

「日本における傷病名を中心とするレセプト情報から得られる指標のバリデーションに関するタスクフォース」報告書

岩上将夫¹ 青木事成² 赤沢学³ 石黒智恵子⁴ 今井志乃ぶ⁵ 大場延浩⁶ 草間真紀子⁷[†] 小出大介⁸ 後藤温⁹ 小林典弘¹⁰ 佐藤泉美¹¹ 中根早百合² 宮崎真¹² 久保田潔^{*13}

*日本における傷病名を中心とするレセプト情報から得られる指標のバリデーションに関するタスクフォース委員長

1 London School of Hygiene and Tropical Medicine

2 中外製薬株式会社 安全性リアルワールドデータサイエンス部

3 明治薬科大学薬学部

4 医薬品医療機器総合機構 医療情報活用部 疫学課

5 国立病院機構本部総合研究センター診療情報分析部

6 日本大学薬学部

7 東京大学大学院薬学系研究科

8 東京大学大学院医学系研究科

9 国立がん研究センター社会と健康研究センター 疫学研究部

10 塩野義製薬株式会社 経営戦略本部 デジタルインテリジェンス部

11 京都大学大学院医学研究科

12 MSD 株式会社 グローバル研究開発本部 薬剤疫学グループ

13 NPO 日本医薬品安全性研究ユニット

[†]現所属: 国立研究開発法人日本医療研究開発機構

キーワード: バリデーション、アウトカム、感度、特異度、陽性的中度、陰性的中度

連絡先: NPO 日本医薬品安全性研究ユニット 〒113-0034 東京都文京区湯島 1-2-13 御茶ノ水明神ビル 4F 久保田潔

E-mail: kubotape-ky@umin.net

〈抄録〉

【目的】医療情報データベース(DB)の活用が薬剤疫学、その他の分野で進んでおり、2017年10月26日に公布され2018年4月1日に実施されるGPSP省令の改正でも医療情報DBを用いた調査が、医薬品の製造販売後データベース調査として具体的に明記された。医療情報DBのデータは通常、調査研究以外の目的で収集されている。したがって、医療情報DBを用いる薬剤疫学研究や製造販売後調査の結果が信頼に足るかを判断する上で、医療情報DBから得られる情報、特に傷病名の正確性をより適切に評価するためのバリデーション研究の実施が必要となる場合も想定される。日本においては、これまでバリデーション研究の報告例は限られており、その経験は不足している。そのため、日本薬剤疫学会の「日本における傷病名を中心とするレセプト情報から得られる指標のバリデーションに関するタスクフォース」(以下「本TF」)は日本の環境下でいかにバリデーション研究を実施するかについての指針を与えることを目的に2016年7月に設立された。対象とするのは、より真に近いと考えられる情報(ゴールドスタンダード)と比較して傷病名などの正確性を評価する研究であり、システムバリデーションや疾患予後予測モデルの妥当性研究は含まれない。

【方法】バリデーション研究に言及した教科書やガイドライン、国内外で実施されたバリデーション研究を検討した。過去の研究の検討では、網羅的な系統的レビューは意図せず、タイトルに“Valid*”が含まれる論文を中心に、これまでにどのような研究が実施されていたのかを俯瞰できる程度にとどめ、(1)利用するDBの種類、(2)バリデーション研究のセッティング、(3)アウトカムの定義、(4)リンケージの有無・方法、(5)ゴールドスタンダードの定義、(6)サンプリング方法とサンプルサイズ、(7)妥当性を測定する指標、(8)指標の閾値・利用法の8項目を検討した。これらに基づき、主に上記8項目に関してバリデーション研究でとるべき手順・チェック項目を検討した。さらに日本において可能なバリデーション研究の有り方を考察した。

【結果】教科書では Strom の Pharmacoepidemiology 第5版41章にバリデーション研究に関する記載が見られた。ガイドラインでは、バリデーション研究の方法が部分的に言及されているものは見出されたが、バリデーション研究それ自体を主たる対象にするものは見出されなかった。検討した国内外100余りのバリデーション研究については、(1)診療報酬請求データ、電子カルテ情報のほか多様なDBが対象にされていた。(2)国・地域・被保険者集団を包括するDBからランダムサンプリングした患者を対象に行った“population-based”と考えられるバリデーション研究も見られたが、少数の医療機関を受診した患者集団や研究者が独自に設定した集団において実施されたバリデーション研究も散見された。(3)ICDコード単独、または薬・検査・診療行為との組み合わせなどがアウトカム定義に用いられていた。(4)リンケージは多くの研究でアウトカム定義の評価に必須の要素であり、多くの研究で利用されていた。(5)ゴールドスタンダードにはカルテレビューのほか、疾患登録などが使われていた。(6)サンプル方法、サンプルサイズは何をゴールドスタンダードとするかに強く依存していた。(7)感度、特異度、陽性的中度、陰性的中度の全てを求めているもののほか、陽性的中度のみを評価したものも多かった。(8)指標の閾値を示したものはほとんどなく、少数、特定の研究の目的を達成する上で必要な「感度〇%以上の指標を見出すことを目標とした」などの記載が見られた。

バリデーション研究の手順としては、「〇〇のデータベースにおいて、□□の疾患を特定する△△の傷病名(またはアルゴリズム)が妥当であることを示す」などのリサーチクエスチョンを明確にすることがまず重要である。今後の研究(バリデーション研究で評価するアウトカム定義を用いて実施する研究)に利用するDBとしては、日本では診療報酬明細書(レセプト)のDBのほか、Diagnosis Procedure Combination (DPC)データ、電子カルテ情報、疾患レジストリーなどが想定される。日本ではDBから元データに遡ることが困難であるため、今後の研究に利用するDBのうち一部の偏ったサンプルを対象にバリデーション研究を行わなければならないことが多い。そのような状況では、バリデーション研究における対象集団の特徴を明確にする、対象集団を重症度で分類したサブグループごとに指標の値を示す、また偽陽性・偽陰性の例の特徴を明確にすることなどにより、バリデーション研究における対象集団と異なる集団においても参考情報として利用できるよう工夫することが求められる。アウトカム定義としては傷病名マスターやICD-10コードを中心とし、複数の定義を評価すべきである。ゴールドスタンダードをカルテレビューで実施する時は、適切な診断基準に沿って判定が可能となるような調査票を設計する必要があり、可能な限り複数の評価者が独立に評価し、その一致度を報告すべきである。求めるべき指標として、陽性的中度だけか、感度、特異度などの評価も必要かはリサーチクエスチョンに依存する。サンプリング方法、サンプルサイズは評価すべき指標に感度が含まれ、ゴールドスタンダードがカルテレビューである時に特に問題となる。対象集団全体またはそのランダムサンプルの全例をレビューすることが基本であるが、特定のアウトカム定義を満たすか否かにより異なる抽出割合でサンプリングする方法、“all possible cases”を想定したサンプリング方法(真のケースを含むと考えられる複数のサブグループからカルテレビューで真のケースを見出し、それを対象集団に含まれる全てのケースとみなす方法)も選択可能である。指標の閾値の絶対的基準を定義することは難しく、各研究の目的により、得られた指標の値を参考にアウトカム定義が適切であるかを判断すべきである。

日本においては、バリデーション研究において重要な役割を果たすリンケージ(照合)は一般に困難であるため、病院(診療所)単位で、保管されている過去の診療報酬明細書(レセプト)を院内のカルテ情報や院内疾患登録などとリンケージ(照合)する病院(診療所)ベースのバリデーション研究が中心とならざるをえない。バリデーション研究にあたっては、DPCデータまたはDPCレセプトの傷病名については医科レセプトの傷病名よりも特異性が高いことが期待され、今後のバリデーション研究でその妥当性が評価されるべきであろう。また、病院(診療所)単位のバリデーション研究以外の方法として、大規模コホート研究における国保レセプトを用いたバリデーション研究も可能な方法と考えられる。

【考察・結論】バリデーション研究はデータベース研究の質を大きく高める(特に、より適切な比較を可能にすることを通じて内的妥当性を高める)ことを可能にする。質の高いデータベース研究は、医薬品などの安全性・有効性に関して、信頼性が高くかつ他では得難い知見の獲得を可能とし、その結果は国民の健康と福祉に貢献しうる。現状、国内でのバリデーション研究の実施は困難を伴う場合が多いが、今後、医療情報DBの適切な利用を進めようとする者は、データベース研究の有用性を高めるためにも、医療情報DBの二次利用に対する国民的理解を得ながら、バリデーション研究を実施しやすい環境の向上に寄与することが求められる。本報告書がバリデーション研

究の意義や基本的な考え方の理解を促し、与えられた条件下で目的を達成するための最適な方法を見出す手掛かりになることを期待する。

目次

1. はじめに
2. バリデーション研究の概要
 - (1) バリデーション研究とは
 - (2) アウトカム定義のバリデーション研究の重要性
 - (3) 代表的な教科書におけるバリデーション研究に関する記載の要約
3. バリデーション研究実施の手順・チェック項目
 - (1) 利用する医療情報データベースの種類を理解
 - (2) バリデーション研究のセッティングの理解
 - (3) アウトカムの定義
 - (4) リンケージの有無・方法の確認
 - (5) ゴールドスタンダードの定義
 - (6) サンプルング方法の選択とサンプルサイズの設定
 - (7) 妥当性を測定する指標の計算
 - (8) 指標の閾値・利用法に関する考察
4. 今後日本で、より望ましいバリデーション研究を実施するために
 - (1) 日本におけるバリデーション研究の現状
 - (2) DPC レセプトの理解とバリデーション研究への利用
 - (3) 大規模コホート研究における国保レセプトを用いたバリデーション研究の可能性
 - (4) 産官学・医療機関における理解の促進
 - (5) 個人を識別するためのリンケージ(照合)
5. 結論

表 1. 妥当性を測定する 4 つの指標

表 2. アウトカムの定義を満たす患者群と満たさない患者群からのサンプルング割合が異なる場合の感度・特異度の求め方

表 3. “all possible cases”を想定したサンプルングを行った場合の 2×2 テーブルの作り方

図. ゴールドスタンダードの定義、サンプルング方法とサンプルサイズの設定、妥当性を測定する指標の計算、の対応関係

付録 1: 同じ感度・特異度であっても研究集団の曝露の割合によって研究結果が変わるケース・コントロール研究の例

付録 2: データベース研究に関連するガイドライン・ガイダンスにおけるバリデーション研究に関する記載の要約

付録 3: 過去に実施されたバリデーション研究のまとめ

付録 4: バリデーション研究におけるサンプルサイズ

付録 5: バリデーション研究の手順・チェック項目の一覧表

1. はじめに

診療報酬請求データ(claims data)や電子カルテ(electronic health records, EHR または electronic medical records, EMR)のような日常診療における医療情報(routinely collected data)は圧倒的な情報量を持つことが最大の強みである⁽¹⁾。こうした医療情報、医療情報データベースを用いた薬剤疫学研究は海外では 1990 年頃より盛んに行われており、昨今、日本においても徐々に見られるようになってきている。

また、国内の薬事規制においても医療情報データベースの活用が進められており、2012 年 4 月 11 日付で公表された「医薬品リスク管理計画指針について(薬食安発 0411 第 1 号、薬食審査発 0411 第 2 号)」の中では、安全性監視活動において医療情報データベースの利用可能性について言及があり、また、2017 年 10 月 26 日付で交付された「医薬品の製造販売後の調査及び試験の実施の基準に関する省令(平成 16 年厚生労働省令第 171 号)」(Good Post-marketing Study Practice, GPSP)の改正によって医薬品の再審査及び再評価の申請資料として医療情報データベースの利用が認められるようになった。こういった薬事規制上の変化に伴い、医薬品の製造販売後調査における医療情報データベースの利活用が現実的課題として上りつつある。

薬剤疫学の分野において医療情報データベースを用いた観察研究は、ランダム化比較試験のサンプルサイズ・観察期間では検出されないレベルの稀なアウトカムと薬剤の関連の検討や日常診療下における薬物使用実態の調査に適している。一方、明確なリサーチクエスチョンを背景に計画されたコホート研究や疾患レジストリーと比べると、情報の正確さや詳細さが不十分であったり、必要な交絡因子の情報が不足していたりすることが弱点である⁽¹⁾。日本薬剤疫学会では、2012 年の「医薬品リスク管理計画指針」に関連し、日本語/英語版の「より良い Pharmacovigilance Plan 策定に向けての提言」⁽²⁾を発表し、日本の市販後の調査の在り方について一石を投じてきた。その提言の中で「データベース研究は安全対策の意思決定につながるエビデンスを与えることもありうるが、アウトカムなどに関する指標のバリデーションが不十分であるか、他のデータ源との連結が困難で、医療機関内の元データにもどっての確認ができない場合には次のアクション (Primary Data Collection を含む)を行うべきかを判断するためのスクリーニングの役割にとどまるだろう。」と言及している通り、データベース研究の信頼性を高める上でバリデーション研究の実施は重要なステップである。しかし、我々の知る限り、バリデーション研究についての方法論・マニュアルは国内外ともに存在しない。

そこで日本薬剤疫学会では「日本における傷病名を中心とするレセプト情報から得られる指標のバリデーションに関するタスクフォース」を立ち上げ、日本・海外において実施されてきたアウトカムなどに関するバリデーション研究及び関連するガイドライン・教科書をレビューし、日本において実施可能なバリデーション研究の手順・問題点・将来的な課題を明確にすることを目標に 2016 年 7 月から活動を行ってきた。その総括としての本報告書は、医療情報データベースを用いて薬剤疫学研究及びそのためのバリデーション研究を行う研究者を対象とするだけでなく、研究の結果を解釈し行政的な判断を行う方々が参照することも想定している。さらに本報告書に書かれたバリデーションに関する方法論・考え方は、日本のレセプト(正式名称は診療報酬明細書または調剤報酬明細書)を用いた薬剤疫学研究だけでなく、医療情報データベースを用いた臨床疫学研究に幅広く適用可能である。

2. バリデーション研究の概要

(1) バリデーション研究とは

バリデーション(validation)は「妥当性(validity)の確認」を意味する。医療情報データベースを用いて薬剤疫学研究を行う際には、研究対象集団・曝露・アウトカム・交絡因子をデータベースから得られる情報(入力された傷病名、処置名、処方内容、検査結果など)を用いて定義する必要があり、本タスクフォースではこの定義を「より真に近いと考えられる情報(ゴールドスタンダード)と比べて妥当性を確認すること」をバリデーションとし、そのために行う研究をバリデーション研究と呼ぶ。

なお、医療情報データベースなどを用いる研究においてバリデーションという言葉は、本タスクフォースの指すバリデーションとは異なる概念で使われることがあるため、注意が必要である。例えば、予後予測モデルの評価もまたバリデーションと呼ばれることがある^(3, 4)。また、薬事規制下において、医薬品承認申請等の資料作成に申請者が使用するコンピューター化システムに対して期待される結果を与えることを検証することはコンピューター化システムバリデーション(CSV)と呼ばれている⁽⁵⁾。

(2) アウトカム定義のバリデーション研究の重要性

データベース研究の際に定義しなければならない要素の中でも、特にアウトカム定義は重要であり研究の質を左右する。アウトカムが何らかの疾病である場合、アウトカム定義にはInternational Classification of Diseases Version 10 (ICD-10)コードなどの傷病名を用いることが多いが、その傷病名が真の疾患発症を意味しているか、また真に疾患発症している患者に対して適切な傷病名が入力されているかどうかは確認してみなければわからない。妥当性が確認されていないアウトカム定義を用いた研究結果は信頼できない可能性がある。入力された傷病名の妥当性は疾患毎並びに医療機関毎に異なる可能性があり、また国の医療制度やデータベース本来の目的(例:診療報酬請求を目的としているか否か)によっても影響を受ける⁽⁶⁾。このため、一概に医療情報データベースの傷病名を妥当と考えることは適切ではないし、逆に日常診療データの傷病名は全て研究で使用できないと決めつけることも適切ではない。より高いエビデンスレベルを求めらるならば、利用するデータベースの中でアウトカム定義の妥当性を疾患毎に適切な方法で得られたサンプルを用いて検討し、妥当性が高いことが確認できたものに絞って薬剤疫学研究を行うことが望ましい。

なお、アウトカム定義だけでなく交絡因子(confounder)や曝露(exposure)の定義の妥当性も重要である。交絡因子については、最終的な研究結果(例:曝露とアウトカムの関連を示す相対リスク)に決定的な影響を与える交絡因子が存在すると考えられるのなら、その因子についてはバリデーション研究の対象となりうる。しかし、交絡因子に固有の研究方法があるわけではなく、通常、アウトカム定義のバリデーション研究と同様に行う。また、曝露については、データベースを利用する薬剤疫学研究において、特定の薬剤の処方または調剤の電子記録によって定義される場合がほとんどであり、電子記録が標準化され同じようにデータに格納されている限りは、高い妥当性を持つと考えられる。処方または調剤の記録が正しくとも、患者が実際に薬を服用しなければバイアスを起こしうるが、コンプライアンス(服薬遵守)はデータベース研究以外の観察研究や介入研

究でも起こりえる問題であり、本タスクフォースが扱うバリデーション研究の対象外である。

(3) 代表的な教科書におけるバリデーション研究に関する記載の要約

バリデーション研究に特化した具体的な方法論・マニュアルではないものの、薬剤疫学の代表的な教科書である *Pharmacoepidemiology, 5th Edition* (Brian L. Strom ら) の中の、Part IV Chapter 41 Validity of Pharmacoepidemiologic Drug and Diagnosis Data には、妥当性を論じるにあたっての留意事項が多く記載されている。以下は、その重要な点についてまとめたものである。

① バリデーション研究について

- (i) バリデーションはデータベースから得られる情報とゴールドスタンダードを比べることにより行うものである。つまり、ゴールドスタンダードではない情報源と比較した結果を「バリデーション(研究)」と呼ばないように明確に区別し、研究者は用語の使用に関し細心の注意を払うべきである。
- (ii) ゴールドスタンダードは、(バリデーション研究の先にある最終目的となる研究の)リサーチクエスション、アウトカム及びゴールドスタンダードの設定に利用できるデータ、他のデータソースの利用可能性等で選択される。
- (iii) バリデーション研究のベストプラクティスには、診療記録番号のような 2 つのデータセットを正確にリンケージ(照合)する情報が必要不可欠である。

② 妥当性を測定する指標について

妥当性(validity)を測定する指標には、感度・特異度・陽性的中度(positive predictive value, PPV)・陰性的中度(negative predictive value, NPV)がある(表 1)。感度と特異度は一般的にトレードオフの関係(感度の高い指標では特異度が低く、特異度の高い指標では感度は低くなる関係)があり、また的中度(PPV と NPV)は感度・特異度に依存する。感度と特異度を重要視しつつ、バリデーション研究においては、感度・特異度・陽性的中度・陰性的中度の 4 指標を算出できるようデザインすることが理想的である。また、これらの指標のうち、いずれの指標が最も重要な指標であるかは、研究のセッティングにより変わる。さらに、同じ感度・特異度であっても、それらが研究結果に与えるバイアスの程度は、研究集団における曝露の割合等の影響を受けるため、指標の絶対値にはあまり意味がない(具体例は付録1を参照)。

表 1.妥当性を測定する 4 つの指標

		ゴールドスタンダードの定義を満たすか (真のケースか)			合計	
		Yes	No			
アウトカムの定義を満たすか (陽性のケースか)	Yes	a (真の陽性)	b (偽陽性)	a+b	陽性的中度(PPV) =a/a+b	
	No	c (偽陰性)	d (真の陰性)	c+d	陰性的中度(NPV) =d/c+d	
合計		a+c	b + d	a+b+c+d		
		感度=a/a+c	特異度=d/b+d			

③ 誤分類に関する考慮事項

曝露(医薬品)に関する administrative data(管理データ)や medical record data(医療記録データ)の利用は、未請求処方、アドヒアランス、OTC 医薬品の併用などの情報を得ることが出来ない等、考慮すべき点もあるものの、(思い出しバイアスが懸念事項となる)質問票等での情報収集と比較しアドバンテージがあり、誤分類については一般に大きな問題とは考えられていない。これに対し、診断に関する administrative data の利用においてはよく懸念となる。特に、外来患者の診断データは、確からしさが低いため注意が必要である。医学的に定義が不十分な疾患や、保険者からの医療費償還時のルールや保険の範囲等、様々な要因が ICD-10 コード化時のバイアス(系統誤差)を起こしている可能性もある。さらに、既に評価されたデータであっても、医療の進歩等や医療環境の変化に伴い、医療行為やコード分類が変化することもあり、再度の妥当性評価が必要となる事がある。

なお本タスクフォースでは、データベース研究に関連する国内外の代表的なガイドライン・ガイドランスにおいても、アウトカムのバリデーション研究に関する言及の有無及び内容の確認を行ったが、研究方法や結果の解釈の仕方についての包括的な記述は確認できなかった(詳細は付録 2 を参照)。

3. バリデーション研究実施の手順・チェック項目

まず本タスクフォースでは、過去のバリデーション研究のレビューを行った(詳細は付録 3 を参照)。その結果を踏まえて、本タスクフォースから日本で行うバリデーション研究の手順及びチェック項目について提案する。

