



KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI



UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY



Studia podyplomowe "Inżynieria oprogramowania" współfinansowane przez Unię Europejską w ramach Europejskiego Funduszu Społecznego

Projekt "Studia podyplomowe z zakresu wytwarzania
oprogramowania oraz zarządzania projektami w firmach
informatycznych" realizowany w ramach
Programu Operacyjnego Kapitał Ludzki

Konstruowanie Baz Danych

Relacyjny Model Danych II

Antoni Ligeza
ligeza@agh.edu.pl
<http://home.agh.edu.pl/~ligeza>
<http://home.agh.edu.pl/~ligeza/wiki>

Logiczny model danych

Formuły atomowe

Reprezentacji danych za pomocą atrybutów polega na wyspecyfikowaniu wartości każdego atrybutu przyjętego do opisu pewnego obiektu. Przyjmuje się, że określona jest funkcja ρ przypisująca te wartości, taka że:

$$\rho : E \times \mathbf{A} \longrightarrow D_{\mathbf{A}},$$

gdzie E jest zbiorem obiektów, $\mathbf{A} = \{A_1, A_2, \dots, A_n\}$ – zbiorem atrybutów wybranych do opisu własności tych obiektów, a $D_{\mathbf{A}}$ – łączną dziedziną wszystkich atrybutów, tj. $D_{\mathbf{A}} = \cup_{i=1}^n D_i$; ponadto, przyjmuje się, że dla każdego atrybutu A_i , $i = 1, 2, \dots, n$, $\rho(e, A_i) \in D_i$, gdzie $e \in E$ jest pewnym obiektem, $A_i \in \mathbf{A}$ – atrybutem, a D_i – dziedziną atrybutu A_i .

Zamiast pisać $\rho(e, A_i) = d$, gdzie $d \in D_i$ jest pewnym elementem dziedziny atrybutu A_i , stosowane są zapisy:

$$A_i(e) = d \tag{1}$$

$$\langle e, A_i, d \rangle \tag{2}$$

Oba powyższe zapisy mówią, że wartość atrybutu A_i dla obiektu e jest równa d . W przypadku gdy dokładnie wiadomo o jaki obiekt chodzi (np. zajmujemy się pojedynczym obiektem, powyższe zapisy mogą przybrać postać uproszczoną:

$$A_i = d \tag{3}$$

Powyższe wyrażenia określa się jako *formuły atomowe* (atomiczne), *atomy* lub *fakty*.

$$\langle A_i, d \rangle \tag{4}$$

$$[A_i, d] \tag{5}$$

Parę (A_i, d) nazywa się *deskryptorem* lub *selektorem*.

Logiczny model danych

Formuły atomowe – model ogólny

Niech:

- $E = \{e_1, e_2, \dots, e_m\}$ oznacza zbiór m obiektów reprezentowanych w modelu,
- $\mathbf{A} = \{A_1, A_2, \dots, A_n\}$ oznacza zbiór n atrybutów reprezentujących wybrane własności,
- D_1, D_2, \dots, D_n oznacza dziedziny powyższych atrybutów.

Ogólna postać atomu przybiera formę:

$$A_j(e_i) = d_{ij}, \quad (6)$$

kompleksu – formę:

$$\langle e_i, A_j, d_{ij} \rangle \quad (7)$$

a analogiczną do selektora formę:

$$[A_j(e_i) = d_{ij}] \quad (8)$$

W wyrażeniach (6), (7), (8) d_{ij} jest wartością j -tego atrybutu dla i -tego obiektu, $i = 1, 2, \dots, m$, $j = 1, 2, \dots, n$.

Powyższy sposób reprezentacji danych, a także w ogólniejszym przypadku – wiedzy, określa się także jako *O-A-V* (od ang. *Object-Attribute-Value*). Model ten może zostać zawężony do postaci *A-V*, zawierającej jedynie parę atrybut-wartość lub też rozszerzony o dodatkowe parametry. Typowym przykładem rozszerzenia jest uzupełnienie o tzw. *współczynnik wiarygodności* (ang. *Certainty Factor*); przybiera on wówczas postać *O-A-V-CF*.

Logiczny model danych – reprezentacja ekstensjonalna

Przy pomocy formuł atomowych można konstruować bardziej złożone wyrażenia reprezentujące określone dane. W szczególności mogą to być formuły opisujące wartości wszystkich atrybutów z zadanego zbioru dla określonego obiektu. Formuły takie odpowiadają logicznie *koniunkcji* atomów, z których każdy opisuje wartość pojedynczego atrybutu.

Definicja 1 *Formułę postaci*

$$\phi_i = [A_1(e_i) = d_{i,1}] \wedge [A_2(e_i) = d_{i,2}] \wedge \dots \wedge [A_n(e_i) = d_{i,n}] \quad (9)$$

nazywamy (pełną) formułą (prostą) opisującą własności obiektu e_i .

