

公的統計マイクロデータ二次的利用推進の ための取り組み

情報・システム研究機構 統計数理研究所

山下智志

CONTENTS

1. 取り組みの背景
 2. オンサイト拠点の仕組みと運営
 3. 利用活性化への取り組み：
公的マイクロデータ研究コンソーシアム
 4. 今後の展開
- (付属資料) データ構造化の手法

取り組みの背景(1)

- 我が国では主要先進国に比べて、社会・経済に関する実証研究が非常に遅れており、状況の改善が急務である
 - 研究や高等教育の場で利用可能なマイクロデータが限られている
 - 研究者にマイクロデータ実証分析に関する経験や研究力量が決定的に不足している
- 人文社会科学分野の国際的なジャーナルでは、マイクロデータを用いた実証分析が主流
 - 実証分析なしに論文が掲載されることはほぼ不可能
 - 我が国の研究者は、海外のマイクロデータに依存するなど著しく不利な状況にある
- 欧米、オーストラリア、韓国などでは政府の公的統計のマイクロデータの研究の利用環境が整えられ、研究の知見に基づいてEvidence-Based Policy Makingが実践されている

取り組みの背景(2)

これらの課題の解決には、国の公的統計調査のマイクロデータを、個人の秘密を厳守しつつ、実証分析を有効に活用する体制の確立が必要

- 我が国研究者の研究力量向上
- 人文社会科学分野の国際競争力の向上
- 国民生活の向上
- 社会経済の発展に資する政策科学研究の促進
- エビデンスに基づく科学的な施策の立案・評価の推進

平成19年の統計法の改正に伴い、**公的統計のマイクロデータ(個票)の二次的利用**(研究・高等教育利用)が認められることとなり、我が国でも、集計データではなく、大規模マイクロデータに基づく人間・社会政策科学研究を推進する機運が高まっている

平成23年からは厚生労働省のナショナルレセプトデータベースの研究利用も開始された

データ基盤整備、ならびに、政府情報の研究者利用の推進に資する環境を整備

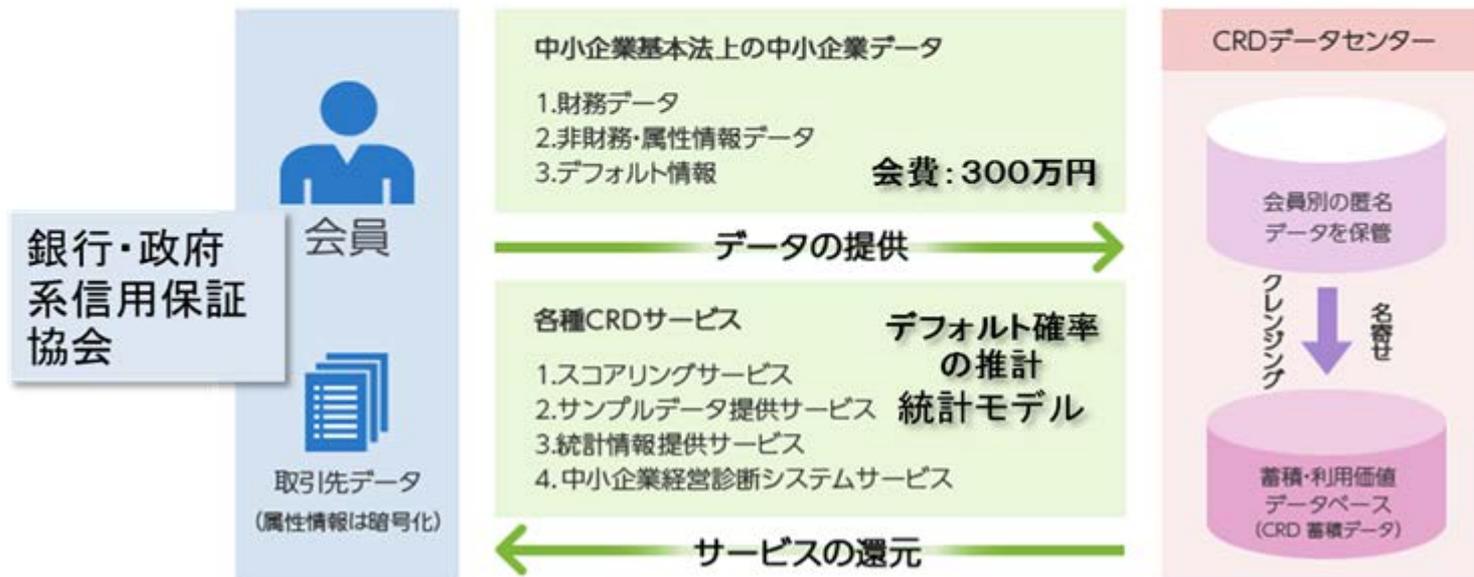
取り組みの背景(3)

情報・システム研究機構 統計数理研究所における 秘匿性マイクロデータ相互利用プラットフォーム開発事例



CRDの仕組み

1. 会員（信用保証協会及び金融機関）は、取引先中小企業の財務データ・非財務データ・デフォルトデータをCRD協会に対して定期的に提供します。
（企業名は全て暗号化され、個別企業名の特定はできない仕組みとなっています。）
2. CRD協会から会員に対しては、蓄積されたデータを加工して各種サービスを還元します。



取り組みの背景(4)

