

Prediction of **U**nderlying **L**atent **C**lasses via **K**-means and **H**ierarchical **C**lustering **A**lgorithm

Guan-Hua Huang, Su-Mei Wang and
Chung-Chu Hsu

07/07/2010

Breast cancer data

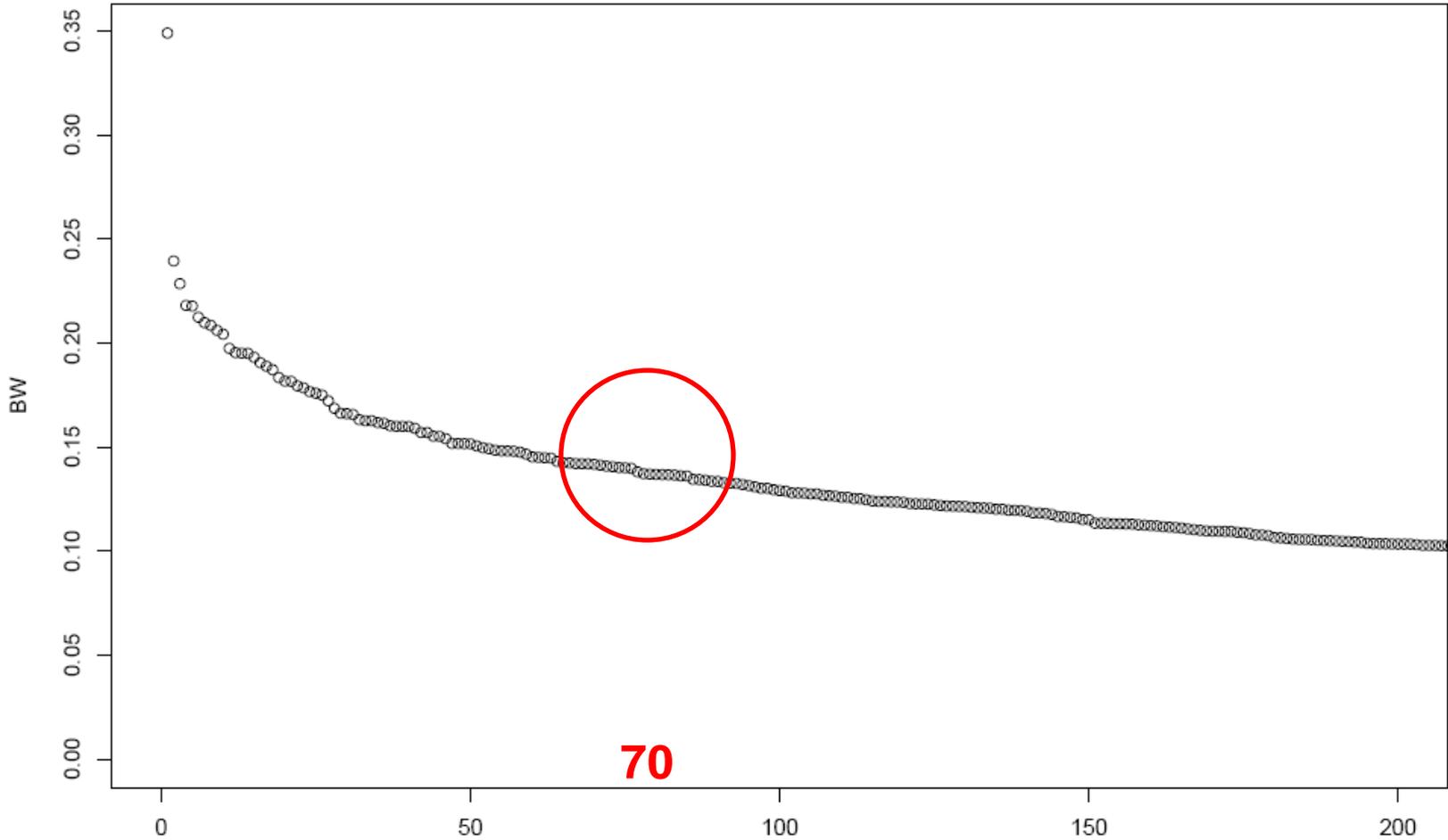
- van't Veer *et al.* Nature 2002
- The 78 sporadic lymph-node-negative breast cancer patients
 - 44 remained free of disease for an interval of at least 5 years (**good prognosis group**)
 - 34 had developed distant metastases within 5 years (**poor prognosis group**).
- Aim to predict good and poor prognostic patients through gene expression profiling

Breast cancer data (cont'd)

- A preliminary two-step gene selection process (from 24481 genes):
 - **4741** genes with the intensity ratio more than two-fold difference and the significance of regulation p-value < 0.01 in more than 3 patients
- Apply a selection of genes based on the ratio of their between-group to within-group sums of squares

$$BW(m) = \frac{\sum_i \sum_c I(d_i = c) (\bar{y}_{cm} - \bar{y}_{.m})^2}{\sum_i \sum_c I(d_i = c) (y_{im} - \bar{y}_{cm})^2}$$

BW plot



Breast cancer data (cont'd)

- Using 70 selected gene expression ratios as observed surrogates, a finite mixture model was fitted.

