



(12) **Gebrauchsmusterschrift**

(21) Aktenzeichen: **20 2017 102 238.2**
(22) Anmeldetag: **13.04.2017**
(47) Eintragungstag: **12.05.2017**
(45) Bekanntmachungstag im Patentblatt: **22.06.2017**

(51) Int Cl.: **G05B 13/02 (2006.01)**
G06N 3/02 (2006.01)
G06F 15/18 (2006.01)

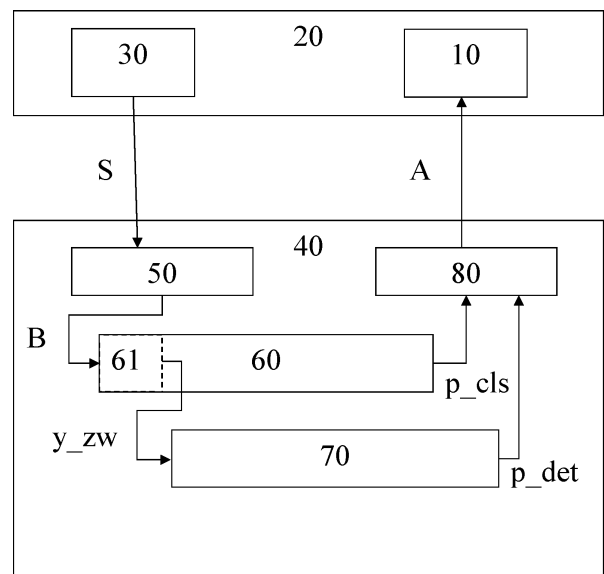
(73) Name und Wohnsitz des Inhabers:
Robert Bosch GmbH, 70469 Stuttgart, DE

Die folgenden Angaben sind den vom Anmelder eingereichten Unterlagen entnommen

(54) Bezeichnung: **Aktorsteuerungssystem**

(57) Hauptanspruch: Aktorsteuerungssystem (40) zum Steuern eines Aktors (10), umfassend mindestens einen Computer und mindestens ein maschinenlesbares Speichermedium, auf dem Befehle gespeichert sind, die beim Ausführen durch den mindestens einen Computer bewirken, dass der mindestens eine Computer ein Verfahren mit den folgenden Schritten ausführt:

- Empfangen eines ermittelten Beobachtungswerts (B), der einen Zustand eines Aktorsystems umfassend den Aktor (10) und eine Umgebung (20) des Aktors (10) charakterisiert;
- Ermitteln eines ersten Ausgabewerts (p_cls) eines ersten maschinellen Lernsystems (60) abhängig von dem Beobachtungswert (B), wobei der erste Ausgabewert (p_cls) wenigstens einen Teil des Beobachtungswerts (B) charakterisiert,
- Ermitteln eines zweiten Ausgabewerts (p_det) eines zweiten maschinellen Lernsystems (70), wobei der zweite Ausgabewert (p_det) eine Wahrscheinlichkeit charakterisiert, dass der Beobachtungswert (B) derart manipuliert wurde, dass der erste Ausgabewert (p_cls) den wenigstens einen Teil des ersten Beobachtungswerts (B) nicht korrekt charakterisiert, und
- Ansteuern des Aktors (10) abhängig von dem ermittelten ersten Ausgabewert (p_cls) und dem ermittelten zweiten Ausgabewert (p_det), wobei das zweite maschinelle Lernsystem (70) den zweiten Ausgabewert (p_det) abhängig von Zwischenergebnissen (y_zw) des ersten maschinellen Lernsystems (60) ermittelt.



Beschreibung

[0001] Die Erfindung betrifft ein Aktorsteuerungssystem.

Stand der Technik

[0002] Aus der DE 10 2005 050 577 A1 ist ein neuronales Netz für eine Steuerungsvorrichtung bekannt. Die Erfindung prüft ein neuronales Netz 1 für eine Steuerungsvorrichtung. Das neuronale Netz weist eine Mehrzahl erster Neuronen N_1, N_2, \dots, N_n in einer ersten Schicht und einem zweiten Neuron M in einer auf die erste Schicht folgenden zweiten Schicht auf. Aus einer vorgegebenen Mehrzahl von Testsignal-Kombinationen wird jede Testsignal-Kombination ausgewählt. Jede Testsignal-Kombination ordnet jedem ersten Neuron N_1, N_2, \dots, N_n einen Test-Eingangssignalvektor ut_1, ut_2, \dots, ut_k zu, welcher entweder ein Nullsignal ist oder das zugehörige erste Neuron N_1, N_2, \dots, N_n derart sättigt, dass das erste Neuron N_1, N_2, \dots, N_n einen unteren Sättigungswert ϕ_{min} ausgibt, oder das zugehörige erste Neuron N_1, N_2, \dots, N_n derart sättigt, dass das erste Neuron N_1, N_2, \dots, N_n einen oberen Sättigungswert ausgibt. Die Testsignal-Kombination wird an die ersten Neuronen N_1, N_2, \dots, N_n angelegt und das Ausgangssignal p des zweiten Neurons M wird erfasst. Ein Teil-Prüfungssignal wird gespeichert, wenn das Ausgangssignal p größer als ein vorgegebener Schwellwert ist. Ein positives Gesamt-Prüfsignal wird ausgegeben, nachdem jede der Testsignal-Kombinationen angelegt wird und wenn kein Teil-Prüfsignal gespeichert ist für die vorgegebene Mehrzahl der Testsignal-Kombinationen gespeichert wird

Vorteil der Erfindung

[0003] Ein Aktorsteuerungssystem nach Anspruch 1 hat demgegenüber dem Vorteil, dass es ein maschinelles Lernverfahren aufweist, das besonders robust gegenüber Irreführungsbeispielen (engl. „Adversarial Examples“) ist. Adversarial Examples sind geringfügig manipulierte Eingangsdaten des maschinellen Lernverfahrens (die bei Bilddaten den unmanipulierten Eingangsdaten so ähnlich sind, dass sie für menschliche Experten praktisch nicht zu unterscheiden sind), die zu einer erheblichen Änderung der Ausgabe des maschinellen Lernverfahrens führen.

[0004] Vorteilhafte Weiterbildungen sind Gegenstand der abhängigen Ansprüche.

Offenbarung der Erfindung

[0005] In einem ersten Aspekt betrifft die Erfindung daher ein Aktorsteuerungssystem zum Steuern eines Aktors mit den Merkmalen des unabhängigen Anspruchs 1. Das Aktorsteuerungssystem ist eingerichtet, einen Beobachtungswert zu empfangen und mit-

tels eines ersten und zweiten maschinellen Lernsystems eine Ansteuerung für den Aktor zu generieren. Der ermittelte Beobachtungswert kann insbesondere ein Ausgangssignal eines Sensors umfassen oder abhängig von einem Ausgangssignal eines Sensors ermittelt worden sein.

[0006] In einem weiterführenden Aspekt ist hierbei das erste maschinelle Lernsystem ein erstes neuronales Netzwerk mit einer ersten Verkettung neuronaler Schichten, und das zweite maschinelle Lernsystem ein zweites neuronales Netzwerk mit einer zweiten Verkettung neuronaler Schichten, wobei die neuronalen Schichten der ersten und zweiten Verkettung neuronaler Schichten in Signalverarbeitungsrichtung bis zu einer letzten gemeinsamen Zwischenschicht gemeinsame neuronale Schichten sind und in Signalverarbeitungsrichtung ab dieser letzten gemeinsamen Zwischenschicht separat (d.h. untereinander nicht verknüpfte) sind.

[0007] Dies bedeutet, dass erstes und zweites neuronales Netz ein gemeinsames neuronales Netz bilden, das eine Verkettung gemeinsamer neuronaler Schichten umfasst, wobei diese Verkettung den Beobachtungswert empfängt. Das neuronale Netz ist hierbei derart aufgebaut ist, dass sowohl die Ermittlung des ersten Ausgabewerts als auch die Ermittlung des zweiten Ausgabewerts ausschließlich („ausschließlich“ im Sinne von „bei fixierten Parametern des neuronalen Netzes ausschließlich“) von Ausgabewerten der gemeinsamen neuronalen Schichten abhängig sind, und dass das erste und zweite neuronale Netz ferner so aufgebaut sind, dass abhängig von den Ausgabewerten der gemeinsamen neuronalen Schichten ein erstes Zwischenergebnis und ein zweites Zwischenergebnis ermittelt werden, wobei der erste Ausgabewert abhängig vom ersten Zwischenergebnis und unabhängig vom zweiten Zwischenergebnis ermittelt wird, und wobei der zweite Ausgabewert abhängig vom zweiten Zwischenergebnis und unabhängig vom ersten Zwischenergebnis ermittelt wird.

