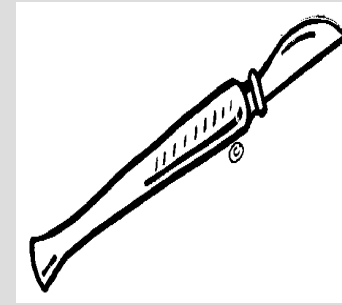

Diagnosis using computers

One disease



Three therapies



Clinical Studies

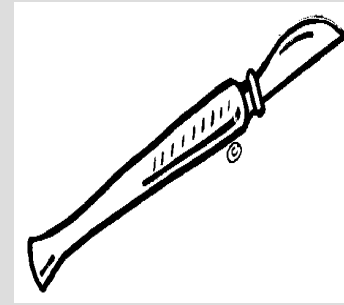
In average



75%



55%



35%

Success

Three subtypes of the disease



A



B



C



A



B



C



100%

60%

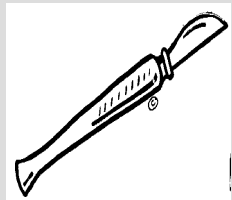
65%



40%

40%

85%



10%

90%

5%



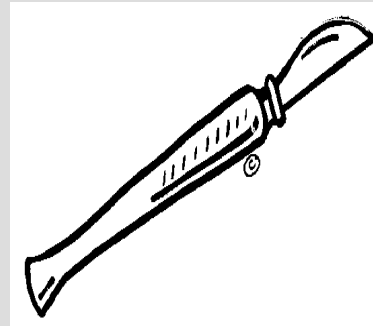
A



100%



B



90%

91,7%

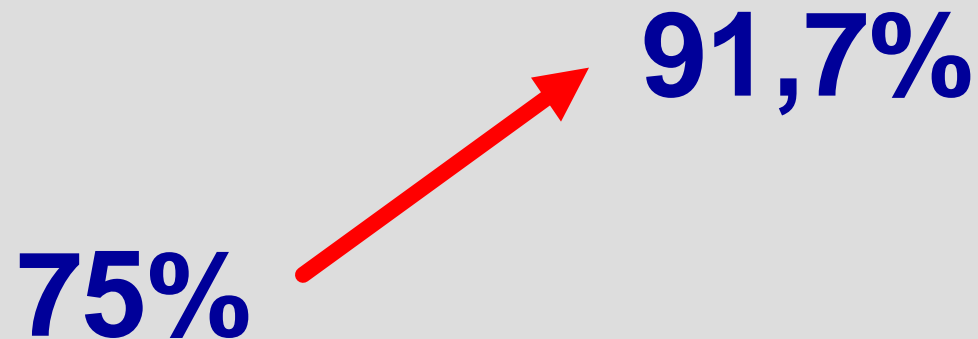


C



85%

Therapeutic success improved because of the refined diagnosis



Without developing any new
therapies

A higher resolution of dividing a disease into subtypes improves therapeutic success rates

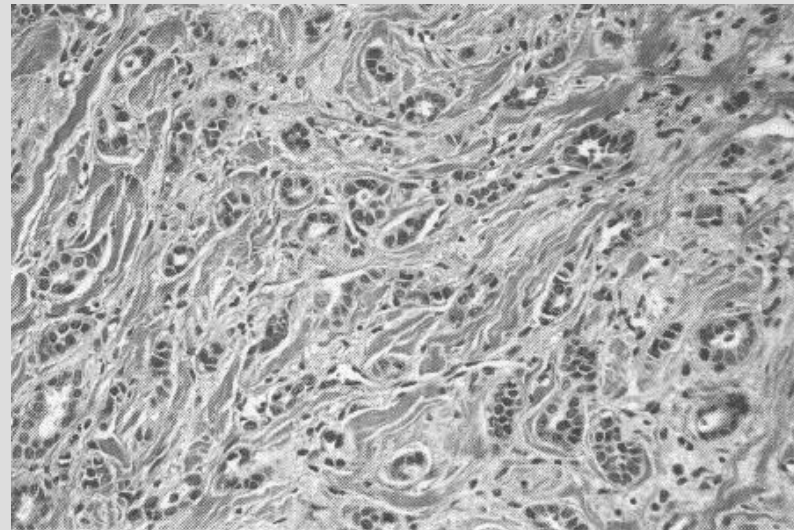


How do we obtain a higher resolution of diagnosis that is clinically relevant?



Looking at cells from outside

The microscope



Details of Metabolism:

The hemogram



Depression Patient		Female / Age: 29		Basic Status Report		Blood Test Date: 12/18/96	
Patient ID: 2373 (1047)		Request: 5716929				Practitioner: Dr. Patricia Kane (0)	
The % Status is the weighted deviation of the laboratory result.							
% Status	Result	Low	High				
98.25	H	2.88	0.80	2.00	A/G Ratio		
44.44	H	4.80	3.20	5.00	Albumin		
-26.19	L	45.00	20.00	125.00	Alkaline Phosphatase		
18.00		10.80	4.00	14.00	Anion Gap		
-33.33	L	10.00	7.00	25.00	B.U.N.		
-23.10		11.11	6.00	25.00	B.U.N./Creatinine Ratio		
-50.00	L	0.00	0.00	200.00	Basophil Count		
-50.00	L	0.00	0.00	3.00	Basophils		
-11.54		0.50	0.00	1.30	Bilirubin, Total		
-5.56		9.30	8.50	10.30	Calcium		
-28.65	L	2.51	2.30	3.30	Calcium/Phosphorus Ratio		
57.89	H	109.00	96.00	108.00	Chloride		
-80.00	L	104.00	140.00	200.00	Cholesterol		
-33.33	L	22.00	20.00	32.00	CO2		
-21.43		0.90	0.70	1.40	Creatinine		
32.40	H	462.00	50.00	550.00	Eosinophil Count		
50.00	H	6.00	0.00	6.00	Eosinophils		
4.17		2.70	1.40	3.80	Free T4 Index (T7)		
-36.67	L	6.00	0.00	45.00	GGT		
-65.00	L	1.90	2.20	4.20	Globulin		
-7.78		89.00	70.00	115.00	Glucose		
-48.95	L	36.00	35.00	130.00	HDL		
-36.36	L	36.50	35.00	46.00	Hematocrit		
-32.86	L	12.60	12.00	15.50	Hemoglobin		
-37.59	L	43.00	25.00	170.00	Iron, Total		
-6.80		198.00	0.00	250.00	LDH		
-58.82	L	96.00	62.00	130.00	LDL		
2.03		2541.00	850.00	4100.00	Lymphocyte Count		
0.00		33.00	18.00	48.00	Lymphocytes		
46.88	H	32.81	27.00	33.00	MCH		
13.01		34.52	32.00	36.00	MCHC		
25.26	H	85.05	80.00	100.00	MCV		
4.78		693.00	200.00	1100.00	Monocyte Count		
50.00	H	6.00	0.00	6.00	Monocytes		
-10.25		4004.00	1500.00	7800.00	Neutrophil Count		
-34.00	L	62.00	48.00	73.00	Neutrophils		
10.00		3.70	2.50	4.50	Phosphorus		
-33.33	L	3.80	3.50	5.30	Potassium		
-18.00		6.80	6.00	8.60	Protein, Total		
-54.62	L	3.84	3.90	5.20	R.B.C.		
-19.05		13.00	0.00	42.00	SGOT		
-25.00	L	12.00	0.00	48.00	SGPT		
-22.73		138.00	135.00	146.00	Sodium		
35.96	H	36.32	26.00	38.00	Sodium/Potassium Ratio		
26.92	H	32.00	22.00	35.00	T-3 Uptake		
-2.50		8.30	4.50	12.50	Thyroxine (T4)		
-26.00	L	60.00	0.00	250.00	Triglycerides		
-28.43	L	1.50	0.40	5.50	Ultra-Sensitive TSH		
-64.00	L	1.80	2.50	7.50	Uric Acid		
5.71		7.70	3.80	10.80	W.B.C.		
31.15		Total Status Deviation					
-10.13		Total Status Skew					

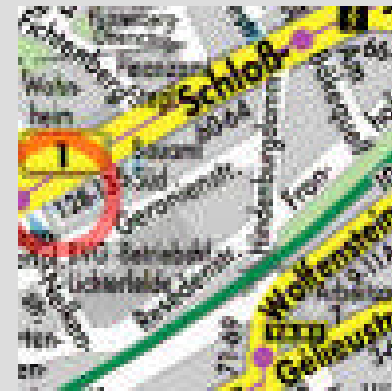
Diagnostics crabwise

- **Deregulation of metabolism causes disease**
- **Occasionally, they also lead to characteristic changes in tissue morphology or the hemogram.**



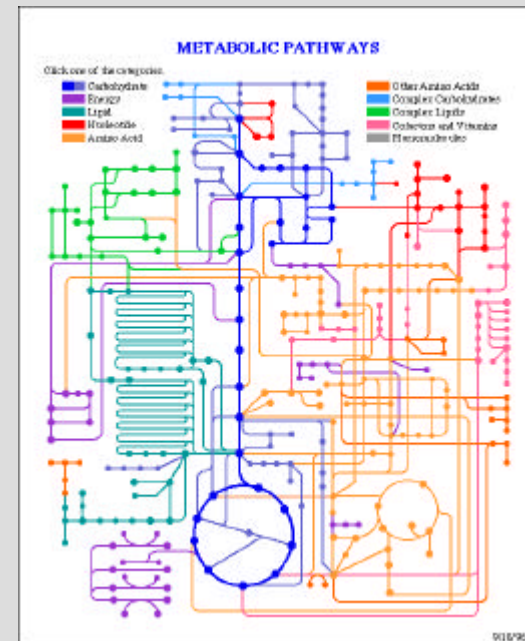
Diagnosics based on details

- A small number of genetic variations, transcription levels, and protein expression levels are routinely measured in single assays.

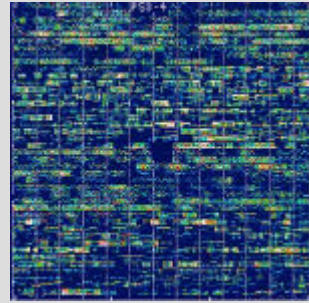
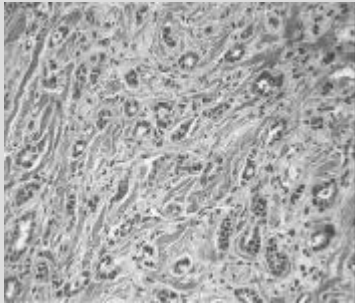


In addition desirable

- A patients metabolism in a bird's eye view



DNA Chip



Tissue

genome:~/ISIBC/original		
ER+Nevins4	d31628_s_at	253.3
ER+Nevins4	d31628_s_at	1386.0
ER+Nevins4	d31628_s_at	209.5
ER+Nevins4	d31716_at	655.3
ER+Nevins4	d31716_at	116.5
ER+Nevins4	d31716_at	596.3
ER+Nevins4	d31716_at	119.5
ER+Nevins4	d31762_at	573.3
ER+Nevins4	d31762_at	104.7
ER+Nevins4	d31762_at	507.8
ER+Nevins4	d31762_at	88.1
ER+Nevins4	d31763_at	698.0
ER+Nevins4	d31763_at	149.9
ER+Nevins4	d31763_at	593.3
ER+Nevins4	d31763_at	115.8
ER+Nevins4	d31764_at	2993.5
ER+Nevins4	d31764_at	426.6
ER+Nevins4	d31764_at	2882.8
ER+Nevins4	d31764_at	508.0
ER+Nevins4	d31765_at	846.5
ER+Nevins4	d31765_at	140.1
ER+Nevins4	d31765_at	1039.5
ER+Nevins4	d31765_at	207.3

**Expression
profile**

**Ok, what is the
problem ?**



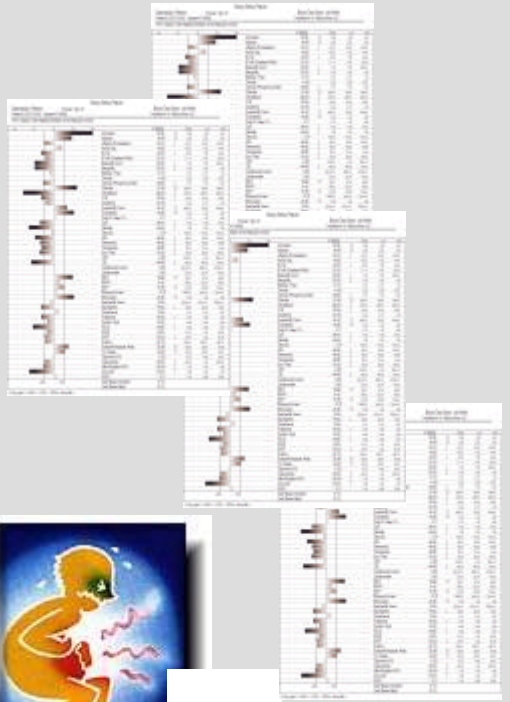
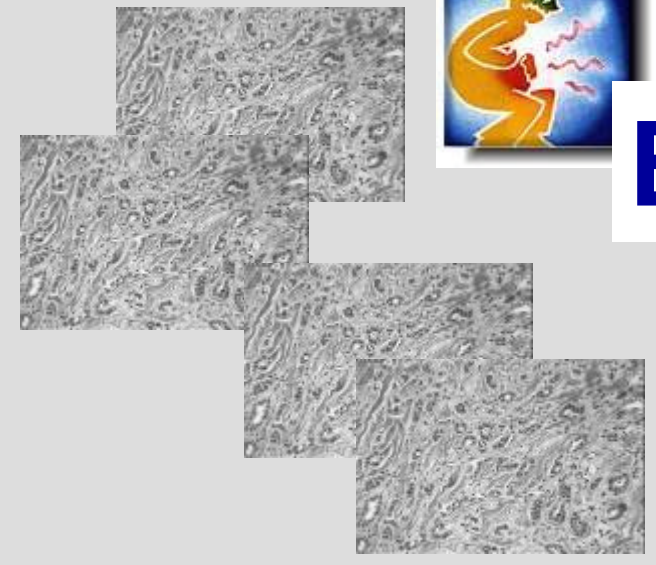
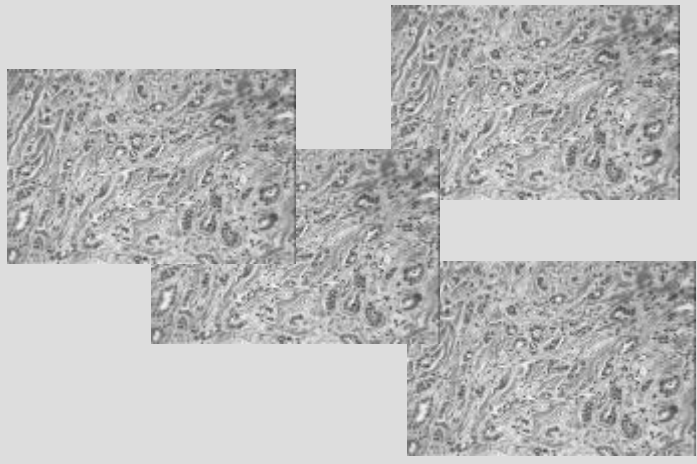
Morphological differences and differences in single assay measurements are the basis of classical diagnosis



A



B



**Are there any differences
between the gene
expression profiles of type
A patients and type B
patients?**

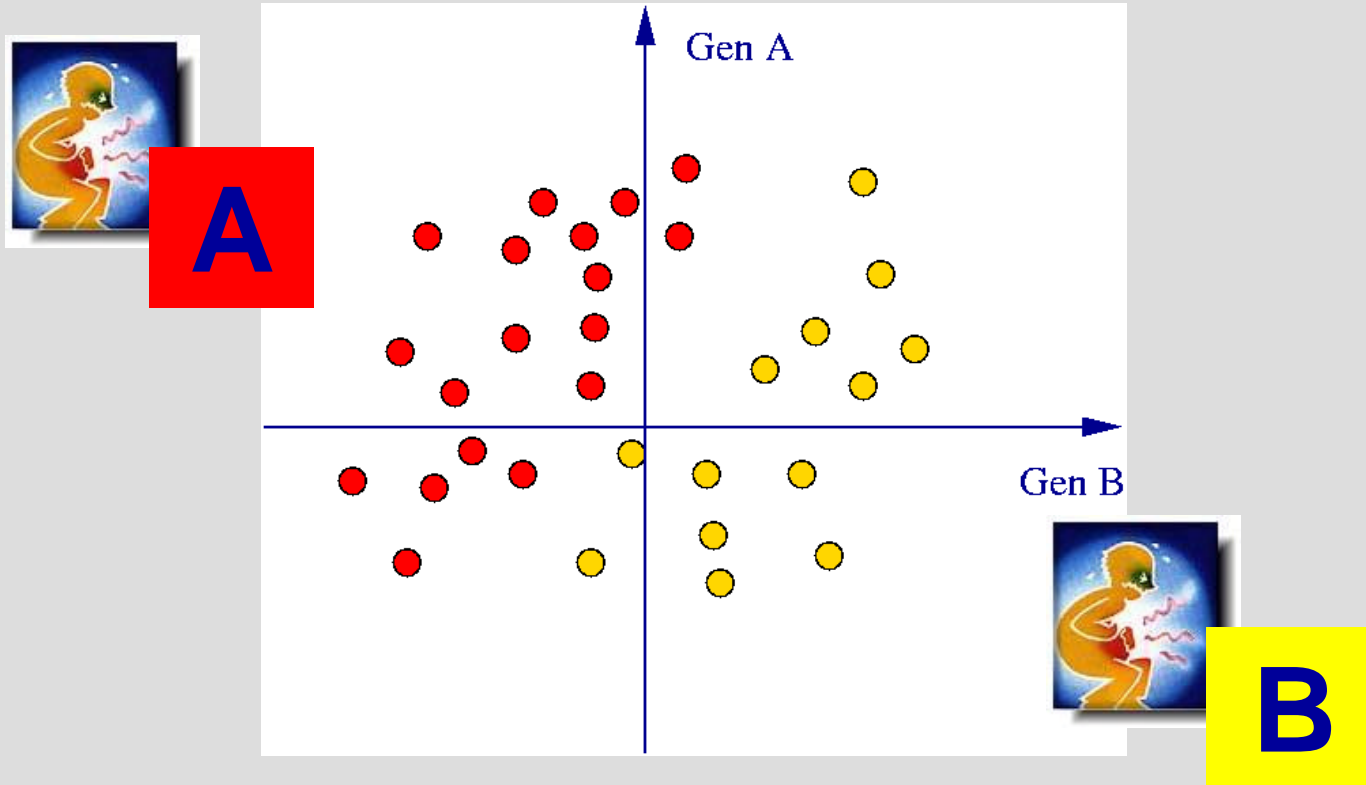
**30.000 genes are a lot.
That's too complex to start
with**

