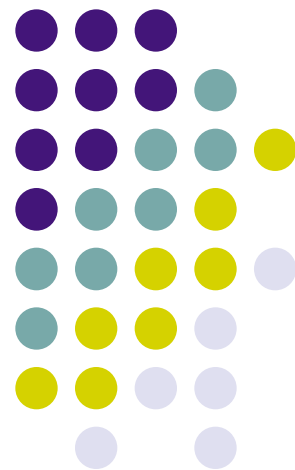


# 数字图像处理

## 第A章

### 机器学习/模式识别基础



# 机器学习/模式识别基础



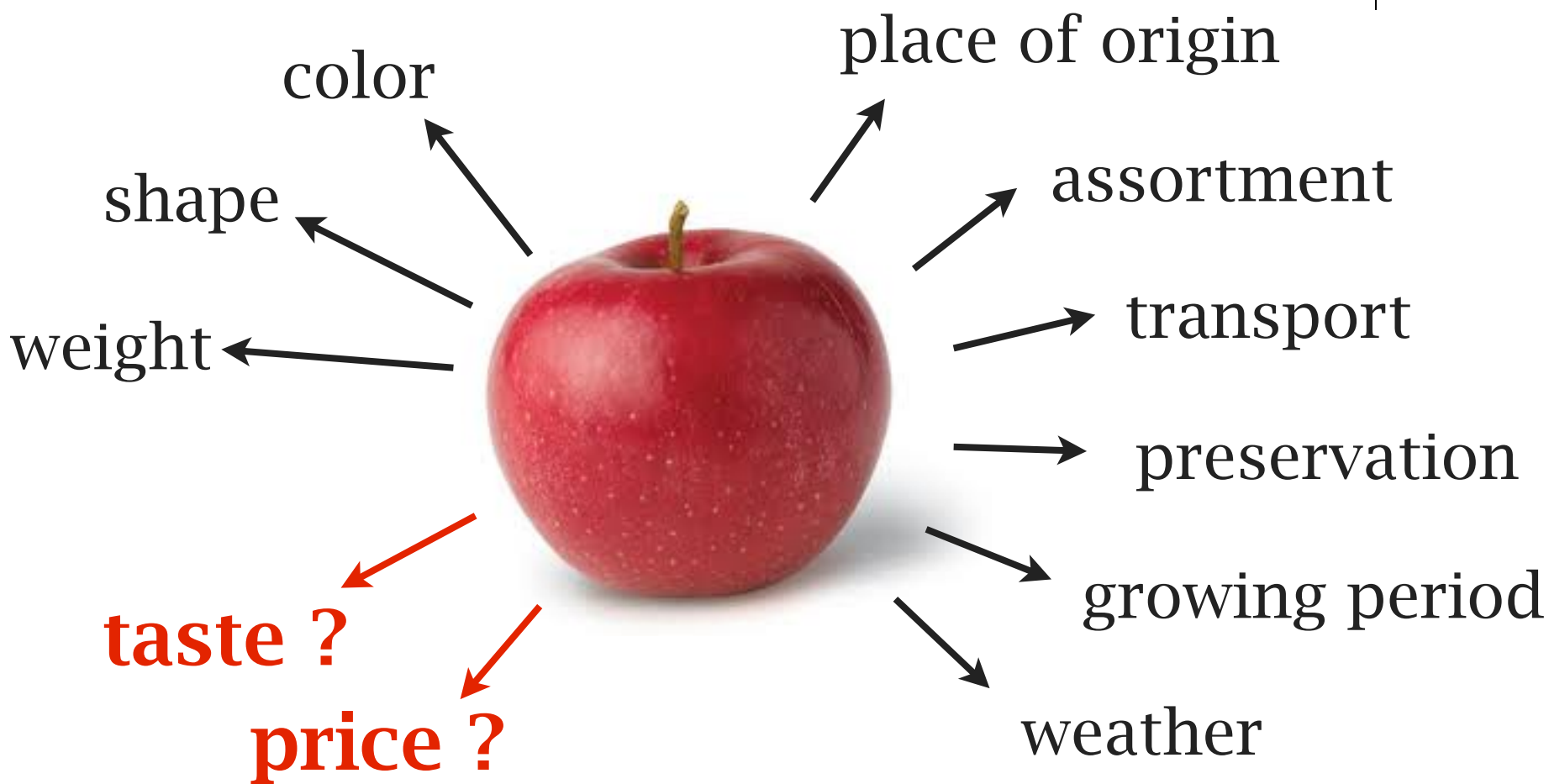
- 预测与识别
- 预测算法
- 特征提取

# 1 预测与识别



- 预测：根据当前的观测，预测未观测事件
- 识别：根据当前的观测，判断是否是预定模式
- 如何定义“事件”或“模式”？
- 机器学习方法：基于数据的定义

# 1 预测与识别



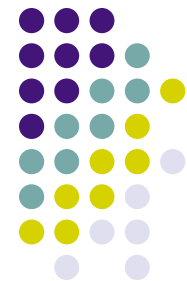
# 1 预测与识别



color={0,1,2,3} weight={0,1,2,3,4}

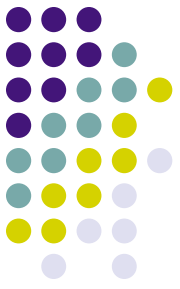
color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

# 机器学习/模式识别基础



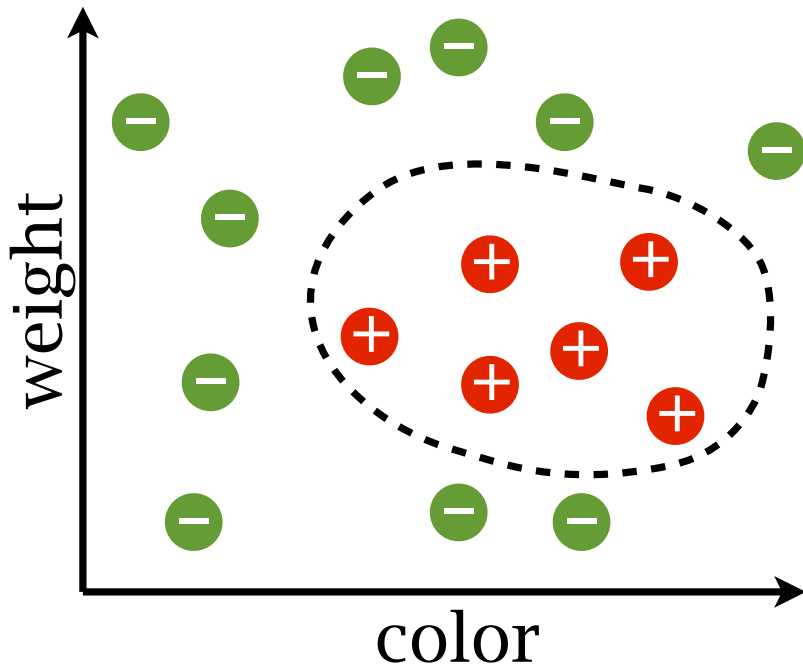
- 预测与识别
- 预测算法
- 特征提取

# Classification



**Features:** color, weight

**Label:** taste is sweet (positive/+) or not (negative/-)



(color, weight)  $\rightarrow$  sweet ?

$$\mathcal{X} \rightarrow \{-1, +1\}$$

ground-truth function  $f$

examples/training data:

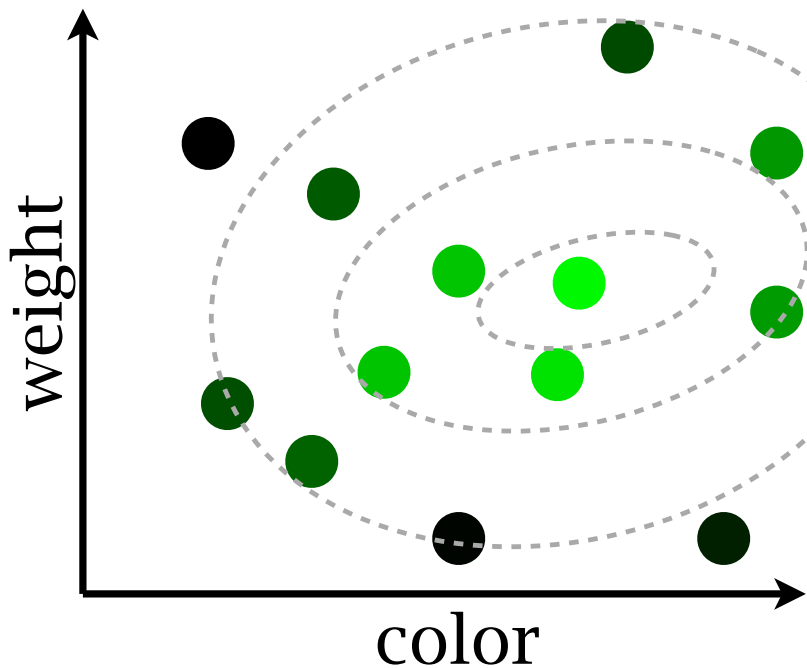
$$\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$$

$$y_i = f(\mathbf{x}_i)$$

# Regression



**Features:** color, weight  
**Label:** sweetness [0,1]



(color, weight)  $\rightarrow$  sweetness  
 $\mathcal{X} \rightarrow [-1, +1]$

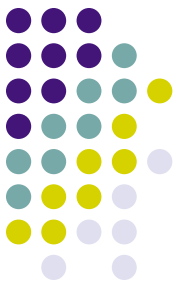
ground-truth function  $f$

examples/training data:  
 $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$

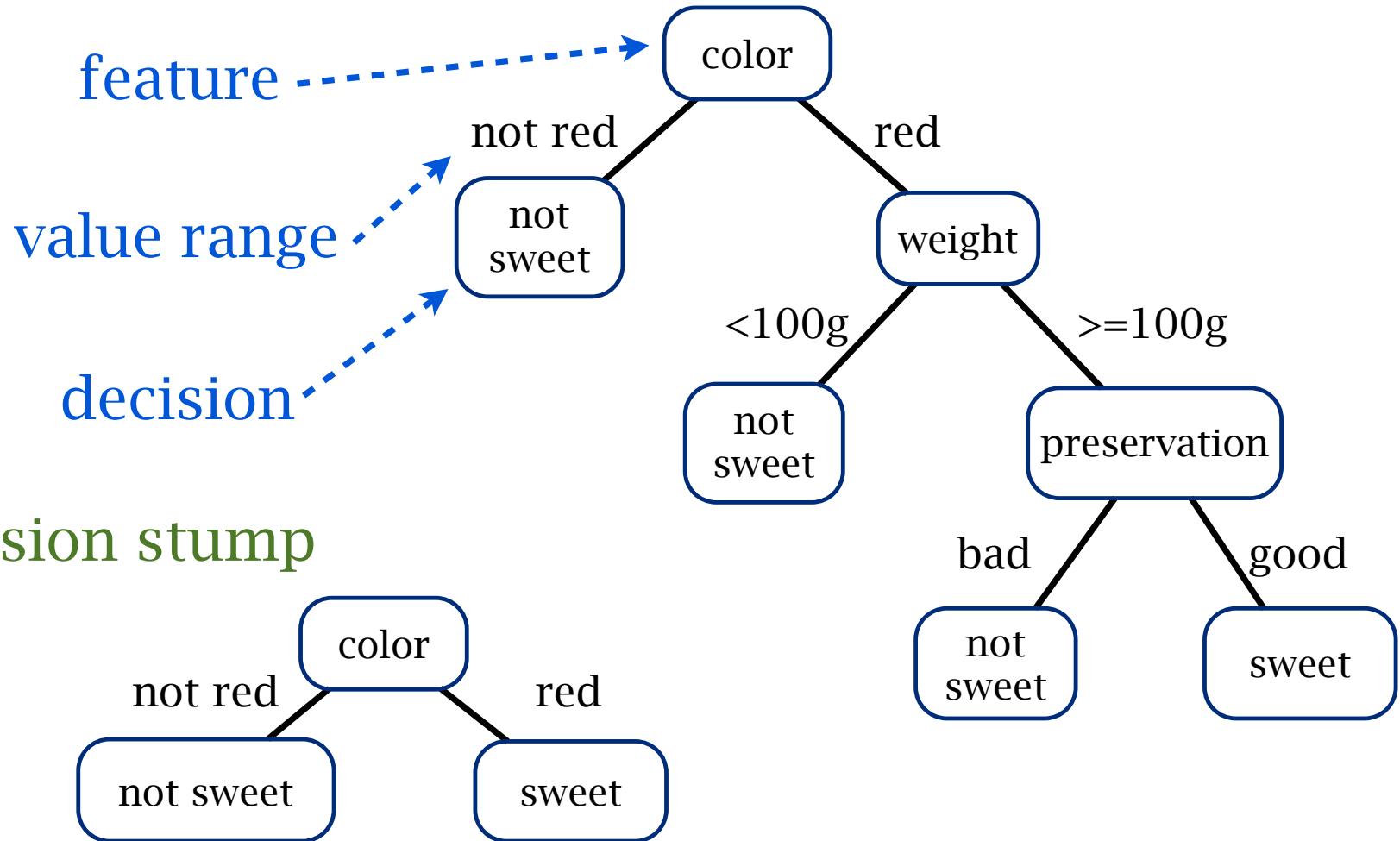
$$y_i = f(\mathbf{x}_i)$$



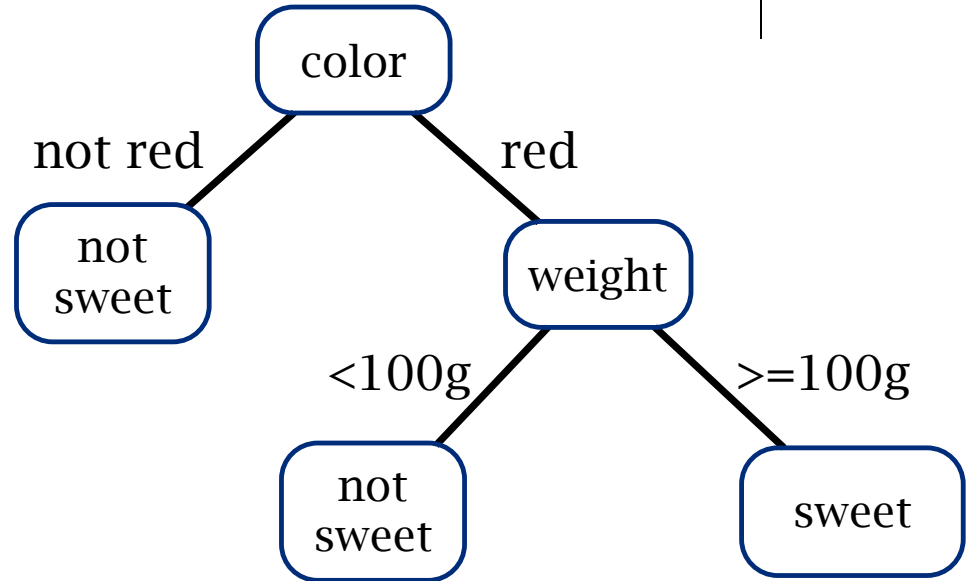
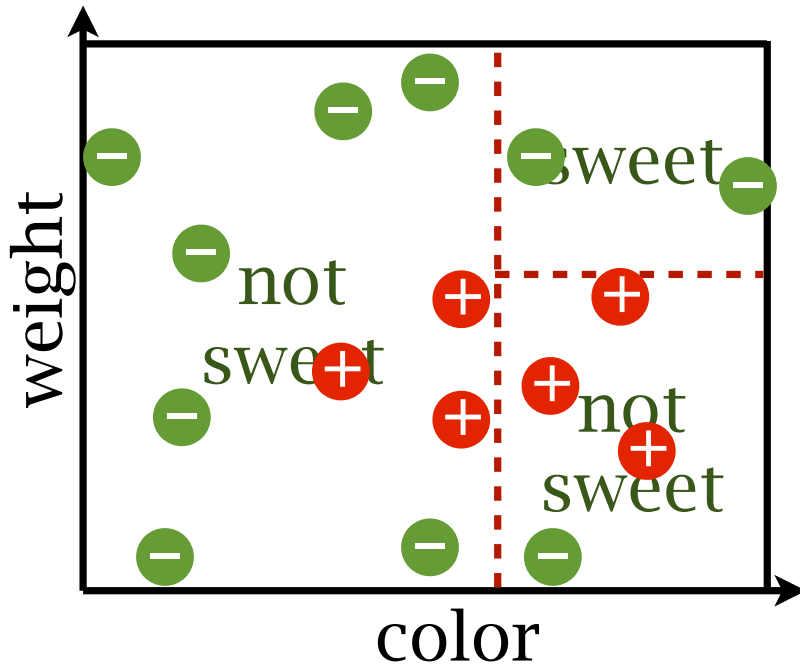
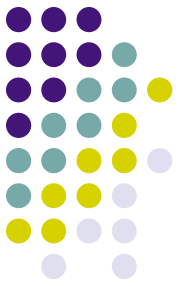
# Decision tree model



decision process  
with a tree structure



# Decision tree model

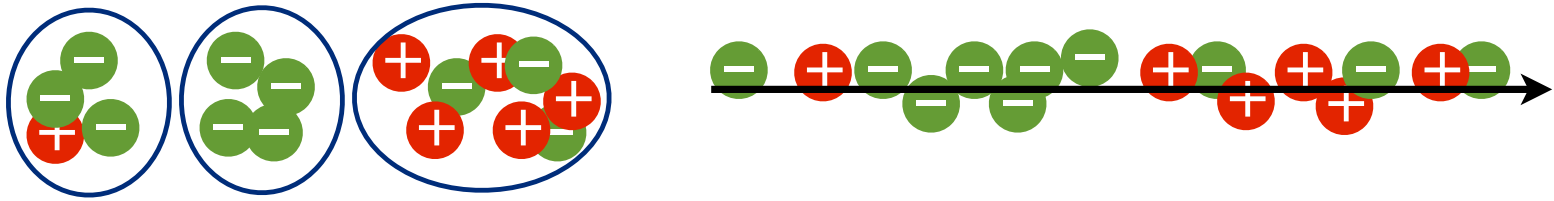


find a decision tree that matches the data?

