși de metoda de selectare a trăsăturilor χ^2 nu au fost regăsiți în Wikipedia, avem aici două linii de bază: una cu 5209 și cealaltă cu 4601 termeni. Metoda noastră a întrecut ușor numai metoda de selectare a trăsăturilor χ^2 cu 4601 termeni. Motivul rezultatelor inferioare în acest caz poate fi specificitatea metodei de selectare a termenului de bază, diferențele de distribuție ale cuvintelor în Reuters-21578 și Wikipedia, etc. Posibilele căi de cercetare în această direcție includ alegerea altor metode de selecție a termenilor, găsirea reprezentărilor documentelor mai potrivite și efectuarea experimentelor și pe alte corpusuri TC.

Kernelurile cluster ierarhice sunt construite prin aplicarea proprietății ultrametrică a unor distanțe de legături utilizate în metodele de clustering aglomerativ. În capitolul 6 am propus două kerneluri de acest fel: unul care este bazat pe prezumția cluster și un altul care utilizează și prezumția de varietate. Desigur – ca și în cazul tehnicilor semisupervizate – acestea sunt așteptate să îmbunătățească performanța chiar dacă doar una dintre prezumții se dovedește a fi corectă. Am propus în același timp un cadru general pentru construcția kernelurilor cluster. Kernelurile au fost testați pe diferite seturi de date folosite pentru benchmark în [Chapelle et al., 2006]. HCK și gHCK au întrecut în mod semnificativ kernelurile lineari și gaussiani utilizați ca măsuri de referință; cea mai semnificativă îmbunătățire măsurată între cea mai bună linie de bază și rezultatele HCK/gHCK a fost de 13.5% în cazul setului de date COIL2. În același timp am comparat kernelurile propuse cu alte kerneluri dependente de date și le-am găsit comparabile cu acestea: HCK și gHCK au produs de mai multe ori rezultate mai bune decât celelate kerneluri. Pe baza experimentelor putem concluziona că aplicarea HCK și gHCK este utilă pentru aproape toate seturile de date. Este loc pentru perfecționare aici, de asemenea: testarea metodei propuse pentru calcularea valorilor kernel pentru punctele nevăzute.

În Capitolul 7 am prezentat trei scheme de reponderare pentru construcția kernelurilor cluster de reponderare (RCK). Ideea de bază este împrumutată de la kernelul bagged cluster: este construită o matrice (un kernel) pozitiv semidefinită prin care valorile de bază ale kernelului sunt reponderate: de exemplu un kernel linear sau gaussian. Am studiat două tipuri de kerneluri de reponderare: kernelul de reponderare gaussian și kernelurile de reponderare care utilizează produsele scalare ale vectorilor de apartenență cluster. Experimentele arată rezultate similare și în

cazul acestor kerneluri, ceea ce întrec în mod evident rezultatele de referință, și prezintă aceste kerneluri cluster comparabili cu kernelurile existente dependente de date.

Rezultatele cuprinse în această teză arată că prin utilizarea kernelurilor dependente de date metodele convenționale kernel pot fi depășite în mod semnificativ. Experimentele arată faptul că kernelurile propuși pot fi utilizate pe o varietate largă a seturilor de date. Din această cauză rezultatele cercetării întreprinse în prezenta teză pot servi ca mijloace utile și eficiente în multe domenii în care sunt folosiți algoritmi de clasificare.

Bibliografia tezei de doctorat

- K. M. Abadir and J. R. Magnus. *Matrix Algebra*. Cambridge University Press, 2005.
- A. Aizerman, E. M. Braverman, and L. I. Rozonoer. Theoretical foundations of the potential function method in pattern recognition learning. *Automation and Remote Control*, 25:821–837, 1964.
- R. Baeza-Yates and B. Ribeiro-Neto. *Modern Information Retrieval*. Addison-Wesley, 1999.
- S. Banerjee and T. Pedersen. An adapted lesk algorithm for word sense disambiguation using wordnet. In *CICLing*, pages 136–145, 2002.