Określenie *formuła pełna* odnosi się do specyfikacji wartości wszystkich atrybutów. Jeżeli w danej formule nie wszystkie wartości atrybutów są sprecyzowane, to formuła taka nie jest formułą pełną. Określenie *formuła prosta* odnosi się do formuł o postaci koniunkcji atomów. Jeżeli wiadomo jaki obiekt jest przedmiotem opisu, specyfikację tego obiektu można pominąć; formuła (9) przyjmie wówczas postać $\psi = [A_1 = d_1] \wedge [A_2 = d_2] \wedge \dots \wedge [A_n = d_n]$.

Jeżeli danych jest m obiektów e_1, e_2, \dots, e_m opisywanych tymi samymi atrybutami A_1, A_2, \dots, A_n , to łączny opis **zbioru** tych obiektów (formuła opisująca relację) może przyjąć postać formuły opisującej własności wszystkich obiektów e_1, e_2, \dots, e_m łącznie,

$$\Phi = \phi_1 \wedge \phi_2 \wedge \dots \wedge \phi_m. \quad (10)$$

Zapis (10) jest zapisem logicznym; z uwagi na obszerność tej formuły oraz jej strukturalny charakter wygodnie jest prezentować taką formułę w postaci tabelarycznej, w odpowiednio skonstruowanej tabeli, gdzie poszczególne wiersze zawierają pełne opisy obiektów a kolumny odpowiadają wartościom danego atrybutu dla wszystkich obiektów.

Logiczny model danych – reprezentacja intensjonalna

Przy pomocy formuł atomowych nie zawierających specyfikacji obiektów można konstruować bardziej złożone wyrażenia reprezentujące określone dane *intensjonalnie*. W szczególności mogą to być formuły specyfikujące wartości (wszystkich) atrybutów z zadanego zbioru. Formuły takie odpowiadają logicznie *koniunkcji* selektorów, z których każdy specyfikuje wartość pojedynczego atrybutu.

Definicja 2 *Formułę postaci*

$$\psi_i = [A_1 = d_{i,1}] \wedge [A_2 = d_{i,2}] \wedge \dots \wedge [A_n = d_{i,n}] \quad (11)$$

nazywamy (pełną) formułą (prostą) specyfikującą własności pewnego obiektu.

Określenie *formuła pełna* odnosi się do specyfikacji wartości wszystkich atrybutów. Jeżeli w danej formule nie wszystkie wartości atrybutów są sprecyzowane, to formuła taka nie jest formułą pełną. Określenie *formuła prosta* odnosi się do formuł o postaci koniunkcji selektorów. Obiektów takich może być więcej niż jeden (obiekty nierozróżnialne) lub może nie istnieć obiekt o wybranych własnościach.

Jeżeli konstruowany jest opis większej liczby obiektów (np. m obiektów e_1, e_2, \dots, e_m opisywanych tymi samymi atrybutami A_1, A_2, \dots, A_n), to łączny opis **zbioru** tych obiektów (formuła opisująca relację) może przyjąć postać formuły opisującej własności wszystkich obiektów e_1, e_2, \dots, e_m łącznie,

$$\Psi = \psi_1 \vee \psi_2 \vee \dots \vee \psi_m. \quad (12)$$

Zapis (12) jest zapisem logicznym; z uwagi na obszerność tej formuły oraz jej strukturalny charakter wygodnie jest prezentować taką formułę w postaci tabelarycznej, w odpowiednio skonstruowanej tabeli, gdzie poszczególne wiersze zawierają pełne opisy obiektów a kolumny odpowiadają wartościom danego atrybutu dla wszystkich obiektów.

Logiczny model danych

Postać tabelaryczna

Logiczną formułę 10 opisującą informację w bazie danych można przedstawić w postaci tabelarycznej:

$$\begin{array}{cccc}
 [A_1(e_1) = d_{1,1}] & [A_2(e_1) = d_{1,2}] & \dots & [A_n(e_1) = d_{1,n}] \\
 [A_1(e_2) = d_{2,1}] & [A_2(e_2) = d_{2,2}] & \dots & [A_n(e_2) = d_{2,n}] \\
 \cdot & \cdot & \dots & \cdot \\
 \cdot & \cdot & \dots & \cdot \\
 \cdot & \cdot & \dots & \cdot \\
 [A_1(e_m) = d_{m,1}] & [A_2(e_m) = d_{m,2}] & \dots & [A_n(e_m) = d_{m,n}],
 \end{array} \tag{13}$$

W schemacie podstawowym przyjmuje się, że:

- obiekty są jednorodne, tzn. wszystkie atrybuty są stosowalne do opisu każdego obiektu,
- wszystkie wartości dowolnego atrybutu mogą być zastosowane do opisu każdego obiektu (choć oczywiście nie muszą być, w danej chwili, prawdziwe),
- żaden obiekt nie wymaga charakteryzowania poprzez dodatkowe atrybuty.

Założenia te mogą być w praktyce osłabiane.