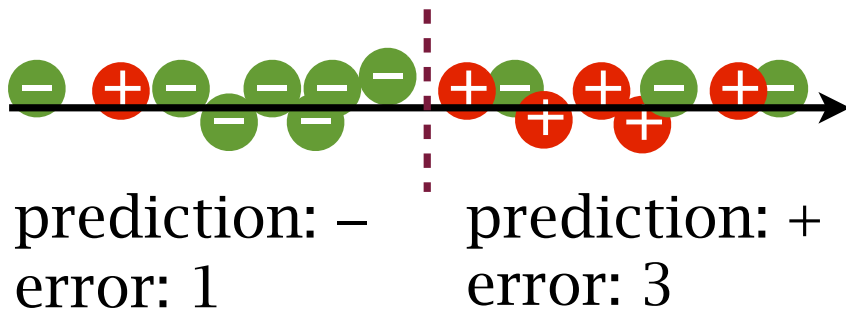
# Split-criterion: classification



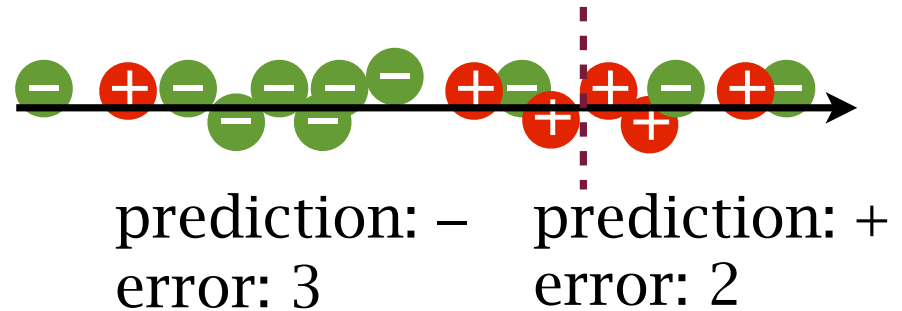
for every possible split of every feature:



## Training error:

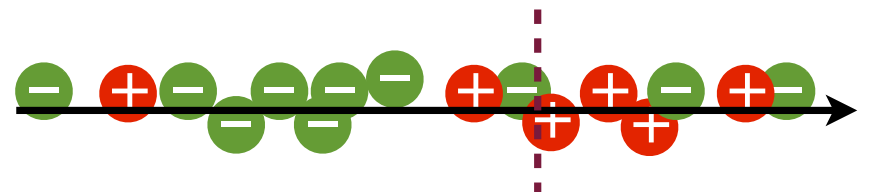


total error: 4

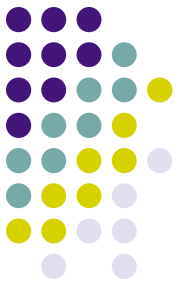


total error: 5

total error: 4



# Split-criterion: classification

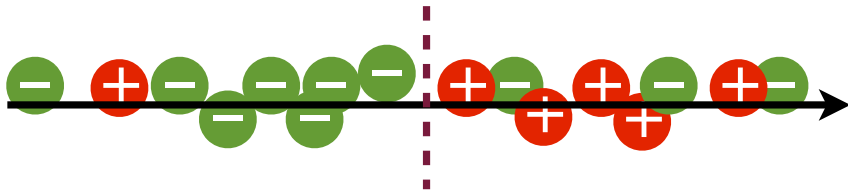


## Information gain (ID3):

$$\text{Entropy: } H(X) = - \sum p_i \ln(p_i)$$

$$\text{Entropy after split: } I(X; \text{split}) = \frac{\# \text{left}}{\# \text{all}} H(\text{left}) + \frac{\# \text{right}}{\# \text{all}} H(\text{right})$$

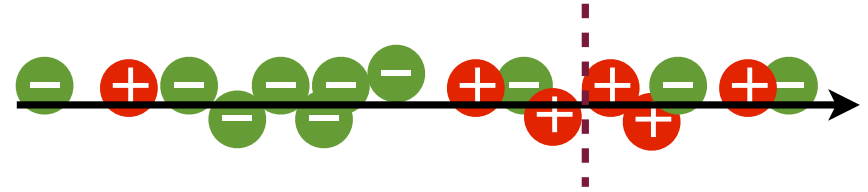
$$\text{Information gain: } H(X) - I(X; \text{split})$$



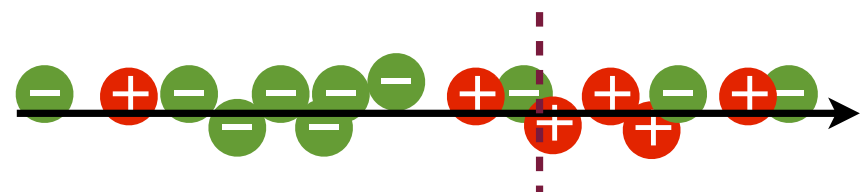
$$H(\text{left}) = -\frac{1}{8} \ln \frac{1}{8} - \frac{7}{8} \ln \frac{7}{8} = 0.3768$$

$$H(\text{right}) = -\frac{5}{8} \ln \frac{5}{8} - \frac{3}{8} \ln \frac{3}{8} = 0.6616$$

$$\begin{aligned} \text{IG} &= H(X) - (0.5 \times 0.3768 + 0.5 \times 0.6616) \\ &= H(X) - 0.5192 \end{aligned}$$



$$\text{IG} = H(X) - 0.6132$$



$$\text{IG} = H(X) - 0.5514$$

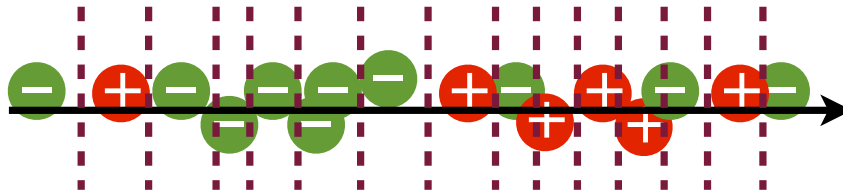


# Split-criterion: classification

## Gain ratio (C4.5):

$$\text{Gain ratio}(X) = \frac{H(X) - I(X; \text{split})}{IV(\text{split})}$$

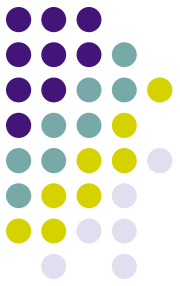
$$IV(\text{split}) = H(\text{split})$$



e.g. student ID

$$IG = H(X) - 0$$

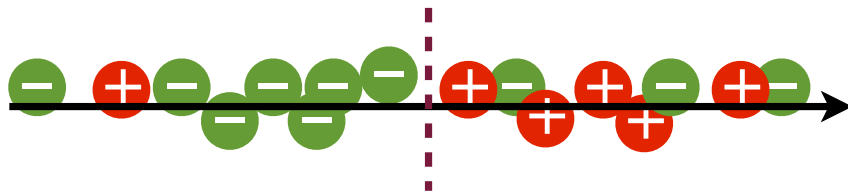
# Split-criterion: classification



## Gini index (CART):

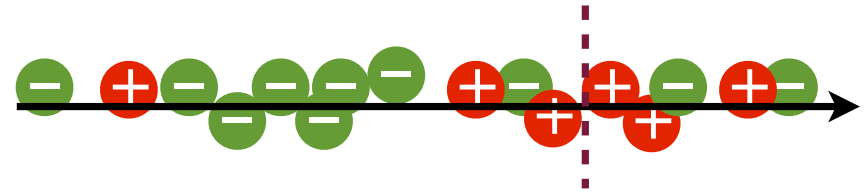
$$\text{Gini: } Gini(X) = 1 - \sum_i p_i^2$$

$$\text{Gini after split: } \frac{\#left}{\#all} Gini(left) + \frac{\#right}{\#all} Gini(right)$$



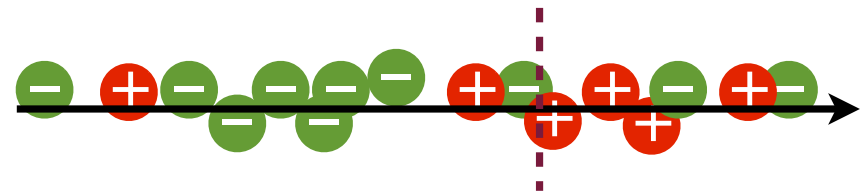
$$IG = H(X) - 0.5192$$

$$Gini = 0.3438$$



$$IG = H(X) - 0.6132$$

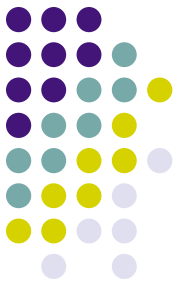
$$Gini = 0.4427$$



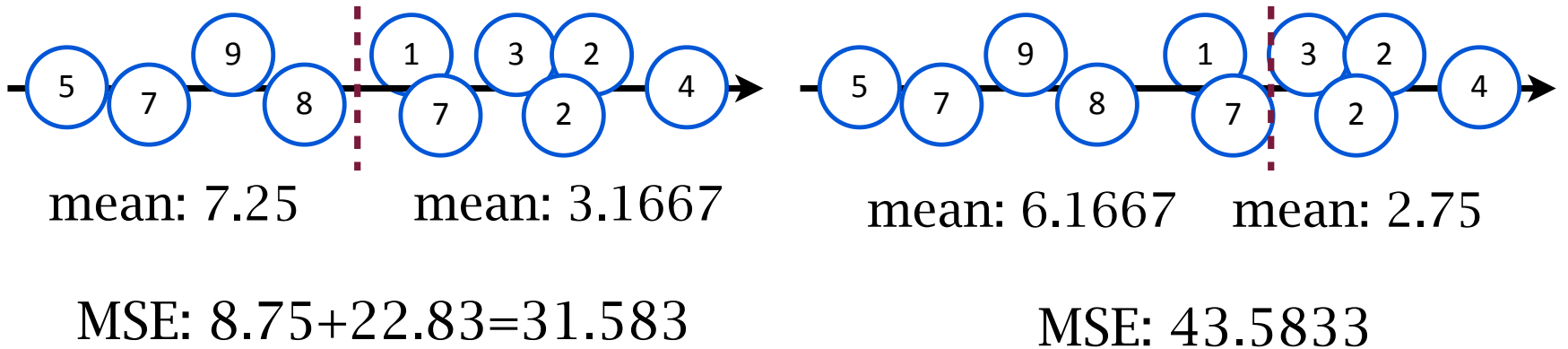
$$IG = H(X) - 0.5514$$

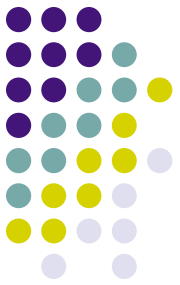
$$Gini = 0.3667$$

# Split-criterion: regression



## Training error:





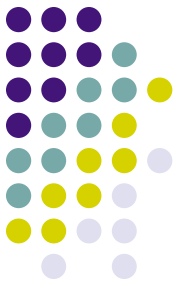
## Split-criterion: stop

Stop criterion:  
no feature to use

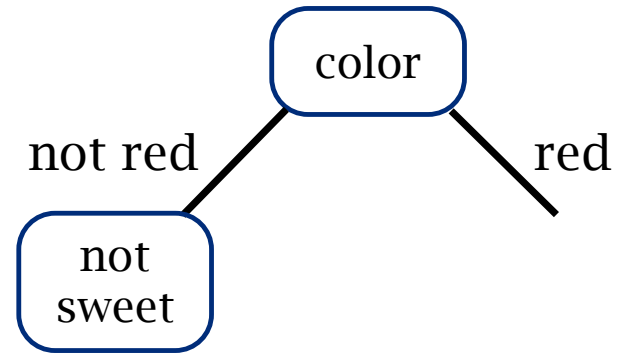
Classification: examples are pure of class

Regression: variance small enough





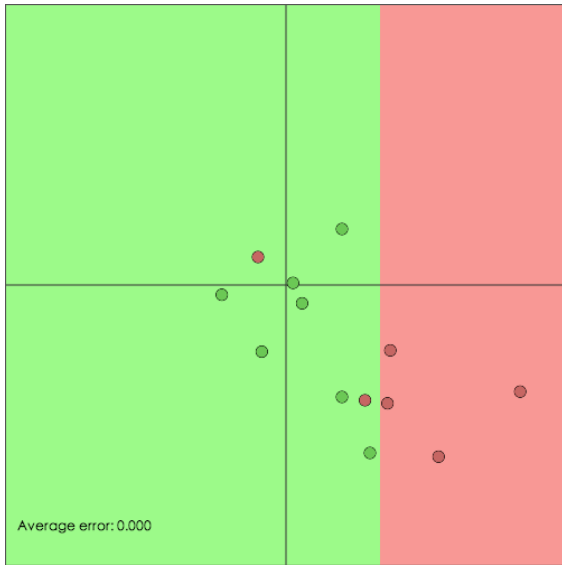
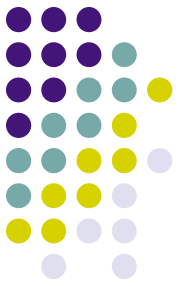
# Make-leaf



