

令和元年6月24日現在

機関番号：33916

研究種目：研究活動スタート支援

研究期間：2017～2018

課題番号：17H07421

研究課題名（和文）市町村が収集・蓄積しているデータベース活用による、認知機能低下の予測モデル構築

研究課題名（英文）Predicting cognitive decline from routinely collected data

研究代表者

尾形 宗士郎（Ogata, Soshiro）

藤田医科大学・医療科学部・講師

研究者番号：00805012

交付決定額（研究期間全体）：（直接経費） 2,100,000円

研究成果の概要（和文）：本研究は宮崎県延岡市と疫学研究を共同で実施し地域在住後期高齢者約480名の認知機能に関するデータを収集した。機械学習法の一つであるスパースモデリングを活用し、市が日常業務として収集・蓄積している情報のみを使用し、認知機能障害に対する予測モデルを作成した（AUC：0.70-0.71）。また本研究調査で追加収集した教育歴、言語流暢性テスト、短期記憶テスト等々（現状の健康診査では収集されていない情報）を加えることで予測精度が向上した（AUC：0.74-0.96）。本研究は既存情報＋簡便な追加検査の情報のみであっても十分な予測精度を有する地域在住後期高齢者を対象とした認知機能障害予測モデルを提示した。

研究成果の学術的意義や社会的意義

認知症対策として認知機能低下を早期発見し予防につなげることが重要な戦略である。本研究で開発した地域在住後期高齢者を対象とした認知機能障害予測モデルは、日本の市町村が既に業務の一環として収集している情報を主として使用している。そのため、市町村は本研究で開発した予測モデルを低コストで今すぐ運用することが可能である。以上のことから、本研究は認知症早期発見システム構築に貢献すると考える。

研究成果の概要（英文）：The present study conducted an epidemiological study collaborating with Nobeoka city, and collected data related to cognitive function from 480 community-dwelling people aged 75 years old and over. By utilizing a sparse modeling as a machine learning method, we developed a diagnostic prediction model of cognitive impairment (i.e., the Routine prediction model) only based on information routinely collected by Japanese cities. This had moderate predictability (AUC: 0.70-0.71). Additionally, we developed an additional diagnostic prediction model (i.e., the Additional prediction model) by adding education level, category fluency test, and short-term memory test to the Routine prediction model. This had relatively high predictability (AUC: 0.74-0.96). Thus, the present study developed the diagnostic prediction models of cognitive impairment with moderate to high predictability, which just required information collected routinely by cities and/or additional information collected easily.

研究分野：応用健康科学

キーワード：認知機能 予測モデル 自治体保有情報 ビッグデータ 循環器病リスク要因 機械学習 AI

様式 C-19、F-19-1、Z-19、CK-19（共通）

## 1. 研究開始当初の背景

認知症者の有病数は世界中で約 3600 万人、日本では 65 歳以上で約 462 万人（有病割合推定値 15%）と推定されている。認知症は発症すると治療・回復が困難なため、認知機能低下を早期発見し予防につなげる必要がある。そのため、認知機能低下者の早期発見を日常業務として実施できる認知機能スクリーニングシステムの構築は急務課題である。

認知機能スクリーニングを日常業務として実施するには、国・市町村が既に日常業務で収集している情報を活用することが現実的かつ効果的であると考えられる。また、認知機能低下のリスク要因に介入可能な循環器病リスク要因（高血圧、肥満、糖尿病等々）が報告されている。特に世界の認知症ケースの約 30% は介入可能なリスク要因によると推定された<sup>1</sup>。

しかし、日本において国・市町村が既に日常業務として収集している情報を活用した認知機能低下の予測モデルは構築されていない。加えて、循環器病リスク要因や既往歴の変遷を、高齢者、特に認知機能が低下した者から問診や質問紙等で正確に収集するのは困難である。

そこで本研究は、市町村が既に日常業務の一環として 2010 年以降収集している診療報酬明細書（レセプトデータ）、健康診査データ、介護認定調査時のデータに着目する。特に、健康診査では循環器病リスク要因が検査項目となっていてデータが収集されている。加えて、介護認定調査結果は介護保険支給のために、自立度や認知機能レベルを評価する項目がありデータが収集されている。さらに、予測に有用・無用な変数を判断し、高精度の予測モデルを作成する機械学習に、本研究は着目する。

