

【公報種別】特許法第17条の2の規定による補正の掲載  
 【部門区分】第6部門第3区分  
 【発行日】令和6年9月13日(2024.9.13)

【国際公開番号】WO2023/132029  
 【出願番号】特願2023-572292(P2023-572292)

【国際特許分類】

G 0 6 F 1 6 / 2 1 5 ( 2 0 1 9 . 0 1 )

G 0 6 F 1 6 / 2 4 ( 2 0 1 9 . 0 1 )

G 0 6 N 2 0 / 0 0 ( 2 0 1 9 . 0 1 )

10

【F I】

G 0 6 F 1 6 / 2 1 5

G 0 6 F 1 6 / 2 4

G 0 6 N 2 0 / 0 0 1 6 0

【手続補正書】

【提出日】令和6年7月1日(2024.7.1)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

20

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

レコード対を取得する取得手段と、

前記レコード対について、複数の類似度関数を用いて複数の類似度を算出する類似度算出手段と、

前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行う予測手段と、

前記予測手段による予測結果を出力する出力手段と、

を備えている情報処理装置。

30

【請求項2】

前記取得手段は、補助データを更に取得し、

前記予測手段は、前記レコード対と、前記複数の類似度と、前記補助データとを参照して、前記レコード対と前記補助データとに応じて定まる重要度を用いて前記レコード対の同一性予測を行う

請求項1に記載の情報処理装置。

【請求項3】

前記予測手段は、前記レコード対を参照して前記重要度を算出する重要度算出手段を備えている

請求項1又は2に記載の情報処理装置。

40

【請求項4】

前記重要度算出手段は、前記複数の類似度の各々に関する重要度を算出し、

前記予測手段は、前記複数の類似度に関する線形和であって、前記各重要度を重み係数とする線形和を用いて、前記同一性予測を行う

請求項3に記載の情報処理装置。

【請求項5】

レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得する取得手段と、

予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が

50

有する 1 又は複数のパラメータ、及び

前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータ

の少なくとも何れかのパラメータを、前記訓練データを参照して生成するパラメータ生成手段と

を備えている情報処理装置。

【請求項 6】

レコード対を取得することと、

前記レコード対について、複数の類似度関数を用いて複数の類似度を算出することと、

前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行うことと、

前記レコード対の同一性予測による予測結果を出力することと、

を含む情報処理方法。

【請求項 7】

レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得することと、

予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する 1 又は複数のパラメータ、及び

前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータ

の少なくとも何れかのパラメータを、前記訓練データを参照して生成することと、

を含む情報処理方法。

【請求項 8】

レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得することと、

予測対象のレコード対について複数の類似度を算出するための複数の類似度算出モデル、及び

前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデル

の少なくとも何れかのモデルを、前記訓練データを参照して生成することと、

を含む学習済モデルの製造方法。

【請求項 9】

コンピュータに、

レコード対を取得する取得処理と、

前記レコード対について、複数の類似度関数を用いて複数の類似度を算出する類似度算出処理と、

前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行う予測処理と、

前記予測処理による予測結果を出力する出力処理と、

を実行させるプログラム。

【請求項 10】

コンピュータに、

レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得する取得処理と、

予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が

10

20

30

40

50

有する 1 又は複数のパラメータ、及び

前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータ

の少なくとも何れかのパラメータを、前記訓練データを参照して生成するパラメータ生成処理と、

を実行させるプログラム。

【手続補正 2】

【補正対象書類名】明細書

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、レコード対の同一性予測を行う技術に関する。

【背景技術】

【0002】

異なるテーブルに格納されたレコードから同一の又は類似するレコードの組み合わせを特定して対応付ける処理が行われている。このような処理は名寄せ処理とも呼ばれる。名寄せ処理によりテーブルの一元管理及びデータの拡張が可能となる。名寄せ処理を行う技術として、機械学習又はルールベースによるマッチングを行う技術が存在する。例えば、特許文献 1 及び非特許文献 1 には、機械学習により名寄せ処理を行う技術が記載されている。特に、特許文献 1 に記載の名寄せ処理装置は、情報処理装置と、記憶部と操作端末とから構成されている。この名寄せ処理装置は、レコード対の類似度を計算する類似度関数を複数用いてレコード対の類似度を計算し、訓練データを用いた機械学習により類似度の重みを学習する。

【先行技術文献】

【特許文献】

【0003】

【特許文献 1】日本国特開 2019 - 185244 号公報

【非特許文献】

【0004】

【非特許文献 1】Pradap Konda, et. al., Magellan: Toward Building Entity Matching Management Systems, Proceedings of the VLDB Endowment, 2016

【発明の概要】

【発明が解決しようとする課題】

【0005】

ところで、レコード対の同一性を判定する手法として様々な手法が存在する。例えば、「アイス」と「あいす」のレコード対はカタカナをひらがなに表記変更することで同一性の判定精度を高くすることができる。また、「ポテトチップス」と「ポテチ」のレコード対は部分文字列を抽出することで同一性の判定精度を高くすることができる。このように、レコード対の同一性の判定に適した手法はレコード対のそれぞれで異なる場合がある。特許文献 1 及び非特許文献 1 に記載の技術では、レコード対によっては同一性の判定を適切に行えないという問題があった。

【0006】

本発明の一態様は、上記の問題に鑑みてなされたものであり、その目的の一例は、レコード対の同一性をより好適に予測できる技術を提供することである。

【課題を解決するための手段】

10

20

30

40

50

## 【0007】

本発明の一側面に係る情報処理装置は、レコード対を取得する取得手段と、前記レコード対について、複数の類似度関数を用いて複数の類似度を算出する類似度算出手段と、前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行う予測手段と、前記予測手段による予測結果を出力する出力手段と、を備えている。

## 【0008】

また、本発明の一側面に係る情報処理装置は、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得する取得手段と、予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する1又は複数のパラメータ、及び前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルが有する1又は複数のパラメータの少なくとも何れかのパラメータを、前記訓練データを参照して生成するパラメータ生成手段と、を備えている。

10

## 【0009】

また、本発明の一側面に係る情報処理方法は、レコード対を取得することと、前記レコード対について、複数の類似度関数を用いて複数の類似度を算出することと、前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行うことと、前記レコード対の同一性予測による予測結果を出力することと、を含む。

20

## 【0010】

また、本発明の一側面に係る情報処理方法は、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得することと、予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する1又は複数のパラメータ、及び前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルが有する1又は複数のパラメータの少なくとも何れかのパラメータを、前記訓練データを参照して生成することと、を含む。

30

## 【0011】

また、本発明の一側面に係る製造方法は、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得することと、予測対象のレコード対について複数の類似度を算出するための複数の類似度算出モデル、及び前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルの少なくとも何れかのモデルを、前記訓練データを参照して生成することと、を含む。

## 【0012】

また、本発明の一側面に係るプログラムは、コンピュータに、レコード対を取得する取得処理と、前記レコード対について、複数の類似度関数を用いて複数の類似度を算出する類似度算出処理と、前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行う予測処理と、前記予測処理による予測結果を出力する出力処理と、を実行させる。

40

## 【0013】

また、本発明の一側面に係るプログラムは、コンピュータに、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得する取得処理と、予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する1又は複数のパラメータ、及び前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコー

50

ド対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータの少なくとも何れかのパラメータを、前記訓練データを参照して生成するパラメータ生成処理と、を実行させる。

【発明の効果】

【0014】

本発明の一態様によれば、レコード対の同一性をより好適に予測できる。

【図面の簡単な説明】

【0015】

【図1】例示的实施形態1に係る情報処理装置の構成を示すブロック図である。

【図2】例示的实施形態1に係る情報処理方法の流れを示すフロー図である。

10

【図3】例示的实施形態1に係る情報処理装置の構成を示すブロック図である。

【図4】例示的实施形態1に係る情報処理方法の流れを示すフロー図である。

【図5】例示的实施形態2に係る情報処理装置の構成を示すブロック図である。

【図6】例示的实施形態2に係る第1のデータと第2のデータの具体例を示す図である。

【図7】例示的实施形態2に係る情報処理方法の流れを示すフロー図である。

【図8】例示的实施形態2に係る統合済データの具体例を示す図である。

【図9】例示的实施形態3に係る情報処理装置の構成を示すブロック図である。

【図10】例示的实施形態3に係る情報処理方法の流れを示すフロー図である。

【図11】例示的实施形態4に係る情報処理装置の構成を示すブロック図である。

【図12】例示的实施形態4に係る画面表示例を示す図である。

20

【図13】各例示的实施形態に係る情報処理装置として機能するコンピュータの構成を示すブロック図である。

【発明を実施するための形態】

【0016】

〔例示的实施形態1〕

本発明の第1の例示的实施形態について、図面を参照して詳細に説明する。本例示的实施形態は、後述する例示的实施形態の基本となる形態である。

【0017】

< 情報処理装置1の構成 >

本例示的实施形態に係る情報処理装置1の構成について、図1を参照して説明する。図1は、情報処理装置1の構成を示すブロック図である。情報処理装置1は、レコード対の同一性予測を行う装置である。情報処理装置1は、取得部11、類似度算出部12、予測部13及び出力部14を備える。