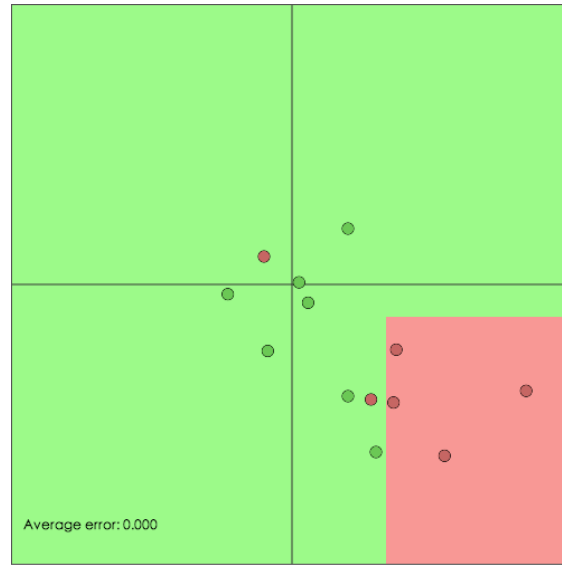
Classification: major class

Regression: mean value

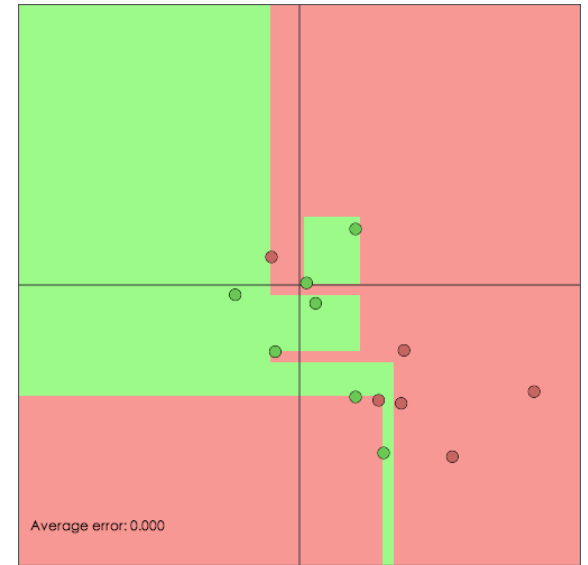
# DT boundary visualization



decision stump

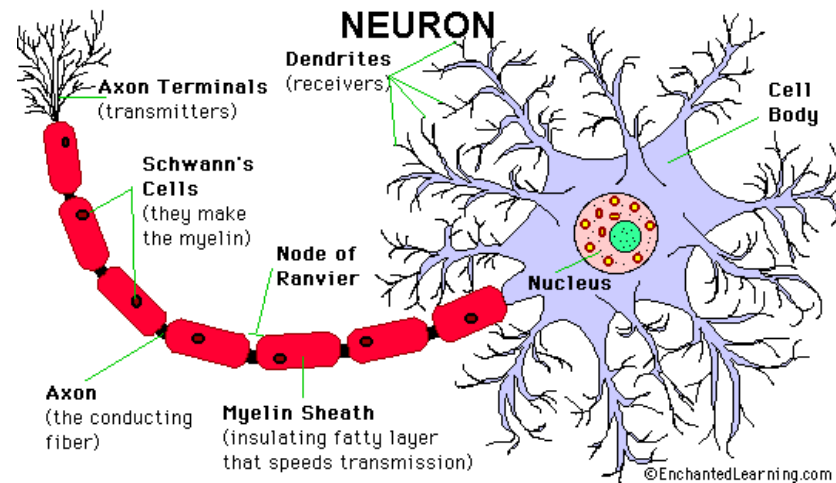
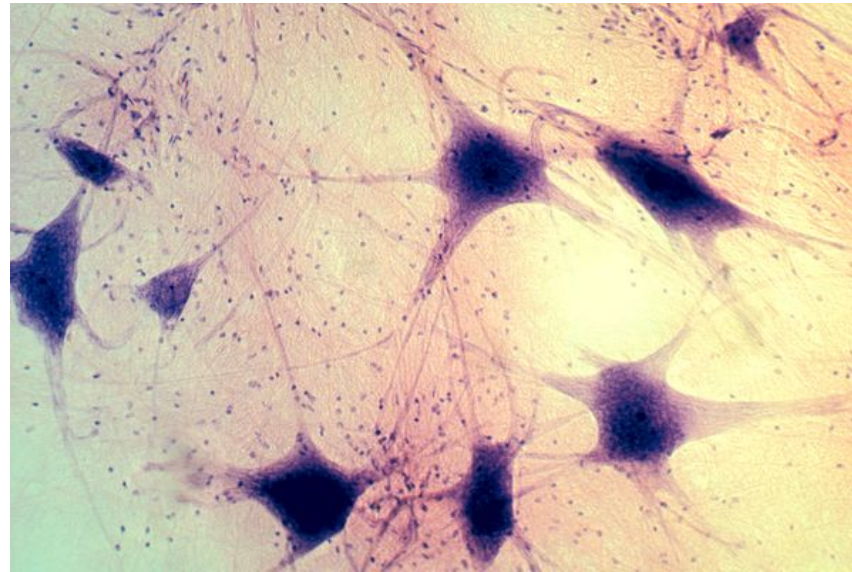
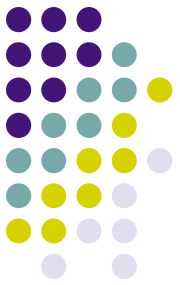


max depth=2

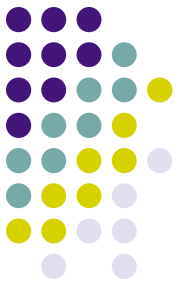


max depth=12

# Neural networks



# Neuron / perceptron



output a function of sum of input

linear function:

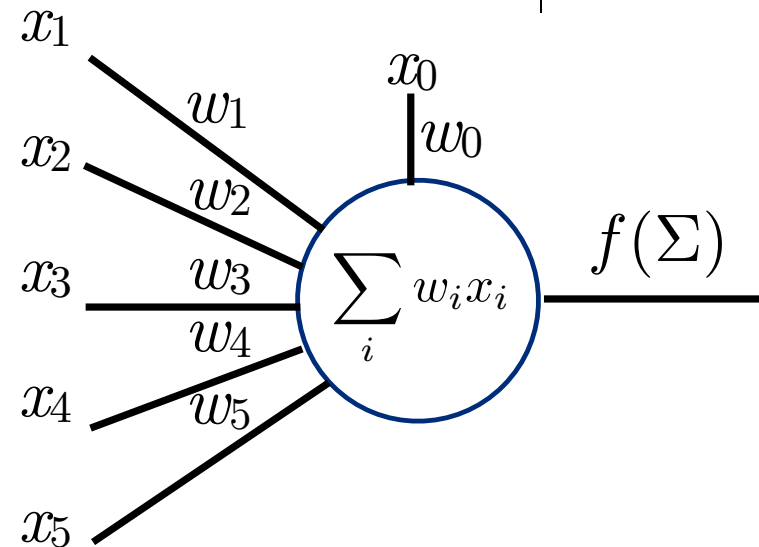
$$f\left(\sum_i w_i x_i\right) = \sum_i w_i x_i$$

threshold function:

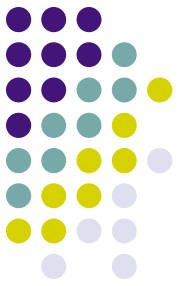
$$f\left(\sum_i w_i x_i\right) = I\left(\sum_i w_i x_i > 0\right)$$

sigmoid function:

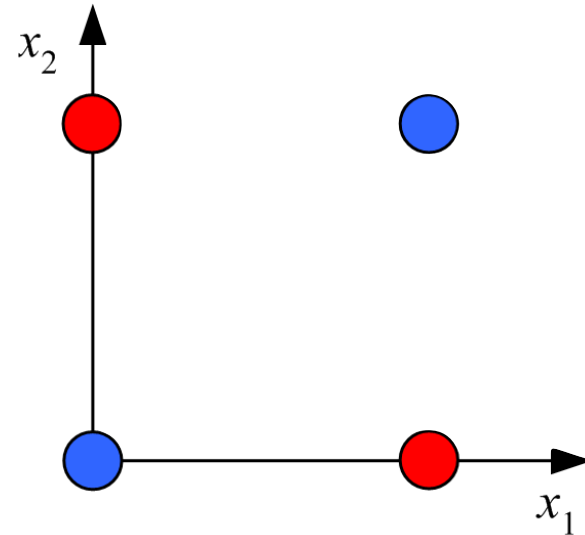
$$f\left(\sum_i w_i x_i\right) = \frac{1}{1 + e^{-\Sigma}}$$



# Limitation of single neuron



$x_1$	$x_2$	$r$
0	0	0
0	1	1
1	0	1
1	1	0



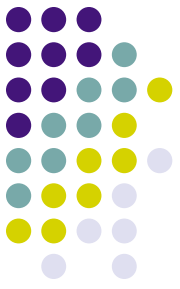
[Minsky and Papert, *Perceptrons*, 1969]



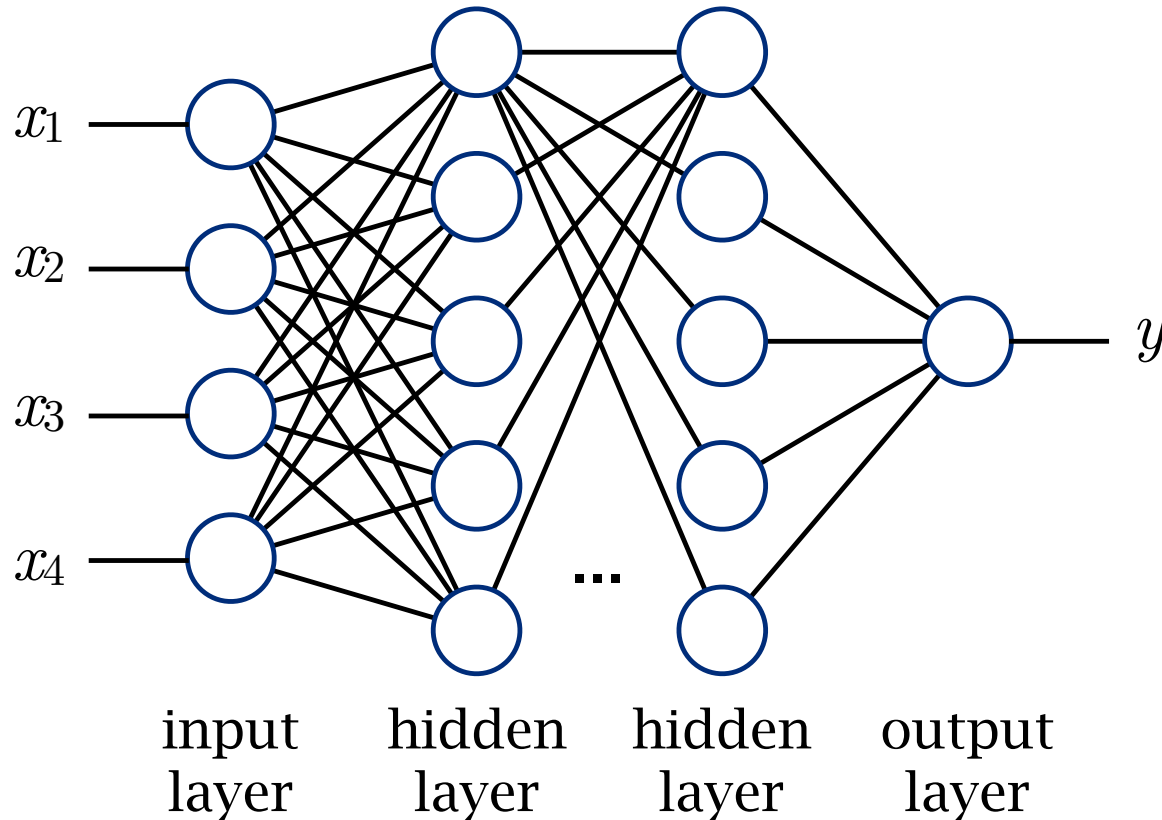
Marvin Minsky  
Turing Award 1969

AI Winter

# Multi-layer perceptrons

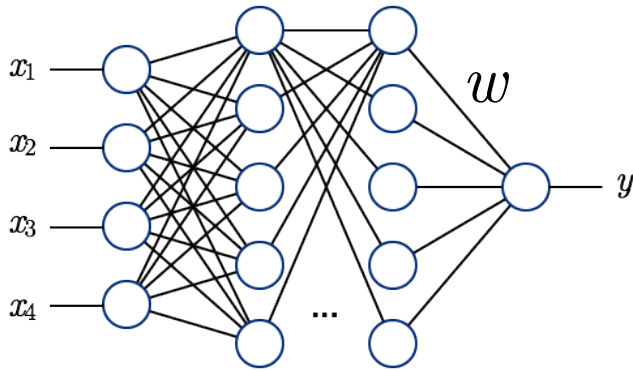
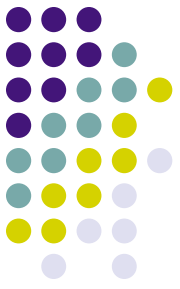


feed-forward network



sigmoid network with one hidden layer can approximate arbitrary function [Cybenko 1989]

# Back-propagation algorithm



$$\hat{y} = F(\mathbf{x})$$

**gradient descent**

$$\text{error: } E(\mathbf{w}) = (F(\mathbf{x}) - y)^2$$

$$\text{update one weight: } \Delta w_{i,j} = -\eta \frac{\partial E(\mathbf{w})}{\partial w_{i,j}}$$

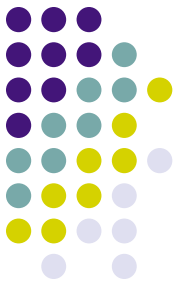
weight of the laster layer

$$\frac{\partial E(\mathbf{w})}{\partial w_{i,j}} = \frac{\partial E(\mathbf{w})}{\partial F(\mathbf{x})} \frac{\partial F(\mathbf{x})}{\partial w_{i,j}}$$

weight of the first layer

$$\frac{\partial E(\mathbf{w})}{\partial w_{i,j}} = \frac{\partial E(\mathbf{w})}{\partial F(\mathbf{x})} \frac{\partial F(\mathbf{x})}{\partial \text{HL2}} \frac{\partial \text{HL2}}{\partial \text{HL1}} \frac{\partial \text{HL1}}{\partial w_{i,j}}$$

# Back-propagation algorithm



**For each** given training example  $(\mathbf{x}, \mathbf{y})$ , **do**

1. Input the instance  $\mathbf{x}$  to the NN and compute the output value  $o_u$  of every output unit  $u$  of the network

2. **For each** network output unit  $k$ , calculate its error term  $\delta_k$

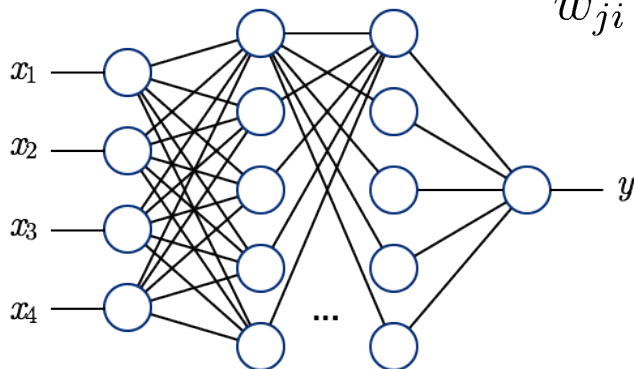
$$\delta_k \leftarrow o_k(1 - o_k)(y_k - o_k)$$

3. **For each** hidden unit  $k$ , calculate its error term  $\delta_h$

$$\delta_h \leftarrow o_k(1 - o_k) \sum_{k \in \text{outputs}} w_{kh} \delta_k$$

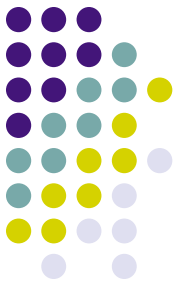
4. Update each network weight  $w_{ji}$  which is the weight associated with the  $i$ -th input value to the unit  $j$