[0008] Nachfolgend werden Ausführungsformen der Erfindung unter Bezugnahme auf die beiliegenden Zeichnungen näher erläutert. In den Zeichnungen zeigen:

[0009] Fig. 1 schematisch eine Interaktion zwischen Aktor und Aktorsteuerungssystem;

[0010] Fig. 2 schematisch eine Interaktion zwischen Trainingssystem und Aktorsteuerungssystem;

[0011] Fig. 3 eine Ausführungsform eines Trainingsverfahrens;

[0012] Fig. 4 eine weitere Ausführungsform des Trainingsverfahrens;

[0013] Fig. 5 einen beispielhaften Aufbau des ersten und zweiten maschinellen Lernsystems;

[0014] Fig. 6 schematisch den Aufbau eines Residualblocks;

[0015] Fig. 7 schematisch den Aufbau eines multiplen Residualblocks.

Beschreibung der Ausführungsbeispiele

[0016] Fig. 1 zeigt einen Aktor **10** in seiner Umgebung **20** in Interaktion mit einem Aktorsteuerungssystem **40**. Aktor **10** und Umgebung **20** werden gemeinschaftlich nachfolgend auch als Aktorsystem bezeichnet. Ein Zustand des Aktorsystems wird mit einem Sensor **30** erfasst, der auch durch eine Mehrzahl von Sensoren gegeben sein kann. Ein Ausgangssignal S des Sensors **30** wird an das Aktorsteuerungssystem **40** übermittelt. Das Aktorsteuerungssystem **40** ermittelt hieraus ein Ansteuersignal A , welches der Aktor **10** empfängt.

[0017] Bei dem Aktor **10** kann es sich beispielsweise um einen (teil-)autonomen Roboter, beispielsweise ein (teil-)autonomes Kraftfahrzeug handeln. Bei dem Sensor **30** kann es sich beispielsweise um einen oder mehrere Videosensoren und/oder einen oder mehrere Radarsensoren und/oder einen oder mehrere Ultraschallsensoren und/oder einen oder mehrere Positionssensoren (beispielsweise GPS) handeln. Alternativ oder zusätzlich kann der Sensor **30** auch ein Informationssystem umfassen, das eine Information über einen Zustand des Aktorsystems ermittelt, wie beispielsweise ein Wetterinformationssystem, das einen aktuellen oder zukünftigen Zustand des Wetters in der Umgebung **20** ermittelt.

[0018] In einem anderen Ausführungsbeispiel kann es sich bei dem Aktor **10** um einen Fertigungsroboter handeln, bei dem Sensor **30** dann beispielsweise um einen optischen Sensor handelt, der Eigenschaften von Fertigungserzeugnissen des Fertigungsroboters erfasst.

[0019] In einem weiteren Ausführungsbeispiel kann es sich bei dem Aktor **10** um ein Freigabesystem handeln, welches eingerichtet ist, die Aktivität eines Geräts freizugeben oder nicht. Bei dem Sensor **30** kann es sich beispielsweise um einen optischen Sensor (beispielsweise zur Erfassung von Bild- oder Videodaten) handeln, der eingerichtet ist, ein Gesicht zu erfassen. Der Aktor **10** ermittelt abhängig vom Ansteuersignal A ein Freigabesignal, das benutzt werden kann, um abhängig vom Wert des Freigabesignals das Gerät freizugeben. Bei dem Gerät kann es sich beispielsweise um eine physische oder logische Zugangskontrolle handeln. Abhängig vom Wert des Ansteuersignals A kann die Zugangskontrolle dann vorsehen, dass Zugang gewährt wird, oder nicht.

[0020] Das Aktorsteuerungssystem **40** empfängt das Ausgangssignal S des Sensors in einer optionalen Empfangseinheit **50**, die das Ausgangssignal S in einen Beobachtungswert B umwandelt (alternativ kann auch unmittelbar das Ausgangssignal S als Beobachtungswert B übernommen werden). Der Beobachtungswert B kann beispielsweise ein Ausschnitt oder eine Weiterverarbeitung des Ausgangssignals S sein. Der Beobachtungswert B wird einem ersten maschinellen Lernsystem **60** zugeführt.

[0021] Das erste maschinelle Lernsystem **60** ermittelt aus dem Beobachtungswert B einen ersten Ausgabewert p_cls . Das erste maschinelle Lernsystem **60** kann in einer Ausführungsform als Klassifikator eingesetzt werden. In diesem Fall kann der erste Ausgabewert p_cls eine Wahrscheinlichkeit charakterisieren, dass der Beobachtungswert B oder ein Teil des Beobachtungswerts B als einer Klasse (aus einer Mehrzahl von Klassen) zugehörig klassifiziert wird. Der erste Ausgabewert kann beispielsweise eine vektorwertige Größe sein, die für jede Klasse der Mehrzahl von Klassen mittels einer zugeordneten Zahl im Wertebereich $[0; 1]$ angibt, wie hoch die Wahrscheinlichkeit ist, dass der Beobachtungswert der jeweiligen Klasse zuzuordnen ist. Der erste Ausgabewert p_cls kann auch eine skalare Größe sein, beispielsweise ein Bezeichner derjenigen Klasse, deren oben beschriebene Wahrscheinlichkeit den größten Wert annimmt.

[0022] Es ist auch möglich, dass ein Beobachtungswert B eine Vielzahl von Größen umfasst, beispielsweise Pixel eines Bilds. Es ist dann möglich, dass eine semantische Segmentierung des Bilds durchgeführt wird, d.h. dass die oben beschriebene Zuordnung zu Klassen für jede dieser Größen, also beispielsweise jeden Pixel des Bilds, einzeln durchgeführt wird, sodass p_cls auch in diesem Sinne eine vektorwertige Größe sein kann.

[0023] In einer alternativen Ausführungsform ist es auch möglich, dass das erste maschinelle Lernsystem **60** eine Regression ausführt, d.h. der erste Ausgabewert p_cls ist eine (im Rahmen der für die digitale Datenverarbeitung notwendigen Diskretisierung) kontinuierliche Größe.

[0024] Das erste maschinelle Lernsystem **60** umfasst ein Block **61**, der aus dem Beobachtungswert B ein Zwischenergebnis y_zw ermittelt (y_zw kann wieder eine vektorwertige Größe sein). Die weiteren Signalverarbeitungsschritte des ersten maschinellen Lernsystems **60** erfolgen in Abhängigkeit dieses Zwischenergebnisses y_zw , d.h. der erste Ausgabewert p_cls wird abhängig von diesem Zwischenergebnis y_zw ermittelt.

[0025] Das Aktorsteuerungssystem **40** umfasst ferner ein zweites maschinelles Lernsystem **70**. Das

zweite maschinelle Lernsystem **70** empfängt das Zwischenergebnis y_{zw} und ermittelt, hieraus einen zweiten Ausgabewert p_{det} , der beispielsweise eine Zahl im Wertebereich $[0; 1]$ sein kann und eine Wahrscheinlichkeit charakterisieren kann, dass der Beobachtungswert B derart manipuliert wurde, dass der erste Ausgabewert p_{cls} wenigstens einen Teil des ersten Beobachtungswerts B nicht korrekt charakterisiert. Im Ausführungsbeispiel wird dies dadurch erreicht, dass das zweite maschinelle Lernsystem derart eingerichtet ist, dass der erste Ausgabewert p_{det} eine Wahrscheinlichkeit charakterisiert, dass es sich bei dem Beobachtungswert B um ein Irreführungsbeispiel handelt.

[0026] Die Manipulation des Beobachtungswerts B kann hierbei durch eine unmittelbare Manipulation des Beobachtungswerts B geschehen, aber auch durch eine Störung des Ausgangssignals S des Sensors **30** oder durch eine Störung der Umgebung **20**.