平成19年統計法改正により、公的統計データの二次的利用(研究・高等教育利用)開放

新しい条項の追加

- 第34条 委託による統計の作成等(オーダーメイド集計)
- 第36条 匿名データの提供
 - 学術研究の発展に資すると認める場合
 - 高等教育の発展に資すると認められる場合

独立行政法人統計センターによるオーダーメイド集計、匿名データ提供の開始

平成22年7月、機構－統計センターと 連携協力協定締結

- 統計数理研究所内に「新領域融合研究センター
統計数理研究所・オンサイト解析室」(当時)を設置
- 同年9月よりサテライト機関としての業務を開始

調印式で握手を交わす堀田凱樹 情報・システム研究機構 機構長(当時)と
戸谷好秀 統計センター 理事長(当時)



CONTENTS

1. 取り組みの背景

2. オンサイト拠点の仕組みと運営

3. 利用活性化への取り組み：
公的マイクロデータ研究コンソーシアム

4. 今後の展開

(付属資料) データ構造化の手法

オンサイト施設の整備・運用

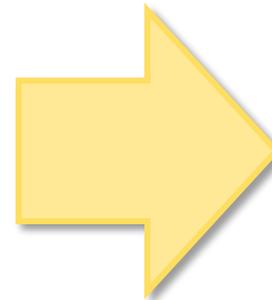
統計法第33条に基づく調査票情報の提供(目的外利用)を受けるためには、一定のセキュリティ要件を満たす必要がある

- データ分析を実施するエリアの入退管理
- 施錠可能なキャビネット等にデータを保管
- 外部ネットワークへの接続禁止



全ての研究者が要件を整えられるわけではない

- 共同研究室を割り当てられている若手研究者
- 施設の都合上、セキュリティを強化することができない組織に属する研究者 など



研究者のデータアクセシビリティを高めるためには、セキュリティ要件を満たす施設の供給が不可欠

→条件を満たすことができない研究者にとって
データ利用申請は困難

オンサイト施設の整備・運用

情報セキュリティが確保された施設と使用を厳重に管理することが可能な仕組みを整備し、提供することにより、33条の適切かつ円滑な運用、調査票情報の適正管理、行政機関等の事務の負担軽減及び手続の迅速化に寄与

高度なセキュリティ環境が整った分析拠点を研究者に提供し
データ利用環境整備の負担を軽減

統計センターとの関係協力協定に基づき、情報・システム研究機構でも統計数理研究所内にオンサイト施設を設置し整備と運用を開始

- 情報・システム研究機構 データサイエンス共同利用基盤施設
社会データ構造化センター オンサイト解析室

オンサイト利用に係る施設基準(試行運用版)について(概要)

1 運営・管理体制について

- ・オンサイト施設管理責任者や、オンサイト施設管理運用窓口を置いていること など

2 オンサイト施設について

- ・施設管理責任者が認めた者や施設利用者以外の立入りを制限し、機密情報を安全に利用するための措置を講じた施設であること
- ・オンサイト施設への入退室に際し、氏名、所属、日付、時刻の記録を行う措置が講じられていること
- ・施設利用者に対し、施設利用許可証を貸与し、施設利用時において常に見やすいところに着用させる措置が講じられていること
- ・利用者端末が複数ある場合、それぞれをパーティションで区切る等の措置が講じられていること
- ・統計センターが指定する回線(SINET)により、利用者端末から中央データ管理施設へのアクセス可能なネットワークが構築されていること
- ・利用者端末及び機器等について、定期的にメンテナンスを行い、正常な状態を維持する体制が整えられていること
- ・オンサイト施設内に、インターネットにアクセス可能なパーソナルコンピュータの設置を行う場合は、大学のセキュリティポリシー等に準拠したものであれば設置可能とし、利用規約等を定め、利用者端末から離れた場所に設置されていること など

3 利用者端末について

- ・利用者端末の設置は、3台までとすること。なお、利用者端末を増設したい場合は、統計センターと協議すること
- ・リモートアクセスのために使用する起動用USBメモリの盗難や紛失を防止するため、管理を徹底すること
- ・利用者端末の盗難、外部への持ち出しを防止する措置を講じること

オンサイト利用に係る施設基準(試行運用版)について(概要)(続き)

4 施設利用者への規制及び監視措置について

- ・パーソナルコンピュータ、カメラ、レコーダ等の記録機器類、無線LAN端末、携帯電話等の通信機器類、その他これらに類する機器について、施設利用者に対し、オンサイト施設内への持込み禁止等の措置が講じられていること
- ・監視カメラ及びカメラ映像を記録する機器を設置するなど、オンサイト施設内における施設利用者の行動を監視できるようにすること など

5 報告・検査措置について

- ・施設管理責任者等の名簿、施設内の利用者端末、その他の機器等の構成及び配置について、統計センターに提出すること
- ・定期的にオンサイト施設の利用状況を統計センターへ報告すること。また、統計センターが施設の利用状況について報告を求めた場合は、速やかに求めに応じること
- ・施設管理者は、統計センターの求めに応じ、オンサイト施設の検査を受けられる体制をとること

6 その他

- ・上記で定めることのほか、オンサイト施設におけるセキュリティの確保に努めること
- ・オンサイト利用に障害等が発生した場合には、その対応を行うこと
- ・調査票情報の漏洩等の事故や不適切な利用等の事象が生じたとき又はその疑いが生じたときは、統計センターと共同で原因究明等に当たること
- ・施設の利用時間は、統計センターの運用体制を踏まえ、施設管理者が決定すること。決定した利用時間は、統計センターへ報告すること

