# Tackling the topology and geometry underlying big data

## Monica Nicolau

Department of Mathematics

Department of Microbiology Immunology

&

Center for Cancer Systems Biology
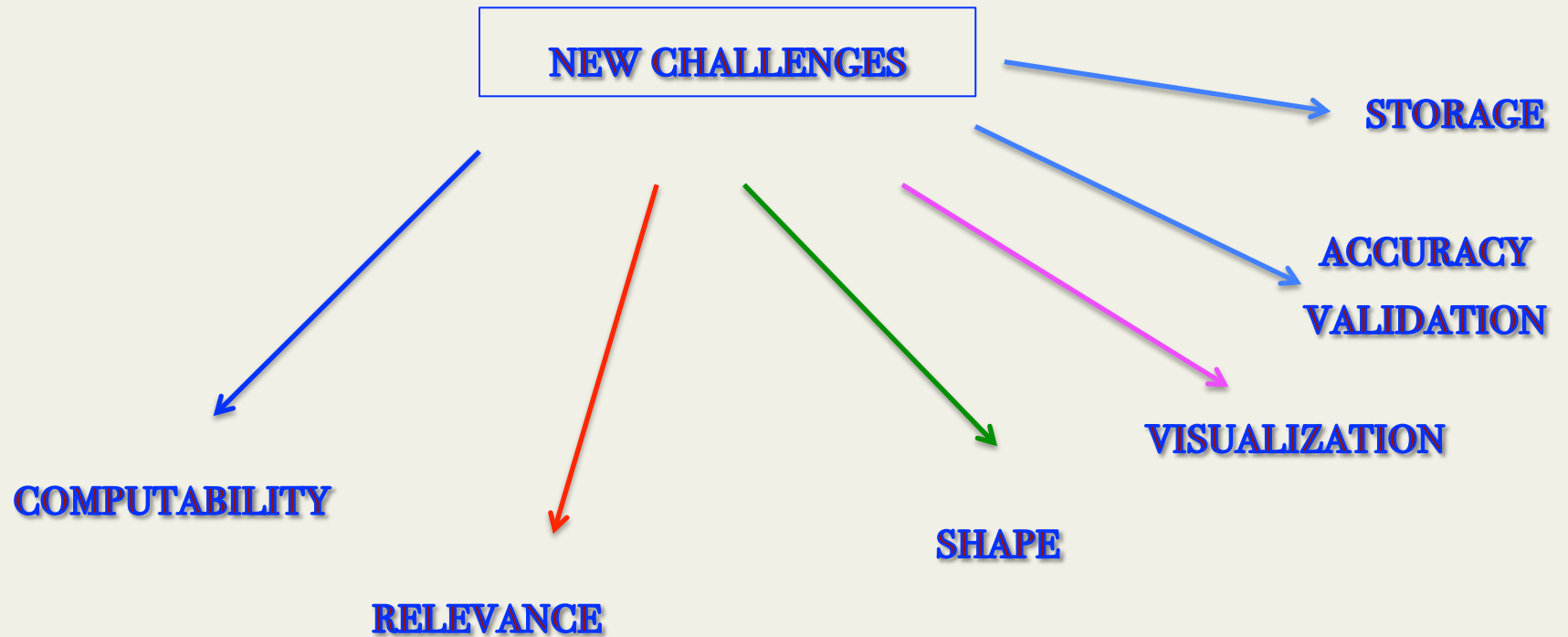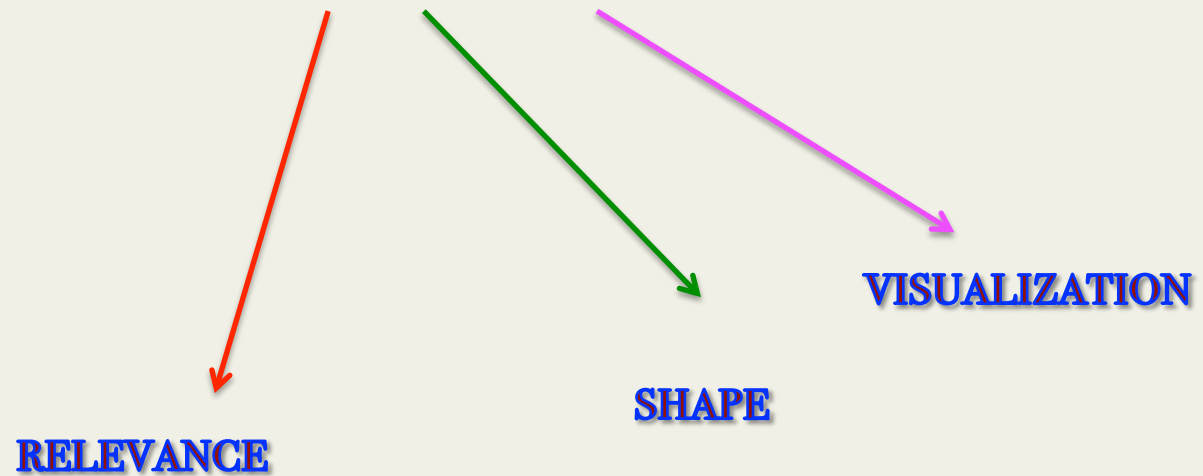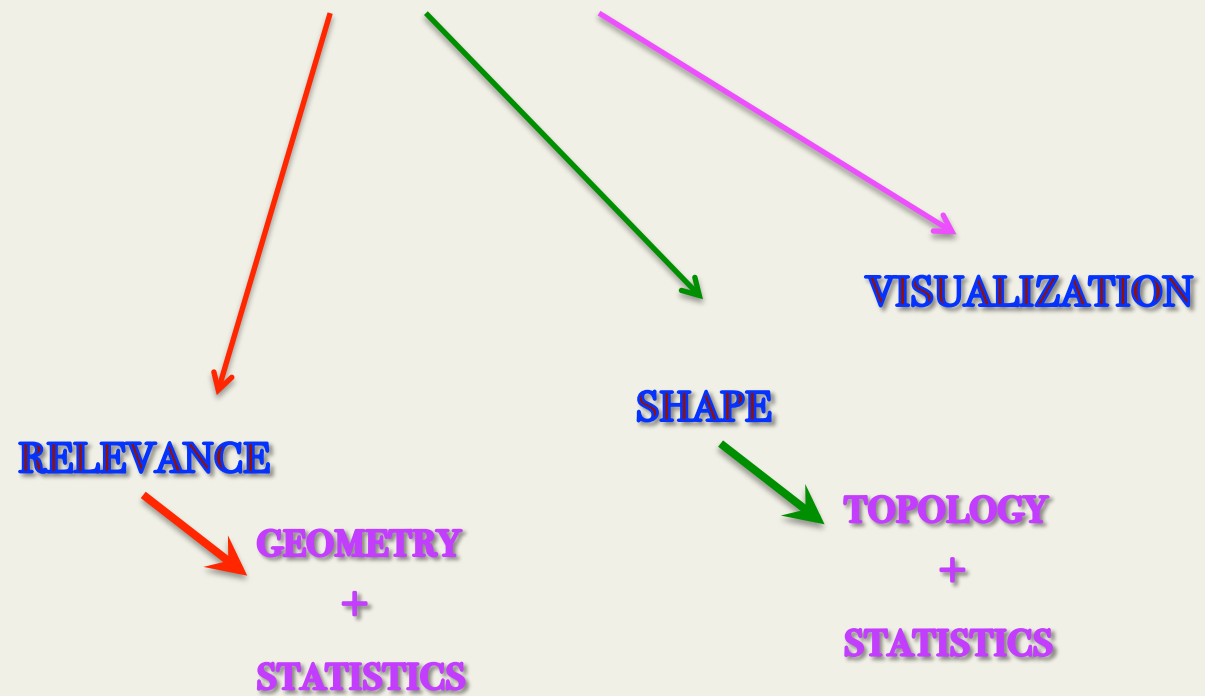
Stanford University

# Big data

# Big data

# Big data

# Big data

NEW CHALLENGES

VISUALIZATION

RELEVANCE

GEOMETRY

+

STATISTICS

SHAPE

TOPOLOGY
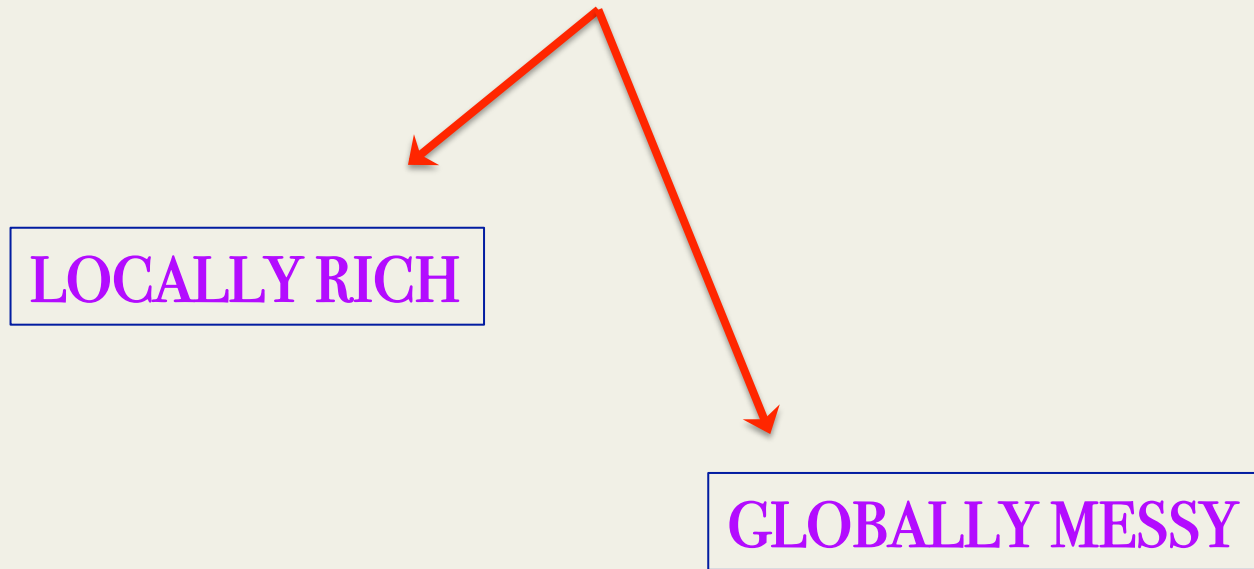
+

STATISTICS

# A little bit of history...

**The Human Genome Project:** *1990 – 2003*   DOE & NIH

international effort to discover all the estimated **20,000-25,000** human genes.

