## Deconvolution

Lou Shaoke
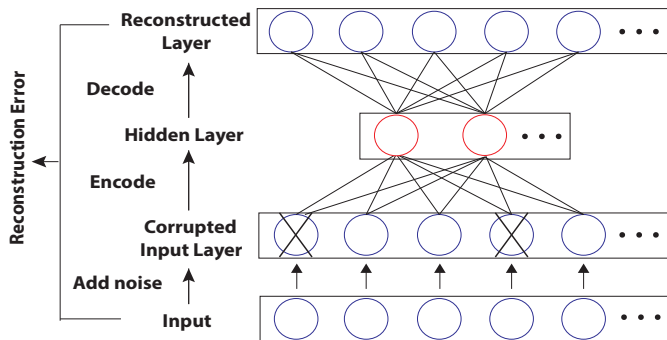
Department of Molecular Biophysics and Biochemistry *loushaoke@gmail.com*

August 12, 2015

Yale

Diagram



$y = sigmod(Wx + b)$
$z = sigmod(W^T y + b\prime)$
$L = \sum(-xlog(z) - (1-x)log(1-z))$
Then use stochastic gradient descent to find local minimum.
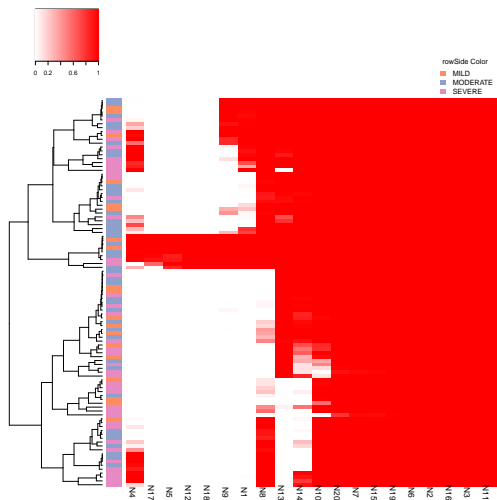$\#visible(v) = \#gene, \#hidden(h) \in 20, 30, 50, \ldots$, learning rate=0.01, corruption
level=0.05; cycle $= 100$

Weight matrix(W, $h \times v$), represent the contribution of each gene for each node. The hidden value(Y, $hs$), can be thought as the activity value of each node in each sample(s)
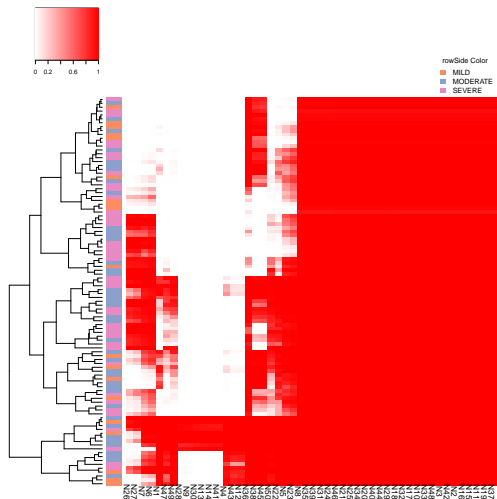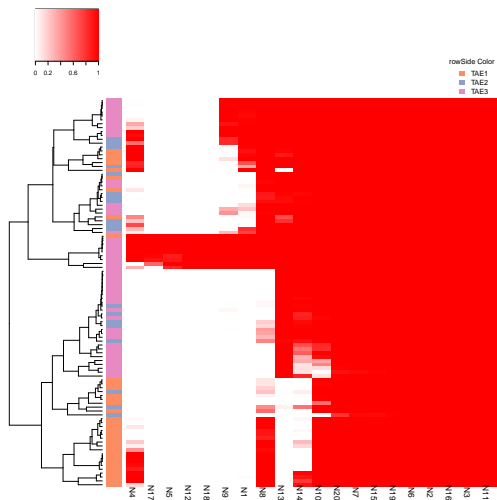
100 cycle, 20 hidden node,learning rate 0.2