統計数理研究所 オンサイト解析室の概要

強固なセキュリティにより、データを保護

- 管理側でデータの持ち出しをコントロール
 - PCへの外部メディア接続の全面制限
 - データ保存用端末による集中管理
 - 外部ネットワークからの遮断によるネットワーク経由のデータ流出防止
- 監視カメラによる常時監視
- パーティションにより区切られた作業環境
 - 意図するor意図しない窃視の防止

統計法第33条による公的統計の目的外利用におけるセキュリティ要件を満たした環境

- 統計センターによる検査を受け、認可

最終生成物データを持ち出す際には、解析室、データ提供者による内容チェックを受ける必要あり。原則として即日持ち出しは出来ない。

本年8月、立川新棟に移転・拡充



統計数理研究所 オンサイト解析室の利用可能システム

解析用端末:2台

- OS: Microsoft Windows7 Professional SP1 64bit
- CPU: Intel Xeon E5-2630 v3 (2.4GHz 8Core)
- GPU: NVIDIA GeForce GTX980 (GDDR5 4GB)
- メモリー: DDR4-SDRAM 32GB (8GB × 4) (ECC/PC-17000)
- ドライブ: SSD SATA 6GB/s 512GB × 2 (RAID1)

光学ドライブ無し、USBポート使用不可

- データは、室内ネットワークを通じてデータ保存用NASにアクセスし、利用する。
- **最終生成物データを持ち出す際には、解析室、データ提供者による内容チェックを受ける必要あり。原則として即日持ち出しは出来ない。**

使用可能なソフトウェア

- IBM SPSS Statistics 24 (日本語版)
Statistical Base, Advanced Statistics, Regression, Custom Tables, Missing Values, Categories, Amos
- R (日本語版) + RStudio
- Microsoft Office 2016 Professional (64bit) (日本語版)
- Adobe Acrobat DC Pro (日本語版)

追加可能なソフトウェア

- SAS 9.4 (32bit/64bit)