[0027] Erster Ausgabewert p_{cls} und zweiter Ausgabewert p_{det} werden einer Ausgabeinheit **80** übermittelt, die hieraus das Ansteuersignal A ermittelt. Beispielsweise ist es möglich, dass die Ausgabeinheit zunächst überprüft, ob der zweite Ausgabewert p_{det} kleiner ist als ein vorgebbare Schwellenwert. Ist dies der Fall, wird abhängig vom ersten Klassifikationsergebnis p_{cls} das Ansteuersignal A ermittelt. Dies ist der Normalfall. Wird hingegen ermittelt, dass der zweite Ausgabewert p_{det} nicht kleiner ist als der vorgebbare Schwellenwert, so kann vorgesehen sein, dass das Ansteuersignal A derart ausgebildet ist, dass es den Aktor A in einen abgesicherten Modus überführt.

[0028] Das Aktorsteuerungssystem **40** umfasst in einer Ausführungsform einen Computer und ein maschinenlesbares Speichermedium (nicht dargestellt), auf dem ein Computerprogramm gespeichert ist, das, wenn es vom Computer ausgeführt wird, diesen veranlasst, die beschriebenen Funktionalitäten des Aktorsteuerungssystems **40** auszuführen. Erstes maschinelle Lernsystem **60** und zweites maschinelle Lernsystem **70** können hier insbesondere als separate oder gemeinsame Computerprogramme implementiert sein.

[0029] Fig. 2 illustriert die Interaktion zwischen einem Trainingssystem **90** und dem ersten maschinellen Lernsystem **60** und dem zweiten maschinellen Lernsystem **70**. Das Trainingssystem **90** hält einen Satz von Trainingsdaten und zugehörigen gewünschten Ergebnissen vor. Aus dem Satz von Trainingsdaten werden ein oder mehrere Beobachtungswerte B ausgewählt und dem ersten maschinellen Lernsystem **60** zur Verfügung gestellt. Es kann sich hierbei um einzelne Beobachtungswerte handeln, also solchen, die dem ersten maschinellen Lernsystem **60** auch bei Interaktion des Aktorsteuerungssystems **40**

mit Aktor **10** und Sensor **30** zugeführt werden. Es kann sich aber auch um einen Stapel (engl. „batch“), also eine Mehrzahl von solchen Beobachtungswerten handeln.

[0030] Das erste maschinelle Lernsystem **60** ermittelt aus diesen ihm zugeführten Beobachtungswerten einen ersten Ausgabewert p_{cls} . Ebenso ermittelt das zweite maschinelle Lernsystem **60** analog zu Fig. 1 einen zweiten Ausgabewert p_{det} . Erster Ausgabewert p_{cls} und zweiter Ausgabewert p_{det} werden wieder dem Trainingssystem **90** zugeführt. Das Trainingssystem **90** ermittelt hieraus ein Parameteranpassungssignal P , das kodiert, welcher Parameter des ersten maschinellen Lernsystems **60** und welcher Parameter des zweiten maschinellen Lernsystems **70** wie ihren Wert ändern sollen. Diese gewünschte Anpassung erfolgt beispielsweise durch die Vorgabe gewünschter Werte für den ersten Ausgabewert p_{cls} und den zweiten Ausgabewert p_{det} und Backpropagation. Das Trainingssystem **90** führt zu diesem Zweck das Parameteranpassungssignal P einem Anpassungsblock **95** zu, der die Parameter im ersten maschinellen Lernsystem **60** und im zweiten maschinellen Lernsystem **70** entsprechend anpasst.

[0031] Trainingssystem **90** umfasst in einer Ausführungsform einen Computer und ein maschinenlesbares Speichermedium (nicht dargestellt), auf dem ein Computerprogramm gespeichert ist, das, wenn es vom Computer ausgeführt wird, diesen veranlasst, die beschriebenen Funktionalitäten des Lernsystems **90** auszuführen.

[0032] Fig. 3a zeigt in einem Flussdiagramm eine Ausführungsform eines Verfahrens zum Trainieren der Parameter des ersten maschinellen Lernsystems **60** und des zweiten maschinellen Lernsystems **70** durch das Trainingssystem **90**.

[0033] Zunächst trainiert das Trainingssystem **90** in einer ersten Phase **1000** die Parameter des ersten maschinellen Lernsystems **60** mit einer Trainingsmenge für das erste maschinelle Lernsystem (**60**) an Beobachtungswerten und dazugehörigen gewünschten ersten Ausgabewerten. Die Parameter des zweiten maschinellen Lernsystems **60** werden in diesem Schritt konstant gehalten, die zweiten Ausgabewerte p_{det} ignoriert.

[0034] Die erste Phase **1000** muss nicht zwingend vom Trainingssystem **90** durchgeführt werden. Es ist auch möglich, dass zu Beginn des Verfahrens die Parameter des ersten maschinellen Lernsystems **60** bereits vollständig trainiert sind.

[0035] Nun werden in einer zweiten Phase **1100** die trainierten Parameter des ersten maschinellen Lernsystems **60** (optional mit Ausnahme der Parameter des Blocks **61**) eingefroren und die Parameter des

zweiten maschinellen Lernsystems **70** (optional inklusive der Parameter des Blocks **61**) trainiert.

[0036] Der Ablauf der zweiten Phase **1100** ist in **Fig. 3b** genauer gezeigt. Zunächst (**1110**) für jeden Beobachtungswert B , der in einer Trainings-Menge für das zweite maschinelle Lernsystem (**70**) enthalten ist, entschieden, ob er manipuliert wird oder nicht. Dies kann beispielsweise zufällig mit einer vorgebbaren Wahrscheinlichkeit, beispielsweise 50%, erfolgen.

[0037] Anschließend (**1120**) werden die gewünschten zweiten Ausgabewerte p_{det} von Beobachtungswerten B , die manipuliert werden sollen, den Wert „1“ gesetzt, ansonsten auf den Wert „0“.

[0038] Im folgenden Schritt **1130** werden Parameter σ und α vorgegeben. σ ist ein vorgebbarer Parameter im Wertebereich $[0; 1]$, α ist ebenfalls ein vorgebbarer Parameter, vorzugsweise im Wertebereich $[0; 1]$, besonders bevorzugt im Bereich $[0.2; 0.3]$, noch weiter bevorzugt im Bereich $[0.24; 0.26]$. Im Ausführungsbeispiel wird α auf den Wert 0.25 gesetzt.

[0039] Dann (**1140**) werden Beobachtungswerte B , die manipuliert werden sollen, durch ihre manipulierte Form B^{adv} ersetzt, die mit dem unten angegebenen Algorithmus ermittelt werden können.

[0040] Für Beobachtungswerte B , für die entschieden wurde, dass sie manipuliert werden sollen, wird eine Manipulation B^{adv} beispielsweise wie folgt durchgeführt: Es wird ein Initialwert $B^{\text{adv}}_0 = B$ initialisiert, und dann iterativ Werte B^{adv}_n gemäß der folgenden Formel ermittelt:

$$B^{\text{adv}}_{n+1} = \text{Clip}_B^\epsilon \{ B^{\text{adv}}_n + \alpha [(1 - \sigma) \text{sgn}(\nabla_B J_{\text{cls}}(B^{\text{adv}}_n, y_{\text{true}}(B))) + \sigma \text{sgn}(\nabla_B J_{\text{det}}(B^{\text{adv}}_n, 1))] \} \quad (\text{Formel 1})$$

[0041] Hierbei ist J_{cls} eine Kostenfunktion des ersten maschinellen Lernsystems **60**. Wird das erste maschinelle Lernsystem **60** als Klassifikator eingesetzt, ist dies bevorzugt eine Kreuzentropie des ersten maschinellen Lernsystems **60**, also die Kreuzentropie der ersten Ausgabewerte p_{cls} und der gewünschten Bezeichner der Klassifikation.

[0042] J_{det} ist eine Kostenfunktion des zweiten maschinellen Lernsystems **70**. Bevorzugt ist dies eine Kreuzentropie des zweiten maschinellen Lernsystems **70**, also die Kreuzentropie der zweiten Ausgabewerte p_{det} und dem Label „1“, das für manipulierte Beobachtungen B^{adv} vorgesehen ist.