## 2. 研究の目的

本研究の目的は下記 3 つとした。市町村が日常業務として収集・蓄積データ（レセプト、健康診査、介護認定調査）と機械学習法を活用することで、

- 1) 地域在住後期高齢者を対象とした現在の認知機能障害の予測モデルを作成する
- 2) 地域在住後期高齢者を対象とした半年から一年後の認知機能障害の予測モデルを作成する。
- 3) これらの予測モデルに、認知機能を簡便に把握する項目（教育歴や認知機能検査の一部など）を加えることで予測精度が向上するか検討する。

## 3. 研究の方法

### 研究デザインと対象者

宮崎県延岡市にて地域在住の後期高齢者（75 歳以上）を対象に、コホートデザインにて疫学研究を実施した。なお、追跡調査はベースライン調査から平均 9 カ月（sd=1 カ月）後に実施した。本研究の参入基準は、1) 75 歳以上の後期高齢者であること、2) 延岡市が実施した 2017 年 4 月から 11 月の後期高齢者健診に参加したこと、3) 本研究の電話調査への参加に同意した者とした。なお、本研究の除外基準は、1) 研究に参加することが困難なほど認知機能が低下した者、2) 研究参加が困難なほど聴力が低下した者とした。延岡市から研究調査依頼をしたところ 568 名の市民から研究参加意向の返信があり、そのうち 480 名が本研究で実施した電話調査に参加した。電話調査参加が 2018 年 2 月中旬から 4 月中旬の参加者を予測モデル作成の training dataset とし、4 月下旬から 8 月中旬の参加者を予測モデルの精度評価をする validation dataset とした。なお、本研究は国立循環器病研究センターと藤田医科大学の倫理委員会から承認を得て実施した。

### 評価項目

認知機能：本研究ではトレーニングを受けた者が研究対象者に電話調査を実施し、認知機能を Telephone Interview for Cognitive Status (TICS-J) で評価した。TICS は Mini-Mental State Exam (MMSE) の電話調査バージョンであり、世界各国の疫学調査で認知機能測定に使用されている<sup>2</sup>。TICS-J は認知症に対して 98.0% の感度と 90.7% の特異度を有すると報告されている<sup>2</sup>。本研究では先行研究で報告された TICS-J のカットオフ値である 32 点以下のものを認知機能障害ありとした。

予測変数群：本研究の予測変数候補として下記変数を検討した。基本情報として、年齢（80 歳以上/未満）、性別、健康診査参加日と電話調査実施日の日間差（30 日単位）、介護保険使用（有/無）、認知症高齢者自立度（I 度以上/未満）とした。循環器病リスク要因として、Body mass index (BMI) 高値（25 kg/m<sup>2</sup> 以上/未満）、循環器病の既往歴（有/無）、高血圧に対する服薬（有/無）、糖尿病に対する服薬（有/無）、脂質異常症に対する服薬（有/無）、血圧高値（SBP 135 mmHg 以上/未満）、HbA1c 高値（5.5%/以上/未満）、high density lipoprotein cholesterol (HDL-c) 低値（65 mg/dL/未満/以上）、low density lipoprotein cholesterol (LDL-c) 高値（120 mg/dL 以上/未満）、estimated glomerular filtration rate (eGFR) 低値（55 mL/min/1.73m<sup>2</sup> 未満/以上）とした。加えて、血圧高値と高血圧服薬、HbA1c 高値と糖尿病服薬、HDL-c 低値と脂質異常症服薬、そして LDL-c 高値と脂質異常症服薬の交互作用項も予測変数として加えた。これらの情報はすべて延岡市のデータベースから抽出したものである。

認知症の古典的リスク要因として教育歴（10 年以上/未満）と言語流暢性テスト（14 点以上/以下）、短期記憶テスト（10 点満点の連続値）として、予測変数として使用した。これらの情

報は電話調査時に収集した。なお、言語流暢性テストでは、1分間にできるだけ多くの動物の名前を対象者に言っただき、その数を点数とした。短期記憶テストはTICS-Jに含まれている10単語の即時再生を利用した。

