

Turinys

IVADAS	1
1. DIMENSIJŲ ATRINKIMAS	4
1.1. Baziniai dimensijų atrinkimo metodai	4
1.1.1. Fisher'io įvertis	5
1.1.2. Atpalaidavimo metodas	5
1.1.3. Asimetrinis priklausomybės koeficientas	6
1.1.4. Absoliučių svorių SVM	6
1.1.5. Rekursyvus dimensijų eliminavimas pagal SVM	7
1.2. Multikriterinis dimensijų atrinkimas	7
1.2.1. Svoriais grįstas multikriterinis suliejimas	8
1.2.2. Reitingais grįstas multikriterinis suliejimas	9
1.2.3. Svoriais ir reitingais grįstas multikriterinis suliejimas	10
1.2.4. Multikriterinis rekursyvus dimensijų eliminavimas	11
2. DIMENSIJŲ ATRINKIMO STABILUMAS	12
2.1. Stabilumo matavimas	13
2.2. Pavienių dimensijų atrinkimo metodų stabilumas	14
2.2.1. Fisher'io dimensijų atrinkimo metodas	14
2.2.2. Atpalaidavimo dimensijų atrinkimo metodas	14
3. Eksperimentai	14
3.1. Naudoti duomenys	14
3.2. Metodologija	15
3.3. Dimensijų atrinkimo metodų sparta	15
LITERATŪRA	15

ĮVADAS

Atsiranda vis tikslesni būdai, pvz. genomo sekvenavimas, gauti biomedicininis duomenis. Tikslumas šiame kontekste reiškia, kad didėja atributų, dažniau vadinamų duomenų dimensijomis, skaičius. Sekvenuojant genomą moderniomis priemonėmis galima gauti net keletą milijonų dimensijų vienam genui(!). Duomenys turintys daug dimensijų yra vadinami daugiamačiais duomenimis.

Norint daugiamačius duomenis suskirstyti į pageidaujamas kategorijas, pvz. atskirti sergančius nuo nesergančių pacientų, pagal vidinę duomenų struktūrą (šis procesas vadinamas klasifikavimu) reikia specifinių klasifikavimo strategijų. To reikia todėl, kad tradiciniai klasifikavimo būdai yra nepajėgūs dirbti su daugiamačiais duomenimis. Taikant tradicines klasifikavimo strategijas daugiamačiams duomenims skaičiavimo laikas pasidaro nebepriimtinas, mažėja klasifikavimo tikslumas.

Viena iš strategijų dirbti su daugiamačiais duomenimis yra mažinti duomenų dimensijų skaičių - naudoti dimensijų atrinkimo (*angl. feature selection*) metodus. Naudojant dimensijų atrinkimo metodus galima supaprastinti duomenis atsirenkant tik tas dimensijas, kurios yra svarbios konkrečiai tiriamai problemai, pvz. nustatant ligos priežastis. Dimensijų atrinkimas yra svarbi duomenų apdorojimo (*angl. preprocessing*) etapo dalis. Pasirinkus našius dimensijų atrinkimo metodus galima sukurti tikslesnius klasifikavimo modelius, sumažinti skaičiavimams reikalingų resursų poreikį, pagreitinti klasifikavimo modelio kūrimo - mokymosi - procesą, taip pat jis padeda vizualizuoti bei geriau suprasti tais duomenimis apibūdinamus procesus.

Norint geriau suprasti biomedicininis duomenis itin svarbu fokusuoti dėmesį į sąlyginai nedidelį dimensijų poaibį. Dimensijų aibės sumažinimas paspartina biomedicininis duomenų tyrimus - tyrėjams reikia atlikinėti bandymus su mažesniu mėginių skaičiumi. Mažesnio skaičiaus mėginių tyrimas kainuoja mažiau, nes mažiau reikia žmonių darbo laiko, mažiau reikia ir cheminių reagentų. Tačiau tyrėjams ne mažiau svarbu yra žinoti kaip tų atrinktųjų dimensijų aibė varijuoja priklausomai nuo nedidelių pasikeitimų pačioje duomenų aibėje - šis kriterijus yra vadinamas dimensijų atrinkimo stabilumu (*angl. robustness*).

Jei atrinktųjų dimensijų poaibis, kurį naudojant duomenų objektai yra išskirstomi teisingoms kategorijoms (klasifikavimo rezultatai yra tikslūs), yra stabilus, tai reiškia, kad tikslinga detalius bandymus su duomenimis pradėti nuo atrinktojo stabilaus dimensijų poaibio. Kitu atveju bandymus atlikinėti reikia su visa duomenų aibe, kas yra mažiau efektyvu. Pastebėtina, kad stabilumo matavimus tikslinga atlikinėti tik atsižvelgiant į klasifikavimo tikslumą.

Dar viena problema dirbant su daugiamačiais duomenimis yra tai, kad dažnai turimas labai ribotas skaičius duomenų objektų (*angl. tuple*). Objektų - dimensijų santykis (ODS) gali skirtis keliais šimtais kartų. Atrodytų, tik laiko klausimas, kada bus paruošta daugiau duomenų objektų, bet žvelgiant

į duomenų gavybos tendencijas pasidaro aišku, kad objektų skaičius niekada nepavys dimensijų skaičiaus. Todėl reikia apgalvoti, kaip bus į tai atsižvelgta kuriant klasifikavimo modelius. Nes esant mažam objektų skaičiui kyla grėsmė, kad klasifikatorius bus sukurtas toks, kuris gerai veiks tik, su tais duomenimis, kuriais remiantis jis buvo sukurtas, tačiau netiksliai klasifikuos naujus duomenis. Tokia problema yra vadinama persimokymo problema (angl. overfitting).

