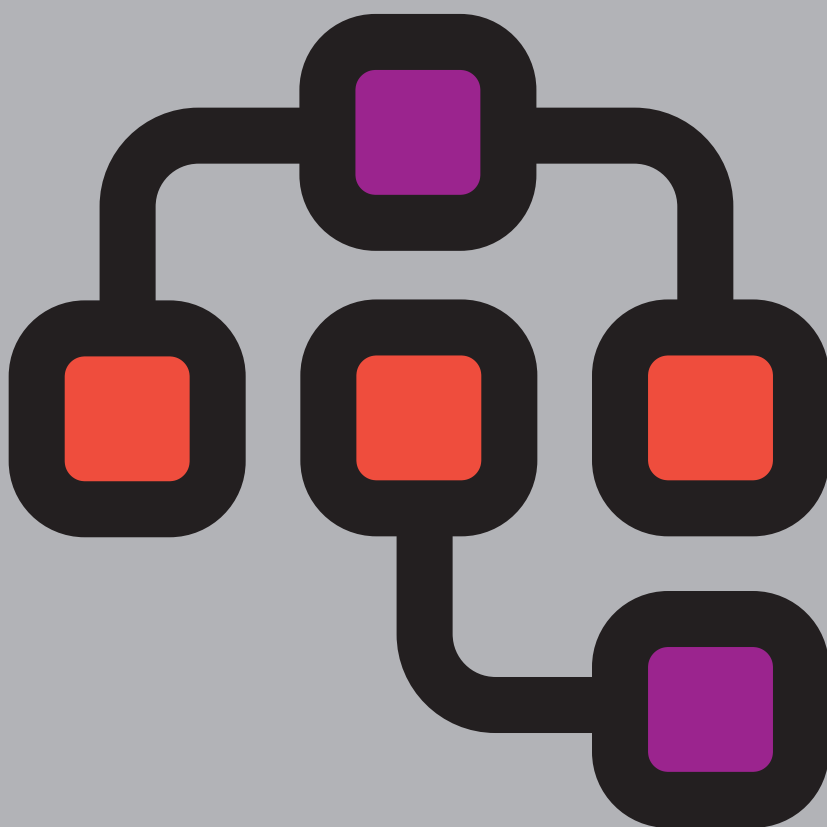


Mirostawa Lasek, Marek Pęczkowski

Enterprise Miner

Wykorzystywanie narzędzi

Data Mining w systemie *SAS*



Enterprise Miner

Mirostawa Lasek, Marek Pęczkowski

Enterprise Miner

Wykorzystywanie narzędzi

Data Mining w systemie *SAS*



Warszawa 2013

Recenzent
prof. dr hab. Witold Chmielarz

Redaktor prowadzący
Małgorzata Yamazaki

Redakcja i korekta
Izabela Mika

Redakcja techniczna
Zofia Kosińska

Projekt okładki i stron tytułowych
Wojciech Markiewicz

Skład i łamanie
Krzysztof Biesaga

ISBN 978-83-235-2770-1 (PDF)

© Copyright by Wydawnictwa Uniwersytetu Warszawskiego, Warszawa 2013

Publikacja dofinansowana przez Wydział Nauk Ekonomicznych
i Rektora Uniwersytetu Warszawskiego

Nowa od wydawcy

W książce przedstawiono zasady posługiwania się programem *SAS Enterprise Miner 6.2* bez oznaczania go symbolem znaku towarowego. Autorzy i wydawca oświadczają, że pojawiające się w tej książce nazwy tego i wszystkich innych produktów ze znakami towarowymi zostały wykorzystane w sposób przyjęty w wydawnictwach, korzystny dla właścicieli znaków towarowych i niebędący pod żadnym względem zamierzonym naruszeniem praw do tych znaków.

Wydawnictwa Uniwersytetu Warszawskiego
00-497 Warszawa, ul. Nowy Świat 4
www.wuw.pl; e-mail: wuw@uw.edu.pl
Dział Handlowy: tel +48 22 55-31-333
e-mail: dz.handlowy@uw.edu.pl
Księgarnia internetowa: www.wuw.pl/ksiegarnia

Wydanie 1

Spis treści

Wstęp	7
1. Zasady posługiwania się programem Enterprise Miner	17
1.1. Rozpoczęcie pracy z programem i tworzenie projektów eksplo- racji danych	17
1.2. Organizacja zbiorów danych wykorzystywanych w analizach eksploracji danych	19
1.3. Tworzenie diagramów i zarządzanie diagramami analizy danych	22
2. Zagadnienia metodyczne przeprowadzania eksploracji danych	24
2.1. Metodyki wykorzystywane na potrzeby analizy danych	24
2.2. Narzędzia eksploracji danych wspierające analizy w poszczególnych etapach metodyki SEMMA programu Enterprise Miner	28
2.3. Ogólne reguły budowania diagramów wg metodyki SEMMA i przykład diagramu	34
3. Przygotowanie danych do analiz Data Mining	37
3.1. Wstępna analiza danych wejściowych	37
3.2. Narzędzia Multiplot i StatExplore w analizie danych	46
3.3. Losowanie próby danych (wykorzystanie węzła Sample)	54
3.4. Podział danych (Data Partition)	56
3.5. Filtrowanie danych (węzeł <i>Filter</i>)	59

3.6.	Wybór zmiennych na potrzeby budowy modeli	62
3.7.	Przeprowadzanie transformacji zmiennych	68
3.8.	Rozwiązywanie problemu brakujących wartości	70
3.9.	Wartości nietypowe	75
4.	Metody prognozowania	77
4.1.	Kryteria oceny modeli predykcyjnych	77
4.2.	Regresja, sieci neuronowe i drzewa decyzyjne jako narzędzia prognozowania	97
4.3.	Sporządzanie prognoz – generowanie i wykorzystywanie kodu skoringowego	173
5.	Metody grupowania	180
5.1.	Analiza skupień	180
5.2.	Sieci neuronowe Kohonena	216
6.	Analiza asocjacji i sekwencji	236
7.	Grupowanie zmiennych	255
	Literatura	264

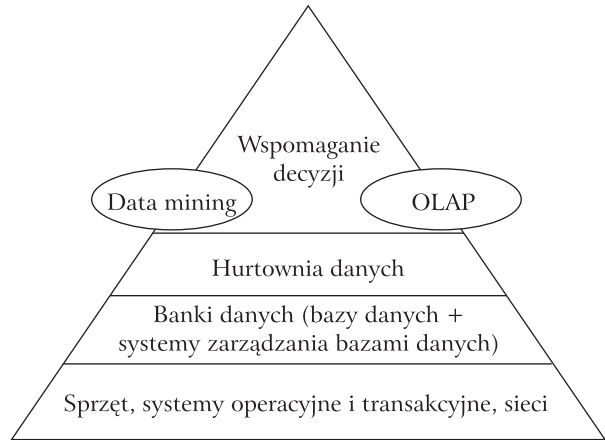
Wstęp

Praca przedstawia podstawy użytkowania programu *Enterprise Miner* firmy SAS oraz krótki opis metod eksploracji danych (*data mining*), których wykorzystanie umożliwia ten program. Pod pojęciem eksploracji danych rozumie się metody statystyczne i metody sztucznej inteligencji umożliwiające odkrywanie nieznanych zależności między danymi w nagromadzonych zbiorach danych [M. Lasek, 2004]. Są to metody, które pozwalają z danych tworzyć wiedzę, tzn. znajdować zależności, wzorce, trendy „ukryte” w danych.

Rozwój metod eksploracji danych związany jest z rozwojem wielu dziedzin: informatyki, statystyki, ekonometrii, ekonomiki i organizacji przedsiębiorstw, zarządzania finansami, teorii i narzędzi wnioskowania w warunkach niepewności. To właśnie ich rozwój, a w szczególności burzliwy rozwój technologii informatycznych, spowodował, że coraz częściej punktem wyjścia w procesach decyzyjnych są dane, a nie hipotezy statystyczne lub koncepcje klasycznych modeli ekonometrycznych. Współczesny sprzęt i oprogramowanie umożliwiają gromadzenie i analizę olbrzymich zbiorów danych, obejmujących bazy mierzone dziesiątkami lub więcej gigabajtów [M. Lasek, M Pęczkowski, 2010(c)].

