

ランダムイズ加工

#なんちゃってOLH #RR

OLHを元とした”なんちゃってOLH”を用いた加工を実施
エンコード

- $v \rightarrow (H, h)$ H :ハッシュ関数, h :ハッシュ値
- $h \in [d']$, $d' = e^\varepsilon + 1$
 - 本来OLHでは h の値域は d' で指定されるが
 - 今回は各属性の値域と等しい

摂動: h に限定して行う

- 維持確率: $y = (H, h)$
- 遷移確率: $y = (H, a)$, $a \neq h$

$$p' = \frac{e^\varepsilon}{e^\varepsilon + d' - 1}$$
$$q' = \frac{1}{e^\varepsilon + d' - 1}$$

	Gender	Age	Zipcode	56
田代	F	15	27	0
山本	M	20	390	2
三浦	F	10	530	3
関口	F	20	695	1
茂呂	M	70	910	0

	Gender	Age	Zipcode	56
田代	F	15	910	1
山本	M	15	530	2
三浦	F	20	530	3
関口	M	20	27	5
茂呂	M	70	390	0

p'	$\frac{e^2}{e^2+1}$	$\frac{e^2}{e^2+6}$	$\frac{e^2}{e^2+494}$	$\frac{e^2}{e^2+5}$
	≈ 0.88	≈ 0.55	≈ 0.01	≈ 0.60

例)なんちゃってOLH, $\varepsilon=2$

	$\varepsilon=1$	$\varepsilon=2$	$\varepsilon=3$
有用性	55.2	54.57	55.68
個人特定攻撃成功数	3/50	1/50	8/50
DP再構築攻撃成功数	23/50	12/50	31/50
合計攻撃成功数	26/100	13/100	39/100
匿名性スコア	74	87	61

$\varepsilon=1, 2, 3$ で実験を行い, ε を決定

DB再構築攻撃

#協調フィルタリング #CF

00

マスクを同アイテムの
最頻値で仮埋め

	2	56	247	260	653
田代	*	1	4	2	3
山本	1	4	0	0	*
三浦	2	*	2	0	3
関口	3	1	*	1	4
茂呂	1	5	0	*	1

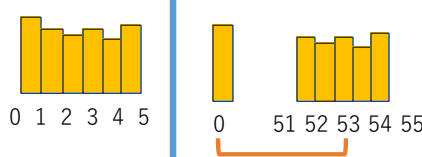
02

マスクを類似度各Top5の計
10個の評価値の
最頻値として予想
(最頻値が複数の場合は類
似度の高い値を選択)

	2	56	247	260	653
田代	*	1	4	2	3
山本	1	4	0	0	*
三浦	2	*	2	0	3
茂呂	3	1	*	1	4
関口	1	5	0	*	1

+50

未視聴(0)の値値と視聴済(1-5)
の値値に差をつける



ユーザ間のコサイン類似度表+50

	田代	山本	三浦	茂呂	関口
田代	1.00	0.76	0.89	0.88	0.76
山本	0.76	1.00	0.86	0.87	0.99
三浦	0.89	0.86	1.00	0.75	0.86
茂呂	0.88	0.87	0.75	1.00	0.87
関口	0.76	0.99	0.86	0.87	1.00

01

ユーザ同士・アイテム同士の
コサイン類似度を計算

	2	56	247	260	653
2	1.00	0.57	0.45	0.56	0.94
56	0.57	1.00	0.20	0.20	0.61
247	0.45	0.20	1.00	0.80	0.61
260	0.56	0.20	0.80	1.00	0.67
653	0.94	0.61	0.61	0.67	1.00

ユーザ間のコサイン類似度表

	田代	山本	三浦	茂呂	関口
田代	1.00	0.49	0.85	0.62	0.31
山本	0.49	1.00	0.69	0.72	0.91
三浦	0.85	0.69	1.00	0.87	0.45
茂呂	0.62	0.72	0.87	1.00	0.44
関口	0.31	0.91	0.45	0.44	1.00

映画間のコサイン類似度表+50

	2	56	247	260	653
2	1.00	0.98	0.63	0.63	0.99
56	0.98	1.00	0.61	0.61	0.99
247	0.63	0.61	1.00	0.51	0.63
260	0.67	0.61	0.51	1.00	0.64
653	0.99	0.99	0.63	0.64	1.00

	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	Σ	AVG
All0	23	22	21	18	11	13	23		16	19	17	14		21	14	20	21	20		20	16	14	343	18.1
User	18	14	18	18	10	9	18		13	19	13	16		17	9	13	17	19		18	17	18	294	15.5
Item	23	17	16	20	11	10	19		14	14	15	9		20	14	16	16	17		19	15	16	301	15.8
W	24	14	19	19	12	10	22		15	17	14	16		19	11	14	18	21		20	18	17	320	16.8
W+50	26	14	18	15	12	12	21		15	19	13	16		20	13	14	17	20		18	20	20	323	17.0