$$w_{ji} \leftarrow w_{ji} + \eta \delta_j x_{ji}$$

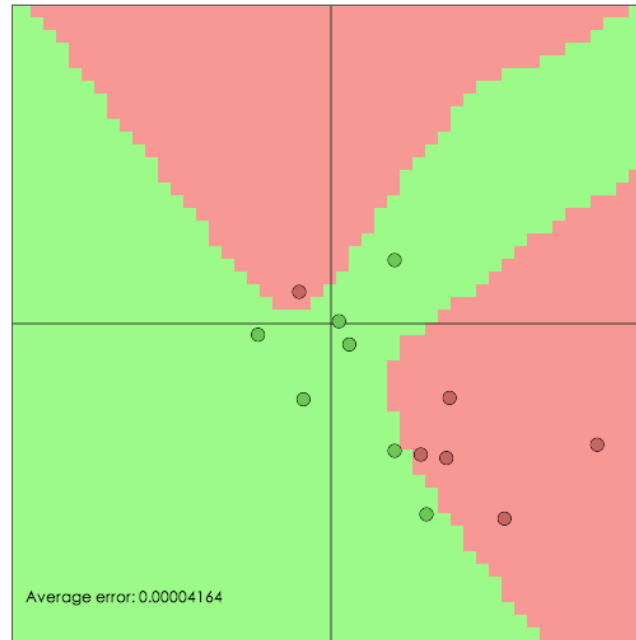




# Advantage and disadvantages



Smooth and nonlinear  
decision boundary

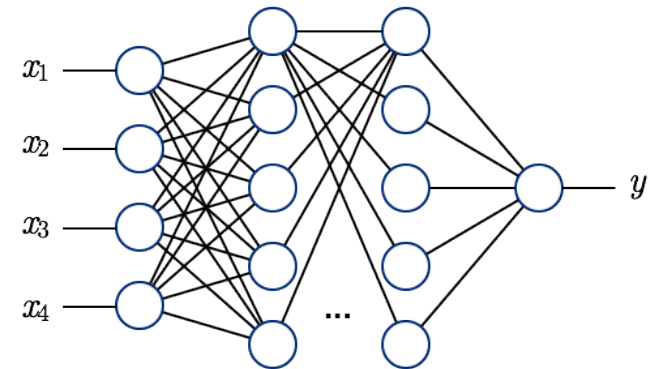


Slow convergence

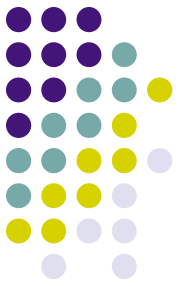
Many local optima

Best network structure unknown

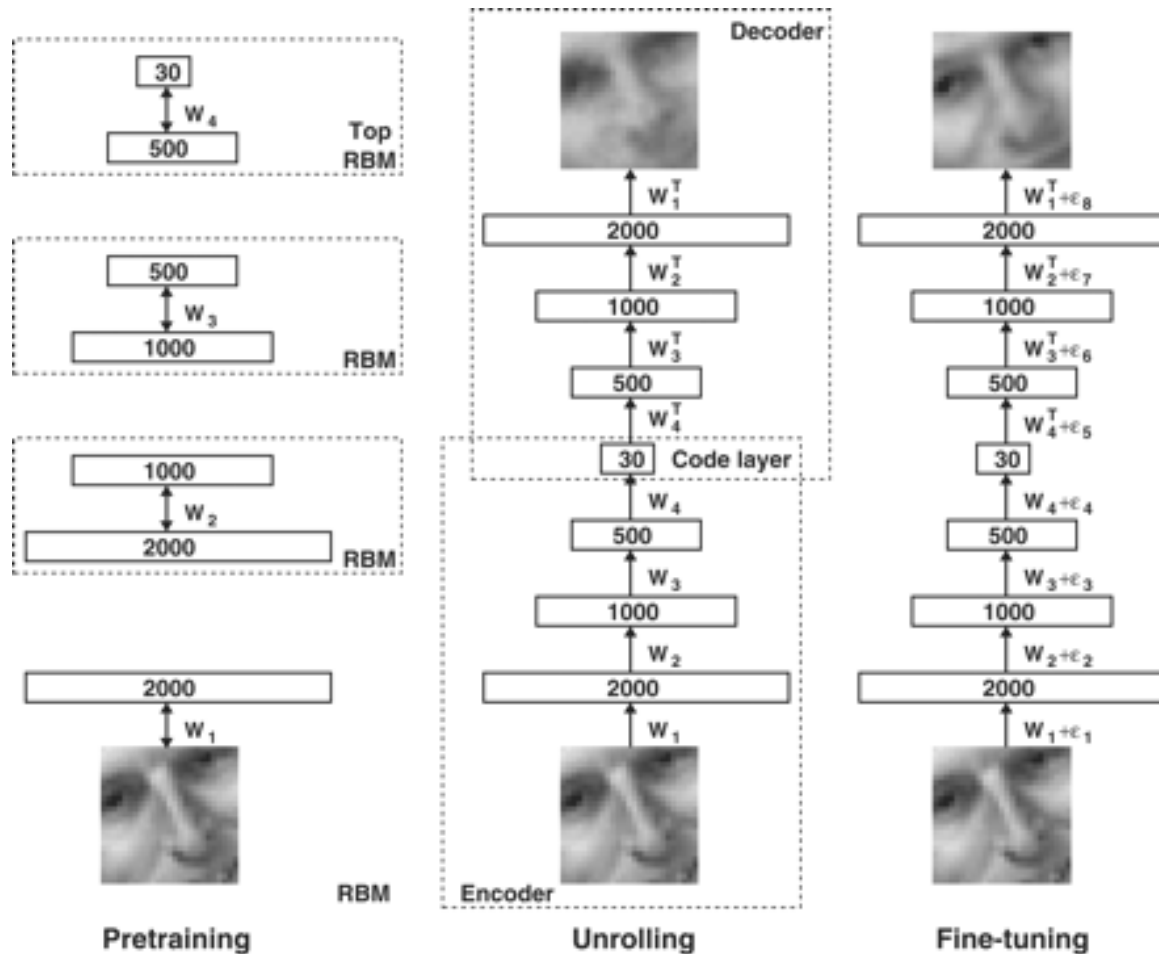
Hard to handle nominal features



# Deep network

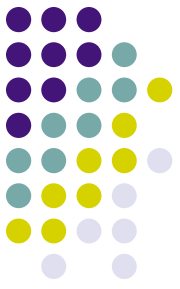


autoencoder:



[Hinton and Salakhutdinov, Science 2006]

# Bayes rule



classification using posterior probability

for binary classification

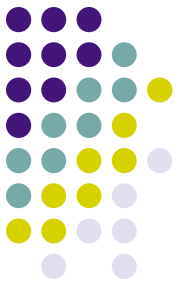
$$f(x) = \begin{cases} +1, & P(y = +1 | \mathbf{x}) > P(y = -1 | \mathbf{x}) \\ -1, & P(y = +1 | \mathbf{x}) < P(y = -1 | \mathbf{x}) \\ \text{random,} & \textit{otherwise} \end{cases}$$

in general

$$\begin{aligned} f(x) &= \arg \max_y P(y | \mathbf{x}) \\ &= \arg \max_y P(\mathbf{x} | y)P(y)/P(\mathbf{x}) \\ &= \arg \max_y P(\mathbf{x} | y)P(y) \end{aligned}$$

how the  
probabilities be  
estimated

# Naive Bayes



$$f(x) = \arg \max_y P(\mathbf{x} | y)P(y)$$

estimation the a priori by frequency:

$$P(y) \leftarrow \tilde{P}(y) = \frac{1}{m} \sum_i I(y_i = y)$$

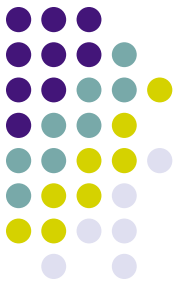
assume features are conditional independence given the class (**naive assumption**):

$$\begin{aligned} P(\mathbf{x} | y) &= P(x_1, x_2, \dots, x_n | y) \\ &= P(x_1 | y) \cdot P(x_2 | y) \cdot \dots \cdot P(x_n | y) \end{aligned}$$

**decision function:**

$$f(x) = \arg \max_y \tilde{P}(y) \prod_i \tilde{P}(x_i | y)$$

# Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

$$P(y = \text{yes}) = 2/5$$

$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

...

$$f(y \mid \text{color} = 3, \text{weight} = 3) \rightarrow$$

$$P(\text{color} = 3 \mid y = \text{yes})P(\text{weight} = 3 \mid y = \text{yes})P(y = \text{yes}) = 0.5 \times 0.5 \times 0.4 = 0.1$$

$$P(\text{color} = 3 \mid y = \text{no})P(\text{weight} = 3 \mid y = \text{no})P(y = \text{no}) = 0.33 \times 0.33 \times 0.6 = 0.06$$

$$f(y \mid \text{color} = 0, \text{weight} = 1) \rightarrow$$

$$P(\text{color} = 0 \mid y = \text{yes})P(\text{weight} = 1 \mid y = \text{yes})P(y = \text{yes}) = 0$$

$$P(\text{color} = 0 \mid y = \text{no})P(\text{weight} = 1 \mid y = \text{no})P(y = \text{no}) = 0$$

# Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

+

color	sweet?
0	yes
1	yes
2	yes
3	yes

smoothed (Laplacian correction) probabilities:

$$P(\text{color} = 0 \mid y = \text{yes}) = (0 + 1)/(2 + 4)$$

$$P(y = \text{yes}) = (2 + 1)/(5 + 2)$$

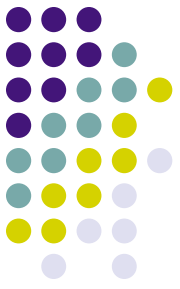
for counting frequency,  
assume every event  
has happened once.

$$f(y \mid \text{color} = 0, \text{weight} = 1) \rightarrow$$

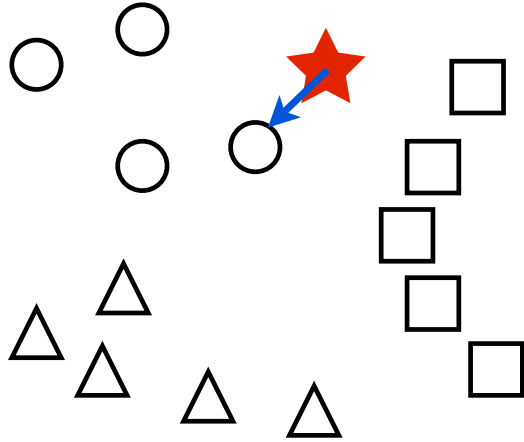
$$P(\text{color} = 0 \mid y = \text{yes})P(\text{weight} = 1 \mid y = \text{yes})P(y = \text{yes}) = \frac{1}{6} \times \frac{1}{7} \times \frac{3}{7} = 0.01$$

$$P(\text{color} = 0 \mid y = \text{no})P(\text{weight} = 1 \mid y = \text{no})P(y = \text{no}) = \frac{2}{7} \times \frac{1}{8} \times \frac{4}{7} = 0.02$$

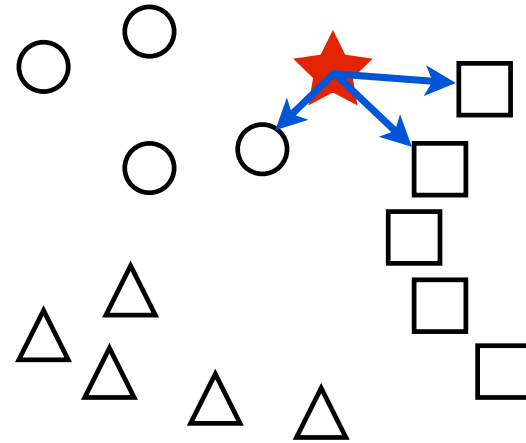
# Nearest neighbor classifier



1-nearest neighbor:

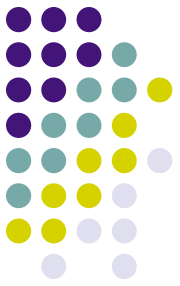


$k$ -nearest neighbor:



- ▶ asymptotically less than 2 times of the optimal Bayes error
- ▶ naturally handle multi-class
- ▶ no training time
- ▶ nonlinear decision boundary
- ▶ slow testing speed for a large training data set
- ▶ have to store the training data
- ▶ sensitive to similarity function

# Linear model

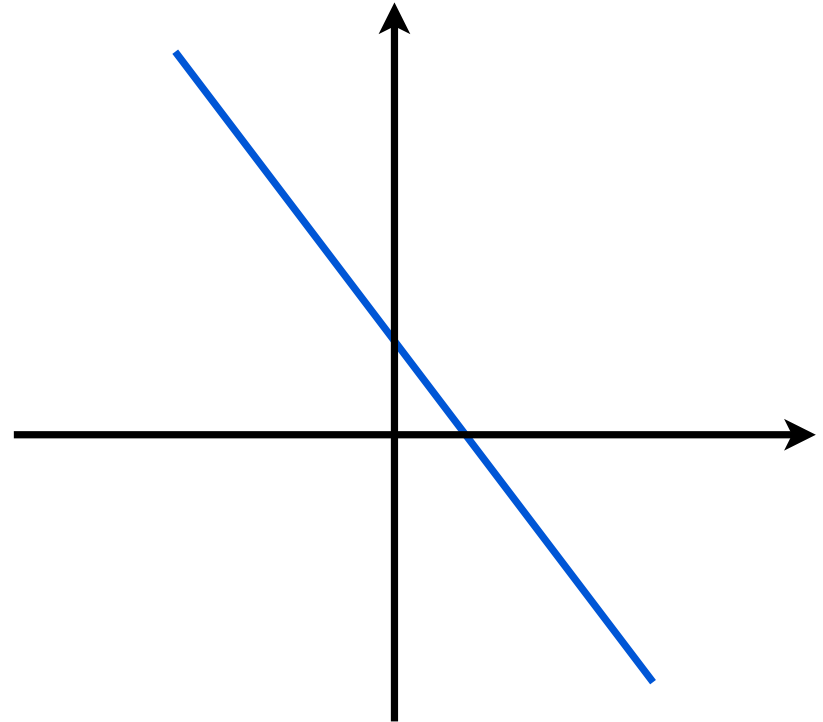


model space:  $\mathbb{R}^{n+1}$

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$$

we sometimes omit the bias

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$$



1.  $w$  with a constant element
2. practically as good as with bias (centered data)



# Least square regression



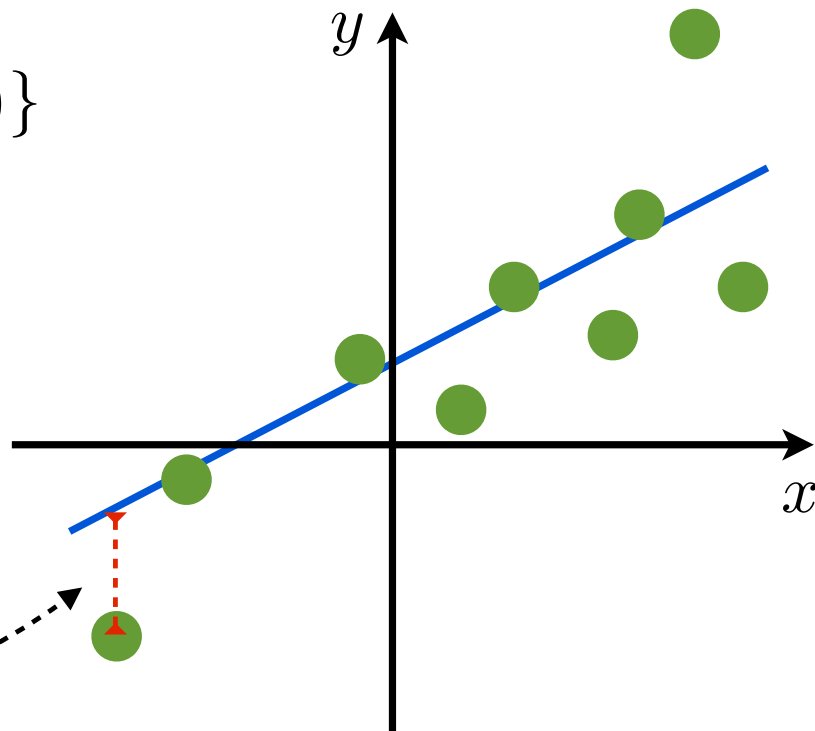
Regression:  $y \in \mathbb{R}$

Training data:

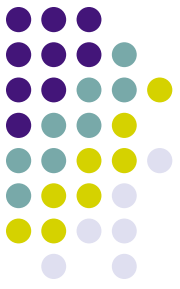
$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_m, y_m)\}$$

Least square loss:

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2$$



# Least square regression



$$L(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2$$

$$\frac{\partial L(\mathbf{w}, b)}{\partial b} = \frac{1}{m} \sum_{i=1}^m 2(\mathbf{w}^\top \mathbf{x}_i + b - y_i) = 0$$

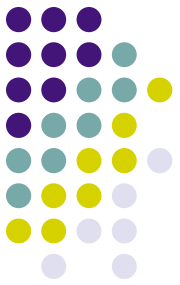
$$\frac{\partial L(\mathbf{w}, b)}{\partial \mathbf{w}} = \frac{1}{m} \sum_{i=1}^m 2(\mathbf{w}^\top \mathbf{x}_i + b - y_i) \mathbf{x}_i = 0$$

$$b = \frac{1}{m} \sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{x}_i) = \bar{y} - \mathbf{w}^\top \bar{\mathbf{x}}$$

$$\mathbf{w} = \left( \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top - \bar{\mathbf{x}} \bar{\mathbf{x}}^\top \right)^{-1} \left( \frac{1}{m} \sum_{i=1}^m (y_i \mathbf{x}_i) - \bar{y} \bar{\mathbf{x}} \right)$$

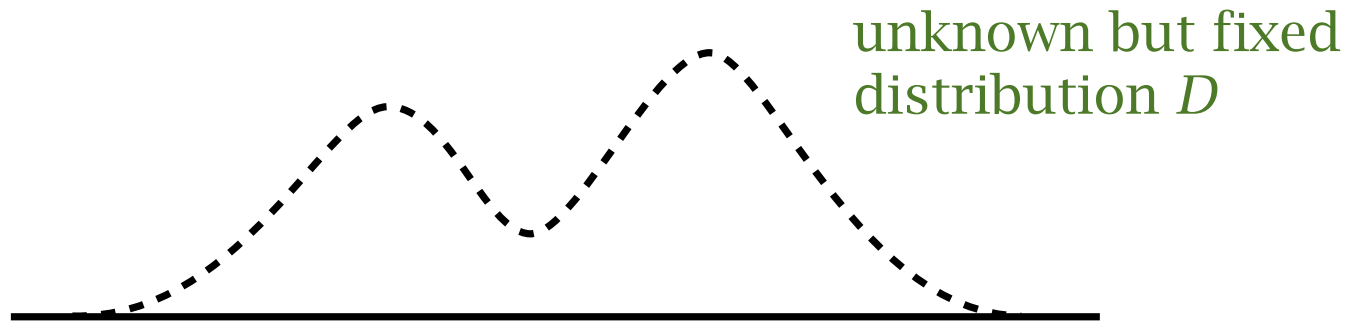
$$= \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, y) = (X^\top X)^{-1} X^\top Y$$

closed  
form  
solution

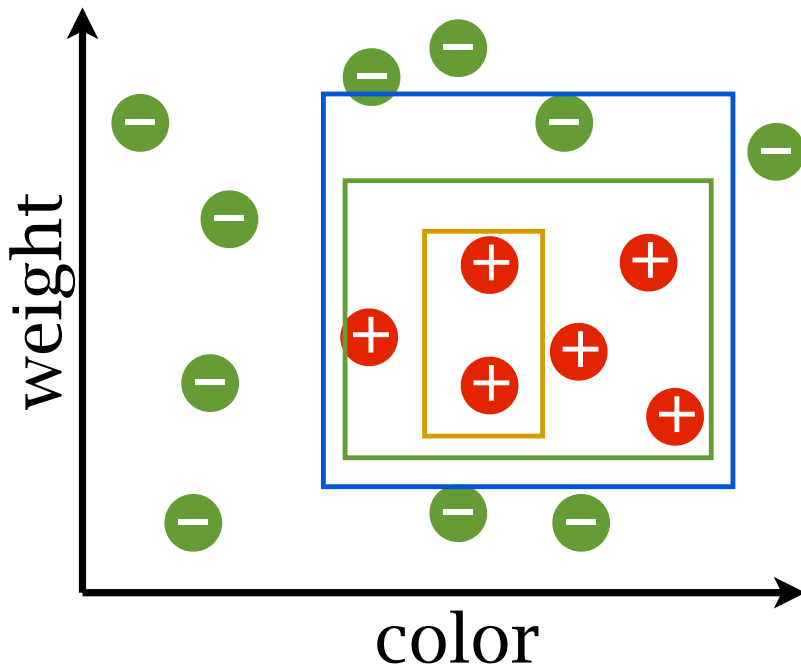
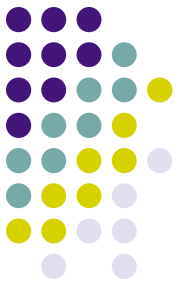


