

# Statistical Pattern Recognition

---

- Review

---

# Statistical Pattern Recognition - Review

- Illustration of a classifier Problem

$$\boxed{a} \rightarrow \boxed{b}$$

- Classification vs Regression
- Preprocessing & Feature extraction
- Curse of dimensionality
- Function approximation

Polynomial curve fitting

Degrees of Freedom (# weights)

Bias & Variance, Regularization

Bayes' Theorem: Probabilistic formulation for the cloumplexity problem

## Decision Making

### Probability Density Estimation

- parametric models  
(Maximum Likelihood & Bayesian learning)
- non parametric
  - \* Histograms
  - \* Kernel-based Methods
  - \* k-nearest neighbour

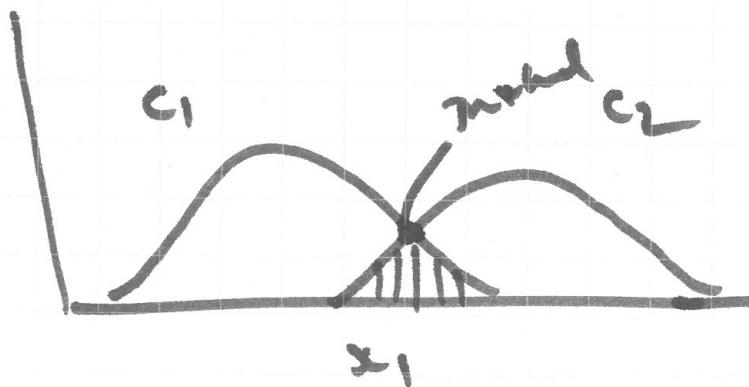


