

---

# **Hurtownie danych**

## **– czyli jak zapewnić dostęp do wiedzy tkwiącej w danych**





**Rodzaj zajęć:** Wszechnica Popołudniowa

**Tytuł:** Hurtownie danych – czyli jak zapewnić dostęp do wiedzy tkwiącej w danych

**Autor:** mgr inż. Andrzej Ptasznik

**Redaktor merytoryczny:** prof. dr hab. Maciej M Sysło

Zeszyt dydaktyczny opracowany w ramach projektu edukacyjnego **Informatyka+** – ponadregionalny program rozwijania kompetencji uczniów szkół ponadgimnazjalnych w zakresie technologii informacyjno-komunikacyjnych (ICT).

**[www.informatykaplus.edu.pl](http://www.informatykaplus.edu.pl)**

**[kontakt@informatykaplus.edu.pl](mailto:kontakt@informatykaplus.edu.pl)**

**Wydawca:** Warszawska Wyższa Szkoła Informatyki

ul. Lewartowskiego 17, 00-169 Warszawa

**[www.wysi.edu.pl](http://www.wysi.edu.pl)**

**[rektorat@wysi.edu.pl](mailto:rektorat@wysi.edu.pl)**

Projekt graficzny: FRYCZ I WICHA

Warszawa 2010

Copyright © Warszawska Wyższa Szkoła Informatyki 2010

Publikacja nie jest przeznaczona do sprzedaży.



**KAPITAŁ LUDZKI**  
NARODOWA STRATEGIA SPÓJNOŚCI



WARSZAWSKA  
WYŻSZA SZKOŁA  
INFORMATYKI

**UNIA EUROPEJSKA**  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



Projekt współfinansowany ze środków Unii Europejskiej w ramach Europejskiego Funduszu Społecznego.

---

# **Hurtownie danych – czyli jak zapewnić dostęp do wiedzy tkwiącej w danych**



**Andrzej Ptasznik**

Warszawska Wyższa Szkoła Informatyki

aptaszni@wwsi.edu.pl

---

**Streszczenie**

Przedmiotem wykładu są podstawy teorii hurtowni danych i aspekty ich wykorzystania. W pierwszej części zostaną omówione podstawowe cechy systemów OLTP (ang. *On-Line Transaction Processing*) oraz systemów OLAP (ang. *On-Line Analytical Processing*). Omówione zostaną podstawowe pojęcia i przykłady projektów hurtowni danych. Przedstawione zostaną podstawowe zagadnienia związane z integracją danych. Omówione zostanie pojęcie analitycznej kostki wielowymiarowej. Zaprezentowane zostaną elementy technologii usług analitycznych i ich znaczenie w systemach typu *Business Intelligence*. W części końcowej wykładu omówione zostaną krótko podstawowe pojęcia związane z eksploracją danych (ang. *Data Mining*).

**Spis treści**

1. Wprowadzenie .....	5
2. Systemu OLTP i OLAP .....	6
3. Podstawy hurtowni danych .....	7
4. Problemy integracji danych .....	10
5. Kostka wielowymiarowa .....	11
6. Systemy Business Intelligence .....	13
7. Eksploracja danych .....	14
8. Podsumowanie.....	15
Literatura .....	16



## 1 WPROWADZENIE

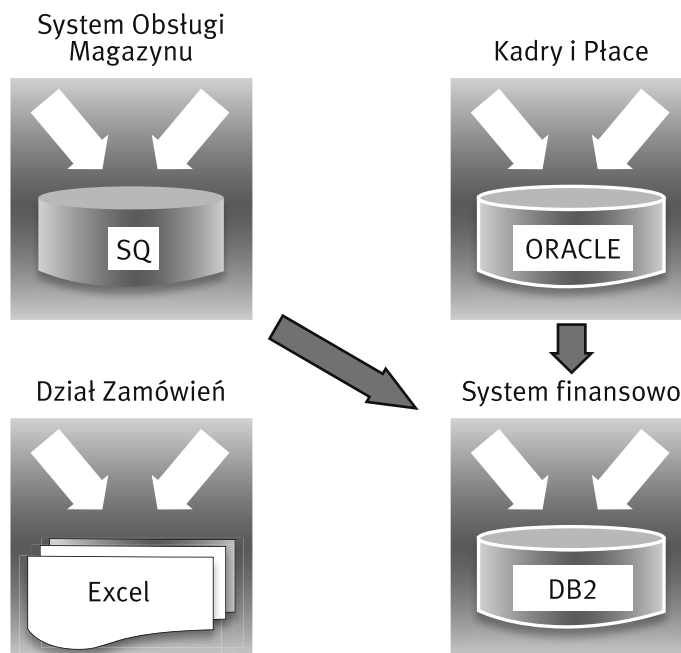
Burzliwy rozwój technologii informatycznych, a w szczególności baz danych, spowodował, że w każdej firmie czy instytucji gromadzone są różne dane na różnych etapach działalności. Bardzo często dane są gromadzone w różnych technologiach – od plików tekstowych poprzez arkusze kalkulacyjne do baz danych. W okresie początkowego rozwoju systemy informatyczne wspomagające działalność firm koncentrowały się na wsparciu działalności operacyjnej. Powstawały różne systemy ukierunkowane na konkretny aspekt działania, przykładowo:

- wystawianie faktur,
- obsługa magazynu,
- systemy kadrowe,
- systemy księgowo,
- obsługa klientów.

Bardzo często systemy takie nie były z sobą w żaden sposób powiązane i wykonywane były przez różnych producentów w różnych technologiach. Stosowanie technologii informatycznych w codziennej działalności firm i instytucji było związane z gromadzeniem danych na potrzeby konkretnego typu działania. Dane gromadzone w różnych systemach, oprócz wspomagania codziennych działań, były wykorzystywane także do celów raportowania i informowania kierownictwa. Podstawowymi problemami takiej działalności były:

- dane po pewnym czasie stawały się niepotrzebne, ponieważ obsługa działalności codziennej nie musiała korzystać z danych historycznych (w systemie obsługi magazynu istotny był aktualny stan towaru w magazynie a nie jaki był ten stan w zeszłym roku) – często w tego typu systemach usuwano starsze dane;
- przetrzymywano bardzo często te same dane w różnych formatach;
- przetwarzanie danych na potrzeby inne niż wsparcie działalności codziennej znacząco wpływało na wydajność tych systemów;

Na rysunku 1 przedstawiony został schemat organizacji instytucji z wykorzystaniem różnych systemów informatycznych.



Rysunek 1.

Przykładowa organizacja firmy z wykorzystaniem różnych systemów informatycznych

Duże ilości gromadzonych danych stają się kopalnią wiedzy, która może zostać wykorzystana do właściwego kierowania firmą i osiągnięcia przewagi konkurencyjnej na rynku.