30

【0018】

(取得部11)

取得部11は、レコード対を取得する。

【0019】

(レコード対・レコード)

レコード対は複数のレコードのセットである。レコードは、一例として、テーブルの行であり、テーブルの列に対応する1又は複数の属性名及び属性値のセットを含む。レコード対に含まれるレコードの数は2であってもよく、また、3以上であってもよい。レコード対は、一例として、第1のテーブルに含まれるレコードと、第2のテーブルに含まれるレコードとのセットである。第1のテーブル及び第2のテーブルは、一例として、事業者の顧客情報を保存したテーブル、又は、商品情報を保存したテーブルである。ただし、第1のテーブル及び第2のテーブルは上述した例に限られず、他のテーブルであってもよい。また、第1のテーブルと第2のテーブルとは同じであってもよく、また、異なってもよい。

40

【0020】

(類似度算出部12)

類似度算出部12は、取得部11が取得したレコード対について、複数の類似度関数を

50

用いて複数の類似度を算出する。換言すると、類似度算出部 1 2 は、 $k$  個 ( $k$  は 2 以上の整数) の類似度関数  $f_i$  ( $1 \leq i \leq k$ ) を用いて、1 つのレコード対について  $k$  個の類似度を算出する。

【0021】

(類似度関数)

類似度関数  $f_i$  は、レコード対に含まれるレコード同士の類似度を算出するための関数である。以下では、類似度関数  $f_i$  を「類似度算出モデル」とも呼ぶ。類似度関数  $f_i$  の入力レコード対であり、類似度関数  $f_i$  の出力はレコード対に含まれるレコード同士の類似度である。複数の類似度関数  $f_i$  は、後述する情報処理装置 2 による学習の対象であり得る。類似度関数  $f_i$  が機械学習により生成される場合、類似度関数  $f_i$  の機械学習の手法は限定されず、一例として、決定木ベース、線形回帰、又はニューラルネットワークの手法が用いられてもよく、また、これらのうちの 2 以上の手法が用いられてもよい。決定木ベースとしては、例えば、LightGBM (Light Gradient Boosting Machine)、ランダムフォレスト、及び XGBoost が挙げられる。線形回帰としては、例えば、ベイズ回帰、サポートベクター回帰、Ridge 回帰、Lasso 回帰、及び Elastic Net が挙げられる。ニューラルネットワークとしては、例えばディープラーニングが挙げられる。

10

【0022】

類似度関数  $f_i$  は、一例として、0 ~ 1 の数値を類似度として出力する。類似度関数  $f_i$  としては、例えば、Jaccard 係数を用いることができる。Jaccard 係数は、集合  $A = \{a_1, a_2, \dots\}$  と集合  $B = \{b_1, b_2, \dots\}$  に対し、 $|A \cap B| / |A \cup B|$  を計算するものである。また、類似度関数  $f_i$  としては例えば、非特許文献 1 に記載された手法が用いられてもよい。また、他の例として、類似度関数  $f_i$  として、例えば文献「Yuliang Li, Jinfeng Li, Yoshihiko Suhara, AnHai Doan, Wang-Chiew Tan, Deep Entity Matching with Pre-Trained Language Models, Proceedings of the VLDB Endowment, 2016」(以下「非特許文献 2」という) に記載された手法が用いられてもよい。ただし、類似度関数  $f_i$  は上述した例に限られず、他の手法によりレコード対の類似度を算出するものであってもよい。

20

【0023】

(予測部 1 3)

予測部 1 3 は、レコード対と、複数の類似度とを参照して、レコード対に応じて定まる重要度を用いてレコード対の同一性予測を行う。

30

【0024】

(重要度)

重要度は、レコード対に応じて定まる情報である。重要度は、一例として、レコード対を参照して算出される。より具体的には、一例として、重要度を算出するための重要度算出モデルを用いて予測部 1 3 が重要度を算出する。この場合、重要度算出モデルの入力はレコード対である。また、重要度算出モデルの出力は重要度である。重要度算出モデルは、後述する情報処理装置 2 による学習の対象であり得る。重要度算出モデルが機械学習により生成される場合、重要度算出モデルの機械学習の手法は限定されず、一例として、決定木ベース、線形回帰、又はニューラルネットワークの手法が用いられてもよく、また、これらのうちの 2 以上の手法が用いられてもよい。

40

【0025】

重要度算出モデルは、一例として、BERT (Bidirectional encoder representations from Transformers)、fastText、word2vec、tf-idf、BM25、等の言語モデルを用いて生成される。また、重要度算出モデルは言語モデルを含んでもよい。言語モデルを用いる場合の重要度の算出処理の具体例について説明する。予測部 1 3 は、一例として、言語モデルを用いてレコード対をベクトルに変換し、このベクトルを更に別の特徴量空間上のベクトルに変換する。更に、予測部 1 3 は、このベクトルを  $k$  クラス分類器 (ソフトマックス関数、等) に入力することで、 $k$  個の重要度を算出する。算出される  $k$

50

個の重要度のそれぞれは、 $k$  個の類似度関数  $i$  のそれぞれに対応する。

【0026】

ただし、重要度を算出する手法は上述した例に限られず、予測部 13 は他の手法により重要度を算出してもよい。予測部 13 は、一例として、ルールベースの処理により重要度を算出してもよい。例えば、予測部 13 は、重要度とレコード対に関する情報とを対応付けたテーブルを参照することにより重要度を算出してもよい。ここで、レコード対に関する情報は、一例として、レコード対に含まれるレコードの特徴量、レコードの分類結果、又はレコードの名称、等を含んでもよい。

【0027】

予測部 13 は、一例として、類似度算出部 12 が算出した複数の類似度に関する線形和であって、各重要度を重み係数とする線形和を用いて、レコード対の同一性予測を行う。ただし、予測部 13 が同一性予測を行う手法は線形和を用いる手法に限られず、予測部 13 は他の手法によりレコード対の同一性予測を行ってもよい。

【0028】

予測部 13 は、一例として、機械学習により生成される予測モデルにレコード対と類似度とを入力することによりレコード対の同一性予測を行ってもよい。この場合、予測モデルの入力は、一例として、 $k$  個の類似度のセットとレコード対とを含む。また、予測モデルの出力は、一例として、同一性の予測結果を含む。また、予測部 13 は、予測モデルが有するパラメータを重要度として算出する。予測モデルの機械学習の手法は限定されず、一例として、決定木ベース、線形回帰、又はニューラルネットワークの手法が用いられてもよく、また、これらのうちの 2 以上の手法が用いられてもよい。

【0029】

(出力部 14)

出力部 14 は、予測部 13 による予測結果を出力する。予測結果は、一例として、レコード対に含まれるレコードが同一であることを示す情報、又は、レコード対に含まれるレコードの類似度を示す情報を含む。

【0030】

予測部 13 による予測結果は、例えばテーブルの統合処理、又は情報検索処理に用いられる。テーブルを統合する場合、予測部 13 により同一であると予測されたレコードを連携することで、複数のテーブルを統合しデータの一元管理を行うことができる。また、情報検索において、検索キーとするレコード（例えば、ユーザにより指定されたレコード）と、所定のテーブルに登録された他の任意のレコードとのレコード対について予測部 13 が同一性予測を行ってもよい。この場合、予測部 13 により同一であると予測されたレコード対に含まれるレコードを、検索結果として情報処理装置 1 が出力してもよい。これにより、検索キーであるレコードと連携されていないテーブルにおける検索処理が可能となる。

【0031】

< 情報処理装置 1 の効果 >

以上のように、本例示的实施形態に係る情報処理装置 1 においては、レコード対について複数の類似度関数を用いて複数の類似度を算出し、レコード対と、複数の類似度とを参照して、レコード対に応じて定まる重要度を用いてレコード対の同一性予測を行う構成が採用されている。ここで、重要度はレコード対に応じて定まるため、複数の類似度に基づく同一性予測の結果は、画一的な手法によるものではなく、レコード対毎の重要度が反映されたものとなる。このため、本例示的实施形態に係る情報処理装置 1 によれば、レコード対の同一性をより好適に予測できるという効果が得られる。

【0032】

< 情報処理方法 S1 の流れ >

本例示的实施形態に係る情報処理方法 S1 の流れについて、図 2 を参照して説明する。図 2 は、情報処理方法 S1 の流れを示すフロー図である。ステップ S11 において、取得部 11 はレコード対を取得する。ステップ S12 において、類似度算出部 12 は、レコー

10

20

30

40

50

ド対について、複数の類似度関数を用いて複数の類似度を算出する。ステップ S 1 3 において、予測部 1 3 は、レコード対と複数の類似度とを参照して、レコード対に応じて定まる重要度を用いてレコード対の同一性予測を行う。ステップ S 1 4 において、出力部 1 4 は予測部 1 3 による予測結果を出力する。