Schizophrenia data

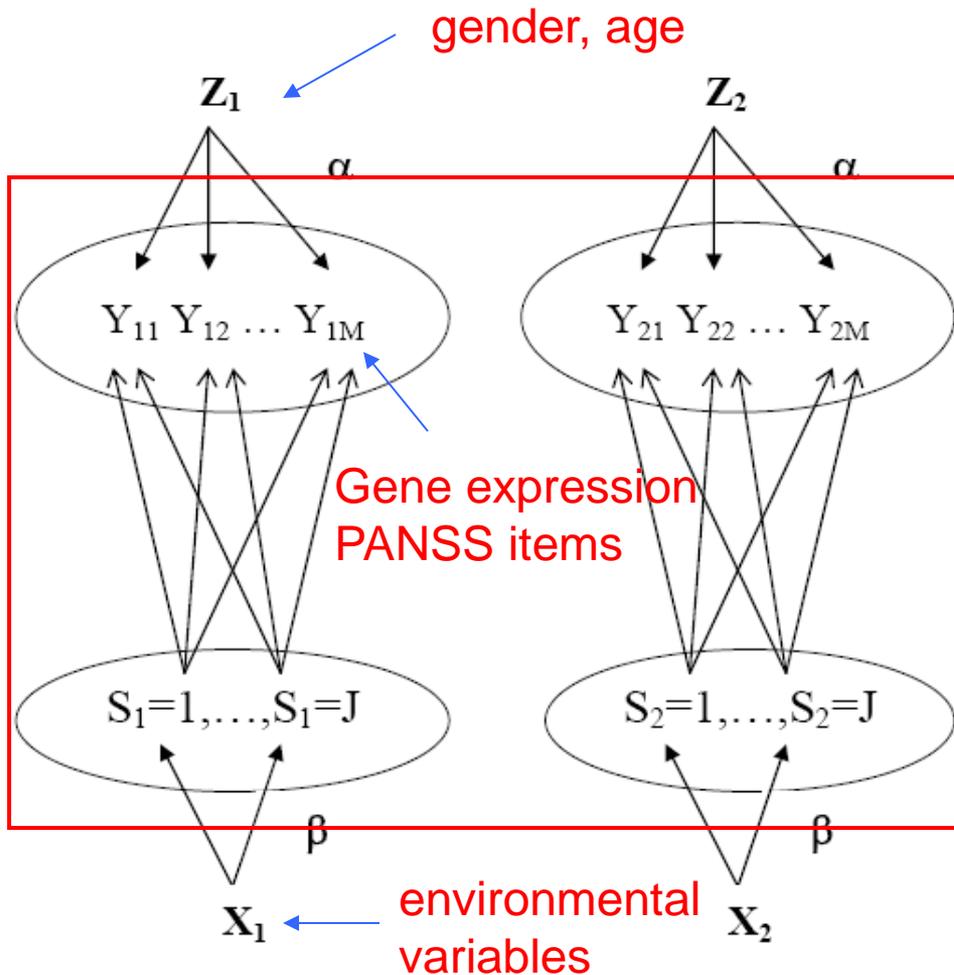
- The data were collected from a series of projects for schizophrenia (Dr. Hai-Gwo Hwu).
- The analyzed data include
 - **169 acute patients** of schizophrenia who were recruited within one week of index admission
 - **160 subsided state patients** who were living with community and under family care
- Aim to
 - explore the subtypes of schizophrenia patients
 - predict patients' phases of chronicity

Schizophrenia data (cont'd)

- Schizophrenia symptoms were assessed by the PANSS:
 - 30 items and consists of three subscales: positive, negative and general psychopathology
 - Each item was originally rated on a 7-point scale (1=absent, 7=extreme), but we reduced the 7-point scale by merging the points that had the response percentages less than 10%

Models

POPULATION (Size = N)



secondary
covariates

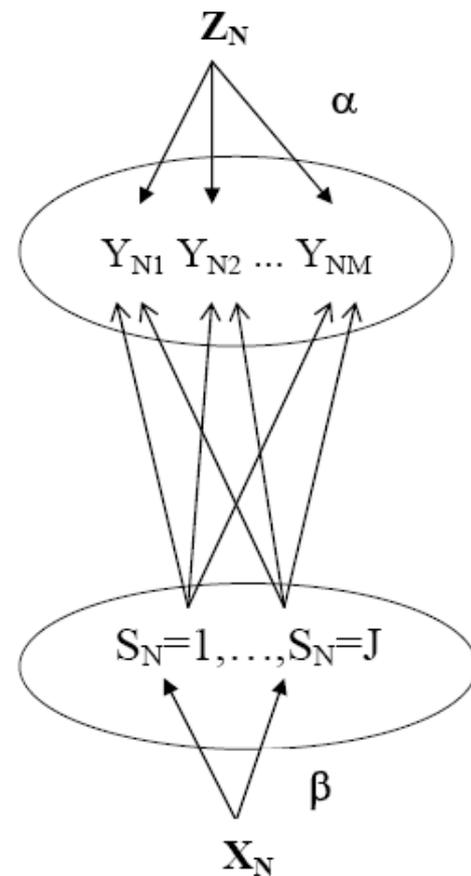
observed
indicators

...

latent class

...

primary
covariates



Introduction

- Finite mixture model is an **analogy of cluster analysis**.
- Finite mixture model **classifies objects** based on their responses to a set of surrogates.
- Measured surrogates are assumed **independent** of one another within any category of the underlying latent variable.
- Use **k-means** and **hierarchical** clustering methods with **covariance** among surrogates as the **distance measure**.

