



AKADEMIA GÓRNICZO-HUTNICZA IM. STANISŁAWA STASZICA W KRAKOWIE

**WYDZIAŁ ELEKTROTECHNIKI, AUTOMATYKI,
INFORMATYKI I INŻYNIERII BIOMEDYCZNEJ**

KATEDRA INFORMATYKI STOSOWANEJ

Praca dyplomowa inżynierska

Sztuczna inteligencja na przykładzie symulacji komputerowej
Artificial intelligence in computer simulation

Autor:	<i>Bartłomiej Konieczny</i>
Kierunek studiów:	<i>Informatyka</i>
Opiekun pracy:	<i>dr inż. Mirosław Gajer</i>

Kraków, 2015

Uprzedzony o odpowiedzialności karnej na podstawie art. 115 ust. 1 i 2 ustawy z dnia 4 lutego 1994 r. o prawie autorskim i prawach pokrewnych (t.j. Dz.U. z 2006 r. Nr 90, poz. 631 z późn. zm.): „Kto przywłaszcza sobie autorstwo albo wprowadza w błąd co do autorstwa całości lub części cudzego utworu albo artystycznego wykonania, podlega grzywnie, karze ograniczenia wolności albo pozbawienia wolności do lat 3. Tej samej karze podlega, kto rozpowszechnia bez podania nazwiska lub pseudonimu twórcy cudzy utwór w wersji oryginalnej albo w postaci opracowania, artystycznego wykonania albo publicznie zniekształca taki utwór, artystyczne wykonanie, fonogram, wideogram lub nadanie.”, a także uprzedzony o odpowiedzialności dyscyplinarnej na podstawie art. 211 ust. 1 ustawy z dnia 27 lipca 2005 r. Prawo o szkolnictwie wyższym (t.j. Dz. U. z 2012 r. poz. 572, z późn. zm.) „Za naruszenie przepisów obowiązujących w uczelni oraz za czyny uchybiające godności studenta student ponosi odpowiedzialność dyscyplinarną przed komisją dyscyplinarną albo przed sądem koleżeńskim samorządu studenckiego, zwanym dalej „sądem koleżeńskim”, oświadczam, że niniejszą pracę dyplomową wykonałem(-am) osobiście i samodzielnie i że nie korzystałem(-am) ze źródeł innych niż wymienione w pracy.

Spis treści

1. Wstęp	7
1.1. Cele pracy	7
1.2. Zawartość pracy	8
2. Uczenie maszynowe - rozdział teoretyczny	9
2.1. Podejścia do uczenia maszynowego	9
2.1.1. Uczenie nadzorowane	9
2.1.2. Uczenie nienadzorowane	10
2.1.3. Uczenie ze wzmocnieniem	12
2.2. Podsumowanie	14
3. Implementacja	15
3.1. Tablica stan-akcja $Q(S, A)$	15
3.2. Polityka wyboru akcji	15
3.2.1. Równowaga pomiędzy eksploracją, a eksploatacją	16
3.2.2. ϵ -greedy policy	16
3.2.3. Optimal policy	16
3.3. Symulacja graficzna	16
3.4. Reprezentacja stanu	17
3.4.1. Akcje	17
3.4.2. Wybór algorytmów	19
3.5. Wykorzystane technologie	20
3.6. Napotkane problemy	20
3.7. Wynik działania	20
4. Podsumowanie	21

1. Wstęp

W ostatniej dekadzie można zauważyć zwiększone zainteresowanie rozwiązaniami z dziedziny sztucznej inteligencji. Innowacyjne pomysły z użyciem tych algorytmów pozwalają nie tylko na interpretację ogromnych ilości danych, których człowiek nie jest w stanie przetworzyć ale również, między innymi, na rozwój autonomicznych pojazdów, jeżdżących bez nadzoru kierowcy.

Według [1] sztuczną inteligencją możemy nazwać “badanie i rozwój inteligentnych maszyn, w szczególności programów komputerowych”.

Inteligentne zachowanie agenta możemy zdefiniować, gdy agent[2]

- dostosowuje swoje zachowanie do aktualnych warunków i celów,
- ma zdolność zmiany otoczenia i celów,
- uczy się z doświadczenia,
- wykonuje odpowiednie do swoich ograniczeń akcję.

