Mining personal banking data to detect fraud

David J. Hand Imperial College London

September 2007

1

Imperial College London Workshop on Data Analysis and Classification In honour of Edwin Diday My research group:

Niall Adams, Adam Brentnall, Martin Crowder, Nick Heard, Dave Weston, Chris Whitrow, Piotr Juszczak, Kiriaki Platanioti, Dimitris Tasoulis, Nicos Pavlidis, Matt Turnbull, James Bentham, Iding Wu, Fanyin Zhou, Christoforos Anagnostopoulos, Daniel Balabanoff, Ed Tricker, Gordon Blunt, ...

Three parts:

- I: Introduction
- II: How big is fraud?
- III: Fraud in banking

I: Introduction

What is fraud?

Criminal deception; the use of false representations to gain an unjust advantage

Concise Oxford Dictionary

Older than humanity itself.

- even animals are known to try to deceive others
- camouflage

The economic imperative

1) Not worth spending \$200m to stop \$20m fraud

e.g. Letter from London Times, August 13, 2007 "Sir, I was recently the victim of an internet fraud. The sum involved was several hundred pounds. My local police refused to investigate, stating that their policy was to investigate only for sums over £5000."

2) The Pareto principle

the first 50% of fraud is easy to stop; next 25% takes the same effort; next 12.5% takes the same effort; ...

3) Resources available for fraud detection are always limited

- in the UK around 3% of police resources go on fraud
- this will not significantly increase

II: How big is fraud?

e.g. In the USA

"Participants in our study estimate U.S. organizations lose 5% of their annual revenues to fraud.

Applied to the estimated 2006 United States Gross Domestic Product, this 5% figure would translate to approximately \$652 billion in fraud losses."

Association of Certified Fraud Examiners

Imperial College London Workshop on Data Analysis and Classification6In honour of Edwin Diday

Cost of fraud

- = immediate direct loss due to fraud
- + cost of fraud prevention and detection
- + cost of lost business (when replacing card)
- + opportunity cost of fraud prevention/detection
- + deterrent effect on spread of e-commerce

7

Does this matter to you?

Identity theft

Fraudsters uses your name and identifying information to

- obtain credit cards
- phone and telecoms
- bank loans
- mortgages
- rent appartments
- if stopped for speeding, or charged with crime, etc.

leaving you with the debts and problems

Identity theft in the USA:

10 million victims in 2003 Average individual loss \approx \$5,000 Total loss to individuals and businesses in 2003 \approx \$50 bn (Federal Trade Commission survey)

+ time to sort out

 \Rightarrow Americans spent nearly 300 million hours resolving ID theft issues in 2003

Typically takes up to two years to sort out the problems, reinstate credit rating, reputation, etc, after detection

III: Fraud in banking

Banking fraud has many aspects

My main focus here is *retail* or *consumer* banking fraud

- personal banking
- credit cards
- home mortgages
- car finance
- personal loans
- current accounts
- savings accounts

Nature of plastic card fraud data

- many transactions billions algorithms must be efficient
- mixed variable types (generally not text, image)
- large number of variables
- incomprehensible variables, irrelevant variables
- different misclassification costs
- many ways of committing fraud
- unbalanced class sizes (c. 0.1% transactions fraudulent)
- delay in labelling
- mislabelled classes
- random transaction arrival times
- (reactive) population drift

Credit card data:

Transaction ID Transaction type Date and time of transaction (to nearest second) Amount Currency Local currency amount Merchant category Card issuer ID ATM ID POS type Cheque account prefix Savings account prefix

Acquiring institution ID Transaction authorisation code Online authorisation performed New card Transaction exceeds floor limit Number of times chip has been accessed Merchant city name Chip terminal capability Chip card verification result

.

Workshop on Data Analysis and Classification12In honour of Edwin Diday

A commercial example of fraud data

US Patent 5,819,226 (see USPTO website) on *Fraud detection and modeling*, (HNC Software in 1992) lists the following variables:

Customer usage pattern profiles representing time-of-day and day-of-week profiles; Expiration date for the credit card; Dollar amount spent in each SIC (Standard Industrial Classification) merchant group category during the current day; Percentage of dollars spent by a customer in each SIC merchant group category during the current day; Number of transactions in each SIC merchant group category during the current day; Percentage of number of transactions in each SIC merchant group category during the current day; Categorization of SIC merchant group categories by fraud rate (high, medium, or low risk); Categorization of SIC merchant group categories by customer types (groups of customers that most frequently use certain SIC categories); Categorization of geographic regions by fraud rate (high, medium, or low risk); Categorization of geographic regions by customer types; Mean number of days between transactions: Variance of number of days between transactions; Mean time between transactions in one day; Variance of time between transactions in one day; Number of multiple transaction declines at same merchant; Number of out-of-state transactions: Mean number of transaction declines: Year-to-date high balance; Transaction amount; Transaction date and time: Transaction type.

