

死因統計の精度及び効率性の向上に資する機械学習の検討に関する研究

研究代表者 今井 健 (東京大学大学院医学系研究科 准教授)

研究要旨

人口動態調査は国勢調査と並ぶ国の主要統計で公衆衛生施策の中心的資料である。本研究は原死因確定に関する調査を行い、我が国での原死因データ収集における課題を抽出し、ICD-11における死亡診断書や死亡統計ルールの変遷を調査すると共に、原死因確定作業に対する機械学習の適用可能性について調査・検討を行うことを目的とした。まず、死亡票の実データを対象に、文字列処理と自動 ICD-10 コード付与を行った上で、オートコーディングツール IRIS を適用し、約 80%の死亡票に対し、仮原死因を確定した。またこの付帯情報による仮原死因の変更の有無、外因や母側病態のコード追加の有無の割合について明らかにすると共に、機械学習による支援ターゲットとして「何らかの付帯情報による仮原死因の変更の有無」が有効であることを明らかにした。また機械学習により I 欄 II 欄病名と付帯情報からこの変更の有無を予測し、確信度と共に導出する2値分類器を複数のアルゴリズムで開発し、何らかの付帯情報が存在する死亡票を対象に **Accuracy 95%, ROC-AUC 0.953, PR-AUC 0.857** という非常に高い分類精度を達成した。機械学習の性質上、100%の精度実現は困難であるが、この分類器は付帯情報の影響による原死因コードの変更の有無のみならず、確信度も出力できることから、従来の人手による確認作業の正確性・効率性向上に大いに寄与する支援ツールとなると期待される。

将来的な ICD-11 導入後の原死因確定ツールとしては、Iris、WHO cause of death identification tool、国内のオートコーディングツールの更新の3種が考えられるが、本支援手法はこれらのどのオートコーディングツールとも組み合わせて利用することが可能であり、我が国における ICD-11 適用後においても死因統計の精度・効率性向上に資する極めて有力な手法と考えられる。

研究分担者

香川璃奈

筑波大学医学医療系 講師

明神大也

奈良県立医科大学公衆衛生学講座助教

研究協力者

大井川仁美

東京大学大学院医学系研究科客員研究員

大江和彦

東京大学大学院医学系研究科 教授

今村知明

奈良県立医科大学公衆衛生学講座 教授

A. 研究目的

我が国において人口動態調査は国勢調査と並ぶ国の基幹統計であり、中でも死因統計は最も重要な情報の一つである。今後 ICD-11 を国内適用するにあたっては原死因データを適切に収集・分析し、国際比較可能なデータを提供することが求められている。レセプトや現在普及が進む電子カルテでは標準病名の採用が進められているが、人

口動態調査の死因は自由入力病名が元となっており完全な自動集計は困難である。また我が国では高齢化が進み死亡者数の増加が見込まれることから、より正確で効率の高いデータ収集の方法の検討が求められている。

そこで、本研究は、原死因確定に関する調査を行い、我が国での原死因データ収集における課題を抽出し、ICD-11における死亡診断書や死亡統計規則の動向を調査すると共に、原死因確定作業に対する機械学習の適用可能性について調査・検討を行うことを目的とした。

B. 研究方法

B-1) 原死因確定プロセスにおける課題の抽出

原死因データ収集における課題については、既に本研究班メンバーらは平成30年度厚生労働統計協会調査研究委託事業において、ヒアリング調査・関連資料分析などを通じ基礎調査を行ってきた。本研究ではこれを発展させ、より詳細な分析を行うため、厚生労働省関係者へのヒアリングと共に、統計法第33条に基づく目的外利用申請によって死亡票・死亡個票データの提供を受け、実データを元にした分析を行った。

1年目は2015～2017年データ、2年目には2018年データ、さらに3年目は2019-2010年データの追加提供を受けて、毎年分析結果のアップデートを行い、最終的には6年分約800万件のデータに対する分析を行った。また同時に統計法22条に基づき、このうち1ヶ月分の抽出データに対する処理の内訳、並びに一部について厚生労働省にて人手チェックに回った300件のランダムサンプリング調査票情報の提供を受け、このサンプリング結果に対する分析と合わせることで、原死因確定プロセスの流れを明らかにした。

以前の基礎調査では、死亡票死亡個票において、病名以外の何らかの付帯情報(I欄II欄病名以外の何らかの補足的情報=手術解剖の有無と所見、外因死の追加事項、生後1年未満死の追加事項、その他付言すべきことがら等)が存在する場合、あるいはオートコーディングで疑義が有った

場合については職員により人手確認が行われており、その件数は月に4万件に及ぶことが判明していた。また、人手確認作業においては、厚生労働省内のオートコーディングシステムが決定した仮の原死因(各自由入力病名に対しICD-10コーディングを行い、その中からルールテーブルに基づき原死因を選択した結果)について、付帯情報を参照しながら必要に応じて原死因の変更を行うということも判明していた。つまり、何らかの付帯情報がある場合、機械学習の援用により、原死因の変更の有無を予め高精度に予測することができれば、人手確認処理を大幅に効率化することができる。しかし、その詳細や実際の件数の割合については不明で、どの程度の効用があるのか見積もることが困難であった。

本調査で原死因確定プロセスの流れを明らかにすることによって、機械学習により支援すべきポイントを明確になると共に、その効果についても論じることができるようになる。

B-2) 機械学習の適用可能性調査

オートコーディングシステムの出力に対し、人手の確認処理が行われている以上、この出力内容を知る必要がある。しかし、本研究では厚生労働省内部で使用されているオートコーディングシステムを用いることができない制限があり、自らこのシステムを模す必要があった。そのため様々な国で広く採用されているフリーの死亡票オートコーディングシステムであるIRISを用いることとした。

本研究で行った「原死因確定支援のための機械学習適用」は、以下の5つのステップからなる。全体構成の詳細については、「別添資料1:本研究で構築したシステムの詳細」を参照されたい。

STEP1) 死亡票・死亡個票からの IRIS 入力用データの作成

初年度ではまず IRIS がどの程度原死因確定に利用できるのか ICD-10 分類提要第2巻の事例、並びに死亡票実データを対象にフィージビリティスタディを行った。IRIS は日本語で利用することがで

きないが、I 欄 II 蘭病名を ICD-10 コーディングで
できれば、ICD-10 コードを入力することで原死因を
十分に決定することができることが判明した。この
詳細については、「別添資料2:インストラクション
マニュアル事例を対象とした IRIS による原死因確
定実験」、並びに「別添資料3実データを対象とし
た IRIS による原死因確定実験」を参照されたい。

IRIS をオートコーディングシステムの代わりに用
いることができるとなると、次は死亡票の自由入力
病名を ICD-10 コーディングする必要がある。

自由入力病名の記載は実に多様であるが、本
研究では多くの文字列処理と標準病名マスターを
用いて病名の ICD-10 コーディングを行い IRIS 処
理へかける方法を取ることにした。死亡票中の病
名で標準病名マスターとの完全一致で ICD コー
ディングできるものは 40%程度であったが、研究 1
~2 年目で多くの文字列処理を入れることにより、
最終的に「全ての病名が ICD コーディング可能な
死亡票」の割合は 80%まで増加した。この詳細に
ついては「別添資料4:死亡に関わる調査票情報
提供に基づいた ICD10 コード自動付与ツールの
作成(香川璃奈・令和2年度分担研究報告書)」を
参照されたい。

研究最終年度では、死亡票・死亡個票を突合し
た実データ(突合 DB: 平成 27~令和 2 年、約
800 万件)に対し、各種の前処理を行った上で、病
名に対し上記の自動 ICD-10 コーディング処理を
行った。

STEP2) IRIS での仮原死因確定処理

次に、死亡個票中の全病名に ICD-10 コードが
振られたものについて IRIS に入力し、仮原死因コ
ードを決定すると共に、確定原死因コードと比較を
行った。その際に、IRIS と国内のコーディングルー
ルの差異を吸収する処理や、期間の処理など各
種前処理も行った。例えば国内で用いられている
5 桁分類コードは IRIS が処理できないため、予め
桁を落として入力する必要がある。また、病名に付
随する期間も非常に多彩な表現バリエーションが
存在するが、これを適切に解析することで「急性」

と書いていなくても自動的に急性のコードを付与
する、といった処理を IRIS が行うことができる。

これらの詳細については、「別添資料 5:IRIS に
入力するための各種前処理」を参照されたい。こ
れらの前処理は全て網羅することが困難である
が、頻度が多い順に確認し、ルール化を行った。

STEP3) IRIS 処理結果の解析

STEP2) と逆に、IRIS 処理の結果確定された仮
原死因コードが国内では使わない/修正が必要な
コードであることもある。そこで、必要に応じて適切
な修正処理を行うとともに、この処理結果を Step1)
で死亡票・死亡個票を結合した「突合 DB」と合わ
せ、この後の解析用の「統合テーブル」を作成し
た。

STEP4) 機械学習用データセットの作成

次に、付帯情報が存在する死亡票のみを対象
に「IRIS が決定した仮原死因が付帯情報の影響
により変更されるか否か」を分類するモデルを学習
するため、統合テーブルの情報を元に、複数のアル
ゴリズムでの機械学習用データの作成を行っ
た。想定した機械学習アルゴリズムは3種類である
・**XGBoost-Simple** が最も単純なベースラインで、
(共通ベクトル) = 「性別・年齢、病名の ICD コー
ド、付帯情報の項目の有無」だけで予測を行う。
・**XGBoost-Embed** は共通ベクトルに、付帯情報
の内容を各手法でベクトル化したものを結合したも
のを使って予測を行う。
・**BERT** は、付帯情報の内容を BERT 言語モデル
に通した出力を、共通ベクトルと上位層で統合して
予測を行う DNN (Deep Neural Network) である。

ここではそれぞれのアルゴリズムに合わせた機械
学習用データセットの作成を行った。

(1) [XGBoost-Simple]

I 欄 II 蘭病名と付帯情報の有無(2 値) を用
いた勾配ブースティング回帰木(XGBoost)

- I 蘭 II 蘭病名については、使用された ICD10 コード、付帯情報については 22 種の項目に対する記載の有無(22 次元)、これに年齢・性別を加えたベクトルを用いた。
- 最も単純なものであるが、既に昨年度研究において 90.3%の正解率を実現しており、本研究ではその他の手法の比較におけるベースラインとなるものである。
- 後の比較では“BASELINE”として参照されている。

(2) [XGBoost-Embed]

(1) のベクトルに加え、付帯情報の文字列の内容を様々な手法で分散表現(ベクトル)に埋め込んだ(Embedding)結果を統合したベクトルに対する勾配ブースティング回帰木(XGBoost)

- (1)は付帯情報について項目の記載の「有無」だけであったが、内容を加味するため、各項目の記載文字列を分散表現に変換し、(1)のベクトルに加えたものである。
- 付帯情報に対する分散表現の獲得手法により、以下に細分化される。

➤ TF・IDF

文章中に出現した単語の重要度を TF (Term Frequency), IDF(Inverse Document Frequency)から算出する古典的な方法である。出現頻度 100 以上の単語だけを用いた。

