

Geeniexpressiohahmojen haku moniulotteisesta datasta

Aija Niissalo

Helsinki 1.2.2008

Tiedon louhinnan seminaari, Ohjaaja: Hannu Toivonen

HELSINGIN YLIOPISTO

Tietojenkäsittelytieteen laitos

Sisältö

| | |
|-------------------------------------------------------|-----------|
| 1 Johdanto | 1 |
| 2 Geeniekspressio ja hakutehtävä | 2 |
| 2.1 Geeniekspressio | 2 |
| 2.2 Hakutehtävän formalisointi | 4 |
| 3 Etäisyysmitat ja moniulotteisuus | 6 |
| 3.1 Euklidinen etäisyys ja korrelaatiomitat | 8 |
| 3.2 Matriisin singulaariarvohajotelma, SVD | 9 |
| 3.3 Moniulotteisuus ja Simpsonin paradoksi | 10 |
| 4 Hakualgoritmit | 11 |
| 4.1 Yksiulotteinen tapaus | 12 |
| 4.2 Yleinen tapaus: heuristinen ratkaisu | 13 |
| 4.3 Pienimmän neliösumman sovitus SVD:llä | 14 |
| 5 Hakutulosten vertailu ja loppupäätelmät | 15 |
| Lähteet | 18 |
| Liitteet | |
| 1 Sanasto | |
| 2 DNA-mikrosiru | |

1 Johdanto

Emäs- ja aminohapposekvenssien samankaltaisuutta voidaan hakea tehokkaasti eliöiden genomien DNA-molekyyleistä ja solujen proteiineista. Mm. BLAST-hakutuloksia [ASG90] hyödyntäen on voitu päätellä sekvenssien välisiä toiminnallisia ja evoluutioon liittyviä suhteita. Geeniekspressio tarkoittaa geenien ilmentymistä eli solun tapahtumasarjaa, jossa DNA:n sisältämä emäs-koodi kopioidaan RNA:ksi ja edelleen tämän RNA-emäsjärjestyksen ohjaamaksi proteiiniksi. (Geeni)ekspressiota mitataan tehokkaasti mm. DNA-mikrosirujen (DNA-array, DNA-chip, microarray) laboratorikokeilla. Mikrosirumittaus on ns. tehoseulontamenetelmä, joka tuottaa kerralla monta tuhatta ekspresioarvoa reaali-lukuna. Hahmonhaku reaaliarvoisesta geeniekspressiodatasta on luonteeltaan erilaista kuin merkkijonosekvenssien haku. Ekspressiodata on moniulotteista: geenejä on tyypillisesti tuhansia ja olosuhteita jopa satoja erilaisia samassa laboratoriomittausarjassa. Kuitenkin helppoudessaan BLAST:in kaltaiselle joka-biologin hahmonhakuvälineelle olisi käyttöä valtavan ekspresiodatamäärän analysoinnissa.

Ekspresiodataa on varastoituna julkisissa tietokannoissa nyt jo yli sadalle eliölle, monille tuhansille mittausarjoille ja sadoille tuhansille mikrosiruille. Määrät kasvavat jatkuvasti. Samanlaisten laboratoriomittauksen ja mittausarjojen etsiminen metatiedon avulla voi viedä kohtuuttomasti aikaa ja käyttäjän pitää tietää, millaista mittausarjaa hän on hakemassa. Geeniekspressiotietoja voidaan käyttää mm. solun toimintamekanismien mallintamisessa, sairauksien luokittelussa, sairausennusteissa sekä lääkkeiden vaikutusten arvioinnissa.

Ajankohtainen ja kiinnostava kysymys on, onko pelkästään ekspresiodatan lukuarvojen suuruuteen perustuva yleiskäyttöinen hakukone mahdollinen. Tietokannoissa olevien mittausarjojen hakutuloksiin saataisiin silloin paremmuusjärjestys. Haku halutaan tehokkaaksi, jolloin hakuun ei annetta kaikkien geenien ekspresioarvoja vaan käyttäjän mielestä tärkeä tai muuten kiinnostava joukko. Olosuhteita ja muuta metatietoa ei hakutilanteessa tarvittaisi. Samansuuntainen arvojakauma, geeniekspressiohahmo eli profiili, kysely- ja tietokantamittausarjan välillä voisi paljastaa yllättävää samankaltaisuutta eri koeolosuhteiden synnyttämissä tiloissa. Hakutuloksen biologinen tulkinta ja tilastollinen analyysi voitaisiin tehdä vastauksena saatavien kokeiden tarkkoja tietoja tulkitsemalla.

Tässä kirjoituksessa esittelen geeniekspressiohahmojen hakualgoritmit, jotka perustuvat reaaliarvoisten mittauspisteiden (1) Euklidiseen etäisyyteen, (2) Pearso-

nin korrelaatioon, (3) Spermanin rank-korrelaatioon tai (4) pienimmän neliösumman sovittukseen singulaariarvohajotelmaa käyttäen (Singular Value Decomposition, SVD). Hakualgoritmien kokeilutietokanta on ladattu opetusministeriön hallinnoiman Tieteen tietotekniikan keskuksen, CSC:n YeastBASE-hiivatietokannan käyttöliittymästä: <http://yeastbase.csc.fi/>. Kokeilutietokantaan otin YeastBASE:n cDNA-sirujen (kokeiluun sopivan) ekspressiodatan.

Luvussa 2 kerrotaan, mitä geeniekspressio on ja miten sitä mitataan. Lisäksi formalisoidaan hakutehtävä ja mainitaan, millaisia hakukoneita on nykyisin käytössä. Luvussa 3 esitellään hakualgoritmeissa käytetyt reaaliarvoisen datan etäisyysmitat sekä moniulotteisen ekspressiodatan analyysiongelmia. Itse hakualgoritmit kuvataan luvussa 4. Loppupäätelmät ja hakutulosten vertailu ovat viimeisessä luvussa 5. Liitteissä annetaan biologisten käsitteiden sanasto ja lyhyt DNA-mikrosirujen kuvaus.

2 Geeniekspressio ja hakutehtävä

Tässä luvussa kuvaan, mitä geeniekspressio on ja miten sitä mitataan sekä mitä laboratoriomittauksessa kullekin geenille saatu ekspressioarvo tarkoittaa. Lisäksi formalisoin hahmonhakutehtävän, joka kohdistuu yhden tai monen mittauksen tuottamaan reaaliarvoiseen datamatriisiin. Tavoitteena on löytää hahmoja, joiden perusteella kyselynä annetun ja tietokannassa olevien mittausten samankaltaisuutta voidaan arvioida. Jos samanlaisia hahmoja löytyy, voidaan tehdä hypoteeseja solujen tiloista. Nämä testataan uusilla laboratorionkokeilla.

2.1 Geeniekspressio

Geeniekspressio on (DNA-molekyylissä olevien) geenien koodaamien proteiinien tuottamista solussa. 1990-luvulta lähtien ekspression suuruutta on voitu mitata kvantitatiivisesti ns. tehoseulontamenetelmillä, kuten DNA-mikrosirujen laboratorionkokeilla [SSD95, LDB96, EiB99]. Mikrosirumittaukset perustuvat Edvin Southernin tekemään keksintöön, jolla solunäytteen sekvenssit löytyvät DNA-fragmentteja ja komplementaarista RNA:ta hybridisoimalla [Sou75]. cDNA-siru mittaa näytesolujen (esim. $10^8 - 10^9$ kpl) hetkellisen tilan poikkeamaa kontrollitilasta ja antaa tuloksena eliön geenien ilmentymisen suhteelliset reaaliarvot yhdessä mittauksessa. Geenejä on tuhansia, ja näiden samanaikaisen ilmentymisen mittaustulos voi kertoa jotain oleellista näytteen biologisesta tilasta, voidaan esimerkiksi ennustaa syöpäpo-

tilaan selviytymismahdollisuus. Koska näyte on peräisin suuresta solupopulaatiosta, yhdellä sirulla mitataan itse asiassa monen tilan yhteisvaikutusta. Tämä on tärkeä muistaa, kun hahmonhakua mallinnetaan ja tuloksia verrataan.

Näyteliuoksessa olevat ekspression aikaansaamat molekyylit on leimattu eri fluori-soivalla värillä kuin kontrolliliuoksen. Leima-aineiden lähettämän valon intensiteetti mitataan laserlaitteella. (Tarkempi sirukuvaus löytyy liitteestä 2, DNA-mikrosiru, jos tämä teknologia ei ole tuttu lukijalle. Ainakin liitteen viimeiseen kuvaan kannattaa tutustua.) Olkoon näytteen ekspression aikaansaama (laserlaitteella mitattu ja normalisoitu) valon intensiteetti geenille i ja sirulle j I_{red} ja vastaavan kontrollinäytteen intensiteetti I_{green} . Tällöin yksittäinen ekspressioarvo datassa on:

$$g_{ij} = \log_2 \frac{I_{red}}{I_{green}}.$$

Geeniekspressiolla tarkoitetaan jatkossa näitä arvoja. Arvot esitetään datamatriisina, jossa riveillä yksilöidään geenit ja sarakkeilla sirut.

Ekspressiodata-arvojen muutosta kutsutaan geeniekspression profiiliksi (gene expression profile). Geenisuuntainen muutos, geeniprofiili (gene profile), kuvaa yhden geenin suhteellisen ekspressiotason vaihtelua monessa eri sirulla tehdyssä mittauksessa. Sirusuuntainen muutos, siruprofiili (array profile), kuvaa yhdellä sirulla mitattua kaikkien geenien suhteellisen ekspressiotason vaihtelua toisiinsa nähden. Näiden geeniekspressiohahmojen etsintä formalisoidaan seuraassa aliluvussa moniulotteisen ekspressiodatan hakutehtävänä ekspressioarvojen muodostamassa vektoriavaruudessa, joka esitetään matriisiesityksenä. Näin hahmoja voidaan vertailla ja hyödyllisiä vastaavuuksia louhia.

Muutama mikrosirujen hahmonhakukone on toteutettu. Hughesin ja kumppaneiden tutkimus [HMJ00] on ensimmäinen suuri koosteisen ekspressiodatan tutkimus. Koska data oli tuotettu homogeenisissa koeolosuhteissa, kokoava hierarkkinen klusterointi [SSZ98, ESB98] Pearsonin korrelaatioetäisyyden kanssa ryhmitteli yhteen biologisesti samankaltaisia tiloja. Lamb ja kumppanit [LCP06] ovat kehittäneet lääkeaine-geeni-sairaus-yhteyden etsivän hakukoneen, Connectivity Map, ihmisen heterogeeniselle geeniekspressiodatalle. Justin Lamb luokittelee [Lam07] esimerkiksi eurooppalaisen ArrayExpress ja vastaavan amerikkalaisen Gene Expression Omnibus (GEO) ekspressiotietokannat super-koosteiksi, joilla Connectivity Mapin kaltainen datan koostaminen on ratkaisematta. Helppo käyttäjäkyselyiden haku (ekspressioarvoja käyttäen) ei niissä ole mahdollista. The Connectivity Mapin hakualgoritmi noudattaa parametritonta, järjestykseen perustuva hahmonhakustrategiaa, joka pe-

rustuu muunneltuun Kolmogorov-Smirnov statistiikkaan. Hakualgoritmi on Subramanian ja kumppaneiden [STM05] formalisoima GSEA-algoritmi, Gene Set Enrichment Analysis. CellMontage-sovelluksella [FKT⁺07] voi hakea kyselyekspressoarvojen samanlaisuutta esimerkiksi osasta GEO-tietokannan dataa. Menetelmä käyttää Spearmanin rank-korrelaatiokerrointa samanlaisuuden mittana. Alusta (platform, esim. tietty mittausinstrumentti ja eliö) pitää kyselyssä valita. Haku tehdään sirukohtaisesti, koko mittaussarjan hahmoa ei haeta (kuten ei Connectivity Mapissäkään).

