

(19)日本国特許庁(JP)

## (12)特許公報(B2)

(11)特許番号  
特許第7529159号  
(P7529159)

(45)発行日 令和6年8月6日(2024.8.6)

(24)登録日 令和6年7月29日(2024.7.29)

(51)国際特許分類		F I			
G 0 6 N	20/00	(2019.01)	G 0 6 N	20/00	1 3 0
G 0 6 N	3/094	(2023.01)	G 0 6 N	3/094	

請求項の数 5 (全12頁)

(21)出願番号	特願2023-528902(P2023-528902)	(73)特許権者	000004226 日本電信電話株式会社 東京都千代田区大手町一丁目5番1号
(86)(22)出願日	令和3年6月17日(2021.6.17)	(74)代理人	110002147 弁理士法人酒井国際特許事務所
(86)国際出願番号	PCT/JP2021/023123	(72)発明者	山田 真徳 東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内
(87)国際公開番号	WO2022/264387	審査官	渡辺 一帆
(87)国際公開日	令和4年12月22日(2022.12.22)		
審査請求日	令和5年10月24日(2023.10.24)		

最終頁に続く

(54)【発明の名称】 学習装置、学習方法、および、学習プログラム

## (57)【特許請求の範囲】

## 【請求項1】

Adversarial Exampleを含む入力データのラベルを予測するためのモデルの学習データを取得するデータ取得部と、

前記モデルのパラメータにノイズを加えた場合とノイズを加えなかった場合とで、前記モデルにおけるloss値のKLダイバージェンスが最大になるようなノイズをパラメータに加え、前記パラメータに対するloss landscapeを平らにしたloss関数と、前記Adversarial Exampleを含む学習データとを用いて、前記モデルの学習を行う学習部とを備えることを特徴とする学習装置。

## 【請求項2】

前記学習部は、  
前記学習データを用いて、前記loss関数により算出されるlossを最小化する前記モデルのパラメータを求める  
ことを特徴とする請求項1に記載の学習装置。

## 【請求項3】

学習された前記モデルを用いて、入力データのラベルを予測する予測部をさらに備えることを特徴とする請求項1に記載の学習装置。

## 【請求項4】

学習装置により実行される学習方法であって、  
Adversarial Exampleを含む入力データのラベルを予測するためのモデルの学習デー

10

20

タを取得する工程と、

前記モデルのパラメータにノイズを加えた場合とノイズを加えなかった場合とで、前記モデルにおけるloss値のKLダイバージェンスが最大になるようなノイズをパラメータに加え、前記パラメータに対するloss landscapeを平らにしたloss関数と、前記Adversarial Exampleを含む学習データとを用いて、前記モデルの学習を行う工程と

を含むことを特徴とする学習方法。

【請求項5】

Adversarial Exampleを含む入力データのラベルを予測するためのモデルの学習データを取得する工程と、

前記モデルのパラメータにノイズを加えた場合とノイズを加えなかった場合とで、前記モデルにおけるloss値のKLダイバージェンスが最大になるようなノイズをパラメータに加え、前記パラメータに対するloss landscapeを平らにしたloss関数と、前記Adversarial Exampleを含む学習データとを用いて、前記モデルの学習を行う工程と

をコンピュータに実行させるための学習プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、モデルの学習装置、学習方法、および、学習プログラムに関する。

【背景技術】

【0002】

従来、分類対象のデータにノイズをのせることで、分類器に誤判定をさせるAdversarial Exampleという攻撃がある。このAdversarial Exampleに対する対策として、例えば、Adversarial Exampleを用いてモデル（分類器）の学習を行うAdversarial Trainingがある。

【0003】

しかし、Adversarial Trainingで学習したモデルは汎化性能が低いという問題がある。これは、Adversarial Trainingで学習したモデルのweightに対するloss landscape（loss関数の形）が尖っていることに起因する。そこで、上記のloss landscapeを平らにするため、モデルのlossを最大化する方向にweightにノイズ（摂動）を加える技術がある。