まず研究の前提事項として、バリデーション研究も数ある研究の1つであり、明確なリサーチクエスションを立てて研究を始める必要がある⁽⁷⁾。その際、「〇〇のデータベースにおいて、□□の疾患を特定する△△の傷病名(またはアルゴリズム)が妥当であることを示す」(仮説検証的研究)あるいは「〇〇のデータベースにおいて、□□の疾患を妥当に特定できるアルゴリズムを探す」(探索的研究)旨の記載が基本となるが、さらに一步踏み込んで、その研究結果が今後どのような研究で使われるのか、最終目的を明確にすることが望ましい。何故ならば、バリデーション研究後の使用目的によって、求めるべき指標やサンプリング方法が異なってくるからである。例えば、後述する通り、薬剤Aと薬剤Bのあるアウトカムに対する相対リスクを正しく求めることが最終目的であれば、アウトカム定義を満たす患者群からのみランダムサンプリングを行い高い陽性的中度を示せることができれば十分かもしれない。一方で、アウトカムの発生率を検討する研究で使われるのなら、感度も求める必要がある。また、他の臨床研究同様、ガイドライン(Guidelines for Good Pharmacoepidemiology Practices など)に沿ったプロトコルの作成と倫理審査のステップを経ることが望ましい。また近年は観察研究であっても、研究プロトコルを UMIN 臨床試験登録システムなどに事前登録することが求められるようになってきている⁽⁸⁾。さらに可能であれば、少数を対象としたパイロット研究により、研究の実施可能性や必要な調査項目の記載漏れ等について確認することが望ましい。

ここから以下の(1)~(8)の項目は、過去のバリデーション研究のレビューの際に用いた項目(付録 3)と対応している。

(1)利用する医療情報データベースの種類を理解

今後の研究、すなわちバリデーション研究で評価する(データベースから得られる情報による)アウトカム定義を用いて実施する研究に利用する医療情報データベースが、入院・外来の診療報酬請求データ・退院時記録・電子カルテ・疾患レジストリー、またはそれらの組み合わせ、のどれに当たるか良く理解しておく。日本のレセプトデータ(DPC レセプト含む)は診療報酬請求データに当たり、Diagnosis Procedure Combination (DPC)データは主に厚生労働省による DPC/PDPS (Diagnosis Procedure Combination / Per-Diem Payment System: 診断群分類包括評価制度)の評価のための退院時記録からなる administrative data である。また厚生労働省と PMDA が構築中の医療情報データベース(Medical Information Database Network、以下 MID-NET®)は当該病院における診療報酬請求データ及び電子カルテを組み合わせたデータベースである。

また、近年多くの日本の学会が疾患レジストリーを開設しているが、海外のそれと異なり、疾患レジストリーを他のデータベースとリンケージ(照合)させて研究に使用することは現状では困難である。このため、疾患レジストリー自体に曝露やアウトカムの情報を蓄積して研究を行う、自己完結型のコホートの構築が志向されている。疾患レジストリーについてもバリデーション研究の実施は重要であり、倫理上の問題の整理を含めて将来的な課題と考えられる。

バリデーション研究を行う前にまず、今後の研究に利用する医療情報データベースがどのような集団を対象に実施され、また”population-based”と言えるのか、について理解しておくことが重要である。この報告書では、「居住地域や加入する保険などが共通する健康人を含む集団」を対象とする研究を”population-based”の研究と呼ぶ。日本において利用可能なデータベースのうち、レセプト情報・特定健診等情報データベース(いわゆる National Database、NDB)のレセプトデータは日本の全国民から得られたデータであり、”population-based”と言える。また、複数の健康保険組合が保有するレセプトデータを蓄積したデータベースは、特定の種類の保険者のデータであり、性・年齢分布は日本の全国民を代表してはいないが、特定の集団を代表する”population-based”のデータベースと考えることができる。一方、特定の医療機関のデータを蓄積したデータベースについては、日本では患者が受診する医療機関を居住地域に関係なく自由に選択できるため、”population-based”とみなすことはできない。また、過去のバリデーション研究の中には、一定期間にある病院に入院した、またはクリニックで診療を受けた全員を対象として実施された研究⁽⁹⁻¹⁴⁾や、あるデータベースの中で特定の疾患をもつ患者集団を対象とする研究⁽¹⁵⁻¹⁹⁾もみられた。

(2)バリデーション研究のセッティングの理解

理想的なバリデーション研究は、今後の研究に利用する医療情報データベースそのもの(からランダムサンプリングを行った集団)を対象とした研究である。この場合、バリデーション研究で求めた妥当性の指標(感度・特異度・陽性的中度・陰性的中度)を、利用する医療情報データベースに直接外挿可能である。しかし現実的には、今後の研究に利用する医療情報データベースのうち一部の偏ったサンプルを対象にバリデーション研究を行わなければならないことが多い。

例えば、上述の通り、NDB のレセプトデータは日本の全国民から得られたデータであり、”population-based”である。しかし日本の現状では、レセプト側から患者または医療機関をランダムサンプリングして元データに遡ってバリデーション研究を実施することは困難である。このため、まずバリデーション研究が実行可能な医療機関を任意に決めるところから開始しなければならないが、その(単数または複数の)医療機関が NDB などレセプトデータを得た集団を代表し、”population-based”と呼べるバリデーション研究になっていることはまずありえない。

このような条件下で実施されるバリデーション研究、特にバリデーション研究で評価したアウトカム定義をレセプトデータで用いることを意図している場合には、バリデーション研究のセッティング・対象集団の特徴の記述がより重要となってくる。例えば、バリデーション研究が実施された地域と期間、医療機関の数と特徴(例:大学病院、救急指定病院、がん拠点病院、受診患者数、病床数、等)および対象集団の特徴(年齢・性別構成、重症度の分布、等)を記述して、これを利用するレセプトデータにおける平均的な特徴と比べることなどにより、どれだけバリデーション研究の結果がレセプトデータ全体に外挿可能かを議論することが望ましい。

また、レセプトデータベース研究での利用を意図して行う、日本における病院単位のバリデーション研究では、バリデーション研究における対象集団の定義への配慮が必要である。例えば、特定の病院の整形外科で診療を受けている患者が糖尿病にも罹患しているものの、糖尿病については他の医療機関で診療を受けているのなら、当該病院に保管されているレセプトデータには糖尿病の傷病名が含まれていない可能性がある。しかし、この患者を糖尿病の傷病名を含むアウト

カム定義の偽陰性例に分類することには慎重でなければならない。他の医療機関で糖尿病の診療を受けていれば、その傷病コードはレセプトデータベースには存在するからである。一般に他の医療機関で診療を受けていることが明らかな患者は病院単位のバリデーション研究では対象集団から除外するか、除外する／しない場合の両方の結果を示すことが適切である。たとえば、先天異常に関するバリデーション研究では、出産後、他の医療機関に移った症例を除くと陽性的中度が高くなることが報告されている⁽²⁰⁾。

一方、特定の医療機関のデータを蓄積したデータベースを用いた研究を意図している場合には、データベースそのもの(からランダムサンプリングを行った集団)を対象にバリデーション研究が実施可能である場合が多い。しかし、データベースを構成している医療機関の一部の医療機関のみでバリデーション研究を行わなければならない場合は、バリデーション研究が実施された医療機関・対象集団の特徴とデータベース全体の平均的な特徴の対比により、結果の外挿性の議論が必要である。

なお、“population-based”ではない医療情報データベースを用いてバリデーション研究を行った場合、または、ある医療情報データベースのうち特定の疾患をもつ患者集団を対象としてバリデーション研究を行った場合には、求めた感度や特異度の解釈の仕方には注意が必要である。例えば、一定期間にある病院に入院した、またはクリニックで診療を受けた全員を対象とする研究⁽⁹⁻¹⁴⁾では、感度・特異度は「その病院に受診したという条件付きの感度・特異度」を意味する。ある疾患をもつ患者集団を対象とする研究⁽¹⁵⁻¹⁹⁾では、感度・特異度は「その疾患を持つという条件付きの感度・特異度」を意味し、その疾患を持つ患者集団における(その疾患の続発症などの)アウトカムの有病割合(prevalence)または発生割合/率(incidence proportion/rate)を推定するのに使うことができる。また、後述するように、感度、陽性的中度などの妥当性の指標の値が異なるサブグループを見出すこと、偽陽性、偽陰性の患者の特徴を明らかにすることなどが、バリデーション研究で評価したアウトカム定義が、それを用いて実施するデータベース研究においてどの程度有用かの情報を与えることにつながると考えられる^(13, 20-29)。

一方で、研究者がアウトカム定義を満たす集団と満たさない集団から異なる割合でサンプリングしたものを単純に足し合わせることによって作り出した集団において感度・特異度を評価した研究も見られた^(26, 30-34)が、このような恣意的な集団はその後の研究で再現できないため、求めた感度・特異度に意味を見出すことは難しい。また、あるアウトカム定義を満たす(したがって本来、陽性的中度しか求めることができない)集団において感度、特異度などを求めている研究^(35, 36)も見られたが、本タスクフォースとしては、アウトカム定義を満たす集団と満たさない集団から異なる割合でサンプリングした場合には、後述する適切な補正を行うなどにより、サンプルを得た集団全体における感度や陽性的中度などを示すことを推奨する。

(3)アウトカムの定義

医療情報データベースにおけるアウトカムは、1つの要素(傷病名のみ、薬剤名のみ、処置名のみ、検査結果のみ、など)、または2つ以上の要素を組み合わせたアルゴリズムによって定義されるが、傷病名、または傷病名とそれ以外の情報の組み合わせが用いられることが多い。本タスクフォースは必要に応じて複数のアルゴリズムについて陽性的中度などの指標を求め、比較を行う

ことを推奨する。傷病名はコード化されたものを使うことが多く、日本のレセプトでは独自の(ただし、ICD-10 と対応可能な)傷病名マスターが使用され、DPC データベースでは ICD-10 コードも使用される。傷病名によるアウトカム定義では、傷病名コードのリスト作り、および傷病名コードが入力された期間と期間内の入力頻度の選択、が必要である。まず傷病名コードのリスト作りについて、同じ疾患の定義を意図していても、異なる研究者が選んだ傷病名コードのリストは異なる可能性がある(例: 上部消化管出血を定義するための傷病名コードのリストに吐血や黒色便を含めるかどうか)。よって、その疾患の専門家も含めた 2 人以上の研究者が議論してリストを作ることが好ましい⁽³⁷⁾。また、ICD-10 コードのように海外でも使われている傷病名コードを用いる場合は、過去の研究で使われた傷病名コードのリストを各文献から参考にできることがある。作成または参照した傷病名コードのリストは論文内または Appendix に明記すべきである。

次に傷病名コードが入力された期間と期間内の入力頻度の選択について考察する。日常診療、特に外来では一度だけの入力病名は検査目的で入力されている可能性があり、あるいは医師が疑い病名を意図して入力した際に「疑い」を付け忘れた可能性も否定できない。そこで、カナダのバリデーション研究で見られるように、外来で一定期間内に 2 回以上の傷病名が入力された時に初めてアウトカムと定義するなど、一定期間内の受診回数をアウトカム定義に含める場合がある。ただし、カナダでは外来診療であっても primary diagnosis の指定が医療費償還の条件になるのに対し、日本の外来診療は出来高払いであり、「主傷病」の指定は必須ではないなど違いがあることを認識しておくことは重要である。また、日本の診療報酬請求制度を考慮すると、「異なる月に」2 回以上の傷病名入力という条件を付けたほうが良いかもしれない。一方、カナダのバリデーション研究では、入院した場合は 1 回の主診断(primary diagnosis)でアウトカム定義が満たされるとしていることが多い⁽³⁸⁻⁴²⁾。日本のレセプトデータでも入院レセプトと外来(入院外)レセプトは判別できるため同様の定義は可能である。ただし、どのような定義が適切かは疾患に大きく依存する(例: 慢性疾患や癌は 2 回以上の傷病名入力が好ましいかもしれないが、急性肝炎や転倒・骨折などの一過性のアウトカムについては 1 回の傷病名でも十分かもしれないし、心筋梗塞など入院することが確実なアウトカムでは入院病名のみを対象としたほうがいいかもしれない)。そのため、その疾患の専門家も含めた 2 人以上の研究者が議論し、複数のパターンが提案された場合はバリデーション研究で比較することが望ましい。

2 つ以上の要素を用いてアルゴリズムを作成する場合^(11, 12, 37, 43, 44)も、複数のパターンを作ってバリデーション研究で比較し、最終的にどのアルゴリズムを用いることが適切か論文の中で議論すべきである。なお、傷病名以外の情報を使う場合、医療機関によって処置や検査の方針が異なる可能性がある(例: 上部消化管内視鏡処置を行う患者に対して全例生検を行う病院もあれば、癌が疑われた場合にのみ生検を行う病院もある)ため、他施設でも処置・検査の方針が同様に注意しながらアルゴリズムを作成することが、結果の一般化可能性を高めることにつながる。

(4)リンケージの有無・方法の確認

海外のバリデーション研究では、評価されるデータベース(全体、またはサンプリングされた一部の患者)をその他のデータベースまたは電子カルテとリンケージ(照合)させて、ゴールドスタンダードの情報を得ることが多い。またカナダの(入院の)Discharge Abstract Database と(外来の)

claims dataなどを組み合わせた administrative data のように評価されるデータ自体が複数のデータベースのリンケージ(照合)によって成り立っている場合もある^(39, 40)。その際、リンケージ(照合)させるデータベース同士に共通の個人レベルの識別子(social security number など)がある場合は1:1 対応させることができ、通常 deterministic(一方のデータベースに含まれる特定の1個人が他方のデータベースに含まれる1個人と同一であることが確実である時に行う)リンケージ(照合)が可能となる^(23, 25)。それができない場合は患者の氏名・生年月日・居住地域などの情報を用いて deterministic または probabilistic(一致する確率が一定程度以上であれば行う)リンケージ(照合)を行う⁽⁴⁵⁾。Deterministic リンケージ(照合)では入力データにエラーが多いとリンケージ(照合)ができないため除外される人の割合が高くなり、Probabilistic リンケージの場合は異なる2人の個人が同一とみなされる可能性が排除できないため、リンケージ(照合)の詳細について記述することが望ましい⁽⁴⁶⁾。

しかし日本の法律・指針および制度では現在、医療情報データベース同士を照合することやレセプトや DPC データベースから各医療機関の電子カルテに遡ることは困難である。特に改正個人情報保護法における「匿名加工情報」および次世代医療基盤法における「匿名加工医療情報」に該当するデータについては、カルテレビューにより行うバリデーション研究で必要となる特定の個人を識別リンケージ(照合)することが法的に禁止されている。匿名加工情報や匿名加工医療情報でないデータについても(例:NDB や疾患レジストリー)、個人レベルでのリンケージ(照合)は大きく制約されていることが多い。よって当面、日本で実施されるバリデーション研究の多くでは、各医療機関の中で保存しているレセプト・DPC データと電子カルテや院内の疾患レジストリーへの登録情報を比べる作業に頼らざるを得ない。なお、各医療機関の中では同じ患者の情報が元々まとめて保存されてあるという前提があるため、この作業はリンケージ(照合)には当たらない。

(5)ゴールドスタンダードの定義

ゴールドスタンダードの設定は、まずカルテレビューをするか否かに大別される。カルテレビューをしない状況は比較的稀で、疾患レジストリーが存在する場合か、臨床検査(例:低ナトリウム血症⁽⁴⁷⁾)や病理(例:癌の確定診断⁽⁴⁴⁾)など検査結果のみでゴールドスタンダードが設定できる場合が含まれる。疾患レジストリーをゴールドスタンダードとする場合^(25, 29, 45, 48, 49)、疾患レジストリーがそもそも高い妥当性を持ち合わせているかを確認しておく必要がある。まだ開設されたばかりで網羅的に症例が登録されていない、あるいは疾患の正確性が十分確認されていない症例が多く混ざっているような疾患レジストリーはゴールドスタンダードとして使うことはできない。疾患レジストリーの妥当性に関する研究(疾患レジストリーのバリデーション研究)が過去に実施されている場合はその結果を引用すべきであるし、実施されていない場合は疾患レジストリーがゴールドスタンダードとして妥当であるとの判断が合理的である旨記述することが望ましい。検査結果をゴールドスタンダードにする場合は、検体が正確に測定・標準化(例:尿蛋白の±、+、2+の分類)・記録されているか、また病理診断や画像診断の場合は診断結果が標準化・電子化されているかを確認しておく必要がある。もし標準化・電子化されていない場合は、標準化・電子化の作業をまず実施するか、実質的に医療機関に戻って行うカルテレビューと同じ方法でバリデーション研究を実施しなければならない。

カルテレ뷰をする場合、その質を可能な限り担保する努力が必要である。過去の文献の中には「カルテレ뷰により目的の疾患の有無を判断しゴールドスタンダードとした」と一文のみ記述しているものが見られたが⁽⁵⁰⁾、これは好ましくない。まず、対象疾患の判断基準について、世界的に確立した診断基準がある場合には、その診断基準にそった情報収集が可能になるような調査票を作成・利用し、評価することが望ましく、報告時には用いた基準を明記する。明確な診断基準がない場合には、その研究でどのような基準を用いたかを明示するべきである。いずれの場合も使用した調査票については、論文内またはAppendixなどの形で示すことが望ましい^(51, 52)。そして、その診断基準または調査票に沿って、1人の患者につき2人(以上の少人数)の評価者が別々に判断することが望ましく、さらに1人の評価者が同一の患者について異なる時点で判断することも考慮する。研究計画書と報告書の中では、評価者間(内)の一致度を表す κ (カッパ)係数を信頼性(reliability)の指標として算出する^(41, 53, 54)。レビューするカルテの数が多い場合には、その一部についてのみ評価者間(内)の一致度を算出することも認められる⁽⁴⁰⁾。なお、評価者の専門性(例:対象疾患の専門医、訓練を受けた医療関係者、等)も記載しておくことが望ましい。評価者の中に対象疾患の専門家がいない場合は、(複数の)スクリーニング担当と(1人の)最終判断担当に分ける方法^(41, 51)も許容されるかもしれない。その場合も、スクリーニング担当者間(内)の一致度は求めることを推奨する。 κ (カッパ)係数の他に、一致しなかった評価結果の取り扱い(例:2人で話し合っって最終的に決めるか、3人目の判断を仰ぐかなど)についても記述する。

(6) サンプルング方法の選択とサンプルサイズの設定

サンプルング方法は、以下(A)~(E)の方法に大別される。それぞれの方法に長所・短所があり、目的と状況に応じて適切に選択することが望ましい。なお、この分類は本タスクフォースで行った過去の約100文献のレビューにもとづいたものであり、(A)~(E)と異なる方法でサンプルングを行うことは必ずしも誤りとは言えない。

(A) 全患者を対象にする場合

まず、カルテレ뷰をしない場合、すなわち疾患レジストリーや検査結果を用いてゴールドスタンダードを定義する場合、データベースまたは対象の医療機関に登録されている全患者を検討することが多い。つまり、データベースの中でアウトカムの定義を満たす患者を全て特定し、一方でゴールドスタンダードの定義を満たす患者も全て特定し、全患者を対象に2×2テーブルを作成する。調査症例数を制限なく検討できる理由は、ひとえにアウトカムとゴールドスタンダードの定義を満たす患者を、それぞれコンピューターを用いて簡単に特定できるからである。もし疾患レジストリーや検査結果が標準化・電子化されておらず手作業で判断する必要がある場合は、まず標準化・電子化による作業の効率化が必要である。

(B) 全患者からランダムサンプルングを行う場合

カルテレ뷰をする場合は、大きな労力・時間がかかるため全例調査は現実的でなく、通常調査症例数を制限せざるを得ない。その際、利用するデータベースまたは対象の医療機関の全患者の中から調査症例をランダムサンプルングすることは可能であり、実際、数千から数万例のランダムサンプルのカルテレ뷰によるバリデーション研究がこれまでにいくつか実施されている^(40, 41, 55)。しかし、頻度の低い疾患では、例え数千例のランダムサンプルングをしたとしてもアウトカ

ムとゴールドスタンダードの定義を満たす患者がそれぞれ数例しか含まれず、妥当性の指標を十分に計算できない⁽³⁸⁾。望ましい精度で感度・特異度・陽性的中度・陰性的中度の全てを求めるためには、全例を代表する①アウトカム定義をみたす集団(陽性のケース)、②アウトカム定義をみたさない集団(陰性のケース)、③ゴールドスタンダードの定義を満たす集団(真のケース)、④ゴールドスタンダードの定義を満たさない集団(偽のケース)を一定数(例:最低 100 人、根拠は後述の段落を参照)特定しなければならない。アウトカムが稀な疾患の発生である場合に、このうち最も困難なのが③にあたる真のケースを一定数以上特定することである。例えば、研究対象集団に含まれる真のケースの期待される割合が 1%であるのなら、100 人の真のケースを特定するためには 10,000 人(100 人÷0.01)のランダムサンプルが必要である。一方、アウトカムの頻度が比較的高いセッティングにおけるバリデーション研究(例:集中治療室でしばしば起こるアウトカム⁽⁵⁶⁾)では、比較的小さなサンプル数でも十分な数の真のケースが含まれる可能性がある。

(C) アウトカム定義を満たす患者群からのみランダムサンプリングを行う場合

データベースの中でアウトカム定義を満たす患者群のみからランダムサンプリングして、その症例に限ってカルテレビューをする方法である^(37, 53, 57, 58)。ただしこの方法では、アウトカム定義を満たさない患者(陰性のケース)が含まれないため、陽性的中度しか求められないことに留意する。その際サンプルサイズの設定としては、「求めたい指標の 95%信頼区間が最大±10%を目標とするならば 100 例、±5%を目標とするならば 400 例必要」という計算式があり(詳細は付録 4 を参照)、これを目安として米国の Mini Sentinel 調査では同サンプリング方法で陽性的中度を求める際には 100 例あれば十分と考えている⁽⁵²⁾(前述した最低 100 人の根拠)。ただし、この見積りは、得られる陽性的中度の大きさが不明であり、得られる陽性的中度の大きさに関わらず±10%または±5%以内になるように陽性的中度 50%(最も信頼区間の幅が広がる状況)を仮定した時のものである。より高い陽性的中度が得られることが相当確実に期待できる場合は、より少ない症例数を設定することも許されるかも知れない。

