

УДК 004.8:004.93

Н.А. Костромов, И.В. Бекетова, С.Л. Каратеев

ФГУП «Государственный научно-исследовательский институт авиационных систем», г. Москва, Россия
Россия, 125319, г. Москва, ул. Викторенко, 7

Формирование пространства признаков и их извлечение из изображений для систем индексации документов

N.A. Kostromov, I.V. Beketova, S.L. Karateev

*The Federal State Unitary Enterprise "State Research Institute of Aviation Systems",
c. Moscow, Russia
Russia, 125319, Moscow, Viktorenko street, 7.*

Formation of Feature Spaces and their Extraction from the Images for Documents Indexing

Н.А. Костромов, И.В. Бекетова, С.Л. Каратеев

ФГУП «Державний науково-дослідний інститут авіаційних систем»,
Росія, м. Москва
Росія, 125319, м. Москва, вул. Вікторенко, 7

Формування простору ознак і їх добування із зображень для систем індексації документів

Предложены оригинальные алгоритмы анализа изображения для формирования вектора признаков. Вектор признаков сформирован из геометрических признаков и составленных на основе системы правил комбинированных признаков. Вектор признаков используется для поиска областей со специфической графической информацией на изображениях документов. Моделирование на тестовой выборке изображений документов показало устойчивость и эффективность работы алгоритмов.

Ключевые слова: индексация, извлечение признаков, машинное обучение

The original algorithms of image analysis for formation of feature vector are proposed. Feature vector is formed from the both type of features, geometric features and secondary complex features compiled on the basis of the logical rules. Feature vector is used to find areas with specific information on the graphic images of documents. Simulation on a test set of images of documents has shown the stability and efficiency of the algorithms.

Key Words: indexing, feature extraction, machine learning.

Запропоновані оригінальні алгоритми аналізу зображень для формування вектора ознак. Вектор ознак сформований з геометричних ознак і складених на основі системи правил комбінованих ознак. Вектор ознак використовується для пошуку областей зі специфічною графічною інформацією на зображеннях документів. Моделювання на текстовій вибірці зображень документів показало стійкість й ефективність роботи алгоритмів.

Ключові слова: індексація, добування ознак, машинне навчання

Введение

При использовании систем автоматического анализа документов возникают задачи обнаружения документов, содержащих специфическую графическую информа-

цию. Существует ряд документов, в которых содержание подкрепляется наглядными образами при помощи графических моделей в форме структурных схем или блок-схем. Сложность задачи обнаружения подобных документов состоит в том, что в отличие от изображений таблиц или других стандартных форм изображения графических моделей рассматриваемого типа, как правило, не имеют стандартной структуры, специфических ключевых слов, символов и графических объектов, по которым можно определить, что данное графическое изображение является структурной схемой или блок-схемой.

Структурные схемы и блок-схемы отличаются широким разнообразием состава и параметров отображающих их геометрических объектов – геометрических фигур, линий и графических символов. Примеры изображений документов, содержащих блок-схемы и структурные схемы, приведены на рис. 1.

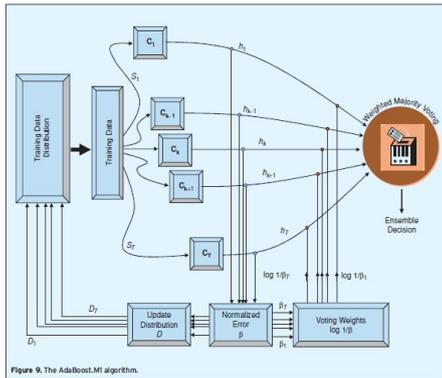


Figure 9. The AdaBoost.M1 algorithm.

where E is the ensemble error [25]. Since $\epsilon_i < 1/2$, E is guaranteed to decrease with each new classifier. In most practical cases, the error decreases very rapidly in the first few iterations, and approaches zero as new classifiers are added. While this is remarkable on its own account, the surprising resistance of AdaBoost to overfitting is particularly noteworthy. Overfitting is a commonly observed phenomenon where the classifier performs poorly on test data, despite achieving a very small training error. Overfitting is usually attributed to memorizing the data, or learning the noise in the data. As a classifier's capacity increases (for example, with the complexity of its architecture), so does its tendency to memorize the training data and/or learn the noise in the data. Since the capacity of an ensemble increases with each added classifier, one would expect AdaBoost to suffer from overfitting, particularly if its complexity exceeds what is necessary to learn the underlying data distribution. Yet, AdaBoost performance usually levels off with increasing number of classifiers with no indication of overfitting.