## 2 SYSTEMY OLTP I OLAP

Tradycyjne systemy baz danych ukierunkowane są na realizację wielu małych i prostych zapytań i mają zapewnić wsparcie dla realizacji codziennych działań pracowników danej firmy lub instytucji. Dla tego typu systemów Edgar Frank „Ted” Codd (brytyjski informatyk, znany przede wszystkim ze swojego wkładu do rozwoju teorii relacyjnych baz danych) wprowadził pojęcie **systemów OLTP** (ang. *On-Line Transaction Processing*) i zdefiniował zbiór zasad, które powinny spełniać systemy tego typu. Podstawowe cechy systemów typu OLTP to:

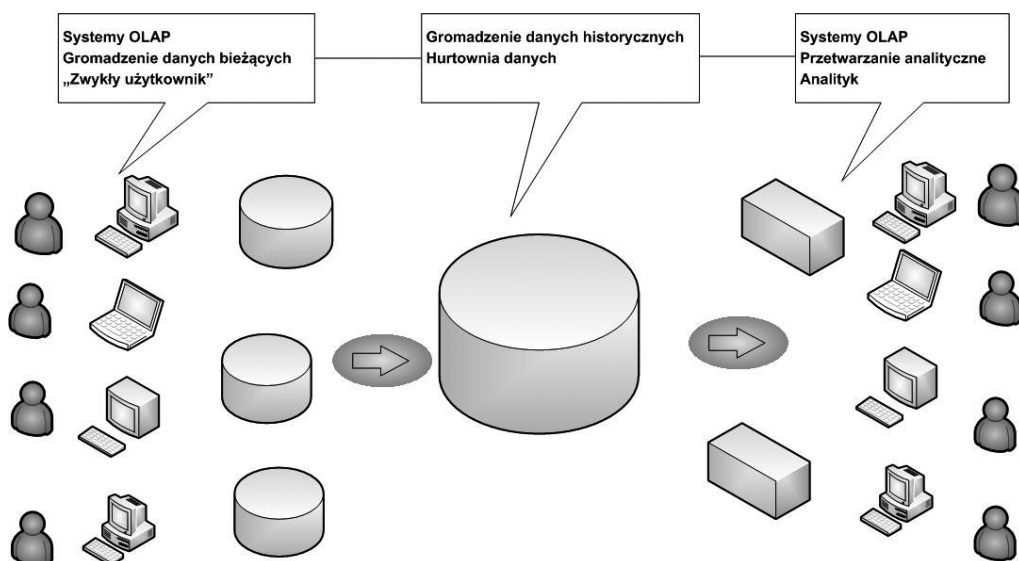
- przechowywane dane są zorientowane procesowo, np. wystawione faktury, otrzymane zamówienia, złożone reklamacje, wykonane przelewy itp.;
- stosunkowo niewielkie rozmiary baz danych (kilka gigabajtów);
- przechowywane są dane bieżące bez konieczności gromadzenia danych historycznych;
- realizowana jest duża ilość w miarę prostych zapytań;
- przechowywane są dane elementarne;
- realizowane są operacje wstawiania, modyfikowania i usuwania danych.

Zbiory danych tworzone w systemach OLTP stają się przydatne do pozyskiwania dodatkowych informacji potrzebnych kierownictwu firmy do podejmowania decyzji. Pojawiają się tu jednak pewne problemy:

- w ramach jednej firmy może istnieć wiele systemów typu OLTP,
- realizowanie dodatkowych czynności w ramach systemu OLTP wpływa na jego wydajność, tym bardziej dlatego, że pozyskiwanie danych analitycznych wymaga wykonywania złożonych zapytań operujących na dużej ilości danych,
- klasyczne zapytania SQL dostarczają danych w postaci dwuwymiarowych tabel, co często jest niewystarczające dla tego typu zastosowań.

Rozwiązaniem tych problemów stała się koncepcja wydzielonych systemów informatycznych świadczących usługi analityczne. Wspomniany wyżej Edgar Codd nazwał systemy tego typu **OLAP** (ang. *On-Line Analytical Processing*) i również dla systemów tego typu sformułował zbiór zasad, które systemy tego typu powinny spełniać. Podstawowe cechy systemów OLAP to:

- przechowywane dane są zorientowane tematycznie, np. sprzedaż produktów, stany zapasów, wydatki, akcje promocyjne itp.;
- bardzo duże ilości gromadzonych danych (rzędu wielu terabajtów);
- przechowywane są dane bieżące i historyczne;
- realizowana są bardzo złożone zapytania operujące na wielkich ilościach danych;
- przechowywane są dane elementarne i zagregowane (sumy, średnie itp.);
- wykonywane są głównie operacje dopisywania nowych danych – praktycznie nie wykonuje się operacji modyfikowania danych.



Rysunek 2.  
Schemat architektury powiązania systemów OLTP i OLAP

Elementem łączącym systemy OLTP i OLAP są wyspecjalizowane bazy danych, gromadzące w specjalnie zaprojektowanych strukturach dane historyczne zwane **hurtowniami danych**. Na rys. 2 przedstawiono schemat architektury systemów OLTP i OLAP z hurtownią danych. Pokazuje on w sposób symboliczny ideę centralnej zbiornicy danych łączącej systemy OLTP i systemy OLAP.

### 3 PODSTAWY HURTOWNI DANYCH

Potrzeba analizy danych dotyczących bieżącej i przyszłej działalności organizacji była podstawowym impulsem do powstania nowych systemów informatycznych. Analiza taka stanowi podstawę do podejmowania decyzji dotyczących zarządzania przedsiębiorstwem i wspomagania podejmowania decyzji. Istniejące dotychczas systemy informatyczne (głównie klasy OLTP) nie mogą dostarczyć potrzebnych danych, gdyż są oparte na operacyjnych bazach danych realizujących codzienne procesy, mogą być rozproszone (dane znajdują się w wielu różnych źródłach), niejednorodne a często nie są z sobą powiązane. Struktury danych są dostosowane do działań operacyjnych, dane są poddawane operacjom modyfikacji. W operacyjnych bazach danych przechowuje się dane odzwierciedlające jedynie aktualny stan lub najnowszą historię, tymczasem do analiz i porównań potrzebne są długookresowe dane historyczne. Rozwiązaniem tego problemu okazała się **hurtownia danych** (ang. *Data Warehouse*). Hurtownia danych jest wydzieloną centralną bazą danych zbierającą informacje służące do zarządzania organizacją. Jest ona odizolowana od baz operacyjnych a jej struktura i użyte do jej budowy narzędzia powinny być zoptymalizowane pod kątem przetwarzania analitycznego. Prosta, najczęściej cytowaną, definicję pojęcia hurtowni danych zaproponował W.H. Inmon (jeden z czołowych teoretyków hurtowni danych i systemów OLAP – autor książki *Building the Data Warehouse*).