## I.I.D. assumption

all training examples and future (test) examples are drawn *independently* from an *identical distribution*



# Hypothesis class

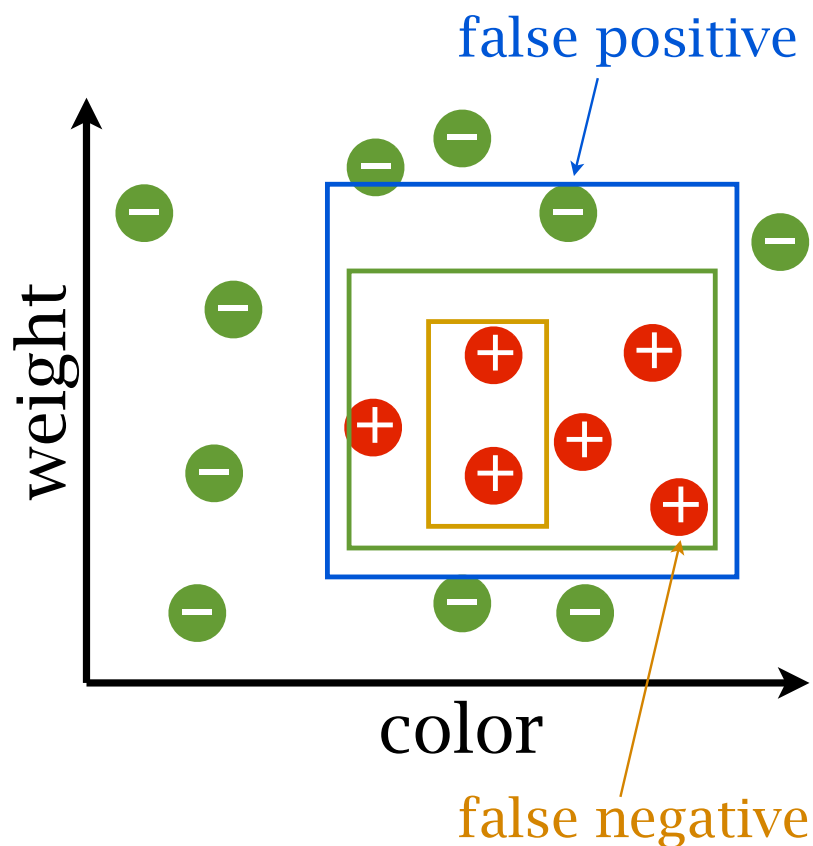


box hypothesis class  $\mathcal{H}$   
contains all boxes

$h \in \mathcal{H}$  is a hypothesis

$$h(\mathbf{x}) = \begin{cases} +1, & \text{if } x \text{ is inside the box} \\ -1, & \text{if } x \text{ is outside the box} \end{cases}$$

# Training and generalization errors



training error

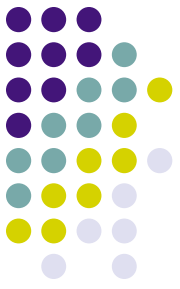
$$\epsilon_t = \frac{1}{m} \sum_{i=1}^m I(h(\mathbf{x}_i) \neq y_i)$$

generalization error

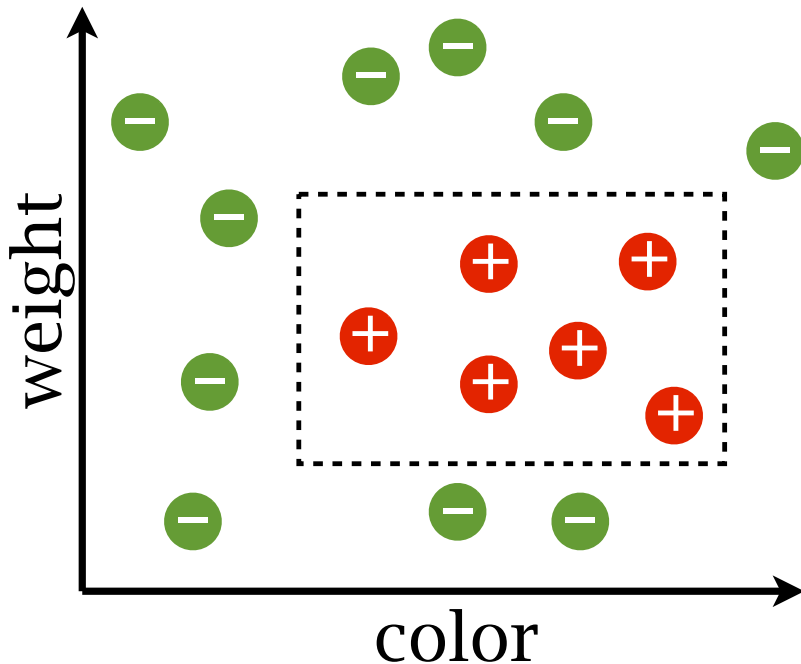
$$\begin{aligned} \epsilon_g &= \mathbb{E}_{\mathbf{x}} [I(h(\mathbf{x}) \neq f(\mathbf{x}))] \\ &= \int_{\mathcal{X}} p(\mathbf{x}) I(h(\mathbf{x}) \neq f(\mathbf{x})) d\mathbf{x} \end{aligned}$$

**find a hypothesis minimizes the generalization error**

# Generalization error



assume i.i.d. examples, and the ground-truth hypothesis is a box



the error of picking a consistent hypothesis:

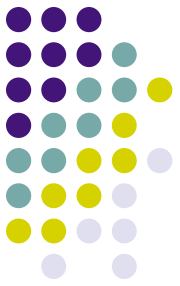
with probability at least  $1 - \delta$

$$\epsilon_g < \frac{1}{m} \cdot (\ln |\mathcal{H}| + \ln \frac{1}{\delta})$$

smaller generalization error:

- ▶ more examples
- ▶ smaller hypothesis space

# Generalization error



for one  $h$

What is the probability of  $h$  is consistent  
 $\epsilon_g(h) \geq \epsilon$

assume  $h$  is **bad**:  $\epsilon_g(h) \geq \epsilon$

$h$  is consistent with 1 example:

$$P \leq 1 - \epsilon$$

$h$  is consistent with  $m$  example:

$$P \leq (1 - \epsilon)^m$$

# Generalization error

$h$  is consistent with  $m$  example:

$$P \leq (1 - \epsilon)^m$$

There are  $k$  consistent hypotheses

Probability of choosing a bad one:

$h_1$  is chosen and  $h_1$  is bad  $P \leq (1 - \epsilon)^m$

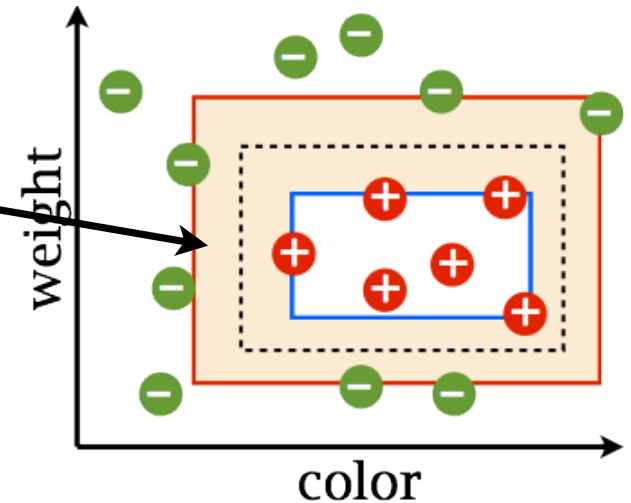
$h_2$  is chosen and  $h_2$  is bad  $P \leq (1 - \epsilon)^m$

...

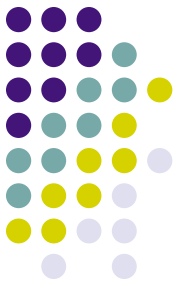
$h_k$  is chosen and  $h_k$  is bad  $P \leq (1 - \epsilon)^m$

overall:

$\exists h$ :  $h$  can be chosen (consistent) but is bad







# Generalization error

$h_1$  is chosen and  $h_1$  is bad  $P \leq (1 - \epsilon)^m$

$h_2$  is chosen and  $h_2$  is bad  $P \leq (1 - \epsilon)^m$

...

$h_k$  is chosen and  $h_k$  is bad  $P \leq (1 - \epsilon)^m$

overall:

$\exists h$ :  $h$  can be chosen (consistent) but is bad

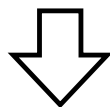
Union bound:  $P(A \cup B) \leq P(A) + P(B)$

$$P(\exists h \text{ is consistent but bad}) \leq k \cdot (1 - \epsilon)^m \leq |\mathcal{H}| \cdot (1 - \epsilon)^m$$

# Generalization error



$$P(\exists h \text{ is consistent but bad}) \leq k \cdot (1 - \epsilon)^m \leq |\mathcal{H}| \cdot (1 - \epsilon)^m$$



$$P(\epsilon_g \geq \epsilon) \leq \frac{|\mathcal{H}| \cdot (1 - \epsilon)^m}{\delta}$$

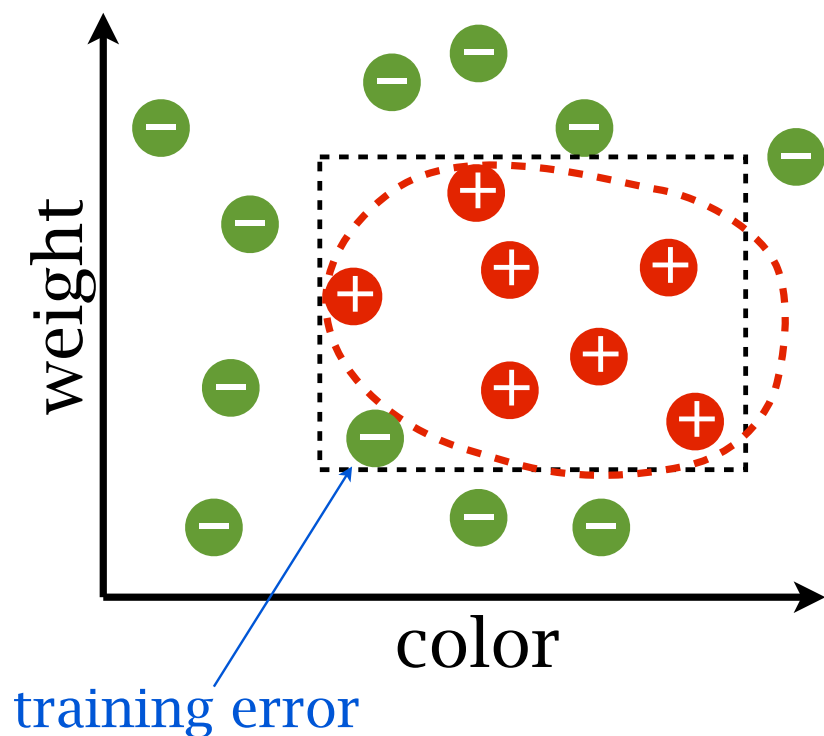
with probability at least  $1 - \delta$

$$\epsilon_g < \frac{1}{m} \cdot \left( \ln |\mathcal{H}| + \ln \frac{1}{\delta} \right)$$

# Inconsistent hypothesis



What if the ground-truth hypothesis is NOT a box: **non-zero training error**



with probability at least  $1 - \delta$

$$\epsilon_g < \epsilon_t + \sqrt{\frac{1}{m} (\ln |\mathcal{H}| + \ln \frac{1}{\delta})}$$

smaller generalization error:

- ▶ more examples
- ▶ smaller hypothesis space
- ▶ **smaller training error**

# 机器学习/模式识别基础



- 预测与识别
- 预测算法
- 特征提取

# Feature extraction



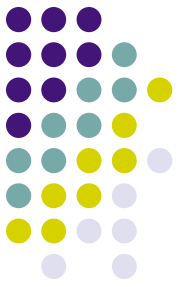
disclosure the inner structure of the data  
to support a better mining performance

feature extraction construct new features  
commonly followed by a feature selection  
usually used for low-level features

digits bitmap:

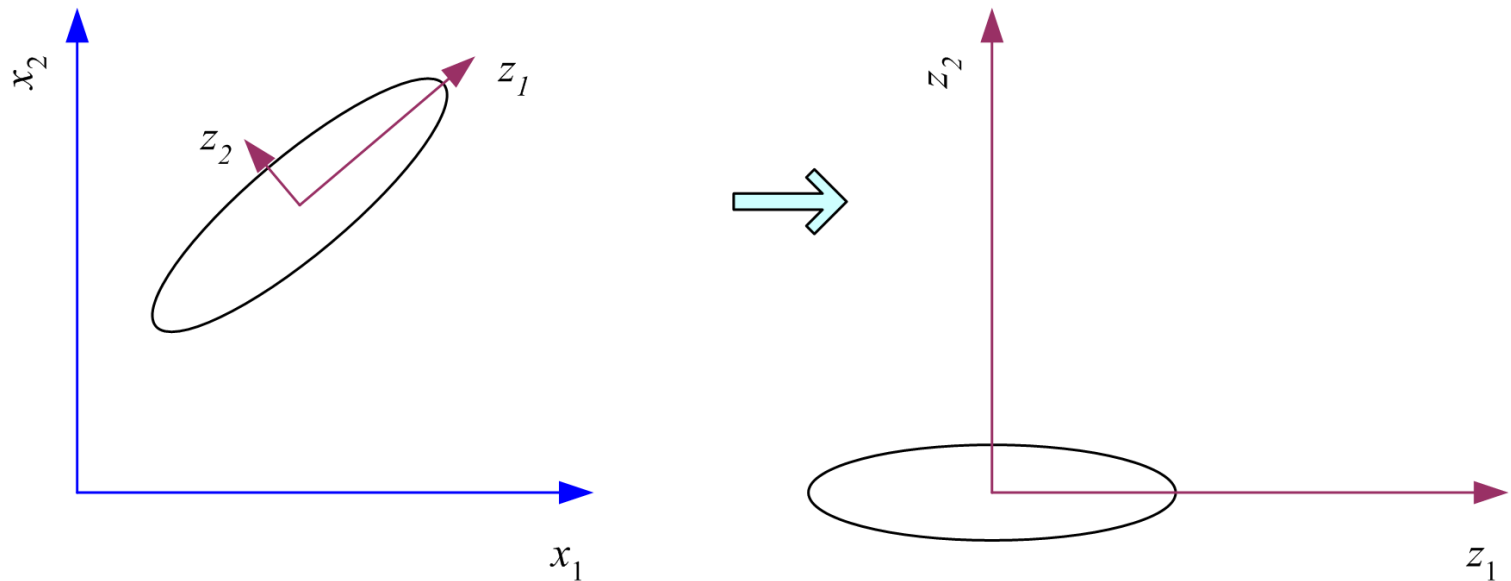
1	2	3	4	0	4
2	3	4	0	1	2
3	4	0	0	4	1
3	1	0	0	2	2
2	0	1	2	3	3
3	3	4	4	1	0