その他のソフトウェア利用については応相談

オンサイト利用施設整備の意義

高度なセキュリティ環境・分析用端末を自前で整える必要がなく、
データ利用者の負担を軽減

あらかじめセキュリティ要件が満たされている環境を利用してもらう
ことで、セキュリティ審査負担・利用者管理負担、漏えいリスクを軽減



これまで環境を整えられなかった研究者が
データ利用の機会を得られる可能性を提供



公的統計データの二次的利用を促進

CONTENTS

1. 取り組みの背景

2. オンサイト拠点の仕組みと運営

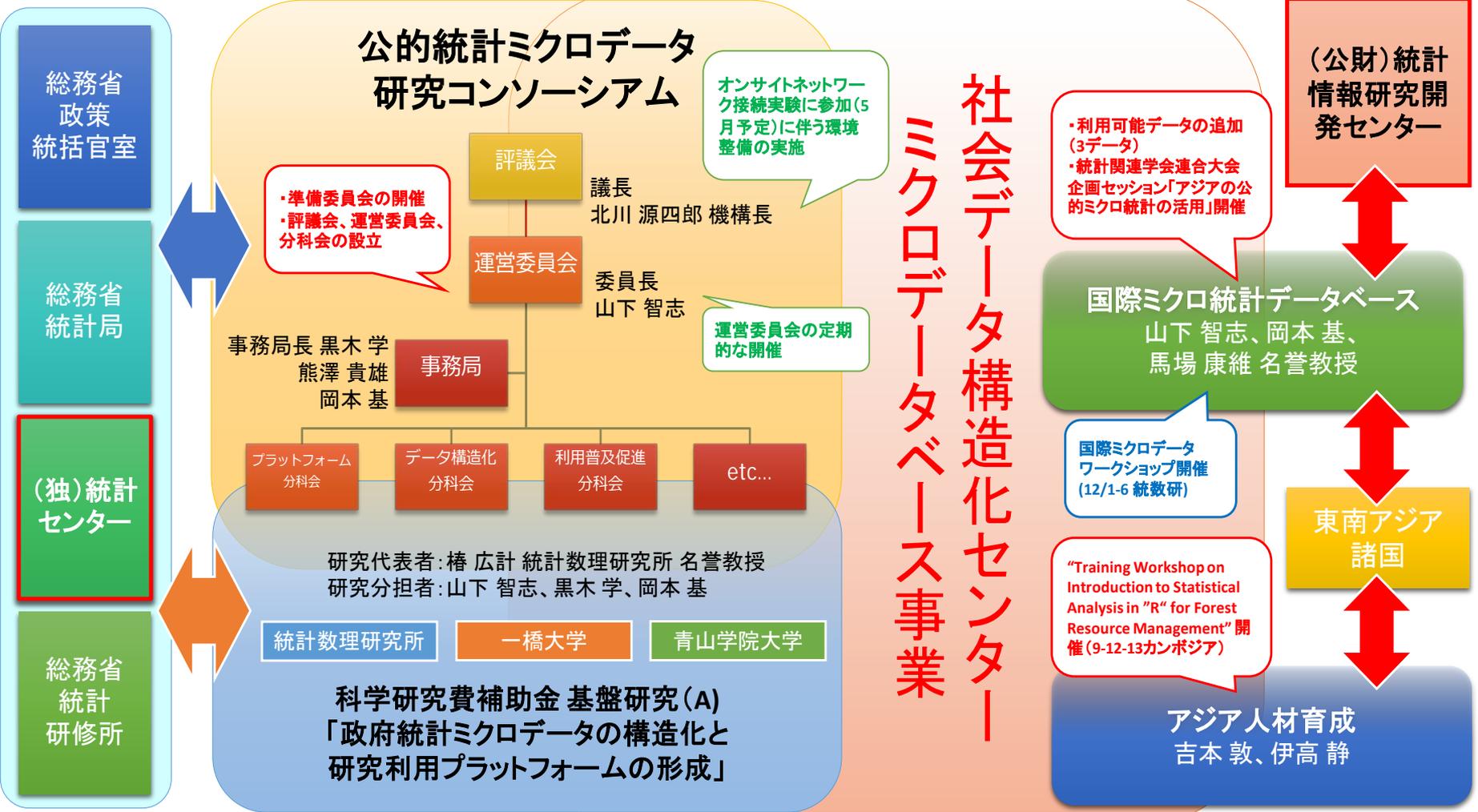
3. 利用活性化への取り組み：
公的マイクロデータ研究コンソーシアム

4. 今後の展開

(付属資料) データ構造化の手法

マイクロデータ事業の運営体制 普及活動と必要技術の開発

※赤太枠: 機構と関係協力協定締結



公的統計マイクロデータ研究コンソーシアムの設立

公的統計の利活用を推進めるために、全国の研究者に向けて、利用しやすい環境を整える必要性

- 必要なセキュリティを確保した「オンサイト施設」を全国的に整備し、潜在的な研究の裾野を広げる

学官産が協力して、公的統計の二次的利用、オンサイトネットワーク整備の支援と意義を広く周知し、諸課題の検討と解決にあたる「公的統計マイクロデータ研究コンソーシアム」を設立

事務局を社会データ構造化センターに設置

平成28年3月29日、

一橋講堂にて設立記念シンポジウムを開催

平成28年8月28日、第1回評議会開催

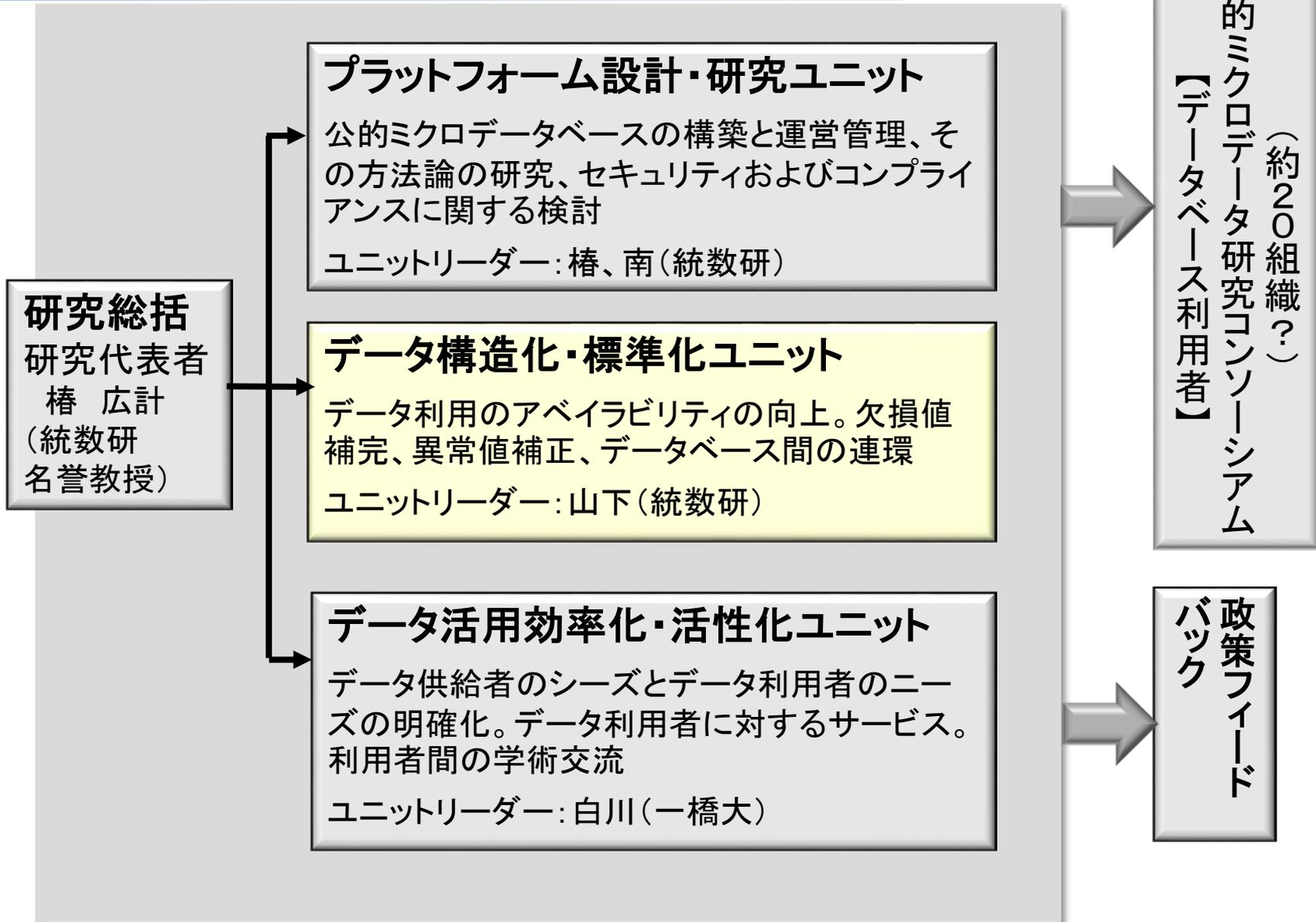
...