determine the complete sequence of the **3 billion** *DNA* subunits
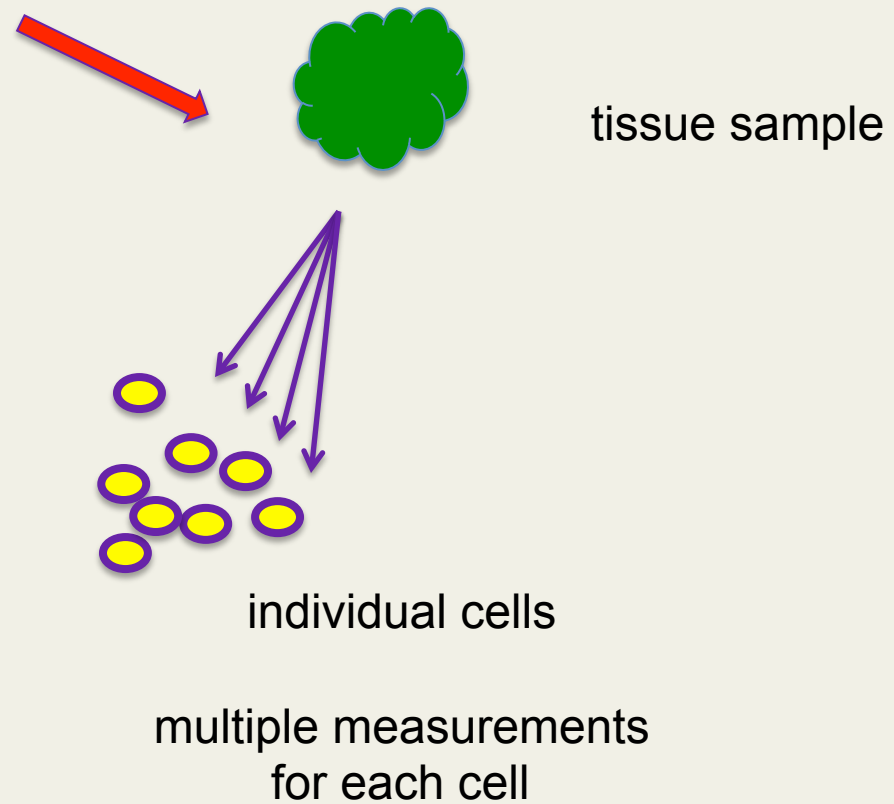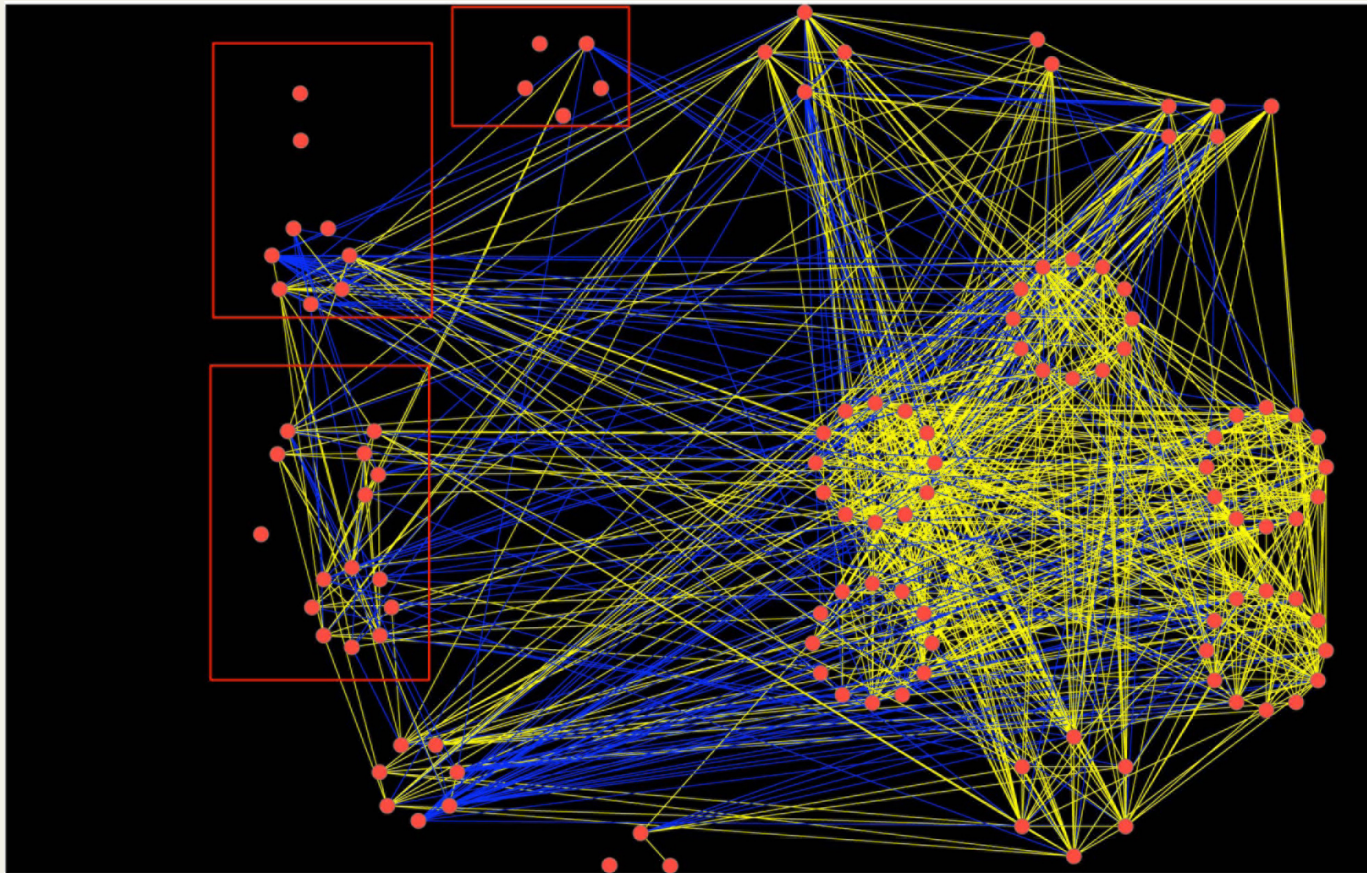
**Data is large**

# Big data

LOCALLY RICH

GLOBALLY MESSY

**example:**

Smoothing a hairball

Nolan Lab
Monack Lab
Stanford

tissue sample

individual cells

multiple measurements
for each cell
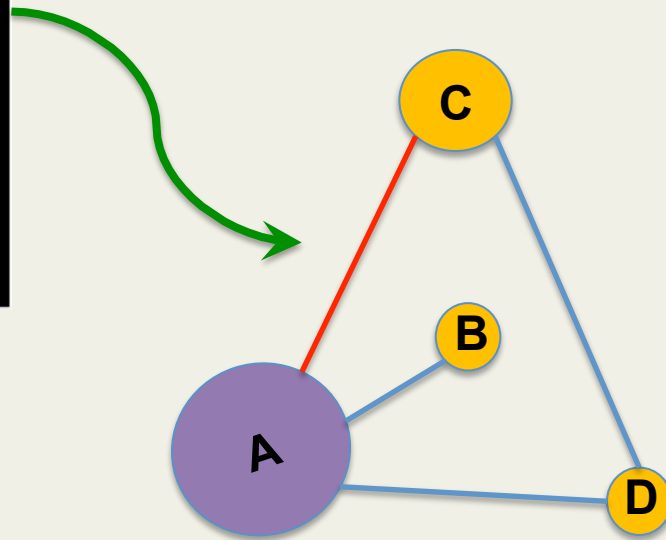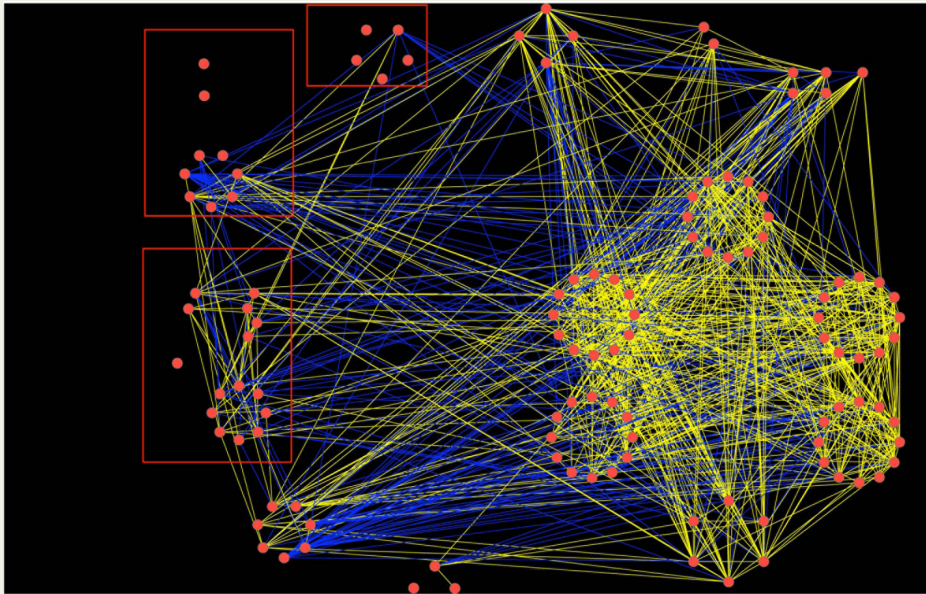
# "Hairball"
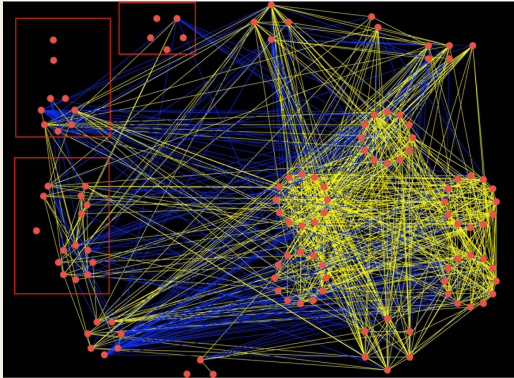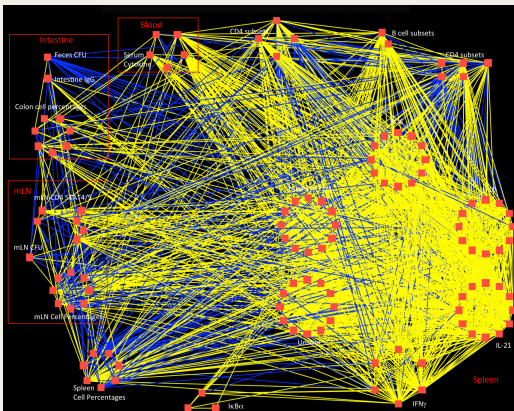


