

Matching document images with ground truth

John D. Hobby

Bell Laboratories, Lucent Technologies, Murray Hill, NJ, USA
e-mail: hobby@bell-labs.com

Received 25 June, 1997 / Revised August 20, 1997

Abstract. Since optical character recognition systems often require very large amounts of training data for optimum performance, it is important to automate the process of finding ground truth character identities for document images. This is done by finding a transformation that matches a scanned image to the machine-readable document description that was used to print the original. Rather than depend on finding feature points, a more robust procedure is to follow up by using an optimization algorithm to refine the transformation. The function to optimize can be based on the character bounding boxes – it is not necessary to have access to the actual character shapes used when printing the original.

Key words: Optical character recognition – Ground truth – Nelder-Mead algorithm

1 Introduction

Training and testing optical character recognition systems often requires large numbers of realistic character and word images. If the correct “ground truth” character identities have to be entered or corrected by hand, it is difficult to gather enough accurate data. Semiautomatic methods are popular now, but they still require a lot of labor (Phillips et al., [14,15]; Rogers et al., [16]; Nagy et al., [12]). Much of this labor can be avoided by automatic ground-truthing where a scanned page image is matched with the original document description that was used to print the page. The procedure suggested by Kanungo et al. [8,9] is to look for certain feature points in the scanned image, transform the original page description to make the feature points match, and then use template matching to look for character images near where the original page description says they should be.

The big advantage of automatic ground truthing is that it requires little or no human labor. Of course this comes at the expense of requiring the original page description, but this limitation is getting less serious as electronic documents become more prevalent. The primary goal of this work is to deal with one of the other

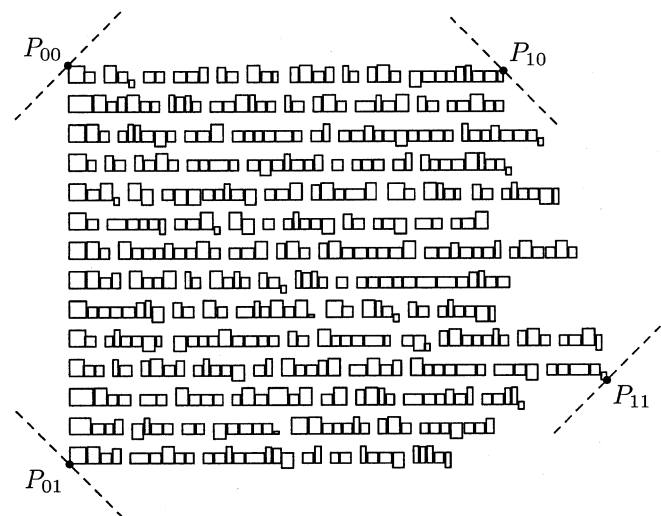


Fig. 1. The bounding boxes for all characters on a small sample page with Kanungo’s feature points indicated. The feature points are chosen so that the 45° dashed lines are supporting lines

limitations of Kanungo’s approach, namely the dependence on finding feature points. It would be better to have a more robust and more accurate way of transforming the original page description to match the image. In addition, our secondary goal is to try to avoid the need to have fonts available for template matching.

The original page description identifies each character on the page and gives its bounding box. Figure 1 shows how Kanungo uses this information to find feature points. The feature points are the bounding box corners where $x + y$, $x - y$, $-x + y$ and $-x - y$ are maximized. Performing a similar computation for connected components in the scanned page image yields another set of four feature points. Kanungo then finds a geometric transformation that aligns the two sets of feature points and hopefully aligns the character bounding boxes with the corresponding character images.

This was not very reliable in our tests because scanned images often have speckles and other material not in the original page description. Kanungo and Haralick [9] give

a heuristic for dealing with this by checking bounding box sizes, but it is always possible to wind up with the wrong feature points. If the scanned image is from a bound book, noise around the edge of the page exacerbates the problem. Even when we find the right feature points, transformations based on just a few feature points can yield suboptimal matches. For these reasons, we attempt to find a transformation that makes the entire page description fit the image as well as possible. We define a function that measures this fit, and then use standard optimization techniques. Since it is still important to have a good starting point for the optimization, the new approach does not replace feature point analysis, it complements it. The optimization can be started from the identity transformation if feature point analysis fails completely, but it is best to avoid this if possible.

Section 2 defines the optimization problem and explains what standard techniques are best for solving it. Section 3 explains how careful bucketing strategies can speed up the optimization enough to make it practical. Section 4 investigates the problem of assigning ground truth to character and word images without *a priori* knowledge of the character shapes or non-geometric distortions in the printing and scanning process. Finally, Sect. 5 gives results and discusses applications, and Sect. 6 gives some concluding results.

2 Finding the transformation

The task is to take character bounding boxes such as those in Fig. 1, and find a geometric transformation that matches them to a similar set of bounding boxes derived from the scanned image. Specifically, we use the bounding boxes for 8-connected sets of black pixels in the image and choose an affine transformation

$$T(x, y) = \begin{pmatrix} t_{xx} & t_{xy} \\ t_{yx} & t_{yy} \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} + \begin{pmatrix} t_x \\ t_y \end{pmatrix}. \quad (1)$$

An initial estimate for T can be computed from the feature points illustrated in Fig. 1. We need to improve this estimate by minimizing a function that depends on the six parameters that appear in (1), and measures the degree by which the page image fails to match the original document description when the transformation T is applied.

In order to have axis-aligned rectangles, it is best to start by applying T to each connected component in the image and then compute bounding boxes as indicated by dashed lines in Fig. 2a. Hence we have two types of boxes:

Image boxes are derived from the result of applying an affine transformation (1) to the page image. They are the bounding boxes of the 8-connected components, where each component is transformed by T before computing the bounding box. The synonymous term *transformed image box* emphasizes the presence of this transformation.

Ideal boxes are the character bounding boxes from the original page description. These are indicated by solid lines in Fig. 2b.

Since characters can break up and/or merge together during printing and scanning, there is not necessarily a one-to-one correspondence. In the example of Fig. 2b, there are two image boxes for the ideal box “r,” and ideal boxes “a” and “l” both correspond to the same image box.

