

Nazwa projektu: **Narzędzie eksploracji danych**

Autorzy specyfikacji:

<i>imię i nazwisko</i>	<i>nr indeksu</i>	<i>grupa</i>	<i>adres email</i>
Kornel Lewandowski	inf94334	czw1145	kornel.lewandowski@gmail.com
Małgorzata Napieraj	inf94347	czw1145	napieraj.malgorzata@gmail.com

Opis projektu

Celem projektu jest napisanie prostego w obsłudze narzędzia eksploracji danych. Narzędzie ma pozwalać na klasyfikację czy ekstrakcję reguł decyzyjnych ze zbioru danych wybranym algorytmem. Dodatkowo, ma umożliwiać wstępne przetwarzanie zbioru wejściowego – dyskretyzację, normalizację, selekcję atrybutów. Ze względu na to, że docelowymi użytkownikami mają być osoby nie specjalizujące się w dziedzinie *data mining*, interfejs użytkownika ma pozwalać na jak najłatwiejszą obsługę, wliczając w to rozbudowaną pomoc kontekstową oraz informowanie o konsekwencjach wykonywanych akcji czy bieżącym stanie procesu eksploracji. System będzie udostępniał kilka wybranych algorytmów eksploracji danych, włączając w to możliwość zmiany parametrów dotyczących danego algorytmu. Ponadto, narzędzie ma umożliwiać import oraz eksport danych w kilku najpopularniejszych formatach plików (XLS, CSV).

Funkcjonalność projektu

Spis wymaganych funkcjonalności z podziałem na obowiązkowe (min. 10) i opcjonalne (min. 5).

<i>Lp</i>	<i>Opis funkcjonalności</i>	<i>Typ</i>	<i>Status</i>
1	Normalizacja i dyskretyzacja danych	obowiązkowa	
2	Selekcja atrybutów	obowiązkowa	
3	Budowa klasyfikatora na podstawie zbioru uczącego (algorytm <i>NaiveBayes</i>)	obowiązkowa	
4	Budowa klasyfikatora na podstawie zbioru uczącego (algorytm <i>J48</i>)	obowiązkowa	
5	Budowa klasyfikatora na podstawie zbioru uczącego (algorytm <i>IBk</i>)	obowiązkowa	
6	Klasyfikacja przykładów na zbudowanym klasyfikatorze	obowiązkowa	
7	Manipulacja parametrami algorytmów klasyfikacji	obowiązkowa	
8	Import danych w formacie XLS/CSV/ARFF	obowiązkowa	
9	Eksport danych w formacie XLS/CSV/ARFF	obowiązkowa	
10	Pomoc kontekstowa	obowiązkowa	
11	Import/eksport modelu klasyfikatora	opcjonalna	
12	Import danych w formacie XML	opcjonalna	
13	Eksport danych w formacie XML	opcjonalna	
14	Udostępnienie dodatkowego algorytmu selekcji atrybutów (<i>GainRatio</i>)	opcjonalna	
15	Budowa klasyfikatora na podstawie zbioru uczącego (algorytm <i>DTNB</i>)	opcjonalna	

Koncepcja realizacji

Koncepcja realizacji projektu polega na wykorzystaniu istniejącej już biblioteki stworzonej do eksploracji danych oraz stworzenia dla niej wygodnego interfejsu użytkownika. Projekt zakłada wykorzystanie narzędzi eksploracji danych zawartych w pakiecie *Weka*, który udostępnia takie funkcje jak wstępne przetwarzanie zbioru wejściowego, zbudowanie klasyfikatora oraz samą klasyfikację. Pozostałe funkcje, czyli import oraz eksport danych, interfejs użytkownika z pomocą kontekstową i użycie klasyfikatorów zostaną zaimplementowane w ramach projektu.

Rdzeń aplikacji napisany zostanie w języku Java i wykorzystywał będzie biblioteki pakietu *Weka*. Interfejs użytkownika będzie wykorzystywał przeglądarkę internetową, co pozwoli stworzyć szybko i łatwo atrakcyjne graficznie i funkcjonalne ekrany aplikacji. W celu zapewnienia wysokiej jakości interfejsu wykorzystane zostaną technologie związane z aplikacjami przeglądarkowymi: HTML, CSS, JavaScript oraz biblioteka jQuery (ewentualnie jQueryUI).

Połączenie między częścią interfejsu uruchamianą w przeglądarce a rdzeniem aplikacji (programem działającym w tle) zapewni zestaw narzędzi Jetty, który będzie w stanie obsłużyć wysyłane z przeglądarki żądania i odpowiedzieć na nie komunikatami w odpowiednim formacie (JSON).

System w wersji podstawowej będzie obsługiwał 3 różne formaty danych (ARFF, CSV, XLS) oraz udostępniał 3 różne algorytmy eksploracji (*NaiveBayes*, *J48*, *Ibk*) i 1 algorytm selekcji atrybutów (*InfoGain*).

Ogólny opis poszczególnych grup funkcjonalności

Normalizacja i dyskretyzacja danych – na etapie wczytywania danych użytkownik będzie mógł dokonać tych operacji na wybranych przez niego atrybutach.

Selekcja będzie polegała na wybraniu przez użytkownika atrybutów, które mają być wyłączone ze zbioru cech opisujących dane. Dodatkowo, będzie istnieć możliwość uszeregowania atrybutów względem zadanej miary jakości i na jej podstawie usunięcia podzbioru atrybutów.

Budowa klasyfikatorów oparta będzie o gotowe rozwiązania znajdujące się w pakiecie *Weka*. Każdy z nich sterowany będzie za pomocą parametrów dostępnych w interfejsie użytkownika (oczywiście to, jakie parametry będą ustalane, jest zależne od wybranego algorytmu klasyfikującego). Klasyfikatory zaprojektowane będą w ten sposób, że dane wejściowe (wstępnie przetworzone przez selekcję atrybutów, dyskretyzację wartości i tym podobne) przekonwertowane zostaną do postaci pliku ARFF

(domyślny format pliku pakietu *Weka*) i poddane budowaniu modelu klasyfikatora. Na tej podstawie otrzymany zostanie model potrafiący zaklasyfikować nowe przypadki podane przez użytkownika.

Import i eksport danych oparty będzie na prostych konwerterach zrealizowanych za pomocą operacji tekstowych. W zależności od skomplikowania konwersji, przewiduje się zarówno wdrożenie własnych rozwiązań, jak i wykorzystanie zewnętrznych bibliotek.

Pomoc kontekstowa będzie oparta przede wszystkim o komunikaty uświadamiające użytkownikowi konsekwencje podejmowanych akcji oraz informujące o obecnym stanie eksploracji.

Materialy źródłowe

- <http://www.cs.waikato.ac.nz/ml/weka/documentation.html>
- <http://api.jquery.com/>
- https://wiki.eclipse.org/Jetty/Tutorial/Embedding_Jetty
- <http://www.cs.put.poznan.pl/jstefanowski/aed/preprocessing1.pdf>
- <http://www.cs.put.poznan.pl/jstefanowski/aed/TPDocenaklasyfikatorow.pdf>
- <http://www.cs.put.poznan.pl/jstefanowski/aed/TPD-grupowanie2011.pdf>
- <http://www.cs.put.poznan.pl/jstefanowski/ml/W2drzewaJS2010.pdf>