

Character templates learning for textual images recognition as an example of learning in structural recognition

Bogdan Savchynskyy, Olexander Kamotskyy
International Research and Training Center
of Information Technologies and Systems,
Acad. Glushkova Ave., 40, Kiev, Ukraine
bogdan@image.kiev.ua

Abstract

¹ Document recognition for digital libraries is characterized by high requirements to a recognition quality and processing of significant amount of single-type documents. So this is a perfect area for single-font approaches because they provide a smaller error rate comparing to multifont approaches and a learning of the font is carried out relatively rarely, because of significant amount of single-type documents.

Traditionally character templates learning is performed for separated characters on a basis of the set of character examples. It leads to recognition errors like in situations when closely placed parts of neighbouring characters are recognized as a single, separate character.

We propose another approach to character templates learning. Namely such templates must be constructed that the result of recognition of a text line image as a whole must match to a text string specified by a teacher. The approach guarantees that not only images of separate characters will be recognized correctly, but also the segmentation of the whole text image into characters will be performed without errors. So in our approach a learning sample consists not from labelled images of separated characters, but from text line images with corresponding text strings.

1. Introduction

Document recognition for digital libraries is characterized by high requirements to a recognition quality and processing of significant amount of single-type documents. So this is a perfect area for single-font approaches because:

- the single-font approaches provide smaller error rate comparing to multifont ones (see [2, 11]);

¹This paper was supported by the EU INTAS project PRINCESS 04-77-7347

- a learning of the font is carried out relatively rarely, because of significant amount of single-type documents.

A striking example of a single-font approach is a Document Image Decoding approach described in papers [1, 2, 3].

Traditionally character templates learning is performed for separated characters on a basis of the character examples set. It leads to recognition errors like in situations when closely placed parts of neighbouring characters are recognized as a single character.

We propose a novel approach to character templates learning. Namely such templates must be constructed that the result of recognition of a text line image as a whole must match to a text string specified by a teacher. The approach guarantees that not only images of separate characters will be recognized correctly, but also the segmentation of the whole text image into characters will be performed without errors. So in our approach a learning sample consists not from labelled images of separated characters, but from text line images with corresponding text strings.

This learning approach can be used only along with a proper recognition algorithm. Namely the algorithm must perform a segmentation and a recognition of text line image simultaneously via dynamic programming like it was discovered in [5] for the first time and like it was rediscovered in modern papers(see for ex. [2, 7, 8]).

Our approach of templates learning is significantly based on the approach of tuning autonomous stochastic automata parameters, described in [10]. Namely the idea of learning problem formulation and solution is the development of the ideas presented in [10].

This paper consists of six sections. We formulate main definitions and a statement of the templates learning problem in the second section. The solution of the problem is presented in the third and the fourth sections. The fifth section contains results of experimental testing and the last one is devoted to conclusions.

2. Learning problem formulation

First we will formulate the problem of a single text line image recognition, because its exact formulation is absolutely necessary for a formulation of the learning problem.

We call a rectangular subset of two-dimensional integer lattice *a field of vision*. Its height and width will be denoted H and W correspondingly:

$$T = \{(i, j) \mid i \in \overline{1, H}, j \in \overline{1, W}\}.$$

Elements of the field of vision T will be called *pixels*.

A denotation V will be used for the set of pixel brightnesses. This set will be later called *the set of signals*. An *image* is a function of the form $x : T \rightarrow V$. Brightness of the pixel t of the image x is denoted with $x(t)$ or x_t . Denotation V^T is used for the set of all possible images.

Let A be a finite set consisting of names of characters and $\alpha \in A$ be the name of a character. A *text line* is considered as a sequence of character names.

We consider the case when images of all the characters have the same height but their width depends just on the character's name. Thus the function $d : A \rightarrow \mathbb{N}$ will be regarded as a *function of a character width*. We consider this function to be known, so its values will not change during learning. We assume that the alphabet A contains also a special character *space* having width equal to one, so any spaces between characters and words can be modelled with a sequence of *spaces*.

A pair (α, j) , where α is a name from the alphabet A and $j \in (1, 2, \dots, W)$ is a column number of the field of vision, is called a *segment*. Each segment $s = (\alpha, j)$ determines a fragment $T(s)$ of the field of vision, which starts with a column number $j(s)$, has the name $\alpha(s)$, the height H and the width $d(\alpha(s))$:

$$T(s) = \{(i, j) \mid i \in \overline{1, H}, j \in \overline{j(s), j(s) + d(\alpha(s)) - 1}\}.$$

Thus segments are fragments of the field of vision labelled by names. The set of all possible segments is denoted with S .

