

# Computational Intelligence: Methods and Applications

## Lecture 26 Density estimation, Expectation Maximization.

Włodzisław Duch  
SCE, NTU, Singapore  
Google: Duch

## Cognitive inspirations

How do we recognize objects? Nobody really knows ...

Objects have features, combinations of features, or rather distributions of feature values in Feature Spaces (FS), characterize objects.

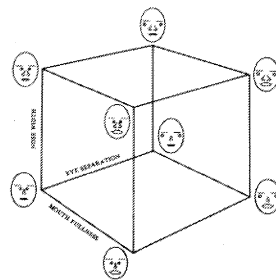
A single object is a point in the FS;

similar objects create a category, or a concept: for ex. happy or sad face, corresponding to some area of the feature space.

$P(\text{Angry}|\text{Face features})$  will have maximum around one of the corners.

In cognitive psychology FS are called “psychological spaces”.

The shape of the  $P(X|C)$  distribution may be quite complex, estimated using known samples to create a fuzzy prototype.



## Density estimation

Knowledge of joint probability density  $P(C,X)$  or just  $P(X)$  allows to do much more than just discrimination!

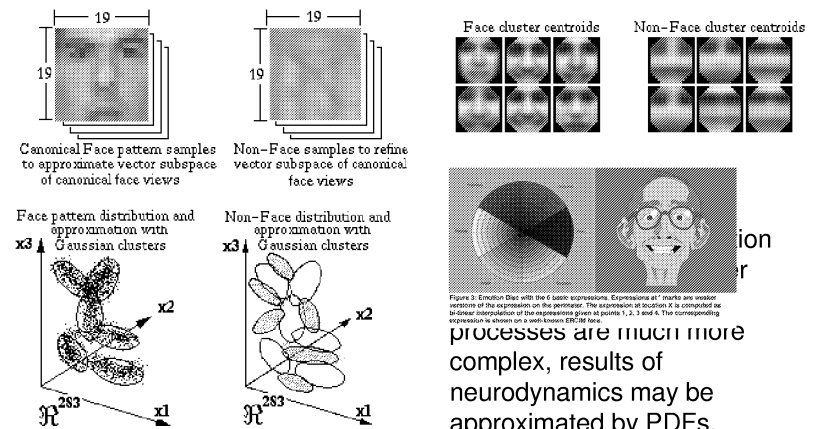
Local maxima of probability density functions (PDFs) correspond to combination of features defining objects in feature spaces.

Estimating PDFs we may create adaptive systems learning from data with or without supervision. They are useful for:

- Auto-association and hetero-association.
- Completion of unknown parts of the input vector (content-addressable memory), prediction of missing values.
- Extraction of logical rules, classical and probabilistic (or fuzzy).
- Finding prototypes for objects or categories in feature spaces.
- Using density functions as heuristics for solution of complex problems, learning from partial info & solving complex problems.

## Object recognition

Population of neural columns, each acting as a weak classifiers to recognize some features, working in chorus – similar to “stacking”. Second-order similarity in low-dimensional (<300) space is sufficient.



## Missing features

Suppose that one of the features  $X = (X_1, X_2, \dots, X_d)$ , for example  $X_1$ , is missing. What is the most likely value for this feature?

Frequently an average value  $E(X_1)$  is used, but is this a reasonable idea? The average may fall in an area where there is no data!

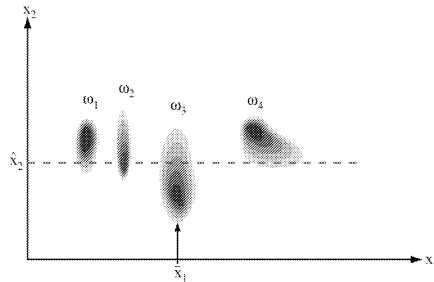


Fig. 2.22,  
Duda, Hart &  
Stork, 2000

In this case if  $X_2$  is known the best answer is the value corresponding to the maximum density at  $\omega_2$ .

Recover missing values searching for maximum density!

## Maximum likelihood

Suppose that the density  $P(X; \theta)$  is approximated using a combination of some parameterized functions. Given a set of observations (data samples)  $D = \{\mathbf{X}^{(i)}, i=1..n$ , what parameters should one choose?

Parameters  $\theta$  may include also missing values, as a part of the model.

A reasonable assumption is that the observed data  $D$  should have high chance of being generated using the model  $P(D; \theta)$ . Assuming that the data vectors  $X^{(i)}$  are independent, the likelihood of obtaining the dataset  $D$  is:

$$l(\theta; D) = \prod_{i=1}^n P(\mathbf{X}^{(i)}; \theta)$$

The most probable parameters of the model (including missing values) maximize likelihood. To avoid products use logarithm and minimize  $-L$

$$\min_{\theta} L(\theta; D) = -\min_{\theta} \ln l(\theta; D) = -\sum_{i=1}^n \ln P(\mathbf{X}^{(i)}; \theta)$$

## Solution

Maximum is found by setting the derivative of the log-likelihood to 0:

$$\frac{\partial L(\theta)}{\partial \theta} = -\sum_{i=1}^n \frac{1}{P(\mathbf{X}^{(i)}; \theta)} \frac{\partial P(\mathbf{X}^{(i)}; \theta)}{\partial \theta} = 0$$

Depending on the parameterization, sometimes this can be solved analytically, but for almost all interesting functions (including Gaussians) iterative numerical minimization methods are used.