Schapire et al. later provided an explanation to this phenomenon based on the so-called margin theory [52]. The details of margin theory are beyond the scope of this paper; however, the margin of an instance x , in loose terms, is its distance from the decision boundary. The further away an instance is from the boundary, the larger its margin, and hence the higher the confidence of the classifier in correctly classifying this instance. Margins are usually described within the context of support vector machines (SVM) type classifiers, which are based on the idea of maximizing the margins between the instances and the decision boundary. In fact, SVMs find those instances, called support vectors, that lie closest to the decision boundary. The support vectors are said to define the margin that separates the classes. The concept of a margin is conceptually illustrated in Figure 10. By using a slightly different definition of a margin—essentially modified for AdaBoost—Schapire et al. show that AdaBoost also boosts the margins, that is, it finds the decision boundary that is further away from the instances of all classes. In the context of AdaBoost, the margin of an

THIRD QUARTER 2006

IEEE CIRCUITS AND SYSTEMS MAGAZINE 31

V. Anish Natarajan, K.M.M. Prabhak / Microprocessors and Microsystems 27 (2013) 513–521

517

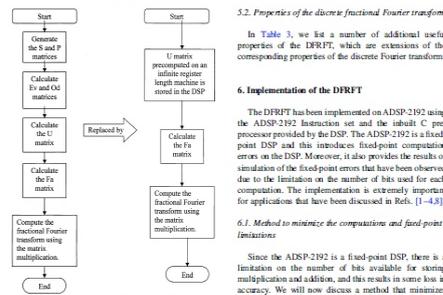


Fig. 4. Flow chart of the original and revised method to compute DFRFT.

(e) The FRFT operator matrix is defined as

$$F^{\alpha}[m,n] = \sum_{k=0}^{N-1} U_k[m] e^{j2\pi k n \alpha} U_k[n], \quad (15)$$

$\mu = \{0, 1, 2, \dots, N-2, (N-N_k)\}$.

(f) The FRFT of order α is then calculated by a matrix multiplication

$$f_x = F^{\alpha} f. \quad (16)$$

The flow chart for the implementation of the DFRFT is shown in Fig. 4. An $O(N \log N)$ algorithm for the digital computation of FRFT is given in Ref. [7], which computes the samples of the continuous FRFT and hence is strictly not the discrete FRFT (DFRFT).

Table 3 Properties of DFRFT

Order	FRFT	DFRFT
0	$f(x)$	$f(x)$
1	$f(x) + j\alpha f'(x)$	$f(x) + j\alpha f'(x)$
2	$f(x)$	$f(x)$
3	$f(x) - j\alpha f'(x)$	$f(x) - j\alpha f'(x)$
4	$f(x)$	$f(x)$
5	$f(x) + j\alpha f'(x)$	$f(x) + j\alpha f'(x)$
6	$f(x)$	$f(x)$
7	$f(x) - j\alpha f'(x)$	$f(x) - j\alpha f'(x)$
8	$f(x)$	$f(x)$

Рисунок 1 – Примеры изображений документов, содержащих блок-схемы и структурные схемы

Важно определить наиболее существенные признаки рассматриваемой графической модели. Данная графическая модель строится при помощи ограниченного набора графических элементов (линии, прямоугольники, ромбы, овалы и т.д.), часто дополняемых короткими текстовыми строками или отдельными графическими символами. Часто изображения содержат шумы и искажения, вызванные процессами получения изображений (сканирование, видеосъёмка) и сжатия полученных изображений.

Достаточно полный обзор современных методов анализа структуры документов приведён в [1], [2]. Методы, используемые для обнаружения изображений графических объектов, можно разделить на три основные категории:

- обнаружение на основе априорной информации о структуре;
- эвристические методы;
- статистические или оптимизационные методы.

Для методов обнаружения на основе априорной информации, как правило, характерны удовлетворительные показатели обнаружения, однако обобщающие способности методов данного типа ограничены.

Эвристические алгоритмы используют заранее определённые наборы логических или предварительно определённых синтаксических правил для получения решения об обнаружении. Наборы эвристических правил обычно строятся на базе системы локальных признаков. Использование эвристик нередко требует сложной пост-обработки для окончательного принятия решений.

Статистические или оптимизационные алгоритмы, как правило, обходятся без априорной информации или используют свободные параметры, значения которых формируются в процессе обучения [3], [4].

Структура алгоритма

Общая структура алгоритмов представлена на рис. 2. Таким образом, выделение признаков можно разделить на следующие этапы:

- 1) выделение текста в документе;
- 2) создание бинарного препарата документа;
- 3) выделение признаков;
- 4) генерация пространства признаков.