Finite mixture model

$\mathbf{Y}_i = (Y_{i1}, \dots, Y_{iM})^T$: M observable surrogates

$$\begin{aligned} f(y_{i1}, \dots, y_{iM}) &= \sum_{j=1}^J \left\{ \Pr(S_i = j) f(y_{i1}, \dots, y_{iM} \mid S_i = j) \right\} \\ &= \sum_{j=1}^J \left\{ \Pr(S_i = j) \prod_{m=1}^M f(y_{im} \mid S_i = j) \right\} \end{aligned}$$



Latent Class Membership Estimation

Background

- The key is to estimate the **latent class membership**.
- Use **K-means** and **hierarchical** clustering methods to group the objects such that observed variables are **statistically independent** within latent classes.
- Use **sample covariance** matrix as the independence measurement.

Independence measurement

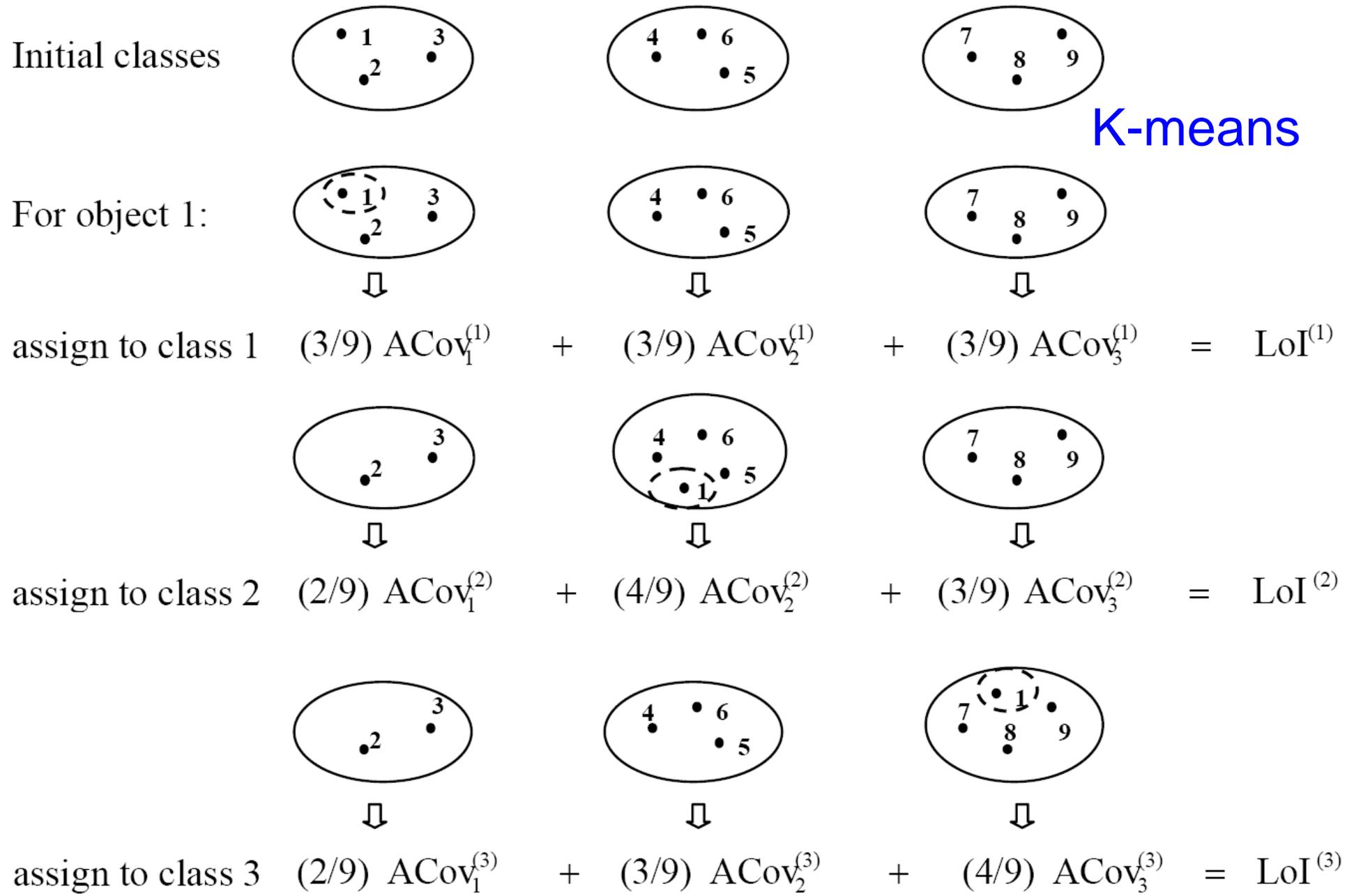
- Supposed $\tilde{\mathbf{Y}}_i = (Y_{i1}, Y_{i2}, \dots, Y_{iM})$

