

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2016-71839

(P2016-71839A)

(43) 公開日 平成28年5月9日(2016.5.9)

(51) Int.Cl.

G06F 21/62 (2013.01)

F I

G06F 21/62 345

テーマコード (参考)

審査請求 有 請求項の数 20 O L (全 18 頁)

(21) 出願番号 特願2015-20104 (P2015-20104)
 (22) 出願日 平成27年2月4日 (2015.2.4)
 (31) 優先権主張番号 103134231
 (32) 優先日 平成26年10月1日 (2014.10.1)
 (33) 優先権主張国 台湾 (TW)

(71) 出願人 502003596
 財団法人 資訊工業策進会
 INSTITUTE FOR INFORMATION INDUSTRY
 台湾台北市和平東路2段106号11樓
 11F, NO. 106, SEC. 2,
 HEPING E. RD., TAIPEI, TAIWAN,

(74) 代理人 100108453

弁理士 村山 靖彦

(74) 代理人 100110364

弁理士 実広 信哉

(74) 代理人 100133400

弁理士 阿部 達彦

最終頁に続く

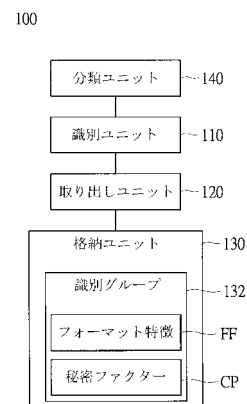
(54) 【発明の名称】 秘密データを識別する方法、電子装置及びコンピュータ読み取り可能な記録媒体

(57) 【要約】 (修正有)

【課題】データの漏れを回避することができる秘密データを識別する方法、電子装置及びコンピュータ読み取り可能な記録媒体を提供する。

【解決手段】電子装置100は、特定フォーマットを表すフォーマット特徴に基づいてデータにおいて特定フォーマットがあるか否かを判定し、特定フォーマットを秘密データとして表す複数の秘密ファクターCPに基づいて、データにおける特定フォーマットが秘密データであるか否かをさらに判定する。これにより、カウント数が多くないが機密記述が大量に含まれたデータの正しい機密レベルを提供するとともに特定フォーマットを有する秘密データを識別する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

特定フォーマットを表すフォーマット特徴と前記特定フォーマットを秘密データとして表す複数の機密ファクターとをそれぞれ有すると共に前記特定フォーマットにそれぞれ対応する複数の識別グループが格納された電子装置に適用される、秘密データを識別する方法であって、

複数のデータのいずれか 1 つを取り出し、それを取り出しデータと定義する工程と、それらのフォーマット特徴のいずれか 1 つを取り出し、それを取り出し特徴と定義する工程と、

前記電子装置が、前記取り出し特徴に基づいて前記取り出しデータが対応する前記特定フォーマットを有するか否かを判定し、前記取り出しデータが対応する前記特定フォーマットを有すると判定した場合に、前記特定フォーマットに対応する複数の秘密ファクターが前記取り出しデータにおける出現頻度が秘密閾値以上であるかを判定し、前記出現頻度が前記秘密閾値以上であると判定した場合に、前記取り出しデータにおける前記特定フォーマットが前記秘密データであることを表し、前記出現頻度が前記秘密閾値よりも小さいと判定した場合に、前記取り出しデータにおける前記特定フォーマットが前記秘密データではないことを表すようにする工程と、

前記電子装置が、複数のフォーマット特徴において取り出されていない前記フォーマット特徴があるか否かを判定し、複数のフォーマット特徴において取り出されていない前記フォーマット特徴があると判定した場合に、取り出されていない前記フォーマット特徴を取り出し、取り出されていない前記フォーマット特徴を前記取り出し特徴と定義することで、改めて前記取り出し特徴に基づいて前記取り出しデータが対応する前記特定フォーマットを有するか否かを判定し、複数のフォーマット特徴において取り出されていない前記フォーマット特徴がないと判定した場合に、複数のデータの次のデータを取り出し、前記次のデータを前記取り出しデータと定義することで、改めて前記取り出しデータが対応する前記特定フォーマットを有するか否かを判定する工程と、

を備えることを特徴とする秘密データを識別する方法。

【請求項 2】

前記電子装置は、前記取り出しデータが対応する前記特定フォーマットを有しないと判定した場合に、それらのフォーマット特徴において取り出されていない前記フォーマット特徴があるか否かを判定することを特徴とする請求項 1 に記載の秘密データを識別する方法。

【請求項 3】

前記電子装置は、それらのフォーマット特徴において取り出されていない前記フォーマット特徴がないと判定した後に、さらに、それらの秘密ファクターとそれらの秘密ファクターがそれらのデータに出現する回数とに基づいて前記取り出しデータに対して分類を行うことを特徴とする請求項 1 に記載の秘密データを識別する方法。

【請求項 4】

前記取り出し特徴に基づいて前記取り出しデータが対応する前記特定フォーマットを有するか否かを判定する工程において、前記取り出し特徴が同一列に 2 つの列終了位置を有し、かつ前記電子装置が前記特定フォーマットにおいて同一列に 2 つの列終了位置を有する数がフォーマット閾値以上であると判定した場合に、前記電子装置は前記取り出しデータが前記特定フォーマットを有すると判定することを特徴とする請求項 1 に記載の秘密データを識別する方法。

【請求項 5】

前記取り出し特徴に基づいて前記取り出しデータが対応する前記特定フォーマットを有するか否かを判定する工程において、前記フォーマット特徴が特定鍵からのメッセージを含み、かつ前記特定フォーマットにおいて前記メッセージを有する数がフォーマット閾値以上である場合に、前記取り出しデータが前記特定フォーマットを有すると判定することを特徴とする請求項 1 に記載の秘密データを識別する方法。

10

20

30

40

50

【請求項 6】

前記取り出し特徴に基づいて前記取り出しデータが対応する前記特定フォーマットを有するか否かを判定する工程において、前記フォーマット特徴がカスタマイズ特徴を含み、かつ前記特定フォーマットにおいて前記カスタマイズ特徴を有する数がフォーマット閾値よりも大きい場合に、前記取り出しデータが前記特定フォーマットを有すると判定することを特徴とする請求項 1 に記載の秘密データを識別する方法。

【請求項 7】

各前記識別グループのそれらの秘密ファクターは、少なくとも 1 つのワード、少なくとも 1 つのストリング、少なくとも 1 つの符号、少なくとも 1 つの数字、少なくとも 1 つの実行指令、及び少なくとも 1 つのフォーマットのいずれか 1 つまたはそれらの組み合わせであることを特徴とする請求項 1 に記載の秘密データを識別する方法。

