

Webデータ対応リアルタイムデータマイニングツールの開発 データマイニングを身近なものに

1. 背景

大量のデータから、隠れた知識を生成する方法として、最近データマイニングという手法が注目されているが、一般にはあまり普及していない。

普及を阻んでいる理由としては、費用的な面と、ツールを扱うには統計の専門家やプログラミング力が必要とされることなどが考えられる。

また、Web上の膨大なデータの中には、今まで気が付かなかった知識が隠れている可能性が高いが、Webから<Table>タグ内の表データを含め多様な形式のデータを取り込めるデータマイニングツールは見当たらない。

2. 目的・ねらい

統計の専門化だけにとどまっているデータマイニングの市場を、一般層まで広げたい。また、従来のデータマイニングツールには無い、Web上のデータを取り込める機能や、マクロやプログラミングを使わなくてもWebアプリと連携して予測モデルを自動更新できる機能などを実現し、データマイニングの応用範囲を広げたい。

本ソフトはSAS社やSPSS社などから出されているような統合的なツールを目指すのではなく、一般の人が気軽に使えて、データマイニングのメリットを享受できるものを目指す。しかも、実用に耐えるデータ量(エクセルレベル以上)を扱えて、効率的な分析アルゴリズムを備えるソフトにしたい。データマイニングの手法としては、当面は利用頻度の高い予想系の重回帰分析を採用し、将来はバスケット分析など他の分析機能を追加する。

3. 開発の内容

(1) HTML, CSV, TSV など様々なファイル形式で存在する表データを、インターネットやディスクなどから容易に入力し、分析結果を得ることができる。

分析対象データは簡単な操作で、格納場所(Web、PC)やフォーマットなどをあまり意識せずに取り込むことができる。取り込み後の分析から予想モデルの作成、分析結果の簡単な解説文出力までを本ソフトが一気に処理する。そのため、ユーザはマウスで予想対象列を指定するだけで、膨大なデータから傾向や法則などを容易に得ることができる。

なお、データ入力方法は次のとおりである。

ファイルを画面上にドラッグ&ドロップ (D&D)
CSV ファイルと TSV ファイル (タブ区切り) は拡張子によらず自動認識
で読み込める。特殊な区切り文字にも対応している。

インターネットショートカット (拡張子: '.URL') を D&D
ブラウザのインターネットショートカットから本画面上に D&D すると、
本ソフトの通信機能が自動的に働いてネット上のデータを収集できる。

「貼り付け」機能を利用
他のソフトからコピー&ペーストで表データをダイレクトに貼り付け
ることができる。

プロジェクトファイルを利用
本ソフト専用のプロジェクトファイル (拡張子: '.DPR') に、ネット
から収集すべきファイル名 (複数可) や FTP パスワード等を事前に記述
しておくことにより、バッチ的にデータを取り込める。HTTP と FTP の
両方に対応している。

既存の表収集ソフト「<Table>バインダ」との連携
開発者作成の他のソフト「<Table>バインダ」との連携機能があり、
これを利用して、HTML ファイル内の<Table>タグ内の表データを分析デ
ータとして取り込める。

(2) マクロやプログラミングを使わずに、容易に予測モデルの作成とそ
れを使ったシステムを構築できる。また、予測モデルの更新作業を自動化
できる。

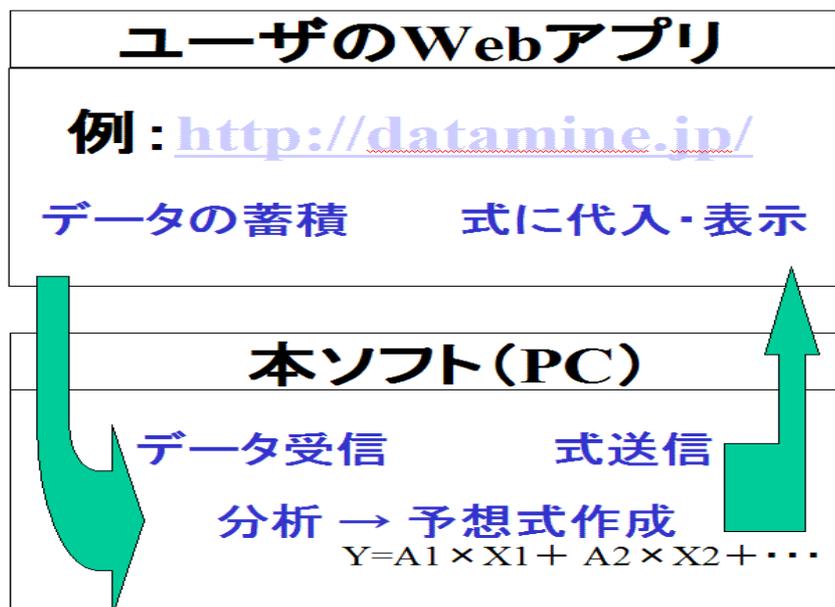


図 1

Web アプリとの連携例として今回作成したのが、「当たるも八卦」(Web 側 TOP 画面 <http://datamine.jp/>) というシステムである。これは Web アプリ側 (データ蓄積) と、リモート PC 側 (予想モデル作成) の連携動作で動いている。(図 1 を参照)