Nicolau, Hotson, Gopinath

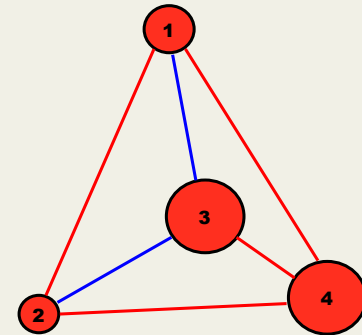# Smoothing the "Hairball"



Nicolau, Hotson, Gopinath
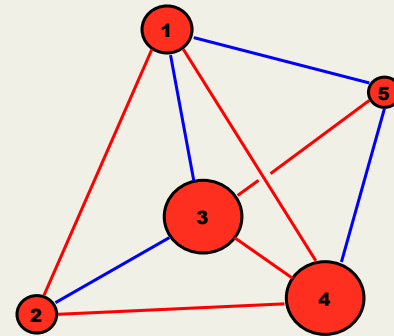
uninfected mice data

infected mice data

# A little bit of history...

**The Human Genome Project:** *1990 – 2003*   **DOE & NIH**

international effort to discover all the estimated **20,000-25,000** human genes.

determine the complete sequence of the **3 billion** *DNA* subunits

**What do these genes do?**

gene expression microarrays –  1995 (Science) & 1996 (Nature Biotechnology)

# Acute Myeloid Leukemia (AML)



AML gene
expression
heatmap

# Networks



AML network

Legend:

● node = gene

— edge = positive correlation

— edge = negative correlation

# Simplified AML Hairball

# Data smoothing: local smoothing of large data



**DATA**

**FLAT DATA**

Nicolau et al – Bioinformatics 2007

# Data smoothing: local smoothing of large data

**1. FLAT CONSTRUCTION** – DATA DE-SPARCING

$$N_1, N_2, \ldots N_k \longrightarrow \hat{N}_1, \hat{N}_2, \ldots \hat{N}_k$$

$$\hat{N}_i \quad \text{FIT TO LINEAR MODEL IN} \quad N_1, N_2, \ldots, N_{i-1}, N_{i+1}, \ldots N_k$$

# Analysis of high throughput data

RELEVANCE

geometric transformations
hypothesis
definition &  testing

SHAPE OF DATA

applied topology and
persistence - robustness

# High throughput data & relevance

- What in the data is relevant to **my study**?

# RELEVANCE

**Understand <u>disease</u> processes from data**

**what is <u>relevant</u> to the disease?**

**transform data to emphasize <u>aberrant</u> patterns compared to healthy tissue data**

**Disease specific genomic analysis (*DSGA*)**
*Nicolau M, Tibshirani R, Børresen-Dale AL, Jeffrey SS* Bioinformatics 2007

# WHAT DOES DISEASE LOOK LIKE?

DSGA – Nicolau et al – Bioinformatics 2007

# WHAT DOES DISEASE LOOK LIKE?

**cancer cells retain memory of their (healthy) cell type signature**

DSGA – Nicolau et al – Bioinformatics 2007

# RELEVANCE –
## *Disease specific genomic analysis - DSGA*

**TUMORS**

**HEALTHY TISSUE**

**microarrays**

DSGA – Nicolau et al – Bioinformatics 2007

# RELEVANCE –
## *Disease specific genomic analysis - DSGA*

**points (HEALTHY)**

**MATHEMATICAL MODEL
FOR HEALTHY STATE**

DSGA – Nicolau et al – Bioinformatics 2007

# RELEVANCE –
## *Disease specific genomic analysis - DSGA*

**points (TUMORS)**



TRANSFORMATION

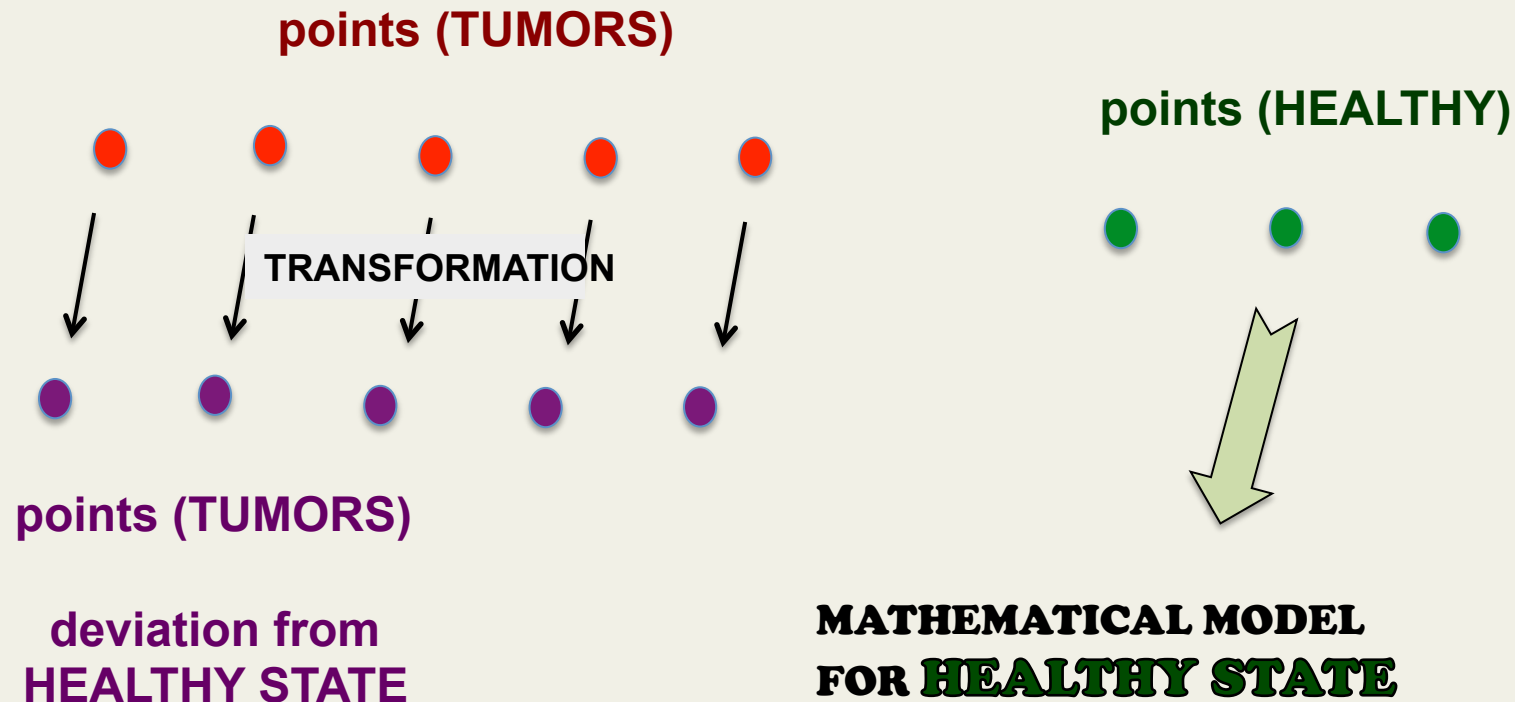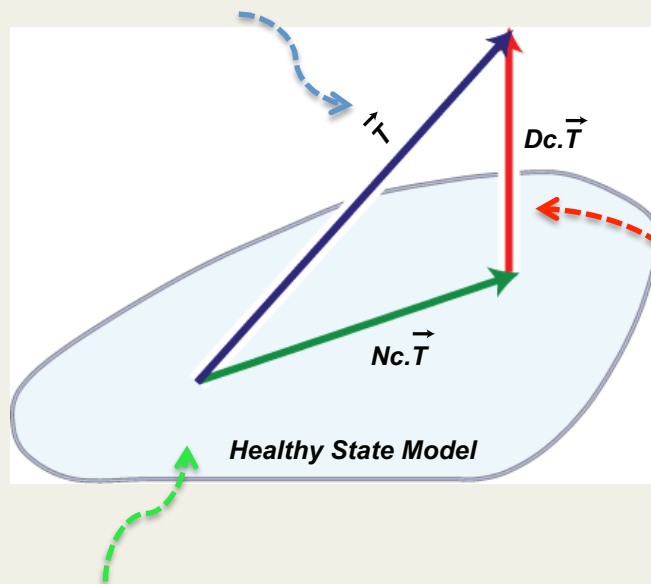**points (HEALTHY)**

**points (TUMORS)**

**deviation from HEALTHY STATE**

**MATHEMATICAL MODEL FOR HEALTHY STATE**

DSGA – Nicolau et al – Bioinformatics 2007

# RELEVANCE –
## *Disease specific genomic analysis - DSGA*



**Tumor data**

$Dc.\vec{T}$

$Nc.\vec{T}$

Healthy State Model

transformed
tumor data
vector of residuals

[Null Hypothesis Space]

Normal tissue data

DSGA – Nicolau et al – Bioinformatics 2007

# RELEVANCE –
## *Disease specific genomic analysis - DSGA*

**Benefits from *Disease component***
**of tumor data:**

1. Highlight **extent of deviation** from normal – aberrant
   behavior

2. **Cleaner** identification of distinct classes

3. **Biology** you highlight is **different** from using original
   data.

