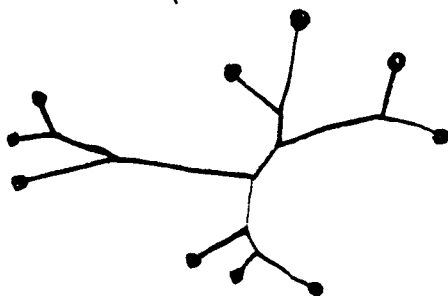


# Clusteringverfahren

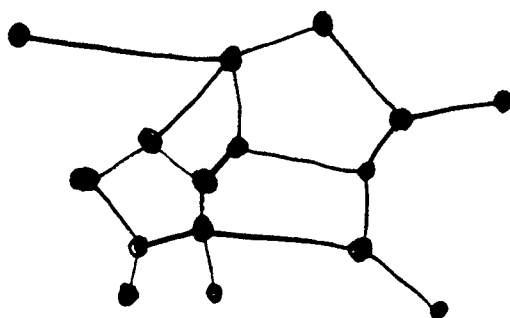
- o Cluster analysis of gene expression dynamics  
 Marco Ramoni, Paola Sebastiani, Isaac Kohane  
 PNAS 14, 2002.

- o Cluster inference methods and graphical models  
 evaluated on NC160 microarray gene expression data  
 Peter Waddell, Hirohisa Kishino  
 Genome Informatics 11, 2000.

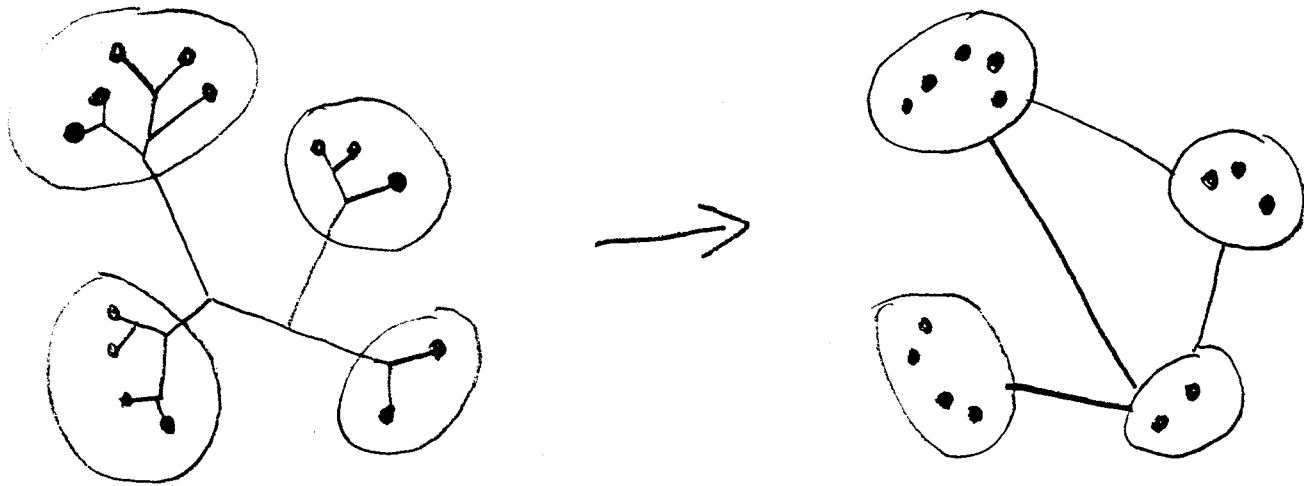
- o Clusterverfahren erzeugen Bäume;



keine Netze im eigentlichen Sinne.



- Clusterverfahren erzeugen Gruppen ähnlicher Gene, die zur weiteren Analyse als Einheit gehandhabt werden können. (5)



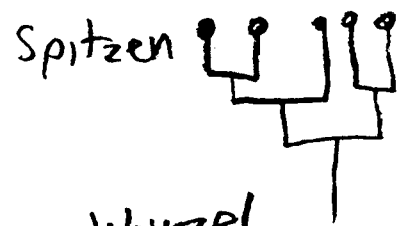
Da durch vereinfacht sich die weitere Analyse erheblich.