#### 統計解析

認知機能障害の予測モデルを作成するために、本研究では least absolute shrinkage and selection operator (lasso) logistic regression analyses を使用した。Lasso の特徴として予測に不要な変数の回帰係数を0にして推定可能な点であり、これにより多変数の中から予測に有用な少数の変数とその回帰係数を選択することができる。本研究では lasso logistic regression analyses を training dataset にて 10-fold cross validation を実施し予測モデルを作成した。なお、10-fold cross validation とはデータを10分割しそのうちの9つのデータで予測モデルを作成し残りの1つのデータで検証を行う作業を10回繰り返す手法である。これにより新しいデータセットに対する頑健さをもつ予測モデルを作成可能である。Training dataset で作成した予測モデルを validation dataset にて予測力の精度を Area under the curve にて評価した。

#### 4. 研究成果

本研究の参入基準をみたま者に対して、延岡市から研究調査依頼をしたところ568名の市民から研究参加意向の返信があり、そのうち480名が本研究で実施した電話調査に参加した。そのうち363名(76%)が training dataset とし、114名(24%)が validation dataset とした。平均年齢(SD)は、training dataset で80.6(4.0)歳、validation dataset で80.7(4.3)歳であった。対象者の特徴の詳細を Table 1 に記載した。

Table 1. Characteristics of the present participants aged 75 years old and over in Nobeoka city, Japan.

Datasets	Training dataset	Validation dataset
N (%)	363 (76%)	114 (24%)
<b>Continuous variables, means (sd)</b>		
Age at telephone interview (years old)	80.64 (4.01)	80.68 (4.34)
TICS-J score	33.80 (3.11)	33.72 (3.43)
Education years	10.86 (2.15)	10.80 (2.41)
Body mass index (kg/m <sup>2</sup> )	22.86 (2.97)	23.44 (3.02)
SBP (mmHg)	133.65 (16.22)	135.60 (15.93)
DBP (mmHg)	71.38 (10.05)	71.62 (10.06)
HbA1c (%)	5.81 (0.53)	5.89 (0.57)
HDL-c (mg/dl)	59.62 (14.65)	57.90 (15.38)
LDL-c (mg/dl)	119.70 (28.62)	109.95 (28.02)
eGFR	62.38 (15.49)	60.82 (14.12)
Category fluency test's score	13.37 (3.84)	13.38 (3.64)
Immediate recall of 10 words	5.46 (1.91)	5.73 (2.02)
<b>Categorical variables, n (%)</b>		
Women	236 (65.0)	36 (31.6)
Dementia grade ≥I	17 (4.7)	4 (3.5)
Use of long-term care insurance	37 (10.2)	13 (11.4)
Medication use for hypertension	187 (51.5)	71 (62.3)
Medication use for diabetes	33 (9.1)	14 (12.3)
Medication use for dyslipidemia	101 (27.8)	35 (30.7)
Medical history of cardiovascular diseases	81 (22.3)	21 (18.4)

Abbreviations: sd, standard deviation; TICS-J, Telephone Interview for Cognitive Status in Japanese; SBP, Systolic blood pressure; DBP, Diastolic blood pressure; HDL-c, high density lipoprotein cholesterol; LDL-c, low density lipoprotein cholesterol; eGFR, estimated glomerular filtration rate.

Tables 2. Coefficients<sup>1</sup> and performance of the prediction models for baseline cognitive impairment obtained by lasso logistic regression analyses