➤ LSI

LSI (Latent Semantic Index) ではトピックという潜在変数を仮定する。各文書の BOW(Bag of Words)あるいは TF・IDF ベクトルを行とした、文書数×単語数の行列を特異値分解することで、文書数×トピック数に次元削減するものである。トピックを間に挟む分、TF・IDF よりも類義語多義語にある程度対応できるとされている。

➤ WORD2VEC

Word2Vec は文書中の「単語」の意味を分散表現として自動獲得するた

めのニューラルネットワークを用いた学習方法である。本研究ではある単語の前後の語から対象単語を復元するタスクにより学習する手法 (CBOW) を用い、得られた付帯情報文字列中の各単語の分散表現の平均を「付帯情報文字列全体」の分散表現とした。

➤ DOC2VEC (PV-DM)

➤ DOC2VEC (PV-DBOW)

Doc2Vec は Word2Vec の手法を文章に拡張した方法である。Word2Vec により得られた単語の分散表現を平均する方法とは異なり、文章自体の分散表現を直接得る手法である。

尚、DOC2VEC については本研究分担研究者の香川璃奈が担当した。

詳細は、別添資料6:「令和3年度分担研究報告書 死亡に関わる調査票情報提供に基づいた ICD10 コード自動付与ツールの作成(香川璃奈)」を参照されたい。

(3) [BERT]

BERT とは 2018 年に Google から発表された、大量の文書リソース(コーパス)から汎用言語モデルを自動獲得するための手法である。このモデルはさらにファインチューニングを行うことによって様々な自然言語処理タスクに汎用的に用いることができ、従来の自然言語処理タスクの多くにおいて SOTA (State of the Art) を達成したことで近年大きな注目を集めている。一般的な使い方としては、何らかのコーパスにて学習された BERT モデルを(必要に応じてさらに追加で事前学習を行い)、その上で特定のタスクを解くための Deep Neural Network アーキテクチャに組み込んでファインチューニングするという方法が用いられる。本研究でもこれら一般的な方法に則り、付帯情報の中の文字列情報を BERT モデルに通した

結果と (1)の情報を全結合層で統合した深層学習モデルを採用した。

- (2) の手法とは異なり、分散表現の獲得には日本語 Wikipedia を元にした BERT 言語モデルを採用し、(1)のベクトルと全結合層で統合したモデルである。

STEP5) 各手法での仮原死因変更有無予測モデルの学習

我が国の原死因確定プロセスにおいては、付帯情報がない場合は、基本的にオートコーディングツールで原死因コードがそのまま確定されるため、これらのモデルに通す必要がない。従って、**何らかの付帯情報が存在する死亡票のみを対象に STEP4 での各手法を用いて、『I 欄 II 欄病名と付帯情報から仮原死因の変更有無を予測し、確信度と共に導出する分類モデル』**を学習した。

また以上の一連の処理は自動化し、Docker 並びに仮想マシンにおいて実行可能なシステムとして実装した。

B-3) ICD-11 における死亡診断書や死亡統計ルールの動向調査

我が国の現行の死亡統計では ICD-10 を元にした WHO による原死因選択ルールが適用されている。しかし 2018 年 6 月に WHO が ICD-11 をリリースした今、ICD-11 における死亡統計の動向は今後の我が国への ICD-11 適用に際し重要である。本研究では初年度から3年度まで継続的に WHO 並びに日本 WHO-FIC 協力センターの関係者へのヒアリング、また WHO-FIC 会議 ITC (Informatics and Terminology Committee) などへの参加によってこの動向調査を行った。

倫理面への配慮

本研究では個人情報や動物愛護に関わる調査・実験は行わない。但し研究の遂行に当たっては、各種法令や「人を対象とする医学系研究に関する倫理指針」を含め各種倫理指針を遵守した。

C. 研究結果

C-1) 原死因確定プロセスにおける課題

まず本研究で関係者へのヒアリングを通じ、原死因確定プロセスにおける課題として、大きく次の 2 つが挙げられた。

(1) オートコーディングシステムで原死因がルールベースで決定できない事例

- 死亡票の I・II 欄傷病名が自由入力であるため、辞書マッチングで ICD-10 コードが付与できないことがある
- 原死因選択ルールに合致しない、疑義がある（「老衰」が年齢と合っていない、希少疾患である等）

(2) 何らかの付帯情報が存在する場合

- 付帯情報とは、I 欄 II 欄病名以外の何らかの補足的情報（手術解剖の有無と所見、外因死の追加事項、生後 1 年未満死の追加事項、その他付言すべきことがら等）のことである。
- 付帯情報が存在する場合は必ず人手での確認処理に回され、必要があればオートコーディングシステムが出力した仮の原死因コードを修正している。
- これには以下のような多様な事例が存在する。
 - ◇ I 欄(ア)に「肺炎」、手術欄に「胃悪性腫瘍切除術・1 週間前」とある。
→本来 I 欄(イ) に「胃癌」と書くべきとみなしてこれを選択する。
 - ◇ I 欄(ア)で「損傷」、手段・状況欄で、自殺、飛び降り、あるいは交通事故とわかるとそれを優先（外因等）。
 - ◇ 「細菌性肺炎」とあるが解剖欄の情報で菌の種類が分かる場合詳細化する。
 - ◇ 「飛び降り自殺」とあるが、I 欄内で産後うつによる影響であると分かることと妊産婦死亡のフラグを付与する。

これらの処理は複雑な事例もあり、非常に労力がかかるものである。一方で、人手確認に回ってきたものが全て何らかの変更を必要としている訳ではない。自動で判断できるもの／変更が無いものはなるべく事前に排除することで効率化が図れると考えられた。

次に、原死因確定プロセスの流れを明らかにするため、統計法33条に基づく目的外利用申請によって提供を受けた死亡票・死亡個票の実データを用い、両者を突合した。以下これを「**突合死亡票データベース (以下突合DB)**」と呼ぶ。突合方法の詳細については「**令和元年度総括研究報告書**」を参照されたい。

さらに、本研究では統計法22条に基づき、上記の一部について、厚生労働省にて人手確認に回った調査票情報の提供を受けた。これは死亡票の1ヶ月分の集計データに対し、人手確認に回った原因を集計したもの(表1) とさらにその中から300件をランダム抽出し、実際に人手確認でどのような対処が行われたかを集計したデータ(表2) である。またこれらの結果の分析に基づいて推計し、昨年までの結果をアップデートした原死因確定プロセスの流れ(最終版)を図1に示す。

			コーディング疑義	
			あり	なし
			16,631	96,308
付帯情報	あり	32.40%	12,973	23,592
	なし	67.60%	3,658	72,716

表1: 人手確認に回った死亡票の内訳

まず、表1の抽出データは全件で 112,939 件あり、「コーディング疑義」とは厚生労働省内のオートコーディングツール(病名に ICD-10 コードを付与し、複数の傷病名から原死因コードを確定するツール)にて、コード化プロセス中に疑義がある、とされたもの

である。また「付帯情報」とは死亡票死亡個票の I 欄 II 欄病名以外の何らかの補足的情報のことである。表中薄灰色の部分(コーディング疑義あり、もしくは付帯情報あり)の部分が**人手確認に回ったもので、合わせて 40,223 件 (35.6%) 存在した**。毎月約 4 万件程度が人手確認されていることになる。一方、残りの約 64% はオートコーディングシステムが決定した原死因コードがそのまま確定される。

次に人手で確認されたもの(Cリスト)からランダムに 300 件サンプリングされたデータに対する「対処内容の内訳」を表2に示す。

		原死因コード修正	
		あり	なし
追加コード (外因・母側 病態)の 付与	あり	11 (3.6%)	36 (12.0%) (うち母側:2)
	なし	30 (10.0%)	223 (74.3%)
計		41 (13.6%)	259 (76.3%)

表2: 人手確認対処内容の内訳

人手確認の後、原死因が変更されたものは 41 件 (13.6%)、変更なしが 259 件 (86.3%) であった。昨年度までは 100 件ランダム抽出したデータの提供を受け分析していたが、少数であるものの、300 件になったことで統計的信頼性が多少向上している。

以上の分析を元に推計した原死因確定プロセスの流れを図1に示す。図1中の「**実線四角**」は、実際のデータから算出されたもので確定値である。付帯情報があるものは全体の 32.3% となっている。

一方、「**破線四角**」は表1の結果からの推定値である。表1の結果から 35.6% が人手確認に回っているため、差し引き 3.2% がコーディング疑義により人手確認に回ったと推計される。

さらに、「**点線四角**」は表2の結果からの推定値である。表2の結果から、人手確認後 259:41 の割合、

- 表 A: 突合 DB の項目・記載内容の詳細
- 表 B: 突合 DB(ユニークキーのみ)の仕様とレコード件数
- 表 C: IRIS 処理結果と死亡票データを合わせた「統合テーブル」仕様
- 表 D: 各処理段階の件数詳細
- 表 E: BERT モデルを用いた各学習手法の予測精度一覧
- 表 F: 機械学習の結果一覧

STEP1) 死亡票・死亡個票からの IRIS 入力用データの作成

突合 DB とその詳細

平成 27～令和 2 年の 6 年間の死亡票・死亡個票の突合 DB の項目と記載内容の詳細について表 A に示す。突合 DB は 63 列からなるテーブルで、8,004,708 件存在した。各項目の詳細仕様は表 A に示す通りである。

突合 DB(ユニークキーのみ)

突合 DB は、「届出地、事件簿番号、処理年月」の組み合わせをキーとして行ったが、このキーが複数回存在するものが 92,768 件存在した。これを除いた結果、突合 DB(ユニークキーのみ)は 7,911,940 件となった。この複数存在するキーはヒアリングの結果、早期提出に起因するものとのことであった。この重複削除処理の結果について各年の件数と突合 DB に対する割合を本報告書末尾の表 D に示す。表 D の L10 を見ると、重複キーの削除によって、元の突合 DB の 98.84% の件数となっているが、微減でありほぼ影響はない。

備考欄の前処理

死亡票・死亡個票の各欄からは、入力時の文字数制限により、入り切らない文字列が備考欄に溢れて記入されることがある。この「備考欄に溢れた文字列」を多くの正規表現ルールにより、元の然るべき項目へ可能な限り復元して結合する処理を行った。

内容は、I 欄 II 欄病名とそれぞれの期間、解剖・手術の詳細、傷害が発生したところ・手段及び状況、その他付言すべき事柄、生後 1 年未満での病死に関する詳細、である。元の項目に復元できなかった文字列のみを「備考欄の文字列」として残した。正規表現処理の詳細は、本報告書全体末尾の「別添資料：備考欄前処理プログラムソース」を参照されたい。

またこの結果得られた「突合 DB(ユニークキー)」の各項目の件数を表 B に示す。

本研究で「付帯情報あり」と判断した基準は表 B に示す通りで、22 項目に対する条件で決定している。しかしながら、備考欄に溢れた文字列を復元することでこの条件に合致する件数は変化してしまう。表 B の件数中の「太字・背景薄灰色」の箇所は、この備考欄処理によって、溢れ文字列が元の項目へ復元された結果、件数に変化があった箇所を示しており、特に項番 41 は、備考欄の文字列が大幅に「本来あるべき項目」へ戻されていることを示している。