Reprezentacja danych w modelu relacyjnym: tabele

W relacyjnych bazach danych informacja zapisywana jest w tabelach. Pierwowzór takiej tabeli może być postrzegany jako odpowiednio zaaranżowana przestrzennie formuła (13).

Interpretacja informacji w tablicy: każdy wiersz (krotka) odpowiada specyfikacji wartości wszystkich atrybutów dla danego obiektu, poszczególne wiersze zawierają opisy kolejnych obiektów, a w dowolnej kolumnie umieszczone są wartości pojedynczego atrybutu dla kolejnych obiektów. Niech:

- A_1, A_2, \dots, A_n będą wybranymi atrybutami o dziedzinach odpowiednio
- D_1, D_2, \dots, D_n ; ponadto niech
- $E = \{e_1, e_2, \dots, e_m\}$ będzie zbiorem pewnych obiektów opisywanych wskazanymi wyżej atrybutami. Reprezentacja tablicowa opisu wszystkich obiektów zbioru E może mieć następującą postać:

$$\begin{array}{c|cccccc}
 E & A_1 & A_2 & \dots & A_j & \dots & A_n \\
 \hline
 e_1 & d_{1,1} & d_{1,2} & \dots & d_{1,j} & \dots & d_{1,n} \\
 e_2 & d_{2,1} & d_{2,2} & \dots & d_{2,j} & \dots & d_{2,n} \\
 \vdots & \vdots & \vdots & & \vdots & & \vdots \\
 e_i & d_{i,1} & d_{i,2} & \dots & d_{i,j} & \dots & d_{i,n} \\
 \vdots & \vdots & \vdots & & \vdots & & \vdots \\
 e_m & d_{m,1} & d_{m,2} & \dots & d_{m,j} & \dots & d_{m,n}
 \end{array} \tag{14}$$

Powyższy sposób reprezentacji jest oszczędniejszy niż poprzedni: zarówno nazwa każdego atrybutu jak i obiektu jest wypisana tylko jeden raz. Równocześnie, reprezentacja taka wydaje się być bardziej przejrzysta i naturalna. Ma ona jednak *model logiczny* (semantykę) i może być interpretowana i przetwarzana logicznie. Może też być oceniana wartość logiczna tablicy jako formuły i jej elementów składowych (wierszy).

Reprezentacja danych w modelu relacyjnym: relacje

Z algebraicznego punktu widzenia, tabela przedstawiająca dane reprezentuje pewną relację R_E określoną w iloczynie kartezjańskim $E \times D_1 \times D_2 \times \dots \times D_n$:

$$R_E \subseteq E \times D_1 \times D_2 \times \dots \times D_n.$$

R_E jest relacją $n + 1$ argumentową, przy czym poszczególne argumenty krotek relacji odpowiadają elementom zbioru E oraz wartościom kolejnych atrybutów.

W rzeczywistości, obiekty e_1, e_2, \dots, e_n istnieją w świecie rzeczywistym; w modelu są one identyfikowane poprzez pewien wyróżniony atrybut (atrybuty), będący nazwą własną danego obiektu, jego identyfikatorem symbolicznym, numerem kolejnym, etc. Powinien to być identyfikator *jednoznaczny*.

Tablicowa reprezentacja relacji przyjmuje ostatecznie jednorodną postać

A_1	A_2	\dots	A_j	\dots	A_n
$d_{1,1}$	$d_{1,2}$	\dots	$d_{1,j}$	\dots	$d_{1,n}$
$d_{2,1}$	$d_{2,2}$	\dots	$d_{2,j}$	\dots	$d_{2,n}$
\vdots	\vdots		\vdots		\vdots
$d_{i,1}$	$d_{i,2}$	\dots	$d_{i,j}$	\dots	$d_{i,n}$
\vdots	\vdots		\vdots		\vdots
$d_{m,1}$	$d_{m,2}$	\dots	$d_{m,j}$	\dots	$d_{m,n}$

(15)

Powyższa tablica reprezentuje pewną relację R , gdzie:

$$R \subseteq D_1 \times D_2 \times \dots \times D_n. \quad (16)$$

W tabeli (15) (relacja R jest zadana *ekstensjonalnie*. Jeżeli opis każdej encji jest unikalny (zapewniona jest jednoznaczna identyfikacja), to $R = \pi_{2,3,\dots,n+1}(R_E)$

Schemat relacji

Zgodnie z matematyczną definicją relacja R jest pewnym zbiorem n -krotek postaci (d_1, d_2, \dots, d_n) , tzn.

$$R \subseteq \{(d_1, d_2, \dots, d_n) : d_1 \in D_1, d_2 \in D_2, \dots, d_n \in D_n\}.$$

W powyższym zapisie zastosowane do opisu atrybuty nie są uwidocznione. Aby móc odczytać znaczenie poszczególnych elementów składowych każdej krotki należy znać nie tylko zbiór wybranych atrybutów, ale także ich przypisanie poszczególnym składowym krotek (tzn. kolejność atrybutów). W tym celu definiuje się pojęcie *schematu relacji*.

