

# Tiedonlouhinta hiirien ja rottien geeniekspressiodatasta

Tiedonlouhinta 2013

Antti Kurronen ja Irene Pöllänen

# Data: Rotat

- Ekspressiodataa, joissa geenien suhteellisten ilmenemisten tasoja\*
- ~ 30 000 riviä -> ~ 1 rivi / geeni

## **ZL\_E18\_SD\_1\_150405 - linear (x3)**

–18 päivän alkio

## **ZL\_neo\_SD\_1\_080405 - linear (x3)**

-vastasyntynyt

## **ZL\_d10-rat-heart\_#3\_130407 - linear (x3)**

–10 päivän vanha

## **ZL\_adult\_SD\_1\_110505 - linear (x3)**

–nuori aikuinen 3kk

## **ZL\_adult\_WS\_1\_100605 - linear (x3)**

–nuori aikuinen, eri kantaa (Wistar, muut Sprague-Dawley)

## **ZL\_old\_WS\_1\_100605 - linear (x3)**

–vanha Wistar 24kk

\*(Kustakin rotasta kolme replikaattia)

# Data:Hiiret

- Ekspressiodataa, joissa geenien suhteellisten ilmenemisten tasoja
- ~ 25 000 riviä -> ~ 1 rivi / geeni

## **JP\_control\_1\_050511\_(MoGene-1\_0-st-v1) - linear (x3)**

–Kontrolli villityypin hiiri

## **JP\_mutant\_1\_050511\_(MoGene-1\_0-st-v1) - linear (x3)**

–Twinkle-transgeeninen hiiri

## **JP\_S1\_SOD het\_210512\_(MoGene-1\_0-st-v1) - linear (x3)**

–SOD2 heterotsygootti poistogeeninen

## **JP\_ST1\_SOD\_het\_Tw\_pos\_251012\_wdh\_(MoGene-1\_0-st-v1) - linear (x3)**

– SOD2 heterotsygootti poistogeeninen x Twinkle-transgeeni risteymä

\*(Kustakin hiirestä kolme replikaattia)





# Tiedonlouhinta tavoitteita:

## **Geenien ekspression voimistuminen ja pienentyminen geeniontologialuokittain**

- Tietyissä biologisissa prosesseissa (Biological Process)
  - Molekulaarisissa toiminnoissa (Molecular Function)
  - Solun osissa (Cellular Component)
- > Samankaltaisuudet ekspressiotasojen muutoksissa hiirillä ja rotilla?

Esim. Proteolyysi on yksi biologisen prosessin luokka, johon kuuluvat geenit hajoittavat proteiineja

## **Hiiren ja rotan ekspressiotasojen vertailu uusilla muuttujilla**

- Vaikuttaako hiirien mutaatiot geenien ekspressiotasoihin samalla tavoin kuin rotilla vanheneminen

# Toteutus

## **Datan esikäsittelyyn**

- R + omia funktioita
- Excel + VBA skriptejä
- Gnumeric

## **Tilastolliseen analysointiin**

- Kingfisher
- Jokin rinnakkainen menetelmä (R tai Matlab)

# Datan esikäsittely:

## Poistetaan

- Datarivit, joista puuttuu geenitiedot
- Identtiset rivit (poistetaan kopiot)

## Yhdistetään

- Datarivit, jotka ovat saman geenin tuloksia  
-> valitaan arvot, joilla hajonta pientä