本処理の結果、最終的に「付帯情報あり」と判断された件数を表 B の最下段、及び表 D の L13 に示す。年によってバラツキはあるが全体として 32.3% であり、図 1 の通りである。

I 欄 II 欄病名に対する ICD コーディング

次に、死亡票の全ての I 欄 II 欄病名に対し、読点の削除、複数病名列挙の展開、文字の正規化、などの文字列処理を施した上で標準病名マスターを用いて ICD-10 コーディングを行った。要素技術は既に昨年度開発しており、本年度は 6 年間全てのデータに対しこれを適用した。

標準病名マスターを用いて、全てのI欄・II欄病名に対しほぼ原記載のままICD-10コーディング可能だったのは死亡票の約44%であるが、上記の多様な前処理を十分に行うことで、全ての病名がICD-10コード化された死亡票は約80%にまで増加した。表DのL17, L18にその件数と割合の詳細を示す。この80%の死亡票がIRISでの仮原死因決定処理へ入力可能なものとなった。

STEP2) IRIS での仮原死因確定処理

IRISは他のSTEPと異なりWindows上で動作するため、解析サーバーのVirtualBox上の仮想マシン(Windows10)で処理を行い、他のSTEPとは共有データ領域を介してデータやりとりを行った。詳細は「【別添資料1】本研究で構築したシステムの詳細」を参照されたい。

IRISによる仮原死因確定処理は、多くのデータに対して一度に行くと極端に遅くなるため、高速化のため各年のデータを50分割して処理を行った。この処理により1年のデータにつき約1日で処理可能であった。また、IRISはICD-10の国際版に準拠しており、日本国内で適用されている独自コード(詳細5桁目分類など)は実装されていない。このような事例についてはIRISが出力する仮原死因コード、死亡票における確定原死因コードのそれぞれについて修正処理を行った。この処理は昨年度と同様のため割愛する。(令和2年度統括研究報告書を参照のこと)

IRIS処理の結果、仮原死因が決定できた件数を表DのL19に、また割合をL20に示す。突合DB全年で99.26%の死亡票に対し、仮原死因が決定可能であった。決定できなかったものはIRISがRejectコードを出力するが(0.74%)、その原因の内訳は、表DのL21~L30に示す通りである。

IRISはRejectコードを出した場合も原死因コードを出力することがある(L31に件数示す)。これはRejected(Maybe)などのケースであるが、非常に稀なことから、後の機械学習での悪影響を考慮し、以降の処理では、Rejected無しのもの(L19)を対象とした。

今回使用した死亡票は、年によって準拠するICD-10コードが異なっており、2015,16年は、ICD-10 2003年版、2017~2020年はICD-10 2013年版となっている。一方、本研究の自動ICD-10コーディングに用いた標準病名マスターはV5.04(2013年版準拠)である。表Dを見ると、この影響により、2015,16年に対するコーディング結果が2017~2020年に比べ若干悪い結果となっていることが判明した。このような項目について薄灰色で示している。また、L38やL48を見ると、仮原死因の変更割合が、2017~2020と比べ明らかに高くなっており、後の機械学習への悪影響が考えられた。そこで、以降の処理では使用するICD-10のバージョンを揃え、**2017~2020年のデータのみにて実験を行うこととした。**

STEP3) IRIS 処理結果の解析

IRIS処理により確定された仮原死因と、国内での確定原死因を比較するため、IRISの出力情報と元の突合DBの情報を合わせ、この後の解析用の「統合テーブル」を作成した。仕様を表Cに示す。統合テーブルは性別、生年、没年、に加え、IRISの処理結果、死亡票の確定原死因、各付帯情報の項目の記載の有無などが全て含まれている。

統合テーブルのコード修正処理

表Cの項番16は、仮原死因と確定原死因の一致・不一致であり、機械学習の際には正解データとして用いられるものである。しかしながら、国内の原死因の方がIRISよりも粒度が細かい(桁が多い)、あるいは逆にIRIS側の方が国内よりも粒度が細かい、また国内では原死因として採用しない

コードが仮原死因として採用されている、などの理由により、本来「一致」として良いものが「不一致」となっているケースがあった。この対処のため、粒度が細かい方のコードは桁を落とす、コードを修正する、あるいは対象レコードから外す、など「**IRISの仮原死因コードと国内原死因コードを合わせる修正処理**」を、頻度が多いパターンを対象に可能な限り行った。このルールは昨年度と同様であるため、詳細は割愛する。

一部対象レコードから外したのものもあるが、表DのL34の通り、元データの99.99%が残っており、大きな影響はない。

この結果、再度「一致・不一致」を計算し、仮原死因の変更有無をまとめたものが、表DのL35～L54である。

IRISが仮原死因を決定したもの(L19)に修正処理を加え、対象とされたレコード数がL33である。この中で付帯情報があったものは2017～2020年データで29.2%存在した。全体では「付帯情報あり」のものは32.3%であったが、対象を「全病名にICD-10コードが付与できた80%」に絞っているため、この差が生まれたと考えられる。例えば、ICD-10コーディング不可能な複雑なケースは付帯情報が含まれている可能性が高い、という理由が考えられる。

一方、付帯情報がある29.2%の中で、**仮原死因の変更があったものは12.9%、無しは87.1%**であった。C-1で述べた「300件抽出データに対する分析結果（原死因変更:13.6%）とも大きく矛盾しない。これら2017～2020年データでの仮原死因変更割合を表Dから抜粋して表3に記す。

仮原死因	変更あり	変更なし	計
付帯あり (29.2%)	162,654 (12.9%)	1,100,078 (87.1%)	1,262,732
付帯なし (70.8%)	164,008 (5.4%)	2,892,357 (94.6%)	3,056,365

表3:2017～2020年データの仮原死因変更の割合

また、原死因の変更以外にも、人手によるチェックで修正が行われるものがある。これが損傷や外因の影響(ICD-10第19章)に対する「外因コード」(V01-Y98)の追加、周産期における母側病態コードの追加などの「コード追加」である。つまり原死因には変更がなくてもこのような補足コードが追加されることがあり得る。「付帯情報あり」についてこれらを細分化したものが以下の表4である。C-1)節での300件抽出データに対する同様の分析「**表2:人手確認対処内容の内訳**」とも似たような分布となっている。

国内のオートコーディングシステムとは異なり、IRISによるシミュレーションではあるが、大規模実データを対象とした解析においても、やはり「付帯情報あり」のうち、約8割のものは「原死因の変更もコードの追加も行う必要がない」ことが確認された。前述の通りコード追加が必要と考えられるケースは現状のオートコーディングシステムでも自動抽出できていることから、後は仮原死因の変更の有無が自動判別できれば、人手対処不要な約8割のデータを自動排除することが可能となる。

仮原死因	コード追加 必要	コード追加 不要
変更あり	29,279 (2.3%)	133,375 (10.6%)
変更なし	112,185 (8.9%)	987,893 (78.2%)

表4:付帯情報有りのものに対するコード変更・追加の割合

STEP4) 機械学習用データセットの作成

次に、「仮原死因の変更あり/無し」を教師データとし、「I欄II欄病名のICD-10コードと付帯情報」から仮原死因変更の有無を予測する分類器を開発した。

研究2年目に、50万件の死亡票サンプリングデータから全病名ICD-10コードが付与された32万件、さらにこの中で付帯情報が有る8万件のデータに対して、最もシンプルな[XGBoost-Simple]を適用し、Accuracy 90.3%を得ている。この結果の詳細は「別添資料8 原死因変更有無予測結

果」を参照されたい。この [XGBoost-Simple] は 90%とそこそこの正解率を達成しているが、付帯情報の各項目について記載の有無しか用いておらず、内容を考慮していない。この「付帯情報の内容」をいかに学習に組み込むか、が最終年度の大きなテーマであった。

BERT を用いた予測実験

最終年度での機械学習手法で最も期待されたのが [BERT] である。近年様々な自然言語処理タスクで従来の精度を塗り替えており、本研究でも大幅な精度向上が見込まれたため、まず BERT を用いた予測実験を行った。

比較のために2年目の[XGBoost-Simple] と同じ 8 万件のデータを用い、年齢、性別、使用された ICD-10 コード (1553 列)、各付帯情報の有無を (22 列) からなる 1577 次元の「共通ベクトル」をそのまま用いた。BERT を用いた本研究のアーキテクチャは、「共通ベクトル(1577 次元)」と「BERT モデルに付帯情報の文字列を通した結果得られる意味情報の分散表現(768 次元)」を上位の全結合層で統合し、最終的に2値分類の結果を得る」というもので、他のタスクでも良く用いられる手法である。

これに対し学習率の減衰方法や、インバランスデータへの対応方法など8種類の細分化を行い、学習を行った。詳細は「別添資料 2: BERT を用いた予測モデルの学習実験」を参照されたい。また、8種類の細分化手法の結果を「表 E: BERT モデルを用いた各学習手法の予測精度一覧」に示す。

結果として事前の期待に反し、最も良かった RedLR でもベースラインである昨年度の Accuracy から大きな向上は見られず、他の細分化手法はベースラインを下回る結果となった。このことから、以降の実験では、BERT 以外の方法 ([XGBoost-Embed]) に絞って本実験を行うこととした

XGBoost 系手法の機械学習用データセット作成

以降では、BERT 以外の方法に絞って実験を行うため、2年目と同様に [XGBoost-Simple]、また3年目に加えた [XGBoost-Embed] の5手法 (TF-IDF、LSI、WORD2VEC、DOC2VEC (PV-DM)、DOC2VEC (PV-DBOW)) についてそれぞれ付帯情報の意味ベクトルを作成/学習し、XGBoost で2値分類するための機械学習用データセットを作成した。詳細を「表 F: 各種機械学習手法の結果一覧」に記す。

全手法に用いた「共通ベクトル」は年齢、性別、出現した病名に対する ICD10 コード、各付帯情報の項目の記載の有無、からなるもので、これに加えた付帯情報の「意味ベクトル」が各手法で異なっている。

[BASELINE] (XGBoost-Simple) では意味ベクトルが存在せず、0 次元である。[TF-IDF] では出現頻度 100 以上(ストップワード除く) の単語を使っているため、2938 次元、[LSI] は 195 次元、残りは 200 次元ベクトルで表現されている。

共通ベクトルと付帯情報の意味ベクトルを結合したものが学習用ベクトルとなり、次元数は表 F に示した通りである。

STEP5) 各手法での仮原死因変更有無予測モデルの学習

最後に、STEP4 で作成されたデータを用いて各手法での仮原死因変更有無予測モデルの学習を行った。この予測モデルは共通して XGBoost を用いている。ハイパーパラメータは事前に探索を行い決定した。パラメータ詳細は本報告書全体の「別添資料: 分類器学習プログラムソース」(fit_and_predict_xgboost.py) を参照されたい。またトレーニングセットとテストセットは 4:1 に分割した。

各手法での結果を「表 F: 各種機械学習手法の結果一覧」に記す。

最も精度が高かったものは [XGBoost-Embed] の TF-IDF であり、Best Accuracy (正解率) 0.949 を実現した。BASELINE (XGBoost-Simple) は 0.915~0.920 であったことから、大きく精度が向上