# Linear methods

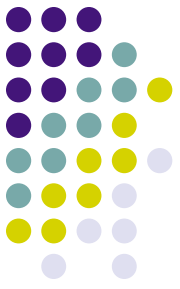


## Principal components analysis (PCA)

rotate the data to align the directions of the variance



# Linear methods



## Principal components analysis (PCA)

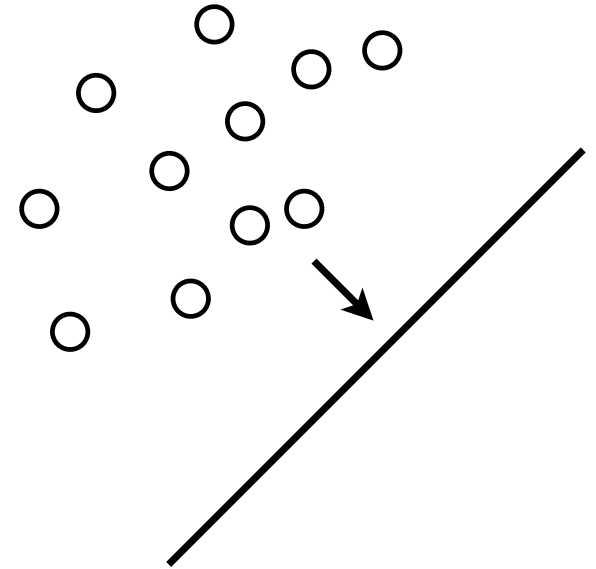
the first dimension = the largest variance direction

$$z = \mathbf{w}^T \mathbf{x}$$

$$\text{Var}(z_1) = \mathbf{w}_1^T \Sigma \mathbf{w}_1$$

find a unit  $\mathbf{w}$  to maximize the variance

$$\max_{\mathbf{w}_1} \mathbf{w}_1^T \Sigma \mathbf{w}_1 - \alpha (\mathbf{w}_1^T \mathbf{w}_1 - 1)$$

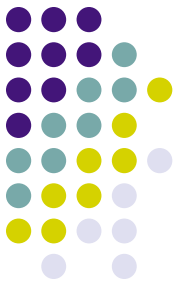


$2\Sigma\mathbf{w}_1 - 2\alpha\mathbf{w}_1 = 0$ , and therefore  $\Sigma\mathbf{w}_1 = \alpha\mathbf{w}_1$

$$\mathbf{w}_1^T \Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1^T \mathbf{w}_1 = \alpha$$

$\mathbf{w}$  is the eigenvector with the largest eigenvalue

# Linear methods



## Principal components analysis (PCA)

the second dimension = the largest variance direction orthogonal to the first dimension

$$\max_{\mathbf{w}_2} \mathbf{w}_2^T \Sigma \mathbf{w}_2 - \alpha (\mathbf{w}_2^T \mathbf{w}_2 - 1) - \beta (\mathbf{w}_2^T \mathbf{w}_1 - 0)$$

$$2\Sigma \mathbf{w}_2 - 2\alpha \mathbf{w}_2 - \beta \mathbf{w}_1 = 0$$

$$\beta = 0 \quad \Sigma \mathbf{w}_2 = \alpha \mathbf{w}_2$$

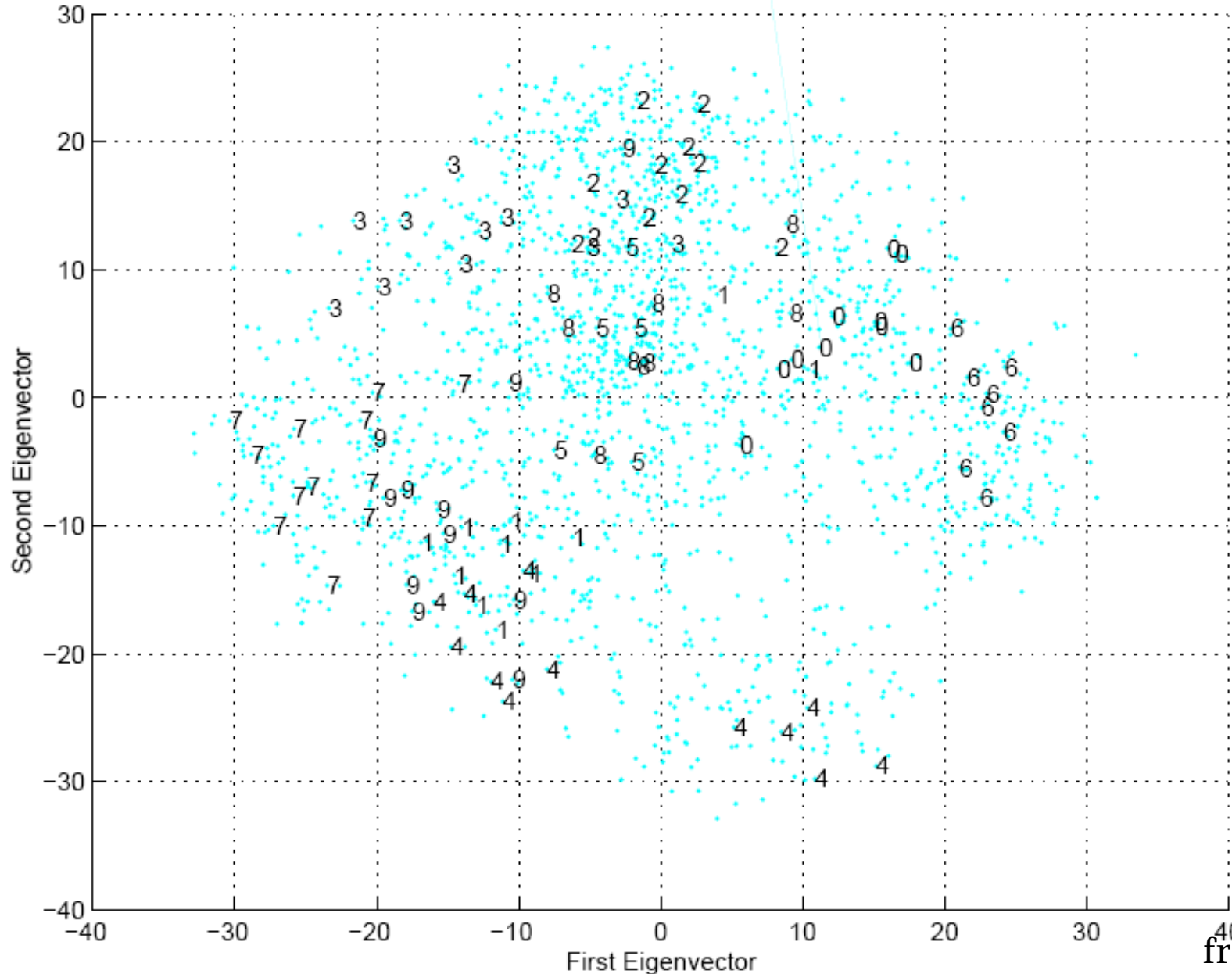
$\mathbf{w}$ 's are the eigenvectors sorted by the eigenvalues



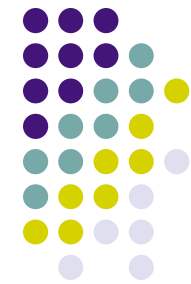
# Linear methods



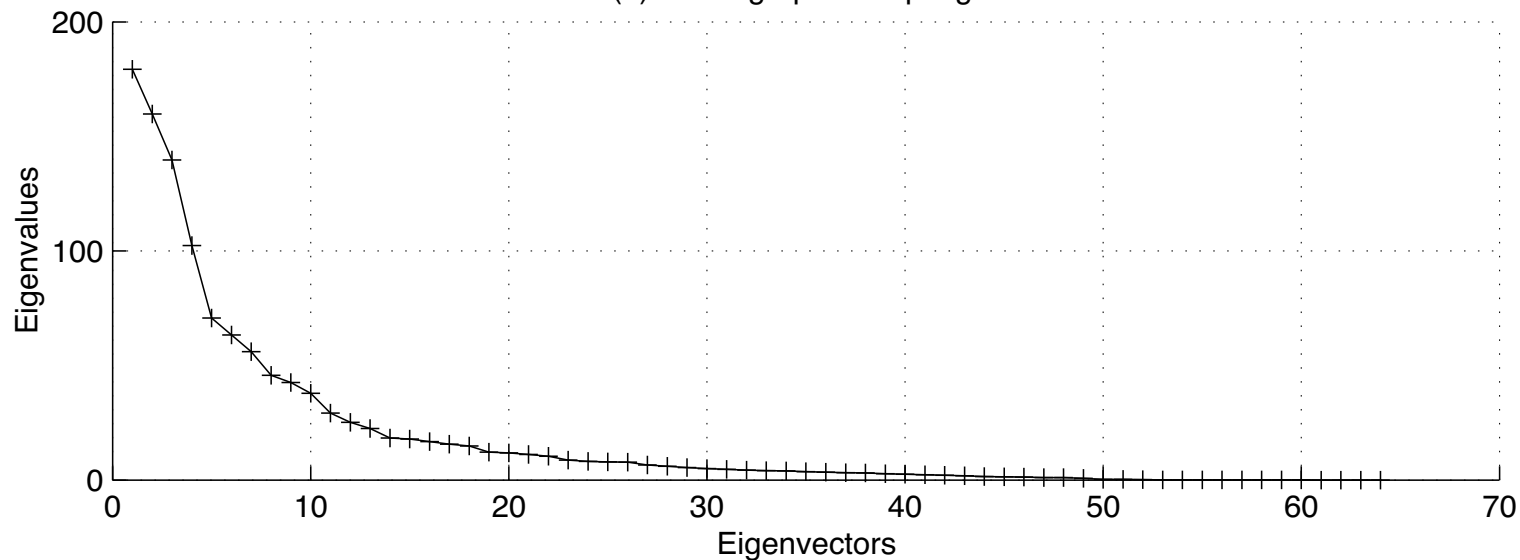
Optdigits after PCA



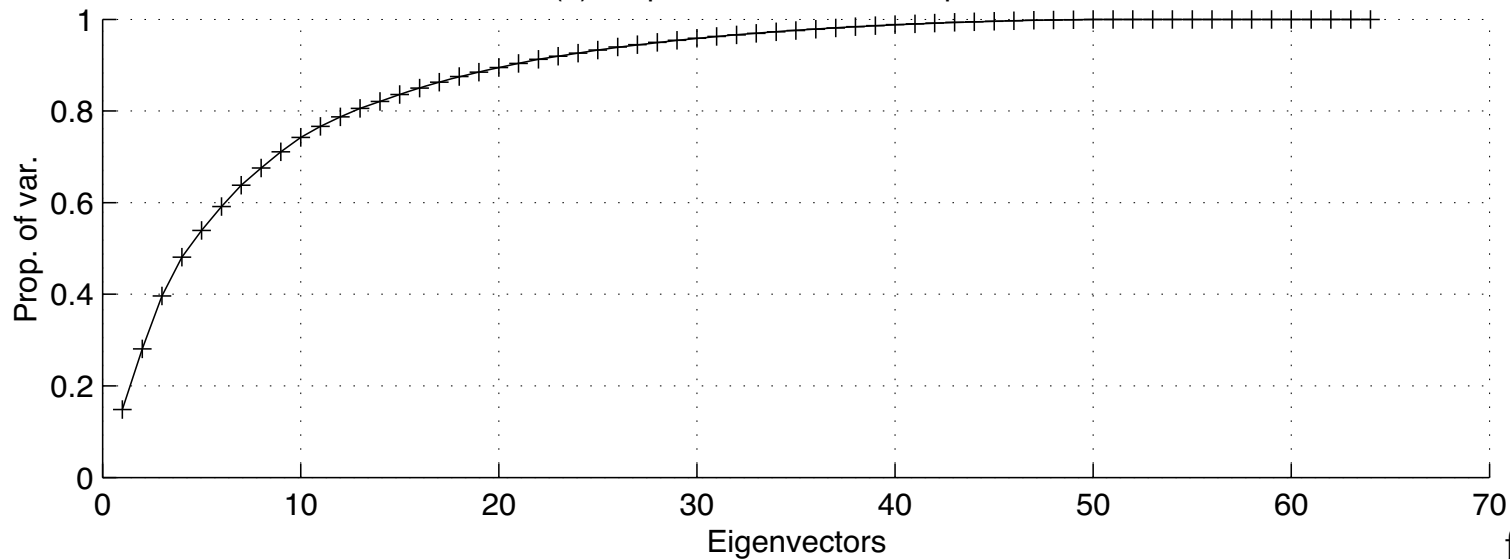
# Linear methods

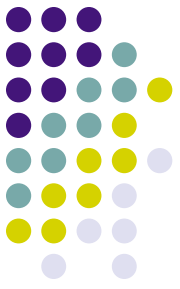


(a) Scree graph for Optdigits



(b) Proportion of variance explained





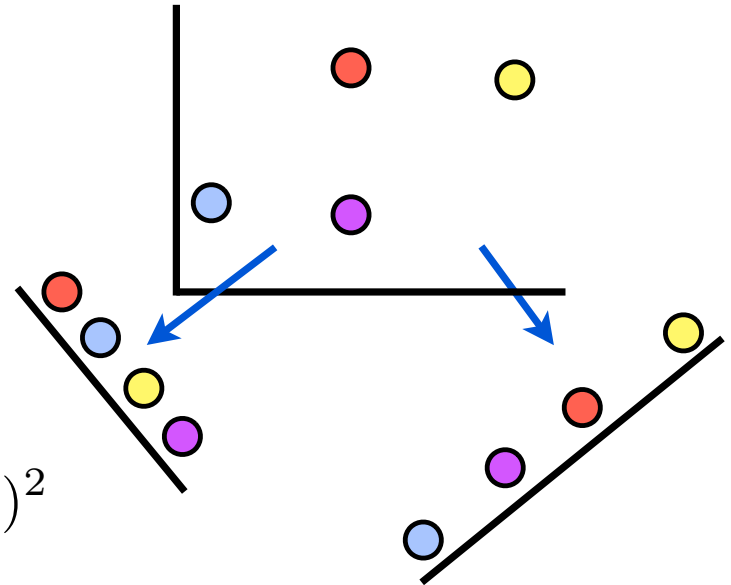
# Linear methods

## Multidimensional Scaling (MDS)

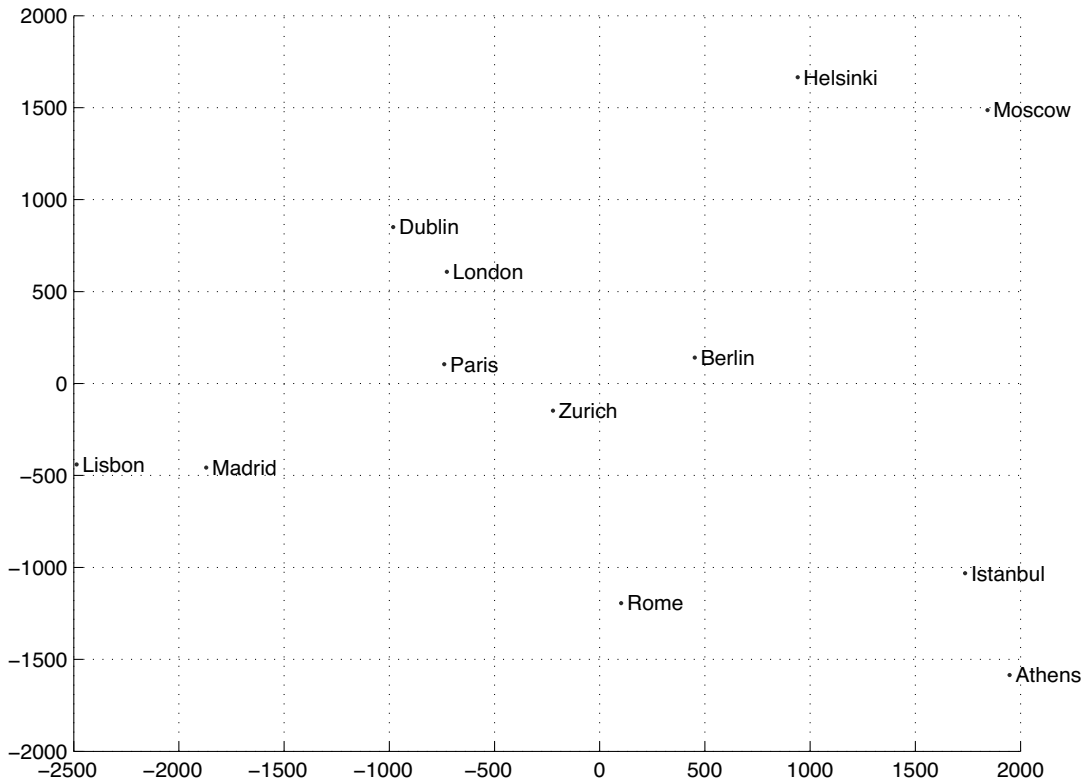
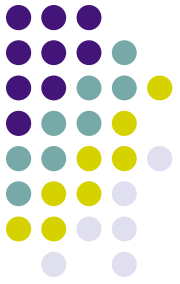
keep the distance into a lower dimensional space

for linear transformation,  
 $W$  is an  $n \times k$  matrix

$$\arg \min_W \sum_{i,j} (\|\mathbf{x}_i^\top W - \mathbf{x}_j^\top W\| - \|\mathbf{x}_i - \mathbf{x}_j\|)^2$$