なお、アウトカムの定義が複数パターン考慮された場合は、それぞれのパターンを満たす患者から 100 例ずつランダムサンプリングをすることが確実ではあるが、各パターンに重複があるとすれば、その分だけ調査症例数を多少減らせるかもしれない。

(D) アウトカム定義を満たす患者群と満たさない群からランダムサンプリングを行う場合

アウトカム定義を満たす患者群からだけでなく、満たさない患者群からも異なる抽出割合でランダムサンプルを実施し、陽性的中度のみならず感度(および特異度・陰性的中度)の推定につなげる方法である⁽⁵⁹⁾。この方法は、アウトカム定義を満たす患者群を予め定めた一定数サンプリングすることにより、得たい精度で陽性的中度を求めた上で、さらに精度は低くなっても感度についてもある程度の情報を得ようとする時に適している。ただし、アウトカムが稀な疾患の場合には、アウトカム定義を満たさない患者群に必要なとされるサンプルサイズは相当大きくなる。どの程度のサンプル数とするかは、研究の目的により異なりうる。この方法を用いる場合には、アウトカム定義を満たす患者群からは最低 100 例、アウトカム定義を満たさない患者群のサンプルサイズについては付録 4 を参照。なお母集団が大きい場合には、アウトカム定義を満たす患者群と満たさない群からそれぞれランダムサンプリングを行う前に、まず全集団からのランダムサンプリングを行う(または適切と考えられる範囲で、特定の期間など全集団の一部の情報を使う)ことで取り扱う母

集団のサイズを減らしデータのハンドリングを良くしておく作業も考慮される。

(E) “all possible cases”を想定したサンプリングを行う場合

研究対象集団(全体あるいはそのランダムサンプリング)の中で真のケースを多く含むことが期待されるサブグループを 2 種類以上(例:「レセプトに該当病名がある」サブグループと「関連する検査結果が正常値の○倍以上」のサブグループ)特定し、これらのサブグループの集合体を対象にカルテレビューで真のケースを見出し、これを「全ての真のケース」と仮定して、陽性的中度のみならず感度(および特異度・陰性的中度)も推定する方法である。この方法において、真のケースを全て含んでいると考えられるサブグループの集合体を Widdifield は“all possible cases”と呼んでいる⁽³⁸⁾。同様の方法で真のケースを全て(または、その大半を)特定しようとする試みはこれまでのバリデーション研究でもみられる^(13, 25, 60-64)。この方法では、複数のサブグループの一つに「アウトカム定義を満たす」グループを通常含めるが、まず、集団全体に含まれる真のケース(と考えられる)全例を特定し、その上でアウトカムの定義を満たすか否かで全患者を対象に 2×2 テーブルを作成するという点では(A)または(B)と共通する。“all possible cases”の方法では、いずれかのサブグループに属する全例についてカルテレビューを実施する。母集団が大きい場合、適切と判断されれば、まず母集団を代表するより小さな集団を(調査期間を短くするか、ランダムサンプリングによって)設定し、その集団に含まれる該当者全員に関してカルテレビューを実施し、最終的に 100 例以上の真のケースを特定することを目標とする。

なお、“all possible cases”の方法は、集団全体またはそのランダムサンプルの調査ではないので、「検討の対象となったサブグループの集合体に含まれない真のケース」の存在を否定することは通常できない。特に、設定可能なサブグループに含まれない真のケースが相当数存在することが確実であるような場合には、“all possible cases”の方法は、感度が過大評価され、アウトカム定義が「実際よりもよい」との誤った印象を与えることにつながるため推奨できない。設定されたサブグループから漏れる真のケースがそれほど多いとは考えられない場合にも、“all possible cases”の方法を用いた場合には、「検討の対象となったサブグループの集合体に含まれない真のケース」の程度が、推定された指標にどのような影響を与えるかについての感度解析をしておくことが望ましい⁽⁶²⁾。

(7)妥当性を測定する指標の計算

妥当性(validity)を測定する指標には、感度・特異度・陽性的中度・陰性的中度が含まれる(前述の表 1 を参照)。臨床疫学の教科書でとりあげられる医学検査の妥当性の検討では、疾患のスクリーニングから確定診断へのプロセスなどを念頭に、感度と特異度が重要視されるが、医療情報データベースにおいて相対リスクの推定を主な目的として実施される薬剤疫学研究においてアウトカム定義の妥当性を検討する際には、まず高い陽性的中度を証明することを最優先にしている研究が多い。その理由は、陽性的中度が高いアウトカム定義の特異度は通常高く(理由:表 1 において陽性的中度 $[a/a+b]$ が高い=偽陽性 $[b]$ が少ない=特異度 $[d/b+d]$ が高い)、また「比較される 2 群(例:糖尿病薬 A 群 vs.糖尿病薬 B 群)間のアウトカムに対する相対リスクは、アウトカム定義の特異度(陽性的中度)が高く、感度が(その大小に関わらず)2 群間で同じ(非差異誤分類の場合)に正しく求まる」(“with perfect specificity, nondifferential sensitivity of disease

misclassification will not bias the risk ratio”(Modern epidemiology 3rd edition)⁽⁶⁵⁾)からである。アウトカム定義の感度が2群間で大きく異なることが考えにくいのであれば、陽性的中度が高いアウトカム定義を用いることが正しい結果を得るための必要条件となる。一方、陽性的中度が高いのみならず感度も高い指標が見いだされれば、研究対象集団における有病割合(prevalence)または発生割合/率(incidence proportion/rate)が正しく推定でき、2群間の全体リスク差/率差を求めることができる。なお、一般的に感度と特異度・陽性的中度はトレードオフの関係にあり、高い陽性的中度を目指してアウトカム定義を厳しく設定する場合には、偽陰性が増え感度が下がる可能性があるため、感度の評価は可能な限り同時に行うことが好ましい。また、陽性的中度は対象集団における有病率の影響を受けるため、特異度の評価も行うことが理想である。

また、本タスクフォースでは、レビューを行った過去の約100文献にもとづき、妥当性の指標として感度・特異度・陽性的中度・陰性的中度を求めることを基本と捉えたが、診断疫学の分野では併せて診断オッズ比(=陽性尤度比/陰性尤度比=(感度×特異度)/((1-感度)×(1-特異度)))やROC曲線下面積も計算されることが多い。バリデーション研究においてもこれらの指標は有用である可能性があり⁽⁶⁶⁾、実際に計算している研究も存在する⁽⁶⁷⁾。よって、これらの指標も感度・特異度・陽性的中度・陰性的中度と同時に計算し、その有用性について検討することは、今後のバリデーション研究の課題の1つと考えられる。

以下に「(6)サンプリング方法の選択とサンプルサイズの設定」の項で示した各サンプリング方法(A)~(E)に対応した各指標の計算方法を示す。

(A) 全患者を対象にした場合、または(B)全患者からランダムサンプリングを行った場合

全患者またはランダムサンプリングされた患者を対象に表1のような2×2テーブルを作成することにより指標の全てを求めることができる。

(C)アウトカム定義を満たす患者群からのみランダムサンプリングを行った場合

2×2テーブルは作成できず、陽性的中度(=アウトカム定義及びゴールドスタンダードの定義を満たす患者数/アウトカム定義を満たす患者数)のみ求めることができる。

(D)アウトカム定義を満たす患者群と満たさない群からランダムサンプリングを行った場合

サンプリングのウェイトを考慮した2×2テーブルを作成することにより感度・特異度を正しく計算できる。つまり、サンプリングした患者のカルテレビュー結果をそのまま2×2テーブルにするのではなく、アウトカム定義を満たす患者群からのサンプリング割合を逆数にして(例:20%のサンプリングをしていたとしたら5倍)掛けた患者数、アウトカム定義を満たさない患者群からのサンプリング割合を逆数にして(例:5%のサンプリングをしていたとしたら20倍)掛けた患者数から2×2テーブルを作成することにより、正しい感度・特異度が求まる(詳細は表2を参照)。

表 2. アウトカム定義を満たす患者群からのサンプリング割合と満たさない患者群からのサンプリング割合が異なる場合の感度・特異度の求め方

シナリオ:全5,000症例のうち、アウトカム定義を満たす患者500人から20%(100人)のサンプリングを行い、満たさない患者4500人から5%(225人)のサンプリングを行い、カルテレビューを行った。その結果、アウトカム定義を満たす患者100人のうち90人、満たさない患者225人のうち45人がそれぞれ真のケースであった。

1) サンプルングのウェイトを考慮していない誤った感度・特異度の求め方

		ゴールドスタンダードの定義を満たすか (真のケースか)		合計	
		Yes	No		
アウトカムの定義を満たすか (陽性のケースか)	Yes	90 (a)	10 (b)	100	陽性的中度 (PPV) = 90/100 = 90% (95% CI: 82.4-95.1%)
	No	45 (c)	180 (d)	225	陰性的中度 (NPV) = 180/225 = 80% (95% CI: 74.2-85.0%)
合計		135	190	325	
		感度 = 90/135 = 66.7%	特異度 = 180/190 = 94.7%		

2) サンプルングのウェイトを考慮した正しい感度・特異度の求め方

		ゴールドスタンダードの定義を満たすか (真のケースか)		合計	
		Yes	No		
アウトカムの定義を満たすか (陽性のケースか)	Yes	450	50	500	陽性的中度 (PPV) = 450/500 = 90% (95% CI: 82.4-95.1%) **
	No	900	3,600	4,500	陰性的中度 (NPV) = 900/4,500 = 80% (95% CI: 74.2-85.0%) ***
合計		1,350	3,650	5,000	
		感度 = 450/1350 = 33.3% (95% CI: 26.3-42.3%)* (95% CI: 26.2-42.2%) [¶]	特異度 = 3,600/3,650 = 98.6% (95% CI: 97.8-99.5%)* (95% CI: 97.7-99.4%) [¶]		

CI = 信頼区間

*以下の近似式による推定

感度 (sen) の 95% 信頼区間: $sen * \exp[\pm 1.96(1-sen)\sqrt{(1/a+1/c)}]$

特異度 (spe) の 95% 信頼区間: $spe * \exp[1.96(1-spe)\sqrt{(1/b+1/d)}]$

詳細については付録 4 参照

[¶] Bootstrap 法による推定 (10,000 bootstrap sample から求めた percentiles)

**実際の検討数 100 例 (全 500 例の 20%) から計算

***実際の検討数 225 例 (全 4,500 例の 5%) から計算

(E) “all possible cases”を想定したサンプリングを行った場合

最終的に全患者を想定した 2×2 テーブルを作成することを念頭に、求められるところから順に各マスを埋めていく作業が必要である(詳細は表 3 を参照)。

表 3. “all possible cases”を想定したサンプリングを行った場合の 2×2 テーブルの作り方

シナリオ: 全 10,000 症例から 2000 例をランダムサンプリングした。真のケースを多く含むことが期待されるサブグループの集合体として 300 人特定された(その中には、アウトカム定義を満たす患者群をサブグループの1つとして含めており、200 人いた)。300 人をカルテレ뷰した結果、240 人が真のケースであり、このうち 180 人がアウトカム定義を満たす患者であった。

		ゴールドスタンダードの定義 を満たすか (真のケースか)		合計	
		Yes	No		
アウトカムの定義を 満たすか (陽性のケースか)	Yes	180***	20	200*	陽性的中度(PPV) = 180/200= 90% (95% CI: 85.0–93.8%)
	No	60***	1,740	1800*	陰性的中度(NPV) = 1,740/1,800 = 96.7% (95% CI: 95.7–97.4%)
合計		240**	1,760**	2,000	
		感度 = 180/240 = 75.0% (95% CI: 69.0–80.3%)	特異度 = 1,740/1,760 = 98.9% (95% CI: 98.3–99.3%)		

*~***各マスを埋めていく順序:

*2,000 症例のうち、アウトカム定義を満たす患者が 200 人。よって満たさない患者は 1,800 人と計算される。

**300 人につきカルテレ뷰をしたところ真のケースは 240 人。よって真でない(偽の)ケースは 1,760 人と計算される。

***カルテレ뷰で真のケースであった 240 人のうち 180 人が陽性のケースであったことから、60 人が偽陰性のケースと計算。(その結果、残りのマスも自動的に埋まる。)

なお、本項目に関する注意点として、すでに何度か言及してきたことではあるが、現状の日本で可能な病院単位のバリデーション研究の患者集団における指標の値は、利用するデータベース(例:レセプトデータ)にそのまま当てはまらない可能性が高いことに留意すべきである。一般に、集団に含まれる特定の疾患の割合(有病割合)が異なれば、感度、特異度が同一であっても陽性的中度や陰性的中度が異なることが知られているが、感度、特異度自体も集団によって異なりうる⁽⁶⁸⁾。たとえば疾患の重症度が偽陰性の割合に影響するのなら、感度(および陽性的中度と陰性的中度)は重症度の分布が異なる二つの集団では異なりうる。また、当該疾患の診断コードを誤つ

て与えられる患者(例:類似する疾患、疾患の初期で診断基準を満たさないが将来当該疾患の患者として診断される可能性が否定できない患者)の分布が集団によって異なるのなら、特異度(および陽性的中度と陰性的中度)も集団によって異なりうる。したがって、病院単位のバリデーション研究では、単に感度や特異度を求めるにとどまらず、感度などに影響を与える疾患の重症度などの因子を明らかにし、その分布と重症度ごとの指標の値を示すべきである。同様に、特異度に影響を与える類似疾患などの存在を明らかにし、その分布と特異度などへの影響を示すべきである。また、重症度については、データベース研究でも利用可能な、治療薬、合併症などで区分した重症度を示し、データベース研究における感度などが、バリデーション研究において求められた感度などどどのように異なるかを予測できるようにすべきである。同様に偽陽性となりやすい疾患などについても、データベース研究においても特定可能な特徴を示すことにより、データベース研究における指標の値の予測に利用可能とすべきである。

理解促進のために架空の例を以下に示す。ある1病院で、レセプトの関節リウマチ病名の陽性的中度を求めるバリデーション研究を行ったと仮定する。カルテレビューをゴールドスタンダードとし、アウトカム定義を満たす群(関節リウマチ病名が付いている患者)から100例のランダムサンプリングを行った。その結果、100人中70人が真のケースと特定され、陽性的中度は70%と計算された。しかし、当病院における関節リウマチ患者の重症度の分布は、レセプトデータベースにおける関節リウマチ患者の重症度の分布と異なる可能性が懸念された。そこでサブグループ解析として、(i)生物学的製剤使用患者、(ii)(生物学的製剤以外の)抗リウマチ薬使用患者、(iii)抗リウマチ薬非使用患者の3群に分けた。その結果、サブグループの構成は(i)群:(ii)群:(iii)群 = 20:60:20であり、陽性的中度は(i)群で90%、(ii)群で70%、(iii)群で50%であった。一方、レセプトデータベースにおけるサブグループの構成は(i)群:(ii)群:(iii)群 = 5:50:45と判明した。つまり、陽性的中度の高かった(i)群よりも、陽性的中度の低かった(iii)群の比率がレセプト全体において高かった。ここで、各サブグループの陽性的中度がバリデーション研究を行った病院とレセプト全体の間で等しかったという(大きな)仮定を置くと、レセプトデータベースでの陽性的中度 $= (5 \times 90 + 50 \times 70 + 45 \times 50) / 100 = 62\%$ と、当病院で求めた陽性的中度70%よりも低いことが予測できる。

このようなサブグループ解析は有用ではあるが、特にカルテレビューの結果をゴールドスタンダードとするバリデーション研究では、サンプル数は通常多くなく、サブグループの数も2つまたは3つ程度にとどめることが現実的であり、1グループに相当数(最低10例以上程度)含まれていることが必要である。したがって、このようなサブグループ解析とともに、偽陽性、偽陰性の患者について詳細に検討し、結果を報告することを本タスクフォースとしては推奨する^(13, 20-29)。偽陽性、偽陰性の患者がバリデーション研究を実施した施設に特異的なのか、あるいは、どのような医療機関でも起こりうると思われるかなどを考察することにより、特定の医療機関において実施されたバリデーション研究の結果をその医療機関の患者が代表しているとは考えられないデータベース研究において利用する上で有用な示唆を与えらる。

(8) 指標の閾値・利用法に関する考察

最後に「求めた陽性的中度や感度・特異度がいくつ以上であれば、そのアウトカムの定義が今後の薬剤疫学研究に使えるか?」という疑問について絶対的な基準の設定は困難である。主な理由として、各研究の目的により閾値は異なるからである。例えば、仮説の検証のための研究では

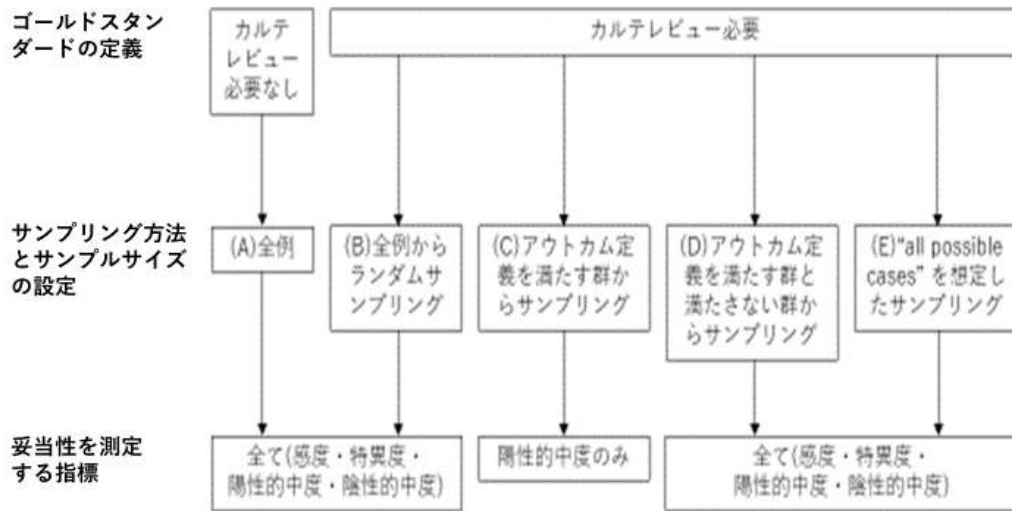
正しい相対リスクが求まるよう高い陽性的中度が要求される一方で、仮説の生成のための(探索的な)研究では多くのアウトカムを把握できるよう高い感度が望まれるかもしれない。また、2 群間の絶対リスク差やアウトカムの有病割合 (prevalence) または発生割合 / 率 (incidence proportion/rate) を求めたい場合も高い感度が要求される。

海外のバリデーション研究では、末期肝疾患の患者を見出すためのアルゴリズムの陽性的中度が 80%以上であることが望ましいとする論文⁽⁶⁹⁾、COPD の患者を見出すためのアルゴリズムの陽性的中度が 85%以上であることが望ましいとする論文⁽⁵⁴⁾が見られた。ただし、前者の「80%」は末期肝疾患の患者における薬の効果の評価するために患者群を特定するという目的に沿って設定された値であり、後者の「85%」は COPD 患者の治療状況や治療の質、その他 COPD 患者に関連する臨床的なクエスチョンを検討するために患者群を特定するという目的に沿って設定された値であることに留意する必要がある。一方、ある値以下であれば、指標として使うことはできないとの言及も見られた。例えば、アナフィラキシーのバリデーション研究において、陽性的中度 70%以下の結果に対して「今回のアウトカムの定義は、あくまで候補を絞りこむ手段に留まり、真のケースの確定にはその後カルテレビューを要する」と考察しているものが見られた⁽³⁷⁾。

一つのバリデーション研究の中で複数のアウトカム定義のパターンまたはアルゴリズムを検討する際は、陽性的中度と感度の相対的なバランスが、最適なアルゴリズムを決定する一助となる。一般的には、アルゴリズムの条件を厳しくするほど陽性的中度(または特異度)が上がるが、感度(あるいはデータベースの中で特定できる患者数)が下がるという、トレードオフ関係が認められる。過去のバリデーション研究では陽性的中度を下げない範囲で感度ができるだけ高くなるアルゴリズムを推奨していることが多いが^(39, 40)、研究者らが考える最適なアルゴリズムについての考察を論文内に記載しておくことが望ましい。

まとめとして、以上で述べたバリデーション研究の手順・チェック項目を一覧表の形にした(付録 5)。バリデーション研究のプロトコール作成・研究の遂行・論文執筆の際に参照されたい。また、上記(1)~(8)の項目のうち、「(5)ゴールドスタンダードの定義」、「(6)サンプリング方法とサンプルサイズの設定」、「(7)妥当性を測定する指標の計算」には対応関係があることはあらためて意識しておきたい(図 1)。

図 1.ゴールドスタンダードの定義、サンプリング方法とサンプルサイズの設定、妥当性を測定する指標、の対応関係。3(バリデーション研究実施の手順・チェック項目)の(5)、(6)、(7)の相互関係を示す。