している。また、例えば 2020 年データで **ROC-AUC で 0.953、PR-AUC で 0.857 と非常に高い精度での分類が可能**であった。例として 2020 年データの「変更あり」のものを対象とした時の ROC Curve を図2に、Precision(適合率)-Recall(再現率) Curve を図3に示す。

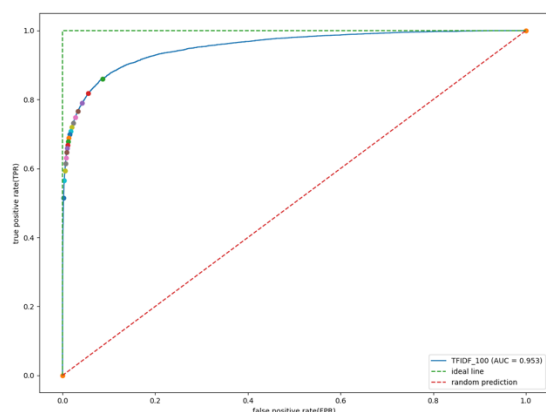


図2 ROC Curve (TF-IDF, AUC 0.953)

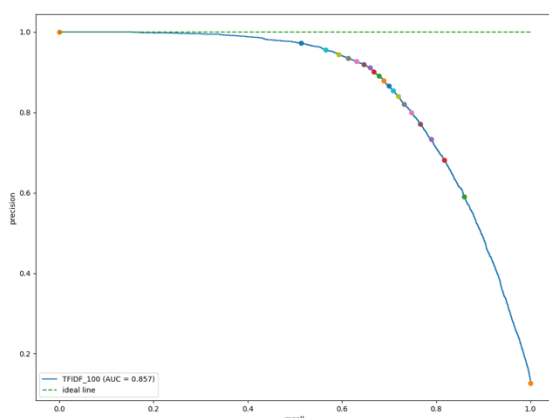


図3 Precision-Recall Curve (TF-IDF, AUC 0.857)

原死因の変更があるものが 12.9% しかないという非常に不均衡 (インバランス) なデータであり、AUC-ROC は高い値が出やすいが、「変更あり」に注目した時の Precision-Recall 曲線の Area under curve (PR-AUC) も非常に高い値であったことは本手法の有効性を示す大きな特徴である。

2017~2020 年はいずれも ICD-10 2013 年版で統一されており、年による大きな差はない。BASELINE と比較すると、[XGBoost-Embed] の

各手法はいずれも大きな精度向上を果たしており、DOC2VEC (PV-DM) が若干悪く、それ以外の LSI, WORD2VEC, DOC2VEC (PV-DBOW) も TF-IDF より僅かに劣るものの、いずれも大差なく同程度に高い精度を実現していた。

上記の手法に共通することとして、システムの最終的な出力は「**変更あり (1)」「なし(0)の2値ではなく、その範囲の動的な数値**であるため、現場で使用する目的に応じて閾値を設定することが可能であると同時に、この数値を適切に変換することで**出力に対する確信度**と解釈することも可能である。

また、閾値の設定によって変化するが、実際の現場での導入に際しては、「変更あり」のものをなるべく高精度にスクリーニングしたいことから、**Recall を重視した設定 (閾値) が適している**と考えられる。

C-3) ICD-11 における死亡診断書や死亡統計ルールの動向

ICD-11 における死亡診断書の動向

2021 年 10 月に開催された WHO-FIC 会議において、WHO が電子的な標準死亡診断書形式の普及を考えているという報告がなされた。WHO の考えでは、各国内で既に同様の電子的な報告システムが導入されている場合、国内のフォーマットを変える必要はないが、この標準形式で出力できるようにすることで、各種ソフトウェアからの可用性の向上を図りたいということであった。各国に対するアンケート結果によると、いくつか各国独自の死亡診断書項目が存在している場合があるが、大まかに WHO 標準形式とそぐっているということであった。今後電子化が十分に浸透すること、また標準形式に含まれる項目について各国のコンセンサスが得られることが重要で、現状すぐにこの標準死亡診断書形式へ移行する訳ではないが、我が国でもこの動きを注視して国内での報告形式の見直しをする必要があると考えられた。

また、我が国では手書き入力された死亡診断書を元に報告の過程で電子化入力されるという手順を踏むが、フランスやドイツにおいては、既に Web あるい

はモバイルアプリにて電子的に入力するシステムが導入されており、その過程でスペルチェックや入力項目間の不整合のチェックがなされている。業務効率化の観点から我が国においても、このような「報告時点」からの電子化と簡易チェックを考えていくことが重要な方向性であると考えられた。

死亡統計ルールの動向

ICD-11における死亡統計の考え方とルールについては、Web 公開されている WHO ICD-11 Reference にて記載がなされた。現時点では詳細情報について主に第2章 2.17 Mortality Statistics から 2.23 Annex for Mortality Coding までに記載されている。根本となる考え方については、ICD-10 の時を踏襲しているものの、細かなコードや記述法については大幅にアップデートされており、これまでの原死因選択ルールベースについても大幅な更新が必要となっている。図4に ICD-11 にて加えられた全体のワークフロー図を示す。

また、2021年10月に開催された WHO-FIC 会議においては、ITC (Informatics and Terminology Committee)にて、WHO cause of death identification tool についての進捗説明が行われた。これは ICD-10-SMoL (Startup Mortality List) と ICD11 reference guide から得られたルールを用いた、Iris とは異なる WHO 製の「原死因確定ツール」である。SMOL 自体は ICD-10 の時代から存在しており、ICD-10 の 3,4桁コードを使ったコーディングができない国を対象とした死因統計報告の第一段階として位置付けられる、粒度の粗いコードセットとなっている。本ツールは ICD-10 の粒度、SMOL の粒度、ICD-11 の粒度それぞれで原死因選択ができる Web ベースのツール提供を目指しているということであった。現時点では Reference Guide から抽出されたルールについては一部しか実装できておらず、CDC death certificate を用いたテストにおいては、SMOL ルールを用いた時で 73.5%の精度、Reference Guide からのルールも

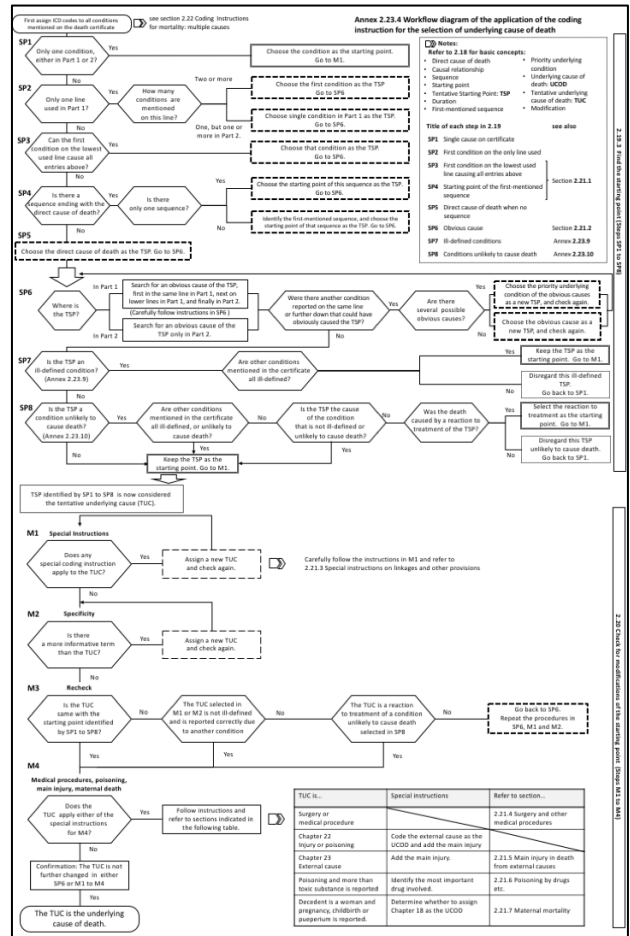


図 4. ICD-11 原死因決定ワークフロー

用いて 84.9%の精度、また ICD-11 でコードされた死亡票について 79.2%の精度と報告があった。但し病名以外の付帯情報(各所見・外因の手段状況・備考欄等)は一切考慮していない。

WHO cause of death identification tool はその元となった粗い粒度である SMOL から発展し今や ICD-10, ICD-11 の粒度での原死因選択ルール全体が実装されようとしている。一方、Iris はフル ICD-10 を用いた原死因選択ツールであり、ICD-11 への移行が進められているがまだ完成には至っていない。WHO によると、WHO cause of death identification tool はツール間の架け橋になり得るもので、本ツールで実装されるルールは他のツールでも用いることができる、ということであった。このことから、WHO の当該ツールでは、単に SMOL の粒度だけではなく、ICD-10,11 の粒度でも原死因選択が可能、Iris と競合

するツールに発展させていく意向があると考えられた。但し、WHO の当該ツールも Iris(ICD-11 対応版)も、現時点ではいずれも ICD-11 における原死因選択ツールを完全に網羅している訳ではなく開発途中である。また病名以外の付帯情報を考慮している訳ではないことには注意が必要である。

D. 考察

本研究の成果で、死亡票実データの約 80%を対象として、IRIS による仮原死因確定処理を行い、確定原死因と比較することで、原死因コードの変更の割合、また外因や母側病態コードの追加割合が明らかになった。

何らかの対処が必要なもの(原死因コードの変更、外因や母側病態コードの追加)は両者を合わせて 2 割であり、最初に 1:4 の分類タスクにより「対処の必要がない 8 割」を除去した後に、細かな対処内容の分類を行う 2 段階処理が適していると考えられた。最初の分類タスクが高精度に行えるだけでまず 8 割の人手処理を削減できることが期待でき、また 2 段階目の処理によって残りの 2 割に対する人手作業の効率化が図れると期待できる。

また、これを開発ターゲットに定め、機械学習により I 欄 II 欄病名と付帯情報からこの変更の有無を予測し、確信度と共に導出する 2 値分類器を複数のアルゴリズムで開発した。何らかの付帯情報が存在する死亡票を対象として、Accuracy 95%, ROC-AUC 0.953, PR-AUC 0.857 という非常に高い分類精度を実現し、本研究の手法の有効性が示された。

一方、表3の結果より「付帯情報がない」にも関わらず「仮原死因に変更があった」ものが 5.4%存在している。これは本来存在しないはずのものであり、ICD-10 コード化する際の手法(本研究では標準病名マスター、厚生労働省内では独自の辞書)、あるいは原死因を決定する際のルールテーブルの違い(IRIS /

国内のオートコーディングツール)の影響と考えられ、本研究の枠組みでは除去できない限界である。

本研究の手法での Accuracy 95% という値は、付帯情報の有無に関係なく、オートコーディングツールの仕組みの違いが持っている本質的な“ずれ”がそもそも 5.4% ある、という事実が大きく影響されていると考えられ、達成しうる上限値なのかもしれない。

本研究では使用することができなかったが、厚生労働省内での本格的な導入時には、国内のオートコーディングツールを元に仮の原死因を決定し、これを教師データとして本研究の手法で学習をやり直すことでさらに精度が向上すると考えられる。