W komputerowych magazynach danych zwanych hurtowniami danych (*data warehouses*), tworzonych w przedsiębiorstwach przemysłowych, bankach, agencjach ubezpieczeniowych, firmach usługowych, znajdują się olbrzymie ilości danych. Dane te trudno poddają się znanym metodom analiz statystycznych i ekonometrycznych, tak aby służyło to budowaniu wiedzy, która mogłaby być przydatna do wspomaganie w podejmowaniu decyzji w zarządzaniu, dokonywaniu wyborów ekonomicznych, znajdowaniu reguł i uogólnień. Narzędzia raportowania *OLAP* (*OnLine Analytical Processing*), pomimo swoich niewątpliwych zalet: wielowymiarowości i wielopoziomowości (umożliwiającej

przechodzenie od ogółu do szczegółu i z powrotem) oraz pomimo analiz zapewniających wgląd w dane z różnych perspektyw badawczych, wynikających z potrzeb różnych użytkowników, okazały się nie w pełni wystarczające. Metody eksploracji danych stanowią analityczne rozszerzenie technik *OLAP* i prezentują podejście do analizy danych służące zasadniczo innym celom niż *OLAP*. Polegają one raczej nie na raportowaniu, ale na zdobywaniu nowej wiedzy. Miejsce *OLAP* i eksploracji danych w piramidzie systemów wspomagania decyzji ilustruje rys. 0.1.



RYSUNEK 0.1. Miejsce *OLAP* i metod eksploracji danych w piramidzie systemów wspomagania decyzji

Źródło: [Z. Chen, 2001, s. 360].

W języku polskim angielski termin *data mining methods* jest tłumaczony jako metody eksploracji danych, odkrywania wiedzy w bazach danych, zgłębiania danych, eksploatacji danych, drażenia danych. Jednak chyba najczęściej, zarówno w polskiej literaturze, na konferencjach, jak i wśród praktyków używana jest nazwa angielska: *data mining*.

Na temat metod eksploracji danych istnieje już od dość dawna bogata literatura (por. np. [M.J.A. Berry, G.S. Linoff, 2000; M.J.A. Berry, G.S. Linoff, 2004; P. Cabena i in., 1998; J. Han, M. Kamber, J. Pei, 2012; D.T.Larose, 2006; D.T.Larose, 2008; O. Maimon, L. Rokach (eds.), 2005; Z. Markov, D.T. Larose, 2009; R. Matignon, 2007; I.H. Witten, E. Frank, M.A. Hall, 2011; N. Ye (ed.), 2003]) dotycząca ich podstaw oraz wykorzystania w różnych zastosowaniach. Ukazują się także prace polskich autorów w języku polskim, głównie w postaci monografii, w polskich czasopismach informatycznych i w Internecie. Przygotowywane są materiały kursowe w języku polskim (por. np. [M. Pęczkowski, 2011 (b)]).

Należy podkreślić, że w Polsce prace dotyczące metod eksploracji danych są coraz intensywniej rozwijane. Są one prowadzone w licznych polskich ośrodkach. Należą do nich m.in. Szkoła Główna Han-

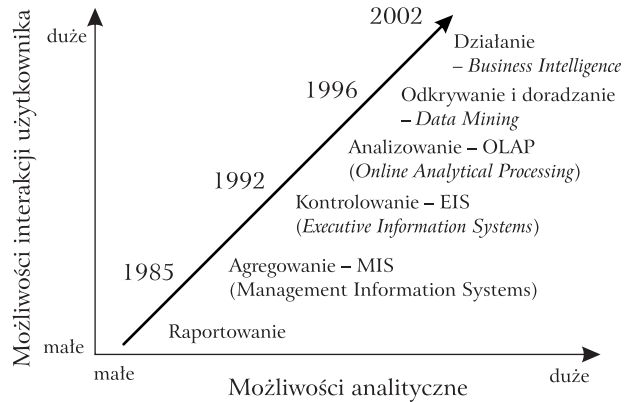
dłowa w Warszawie, Uniwersytet Gdański, Uniwersytet Ekonomiczny we Wrocławiu, Uniwersytet Ekonomiczny w Katowicach, Uniwersytet Technologiczno-Przyrodniczy w Bydgoszczy, Politechnika Białostocka, Uniwersytet w Białymstoku, Uniwersytet Ekonomiczny w Krakowie.

Nasze próby wykorzystywania metod eksploracji danych wskazały na ich dużą przydatność w przeprowadzaniu analizy danych [M. Lasek, M. Pęczkowski, 2010 (c)]. Dotyczyły m.in. analizy zróżnicowania pięciuset największych firm Rzeczypospolitej [M. Lasek, M. Pęczkowski, 2008 (b)], przeprowadzania segmentacji klientów na potrzeby prowadzenia kampanii reklamowych [M. Kutera, M. Lasek, 2010], analizowania i prognozowania kondycji ekonomicznej przedsiębiorstw [M. Lasek, 2007 (a)], analizy działalności inwestycyjnej gospodarstw agroturystycznych [M. Lasek, E. Nowak, M. Pęczkowski, 2008], analiz finansowych przedsiębiorstw [M. Lasek, M. Pęczkowski, 2008 (a)], przewidywania groźby upadłości lub konieczności prowadzenia postępowania układowego przedsiębiorstw [M. Lasek, M. Pęczkowski, D. Wierzbą, 2009]. Wymienione prace prezentują wyniki prowadzonych przez nas analiz i różnorodnych badań, chociaż oczywiście w literaturze (zwłaszcza anglojęzycznej) można znaleźć wiele prac opisujących zarówno przydatność poszczególnych metod eksploracji danych do różnych celów, jak i całościowe projekty zastosowań. Często powołujemy się na własne artykuły dotyczące zagadnień eksploracji danych także z tego względu, że czytelnik może tam znaleźć wskazane przez nas dość liczne pozycje literatury odnoszące się do szczegółowych kwestii i problemów przedstawianych w oddzielnych, podejmujących odrębne tematy artykułach, jak np. wyznaczanie liczby skupień w grupowaniu obiektów, graficzna ocena jakości modeli, grupowanie zmiennych, sporządzanie prognoz. Uznaliśmy, że nie byłoby celowe powtarzanie już tam przytaczanych czy też cytowanych pozycji.

Obecnie metody eksploracji danych znajdują zastosowania praktyczne w ramach systemów oprogramowania zwanych systemami *Business Intelligence* – rys. 0.2.

Metody eksploracji danych dzielone są w różny sposób. W przypadku kryterium celu stosowania można wyróżnić metody:

- klasyfikacji – umożliwiające przydział obiektów do z góry zdefiniowanych klas (podzbiorów); należy do nich np. analiza dyskryminacyjna, w której dąży się do znalezienia funkcji umożliwiającej przewidywanie przynależności nowego obiektu do danej klasy;
- regresji – obejmujące znajdowanie związków opisujących wpływ jednej lub większej liczby cech (zmiennych objaśniających) na wybraną cechę (zmienną objaśnianą);



RYSUNEK 0.2. Ewolucja od statycznych raportów do metod eksploracji danych i systemów *Business Intelligence*

Źródło: [N. Rasmussen i in., 2002, s. 5].

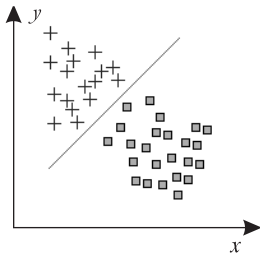
- grupowania (analizy skupień) – dotyczące podziału zbioru obiektów na skończoną liczbę grup w ten sposób, aby obiekty podobne do siebie znalazły się w tej samej grupie;
- odkrywania charakterystyk – polegające na znajdowaniu opisu grup obiektów za pomocą skończonej, możliwie małej liczby cech (charakterystyk) określanych często mianem profili (np. klientów bankowych, osób kupujących określone towary, użytkowników określonych typów komputerów);
- odkrywania asocjacji – dotyczące odkrywania związków między obiektami lub grupami obiektów opisanych przez wiele cech ilościowych lub jakościowych;
- odkrywania sekwencji – służące do odnajdywania kolejności następowania zdarzeń lub pojawiania się obiektów;
- wykrywania zmian i odchyłeń – polegające na poszukiwaniu wartości nietypowych (odstających, skrajnych), a także systematycznych błędów pomiaru.