4. 今後日本で、より望ましいバリデーション研究を実施するために

最後に、日本におけるバリデーション研究の現状を総括し、今後より望ましいバリデーション研究を実施するために必要な知識・事項・提言についてまとめた。

(1) 日本におけるバリデーション研究の現状

日本のバリデーション研究は 2010 年頃から実施されるようになり、2013 年以後に論文がいくつか発表されている。

2013 年の大場らの論文⁽⁷⁰⁾は、(株)日本医療データセンターのレセプトデータと資格喪失理由を含む被保険者のデータを用い、レセプト上の転帰「死亡」とレセプト情報から求めた Charlson Comorbidity Index を組み合わせ、保険者の保有する資格喪失理由「死亡」をゴールドスタンダードとするバリデーション研究であり、レセプト情報から得られる情報による「死亡」の定義が感度 60%程度、陽性的中度 95%をもつことなどを明らかにした。

2015 年の佐藤らの論文⁽¹¹⁾は、単一の病院におけるレセプトデータから新規に発生した乳がんを定義し、院内がん登録をゴールドスタンダードとしてその妥当性を検討するバリデーション研究であり、傷病名と治療に関する情報を組み合わせた定義が感度 90%、陽性的中度 87%などをもつことを明らかにした。

2015 年の山口らの論文⁽⁷¹⁾では、DPC 参加病院で EBM Provider((株)メディカル・データ・ビジョン)のうち、検査値の得られる 16 病院のデータを用い、データベースから利用可能な患者情報、病名一覧、薬剤、処置、検査による専門医の判断をゴールドスタンダードとして、ICD10 コードと薬、診療行為の組み合わせによる静脈血栓塞栓症と出血性イベントのアウトカム定義の陽性的中度がそれぞれ 75.0%、73.3%であることを明らかにした。

2016 年の田中ら⁽⁷²⁾の論文は 3 つの病院でレセプト情報から得られる非定型大腿骨折の診断コードを電子カルテから得られる情報をゴールドスタンダードとして評価するバリデーション研究であり、感度は 82%前後、特異度は 100%であったことを明らかにした。

2017 年の山名ら⁽⁷³⁾の論文は、4 つの DPC 病院でランダムサンプルした 315 人の患者の DPC データの 16 の疾患、10 の診療行為、SS-MIX データの 13 の臨床検査結果を、カルテレビューをゴールドスタンダードとして評価するバリデーション研究であり、DPC データ、SS-MIX のデータの妥当性は将来の研究に使う上で十分な妥当性を有していることを明らかにした。

(2) 行政が主導するバリデーション研究

MID-NET®は、医薬品の安全対策を目的として、厚生労働省と PMDA の医療情報データベース基盤整備事業において構築されたデータベースシステム及び関連ネットワークの総称である。厚生労働省が公募により選定した協力医療機関 10 拠点 23 病院における、2009 年以降のレセプトデータ、DPC データ、及び検査結果値を含む電子カルテデータ(以下、あわせて「HIS (Hospital Information System) データ」という。)が含まれており、データベース規模としては 2018 年までに約 400 万人になる見込みである。MID-NET®は分散型データベース構造であり、各拠点にデータベース(以下、「統合 DS」(統合データソース))が設置されている。各拠点の統合 DS には、SS-MIX2 の規約に基づき送信された HIS データが移行されており、この移行の過程で、傷病名、検体検査名、

薬剤名等に対してマッピング・標準化処理(医療機関で利用されている院内コード及び名称から標準コード及び名称への変換、臨床検査値と単位等の標準化など)が行われている。利用者がMID-NET®を解析する場合、作成したプログラムをオンサイトセンターから専用回線を通じて対象医療機関に送付すると、医療機関の担当者が依頼内容を承認することで対象となるデータが連結不可能匿名化処理を経て抽出され、解析処理が実行される仕組みとなっている。この抽出データ及び解析結果は医療機関の担当者による承認を経て複数施設統合処理センター(以下、「センター」)に送付される。利用者はセンターに送られてきたデータをシステム上で更に解析することが可能であるが、結果の持ち出しについては統計処方を実施した集計結果についてのみ可能である。

2017年度より、日本医療研究開発機構の医薬品等規制調和・評価研究事業として、MID-NET®協力医療機関におけるアウトカムバリデーション研究が実施されている。各医療機関内において、候補とするアウトカム定義に基づき抽出された集団からランダムサンプリングした患者について、MID-NET®には含まれていない診療録情報を確認し、各拠点共通の判定基準に基づいて、その領域の専門医が真偽を判定する。対象となるアウトカムは、医薬品のリスク管理計画に記載されている安全性検討事項の上位のものから選定され、3年間で約20のアウトカムに対して実施される予定である。

(3) DPC レセプトのアウトカム定義への利用

日本のレセプトデータの利用を意図して行うバリデーション研究において、DPC レセプトの傷病名は医科レセプトの傷病名よりも特異性が高いことが期待され、今後のバリデーション研究でその妥当性が評価されるべきである。その際、DPC レセプトの構造についてよく理解しておくことが重要と考えられる。

本邦では2003年よりDPCに基づく、包括支払制度(DPC/PDPS)が導入された^(74,75)。DPC参加病院ではDPC導入の影響評価に関わる調査のために3カ月に1度、厚生労働省に提出されるDPCデータおよび医療費償還のためのDPCレセプトとDPCに包括されなかった患者の医科レセプトが作成される。DPC準備病院ではDPCデータおよび医科レセプトが作成され、それ以外の病院ではそれまで通り医科レセプトのみ作成される。DPCレセプトはDPCデータのDファイルに相当し、傷病コード(7桁)、ICD-10コード、傷病名区分コードが記録され、DPCデータの様式1と同等の疾病情報が得られる。さらに、医科レセプトと同様、主傷病フラグ、診療開始日が記録される。但し、主傷病フラグは任意の入力項目であることから、バリデーション研究のアウトカム定義では、傷病名区分がより有用と考えられる。DPCレセプトにおける傷病名区分は、「傷病名(医療資源を最も投入した傷病名)」「副傷病名」「主傷病名」「入院の契機となった傷病名」「医療資源を2番目に投入した傷病名」「入院時併存傷病名」「入院後発症傷病名」の7項目で構成されるため、バリデーション研究のアウトカム定義で有用と考えられる。また、DPCレセプト、医科レセプトにはその病院における当該疾病に係る診療開始日が記録されるため、曝露との時間的前後関係を特定できる。但し、複数の病院の電子カルテやその病院で発行されたレセプト情報を集積したデータベースでは、病院が変わると患者の追跡ができないため、当該病院における診療開始日より前に疾患が発生している可能性があり、疾患の発生日の推定には注意が必要である。さらに、DPC/PDPSは急性期の疾患を対象とし、退院時に医療資源病名が確定し、適応される入院期間に上限があ

るという特徴がある。よって、長期の入院を伴う疾患など慢性期の傷病名のバリデーションにはDPCレセプトは不向きであり、医科レセプト情報との併用が必要である。一方、バリデーション研究においてアウトカムに対するアルゴリズムが組まれる場合には、手術、検査、処置の情報も重要となる。DPCレセプトにはこれら診療行為の情報が含まれるので、それぞれの情報をDPCコードから特定することができる。また、DPCコード自体をアウトカムや研究対象集団を決めるアルゴリズムとして利用することも可能かも知れない。DPCコードは、疾病と治療内容の組み合わせで患者を分類したもので、14桁のコードで成り立っており、同一のコードは同一の疾患であるだけでなく、治療内容も同一であることを意味するからである⁽⁷⁶⁾。

(4) 大規模コホート研究における国保レセプトを用いたバリデーション研究の可能性

バリデーション研究は、単施設もしくは複数施設で実施されることもあるが、外的妥当性の観点からは、研究の母集団を代表する集団でバリデーション研究を行うことが望ましい。国内の複数の地域住民を対象としたpopulation-basedコホート研究で、国民健康保険(国保)レセプト、後期高齢者医療制度(後期高齢医療)レセプトなどのレセプト情報を収集し、かつ、悉皆的な疾病登録が実施されている場合には、外的妥当性の高いバリデーション研究が実施可能である。

2010年以降に開始された、第3世代の大規模コホート研究の次世代多目的コホート研究(JPHC-NEXT)や東北メディカル・メガバンク機構(ToMMo)では、研究参加者のレセプト情報も収集している⁽⁷⁷⁾。このようなコホート研究では、レセプト情報と全国がん登録などの悉皆的な疾病登録とリンケージ(照合)させることにより、バリデーション研究が可能となる。全国がん登録の研究利用では、原則本人からの同意が必要であるため、全国がん登録を用いたアウトカム定義のバリデーション研究実施は困難だが、JPHC-NEXTやToMMoではこの同意の有無に関する情報も含まれている。代表性のある大規模なコホート研究でバリデーション研究が実施されれば、レセプト情報でがんなどのアウトカム定義の評価において有用な情報となることが期待される。悉皆的な疾病登録がない場合であっても、基幹病院への入院を要するような重篤疾患であれば、コホート調査対象地域の基幹病院で入退院サマリ、病歴(紙カルテや電子カルテ)を参照することにより、バリデーション研究が実施可能と考えられる。また、血液検査値や血圧値などで疾病の有無が判別可能な糖尿病、高血圧、脂質異常症、慢性腎機能障害などの慢性疾患では、健康診断結果や病院受診時の臨床検査値を参照すればバリデーション研究ができると考えられる

大規模コホート研究におけるバリデーション研究にも限界がある。第1に、現時点では、国保レセプトや後期高齢者医療レセプト以外のレセプト情報の収集は困難であり、その他(例:協会けんぽ)のレセプト情報を用いたバリデーション研究は実施できない。国保加入者や後期高齢者医療受給者とその他の保険加入者では、背景因子が大きく異なる可能性がある点に注意が必要である。第2に、上記のようなコホート研究の対象者数は約10万人で「大規模」であるが、希少疾患のバリデーション研究は困難である。第3に、背景因子やレセプト病名登録に地域差があることが予想される場合には、コホート対象地域のバリデーション研究結果をその他の集団への一般化可能性には限界がある。

このような限界に留意しつつも、大規模コホート研究におけるバリデーション研究が実施され、レセプト情報の妥当性に関するエビデンスが集積することに期待したい。

(5) 産官学・医療機関における理解の促進

2017年のGPSP省令の改正に伴い、医療情報データベースの活用が製造販売後調査の一つの方法として認められるに至った。バリデーション研究の必要性については、PMDAによって平成26年3月31日に発表された「医療情報のデータベース等を用いた医薬品の安全性評価における薬剤疫学研究の実施に関するガイドライン」や、平成29年6月9日に厚生労働省の2課長通知として発出された「製造販売後の医薬品安全性監視における医療情報データベースの利用に関する基本的考え方について」に一部盛り込まれている。しかし、より高いエビデンスレベルが求められる場合にはバリデーション研究の実施が必要であるという認識はおろか、バリデーション研究という用語自体、馴染みがない産官の関係者が多い。本報告書が、バリデーション研究の理解の促進の一助となることを期待したい。製薬産業をはじめデータベース研究を実施する民間企業においても、バリデーション研究の意義を理解し、その推進に協力するべく、業界内での啓発活動が求められる。医療機関において、バリデーション研究の意義、特に診療報酬請求のために発行されるレセプトのバリデーション研究がデータベース研究の質向上に必要であることへの理解促進も重要である。レセプトは本来医療費償還のための請求書であり、「請求書として正確」であることは「医学的に正確」であることとは合致しない。たとえば claims database を含む administrative database では低ナトリウム血症⁽⁴⁷⁾や血清クレアチニン値によって定義される急性腎障害⁽¹⁸⁾の感度が低いことが知られているが、これは低ナトリウムや血清クレアチニン上昇自体が主たる診断(したがって医療費償還の対象)となることが多くないことを反映している。たとえ「請求書として正確」であっても、バリデーション研究における特異度や感度は100%にはなりえないことへの理解促進が必要である。学(アカデミア)にはバリデーション研究への医療機関の協力を得つつ、社会的重要性も考慮しながらバリデーション研究を現場で推し進める役割を果たすことが求められる。

(6) 個人を識別するためのリンケージ(照合)

米国における Medicare、Medicaid、民間保険由来の医療費償還請求書(claims)や Health maintenance organization (HMO)などが有する Electronic Health Record のデータベースのデータは、(研究のために提供されるデータセットに含まれるかどうかは別として) Social Security Number などの個人ごとに付与された ID (個人 ID) とともに管理されている^(20, 53, 78)。北欧⁽⁴⁷⁾やカナダ^(49, 54)における各種レジストリーや退院時サマリのデータベースのデータも個人 ID とともに管理されており、バリデーション研究を目的に個人 ID を用いたリンケージ(照合)が行われることがある。特に北欧では、バリデーション研究に限らず、個別の(薬剤疫学)研究ごとに、複数のデータベースを個人 ID でリンケージ(照合)し、研究用の新たなデータセットを作成して研究者に提供することが恒常的に実施されている。これに対し、日本ではデータベースのリンケージ(照合)にもとづくバリデーション研究の実施は困難である。理由の一つとして、研究者等が利用可能なデータベースの多くはそれぞれで匿名化されたデータであり、また日本には、異なるデータベース間で同一人物を紐付けるために利用可能な個人 ID が存在しないことが上げられる。また、2017年5月30日に全面施行された改正個人情報保護法にて新たに導入された匿名加工情報の利用にあたっては、同法第36条および第38条にて「本人を識別するために、当該匿名加工情報を他の情報と照合してはならな

い」旨が定められている。

2017年に公布された「医療分野の研究開発に資するための匿名加工医療情報に関する法律」(次世代医療基盤法)は通常では許されない個人情報連結を可能とするものであり、国の認定を受けた「認定匿名加工医療情報作成事業者」が異なるタイプのデータベース同士をリンケージ(照合)できる可能性がある。たとえばレセプトデータと疾患レジストリーをリンケージ(照合)することが可能なら、研究の質の大幅な向上が期待できる。特にレセプトデータと複数の疾患レジストリーをリンケージ(照合)できれば、疾患レジストリーAから特定された疾患Aをもつ患者集団において(レセプトデータで特定可能な)薬の曝露がどのように疾患Bの発生というアウトカムを起こすかを疾患レジストリーBで確認する、などの研究が可能になるかもしれない。

しかし、次世代医療基盤法の下においても、バリデーション研究、特にカルテレビューによるバリデーション研究は必ずしも容易ではないと考えられる。たとえば、その第18条では、複数の医療機関などから個人情報を得て匿名加工医療情報を作成する「認定匿名加工医療情報作成事業者」について、「本人を識別するために、当該匿名加工医療情報を他の情報と照合してはならない」旨が定められている。ただし、運用のあり方によっては、次世代医療基盤法はカルテレビューによるバリデーション研究を含む様々な形態のレコードリンケージの根拠となりうるとも考えられる。今後、次世代医療基盤法の下でどのようなリンケージ(照合)が行われていくのかに注目したい。

その他、「人を対象とする医学系研究に関する倫理指針」の「第17 匿名加工情報の取り扱い」では匿名加工情報を作成またはその提供を受けた研究者は「本人を識別するために、当該匿名加工情報を他の情報と照合してはならない」旨が定められている。ただし、「人を対象とする医学系研究に関する倫理指針」で言及されている「匿名加工情報」は改正個人情報保護法における匿名加工情報取扱事業者による商業的データベースなどの「匿名加工情報」に限定されることに留意する必要がある。しかし、現状に日本では匿名加工情報取扱事業者による「匿名加工情報」以外のデータベースについても、リンケージ(照合)は強く制限されている。たとえば「レセプト情報・特定健診等情報(NDB)の提供に関するガイドライン」でも、その「第5 レセプト情報等の提供依頼申出手続」において「有識者会議が特に認めた場合を除き、提供されたその他の個体識別が可能となる可能性があるデータ(別の利用目的で提供されたその他のレセプト情報等を含む)とのリンケージ(照合)を行わないこと」とされており、「第6 提供依頼申出に対する審査」の「4 審査基準」でも「特定個人を識別することを内容とする分析方法、手法も認めない」とされている。

バリデーション研究、特にカルテ情報をゴールドスタンダードとして実施するバリデーション研究は、匿名化されたデータベースに含まれる情報の正確性を評価するために、データベース内の個人の特定を不可欠のプロセスとする研究である。しかし現在の日本では、データベース内の個人の特定を不可欠とする研究が存在すること、さらに、その研究がデータベース研究そのものの質の向上に必須であるとの認識が十分とは言い難い。我が国のデータベース研究の質を担保するためには、本人を識別するために他の情報とリンケージ(照合)することを認め、適切なバリデーション研究を実施できる環境の整備は必要不可欠である。

5. 結論

バリデーション研究その他により、データベース研究を十分質の高い研究とすることは可能であ

り、そのようなデータベース研究は医薬品などの安全性・有効性に関するその他の方法では困難な信頼できる知見を得て、国民の健康と福祉に貢献しうる。また厳密なプライバシー保護と両立しうる方法でリンケージ(照合)を認める仕組み作りは可能である。国では 2017 年に厚生労働大臣を本部長とする「データヘルス改革推進本部」を立ち上げ、「健康医療介護分野の工程表」をまとめ、医療等 ID の導入を踏まえ、データの連結を推し進める方針でデータ連結の調査・研究も計画されている。その成果としてカルテなどを用いるバリデーション研究を含むデータのリンケージ(照合)の可能性が広がることにも期待したい。今後、データベースの適切な利用を進めようとする者は、これらの点に関する国民的理解を深め、バリデーション研究を実施しやすい環境を実現してデータベース研究の有用性を高めることに関しても努めるべきであろう。

バリデーション研究は、データベース研究のより適切な結果の解釈を可能にするが、その方法論はデータベース研究と同様に固定なものではない。現実的には、研究に用いたデータベースの種類、そのデータベースに課せられている法的および技術的制限、研究で利用できる人的および金銭的リソースなどの条件を考慮しつつ、最適な方法はそれぞれの研究ごとに熟慮されるべきものである。

この報告書がバリデーション研究の意義や基本的な考え方の理解を促し、与えられた条件下で目的を達成するための最適な方法を見出す手掛かりになることを期待する。

参考文献

1. Strom BL. Overview of Automated Databases in Pharmacoepidemiology. In Textbook of Pharmacoepidemiology. Eds Strol BL, Kimme SE. John Wiley & Sons, Ltd. NJ, 2006, 167-71.
2. 小宮山靖、青木事成、古閑晃、久保田潔. より良い Pharmacovigilance Plan 策定に向けての提言. 薬剤疫学 2015;20:73-83.
3. Uneno Y, Taneishi K, Kanai M, Okamoto K, Yamamoto Y, Yoshioka A, et al. Development and validation of a set of six adaptable prognosis prediction (SAP) models based on time-series real-world big data analysis for patients with cancer receiving chemotherapy: A multicenter case crossover study. PloS one. 2017;12(8):e0183291. Epub 2017/08/25. doi: 10.1371/journal.pone.0183291. PubMed PMID: 28837592; PubMed Central PMCID: PMC5570326.
4. op den Winkel M, Nagel D, Sappl J, op den Winkel P, Lamerz R, Zech CJ, et al. Prognosis of patients with hepatocellular carcinoma. Validation and ranking of established staging-systems in a large western HCC-cohort. PloS one. 2012;7(10):e45066. Epub 2012/10/17. doi: 10.1371/journal.pone.0045066. PubMed PMID: 23071507; PubMed Central PMCID: PMC3465308.
5. 厚生労働省医薬食品局監視指導・麻薬対策課長. 医薬品・医薬部外品製造販売業者等におけるコンピュータ化システム適正管理ガイドラインについて. 平成 22 年 10 月 21 日. <http://ecompliance.co.jp/CSV/csguideline.pdf>.
6. Strom BL. Validity of Pharmacoepidemiologic Drug and Diagnosis Data. In Textbook of Pharmacoepidemiology. Eds Strol BL, Kimme SE. John Wiley & Sons, Ltd. NJ, 2006, 239-58.
7. ISPE. Guidelines for Good Pharmacoepidemiology Practices (GPP). <https://www.pharmacoepi.org/resources/policies/guidelines-08027/>.
8. Berger ML, Sox H, Willke RJ, Brixner DL, Eichler HG, Goettsch W, et al. Good practices for real-world data studies of treatment and/or comparative effectiveness: Recommendations from the joint ISPOR-ISPE Special Task Force on real-world evidence in health care decision making. Pharmacoepidemiology and drug safety. 2017;26(9):1033-9. Epub 2017/09/16. doi: 10.1002/pds.4297. PubMed PMID: 28913966; PubMed Central PMCID: PMC5639372.
9. Semins MJ, Trock BJ, Matlaga BR. Validity of administrative coding in identifying patients with upper urinary tract calculi. The Journal of urology. 2010;184(1):190-2. Epub 2010/05/19. doi: 10.1016/j.juro.2010.03.011. PubMed PMID: 20478584.
10. Gonzalez-Fernandez M, Gardyn M, Wyckoff S, Ky PK, Palmer JB. Validation of ICD-9 Code 787.2 for identification of individuals with dysphagia from administrative databases. Dysphagia. 2009;24(4):398-402. Epub 2009/04/29. doi: 10.1007/s00455-009-9216-1. PubMed PMID: 19399554.
11. Sato I, Yagata H, Ohashi Y. The accuracy of Japanese claims data in identifying breast cancer cases. Biological & pharmaceutical bulletin. 2015;38(1):53-7. Epub 2015/03/07. doi: 10.1248/bpb.b14-00543. PubMed PMID: 25744458.
12. MacFarlane LA, Liu CC, Solomon DH, Kim SC. Validation of claims-based algorithms for gout flares. Pharmacoepidemiology and drug safety. 2016;25(7):820-6. Epub 2016/05/28. doi: 10.1002/pds.4044.

