

Pattern Recognition

(Αναγνώριση Προτύπων)

Nonparametric Techniques

(Μη-Παραμετρικές Τεχνικές)

Panos Trahanias

UNIVERSITY OF CRETE DEPARTMENT of COMPUTER

SCIENCE



Issues in Parametric Techniques:

- Often the shape of the distribution is not known.
- In practice, many distributions are <u>multimodal</u> (contain more than one modes), whereas most parametric models are <u>unimodal</u>.
- Approximation of multimodal distributions as a product of unimodal ones is not appropriate.
- Non-parametric techniques: Estimation of distribution function from scratch.
 - Settimation of pdf's $p(x/\omega_i)$ from the data via generalization of the <u>multi-dimensional histogram</u>.
 - Solution Ψ Direct Estimation of the posterior probabilities $P(\omega_i / x)$ and the discriminant functions.



Based on the fact that the probability of sample x is found within area R results from

$$P = P(\mathbf{x} \in R) = \int_{\mathbf{x}' \in R} p(\mathbf{x}') d\mathbf{x}'$$

The above integral can be approximated either by the product of $p(\mathbf{x})$ times the area R, or by the number of samples that are found within R

$$p(x|\omega_i)$$

$$0.4$$

$$p(x^*)$$

$$0.3$$

$$0.2$$

$$0.2$$

$$0.2$$

$$0.2$$

$$0.2$$

$$0.2$$

$$R \xrightarrow{\omega_2}$$

$$\omega_1$$

$$P = \int_{\mathbf{x}' \in R} p(\mathbf{x}') d\mathbf{x}' \approx p(\mathbf{x}^*) \cdot V \approx k / n$$

V: Area R. In one-dimensional case, V= length of R

k: Number of samples within R

n: Total number of samples



Density (Distribution) Estimation

In order to estimate the density at *x*, we choose a series of areas *R*1, *R*2,...,*Rn*,

that contain x, whereby Ri is used for i samples. Let Vn the volume of Rn, kn the number of samples within the n-th area, and pn(x) the n-th estimation of p(x). Then,





- There are two ways to create series of areas R_i for $p_n(x)$ to converge to p(x):
 - The volume of an initial area is reduced by defining a sequence of volumes V_n as Yunt Vions of n, e.g.
 Density estimation via the Parzen Windows

<u>method</u>

♥ We define k_n as a function of n, hence V_n increases until it contains k_n samples. e.g.

 □ Density estimation via the k_n-Nearest

 Neighbor method

Two Approaches







- Method based on finding the number of samples within a given area, where the area shrinks as the number of samples grows larger.
- The number of samples within the area is computed via the use of a "window function", the Parzen Window.



Parzen Windows



$$p_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{V_n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right) = \frac{1}{n} \sum_{i=1}^n \delta_n(\mathbf{x} - \mathbf{x}_i) \quad \text{ónov} \quad \delta_n(\mathbf{x}) = \frac{1}{V_n} \varphi\left(\frac{\mathbf{x}}{h_n}\right)$$

Window (.) can be a general functional, not necessarily hypercube. For $p_n(x)$ to be a valid p.d.f. for every *n*, it should hold

$$\varphi(\mathbf{u}) \ge 0$$
 και $\int_{\mathbf{u}} \varphi(\mathbf{u}) d\mathbf{u} = 1$

- ▶ p(x) is computed as a superposition of Gaussians, where each Gaussian is centered in the corresponding training sample. Parameter h_n is the Gaussian standard deviation!!!







Training Samples



Effect of Window Width h_n

$$\delta_n(\mathbf{x}) = \frac{1}{V_n} \varphi\left(\frac{\mathbf{x}}{h_n}\right), \quad V_n = h_n^d$$



Parameter *hn* has an effect on both the window width and its amplitude:

When *hn* is large (small), window becomes wide (narrow), window amplitude is small (large) and x must be far from (close to) xi for function $\delta n(x-xi)$ to change drastically compared to its value at $\delta n(0)$.



Effect of Window Width h_n

Hence



How the window width affects the estimation of p.d.f. p(x):

 \checkmark

✓ If *hn* is <u>large</u>, estimate pn(x) results as superposition of *n* "wide" functions centered at the training samples and constitutes a <u>smooth</u>, "out-of-focus" estimation of p(x), <u>with low resolution</u>.

