

COM実験(2008/10/9)

櫻井彰人

実験の目的

- データの(客観的、自動的)分析を試みる、その基礎を知る

レポートについて1

- レポートに報告して戴きたいことは次の点です。
- スライド中、実験1及び実験2の「実験手順」となっている部分に書かれている手順に従い、実験を行い、その結果(自分が行ったこと、得られた結果)とそれに関する考察
 - 実験1「文字認識」
 - 実験2「歌詞の分類」
 - 実験3「ドル円レートの予測」は、ご参考
- 感想

レポートについて2

- 言うまでもないことですが、他人のレポート・著作物等を写してはいけません。自分の独力で作成してください。
- 他人の著作物からの引用には、出典を明記してください。
- 締め切りは、10月29日です。

では、本論に

データには裏がある

- 申し訳ない、ちょっと表現が悪い。
- 正しくは、「データには構造がある」と言うべき。
- データというのは、数字の並び、または、文字の並び。一列に並んだもの。だから構造はない。
- そして、紙の上にな書かれたり、コンピュータのメモリに入れられているだけ
- しかし、それが何か実世界のもの(自然物、人工物)によって生成されたものなら、それらの間に何らかの関係があろう。

データの裏を知りたい、知ろう

- となれば、データの背後関係を知りたくなろう。
- 単なる興味の問題ではなく、世の中を生きていく上で(少々大げさか)必要なこと
- どうしてそのデータは出てきたのか？
- データにはどういう関係があるのか？
- (そのデータに関して)将来どうなるのか？

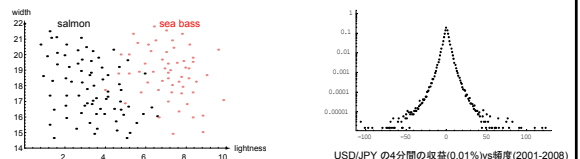
例えば

- 最近、「app store」という言葉を耳にしたのだが、これは今後流行るのだろうか？
- いやいや、これほど難しい問いに答えられるわけではない。が、例えば、
- 最近、体重の変化がこんな(日々の体重がずら〜っと並んだ表)具合だが、おかしくないか？病院に行くべきか？
- なんて気が利かないんだ、このブラウザは？私の過去の行動を見ているんだから、次やりたいことぐらい、察してくれよ。
- 写真の数が増えすぎて困った。見たい写真をいくつか教えるから、それに似た写真を集めて、時間順に並べてくれると助かるんだが。

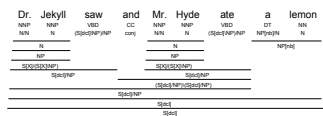
ここで、データとは

- 数字の列
- 文字の列(言葉というべき)。言語
- 写真の集まり
- 絵画の集まり
- 行動の表現(って何だろう?)の集まり
- 音の表現(って何だろう?)の集まり

構造とは



R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern Classification*



J. Curran and S. Clark. C&C tools.

その構造1

- 数値データ(最も無味乾燥)
 - 正常と考えられる個体群、異常と考えられる個体群からデータが得られる。
 - 時間とともに変化する(時系列)。その変化の原因。
- 文字データ(少し構造があるように見える)
 - どのような文字を使うかという癖(人により違う)
 - どのような単語を使うかという癖(人により違う)
- 画像(写真、絵画等)
 - 対象(意図、目的に依存)
 - 構図、描き方、撮り方(作成者、状況等に依存)



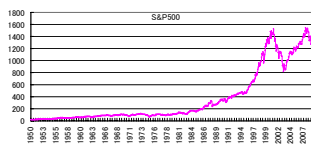
http://cert.yahoo.co.jp/text/digicame/chap2/c2_0302.html

その構造2

- 音(音楽、声、虫の音、、、)
 - 道具(楽器、性別・年齢・健康、虫の種類)
 - 音楽の場合に、楽曲の種類、演奏等
- 行動
 - (買物) どのような目的。誰のために。何のために。どこで。
 - (web閲覧) 目的。何(ネットワーク、PC等)を使って。