2.1 The mismatch function

The above example suggests that we measure the degree of mismatch by developing a function d that takes two boxes and measures how much they would have to be changed in order for one of the boxes to contain the other. For each ideal box A , we can find the image box B that minimizes $d(A, B)$, and then apply a standard vector norm to the resulting list of d values. This approach allows characters to break up or merge together and does not penalize for noise or non-text material for which there are image boxes but no ideal boxes.

Hence we define $d(A, B)$ to be

$$\begin{aligned} \min & (d_f(A_{x1}, A_{x2}, B_{x1}, B_{x2}) + d_f(A_{y1}, A_{y2}, B_{y1}, B_{y2}), \\ & d_f(B_{x1}, B_{x2}, A_{x1}, A_{x2}) + d_f(B_{y1}, B_{y2}, A_{y1}, A_{y2})) \\ & + d_p(A_{x2} - A_{x1}, B_{x2} - B_{x1}) + d_p(A_{y2} - A_{y1}, B_{y2} - B_{y1}), \end{aligned}$$

where $x1$ and $x2$ subscripts refer to a box’s minimum and maximum x coordinate, $y1$ and $y2$ subscripts refer to the minimum and maximum y coordinate, d_p is a penalty term that is nonzero when the ratio of its arguments is too small or too large, and $d_f(x_1, x_2, x_3, x_4)$ equals

$$\begin{cases} 0 & \text{if } x_3 \leq x_1 \\ & \text{and } x_2 \leq x_4; \\ \min(|x_3 - x_1|, |x_4 - x_2|) \\ + \max(0, x_2 - x_1 - (x_4 - x_3)) & \text{otherwise.} \end{cases}$$

Function d_f is the amount of translation and stretching necessary to fit interval $[x_1, x_2]$ within interval $[x_3, x_4]$. The d_p terms ensure that unreasonably large objects are not treated as run-together characters and that broken characters are required to contain reasonably large connected components. A generous choice is

$$d_p(a, b) = \max(0, \max(a, b) - 8 \cdot \min(a, b)).$$

In other words, we allow the height or width of an image box to differ from the corresponding ideal box dimension by a factor of 8 before penalties are accessed. The number 8 was determined empirically by looking for large contributions to the mismatch function after optimization and manually verifying that they were not due to inadequate tolerance for the size ratio in $d_p(a, b)$.

We can now define the mismatch function as follows: take the six parameters that appear in (1); use them to find bounding boxes of the connected components in the transformed image, choose among these transformed image boxes, an image box B that minimizes $d(A, B)$ for each ideal box A ; and apply a vector norm to the resulting $d(A, B)$ values. We choose the L_4 norm (fourth root of sum of fourth powers) since experiments suggested that the L_∞ norm is too sensitive to noise and the L_2 norm causes the mismatch function to have too many local minima.

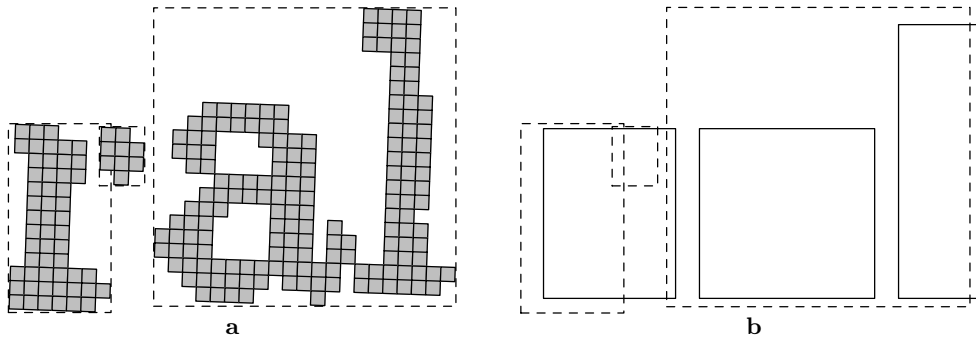


Fig. 2a. Portions of a slightly rotated 200 dpi page image with bounding boxes; **b.** the same image boxes (dashed lines) superimposed on the ideal boxes from the original page description (solid lines)

2.2 Minimizing the mismatch

Most standard optimization methods do not work well on the mismatch function because it lacks continuity and smoothness properties. Many comparisons are needed to evaluate the function and it does not take much of a change in the transformation parameters (1) to reverse one of the comparisons. This means that the mismatch function has an extremely large number of slope discontinuities. As explained by Wright [20], optimization problems of this kind are best solved by direct search methods such as Nelder–Mead [13] and Torczon’s algorithm [19].

The basic idea of a direct search method is to perform a sequence of function evaluations, using the previous history to decide what input to give for the next function evaluation. Figure 3 illustrates the first and last function evaluations when using Nelder–Mead to minimize the mismatch function for a sample page scanned at 200 dpi. In this case there were 219 function evaluations and the mismatch was reduced from 11.63 to 4.13. Superimposing the transformed image boxes on the original page description gives a visual impression of how much the mismatch was reduced.

If Fig. 3 were based on Torczon’s algorithm, it would appear almost the same, but more function evaluations would be needed. Torczon’s algorithm is well suited to parallel computation because its outer loop does several function evaluations per iteration and these could all be done in parallel. Torczon’s algorithm has guaranteed convergence properties that Nelder–Mead lacks, but the results in Sect. 5 will show that it requires significantly more function evaluations. In addition, it is hard to get enough information about the mismatch function to make use of the convergence properties.