Остановимся подробнее на каждом из этапов.



Рисунок 2 – Общая схема формирования вектора вторичных признаков для индексации документа

На первом этапе происходит выделение областей (маски) текста на изображении. Выделение текстовых строк и формирование маски текста на изображении выполняется путем классификации скользящим окном на основе разработанного авторами оригинального алгоритма адаптивного бустинга.

На втором этапе происходит бинаризация изображения (рис. 3). Перед бинаризацией изображение преобразуется в оттенки серого. Затем по всему изображению проходит окно (ядро низкочастотного фильтра), после чего из усредненных соответствующих пикселей вычитаются значения соответствующих пикселей исходного изображения (1). Если их разница отрицательна, то значение пикселя считается равным 0, иначе пиксель принимает значение, получившееся при вычитании (2). Затем выбирается порог бинаризации. В качестве порога бинаризации может использоваться комбинированный порог:

$$I'(x, y) = \sum_i^M \sum_j^M G(i, j)I(x+i, y+j) - I(x, y) \quad (1)$$

$$I''(x, y) = \begin{cases} I'(x, y), & \text{если } I'(x, y) > 0 \\ 0, & \text{если } I'(x, y) \leq 0 \end{cases} \quad (2)$$

$$I_{bin}(x, y) = \begin{cases} 1, & \text{если } I''(x, y) \geq I_t \\ 0, & \text{если } I''(x, y) < I_t \end{cases} \quad (3)$$

где I – изображение в оттенках серого, $I(x, y)$ – пиксель изображения, M – размер ядра свертки, G – ядро свертки фильтра, $I_{bin}(x, y)$ – бинаризованный пиксель, I_t – порог бинаризации.

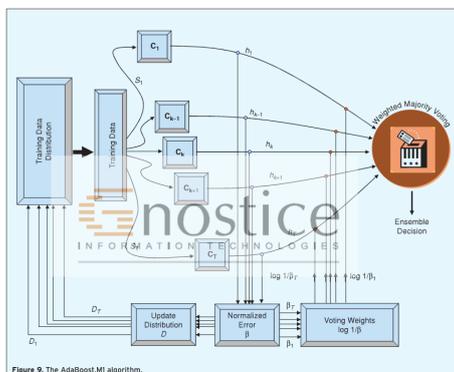


Figure 9. The AdaBoost.M1 algorithm.

where E is the ensemble error [26]. Since $\alpha_i < 1/2$, E is guaranteed to decrease with each new classifier. In most practical cases, the error decreases very rapidly in the first few iterations, and approaches zero as new classifiers are added. While this is remarkable on its own account, the surprising resistance of AdaBoost to overfitting is particularly noteworthy. Overfitting is a commonly observed phenomenon where the classifier performs poorly on test data, despite achieving a very small training error. Overfitting is usually attributed to memorizing the data, or learning the noise in the data. As a classifier's capacity increases (for example, with the complexity of its architecture), so does its tendency to memorize the training data and/or learn the noise in the data. Since the capacity of an ensemble increases with each added classifier, one would expect AdaBoost to suffer from overfitting, particularly if its complexity exceeds what is necessary to learn the underlying data distribution. Yet, AdaBoost performance usually levels off with increasing number of classifiers with no indication of overfitting.

Schapire *et al.* later provided an explanation to this phenomenon based on the so-called *margin theory* [52]. The details of margin theory are beyond the scope of this paper; however, the margin of an instance x , in loose terms, is its distance from the decision boundary. The further away an instance is from the boundary, the larger its margin, and hence the higher the confidence of the classifier in correctly classifying this instance. Margins are usually described within the context of support vector machine (SVM) type classifiers, which are based on the idea of maximizing the margins between the instances and the decision boundary. In fact, SVMs find those instances, called support vectors, that lie closest to the decision boundary. The support vectors are said to define the margin that separates the classes. The concept of a margin is conceptually illustrated in Figure 10.

By using a slightly different definition of a margin—suitably modified for AdaBoost—Schapire *et al.* show that AdaBoost also boosts the margins, that is, it finds the decision boundary that is further away from the instances of all classes. In the context of AdaBoost, the margin of an

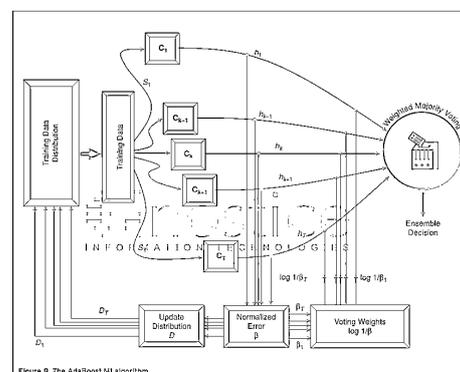


Figure 9. The AdaBoost.M1 algorithm.