[0043] Die Funktion $\text{Clip}_B^\epsilon(z)$ normiert Werte einer Variablen z auf eine ϵ -Kugel um B . Die Norm kann hierbei eine L^2 -Norm sein, oder auch eine L^∞ -Norm.

[0044] Die Anzahl der Iterationen kann beispielsweise mittels eines Konvergenzkriteriums begrenzt werden, oder auf einen festen Wert vorgegeben werden, beispielsweise **10**.

[0045] Im nun folgenden Schritt **1150** werden die trainierten Parameter des ersten maschinellen Lernsystems **60** (optional mit Ausnahme der Parameter des Blocks **61**) eingefroren und die Parameter des zweiten maschinellen Lernsystems **70** (optional inklusive der Parameter des Blocks **61**) trainiert.

[0046] Das Trainieren der Parameter erfolgt mittels der beschriebenen Kostenfunktion J_{det} des zweiten maschinellen Lernsystems **70** und Backpropagation.

[0047] Damit endet die zweite Phase **1100**.

[0048] Diese Form des Lernens kann derart durchgeführt werden, dass der Trainingsmengen so groß gewählt wird, dass die erste Phase **1000** und die zweite Phase **1100** jeweils nur einmalig durchgeführt wird. Um ein Trainieren der Parameter des zweiten maschinellen Lernsystems **60** zu gewährleisten, ist vorteilhafterweise vorgesehen, den Wert für die Wahrscheinlichkeit, dass die Beobachtungen B manipuliert werden sollen, auf den Wert „1“ zu setzen.

[0049] **Fig. 4** zeigt in einem Flussdiagramm eine weitere Ausführungsform des Verfahrens zum Trainieren der Parameter des ersten maschinellen Lernsystems **60** und des zweiten maschinellen Lernsystems **70** durch das Trainingssystem **90**. Dieses Verfahren macht das Aktorsteuerungssystem **40** besonders robust gegen Angreifer, die Beobachtungsdaten auch mittels Kenntnissen über den inneren Aufbau des zweiten maschinellen Lernsystems **70** manipulieren.

[0050] Zunächst wird in einem ersten Schritt **2000** eine Untermenge von Beobachtungswerten B aus einer Trainings-Menge ausgewählt. Der Trainings-Menge umfasst erneut Paare von Beobachtungswerten B und zugehörige gewünschte ersten Ausgabewerte.

[0051] Nun (**2100**) wird mit einer vorgebbaren Wahrscheinlichkeit von beispielsweise 50% entschieden, ob die Beobachtungswerte B dieser Untermenge manipuliert werden sollen oder nicht. Wenn ja, folgt Schritt **2200**, andernfalls Schritt **2300**.

[0052] In Schritt **2200** wird für die Datenpunkte der ausgewählten Untermenge dann der zweiten Ausgabewert auf den Wert „1“ gesetzt. Es folgt Schritt **2210**, in dem die Variable σ optional mittels eines (Pseudo-)Zufallszahlengenerators zufällig auf einen Wert aus dem Wertebereich $[0; 1]$ gesetzt und gespeichert. Eine derartige zufällige Wahl von σ macht das Aktorsteuerungssystem **40** besonders robust gegen eine breite Klasse möglicher Manipulationen der

Beobachtungswerte B . (Alternativ kann die Variable σ auch auf einen festen Wert festgelegt werden, oder einmalig zufällig ausgewählt werden). Die Variable α wird wie in Schritt **1130** beschrieben gewählt.

[0053] Im folgenden Schritt **2210** werden für alle Beobachtungswerte B der ausgewählten Untermenge Manipulationswerte B^{adv} wie in Schritt **1150** beschrieben ermittelt und die ursprünglichen Beobachtungswerte B durch die zugehörigen ermittelten manipulierten Beobachtungswerte B^{adv} ersetzt. Es folgt Schritt **2400**.

[0054] In Schritt **2300** wird für die Datenpunkte der ausgewählten Untermenge dann der zweiten Ausgabewert auf den Wert „0“ gesetzt.

[0055] Wurde hingegen entschieden, dass keine Manipulation der Beobachtungswerte B durchgeführt werden soll, werden in Schritt **2300** die Beobachtungswerte B in der Untermenge belassen, und es folgt ebenfalls Schritt **2400**.

[0056] Nun werden in Schritt **2400** die Parameter des ersten maschinellen Lernsystems **60** und des zweiten maschinellen Lernsystems **70** mit der (ggf. manipulierten) Untermenge von Beobachtungswerten und den zugehörigen ersten und zweiten gewünschten Ausgabewerten trainiert.

[0057] Soll das Verfahren mehrfach iteriert durchgeführt werden, können die Datenpunkte der ausgewählten Untermenge aus der Trainings-Menge entfernt werden (**2500**) und es kann an dieser Stelle zurückverzweigt werden zu Schritt **2000**, und das Verfahren beispielsweise so lange durchgeführt werden, bis die Trainings-Menge keine Datenpunkte mehr enthält.

[0058] Optional kann das so beschriebene Verfahren für die gleichen Trainings-Mengen mehrfach durchgeführt werden, beispielsweise mit einer vorgebbaren Häufigkeit.

[0059] Fig. 5 zeigt den Aufbau des ersten maschinellen Lernsystems **60** und des zweiten maschinellen Lernsystems **70** gemäß einer Ausführungsform der Erfindung, in dem beide durch künstliche neuronale Netze gegeben sind.

[0060] Das erste maschinelle Lernsystem **60** besteht aus einer Eingabeschicht **100**, dem der Beobachtungswert B zugeführt wird, gefolgt von einer Konvolutionsschicht **110**, multiplen (hier beispielhaft: fünffachen) Residualblöcken **120a**, **120b**, **120c** einem global-average Pooling Schicht **130** und einer Fully-connected Schicht **140**, dessen Ausgabewert der erste Ausgabewert p_{cls} ist. Ein Ausgangssignal einer vorhergehenden Schicht des neuronalen Netzes wird in üblicher Weise der folgenden Schicht als Eingabesi-

gnal zugeführt. Es kann vorgesehen sein, dass die Dimensionalität dieses Signals von einem oder mehrerer der multiplen Residualblöcke **120a**, **120b**, **120c** reduziert wird. Der globalaverage Pooling Layer ist in üblicher Weise ebenfalls eingerichtet, eine dimensionale Reduktion durchzuführen.

[0061] Konvolutionsschicht **110** führt vorteilhafterweise eine 3×3 -Konvolution durch.

[0062] Das zweite maschinelle Lernsystem **70** besteht aus einer Konvolutionsschicht **170**, gefolgt von einer optionalen max-pooling Schicht **155**, einer Konvolutionsschicht **160**, einer weiteren, ebenfalls optionalen max-pooling Schicht **165**, zwei weiteren Konvolutionsschichten **170**, **175** und einer global-average Pooling Schicht **180**. Das Ausgangssignal der global-average Pooling Schicht **180** ist das zweite Ausgabesignal p_{det} .

[0063] Konvolutionsschicht **175** ist hierbei vorteilhafterweise eine Konvolution mit einer geringeren Schrittweite als die Konvolutionsschichten **150**, **160** und **170**. Beispielsweise führt die Konvolutionsschicht **175** eine 1×1 -Konvolution durch, die Konvolutionsschichten **150**, **160**, **170** eine 3×3 -Konvolution.

[0064] Es sei angemerkt, dass die in Fig. 5 illustrierten Konvolutionsschichten **110**, **150**, **160**, **170**, **175** jeweils die eigentliche Konvolution umfassen, sowie ferner eine Stapel-Normalisierung (engl. „batch normalization“), in der die Werte des Eingangssignals x für den Fall, dass es ein hochdimensionales Signal ist, das aus einem Stapel von Beobachtungswerten B gewonnen wurde, hinsichtlich dieses Stapels von Beobachtungswerten normalisiert werden, und einen Aktivierungsblock, in dem das Ausgangssignal als Aktivierungsfunktion des Eingangssignals verwendet wird, beispielsweise eine Rectified Linear Unit (ReLU).