If hn is small, pn(x) results as superposition of n "narrow" functions, i.e. <u>"erratic or noisy"</u> estimation of p(x).



Convergence Properties of $p_n(\mathbf{x})$

- Mean value of random variable $p_n(x)$: $\overline{p}_n(\mathbf{x}) = E[p_n(\mathbf{x})] = \int \delta_n(\mathbf{x} - \mathbf{v})p(\mathbf{v})d\mathbf{v}$ is the convolution of p(x) with the window function, hence it is a blurred version of p(x).
- It holds,

$$\delta_n(\mathbf{x} - \mathbf{v}) - \underset{V_n \to 0}{\longrightarrow} \delta(\mathbf{x} - \mathbf{v}) \quad \text{prote} \quad \overline{p}_n(\mathbf{x}) - \underset{n \to \infty}{\longrightarrow} p(\mathbf{x})$$

Variance of random variable $p_n(x)$:

$$\operatorname{var}(p_n(\mathbf{x})) = \frac{1}{nV_n} \int \frac{1}{V_n} \varphi^2 \left(\frac{\mathbf{x} - \mathbf{v}}{h_n}\right) p(\mathbf{v}) d\mathbf{v} - \frac{1}{n} \overline{p}_n^2(\mathbf{x}) \leq \frac{\sup(\varphi(\cdot))\overline{p}_n(\mathbf{x})}{nV_n}$$

- Accordingly, we have $\frac{n}{n}$ Accordingly, we have $\frac{n}{n}$
- Valid choices:

$$V_n = V_1 / \sqrt{n}$$
 $\dot{\eta}$ $V_n = V_1 / \ln n$



Parzen Windows Examples

$$p_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h_n} \varphi \left(\frac{x - x_i}{h_n} \varphi \left(\frac{x - x_i}{h_n} \varphi \left(\frac{x - x_i}{h_n} \right) \right) \right)$$

• ,

As n approaches infinity, estimation becomes accurate,

regardless of the window width.





Parzen Windows Examples



Parzen Windows Examples

As n approaches infinity, estimation becomes

accurate, regardless of the window width.





Parzen Windows for Pattern Classification

- Estimation of likelihood $p(\mathbf{x}|\omega_i)$ from training data via the Parzen windows approach, and use of Bayes rule for classification, e.g. computation of posterior probabilities and classification based on largest posterior probability.
- <u>Pros</u>: No need for problem specific assumptions, only the existence of training data set!
- <u>Cons</u>: Requires (many)^d data so that estimate converges to actual distribution.
 - Moreover, as the dimension increases, the requirement for (many)^d data becomes ((many)^{many (n)})ⁿ !!!! □ Curse of dimensionality !
 - The only way to overcome the above is the prior knowledge about "good" training data!
- The training error can become arbitrarily small (even zero), by choosing small windows! Still this is not desirable, since it'll result in <u>overfitting</u> and it'll reduce performance in the classification of test data.



Parzen Windows for Pattern Classification

Very small window [] very small/detailed separation in the feature space, not desirable! Larger window [] Larger training error, but better generalization performance! Desirable property.





In practice we're looking at small windows in areas with many (dense) training samples, and large windows in areas with little (sparse) training samples! How can this be accomplished...?



Probabilistic Neural Networks (PNN)

Input: {xk; k=1,...,d} d nodes, each corresponds to a feature.

wjk: weights that link *k*–*th* input node with *j*–*th* node of category ω_2 ω_{I} hidden layer (pattern node). aji Sparsely connected Hidden layer: *n* nodes, each corresponds to a pattern, i.e. pattern training sample, j=1,2,...,n. **Fully connected** wjk Output layer: *c* nodes, each represents a class. input x_1 x_{2} x_d *aji*: weights that link *j*–*th* hidden node with *i*–*th* output node, *i*=1,2,...,*c*



d-dimensional input vector x



PNN-Training





PNN-Classification

- Each class node sums the results of the pattern nodes that are connected to it. Accordingly, the activation of each class represents the p.d.f. estimate with a Gaussian Parzen window with covariance matrix σ²L_{4×d}, where I is the ide Algorithm 2 (PNN classification)

$$\begin{array}{ll}
 \underline{l} \ \underline{login \ initialize} \ k = 0, \mathbf{x} = \text{test pattern} \\
 \underline{l} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{login \ k \leftarrow k + 1} \\
 \underline{s} \ \underline{s} \\
 \underline{s} \ \underline{s} \$$