$$256 \times 256 \times 8 \longrightarrow$$

$$\frac{256 \times 256 \times 8}{2} \approx 158000$$

$$\approx 10$$

$$x_1 \rightarrow \frac{\text{height}}{\text{width}} \rightarrow$$

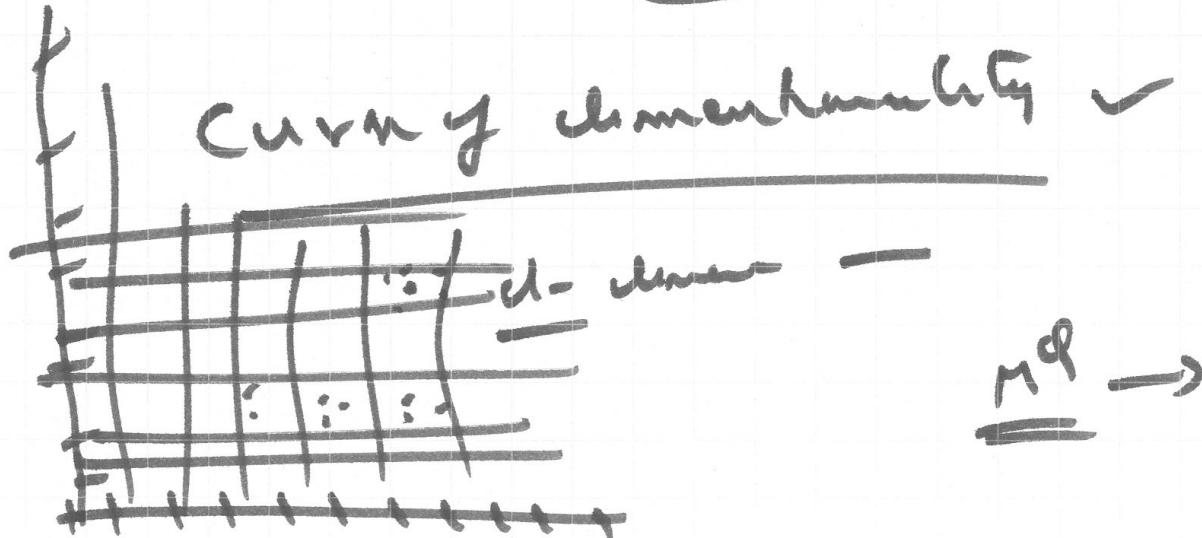


## Classification vs Regression

Prediction ✓



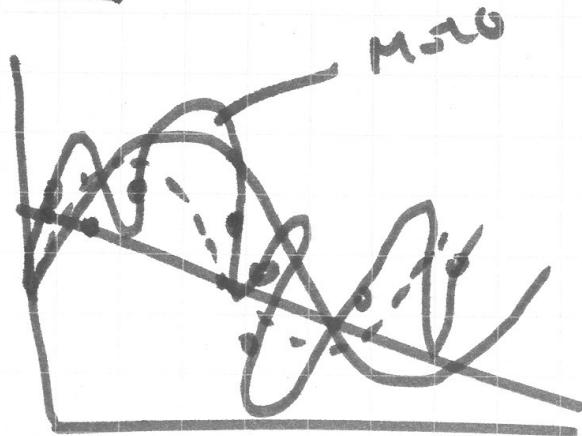
Curves of dimensionality ✓



Curve fitting : Polynomial Model

$$y(x) = w_0 + w_1 x + \dots + w_m x^m \\ = \sum_{n=0}^m w_n x^n$$

$$y(x) = \underbrace{\sin(0.2x) + 0.5}_{n=0}$$

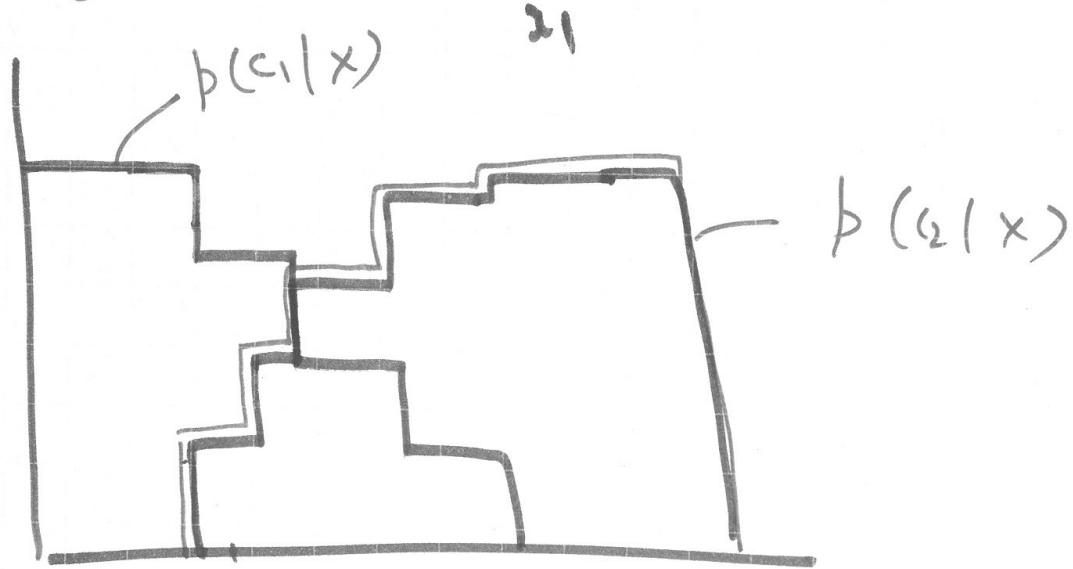


$$p(c_k | x_L) = \frac{\sum p(x_L | c_k) \cdot p(c_k)}{p(x_L)} = 1$$

↓  
 posterior  
 ↓  
 class considered  
 p<sub>post</sub>  
 ↓  
 unlabeled p<sub>post</sub>