- S. Basu. *Semi-supervised Clustering: Probabilistic Models, Algorithms and Experiments*. PhD thesis, The University of Texas at Austin, 2005.
- M. Belkin, P. Niyogi, and V. Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of Machine Learning Research*, 7:2399–2434, 2006.
- A. Berger. Error-correcting output coding for text classification. In *Proceedings of IJCAI-99 Workshop on Machine Learning for Information Filtering*, 1999.
- P. Berkhin. Survey of clustering data mining techniques. Technical report, Accrue Software, San Jose, CA, 2002.
- M. Bernstein, V. de Silva, J. C. Langford, and J. B. Tenenbaum. Graph approximations to geodesics on embedded manifolds, 2000.
- T. D. Bie. *Semi-Supervised Learning Based On Kernel Methods And Graph Cut Algorithms*. PhD thesis, Katholieke Universiteit Leuven, Kasteelpark Arenberg 10, 3001 Leuven (Heverlee), 2005.
- T. D. Bie, J. A. K. Suykens, and B. D. Moor. Learning from general label constraints. In A. L. N. Fred, T. Caelli, R. P. W. Duin, A. C. Campilho, and D. de Ridder, editors, *Structural, Syntactic, and Statistical Pattern Recognition, Joint IAPR International Workshops, SSPR 2004 and SPR 2004, Lisbon, Portugal, August 18-20, 2004 Proceedings*, volume 3138 of *Lecture Notes in Computer Science*, pages 671–679. Springer, 2004. ISBN 3-540-22570-6.
- C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer Verlag, New York, 2006.
- A. Blum and S. Chawla. Learning from labeled and unlabeled data using graph mincuts. In *Proc. 18th International Conf. on Machine Learning*, pages 19–26. Morgan Kaufmann, San Francisco, CA, 2001.
- Z. Bodó. Hierarchical cluster kernels for supervised and semi-supervised learning. In *Proceedings of the 4th International Conference on Intelligent Computer Communication and Processing (ICCP 2008)*, pages 9–16. IEEE, August 28–30 2008.
- Z. Bodó and Z. Minier. On supervised and semi-supervised k-nearest neighbor algorithms. In *Proceedings of the 7th Joint Conference on Mathematics and Computer Science*, volume LIII, pages 79–92. Studia Universitatis Babeş-Bolyai, Series Informatica, July 2008.
- Z. Bodó and Z. Minier. Semi-supervised feature selections with svms. In *Proceedings of the conference Knowledge Engineering: Principles and Techniques (KEPT 2009)*, pages 159–162. Presa Universitară Clujeană, July 2–4 2009. Special Issue of Studia Universitatis Babeş-Bolyai, Series Informatica.
- Z. Bodó, Z. Minier, and L. Csató. Text categorization experiments using Wikipedia. In *Proceedings of the conference Knowledge Engineering: Principles and Techniques (KEPT 2007)*, pages 66–72. Presa Universitară Clujeană, June 6–7 2007. Special Issue of Studia Universitatis Babeş-Bolyai, Series Informatica.
- I. Borg and P. J. F. Groenen. *Modern multidimensional scaling, 2nd edition*. Springer-Verlag, New York, 2005.
- B. E. Boser, I. Guyon, and V. N. Vapnik. A training algorithm for optimal margin classifiers. *Computational Learning Theory*, 5:144–152, 1992.
- C. J. C. Burges. A tutorial on support vector machines for pattern recognition. In *Knowledge Discovery and Data Mining*, volume 2, pages 121–167. 1998.
- R. Busa-Fekete and A. Kocsor. Locally linear embedding and its variants for feature extraction. In *IEEE*

- International Workshop on Soft Computing Applications, SOFA 2005*, 2005.