Wykorzystując powyższe definicję, sztuczną inteligencję określamy jako dziedzinę naukową zajmującą się badaniem, rozwojem i implementacją programów i maszyn wykazujących cechy inteligencji, tzn. takie które ucząc się z doświadczenia, będąc zmiennym w stosunku do otaczającego ich otoczenia i warunków dążą do wykonania swoich celów uwzględniając obowiązujące je ograniczenia.

1.1. Cele pracy

Celem pracy jest opis i implementacja inteligentnego agenta. Agent wykorzystując algorytmy uczenia ze wzmocnieniem wyciąga wnioski z podejmowanych akcji i dostosowuje swoje zachowanie. Aby lepiej zaprezentować wyniki jego działania, stworzone zostanie proste środowisko graficzne, które będzie przedstawiać wyniki przeprowadzanych przez robota wyborów.

W pracy zostanie wyjaśnione pojęcie uczenia maszynowego i podstawowe podejścia do rozwiązania problemów z tej dziedziny. Dzięki temu możliwe będzie podsumowanie różnic między uczeniem ze wzmocnieniem, a pozostałymi podejściami do uczenia maszynowego.

1.2. Zawartość pracy

Pierwszy rozdział stanowi wprowadzenie do pracy, określający jej zakres. W rozdziale drugim przedstawiono wyjaśnienie podstawowych pojęć i podstaw teoretycznych. W rozdziale trzecim zostały wyjaśnione szczegóły implementacji inteligentnego agenta oraz opis wybranego algorytmu. Pracę zamyka rozdział podsumowujący jej treść.

2. Uczenie maszynowe - rozdział teoretyczny

Jednym z najistotniejszych zagadnień z dziedziny sztucznej inteligencji jest uczenie maszynowe. Uczenie maszynowe jest metodą analizy danych, która automatyzuje budowę modelu analitycznego na podstawie nauki z danych. W wielu zastosowaniach ich użycie jest znacznie bardziej efektywne od manualnego programowania, w wyniku czego uczenie maszynowe znalazło szerokie zastosowanie w informatyce i innych dziedzinach. W ostatniej dekadzie można zauważyć zwiększone użycie metod uczenia maszynowego[3].

2.1. Podejścia do uczenia maszynowego

- uczenie nadzorowane,
- uczenie nienadzorowane,
- uczenie ze wzmocnieniem,

2.1.1. Uczenie nadzorowane

Uczenie nadzorowane polega na wnioskowaniu funkcji z określonych danych trenujących. Wykorzystując dostarczone przykłady algorytmy potrafią estymować wartości danych, które mogą nie występować w podanym zbiorze wejściowym. Dzięki generalizowaniu z przykładów, metody uczenia nadzorowanego są w stanie wyznaczać przewidywane wartości na podstawie danych trenujących.

Ważną cechą danych trenujących w uczeniu nadzorowanym jest konieczność ich oznaczenia. Algorytm, aby móc szacować pożądane wartości funkcji, musi posiadać wiedzę o ich cechach.

Przykładem zastosowania algorytmów uczenia nadzorowanego jest system rozpoznawania niechcianych wiadomości w klientach pocztowych. Danymi wejściowymi są w tym przypadku kategoryzowane na pożądane lub niepożądane wiadomości e-mail. System generalizując podane mu przykłady jest w stanie zidentyfikować kolejne wiadomości i wykonać odpowiednią akcję, zależnie od preferencji użytkownika (może to być na przykład usunięcie lub przeniesienie do zdefiniowanego folderu).

Wiele różnych algorytmów uczenia nadzorowanego zostało wykorzystanych by rozwiązać problem klasyfikacji wiadomości e-mail. Użyto między innymi algorytmów k-nearest neighbor[4], Naive Bayes[5][6] czy Random Forest[7], jednak wiąże się to z istotnymi wadami[8]:

- **Wymagane oznaczenie danych testowych.** Metody uczenia nadzorowanego wymagają, aby dane trenujące były oznaczone. W przypadku klasyfikacji wiadomości e-mail, konieczne jest ich oznaczenie w zależności od tego czy są szkodliwe czy nie. Problem stwarza tutaj wielkość danych. Ilość wiadomości, która jest wymieniana w sieci jest bardzo duża. W związku z czym, żeby klasyfikacja miała sens, wymagane też jest oznaczenie sporej ilości przykładów, co nie zawsze jest możliwe i opłacalne do zrealizowania.
- **Mała liczba danych testowych.** W związku z niewielką (w stosunku do wszystkich możliwych) ilością danych trenujących, algorytm jest mało odporny na modyfikowane dane. Osoby roszylające niechciane wiadomości bardzo często będą zmieniać ich treść i strukturę, na taką, która nigdy nie pojawiła się wśród danych trenujących. Może mieć to negatywny wpływ na wynik działania algorytmu.