10

【請求項 8】

各前記フォーマット特徴は、少なくとも 1 つのワード、少なくとも 1 つのストリング、少なくとも 1 つの符号、少なくとも 1 つの数字、少なくとも 1 つの実行指令、及び少なくとも 1 つのフォーマットのいずれか 1 つまたはそれらの組み合わせであることを特徴とする請求項 1 に記載の秘密データを識別する方法。

【請求項 9】

特定フォーマットを表すフォーマット特徴と前記特定フォーマットを秘密データとして表す複数の機密ファクターとをそれぞれ有すると共に前記特定フォーマットにそれぞれ対応する複数の識別グループを格納するための格納ユニットと、

20

前記格納ユニットに電氣的に接続され、それらのデータ及びそれらの識別グループを取り出すための取り出しユニットと、

前記取り出しユニットに電氣的に接続される識別ユニットであって、

前記取り出しユニットを介して、それらのデータのいずれか 1 つを取り出し、それを取り出しデータと定義する工程と、

前記取り出しユニットを介して、それらのフォーマット特徴のいずれか 1 つを取り出し、それを取り出し特徴と定義する工程と、

前記取り出し特徴に基づいて前記取り出しデータが対応する前記特定フォーマットを有するか否かを判定し、前記取り出しデータが対応する前記特定フォーマットを有すると判定した場合に、前記特定フォーマットに対応するそれらの秘密ファクターが前記取り出しデータにおける出現頻度が秘密閾値以上であることを判定し、前記出現頻度が前記秘密閾値以上であると判定した場合に、前記取り出しデータにおける前記特定フォーマットが前記秘密データであることを表し、前記出現頻度が前記秘密閾値よりも小さいと判定した場合に、前記取り出しデータにおける前記特定フォーマットが前記秘密データではないことを表すようにする工程と、

30

それらのフォーマット特徴において取り出されていない前記フォーマット特徴があるか否かを判定し、それらのフォーマット特徴において取り出されていない前記フォーマット特徴があると判定した場合に、前記取り出しユニットを介して取り出されていない前記フォーマット特徴を取り出し、取り出されていない前記フォーマット特徴を前記取り出し特徴と定義することで、改めて前記取り出し特徴に基づいて前記取り出しデータが対応する前記特定フォーマットを有するか否かを判定し、それらのフォーマット特徴において取り出されていない前記フォーマット特徴がないと判定した場合に、前記取り出しユニットを介して複数のデータの次のデータを取り出し、前記次のデータを前記取り出しデータと定義することで、改めて前記取り出しデータが対応する前記特定フォーマットを有するか否かを判定する工程と、を実行する識別ユニットと、

40

を備えることを特徴とする秘密データを識別する電子装置。

【請求項 10】

前記識別ユニットは、前記取り出しデータが対応する前記特定フォーマットを有しないと判定した場合に、複数のフォーマット特徴において取り出されていない前記フォーマット特徴があるか否かを判定することを特徴とする請求項 9 に記載の秘密データを識別する

50

電子装置。

【請求項 1 1】

前記識別ユニットに電氣的に接続される分類ユニットであって、前記識別ユニットがそれらのフォーマット特徴において取り出されていない前記フォーマット特徴がないと判定した場合に、それらの秘密ファクターとそれらの秘密ファクターがそれらのデータに出現する回数とに基づいて前記取り出しデータに対して分類を行う分類ユニットをさらに備えることを特徴とする請求項 9 に記載の秘密データを識別する電子装置。

【請求項 1 2】

前記取り出し特徴が同一列に 2 つの列終了位置を有し、かつ前記識別ユニットが前記特定フォーマットにおいて同一列に 2 つの列終了位置を有する数がフォーマット閾値以上であると判定した場合に、前記識別ユニットは前記取り出しデータが前記特定フォーマットを有すると判定することを特徴とする請求項 9 に記載の秘密データを識別する電子装置。

10

【請求項 1 3】

前記フォーマット特徴が特定鍵からのメッセージを含み、かつ前記識別ユニットが前記特定フォーマットにおいて前記メッセージを有する数がフォーマット閾値以上であると判定した場合に、前記識別ユニットは前記取り出しデータが前記特定フォーマットを有すると判定することを特徴とする請求項 9 に記載の秘密データを識別する電子装置。

【請求項 1 4】

前記フォーマット特徴がカスタマイズ特徴を含み、かつ前記識別ユニットが前記特定フォーマットにおいて前記カスタマイズ特徴を有する数がフォーマット閾値よりも大きいと判定した場合に、前記識別ユニットは前記取り出しデータが前記特定フォーマットを有すると判定することを特徴とする請求項 9 に記載の秘密データを識別する電子装置。

20

【請求項 1 5】

各前記識別グループのそれらの秘密ファクターは、少なくとも 1 つのワード、少なくとも 1 つのストリング、少なくとも 1 つの符号、少なくとも 1 つの数字、少なくとも 1 つの実行指令、及び少なくとも 1 つのフォーマットのいずれか 1 つまたはそれらの組み合わせであることを特徴とする請求項 9 に記載の秘密データを識別する電子装置。

【請求項 1 6】

各前記フォーマット特徴は、少なくとも 1 つのワード、少なくとも 1 つのストリング、少なくとも 1 つの符号、少なくとも 1 つの数字、少なくとも 1 つの実行指令、及び少なくとも 1 つのフォーマットのいずれか 1 つまたはそれらの組み合わせであることを特徴とする請求項 9 に記載の秘密データを識別する電子装置。

30

【請求項 1 7】

ユーザコンピュータと遠隔サーバとの間に設けられ、前記ユーザコンピュータと前記遠隔サーバとの間に伝送される各前記データにおける前記特定フォーマットが秘密データであるか否かを識別することを特徴とする請求項 9 に記載の秘密データを識別する電子装置。

【請求項 1 8】

ユーザコンピュータに接続され、ネットワーク接続を介して前記ユーザコンピュータのそれらのデータを取り出し、各前記データにおける前記特定フォーマットが秘密データであるか否かを識別することを特徴とする請求項 9 に記載の秘密データを識別する電子装置。

40

【請求項 1 9】

ユーザコンピュータの内部に設けられ、前記ユーザコンピュータからそれらのデータが出力された場合に、それらのデータを取り出し、各前記データにおける前記特定フォーマットが秘密データであるか否かを識別することを特徴とする請求項 9 に記載の秘密データを識別する電子装置。

【請求項 2 0】

コンピュータによって実行可能なプログラムが記録され、プロセッサによって読み取られた場合に、前記プロセッサは、前記コンピュータによって実行可能なプログラムを

50

実行することで、請求項 1 に記載の秘密データを識別する方法を実施可能であることを特徴とするコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、秘密データを識別する方法、電子装置及びコンピュータ読み取り可能な記録媒体に関し、特にファイルにおける特定フォーマットが秘密データであるか否かを識別する方法、電子装置及びコンピュータ読み取り可能な記録媒体に関するものである。

