

追跡終了後コホート研究を用いた共通化データベース基盤整備と その活用に関する研究：進捗報告

研究代表者 玉腰暁子（北海道大学大学院医学研究科・教授）

研究要旨

国内で実施され追跡を終了した複数のコホート研究情報を共通化し、その利用環境を整え、将来にわたって向後終了するコホート研究も組み入れ可能な体制を構築するため、データを広く共有化し二次利用を促進するためのシステムであるデータアーカイブ化が進められている社会科学系ならびにライフサイエンス系のデータアーカイブセンターの現状を把握した。疫学研究データには機微情報を含むのみならず、収集する項目数も多いことから、80%程度のレコードは一意性があるものとして対応することが必要であり、研究開始時点での対象者への説明のあり方、完全に連結不可能匿名化するタイミングやその方法などの検討とともに、研究内容によっては共同研究を締結して行うなどの対応が必要と考えられた。また、今後、疫学研究のデータアーカイブを構築し利用を進めていくためには、教育も含めた利用環境整備も必要である。一方、統計法の規定から、現状では人口動態統計資料から得られた死因情報をコホート研究のアウトカムとして公開利用することはできない。そのため、がん登録推進法による死亡者情報票の活用、ならびに現制度化で運用するために必要に応じて死因情報を入手・利用する方法を提案した。

分担研究者

磯 博康（大阪大学大学院医学系研究科・教授）
大橋靖雄（中央大学理工学部・教授）
祖父江友孝（大阪大学大学院医学系研究科・教授）
辻 一郎（東北大学大学院医学系研究科・教授）

A. 研究目的

本研究では、国内で実施され追跡を終了した複数のコホート研究情報を共通化し、その利用環境を整え、将来にわたって向後終了するコホート研究も組み入れ可能な体制を構築するために必要な事項を検討する。今年度は、国内における社会学分野ならびにライフサイエンス分野のデータアーカイブ

の現状、各コホートが持つデータを公開した場合の一意性の問題について検討した。さらに、人口動態統計情報を利用して把握されている死因情報は、統計法の規定により、データアーカイブ時には制限つきであれ公開情報とできないことから、その対応として代替案として、がん登録推進法において示された死亡者情報票の活用の提案、ならびに現制度化での利用方法の提案を行った。

B. 研究方法

各分野の専門家から現状を伺い、追跡を終了したコホート研究データアーカイブを公開する場合の課題を検討した。

三府県コホートデータを用い、分類数、標本数を
変化させた場合にそれぞれどのような頻度で一意
性が見られるかを比較した。

がん登録推進法について、条文や政令をもとに、
追跡終了後コホート研究を用いた共通化データベ
ースが構築された場合における、死亡者情報活用を
めぐる諸問題について検討を行った。

現制度下で二次的に死因情報を利用するため、研
究の都度、コホートデータに死因情報を付与する方
法について、JALS データを用い検討した。

C. 研究結果

-1 社会学分野のデータアーカイブの現状

データアーカイブセンターの意義は、統計調査、
社会調査の個票データを収集・保管し、その散逸を
防ぐとともに、学術目的での二次的な利用のために
提供することにある。社会科学系では、特に若手の
研究者がデータアーカイブを利用して、オリジナル
な枠組みで分析を行い、新たな知見を出していくこ
とがより一般的になってきている。そのためのセン
ターの1つであるSSJDA (Social Science Japan Data
Archive) には現在、約1600件のデータが寄託され
ており、2013年度は2700件の利用があった。このよ
うに活用が進んでいる背景には、データアーカイブ
センターが設立されたこと、ならびに二次分析のメ
リットが広く研究者に認識されたことがある。この
ようにデータが収集・公開され第三者が分析するこ
とは、データの再現性を確認することにつながる。
また、特に公的資金が投入され実施された調査デー
タに関しては、調査者個人のものではないという認
識も広まりつつある。データアーカイブを二次利用
するメリットは、既に行われている調査を繰り返さ
ずに済み労力、資金とも無駄な投入を避けることが
できること、特に多くの変数を得るような調査では
得られたすべての情報を調査者が解析することは
ないため、利用されていない変数について独自のア
イディアで解析することで、新たな知見を得ること
ができること、若手研究者にとっては、自身で小規
模な回収率の高くない調査を行うことに比べ、質の
よい調査データにアクセスできること、学生教育の