【先行技術文献】

【非特許文献】

【0004】

【文献】Diederik P. Kingma, Max Welling, "Auto-Encoding Variational Bayes"、[2021年6月4日検索]、インターネット<URL:https://arxiv.org/pdf/1312.6114.pdf>

【文献】Dongxian Wu, Shu-Tao Xia, Yisen Wang, "Adversarial Weight Perturbation Helps Robust Generalization"、[2021年6月4日検索]、インターネット<URL:https://arxiv.org/pdf/2004.05884>

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかし、上記の技術は、ノイズがのっていないデータに対する予測性能が低下するという問題がある。そこで、本発明は、前記した問題を解決し、Adversarial exampleに対する頑健さを確保しつつ、ノイズがのっていないデータに対しても精度よく予測できるモデルを学習することを課題とする。

【0006】

前記した課題を解決するため、本発明は、Adversarial Exampleを含む入力データのラベルを予測するためのモデルの学習データを取得するデータ取得部と、前記モデルのパラメータにノイズを加えた場合とノイズを加えなかった場合とで、前記モデルにおけるlos

10

20

30

40

50

s値のKLダイバージェンスが最大になるようなノイズをパラメータに加え、前記パラメータに対するloss landscapeを平らにしたloss関数と、前記Adversarial Exampleを含む学習データとを用いて、前記モデルの学習を行う学習部とを備えることを特徴とする。

【発明の効果】

【0007】

本発明によれば、Adversarial Exampleに対する頑健さを確保しつつ、ノイズがのっていないデータに対しても精度よく予測できるモデルを学習することができる。

【図面の簡単な説明】

【0008】

【図1】図1は、学習装置の構成例を示す図である。

10

【図2】図2は、式(10)におけるMAX vを求めるためには、Fisher情報行列Gの最大固有値に対応する固有ベクトルhを求めればよい理由を説明するための式である。

【図3】図3は、学習装置の処理手順の例を示すフローチャートである。

【図4】図4は、学習装置の処理手順の例を示すフローチャートである。

【図5】図5は、学習装置の適用例を説明するための図である。

【図6】図6は、学習装置により学習されたモデルに対する実験結果を示す図である。

【図7】図7は、学習プログラムを実行するコンピュータの構成例を示す図である。

【発明を実施するための形態】

【0009】

以下、図面を参照して、本発明の実施の形態(実施形態)を説明する。なお、本発明は以下に説明する実施形態に限定されない。

20

【0010】

[学習装置の概要]

本実施形態の学習装置は、Adversarial Example(ノイズが付加されたデータ)を含むデータを用いて、入力されたデータのラベルを予測するモデルの学習を行う。ここで、学習装置は、モデルの学習に用いる損失関数(loss関数)として、モデルのパラメータにノイズを加えた場合とノイズを加えなかった場合とで、モデルにおけるloss値のKLダイバージェンスが最大になるようなノイズをパラメータに加え、パラメータに対するloss landscapeを平らにしたloss関数を用いる。

【0011】

30

これにより、学習装置は、Adversarial Exampleに対する頑健さを確保しつつ、ノイズがのっていないデータに対しても精度よくラベルを予測できるモデルを学習することができる。

【0012】

[学習装置の構成例]

図1を用いて、学習装置10の構成例を説明する。学習装置10は、例えば、入力部11、出力部12、通信制御部13、記憶部14、および、制御部15を備える。

【0013】

入力部11は、各種データの入力を受け付けるインターフェースである。例えば、入力部11は、後述する学習処理および予測処理に用いるデータの入力を受け付ける。出力部12は、各種データの出力を行うインターフェースである。例えば、出力部12は、制御部15により予測されたデータのラベルを出力する。

40

【0014】

通信制御部13は、NIC(Network Interface Card)等で実現され、ネットワークを介したサーバ等の外部の装置と制御部15との通信を制御する。例えば、通信制御部13は、学習対象のデータを管理する管理装置等と制御部15との通信を制御する。