$$\sum_{c_k} p(c_k | x_L) = 1$$

$$\sum_k p(x_L | c_k) p(c_k) = p(x_L)$$

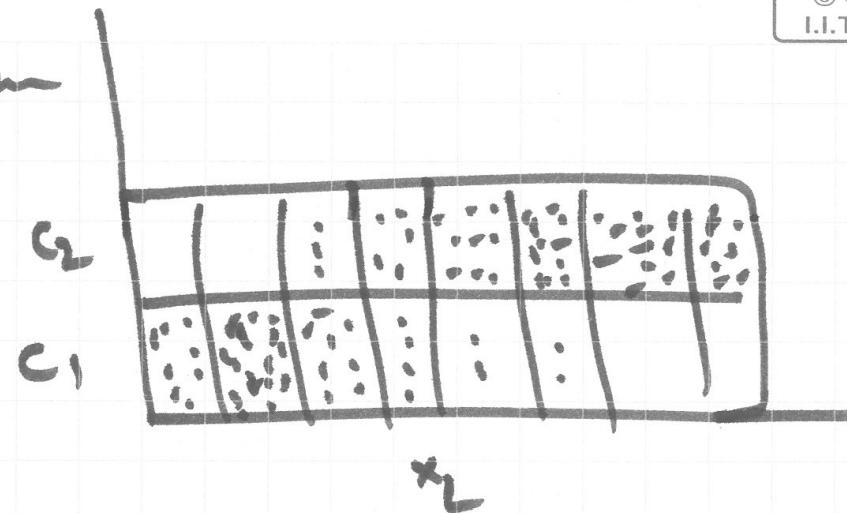


Buyee's Problem

$$N = 100$$

$$C_1 = 70$$

$$C_2 = 30$$



$$P(C_k | x_L) = \frac{4}{100}$$

$$\underbrace{P(C_k | x_L) P(x_L)}_{\frac{4}{9} \times \frac{1}{10}} ; P(x_L | c_k) P(c_k)$$

$$= \frac{4}{20} \cdot \frac{7}{10} = \frac{14}{100}$$

## Decision Rules

$$p(c_{ik}|x) > p(c_j|x) \quad i \neq j$$

$$p(x|c_k) p(c_k) > p(x|c_j) p(c_j)$$

$$\begin{aligned} p(c) = & \sum p(x \in R_2, c_i) + p(x \in R_1, c_j) \\ & \cdot \int_{R_2} p(c_i|x) p(c_i) + \int_{R_1} p(c_j|x) p(c_j) \end{aligned}$$

# Probability Density Estimation

## Parametric Methods

- Maximum likelihood
- Bayesian inference

## Non parametric Methods

- Histogram
  - Kernel based approaches
  - K-nearest neighbor ✓
- K-L distance (Kullback - Leibler distance) ✓

Mixture Models :

Description of the Models

Training Methods

- ✓ — Nonlinear optimization
- ✓ — Re-estimation [EM-Algorithm]
- ✓ — Stochastic sequential estimation

## Parametric Methods

### Multivariate Gaussian

$$p(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left[ -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right]$$

### Discriminant function

$$y_k(x) = \ln p(x|c_k) + \ln p(c_k)$$

### Maximum Likelihood

$$p(c_k|x) = \frac{p(x|c_k) p(c_k)}{p(x)}$$

## Maximum Likelihood

$$p(x|\theta); \quad x = \{x^1, x^2, \dots, x^N\} \quad \checkmark$$

$$\theta = \{\theta_1, \theta_2, \dots, \theta_M\}$$

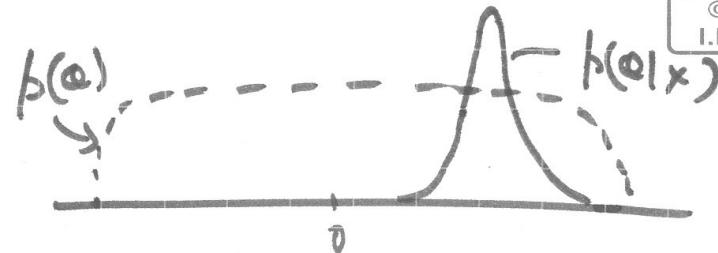
$$p(x|\theta) = \prod_{m=1}^N p(x^m|\theta) = L(\theta) \quad \checkmark$$

$$\mathcal{E} = -n L(\theta) = - \sum_{m=1}^N \ln p(x^m|\theta) \quad \checkmark$$

