V12 Multi-variate analysis

Program for today:

- *Staphylococcus aureus* Africa project analysis for confounding variables
- Overview of multivariate analysis for omics projects
- Case study: gene-regulatory network for breast cancer
- Case study: single cell methylation and expression data

Review (V1): S. aureus in Germany vs. Africa: StaphNet

6 study sites each collected 100 isolates of healthy volunteers and 100 of blood culture or clinical infection sites.

Aim

microbiological and molecular characterization of African *S. aureus* isolates by DNA microarray analysis including clonal complex analysis

supplemented by Whole Genome Sequencing



Review (V1): Distribution of clonal complexes



Some clonal complexes more prevalent in Africa,

others predominant in Germany.

V12

Processing of Biological Data

Review (V1): Activitity of individual probes for CCs

	K R = 9 • C • − = S2 Excel Book - Microsoft Excel																				
Da	Datei Start Einfügen Seitenlayout Formeln Daten Überprüfen Ansicht ABBYY FineReader 12 Acrobat																				
	Ausschneiden	Calibri \cdot 11 \cdot A^{*} $\stackrel{*}{=}$ $=$	≫~~ {	Zeilen	umbruch		Standa	rd	+		< <u>₹</u>		Stand	ard	Gut		-	÷	*	Σ	
Einf	igen ≪ Sormat übertr	ngen F K U - 🔄 - 🖄 - 🔺 = = = =		a+ Verbin	den und ze	entrieren	- 19-	% 000	€,0 ,00 0,€ 00,	Bedir	≡zi ngte A	ls Tabelle	Neutr	ral	Schle	echt	E	infügen l	.öschen F	ormat	2
_	Zwischenablage	ा Schriftart ा	Ausri	thtung			a l	Zahl	G.	Formatio	erung + ior	matieren	Formatvo	orlagen					Zellen		
	B80 •	f _x tet.K.																			1
	А	В	С	D	E	F	G	Н	1	J	K	L	М	N	0	Р	Q	R	S	Т	
1		Groups	African	German	African	German	African (Serman	African	German	African (German	African	German 3.	African	German	African (German	African	German (Δfi
3		Numbers (n)	600	600	109	57	51	88	105	25	48	53	44	51	11	75	83	2	11	48	-
4		Clonal complex (CC)	all	CCs	CC	15	CC4	5	CC1	.21	CC	8	CC	5	CC	30	CC15	52	CC	.7	
5	SPECIES MARKERS	rrnD1Saureus.	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	
6		gapA kato	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	
8		rata coA	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	
9		nucl	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	
10		spa	100%	100%	100%	100%	98%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	
11		sbi	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	
12	REGULATORY GENES	sarA	100%	100%	100%	98%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	
13		saes	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	
14		agri total	35%	55%	0%	0%	41%	99%	99%	0%	100%	100%	0%	0%	0%	0%	100%	100%	100%	100%	
16		agrB.I	54%	60%	0%	0%	100%	100%	84%	92%	100%	100%	0%	0%	0%	0%	100%	100%	100%	100%	
17		agrC.I	57%	59%	0%	2%	96%	92%	99%	100%	100%	100%	0%	0%	0%	0%	100%	100%	100%	98%	
18		agrD.I	35%	55%	0%	0%	41%	99%	0%	0%	100%	100%	0%	0%	0%	0%	100%	100%	100%	100%	
19		agriltotal.	27%	25%	99%	100%	0%	0%	0%	0%	0%	0%	100%	100%	0%	0%	0%	0%	0%	0%	
20		agrB.II	27%	25%	99%	100%	0%	0%	0%	0%	0%	0%	98%	100%	0%	0%	0%	0%	0%	0%	
21		agrc.II	27%	25%	99%	100%	0%	0%	0%	0%	0%	0%	100%	100%	0%	0%	0%	0%	0%	0%	
23		agrilltotal.	16%	14%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	100%	100%	0%	0%	0%	0%	
24		agrB.III	16%	14%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	100%	100%	0%	0%	0%	0%	
25		agrC.III	15%	14%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	91%	97%	0%	0%	0%	0%	
26		agrD.III	16%	14%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	100%	100%	0%	0%	0%	0%	
27		agrIV.total.	37%	6%	0%	0%	59%	1%	100%	100%	6%	2%	0%	0%	0%	0%	100%	100%	0%	0%	
28		agrCIV	23%	41%	0%	0%	59%	1%	100%	100%	96%	98%	0%	0%	0%	0%	0%	100%	100%	98%	
30		hld	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	
31	METHICILLIN RESISTANCE	mecA	3%	4%	0%	0%	2%	2%	0%	0%	13%	2%	5%	16%	0%	0%	0%	0%	0%	0%	
32	AND SCCmer TVPING	delta_mecR	2%	3%	0%	0%	0%	2%	0%	0%	10%	2%	5%	16%	0%	0%	0%	0%	0%	0%	
33	Soomee III Ind	ugpQ	3%	4%	0%	0%	2%	2%	0%	0%	13%	2%	5%	16%	0%	0%	0%	0%	0%	0%	
34		ccrA.1	1%	0%	0%	0%	0%	0%	0%	0%	0%	4%	0%	0%	0%	0%	0%	0%	0%	0%	
36			0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	
37		Q9XB68.dcs	1%	3%	0%	0%	0%	2%	0%	0%	10%	0%	5%	12%	0%	0%	0%	0%	0%	0%	
38		ccrA.2	3%	4%	0%	0%	0%	2%	0%	0%	10%	2%	5%	16%	0%	0%	0%	0%	0%	0%	
39		ccrB.2	3%	4%	0%	0%	0%	2%	0%	0%	10%	2%	5%	16%	0%	0%	0%	0%	0%	0%	
40		kdpA	1%	1%	0%	0%	0%	0%	0%	0%	0%	0%	0%	12%	0%	0%	0%	0%	0%	0%	
41		Kaps Kaps	1%	1%	0%	0%	0%	0%	0%	0%	0%	0%	0%	12%	0%	0%	0%	0%	0%	0%	
42	/12	kdpD SCC	1%	1%	Proce	ssing	of Bio	loaie	al Dat	a 0%	0%	0%	0%	12%	0%	0%	0%	0%	0%	0%	
44		kdpE.SCC	1%	1%	0%	0%	0%	.0948 0%	0%	0%	0%	0%	0%	12%	0%	0%	0%	0%	0%	0%	
45		mecl	0%	1%	0%	0%	0%	0%	0%	0%	0%	0%	0%	12%	0%	0%	0%	0%	0%	0%	