- Clusterverfahren benötigen ein Distanzmaß<sup>4</sup> zwischen den Genen, eine sog. Metrik. Diese kann auf vielfältige Weise definiert werden. (4)

- Clustering erfordert also zwei Schritte:

1. Berechnung aller paarweisen Distanzen
2. Erzeugung des Baums

- Die Menge der Distanzen kann bestimmte Eigenschaften haben, z.B. ultrametrisch sein:



- Oder auch nicht.

# Partielle Korrelation

(4.5)

$$r_{xy.z} = (r_{xy} - r_{xz}r_{yz}) / \left( (1 - r_{xz}^2)(1 - r_{yz}^2) \right)^{1/2}$$

Für mehr als 3 Variablen:

$$r_{xy.g} = -w_{xy} / (w_{yx} w_{xy})^{1/2}$$

$$W = V^{-1} \quad (\text{inverse Kovarianzmatrix})$$

# Distanzen

(4.6)

$$1. \quad \delta_r = 1 - r \quad \left. \vphantom{\delta_r} \right\} 0 \dots 2$$

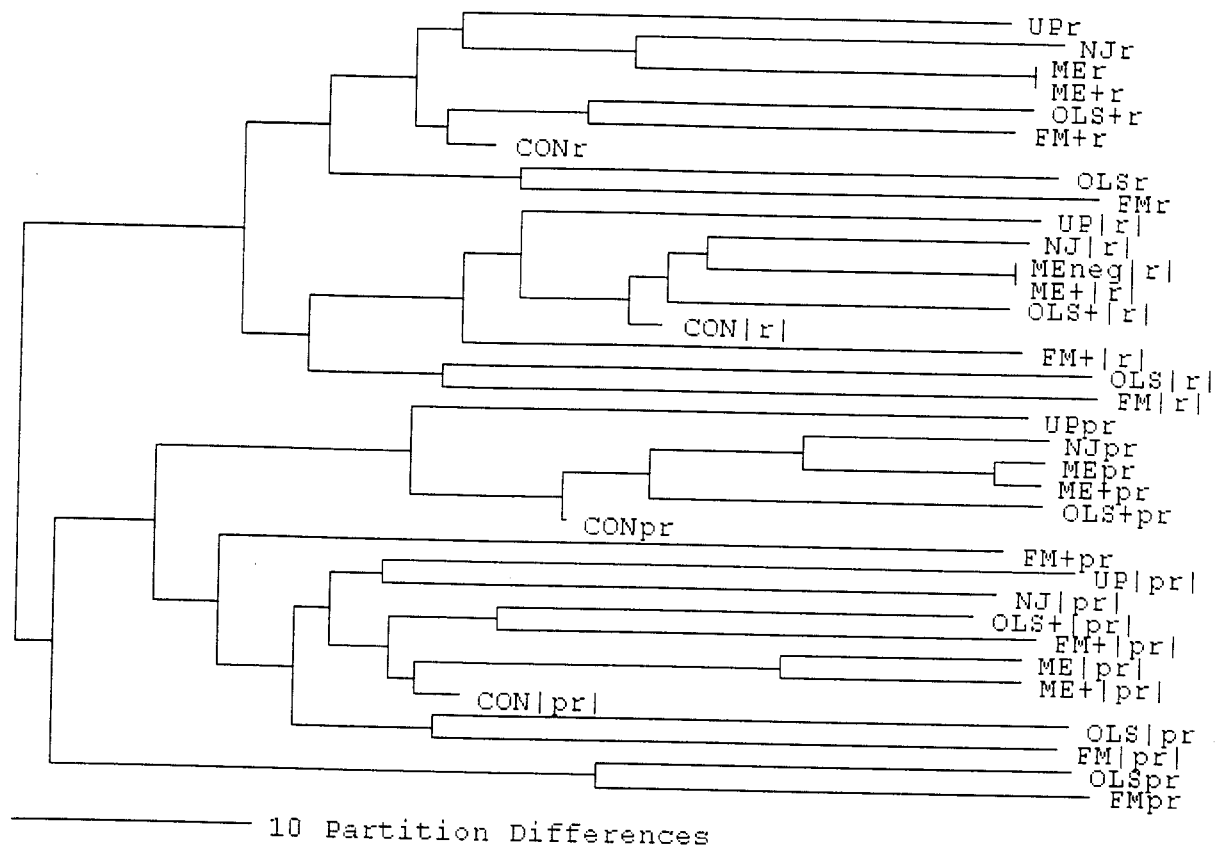
$$2. \quad \delta_{pr} = 1 - pr$$

$$3. \quad \delta_{|r|} = 1 - |r| \quad \left. \vphantom{\delta_{|r|}} \right\} 0 \dots 1$$

$$4. \quad \delta_{|pr|} = 1 - |pr|$$



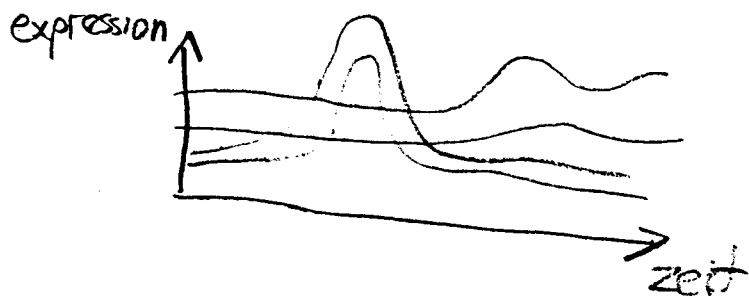
o Das Verhalten der Methoden kann in einem Baum dargestellt werden.



Ergebnis:

1. Gleiche Distanzmasse ergeben "ähnlichere" Bäume als verschiedene bei gleichem Clusterverfahren
2. ME - Methoden sind untereinander "ähnlicher" als gegenüber NJ.
3. Das Spektrum der erhaltenen Bäume ist so gross, dass wohl für jede Anwendung ein passender dabei ist.

- Ramoni et. al. clustern gene aufgrund ihrer expressionsdynamik. (7)



- Ein Bayesianischer Ansatz liefert ein Distanzmass welches die Ähnlichkeit der zeitlichen Expressionsverläufe beschreibt
- Aufgrund dieser Distanzen erzeugt ein Programm namens CAGED Baumstrukturen

- Expressionsverlauf wird beschrieben durch einen autoregressiven Ansatz. (8)

$$\vec{x}_t = \underline{A}_1 \vec{x}_{t-1} + \underline{A}_2 \vec{x}_{t-2} + \dots + \underline{A}_p \vec{x}_{t-p} + \vec{c} + \vec{\varepsilon}$$

Expression zum Zeitpunkt  $t$  hängt von den vorhergehenden Expressionsschritten ab. Oft verwendet man nur den letzten.  $p$  ist die Ordnung.

- Der zufällige Fehler  $\vec{\varepsilon}$  erlaubt noch Annahme einer geeigneten Verteilung die Berechnung der Wahrscheinlichkeit, das zwei Expressionsverläufe in Wirklichkeit identisch sind.

o Die Likelihood Funktion:

(8.5)