$\frac{\partial \mathcal{E}}{\partial \theta} = 0 \Rightarrow$  optimal parameters

$$\hat{\mu} = \frac{1}{N} \sum_{m=1}^N x^m; \quad \hat{\Sigma} = \frac{1}{N} \sum_{m=1}^N (x^m - \hat{\mu})(x^m - \hat{\mu})^T$$

Bayesian Inference



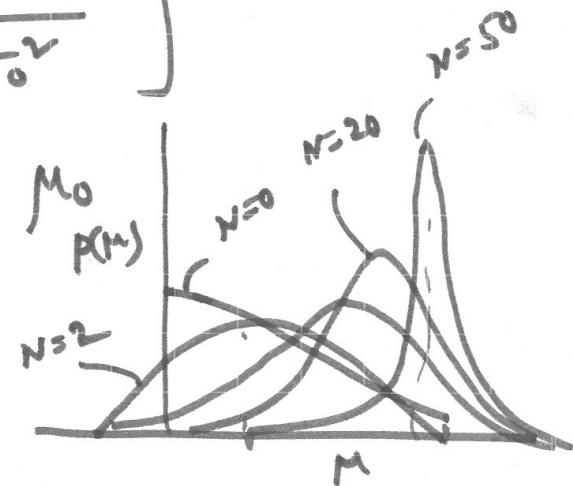
$$p(x|\theta) = \prod_{m=1}^N p(x_m|\theta)$$

$$\underline{p(\theta|x)} = \frac{\underline{p(x|\theta) f(\theta)}}{p(x)} = \frac{\underline{p(\theta)}}{\underline{f(x)}} \prod_{m=1}^N p(x_m|\theta)$$

$$p_0(\mu) = \frac{1}{(2\pi\sigma_0^2)^{\frac{1}{2}}} \exp \left[ -\frac{(\mu-\mu_0)^2}{2\sigma_0^2} \right]$$

$$\mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \bar{x} + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0$$

$$\sqrt{\frac{1}{\sigma_N^2}} = \frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}$$



# Nonparametric Methods

## Histograms

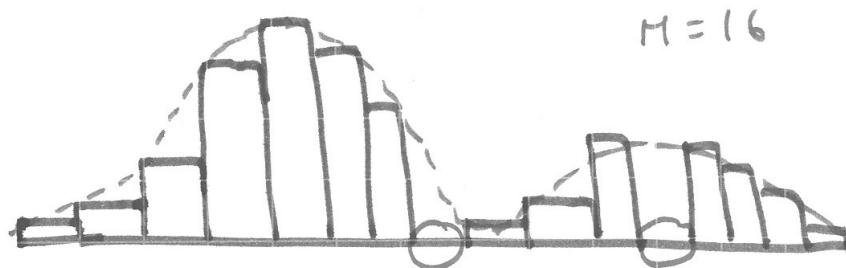
Smoothing factor  
= width of one bin

### Disadv

- Discontinuous ✓
- Curves dimensionality

Ave  $\bar{x}$  can be constructed incrementally ✓  
(smooth out the crucial details)

# bins  $\begin{cases} \text{low} & (\text{smooth & noisy}) \\ \text{high} & (\text{sparse & noisy}) \end{cases}$  ✓



## Density Estimation (in General)

$$\checkmark P = \int_R p(x) dx ; \quad P(k) = \frac{N!}{k!(N-k)!} p^k (1-p)^{N-k}$$

$$\text{Mean } P = \frac{k}{N} ; \quad P = \int_R p(x) dx \approx \frac{p(x)}{N}$$

$$p(x) = \frac{1}{\sqrt{Nv}}$$

Kernel-based Method

$$k = \sum_{m=1}^N K\left(\frac{x-x_m}{v}\right) ; \quad \tilde{p}(x) = \frac{1}{N} \sum_{m=1}^N \frac{1}{v} K\left(\frac{x-x_m}{v}\right)$$