機械学習の性質上、100%の精度実現は困難であるが、この本研究の手法は付帯情報の影響による原死因コードの変更の有無のみならず、確信度も出力することが可能である。これも参考情報としてユーザーへ提示することで従来の人手による確認作業の正確性・効率性向上に大いに寄与する支援ツールとなると期待される。

BERT が事前の予想に比してあまり効果を発揮しなかった理由としては、付帯情報の文字列は総じて長いものが少なく、BERT モデルへの入力に適していないということが考えられる。寧ろ TF・IDF のような古典的手法が奏功していたことを考えると、高度な言語モデルで「文脈を捉える」処理よりも、「特定の重要な単語の出現をきちんと捉えられているか」という点の方が効いているのではないかと考えられた。また、[XGBoost-Embed] の TF・IDF、LSI、WORD2VEC、DOC2VEC(PB-DBOW)の4手法の間にはそこまで大きな差はないが、唯一 TF・IDF が異なっているのはベクトルの次元数が他と比べて多いことである。2段階目で XGBoost を用いることを考えると入力となるベクトルの次元数はある程度大きい方が良い、という可能性も考えられる。

ICD-11における死亡診断書の動向からは、我が国においても入力時点からの電子化を推進することで作業効率化とスペルチェックや入力項目間の不整合のチェックなどによる質向上を図る方向性が今後重要と考えられた。項目についても我が国独自の項目(あるいは選択肢の粒度)とWHO標準形式との整合性を確保していく観点が重要と思われる。多少本研究の本筋から外れるが、そもそも入力欄に文字数限定があり、溢れた文字が備考欄に記載される、という方式自体が前時代的に感じられる。病名や期間、手術、解剖の詳細などの情報が途切れてしまい、AIなどを用いた自動解析に極めて適さない仕組みである。海外で、死亡診断書の入力を電子デバイスから行っているシステムで同様の仕様は見当たらなかった。本研究では「備考欄に溢れた文字を元に復元する前処理」を適用したが、その開発に膨大な時間がかかったことを考えると、今後この点を改善することも重要であろう。

さらに、このように一層進んだ電子化方式にて収集された死亡票データに対し原死因を確定するプロセスについては、効率化の観点から”Iris ICD-11 対応版”もしくは”WHO cause of death identification tool”を採用することも有力な選択肢であることには変わりない。しかし、両者とも開発途中であり今後動向を注視することが必要である。WHOによると当該ツールに導入される「原死因選択ルールセット」は、他のツールでも用いることができるよう共有するということであったため、**将来仮にどちらを用いることとなっても実用上の問題は生じない**と思われた。

これらの2つはこれまで国内(厚生労働省内)で導入されてきたオートコーディングツールをICD-11対応とするときの代替候補ということであり、もちろん我が国独自のオートコーディングツールをICD-11対応へ更新するという方向性も考えられる。しかし、以上3つのどの方式を取るとしてもそのツール単体では「付

帯情報を考慮するような高度な原死因確定処理」は実現できておらず、現在病名以外の各種付帯情報の存在等によって人手による確認が行われている(約35%の死亡票)ステップについては何らかの支援が必要である。本研究のAI技術による支援手法は、選択ルールセットを実装した何らかのオートコーディングツールにより仮の原死因を決定した後、「付帯情報を考慮すると仮原死因コードに変更が起こるか否か」をその確信度と共に高精度に提示するものであるため、**どのオートコーディングツールにも組み込むことが可能**である。そのため本研究のような**AIによる支援手法をルールベースのオートコーディングツールと組み合わせることはICD-11においても十分に可能**であり(ICD-11コードにおいて再学習させるプロセスは必要であるが)、原死因確定プロセスの正確性・効率性向上のために極めて重要であると考えられる。本研究での成果は今後我が国における原死因決定プロセスへのAI活用に向けて大いに寄与する貴重な知見であると共に、本成果を元にした現場へのAI支援システム導入が望まれる。

E. 結論

本研究では、死亡票の実データに対する仮原死因確定コーディングツールとしてIRISを使用し、約80%の死亡票に対し仮原死因を決定した上で、これを確定原死因と比較することで得られた教師データを元に、機械学習で「付帯情報の影響により原死因コードが変更されるか否か」を確信度と共に高精度に提示する手法を開発した。結果として**Accuracy 95.0%, ROC-AUC 0.953, PR-AUC 0.857**という非常に高い精度を得た。この支援手法を既存のオートコーディングツールの処理結果に組み合わせることで原死因決定プロセスの正確性・効率性向上に大いに寄与すると期待される。また、この手法は将来の我が国でのオートコーディングツ

ルの選択候補に依存せず、どれであっても組み合わせさせて使えることから、近い将来の ICD-11 適用後においても有力な手法と考えられる。

F. 健康危険情報

なし

G. 研究発表

1. 明神 大也, 大井川 仁美, 香川 璃奈, 今村 知明, 今井 健. 死因統計の精度と効率性の向上に向けた我が国の原死因確定課題の抽出. 医療情報学 40(Suppl.):674-676, 2020.
2. 大井川 仁美, 明神 大也, 香川 璃奈, 今村 知明, 今井 健. 原死因確定プロセスにおける IRIS の国内導入可能性に関する基礎的な検討. 医療情報学 40(Suppl.):677-682, 2020.
3. 今井 健, 明神大也, 大井川仁美, 香川璃奈, 今村知明. 原死因確定作業についての実態・問題点の把握、ならびに正確・効率性向上に向けた機械学習の適用可能性と課題に関する調査研究. 厚生 の 指 標 , 2020;67(3): 17-24, 2020.
4. 大井川仁美, 今井 健, 香川璃奈, 明神大也, 今村知明. 原死因決定プロセスの効率化に資する機械学習による原死因コード変更予測. 医療情報学 41(Suppl.):797-800, 2021.<第 41 回 医療情報学連合大会 研究奨励賞>

H. 知的財産権の出願・登録状況

なし

表A:「突合DB」の項目・記載内容の詳細

項番	項目名	由来	項目日本語名称	バイト数	意味	記載内容の仕様
0	id		仮名ID		Irisに入れる用の通し番号	
1	sub_area	死亡個票	届出地(8)	8	届け出が出された場所の番号	<1,2バイト目(都道府県)> 01-47 北海道～沖縄県 「都道府県別市区町村符号及び保健所符号一覧」参照のこと。昭和47年1月からは行政管理庁(現総務省)の定める「統計に用いる標準地域コード」を採用した。 <3,4バイト目(保健所)> 01-29 指定都市が設置した保健所 31-49 保健所を設置する市(指定都市を除く)が設置した保健所 51-98 道府県が設置した保健所 01-69 区部の保健所(東京都) 71-98 区部以外の保健所(東京都) <5バイト目(支所符号)> A-Z 受け付けた市区町村支所の符号 A支所～Z支所 △ 支所なし <6バイト目(市区町村:種類)> 1 指定都市、特別区 2 市(指定都市を除く) 3-7 町村 <7,8バイト目(市区町村:順位)> 01-99 都道府県における順位
2	jiken_nm		事件簿番号(4)	4	事件簿番号	0001-9999市区町村(支所を含む)における事件簿番号を表す。
3	shori_ym		処理年月(6)	6	処理年月	ex. 201501-202012 2015年1月処理～2020年12月処理
4	bikou_yn		備考欄記述有無(1)	1	備考欄に記述があるかないか	0 記述なし 1 記述あり
5	death_kbn		「死亡した人の住所」のうち届出地区分(1)	1	届出地の区分	1 届出地と同じ 2 届出地以外 3 外国 V 不詳 △△ 未記入または届出地不詳
6	d_tdfk		都道府県(8)	8	死亡したところの都道府県名	漢字
7	d_city		市区町村(12)	12	死亡したところの市、郡、東京都の区	漢字
8	d_town		市区町村(18)	18	死亡したところの町、村、指定都市の区	漢字
9	d_add		コード(7)	7	死亡したところのコード	<1,2バイト目(都道府県)> 01-47 北海道～沖縄県 「都道府県別市区町村符号及び保健所符号一覧」参照のこと。昭和47年1月からは行政管理庁(現総務省)の定める「統計に用いる標準地域コード」を採用した。 △△ 未記入または届出地不詳 <3バイト目(市区町村:種類)> 1 指定都市、特別区 2 市(指定都市を除く) 3-7 町村 △ 未記入または届出地不詳 <4,5バイト目(市区町村:順位)> 01-99 都道府県における順位 △△ 未記入または届出地不詳 <6,7バイト目(保健所)> 01-29 指定都市が設置した保健所 31-49 保健所を設置する市(指定都市を除く)が設置した保健所 51-98 道府県が設置した保健所 01-69 区部の保健所(東京都) 71-98 区部以外の保健所(東京都) △△ 未記入または届出地不詳
10	la_c		「死亡の原因」に含まれるA欄・ア欄の原因(38)	38	A欄記載の傷病名	漢字
11	la_p		期間(16)	16	期間	漢字
12	lb_c		「死亡の原因」に含まれるB欄のイ欄の原因(38)	38	B欄記載の傷病名	漢字
13	lb_p		期間(16)	16	期間	漢字
14	lc_c		「死亡の原因」に含まれるC欄のウ欄の原因(38)	38	C欄記載の傷病名	漢字
15	lc_p		期間(16)	16	期間	漢字
16	ld_c		「死亡の原因」に含まれるD欄のエ欄の原因(76)	76	D欄記載の傷病名	漢字
17	ld_p		期間(32)	32	期間	漢字
18	ll_c		「死亡の原因」に含まれるE欄の原因(76)	76	E欄記載の傷病名	漢字
19	ll_p		期間(32)	32	期間	漢字
20	ope_flg		手術フラグ(1)	1	手術の有無	1なし 2あり △未記入
21	ope_detail		手術の部位及び所見(116)	116	手術の詳細記述	漢字
22	ope_bk_flg		(手術)備考欄への記載(1)	1	手術に関して備考がある場合はフラグ	1記載あり △記載なし
23	ope_ymd		手術日(8)	8	手術年月日	<1～4バイト目> 0000-9999手術をした年(西暦) △△△△未記入もしくは手術無し <5,6バイト目> 01～12手術をした月 △△未記入もしくは手術無し <7,8バイト目> 01～31手術をした日 △△未記入もしくは手術無し

