

垂直分割における通信効率の良い 一致度の秘匿分散計算

青木 良樹[†] 菊池 浩明[†]

寺田 雅之[‡] 石井 一彦^{*} 関野 公彦^{*}

[†] 東海大学 [‡] NTTドコモ先進技術研究所 ^{*} NTTドコモサービス&ソリューション開発部

論文の訂正

3.2 提案方式：一致度の推定

誤

k 人の類似ユーザを作り、そのユーザ組の一致度を求める。

次に、一致度の推移性を利用し、任意の n^2 組の一致度を推定する。必要な通信コストは $nk = O(n)$ であり、効率が良いが、推定される一致度は誤差を含む。

Lathia らの提案した擬似ユーザを利用して $O(n^2)$ を n に削減する方法を提案する。まず、Algorithm?? をつかって、三組織間で擬似ユー

Algorithm 1 一致度の予測方法

入力: A の評価値 $r_{u,i}$, ($i \in I_A$), B の評価値 $r_{u,i}$, ($i \in I_B$), C の評価値 $r_{u,i}$, ($i \in I_C$).

出力: 擬似ユーザと各ユーザの C, T, D, N
 $\hat{C}_{u,c}, \hat{D}_{u,c}, \hat{T}_{u,c}, N$.

4. 式 (3, 4, 5) の区間推定を行い、全アイテム間 u, v について $[\underline{C}_{uv,c_1}, \overline{C}_{uv,c_1}], \dots, [\underline{C}_{uv,c_k}, \overline{C}_{uv,c_k}]$ を求め、その集約区間 $\underline{C}_{uv,*} = \bigvee_k \underline{C}_{uv,c_k}$, $\overline{C}_{uv,*} = \bigvee_k \overline{C}_{uv,c_k}$ を求める。 D と T についても同様に算出する。

3.2 提案方式：一致度の推定

正

k 人の類似ユーザを作り、そのユーザ組の一致度を求める。

次に、一致度の推移性を利用し、任意の n^2 組の一致度を推定する。必要な通信コストは $nk = O(n)$ であり、効率が良いが、推定される一致度は誤差を含む。

Lathia らの提案した擬似ユーザを利用して $O(n^2)$ を $O(n)$ に削減する方法を提案する。まず、Algorithm2 をつかって、三組織間で擬似

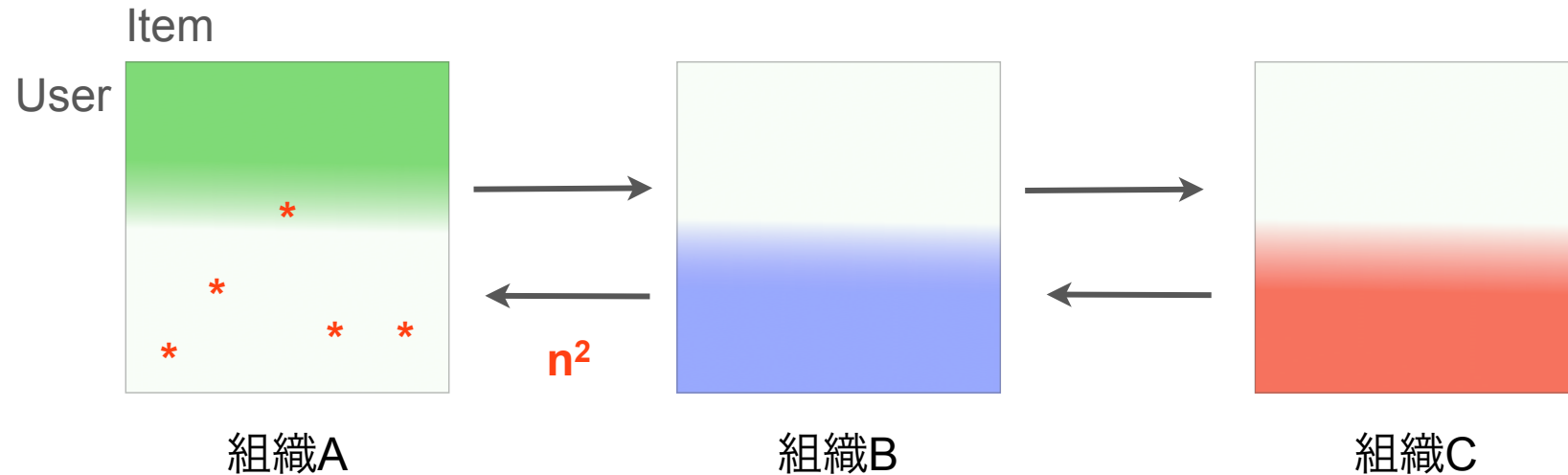
Algorithm 1 一致度区間の予測方法

入力: A の評価値 $r_{u,i}$, ($i \in I_A$), B の評価値 $r_{u,i}$, ($i \in I_B$), C の評価値 $r_{u,i}$, ($i \in I_C$).

出力: 集約区間 $[\underline{C}_{uv,*}, \overline{C}_{uv,*}]$, $[\underline{D}_{uv,*}, \overline{D}_{uv,*}]$, $[\underline{T}_{uv,*}, \overline{T}_{uv,*}]$.

4. 式 (3, 4, 5) の区間推定を行い、全アイテム間 u, v について $[\underline{C}_{uv,c_1}, \overline{C}_{uv,c_1}], \dots, [\underline{C}_{uv,c_k}, \overline{C}_{uv,c_k}]$ を求め、その集約区間 $\underline{C}_{uv,*} = \bigvee_k \underline{C}_{uv,c_k}$, $\overline{C}_{uv,*} = \bigwedge_k \overline{C}_{uv,c_k}$ を求める。 D と T についても同様に算出する。

垂直分割DM



目的

B,Cの評価値を基に，類似度を正確に予測する
それぞれの組織の評価値は秘匿

問題点

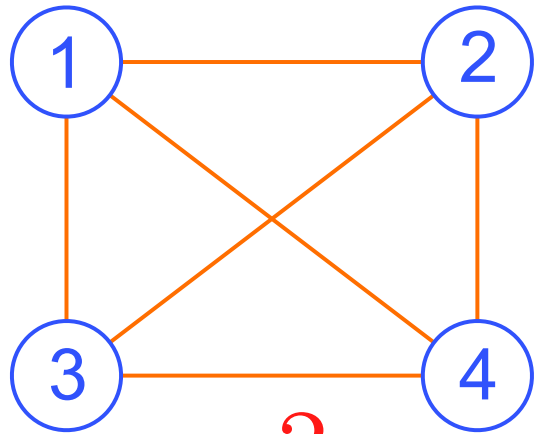
全ユーザ間の組み合わせ $O(n^2)$ の情報交換

既存研究

	手法	通信コスト
[VC03]	k-means	$O(n^2)$
[BVK11]	Slope One+CF	$O(n^2)$
本発表	Concordant	$O(kn)$

