

【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7年3月5日(2025.3.5)

【公開番号】特開2025-20137(P2025-20137A)
 【公開日】令和7年2月12日(2025.2.12)
 【年通号数】公開公報(特許)2025-026
 【出願番号】特願2024-177329(P2024-177329)
 【国際特許分類】

G 0 6 N 3/0475(2023.01)

10

G 0 6 N 3/045(2023.01)

G 0 6 N 3/08(2023.01)

【F I】

G 0 6 N 3/0475

G 0 6 N 3/045

G 0 6 N 3/08

【手続補正書】

【提出日】令和7年2月25日(2025.2.25)

【手続補正1】

20

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

コンピュータ実装される方法であって、

ターゲットニューラルネットワークを展開する本番環境によって複数のセグメントを保持するステップであって、前記複数のセグメントは、

テストケース生成ニューラルネットワークを使用することによって複数のテスト入力を生成し、

30

前記ターゲットニューラルネットワークを使用して前記複数のテスト入力を処理してテスト入力ごとに1つまたは複数の出力を生成し、

テスト入力ごとに前記ターゲットニューラルネットワークによって生成された前記1つまたは複数のテスト出力に基づいて、前記ターゲットニューラルネットワークによって1つまたは複数の基準を満たしていないテスト出力の生成をもたらす1つまたは複数の不合格のテスト入力を識別し、

前記1つまたは複数の不合格のテスト入力に基づいてセグメントを生成することによって識別されている、保持するステップと、

前記本番環境によって、ユーザからオンラインネットワーク入力を受信するステップと

40

、前記本番環境によって、前記オンラインネットワーク入力が前記複数のセグメントの少なくとも1つを含むことを決定するステップと、

それに応答して、前記本番環境によって、前記ターゲットニューラルネットワークを使用して前記オンラインネットワーク入力を処理することなく、前記オンラインネットワーク入力に応答してデフォルト出力を生成し、オンラインネットワーク出力を生成するステップと、

を含む、方法。

【請求項2】

前記本番環境が展開前の環境システムから前記複数のセグメントを受信する、請求項1

50

に記載の方法。

【請求項 3】

前記 1 つまたは複数の不合格のテスト入力を識別するステップが、
テスト入力ごとに前記ターゲットニューラルネットワークによって生成された前記 1 つ
または複数のテスト出力に自然言語処理技法を適用するステップを含む、請求項 1 または
2 に記載の方法。

【請求項 4】

前記 1 つまたは複数の不合格のテスト入力を識別するステップが、
1 つまたは複数の基準に対してテスト入力ごとに前記ターゲットニューラルネットワー
クによって生成された前記 1 つまたは複数のテスト出力を評価するステップを含む、
請求項 1 乃至 3 の何れか一項に記載の方法。

10

【請求項 5】

前記 1 つまたは複数の基準が、テスト出力が攻撃的または偏ったコンテンツ、トレー
ニング中に処理される保護可能な情報、あるいは個人の連絡先情報を含むべきでないとい
うことを指定する、請求項 4 に記載の方法。

【請求項 6】

前記 1 つまたは複数の不合格のテスト入力に基づいて前記セグメントを生成するステッ
プが、
前記 1 つまたは複数の不合格のテスト入力において発生する様々なセグメントを複数の
テキストクラスタにグループ化するためにテキストクラスタリングアルゴリズムを適用す
るステップと、
前記複数のテキストクラスタの 1 つに含まれるセグメントに基づいて、前記セグメント
を生成するステップと、を含む
請求項 1 乃至 5 の何れか一項に記載の方法。

20

【請求項 7】

前記本番環境に展開される前記ターゲットニューラルネットワークを調整するために前
記複数のセグメントを使用するステップを更に含む、
請求項 1 乃至 6 の何れか一項に記載の方法。

【請求項 8】

1 つまたは複数のコンピュータによって実行されると、前記 1 つまたは複数のコンピュ
ータに、
ターゲットニューラルネットワークを展開する本番環境によって複数のセグメントを保
持することであって、前記複数のセグメントは、
テストケース生成ニューラルネットワークを使用することによって複数のテスト入力
を生成し、
前記ターゲットニューラルネットワークを使用して前記複数のテスト入力を処理して
テスト入力ごとに 1 つまたは複数の出力を生成し、
テスト入力ごとに前記ターゲットニューラルネットワークによって生成された前記 1
つまたは複数のテスト出力に基づいて、前記ターゲットニューラルネットワークによって
1 つまたは複数の基準を満たしていないテスト出力の生成をもたらす 1 つまたは複数の不
合格のテスト入力を識別し、
前記 1 つまたは複数の不合格のテスト入力に基づいてセグメントを生成することによ
って識別されている、保持することと、
前記本番環境によって、ユーザからオンラインネットワーク入力を受信することと、
前記本番環境によって、前記オンラインネットワーク入力が入力された前記複数のセグメントの少
なくとも 1 つを含むことを決定することと、
それに応答して、前記本番環境によって、前記ターゲットニューラルネットワークを使
用して前記オンラインネットワーク入力を処理することなく、前記オンラインネットワー
ク入力に応答してデフォルト出力を生成し、オンラインネットワーク出力を生成すること
と、

