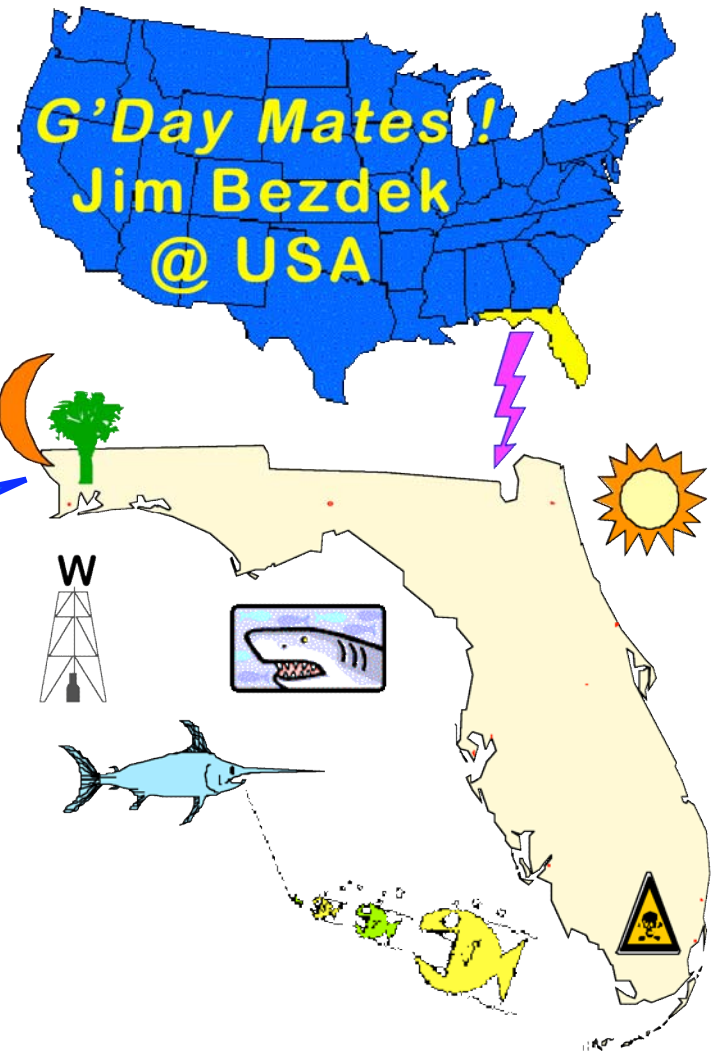


I'm from ...



Pensacola

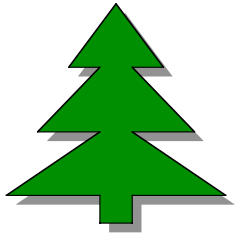


Pensacola
before
Hurricane IVAN

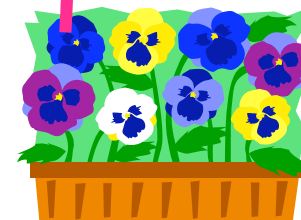
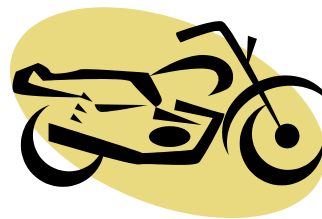
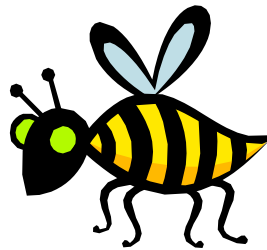




What **IS** Pattern Recognition ?



Where will this bee go?



Pattern Recognition



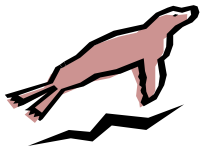
Numerical Pattern Recognition



Feature Analysis



Similarity : Distances and Norms



Label Vectors and Partitions



Cluster Analysis



Classifier Design

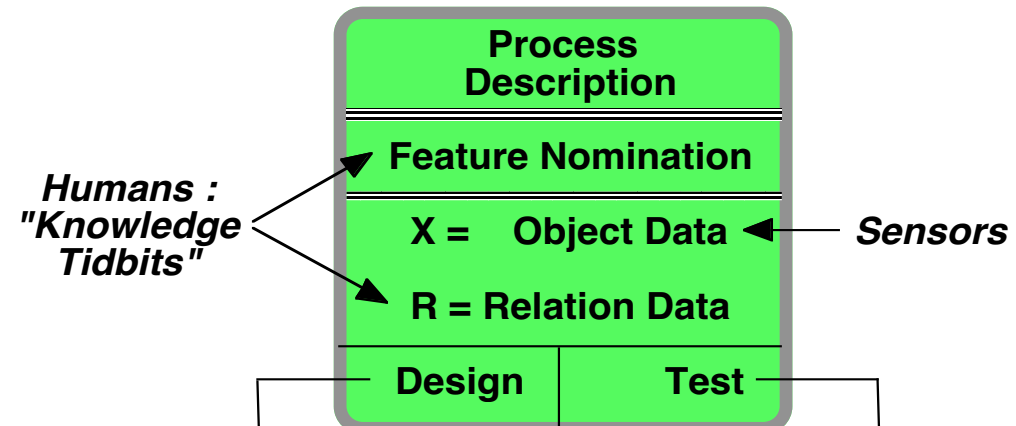
Numerical Pattern Recognition

Process Description

Feature Analysis

Cluster Analysis

Classifier Design



2 kinds of data for Pattern Recognition

Objects $O = \{o_1, \dots, o_n\} : o_i = i\text{-th } \textit{physical} \text{ object}$

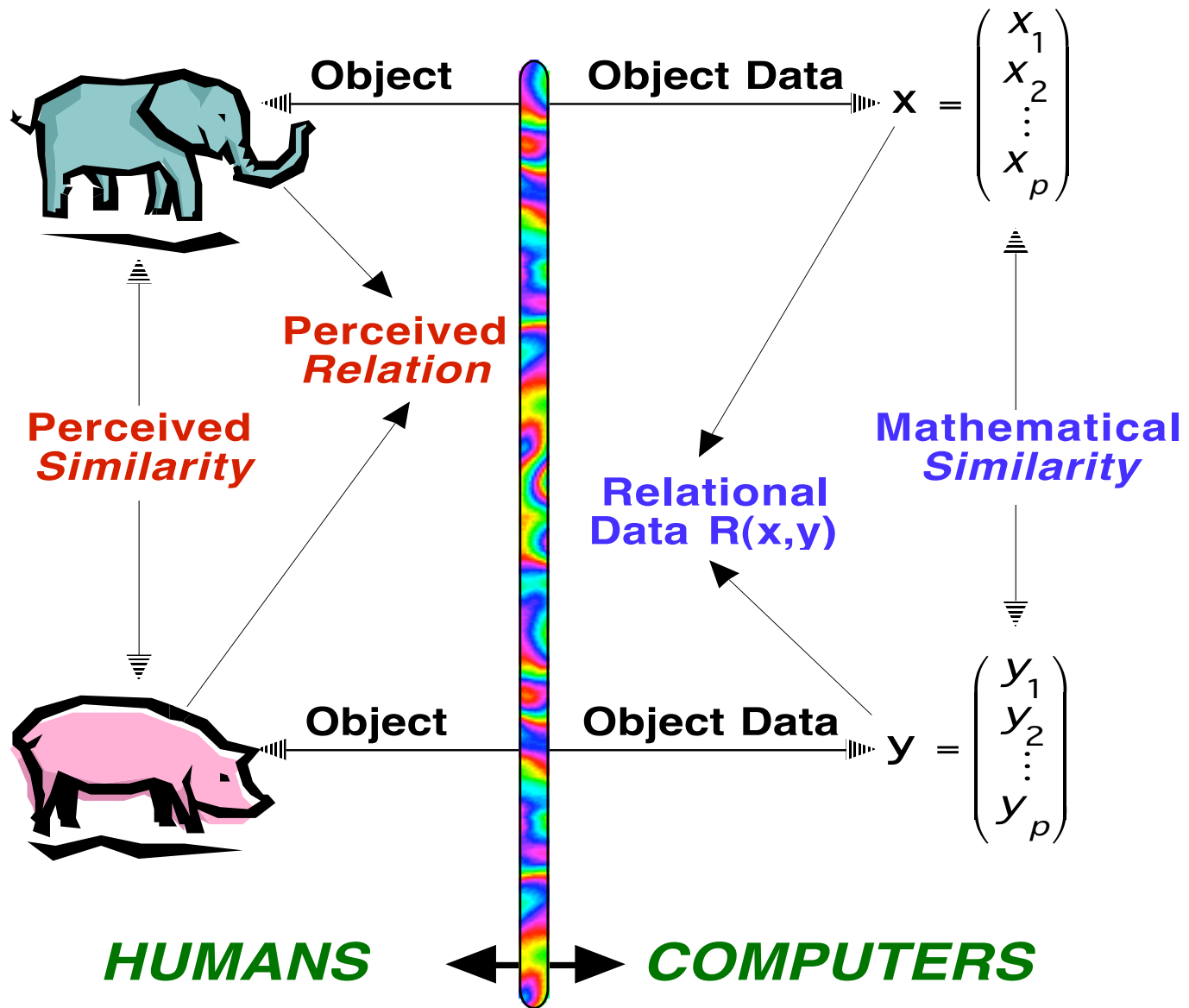
Object Data $X = \{x_1, \dots, x_n\} : x_i = \textit{feature vector} \text{ for } o_i$
 $x_{ji} = j\text{-th } (\textit{measured}) \text{ feature of } x_i : 1 \leq j \leq p$

Relational data $R = [r_{ij}] = \textit{relationship}(o_i, o_j) \text{ or } (x_i, x_j)$
 $s_{ij} = \textit{similarity}(o_i, o_j) \text{ or } (x_i, x_j)$
 $d_{ij} = \textit{dissimilarity}(o_i, o_j) \text{ or } (x_i, x_j)$

Typically
($R = D$)

(i)	$d_{ij} \geq 0$	$1 \leq i \neq j \leq n$
(ii)	$d_{ii} = 0$	$1 \leq i \leq n$
(iii)	$d_{ij} = d_{ji}$	$1 \leq i \neq j \leq n$

When $X \rightarrow D$, $d_{ij} = ||x_i - x_j||_E$, D is “**Euclidean**”



Measuring (Dis)Similarity : The 5 **Good** Norms

Mahalonobis

< , > Norms

Diagonal

$$\delta_{M^{-1}} = \|\mathbf{x} - \mathbf{v}\|_{M^{-1}} = \sqrt{(\mathbf{x} - \mathbf{v})^T \mathbf{M}^{-1} (\mathbf{x} - \mathbf{v})}$$

$$\delta_{D^{-1}} = \|\mathbf{x} - \mathbf{v}\|_{D^{-1}} = \sqrt{(\mathbf{x} - \mathbf{v})^T \mathbf{D}^{-1} (\mathbf{x} - \mathbf{v})}$$

Euclidean

$$\delta_2 = \|\mathbf{x} - \mathbf{v}\|_{I^{-1}} = \sqrt{(\mathbf{x} - \mathbf{v})^T \mathbf{I}^{-1} (\mathbf{x} - \mathbf{v})}$$

City Block

Minkowski Norms

Sup or Max

$$\delta_1 = \|\mathbf{x} - \mathbf{v}\|_1 = \sum_{j=1}^p |\mathbf{x}_j - \mathbf{v}_j|$$

$$\delta_\infty = \|\mathbf{x} - \mathbf{v}\|_\infty = \max_{j=1}^p \{|\mathbf{x}_j - \mathbf{v}_j|\}$$