有用性向上

#クロス集計 #スワップ

01

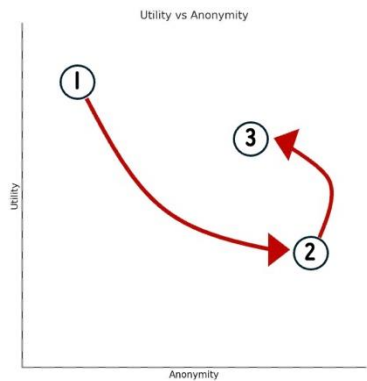
有用性と安全性の関係

右図にデータ加工の推移を示す

- 配布データBXX.csv
- なんちゃってOLH加工後
- 有用性向上プログラム実行後

安全性スコアをなんちゃってOLHにより向上
他チームのレベルまで上げることに成功

一方、有用性スコアは低下したため、ル
ール上の有用性向上のためプログラムを実行



02

有用性を上げるプログラム

以下の4ステップを繰り返す

- 元データBXX.csvとCXX_{i}.csvの
全共通カラムについてクロス集計
の差分を計算
- 1の計算結果の中で絶対誤差の合計が最も
大きい列の組を選択
- 2で選んだ列の組について、8ファイルに
対しての出現回数が最大のペアと最小の
ペアを選択 (右表の赤字)
- 3で選んだペアをもつ行をランダムで1行
選択し、最大のペアの内容を
最小のペアの内容で上書き (選択した列
の組の絶対誤差の合計を2低減)

	0	1	2	3	4	5
M	3096	484	335	560	671	391
F	2006	519	346	482	656	454

	0	1	2	3	4	5
M	1908	598	641	757	869	633
F	1523	502	630	670	671	598

	0	1	2	3	4	5
M	-1188	114	306	197	198	242
F	-483	-17	284	188	15	144

03

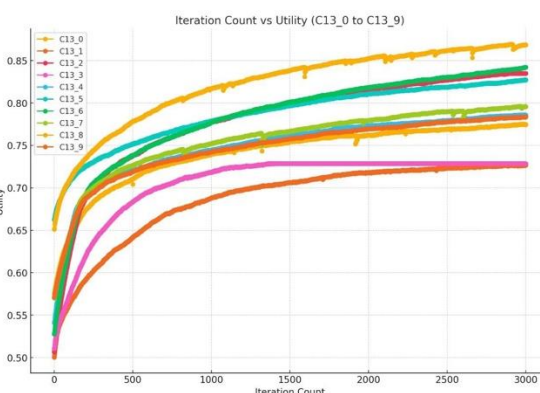
反復回数と有用性のグラフ

有用性向上のプログラムを3000回実行
列数の多いCXX_1, CXX_3は試行回数を
増やしても有用性スコアが72辺りで収
束※1, 2

CXX_1, CXX_3に合わせて有用性スコア
72.2に揃えたファイルをまとめて提出

※1 列数が多いと、有用性向上プログラムで上書きした
列とその他の列の組み合わせの数も多くなり、一回の上
書きで他の組み合わせの絶対誤差の合計が大きくなっ
てしまうことがあるため

※2 右図は、別の加工ファイル (mixed) の有用性向上記
録だが上記と同様に有用性スコアは72辺りに収束する



個人特定攻撃

#ユークリッド距離 #全組み合わせ

00

BXXb_E.csvと
CXX_0~CXX_9.csv
の1-5の評価値に
+50

	2	56	247	260
A	3	0	2	4
B	5	3	0	0
C	5	0	1	4

	2	56	247	260
A	53	0	52	54
B	55	53	0	0
C	55	0	51	54

01

BXXb_E.csvとCXX_0~CXX_9.csv
の各ファイルについて、Ratingとの誤差
を計算し、各ユーザーごとに誤差が少な
いTOP10を算出
※ BXXb_E.csvはBXXb.csvの補完版

	2	56	247	260
田代	52	0	54	55
山本				
三浦				

以下49名

	2	56	247	260
A	53	0	52	54
B	55	53	0	0
C	55	0	51	54

1

23

4

02

TOP10とBXXa.csvを比較し
て基本属性が3つ以上一致する
ものを候補の行として保存

	Gender	Age	Occupat ion	Zip-code
田代	M	25	5	231
山本				

以下49名

	Gender	Age	Occup ation	Zip-code
A	F	18	25	640
C	M	25	19	231

03

各ジャンルの候補行の表

	C_0	C_1	C_2	C_3
田代	C	O	C	G
山本	B	U	F	J
三浦	P	F	O	P

他に、2,500攻撃・・・BXXa.csv, BXXb.csvを全部組み合わせ、全ジャンルの加エデータCとのハミング距離を計算し、各ジャンルのにもっともらしい行が確認できたaのindexを選ぶ攻撃を行いました。