Išgaunamuose biomediciniuose duomenyse didėja dimensijų skaičius, todėl daugėja ir nereikalingų ar klaidingų duomenų - triukšmo. Triukšmas atsiranda dėl įvairių priežasčių, pvz, cheminiai preparatai buvo netinkamai paruošti. Duomenyse esant triukšmui su jais tampa sudėtinga dirbti, prastėja klasifikavimo rezultatai. Atsiranda poreikis duomenų apdorojimo etape identifikuoti triukšmingus duomenis ir juos pašalinti iš duomenų rinkinio.

Taigi, dirbant su daugiamatiais duomenimis, reikia atsižvelgti į keletą kriterijų:

1. Klasifikavimo tikslumą;
2. Dimensijų atrinkimo stabilumą, atsižvelgiant į klasifikavimo rezultatus;
3. Triukšmo lygį duomenyse;
4. Skaičiavimo išteklių naudojimo racionalumą.

Reikalavimas vienu metu atsižvelgti į keletą kriterijų apsunkina užduotį. Klasifikuojant daugiamatius duomenis uždavinys yra surasti geriausius rezultatus duodančią strategiją, kuri geriausiai atsižvelgia į aukščiau minėtus kriterijus.

Darbo eksperimentinei daliai reikalingus skaičiavimo išteklius, suteikė VU MIF ITC. Duomenų apdorojimo algoritmų implementavimui buvo naudojama R programavimo kalba. Eksperimentai atlikti profesinės praktikos MII metu.

Šio darbo tikslas yra išanalizuoti darbo su daugiamatiais duomenis ypatybes.

Šiam darbui yra keliamos tokios užduotys:

1. Apžvelgti esamus klasifikavimo metodus;
2. Išanalizuoti bazinių dimensijų atrinkimo metodų spartą;
3. Išanalizuoti bazinių dimensijų atrinkimo metodų stabilumą;
4. Išanalizuoti kaip dimensijų atrinkimo metodai įtakoja klasifikatorių tikslumą;
5. Pasiūlyti naują metodą daugiamatį duomenų klasifikavimui.

Bakalaurinis darbas suskirstytas į 3 skyrius:

1. Pirmajame skyriuje apžvelgsiu nagrinėjamą dalykinę sritį;
2. Antrajame skyriuje pristatysiu siūlomą metodą daugiamatį duomenų klasifikavimui;

3. Trečiajame skyriuje aprašysiu darbo metu atliktus eksperimentus.

1. DIMENSIJŲ ATRINKIMAS

Klasifikuojant daugiamacių duomenis susiduriame su taip vadinamu dimensiškumo prakeiksmu (angl. curse of dimensionality). Pavyzdžiui, kai dimensijų skaičius įkopia į trečią ar ketvirtą eilę naudoti paprastus klasifikavimo algoritmus tampa nebeefektyvu nei laiko, nei klasifikavimo našumo atžvilgiais. Vienas iš būdų kovoti su dimensiškumo prakeiksmu yra naudoti vienokius ar kitokius dimensijų skaičiaus mažinimo metodus. Dimensijų atrinkimas yra svarbus etapas duomenų apdorojimui daugelyje mašininio mokymosi taikymų, jis dažniausiai yra naudojamas surasti mažiausią dimensijų poaibį, kuris maksimaliai pagerina klasifikavimo modelio našumą.

Pagal tai, kaip dimensijų atrinkimo metodai bendradarbiauja su klasifikatoriumi, dimensijų atrinkimo metodus galima išskirstyti į tris kategorijas:

1. Filtravimo metodai (angl. filter methods). Jie dirba tiesiogiai su duomenimis, o jų rezultatas gali būti dimensijų įvertinimas svoriais, dimensijų reitingavimas ar tiesiog dimensijų poaibis. Tokių metodų pagrindinis privalumas yra tai, kad jie yra greiti ir nepriklausomi nuo klasifikavimo metodo, tačiau gali sumažinti klasifikavimo tikslumą.
2. Įvyniojimo metodai (angl. wrapper methods). Jie atlieka paiešką dimensijų aibėje vadovaudamiesi klasifikavimo modeliu (pvz. klasifikavimo našumu po pakartotinio įvertinimo). Jie dažnai duoda geresnius rezultatus negu filtravimo metodai, bet yra reiklesni resursams.
3. Įdėtiniai metodai (angl. embedded methods). Jie dimensijų atrinkimui naudoja vidinius klasifikatoriaus duomenis (pvz. svoriai gauti pagal SVM). Šie metodai dažnai siūlo gerus mainus tarp klasifikavimo našumo ir skaičiavimų sudėtingumo.

Šiame skyriuje nagrinėsiu bazinius dimensijų atsirinkimo metodus, keletos kriterijų suliejimą metodą (angl. feature selection based on multicriterion fusion)[YM11], bei stabilių dimensijų grupių išskyrimo metodą[LYD09].

1.1. Baziniai dimensijų atrinkimo metodai

Yra pasiūlyta daugybė dimensijų atrinkimo metodų. Šiame skyriuje aptarsiu keletą taip vadinamų bazinių¹ dimensijų atrinkimo metodų:

1. Fisher'io įvertis (angl. Fisher ratio)[PWCG01];

¹Aptarsiu tik bazinius todėl, kad jie yra svarbiausi, nes jų tarpusavio rezultatai yra nekoreliuojantys, o nekoreliuojančius rezultatus galima panaudoti juos apjungiant, taip gauname sinergijos efektą.