際にも、実データを用いた教育を行うことができ
ることである。

現在、SSJDA の運営費用は文部科学省 (2010 年度
より国立大学法人共同利用・共同研究拠点)、東京
大学社会科学研究所から運営費、データアーカイブ
に関わる科学研究費で賄われている。データアーカ
イブセンターの活動として行われている業務の主
なもの、データ寄託の依頼・受付、データ整理、
データ秘匿処理、メタデータの作成、データ利用の
受付・提供、リモート集計の提供、二次利用成果の
公開、データ寄託者の表彰、二次利用促進と適切な
解析のための研究会・セミナーの開催等多岐にわた
っている。

-2 ライフサイエンス分野のデータアーカイブ の現状

バイオサイエンスデータベースセンター (NBDC)
では特にヒトに関する情報に特化した NBDC ヒトデ
ータベースを構築し、2013 年 10 月から運用が開始
された。ヒトを対象とするデータであるため、特に
個人につながる情報の保護対策が重要となる。そこ
で NBDC では、欧米のデータベースを参考に受け皿
づくりが進められている。取り扱われるデータは匿
名化されたもののみで、レベルに応じたアクセス制
限が行われている。データ提供と利用に関する審査
は、NBDC で行われるが、原則として試料提供者か
らデータ共有に関する同意を事前に受けておくこ
とが求められている。ただし、過去に収集された既
存試料・情報で同意を取り直すことが困難な場合に
は、データ共有について倫理委員会で承認されてい
ることが要件である。データ共有に関するガイドラ
イン、セキュリティレベルに関するガイドライン等
が定められ、HP 上で公開されている。2015 年 2 月
現在、研究データ 15 件が HP に公開され、そのうち
制限公開 10 件、オープン 5 件であるが、すべてゲ
ノムに関連するもので、いわゆる疫学研究のデー
タは今までのところ、寄託はされていない。

一意性の検討

三府県コホート対象者約 100,000 例全てを使用

した場合と、無作為抽出により標本数を1万、1,000、100に変化させた場合各々で、分類数 K とユニークセル数 S_1 、分類数に占めるユニークセル数の割合 S_1/K を検討した。全例を使用した場合、分類数 K が増加するとともに、ユニークセル数 S_1 およびユニークセル数の割合 S_1/K は増加するが、ユニークセル数の割合 S_1/K については、分類数 K が約20,000例になるまで急増後80%程度でプラトーに達し、分類数 K が80,000例あたりからさらに増加する、というパターンを示した。この傾向を全例を使用した場合と標本数を減らした場合で比較すると、分布の形状は似通っていたが、急増する部分の勾配が緩やかになり(100例使用の場合は40例程度まで)、プラトー部分が狭くなる傾向があった。

がん登録推進法における死亡者情報票の活用

死亡者情報票を利用したコホートデータ追跡情報入手法を検討するため、癌登録推進法における死亡者情報票の取り扱いについて確認した。がん登録推進法第11条(死亡者情報票の作成及び提出)に死亡者情報票の作成に関する事項は明記されている。それによると、死亡者情報票は、死亡の届書その他の関係書類に基づいて、市町村長が作成するもので、死亡した者に関する氏名、性別、生年月日、死亡の時にける住所、死亡の日、死亡の原因、死亡診断書の作成に係る病院又は診療所の名称及び所在地その他の厚生労働省令で定める情報が含まれる。全死亡者に関する死亡者情報票は、電磁的記録又は書類により作成され、保健所、都道府県を経て、国(国立がん研究センター)に提出される。それを受けて、国立がん研究センターは、死亡者情報票と全国がん登録情報とを照合する。その照合期間は、厚生科学審議会がん登録部会の政令案では100年とされたことから、がん患者の生命予後をほぼ完璧に追跡することが可能になったと思われる。