Then,

$$\text{Cov}(\tilde{\mathbf{Y}}_i) = \begin{bmatrix} \text{cov}(Y_{i1}, Y_{i1}) & \text{cov}(Y_{i1}, Y_{i2}) & \cdots & \text{cov}(Y_{i1}, Y_{iM}) \\ \text{cov}(Y_{i2}, Y_{i1}) & \text{cov}(Y_{i2}, Y_{i2}) & \cdots & \text{cov}(Y_{i2}, Y_{iM}) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(Y_{iM}, Y_{i1}) & \text{cov}(Y_{iM}, Y_{i2}) & \cdots & \text{cov}(Y_{iM}, Y_{iM}) \end{bmatrix}$$

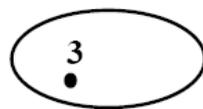
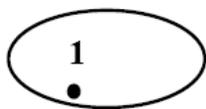
- $\text{ACov} = \text{mean}(| \text{entries in non-diagonal-block} |)$

K-means

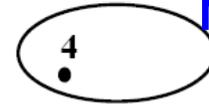
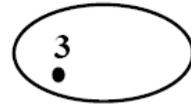


=> Assign object 1 to the class corresponding to minimum LoI

Initial

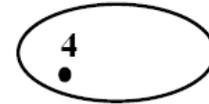
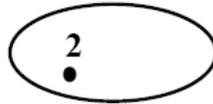


Agglomerative
hierarchical



merge 1, 2

$$(2/4)ACov_{(12)} + (1/4)ACov_3 + (1/4)ACov_4 = LoI^{(12)}$$



merge 1, 3

$$(2/4)ACov_{(13)} + (1/4)ACov_2 + (1/4)ACov_4 = LoI^{(13)}$$

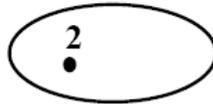
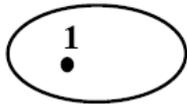
⋮

⋮

⋮

⋮

⋮

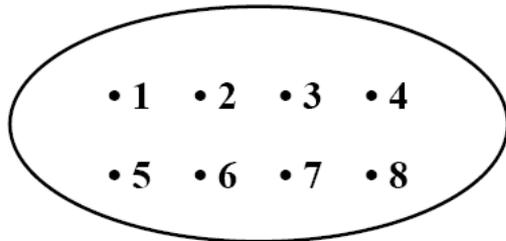


merge 3, 4

$$(1/4)ACov_1 + (1/4)ACov_2 + (2/4)ACov_{(34)} = LoI^{(34)}$$

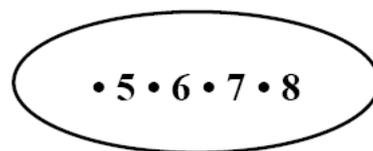
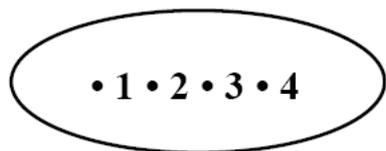
=> Merge the pair of classes whose combination results in the minimum LoI

Preliminary class

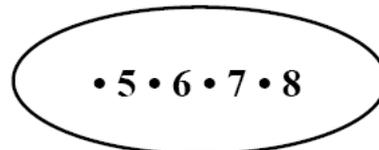


Divisive hierarchical

Split by 2-means

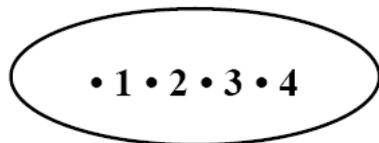


Split class 1 by 2-means



$$(2/8) \text{ACov}_{(1)_1} + (2/8) \text{ACov}_{(1)_2} + (4/8) \text{ACov}_2 = \text{LoI}^{(1)}$$

Split class 2 by 2-means



$$(4/8) \text{ACov}_1 + (2/8) \text{ACov}_{(2)_1} + (2/8) \text{ACov}_{(2)_2} = \text{LoI}^{(2)}$$