2. Atpalaidavimo (ang. relief) metodas[RSK03];
3. Asimetrinis priklausomybės koeficientas[Sha01] (ADC) (angl. Asymmetric Dependency Coefficient);
4. Absoliučių svorių SVM[Vap00] (AW-SVM) (angl. Absolute Weight SVM)
5. Rekursyvus dimensijų eliminavimas pagal SVM[GWBV02] (SVM-RFE) (angl. Recursive Feature Elimination by SVM)

1.1.1. Fisher'io įvertis

Fisher'io įvertis vertina individualias dimensijas pagal jų klasių atskiriamąją galią. Dimensijos įvertis yra sudarytas iš tarpklasinio skirtumo santykio su vidiniu klasės pasiskirstymu:

$$FR(j) = \frac{(\mu_{j1} - \mu_{j2})^2}{\sigma_{j1}^2 + \sigma_{j2}^2}, \quad (1)$$

kur, j - yra dimensijos indeksas, μ_{jc} - dimensijos j reikšmių vidurkis klasėje c , σ_{jc}^2 - dimensijos j reikšmių standartinis nuokrypis klasėje c , kur $c = 1, 2$. Kuo didesnis yra Fisher'io įvertis, tuo geriau ta dimensija atskiria klases.

1.1.2. Atpalaidavimo metodas

Atpalaidavimo metodas iteratyviai skaičiuoja dimensijų „susietumą“. Pradžioje „susietumas“ visoms dimensijoms yra lygus nuliui. Kiekvienoje iteracijoje atsitiktinai² pasirenkamas objektas iš duomenų bazės, surandami artimiausi kaimynai iš tos pačios ir kitos klasės, ir atnaujinamos visų dimensijų „susietumo“ reikšmės. Dimensijos įvertis yra vidurkis visų objektų atstumų iki artimiausių kaimynų iš tos pačios ir kitos klasės:

$$W(j) = W(j) - \frac{diff(j, x, x_H)}{n} + \frac{diff(i, x, x_M)}{n}, \quad (2)$$

kur $W(j)$ - j -osios dimensijos „susietumo“ įvertis, n - objektų aibės dydis, x - atsitiktinai pasirinktas objektas, x_H - artimiausias kaimynas iš tos pačios klasės (angl. nearest-Hit), x_M - artimiausias kaimynas iš kitos klasės, $diff(j, x, x')$ - j -osios dimensijos reikšmių skirtumas tarp laisvai pasirinkto objekto x ir atitinkamo kaimyno, kur skirtumą į intervalą $[0, 1]$ normalizuojanti funkcija yra:

$$diff(j, x, x') = \frac{|x_j - x'_j|}{x_{jmax} - x_{jmin}}, \quad (3)$$

kur x_{jmax} ir x_{jmin} yra maksimali ir minimali j -osios dimensijos reikšmės. „Susietumo“ reikšmių atnaujinimas yra vykdomas n kartų ir kuo didesnė galutinė

²Pastebėtina, kad metodas turi atsitiktinį elementą, todėl klasifikavimo ir dimensijų atrinkimo stabilumo rezultatai dažniausiai šiek tiek varijuoja nekeičiant konfigūracijos.

reikšmė, tuo svarbesnė dimensija. Pastebėtina, kad šis algoritmas veikia tik su dviejomis klasėmis, nors yra ir išplėtimų.

1.1.3. Asimetrinis priklausomybės koeficientas

Asimetrinis priklausomybės koeficientas yra dimensijų reitingavimo metodas, kuris matuoja klasės Y etiketės (angl. label) tikimybinę priklausomybę j -ajai dimensijai, naudodamas informacijos prieaugį (angl. information gain):

$$ADC(Y, j) = \frac{MI(Y, X_j)}{H(Y)}, \quad (4)$$

kur $H(Y)$ - klasės Y entropija, o $MI(Y, X_j)$ - yra bendrumo informacija (angl. mutual information) tarp klasės etiketės Y ir j -osios dimensijos

$$H(Y) = - \sum_y p(Y = y) \log p(Y = y), \quad (5)$$

$$H(X_j) = - \sum_x p(X_j = x) \log p(X_j = x), \quad (6)$$

$$MI(Y, X_j) = H(Y) + H(X_j) - H(Y, X_j), \quad (7)$$

$$H(Y, X_j) = - \sum_{y, x_j} p(y, x_j) \log p(y, x_j), \quad (8)$$

Kuo didesni ADC įverčiai, tuo dimensija yra svarbesnė, nes turi daugiau informacijos apie duomenų klases.

1.1.4. Absoliučių svorių SVM

Atraminių vektorių metodas (SVM) yra vienas populiariausių klasifikavimo algortimų, nes jis gerai susidoroja su daugiamatiais duomenimis. Yra keletas bazinių SVM variantų, bet šiame darbe naudosime tiesinį SVM, nes jis demonstruoja gerus rezultatus dirbant su genų ekspresijos duomenimis. Tiesinis SVM yra hiperplokštuma apibrėžta kaip:

$$\sum_{j=1}^p w_j x_j + b_0 = 0, \quad (9)$$

kur p - dimensijų kiekis, w_j - j -osios dimensijos svoris, x_j - j -osios dimensijos kintamasis, b_0 - konstanta. Dimensijos absoliutus³ svoris w_j gali būti panaudotas dimensijų reitingavimui. Pastebėtina, kad svorių nustatymas yra atliekamas tik vieną kartą⁴.