Review (V1): Principal component analysis of 1200 strains

Input data: binary matrix of MA data; dimension 1200 x 334 probes PCA identifies local clusters that are characteristic

for particular clonal complexes



V12

Staphylococcus aureus data from Africa project (V1)

Site	# of cases below 1 year	# of cases 1 to 5 years	# of cases 6 - 25 years	# of cases 26 – 65 years	# of cases above 66 years
Africa + Germany (clinical)	88	109	90	225	88
Africa + Germany (commensal)	19	34	363	175	9
Africa (clinical)	86	106	53	54	1
Africa (commensal)	17	34	156	89	4
Germany (clinical)	2	3	37	171	87
Germany (commensal)	2	0	207	86	5

Age distribution is **heavily skewed**:

many small kids / babies in Africa – few seniors in Africa

very few small kids / babies in Germany – many seniors in Germany

Does this affect the analysis + interpretation?

Analyze whether age is a confounding variable

To test whether age is a **confounding variable**, one can compare the results from simple linear regression with those from multiple linear regression.

The principal difference between these two types of regression models is the number of explanatory variables:

(1) the simple linear regression (SLR) model uses only one dependent variable y and one explanatory variable x: $y = a + b \cdot x$

In our case, **y** stands for the binary output from the Alere-chip experiment for a particular gene. **y** therefore has values of 0 or 1.

With the binary variable x we could encode the sites Africa (x = 0) / Germany (x = 1). *a* and *b* are weights estimated by the model.

Generally SLR tries to find such weights (values for **a** and **b**) so that the difference between the estimated **y** and actual **y** will be the smallest.

Analyze whether age is a confounding variable

(2) the multiple linear regression model also has one dependent variable **y** but more than one explanatory variables

 $y = a + b_1 \cdot x_1 + b_2 \cdot x_2 + \dots + b_n \cdot x_n$

As above, **y** will be the Alere-chip entry for a gene with value 0 or 1.

The site, clin/com and age categories will be used as explanatory variables .

Steps of testing age categories for confounding

(1) Estimate a linear regression model for the dependent variable and one or more explanatory variables.