構造が分かれば

- 次の行動がとれる
 - 健康データ: 病気だと分かれば、医者に行く
 - 時系列: 次の予測
 - 文学作品: 気に入った作者が分かれば、次の作品を読む
 - web の流行: それに乗る(乗らないという選択もあり)
- 一般に、分類し、それに合わせた行動をとる。



どうしたらできるか？

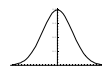
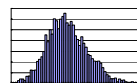
- コンピュータサイエンスでは、昔から、「機械学習」という分野で深く研究されてきた
 - 人工知能の一分野と考えてよい
 - なぜ「学習」か？
 - 人間の学習は、丸暗記を除くと、理解して、記憶して、次に使えるようにすること
 - この「理解する」ということがポイント
 - データの構造を知ることが、理解することの重要な一部
 - なぜ「機械」か？
 - コンピュータのこと。英語では、計算機械と言った。
 - なぜ「ロボットでない」のか？
 - ロボットにも機械学習は必要。しかし、ロボットに限らない。
 - 人間との違いは？
 - 人間ほど賢くはない。コンピュータは「実世界」を知らない。
 - しかし、大量のデータを、文句も言わず、処理してくれる

どうしたらできる？ その2

- 非常にたくさんのアルゴリズムが開発されてきた
- たくさんありすぎて、説明できない(私も全部知っているわけでは、当然、ない)。
 - というわけで、説明しません。
- 実験では、道具を使って、体感することで我慢してください。
 - 申し訳なし
- もっと知りたい場合には、是非、私の講義を聴いてください。

統計との違い

- 昔(20年前)はまったく別のdiscipline
- しかし、今は融合が進んでいる。
- しかし、違いはある
- 統計:
 - 機械学習に比べると簡単な(確率的)構造を考える。そのかわり、検定することが頭にある。主に数値データ
- 機械学習:
 - どんなモデルでも考える。検定したいが、できないことが結構ある。記号データも平気で扱う。



データマイニングとの違い

- 実は、ない。
- 研究者もかなり重なっている
- あえて言えば:
 - 機械学習: アルゴリズムに重点
 - 精度、速さ、軽さ、応用の広さ、表現力
 - データマイニング: 大量のデータに重点
 - 知識、アクション、データ個別的でもOK

話しが抽象的になっているので、

- 簡単な実例を一つ二つ。

データの例

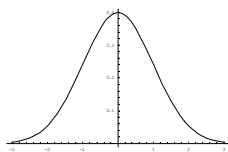
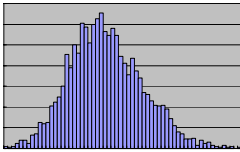
Make	Model	Year	Head inj. c.	Chest decel.	L. Leg R. Leg	D/P	Protection	アクコー ム数	血圧(mmHg)	(LOW) (HIGH)				
Acura	Integra	87	599	35	791	262	Driver	manual belts	3	金	30	少女	87	130
Acura	Integra RS	90	585	1545	1301	Driver	Motorized belts	5	日	20	少女	80	130	
Acura	Legend LS	88	439	80	929	708	Driver	fl airbag	6	月	30	高学	90	130
Audi	90	89	600	49	168	1971	Driver	manual belts	7	火	18	少女	87	130
Audi	100	89	185	35	988	894	Driver	fl airbag	8	水	18	少女	104	140
BMW	325i	90	1038	50	868	Driver	fl airbag	9	木	20	少女	83	130	
Buick	Century	91	810	47	1360	315	Driver	passive belts	10	金	20	高学	94	130
Buick	Black Park	88	1467	8										
Buick	Le Sabre	90		327/11/2000	53.6875	54.5156	51.0312	51.25	40198100	51.250	200049			
Buick	Regal	88	880	28/11/2000	51.8375	53.1875	50.825	51	52637000	51.000	200049			
Buick	Regal	88	880	29/11/2000	51.8125	53	50.3125	51.6875	55318000	51.688	200049			
Buick	Regal	88	880	30/11/2000	50.1875	50.8375	45.1875	47.875	108469000	47.875	200049			
Buick	Regal	88	880	31/12/2000	49.1875	51.625	47.25	45.5	70468000	48.500	200049			
Buick	Regal	88	880	01/12/2000	49.0625	49.5625	45	45.9125	9501200	45.913	200050			
Buick	Regal	88	880	02/12/2000	47.75	52.125	47.3125	52.125	90848900	52.125	200050			
Buick	Regal	88	880	06/12/2000	52	53.5625	51.2656	51.4375	71419200	51.438	200050			
Buick	Regal	88	880	07/12/2000	50.3125	51	49	49.9375	46444400	49.938	200050			
Buick	Regal	88	880	08/12/2000	51.8375	53.25	51	52.375	95400200	52.375	200050			
Buick	Regal	88	880	03/12/2000	52.875	55.75	52.625	54.8125	78621500	54.813	200051			
Buick	Regal	88	880	12/12/2000	54.75	55.125	53.3125	54.375	39485300	54.375	200051			
Buick	Regal	88	880	13/12/2000	55.1875	55.75	50.3125	51.125	54333600	51.125	200051			
Buick	Regal	88	880	14/12/2000	51.0625	52.5625	50.875	50.9375	48244400	50.938	200051			
Buick	Regal	88	880	15/12/2000	50.0625	50.1875	47.125	48.1719	100237800	48.172	200051			
Buick	Regal	88	880	18/12/2000	49	50.125	42.3125	42.9375	126032400	42.938	200052			
Buick	Regal	88	880	19/12/2000	45	46	41.5	41.75	99018900	41.750	200052			