where E is the ensemble error [26]. Since $\alpha_i < 1/2$, E is guaranteed to decrease with each new classifier. In most practical cases, the error decreases very rapidly in the first few iterations, and approaches zero as new classifiers are added. While this is remarkable on its own account, the surprising resistance of AdaBoost to overfitting is particularly noteworthy. Overfitting is a commonly observed phenomenon where the classifier performs poorly on test data, despite achieving a very small training error. Overfitting is usually attributed to memorizing the data, or learning the noise in the data. As a classifier's capacity increases (for example, with the complexity of its architecture), so does its tendency to memorize the training data and/or learn the noise in the data. Since the capacity of an ensemble increases with each added classifier, one would expect AdaBoost to suffer from overfitting, particularly if its complexity exceeds what is necessary to learn the underlying data distribution. Yet, AdaBoost performance usually levels off with increasing number of classifiers with no indication of overfitting.

Schapire *et al.* later provided an explanation to this phenomenon based on the so-called *margin theory* [52]. The details of margin theory are beyond the scope of this paper; however, the margin of an instance x , in loose terms, is its distance from the decision boundary. The further away an instance is from the boundary, the larger its margin, and hence the higher the confidence of the classifier in correctly classifying this instance. Margins are usually described within the context of support vector machine (SVM) type classifiers, which are based on the idea of maximizing the margins between the instances and the decision boundary. In fact, SVMs find those instances, called support vectors, that lie closest to the decision boundary. The support vectors are said to define the margin that separates the classes. The concept of a margin is conceptually illustrated in Figure 10.

By using a slightly different definition of a margin—suitably modified for AdaBoost—Schapire *et al.* show that AdaBoost also boosts the margins, that is, it finds the decision boundary that is further away from the instances of all classes. In the context of AdaBoost, the margin of an

Рисунок 3 – Изображение до и после бинаризации

В результате должно получиться черно-белое изображение, как показано на рис. 3. Постобработка бинарного препарата включает в себя удаление отдельных точек и операции морфологического утоньшения. Далее возможно удаление областей изображения с выделенным на первом этапе текстом.

На третьем этапе происходит выделение признаков, сначала простых (рис. 4), а затем на их основе объединение в комбинированные признаки – прямоугольники и связи между ними (рис. 5). Простые геометрические признаки представляют собой горизонтальные и вертикальные линии. Эти линии обладают следующими атрибутами: длина, стиль (гладкая, штрихпунктирная), связность с другими линиями (пересечение или соединение). Линии обнаруживаются путем перебора пикселей по горизонтали

для горизонтальных линий и пикселей по вертикали для вертикальных линий. Алгоритм анализирует линию на стиль начертания – гладкий или штрихпунктирный. Псевдокод представлен в листинге 1.

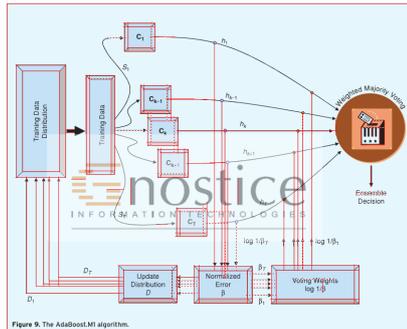


Figure 9. The AdaBoost.M1 algorithm.

where E is the ensemble error [26]. Since $\epsilon_t < 1/2$, E is guaranteed to decrease with each new classifier. In most practical cases, the error decreases very rapidly in the first few iterations, and approaches zero as new classifiers are added. While this is remarkable on its own account, the surprising resistance of AdaBoost to overfitting is particularly noteworthy. Overfitting is a commonly observed phenomenon where the classifier performs poorly on test data, despite achieving a very small training error. Overfitting is usually attributed to memorizing the data, or learning the noise in the data. As a classifier's capacity increases (for example, with the complexity of its architecture), so does its tendency to memorize the training data and/or learn the noise in the data. Since the capacity of an ensemble increases with each added classifier, one would expect AdaBoost to suffer from overfitting, particularly if its complexity exceeds what is necessary to learn the underlying data distribution. Yet, AdaBoost performance usually levels off with increasing number of classifiers with no indication of overfitting.

