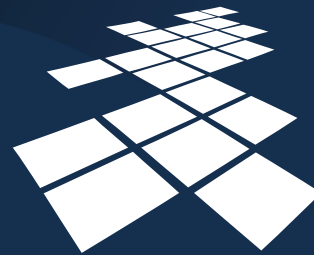


Zaawansowane systemy baz danych - ZSBD

Standard SQL/MM

Wykład prowadzi:
Marek Wojciechowski



UCZELNIA
ONLINE

Standard SQL/MM

Wykład poświęcony jest standardowi SQL/MM, który jest nowym standardem uzupełniającym język SQL o biblioteki do obsługi specjalistycznych danych i aplikacji.

Celem wykładu jest przedstawienie głównych idei standardu SQL/MM, a w szczególności jego części poświęconych przetwarzaniu danych tekstowych, przestrzennych i obrazów w bazach danych.

Do zrozumienia treści wykładu niezbędna jest znajomość systemów baz danych i języka SQL oraz podstawowej problematyki multimedialnych baz danych.



Plan wykładu

- Wprowadzenie do standardu SQL/MM
- Przegląd części standardu SQL/MM
- Omówienie specyfikacji SQL/MM Full Text
- Omówienie specyfikacji SQL/MM Spatial
- Omówienie specyfikacji SQL/MM Still Image
- Podsumowanie

Wykład rozpocznie się od krótkiego wprowadzenia do standardu SQL/MM, przedstawienia jego genezy i zakresu. Następnie szczegółowo omówione i zilustrowane przykładami będą jego specyfikacje składowe dotyczące danych tekstowych, przestrzennych i obrazów. Dla kompletności, krótko przedstawione będą również pozostałe części standardu, dotyczące eksploracji danych i obsługi danych historycznych.



Standard SQL/MM

- SQL/MM: SQL Multimedia and Application Packages
- Opracowywany i publikowany przez ISO
- Standard obejmujący wiele części:
 - części poświęcone szeroko pojmowanym multimediom
 - części poświęcone specjalistycznym zastosowaniom
- Oparty o SQL i jego typy definiowane przez użytkownika (typy obiektowe SQL99)

Standard SQL/MM (3)

Podobnie jak SQL, SQL/MM jest standardem ISO i składa się z wielu części, przy czym części SQL/MM są ze sobą raczej luźno związane w porównaniu ze specyfikacjami SQL.

Pełna nazwa standardu SQL/MM brzmi: SQL Multimedia and Application Packages, co oznacza że jego części niekoniecznie muszą dotyczyć przetwarzania szeroko pojętych danych multimedialnych, ale również specjalistycznych zastosowań systemów baz danych. Należy też zwrócić uwagę, że z punktu widzenia standardu termin multimedia obejmuje również dane przestrzenne i tekstowe.

SQL/MM stanowi uzupełnienie języka SQL i jest oparty o mechanizmy obiektowo-relacyjne, które pojawiły się w standardzie SQL99.



Geneza standardu SQL/MM

- Wnioski z prac nad rozszerzeniami SQL dla danych tekstowych, przestrzennych, multimedialnych
 - konflikty nazw np. CONTAINS
 - przewidywane trudności w implementacji
- Dostępność obiektowych typów danych od SQL99
 - zamiast rozszerzeń SQL – biblioteki typów obiektowych

Ponieważ standard języka SQL nie zawiera konstrukcji do obsługi takich danych jak multimedia, duże obiekty tekstowe, czy też dane przestrzenne, środowiska zajmujące się tworzeniem oprogramowania do przetwarzania tych specjalistycznych typów danych rozpoczęły pracę nad propozycjami rozszerzenia języka SQL o potrzebne im elementy. Niestety okazało się, że ewentualne rozszerzenia SQL dedykowane dla poszczególnych rodzajów danych mogą być niekompatybilne ze sobą. Najłatwiejszym do zauważenia potencjalnym konfliktem był konflikt słów kluczowych. Przykładowo, słowo kluczowe CONTAINS, używane jest zarówno w kontekście tekstowych baz danych (do wskazania, że dane słowo lub fraza zawiera się w danym fragmencie tekstu), jak i baz przestrzennych i multimedialnych (do wskazania, że jeden obiekt zawiera w sobie inny obiekt).

Ze względu na wspomniane wyżej problemy zarzucono koncepcję rozszerzania języka SQL w zakresie wsparcia dla baz danych tekstowych, przestrzennych i multimedialnych. Zwyciężyła koncepcja opracowania nowego standardu, obejmującego specyfikację bibliotek opartych o typy obiektowe SQL99, przeznaczonych do obsługi poszczególnych specjalistycznych rodzajów danych i aplikacji. Nowy standard natychmiast stał się znany pod nazwą SQL/MM („MM” od „MultiMedia”). Dzięki oparciu specyfikacji SQL/MM o obiektowe typy SQL, funkcjonalność bibliotek jest w sposób naturalny dostępna z poziomu poleceń języka SQL, np. poprzez wywołania metod bibliotecznych typów obiektowych w wyrażeniach języka SQL. Z myślą o użytkownikach niechętnie korzystających z mechanizmów obiektowo-relacyjnych, dla użytkowych metod typów SQL/MM standard specyfikuje odpowiadające im funkcje SQL.



Części standardu SQL/MM

- Part 1: Framework (baza dla pozostałych części)
- Part 2: Full-Text (tekstowe bazy danych)
- Part 3: Spatial (przestrzenne bazy danych)
- Part 5: Still Image (obrazy)
- Part 6: Data Mining (eksploracja danych)
- Part 7: History (dane historyczne)

W chwili obecnej (wiosna 2006) standard SQL/MM obejmuje pięć części (1, 2, 3, 5 i 6) o statusie oficjalnego standardu oraz jedną (7) w stadium Working Draft. Nie ma w standardzie SQL/MM części czwartej. Miała ona dotyczyć ogólnych operacji matematycznych (General Purpose Facilities), ale prace nad nią zarzucono kilka lat temu.

Część pierwsza – Framework ma charakter ogólny. Zawiera ona informacje o zakresie standardu oraz definicje i koncepcje wspólne dla pozostałych, specjalistycznych części. Część pierwsza dotyczy między innymi sposobu, w jaki inne części standardu SQL/MM wykorzystują mechanizm obiektowych typów SQL.

Części druga, trzecia i piąta poświęcone są multimediom w rozumieniu standardu SQL/MM, czyli odpowiednio danym tekstowym, przestrzennym i obrazom (nieruchomym).

Części 6 i 7 dotyczą specjalistycznych zastosowań (Application Packages). Część szósta poświęcona jest eksploracji danych, a siódma obsłudze danych historycznych w bazach danych.

Zwraca uwagę brak części poświęconych danym audio i wideo, których można było się spodziewać w standardzie, którego nazwa sugeruje nacisk na multimedia. Nie wyklucza się opracowania tych części w przyszłości, w zależności od odzewu środowiska na części już istniejące.