Parameters of the Statistical Norms

Mean Vector

$$\bar{\mathbf{v}} = \sum_{k=1}^n \mathbf{x}_k / n$$

Covariance
Matrix

$$\mathbf{M} = [\mathbf{m}_{ij}] = \frac{\sum_{k=1}^n (\mathbf{x}_k - \bar{\mathbf{v}})(\mathbf{x}_k - \bar{\mathbf{v}})^T}{n}$$

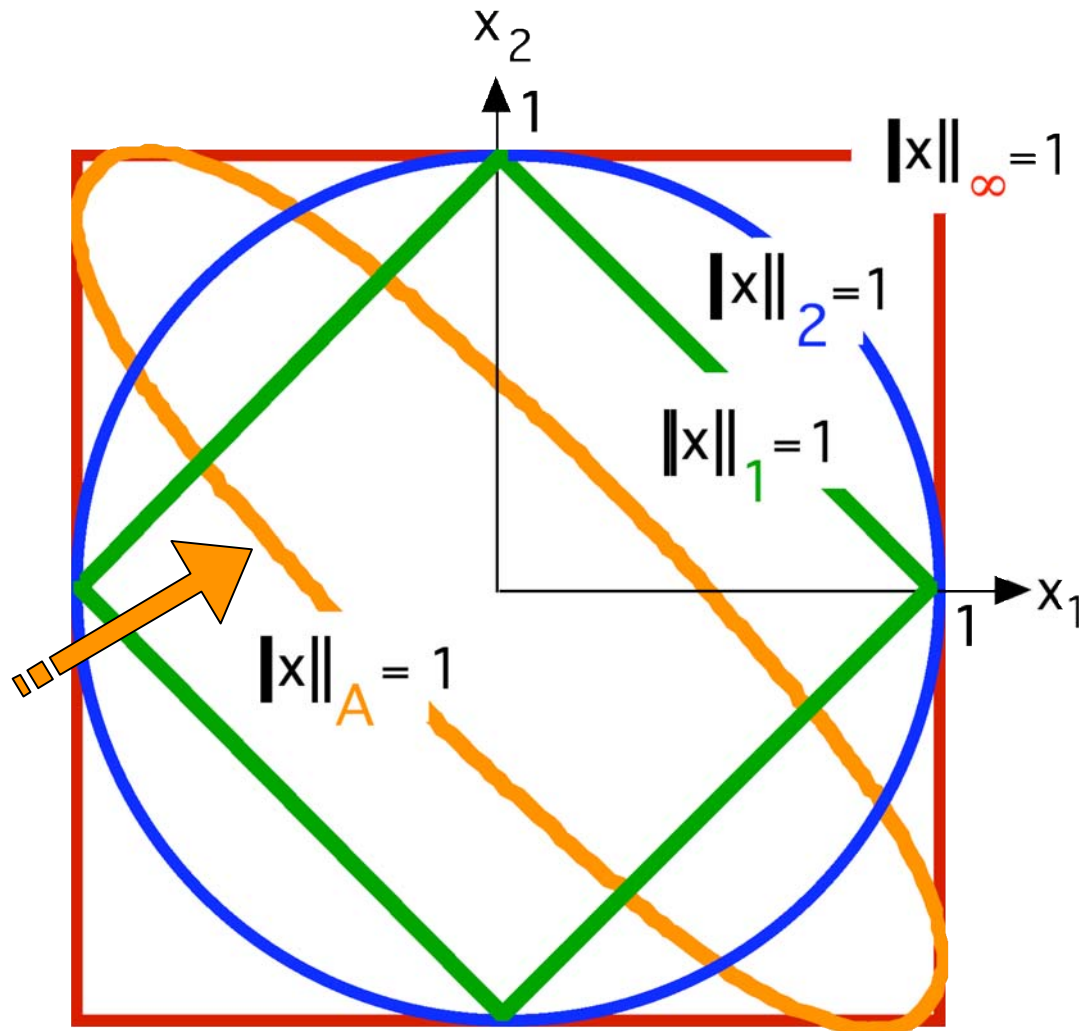
Variance
(Matrix)

$$\mathbf{D} = \begin{bmatrix} m_{11} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & m_{pp} \end{bmatrix} = \begin{bmatrix} s_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & s_p^2 \end{bmatrix}$$

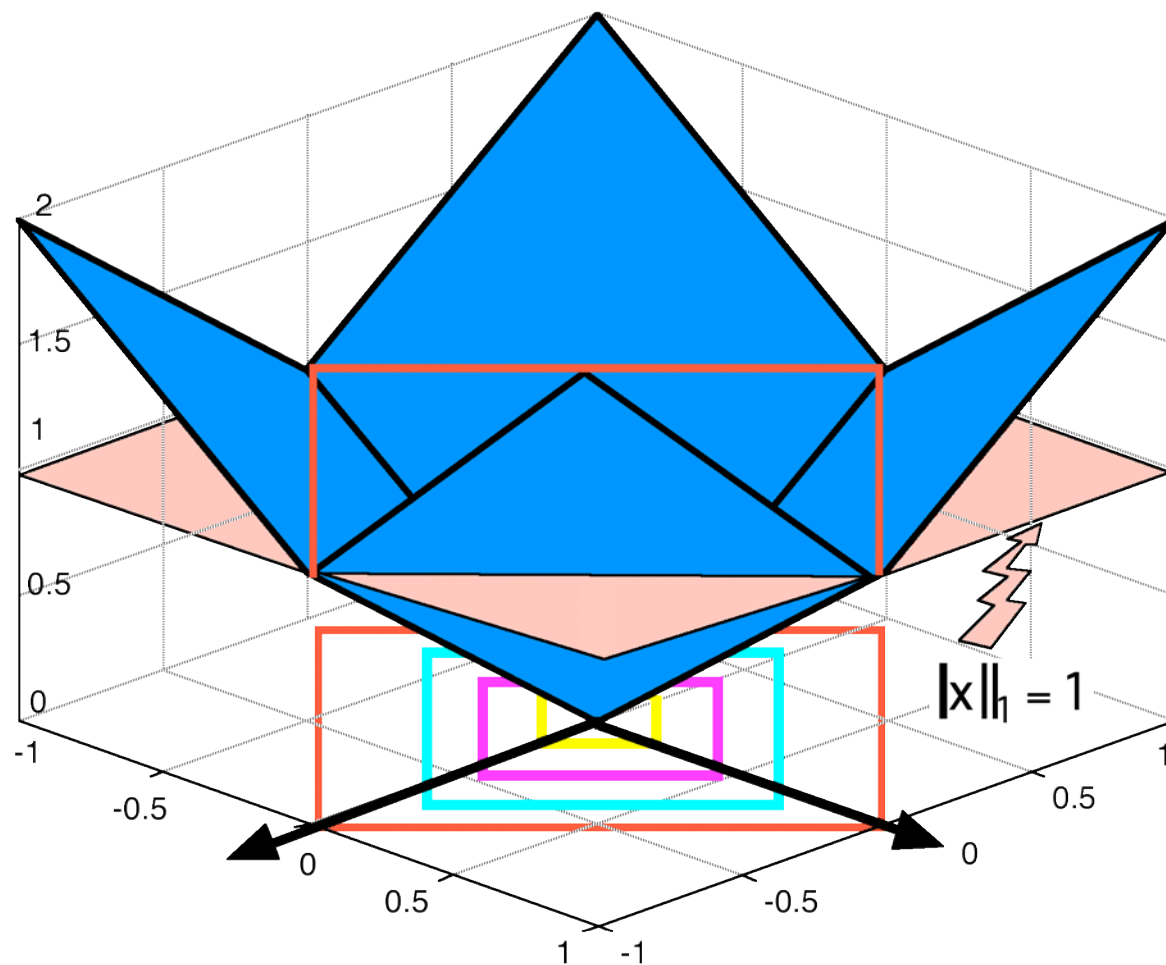
Level Sets and Open balls of Norms in \mathbb{R}^2

Inner Product Norms induce elliptical balls via eigenstructure of A

$$\langle \mathbf{x}, \mathbf{x} \rangle_A = \|\mathbf{x}\|_A^2 = \mathbf{x}^T \mathbf{A} \mathbf{x}$$