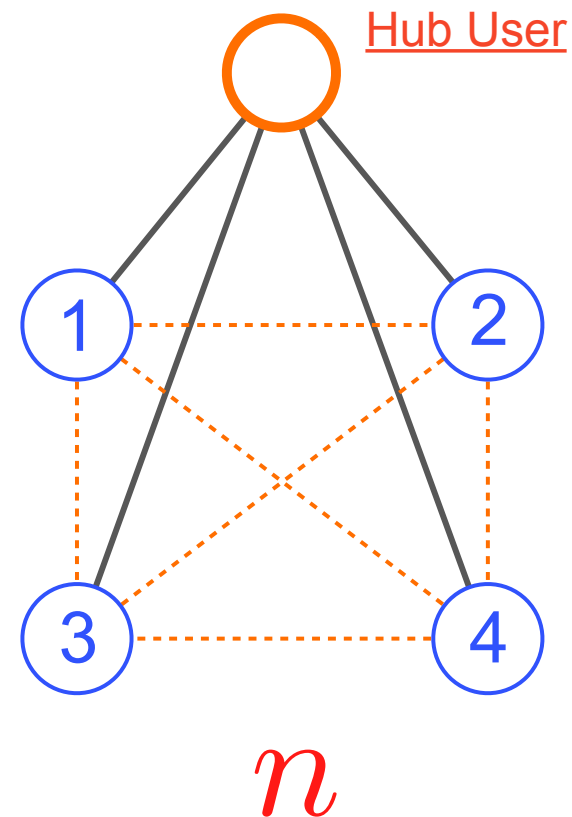
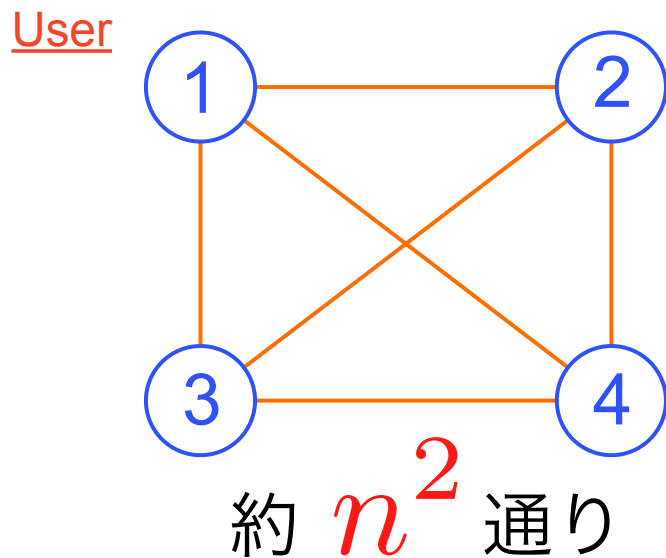
本アプローチの新規性

User



約 n^2 通り

本アプローチの新規性



+

Concordant(一致度)に基づく
類似度推定[Lathia 2007]

一致度の定義 (Somers' d)

一致 (**C**oncordant)

不一致 (**D**iscordant)

同順位 (**T**ied)

	i_1	i_2	i_3	i_4	i_5	
u_1	4	2	5	1	5	
	+	-	+	-	+	平均値より
u_2	1	2	3		4	+ ? - ?
	-	-	+	/	+	
$u_{1,2}$	d	c	c	t	c	

一致度の定義 (Somers' d)

一致 (**C**oncordant)

不一致 (**D**iscordant)

同順位 (**T**ied)

	i_1	i_2	i_3	i_4	i_5	
u_1	4	2	5	1	5	
	+	-	+	-	+	平均値より
u_2	1	2	3		4	+ ? - ?
	-	-	+	/	+	
$u_{1,2}$	d	c	c	t	c	