Schapire et al. later provided an explanation to this phenomenon based on the so-called margin theory [52]. The details of margin theory are beyond the scope of this paper; however, the margin of an instance x , in loose terms, is its distance from the decision boundary. The further away an instance is from the boundary, the larger its margin, and hence the higher the confidence of the classifier in correctly classifying this instance. Margins are usually described within the context of support vector machine (SVM) type classifiers, which are based on the idea of maximizing the margins between the instances and the decision boundary. In fact, SVMs find those instances, called support vectors, that lie closest to the decision boundary. The support vectors are said to define the margin that separates the classes. The concept of a margin is conceptually illustrated in Figure 10.

By using a slightly different definition of a margin—suitably modified for AdaBoost—Schapire et al. show that AdaBoost also boosts the margins, that is, it finds the decision boundary that is further away from the instances of all classes. In the context of AdaBoost, the margin of an

setting. Alternatively, entirely different type of classifiers, such as MLPs, decision trees, nearest neighbor classifiers, and support vector machines, can also be combined for added diversity. However, combining different models, or even different architectures of the same model, is used only for specific applications that warrant them. Correlation Diversity is measured as the correlation between two classifier outputs, defined as

$$A_j = \frac{ad - bc}{\sqrt{(a+b)(c+d)(a+c)(b+d)}} \quad 0 \leq \rho \leq 1. \quad (1)$$

Maximum diversity is obtained for $\rho = 0$, indicating that the classifiers are uncorrelated. Q-Statistic Defined as

$$Q = \frac{ad - bc}{ad + bc} \quad (2)$$

2.3. Measures of Diversity

Several measures have been defined for quantitative assessment of diversity. The most commonly used measures, different from the pair-wise diversity measures, are overall diversity Q and Q -Statistic. If the same instance are correctly classified by both classifiers, and negative values, otherwise. Maximum diversity is, once again, obtained for $Q = 0$.

Divergence and Double Fault Measures: The divergence is the probability that the two classifiers will disagree, whereas the double fault measure is the



Figure 4. k-fold data splitting for generating different, but overlapping, training datasets.

Рисунок 4 – Изображения с обнаруженными линиями

```

procedure GetLineFromImage(
    imgBin2Int[,]: Integer; //Массив пикселей бинаризованного изображения
    W_img, H_img: Integer; //Ширина и высота изображения в пикселях
    LineList: List; //Список горизонтальных линий
    hmin: Integer; // Минимальная длина линии
    hspan: Integer; // Промежуток между линиями
    hmax: Integer; // Максимальная длина линии
    LineColor: Integer; // Цвет линии
    LDensity: Double; // Плотность линии
    DashLength: Integer); // Длина штриха
{
//Выделение горизонтальных линий
//Минимальная длина
lmin := hmin;
//Максимальная длина
lmax := hmax;
//Максимальные промежутки между соседними пикселями линии
lspan:= hspan;
//Цикл обработки по строкам для горизонтальных линий
For i := 0 To H_img - 1 Do
{
//Сбросить счетчики
clen := 0; cden := 0; cseq := 0; cemp := 0; cadd := 0;
ptrLine^.StyleDraw := ldsSolid;
ptrLine^.StyleCross := lscNon;
ptrLine^.StyleOrientation := lsoInclined;
ptrLine^.UseForRect := false;
ptrLine^.CountCross := 0;
ptrLine^.LineLength := 0;
cspan := 0;
//Инициализируем прямую как "несуществующую"
ptrLine^.Line.Left := -1;
ptrLine^.Line.Right := -1;
ptrLine^.Line.Top := i;
ptrLine^.Line.Bottom := i;
}
}
    
```

```

ptrLine^.StyleOrientation := IsoHorizontal;
For j := 0 To Wimg - 1 Do
{
    //Увеличиваем счетчик длины линии
    Inc(clen);
    //Если пиксель совпадает с цветом линии
    if(imgBin2Int[i,j] = LineColor)
    {
        //Увеличиваем счетчик промежутков между кусками линии
        if(cspan > 0)
        {
            Inc(cden);
            if(cadd > DashLength)Then
                cseq := DashLength;
        }
        end;
        Inc(cadd);
        //Увеличиваем счетчик совпадающих пикселей линии
        Inc(cseq);
        // Если это первый совпавший пиксель,
        // то инициализируем начальную точку
        If((ptrLine^.Line.Left = -1)And(cseq > DashLength))
        {
            clen := 0;
            cemp := 0;
            ptrLine^.Line.Left := j - DashLength;
        }
        //Т.к. совпал цвет линии, то присваиваем крайнее положение справа
        if(cseq > DashLength)
        {
            ptrLine^.Line.Right := j + 1;
            //change
            cspan := 0;
        }
        //Т.к. цвет линии совпадает, то обнуляем счетчик "пустоты"
        //cspan := 0;
    }
    //Цвет линии не совпадает
    else
    {
        //Увеличиваем счетчик пустого пространства между отрезками линии
        Inc(cspan);
        Inc(cemp);
        if(cspan > DashLength)
            cadd := 0;
    }
    //Если промежуток между соседними точками больше разрешенного промежутка
    if(cspan >= lspan)
    {
        // Проверяем не является ли линия пустой
        if((ptrLine^.Line.Left <> -1)And(ptrLine^.Line.Right <> -1))
        {
            // Задаем ориентацию линии горизонтальной
            ptrLine^.StyleOrientation := IsoHorizontal;
            // Проверяем линию на допустимость значений и
            // проверяем тип начертания гладкая или штрихпунктирная
            if(CheckLine(ptrLine, LDensity))
            {
                // Добавляем линию
                LineList.Add(ptrLine);
            }
        }
    }
}