#### 【 0 0 3 3 】

< 情報処理方法 S 1 の効果 >

以上のように、本例示的实施形態に係る情報処理方法 S 1 においては、レコード対について複数の類似度関数を用いて複数の類似度を算出し、レコード対と、複数の類似度とを参照して、レコード対に応じて定まる重要度を用いてレコード対の同一性予測を行う構成が採用されている。このため、本例示的实施形態に係る情報処理方法 S 1 によれば、レコード対の同一性をより好適に予測できるという効果が得られる。

10

#### 【 0 0 3 4 】

< 情報処理装置 2 の構成 >

次いで、本例示的实施形態に係る情報処理装置 2 の構成について、図 3 を参照して説明する。図 3 は、情報処理装置 2 の構成を示すブロック図である。情報処理装置 2 は、レコード対の同一性を予測するために用いるパラメータを生成する装置である。情報処理装置 2 は、取得部 2 1 及びパラメータ生成部 2 2 を備える。

#### 【 0 0 3 5 】

取得部 2 1 は、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得する。同一性に関するラベルは、一例として、レコード対に含まれるレコードが同一であるか否かを示す。

20

#### 【 0 0 3 6 】

パラメータ生成部 2 2 は、( i ) 予測対象のレコード対について複数の類似度を算出するための複数の類似度関数  $i$  の各々が有する 1 又は複数のパラメータ、及び ( i i ) 予測対象のレコード対と、複数の類似度とを参照して、予測対象のレコード対に応じて定まる重要度を用いて予測対象のレコード対の同一性予測を行う予測部 1 3 が、重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータ、の少なくとも何れかのパラメータを、訓練データを参照して生成する。

#### 【 0 0 3 7 】

< 情報処理装置 2 の効果 >

以上のように、本例示的实施形態に係る情報処理装置 2 においては、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得し、予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する 1 又は複数のパラメータ、及び予測対象のレコード対と、複数の類似度とを参照して、予測対象のレコード対に応じて定まる重要度を用いて予測対象のレコード対の同一性予測を行う予測手段が、重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータの少なくとも何れかのパラメータを、訓練データを参照して生成する構成が採用されている。このため、本例示的实施形態に係る情報処理装置 2 によれば、レコード対の同一性をより好適に予測可能なパラメータを生成できるという効果が得られる。

30

#### 【 0 0 3 8 】

< 情報処理方法 S 2 の流れ >

本例示的实施形態に係る情報処理方法 S 2 の流れについて、図 4 を参照して説明する。図 4 は、情報処理方法 S 2 の流れを示すフロー図である。ステップ S 2 1 において、取得部 2 1 は、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得する。ステップ S 2 2 において、パラメータ生成部 2 2 は、( i ) 予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する 1 又は複数のパラメータ、及び ( i i ) 予測対象のレコード対と、複数の類似度とを参照して、予測対象のレコード対に応じて定まる重要度を用いて予測対象のレコード対の同一性予測を行う予測手段が、重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータ、の少なくとも何れかのパラメータを、訓練データを参照して生成する

40

50

## 【 0 0 3 9 】

## &lt; 情報処理方法 S 2 の効果 &gt;

以上のように、本例示的实施形態に係る情報処理方法 S 2 においては、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得し、予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する 1 又は複数のパラメータ、及び予測対象のレコード対と、複数の類似度とを参照して、予測対象のレコード対に応じて定まる重要度を用いて予測対象のレコード対の同一性予測を行う予測手段が、重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータの少なくとも何れかのパラメータを、訓練データを参照して生成する構成が採用されている。このため、本例示的实施形態に係る情報処理方法 S 2 によれば、レコード対の同一性をより好適に予測可能なパラメータを生成できるという効果が得られる。

10

## 【 0 0 4 0 】

## &lt; 製造方法 &gt;

情報処理装置 2 は、学習済モデルの製造方法を実行する装置として特定することもできる。ここで、学習済モデルの製造方法は、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得することと、複数の類似度算出モデル及び重要度算出モデルの少なくとも何れかのモデルを、訓練データを参照して生成することと、を含む。

## 【 0 0 4 1 】

## 〔 例示的实施形態 2 〕

本発明の第 2 の例示的实施形態について、図面を参照して詳細に説明する。なお、例示的实施形態 1 にて説明した構成要素と同じ機能を有する構成要素については、同じ符号を付し、その説明を繰り返さない。

20

## 【 0 0 4 2 】

## &lt; 情報処理装置 1 A の構成 &gt;

図 5 は、本例示的实施形態に係る情報処理装置 1 A の構成を示すブロック図である。情報処理装置 1 A は、制御部 1 0 A、記憶部 2 0 A、通信部 3 0 A 及び入出力部 4 0 A を備える。

## 【 0 0 4 3 】

## ( 通信部 3 0 A )

通信部 3 0 A は、情報処理装置 1 A の外部の装置と通信回線を介して通信する。通信回線の具体的構成は本例示的实施形態を限定するものではないが、通信回線は一例として、無線 LAN ( Local Area Network )、有線 LAN、WAN ( Wide Area Network )、公衆回線網、モバイルデータ通信網、又は、これらの組み合わせである。通信部 3 0 A は、制御部 1 0 A から供給されたデータを他の装置に送信したり、他の装置から受信したデータを制御部 1 0 A に供給したりする。

30

## 【 0 0 4 4 】

## ( 入出力部 4 0 A )

入出力部 4 0 A には、キーボード、マウス、ディスプレイ、プリンタ、タッチパネル等の入出力機器が接続される。入出力部 4 0 A は、接続された入力機器から情報処理装置 1 A に対する各種の情報の入力を受け付ける。また、入出力部 4 0 A は、制御部 1 0 A の制御の下、接続された出力機器に各種の情報を出力する。入出力部 4 0 A としては、例えば USB ( Universal Serial Bus ) などのインタフェースが挙げられる。

40

## 【 0 0 4 5 】

## ( 制御部 1 0 A )

制御部 1 0 A は、図 5 に示すように、取得部 1 1、類似度算出部 1 2、予測部 1 3、出力部 1 4、及び統合部 1 5 A を備える。

## 【 0 0 4 6 】

## ( 取得部 1 1 )

50

取得部 1 1 は、レコード対に含まれる第 1 のレコード  $e$  を含む第 1 のデータ  $x$  と、レコード対に含まれる第 2 のレコード  $e'$  を含む第 2 のデータ  $x'$  とを取得する。第 1 のデータ  $x$  及び第 2 のデータ  $x'$  は、例えば複数のレコードを含むテーブルである。第 1 のレコード  $e$   $x$  と、第 2 のレコード  $e'$   $x'$  とは、一例として、以下のように表現される。

$$e = (a_1 : v_1, a_2 : v_2, \dots, a_d : v_d)$$

$$e' = (a_1 : v_1, a_2 : v_2, \dots, a_{d'} : v_{d'})$$

【0047】

ここで、 $a_l \in A_l$  ( $l = 1, 2, \dots, d$ )、及び  $a'_m \in A'_m$  ( $m = 1, 2, \dots, d'$ ) は属性名であり、 $A_l$  及び  $A'_m$  は、例えば文字列空間である。 $v_l \in V_l$  及び  $v'_m \in V'_m$  は属性値であり、 $V_l$  及び  $V'_m$  は、例えば文字列空間又は実数空間である。 $d$  はレコード  $e$  が有する属性の数であり、 $d'$  はレコード  $e'$  が有する属性の数である。換言すると、第 1 のレコード  $e$  及び第 2 のレコード  $e'$  はそれぞれ、属性名と属性値とのセットを複数含む。

【0048】

図 6 は、第 1 のデータ  $x$  及び第 2 のデータ  $x'$  の具体例であるテーブル  $T_1$  及びテーブル  $T_2$  を示す図である。テーブル  $T_1$  及びテーブル  $T_2$  は行と列からなり、行はレコード、列は属性に対応する。換言すると、テーブル  $T_1$  は、複数の第 1 のレコード  $e_1, e_2, \dots$  を含む。また、テーブル  $T_2$  は、複数の第 2 のレコード  $e'_1, e'_2, \dots$  を含む。

【0049】

図 6 の第 1 のレコード  $e_2$  は、 $e_2 = (\text{商品名} : \text{ポテトチップス}, \text{価格} : 198)$  と表される。第 1 のレコード  $e_2$  において、属性名が「商品名」である属性の属性値は「ポテトチップス」であり、また、属性名が「価格」である属性の属性値は「198」である。

【0050】

テーブル  $T_1$  の属性名及び属性値と、テーブル  $T_2$  の属性名及び属性値とは、同じであってもよく、また、異なってもよい。図 6 の例で、取得部 1 1 が取得するレコード対  $(e, e')$  は、テーブル  $T_1$  に含まれる第 1 のレコード  $e_1, e_2, \dots$  のいずれかと、テーブル  $T_2$  に含まれる第 2 のレコード  $e'_1, e'_2, \dots$  のいずれかとの対である。

【0051】

(類似度算出部 1 2)

