**Proteome and Genome Research Unit**

# Non-negative Matrix Factorization for Methylation Data Deconvolution

## Update on the Data Challenge and the related paper: Lutsik et al

**Petr Nazarov**

petr.nazarov@lih.lu

**LIH, Strassen, Luxembourg**

# Outline

- **Concept of NMF**

- **Update on the Data Challenge**

- **Paper**

**METHOD** **Open Access**

# MeDeCom: discovery and quantification of latent components of heterogeneous methylomes
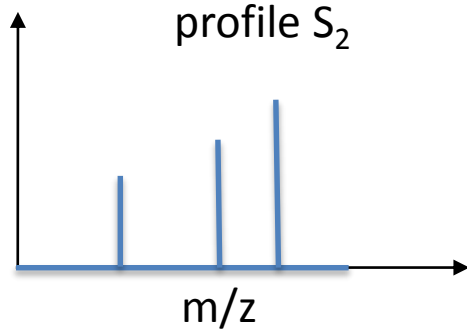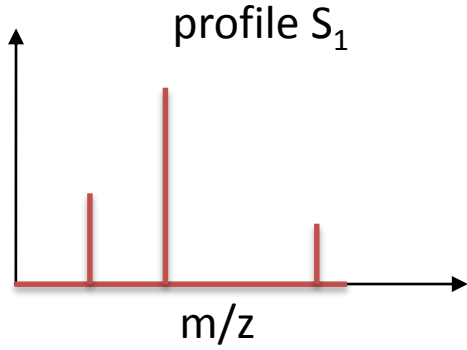
Pavlo Lutsik[1,4†], Martin Slawski[2,3,5†], Gilles Gasparoni[1], Nikita Vedeneev[2], Matthias Hein[2*] and Jörn Walter[1*]

**Abstract**

It is important for large-scale epigenomic studies to determine and explore the nature of hidden confounding variation, most importantly cell composition. We developed MeDeCom as a novel reference-free computational framework that allows the decomposition of complex DNA methylomes into latent methylation components and their proportions in each sample. MeDeCom is based on constrained non-negative matrix factorization with a new biologically motivated regularization function. It accurately recovers cell-type-specific latent methylation components and their proportions. MeDeCom is a new unsupervised tool for the exploratory study of the major sources of methylation variation, which should lead to a deeper understanding and better biological interpretation.
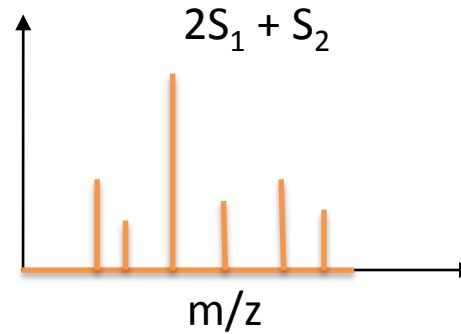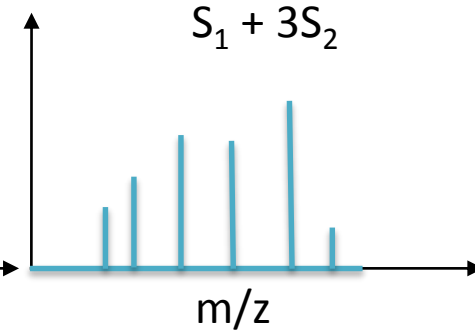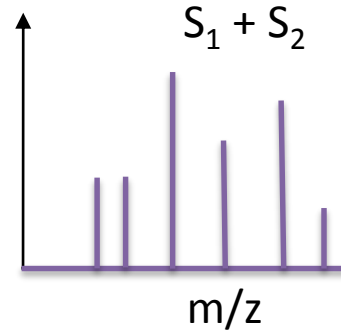
**Keywords:** DNA methylation, DNA methylome, Cell heterogeneity, Deconvolution, Matrix factorization, Epigenetics

# Concept: NMF

profile $S_1$

m/z

profile $S_2$

m/z

Mixing

$$X = S \times M$$

$$M = \begin{array}{|c|c|c|} \hline 1 & 2 & 1 \\ \hline 1 & 1 & 3 \\ \hline \end{array}$$