**Hurtownia danych** to zbiór zintegrowanych, nieulotnych, ukierunkowanych baz danych, wykorzystywanych w systemach wspomagania decyzji.

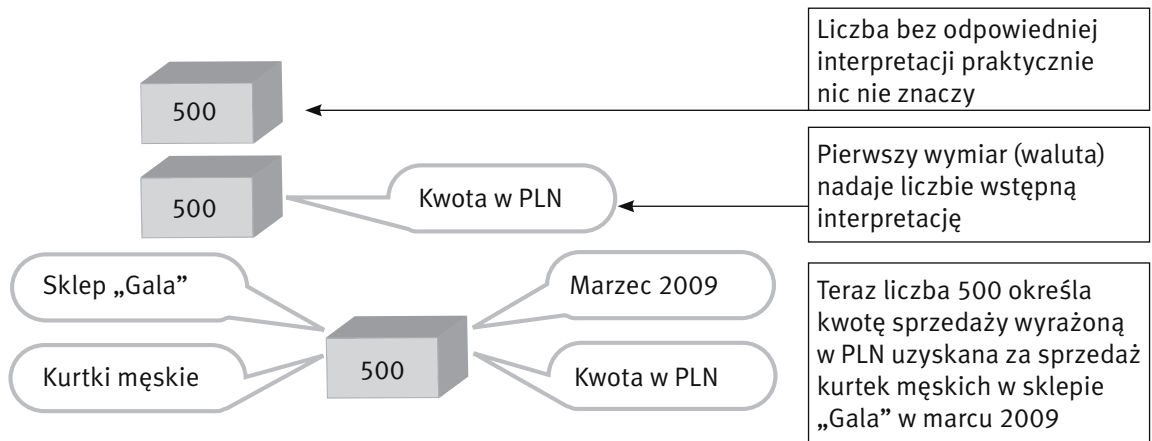
Podstawowe cechy hurtowni danych to:

- **Jest scentralizowaną bazą danych** – gromadzi dane z różnych źródeł i przechowuje je w specjalnie zaprojektowanych strukturach.
- **Jest oddzielona od baz operacyjnych** – tym samym operacje wykonywane na danych gromadzonych w hurtowniach nie wpływają na wydajność operacji realizowanych w systemach OLTP.
- **Scala informacje z wielu źródeł** – ponieważ dane dotyczące jednego procesu mogą być w konkretnej firmie tworzone i przechowywane w różnych bazach danych lub nawet w plikach czy arkuszach kalkulacyjnych.
- **Jest zorientowana tematycznie** – gromadzi dane opisujące różne aspekty działalności firmy.
- **Przechowuje dane historyczne** – hurtownie mają niezaspokojony „apetyt” na dane, im dłuższa historia przechowywanych danych tym większe możliwości analizy.
- **Utrzymuje wielką ilość informacji** – w hurtowniach danych praktycznie nie wykonuje się operacji usuwania danych, czyli ilość danych tylko rośnie wraz z dostarczaniem nowych porcji danych.
- **Agreguje informacje** – z punktu widzenia analizy najczęściej interesują nas podsumowania, obliczenia średnich i inne działania matematyczne wykonywane na grupach danych.

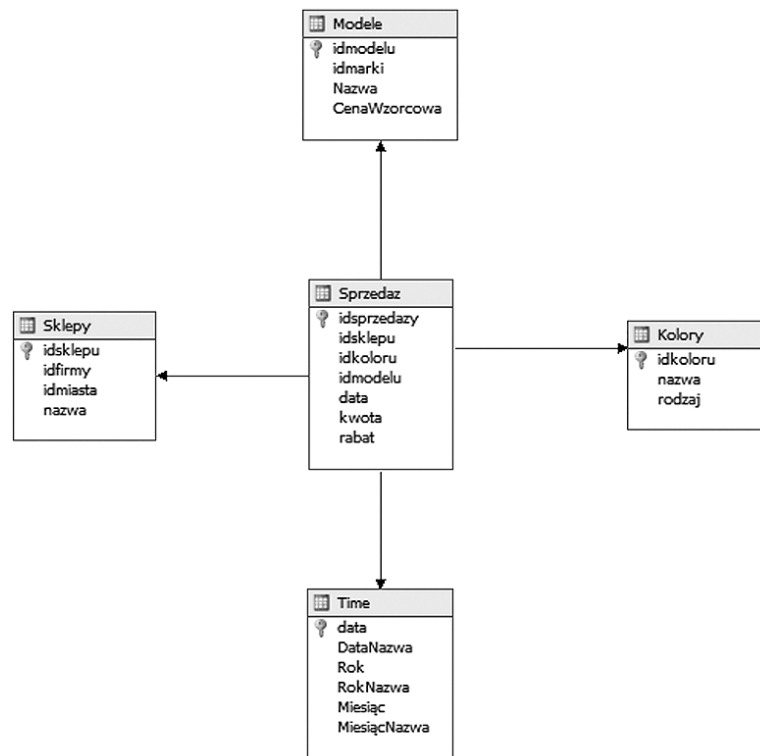
Najczęściej hurtownie danych są tworzone jako bazy relacyjne, w których są projektowane tabele faktów i tabele wymiarów. **Fakt** to pojedyncze zdarzenie będące podstawą analiz (np. sprzedaż produktów, udzielone kredyty itp.). Fakty są opisane przez wymiary i miary. **Miara** to wartość liczbową dowiązana do danego faktu, np. kwota sprzedaży, ilość sztuk, a **wymiar** to cecha opisująca dany fakt, np. data, klient, produkt, lokalizacja. Dodatkowo, wymiary zawierają atrybuty, które są dodatkowymi cechami wymiaru, np. dla wymiaru czas atrybutami mogą być miesiąc, kwartał i rok. Istotę pojęć miar i wymiarów omówimy na przykładzie. Podstawowymi elementami gromadzonymi w hurtowniach są wartości liczbowe, czyli miary pewnych faktów.

Jak pokazano na rys. 3, wymiary są cechami opisującymi wartość miar, czyli nadają wartościom liczbowym odpowiedni sens. Najczęściej stosowanym wzorcem przy projektowaniu hurtowni jest tak zwany **schemat gwiazdy**. Na rysunku 4 przedstawiono przykładowy projekt hurtowni danych opisujący sprzedaż samochodów.