Innym kryterium podziału metod eksploracji danych jest charakter związku lub podziału obiektów (cech), który może być liniowy lub nieliniowy.¹ Wyróżnia się związki liniowe oraz związki nieliniowe między obiektami, jak widać to na rys. 0.3. Przypadki *a*, *b* i *d* przedstawiają ujęcie liniowe w modelach analizy danych, chociaż linie z rys. 0.3*b* oraz 0.3*d* nie są liniami prostymi. Podział na metody liniowe i nieliniowe wynika z zależności stwarzanych przez parametry modeli (charakter dopasowania parametrów modeli).

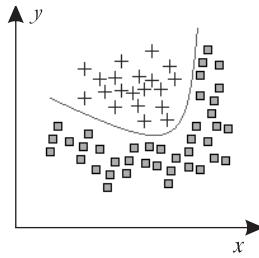
Celem naszej pracy jest umożliwienie czytelnikowi w możliwie krótkim czasie opanowania zasad użytkowania programu *Enterprise Miner* oraz poznania podstawowych metod eksploracji danych. Praca nie

¹ Przyjęta tu terminologia; liniowy, nieliniowy, wynika z dosłownego tłumaczenia z języka angielskiego. Wydaje się, że poprawniejsze byłyby terminy: związki jednowymiarowe oraz związki wielowymiarowe.

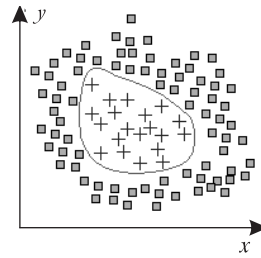
a) liniowość



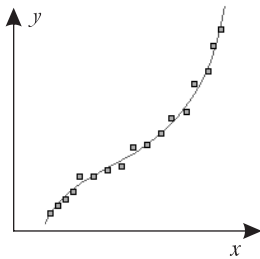
b) liniowość



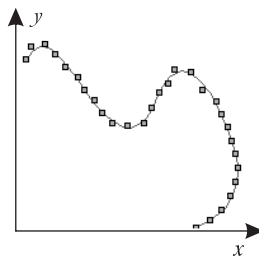
c) nieliniowość



d) liniowość



e) nieliniowość



RYSUNEK. 0.3. Liniowe i nieliniowe modele eksploracji danych

Źródło: na podstawie [H. Lohninger, 1999].

jest szczegółową instrukcją użytkownika programu *Enterprise Miner* ani podręcznikiem z zakresu metod eksploracji danych (takie podręczniki i materiały kursowe są opracowywane i udostępniane przez *SAS Institute*). Przedstawiamy jedynie podstawowe opcje programu oraz ogólną charakterystykę wybranych metod, zwracając uwagę na potrzeby w zakresie odpowiedniego przygotowania danych oraz na właściwą interpretację uzyskiwanych wyników. Niektóre możliwości programu, jak nam się wydaje rzadziej używane przez początkującego użytkownika, który dopiero zapoznaje się z programem i metodami, zostały jedynie zasygnalizowane – zainteresowany czytelnik może się z nimi zapoznać w dokumentacji programu dostarczanej przez firmę *SAS* lub uczestnicząc w specjalistycznych kursach organizowanych przez tę firmę. Aby umożliwić czytelnikowi wykorzystywanie tylko wybranych przez niego rozdziałów (bez potrzeby czytania całej pracy), musieliśmy powtórzyć niektóre informacje. W takich przypadkach staraliśmy się o dostosowanie opisów do potrzeb tematyki rozdziałów, w których są zamieszczone.

Możliwości programu i metod ilustrujemy przykładami analiz zbiorów danych z różnych dziedzin i źródeł:

- zbioru dotyczącego charakterystyki polskich gospodarstw domowych *F2007BD*, liczącego 37 121 obiektów (czyli gospo-

darstw domowych) oraz utworzonego na jego podstawie zbioru *F2007ZYWN*, przedstawiającego wydatki na żywność, alkohole i tytoń w gospodarstwach domowych, uwzględniającego 31 pozycji takich wydatków (zbiory danych pochodzą z badania budżetów gospodarstw domowych przeprowadzanego przez GUS i dotyczą danych z 2007 r.) [Budżety gospodarstw domowych ..., 2008; Metodologia badań budżetów ..., 2011; M. Pęczkowski, 2011 (a)],

- zbioru *HMEQ* dotyczącego udzielania kredytów przez bank, liczącego 5960 obiektów (klientów banku), zbioru *HMEQ_Score* o wybieranej liczbie obserwacji (klientów lub potencjalnych klientów) dla prowadzenia na bieżąco badań skoringowych analizy ryzyka niewywiązania się klientów z płatności (badania ich ewentualnej niewypłacalności) oraz zbioru danych *BANK* o usługach świadczonych przez bank, liczącego 32 367 obserwacji reprezentujących transakcje bankowe – usługi świadczone klientom przez bank, realizowane przez około 8000 klientów banku (te dwa zbiory danych są używane na kursach prowadzonych przez *SAS Institute Inc.*, np. [Applied Analytics Using ..., 2008]),
- zbioru *CHURN* z danymi dotyczącymi wykorzystywania usług przez klientów firmy telefonicznej, liczącego 3333 obiekty (klientów, którzy zrezygnowali lub nie zrezygnowali z usług firmy) oraz zbioru transakcji sprzedaży i zakupu towarów, którymi w analizowanym przypadku były rozmaite warzywa (w [D.T. Larose, 2006] podane są adresy do stron internetowych, gdzie można znaleźć zbiór *CHURN*; dane dotyczące handlu warzywami zamieszczone są bezpośrednio w tej książce).

Praca składa się z siedmiu rozdziałów.

W rozdziale 1. zapoznajemy czytelnika z podstawami posługiwania się programem *Enterprise Miner*. Opisujemy, jak rozpocząć pracę z programem, jak przygotować dane, które będą wykorzystywane w analizach, oraz jak utworzyć projekt analizy danych, który będzie złożony z diagramów definiujących poszczególne kroki składające się na analizę danych, począwszy od wprowadzenia danych i wstępnej ich analizy oraz obróbki, poprzez zastosowanie wybranych metod eksploracji danych, kończąc na interpretacji uzyskanych wyników.

W rozdziale 2. przedstawiamy metody przydatne do przeprowadzania analiz danych. Oddzielną część tego rozdziału (podrozdz. 2.2) poświęcamy metodyce analizy danych *SEMMA* (*Sample, Explore, Modify, Model, Assess*). Jest to oryginalna metodyka opracowana przez firmę *SAS Institute*, definiująca kolejne kroki i narzędzia analizy danych. Opisujemy udostępniane w *Enterprise Miner* narzędzia analizy danych, które są zalecane w ramach poszczególnych kroków metodyki. Narzędzia

te stanowią węzły diagramów analizy danych. Węzły na diagramach łączone są liniami zakończonymi strzałkami wskazującymi kolejność kroków analizy danych. W tym rozdziale omawiamy także sposoby budowania diagramu analizy danych i wskazujemy, jakich zasad budowy diagramów należy przestrzegać, aby zapewnić poprawność przeprowadzanych analiz.

Cały rozdział 3. poświęcony jest problematyce przygotowywania danych na potrzeby analizy związanej z eksploracją danych. Obejmuje on dostępne w *Enterprise Miner* metody wstępnej statystycznej analizy danych i opis narzędzi udostępnianych w celu przeprowadzenia takiej analizy. Przedstawiamy także zagadnienia losowania próby do przeprowadzania analiz, gdy uznamy, że nie ma potrzeby posłużenia się całym zbiorem. Omawiamy problem podziału danych na dane treningowe, walidacyjne i testowe (stosowanego na potrzeby budowania modeli) i przedstawiamy celowość jego wykonania. W tym rozdziale poruszamy także zagadnienia filtrowania danych, wyboru zmiennych, ze względu na które będziemy analizować dane, przeprowadzania transformacji zmiennych, zastępowania brakujących wartości wartościami wyliczonymi przez program lub przyjętymi przez użytkownika wartościami (tzw. imputacja danych) oraz wpływu wartości nietypowych (*outliers*) na wyniki analiz. Problematyce wyboru zmiennych, z uwagi na jej znaczenie w przypadkach wykorzystywania metod eksploracji danych, poświęcamy szczególną uwagę. Przedstawiamy specjalne narzędzie programu, które może być wykorzystywane, aby właściwie dobrać zmienne do tworzonego modelu – narzędzie o nazwie *Variable Selection*, które udostępnia metody doboru zmiennych. Ilustrujemy tu m.in. możliwość posłużenia się kryterium współczynnika determinacji R^2 oraz kryterium Chi^2 dla doboru zmiennych. W dalszej części pracy przedstawimy także inne metody, które mogą być pomocne w doborze zmiennych, mianowicie możliwości wykorzystywania drzew decyzyjnych oraz dokonywania selekcji zmiennych za pomocą zbudowanego przez użytkownika modelu regresji (liniowej lub logistycznej).