【0015】

記憶部14は、RAM(Random Access Memory)、フラッシュメモリ(Flash Memory)等の半導体メモリ素子、または、ハードディスク、光ディスク等の記憶装置によって実現され、後述する学習処理により学習されたモデルのパラメータ等が記憶される。

50

## 【 0 0 1 6 】

制御部 1 5 は、例えば、CPU (Central Processing Unit) 等を用いて実現され、記憶部 1 4 に記憶された処理プログラムを実行する。これにより、制御部 1 5 は、図 1 に例示するように、取得部 1 5 a、学習部 1 5 b および予測部 1 5 c として機能する。

## 【 0 0 1 7 】

取得部 1 5 a は、後述する学習処理および予測処理に用いるデータを、入力部 1 1 あるいは通信制御部 1 3 を介して取得する。

## 【 0 0 1 8 】

学習部 1 5 b は、Adversarial Example を含むデータを学習データとして用いて、入力されたデータのラベルを予測するモデルの学習を行う。ここで、学習部 1 5 b は、モデルの学習に用いる loss 関数として、モデルのパラメータにノイズを加えた場合とノイズを加えなかった場合とで、モデルにおける loss 値の KL ダイバージェンスが最大になるようなノイズをパラメータに加えて、パラメータに対する loss landscape を平らにした loss 関数を用いる。

10

## 【 0 0 1 9 】

ここで、学習部 1 5 b による、モデルの学習方法の基本的な考え方を説明する。例えば、学習対象のモデルは、データ  $x$  のラベル  $y$  の確率分布を表すモデルであり、パラメータを用いて、式 ( 1 ) により表される。なお、式 ( 1 ) における  $f$  は、モデルが出力するラベルを表すベクトルである。

## 【 0 0 2 0 】

## 【数 1】

$$p_{\theta}(y_k|x) = \frac{\exp f_k(x; \theta)}{\sum_i \exp f_i(x; \theta)} \quad \dots (1)$$

20

## 【 0 0 2 1 】

そして、学習部 1 5 b は、式 ( 2 ) で表される loss 関数の値が小さくなるように、モデルのパラメータ を決定することにより、モデルの学習を行う。ここで、 $p(y|x)$  は、真の確率を表す。

30

## 【 0 0 2 2 】

## 【数 2】

$$l(x, y; \theta) = p(y|x) \log p_{\theta}(y|x) \quad \dots (2)$$

## 【 0 0 2 3 】

ここで、学習部 1 5 b は、データ  $x$  にノイズ が乗せられた Adversarial Example ( 式 ( 3 ) 参照) に対しても正しくラベルを予測できるようにモデルの学習を行う。つまり、学習部 1 5 b は、式 ( 4 ) に示す Adversarial Training を行う。

## 【 0 0 2 4 】

## 【数 3】

$$\max_{\eta} E_{x, y \sim p(x, y)} [l(x + \eta, y; \theta)] \quad \dots (3)$$

40

## 【数 4】

$$\min_{\theta} \left( \max_{\eta} E_{x, y \sim p(x, y)} [l(x + \eta, y; \theta)] \right) \quad \dots (4)$$

50

## 【 0 0 2 5 】

ここで、従来、上記のモデルのAdversarial Training (AT)において、モデルの汎化性能を上げるためにweight (モデルのパラメータ) にノイズ (摂動) を加えることにより、weightに対するloss landscapeを平らにする方法がある。この方法 (Adversarial Weight Perturbation、AWP) における、loss関数は、式 (5) および式 (6) により表される。なお、 $w$  (weight) は、学習対象のモデルのパラメータであり、前記した に相当する。  $\alpha$  は、ノイズ ( $v$ ) の大きさを調節する係数で、 $w$  のフロベニウスノルムではかったときのスケールに合うような値が設定される。つまり、パラメータはスケール不変性があるため、  $\alpha$  はそのスケールの変化を吸収する役割を担う。

## 【 0 0 2 6 】

## 【数 5】

$$\rho(w) = \frac{1}{N} \sum_n \max_{\|x'_n - x_n\|_p \leq \epsilon} l(x_n, y_n, w) \quad \dots (5)$$