Istniejące specyfikacje są ciągle rozwijane. Niektóre części doczekały się już drugiej edycji, a dla niektórych z nich trwają prace nad trzecią edycją.



SQL/MM Full-Text

- Dotyczy danych tekstowych różniących się od znakowych:
 - rozmiarem
 - strukturą (zdania, akapity)
 - typami operacji
- Zakres SQL/MM Full-Text:
 - typy danych dla dokumentów tekstowych oraz wzorców wyszukiwania
 - metody dopasowywania dokumentów do wzorca
 - konwersja do/z typów znakowych
- SQL/MM Full-Text zakłada obsługę wielu języków

SQL/MM Full-Text dotyczy przetwarzania dokumentów tekstowych, które od tradycyjnych łańcuchów znaków odróżnia przede wszystkim znacznie większy rozmiar, ale również obecność struktury (podział na zdania i akapity) oraz charakter operacji wyszukiwania. O ile dla prostych łańcuchów znaków wyszukiwanie sprowadza się do prostego dopasowywania do wzorca, w przypadku danych tekstowych kryteria wyszukiwania dotyczą obecności słów kluczowych lub fraz i uwzględniają odmianę, wymowę, a nawet znaczenie poszczególnych słów czy fraz w kontekście konkretnego języka.

Standard SQL/MM Full-Text definiuje typy danych dla dokumentów tekstowych i wzorców wyszukiwania, a także metody dopasowywania dokumentów tekstowych do wzorca. W zakres standardu wchodzi również metody konwersji danych tekstowych do/z typów znakowych.

SQL/MM Full-Text zakłada obsługę wielu języków. Niemniej, mechanizmy przeszukiwania tekstów objęte standardem SQL/MM są szczególnie odpowiednie dla języków, które łatwo poddają się analizie komputerowej w zakresie wyodrębniania poszczególnych słów i zdań. Do tej klasy języków należą języki zachodnie (w tym np. angielski, niemiecki, polski), w przypadku których słowa oddzielone są odstępami, a zdania znakami interpunkcyjnymi. Inaczej jest np. w języku japońskim, gdzie wyodrębnienie słów z tekstu wymaga analizy kontekstu.



Typy danych SQL/MM Full-Text

- Dokumenty tekstowe reprezentuje typ FullText
- Metody typu FullText:
 - Contains – test binarny (0/1) czy dokument pasuje do wzorca
 - Score – miara zgodności ze wzorcem
- Wzorce reprezentowane jako
 - wystąpienia typu FT_Pattern
 - łańcuchy znaków (CHARACTER VARYING)

SQL/MM Full-Text definiuje kilka typów obiektowych, z których podstawowe znaczenie ma typ FullText, służący do reprezentacji dokumentów tekstowych. Spośród metod typu FullText największe znaczenie mają: Contains i Score. Contains umożliwia sprawdzenie czy dokument pasuje do wzorca i zwraca wartość 1 gdy pasuje a 0 w przeciwnym wypadku. Metoda Score dla danego dokumentu i wzorca zwraca miarę zgodności (ang. relevance) dokumentu z wzorcem w formie liczby zmiennoprzecinkowej. Metoda Score może być wykorzystana do generowania rankingów dokumentów wg dopasowania do zadanego wzorca. Parametrem obu metod jest wzorec. Wzorce mogą być reprezentowane w postaci wystąpień dedykowanego do tego celu typu FT_Pattern albo jako zwykłe łańcuchy znaków (typu CHARACTER VARYING).



Klasy wzorców (1/3)

- Wystąpienie konkretnego słowa lub frazy
 - `tresc.Contains(' "standard" ')`
 - `tresc.Contains(' "standard języka" ')`
- Wzorce ze znakami `_` i `%`
 - `tresc.Contains(' "sport_" ')`
 - `tresc.Contains(' "standard%" ')`
- Wystąpienie co najmniej jednego słowa/frazy z listy
 - `tresc.Contains(' "standard", "język%" ')`
- Wystąpienie formy danego słowa/frazy wg reguł odmiany danego języka (domyślny jest angielski)
 - `tresc.Contains('STEMMED FORM OF POLISH "standard języka" ')`

Standard SQL/MM (8)

Ten slajd i dwa kolejne przedstawiają klasy wzorców, które wg standardu SQL/MM Full-Text mogą być wykorzystane do przeszukiwania kolekcji dokumentów. Dla każdej klasy wzorców przedstawiono przykład (przykłady) jego wykorzystania w metodzie `Contains` na rzecz dokumentu reprezentowanego przez zmienną „`tresc`”.

Pierwsza klasa wzorców żąda wystąpienia konkretnego słowa lub frazy. Druga przy dopasowaniu wzorca do dokumentu wykorzystuje znaki specjalne „`_`” i „`%`” o takim samym znaczeniu jak dla operatora `LIKE` w języku SQL. Trzecia klasa wzorców to wzorce wymagające wystąpienia co najmniej jednego słowa/frazy z listy słów/fraz. Ostatnia przedstawiona na slajdzie klasa wzorców to wzorce wymagające wystąpienia słowa lub frazy lub jego/jej formy gramatycznej, będące(-go)(-j) wynikiem odmiany zgodnie z regułami danego języka.

Standard SQL/MM Full-Text zakłada, że każdy dokument, słowo, czy fraza posiada przypisany język. W definicji wzorców, każde słowo lub fraza mogą być poprzedzone nazwą języka (`ENGLISH`, `GERMAN`, `POLISH`, ...). W przypadku gdy słowo (lub fraza) nie jest poprzedzone we wzorcu nazwą języka, przyjmowany jest język domyślny, jakim wg standardu jest język angielski. W ramach jednego złożonego wzorca mogą występować słowa lub frazy z różnych języków.



Klasy wzorców (2/3)

- Wystąpienie słów/fraz w sąsiedztwie
 - `tresc.Contains(' "standard" NEAR ("SQL", "MPEG") WITHIN 2 PARAGRAPHS ')`
- Wystąpienie słów/fraz w tym samym kontekście
 - `tresc.Contains(' "standard" IN SAME SENTENCE AS "SQL" ')`
- Wzorce złożone za pomocą operatorów logicznych AND, OR i NOT
 - `tresc.Contains(' "SQL/MM" & ("Text" | "Spatial") & NOT "Image" ')`

Standard SQL/MM (9)

Pierwsza z klas wzorców wymienionych na slajdzie wymaga wystąpienia dwóch słów lub fraz w sąsiedztwie ograniczonym podaną liczbą znaków (CHARACTERS), słów (WORDS), zdań (SENTENCES) lub akapitów (PARAGRAPHS) z możliwością wskazania, że podane słowa/frazy muszą wystąpić w wyspecyfikowanej we wzorcu kolejności (IN ORDER). Druga klasa wzorców wymaga wystąpienia dwóch słów/fraz w tym samym kontekście, przy czym kontekstem może być zdanie (SENTENCE) lub akapit (PARAGRAPH). Trzecia klasa wzorców to wzorce złożone, zbudowane z prostych wzorców połączonych logicznymi operatorami AND („&”), OR („|”) i NOT („NOT”). Budując wzorce złożone można wykorzystać nawiasy do wskazania kolejności ewaluacji operatorów. W przypadku braku nawiasów, obowiązuje priorytet operatorów logicznych w następującym porządku (od najwyższego): „NOT”, „&”, „|”.