(2) Repeat step 1 with age categories added as further explanatory variable.

(3) Compare the weights obtained in steps 1 and 2.

As a rule of thumb, if the weight (-s) (regression coefficient(-s)) from step 1 changes by more than 10%, then the additional variable (here: age) may be considered as a **confounder**.

By following these steps, one can test for every significant finding (for example, gene association) whether age is a confounder.

Reasons for this could be e.g. a significant imbalance in the distribution of age among samples.

Case study

Case study: test whether age categories are a confounding variable for the 2 genes lukS.PV and sdrC..total

Reason for selecting these genes:

these 2 genes have very different frequencies in African vs German sites as well as in clinical vs commensal samples.

Africa was encoded as $x_1 = 0$ and Germany as $x_1 = 1$.

Clinical samples were encoded as $x_2 = 1$ and commensal with $x_2 = 0$.

Age categories were encoded from $x_3 = 1$ to 5.

Multiple linear regression model for the lukS.PV gene

The Alere result for this gene for different samples is the dependent variable. The site affiliation + clin/com values are explanatory variables.

The table lists the dependent (lukS.PV) and explanatory (Africa_value, clin_com_value) variables for 10 samples out of 1200 samples.

#	samples	lukS.PV	Africa_value	clin_com_value
1	FR-B001	0	0	1
2	FR-B003	0	0	1
3	FR-B004	0	0	1
4	FR-B005	0	0	1
5	FR-B007	0	0	1
6	FR-B008	0	0	1
7	FR-B009	0	0	1
8	FR-B010	0	0	1
9	FR-B011	0	0	1
10	FR-B012	0	0	1

Since all these samples are from a German site, the Africa_value = 0. Also, all samples are clinical (clin_com_value = 1).

lukS.PV

Application of linear regression determines optimal weights w_1 , w_2 , w_3 .

For every sample we get $lukS.PV = w_1 + w_2 \cdot Africa.value + w_3 \cdot clin com value .$

For the first sample FR-B001, the formula would be $0 = w_1 + w_2 \cdot 0 + w_3 \cdot 1$.

Results from multiple linear regression (coefficients marked in bold):

	Estimate	Std. Error	t value	Pr(> t)
w ₁ for intercept	-0.07250	0.01781	-4.070	5e-05 ***
w ₂ for Africa_value	0.42833	0.02057	20.825	<2e-16 ***
w ₃ for clin_com_value	0.19500	0.02057	9.481	<2e-16 ***

In other words, the following model is estimated: IukS.PV = -0.07 + 0.42833 · Africa_value + 0.195 · clin_com_value

t value : equal to coefficient (estimate) divided by the standard error. Pr(>|t|) : p-value = probability of seeing a result as extreme in random data.

lukS.PV

We then added a further variable "age category" with weight w_4 to the model.

lukS.PV = $w_1 + w_2 \cdot Africa.value + w_3 \cdot clin com value + w_4 \cdot age$

	Estimate	Std. Error	t	value Pr(> t)
(Intercept)	0.06211	0.04559	1.362	0.17333
Africa_value	0.39077	0.02360	16.556	< 2e-16 ***
clin_com_value	0.19470	0.02049	9.503	< 2e-16 ***
age	-0.03618	0.01129	-3.206	0.00138 **

lukS.PV = 0.06211 + 0.39077 · Africa value + 0.19470 · clin com value - 0.03618 · age

lukS.PV

This result shows

(a) that the age category has a very small impact (its own weight is close to 0) and(b) the two other weights (for the site and clin/com) did not change much.

E.g. the weight of the Africa_values changed in relative terms by :

 $\frac{(0.42833 - 0.39077)}{0.42833} \cdot 100\% = 8.8\%$

The weight of clin_com_value changed by only 0.15%.

Both values are smaller than 10% (rule of thumb).

Conclusion:

There is no statistical evidence that age acts as a confounding variable.

Same analysis for gene sdrC_total

Before adding age categories:

Coefficients:	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.02083	0.0125	81.711	< 2e-16 ***
Africa_value	-0.12833	0.0144	-8.896	< 2e-16 ***
clin_com_value	-0.05833	0.0144	-4.044	5.6e-05 ***

After adding age categories:

Coefficients:	Estimate Std. Error	t value	Pr(> t)
(Intercept)	0.975445 0.0321	30.407	< 2e-16 ***
Africa_value	-0.115667 0.0166	-6.964	5.44e-12 ***
clin_com_value	-0.058232 0.0144	-4.039	5.71e-05 ***
age-category	0.012198 0.0079	1.536	0.125

Weight of Africa_value changed by **9.87%**, weight of clin_com_value changed by **0.17%**

Conclusion

There is no evidence from our preliminary analysis for the genes lukS.PV and sdrC..total that age acts as a confounder in the association of genes with invasiveness and site affiliation.