- C.-C. Chang and C.-J. Lin. *LIBSVM: a library for support vector machines*, 2001. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- O. Chapelle, J. Weston, and B. Schölkopf. Cluster kernels for semi-supervised learning. In S. Becker, S. Thrun, and K. Obermayer, editors, *NIPS*, pages 585–592. MIT Press, 2002. ISBN 0-262-02550-7.
- O. Chapelle, B. Schölkopf, and A. Zien. *Semi-Supervised Learning*. MIT Press, Sept. 2006. Web page: <http://www.kyb.tuebingen.mpg.de/ssl-book/>.
- Chung. Spectral graph theory (reprinted with corrections). In *CBMS: Conference Board of the Mathematical Sciences, Regional Conference Series*, 1997.
- D. A. Cohn, Z. Ghahramani, and M. I. Jordan. Active learning with statistical models. *Journal of Artificial Intelligence Research*, 4:129–145, 1996. ISSN 1076-9757.
- C. Corley, A. Csomai, and R. Mihalcea. Text semantic similarity, with applications. In *Proceedings of International Conference Recent Advances in Natural Language Processing (RANLP 2005)*, September 2005.
- T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*. MIT Press, Cambridge, MA, 2001.
- T. Cover and J. Thomas. *Elements of Information Theory, Second Edition*. Wiley-Interscience, 2006.
- T. M. Cover and P. E. Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, IT-13, 1967.
- K. Crammer and Y. Singer. A new family of online algorithms for category ranking. In *The 25th Annual International ACM SIGIR Conference*. ACM, 2002.
- N. Cristianini, J. Shawe-Taylor, A. Elisseeff, and J. S. Kandola. On kernel-target alignment. In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14 [Neural Information Processing Systems: Natural and Synthetic, NIPS 2001, December 3-8, 2001, Vancouver, British Columbia, Canada]*, pages 367–373. MIT Press, 2001.
- N. Cristianini, J. Shawe-Taylor, and H. Lodhi. Latent semantic kernels. *J. Intell. Inf. Syst*, 18(2-3): 127–152, 2002.
- N. Cristianini, J. Kandola, A. Vinokourov, and J. Shawe-Taylor. Kernel methods for text processing. In J. A. K. Suykens, G. Horváth, S. Basu, C. Micchelli, and J. Vandewalle, editors, *Advances in Learning Theory: Methods, Models and Applications*, pages 197–221. IOS Press, 2003.
- L. Csátó and Z. Bodó. *Neurális hálók és a gépi tanulás módszerei (Rețele neurale și metode de instruire automată)*. Presa Universitară Clujeană, Cluj-Napoca, 2008.
- L. Csátó and Z. Bodó. Decomposition methods for label propagation. In *Proceedings of the conference Knowledge Engineering: Principles and Techniques (KEPT 2009)*, pages 127–130. Presa Universitară Clujeană, July 2–4 2009. Special Issue of Studia Universitatis Babeş-Bolyai, Series Informatica.
- F. Debole and F. Sebastiani. An analysis of the relative hardness of reuters-21578 subsets. *Journal of the American Society for Information Science and Technology*, 56:971–974, 2004.
- S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman. Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41, June 1990.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39:1–38, 1977.
- I. S. Dhillon, Y. Guan, and B. Kulis. Kernel k-means: spectral clustering and normalized cuts. In *ACM SIGKDD – Knowledge discovery and data mining*, pages 551–556, 2004.
- T. G. Dietterich and G. Bakiri. Solving multiclass learning problems via error-correcting output codes. *Journal of Artificial Intelligence Research*, 2:263–286, 1995.
- L. Diosan, M. Oltean, A. Rogozan, and J.-P. Pécuchet. Improving SVM performance using a linear combination of kernels. In B. Beliczynski, A. Dzielinski, M. Iwanowski, and B. Ribeiro, editors, *Adaptive and Natural Computing Algorithms, 8th International Conference, ICANNGA 2007, Warsaw, Poland, April 11-14, 2007, Proceedings, Part II*, volume 4432 of *Lecture Notes in Computer Science*, pages

- 218–227. Springer, 2007. ISBN 978-3-540-71590-0.