Imperial College London

Workshop on Data Analysis and Classification13In honour of Edwin Diday

"Additional fraud-related variables which may also be considered are listed below"

Imperial College London

Workshop on Data Analysis and Classification14In honour of Edwin Diday

Current Day Cardholder Fraud Related Variables bweekend current day boolean indicating current datetime considered weekend cavapvdl current day mean dollar amount for an approval cavapvdl current day mean dollar amount for an approval cavaudl current day mean dollars per auth across day ccoscdoni current day cosine of the day of month i.e. cos(day ((datepart(cst.sub.-- dt) * &TWOPI)/30)); ccoscdow current day cosine of the day of week i.e. cos(weekday ((datepart(cst.sub.-- dt) * &TWOPI)/7)); ccoscmoy current day cosine of the month of year i.e. cos(month ((datepart(cst.sub.-- dt) * &TWOPI)/12)); cdom current day of month cdow current day day of week chdzip current cardholder zip chibal current day high balance chidcapy current day highest dollar amt on a single cash approve chidcdec current day highest dollar amt on a single cash decline chidmapy current day highest dollar amt on a single merch approve chidmdec current day highest dollar amt on a single merch decline chidsapy current day highest dollar amount on a single approve chidsau current day highest dollar amount on a single auth chidsdec current day highest dollar amount on a single decline cmoy current day month of year cratdcau current day ratio of declines to auths csincdom current day sine of the day of month i.e. sin(day ((datepart(cst.sub.-- dt) * &TWOPI)/30)); csincdow current day sine of the day of week i.e. sin(weekday ((datepart(cs.sub.-- dt) * &TWOPI)/7)), csincmoy current day sine of the month of year i.e. sin(month ((datepart(cs.sub.-- dt) * &TWOPI)/12)); cst.sub.-- dt current day cst datetime derived from zip code and CST auth time ctdapy current day total dollars of approvals ctdau current day total dollars of auths ctdcsapy current day total dollars of cash advance approvals ctdcsdec current day total dollars of cash advance declines ctddec current day total dollars of declines ctdmrapy current day total dollars of merchandise approvals ctdmrdec current day total dollars of merchandise declines ctnapy current day total number of approves ctnau current day total number of auths ctnau10d current day number of auths in day <= \$10 ctnaudy current day total number of auths in a day ctncsapy current day total number of cash advance approvals ctncsapy current day total number of cash approves ctncsdec current day total number of declines cmmrapy current day total number of merchandise approvals ctnmrdec current day total number of merchandise declines ctnsdapy current day total number of approvals on the same day of week as current day ctnwdaft current day total number of weekday afternoon approvals ctnwdapv current day total number of weekday approvals ctnwdeve current day total number of weekday evening approvals ctnwdmor current day total number of weekday morning approvals ctnwdnit current day total number of weekday night approvals ctnweaft current day total number of weekday morning approvals ctnweapy current day total number of weekend approvals ctnweeve current day total number of weekend evening approvals ctnwemor current day total number of weekend morning approvals ctnwenit current day total number of weekend night approvals current day current balance cyraud1 current day variance of dollars per auth across day czrate1 current day zip risk group 1 Zip very high fraud rate czrate2 current day zip risk group 2 'Zip high fraud rate' czrate3 current day zip risk group 3 'Zip medium high fraud rate' czrate4 current day zip risk group 4 'Zip medium fraud rate' czrate5 current day zip risk group 5 'Zip medium low fraud rate' czrate8 current day zip risk group 6 'Zip low fraud rate' czrate7 current day zip risk group 7 'Zip very low fraud rate' czrate8 current day zip risk group 8 'Zip unknown fraud rate` ctdsfa01 current day total dollars of transactions in SIC factor group 01 ctdsfa02 current day total dollars of transactions in SIC factor group 02 ctdsfa03 current day total dollars of transactions in SIC factor group 03 ctdsfa04 current day total dollars of transactions in SIC factor group 04 ctdsfa05 current day total dollars of transactions in SIC factor group 05 ctdsfa06 current day total dollars of transactions in SIC factor group 06 ctdsfa07 current day total dollars of transactions in SIC factor group 07 ctdsfa08 current day total dollars of transactions in SIC factor group 08 ctdsfa09 current day total dollars of transactions in SIC factor group 09 ctdsfa10 current day total dollars of transactions in SIC factor group 10 ctdsfa11 current day total dollars of transactions in SIC factor group 11 ctdsra01 current day total dollars of transactions in SIC fraud rate group 01 ctdsra02 