Each text line $(\alpha_1, \dots, \alpha_L)$ of the length L , depicted in the image x , uniquely corresponds to the sequence of closely fitted segments (s_1, \dots, s_L) which cover the whole field of vision T :

$$\begin{cases} \alpha(s_l) = \alpha_l & \forall l \in \overline{1, L} \\ j(s_1) = 1 \\ j(s_{l+1}) = j(s_l) + d(\alpha(s_l)) & \forall l \in \overline{1, L-1} \\ j(s_L) + d(\alpha(s_L)) = W + 1 \end{cases} \quad (1)$$

We call the sequence $\bar{s} = (s_1, \dots, s_L)$ satisfying conditions (1) *a segmentation*. Its length will be denoted as $L(\bar{s})$. The set of all possible segmentations of the field of vision will be denoted with \overline{S} .

Let us introduce the set E of parameter vectors. Each parameter vector \bar{e} consists of templates of all characters: $\bar{e} = \{e(\alpha), \alpha \in A\}$. A template $e(\alpha)$ determines an ideal view of the character α . A template of a character can be, for example, a noise-free image of the character. Each template $e(\alpha)$ can be regarded as a vector itself. Thus we can interpret the parameter vector \bar{e} as a vector in a multidimensional vector space, meaning that coordinates of templates are coordinates of the vector.

We will assume that *a local dissimilarity function* of the form $f : E \times S \times V^T \rightarrow \mathbb{R}$ is defined. Its value $f_{\bar{e}}(s, x)$ determines a similarity between a segment s of an image x and the character $\alpha(s)$. The greater the similarity, the less value the function takes. A value of the function can be for example equal to a sum of squared differences of image pixels and pixels of character's $\alpha(s)$ template.

The similarity of the whole image x and a text line $(\alpha_1, \dots, \alpha_L)$, having a corresponding segmentation $\bar{s} = (s_1, \dots, s_L)$, is defined as a total similarity of all the segments, i.e. it is equal to $\sum_{l=1}^{L(\bar{s})} f_{\bar{e}}(s_l, x)$.

Recognition of an image x representing a text line consist in a search for the segmentation $\bar{s}^* = (s_1^*, s_2^*, \dots, s_L^*)$, which is the most similar to the image:

$$\bar{s}^* = \arg \min_{\bar{s} \in \overline{S}} \sum_{l=1}^{L(\bar{s})} f_{\bar{e}}(s_l, x). \quad (2)$$

The problem (2) can be solved via dynamic programming, which was done for the first time by Kovalevsky [5].

Generally speaking, the function $\sum_{l=1}^{L(\bar{s})} f_{\bar{e}}(s_l, x)$ takes its minimal value on some set of optimal segmentation \overline{S}^* . In this case we will consider that \bar{s}^* denotes any element of this set. A situation when the solution (2) must be unique will be distinguished with a sign \doteq .

Further we will assume that a procedure for the search of an optimal segmentation \bar{s}^* is available. This procedure solves the problem (2) for any triple $P = \langle A, d, f_{\bar{e}} \rangle$, which is the set of characters A , the width function d and a local disparity function $f_{\bar{e}}$, defined up to the vector of parameters \bar{e} . We will call this procedure *a recognition algorithm* $\bar{s}_P^*(x, \bar{e})$. Its implementation is not important for us, but we will assume that it can track the cases when the solution \bar{s}^* is not unique.

Templates learning algorithm, corresponding to the recognition problem (2), consists in choosing a parameter vector \bar{e} (i.e. choosing character templates) for a given triple $P = \langle A, d, f_{\bar{e}} \rangle$.

Let $X^o = \{x_r^o \mid r \in \overline{1, R}\}$ be a learning set of images and $\overline{S}^o = \{\bar{s}_r^o \mid r \in \overline{1, R}\}$ be the set of corresponding segmentations.

Problem 1 (*character templates learning*) consists in search for such a parameter vector \bar{e}^* that the result of

recognition of each learning image x_r^o must be the only corresponding segmentation \bar{s}_r^o :

$$\bar{s}_r^o \stackrel{!}{=} \arg \min_{\bar{s} \in \bar{S}} \sum_{l=1}^{L(\bar{s})} f_{\bar{e}^*}(s_l, x_r^o) \quad \forall r \in \overline{1, R}.$$

Complexity of the problem is determined by the fact that the minimum must be found on the set of segmentations \bar{S} , and the cardinal number of this set is proportional to the number of *all possible text lines*, which can be depicted in the field of vision.

3. Solution of the learning problem

We propose the solution of the learning problem 1 in the case, when a dissimilarity function is linear with respect to parameters $\bar{e} = \{e_1, \dots, e_K\}$, i.e. it can be presented in the form:

$$f_{\bar{e}}(s, x) = \sum_{k=1}^K e_k \cdot \varphi_k(s, x) = \langle \bar{e}, \bar{\varphi}(s, x) \rangle. \quad (3)$$

Here $\bar{\varphi}(s, x) = \{\varphi_1(s, x), \dots, \varphi_K(s, x)\}$ is a vector, which directly depends on an image x and segment s .