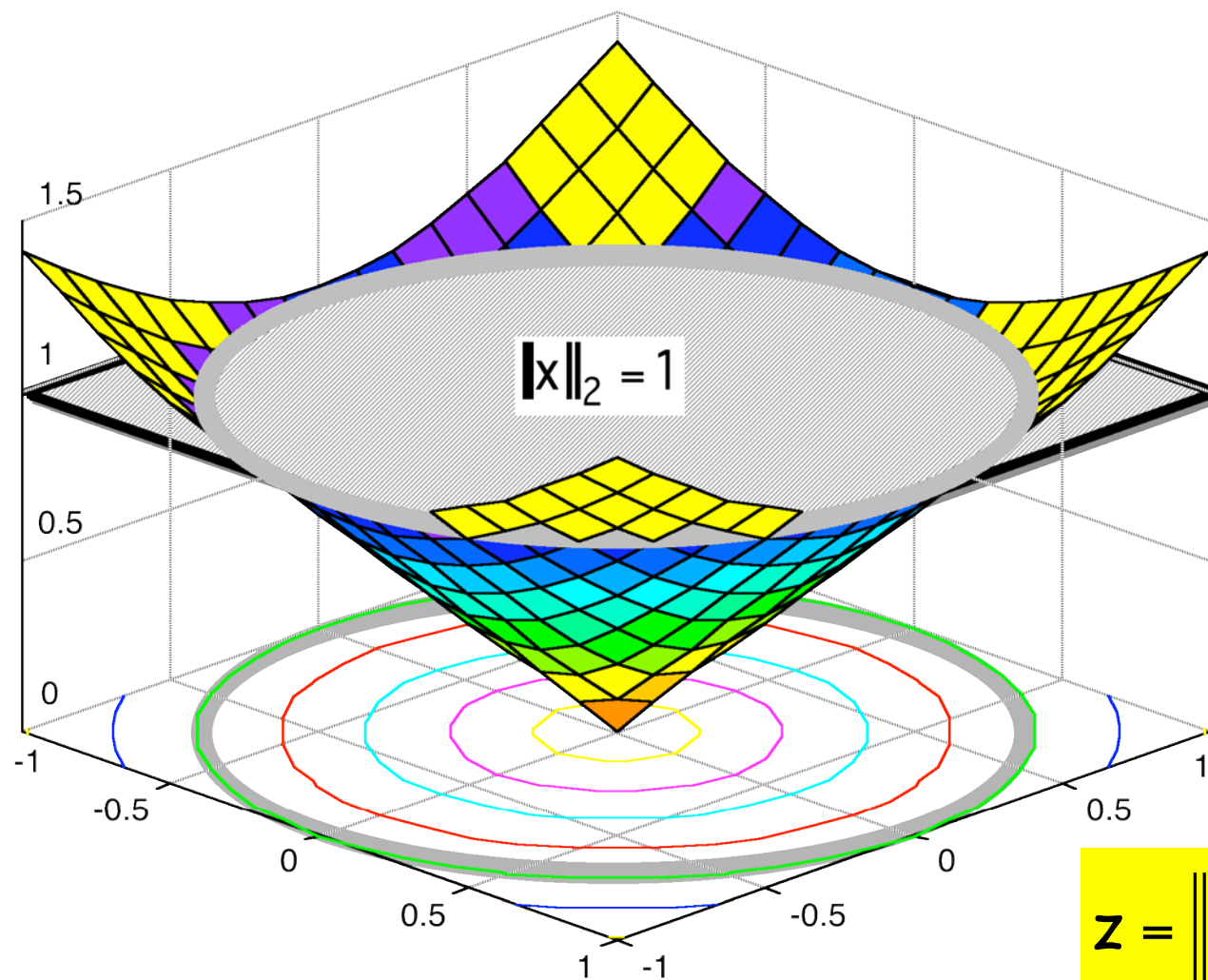


Level Sets and Open Balls of 1-norm



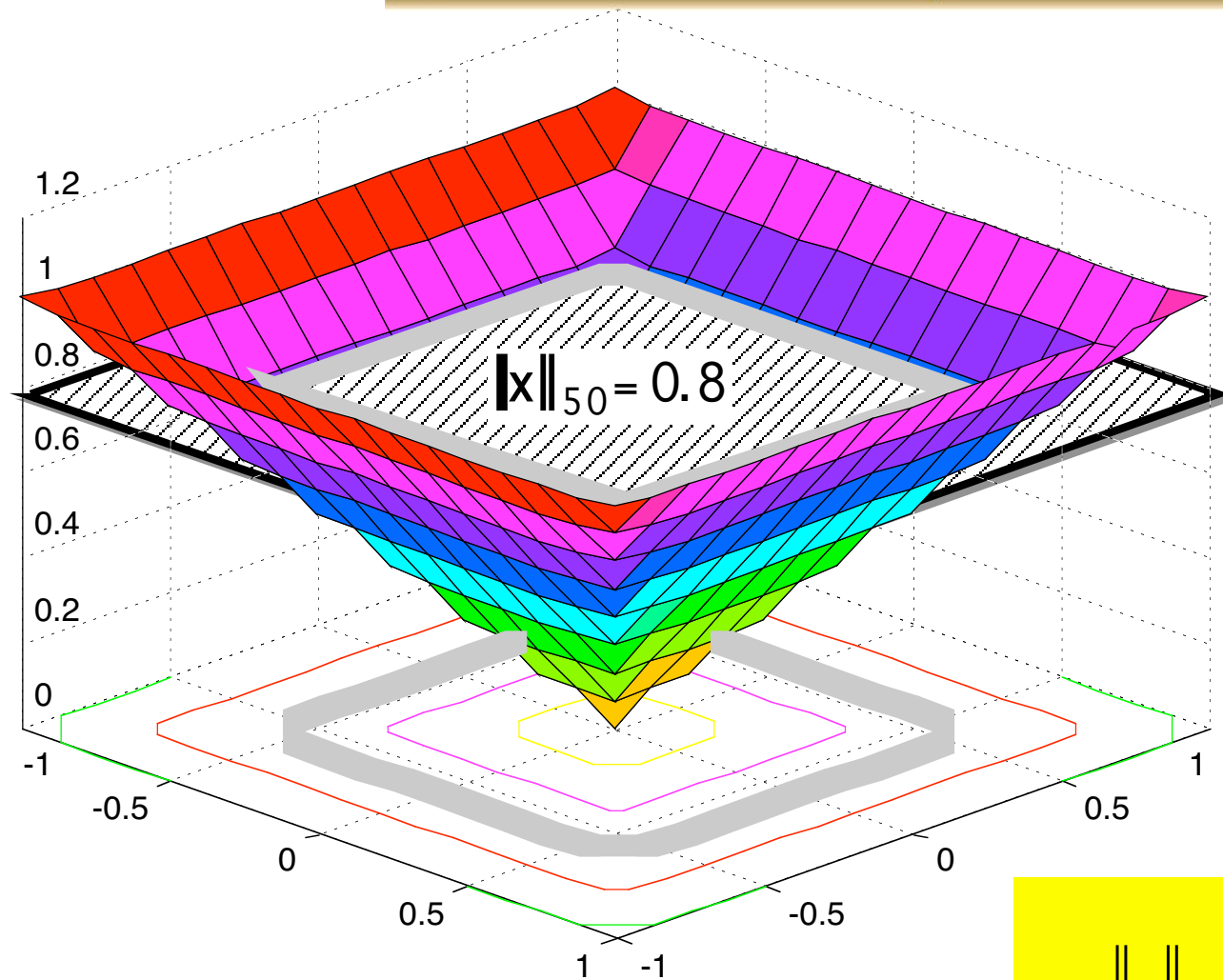
$$z = \|x\|_1 = |x| + |y|$$

Level Sets and Open Balls of 2-norm



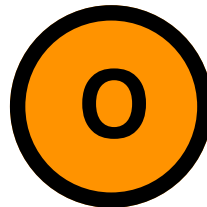
$$z = \|x\|_2 = \sqrt{x^2 + y^2}$$

Level Sets and Open Balls of 50-norm



$$z = \|x\|_{50} = \left(x^{50} + y^{50} \right)^{\frac{1}{50}}$$

OBJECTS



Feature Nomination : Issues



What are the problems to solve



What variables are important



Which variables *can be* measured



What data analysis will be done



Are *more* features better than *less*



Propose and collect raw (object) data

Feature Extraction : Why?

$$\text{Find } \underbrace{\phi : P(\mathbb{R}^p)}_X \mapsto \underbrace{P(\mathbb{R}^q)}_{Y=\phi[X]} \text{ to:}$$

1. Increase dimension : $p < q$

Functional link networks
Support Vector Machines

2. Decrease dimension : $p > q$

Decrease time and space
Eliminate "redundant" features
Assess tendencies visually ($q=2$)

3. Y should have the "same information" as X for :

Clustering and Classification
Prediction and Control

Feature Extraction : **How?**

1. Functions $f : \underbrace{\mathbb{R}^p}_X \mapsto \underbrace{\mathbb{R}^q}_{Y=f(X)}$

Linear

Axial ON Projection
Derived ON Projection (PCA)
Derived non-OG Projection

Non-Linear

3. Use your



2. Algorithms $\underbrace{A : \mathcal{P}(\mathbb{R}^p)}_X \mapsto \underbrace{\mathcal{P}(\mathbb{R}^q)}_{Y=A[X]}$

Sammon's Algorithm
FFBP Neural Nets
SVMs (in theory !)

Feature Analysis Example : Iris Data

Each object vector represents an Iris flower

$n = 150$ vectors

$p = 4$ features

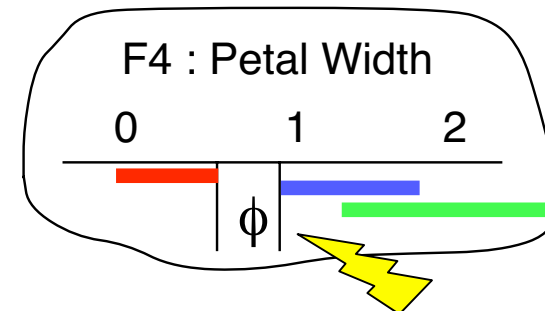
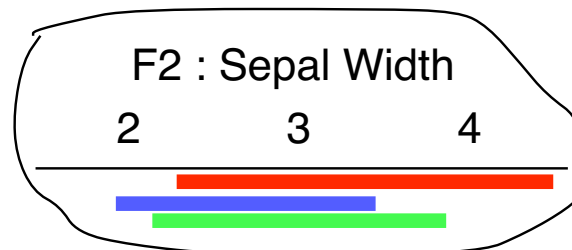
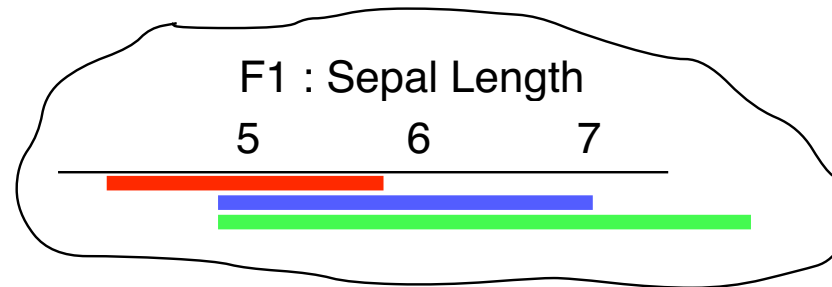
$c = 3$ *physical* subspecies

($n_i = 50$ vectors each)

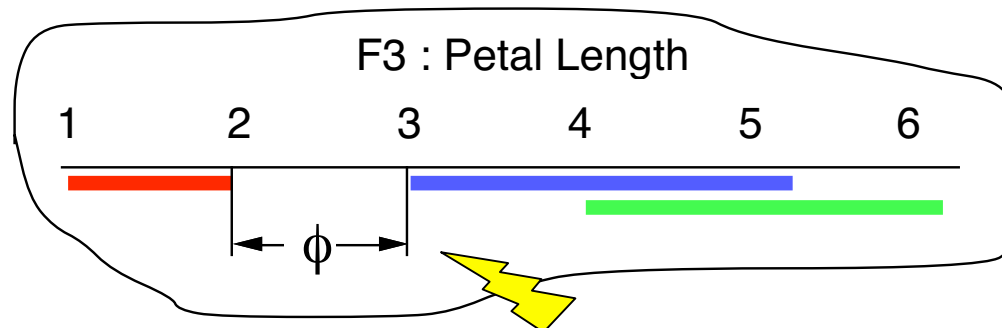
$$\mathbf{x}_k = \begin{pmatrix} x_{1k} \leftarrow \text{sepal length} \\ x_{2k} \leftarrow \text{sepal width} \\ x_{3k} \leftarrow \text{petal length} \\ x_{4k} \leftarrow \text{petal width} \end{pmatrix}$$