Definicja 3 *Schematem relacji R o danych atrybutach A_1, A_2, \dots, A_n , takiej że $R \subseteq D_1 \times D_2 \times \dots \times D_n$ nazywamy ciąg (A_1, A_2, \dots, A_n) . W celu jawnej specyfikacji schematu relacji R piszemy $R(A_1, A_2, \dots, A_n)$.*

Wyrażenie $R(A_1, A_2, \dots, A_n)$ czytamy „relacja R o schemacie A_1, A_2, \dots, A_n ” lub też „relacja R ma schemat A_1, A_2, \dots, A_n ”. Znajomość schematu relacji pozwala dokonać właściwej interpretacji jej elementów, odczytać znaczenie poszczególnych argumentów krotek tworzących tą relację. W celu jawnej specyfikacji dziedzin atrybutów danej relacji jej schemat zapisuje się też jako $R(A_1 : D_1, A_2 : D_2, \dots, A_n : D_n)$.

Uwaga:

Wielu Autorów podręczników z zakresu baz danych definiuje schemat relacji jako *zbiór* atrybutów $\{A_1, A_2, \dots, A_n\}$; podejście takie wydaje się niewystarczające dla właściwej interpretacji relacji reprezentowanej jako zbiór krotek – dla dowolnej relacji n -argumentowej istniałoby wówczas $n!$ możliwych interpretacji jej argumentów.

Własności relacji

Tablice reprezentujące relacje posiadają następujące własności:

- liczba kolumn jest z góry ustalona (na etapie projektowania relacyjnej bazy danych) i jest taka sama dla wszystkich wierszy (krotek),
- kolumny etykietowane są atrybutami; z każdą kolumną związany jest dokładnie jeden atrybut (o określonej dziedzinie),
- liczba wierszy zależy od liczby opisywanych obiektów i jest dowolna (można nawet zdefiniować pustą relację o zadanym schemacie),
- kolejność kolumn nie jest w zasadzie istotna, o ile przypisze się im odpowiednie atrybuty (dotyczy to tabeli; w relacji będącej zbiorem krotek, których elementy nie są etykietowane atrybutami, nie wolno zmieniać kolejności argumentów),
- kolejność wierszy nie jest istotna; relacja jest *zbiorem* krotek,
- na przecięciu wierszy i kolumn znajdują się wartości atomiczne będące elementami dziedziny właściwej dla atrybutu przypisanego danej kolumnie,
- wszystkie wiersze tabeli reprezentującej relację są różne (różnią się przynajmniej na jednej pozycji – dwa identyczne wiersze nie mogą się powtarzać w tabeli, o ile jest ona zapisem relacji); w rzeczywistej tabeli mogą wystąpić identyczne wiersze, tzw. *duplikaty* (o ile nie określono klucza ani indeksu jednoznacznego).

Wymaganie dotyczące rozróżnialności rekordów jest konsekwencją definicji relacji jako zbioru; w tabeli stanowiącej reprezentację relacji żadna krotka nie może wystąpić dwa lub więcej razy. W konsekwencji, żaden wiersz nie powinien się powtórzyć.

Problem duplikatów

W realnych bazach danych gdzie tabele zapisywane są w pewnych plikach, nie ma w zasadzie przeszkód aby pewne wiersze powtarzały się. Jeżeli przy definiowaniu schematu relacji (atrybutów i formatu tabeli) nie zdefiniuje się żadnego *klucza*, ani nie zażąda się aby rekordy były *indeksowane unikatowo*, to wprowadzenie dwóch lub więcej identycznych wierszy do tabeli jest możliwe, a nawet czasem celowe; co więcej, tabele zawierające duplikaty mogą się pojawić w wyniku realizacji pewnych operacji algebraicznych na istniejących już tabelach bez powtórzeń, np. w wyniku zastosowania operacji projekcji. Powtarzające się wiersze tabeli nazywa się *duplikatami*. Powoduje to jednak określone konsekwencje, które należy uwzględnić już na etapie projektu.

- obiekty opisywane takimi wierszami nie są rozróżnialne, a więc celowość reprezentowania więcej niż jednego takiego obiektu jest wątpliwa,
- usuwając informację o takim obiekcie z bazy należy wyraźnie sprecyzować, czy chodzi nam o usunięcie jednego jej wystąpienia, dwóch, trzech, itd., czy też wszystkich,
- mogą pojawić się trudności związane z przechowywaniem takich tabel, jeżeli dostęp do ich elementów (wierszy) realizowany jest za pomocą indeksu unikatowego,
- w wyniku realizacji pewnych operacji na takich tabelach (np. złożenia lub iloczynu kartezjańskiego) następuje niekontrolowane i nie zawsze celowe powielanie informacji,
- niektóre operacje na tabelach (np. łączenie) wymagają aby tzw. tabela główna (tabela nadrzędna) była indeksowana jednoznacznie względem pól łączących (lub aby pola te definiowały klucz) i bez spełnienia tego warunku operacje te nie mogą być realizowane; dlatego w większości przypadków, już na etapie projektowania tabel bazy danych, definiuje się dla każdej tablicy klucz podstawowy.