**Let's start with considering
only two genes:**

gene A und gene B

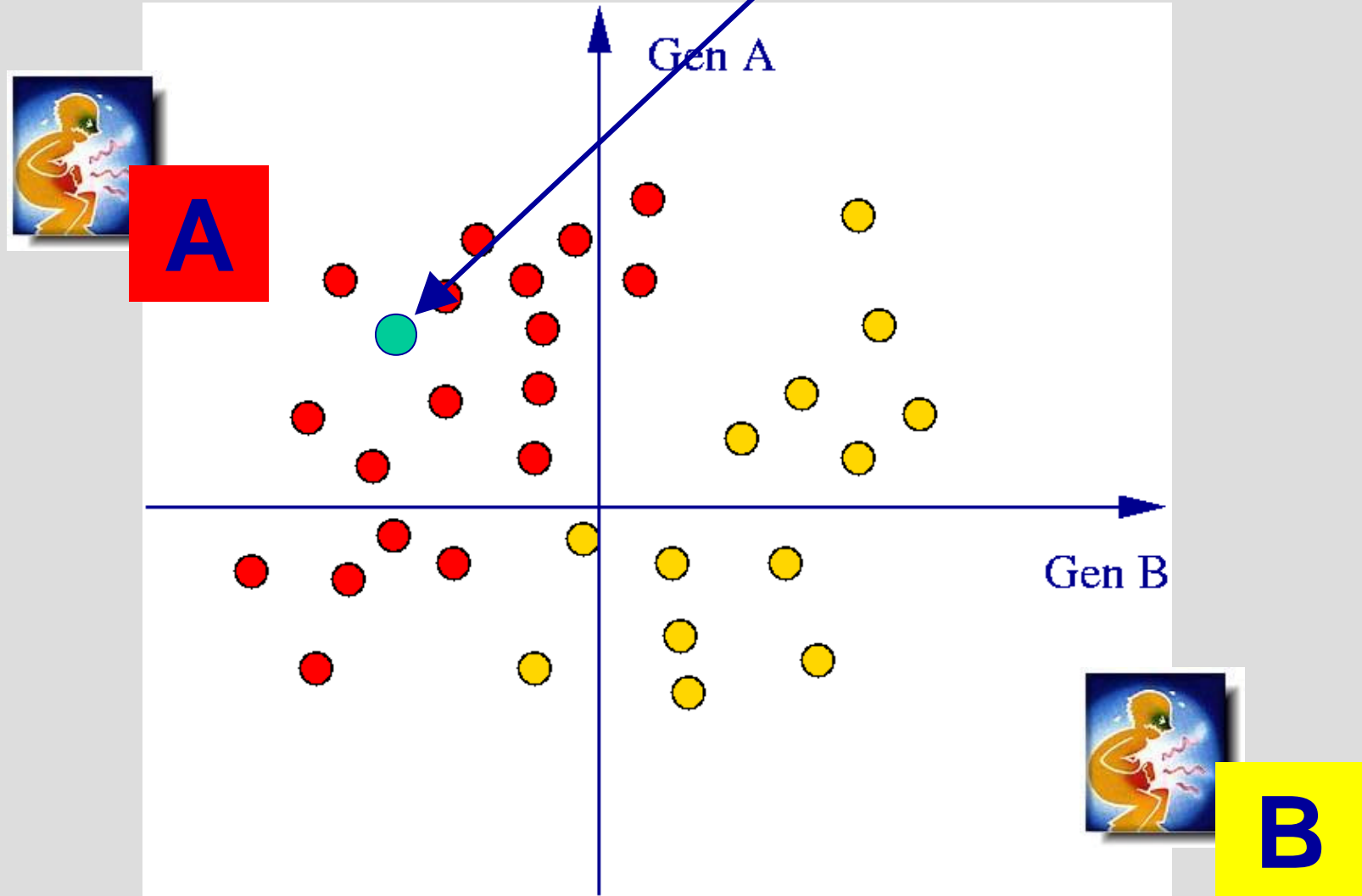


In this situation we can see that ...



... there is a difference.

A new patient



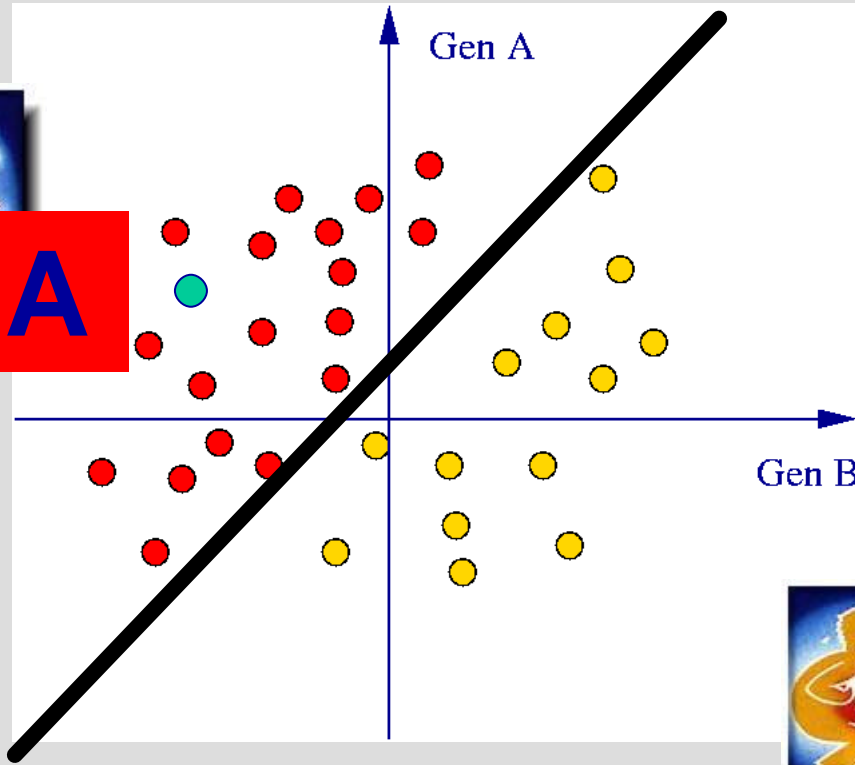


The new patient

A

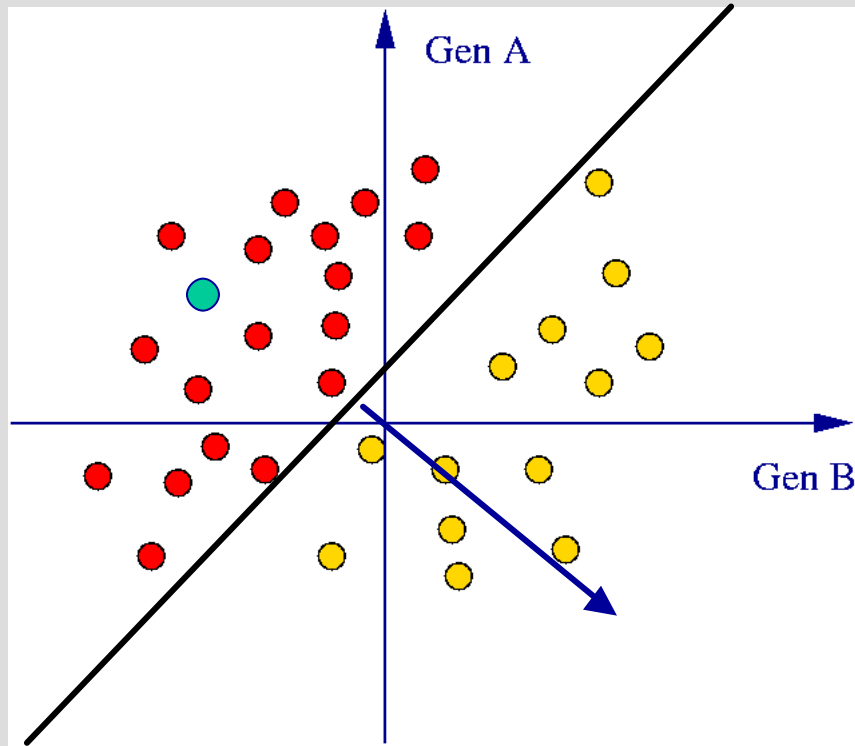


A



B

Here everything is clear.



The normal vector of the separating line can be used as a signature

.... the separating line is not unique

What exactly do we mean if we talk about signatures?

x_1, \dots, x_{30000} : expression levels

$f(x_1, \dots, x_{30000})$: Mapping that assigns one number to the expression levels

High values of f indicate class 1

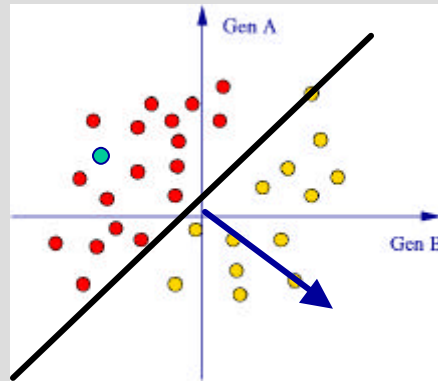
low values class 2

Example:

$$f(x_1, \dots, x_{30000}) = x_1 \quad \text{gene 1 is the signature}$$

Or a normal vector is the signature:

$$f(x_1, \dots, x_{30000}) = b_0 + b_1x_1 + b_2x_2$$

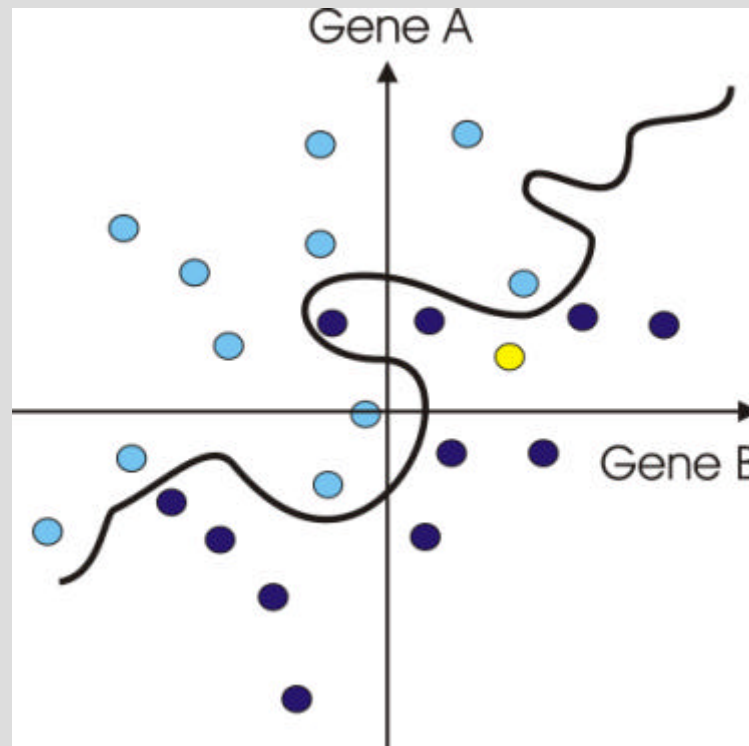


if x_1 and x_2 are the two genes in the Diagram

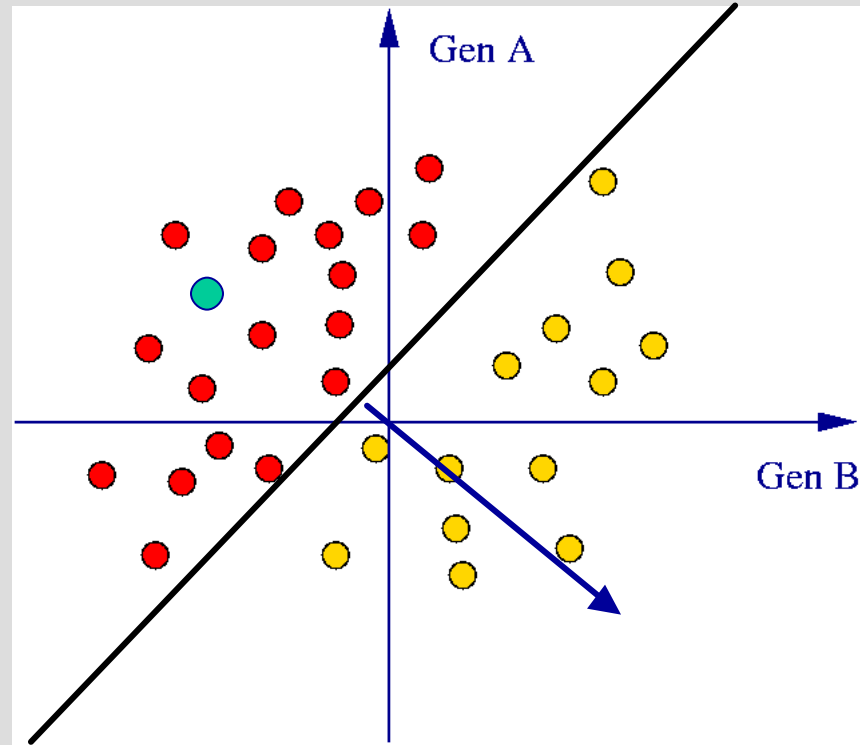
Using all genes yields:

$$f(x_1, \dots, x_{30000}) = b_0 + \sum_{i=1}^{30000} b_i x_i$$

Or you choose a very complicated signature



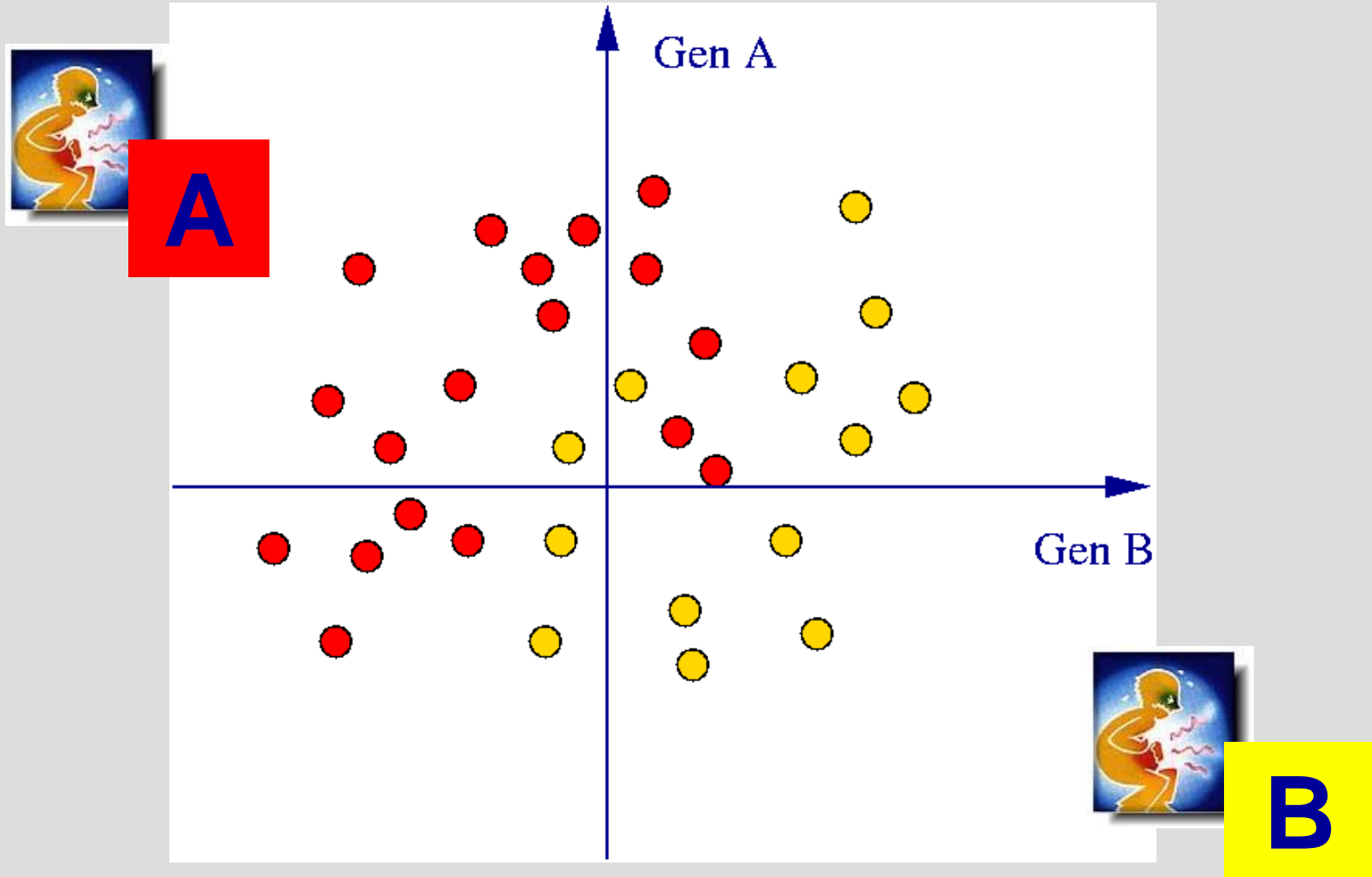
$$f(x_1, \dots, x_{30000}) = \text{complicated}$$



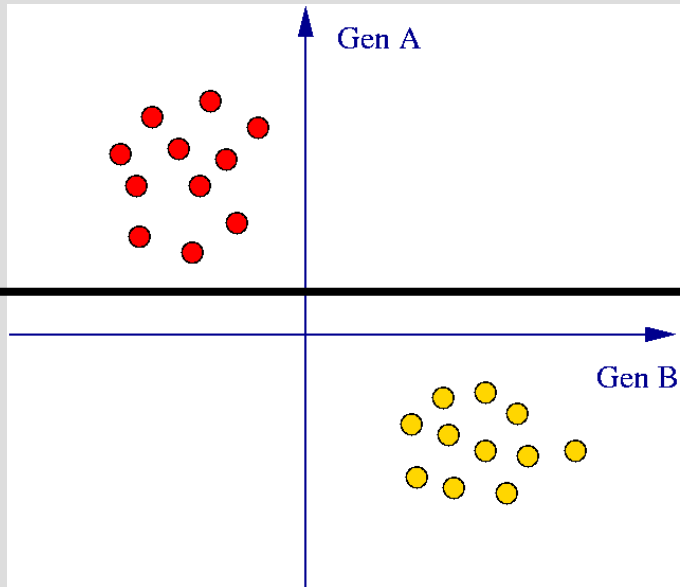
Unfortunately, expression data is different.