W rozdziale 4. omawiamy możliwości wykorzystania metod eksploracji danych na potrzeby prognozowania. Przedstawiamy zastosowanie w prognozowaniu metod: regresji liniowej i logistycznej, sieci neuronowych oraz drzew decyzyjnych. Podajemy przykłady zastosowania tych metod, posługując się programem *Enterprise Miner* i wykorzystując w możliwie szerokim zakresie domyślne ustawienia parametrów budowanych modeli regresji, sieci neuronowych i drzew decyzyjnych, tak aby umożliwić czytelnikowi opanowanie zasad budowy i wykorzystywania tych modeli. W praktycznych zastosowaniach potrzebne jest zwykle stopniowe udoskonalanie modeli metodą prób i błędów aż do dobrania takich parametrów, które pozwolą zbudować jak najlep-

szy model uwzględniający jego dostosowanie do potrzeb analizy, ale też i biorąc pod uwagę fakt, że budując model jesteśmy uzależnieni od właściwości i jakości posiadanego zbioru danych, który służy do jego utworzenia. Brak odpowiednich danych może stanowić poważną przeszkodę w zbudowaniu modelu przydatnego dla praktycznych zastosowań.

Opisujemy sposób przeprowadzania oceny poprawności prognoz oraz generowania tzw. kodu skoringowego (*score code*), który służy do przewidywania wartości cechy nowych obiektów, niewystępujących w dotychczas badanym zbiorze.

Skoring (*scoring; score*) jest to ocena przypisywana jakiemuś obiektowi (podobnie jak stopnie w szkole, ocena punktowa kredytobiorcy w banku lub punkty w konkurencji sportowej), która pozwala porównywać obiekty ze względu na daną cechę (zmienną wynikową). Liczbowa wartość tej oceny jest obliczana na podstawie modelu uzyskanego z treningowego zbioru danych (*trained model*). Skoring (ocenie, punktacja) jest uogólnieniem pojęcia wartości teoretycznej zmiennej wynikowej i może być obliczany dla obiektów, dla których znamy wartości zmiennych objaśniających, ale nie znamy wartości zmiennej objaśnianej (*target*). Jest to więc przewidywana (prognozowana) wartość zmiennej objaśnianej dla tego obiektu. Skoring nowego zbioru danych (tzn. niebiorącego udziału w treningu) jest końcowym wynikiem większości problemów eksploracji danych. Na podstawie uzyskanej metody (wzoru) – tzw. *scoring formula* – obliczamy skoring (*ocenę skoringową*) nowego obiektu (np. respondenta ocenianego ze względu na prawdopodobieństwo pozytywnej odpowiedzi na ofertę marketingową).

W naszej pracy metodę zilustrujemy przykładem dotyczącym przewidywania spłaty kredytu bankowego przez klienta.

Rozdział 5. poświęcamy zastosowaniu metod eksploracji danych i programu *Enterprise Miner* w grupowaniu obiektów. Rozważamy zastosowanie metody grupowania należącej do niehierarchicznych metod tworzenia skupień oraz przedstawiamy możliwości tworzenia grup za pomocą sieci neuronowych Kohonena.

W rozdziale 6. opisujemy, jak przeprowadzić analizę asocjacji, jak wygenerować reguły przedstawiające zależności między zmiennymi, co umożliwi wyciąganie wniosków dotyczących współwystępowania określonych wartości zmiennych (reguły asocjacyjne) oraz następowania zjawisk w czasie (reguły sekwencji).

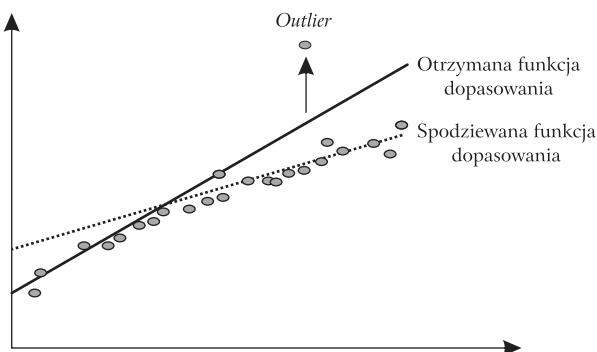
Rozdział 7. dotyczy grupowania zmiennych, które ma na celu ograniczenie redundancji informacji wprowadzanej przez zmienne, a także zmniejszenie wymiaru problemu, co ułatwia użytkowanie i interpretację modeli.

Metody eksploracji danych znajdują zastosowanie w prowadzeniu kampanii reklamowych, utrzymaniu klientów i zdobywaniu nowych klientów, analizach i ocenach kredytobiorców bankowych oraz ocenach ryzyka kredytowego, promocji produktów i usług, prognozowaniu i planowaniu sprzedaży, przewidywaniu sukcesów kampanii reklamowych oraz sprzedaży produktów, badaniach rynku.

Należy pamiętać, że proponowane metody wymagają dogłębnego ich zrozumienia przed zastosowaniem, aby móc prawidłowo odczytać wygenerowane przez nie wyniki. Wyniki te są na ogół silnie uzależnione od danych wejściowych. Narzuca to konieczność krytycznego podejścia do uzyskanych wyników przy ich interpretacji. Potrzebne jest też nabycie umiejętności doboru proponowanych w ich ramach algorytmów i wyznaczania parametrów wejściowych (początkowych) określających sposób realizacji algorytmów specyficznych dla poszczególnych proponowanych metod, np. przyjęcia niektórych wartości początkowych lub warunków zakończenia działania algorytmów z iteracjami. Trzeba pamiętać, że wnioski w przypadku metod eksploracji danych powinny być formułowane raczej w postaci domniemań niż kategoriycznych stwierdzeń.

Należy mieć także na uwadze, że metody eksploracji danych przed ich zastosowaniem wymagają wstępnej analizy danych, podczas której rozpatrywany jest charakter danych, np. występowanie wartości odstających (*outliers*) – por. rys. 0.4. O wynikach decyduje dobór cech, według których oceniane są obiekty. Aby uniknąć redundancji cech, wystarczająca okazuje się często analiza korelacji. Do niektórych analiz należy wybierać najslabiej skorelowane zmienne. W wielu przypadkach należy dążyć do ograniczenia liczby cech, gdy ich liczba jest zbyt duża z uwagi na potrzeby analizy, charakter zastosowanej metody lub ilość dostępnych danych (obiektów ze znaną charakterystyką).

