Conference on Phenotype MicroArray Analysis of Cells, Florence/Italy 2015: opm Workshop

Analysing Phenotype MicroArray data with *opm*

Benjamin Hofner, FAU, Erlangen/Germany Lea A.I. Vaas, Fraunhofer Institute, Hamburg/Germany María del Carmen Montero Calasanz, Newcastle University, Newcastle/UK Markus Göker, DSMZ, Braunschweig/Germany









Why did we implement opm?

Part 1:

- Free, easily extendable software
- Flexible production of high-quality graphics and robust statistical analysis
- Interactive or fully automated usage possible

Part 2:

- Flexible **metadata** management
- Reproducible research
- Easy interaction with other software (and databases)

Why R?

- *de facto* standard for freesoftware statistical computing
- all operating systems
- flexible and clean coding
- non-interactive and interactive use
- powerful GUIs/IDEs (e.g. RStudio[™])







www.rstudio.com

Parts of this workshop

- 1. Reading PM data into opm
- 2. Adding and manipulating **metadata**
- 3. Visualizing PM data in opm
- 4. Estimating and visualizing **curve parameters**
- 5. Statistical comparisons of experimental groups

Metadata

- "data about data"
- structural or (here) descriptive: organism measured, growth conditions etc.

Goal of opm:

- **self-describing objects** that carry all user-relevant information (data and their metadata)
- **sufficient** for subsequent plots and statistical tests
- can be written to files and **distributed** between systems
- use a **light-weight**, flexible file format

The data: multi-dimensional, need aggregation



Six wells, 32 plates

Aggregating: estimating curve parameters



Options

- **parametric models** (Gompertz, logistic etc.): explicitly estimate the four parameters or parameters from which the four can be derived
- **splines**: estimate smooth curves non-parametrically and analytically derive the four parameters from the smooth curves

Aggregating: estimating curve parameters



 parametric models (Gompertz, logistic etc.) work only well for rather regular curve shapes

 splines usually work as well in these situations and outperform them in other cases

Vaas et al. PLoS ONE 7: e34846, 2012

Aggregating: estimating curve parameters

Some spline methods perform better than others:



Hour

Hour

What is the correct statistical procedure to compare four groups pairwise to each other?

Data situation:		Typical approach: six separate t-tests	
7 plates: 10 plates: 10 plates: 5 plates:	Ax1 Ax2 Ax4 Ax6	Ax2 - Ax1 ,	p =
		Ax4 - Ax1 ,	p =
		Ax6 - Ax1 ,	p =
		Ax4 - Ax2 ,	p =
		Ax6 - Ax2 ,	p =
		Ax6 - Ax4.	p =

What is the correct statistical procedure to compare four groups pairwise to each other?



Why the t-tests fail:



The more pairwise comparisons of multiple groups from the same sample, the larger the risk of identifying **falsepositive** differences.

One must correct for the **family-wise error rate**

95% family-wise confidence level

Multiple comparison of means

- user-defined comparisons
- inherent multiplicity adjustment
- significance of difference
- and effect size visible
- on original scale

Hothorn, T. et al. (2008) Simultaneous inference in general parametric models. Biometr. J. 50. 346-363.

A02 (Dextrin) - A01 (Negative Control) A03 (D-Maltose) - A01 (Negative Control) A04 (D-Trehalose) - A01 (Negative Control) A05 (D-Cellobiose) - A01 (Negative Control) A06 (b-Gentiobiose) - A01 (Negative Control) A07 (Sucrose) - A01 (Negative Control) A08 (D-Turanose) - A01 (Negative Control) A09 (Stachyose) - A01 (Negative Control) A10 (Positive Control) - A01 (Negative Control) A03 (D-Maltose) - A02 (Dextrin) A04 (D-Trehalose) - A02 (Dextrin) A05 (D-Cellobiose) - A02 (Dextrin) A06 (b-Gentiobiose) - A02 (Dextrin) A07 (Sucrose) – A02 (Dextrin) A08 (D-Turanose) - A02 (Dextrin) A09 (Stachvose) - A02 (Dextrin) A10 (Positive Control) – A02 (Dextrin) A04 (D-Trehalose) - A03 (D-Maltose) A05 (D-Cellobiose) - A03 (D-Maltose) A06 (b-Gentiobiose) - A03 (D-Maltose) A07 (Sucrose) - A03 (D-Maltose) A08 (D-Turanose) - A03 (D-Maltose) A09 (Stachyose) - A03 (D-Maltose) A10 (Positive Control) - A03 (D-Maltose) A05 (D-Cellobiose) - A04 (D-Trehalose) A06 (b-Gentiobiose) - A04 (D-Trehalose) A07 (Sucrose) - A04 (D-Trehalose) A08 (D-Turanose) - A04 (D-Trehalose) A09 (Stachyose) - A04 (D-Trehalose) A10 (Positive Control) - A04 (D-Trehalose) A06 (b-Gentiobiose) - A05 (D-Cellobiose) A07 (Sucrose) – A05 (D–Cellobiose) A08 (D-Turanose) - A05 (D-Cellobiose) A09 (Stachyose) - A05 (D-Cellobiose) A10 (Positive Control) - A05 (D-Cellobiose) A07 (Sucrose) - A06 (b-Gentiobiose) A08 (D-Turanose) - A06 (b-Gentiobiose) A09 (Stachyose) - A06 (b-Gentiobiose) A10 (Positive Control) – A06 (b–Gentiobiose) A08 (D-Turanose) - A07 (Sucrose) A09 (Stachyose) - A07 (Sucrose) A10 (Positive Control) - A07 (Sucrose) A09 (Stachyose) - A08 (D-Turanose) A10 (Positive Control) - A08 (D-Turanose) A10 (Positive Control) - A09 (Stachyose)



Data manipulation in *opm*



CSV := Character Separated Values; YAML, a data serialization format

Data analysis in opm



Further reading

• **Project homepage:** http://opm.dsmz.de

• **Package and further resources:** http://www.goeker.org/opm

• **Development version of the package:** https://r-forge.r-project.org/R/?group_id=1573

Overview on opm:

Lea A.I. Vaas, J. Sikorski, B. Hofner, N. Buddruhs, A. Fiebig, H.-P. Klenk and M. Göker. "opm: An R package for analysing OmniLog® Phenotype MicroArray Data". Bioinformatics 29 (14): 1823-1824, 2013.

Parameter comparison and visualization with opm:

Lea A.I. Vaas, J. Sikorski, V. Michael, M. Göker and H.-P. Klenk. "Visualization and curve-parameter estimation strategies for efficient exploration of phenotype microarray kinetics". PLoS ONE 7 (4): e34846, 2012.

Advanced usage of opm to detect differential expressions:

B. Hofner, L. Boccuto and M. Göker. "Controlling false discoveries in highdimensional situations: Boosting with stability selection". BMC Bioinformatics 16 (6): 144, 2015.