$$C_{1,2} = 3$$

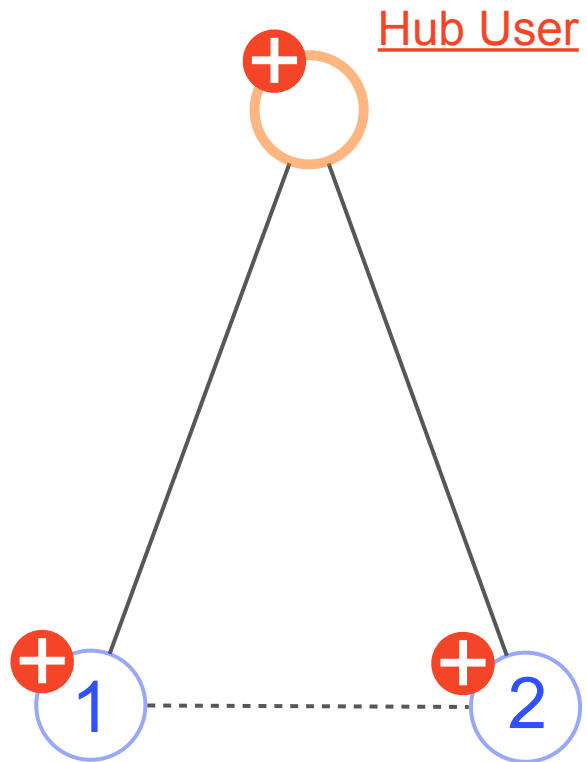
$$D_{1,2} = 1$$

$$T_{1,2} = 1$$

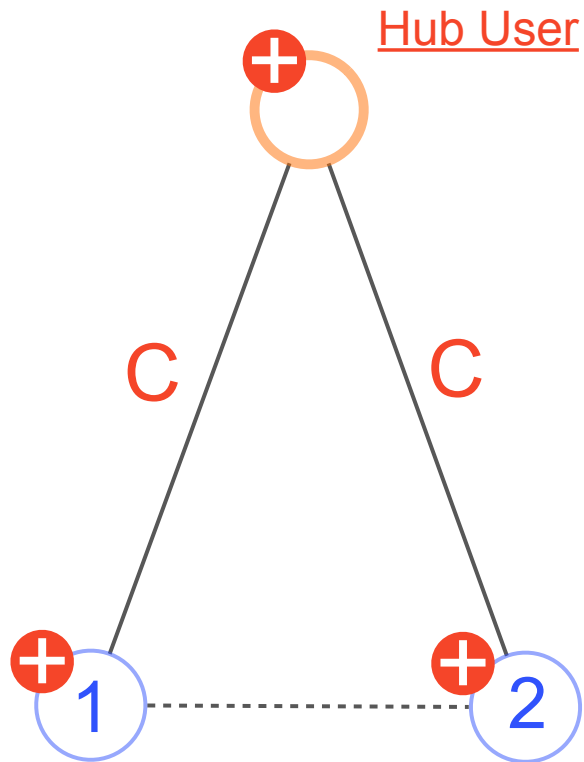
$$d_{a,b} = \frac{C - D}{N - T}$$

$$d_{1,2} = \frac{3 - 1}{5 - 1} = 0.5$$

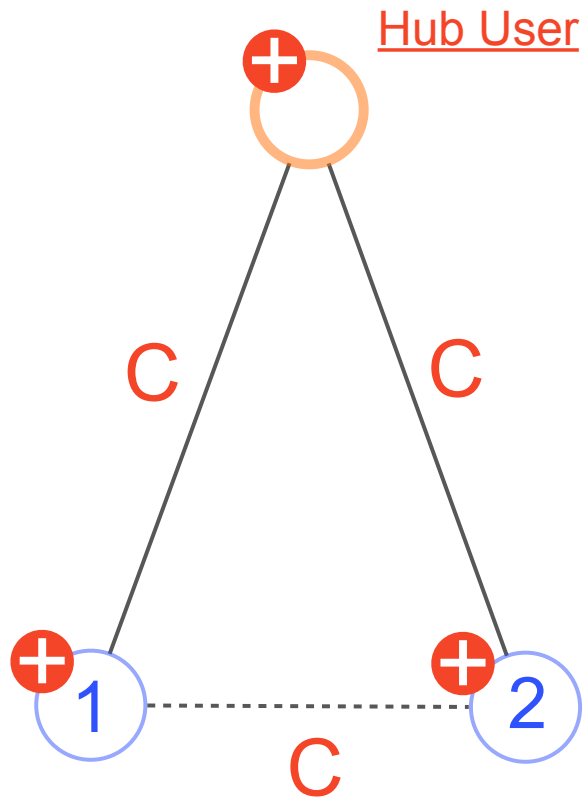
一致度の性質



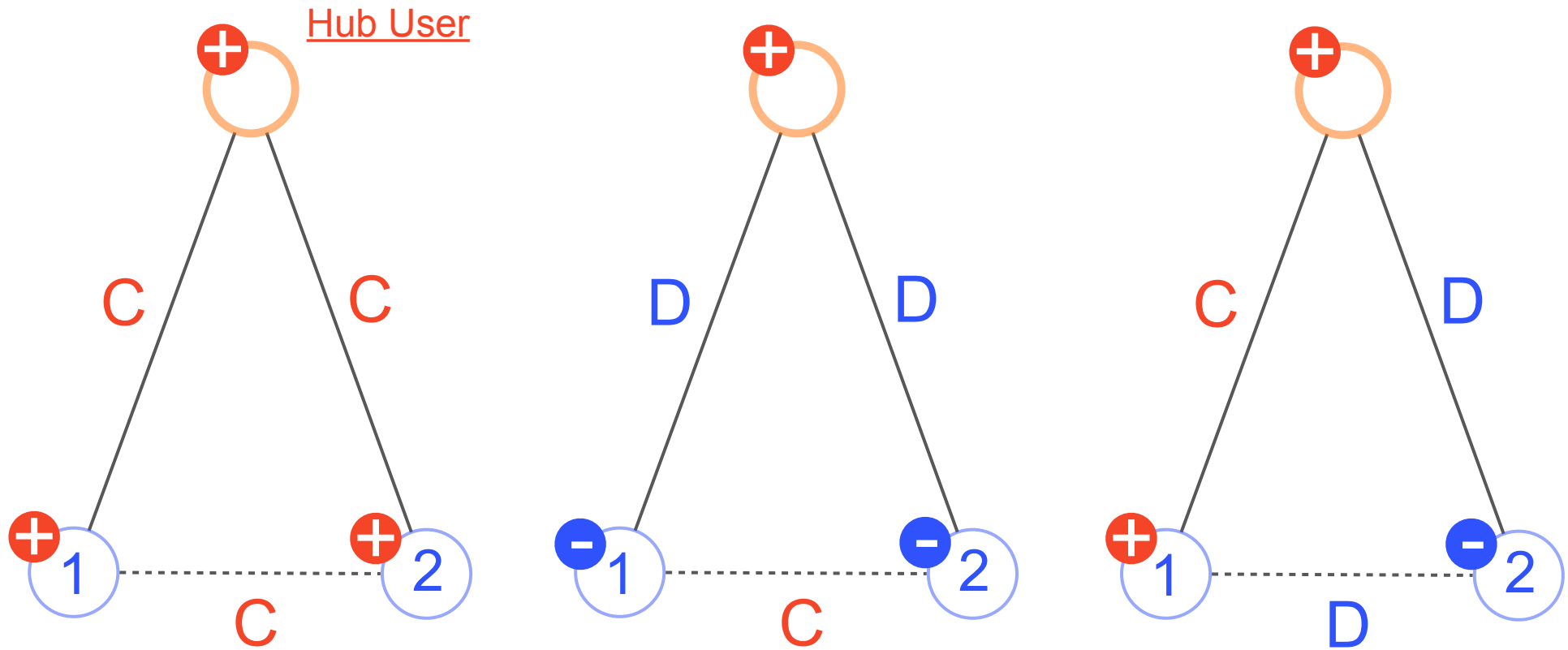
一致度の性質



一致度の性質



一致度の性質



一致度の推定 (Concordantの推定) [Lathia07]

最悪のケース	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	
$U_{1,hub}$	c	c	c	c	d	d	d	t	下限 $\underline{C}_{12} = 0$
$U_{2,hub}$	d	d	d	d	c	c	t	t	
$U_{1,2}$	d	d	d	d	d	d	t	t	

一致度の推定 (Concordantの推定) [Lathia07]