current day total dollars of transactions in SIC fraud rate group 02 ctdsra03 current day total dollars of transactions in SIC fraud rate group 04 ctdsra03 current day total dollars of transactions in SIC fraud rate group 04 ctdsra03 current day total dollars of transactions in SIC fraud rate group 04 ctdsra03 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra04 current day total dollars 03 ctdsra04 current day total dollars of transactions in SIC fraud rate group 04 ctdsra05 current day total dollars of transactions in SIC fraud rate group 05 ctdsra06 current day total dollars of transactions in SIC fraud rate group 06 ctdsra07 current day total dollars of transactions in SIC fraud rate group 07 ctdsva01 current day total dollars in SIC VISA group 01 ctdsva02 current day total dollars in SIC VISA group 02 ctdsva03 current day total dollars in SIC VISA group 03 ctdsva04 current day total dollars in SIC VISA group 04 ctdsva05 current day total dollars in SIC VISA group 05 ctdsva06 current day total dollars in SIC VISA group 06 ctdsva07 current day total dollars in SIC VISA group 07 ctdsva08 current day total dollars in SIC VISA group 08 ctdsva09 current day total dollars in SIC VISA group 09 ctdsva10 current day total dollars in SIC VISA group 10 ctdsva11 current day total dollars in SIC VISA group 11 ctnsfa01 current day total number of transactions in SIC factor group 01 ctnsfa02 current day total number of transactions in SIC factor group 02 ctnsfa03 current day total number of transactions in SIC factor group 03 ctnsfa04 current day total number of transactions in SIC factor group 04 ctnsfa05 current day total number of transactions in SIC factor group 05 ctnsfa06 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 05 ctnsfa06 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 05 ctnsfa06 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 current day total number of transactions in SIC factor group 06 ctnsfa07 cu 07 ctnsfa08 current day total number of transactions in SIC factor group 08 ctnsfa09 current day total number of transactions in SIC factor group 09 ctnsfa10 current day total number of transactions in SIC factor group 10 ctnsfa11 current day total number of transactions in SIC factor group 11 ctnsra01 current day total number of transactions in SIC fraud rate group 01 ctnsra02 current day total number of transactions in SIC fraud rate group 02 ctnsra03 current day total number of transactions in SIC fraud rate group 03 ctnsra04 current day total number of transactions in SIC fraud rate group 04 ctnsra05 current day total number of transactions in SIC fraud rate group 05 ctnsra06 current day total number of transactions in SIC fraud rate group 06 ctnsra07 current day total number of transactions in SIC fraud rate group 07 ctnsva01 current day total number in SIC VISA group 01 ctnsva02 current day total number in SIC VISA group 02 ctnsva03 current day total number in SIC VISA group 03 ctnsva04 current day total number in SIC VISA group 04 ctnsva05 current day total number in SIC VISA group 05 ctnsva06 current day total number in SIC VISA group 06 ctnsva07 current day total number in SIC VISA group 07 ctnsva08 current day total number in SIC VISA group 08 ctnsva09 current day total number in SIC VISA group 09 ctnsva10 current day total number in SIC VISA group 10 ctnsva11 current day total number in SIC VISA group 11 7 Day Cardholer Fraud Related Variables raudymdy 7 day ratio of auth days over number of days in the window ravapvdl 7 day mean dollar amount for an approval ravaudl 7 day mean dollars per auth across window rddapv 7 day mean dollars per day of approvals rddapv2 7 day mean dollars per day of approvals on days with auths rddau 7 day mean dollars per day of auths on days with auths rddauall 7 day mean dollars per day of auths on all days in window rddcsapv 7 day mean dollars per day of cash approvals rddcsdec 7 day mean dollars per day of cash declines rdddec 7 day mean dollars per day of declines rdddec2 7 day mean dollars per day of declines on days with auths rddmrapy 7 day mean dollars per day of merchandise approvals rddmrdec 7 day mean dollars per day of merchandise declines rdnapu 7 day mean number per day of approvals rdnau 7 day mean number per day of auths on days with auths rdnauall 7 day mean number per day of auths on all days in window rdncsapv 7 day mean number per day of cash approvals rdncsdec 7 day mean number per day of cash declines rdndec 7 day mean number per day of declines rdnmrapv 7 day mean number per day of merchandise approvals rdnmrdec 7 day mean number per day of merchandise declines rdnsdap2 7 day mean number per day of approvals on same day of week calculated only for those days which had approvals rdnsdapy 7 day mean number per day of approvals on same day of week as current day rdnwdaft 7 day mean number per day of weekday afternoon approvals rdnwdapy 7 day mean number