## 【数 6】

$$\min_w \max_v \{ \rho(w) + \rho(w + \alpha \odot v) - \rho(w) \} = \min_w \max_v \rho(w + \alpha \odot v) \quad \dots (6)$$

$\alpha$ : ノイズの大きさを調整する係数

$w$ : モデルのパラメータ

$v$ : モデルのパラメータに対するノイズ

$\odot$ : アダマール積

## 【 0 0 2 7 】

ここで、filter normalizationで可視化されたweight loss landscapeを平らにしたいので、  $\alpha$  は、フィルターごとの $w$ のスケールにあったノイズ (摂動) になるよう以下の式 (7) のように定義される。なお、 $k$  はフィルターのインデックスである。

## 【 0 0 2 8 】

## 【数 7】

$$\alpha^k = \frac{\|w^k\|}{\|\nabla_v \rho(w^k + v^k)\|} \quad \dots (7)$$

## 【 0 0 2 9 】

よって、 $v$  を最大にするための更新式は、式 (8) のように表される。

## 【 0 0 3 0 】

## 【数 8】

$$v \leftarrow \prod_{\gamma} \left( v + \eta_2 \alpha \odot \nabla_v \frac{1}{N} \sum_n l(x'_n, y_n, w + v) \right) \quad \dots (8)$$

## 【 0 0 3 1 】

ちなみに先行研究では、上記の $v$ を最大にするための更新は1回で充分であることが確認されている。また、 $w$ の更新式は、以下の式 (9) のように表される。

## 【 0 0 3 2 】

10

20

30

40

50

【数 9】

$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v} - \eta_3 \nabla_{\mathbf{v}+\mathbf{w}} \frac{1}{N} \sum_n l(\mathbf{x}'_n, \mathbf{y}_n, \mathbf{w} + \mathbf{v}) - \mathbf{v} \quad \dots (9)$$

【0033】

ここで、上記のAWPでは、loss値を最大化するようにwにノイズを加えていたが、学習部15bは、loss値のKLダイバージェンスを最大化するように、wにノイズを加える。このloss関数は、以下の式(10)により表される。なお、式(10)における(w)は、式(5)に示す(w)に相当する。

【0034】

【数10】

$$\min_{\mathbf{w}} \max_{\mathbf{v}} \{ \rho(\mathbf{w}) + D_{\text{KL}}(\rho(\mathbf{w}) \parallel \rho(\mathbf{w} + \alpha \odot \mathbf{v})) \}$$

$$\alpha^k = \frac{\|\mathbf{w}^k\|}{\|\mathbf{v}_{\text{KL}}\|} = \|\mathbf{w}^k\| \quad \dots (10)$$

$\mathbf{v}_{\text{KL}}$ : KLダイバージェンスを最大化するよう作ったノイズ

【0035】

式(10)におけるMAX vを求めるためには、Fisher情報行列Gの最大固有値に対応する固有ベクトルhを求めればよい。その理由となる式を図2に示す。

【0036】

よって、vを最大にするための更新式は、式(11)のようになる。

【0037】

【数11】

$$\mathbf{v} \leftarrow \prod_{\gamma} (\mathbf{v} + \eta_2 \mathbf{h}_1) \quad \dots (11)$$

【0038】

なお、上記のFisher情報行列は巨大なので、負数に固有値分割を行うと時間がかかりすぎる。そのため、例えば、power iterationを使い最大固有値を計算する。また、Fisher情報行列を計算するとき、

【数12】

$$\frac{\partial}{\partial \mathbf{w}} \log p_{\theta}(y|x)$$

を計算する必要があるが、これは入力より出力の次元が大きいため、通常のDeep Learningで用いられるback propagationを使うと計算効率がよくない。そのためforward propagationで勾配を計算したいが、Pytorch等の既存の深層学習ライブラリにforward propagationのモードは用意されていない。そのため、以下の文献1に記載のROP trickを使ってforward propagationを実現する。