[0065] Das Eingangssignal y_{zw} des zweiten maschinellen Lernsystems **70**, das der Konvolutionsschicht **150** zugeführt wird, wird aus dem ersten maschinellen Lernsystem **70** als Ausgangssignal des Eingangsblocks **100** oder des Konvolutionsblocks **110** oder eines der multiplen Residualblöcke **120a**, **120b**, **120c** abgezweigt.

[0066] Fig. 6 illustriert den Aufbau eines Residualblocks, aus dem wie in Fig. 7 illustriert multiple Residualblöcke zusammengesetzt sind. Ein Eingangssignal x des Residualblocks wird zunächst einem Konvolutionsblock **200** zugeführt. Auf den Residualblock folgt eine Stapel-Normalisierung **210**.

[0067] Es folgt ein Aktivierungsblock **220**, in dem das Ausgangssignal als Aktivierungsfunktion des Eingangssignals verwendet wird, beispielsweise eine

Rectified Linear Unit (ReLU). Es folgt ein weiterer Konvolutionsblock **230**, eine weiterer Stapel-Normalisierungs-block **240** und ein Additionsblock **260**. Das Eingangssignal x wird parallel zum Konvolutionsblock **200** einem optionalen noch weiteren Konvolutionsblock **270** und einem noch weiteren Stapel-Normalisierungs-Block **280**. Das Ausgangssignal des noch weiteren Stapel-Normalisierungs-Blocks **280** wird dem Additionsblock **260** zugeführt. Noch weiterer Konvolutionsblock **270** und noch weiterer Stapel-Normalisierungs-Block **280** werden nur benötigt, wenn der Residualblock x eine dimensionale Reduktion des Eingangssignals x durchführen soll. In diesem Fall führen auch Konvolutionsblock **200** und Stapel-Normalisierungs-Block **210** die dimensionale Reduktion durch.

[0068] Additionsblock **260** addiert das Ausgangssignal des weiteren Stapel-Normalisierungs-Blocks **240** bzw. des weiteren Aktivierungsblocks **250** und das Ausgangssignal des Stapel-Normalisierungs-Blocks **280** bzw. das Eingangssignal x , und führt es optional einem noch weiteren Aktivierungsblock **290** zu, um das Ausgangssignal y zu gewinnen.

[0069] Fig. 7 illustriert den Aufbau eines multiplen Residualblocks. Eingangssignal x wird einem ersten Residualblock **300** zugeführt, dem weitere Residualblöcke **310**, **320**, **330**, **340** folgen, um das Ausgangssignal y zu gewinnen. Soll der multiple Residualblock eine dimensionale Reduktion durchführen, wird der erste Residualblock **300** so gewählt, dass er diese dimensionale Reduktion durchführt (d.h. er weist Konvolutionsblock **270** und Stapel-Normalisierungsblock **280** auf, vgl. Fig. 4). Andernfalls wird der erste Residualblock **300** ebenso wie die weiteren Residualblöcke **310**, **320**, **330**, **340** so gewählt, dass er keine dimensionale Reduktion durchführt.

ZITATE ENTHALTEN IN DER BESCHREIBUNG

Diese Liste der vom Anmelder aufgeführten Dokumente wurde automatisiert erzeugt und ist ausschließlich zur besseren Information des Lesers aufgenommen. Die Liste ist nicht Bestandteil der deutschen Patent- bzw. Gebrauchsmusteranmeldung. Das DPMA übernimmt keinerlei Haftung für etwaige Fehler oder Auslassungen.

Zitierte Patentliteratur

- DE 102005050577 A1 [0002]

Schutzansprüche

1. Aktorsteuerungssystem (**40**) zum Steuern eines Aktors (**10**), umfassend mindestens einen Computer und mindestens ein maschinenlesbares Speichermedium, auf dem Befehle gespeichert sind, die beim Ausführen durch den mindestens einen Computer bewirken, dass der mindestens eine Computer ein Verfahren mit den folgenden Schritten ausführt:

– Empfangen eines ermittelten Beobachtungswerts (B), der einen Zustand eines Aktorsystems umfassend den Aktor (**10**) und eine Umgebung (**20**) des Aktors (**10**) charakterisiert;

– Ermitteln eines ersten Ausgabewerts (p_cls) eines ersten maschinellen Lernsystems (**60**) abhängig von dem Beobachtungswert (B), wobei der erste Ausgabewert (p_cls) wenigstens einen Teil des Beobachtungswerts (B) charakterisiert,

– Ermitteln eines zweiten Ausgabewerts (p_det) eines zweiten maschinellen Lernsystems (**70**), wobei der zweite Ausgabewert (p_det) eine Wahrscheinlichkeit charakterisiert, dass der Beobachtungswert (B) derart manipuliert wurde, dass der erste Ausgabewert (p_cls) den wenigstens einen Teil des ersten Beobachtungswerts (B) nicht korrekt charakterisiert, und

– Ansteuern des Aktors (**10**) abhängig von dem ermittelten ersten Ausgabewert (p_cls) und dem ermittelten zweiten Ausgabewert (p_det), wobei das zweite maschinelle Lernsystem (**70**) den zweiten Ausgabewert (p_det) abhängig von Zwischenergebnissen (y_zw) des ersten maschinellen Lernsystems (**60**) ermittelt.

2. Aktorsteuerungssystem (**40** nach Anspruch 1, wobei der erste Ausgabewert (p_cls) wenigstens einen Teil des Beobachtungswerts (B) als einer Klasse aus einer Mehrzahl von Klassen zugehörig klassifiziert, und der zweite Ausgabewerts (p_det) eine Wahrscheinlichkeit charakterisiert, dass der Beobachtungswert (B) derart manipuliert wurde, dass der wenigstens eine Teil des ersten Beobachtungswerts (B) vom ersten maschinellen Lernsystem (**60**) fehlklassifiziert wurde.

3. Aktorsteuerungssystem (**40**) nach Anspruch 1 oder 2, wobei der zweite Ausgabewert (p_det) eine Wahrscheinlichkeit charakterisiert, dass der Beobachtungswert (B) ein Irreführungsbeispiel, des ersten maschinellen Lernsystems (**60**) ist;

4. Aktorsteuerungssystem (**40**) nach einem der Ansprüche 1 bis 3, wobei das erste maschinelle Lernsystem (**60**) ein erstes neuronales Netzwerk mit einer ersten Verkettung neuronaler Schichten (**100**, **110**, **120a**, **120b**, **120c**, **130**, **140**) ist, und das zweite maschinelle Lernsystem (**70**) ein zweites neuronales Netzwerk mit einer zweiten Verkettung neuronaler Schichten (**100**, **110**, **120a**, **120b**, **120c**, **150**, **155**, **160**, **165**, **170**, **175**, **180**) ist, wobei die neuronalen Schichten der ersten (**100**, **110**, **120a**, **120b**,

120c, **130**, **140**) und zweiten (**100**, **110**, **120a**, **120b**, **120c**, **150**, **155**, **160**, **165**, **170**, **175**, **180**) Verkettung neuronaler Schichten bis zu einer letzten gemeinsamen Zwischenschicht (**100**, **110**, **120a**, **120b**, **120c**) gemeinsame neuronale Schichten sind und ab dieser letzten gemeinsamen Zwischenschicht (**100**, **110**, **120a**, **120b**, **120c**) separat sind.

5. Aktorsteuerungssystem (**40**) nach Anspruch 4, wobei die gemeinsamen neuronalen Schichten eine Konvolutionsschicht (**110**) umfassen.

6. Aktorsteuerungssystem (**40**) nach Anspruch 4 oder 5, wobei die gemeinsamen neuronalen Schichten mindestens einen multiplen, insbesondere 5-fachen, Residualblock (**120a**, **120b**, **120c**) umfassen.

7. Aktorsteuerungssystem (**40**) nach einem der Ansprüche 4 bis 6, wobei das zweite neuronale Netz neben den gemeinsamen neuronalen Schichten noch mindestens drei, insbesondere genau vier, weitere Konvolutionsschichten (**150**, **160**, **170**, **175**) umfasst.

8. Aktorsteuerungssystem (**40**) nach einem der Ansprüche 1 bis 7, wobei der Aktor (**10**) ein autonomer oder teilautonomer Roboter, insbesondere ein Kraftfahrzeug oder ein Fertigungsroboter, ist.