We wrote in our manuscript:

"The discrepancy in population age between the German and African cohort potentially biases the 'true' distribution of clones and genes between isolates from the different geographic regions ...

[but] application of a multiple linear regression model for the detection rate of Panton-Valentine leucocidin genes failed to provide evidence that age acts as a confounding variable"

Ruffing et al. Sci. Rep. 7, 154 (2017)

Diabetes/HIV as confounding variables

Next, we tested using Fisher's exact test whether

(a) diabetes and HIV have similar frequencies in the total groups of African and German samples and

(b) whether diabetes and HIV have similar frequencies in selected groups of African and German individuals carrying particular clonal complexes.

The Fisher test considers the distribution provided in a 2×2 table.

	Africa	Germany	Row Total
HIV+	а	b	a + b
HIV-	С	d	c + d
Column Total	a + c	b + d	a + b + c+ d = n

The formula for the (exact) p-value calculation is :

$$p-value = \frac{(a+b)!(c+d)!(a+c)!(b+d)!}{a!b!c!d!n!}$$

Explanation: these are the number of possible combinatoric combinations for these fields.

Processing of Biological Data

Analysis of HIV co-infection

First, we will test the null hypothesis that "HIV is equally distributed in African and German samples".

(a) For all African samples and all German samples we obtain the following dependencies of HIV carriers (HIV+) and of individuals without approved HIV status (you may say non-carriers) (HIV-):

	Africa	Germany
HIV+	41	0
HIV-	315	586

The p-value obtained for this table can be interpreted as the sum of evidence provided by the observed data—or any more extreme table—for the null hypothesis that "there is no difference in the proportions of HIV carriers among the African and German individuals tested in our study".

The smaller the value of p, the greater the evidence for rejecting the null hypothesis.

Analysis of HIV co-infection

For the data shown above,

$$p-value = \frac{(41+0)!(315+586)!(41+315)!(0+586)!}{41!0!315!586!942!} = 1.03838e-18$$

Thus, there is very strong evidence from the observed frequencies that African and German individuals **are not equally likely to be HIV carriers**.

Analysis of diabetes co-infection

Similarly, we can obtain Fisher's exact p-value for the distribution of diabetes among African and German samples.

Africa Germany diab+ 4 68 diab- 475 526 p-value = 3.73425e-14

Also, here, the null hypothesis of a similar distribution is **strongly rejected** suggesting the prevalence of diabetes in individuals from Germany compared to individuals from Africa.

Of course, we can trace this imbalance back to the **difference in age categories** of the two groups.

HIV/diabetes in individuals with selected CCs

Next, we tested the distribution of HIV/diabetes in individuals carrying *S. aureus* from selected clonal complexes (CC15, CC45, CC121, CC30 which showed significant imbalance in German/African samples).

These are the results (tables + p-values from Fisher's exact test)

RF_HIV

CC15	Africa	Germany
hiv+	4	0
hiv-	65	57
p-value 0.126		0.25 (after correction for false discovery rate (FDR))
CC45	Africa	Germany
hiv+	1	0
hiv-	40	87
p-value	0.320	0.42 (FDR-corrected)

HIV/diabetes in individuals with selected CCs

CC121	Africa	Germany
hiv+	11	0
hiv-	40	24
p-value	0.0132	0.05 (FDR-corrected)

CC30	Africa	Germany
hiv+	0	0
hiv-	7	75
p-value	1	1 (FDR-corrected)

HIV/diabetes in individuals with selected CCs

RF_CCS	SI_Diab_	mel
CC15	Africa	Germany
diab+	0	1
diab-	88	56
p-value	0.393	0.52 (FDR-corrected)
CC45 diab+	Africa 0	Germany 12 75
p-value	47 0.0081	0.03 (FDR-corrected)
CC121 diab+ diab- p-value	Africa 0 57 0.305	Germany 1 24 0.52 (FDR-corrected)
CC30 diab+ diab- p-value	Africa 0 9 1	Germany 7 68 1 (FDR-corrected)
V12		Processing of Biological Data

Interpretation

In most cases, there is no evidence based on our data to reject the null hypothesis of assuming a similar distribution of HIV and diabetes carriers among African and German samples belonging to **particular clonal complexes**.