最悪のケース	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8
$U_{1,hub}$	c	c	c	c	d	d	d	t
$U_{2,hub}$	d	d	d	d	c	c	t	t
$U_{1,2}$	d	d	d	d	d	d	t	t

下限
 $\underline{C}_{12} = 0$

最良のケース	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8
$U_{1,hub}$	c	c	c	c	d	d	d	t
$U_{2,hub}$	c	c	t	t	d	d	d	d
$U_{1,2}$	c	c	t	t	c	c	c	t

上限
 $\overline{C}_{12} = 5$

一致度の推定 (Concordantの推定) [Lathia07]

最悪のケース	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	
$U_{1,hub}$	c	c	c	c	d	d	d	t	下限 $\underline{C}_{12} = 0$
$U_{2,hub}$	d	d	d	d	c	c	t	t	
$U_{1,2}$	d	d	d	d	d	d	t	t	
									— 推定値
最良のケース	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	
$U_{1,hub}$	c	c	c	c	d	d	d	t	上限 $\overline{C}_{12} = 5$
$U_{2,hub}$	c	c	t	t	d	d	d	d	
$U_{1,2}$	c	c	t	t	c	c	c	t	

一致度の推定 (Tied) [Lathia07]

最悪のケース	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8
$U_{1,hub}$	t	t	t	t	c	c	d	d
$U_{2,hub}$	t	t	t	t	t	c	d	d
$U_{1,2}$	t	t	t	t	t	c	c	c

下限
 $\underline{T}_{12} = 5$

最良のケース	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8
$U_{1,hub}$	t	t	t	t	c	c	d	d
$U_{2,hub}$	d	d	c	t	t	t	t	t
$U_{1,2}$	t	t	t	t	t	t	t	t

上限
 $\overline{T}_{12} = 8$

$$\max(T_{uc}, T_{vc}) \leq T_{uv} \leq \min(T_{uc} + T_{vc}, N)$$

一致度の推定式

- 一致 (**C**oncordant)

$$\begin{aligned} & \max(C_{uc} + C_{vc} - N, 0) + \max(D_{uc} + D_{vc} - N, 0) \\ & \leq C_{uv} \leq \min(C_{uc}, C_{vc}) + \min(D_{uc}, D_{vc}) \end{aligned}$$

- 不一致 (**D**iscordant)

$$\max(N - (\overline{C_{uv}} + \overline{T_{uv}}), 0) \leq D_{uv} \leq N - (\underline{C_{uv}} + \underline{T_{uv}})$$

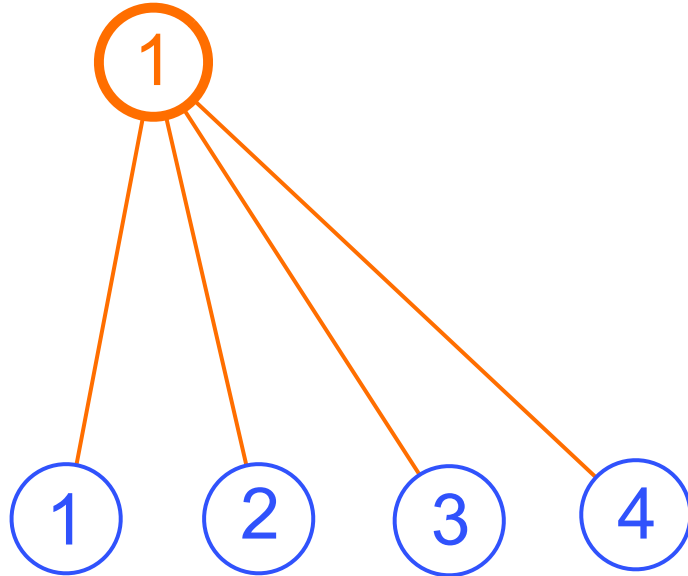
- 同順位 (**T**ied)