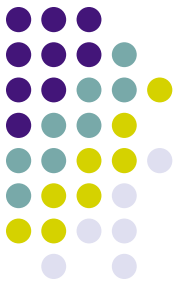


# Linear methods



from [Intro. ML]

# Linear methods



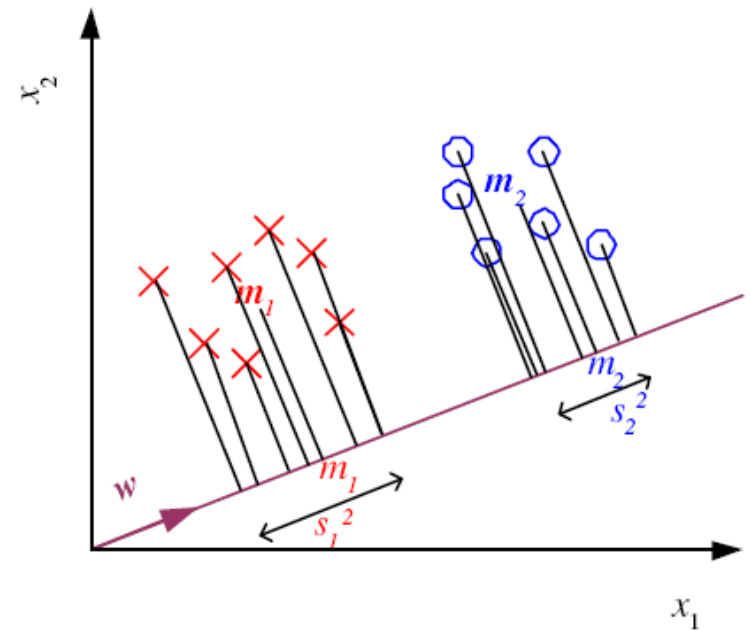
## Linear Discriminant Analysis (LDA)

find a direction such that the two classes are well separated

$$z = \mathbf{w}^T \mathbf{x}$$

$m$  be the mean of a class

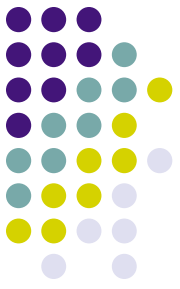
$s^2$  be the variance of a class



maximize the criterion

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

# Linear methods



## Linear Discriminant Analysis (LDA)

$$\begin{aligned}(m_1 - m_2)^2 &= (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 \\ &= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} \\ &= \mathbf{w}^T \mathbf{S}_B \mathbf{w}\end{aligned}$$

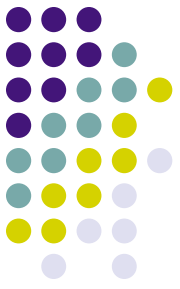
$$\begin{aligned}s_1^2 &= \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t \\ &= \sum_t \mathbf{w}^T (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T \mathbf{w} r^t \\ &= \mathbf{w}^T \mathbf{S}_1 \mathbf{w}\end{aligned}$$

$$s_1^2 + s_2^2 = \mathbf{w}^T \mathbf{S}_W \mathbf{w} \quad \mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$

The objective becomes:

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} = \frac{|\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)|^2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

# Linear methods



## Linear Discriminant Analysis (LDA)

The objective becomes:

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} = \frac{|\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)|^2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

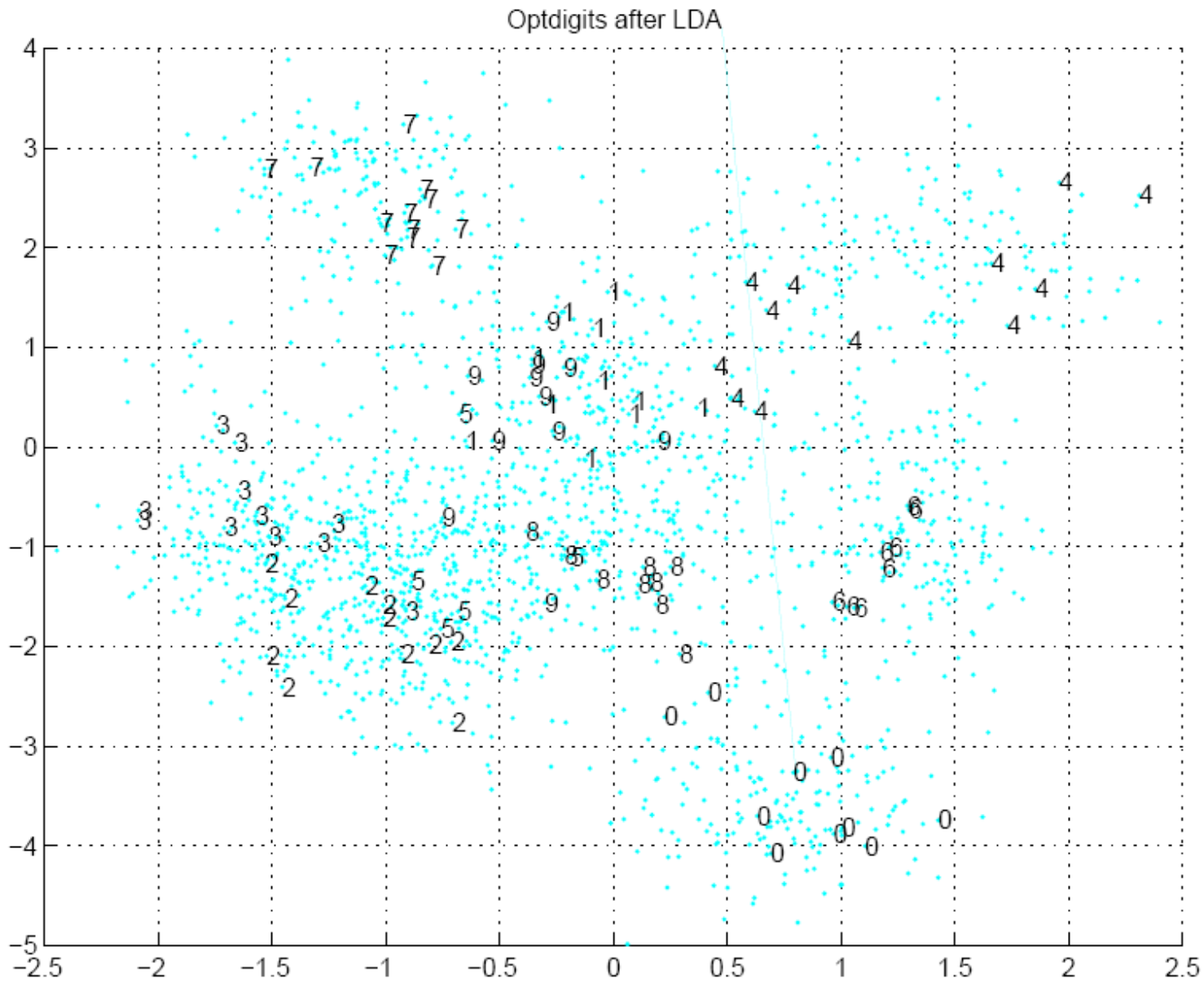
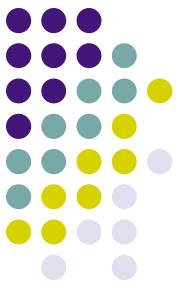
$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \left( 2(\mathbf{m}_1 - \mathbf{m}_2) - \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \mathbf{S}_W \mathbf{w} \right) = 0$$

Given that  $\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) / \mathbf{w}^T \mathbf{S}_W \mathbf{w}$  is a constant, we have

$$\mathbf{w} = c \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

just take  $c = 1$  and find  $\mathbf{w}$

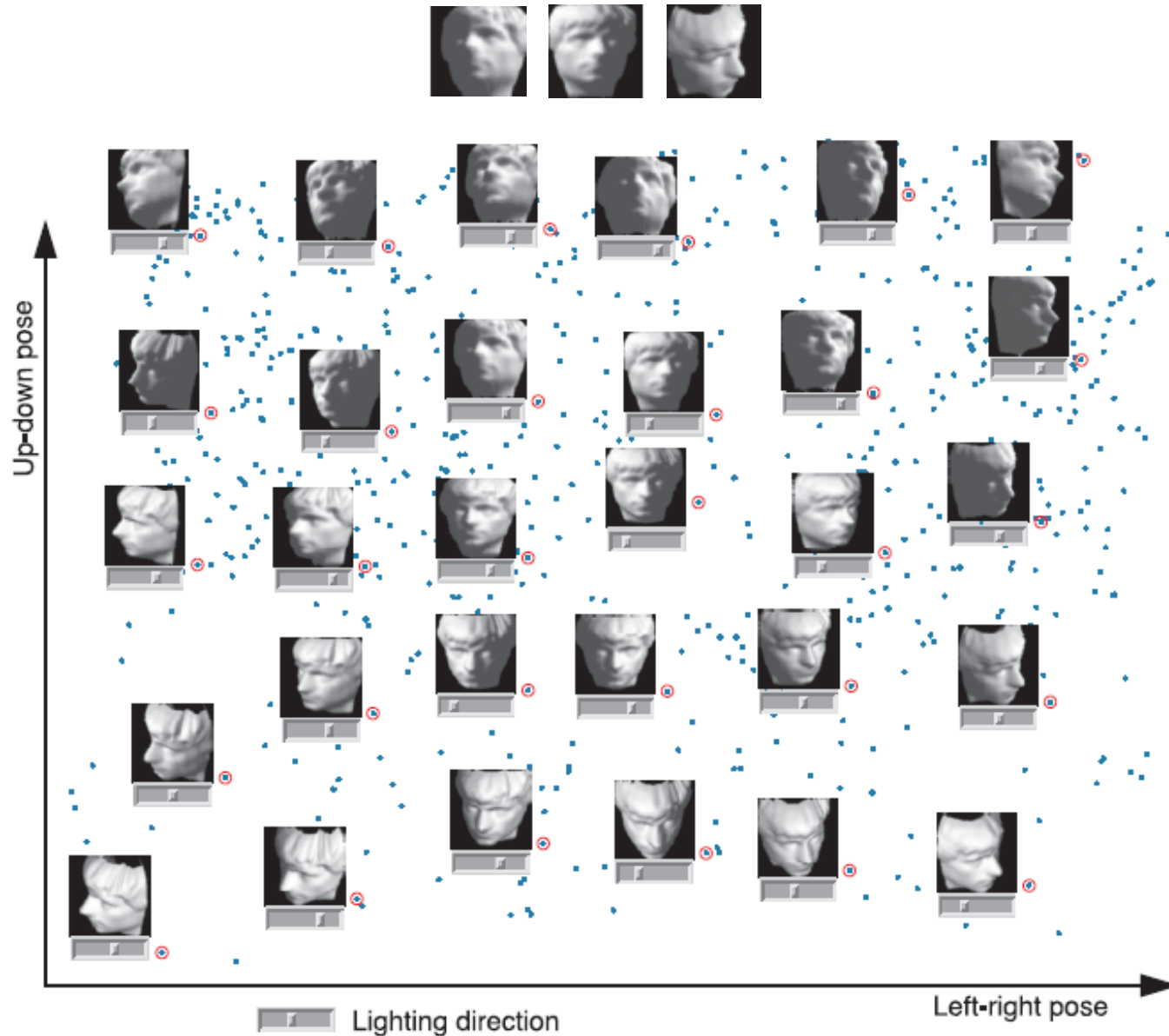
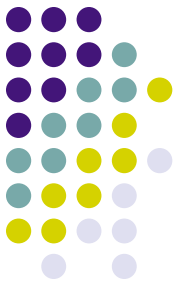
# Linear methods



from [Intro. ML]



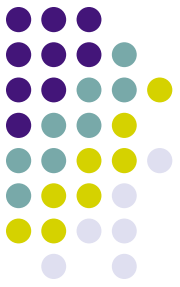
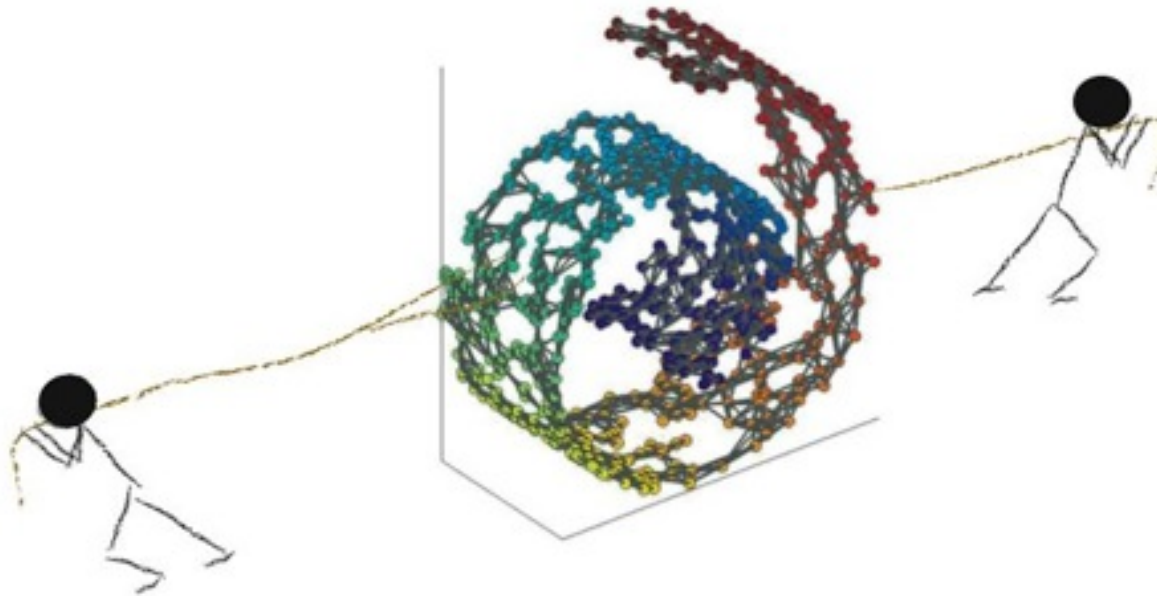
# Manifold learning