2.1.2. Uczenie nienadzorowane

Podobnie jak w uczeniu nadzorowanym, algorytmy uczenia nienadzorowanego wyznaczają funkcje na podstawie danych wejściowych, jednak są w stanie odkryć niewidoczne zależności między nimi. Konsekwencją wynikającą z charakterystyki algorytmów uczenia nadzorowanego jest niemożność określenia błędu lub poprawności rozwiązania. Celem działania algorytmu może być na przykład kategoryzowanie informacji (klasteryzacja).

W odróżnieniu od uczenia nadzorowanego, metody uczenia nienadzorowanego są w stanie wykryć ukryte wzorce w danych wejściowych.

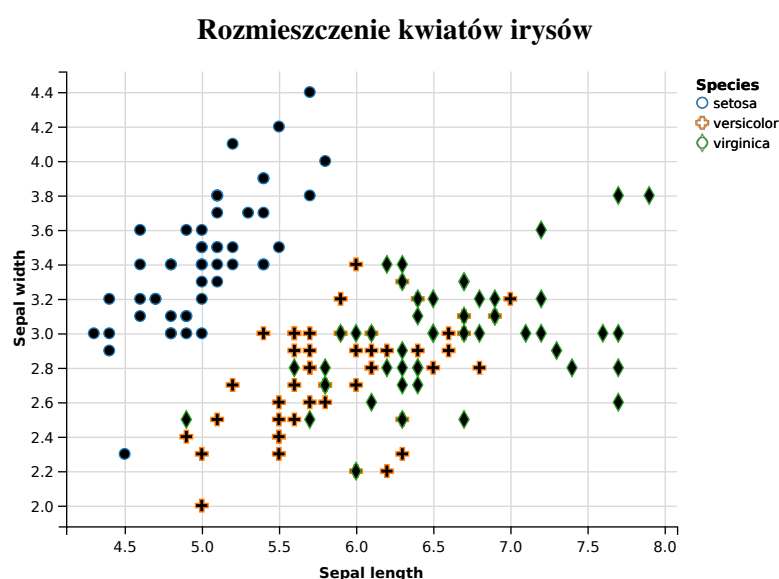
Używając jako wejścia informacji dotyczących kwiatów irysów w postaci przedstawionej na listingu 2.1, algorytmy uczenia nienadzorowanego mogą być zastosowane w celu wywnioskowania gatunku kwiatu (*setosa*, *versicolor*, *virginica*) na podstawie długości i szerokości płatków (*sepal*) i listka kielicha (*petal*).

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
53	6.9	3.1	4.9	1.5	versicolor
54	5.5	2.3	4.0	1.3	versicolor
55	6.5	2.8	4.6	1.5	versicolor
56	5.7	2.8	4.5	1.3	versicolor
101	6.3	3.3	6.0	2.5	virginica
102	5.8	2.7	5.1	1.9	virginica
103	7.1	3.0	5.9	2.1	virginica
104	6.3	2.9	5.6	1.8	virginica

Listing 2.1. Przykład danych o kwiatach irysów

Na rys. 2.1, przedstawione zostało rozmieszczenie gatunków kwiatów w zależności od długości i szerokości płatków kwiatu. Wyraźnie widać podział na dwa podstawowe klastry

- gatunek setosa,
- gatunek versicolor i virginica.



Rys. 2.1. Populacja kwiatów irysów w zależności od szerokości i długości płatków kwiatu. Źródło: Opracowanie własne