```

```

// Создаем новую линию
New(ptrLine);
}
// Сбрасываем счетчики, т.к. превышено значение
// пустоты между отрезками линии
clen := 0; cden := 0; cseq := 0;   cemp := 0; cadd := 0;
ptrLine^.StyleDraw := ldsSolid;
ptrLine^.StyleCross := lscNon;
ptrLine^.StyleOrientation := lsoInclined;
ptrLine^.UseForRect := false;
ptrLine^.CountCross := 0;
ptrLine^.LineLength := 0;
//Инициализируем прямую как "несуществующую"
ptrLine^.Line.Left := -1;
ptrLine^.Line.Right := -1;
ptrLine^.Line.Top := i;
ptrLine^.Line.Bottom := i;
}
}
// Проверяем линию на допустимость значений и
// проверяем тип начертания (гладкая или штрихпунктирная линия)
if(CheckLine(ptrLine, LDensity))
{
// Добавляем линию
LineList.Add(ptrLine);
// Создаем новую линию
New(ptrLine);
}
}

```

Листинг 1 – Псевдокод нахождения горизонтальных линий в бинарном изображении

Далее эти линии вместе с маской текста образуют более сложные признаки. Выделяются прямоугольники из найденных линий, определяется включение в прямоугольниках других линий, других прямоугольников, текста, наличие и количество связей с другими прямоугольниками или линиями.

setting. Alternatively, entirely different type of classifiers, such as MLPs, decision trees, nearest neighbor classifiers, and support vector machines can also be combined for added diversity. However, combining different models, or even different architectures of the same model, is used only for specific applications that warrant them. Diversity is typically obtained through resampling of the training data, as this procedure is theoretically more tractable. Finally, diversity can also be achieved by using different features. In fact, generating different classifiers using random feature subsets is known as the *random subspace method* [44], and it has found widespread use in certain applications, which are discussed later in future research areas.

2.2. Measures of Diversity
Several measures have been defined for quantitative assessment of diversity. The most commonly used measures, referred to as the *confusion matrix* [45] (see Fig. 7) are: P is the overall probability of correct classification, and an overall diversity measure is obtained by averaging these pairwise measures. Given two hypotheses h_1 and h_2 , we use the notations

	h_1 is correct	h_1 is incorrect
h_1 is correct	a	b
h_1 is incorrect	c	d

Disagreement and Double Fault Measures. The disagreement is the probability that the two classifiers will disagree, whereas the double fault measure is the

$$P_d = \frac{ad - bc}{\sqrt{(a+b)(c+d)(a+c)(b+d)}}, \quad 0 \leq P_d \leq 1. \quad (1)$$

Maximum diversity is obtained for $P_d = 0$, indicating that the classifiers are uncorrelated. Q-Statistic Defined as

$$Q = \frac{ad - bc}{ad + bc} \quad (2)$$

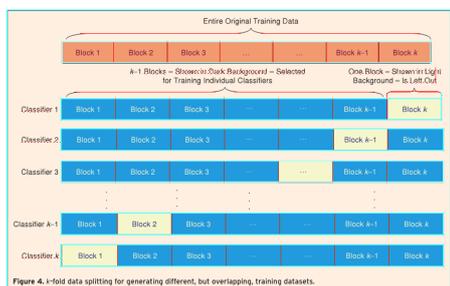


Figure 4. k-fold data splitting for generating different, but overlapping, training datasets.

pairs for the second level classifier, C_{T+1} . Traditionally, the k-fold selection process described in Section 2 is used to obtain the training data for classifier C_{T+1} . Specifically, the entire training dataset is divided into T blocks, and each first-level classifier C_1, \dots, C_T is first trained on a different set of $T-1$ blocks of the training data. Therefore, there is one block of data not seen by each of the classifiers C_1 through C_T . The outputs of each classifier for the block of instances on which it was not trained, along with the correct labels of those instances, constitute the training data for the second level meta-classifier C_{T+1} . Once C_{T+1} is trained, all data are pooled, and individual classifiers C_1, \dots, C_T are retrained on the entire database, using a suitable resampling method.