³Svorį reikia imti absoliutaus dydžio, nes neigiamas svoris implikuoja priklausomybę vienai klasei, o teigiamas kitai klasei.

⁴SVM-RFE - metodas svorius nustato daug kartų.

1.1.5. Rekursyvus dimensijų eliminavimas pagal SVM

Rekursyvus dimensijų eliminavimas pagal SVM yra vienas populiariausių dimensijų atrinkimo algoritmų. Todėl, jis yra naudojamas, kaip atskaitos taškas (angl. threshold) vertinant kitus dimensijų atrankos metodus. Iš esmės šis metodas yra daugkartinis absoliučių svorių SVM metodo taikymas nuolat išmetinėjant dimensijas su mažiausiais svoriais. Rekursyvus dimensijų eliminavimas mums padeda surasti tinkamą dimensijų poaibį, kas nevisada pavyksta su dimensijų reitingavimo metodais. Bendroji rekursyvaus dimensijų eliminavimo procedūra: Jei trečiajame algoritmo žingsnyje yra pašalinama tik viena

Algorithm 1 Rekursyvus dimensijų eliminavimas

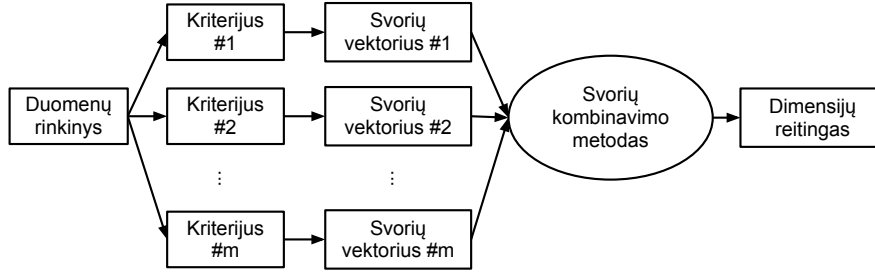
1. Turime pilną dimensijų rinkinį F_0 , nustatome $i = 0$.
 2. Įvertiname kiekvienos dimensijos kokybę dimensijų aibėje F_i .
 3. Išmetame mažiausiai svarbią dimensiją iš F_i tam, kad gautume dimensijų rinkinį F_{i+1} .
 4. Nustatome $i = i + 1$ ir grįžtame į antrąjį žingsnį kol nėra patenkinta algoritmo pabaigos sąlyga.
-

dimensija, tai gauname dar ir dimensijų reitingavimą, o jei pašalinamos kelios dimensijos ar jų dalis (pvz. 50%) tai reitingavimo negauname. Pastebėtina, kad rekursyvus dimensijų eliminavimas gali labai padidinti algoritmo sudėtingumą skaičiavimo resursų atžvilgiu. Algoritmo pabaigos sąlyga gali būti koks nors konkretus dimensijų skaičius arba tiesiog dimensijų aibę mažinti tol, kol dimensijų visai nebeliks.

1.2. Multikriterinis dimensijų atrinkimas

Multikriterinio dimensijų atrinkimo metodų esmė yra panaudoti kelis dimensijų atrinkimo metodus suliejant jų rezultatus į vieną bendrą rezultatą. Kodėl naudinga sulieti keletą dimensijų atrinkimo metodų rezultatų? Sulieti keletą dimensijų atrinkimo metodų rezultatų naudinga, nes pavieniai dimensijų atrinkimo metodai be to, kad turi savitų privalumų, visada turi ir savo silnybių, pavyzdžiui, jautrumas išimtims (angl. outliers), negali rasti dimensijų tarpusavio priklausomybių, etc. Yra skiriamos trys priežastys, kodėl keletas kombinuotų silpnų ir nestabilių dimensijų atrinkimo metodų gali duoti geresnius rezultatus[Die00]:

1. Keletas skirtingų, bet vienodai optimalių hipotezių gali būti teisingos, ir kriterijų kombinavimas sumažina tikimybę, kad bus pasirinkta neteisinga hipotezė;
2. Atskiri kriterijai gali dirbti skirtinguose lokaliuose optimumuose, tuo tarpu kombinavimas gali geriau reprezentuoti tikrąją duomenų funkciją;



1 pav.: Svoriais grįstas multikriterinis suliejimas.

3. Tikroji duomenų funkcija negali būti reprezentuojama jokia hipoteze paskiro algoritmo hipotezių erdvėje ir agreguojant pavienių metodų rezultatus galima praplėsti hipotezių erdvę.

Suliejant keletą skirtingų metodų suliejamose gerosios pavienių dimensijų atrinkimo metodų savybės, taip kompensuojant algoritmų silpnybes.

Galima pavienių dimensijų atrankos rezultatus sulieti pagal šias suliejimo strategijas:

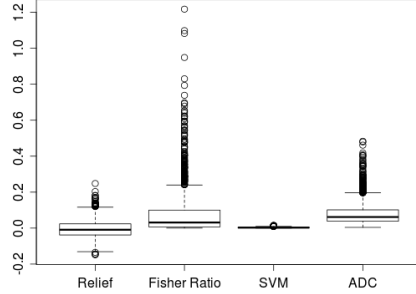
1. Svoriais grįstas multikriterinis suliejimas;
2. Reitingais (angl. rank) grįstas multikriterinis suliejimas;
3. Svoriais ir reitingais grįstas multikriterinis suliejimas.