=> Split the class whose division results in the minimum LoI

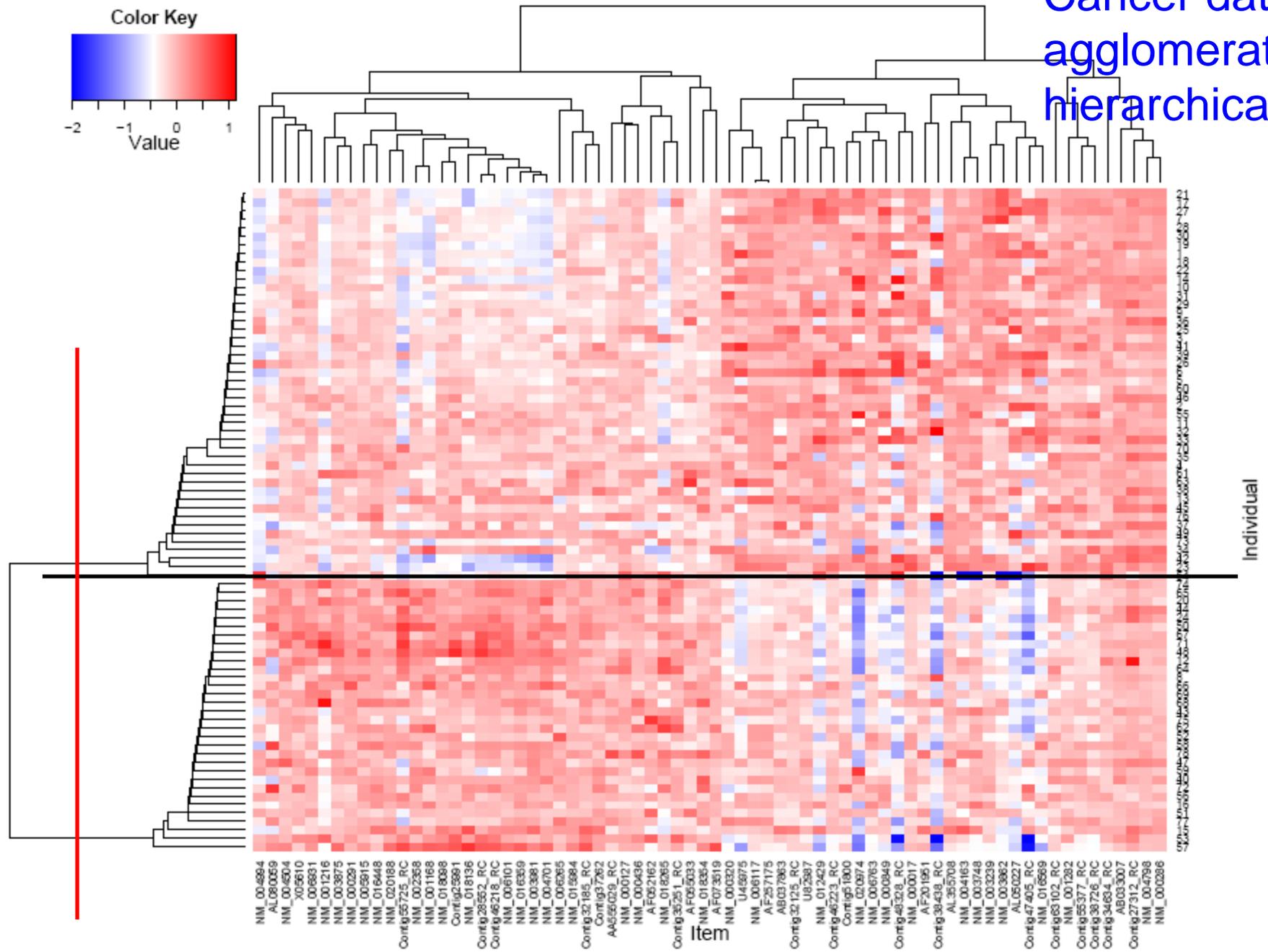
Classification using finite mixture models

