宇宙航空研究開発機構 宇宙科学研究本部

科学衛星運用・データ利用センタ-

Center for Science-satellite Operation and Data Archive

2009.11.27 / No. 193

ISAS/JAXA

衛星データとデータマイニング 第2回 データマイニングの実際

本田 理恵(高知大学自然科学研究系理学部門 / JAXA 宇宙科学情報解析研究系 客員)

## 1. はじめに

前回のデータマイニングの概略の説明に引き続き、今回 はオープンソースソフトウェアの Weka を用いた学習事例 を交えながら代表的な手法について具体的に紹介します。 事例としては、はじめての方に衛星データへの活用のイ メージをつかんでいただくために、簡潔な例を用いました が、教科書的記述が多くなってしまったことにはお詫びい たします。また、紙面の制約のため図が小さく見にくくなっ てしまいましたが, web 版には大きめに掲載されています のでそちらもご参照ください。

PLAIN News

http://www.isas.jaxa.jp/docs/PLAINnews/

# 決定木学習

機械学習やデータマイニングでは、"分類"に属する問 題が数多く扱われます。"分類"は,複数の属性を持つデー タにおいて、ある属性を目的属性、それ以外の属性を従属 属性として、従属属性だけから目的属性を決定する問題と してとらえることができます。機械学習の分野では、目的 属性をラベルまたはクラス、従属属性をパターンと呼ぶこ ともあります。

表1に、分類で扱われるデータセットの一例を示します。 ここでは、衛星のテレメトリデータを想定して、搭載セン サの誤動作の有無, 電圧, 温度, 他のセンサの動作の有無 を属性とします (注1)。このようなデータセットから、セ ンサの誤動作がおこる条件を知りたい、という要求がある ものとします。この場合,衛星の誤動作の有無を目的属性, それ以外を従属属性とすることになります。

11	/ / Ľ		ビング研究所用	1 FIX
温度	電圧 A	電圧 B	他センサ	誤動作
高	高	高	無	有
高	高	高	有	有
中	高	高	無	無
低	中	高	無	無
低	低	中	無	無
低	低	中	有	有
中	低	中	有	無
高	中	高	無	有
高	低	中	無	無
低	中	中	無	無
高	中	中	有	無
中	中	高	有	無
中	高	中	無	無
低	中	高	有	有

1 データセットの例-	センサ誤動作情報-
-------------	-----------

"分類"を扱う手法には、ナイーブベイズ、サポートベクター マシンなど種々のものがありますが、ここでは決定木学習 を取り上げます。表1のデータに対して求められた決定木 の例を図1に示します。決定木の根または節には従属属性 に対する条件式が割り当てられます。一方、葉には目的属 性が割り当てられます。よって、目的属性の値が未知で従 属属性が既知のデータが与えられた場合,従属属性の値を 参照しながら決定木を根から葉へたどっていくことによっ て、この事例の目的属性値を推定することができます。



図1.表1のデータに対して得られた得られた決定木の例。 括弧内数字は実際に分類された事例数。

このような木構造を、データからの学習によって形成し ようとするのが決定木学習です。決定木学習では、まず全 データに対して考えられる全ての分割条件を用いてそれぞ れ分割を試行し、目的属性に関しての分割前と分割後の データセットの多様性指標の変化量を計算します。ここで, 多様性指標の減少量がもっとも大きくなる条件を最初の分 割条件として採用してデータセットを分割します。この過 程を分割後の各データセットが目標属性について均一にな るか、それ以上の分割条件がなくなるまで再起的に実施し ていきます。多様性指標には様々な物が有りますが、代表 的なアルゴリズムである Quinlan の C4.5 や C5.0 では情報 エントロピーに基づく情報利得比が用いられています。