1.2.1. Svoriais grįstas multikriterinis suliejimas

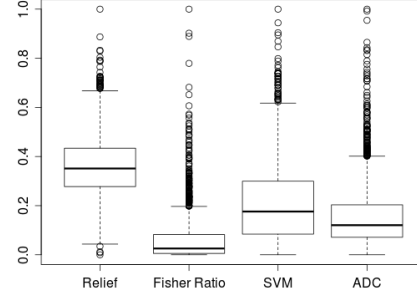
Svoriais grįsto multikriterinio dimensijų atrinkimo suliejimo pagal svorius algoritmo pirmajame žingsnyje kiekvienas bazinis metodas priskiria duomenų rinkinio dimensijoms svorius, tada tie svoriai yra kombinuojami į vieną sutarties (angl. consensus) svorių vektorių, kurio pagrindu yra gaunami dimensijų reitingai. Algoritmas yra pavaizduotas 1 pav.

Suliejant svorius svarbu yra užtikrinti, kad svoriai, gauti naudojant skirtingus bazinius kriterijus, būtų palyginami. Todėl svorių normalizavimas turi būti atliekamas prieš svorių kombinavimą. Kitu atveju dimensijų įvertinimo metodai bus nepalyginami. Paveikslėlyje 2 pav. nenormalizuotų pavienių dimensijų vertinimo metodų skiriasi netgi suteiktų svorių intervalai. Paveikslėlyje 3 pav. matome, kad net ir normalizavus svorius gana stipriai skiriasi svorių kvartilai - į tai reikia atkreipti dėmesį interpretuojant galutinius dimensijų vertinimo rezultatus. Šiame darbe svoriai yra normalizuoti intervale $[0, 1]$ pagal formulę:

$$u'_i = \frac{u_i - u_{imin}}{u_{imax} - u_{imin}}, \quad (10)$$



2 pav.: Pavienių dimensijų atrinkimo metodų nenormalizuotas svorių pasiskirstymas.



3 pav.: Pavienių dimensijų atrinkimo metodų normalizuotas svorių pasiskirstymas.

kur u_i - dimensijų svorių vektorius pagal i kriterijų, $u_{i\min}$ - minimali u_i svorių vektoriaus reikšmė, $u_{i\max}$ - maksimali u_i svorių vektoriaus reikšmė, u'_i - normalizuotų svorių vektorius.

Sutarties svorių vektorius u yra vidurkis normalizuotų svorių vektorių:

$$u = \frac{1}{m} \sum_{i=1}^m u'_i, \quad (11)$$

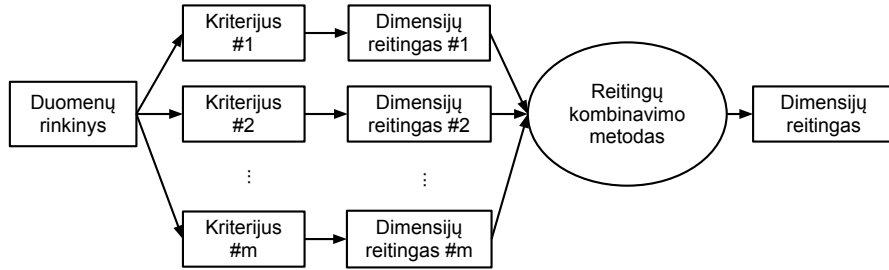
kur m yra bazinių kriterijų skaičius. Reikia paminėti, kad didesnė svorio reikšmė reiškia, kad dimensija yra geresnė.

1.2.2. Reitingais grįstas multikriterinis suliejimas

Reitingais grįsto multikriterinio suliejimo pagal reitingus metodas gauna duomenų rinkinio dimensijų reitingą, pagal keletą bazinių dimensijų reitingavimo kriterijų. Algoritmo pirmajame žingsnyje keletas dimensijų atrinkimo kriterijų grąžina dimensijų reitingu, paskui tie reitingai yra kombinuojami į vieną bendrą dimensijų reitingą. Algoritmas yra pavaizduotas 4 pav. Suliejimo pagal reitingus metodas nereikalauja dimensijų atrinkimo metodų rezultatų normalizavimo, nes tiesiog imame dimensijoms priskirtus reitingus ir juos kombinuojame. Skirtingai nei suliejimo pagal svorius algoritme, baziniai dimensijų atrinkimo kriterija dimensijų eliminavimas[YM11] susideda iš dviejų dalių: keletos dimensijų atrinkimi turi gražinti dimensijų reitingus, o ne svorius.

Dimensijų reitingų kombinavimui yra keletas metodų[DKNS01], tačiau paprastumo dėlei šiame darbe naudosiu Borda balsavimą⁵ (angl. Borda count).

⁵Dar žinomas kaip „Pažymių metodas“. Jis buvo pasiūlytas prancūzų matematiko ir fiziko Jean-Charles de Borda 1770 metais.



4 pav.: Reitingais grįstas multikriterinis suliejimas.

Tarkime, kad turime m balsuotojų ir p kandidatų aibę. Tada Borda balsavimo metodas kiekvienam i -ajam balsuotojui sukuria balsų vektorių v_i tokiu būdu: geriausiai įvertintam kandidatui suteikiama p taškų, antrajam kandidatui $p - 1$, ir t.t. Galutiniai taškai yra gaunami sudedant visų balsuotojų taškus

$$v = \sum_{i=1}^m v_i, \quad (12)$$

kur v yra suminių taškų vektorius, o iš jo galime gauti ir dimensijų reitingus.