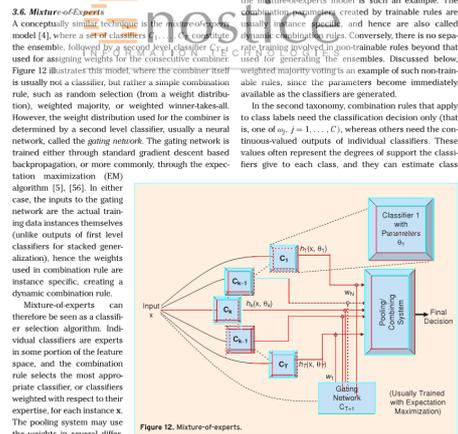


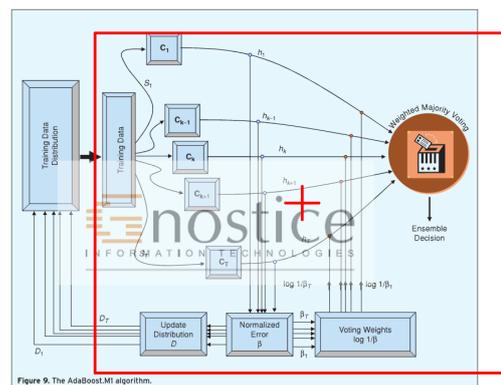
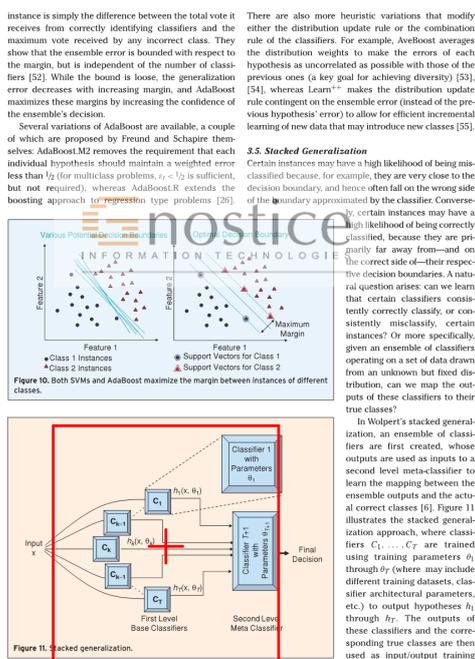
Figure 12. Mixture-of-experts.

Рисунок 5 – Изображение с обнаруженными линиями и прямоугольниками

6. Количество полусвязных линий (линия имеет соединение с прямоугольником, только с одной стороны).
7. Количество полусвязанных прямоугольников (прямоугольник связан только с линией, но не с другими прямоугольниками).
8. Количество полностью связанных линий (линии, которые своими концами соединены с прямоугольниками)
9. Количество полностью связанных прямоугольников (прямоугольники, которые связаны друг с другом линией)
10. Количество прямоугольников внутри области и для каждого следующего признака:
 - а) количество всех линий этого прямоугольника,
 - б) количество всех горизонтальных линий,
 - в) наличие текста,
 - г) нормализованная площадь к площади области,
 - д) соотношение ширины к высоте.

Выводы

Для оценки эффективности предложенного метода было проведено моделирование работы оригинального классификатора на базе изображений документов. База изображений документов была разбита на две части: обучающую выборку и тестовую выборку. Обучающая база изображений образована из 600 изображений (300 – из положительных и 300 – из отрицательных примеров). Тестовая база изображений образована из 300 изображений (150 – из положительных и 150 – из отрицательных примеров).



where E is the ensemble error [26]. Since $\epsilon_i < 1/2$, E is guaranteed to decrease with each new classifier. In most practical cases, the error decreases very rapidly in the first few iterations, and approaches zero as new classifiers are added. While this is remarkable on its own account, the surprising resistance of AdaBoost to overfitting is particularly noteworthy. Overfitting is a commonly observed phenomenon where the classifier performs poorly on test data, despite achieving a very small training error. Overfitting is usually attributed to memorizing the data, or learning the noise in the data. As a classifier's capacity increases (for example, with the complexity of its architecture), so does its tendency to memorize the training data and/or learn the noise in the data. Since the capacity of an ensemble increases with each added classifier, one would expect AdaBoost to suffer from overfitting, particularly if its complexity exceeds what is necessary to learn the underlying data distribution. Yet, AdaBoost performance usually levels off with increasing number of classifiers with no indication of overfitting.