【背景技術】

【0002】

秘密データを識別する技術は、データ保護の関連分野に用いられる。秘密データの識別メカニズムを通じて、高機密性が潜在する秘密データをさらに識別することができる。

【0003】

従来の秘密データ識別技術は、個人データまたは秘密ストリングに対してのみ分析識別するものであり、機密レベルが、見つけ出した類型またはカウント数に比例するのが一般的である。カウント数が多くないが機密記述が大量に含まれたデータ（例えば履歴、カルテ等）に対して、正しい機密レベルを提供することができない。また、従来の秘密データ識別技術は、大量の既知データ全体の内容について学習を行い、既知データの特徴を取得した後、上記特徴を識別すべきデータの特徴と比較対照することで、識別データが秘密データであるか否かを判定する。従って、従来の秘密データ識別技術は、既知データと同一または類似する秘密データしか見つけ出すことができず、既知データと同一のテンプレートまたはフォーマットを使用した秘密データを見つけ出すことができない。

【発明の概要】

【発明が解決しようとする課題】

【0004】

カウント数が多くないが機密記述が大量に含まれたデータの正しい機密程度を提供するとともに特定のテンプレートまたはフォーマットを有する秘密データを識別することによってデータの漏れを回避することができる秘密データを識別する方法、電子装置及びコンピュータ読み取り可能な記録媒体を提供する。

【課題を解決するための手段】

【0005】

本発明は、特定フォーマットを表すフォーマット特徴と前記特定フォーマットを秘密データとして表す複数の機密ファクターとをそれぞれ有すると共に前記特定フォーマットにそれぞれ対応する複数の識別グループが格納された電子装置に適用される、秘密データを識別する方法であって、複数のデータのいずれか 1 つを取り出し、それを取り出しデータと定義する工程と、複数のフォーマット特徴のいずれか 1 つを取り出し、それを取り出し特徴と定義する工程と、電子装置が、取り出し特徴に基づいて取り出しデータが対応する特定フォーマットを有するか否かを判定し、取り出しデータが対応する特定フォーマットを有すると判定した場合に、特定フォーマットに対応する複数の秘密ファクターの取り出しデータにおける出現頻度が秘密閾値以上であるかを判定し、出現頻度が秘密閾値以上であると判定した場合に、取り出しデータにおける特定フォーマットが秘密データであることを表し、出現頻度が秘密閾値よりも小さいと判定した場合に記取り出しデータにおける特定フォーマットが秘密データではないことを表すようにする工程と、電子装置が、複数のフォーマット特徴において取り出されていないフォーマット特徴があるか否かを判定し、複数のフォーマット特徴において取り出されていないフォーマット特徴があると判定した場合に、取り出されていないフォーマット特徴を取り出し、取り出されていないフォーマット特徴を取り出し特徴と定義することで、改めて取り出し特徴に基づいて取り出しデータが対応する特定フォーマットを有するか否かを判定し、複数のフォーマット特徴において取り出されていないフォーマット特徴がないと判定した場合に、複数のデータの次のデータを取り出し、次のデータを取り出しデータと定義することで、改めて取り出しデー

10

20

30

40

50

タが対応する特定フォーマットを有するか否かを判定する工程と、を備えることを特徴とする秘密データを識別する方法を提供する。

【0006】

また、本発明は、特定フォーマットを表すフォーマット特徴と特定フォーマットを秘密データとして表す複数の機密ファクターとをそれぞれ有すると共に特定フォーマットにそれぞれ対応する複数の識別グループを格納するための格納ユニットと、格納ユニットに電氣的に接続され、複数のデータ及び複数の識別グループを取り出すための取り出しユニットと、取り出しユニットに電氣的に接続される識別ユニットであって、取り出しユニットを介して、複数のデータのいずれか1つを取り出し、それを取り出しデータと定義する工程と、取り出しユニットを介して、複数のフォーマット特徴のいずれか1つを取り出し、それを取り出し特徴と定義する工程と、取り出し特徴に基づいて取り出しデータが対応する特定フォーマットを有するか否かを判定し、取り出しデータが対応する特定フォーマットを有すると判定した場合に、特定フォーマットに対応する複数の秘密ファクターの取り出しデータにおける出現頻度が秘密閾値以上であるかを判定し、出現頻度が秘密閾値以上であると判定した場合に、取り出しデータにおける特定フォーマットが秘密データであることを表し、出現頻度が秘密閾値よりも小さいと判定した場合に、取り出しデータにおける特定フォーマットが秘密データではないことを表すようにする工程と、複数のフォーマット特徴において取り出されていないフォーマット特徴があるか否かを判定し、複数のフォーマット特徴において取り出されていないフォーマット特徴があると判定した場合に、取り出しユニットを介して取り出されていないフォーマット特徴を取り出し、取り出されていないフォーマット特徴を取り出し特徴と定義することで、改めて取り出し特徴に基づいて取り出しデータが対応する特定フォーマットを有するか否かを判定し、複数のフォーマット特徴において取り出されていないフォーマット特徴がないと判定した場合に、取り出しユニットを介して複数のデータの次のデータを取り出し、次のデータを取り出しデータと定義することで、改めて取り出しデータが対応する特定フォーマットを有するか否かを判定する工程と、を実行する識別ユニットと、を備えることを特徴とする秘密データを識別する電子装置を提供する。

10

20

【0007】

また、本発明は、コンピュータによって実行可能なプログラムが記録され、プロセッサによって読み取られた場合に、プロセッサは、上記秘密データを識別する方法における工程を実行可能であることを特徴とするコンピュータ読み取り可能な記録媒体を提供する。

30

【発明の効果】

【0008】

上記のように、本発明に係る秘密データを識別する方法、電子装置及びコンピュータ読み取り可能な記録媒体によれば、特定フォーマットを有するデータが秘密データであるか否かを判定することができる。これにより、本発明に係る秘密データを識別する方法、電子装置及びコンピュータ読み取り可能な記録媒体は、カウント数が多くないが機密記述が大量に含まれたデータの正しい機密レベルを提供するとともに特定フォーマットを有する秘密データを識別することができ、データの漏れを回避することができる。