Imperial College London

Workshop on Data Analysis and Classification15In honour of Edwin Diday

per day of weekday approvals rdnwdeve 7 day mean number per day of weekday evening approvals rdnwdmor 7 day mean number per day of weekday morning approvals rdnwdnit 7 day mean number per day of weekday night approvals rdnweaft 7 day mean number per day of weekend afternoon approvals rdnweapv 7 day mean number per day of weekend approvals rdnweater 7 day mean number per day of weekend evening approvals rdnwemor 7 day mean number per day of weekend morning approvals rdnwenit 7 day mean number per day of weekend night approvals rdnwemit 7 day mean number per day of weekend morning approvals rdnwenit 7 day mean number per day of weekend morning approvals rdnwenit 7 day mean number per day of weekend morning approvals rdnwenit 7 day mean number per day of weekend morning approvals rdnwenit 7 day mean number per day of weekend morning approvals rdnwenit 7 day mean number per day of weekend morning approvals rdnwenit 7 day mean number per day of weekend morning approvals rdnwenit 7 day mean number per day of weekend morning approvals rdnwenit 7 day mean number per day of weekend rhidcapy 7 day highest dollar amt on a single cash approve rhidcdec 7 day highest dollar amt on a single cash decline rhidmapy 7 day highest dollar amt on a single merch approve rhidmdec 7 day highest dollar amt on a single merch decline rhidsapv 7 day highest dollar amount on a single approve rhidsam 7 day highest dollar amount on a single auth rhidsdec 7 day highest dollar amount on a single decline rhidtapy 7 day highest total dollar amount for an approve in a single day rhidtau 7 day highest total dollar amount for any auth in a single day rhidtdec 7 day highest total dollar amount for a single day rhinapy 7 day highest number of approves in a single day rhinau 7 day highest number of auths in a single day rhindec 7 day highest number of declines in a single day rnaudy 7 day number of days in window with any auths rnausd 7 day number of same day of week with any auths rnauwd 7 day number of weekdays days in window with any auths rnauwe 7 day number of weekend days in window with any auths rncsaudy 7 day number of days in window with cash auths rnmraudy 7 day number of days in window with merchant auths rtdapy 7 day total dollars of approvals rtdau 7 day total dollars of auths rtdcsapv 7 day total dollars of cash advance approvals rtdcsdec 7 day total dollars of cash advance declines rtddec 7 day total dollars of declines rtdmrapv 7 day total dollars of merchandise approvals rtdmrdec 7 day total dollars of merchandise declines rtnapv 7 day total number of approvals rtnapvdy 7 day total number of approvals in a day rtnan 7 day total number of auths rtnau10d 7 day number of auths in window <= \$10 rtncsapy 7 day total number of cash advance approvals rtncsdec 7 day total number of cash advance adeclines rtndec 7 day total number of merchandise approvals rtnmrdec 7 day total number of merchandise declines rtnsdapv 7 day total number of approvals on the same day of week as current day rtnwdaft 7 day total number of weekday afternoon approvals rtnwdapv 7 day total number of weekday approvals rtnwdeve 7 day total number of weekday evening approvals rtnwdmor 7 day total number of weekday morning approvals rtnwdnit 7 day total number of weekend approvals rtnweaft 7 day total number of weekend afternoon approvals rtnweapv 7 day total number of weekend evening approvals rtnwemor 7 day total number of weekend morning approvals rtnwenit 7 day total number of weekend night approvals rvnaudl 7 day variance of dollars per auth across window Profile Cardholder Fraud Related Variables paudymdy profile ratio of auth days over number of days in the month payapydl profile mean dollar amount for an approval payaudl profile mean dollars per auth across month pchdzip profile the last zip of the cardholder pdbm profile value of `date became member` at time of last profile update pddapv profile daily mean dollars of approvals pddapv2 profile daily mean dollars of approvals on days with auths pddau profile daily mean dollars of auths on days with auths pddau30 profile daily mean dollars of auths on all days in month pddcsapy profile daily mean dollars of cash approvals pddcsdec profile daily mean dollars of cash declines pdddec profile daily mean dollars of declines pdddec2 profile daily mean dollars mean dollars of merchandise approvals pddmrdec profile daily mean dollars of merchandise declines pdnapy profile daily mean number of approvals pdnau profile daily mean number of auths on days with auths pdnau30 profile daily mean number of auths on all days in month pdncsapy profile daily mean number of cash approvals pdncsdec