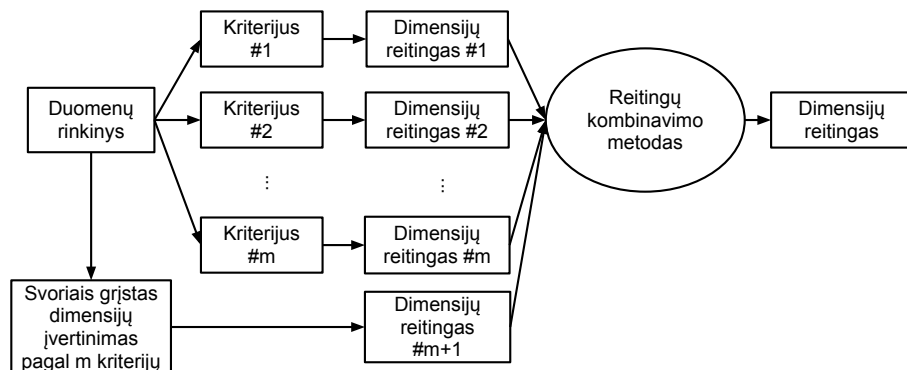
1.2.3. Svoriais ir reitingais grįstas multikriterinis suliejimas

Svoriais ir reitingais grįsto multikriterinio suliejimo metodas nuo reitingais grįsto multikriterinio suliejimo metodo skiriasi tuo, kad kaip dar vienas reitingas yra panaudojamas svoriais grįsto multikriterinio dimensijų atrinkimo metu gautas reitingas. Multikriterinio dimensijų įverčių ir pagal svorius, ir pagal reitingus metodas vyksta trimis žingsniais:

1. Gauname dimensijų reitingus pagal m pavienių dimensijų atrinkimo metodus;
2. Suliejame dimensijų įverčius pagal svorius ir taip gauname vieną dimensijų reitingą;
3. Reitinguojame dimensijas pagal visus turimus $m + 1$ pavienius reitingus.

Algoritmas yra pavaizduotas 5 pav.

Kadangi yra suliejami keli mažai koreliuojantys dimensijų reitingavimo metodai, yra pasiekiamas didesnis dimensijų atrinkimo stabilumas, kai varijuoja treniravimosi duomenų poaibis (angl. subsampling).



5 pav.: Svoriais ir reitingais grįstas multikriterinis suliejimas.

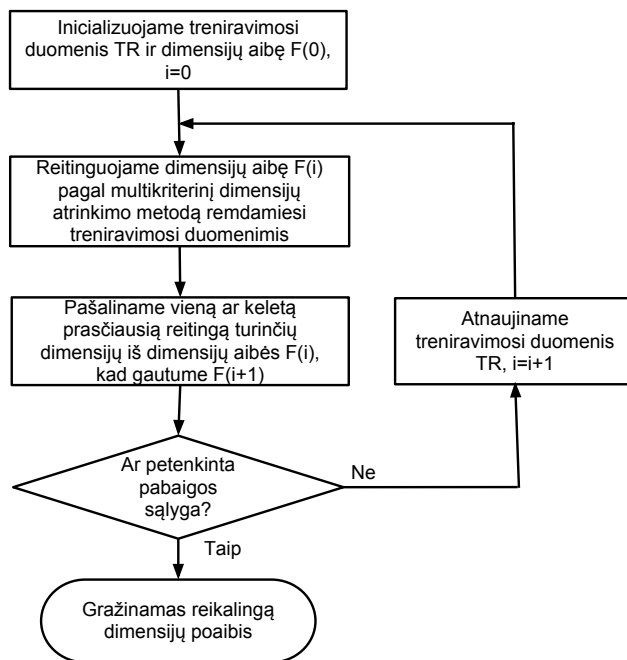
1.2.4. Multikriterinis rekursyvus dimensijų eliminavimas

Jei dimensijų atrinkimo tikslas yra pagerinti klasifikavimo rezultatus, tai taikymas multikriterinio dimensijų atrinkimo metodų nebūtinai duos pageidaujamą rezultatą, nes yra pastebėta, kad vien dimensijų reitingavimas nebūtinai suranda geriausią dimensijų poaibį. Tam, kad būtų surastas geriausias dimensijų poaibis reikia kombinuoti multikriterinį dimensijų reitingavimą su paieškos strategija. Rekursyvus dimensijų eliminavimas yra dažnai naudojama paieškos strategija dimensijų atrinkimui. Todėl yra kombinuojamas multikriterinis dimensijų reitingavimas ir rekursyvus dimensijų eliminavimas.

Multikriterinis rekursyvus dimensijų eliminavimas[YM11] susideda iš dviejų dalių: keletos dimensijų atrinkimo kriterijų suliejimo ir pagal svorius, ir pagal reitingus, ir rekursyvaus dimensijų eliminavimo aprašyto algoritme nr. 1. Algoritmas pavaizduotas 6 pav.

Yra pastebėta, kad standartinis rekursyvus dimensijų eliminavimas, kai vienos iteracijos metu yra eliminuojama viena dimensija, gali labai padidinti algoritmo sudėtingumą. Todėl genų ekspresijos duomenims yra rekomenduotina eliminuoti keletą dimensijų vienu metu.