2.2 Hakutehtävän formalisointi

Ekspressiomittaussarja sisältää yhden tai monta (DNA-siru)mittausta. Mittaussarjan tulos on reaalitylukuarvoinen datamatriisi, jossa geenit ovat rivejä ja sarakkeet mittauksia, yhden sirun tuottamia ekspressoarvoja. Tehtävä on tiedonhaku näistä vaihtelevaulotteisista datapistejoukoista (koe- eli mittaussarjoista), joita datamatriisit kuvaavat ja joissa yksi piste vastaa yhtä mittausta.

Kyselyn kohteena oleva tietokanta koostuu datapistejoukoista, jotka sijaitsevat eriulotteisissa reaalitylukuavaruuksissa ja joiden koordinaattien vastaavuus ei ole triviaalia. Datapistejoukkojen vastaavuusongelman aiheuttaa mittaussarjojen erilainen koko; mittauksia voi olla yhdessä mittaussarjassa jopa yli sata tai vain yksi, jolloin datamatriisi kutistuu yhdeksi pystyvektoriksi. Geeniavaruuden ulottuvuus oletetaan samaksi kaikille mittauksille. Geeniavaruuden moniulotteisuuteen liittyy kuitenkin ongelmia.

Kyselyä annetaan uusi datapistejoukko ja vastauksena halutaan tietokannan datapistejoukot järjestettynä sen mukaan, kuinka samanlaisia ne ovat kyselypistejoukon kanssa. Jotta vastaaminen onnistuu, pitää ratkaista eriulotteisten ja toisiaan vastamattomien avaruuksien kohdentaminen toistensa suhteen siten, että vertailukelpoisuus toteutuu. Myös yhden datapisteen vastaavuuden etsiminen on mielenkiintoinen tapaus. Seuraavaksi formalisoidaan edellä kuvattu tilanne. Samoja merkintöjä käytetään jatkossa.

Vektorit, matriisit ja joukot merkitään lihavoidulla fontilla; reaalityluvut, kuten lukumäärät, indeksit ja arvot merkitään kursivilla. Reaalitylukujen joukolle on oma merkintä: \mathbb{R} . Vektoreita merkitään myös alkioittain: $\mathbf{x} = (x_1, x_2, \dots, x_p)$.

Jokainen ekspressiomittaus antaa ekspressoarvon kullekin mitatulle geenille. Näin saadaan ekspressoarvoja, g_i , missä i kuuluu järjestettyyn joukkoon geenit, $\mathbf{G} =$

$\{1, 2, \dots, m\}$. Voidaan olettaa että joukko \mathbf{G} on aina sama ja sen alkioiden järjestyks on määrätty. Käytännön nimenvastaavuusongelmat ja puuttuva data ratkaistaan erikseen. Kukin mittaus antaa datan $\mathbf{g} = (g_1, g_2, \dots, g_m)$, joka on mittauksen ekspressioarvoista muodostettu vektori. Mittaussarja sisältää n kappaletta edellä määriteltyjä mittauksia \mathbf{g}_j , missä j kuuluu joukkoon mittaukset, $\mathbf{A} = \{1, 2, \dots, n\}$. Mittaussarja voi sisältää mm. kontrollimittauksia, yksittäisiä eri olosuhteissa tehtyjä mittauksia tai/ja aikasarjamittauksia (esim. 'control/control', 'treatment/control', 'time1/control', 'time2/control'). Mittaussarja antaa n datavektoria. Ne yhdessä muodostavat mittaussarjan \mathbf{D} , joka on alussa mainittu datapistejoukko. Mittaukset vaihtelevat eri mittaussarjoissa sekä tulkinnaltaan että lukumäärältään.

Ekspressiodatan tietokanta \mathbf{DB} , johon haku kohdistuu, on kokoelma tällaisia $m \times n$ -datamatriiseja \mathbf{D} . Kukin tietokannan datamatriisi \mathbf{D}_k kuuluu joukkoon $\mathbb{R}^{m \times n}$, jossa m on geenien lukumäärä ja n on sirujen lukumäärä siten, että m on kaikille sama, mutta n voi vaihdella. k on mittaussarjan indeksi tietokannassa. Datamatriisi \mathbf{D}_k koostuu mittausarvoista g_{ij} siten, että i kuuluu joukkoon \mathbf{G} ja j kuuluu joukkoon \mathbf{A} . Silloin tietokanta \mathbf{DB} kuuluu joukkoon $\mathbb{R}^{m \times \sum n}$ ja \mathbf{D}_k kuuluu joukkoon $\mathbf{DB} = \{\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_K\}$, missä K on mittaussarjojen lukumäärä tietokannassa. Koko tietokannan sirujen määrää, $\sum n$, merkitään jatkossa p :llä.

Kyselyn lähtökohtana on tilanne, jossa on suoritettu uusi mittaussarja \mathbf{Q} (erikoistapauksena yksittäinen mittaus), ja halutaan tietää, onko tietokannassa mittaussarjaa, jossa suhteellinen ekspressio on muuttunut samalla tavoin kuin uudessa mittaussarjassa. Lisäksi ollaan kiinnostuneita vain joistakin geneistä, ei koko joukosta \mathbf{G} . Kyselyssä on annettava ne geenit joita kysely koskee (= jokin joukon \mathbf{G} osajoukko) ja näille uudessa laboratoriokeessa saadut mittausarvot. Genejä valitaan M kpl. Mittausarvoja on kullakin geenillä N kappaletta, jos kokeessa tehtiin N mittauksia. Kyselyssä on $M \times N$ mittausarvoa. Näin kysely \mathbf{Q} kuuluu joukkoon $\mathbb{R}^{M \times N}$, jossa M on kyselygeenien lukumäärä ja N on kyselysirujen lukumäärä ja \mathbf{G}_M , joka sisältyy joukkoon \mathbf{G} , on kyselygeenien joukko. Kysely koostuu mittausarvoista q_{ij} siten, että i kuuluu joukkoon $\{1, 2, \dots, M\}$ ja j kuuluu joukkoon $\{1, 2, \dots, N\}$.

Esimerkki: Alla olevassa taulukossa on kolmen geenin ekspressioarvot kyselyn datamatriisille \mathbf{Q} , jossa on neljä sirua, tietokantamatriisille $\mathbf{D1}$, jossa on myös neljä sirua ja tietokantamatriisille $\mathbf{D2}$, jossa on vain kaksi sirua. Ensimmäisellä kyselysirulla s1 on mitattu aineella A käsiteltyjen solujen ekspression suhdetta käsittelemättömien solujen ekspression. Toisella s2, kolmannella s3 ja neljännellä s4 sirulla vastaavasti aineilla B, C ja D käsiteltyjen solujen ekspression suhdetta kontrolliin.

| | s1 | s2 | s3 | s4 | s5 | s6 | s7 | s8 | s9 | s10 |
|--------------|----|----|----|----|----|----|----|----|----|-----|
| mittaussarja | Q | Q | Q | Q | D1 | D1 | D1 | D1 | D2 | D2 |
| geeni1 | -1 | -2 | -1 | -2 | -4 | -3 | -4 | -3 | -4 | -4 |
| geeni2 | 2 | 6 | 5 | 1 | 3 | 2 | 4 | 1 | 5 | 2 |
| geeni3 | 5 | 1 | 2 | 3 | 2 | 4 | 1 | 2 | 2 | 3 |

Kyselyn geeniprofili geenille 1 on $(-1, -2, -1, -2)$, jossa $q_{11} = -1$, $q_{12} = -2$, $q_{13} = -1$ ja $q_{14} = -2$. Siruprofili sirulle s1 on $(-1, 2, 5)$, jossa $q_{11} = -1$, $q_{21} = 2$, $q_{31} = 5$. Geeni 1 on alaspäin ekspressoitunut, sen ekspressio on vähäisempää näytteessä kuin kontrollinäytteessä. Muut geenit ovat ylöspäin ekspressoituneita.

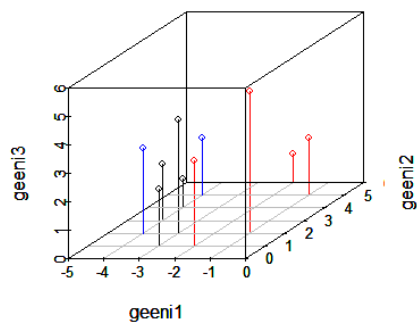
Kuvassa¹ 1 on (a) geenien suhteen kolmiulotteisessa avaruudessa kuvattu esimerkin kyselymittaussarja \mathbf{Q} (neljä mittausta, punainen) ja kaksi tietokannan mittaussarjaa (\mathbf{D}_1 : neljä mittausta, musta; \mathbf{D}_2 : kaksi mittausta, sininen) sekä (b) geeniekspresioarvojen säännöllisiä hahmoja havainnollistava esimerkki. Jos kuvan (a) kolmen datapistejoukon kuvan kääntäisi niin, että pisteet esitettäisiin mittausavaruudessa (eri sirut avaruuden suuntina), datapisteitä olisi jokaisella mittaussarjalla kolme, yksi jokaista geeniä kohden. Avaruudet olisivat eriulotteiset (mutta tehollinen avaruus olisi kolmiulotteinen, koska genejä on vain kolme: kolme vektoria, tässä datapistettä, voi virittää vain kolmiulotteisen avaruuden). Silloin eriulotteisten avaruuksien kohdentaminen toistensa suhteen ei ole triviaalia. SVD-matriisihajotelma (Singular Value Decomposition, SVD) tarjoaa linjaus- ja skaalausoperaatiot (ts. kääntö ja venytys) niin, että rivi- ja sarakeavaruuksien linjaus onnistuu hallitusti. Kun data-matriisilla operoidaan vektoriin, se voidaan SVD:n avulla esittää linjaus-skaalaus-linjaus-operaationa (tarkemmin seuraavassa luvussa). Seuraavassa kuvassa 2 ovat esimerkin kaikki siruprofiilit.