【0039】

(文献1) [Adding functionality] Hessian and Fisher Information vector products、<https://discuss.pytorch.org/t/adding-functionality-hessian-and-fisher-information-vector-products/23295/2>

【0040】

10

20

30

40

50

学習部 15 b は、Adversarial Example を含む学習データと、上記の loss 関数とを用いて、入力されたデータのラベルを予測するモデルの学習を行う。つまり、学習部 15 b は、学習データを用いて、上記の loss 関数により算出される loss を最小化するようなモデルのパラメータを求める。

#### 【0041】

予測部 15 c は、学習されたモデルを用いて、入力データのラベルを予測する。例えば、予測部 15 c は、学習されたパラメータを上記式(1)に適用することにより、新たに取得されたデータの各ラベルの確率を算出し、最も確率が高いラベルを出力する。これにより、学習装置 10 は、例えば、入力データが Adversarial Example であった場合にも、正しいラベルを出力することができる。

10

#### 【0042】

##### [学習処理]

次に、図3を参照して、学習装置 10 による学習処理手順の例について説明する。図3に示す処理は、例えば、学習処理の開始を指示する操作入力があったタイミングで開始される。

#### 【0043】

まず、取得部 15 a が、Adversarial Example を含む学習データを取得する(S1)。次に、学習部 15 b が、学習データと、loss 関数とを用いて、入力データのラベルの確率分布を表すモデルを学習する(S2)。なお、この loss 関数は、上記の通り、モデルのパラメータにノイズを加えた場合とノイズを加えなかった場合とで、当該モデルにおける loss 値の KL ダイバージェンスが最大になるようなノイズをパラメータに加え、パラメータに対する loss landscape を平らにした loss 関数である。学習部 15 b は、S2 で学習されたモデルのパラメータを記憶部 14 に記憶する。

20

#### 【0044】

##### [予測処理]

次に、図4を参照して、学習装置 10 による入力データのラベルの予測処理の例について説明する。図4に示す処理は、例えば、予測処理の開始を指示する操作入力があったタイミングで開始される。

#### 【0045】

まず、取得部 15 a は、ラベルの予測対象のデータを取得する(S11)。次に、予測部 15 c は、学習部 15 b により学習されたモデルを用いて、S11 で取得されたデータのラベルを予測する(S12)。例えば、予測部 15 c は、学習されたパラメータを上記の式(1)に適用することにより、S11 で取得されたデータ  $x'$  の  $p(x')$  を算出し、最も確率が高いラベルを出力する。これにより、例えば、データ  $x'$  が Adversarial Example であった場合でも、学習装置 10 は、正しいラベルを出力することができる。

30

#### 【0046】

##### [学習装置の適用例]

上記の学習装置 10 を、データの異常検知に適用してもよい。この場合の適用例を、図5を参照しながら説明する。ここでは、前記した予測部 15 c の機能が、検知装置 20 に装備される場合を例に説明する。

40

#### 【0047】

例えば、学習装置 10 は、データ取得装置から取得した教師データ(学習データ)と、前記した loss 関数とを用いて、モデルの学習(Adversarial Training)を行う。その後、検知装置 20 は、データ取得装置から新たなデータ  $x'$  を取得すると、学習済みモデルを用いて、データ  $x'$  の  $p(x')$  を算出する。そして、検知装置 20 は、確率が最も高いラベルに基づき、データ  $x'$  が異常なデータか否かのレポートを出力する。

#### 【0048】

##### [実験結果]

次に、本実施形態の学習装置 10 により学習されたモデルによる、ラベルの予測精度の評価実験の結果を図6に示す。本実験では、本実施形態の学習装置 10 により学習された

50

モデルについて、robust accとnatural accを評価した。

【 0 0 4 9 】

robust accは、Adversarial Exampleがのったデータの分類精度（データのラベルの予測精度）を示す値である。また、natural accは、ノイズなしのデータの分類精度を示す値である。robust accもnatural accも、0～100の値をとる。比較対象としたのは、ATにより学習したモデルと、AWPにより学習したモデルである。実験条件は以下の通りである。