類似度算出部 1 2 は、 $k$  個 ( $k$  は 2 以上の整数) の類似度関数  $s_i$  ( $1 \leq i \leq k$ ) を用いて、1 つのレコード対  $(e, e')$  について  $k$  個の類似度  $s_i$  を算出する。類似度算出部 1 2 が  $k$  個の類似度  $s_i$  を算出する処理の詳細については後述する。

【0052】

(予測部 1 3)

予測部 1 3 は、レコード対  $(e, e')$  と、複数の類似度  $s_i$  とを参照して、レコード対  $(e, e')$  に応じて定まる重要度を用いてレコード対の同一性予測を行う。本例示の実施形態において、予測部 1 3 は、レコード対  $(e, e')$  を参照して重要度を算出する重要度算出部 1 3 1 A を備えている。予測部 1 3 が行う同一性の予測処理、及び重要度算出部 1 3 1 A が行う重要度の算出処理の詳細については後述する。

【0053】

(出力部 1 4)

出力部 1 4 は、予測部 1 3 による予測結果を出力する。予測結果は、一例として、レコード対に含まれるレコード同士が同一であるか否かを示す情報を含む。また、予測結果は、レコード対に含まれるレコード同士の類似の程度を示す情報を含んでもよい。出力部 1 4 は、予測結果を記憶部 2 0 A 又は外部記憶装置に書き込むことにより出力してもよく、また、入出力部 4 0 A に接続された出力装置 (表示装置、印刷装置、等) に出力してもよい。また、出力部 1 4 は、通信部 3 0 A を介して予測結果を他の装置に送信することにより予測結果を出力してもよい。

【0054】

(統合部 1 5 A)

10

20

30

40

50

統合部 15 A は、出力部 14 が出力する予測結果を参照して、第 1 のデータと第 2 のデータとから、統合済データを生成する。統合部 15 A が行う統合済データの生成処理の詳細については後述する。

## 【0055】

(記憶部 20 A)

記憶部 20 A には、取得部 11 が取得する第 1 のデータ  $x$  及び第 2 のデータ  $x'$  が記憶されるとともに、予測部 13 の予測結果  $PR$  が記憶される。また、記憶部 20 A には、複数の類似度関数  $i$ 、重要度算出モデル  $g$ 、及びパラメータ  $P$  が記憶される。

## 【0056】

類似度関数  $\{i_1, \dots, i_k\}$  は、上述の例示的实施形態 1 で示したように、例えば Jaccard 係数、又は非特許文献 1 若しくは非特許文献 2 に記載された手法により類似度を算出する関数である。類似度関数  $i_i$  は、一例として、情報処理装置 1 A のユーザ等により入力される。類似度関数  $i_i$  は、一例として、レコード対  $(e, e')$  に対して 0 から 1 の数値を類似度として出力する。この場合、例えば、出力値が 1 に近いほど類似性が高く、0 に近いほど類似性が低い。類似度関数  $i_i$  は、一例として、学習可能なパラメータを備えた関数である。

10

## 【0057】

重要度算出モデル  $g$  は、重要度算出部 131 A が重要度を算出するために用いるモデルである。重要度算出モデル  $g$  は、上述の例示的实施形態 1 で示したように、例えば BERT、fastText、word2vec、tf-idf、BM25、等の言語モデルを用いて生成される。また、重要度算出モデル  $g$  は言語モデルを含んでもよい。

20

## 【0058】

記憶部 20 A に記憶されたパラメータ  $P$  は、 $k$  個の類似度関数  $i_i$  の各々が有する 1 又は複数のパラメータ  $i_i$ 、及び、重要度算出モデル  $g$  が有する 1 又は複数のパラメータ  $w$ 、の少なくとも何れかのパラメータを含む。

## 【0059】

< 情報処理方法 S1 A の流れ >

図 7 は、情報処理装置 1 A が実行する情報処理方法の一例である情報処理方法 S1 A の流れを示すフロー図である。なお、一部のステップは並行して、又は順序を換えて実行されてもよい。また、既に説明した内容についてはその説明を繰り返さない。

30

## 【0060】

(ステップ S101)

ステップ S101 において、取得部 11 は、第 1 のデータ及び第 2 のデータを取得する。取得部 11 は、一例として、情報処理装置 1 A のユーザ等が入出力部 40 A に接続された入力装置を用いて入力した第 1 のデータ及び第 2 のデータを取得する。また、取得部 11 は、通信部 30 A を介して他の装置から第 1 のデータ及び第 2 のデータを受信することにより、第 1 のデータ及び第 2 のデータを取得してもよい。また、取得部 11 は、外部接続された記憶装置から第 1 のデータ及び第 2 のデータを読み出すことにより、第 1 のデータ及び第 2 のデータを取得してもよい。取得部 11 は、取得した第 1 のデータ及び第 2 のデータを記憶部 20 A に記憶する。

40

## 【0061】

(ステップ S102)

ステップ S102 において、取得部 11 は、記憶部 20 A に記憶されたパラメータ  $P$  を取得する。

## 【0062】

(ステップ S103)

ステップ S103 において、取得部 11 は、予測対象であるレコード対  $(e, e')$  を取得する。

## 【0063】

(ステップ S104)

50

ステップ S 1 0 4 において、類似度算出部 1 2 は、k 個の類似度関数  $s_i$  を用いて、レコード対 ( e , e ' ) について k 個の類似度  $s_i$  を算出する。k 個の類似度関数  $s_i$  がそれぞれ異なるため、算出される k 個の類似度  $s_i$  もそれぞれ異なった値となり得る。例えば、「アイス」と「あいす」のレコード対の場合、表記変更を行って算出される類似度  $s_i$  は、類似性が高いことを示す値となる一方、部分文字列を抽出して算出される類似度  $s_i$  は、類似性が低いことを示す値となる。また、「ポテトチップス」と「ポテチ」のレコード対の場合、表記変更を行って算出される類似度  $s_i$  は、類似性が低いことを示す値となる一方、部分文字列を抽出して算出される類似度  $s_i$  は、類似性が高いことを示す値となる。

【 0 0 6 4 】

(ステップ S 1 0 5 )

ステップ S 1 0 5 において、重要度算出部 1 3 1 A は、レコード対 ( e , e ' ) を参照して、複数の類似度  $s_i$  の各々に関する重要度  $g_i$  を算出する。重要度算出部 1 3 1 A は、一例として、重要度算出モデル  $g$  を用いて重要度  $g_i$  を算出する。

【 0 0 6 5 】

重要度算出モデル  $g$  は、複数の類似度  $s_i$  のそれぞれについて重要度  $g_i$  を算出するためのモデルである。重要度算出モデル  $g$  は、一例として、

【数 1】

$$g: X \times X' \rightarrow [0,1]^k, \sum_{i=1}^k \{g(e, e')\}_i = 1$$

と表される。換言すると、重要度算出モデル  $g$  により算出される k 個の重要度  $\{g(e, e')\}_i$  の総和は 1 である。

【 0 0 6 6 】

重要度算出部 1 3 1 A が行う重要度  $g_i$  の算出処理の具体例について説明する。まず、重要度算出部 1 3 1 A は、言語モデルで第 1 のレコード e 及び第 2 のレコード e ' の各属性値の文字列をベクトルに変換する。具体的には、例えば、重要度算出部 1 3 1 A は、レコード対 ( e = ( 商品名 : ポテトチップス , 価格 : 1 9 8 ) , e ' = ( 商品名 : ポテチ , 評価 : 5 ) ) を、レコード対 ( e , e ' ) を文字列に変換する関数  $serialize(e, e')$  により、「[CLS][COL]商品名[VAL]ポテトチップス[COL]価格[VAL]198[SEP][COL]商品名[VAL]ポテチ[COL]評価[VAL]5[SEP]」という文字列に変換する。ここで、[CLS]と[COL]と[VAL]と[SEP]は、それぞれ文章の始まりと、属性名と、属性値と、レコードの区切りを示す記号である。

【 0 0 6 7 】

更に、重要度算出部 1 3 1 A は、生成した文字列を言語モデル (例えば、BERT) によりベクトルに変換する。続いて、重要度算出部 1 3 1 A は、言語モデルにより得られたベクトルに対し、連結、和、深層学習等を適用することで、新たな L 次元ベクトル z に変換する。

【 0 0 6 8 】

更に、重要度算出部 1 3 1 A は、変換した L 次元ベクトル z を、k クラス分類器に入力することで、k 個の重要度  $\{g(e, e')\}_i$  を算出する。k クラス分類器としては、例えば線形分類器、深層学習等の技術が用いられる。k クラス分類器として、例えば文献「Robert A. Jacobs, Michael Jordan, Geoffrey Hinton: Adaptive Mixtures of Local Experts, Neural Computation 3, 79-87 (1991)」の文献に記載された技術、又は、「Noam Shazeer, Quoc Le, Geoffrey Hinton: Jeffrey Dean: OUTRAGEOUSLY LARGE NEURAL NETWORKS: THE SPARSELY-GATED MIXTURE-OF-EXPERTS LAYER, ICLR 2017」の文献に記載された技術が用いられてもよい。