3 Etäisyysmitat ja moniulotteisuus

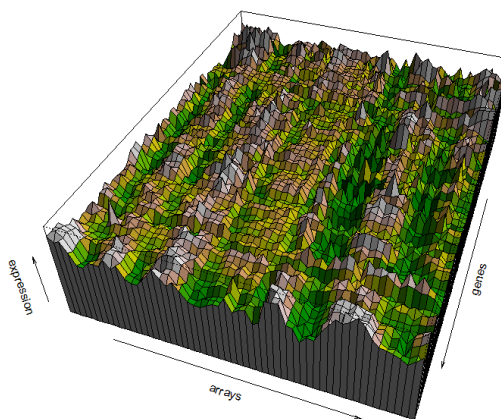
Alla esitellään etäisyysmitat, joita on käytetty luvun 4 hahmonhakualgoritmeissa. Lisäksi on muutama huomio ekspressidatan ominaisuuksista.

¹Kuvat on tuotettu R-ohjelmistolla [IG96]

Kysely: punainen, D1: musta, D2: sininen



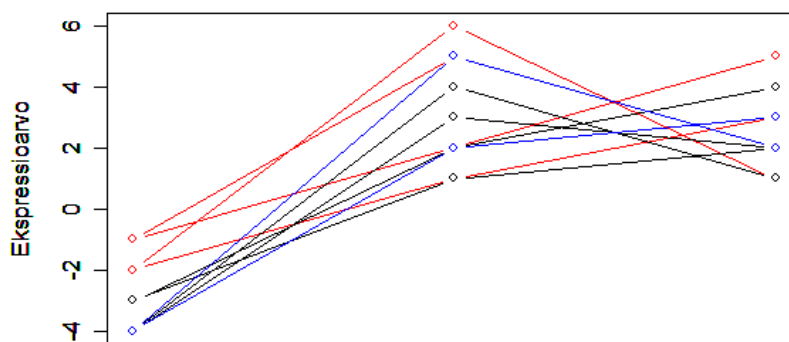
(a) Kysely ja kaksi tietokantadatapistejoukkoa



(b) Osa yhden mittausarja datasta

Kuva 1: (a) Vastaus kysymykseen, kumpi datapistejoukko vastaa kyselyä, ei ole triviaalia näinkään yksinkertaisessa esimerkissä. (b) Spellmanin ja kumppaneiden [SSZ98] solukiertotutkimuksissa käytetty ekspressiodata osoitti solukiertoon liittyvää säännöllistä geeniekspressiovaihtelua. Kuvassa on 70 geenin ekspressioarvojen suuruudet yhtenä, geenien identiteetti toisena ja siruidentiteetti kolmantena suunnana.

Esimerkkitaulukon siruprofiilit



Geenit vasemmalta oikealle: geeni1, geeni2 ja geeni3
Kysely: punainen, D1: musta, D2: sininen

Kuva 2: Siruprofiili kuvaa yhdellä sirulla olevien geenien ekspressioarvojen muutosta. Jokaisella sirulla on mitattu kolmen geenin ekspressiotaso. Värit vastaavat kutakin mittausarjaa. Sinisellä sarjalla on vain kaksi sirua. Kahdella muulla on neljä.

3.1 Euklidinen etäisyys ja korrelaatiomitat

Euklidinen etäisyys on geometrinen etäisyysmitta ja Minkovskin etäisyyden erikoistapaus. Euklidisessa etäisyydessä yksittäisen koordinaatin suuri numeroarvo voi hallita etäisyyden numeroarvoa. Etäisyys voi kasvaa miten suureksi tahansa. Jos mitaukset ovat tarkkoja ja poikkeavia havaintoja (outlier) ei ole, Euklidinen etäisyys mittaa hyvin mittausten vastaavuutta. Sitä käytetään monissa geeniekspression analyysissä mittausten epätarkkuudesta huolimatta.

Pearsonin korrelaatiokerroin on numeerinen mitta satunnaismuuttujien väliselle lineaariselle riippuvuudelle. Arvo on $-1:n$ (vastakkaissuuntainen korrelaatio) ja $1:n$ (samansuuntainen korrelaatio) välillä, 0 merkitsee, että muuttujat eivät korreloi. Geometrisesti tulkittuna korrelaatio on verrannollinen vektoreiden välisen kulman kosiniin. Pearsonin korrelaatio on kohinaherkkä. Pearsonin korrelaatiota käytetään paljon geeniekspression klusteroinneissa.

Spearmanin ρ on ei-parametrinen järjestyskorrelaatio, joka on Pearsonin korrelaatio siten, että havaintoarvojen (esimerkiksi ekspressioarvojen) suuruuden järjestysnumeroita käytetään korrelaatiolaskennan arvoina. Spearmanin järjestyskorrelaatiokerroin, ρ , mittaa kahden muuttujan havaintoarvojen suuruusjärjestyksien yhteensopivuutta. Se sopii esimerkiksi tilanteeseen, jossa ekspressiomittaukset on tehty eri teknologialla.

Laskentakaavoissa vektoreina ovat $\mathbf{x} = (x_1, x_2, \dots, x_m)$ ja $\mathbf{y} = (y_1, y_2, \dots, y_m)$.

Kaava 1 *Vektoreiden Euklidinen etäisyys, tason pisteiden etäisyyden yleistys:*

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2}$$

Kaava 2 *Vektoreiden Pearsonin korrelaatiokerroin, kun s_x, s_y ovat $\mathbf{x}:n$ ja $\mathbf{y}:n$ koordinaattien keskihajonnat (viiva merkitsee keskiarvoa):*

$$\text{cor}(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{m s_x s_y}$$

Kaava 3 *Vektoreiden Spearmanin ρ , kun RANK merkitsee vektorin koordinaattiarvon suuruuden järjestysnumeroa (1, 2 jne.):*

$$RE_i = \text{RANK}(x_i) - \text{RANK}(y_i)$$

$$\rho(\mathbf{x}, \mathbf{y}) = 1 - \frac{6 \sum_{i=1}^m RE_i^2}{m(m^2 - 1)}$$

Esimerkkitaulukossa: Pienin Euclidinen etäisyys (tietokantasiruille) on kyselys-
 irun s_1 ja s_6 välillä: $d(\mathbf{s1}, \mathbf{s6}) = \sqrt{(-1 - (-3))^2 + (2 - 2)^2 + (5 - 4)^2} = \sqrt{5}$. Paras
 Pearsonin korrelaatio s_1 :lle on myös s_6 :lla: $cor(\mathbf{s1}, \mathbf{s6}) = \frac{\sum_{i=1}^3 (s_{1i}-2)(s_{6i}-1)}{3*3*3.6} = 0.97$.
 Paras Spearmanin korrelaatio s_1 :lle tietokantariruista on kolmella siruilla s_6 , s_8 ja
 s_{10} , kaikilla korrelaatio on yksi ja ekspressiotason suuruusjärjestys on kaikille
 $(geeni1^{RANK}, geeni2^{RANK}, geeni3^{RANK}) = (1, 2, 3)$.

3.2 Matriisin singulaariarvohajotelma, SVD

SVD:tä on käytetty mm. tiedonhaussa (information retrieval) ja kuvahahmojen et-
 sinnässä sekä erilaisissa signaalintunnistustehtävissä. Myös PageRank [BrP98] ja
 HITS-algoritmit [Kle99] hakevat ominaisvektoreita (ks. alla). Pääkomponenttiana-
 lyysillä (Principal Component Analysis, PCA) visualisoidaan myös ekspressiodataa.
 PCA on SVD keskitetyille datamatriisille. Alter ja kumppanit [ABB00] ovat käyt-
 täneet menestyksekkäästi SVD:tä solujen eri tilojen identifioimiseen ja luokitelleet
 geenejä esimerkiksi sen mukaan, kuuluko niiden toiminta oleellisesti hiivan solu-
 kierron tiettyihin vaiheisiin. He ovat laajentaneet menetelmän käyttöä yleistetyllä
 SVD:llä niin, että eri lajienkin geeniekspressiota on voitu verrata.

SVD [PFT86, GoV83, Eld07] esittää matriisin \mathbf{D} , joka kuuluu joukkoon $\mathbb{R}^{m \times n}$, ha-
 joitettuna kolmeen osaan: $\mathbf{D} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ siten, että matriisin \mathbf{U} sarakkeet ovat mat-
 riisin $\mathbf{D} \mathbf{D}^T$ ominaisvektoreita ja matriisin \mathbf{V} sarakkeet matriisin $\mathbf{D}^T \mathbf{D}$ ominais-
 vektoreita ja $\mathbf{\Sigma}$:n diagonaalien, $(\sigma_1, \sigma_2, \dots, \sigma_n)$, alkiot ovat matriisin $\mathbf{D} \mathbf{D}^T$ kuten
 myös matriisin $\mathbf{D}^T \mathbf{D}$ ominaisarvojen neliöjuuria, eli singulaariarvoja. Singulaariar-
 vot ovat diagonaalilla järjestyksessä suurimmasta pienimpään. Ominaisvektorit (gee-
 niekspression yhteydessä: ominaisgeenit, \mathbf{V} , ja ominaissirut, \mathbf{U}) ja singulaariarvot
 vastaavat kolmikkona toisiaan niin, että matriisin ominaisgeenit kuvautuvat (pro-
 jisoituvat) mittausmatriisin kuvauksessa skaalatuille ominaissiruille ja päinvastoin:
 $\mathbf{D} \mathbf{v}_i = \sigma_i \mathbf{u}_i$ ja $\mathbf{D}^T \mathbf{u}_i = \sigma_i \mathbf{v}_i$. Matriisi koostuu ominaisgeeneistään ja -siruistaan:
 $\mathbf{D} = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^T$, ominaisgeenien ja -sirujen alkiot sekä singulaariarvot antavat ker-
 toimet lineaarikombinaatioon, jolla matriisin voi muodostaa. Singulaariarvot kerto-
 vat summassa olevien osamatriisien normin ja ominaisvektorit rivi- ja sarakeavaruus-
 den ortonormaalit kantavektorit, kukin summan jäsen on rank-1 matriisi.

Ominaisvektoreiden alkioiden lukuarvot kertovat, mitkä geenit ja sirut ovat voi-
 makkaimmin vaikuttaneet ko. lineaarikombinaation syntymiseen. Singulaariarvojen
 suhteet kuvaavat eri ulottuvuussuuntien merkitystä. Singulaariarvot pienenevät no-
 peasti, jos matriisissa on paljon lineaarista riippuvuutta. Jos kahden mittausarjan

(muutammat) ensimmäiset ominaissirut (\sim pääkomponentit) ja singulaariarvot vastaavat hyvin toisiaan ja singulaariarvot pienenevät nopeasti, mittaussarjojen voisi olettaa kuvaavan (jossain määrin) solujen samojen tilojen joukkoa.