本年9月より、会員募集開始(予定)



設立記念シンポジウムの様子

データ活用のための研究体制



データ構造化の目的と手法

データの供給を受けた研究者がストレスなく研究を開始できるよう、データベースを供給段階で整えておく必要がある

・欠損値補間

経済統計、医療統計を中心に研究成果が多い。
経済(Hot Deck, Cold Deck)と医療(ICE, Knn, MICE)とは系統がやや異なる
単純は「削除」「平均値補間」は少なくなりつつある。

・異常値補正

研究成果は少ない。
3シグマ折り返し処理、数値情報の順位化、数値情報のカテゴリー化などの
経験的処理が一般的である。
一般的にモデルパラメータに与える影響が大きい

・リレーション、データ・リンケージ

いわゆる名寄せ。
情報が複数のデータベースに分散している場合、必須の処理
特定フィールドの完全一致から、確率一致性へ進化

・テキストデータの数量化

CONTENTS

1. 取り組みの背景
2. オンサイト拠点の仕組みと運営
3. 利用活性化への取り組み:
公的マイクロデータ研究コンソーシアム

4. 今後の展開

(付属資料) データ構造化の手法

今後の展開(1)

官学データ利用(政府データの学術利用)

- ・オンサイト拠点の全国展開によるデータ・アクセシビリティの向上
→オンサイト拠点設置・運営のための費用負担
- ・公的マイクロデータ研究コンソーシアム活動による広報活動と利用者ニーズの把握
- ・データ構造化によるクレンジングされたデータの供給、データ・リライアビリティの向上
→データ構造化研究の推進
- ・データ利用のための人材育成、情報提供

学学データ利用(研究活動データの学術相互利用)

- ・社会データアーカイブズの開発と利用(東大、慶應大などの実績)
→データ提供者のメリット、利用者フレンドリーなインターフェースなどが課題

民学データ利用(民間データの学術利用)

- ・民間との共同研究などでデータを入手可能
→知財契約、守秘義務契約などのコンプライアンス、セキュリティ体制などなれていない研究者にはハードルが高い

今後の展開(2)

民官データ利用(民間データの政府利用)

→所轄官庁が企業よりデータを取得することは一部行われている(例 銀行法24条)

ただし省庁間で利用する仕組みはこれから?

cf. 銀行法24条のデータ提出先は内閣総理大臣であるが実質的には金融庁が独占利用に近い状態

民民データ利用(民間データの民間利用)

→コンプライアンスが整理できた業界からデータの共有化は進められている。
(データベンダーの存在)

学民データ利用、学官データ利用については学学のアーカイブズを通して可能

→アーカイブズの改良が必須?

・統計改革推進会議

官官データ利用の早期実現、
官学データ利用については早急推進、
官民データ利用については重要課題として積極的に検討

・骨太の方針

官民が保有するデータの徹底した利活用を図るべく、新しい社会インフラとなるデータ利活用基盤を構築する。「官民ラウンドテーブル」等を通じた公共データのオープン化、安心してデータ流通を促進させるための法制度整備等を進める。

CONTENTS

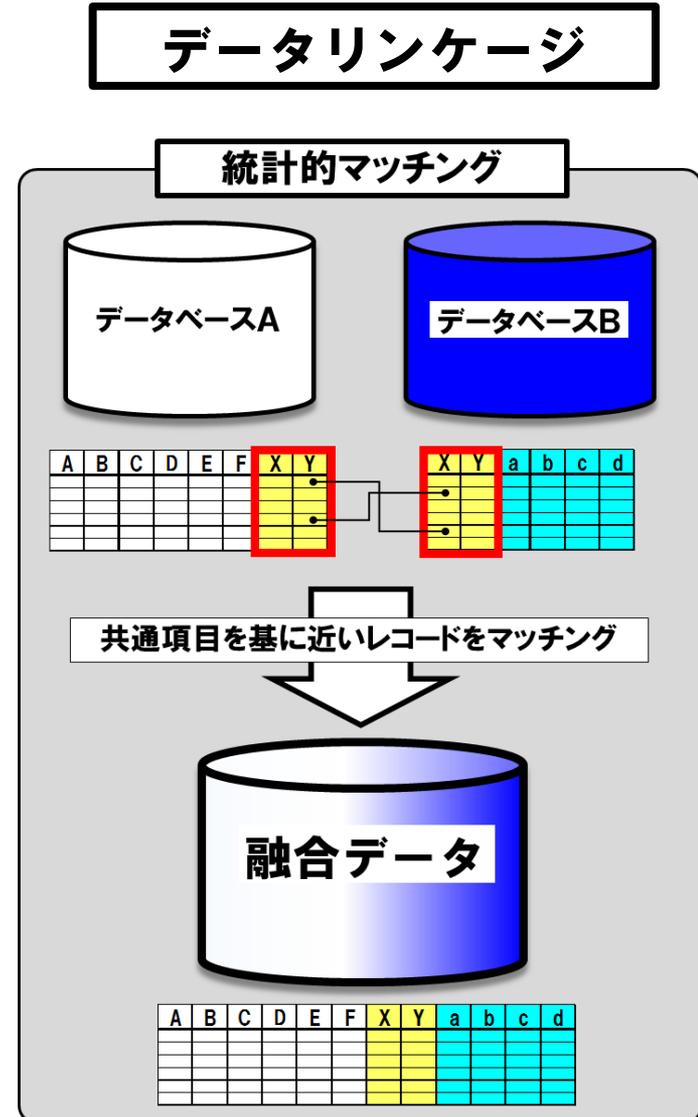
1. 取り組みの背景
2. オンサイト拠点の仕組みと運営
3. 利用活性化への取り組み:
公的マイクロデータ研究コンソーシアム
4. 今後の展開