統計的な表現方法

- 一つ一つのデータには規則性はないのに、集団でみると(出現頻度等に)規則性が見られる。
 - 一様分布
 - 公平なサイコロ、コイン投げ
 - 正規分布
 - たくさんの独立な要因がからむとき
 - Zipfの法則、80/20の法則、1/f ゆらぎ

正規分布

- 多数の独立な要因の結果発生する
 - コイン100回投げの表の回数の1000回実験
 - 例: (本当は違うのだが) 試験の成績

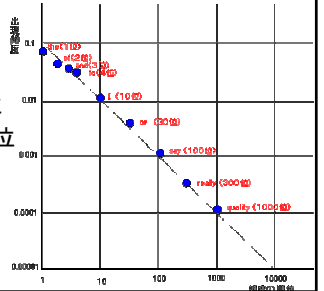


<http://www.cstep.bc.edu/TIMSSI/databas.html> 49 (calculus)

Plot[(1/Sqrt[2 Pi]) Exp[-x^2/2]], {x, -3, 3}

Zipfの法則

- 「単語の出現頻度がその順位に反比例する」という経験則
 - 都市の人口とその順位
 - HPのヒット数とその順位
 - 資産とその順位



機械学習での規則の表現

- 条件文
 - IF むにゃむにゃ THEN あれこれ
 - 信頼性は?? %
- 決定木
 - 次のスライド
- ニューラルネットワーク
- ほかにいろいろ

If-then 規則

- If 涙産生 = 少 then 推薦 = しない
- If 年齢 = 若い and 乱視 = なし and 涙産生 = 通常 then 推薦 = ソフト
- If 年齢 = 老眼以前 and 乱視 = なし and 涙産生 = 通常 then 推薦 = ソフト
- If 年齢 = 老眼 and めがね = 近視 and 乱視 = なし then 推薦 = なし
- If めがね = 遠視 and 乱視 = なし and 涙産生 = 通常 then 推薦 = ソフト
- If めがね = 近視 and 乱視 = あり and 涙産生 = 通常 then 推薦 = ハード
- If 年齢 = 若い and 乱視 = あり and 涙産生 = 通常 then 推薦 = ハード
- If 年齢 = 老眼以前 and めがね = 遠視 and 乱視 = あり then 推薦 = しない
- If 年齢 = 老眼 and めがね = 遠視 and 乱視 = あり then 推薦 = しない