$$\max(T_{uc}, T_{vc}) \leq T_{uv} \leq \min(T_{uc} + T_{vc}, N)$$

提案方式：区間の集約方式

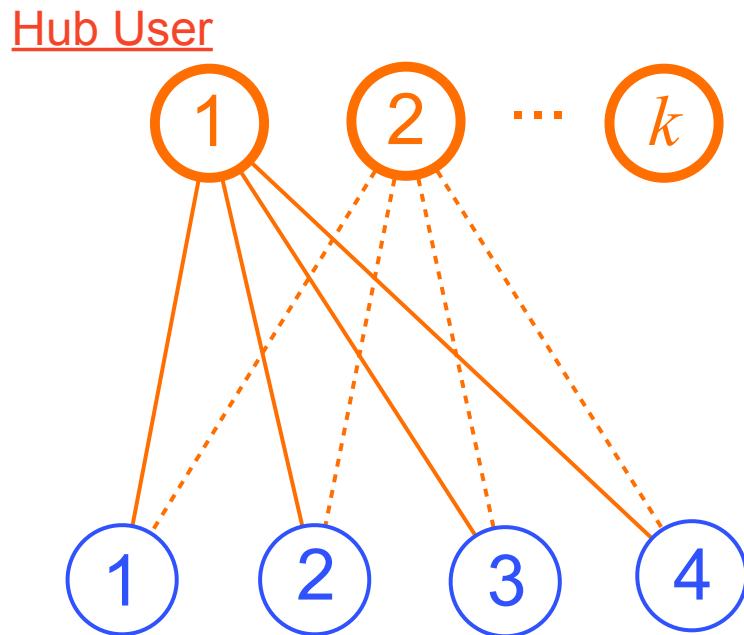
Hubユーザ c_j を独立に k 人作成

Hub User



提案方式：区間の集約方式

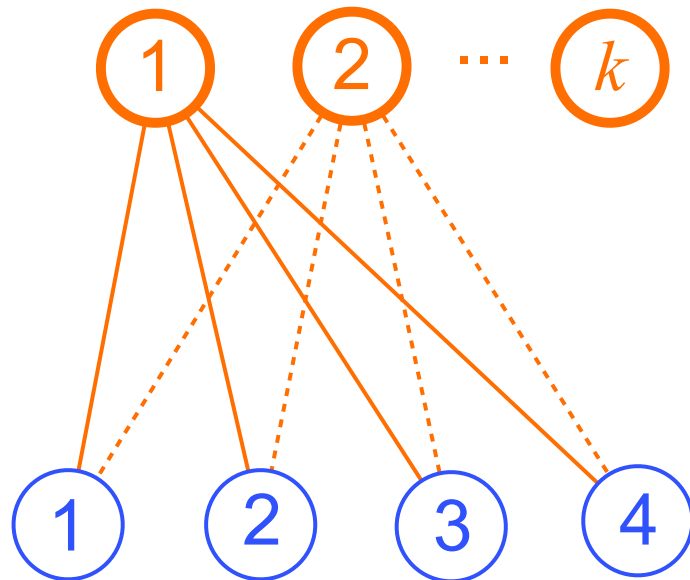
Hubユーザ c_j を独立に k 人作成



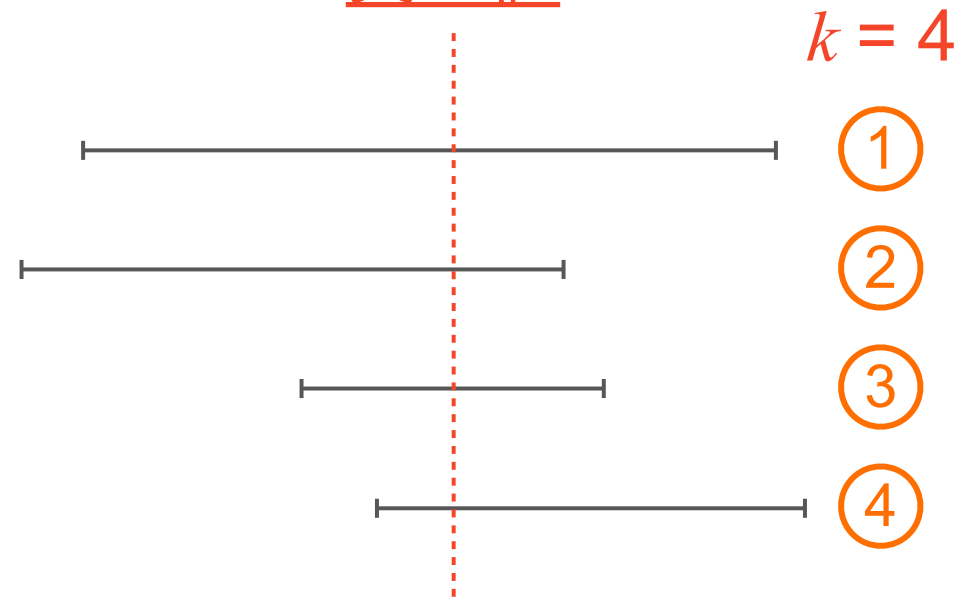
提案方式：区間の集約方式

Hubユーザ c_j を独立に k 人作成

Hub User

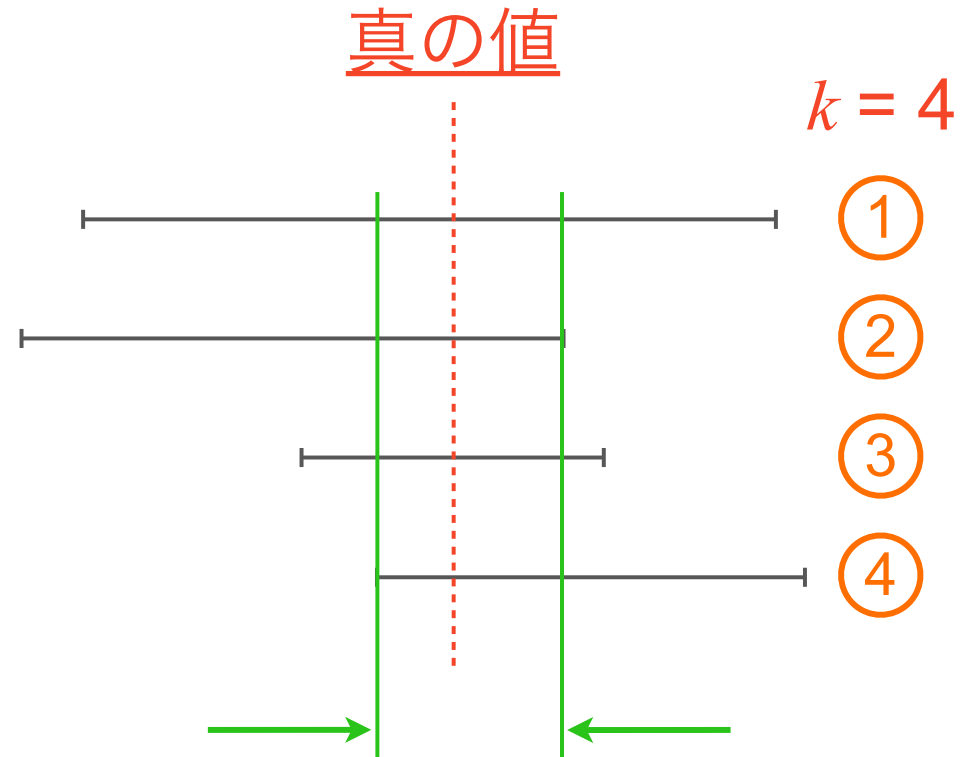
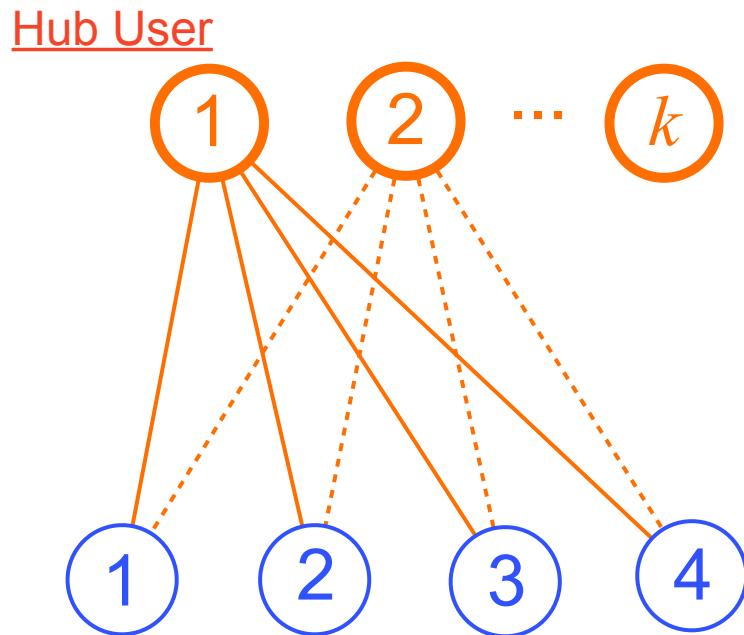