(付属資料) データ構造化の手法

データリンケージ、リレーションとは？

- 複数のデータベースをレコード単位で結合 (Matching) することにより、豊富な情報を持つ単一のデータベースを構築する方法。
- 新たな調査やデータ収集を行うことなく有益な情報を持つデータベースを構築することが可能。

マーケティング・サイエンスの分野では「データ融合」(Data Fusion) とも呼ばれており、POSデータと消費者アンケートデータを統合することにより、消費者行動に関する詳細な研究・分析が行われている(星野(2009))



特定フィールドの完全一致から、統計的一致性評価へ

これまで

特定フィールド(名称、住所、設立年など)の完全一致

→ 名称・住所表記のぶれなどから、不完全一致の対応アルゴリズムはある
完全一致しないデータは基本的にリレーションしない

これから

類似度を指標化し、ある閾値を超えたものはリレーションを行う

→ 統計的マッチング

0. 統計的マッチングの際の前提条件

(1) 「条件付き独立性」

- ・ $[X, Z]$ と $[Y, Z]$ をマッチングする場合 (共変量は Z)、 X と Y に関する Z の条件付き独立性 (CIA: Conditional Independence Assumption) が成立していることが前提。

(2) 「強く無視できる割り当て」

- ・ どちらの群に割り当てられるかは共変量の値に依存し、従属変数による割り当ての影響はあくまで「共変量と従属変数の関係」を通じてのみ間接的に依存している。

高部勲(2017)より抜粋

傾向スコア（propensity score）マッチング

共変量マッチングはすべてのデータペアの距離を計算する必要があるため、データ量が多いと計算負荷が現実的でない。

→ より簡便的なマッチングを考える（傾向スコア）

傾向スコアの定義は安定していない。

2群判別のロジットZスコアの1次元で定義しているものが多い。

医療の死亡、
金融の倒産など
外的な2値変数が存在

汎用データベースには
用いることが困難

(1) 最近隣法（nearest neighbor matching）

- ・ 傾向スコアでみて、最も距離の近いレコードとマッチングする。

(2) キャリパーマッチング

- ・ ある特定の距離以上になるときは、マッチングしない。
- ・ マッチングできないレコードが生ずる可能性がある。

(3) カーネルマッチング

- ・ 対象群のすべてのレコードの値を、カーネルの重みで利用する。
- ・ 処置群と対象群を1対1でマッチングする際に生じる問題点（マッチングできないレコードが生じる可能性等）に対処するためにヘックマンが提案。

$$\hat{y}_{i0} = \frac{\sum_{j=1}^N (1 - Z_j) K_{ij} y_{j0}}{\sum_{j=1}^N (1 - Z_j) K_{ij}}$$

