

How to prepare the percentages for a heatmap?  
What is the meaning of the color scale?

Questions

Tab Delimited, Multiple Sample (TDMS)  
File Format

NAME	GenBank	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5
Treatment	placebo	placebo	placebo	placebo	placebo	placebo
Patient ID	101	102	103	104	105	
Gender	F	F	M	M	M	
A1	AA714409	0.245457	0.226136	-0.273932	0.273945	0.416614
A2	AA156571	-0.32316	0.007095	0.354105	-0.32402	-0.40812
A3	AI219422	0.202585	-0.08818	-0.32919	-0.36248	-0.29789
A4	AA521392	-0.35597	0.253059	0.359742	-0.14589	0.032914
A5	AI344051	0.033734	0.341731	0.371982	-0.11216	-0.48239
A6	AA705217	0.104424	-0.19359	0.426171	0.275234	-0.07636
A7	AA595793	-0.02515	0.047777	-0.30558	-0.15541	0.401895
A8	HB1621	0.024793	-0.10033	0.137568	-0.10322	0.287513
A9	T48924	0.483648	0.429508	-0.23373	0.369819	0.204742
A10	W24076	-0.03915	-0.00601	0.479825	0.085283	0.338572
A11	AA596275	0.265352	0.067	0.027318	-0.05708	0.463785
A12	AA911076	-0.21338	0.229778	-0.18721	0.052883	0.119687
A13	CA11036	-0.43006	0.067805	0.703913	-0.45023	-0.34687

Matrix values are percentages  
Lines are populations of cells  
Columns are samples  
Samples are arranged as groups

**Usual "absolute difference" use:**  
log transform  
overall centering

**"Scaled discovery" use:**  
log transform with soft floor  
overall centering using median  
standardize range using percentiles

**"Supervised" use:**  
log transform with soft floor  
overall centering using median  
standardize average within group dispersion

**"Reference group" option:**  
discovery or supervised use then  
fix the reference group level

**"Summary" use:**  
discovery or supervised use  
reduce to central tendency of each group

**"Noisy differences" use:**  
discovery (or supervised) use then  
sort values within group

**"Robust noisy differences" use:**  
discovery or supervised use  
sort values within group  
substitute values using one gaussian per group

Typical uses

Definitions

Main formula

transform, center, then scale

$$display_{ij} = \frac{transform(percent_{ij}) - center_{ij}}{scale_{ij}}$$

CytoPreparePercentage

Steps

- 1/ transform
  - none
  - log2
  - asinh
  - option: soft floor**
- 2/ centering
  - to: overall.line, ref.group, line.but.ref
  - method: trimmed mean, median
- 3/ standardize
  - to: overall dispersion = 1
  - to: low-high percentile range = 1
  - discovery = no grouping
  - to: average within group dispersion = 1
  - supervised = groups defined**
  - method: dispersion = sd or mad
- 4/ sort within group
  - imputation: NA = median
- 5/ substitute
  - overall rank
  - overall gaussian
  - gaussian within group
- 6/ within group summarized value
  - median
  - median plus percentiles

Step benefits

- transform:**
- \* compress overall dynamical range
  - \* make dynamical range similar between populations
  - \* soft floor avoids over dispersion of lowest values
- centering:**
- \* remove base line for each population
  - \* allow comparing differences within each population
  - \* color scale relates differences to the base line
- standardize:**
- \* relate differences within each population to the same standard
  - \* stretch or squeeze differences in proportion to the standard
  - \* color scale relates differences to the same standard
  - \* within group dispersion standard enlightens statistical significance
  - \* wgds avoids confusing within and between group dispersions
- sort within group:**
- \* visual check that values between groups overlap or not
- substitute:**
- \* reshape data in a robust approach trying to weaken outliers
  - \* grouping is kept between population, but sample order is lost
- reference group:**
- \* set a group as base line; any further scaling should not change it
  - \* define the value for this group
- summarize:**
- \* give a synthetic view, possibly with dispersion