For these reasons, we shall concentrate on Nelder–Mead. For our six-dimensional problem, the Nelder–Mead algorithm maintains a set of seven

$$\left(t_{xx}, t_{xy}, t_{yx}, t_{yy}, \frac{t_x}{W}, \frac{t_y}{H} \right)$$

values that are the vertices of a simplex in Euclidean 6-space, where W and H are constant scale factors. (It suffices to let W and H be the width and height of the page). If T_7 is the vertex where the mismatch function is greatest (i.e., worst) and \bar{T} is the average of the other

To be, or not to be: that is the question:
Whether 'tis nobler in the mind to suffer
The slings and arrows of outrageous fortune,
Or to take arms against a sea of troubles,
And, by opposing end them? To die: to sleep;
No more; and, by a sleep to say we end
The heartache and the thousand natural shocks
That flesh is heir to, 'tis a consummation
Devoutly to be wished. To die, to sleep;
To sleep; perchance to dream: ay, there's the rub;
For in that sleep of death what dreams may come,
When we have shuffled off this mortal coil,
Must give us pause. Where's the respect
That makes calamity of so long life;

$$\text{mismatch } 11.63 \text{ at } (t_{xx}, t_{xy}, t_{yx}, t_{yy}, t_x, t_y) = (1.002, 0.002, -0.002, 1.002, -9.88, 2.91)$$

To be, or not to be: that is the question:
Whether 'tis nobler in the mind to suffer
The slings and arrows of outrageous fortune,
Or to take arms against a sea of troubles,
And, by opposing end them? To die: to sleep;
No more; and, by a sleep to say we end
The heartache and the thousand natural shocks
That flesh is heir to, 'tis a consummation
Devoutly to be wished. To die, to sleep;
To sleep; perchance to dream: ay, there's the rub;
For in that sleep of death what dreams may come,
When we have shuffled off this mortal coil,
Must give us pause. Where's the respect
That makes calamity of so long life;

$$\text{mismatch } 4.13 \text{ at } (t_{xx}, t_{xy}, t_{yx}, t_{yy}, t_x, t_y) = (0.999, 0.007, -0.003, 0.998, -11.14, 22.17)$$

Fig. 3. The transformed image boxes superimposed on the original page description before and after using Nelder–Mead to minimize the mismatch. The corresponding transformation parameters and mismatch values appear below each image

six vertices, a typical Nelder–Mead step consists of trying one or two points along the $\bar{T}T_7$ line and trying to update the simplex by replacing one of the existing vertices.

A Nelder–Mead iteration begins by trying the “reflection point” $T_{\text{ref}} = 2\bar{T} - T_7$. If this is better than the

previous best point T_1 , try $\bar{T} + 1.3(T_{\text{ref}} - \bar{T})$ and replace T_7 by the better of this point and T_{ref} .

If T_{ref} is not as good as T_1 but is better than the point T_6 that was second-worst, replace T_7 by T_{ref} . Otherwise, T_{ref} is either worse than T_7 or worse than T_6 but better than T_7 . Let the “contraction point” T_{con} be $(\bar{T} + T_7)/2$ in the former case and $(\bar{T} + T_{\text{ref}})/2$ in the latter case. If T_{con} is better than T_7 , replace T_7 by T_{con} .

If T_{con} fails to improve on T_7 , Nelder–Mead shrinks the simplex toward the best point T_1 by replacing T_i with $(T_i + T_1)/2$ for each of the other 6 vertices T_i . This completes one iteration. Refer to [13] for explanations of numeric parameters and details. Nelder–Mead proceeds as outlined above until some convergence test is satisfied; e.g., all 7 vertices are within Euclidean distance 3×10^{-5} of T_1 or the 7 function values are identical to within one part in 10^4 .

We modify this convergence test by defining a *critical value* for the mismatch function, and using very tight tolerance in the other tests if the critical value has not been achieved. The purpose of the critical mismatch value is to ensure that there is a good chance of assigning appropriate ground truth. The critical value is what the mismatch function would return if each $d(A, B)$ value were half the median character width.

3 Efficient implementation

The Nelder–Mead algorithm can require 300 function evaluations to find the transformation parameters that minimize the mismatch function. For a typical page containing 1800 characters and 1800 connected components, Sect. 2.1 apparently requires 1800^2 evaluations of $d(A, B)$ for each invocation of the mismatch function. The resulting 1 billion $d(A, B)$ evaluations would probably make the overall algorithm unacceptably slow.

Two strategies result in significant speed-ups:

1. Use a carefully chosen subset of the ideal boxes – the ones where $d(A, B)$ is likely to be largest if the transformation parameters are wrong.
2. Preprocess the transformed image boxes using a bucketing algorithm that allows $d(A, B)$ to be evaluated only for the most reasonable (A, B) pairs.

Strategy 1 is simplest, but it needs to be used carefully since it involves changing the mismatch function rather than just evaluating it more efficiently.

3.1 Choosing a subset of the ideal boxes

The basic idea of choosing a subset of the boxes is relevant to a wide range of problems. For instance some skew detection methods are based on fiducial points derived from connected components or other image features [1, 17]. Taking a random sample of the fiducial points would certainly result in a speed up.

The mismatch function is especially well suited to random sampling because it is designed to allow for material in the page image that does not correspond to any

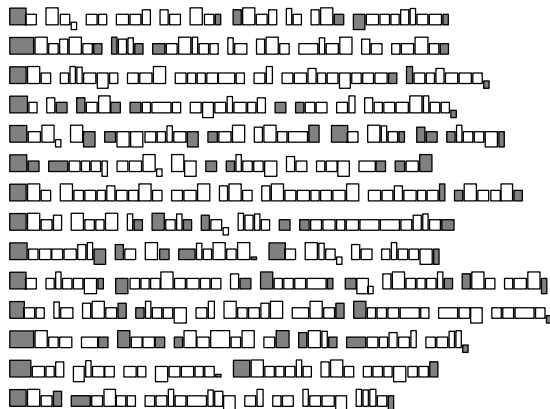


Fig. 4. The character boxes from a sample page with a subset shaded to indicate which of them might be used as ideal boxes for the mismatch function

characters in the original page description. We can base the mismatch function’s ideal boxes on just a subset of the characters as indicated by the shaded boxes in Fig. 4.

The selected characters in Fig. 4 are at the beginnings and ends of lines, around any other large empty spaces, and around some of the word breaks. If a poorly chosen transform is applied to the image, it is likely that many of the selected boxes will not match any of the resulting image boxes. Hence the mismatch function will tend to return large values when given non-optimal transformation parameters.