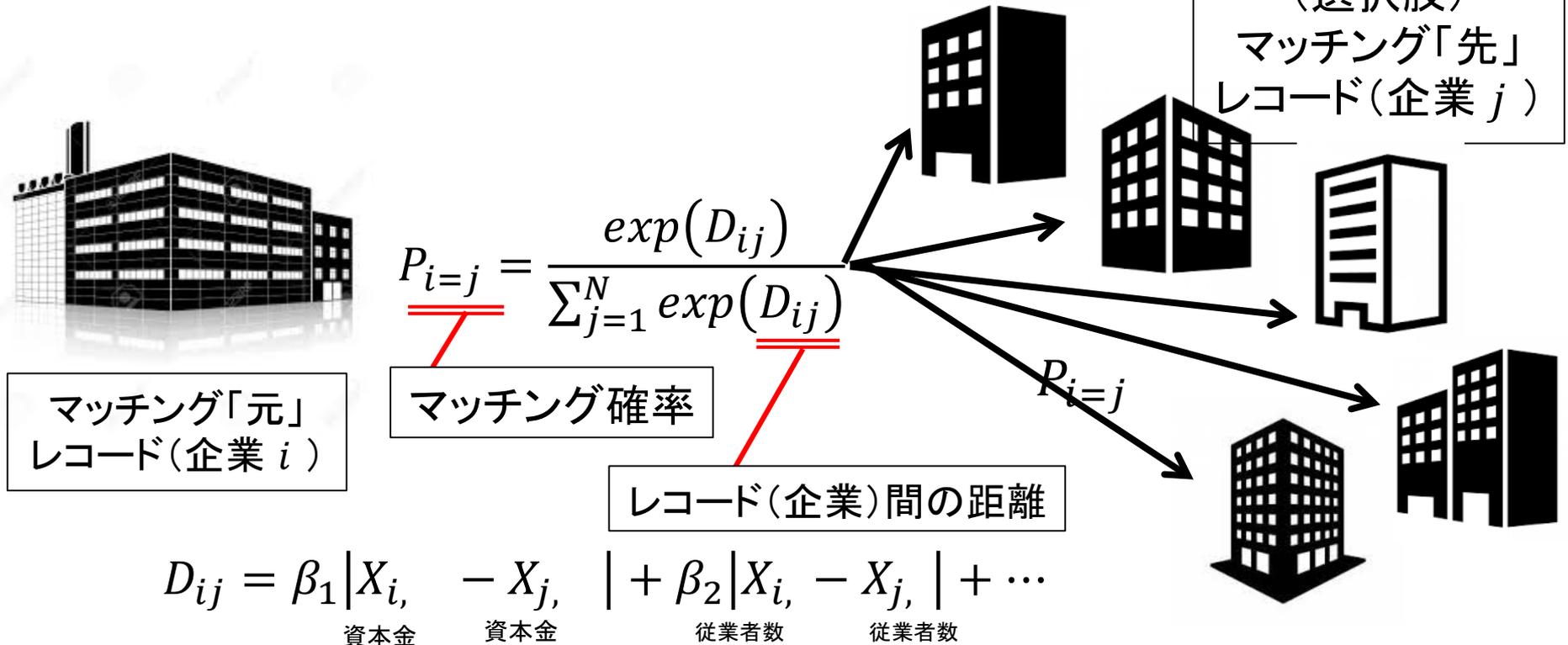
※傾向スコアについては、主に経済学分野で、企業データに対する適用事例がある。

※あるイベントが起こった地域の企業と、起こっていない地域の類似企業をマッチングし、その差をみることで、イベントの影響について分析する事例がほとんど。

多項ロジットモデルに基づく統計的マッチング

- 多項ロジットモデルの枠組みを統計的マッチングの問題に適用]

Cf. マーケティングにおける消費者選択確立推計のモデル



マッチング確率により、レコード間リレーションを作る。

リレーション作成ルールは様々:

例) レコード間で最高確率のペアをマッチング

→ 残ったレコード間で最高確率のペアをマッチング → 繰り返す...

→ ある一定の確率ではマッチングさせない

リレーション実験例

○マッチング元:

「帝国データバンク」データ(平成24年2月分)

商品名: COSMOS II 企業概要ファイル・レイアウトC

産業分類は、「日本標準産業分類」を基に類似のコードを付与。

地域コード(都道府県・市区町村)は政府統計のものを使用。

⇒ 質的データ項目の内容は類似している

○マッチング先:

「平成24年経済センサス - 活動調査」 ミクロデータ(※)

(※統計法第33条による二次的利用の制度に基づき提供を受けたもの。)

民間企業データベースと公的ミクロ企業データベースがくつつくかの実験

距離の計算のイメージ

		帝国データバンク【マッチング元】 i			
		TDB企業 1	TDB企業 2	...	TDB企業 M
経済センサス【マッチング先】 j	EC企業 1	dist(1, 1)	dist(1, 2)		dist(1, M)
	EC企業 2	dist(2, 1)	dist(2, 2)	...	距離(2, M)
	⋮	⋮	⋮		⋮
	EC企業 N	dist(N, 1)	dist(N, 2)	...	

民間データは主に借入をしている企業のサンプルデータ、
政府データはセンサスデータ、

お互いに独自の変数があり、統合することによって情報量が増大する。

$$P_{ij} = \frac{\exp(D_{ij})}{\sum_{j=1}^N \exp(D_{ij})} \quad \text{最尤法で推計}$$

$$D_{ij} = \beta_1 |X_{i, \text{資本金}} - X_{j, \text{資本金}}| + \beta_2 |X_{i, \text{従業員数}} - X_{j, \text{従業員数}}| + \dots$$