What can go wrong?

There is no separating straight line



Gene A is important



A

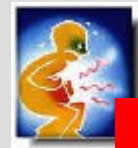
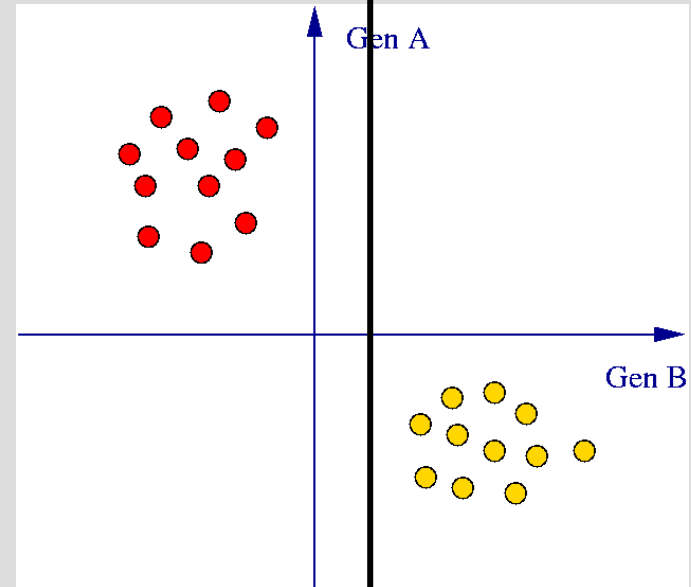
Gene A high



B

Gene A low

Gene B is important



A

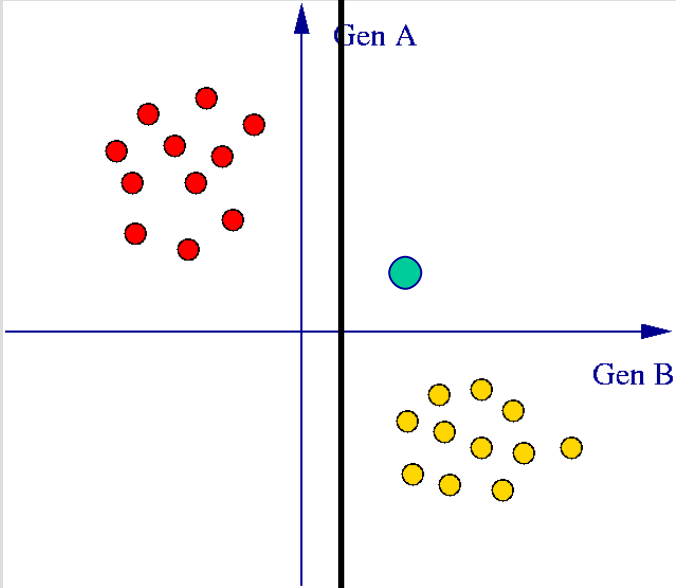
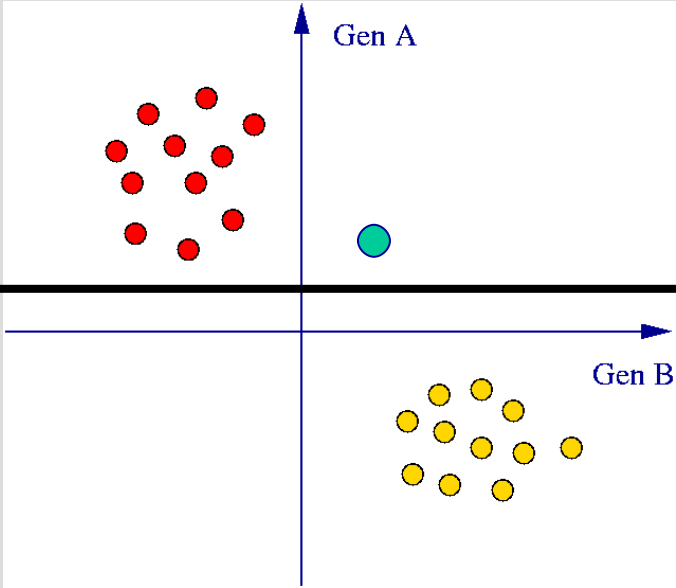
Gene B low



B

Gene B high

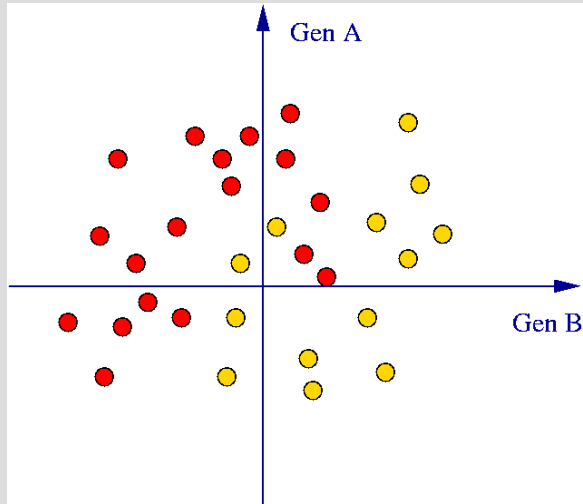
New patient ?



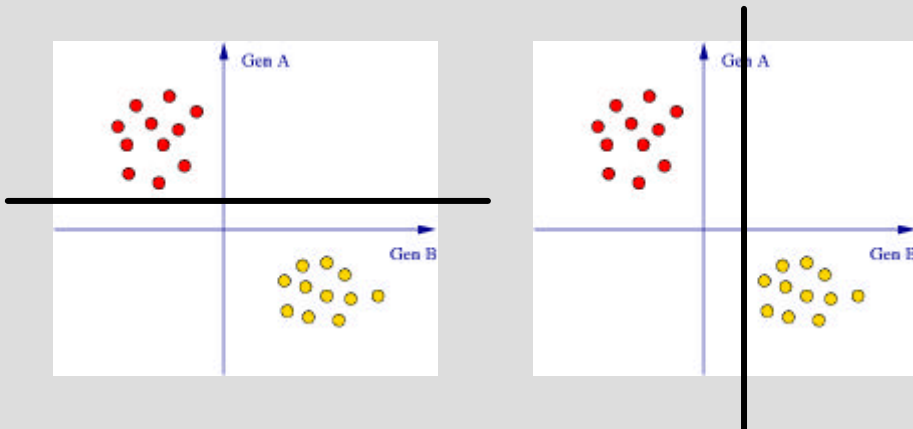
A



B

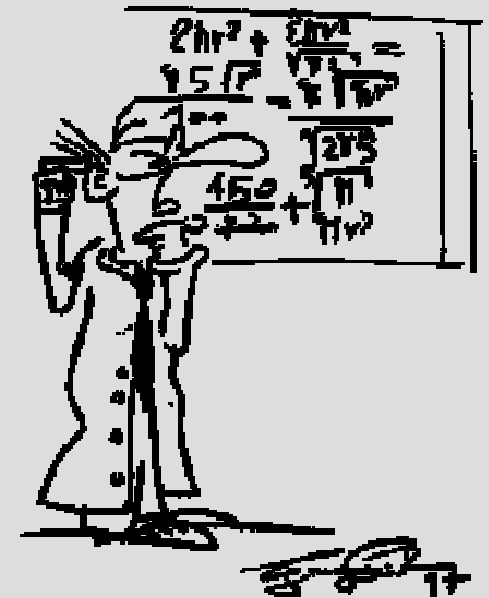
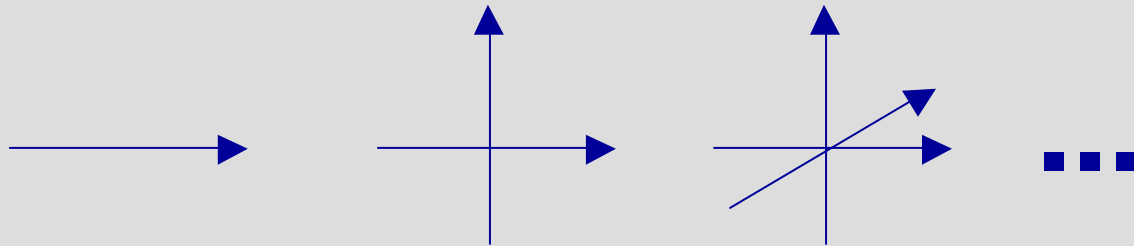


Problem 1:
No separating line

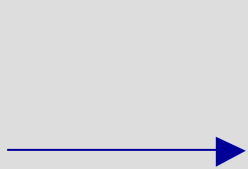


Problem 2:
To many separating lines

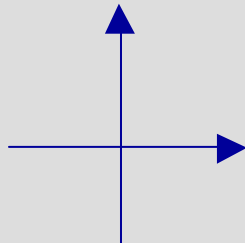
In praxis we look at thousands of genes, generally more genes than patients



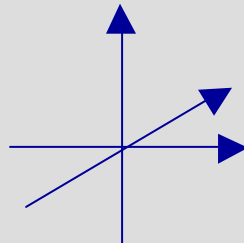
And in 30000 dimensional spaces different laws apply



1



2



3

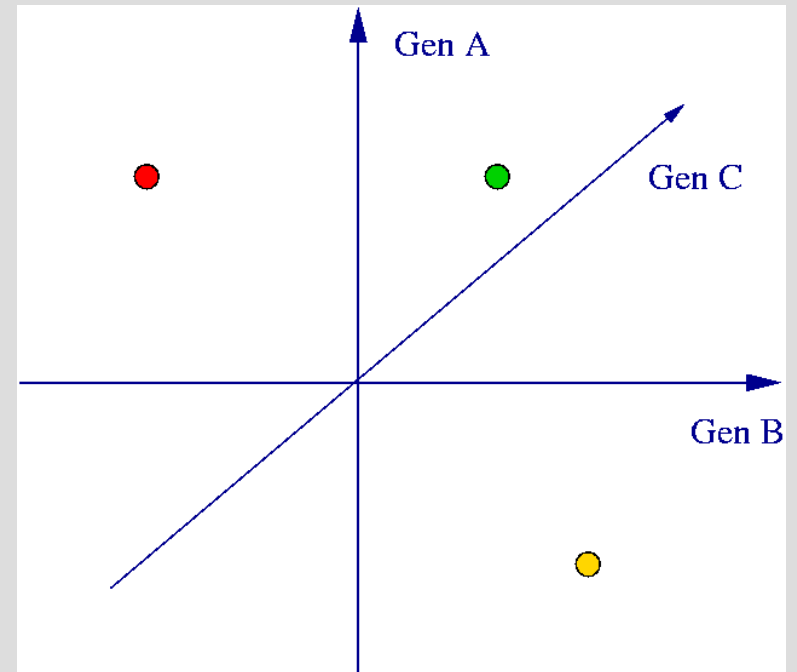
...

30000

- **Problem 1 never exists!**
- **Problem 2 exists almost always!**

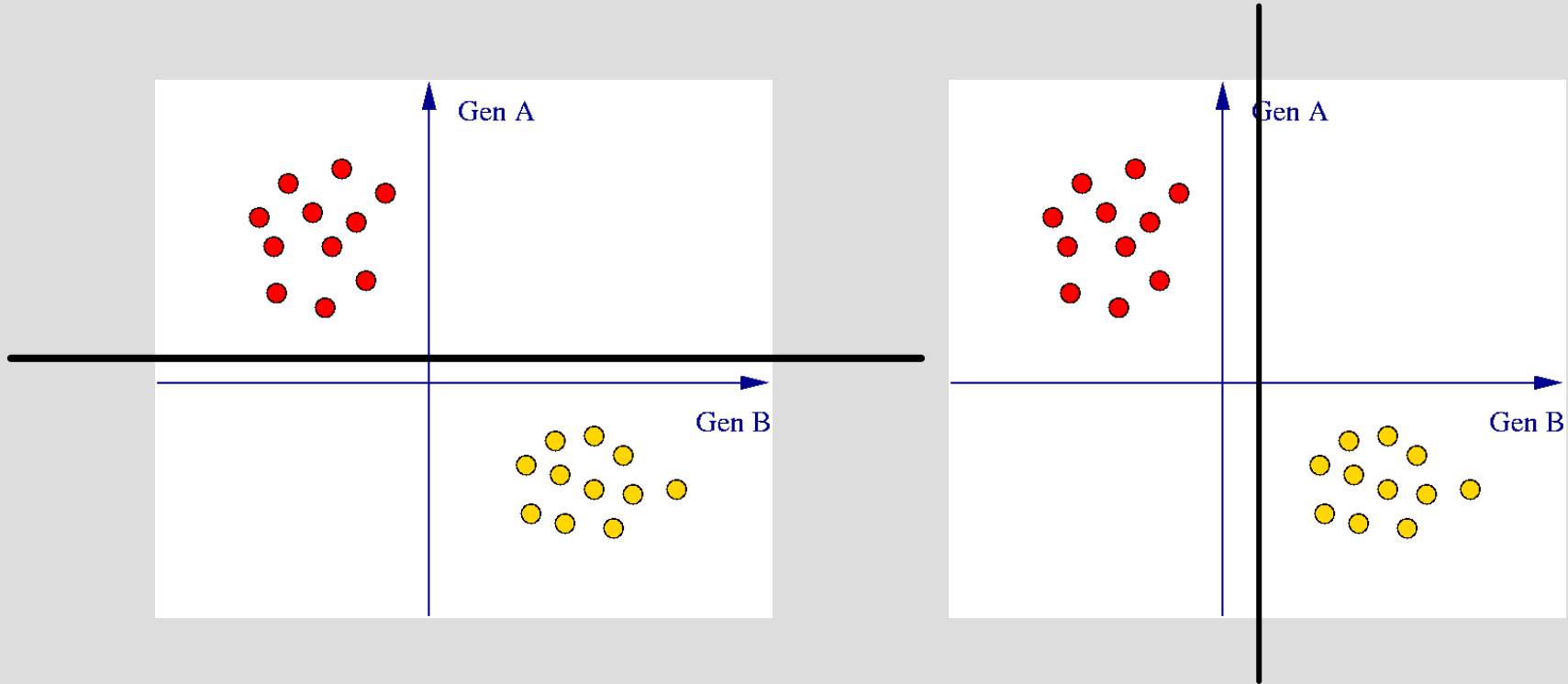
Spent a minute thinking about this in three dimensions

Ok, there are three genes, two patients with known diagnosis, one patient of unknown diagnosis, and separating planes instead of lines



OK! If all points fall onto one line it does not always work. However, for measured values this is very unlikely and never happens in praxis.