真の値



提案方式：区間の集約方式

Hubユーザ c_j を独立に k 人作成



提案通信プロトコル (Algorithm 1)

Step1

	組織A				組織B				組織C	
	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9	i_{10}
U_1	1	2		4	3	4	1	2	3	2
U_2	3	5	1		2		3	4	2	3
C_1	1	2	4	2	4	5	3	2	1	2

提案通信プロトコル (Algorithm 1)

Step1

	組織A				組織B				組織C	
	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9	i_{10}
U_1	1	2		4	3	4	1	2	3	2
U_2	3	5	1		2		3	4	2	3
C_1	1	2	4	2	4	5	3	2	1	2

Step2

	C	D	T	N	C	D	T	N	C	D	T	N
$U_1:C_1$	2	1	1	4	4	0	0	4	0	2	0	2
$U_2:C_1$	0	2	2	4	0	2	2	4	2	0	0	2

提案通信プロトコル (Algorithm 1)

Step1

	組織A				組織B				組織C	
	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9	i_{10}
U_1	1	2		4	3	4	1	2	3	2
U_2	3	5	1		2		3	4	2	3
C_1	1	2	4	2	4	5	3	2	1	2

Step2

	C	D	T	N	C	D	T	N	C	D	T	N
$U_1:C_1$	2	1	1	4	4	0	0	4	0	2	0	2
$U_2:C_1$	0	2	2	4	0	2	2	4	2	0	0	2

Step3

	C	D	T	N
$U_1:C_1$	6	3	1	10
$U_2:C_1$	2	4	4	10

Alg. 2 秘匿総和プロトコルで共有

提案通信プロトコル (Algorithm 1)

Step1

	組織A				組織B				組織C	
	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9	i_{10}
U_1	1	2		4	3	4	1	2	3	2
U_2	3	5	1		2		3	4	2	3
C_1	1	2	4	2	4	5	3	2	1	2

Step2

	C	D	T	N	C	D	T	N	C	D	T	N
$U_1:C_1$	2	1	1	4	4	0	0	4	0	2	0	2
$U_2:C_1$	0	2	2	4	0	2	2	4	2	0	0	2

Step3

	C	D	T	N
$U_1:C_1$	6	3	1	10
$U_2:C_1$	2	4	4	10

Alg. 2 秘匿総和プロトコルで共有

Step4

$$[\underline{C}_{12,c_1}, \overline{C}_{12,c_1}] = [0, 5]$$

$$[\underline{T}_{12,c_1}, \overline{T}_{12,c_1}] = [4, 5]$$

$$[\underline{D}_{12,c_1}, \overline{D}_{12,c_1}] = [0, 6]$$

区間推定, 出力

実験

目的

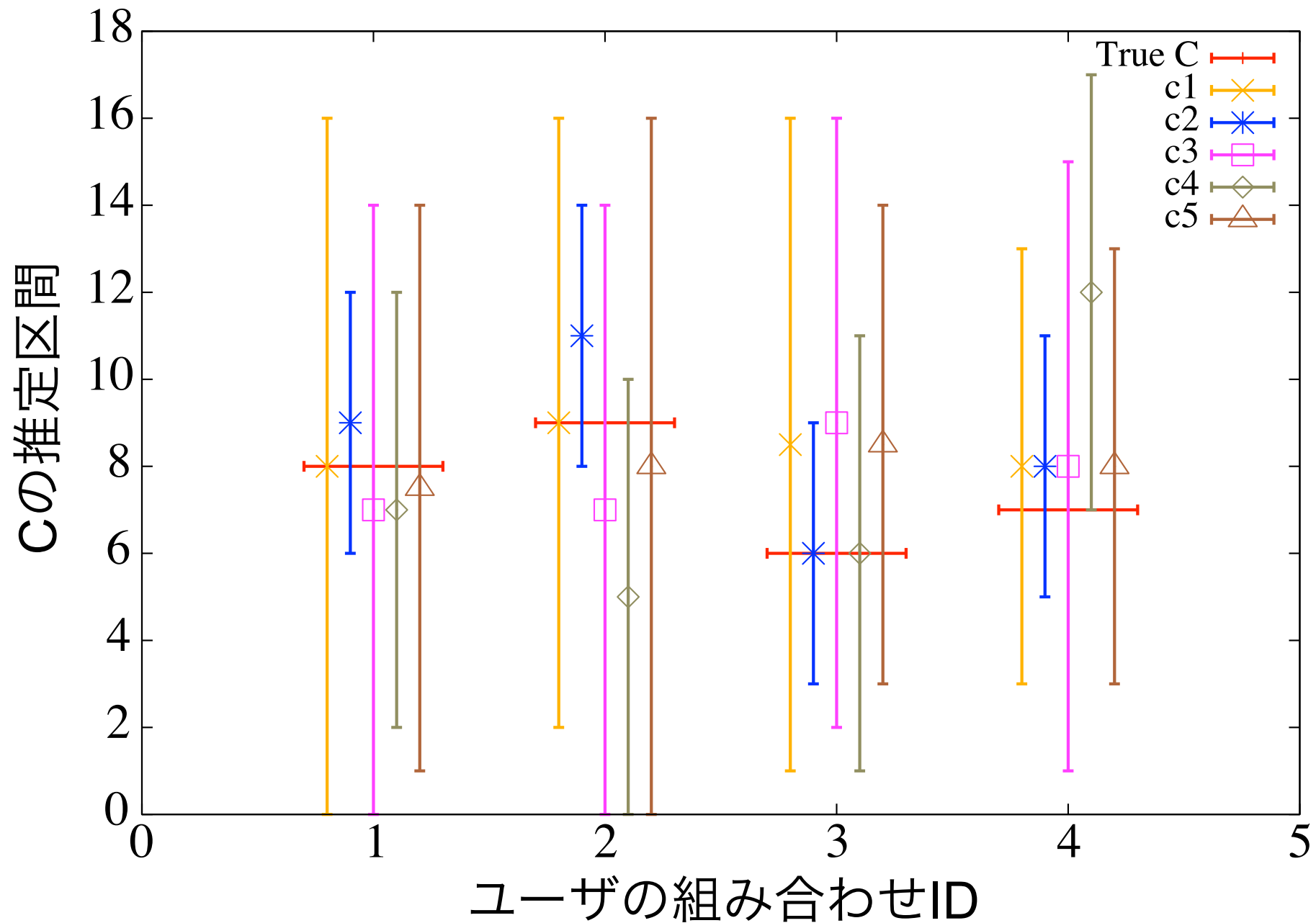
- Hubユーザ数 k の増加に伴う区間の変化
- 予測した類似度と真の類似度の差