Multivariate Normal Kernel

$$\tilde{p}(x) = \frac{1}{N} \sum_{m=1}^N \frac{1}{(2\pi v)^d} e^{-\frac{\|x-x_m\|^2}{2v}}$$

$$\tilde{p}(x) = h(x) * p(x)$$

$h(x)$  = impulse fn.

$$\tilde{p}(x) \approx p(x)$$

$$h \rightarrow 0$$

$$\boxed{\tilde{p}(x) = p(x)}$$

∴ due to finite points

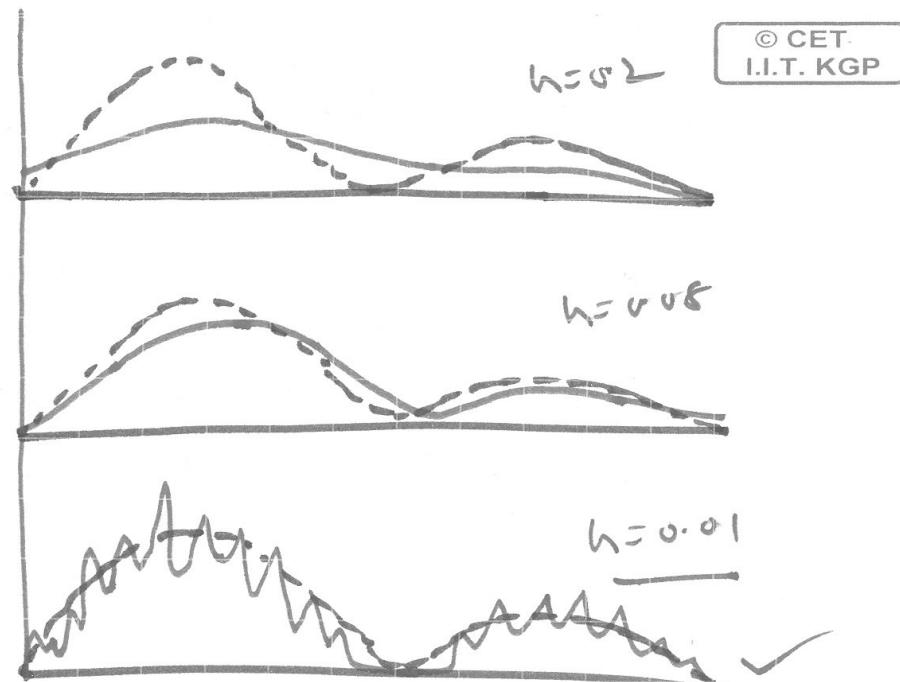
$$\tilde{p}(x) \text{ is noisy}$$

( $K$ -nearest neighbor)