24	ap_flg	解剖フラグ(1)	1	解剖の有無	1解剖なし 2解剖あり △未記入
25	ap_detail	解剖の部位及び所見(116)	116	解剖の詳細記述	漢字
26	ap_bk_flg	(解剖)備考欄への記載(1)	1	解剖に関して備考がある場合はフラグ	1記載あり △記載なし
27	dc_type	死因の種類(2)	2	自然死・不慮の外因死・その他の不詳の死などの種類を表す数字	01 病死・自然死 02 交通 03 点等 04 溺水 05 火災 06 窒息 07 中毒 08 その他 09 自殺 10 他殺 11 不詳の外因死 12 不明の死 △△,00 未記入
28	inj_ymd	傷害が発生したとき(8)	8	障害発生年月日	<1~4バイト目> 0000-9999傷害が発生した年(西暦) △△△△未記入もしくは傷害ではない <5,6バイト目> 01-12傷害が発生した月 △△未記入もしくは傷害ではない <7,8バイト目> 01-31傷害が発生した日 △△未記入もしくは傷害ではない
29	inj_hm	傷害が発生したとき(5)	5	障害発生時分	<1,2バイト目> 00-11傷害が発生した時 △△未記入もしくは傷害ではない <3,4バイト目> 00-59傷害が発生した日 △△未記入もしくは傷害ではない
30	inj_p_type	傷害が発生したところの種別(1)	1	住居・工場・道路などの障害が発生した場所の種類を表す数字	1住居 2工場及び建設現場 3道路 4その他 △未記入もしくは傷害ではない
31	inj_detail	傷害が発生したところその他の記述(40)	40	傷害が発生したところの種別がない場合の記述(その他)	漢字
32	inj_tdfk	傷害発生場所(8)	8	傷害発生都道府県	漢字
33	inj_city	傷害発生場所(12)	12	傷害発生市郡	漢字
34	inj_town	傷害発生場所(18)	18	障害発生区町村	漢字
35	inj_method	手段及び状況(120)	120	障害発生時の手段及び状況に関する詳細記述	漢字
36	inj_biko	(傷害)備考欄への記載(1)	1	傷害発生に関して備考がある場合はフラグ	1記載あり △記載なし
37	inf_detail	「生後1年未満での病死」の病態・異状の詳細	84	妊娠・分娩時における母体の病態または以上に関する詳細記述	漢字
38	inf_biko_yn	備考欄への記載	1	「生後1年未満での病死」に関して備考がある場合はフラグ	1記載あり △未記入もしくは1歳未満ではない
39	sonota	その他付言すべき事柄	60	その他に付言すべきことに関する詳細記述	漢字
40	biko_yn	備考欄外字有無	1	備考欄に外字がある場合はフラグ	0備考欄あり、外字なし 1備考欄あり、外字あり △備考欄なし
41	biko_detail	備考欄	1024	備考に関する詳細記述	漢字
42	ym	調査年	2	調査年	西暦の下2桁
43	sub_ym	提出年月	4	死亡票提出年月日	元号2桁
44	jiken_sb	事件簿番号サブ	1	事件簿番号の補助番号	1-9 事件簿番号は原則として同じ番号はないはずであるが、万一同じ番号があったときは1~9をつけて区分する。 △ 通常は△
45	sex	性別	1	性別	1男 2女
46	birth_yn	出生年月日フラグ	1	出生年月日あるか	V不詳(出生年月日不詳。以下33~39カラムはVVVVVV) △出生年月日の記入があるもの
47	birth_ymd	出生年月日時分	7	出生年月日	<1バイト目> 1 明治 2 大正 3 昭和 4 平成 5 令和 V 不詳 <2,3バイト目> 01-64 元号に対して1年~終止年 VV 不詳 <4,5バイト目> 01-12 1月~12月 VV 不詳 <6,7バイト目> 01-31 月に対して1日~終止日 VV 不詳
48	birth_hm	出生年月日時分	4	出生時分	<1,2バイト目> 00-23午前0時~午後11時 △△生存期間が31日以上とき(調査票記入要領による) VV不詳 <3,4バイト目> 00-59分~59分 △△生存期間が31日以上とき(調査票記入要領による) VV不詳
49	death_yn	死亡年月日フラグ	1	死亡年月日あるか	V不詳(死亡年月日不詳。以下45~51カラムはVVVVVV) △死亡年月日の記入があるもの

50	death_ymd	死亡年月日時分	7	死亡年月日	<1バイト目> 3 昭和 4 平成 V 不詳 <2,3バイト目> 01-64 元号に対して1年～終止年 VV 不詳 <4,5バイト目> 01-12 1月～12月 VV 不詳 <6,7バイト目> 01-31 月に対して1日～終止日 VV 不詳
51	death_hm	死亡年月日時分	4	死亡時分	<1,2バイト目> 00-23午前0時～午後11時 VV不詳 <3,4バイト目> 00-590分～59分 VV不詳
52	ori_dc	原死因(5)	5	原死因コード	<1～3バイト目> A00-U04 外部参照。原死因の詳細については「疾病、傷害および死因統計分類提要 ICD-10(2003年版)準拠 第2巻内容例示表」及び「人口動態統計用として追加設定した基本分類細分項目一覧」(平成27年)を参照のこと <4バイト目> 0-9原死因4桁目設置有りの場合 △原死因4桁目設置無しの場合 <5バイト目> A-1原死因5桁目設置有りの場合 △原死因5桁目設置無しの場合
53	gaiin	外因符号	4	外因コード	<1～3バイト目> V01-Y89外部参照。原死因S00.0～T98.31についての外因符号。外因符号の詳細については「疾病、傷害および死因統計分類提要 ICD-10(2003年版)準拠 第2巻内容例示表」を参照のこと <4バイト目> △△△外因死以外(原死因A00.0～R99-U04.9)の場合 0-9外因符号(上3桁)がV01～Y89で4桁目設置有りの場合の発生場所等、外因の状況コード。○交通事故以外の不慮の事故(外因符号W00～X59)による死亡の場合、発生場所コード。0 家(庭)、1 居住施設、2 学校、施設及び公共の地域、3 スポーツ施設及び競技施設、4 街路及びハイウェイ、5 商業及びサービス施設、6 工業地域及び建築現場、7 農場、8 その他の明示された場所、9 詳細不明の場所。○それ以外の外因(外因符号V01～V98、X60～Y89)による死亡の場合、発生場所等、外因の状況コード。詳細については「疾病、傷害および死因統計分類提要 ICD-10(2003年版)準拠 第2巻内容例示表」を参照のこと △△外因死以外(原死因A00.0～R99-U04.9)、外因符号4桁目設置無しの場合
54	mom_sy	母側病態	5	母側病態コード	<1～3バイト目> P00-P99外部参照。母側病態の詳細については「疾病、傷害および死因統計分類提要 ICD-10準拠 第2巻内容例示表」及び「人口動態統計用として追加設定した基本分類細分項目一覧」(平成27年)を参照のこと △△△スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合) <4バイト目> 0-9母側病態4桁目設置有りの場合 △スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡、母側病態4桁目設置無しの場合) <5バイト目> A-1母側病態5桁目設置有りの場合 △スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡、母側病態5桁目設置無しの場合)
55	twin_yn	単多胎別	2	単胎か多胎か数字	<1バイト目> 1単胎 2-7双子～7つ子 88つ子以上 △スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合) <2バイト目> 1-71番目～7番目 88番目以上 V不詳 △スペース(単胎、1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合)
56	preg_yn	妊娠週数	1	妊娠の有無	0 妊娠週数記入のあるもの V 不詳 △ スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合)
57	preg_w	妊娠週数	2	妊娠週数	12-9812週～98週 VV不詳 △△スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合)
58	mon_b_ymd	母の生年月日	7	母の生年月日	<1バイト目> 3昭和 4平成 V不詳 △スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合) <2,3バイト目> 01-64元号に対して1年～終止年 VV不詳 △△スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合) <4,5バイト目> 01-121月～12月 VV不詳 △△スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合) <6,7バイト目> 01-31月に対して1日～終止日 VV不詳 △△スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合)
59	p_preg_bir	前回の妊娠	2	前回までの出生児数	00-180人～18人 1919人以上 VV不詳 △△スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合)
60	p_preg_dea	前回の妊娠	2	前回までの死産児数	00-190胎～19胎 2020胎以上 VV不詳 △△スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合)

61	child_1	子の数	2	今回の出生子含む出生子	01-191人～19人 2020人以上 △△スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合) VV不詳
62	child_2	子の数	2	出産児	01-391～39胎・人(出産児＝出生子(155～156カラム)＋死産児(妊娠満22週以後)) 4040胎・人以上(出産児＝出生子(155～156カラム)＋死産児(妊娠満22週以後)) △△スペース(1歳以上及び年齢不詳の死亡、0歳で病死以外の死亡の場合) VV不詳

表B:「突合DB(ユニークキーのみ)」の項目仕様とレコード件数

項番	項目名	意味	Bytes	付帯情報の有無 の判定基準	2015年		2016年		2017年		2018年		2019年		2020年	
					なし	処理後	なし	処理後	なし	処理後	なし	処理後	なし	処理後	なし	処理後
備考欄の処理(溢れ文字列の元項目への結合)																
0	id	Iris処理用の仮名ID	7		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
1	sub_area	届け出が出された場所の番号	8		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
2	jiken_nm	事件簿番号	4		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
3	shori_ym	処理年月	6		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
4	bikou_yn	備考欄に記述があるかないか	1		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
5	death_kbn	届出地の区分	1		1236015	1236015	1265300	1265300	1309518	1309518	1358079	1358079	1376665	1376665	1366162	1366162
6	id_idfk	死亡したところの都道府県名	8		1156709	1156709	1183673	1183673	1223805	1223805	1262804	1262804	1282413	1282413	1275576	1275576
7	d_city	死亡したところの市、郡、東京都の区	12		1158995	1158995	1185824	1185824	1226189	1226189	1265079	1265079	1284561	1284561	1278077	1278077
8	d_town	死亡したところの町、村、指定都市の区	18		1124819	1124819	1152259	1152259	1194588	1194588	1236935	1236935	1255436	1255436	1247524	1247524
9	d_add	死亡したところのコード	7		1232357	1232357	1261446	1261446	1305700	1305700	1354030	1354030	1372307	1372307	1363004	1363004
10	la_c	ア欄記載の傷病名	38		1233909	1233909	1263174	1263174	1307433	1307433	1355898	1355898	1374327	1374327	1364810	1364811
11	la_p	期間	16		1214814	1215632	1244394	1245571	1288530	1289757	1337818	1339021	1355927	1355927	1348402	1349610
12	lb_c	イ欄記載の傷病名	38		450934	450934	457226	457226	462464	462464	471565	471565	459618	459618	440092	440092
13	lb_p	期間	16		347244	350195	354441	356991	359010	361373	369950	372359	361549	364192	348866	351381
14	lc_c	ウ欄記載の傷病名	38		96485	96485	97574	97574	97091	97091	98826	98827	93878	93878	88240	88240
15	lc_p	期間	16		75853	76334	77120	77503	76986	77327	79260	79566	75721	76110	71628	71993
16	ld_c	エ欄記載の傷病名	76		18173	18173	18223	18223	18038	18038	18459	18459	17259	17259	16345	16345
17	ld_p	期間	32		13544	13642	13495	13599	13332	13409	13804	13886	13069	13151	12473	12563
18	ll_c	II欄記載の傷病名	76		419695	419695	425148	425149	435249	435249	445932	445932	443940	443940	436382	436382
19	ll_p	期間	32		388944	390494	395980	397590	406795	408483	419622	419622	419840	419840	413205	414629
20	ope_flg	手術の有無	1	2の場合「有」	1225149	1225149	1255342	1255342	1300736	1300736	1349277	1349277	1366293	1366293	1357354	1357354
21	ope_detail	手術の詳細記述	116	記載あれば「有」	182103	183651	183126	184509	185573	186793	188201	189459	186746	188132	183044	184462
22	ope_bk_flg	手術に関する備考がある場合は「有」	1	記載あれば「有」	2	2	1	1	2	2	3	3	3	3	3	3
23	ope_ymnd	手術年月日	8	記載あれば「有」	162250	170428	162787	171006	164631	172479	166800	176037	165743	174722	162755	172148
24	ap_flg	解剖の有無	1	2の場合「有」	1224106	1224106	1254352	1254352	1299954	1299954	1348571	1348571	1365596	1365596	1356691	1356691
25	ap_detail	解剖の詳細記述	116	記載あれば「有」	29336	29644	30341	30655	30135	30453	29610	29869	28774	29019	24894	25136
26	ap_bk_flg	解剖に関する備考がある場合は「有」	1	記載あれば「有」	0	0	1	1	0	0	4	4	7	7	7	7
27	dc_type	自然死・不慮の外因死・その他の不詳の死などの種類を表す数字	2	「空.00.01」以外の記載あれば「有」	1235172	1235172	1264496	1264496	1308575	1308575	1357129	1357129	1375599	1375599	1365430	1365430
28	inj_ymnd	傷害発生年月日	8	記載あれば「有」	49959	49959	47870	47870	47829	47829	48061	48061	45781	45781	45373	45373