100 cycle, 50 hidden node,learning rate 0.1 and corruption level 0.01

100 cycle, 20 hidden node,learning rate 0.2

100 cycle, 50 hidden node,learning rate 0.1 and corruption level 0.01

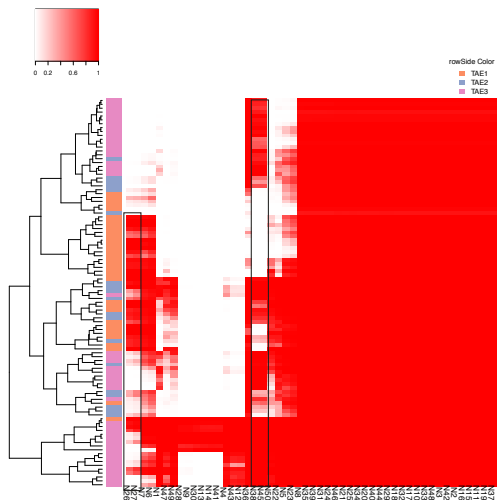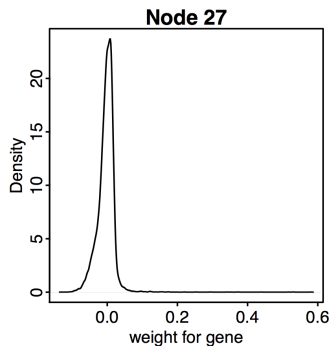We define genes contribution into three classes: extreme_high(H), normal, extreme_low(L) by $\mu \pm 2\sigma$

| xx | node26 | node27 | both |
|----|--------|--------|------|
| H | 232 | 246 | 232 |
| L | 713 | 681 | 676 |

The same way, we characterize node 38,45, 50.

Yale

# Function annotation

1. Highly Enrichment in keratinization for TAE1 specific high weighted gene. Asthma associate with keratinization.
2. Fatty Acid related process
3. Hormone response related process

Yale

Clustering using genes expression from above nodes

Clustering using random selected genes

1. purity
   $Purity(\Omega, C) = \frac{1}{N} \sum_k max_j |\omega_k \cap c_j|$

2. normalized mutual information
   $NMC(\Omega, C) = \frac{I(\Omega;C)}{(H(\Omega)+H(C))/2}$

3. Rand Index:

4. F-value:

Yale

Eosinophilic and Neutrophilic Inflammation in Asthma. eosinophilic and neutrophilic inflammation in asthma, and they are not mutually exclusive subtype.

Neutrophils are prominent in airway secretions during acute severe asthma

Macrophages exert prominent effects in the defense of the respiratory tract from airborne pathogens

Distinct cellular subtypes of asthma based on the presence or absence of sputum granulocytesnamely, eosinophilic asthma (EA), neutrophilic asthma (NA), mixed eosinophilic and neutrophilic asthma (ME/NA) and paucigranulocytic asthma (PGA)

Yale

1. 5 cell types, all sample: control vs astham
2. 5 cell types, samples: control vs severe asthma
3. 3 cell types, all sample
4. 3 cell type, control versus severe asthma

Given FDR<0.05, all above comparison have no significant genes, if FDR ¡ 0.1, the last comparison finds :DEFA3, DEFA1 and RPS4Y1(robosomal protein) Microphase cell line.

Yale

# Followup

- Expression value, log-based?

- Conventionally, use cell count to determine subtypes, know more about the experiment design

- other medical information: disease duration, medication etc

- cellular changes $<$ cell proportion changes, inflammatory response, proportion or count change is the best and easiest way to do it.

- meta analysis and supervised learning

Sources

Mixtures

Separated
Sources

http://research.ics.aalto.fi/ica/cocktail/cocktail_en.cgi

Yale

Given a set of mixed gene expression sets $X_{gs}$ for gene $g \in 1, 2, ..., G$ and sample $s \in 1, 2, ..., m$. The samples can be from case-control tissue, different tissues type and blood sample etc. Due to the sample hetergeneous, the gene expression should be a mixture of expression of different cell type/condition $w \in 1, 2, ..., W$.

The motivation:

Given a set of mixed gene expression sets $X_{gs}$ for gene $g \in 1, 2, ..., G$ and sample $s \in 1, 2, ..., m$. The samples can be from case-control tissue, different tissues type and blood sample etc. Due to the sample hetergeneous, the gene expression should be a mixture of expression of different cell type/condition $w \in 1, 2, ..., W$.

The motivation:

1) Can we deconvolve the expression to cell-type specific expr? csSAM, PERT

Given a set of mixed gene expression sets $X_{gs}$ for gene $g \in 1, 2, ..., G$ and sample $s \in 1, 2, ..., m$. The samples can be from case-control tissue, different tissues type and blood sample etc. Due to the sample hetergeneous, the gene expression should be a mixture of expression of different cell type/condition $w \in 1, 2, ..., W$.