The only exceptions to this areCC45 (diabetes - p=0.008/q=0.03)andCC121 (HIV - borderline p=0.013/q=0.05).

Therefore, we concluded in the paper:

"we observed statistically significant imbalances in the frequencies of all these clonal complexes XXX, YYY ... between African and Germany.

We tested based on Fisher's exact test that these imbalances were not due to an imbalance of HIV and diabetes carriers in both groups.

The only exceptions to this are CC45 (diabetes) and CC121 (HIV) where such associations cannot be ruled out."

Interpretation

On the other hand, the FDR-corrected p-values for CC45 and CC121 are borderline (0.03 and 0.05).

Therefore, there only exists weak statistical evidence for a significant association between CC45 and diabetes or between CC121 and HIV.

integration of multi-omics data

Overview of methods for multivariate analysis:

https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-015-0857-9

Different types of high-throughput technologies allow us to collect information on the molecular components of biological systems

- e.g. nucleotide sequencing,
- DNA-chips measuring gene expression and
- protein mass spectrometry measuring protein abundances).

Therefore, in order to draw a more comprehensive view of biological processes, experimental data made on different layers have to be integrated and analyzed.

The development of methods for the integrative analysis of multi-layer datasets is one of the most relevant problems computational scientists are addressing nowadays.

Bersanelli et al. BMC Bioinformatics (2016) **17(Suppl 2)**:S15

Graph-based integration of multi-omics data

Some group of approaches use graphs to model the interactions among variables.

These approaches, designated as "network-based" (NB), take into account currently known (e.g. protein-protein interactions) or predicted (e.g. from correlation analysis) relationships between biological variables.

Then, graph measures (e.g. degree, connectivity, centrality) and graph algorithms (e.g. sub-network identification) are used to identify valuable biological information.

Importantly, networks are used in the modeling of the cell's intricate wiring diagram and suggest possible mechanisms of action at the basis of healthy and pathological phenotypes

Bersanelli et al. BMC Bioinformatics (2016) **17(Suppl 2)**:S15

Bayesian integration of multi-omics data

The second criterion is whether the approach is Bayesian (BY).

These approaches use a statistical model in which, starting from an *a priori* reasonable assumption about the data probability distribution (*parametric* or *non-parametric*)

it is possible to compute the updated posterior probability distribution making use of the Bayes' rule.

In the network-based area, Bayesian networks are another promising framework for the analysis multi-omics data.

4 classes of methods:

- network-free non-Bayesian (NF-NBY),
- network-free Bayesian (NF-BY),
- network-based non-Bayesian (NB-NBY) and
- network-based Bayesian (NB-BY) methods

Bersanelli et al. BMC Bioinformatics (2016) **17(Suppl 2)**:S15

Overview of multi-omics methods



Grey: network-free, non-Bayesian methods;

yellow: network-free, Bayesian methods;

blue: network-based, non-Bayesian methods;

green: network-based **Bayesian methods**

Method	Multi-omics approach	Impleme
Camelot	Bivariate predictive regression model	NA
CNAmet	Multi-omics gene-wise scores	R
FALDA	FA + LDA of a joint matrix	NA
Integromics	Regularized CCA, sparse PLS	R
iPAC	Sequential	NA
MCD	Sequential	NA
MCIA	Multiple co-inertia analysis	R
sMBPLS	Sparse Multi-Block PLS regression	Matlab
Coalesce	Multi-omics probabilities	C ++
iCluster	Joint Gaussian latent variable models	R
MDI	DMA mixture models	Matlab
PSDF	Hierarchical DMA mixture models	Matlab
TMD	Hierarchical DMA mixture models	Matlab
Kernel Fusion	Integration of omics-specific kernels	Matlab
Endeavour	Integration of omics-specific ranks with order statistics	Webserv
MOO	Sub-network extraction on MWG	R
Multiplex	Joint analysis of multi-layered networks	NA
NuChart	Analysis of a MWG	R
SNF	Similarity network fusion	Matlab, F
SteinerNet	Sub-network extraction on MWG	Webserv
stSVM	MWG	R
Paradigm	Multi-omics bayesian factor graphs	C ++
Conexic	Sequential	Java