Visual Examination of Feature Ranges

— A = Sestosa
— B = Versicolor
— C = Virginica



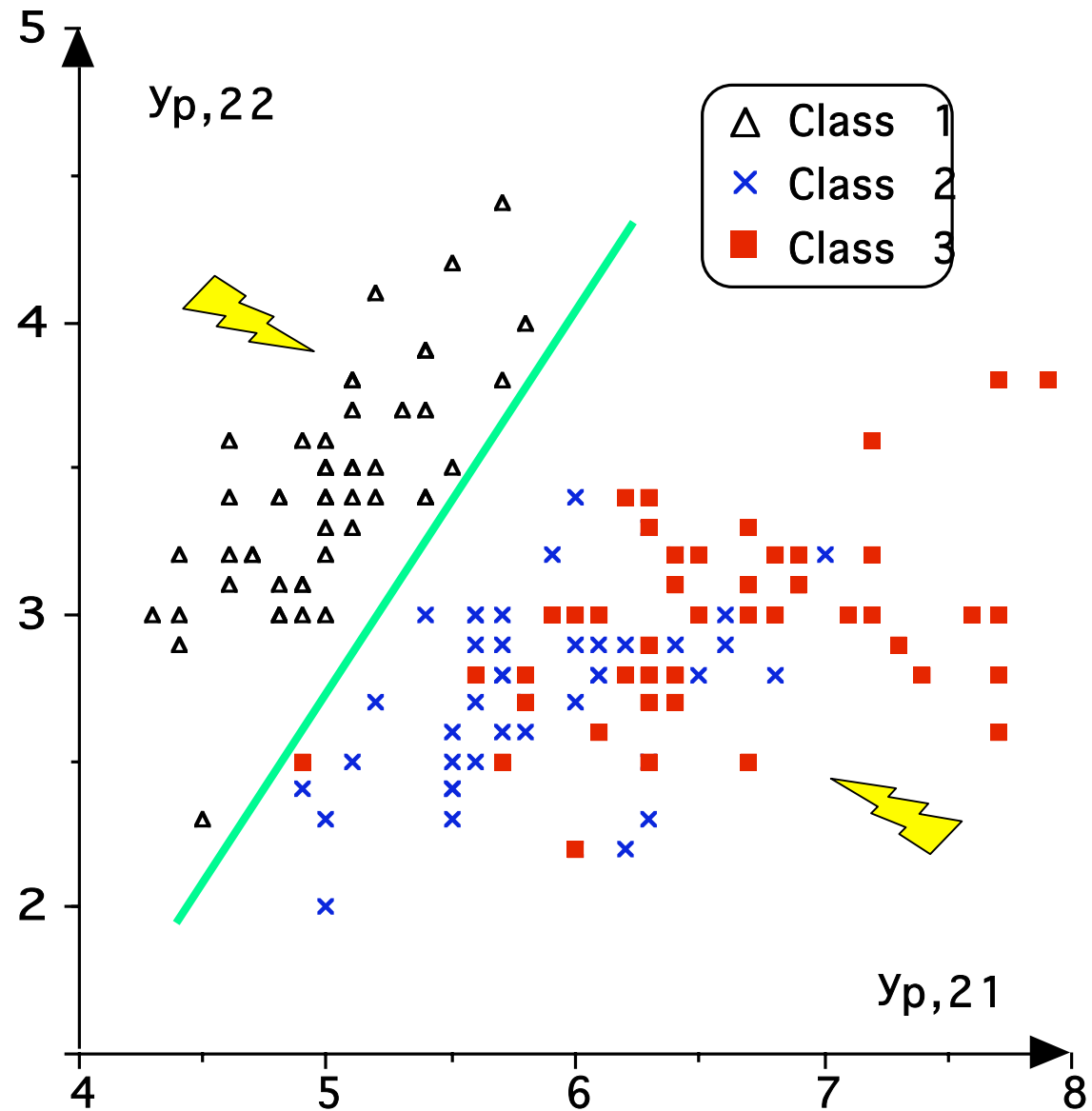
**F3 or F4
separates A
from BUC**



Linear
Selection
(Projection)

F1 vs F2

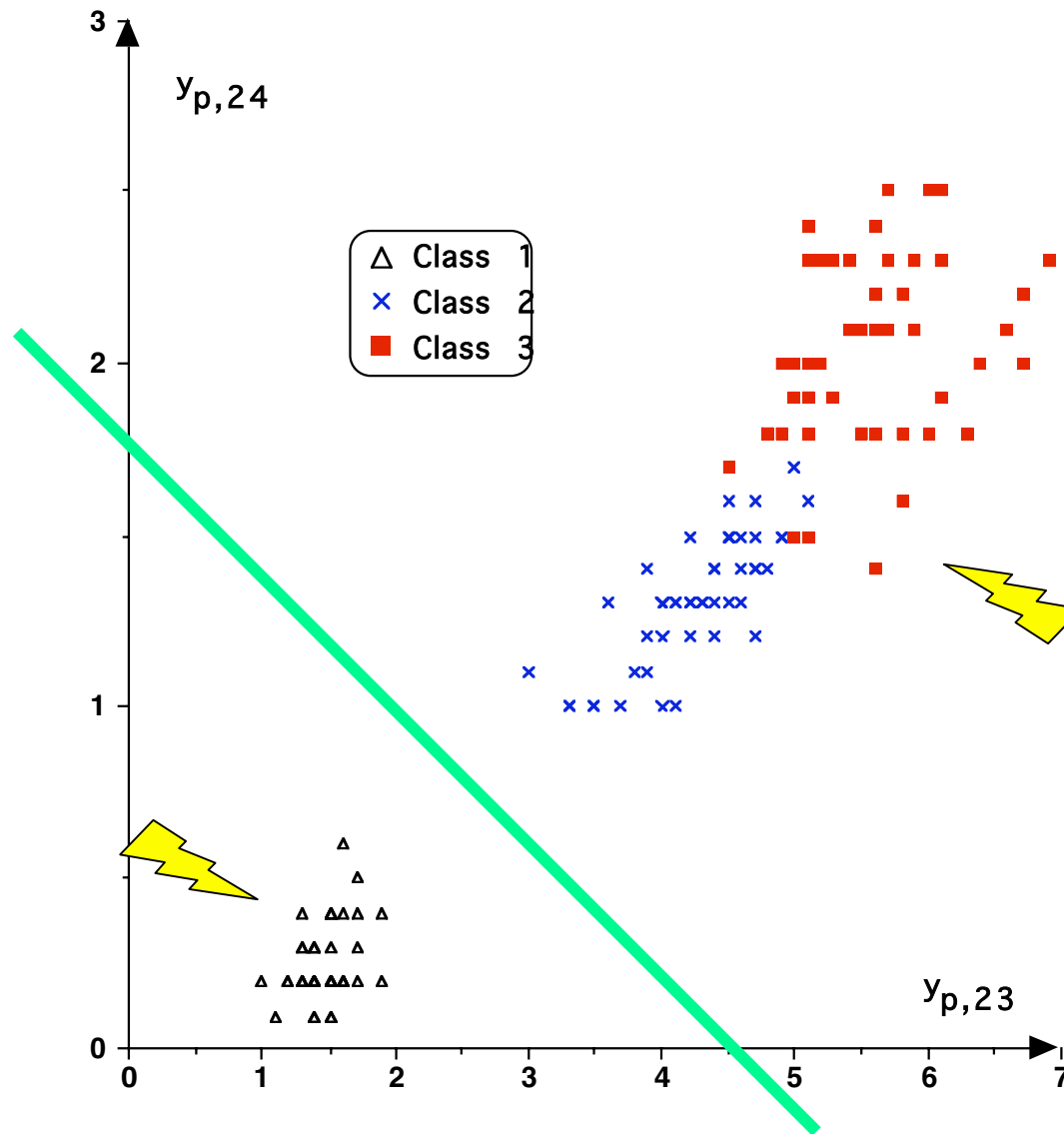
(F1, F2) *linearly*
separates A
from BUC
-----and-----
suggests that Iris
has only $c = 2$
geometric clusters



Linear Selection (Projection)

F3 vs F4

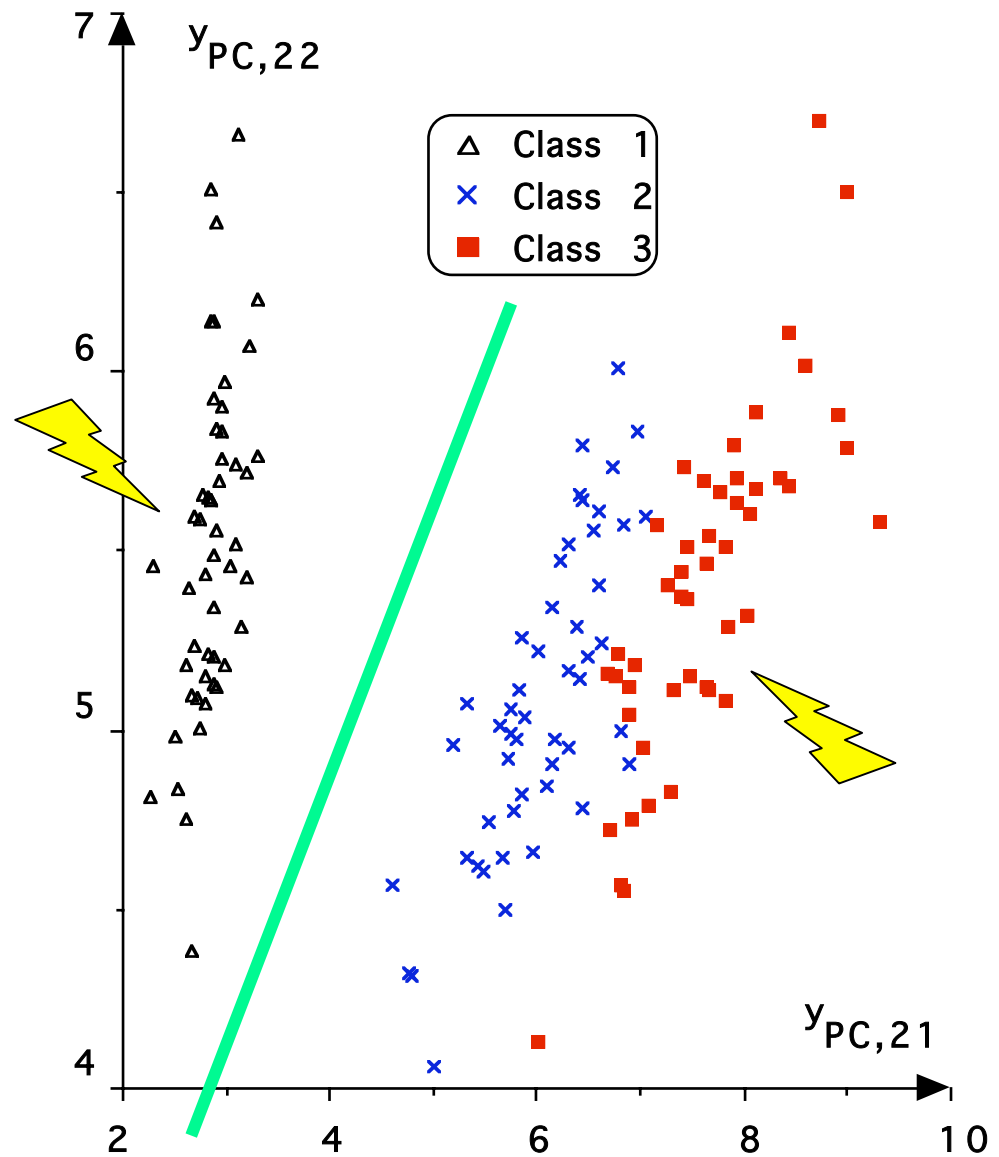
(F1, F2) *linearly* separates A from BUC
 -----and-----
 suggests that Iris has only $c = 2$
geometric clusters



Linear Extraction
(PCA Projection)

First 2 Principle
Components (y_1, y_2)