PubMed PMID: 27230083; PubMed Central PMCID: PMC4930384.

13. Hsieh CY, Chen CH, Li CY, Lai ML. Validating the diagnosis of acute ischemic stroke in a National Health Insurance claims database. *Journal of the Formosan Medical Association = Taiwan yi zhi*. 2015;114(3):254-9. Epub 2013/10/22. doi: 10.1016/j.jfma.2013.09.009. PubMed PMID: 24140108.
14. de Achaval S, Feudtner C, Palla S, Suarez-Almazor ME. Validation of ICD-9-CM codes for identification of acetaminophen-related emergency department visits in a large pediatric hospital. *BMC health services research*. 2013;13:72. Epub 2013/02/26. doi: 10.1186/1472-6963-13-72. PubMed PMID: 23433397; PubMed Central PMCID: PMC4930384.
15. Vanderloo SE, Johnson JA, Reimer K, McCrea P, Nuernberger K, Krueger H, et al. Validation of classification algorithms for childhood diabetes identified from administrative data. *Pediatric diabetes*. 2012;13(3):229-34. Epub 2011/07/21. doi: 10.1111/j.1399-5448.2011.00795.x. PubMed PMID: 21771232.
16. Newton KM, Wagner EH, Ramsey SD, McCulloch D, Evans R, Sandhu N, et al. The use of automated data to identify complications and comorbidities of diabetes: a validation study. *Journal of clinical epidemiology*. 1999;52(3):199-207. Epub 1999/04/21. PubMed PMID: 10210237.
17. Youngson E, Welsh RC, Kaul P, McAlister F, Quan H, Bakal J. Defining and validating comorbidities and procedures in ICD-10 health data in ST-elevation myocardial infarction patients. *Medicine*. 2016;95(32):e4554. Epub 2016/08/12. doi: 10.1097/md.0000000000004554. PubMed PMID: 27512881; PubMed Central PMCID: PMC4985336.
18. Molnar AO, van Walraven C, McArthur E, Fergusson D, Garg AX, Knoll G. Validation of administrative database codes for acute kidney injury in kidney transplant recipients. *Canadian journal of kidney health and disease*. 2016;3:18. Epub 2016/04/09. doi: 10.1186/s40697-016-0108-7. PubMed PMID: 27057318; PubMed Central PMCID: PMC4823855.
19. Ingeman A, Andersen G, Hundborg HH, Johnsen SP. Medical complications in patients with stroke: data validity in a stroke registry and a hospital discharge registry. *Clinical epidemiology*. 2010;2:5-13. Epub 2010/09/25. PubMed PMID: 20865097; PubMed Central PMCID: PMC2943185.
20. Palmsten K, Huybrechts KF, Kowal MK, Mogun H, Hernandez-Diaz S. Validity of maternal and infant outcomes within nationwide Medicaid data. *Pharmacoepidemiology and drug safety*. 2014;23(6):646-55. Epub 2014/04/18. doi: 10.1002/pds.3627. PubMed PMID: 24740606; PubMed Central PMCID: PMC4205050.
21. Bowker SL, Savu A, Donovan LE, Johnson JA, Kaul P. Validation of administrative and clinical case definitions for gestational diabetes mellitus against laboratory results. *Diabetic medicine : a journal of the British Diabetic Association*. 2017;34(6):781-5. Epub 2016/10/16. doi: 10.1111/dme.13271. PubMed PMID: 27743395.
22. Franchi C, Giussani G, Messina P, Montesano M, Romi S, Nobili A, et al. Validation of healthcare administrative data for the diagnosis of epilepsy. *Journal of epidemiology and community health*. 2013;67(12):1019-24. Epub 2013/09/12. doi: 10.1136/jech-2013-202528. PubMed PMID: 24022813.
23. Metcalfe A, Sibbald B, Lowry RB, Tough S, Bernier FP. Validation of congenital anomaly coding in Canada's administrative databases compared with a congenital anomaly registry. *Birth defects*

- research Part A, Clinical and molecular teratology. 2014;100(2):59-66. Epub 2013/12/07. doi: 10.1002/bdra.23206. PubMed PMID: 24307632.
24. Stringer E, Bernatsky S. Validity of juvenile idiopathic arthritis diagnoses using administrative health data. *Rheumatology international*. 2015;35(3):575-9. Epub 2014/10/02. doi: 10.1007/s00296-014-3142-8. PubMed PMID: 25270916.
25. Dregan A, Moller H, Murray-Thomas T, Gulliford MC. Validity of cancer diagnosis in a primary care database compared with linked cancer registrations in England. Population-based cohort study. *Cancer epidemiology*. 2012;36(5):425-9. Epub 2012/06/26. doi: 10.1016/j.canep.2012.05.013. PubMed PMID: 22727737.
26. Furlan JC, Fehlings MG. The National Trauma Registry as a Canadian spine trauma database: a validation study using an institutional clinical database. *Neuroepidemiology*. 2011;37(2):96-101. Epub 2011/09/17. doi: 10.1159/000330835. PubMed PMID: 21921642.
27. Schwartz KL, Jembere N, Campitelli MA, Buchan SA, Chung H, Kwong JC. Using physician billing claims from the Ontario Health Insurance Plan to determine individual influenza vaccination status: an updated validation study. *CMAJ open*. 2016;4(3):E463-e70. Epub 2016/10/13. doi: 10.9778/cmajo.20160009. PubMed PMID: 27730110; PubMed Central PMCID: PMC5047797.
28. Coleman N, Halas G, Peeler W, Casaclang N, Williamson T, Katz A. From patient care to research: a validation study examining the factors contributing to data quality in a primary care electronic medical record database. *BMC family practice*. 2015;16:11. Epub 2015/02/05. doi: 10.1186/s12875-015-0223-z. PubMed PMID: 25649201; PubMed Central PMCID: PMC4324413.
29. Joseph KS, Fahey J. Validation of perinatal data in the Discharge Abstract Database of the Canadian Institute for Health Information. *Chronic diseases in Canada*. 2009;29(3):96-100. Epub 2009/06/17. PubMed PMID: 19527567.
30. Rudmik L, Xu Y, Kukec E, Liu M, Dean S, Quan H. A validated case definition for chronic rhinosinusitis in administrative data: a Canadian perspective. *International forum of allergy & rhinology*. 2016;6(11):1167-72. Epub 2016/11/04. doi: 10.1002/alr.21801. PubMed PMID: 27228224.
31. Ramalle-Gomara E, Ruiz E, Serrano M, Bartulos M, Gonzalez MA, Matute B. Validity of discharge diagnoses in the surveillance of stroke. *Neuroepidemiology*. 2013;41(3-4):185-8. Epub 2013/09/21. doi: 10.1159/000354626. PubMed PMID: 24051447.
32. Gabriel SE, Crowson CS, O'Fallon WM. A mathematical model that improves the validity of osteoarthritis diagnoses obtained from a computerized diagnostic database. *Journal of clinical epidemiology*. 1996;49(9):1025-9. Epub 1996/09/01. PubMed PMID: 8780612.
33. Husain N, Blais P, Kramer J, Kowalkowski M, Richardson P, El-Serag HB, et al. Nonalcoholic fatty liver disease (NAFLD) in the Veterans Administration population: development and validation of an algorithm for NAFLD using automated data. *Alimentary pharmacology & therapeutics*. 2014;40(8):949-54. Epub 2014/08/27. doi: 10.1111/apt.12923. PubMed PMID: 25155259; PubMed Central PMCID: PMC4331854.
34. Verkooijen HM, Fioretta G, Chappuis PO, Vlastos G, Sappino AP, Benhamou S, et al. Set-up of a

- population-based familial breast cancer registry in Geneva, Switzerland: validation of first results. *Annals of oncology : official journal of the European Society for Medical Oncology*. 2004;15(2):350-3. Epub 2004/02/05. PubMed PMID: 14760133.
35. Ogdie A, Alehashemi S, Love TJ, Jiang Y, Haynes K, Hennessy S, et al. Validity of psoriatic arthritis and capture of disease modifying antirheumatic drugs in the health improvement network. *Pharmacoepidemiology and drug safety*. 2014;23(9):918-22. Epub 2014/07/22. doi: 10.1002/pds.3677. PubMed PMID: 25044030; PubMed Central PMCID: PMC4149813.
 36. Harboe KM, Anthonen K, Bardram L. Validation of data and indicators in the Danish Cholecystectomy Database. *International journal for quality in health care : journal of the International Society for Quality in Health Care*. 2009;21(3):160-8. Epub 2009/03/24. doi: 10.1093/intqhc/mzp009. PubMed PMID: 19304994.
 37. Walsh KE, Cutrona SL, Foy S, Baker MA, Forrow S, Shoaibi A, et al. Validation of anaphylaxis in the Food and Drug Administration's Mini-Sentinel. *Pharmacoepidemiology and drug safety*. 2013;22(11):1205-13. Epub 2013/09/17. doi: 10.1002/pds.3505. PubMed PMID: 24038742; PubMed Central PMCID: PMC4113322.
 38. Widdifield J, Ivers NM, Young J, Green D, Jaakkimainen L, Butt DA, et al. Development and validation of an administrative data algorithm to estimate the disease burden and epidemiology of multiple sclerosis in Ontario, Canada. *Multiple sclerosis (Houndmills, Basingstoke, England)*. 2015;21(8):1045-54. Epub 2014/11/14. doi: 10.1177/1352458514556303. PubMed PMID: 25392338.
 39. Butt DA, Tu K, Young J, Green D, Wang M, Ivers N, et al. A validation study of administrative data algorithms to identify patients with Parkinsonism with prevalence and incidence trends. *Neuroepidemiology*. 2014;43(1):28-37. Epub 2014/10/18. doi: 10.1159/000365590. PubMed PMID: 25323155.
 40. Widdifield J, Bombardier C, Bernatsky S, Paterson JM, Green D, Young J, et al. An administrative data validation study of the accuracy of algorithms for identifying rheumatoid arthritis: the influence of the reference standard on algorithm performance. *BMC musculoskeletal disorders*. 2014;15:216. Epub 2014/06/25. doi: 10.1186/1471-2474-15-216. PubMed PMID: 24956925; PubMed Central PMCID: PMC4078363.
 41. Tu K, Wang M, Jaakkimainen RL, Butt D, Ivers NM, Young J, et al. Assessing the validity of using administrative data to identify patients with epilepsy. *Epilepsia*. 2014;55(2):335-43. Epub 2014/01/15. doi: 10.1111/epi.12506. PubMed PMID: 24417710.
 42. Schultz SE, Rothwell DM, Chen Z, Tu K. Identifying cases of congestive heart failure from administrative data: a validation study using primary care patient records. *Chronic diseases and injuries in Canada*. 2013;33(3):160-6. Epub 2013/06/06. PubMed PMID: 23735455.
 43. Williamson T, Green ME, Birtwhistle R, Khan S, Garies S, Wong ST, et al. Validating the 8 CPCSSN case definitions for chronic disease surveillance in a primary care database of electronic health records. *Annals of family medicine*. 2014;12(4):367-72. Epub 2014/07/16. doi: 10.1370/afm.1644. PubMed PMID: 25024246; PubMed Central PMCID: PMC4096475.

44. Kim SC, Gillet VG, Feldman S, Lii H, Toh S, Brown JS, et al. Validation of claims-based algorithms for identification of high-grade cervical dysplasia and cervical cancer. *Pharmacoepidemiology and drug safety*. 2013;22(11):1239-44. Epub 2013/09/13. doi: 10.1002/pds.3520. PubMed PMID: 24027140; PubMed Central PMCID: PMC3855630.
45. Aboa-Eboule C, Mengue D, Benzenine E, Hommel M, Giroud M, Bejot Y, et al. How accurate is the reporting of stroke in hospital discharge data? A pilot validation study using a population-based stroke registry as control. *Journal of neurology*. 2013;260(2):605-13. Epub 2012/10/19. doi: 10.1007/s00415-012-6686-0. PubMed PMID: 23076827; PubMed Central PMCID: PMC3566387.
46. Rezaie A, Quan H, Fedorak RN, Panaccione R, Hilsden RJ. Development and validation of an administrative case definition for inflammatory bowel diseases. *Canadian journal of gastroenterology = Journal canadien de gastroenterologie*. 2012;26(10):711-7. Epub 2012/10/13. PubMed PMID: 23061064; PubMed Central PMCID: PMC3472911.
47. Holland-Bill L, Christiansen CF, Ulrichsen SP, Ring T, Jorgensen JO, Sorensen HT. Validity of the International Classification of Diseases, 10th revision discharge diagnosis codes for hyponatraemia in the Danish National Registry of Patients. *BMJ open*. 2014;4(4):e004956. Epub 2014/04/25. doi: 10.1136/bmjopen-2014-004956. PubMed PMID: 24760354; PubMed Central PMCID: PMC4010845.
48. Vasta R, Boumediene F, Couratier P, Nicol M, Nicoletti A, Preux PM, et al. Validity of medico-administrative data related to amyotrophic lateral sclerosis in France: A population-based study. *Amyotrophic lateral sclerosis & frontotemporal degeneration*. 2017;18(1-2):24-31. Epub 2016/11/01. doi: 10.1080/21678421.2016.1241280. PubMed PMID: 27797285.
49. Jiang J, Southern D, Beck CA, James M, Lu M, Quan H. Validity of Canadian discharge abstract data for hypertension and diabetes from 2002 to 2013. *CMAJ open*. 2016;4(4):E646-e53. Epub 2016/12/27. doi: 10.9778/cmajo.20160128. PubMed PMID: 28018877; PubMed Central PMCID: PMC5173472.
50. Hux JE, Ivis F, Flintoft V, Bica A. Diabetes in Ontario: determination of prevalence and incidence using a validated administrative data algorithm. *Diabetes care*. 2002;25(3):512-6. Epub 2002/03/05. PubMed PMID: 11874939.
51. Tu K, Wang M, Young J, Green D, Ivers NM, Butt D, et al. Validity of administrative data for identifying patients who have had a stroke or transient ischemic attack using EMRALD as a reference standard. *The Canadian journal of cardiology*. 2013;29(11):1388-94. Epub 2013/10/01. doi: 10.1016/j.cjca.2013.07.676. PubMed PMID: 24075778.
52. Cutrona SL, Toh S, Iyer A, Foy S, Cavagnaro E, Forrow S, et al. Design for validation of acute myocardial infarction cases in Mini-Sentinel. *Pharmacoepidemiology and drug safety*. 2012;21 Suppl 1:274-81. Epub 2012/01/25. doi: 10.1002/pds.2314. PubMed PMID: 22262617; PubMed Central PMCID: PMC3679667.
53. Wahl PM, Terrell DR, George JN, Rodgers JK, Uhl L, Cataland S, et al. Validation of claims-based diagnostic codes for idiopathic thrombotic thrombocytopenic purpura in a commercially-insured population. *Thrombosis and haemostasis*. 2010;103(6):1203-9. Epub 2010/03/31. doi:

- 10.1160/th09-08-0595. PubMed PMID: 20352159.
54. Lacasse Y, Daigle JM, Martin S, Maltais F. Validity of chronic obstructive pulmonary disease diagnoses in a large administrative database. *Canadian respiratory journal*. 2012;19(2):e5-9. Epub 2012/04/27. PubMed PMID: 22536584; PubMed Central PMCID: PMC3373291.
55. Baghestan E, Bordahl PE, Rasmussen SA, Sande AK, Lyslo I, Solvang I. A validation of the diagnosis of obstetric sphincter tears in two Norwegian databases, the Medical Birth Registry and the Patient Administration System. *Acta obstetrica et gynecologica Scandinavica*. 2007;86(2):205-9. Epub 2007/03/17. doi: 10.1080/00016340601111364. PubMed PMID: 17364284.
56. Rishi MA, Kashyap R, Wilson G, Hocker S. Retrospective derivation and validation of a search algorithm to identify extubation failure in the intensive care unit. *BMC anesthesiology*. 2014;14:41. Epub 2014/06/04. doi: 10.1186/1471-2253-14-41. PubMed PMID: 24891838; PubMed Central PMCID: PMC4041644.
57. Mamtani R, Haynes K, Boursi B, Scott FI, Goldberg DS, Keefe SM, et al. Validation of a coding algorithm to identify bladder cancer and distinguish stage in an electronic medical records database. *Cancer epidemiology, biomarkers & prevention : a publication of the American Association for Cancer Research, cosponsored by the American Society of Preventive Oncology*. 2015;24(1):303-7. Epub 2014/11/13. doi: 10.1158/1055-9965.epi-14-0677. PubMed PMID: 25389114; PubMed Central PMCID: PMC4294969.
58. Fombonne E, Heavey L, Smeeth L, Rodrigues LC, Cook C, Smith PG, et al. Validation of the diagnosis of autism in general practitioner records. *BMC public health*. 2004;4:5. Epub 2004/04/29. doi: 10.1186/1471-2458-4-5. PubMed PMID: 15113435; PubMed Central PMCID: PMC394332.
59. Klungsoyr K, Harmon QE, Skard LB, Simonsen I, Austvoll ET, Alsaker ER, et al. Validity of pre-eclampsia registration in the medical birth registry of Norway for women participating in the Norwegian mother and child cohort study, 1999-2010. *Paediatric and perinatal epidemiology*. 2014;28(5):362-71. Epub 2014/07/22. doi: 10.1111/ppe.12138. PubMed PMID: 25040774; PubMed Central PMCID: PMC4167249.
60. Ducharme R, Benchimol EI, Deeks SL, Hawken S, Fergusson DA, Wilson K. Validation of diagnostic codes for intussusception and quantification of childhood intussusception incidence in Ontario, Canada: a population-based study. *The Journal of pediatrics*. 2013;163(4):1073-9.e3. Epub 2013/07/03. doi: 10.1016/j.jpeds.2013.05.034. PubMed PMID: 23809052.
61. Hall GC. Validation of death and suicide recording on the THIN UK primary care database. *Pharmacoepidemiology and drug safety*. 2009;18(2):120-31. Epub 2008/12/26. doi: 10.1002/pds.1686. PubMed PMID: 19109801.
62. Kumamaru H, Judd SE, Curtis JR, Ramachandran R, Hardy NC, Rhodes JD, et al. Validity of claims-based stroke algorithms in contemporary Medicare data: reasons for geographic and racial differences in stroke (REGARDS) study linked with Medicare claims. *Circulation Cardiovascular quality and outcomes*. 2014;7(4):611-9. Epub 2014/06/26. doi: 10.1161/circoutcomes.113.000743. PubMed PMID: 24963021; PubMed Central PMCID: PMC4109622.

63. Jorgensen LK, Dalgaard LS, Ostergaard LJ, Andersen NS, Norgaard M, Mogensen TH. Validity of the coding for herpes simplex encephalitis in the Danish National Patient Registry. *Clinical epidemiology*. 2016;8:133-40. Epub 2016/06/23. doi: 10.2147/clep.s104379. PubMed PMID: 27330328; PubMed Central PMCID: PMC4896464.
64. Tairou F, De Wals P, Bastide A. Validity of death and stillbirth certificates and hospital discharge summaries for the identification of neural tube defects in Quebec City. *Chronic diseases in Canada*. 2006;27(3):120-4. Epub 2007/02/20. PubMed PMID: 17306063.
65. Greenland S, Lasth TL. Disease Misclassification. In *Modern Epidemiology 3rd ed.* Eds Rothman KJ, Greenland S, Lasth TL, Lippincott Williams & Wilkins PA, 2008. 358-60.
66. Glas AS, Lijmer JG, Prins MH, et al. The diagnostic odds ratio: a single indicator of test performance. *Journal of Clinical Epidemiology* 2003;56(11):1129-35.
67. Turner RM, Chen YW, Fernandes AW. Validation of a Case-Finding Algorithm for Identifying Patients with Non-small Cell Lung Cancer (NSCLC) in Administrative Claims Databases. *Front Pharmacol*. 2017 Nov 30;8:883.
68. Greenland S, Lasth TL. Relation of Predictive Values to Sensitivity and Specificity. In *Modern Epidemiology 3rd ed.* Eds Rothman KJ, Greenland S, Lasth TL, Lippincott Williams & Wilkins PA, 2008. 357-8.
69. Goldberg DS, Lewis JD, Halpern SD, Weiner MG, Lo Re V, 3rd. Validation of a coding algorithm to identify patients with hepatocellular carcinoma in an administrative database. *Pharmacoepidemiology and drug safety*. 2013;22(1):103-7. Epub 2012/11/06. doi: 10.1002/pds.3367. PubMed PMID: 23124932; PubMed Central PMCID: PMC3540172.
70. Ooba N, Setoguchi S, Ando T, Sato T, Yamaguchi T, Mochizuki M, et al. Claims-based definition of death in Japanese claims database: validity and implications. *PloS one*. 2013;8(5):e66116. Epub 2013/06/07. doi: 10.1371/journal.pone.0066116. PubMed PMID: 23741526; PubMed Central PMCID: PMC3669209.
71. 山口拓洋、富士武史、赤木将男、安部靖之、中村真潮、山田典一ほか. 医療情報データベースを用いた静脈血栓塞栓症発症、出血性イベントのバリデーション研究. *医薬品情報学* 2015; 17:87-93.
72. Tanaka S, Hagino H, Ishizuka A, Miyazaki T, Yamamoto T, Hosoi T. Validation Study of Claims-based Definitions of Suspected Atypical Femoral Fractures Using Clinical Information. *薬剤疫学* 2016; 21:13-9.
73. Yamana H, Moriwaki M, Horiguchi H, Kodan M, Fushimi K, Yasunaga H. Validity of diagnoses, procedures, and laboratory data in Japanese administrative data. *Journal of epidemiology*. 2017;27(10):476-82. Epub 2017/02/01. doi: 10.1016/j.je.2016.09.009. PubMed PMID: 28142051; PubMed Central PMCID: PMC5602797.
74. Ishii M. DRG/PPS and DPC/PDPS as Prospective Payment Systems. *Japan Medical Association journal : JMAJ*. 2012;55(4):279-91. Epub 2012/07/01. PubMed PMID: 25237234.
75. Okamura S, Kobayashi R, Sakamaki T. Case-mix payment in Japanese medical care. *Health policy*. 2005;74(3):282-6.