Implementation

Existing tools

	legend
	MWG = multi-weighted
	graph;
	FA = factor analysis;
	LDA = linear discriminant
	analysis;
	CCA = canonical correlation
	analysis;
	PLS = partial least squares;
	DMA = Dirichelet
/er	multinomial allocation
D	
ver	
	Bersanelli et al. BMC Bioinformatics (2016) 17(Suppl 2) :S15

V12

Multi-omics analysis of breast cancer network



Processing of Biological Data

Breast cancer network from TCGA data

ca. 1300 differentially expressed genes.

Hierarchical clustering of coexpression network: 10 **modules** with 26 - 295 genes.

Regulatory info from databases Jaspar, Tred, MSigDB.

Shown are 3 modules.

Squares are known drug targets.



Drug Targets in breast cancer network

Table S4. The identified key gene nodes in the breast cancer network (12) whose protein products are targeted by anti-cancer drugs. (1) means that at least one drug that targets this gene product is reported in this database, and (0) means no drugs are reported for the respective gene in this database. Not included are substances that are known to be cancerogenous or mutagenic.

Target gene	Drug and antineoplastic agents		PharmGKB	Cancer Resource
AKT1	U 0126;tyrphostin AG 1478; Ursodeoxycholic Acid;Valproic Acid;tyrphostin AG 1024; trametinib; Tretinoin	1	0	1
BRCA2	Tretinoin; trichostatin A; Estradiol; transplatin; troglitazone; Tunicamycin; fulvestrant		0	1
ESR1	exemestane;tamoxifen		. 1	. 1
TGFB1	Doxorubicin; Fluorouracil; Thalidomide; Entinostat; Hyaluronidase		0	1
TP53	4-biphenylmine; alliin; Apigenin; Atropine;bicalutamide;butylidenephthalide		0	1

Some key genes are protein targets of known anti-cancer drugs,

 \rightarrow relevance of key genes is validated

Hamed et al. Nucl Acids Res 43: W283-W288 (2015)

Case study: single-cell analytics

http://www.nature.com/nmeth/journal/v13/n3/full/nmeth.3728.html

In the presence of serum, mouse embryonic stem cells (ESCs) epigenetic heterogeneity constitute a metastable population with stochastic switching between transcriptional states.

Parallel single-cell sequencing links transcriptional and

Christof Angermueller^{1,7}, Stephen J Clark^{2,7}, Heather J Lee^{2,3,7}, Jain C Macaulav^{3,7}, Mabel J Teng³, Tim Xiaoming Hu^{1,3,4}, Felix Krueger⁵, Sébastien A Smallwood², Chris P Ponting^{3,4}, Thierry Voet^{3,6}, Gavin Kelsey², Oliver Stegle¹ & Wolf Reik^{2,3}

This transcriptional heterogeneity has been linked to the differentiation potential of ESCs. E.g. NANOG¹⁰ cells have an increased propensity to differentiate and elevated expression of differentiation markers compared with NANOG^{hi} cells.

Sorted populations of cells show different levels of DNA methylation between transcriptional states, such as gains in DNA methylation in NANOG¹⁰ and REX1/ZFP42^{lo} cells compared with, respectively, NANOG^{hi} and REX1^{hi} cells.

To investigate the link between epigenetic and transcriptional heterogeneity in ESCs, Reik et al. performed scM&T-seq on 76 individual serum ESCs.

scMT & T-seq protocol

Single cells are collected and lysed.

Then poly-A RNA is captured on magnetic beads and physically separated from DNA.

Single cell isolation 99 Lysis SMARTer Oligo dT-VN Streptavidin magnetic bead Poly-A mRNA capture AAAAAAA Separation of poly-A mRNA and DNA 9E)

Angermüller et al. Nature Methods (2016) 13, 229

scMT & T-seq protocol

Amplified cDNA is generated from mRNA on beads.

DNA is bisulfite converted and Illumina sequencing libraries are prepared from both components in parallel.



Angermüller et al. Nature Methods (2016) 13, 229

How good is the protocol: check against single cell bisulfite sequencing (scBS-seq)



Methylome coverage in scM&T-seq libraries was lower than that in scBS-seq libraries.