$S_1 + S_2$

m/z

$S_1 + 3S_2$

m/z

$2S_1 + S_2$

m/z

**Non-negative**          **Non-negative**          **Non-negative**

$$X \approx T \times A$$

~S        ~M
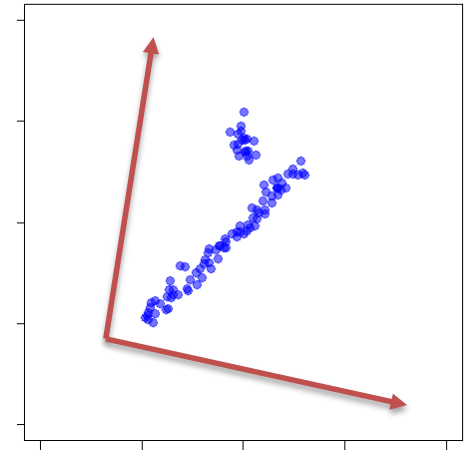
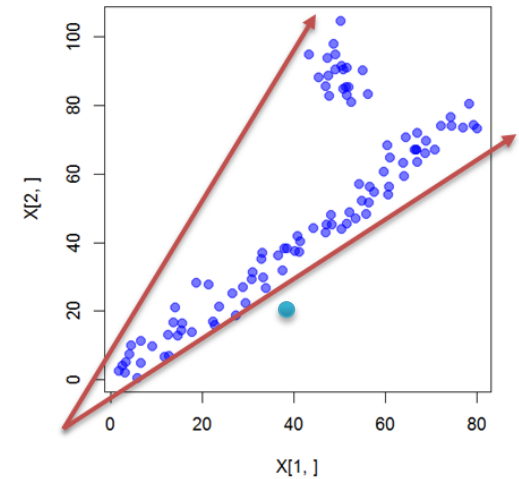# Concept: NMF

**PCA**



**NMF**



**NMF: issue 1**



**NMF: issue 2**



- Multiple solutions
- Is the minimal description stable?

$\Rightarrow$ we need:
➢ additional restrictions
➢ regularizations during fitting

# Data Challenge

## Place & Participants



## Aussois, Dec 2018

**Invited speakers:**

➤ E. Andres **Houseman**, independent data scientist, USA, **RefFreeEWAS**

➤ Pavlo **Lutsik**, from DKFZ, Heildeberg, Germany, **MeDeCom**

➤ Eugene **Lurie**, from BCM, Houston, USA, **Edec**

**Participants:**

➤ 9 (10) commands 3-4 members: FR, DE, US, RU, LU, NL, … ?

# Data Challenge

## Structures

**2 sub-challenges:**    **X = S x M + noise**

➢ **Training**: 3 cell types, 100 synthetic samples, no confounding variables
➢ **Main**: $k$ cell types ($k$=5), 100 synthetic samples, $y$ confounding variables

**Our team:**

➢ Fabian Bergmann (MSc student) – IT, submits, fine tuning, RefFreeEWAS
➢ Tony Kaoma – wide search for alternative algorithms, MeDeCom
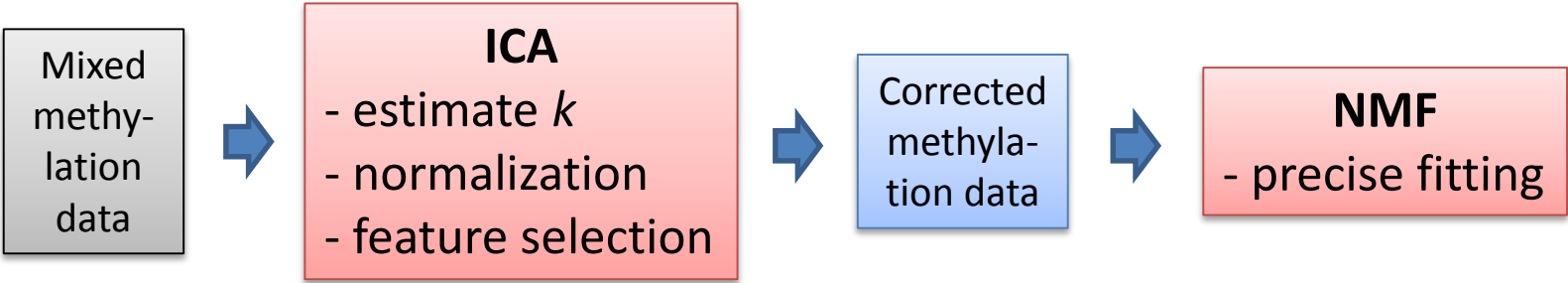➢ Petr Nazarov – ICA, moderating FB ☺

**Reasons why we won sub-challenge 1:**

➢ our tuning of the parameters was more efficient – RefFreeEWAS overfits!
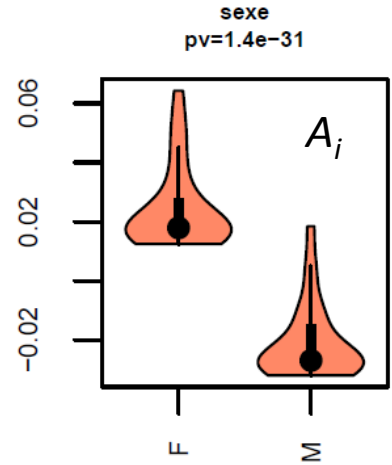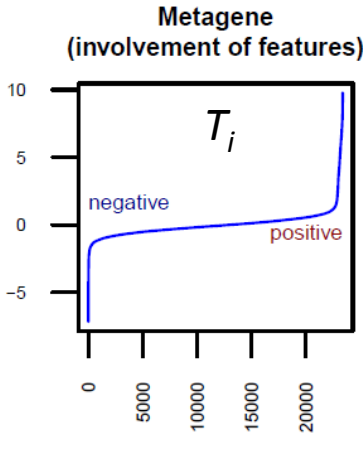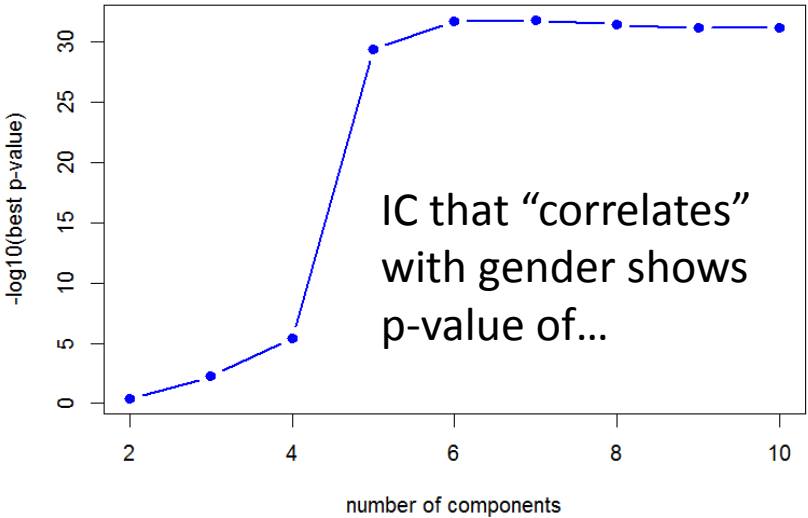➢ search for methods for initial estimation by TK helped