- For a new object $Y^* = (Y_1^*, \dots, Y_M^*)$ with the disease status D^*

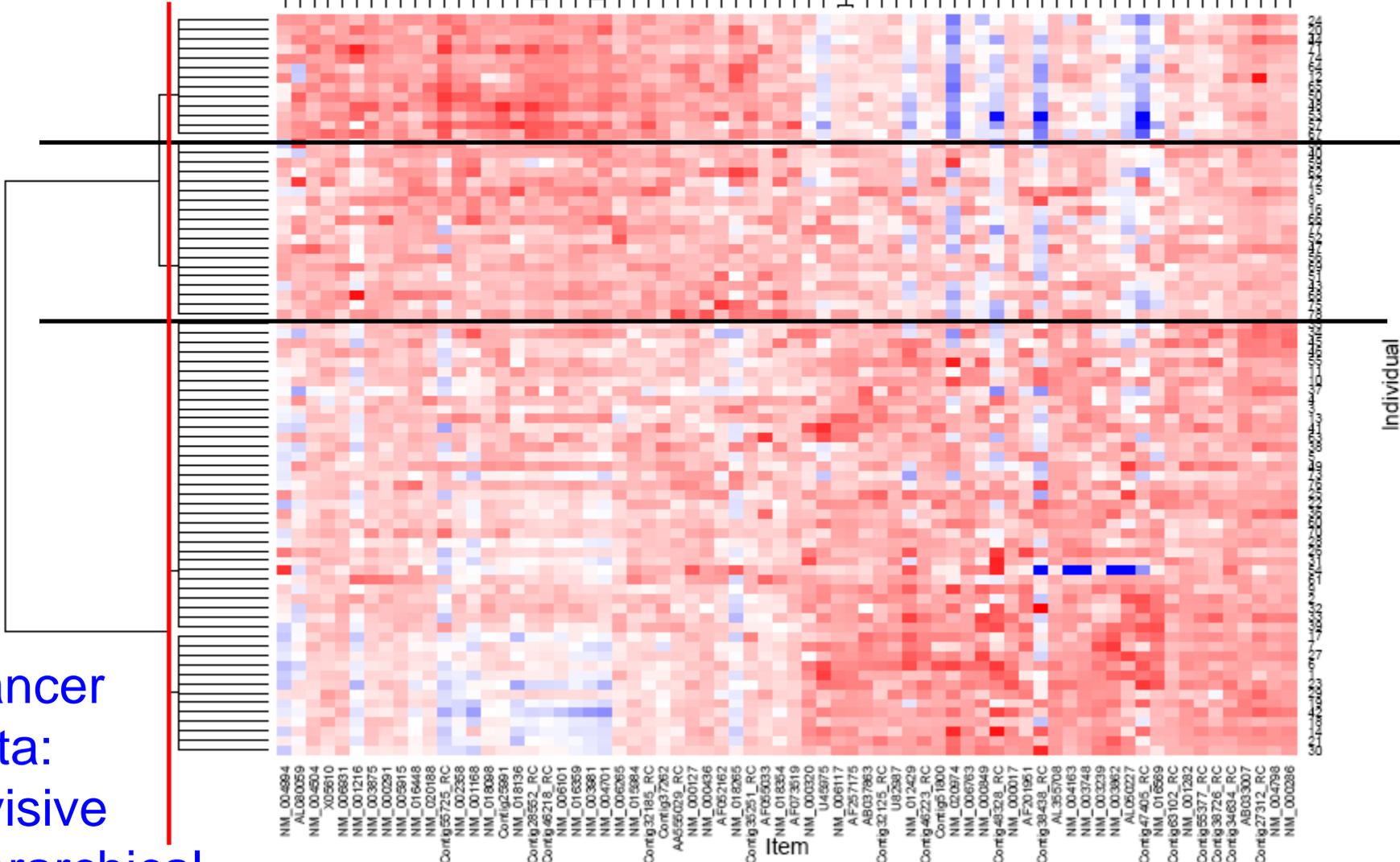
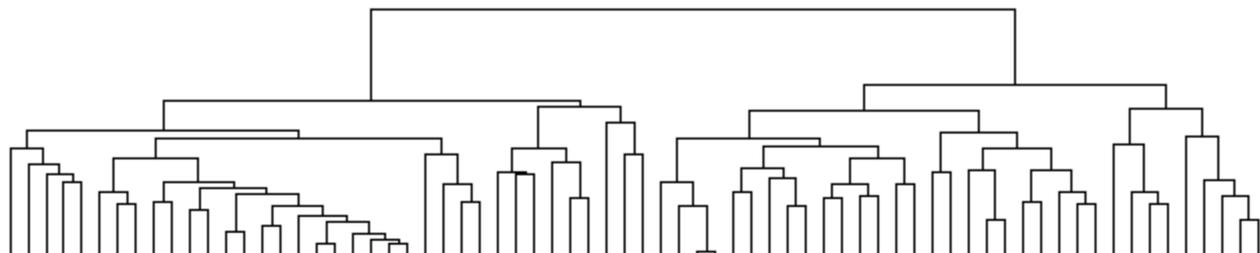
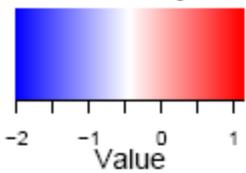
$$\Pr(D^* = c | Y^*) = \sum_{j=1}^J \left\{ \Pr(D^* = c | S^* = j, Y^*) \times \Pr(S^* = j | Y^*) \right\}$$

- Allocate Y^* to $D^*=c^*$ at which the maximum estimated posterior probability is reached

Cancer data:
agglomerative
hierarchical



Color Key



NM_004694
AL080059
NM_004504
X05610
NM_006831
NM_001216
NM_003875
NM_000291
NM_005915
NM_016448
NM_020188
Corrig55725_RC
NM_002358
NM_001168
NM_018088
Corrig25991
NM_018136
Corrig28552_RC
Corrig46218_RC
NM_006101
NM_016359
NM_003881
NM_004701
NM_006265
NM_016984
Corrig32185_RC
Corrig37262
AA055025_RC
NM_000127
NM_000436
AF052162
NM_018265
Corrig35251_RC
AF055033
NM_018354
AF073619
NM_000320
U45975
NM_006117
AF257175
AB037863
Corrig32125_RC
U82887
NM_012429
Corrig46223_RC
Corrig51800
NM_020974
NM_006763
NM_000949
Corrig49328_RC
NM_000017
AF201951
Corrig39439_RC
AL355708
NM_004163
NM_003748
NM_003239
NM_003862
AL050227
Corrig47405_RC
NM_016589
Corrig63102_RC
NM_001282
Corrig55377_RC
Corrig39726_RC
Corrig34683_RC
AB033007
Corrig27312_RC
NM_004798
NM_000286

Individual

Cancer
data:
divisive
hierarchical

Leave-one-out cross-validation

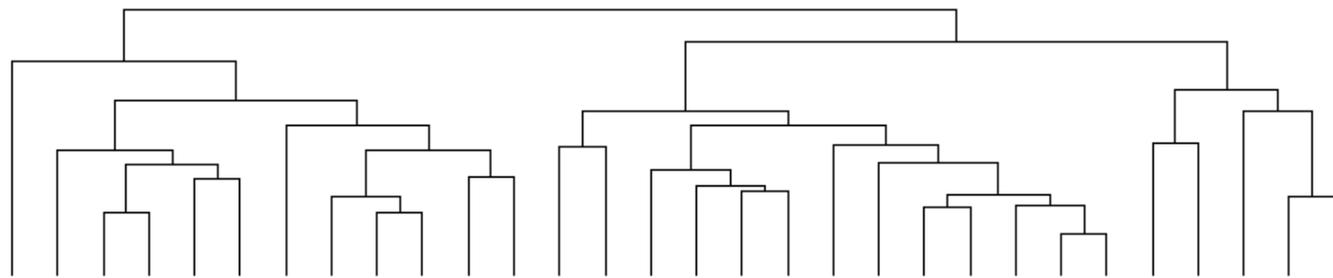
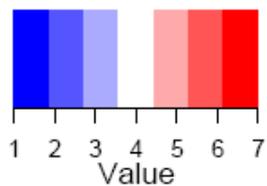
- Misclassification rates in predicting poor vs. good prognosis
 - k-means: 24.36%
 - agglomerative hierarchical: 26.92%
 - divisive hierarchical: 29.49%