## Muuttujat

Luodaan ryhmien välisiä eroja kuvaavia muuttujia



# Datan esikäsittely:

## **Normalisointi**

- Lowessin menetelmällä

## **Numeerisille arvoille log<sub>2</sub> -muunnos**

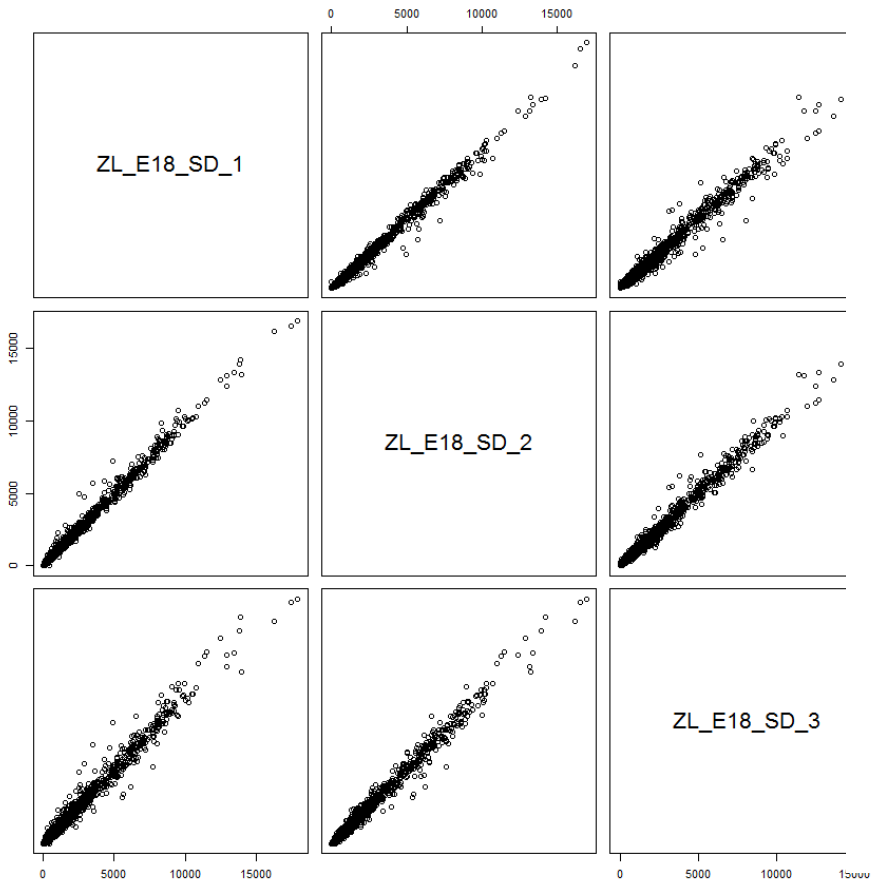
-> mukavemman näköinen data

## **Korkeimman 5% poisto – tarvitaanko**

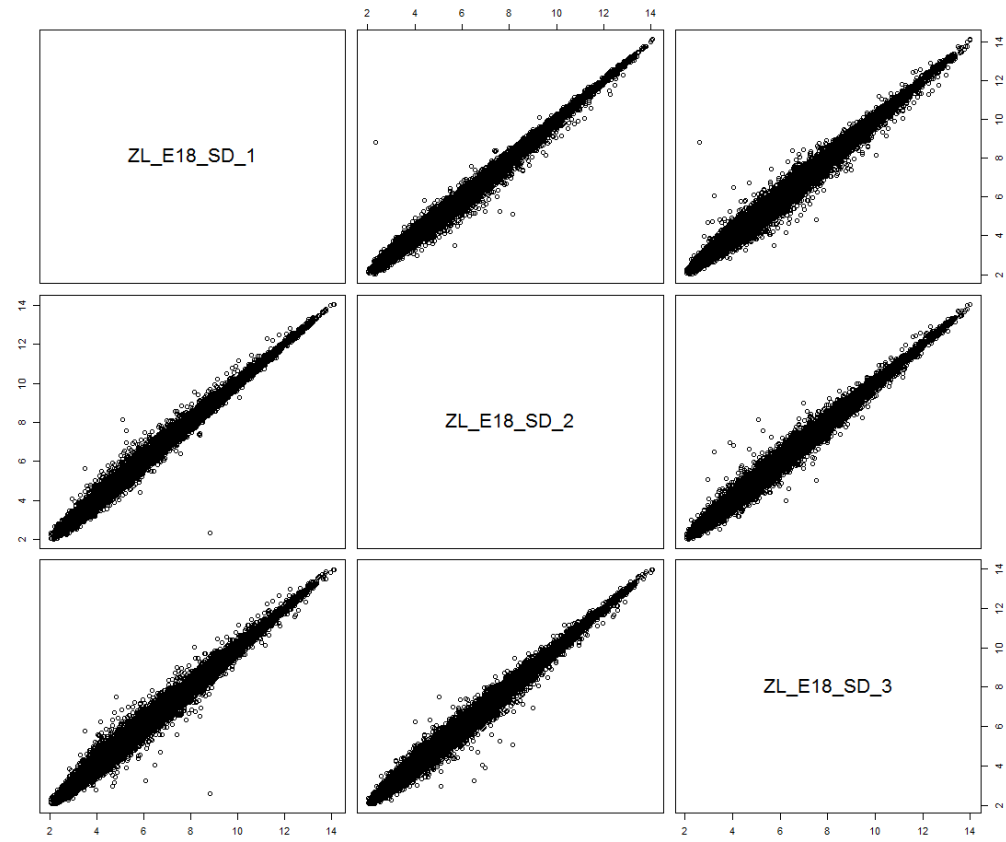
-laitteiston aiheuttamaa virhettä?