	Model 1	Model 2	Model 3
<b>Predictors</b>			
Intercept	-2.256	-3.327	2.155
Age >=80	0.163	0.210	-
Women	-	0.051	0.639
Dementia grade >=I in long-term care insurance	1.069	1.533	0.725
Use of long-term care insurance	0.427	0.252	0.176
Difference days between the health check-up and the telephone interview (per 30 days)	0.106	0.160	0.142
High level of body mass index >=25 kg/m <sup>2</sup>	0.187	0.253	-
High level of SBP >=135 mmHg	0.391	0.409	0.319
High level of HbA1c >=5.5%	-	0.079	0.011
Low level of HDL-c <65 mg/dL	0.168	0.314	0.015
High level of LDL-c >=120 mg/dL	-0.330	-0.519	-0.495
Low level of eGFR <55 mL/min/1.73m <sup>2</sup>	0.263	0.498	0.216
Medication use for hypertension	-	-0.387	-
Medication use for diabetes	-	1.916	-
Medication use for dyslipidemia	-	-	-
Medical history of cardiovascular diseases	-	0.187	-
Interaction between high level of SBP and medication use for hypertension	0.161	0.527	0.039
Interaction between high level of HbA1c and medication use for diabetes	-	-1.995	-
Interaction between low level of HDL-c and medication use for dyslipidemia	0.232	0.245	0.551
Interaction between high level of LDL-c and medication use for dyslipidemia	-	-0.132	-
Education years <10 years	NA	1.050	0.381
Category fluency test's score < 14	NA	NA	0.722
Immediate recall of 10 word (continuous scores)	NA	NA	-1.101
<b>Performance of the prediction models</b>			
<b>Training dataset</b>			
Area Under the Curve (95% confidence interval)	0.70 (0.64-0.76)	0.75 (0.69-0.80)	0.91 (0.88-0.94)
Positive predictive values	0.77	0.68	0.82
Negative predictive values	0.74	0.77	0.88
Accuracy	0.75	0.76	0.87
<b>Validation dataset</b>			
Area Under the Curve (95% confidence interval)	0.71 (0.60-0.82)	0.74 (0.64-0.85)	0.96 (0.92-0.99)
Positive predictive values	0.67	0.56	0.89
Negative predictive values	0.74	0.81	0.88
Accuracy	0.73	0.73	0.89

Abbreviations: SBP, Systolic blood pressure; HDL-c, high density lipoprotein cholesterol; LDL-c, low density lipoprotein cholesterol; eGFR, estimated glomerular filtration rate; CVD, cardiovascular diseases.

1 Dashes (i.e., -) meant that the variables were not selected as the predictors by the lasso logistic regression analyses.

Tables 3. Coefficients<sup>1</sup> and performance of prediction models for 6-12 months later cognitive impairment obtained by lasso logistic regression analyses.

	Model 1	Model 2
<b>Predictors</b>		
Intercept	-1.545	-0.271
Age >=80	0.583	0.460
Women	-0.075	-
Dementia grade >=I in long-term care insurance	-	-
Use of long-term care insurance	-	-
Difference days between the health check-up and the telephone interview (per 30 days)	-	-
High level of body mass index >=25 kg/m <sup>2</sup>	-	-
High level of SBP >=135 mmHg	-	-
High level of HbA1c >=5.5%	-	-
Low level of HDL-c <65 mg/dL	-	-
High level of LDL-c >=120 mg/dL	-0.455	-0.438
Low level of eGFR <55 mL/min/1.73m <sup>2</sup>	-	-
Medication use for hypertension	0.488	0.421
Medication use for diabetes	0.153	0.104
Medication use for dyslipidemia	-	-
Medical history of cardiovascular diseases	-	-
Interaction between high level of SBP and medication use for hypertension	-	-
Interaction between high level of HbA1c and medication use for diabetes	-	-
Interaction between low level of HDL-c and medication use for dyslipidemia	-	-
Interaction between high level of LDL-c and medication use for dyslipidemia	-0.272	-0.160
Education years <10 years	0.533	0.340
Category fluency test's score < 14	-	-
Immediate recall of 10 word (continuous scores)	-	-0.211
<b>Performance of the prediction models</b>		
<b>All dataset with 10-fold cross validation</b>		
Area Under the Curve (95% confidence interval)	0.72 (0.66-0.78)	0.73 (0.67-0.80)
Positive predictive values	0.58	0.92
Negative predictive values	0.74	0.75
Accuracy	0.73	0.76

Abbreviations: SBP, Systolic blood pressure; HDL-c, high density lipoprotein cholesterol; LDL-c, low density lipoprotein cholesterol; eGFR, estimated glomerular filtration rate; CVD, cardiovascular diseases.

<sup>1</sup> Dashes (i.e., -) meant that the variables were not selected as the predictors by the lasso logistic regression analyses.