DSGA – Nicolau et al – Bioinformatics 2007

# Analysis of **high throughput** data

RELEVANCE
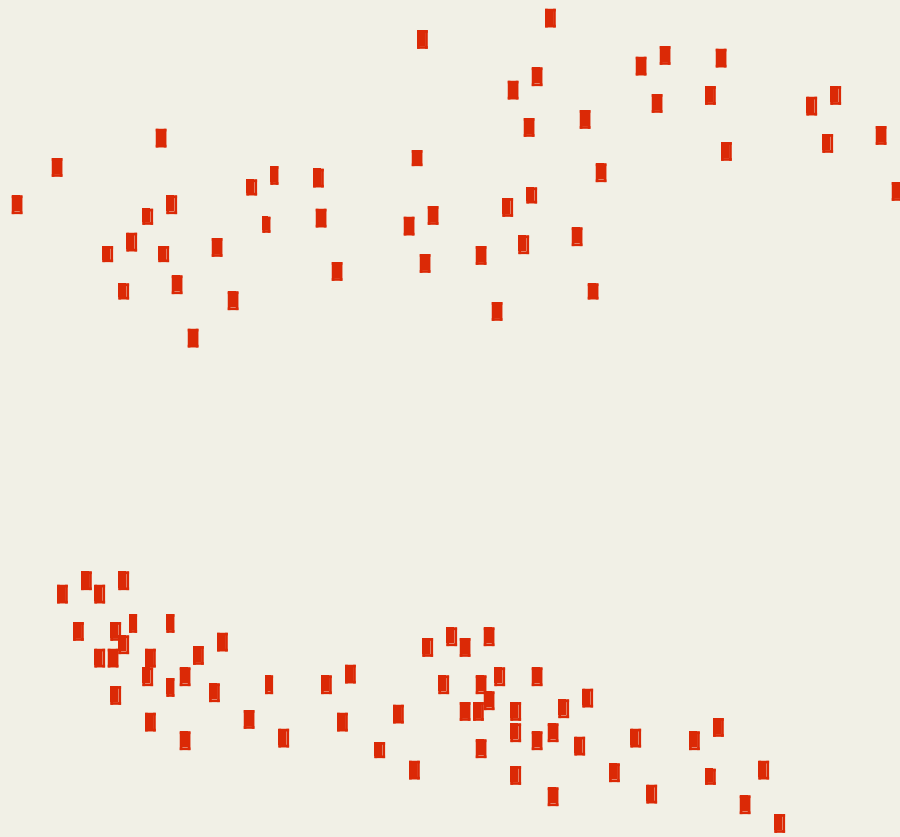
**geometric** transformations
hypothesis
definition & testing

SHAPE OF DATA

applied **topology** and
persistence - robustness
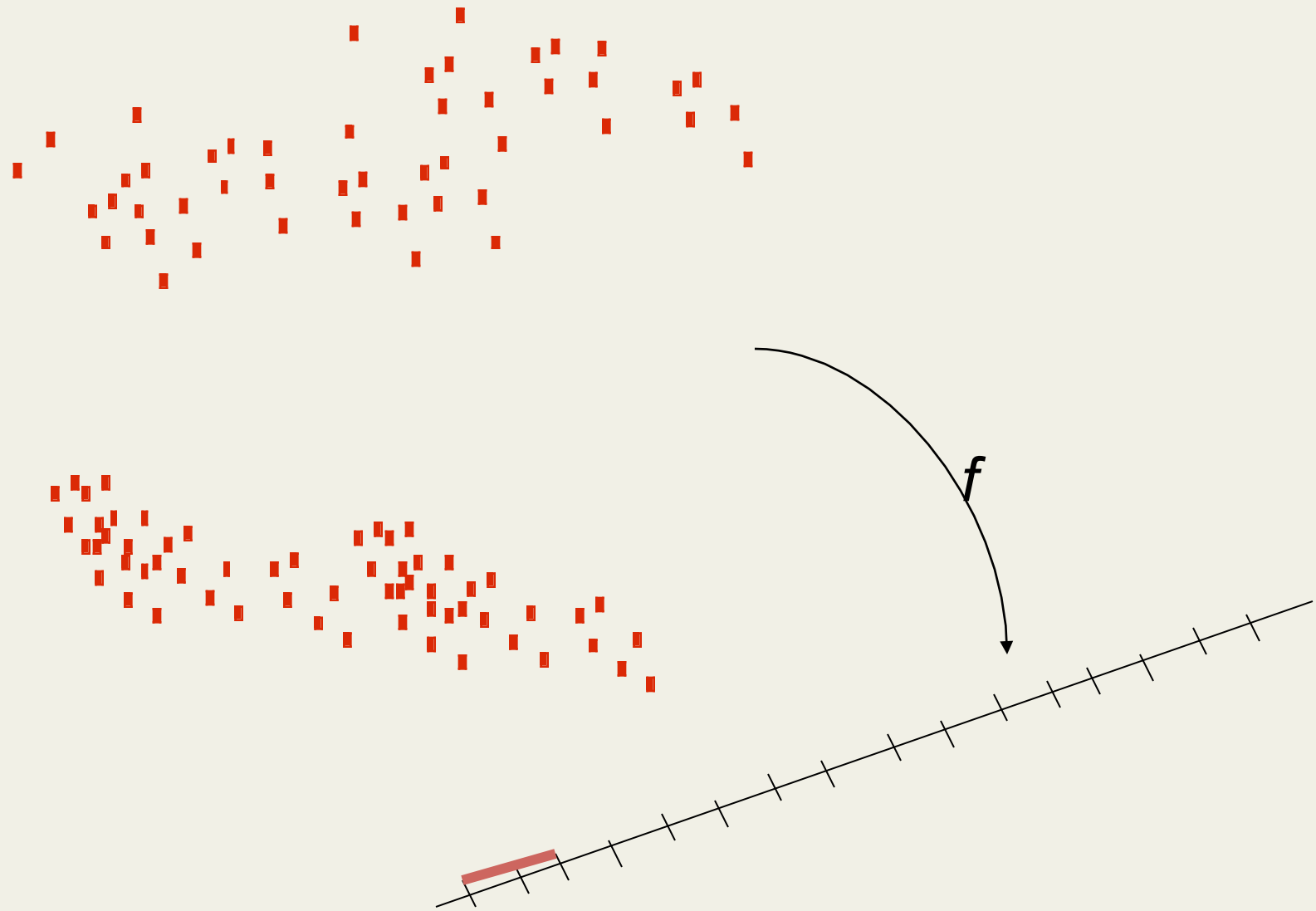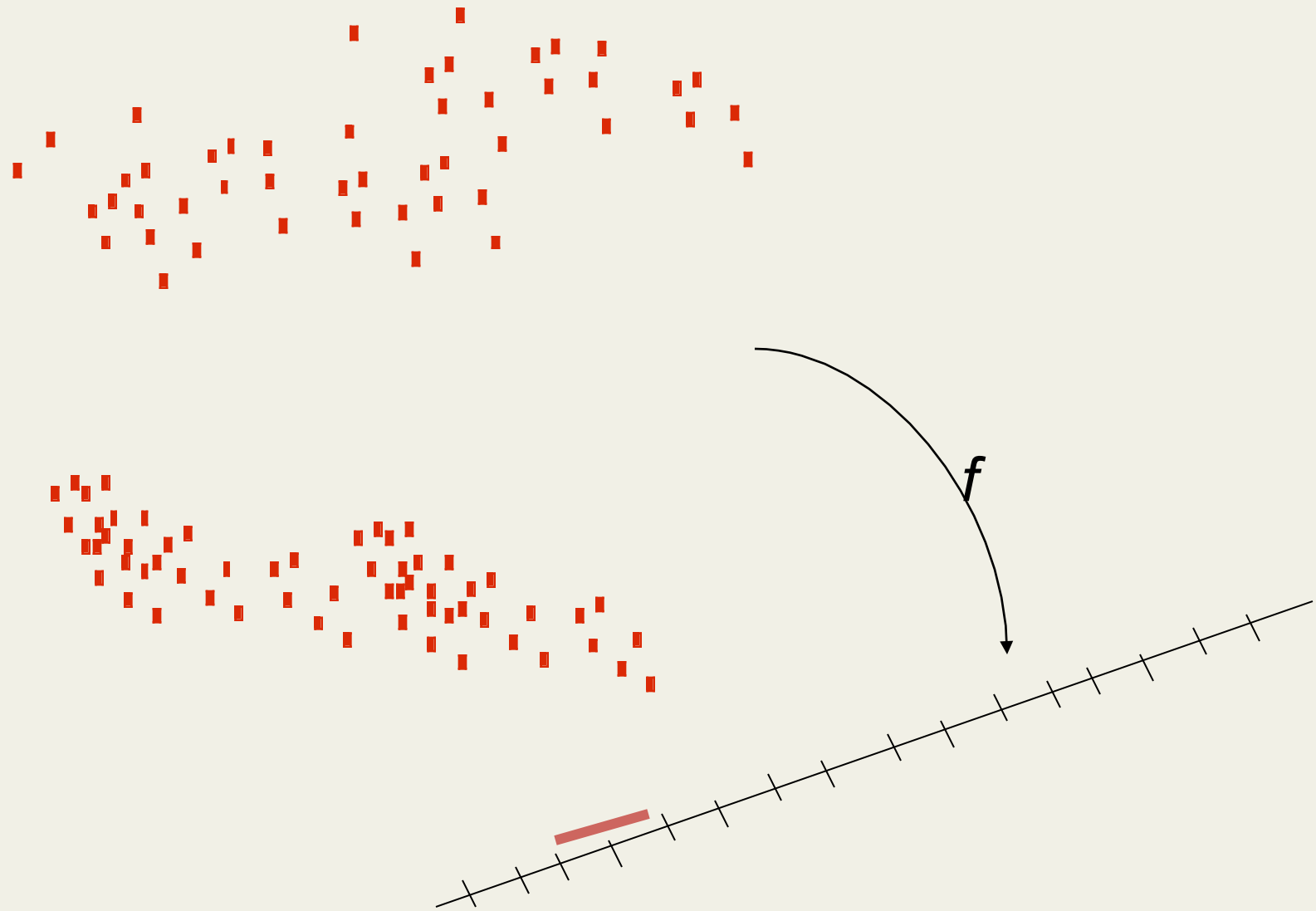
# SHAPE OF DATA
## *TOPOLOGY  &  Mapper*