Rysunek 3.  
Interpretacja miary



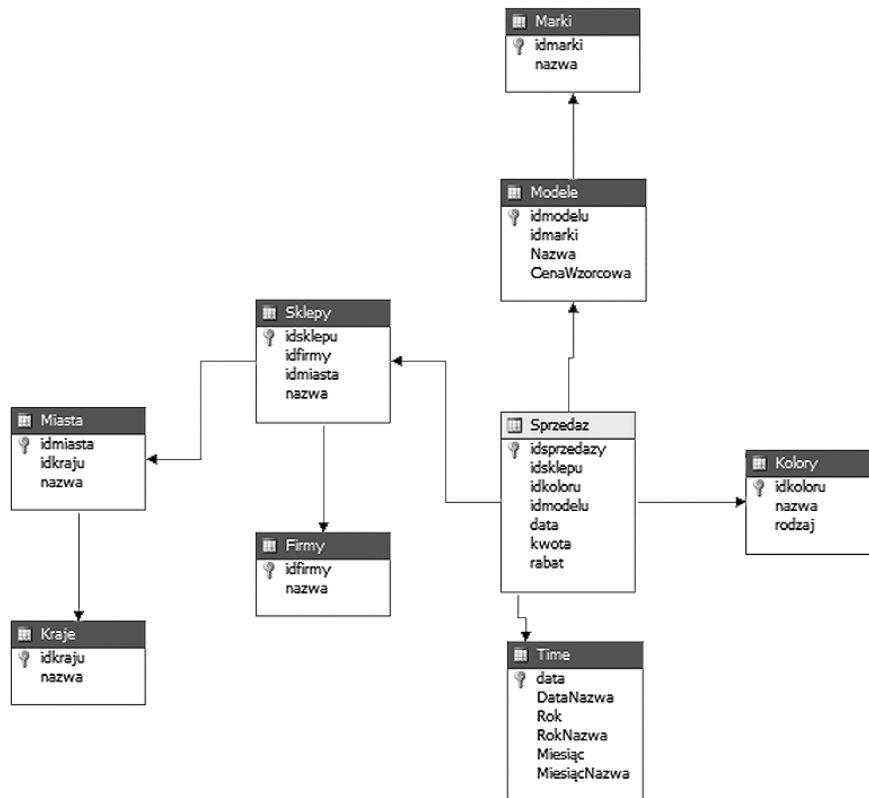
Rysunek 4.  
Przykładowy projekt hurtowni danych w schemacie gwiazdy

Centralną tabelą jest tabela o nazwie Sprzedaz, w której są zapisywane fakty opisujące kwoty uzyskane za sprzedaż samochodów. Tabela faktów łączy się z czterema tabelami opisującymi różne wymiary (kolor, model, sklep i czas). Połączenia tabel wymiarów z tabelą faktów są realizowane za pomocą odpowiednich kluczy obcych.

Do podstawowych cech schematu gwiazdy należy zaliczyć:

- prosta struktura, dzięki czemu schemat jest łatwy do zrozumienia;
- duża efektywność zapytań ze względu na niewielką liczbę połączeń tabel;
- dominująca struktura dla hurtowni danych, wspierana przez wiele narzędzi.

Rozwinięciem schematu gwiazdy jest schemat **płatka śniegu**, który występuje wtedy, gdy wymiary są powiązane z innymi tabelami. Na rysunku 5 przedstawiono przykładowy projekt hurtowni w schemacie płatka śniegu, który jest rozszerzeniem projektu z rysunku 4.



Rysunek 5. Przykładowy projekt hurtowni danych w schemacie płątka śniegu

Do podstawowych cech schematu płątka śniegu należy zaliczyć:

- spadek wydajności zapytań w porównaniu ze schematem gwiazdy ze względu na większą liczbę połączeń tabel;
- struktura łatwiejsza do modyfikacji;
- wykorzystywany rzadziej niż schemat gwiazdy, gdyż efektywność zapytań jest ważniejsza niż efektywność ładowania danych do tabel wymiarów;

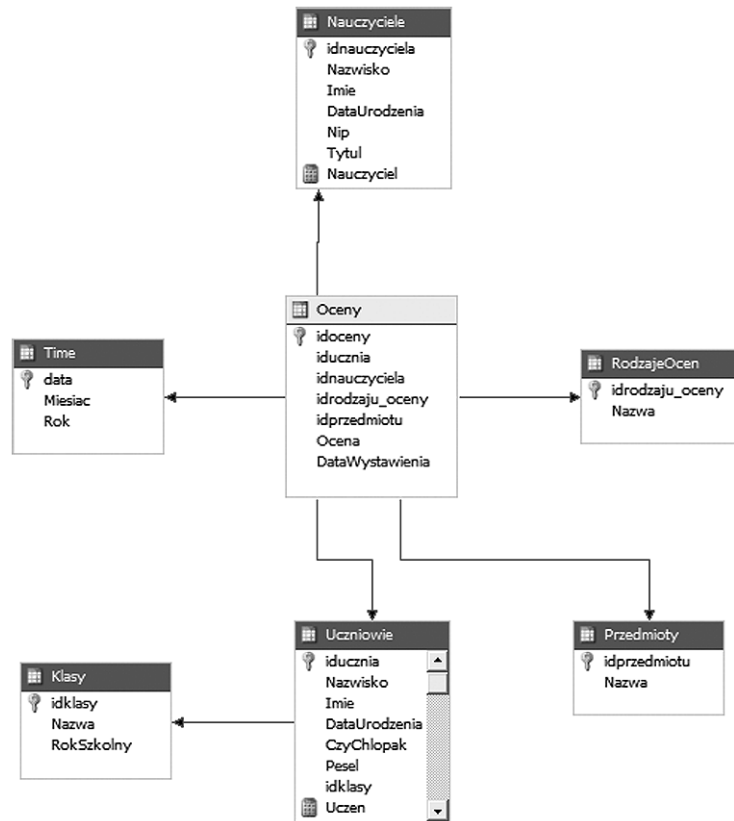
Hurtownie danych stanowią podstawowe źródło zasilające procesy analizy danych. Przedstawione przykłady projektów hurtowni są jedynie wycinkiem, gdyż w rzeczywistości hurtownie składają się z wielu podobnych struktur danych, opisujących różne fakty i korzystające z różnych wymiarów. Na rys. 6 został przedstawiony jeszcze jeden przykład projektu struktury hurtowni danych, w którym faktami są oceny wystawione uczniom. Każda ocena jest charakteryzowana przez :

- datę jej wystawienia – wymiar Time;
- ucznia, który otrzymał ocenę – wymiar Uczniowie, który jest dodatkowo opisywany przez wymiar Klasy;
- nauczyciela, który ocenę wystawił – wymiar Nauczyciele;
- przedmiot, z którego ocena została wystawiona – wymiar Przedmioty;
- rodzaj wystawionej oceny – wymiar RodzajeOcen.

Tworzenie hurtowni danych dla jednej szkoły wydaje się niecelowe ze względu na stosunkowo niewielką ilość danych, ale można sobie wyobrazić istnienie takiej hurtowni w skali kraju i wtedy stanowiłaby podstawę do analizy skuteczności nauczania.