Klasy wzorców (3/3)

- Wzorce odwołujące się do tematyki tekstu
 - `tresc.Contains('IS ABOUT "język zapytań"')`
- Wzorce odwołujące się do brzmienia w danym języku (domyślny jest angielski)
 - `tresc.Contains(' SOUNDS LIKE "sequel" ')`
- Dopasowanie przybliżone (odporne na „literówki”)
 - `tresc.Contains(' FUZZY FORM OF "teacher" ')`
- Wzorce dopuszczające synonimy
 - `tresc.Contains(THESAURUS "computer science" EXPAND SYNONYM TERM OF "list" ')`

Pierwsza z klas wzorców wymienionych na slajdzie służy do wyszukania dokumentów na dany, ogólnie wyspecyfikowany temat. Druga klasa wzorców odwołuje się do brzmienia słowa/frazy w danym języku (domyślny jest angielski). Kolejna klasa wzorców to dopasowanie odporne na tzw. „literówki”, które może odnaleźć teksty, w których popełniono błędy ortograficzne lub drobne błędy edycyjne. Ostatnia z przedstawionych na slajdzie klas wzorców służy do wyszukania dokumentów zawierających słowa/frazy wywiedzione z podanego słowa lub frazy (np. synonimy). W definicji wzorca tego typu przede wszystkim należy wskazać nazwę słownika wyrazów bliskoznacznych (THESAURUS), który ma być wykorzystany do generacji związanych terminów. Dostępne klauzule wskazujące metodę generacji wyszukiwanych słów i fraz to:

- (a) EXPAND SYNONYM TERM OF – synonimy,
- (b) EXPAND BROADER TERM OF – słowa/frazy o szerszym znaczeniu,
- (c) EXPAND NARROWER TERM OF – słowa/frazy o węższym znaczeniu,
- (d) EXPAND TOP TERM OF – najbardziej ogólne terminy z słów/fraz o szerszym znaczeniu,
- (e) EXPAND PREFERRED TERM OF – preferowane terminy ze zbioru synonimów,
- (f) EXPAND RELATED TERM OF – terminy powiązane.



SQL/MM Full-Text – Przykład

```
CREATE TABLE raporty (  
  numer INTEGER,  
  tresc FULLTEXT  
)
```

```
SELECT numer  
FROM raporty  
WHERE tresc.CONTAINS('STEMMED FORM OF "standard"  
  IN SAME PARAGRAPH AS SOUNDS LIKE "sequel"') = 1
```

Na slajdzie pokazano przykład wykorzystania możliwości SQL/MM Full-Text z poziomu języka SQL. Pierwsze polecenie tworzy tabelę RAPORTY, której pierwsza kolumna typu INTEGER zawiera numery raportów, a druga kolumna typu FULLTEXT treści raportów. Drugie polecenie to zapytanie zwracające numery raportów zawierających w tym samym akapicie dowolną formę słowa „standard” wg reguł odmiany języka angielskiego (język nie został podany jawnie, a angielski jest domyślny) i słowo brzmiące jak „sequel”.



SQL/MM Spatial

- Dotyczy danych przestrzennych
 - figury geometryczne, lokalizacja i topologia obiektów
- Zakres SQL/MM Spatial:
 - typy danych dla danych przestrzennych
 - tworzenie i porównywanie geometrycznych obiektów przestrzennych, obliczanie wartości miar
 - konwersje do/z typów znakowych i binarnych
- Przede wszystkim zastosowania geograficzne
- Jeden z kilku standardów dla danych przestrzennych

SQL/MM Spatial definiuje obiektowe typy danych i ich metody do przetwarzania danych przestrzennych tj. dotyczących geometrii, lokalizacji i topologii obiektów. Przetwarzanie danych przestrzennych obejmuje tworzenie i porównywanie obiektów geometrycznych oraz obliczanie wartości miar takich jak np. pole powierzchni. W zakres standardu wchodzi również metody konwersji między typami SQL/MM Spatial a innymi znakowymi i binarnymi reprezentacjami danych przestrzennych.

Specyfikacja SQL/MM Spatial jest szczególnie zorientowana na przetwarzanie danych w systemach informacji geograficznej (GIS), ale obszar jej zastosowań wykracza poza przetwarzanie informacji o obiektach na powierzchni Ziemi i obejmuje również np. projektowanie układów elektronicznych.

SQL/MM Spatial jest związany z dwoma innymi, rozwijanymi równolegle, standardami dla danych przestrzennych opracowywanymi przez ISO Technical Committee TC 211 i Open GIS Consortium.



Możliwości SQL/MM Spatial

- W chwili obecnej standard wspiera dane:
 - 0-wymiarowe (punkty)
 - 1-wymiarowe (linie)
 - 2-wymiarowe (figury płaskie)
- Wsparcie dla testów przecinania, pokrywania, itp.
- Wsparcie dla różnych przestrzennych układów odniesienia
 - przede wszystkim opisujących geografie naszej planety i jej regionów (krzywizna Ziemi)

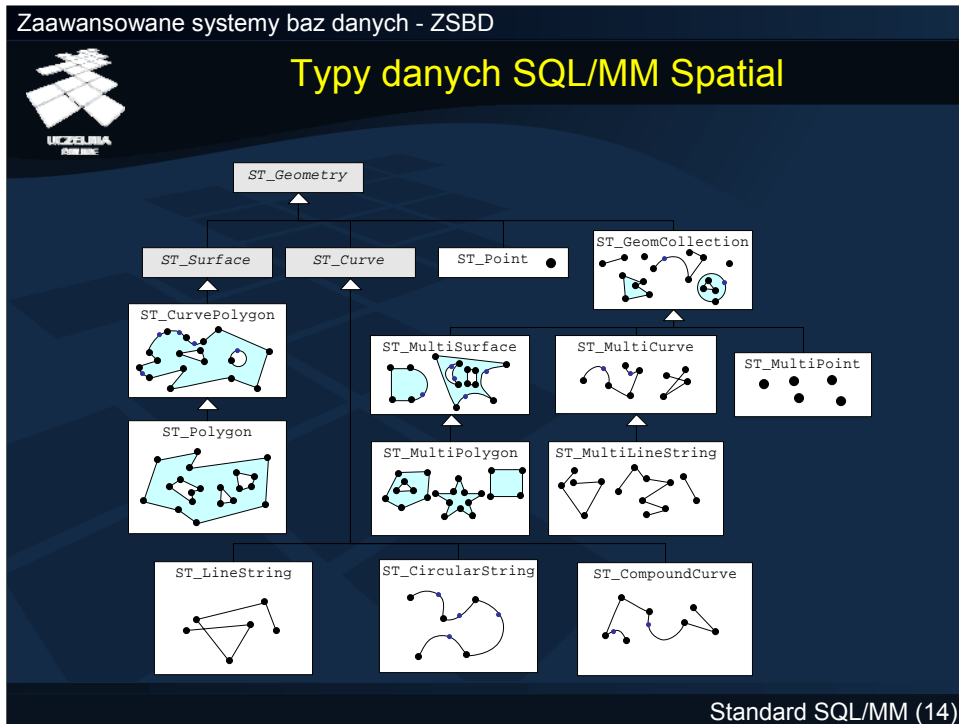
Standard SQL/MM (13)