# Datan esikäsittely:

Rat original data: ZL\_E18\_SD

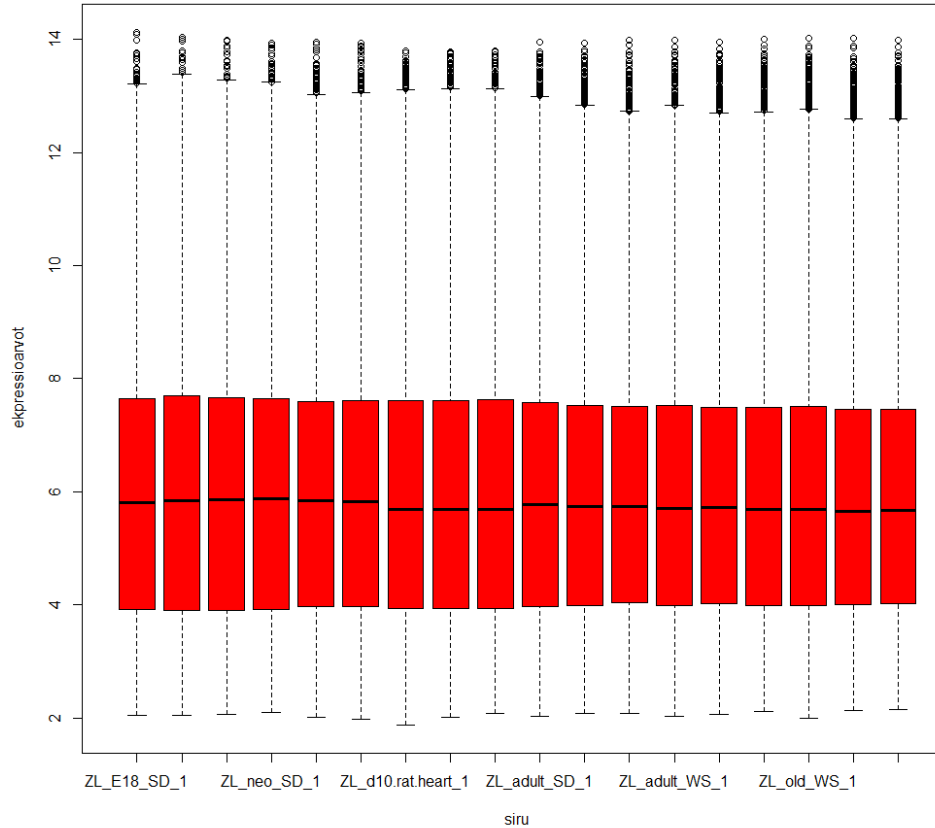


rat log2 transformed data: ZL\_E18\_SD

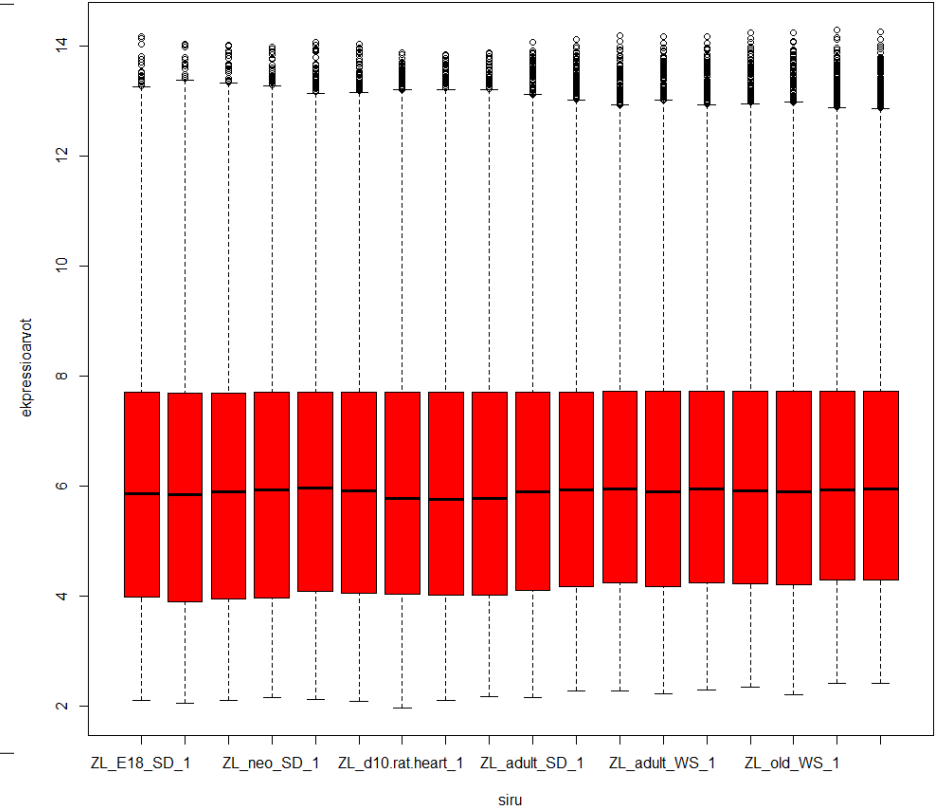


# Datan esikäsittely: Normalisointi?

Boxplot: log2 transformed data

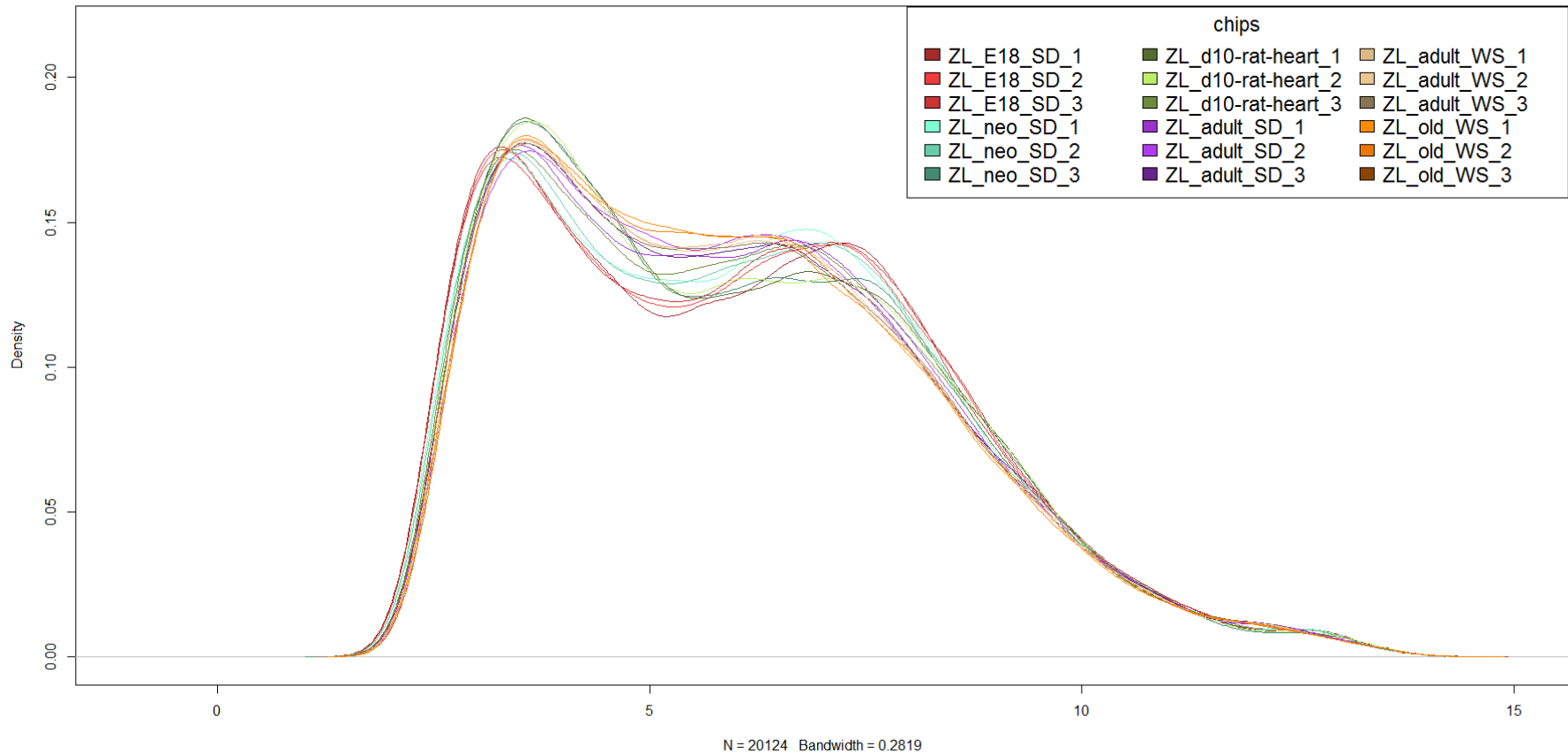


Boxplot: log2 transformed normalized data