With more gene than patients the following problem exists:

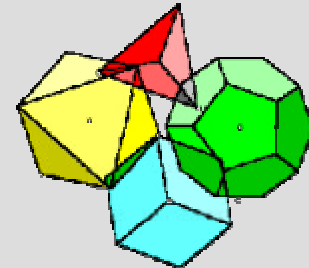
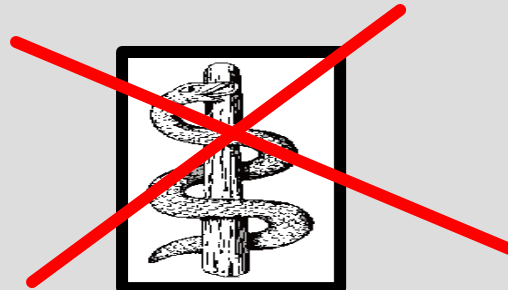


Hence for microarray data it always exists

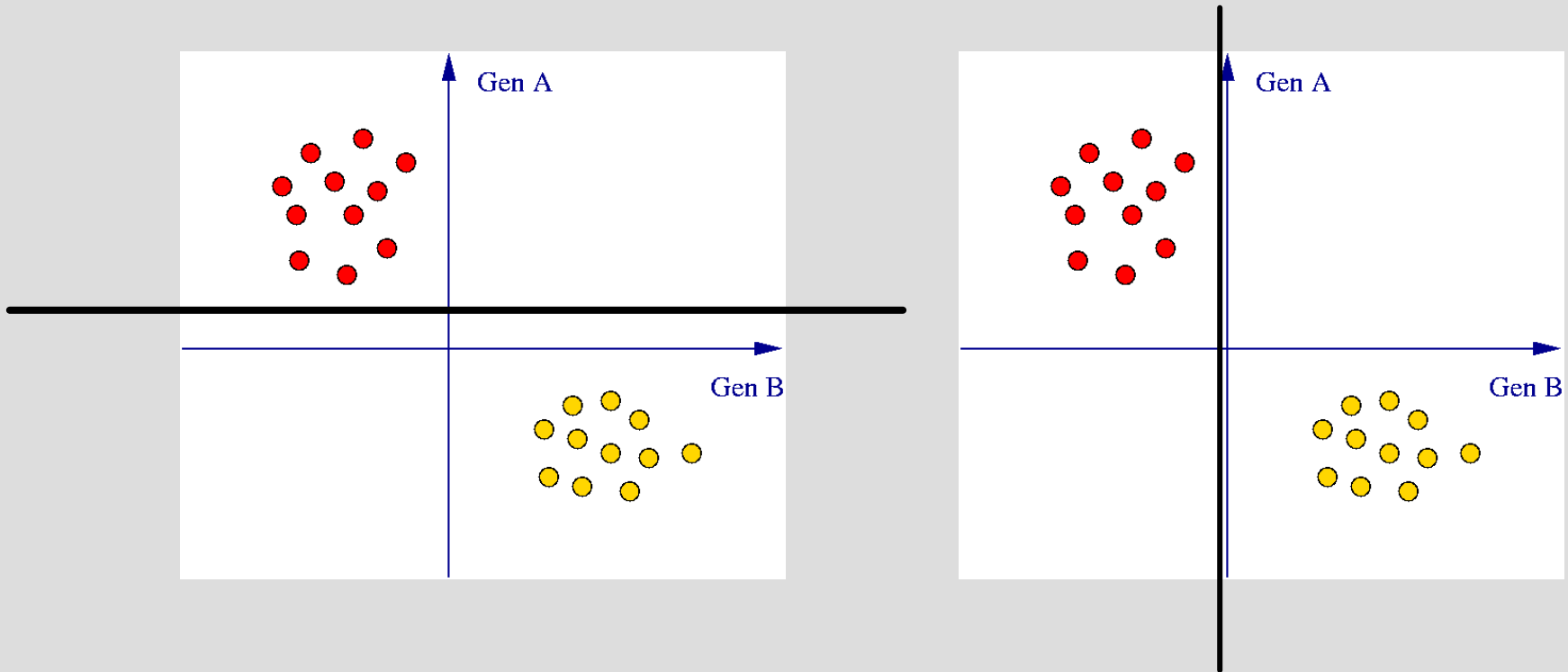
From the data alone we can not decide which genes are important for the diagnosis, nor can we give a reliable diagnosis for a new patient



This has little to do medicine. It is a geometrical problem.



Whenever you have expression profiles from two groups of patients, you will find differences in their genes expression ...



... no matter how the groups are defined .

There is a guarantee that you find a signature:

- which separates malignant from benign tumors

- but also



Müllers from Schmidts



- or using an arbitrary order of patients odd numbers from even numbers

In summary:

If you find a separating signature, it does not mean (yet) that you have a nice publication ...

... in most cases it means nothing.



Wait! Believe me!

There are meaningful differences in gene expression. And these must be reflected on the chips.



Ok,OK...

On the one hand we know that there are completely meaningless signatures and on the other hand we know that there must be real disorder in the gene expression of certain genes in diseased tissues

How can the two cases be distinguished?



**What are
characteristics of
meaningless
signatures?**

They come in large numbers

Parameters have high variances

Under-determined models

We have searched in a huge set of possible signatures

No regularization

They reflect details and not essentials

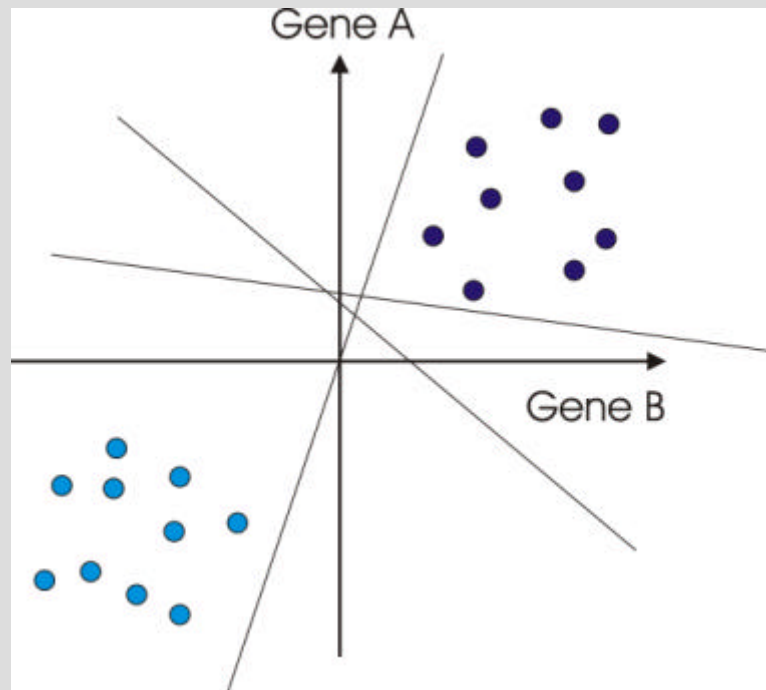
Overfitting



Under-determined models

They come in large numbers

Parameters have high variances



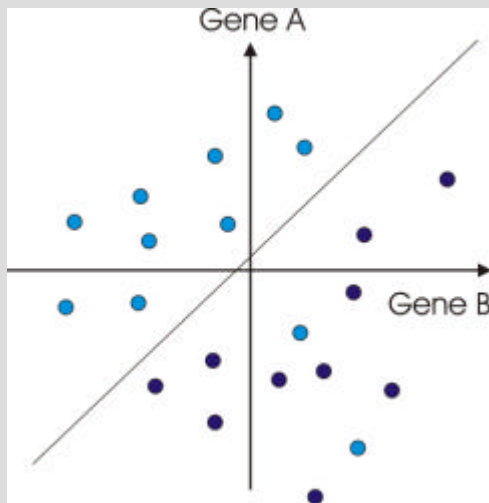
No regularization

We have searched in a huge set of possible signatures

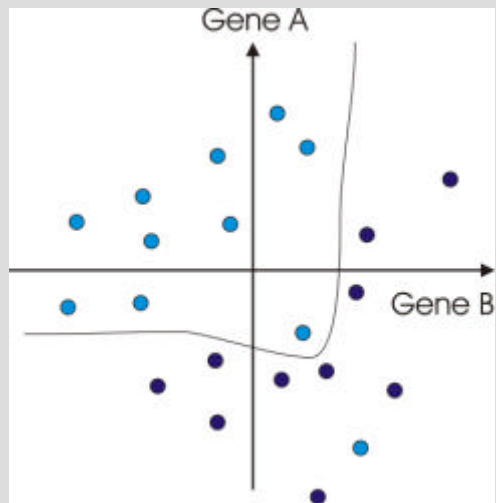
When considering all possible separating planes there must always be one that fits perfectly, even in the case of no regulatory disorder

Overfitting

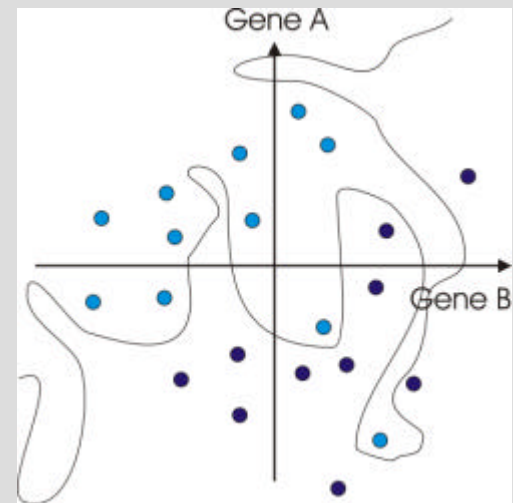
They reflect details and not essentials



2 errors



1 error



no errors

Signatures do not need to be perfect

Examples for sets of possible signature

- All quadratic planes
- All linear planes
- All linear planes depending on at most 20 genes
- All linear planes depending on a given set of 20 genes

High probability for finding a fitting signature

Low probability that a signature is meaningful



Low probability for finding a fitting signature

High probability that a signature is meaningful

What are strategies for finding meaningful signatures?

Later we will discuss 2 possible approaches

- Gene selection followed by linear discriminant analysis, and the PAM program
- Support Vector Machines

What is the basis for this methods?

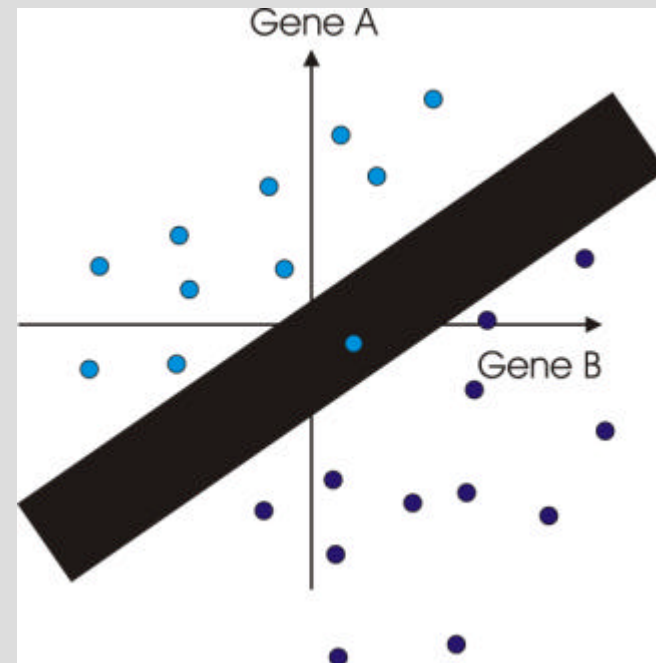
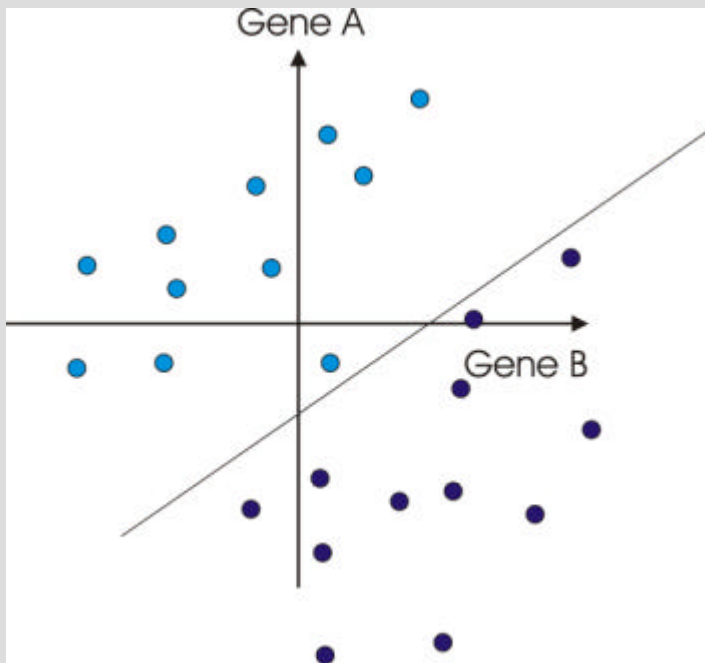


Gene selection

When considering all possible linear planes for separating the patient groups, we always find one that perfectly fits, without a biological reason for this.

When considering only planes that depend on maximally 20 genes it is not guaranteed that we find a well fitting signature. If in spite of this it does exist, chances are good that it reflects transcriptional disorder.

Support Vector Machines



Fat planes: With an infinitely thin plane the data can always be separated correctly, but not necessarily with a fat one.

Again if a large margin separation exists, chances are good that we found something relevant.

Large Margin Classifiers

Both gene selection and Support Vector Machines confine the set of a priori possible signatures. However, using different strategies.

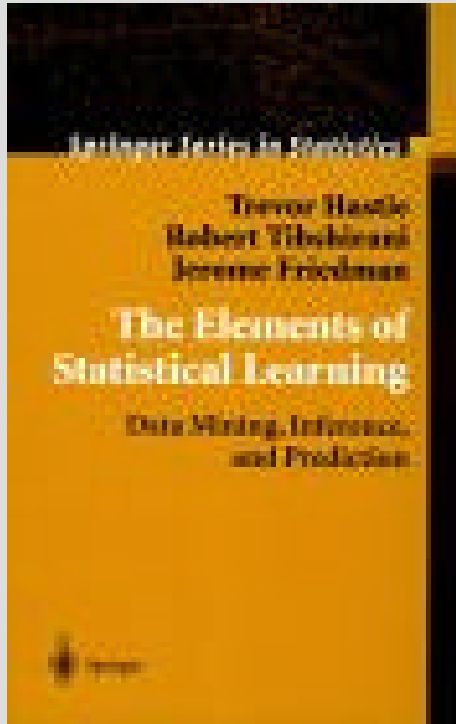
Gene selection wants a small number of genes in the signature (sparse model**)**

SVMs want some minimal distance between data points and the separating plane (large margin models**)**

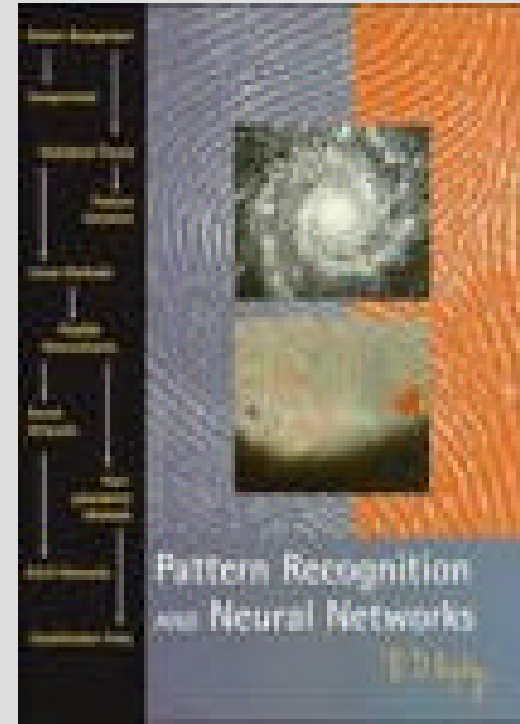
There is more than you could do ...

Learning Theory

Ridge Regression, LASSO, Kernel based methods, additive Models, classification trees, bagging, boosting, neural nets, relevance vector machines, nearest-neighbors, transduction etc. etc.



**The Elements of
Statistical Learning**
Hastie, T. Tibshirani,
R. Friedman, J



**Pattern
Recognition and
Neural Networks**
Brian D. Ripley

Questions



Coffee



Learning Methods

Setup

We have 200 patient profiles and 30000 genes on the chip

Patients can be divided into two groups according to some clinical or pathological criterion. There are 100 patients in each group.

The group distinction is not derived from the expression data

Problem: Can we reconstruct the group assignments from the expression profiles?