Let us denote a learning sample consisting of images and corresponding segmentations as Π :

$$\Pi = \{(x_r^o, \bar{s}_r^o) \mid r \in \overline{1, R}\}.$$

Now the learning problem 1 takes the form: find such \bar{e}^* , that

$$\bar{s}^o \stackrel{!}{=} \arg \min_{\bar{s} \in \bar{S}} \sum_{l=1}^{L(\bar{s})} \sum_{k=1}^K e_k^* \varphi_k(s_l, x^o) \quad \forall (x^o, \bar{s}^o) \in \Pi. \quad (4)$$

An equality (4) means, that the system of inequalities:

$$\left\{ \begin{array}{l} \sum_{l=1}^{L(\bar{s}^o)} \sum_{k=1}^K e_k^* \varphi_k(s_l^o, x^o) < \sum_{l=1}^{L(\bar{s})} \sum_{k=1}^K e_k^* \varphi_k(s_l, x^o) \\ \forall \bar{s} \neq \bar{s}^o, \end{array} \right.$$

must hold for each pair $(x^o, \bar{s}^o) \in \Pi$.

These inequalities are equivalent to

$$\left\{ \begin{array}{l} \sum_{l=1}^{L(\bar{s})} \sum_{k=1}^K e_k^* \varphi_k(s_l, x^o) - \sum_{l=1}^{L(\bar{s}^o)} \sum_{k=1}^K e_k^* \varphi_k(s_l^o, x^o) > 0 \\ \forall \bar{s} \neq \bar{s}^o. \end{array} \right.$$

After grouping components of parameter vector \bar{e}^* and denoting all coefficients of e_k^* as $\varphi'_k(\bar{s}, \bar{s}^o, x^o)$ we will get:

$$\left\{ \begin{array}{l} \sum_{k=1}^K e_k^* \varphi'_k(\bar{s}, \bar{s}^o, x^o) > 0 \\ \forall \bar{s} \neq \bar{s}^o. \end{array} \right. \quad (5)$$

Thus the learning problem 1 consists in solving (5) with respect to the vector $\bar{e}^* = (e_1^*, \dots, e_K^*)$.

Since the linear system (5) contains $|\bar{S}| - 1$ inequalities, the classical linear optimization methods cannot be used for its solution, because their complexity directly depends on the number of inequalities. But the solution of (5) can be found by means of linear discriminant analysis, containing algorithms which complexity does not depend on the cardinal number of the system, but on other its properties.

We used algorithms of Kozinec and perceptron (see, for ex. [10]) for the solution of the system (5). These algorithms are iterative and at each iteration parameter vector is corrected on the basis of a single unfulfilled inequality of the system (5). Thus they need only a method of effective search for such an inequality. The proof of a finite step convergence of these algorithms can be found in [10].

In our case, if any inequality of the system (5) does not hold for the current value of parameters \bar{e} , then it means that such a learning pair (x^o, \bar{s}^o) exists that equality (4) does not hold and in its turn it means that the learning image x^o is recognized incorrectly. And vice versa, in order to check whether (5) is fulfilled we are to check whether (4) is fulfilled, in other words we must recognize all the learning images $x^o \in X^o$.

Thus Kozinec and perceptron algorithms iteratively refine parameter vector \bar{e} by means of recognition procedure $\bar{s}_P^*(x, \bar{e})$. Under given $P = \langle A, d, f_{\bar{e}} \rangle$ the procedure finds the solution of the recognition problem (2). This solution determines such an inequality of the system (5) which is not fulfilled for a current value of parameters \bar{e} .

Algorithm 1 Kozinec algorithm for solving (5)

- 1: Choose any vector $\bar{e} \neq 0$ from the convex hull of the vectors set $\{\bar{\varphi}'\}$. For example, we can take it equal to $\bar{\varphi}'(\bar{s}, \bar{s}_r^o, x_r^o) \neq 0, \bar{s} = \bar{s}_P^*(x_r^o, 0)$ for any $r \in \overline{1, R}$.
- 2: Find an inequality which does not hold for the current value of \bar{e} :

$$\text{find } r \in \overline{1, R}, \text{ such that } \bar{s}_r^o \stackrel{!}{\neq} \bar{s}_P^*(x_r^o, \bar{e}).$$
- 3: If there is no such r then **end**. Parameter vector \bar{e}^* is found.
- 4: Else calculate new parameter vector as a perpendicular dropped to a segment connecting vectors \bar{e} and $\bar{\varphi}' = \bar{\varphi}'(\bar{s}_P^*(x_r^o, \bar{e}), \bar{s}_r^o, x_r^o)$:

$$\bar{e} := k \cdot \bar{e} + (1 - k) \cdot \bar{\varphi}', \quad (6)$$

$$\text{where } k = \frac{(\bar{\varphi}', \bar{\varphi}')}{(\bar{e}, \bar{e}) + 2(\bar{e}, \bar{\varphi}') + (\bar{\varphi}', \bar{\varphi}')}. \quad (7)$$

Proceed to step 2.