9. Aktorsteuerungssystem (**40**) nach einem der Ansprüche 1 bis 7, wobei der Aktor (**10**) ein Freigabesystem und die Beobachtung (B) ein Ausgangssignal eines Bild- und/oder Videoerfassungssystems ist.

10. Aktorsteuerungssystem (**40**) nach einem der Ansprüche 1 bis 9, dessen zweites (**60**) maschinelle Lernsystem mit einem Trainingssystem (**90**) trainiert wurde,

wobei das Trainingssystem (**90**) zum Trainieren des Aktorsteuerungssystems (**40**) nach einem der Ansprüche 1 bis 9 eingerichtet ist,

wobei das Trainingssystem (**90**) mindestens einen zweiten Computer und mindestens ein zweites maschinenlesbares Speichermedium umfasst, auf dem Befehle gespeichert sind, die beim Ausführen durch den mindestens einen zweiten Computer bewirken, dass der mindestens eine zweite Computer ein Verfahren zum Trainieren des zweiten maschinellen Lernsystems (**70**) mit den folgenden Schritten ein- oder mehrfach ausführt:

b) Auswahl einer Untermenge von Beobachtungswerten (B) aus einer Trainings-Menge für das zweite maschinelle Lernsystem (**70**), die Paare von Beobachtungswerten (B) und zugehörigen gewünschten ersten Ausgabewerten umfasst,

c) Entscheiden, ob die Beobachtungswerte (B) dieser Untermenge manipuliert werden sollen oder nicht,

d) Setzen des gewünschten zweiten Ausgabewerts auf den Wert eines vorgebbaren ersten Zahlenwerts, wenn die Beobachtungen (B) nicht manipuliert werden bzw. auf einen zweiten, vom ersten Zahlen-

wert verschiedenen, vorgebbaren zweiten Zahlenwert, wenn die Beobachtungen (B) manipuliert werden,

e) Manipulation der Beobachtungswerte (B) dieser Untermenge, sofern entschieden wurde, dass die Manipulation durchgeführt werden soll,

f) Ersetzen der ursprünglichen Beobachtungswerte (B) durch die zugehörigen manipulierten Beobachtungswerte (B^{adv}), und

g) Trainieren der Parameter des ersten maschinellen Lernsystems (**60**) und des zweiten maschinellen Lernsystems (**70**) mit der (ggf. manipulierten) Untermenge von Beobachtungswerten (B, B^{adv}) und zugehörigen ersten und zweiten gewünschten Ausgabewerten.

relativ zum Gradienten der zweiten Kostenfunktion (J_{det}) gewichtet wird zufällig gewählt wird.

Es folgen 7 Seiten Zeichnungen

11. Aktorsteuerungssystem (**40**) nach Anspruch 10, wobei die auf dem zweiten maschinenlesbaren Speichermedium gespeicherten Befehle des Trainingssystem (**90**) bewirken, dass der mindestens eine zweite Computer ein Verfahren nicht nur zum Trainieren des zweiten (**70**), sondern auch zum Trainieren des ersten maschinellen Lernsystems (**60**) durchführt, bei welchem vor Durchführung der Schritte b) bis g) der folgende Schritt durchgeführt wird: a) Trainieren des ersten maschinellen Lernsystems (**60**) mit einer Trainings-Menge für das erste maschinelle Lernsystem (**60**) aus Beobachtungswerten (B) und zugehörigen gewünschten ersten Ausgabewerten.

12. Aktorsteuerungssystem (**40**) nach Anspruch 10 oder 11, bei dem die manipulierten Beobachtungswerte (B^{adv}) abhängig von einem Wert einer ersten Kostenfunktion (J_{cls}) des ersten maschinellen Lernsystems (**60**) für den jeweiligen ursprünglichen Beobachtungswert (B) generiert werden.

13. Aktorsteuerungssystem (**40**) nach Anspruch 12, bei dem die manipulierten Beobachtungswerte (B^{adv}) auch abhängig von einem Wert einer zweiten Kostenfunktion (J_{det}) des zweiten maschinellen Lernsystems (**70**) für den jeweiligen ursprünglichen Beobachtungswert (B) generiert werden.

14. Aktorsteuerungssystem (**40**) nach Anspruch 13, wobei die manipulierten Beobachtungswerte (B^{adv}) von Beobachtungswerten (B) abhängig von einem Gradienten der ersten Kostenfunktion (J_{cls}), insbesondere einer ersten Kreuzentropie, des ersten maschinellen Lernsystems (**60**) ist und abhängig von einem Gradienten der zweiten Kostenfunktion (J_{der}), insbesondere einer zweiten Kreuzentropie, des zweiten maschinellen Lernsystems (**70**) ist.

15. Aktorsteuerungssystem (**40**) nach Anspruch 14, wobei ein Parameter (σ), der charakterisiert, wie stark der Gradient der ersten Kostenfunktion (J_{cls})