*...and ICA was not needed at al ☺ !*

## Winning strategy for sub-challenge 2

Mixed methy-lation data → **ICA**
- estimate $k$
- normalization
- feature selection → Corrected methyla-tion data → **NMF**
- precise fitting

Gender was one of the confounders



IC that "correlates" with gender shows p-value of…

Metagene (involvement of features)

$T_i$

negative
positive

sexe
pv=1.4e-31

$A_i$

F    M

**X = T x A**

# Data Challenge

## Teams



Pavlo Lutsik: MeDeCom

Eugene Lurie: Edec

Almost the only biologist ☺

# Data Challenge

## Conclusions & Results

➢ Even the top methods on NMF must be fine-tuned in order to give good results

➢ ICA is of no help for simple tasks. But can be useful in more complex situations (e.g. confounders)

➢ Pavlo Lutsik, the developer of MeDeCom and co-author of RnBeads, was "as**toni**shed" by Tony's results on his own tool and proposed to work together on the protocols paper

➢ The general paper based on the challenge was planned, but i.m.h.o., the chances are vague

➢ It brings new knowledge and simply… a lot of fun ☺

E.Andreas Houseman: RefFreeEWAS

**METHOD**　　　　　　　　　　　　　　**Open Access**

CrossMark

# MeDeCom: discovery and quantification of latent components of heterogeneous methylomes

Pavlo Lutsik[1,4†], Martin Slawski[2,3,5†], Gilles Gasparoni[1], Nikita Vedeneev[2], Matthias Hein[2*] and Jörn Walter[1*]

**Abstract**

It is important for large-scale epigenomic studies to determine and explore the nature of hidden confounding variation, most importantly cell composition. We developed MeDeCom as a novel reference-free computational framework that allows the decomposition of complex DNA methylomes into latent methylation components and their proportions in each sample. MeDeCom is based on constrained non-negative matrix factorization with a new biologically motivated regularization function. It accurately recovers cell-type-specific latent methylation components and their proportions. MeDeCom is a new unsupervised tool for the exploratory study of the major sources of methylation variation, which should lead to a deeper understanding and better biological interpretation.

**Keywords:** DNA methylation, DNA methylome, Cell heterogeneity, Deconvolution, Matrix factorization, Epigenetics

# MeDeCom Paper

**Main idea (sorry, it is simple but… ☺)**

**D = T×A + E**

**Standard NMF:**

$$\min_{T,A} ||D - TA||_F^2 = \sum_{i=1}^m \sum_{j=1}^n (D_{ij} - (TA)_{ij})^2$$
$$\text{subject to}$$
$$0 \leq T_{is} \leq 1 \ \forall i, s$$
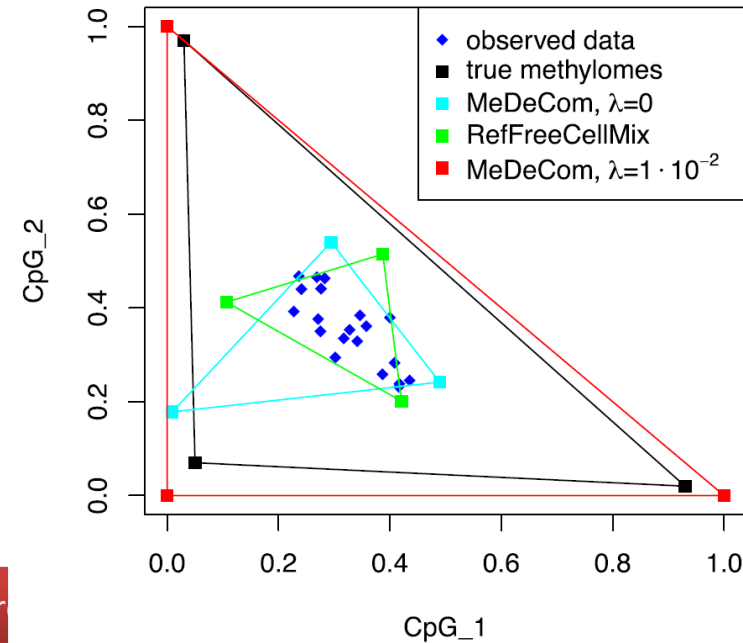$$A_{sj} \geq 0 \ \forall s, j$$
$$\sum_{s=1}^k A_{sj} = 1 \ \forall j.$$

