

令和7年国勢調査におけるCANCEIS補完の実装に向けて （参考資料）

令和6年3月
総務省統計局

○ **ホットデック法**は、欠測値や矛盾値を同一のデータセット内のデータから一定の方法で補完する方法

○ **単純なホットデック法の運用では、エディットルールを充足する保証がない**

(補完前後でデータが変わり、補完前にパスしたチェックを補完後にパスしない可能性がある)

例) 年齢、配偶関係、産業が不詳の場合 ⇒ 慎重にエディットしないと、年齢を14歳に補完した場合などに配偶関係や産業に矛盾が生じる

<ホットデック法で性別のみ補完する例>

ID	世帯主との続き柄	性別	年齢	性別の補完列
1	1 世帯主	1 男	39	性別 = 1
2	2 配偶者	2 女	35	性別 = 2
3	3 子	1 男	13	性別 = 1
4	3 子		10	性別 = 1
5	4 その他の親族	2 女	40	性別 = 2
6	4 その他の親族	1 男	空白	性別 = 1
7	4 その他の親族	2 女	13	性別 = 2
8	5 親族以外		空白	性別 = 2
9	5 親族以外	1 男	44	性別 = 1
10	5 親族以外	2 女	36	性別 = 2

<ID順に処理>

① 処理の都度、性別の補完列を置換

- ※ 例えば、ID 2 に到達した際、ID 2 の「性別」欄のデータ (性別 2) に性別の補完列を置換
- ※ 3 人目は男性なので、性別の補完列は再び 1 になる。

② 空白のレコードに遭遇した際、その直前の性別の補完列の値を補完

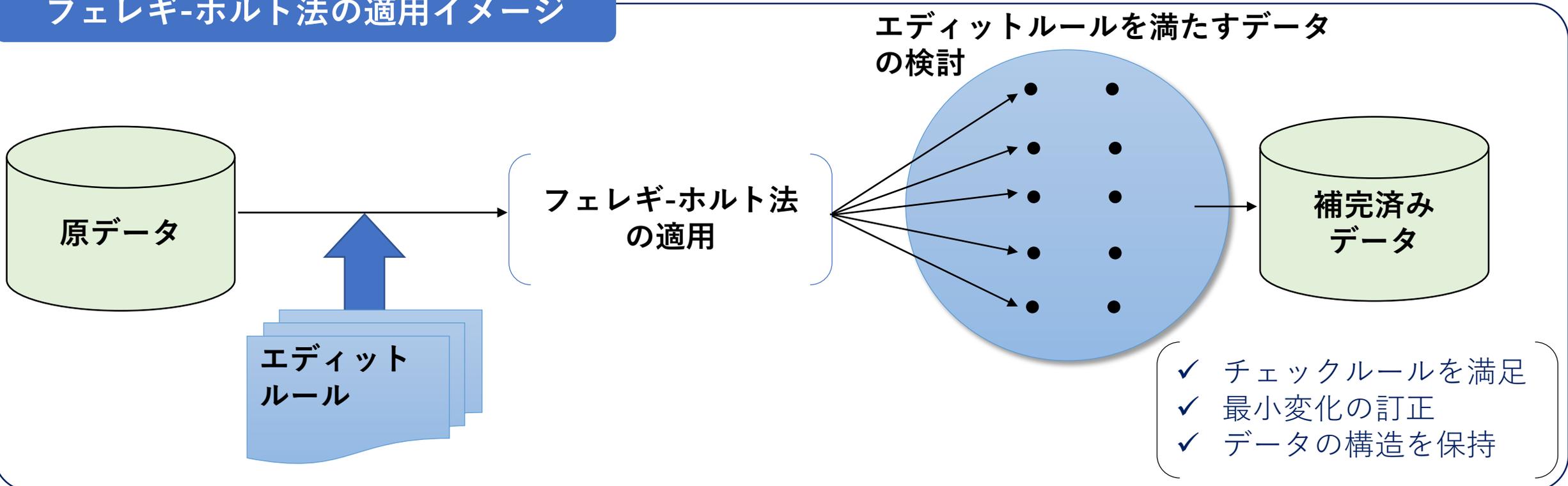
- ※ 4 人目の性別は空白のため、性別の補完列の値 (この場合は「1 男性」) から、空白に「1 男性」を補完