SQL/MM Spatial obecnie wspiera dane 0-wymiarowe (punkty), 1-wymiarowe (linie) i 2-wymiarowe (figury płaskie). W rzeczywistych zastosowaniach pojawiają się oczywiście również obiekty 3-wymiarowe, niestety obecna specyfikacja SQL/MM (edycja druga) ich nie uwzględnia.

Dla obiektów geometrycznych standard przewiduje testy wzajemnego położenia obiektów np. przecinania, pokrywania, zawierania.

Dużą wagę specyfikacja SQL/MM Spatial przykładła do przestrzennych układów odniesienia (ang. spatial reference system). Większość z nich opisuje geografie naszej planety i jej regionów, co jest konsekwencją faktu, że typowymi użytkownikami wymagającymi przetwarzania danych przestrzennych są władze państwowe i lokalne oraz duże korporacje, operujące na danych geograficznych. Wybór układu wpływa m.in. na metodę wyznaczania odległości.

Typy danych SQL/MM Spatial



Podstawowy zbiór typów danych SQL/MM Spatial to typy reprezentujące poszczególne figury (kształty) geometryczne. Typy te tworzą hierarchię przedstawioną na slajdzie. Korzeniem tej hierarchii jest abstrakcyjny typ `ST_Geometry`. Inne typy abstrakcyjne, pełniące rolę węzłów pośrednich w hierarchii to `ST_Curve` (krzywa) i `ST_Surface` (2-wymiarowa figura). Typy `ST_MultiCurve` (kolekcja krzywych) i `ST_MultiSurface` (kolekcja figur 2-wymiarowych) mogą być abstrakcyjne lub nie - zależnie od implementacji. W hierarchii występują ponadto następujące nieabstrakcyjne podtypy: `ST_Point` (punkt), `ST_LineString` (sekwencja odcinków prostych), `ST_CircularString` (sekwencja łuków – odcinków okręgu), `ST_CompoundCurve` (sekwencja połączonych krzywych), `ST_CurvePolygon` (figura 2-wymiarowa o jednym brzegu i opcjonalnych dziurach w formie zamkniętej sekwencji połączonych krzywych), `ST_Polygon` (wielokąt z opcjonalnymi dziurami w kształcie wielokątów), `ST_GeomCollection` (kolekcja figur geometrycznych), `ST_MultiPoint` (kolekcja punktów), `ST_MultiLineString` (kolekcja `ST_LineString`) i `ST_MultiPolygon` (kolekcja `ST_Polygon`).

Spośród pozostałych typów SQL/MM Spatial należy wymienić `ST_SpatialRefSystem` służący do opisu przestrzennych układów odniesienia oraz typy `ST_Angle` i `ST_Direction` reprezentujące odpowiednio kąty i kierunki.



Metody typów SQL/MM Spatial

- Metody do odczytu właściwości i miar obiektów
 - np. ST_Boundary, ST_Length, ST_Area
- Metody do porównywania obiektów geometrycznych
 - ST_Equals, ST_Disjoint, ST_Intersects, ST_Crosses, ST_Overlaps, ST_Touches, ST_Within, ST_Contains i ST_Distance
- Metody do tworzenia nowych obiektów geometrycznych
 - ST_Difference, ST_Intersection i ST_Union
- Metody do konwersji z/do zewnętrznych formatów danych

Metody typów SQL/MM Spatial można podzielić na 4 grupy:

- (a) metody do odczytu właściwości i miar dla obiektów geometrycznych,
- (b) metody do porównywania obiektów geometrycznych,
- (c) metody do tworzenia nowych obiektów geometrycznych,
- (d) metody do konwersji z/do zewnętrznych formatów danych (tekstowych, binarnych).

Przykłady metod do odczytu właściwości i miar dla obiektów geometrycznych to:

ST_Boundary zwracająca brzeg figury geometrycznej, ST_Length zwracająca długość krzywej, ST_Area zwracająca pole powierzchni figury 2-wymiarowej.

Metody do porównywania obiektów geometrycznych to ST_Equals (do testowania czy figury są równe), ST_Disjoint (do testowania czy figury są rozłączne), ST_Intersects, ST_Crosses i ST_Overlaps (trzy bardzo podobne metody do testowania czy wnętrza figur mają część wspólną), ST_Touches (do testowania czy figury stykają się brzegiem), ST_Within i ST_Contains (do testowania czy jedna figura zawiera się w drugiej) oraz ST_Distance (do wyznaczania minimalnej odległości między punktami dwóch figur).

Metody do tworzenia nowych obiektów geometrycznych to przede wszystkim ST_Difference, ST_Intersection i ST_Union zwracające odpowiednio różnicę, część wspólną i sumę figur.

Przykładowe metody do konwersji danych to metody do konwersji między typami SQL/MM Spatial a formatem GML (Geography Markup Language), np. ST_MPointFromGML, ST_LineFromGML, ST_AsGML.



SQL/MM Spatial – Przykład

```
CREATE TABLE kraje (  
  nazwa_kraju VARCHAR(30),  
  lokalizacja ST_GEOMETRY )
```

```
SELECT lokalizacja.area  
FROM kraje  
WHERE nazwa_kraju = 'POLSKA'
```

```
SELECT nazwa_kraju FROM kraje  
WHERE lokalizacja.ST_Touches(  
  SELECT lokalizacja FROM kraje  
  WHERE nazwa_kraju = 'POLSKA')
```

Na slajdzie pokazano przykład wykorzystania możliwości SQL/MM Spatial z poziomu języka SQL. Pierwsze polecenie tworzy tabelę KRAJE, której pierwsza kolumna zawiera nazwy krajów, a druga kolumna figury geometryczne reprezentujące ich lokalizacje i kształty w postaci obiektów ST_GEOMETRY. Drugie polecenie to zapytanie zwracające pole powierzchni Polski. Trzecie polecenie to zapytanie zwracające nazwy krajów graniczących z Polską.