- R. Duda, P. Hart, and D. Stork. *Pattern Classification*. John Wiley and Sons, 2001. 0-471-05669-3.
- W. K. Estes. *Classification and Cognition*. Oxford University Press, 1994. ISBN 9780195073355.
- B. Fischer, V. Roth, and J. M. Buhmann. Clustering with the connectivity kernel. In S. Thrun, L. K. Saul, and B. Schölkopf, editors, *NIPS*. MIT Press, 2003. ISBN 0-262-20152-6.
- N. Fuhr, S. Hartmann, G. Knorz, G. Lustig, M. Schwantner, and K. Tzeras. AIR/X – a rule-based multistage indexing system for large subject fields. In A. Lichnerowicz, editor, *Proceedings of RIAO-91, 3rd International Conference “Recherche d’Information Assistée par Ordinateur”*, pages 606–623, Barcelona, ES, 1991. Elsevier Science Publishers, Amsterdam, NL.
- J. Fürnkranz. Round robin classification. *Journal of Machine Learning Research*, 2:721–747, 2002.
- E. Gabrilovich and S. Markovitch. Computing semantic relatedness using wikipedia-based explicit semantic analysis. In *Proceedings of The 20th International Joint Conference on Artificial Intelligence (IJCAI)*, January 2007.
- I. P. Gent, P. Prosser, B. M. Smith, and W. Wei. Supertree construction with constraint programming. In *ICCP: International Conference on Constraint Programming (CP), LNCS*, 2003.
- A. Gliozzo and C. Strapparava. Domain kernels for text categorization. In *Proceedings of the Ninth Conference on Computational Natural Language Learning (CoNLL-2005)*, pages 56–63, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics.
- G. H. Golub and C. F. Van Loan. *Matrix Computations, 3rd Edition*. The Johns Hopkins University Press, Baltimore, MD, 1996.
- C. M. Grinstead and J. L. Snell. *Introduction to Probability*. AMS, 2003.
- I. Guyon and A. Elisseeff. An introduction to variable and feature selection. *Journal of Machine Learning Research*, 3:1157–1182, 2003.
- J. Ham, D. D. Lee, S. Mika, and B. Schölkopf. A kernel view of the dimensionality reduction of manifolds. In C. E. Brodley, editor, *Machine Learning, Proceedings of the Twenty-first International Conference (ICML 2004), Banff, Alberta, Canada, July 4-8, 2004*, volume 69 of *ACM International Conference Proceeding Series*. ACM, 2004.
- G. Hirst and D. St-onge. Lexical chains as representations of context for the detection and correction of malapropisms, Aug. 31 1997.
- P. Jackson and I. Moulinier. *Natural Language Processing for Online Applications: Text Retrieval, Extraction & Categorization*. John Benjamins, 2002. ISBN 90-272-4989-X.
- A. K. Jain and R. C. Dubes. *Algorithms for Clustering Data*. Prentice-Hall, 1988.
- A. K. Jain, M. N. Murty, and P. J. Flynn. Data clustering: A review. *CSURV: Computing Surveys*, 31, 1999.
- J. J. Jiang and D. W. Conrath. Semantic similarity based on corpus statistics and lexical taxonomy. *CoRR*, cmp-lg/9709008, 1997. informal publication.
- T. Joachims. Text categorization with support vector machines: Learning with many relevant features. Technical Report LS VIII-Report, Universität Dortmund, Dortmund, Germany, 1997.
- T. Joachims. Text categorization with support vector machines: Learning with many relevant features. In C. Nedellec and C. Rouveirol, editors, *Machine Learning: ECML-98, 10th European Conference on Machine Learning, Chemnitz, Germany, April 21-23, 1998, Proceedings*, volume 1398 of *Lecture Notes in Computer Science*, pages 137–142. Springer, 1998. ISBN 3-540-64417-2.
- I. T. Jolliffe. *Principal Component Analysis*. Series in Statistics. Springer Verlag, 2002.