profile daily mean number of cash declines pdndec profile daily mean number of declines pdnmrapy profile daily mean number of merchandise approvals pdnmrdec profile daily mean number of merchandise declines pdnw1ap2 profile mean number of approvals on Sundays which had auths pdnw1apv profile mean number of approvals on Sundays (day 1 of week) pdnw2ap2 profile mean number of approvals on Mondays which had auths pdnw2apv profile mean number of approvals on Mondays (day 2 of week) pdnw3ap2 profile mean number of approvals on Tuesdays which had auths pdnw3apv profile mean number of approvals on Tuesdays (day 3 of week) pdnw4ap2 profile mean number of approvals on Wednesdays which had auths pdnw4apv profile mean number of approvals on Wednesdays (day 4 of week) pdnw5ap2 profile mean number of approvals on Thursdays which had auths pdnw5apy profile mean number of approvals on Thursdays (day 5 of week) pdnw6ap2 prdfile mean number of approvals on Fridays which had auths pdnw6apy profile mean number of approvals on Fridays (day 6 of week) pdnw7ap2 profile mean number of approvals on Saturdays which had auths pdnw7apy profile mean number of approvals on Saturdays (day 7 of week) pdnwdaft profile daily mean number of weekday afternoon approvals pdnwdapv profile daily mean number of weekday approvals pdnwdeve profile daily mean number of weekday evening approvals pdnwdmor profile daily mean number of weekday morning approvals pdnwdnit profile daily mean number of weekday night approvals pdnweaft profile daily mean number of weekend afternoon approvals pdnweapy profile daily mean number of weekend approvals pdnweeve profile daily mean number of weekend evening approvals pdnwemor profile daily mean number of weekend morning approvals pdnwenit profile daily mean number of weekend night approvals pexpir profile expiry date stored in profile; update if curr date>pexpir phibal profile highest monthly balance phidcapy profile highest dollar amt on a single cash approve in a month phidcdec profile highest dollar amt on a single cash decline in a month phidmapy profile highest dollar amt on a single merch approve in a month phidmdec profile highest dollar amt on a single merch decline in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dollar amount on a single approve in a month phidsapy profile highest dolla profile highest dollar amount on a single decline in a month phidtapy profile highest total dollar amount for an approve in a single day phidtau profile highest total dollar amount for any auth in a single day phidtdec profile highest total dollar amount for a decline in a single day phinapy profile highest number of approves in a single day phinau profile highest number of auths in a single day phindec profile highest number of declines in a single day pm1avbal profile average bal. during 1st 10 days of mo. pm1nauths profile number of auths in the 1st 10 days of mo. pm2avbal profile average bal. during 2nd 10 days of mo. pm2nauths profile number of auths in the 2nd 10 days of mo. pm3avbal profile average bal. during remaining days pm3nauths profile number of auths in the last part of the month. pmovewt profile uses last zip to determine recent residence move; pmovewt =2 for a move within the previous calendar month; pmovew pnaudy profile number of days with auths pnauw1 profile number of Sundays in month with any auths pnauw2 profile number of Mondays in month with any auths pnauw3 profile number of Tuesdays in month with any auths pnauw4 profile number of Wednesdays in month with any auths pnauw5 profile number of Thursdays in month with any auths pnauw6 profile number of Fridays in month with any auths pnauw7 profile number of Saturdays in month with any auths pnauw6 profile number of weekday days in month with any auths pnauwe profile number of weekend days in month with any auths pncsaudy profile number of days in month with cash auths pnmraudy profile number of days in month with merchant auths pnweekday profile number of weekday days in the month pnweekend profile number of weekend days in the month pratdcau profile ratio of declines to auths profage profile number of months this account has had a profile (up to 6 mo.) psdaudy profile standard dev. of # days between transactions in a month psddau profile standard dev. of \$ per auth in a month ptdapy profile total dollars of approvals in a month ptdau profile total dollars of auths in a month ptdaudy profile total dollars of auths in a day ptdcsapy profile total dollars of eash advance approvals in a month ptdcsdec profile total dollars of cash advance declines in a month ptddec profile total dollars of declines in a month ptdmrapy profile total dollars of merchandise approvals in a month ptdmrdec profile total