Kyselymittaus voidaan projisoida tietokannan mittaussarjaan ominaissirujen avulla niin, että kyselysirujen lyhimmat etäisyydet tietokannan mittaussarjaan voidaan määrätä (pienimmän neliösumman sovitus, Least Square, LS), jos ekspressioarvojen oletetaan olevan lineaarisesti riippuvaisia. Projektio tehdään pseudokäänteismatriisin $\mathbf{D}^{-1} = \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^T$, avulla, missä $\mathbf{\Sigma}^+$ on singulaariarvojen käänteisluvuista muodostettu diagonaalimatriisi. Vektori $\hat{\mathbf{q}} = \mathbf{U}_k^T \mathbf{q}_j$ on projektio, joka linjaa kyselysirun tietokantamittaussarjan ominaissiruille. Operaatio $\mathbf{V}\mathbf{\Sigma}^+\hat{\mathbf{q}}$ skaalaa ja linjaa em. vektorin tietokantamittaussarjan ominaisgeeneille ja tuloksena saadaan vektori \mathbf{fitX} . Tämä projisoidaan edelleen tietokantamittaussarjan mittausavaruuteen. Yhden kyselysirun etäisyys mittaussarjasta on edellä kuvatun projektiosovituksen ja ko. sirun etäisyys eli jäännöksen normi, joka merkitään $\|\mathbf{D}_k \mathbf{fitX} - \mathbf{q}_j\|_2$. Koko kyselymittausarjan etäisyys määritellään em. normien summaksi: $\sum_{j=1}^n \|\mathbf{D}_k \mathbf{fitX} - \mathbf{q}_j\|_2$.

Kokeellisen datan laatu voi olla erittäin vaikea havaita etukäteen. Pienimmän neliösumman ongelmat voivat olla sekä ylimäärätyjä (overdetermined, datapisteitä on enemmän kuin parametrejä) että alimäärätyjä (underdetermined, datassa esiintyy monitulkintaisia parametrikombinaatioita). SVD:llä LS-sovitusta laskettaessa tätä ei tarvitse etukäteen tietää [PFT86]. Myös SVD ovat herkkä poikkeaville havainnoille, siksi ne on poistettu tietokannasta.

3.3 Moniulotteisuus ja Simpsonin paradoksi

Euklidisen avaruuden tilavuus kasvaa eksponentiaalisesti, kun uusia ulottuvuuksia lisätään². Moniulotteisen avaruuden tilavuus on keskittynyt avaruuden laiduille. Avaruus on enimmäkseen tyhjä, mikä merkitsee sitä että datan sisäinen rakenne on yleensä pienempiulotteinen. Siksi luokittelualgoritmeissa datavektorit voidaan usein projisoida aliavaruuteen ja luokitella siellä. Parhaan aliavaruuden löytäminen riippuu louhintatehtävästä. Hakuavaruuden kokoa rajoitetaan luvun 4 hahmonhakualgoritmeissa rajaamalla haku vain kyselyn sisältämiin geeneihin. Ekspressiodata on kuitenkin melko tasaisesti jakautunutta origon ympärille, joten vaikka geenejä otettaisiin hakuun tuhansia, hakutulokset eivät yleensä vääristy liikaa. Tämä voi olla

²Tämä on dynaamisen ohjelmoinnin keksijän Richard Bellmanin nimeämä 'moniulotteisuuden kirous'

hiivan datan ominaisuus, jonka olemassa olo mm. ihmisen datan kohdalla pitää testata erikseen.

Yksi mikrosirudatan analyysituloksia arvioitaessa huomioon otettava asia on riippuvuuksien ja riippumattomuuksien hienovaraisuuksien huomioiminen. Tämä liittyy 'Simpsonin paradoksina' tunnettuun tilastolliseen ilmiöön. Seuraava *esimerkki* [HMS01] tuntuu ristiriitaiselta. Lääke A parantaa vanhat, suhteessa 30 parantunutta 90 lääkkeen käyttäjästä (30/90), ja nuoret, suhteessa kymmenen parantunutta kymmenestä lääkkeen käyttäjästä (10/10), paremmin kuin lääke B, jolle vastaavat luvut ovat 2/10 vanhoille ja 48/90 nuorille. Kuitenkin, kun ryhmät yhdistetään ja katsotaan kuinka tilastollisesti on käynyt, lääke B (, jolle vastaavat luvut ovat 50/100) parantaa yhdistetyssä ryhmässä paremmin kuin lääke A (, jolle vastaavat luvut ovat 40/100). Ryhmien koon erilaisuus ja parantumismittausten suhteellisuus aikaansaa paradoksin. Koska solunäytteessä mitatut RNA-suhteet voivat olla peräisin (esimerkiksi) kahdesta eri solupopulaatiosta ja populaatioiden kokoero ei ole tiedossa, saattaa mittauksen näyttää siltä, että joku geeni säätelee toista geeniä alentamalla ekspressiota, vaikka itse asiassa molemmissa solupopulaatioissa erikseen geeni säätelee toista geeniä ekspressiota lisäämällä³. Myös (yllä olevan esimerkin tapaan) solunäytteiden solupopulaatiot ovat eri kokoisia ja mittaukset ovat suhteellisia. Tämän tyyppiset rajoitukset ekspressiomittauksissa edellyttävät hahmonetsintätulosten varmentamista muulla tavoin ja biologisen asiantuntemuksen käyttöä tulosten tulkinnassa.

4 Hakualgoritmit

Haku pyrkii löytämään jotakuinkin samansuuntaiset ja normiltaan samanlaiset vektorijoukot. Tätä varten määritellään yleiseen tapaukseen heuristinen laatuluku L , jolla samanlaisuutta arvioidaan. Hakuongelman ratkaisut käyn läpi yksinkertaisemmasta tilanteesta hankalampaan. Yksiulotteisessa tapauksessa käytän samanlaisuuden kriteerinä Euklidista etäisyyttä tai korrelaatiomittoja. Yleisen tapauksen koostan yksiulotteisista tapauksista em. laatuluvun avulla. SVD ratkaisee ulottuvuusongelmat geenien ja sirujen suuntaan samanaikaisesti. SVD-LS haussa lasken LS-sovituksen jäännöksen (residual) normin kullekin kyselysirulle erikseen (jokaista tietokantamittaussarjaa kohden) ja summaan sen jokaiselle tietokantamittaussarjalle

³Tommi Jaakkola kurssilla: Intensive course on modeling biological networks, Otaniemi elokuu 2007

erikseen. Pienin summa voittaa.

Haku tehdään vain hakugeenin, \mathbf{G}_M , sisältävään osatietokantaan, $\mathbf{DB}_M \in \mathbb{R}^{M \times p}$, missä p on kaikkien sirujen lukumäärä tietokannassa. Geenisuuntaista profiilia ei etsitä; se tulee implisiittisesti mukaan, jos siruprofiilit ovat riittävän samanlaisia. Kunkin mittaussarjan kyselysiruja lähinnä olevat mittaukset määräävät mittaussarjan vastaavuuden kyselyn kanssa. Kaukana olevat mittaukset eivät huononna tulosta, jos lähellä olevia löytyy.

Selityksissä käytän merkintöjä: (1) *SORT* merkitsee järjestämistä, pienimmästä suurimpaan. (2) *SORT^{desc}* merkitsee järjestämistä, suurimmasta pienimpään. (3) *RANK* ja *RANK^{rev}* merkitsevät järjestysnumeron antamista arvon suuruuden perusteella - *rev* merkitsee, että suurin arvo saa pienimmän järjestysnumeron. (4) \bar{x} on lukujen x_i keskiarvo. (5) K on tietokannassa \mathbf{DB} olevien mittaussarjojen \mathbf{D}_k lukumäärä ja (6) p on kaikkien tietokantasirujen lukumäärä.

4.1 Yksiulotteinen tapaus

Kyselysyötteenä on yksi vektori, jos kyselyssä ei ole kuin yksi mittausta: $\mathbf{Q} = (q_{11}, \dots, q_{M1})$. Tietokannan mittaussarja on (mallissa) samanlainen kuin kysely, jos jonkin mittaussarjasirun Euklidinen etäisyys kyselysirusta on lähellä nollaa geenivaruudessa mitattuna tai jos kyselyn siruprofiili korreloi hyvin mittaussarjan jonkun siruprofiilin kanssa. Ensimmäisenä on triviaali tapaus. Jos vain yhden geenin ekspressiota on mitattu, vektori on itse asiassa skalaari ja vertailussa käytetään vain arvojen erotusta. Seuraavaksi oletetaan, että geenejä on useampia.

Edellisen luvun etäisyyksiä käyttäen voi mittausten vastaavuuden määrätä kolmella eri kriteerillä. (1) Kuinka kaukana ekspressioarvot ovat toisistaan? (2) Kuinka hyvin mittaukset korreloivat, kun negatiivinen korrelaatio katsotaan huonommaksi kuin ei korrelaatiota ollenkaan? (3) Kuinka hyvin eri geenien ekspression keskinäinen vaihtelu mittaustilanteessa korreloi, kun vain ekspressiotason suuruuden järjestys geenien kesken merkitsee? Näistä saamme kolme yksinkertaista algoritmia, Algoritmi 1, Algoritmi 2 ja Algoritmi 3. Jokainen tietokannan mittaussarja, $k \in \{1, 2, \dots, K\}$, saa järjestysnumeron perustuen siihen, kuinka etäällä lyhimmillä etäisyydellä kyselystä oleva siru on sijoittunut tai kuinka hyvin mittaussarjan parhaiten korreloiva siru on sijoittunut.

Kullekin algoritmille tehdään sama **alustus**: Erotta tietokannasta \mathbf{DB} geenejä $\{1, \dots, M\}$ vastaavat rivit. Osatietokanta, johon kysely kohdistuu, on nyt $\mathbf{DB}_M^{M \times p}$,

missä p on tietokannan sirujen lukumäärä. Yksittäinen mittausarvo osatietokannassa on g_{ij} , missä $i \in \{1, 2, \dots, M\}$ ja $j \in \{1, 2, \dots, p\}$.

Algoritmi 1 Ehaku($\mathbf{Q} = \mathbf{q}_1 = (q_{11}, \dots, q_{M1}), \mathbf{DB}_M$)

Hakualgoritmi yksiulotteiseen hahmon hakuun Euklidisen etäisyyden avulla.