- Window width: instead of choosing it as a function of the number of training samples ($V_n = V_1 / \sqrt{n}$), why not choosing it as a function of the training data?
- \blacktriangleright Note: a large window is desirable in areas with small number of data, whereas a small window is appropriate in dense (with data) areas!
- k-nearest neighbor estimation algorithm:
 - \clubsuit Choose an initial area centered at **x** where $p(\mathbf{x})$ is estimated.
 - \checkmark Increase window until a predefined number of k_n samples falls inside the window <u> k_n </u> those are the k_n nearest V_n **neighbors** of *x*. ♥ Pdf (density) is estimated $\lim_{n \to \infty} k_n = \infty$; $\lim_{n \to \infty} k_n / n = 0$

Solutions of $p_n(x)$:











How to Choose k_n



Note that as *kn* increases, estimation accuracy increases...!

In classification problems, we usually tune *kn* (or *hn* for Parzen windows), so that the classifier operates with the lowest error for the <u>validation test dataset</u>.



- → KNN can be employed to estimate the posterior probabilities: Actually, posterior probabilities in an area around *x* are calculated as the ratio of samples within the area with class label ω_i .
 - \checkmark *n*: total number of patterns from all classes
 - \mathbf{i} *k*_{*i*}: number of patterns from class *i* in the area around *x*
 - k: total number of patterns from all classes in the area around x

$$p_n(\mathbf{x}, \omega_i) = \frac{\frac{k_i}{N}}{V}$$
$$P_n(\omega_i | \mathbf{x}) = \frac{p_n(\mathbf{x}, \omega_i)}{\sum_{j=1}^c p_n(\mathbf{x}, \omega_j)} = \frac{k_i}{k}$$



KNN Classification

\succ Classification of a pattern *x*,

- ✤ From the *n* training samples, the *k* nearest neighbors are identified regardless of their class, (where *k* is odd for a two-class problem). $\sum_{i=k}^{n} k_i = k$
- $\stackrel{\bullet}{\rightarrow}$ Let k_i the number of samples from class *i*,
- $\mathbf{\mathfrak{S}}$ Classify **x** to the class with the largest k_i !





Nearest Neighbor Classifier

>

The simplest KNN classifier is for k=1! In this case x assumes the class of its **nearest neighbor**! Is this a good classifier...?



NN classifier partitions the feature space as a Voronoi tessellation, where each cell assumes the label of the class within it.

Given infinite number of training samples, the classification error probability is bounded by twice the error of the Bayesian classifier (valid for small error probabilities).

Error Rates





NN Classifier

K-NN Classifier



- Partial distances
 - Use a subset r of d dimensions, to compute the distance of the unknown pattern from the training samples:

$$D_r(\mathbf{a},\mathbf{b}) = \left(\sum_{k=1}^r (a_k - b_k)^2\right)^{1/2}$$

- Search trees
 - Establish "search trees" where training samples are selectively linked so that, for the classification of a new pattern, distance computation is required only to some entry/root patterns and the patterns linked to them.
- Editing, pruning, or condensing
 - Prune patterns that are surrounded from patterns of the same class.
 - 1. Compute the Voronoi diagram of the training patterns
 - 2. For each pattern, if any of its neighbors is not of the same class, mark it.
 - 3. Prune non-marked patterns and compute the revised



Metrics and KNN Classification

Properties:

- ♥ Nonegativity: $D(\mathbf{a}, \mathbf{b}) \ge 0$
- \checkmark reflexivity: $D(\mathbf{a}, \mathbf{b}) = 0$ iff $\mathbf{a} = \mathbf{b}$
- \checkmark symmetry: $D(\mathbf{a},\mathbf{b}) = D(\mathbf{b},\mathbf{a})$
- ^t triangle inequali $D(\mathbf{a}, \mathbf{b}) + D(\mathbf{b}, \mathbf{c}) ≥ D(\mathbf{a}, \mathbf{c})$
- Euclidean Distance:
- Minkowski Metric:
- Manhatan Distance:

$$L_{2}(\mathbf{a}, \mathbf{b}) = \left(\sum_{k=1}^{d} (a_{k} - b_{k})^{2}\right)^{1/2}$$
$$L_{p}(\mathbf{a}, \mathbf{b}) = \left(\sum_{k=1}^{d} |a_{k} - b_{k}|^{p}\right)^{1/p}$$
$$L_{1}(\mathbf{a}, \mathbf{b}) = \sum_{k=1}^{d} |a_{k} - b_{k}|$$



Distance 1 from the center via the use of each one of the metrics *Lp*.