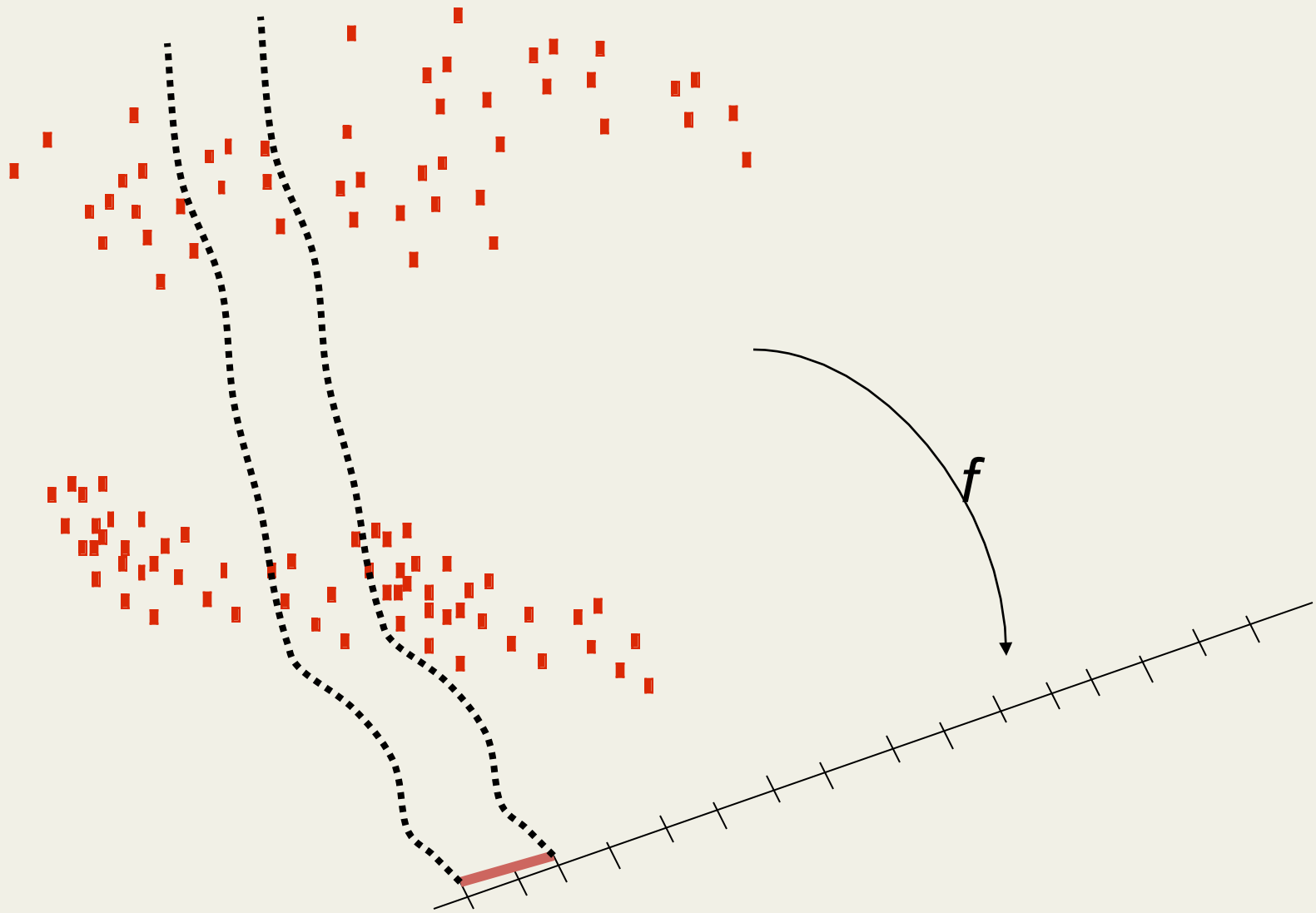Singh, Memoli, Carlsson **Point Based Graphics** 2007

# SHAPE OF DATA - mapper



*f*

*Mapper*

*f*

# *Mapper*

*Mapper*

clusters

*f*
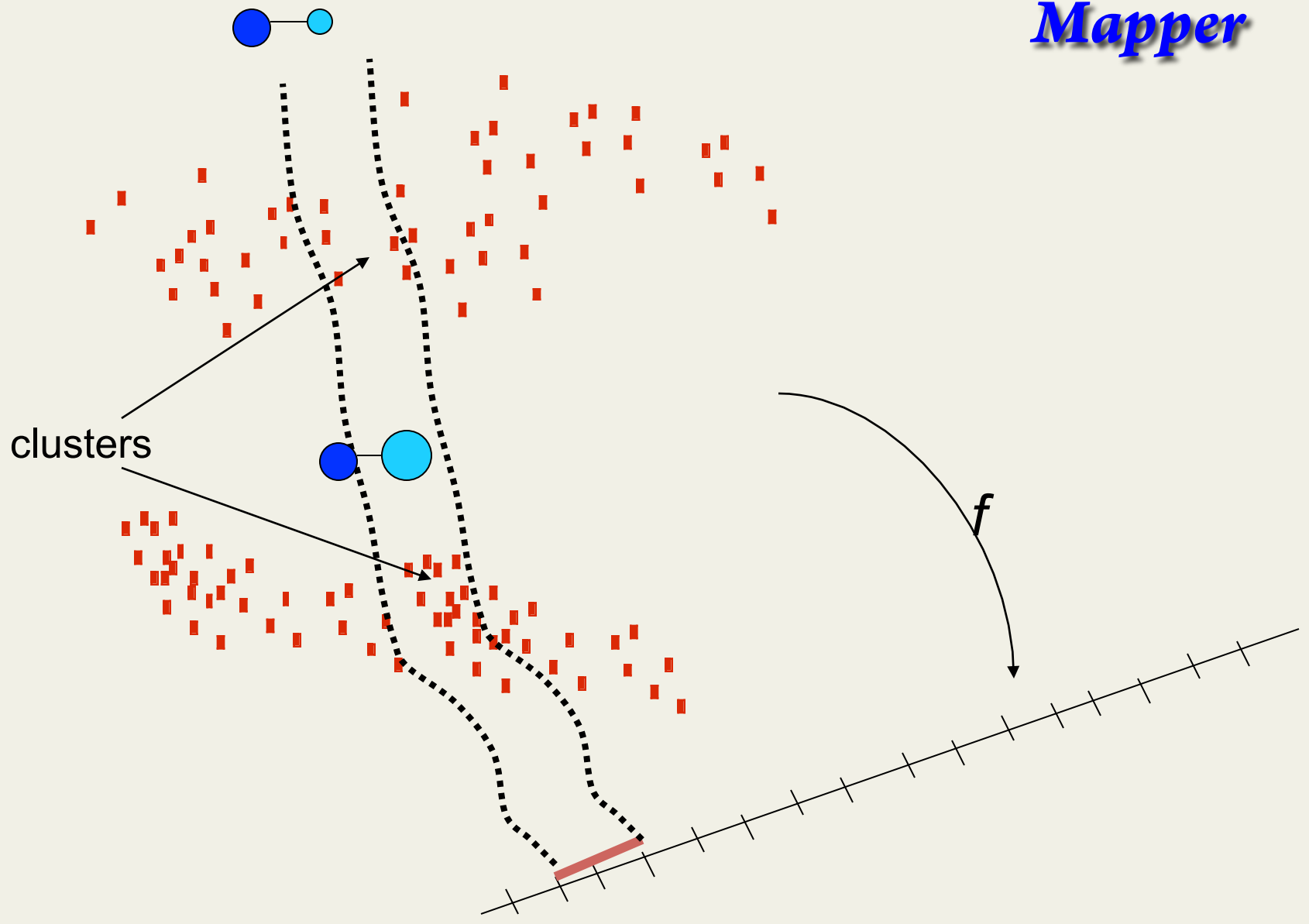
*Mapper*

clusters

*f*

*Mapper*

clusters

*f*

*Mapper*

nicolau-AMS

*Mapper*

*Mapper*

*Mapper*

*Mapper*

*Mapper*

*f*

# Mapper filter function: overall deviation from Healthy State Model



Tumor data

$Dc.\vec{T}$

$Nc.\vec{T}$

**Healthy State Model**

[Null Hypothesis Space]

Normal tissue data

**transformed tumor data vector of residuals**
*DcTumor*
**vector magnitude of**
*Disease Component*

# Progression Analysis of Disease – PAD

RELEVANCE

**geometric** transformations
**DSGA**

SHAPE OF DATA

applied **topology**
**Mapper**

**Topology based data analysis identifies subgroup of breast cancers with unique mutational profile and excellent survival**

*Nicolau M, Levine AJ, Carlsson G*   Proc Nat Acad Sci 2011

# Progression Analysis of Disease: PAD
## running Mapper on DSGA-transformed data

DSGA – transformed data from tumors & normals:
disease component

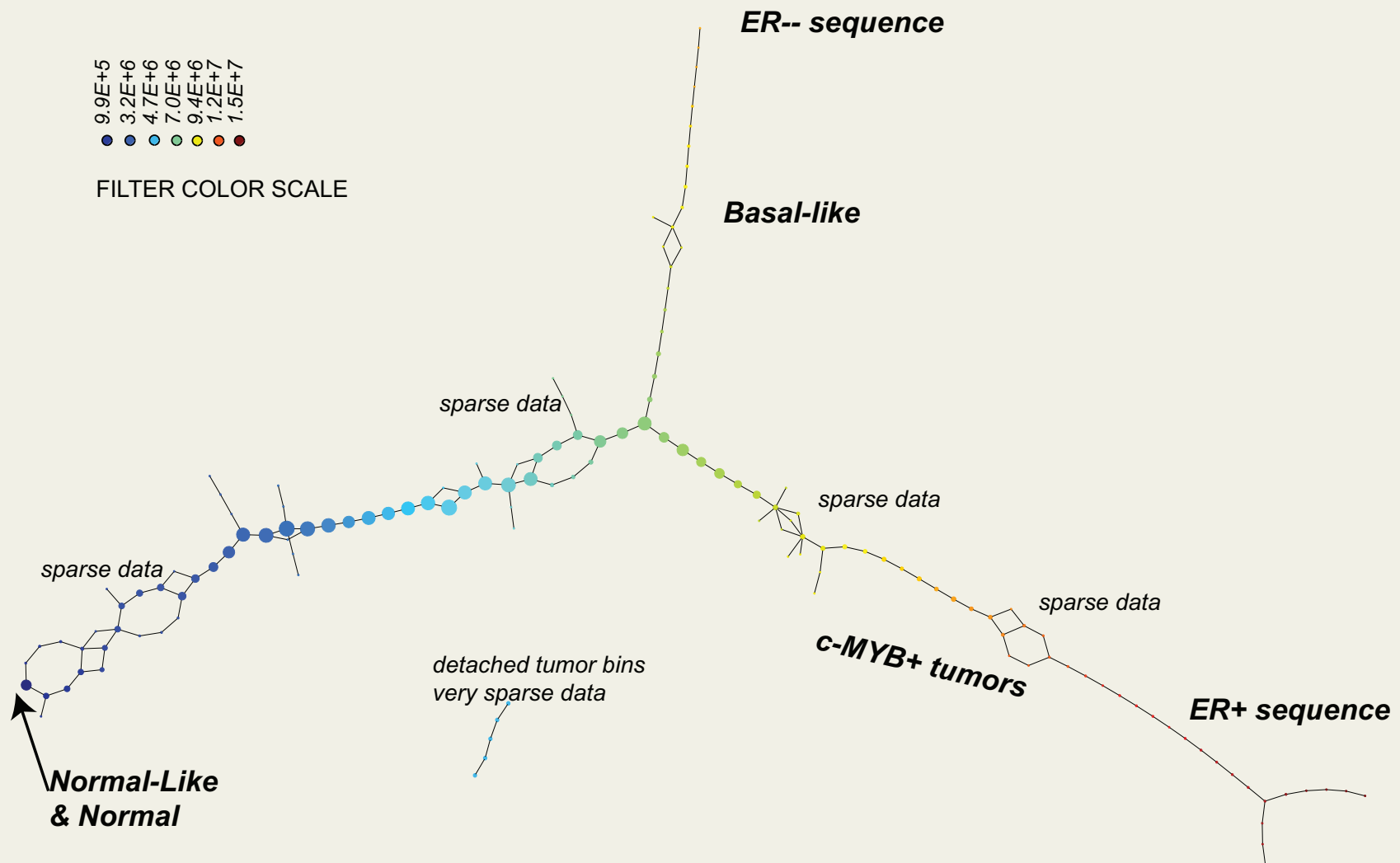Mapper filter: ($L^p$-vector magnitude)$^q$