Nie jesteśmy w stanie, w ramach tego wykładu, omówić wszystkich aspektów tworzenia hurtowni danych, gdyż są to zagadnienia złożone i praktycznie każdy projekt ma swoją specyfikę i może wyglądać zupełnie inaczej w zależności od swojego przeznaczenia i założeń jakie dana firma przyjęła przy realizacji. Przedstawione zasady stanowią punkt wyjścia przy realizowaniu konkretnego projektu.

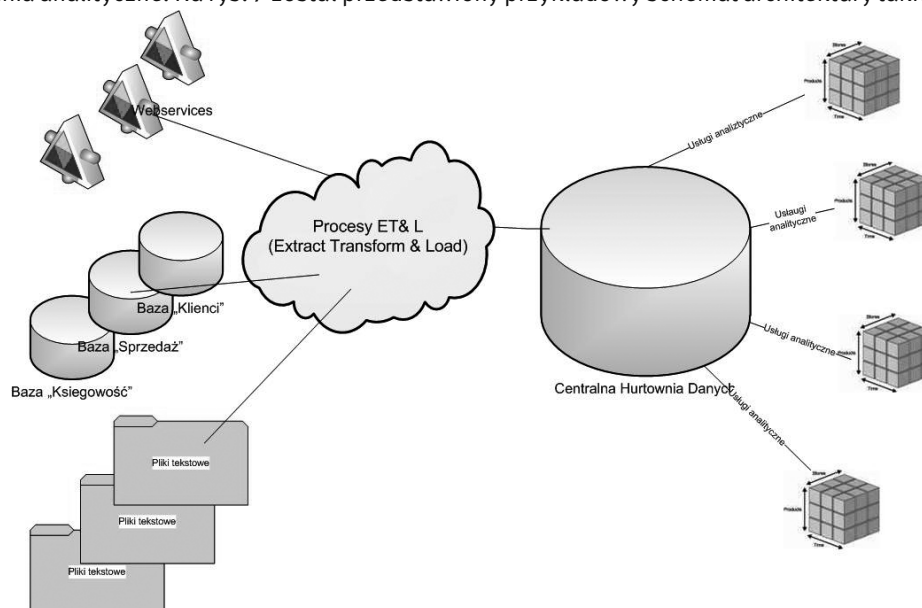




Rysunek 6. Przykładowy projekt hurtowni dla wystawianych ocen w szkołach

#### 4 PROBLEMY INTEGRACJI DANYCH

Hurtownie danych są zasilane danymi pobieranymi z systemów OLTP, które mogą być wykonane w różnych technologiach oraz innych źródeł danych dostępnych w konkretnej firmie. Na bazie hurtowni są realizowane różne zadania analityczne. Na rys. 7 został przedstawiony przykładowy schemat architektury takiego systemu.

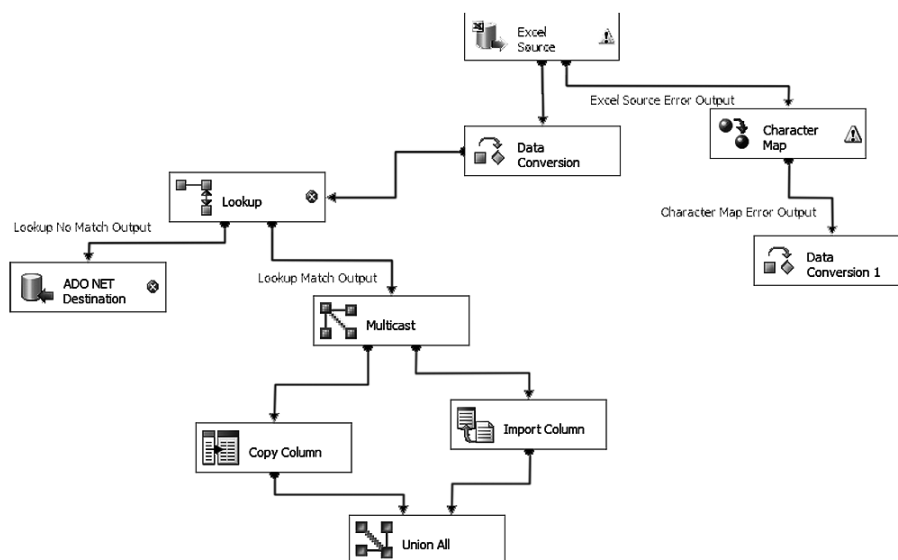


Rysunek 7. Architektura otoczenia hurtowni danych

Przedstawiony na rys. 7 schemat pokazuje warstwę (nazwaną Procesy ET&L), która występuje pomiędzy systemami OLTP i innymi źródłami danych a hurtownia danych. Problemy związane z pozyskiwaniem danych dla hurtowni są jednymi z najtrudniejszych zadań przy jej tworzeniu. W ramach warstwy ET&L (ang. *Extract Transform & Load* – pobierz, przekształć i zapisz) są realizowane następujące zadania:

- standaryzacja danych – ponieważ dane pobierane mogą być z wielu różnego typu źródeł, to należy doprowadzić je do jednakowej postaci;
- konwersja typów danych – różne systemy mogą w inny sposób zapisywać dane i dlatego należy je doprowadzić do tego samego typu;
- transformacje danych – dane w systemach roboczych mogą być przechowywane w innej postaci niż postać ich zaprojektowana w hurtowni, dlatego należy je odpowiednio przekształcić;
- agregacja danych – w hurtowniach nie musimy zapisywać każdej elementarnej danej z systemów operacyjnych a jedynie pewne zbiorcze wartości;
- integracja danych z różnych źródeł – dane tego samego rodzaju z punktu widzenia hurtowni (np. opis klienta) mogą być zapisywane w różnych źródłach danych i przed zapisaniem w hurtowni należy je odpowiedni powiązać,
- czyszczenie danych i kontrola poprawności – ponieważ w systemach operacyjnych mogą być przechowywane dane błędne, dlatego przed zapisaniem w hurtowni należy je sprawdzić i usunąć dane błędne;
- dodatkowe przekształcenia, np. przeliczenie wartości różnych walut.

Zadania warstwy ET&L są wspierane przez różne technologie, w ramach których projektuje się i programuje działanie odpowiednich procesów. Na rysunku 8 został przedstawiony przykładowy fragment schematu procesu ET&L wykonany w MS SQL Server 2008 Integration Services.



Rysunek 8.