【 0 0 5 0 】

画像のデータセット:Cifar10

Deep learning model: Resnet18

Adversarial Example: PGD

PGDのパラメータ: eps=8/255, train\_iter=7, eval\_iter=20, eps\_iter=0.01, rand\_init=True, clip\_min=0.0, clip\_max=1.0

【 0 0 5 1 】

図6に示すように、学習装置10により学習されたモデルは、ATにより学習されたモデルと比べ、robust acc、natural accともに高い値になっている。また、本実施形態の学習装置10により学習されたモデルは、AWPにより学習されたモデルと比べ、robust accは少し低い値となっているが、natural accは大幅に高い値になっている。

【 0 0 5 2 】

したがって、学習装置10により学習されたモデルは、Adversarial exampleに対する頑健さを確保しつつ、ノイズがのっていないデータに対しても精度よく予測できるモデルであることが確認できた。

【 0 0 5 3 】

[ システム構成等 ]

また、図示した各部の各構成要素は機能概念的なものであり、必ずしも物理的に図示のように構成されていることを要しない。すなわち、各装置の分散・統合の具体的な形態は図示のものに限られず、その全部又は一部を、各種の負荷や使用状況等に応じて、任意の単位で機能的又は物理的に分散・統合して構成することができる。さらに、各装置にて行われる各処理機能は、その全部又は任意の一部が、CPU及び当該CPUにて実行されるプログラムにて実現され、あるいは、ワイヤードロジックによるハードウェアとして実現され得る。

【 0 0 5 4 】

また、前記した実施形態において説明した処理のうち、自動的に行われるものとして説明した処理の全部又は一部を手動的に行うこともでき、あるいは、手動的に行われるものとして説明した処理の全部又は一部を公知の方法で自動的に行うこともできる。この他、上記文書中や図面中で示した処理手順、制御手順、具体的名称、各種のデータやパラメータを含む情報については、特記する場合を除いて任意に変更することができる。

【 0 0 5 5 】

[ プログラム ]

前記した学習装置10は、パッケージソフトウェアやオンラインソフトウェアとしてプログラムを所望のコンピュータにインストールさせることによって実装できる。例えば、上記のプログラムを情報処理装置に実行させることにより、情報処理装置を学習装置10として機能させることができる。ここで言う情報処理装置には、デスクトップ型又はノート型のパーソナルコンピュータが含まれる。また、その他にも、情報処理装置にはスマートフォン、携帯電話機やPHS(Personal Handyphone System)等の移動体通信端末、さらには、PDA(Personal Digital Assistant)等の端末等がその範疇に含まれる。

【 0 0 5 6 】

また、学習装置10は、ユーザが使用する端末装置をクライアントとし、当該クライアントに上記の処理に関するサービスを提供するサーバ装置として実装することもできる。この場合、サーバ装置は、Webサーバとして実装することとしてもよいし、アウトソーシ

10

20

30

40

50

ングによって上記の処理に関するサービスを提供するクラウドとして実装することとしてもかまわない。

【0057】

図7は、学習プログラムを実行するコンピュータの一例を示す図である。コンピュータ1000は、例えば、メモリ1010、CPU1020を有する。また、コンピュータ1000は、ハードディスクドライブインタフェース1030、ディスクドライブインタフェース1040、シリアルポートインタフェース1050、ビデオアダプタ1060、ネットワークインタフェース1070を有する。これらの各部は、バス1080によって接続される。

【0058】

メモリ1010は、ROM(Read Only Memory)1011及びRAM(Random Access Memory)1012を含む。ROM1011は、例えば、BIOS(Basic Input Output System)等のブートプログラムを記憶する。ハードディスクドライブインタフェース1030は、ハードディスクドライブ1090に接続される。ディスクドライブインタフェース1040は、ディスクドライブ1100に接続される。例えば磁気ディスクや光ディスク等の着脱可能な記憶媒体が、ディスクドライブ1100に挿入される。シリアルポートインタフェース1050は、例えばマウス1110、キーボード1120に接続される。ビデオアダプタ1060は、例えばディスプレイ1130に接続される。