Imperial College London

Workshop on Data Analysis and Classification16In honour of Edwin Diday

dollars of merchandise declines in a month ptdsfa01 profile total dollars of transactions in SIC factor group 01 ptdsfa02 profile total dollars of transactions in SIC factor group 02 ptdsfa03 profile total dollars of transactions in SIC factor group 03 ptdsfa04 profile total dollars of transactions in SIC factor group 04 ptdsfa05 profile total dollars of transactions in SIC factor group 05 ptdsfa06 profile total dollars of transactions in SIC factor group 06 ptdsfa07 profile total dollars of transactions in SIC factor group 07 ptdsfa08 profile total dollars of transactions in SIC factor group 08 ptdsfa09 profile total dollars of transactions in SIC factor group 09 ptdsfa10 profile total dollars of transactions in SIC factor group 10 ptdsfa11 profile total dollars of transactions in SIC factor group 11 ptdsra01 profile total dollars of transactions in SIC fraud rate group 01 ptdsra02 profile total dollars of transactions in SIC fraud rate group 02 ptdsra03 profile total dollars of transactions in SIC fraud rate group 03 ptdsra04 profile total dollars of transactions in SIC fraud rate group 04 ptdsra05 profile total dollars of transactions in SIC fraud rate group 05 ptdsra06 profile total dollars of transactions in SIC fraud rate group 06 ptdsra07 profile total dollars of transactions in SIC fraud rate group 07 ptdsva01 profile total dollars in SIC VISA group 01 ptdsva02 profile total dollars in SIC VISA group 02 ptdsva03 profile total dollars in SIC VISA group 03 ptdsva04 profile total dollars in SIC VISA group 04 ptdsva05 profile total dollars in SIC VISA group 05 ptdsva06 profile total dollars in SIC VISA group 06 ptdsva07 profile total dollars in SIC VISA group 06 ptdsva06 profile total dollars in SIC VISA group 06 ptdsva06 profile total dollars in SIC VISA group 06 ptdsva06 ptd 07 ptdsva08 profile total dollars in SIC VISA group 08 ptdsva09 profile total dollars in SIC VISA group 09 ptdsva10 profile total dollars in SIC VISA group 10 ptdsva11 profile total dollars in SIC VISA group 09 ptdsva09 profile total dollars in SIC VISA group 09 ptdsva09 profile total dollars in SIC VISA group 09 ptdsva09 profile total dollars in SIC VISA group 09 ptdsva09 profile total dollars in SIC VISA group 09 ptdsva09 ptdsva0 11 ptnapy profile total number of approvals in a month ptnapydy profile total number of approves a day ptnau profile total number of auths in a month ptnau10d profile number of auths in month <= \$10 ptnaudy profile total number of auths in a day ptncsapy profile total number of cash advance approvals in a month ptncsdec profile total number of cash advance declines in a month ptndec profile total number of declines in a month ptndecdy profile total number of declines in a day ptnmrapy profile total nurnher of merchandise approvals in a month ptnmrdec profile total number of merchandise declines in a month ptnsfa01 profile total number of transactions in SIC factor group 01 ptnsfa02 profile total number of transactions in SIC factor group 02 ptnsfa03 profile total number of transactions in SIC factor group 03 ptnsfa04 profile total number of transactions in SIC factor group 04 ptnsfa05 profile total number of transactions in SIC factor group 05 ptnsfa06 profile total number of transactions in SIC factor group 06 ptnsfa07 profile total number of transactions in SIC factor group 07 ptnsfa08 profile total number of transactions in SIC factor group 08 ptnsfa09 profile total number of transactions in SIC factor group 09 ptnsfa10 profile total number of transactions in SIC factor group 10 ptnsfa11 profile total number of transactions in SIC factor group 11 ptnsra01 profile total number of transactions in SIC factor group 10 ptnsfa11 profile total number of transactions in SIC factor group 10 ptnsfa10 ptnsf group 01 ptnsra02 profile total number of transactions in SIC fraud rate group 02 ptnsra03 profile total number of transactions in SIC fraud rate group 03 ptnsra04 profile total number of transactions in SIC fraud rate group 04 ptnsra05 profile total number of taansactions in SIC fraud rate group 05 ptnsra06 profile total number of transactions in SIC fraud rate group 06 ptnsra07 profile total number of transactions in SIC fraud rate group 07 ptnsva01 profile total number in SIC VISA group 01 ptnsva02 profile total number in SIC VISA group 02 ptnsva03 profile total number in SIC VISA group 03 ptnsva04 profile total number in SIC VISA group 04 ptnsva05 profile total number in SIC VISA group 05 ptnsva06 profile total number in SIC VISA group 06 ptnsva07 profile total number in SIC VISA group 07 ptnsva08 profile total number in SIC VISA group 08 ptnsva09 profile total number in SIC VISA group 09 ptnsva10 profile total number in SIC VISA group 10 ptnsva11 profile total number in SIC VISA group 11 ptnw1apv profile total number of approvals on Sundays (day 1 of week) ptnw2apv profile total number of approvals on Mondays (day 2 of week) ptnw3apv profile total number of approvals on Tuesdays (day 3 of week) ptnw4apv profile total number of approvals on Wednesdays (day 4 of week) ptnw5apv profile total number of approvals on Thursdays (day 5 of week) ptnw6apv profile total number of approvals on Fridays (day 6 of week) ptnw7apv profile total number of approvals on Saturdays (day 7 of week) ptnwdaft profile total number of weekday afternoon approvals in a month ptnwdapv profile total number of weekday approvals in a month ptnwdeve profile total number of weekday evening approvals in a month ptnwdmor profile total number of weekday morning approvals in a month ptnwdnit profile total number of weekday night approvals in a month ptnweaft profile total number of weekend afternoon approvals in a month ptnweapy profile total number of weekend approvals in a month ptnweeve profile total number of weekend evening approvals in a month ptnwemor profile total number of weekend morning approvals in a month ptnwenit profile total number of weekend morning approvals in a month pydaybtwn profile variance in number of days between trx's (min of 3 trx) pyraudl profile variance of dollars per auth accoss month MERCHANT FRAUD VARIABLES mtotturn Merchant Total turnover for this specific merchant msicturn Merchant Cumulative SIC code turnover mctrtage Merchant Contract age for specific merchant maagsic Merchant Average contract age for this SIC code mavgnbtc Merchant Average number of transactions in a batch maamttrx Merchant Average amount per transaction (average amount per authorizations) myaramt Merchant Variance of amount per transaction may at the second s Merchant Average time between batches mavgtaut Merchant Average time between authorizations for this merchant mratks Merchant Ratio of keyed versus swiped transactions mnidclac Merchant Number of identical customer accounts mnidcham Merchant Number of identical charge amounts mtrxsrc Merchant What is the source of transaction (ATM, merchant, etc.) mtrxtrsp Merchant How is the transaction transported to the source (terminal, non-terminal, voice authorization) mfloor Merchant Floor limit mchgbks Merchant Charge-backs received mrtrvs Merchant Retrievals received (per SIC, merchant, etc.). The issuer pays for a retrieval, macgrat Merchant Acquirer risk managment rate (in Europe one merchant can have multiple acquires, but they dont have records about how many or who.) morevrsk Merchant Previous risk management at this merchant? Yes or No mtyprsk Merchant Type of previous risk management (counterfeit, multiple imprint, lost/stolen/not received) msicrat Merchant SIC risk management rate mpctaut Merchant Percent of transactions authorized

Imperial College London

Workshop on Data Analysis and Classification17In honour of Edwin Diday

Unbalanced classes

Detectorcorrectly identifies 99 in 100 legitimate transactionsandcorrectly identifies 99 in 100 fraudulent transactions

Pretty good?