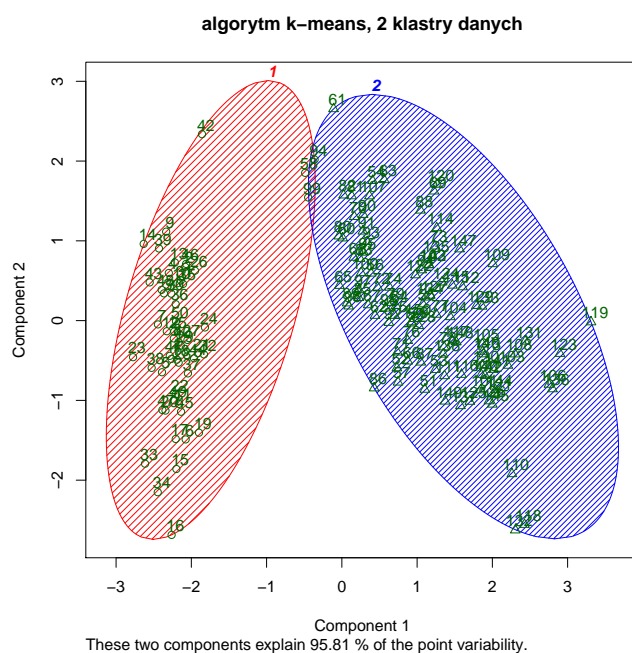
Stosując algorytmy uczenia nienadzorowanego, nie wiemy w jaki sposób klasyfikować dane trenujące. W tym celu możemy wykorzystać algorytm k-means[9], który na podstawie określonych cech grupuje dane w podaną liczbę klastrów. Poddając przedstawione dane działaniu algorytmu k-means, uzyskujemy wynik przedstawiony na rys. 2.2 oraz rys. 2.3.

Na rys. 2.2 można zaobserwować, że algorytm zadziała z dużą skutecznością rozróżniając pomiędzy dwoma gatunkami dla dwóch klastrów, tj. pomiędzy gatunkiem setosa, a gatunkami versicolor i virginica, jednak algorytm będzie miał problem rozróżnić pomiędzy gatunkami versicolor i virginica. W tym celu można zbiór podzielić na trzy klastry jak na rys.2.3, ale nie gwarantuje to wystarczającej skuteczności przydziału.

Tym samym można zaobserwować, że gatunek versicolor i virginica nie są rozróżnialne, używając przedstawionych danych wejściowych.

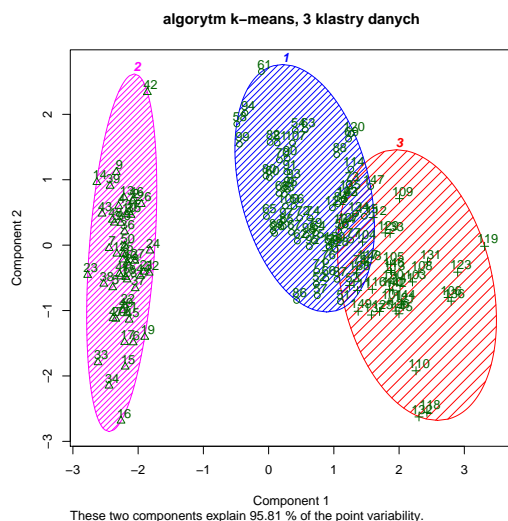
Rozważając sytuację, w której dodana by została duża liczba dodatkowych wpisów, różnych od aktualnych danych, można by było zaobserwować pojawienie się kolejnego klastra danych, co oznacza pojawienie się nowego gatunku kwiatu.

Wynik działania algorytmu k-means



Rys. 2.2. Wynik działania algorytmu klasyfikacji k-means dla 2 klastrów kwiatów irysów w zależności od szerokości i długości płatków kwiatu. Źródło: Opracowanie własne

Wynik działania algorytmu k-means

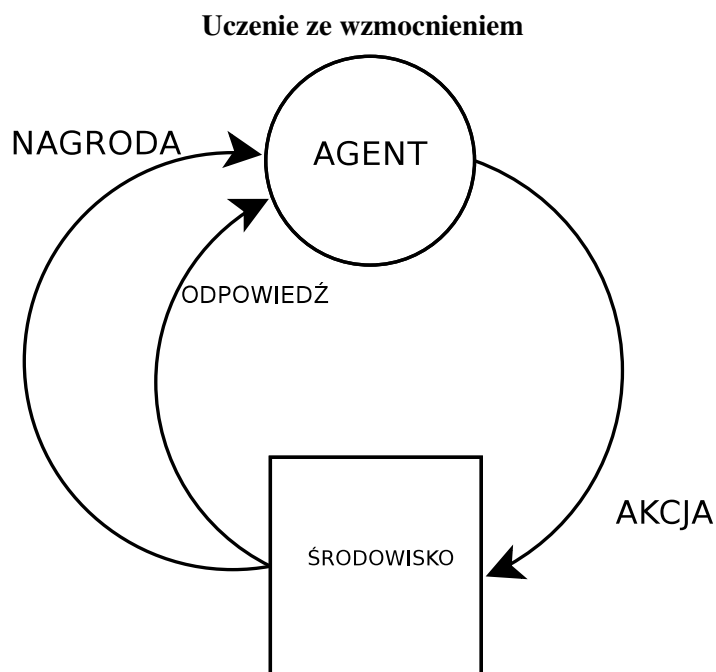


Rys. 2.3. Wynik działania algorytmu klasyfikacji k-means dla 3 klastrów kwiatów irysów w zależności od szerokości i długości płatków kwiatu. Źródło: Opracowanie własne