Problem jednoznacznej identyfikacji encji

Jeżeli żaden z atrybutów występujących w schemacie tablicy nie pozwala na jednoznaczną identyfikację obiektów, to zwykle wprowadza się (sztucznie) nowy atrybut jednoznacznie identyfikujący każdy obiekt (identyfikator).

Przykładami takich identyfikatorów mogą być: numer PESEL, numer NIP, identyfikator pracownika w zakładzie pracy, sygnatura książki w bibliotece, numer rejestracyjny samochodu, itp. Wszystkie one mają pewną istotną własność, a mianowicie w sposób *jednoznaczny* identyfikują dany obiekt. Formalnie, jeżeli A jest takim wybranym atrybutem o dziedzinie D_A , to musi istnieć wzajemnie jednoznaczna *funkcja identyfikująca* μ , postaci:

$$\mu : E \longrightarrow D_A. \quad (17)$$

Warunek (17) oznacza, że każdy obiekt ze zbioru E może być reprezentowany przez pewien identyfikator, przy czym żądanie wzajemnej jednoznaczności funkcji identyfikującej czyni zadość wymaganiu *rozróżnialności* obiektów, tzn. mając daną wartość funkcji jesteśmy w stanie jednoznacznie zidentyfikować obiekt. Oczywiście nie wszystkim potencjalnym identyfikatorom muszą być w danej chwili przypisane relacje istniejące obiekty; zazwyczaj definicja dziedziny atrybutu identyfikującego jest tak skonstruowana, aby zawsze można było opisać nowe obiekty, jeżeli takie pojawią się i zaistnieje potrzeba ich reprezentacji w bazie danych.

Jeżeli w danym zestawie atrybutów nie istnieje (nie jest wprowadzony) żaden pojedynczy atrybut czyniący zadość wymaganiu jednoznacznej identyfikacji, można wybrać odpowiedni zestaw atrybutów $A^1, A^2, \dots, A^k \in \{A_1, A_2, \dots, A_n\}$ spełniających żądanie jednoznacznej identyfikacji. W takim przypadku funkcja identyfikująca przybiera postać:

$$\mu : E \longrightarrow D^1 \times D^2 \times \dots \times D^k, \quad (18)$$

gdzie, D^1, D^2, \dots, D^k są odpowiednio dziedzinami atrybutów A^1, A^2, \dots, A^k wchodzących w skład identyfikatora.

Pojęcie klucza

Pojęcie klucza i definicje pochodne

- **Zbiór identyfikujący:** dowolny atrybut lub zestaw atrybutów pozwalających na jednoznaczną identyfikację obiektu w bazie danych.
- **Klucz:** każdy minimalny zbiór identyfikujący, tzn. taki, że żaden jego właściwy podzbiór nie jest wystarczający do jednoznacznej identyfikacji obiektu; klucz stanowi więc możliwie najprostszy zestaw atrybutów wystarczający do rozróżnienia wszystkich obiektów,
- **Klucz prosty:** klucz złożony z pojedynczego atrybutu,
- **Klucz złożony:** klucz składający się z więcej niż jednego atrybutu,
- **Klucz podstawowy:** wybrany klucz preferowany przez użytkownika; nazywa się go też lub *kluczem głównym*,
- **Nadklucz:** każdy zbiór atrybutów zawierający w sobie pewien klucz,
- **Podklucz:** dowolny właściwy podzbiór klucza.
- **Klucz obcy:** atrybut lub ich zestaw występujący w danej tabeli a będący kluczem określony w innej tabeli, stosowany w celu identyfikacji elementów dla ewentualnego złączenia; klucz obcy nie musi być kluczem w tabeli dołączanej,
- **Atrybut kluczowy:** atrybut wchodzący w skład przynajmniej jednego klucza.

Klucz jest określany dla schematu relacji; co jest kluczem a co nie nie zależy od aktualnej zawartości tabeli. Pojęcie klucza jest związane z istnieniem *zależności funkcyjnych* w tabelach. Dana tabela może mieć wiele kluczy. Wyznaczenie wszystkich kluczy stanowi istotne zadanie przy normalizacji.

Relacje, tabele logiczne, tabele fizyczne

Relacje a tabele logiczne

Tabela logiczna (wirtualna) stanowi (wyidealizowaną) reprezentację relacji; może ona jednak różnić się od relacji definiowanej matematycznie:

- Dziedzina relacji (ze względu na i -ty argument; zadawana ekstensjonalnie) reprezentowanej w tabeli jest zazwyczaj podzbiorem dziedziny D_i i -tego atrybutu (zadawanej intensjonalnie),
- Dla relacji zawsze istnieje klucz; dla tabeli – niekoniecznie,
- Tabela może zawierać duplikaty (nie dopuszczalne w relacji),
- W tabeli można zmieniać kolejność kolumn (przy zachowaniu etykietowania atrybutami), w relacji – nie,
- Istnieje dwoista terminologia (np. krotka – rekord, relacja – tabela).