③ 5 人目は女性なので性別の補完列を「2 女性」に置換

男女はほぼ同頻度で出現するため、最終的な補完回数は男女で半数ずつになると想定

- フェレギ・ホルト法は、補完により変更されたレコードが**エディットに失敗しないことを保証する全体的なモデル**であり、以下の原理により、データの欠測値を修正するためのアルゴリズムを提供
 - ✓ 各レコードが**全てのチェックルールを満足する**
 - ✓ **できるだけ少ない変更で訂正が達成される**
 - ✓ **データの構造を保持する補完手順**となっている
- これにより、上記の利点を保持した上で、ドナー候補を探索し、欠測値を補完可能
- **ただし、エディットする変数を特定し、ドナーを探索する流れはデータ処理上、非効率**

フェレギ-ホルト法の適用イメージ



CANCEISのドナー検索・補完の流れ

○ CANCEISは、**要補完ユニット**のドナー候補として、要補完ユニットと距離が近く、類似性の高い複数のユニット（**NMCIA** ^(注1)）をリスト化し、その中から最終ドナーをランダムに選択して補完

注1) NMCIA・・・Near Minimum Change Imputation Actions

ドナー検索・補完の流れ

※ 単身世帯の例

要補完ユニット (ID=1)

ユニット(世帯)番号	世帯員番号	年齢	配偶関係
1	1	25	

第1ステージ

✓ 要補完ユニットの近隣の500ユニットからドナーを検索

ユニット(世帯)番号	世帯員番号	年齢	配偶関係
5	1	34	離別
8	1	27	未婚
⋮	⋮	⋮	⋮

✓ ドナー候補の中から、NMCIAを10ユニット順次選定

ユニット(世帯)番号	世帯員番号	年齢	配偶関係
20	1	23	未婚
31	1	25	未婚
⋮	⋮	⋮	⋮

✓ NMCIAの中で最適な補完を実現する補完アクション (IA) を特定

ユニット(世帯)番号	世帯員番号	年齢	配偶関係
31	1	25	未婚

第2ステージ (検索範囲の拡大)

✓ 第1ステージとは別の500ユニットからドナーを検索

ユニット(世帯)番号	世帯員番号	年齢	配偶関係
510	1	54	離別
513	1	47	未婚
⋮	⋮	⋮	⋮

✓ NMCIAを更新 (第1ステージより優れたドナーがある場合)

ユニット(世帯)番号	世帯員番号	年齢	配偶関係
520	1	25	未婚
31	1	25	未婚
⋮	⋮	⋮	⋮

✓ IAより有意に(注2)優れたIAがなければ検索を終了し、第2ステージ終了時点のNMCIAからランダムに最終ドナーを選定

✓ そうでなければ第3ステージに進み、1,000ユニットを検索

✓ 以下同様に最大10ステージまで検索 (3ステージ以降、検索数はステージごとに倍増)

注2) 第1ステージの最適なIAに比した第2ステージのIAの品質 (P8の品質評価式の値) の改善が10%未満の場合、それ以上のステージの追加は無意味と判断し、検索を終了

「距離が近い」ユニット（距離関数）

○ CANCEISの距離関数は、あるユニット（世帯）の調査票への回答内容と、別のユニットの回答内容との近さを定量化するための指標であり、年齢、国籍、世帯主との続き柄などの質的な相違も加味して算出※

※ ユニット間の地理的な近さも加味

D_{fp} : 要補完ユニット(FU)と通過ユニット(PU)の距離

V_{fi} : 要補完ユニットにおける*i*番目の変数

V_{pi} : 通過ユニットにおける*i*番目の変数

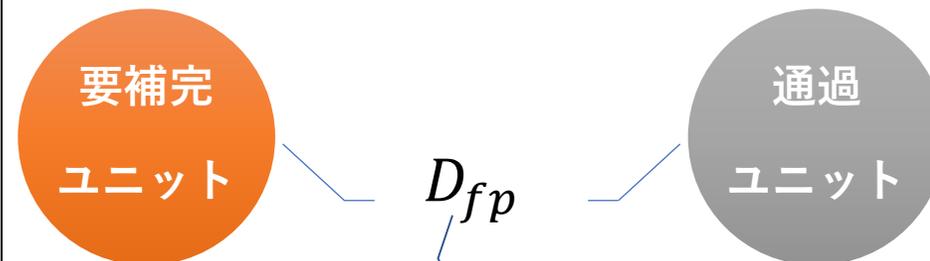
$D_i(V_{fi}, V_{pi})$: V_{fi} と V_{pi} の距離

w_i : *i*番目の変数におけるウエイト

$$D_{fp} = \sum_i w_i D_i(V_{fi}, V_{pi})$$

具体例

$$D_i = \begin{cases} 0, & \text{if } V_{fi} = V_{pi} \\ 1, & \text{otherwise} \end{cases}$$