76. Matsuda S, Ishikawa KB, Kuwabara K, Fujimori K, Fushimi K, Hashimoto H. Development and use of the Japanese case-mix system. *Eurohealth*. 2008; 14:25-30.
77. 井上真奈美. 日本におけるがんの大規模コホート研究：歴史と展開. *日本保険医学会誌* 2015;113:147-57.
78. Setoguchi S, Solomon DH, Glynn RJ, Cook EF, Levin R, Schneeweiss S. Agreement of diagnosis and its date for hematologic malignancies and solid tumors between medicare claims and cancer registry data. *Cancer causes & control : CCC*. 2007;18(5):561-9. Epub 2007/04/21. doi: 10.1007/s10552-007-0131-1. PubMed PMID: 17447148.

「日本における傷病名を中心とするレセプト情報から得られる指標のバリデーションに関するタスクフォース」報告書

付録 1:同じ感度・特異度であっても研究集団の曝露の割合によって研究結果が変わるケース・コントロール研究の例

Strom 編の Pharmacoepidemiology 第 5 版の第 41 章に示された、曝露の誤分類がケース・コントロール研究の結果に与える影響について解説する。ここで解説するのと同様の問題はアウトカムの誤分類でも起こりうる。

以下のケース・コントロール研究で示すように、同じ感度・特異度であっても、それらの誤分類が結果に与えるバイアスの程度は、研究集団における曝露の割合等の影響を受ける。従って、これらの指標の絶対値にはあまり意味が無く、究極の基準は効果の測定に影響するバイアスの程度である。

【セッティング】

- 真のオッズ比(OR):3.0、ケースとコントロール各 1,000 例
 - 曝露の測定への妥当性尺度:
 - 感度:(ケース)0.9、(コントロール)0.8
 - 特異度:(ケース)0.95、(コントロール)0.99

ソース集団(およびコントロール)の曝露割合 = 10%	ソース集団(およびコントロール)の曝露割合 = 90%																																				
<p>真)</p> <table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <thead> <tr> <th></th> <th>曝露</th> <th>非曝露</th> </tr> </thead> <tbody> <tr> <th>Case</th> <td>250</td> <td>750</td> </tr> <tr> <th>Control</th> <td>100</td> <td>900</td> </tr> </tbody> </table> <p style="text-align: right;">OR =3.0</p> <p>測定値)</p> <table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <thead> <tr> <th></th> <th>曝露</th> <th>非曝露</th> </tr> </thead> <tbody> <tr> <th>Case</th> <td>262^a</td> <td>738^b</td> </tr> <tr> <th>Control</th> <td>89^c</td> <td>911^d</td> </tr> </tbody> </table> <p style="text-align: right;">OR =3.6</p> <p style="font-size: small;">a: $250 \cdot 0.9 + 750 \cdot (1 - 0.95)$; b: $250 \cdot (1 - 0.9) + 750 \cdot 0.95$; c: $100 \cdot 0.8 + 900 \cdot (1 - 0.99)$; d: $100 \cdot (1 - 0.8) + 900 \cdot 0.99$</p>		曝露	非曝露	Case	250	750	Control	100	900		曝露	非曝露	Case	262 ^a	738 ^b	Control	89 ^c	911 ^d	<p>真)</p> <table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <thead> <tr> <th></th> <th>曝露</th> <th>非曝露</th> </tr> </thead> <tbody> <tr> <th>Case</th> <td>964</td> <td>36</td> </tr> <tr> <th>Control</th> <td>900</td> <td>100</td> </tr> </tbody> </table> <p style="text-align: right;">OR =3.0</p> <p>測定値)</p> <table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <thead> <tr> <th></th> <th>曝露</th> <th>非曝露</th> </tr> </thead> <tbody> <tr> <th>Case</th> <td>869^e</td> <td>131^f</td> </tr> <tr> <th>Control</th> <td>721^g</td> <td>279^h</td> </tr> </tbody> </table> <p style="text-align: right;">OR =2.6</p> <p style="font-size: small;">e: $964 \cdot 0.9 + 36 \cdot (1 - 0.95)$; f: $964 \cdot (1 - 0.9) + 36 \cdot 0.95$; g: $900 \cdot 0.8 + 100 \cdot (1 - 0.99)$; h: $900 \cdot (1 - 0.8) + 100 \cdot 0.99$</p>		曝露	非曝露	Case	964	36	Control	900	100		曝露	非曝露	Case	869 ^e	131 ^f	Control	721 ^g	279 ^h
	曝露	非曝露																																			
Case	250	750																																			
Control	100	900																																			
	曝露	非曝露																																			
Case	262 ^a	738 ^b																																			
Control	89 ^c	911 ^d																																			
	曝露	非曝露																																			
Case	964	36																																			
Control	900	100																																			
	曝露	非曝露																																			
Case	869 ^e	131 ^f																																			
Control	721 ^g	279 ^h																																			

(教科書を元に著者ら作成)

「日本における傷病名を中心とするレセプト情報から得られる指標のバリデーションに関するタスクフォース」報告書

付録 2: データベース研究に関連するガイドライン・ガイダンスにおけるバリデーション研究に関する記載の要約

本タスクフォースでは、データベース研究に関連する国内外のガイドライン・ガイダンス等から、アウトカムのバリデーション研究に関する言及の有無や言及されている内容について抜き出した。以下は、その重要な点についてまとめたものである。

(ア) 医療情報のデータベース等を用いた医薬品の安全性評価における薬剤疫学研究の実施に関するガイドライン(初版 平成 26 年 3 月 31 日:PMDA)

本ガイドラインは、PMDA 及び製薬企業等が医療情報のデータベースを二次利用して医薬品の安全性評価を行う際に、適切な薬剤疫学研究が実施されるよう留意事項をまとめたものであり、「3. データソース」の「(3)バリデーション」において、バリデーション研究の目的や方法論、実施する意義、課題などについて網羅的に解説している。

① バリデーション研究の必要性について

1. バリデーション研究を実施することを基本推奨する。特に診療報酬データの傷病名コードのみでアウトカムを定義する場合には、バリデーション研究が必要となることがある。また、研究実施計画書および研究結果報告書(もしくはそれに代わる記録)では、バリデーション研究の実施について下記の点について言及する: 対象集団、曝露、アウトカム、または共変量等の各種定義に対するバリデーション研究について、予定も含めた実施の有無について記載(実施する場合はその研究の計画等を添付、実施しない場合にはその理由を記載)。
2. 既存のバリデーション研究の結果を利用する場合には、そうすることが適切だとした理由を記載し、参考とした資料を研究実施計画書に添付。

② アウトカム判定(基準)について

1. アウトカムの判定基準は、関連する臨床ガイドライン等に定められた標準的な診断基準があればそれを参考に作成し、関連する臨床ガイドライン等がない場合であっても、可能な限り一般的に確立された診断基準を参考にした判定基準を設定することが望ましい。
2. 臨床的知識が必要とされるため臨床医が判定者となることが多いが、結果の客観性を担保するために複数の判定者によって実施し、判定結果のバラツキを抑えるために少数の判定者によって実施することが望ましい。

③ 妥当性を示す指標について

1. 感度、特異度はバリデーション研究の対象とした集団の特性の構成に影響されない指標値である。しかし、実施可能性の観点から多くの場合、陽性的中度しか求められないことも多い。
2. 陽性的中度は、研究のために用いる特定のデータベースと定義に特異的な指標

値である。しかし、研究に使用するデータの対象者が連結不可能匿名化されている場合等は、診療録等の情報を入手し妥当性の判定を行うことができないため、他のデータベースで同様の定義に関して実施したバリデーション研究の結果を参考にすることもある。

(イ) 製造販売後データベース調査実施計画書の記載要領(平成 30 年 1 月 23 日:PMDA)

本記載要領は、「医薬品の製造販売後の調査及び試験の実施の基準に関する省令等の一部を改正する省令」(平成 29 年厚生労働省令第 116 号)に製造販売後データベース調査が規定されたことを踏まえ、医薬品の製造販売業者等が製造販売後データベース調査実施計画書を作成する際の参考として提示されている。そのため、本要領におけるバリデーション研究に関する記述も、計画書への記載方法に特化している。以下にバリデーション研究に関連する箇所を引用する。

13.8. バリデーション

バリデーションの記載要領

- 調査で設定した全ての定義について、バリデーションスタディに関する既存の研究報告の有無及びその概要、及びバリデーションスタディ実施予定の有無を記載すること。
- バリデーションスタディが求められる対象として、対象集団の定義、曝露の定義、アウトカムの定義、共変量等が挙げられるが、特にアウトカムの定義については、既存の研究報告の有無及びバリデーションスタディの実施の有無について明記すること。
- バリデーションスタディを実施する場合は、計画の概要を記載するか、実施計画書等を付録とすること。付録とする場合、本章にはその旨を記載すること。
- バリデーションスタディを実施しない場合は、実施しない場合であっても適切な調査が実施可能と判断した理由を記載すること。
- 既存のバリデーションスタディの結果を利用してアウトカム等を定義する場合には、そのバリデーションスタディの結果を利用することが適切な理由及びバリデーションスタディの結果を記載するとともに、参考とした資料を付録とすること。

18. 付録

- 付録は以下のようなものが想定される。
 - バリデーションスタディ実施計画書又は結果報告書、もしくはアウトカム定義等の妥当性に関する資料、関連文献等

(ウ) Best Practices for Conducting and Reporting Pharmacoepidemiologic Safety Studies Using Electronic Healthcare Data (May 2013:FDA)

本ガイダンスは、診療報酬請求データ・電子カルテ等の医療情報データを用いた薬剤疫学的な安全性研究の実施・報告に関するベストプラクティスを記述している。本ガイダンス中、バリデーション研究の実施方法については、紹介の要素が強い以下の一文を除き具体的な記載はなかった: 「ICD のようなコードに基づくアウトカム定義での陽性的中度の評価は、データソースからそのコードを持つ症例の全て、もしくはサンプルを抽出し、その症例が本当にそのコードされたイベントを経

験したか確認するために、その症例の一次医療データ(通常はカルテ)をレビューする」。一方、アウトカムの妥当性に言及した「C. Study Design: Outcome Definition and Ascertainment」の章中の項(2. Validation of Outcomes)を中心に、いくつかバリデーション研究に関連する記載が下記のとおり存在した。

① バリデーション研究の必要性について

データベース研究を実施する際のデータソースの選択においては、アウトカムの妥当性に関する以下の点を検討の上、プロトコルへの記載も考慮すべきである。

1. アウトカムや他の研究因子(例えば、曝露、重要な共変数、選択/除外基準)の妥当性検証の実施可能性。
2. 選択したデータベースで実施されたアウトカムに関するバリデーション研究を含むあらゆる関連論文とその詳細(バリデーション研究が実施された集団、データベースや時期、及びその性能特性)。
3. 事前のアウトカムのバリデーションが無い研究では、研究者が使用したアウトカム定義を適切と考えた根拠。
4. アウトカムの重症度が把握できるか(患者がアウトカムについて異なった表現で訴えたり、重症度が様々なものについては特に重要である。例えば、アレルギー反応は、皮膚の発疹からアナフィラキシーショックのような生命を脅かすものまで様々である)。

② 妥当性を示す指標について

1. 研究のデザインにおいて、研究者は選択した臨床アウトカムの妥当性が検証されているかを確認すべきである。具体的には、臨床的に適切なアウトカム定義を確立し、その定義の陽性的中度を評価する。
2. 感度も重要であり、たとえば、あるアウトカム定義が高い陽性的中度を示しても、もしそのアウトカムが臨床的な関心をあまり引かず、それゆえ見逃されるのであれば、感度は低くなりうる。
3. したがって、研究者は、(1)そのアウトカム定義の陽性的中度と研究の内部妥当性に与える潜在的な影響、(2)アウトカム定義の感度、すなわち選択されたデータソースと研究対象集団においてアウトカムが特定される程度、について議論すべきである。
4. さらに、実社会での患者集団への外的妥当性(研究結果の一般化可能性)に与える潜在的な影響も議論すべきである。

③ 診断名をアウトカム定義に用いる際の考慮事項(米国でのプラクティス)

1. 退院時診断の主病名、副病名の区分はしばしば恣意的になされる。よって、退院時診断名の順番が、それらの医学的な重要性と一致しているとは限らない。
2. 入院時の保険請求データにおける ICD コードは、一般的により信頼性が高く、より重篤な疾病を反映しやすいことから、請求データを用いてアウトカムを定義する際には、可能な限り入院データのコードを使用しバリデーションするのが良い。
✓ 外来の請求データにおける ICD コードは、保険償還を最大化するため、しば

しば「upcoding」や「downcoding」がされてきた(Strom 2005)。

3. 患者が療養中に死亡した時は多くの場合でシステムに記録されるが、療養中ではない時に起こった死亡は電子的なヘルスケアデータシステムに記録されないことがある。よって死亡の信頼性のある確認は、他のシステムとのリンケージを通じて実施されることが好ましい。

(エ) Guideline on good pharmacovigilance practices (GVP): Module VIII – Post-authorisation safety studies (Rev 3)

European Medicines Agency (EMA)は、Post-authorisation safety studies (PASS)の実施についてガイダンスを発出している。このガイダンスには、バリデーション研究そのものの説明や実施方法に関する記述はないが、バリデーション研究に関連する記述はある。

バリデーション研究の必要性については以下のように記されている:「利用するデータベースが研究で求められる程度に詳細で正確な情報を有していない場合、カルテに戻り診断名の確認を考慮する。関心のあるアウトカムによっては、症例ごとに行う判断もしくはランダムサンプリングした症例レビューによりバリデーションをしてもよい。」

また、研究計画書には、データソース(1, 2)と品質管理(3)に関連して以下を記載すべきとある。

1. 利用するデータソースに含まれる記録やコーディングの妥当性に関するあらゆる情報
2. 診断名のバリデーションのための専門家委員会の関与の有無
3. 品質管理に関して、エンドポイントのバリデーションを含めたデータの質を確認する方法や手順

最終研究報告書にも、コーディングの妥当性に関して記述し、それらが研究に及ぼす限界について考察するべきである。

(オ) Scientific guidance on post-authorisation efficacy studies

EMAはPost-authorisation efficacy studiesのガイダンスも公開しているが、その中でバリデーション研究に直接的に言及する記述は存在しない。医薬品のベネフィットを評価する観察研究においては、高い精度の曝露とアウトカムに関する情報が必要であるとの記述にとどまっている。

(カ) ENCePP: Guide on Methodological Standards in Pharmacoepidemiology (Revision 6)

本ガイダンスでは、以下のようなバリデーション研究に関連する記述が見られた。

1. 医療データベースを用いたアウトカムの適切な評価は、研究自体の正当性を示す上で極めて重要である。
2. アウトカム定義で用いられる因子の完全性と妥当性について十分な検討が必要であり、症例定義やアルゴリズムに用いられた仮説(前提)は検証されるべきである。バリデーションの実施方法については、専門家によるカルテレビューだけで

なく、がん登録や死亡レジストリー、もしくはそれらから算出された発生率や有病割合との一致性を確認することによる方法もある。

3. 研究によく用いられるデータベースでは、過去に重要な因子(key variables)のバリデーション研究が実施されその結果が文書化されていることがあるが、そのような過去のバリデーション研究の結果を外挿する場合には、(バリデーション研究における)因子や分析内容との差異、その後の医療、診療行為、コーディングの変化の影響を考慮する必要がある。つまり、医療制度/環境とデータが生み出される手順の両面を熟知することが好ましい。

なお、ENCePP による研究計画書のチェックリスト(Revision 3)では、Section 6: Outcome definition and measurement の中に、「6.3 Does the protocol address the validity of outcome measurement? (e.g. precision, accuracy, sensitivity, specificity, positive predictive value, prospective or retrospective ascertainment, use of validation sub-study)」というチェック項目が存在する。

(キ) Guidelines for Good Database Selection and use in Pharmacoepidemiology Research

このガイドラインでは、薬剤疫学研究におけるデータベースの選択と使用に関するチェックリストを示し、それぞれのチェック項目について説明している。バリデーション研究に関連するものとして、以下の記述が確認できる。

1. レセプトデータ(claims data)の診断名は、最終的な診断名ではなく、“working”な診断名であることがあるため、偽陽性や偽陰性を考慮すべきである。
2. データのバリデーションの方法として、カルテ(clinical note)や死亡診断書のような外部資料との突合せがある。電子カルテ(EMR)がデータソースの場合では、死亡統計や大規模臨床試験のメタ解析などの外部の数値データとの比較といった方法によるバリデーションを実施すべきである。
3. 妥当性を示す指標については、主要アウトカムを特定するために使われたアルゴリズム/方法の陽性的中度を算出することがある。

(ク) GPS – Good Practice in Secondary Data Analysis: Revision after Fundamental Reworking

二次データを用いた分析を実施するための基本的基準を確立することを目的に作られた文書である。「Recommendation 6.4 – Data Quality」の項に、利用可能な情報を用いデータの信頼性と妥当性を評価すべき、といったバリデーションに関する概念的な記述はあるものの、それを実行するための具体的な手段に関する記述は存在しない。

(ケ) Guidelines for Good Pharmacoepidemiology Practices (GPP)

薬剤疫学研究の質と完全性の確保に向けた最低限のプラクティスと手順を記すことを目的としている文書である。バリデーション研究に特化した記述は無いが、使用したデータベース、アウトカムの定義、妥当性の検証方法、診断のバリデーションのための専門家委員会の有無、評価手順などについて研究計画書に記述すること、としている。

(コ) A Checklist for Retrospective Database Studies – Report of the ISPOR Task Force on Retrospective Databases

レトロスペクティブなデータベース研究を実施する上でのチェックリストとして利用されることを想定した文書である。随所でデータの質の確保についての言及がなされているが、アウトカムのバリデーション研究についての詳細な説明はない。しかし、データの質の確保の一環として、アウトカムを特定するためのコードやアルゴリズムが妥当であるという根拠は提示されるべきとしている。また、その根拠としてしばしば過去の研究が引用されるが、理想的には一次データソースに対してのバリデーションがなされるべき、と記されている。

(サ) Developing a Protocol for Observational Comparative Effectiveness Research: A User’s Guide

The Agency for Healthcare Research and Quality による比較効果研究の研究計画立案に関するガイダンスであり、データベース研究を実施する上での注意点などについて記載されている。コーディングの妥当性検証の重要性については訴えているものの、バリデーション研究自体の方法については記述されず、バリデーション研究を実施する上で参考になる論文を紹介するにとどまっている。

(シ) The REporting of studies Conducted using Observational Routinely-collected health Data (RECORD) Statement

STrengthening the Reporting of OBservational studies in Epidemiology (STROBE) を基にした”routinely collected health data”を用いた研究報告に関する基準である。バリデーション研究自体についての方法論などの解説はないが、関係する記載として、集団の選択に用いられたコードやアルゴリズムに関するあらゆるバリデーション研究を引用すべきであり、また、本研究のためのバリデーションを実施したがその結果を論文化していない場合は、そのバリデーションの詳細な方法や結果を(本研究の論文に)記述すべきとしている。

以上、国内外のガイドラインなどでは「アウトカムのバリデーションは重要である」という考えは共通していたが、そのための研究方法や結果の解釈の仕方についての包括的な記述は確認できなかった。その理由の一つとして、研究の目的、用いるデータソース、注目するアウトカムの医学的性質(重症度や珍しさなど)、医療制度・医事会計制度を含めたデータが収集される環境などにより最適なバリデーションの方法が異なり、画一的なプラクティスを提供することが難しいからではないかと推察される。

「日本における傷病名を中心とするレセプト情報から得られる指標のバリデーションに関するタスクフォース」報告書

付録 3: 過去に実施されたバリデーション研究のまとめ

本タスクフォースでは、バリデーション研究の要点について先例から学び、今後の日本の研究に活かすことを1つの目標として過去のバリデーション研究をレビューした。ただし、網羅的なシステムティックレビューではなく、どのような研究があるかを俯瞰することを目的とし、対象文献数は105に留めた。具体的には、2016年12月にPubMedで(網羅性はないが重要な論文ができるだけ含まれることが期待できる)検索式

((sensitivity[Title/Abstract]) OR predictive value[Title/Abstract]) AND database[Title/Abstract]) AND valid*[Title] でヒットした412論文から、タイトル・アブストラクトで論文を絞りこんだ。その際、(特定のアウトカムやデータベースについての)システムティックレビュー論文は除外したが、これから行う研究のためのプロトコル論文は残した。また、バリデーション研究の目的(例:バリデーションされた傷病名などを今後のデータベース研究におけるアウトカム、曝露、研究対象集団の特定に使うことが想定されるか)に制限は加えなかった。その後、検索式ではヒットしなかったがタスクフォースのメンバーが把握している重要な論文をいくつか追加した。各研究の詳細は別途エクセルファイル(付録3 文献リスト.xlsx)に示しているが、ここでは下記の8個の項目に焦点を当て、各項目において重要な研究を選択して紹介することとした。なお、引用文献番号は別途エクセルファイル(付録3 文献リスト.xlsx)と対応している。