But suppose only 1 in 1000 transactions are fraudulent

Imperial College London

Workshop on Data Analysis and Classification18In honour of Edwin Diday

		True class		
		Legit	Fraud	
Predicted	Legit	99%	1%	
class	Fraud	1%	99%	
Numbers		999	1	

		True cl	ass	
		Legit	Fraud	
Predicted	Legit	989.01	0.01	
class	Fraud	9.99	0.99	0.99 / (9.99+0.99) = 0.09
Numbers		999	1	

Imperial College London

Workshop on Data Analysis and Classification19In honour of Edwin Diday

91% of suspected frauds are in fact legitimate

This matters because:

- operational decisions must be made (stop card?)
- good customers must not be irritated

Customers are pleased you care: up to a point

Workshop on Data Analysis and Classification In honour of Edwin Diday



Delay in learning class labels

- if fraud alarm is raised, then true class quickly known
- if no alarm, then not detected until statement

This makes it different from the standard supervised classification paradigm

Banks cannot always say for sure when a fraud commences

Mislabelled classes

Not all fraudulent transactions are labelled as fraud (account holder fails to check carefully)

Not all legitimate transactions are labelled as legitimate

There may be subtleties

e.g. account holder makes transactions and then claims card was stolen

Such transactions are fraudulent because the holder declares them as such

Reactive population drift

- banks implement detection/prevention strategies
- fraudsters don't generally give up! but change strategies

Reactive population drift example 1: *Chip and PIN*

Chip and PIN intended/predicted to end card fraud

After UK rollout on 14 Feb 06, CC fraud in UK did decline How much was a consequence of the publicity?

but

- predicted to lead to increase in identity theft

and

- Lloyds TSB observed increase in fraudulent use of UK cards in Europe (no C&P – mag stripe still counterfeited)
- observed increase in ATM and cardholder not present fraud
- in fact, crooks installed data skimmers into C&P terminals (such devices can be purchased for < £100), over £1m stolen from Shell gas stations

Imperial	College
London	

Workshop on Data Analysis and Classification24In honour of Edwin Diday

Plastic card fraud in the UK (Gordon Blunt)



Imperial College London Workshop on Data Analysis and Classification25In honour of Edwin Diday

What is a good system?

'Classifies fraudulent transactions as fraudulent, and legitimate transactions as legitimate' ?

But: no method is perfect Need: criteria for assessing effectiveness

Timeliness: time scale: count of fraud transactions misclassified

Standard two class classification criteria inadequate:

- misclassification rate: treats two types of misclassification equally
- Gini coefficient (AUC): averages over all misclassification cost ratios
- Kolmogorov-Smirnov statistic: data driven cost ratio

Unbalanced classes

		True class		
Predicted class		Fraud	Legitimate	
	Fraud	A	В	
	Legitimate	С	D	

A very well known consumer credit organisation evaluates fraud using the two ratios

$$R_1 = A/(A+C)$$
 (= Sensitivity)
 $R_2 = B/(A+B)$ (= 1- Precision)

Imperial College London

Workshop on Data Analysis and Classification27In honour of Edwin Diday

In itself, this would appear to be fine

But in fact, the units of assessment they use are *accounts*

An account is flagged as potentially fraudulent if *at least one transaction is so flagged*

Problem 1: This means that one can make the probability of flagging an account as fraudulent as near to 1 as one wishes by examining enough transactions

Problem 2: Fails to include *timeliness* in the measure

A superior measure

Consider each series of transactions ending with either (i) a *fraud flag* on a true fraud Or

(ii) or end of observed sequence

nnnnfnnfnn**n**nnfnnnnn**n**nnnn**f**

		True class	
Predicted class		Fraud	Legitimate
	Fraud	$m_{f/f}$	$m_{n/f}$
	Legitimate	$m_{f/n}$	$m_{n/n}$

Imperial College London

Workshop on Data Analysis and Classification29In honour of Edwin Diday

nnnfnnfnn**n**nnfnnnnn**n**nnnn**f**

		True class		
Predicted class		Fraud	Legitimate	
	Fraud	1	2	
	Legitimate	3	21	

Imperial College London

Workshop on Data Analysis and Classification30In honour of Edwin Diday

		True class		
		Fraud	Legitimate	
Predicted class	Fraud	$m_{f/f}$	$m_{n/f}$	
	Legitimate	$m_{f/n}$	$m_{n/n}$	

Overall performance measure for given threshold:

$$T_1 = \left(m_{f/f} + m_{n/f} + k m_{f/n} \right) / \left(k m_f + m_n \right)$$

where *k* is the estimated relative cost of misclassifying a fraud as legitimate compared to misclassifying a legitimate as fraud

Or, if the bank can afford to investigate C cases

$$T_2$$
: minimise $m_{f/n}$ subject to $\left(m_{f/f} + m_{n/f}\right) = C$

Imperial College London Workshop on Data Analysis and Classification31In honour of Edwin Diday

Constructing suspicion scores

Rules: detect known suspicious behaviour

e.g. most office workers do not shop in working hours

- Supervised classification: learn to distinguish fraud from legit
- Anomaly detection: detect anomalous behaviour
- Precursors: identify behaviour which precedes fraud
- Change points: change of state in account
- Multilevel: transaction/account/merchant
- Link analysis: networks