D_{fp} が小さいほど、FUとPUが類似している。

※ V_{fk} と V_{pk} の値が一致するかしないかにより、 $D_k(V_{fk}, V_{pk})$ を0～1の数値に設定する。

$$D_{fp} = w_1 D_1(V_{f1}, V_{p1}) + w_2 D_2(V_{f2}, V_{p2}) + \dots + w_n D_n(V_{fn}, V_{pn})$$

ユニット	年齢	国籍	...	世帯の種類
要補完ユニット	V_{f1}	V_{f2}	...	V_{fn}
通過ユニット	V_{p1}	V_{p2}	...	V_{pn}

✓ 例えば、国籍が両ユニットとも日本人なら D_2 は0で、それにウエイトを掛ける。

距離のイメージ

- ユニット番号 1 を要補完ユニットとすると、ドナー候補は、世帯員数が同一のユニット 4 及び 5 がドナー候補
- これを要補完ユニットとの距離で表すと、ユニット 4 との距離 (0) の方がユニット 5 との距離 (7) より近いと判断

※ 簡単化のため、いずれのウエイトも 1 とした

回答データ

距離関数で出した距離

ユニット (世帯)番号	世帯員番号	年齢	国籍	続き柄	配偶関係	年齢	国籍	続き柄	配偶関係
1	1	45	日本人	世帯主	配偶者あり	-	-	-	-
1	2	46	日本人	配偶者		-	-	-	-
1	3	18	日本人	子	未婚	-	-	-	-
1	4	16	日本人	子		-	-	-	-
2	1	55	日本人	世帯主	配偶者あり	-	-	-	-
2	2	20	日本人	配偶者	配偶者あり	-	-	-	-
3	1	19	日本人	世帯主	未婚	-	-	-	-
4	1	45	日本人	世帯主	配偶者あり	0	0	0	0
4	2	46	日本人	配偶者	配偶者あり	0	0	0	0
4	3	18	日本人	子	未婚	0	0	0	0
4	4	16	日本人	子	未婚	0	0	0	0
5	1	45	アメリカ	世帯主	配偶者なし	0	1	0	1
5	2	46	アメリカ	兄弟姉妹	配偶者なし	0	1	1	0
5	3	18	アメリカ	子	未婚	0	1	0	0
5	4	16	アメリカ	他の親族	未婚	0	1	1	0

ユニット 1 と 4 の距離
 $D_{f4} = 0$

ユニット 1 と 5 の距離
 $D_{f4} = 7$

ドナー候補による補完の品質 (Quality of Imputation Action)

- CANCEISにおける各ドナー候補による補完の品質は、(1)原データからの最小変化性(D_{fa} : ③と①の距離)及び(2)実在可能性(D_{ap} : ③と②の距離)について下式で定義され、ドナー検索の結果、下式の値が小さい10ユニットがNM CIAに登録される。
- NM CIAの中からランダムに1ユニットが、最終的なドナーに選定される。

※ なお、我が国のように、CANCEIS補完において欠測値補完のみ行う場合は、原データからの最小変化性は担保されるため、 D_{fa} はゼロとなり、 D_{fpa} は D_{ap} のみに依存する。

<品質評価式>

$$D_{fpa} = \alpha D_{fa} + (1 - \alpha) D_{ap}$$
 加重平均
 ①要補完ユニットと ③仮補完ユニットと
 ③仮補完ユニットとの距離 ②通過ユニットとの距離

*) α : ユーザー定義システムパラメータ ($0.5 < \alpha \leq 1.0$)

単身世帯のユニット1 (①要補完ユニット) に対し、近隣の単身世帯500ユニットの品質を確認する事例

①補完前のデータ(要補完ユニット)