It is easy to choose such a subset of the characters if we assume that the original page description lists boxes approximately in reading order. Suppose there are n characters numbered $0, 1, 2, \dots, n - 1$, where each character i has bounding box A_i . Let σ be a pseudorandom function that maps $0, 1, \dots, n - 1$ into the interval $[1, 1.3]$, and consider the positive integers less than n in order of decreasing

$$\sigma(i) \cdot d(A_{i-1}, A_i).$$

For each such integer i , make i and $i - 1$ part of the chosen subset until the chosen subset reaches some predetermined size. This predetermined size will typically be something like

$$\min(\max(300, 0.15n), n).$$

The purpose of the function σ is to ensure that if there are a lot of roughly equal $d(A_{i-1}, A_i)$ values, the chosen subset samples them in a roughly uniform manner with respect to the index i . The number 1.3 in the definition of σ is a fairly arbitrary value that quantifies what we mean by “roughly equal $d(A_{i-1}, A_i)$ values.” If inter-word spaces vary by more than a factor of 1.3, we trust that it is all right to concentrate on the largest spaces.

3.2 Bucketing strategies

Transformed image boxes B_0, B_1, \dots, B_{m-1} must be preprocessed so that, for any given box A , we can quickly

find the B_i that minimizes $d(A, B_i)$. It is natural to use bucketing because the B_i are distributed in a fairly uniform fashion across the non-blank areas of the page, and we can find a reasonable upper bound on the optimal $d(A, B_i)$. This is because the ideal boxes are easily ordered so that if A' follows A , the optimal B_i for A is likely to produce a low value for $d(A', B_i)$ as well.

Another way to guess a good B_i is to use information from a previous invocation of the mismatch function. If the transformation T has not changed much, the transformed image boxes B_0, B_1, \dots, B_{m-1} will not differ much from one invocation to the next. Hence we can try the B_i that corresponds to the one chosen last time for ideal box A . Thus we have two guesses for image boxes that are supposed to give low $d(A, B_i)$ values and these provide an upper bound on the actual minimum d value.

There are many ways to set up buckets so that a good upper bound on the minimum d value and a sufficiently uniform distribution of image boxes sharply limit the number of buckets and image boxes to be examined. One method is to classify boxes according to (x_1, y_1) and

$$\max \left(\frac{x_2 - x_1}{W}, \frac{y_2 - y_1}{H} \right), \quad (2)$$

where W and H are the width and height of the page and (x_1, y_1) and (x_2, y_2) are the minimum and maximum x, y coordinates for the box. The observed values of (2) are grouped so that one group has the small values and each other group covers at most a factor of two. Within each group, the bucket is determined by (x_1, y_1) in the natural fashion. Each time we go from one group to the next, values of (2) double and we reduce the number of (x_1, y_1) buckets by a factor of four.

4 Assigning ground truth

We have seen how to find a transformation (1) that makes the image boxes approximately match the ideal boxes that carry ground truth character identities from the original document description. If this were an exact one-to-one correspondence, it would be trivial to assign ground truth to the image boxes. When the correspondence is far from one-to-one, the most reliable approach is probably Kanungo's procedure of comparing against the expected character shapes with a small range of x and y displacements. Character shapes for template matching can also be extracted from the scanned image as suggested by Nagy et al. [12]. The purpose of this section is to consider what can be done if the ideal character shapes are not available.

The pleasant situation of a one-to-one correspondence is most relevant to the problem of assigning ground truth on a word-by-word basis, so we start with this problem. We then use information from the word-by-word problem to attack the more difficult problem of assigning character-level ground truth.

4.1 Ground truth at the word level

Consider the problem of identifying words in the original document description with the corresponding parts of the page image. We can assume that a word is described by consecutive records from the document description, where each record gives a character identity and the corresponding bounding box. Hence the task is to find word breaks in the document description, combine character bounding boxes to get a bounding box for each word, and then expect the transformed image boxes for each word to be approximately contained in the word's bounding box.

One way to decide if there should be a word break between two records in the document description is just to measure the distance between the character bounding boxes and test it against a threshold as do Chen et al. [3]. In our case, the original document description was derived from T_EX output so we use a version of Knuth's rule for finding word breaks in T_EX output [10,11]. Insert a word break between document description records A and B if

$$B_{x1} - A_{x2} > \frac{s}{6}, \quad B_{x1} - \bar{x} < -\frac{2s}{3} \quad \text{or} \quad |B_{y2} - A_{y2}| > \frac{5s}{6},$$

where $x1$ and $x2$ subscripts refer to the minimum and maximum x coordinates, $y2$ subscripts refer to the bottom y coordinate, s is the average of A 's font size and B 's font size, and \bar{x} is the maximum of A_{x2} and the C_{x2} values for all records C that precede A and follow the last word break. Another way to say this is that we break the document description into words by checking for horizontal spaces of more than $\frac{1}{6}$ of the font size or vertical spaces of more than $\frac{5}{6}$ of this size, while allowing $\frac{2}{3}$ of this size for backspacing when building up accented characters.

Refer to the bounding box of all the ideal boxes in a given word as a *word box*. We can assign transformed image boxes to words by just testing for approximate containment. An image box B is approximately contained in word box W if

$$d_f(B_{x1}, B_{x2}, W_{x1}, W_{x2}) + d_f(B_{y1}, B_{y2}, W_{y1}, W_{y2}) \leq \epsilon, \quad (3)$$

where d_f is as given in Sect. 2.1 and ϵ is a suitable threshold. This test is fast enough that we can try each transformed image box against each word box until finding a word box that approximately contains the image box.

4.2 Character level ground truth

If we want to find the portion of the page image that corresponds to each character in the original document description, the word-level ground truth allows the problem to be solved separately for each word. The correspondence between ideal boxes and transformed image boxes may be imprecise as illustrated in Fig. 5a, but this can be improved by doing an additional transformation specific to the current word. The bounding boxes of the

transformed image boxes assigned to the word (dashed lines in Fig. 5b) will not precisely match the word box (solid lines in Fig. 5b), but a transformation of the form

$$\bar{T}(x, y) = (\bar{t}_{xx}x + \bar{t}_x, \bar{t}_{yy}y + \bar{t}_y) \quad (4)$$

can fix this. The result is a better correspondence between ideal boxes and image boxes as shown in Fig. 5c.

The remaining task is to use this correspondence to map image boxes and the associated artwork to the ideal boxes. In case this cannot be done by comparing the observed images against the expected character shapes as Kanungo suggests, we now give an algorithm for doing this: Algorithm 1 uses an approximate containment test like (3), but with a tighter tolerance. We omit the details in Step 3 where run-together character images get cut up since this cannot be done reliably without knowledge of the character shapes. See Hobby [6] for an example of an application where Algorithm 1 was used successfully.