**MeDeCom's regularization:**

$$\min_{T,A} ||D - TA||_F^2 + \boxed{\lambda \sum_{i=1}^m \sum_{s=1}^k \omega(T_{is})}, \text{with } \omega(x) = x(1-x)$$

$$\text{subject to } 0 \leq T_{is} \leq 1 \ \forall i, s$$
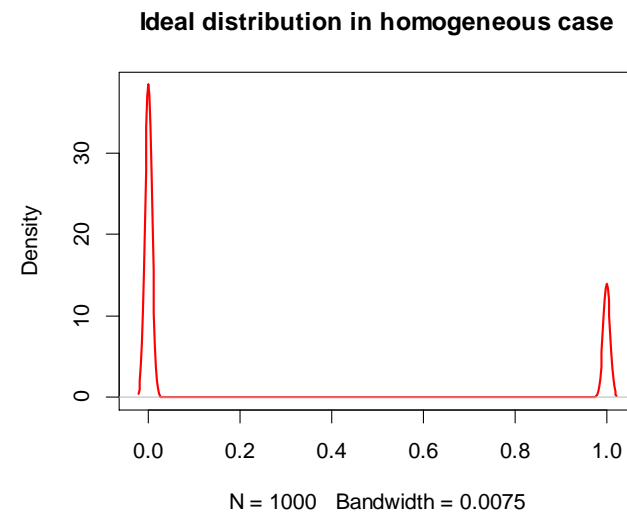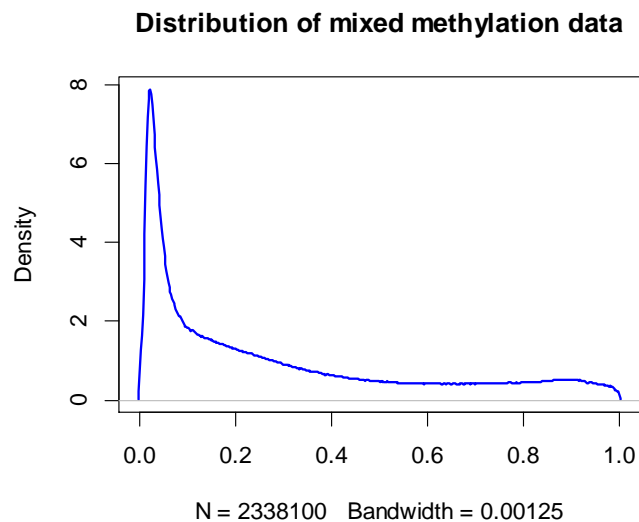$$A_{sj} \geq 0 \ \forall s, j$$
$$\sum_{s=1}^k A_{sj} = 1 \ \forall j,$$



Legend:
- ◆ observed data
- ■ true methylomes
- ■ MeDeCom, λ=0
- ■ RefFreeCellMix
- ■ MeDeCom, λ=1·10⁻²