Total data points =  $N$

# points of class  $K$  =  $\underline{N_K}$

$$\sum_K N_K = N$$



$$p(x|c_k) = \frac{N_k}{N} \checkmark$$

$$p(x) = \frac{K}{N} \checkmark$$

$$p(c_k) = \frac{N_k}{N} \checkmark$$

$$P(c_a|x) = \frac{p(x|c_a) p(c_a)}{p(x)} = \frac{p(c_a)}{K}$$

Smoothing parameters

Mistogram - width of bins

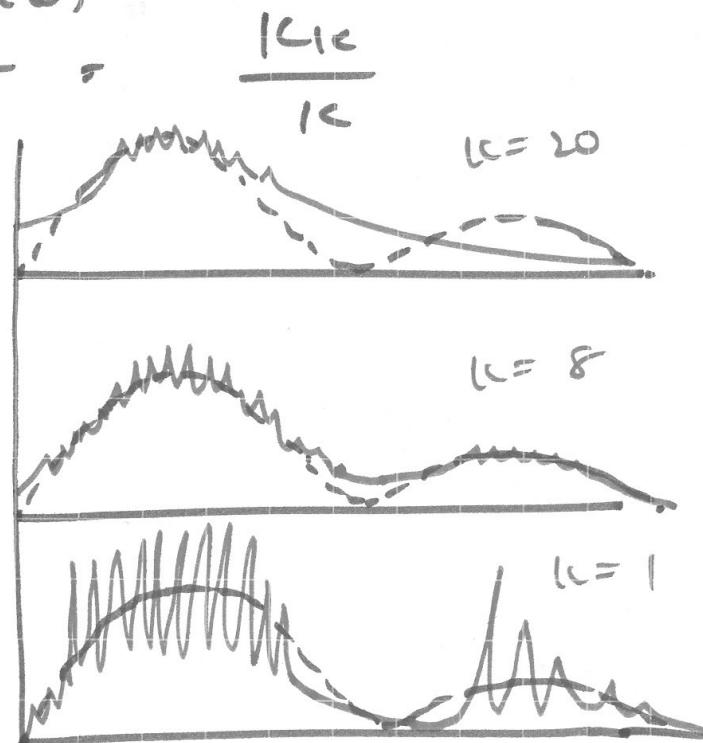
(kernel) - width of the kernel "h"

K-nearest - number of k  
neighbors

$$L = - \int \tilde{p}(z) \ln \frac{\tilde{p}(z)}{p(z)} dz$$

(K-L distance)

Cumulative - Levenshtein distance



$\Rightarrow$  Distance between true and  
the estimated distance

GMM : EM Algorithm

---

## Mixture Models

### Parametric form

Specify functional form is easier ✓

Evaluation with new data is faster

### Non-parametric

No specific functional form ✓

# variables grow with # training points  
very slow for evaluation of new data

### Semi-parametric ✓

No specific functional form

# variables grows with complexity  
Computationally intensive ✓

### Training Details

Non-linear optimization

Re-estimation [EM algorithm]