- J. Kandola, J. Shawe-Taylor, and N. Cristianini. Optimizing kernel alignment over combinations of kernels. Technical Report 2002-121, Department of Computer Science, Royal Holloway, University of London, UK, 2002a.
- J. S. Kandola, J. Shawe-Taylor, and N. Cristianini. Learning semantic similarity. In S. Becker, S. Thrun, and K. Obermayer, editors, *NIPS*, pages 657–664. MIT Press, 2002b. ISBN 0-262-02550-7.
- C. Leacock and M. Chodorow. Combining local context and WordNet similarity for word sense identi-

- fication. In C. Fellbaum, editor, *WordNet: An Electronic Lexical Database*, pages 265–283. The MIT Press, Cambridge, Massachusetts, 1998.
- M. D. Lee, B. Pincombe, and M. Welsh. An empirical evaluation of models of text document similarity. In *CogSci2005*, pages 1254–1259, 2005.
- C. Leslie and R. Kuang. Fast kernels for inexact string matching. In *COLT: Proceedings of the Workshop on Computational Learning Theory*, Morgan Kaufmann Publishers, 2003.
- C. S. Leslie, E. Eskin, J. Weston, and W. S. Noble. Mismatch string kernels for SVM protein classification. In S. Becker, S. Thrun, and K. Obermayer, editors, *NIPS*, pages 1417–1424. MIT Press, 2002. ISBN 0-262-02550-7.
- D. D. Lewis. Evaluating and optimizing autonomous text classification systems. In *SIGIR '95*, pages 246–254, New York, NY, USA, 1995. ACM Press.
- D. Lin. An information-theoretic definition of similarity. In *ICML*, pages 296–304, 1998.
- H. Lodhi, C. Saunders, J. Shawe-Taylor, N. Cristianini, and C. Watkins. Text classification using string kernels. *Journal of Machine Learning Research*, 2:419–444, 2002.
- H. Lütkepohl. *Handbook of matrices*. John Wiley & Sons Ltd., Chichester, 1996. ISBN 0-471-97015-8.
- G. A. Miller, R. Beckwith, C. Fellbaum, D. Gross, and K. Miller. Five papers on wordnet. Technical Report, Cognitive Science Laboratory, Princeton University, 1993.
- Z. Minier, Z. Bodó, and L. Csató. Segmentation-based feature selection for text categorization. In *Proceedings of the 2nd International Conference on Intelligent Computer Communication and Processing (ICCP 2006)*, pages 53–59. IEEE, September 1–2 2006.
- Z. Minier, Z. Bodó, and L. Csató. Wikipedia-based kernels for text categorization. In *Proceedings of the 9th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC 2007)*, pages 157–164. IEEE, September 26–29 2007.
- T. M. Mitchell. The discipline of machine learning. Technical Report CMU-ML-06-108, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, 2006.
- R. Morelos-Zaragoza. *The Art of Error Correcting Coding*. John Wiley and Sons, Inc., pub-WILEY:adr, 2002. ISBN 0-471-49581-6.
- A. Moschitti. A study on optimal parameter tuning for rocchio text classifier. In F. Sebastiani, editor, *Advances in Information Retrieval, 25th European Conference on IR Research, ECIR 2003, Pisa, Italy, April 14-16, 2003, Proceedings*, volume 2633 of *Lecture Notes in Computer Science*, pages 420–435. Springer, 2003. ISBN 3-540-01274-5.
- A. Y. Ng, M. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14*, Cambridge, MA, 2002. MIT Press.
- K. Nigam. *Using unlabeled data to improve text classification*. PhD thesis, Carnegie Mellon University, 2001.
- K. Nigam, A. McCallum, S. Thrun, and T. M. Mitchell. Text classification from labeled and unlabeled documents using EM. *Machine Learning*, 39(2/3):103–134, 2000.
- L. Page, S. Brin, R. Motwani, and T. Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford Digital Library Technologies Project, 1998.