Anhängende Zeichnungen

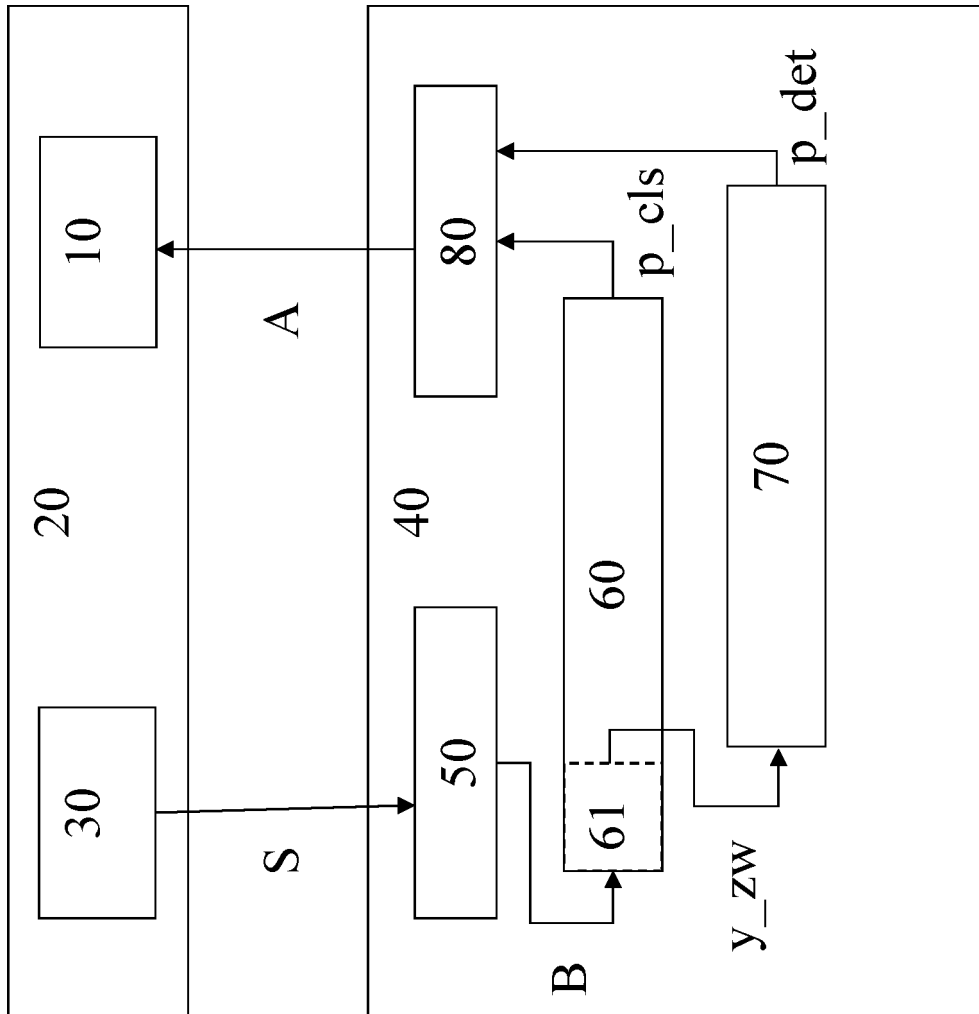


Fig. 1

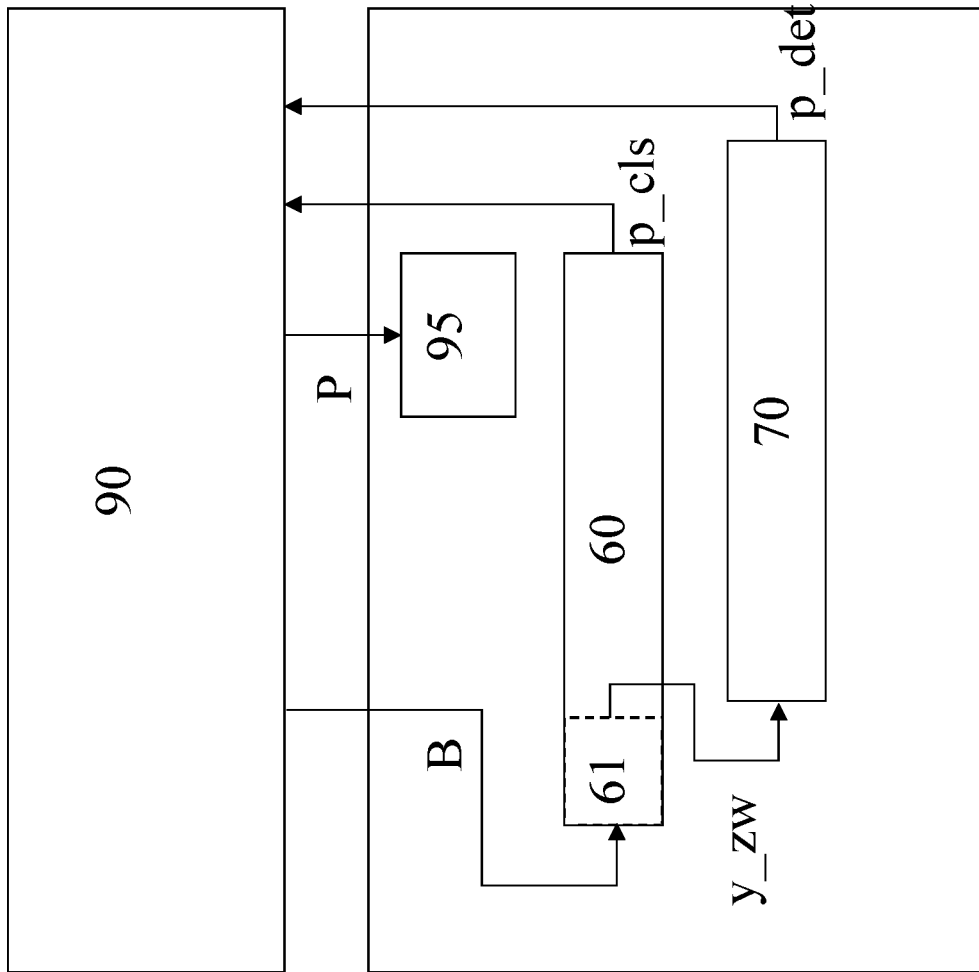


Fig. 2

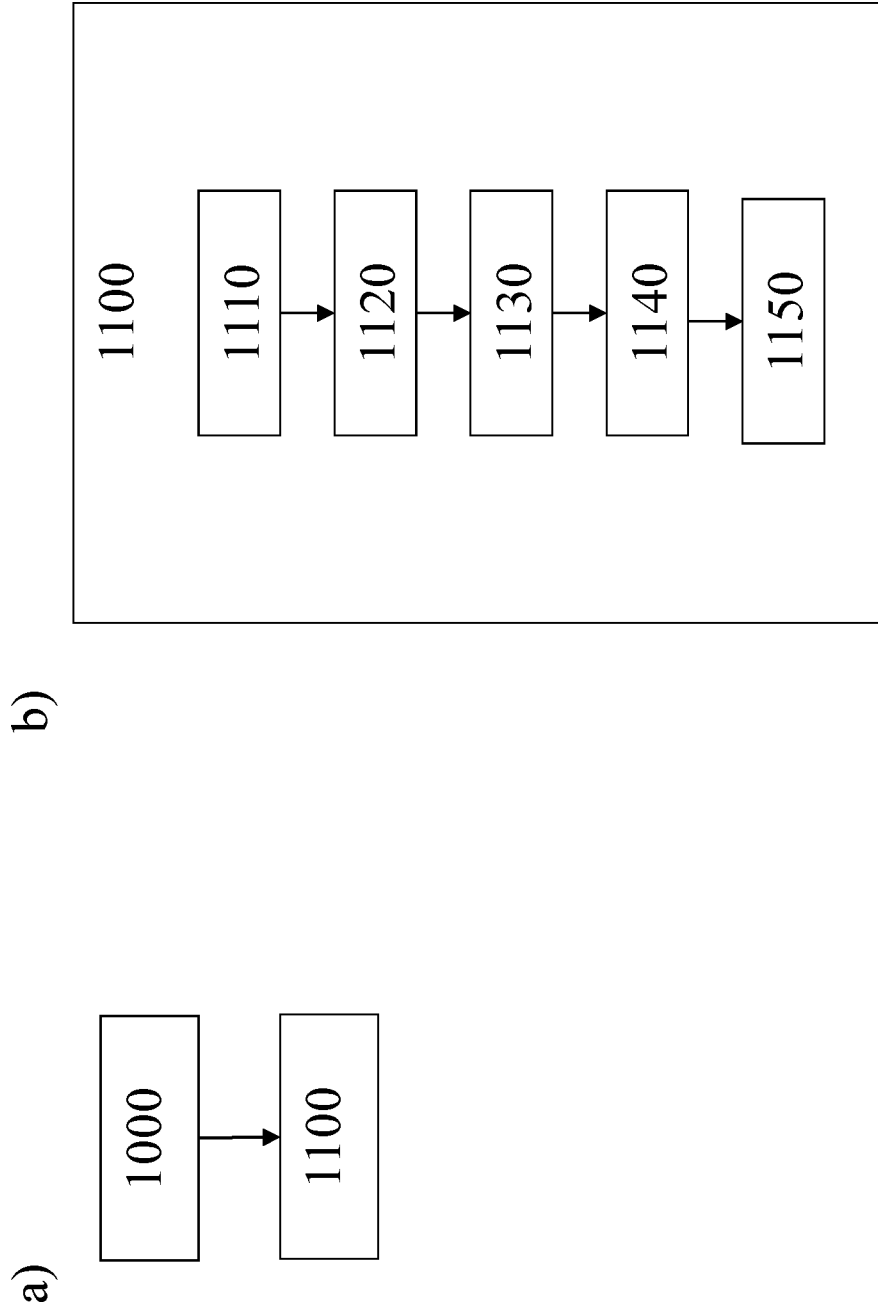


Fig. 3