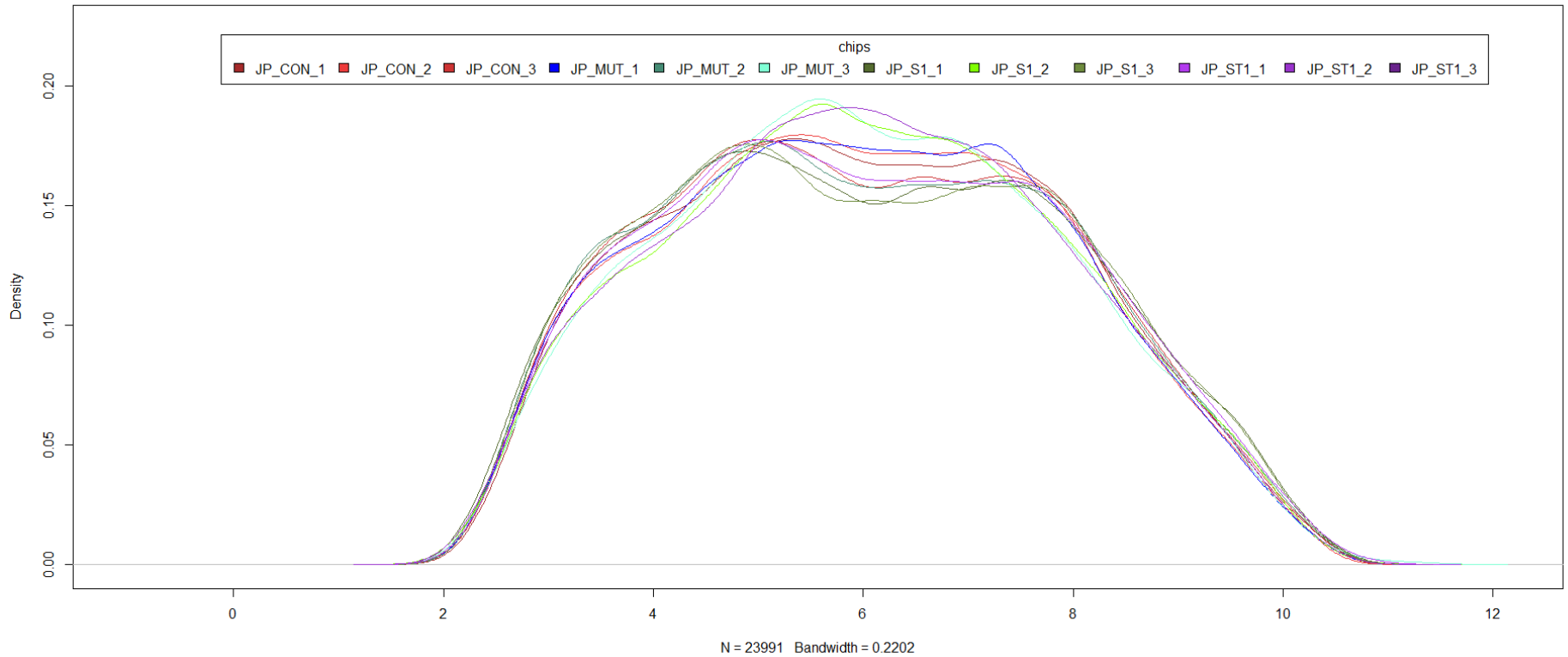
# Datan esikäsittely: jakauma

Density: Rat log2 transformed data



# Datan esikäsittely: jakauma

Density: Mouse log2 transformed data (normalized)



# Alustavia tuloksia

## **Tilastolliset analyysit vasta alkaaneet**

- Esitarkasteltu dataa, jota ei ole täysin esiprosessoitu
- GO-luokkien ja ekspressiotasojen välisiä korrelaatioita ei vielä todettu hiiridatalla (Kingfisher ohjelmalla)

# Loppupäätelmät

- Työläintä oli ymmärtää itse ongelmaa samalla tavalla muiden kanssa
- Datassa oli paljon esikäsiteltävää
- Data on hyvää -> kaikki replikaatit saadaan valittua analyysiin.
- Myöhään löydetyt hiiren datassa olevat identtiset rivikopiot viivästyttivät analyysia