- S. Patwardhan, S. Banerjee, and T. Pedersen. Using measures of semantic relatedness for word sense disambiguation. In A. F. Gelbukh, editor, *Computational Linguistics and Intelligent Text Processing, 4th International Conference, CICLing 2003, Mexico City, Mexico, February 16-22, 2003, Proceedings*, volume 2588 of *Lecture Notes in Computer Science*, pages 241–257. Springer, 2003. ISBN 3-540-00532-3.
- J. C. Platt, N. Cristianini, and J. Shawe-Taylor. Large margin DAGs for multiclass classification. In S. A. Solla, T. K. Leen, and K.-R. Müller, editors, *NIPS*, pages 547–553. The MIT Press, 1999. ISBN 0-262-19450-3.
- P. Resnik. Using information content to evaluate semantic similarity in a taxonomy. In *IJCAI*, pages

- 448–453, 1995.
- C. J. V. Rijsbergen. *Information Retrieval*. Butterworths, 1979.
- S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290:2323–2326, 2000.
- S. Russel and P. Norvig. *Artificial Intelligence: a Modern Approach*. Prentice-Hall, 1995.
- G. Salton, A. Wong, and A. C. S. Yang. A vector space model for automatic indexing. *Communications of the ACM*, 18:229–237, 1975.
- L. K. Saul and S. T. Roweis. An introduction to locally linear embedding, 2001.
- B. Schölkopf and A. J. Smola. *Learning with Kernels*. The MIT Press, Cambridge, MA, 2002.
- B. Schölkopf, A. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. Technical Report 44, Max-Planck Institut für biologische Kybernetik, Arbeitsgruppe Bülthoff, Spemannstrasse 38, 2076 Tobingen, Germany, Dec. 1996.
- B. Schölkopf, A. J. Smola, and K.-R. Müller. Kernel principal component analysis. *Advances in kernel methods: support vector learning*, pages 327–352, 1999.
- H. Schütze. Automatic word sense discrimination. *Computational Linguistics*, 24(1):97–123, 1998.
- F. Sebastiani. Machine learning in automated text categorization. *ACM Computing Surveys*, 34(1):1–47, 2002.
- J. Shlens. A tutorial on principal component analysis, 2005.
- V. Sindhwani, P. Niyogi, and M. Belkin. Beyond the point cloud: from transductive to semi-supervised learning. In L. D. Raedt and S. Wrobel, editors, *Machine Learning, Proceedings of the Twenty-Second International Conference (ICML 2005), Bonn, Germany, August 7-11, 2005*, volume 119 of *ACM International Conference Proceeding Series*, pages 824–831. ACM, 2005. ISBN 1-59593-180-5.
- D. Sloughter. The calculus of functions of several variables, 2001.
- J. B. Tenenbaum, V. de Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, Dec. 2000.
- H. tien Lin and C. jen Lin. A study on sigmoid kernels for SVM and the training of non-PSD kernels by SMO-type methods, 2003.
- D. Tikk. *Szövegbányászat*. Typotex, Budapest, 2007.
- V. Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- G. Varelas. Semantic similarity methods in wordnet and their application to information retrieval on the web. Technical Report TR-TUC-ISL-01-2005, Technical Univ. of Crete (TUC), Dept. of Electronic and Computer Engineering, Chania, Crete, Greece, 2005.
- S. Vishwanathan and N. M. Murty. Kernel enabled K-means algorithm. Technical report, The Indian Institute of Science, Bangalore, 2002.
- S. V. N. Vishwanathan and A. J. Smola. Fast kernels for string and tree matching. In S. Becker, S. Thrun, and K. Obermayer, editors, *NIPS*, pages 569–576. MIT Press, 2002. ISBN 0-262-02550-7.
- S. V. N. Vishwanathan, K. M. Borgwardt, O. Guttman, and A. J. Smola. Kernel extrapolation. *Neurocomputing*, 69(7-9):721–729, 2006.