年齢	めがね	乱視	涙産生	推薦
若い	近視	なし	少	しない
若い	近視	なし	通常	ソフト
若い	近視	あり	通常	ハード
若い	遠視	なし	少	しない
若い	遠視	なし	通常	ソフト
若い	遠視	あり	通常	ハード
老眼以前	遠視	なし	少	しない
老眼以前	遠視	なし	通常	ソフト
老眼以前	遠視	あり	通常	ハード
老眼以前	近視	なし	少	しない
老眼以前	近視	なし	通常	ソフト
老眼以前	近視	あり	通常	ハード
老眼	遠視	なし	少	しない
老眼	遠視	あり	通常	ハード

実験1:文字認識

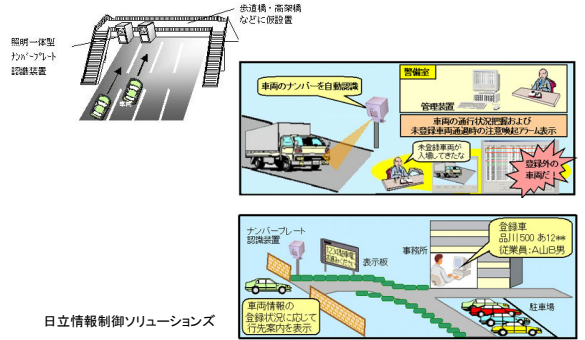
・実世界で活躍中

- 高速道路の料金所でナンバープレートを読む
 - ・Nシステムもそうらしい
- Wikipediaより



- 郵便番号自動読み取り(「区分け」と「配達順序」が重要)。手書きの住所も読む。高速。
- 蛇足ですが、「正解率99% ネット認証技術、書籍のデジタル化に威力」という記事を見て下さい。

蛇足



文字認識は結構難しい

- ・人間なら、崩し字でなければ、簡単だと思う
- ・けれども、平仮名の読み方を、日本語を勉強したことの無い人に教えることを想像してみてください。
- ・何が難しいか？
- ・規則が書けない！(実は分からないのかも、実は「ない」のかも)

実験データの説明

- ・簡単な文字認識
 - 数字のみ。
- ・データの前処理(これが大変)済み
 - 分離(他の文字から分離)済み
 - 整形(大きさ、傾き、重心等)済み
- ・それでも、結構、難しそう。
 - 「数字」を知らない(!)人に区別の仕方を説明してみよう
- ・本質「分類規則が表現できない」
 - データから得る
- ・データのもと:
 - UCI Machine Learning Repository
 - Optical Recognition of Handwritten Digits Data Set

データの前処理

This is a diagonal image that will permit you the corner pattern of 012

This is a diagonal image that will permit you the corner pattern of 012

The non-parallel hand line is a very unstable and pattern making difficult the slow angle extraction.

The hand and data writing is also not as well as the standard and connected character.

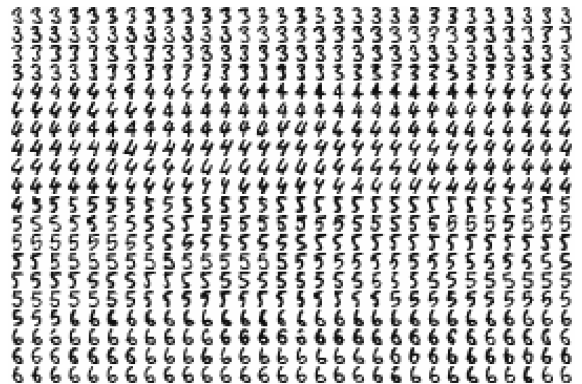
The non-parallel hand line is a very unstable and pattern making difficult the slow angle extraction.

The hand and data writing is also not as well as the standard and connected character.