Wymogiem użycia określonej metody eksploracji danych może być posiadanie danych dla odpowiednio dużej liczby obiektów w stosunku



RYSUNEK 0.4. Efekt wpływu wartości odstających (*outliers*)

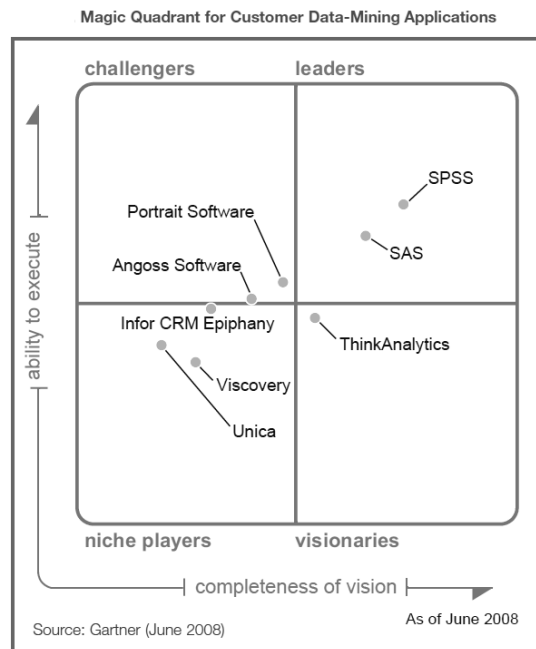
Źródło: na podstawie [A.P. Engelbrecht, 2002, s. 92].

do liczby cech. Możemy też po prostu uznać, że liczba cech opisujących dane dla obiektów jest zbyt duża (wręcz zastraszająco duża, np. wynosi kilkaset), co utrudni zbudowanie właściwego modelu predykcyjnego. Istnieją algorytmy doboru cech do analizy, które ograniczają liczbę zmiennych. Możliwe są dwie strategie postępowania przy stosowaniu tych algorytmów:

- (1) rozpoczyna się od małej liczby cech, po czym są dodawane kolejne cechy, aż do momentu uzyskania pożądanej jakości modelu (określane mianem algorytmów doboru cech w przód);
- (2) rozpoczyna się od zestawu wszystkich cech i kolejno eliminuje cechy, które mają najmniejszy wkład w uzyskanej jakości modelu, tj. z których rezygnacja w jak najmniejszym stopniu pogorszy jakość modelu (określane mianem algorytmów eliminacji wstecznej cech).

Kombinacją wymienionych powyżej strategii jest tzw. metoda krokowa, którą należy tu wymienić dla kompletności algorytmów i o której będzie jeszcze mowa w dalszej części tej pracy.

Zachętą do zapoznania się z programem *Enterprise Miner* firmy SAS może być wysoka pozycja, jaką zajmuje ta firma na rynku producentów narzędzi eksploracji danych. Przykładem potwierdzającym tę tezę są wyniki badań firmy *Gartner* prezentowane w Internecie w serii słynnych „magicznych kwadratów” (*magic quadrant*), co np. ilustruje rys. 0.5, ukazujący, że SAS znajduje się wśród liderów aplikacji do eksploracji danych wykorzystywanych do analizy danych o klientach.



RYSUNEK 0.5 Miejsce firmy SAS wśród dostawców narzędzi eksploracji danych

Źródło: wyniki badań Gartner, Inc., 2008 opublikowane w Internecie (*Gartner RAS Core Research Note G00158953*, podano za: <http://bi.pl/publications/art>, dostęp w dniu 18.04.2011 lub też http://www.spss.com.hk/PDFs/Gartner_Magic_Quadrant, dostęp w dniu 7.11.2011).

1

Zasady postępowania się programem *Enterprise Miner*

1.1. Rozpoczęcie pracy z programem i tworzenie projektów eksploracji danych

Program *SAS Enterprise Miner* (wersja 6.2) wywołuje się bezpośrednio z poziomu Windows, ale korzystanie z programu wymaga zainstalowania systemu *SAS 9.2*.



RYSUNEK 1.1. Okno programu *SAS Enterprise Miner* 6.2 po jego uruchomieniu

Źródło: opracowanie własne przy użyciu programu *SAS Enterprise Miner* 6.2.

Po uruchomieniu programu *SAS Enterprise Miner* w pojawiającym się oknie należy podać nazwę użytkownika oraz hasło (rys. 1.1). Można wybrać tryb pracy: lokalny (*Personal Workstation*) albo klient-serwer

(po odznaczeniu opcji *Personal Workstation* widoczne są dostępne serwery).

Aby kontynuować, trzeba utworzyć *projekty*, a w nich *diagramy*, za pomocą których będziemy opisywać proces przetwarzania. Diagramy składają się z *węzłów* określających sposób przetwarzania danych. Są one połączone strzałkami, które wskazują kolejność przetwarzania. Każdy węzeł służy do wywołania (uruchomienia) procedury obliczeniowej i podania informacji o parametrach realizowanej procedury. Należy utworzyć tzw. *biblioteki* i skojarzyć je z katalogami, w których znajdują się źródłowe zbiory danych.

Po zalogowaniu się do aplikacji (przycisk *Log On* – na rys. 1.1) ukazuje się okno *Enterprise Miner* (rys. 1.2). Można pracować, korzystając z już istniejącego projektu lub utworzyć nowy projekt.

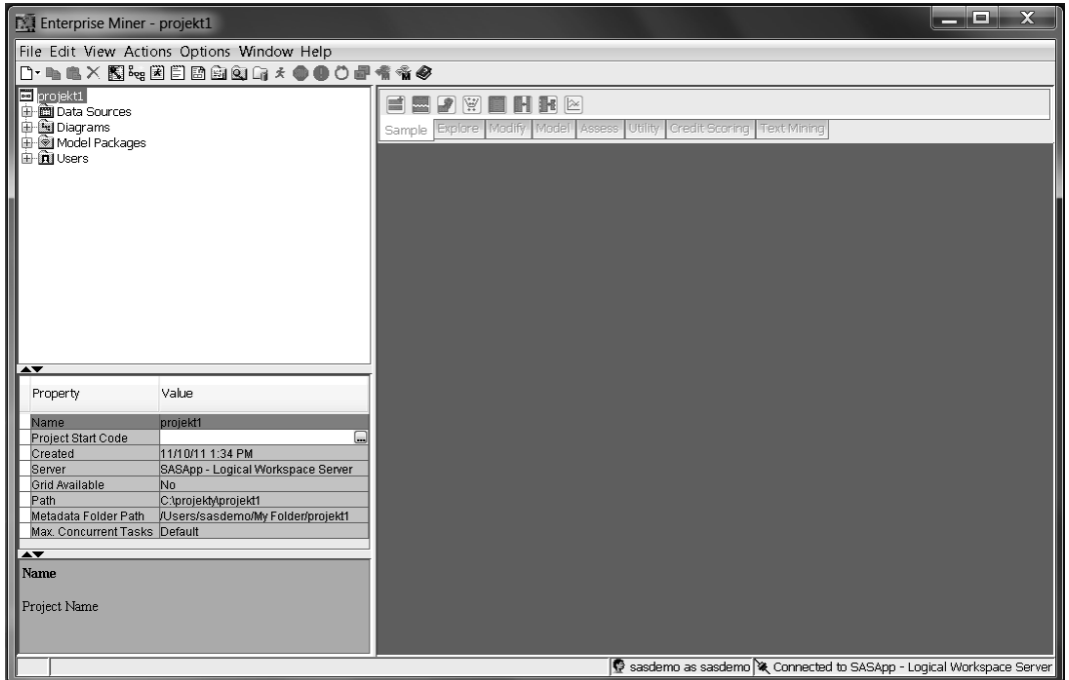


RYSUNEK 1.2. Okno powitalne w SAS Enterprise Miner

Źródło: opracowanie własne przy użyciu programu SAS Enterprise Miner 6.2.

Projekty budujemy w celu efektywnego wykorzystywania narzędzi eksploracji danych. W projektach przedstawiamy, z jakich zbiorów danych będziemy korzystać, jakie i w jakiej kolejności algorytmy będziemy stosować oraz jak będziemy przedstawiać i gdzie umieszczać wyniki przetwarzania.

Aby utworzyć nowy projekt, możemy wybrać opcję *New Project* (rys. 1.2) lub z głównego menu kolejno *File|New|Project* (lub skorzystać z odpowiedniej ikony). W pojawiającym się oknie wpisujemy nazwę nowego projektu i określamy katalog, w którym go zapiszemy. W podanym katalogu projektów jest tworzony podkatalog o podanej nazwie projektu, a w nim podkatalogi zawierające informacje o projekcie. Po zapisaniu wyświetlane jest okno projektu, wyglądające tak jak przed-



RYSUNEK 1.3. Okno projektu

Źródło: opracowanie własne przy użyciu programu SAS *Enterprise Miner* 6.2.

stawiono na rys. 1.3 (projektowi przy jego tworzeniu nadano nazwę *projekt1*).