The perceptron algorithm looks quite similar to Kozinec one, but a bit simpler. Applied to the system (5) it takes the form:

Algorithm 2 Perceptron algorithm for solving (5)

- 1: Assign $\bar{e} = 0$.
- 2: Find an inequality which does not hold for the current value of \bar{e} :
find $r \in \overline{1, R}$, such that $\bar{s}_r^o \neq \bar{s}_P^*(x_r^o, \bar{e})$.
- 3: If there is no such r then **end**. Parameter vector \bar{e}^* is found.
- 4: Else calculate new parameter vector \bar{e} by the formula:

$$\bar{e} := \bar{e} + \bar{\varphi}'(\bar{s}_P^*(x_r^o, \bar{e}), \bar{s}_r^o, x_r^o).$$

Proceed to step 2.

To finish the consideration of the learning problem we must only concretize a local dissimilarity function $f_{\bar{e}}$. Its form determines a way of calculating vector $\bar{\varphi}'(\bar{s}, \bar{s}^o, x)$, which is used in algorithms 1 and 2.

4. Selection of the local dissimilarity function

In practice an input image of the character α is usually regarded as an ideal one $\chi(\alpha)$ distorted by some noise r . Most typically, the noise is considered as an additive one with spherically symmetric probability distribution $p_n(r)$:

$$x = \chi(\alpha) + r. \quad (8)$$

Here x is the distorted character image. In this case correlation approach [6, 4] to recognition can be used and a likelihood function can be used as a local dissimilarity function $f_{\bar{e}}$. Under assumption that a probability density $p_n(r)$ of the noise is a monotone decreasing function $f(\cdot)$ of the sum of its squared components:

$$p_n(r) = f(r^2), \quad (9)$$

function $f_{\bar{e}}$ takes the form of the Euclidean distance between vectors $x = \{x_{i,j} \mid (i,j) \in T(s)\}$ and $\bar{e} = \chi(\alpha) = \{\chi_{i,j}^\alpha \mid (i,j) \in T(s)\}$:

$$f_{\bar{e}}(s, x) = \sum_{\substack{i \in \overline{1, H} \\ j=1, d(\alpha(s))}} (\chi_{i,j}^{\alpha(s)} - x_{i,j(s)+j})^2. \quad (10)$$

Formula (10) can be presented in another form:

$$\begin{aligned} f_{\bar{e}}(s, x) &= \sum_{\substack{i \in \overline{1, H} \\ j=1, d(\alpha(s))}} (\chi_{i,j}^{\alpha(s)} - x_{i,j(s)+j})^2 = \\ &= \sum_{\substack{i \in \overline{1, H} \\ j=1, d(\alpha(s))}} (x_{i,j(s)+j})^2 + \\ &+ \sum_{\substack{i \in \overline{1, H} \\ j=1, d(\alpha(s))}} \left((\chi_{i,j}^{\alpha(s)})^2 - 2\chi_{i,j}^{\alpha(s)} \cdot x_{i,j(s)+j} \right) = \\ &= \tilde{f}_{\bar{e}}(s, x) + \sum_{\substack{i \in \overline{1, H} \\ j=1, d(\alpha(s))}} (x_{i,j(s)+j})^2, \end{aligned} \quad (11)$$

where $\tilde{f}_{\bar{e}}(s, x)$ is used to denote a quantity $\sum_{\substack{i \in \overline{1, H} \\ j=1, d(\alpha(s))}} \left((\chi_{i,j}^{\alpha(s)})^2 - 2\chi_{i,j}^{\alpha(s)} \cdot x_{i,j(s)+j} \right)$.

To avoid piling indexes we will use a nonstandard designation starting from this section. Namely, we will denote with $s \in \bar{s}$ an enumeration of all segments of segmentation \bar{s} . For example the designation $\sum_{l=1}^{L(\bar{s})} s_l$ is exactly the same as the designation $\sum_{s \in \bar{s}} s$.