【0059】

ハードディスクドライブ1090は、例えば、OS1091、アプリケーションプログラム1092、プログラムモジュール1093、プログラムデータ1094を記憶する。すなわち、上記の学習装置10が実行する各処理を規定するプログラムは、コンピュータにより実行可能なコードが記述されたプログラムモジュール1093として実装される。プログラムモジュール1093は、例えばハードディスクドライブ1090に記憶される。例えば、学習装置10における機能構成と同様の処理を実行するためのプログラムモジュール1093が、ハードディスクドライブ1090に記憶される。なお、ハードディスクドライブ1090は、SSD(Solid State Drive)により代替されてもよい。

【0060】

また、上述した実施形態の処理で用いられるデータは、プログラムデータ1094として、例えばメモリ1010やハードディスクドライブ1090に記憶される。そして、CPU1020が、メモリ1010やハードディスクドライブ1090に記憶されたプログラムモジュール1093やプログラムデータ1094を必要に応じてRAM1012に読み出して実行する。

【0061】

なお、プログラムモジュール1093やプログラムデータ1094は、ハードディスクドライブ1090に記憶される場合に限らず、例えば着脱可能な記憶媒体に記憶され、ディスクドライブ1100等を介してCPU1020によって読み出されてもよい。あるいは、プログラムモジュール1093及びプログラムデータ1094は、ネットワーク(LAN(Local Area Network)、WAN(Wide Area Network)等)を介して接続される他のコンピュータに記憶されてもよい。そして、プログラムモジュール1093及びプログラムデータ1094は、他のコンピュータから、ネットワークインタフェース1070を介してCPU1020によって読み出されてもよい。

【符号の説明】

【0062】

- 10 学習装置
- 11 入力部
- 12 出力部
- 13 通信制御部
- 14 記憶部
- 15 制御部

10

20

30

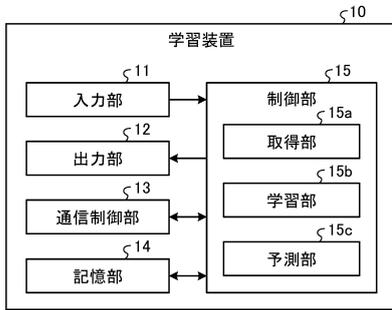
40

50

- 1 5 a 取得部
- 1 5 b 学習部
- 1 5 c 予測部
- 2 0 検知装置

【図面】

【図 1】



【図 2】

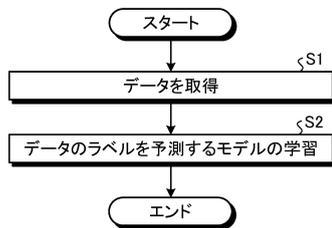
$$\begin{aligned}
 D_{\text{KL}}(p(y|x;w) \| p(y|x;w + \alpha \circ v)) &= \int p(y|x;w) p(x) (\log p(y|x;w) - \log p(y|x;w + \alpha \circ v)) dy dx \\
 &= \int p(y|x;w) p(x) \left( \log p(y|x;w) - \log p(y|x;w) - \sum_{\theta_i} \frac{\partial}{\partial \theta_i} \log p(y|x;w) (\alpha^{\theta_i}) - \sum_{\theta_i} \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p(y|x;w) (\alpha^{\theta_i}, \alpha^{\theta_j}) + O(\alpha^3) \right) dx dy \\
 &\approx \int p(x) p(y|x;w) \left( - \sum_{\theta_i} \frac{\partial}{\partial \theta_i} \log p(y|x;w) (\alpha^{\theta_i}) (\alpha^{\theta_i}) \right) dx dy \\
 &= \int p(x) p(y|x;w) \left( \sum_{\theta_i} \frac{\partial}{\partial \theta_i} \log p(y|x;w) \frac{\partial}{\partial \theta_i} \log p(y|x;w) (\alpha^{\theta_i}) (\alpha^{\theta_i}) \right) dx dy \\
 &= \int p(x) \sum_{\theta_i} (\alpha^{\theta_i})^2 C_{ij} (\alpha^{\theta_i}) dx \quad , \Delta \theta_i = C_{ij} \theta_i \\
 &= \int p(x) \sum_{\theta_i} v_i C_{ij}^{(v)} dx \\
 &\approx \int p(x) \lambda_i \| \mathbf{h}_i \|^2 dx \quad , \lambda_i \mathbf{h}_i = C_{ij} \theta_i \\
 C_{ij} &\equiv \int p(y|x;w) \frac{\partial}{\partial w_i} \log p(y|x;w) \frac{\partial}{\partial w_j} \log p(y|x;w) dy \\
 C_{ij}^{(v)} &\equiv \int p(y|x;w) \alpha_i^k \frac{\partial}{\partial w_i} \log p(y|x;w) \alpha_j^k \frac{\partial}{\partial w_j} \log p(y|x;w) dy
 \end{aligned}$$

