

Faculty of Life Sciences



Getting a little from a lot, reduction and modeling of -omics data

Kasper Brink

February 3, 2012 Slide 1/29

Intro

Cassava (Manihot esculenta)



- Important source of nutrition in Africa, Asia and South America
- Contains Cyanogenic Glucosides



Intro

Cassava (Manihot esculenta)



- Important source of nutrition in Africa, Asia and South America
- Contains Cyanogenic Glucosides

<u>3/8/09</u>

By:"YoshiMa



Intro

- 32 plants
- LC-MS spectra from the 32 plants
- Arrays for 13865 genes from the same plants
- Construct a model on the form
 Spectra = expression × coefficients + noise
- Problems: dimension reduction, variable selection, validation







• 109431241



Slide 4/29 — Kasper Brink — Getting a little from a lot, reduction and modeling of -omics data — February 3, 2012



	time	mz	intensity	
[1,]	5.205	124.7990	20	
[2,]	5.205	125.1057	17	
[3,]	5.205	125.5501	12	
[4,]	5.205	126.6448	27	
[5,]	5.205	127.3360	13	
[6,]	5.205	127.7387	22	
[3682	2288,]	1802.489	1202.724	398
[3682	2289,]	1802.489	1203.069	309
[3682	2290,]	1802.489	1203.419	178
[3682	2291,]	1802.489	1204.025	70
[3682	2292,]	1802.489	1204.426	48
[3682	2293,]	1802.489	1204.741	80







Slide 7/29 — Kasper Brink — Getting a little from a lot, reduction and modeling of -omics data — February 3, 2012



■ 109431241■ 17463549



Extract signals from 3D tensor

Slide 9/29 — Kasper Brink — Getting a little from a lot, reduction and modeling of -omics data — February 3, 2012



The easy reduction: Average in the required direction





109431241 17463549 11872



Slide 11/29 — Kasper Brink — Getting a little from a lot, reduction and modeling of -omics data — February 3, 2012

The more meaningful reduction: Decomposition into mixing matrices and underlying components with some restraints (non-negativity, sparsity)



This requires that the individual spectra are exactly the same dimension to construct the tensor $\underline{\mathbf{Y}}$



Slide 12/29 — Kasper Brink — Getting a little from a lot, reduction and modeling of -omics data — February 3, 2012





109431241 17463549 11872 1855











		151	_1		151_2	2	151_3		151_4	1	152_1	
[1,]	24	13.51	56	170	.1290) 137	.9091	131	.3433	138.	5385	
[2,]	181	3.11	80	2502	.2420	0 1337	.7850	1903	.2860	3155.	5690	
[3,]	647	1.47	70	1927	.5710) 1923	.8660	1023	.3690	1088.	7270	
[4,]	1279	94.19	00	5849	.8330	8092	.0160	12470	.7900	17349.	1100	
[5,]	87	75.24	62	649	.8030	516	.7879	451	.6719	694.	3333	
[6,]	32	24.43	75	139	.4516	5 157	.3231	121	.8254	260.	1695	
[1386	0,]	257	.98	460	265	.31250	299	.25000	399	.62120	149	.66100
[1386	1,]	96	.68	852	257	.23440	291	.28130	178	.21540	108	.80650
[1386	2,]	596	.07	690	264	.51560	143	.07940	66	.92424	52	.63768
[1386	3,]	56	.67	742	71	.33871	60	.19672	88	.40984	52	.62903
[1386	4,]	2112	.12	500	1995	.82100	1702	.69400	2240	.42200	1716	.50000
[1386	5,]	1259	.39	400	3086	.60700	1551	.72900	881	.36360	2081	.33300

Numbers

n 109431241 **17463549 11872 a** 1855 **11872+443680**

How to choose relevant genes? The easy way: Fit linear model with multiple testing correction (limma) and select the genes with the lowest p-values. Very fast and simple, but is it meaningful?



A more interesting approach, clustering by profile. Use NMF algorithm



Modelling

At some point data will look like this



Multivariate regression with PLS



Slide 22/29 — Kasper Brink — Getting a little from a lot, reduction and modeling of -omics data — February 3, 2012

Numbers

109431241 17463549 11872 a 1855 **6** 11872+443680 o 11872+3200



Biological results

- The coefficient matrix can be seen as the effect of each gene at each time point
- At a given time point (peak) this can be used to construct a ranked list of genes
- Compare with existing knowledge of genes



Numbers

109431241 17463549 11872 a 1855 **11872+443680** ⁶ 11872+3200 **a** 300

Genetic algorithms

- How to determine the number of genes that go into the model?
- Possible solution is to use a genetic algorithm for the selection
- Allows selection from a larger set of genes, not just the top 100
- Fitness measure is the same as before

Genetic algorithms

- GA is an optimization algorithm
- Chromosomes 111100101010101111000 each have a fitness measure
- High fitness chromosomes are more likely to have offspring
- Crossover 111100101010101111000 001010101110101111000
- Mutation 011100101010101111000