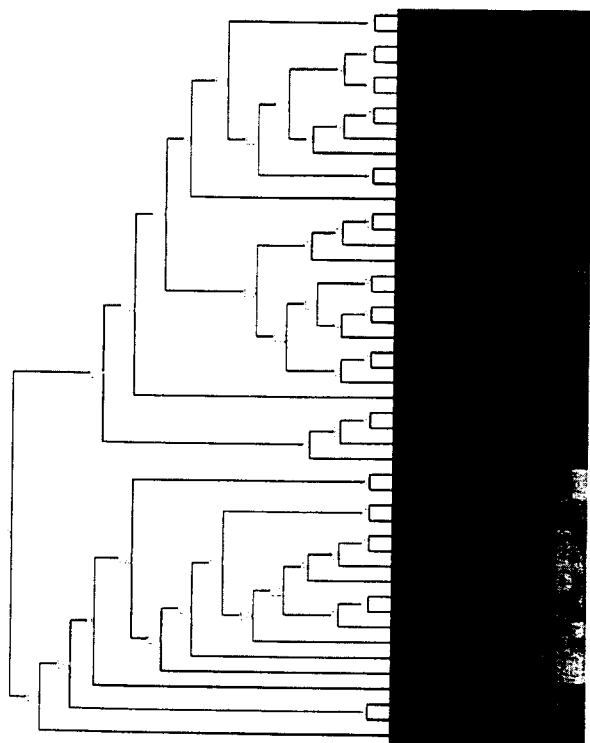
$$P(M_c|y) \propto P(M_c) f(y|M_c)$$

$$f(y|M_c) = \frac{\Gamma(1)}{\Gamma(1+m)} \times \prod_{k=1}^c \frac{\Gamma(m_k/m + m_k) \left(\frac{RSS_k}{2}\right)^{(q-n_k)/2} \Gamma\left(\frac{m_k-1}{2}\right)}{\Gamma(m_k/m) (2\pi)^{(n_k-1)/2} \det(X_k^T X_k)}$$

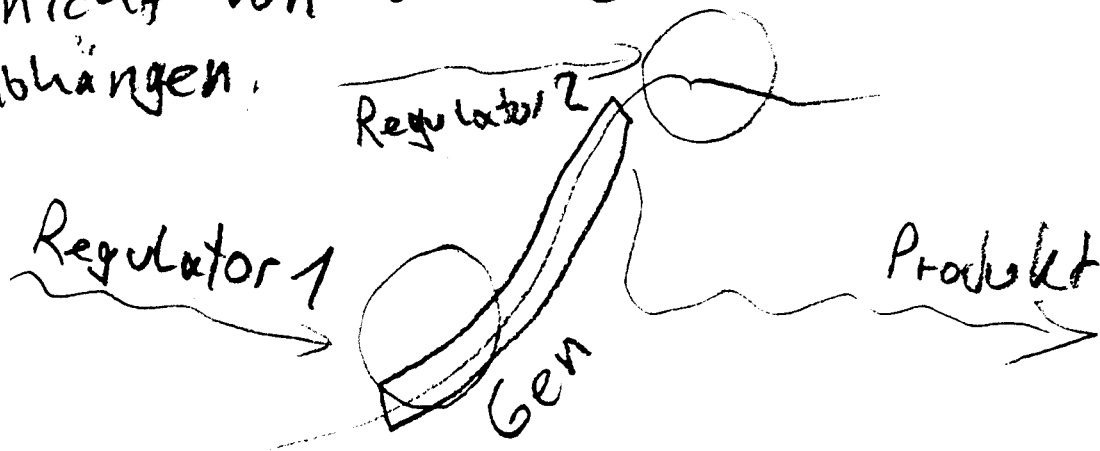
o Distanzmass:

$$De(s_i, s_j) = \sqrt{\sum_{t=1}^n (y_{it} - y_{jt})^2}$$

o Ergebnis ist ein Baum mit einer Darstellung der Expressionsverläufe an den Spitzen



- o Ramoni et al. berücksichtigen erstmals die Abhängigkeiten in einer Zeitreihe. (10)  
Andere Verfahren betrachteten die Messpunkte als unabhängig und machten Fehler.
- o Der autoregressive Ansatz benötigt unnötig viele unbekannte Parameter: Ein Gen kann nicht von beliebig vielen anderen Genen abhängen.



## Zusammenfassung:

(11)

- o Der Ansatz berücksichtigt die dynamische Struktur der Genexpression
- o Er reduziert die Zahl der Dimensionen durch Zusammenfassen ähnlicher Gene
- o Besser wäre das direkte Ableiten eines Netzwerks