- U. von Luxburg. A tutorial on spectral clustering. Technical Report 149, Max Planck Institute for Biological Cybernetics, August 2006.
- J. Voss. Measuring wikipedia, 2005.
- J. H. Ward, Jr. Hierarchical grouping to optimize an objective function. *Journal of the American Statistical Association*, 58:236–244, 1963.
- C. wei Hsu and C. jen Lin. A comparison of methods for multi-class support vector machines, 2001.
- J. Weston and C. Watkins. Support vector machines for multiclass pattern recognition. In *Proceedings of the Seventh European Symposium On Artificial Neural Networks*, 4 1999.
- J. Weston, A. Elisseeff, B. Schölkopf, and M. Tipping. Use of the zero-norm with linear models and kernel methods. *Journal of Machine Learning Research*, 3:1439–1461, 2003.
- J. Weston, C. Leslie, D. Zhou, A. Elisseeff, and W. S. Noble. Semi-supervised protein classification

- using cluster kernels, 2004.
- J. Weston, C. Leslie, E. Ie, and W. S. Noble. Semi-supervised protein classification using cluster kernels. In O. Chapelle, B. Schölkopf, and A. Zien, editors, *Semi-Supervised Learning*, chapter 19, pages 343–360. MIT Press, 2006.
- S. K. M. Wong, W. Ziarko, and P. C. N. Wong. Generalized vector spaces model in information retrieval, proceedings of the 8th annual international ACM SIGIR conference on research and development in information retrieval. In J. M. Tague, editor, *Proceedings of the 8th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 18–25, Montreal, Quebec, Canada, June 1985. ACM Press.
- S. K. M. Wong, W. Ziarko, V. V. Raghavan, and P. C. N. Wong. On modeling of information retrieval concepts in vector spaces. *ACM Trans. on Database Sys.*, 12(2):299, June 1987.
- B. Y. Wu and K.-M. Chao. *Spanning Trees and Optimization Problems*. Chapman and Hall/CRC, Boca Raton, Florida, 2004.
- Z. Wu and M. S. Palmer. Verb semantics and lexical selection. In *ACL*, pages 133–138, 1994.
- Y. Yang and C. G. Chute. A linear least squares fit mapping method for information retrieval from natural language texts. In *COLING*, pages 447–453, 1992.
- Y. Yang and C. G. Chute. An example-based mapping method for text categorization and retrieval. *ACM Transactions on Information Systems*, 12(3):252–295, July 1994.
- Y. Yang and J. O. Pedersen. A comparative study on feature selection in text categorization. In *International Conference on Machine Learning*, pages 412–420, 1997.
- Z.-P. Yang, X. Zhang, and C.-G. Cao. Inequalities involving khatri-rao products of hermitian matrices. *The Korean Journal of Computational & Applied Mathematics*, 9(1):125–133, 2002. ISSN 1229-9502.
- Q. Yong and Y. Jie. Geodesic distance for support vector machines. *Acta Automatica Sinica*, 31(2): 202–208, 2005.
- J. A. Zdziarski. *Ending spam: Bayesian content filtering and the art of statistical language classification*. No Starch Press, pub-NO-STARCH:adr, 2005. ISBN 1-59327-052-6.
- D. Zhou, B. Schölkopf, and T. Hofmann. Semi-supervised learning on directed graphs. In *NIPS*, 2004.
- X. Zhu. *Semi-supervised learning with graphs*. PhD thesis, Carnegie Mellon University, Pittsburgh, PA, USA, 2005. Chair-John Lafferty and Chair-Ronald Rosenfeld.
- X. Zhu and Z. Ghahramani. Learning from labeled and unlabeled data with label propagation. Technical Report CMU-CALD-02-107, Carnegie Mellon University, 2002.
- X. Zhu, T. J. Rogers, R. Qian, and C. Kalish. Humans perform semi-supervised classification too. In *AAAI*, page 864. AAAI Press, 2007. ISBN 978-1-57735-323-2.