Algorithm 1 How to use an approximate containment test to take the character images assigned to a word and assign them to ideal boxes within the word, cutting up images if necessary.

1. Test each transformed image box for approximate containment in each of the word’s ideal boxes and assign the image boxes accordingly. If a transformed image box is approximately contained in more than one ideal box, choose the one where the d_f values are smallest.
2. For each unassigned image box B , find the set S_B of ideal boxes that intersect B and label B as a conglomeration of material from members of S_B .
3. If desired, try to cut up each image for box B from Step 2 cookie-cutter fashion, and assign the pieces to the appropriate ideal boxes.
4. Use the character identity and font labels from each ideal box to give ground truth labels for the images assigned that ideal box.

Another option is to use Algorithm 1 to extract prototypes for template matching. Combining everything the algorithm assigns to a given ideal box should produce the character image corresponding to that ideal box. For any character common enough to occur several times per page, throw out any unreliable instances and average the rest of them together. Character images containing the results of Step 3 should be considered unreliable and the same goes for cases where the approximate containment test in Step 1 would have failed if the tolerances were a little tighter.

The next step is to average the reliable character instances together. See [5] for a full discussion of how to do this. The result is a set of high-quality templates for all but the most uncommon characters. These can then be used in Kanungo’s template matching process.

5 Results

The algorithm was implemented in C++ and tested on 28 pages from six different documents. Each page was

printed at 600 dpi, photocopied once, and scanned or faxed at various resolutions. Table 1 gives basic information about the documents and the scanned images. The original document descriptions consisted of L^AT_EX output that was post processed to produce ASCII files listing character identities and bounding boxes for use by the program. This ASCII format was not specific to T_EX or L^AT_EX and could have been generated from a PDF or PostScript document description if a suitable conversion program were available.

In order to test the optimization procedures as thoroughly as possible, we start with a simple estimate for the transformation that matches the image to the document description, and then use all available means to reduce the mismatch. The goal is to get the mismatch below a critical value based on half the width of a typical character as suggested in Sect. 2.2.

The simple estimate for the transformation differs from Kanungo’s approach in how it tries to cope with bad feature points from speckles and miscellaneous noise in the page images. We find the four feature points illustrated in Fig. 1 for both the ideal boxes and the image boxes, then try to transform any three of the image box features into the corresponding ideal box features. This gives four possible transformations from which we can select the one closest to the identity transform. If none of these appear reasonable, we try again using additional constraints with pairs of feature points instead of triples. An additional heuristic attempts to cope with fax header lines that appear in the images but not in the document description.

Once the initial transformation is chosen, we can pick a subset of the ideal boxes as explained in Sect. 3.1 and then use the Nelder–Mead algorithm to minimize the mismatch as explained in Sect. 2.2. If the mismatch remains above the critical value, a second round of minimization uses the same starting point, but bases the mismatch on all of the ideal boxes instead of just a subset of them. If this fails to reach the critical mismatch value, it may be that certain ideal boxes cannot be matched.¹ These ideal boxes can be found by examining the last function evaluation in Round 2, and finding the ideal boxes that contribute the most to the mismatch. For instance one 200 × 100 dpi page image from Document A had mismatch contributions of 15.2 and 14.9 for the 85th and 84th ideal boxes, while no other ideal box had a mismatch contribution of more than 3.9. Round 3 consists of dropping as many as three such high-contributing boxes from the mismatch function and repeating the minimization process.

Figure 6 shows how this three-round minimization process performs on the test pages. Since the mismatch function is the L_4 norm of the contributions for the ideal boxes, the mismatch values used in the figure are normalized by dividing by the fourth root of the number of ideal boxes. This means that the coordinates in the figure are essentially in pixel units, so the higher “after minimization” values for the 400 dpi test pages just

¹ This can happen when character images merge with large connected components such as rule lines in a table.

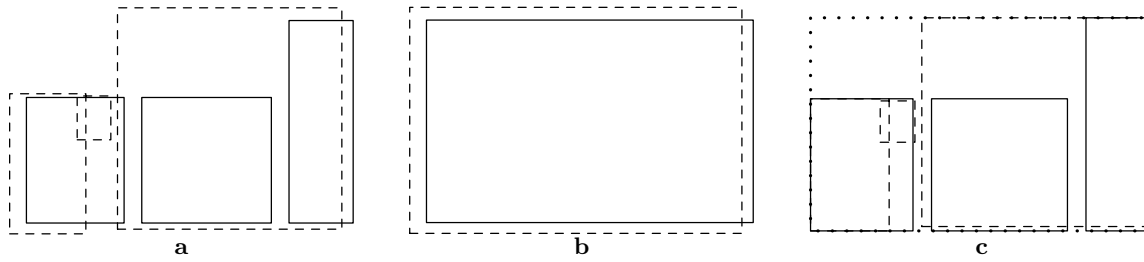


Fig. 5a. Transformed image boxes (*dashed lines*) superimposed on the corresponding ideal boxes (*solid lines*); **b** the word box (*solid lines*) and the bounding box of the transformed image boxes (*dashed lines*); **c** The boxes from **a** with the image boxes transformed so that the bounding boxes from **b** coincide to form the *dotted* box

Table 1. Statistics about the test pages and the character images extracted by page layout analysis. Documents *A–C* are preprints of journal articles and documents *E–G* are software manuals. The “Truth” column lists the average characters per page from the original document description and the last three columns count connected components in the scanned images

Document				Average characters or components			
Id	Font	Language	Pages	Truth	200 × 100	200 dpi	400 dpi
<i>A</i>	cmr 10pt	English	5	1715	2100	1819	1838
<i>B</i>	Times 10pt	English	6	2371	2655	2453	2519
<i>C</i>	cmr 12pt	English	5	1272	1839	1423	1379
<i>E</i>	cmr 10pt	English	4	1879	2533	2026	2016
<i>G</i>	cmr 11pt	German	4	1643	2257	1831	1821
<i>S</i>	cmr 11pt	Spanish	4	1453	1916	1564	1540