方法

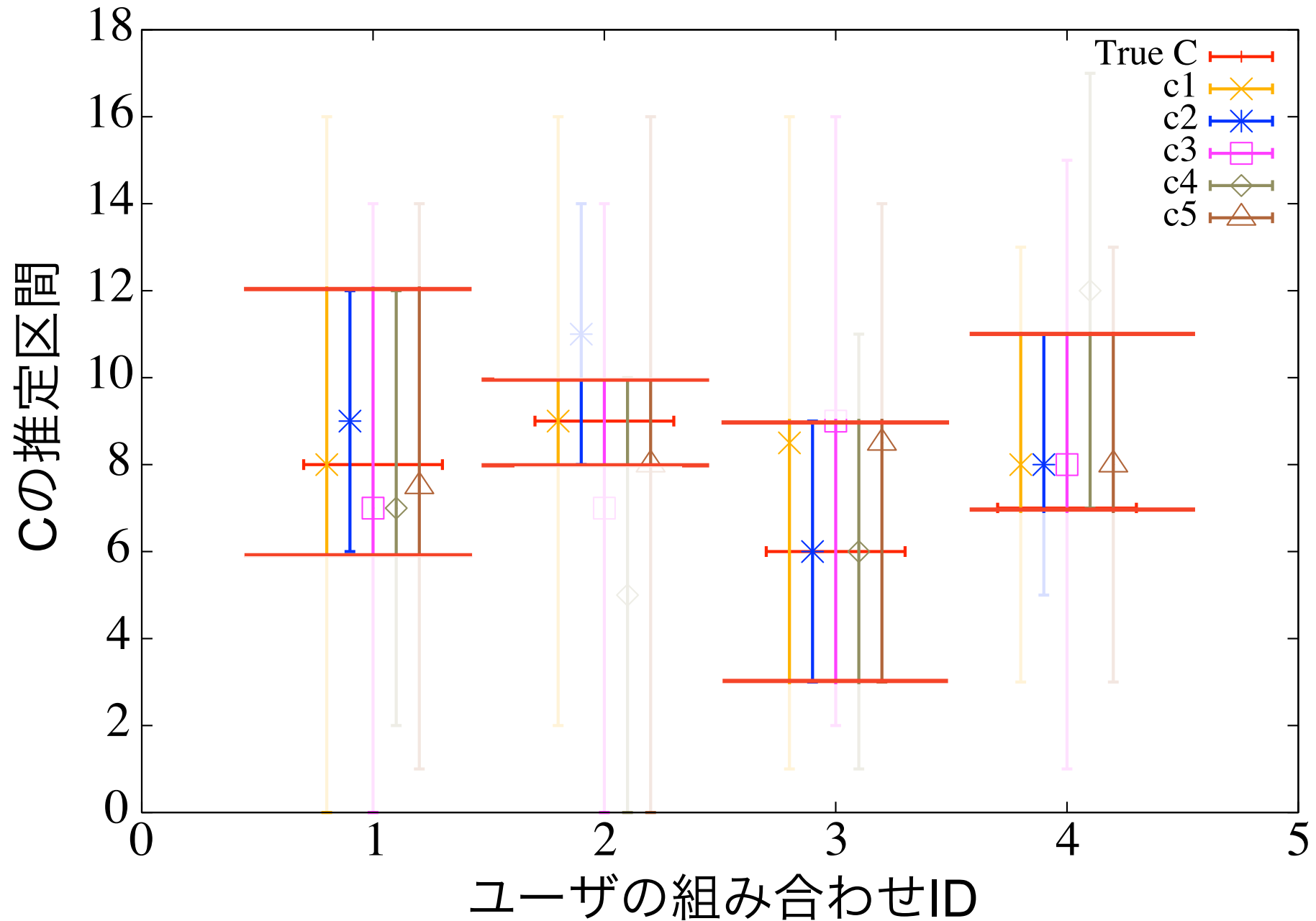
ユーザ数 6, アイテム数 17, 評価数 102, のランダムデータセットを使い提案方式で区間を推定.