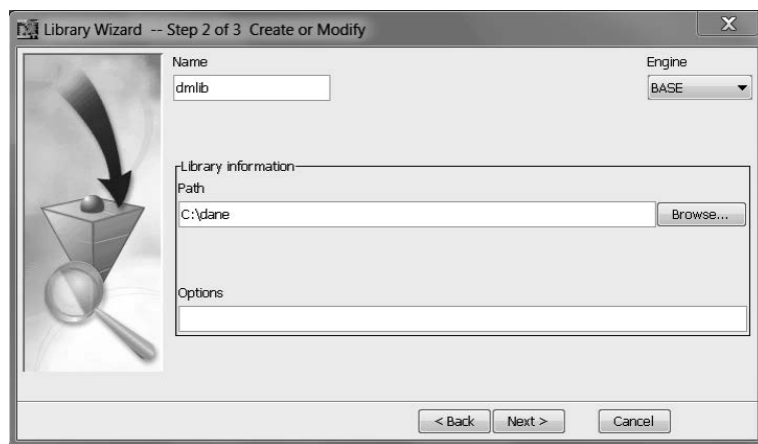
Aby korzystać z już uprzednio utworzonego projektu, wybieramy opcję *Open Project* (rys. 1.2) lub polecenia *File | Open Project* z głównego menu.

1.2. Organizacja zbiorów danych wykorzystywanych w analizach eksploracji danych

Zbiory danych przechowywane są w bibliotekach danych. Pełna specyfikacja zbioru danych w systemie SAS obejmuje nazwę zbioru poprzedzoną nazwą biblioteki i kropką oddzielającą obie nazwy, np. *Bibli1.dane*, oznacza zbiór danych o nazwie *dane* znajdujący się w bibliotece o nazwie *Bibli1*.

Dostęp do poleceń zarządzania bibliotekami i zbiorami danych możemy uzyskać po kliknięciu ikony *Explorer*. Jeżeli ikona nie jest widoczna na ekranie, to musimy wybrać z menu głównego polecenie *View | Explorer*, aby ją uaktywnić. Pojawia się okno *Explorer*. Aby utworzyć nową bibliotekę, można wówczas wykorzystać menu podręczne uruchamiane prawym klawiszem myszy w lewej części okna *Explorer*.

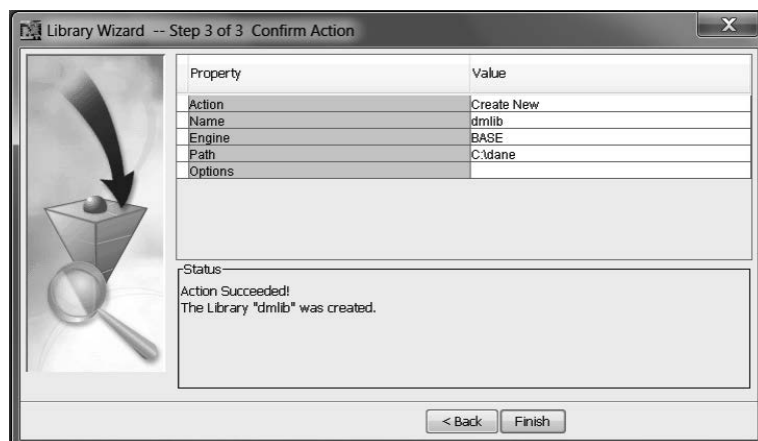
Jeżeli chcemy utworzyć nową bibliotekę, to możemy zamiast ikony *Explorer* wybrać z głównego menu polecenia *File | New | Library* (lub skorzystać z odpowiedniej ikony z grupy ikon znajdujących się poniżej głównego menu). Zostaje otwarty *Kreator Biblioteki (Library Wizard)*. Po kliknięciu przycisku *Next* otwiera się okno tworzenia biblioteki (rys. 1.4), w którym możemy wpisać nazwę biblioteki (przyjęliśmy nazwę *dmlib*) i wpisać lub wskazać (posługując się przyciskiem *Browse*) ścieżkę dostępu do katalogu zawierającego dane, które będziemy wykorzystywać (przyjmijmy, że nasze zbiory danych zostały umieszczone na dysku w katalogu *c:\dane*).



RYSUNEK 1.4. Tworzenie dostępu do zbioru danych na potrzeby eksploracji

Źródło: Opracowanie własne przy użyciu programu SAS Enterprise Miner 6.2.

Po przejściu do następnego okna (przycisk *Next*) potwierdzamy wybrane opcje przyciskiem *Finish*. Biblioteka zostaje utworzona i wszystkie zbiory SAS przypisane bibliotece stają się dostępne w programie *Enterprise Miner* (rys. 1.5).



RYSUNEK 1.5. Ostatnie okno kreatora; tu potwierdzamy utworzenie biblioteki ze zbiorami danych na potrzeby eksploracji

Źródło: opracowanie własne przy użyciu programu SAS Enterprise Miner 6.2.

	Numer gospodarstwa	Miesiąc badania	Województwo	Klasa miejscowości	Region	Liczba osób	Powierzchnia użytkowa mieszkania (gdy własność lub najem)	Telefon stacjonarny	Ocena sytuacji materialnej	Grupa społeczno-ekonomiczna (grs)
1	100260821	8 06		2	3	1	35	1	5	4
2	222970321	3 04		6	6	1	32	2	5	5
3	221800121	1 14		6	1	1	20	2	5	5
4	214070621	6 24		3	2	3	34	2	5	5
5	122340421	4 30		6	4	1	60	2	3	4
6	122780911	9 06		6	3	1	50	2	5	5
7	111190121	1 24		4	2	1	60	2	4	4
8	214411221	12 24		5	2	1	76	2	5	5
9	222390811	8 10		6	1	1	45	2	5	3
10	215090311	3 14		4	1	2	55	2	5	3
11	206730521	5 20		2	3	2	50	2	5	5
12	120110611	6 10		6	1	1	60	2	5	5
13	122340311	3 30		6	4	1	50	2	5	5
14	111370621	6 22		4	6	1	52	2	5	5
15	213600711	7 30		5	4	2	60	2	4	5
16	213990811	8 24		3	2	1	31	2	5	5
17	220161021	10 30		6	4	1	52	2	3	4
18	121430721	7 20		6	3	1	60	2	5	4
19	120190211	2 10		6	1	1	50	2	5	2
20	220330211	2 30		6	4	1	36	2	5	5
21	220561211	12 26		6	3	1	80	2	3	3

RYСУNEK 1.6. Ekran z widocznym fragmentem zawartości jednego z wykorzystywanych w pracy zbiorów (F2007BD)

Źródło: opracowanie własne przy użyciu programu SAS 9.2.

Zbiory danych są przedstawiane w postaci tablic. W kolejnych wierszach tablicy są opisane badane obiekty, a w kolumnach – analizowane zmienne. Aby obejrzeć zawartość zbioru, najprościej dwukrotnie kliknąć na jego nazwie (domyślnie zbiór jest otwierany za pomocą SAS *Enterprise Guide*). Można też kliknąć jego nazwę prawym klawiszem myszy i wybrać z podręcznego menu jedną z opcji *Open* (możemy obejrzeć zawartość zbioru za pomocą SAS *Enterprise Guide* lub SAS 9.2).

W naszej pracy będziemy analizować różne zbiory danych, z wielu dziedzin, aby możliwie ciekawie zilustrować przydatność eksploracji danych do różnych celów. Na rysunku 1.6 przedstawiamy fragment jednego z wykorzystywanych zbiorów.

Aby zbudować modele predykcyjne za pomocą metod eksploracji danych, zbiór danych jest zwykle dzielony na trzy zbiory zawierające:

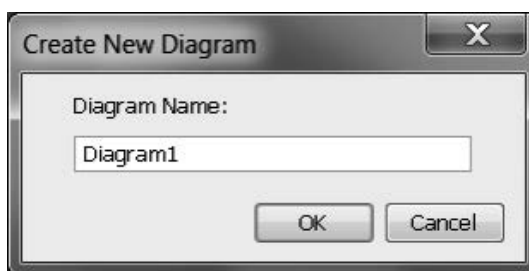
- (1) dane treningowe,
- (2) dane walidacyjne,
- (3) dane testowe.