## Assumptions / Requirments
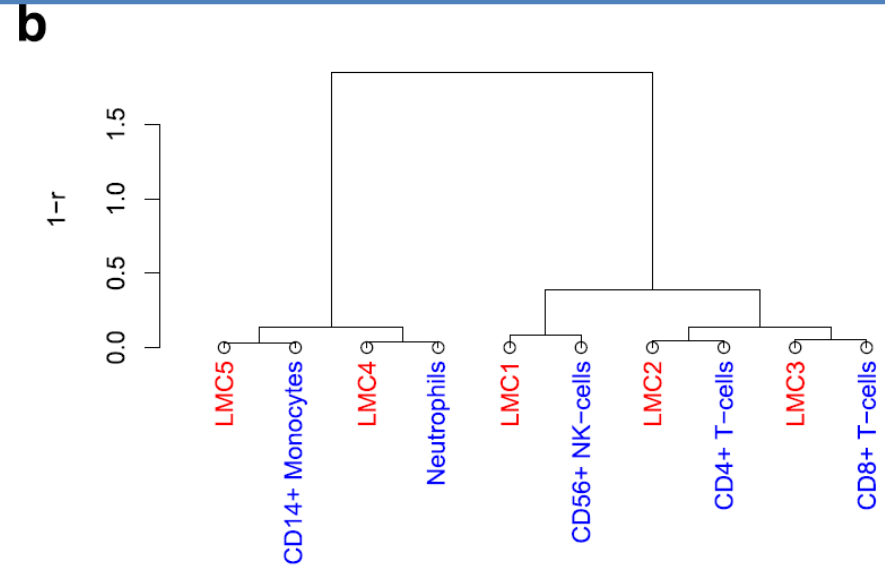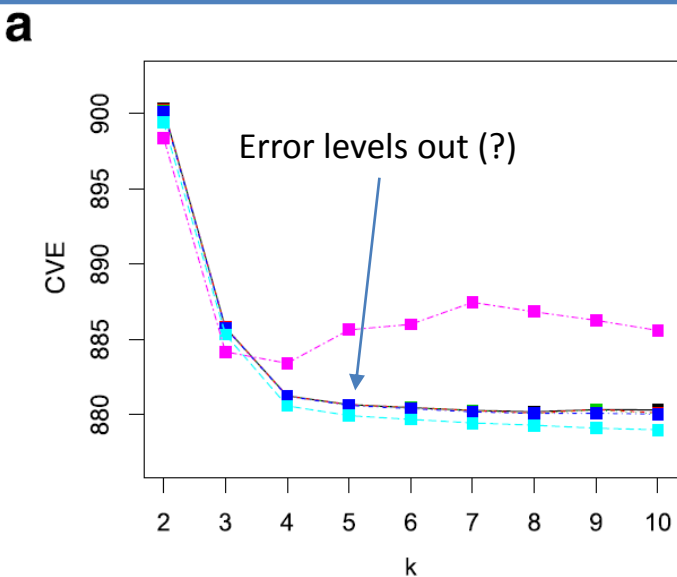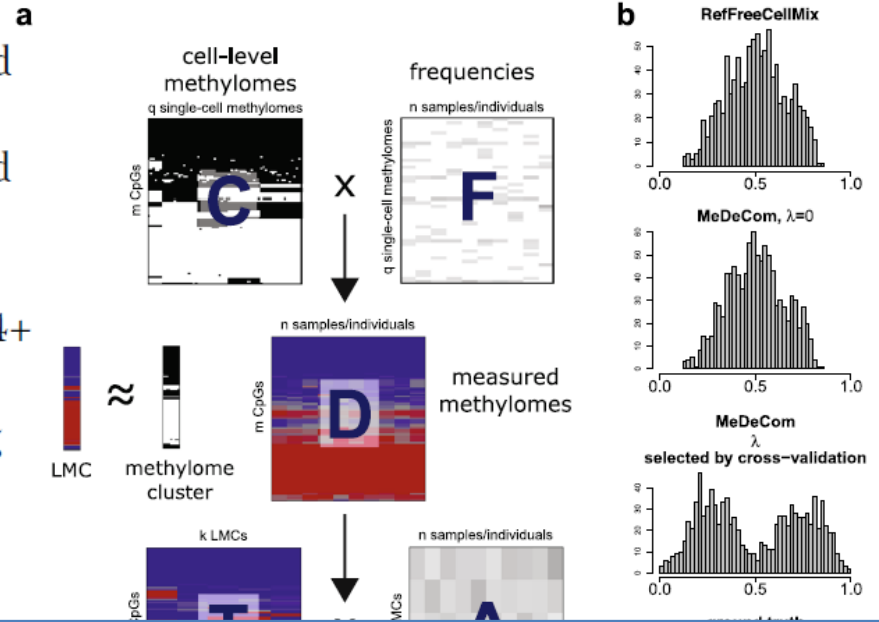
(1) Cell population consists of finite (and small) number of sub-populations.

(2) Each cell subpopulation have homogenous methylome profile => $\forall$CpG can be either 0 or 1.

(3) Population mixtures are variable b/w samples.

(4) Low level of technical noise and high level of biological variability.

