

系統図自動入力における文字列検出方式

5M-9

七野 剛, 二俣 隆, 柴山 純一, 前田 暲
三菱電機(株) 情報電子研究所

1. はじめに

既存の系統図を認識し、効率よく設備情報をデータベース化するシステムを開発中である[1]。文字列と図形は接触して記述されることがあり、このような図面からでも正しく文字列の位置情報を検出する方式について報告する。

2. 文字列検出における問題点

対象とした系統図は、計算機入力を目的として作成されたものではなく、特に、文字については厳密な記入制限がない。このため、文字列検出を行う場合、次に示す問題が発生する。

①文字要素と図形要素(ライン・シンボル)の接触

文字要素とライン、文字要素とシンボルが、接触することがある(図1 a, b)。孤立した文字要素を統合する方法では、文字列を正しく検出することができない。

②異なる文字列との近接あるいは接触

水平方向に近接して記述されることがある(図1 c)。文字列間隔よりも文字間隔が大きいことがあるので、文字列の境界を検出するのが困難である。

垂直方向に接触して記述されることがある(図1 d)。接触の度合いが様々であり、また、文字サイズが1種ではないため、接触の境界を検出するのが困難である。

3. 処理方式の特長

上記の問題を解決するために開発した文字列検出方式の特長を示す。

①文字列の検出に先立ってラインおよびシンボルの認識を行い、認識できた図形要素を除いたベクトルデータの中から文字列の検出を行う。ラインおよびシンボルの認識は、水平垂直の長い線(ライン)・特定の大きさの黒領域が近接して存在する(シンボル)など、確実な特徴に基づいて行われる。このため、文字要素と図形要素が接触していても分離が正しく行なえる。

②文字列の検出は、近接する文字要素を統合して行なうが、統合の後、垂直方向に接触した文字列を分割する。また、統合・分割の際、シンボル認識の結果を有効に利用する。このため、異なる文字列が近接あるいは接触していても文字列単位に検出できる。

4. 文字列検出

処理は、図面の画像を細線化およびベクトル化したベクトルデータを用い、以下の順に進められる。

1)要素分離; 先に認識されている図形要素と、明らかに文字となり得ない短い線(ノイズ)、長い線(認識仕様でない引き出し線など)を図面全体のベクトルデータから除き、これを文字要素とする。

2)統合処理; 近接する文字要素を統合して文字列の候補を作成するが、近接していても「くいちがい」の度合いが大きい場合(図1 c)は統合しない。また、シンボル内の文字列と外の文字列が近接して記述されていることがある(図2 a)が、一個の文字列に統合されてしまうを防ぐため、破線で示したシンボルが認識されている場合は内外の要素を統合しない。これにより、異なる文字列が水平方向に近接していても文字列単位に検出できる。

3)分割処理; 文字列候補が規定の高さを越えた場合、水平方向に文字列候補を構成する線の数を加算したヒストグラムを作成する。そして、ヒストグラムの値が最も小さいところで分割を行う。また、「弁記号には、規定の文字高さがh1の文字列とh2の文字列の2行が近接している」などの規則がある。従って、文字列候補が認識された弁記号に近接していれば、規定の文字サイズにより境界を決定し分割する(図2 b)。これにより、異なる文字列が垂直方向に接触していても文字列単位に検出できる。

5. おわりに

文字列が図形と接触、異なる文字列と近接あるいは接触していても正しく検出する方式を報告した。これにより、系統図入力が効率的に行えるシステムを開発できる。

[1] 柴山他, 「系統図自動入力における対話修正方式」
情報処理学会第39回全国大会 3E-3

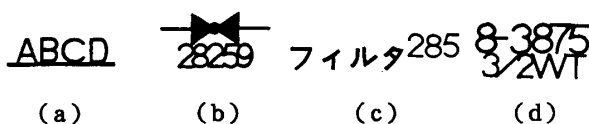


図1. 文字列検出において問題となるパターン

Character Strings Segmentation for Automatic Input of Plant Diagrams

Go SHICHINO, Takashi FUTATSUMATA, Jun-ichi SHIBAYAMA and Akira MAEDA Mitsubishi Electric Corporation

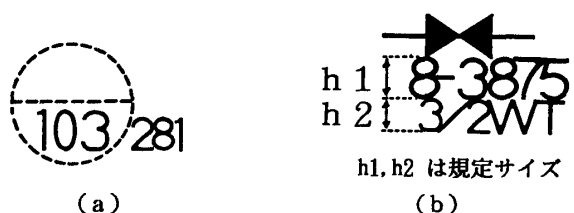


図2. 文字列の分割