- 結果1：擬似ユーザによる推定区間の変化
- 結果2：擬似ユーザ数 k の変化による区間の変化
- 結果3：真の類似度と推定した類似度の分布

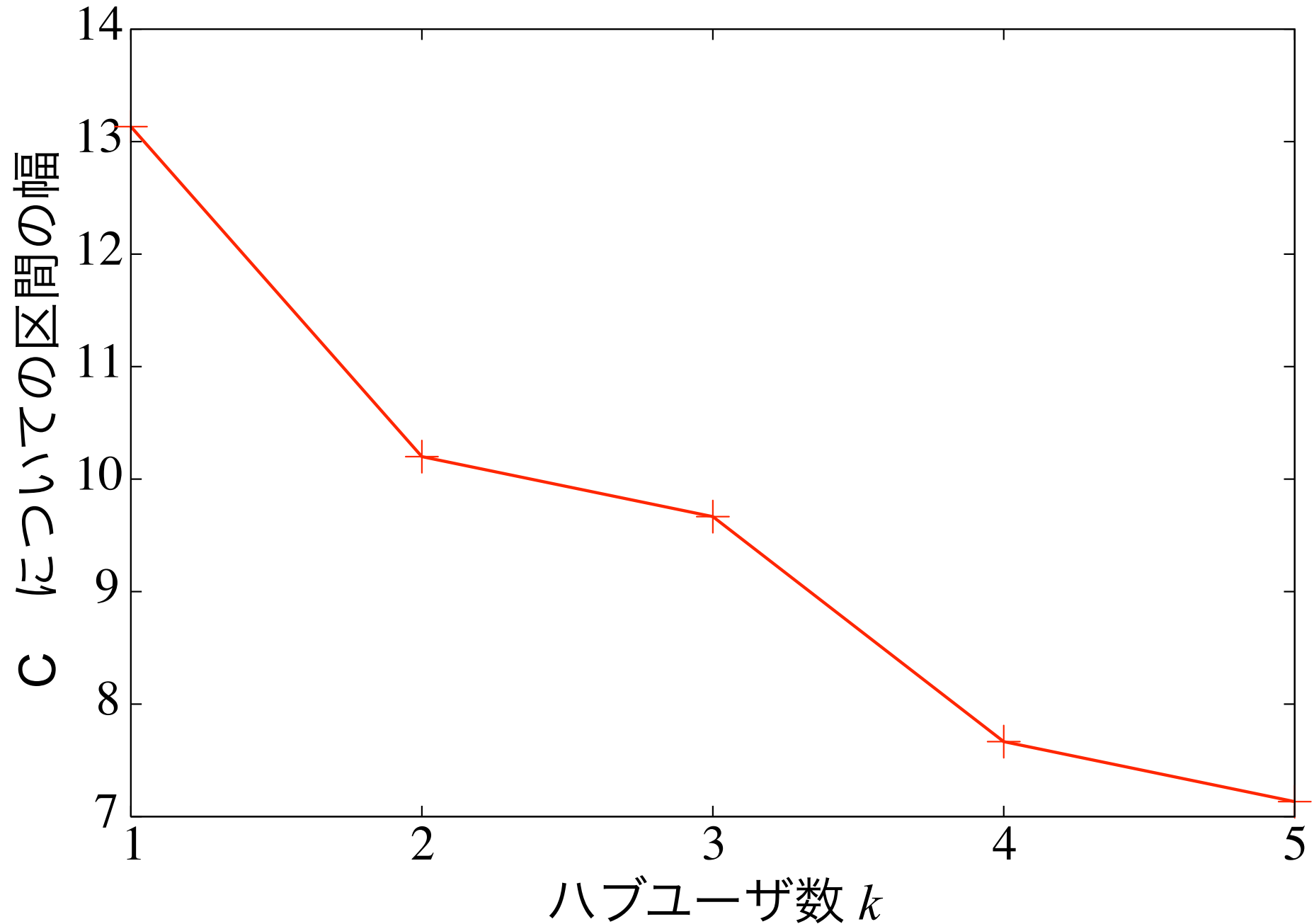
結果1：擬似ユーザによる推定区間の変化



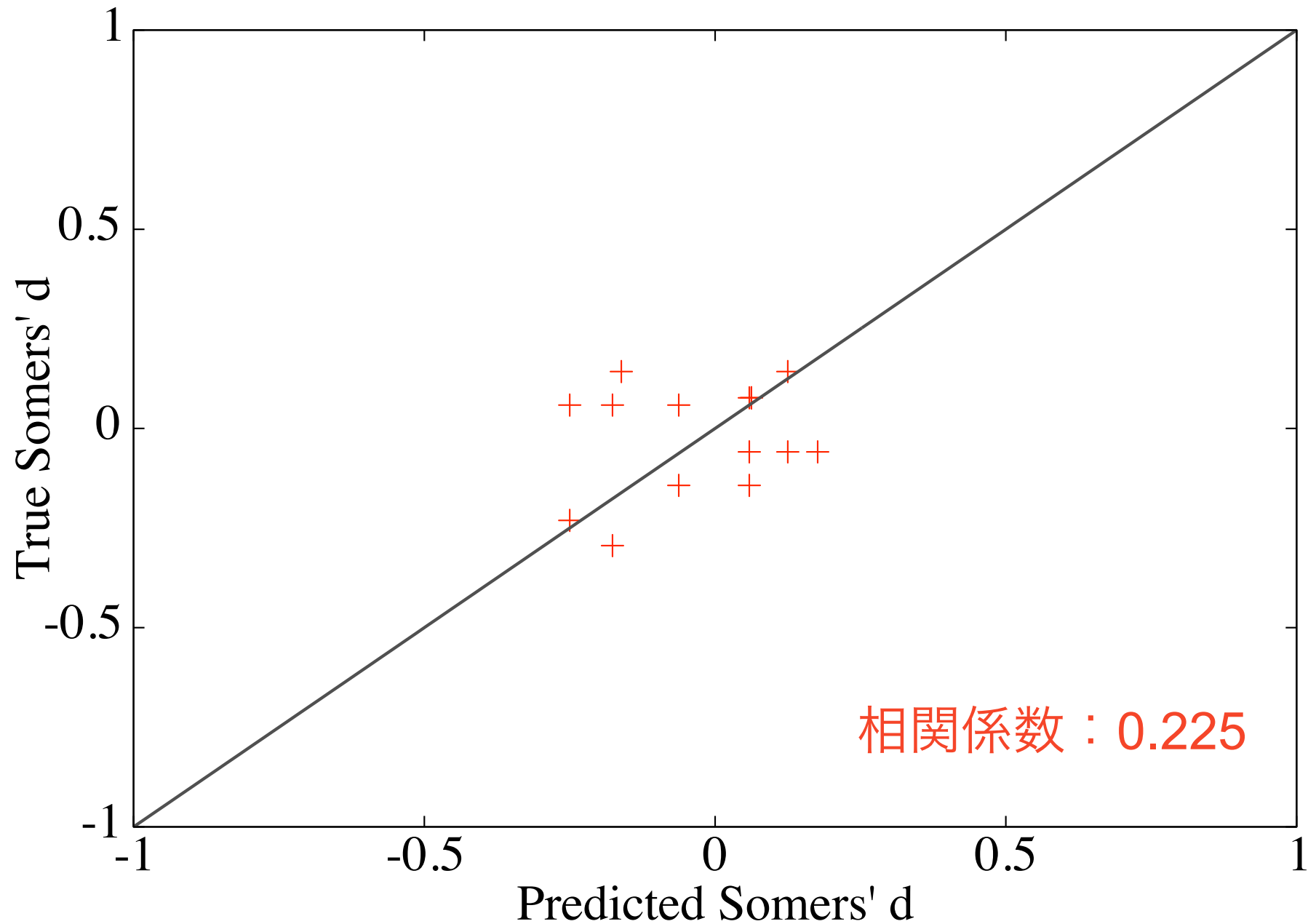
結果1：擬似ユーザによる推定区間の変化



結果2： 擬似ユーザ数 k の変化による区間の変化



結果3： 真の類似度と推定した類似度の分布



まとめ

- 3 組織間, 垂直分割の環境における, ユーザ間類似度の計算方法を提案
- 通信コストを n^2 から kn に改善
- 複数のハブユーザを使い推定区間を集約