Many local minima of the likelihood function are expected, so the minimization problem may be difficult.

Likelihood estimation may be carried for samples from a given class  $P(X|\omega, \theta)$ , assuming that the probability of generating  $n$  such samples is equal to  $P(X|\omega, \theta)^n$ , and the *a priori* class probabilities are estimated from their frequencies.

Such parametric models are called “generative” models.

## Example

Example from “Maximum likelihood from incomplete data via the EM algorithm”, Dempster, Laird, Rubin 1977, data by Rao, from population genetics. There are 197 observation of 4 types of bugs:

$n_1=125$  times species (class)  $\omega_1$ ,  $n_2 = 18$  from class  $\omega_2$ ,  $n_3 = 20$  from class  $\omega_3$ , and  $n_4 = 34$  from class  $\omega_4$ . An expert provided the following parametric expressions for the probabilities to find these bugs:

$$P(\mathbf{X} | \omega_1; \theta) = (2 + \theta) / 4; P(\mathbf{X} | \omega_2; \theta) = (1 - \theta) / 4$$

$$P(\mathbf{X} | \omega_3; \theta) = (1 - \theta) / 4; P(\mathbf{X} | \omega_4; \theta) = \theta / 4$$

Find the value of parameter that maximize the likelihood:

$$l(\theta) = P(\mathbf{X} | \omega_1; \theta)^{n_1} P(\mathbf{X} | \omega_2; \theta)^{n_2} P(\mathbf{X} | \omega_3; \theta)^{n_3} P(\mathbf{X} | \omega_4; \theta)^{n_4}$$

Multiplicative constants  $n!/(n_1!n_2!n_3!n_4!)$  are not important here.

## Solution

Log-likelihood: 
$$L(\theta) = -n_1 \ln P(\mathbf{X} | \omega_1; \theta) - n_2 \ln P(\mathbf{X} | \omega_2; \theta) - n_3 \ln P(\mathbf{X} | \omega_3; \theta) - n_4 \ln P(\mathbf{X} | \omega_4; \theta)$$

Derivative: 
$$\frac{\partial L(\theta)}{\partial \theta} = - \left( \frac{n_1}{2 + \theta} - \frac{n_2 + n_3}{1 - \theta} + \frac{n_4}{\theta} \right) = 0$$

Quadratic equation for  $\theta$  allows for analytical solution:  $\theta = 0.6268$ ; now the model may provide estimations of expected frequencies:

$$\langle n_1 \rangle = N(2 + \theta) / 4 = 129.4$$

For all 4 classes, expected (real) number of observation:

$$\langle n_1 \rangle = 129 \text{ (125)}, \langle n_2 \rangle = 18 \text{ (18)}, \langle n_3 \rangle = 18 \text{ (20)}, \langle n_4 \rangle = 31 \text{ (34)}$$

In practice analytic solutions are rarely possible.

## General formulation

Given data vectors  $D = \{\mathbf{X}^{(i)}\}$ ,  $i=1..n$ , and some parametric functions  $P(\mathbf{X}|\theta)$  that model the density of the data  $P(\mathbf{X})$  the best parameters should minimize log-likelihood for all data samples:

$$\theta^* = \arg \min_{\theta} L(\theta | D) = - \sum_{i=1}^n \ln P(\mathbf{X}^{(i)}; \theta)$$

$P(\mathbf{X}|\theta)$  is frequently a Gaussian mixture; for a single Gaussian standard solution will give the formula for mean and variance.

Assume now that  $\mathbf{X}$  is not complete – features, or whole parts of the vector are missing. Let  $\mathbf{Z}=(\mathbf{X},\mathbf{Y})$  be the complete vector. Joint density:

$$P(\mathbf{Z} | \theta) = P(\mathbf{X}, \mathbf{Y} | \theta) = P(\mathbf{Y} | \mathbf{X}, \theta) P(\mathbf{X} | \theta)$$

Initial joint density may be formed analyzing cases without missing values; the idea is to maximize the complete data likelihood.

## What to expect? E-step.

Original likelihood function  $L(\theta|\mathbf{X})$  is based on incomplete information, and since  $\mathbf{Y}$  is unknown it may be treated as a random variable that should be estimated.

Complete-data likelihood function  $L(\theta|\mathbf{Z})=L(\theta|\mathbf{X},\mathbf{Y})$  may be evaluated calculating the expectation of incomplete likelihood over  $\mathbf{Y}$ . This is done iteratively, starting from initial estimation  $\theta^{i-1}$  new estimation  $\theta^i$  of parameters and missing values is generated:

$$Q(\theta | \theta^{i-1}) = E_{\mathbf{Y}} [\ln P(\mathbf{X}, \mathbf{Y} | \theta) | \mathbf{X}, \theta^{i-1}]$$

where  $\mathbf{X}$  and  $\theta^{i-1}$  are fixed,  $\theta$  is a free variable, and the conditional expectation is calculated using the joint distribution of the  $\mathbf{X}, \mathbf{Y}$  variable with fixed  $\mathbf{X}$

$$E[Y | X = x] = \int y P_{Y|X}(x, y) dy$$

See detailed ML discussion in Duda, Hart & Stork, Chapter 3