SQL/MM Still Image

- Dotyczy obrazów nieruchomych (np. fotografii)
 - obrazy bitmapowe, wsparcie dla różnych formatów
- Zakres SQL/MM Still Image:
 - typy danych dla obrazów i ich właściwości
 - metody do modyfikacji obrazów
 - wyszukiwanie w oparciu o zawartość
- Nie obejmuje semantycznych opisów zawartości
- Specyfikacja zaimplementowana w Oracle10g

Specyfikacja SQL/MM Still Image dotyczy obrazów nieruchomych, takich jak np. fotografie. SQL/MM Still Image przyjmuje, że obraz jest dwuwymiarową tablicą pikseli (obraz bitmapowy). Zakłada się, że implementacje standardu powinny wspierać wiele różnych formatów graficznych np. JPEG, GIF, TIFF.

Standard definiuje strukturalne typy SQL umożliwiające składowanie obrazów w bazie danych, reprezentację właściwości wizualnych obrazów, podstawowe operacje przetwarzania obrazów (skalowanie, obroty) oraz wyszukiwanie obrazów w oparciu o zawartość.

SQL/MM dla obrazów nie standaryzuje metod opisu zawartości semantycznej. Nie stanowi więc w tym zakresie konkurencji dla standardu MPEG-7.

Still Image jest obecnie (wiosna 2006) jedyną częścią standardu SQL/MM, która doczekała się implementacji w rzeczywistym systemie zarządzania bazą danych. SQL/MM Still Image jest wspierany przez system Oracle od wersji 10g.



Wyszukiwanie w oparciu o zawartość w SQL/MM Still Image

- Content-Based Image Retrieval
- Realizowane poprzez wyszukiwanie obrazów podobnych do wzorca cech wizualnych
- Uwzględniane właściwości wizualne obrazów:
 - średni kolor
 - histogram kolorów (udział kolorów w obrazie)
 - lokalizacja kolorów
 - tekstura
- Możliwość przypisania wag poszczególnym właściwościom
- Miara odległości obrazów

SQL/MM Still Image specyfikuje typy danych i ich metody umożliwiające wyszukiwanie obrazów w oparciu o zawartość (ang. Content-Based Image Retrieval). Wyszukiwanie jest realizowane względem wzorca cech wizualnych, typowo wywiedzionego z obrazu podanego jako przykład. SQL/MM Still Image umożliwia uwzględnienie następujących właściwości wizualnych przy wyszukiwaniu:

- (a) średni kolor,
- (b) histogram kolorów (udział kolorów w obrazie),
- (c) lokalizacja kolorów,
- (d) tekstura.

Standard przewiduje możliwość przypisywania wag poszczególnym właściwościom i specyfikuje miarę podobieństwa (a właściwie odległości) obrazów.



Typy danych SQL/MM Still Image

- SI_StillImage – reprezentuje obraz
- SI_AverageColor – reprezentuje średni kolor obrazu
- SI_ColorHistogram – reprezentuje histogram kolorów
- SI_PositionalColor – reprezentuje lokalizację kolorów obrazu
- SI_Texture – reprezentuje teksturę obrazu
- SI_Color – reprezentuje kolor
- SI_FeatureList – reprezentuje listę właściwości wizualnych obrazu

W standardzie SQL/MM obrazy są reprezentowane za pomocą typu SI_StillImage. Atrybuty typu SI_StillImage to: SI_content (obraz w postaci binarnej, typu Binary Large Object), SI_contentLength (rozmiar w bajtach), SI_reference (referencja typu DATALINK do zewnętrznego źródła danych przechowującego obraz), SI_format (format graficzny obrazu), SI_height (wysokość obrazu w pikselach), SI_width (szerokość obrazu w pikselach). Metody typu SI_StillImage pozwalają m.in. na skalowanie i obroty obrazu, konwersje między formatami graficznymi oraz generację miniaturki obrazu w mniejszej rozdzielczości (ang. thumbnail).

Oprócz podstawowego typu SI_StillImage, SQL/MM Still Image definiuje również kilka typów służących do reprezentacji właściwości obrazu. Typ SI_AverageColor reprezentuje średni kolor obrazu, SI_ColorHistogram reprezentuje udział kolorów w obrazie w formie histogramu, SI_PositionalColor reprezentuje lokalizację kolorów na obrazie, a SI_Texture służy do zapamiętania informacji o teksturze obrazu. Właściwości odnoszące się do kolorystyki obrazu reprezentują poszczególne kolory jako wartości pomocniczego typu danych SI_Color.

Ponadto, SQL/MM Still Image przewiduje jeszcze typ danych SI_FeatureList do reprezentacji listy właściwości wizualnych obrazu wraz z przypisanymi im wagami. Typ ten jest wykorzystywany do testów podobieństwa uwzględniających więcej niż jedną właściwość obrazu.



Metody typu SI_StillImage

- Konstruktory
- Metody do odczytu atrybutów
- Metody do zmiany zawartości binarnej obrazu
- Metody przetwarzania obrazu
- Metody do testów podobieństwa

Metody typu SI_StillImage można podzielić na następujące grupy:

1. Konstruktory, m.in.: SI_StillImage(BINARY LARGE OBJECT) – tworzący obiekt SI_StillImage, dla podanego obiektu binarnego zawierającego obraz; SI_StillImage(DATALINK) – tworzący obiekt SI_StillImage, dla podanej referencji do obiektu binarnego zawierającego obraz.
2. Metody do odczytu atrybutów: SI_Content – zwracająca zawartość binarną w postaci BINARY LARGE OBJECT, SI_ContentLength – zwracająca rozmiar zawartości binarnej w bajtach, SI_Height – zwracająca wysokość obrazu w pikselach, SI_Width – zwracająca szerokość obrazu w pikselach, SI_Format – zwracająca format graficzny obrazu.
3. Metody do zmiany zawartości binarnej obrazu: SI_SetContent – ustawiająca nową zawartość binarną (podmiana obrazu), SI_ChangeFormat – dokonująca konwersji formatu obrazu.
4. Metody do przetwarzania obrazu: SI_Scale, SI_Resize, SI_Thumbnail, SI_Rotate (omówione na następnym slajdzie).
5. Metody do testów podobieństwa: przeciążona metoda SI_Score (omówiona na jednym z kolejnych slajdów).