【 0 0 6 9 】

10

20

30

40

50

例えば、 $i = 1, \dots, k$ において、 $L$ 次元ベクトル $w_i$ に対し、線形ソフトマックス関数の $i$ 次元の出力である重要度 $\{g(e, e')\}_i$ は、  
 $\exp(w_i^T \cdot z) / (\exp(w_1^T \cdot z) + \exp(w_2^T \cdot z) + \dots + \exp(w_k^T \cdot z))$

により算出される。ここで、 $L$ 次元ベクトル $w_i$ は、重要度算出モデル $g$ の学習可能なパラメータ $w$ の一例である。また、「 $w_i^T \cdot z$ 」は $L$ 次元ベクトル $w_i$ と $L$ 次元ベクトル $z$ の内積である。

【0070】

(ステップS106)

ステップS106において、予測部13は、類似度算出部12が計算した類似度 $s_i$ とレコード対 $(e, e')$ とを用いて、レコード対 $(e, e')$ の同一性を予測する。予測部13は、一例として、 $k$ 個の類似度 $s_i$ を用いてレコード対 $(e, e')$ に含まれるレコード同士の類似度を算出し、算出した類似度が閾値 $q$ (例えば、 $q = 0.5$ )より大きい場合に、レコード $e$ とレコード $e'$ とが同一であると予測し、算出した確率が閾値 $q$ 以下である場合に同一でないと予測する。

10

【0071】

予測部13が算出する確率は、レコード対 $(e, e')$ について $k$ 個の類似度 $s_i$ を統合し予測した結果を示すものであり、一例として、 $0 \sim 1$ の数値である。本例示的实施形態において、予測部13は、レコード対 $(e, e')$ と類似度 $s_i$ とを入力とする確率関数 $h$ により、確率を算出する。確率関数 $h$ は、一例として、 $k$ 個の類似度 $s_i = \phi_i(e, e')$

20

【数2】

$$h(e, e' | \{\phi_i\}_{i=1}^k) = \sum_{i=1}^k \{g(e, e')\}_i \phi_i(e, e') \quad \dots \text{(数式1)}$$

【0072】

上述の(数式1)において、重要度 $\{g(e, e')\}_i$ は重要度算出部131Aが算出する重要度であり、類似度 $s_i = \phi_i(e, e')$ は、類似度関数 $\phi_i$ によりレコード対 $(e, e')$ について算出された類似度である。(数式1)を用いる場合、換言すると、予測部13は、複数の類似度 $s_i$ に関する線形和であって、各重要度 $\{g(e, e')\}_i$ を重み係数とする線形和を用いて、同一性予測を行う。

30

【0073】

本例示的实施形態では、異なる複数のレコード対 $(e, e')$ のそれぞれについて算出された $k$ 個の類似度 $s_i$ が同じであっても、重要度 $\{g(e, e')\}_i$ はレコード対のそれぞれで異なり得る。換言すると、予測部13による予測結果には、類似度 $s_i$ だけでなく、レコード対により定まる重要度 $g_i$ が反映される。このように、予測部13が同一性を予測する手法はレコード対によって異なり得る。

【0074】

(ステップS107)

ステップS107において、出力部14は、予測部13の予測結果を出力する。一例として、出力部14は、予測結果を記憶部20Aに記憶する。

40

【0075】

(ステップS108)

ステップS108において、予測部13は、予測対象である全てのレコード対 $(e, e')$ について同一性の予測を行ったかを判定する。予測対象である全てのレコード対 $(e, e')$ について予測処理が完了した場合(ステップS108; YES)、予測部13はステップS109の処理に進む。一方、また予測対象であるレコード対 $(e, e')$ が残っている場合(ステップS108; NO)、予測部13はステップS103の処理に戻り

50

、次のレコード対  $(e, e')$  について同一性の予測を行う。すなわち、情報処理装置 1 A は、予測対象である全てのレコード対  $(e, e')$  について、ステップ S 1 0 3 ~ S 1 0 7 処理を実行する。

【0076】

(ステップ S 1 0 9)

ステップ S 1 0 9 において、統合部 1 5 A は、出力部 1 4 が出力する予測結果を参照して、第 1 のデータと第 2 のデータとから、統合済データを生成する。統合済データは、一例として、統合部 1 5 A は、予測部 1 3 が同一であると予測したレコード対に含まれるレコード同士を統合したレコードを含む。

【0077】

図 8 は、統合済データの一例であるテーブル T 3 を示す図である。テーブル T 3 は、複数のレコード  $f_1, f_2, \dots$  を含む。レコード  $f_1$  は、図 6 の第 1 のレコード  $e_1$  と第 2 のレコード  $e'_2$  を統合したレコードである。レコード  $f_2$  は、図 6 の第 1 のレコード  $e_2$  と第 2 のレコード  $e'_3$  を統合したレコードである。レコード  $f_3$  は、図 6 の第 1 のレコード  $e_3$  と第 2 のレコード  $e'_1$  を統合したレコードである。

【0078】

<実施例>

次に、本例示の実施形態の具体的な実施例を説明する。この例で、類似度関数  $\{i\}$  として、類似度関数  $1 \sim 3$  を用いる。類似度関数  $1$  は、レコード対の商品名の Jaccard 係数を算出する関数である。類似度関数  $2$  は、レコード対の商品名がひらがなであればカタカナに変換してから Jaccard 係数を算出する関数である。類似度関数  $3$  は、上述の非特許文献 2 に記載された手法により類似度を算出する関数である。ここで、類似度関数  $3$  は学習可能なパラメータ  $3$  を持つ。

【0079】

(ステップ S 1 0 1)

図 7 のステップ S 1 0 1 において、取得部 1 1 は、同一性が未知のレコード対の集合であるテストデータ  $D_{test} = \{((商品名: しょうゆせんべい, 価格: 268), (商品名: ショウユセンベイ, 評価: 4)), \dots, ((商品名: ヨモギ団子, 価格: 190), (商品名: みたらしだんご, 評価: 3))\}$  を取得する。

【0080】

(ステップ S 1 0 2 ~ S 1 0 4)

類似度算出部 1 2 は、類似度  $S = (s_1, s_2, s_3)$  を算出する。ここで、類似度算出部 1 2 は、パラメータ  $3$  を記憶部 2 0 A から読み取り、読み取ったパラメータ  $3$  を用いて類似度  $s_3$  を算出する。具体的には、類似度算出部 1 2 は、テストデータ  $D_{test}$  のレコード対  $\{e = (商品名: しょうゆせんべい, 価格: 268), e' = (商品名: ショウユセンベイ, 評価: 4)\}$  の類似度  $S = (\varphi_1(e, e'), \varphi_2(e, e'), \varphi_3(e, e'))^T = (0, 1, 0.7)^T$  を計算する。

【0081】

(ステップ S 1 0 5)

予測部 1 3 は、レコード対  $(e, e')$  の属性名と属性値を連結する関数  $serialize(e, e')$  を用いて、レコード対  $(e, e')$  から文字列「[CLS][COL]商品名[VAL]しょうゆせんべい[COL]価格[VAL]268[SEP][COL]商品名[VAL]ショウユセンベイ[COL]評価[VAL]4[SEP]」を作成する。また、予測部 1 3 は、事前学習済み言語モデルである BERT によりこの文字列のベクトル表現である L 次元ベクトル  $v$  を得る。更に、予測部 1 3 は、線形ソフトマックス関数を用いて  $i = 1, 2, 3$  に対し類似度関数  $i$  の重みである重要度  $g_i$  を、

$$g_i = e^{(w_i^T \cdot v)} / (e^{(w_1^T \cdot v)} + e^{(w_2^T \cdot v)} + e^{(w_3^T \cdot v)})$$

と計算し、 $(g_1, g_2, g_3) = (0.1, 0.6, 0.3)$  を得る。ここで、 $w_1, w_2$  及び  $w_3$  は実数ベクトルであり、重要度算出モデル  $g$  の学習可能なパラメータ  $w$  の一

10

20

30

40

50

例である。

【0082】

(ステップS106)

ステップS106において、予測部13は、類似度算出部12が計算した類似度 $S$ に重要度 $g_i$ をかけた和を確率として算出する。類似度 $S = (0, 1, 0.7)^T$ であり、重要度 $(g_1, g_2, g_3) = (0.1, 0.6, 0.3)$ であるから、

$$h(e, e') = 0.1 \times 0 + 0.6 \times 1 + 0.3 \times 0.7 = 0.81$$

となる。算出された値「0.81」が予め定められた閾値 $q = 0.5$ よりも大きいため、予測部13は、レコード対 $(e, e')$ に含まれるレコード $e$ とレコード $e'$ とが同一であると予測する。

10

【0083】

(ステップS107)

ステップS107において、出力部14がレコード対 $(e, e')$ の同一性予測結果を出力する。以上の同一性予測と出力をテストデータ $D_{test}$ の全てのレコード対に対して適用する。