(1) 利用する医療情報データベースの種類

アウトカム定義のために使われたデータベースは下記のように大別された。①②③はadministrative dataとも呼ばれる。Administrative dataは患者診療に用いるのではなく、主に病院の管理目的に集めたデータである。

① 診療報酬データ(claims data):

医療機関が診療報酬を請求書するためのデータであり、外来と入院の両方を対象にしているものが多い。代表的なものとしては米国の Medicare⁽¹⁻³⁾、Medicaid⁽⁴⁾、民間の診療請求データ⁽⁵⁾や、台湾の National Health Research Database^(6,7)、カナダの Ontario Health Insurance Plan^(8,9)、Alberta州の Data Integration, Measurement, and Reporting (DIMR)⁽¹⁰⁾などがある。

② 退院時記録(discharge data):

患者退院時に主に担当医が作成する退院時記録のデータであり、カナダの複数の州の discharge abstract database (DAD)⁽¹¹⁻¹³⁾、デンマークの Danish National Hospital Register^(14,15)、フランスの French Hospital Discharge Database^(16,17)などがある。

③ 外来の administrative data:

カナダの National Ambulatory Care Reporting System (NACRS)⁽¹⁸⁾は、診療報酬請求目的ではない一般外来と救急外来の administrative data である。

④ 電子カルテ(Electronic Medical Records):

日常診療の電子カルテをデータベース化したものであり、アメリカの退役軍人データベース⁽¹⁹⁻²¹⁾、

オランダの Integrated Primary Care Information (IPCI)⁽²²⁾、イタリアの Health Search/CSD Patient Database(HSD)^(22,23)、カナダの Canadian Primary Care Sentinel Surveillance Network(CPCSSN)の電子医療記録⁽²⁵⁾、ノルウェーの Patient Administration System(PAS)⁽²⁶⁾、スペインの Minimum Basic Data Set (MBDS)⁽²⁷⁾のほか、イギリスのプライマリーケアデータベースの The Health Improvement Network(THIN)⁽²⁸⁻³³⁾などが代表的である。イギリスの General Practice Research Database(GPRD)⁽³⁴⁾もプライマリーケアデータベースの 1 つであるが、GPRD は現在他のデータベースとのリンクが行える Clinical Practice Research Datalink に進化しており、GPRD については下記の⑥に含めた。

⑤ 疾患レジストリー:

特定の疾患を持つ患者を登録しているデータベースであり、日常診療データとは特徴が異なる。ノルウェーの Medical Birth Registry of Norway^(26,35,36)やカナダの糖尿病レジストリー⁽³⁷⁾や Trauma Registry⁽³⁸⁾、米国の Colorado Violent Death Reporting System⁽³⁹⁾、オランダの口唇口蓋裂のレジストリー⁽⁴⁰⁾、英国の National Hip Fracture Database⁽⁴¹⁾、など様々な疾患を対象としたレジストリーがある。疾患レジストリーはゴールドスタンダードの定義に使われることが多いが、疾患レジストリーそのものをバリデーションする研究も散見された。

⑥ 上記の組み合わせ:

以上の①～⑤を 2 つ以上組み合わせたデータも存在した。特にカナダでは、オンタリオ州の外来の診療報酬データと退院時記録の組み合わせ⁽⁴²⁻⁴⁴⁾、さらに日帰り手術データ Same Day Surgery Database をも組み合わせたデータ⁽⁴⁵⁾が評価されている。同様にアルバータ州⁽¹⁸⁾、ブリティッシュコロンビア州⁽²⁴⁾、カルガリー州⁽⁴⁶⁾、マニトバ州⁽⁴⁷⁾、ノバスコシア州⁽⁴⁸⁾でも Discharge Abstract Database (DAD)と Physician Claims や National Ambulatory Care Reporting System (NACRS)の組み合わせ、DAD、Medical Services Insurance (MSI)、Mental Health Outpatient Information System (MHOIS)の組み合わせ⁽⁴⁸⁾など、リンクされたデータを評価した研究が多い。イギリスのプライマリーケアデータベースの 1 つである Clinical Practice Research Datalink(CPRD)^(49,50)も入院・外来の Hospital Episode Statistics(HES)や Cancer Registry などとリンクさせて研究に使用することができる。また、死亡診断書、警察からの報告、監察医、検視官などかの情報を集めて作成されているアメリカの National Violent Death Reporting System⁽³⁹⁾、薬、退院時データと検査に関するデータなどからなるイタリアの Drug Administrative Database (DAD)⁽⁵¹⁾、フランスの退院時記録 (hospital discharge data)と診療報酬データ (health insurance data)の組み合わせ⁽¹⁶⁾などもバリデーション研究によって評価されている。

⑦ 薬剤疫学研究のために作られたデータベース:

少数ではあったが、米国の(Mini-)Sentinel Distributed Database⁽⁵²⁻⁵⁴⁾やヨーロッパ数カ国による European Commission-funded exploring and understanding adverse drug reactions (EU-ADR) project のデータベース⁽²²⁾、アメリカの Brigham and Women's Hospital's における Medicare とリンクした EMR データベース⁽⁵⁵⁾など、薬剤疫学研究のために作られたデータベースも存在した。

(2) バリデーション研究のセッティング

① 国:

今回レビューした研究の中では、カナダと米国が最も多く、英国、デンマークなども多かった。

② 研究対象期間:

データの研究対象開始年は2000年以降に集中し、2000年以降バリデーション研究は年々増加する傾向が見られた。データの研究対象期間は1年未満のものが最も多く、5年未満が過半数を占めた。

③ 研究の”population”:

研究が対象としている”population”については、国や州といった地域の住民、Health maintenance organization (HMO)などの被保険者、グループ病院などの複数施設⁽²³⁾、大学病院などの単施設で診療を受けた患者^(56,57)、特定の疾患を有する患者^(18,42,48,58-60)など様々であった。一方で、研究において想定される population は国や地域の医療体制に依存することも多く、論文内の記載だけでは詳細不明であるものもあった。一般的には、居住地域や加入する保険などが共通する健康人を含む集団を研究対象としており、さらにそこから検討症例または研究施設の全てまたはランダムサンプリングを用いるバリデーション研究は”population-based”とみなせる。一方で、複数/単施設を対象としている研究は”population-based”ではないと判断されることが多い。しかし、ある地域において特定の疾患を持つ患者が全例その複数/単施設に集まってくる状況では、研究が”population-based”とみなせる場合があり⁽⁶¹⁾、データベースの特徴を論文もしくは他の資料を参照して判断する必要があった。

(3) アウトカムの定義

アウトカムの定義として、疾患コードである International Classification of Diseases 9th Revision (ICD-9)もしくは ICD-10 が主に使用されていたが、英国プライマリーケアデータの Read code^(29-32,62,63)のようなデータベース特有の疾患コードも見られた。また、薬剤処方コード^(64,65)や処置・手術コード⁽⁶⁶⁻⁶⁸⁾、検査値^(15,32)を用いてアウトカムを定義している研究も存在した。疾患コードに薬剤や手術のコードを組み合わせる^(55,58,69)、あるいは疾患コードと入院記録、受診回数情報を組み合わせる^(53,58,70-72)ことで、陽性的中度や感度、特異度の値を高める工夫をしている研究もあった。

(4) リンケージの有無・方法

今回レビューした研究の中では、複数(2つ以上)のデータベースがリンクされていることが多かった。すでにリンク済みのデータベースを使用していた研究と、その研究のために後述する方法によってリンケージさせた研究があった。

リンケージの目的は主に2つあり、1つはアウトカム定義に用いるデータベースとゴールドスタンダードの定義に用いるデータベースなどとのリンケージであり、もう1つは(1)の⑥で言及したカナダの複数の州で実施されている入院と外来の記録を組み合わせるなどのより質の高いアウトカム定義を作成するための複数のデータベース同士のリンケージである。日本における個別の病院におけるバリデーション研究では、バリデーション研究を行う外部研究者に提供するデータと病院内部のデータを研究に固有の番号でリンケージしたものも見られた⁽⁶⁷⁾。

リンケージの方法は以下の2つであった。

① Deterministic linkage:

2つ(以上)のファイルに共通で存在する1つ以上の個人固有の識別子(社会保障番号, 住民登録番号, 姓名等)を用いてマッチングさせる方法である。単独あるいは複数のマッチングに利用した識別子の全てが一致した場合にだけ同一個人とみなす。そのため, 例えば番号や姓名のスペルに1つでも違いがあれば, 同一人物であってもマッチされることはない。よって, この方法では, 誤って他人とリンケージされる事は少なく同一人物との一致割合は高いが, 不正確なデータが多ければ一致割合が低くなる。また, 特異性の高い識別子が使用できない(データベース上に存在しない)場合にも, 同一人物の一致割合が低くなる事が起こりうる予想される。

使用された識別子としては, 国や地域で付与された健康情報を管理するための番号(カナダ^{11,12,37,46,61,73-80})やイギリス⁽⁴⁹⁾), 住民登録番号(デンマーク^(14,15,81))やノルウェー⁽³⁵⁾), 社会保障番号(米国^(4,39,82))などが使われていた。他に一意の識別子(“unique identification code”, “unique patient number”など)(イタリア⁽⁵¹⁾, 日本^(43,67), オランダ⁽⁶⁴⁾, 台湾⁽⁷⁾, カナダ^(8,10,18,24,43,59,83), 米国⁽⁸⁴⁾)を使用したと記載している研究も見られた。これらの研究では, これら一意の識別子の補足として他の識別子が使われている例もあった。複数の識別子を用いた deterministic な方法と考えられる。

② Probabilistic linkage:

複数の識別子を用いてマッチングさせる方法であり, 全ての識別子が一致する必要はない。この方法では, マッチングに用いる識別子毎に(同一人物と一致する確率に応じて)マッチングへ関与する重みを変え, その重みを足し合わせて一定の閾値を超した場合に, 同一個人とみなす。使用する識別子及び閾値を決定するのはリンケージを実施する研究者である。Deterministic な方法に比べて複雑なことから, probabilistic linkage を使用している研究は少なかった(台湾⁽⁶⁾, 米国⁽¹⁷⁾)。

リンケージの方法については, 詳細が記載されていない研究も多かった。特に単施設で実施された研究は明確な記載がない場合が多かったが^(19,21,26,38,56,57,63,84-89), 単施設ではカルテ番号や姓名等の匿名化されていない個人情報を用いることを前提にしている(ために敢えて記載しない)可能性が考えられる。また, 米国の同一 HMO 内で実施された研究^(19,43,52,90), 電子カルテ(Electronic Medical Records)における診断コードなどを, 同一のデータ源の検査値などによって評価した研究^(32,65)でも, リンケージの方法に特に言及されていないものが見られたが同様の理由によると考えられる。

(5) ゴールドスタンダードの設定

アウトカム定義の妥当性を確認するために参照となるゴールドスタンダードの種類は下記のように大別された。

① 疾患レジストリー:

国全体または地域単位で, がん^(49,67,82,91), 心筋梗塞^(11,18), 脳卒中^(7,17)などのよく見られる疾患から, 筋萎縮性側索硬化症⁽¹⁶⁾, 先天性奇形⁴⁶⁾, 自閉症⁽⁴⁸⁾, 小児糖尿病^(42,47,76), など比較的稀な疾患まで, レジストリーをゴールドスタンダードとした研究が見られた。

② 検査結果:

血清クレアチニン値に基づいた慢性腎不全⁽³²⁾及び急性腎障害⁽⁵⁹⁾, 血清ナトリウム値に基づいた

低ナトリウム血症⁽¹⁵⁾、経口ブドウ糖負荷(OGTT)試験に基づいた妊娠糖尿病⁽⁷³⁾をゴールドスタンダードとしている研究が見られた。また、子宮頸がんのバリデーション研究において、病理診断をゴールドスタンダードとしているものがあつた⁽⁹²⁾。

③ カルテレビュー:

最も頻度が多く、様々な疾患のバリデーション研究で行われていた。医師、特に対象疾患の専門家が直接カルテレビューをして判断しているものと、chart abstractor と呼ばれるカルテレビューのトレーニングを受けた人達が専門家の作ったチェック項目に沿って判断しているものに大別された。前者は診断が比較的難しい疾患や専門的知識を要する疾患(例:特発性肺線維症⁽⁹⁰⁾、多発性硬化症⁽⁷⁵⁾、化膿性汗腺炎⁽⁸⁵⁾、うつ⁽⁶⁵⁾の緩解、等)に多く、後者は比較的よく見られる疾患や診断基準が確立している疾患(例:脳卒中⁽⁴⁵⁾、上部消化管出血^(19,22)、てんかん^(71,93)、等)に多く見られた。2人の担当者が独立してカルテレビューをおこなっている研究が多かつたが、1人でおこなつた研究もあり、また記載がない論文も見られた。また、chart abstractor がスクリーニングを行い、可能性の高いと判断された症例を専門家がさらにレビューし最終判断した研究もあつた⁽⁷¹⁾。

その他の詳細事項についての記載も論文によって様々で、ゴールドスタンダードの詳細(実際に用いた調査票)を論文内または Appendix に示している模範的な論文^(4,53)もあれば、「カルテレビューにより目的の疾患の有無を判断しゴールドスタンダードとした」と一文のみしか記載していない論文も散見された。

④ 医師への質問票:

イギリスのプライマリーケアデータでは、傷病名を入力した医師に質問票を送り、その解答をゴールドスタンダードとする研究が多く見られた^(31,33,50)。これは、一応カルテレビューになるが、傷病名を入力したプライマリーケア医自身がそれを行っているため、判断にバイアスが生じている可能性が高い。

⑤ その他:

妊娠登録における妊娠中の糖尿病や喘息を薬の使用をゴールドスタンダードとして評価した研究のほか⁽³⁵⁾、性質の違う2つ以上の情報を組み合わせてゴールドスタンダードを設定した研究(例:糖尿病薬の処方歴、採血結果、または糖尿病病名のついた退院歴を糖尿病のゴールドスタンダードとした研究⁽⁵⁸⁾)や、インフルエンザ^(8,94)やうつ⁽⁶²⁾について住民に直接質問をし得られた解答をゴールドスタンダードとした研究、データベース上の診断コードをゴールドスタンダードとし薬剤の使用によるアウトカム定義を評価した研究⁽⁶⁴⁾、異なるデータベースの登録情報をゴールドスタンダードとした研究(例:保険者が有するデータに資格喪失理由として入力された死亡の情報をゴールドスタンダードとした日本のレセプトのバリデーション研究⁽⁷²⁾や、Nursing home Minimum Data Set への登録をケアホーム居住のゴールドスタンダードとした研究⁽¹⁾)もあつた。また、ほとんどのバリデーション研究ではゴールドスタンダードの定義は単一であつたが、敢えて複数のゴールドスタンダードの定義を示し、それぞれのゴールドスタンダードに対してアウトカム定義の妥当性を検討している研究が見られた^(34,44,94)(例:専門医または心臓カテーテルによる虚血性心疾患の診断を”hard”なゴールドスタンダード、一般医が診断した虚血性心疾患を”soft”なゴールドスタンダードとした⁽⁴⁴⁾)、疾患登録または(疾患登録にない時は)カルテレビューを行つて判断するとして複雑なゴールドスタンダードを用いた研究も見られた⁽⁷⁾。

(6) サンプリング方法とサンプルサイズ

サンプリング方法とサンプルサイズは、何をゴールドスタンダードとするかで大きく異なるため、ここでは下記の 6 つに分類した。

① カルテレビューをゴールドスタンダードとする場合：

陽性的中度のみを求める研究ではサンプルサイズは比較的小さいものが多い。米国の Mini-sentinel における陽性的中度のみを求めるバリデーション研究のデザインでは 95%信頼区間を±10%の精度で求めるためには 100 例程度のカルテレビューで十分であることが示されている⁽⁵⁴⁾。予測される陽性的中度を前提に求めたい 95%信頼区間から必要なサンプルサイズを計算した論文でも、必要症例数は 77 から 250 例で十分であることが示されている^(50,95)。

これに対し、陽性的中度のみならず、感度、特異度、陰性的中度をも求めようとする研究では、サンプルサイズは格段と大きくなる。カナダで行われた、ランダムサンプルで妥当性に関する 4 つの指標を求めるバリデーション研究には合計 9,500 例のカルテレビューを実施した関節リウマチのバリデーション研究⁽⁹⁶⁾、7,500 例のカルテレビューをした epilepsy のバリデーション研究⁽⁷¹⁾、5,000 例のカルテレビューをした脳卒中/TIA のバリデーション研究⁽⁴⁵⁾が含まれる。また単一の病院においてではあるが、25,761 例のカルテレビューにより、会陰裂傷のバリデーション研究がノルウェーで実施されている⁽²⁶⁾。

カルテレビューを前提にサンプルサイズを小さくする試みも行われており、ノルウェーの出生レジストリーの、カルテレビュー（および退院時の診断コード）をゴールドスタンダードとする妊娠高血圧腎症(preeclampsia)のバリデーション研究では、アウトカム定義を満たす 3,500 例全例とアウトカム定義を満たさない 75,311 例からのランダムサンプル 1,840 が検討され、サンプルの割合の異なりを適切に補正して妥当性に関する 4 つの指標を求めている⁽³⁶⁾。アウトカム定義を満たす集団と満たさない集団から異なる割合でサンプルを求める研究は、他にもいくつか見られる。例えば、カルガリー大学耳鼻科の慢性鼻副鼻腔炎のバリデーション研究ではアウトカム定義を満たす症例と満たさない症例各 100 例⁽¹⁰⁾、スペインの脳卒中のバリデーション研究ではアウトカム定義を満たす 300 例と満たさない 100 例⁽²⁷⁾、アメリカの全米の退役軍人 500 万人のデータベースを用いて行われた非アルコール性脂肪肝疾患のバリデーション研究ではアウトカム定義を満たす 450 例と満たさない 150 例⁽²⁰⁾をカルテレビューで検討している。これらの研究では、感度と特異度を計算する際の補正が適切に行われておらず、示された妥当性に関する指標は陽性的中度と陰性的中度以外は適切とは思われないが、カルテレビューをする症例数を減らすために、アウトカム定義を満たす症例と満たさない症例を異なる割合でサンプルした例と考えられる。

また、カナダにおける有病割合の低い多発性硬化症のバリデーション研究では、ランダムサンプルでは妥当性に関する指標を検討する上で十分なケースが得られないとの理由から、①EMR の cumulative patient profile (CPP)、②free text entries for MS-related phrases、③physician billing codes)から多発性硬化症の全症例を含むと考える”all possible cases”943 例を特定し、カルテレビューにより 247 例の多発性硬化症の症例を特定している⁽⁷⁵⁾。同様のアプローチは、ICD コードとテキスト検索により腸重積症例 565 例を特定したカナダにおける研究⁽⁶¹⁾、アメリカの Reasons for Geographic and Racial Differences in Stroke (REGARDS)における脳卒中のバリデーション研究で

Medicare の診断コードと脳卒中を疑わせる症状の自己報告のいずれかに該当する例のカルテレビューで真のケースを特定した研究⁽²⁾、3つのデータベースのいずれかで特定された Neural Tubal Defect による死亡、死産をさらにカルテレビューで真のケースを特定した研究⁽⁹⁷⁾、などでもとられている。特に REGARDS コホートを利用した研究では、primary analysis では特定されたケースを全ケースとしているが、特定されていない真のケースがどの程度存在すると、妥当性の指標にどのような影響を与えるかの感度解析が行われている⁽²⁾。”all possible cases”のアプローチに区分可能な研究としては、複数のサブグループから異なる抽出割合で得たサンプルのカルテレビューによって急性心不全を特定した研究もあるが⁽⁹⁸⁾、抽出割合が異なるために、指標の求め方は複雑である。

② 調査票またはインタビューの結果をゴールドスタンダードとする場合：

この方法を用いた場合、調査されたサンプルサイズは 500 例以下がほとんどであった。郵送による調査票の回答は英国で多く用いられている。英国ではプライマリ医によって登録される医療情報データベース(The Health Improvement Network、Clinical Practice Research Datalink など)が多くの医学研究に利用されており、これらのデータベース上の傷病名の妥当性研究では、該当のプライマリ医に送付された調査票の結果が用いられている。大腸がん⁽²⁹⁾、膀胱がん⁽³⁰⁾、乾癬性関節炎⁽³¹⁾、c 型肝炎⁽³³⁾などの疾患が対象で、いずれも 100 から 200 例程度のサンプルサイズとなっている。

またインタビュー調査の結果は、インフルエンザワクチン接種の妥当性研究で使われている。オーストラリアの妊婦を対象にした妥当性研究⁽⁹⁴⁾では、協力の得られた出産した女性に対して実施された電話インタビューによるセルフレポートが使われ、サンプルサイズは 563 人であった。一方、カナダのオンタリオ州で実施された研究⁽⁸⁾では、既に実施されていた地域住民を対象としたインタビュー調査の結果が使われたため、サンプルサイズは 47,301 人と大きかった。

③ 電子カルテなどの診療記録をゴールドスタンダードとする場合：

電子的に記録された情報からは、コードなどで該当の情報を特定する事が容易であることから、手間と時間を省くことが可能であり、カルテレビューに比べてサンプルサイズを大きく設定することが可能である。事実、数万～数十万の規模のサンプルサイズを扱った研究もあった^(9,24)。