ユニット(世帯)番号	世帯員番号	年齢	住宅の建て方	性別	続き柄	配偶関係
1	1	25	共同住宅	男	世帯主	
5	1	34	一戸建て	男	世帯主	離別
8	1	27	一戸建て	男	世帯主	未婚
12	1	35	一戸建て	女	世帯主	配偶者あり
15	1	28	共同住宅	女	世帯主	未婚
20	1	23	共同住宅	男	世帯主	未婚
⋮	⋮	⋮	⋮	⋮	⋮	⋮

③各ドナー候補の配偶関係を仮に①に補完したデータ(仮補完ユニット)から、 D_{ap} を算出

D_{fa} : 最小変化性 ③と①の距離	D_{ap} : 実在可能性 ③と②の距離
0.00	2.00
0.00	1.10
0.00	3.00
0.00	1.15
0.00	0.10
⋮	⋮

D_{fa} と D_{ap} の加重平均 D_{fpa} が最小となる10ユニットが NM CIAに入る

②実在するデータ(通過ユニット)

※要補完ユニットとのレコードの位置が近接する500ユニット

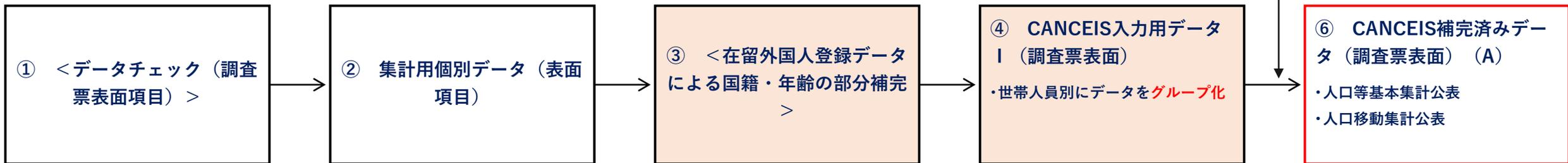
※) 年齢の距離は、両者の年齢差が1~2歳程度であれば、ほぼゼロとなり、年齢差が6歳以上であれば1となる関数

- CANCEISは、国勢調査における既存のデータチェックと演繹的補完の拡張として、ドナー補完を体系的に導入するためのプログラムとして位置付けることが可能
- 従前のデータチェックと整合したエディットルールをCANCEISのモジュール（⑤、⑪）に組み込むことにより、既存の集計プロセスと整合した集計プロセスの一環としてCANCEISを実装することが可能

集計の流れ（概略）

※下図は、令和7年国勢調査にCANCEISを実装することを想定した集計プロセスの概念図

◆調査票表面項目



◆調査票裏面項目



※ 国勢調査は、集計区分別に集計・公表時期が異なるため、

- 「人口等基本集計」は、調査票表面項目②を使用しCANCEISで補完した結果を公表（裏面項目は不使用）
- 「就業状態等基本集計」は、表面の補完結果⑥と調査票裏面項目⑧を統合したデータを用いCANCEISで補完

（ 学齢6歳未満の労働力状態を演繹的に補完
緯度・経度情報の付与 ）

〔 ⑪ 都道府県別モジュール 〕

国籍データの事前補完

- 国籍データは、人口構成比（日本人の比率と外国人の比率）が不均衡であり、事前処理なしにCANCEISを実行すると、外国人の過少補完となるため、CANCEIS実行前の事前補完が必須
- 事前処理は、令和2年調査の不詳補完値作成方法に準じつつ、CANCEISで代替可能な処理（A. 二人以上の世帯及びC. 単身世帯のうち民営賃貸共同住宅に居住している年齢不詳の者に係る処理）は除外