**Distribution of mixed methylation data**



N = 2338100   Bandwidth = 0.00125

**Ideal distribution in homogeneous case**



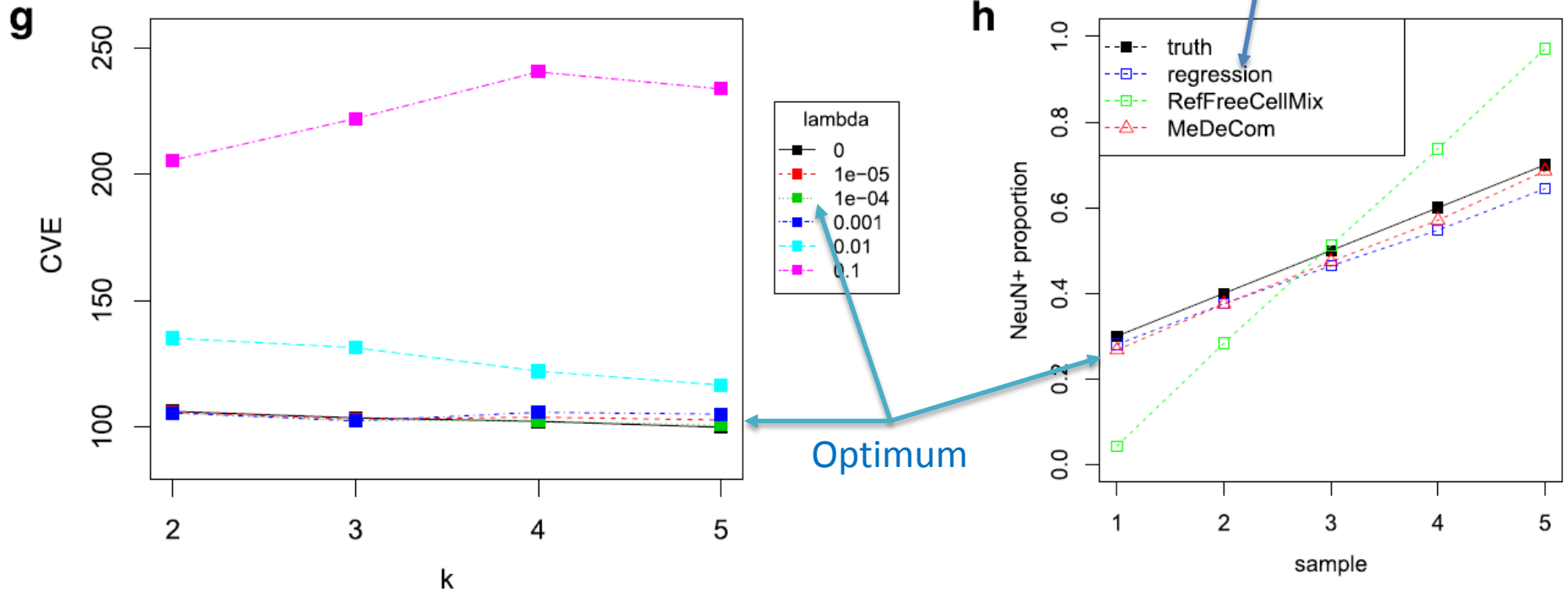N = 1000   Bandwidth = 0.0075

## Synthetic data

- $k_{sim} = 2$ with two distant cell types (neutrophils and CD4+ T cells).
- $k_{sim} = 2$ with two similar cell types (neutrophils and monocytes).
- $k_{sim} = 3$ with two similar cell types and one distant from the first two (neutrophils, monocytes and CD4+ T cells).
- $k_{sim} = 5$ with all major blood cell types, excluding eosinophils and B cells.

**a**

cell-level methylomes
q single-cell methylomes

frequencies
n samples/individuals
q single-cell methylomes

C X F

LMC ≈ methylome cluster

measured methylomes
n samples/individuals

D

k LMCs
n samples/individuals

T A

**b**

RefFreeCellMix

MeDeCom, λ=0

MeDeCom λ selected by cross-validation

**a**

Error levels out (?)

| lambda | |
|---|---|
| 0 | ■ |
| 1e−05 | ■ |
| 1e−04 | ■ |
| 0.001 | ■ |
| 0.01 | ■ |
| 0.1 | ■ |

CVE vs k

**b**

1−r

LMC5 — CD14+ Monocytes — LMC4 — Neutrophils — LMC1 — CD56+ NK-cells — LMC2 — CD4+ T-cells — LMC3 — CD8+ T-cells