1. **Laske** $d_j \leftarrow d(\mathbf{q}_1, \mathbf{g}_j) = \sqrt{\sum_{i=1}^M (q_{i1} - g_{ij})^2}, \forall j \in \{1, 2, \dots, p\}$
2. **Laske** $d_k \leftarrow \min_{\mathbf{g}_j \in \mathbf{D}_k} (d_j), \forall k \in \{1, 2, \dots, K\}$
3. **return** $SORT(d_k); k \in \{1, 2, \dots, K\}$

Algoritmi 2 Phaku($\mathbf{Q} = \mathbf{q}_1 = (q_{11}, \dots, q_{M1}), \mathbf{DB}_M$)

Hakualgoritmi yksiulotteiseen hahmon hakuun Pearsonin korrelaation avulla

1. **Laske** $cor_j \leftarrow cor(\mathbf{q}_1, \mathbf{g}_j) = \frac{\sum_{i=1}^M (q_{i1} - \bar{q})(g_{ij} - \bar{g}_j)}{Ms_{q_1} s_{g_j}}, \forall j \in \{1, 2, \dots, p\}$
2. **Laske** $cor_k \leftarrow \max_{\mathbf{g}_j \in \mathbf{D}_k} (cor_j), \forall k \in \{1, 2, \dots, K\}$
3. **return** $SORT^{desc}(cor_k); k \in \{1, 2, \dots, K\}$

Algoritmi 3 Shaku($\mathbf{Q} = \mathbf{q}_1 = (q_{11}, \dots, q_{M1}), \mathbf{DB}_M$)

Hakualgoritmi yksiulotteiseen hahmon hakuun Spearmanin korrelaation avulla

1. **Laske** $\rho_j \leftarrow \rho(\mathbf{q}_1, \mathbf{g}_j) = 1 - \frac{6 \sum_{i=1}^M RE_i^2}{M(M^2 - 1)}, \forall j \in \{1, 2, \dots, p\}$
2. **Laske** $\rho_k \leftarrow \max_{\mathbf{g}_j \in \mathbf{D}_k} (\rho_j), \forall k \in \{1, 2, \dots, K\}$
3. **return** $SORT^{desc}(\rho_k); k \in \{1, 2, \dots, K\}$

4.2 Yleinen tapaus: heuristinen ratkaisu

Yleisessä hakutapauksessa mittausarvojen vastaavuus määritellään heuristisen laatuluvun avulla. Parhaat mittausarvojen vastaavuudet haetaan ensin joka kyselysirulle erikseen ja tietokannan sirut saavat laatuluvun kyselysirun ja tietokantasirun parittaisen etäisyysmitta-arvon suuruuden järjestysnumerona (pienin etäisyys voittaa). Tämän voin laskea yksiulotteisen tapauksen kolmea eri kriteeriä käyttäen. Jokaisen tietokantakokeen k lopulliseen laatulukuun L_k summataan *kutakin kyselysirua* parhaiten vastaavan sirun laatuluku tässä mittausarjassa (eli N kpl). Laatulukujen summan arvo määrää mittausarjan sijoituksen vastauksessa. Kyselysyytteenä

on matriisi, jos kyselyssä on monta mittausta: $\mathbf{Q}^{M \times N}$, joka koostuu vektoreista q_{ij} siten, että i kuuluu joukkoon $\{1, \dots, M\}$ ja j kuuluu joukkoon $\{1, \dots, N\}$, M on geenien lukumäärä kyselyssä ja N sirujen. Vastaavuudessa ei tietokantamittausten mittaussarjakohtainen lukumäärä vaikuta hakuun, vain kyselymittausta (kutakin mittaussarjaa kohden) lähinnä olevat mittaussarjan mittaukset määräävät samanlaisuuden. Jos kyselysirulle löytyy mittaussarjassa hyvin samanlainen siru, se riittää. Silloin mittaussarjassa on tehty ainakin yksi mittaus, jossa samanlaisuus toteutuu. Vaikka muita mittauksia on lisäksi tehty, siitä ei rankaista. Tietokannan yksittäisellä sirulla voi olla paras vastaavuus moneen kysely-siruun nähden. Näin käy välttämättä, jos tietokannan mittaussarjassa ei ole ollut kuin yksi siru. Algoritmissa 4 käytetään yhtä yksiulotteisen tapauksen hakutavoista **Ehaku**, **Phaku** tai **Shaku** samassa haussa. Tietokannan alustus osatietokannaksi tehdään kuten yksiulotteisessa tapauksessa.

Algoritmi 4 **Lhaku**($\mathbf{Q}^{M \times N}$, \mathbf{DB}_M)

Hahmon haku laatuluvun avulla

1. Alusta $L_k \leftarrow 0, \forall k \in \{1, 2, \dots, K\}$

2. **for** i *in* $\{1, 2, \dots, N\}$

$$L_i^{1\text{-ulotteinen}} \leftarrow \begin{cases} \text{RANK}(\mathbf{Ehaku}(q_i, \mathbf{DB}_M))\text{tai} \\ \text{RANK}^{rev}(\mathbf{Phaku}(q_i, \mathbf{DB}_M))\text{tai} \\ \text{RANK}^{rev}(\mathbf{Shaku}(q_i, \mathbf{DB}_M)) \end{cases}$$

$$L_k \leftarrow L_k + L_{ki}^{1\text{-ulotteinen}}, \forall k \in \{1, 2, \dots, K\}$$

3. **return** $\text{SORT}(L_k)$; $k \in \{1, 2, \dots, K\}$

4.3 Pienimmän neliösumman sovitus SVD:llä

Pienimmän neliösumman sovitus SVD:llä lasketaan luvussa 3 kuvatulla tavalla:

Jokaiselle kyselysirulle lasketaan pienimmän neliösumman jäännöksen (residual) normi osatietokantamittaussarjojen pseudokäänteismatriisin $\mathbf{D}^{-1} = \mathbf{V}\Sigma^+\mathbf{U}^T$ avulla (algoritmissa oleva tietokantaindeksi k identifioi tietokantamittaussarjan). Jos kyselysiru \mathbf{q} olisi (singulaarisen) \mathbf{D} :n rangilla, projektiolla $\mathbf{DfitX} = \mathbf{q}$ on ääretön määrä ratkaisuja. Jos taas kyselysiru ei olisi (singulaarisen) \mathbf{D} :n rangilla, ratkaisuja ei olisi ollenkaan. Kummassakin tapauksessa lähin mahdollinen (pienimmän

neliösumman mielessä) \mathbf{fitX} löytyy SVD:n avulla, mikä minimoi jäännöksen normin $\|\mathbf{DfitX} - \mathbf{q}\|_2$ (tässä 2-normi) [PFT86, GoV83]. Algoritmikuvassa 5 on tämä esitetty pseudokoodina. Alustus sisältää taas hakugeenien osatietokantaan rajoittumisen, kuten muissakin algoritmeissa. Itse SVD-hajotelma lasketaan esimerkiksi R-ohjelmistolla. Rutiini SVD:n laskuun löytyy useista ohjelmistoista. SVD:n voi laskea $O(mn^2)$ liukulukuoperaatiolla.

Algoritmi 5 *Hahmon haku SVD-LS-menetelmällä*

```

1. for  $k$  in  $\{1, \dots, K\}$ 
     $[\mathbf{U}_k, \mathbf{\Sigma}_k, \mathbf{V}_k] \leftarrow SVD(\mathbf{D}_k)$ 
    Alusta  $Residual_k \leftarrow 0$ 
    for  $j$  in  $\{1, 2, \dots, N\}$ 
         $\hat{\mathbf{q}} \leftarrow \mathbf{U}_k^T \mathbf{q}_j$  ; kyselysirun projektio ominaissirukannassa
         $\mathbf{fitX} \leftarrow \mathbf{V}_k \mathbf{\Sigma}^+ \hat{\mathbf{q}}$  ; sovitus
         $Residual_k \leftarrow Residual_k + \|\mathbf{D}_k \mathbf{fitX} - \mathbf{q}_j\|_2$  ; normin summaus
2. return  $SORT(Residual_k)$ ;  $k \in \{1, 2, \dots, K\}$ 

```

Esimerkkitaulukossa: Euklidisen etäisyyden avulla laskettaessa \mathbf{D}_1 on niukasti lähempänä kyselyn mittausarjaa verrattuna \mathbf{D}_2 :een. Siruja $s1$ ja $s4$ lähinnä ovat \mathbf{D}_1 :n sirut $s6$ ja $s8$. Sekä siruja $s2$ ja $s3$ lähinnä on \mathbf{D}_2 :n siru $s9$. Vaa'ankielenä on \mathbf{D}_2 :n siru $s10$ - se ei ole kuin vasta kolmanneksi lähinnä siruun $s4$ nähden. Jos se olisi ollut toisena, tulos olisi ollut 'tasapeli', koska \mathbf{D}_1 :n sirut olivat toisena molemmissa 'ei-ykkösenä'-tapauksissa. Pearsonin korrelaation avulla laskettaessa \mathbf{D}_1 voittaa selvästi. Sen siru $s6$ korreloi parhaiten sirujen $s1$ ja $s4$ kanssa, vastaavasti $s7$ korreloi parhaiten sirujen $s2$ ja $s3$ kanssa. \mathbf{D}_2 :n siru $s9$ on toisena sirujen $s1$ ja $s4$ suhteen, siru $s10$ on vasta kolmantena sirujen $s2$ ja $s3$ suhteen. Siis kaksi on sekin \mathbf{D}_1 :llä. Spermanin korrelaatiolla laskettaessa tulos on tasapeli. Näin käy helposti, jos vain kolme geeniä on haussa. SVD-LS menetelmällä \mathbf{D}_1 voittaa selvästi. Sille kaikkien kyselysirujen jäännöksen normi on käytännössä nolla. \mathbf{D}_2 :lle kyselysirun $s2$ jäännöksen normi on alle yksi, muille yli yksi. Tarkemmin: $\|\mathbf{D}_2x - s1\|_2 = 3.40$, $\|\mathbf{D}_2x - s2\|_2 = 0.84$, $\|\mathbf{D}_2x - s3\|_2 = 1.97$ ja $\|\mathbf{D}_2x - s4\|_2 = 1.07$.

5 Hakutulosten vertailu ja loppupäätelmät

Menetelmieni hakutuloksia olen verrannut (yksittäistulosten ja biologisen relevanssiarvion lisäksi) erillisellä ohjelmalla. Ohjelma vertaa hakutulosten mittausarjojen

tulosparemmuusjärjestyksen alku- ja loppupääleikkausten sisältämien alkiodien lukumäärää. Näin voi arvioida, löytyykö hakutuloksille yhteistä konsensusta paremmuusjärjestykseen. Vertailu koskee hiivatietokannan, jossa on 31 mittaussarjaa, 479 sirua ja 5939 geeniä, hakutuloksia. Pearsonin ja Spearmanin korrelaatiot tuottavat hyvin samanlaisia tulosjärjestyksiä. Seuravaksi eniten samankaltaisuutta on SVD-LS- ja Pearsonin korrelaatio -menetelmillä. Euklidinen haku on tuottaa samanlaisia tulosjärjestyksiä Pearsonin korrelaation kanssa, mutta harvemmin kuin kaksi edellistä paria. Kaikkien eri menetelmien yhteiskonsensus löytyy myös usein. Pienet vaihtelut paremmuusjärjestyksessä tasottuvat leikkausvertailulla, Spearmanin ρ ja Kendallin τ ovat kohinaherkempiä konsensusarvioinnissa kuin em. leikkauksen käyttö.

Euklidisen etäisyyden ja laatuluvun käyttö näyttää löytävän mittaussarjojen samanlaisuutta, jos poikkeavia havaintoja ei tietokannassa ole ja haku sisältää riittävän paljon lukuarvoltaan erilaisesti ekspressoituneita genejä. Euklidisessa haussa merkittävästi ekspressoituneet geenit yleensä määräävät samanlaisuuden eri mittausten välille.