Stochastic sequential estimate

# Mixture Model (cont..)

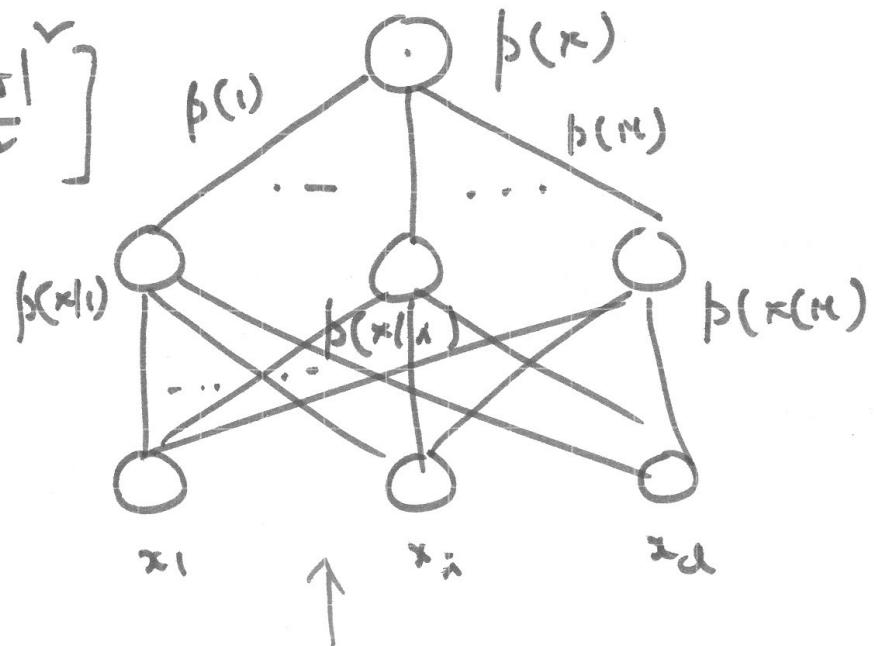
$$p(x) = \sum_{j=1}^M p(x|j) p(j)$$

Mixture distribution  
Component density

Mixing parameters  
Linear combination of component densities

$$p(x|j) = \frac{1}{(2\pi\sigma_j^2)^{\frac{d}{2}}} \exp\left[-\frac{|x - \mu_j|^2}{2\sigma_j^2}\right]$$

$p(x|c_k)$  - class specific  
MMP model



## Maximum likelihood formulation (Mixture model)

$$\begin{aligned}
 E &= -\ln L = -\sum_{n=1}^N \ln p(x^n) = -\sum_{n=1}^N \ln \left[ \sum_{j=1}^M p(x^n|j) p(j) \right] \\
 &= -\sum_{n=1}^N \ln \left[ \sum_{j=1}^M p(x^n|j) p(j) \right]
 \end{aligned}$$

Complexity in minimizing the error function (E)

- Singularity problems
- Avoiding local minima

M ✓

$$\frac{\partial E}{\partial \mu_J} = 0 ; \quad \frac{\partial E}{\partial \sigma_T} = 0 ; \quad \frac{\partial E}{\partial \rho_J} = 0 ;$$

$$\hat{\mu}_J = \sum_{x^n} P(J|x^n) \times x^n / \sum_{x^n} P(J|x^n)$$

$$\hat{\sigma}_J^2 = \frac{1}{n} \sum_{x^n} P(J|x^n) \| x^n - \hat{\mu}_J \|^2$$

$$\hat{P}(J) = \frac{1}{n} \sum_{x^n} P(J|x^n)$$

# The EM Algorithm

$\hat{\mu}_j$ ,  $\hat{\sigma}_j^2$ ,  $\hat{p}(j)$   $\rightarrow$  do not provide direct soln  
Nonlinear coupled eq

Sol'n :- Iterative scheme to find min of ' $E'$ '

1. initial guess of the parameter  $\Rightarrow$  old parameter

2. Evaluate RHS of  $\hat{\mu}_j$ ,  $\hat{\sigma}_j^2$ ,  $\hat{p}(j)$  ✓

3. Re estimate of the parameter  $\Rightarrow$  new parameter ✓

4. Verify the curr ' $E'$  with new parameter ✓

$$E^{\text{new}} - E^{\text{old}} = - \sum_{\mathcal{J}} \ln \left[ \frac{p^{\text{new}}(x^n)}{p^{\text{old}}(x^n)} \right]$$

$$= - \sum_{\mathcal{J}} \frac{\sum_{\mathcal{J}} p^{\text{new}}(\mathcal{J}) p^{\text{new}}(x^n | \mathcal{J})}{p^{\text{old}}(x^n)} \cdot \frac{p^{\text{old}}(\mathcal{J}(x^n))}{p^{\text{old}}(\mathcal{J}(x^n))}$$

$$\leq - \sum_{\mathcal{J}} \left\{ p^{\text{old}}(\mathcal{J}(x^n)) \ln \left[ \frac{p^{\text{new}}(\mathcal{J}) p^{\text{new}}(x^n | \mathcal{J})}{p^{\text{old}}(\mathcal{J}) p^{\text{old}}(x^n | \mathcal{J})} \right] \right\}$$

$$E^{\text{new}} \leq E^{\text{old}} + Q$$

$$Q = - \sum_{\mathcal{J}} \left\{ p^{\text{old}}(\mathcal{J}(x^n)) \ln \left[ p^{\text{new}}(\mathcal{J}) p^{\text{new}}(\mathcal{J}(x^n)) \right] \right\}$$

$$Q_{\text{true}} = - \sum_{\mathcal{J}} \left\{ p^{\text{old}}(\mathcal{J}(x^n)) \left[ \ln p^{\text{new}}(\mathcal{J}) - \alpha \ln \frac{1}{f_{\mathcal{J}}} - \frac{\|x^n - \mu_{\mathcal{J}}^{\text{new}}\|^2}{2(f_{\mathcal{J}}^{\text{new}})^2} \right] \right\} \text{true}$$