【0084】

<情報処理装置1Aの効果>

以上のように、本例示的实施形態に係る情報処理装置1Aにおいては、レコード対 $(e, e')$ を参照して重要度 $g_i$ を算出し、算出した重要度 $g_i$ を用いてレコード対の同一性予測を行う構成が採用されている。このため、本例示的实施形態に係る情報処理装置1Aによれば、例示的实施形態1に係る情報処理装置1の奏する効果に加えて、レコード対 $(e, e')$ を用いて算出される重要度 $g_i$ を加味した同一性予測を行うことができ、レコード対 $(e, e')$ の同一性をより適切に予測できるという効果が得られる。

20

【0085】

<変形例>

上述の例示的实施形態において、取得部11は補助データ $u$ を更に取得し、予測部13は、レコード対 $(e, e')$ と、複数の類似度 $s_i$ と、補助データ $u$ とを参照して、レコード対 $(e, e')$ と補助データ $u$ とに応じて定まる重要度 $g_i$ を用いてレコード対 $(e, e')$ の同一性予測を行ってもよい。

【0086】

補助データ $u$ は、一例として、レコードの名前、レコードの特徴量、及び/又はレコードの分類結果(お菓子、人名、等)、を示す情報を含む。ここで、補助データ $u$ は、一例として、Wikipedia(登録商標)等の外部データから得られるレコードに関する情報を含んでもよい。また、補助データ $u$ は、一例として、類似度関数 $s_i$ のパラメータ及び/又は重要度算出モデル $g$ のパラメータ $w$ の学習で用いられた訓練データの数を含んでもよい。ただし、補助データ $u$ は上述した例に限られず、他の情報を含んでもよい。補助データ $u$ は、一例として、離散的な情報を表すワンホットベクトルである。

30

【0087】

この場合、重要度算出モデル $g$ には、レコード対 $(e, e')$ に加えて補助データ $u$ が入力される。一例として、ベクトルである補助データ $u$ は、上述の $L$ 次元ベクトル $z$ に連結され、連結されたベクトルとパラメータ $w$ を用いて重要度 $g_i$ が算出される。

40

【0088】

本変形例では、予測部13は、レコード対 $(e, e')$ と、複数の類似度 $s_i$ と、補助データ $u$ とを参照して、レコード対 $(e, e')$ と補助データ $u$ とに応じて定まる重要度 $g_i$ を用いてレコード対 $(e, e')$ の同一性予測を行う。これにより、予測部13はレコード対 $(e, e')$ の同一性の予測精度をより高くすることができる。

【0089】

[例示的实施形態3]

本発明の第3の例示的实施形態について、図面を参照して詳細に説明する。なお、例示的实施形態1~2にて説明した構成要素と同じ機能を有する構成要素については、同じ符

50

号を付記し、その説明を繰り返さない。

【0090】

< 情報処理装置 1 B の構成 >

図 9 は、本例示の実施形態に係る情報処理装置 1 B の構成を示すブロック図である。情報処理装置 1 B の制御部 1 0 A は、取得部 1 1、類似度算出部 1 2、予測部 1 3、出力部 1 4、統合部 1 5 A に加えて、学習部 1 6 B を備える。

【0091】

本例示の実施形態に係る取得部 1 1 は、レコード対  $(e_j, e'_j)$  と、当該レコード対  $(e_j, e'_j)$  の同一性に関するラベル  $y_j$  との組を複数含む訓練データ  $D_{tr}$  を更  
10  
に取得する。訓練データ  $D_{tr}$  は、上述のパラメータ  $P$  を学習するために用いられる。訓練データ  $D_{tr}$  は、一例として、

【数 3】

$$D_{tr} = \{(e_j, e'_j, y_j)\}_{j=1}^n$$

と表現される。ここで、 $n$  は、レコード対  $(e_j, e'_j)$  の総数である。ラベル  $y_j$  は、一例として、「0」又は「1」である。「1」は、第 1 のレコード  $e_j$  と第 2 のレコード  $e'_j$  とが同一である旨を示し、「0」は、第 1 のレコード  $e_j$  と第 2 のレコード  $e'_j$  とが同一でない旨を示す。  
20

【0092】

学習部 1 6 B は、( i ) 類似度算出部 1 2 が類似度  $s_i$  を算出するために用いる複数の類似度関数  $f_i$  の各々が有する 1 又は複数のパラメータ  $\theta_i$ 、及び ( i i ) 重要度算出部 1 3 1 A が重要度を算出するために用いる重要度算出モデル  $g$  が有する 1 又は複数のパラメータ  $w$ 、の少なくとも何れかのパラメータ  $P$  を、前記訓練データを参照して生成する。学習部 1 6 B は、本明細書に係る「パラメータ生成手段」の一例である。

【0093】

< 情報処理方法 S 2 B の流れ >

図 1 0 は、情報処理装置 1 B が実行する情報処理方法の一例である情報処理方法 S 2 B の流れを示すフロー図である。なお、一部のステップは並行して、又は順序を換えて実行  
30  
されてもよい。また、既に説明した内容についてはその説明を繰り返さない。

【0094】

( ステップ S 2 0 1 ・ S 2 0 2 )

ステップ S 2 0 1 において、取得部 1 1 は、訓練データ  $D_{tr}$  を取得する。訓練データ  $D_{tr}$  は、一例として、情報処理装置 1 B のユーザにより入力される。また、ステップ S 2 0 2 において、取得部 1 1 は、複数の類似度関数  $f_i$  を取得する。類似度関数  $f_i$  は、一例として、情報処理装置 1 B のユーザにより入力される。

【0095】

( ステップ S 2 0 3 )

ステップ S 2 0 3 において、学習部 1 6 B は、訓練データ  $D_{tr}$  を用いて、パラメータ  $\theta_i$  及びパラメータ  $w$  の少なくとも何れかを学習する。ここで、パラメータ  $\theta_i$  は、類似度関数  $f_i$  が有するパラメータの集合である。また、パラメータ  $w$  は、重要度算出モデル  $g$  が有するパラメータの集合である。  
40

【0096】

学習部 1 6 B は、一例として、目的関数  $L$  によりパラメータ  $\theta_i$  とパラメータ  $w$  とを最適化する。この最適化は、一例として、

【数 4】

$$\min_{w, \{\theta_i\}_{i=1}^k} \frac{1}{n} \sum_{j=1}^n l(h_w(e_j, e'_j | \{\phi_{i,\theta_i}\}_{i=1}^k), y_j) + \alpha \Omega(w, \{\theta_i\}_{i=1}^k)$$

と表される。ここで、評価指標  $l$  は、

【数 5】

$$l: [0,1] \times \{0,1\} \rightarrow [0, \infty]$$

10

である。すなわち、評価指標  $l$  は、

訓練データ  $D_{tr}$  のレコード対  $(e_j, e'_j)$  に含まれるレコード同士が同一である確率（確率関数  $h_w$  の出力）と、

「0」又は「1」のラベル  $y_j$  と、

を入力とし、0以上の値を出力する損失関数である。評価指標  $l$  としては、例えばクロスエントロピー誤差を用いることができる。

【0097】

また、目的関数  $L$  において、 $\alpha$  は非負値のハイパーパラメータである。ハイパーパラメータ  $\alpha$  は、情報処理装置 1 B のユーザ等が定めてもよいし、訓練データ  $D_{tr}$  とは別の同一性が既知のレコード対の集合を用いて自動的に決定された値であってもよい。 $\alpha$  はパラメータに対する正則化項であり、 $L_2$  ノルムを用いてもよい。上の式においてパラメータ  $\theta_i$  を固定してパラメータ  $w$  のみを最適化してもよい。

20

【0098】

学習部 16 B は、生成したパラメータ  $w$  及びパラメータ  $\theta_i$  を記憶部 20 A に保存する。学習部 16 B が生成したパラメータ  $w$  及びパラメータ  $\theta_i$  は、類似度算出部 12 による類似度  $s_i$  の算出処理、及び / 又は予測部 13 による同一性の予測処理において用いられる。

【0099】

< 実施例 >

30

次に、本例示の実施形態の具体的な実施例について説明する。例えば、テーブル T1 の第1のレコード  $e_1 = (\text{商品名: ポテトチップス, 価格: } 198)$ 、第1のレコード  $e_2 = (\text{商品名: アイス, 価格: } 148)$  と、テーブル T2 の第2のレコード  $e'_1 = (\text{商品名: ポテチ, 評価: } 5)$ 、第2のレコード  $e'_2 = (\text{商品名: あいす, 評価: } 4)$  について、訓練データ  $D_{tr}$  を、

$$D_{tr} = \{ (e_1, e'_1, 1), (e_2, e'_2, 1), (e_1, e'_2, 0), (e_2, e'_1, 0) \}$$

とする。

【0100】

また、類似度関数  $\{s_i\}$  として、類似度関数  $s_1 \sim s_3$  を用いる。類似度関数  $s_1 \sim s_3$  は、上述の例示の実施形態 1 の実施例で示した類似度関数  $s_1 \sim s_3$  と同様である。類似度関数  $s_3$  は学習可能なパラメータ  $\theta_3$  を有する。

