

On the Approximate Communal Fraud Scoring of Credit Applications

Clifton Phua
 PHD Candidate
clifton.phua@infotech.monash.edu.au
<http://www.bsys.monash.edu.au/people/cphua>
 2005

Other contributors:
 Dr. Ross Gayler
 Assoc. Prof. Vincent Lee
 Assoc. Prof. Kate Smith

Overview

- Credit applications, fraud detection, and existing techniques
- Pair-wise matching, approximation of pair-wise communal score, smoothed k -wise score, and temporal and spatial weights
- Experiments: synthetic data, training and scoring phases, and results
- Discussion: directed graph, types of subgraphs, prediction scores, and related work

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 2

Importance of credit application fraud detection

- Each credit application contains identity and non-identity information
- Credit bureaux collect millions of enquiries related to credit applications
- Identity crime
- First stage of credit life cycle

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 3

Existing techniques

- Attribute validation / verification rules
- Credit history matching
- Black lists
- Classification

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 4

Pair-wise matching (Part 1)

- Better timeliness, sensitivity, and specificity
- Pair-wise matching with application / time window, training and scoring phases

app_id	date_received	...	empidm.com
leg4-6	31/12/2004	...	
leg4-5	30/12/2004	...	
leg4-4	30/12/2004	...	
leg4-3	29/12/2004	...	
leg4-2	28/12/2004	...	
leg4-1	28/12/2004	...	
...
leg4-0	14/02/2004	...	0.035
leg4-iii	6/02/2004	...	0
leg4-i	26/01/2004	...	0.9
leg4-ii	26/01/2004	...	1.55
leg4-i	25/01/2004	...	0

Before Training

After Training

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 5

Pair-wise matching (Part 2)

- Pair-wise relationship with case-sensitive exact / fuzzy and cross matching
- Permanent, stable, and transient identity attributes

date received	given_name	surname	street_number	current_address	previous_address	radius	postcode	state	home_phone
31/12/2004	jeanette	pan	3	jean avenue	mariah avenue	clayton	6000	vic	85777756
7/1/2004	chuan wei	pan	2	erwin jean	mariah ave	clayton	3168	vic	85777756
6/1/2004	chuan wei	pan	2	erwin jean	mariah ave	clayton	3168	vic	85777756
5/1/2004	jeanette	pan	3	jean ave	mariah avenue	clayton	3168	vic	85777756
4/1/2004	jeanette	pan	3	jean avenue	mariah ave	clayton	3168	vic	85777756
3/1/2004	jeanette	pan	3	jean avenue	mariah avenue	clayton	3168	vic	85777756
2/1/2004	jeanette	pan	3	jean avenue	mariah avenue	clayton	3168	vic	85777756

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 6

Approximation of the pair-wise communal score (Part 1)

- Black list
 - Outlinks to known fraud applications
 - Outlinks to applications with excessive inlinks
- White list
 - Legitimate relationships
- Anomalous relationships
- No relationships

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 7

Approximation of the pair-wise communal score (Part 2)

- Score range, advantage, disadvantage, and secondary use

Application Pairs	Black List	White List	Anomalous Links	Unlinked
Score range	Highest	Lowest	Medium to high	None
Main advantage	Reasonably accurate	Filters normal relations	Finds irregular relations	Contains weak or no relations
Main disadvantage	Time delay	Prone to manipulation	Reliant on attribute weights	Unlinked (once-off) fraud
Secondary use other than real-time fraud detection	Update attribute, normal, temporal, spatial weights	Viral marketing	Error detection	-

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 8

Smoothed k -wise scoring function

$$S(v_i) = \sum_{v_j \in N(v_i)} \left[(1 - \alpha) * w_{ij} + \alpha * \frac{S(v_j)}{E_o(v_j)} \right]$$

- v_i refers to the current application which is being scored; v_j refers to the previous application which has been scored
- High suspicion score roughly captures the strong connection between v_i and v_j (high link weight), and/or the quality of v_j (high average suspicion scores), and/or the quantity of v_j (large number of connected applications); and vice versa
- α controls the contribution of previous scores of v_j

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 9

Communal, temporal, and spatial weights

$$\tilde{w}_{ij} = w_{communal_{ij}} * w_{temporal_{ij}} * w_{spatial_{ij}}$$

- \tilde{w}_{ij} roughly captures the relative strength of the communal links over time and space
- $w_{temporal}$ is the link weight derived from pair-wise time difference of *date_received*, $w_{spatial}$ is the link weight derived from pair-wise geographical distance of *postcode*
- High link weight roughly indicates strong connection, and/or no time difference, and/or maximum geographic distance between v_i and v_j ; and vice versa