Pearsonin korrelaatiolla voi hakea lineaarista riippuvuutta kyselyn ja tietokannan mittausten välille. Lukuarvojen suuruudella ei silloin ole väliä, kunhan arvovaihtelu on ollut samansuuntaista. Jos kyselyllä tai tietokannan mittaussarjalla on tasaista numeroarvojen vääristymää tai mittausten kontrollitaso on eri, Pearsonin korrelaatio voi löytää samantapaiset kokeet, joita Euklidisella etäisyydellä tehty haku ei löydä. Pearsonin korrelaatio on herkkä datassa esiintyvälle kohinalle ja lineaarisuuden löytyminen on osittain sattumanvaraista, jos datassa lähes kaikki arvot ovat lähellä nolla-tasoa olevia arvoja. Euklidista etäisyysmittaa käytettäessä pienillä lukuarvoilla ei tule suuria eroja, mutta Pearsonin korrelaatiota käytettäessä normiltaan pienten vektoreiden suunnat voivat vaihdella merkittävästi, vaikka normi itsessään olisi merkityksetön eroavaisuutta etsittäessä. Jos mittaukset olisivat tarkkoja, silloin suunnan vaihtelu on tietenkin merkitsevä eroja etsittäessä.

Spearmanin järjestyskorrelaation ongelma on esimerkiksi mittausten lukuarvojen tarkkuusraja ja geenijoukon koko. Läheisillä arvoilla olevien geenien järjestykseen voi tulla mittaustarkkuusrajan määräämää sattumanvaraisuutta. Liian pieni määrä genejä ei erottele mittauksia tarpeeksi. Järjestykseen perustuvan korrelaation hyöty tulee sen parametrattomuudesta. Mm. eri alustojen (mittaustekniikoita on erilaisia ja niiden tuottamat ekspressioarvot eivät suoraan ole vertailukelpoisia) mittauksia voidaan silloin vertailla keskenään.

SVD:hen perustuva menetelmä ei etsi suoraan kahden mittauksen välistä lineaarista riippuvuutta, vaan koko mittaussarjaan liittyvää lineaarisuutta, jolloin koko mittaussarjan hahmo erottuu. SVD:n etsimän lineaarisuuden merkitys ei ole selvää solubiologian epälineaarisissa prosesseissa - mikä on menetelmän huono puoli. Tietokannan esianalyysit kuitenkin osoittivat, että mittauservojen välillä on lähes lineaarista riippuvuutta suuressakin geenijoukossa (mittausarvojen lineaarisuutta testattaessa). Press ja kumppaneiden mukaan [PFT86] SVD:n käyttö pienimmän neliösumman parametrien arvioinnissa ja itse asiassa koko LS-menetelmä voidaan tulkita tilastolliseksi maksimaalisen todennäköisyyden arvioksi (maximum likelihood estimation) siten, että LS-menetelmällä löydetyt parametrit maksimoivat mitatun datan todennäköisyyden. Luulen, että yksi syy SVD:n tehoon niin monella sovellusallalla perustuu osittain siihen, että SVD:llä LS-laskenta voidaan tehdä mahdollisimman häiriöttömästi. SVD tasoittaa vertailuongelmaa, joka syntyy paljon ja vähän siruja sisältävien mittaussarjojen välille. Laatuluville ja Euklidisella etäisyydellä tai korrelaatioilla vertailtaessa monta sirua sisältävästä mittaussarjasta löytyy tilastollisesti helpommin (ehkä liiankin helposti) samanlainen siru kyselysirun kanssa kuin SVD:llä vertailtaessa (tämä koskee esimerkiksi aikasarjadataa). Laatuluville käyttöä voisi kokeilla niin, että lisäisi sille sopivan painotuksen, joka huonontaisi laatua, jos kyselyllä ja tietokantamittaussarjalla on paljon kokoeroa. Perustavanlainen kysymys on silloin, millaista vastaavuutta halutaan.

[PFT86] on hyvä lähde SVD:n soveltamisessa, jos SVD:n käyttö kiinnostaa lukijaa. Hakukoneen kehittämisessä voisi seuraavaksi kokeilla ICA:n (Independent Component Analysis) ja NMF:n (Nonnegative Matrix Factorization) käyttöä. ICA:n etu SVD:hen nähden on, että sillä voi analysoida tilastollista riippumattomuutta. NMF löytää tunnetusti paikallisia hahmoja.

Lähteet

- ABB00 Alter, O., Brown, P. O. ja Botstein, D., Singular value decomposition for genome-wide expression data processing and modelling. *Proc. Natl. Acad. Sci. USA*, 97,18(2000), sivut 10101–10106.
- ASG90 Altschul, S. F., Gish, W., Miller, W., Myers, E. W. ja Lipman, D. J., Basic local alignment search tool. *Journal of Molecular Biology*, 215, sivut 403–410.
- BrP98 Brin, S. ja Page, L., The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems*, 30, sivut 107–117.
- EiB99 Eisen, M. B. ja Brown, P. O., Dna arrays for analysis of gene expression. *Methods Enzymol*, 303, sivut 179–205.
- Eld07 Eldén, L., *Matrix Methods in Data Mining and Pattern Recognition*. SIAM, 2007.
- ESB98 Eisen, M. B., Spellman, P. T., Brown, P. O. ja Botstein, D., Cluster analysis and display of genome-wide expression patterns. *Proc. Natl. Acad. Sci. USA*, 95, sivut 14863–14868.
- FKT⁺07 Fujibuchi, W., Kiseleva, L., Taniguchi, T., Harada, H. ja Horton, P., Cellmontage: Similar expression profile search server. *Bioinformatics, on line version*. [<http://bioinformatics.oxfordjournals.org/cgi/content/short/23/22/3103>, on line version].
- GoV83 Golub, G. ja Van Loan, C., *Matrix computations*. The Johns Hopkins University Press, Oxford, 1983.
- HMJ00 Hughes, T., Marton, M., Jones, A., Roberts, C., Stoughton, R., Armour, C., Bennett, H., Coffey, E., Dai, H., He, Y., Kidd, M., King, A., Meyer, M., Slade, D., Lum, P., Stepaniants, S., Shoemaker, D., Gachotte, D., Chakraburttty, K., Simon, J., Bard, M. ja Friend, S., Functional discovery via a compendium of expression profiles. *Cell*, 102,1(2000), sivut 109–126.
- HMS01 Hand, D., Mannila, H. ja Smyth, P., *Principles of Data Mining*. The MIT Press, 2001.

- IG96 Ihaka, R. ja Gentleman, R., R: A language for data analysis and graphics. *Journal of Computational and Graphical Statistics*, 5,3(1996), sivut 299–314.
- Kle99 Kleinberg, J., Authoritative sources in a hyperlinked environment. *Journal of the ACM*, 46, sivut 604–632.
- Lam07 Lamb, J., The connectivity map: a new tool for biomedical research. *Genome Res.*, 7,1(2007), sivut 54 – 60.
- LCP06 Lamb, J., Crawford, E., Peck, D., Modell, J. W., Blat, I. C., Wrobel, M. J., Lerner, Jim Brunet, J.-P., Subramanian, A., Ross, K. N., Reich, M., Hieronymus, H., Wei, G., Armstrong, S. A., Haggarty, S. J., Clemons, P. A., Wei, R., Carr, S. A., Lander, E. S. ja Golub, T. R., The connectivity map: using gene-expression signatures to connect small molecules, genes, and disease. *Science*, 313,5795(2006), sivut 1929–1935.
- LDB96 Lockhart, D. J., Dong, H., Byrne, M. C., Follettie, M. T., Gallo, M. V., Chee, M. S., Mittmann, M., Wang, C., Kobayashi, M., Norton, H. ja Brown, E. L., Expression monitoring by hybridization to high-density oligonucleotide arrays. *Nat Biotech*, 14,12(1996), sivut 1675–1680.
- PFT86 Press, W., Flannery, B., Teukolsky, S. ja Vetterling, W., *Numerical recipes, the art of scientific computing*. Cambridge University Press, Cambridge, 1986.
- Sou75 Southern, E., Detection of specific sequences among dna fragments separated by gel electrophoresis. *J Mol Biol.*, 98, sivut 503–517.
- SSD95 Schena, M., Shalon, D., Davis, R. W. ja Brown, P. O., Quantitative monitoring of gene expression patterns with a complementary dna microarray. *Science*, 270,5235(1995), sivut 467–470.
- SSZ98 Spellman, P. T., Sherlock, G., Zhang, M. Q., Iyer, V. R., Anders, K., Eisen, M. B., Brown, P. O., Botstein, D. ja Futcher, B., Comprehensive identification of cell cycle–regulated genes of the yeast *saccharomyces cerevisiae* by microarray hybridization. *Molecular Biology of the Cell*, 9,12(1998), sivut 3273–3297. [Myös <http://yeastbase.csc.fi/>, 3.5.2007].

- STM05 Subramanian, A., Tamayo, P., Mootha, V. K., Mukherjee, S., Ebert, B. L., Gillette, M. A., Paulovich, A., Pomeroy, S. L., Golub, T. R., Lander, E. S. ja Mesirov, J. P., Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. USA*, 102,43(2005), sivut 15545–15550. [GSEA <http://www.broad.mit.edu/gsea/>, 27.7.2007].
- TuL05 Tuimala, J. ja Laine, M., toimittajat, *DNA Microarray Data Analysis*. Picaset Oy, Helsinki, 2005. [Myös <http://www.csc.fi/english/research/sciences/bioscience/books>].

Liite 1. Sanasto

Sanasto sisältää tämän kirjoituksen biologista käsitteisanastoa ja sellaisia asioita, joita saattaa tulla vastaan, kun kirjoituksen lähteisiin tutustuu. Biologiassa on paljon poikkeuksia ja muuttuvaa tietoa, joten sanojen selityksetkin muuttuvat ja riippuvat asiayhteydestä. Yleisimmät käsitteet on selitetty ensin.

- *Eliöt - eukaryootit ja prokaryootit*

Taksonomisesti eliöt voidaan jakaa kolmeen ryhmään: arkkeliöihin (arkit, *Archaea*), bakteereihin (*Bacteria*) ja aitotumallisiin (*Eukarya*). Solutyypin mukaan jaoteltuna arkit ja bakteerit kuuluvat prokaryootteihin (esitumallisiin), joilla ei ole solukalvon erottamaa tumaa, ja eukaryootteihin (aitotumallisiin), joilla on sekä tuma että muita kalvopäällysteisiä DNA:ta sisältäviä rakenteita (esimerkiksi: eläimillä mitokondriot ja kasveilla viherhiukkaset). Taksonomisena ryhmänä voidaan käyttää myös kliinejä (clade). Kliinin jäsenet jakavat piirteitä, jotka ne ovat perineet yhteiseltä esi-isältä. *Neomura*-kliiniin on ehdotettu kuuluvaksi sekä arkit että eukaryootit. Samoin on ehdotettu, että *Neomura* olisi kehittynyt bakteereista, tarkemmin nykyäänkin maaperässä runsaana esiintyvistä aktinobakteereista (*Actinobacteria*). Hiiva ja ihminen kuuluvat eukaryootteihin. Leivontahiiva (*Saccharomyces cerevisiae*) on molekyyli- ja solubiologiassa laajalti käytetty eukaryoottinen mallieliö (model organism), jolle on olemassa lajikohtaisia tietokantoja, mm. Standfordin yliopiston yhteydessä oleva *Saccharomyces* Genome Database, SGD. Viruksia ei lueta eliöiden taksonomiaan, koska ne ovat soluttomia ja niillä ei ole aineenvaihduntaa. Virukset käyttävät eliöitä isäntäsoluinaan. Muita soluttomia elämänmuotoja ovat esimerkiksi prionit.