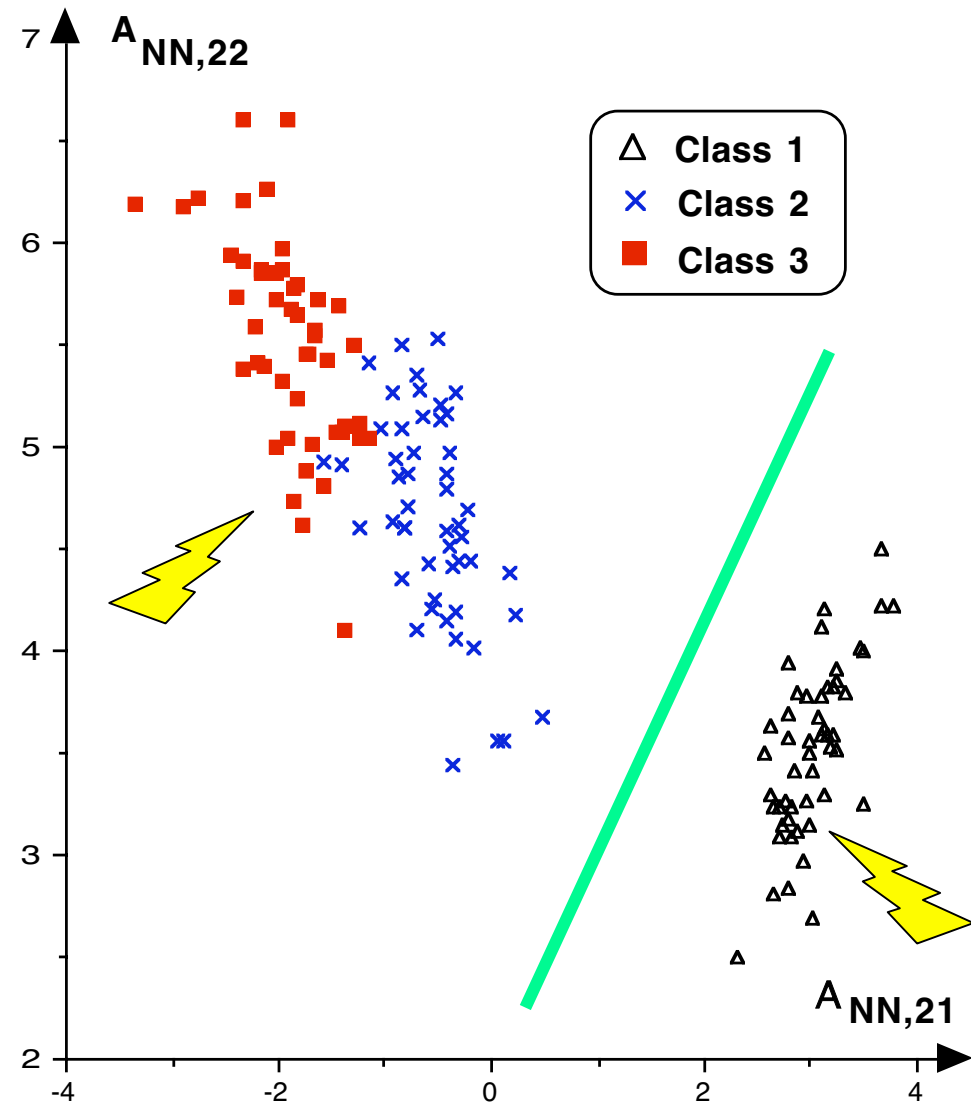
(y_1, y_2) *linearly*
separates A
from BUC
-----and-----
suggests that Iris
has only $c = 2$
geometric clusters



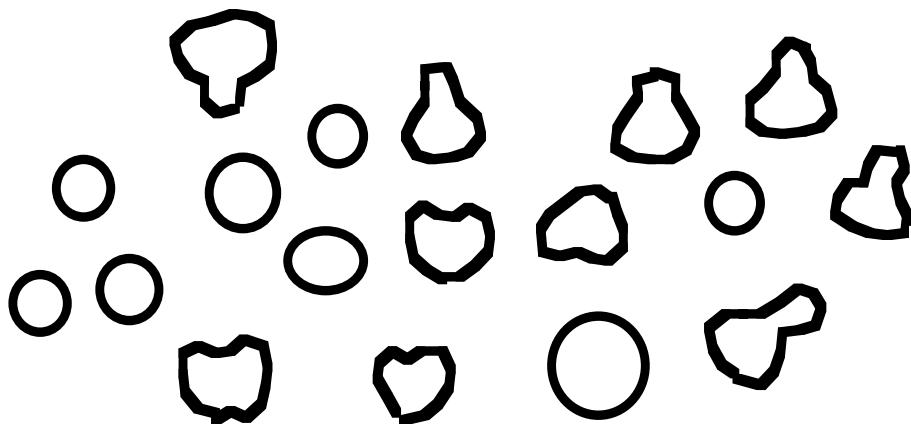
NL Extraction by
4:2:4 FFBP NN

Extracted
feature pair
(NN_{21} , NN_{22})

(NN_{21} , NN_{22})
linearly
separates A
from BUC
-----and-----
suggests that Iris
has only $c = 2$
geometric clusters



Label the similar objects

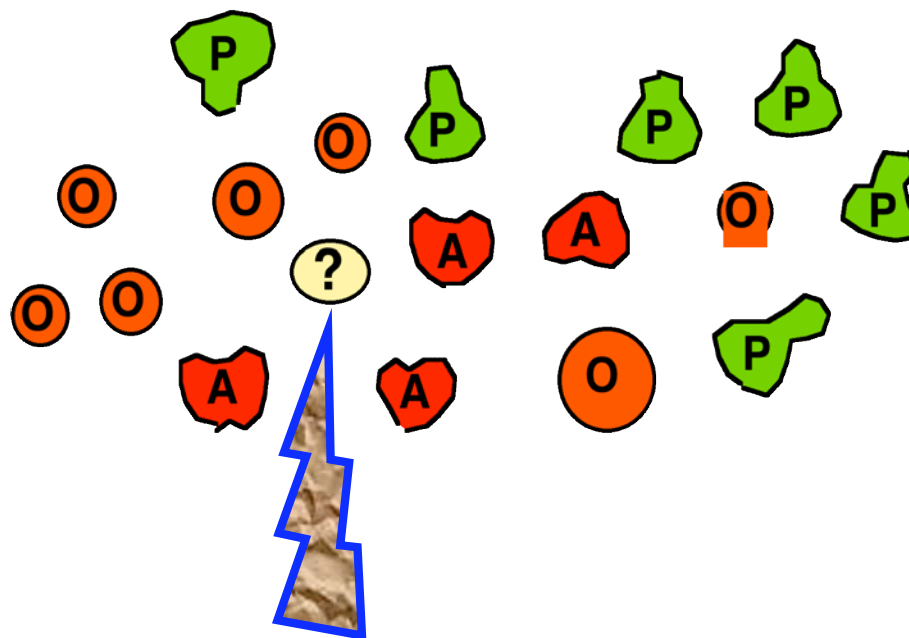


What is
Clustering ?

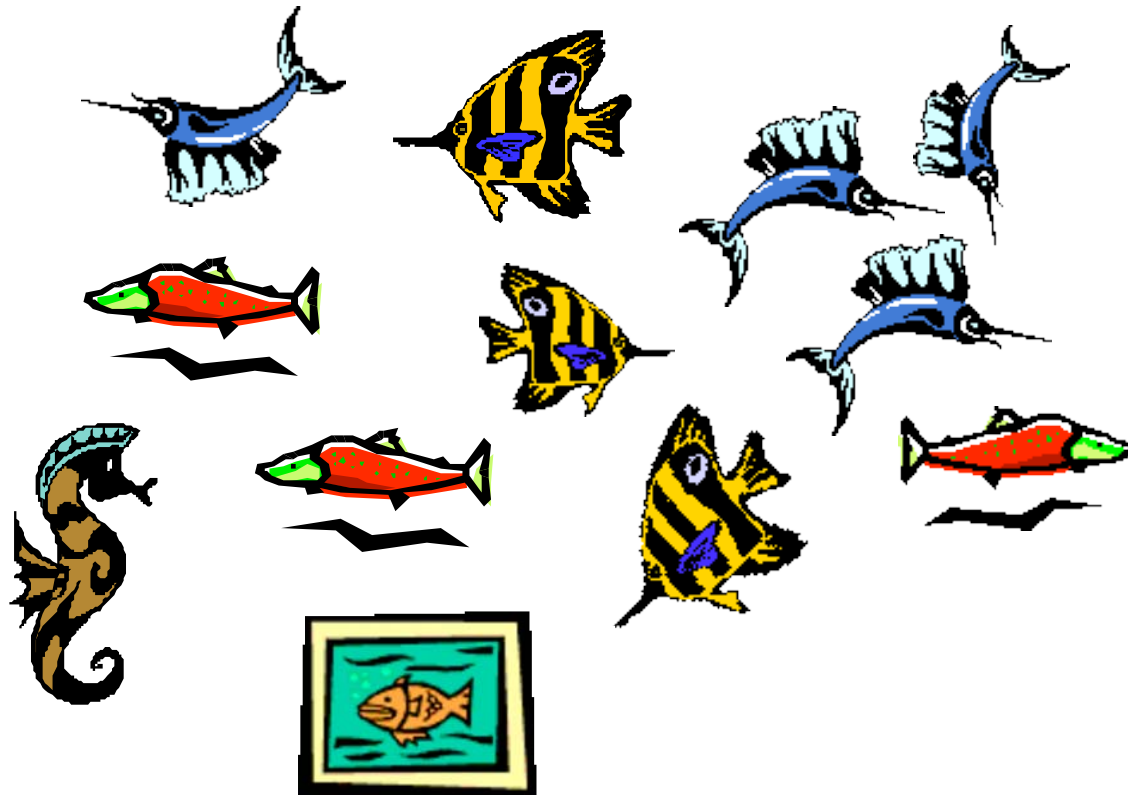
(Unsupervised Learning)

What is "similar"?

How many groups (c) ?



Data sets usually contain **mixtures** of objects



Clustering **groups** and **labels** objects

Easy for Humans !

Hard for computers !

Confusing Objects

Similarity Mistakes



What is "similar"?



Wrong features?

Mystery Data (Noise)



Incomplete data

Non-linear boundaries

How many groups (c)?

What kind of labels?

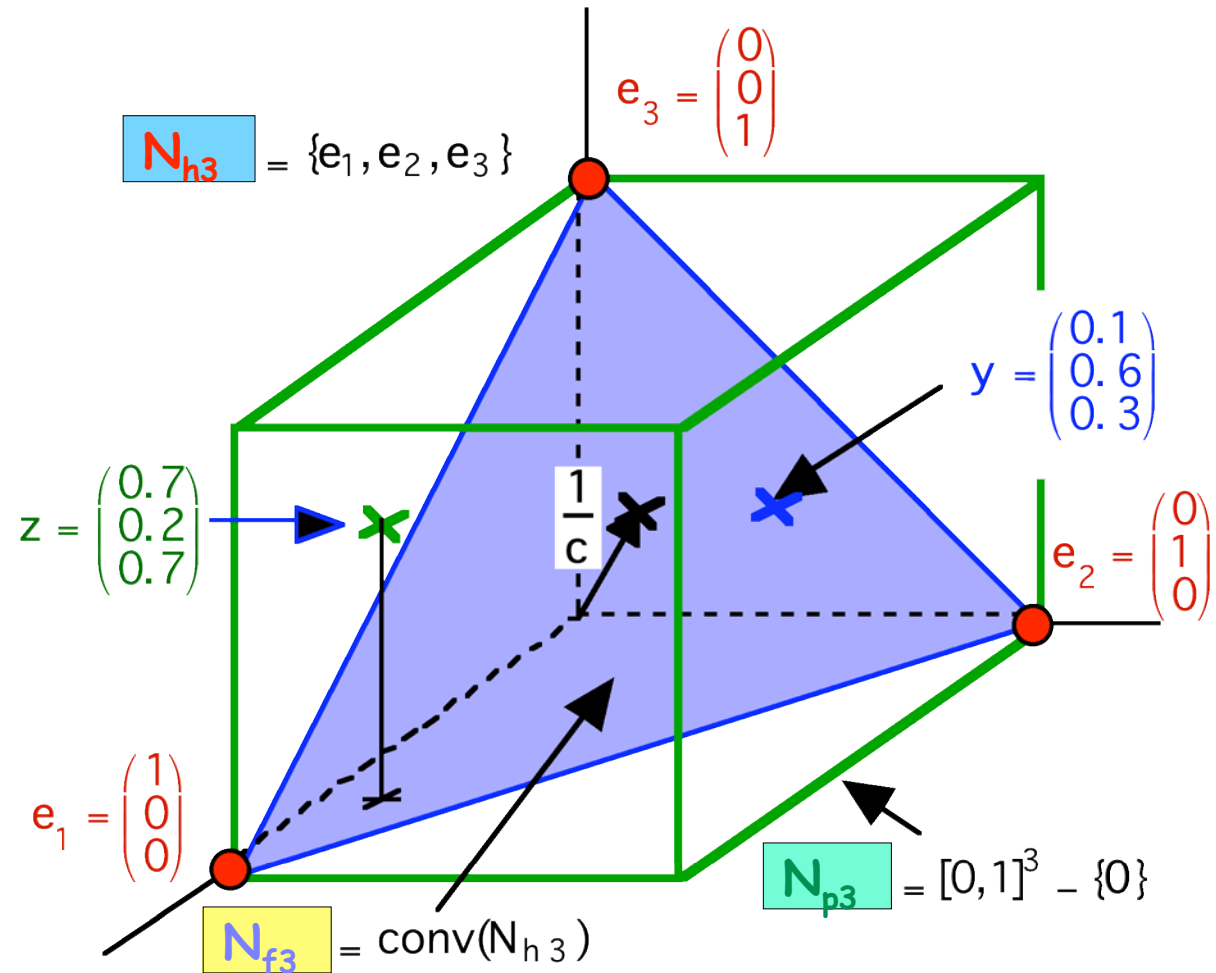
Yipes !