Chess-board Distance:

 $L_{\infty}(\mathbf{a},\mathbf{b}) = \max_{k=1,\dots,d} (|a_k - b_k|)$

k = 1



Metrics and KNN Classification



When the pattern space is transformed after multiplying each feature by a constant, distances in the transformed space may vary significantly from the initial ones. Evidently, this effects the performance of the KNN

It is of utmost importance to employ metrics that are invariant under basic transformations, such as translation (shift), scaling, rotation, line thickness, shear. Example case: optical character



 \succ

Tangent Distance



Let r transformations, αi .

Let **x**' a pattern.

Transformed pattern, *Fi(x';αi)*.

Tangent vector: $TV_i = F_i(x'; \alpha_i) - x'$.

Tangent matrix: *Tdxr=[TV1,..., TVr]*.

Tangent space: space defined by the r linearly independent tangent vectors TV_i that pass

from \boldsymbol{x}' . Forms a linear approximation of the space of transformed \boldsymbol{x}' .

Tangent distance:

 $D_t D_t \mathbf{A}_{\mathrm{tan}}(\mathbf{x}_{\mathbf{X}}, \mathbf{x}) \equiv \min_{\mathbf{w} \in \mathbf{W}} [\mathbf{x}_{\mathbf{W}}(\mathbf{x}_{\mathbf{T} \mathbf{W}}) \mathbf{T}_{\mathbf{A}}] - \mathbf{x} \|]$

Actually, it is Euclidean distance of \boldsymbol{x} from the tangent space of \boldsymbol{x}'

Tangent Distance







We find the optimal \mathbf{a} via iterative gradient descent. For gradient descent we need the derivative of the (squared) Euclidean distance. The Euclidean distance in Eq. 61 obeys

$$D^{2}(\mathbf{x}' + \mathbf{T}\mathbf{a}, \mathbf{x}) = \|(\mathbf{x}' + \mathbf{T}\mathbf{a}) - \mathbf{x}\|^{2},$$
(62)

and we compute the gradient with respect to the vector of parameters \mathbf{a} — the projections onto the tangent vectors — as

$$\nabla_{\mathbf{a}} D^2(\mathbf{x}' + \mathbf{T}\mathbf{a}, \mathbf{x}) = 2\mathbf{T}^t(\mathbf{x}' + \mathbf{T}\mathbf{a} - \mathbf{x}).$$
(63)

Thus we can start with an arbitrary **a** and take a step in the direction of the negative gradient, updating our parameter vector as

$$\mathbf{a}(t+1) = \mathbf{a}(t) - \eta \mathbf{T}^t (\mathbf{T}\mathbf{a}(t) + \mathbf{x}' - \mathbf{x}), \tag{64}$$



 \succ

Reduced Coulomb Energy (RCE) Networks

An RCE network, during training tunes the window width around each pattern according to its distance from the nearest pattern in a different class.

During training, each pattern initiates a new circle and circle radii are adapted so that they do not contain patterns of different classes.

Black circles represent class 1, pink circles class 2, whereas dark red areas represent ambiguous regions where no classification decision can be made.







Algorithm 4 (RCE training)

1	<u>begin</u> initialize $j = 0, n = \#$ patterns, $\epsilon = $ small param, $\lambda_m = $ max radius
2	$\underline{\mathbf{do}} \ j \leftarrow j+1$
3	train weight: $w_{jk} \leftarrow x_k$
4	find nearest pt not in ω_i : $\hat{\mathbf{x}} \leftarrow \arg \min D(\mathbf{x}, \mathbf{x}')$
	$\mathbf{x} ot \in \omega_i$
5	$\textbf{set radius: } \lambda_j \leftarrow Min[D(\hat{\mathbf{x}}, \mathbf{x}') - \epsilon, \lambda_m]$
6	$\underline{\mathbf{if}} \ \mathbf{x} \in \omega_i \ \underline{\mathbf{then}} \ a_{ic} \leftarrow 1$
7	$\underline{\mathbf{until}} \ j = n$
8	end

Algorithm 5 (RCE classification)