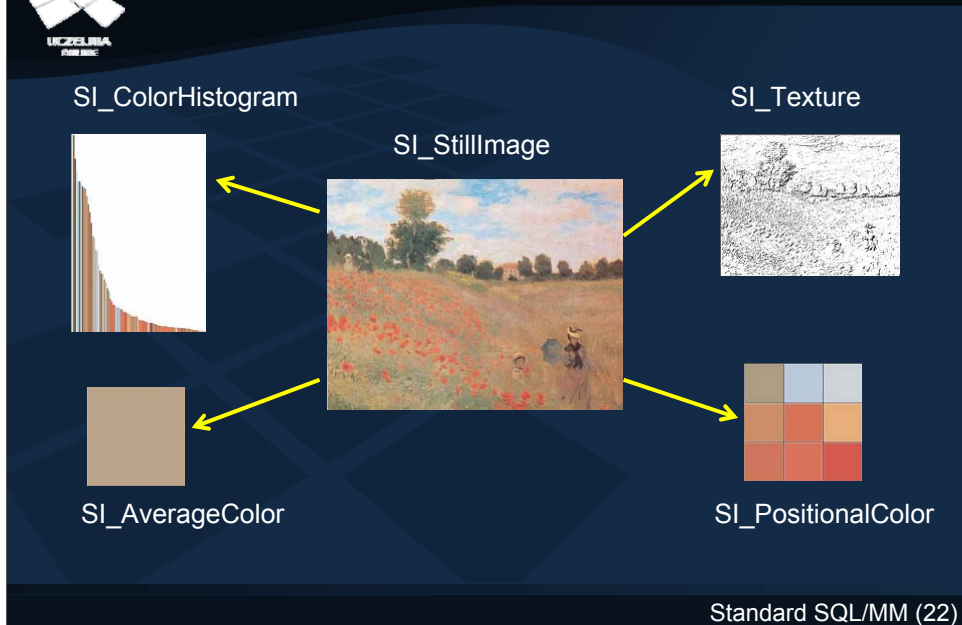


Przetwarzanie obrazów



Slajd ilustruje działanie metod typu `SI_StillImage` do przetwarzania obrazu: `SI_Resize` – dokonującej zmiany rozmiaru obrazu, bez zachowania proporcji; `SI_Scale` – dostępnej w dwóch wersjach; skalującej obraz z zachowaniem proporcji w ramach podanego okna lub wg podanego współczynnika skalowania; `SI_Thumbnail` – generującej miniaturę obrazu; `SI_Rotate` – obracającej obraz o zadany kąt.

Właściwości wizualne obrazu



Slajd ilustruje znaczenie czterech dostępnych w SQL/MM Still Image typów właściwości wizualnych obrazu.

Każdą z właściwości można wyznaczyć dla obrazu za pomocą konstruktora typu właściwości, przyjmującego jako parametr obiekt `SI_StillImage`: `SI_AverageColor(StillImage)`, `SI_ColorHistogram(StillImage)`, `SI_PositionalColor(StillImage)`, `SI_Texture(StillImage)`.

Jako alternatywę dla powyższych konstruktorów (wygodną szczególnie w zapytaniach SQL) standard definiuje poniższe cztery funkcje SQL, których działanie sprowadza się do wywołania stosownych konstruktorów: `SI_findAvgColor(StillImage)`, `SI_findClrHstgr(StillImage)`, `SI_findPstnlClr(StillImage)`, `SI_findTexture(StillImage)`.

Sposób opisu tekstury, algorytm jej wyznaczania oraz algorytm określania podobieństwa względem tekstury nie są zawarte w standardzie i będą zależne od implementacji. Dla pozostałych trzech właściwości standard wskazuje sposób ich reprezentacji i ogólne algorytmy ich ekstrakcji z obrazu.



Algorytm wyznaczania średniego koloru

- Obraz jest próbkowany za pomocą siatki referencyjnej
- Komponenty koloru R, G i B z poszczególnych próbek są niezależnie od siebie sumowane i dzielone przez liczbę próbek

$$(\bar{R}, \bar{G}, \bar{B}) = \left(\frac{\sum_{i=1}^{SI_width} \sum_{j=1}^{SI_height} R(i, j)}{SI_width * SI_height}, \frac{\sum_{i=1}^{SI_width} \sum_{j=1}^{SI_height} G(i, j)}{SI_width * SI_height}, \frac{\sum_{i=1}^{SI_width} \sum_{j=1}^{SI_height} B(i, j)}{SI_width * SI_height} \right)$$

- Wynikowy średni kolor jest reprezentowany jako jedna właściwość typu SI_Color

SQL/MM proponuje następujący algorytm wyznaczania średniego koloru dla obrazu. Najpierw obraz jest próbkowany tzw. siatką referencyjną, wskazującą piksele, z których będzie uśredniany kolor. Następnie komponenty koloru z poszczególnych próbek są niezależnie uśredniane poprzez zsumowanie ich i podział sum przez liczbę próbek. Na slajdzie przedstawiono wzór na średni kolor przy założeniu, że siatka referencyjna jest tworzona przez wszystkie piksele obrazu. Wynik uśredniania jest zapisany w obiekcie SI_AverageColor w jego właściwości typu SI_Color.



Algorytm generacji histogramu kolorów

- Przestrzeń kolorów dzielona jest na obszary
- Każdy obszar przestrzeni kolorów obejmuje pewien zbiór kolorów i jest reprezentowany przez jeden kolor C_i
- Dla każdego obszaru kolorów określana jest jego częstotliwość F_i w ramach obrazu poprzez iterację po wszystkich pikselach
- Każdy piksel zwiększa wartość F_i o 1 dla tego zakresu, do którego należy jego kolor
- Wartości F_i są normalizowane, aby zawierały się w zakresie od 0 do 100
- Wynikowy histogram jest sekwencją par (kolor, częstotliwość) w postaci dwóch tablic (ARRAY)

SQL/MM proponuje następujący algorytm generacji histogramu kolorów dla obrazu. Przestrzeń kolorów dzielona jest na pewną liczbę obszarów, zależną od implementacji. Każdy obszar przestrzeni kolorów obejmuje pewien zbiór kolorów i jest reprezentowany przez jeden kolor C_i . Dla każdego obszaru kolorów określana jest jego częstotliwość F_i w ramach obrazu, poprzez iterację po wszystkich pikselach. Każdy piksel zwiększa wartość licznika F_i o 1 dla tego obszaru kolorów, do którego należy jego kolor. Po zliczeniu częstotliwości, wartości F_i są normalizowane, aby zawierały się w zakresie od 0 do 100. Wynikowy histogram, logicznie będący sekwencją par (kolor, częstotliwość), fizycznie zapamiętywany jest w postaci dwóch tablic (ARRAY). Pierwsza z tablic zawiera kolory, a druga odpowiadające im częstotliwości. Tablice mają oczywiście taką samą liczbę elementów, zależną od implementacji, odpowiadającą maksymalnej dostępnej w implementacji długości histogramu.



Algorytm wyznaczania lokalizacji kolorów

- Obraz jest dzielony na „siatkę” $m \times n$ prostokątów
- Dla każdego z prostokątów wyznaczany jest dominujący kolor
- Wynikowa lokalizacja kolorów jest reprezentowana w postaci tablicy (ARRAY) obiektów `SI_Color`

SQL/MM proponuje następujący algorytm ekstrakcji lokalizacji kolorów dla obrazu. Najpierw obraz dzielony jest na siatkę $m \times n$ prostokątów, o wymiarach zależnych od implementacji. Następnie dla każdego z prostokątów wyznaczany jest kolor dominujący, jako kolor o największej częstotliwości w histogramie kolorów wygenerowanym dla prostokąta (algorytmem przedstawionym na poprzednim slajdzie). Wynikowa lokalizacja kolorów jest reprezentowana w postaci tablicy (ARRAY) obiektów `SI_Color`.