Additional independent test set

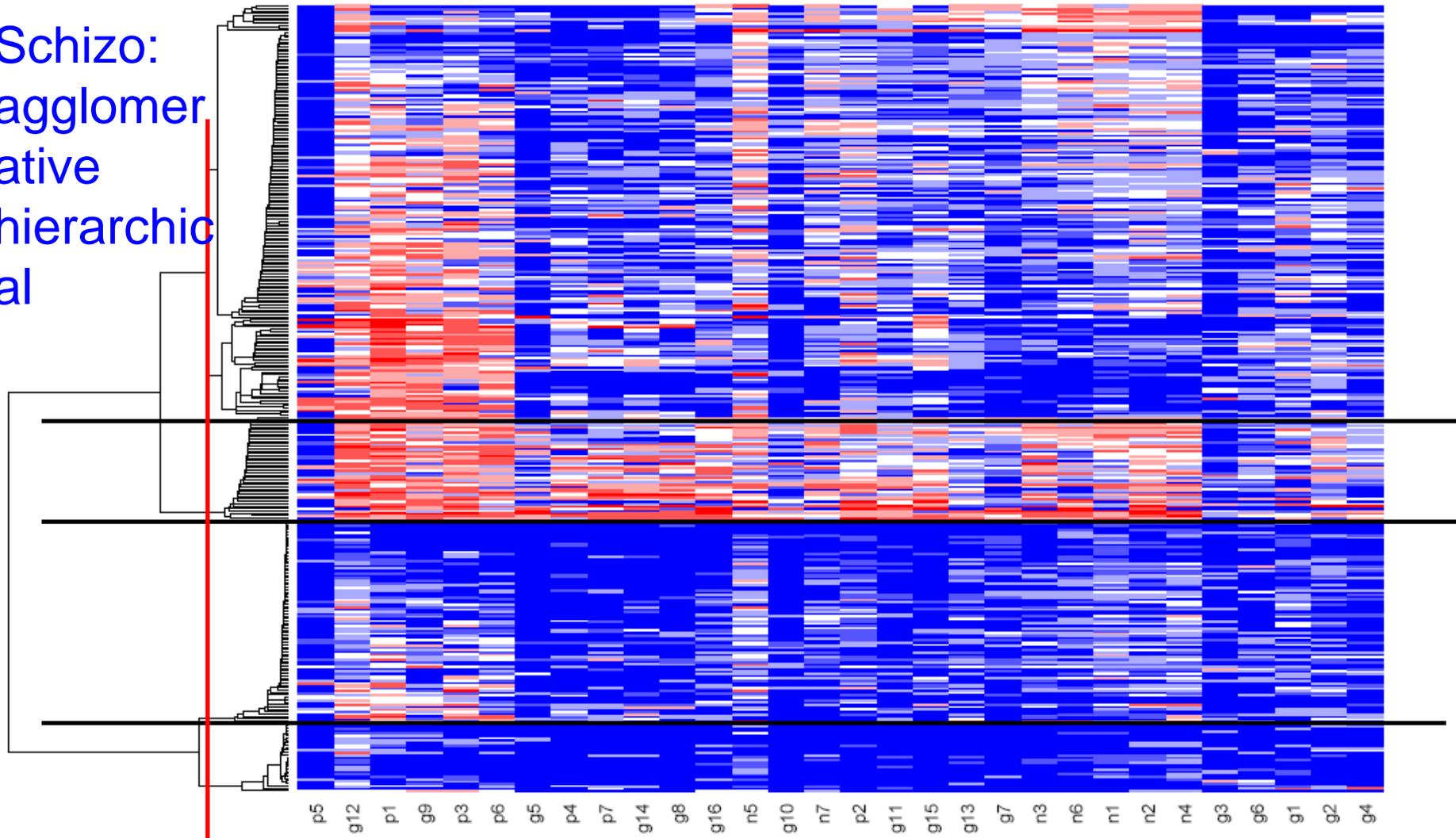
- Independent 19 young, lymph-node-negative breast cancer patients:
 - 12 poor prognosis
 - 7 good prognosis

No.	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
True	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1
KM	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	0	1
AH	0	0	0	1	0	0	0	0	0	0	0	1	1	1	1	1	1	0	1
DH	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	0	1	0	1