データアーカイブ化における死因情報の利用：コホートデータに死因情報を付与する方法

統計法の規定から、人口動態統計資料から得られた死因情報をコホート研究のアウトカムとして公

開利用することは不可能な状態にある。そこで、代替案の一つとして、死因を連結した形でのデータセット構築・アーカイブ化ではなく、必要時に中央(アーカイブデータを保持するセンター等)で死因照合作業を行って解析用データセットを作成する方法とその妥当性をJALSの実データを用いて検討した。

提案するデータ利用基盤の概略は次のとおりである。研究コンソーシアムに参加する各研究が、基本データ(生活習慣、検査データなど)と死因を除いた追跡データをアーカイブセンターに提供する。アーカイブセンターでは、基本データベースと追跡データベースを分けて構築する。その際、基本データベースは原則登録時から修正なしの状態、追跡データベース(その後の死因照合作業で必要となる「死亡地(市町村)」、「死亡日」、「生年月日」、「性別」を含む)は、今後の追跡継続に応じて更新できる構造とする。このデータベースを利用した研究を行いたい研究者は、死因情報を得るために厚生労働省に対し人口動態調査二次利用申請を行う。

承認後に提供を受けた死因情報をアーカイブセンター内で、「死亡地(市町村)」、「死亡日」、「生年月日」、「性別」をキー変数として、保有する追跡情報と照合する。死因を付与した一時的な解析データセットを作成し、研究計画に基づいた解析に使用する。研究終了後は死因情報を削除(抹消)し、厚生労働省に利用後報告を行う。

この方法の妥当性を確認するため、JALS対象地域の市町村で1999年1月1日から2012年12月31日までに発生した死亡の調査票情報を厚生労働省に申請、入手した。JALSの対象者で、当該期間中に死亡が特定できていた7,137件(職域コホートと死亡調査データが確定していないコホートを除く)のうち、99.5%が、性別、生年月日、死亡年月日、死亡時の居住市町村名をキー変数として人口動態調査データと一致した。なお、不一致のうち14件は以前にJALSが行った死因照合作業において既に未照合が判明しており、各コホートに対して死亡時情報を確認したがいずれも情報に誤りがなく、人口動態統計作成の過程で入力間違い等が発生した事

例と判断した。このため、今回照合出来なかった例は、実質として19件(0.27%)であった。

D. 考察

国民の税金を投入し、多くの人手と長期の追跡を経て構築されたコホート研究データをアーカイブ化、広く利用可能にすることは、研究の透明性確保、第三者による研究結果の検証、若手の育成に寄与するのみならず、研究の無駄・重複を減らし、必要な公費・労力を新しい有意義な研究に向けるという意義もある。実際、社会科学系分野では、そのためのセンターが設立され、現在では多くのデータが二次利用されている。しかし、今の形になるまでに、10~20年の年月を要しており、データ寄託がある一定数に達するまで、利用のメリットが十分に浸透するよう働きかけるとともに、利用のための環境整備も必要と考えられた。ライフサイエンス分野では、NBDCがデータアーカイブセンターの役割を担い始めた。しかし、いまだ緒についたところで、特にヒトを対象としたデータに関して実績があがるのはまだこれからと考えられた。加えて、疫学研究は単に生体情報のみならず、生活習慣や心情等に関する情報も収集されることが多く、追跡結果も死亡・疾病罹患など機微情報を含む。さらに、収集する項目数も多いため、その組み合わせにより80%程度のレコードは一意性があるものとして対応する必要がある。したがって、研究開始時点での対象者への説明のあり方、完全に連結不可能匿名化にするタイミングやその方法など、今後の検討課題である。また、単に塩基配列などの公開と異なり、疫学研究で構築されたデータセットの公開内容は一部に制限し、機微情報を取り扱う場合には共同研究を締結するなどの対応が必要と思われる。