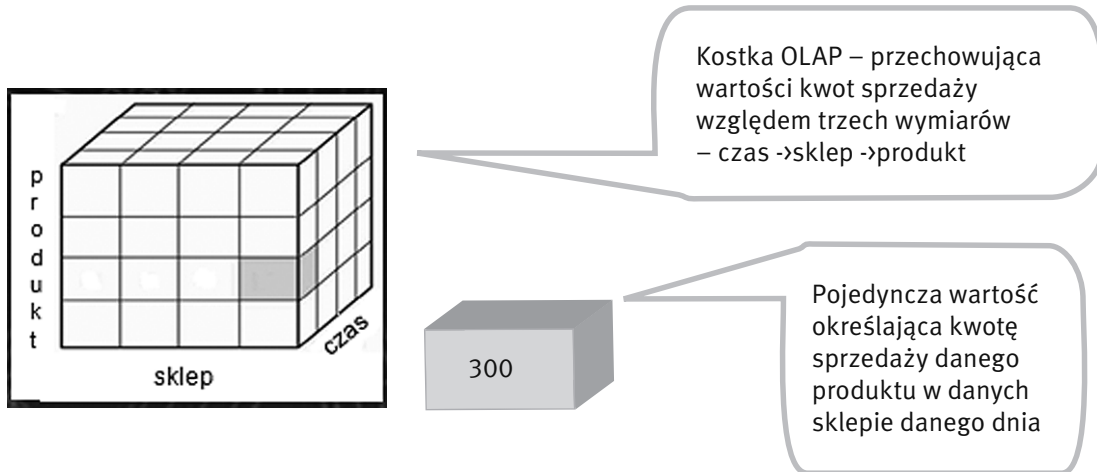
Przykładowy pakiet usługi MS SQL Server 2008 Integration Services

Technologia MS SQL Server 2008 Integration Services umożliwia definiowanie złożonych procesów pozyskiwania, przekształcania i zapisywania danych z różnych źródeł. Projektowany schemat przetwarzania prezentowany jest za pomocą ikon opisujących różne etapy i zadania procesu.

## 5 KOSTKA WIELOWYMIAROWA

Hurtownie danych stanowią punkt wyjścia do realizacji usług analitycznych. Najczęściej stosowanym elementem usług analitycznych jest **wielowymiarowa kostka OLAP**, która przechowuje dane w sposób bardziej przypominający wielowymiarowe arkusze kalkulacyjne niż tradycyjną, relacyjną bazę danych.

Kostka umożliwia wyświetlanie i oglądanie danych z różnych punktów widzenia. Do jej budowy potrzeba dowolnego źródła danych opartego na tabelach relacyjnych – czyli najczęściej kostki wielowymiarowe buduje się w oparciu o hurtownie danych. Kostka składa się z miar, wymiarów oraz poziomów i jest zoptymalizowana pod kątem szybkiego i bezpiecznego dostępu do danych wielowymiarowych. **Miary** to wskaźniki numeryczne (ile?), natomiast **wymiary** reprezentują dane opisowe (kto? co? kiedy? gdzie? jak?). Wymiary są pogrupowane za pomocą **poziomów**, które odzwierciedlają hierarchię i umożliwiają użytkownikom zwiększanie lub zmniejszanie poziomu szczegółowości analizowanego wymiaru. Jak widać, kostka OLAP oparta jest o te same pojęcia (miary i wymiary) co schematy hurtowni danych. Trudno graficznie zaprezentować strukturę wielowymiarową – dlatego najczęściej kostka jest pokazywana w postaci sześcianu, czyli kostki złożonej z trzech wymiarów.



Rysunek 9.  
Kostka OLAP

Podczas analizy z wykorzystaniem kostek wielowymiarowych, dane są poddawane typowym operacjom, do których zaliczamy m.in.:

- **zwijanie** – podnoszenie poziomu agregacji czyli, uogólnianie danych;
- **rozwijanie** – zmniejszanie poziomu agregacji, dane stają się bardziej szczegółowe;
- **selekcja** – wybór interesujących elementów wymiarów;
- **projekcja** – zmniejszanie liczby wymiarów.

Obsługę tworzenia i eksploatacji kostek wielowymiarowych wspierają różne technologie, między innymi MS SQL Server 2008 Analysis Services. Na rysunku 10 zostało przedstawione przykładowe zestawienie na bazie kostki OLAP, opisującej sprzedaż samochodów. Zestawienie pokazuje wartość sprzedaży poszczególnych marek samochodów w kolejnych latach.

		Rok Nazwa   Miesiąc Nazwa								
		Rok 2000	Rok 2001	Rok 2002	Rok 2003	Rok 2004	Rok 2005	Rok 2006	Rok 2007	Rok 2008
Marka	Model	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota
Fiat		87160932	83147859	88051496	84921857	86334394	83859853	80762763	84165778	87452700
Ford		108975594	106558654	110131873	101560521	104506206	103569152	102243478	106147122	105850450
Honda		103068833	100467684	106067324	103753240	104024481	98014402	97112774	100667562	98499940
Kia		55557236	53813716	56930702	50686638	53099688	50090456	51065556	55321341	56766143
Opel		119280487	125054961	117615346	112393332	120499921	112738729	113487766	122608813	123282674
Grand Total		474043082	469042874	478796741	453315588	468464690	448272592	444672337	468910616	471851907

Rysunek 10.  
Przykładowe zestawienie zbiorcze na bazie kostki OLAP

Kolejne zestawienie na rys. 11 pokazuje elementy uszczegółowienia, polegające na rozbiściu kwot rocznych na poszczególne miesiące oraz rozbiściu kwot sprzedaży marki Fiat na poszczególne modele.

		Rok Nazwa   Miesiąc Nazwa											
		Rok 2000											
		April 2000	August 2000	December 2000	February 2000	January 2000	July 2000	June 2000	May 2000	April 2000	March 2000	February 2000	January 2000
Marka	Model	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota	Kwota
Fiat	Brava	2866953	1874901	1988477	2200225	2190336	2170497	2515317	2315				
	Fiat 500	1582491	1743247	1425258	1006753	1623506	1484610	1602509	1280				
	Grande Punto	2174695	2018129	1627548	1189492	2090553	1950285	2320150	1887				
	Panda	1629036	1263054	1560842	1123210	1761087	1308147	1804285	1175				
	Total	8253175	6899331	6602125	5519680	7665482	6913539	8242261	6660				
Ford		9498183	9551630	8166704	6978720	8626562	10254700	10365812	8557				
Honda		8352079	7671737	9034879	7110533	9081817	8114449	8690739	7634				
Kia		4408343	5342258	5416620	3208120	5170567	3824752	5032367	4368				
Opel		9558239	11524474	10282450	7072156	11777633	8844828	11948671	8134				
Grand Total		40070019	40989430	39502778	29889209	42322061	37952268	44279850	3535				

Rysunek 11.

Przykładowe zestawienie zbiorcze na bazie kostki OLAP z elementami uszczegółowienia

Do obsługi i pozyskiwania danych z kostek wielowymiarowych istnieje specjalny język MDX (ang. **Multi Dimensional eXpressions** – wyrażenia wielowymiarowe) – opis tego języka wykracza zdecydowanie poza ramy naszego wykładu. Wielowymiarowe kostki OLAP są przechowywane w specjalizowanych strukturach zoptymalizowanych pod kątem szybkości pozyskiwania danych.