Nors SVM-RFE dimensijų atrinkimo algoritmas ir yra labai populiarus, tačiau yra žinoma, kad jam trūksta stabilumo. Todėl kombinuodami didesnę stabilumą turintį multikriterinį dimensijų atrinkimą su rekursyvaus dimensijų eliminavimo paieškos strategija, turėtume gauti stabilesnį dimensijų atrinkimo algoritmą.



6 pav.: Multikriterinio rekursyvaus dimensijų eliminavimo algoritmas.

2. DIMENSIJŲ ATRINKIMO STABILUMAS

Dimensijų atrinkimo technikų stabilumas gali būti apibrėžtas kaip dimensijų atrinkimo rezultatų variacijos dėl mažų pakeitimo duomenų rinkinyje. Pakeitimai duomenų rinkinyje gali būti duomenų objektų lygio (pvz. pridedami ar atimami duomenų objektai), dimensijų lygio (pvz. pridedant dimensijoms triukšmo) ar abiejų lygių kombinacija.

Dimensijų atrinkimo technikų stabilumas yra vis didesnę svarbą įgaunanti tyrimų kryptis. Stabilumo aktualumas yra sąlygotas to, kad biologiniuose duomenyse galima gana užtikrintai daryti prielaidą, kad konkrečiai problemai yra aktualios tik tam tikros dimensijos. Todėl dalykinės srities ekspertams yra aktualu naudoti tik tuos dimensijų atrinkimo metodus, kurie yra stabilūs ir relevantiški modeliuojamai problemai, nes tai atpigina tolimesnę duomenų analizę.

Šiame skyriuje apžvelgsime teorinį stabilumo matavimų modelį. Taip pat įvertinsime pavienių dimensijų atrinkimo metodų stabilumą taikant juos įvairiems duomenų rinkiniams. Taip pat išanalizuosime situaciją, kai kombinuojami keletos dimensijų atrinkimo metodų rezultatai.

2.1. Stabilumo matavimas

Vertinant dimensijų atrinkimo metodų stabilumą yra svarbu kaip panašiai yra atrenkamos dimensijos, kai yra atliekamas dimensijų atrinkimas su vis kitu duomenų poaibiu. Kuo panašesnius dimensijų atrinkimo rezultatus gausime, tuo stabilumas yra didesnis. Vidutinis stabilumas gali būti apibrėžtas kaip vidurkis visų reitingavimo metu gautų sąrašų porų tarpusavio panašumo įverčių:

$$S_{tot} = \frac{2 \sum_{i=1}^{k-1} \sum_{j=i+1}^k S(f_i, f_j)}{k * (k - 1)}, \quad (13)$$

kur k žymi kiek kartų buvo imtas skirtingas poaibis objektų dimensijų atrinkimui, f_i, f_j - dimensijų atrinkimo rezultatas - reitingai, $S(f_i, f_j)$ - yra kokia nors panašumo matavimo funkcija.

Kaip matome dimensijų atrinkimo stabilumas priklauso nuo to, kokią panašumo funkciją naudosime. Tradicinės panašumo funkcijos (persidengimo procentas, Pearson'o koreliacija, Spearman'o koreliacija, Jaccard indeksas) gali būti taikomos, bet jos yra linkusios priskirti didesnes panašumo reikšmes, kai pasirenkamas didesnis dimensijų poaibis. Taip yra dėl padidėjusio sisteminio nuokrypio (ang. bias), nes imant didesnę poaibį padidėja tikimybė tiesiog atsitiktinai pasirinkti dimensiją. Kad išvengtume šios problemos panašumui vertinti buvo pasirinktas Kunchevos [Kun07] indeksas:

$$KI(f_i, f_j) = \frac{r * N - s^2}{s * (N - s)} = \frac{r - (s^2/N)}{s - (s^2/N)}, \quad (14)$$

kur $s = |f_i| = |f_j|$ yra atrinktų dimensijų aibės dydis, $r = |f_i \cap f_j|$ - abiem atrinktiems dimensijų poaibiems bendrų dimensijų skaičius, N - bendras duomenų aibės dimensijų skaičius. Pastebėtina, kad formulėje esantis atėminys s^2/N ištaiso sisteminį nuokrypį atsirandantį dėl galimybės atsitiktinai pasirinkti dimensijas. Kunchevos indeksas gali įgyti reikšmes iš intervalo $[-1, 1]$, kur didesnė reikšmė reiškia didesnę panašumą, o artimos nuliui reikšmės reiškia, kad dimensijos atrenkamos daugiausia atsitiktinai. Kunchevos indekso ypatybė yra ta, kad jis atsižvelgia tik į persidengiančias, tačiau visiškai nekreipia dėmesio į koreliuojančias dimensijas.

Vertinant stabilumą tarp skirtingų metodų gali iškilti problemų, nes ne visi dimensijų atrinkimo metodai gražina rezultatą tokiu pačiu formatu. Šiame darbe dimensijų atrinkimo metodų rezultatas yra ne dimensijai priskirtas svoris, bet dimensijos reitingas. Todėl f_i yra sąrašas, kurio ilgis yra N , kur pirmas sąrašo elementas yra geriausią reitingą turinčios dimensijos numeris, o paskutinis sąrašo elementas yra blogiausią reitingą turinčios dimensijos numeris.