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 10

Experiments - synthetic data

- External concerns: Privacy, confidentiality, ethical, and competitive; Technical reasons: Data set size, population drift, concept drift, adversarial countermeasures, and data entry error rates
- Original records, duplicate records (uniform and poisson distributions), error rates
- 52,750 applications, 19 attributes which span the entire year 2004. ~10% are fraudulent (regular frauds, occasional frauds, seasonal frauds, and once-off frauds). ~90% are legal applications and 50 are hand-crafted applications with both fraudulent and legal examples

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 11

Experiments - training & scoring phases

- All attributes are included except *rec_id* (all unique), *date_received* (for calculating $w_{temporal}$), *postcode* (for calculating $w_{spatial}$) *street_number*, *suburb*, and *state* (too dense). Also, three additional attributes are created to store *suspicion_score*, *outlinks*, and *inlinks*
- Train 10,000 applications (01/01/2004 to 10/03/2004); score 42,750 applications from (10/03/2004 to 31/12/2005). The scoring phase retains trained applications which have non-zero scores

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 12

Experiments - results

- Parameter tuning and scalability issues
- Significant change of link activity in the scoring phase over the training phase which is probably due to seasonal frauds

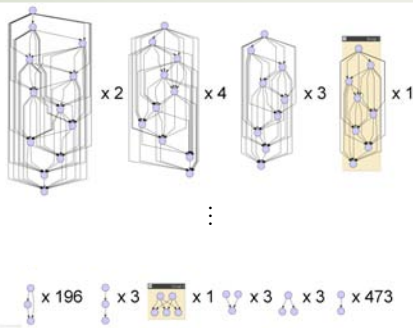
	Training	Scoring
g	10,000	44,238
h	10,000	10,000
$T_{max,org}$	0.5	0.5
T_g	10	10
$T_{max} = T_{max,org}$	3	3
W_{ij}	126	126
α	0.8	0.8
h	2,917	17,968
$computation_{org}$	2.7 min-sec 4.2 min	4.3 min-sec 6.68 min
$\bar{d}_1 = \bar{d}_D = \frac{h}{g}$	0.292	0.418
$\Delta = \frac{h}{g(g-1)}$	0.000029	0.00001
$\frac{\sum_{i=1}^h S(v_i)}{g}$	0.116	0.197

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 13

Discussion - directed graph

- The descriptive directed graph allow the analyst to visually inspect / explore the subgraphs
- The following slide presents the compressed hierarchical subgraph structures of linked applications after training.
- The synthetic data generation process causes all the 27 different types of subgraphs to be disconnected although this is unlikely with real applications

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 14

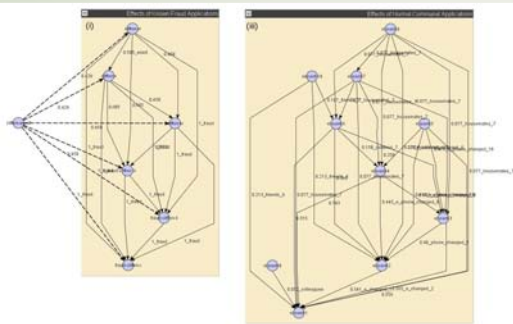


On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 15

Discussion - types of directed subgraphs

- Subgraph (i) consists of known fraud applications and the subsequent linked applications.
- Subgraph (ii) is made up of linked applications which have frequent address changes by the same identity.
- In addition to subgraph (ii), subgraph (iii) includes linked applications submitted by an identity's social network.
- Subgraph (iv) consists of linked applications with data entry errors which may also turn out to be frauds.
- Subgraph (v) illustrates synthetic applications which "mix and match" attributes from some other previous applications.
- Subgraph (vi) have all exact applications to demonstrate the effects of temporal and spatial weights.
- Subgraph (vii) shows the effects of exponential smoothing.

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 16



On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 17

Discussion - prediction scores

- The predictive suspicion scores allow the analyst to rank the most recent applications and investigate the most suspicious applications
- The decision thresholds can be defined as:
 - If $0 < S(v_i) \leq T_{lower}$, then current application v_i is unusual
 - If $T_{lower} < S(v_i) < T_{upper}$, then v_i is suspicious
 - If $S(v_i) \geq T_{upper}$, then v_i should be investigated
- If the investigated v_i is fraudulent, feedback should be provided to the fraud detection system
- If $T_{lower} = 0.8$ and $T_{upper} = 1$, then out of 50 applications, thirteen are be investigated (strong fraud symptoms) and three are considered suspicious (some fraud symptoms)

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 18

Related work

- There is no academic research, to the best of our knowledge, into the scoring of dynamic credit applications which accounts for its sparse-identifiers, communal, temporal, and spatial aspects
- However, there are other related and established application fields in multi-attribute pair-wise matching (for example, record linkage/deduplication detection), and single-attribute communal scoring/directed graphs with explicit links (for example, telecommunications fraud detection, terrorist detection, social network analysis, and webpage ranking)

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 19

The End of Presentation 😊

Questions?

On the Approximate Communal Fraud Scoring of Credit Applications (CSCC9) 20