Dane z tych zbiorów zawierają zmienną objaśnianą, której wartości staramy się przewidzieć, oraz wartości zmiennych objaśniających, na podstawie których będziemy opierać nasze przewidywania. Zbiór treningowy jest zbiorem danych, na podstawie których wykrywamy możliwe zależności między zmiennymi. Zbiór walidacyjny jest używany do poprawy jakości modelu, aby lepiej go dostosować bądź dopasować

do celów, dla jakich jest tworzony, to jest aby był bardziej poprawny i dokładny, czyli bardziej precyzyjny dla ich realizacji¹. Zbiór testowy jest zbiorem, który służy nam do zbadania, na ile wykryte zależności są prawdziwe dla innych danych niż użyte do „ich wykrycia”. Nim dokonamy podziału zbioru i zbudujemy modele do analizy danych za pomocą metod udostępnianych przez *Enterprise Miner*, musimy utworzyć diagram wyznaczający przebieg analiz.

1.3. Tworzenie diagramów i zarządzanie diagramami analizy danych

Nowy diagram możemy utworzyć, wybierając z głównego menu polecenia *File | New | Diagram* lub klikając odpowiednią ikonę, lub wybierając z menu podręcznego w oknie projektu prawym klawiszem myszy polecenie *Diagrams*. Diagramowi możemy nadać własną nazwę. Wpisujemy nazwę diagramu (tu – rys. 1.7: *Diagram1*) i klikamy przycisk *OK*.



RYSUNEK 1.7. Tworzenie diagramu eksploracji danych

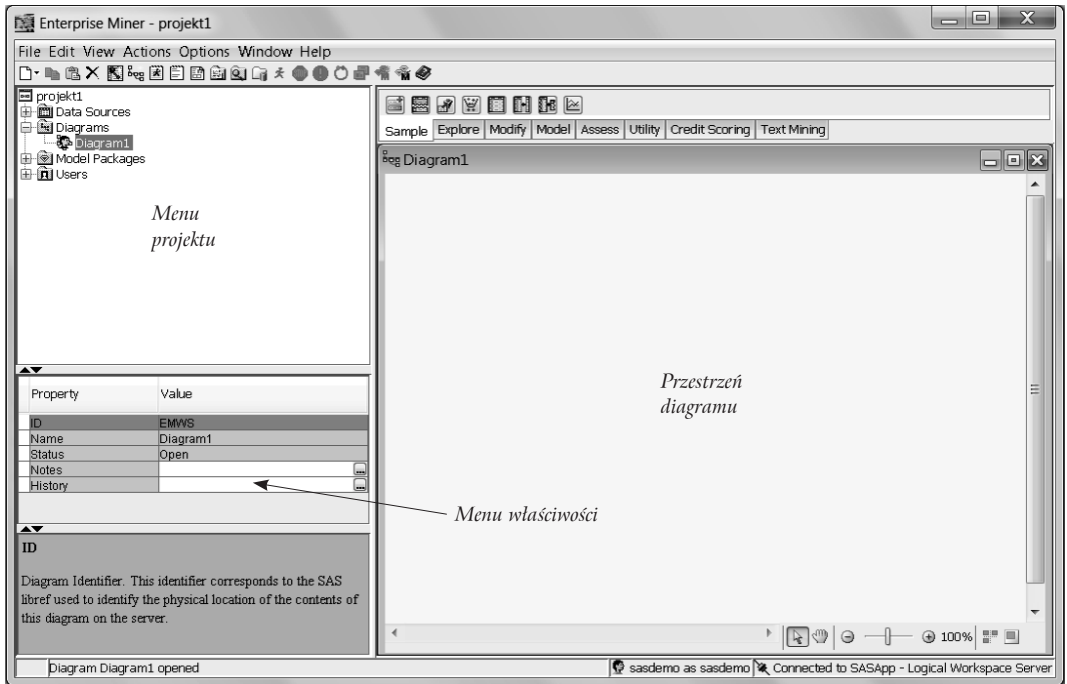
Źródło: opracowanie własne przy użyciu programu SAS *Enterprise Miner 6.2*.

W projekcie powstaje diagram o podanej nazwie zawierający puste pole robocze, w którym będziemy tworzyć schemat analizy danych (rys. 1.8 – por. *Przestrzeń diagramu*). Powyżej tego pola widoczne są zakładki ikon węzłów, z których będziemy budować diagram.

Menu projektu (rys. 1.8) umożliwia śledzenie struktury projektu. Służy też do dodawania nowych komponentów do projektu (np. źródeł danych) i nawigacji pomiędzy diagramami. W *Menu właściwości* (rys. 1.8) określa się parametry wykorzystywanych metod, a w *Przestrzeni diagramu* (rys. 1.8) – buduje procesy eksploracji danych, wstawiając węzły diagramu realizujące poszczególne procedury i łącząc je, aby wskazać kierunek przepływu przetwarzania.

Czynności przetwarzania określamy, przeciągając na diagram odpowiednie ikony węzłów. Aby wskazać kierunek analizy, łączymy wpro-

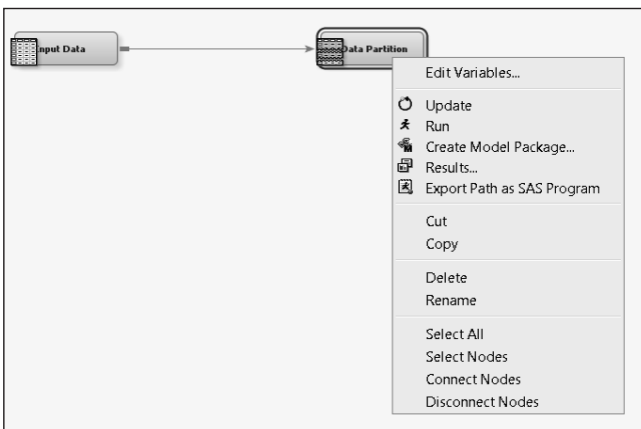
¹ W *Enterprise Miner* statystyki pozwalające ocenić jakość modelu są domyślnie obliczane na podstawie danych ze zbioru walidacyjnego (także już podczas budowy modelu „z danych treningowych”).



RYСУNEK 1.8. Okno *Enterprise Miner*; po prawej obszar *Diagram1*, w którym powstanie schemat eksploracji danych

Źródło: opracowanie własne przy użyciu programu SAS *Enterprise Miner* 6.2.

wadzone węzły, przeciągając linię ze strzałką skierowaną w kierunku przetwarzania, od węzła do węzła diagramu – jak ilustruje to rys. 1.9, na którym są widoczne dwa węzły: *Input Data* i *Data Partition* połączone strzałką.



RYСУNEK 1.9. Przestrzeń diagramu z wprowadzonymi dwoma przykładowymi węzłami określającymi zakres analizy danych i strzałką wskazującą kierunek przetwarzania; widoczne także przykładowe menu podręczne „obsługi” węzła uzyskiwane po kliknięciu na wybranym węźle prawym przyciskiem myszy

Źródło: opracowanie własne przy użyciu programu SAS *Enterprise Miner* 6.2.

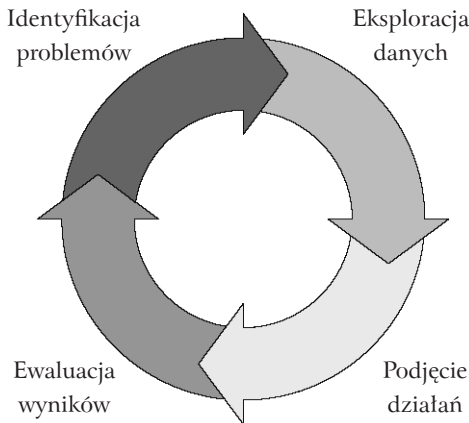
Zagadnienia metodyczne przeprowadzania eksploracji danych

2.1. Metodyki wykorzystywane na potrzeby analizy danych

W literaturze są opisywane i stosowane w praktyce różne metodyki wspomagające organizację eksploracji danych. Do najbardziej znanych należą:

- *Virtuous Cycle of Data Mining*,
- *CRISP-DM*,
- *Six Sigma*,
- *SEMMA*.