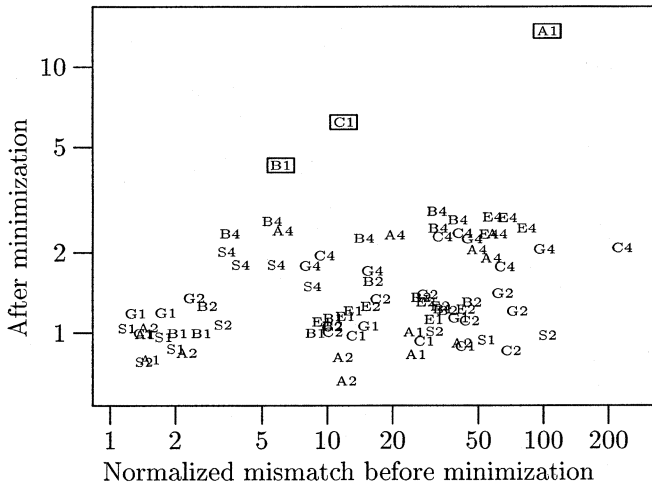


Fig. 6. A scatter plot of the normalized mismatch before and after the three-round minimization process for pages from the test documents. Each label gives a document id from Table 1 followed by a digit that gives the resolution: 1 means 200 × 100 dpi; 2 means 200 dpi; 4 means 400 dpi. Boxes identify pages where the critical mismatch value could not be achieved. All data is for Nelder–Mead minimization

mean that the actual distances involved in the mismatch computation are staying fairly constant. (The mismatch values in Fig. 3 would have to be divided by $\sqrt[3]{300} \approx 4.16$ in order to have the same units.)

The three boxed labels in Fig. 6 are for 200 × 100 dpi test pages where Nelder–Mead minimization failed to reach the critical mismatch value. Since all other test pages led to less than critical mismatch values, the 200 × 100 dpi resolution images appear to be the primary trouble spot for Nelder–Mead. This is relatively encouraging

in view of the high mismatch values before minimization. Except for the small group of 100 × 200 and 200 dpi labels at horizontal positions < 3.5 and the small group of 400 dpi labels at horizontal positions < 7, it was apparently not possible to find three appropriate feature points on which to base the initial transformation. Note that minimization often reduced the normalized mismatch by a factor of two even when the initial mismatch value was good enough to suggest that the feature points were appropriate.

If the mismatch values after minimization scale with resolution, is there some underlying cause? We test this by viewing the mismatch graphically as indicated by the transformation (4) for each word. Figure 7 shows how the center of each word in the 200 dpi scanned image gets shifted for the two pages from Document *A*; i.e., it displays the additional transformations needed after the best affine transformation has been applied. These transformations are complicated but they do not look like random noise. If the similarities between parts **a** and **b** of the figure are due to the fact that they both originated from the same copier and fax machine, it may be possible to measure the nonlinearities and compensate for them in future experiments as suggested by Kanungo et al., [8,9]. On the other hand, Sect. 4.2 explained how to find the transformations (4) without such prior calibration if the mismatch is not too great.

In view of the importance of getting the mismatch below the critical value, it is worth checking if alternative minimization strategies can reduce or eliminate instances of failure to reach the critical value. One approach is to try other initial transformations. For instance, the failures in Fig. 6 can be eliminated by starting at the transformation that turned out to be optimal for the previous page of the document in question. Of course it would be

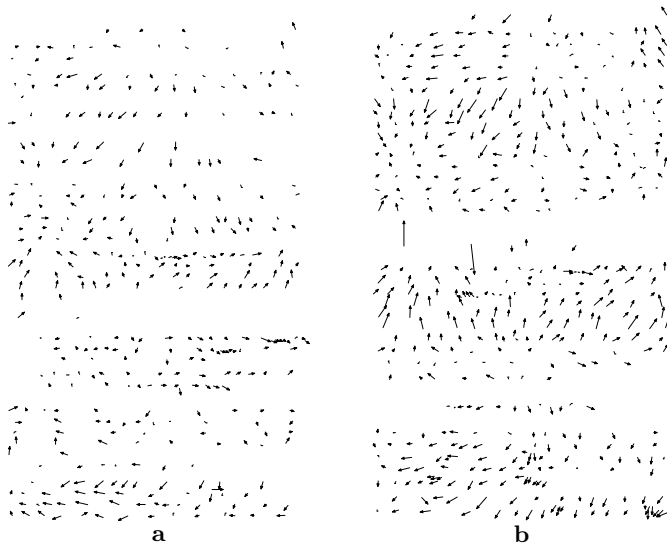


Fig. 7a,b. Displacements for the center of each word's bounding box due to the transformation (4) when matching 200 dpi images to the original document descriptions for two pages of Document A. The displacement vectors are exaggerated by a factor of 20

better to find an optimization strategy that outperforms Nelder–Mead on the 200×100 dpi pages.

The guaranteed convergence properties of Torczon's algorithm make it a prime candidate for these difficult, low-resolution test pages. Figure 8 compares the effectiveness of Nelder–Mead minimization with three versions of Torczon's algorithm. The three versions differ in the setting of a parameter that controls the number of function evaluations per iteration of the optimization routine. Torczon refers to this parameter as “the number of search directions.” The minimum allowable value of 12 gives rise to the solid round dots in Fig. 8; the small open circles are for 24 search directions; and the large open circles are for 96 search directions. The Nelder–Mead data indicated by + signs in the figure do not show a dramatically different distribution.

Figure 8 shows that the mismatch after minimization was always either less 1.5 or more than 4. This separates the data points into those where the mismatch was successfully reduced to the critical value, and those where the minimization process should be considered unsuccessful. Since the most notable feature of this latter group is that dots (Torczon's with 12 search directions) are overrepresented, we tentatively conclude that it is probably best to use more than 12 search directions.

Table 2 shows how Torczon's algorithm compares with Nelder–Mead in terms of the number of function evaluations needed to reach the minimum. The “Overall” row shows that the overall average number of function evaluations is much less for Nelder–Mead. Torczon's algorithm performs best with 24 search directions, but even in that case, Nelder–Mead does fewer than half as many function evaluations. It also helps to increase the resolution. The overall averages are reduced at 200 and 400 dpi because it is seldom necessary to resort to more than one round of minimization.