## 6 SYSTEMY BUSINESS INTELLIGENCE

**Business Intelligence (BI)** – **analitka biznesowa** – jest pojęciem bardzo szerokim. Do dzisiaj nie istnieje powszechnie przyjmowana definicja systemów tej klasy. Najbardziej ogólnie można przedstawić je jako proces przekształcania danych w informacje, a informacji w wiedzę, która może być wykorzystana do zwiększenia konkurencyjności przedsiębiorstwa. Systemy BI są mocno uzależnione od utworzenia hurtowni danych, które umożliwiają ujednoczenie i powiązanie danych zgromadzonych z różnorodnych systemów informatycznych przedsiębiorstwa. Utworzenie hurtowni danych zwalnia systemy transakcyjne od tworzenia raportów i umożliwia równoczesne korzystanie z różnych systemów BI. System BI opierają się na następującej koncepcji:

- system BI generuje standardowe raporty lub wylicza kluczowe wskaźniki efektywności działania przedsiębiorstwa (ang. *Key Performance Indicators*);
- na podstawie standardowych raportów i wskaźników stawia się hipotezy;
- postawione hipotezy weryfikuje się poprzez wykonywanie szczegółowych analiz danych z wykorzystaniem różnego rodzaju narzędzi analitycznych (np. OLAP, *data mining*).

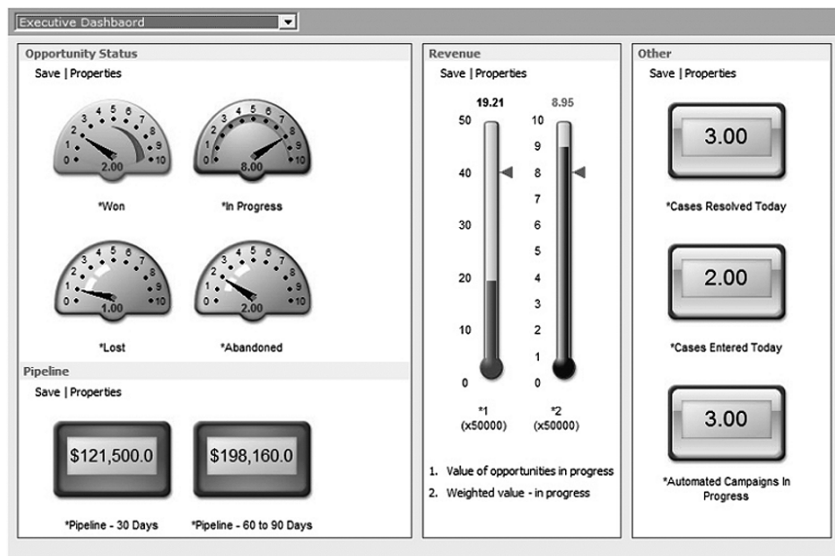
Najczęściej spotykane odmiany systemów zaliczanych do BI to:

- EIS – systemy powiadamiania kierownictwa (ang. *Executive Information Systems*).
- DSS – systemy wspomagania decyzji (ang. *Decision Support Systems*).
- MIS – systemy wspomagania zarządzania (ang. *Management Information Systems*).
- GIS – systemy informacji geograficznej (ang. *Geographic Information Systems*).

Systemy BI są narzędziem dla menedżerów i specjalistów zajmujących się analizami i strategią. Dla menedżerów niższych szczebli, którzy oczekują informacji o aktualnym stanie procesów, są przeznaczone rozwiązania **Business Activity Monitoring (BAM)**, umożliwiające przetwarzanie napływających na bieżąco danych. Techniki prezentacyjne są dobierane odpowiednio do potrzeb użytkownika. Jednym ze sposobów prezentowania wyników wstępnej analizy i sygnalizowania przekroczenia założonych wartości w działalności firmy jest koncepcja **kokpitu menadżera**. Idea kokpitu jest oparta na założeniu, aby bardzo szybko informować menadżera o wartościach podstawowych wskaźników oraz sygnalizować niekorzystne zjawiska zachodzące w jego dziedzinie odpowiedzialności. Do graficznej prezentacji takich faktów są używane proste gadżety (wskaźniki, sy-



gnalizatory świetlne, liczniki). Elementy kokpitu powinny dać ogólny obraz procesów zachodzących w firmie. Na rysunku 12 został pokazany przykładowy kokpit menadżera.



Rysunek 12. Przykładowa postać kokpitu menadżera

Jeżeli z obrazu wskaźników kokpitu wynika problem, to należy uruchomić inne, przeważnie bardziej złożone procesy analizy.

## 7 EKSPLOACJA DANYCH

**Eksploracja danych** (spotyka się również określenie drążenie danych, pozyskiwanie wiedzy, wydobywanie danych, ekstrakcja danych) (ang. *data mining*) – jest jednym z etapów procesu, który bywa nazywany **odkrywaniem wiedzy z baz danych** (ang. *Knowledge Discovery in Databases, KDD*). Idea eksploracji danych jest oparta na wykorzystaniu komputerów i ogromnych zbiorów danych do znajdowania ukrytych dla człowieka prawidłowości w danych zgromadzonych w hurtowniach danych. Istnieje wiele technik eksploracji danych, które są oparte na zaawansowanej statystyce (statystyczna analiza wielowymiarowa) oraz technikach i metodach wywodzących się z obszaru badań nad sztuczną inteligencją. Główne przykłady stosowanych rozwiązań to:

- wizualizacje na wykresach,
- metody statystyczne,
- sieci neuronowe,
- metody uczenia maszynowego,
- metody ewolucyjne,
- logika rozmyta,
- zbiory przybliżone.

Motywację dla rozpatrywania tego typu narzędzi stanowi ciągły wzrost technicznych możliwości gromadzenia i analizy danych, w których ukryte są potencjalnie cenne informacje dopełniające wiedzę. Zastosowanie technik KDD daje szczególnie dobre wyniki w nowych dziedzinach, gdzie tak zwana wiedza ekspercka jest jeszcze w dużej mierze niepełna i nieugruntowana. Do takich dziedzin można przykładowo zaliczyć:

- Analizę różnych aspektów ruchu internetowego.
- Marketing z wykorzystaniem Internetu.
- Rozpoznawanie obrazu, pisma, mowy itd.
- Wspomaganie diagnostyki medycznej.