## Cell mixture

**Dataset ArtMixN**: cell sorting into NeuN+ and NeuN- cells.

*NeuN = RBFOX3 protein*



Regression – if S matrix is known

Optimum
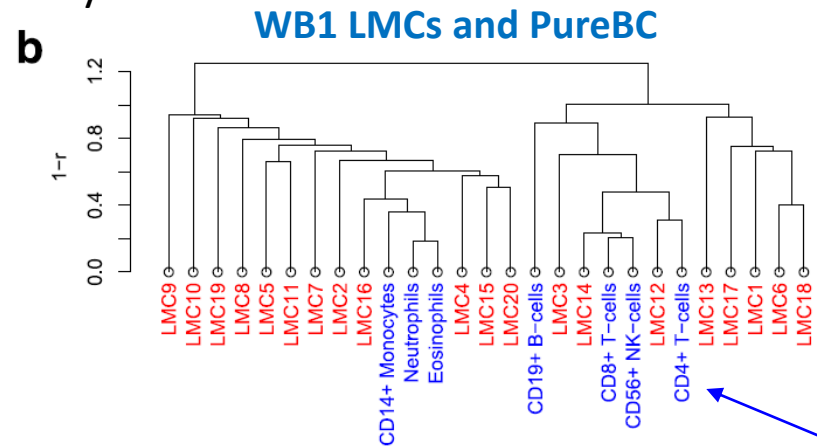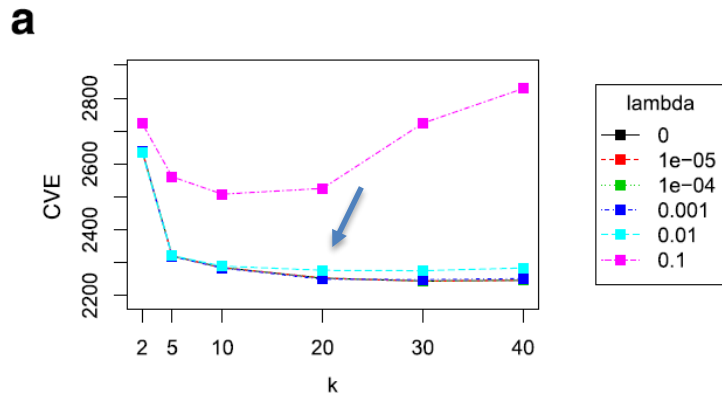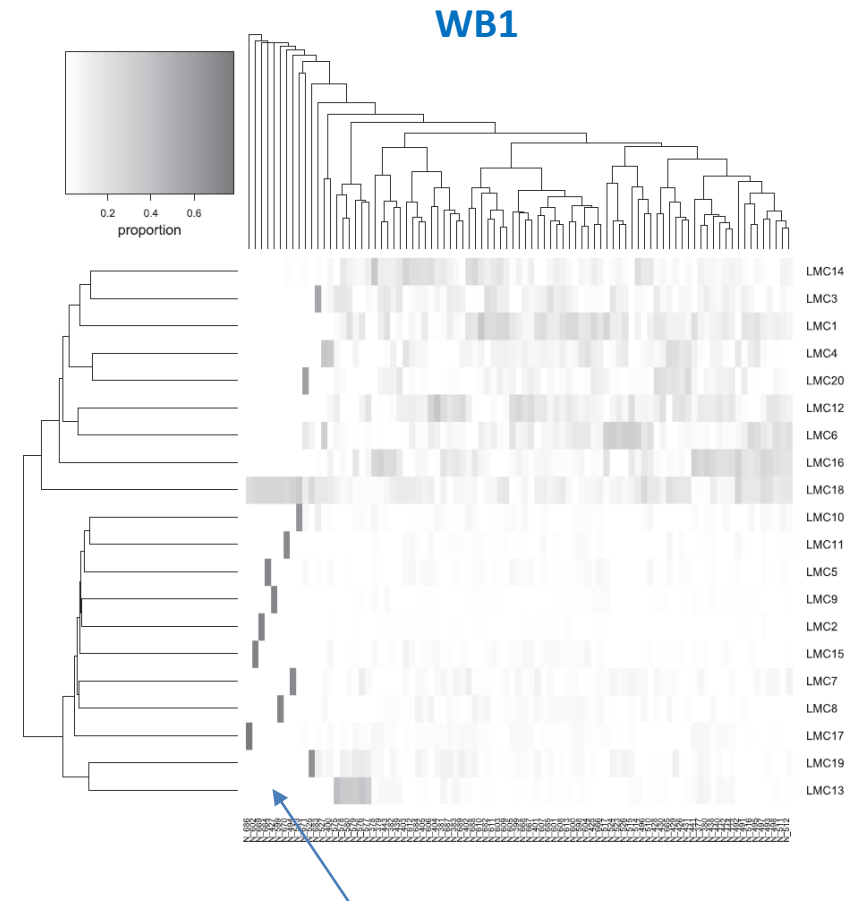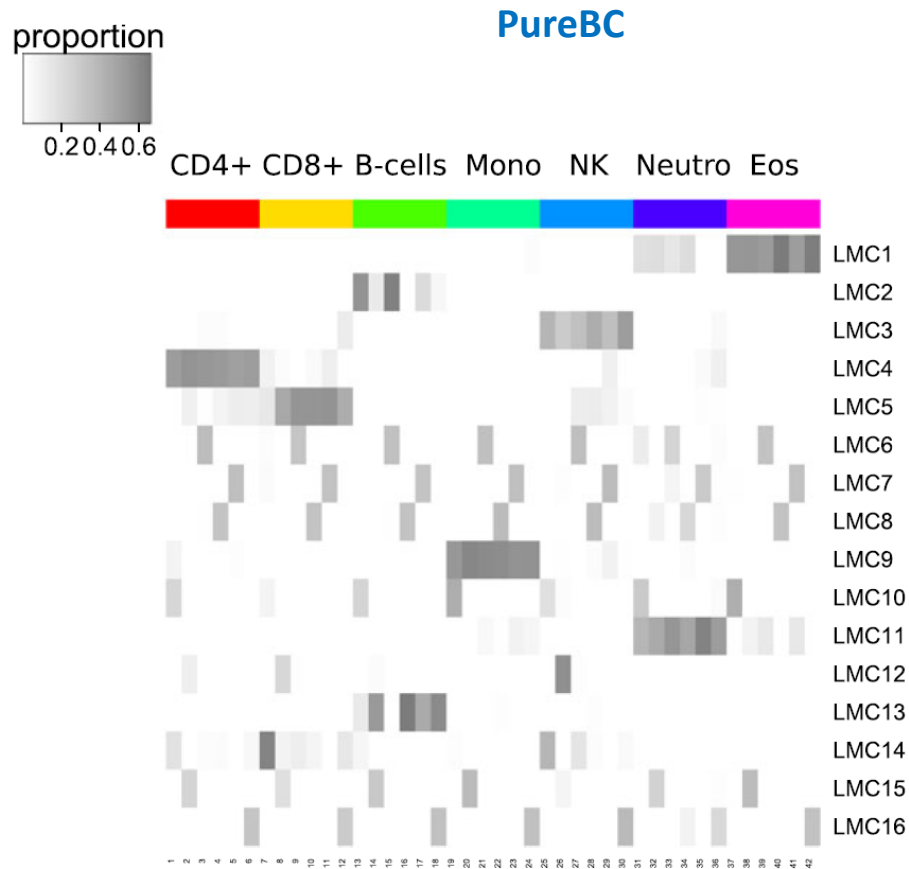
## Blood samples

Whole blood samples were used.