W przypadku *Virtuous Cycle of Data Mining* postuluje się przeprowadzanie badania danych pod względem stosowanej metodyki w sposób zbliżony do sterowania *procesami biznesowymi*, gdzie przez *proces biznesowy* jest rozumiana realizacja ciągu powiązanych ze sobą etapów postępowania prowadzących do rozwiązania określonego problemu lub do osiągnięcia określonego efektu. Proces *Virtuous Cycle of Data Mining*



RYSUNEK 2.1. Metodyka postępowania zgodnie z *Virtuous Cycle of Data Mining*

Źródło: opracowanie własne na podstawie [M.J.A. Berry, G.S. Linoff, 2000; M.J.A. Berry, G.S. Linoff, 2004].

składa się z czterech etapów: identyfikacji problemów, eksploracji danych, podjęcia działań, oceny (ewaluacji) wyników (rys. 2.1).

Etapy metodyki *Virtuous Cycle of Data Mining* realizuje się w sposób powtarzalny (iteracyjnie). Obejmują one:

- (1) identyfikację problemów:
 - analizę działalności podmiotu i znalezienie tych jej aspektów, które potencjalnie mogą zostać zoptymalizowane,
 - identyfikację działalności, które mogą wpłynąć na dostępność danych i możliwość podejmowania działań,
 - ocenę wiarygodności stosowanych źródeł danych oraz możliwości pozyskiwania danych,
 - zebranie wiedzy o problemie wynikającej z doświadczenia i intuicji praktyków;
- (2) eksplorację danych:
 - identyfikację i pozyskanie danych,
 - sprawdzenie i czyszczenie danych,
 - przekształcenie danych oraz uzyskanie właściwego układu danych,
 - wybranie próby uczącej,
 - wybranie metody modelowania,
 - ocenę jakości modelu;
- (3) podjęcie zamierzonych działań na podstawie wyników uzyskanych w kroku eksploracji danych;
- (4) ewaluację (ocenę) wyników oraz dokonanie zmian i ulepszeń w budowanym modelu.

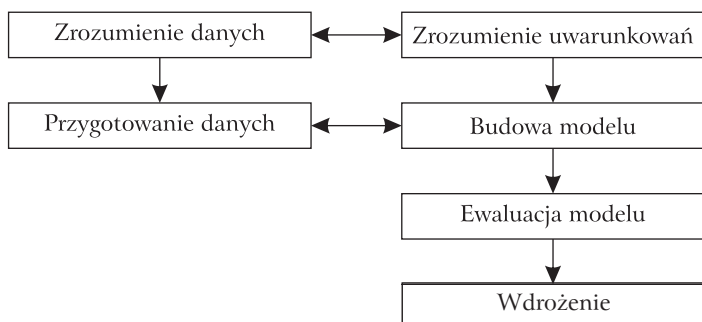
Kolejna z wymienionych metodyk *CRISP-DM* (*Cross-Industry Standard Process for Data Mining*) została opracowana w 1996 r. przez analityków z *DaimlerChrysler*, *SPSS* (*Statistical Package for the Social Science*) i *NCR*.

CRISP-DM (rys. 2.2) proponuje zastosowanie standardowego procesu dopasowania eksploracji danych do ogólnej strategii rozwiązywania problemów komórki biznesowej lub badawczej. Składa się z sześciu faz, które obejmują realizację wymienionych poniżej czynności:

- (1) zrozumienie uwarunkowań biznesowych:
 - sformułowanie celów i wymagań projektu zgodnie z ich rozumieniem w komórce biznesowej lub badawczej, której dotyczą,
 - wykorzystanie sformułowanych celów i wymagań do opracowania definicji problemu eksploracji danych,
 - stworzenie wstępnego planu działań zmierzających do osiągnięcia sformułowanych celów;
- (2) zrozumienie danych:
 - zebranie danych,

- wstępną analizę danych i ocenę jakości danych,
 - (jeżeli trzeba), wybranie interesujących podzbiorów, które mogą zawierać wzorce;
- (3) przygotowanie danych:
- przygotowanie ze wstępnych danych ostatecznego zbioru danych, który będzie wykorzystywany we wszystkich następnych fazach,
 - wybór obserwacji i zmiennych, które będą analizowane,
 - wykonanie przekształceń zmiennych (jeżeli jest to konieczne),
 - czyszczenie danych;
- (4) budowę modelu:
- wybór i zastosowanie odpowiednich technik modelowania,
 - skalowanie parametrów modelu,
 - jeżeli trzeba, należy wrócić do etapu przygotowania danych, by przybrały one postać odpowiadającą specyficznym wymaganiom danej techniki eksploracji danych;
- (5) ewaluacja modelu:
- ocenę modelu (lub kilku modeli) pod względem jakości i efektywności, przed jego wdrożeniem,
 - ustalenie, czy model rzeczywiście spełnia wszystkie założenia ustalone w pierwszej fazie,
 - ocenę, czy są jakieś ważne cele biznesowe lub badawcze, które nie zostały w należyty sposób uwzględnione,
 - podjęcie decyzji co do wykorzystania wyników eksploracji;
- (6) wdrożenie:
- wykorzystanie stworzonego modelu (względnie modeli),
 - wykorzystywanie wniosków z poprzednio zrealizowanych projektów w kolejnych, nowych projektach.

Kolejna faza i realizowane czynności zależą od wyników poprzednich faz i czynności. Niekiedy jesteśmy zmuszeni do powrotu do faz wcześniejszych, zanim przejdziemy do kolejnej fazy (rys. 2.2).

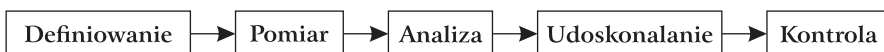


RYСУNEK 2.2. Metodyka CRISP-DM

Źródło: opracowanie własne na podstawie [D.T. Larose, 2006, s. 4–9].

Następna z wymienionych na wstępie rozdziału metodyk przeprowadzania analizy danych to metodyka *Six-Sigma*. Metodyka ta wypracowana w koncernie *Motorola* w połowie lat osiemdziesiątych XX w. stanowi zbiór „dobrych praktyk” zarządzania danymi dla doskonalenia procesów stosowanych na potrzeby sterowania jakością. Nazwa pochodzi od sześciokrotnego odchylenia standardowego w rozkładzie normalnym. Sześć Sigma oznacza w tym przypadku 3,4 wadliwych produktów na 1 milion wykonanych i jest ogólnym wskazaniem co do poziomu jakości w całej działalności firmy lub instytucji. Celem wdrożenia programu *Six-Sigma* jest zmniejszenie liczby defektów do 3,4 defekta na milion egzemplarzy produktu [M. Lasek, M. Pęczkowski, 2005].

Założono, że wdrażanie metodyki *Six-Sigma* będzie odbywać się zgodnie z procesem zarządzania danymi w toku realizacji pięciu etapów na potrzeby sterowania jakością. Jest on w skrócie nazywany *DMAIC*, od pierwszych liter słów *Define Measure Analyze Improve Control* (rys. 2.3).



RYSUNEK 2.3. Realizacja procesu *Six-Sigma*

Źródło: opracowanie własne na podstawie [Ch.B. Tayntor, 2003].

Definiowanie obejmuje określenie celów związanych z procesami, które mają być udoskonalone z jednoczesnym uwzględnieniem wymagań klienta oraz strategii przedsiębiorstwa. Pomiar jest rozumiany jako pomiar wyników działania aktualnie istniejących rozwiązań oraz gromadzenia danych w celu dokonania porównań w przyszłości. Analiza dotyczy znajdowania związków pomiędzy różnymi czynnikami wpływającymi na dany proces. Udoskonalenie to wprowadzanie udoskonaleń i poprawek eliminujących wcześniej wykryte problemy. Kontrola oznacza zaś monitorowanie wyników działania zastosowanych udoskonaleń. Postępowanie może zostać przejęte dla potrzeb eksploracji danych. Między innymi na podstawie postępowania założonego w metodyce *Six-Sigma* została zaproponowana przez *SAS Institute* metodyka eksploracji danych o nazwie *SEMMA* (rys. 2.4).

Metodyka *SEMMA* składa się z pięciu etapów: *Sample Explore Modify Model Assess* i od pierwszych liter wyrazów oznaczających kolejne etapy utworzono jej nazwę *SEMMA* [M. Lasek, M. Pęczkowski, 2010 (b)].



RYSUNEK 2.4. Kolejne etapy metodyki *SEMMA*

Źródło: opracowanie własne na podstawie [Reference, 2009].