Giorgos Vamvakas

データを図にして並べてみました



UCI Machine Learning Repository

実験3 ドル円レートの予測

- FX: 外国為替証拠金取引
 - 証拠金(保証金)を業者に預託し、主に差金決済による通貨の売買を行なう取引
- FXで利益を上げることができるのだろうか?
 - 仲介業者の取り分(スプレッド)は小さい。宝くじとは違う
 - しかし、典型的なギャンブル。ゼロサムゲーム。勝者がいれば、敗者がいる。そして、敗者が圧倒的に多い(80-20の法則、幕法則)と思う。
- 値動きは、基本的には、ランダムウォークのはず
 - すなわち、予測不能のはず。
- しかし、少し、試してみよう

実験データの説明

- 米ドル(USD)を日本円(JPY)で売買する
- 価格(?)の単位は 0.01 円
- 分単位の値動きを見る
 - 分足という



東京金融取引所

実験データの説明

- USD/JPYの分足を用いる
 - まず、(いつでもよいのだが)2008年9月1日にしよう
 - Forexiteというサイトのものを用いる。時刻は GMT+1 (Central European Time)。データの正確性の保証はない。
 - 24時間中の、一分毎、Open (始値)、High (高値)、Low (安値)、Close (終値)が時系列に記されている。
 - ある「分」の「終値 - 始値」(収益)を予測する
 - 難しい: 何を、予測の根拠に用いるか? i.e. どんな特徴量を用いるか
- (まずは)5分前からの各「分」の収益を用いよう

010908.zip

010908.zip 中の
010908.txt

```
1 <TICKER>,<DTYYYYMMDD>,<TIME>,<OPEN>,<HIGH>,<LOW>,<CLOSE>
2 EURUSD,20080901,000000,1.4698,1.4698,1.4697,1.4698
3 EURUSD,20080901,000100,1.4699,1.4700,1.4698,1.4699
...
4320 USDCHF,20080901,235800,1.1017,1.1017,1.1017,1.1017
4321 USDCHF,20080901,235900,1.1017,1.1024,1.1017,1.1021
4322 USDJPY,20080901,000000,108.31,108.31,108.31,108.31
4323 USDJPY,20080901,000100,108.30,108.31,108.29,108.29
4324 USDJPY,20080901,000200,108.28,108.30,108.28,108.29
4325 USDJPY,20080901,000300,108.30,108.30,108.30,108.30
```

ここから

作業手順

- 2008年9月1日のデータを得る。
- Excelファイルとし、各分について5分前からの収益を計算する。当「分」の収益も求める。
- 予測問題を簡単にするために、up or down を示す値をつける(上昇していたら +1, 下降していたら -1)
 - このファイルは用意してある。

010908.zip 中の
010908.txt から
USDJPY を全部

USDJPY080901.xls

	B	C	D	E	F	G	H	I	J	K	L	M	N
1	<TIME>	<OPEN>	<HIGH>	<LOW>	<CLOSE>	5 before	4 before	3 before	2 before	1 before	predict	updown	
2	0	108.31	108.31	108.31	108.31								
3	100	108.3	108.31	108.29	108.29								
4	200	108.28	108.3	108.28	108.29								
5	300	108.3	108.3	108.3	108.3								
6	400	108.29	108.31	108.29	108.29								
7	500	108.3	108.3	108.25	108.25	0	-0.01	0.01	0	0	-0.05	-1	
8	600	108.25	108.25	108.25	108.25	-0.01	0.01	0	0	-0.05	0	0	
9	700	108.25	108.26	108.24	108.24	0.01	0	0	-0.05	0	-0.01	-1	

実験手順: ファイルの準備

- 収益の部分のみを取り出し、csv ファイルを作る。
- arff のヘッダーをつけ、Weka用の arff ファイルとする。メモ帳なりエディタなりを使う方が間違いが少ない

```
1 @relation USDJPYReturns080901;
2 @attribute 5m numeric;
3 @attribute 4m numeric;
4 @attribute 3m numeric;
5 @attribute 2m numeric;
6 @attribute 1m numeric;
7 @attribute this numeric;
8 @attribute updown {-1,0,1};
9 @data;
10 0,-0.01,0.01,0.0,-0.05,-1;
11 -0.01,0.01,0.0,-0.05,0,0;
```

実験手順: 味見

- Weka を使う
 - 使うときに、当「分」の収益という属性は "Preprocess" で remove して下さい(次スライド)
 - trees にある J48: 決定木
 - functions にある neural network
 - functions にある SMO: support vector machine の一つ
 - Bayes にある naïve Bayes
- どれも正解率は 1/3 ぐらいであろう。
 - 気軽に、なんだあー様ランダムなんだ、と思わないで下さい。データを見てみると、収益が正、零、負になる割合がほぼ 1/3 なのです。信じられますか?