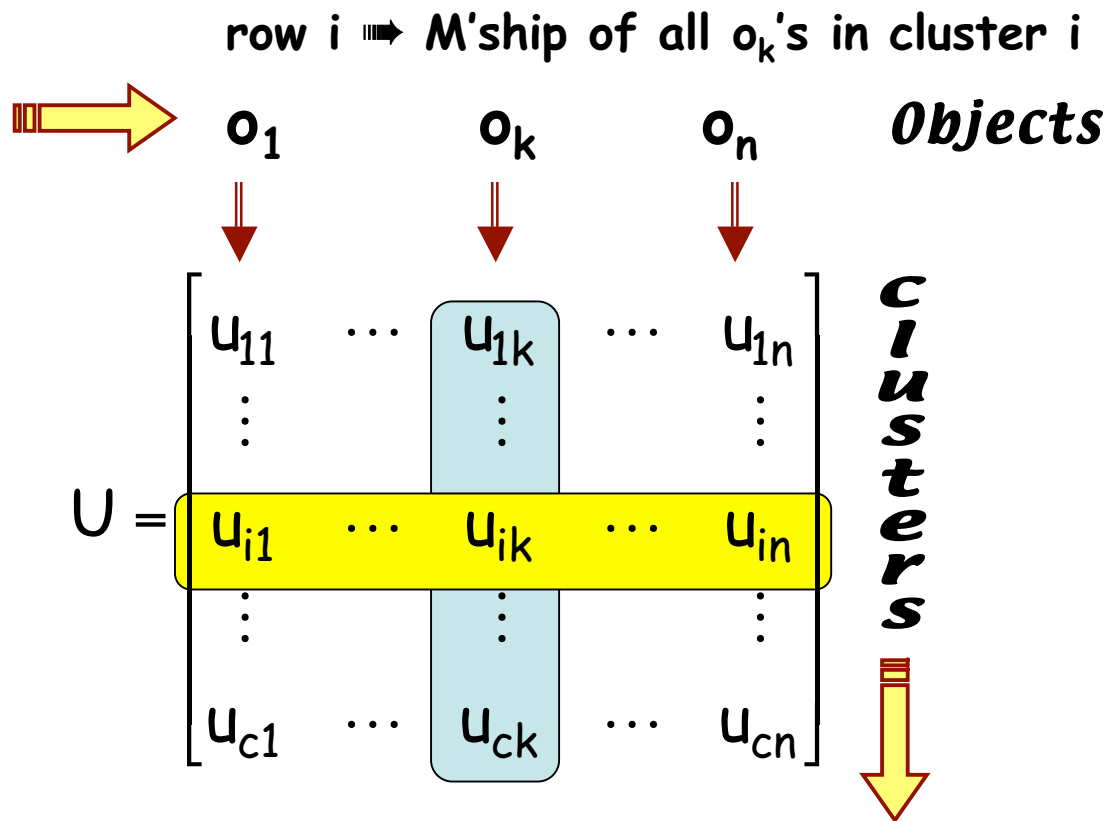
Unipolar Label Vectors @ $c = 3$

Crisp
= Vertices

Fuzzy or
Probabilistic
= Triangle

Possibilistic
= Cube - $\{0\}$





col $j \Rightarrow$ M'ship of o_k in each cluster

Partition Matrices

Membership Functions

$$u_i: O \rightarrow [0, 1]$$

$u_i(o_k) = u_{ik} =$ M'ship of o_k in cluster i

Crisp

Fuzzy/Prob

Possibilistic

Row sums

$$\sum_k u_{ik} > 0$$



Col sums

$$\sum_i u_{ik} = 1$$



$$\sum_{i=1}^c u_{ik} \leq c$$

M'ships

$$u_{ik} \in \{0, 1\}$$

$$u_{ik} \in [0, 1]$$



Set Name

M_{hcn}

\subset

M_{fcn}

\subset

M_{pcn}

Example

1 0 0 0
0 1 0 0
0 0 1 1

1 .07 0 .44
0 .91 0 .06
0 .02 1 .50

1 .07 1 .44
0 .91 0 .52
0 .02 1 .38

Take a 2nd



Partition Set Theory

$$M_{fcn} = \text{conv}(M_{fcn0}) = \text{convex polytope}$$

Each face of M_{fcn0} looks like N_{fc}

$$U = [1/c] = \text{centroid of } M_{fcn}$$

$$|M_{hcn}| = \left(\frac{1}{c!} \right) \sum_{j=1}^c \binom{c}{j} (-1)^{c-j} j^n \approx \left(\frac{c^n}{c!} \right)$$

$$\dim(M_{fcn}) = n(c-1)$$

Non Degenerate

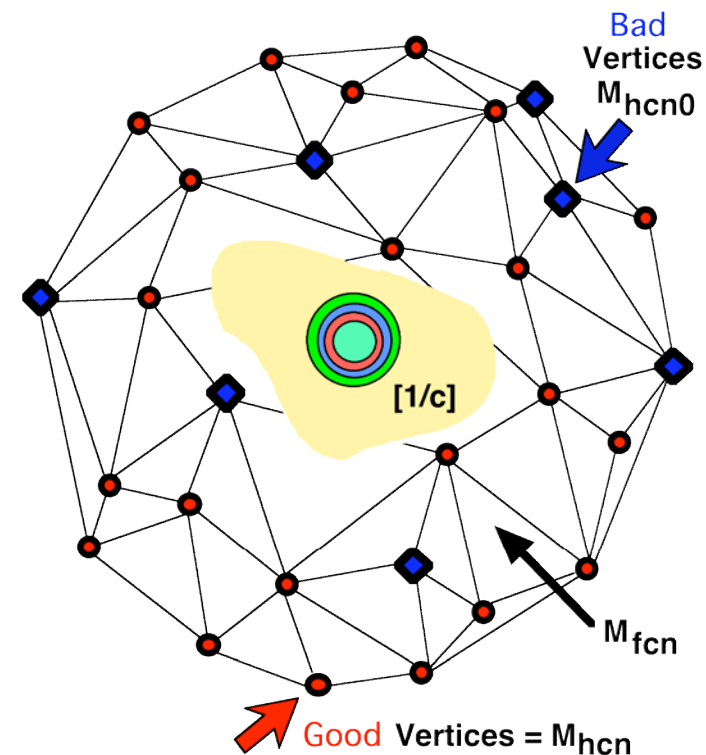
$$\left(0 < \sum_{k=1}^n u_{ik} \right)$$

$$M_{hcn} \subset M_{fcn}$$

Degenerate (0 rows)

$$\left(0 \leq \sum_{k=1}^n u_{ik} \right)$$

$$M_{hcn0} \subset M_{fcn0}$$



Hardening Partitions

Defuzzification (**Deprobabilization = Bayes Rule** for $U \in M_{fcn}$!)

$$U = \begin{bmatrix} 1 & .07 & .6 & .5 \\ 0 & .81 & .2 & 0 \\ 0 & .12 & .7 & .5 \end{bmatrix} \in M_{pcn} \quad H_1(U) = U_{MM} = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \in M_{hcn}$$

↑ ↑

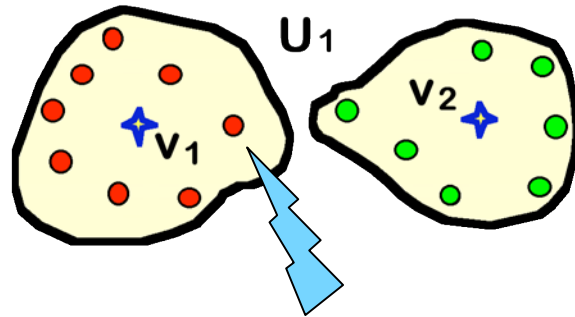
resolve ties randomly

α -cut thresholding ($0 < \alpha < 1$) **filters** chosen levels of m'ship

$$H_{.9}(U) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \in M_{hcn} \quad H_{.8}(U) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \in M_{hcn}$$

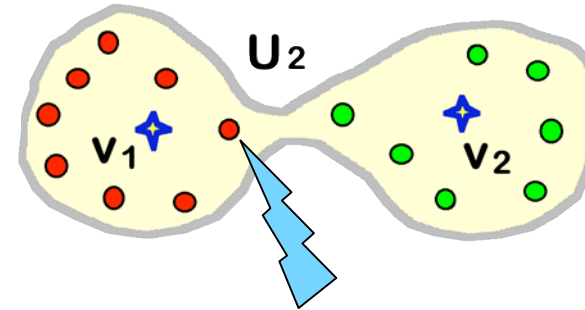
↑ ↑

Clusters represented by partition matrix (U) *or* prototypes $V = \{v_k\}$



$$U_c = \begin{array}{cccc|cccc} 1 & 1 & \dots & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 1 & 1 & \dots & 1 \end{array}$$

Crisp : hard boundaries



$$U_f = \begin{array}{cccc|cccc} .9 & 1 & \dots & .6 & .3 & 0 & .2 & \dots & 0 \\ .1 & 0 & \dots & .4 & .7 & 1 & .8 & \dots & 1 \end{array}$$