However, genome-wide CpG coverage at matched sequencing depth (c) and coverage of different regions (d) was consistent across protocols.



Angermüller et al. Nature Methods (2016) 13, 229

V12

Clustering based on DNA methylation data



Shown is a hierarchical clustering analysis of gene-body methylation for the 300 most variable genes in terms of DNA methylation.

 \rightarrow 3 main clusters (green, mauve, orange)

Angermüller et al. Nature Methods (2016) 13, 229

gene.-body methylation levels

е

Clustering based on expression data



Shown is a hierarchical clustering analysis of gene expression for the 300 most variable genes (on the basis of DNA-methylation variance (same coloring scheme)).

 \rightarrow Both data yield distinct clustering of cells.

This suggests that global methylome and transcriptome profiles reveal complementary, but distinct, aspects of cell state.

This is also consistent with previous observations that the transcriptome and methylome are partially uncoupled in serum ESCs.

Angermüller et al. Nature Methods (2016) 13, 229

V12

Associations of expression and DNA-methylation variation



1,493 associations were found between the expression of individual genes and DNA-methylation variation in several genomic contexts (FDR) < 10%).

There exist **both positive and negative associations**, highlighting the complexity of interactions between the methylome and the transcriptome.

Pearson correlation of methylation and expression of individual genes

Also distal regulatory elements including low-methylation regions (LMRs) had a fair balance of positive and negative

associations.



Angermüller et al. Nature Methods (2016) 13, 229

Processing of Biological Data





Negative correlations between DNA methylation and gene expression were predominant for non-CGI promoters.

Angermüller et al. Nature Methods (2016) 13, 229

V12

Zoomed-in analysis for Esrrb



Dot size : CpG coverage, Dot colors correspond to single cells.

Angermüller et al. Nature Methods (2016) 13, 229

V12

Esrrb is a known hub gene in pluripotency networks. Its expression negatively correlates with the methylation of several LMR and p300 sites overlapping 'super-enhancers' in the genomic neighborhood.

Zoomed-in analysis for Esrrb

Two scatter plots at KII2 Zlp42 the top right: association between DNA 25 50 Methylation (%) 1.0 -0.5 0.5 0 25 50 Methylation (%) 75 75 100 methylation at a Variance 1,000 p300 region 500 (outlined in yellow) 1.0 -Correlation (r) and at an LMR 0.5 -0. -0.5 (outlined in blue) -1.0 100 and Esrrb 75 expression. Methylation (%) 50 25 0 Esrrb LMR p300 L 1 I I 1 Super-enhancer CGI . . 86,350,000 86,400,000 86,450,000 86,500,000 Esrrb GRCm38 chr12: 86,361,117 - 86,521,628 (161 kb) TSS

Angermüller et al. Nature Methods (2016) 13, 229

V12

Processing of Biological Data

100

Correlations of DNA methylation and expression



Gene-specific association analysis of correlations between DNA methylation in different genomic contexts and gene expression in individual cells.

Shown are methylation-expression correlations for all variable genes in single cells, for each annotation, with the correlation obtained from matched RNA-seq and BS-seq of a bulk cell population superimposed (orange circles). Prom = promoter.

Angermüller et al. Nature Methods (2016) 13, 229

V12

Summary

- Multi-variate vs. single-variate analysis reveals possible confounding effects
- Multi-omics methods: graph-based and/or Bayesian methods for data integration
- Single cell analysis showed:
- non–CGI promoter methylation and transcription are negatively associated in single cells / both positive and negative associations at distal regulatory regions.
- expression levels of many pluripotency factors, such as *Esrrb* are negatively associated with DNA methylation → an important mechanistic component of fluctuating pluripotency in serum ESCs is epigenetic heterogeneity
- the strength of the connection between the methylome and the transcriptome can vary from cell to cell
- Q: is our understanding / data generation ready for multi-omics analysis?

Tools for meta-dimensional analysis



a | Concatenation-based integration involves combining data sets from different data types at the raw or processed data level before modelling and analysis.

V12

Ritchie et al. *Nature Reviews Genetics* **16**, 85–97 (2015) Processing of Biological Data

b | Transformation-based integration involves performing mapping or data transformation of the underlying data sets before analysis, and the modelling approach is applied at the level of transformed matrices. c | Model-based integration is the process of performing analysis on each data type independently, followed by integration of the resultant models to generate knowledge about the trait of interest.