Tabele logiczne a tabele fizyczne

Tabela fizyczna (pamiętana) stanowi zapis tabeli logicznej; może się od niej jednak istotnie różnić:

- rekord logiczny może być różny od rekordu fizycznego (przetwarzanie danych, pola wyliczane, pola segmentowe),
- dane mogą być kodowane,
- dane często mają wymiar (pomijany w tabeli fizycznej),
- dla każdego zagadnienia można zbudować wiele reprezentacji fizycznych (różna liczba i schematy tabel).

Problem wartości *NULL*

Wartości *NULL* to tzw. wartości puste (albo zerowe); reprezentują one brak wartości danego atrybutu (chwilowy lub permanentny), który może być interpretowany jako:

- brak znajomości wartości danego atrybutu (dane niedostępne, wartość nieznana; np. data urodzenia, numer NIP czy PESEL),
- niemożliwość wyznaczenia wartości danego atrybutu w rekordzie (atrybut niestosowalny; np. pojemność skokowa dla samochodów elektrycznych),
- niepewność co do istnienia tej wartości (np. gdy nie wiadomo, czy osoba posiada telefon).

Wartości *NULL* są różne od spacji, zera czy też stringu pustego. W pewnych okolicznościach, wartości *NULL* są traktowane jak każde inne (np. wyświetlanie, porządkowanie); w innych okolicznościach nie podlegają przetwarzaniu (dają nieokreślone wartości funkcji, dwie wartości *NULL* nie są traktowane jako równe – brak połączenia). Porównanie wartości *NULL* z dowolną wartością daje wartość logiczną *UNKNOWN* lub *NULL* (różną od *TRUE* i *FALSE*). Prowadzi to do logiki trówartościowej:

AND	<i>TRUE</i>	<i>FALSE</i>	<i>NULL</i>	OR	<i>TRUE</i>	<i>FALSE</i>	<i>NULL</i>
<i>TRUE</i>	<i>TRUE</i>	<i>FALSE</i>	<i>NULL</i>	<i>TRUE</i>	<i>TRUE</i>	<i>TRUE</i>	<i>TRUE</i>
<i>FALSE</i>	<i>FALSE</i>	<i>FALSE</i>	<i>FALSE</i>	<i>FALSE</i>	<i>TRUE</i>	<i>FALSE</i>	<i>NULL</i>
<i>NULL</i>	<i>NULL</i>	<i>FALSE</i>	<i>NULL</i>	<i>NULL</i>	<i>TRUE</i>	<i>NULL</i>	<i>NULL</i>

NOT	<i>TRUE</i>	<i>FALSE</i>	<i>NULL</i>
<i>TRUE</i>	<i>FALSE</i>	<i>TRUE</i>	<i>NULL</i>

Operatory *IS NULL* oraz *IS NOT NULL* pozwalają sprawdzić, czy dana wartość jest *NULL*. Użycie wartości *NULL* może prowadzić do niezamierzonych efektów (np. $0 \cdot x = 0$ dla dowolnej liczby x , ale $0 \cdot \text{NULL} = \text{NULL}$).

Sortowanie i indeksowanie

Sortowanie

Operacja *sortowania* oznacza uporządkowanie rekordów w tabeli względem jednego lub kilku kryteriów (ustalonych hierarchicznie). Najczęściej ma miejsce *uporządkowanie leksykograficzne* (alfabetyczne). **Porządek leksykograficzny:**

$$X_1X_2 \dots X_k < Y_1Y_2 \dots Y_m$$

wtw. gdy

- $k < m$ oraz $X_1X_2 \dots X_k = Y_1Y_2 \dots Y_k$, albo
- dla pewnego $i \leq \min(k, m)$ $X_1X_2 \dots X_{i-1} = Y_1Y_2 \dots Y_{i-1}$ oraz $X_i < Y_i$.

Sortowanie może odbywać się według jednego lub kilku pól; *ASC* oznacza porządek rosnący (default), *DESC* – porządek malejący. Uwaga: porządkowanie liczb zadeklarowanych jako tekst daje niezamierzony efekt!

Indeksowanie

Indeksy to pomocnicze struktury danych pozwalające na szybkie odszukiwanie rekordów o zadanych wartościach atrybutów. Indeksy zazwyczaj mają strukturę hierarchiczną (blokową) i są realizowane za pomocą tzw. B-drzew. Użycie indeksów znakomicie przyspiesza wyszukiwanie informacji oraz operacje na danych, ale ich obsługa (aktualizacja) jest pracochłonna.