Web 上で蓄積されたデータを、本ソフトに取り込み、分析して予想式(つまり偏回帰係数と定数項)を決定し、その係数ファイルを Web サーバーに送信する。この取り込みから送信までの一連の流れは、プロジェクトファイルで簡単に設定できる。ちなみに「当たるも八卦」は、OS のタスクスケジューラに設定されており、2 日毎に予想式が自動更新されている。

この Web 連携機能を利用すると、Web 対応の分析予想系のシステムを構築する場合に、ユーザは Web 側のアプリを作成するだけで良い。但し、前処理がほとんど不要なデータを吐き出すように Web アプリを作成しておく必要がある。

(3) 予測因子を A I C 法等により適切にかつ自動的に選択し、予測信頼性の高い最適なモデルを作成できる。

説明変数同士の相関が強いデータが存在する場合、多重共線性の関係でそのまま分析すると予測信頼性が落ちるが、本ソフトでは重回帰分析処理の前に相関の強いデータの一方を自動的に外し、さらに A I C 法と変量増加法の組み合わせで多重共線性を回避している。

(4) 解析結果に悪影響を及ぼす異常データを除去する

数字以外のデータが入っている行を自動的に除外するようにした。

4. 従来の技術 (または機能) との相違

(1) 多様な形式の Web データを取り込める。

Web から <Table> タグ内の表データを含め多様な形式のデータを簡単な操作で収集できる。

(2) Web アプリとの連携機能がある。

本ソフトを分析、通信部品として利用すると、Web 側のアプリ開発とプロジェクトファイルの作成だけで分析・予想系の Web システムを容易に構築できる。

(3) 操作が簡単である。

データ入力から分析結果の出力までの操作は、D & D、[編集] ボタンクリック、予想対象列クリック、及び [個別分析] ボタンクリック

だけである。

の表画面では、必要に応じて表やデータを修正できる。

(4) 分析結果を文章で出力する。

重回帰分析の結果を簡単な文章で出力する。

例：<http://datamine.jp/report2-1.htm>

(5) 分析計算が高速である。

最適なモデルを自動選択するための独自のロジックを採用し、今回これをさらに改善した。図2において、縦軸は所要時間、横軸はデータ量である。左図は従来手法、右図は独自の手法を使ってそれぞれ処理時間を計測したものである。

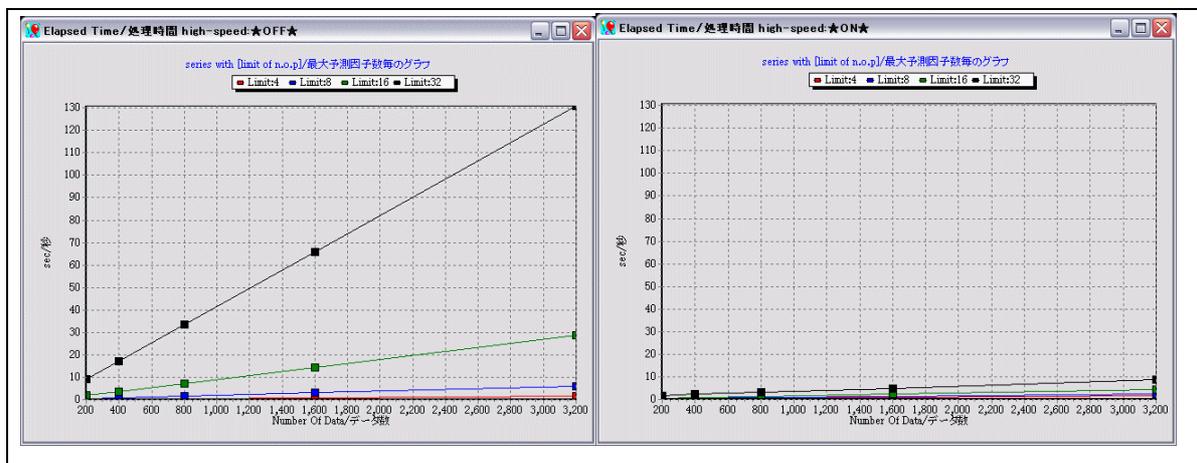


図2

5. 期待される効果

統計の専門化だけにとどまっているデータマイニングの市場が、一般層まで広がることが期待できる。また、Web アプリとの連携システムなど、データマイニングの応用範囲が今よりも広がると思われる。

6. 普及（または活用）の見通し

Web 上にシェアウェアとして公開し、だれでも気軽に本ソフトを利用できるようにしたい。本ソフトのダウンロード数は、開発者が作成した既存ソフトのダウンロード数の5万程度を見込む。

7. 開発者名（所属）

若松 桜男（有限会社ソフトポート）

(参考) 開発者URL <http://www.softport.co.jp/>