また、木構造は, if then 形式のルールで表現する こともでき、このルール自体が現象に関する知識で あって、決定木学習によってこの知識を発掘するこ とができた、と解釈することもできます。 図1の決 定木では,"if temperature=high and V\_B=high then sensor failure=true, if temperature = low and otherSensor=active then failure=true となり、このセ ンサが不具合を起こすのは、温度が高くB点の電圧が高い とき、または、温度が低く他のセンサが動作しているとき、 ということになります。

# 3. 相関分析

発生した事象の共起関係を調べたい場合には相関分析が 有用です。相関分析は、トランザクションに現れるアイテ ム間の共起関係を分析します。顧客の購買分析とのアナロ ジーで説明すると、トランザクションは顧客の買い物籠(バ スケット). アイテムはバスケットにはいった商品として とらえることができます。

このようなデータから相関ルールと呼ばれる共起関係を 抽出します。相関ルールはA→Bという形で表記され,A という事象が発生すれば高い確率でBという事象も発生す ることを示します。相関ルールの重要性の指標には支持度 (support), 確信度 (confidence) があります。Support(A  $\rightarrow$  B) は全トランザクション数に対するAとBが共起す るトランザクション数の割合を表し, confidence(A→ [裏へ続く]

B) はAを含むトランザクション数に対するAとBが共 起するトランザクション数の割合を示します。よって confidence(A  $\rightarrow$  B) と Support(A  $\rightarrow$  B) が大きいルールが 重要なルールということになります。

ルールの抽出にあたっては、アイテムの数が増加する ことによって組み合わせ数が膨大になって計算が困難に なるという問題がありますが、IBM アルマディン研究所の Agrawal によって開発された Apriori アルゴリズムによっ て大規模データから効率よくルールを抽出することが可能 になっています。

#### 4. クラスタリング

この他,一般的な用途によく用いられる手法としてクラス タリングがあります。クラスタリングは,数値やカテゴリ 値からなる多次元ベクトルを対象とし,類似度に基づいたグ ループ分けを行う手法です。数値属性のベクトルについて は,類似度の代わりに相違度としてベクトル間のユークリッ ド距離がよく使用されます。クラスタリングは,決定木学習 のようにあらかじめ目標属性のようなお手本が与えられない ため,教師無し学習の一つです。

具体的な手法には、階層的な手法と非階層的な手法があ ります。階層併合的クラスタリングは、1データを1クラ スタに割り当てた状態を初期状態とし、データを相違度の 小さいものから順番に併合して最終的に一つのグループに なるまで反復します。併合過程のデータ集合をまとめた データ集合は樹状図をなし、様々な階層レベルでグループ を評価することができます。

一方,非階層的な手法には K-means 法などがあります。 K-means 法ではあらかじめグループの個数 k を定め,K 個 のクラスタの重心の初期値をランダムに与えます。次に各 データについて,全てのクラスタ重心との距離を計算して 最も距離のクラスタにデータを割り付けます。全てのデー タの割り付けの完了後,各クラスタの属するデータセット からクラスタの重心を決定し直します。この過程を決めら れた回数,またはデータの重心の変化がなくなるまで反復 することによって,尤もらしいグループ分けを得ることが できます。

K-means 法は非常に明快なアルゴリズムですが,K をあ らかじめ与えなければならない,初期値によって結果が大 きく異なる,といった欠点も持ちます。使用時にはこうし た特徴に対する注意が必要です。

この他,多変量正規分布の混合モデルでの近似もよく用いられます。学習(パラメータ推定)には EM アルゴリズムが用いられます。

なお,数値属性データではクラスタリングの実施前に, 異なる属性値間の寄与を等しくするために各属性の間で正 規化を行うことが必要な場合も有ります。

## 5. Weka によるデータマイニング実例

それでは前回紹介したニュージーランド大で開発された データマイニングや機械学習のための統合的ソフトウェア Wekaを使用して,前述の代表的な手法によるデータの学 習の実例を紹介します。