- **PureBC** – 7 MACS-purified cell types: neutro, mono, B cells ,CD4+, CD8+, NK, eosinophils
- WB1 – 87 rheumatoid arthritis patients
- WB2 – 442 cancer-free patients from EPIC Italy

## PureBC samples



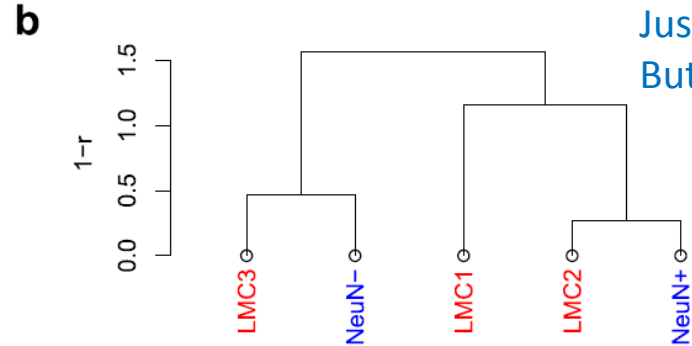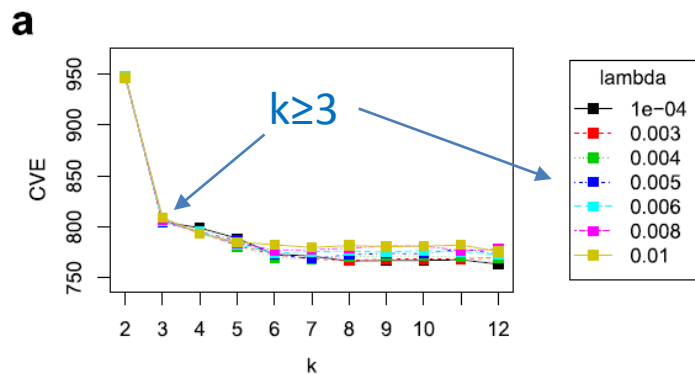Specific to a single donor? Mutations?

# MeDeCom Paper

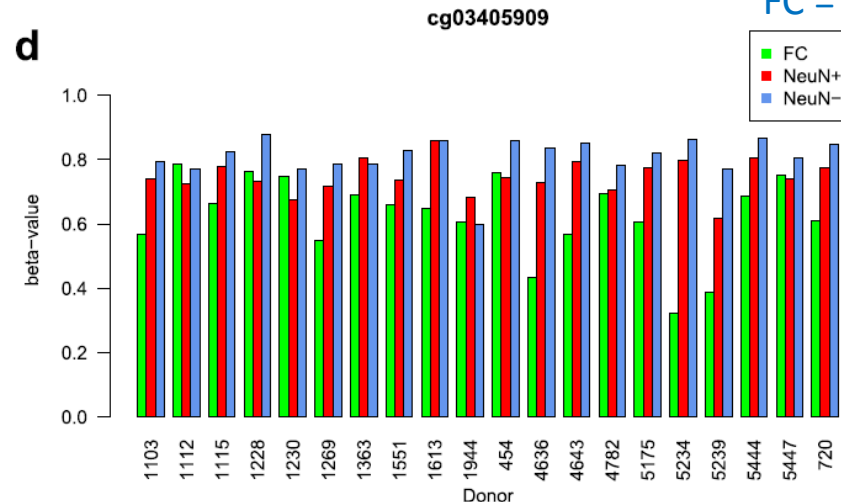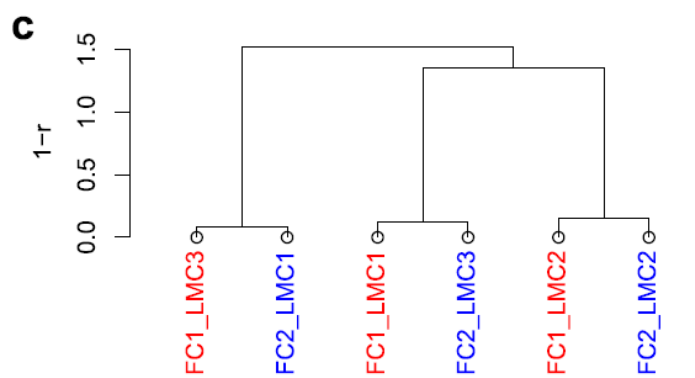## Brain samples

Whole blood samples were used.

- **PureN** – 2x29 NeuN+/- fractions of 29 healthy controls
- FC1 – 2x10 frontal cortex of MDD (major depression disorder) patients
- FC1 – 114 frontal cortex of AD (Alzheimer's disease) patients

Why k=3 ?
Just to have stable result...
But it is not enough.

frontal cortex
FC = NeuN$^+$ + NeuN$^-$

# MeDeCom Paper

## Conclusions

MeDeCom
(1) provides significant advances compared to other methods;
(2) uses with biologically relevant constrains and its LMCs are more interpretable;
(3) acts robustly on artificial and real data;
(4) identifies key methylation signatures;

(5) (IMHO) Low k is more dangerous than high k.

But: Separation of specific blood cell subtypes (similar methylomes) becomes challenging