Fuzzy : soft boundaries

4 types of clustering models

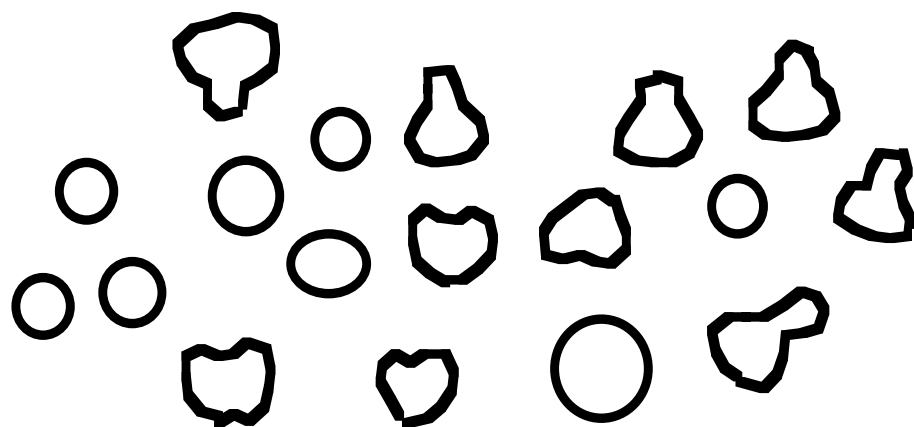
U only models

V only models

(U, V, +) models

(U, V) models

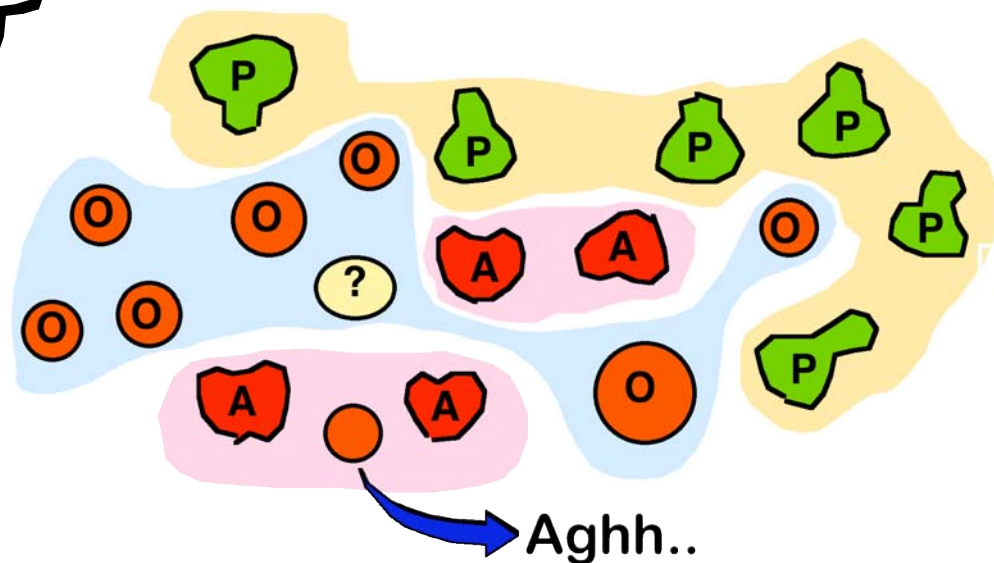
Find regions that
contain similar objects



What are
" (decision) regions"?

What is
Classification ?

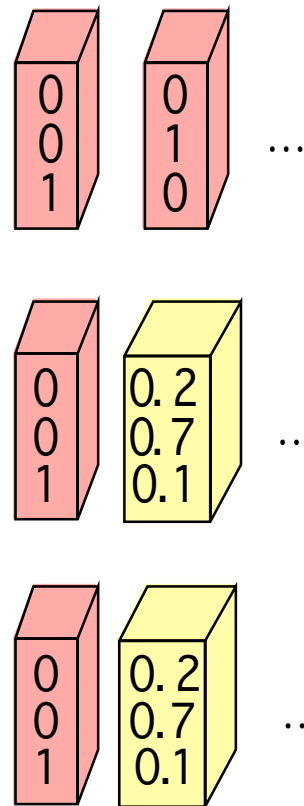
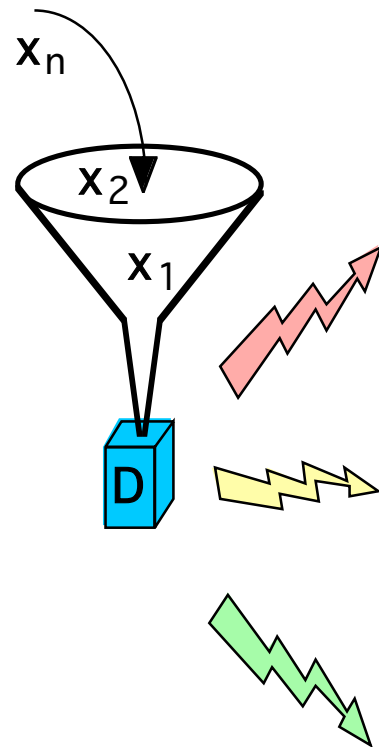
(Supervised Learning)



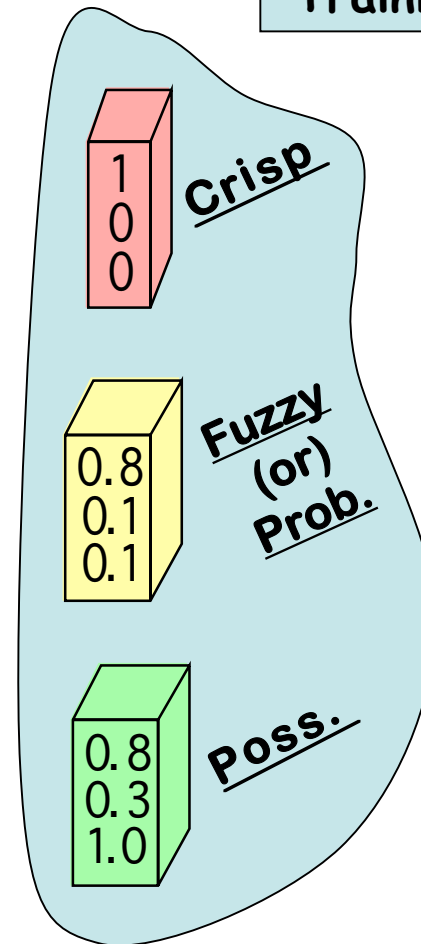
Classifier Functions

$$D : \mathcal{R}^p \mapsto N_{pc} = [0, 1]^c$$

Feature
Vectors
in \mathcal{R}^p



3 Kinds of
Training Data !



Label
Vectors
in N_{pc}

Classifier Training

1

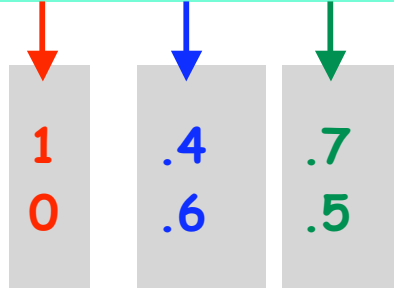
Choose a family of models $\mathcal{F} = \{ D_1(\theta) \dots D_i(\theta) \dots \}$

2

Training Data $X_t = \{ \mathbf{x}_{t,1} \dots \mathbf{x}_{t,k} \dots \mathbf{x}_{t,n} \}$

3

Training Labels $U_t =$



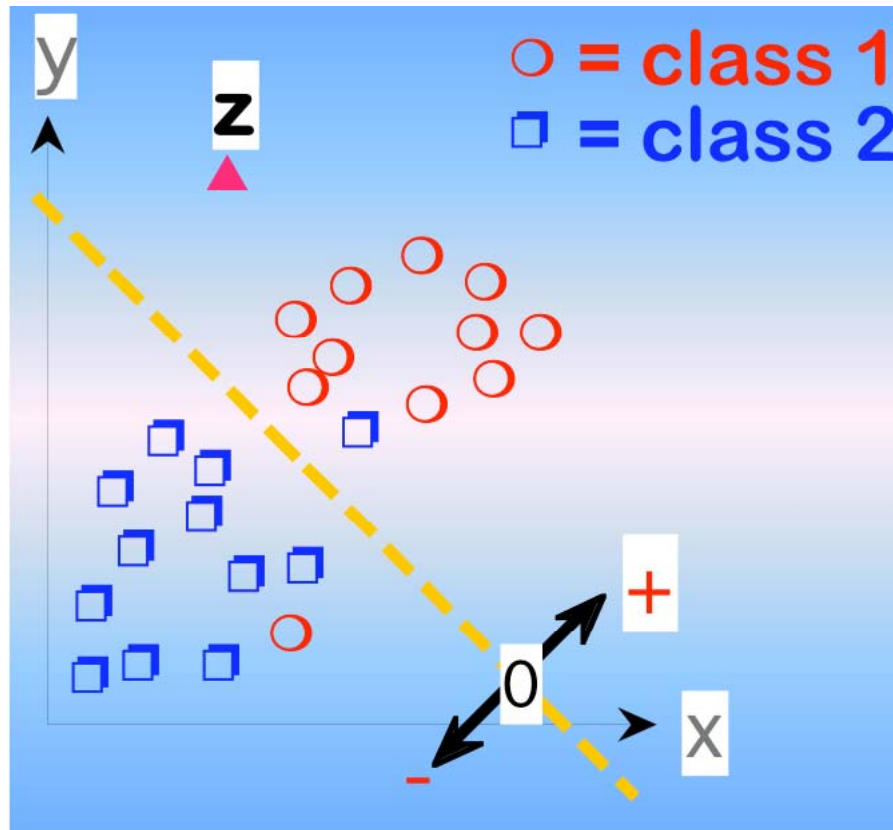
4

Algorithm A “looks for” an optimal (θ^*) of D

5

Use (X, U_x, A) to find (θ^*) of D (“learning”)

$$d_{HP}(x) = \underbrace{\langle x, \vec{w} \rangle + \alpha}_{\text{(Linear) HP Model}} + \underbrace{X_+}_{\text{Data}} + \underbrace{U_+}_{\text{Labels}} + \underbrace{\text{Hyperplane}}_{\text{Algorithm}}$$



$$\theta = \left[\begin{pmatrix} w_1 \\ w_2 \end{pmatrix}, \alpha \right]$$

$$d_{HP}(x) = x + y - 1$$

(e.g. SVM) Classifier

$$\theta^* = \left[\begin{pmatrix} 1 \\ 1 \end{pmatrix}, -1 \right]$$

D_{HP}^* *crisply* labels *every* point in \mathbb{R}^p

Classifier Function

$$D_{HP}^*(x) = \left\{ \begin{array}{lll} e_1 = (1 \ 0)^T & ; d_{HP}(x) > \alpha^* & \Rightarrow x \in 1 \\ e_1 \text{ or } e_2 & ; d_{HP}(x) = \alpha^* & \Rightarrow \text{Tie} \\ e_2 = (0 \ 1)^T & ; d_{HP}(x) < \alpha^* & \Rightarrow x \in 2 \end{array} \right\}$$

Decision (or Discriminant) function

Mistakes can be costly : a doctor treats patient x with A or B

Use (A| $x \in 1$) \Rightarrow cures x

Use (B| $x \in 1$) \Rightarrow no cure, no harm

Use (B| $x \in 2$) \Rightarrow cures x

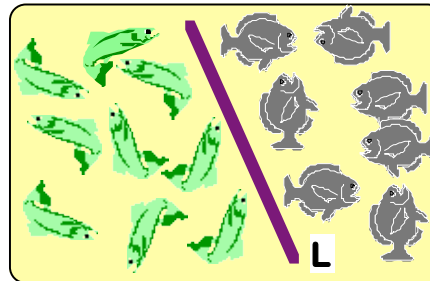
Use (A| $x \in 2$) \Rightarrow ☠️ - a cure?