29	inj_hm	傷害発生時分	5	記載あれば「有」	38950	38950	37433	37433	37525	37525	36994	36994	35044	35044	34479	34479
30	inj_p_type	住居・工場・道路などの傷害が発生した場所の種類を表す数字	1	記載あれば「有」	50014	50014	47978	47978	47954	47954	48179	48179	45897	45897	45593	45593
31	inj_detail	傷害が発生したところの種類にない場合の記述(その他)	40	記載あれば「有」	14020	14020	13527	13527	13435	13435	13450	13450	13359	13359	12951	12951
32	inj_tdfk	傷害発生都道府県	8	記載あれば「有」	48211	48211	46278	46278	46459	46459	46897	46897	44569	44569	44499	44499
33	inj_city	傷害発生市郡	12	記載あれば「有」	47783	47783	45857	45857	46155	46155	46554	46554	44271	44271	44167	44167
34	inj_town	傷害発生区町村	18	記載あれば「有」	20275	20275	19414	19414	19618	19618	19628	19628	18805	18805	18605	18605
35	inj_method	傷害発生時の手段及び状況に関する詳細記述	120	記載あれば「有」	49125	49894	47148	47882	47111	47850	47669	48409	45524	46229	45230	45951
36	inj_biko	傷害発生に関して備考がある場合はフラグ	1	記載あれば「有」	0	0	0	0	0	0	2	2	1	1	0	0
37	inf_detail	妊娠・分娩時における母体の病態または以上に関する詳細記述	84	記載あれば「有」	621	621	661	661	563	563	582	582	548	548	547	547
38	inf_biko_yn	「生後1年未満での病死」に関する備考がある場合はフラグ	1	記載あれば「有」	0	0	0	0	0	0	1	1	0	0	0	0
39	sonota	その他に付言すべきことに関する詳細記述	60	記載あれば「有」	84312	84105	83150	83159	84982	85135	87383	87572	84598	84844	86204	86465
40	biko_yn	備考欄に外字がある場合はフラグ	1	記載あれば「有」	203493	203493	211584	211584	223953	223953	242542	242542	249490	249490	249037	249037
41	biko_detail	備考に関する詳細記述	1024	記載あれば「有」	203492	34040	211584	36226	223952	36866	242542	40277	249490	40039	249036	40729
42	ym	調査年	2		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
43	sub_ym	死亡票提出年月日	4		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
44	jiken_sb	事件簿番号の補助番号	1		7	7	30	30	3	3	0	0	0	0	0	0
45	sex	性別	1		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
46	birth_yn	出生年月日があるか	1		524	524	489	489	578	578	522	522	586	586	584	584
47	birth_yumd	出生年月日	7		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
48	birth_hm	出生時分	4		915	915	912	912	868	868	858	858	801	801	743	743
49	death_yn	死亡年月日があるか	1		21	21	19	19	24	24	30	30	27	27	19	19
50	death_yumd	死亡年月日	7		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
51	death_hm	死亡時分	4		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
52	ori_dc	原死因コード	5		1236039	1236039	1265327	1265327	1309551	1309551	1358114	1358114	1376694	1376694	1366215	1366215
53	gain	外因コード	4		65283	65283	63794	63794	66994	66994	68948	68948	66577	66577	66067	66067
54	mom_sy	母親病態コード	5		1783	1783	1815	1815	1675	1675	1703	1703	1592	1592	1445	1445
55	twin_yn	単胎か多胎か数字	2		1783	1783	1815	1815	1675	1675	1703	1703	1592	1592	1445	1445
56	preg_yn	妊娠の有無	1		1783	1783	1815	1815	1675	1675	1703	1703	1592	1592	1445	1445
57	preg_w	妊娠週数	2		1783	1783	1815	1815	1675	1675	1703	1703	1592	1592	1445	1445
58	mon_b_yumd	母の生年月日	7		1783	1783	1815	1815	1675	1675	1703	1703	1592	1592	1445	1445
59	p_preg_bir	前回までの出生児数	2		1783	1783	1815	1815	1675	1675	1703	1703	1592	1592	1445	1445
60	p_preg_dea	前回までの死産児数	2		1783	1783	1815	1815	1675	1675	1703	1703	1592	1592	1445	1445
61	child_1	今回の出生子含む出生子	2		1783	1783	1815	1815	1675	1675	1703	1703	1592	1592	1445	1445

62	child_2	出産児	2	1783	1783	1815	1815	1675	1675	1703	1703	1592	1592	1445	1445
		付帯情報有の個数		402535	402270	409253	409065	421472	421300	440264	440176	443578	443521	438377	438271

注: (1) 本「突合DB(ユニークキーのみ)」とは死亡票・死亡個表の突合結果(「突合DB」)に対し、**届出地(sub_area)・事件簿番号(jiken_nm)・処理年月(shori_ym)**の組み合わせ

キーとし、重複がないものだけに絞ったものである。(元の「突合DB」の98.84%)

(2) 各欄からは、入力時の文字数制限により、入り切らない文字列が備考欄(biko_detail)に溢れて記入されることがあるが、

本研究ではこの「備考欄に溢れた文字列を正規表現ルールにより、元の**残るべき項目へ可能な限り復元して結合する処理**を行った。

(3) 各年の「なし」はこの**処理を行う前の「突合DB(ユニークキーのみ)」**の各項目件数、「**処理後**」はこの**処理を行った後の件数**を示している。

(4) 5列目「付帯情報の有無の判定基準」とは、「**付帯情報として用いた項目**」について「**どういう条件であれば、付帯情報ありとして判定したか**」の基準が入力されている

(5) 件数中の「**太字・背景薄灰色**」の箇所は、この備考欄処理によって、**溢れ文字列が元の項目へ復元された結果、件数に変化があった箇所**を示す。

特に項番41は、**備考欄の文字列が大幅に「本来あるべき項目」へ戻されている**ことを示している。

表C: IRIS処理結果と死亡票データ(突合 DBユニークキーのみ)を合わせた「統合テーブル」仕様

配列番号	項目	意味	機械学習用基本ベクトルでの使用	備考
0	CertificateKey	仮名ID	YES	突合DB(ユニークキー)における仮名ID
1	DateBirth	生年	YES	機械学習用基本ベクトルでは、これらを合わせて、死亡時年齢に変換。(小数点以下2位で四捨五入)
2	DateDeath	没年	YES	
3	Sex	性別	YES	1: 男性 2: 女性
4	Iris(反転)	IRIS による仮原死因コード	YES	外因コードがある場合、日本仕様に合わせ死因コードと順番を入れ替え済み
5	MainInj	IRIS が付与した外因コード	NO	(例) V494 など
6	Reject	IRIS: Reject されたか否か & 種別	NO	No, Maybe, Code, Syntax, MainInjury, Muse
7	SelectedCodes	IRIS: コード置換・修正処理前の ICD10コード	NO	(例) Ia: R54(2Hours)/II: D531(30Months) など
8	SubstitutedCodes	IRIS: IRIS内辞書により置換されたICD10コード	NO	(例) R54(2Hours)*D531(30Months) など
9	ErnCodes	IRIS: N/A	NO	IRIS Ver5 以降使用されていない
10	AcmeCodes	IRIS: MUSE 処理後の ICD10コード	YES	(例) R54*D531 など。 機械学習用の「用いられたICD-10コードに基づいた数値ベクトル」には仮原死因コードと合わせて、この欄の処理結果を用いる。
11	MultipleCodes	IRIS: 多重分類分類用の ICD10コード	NO	(例) D531 R54 など。アルファベットソートされている。
12	ToDoList	IRIS: GUI インタフェースで入力されたTODO	NO	使っていない。
13	確定	日本: 死亡票に基づく確定原死因コード	NO	(例) D531 など
14	外因	日本: 死亡票に基づく外因コード	NO	(例) V494 など
15	母側	日本: 死亡票に基づく母側病態コード	NO	(例) P008C など
16	一致/不一致	仮原死因コードと確定原死因コードの一致	YES	0: 一致、1: 不一致。機械学習では正解データとして使用
17	Ia_c	A欄記載の傷病名	NO	記載あり:1、記載なし:0
18	Ia_p	期間	NO	記載あり:1、記載なし:0
19	Ib_c	I欄記載の傷病名	NO	記載あり:1、記載なし:0
20	Ib_p	期間	NO	記載あり:1、記載なし:0
21	Ic_c	U欄記載の傷病名	NO	記載あり:1、記載なし:0
22	Ic_p	期間	NO	記載あり:1、記載なし:0
23	Id_c	E欄記載の傷病名	NO	記載あり:1、記載なし:0
24	Id_p	期間	NO	記載あり:1、記載なし:0
25	Il_c	II欄記載の傷病名	NO	記載あり:1、記載なし:0
26	Il_p	期間	NO	記載あり:1、記載なし:0
27	ope_flg	手術の有無	YES	付帯情報項目
28	ope_detail	手術の詳細記述	YES	付帯情報項目
29	ope_bk_flg	手術に関して備考がある場合はフラグ	YES	付帯情報項目
30	ope_ymd	手術年月日	YES	付帯情報項目
31	ap_flg	解剖の有無	YES	付帯情報項目
32	ap_detail	解剖の詳細記述	YES	付帯情報項目
33	ap_bk_flg	解剖に関して備考がある場合はフラグ	YES	付帯情報項目
34	dc_type	自然死・不慮の外因死・その他の不詳の死などの種類を表す数字	YES	付帯情報項目
35	inj_ymd	障害発生年月日	YES	付帯情報項目
36	inj_hm	障害発生時分	YES	付帯情報項目
37	inj_p_type	住居・工場・道路などの障害が発生した場所の種類を表す数字	YES	付帯情報項目
38	inj_detail	障害が発生したところの種類にない場合の記述(その他)	YES	付帯情報項目
39	inj_tdfk	傷害発生都道府県	YES	付帯情報項目
40	inj_city	傷害発生市郡	YES	付帯情報項目
41	inj_town	障害発生区町村	YES	付帯情報項目
42	inj_method	障害発生時の手段及び状況に関する詳細記述	YES	付帯情報項目
43	inj_biko	傷害発生に関して備考がある場合はフラグ	YES	付帯情報項目
44	inf_detail	妊娠・分娩時における母体の病態または以上に関する詳細記述	YES	付帯情報項目
45	inf_biko_yn	「生後1年未満での病死」に関して備考がある場合はフラグ	YES	付帯情報項目
46	sonota	その他に付言すべきことに関する詳細記述	YES	付帯情報項目
47	biko_yn	備考欄に外字がある場合はフラグ	YES	付帯情報項目
48	biko_detail	備考に関する詳細記述	YES	付帯情報項目
49	付帯有無	付帯情報の有り無し	NO	27~48に「付帯情報有りの基準」を満たすものがあれば1、それ以外0。これが1のものだけを機械学習の対象とした。