Consider a single gene first

a_1, \dots, a_{100} expression levels in group a

b_1, \dots, b_{100} expression levels in group b

$$\bar{a} = \frac{1}{100} (a_1 + \dots + a_{100})$$

$$\bar{b} = \frac{1}{100} (b_1 + \dots + b_{100})$$

c expression level of a patient
with unknown diagnosis

Compare $|c - \bar{a}|$ and $|c - \bar{b}|$

Diagnosis : a if $|c - \bar{a}| < |c - \bar{b}|$

b if $|c - \bar{a}| \geq |c - \bar{b}|$

Both groups are summarized by the mean gene expression in this

Diagnosis is according to the closest mean

Consider two genes:

$a_{1,1}, \dots, a_{1,100}, a_{2,1}, \dots, a_{2,100}$ group a

$b_{1,1}, \dots, b_{1,100}, b_{2,1}, \dots, b_{2,100}$ group b

$$\bar{a} = (\bar{a}_1, \bar{a}_2)$$

$$\bar{b} = (\bar{b}_1, \bar{b}_2)$$

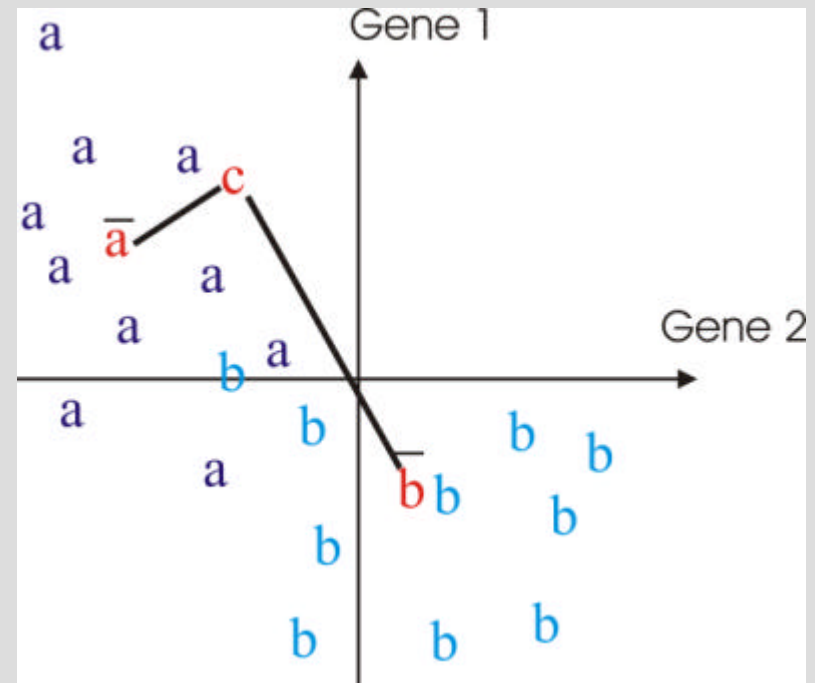
$c = (c_1, c_2)$ Patient without diagnosis

Compare: $d_a = (\bar{a}_1 - c_1)^2 + (\bar{a}_2 - c_2)^2$ and

$$d_b = (\bar{b}_1 - c_1)^2 + (\bar{b}_2 - c_2)^2$$

Diagnosis: a if $d_a < d_b$

b else



Many (N) genes:

$a_{i,j}$ Gene i in Patient j from group a

$b_{i,j}$ Gene i in Patient j from group b

$$\bar{a} = (\bar{a}_1, \dots, \bar{a}_N)$$

$$\bar{b} = (\bar{b}_1, \dots, \bar{b}_N)$$

c_1, \dots, c_N Patient without diagnosis

Compare distances to the centroids :

$$d_a = \sum_{i=1}^N (\bar{a}_i - c_i)^2$$

$$d_b = \sum_{i=1}^N (\bar{b}_i - c_i)^2$$

Diagnosis : a if $d_a < d_b$

b else

Nearest Centroid Method

(Plain Vanilla)

**Patient groups are
modelled separately by
centroids**

**Diagnosis is according
to the nearest centroid
in euclidean distance**

$a_{i,j}$ gene i in patient j from group a

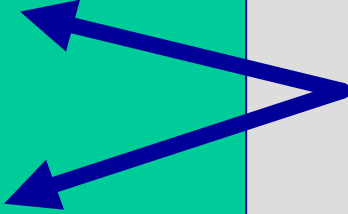
$b_{i,j}$ gene i in patient j from group b

$$d_a = \sum_{i=1}^N (\bar{a}_i - c_i)^2$$

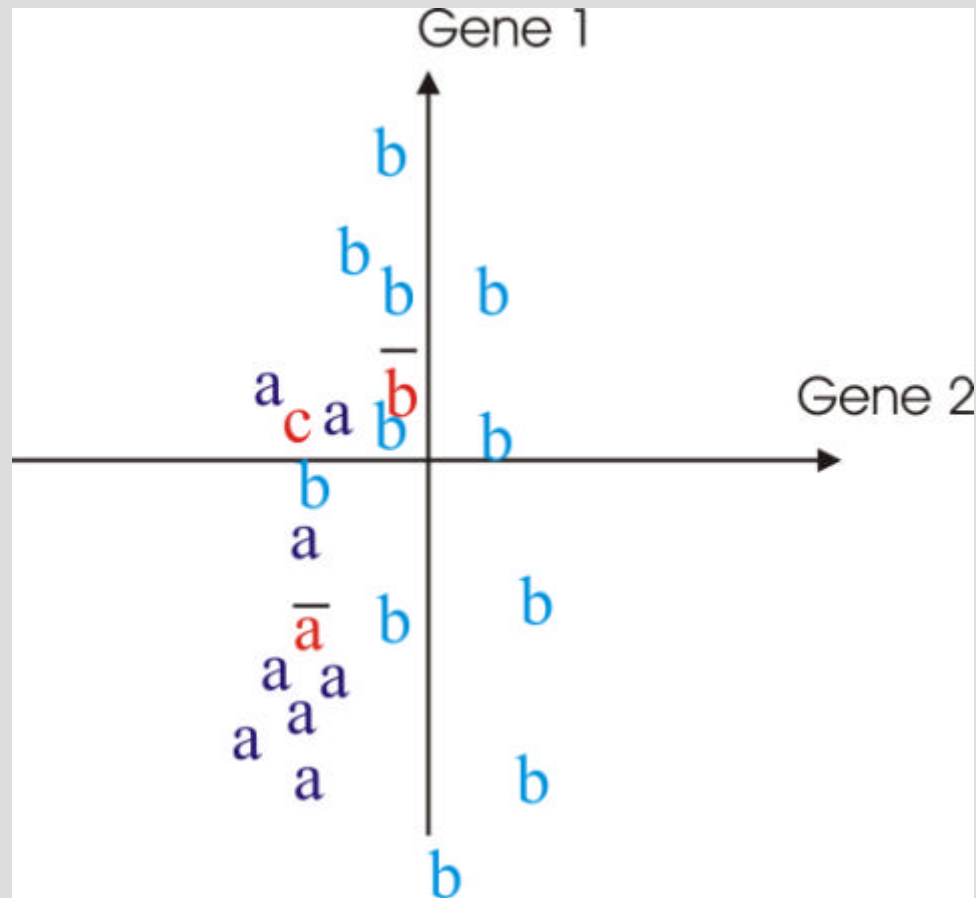
$$d_b = \sum_{i=1}^N (\bar{b}_i - c_i)^2$$

Diagnosis : a if $d_a < d_b$
b else

**All N genes
contribute equally
to the diagnosis ...**

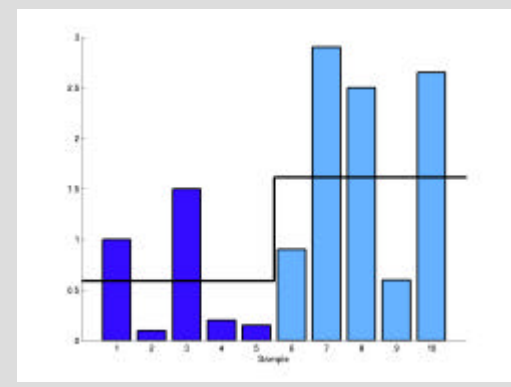
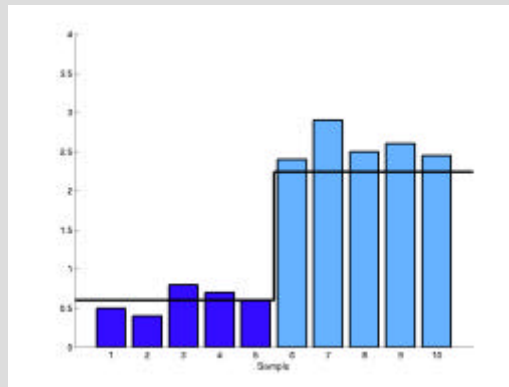
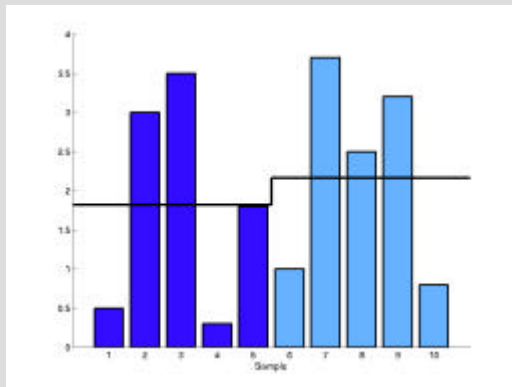


... that is a problem



Genes with a small „variance“ should get more weight than genes with high variance

$$d_a = \sum_{i=1}^N w_i (\bar{a}_i - c_i)^2 \quad d_b = \sum_{i=1}^N w_i (\bar{b}_i - c_i)^2$$



Use the pooled within class variance ... instead of the overall variance

The variances need to be estimated

$$s_i^2 = \frac{1}{n-2} \sum_{j=1}^{n/2} (a_{i,j} - \bar{a}_i)^2 + (b_{i,j} - \bar{b}_i)^2$$

pooled in class variance

In our case :

$$n = 200$$

→ SAM

$$w_i = (s_i + s_0)^2$$

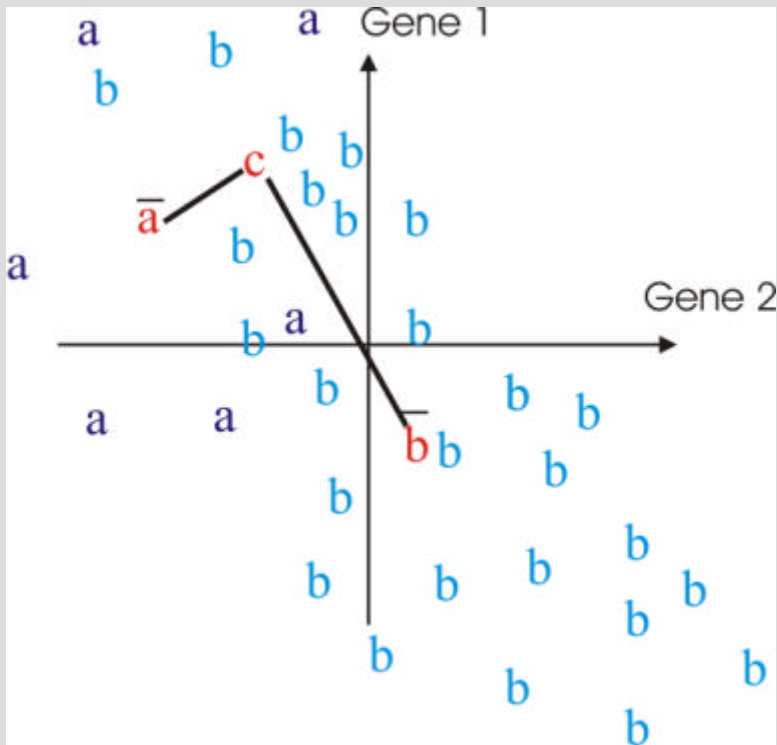
$$s_0^2 = \text{median}(s_1^2, \dots, s_N^2)$$

The estimated variance is not the true variance. It can be higher or lower. If a small variance is underestimated s_i^2

can be very small and w_i is unnaturally high.

While this is a rare event for a fixed gene it happens quite often if we are looking for 30000 genes

Is c an a or a b?



Is closer to the a centroid but there much more b than a samples

If this reflects the true population, than c should be classified as b

Baseline correction

p_a = relative size of group a
i.e. relative frequency of type a
samples in the study, or expert
knowledge

$$p_b = 1 - p_a$$

$$d_a(c) = \sum_{i=1}^N \frac{(\bar{a}_i - c_i)^2}{(\mathbf{s}_i + \mathbf{s}_0)^2} - 2 \log p_a$$

$$d_b(c) = \sum_{i=1}^N \frac{(\bar{b}_i - c_i)^2}{(\mathbf{s}_i + \mathbf{s}_0)^2} - 2 \log p_b$$

Discriminant Score

distance to the
centroid

$$d_a(c) = \sum_{i=1}^N \frac{(\bar{a}_i - c_i)^2}{(\mathbf{s}_i + \mathbf{s}_0)^2} - 2 \log \mathbf{p}_a$$

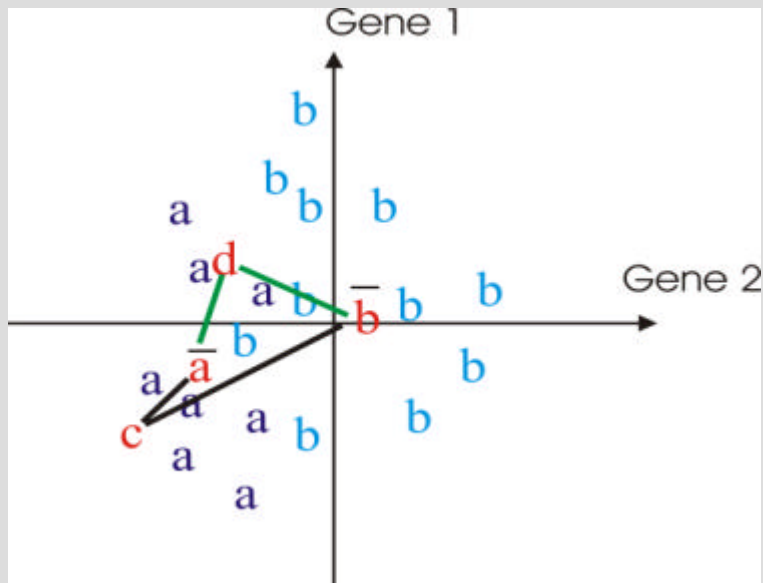
baseline
correction

$$d_b(c) = \sum_{i=1}^N \frac{(\bar{b}_i - c_i)^2}{(\mathbf{s}_i + \mathbf{s}_0)^2} - 2 \log \mathbf{p}_b$$

pooled within
class variance

variance
regularization
parameter

Classification probabilities



Both c and d are diagnosed as group a

But for d that was a close decision

$$\text{Prob} [Group(c) = a] = \frac{e^{-\frac{1}{2} d_a(c)}}{e^{-\frac{1}{2} d_a(c)} + e^{-\frac{1}{2} d_b(c)}}$$

$$\text{Prob} [Group(c) = b] = 1 - \text{Prob} [Group(c) = a]$$

Putting things into context

$d_a(c) = d_b(c)$ is a linear plane

We are still using all the 30000 genes

→ Overfitting problem

The plane is not necessarily optimal in terms of separation

This might be an advantage or a disadvantage

There is already some regularization going on

Variable selection

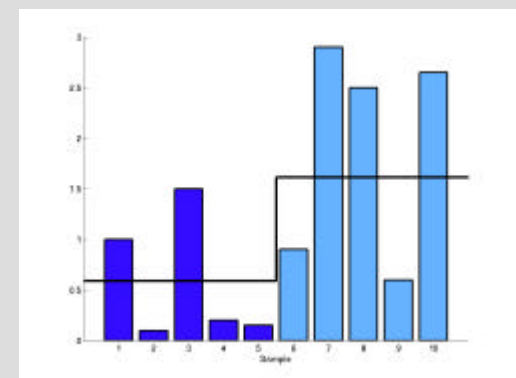
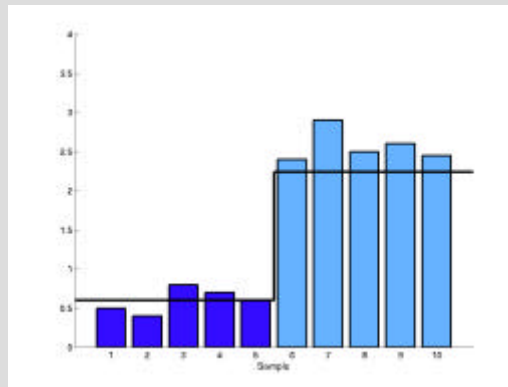
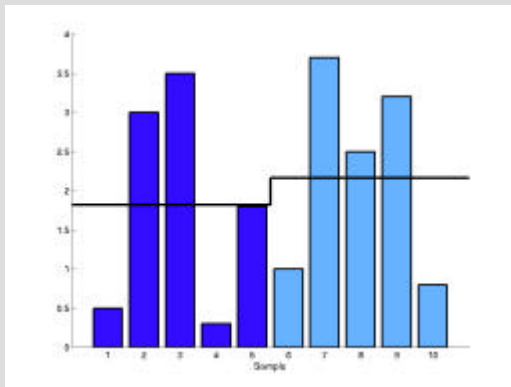
30000 genes are to many

They may cause overfitting

They introduce noise ... there weights are low ... but if there are many ...