Indeksy są implementowane tak, że każdy wierzchołek B-drzewa zajmuje blok na dysku (4096 bajtów), a więc może zawierać kilkaset wskaźników do węzłów potomnych (dla typowego wyszukania wymagane jest np. tylko 3 dostępy do dysku).

Indeksy mogą być pojedyncze (oparte na jednym atrybucie) oraz wielokrotne (dla kilku atrybutów). Można je deklarować jako unikatowe (bez powtórzeń) lub jako dopuszczające powtórzenia.

Problemy modelowania danych

Dane zawarte w bazie danych powinny wiernie oddawać stan rzeczywisty modelowanego systemu, tzn. powinny one być zgodne z rzeczywistością. W praktyce modelowanie zbiorów encji w relacyjnych bazach danych napotyka szereg problemów teoretycznych i praktycznych. Najważniejsze z nich to:

1. Problem zupełności.

- (a) *Zupełność fizyczna*: czy wszystkie encje modelowaneo świata zostały opisane w bazie (czy nic nie zostało pominięte)?
- (b) *Zupełność logiczna*: czy wyspecyfikowano wszystkie rekordy pokrywające zadany (wymagany) obszar uniwersum U ?

2. Problem rozróżnialności.

- (a) *Rozróżnienie fizyczne*: czy wszystkie różne encje reprezentowane są różnymi rekordami?
- (b) *Rozróżnialność potencjalna*: czy wybrany schemat relacji (atrybuty) pozwala rozróżnić wszystkie encje na zadanym poziomie abstrakcji (wyróżnić wszystkie wymagane klasy abstrakcji)?

3. Problem poprawności (fizycznej realizowalności).

- (a) *Realizowalność fizyczna (poprawność potencjalna)*: czy wszystkie rekordy opisują fizycznie realizowalne (potencjalnie) encje?
- (b) *Poprawność*: czy wszystkie rekordy opisują realnie istniejące encje?

4. Spójność danych.

- (a) *Spójność zewnętrzna (administracyjna)*: czy aktualny stan bazy danych dokładnie odzwierciedla aktualny stan systemu ?
- (b) *Spójność wewnętrzna*: czy spełnione są wszystkie ograniczenia, tzw. *więzy spójności*, nałożone na schemat i instancję bazy danych?

Więzy integralności

Dla zapewnienia (względnej) poprawności danych konstruuje się różnego rodzaju *więzy integralności* i nakłada je na strukturę (schemat) oraz realizację (instancję) bazy danych. Typowe więzy spójności obejmują:

Lokalne więzy spójności:

- definiowanie dziedzin atrybutów (typu, długości),
- nakładanie ograniczeń na wartości atrybutu w rekordzie (reguły poprawności operujące na pojedynczym polu, np. `Cena >= 0`),
- nakładanie łącznych ograniczeń na wartości atrybutów w rekordzie (reguły poprawności operujące na polach każdego pojedynczego rekordu, np. `Cena_kupna < Cena_sprzedaży`),
- definiowanie masek wprowadzania oraz formatów wyprowadzania danych.

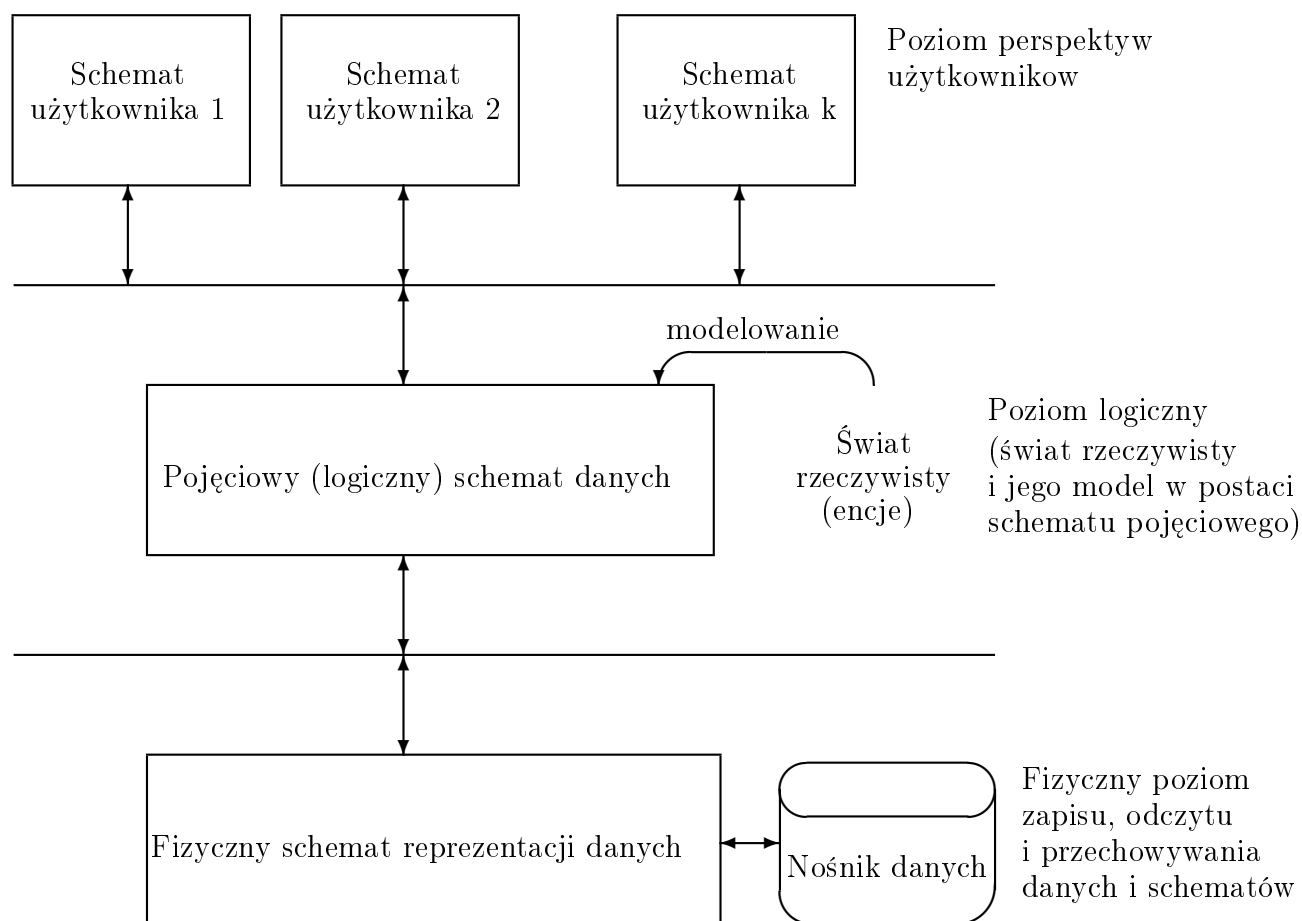
Globalne więzy spójności:

- niedopuszczenie powtarzania się wartości wybranego pola (lub grupy pól) poprzez definiowanie klucza lub indeksu bez powtórzeń,
- definiowanie więzów referencyjnych (wymuszanie więzów integralności); wartości klucza obcego muszą występować w tablicy nadrzędnej (lub być *NULL*), wymagana jest zgodność typów przy złączeniach,
- globalne ograniczenia funkcyjne (weryfikowane okresowo) (np. `Stan_magazynu = Stan_poprzedni - Suma_sprzedaży + Suma_dostaw`),
- globalne ograniczenia numeryczne (np. liczby wybranych rekordów), algebraiczne (np. zawieranie się wybranych zbiorów) i logiczne (spełnianie wybranych formuł definiujących pożądane własności).

Podstawowe założenia RBD – Postulaty Codda

1. **Postulat informacyjny.** Dane są reprezentowane *jedynie* przez wartości atrybutów w wierszach tabel.
2. **Postulat dostępu.** Każda wartość jest dostępna poprzez podanie tabeli, atrybutu i klucza.
3. **Postulat dotyczący wartości NULL.** Dostępna jest specjalna wartość NULL dla reprezentacji wartości nieokreślonej, inna od wszystkich, i podlegająca przetwarzaniu.
4. **Postulat dotyczący katalogu.** Struktura bazy danych jest dostępna w katalogu będącym relacyjną bazą danych.
5. **Postulat języka danych.** System musi dostarczać pełnego języka przetwarzania danych (interakcja, aplikacje, definicje, przetwarzanie).
6. **Postulat modyfikowalności perspektyw.** System musi umożliwiać modyfikowanie perspektyw, o ile jest ono semantycznie realizowalne.
7. **Postulat modyfikowalności danych.** System musi umożliwiać operacje modyfikacji danych (INSERT, UPDATE, DELETE).
8. **Postulat fizycznej niezależności danych.** Zmiany fizycznej reprezentacji danych i organizacji dostępu nie wpływają na aplikacje.
9. **Postulat logicznej niezależności danych.** Zmiany wartości w tabelach nie wpływają na aplikacje.
10. **Postulat niezależności więzów spójności.** Więzy spójności są definiowalne w bazie i nie zależą od aplikacji.
11. **Postulat niezależności dystrybucyjnej.** Działanie aplikacji nie zależy od modyfikacji dystrybucji bazy.
12. **Postulat bezpieczeństwa względem operacji niskiego poziomu.** Operacje niskiego poziomu (na poziomie rekordu) nie mogą naruszać modelu relacyjnego i więzów spójności.

Poziomy reprezentacji danych





KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI



UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY



Studia podyplomowe "Inżynieria oprogramowania" współfinansowane przez Unię Europejską w ramach Europejskiego Funduszu Społecznego

Projekt "Studia podyplomowe z zakresu wytwarzania
oprogramowania oraz zarządzania projektami w firmach
informatycznych" realizowany w ramach
Programu Operacyjnego Kapitał Ludzki