インストール まずは、http://www.cs.waikato.ac.nz/

ml/weka から使用環境に対応したソフトウェアをダウン ロードしてインストールします。2009 年 11 月現在, 最 新の stable バージョンは 3.6.1 で, Windows, Mac OS, Linux などの OS 向けの実行形式が準備されています。実 行にあたっては Java VM 1.5 以降を必要としますので適 宜準備をしてください。なお,上記サイトにはチュート リアル, FAQ,サンプルデータ,メーリングリストへの参 加法など各種情報が掲載されています。

データの準備 Wekaの入力データの標準的なフォーマットはARFFです。図2に表1のARF形式のファイル" sensor.arff"を示します(ファイルの拡張子にarffをつけます)。ARFFファイルはヘッダ部とデータ部からなり、ヘッダ部は@dataまでです。

ヘッダ部の@relationの後にはテープルの名前を,@ attributeの後には属性の名前とタイプを記述しま す。カテゴリ値の場合は {high, mid, low}のように 列挙して示します。数値属性の場合は,"@attribute temperature numeric"のように記述します。なお,数 値属性のタイプにはreal, integer, numericを指定す ることができます。データ部はデータ毎に各行に記述さ れた","で区切られた属性値の羅列であり,CSVファ イルのフォーマットと同じとなっています(属性値の順 番はヘッダに記述された属性の順番に従います)。よっ て,CSVファイルの先頭にヘッダを不可するだけで,入 カファイルを簡単に作成することが可能です。

@relation sensor			
@attribute temperature {high, mid, low}			
@attribute V_A{high, mid, low}			
@attribute V_B{high, mid, low}			
@attribute otherSensor {on, off}			
@attribute failure {true, false}			
@data			
high, high, high, off, true			
high, high, high, on, true			
mid, high, high, off, false			
low, mid, high, off, false			
low, low, mid, off, false			
low, low, mid, on, true			
mid, low, mid, on, false			
high, mid, high, off, true			
high, low, mid, off, false			
low, mid, mid, off, false			
high, mid, mid, on, false			
mid, mid, high, on, false			
mid, high, mid, off, false			
low, mid, high, on, true			

図 2 入力用データ sensor. arff の例 (ARFF 形式)

**起動** Weka はインストールディレクトリにある鳥の イメージの Weka3-6-1というアイコンをクリック することで起動できます。起動画面には"Explorer, Experimenter, KnowledgeFlow, Simple-CLI"の4つの ボタンが表示されますが,ここではGUIを使用して対 話的に処理を行うことができる"Explorer"を選択しま す。各処理は最上部のボタンで選択できるようになって おり,それぞれ, Preprocess(前処理), Classify(分 類,決定木学習を含む), Cluster (クラスタリング), Associate (相関分析), Select attribute (属性選択), visualize(可視化)となっています(図3参照)。なお, 起動時には Preprocess のみ選択可能となっています。

**データ読み込み** まずPreprocess(前処理)を選択して, 上から2段目左側の"Open file..."ボタンをクリッ クして arff(ここでは" sensor. arff")ファイルを指 定して読み込みます。

図3に読み込み直後の画面を示します。この画面で はデータの属性や分布についての特徴が表示されま す。右下に表示されたヒストグラム左上のリストで 目標属性(クラス)を選択して、右側の"visualize all"ボタンをクリックすると、各属性値のデータに 占める目標属性値(クラス)毎のデータ分布を、色付 きヒストグラムで一覧することができます。



図3 Weka 前処理画面例

 決定木学習 決定木学習を実施するには、最上段の "Classify"ボタンを選択して図4のような分類画面 に切り替えます。次に2段目の"Classifier"欄から trees、J48 (C4.5に相当)を選択します。

左側の"Test option"では学習の方法として交叉検 定法などが選択できますが、ここではひとまず全デー タを使用して学習を行う"Use training set"を選択 します。詳細なオプションはその下の"More options" で指定することができます。また、その下のリストで、 目的属性を選択します。ここではセンサの誤動作を示 す"failure"を選択します。