They can not all matter

→ Choose genes:



Choose the genes with the highest weights
regularized t-score a la SAM

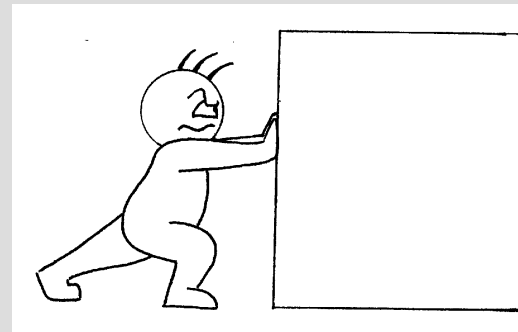
Hard thresholding vs. soft thresholding

Lets say we pick the top 100 genes

Gene Nr. 100 is in but gene Nr. 101 is not,

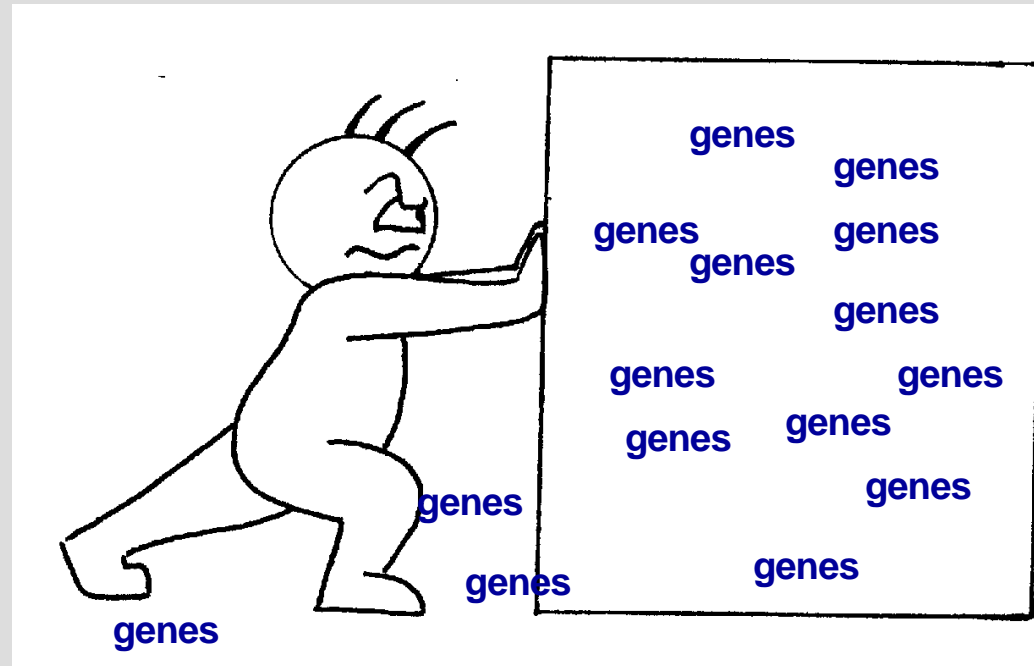
however, both genes are almost equally informative

If you want to get rid of genes you can chop them off or slowly push them out



The shrunken centroid method and the PAM program

Tibshirani et al 2002



Idea

Genes with high weights are influential for diagnosis

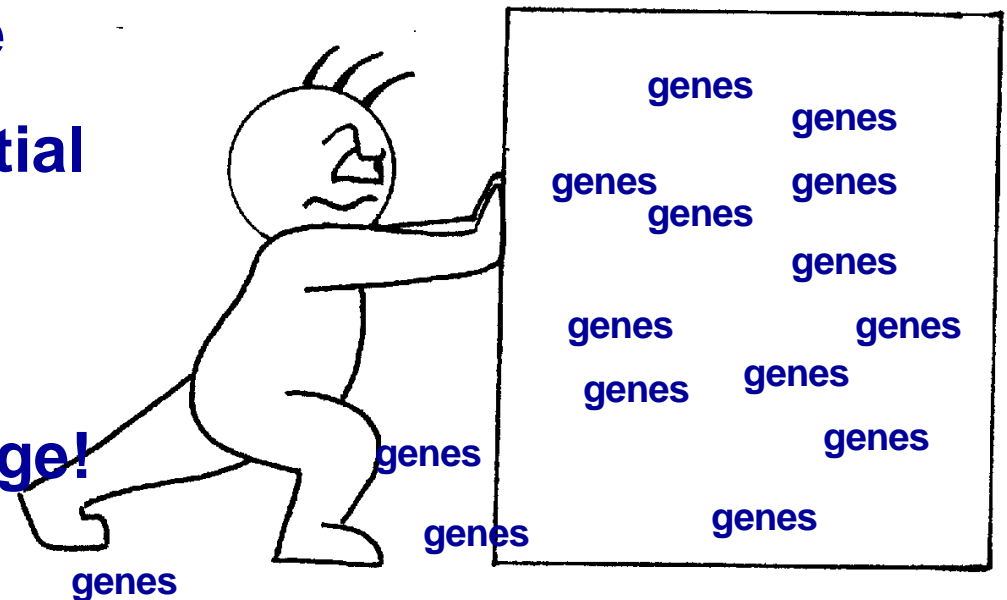
Genes with lower weights are less influential for diagnosis

Genes that are excluded can not be influential for diagnosis at all

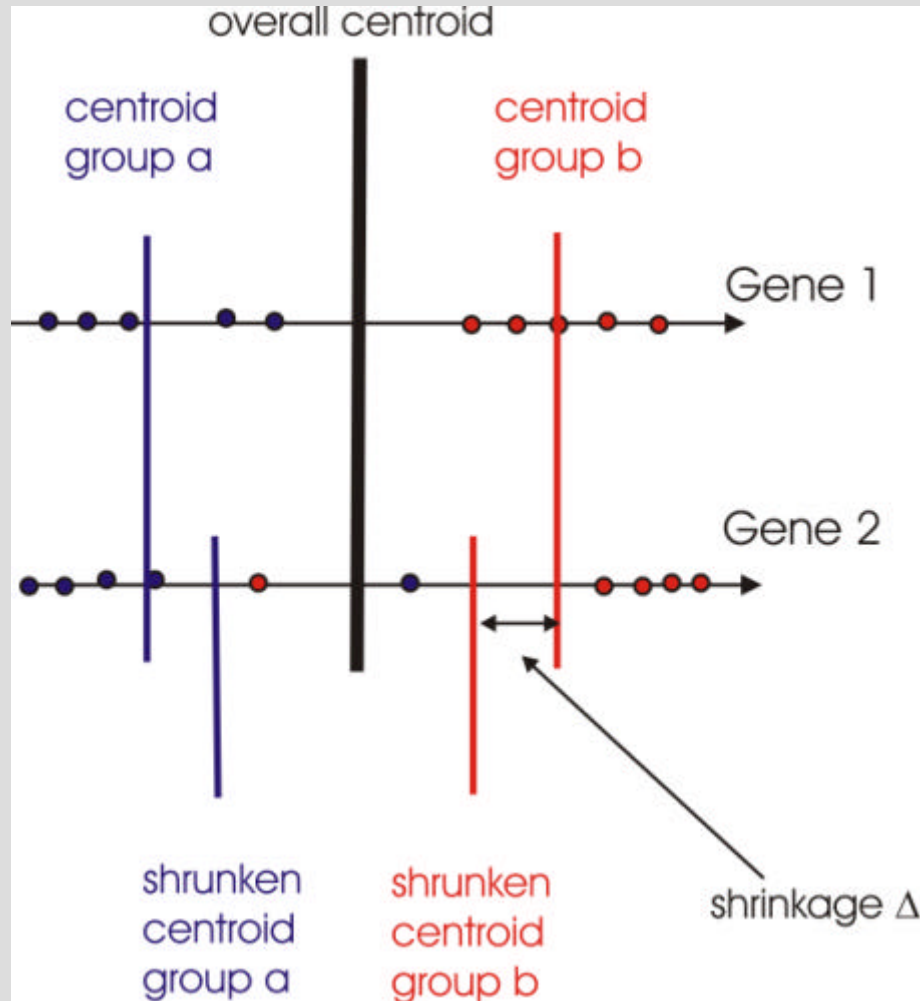
Before you exclude a gene totally from analysis make

it continuously less influential for the diagnosis

How? By centroid shrinkage!



Centroid shrinkage



Notation

\bar{a}_i mean of gene i in group a

\bar{b}_i mean of gene i in group b

\bar{x}_i mean of gene i using all data

Let

$$D_{i,a} = \frac{\bar{a}_i - \bar{x}_i}{m_a (\mathbf{s}_i + \mathbf{s}_0)}, \quad m_a = \sqrt{1/n_a + 1/n}$$

$$D_{i,b} = \dots$$

or

$$\bar{a}_i = \bar{x}_i + m_a (\mathbf{s}_i + \mathbf{s}_0) D_{i,a}$$

$$\bar{b}_i = \dots$$

group centroid

overall centroid

scaling factor

$$\bar{a}_i = \bar{x}_i + m_a (\mathbf{s}_i + \mathbf{s}_0) D_{i,a}$$

offset

$$\bar{a}_i = \bar{x}_i + m_a (\mathbf{s}_i + \mathbf{s}_0) D'_{i,a}$$

shrunk offset

$$D'_{i,a} = \text{sign}(D_{i,a}) (|D_{i,a}| - \Delta)_+$$

shrinkage parameter

$(\dots)_+ = \text{truncation at zero}$

Ok, the same in words for those who do not like formulae



Gene by gene, we shrink the group centroids towards the overall centroids standardized by the within-class standard deviations until the group centroids fall onto the coverall centroid ... then the gene is excluded.

When a group centroid moves towards the overall centroid the corresponding gene becomes continuously less influential for diagnosis until it is finally excluded