10

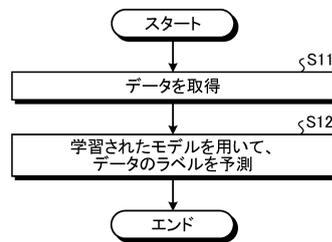
20

30

【図 3】

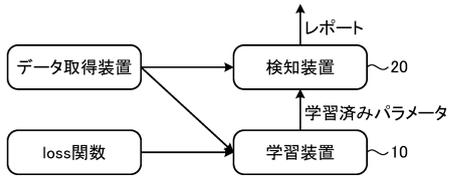


【図 4】



40

【図 5】

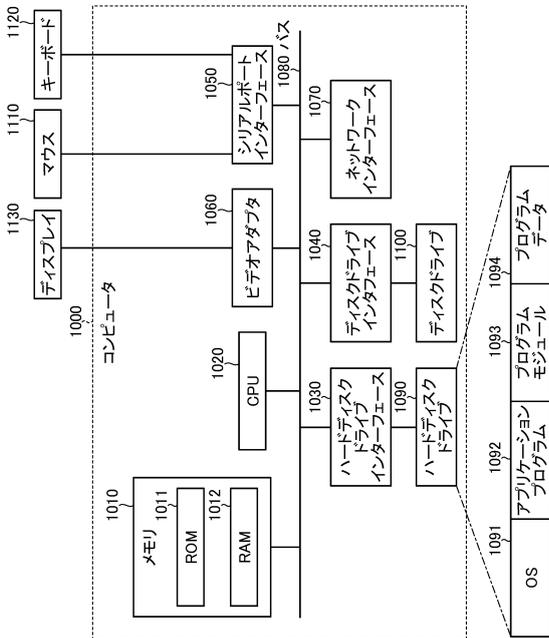


【図 6】

	robust acc	natural acc
AT	52.79	85.57
AWP	55.39	82.00
本実施形態	54.80	86.57

10

【図 7】



20

30

40

50

---

フロントページの続き

- (56)参考文献      WU, Dongxian et al. , "Adversarial Weight Perturbation Helps Robust Generalization" , arXiv.org [online] , 2020年 , pp. 1-20 , [retrieved on 2021.08.16], Retrieved from the Internet:  
                         URL: <https://arxiv.org/abs/2004.05884v2>
- MIYATO, Takeru et al. , "Distributional Smoothing with Virtual Adversarial Training" , arXiv.org [online] , 2016年 , pp. 1-12 , [retrieved on 2021.08.16], Retrieved from the Internet:  
                         URL: <https://arxiv.org/abs/1507.00677v9>
- (58)調査した分野 (Int.Cl. , D B 名)
- G 0 6 N    2 0 / 0 0 - 2 0 / 2 0
- G 0 6 N    3 / 0 2 - 3 / 1 0