Patient x *wants* D_{HP}^* to be *soft*

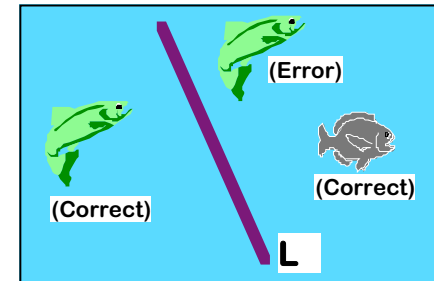
$$D_{HP}^*(x) = \left\{ \begin{array}{ll} \begin{pmatrix} 0.95 \\ 0.05 \end{pmatrix} & \Rightarrow \begin{pmatrix} \text{use} \\ A \end{pmatrix} \\ \begin{pmatrix} 0.53 \\ 0.47 \end{pmatrix} & \Rightarrow \begin{pmatrix} \text{more} \\ \text{tests} \end{pmatrix} \end{array} \right\}$$

3 Kinds of Training and Operation

Non-Adaptive
Off Line **Training**



Non-Adaptive
On Line **Operation**



3 Data Sets

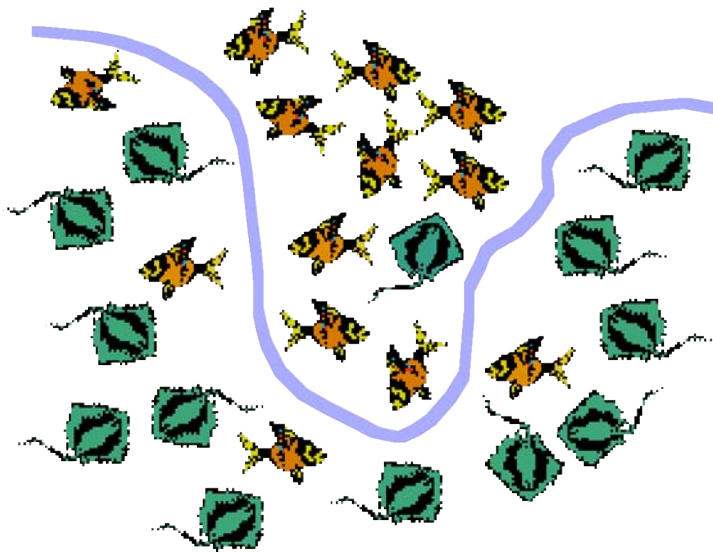
needed to Train and Test a Classifier

Training $X_{tr} = \{\mathbf{x}_{tr,1} \dots \mathbf{x}_{tr,n}\}$ with labels $u_{tr} \in N_{pc}$

Validation $X_v = \{\mathbf{x}_{va,1} \dots \mathbf{x}_{va,s}\}$ with labels $u_{va} \in N_{hc}$

Test $X_t = \{\mathbf{x}_{te,1} \dots \mathbf{x}_{te,q}\}$ with labels $u_{te} \in N_{hc}$

Data can be uncooperative !



Non-Linear Clusters

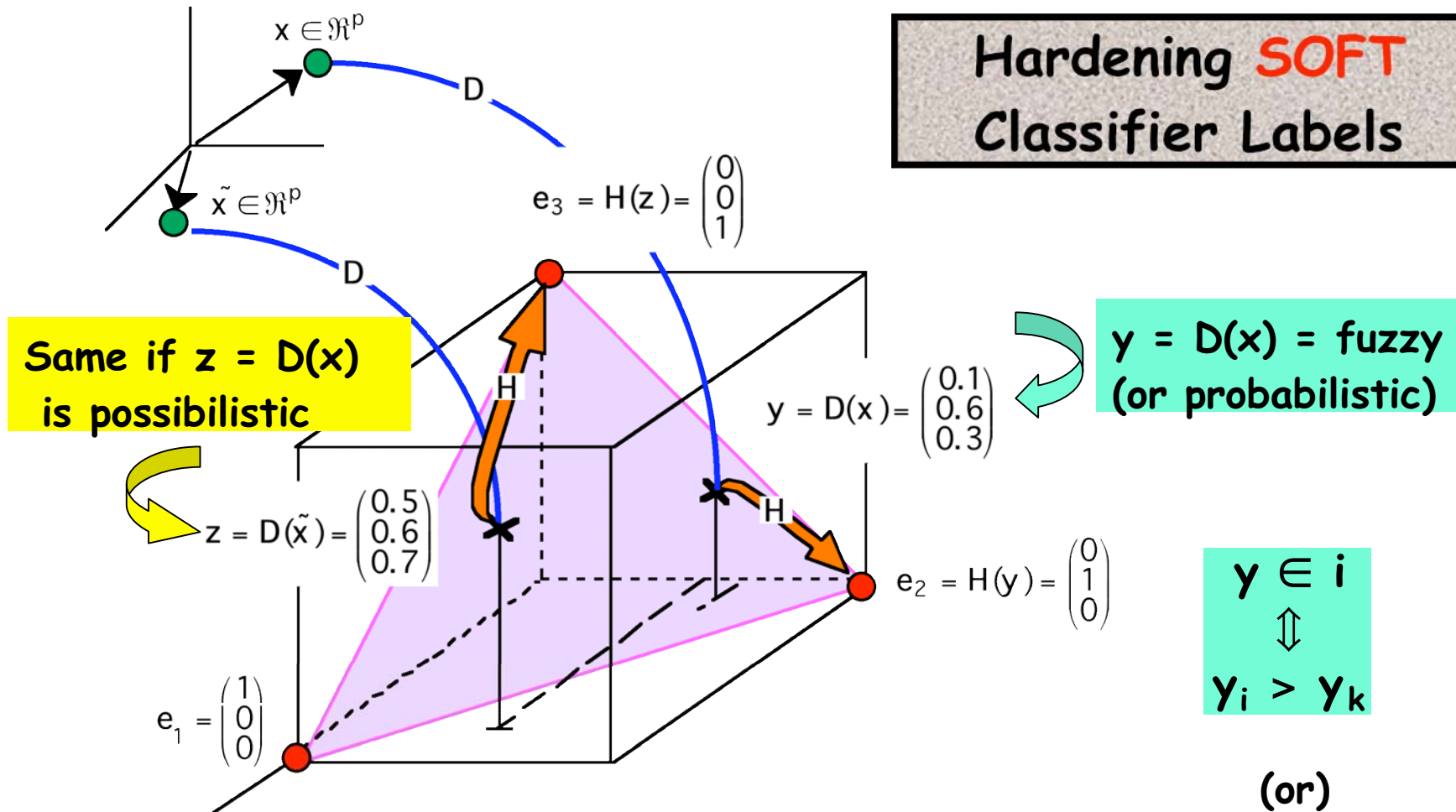
Noise & Outliers

Mixtures of Points

High Dimensional

Wrong Features

Hardening **SOFT** Classifier Labels



Error rate estimation

$D(X)$ = **classifier** learned by applying algorithm A to X

$U_{D(X)}$ = (possibly) **soft labels** given to X by classifier D

$H(U_{D(X)})$ = **crisp labels** given to X by hardening $U_{D(X)}$

U_X = (known) **crisp labels** of data set X (X_{va} or X_{te})

Empirical Error Rate
(generalization error)

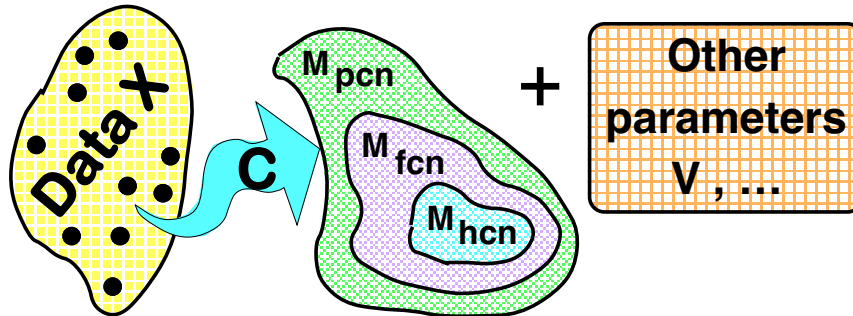
$$E_D(X) = \left(\frac{\# \text{ wrong}}{\# \text{ tried}} \right) = \frac{\|U_X - H(U_{D(X)})\|}{2|X|}$$

Minimize E_D **iteratively** by learning with X_{tr} and then computing $E_D(X_{va})$

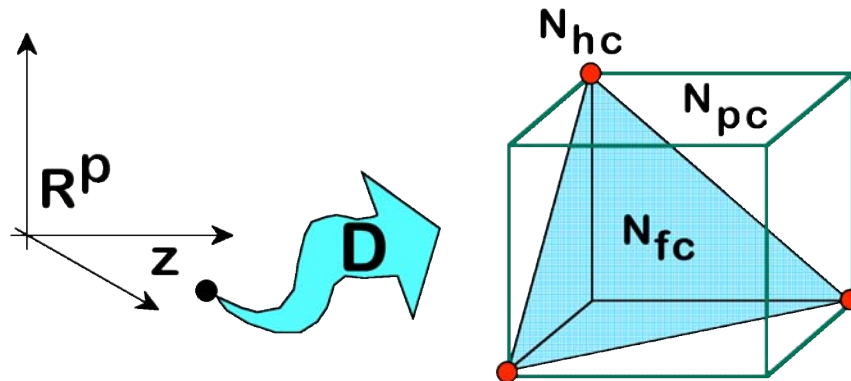


Estimate E_D **ONCE** by computing $E_D(X_{te})$

Clustering vs Classification : Summary



Clustering labels *only* the submitted data



Classifiers can label *all* the data in their domain space

Building a PR system - the *big*



Process Description	Feature Analysis	Cluster Analysis	Classifier Design
Model Type	Operation	Criterion	Model Basis
Numerical Rule-Based ⋮ Syntactic	Extraction Filtering ⋮ 2D-Displays	Dissimilarity Connectivity ⋮ Frequency	Deterministic Fuzzy ⋮ Statistical
Data Needed	Model + Algorithm	Model + Algorithm	Model + Algorithm
Type Structure Nomination Collection ⋮ Sensors	Projection PCA/LDA Sammon SOMs ⋮ VAT-coVAT	c Means SAHN Entropy GMD ⋮ ART NN	FFBP NN 1-np rules K-nn rules SVMs ⋮ Bayes Quad.



Numerical Pattern Recognition



Duda, R. and Hart, P. (1973). Pattern Classification and Scene Analysis, Wiley Interscience, NY, 1973.



Devijver, P. and Kittler, J., Pattern Recognition: A Statistical Approach, Prentice-Hall, Englewood Cliffs, NJ, 1982.



Theodoridis, S. and Koutroumbas, K. (2003). Pattern Recognition, 2nd ed., Academic Press, NY.

Bezdek, J. C. (1981). Pattern Recognition with Fuzzy Objective Function Algorithms, Plenum, NY.

Bezdek, J. C. and Pal, S. K. (1992). Fuzzy Models for Pattern Recognition, IEEE Press, Piscataway, NJ.

Bezdek, J. C., Keller, J. M., Krishnapuram, R. and Pal, N. R. (1999). Fuzzy models and algorithms for Pattern Recognition and Image Processing, Springer, NY.



Machine Learning & SVMs



Scholkopf, B. and Smola, A.J. (2002). Learning with Kernels, MIT Press, Boston

Vidyasagar, M. (2003). Learning and Generalization, Springer, London.

Devroye, L., Györfi, L. and Lugosi, G. (1996). A Probabilistic Theory of Pattern Recognition, Springer, NY,



Syntactic Pattern Recognition

Pavlidis, T. (1980) , Structural pattern recognition, Springer-Verlag, New York.

Fu, K. S. (1982). Syntactic Pattern Recognition and Applications, Prentice Hall, Inc, New Jersey.

Bunke, H., (ed.) (1992) Advances in Structural and Syntactic Pattern Recognition, World Scientific Press, Singapore.