40

【0101】

ステップ S201 において、取得部 11 は訓練データ  $D_{tr}$  を取得する。また、ステップ S203 において、学習部 16 B は、予測部 13 による訓練データ  $D_{tr}$  のレコード対  $(e_j, e'_j)$  の同一性予測がよく正解するように、クロスエントロピー誤差に基づいて、重要度算出モデル  $g$  のパラメータ  $w$  と類似度関数  $s_i$  のパラメータ  $\theta_i$  を、確率的勾配降下法を用いて最適化する。最適化されたパラメータ  $w$  とパラメータ  $\theta_i$  とは、記憶部 20 A に保存される。

【0102】

50

< 情報処理装置 1 B の効果 >

以上のように、本例示的实施形態に係る情報処理装置 1 B においては、重要度算出モデル  $g$  が有するパラメータ  $w$  及び類似度関数  $i$  が有するパラメータ  $i$  の少なくとも何れかのパラメータを、訓練データ  $D_{tr}$  を参照して生成する構成が採用されている。このため、本例示的实施形態に係る情報処理装置 1 B によれば、例示的实施形態 1 に係る情報処理装置 1 の奏する効果に加えて、レコード対の同一性をより好適に予測可能なパラメータを生成できるという効果が得られる。

【 0 1 0 3 】

< 変形例 >

上述の例示的实施形態において、訓練データ  $D_{tr}$  は、補助データ  $u$  を含んでいてもよい。この場合、訓練データ  $D_{tr}$  は、一例として、

【数 6】

$$D_{tr} = \{(e_j, e'_j, u_j, y_j)\}_{j=1}^n$$

と表される。学習部 1 6 B は、補助データ  $u$  を含む訓練データ  $D_{tr}$  を用いてパラメータ  $w$  とパラメータ  $i$  とを最適化する。

【 0 1 0 4 】

〔例示的实施形態 4〕

本発明の第 4 の例示的实施形態について、図面を参照して詳細に説明する。なお、例示的实施形態 1 ~ 3 にて説明した構成要素と同じ機能を有する構成要素については、同じ符号を付記し、その説明を繰り返さない。

【 0 1 0 5 】

< 情報処理装置 1 C の構成 >

図 1 1 は、本例示的实施形態に係る情報処理装置 1 C の構成を示すブロック図である。情報処理装置 1 C の制御部 1 0 A は、取得部 1 1、類似度算出部 1 2、予測部 1 3、出力部 1 4、学習部 1 6 B に加えて、検索結果出力部 1 7 C を備える。

【 0 1 0 6 】

本例示的实施形態に係る取得部 1 1 は、レコード対  $(e, e')$  に含まれる第 1 のレコード  $e$  として、ユーザからの入力データを取得する。ユーザからの入力データは、一例として、入出力部 4 0 A に接続された入力装置（例えば、キーボード、マウス、等）により入力される。

【 0 1 0 7 】

また、取得部 1 1 は、レコード対  $(e, e')$  に含まれる第 2 のレコード  $e'$  として、対象データに含まれる複数のレコードの 1 つを取得する。対象データは、検索対象のデータであり、一例として、1 又は複数のテーブルを含む。

【 0 1 0 8 】

予測部 1 3 は、第 1 のレコード  $e$  と、対象データに含まれる複数のレコードの各々のレコード対に対して同一性予測を行う。検索結果出力部 1 7 C は、出力部 1 4 が出力する各々の予測結果  $PR$  を参照して、入力データに基づく検索結果であって、対象データを検索対象とする検索結果を出力する。検索結果出力部 1 7 C は、一例として、入出力部 4 0 A に接続された出力装置（ディスプレイ、プリンタ、等）に検索結果を出力する。また、検索結果出力部 1 7 C は、通信部 3 0 A を介して接続された他の装置に検索結果を送信することにより、検索結果を出力してもよい。また、検索結果出力部 1 7 C は、検索結果を記憶部 2 0 A 又は外部記憶装置に記憶することにより検索結果を出力してもよい。

【 0 1 0 9 】

図 1 2 は、検索結果出力部 1 7 C が出力する画面表示の具体例を示す図である。図 1 2 の例で、入力データは、ユーザがテキストボックス 5 1 に入力する文字列であり、対象データは、上述の例示的实施形態 1 において図 6 に示したテーブル  $T_1$  及びテーブル  $T_2$  で

ある。予測部 13 は、ユーザの入力データである第 1 のレコード e と、テーブル T 1 に含まれるレコード及びテーブル T 2 に含まれるレコード e' の各々のレコード対に対して同一性予測を行う。予測部 13 が行う同一性の予測処理は、上述の例示的实施形態 2 で説明したため、その説明を繰り返さない。

#### 【0110】

図 12 の例において、検索結果出力部 17C は、予測部 13 の予測結果 PR を参照して、入力データに基づく検索結果 53、及び検索結果 54 を出力する。検索結果 53 は、「ポテチ」の文字列を入力データとして、テーブル T 1 から検索された検索結果である。検索結果 54 は、「ポテチ」の文字列を入力データとして、テーブル T 2 から検索された検索結果である。

10

#### 【0111】

< 情報処理装置 1C の効果 >

以上のように、本例示的实施形態に係る情報処理装置 1C においては、出力部 14 が出力する各々の予測結果を参照して、入力データに基づく検索結果であって、対象データを検索対象とする検索結果を出力する構成が採用されている。このため、本例示的实施形態に係る情報処理装置 1C によれば、例示的实施形態 1 に係る情報処理装置 1 の奏する効果に加えて、入力データに基づく対象データからの検索をより好適に行うことができるという効果が得られる。

#### 【0112】

情報処理装置 1C は、以下のようにも記載され得る。

20

ユーザからの入力データと、対象データに含まれる複数のレコードの 1 つとをレコード対として取得する取得手段と、

前記レコード対について、複数の類似度関数を用いて複数の類似度を算出する類似度算出手段と、

前記入力データと、前記対象データに含まれる複数のレコードの各々のレコード対に対して、前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行う予測手段と、

前記予測手段による予測結果を参照して、前記入力データに基づく検索結果であって、前記対象データを検索対象とする検索結果を出力する出力手段と、  
を備えている情報処理装置。

30

#### 【0113】

〔ソフトウェアによる実現例〕

情報処理装置 1、1A、1B、1C、2（以下「情報処理装置 1 等」という）の一部又は全部の機能は、集積回路（IC チップ）等のハードウェアによって実現してもよいし、ソフトウェアによって実現してもよい。

#### 【0114】

後者の場合、情報処理装置 1 等は、例えば、各機能を実現するソフトウェアであるプログラムの命令を実行するコンピュータによって実現される。このようなコンピュータの一例（以下、コンピュータ C と記載する）を図 13 に示す。コンピュータ C は、少なくとも 1 つのプロセッサ C1 と、少なくとも 1 つのメモリ C2 と、を備えている。メモリ C2 には、コンピュータ C を情報処理装置 1 等として動作させるためのプログラム P が記録されている。コンピュータ C において、プロセッサ C1 は、プログラム P をメモリ C2 から読み取って実行することにより、情報処理装置 1 等の各機能が実現される。

40

#### 【0115】

プロセッサ C1 としては、例えば、CPU（Central Processing Unit）、GPU（Graphic Processing Unit）、DSP（Digital Signal Processor）、MPU（Micro Processing Unit）、FPU（Floating point number Processing Unit）、PPU（Physics Processing Unit）、マイクロコントローラ、又は、これらの組み合わせなどを用いることができる。メモリ C2 としては、例えば、フラッシュメモリ、HDD（Hard Disk Drive）、SSD（Solid State Drive）、又は、これらの組

50

み合わせなどを用いることができる。

【0116】

なお、コンピュータCは、プログラムPを実行時に展開したり、各種データを一時的に記憶したりするためのRAM(Random Access Memory)を更に備えていてもよい。また、コンピュータCは、他の装置との間でデータを送受信するための通信インタフェースを更に備えていてもよい。また、コンピュータCは、キーボードやマウス、ディスプレイやプリンタなどの入出力機器を接続するための入出力インタフェースを更に備えていてもよい。

【0117】

また、プログラムPは、コンピュータCが読み取り可能な、一時的でない有形の記録媒体Mに記録することができる。このような記録媒体Mとしては、例えば、テープ、ディスク、カード、半導体メモリ、又はプログラマブルな論理回路などを用いることができる。コンピュータCは、このような記録媒体Mを介してプログラムPを取得することができる。また、プログラムPは、伝送媒体を介して伝送することができる。このような伝送媒体としては、例えば、通信ネットワーク、又は放送波などを用いることができる。コンピュータCは、このような伝送媒体を介してプログラムPを取得することもできる。

【0118】

〔付記事項1〕

本発明は、上述した実施形態に限定されるものでなく、請求項に示した範囲で種々の変更が可能である。例えば、上述した実施形態に開示された技術的手段を適宜組み合わせ得られる実施形態についても、本発明の技術的範囲に含まれる。