2.1.3. Uczenie ze wzmocnieniem

Uczenie ze wzmocnieniem jest obszarem uczenia maszynowego, w którym inteligentny agent wykonując akcję w otoczeniu otrzymuje odpowiednią nagrodę. Robot decyduje jaka akcja zostanie wykonana

na podstawie istniejącej polityki. Algorytmy uczenia ze wzmocnieniem nie wymagają danych trenujących by wykonać swoje zadanie. Na podstawie wiedzy o stanie otoczenia, możliwych do wykonania akcji w danym stanie oraz funkcji nagrody, która agent otrzymuje za wykonanie danych akcji, robot, ucząc się metodą prób i błędów odkrywa sposoby wybierania kolejnych akcji (polityki) tak, aby osiągnąć swój cel.



Rys. 2.4. Schemat uczenia ze wzmocnieniem. Źródło: Opracowanie własne

Na rys. 2.4 przedstawiono schemat współpracy agenta ze środowiskiem. Agent podejmując akcje w otoczeniu dostaje nagrodę za podjętą decyzję w danym stanie, oraz informacje o nowym, aktualnym stanie środowiska (otoczenie może ulec zmianie po wykonaniu przez robota akcji).

Algorytmy uczenia ze wzmocnieniem dotyczą specyficznych zadań występujących w uczeniu maszynowym. Problem polega na odnalezieniu przez agenta optymalnej akcji do wykonania na podstawie wiedzy o aktualnym jego stanie. W przypadku, gdy to działanie jest powtórzone, możemy mówić o MDP (Procesy decyzyjne Markowa).

Procesem decyzyjnym Markowa nazywamy model środowiska w algorytmach uczenia ze wzmocnieniem, oznaczamy go jako:

$$MDP = (S, A, P(s, s'), R(r, r'), \gamma)$$

gdzie

1. S to zbiór stanów
2. A zbiór możliwych akcji
3. $P(s, s')$ prawdopodobieństwo przejścia do stanu s' ze stanu s po wykonaniu akcji a , w czasie t

4. $R(r, r')$ nagroda po przejściu do stanu s' ze stanu s
5. γ współczynnik regulujący stosunek między nagrodami przewidywanymi w późniejszym czasie i otrzymywanymi aktualnie, $0 < \gamma < 1$

Zadaniem MDP jest uzyskanie takiej polityki robota, która zmaksymalizuje sumę otrzymywanych nagród.

$$\sum_{t=0}^{\infty} \gamma^t R_{a_t}(s_t, s_{t+1})$$

Algorytmy uczenia ze wzmocnieniem znalazły zastosowanie między innymi w implementacji autonomicznego helikoptera[10], udoskonaleniu działania windy[11], usprawnienia sygnalizacji świetlnej[12] czy budowie inteligentnego robota[13].

2.2. Podsumowanie

Uczenie nienadzorowane posiada tę przewagę nad uczeniem nadzorowanym, że nie wymaga oznaczenia danych trenujących oraz potrafi wykryć ukryte zależności pomiędzy cechami przykładów. Jednak w odróżnieniu od obydwu metod, uczenie ze wzmocnieniem jest odmienne w tym, że nie wymaga podawania jakichkolwiek danych przykładowych, a jedynie na podstawie zdefiniowanych nagród i akcji, agent uczy się odpowiedniego działania. By zaimplementować robota uczącego się wykonywania akcji w otoczeniu wybrano implementacje metod uczenia ze wzmocnieniem, ze względu na charakterystykę środowiska odpowiadającą wymaganiom procesów decyzyjnych Markowa.

3. Implementacja

Agent mając wiedzę na temat

- akcji, które może wykonać,
- stanu, w jakim się znajduje,
- zasad, definiujących warunki otrzymania nagrody,
- polityki wybierania akcji

jest w stanie dostosować swoje kolejne działania (wybrane akcje), tak aby uzyskać jak najlepszy wynik w kolejnych iteracjach symulacji. Celem agenta jest dotarcie do punktu końcowego, zdobywając jak największą nagrodę.