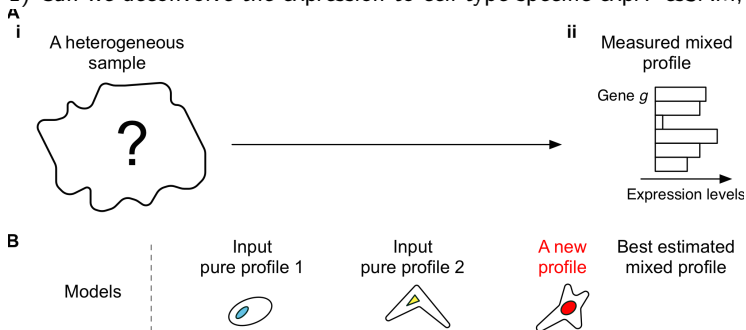The motivation:
1) Can we deconvolve the expression to cell-type specific expr? csSAM, PERT
2) Can we deconvolve the expression to cell-type like value/latent? for example: cancer, control;(DeMix, ISOpure etc)

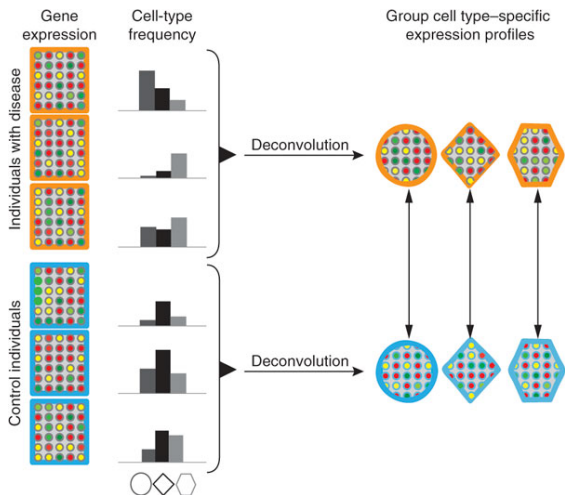Yale

$X = A \times S$

# Linear Deconvolution

$X = A \times S$

SS Shen-Orr (Nature method 2010): Two-group Model

$X_{ij} = \sum_{k=1}^{K} w_{ik} h_{kj}^{(1)} + e_{ij}$ and $X_{ij} = \sum_{k=1}^{K} w_{ik} h_{kj}^{(2)} + e_{ij}$

## Algorithms available in CellMix

The *CellMix* package includes several deconvolution algorithms, which differ in term of input and output data. The following table helps choosing an appropriate algorithm according to the data available and the desired output.

| Description | Basis | Coef | Marker | Iter |
|---|---|---|---|---|
| lsfit Partial deconvolution of proportions using least-squares fits (Abbas et al. (2009)) | red | green | - | - |
| cs-lsfit Partial deconvolution of cell signatures using least-squares fits | green | red | - | - |
| qprog Estimates proportions from known expression signatures using quadratic programming (Gong et al. (2011)) | red | green | - | - |
| cs-qprog Estimates constrained cell-specific signatures from proportions using quadratic programming [experimental] | green | red | - | - |
| DSA Complete deconvolution using Digital Sorting Algorithm (Zhong et al. (2013)) | green | green | red | - |
| csSAM Estimates cell/tissue specific signatures from known proportions using SAM (Shen-Orr et al. (2010)) | green | red | - | - |
| DSection Estimates proportions from proportions priors using MCMC (Erkkila et al. (2010)) | green | red/green | - | 500 |
| ssKL Semi-supervised NMF algorithm for KL divergence, using marker genes (Gaujoux et al. (2011)) | green | green | red | 3000 |
| ssFrobenius Semi-supervised NMF algorithm for Euclidean distance, using marker genes (Gaujoux et al. (2011)) | green | green | red | 3000 |
| meanProfile Compute proportion proxies as mean expression profiles | - | green | red | - |
| deconf Alternate least-square NMF method, using heuristic constraints (Repsilber et al. (2010)) | green | green | - | 1000 |

🟥 Required input  🟩 Estimated output  🟥🟩 Required input and estimated output

**Basis**  Cell-specific signatures
**Coef**  Cell proportions
**Marker** *Input:* cell-specific marker list
        *Output:* cell-specific differential expression (e.g., Case vs. Control)

## Other algorithms not – yet – available in CellMix

- TEMT: A mixture model for expression deconvolution from RNA-seq in heterogeneous tissues (Li et al. (2013))
- DeMix: Deconvolution for Mixed Cancer Transcriptomes Using Raw Measured Data (Ahn et al. (2013))
- ISOpure: Computational purification of individual tumor gene expression profiles leads to significant improvements in prognostic prediction

- 1. technical reasons and data transformation. Yi Zhong et al 2012 response to SS Shen-Orr.
- 2. Theoretical and pratical. How to evelute the results? It is good if more DEGs were found?
  Celltype-wise seperation?
  blind expression seperation? especially for more complex situation. such as metatstasis tissue with adjacent and original tissue.

Yale

- - use known algorithms to explore functional
- - From the practical view: diagnosis and prognosis
  blood test: marker and diagnosis (require clinical information)
- - Metastasis
  Seed and soil
- - Combination?
  Origin site $\rightarrow$ blood $\rightarrow$ Metastasis

Yale