40

【図面の簡単な説明】

【0009】

【図1】本発明の一実施例に係る秘密データを識別する電子装置の模式図である。

【図2A】本発明の一実施例に係る秘密データを識別する方法のフロー図である。

【図2B】本発明の一実施例に係る秘密データを識別する方法のフロー図である。

【図3A】本発明の一実施例に係る電子装置が取り出しデータにフォームがあると判定した様子を示す模式図である。

【図3B】本発明の一実施例に係る電子装置が取り出しデータにフォームがあると判定した様子を示す模式図である。

【図4A】本発明の他の実施例に係る電子装置が取り出しデータにリストがあると判定し

50

た様子を示す模式図である。

【図4B】本発明の他の実施例に係る電子装置が取り出しデータにリストがあると判定した様子を示す模式図である。

【図5A】本発明の他の実施例に係る電子装置が取り出しデータにパターンがあると判定した様子を示す模式図である。

【図5B】本発明の他の実施例に係る電子装置が取り出しデータにパターンがあると判定した様子を示す模式図である。

【図6】本発明の他の実施例に係る電子装置が受信したデータにおける特定フォーマットの内容が秘密データであるか否かを判定する。

【発明を実施するための形態】

【0010】

以下、本発明の各種の例示性実施例について、添付図面を参照しながら詳しく説明する。ここで説明しておきたいのは、本発明の概念は、異なる形式で表現されるため、明細書に述べた例示性実施例に限定されるものではない。また、図面における同一素子には、同一符号を付す。

【0011】

本発明の実施例に係る秘密データを識別する電子装置は、特定フォーマットを表すフォーマット特徴に基づいてデータにおいて特定フォーマットがあるか否かを判定し、次に、さらに特定フォーマットを秘密データとして表す複数の秘密ファクターに基づいて、データにおける特定フォーマットが秘密データであるか否かを判定する。また、本発明の実施例に係る電子装置に対応して実行される秘密データを識別する方法において、ファームウェア、ソフトウェアまたはハードウェア回路の方法により電子装置に実施可能である。

【0012】

まず、図1は、本発明の一実施例に係る秘密データを識別する電子装置の模式図である。図1に示すように、秘密データを識別する電子装置100は、電子装置100によって受信されたデータにおける特定フォーマットの内容が秘密データであるか否かを識別し、データの漏れを回避するためのものである。この実施例において、電子装置100は、スマートフォン、デスクトップコンピュータ、ノートブックコンピュータ、またはその他データを受信可能な電子装置であってもよい。

【0013】

電子装置100は、ユーザコンピュータと遠隔サーバとの間（図示せず）に設けられており、ユーザコンピュータと遠隔サーバとの間に伝送されるデータにおける特定フォーマットが秘密データであるか否かを識別することができる。また、電子装置100は、ユーザコンピュータ（図示せず）に電氣的に接続されることで、ネットワーク接続を介してユーザコンピュータにおけるデータを取り出すとともに、取り出されたデータにおける特定フォーマットが秘密データであるか否かを識別することもできる。さらに、電子装置100は、ユーザコンピュータの内部（図示せず）に設けられることで、ユーザコンピュータからデータが出力された場合に、出力されたデータにおける特定フォーマットが秘密データであるか否かを識別することができる。本発明は、電子装置の設置位置について何ら制限するものではない。これにより、電子装置100は、秘密データが窃取意図のある者によって取得されることを防止し、データの漏れを回避することができる。

【0014】

電子装置100は、演算処理ユニットとしての識別ユニット110と、取り出しユニット120と、格納ユニット130とを含む。格納ユニット130には複数の識別グループ132が格納されている。各識別グループ132は、特定フォーマットに対応し、かつ対応する特定フォーマットを表すフォーマット特徴FFを有する。つまり、各識別グループ132がフォーマット特徴FFを有することで、演算処理ユニットとしての識別ユニット110は、データにおける内容が対応する特定フォーマットを有するか否かを識別することができる。1つの例として、特定フォーマットがフォーム（FORM）である場合、フォームのフォーマット特徴FFは、複数の列において2つの列終了位置（End-of-

10

20

30

40

50

Line)を有する特徴であってもよい。さらに例を挙げれば、特定フォーマットがリスト(LIST)である場合、リストのフォーマット特徴FFは、複数の「TAB」鍵からのメッセージを有する特徴であってもよい。さらに例を挙げれば、特定フォーマットがユーザ自身によって定義されたテンプレート(TEMPLATE)である場合、テンプレートのフォーマット特徴FFは、ユーザ自身によって定義された特徴であってもよい。この実施例において、各フォーマット特徴FFは、少なくとも1つのワード、少なくとも1つのストリング、少なくとも1つの符号、少なくとも1つの数字、少なくとも1つの実行指令、及び少なくとも1つのフォーマットのいずれか1つまたはそれらの組み合わせであってもよく、これらに限定されるものではない。

【0015】

また、各識別グループ132は、対応する特定フォーマットを秘密データとして表す複数の秘密ファクターCPを有する。つまり、各識別グループ132が複数の秘密ファクターCPを有することで、演算処理ユニットとしての識別ユニット110は、データにおける特定フォーマットの内容が秘密データであるか否かを識別することができる。1つの例として、特定フォーマットが履歴フォーム(図3Aを参照)である場合、秘密ファクターCPは、「名前」、「身分証明書」、「携帯電話」、及び「連絡住所」等の名詞であってもよい。さらに例を挙げれば、特定フォーマットが住所録リスト(図4Aを参照)である場合、秘密ファクターCPは、「生年月日」、「身長」、「体重」、「住所」、及び「電話」等の名詞であってもよい。さらに例を挙げれば、特定フォーマットがユーザ自身によって定義されたテンプレート(図5Aを参照)である場合、秘密ファクターCPは、「計画目的」、及び「お客様要求」等、ユーザ自身によって定義された名詞であってもよい。この実施例において、各識別グループ132に対応する複数の秘密ファクターCPは、少なくとも1つのワード、少なくとも1つのストリング、少なくとも1つの符号、少なくとも1つの数字、少なくとも1つの実行指令、及び少なくとも1つのフォーマットのいずれか1つまたはそれらの組み合わせであってもよく、これらに限定されるものではない。

【0016】

電子装置100において複数の識別グループ132が格納ユニット130に格納される方法は、従来の格納方法である。当業者は、電子装置100において複数の識別グループ132が格納ユニット130に格納される方法を理解することができるため、ここでは詳しい説明を省略する。この実施例において、格納ユニット130は、フラッシュメモリチップ、リードオンリーメモリチップまたはランダムアクセスメモリチップ等、揮発性または非揮発性の記憶チップであってもよく、好ましくは非揮発性メモリである。

【0017】

また、電子装置100は、ユーザが識別インタフェースにおいて識別しようとする特定フォーマット(例えばユーザ自身によって定義された名詞)を設定し、かつ受信されたデータにおける特定フォーマットの内容が秘密データであるか否かを識別できるように、識別インタフェース(図示せず)を表示するための表示ユニットをさらに有する。当然ながら、識別しようとする特定フォーマット及びそれに対応する識別グループ132が予め格納ユニット130に設定された場合には、表示ユニットを設けなくてもよく、本発明はこれに限定されるものではない。