Po osiągnięciu celu symulacja jest resetowana, jej wynik jest zapisywany. Robot ucząc się na podstawie poprzednio dokonanych wyborów, buduje tablicę mapującą stan-akcja $Q(S, A)$ do konkretnej wyliczonej wartości. W programie do przedstawienia wartości tablicy stan-akcja $Q(S, A)$ wykorzystano wielowymiarową tablicę.

3.1. Tablica stan-akcja $Q(S, A)$

Tablica $Q(S, A)$ jest to struktura danych przechowująca wszystkie stany i możliwe do wykonania w nich akcje. Każdej z tych akcji przyporządkowany jest odpowiednia wartość, która jest poddawana modyfikacji przez odpowiedni algorytm uczenia ze wzmocnieniem. Ostatecznym celem działania algorytmu jest uzyskanie takich informacji w strukturze $Q(S, A)$, aby robot korzystając z polityki optymalnej (tzn. zawsze wybierał najbardziej korzystne akcje), wykonał cel z jak najlepszym wynikiem.

TODO PRZYKŁADOWA REPREZENTACJA

3.2. Polityka wyboru akcji

Agent mając do dyspozycji jedną z możliwych akcji a ze zbioru wszystkich akcji A , musi podjąć decyzję, która akcję ma wybrać jako następną. W tej implementacji rozważane są dwie polityki

1. ϵ -greedy policy
2. optimal policy

3.2.1. Równowaga pomiędzy eksploracją, a eksploatacją

Od wyboru polityki zależy stopień eksploracji algorytmu. Gdyby agent za każdym razem wybierał najbardziej korzystną według niego akcję, byłby podatny na zakleszczenie w nieoptymalnym rozwiązaniu. Wprowadzając pewną losowość wyboru akcji agent, będzie w stanie odkryć nowe rozwiązania, nawet jeżeli według jego aktualnej wiedzy znalazł to najbardziej optymalne. Powstało dużo prac dokładniej badających to zagadnienie, patrz ??, ??.

3.2.2. ϵ -greedy policy

Polityka ϵ -greedy jest strategią wyboru akcji polegającą na wyborze losowej akcji z prawdopodobieństwem ϵ , a w przeciwnym wypadku wybranie najbardziej korzystnej akcji. Agent wybiera akcję, która według niego będzie najbardziej korzystna z prawdopodobieństwem $1 - \epsilon$, a w przeciwnym wypadku wybierze akcję losową. Parametr ϵ jest regulowany zależnie od potrzeb i musi spełniać zależność $0 < \epsilon < 1$.

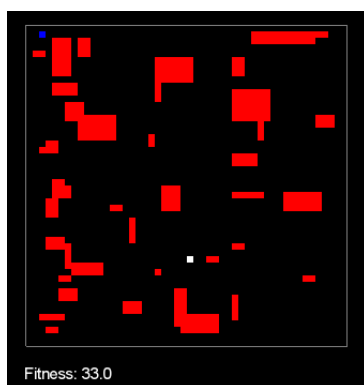
3.2.3. Optimal policy

Polityka optymalna, agent za każdym razem wybiera najkorzystniejszą według jego wiedzy akcję.

3.3. Symulacja graficzna

Jako symulację graficzną działania algorytmu wykorzystano figury geometryczne reprezentujące: agenta, przeszkody, granice i cel. Agent porusza się na dwuwymiarowej przestrzeni o wymiarach $50 * 50$ pikseli.

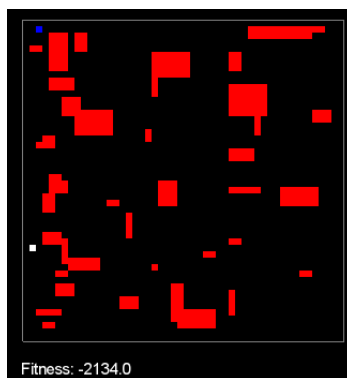
Przykładowy stan środowiska w symulacji przedstawiono na rys. 3.1, 3.2.



Rys. 3.1. Przykładowy stan symulacji. Źródło: Opracowanie własne

Na rys. 3.1, 3.2, można zauważyć następującą reprezentację obiektów jako kolor:

- czerwony - przeszkoda,



Rys. 3.2. Przykładowy stan symulacji. Źródło: Opracowanie własne

- biały - agent,
- niebieski - nagroda,
- szary - granica.