Basing on the specified transformation (11) we rewrite a learning problem 1:

$$\begin{aligned} \bar{s}_r^o \stackrel{!}{=} \arg \min_{\bar{s} \in \bar{S}} \sum_{s \in \bar{s}} f_{\bar{e}}(s, x) = \\ = \arg \min_{\bar{s} \in \bar{S}} \left(\sum_{s \in \bar{s}} \tilde{f}_{\bar{e}}(s, x) + \sum_{s \in \bar{s}} \sum_{\substack{i \in \overline{1, H} \\ j=1, d(\alpha(s))}} (x_{i,j(s)+j})^2 \right) = \\ = \arg \min_{\bar{s} \in \bar{S}} \left(\sum_{s \in \bar{s}} \tilde{f}_{\bar{e}}(s, x) + \sum_{\substack{i \in \overline{1, H} \\ j \in \overline{1, W}}} (x_{i,j})^2 \right). \end{aligned} \quad (12)$$

The second item in the formula (12) depends only on the input image x , not on a segmentation \bar{s} . Thus the item does not affect the argument of minimum search. So the formula (12) can be rewritten in the form:

$$\begin{aligned} \bar{s}_r^o \stackrel{!}{=} \arg \min_{\bar{s} \in \bar{S}} \sum_{s \in \bar{s}} \tilde{f}_{\bar{e}}(s, x) = \\ = \arg \min_{\bar{s} \in \bar{S}} \sum_{s \in \bar{s}} \sum_{\substack{i \in \overline{1, H} \\ j=1, d(\alpha(s))}} \left((\chi_{i,j}^{\alpha(s)})^2 - 2\chi_{i,j}^{\alpha(s)} \cdot x_{i,j(s)+j} \right). \end{aligned} \quad (13)$$

The target function in the formula (13) is nonlinear with respect to parameters $\chi_{i,j}^\alpha$, thus the problem can not be solved directly in the same manner like it was done in previous section.

For solving the problem we will substitute each function of the form $(\chi_{i,j}^\alpha)^2 - 2\chi_{i,j}^\alpha \cdot x_{i,j(s)+j}$ such, that the only parameter $\chi_{i,j}^\alpha$ must be learned, with a function $\dot{e}_{i,j}^\alpha \cdot x_{i,j(s)+j} + e_{i,j}^\alpha$, which is linear in two parameters $\dot{e}_{i,j}^\alpha$ and $e_{i,j}^\alpha$. Obviously the set of such functions includes the set of the first-type functions and coincides with it in the case $e_{i,j}^\alpha \geq 0$ and $\dot{e}_{i,j}^\alpha = \pm 2 \cdot \sqrt{e_{i,j}^\alpha}$. But note that the introduced parameters $e_{i,j}^\alpha$ and $\dot{e}_{i,j}^\alpha$ do not have such a clear physical meaning like $\chi_{i,j}^\alpha$, which defines a colour of a given image pixel.

In this case the problem (13) takes the form:

$$\bar{s}_r^o \stackrel{!}{=} \arg \min_{\bar{s} \in \bar{S}} \sum_{s \in \bar{s}} \sum_{\substack{i \in \overline{1, H} \\ j = \overline{1, d(\alpha(s))}}} \dot{e}_{i,j}^\alpha \cdot x_{i,j(s)+j} + e_{i,j}^\alpha, \quad (14)$$

a parameter vector \bar{e} is composed from all the coefficients:

$$\bar{e} = \left\{ \dot{e}_{i,j}^\alpha, e_{i,j}^\alpha \mid \alpha \in A, i \in \overline{1, H}, j = \overline{1, d(\alpha)} \right\},$$

and the vector $\bar{\varphi}'(\bar{s}, \bar{s}^o, x)$ from the problem (14) consists of components $\dot{\varphi}_{i,j}^{\prime\alpha}$ and $\varphi_{i,j}^{\prime\alpha}$, which depend on the colours of an input image x :

$$\bar{\varphi}'(\bar{s}, \bar{s}^o, x) = \left\{ \dot{\varphi}_{i,j}^{\prime\alpha}, \varphi_{i,j}^{\prime\alpha} \mid \alpha \in A, i \in \overline{1, H}, j = \overline{1, d(\alpha)} \right\}, \quad (15)$$

where

$$\begin{aligned} \dot{\varphi}_{i,j}^{\prime\alpha} &= \sum_{\substack{s \in \bar{s} \\ \alpha(s) = \alpha}} x_{i,j(s)+j}^o - \sum_{\substack{s \in \bar{s}^o \\ \alpha(s) = \alpha}} x_{i,j(s)+j}^o, \\ \varphi_{i,j}^{\prime\alpha} &= \sum_{\substack{s \in \bar{s} \\ \alpha(s) = \alpha}} 1 - \sum_{\substack{s \in \bar{s}^o \\ \alpha(s) = \alpha}} 1. \end{aligned}$$