【0119】

〔付記事項2〕

上述した実施形態の一部又は全部は、以下のようにも記載され得る。ただし、本発明は、以下の記載する態様に限定されるものではない。

【0120】

(付記1)

レコード対を取得する取得手段と、  
前記レコード対について、複数の類似度関数を用いて複数の類似度を算出する類似度算出手段と、  
前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行う予測手段と、  
前記予測手段による予測結果を出力する出力手段と、  
を備えている情報処理装置。

【0121】

上記の構成によれば、レコード対の同一性をより好適に予測できる。

【0122】

(付記2)

前記取得手段は、補助データを更に取得し、  
前記予測手段は、前記レコード対と、前記複数の類似度と、前記補助データとを参照して、前記レコード対と前記補助データとに応じて定まる重要度を用いて前記レコード対の同一性予測を行う、  
付記1に記載の情報処理装置。

【0123】

上記の構成によれば、重要度はレコード対だけでなく補助データの内容を反映した情報となる。このような重要度を用いてレコード対の同一性を予測することにより、レコード対の同一性の予測精度をより高くすることができる。

【0124】

(付記3)

前記予測手段は、前記レコード対を参照して前記重要度を算出する重要度算出手段を備

10

20

30

40

50

えている、

付記 1 又は 2 に記載の情報処理装置。

【 0 1 2 5 】

上記の構成によれば、レコード対を参照して算出される重要度を用いてレコード対の同一性予測を行うことにより、レコード対の同一性の予測精度をより高くすることができる。

【 0 1 2 6 】

( 付記 4 )

前記重要度算出手段は、前記複数の類似度の各々に関する重要度を算出し、

前記予測手段は、前記複数の類似度に関する線形和であって、前記各重要度を重み係数とする線形和を用いて、前記同一性予測を行う、

付記 3 に記載の情報処理装置。

【 0 1 2 7 】

上記の構成によれば、重要度を重み係数とする類似度の線形和を用いて同一性予測を行うことにより、レコード対の同一性の予測精度を高くすることができる。

【 0 1 2 8 】

( 付記 5 )

前記取得手段は、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを更に取得し、

当該情報処理装置は、

前記類似度算出手段が前記類似度を算出するために用いる前記複数の類似度関数の各々が有する 1 又は複数のパラメータ、及び、

前記重要度算出手段が前記重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータ、の少なくとも何れかのパラメータを、前記訓練データを参照して生成するパラメータ生成手段を更に備えている、

付記 3 又は 4 に記載の情報処理装置。

【 0 1 2 9 】

上記の構成によれば、訓練データを参照して生成したパラメータを用いることで、レコード対の同一性をより好適に予測することができる。

【 0 1 3 0 】

( 付記 6 )

前記取得手段は、前記レコード対に含まれる第 1 のレコードを含む第 1 のデータと、前記レコード対に含まれる第 2 のレコードを含む第 2 のデータとを取得し、

当該情報処理装置は、前記出力手段が出力する前記予測結果を参照して、前記第 1 のデータと前記第 2 のデータとから、統合済データを生成する統合手段を備えている、  
付記 1 から 5 の何れか 1 つに記載の情報処理装置。

【 0 1 3 1 】

上記の構成によれば、第 1 のデータと第 2 のデータとをより好適に統合することができる。

【 0 1 3 2 】

( 付記 7 )

前記取得手段は、

前記レコード対に含まれる第 1 のレコードとして、ユーザからの入力データを取得し、

前記レコード対に含まれる第 2 のレコードとして、対象データに含まれる複数のレコードの 1 つを取得し、

前記予測手段は、前記第 1 のレコードと、前記対象データに含まれる複数のレコードの各々とのレコード対に対して前記同一性予測を行い、

当該情報処理装置は、前記出力手段が出力する各々の前記予測結果を参照して、前記入

10

20

30

40

50

カデータに基づく検索結果であって、前記対象データを検索対象とする検索結果を出力する検索結果出力手段を備えている、  
付記 1 から 5 の何れか 1 つに記載の情報処理装置。

【 0 1 3 3 】

上記の構成によれば、入力データに基づく対象データからの検索をより好適に行うことができる。

【 0 1 3 4 】

( 付記 8 )

レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得する取得手段と、

予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する 1 又は複数のパラメータ、及び、

前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータ、

の少なくとも何れかのパラメータを、前記訓練データを参照して生成するパラメータ生成手段と、

を備えている情報処理装置。

【 0 1 3 5 】

上記の構成によれば、レコード対の同一性をより好適に予測可能なパラメータを生成できる。

【 0 1 3 6 】

( 付記 9 )

レコード対を取得することと、

前記レコード対について、複数の類似度関数を用いて複数の類似度を算出することと、

前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行うことと、

前記レコード対の同一性予測による予測結果を出力することと、

を含む情報処理方法。

【 0 1 3 7 】

上記の情報処理方法によれば、上述した情報処理装置と同様の効果を奏する。

【 0 1 3 8 】

( 付記 10 )

レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得することと、

予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する 1 又は複数のパラメータ、及び、

前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルが有する 1 又は複数のパラメータ

の少なくとも何れかのパラメータを、前記訓練データを参照して生成することと、

を含む情報処理方法。

【 0 1 3 9 】

上記の情報処理方法によれば、上述した情報処理装置と同様の効果を奏する。

【 0 1 4 0 】

( 付記 11 )

レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得することと、

10

20

30

40

50

予測対象のレコード対について複数の類似度を算出するための複数の類似度算出モデル、及び、

前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデル、の少なくとも何れかのモデルを、前記訓練データを参照して生成することと、を含む学習済モデルの製造方法。

【0141】

上記の構成によれば、レコード対の同一性をより好適に予測可能なモデルを製造することができる。

【0142】

(付記12)

コンピュータに、

レコード対を取得する取得処理と、

前記レコード対について、複数の類似度関数を用いて複数の類似度を算出する類似度算出処理と、

前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行う予測処理と、

前記予測処理による予測結果を出力する出力処理と、

を実行させるプログラム。

【0143】

上記の構成によれば、上述した情報処理装置と同様の効果を奏する。

【0144】

(付記13)

コンピュータに、

レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得する取得処理と、

予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する1又は複数のパラメータ、及び、

前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルが有する1又は複数のパラメータ

の少なくとも何れかのパラメータを、前記訓練データを参照して生成するパラメータ生成処理と、

を実行させるプログラム。

【0145】

上記の構成によれば、上述した情報処理装置と同様の効果を奏する。

【0146】

[付記事項3]

上述した実施形態の一部又は全部は、更に、以下のように表現することもできる。

【0147】

少なくとも1つのプロセッサを備え、前記プロセッサは、レコード対を取得する取得処理と、前記レコード対について、複数の類似度関数を用いて複数の類似度を算出する類似度算出処理と、前記レコード対と、前記複数の類似度とを参照して、前記レコード対に応じて定まる重要度を用いて前記レコード対の同一性予測を行う予測処理と、前記予測処理による予測結果を出力する出力処理とを実行する情報処理装置。

【0148】

なお、この情報処理装置は、更にメモリを備えていてもよく、このメモリには、前記取得処理と、前記類似度算出処理と、前記予測処理と、前記出力処理とを前記プロセッサに

10

20

30

40

50

実行させるためのプログラムが記憶されていてもよい。また、このプログラムは、コンピュータ読み取り可能な一時的でない有形の記録媒体に記録されていてもよい。

【0149】

上述した実施形態の一部又は全部は、更に、以下のように表現することもできる。

少なくとも1つのプロセッサを備え、前記プロセッサは、レコード対と、当該レコード対の同一性に関するラベルとの組を複数含む訓練データを取得する取得処理と、予測対象のレコード対について複数の類似度を算出するための複数の類似度関数の各々が有する1又は複数のパラメータ、及び前記予測対象のレコード対と、前記複数の類似度とを参照して、前記予測対象のレコード対に応じて定まる重要度を用いて前記予測対象のレコード対の同一性予測を行う予測手段が、前記重要度を算出するために用いる重要度算出モデルが有する1又は複数のパラメータの少なくとも何れかのパラメータを、前記訓練データを参照して生成するパラメータ生成処理とを実行する情報処理装置。

10

【0150】

なお、この情報処理装置は、更にメモリを備えていてもよく、このメモリには、前記取得処理と、前記パラメータ生成処理とを前記プロセッサに実行させるためのプログラムが記憶されていてもよい。また、このプログラムは、コンピュータ読み取り可能な一時的でない有形の記録媒体に記録されていてもよい。

【符号の説明】

【0151】

1、1A、1B、1C、2 情報処理装置

20

10A 制御部

11、21 取得部

12 類似度算出部

13 予測部

14 出力部

15A 統合部

16B 学習部

17C 検索結果出力部

20A 記憶部

22 パラメータ生成部

30

30A 通信部

40A 入出力部

131A 重要度算出部

S1、S1A、S2、S2B 情報処理方法

40

50