3.4. Reprezentacja stanu

Stan w jakim znajduje się aktualnie robot jest przedstawiony w postaci siatki 3×3 . Agent jest świadomy obiektów znajdujących się nad nim, pod nim oraz po jego prawej i lewej stronie. Agent może napotkać następujące obiekty

- przeszkoda (O), robot **może** wejść na przeszkodę, jednak traci za to określoną liczbę punktów;
- nagroda (P), stan końcowy. Gdy agent osiąga cel symulacja zapisuje wynik działania i uruchamia się kolejny epizod nauki robota;
- granica (B), granica jest obiektem wyznaczającym pole symulacji, dlatego też nie jest możliwe przemieszczenie się na nią;
- puste pole (E), pole nie zawierające żadnego z pozostałych obiektów;

Przykładowe stany w jakich może się znaleźć robot przedstawiono na rys. 3.3, 3.4, 3.5 i 3.6.

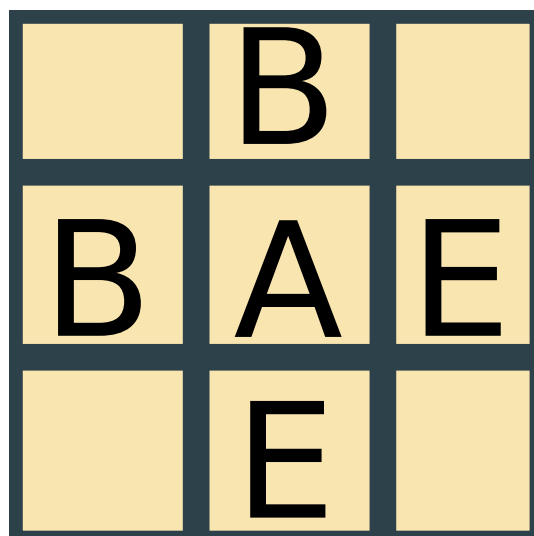
3.4.1. Akcje

Robot znajdując się w dowolnym stanie nie graniczącym z obiektem B (granica), jest w stanie wykonać jedną z czterech akcji:

1. ruch do góry,
2. ruch do dołu,



Rys. 3.3. Agent graniczy od górnej strony z przeszkodą, natomiast w pozostałych kierunkach znajdują się puste pola. Źródło: Opracowanie własne

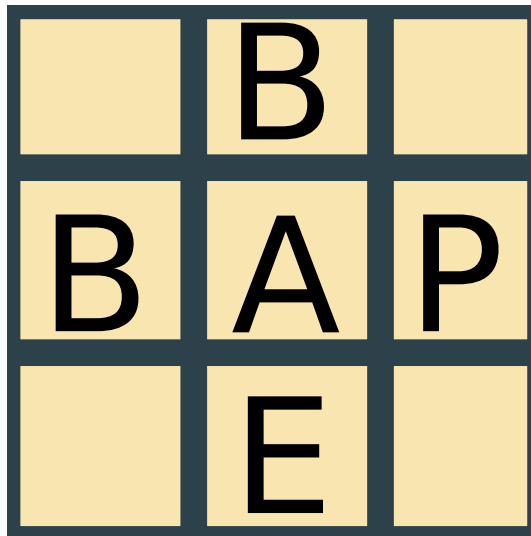


Rys. 3.4. Agent graniczy od góry i lewej strony z granicą, oznacza to, że robot znajduje się w górnym lewym rogu mapy, a jedynymi akcjami jakie może wykonać to ruch w prawo i ruch w dół. Źródło: Opracowanie własne

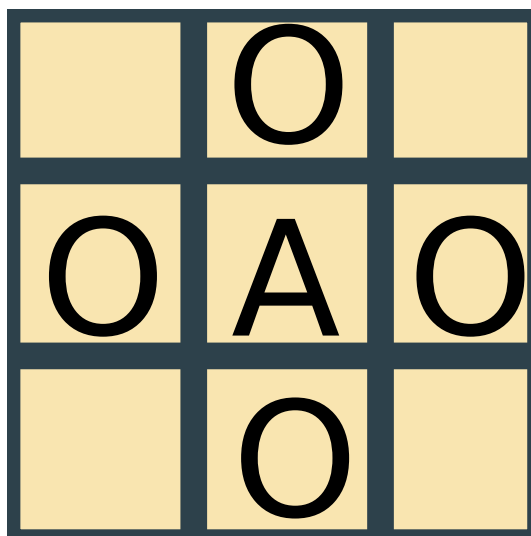
3. ruch w prawo,

4. ruch w lewo.