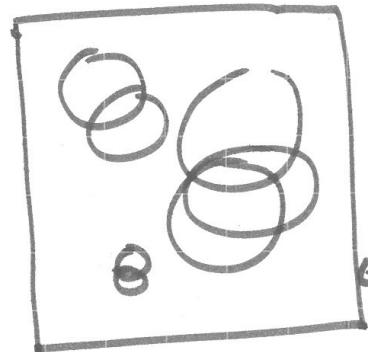
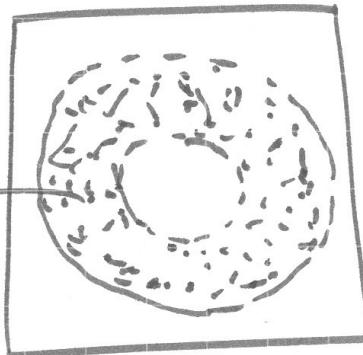
$$\overline{\mu}_T^{\text{new}} = \frac{\sum_{j=1}^3 p^{\text{old}}(\mathcal{T}(x_j)) x_j}{\sum_{j=1}^3 p^{\text{old}}(\mathcal{T}(x_j))} \quad \checkmark$$

$$(\overline{\mu}_T^{\text{new}})^r = \frac{1}{d} \frac{\sum_{j=1}^3 p^{\text{old}}(\mathcal{T}(x_j)) \|x_j - \overline{\mu}_T^{\text{new}}\|^r}{\sum_{j=1}^3 p^{\text{old}}(\mathcal{T}(x_j))} \quad \checkmark$$

$$P(\mathcal{T})^{\text{new}} = \frac{1}{n} \sum_{j=1}^n p^{\text{old}}(\mathcal{T}(x_j)) \quad \checkmark$$

## Example

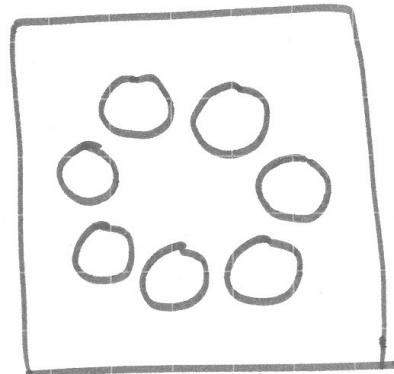
1000  
data points



© CET  
I.I.T. KGP

Initialization

After 20 iterations



GMM vs classification

Incomplete data

hybridized complete data, \*

$$2^m \rightarrow (1, M) \checkmark$$

x3

$$E^{\text{Comb}} = \cancel{-\ln} - \sum_n \ln \underbrace{p^{\text{new}}(x^n, t^n)}_{p^{\text{new}}(x^n | t^n)}$$

$$= - \sum_n \ln \left[ p^{\text{new}}(t^n) p^{\text{new}}(x^n | t^n) \right]$$

1. Given some values to all parameters of GMM ✓
2. Find PD of  $t^n$  using Bayes' Theorem & old formula
3. Compute  $E[E^{\text{Comb}}]$  wrt PD of  $t^n$  ✓ - Est step
4. Compute New parameter by minimizing  $E$  ✓

$$E[E^{\text{Comb}}] = \sum_{t=1}^T \sum_{n=1}^N E^{\text{Comb}} \frac{1}{T} p^{\text{old}}(t^n | x^n)$$

$$= - \sum_n \left\{ \sum_J p^{\text{old}}(J | x^n) \ln \left[ p^{\text{new}}(J) p^{\text{new}}(x^n | J) \right] \right\}$$