$R$

tumors

normal

Health State Model

0

# Progression Analysis of Disease: PAD
## running Mapper on DSGA-transformed data

9.9E+5
3.2E+6
4.7E+6
7.0E+6
9.4E+6
1.2E+7
1.5E+7

FILTER COLOR SCALE

*ER-- sequence*

*Basal-like*

*sparse data*

*sparse data*

*sparse data*

*detached tumor bins
very sparse data*

*c-MYB+ tumors*

*sparse data*

*ER+ sequence*

**Normal-Like
& Normal**

*Nicolau et al* PNAS 2011

# *c-MYB+* group

**survival analysis**

**K-M survival**
cMYB+group

Group is
***homogeneous* &**
***distinct***
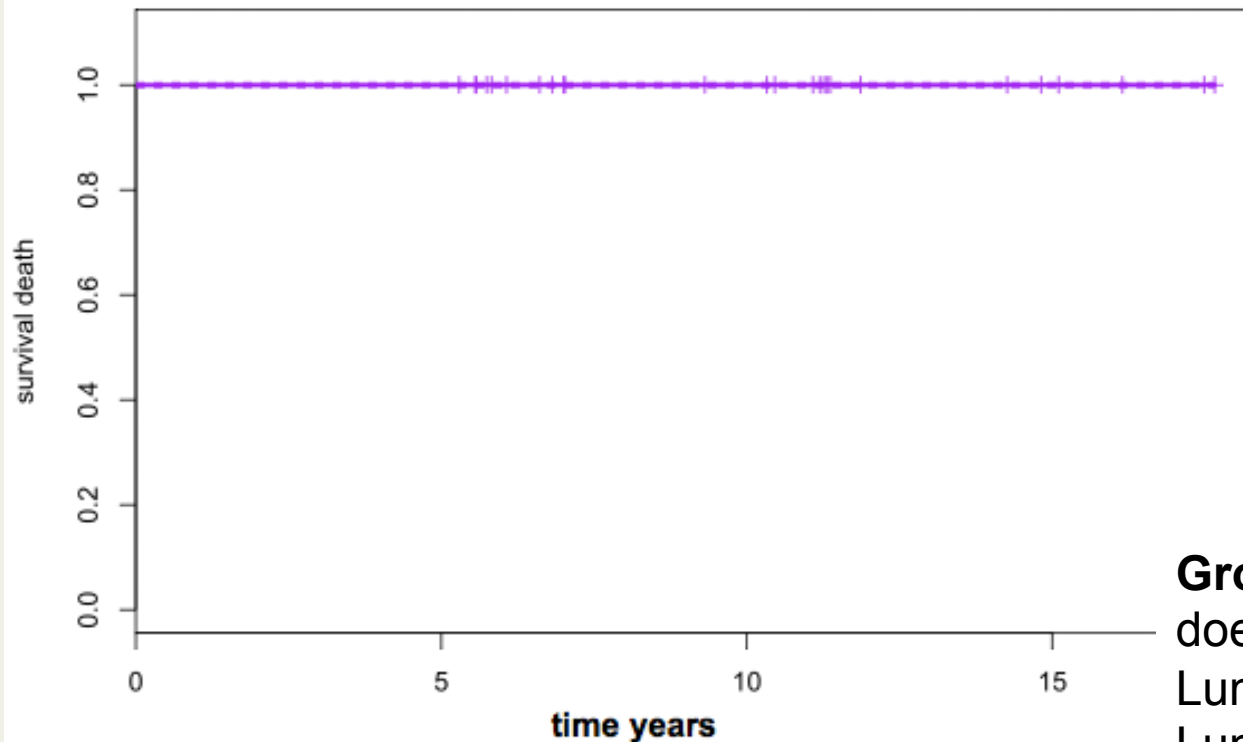   mathematically
   biologically

**predictor variables**
   few (1 or 2)

**significant variables**
   biologically meaningful

**Group is new**:
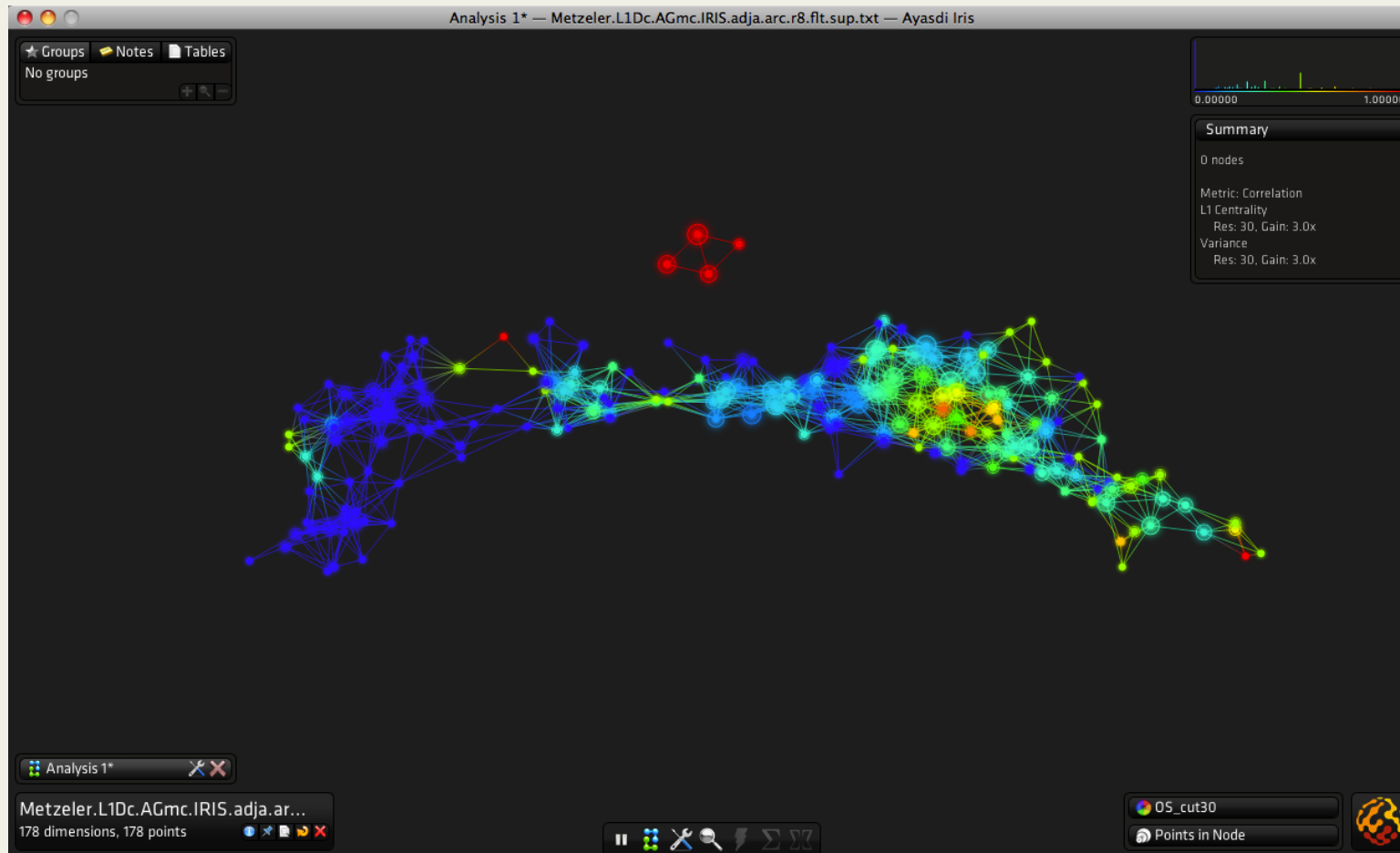doesn't follow old classification
Luminal A
Luminal B
unclassified