ベースライン時認知機能障害に対する予測モデルを、lasso logistic regression analysesにて作成した (Table 2)。Model 1 では市町村が既に通常業務で収集している情報のみで作成した予測モデルである。Table 2 に回帰係数が記載されている変数が、各モデルにおいて認知機能障害を予測するうえで重要な変数であり、その重みづけが回帰係数で示されている。Model 1 の AUC は training dataset で 0.70 (0.64-0.76)、validation dataset で 0.71 (0.60-0.82) と中程度の予測力を保有していた。Model 2 は Model 1 に教育歴のみを追加したモデルであり、AUC が training dataset で 0.75 (0.69-0.80)、validation dataset で 0.74 (0.64-0.85) と予測力が向上した。Model 3 は Model 2 に言語流暢性テストと 10 単語の即時再生のスコアを追加投入したモデルであり、AUC が training dataset で 0.91 (0.88-0.94)、validation dataset で 0.96 (0.92-0.99) と予測力が向上した。ただし、Model 3 は TICS-J の一部課題である 10 単語の即時再生スコアを予測変数として使用したため、データに過剰適合している可能性がある。

フォローアップ時(平均追跡期間9カ月[SD=1カ月])の認知機能障害に対する予測モデルを、lasso logistic regression analysesにて作成した (Table 3)。フォローアップ時のサンプルサイズは 308 名と小さく、training dataset と validation dataset に分けて予測モデルを構築するのが困難であったため、全体のデータセットに対して 10-fold cross validation により予測モデルを構築した。Model 1 と Model 2 の AUC はそれぞれ 0.72 (0.66-0.78) と 0.73 (0.67-0.80) であった。いっぽうで Positive predictive value はそれぞれ、0.58 と 0.92 と Model 2 のほうが精度は高かった。

本研究によって市町村が通常業務で収集している既存情報と簡便な追加検査の情報で、現在の認知機能障害及び、約一年後の認知機能障害を中-高精度に予測するモデルを提示することでできた。本研究で提示した認知機能障害の予測モデルは、全国の市町村が有する既存情報を有効活用しているため、市町村の公衆衛生サービスとしての認知症早期発見システムに組み込むことが容易である。以上から本研究は認知症早期発見システム構築に貢献すると考える。

## References

- 1: Norton S, et al. Potential for primary prevention of Alzheimer's disease: an analysis of population-based data. *Lancet Neurol.* 2014 Aug;13(8):788-94.
2. Konagaya Y, et al. Validation of the Telephone Interview for Cognitive Status (TICS) in Japanese. *Int J Geriatr Psychiatry.* 2007 Jul;22(7):695-700.

## 5. 主な発表論文等

〔雑誌論文〕(計 0 件) 現在、認知機能障害の予測モデルを検討した結果を論文にし、国際雑誌から査読をうけている最中である。

### 〔学会発表〕(計 4 件)

1. 尾形宗士郎、西村邦宏、宮本恵宏. 地域在住後期高齢者における認知機能を予測する要因の検討. 日本公衆衛生学会. 2018年.
2. Soshiro Ogata, Michikazu Nakai, Misa Takegami, Kunihiro Nishimura, Yoshihiro Miyamoto. Developing a Prediction Model to Estimate Low Cognitive Function in Japanese People Aged 75 Years and Over. *Gerontological Society of America.* 2018.
3. Haruka Tanaka, Soshiro Ogata, Chisato Hayashi. Temporal order between depressive symptoms and subjective memory complaints using cross lagged panel model. *Gerontological Society of America.* 2018.
4. 尾形宗士郎, 清重映里, 竹上未紗, 中井陸運, 中尾葉子, 神出計, 西村 邦宏, 宮本恵宏. 地域在住後期高齢者における認知機能と過去数年間の循環器病リスク要因の経時変化の関連. 日本循環器病予防学会. 2019年.

〔図書〕(計 0 件) 〔産業財産権〕 出願状況 (計 0 件) 取得状況 (計 0 件) 〔その他〕 該当なし

## 6. 研究組織

### (1) 研究分担者

研究分担者氏名：該当なし

ローマ字氏名：

所属研究機関名：

部局名：

職名：

研究者番号 (8 桁)：

### (2) 研究協力者

研究協力者氏名：該当なし

ローマ字氏名：

※科研費による研究は、研究者の自覚と責任において実施するものです。そのため、研究の実施や研究成果の公表等については、国の要請等に基づくものではなく、その研究成果に関する見解や責任は、研究者個人に帰属されます。