- *DNA*

DNA on solun makromolekyyli-rakenteinen nukleinihappo, johon on koodattu eliöiden solujen ja joidenkin virusten geneettinen materiaali. Koodin tietoyksikköinä ovat neljä emästä, joiden järjestys on kaiken elollisen toiminnan ja kehityksen takana. Eliön lisääntyessä geneettinen materiaali kopioituu ja välitetään jälkeläisille. Eukaryooteissa DNA on järjestäytynyt kromosomeiksi solun tumassa. Prokaryooteissa DNA on rengasmaisena molekyylinä solun sisällä ilman tumaa nukleoidiksi kutsutulla alueella. Eläin- ja kasvisolujen mitokondriot sekä kasvisolujen kloroplastit (viherhiukkaset) sisältävät oman rengasmaisen DNA:nsa kuten bakteereissa.

- *RNA*

RNA on DNA:sta kopioituva ribonukleiinihappo, jota käytetään solun sisällä DNA-emästen järjestyksen tiedon siirtoon, kopiointiin ja jossain tapauksissa geneettisenä materiaalina. RNA:han kopioituu vain pieni osa DNA:ta ja sen rakennekin on hieman toisenlainen.

RNA-tyyppejä on kolme lähetti-RNA (messenger RNA, mRNA), ribosomi-RNA (ribosomal RNA, rRNA) ja siirtäjä-RNA (transfer RNA, tRNA). Virus-ten geneettinen materiaali voi olla RNA:ta.

- *Geeni* (Gene)

Geenillä on useita määritelmiä riippuen tarkkuudesta ja yhteydestä, jossa asiasta puhutaan. Yleensä geeni on sellainen osa DNA:ta, jota koodaamalla (transkriptiossa) saadaan muodostettua yksi solussa toimiva proteiini. Geenejä eliöillä on tuhansia. Esimerkiksi hiivalla n. 6000 ja ihmisellä yli 20 000 (ennen väitettiin, että 40 000), bakteereillakin useampi tuhat. Eukaryooteilla geenin koodaavien osien, eksonien, välissä on proteiinia koodaamattomia jaksoja, introneja. Intronien alussa ja lopussa on jaksoja, joiden avulla ne voidaan poistaa (silmukointi) ennen proteiinisynteesiä.

- *Genomi*

Genomi on organismin koko perintöainees, joka on koodattu DNA:han.

- *Sekvenssi*

Sekvenssi on DNA:ssa olevien emästen (A, C, G ja T), RNA:ssa olevien emästen, proteiinien aminohappojen (yl. 20 kirjainkoodia) tai muun biologisen koodausosasten merkkijonojen järjestys.

- *Nukleotidi*

Nukleotidi on monomeeri, joista nukleiinihapot DNA ja RNA muodostuvat. Nukleotidi koostuu fosfaatti-, sokeri- ja emäsosasta. RNA:n sokeriosana on riboosi ja DNA:n deoksiriboosi. DNA:n emäsosia ovat guaniini, adeniini, sytosiini ja tymiini. RNA:n emäsosat ovat guaniini, adeniini, sytosiini ja urasiili. Ihmisen DNA sisältää kolme miljardia nukleotidia. Viruksellakin niitä voi olla tuhansia.

- *Aminohappo*

Aminohappo on orgaaninen happo- ja emäspään (korkboksyyli- ja amino-ryhmä) sisältävä yhdiste, joita yhdistämällä saadaan muodostettua elämään tarvittavat proteiinit. Eliöt käyttävät kahtakymmentä (nykytiedon mukaan kaksi lisää - eräällä mikrobilla) erilaista aminohappoa tuhansien erilaisten proteiinin muodostamiseen. Kolmen emäksen järjestys DNA:ssa vastaa yhtä luonnon aminohappoa.

- *Proteiini* (Valkuaisaine)

Proteiinit ovat polymeroituneita aminohappoja (polypeptidejä). Yksi proteiini voi sisältää satoja ja jopa tuhansia aminohappo-osia. PDB:ssä (Protein Data Bank) on lähes 40 000 proteiinin rakenne. Proteiineilla on ensinnäkin primäärinen rakenne, joka kertoo aminohappojärjestyksen, josta se muodostuu. Aminohapot liittyvät toisiinsa peptidisidoksin. Sekundaarisella rakenteella viitataan muun muassa proteiinin muotoon vaikuttaviin rikkisiltoihin ja vetysidoksiin. Tertiäärinen rakenne tarkoittaa proteiinin avaruudellista rakennetta kokonaisuudessaan. Kvaternaarinen rakenne tarkoittaa usean aminohappoketjun ryhmittymää. Esim. hemoglobiini on neljän yhteen liittyneen proteiinin muodostama tetrameeri. Myös useimmat entsyymit ovat proteiineja.

- *Transkriptio*

Solun proteiinisynteesiprosessin eräs osa on transkriptio, jossa DNA:n koodin mukainen emäsjärjestys välittyy aminohappojärjestykseksi syntyvään proteiinimolekyylisiin (RNA:n muodostuminen DNA:sta). Transkriptiotekijät (transcription factor) säätelevät monimutkaisella prosessilla miten ja milloin jokin proteiini solussa aletaan koodata. Transkriptiotekijät mm. paikallistavat DNA:n promootorialueita (alue, jonka jälkeen jonkun proteiinin koodaus voi alkaa) sitoutumalla suoraan DNA:n ympärille sopivalla avaruudellisella tavalla - usein lisäksi muiden transkriptiotekijöiden avustamana.

- *Metabolia* (Aineenvaihdunta)

Metabolia on biologinen prosessi, jossa ravintoaineita muokataan, jotta saataisiin energiaa ja rakennusaineita solun tai eliön käyttöön. Metabolisia prosesseja on anabolisia ja katabolisia. Anabolisissa prosesseissa rakennetaan monimutkaisempia molekyylejä yksinkertaisista rakennusaineista. Katabolisissa prosesseissa monimutkaisemmat yhdisteet hajotetaan pienemmiksi.

- *Metaboliitti* (Aineenvaihduntatuote)

Metaboliitti on aineenvaihdunnassa syntynyt yhdiste, josta usein syntyy edelleen muita yhdisteitä. Metaboliaverkoilla kuvataan solun aineenvaihdunnan tapahtumia suunnattuna verkkona, jossa kukin solmu vastaa reaktiota ja reaktioon liittyy sitä katalysoiva entsyymi. Kaaret yhdistävät reaktiot reaktiopoluiksi.

- *cDNA* (complementary DNA)

cDNA on komplementaarinen DNA-juoste, joka on syntetisoitu käyttäen mallina transkriptiossa syntynyttä RNA:ta.

- *Geeniekspressio* (geenin ilmeneminen - gene expression)

Geeniekspressio⁴ on DNA:n koodaaman proteiinin tuottamista solussa. DNA:ssa olevat geenit eivät ilmene (samankaan eliön) joka solussa samalla tavalla. Ilmeneminen on monimutkaisesti ohjattua ja riippuu solun laadusta ja ympäristön olosuhteista. Ekspressiossa yhdistyy mRNA-molekyylien, DNA-sekvenssien ohjaama, transkriptio ja mRNA-molekyylien ohjaaman proteiinin tuottaminen translaatiossa. Väljemmin puhuttaessa pelkkää transkriptiota kutsutaan ekspressioksi. Eri geenien yhtäaikaista ekspressiota voidaan mitata DNA-siruilla.

- *DNA-siru* (mikrosiru - siru - DNA microarray - chip)

Katso liite 2: DNA-mikrosiru.

- *BLAST* (Basic Linear Alignment Search Tool)

BLAST on eräs käytetyimmistä sekvenssien hakumenetelmistä. Se on heuristinen menetelmä ja sen sisältämä ohjelmistovalikoima mahdollistaa tietokantahakujen tekemisen riippumatta siitä, ovatko hakusekvenssi ja tietokannan sisältämät sekvenssit proteiini- tai nukleotidisekvenssejä. Monipuolisten hakumahdollisuuksien lisäksi BLAST:n etuja ovat käytettyjen algoritmien nopeus, herkkyyys sekä ohjelmiston vapaa saatavuus. Useimmille BLAST on tuttu ennen kaikkea NCBI:n ja EBI:n kaltaisten keskusten www-palvelimista. Näiden kaikille avointen palvelujen lisäksi monet laboratoriot ja palvelukeskukset ovat toteuttaneet omia BLAST-palveluja. Altschul ja kumppaneiden Basic Local Alignment Search Tool [ASG90] on tieteen eniten viitattu artikkeli.

⁴Viime vuosina on pystytty mittaamaan yhden solun ekspressiota. Animaatio yhden solun ekspressiosta löytyy esimerkiksi Yeshiva yliopiston Singer Laboratoriosta osoitteesta http://singerlab.aecom.yu.edu/supplements/natrevmcb_v5p855/movies03.htm.

- *SNP* (Single Nucleotide Polymorphism, snipit)

SNP on mutaatio (emäsosan vaihtuminen) yhdessä nukleotidissa, joka periytyy. Sen täytyy esiintyä vähintään yhdessä prosentissa populaatiota, jotta se hyväksytään snipiksi. Snipit edustavat 90% ihmisten genomien välisistä eroista ja niitä ilmenee 100 - 300 emäksen välein sekä introneissa että eksoneissa. Koska tällaiset perimäerot selittävät suurimman osan yksilöiden välisestä eroavuudesta, luonnonvalinta kohdistuu erityisen voimakkaasti niihin. Snippejä käytetään paljon laboratoriokokeiden suunnittelussa, kun haetaan sopivaa kohtaa DNA:ssa, johon koe tulisi kohdistaa. Eri organismien snipeille on olemassa tietokantoja.

- *PCR* (Polymerase Chain Reaction)

Menetelmä, jolla syntetisoidaan DNA:n jostain segmentistä monta kopiota. Siinä käytetään lämpökäsittelykiertoa, jotta polymeerasireaktio saadaan halutunlaiseksi. Menetelmää käytetään monien geneettisten tutkimusten osana.

- *Elektroforeesi* (Electrophoresis)

Elektroforeesi on analysointimenetelmä, jolla eri painoiset DNA-sekvenssit voidaan erotella. Tulokset kuvataan ja näin syntyy kuvainformaatiota, jota on tarpeen analysoida ja säilyttää. Esimerkiksi PCR-menetelmällä (ks. edellä) tuotettuja DNA-sekvessejä voidaan analysoida elektroforeesin avulla.

- *Proteomiikka ja proteomi*

Proteomiikka tutkii proteiinien rakennetta ja ominaisuuksia huipputehokkaita laitteita ja teknologioita hyödyntäen. Miikka-päätettä käytetään samaan tapaan monilla biotieteen aloilla. Proteomi on solun (tai koko eliön solujen) tuottamien proteiinien kokonaisuus samoin kuin genomi on geenien ja muun perimäaineen kokonaisuus. Transkriptomi on puolestaan solun tai solujoukon mRNA-molekyylien kokonaisuus.