【0018】

取り出しユニット120は、識別ユニット110が受信されたデータをさらに識別できるように、格納ユニット130に電氣的に接続されるとともに、複数のデータ及び複数の識別グループ132を取り出すものである。識別ユニット110は、取り出しユニット120に電氣的に接続され、電子装置100の主要な演算中心としての演算処理ユニットであり、各分析、演算及び制御を行うものである。この実施例において、識別ユニット110は、中央処理器、マイクロ制御器または埋め込み型制御器等の処理チップであってもよい。識別ユニット110及び取り出しユニット120は、中央処理器、マイクロ制御器または埋め込み型制御器等の処理チップに統合されてもよく、本発明は、これに限定されるものではない。

10

20

30

40

50

【0019】

識別ユニット110は、下記の工程を実行することで、受信されたデータにおける特定フォーマットの内容が秘密データであるか否かを識別する。

【0020】

図1、図2Aを同時に参照すると、まず、識別ユニット110は、取り出しユニット120を介して複数のデータのいずれか1つを取り出し、それを取り出しデータと定義することで、取り出しデータにおける特定フォーマットの内容が秘密データであるか否かをさらに識別する(ステップS210)。識別ユニット110は、取り出しユニット120を介して外部装置から上記複数のデータを取り出すか、または格納ユニット130に予め格納された複数のデータを取り出すことができ、本発明はそれに限定されるものではない。

10

【0021】

次に、識別ユニット110は、取り出しユニット120を介して格納ユニット130に格納された複数のフォーマット特徴FFのいずれか1つを取り出し、それを取り出し特徴(ステップS220)と定義する。この場合の取り出し特徴は、ある特定フォーマット(例えばフォームまたはリスト等の特定フォーマット)を表す。さらに、識別ユニット110は、取り出し特徴に基づいて取り出しデータが対応する特定フォーマットを有するか否かを判定する(ステップS230)。即ち、識別ユニット110は、取り出しデータに所定の数量の取り出し特徴があるか否かを判定することにより、取り出しデータに現在取り出されたフォーマット特徴FFの特定フォーマットがあるか否かを判定する。この実施例において、特定フォーマットは、フォーム、リスト、ユーザ自身によって定義されたテンプレート、またはその他規則性特徴を有する特定フォーマットであってもよく、本発明はそれに限定されるものではない。特定フォーマットに対応するフォーマット特徴FFは、特定フォーマットにおいてのみ出現する特徴、例えば特定鍵からのメッセージ、連続ブランク等の特徴から選択されてもよく、本発明はそれに限定されるものではない。

20

【0022】

識別ユニット110が取り出しデータにおいて対応する特定フォーマットがあると判定した場合には、取り出しデータにおいて取り出し特徴に対応する特定フォーマットがあることを表す。この場合、識別ユニット110は、取り出しデータにおける特定フォーマットの内容が秘密データであるか否かをさらに判定する(ステップS240)。逆に、識別ユニット110が、取り出しデータにおいて対応する特定フォーマットがないと判定した

30

。

【0023】

1つの例として、特定フォーマットがフォームである場合、そのフォーマット特徴FFは、図3Aに示すように、同一列に少なくとも2つの列終了位置を有するものである。従って、取り出しユニット120がフォームを表すフォーマット特徴FFを取り出した場合に、識別ユニット110は、フォームの内容において、その同一列に少なくとも2つの列終了位置を有する数がフォーマット閾値以上であるか否かを判定する。YES(はい)と判定した場合に、識別ユニット110は、取り出しデータにフォームを表す特定フォーマットがあると認定する。逆に、識別ユニット110は、取り出しデータにフォームを表す特定フォーマットがないと認定する。上記フォーマット閾値は、実際のフォームに応じて設定することができ、本発明はそれに限定されるものではない。識別ユニット110は、取り出しデータにフォームを表す特定フォーマットがあるか否かを識別した後、取り出しユニット120を介してフォームにおける内容(図3Bを参照)を取り出し、フォームにおける内容が秘密データであるか否かをさらに判定する。

40

【0024】

さらに例を挙げれば、特定フォーマットがリストである場合、そのフォーマット特徴FFは、図4Aに示すように、複数の「TAB」からのメッセージである。従って、取り出

50

しユニット120がリストを表すフォーマット特徴FFを取り出した場合に、識別ユニット110は、リストにおける内容に上記メッセージを有する数がフォーマット閾値以上であるかを判定する。YESと判定した場合に、識別ユニット110は、取り出しデータにリストを表す特定フォーマットがあると認定する。逆に、識別ユニット110は、取り出しデータにリストを表す特定フォーマットがないと認定する。上記フォーマット閾値は、実際のリストに基づいて設定してもよく、本発明はそれに限定されるものではない。識別ユニット110は、取り出しデータにリストを表す特定フォーマットがあるか否かを識別した後、取り出しユニット120を介してリストにおける内容を取り出し(図4Bを参照)、リストにおける内容が秘密データであるか否かをさらに判定する。

【0025】

さらに例を挙げれば、特定フォーマットがユーザ自身によって定義されたテンプレートである場合、そのフォーマット特徴FFは、カスタマイズ特徴である。即ち、フォーマット特徴FFは、ユーザ自身によって定義されてなるものである。図5Aに示すように、カスタマイズ特徴は、「計画目的」及び「お客様要求」等の特徴である。従って、取り出しユニット120がカスタマイズ特徴を表すフォーマット特徴FFを取り出した場合に、識別ユニット110は、テンプレートの内容に上記カスタマイズ特徴を有する数がフォーマット閾値以上であるかを判定する。YESと判定した場合に、識別ユニット110は、取り出しデータにテンプレートを表す特定フォーマットがあると認定する。逆に、識別ユニット110は、取り出しデータにテンプレートを表す特定フォーマットがないと認定する。上記フォーマット閾値は、実際のテンプレートに基づいて設定してもよく、本発明はそれに限定されるものではない。識別ユニット110は、取り出しデータにテンプレートを表す特定フォーマットがあるか否かを識別した後、取り出しユニット120を介してテンプレートにおける内容を取り出し(図5Bを参照)、テンプレートにおける内容が秘密データであるか否かをさらに判定する。

【0026】

上記の3つの例において、当業者は、識別ユニット110が取り出しユニット120を介して特定フォーマット(例えばフォーム、リスト、テンプレート)における内容を取り出す実施方法を理解することができるため、ここでは詳しい説明を省略する。