Example 1: Supervised classification (Chris Whitrow)

Basic principle:

Given a set of known fraudulent and legitimate transactions/accounts, along with descriptive variables for each, condense these to a rule enabling correct classification of new transactions/accounts from only the descriptive variables

Workshop on Data Analysis and Classification33In honour of Edwin Diday

Classification methods compared:

- logistic regression
- quadratic discrimination
- naive Bayes classifier
- decision tree
- k-nearest neighbour
- SVMs with radial basis kernels
- random forests

Bank A:

- 175 million transactions: 1st August 05 to 30th Nov 05
- 16.8 million accounts
- 5,946 accounts experienced fraud at POS terminals in obs'n period
- 76 variables per transaction; mostly categorical
- rolling window activity records 0, 1, 3, 7 days
- activity records sacrifice immediacy of individual transactions
- but potential for more accuracy

Two explorations:

1: *Random*: Train on random 70%, test on remainder

2. Prediction: Train up to 30th Oct 05, test after

Imperial College London

Workshop on Data Analysis and Classification36In honour of Edwin Diday

Random performance



Imperial College London Workshop on Data Analysis and Classification37In honour of Edwin Diday

Example 2: One class modelling: outliers (Piotr Juszczak)

Basic principle: *build a model for the 'norm' for this customer and detect when it deviates*

'Norm' can be based on

- this customer compared with other customers
- this customer compared with self at previous times
- a combination of these

Basic advantage of one-class approach

- can detect new kinds of anomalies, not seen before
- more power in dynamic fraud environment?

Bank B:

- 44,637 accounts
- 2,374,311 transactions
- 3,742 fraudulent accounts
- 53,844 fraudulent transactions
- 3 months data

77 variables, we used

- size of transaction
- difference between current and previous transaction size
- sum of current and previous transaction sizes
- product of current and previous transaction sizes
- time of transaction
- time between current and previous transaction
- merchant category
- ATM ID code



Imperial College London

Workshop on Data Analysis and Classification40In honour of Edwin Diday

Preprocessing the categorical variables (MCC and ATM)

A(j,i) = no. times ATM *j* is accessed from account *i*

$$ATM(j) = (A(j,1), A(j,2), \dots, A(j,K))^{T}$$

- \Rightarrow dissimilarity matrix between ATMs
- \Rightarrow reduce dimensionality of ATMs using MDS
- \Rightarrow combine with continuous variables

Similar for MCCs



Imperial College London Workshop on Data Analysis and Classification42In honour of Edwin Diday

Used several methods for building the pdfs:

- Parzen kernel
- Naive Bayes with Parzen kernel for each variable
- Single multivariate Gaussian
- Mixture of multivariate Gaussian
- 1-nearest neighbour
- Support Vector Data Description
- Self-Organising map
- Minimum spanning tree data description
- Minimax probability machine



Imperial College London

Workshop on Data Analysis and Classification44In honour of Edwin Diday



Imperial College London Workshop on Data Analysis and Classification45In honour of Edwin Diday

Example 3: One class modelling: Peer group analysis (Dave Weston)

Individual account profiles:

Model behaviour and compare new transaction with past

But

Spending behaviour just before Christmas is anomalous Individual profile models may flag such transactions

So

Identify others with similar past behaviour (peer group) Compare new transaction with their new transactions Target account tracks peer group



Peer group quality: dispersion about past target behaviour

A clustering time series problem

Imperial College	Workshop on Data Analysis and Classification	47
London	In honour of Edwin Diday	

Bank C:

4,159 accounts, with at least 80 transactions over a 4 month period with no fraud in first 3 months: 241 had fraud in last month

Build 4,159 peer groups using first 3 months data

Split 3 months into *n* windows and summarise each window by

- total amount withdrawn
- total number of transactions
- entropy of MCCs
- \Rightarrow combine windows to give up to 3*n* dmensional space for finding peer groups

Imperial College London

Workshop on Data Analysis and Classification48In honour of Edwin Diday

Population outlier detection - robustified peer group



Imperial College London

Workshop on Data Analysis and Classification49In honour of Edwin Diday

Conclusions

Fraud detection problems

- may involve high dimensions, messy data, large n
- typically have unbalanced classes
- often have mislabelled classes, delay in labelling
- may involve dynamic, reactive data distributions

There are

- many approaches / different aspects
- issues of how to measure performance

Other, deeper questions

The economic imperative

About methodology

How much do we learn from ad hoc comparisons of methods on particular data sets?

About society

Is society changing? Accepting some degree of fraud? Fraud management requires a holistic approach, blending tactical and strategic solutions with the state-of-the-art technology solutions and best practice in fraud strategy and operations

James Gilmour, Editor Credit Risk International, 2003

Imperial College London

Workshop on Data Analysis and Classification52In honour of Edwin Diday



<u>d.j.hand@imperial.ac.uk</u>

http://stats.ma.ic.ac.uk/djhand/public_html/

Imperial College London

Workshop on Data Analysis and Classification53In honour of Edwin Diday