統計法の規定から、現状では人口動態統計資料から得られた死因情報をコホート研究のアウトカムとして公開利用することはできない。米国ではNational Death Index (NDI) という、厚生省 (U.S. Department of Health and Human Services) の下部機関が、研究目的での生存・死亡確認情報 (死亡時には、死亡年月日や死因などを含む) の提供を行

っている。研究者は、調査対象者リスト (氏名、性、生年月日、住所、社会保障番号など) を提出し、審査にパスすると、有料 (基本料 350 ドル + 対象者 1 人 1 年あたり 15 セント) で、上記情報が提供される。これにより、米国の疫学研究・臨床研究のレベルと即時性は飛躍的に向上し、医学研究や医薬品開発において国際的に有利な地位を確保することができたといえ、今のままでは日本の疫学研究は後塵を拝する。がん登録推進法により提供される死亡者情報票データを活用して、米国の NDI と同様のシステムを作るには法制度の改革が必要であるが、それが実現すれば、わが国の疫学研究・臨床研究や医薬品・医療機器開発は発展すると思われ、それは政府「健康・医療戦略」の目指すところと合致するものであろう。

一方で、現制度化での運用方法を検討するため、死因情報を外したアーカイブ環境を想定し、必要時に死因を人口動態二次利用申請し、アーカイブセンターにて照合・集計・解析を行う運用例を提案した。JALS で実際に死因照合を行ったところ 99.5% で照合が可能であり、照合作業の技術的側面、作業手順化の面で問題はなかった。今後検討すべき課題の一つとしては、人口動態統計二次利用申請に基づくことから、データの保持期間が公的研究費の継続期間に限定されることがあげられる。研究の質や結果の再現性を保証するという点では、解析に使用した死因付きのデータセットが長期に保持できることが望ましい。また、今回提案する方法では、死因付きデータセットの利用場所、すなわち解析場所がアーカイブセンター (あるいは申請書に記載した研究者の所属する機関) に限られる。そのため、アーカイブデータの利用規定も合わせて、アーカイブセンターで対応する場合は、データ解析を行える環境 (物理的な環境、統計家の配置等) について検討する必要がある。また、死因データの申請者の所属機関で実施する場合には、アーカイブデータの外部利用の規約等の整備も必要があるといえる。

E. 結論

疫学研究により得られたデータを広く共有化す

るためのシステムであるデータアーカイブ化に向けた課題を整理するために、社会科学系ならびにライフサイエンス系のデータアーカイブセンターの現状を把握した。疫学研究データには機微情報を含むのみならず、収集する項目数も多いことから、80%程度のレコードは一意性があるものとして対応することが必要である。研究開始時点での対象者への説明のあり方、完全に連結不可能匿名化にするタイミングやその方法などの検討とともに、共同研究を締結するなどの対応が必要と思われる。また、疫学研究のデータアーカイブを構築していくためには、先行する社会科学系データアーカイブの運営システムから学ぶと同時に、利用のための環境整備も必要と考えられた。一方、統計法の規定から、現状では人口動態統計資料から得られた死因情報をコホート研究のアウトカムとして公開利用することはできない。そのため、がん登録推進法による死亡者情報票の活用、ならびに現制度化で運用するために必要に応じて死因情報を入手・利用する方法を提案

した。

F. 健康危機情報

なし

G. 研究発表

1. 論文発表

なし

2. 学会発表

なし

H. 知的財産権の出願・登録状況（予定を含む）

1. 特許取得

なし

2. 実用新案登録

なし

3. その他

なし