【0027】

ステップS240に戻り、識別ユニット110は、この特定フォーマットに対応する複数の秘密ファクターCPの取り出しデータにおける出現頻度が秘密データ閾値以上であるかを判定することにより、取り出しデータにおける特定フォーマットの内容が秘密データであるか否かを判定する。秘密ファクターCPは、対応する特定フォーマットが秘密データである確率を表すものである。従って、特定フォーマットにおいて秘密ファクターCPが多く出現するほど、特定フォーマットが秘密データである確率が高いことを表す。秘密ファクターCPの設定について、前の実施例に記載された通りであるため、ここでは詳しい説明を省略する。これにより、識別ユニット110が、秘密ファクターCPの出現頻度が秘密閾値以上であると判定した場合に、取り出しデータにおける特定フォーマットが秘密データであることを表す(ステップS250)。逆に、識別ユニット110が、秘密ファクターCPの出現頻度が秘密閾値よりも小さいと判定した場合に、取り出しデータにおける特定フォーマットが秘密データではないことを表す(ステップS260)。上記秘密閾値は、実際の複数の秘密ファクターCPの取り出しデータにおける出現頻度に基づいて設定されたものであり、本発明はそれに限定されるものではない。

【0028】

1つの例として、図3A~図3Bに示すように、特定フォーマットがフォームであるとする。このうち、フォームは、秘密ファクターCPの名詞として、「名前」、「身分証明書」、「携帯電話」、及び「連絡住所」を有する。各名詞には、例えば「名前」と同義である「名字」、「名称」、「人名」、「Name」等の同義字が現れる可能性がある。従って、判定の過程において、識別ユニット110は、同義字を同一の字句と見なす。この実施例において、識別ユニット110は、同義字関数STF(i)を介して各字句がフォ

10

20

30

40

50

ームに出現する重要性を算出することで、各字句とフォームとの間の関連性を得ることができる。本実施例における同義字関数 $STF(i)$ は、以下のように示すことができる。

【0029】

【数1】

$$STF(i) = \frac{n_{ij}}{\sum_k N_{kj}} \times \omega_i$$

【0030】

10

ここで、 n_{ij} は、第 i 種の字句が第 j 個のフォームに出現する回数を表す。 ω_i は第 i 種の字句の重みを表す。 N_{kj} は第 j 個のフォームにおけるすべての k 個の字句を表し、かつ $k = 0$ 。

【0031】

ここで注意すべき点は、識別ユニット 110 が同義字を同一の字句と見なす点である。即ち、識別ユニット 110 がフォームにおける「連絡住所」、「名前」、「名称」、「人名」、及び「身分証明書」を見つけ出した場合、識別ユニット 110 は、「連絡住所」を第 1 種の名詞として見なし、「名前」、「名称」、「人名」を第 2 種の字句として見なし、「身分証明書」を第 3 種の字句として見なす。各種の字句の重みについて、 ω_1 が 0.5 であり、 ω_2 が 0.2 であり、 ω_3 が 0.3 であるとする場合、識別ユニット 110 は、同義字関数 STF を介して各字句がフォームに出現する重要性を算出する。第 1 種の字句としては、 $STF(1) = 1/5 * 0.5 = 0.1$ であり、第 2 種の字句としては、 $STF(2) = 3/5 * 0.2 = 0.12$ であり、第 3 種の字句としては、 $STF(3) = 1/5 * 0.3 = 0.06$ である。

20

【0032】

次に、この実施例における識別ユニット 110 は、さらに、情報関数 PIF を介してフォームにおいて秘密ファクター CP の字句として出現する確率を算出する。この実施例における情報関数 PIF は、以下の通りである。

【0033】

【数2】

$$PIF = \frac{P_n}{P_t}$$

30

【0034】

ここで、 P_t は、現在秘密ファクター CP としての名句の数を表す。 P_n は、フォームにおいて秘密ファクター CP の字句として出現する数を表す。上記の例としては、フォームには、秘密ファクター CP の名詞として、「名前」、「身分証明書」、「携帯電話」、及び「連絡住所」の 4 つの名詞がある。識別ユニット 110 は、フォームにおいて「連絡住所」、「名前」、「名称」、「人名」、及び「身分証明書」の 5 つの名詞を見つけ出し、見つけ出した 5 つの名詞を 3 種の字句に分類する。この場合、演算処理ユニットとしての識別ユニット 110 が $PIF = 3/4$ として算出したため、フォームにおいて秘密ファクター CP の名詞として出現する確率が 75% であることを表す。

40

【0035】

次に、識別ユニット 110 は、秘密データ関数 $PIFV$ を介して、フォームに対応する 4 つの秘密ファクター CP の取り出しデータにおける出現頻度を算出する。この実施例における秘密データ関数 $PIFV$ は、以下の通りである。

【0036】

【数 3】

$$PIFV = \left(\sum_n STF(i) \right) \times PIF$$

【0037】

ここで、 $nSTF(i)$ は、各字句がフォームにおいて出現する重要性の総計を表す。PIF は、フォームにおいて秘密ファクターの字句として出現する確率を表す。上記の例に続き、 $PIFV = (0.1 + 0.12 + 0.06) \times 0.75 = 0.21$ であることは、フォームに対応する4つの秘密ファクターCPの取り出しデータにおける出現頻度が0.21であることを表す。

10

【0038】

最後に、識別ユニット110は、出現頻度が秘密閾値以上であるかを判定する。上記の例に続き、この実施例における秘密閾値は0.1とする。従って、識別ユニット110は、秘密ファクターCPの出現頻度(0.21である)が秘密閾値(0.1である)よりも大きいと判定し、取り出しデータにおけるフォームの内容が秘密データであることを表す。これにより、識別ユニット110は、ステップS210~S260を介して、取り出されたデータにおける特定フォーマットが秘密データであるか否かを判定することができる。

【0039】

これにより、識別ユニット110は、特定フォーマットを表す秘密ファクターCPを介して取り出しデータにおける特定フォーマットの秘密性を識別することができ、高秘密性のデータの漏れを回避することができる。

20

【0040】

次に、識別ユニット110は、複数のフォーマット特徴FFにおいて取り出されていないフォーマット特徴FFがあるか否かを判定する(ステップS270)。即ち、識別ユニット110は、取り出しデータにその他の特定フォーマットがあるか否かをさらに判定する。識別ユニット110が、取り出されていないフォーマット特徴FFがあると判定した場合に、ステップS220に戻り、取り出しユニット120を介して取り出されていないフォーマット特徴FFを取り出す。この場合、識別ユニット110は、取り出されていないフォーマット特徴FFを取り出し特徴と定義することで、改めて定義された取り出し特徴に基づいて取り出しデータに対応する特定フォーマットがあるか否かを改めて判定する。上記の例に続き、フォームのフォーマット特徴FFを判定した後、識別ユニット110がリストを表すフォーマット特徴FFが取り出されていないと判定した場合に、識別ユニット110は、リストを表すフォーマット特徴FF(即ちフォーマット特徴FFが複数の「TAB」鍵からのメッセージである)を取り出し特徴として定義することで、改めて取り出し特徴に基づいて取り出しデータにリストのフォーマットがあるか否かを判定する。