Experiments with Kozinec and perceptron algorithms for learning vector containing parameters $\dot{e}_{i,j}^\alpha, e_{i,j}^\alpha$ in the problem (14) showed, that on real images learning takes such an amount of time which is too much for practical use of these algorithms. Significant amount of the algorithm iterations can be explained with such empirical considerations. Different components of the vector $\bar{\varphi}'(\bar{s}, \bar{s}^o, x)$ affect the values of parameter vector in nonuniform way. For example, each increase of coefficient $\dot{e}_{i,j}^\alpha$ by value $x_{i,j(s)+j}$, which is the colour of an image at some point and takes value, for example, in the interval from 0 to 255, requires adequate learning of parameter $e_{i,j}^\alpha$. But this operation takes up to 255 cycles of perceptron algorithm, in each cycle $e_{i,j}^\alpha$ is increased by 1.

A uniform influence on different components of parameter vector can be provided with an orthonormalized basis

for functions $\dot{e}_{i,j}^\alpha \cdot x_{i,j(s)+j} + e_{i,j}^\alpha$. These functions can be represented in such a form:

$$\bar{s}_r^o \stackrel{!}{=} \arg \min_{\bar{s} \in \bar{S}} \sum_{s \in \bar{s}} \sum_{\substack{i \in \overline{1, H} \\ j = \overline{1, d(\alpha(s))}}} \left(\dot{e}_{i,j}^\alpha \cdot \psi_1(x_{i,j(s)+j}) + e_{i,j}^\alpha \cdot \psi_0(x_{i,j(s)+j}) \right), \quad (16)$$

where functions $\psi_i(x)$ are orthonormalized Chebyshev polynomials. Each $\psi_i(x)$ is an i -th order polynomial of the signal value x . Orthonormalization means that such restrictions hold:

$$\sum_{x \in V} \psi_i(x) \cdot \psi_j(x) = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}. \quad (17)$$

In this basis the components (15) of the vector $\bar{\varphi}'(\bar{s}, \bar{s}^o, x)$ take a form:

$$\begin{aligned} \dot{\varphi}_{i,j}^{\prime\alpha} &= \sum_{\substack{s \in \bar{s} \\ \alpha(s) = \alpha}} \psi_1(x_{i,j(s)+j}) - \sum_{\substack{s \in \bar{s}^o \\ \alpha(s) = \alpha}} \psi_1(x_{i,j(s)+j}^o), \\ \varphi_{i,j}^{\prime\alpha} &= \sum_{\substack{s \in \bar{s} \\ \alpha(s) = \alpha}} \psi_0(x_{i,j(s)+j}) - \sum_{\substack{s \in \bar{s}^o \\ \alpha(s) = \alpha}} \psi_0(x_{i,j(s)+j}^o). \end{aligned}$$

Due to the use of such a basis for a local dissimilarity function processing time of algorithms 1 and 2 decreased by two-five times on real images.

5. Experimental testing

The proposed approach was tested on real and artificial images. Two other approaches have been used to compare with our one.

The first approach follows from the maximum likelihood estimation of templates under condition that an image noise is described with formulae (8) and (9). The solution of the estimation problem comes to averaging of learning images' fragments corresponding to the same character. The use of obtained templates as parameter vector of recognition algorithm with the dissimilarity function of the form (10) will be called a *method of templates averaging*.

The second approach is a multifont one. We took a popular commercial OCR program **FineReader** as a representative of this approach.

We present the results of testing learning algorithms 1 and 2 on real images like one presented in the fig. 1. The procedure of images creation was the following: computer text was printed using a dot-matrix printer with an old printing band that caused a nonuniform noise in the image. To avoid errors connected with text lines detection and skew

Table 1. Percentage of errors committed by the proposed algorithm (learning), algorithm of templates averaging and multifont commercial OCR program FineReader.

algorithm	Percentage of errors
learning	3,8
averaging	13
FineReader	17,5

all the image text lines have been separated and unskewed manually.

Then for each line that was used for learning a corresponding segmentation (i.e. names and exact horizontal position of each character) was entered by a teacher. Producing such a segmentation manually is a hard work, so in our experiments it was done automatically. For an automatic creation of the segmentation just a text line without information about character positions was entered by the teacher. The character positions were then produced with an unsupervised learning algorithm like one presented in [9].

A Kozinec algorithm with the local dissimilarity function of the form defined with formulae (16) and (17) was used to learn character templates.

We used 100 segmented text lines for testing. Ten on these lines have been used to learn templates. About thirty lines including all the learning lines labelled by ticks are stacked in a single image in the fig. 1. In the fig. 2 some of input lines and lines containing recognition results are stacked together. Learning lines are also ticked off. It is easy to see that lines used for learning are recognized without errors as it is demanded by the problem 1 formulation.