Agent **nie może** przemieścić się na granicę, w związku z czym w przypadku sąsiedztwa z granicą, robot posiada ograniczony zakres ruchów.



Rys. 3.5. Agent graniczy od góry i lewej strony z granicą, a od prawej z stanem końcowym. Agent wykonując ruch w prawo osiągnie swój cel, powodując tym samym zakończenie i zresetowanie symulacji. Źródło: Opracowanie własne



Rys. 3.6. Agent graniczy z każdej strony z przeszkodą, jakkolwiek ruch powoduje utratę punktów. Źródło: Opracowanie własne

3.4.2. Wybór algorytmów

3.4.2.1. Q-learning

Algorytm Q-learning [watkins1992q] jest off-policy

3.4.2.2. SARSA

3.5. Wykorzystane technologie

Java 8 kod został napisany w języku Java. Do generowania grafiki została wykorzystana biblioteka LibGDX. W celu notowania wyników do plików kalkulacyjnych została wykorzystana biblioteka Apache POI.

3.6. Napotkane problemy

Przy projektowaniu aplikacji, największy problem stanowiło odpowiednie przedstawienie aktualnego stanu środowiska i robota. Nieodpowiednie odwzorowanie stanu, może powodować niepoprawną naukę agenta. Zakładając rozmiar mapy $50 * 50$, można zauważyć, że przy naiwnej reprezentacji stanu, tzn. takiej, w której każda możliwa kombinacja stanu mapy jest reprezentowana osobno, powstaje ogromna ilość możliwych akcji w tablicy wartości $Q(S, A)$. Ostatecznie postanowiono przedstawić stan w postaci siatki $3 * 3$. Reprezentacja stanu Testowanie działania algorytmów

3.7. Wynik działania

4. Podsumowanie

Bibliografia

- [1] John McCarthy. „What is artificial intelligence”. W: *URL: <http://www-formal.stanford.edu/jmc/whatisai.html>* (2007), s. 38.
- [2] Poole David L. i Mackworth Alan K. *Artificial Intelligence: Foundations of Computational Agents*. New York, NY, USA: Cambridge University Press, 2010. ISBN: 0521519004, 9780521519007.
- [3] Pedro Domingos. „A few useful things to know about machine learning”. W: *Communications of the ACM* 55.10 (2012), s. 78–87.
- [4] Loredana Firté, Camelia Lemnaru i Rodica Potolea. „Spam detection filter using KNN algorithm and resampling”. W: *Proceedings of the 2010 IEEE 6th International Conference on Intelligent Computer Communication and Processing*. IEEE. 2010, s. 27–33.
- [5] Muhammad N Marsono, M Watheq El-Kharashi i Fayez Gebali. „Binary LNS-based naïve Bayes inference engine for spam control: noise analysis and FPGA implementation”. W: *Computers & Digital Techniques, IET* 2.1 (2008), s. 56–62.
- [6] R Deepa Lakshmi i N Radha. „Spam classification using supervised learning techniques”. W: *Proceedings of the 1st Amrita ACM-W Celebration on Women in Computing in India*. ACM. 2010, s. 66.
- [7] Irena Koprinska i in. „Learning to classify e-mail”. W: *Information Sciences* 177.10 (2007), s. 2167–2187.
- [8] Wenjuan Li i in. „Towards designing an email classification system using multi-view based semi-supervised learning”. W: *Trust, Security and Privacy in Computing and Communications (Trust-Com), 2014 IEEE 13th International Conference on*. IEEE. 2014, s. 174–181.
- [9] Adam Coates i Andrew Y Ng. „Learning feature representations with k-means”. W: *Neural Networks: Tricks of the Trade*. Springer, 2012, s. 561–580.
- [10] Pieter Abbeel i in. „An application of reinforcement learning to aerobatic helicopter flight”. W: *Advances in neural information processing systems* 19 (2007), s. 1.
- [11] A Barto i RH Crites. „Improving elevator performance using reinforcement learning”. W: *Advances in neural information processing systems* 8 (1996), s. 1017–1023.

- [12] Marco Wiering i in. „Multi-agent reinforcement learning for traffic light control”. W: *ICML*. 2000, s. 1151–1158.
- [13] Hajime Kimura, Tom Yamashita i Shigenobu Kobayashi. „Reinforcement learning of walking behavior for a four-legged robot”. W: *Decision and Control, 2001. Proceedings of the 40th IEEE Conference on*. T. 1. IEEE. 2001, s. 411–416.