30

【0041】

逆に、識別ユニット110が、取り出されていないフォーマット特徴がないと判定した場合に、取り出しデータに判定すべき特定フォーマットがないことを表す。この場合、識別ユニット110は、ステップS210に戻り、複数のデータにおける次のデータを取り出す。さらに、識別ユニット110は、次のデータを取り出しデータとして定義することで、取り出しデータに対応する特定フォーマットがあるか否かを改めて判定する。

40

【0042】

また、図1、図2A、図2Bを同時に参照すると、電子装置100は、分類ユニット140をさらに含む。分類ユニット140は、識別ユニット110に電氣的に接続され、現在の取り出しデータに対して分類を行うものである。さらに詳しくは、識別ユニット110が、取り出されたフォーマット特徴FFがないと判定した場合に、分類ユニット140は、現在の取り出しデータに対してさらに分類することで、取り出しデータにおける特定フォーマットがどの種類であるかをさらに判定することができる(ステップS275)。識別ユニット110は、分類ユニット140が現在の取り出しデータの分類を終了した後

50

に、ステップS 2 1 0に戻り、複数のデータにおける次のデータを取り出す。1つの例として、分類ユニット1 4 0は、フォームを有する取り出しデータを履歴表、給料表、カルテ表、またはその他高秘密性のフォームに分類する。若しくは、分類ユニット1 4 0は、リストを有する取り出しデータを住所録、内線表、またはその他高秘密性のリストに分類する。

【0043】

この実施例において、すべてのデータを関連性を有するため、分類ユニット1 4 0は、特定フォーマットにおける複数の秘密ファクターCPと、上記秘密ファクターCPがすべてのデータにおいて出現する回数とに基づいて現在の取り出しデータに対して分類を行う。例えば、分類ユニット1 4 0は、「履歴」、「名前」、「身分証明書」、「携帯電話」、及び「連絡住所」の5つの字句を秘密ファクターCPとする。分類ユニット1 4 0は、上記の5つの字句と、上記の5つの字句がすべてのデータにおいて出現する回数とに基づいて現在の取り出しデータに対して分類を行う。当然ながら、すべてのデータの間に関連性がない場合には、分類ユニット1 4 0は、特定フォーマットの複数の秘密ファクターCPにのみ基づいて現在の取り出しデータに対して分類を行うこともでき、本発明はそれに限定されるものではない。

10

【0044】

また、本実施例における分類ユニット1 4 0は、例えばTFIDF (term frequency - inverse document frequency)、サポートベクトルマシン (support vector machines、SVM)、ベイジアン分類法 (bayesian classification)、またはバックプロパゲーションニューラルネットワーク (back propagation neural (BPN) network) 等の分類アルゴリズムにより、現在の取り出しデータに対して分類を行うことで、取り出しデータの分類をより正確に行う。当業者は、分類ユニット1 4 0が分類アルゴリズムにより現在の取り出しデータに対して分類を行う実施及び運用方法を理解することができるため、ここでは詳しい説明を省略する。

20

【0045】

これにより、分類ユニット1 4 0は、特定フォーマットの取り出しデータに対して分類を行うことができる。従って、すべてのデータの識別が終了した場合に、ユーザは、すべてのデータにおける特定フォーマットがどの種類であるかを理解することができ、すべてのデータに対して制御を行うことができる。

30

【0046】

以下、ユーザがユーザコンピュータ10を介して1つのデータDAを遠隔サーバ20に伝送することを例にして説明する。図6に示すように、電子装置100は、ユーザコンピュータ10と遠隔サーバ20との間に設けられることで、ユーザコンピュータ10からのデータDAにおける特定フォーマットの内容が秘密データであるか否かを判定する。説明の簡単化のために、本実施例におけるデータDAは、図3Aに示すフォームを有し、この場合に取り出されたフォーマット特徴FFは、フォームを表す特定フォーマットである。

【0047】

図1、図3A、図6を同時に参照すると、ユーザがユーザコンピュータ10を介してデータDAを遠隔サーバ20に伝送する過程において、電子装置100における識別ユニット110は、取り出しユニット120を介してデータDAを取り出す。この場合、電子装置100は、データDAにおける特定フォーマットの内容が秘密データであるか否かを判定し、かつ、秘密データの漏れを回避するために、しばらくの間、データDAを遠隔サーバ20に伝送しない。

40

【0048】

まず、電子装置100における識別ユニット110は、現在取り出されたフォーマット特徴FF (即ちフォームを表す特定フォーマット) に基づいてデータDAにおいてフォームを表す特定フォーマットがあると判定する。識別ユニット110がデータDAにおいてフォームを表す特定フォーマットがあるか否かを判定する方法について、上記の実施例に

50

記載された通りであるため、ここでは詳しい説明を省略する。

【 0 0 4 9 】

次に、電子装置 1 0 0 における識別ユニット 1 1 0 は、フォームを表す特定フォーマットに対応する複数の秘密ファクター C P のデータ D A における出現頻度に基づいて、データ D A におけるフォームの内容が秘密データであると判定する。識別ユニット 1 1 0 がデータ D A においてフォームを表す特定フォーマットの内容が秘密データであるか否かを判定する方法について、上記の実施例に記載された通りであるため、ここでは詳しい説明を省略する。

【 0 0 5 0 】

さらに、電子装置 1 0 0 における識別ユニット 1 1 0 は、まだ識別していないフォーマット特徴 F F があるか否かをさらに判定する。この実施例において、この場合の識別ユニット 1 1 0 には取り出されていないフォーマット特徴 F F が既にある。即ち、識別ユニット 1 1 0 は、データ D A における特定フォーマットを既に判定した。次に、電子装置 1 0 0 における分類ユニット 1 4 0 は、複数の秘密ファクター C P に基づいてデータ D A に対して分類を行うとともに、データ D A を履歴データに分類する。分類ユニット 1 4 0 がデータ D A を履歴データに分類する方法について、上記の実施例に記載された通りであるため、ここでは詳しい説明を省略する。

【 0 0 5 1 】

この場合、電子装置 1 0 0 は、ユーザコンピュータ 1 0 からのデータ D A におけるフォームが履歴データであり、かつこの履歴データが秘密データであると判定する。電子装置 1 0 0 は、データ D A におけるフォームが秘密データであると判定した後、実際の情報安全防護に基づいて後続の処理を行うことができる。例えば、電子装置 1 0 0 は、データ D A が遠隔サーバ 2 0 に伝送されることを許可しないと同時に、システム管理者に対してユーザコンピュータ 1 0 が秘密データを遠隔サーバ 2 0 に伝送中であることを通知する。これにより、電子装置 1 0 0 は、出力されたデータ D A における特定フォーマットが秘密データであるか否かを識別することができ、秘密データが窃取意図のある者によって取得されることを防止し、データの漏れを回避することができる。