図4 決定木学習の実施画面例

次にその下側にある"Start"ボタンをクリックす ると学習が開始されます。学習結果は得られた決定 木に対する評価値とともに右側のウィンドウに表示 されます。得られた決定木の尤もらしさは,TP (true positive rate), FP rate (false positive rate), precision, recall などで表示されます。詳細は省略 しますが,この場合の正答率は100%で性能のよい決定 木が得られたことになります。

なお、作成された決定木を可視化するには、左下の 学習結果のリスト "Result list"に追加された計算結 果のリストを右クリックして "Visualize tree"を選 択すると、図1のような決定木が別のウィンドウに表 示されます。

相関分析 次に,決定木学習と同じデータで相関分析 を実施します。図2のデータは本来の相関分析が扱う データとは少し異なりますが,行をトランザクション, 観測される属性値をアイテムと解釈することによって 相関分析を行うことが可能です。

最上段の"Associate"をクリックして相関分析用の ウィンドウに切り替えます。次の段の"Associator"を クリックして "apriori"を選択します。"apriori……" と表示されたリストの欄をクリックすると、抽出時の support の敷居値などを設定することができます(デ フォルトでは0.2)。"Start"ボタンを押すと計算が実行 され、右ウィンドウに以下のような抽出された相関ルー ルが表示されます。なお、上記のルールの確信度は1、 支持度は14/3 で重要性の高いルールといえます。

- 1. temperature=mid 4 ==> failure=false 4  $\ensuremath{\mathsf{a}}$ 
  - conf:(1)(温度が中程度であれば誤動作なし)
- 2.  $V_A=low 4 \implies V_B=mid 4$  conf:(1)
- (A の電圧が低いとき B の電圧も低い)(以下略)

クラスタリング 最後にクラスタリングを実施します。 データには Weka のインストールディレクトリの下の data ディレクトリに入っている"iris.arff"を利用 します。このデータセットには、アヤメの属性(花弁 の幅,長さなど4種,数値)がアヤメの種類とともに 格納されています。単純なデータではありますが、天 文学,惑星科学,地球観測などの分野で一般的である 各種スペクトルのグループ化も、この実験から容易に イメージすることが可能と思います。

まず" Preprocess" ボタンを押して、" Open file" で data の "iris. arff"を指定します。数値属性の場合, 必要であれば前処理を実施します。属性ごとの正規化 を実施するには、Filter 欄の "Choose" ボタンを押し て、filters, unsupervised, attribute を順に指定し て、" standarize"を選択し、" apply" ボタンを押しま す。これによって、各属性値は属性毎に0を平均値と して標準偏差が1になるように正規化されます。

次に最上部の"Cluster"ボタンをおして、図5の クラスタリング画面に切り替え、"Choose"からア ルゴリズムを選択します。ここでは簡単のため" simpleKmeans"を選択します。選択したアルゴリズム が表示されているリスト欄をクリックするとクラスタ 数 (numClusters, k) などのパラメータを入力するこ [裏へ続く] とができます。今回はデータセットの記述からアヤメの 種類が3個であることが既知なので,numClusters=3と します。また、学習の際にアヤメの種類をマスクして クラスタリング結果と比較するために、左側のCluster mode から"Classes to clusters evaluation"をクリッ クして、リストで class (アヤメの種類)を指定します。 Start ボタンをクリックすると処理が始まり、右画面に その結果が表示されます。最下部に表示された各データ のクラスタ番号とマスクしたアヤメの種類の比較結果か ら、誤分類は 17/150,11%で、教師データを与えなくても、 まずまずのグループ化ができたたことがわかります。

なお、"Result list"を右クリックして、"visualize cluster assignments"を指定すると、属性毎のクラスタ 分布が可視化されます。データ毎のクラスタリング結果 を見るには、この画面の save ボタンをクリックして保存 用のファイル名を指定してエクスポートします。