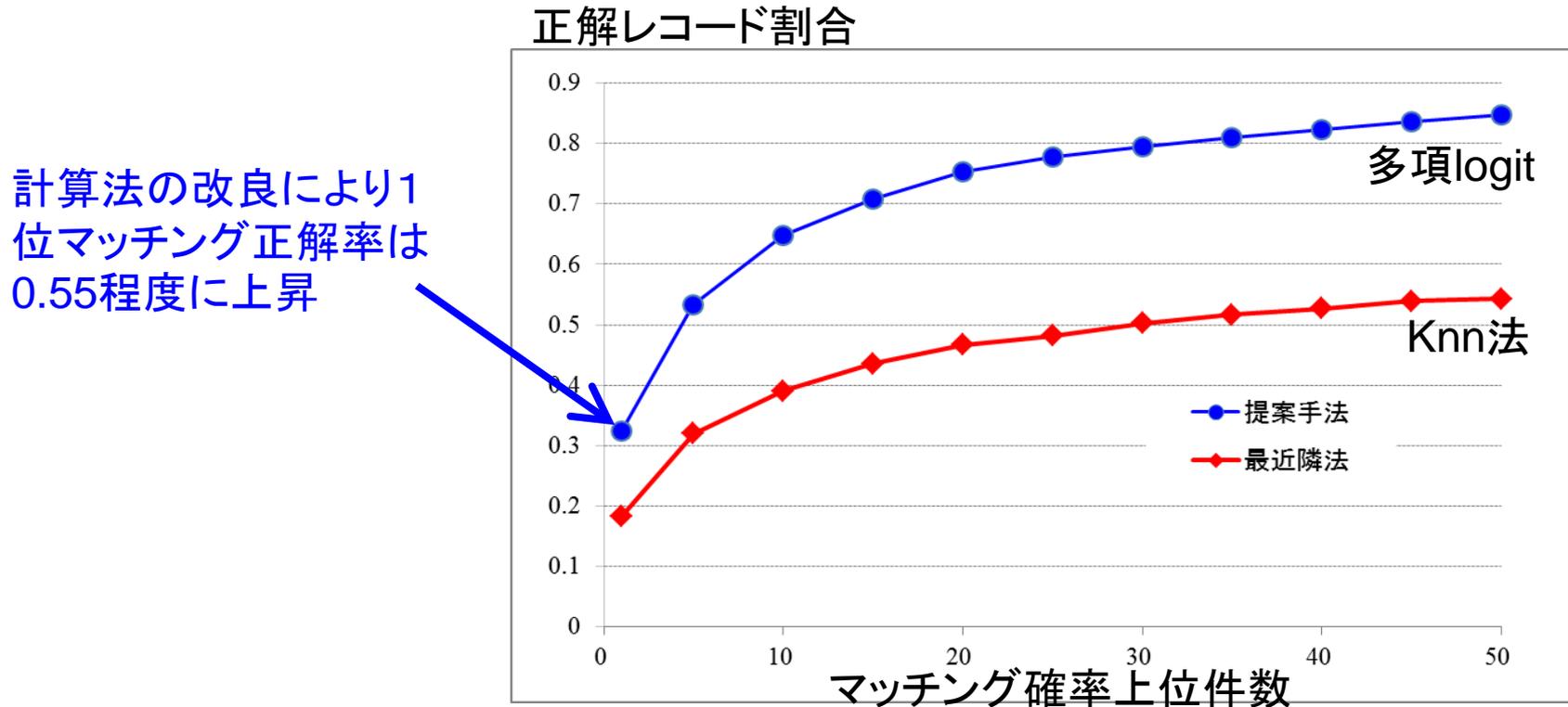
$$LL = \log \left[\prod_{i,j} P_{ij}^{\delta_{ij}} \right] = \sum_{i,j} \delta_{ij} \log [P_{ij}]$$

- LL を最小にするパラメータ β_1 等を求める。
⇒ウエイトを統計的に推定することが可能。

リレーション実験例

2つのデータでマッチングが確認できたものだけで、
正解リレーションデータベースをつくり、
モデルでリレーションを再現できるかを実験。

正解割合はマッチング元から見て、正しいマッチング先を見つけた割合。



社会実装を意識したリレーション実験の例

低質ビッグデータと高質少量データの融合法(アパートローンリスク)

10日ごと●●県の全数データを自動ダウンロード
→約2年分のデータ集積済み

低質ビッグデータ

Web賃貸住宅募集データ(パネル)

空室→占室モデル
(OVモデル)
2項ロジットモデル

高質少量データ

鑑定士によるサーベイ
データ(実地調査・パネル)

占室→空室モデル
(VOモデル)
2項ロジットモデル

銀行の融資情報

部屋の初期状態
 $S_{0,v}(0)$

個別部屋の
状態遷移行列

前 後	占室	空室
占室	$P_{0 \rightarrow 0}$	$P_{0 \rightarrow v}$
空室	$P_{v \rightarrow 0}$	$P_{v \rightarrow v}$

社会実装のための運営会議を
定期的(3ヶ月ごと)に開催

【モデル開発側】
統計数理研究所、国立情報学研究所、長崎大学
社会データ構造化センター

【モデル実装側】
滋賀銀行、CRD協会

CRD協会のデータ提供システムにより全都道府県庁に対して開発モデルを提供を行うことで合意
(すでに全都道府県庁がCRDの会員である)

3ヶ月ごと
滋賀県
の4000軒
のアパートを
パネル調査

7回調査済み
↓
アパートの経
年変化を観
測可能に

t期後の状態
 $P_{0,t}, P_{v,t}$ の
推定

- ・個別部屋のT年後までの
収益シュミレーション
- ・棟全体のポートフォリオ評価
- ・アパートローンの
貸倒確率の計量化

占空モデルの予測力評価
一般的な信用リスクモデルに遜
色ないレベルを達成

(1)占室→空室モデル
上記サーベイのパネルデータを元に、現在空室である部
屋が翌期までに占室になる確率を2項ロジットモデルを最
尤法により推計する

$$P_{i,t} = \frac{1}{1 + \exp Z_{i,t}^P}$$

$$Z_{i,t}^P = \alpha^P + \sum_j \beta_j^P x_{i,j,t}$$

占空モデルの開発
試作版開発済み

$$L^*(\alpha^{*P}, \beta_j^{*P}) = \arg \max_{\alpha^P, \beta_j^P} \prod_{i,t} P_{i,t}^{\delta_{i,t}^P} (1 - P_{i,t}^P)^{1 - \delta_{i,t}^P}$$