Weka手順: 蛇足と補足

選択して
クリックする

信じられないほど左右対称。綺麗。

実験手順: 他の日との比較

- 別の日のデータではどうか？
- 9月2日、9月3日、9月4日で試してみよう
 - xls, csv, arffファイルは自分で作りましょう
- どうですか？
- 予測するための情報が不足でしょうか？
- きっとそうでしょう。では、一分前の「分」の高値安値を含めてみましょう。
 - 9月1日のファイルが作ってあります。他の日ファイルも作り、試しましょう。
 - でもうまいけません。

020908.zip
030908.zip
040908.zip

USDJPY080901A.xls

実験手順

- では、5分足(5分間一区切り)を試してみよう。
 - 1分足では、他の人の動きをみて動く(相関が発生する)ということが少ないので、動きはランダムになり、従って、予測できない。
 - しかし、5分程度みれば、データ間に相関が発生し、従って、予測可能となる可能性がある。
 - 本当か？

実際には逆に、収益の時間相関は20分くらいまでは存在すると報告されている。例えば
P. Gopikrishnan, et al. Scaling of the distribution of fluctuations of financial market indices,
Physical Review E vol. 60, 5305 - 5316 (1999)

実験手順

- まず、5分ごとのデータにする前に、「5分間」の始値、高値、安値、終値を求める
- 次に、5分区切りのデータを抽出する。

コピー ← 最小値

A	B	C	D	E	F	G	H	I	J	K	L
<YMDDD>	<TIME>	<OPEN>	<HIGH>	<LOW>	<CLOSE>	<TIMESOC>	<OPEN>	<HIGH+4m>	<LOW+4m>	<CLOSE+4m>	
20080902	0	10818	10818	10814	10814	0	10818	10818	10807	10808	
20080902	100	10812	10813	10812	10812	100	10814	10814	10807	10808	
20080902	200	10813	10814	10812	10812	200	10814	10814	10807	10807	
20080902	300	10811	10811	10809	1081	300	10811	10811	10807	10808	
20080902	400	10809	10809	10807	10808	400	10809	10809	10807	10808	
20080902	500	10808	10808	10806	10806	0	10808	10808	10807	10808	

<TIME>を500で割った余り
あとで、この値の値をソートすれば、
5分区切りのデータが得られる

コピー

実験手順

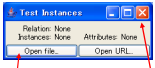
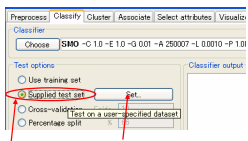
- Weka を用いてみる
 - J48, SMO, naive Bayes など
- 少し正解率が高い
 - しかし、データ分布を見てみると、収益=0 の事例が少ないことが分かる。
 - やはり予測できないのか。
 - しかし、データ数が少ないからかもしれない。
- 9月1日~4日で試してみよう
 - 少しよいか？
 - しかし、別の日、例えば、8月6日~8日でテストをしてみたらどうだろうか？

060808.zip
070808.zip
080808.zip

実験手順

- 別のデータでテストする方法
- 学習データと同じ属性数(並び)のテストデータを用意する。
- 今回は、当「分」の収益という属性を持ったファイルしかないで、そのままではテスト用のデータファイルにはならない。そこで、当「分」の収益という属性を削除したデータファイルを作る。Weka ができる。

Weka補足: テストデータの指定



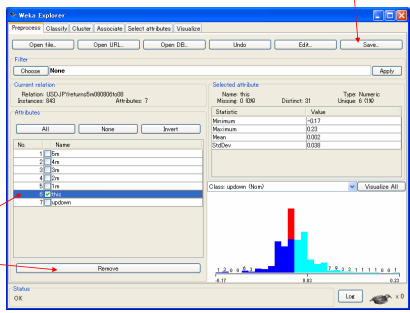
次に、ここをクリックし、

テストデータの
入ったファイルを
指定し

ここをクリックして閉じる

Weka補足: ファイルを作る

② 保存する



① 選択して
クリックする