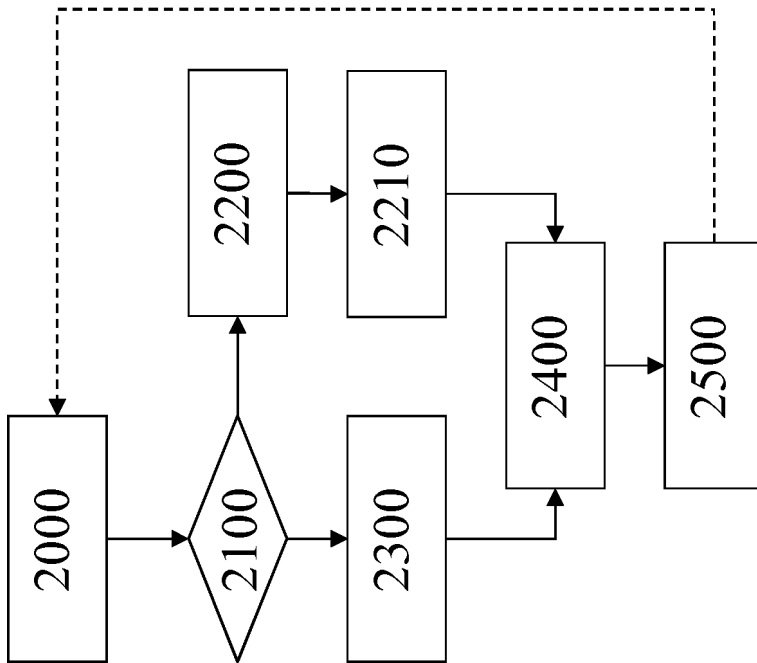


Fig. 4

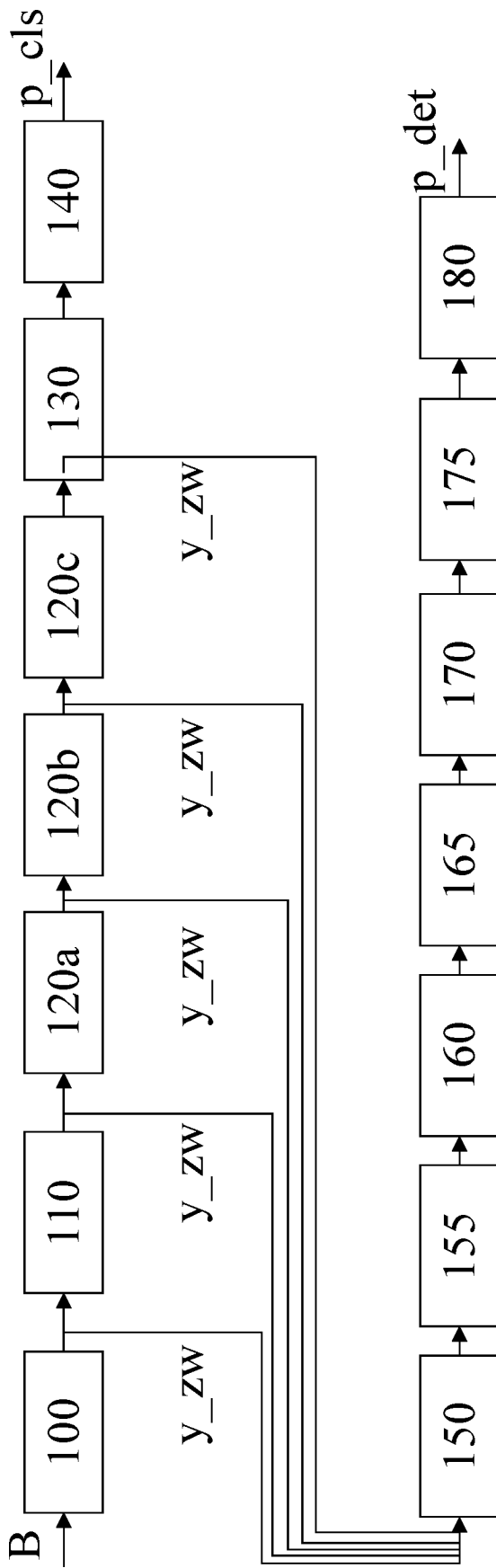


Fig. 5

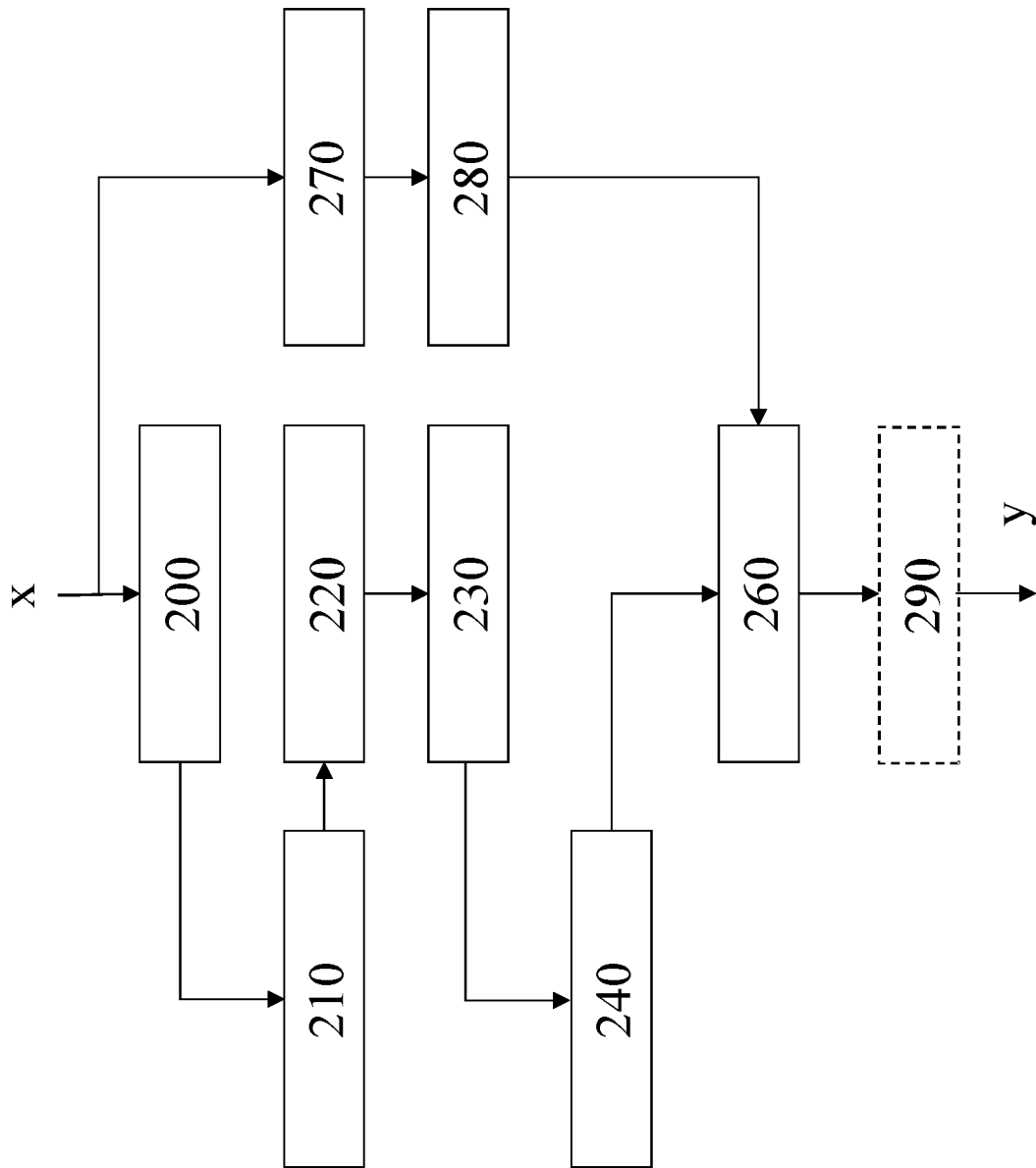


Fig. 6

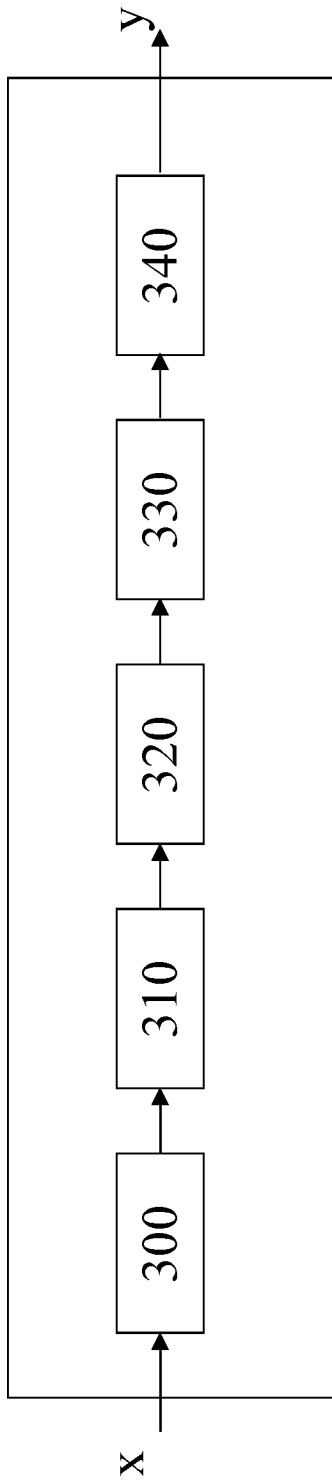


Fig. 7