注:

- (1) 死亡個票中の全病名は項番17-26に相当するが、機械学習用基本ベクトルにはこのままではなく、IRIS でのICD-10コード修正結果である AcmeCodes(項番 10)を用いた。を用いた。IRIS でのコード修正機能により、期間表現を利用して「急性」のコードに変換する、などの処理がなされるためである。

表D:各処理段階の件数詳細

L	西暦年	2020	2019	2018	2017	2016	2015	2017-20	ALL		
1	元号	令和2年	平成31年/ 令和1年	平成30年	平成29年	平成28年	平成27年	N/A	N/A		
2	国内使用 ICD-10 バージョン	2013年版				2003年版		2013年版	N/A		
3	死亡票件数	1384300	1393504	1374469	1351944	1319030	1301379	5504217	8124626		
4	死亡個票件数	1384568	1393900	1374780	1325955	1280853	1249469	5479203	8009525		
5	突合DB	レコード数	1384263	1392578	1373589	1325413	1280808	1248057	5475843	8004708	
6		男性	713650	714048	705560	683577	661501	645711	2816835	4124047	
7		女性	670613	678530	668029	641836	619307	602346	2659008	3880661	
8	重複排除	キー(届出地+事件簿番号+提出年月日)の重複により排除したレコード数 (L5-L9)		18048	15884	15475	15862	15481	12018	65269	92768
9	突合DB (ユニークキー)	レコード数	1366215	1376694	1358114	1309551	1265327	1236039	5410574	7911940	
10		割合(／L5)	98.70%	98.86%	98.87%	98.80%	98.79%	99.04%	98.81%	98.84%	
11		男性	703561	705117	697105	674859	652853	639005	2780642	4072500	
12		女性	662654	671577	661009	634692	612474	597034	2629932	3839440	
13		付帯情報あり(備考欄文字列前処理後)	438271	443521	440176	421300	409065	402270	1743268	2554603	
14		割合(／L5)	32.1%	32.2%	32.4%	32.2%	32.3%	32.5%	32.2%	32.3%	
17	IRIS処理	IRIS処理に回した件数 (全傷病名がICD-10コーディング可能)	1104767	1111603	1088872	1047946	1008941	984617	4353188	6346746	
18		割合(／L9)	80.86%	80.74%	80.18%	80.02%	79.74%	79.66%	80.46%	80.22%	
19		IRIS で仮原因決定可能(Rejectなし)	1095698	1102884	1080590	1040397	1002020	977952	4319569	6299541	
20		割合(／L17)	99.18%	99.22%	99.24%	99.28%	99.31%	99.32%	99.23%	99.26%	
21		IRIS による Reject(と原因の内訳)	9069	8719	8282	7549	6921	6665	33619	47205	
22		割合(／L17)	0.82%	0.78%	0.76%	0.72%	0.69%	0.68%	0.77%	0.74%	
23		Syntax	102	115	119	109	97	97	445	639	
24		Code	1014	1014	949	891	887	842	3868	5597	
25		MultipleCause	0	0	0	0	0	0	0	0	
26		MayBe	7820	7461	7064	6417	5791	5552	28762	40105	
27		Muse	55	62	66	60	63	61	243	367	
28		Interval	0	0	0	0	0	0	0	0	
29		MainInjury	78	67	84	72	83	113	301	497	
30		Coder	0	0	0	0	0	0	0	0	
31	仮原因(UCコード)が付与された件数	1104535	1111357	1088639	1047738	1008741	984422	4352269	6345432		
32	外因コードが付与された件数	36603	36364	37180	35741	34973	36069	145888	216930		
33	解析用統合テーブル (BasicTable)	コード修正・非対象レコード削除後 (L19中)	1095585	1102777	1080458	1040277	1001497	977469	4319097	6298063	
34		割合(／L19)	99.99%	99.99%	99.99%	99.99%	99.95%	99.95%	99.99%	99.98%	
35		付帯情報あり	319434	322559	317598	303141	292776	287358	1262732	1842866	
36		割合(／L33)	29.2%	29.2%	29.4%	29.1%	29.2%	29.4%	29.2%	29.3%	
37		仮原因変更有り	40493	40945	41516	39700	48190	48134	162654	258978	
38		割合(／L35)	12.7%	12.7%	13.1%	13.1%	16.5%	16.8%	12.9%	14.1%	
39		追加コードあり	7288	7307	7696	6988	7601	7699	29279	44579	
40		追加コードなし	33205	33638	33820	32712	40589	40435	133375	214399	
41		仮原因変更なし	278941	281614	276082	263441	244586	239224	1100078	1583888	
42		割合(／L35)	87.3%	87.3%	86.9%	86.9%	83.5%	83.2%	87.1%	85.9%	
43		追加コードあり	27739	27876	28509	28061	26389	27371	112185	165945	
44		追加コードなし	251202	253738	247573	235380	218197	211853	987893	1417943	
45		付帯情報なし	776151	780218	762860	737136	708721	690111	3056365	4455197	
46	割合(／L33)	70.8%	70.8%	70.6%	70.9%	70.8%	70.6%	70.8%	70.7%		
47	仮原因変更有り	40471	41391	41090	41056	63104	63032	164008	290144		
48	割合(／L45)	5.2%	5.3%	5.4%	5.6%	8.9%	9.1%	5.4%	6.5%		
49	追加コードあり	2065	2057	2088	1975	1706	1795	8185	11686		
50	追加コードなし	38406	39334	39002	39081	61398	61237	155823	278458		
51	仮原因変更なし	735680	738827	721770	696080	645617	627079	2892357	4165053		
52	割合(／L45)	94.8%	94.7%	94.6%	94.4%	91.1%	90.9%	94.6%	93.5%		
53	追加コードあり	5537	5386	5396	5149	4113	3821	21468	29402		
54	追加コードなし	730143	733441	716374	690931	641504	623258	2870889	4135651		
55	機械学習用データ	対象レコード数 (=L35)	319434	322559	317598	303141	292776	287358	1262732	1842866	
56		ICDコード種類数 (Acme)	2061	2091	2066	2050	2006	2024	2693	2853	

- 注:
- (1) 2015, 2016 年と、2017~2020年 で原因コードに使用されたICD-10のバージョンが異なる。(それぞれ 2003年版準拠、2013年版準拠)
 - (2) 本研究では全体を通じて、病名のICD-10コーディングに「標準病名マスター-v5.04 (ICD-10 2013年版準拠)」を用いている。
 - (3) そのため、「2015,2016年」は「2017~2020年」と比べて、「全病名にICD-10コーディング可能な割合」「仮原因変更有り割合」など、いくつかの項目で低くなっている。このような使用されるICD-10バージョンによる影響が確認された値は薄灰色の網掛けで表示してある。
 - (4) 機械学習以降の処理については、別表「各種機械学習手法の結果一覧」を参照のこと

表E: BERTモデルを用いた各学習手法の予測精度一覧

手法		RedLR	BalRedLR	BalRedLR_ BertOnly	BalRedLR_ Norm	CwRedLR	CyLR	BalCyLR	CwCyLR
Accuracy		0.914	0.866	0.819	0.863	0.877	0.900	0.858	0.882
F1-score		0.591	0.535	0.391	0.548	0.543	0.564	0.510	0.549
Precision		0.756	0.481	0.343	0.475	0.517	0.641	0.457	0.537
Recall		0.485	0.602	0.455	0.648	0.571	0.504	0.576	0.561
Confusion Matrix (正解, 予測)	(0, 0)	13918	12890	12422	12750	13130	13656	12817	13235
	(0, 1)	328	1356	1824	1496	1116	590	1429	1011
	(1, 0)	1078	833	1140	737	897	1038	888	918
	(1, 1)	1014	1259	952	1355	1195	1054	1204	1174

表F: 各種機械学習手法の結果一覧

対象年		2020	2019	2018	2017	2016	2015	
使用ICD-10バージョン		2013年版				2003年版		
機械学習用データ	対象レコード数	319434	322559	317598	303141	292776	287358	
	ICDコード種類数 (Acme)	2061	2091	2066	2050	2006	2024	
	共通ベクトルの次元数 (ID, 正解を除く)	2085	2115	2090	2074	2030	2048	
	付帯情報意味ベクトルの次元数							
	Embedding手法							
	BASELINE	0	0	0	0	0	0	
	TFIDF	2938	2938	2938	2938			
	LSI	195	195	195	195			
	WORD2VEC	200	200	200	200			
	DOC2VEC (PV-DM)	200	200	200	200			
	DOC2VEC (PV-DBOW)	200	200	200	200			
	学習用ベクトルの次元数							
	BASELINE	2085	2115	2090	2074	2030	2048	
	TFIDF	5023	5053	5028	5012			
LSI	2280	2310	2285	2269				
WORD2VEC	2285	2315	2290	2274				
DOC2VEC (PV-DM)	2285	2315	2290	2274				
DOC2VEC (PV-DBOW)	2285	2315	2290	2274				
学習結果 (予測精度)	評価尺度	Embedding手法						
	Accuracy (最良)	BASELINE	0.920	0.917	0.915	0.915	0.907	0.904
		TFIDF	0.949	0.948	0.949	0.949		
		LSI	0.945	0.944	0.944	0.945		
		WORD2VEC	0.948	0.942	0.942	0.943		
		DOC2VEC (PV-DM)	0.933	0.930	0.931	0.930		
		DOC2VEC (PV-DBOW)	0.944	0.941	0.942	0.943		
	ROC-AUC	BASELINE	0.925	0.923	0.922	0.921	0.927	0.928
		TFIDF	0.953	0.949	0.951	0.950		
		LSI	0.945	0.943	0.942	0.943		
		WORD2VEC	0.943	0.940	0.941	0.942		
		DOC2VEC (PV-DM)	0.933	0.929	0.932	0.930		
		DOC2VEC (PV-DBOW)	0.944	0.943	0.943	0.946		
	PR-AUC	BASELINE	0.742	0.727	0.735	0.735	0.793	0.791
		TFIDF	0.857	0.847	0.859	0.859		
		LSI	0.835	0.828	0.834	0.839		
		WORD2VEC	0.828	0.817	0.826	0.830		
		DOC2VEC (PV-DM)	0.782	0.771	0.785	0.780		
		DOC2VEC (PV-DBOW)	0.830	0.821	0.831	0.837		