**The amount of shrinkage is controlled
by Delta**

**Little shrinkage many genes are still
contributing to the centroids**

**High shrinkage only few genes are
still in the analysis**

**The amount of shrinkage can be
determined by**

**cross validation ... we will discuss
this later**



Estrogen Receptor Status

- **7000** genes
- **49** breast tumors
- **25** ER+
- **24** ER-

ER-

ER+

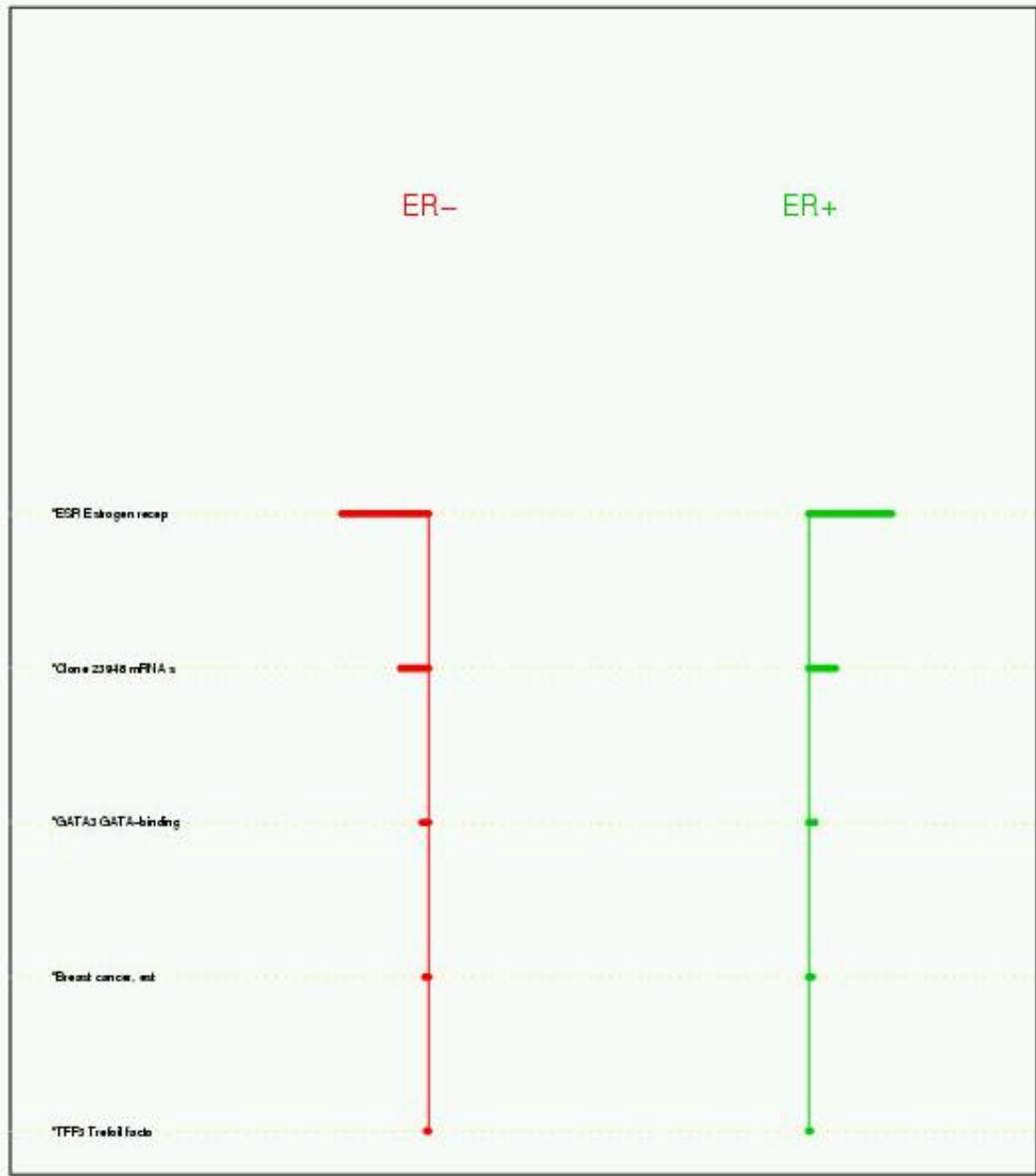
*ESR Estrogen recep

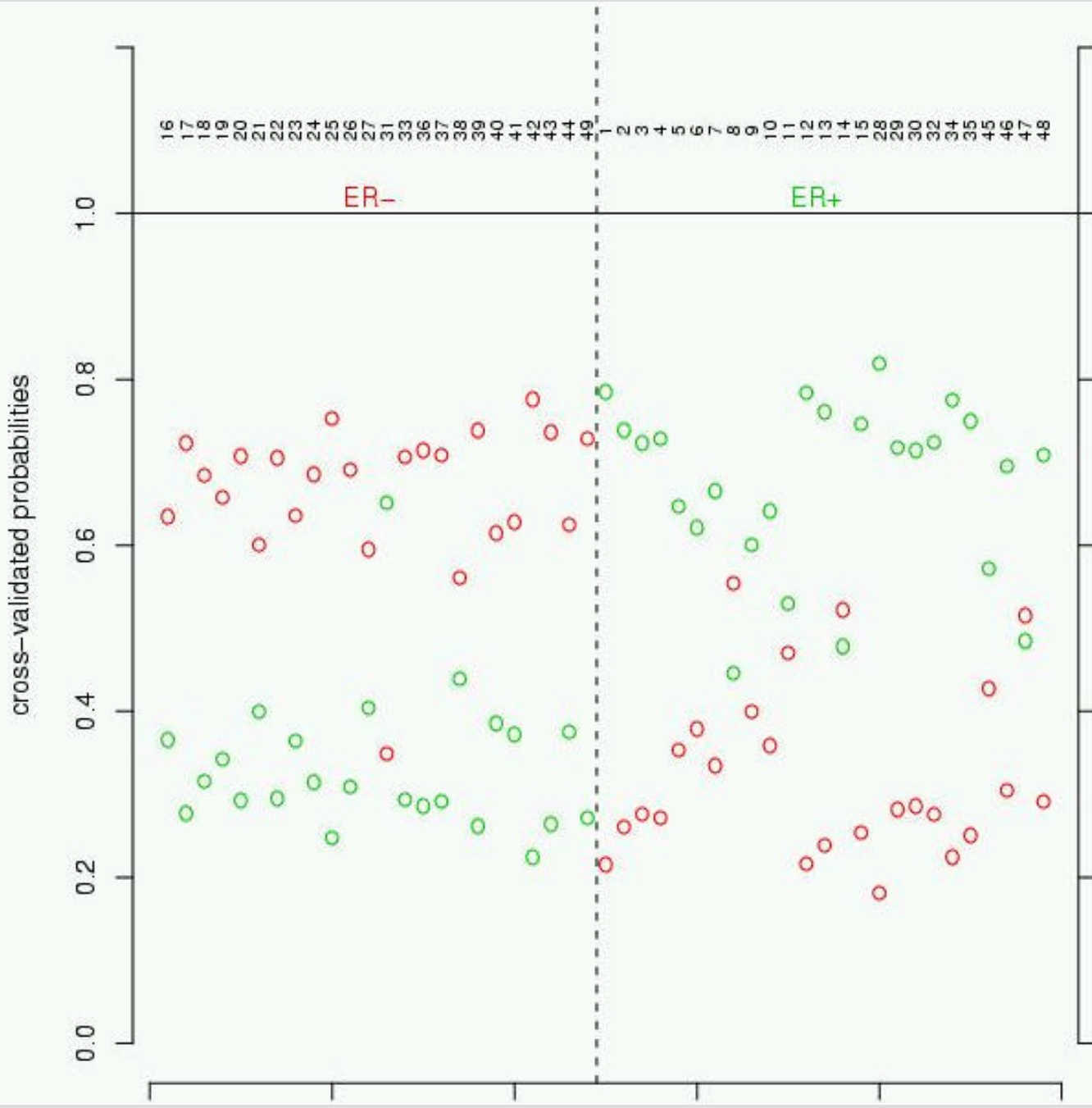
*Gene 23946 mRNA s

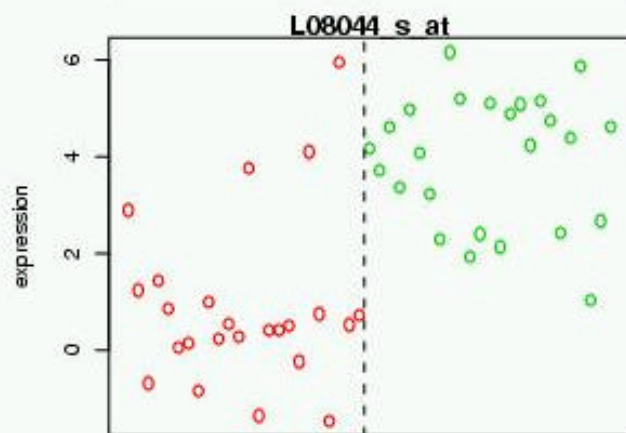
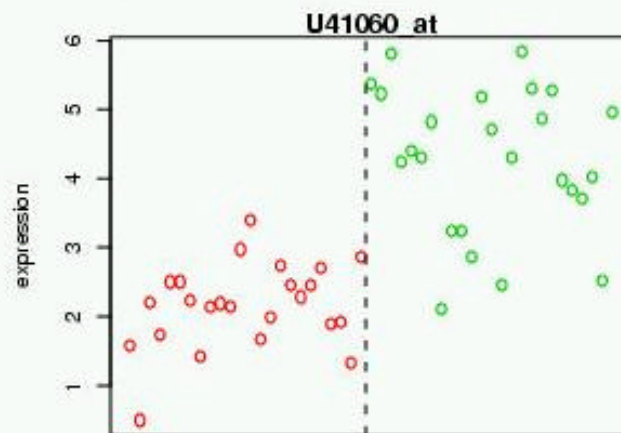
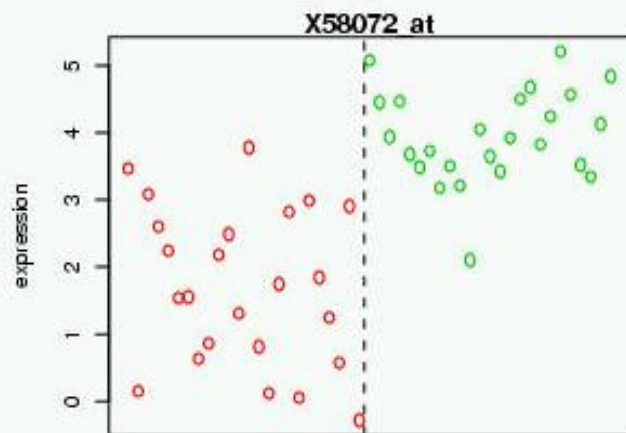
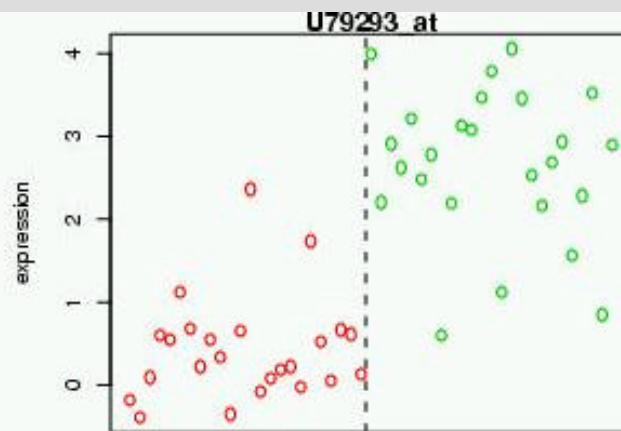
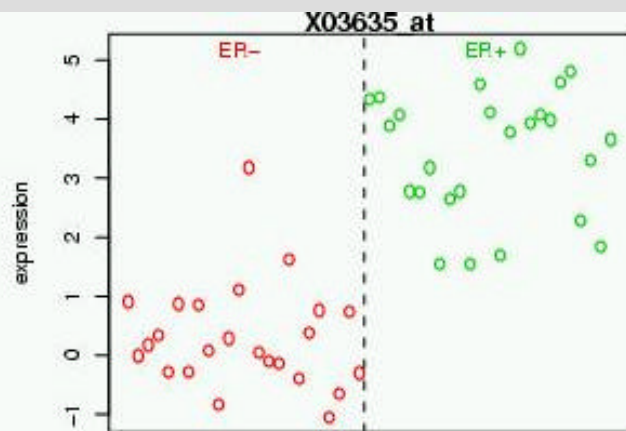
*GATA2 GATA-binding

*Breast cancer, int

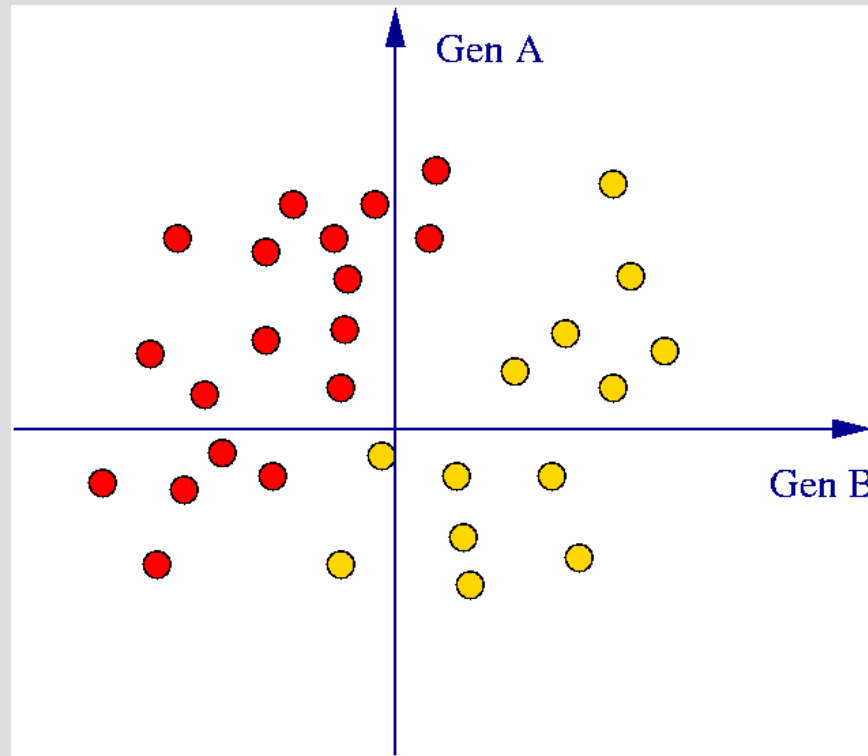
*TFPI1 Tissue factor





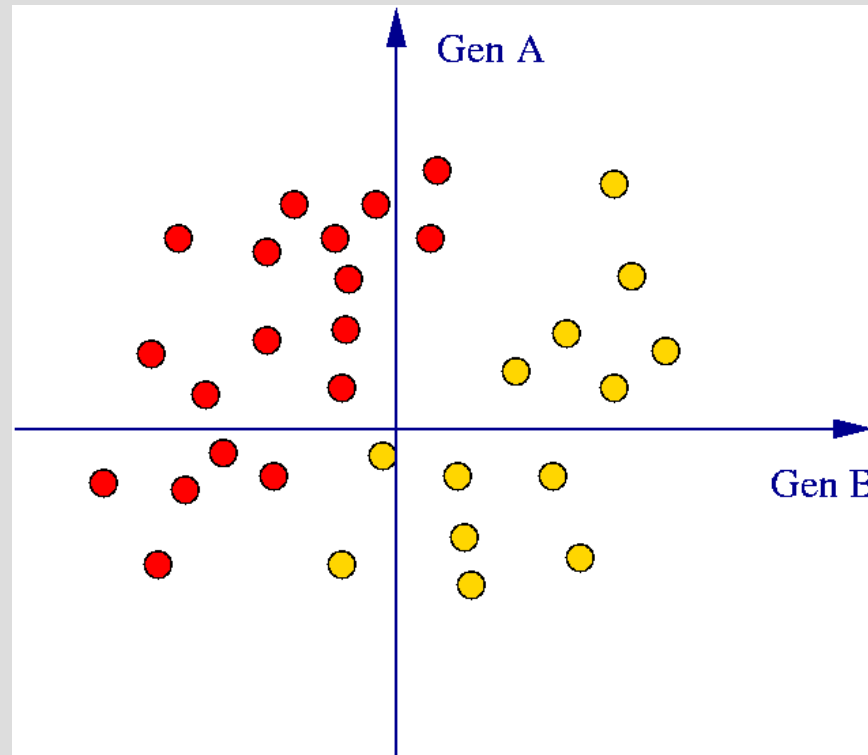


Imagine we have a study with 30000 genes 29998 of them with no biological significance and the 2 below



What would PAM do?

What would PAM do?



Fail

Pam would not find these two genes because their group centroids are too near to the overall centroid

Each of them is a poor classifier, together they are a good one

This is both a bug and a feature of PAM

Again, there is regularization going on

PAM does not find everything, but what it finds has a good chance to be of importance

- **PAM does variable selection by screening one gene after another**
- **The centroids are the signatures**
- **So when we decide whether a gene should go into a signature we only look at this single gene and decide**
- **Interaction of genes is unimportant for the selection**
- **We combine consistently up and down regulated genes into signatures**

Devices of regularization used by PAM

- Gene selection**

- Shrinkage**

- Gene selection by screening (no wrapping)**

- The weight of a gene only depends on the gene and not on its interaction with others**

- Use of a baseline depending on the population size of the groups ... more information in addition to the expression data**

Questions



Coffee



What did we learn so far, and what didn't we?

- The high dimensional data leads to overfitting problems
- There are meaningful signatures and those that mean nothing
- Regularization (PAM,SVM,...) helps finding meaningful signatures
- ...

-... but if I have found one there is still no guarantee

-The patients in my data display differences in a signature between group a and b ... but does this apply to a new patient too?

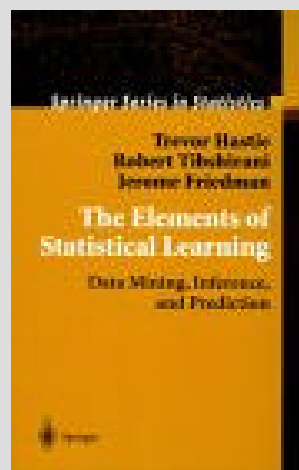
- Is the signature predictive? Can it be used for diagnosis?

Problems:

1. How much regularization is good?

2. If I have found a signature, how do I know whether it is meaningful and predictive or not?

Model Selection & Model Assessment



Chapter 7

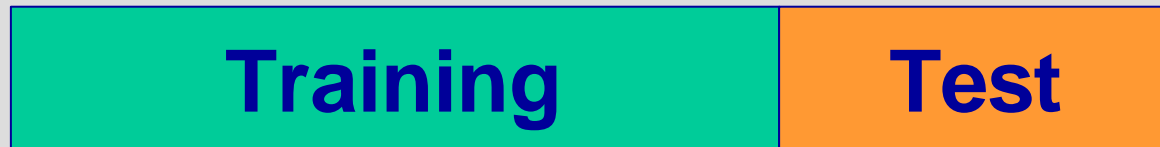
Cross-Validation and Bootstrap

We only discuss Cross-Validation

Test and Training Data

150

50



Split your profiles randomly into a **training set** and a **test set**

Train your model only using the data in the training set

(define centroids, calculate normal vectors for large margin separators, ...)

Apply the model to the test data ...

The setup

$x = (x_1, \dots, x_n)$ profile

$y \in \{a, b\}$ class assignment

$g(x) = y$ true class of x

$\hat{g}(x) = \hat{y}$ predicted class of x

$\hat{p}_y(x)$ estimated probability that x is of class y

(PAM, logistic regression, logistic discrimination etc.)

For example :

$f(x)$ signature, c cutoff

$\hat{g}(x) = a$ if $f(x) > c$

$\hat{g}(x) = b$ if $f(x) \leq c$

Trainings and Test Data

Trainings data :

$$x_j^{train} = (x_{1,j}^{train}, \dots, x_{n,j}^{train})$$

a trainings profile

$$y_j^{train}$$

its true class (used when fitting the model)

$$\hat{y}_j^{train}$$

its predicted class

$$\hat{p}_y(x_j^{train})$$

estimated probability that x_j^{train} is of class y

Test data :

$$x_j^{test} = (x_{1,j}^{test}, \dots, x_{n,j}^{test})$$

a test profile

$$y_j^{test}$$

its true class (NOT used when fitting the model)

$$\hat{y}_j^{test}$$

its predicted class

$$\hat{p}_y(x_j^{test})$$

estimated probability that x_j^{test} is of class y

Errors & Deviances

Notation : Indicator function

$$I(y = \hat{y}) = 1 \text{ if } y = \hat{y}$$

$$I(y = \hat{y}) = 0 \text{ if } y \neq \hat{y}$$

Trainings Error :

$$err^{train} = \sum_{\text{trainingsample}} I(y_j^{train} \neq \hat{y}_j^{train}) \quad \text{Number of misclassifications in the trainings set}$$

Trainings Deviance :

$$dev^{train} = -\frac{2}{N^{train}} \sum_{\text{trainingsample}} I(y = a) \log \hat{p}_a(x_j^{train}) + I(y = b) \log \hat{p}_b(x_j^{train})$$

Test Error :

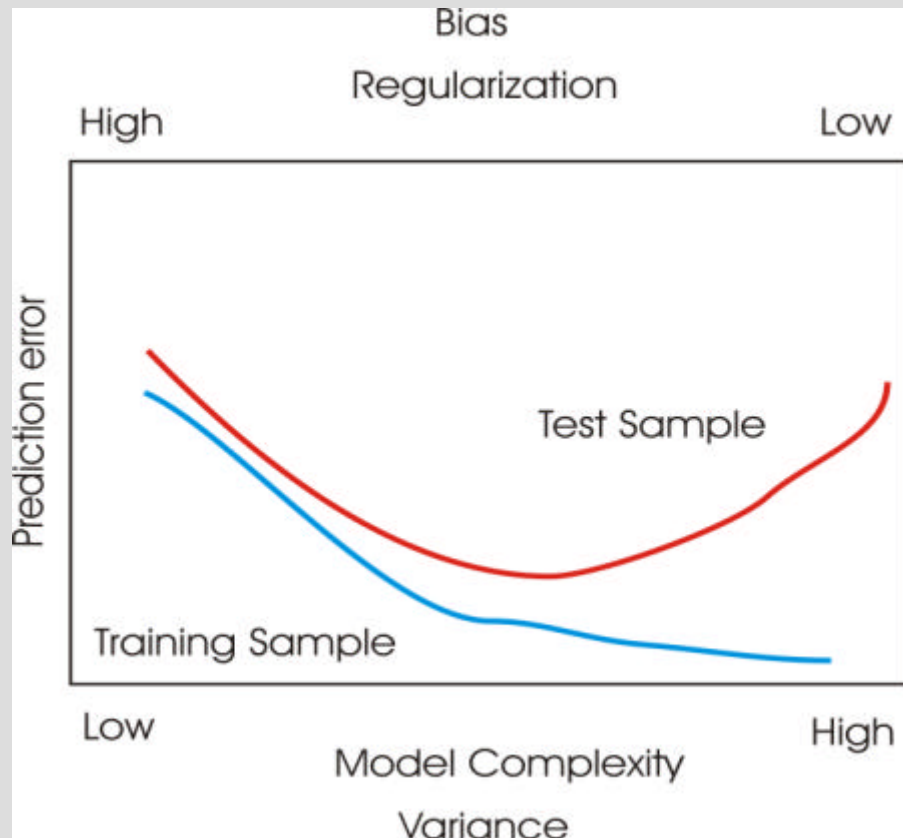
$$err^{test} = \dots$$

Test Deviance :

$$dev^{test} = \dots$$

The deviance is a continuous probabilistic error measure

The bias variance trade off



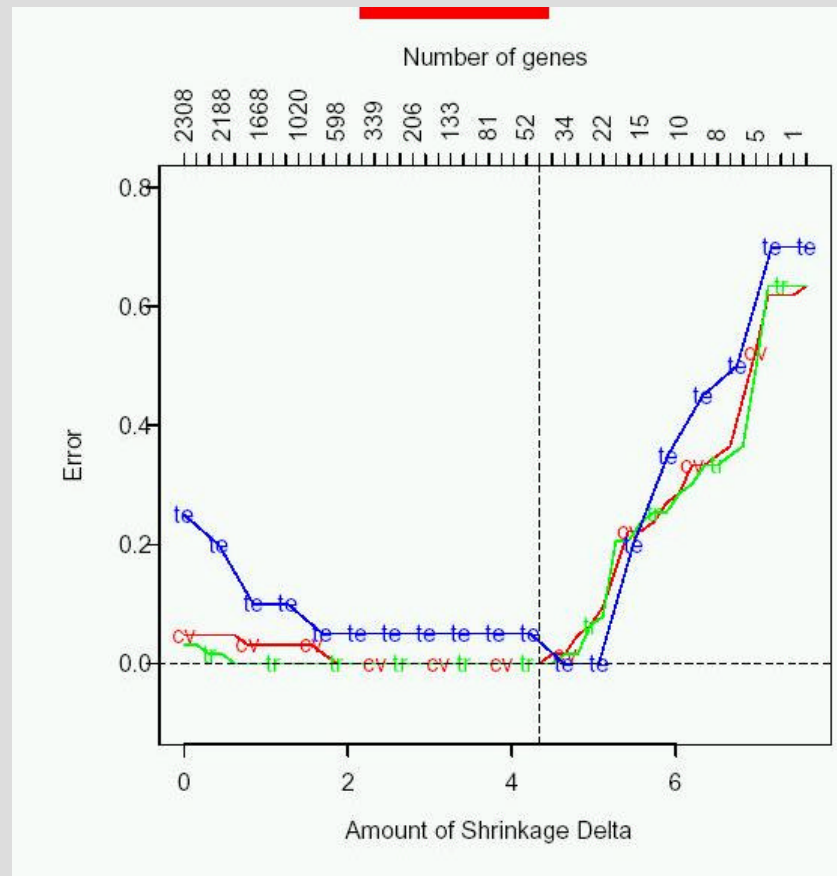
Model Complexity:

- max number of genes
- shrinkage parameter
- minimal margin
- etc

Small round blue cell tumors

4 classes

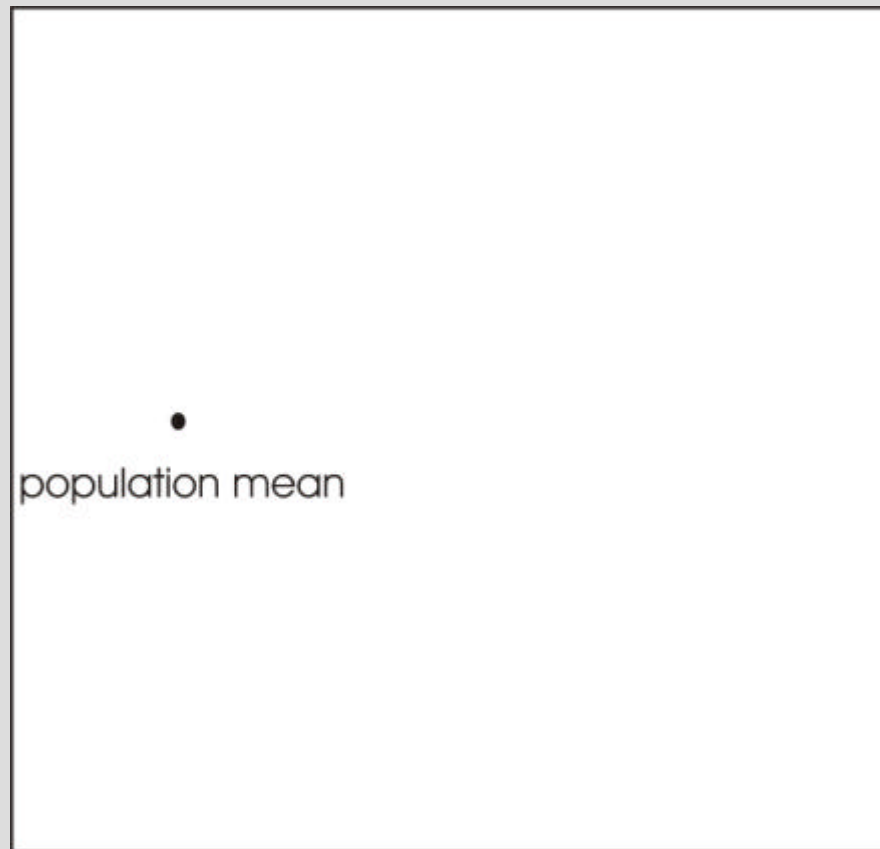
(Data: Khan et al. 2001)
(Analysis (PAM): Hastie et al 2002)



How come?