令和2年不詳補完値における部分補完の方法

A 二人以上の世帯

小地域別、男女・世帯人員の構成別、住宅の建て方別に、**基本項目不詳世帯以外の世帯をドナーとしたホットデッキ法**により、世帯員の年齢及び国籍の不詳を補完

B 単身世帯で国籍不詳の者

小地域別、男女別に、在留外国人登録データを活用した**コールドデッキ法**により国籍及び年齢の不詳を補完

C 単身世帯のうち、民営賃貸共同住宅に居住している年齢不詳の者

市（区）町村別、男女別に年齢を確率的に補完

CANCEISにおける補完処理で代替可能のため、R2年のような事前の部分補完は不要

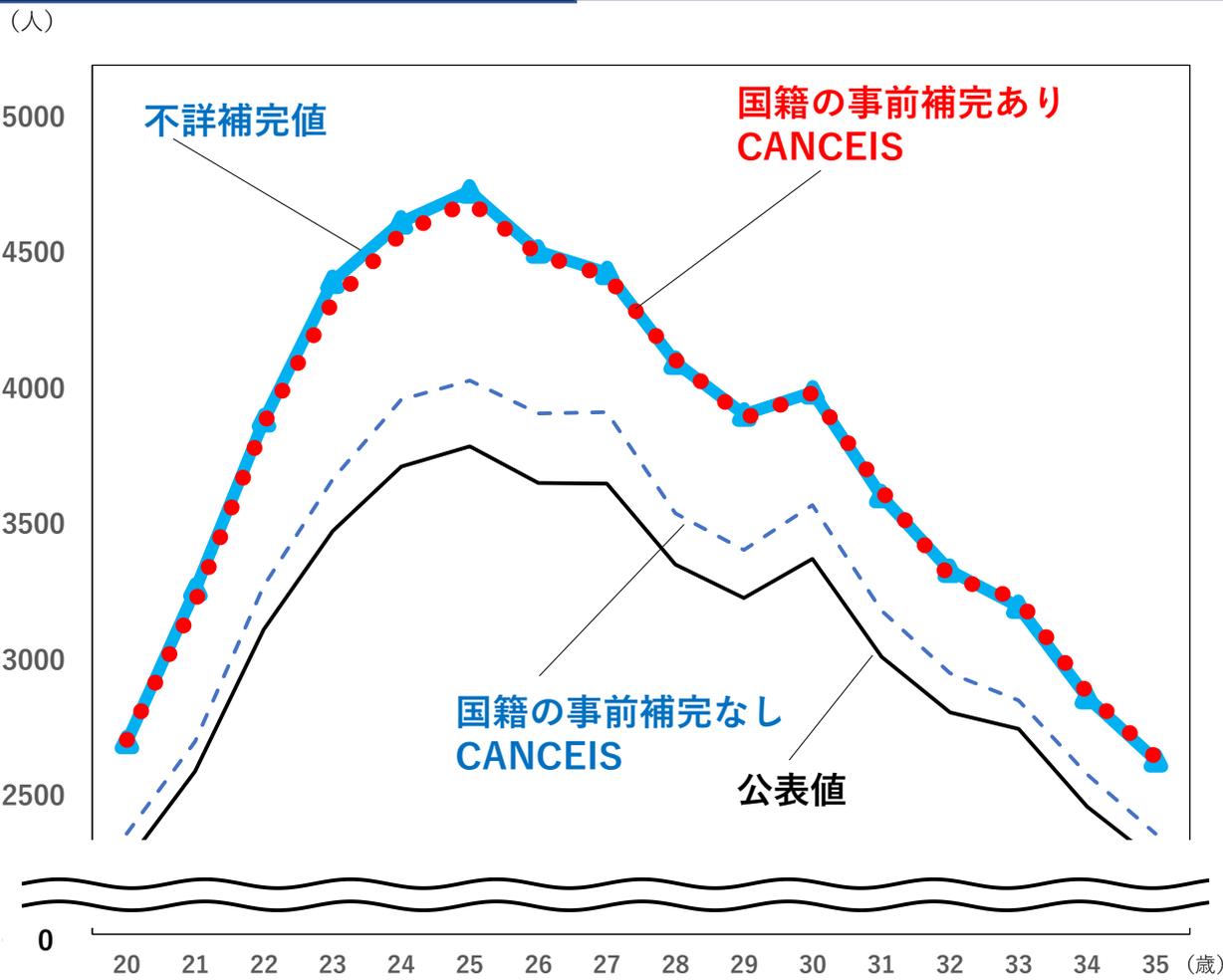
令和7年国勢調査の集計におけるCANCEIS適用の前処理として継承

CANCEISにおける補完処理で代替可能のため、R2年のような事前の部分補完は不要

国籍の事前補完を行ったCANCEIS試算結果

○ 下図は、令和2年国勢調査の人口等基本集計におけるCANCEISの適用結果を国籍の事前補完の有無別に示した
もの。国籍の事前補完ありのCANCEIS適用結果は、国籍の事前補完なしのCANCEIS適用結果よりも不詳補完値
に近い結果となった。

年齢別人口(愛知県)・外国人男性



年齢別人口(愛知県)・外国人女性