Galiausiai yra svarbu paminėti, kad dimensijų stabilumas nėra matuojamas visiškai nepriklausomai - jis yra matuojamas atsižvelgiant į klasifikavimo rezultatus. Dimensijų atrinkimo metodų stabilumas yra matuojamas tik tada

kai atrinktos dimensijos duoda gerus klasifikavimo rezultatus. Taip yra, nes kokios nors dalykinės srities ekspertui, nėra naudingos tos dimensijų atrinkimo strategijos, kurios duoda labai stabilius rezultatus, bet nėra naudingos klasifikavimo modelio kūrime.

2.2. Pavienių dimensijų atrinkimo metodų stabilumas

Šiame skyriuje apžvelgsime pavienių dimensijų atrinkimo metodų stabilumą. Stabilumas visiems metodams buvo matuojamas atsižvelgiant į klasifikavimo rezultatus - buvo matuojamas stabilumas tikra

2.2.1. Fisher'io dimensijų atrinkimo metodas

2.2.2. Atpalaidavimo dimensijų atrinkimo metodas

Šis metodas yra vienas nestabiliausių, nes pasikeitimai duomenų rinkinyje stipriai įtakoja rezultatus.

3. Eksperimentai

3.1. Naudoti duomenys

Šiame darbe eksperimentai buvo atliekami su biomedicininiais viešai prieinamais genų ekspresijos duomenų rinkiniais.

1 lentelė. Darbe naudoti duomenų rinkiniai

Pavadinimas	Šaltinis	Objektų skaičius (+/-)	Dimensijų skaičius	ODS
Gaubtinės žarnos auglys (angl. Colon)	[ABN ⁺ 99]	62 (40/22)	2000	3,1%
Centrinės nervų sistemos auglys (CNS)	[PTG ⁺ 02]	60 (39 AAL / 21 AML)	7129	0.84%
Prostatos auglys	[SFR ⁺ 02]	102 (52/50)	6033	1.7%
Šizofrenija ir maniakinė depresija	[Ins]	90 (bp ⁶ : sz ⁷ : cc ⁸ =30:31:29)	22283	0.403%

⁶bp (angl. Bipolar disorder) - maniakine depresija sergantys pacientai.

⁷sz (angl. Schizophrenia) - šizofrenija sergantys pacientai.

⁸cc (angl. Control Crowd) - kontrolinė grupė.

Duomenų rinkinius apibūdinantis dydis ODS, kuris turimiems duomenims tesiekia nuo 0,403% iki 3,01% procento, parodo, kad turime labai retus (angl. sparse) duomenis, o tai labai apsunkina mokymosi procesą ir gali sukelti persimokymo (angl. overfitting) problemą.

3.2. Metodologija

3.3. Dimensijų atrinkimo metodų sparta

Literatūra

- [ABN⁺99] U. Alon, N. Barkai, D.A. Notterman, K. Gish, S. Ybarra, D. Mack, and A.J. Levine. Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proceedings of the National Academy of Sciences*, 96(12):6745, 1999.
- [Die00] T. Dietterich. Ensemble methods in machine learning. *Multiple classifier systems*, pages 1–15, 2000.
- [DKNS01] C. Dwork, R. Kumar, M. Naor, and D. Sivakumar. Rank aggregation revisited. In *Proceedings of WWW10*, pages 613–622, 2001.
- [GWBV02] Isabelle Guyon, Jason Weston, Stephen Barnhill, and Vladimir Vapnik. Gene selection for cancer classification using support vector machines. *Mach. Learn.*, 46(1-3):389–422, March 2002.
- [Ins] Stanley Medical Research Institute. Online genomics database. [žiūrėta 2012-04-03].
- [Kun07] Ludmila I. Kuncheva. A stability index for feature selection. In *Artificial Intelligence and Applications*, pages 421–427, 2007.
- [LYD09] Steven Loscalzo, Lei Yu, and Chris Ding. Consensus group stable feature selection. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '09, pages 567–576, New York, NY, USA, 2009. ACM.
- [PTG⁺02] S.L. Pomeroy, P. Tamayo, M. Gaasenbeek, L.M. Sturla, M. Angelo, M.E. McLaughlin, J.Y.H. Kim, L.C. Goumnerova, P.M. Black, C. Lau, et al. Prediction of central nervous system embryonal tumour outcome based on gene expression. *Nature*, 415(6870):436–442, 2002.
- [PWCG01] Paul Pavlidis, Jason Weston, Jinsong Cai, and William Noble Grundy. Gene functional classification from heterogeneous data. In *Proceedings of the fifth annual international conference on Computa-*

- tional biology*, RECOMB '01, pages 249–255, New York, NY, USA, 2001. ACM.
- [RSK03] Marko Robnik-Sikonja and Igor Kononenko. Theoretical and empirical analysis of relieff and rrelieff. *Machine Learning*, 53(1-2):23–69, 2003.
 - [SFR⁺02] D. Singh, P.G. Febbo, K. Ross, D.G. Jackson, J. Manola, C. Ladd, P. Tamayo, A.A. Renshaw, A.V. D’Amico, J.P. Richie, et al. Gene expression correlates of clinical prostate cancer behavior. *Cancer cell*, 1(2):203–209, 2002.
 - [Sha01] C. E. Shannon. A mathematical theory of communication. *SIG-MOBILE Mob. Comput. Commun. Rev.*, 5(1):3–55, January 2001.
 - [Vap00] V.N. Vapnik. *The nature of statistical learning theory*. Statistics for engineering and information science. Springer, 2000.
 - [YM11] F. Yang and KZ Mao. Robust feature selection for microarray data based on multicriterion fusion. *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, 8(4):1080–1092, 2011.