Typ danych SI_FeatureList

- Służy do reprezentacji zbioru właściwości wizualnych
- Wykorzystywany do testów podobieństwa obrazów z uwzględnieniem więcej niż jednej właściwości
- Reprezentuje złożony wzorec do testów podobieństwa
- Umożliwia przypisanie wag poszczególnym właściwościom

Zbiór właściwości wizualnych obrazu można zapamiętać jako jeden obiekt typu SI_FeatureList. Podstawowym przeznaczeniem tego typu danych jest umożliwienie przeprowadzania testów podobieństwa obrazów z uwzględnieniem więcej niż jednej właściwości. Obiekty typu SI_FeatureList stanowią złożone wzorce do testów podobieństwa obrazów, w których każdej uwzględnionej właściwości przypisana jest waga. Waga jest liczbą zmiennoprzecinkową z przedziału <0.0, 1.0>.

Złożony wzorec można utworzyć następującym konstruktorem typu SI_FeatureList: SI_FeatureList(SI_AverageColor, DOUBLE PRECISION, SI_ColorHistogram, DOUBLE PRECISION, SI_PositionalColor, DOUBLE PRECISION, SI_Texture, DOUBLE PRECISION). Parametry typu zmiennoprzecinkowego to wagi przypisane poprzedzającym je na liście argumentów właściwościom. W przypadku gdy dana właściwość ma nie być zawarta we wzorcu, należy podać NULL jako wartość właściwości i jej wagi lub przypisać właściwości wagę 0.

Wzorec reprezentowany przez obiekt SI_FeatureList można modyfikować za pomocą metod SI_SetFeature, z których każda ustawia wartość danej właściwości i przypisuje jej wagę: SI_SetFeature(SI_AverageColor, DOUBLE PRECISION), SI_SetFeature(SI_ColorHistogram, DOUBLE PRECISION), SI_SetFeature(SI_PositionalColor, DOUBLE PRECISION), SI_SetFeature(SI_Texture, DOUBLE PRECISION).



Testowanie podobieństwa obrazów

- Metody SI_Score typu SI_StillImage z argumentem typu pojedynczej właściwości lub SI_FeatureList
- Metody SI_Score typów pojedynczych właściwości lub SI_FeatureList z argumentem typu SI_StillImage
- Brak metody do bezpośredniego porównywania dwóch obiektów SI_StillImage ze sobą!
- Wartości SI_Score ≥ 0 (0 – najlepsze dopasowanie)
- Dla testów poprzez SI_FeatureList wynik SI_Score to suma ważona wyników testów dla właściwości składowych

$$\frac{\sum_{i=1}^N F_i \cdot SI_Score(image) \cdot W_i}{\sum_{i=1}^N W_i}$$

Standard SQL/MM (27)

Do realizacji testów podobieństwa na potrzeby wyszukiwania obrazów w oparciu o zawartość służą metody:

1. Przeciążona metoda SI_Score typu SI_StillImage: SI_Score(SI_AverageColor), SI_Score(SI_ColorHistogram), SI_Score(SI_PositionalColor), SI_Score(SI_Texture), SI_Score(SI_FeatureList).

2. Metody SI_Score(SI_StillImage) typów SI_AverageColor, SI_ColorHistogram, SI_PositionalColor, SI_Texture i SI_FeatureList.

Zwraca uwagę brak metody do bezpośredniego porównywania dwóch obiektów SI_StillImage ze sobą. W celu porównania dwóch obrazów należy najpierw z jednego z nich wyekstrahować właściwość lub właściwości wizualne, które mają być uwzględnione w teście podobieństwa. Jeśli test ma dotyczyć więcej niż jednej właściwości, należy w kolejnym kroku utworzyć obiekt SI_FeatureList przypisując poszczególnym właściwościom wagi.

Wynikiem zwracanym przez metody SI_Score jest nieujemna liczba zmiennoprzecinkowa. Im mniejsza wartość tym lepiej wzorzec właściwości reprezentuje porównywany z nim obraz. Wartość 0 oznacza najlepsze dopasowanie do wzorca.

W przypadku testów realizowanych z wykorzystaniem obiektu SI_FeatureList, końcowy wynik metody SI_Score jest średnią ważoną wyników metod SI_Score dla poszczególnych właściwości, co ilustruje wzór przedstawiony na slajdzie. N jest liczbą właściwości zawartych w obiekcie SI_FeatureList, F_i to i-ta właściwość, a W_i jej waga.



SQL/MM Still Image – Przykład 1

```
CREATE TABLE flagi
(panstwo VARCHAR2(40),
flaga SI_StillImage);
```

```
SELECT f1.flaga
FROM flagi f1
WHERE SI_findAvgClr(
    (SELECT f2.flaga FROM flagi f2
    WHERE f2.panstwo = 'Polska')
).SI_Score(f1.flaga) < 30;
```

Na slajdzie pokazano przykład wykorzystania możliwości SQL/MM Still Image z poziomu języka SQL. Pierwsze polecenie tworzy tabelę FLAGI, której pierwsza kolumna zawiera nazwy państw, a druga kolumna obrazki reprezentujące ich flagi w postaci obiektów SI_StillImage. Drugie polecenie to zapytanie zwracające flagi podobne do polskiej w sensie średniego koloru. W podzapytaniu z bazy danych jest odczytywany obrazek z polską flagą, następnie wyznaczany jest dla niego funkcją SQL SI_findAvgClr średni kolor jako wzorec do testów podobieństwa. W warunku WHERE zapytania zewnętrznego testowana jest odległość poszczególnych obrazów z bazy danych od tego wzorca średniego koloru. Za podobne uznawane są obrazy, dla których odległość od wzorca jest mniejsza niż 30.

Należy podkreślić, że dobór progu podobieństwa (w przykładzie: 30) jest zadaniem trudnym. Selektywność danej wartości progowej zależy od natury przeszukiwanej kolekcji obrazów. Dlatego też aby odpowiednio dobrać wartość progu odległości, może być konieczne wykonanie zapytania dla kilku różnych wartości progu aż do uzyskania satysfakcjonujących wyników.



SQL/MM Still Image – Przykład 2 (1/2)

```

DECLARE
    tempimage          SI_StillImage;
    tempAvgColor       SI_AverageColor;
    tempTexture        SI_Texture;
    myFeatureList      SI_FeatureList;
    score              DOUBLE PRECISION;
BEGIN
    SELECT flaga INTO tempimage FROM flagi
    WHERE panstwo = 'Polska';

    tempAvgColor := NEW SI_AverageColor(tempimage);
    tempTexture := NEW SI_Texture(tempimage);
    myFeatureList := NEW SI_FeatureList(
        tempAvgColor, 0.7, NULL, NULL,
        NULL, NULL, tempTexture, 0.3);
    ...