【 0 0 5 2 】

また、本発明は、コンピュータ読み取り可能な記録媒体により、上記秘密データを識別する方法におけるコンピュータプログラムを格納することで上記の工程を行うこともできる。このコンピュータ読み取り可能な記録媒体は、フロッピー(登録商標)ディスク、ハードディスク、光ディスク、USBドライブ、磁気テープ、ネットワークを介してアクセス可能なデータベース、または当業者が容易に想到し得る同一機能を有する記録媒体であってもよい。

【 0 0 5 3 】

上記のように、本発明の実施例に係る秘密データを識別する方法、電子装置及びコンピュータ読み取り可能な記録媒体は、特定フォーマットを有するデータが秘密データであるか否かを判定することができる。これにより、本発明の実施例に係る秘密データを識別する方法、電子装置及びコンピュータ読み取り可能な記録媒体は、カウント数が多くないが機密記述が大量に含まれたデータの正しい機密レベルを提供するとともに特定フォーマットを有する秘密データを識別することができ、データの漏れを回避することができる。

【 0 0 5 4 】

上述したものは、本発明の好ましい実施例に過ぎず、本発明の実施の範囲を限定するためのものではない。

【 符号の説明 】

【 0 0 5 5 】

- 1 0 0 電子装置
- 1 1 0 識別ユニット
- 1 2 0 取り出しユニット
- 1 3 0 格納ユニット

10

20

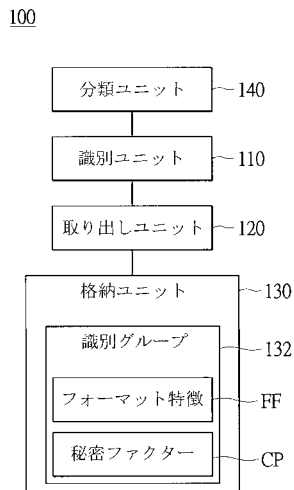
30

40

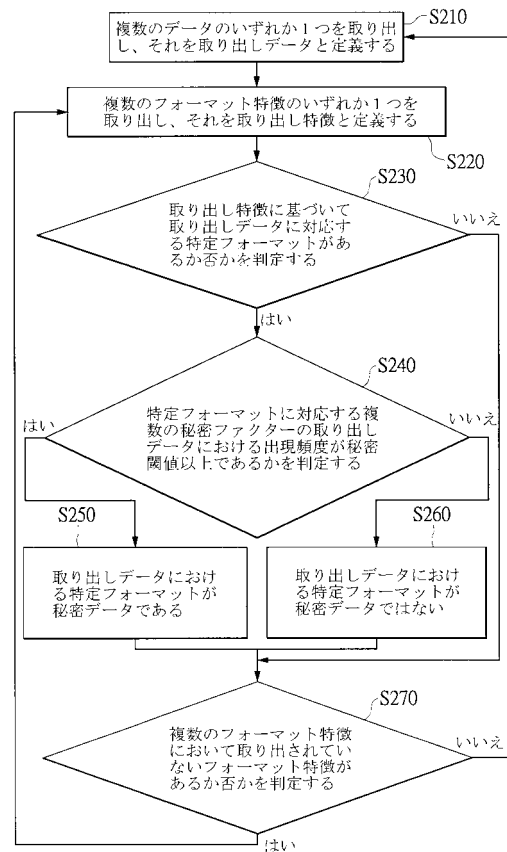
50

- 1 3 2 識別グループ
- 1 4 0 分類ユニット

【 図 1 】



【 図 2 A 】



【 図 4 B 】

学籍番号	座席番号	名前	性別	生年月日	身長(センチ)	体重(キロ)	住所	電話
253001	1	小明	男性	801202	172	67	台北市XXXXXXXX	0228224698
253002	2	小王	男性	810830	171	63	台北市XXXXXXXX	0228543289
253003	3	小強	男性	801010	172	58	台北市XXXXXXXX	0225553281
253004	4	小林	男性	810324	174	68	台北市XXXXXXXX	0225490377
253005	5	小李	男性	810629	180	73	台北市XXXXXXXX	0228973321
.....

【 図 5 A 】

テンプレート

一、計画目的
 本計画は、ネットワークバージョン個人情報検出システム（以下、本システムという）「1」を強化することを目的とする。当該ソフトウェアは、主に、企業秘密データの自動化検索収集に用いられる。

二、お客様要求：
 本プロジェクトのお客様要求の説明には、以下の「応用情景」、「システム要求」、及び「クラウド演算」の3つの部分が含まれる。

1、応用情景：
 2、システム機能：
 本プロジェクトのシステム機能は、それぞれ「機能性要求」及び「非機能性要求」の2つの部分について説明する。

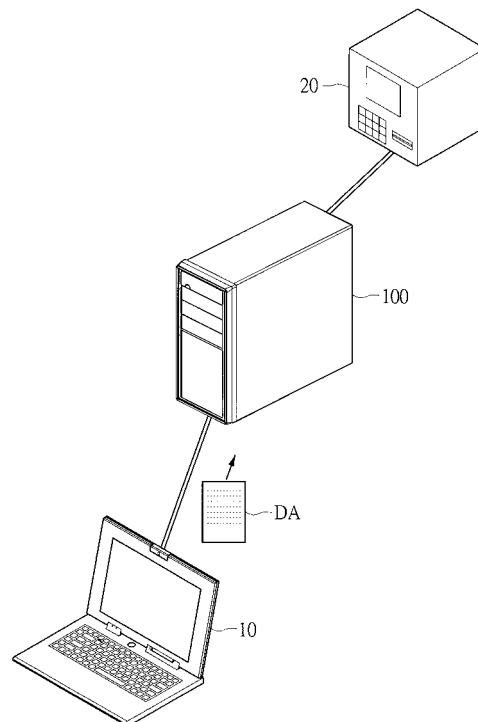
【 図 5 B 】

計画目的
 本計画は、ネットワークバージョン個人情報検出システム（以下、本システムという）「1」を強化することを目的とする。当該ソフトウェアは、主に、企業秘密データの自動化検索収集に用いられる。

お客様要求：
 本プロジェクトのお客様要求の説明には、以下の「応用情景」、「システム要求」、及び「クラウド演算」の3つの部分が含まれる。

応用情景：
 お客様要求：
 本プロジェクトのシステム機能は、それぞれ「機能性要求」及び「非機能性要求」の2つの部分について説明する。

【 図 6 】



フロントページの続き

(72)発明者 葉 信延

台湾新北市三重區三和路四段292巷3號3エフ

(72)発明者 劉 建宗

台湾新北市中和區國光街112巷2弄10號1エフ