図6にクラスタリングの結果をクラスタ毎に色を変え たスペクトル形式で表示してみます(エクスポートした



図5 クラスタリングの実行画面例

ADASS とは、Astronomical Data Analysis Software & Systems の略称で、その名前が示すとおり、天文データを扱う ためのソフトウェアに関する国際カンファレンスです。毎年開催 されており、本年度は札幌で行なわれました。日本での開催 は ADASS 史上初であり、多くの日本人開発者・研究者の参 加がみられました。

ADASS は、自然科学系の研究会等と同様のロ頭発表・ポ スター発表のほか、フォーカスデモやフロアデモといった、ソ フトウェアを実際に動かしながら紹介する発表もあり、その有用 性をアピールする絶好の機会となっています.このほか、BoF といって、あるテーマに関心が高い人達が集まる小ミーティン グもあり、毎年、FITS に関するテーマなどの会合があります. 我々、C-SODA グループからは、宇宙科学データ公開サー ビス "DARTS"、次世代 FITS I/O ライブラリ "SFITSIO"をポ スターで、天体ナビゲーションツールである "JUDO"、「すざく」 XIS 用 QL ツール "UDON" をフォーカスデモで発表しました. "IUDO"と "UDON"のフォーカスデモは、昨年度も行なっ データからプロット)。この図からもスペクトルの形状に 応じた自動的なグループ化がなされていることを見て取 ることができます。



### 6. まとめ

今回は代表的な3つの手法,決定木学習,相関分析,ク ラスタリング(k-means法)について,toy dataに対し てWekaを使用した分析例を含めながら説明しました。使 用したデータはまさに"toy-おもちゃ"レベルですが,こ こで取り上げたような分析のニーズは,様々なシチュエー ションで頻出するものであると考えられます。データマイ ニングは知識の発見という大きな目的を目指すものではあ りますが,分析指針が定まらず眠っているデータや,煩雑 な手作業の自動化などにも活用可能ですので,お試しいた だけると幸いです。

次回は、少し進んだモデルや、実際の活用例などについ て紹介する予定です。

参考文献等(1)元田浩ほか,データマイニングの基礎,オー ム社,2006

注1 文献(1)の掲載例(ゴルフプレイデータ)を修正して利用。

# 2009 札幌 ADASS 報告

山内 千里 (JAXA 宇宙科学情報解析研究系 /C-SODA) たのですが,それでも30人程度の方々に動作している様子 を見てもらう事ができました. "JUDO"は現状ではナビゲーショ ンツールとしてしか使われていませんが,今年度の開発によっ てより汎用的な QL 画像切り出しシステムに発展していく予定 です.ご期待ください.

FITS I/O ライブラリ "SFITSIO" については、いずれ PLAIN ニュースでももっと誌面を使って紹介する予定ですが、 CFITSIO よりも圧倒的に簡単に FITS ファイルが扱える C 言語 プログラマ向きのライブラリです. 今回の ADASS ではポスター の他、FITS の BoF でも座長のご好意で 5 分間の発表をさせ てもらう事ができました. チラシとマニュアルのコピーも在庫ゼ ロとなり、予想以上の関心の高さに驚きました.

来年の ADASS はアメリカのボストンです. 自分のソフトウェア をアピールするのにも, いろいろなソフトウェアに関する情報を 集めるのにも, たいへん有意義なカンファレンスです. データ 処理に深く関わっているそこのあなたも来年は参加してみませ んか?

編集発行:宇宙航空研究開発機構 宇宙科学研究本部 科学衛星運用・データ利用センター 〒229-8510 相模原市由野台 3-1-1 Tel.042-759-8767 住所変更等 e-mail:news@plain.isas.jaxa.jp 本ニュースはインターネットでもご覧になれます.http://www.isas.jaxa.jp/docs/PLAINnews ●編集後記:いつのまにか、相模原キャンパスは落ち葉で埋め尽くされています。ふと気づいたら、もう師走が目の前です!(K.E.)