30

40

50

を含む動作を実行させる命令を記憶する、1つまたは複数のコンピュータ可読記憶媒体。

【請求項9】

前記本番環境が展開前の環境システムから前記複数のセグメントを受信する、請求項8に記載のコンピュータ可読記憶媒体。

【請求項10】

前記1つまたは複数の不合格のテスト入力を識別することが、

テスト入力ごとに前記ターゲットニューラルネットワークによって生成された前記1つまたは複数のテスト出力に自然言語処理技法を適用することを含む、請求項8または9に記載のコンピュータ可読記憶媒体。

【請求項11】

前記1つまたは複数の不合格のテスト入力を識別することが、

1つまたは複数の基準に対してテスト入力ごとに前記ターゲットニューラルネットワークによって生成された前記1つまたは複数のテスト出力を評価することを含む、請求項8乃至10の何れか一項に記載のコンピュータ可読記憶媒体。

【請求項12】

前記1つまたは複数の基準が、テスト出力が攻撃的または偏ったコンテンツ、トレーニング中に処理される保護可能な情報、あるいは個人の連絡先情報を含むべきでないということ指定する、請求項11に記載のコンピュータ可読記憶媒体。

【請求項13】

前記1つまたは複数の不合格のテスト入力に基づいて前記セグメントを生成することが

、前記1つまたは複数の不合格のテスト入力において発生する様々なセグメントを複数のテキストクラスタにグループ化するためにテキストクラスタリングアルゴリズムを適用することと、

前記複数のテキストクラスタの1つに含まれるセグメントに基づいて、前記セグメントを生成することと、を含む

請求項8乃至12の何れか一項に記載のコンピュータ可読記憶媒体。

【請求項14】

1つまたは複数のコンピュータと、命令を記憶する1つまたは複数の記憶媒体とを備えるシステムであって、前記命令が1つまたは複数のコンピュータによって実行されると、前記1つまたは複数のコンピュータに、

ターゲットニューラルネットワークを展開する本番環境によって複数のセグメントを保持することであって、前記複数のセグメントは、

テストケース生成ニューラルネットワークを使用することによって複数のテスト入力を生成し、

前記ターゲットニューラルネットワークを使用して前記複数のテスト入力を処理してテスト入力ごとに1つまたは複数の出力を生成し、

テスト入力ごとに前記ターゲットニューラルネットワークによって生成された前記1つまたは複数のテスト出力に基づいて、前記ターゲットニューラルネットワークによって1つまたは複数の基準を満たしていないテスト出力の生成をもたらす1つまたは複数の不合格のテスト入力を識別し、

前記1つまたは複数の不合格のテスト入力に基づいてセグメントを生成することによって識別されている、保持することと、

前記本番環境によって、ユーザからオンラインネットワーク入力を受信することと、

前記本番環境によって、前記オンラインネットワーク入力が入力された前記複数のセグメントの少なくとも1つを含むことを決定することと、

それに応答して、前記本番環境によって、前記ターゲットニューラルネットワークを使用して前記オンラインネットワーク入力を処理することなく、前記オンラインネットワーク入力に応答してデフォルト出力を生成し、オンラインネットワーク出力を生成することと、

10

20

30

40

50

を含む動作を実行させる、システム。

【請求項 15】

前記本番環境が展開前の環境システムから前記複数のセグメントを受信する、請求項 14 に記載のシステム。

【請求項 16】

前記 1 つまたは複数の不合格のテスト入力を識別することが、
テスト入力ごとに前記ターゲットニューラルネットワークによって生成された前記 1 つまたは複数のテスト出力に自然言語処理技法を適用することを含む、
請求項 14 または 15 に記載のシステム。

【請求項 17】

前記 1 つまたは複数の不合格のテスト入力を識別することが、
1 つまたは複数の基準に対してテスト入力ごとに前記ターゲットニューラルネットワークによって生成された前記 1 つまたは複数のテスト出力を評価することを含む、
請求項 14 乃至 16 の何れか一項に記載のシステム。

【請求項 18】

前記 1 つまたは複数の基準が、テスト出力が攻撃的または偏ったコンテンツ、トレーニング中に処理される保護可能な情報、あるいは個人の連絡先情報を含むべきでないということを指定する、請求項 16 に記載のシステム。

【請求項 19】

前記 1 つまたは複数の不合格のテスト入力に基づいて前記セグメントを生成することが

、
前記 1 つまたは複数の不合格のテスト入力において発生する様々なセグメントを複数のテキストクラスタにグループ化するためにテキストクラスタリングアルゴリズムを適用することと、

前記複数のテキストクラスタの 1 つに含まれるセグメントに基づいて、前記セグメントを生成することと、を含む

請求項 14 乃至 18 の何れか一項に記載のシステム。

【請求項 20】

前記動作がさらに、

前記本番環境に展開される前記ターゲットニューラルネットワークを調整するために前記複数のセグメントを使用することを含む、

請求項 14 乃至 19 の何れか一項に記載のシステム。

10

20

30

40

50