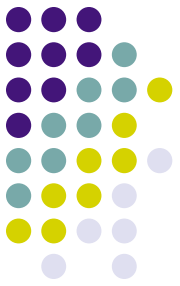
# Manifold learning

A low intrinsic dimensional data embedded in a high dimensional space

cause a bad distance measure

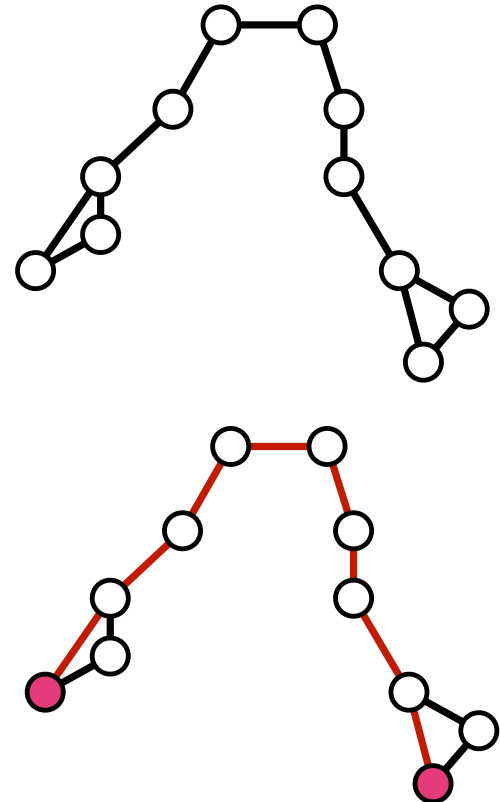


# Manifold learning

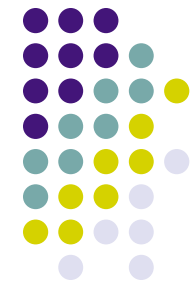


## ISOMAP

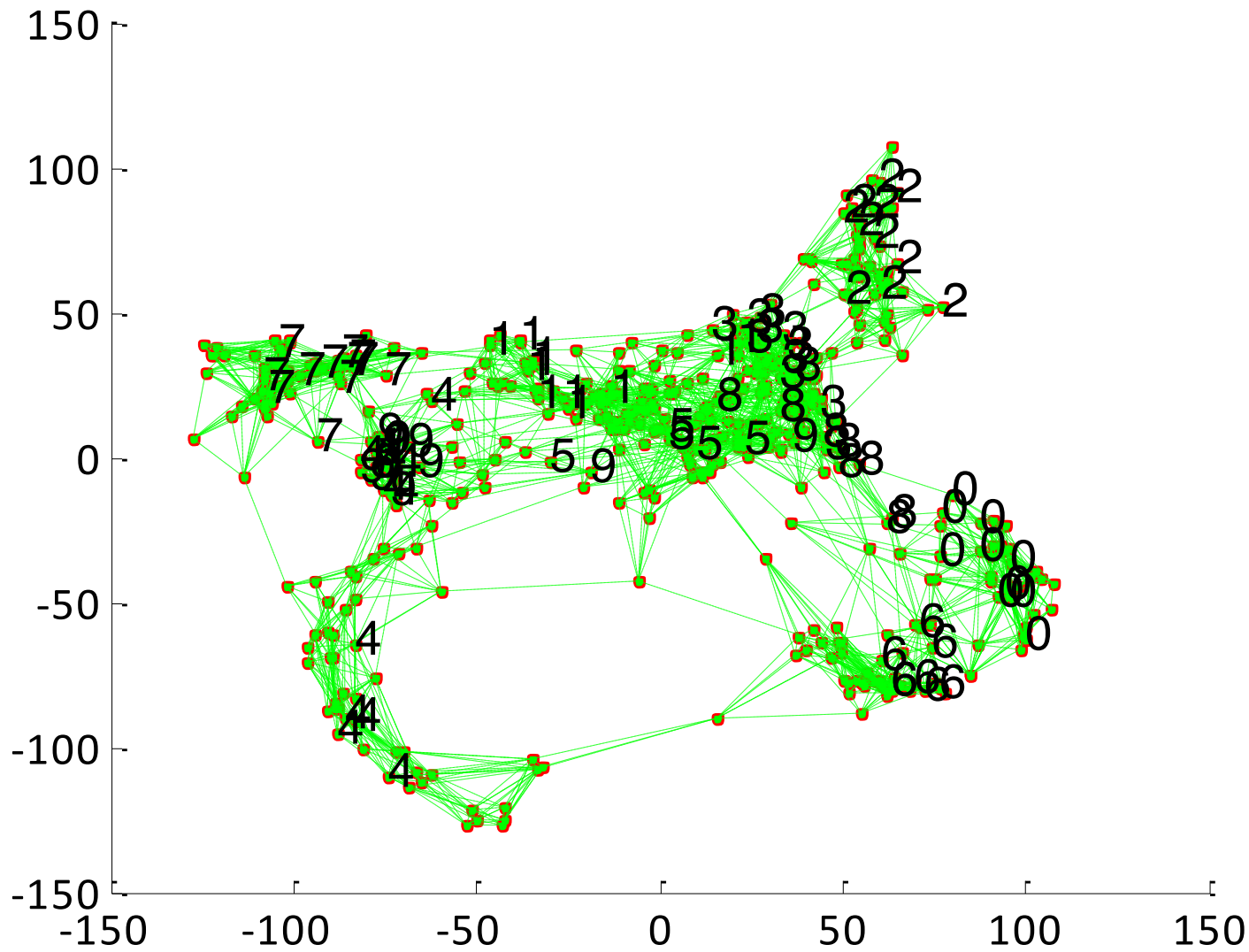
1. construct a neighborhood graph (kNN and  $\epsilon$ -NN)
2. calculate distance matrix as the shortest path on the graph
3. apply MDS on the distance matrix



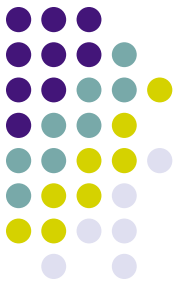
# Manifold learning



Optdigits after Isomap (with neighborhood graph).



# Manifold learning



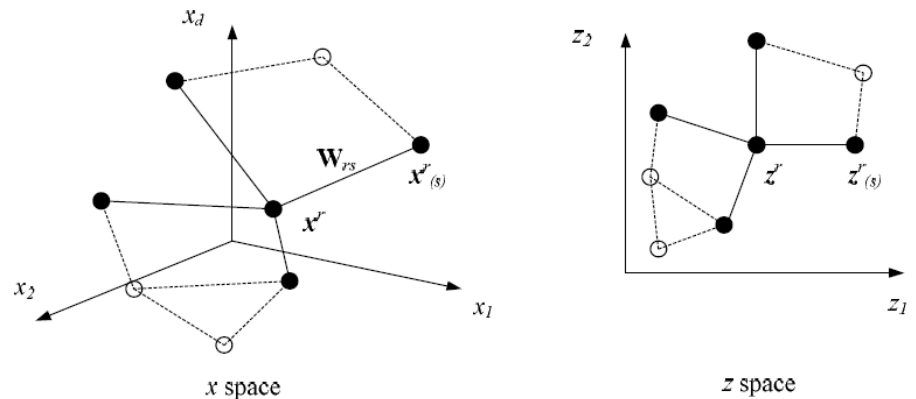
Local Linear Embedding (LLE):

1. find neighbors for each instance
2. calculate a linear reconstruction for an instance

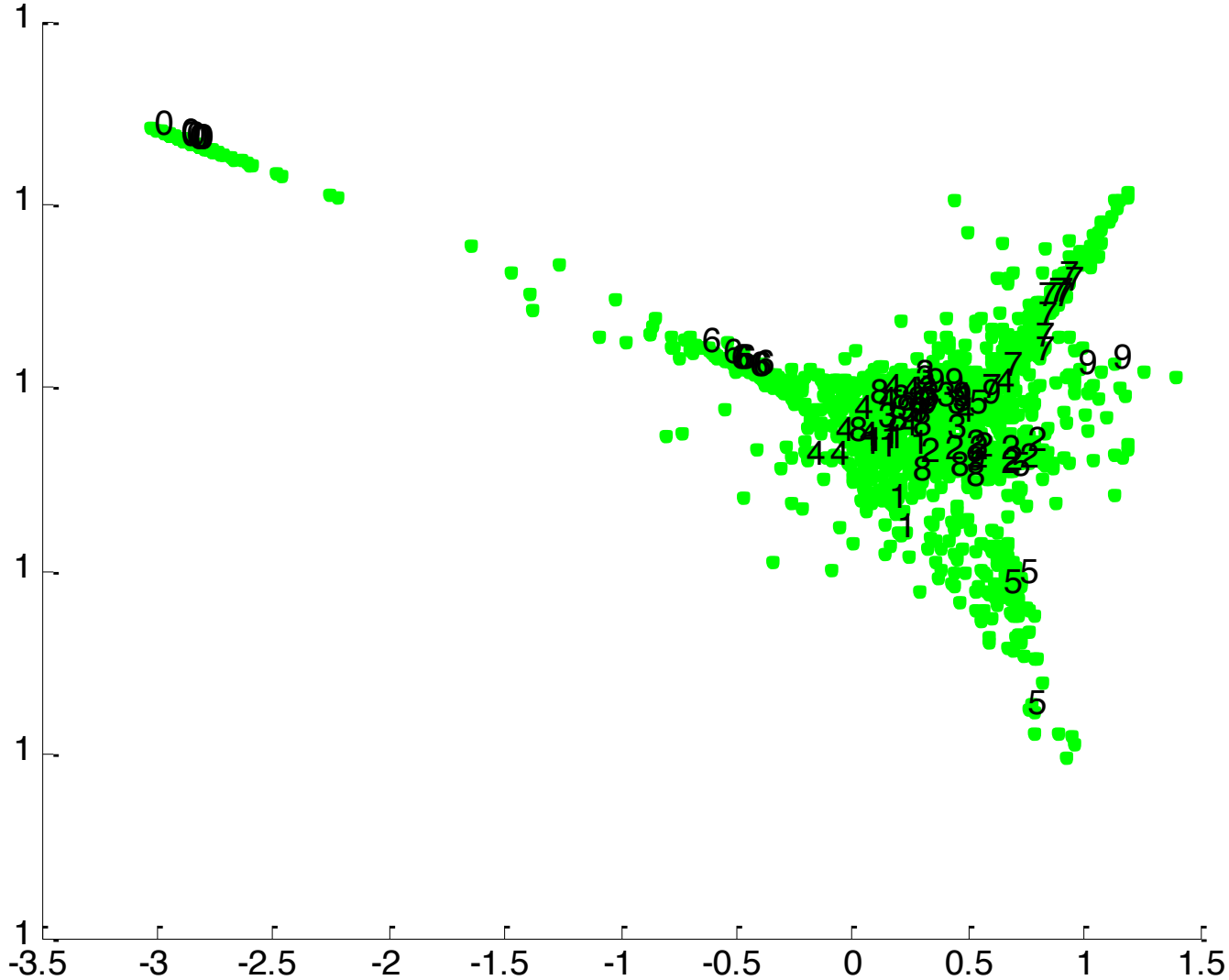
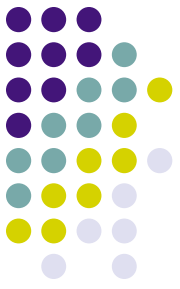
$$\sum_r \|\mathbf{x}^r - \sum_s \mathbf{W}_{rs} \mathbf{x}_{(r)}^s\|^2$$

3. find low dimensional instances preserving the reconstruction

$$\sum_r \|\mathbf{z}^r - \sum_s \mathbf{W}_{rs} \mathbf{z}^s\|^2$$



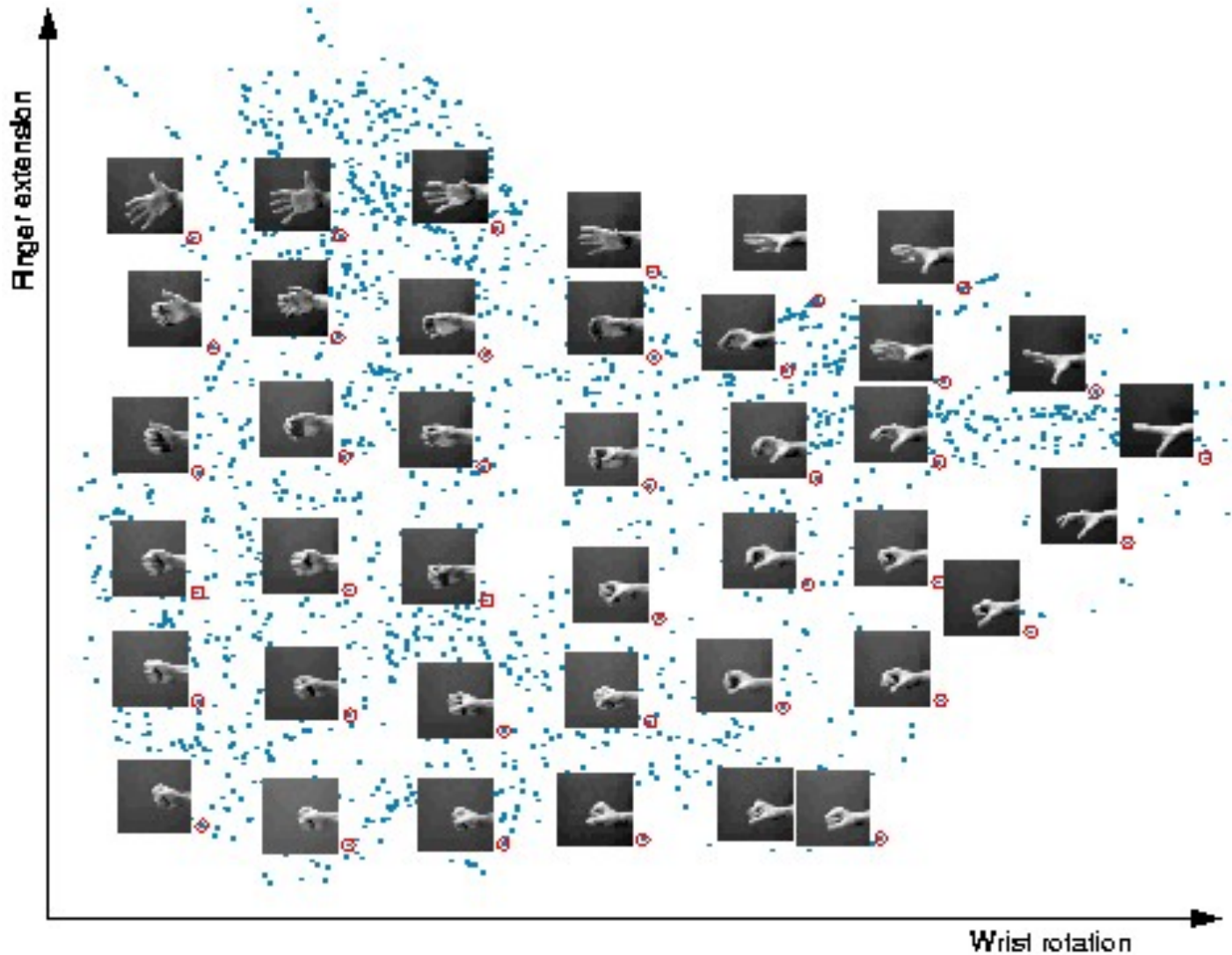
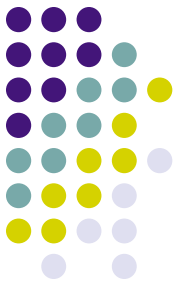
# Manifold learning



from [Intro. ML]

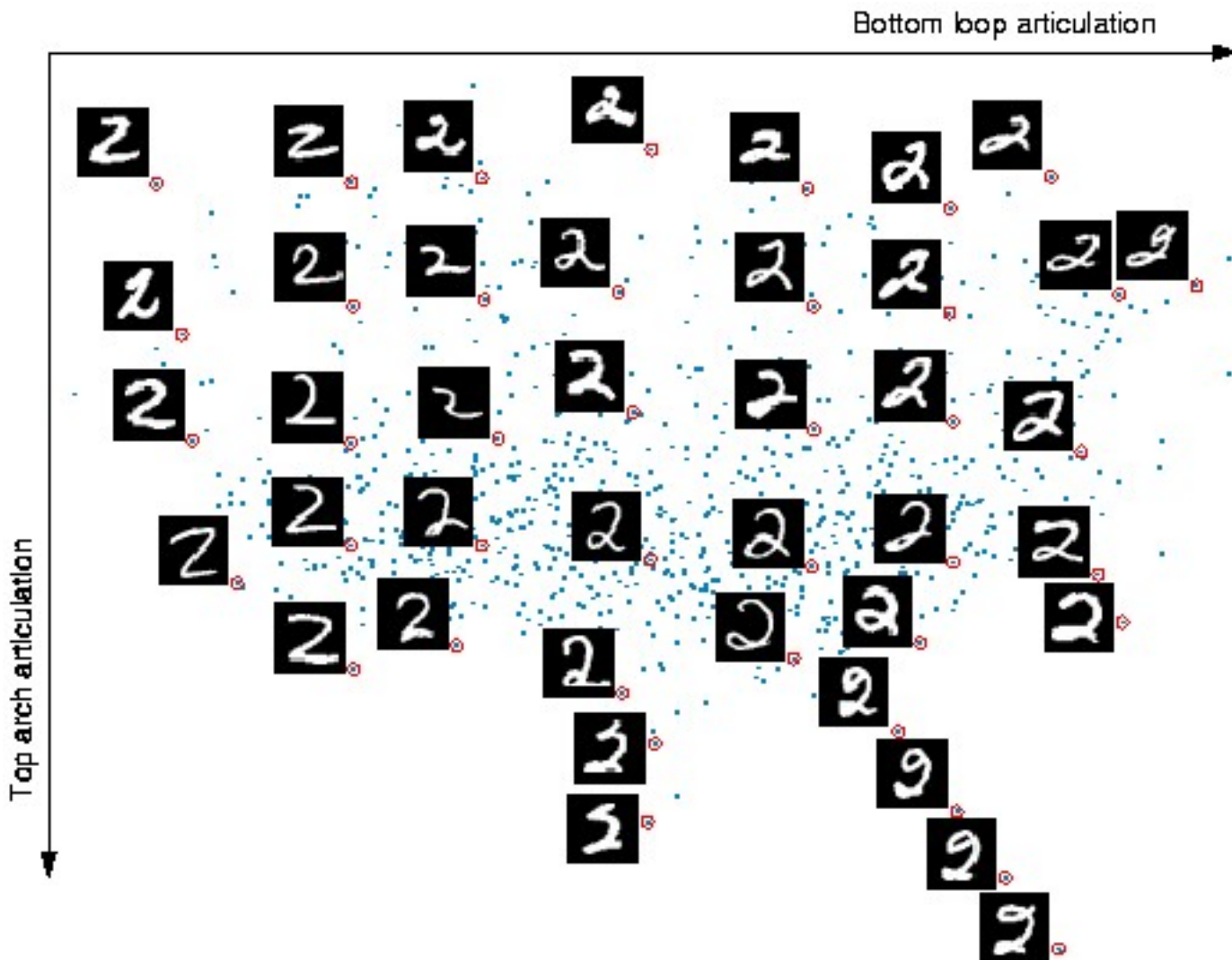
# Manifold learning

more manifold learning examples



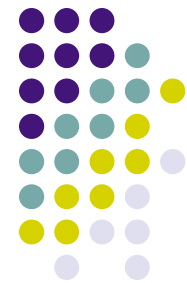
# Manifold learning

more manifold learning examples





# 机器学习/模式识别基础



## 小结

- 预测与识别
  - 用数据定义“模式”
- 预测算法
- 特征提取