The results of testing are presented at the table 1. Lines used for learning have not been considered. You can see that the proposed approach surpasses template averaging approach more then in three times and multifont recognition software more then in four times.

The approach presented in the article has been tested also on artificial images. These images were of two types. Images of the first type were noised according to formulae (8) and (9). But the difference in recognition results between the proposed approach and the approach of averaging templates for such images was negligible.

Images of the second type were noised in completely another manner. Namely parts of character images were shifted or deformed such that the averaging gave incorrect recognition results in contrast to our approach, that gave error-free results.

The simplest example of such an artificial tests is presented in the fig. 3 and 4. After learning on the image 3 with averaging approach the image 4 was recognized as CCE in

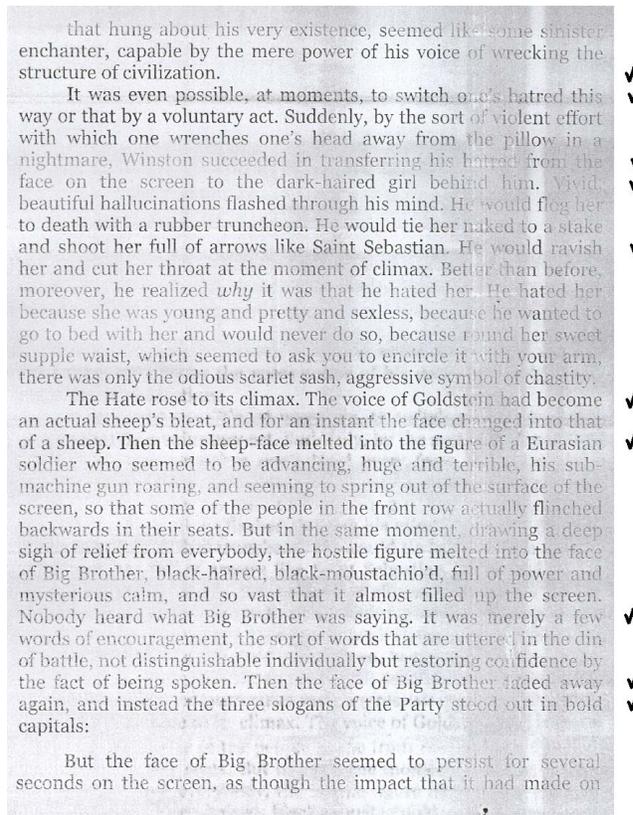


Figure 1. Some of image lines used to test the proposed learning approach. Lines used for learning are ticked off.

contrast to our approach that led to error-free recognition.

Quantitatively but not qualitatively more complicated example is presented in the fig. 5 and 6. Degraded characters are labelled with a grey background. Recognition results after learning with the averaging approach are presented in the fig. 6, while our approach gives error-free results.

6. Conclusions

Theoretical and practical analysis of the proposed approach argues:

- the proposed approach leads to approximately the same recognition quality as an averaging approach in a case of known image noise, described with formulae (8) and (9);
- in cases when noise distribution is far from being described with formulae (8) and (9) the proposed ap-

that hung about his very existence, seemed like some sinister
 that hung about his very existence, seemed like some sinister
 enchanter, capable by the mere power of his voice of wrecking the
 enchanter, capable by the mere power of his voice of wrecking the
 structure of civilization. ✓
 It was even possible, at moments, to switch one's hatred this
 way or that by a voluntary act. Suddenly, by the sort of violent effort ✓
 way or that by a voluntary act. Suddenly, by the sort of violent effort
 with which one wrenches one's head away from the pillow in a
 with which one wrenches one's head away from the pillow in a
 nightmare, Winston succeeded in transferring his hatred from the
 nightmare, Winston succeeded in transferring his hatred from the
 face on the screen to the dark-haired girl behind him. ✓
 face on the screen to the dark-haired girl behind him. ✓
 Vivid, beautiful hallucinations flashed through his mind. He would flog her
 beautiful hallucinations flashed through his mind. He would flog her
 to death with a rubber truncheon. He would tie her naked to a stake ✓
 to death with a rubber truncheon. He would tie her naked to a stake
 and shoot her full of arrows like Saint Sebastian. He would ravish ✓
 and shoot her full of arrows like Saint Sebastian. He would ravish
 her and cut her throat at the moment of climax. Better than before,
 her and cut her throat at the moment of climax. Better than before,
 her and cut her throat at the moment of climax. Better than before,
 moreover, he realized why it was that he hated her. He hated her ✓
 moreover, he realized why it was that he hated her. He hated her
 because she was young and pretty and sexless, because he wanted to
 because -he was young and pretty and sexless, because he wanted to
 go to bed with her and would never do so, because round her sweet
 go to bed with her and would never do so, because round her sweet
 supple waist, which seemed to ask you to encircle it with your arm,
 sunn.e waist, which seemed to ask you to encircle it with your arm,
 there was only the odious scarlet sash, aggressive symbol of chastity.
 there was only the odious scarlet sash, aggressive symeo-of chastity, ✓