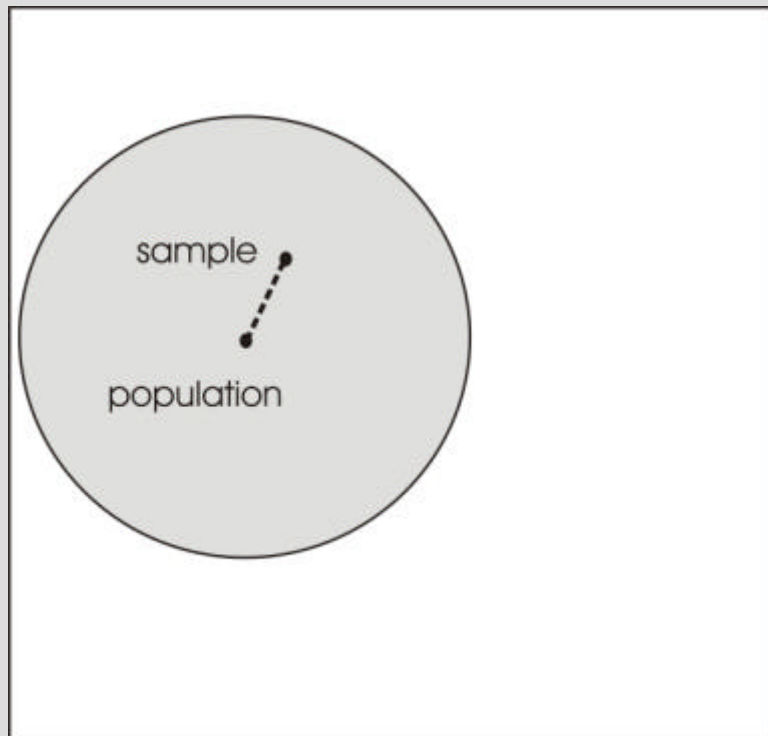
Population mean:

Genes have a certain mean expression and correlation in the population

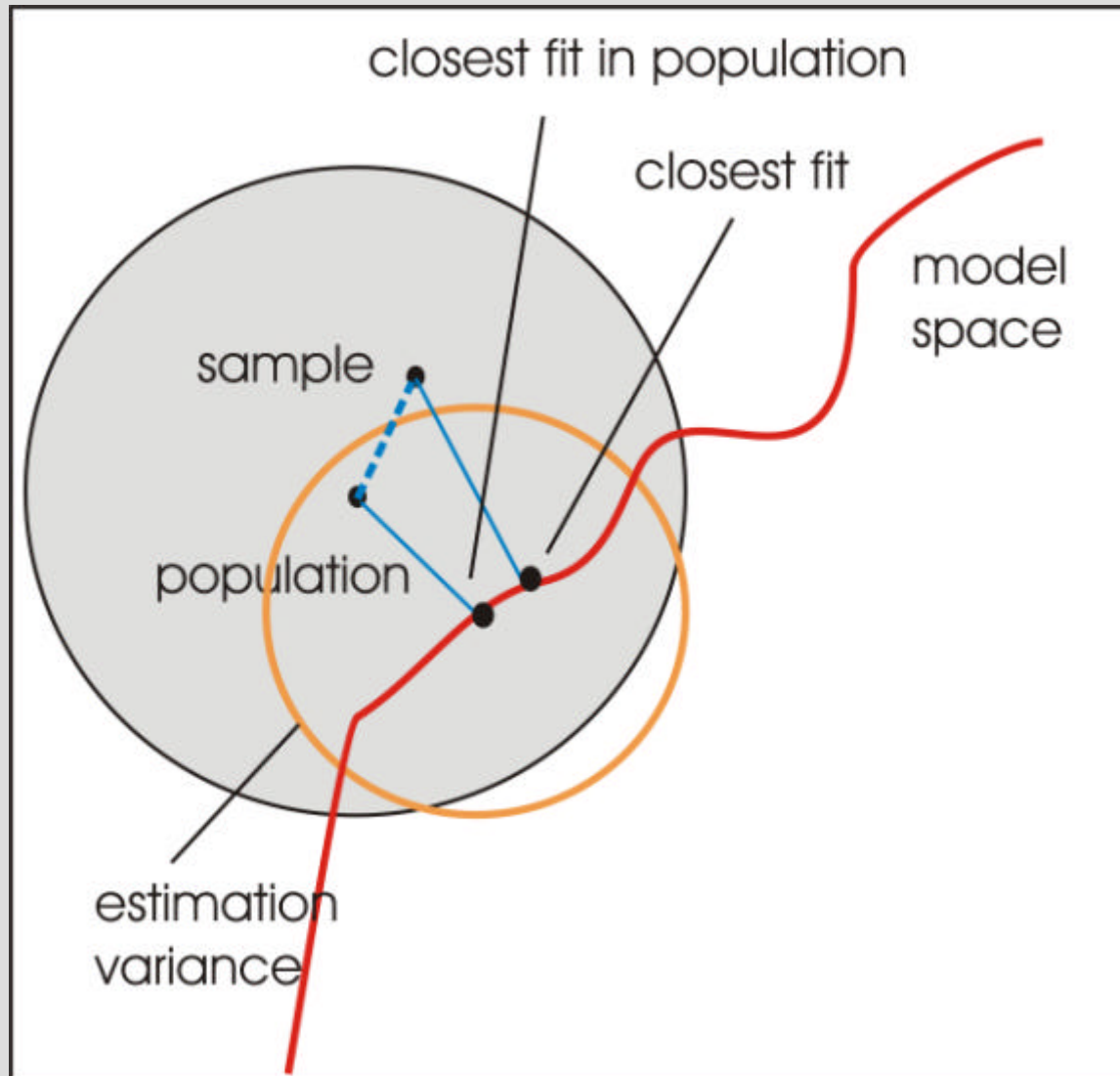


Sample mean:

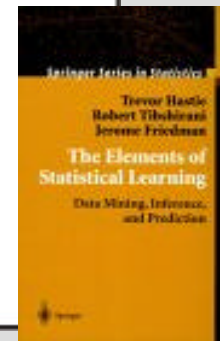
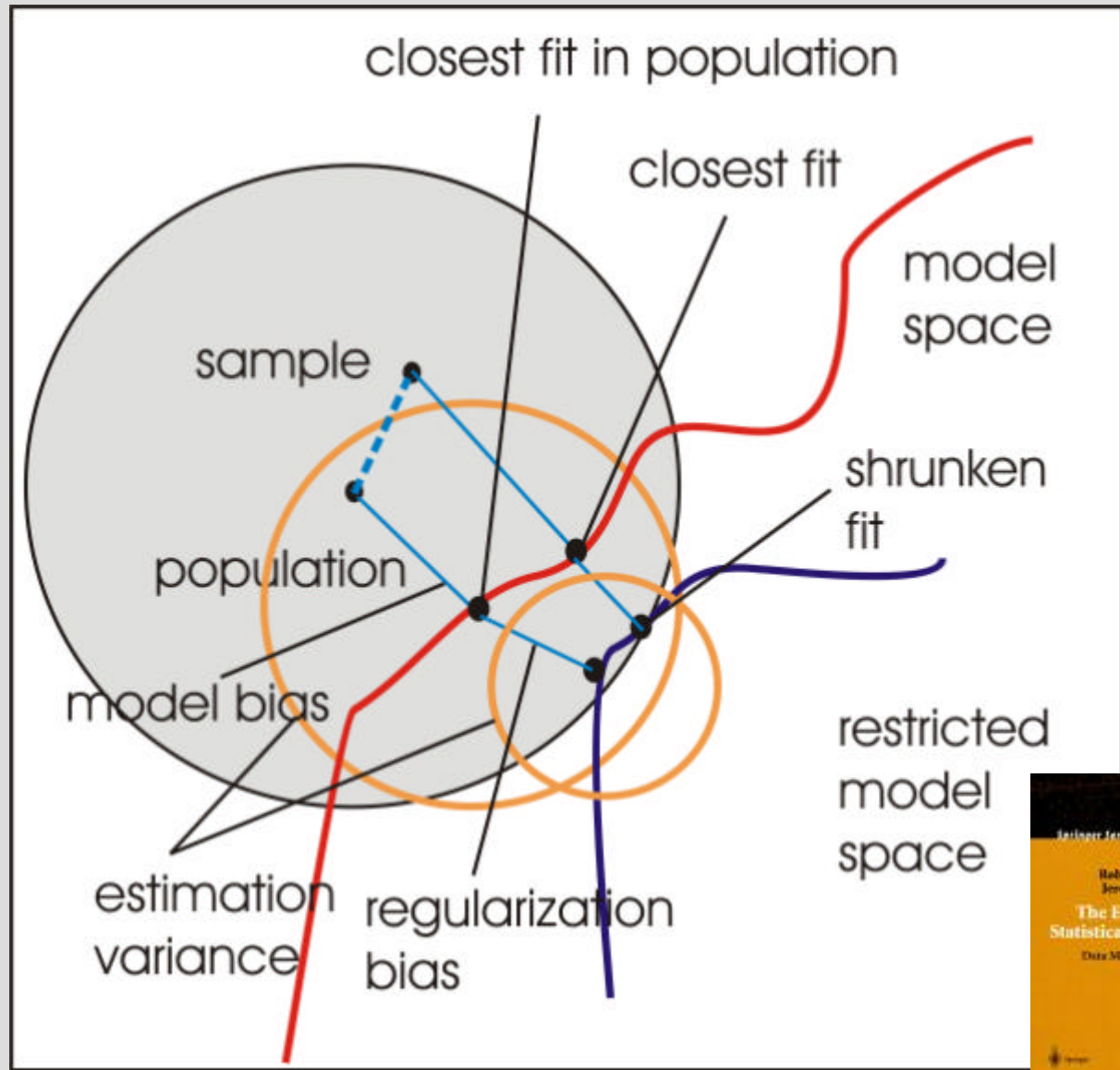
We observe average expression and empirical correlation



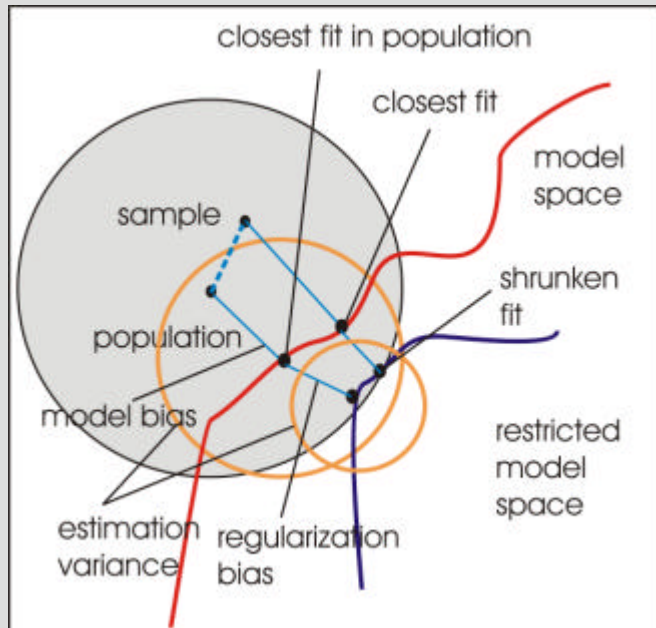
Fitted model:



Regularization



Bias-Variance-Trade-Off in PAM and in general



A lot of shrinkage:

Poor fit & low variance

Little shrinkage

Good fit & high variance

How much shrinkage should I use?

Model Selection with separate data

100

50

50



Split of some samples for Model Selection

Train the model on the training data with different choices for the regularization parameter

Apply it to the selection data and optimize this parameter (Model Selection)

Test how good you are doing on the test data (Model Assessment)

10 Fold Cross-Validation

...



...

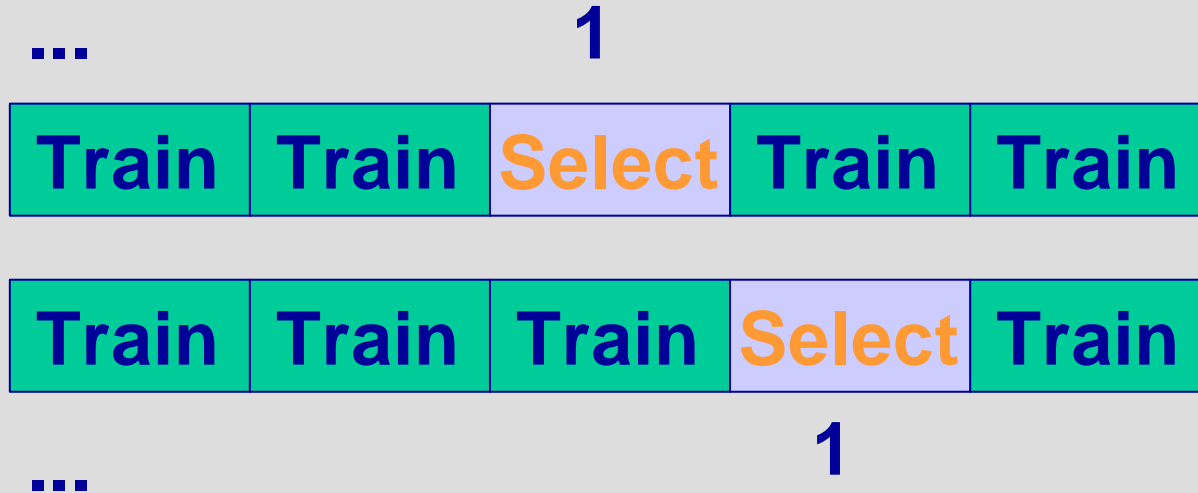
Chop up the training data (**don't touch the test data**) into 10 sets

Train on 9 of them and predict the other

Iterate, leave every set out once

Select a model according to the prediction error (deviance)

Leave one out Cross-Validation



Essentially the same

But you only leave one sample out at a time and predict it using the others

Good for small training sets

Model Assessment

How well did I do?

Can I use my signature for clinical diagnosis?

How well will it perform?

How does it compare to traditional methods?

The most important thing:

Don't fool yourself! (... and others)



This guy (and others) thought for some time he could predict the nodal status of a breast tumor from a profile taken from the primary tumor!

**... there are significant differences.
But not good enough for prediction**

(West et al PNAS 2001)

DOs AND DONTs :

1. Decide on your diagnosis model (PAM,SVM,etc...) and **don't change your mind later on**
2. Split your profiles randomly into a **training set and a test set**
3. Put the data in the test set away.
4. Train your model only using the data in the training set
(**select genes**, define centroids, calculate normal vectors for large margin separators,**perform model selection ...**)
don't even think of touching the test data at this time
5. Apply the model to the test data ...
don't even think of changing the model at this time
6. Do steps 1-5 only once and accept the result ...
don't even think of optimizing this procedure

The selection bias

- You can not select 20 genes using all your data and then with this 20 genes split test and training data and evaluate your method.
- There is a difference between a model that restricts signatures to depend on only 20 genes and a data set that only contains 20 genes
- Your model assessment will look much better than it should

FAQ

**How many
patients do
we need?**

**Do we need
to replicate
patient
profiles?**

**Do we need to
consult a
bioinformatics
expert?**

**When on do
we need to
contact
him/her?**

**Where do we
find him/her?**

Thank
You