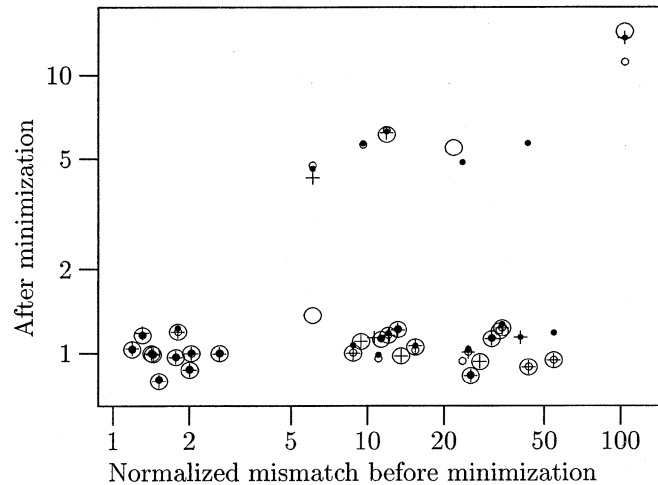


Fig. 8. A scatter plot of the normalized mismatch before and after the three-round minimization process for 200×100 dpi pages from the test documents. The + signs are for Nelder–Mead minimization, and the dots, small circles, and large circles are for Torczon's algorithm with 12, 24, and 96 search directions, respectively

Table 2 may be a little misleading because it weights iterations in Round 1 just like those in Rounds 2 and 3 even though the later rounds make the mismatch function harder to evaluate by basing it on all or almost all the ideal boxes instead of using just 15% of them. Hence, the actual run times in Table 3 show an even bigger speed up at the higher resolutions. Torczon's algorithm still takes more than twice as long, but this could change if parallelism were enabled. (The test machine was an SGI Challenge XL with 12 MIPS R4400 processors running at 150 MHz, but all tests were done on a single processor).

What about the accuracy of the ground truth produced by the algorithm? Kanungo has shown that reliable ground truth can be obtained by template matching against the expected character shapes, once the transformation is known. To test how well the algorithm performs without template matching, it was used to evaluate imperfectly segmented characters derived from the 200 dpi pages described in Table 1. The method of [2, 7] was used to segment each page image into characters, then the resulting character boxes were used as the image boxes and matched against the ideal boxes from the \TeX output. Thus each image box was labeled as corresponding to zero or more ideal boxes or parts thereof, and these labels were checked by hand.

Out of a total of 48 553 image boxes, 2529 were not labeled with any ideal boxes. Most of these rejected image boxes were from fax header lines or other parts of the page images where the \TeX output did not specify any characters. However, 624 of the rejected image boxes could have been labeled. This high number of false rejects was due to conservative tolerances designed to reject questionable cases rather than risk generating bad output.

A total of 2428 of the image boxes were correctly labeled as more than one character and 168 image boxes

Table 2. Average number of times the mismatch function had to be evaluated to process a test page broken down by resolution, optimization strategy, number of rounds of optimization needed, and success in reducing the mismatch below the critical value

Scenario	200 × 100 dpi				200 dpi	400 dpi
	Torczon 12	Torczon 24	Torczon 96	NM	NM	NM
Round 1 works	1596	899	2748	307	335	312
Round 2 works	4706	2566	6302	1229	1226	1366
Round 3 works	9525	2673	8853	1918	–	–
Round 2 fails	4485	2590	9134	1648	–	–
Round 3 fails	5361	4149	12405	–	–	–
Overall	2924	1501	4457	657	399	387

Table 3. Average run time in seconds per test page as a function of resolution, optimization strategy, number of rounds of optimization needed, and success in reducing the mismatch below the critical value

Scenario	200 × 100 dpi				200 dpi	400 dpi
	Torczon12	Torczon24	Torczon96	NM	NM	NM
Round 1 works	57.6	34.0	96.2	12.5	12.8	13.2
Round 2 works	285.9	169.6	631.7	128.1	77.0	89.9
Round 3 works	833.6	251.5	803.4	158.9	–	–
Round 2 fails	275.3	203.6	568.4	124.5	–	–
Round 3 fails	321.8	358.9	1145.2	–	–	–
Overall	169.5	93.9	281.4	44.3	17.3	18.7

were correctly labeled as only part of a character. There were 6 cases where two characters ran together but they were labeled as as a single character plus “something that got rejected.” In other words, the algorithm indicated that the 6 objects were too big to be single characters but it could only find one ideal box that clearly matched the image box.

Only 5 image boxes were truly mislabeled: 2 image boxes that should have been labeled as periods or commas were labeled as part of a neighboring character; and 3 image boxes resulting from a character merging with part of an adjacent character were labeled as single characters. This corresponds to an error rate of about one part in 10^4 if you exclude the rejections and other exceptional conditions described above. Using Algorithm 1 to extract templates for template matching as suggested at the end of Sect. 4.2 would have eliminated all of these errors.

6 Conclusion

When generating ground truth by matching a page image with the original document description, it is of primary importance to find a geometric transformation that maps coordinates appropriately. We have seen how standard optimization techniques can improve the accuracy of such a transformation, even if the initial approximation is way off. In those few cases where optimization does not yield an appropriate transformation, the mismatch function reveals the problem, and restarting the optimization with a different initial transformation can solve the problem.

The mismatch function presented in Sect. 2.1 could be generalized to deal with other applications. This should allow the basic idea of using optimization techniques to match a scanned image with the original document description to apply to any sort of text or graphics. For

instance, if the page contains line graphics and the document description makes it clear where all the ink should be placed, the mismatch function can just look at some of the places where black pixels should be found and measure the distance to the nearest black significant black area in the image. Such a modification in the mismatch function may also yield better results on low-resolution text where the correspondence between characters and bounding boxes of connected components sometimes begins to break down.

Consider the statement by Kanungo and Haralick that their methodology is general enough to handle documents in any language [9, p. 674]. The algorithm presented here works for any language where bounding boxes of connected components are reasonably well correlated with character bounding boxes. In other words, the mismatch function would have to be generalized somehow in order to cope with a language like Arabic where characters are not separated by white space.