Color Key



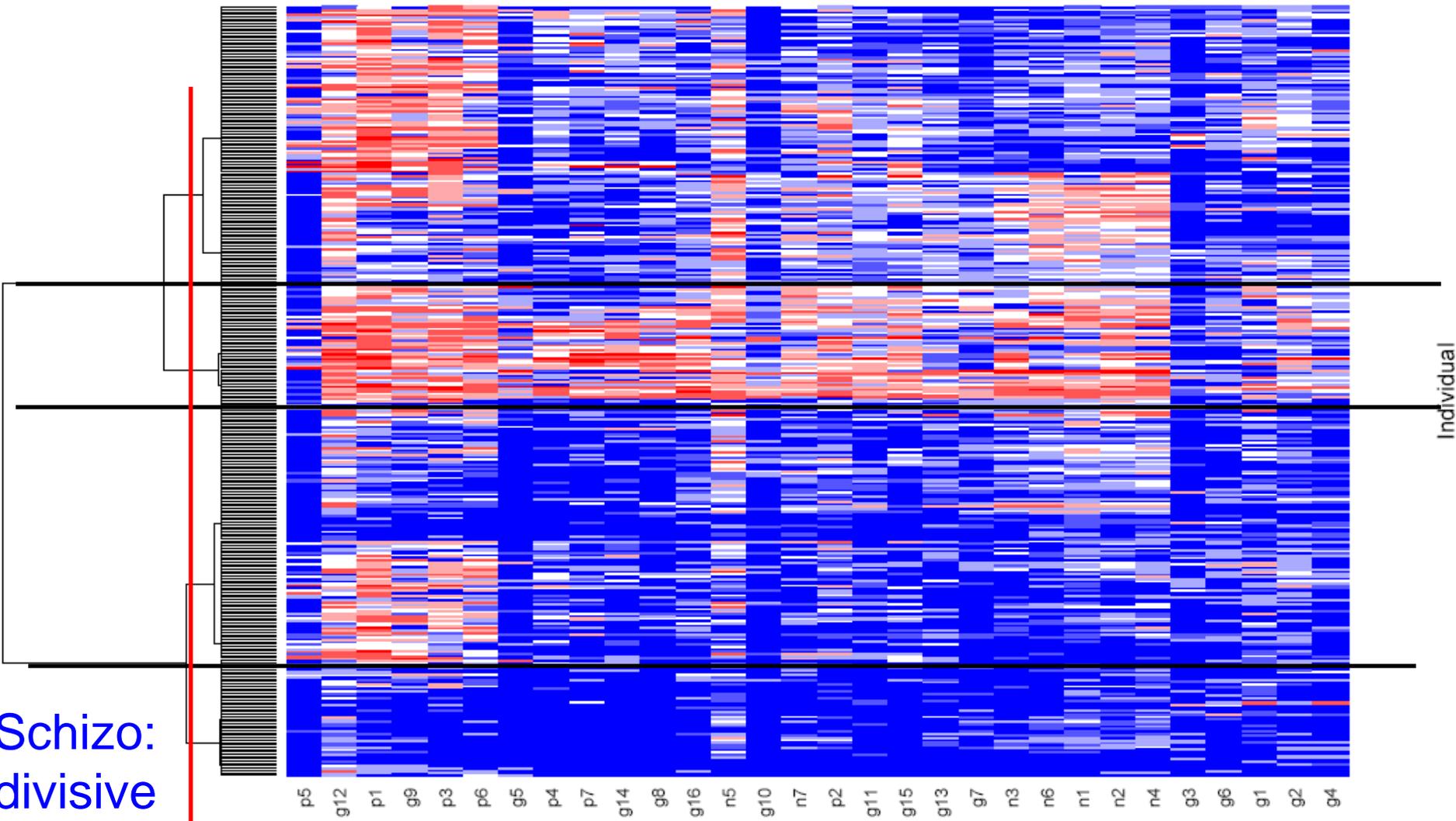
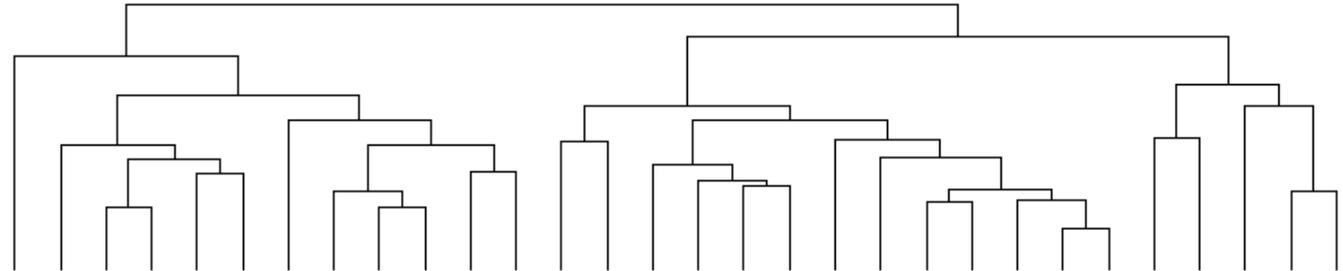
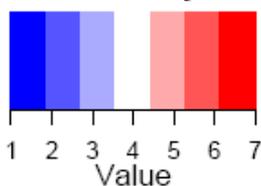
Schizo:
agglomerative
hierarchical



p5 g12 p1 g9 p3 p6 g5 p4 p7 g14 g8 g16 n5 g10 n7 p2 g11 g15 g13 g7 n3 n6 n1 n2 n4 g3 g6 g1 g2 g4

Individual

Color Key



Schizo:
divisive
hierarchical

Leave-one-out cross-validation

- Misclassification rates in predicting acute vs. subsided schizophrenia
 - k-means: 23.10%
 - agglomerative hierarchical: 24.01%
 - divisive hierarchical: 28.27%