```

1

2

3

4

Standard SQL/MM (29)

Ten i następny slajd pokazują przykład wyszukania obrazów podobnych do zadanego z uwzględnieniem więcej niż jednej właściwości. W takim wypadku niezbędne jest najpierw utworzenie złożonego wzorca podobieństwa obrazów w postaci obiektu `SI_FeatureList`, a następnie wyznaczenie wartości odległości poszczególnych obrazów od tego wzorca. Odpowiednią sekwencję operacji wygodnie można zapisać w proceduralnym języku pozwalającym na zagnieżdżanie poleceń SQL, takim jak np. PL/SQL na platformie Oracle.

Niniejszy slajd pokazuje pierwszą część anonimowego bloku PL/SQL, którego zadaniem jest wyświetlenie nazw państw, których flagi są podobne do polskiej w sensie średniego koloru i tekstury, przy czym średni kolor ma mieć większą wagę (0.7) niż tekstura (0.3).

W sekcji deklaracji zmiennych (1) deklarowane są kolejno zmienne, w których pamiętane będą: obrazek z polską flagą, jego średni kolor, jego tekstura, wzorzec do testów podobieństwa, wynik bieżącego testu podobieństwa. Działanie bloku kodu rozpoczyna się pobraniem z bazy danych obrazka z polską flagą do zmiennej w programie (2). Następnie za pomocą konstruktorów `SI_AverageColor` i `SI_Texture` wyznaczane są z obrazka jego właściwości średniego koloru i tekstury (3). Wreszcie z tych dwóch właściwości tworzony jest obiekt `SI_FeatureList` z uwzględnieniem zadanych wag (4).



SQL/MM Still Image – Przykład 2 (2/2)

```
...  
FOR c IN (SELECT panstwo, flaga FROM flagi) LOOP  
  score := myFeatureList.SI_Score(c.flaga);  
  IF score < 10 THEN  
    dbms_output.put_line(c.panstwo);  
  END IF;  
END LOOP;  
END;  
/
```

5

6

7

Standard SQL/MM (30)

Po przygotowaniu złożonego wzorca do testów podobieństwa, w pętli FOR z podzapytaniem przeglądane są kolejne flagi odczytane z bazy danych (5). Dla każdej z nich metodą SI_Score wyznaczana jest jej odległość od wzorca (6). Jeśli wartość ta jest mniejsza niż 10, nazwa państwa wyświetlana jest na konsoli (7).



SQL/MM Data Mining

- Biblioteka typów i ich metod do eksploracji danych
- Ma umożliwiać wymianę modeli między systemami
- Obecnie wspiera cztery metody eksploracji danych:
 - odkrywanie reguł asocjacyjnych
 - klasyfikację
 - analizę skupień
 - regresję
- Jeden z wielu standardów dotyczących eksploracji danych

Standard SQL/MM (31)

Eksploracja danych (ang. data mining) to odkrywanie wcześniej nieznanymi, nietrywialnymi, potencjalnie interesujących wzorców, reguł, modeli z dużych zbiorów danych. SQL/MM Data Mining stanowi próbę dostarczenia standardowych interfejsów do algorytmów eksploracji danych na gruncie obiektowo-relacyjnych baz danych. Ma na celu umożliwienie wymiany modeli eksploracji danych między różnymi systemami. Standard nie określa algorytmów implementujących poszczególne modele ani formatu przechowywanych danych dla eksploracji danych.

Obecnie SQL/MM Data Mining wspiera cztery podstawowe metody eksploracji danych: odkrywanie reguł asocjacyjnych, klasyfikację, analizę skupień i regresję.

SQL/MM Data Mining jest jednym z wielu standardów dotyczących eksploracji danych. Rozważane jest uzgodnienie go w przyszłości ze standardem specyfikującym interfejs do eksploracji danych dla języka Java (Java Data Mining API).



SQL/MM History

- Obsługa historii zmian dla tabel bazy danych
 - składowanie danych historycznych
 - zarządzanie danymi historycznymi
 - udostępnianie danych historycznych
- Historia tabeli rozumiana jako sekwencja operacji INSERT, UPDATE i DELETE wykonanych na niej w przeszłości

Ostatnią obecnie częścią standardu SQL/MM jest część siódma – History. Nie doczekała się ona jeszcze statusu obowiązującego standardu. Aktualnie (wiosna 2006) jest w stadium Working Draft.

SQL/MM History dotyczy obsługi historii zmian dla tabel bazy danych, obejmującej składowanie danych historycznych, zarządzanie nimi i udostępnianie ich z poziomu zapytań SQL.

Historia tabeli rozumiana jako sekwencja operacji INSERT, UPDATE i DELETE wykonanych na niej w przeszłości.

SQL/MM History ma umożliwić np. sprawdzenie jaką wartość miał atrybut wiersza tabeli w danym momencie w przeszłości.



Podsumowanie

- SQL/MM uzupełnia SQL o biblioteki do obsługi specjalistycznych danych i aplikacji
- SQL/MM: SQL Multimedia and Application Packages
- Jeszcze za wcześnie by mówić o akceptacji standardu
 - jedyna istniejąca implementacja to SQL/MM Still Image w Oracle 10g
 - funkcjonalność definiowana przez standard dostępna w istniejących SZBD, ale inna składnia

SQL/MM uzupełnia język SQL o biblioteki do obsługi specjalistycznych danych i aplikacji. Ma z założenia dotyczyć szeroko pojętych danych multimedialnych i zaawansowanych zastosowań systemów baz danych. Obecnie obejmuje obsługę danych tekstowych, przestrzennych i obrazów oraz eksplorację danych i obsługę danych historycznych.

W chwili obecnej (wiosna 2006) trudno jeszcze mówić o sukcesie standardu. SQL/MM jest ciągle rozwijany, ale prawie nie pojawiają się jego implementacje. Jedyną zaimplementowaną częścią standardu jest SQL/MM Still Image, wspierany w Oracle 10g. Dostępne systemy zarządzania bazami danych w dużym stopniu oferują już funkcjonalność przewidywaną przez standard SQL/MM. Pozostaje jednak kwestia dostosowania interfejsu do wyznaczonego przez standard.



Materiały dodatkowe

- ISO/IEC 13249, Information Technology – Database Languages – SQL Multimedia and Application Packages (specyfikacja standardu ISO)
- Melton J., Eisenberg A.: SQL Multimedia and Application Packages (SQL/MM). SIGMOD Record 30(4), 2001
- Stolze K.: SQL/MM Spatial: The Standard to Manage Spatial Data in Relational Database Systems. BTW 2003