Schapire *et al.* later provided an explanation to this phenomenon based on the so-called *margin theory* [52]. The details of margin theory are beyond the scope of this paper; however, the margin of an instance x , in loose terms, is its distance from the decision boundary. The further away an instance is from the boundary, the larger its margin, and hence the higher the confidence of the classifier in correctly classifying this instance. Margins are usually described within the context of support vector machine (SVM) type classifiers, which are based on the idea of maximizing the margins between the instances and the decision boundary. In fact, SVMs find those instances, called support vectors, that lie closest to the decision boundary. The support vectors are said to *define the margin* that separates the classes. The concept of a margin is conceptually illustrated in Figure 10.

By using a slightly different definition of a margin—suitably modified for AdaBoost—Schapire *et al.* show that AdaBoost also *boosts* the margins, that is, it finds the decision boundary that is further away from the instances of all classes. In the context of AdaBoost, the margin of an

Рисунок 6 – Примеры работы алгоритма обнаружения и выделения

При моделировании на тестовой выборке изображений документов алгоритм обнаружения специфической графической информации показал устойчивость и эффективность работы. Получены следующие характеристики классификатора:

- вероятность правильной классификации (по тестовой выборке положительных примеров) составляет: 0.91;

– вероятность ложноположительной классификации (по тестовой выборке отрицательных примеров) составляет: 0.056.

Примеры работы алгоритма обнаружения специфической графической информации показаны на рис. 6, где рамками выделены области изображений документов, классифицированные как содержащие блок-схемы.

В качестве дальнейшего направления развития предполагается адаптация алгоритма к условиям работы с зашумлёнными изображениями, содержащими искажения. Для этого необходимо произвести соответствующие доработки алгоритмов бинаризации, выделения первичных признаков и формирования сложных признаков.

Литература

1. S. Mao. Document structure analysis algorithms: a literature survey / S. Mao, A. Rosenfeld, T. Kanungo // Document Recognition and Retrieval X. – 2003. – Vol. 5011. – P. 197-207.
2. Shazia Akram. Document Image Processing. A Review / Shazia Akram, Mehraj-Ud-Din Dar, Aasia Quyoum // International Journal of Computer Applications. – 2010. – Vol. 10, № 5. – P. 35-40.
3. G. Nagy. A prototype document image analysis system for technical journals / G. Nagy, S. Seth, M. Viswanathan // Computer. – 1992. – Vol. 25. – P. 10-22.
4. A. Lemaitre. Contribution of Multiresolution Description for Archive Document Structure Recognition / A. Lemaitre, J. Camillerapp, B. Couasnon // Ninth International Conference on Document Analysis and Recognition. – 2007. – Vol. 01. – P. 247-251.

Literatura

1. S. Mao, A. Rosenfeld, and T. Kanungo. Document structure analysis algorithms: a literature survey // Document Recognition and Retrieval X. – 2003. – Vol. 5011. – P. 197-207.
2. Shazia Akram, Mehraj-Ud-Din Dar, Aasia Quyoum. Document Image Processing - A Review // International Journal of Computer Applications. – 2010. – Vol. 10, No. 5. – P. 35-40.
3. G. Nagy, S. Seth, M. Viswanathan. A prototype document image analysis system for technical journals // Computer. – 1992. – Vol. 25. – P. 10-22.
4. Lemaitre, J. Camillerapp, B. Couasnon. Contribution of Multiresolution Description for Archive Document Structure Recognition // Ninth International Conference on Document Analysis and Recognition. 2007. Vol. 01. P. 247-251.

RESUME

N.A. Kostromov, I.V. Beketova, S.L. Karateev

Formation of Feature Spaces and their Extraction from the Images for Documents Indexing

The article describes the formation of the original algorithms for generation of feature vectors for the detection of areas with specific the graphic images in the document image.

The process of formation of feature vector consists of two stages. In the first phase, we get a binary document image. Lines and text strings are allocated from the resulting image. The result is a set of primary attributes. In the second stage based on the analysis of a set of primary features, the combined features, i.e. rectangular areas, signs of the text, the signs of connectivity and nesting stand out. Feature vector is formed by combining the primary features and combined ones.

Trained with the given arrays of images, the system can detect the image of the document with the desired graphical information. The following characteristics of the classifier were obtained: the probability of correct classification (as a test set of positive examples) is 0.91; the probability of false positive classifications (on a test set of negative examples) is 0.056.

Статья поступила в редакцию 05.06.2012.