Liite 2. DNA-mikrosiru

DNA-mikrosirussa (tarkemmin cDNA, ks. sanasto) on lasi- tai nailonalusta, johon on kiinnitetty kokoelma cDNA-koettimia (cDNA probe). cDNA-koettimet sisältävät 50-500 emäsparia kukin. Koettimina voidaan käyttää myös olikonukleotideja. Olikonukleotideissa on vain 10-80 nukleotidia.

Mikrosirukokeessa hybridisoidaan sirun koettimia leimattujen näytemolekyylien kanssa. Näytemolekyylit ovat peräisin laboratoriokokeen solunäytteen (usein punainen leima) ja kontrollisolunäytteen (usein vihreä leima) ekspressiokokeesta. Koettimet vastaavat (lähes) kaikkia eliön geenejä ja leimatut molekyylit vastaavat (jossain määrin) koetilanteen transkriptiossa (ks. sanasto) tuotetun mRNA:n määrää. Koettimien ja leimattujen molekyylien nukleotidien emäsosat pyrkivät hybridisoitumaan aina pareittain A T:hen ja C G:hen (jos A, adeniini, on koettimessa niin T, tymiini, on leimatussa molekyylissä - eli esimerkiksi ATACG tarttuisi TATGC:hen).

Koska koettimet ovat spesifisiä, paikallaan siruhilassa ja paikka tiedetään, voidaan solunäytteessä tapahtunutta geeniekspression kvantitatiivista suuruutta arvioida leimattujen molekyylien hybridisaatiosuhteen määrää mittaamalla. Mittaustulokset annetaan usein käsitellyn näytteen cDNA-juosteen ja kontrolli-cDNA-juosteen hybridisaatiomäärien suhteen kaksikantaisena logaritmina. Hybridisaation määrä arvioidaan laser-laitteella mittaamalla leimausaineiden lähettämän valon intensiteetit.⁵ Sirun pinnalla on säännöllinen hila (grid) ja jokaisessa hilapisteessä (spot, näitä on jopa kymmeniätuhansia) on identifioitu koetinmolekyylijoukko siten, että yhtä laatua on aina vain yhdessä hilapisteessä. Hilapisteessä voi olla esimerkiksi 10^7 - 10^8 yhden geenin kloonina ja hilapisteitä voi yhdellä sirulla olla esimerkiksi kymmenentuhatta kappaletta. Hybridisointitapahtumassa näyteliuoksen leimatut molekyylit tarttuvat emäsjärjestyksensä vastaaviin koettimiin. Kun hybridisaatio on stabiloitunut kylliksi, ylimäärä näyteliuoksesta huuhdellaan pois, siru kuivataan ja hilaan tarttuneet leimatut molekyylit kuvataan laserlaitteella.

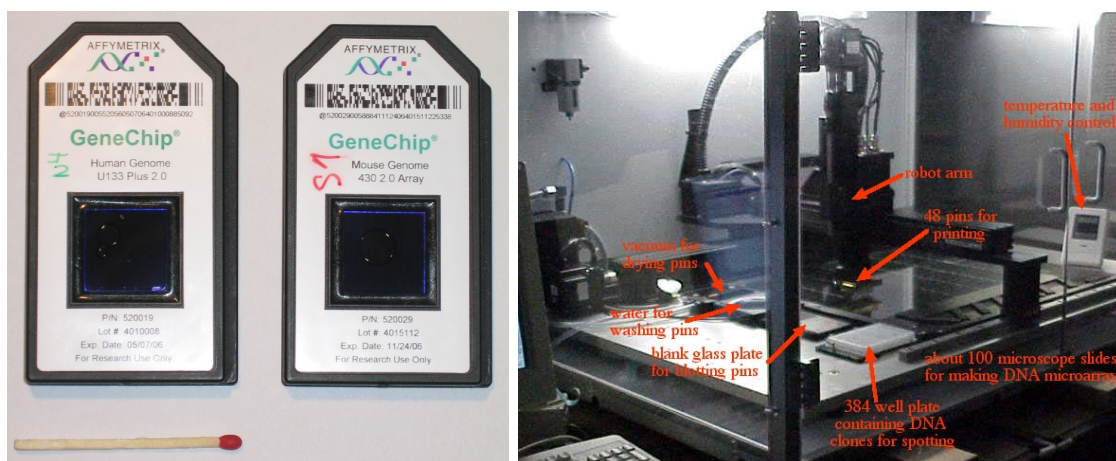
Kuvassa 3 a⁶ on kaksi Affimetrix:in sirua ja kuvassa 3 b⁷ on robotti, joka valmistaa laboratorioon omatekoisia siruja.

Kuvassa 4, on CSC:n julkaiseman DNA Microarray Data Analysis -kirjasta [TuL05]

⁵Samaan tarkoitukseen käytetään fluoresoivien väriaineiden lisäksi radioaktiivisuuteen perustuva leimaamista.

⁶Kuva 1 a on Wikipediasta ja julkaistu 'GNU Free Documentation License' -lisenssillä.

⁷Kuva 1 b on sivulta <http://www.bio.davidson.edu/Courses/genomics/arrays/arrayer.html>; sieltä löytyy myös liikkuvaa kuvaa mikrosirujen valmistuksesta.



(a) Affymetrix GeneChip -siruja

(b) Sirujen printtausrobotti

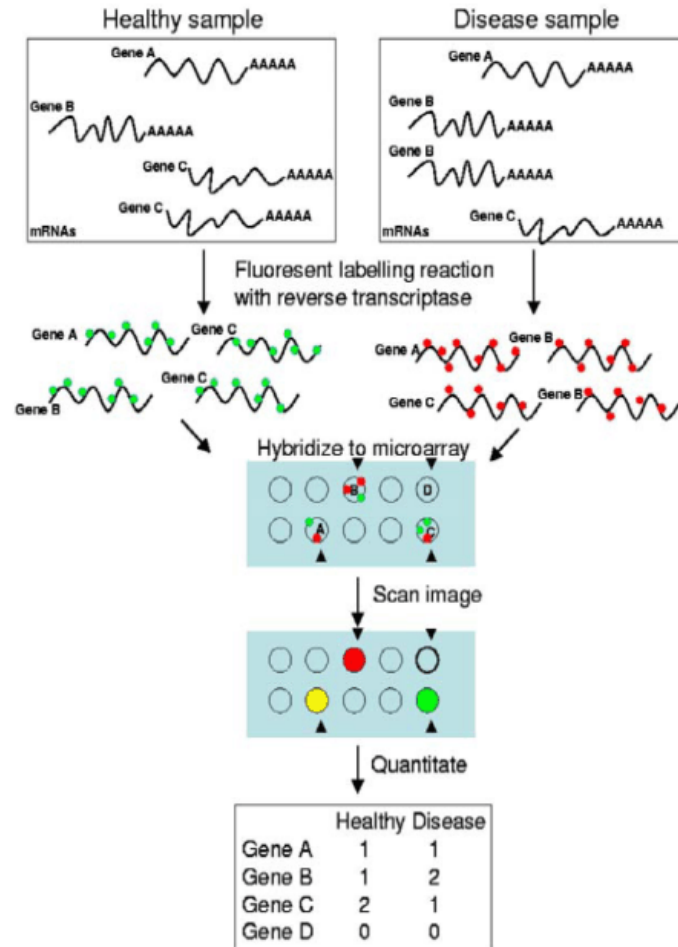
Kuva 3: a-kuvassa on ihmisen ja hiiren ekspressiomittauksiin Affymetrix Inc.:in valmistamat DNA-mikrosirut. Sirulle mahtuu koettimia niin paljon, että eliön lähes kaikkien geenien ilmeneminen voidaan mitata. Affymetrix:in GeneChip-sirujen tuottama ekspressioarvo kuvaa geenin suhteellista ekspressiota saman mittauksen muiden geenien ekspressioon. Siru tuotetaan *in situ* syntetisoinnilla, jolloin koettimet 'rakennetaan' fotolitografisella maskilla suoraan omille hilapaikoilleen. Näin tuotetuilla siruilla genejä vastaavat paikat ovat samassa sirutyypissä identtiset. Koska sirut ovat tietyn standardin mukaisia, eri sirujen tuottamaa dataa on helpompi vertailla keskenään kuin omatekoisten sirujen. Laboratoriot valmistavat kuitenkin siruja itse, koska silloin koetinvalikoima voidaan suunnitella omaa tarvetta vastaavaksi (ja ne ovat myös halvempia). Koettimien näytevalikoima printataan sirun pinnalle robotiikkaa hyväksi käyttäen. b-kuvassa on mikrosiru robotti The Davidson Trust:in Department of Biology:n Davidson College:sta. Tällaisen sirun (kaksiväri-DNA-siru) tuottama ekspressioarvo on suhteellista geenin eri näytteiden välillä - koenäytteen (esimerkiksi: sairas kudos) ja kontrollinäytteen (esimerkiksi: verrokki kudos).

kuva, jossa on esimerkki mikrosirukokeen työnkulusta.

Sirut saadaan kuva-analyysin ja normalisoinnin tuloksena monen tuhannen geenin ekspressioarvot yhtä aikaa. Mittaustulokset annetaan (kaksivärisiruilla) usein sirukokeen cDNA-juosteen ja kontrolli-cDNA-juosteen hybridisaatiomäärien suhteen logaritmina.

Logaritmuunnos tehdään, koska suhteellisessa mittauksessa näytteen kaksinkertaisen ekspression suhde kontrolliin olisi arvoltaan 2 ja kaksi kertaa pienemmän olisi 0.5. Kun arvot muutetaan logaritmisiksi, suhteen arvot $[0, 1]$ päätyvät välille $[-\infty, 0]$ ja suhteen arvot $[1, \infty]$ päätyvät välille $[0, \infty]$. Näin yli- ja aliekspressoituminen (kontrolliin verrattuna) skaalautuu samalla tavalla. Kirjallisuudessa toivotaan usein, että suhteelliseen ekspressiodataan saataisiin yhteinen kontrollistandardi, jolloin arvojen keskinäinen vertailu olisi helpompaa.

Sirut valmistetaan eliökohtaisesti. Siruja voi hankkia teollisesti valmistettuna (esim. Affymetrix, ks. kuva 3) tai niitä voidaan valmistaa laboratoriossa korkean teknologian erityislaitteilla.



Kuva 4: Kuvassa näkyy tyypillisen kaksiväri-DNA-mikrosirun (koenäyte ja kontrollinäyte värjätään eri väreillä) työnkulku. Taulukon arvoista saadaan ekspressiomittauksen arvot (esimerkiksi) ottamalla kaksikantainen logaritmi sairaan ja terveen näytteen intensiteettiarvoista - geeni A: $0 = \log_2(1/1)$, geeni B: $1 = \log_2(2/1)$, geeni C: $-1 = \log_2(1/2)$, geenillä D arvo puuttuu (se voidaan arvioida esimerkiksi KNN-menetelmällä, K-Nearest Neighbors, jos dataa on tarpeeksi). Silloin yli- ja aliekspressoituminen skaalautuu samalla tavalla.