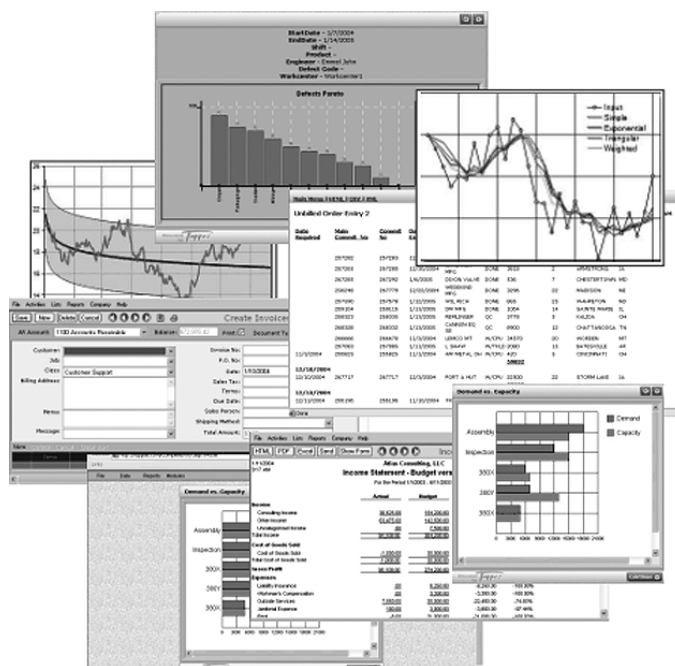
- Badania genetyczne.
- Analizę historii operacji bankowych i zapobieganie wyłudzeniom.
- Optymalizację działań związanych z systemami CRM zajmujących się zarządzaniem relacjami z klientami.

Proces odkrywania wiedzy z danych przebiega według poniższego schematu:

- **Zrozumienie dziedziny problemu** – złożoność danych a także problemów stawianych przy okazji ich analizy, coraz częściej nie umożliwia natychmiastowe sformułowanie pytań, na które użytkownik chce uzyskać odpowiedź. Trzeba dobrze zrozumieć problem, dla którego chcemy stosować techniki KDD.
- **Budowa roboczego zbioru danych** – określenie, z jakich zasobów danych będziemy korzystać w procesie KDD.
- **Oczyszczenie, przekształcanie i redukcja danych** – istotę tego problemu omówiliśmy w rozdziale poświęconym integracji danych.
- **Eksploracja danych (data mining)** – realizacja procesu odkrywania wiedzy przy użyciu bardzo różnorodnych technik, opartych na statystyce, sztucznej inteligencji, czy też odwołujących się do metod uczenia maszynowego.

Podstawowym problemem procesów odkrywania wiedzy tkwiącej w danych jest to, że różnych regularności jest w danych praktycznie „nieskończenie wiele”, zaś dla użytkownika interesujące będą tylko niektóre z nich i to w różnym stopniu. Osiągnięcie dobrych wyników w procesie eksploracji danych jest uzależnione nie tylko od danych i wykorzystywanych technologii ale przede wszystkim od wiedzy i zaangażowania analityków wykonujących te zadania. Przykładowe postaci zobrazowań wyników, które można uzyskiwać w procesie eksploracji danych przedstawiono na rys 13.

Techniki i metody eksploracji danych są w stadium ciągłego rozwoju i należy się spodziewać nowych rozwiązań w tym zakresie.



Rysunek 13.

Przykładowa postać kokpitu menadżera. (źródło [www.shopfloorreporting.com](http://www.shopfloorreporting.com)).

## 8 PODSUMOWANIE

Hurtownie danych są wydzielonymi, specjalizowanymi bazami danych, przeznaczonymi do wspomaganie usług analitycznych. Wdrożenie hurtowni danych może dostarczyć firmie wiele korzyści:

- **Odciążenie systemów transakcyjnych** – przygotowanie analiz i zestawień nie obciąża już systemów transakcyjnych, które mogą obsługiwać bieżące operacje. Zasilenie hurtowni danymi z systemów

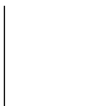
źródłowych wykonywane jest automatycznie i najczęściej odbywa się w cyklu dziennym, z reguły w nocy, gdy użytkownicy nie korzystają z systemu.

- **Poprawa jakości analizowanych danych** – analizując dane w hurtowni danych na zagregowanym poziomie dużo łatwiej wychwycić pewne nieprawidłowości w systemach źródłowych. W hurtowni danych bardzo dobrze widać np., czy koszty są przypisane do odpowiednich nośników, czy wszyscy klienci są przypisani do regionów sprzedaży lub handlowców itd.
- **Przechowywanie danych o długim horyzoncie czasowym** – dzięki temu, że w hurtowni danych mamy łatwy dostęp do danych wieloletnich możemy wykonywać bardziej trafne prognozy, czy też doszukiwać się określonych trendów.
- **Łączenie danych pochodzących z różnych systemów transakcyjnych** – hurtownia danych może pobrać dane z praktycznie każdego źródła danych. Dane te są następnie porządkowywane i dokonywana jest unifikacja pojęć i mierników. Dzięki temu możliwe staje się porównanie niejednorodnych danych.
- **Udostępnienie danych dla wszystkich potrzebujących** – w hurtowni danych możemy zdefiniować poszczególnym użytkownikom uprawnienia do odpowiedniego wycinka danych. Przy pomocy narzędzi analitycznych i wizualizacji danych, użytkownicy mogą wykonywać na ich bazie różne zestawienia, raporty i analizy.

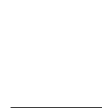
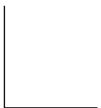
## LITERATURA

1. Hand D., Mannila H., Smyth P., *Eksploracja danych*, WNT, Warszawa 2002
2. Jarke M., Lenzerini M., Vassiliou Z., Vassiliadis P., *Hurtownie danych. Podstawa organizacji funkcjonowania*, WSiP, Warszawa 2003
3. Poe V., Klauer P., Brobst S., *Tworzenie hurtowni danych*, WNT, Warszawa 2000
4. Surma J., *Business Intelligence. Systemy wspomaganie decyzji biznesowych*, WN PWN, Warszawa 2009
5. Todman Ch., *Projektowanie hurtowni danych*, WNT, Warszawa 2003











W projekcie **Informatyka +**, poza wykładami i warsztatami, przewidziano następujące działania:

- 24-godzinne kursy dla uczniów w ramach modułów tematycznych
- 24-godzinne kursy metodyczne dla nauczycieli, przygotowujące do pracy z uczniem zdolnym
  - nagrania 60 wykładów informatycznych, prowadzonych przez wybitnych specjalistów i nauczycieli akademickich
    - konkursy dla uczniów, trzy w ciągu roku
    - udział uczniów w pracach kół naukowych
    - udział uczniów w konferencjach naukowych
      - obozy wypoczynkowo-naukowe.

Szczegółowe informacje znajdują się na stronie projektu

**[www.informatykaplus.edu.pl](http://www.informatykaplus.edu.pl)**