④ 疾患登録をゴールドスタンダードとする場合：

疾患が電子的に登録されている場合には情報の特定が容易であり、大きなサンプルサイズでの研究が実施可能である。カナダで実施された妊娠糖尿病の妥当性研究⁽³⁷⁾では州の周産期レジストリーが用いられ、そのサンプルサイズは 411,390 例であった。該当症例全例が使用できたため、感度・特異度・陰性的中度・陽性的中度の全てが算出可能であった。他にもがん登録が用いられたイギリスの General Practitioner Research Database (GPRD)のがんの妥当性研究⁽⁴⁹⁾のサンプルサイズは 42,556 例と大きい。一般的には疾患登録の精度は高いと考えられるが、地域全体のがん登録情報を用いた研究で、感度・特異度・陽性的中度・陰性的中度を算出できても、陽性的中度が低い場合がある。米国のペンシルバニア州で州のがん登録が使用された研究⁽⁸²⁾では陽性的中度が低く、その理由として疾患登録の不完全性の可能性を指摘し、がん登録で見逃された患者の割合とそれが陽性的中度に与える影響を検討する感度分析を行っていた。また、単施設等の小規模での妥当性研究で、地域の疾患登録ではなく、当該の医療機関の疾患登録を使用している

場合はその分サンプルサイズが小さい。例えば、単施設で実施された小児糖尿病の妥当性研究⁽⁴²⁾では 1,323 名、乳がんの妥当性研究⁽⁶⁷⁾は 633 例であった。その他使われた疾患登録には、小児糖尿病疾患登録⁽³⁸⁾や先天性奇形疾患登録⁽³⁷⁾なども使われていた。

⑤ 検査値や画像検査の結果をゴールドスタンダードとする場合:

画像診断の結果や、血液検査等の測定値が使用されており、単独あるいは少数の医療機関で実施された研究ではこれらの情報が入手しやすかったと考えられる。それらの研究のサンプルサイズは数百例と小さかった^(59,70,81,99,100)。約 25,000 症例の大きなサンプルサイズで実施された研究は、妥当性の確認の対象となるレセプトデータに元々連結している電子カルテがあり、参照が容易であったと考えられる⁽⁹²⁾。

⑥ その他:

レセプトや電子カルテ、検査値などの複合や、アルゴリズムの作成、保険資格喪失理由の記録や、複数のデータベースの記録の有無、疾患スクリーニングの結果などが使用されており、サンプルサイズも様々であった^(41,48,72,101)。

(7) 妥当性を測定する指標

妥当性 (validity) の指標として、陽性的中度、陰性的中度、感度、特異度が用いられていた。その内、求めた指標 (の組み合わせ) は以下の 2 パターンに大別された。

① 陽性的中度

Mini-Sentinel^(53,54)、THIN^(30,33)や GPRD⁽³⁴⁾等のデータベースにおけるバリデーション研究では陽性的中度のみを求めていた。ゴールドスタンダードはカルテレ뷰が多く、検査結果⁽⁹²⁾や General practitioner (GP)⁽³³⁾からの質問票の回答とした研究も見られた。いずれの研究もアウトカム定義を満たす患者群からのみランダムサンプリングを行っており、感度や特異度は求められていない。

② 陽性的中度、陰性的中度、感度、特異度

アウトカム定義を満たす患者群および満たさない患者群からサンプリングを行った研究では、4 つ全ての指標を求めていた。しかし、指標全てを求められるにも関わらず、陰性的中度を求めている研究、陰性的中度および特異度を求めている研究があった⁽²⁹⁾。また、一病院で特定された真のケースとその病院の患者の administrative database での該当コードを照らし合わせて感度と陽性的中度のみを求めている研究もあった⁽⁷⁴⁾。

偽陽性、偽陰性の症例については、偽陽性・偽陰性となった理由などの詳細を表などに記載している研究もみられた^(4, 7,8,38,46,49,51,73,74,102,103)。

4 つの指標以外に、複数人でカルテレ뷰を行っている研究では κ 係数を算出した研究も見られた^(5,12,44,65,80,104,105)。異なる評価者間の一致度 (inter-rater agreement) のみならず、同一の評価者の異なる時点での一致度 (intra-rater agreement) を評価した研究もみられた⁽⁸⁾。

(8) 指標の閾値・利用法

感度や陽性的中度の値としてどのレベルが望ましいか、または指標の利用方法については、殆どの論文では言及はなかったが、ここでは記載のあった少数の研究を紹介する。目標とする陽性的中度を 80%以上⁽⁸⁹⁾や 85%以上⁽⁸⁰⁾と予め定めている研究や、複数のアルゴリズムを比較する際に

陽性的中度 70%を「低い」⁽⁶⁷⁾、67%未満を「低い」⁽⁸¹⁾と判断している研究があった。予防接種有無の自己報告を実際の接種記録と比較した研究⁽⁸⁾では、低い感度(49.8%)かつ高い陽性的中度(88.4%)という集団であっても、自己申告で既接種と回答した部分集団は予防接種の安全性を判断するための自己対照研究(Self-Controlled Study)には用いることができる、と筆者らは考察している。また、パーキンソン病／パーキンソン症候群の疫学調査のためのバリデーション研究で「陽性的中度が高く、かつ感度特異度も高い指標を選択した」旨が記載されている研究もあった⁽⁹⁾。

本報告書の「2. バリデーション研究の概要」の「バリデーション研究に関する教科書のまとめ (Pharmacoepidemiology, 5th Edition (Brian L. Stromら)」の紹介でも言及されている通り、同一の感度や特異度であっても、薬剤疫学研究の結果への影響は研究集団における曝露の割合等で変わりうる。指標の望ましい値に言及されている研究が少ないのは、「望ましい値」は研究対象集団、研究の目的などで異なりうることを反映していると考えられる。

参考文献

別途エクセルファイル(付録3 文献リスト.xlsx)参照

「日本における傷病名を中心とするレセプト情報から得られる指標のバリデーションに関するタスクフォース」報告書

付録 4: バリデーション研究におけるサンプルサイズ

[1]陽性的中度 (PPV) のみを求めるバリデーション研究

一般にアウトカム指標の定義を満たす n 例のサンプルのカルテレビューを行い a 例の真のケースが見出されれば、PPV は

$$PPV = \frac{a}{n} \quad (1)$$

であり、その 95%信頼区間 (95%CI_{PPV})

$$95\%CI_{PPV} = PPV \pm \delta_{PPV} \quad (2)$$

の δ_{PPV} (絶対精度、absolute precision) は二項分布の正規近似によって

$$\delta_{PPV} = 1.96 \sqrt{\frac{PPV(1-PPV)}{n}} \quad (3)$$

と与えられる。 δ_{PPV} の値が最大であるのは $PPV=0.5$ の場合であり、逆に、絶対精度をある δ_{PPV} の値以下にするための n の大きさは

$$n = \left(\frac{0.98}{\delta_{PPV}}\right)^2 \sim \frac{1}{\delta_{PPV}^2} \quad (3)$$

である¹。PPV の 95%信頼区間を ± 0.1 以内としたければ 100 例 ($1/0.1^2$)、 ± 0.05 以内としたければ 400 例 ($1/0.05^2$) が必要である。

[2]陽性的中度 (PPV)、陰性的中度 (NPV)、感度、特異度を求めるバリデーション研究

[2-1]単純なランダムサンプル

バリデーション研究を実施する集団全体から単純なランダムサンプルによって得た n 例のカルテレビューを行い、陽性的中度 (PPV)、陰性的中度 (NPV)、感度、特異度を本報告書の表 1 に示した式を用いて求める場合、95%信頼区間が最も広くなるのは通常は感度である。ある特定の絶対精度 (δ_{sen}) で感度を求める場合には、上記式(3)と同様にサンプルに含まれる真のケース (表 1

の $a+c$) の大きさが $\frac{1}{\delta_{sen}^2}$ となることが求められる。 n 例の単純なランダムサンプルに含まれる真のケースの割合は、集団における真のケースの割合 (有病割合、prevalence、以下 $prev$ と略す) と同一とみなすことができ、 $n \cdot prev$ の大きさが $\frac{1}{\delta_{sen}^2}$ に等しいような n が必要である。すなわち必要なサンプルサイズは

$$n = \frac{1}{prev} \frac{1}{\delta_{sen}^2} \quad (4)$$

と与えられる。たとえば、真のケースを 1%含む集団からの単純なランダムサンプルで感度の 95%信頼区間を ± 0.1 以内としたければ 10,000 例 ($(1/0.01)(1/0.1^2)$) のカルテレビューが必要である。実際、関節リウマチ (RA) に関して、「RA の診断による入院または 2 年以内に 3 回の外来診療における

RA の診断をもち、うち 1 例は専門家による」という定義の感度が±0.1 (0.77、95%CI:0.67-0.87)であることを示したカナダで実施されたランダムサンプルによるバリデーション研究では、7,500 人のカルテレビューが実施されている²。この研究では当該定義による TP のケースは 53 例、FN は 16 例だったので、有病割合は 0.92%程度であったと考えられる。

[2-2]層別サンプル

バリデーション研究を実施する N 人の集団のうち、アウトカム指標”Yes”の者が N_1 人、“No”の者が N_0 人である場合に、これら N_1 人と N_0 人から異なる割合でサンプルを得る層別サンプル (stratified sampling)によるバリデーション研究も可能である。

ただし、この方法への過大な期待は禁物であり、層別サンプルを使って、感度をたとえば±0.1で求めるのに必要なサンプルサイズは、上述の[2-1]の単純なランダムサンプルで感度を±0.1で求める場合のサンプルサイズと通常大きくは変わらない。この方法が有用なのは、たとえば PPV については必要な精度の情報を得る一方、感度についても、絶対精度は若干低くても何らかの情報を得ておきたい、というような場合である。たとえばノルウェーにおける妊娠レジストリにおける子癩前症(preeclampsia)のバリデーション研究では、この疾患コードをもつ 3,500 例全例と疾患コードをもたない 75,311 例の 2.4%の 1,840 例のレビューが antenatal chart と退院時の診断コードをゴールドスタンダードとして実施されている。疾患コードをもつ 3,500 例のうち、真のケースは 2,936 例であり、PPV=0.84、その 95%信頼区間は (0.83, 0.85) であるが、感度については 0.43、その 95%信頼区間は (0.39, 0.48) であり、PPV については 95%信頼区間が±0.01 で求められたのに対し、感度の 95%信頼区間については±0.04 であった。

本 TF のメンバーなどによる最近の研究において、層別サンプルを用いたバリデーション研究では二段階で研究を実施することが推奨されている⁴。第一段階 (Step I) においては[1]と同様にアウトカム定義を満たす者 N_1 人からランダム抽出した n_1 人のサンプルについてカルテレビューを実施し、真のケース a 人と非ケース b 人を得て ($a+b=n_1$)、PPV (a/n_1) を推定する。次いで第二段階 (Step II) において、以下の式で与えられるアウトカム定義を満たさない者 N_0 人からの n_0 人のサンプルについてカルテレビューを実施し、真のケース c 人と非ケース d 人を特定する ($c+d=n_0$)。

$$n_0 = \frac{N_0}{N_1} \frac{1}{PPV} f * \quad (5)$$

式(5)で $f*$ の値は感度の絶対精度として期待される値” δ_{sen} ”と”a”の 2 つによって決まり、本付録末尾に掲げた表にその値を示した。本付録末尾の表における”a”の値が、Step I において n_1 人に含まれる真のケースの数よりも多い場合には、アウトカム定義を満たす者について追加のカルテレビューを行う。たとえば、Step I で PPV の 95%CI を±0.1 で求めるために 100 人のカルテレビューを実施し、真のケースが 80 人 ($a=80, b=20$) で PPV=0.8 と推定されたとする。この場合、得られた PPV の 95%信頼区間は二項分布の正規近似により 95%信頼区間 (0.72-0.88) と計算され±0.1 の範囲内だが、本付録末尾の表を利用し、たとえば a については $a=80$ ではなく $a=100$ を用いるのであれば、 $a=100$ が満たされるように、アウトカム定義を満たす者 25 名程度のカルテレビューを追加で実施する。本付録末尾の表の $f*$ の算出法については元論文⁴を参照されたい。

n_1 例と n_0 例のカルテレビュー実施し、a、b、c、d の値を得た後の PPV と NPV および、その 95%

信頼区間の推定方法は、単純なランダムサンプルの場合と同様であり、標準的な教科書における割合 (proportion) とその 95% 信頼区間の推定方法を参照する。

また、感度 (sen) と特異度 (spe) は以下のように求めることができる。

$$\text{sen} = \text{PPV} \cdot N_1 / [\text{PPV} \cdot N_1 + (1 - \text{NPV}) N_0] \quad (6)$$

$$\text{spe} = \text{NPV} \cdot N_0 / [(1 - \text{PPV}) N_1 + \text{NPV} \cdot N_0] \quad (7)$$

その 95% 信頼区間については、bootstrap 法を用いるか、以下の近似式⁴を用いる。

$$95\% \text{CI}_{\text{sen}} = \text{sen} \cdot \exp \left[\pm 1.96(1 - \text{sen}) \sqrt{\frac{1}{a} + \frac{1}{c}} \right] \quad (8)$$

$$95\% \text{CI}_{\text{spe}} = \text{spe} \cdot \exp \left[\pm 1.96(1 - \text{spe}) \sqrt{\frac{1}{b} + \frac{1}{d}} \right] \quad (9)$$

参考文献

1. McNeil D. Epidemiological research methods. John Wiley & Sons; 1996, 266 p.
2. Widdifield J, Bombardier C, Bernatsky S, et al. An administrative data validation study of the accuracy of algorithms for identifying rheumatoid arthritis: the influence of the reference standard on algorithm performance. BMC Musculoskelet Disord. 2014;15:216.
3. Klungsøyr K, Harmon QE, Skard LB, et al. Validity of pre-eclampsia registration in the medical birth registry of Norway for women participating in the Norwegian mother and child cohort study, 1999-2010. Paediatr Perinat Epidemiol. 2014;28:362-71.
4. Kubota K, Iwagami M, Yamaguchi T. Designing the validation study with the stratified sampling. In submission.

δ_{sen}	a=50	a=60	a=70	a=80	a=90	a=100	a=120	a=150	a=200	a=300	a=500
0.10	80.46	69.26	63.58	60.14	57.84	56.18	53.96	52.01	50.26	48.70	47.56
0.11	58.25	52.69	49.59	47.60	46.22	45.21	43.81	42.55	41.40	40.35	39.57
0.12	45.32	42.11	40.22	38.96	38.07	37.40	36.46	35.60	34.80	34.06	33.51
0.13	36.80	34.77	33.52	32.67	32.05	31.59	30.93	30.32	29.75	29.21	28.80
0.14	30.75	29.38	28.50	27.90	27.46	27.13	26.65	26.20	25.77	25.37	25.06
0.15	26.25	25.26	24.63	24.19	23.87	23.61	23.26	22.91	22.59	22.28	22.04
0.16	22.76	22.04	21.56	21.23	20.98	20.79	20.51	20.25	19.99	19.75	19.56
0.17	19.99	19.44	19.08	18.82	18.63	18.48	18.26	18.05	17.85	17.65	17.50
0.18	17.75	17.32	17.03	16.83	16.67	16.56	16.38	16.21	16.05	15.89	15.77
0.19	15.90	15.56	15.33	15.16	15.04	14.94	14.80	14.66	14.53	14.40	14.30
0.20	14.35	14.07	13.89	13.75	13.65	13.57	13.45	13.34	13.23	13.12	13.04

アウトカム定義を満たす者からのランダムサンプルのカルテレビューにより真のケースが a 人特定されていることを前提に、アウトカム定義を満たさない者から、本表に示された値(f*)を用いて得られる $((N_0/N_1)/PPV \text{ f*})$ のサイズのサンプルをランダムサンプルし、カルテレビューを実施すれば、感度の 95%信頼区間は推定値 $\pm \delta_{sen}$ 以内である。詳しくは元論文⁴参照。

バリデーション研究のプロトコルや論文中に記載すべき情報一覧(チェックリスト)				
チェック項目	内容	確認結果		
		はい	いいえ	不明
リサーチクエストとプロトコル	リサーチクエスト(バリデーション研究後の使用目的も含む)が明記されているか。			
	ガイドライン(Guidelines for Good Pharmacoepidemiology Practices など)にそってプロトコルが作成されているか。			
倫理審査と研究計画の登録	研究実施に関する倫理審査委員会による承認について記載されているか。			
	研究計画の事前登録(UMIN 臨床試験登録システムなど)がなされているか。			
パイロット研究	少数例を対象としたパイロット研究が予定されているか。			
(1)利用する医療情報データベースの種類	データベースの種類(入院・外来の診療報酬請求データ・退院時記録・電子カルテ・疾患レジストリー、またはそれらの組み合わせ)が明記されているか。			
(2)バリデーション研究のセッティング	対象の地域が定義されているか。			
	データの期間が明記されているか。			
	調査された医療機関の数が明記されているか。			
	医療機関の特徴(大学病院、救急指定病院、がん拠点病院、受診患者数、病床数、等)が明記されているか。			
	研究対象集団の特徴(年齢・性別構成、重症度の分布、等)が明記されているか。			
	バリデーション研究の対象集団と、利用するデータベースの想定する“population”の差異についての検討・考察がされているか。			
(3)アウトカムの定義	どの情報を使うか(傷病名のみ、薬剤名のみ、処置名のみ、検査結果のみ、またはその組み合わせによるアルゴリズム、など)。			
	傷病名をアウトカム定義に使う場合、傷病名コードの種類(傷病名マスター、ICD-10 コード、等)、選択した傷病名コードのリスト、傷病名入力の期間・回数に関する条件(例:異なる月に連続して2回、半年以内に2回、等)の記述があるか。			

(4)リンケージの有無・方法の確認	リンケージ(照合)を行うか(注: 各医療機関の中で保存しているデータを比べる作業はリンケージとはみなさない)。			
	リンケージ(照合)を行った場合、リンケージのタイプ(deterministic または probabilistic)、使用した因子(Social security number、年齢・性別・生年月日・居住地域、等)、リンケージが出来なかった人の数や割合などが記載されているか。			
	Probabilistic リンケージの場合、リンケージの質について記載されているか。			
(5)ゴールドスタンダードの定義	ゴールドスタンダードの種類(疾患レジストリー、検査結果、カルテレビュー、等)が明記されているか。			
	疾患レジストリーの場合、それ自体の妥当性に関する研究が過去に実施されている場合はその引用、実施されていない場合は疾患レジストリーがゴールドスタンダードとして妥当であるとの判断が合理的である旨の記述があるか。			
	検査結果の場合、検査結果が正確に測定・標準化・記録されているかの記載があるか。			
	カルテレビューの場合、対象疾患の判断基準または調査票を明示し、評価者の専門性を記載し、2人(あるいは少数)の評価者が別々に判断、評価者間の一致度を表す κ 係数を計算し、一致しなかった評価結果の取り扱いについて明記されているか。			
(6)サンプリング方法とサンプルサイズの設定(右のいずれかを選択する)	[A]全患者を対象にする場合:アウトカムとゴールドスタンダードの定義を満たす患者を、それぞれコンピューターを用いて簡単に同定できるかを確認する(できない場合はデータの標準化・電子化による作業の効率化が必要)。			
	[B]全患者からランダムサンプリングを行う場合:アウトカム定義を満たす集団(陽性のケース)、満たさない集団(陰性のケース)、ゴールドスタンダードの定義を満たす集団(真のケース)、満たさない集団(偽のケース)をそれぞれ最低100人(求める指標の95%信頼区間が最大±10%に対応)ずつ確保できるようサンプル数を計算する(確保が難しい可能性が最も高い「真のケースが最低100人」になるように計算する)。			
	[C]アウトカム定義を満たす患者群からのみランダムサンプリングを行う場合:最低100例			

	が目安となる。			
	[D]アウトカム定義を満たす患者群と満たさない群からランダムサンプリングを行う場合：アウトカム定義を満たす患者群からは最低 100 例、アウトカム定義を満たさない患者群からは付録 4 の近似式などを参考にサンプルサイズを計算する。			
	[E]“all possible cases”を想定したサンプリングを行う場合：研究対象集団（全体またはランダムサンプリングした症例）の中で真のケースを多く含むことが期待されるサブグループを 2 種類以上特定し足し合わせることでサブグループの集合体を作る。複数のサブグループからランダムサンプリングをする際は、抽出率を同じ割合にして、最終的に同定される真のケースが最低 100 例以上になるよう必要サンプル数を見積もる。			
(7)妥当性を測定する指標の計算	妥当性を測定する指標として求めるものは何かの説明されているか。			
	[A]全患者を対象にした場合、または、[B]全患者からランダムサンプリングを行った場合：2×2 テーブルが作成され、感度・特異度・陽性的中度・陰性的中度が計算されているか。			
	[C]:アウトカム定義を満たす患者群からのみランダムサンプリングを行った場合：陽性的中度が計算されているか。			
	[D]アウトカム定義を満たす患者群と満たさない群からランダムサンプリングを行った場合：ウェイトを考慮した 2×2 テーブルが作成され、感度・特異度・陽性的中度・陰性的中度が計算されているか。			
	[E]“all possible cases”を想定したサンプリングを行った場合：得られた情報をもとに 2×2 テーブルが作成され、感度・特異度・陽性的中度・陰性的中度が計算されているか。			
(8)指標の閾値・利用法に関する考察	同じ研究セッティングで複数のアルゴリズムを比較する場合、著者らが考える最適なアルゴリズムについての考察がなされているか。			