*Nicolau et al* PNAS 2011
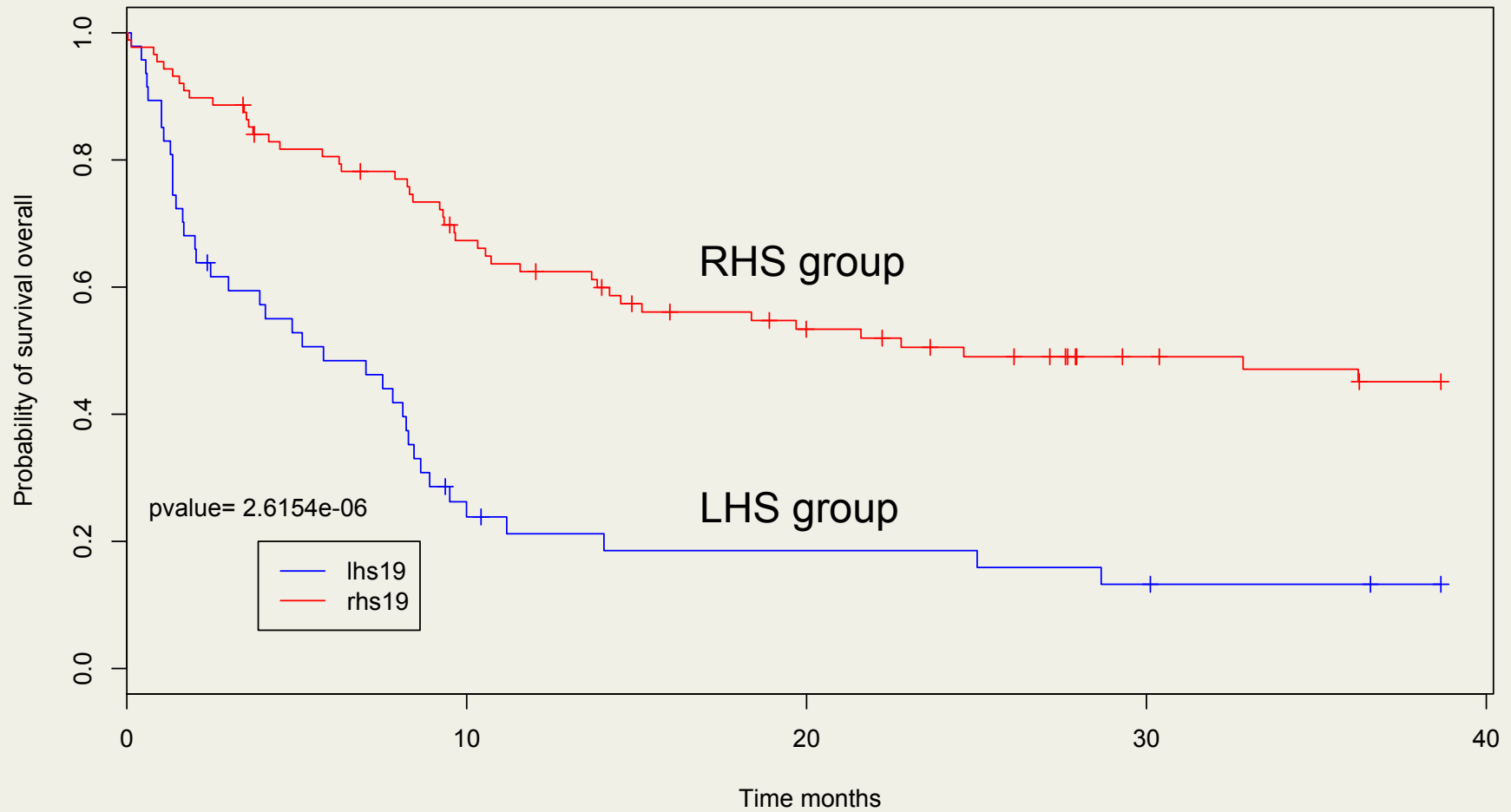
# Acute Myeloid Leukemia

**another example**
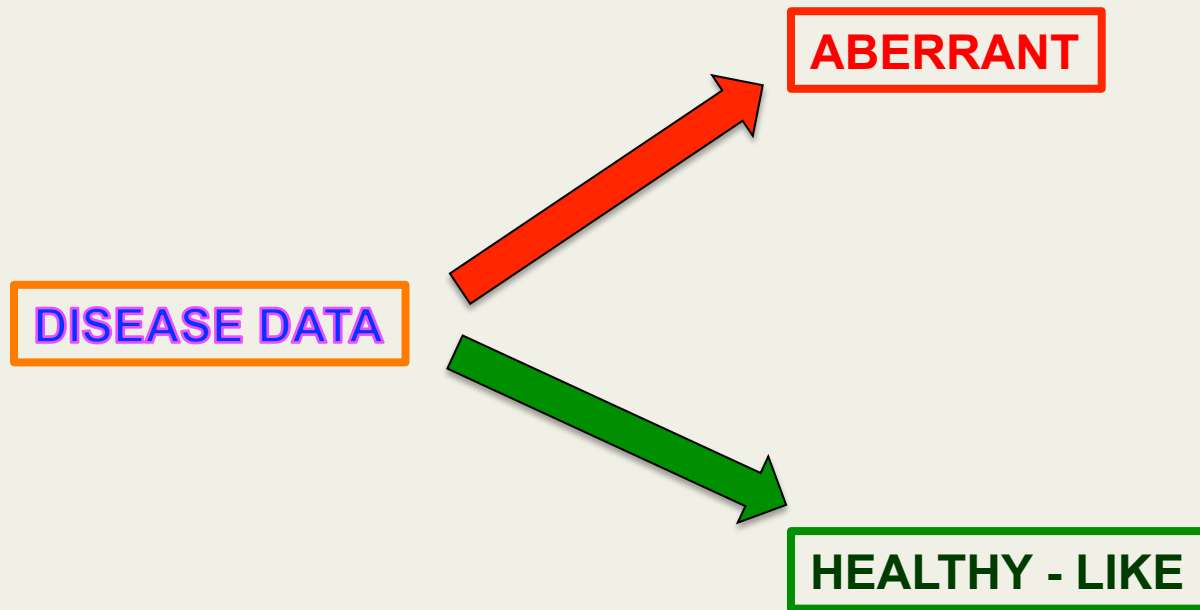
# Disease component IRIS: AML



node color: survival

# survival – RHS vs LHS



**KM survival Metzeler DSGA_Dc**
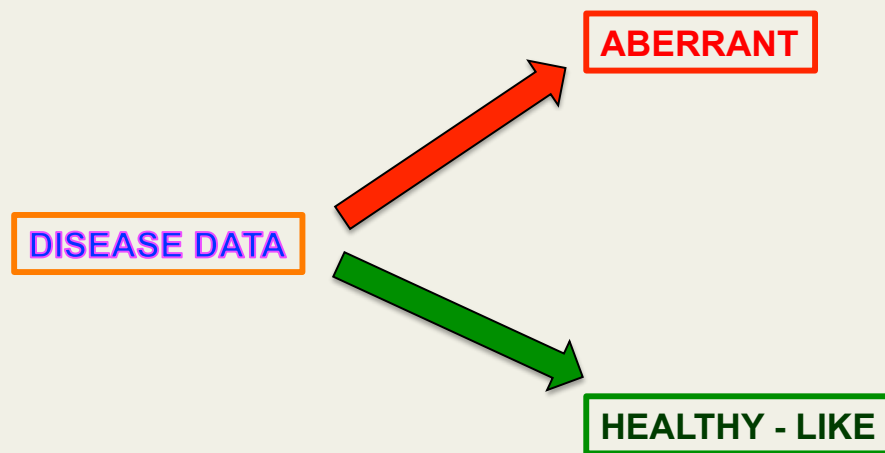**kernel.embed_IRIS_L1cent&var  RHS vs LHS**

# Another look at AML data



ABERRANT

DISEASE DATA

HEALTHY - LIKE

**Disease Specific Genomic Analysis: DSGA**
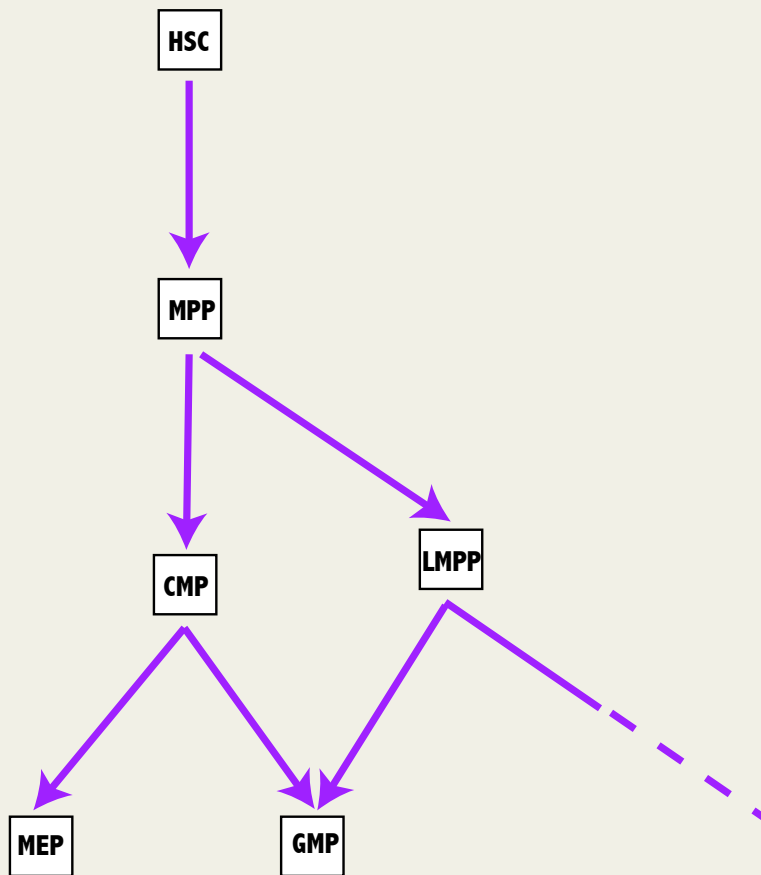
# Another look at AML data
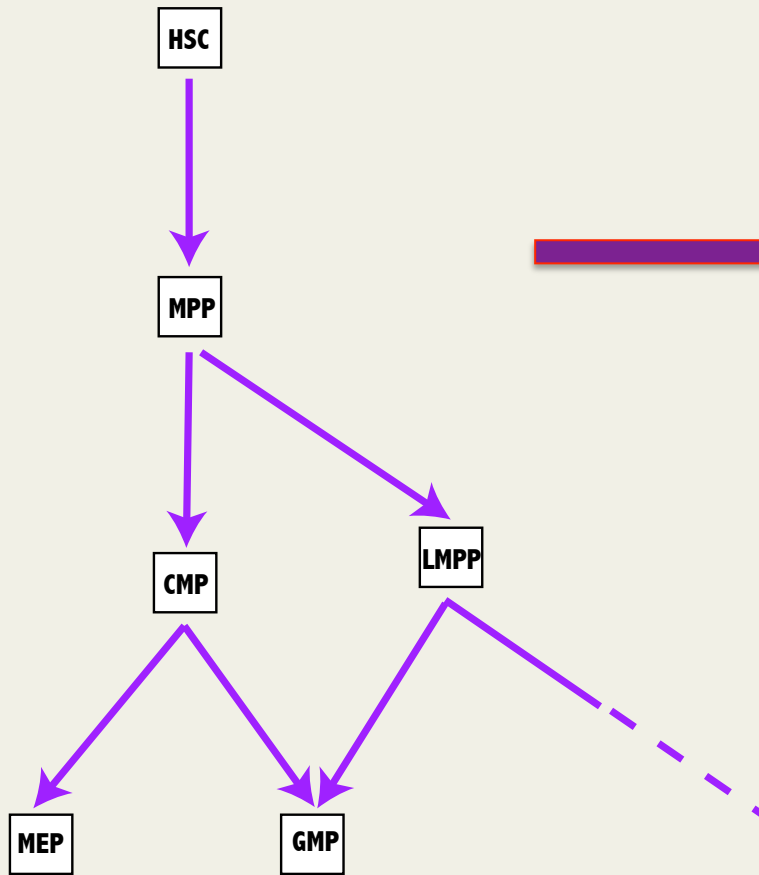