A number of important questions remain to be more fully addressed in future work. Does searching for an affine transformation (1) provide the right number of degrees of freedom? Three degrees of freedom suffice to determine how a piece of paper is positioned on a scanner, yet affine transformations provide six and Fig. 7 suggests that it would help to allow much more complicated transformations. Affine transformations were chosen as a compromise because direct search optimization algorithms have a reputation for behaving poorly when the dimensionality of the search space is too high.

Another question is what optimization algorithm is best. Since Torczon’s algorithm is designed to run in parallel, its run times could be significantly better than Nelder–Mead if parallelism were enabled. On the other hand, the mismatch function itself could be parallelized since each ideal box can be considered separately. The tests reported in Sect. 5 are less than definitive, but the

difficulties that occasionally cause Nelder–Mead to fail to find the desired minimum do not seem to be amenable to Torczon’s superior convergence properties. The mismatch function could well have undesired local minima.

A possible alternative to direct search methods such as Nelder–Mead and Torczon’s algorithm is to try to use differential semblance optimization to construct a smoother or more continuous function to minimize as suggested by Symes [18], Gockenbach [4].

It also seems promising to try a real segmentation algorithm instead of just finding bounding boxes of connected components. The third round of optimization used in the trials in Sect. 5 was specifically designed to cope with missed segmentations, but it would undoubtedly be better to improve the segmentation. This would be especially important when dealing with languages where characters tend to run together or languages where characters can contain many small connected components.

Finally, it would help to be able to extract symbol identities and bounding boxes from a PostScript or PDF file, instead of depending on \TeX or \LaTeX output. This would make it easier to generate a wide range of input for the ground truthing process and it would probably allow for more accurate bounding boxes. We chose \TeX output simply because it is easy to parse – PostScript or PDF would work as well, provided that it contains the required symbol identities.

Acknowledgements. Virginia Torczon and David Serafini provided an implementation of Torczon’s algorithm and adjusted it to be more easily called as a subroutine. Margaret Wright’s careful implementation of the Nelder–Mead algorithm was also very helpful.

References

1. Henry S. Baird. The skew angle of printed documents. In: Proc. SPSE 40th Symp. Hybrid Imaging Systems, pp 21–24, Rochester, NY., May 1987
2. Henry S. Baird. Anatomy of a versatile page reader. Proceedings of the IEEE 80(7): pp 1059–1065, 1992. Special Issue on OCR
3. Su Chen, Robert M. Haralick, Ihsin T. Phillips. Perfect document layout ground truth generation using DVI files and simultaneous text word detection from document images. In: Proceedings of the Fourth Annual Symposium on Document Analysis and Information Retrieval, Las Vegas, NV, April 1995
4. Mark S. Gockenbach. An abstract analysis of differential semblance optimization. Technical Report TR94–18, Dept. of Computational and Applied Mathematics, Rice University, 1994
5. John D. Hobby, Henry S. Baird. Degraded character image restoration. In: Proceedings of the Fifth Annual Symposium on Document Analysis and Image Retrieval, pp 233–245, 1996
6. John D. Hobby, Tin K. Ho. Enhancing degraded document images via bitmap clustering and averaging. In: ICDAR’97: Fourth International Conference on Document Analysis and Recognition, August 1997

7. D. J. Ittner, Henry S. Baird. Language-free layout analysis. In: Proceedings of the Second International Conference on Document Analysis and Recognition, pp 336–340, Tsukuba Science City, Japan, October 1993
8. Tapas Kanungo. Validation of Document Degredation Models. PhD thesis, University of Washington, 1995
9. Tapas Kanungo, Robert M. Haralick. Automatic generation of character groundtruth for scanned documents: a closed-loop approach. In: Proceedings of Int. Conf. on Pattern Recognition, Vol. C, pp 669–675, Vienna, Austria, 1996
10. Donald E. Knuth. The \TeX book. Reading, MA: Addison Wesley 1986. Vol. A of Computers and Typesetting
11. Donald E. Knuth. \TeX ware. Technical Report CS-TR-86-1097, Dept. of Computer Science, Stanford University, April 1986
12. George Nagy, Yihong Xu. Priming the recognizer. In: Jonathan J. Hull, Suzanne Liebowitz Taylor (eds) International Workshop on Document Analysis Systems (DAS’96), pp 263–281, 1996
13. J. A. Nelder, R. Mead. A simplex method for function minimization. Computer Journal 7: 308–313, 1965
14. Ihsin T. Phillips, Su Chen, J. Ha, Robert M. Haralick. English document database design and implementation methodology. In: Proceedings of the Second Annual Symposium on Document Analysis and Information Retrieval, pp 65–104, April 1993
15. Ihsin T. Phillips, J. Ha, Robert M. Haralick. Implementation methodology and error analysis for the University of Washington English Document Image Database-I. In: Proceedings of the SPIE, Vol. 2103, pp 155–173, 1994
16. R.P. Rogers, Ihsin T. Phillips, Robert M. Haralick. Semi-automatic production of highly accurate word bounding box ground truth. In: International Association for Pattern Recognition Workshop on Document Analysis Systems, pp 375–387, October 1996
17. A. Lawrence Spitz. Skew determination in CCITT group 4 compressed document images. In: Proceedings of the Symposium on Document Analysis and Information Retrieval, pp 11–25, April 1992
18. William W. Symes. Velocity inversion: a case study in infinite-dimensional optimization. Mathematical Programming 48: 71–102 (1990)
19. Virginia Torczon. PDS: Direct search methods for unconstrained optimization on either sequential or parallel machines. Technical Report CRPC-TR92206, Rice University Center for Research on Parallel Computation, 1992
20. Margaret H. Wright. Direct search methods: once scorned, now respectable. In: D. F. Griffiths, G. A. Watson (eds) Proceedings of the 1995 Dundee Biennial Conference in Numerical Analysis, pp 191–208, Harlow, UK: Addison Wesley Longman, 1995

John D. Hobby received his B.S. in mathematics and computer science from the University of Washington in 1980, and a Ph.D. in computer science from Stanford University in 1985. Since then, he has been a member of technical staff at Bell Labs in Murray Hill, New Jersey. His research interests involve problems in document analysis including techniques for improving degraded images. He has also published research articles on font generation, splines, graphics algorithms, numerical stability, and computational geometry, and he has created a graphics language called MetaPost.