Figure 2. The result of recognition of the first half of image lines in the fig. 1. Lines used for learning ticked off and according to problem 1 formulation are recognized without errors.

proach surpasses the averaging approach and other approaches based on a maximum likelihood estimation of character templates;

- learning time can be significantly decreased if the proposed learning algorithm starts from a good approximation of the character templates received for example by averaging approach;
- a learning quality in the approach depends not on the number of samples but on that how learning samples cover all the cases that appear during recognition phase. Good recognition results were achieved in the example presented in the fig. 1, because lines from learning sample contained as character images from

CEEEEEEE

Figure 3. A very simple example of artificial image. The last character E is degraded.

CEE

Figure 4. After learning on the image 3 with averaging approach this image was recognized incorrectly as CEE in contrast to our learning approach which led to a error-free result.

мова - це не просто спосіб спілкування,
 а щось більш значуще. мова - це всі глибинні
 пласти духовного життя народу, його
 історична пам'ять, найцінніше надбання віків,
 мова - це ще й музика, мелодика, фарби,
 буття, сучасна, художня, інтелектуальна і
 мисленнева діяльність народу.

Figure 5. An artificial image for learning. Degraded characters labelled with a grey background.

грандіозні речі робляться грандіозними
 засобами, одна природа робить велике даром
 ↓↓
 грайдідзн' реч' рдблч_ься_грайдідзиимн_
 зесдбаки,_ддиа_понрдда_рдбить'ітелике_оардкі

Figure 6. An input image is at the top, the result of recognition is at the bottom. The recognition result was achieved after learning with the averaging approach on the image 6.

the left, relatively "good" part of the whole image, so and character images from the right, relatively "bad" one;

- the basis of the local dissimilarity function must be chosen carefully. In our case the use of orthonormalized Chebyshev polynomials enabled a possibility to work with real images;
- the proposed learning approach can be used not only for text recognition, but also for a wide range of structural recognition problems. The only condition is that the problem must be formulated in a form (2) with a linear dissimilarity function (3) and recognition problem itself must be solvable too.

References

- [1] P. A. Chou and G. E. Kopec. A stochastic attribute grammar model of document production and its use in document image decoding. In H. B. L. Vincent, editor, *Document Recognition II*, volume 2422 of *SPIE Proc.*, pages 66–73, 1995.
- [2] G. E. Kopec and P. A. Chou. Document image decoding using markov source models.
- [3] G. Kopec and L. M. Document-specific character template estimation. In *IS&T/SPIE 1996 Intl. Symposium on Electronic Imaging: Science & Technology*, San Jose, CA, Jan. 27–Feb. 2, 1996.
- [4] V. Kovalevsky. Korreliacionnyj metod raspoznavaniya izobrazhenij; in russian (Correlation method of image recognition). *Journal vychislitelnoj matematiki i matematicheskoi fiziki*, 2(4), 1962.
- [5] V. Kovalevsky. Optimalnyj algoritm raspoznavania nekotorych posledovatelnostej izobrazhenij; in russian (Optimal algorithm recognizing some sequences of images). *Kibernetika*, (4):75–80, 1967.
- [6] V. Kovalevsky. *Metody optimalnykh reshenij v raspoznavanii izobrazhenij; in russian (Methods of optimal solutions in image recognition)*. Nauka, Moskow, 1976.
- [7] U. Marti and H. Bunke. Handwritten sentence recognition. In *Proceedings of 15th International Conference on Pattern Recognition*, volume 3, pages 467–470, Barcelona, 2000.
- [8] R. Plamondon and S. N. Srihari. On-line and off-line handwriting recognition: A comprehensive survey. *IEEE Trans. on PAMI*, 22(1):63–84, January 2000.
- [9] P. Sarkar, H. S. Baird, and X. Zhang. Training on severely degraded text-line images. In *IAPR 7th International Conference on Document Analysis and Recognition (ICDAR03)*, pages 38–43, Edinburgh, Scotland, August 2003.
- [10] M. I. Schlesinger and V. Hlaváč. *Ten lectures on statistical and structural pattern recognition*. Kluwer Academic Publishers, Dordrecht/Boston/London, 2002.
- [11] Y. Xu and G. Nagy. Prototype extraction and adaptive OCR. *IEEE Trans. Pattern Anal. Mach. Intell.*, 21(12):1280–1296, 1999.