**NORMAL COMPONENT OF AML CELLS**

**do AML tumor cells retain a memory of healthy signatures?**

**do differences in this memory have significance for disease?**

# Hematopoiesis

# Hematopoiesis

HSC → MPP → CMP → MEP, GMP
MPP → LMPP → GMP

# AML developmental stages

LSC → LPC → Blast
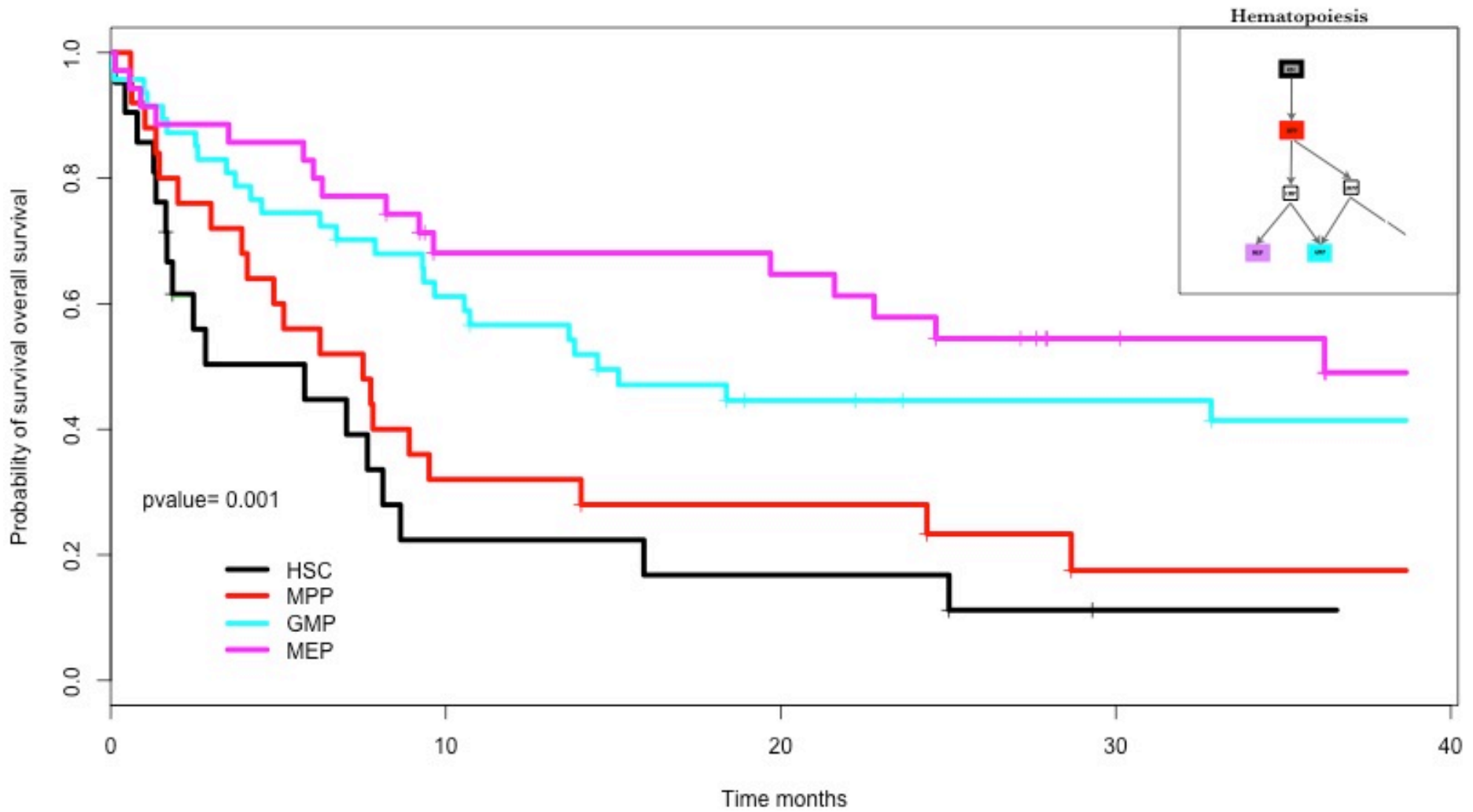
● CD34+/CD38-    *LSC*
● CD34+/CD38+    *LPC*
● CD34-          *Blast*

*Majeti* markers

KM survival Metz global Nc.scores - HSC v MPP v GMP v MEP

# conclusion

***Disease component*** highlights aberrant behavior

association with clinical characteristics
cleaner groups of genes associated with distinct biology
together with Mapper found *novel group of **breast cancer***
             & found strong *association with survival in **AML***
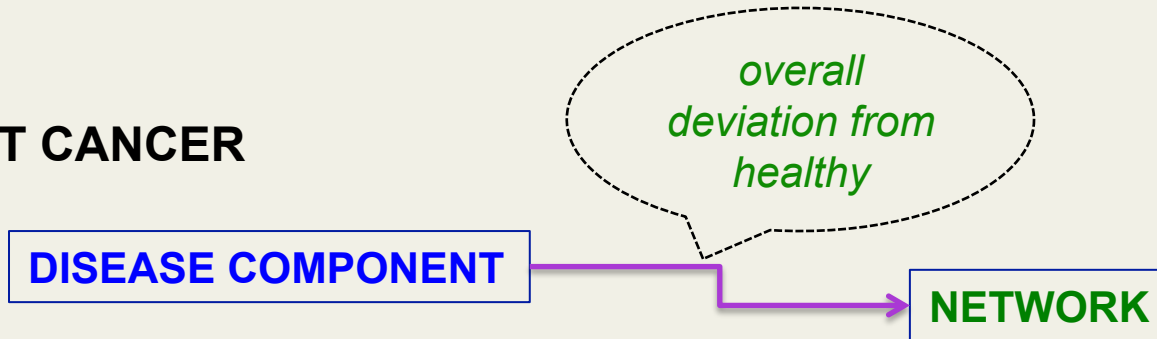
***Normal component*** extracts memory of healthy types

        association with clinical characteristics
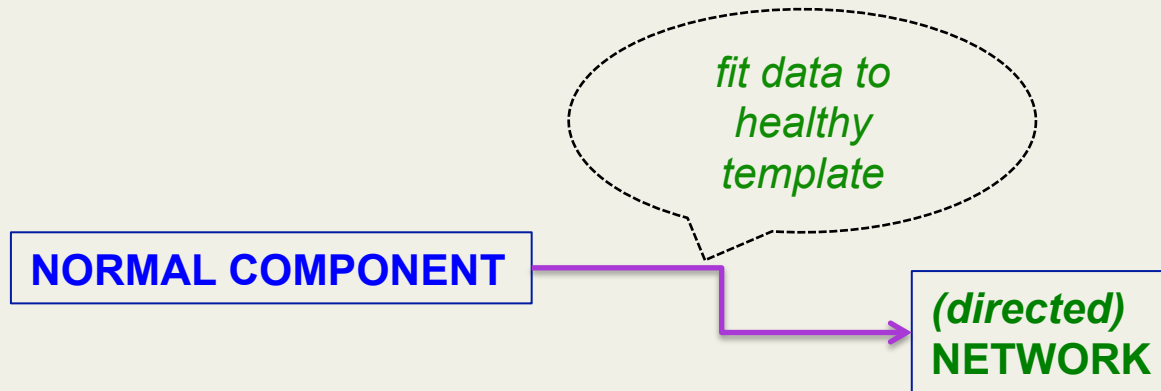        identified *novel groups of **AML***
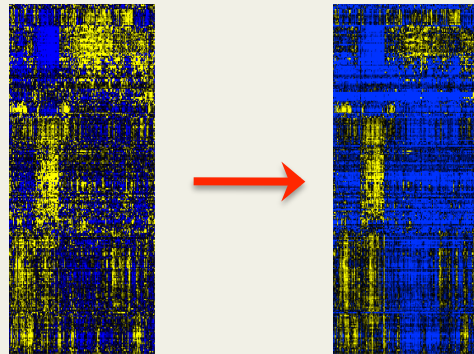
# Networks/hairballs everywhere
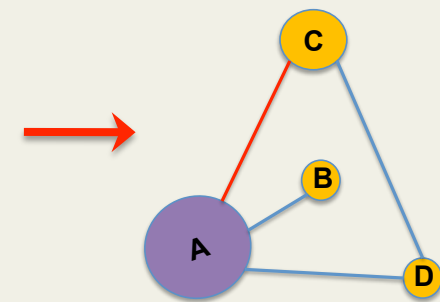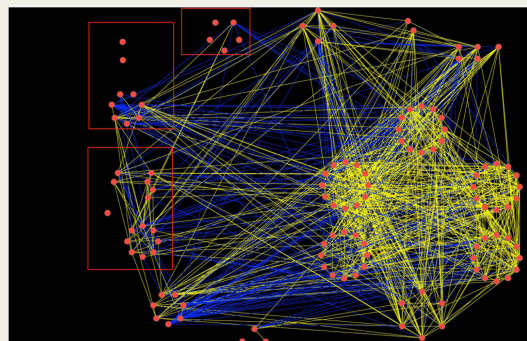
**Large data** → LOCALLY RICH

GLOBALLY MESSY

**locally smooth data**

**locally smooth data similarities (hairball)**

# Thanks -

## Computational:

**Gunnar Carlsson** (Stanford Mathematics)
**Sylvia Plevritis** (Stanford Cancer Center for Systems Biology)
**Rob Tibshirani** (Stanford Statistics)

## Biology:

**Arnold Levine** (Princeton IAS – School of Natural Studies)
**Anne-Lise Børresen-Dale** (Genetics, University of Oslo, Norway)
**Stefanie Jeffrey** (Stanford Surgery)
**Amato Giaccia** (Stanford Radiation Oncology)
**Janine Erler** (Cell and Molecular Biology, Institute of Cancer Research, London, UK)
**Ravindra Majeti** (Stanford Hematology)
**Garry Nolan** (Stanford Immunology)

# *Thanks - funding*

**NIH – National Human Genome Research Institute (NHGRI)**

**California Breast Cancer  Research Program**

**DARPA**

**Air Force Office of Scientific Research**

**National Institutes of Health**