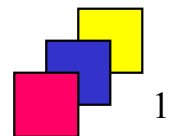


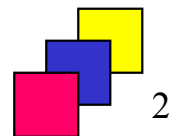
Optymalizacja poleceń SQL

Statystyki



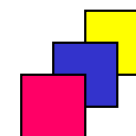
Statystyki (1)

- Informacje, opisujące dane i struktury obiektów bazy danych.
- Przechowywane w słowniku danych.
- Używane przez optymalizator do oszacowania:
 - selektywności predykatów polecenia,
 - kosztu użycia ścieżek dostępu,
 - kosztu operacji I/O i czasu procesora,
 - **kosztu planu wykonania polecenia.**
- Tylko **aktualne** statystyki użyteczne!
 - Statystyki są **statyczne** – nie są automatycznie uaktualniane przy zmianie danych.



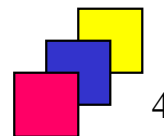
Statystyki (2)

- **Przykłady statystyk:**
 - dla relacji:
 - liczba rekordów,
 - liczba bloków,
 - średnia długość rekordu,
 - dla atrybutu relacji:
 - liczba różnych wartości,
 - liczba rekordów, w których atrybut ma wartość pustą,
 - rozkład wartości (histogram),
 - dla indeksu:
 - liczba bloków-liści,
 - wysokość drzewa,
 - wskaźnik zgrupowania indeksu,
 - statystyki systemowe:
 - wykorzystanie procesora,
 - liczba operacji we/wy.



Statystyki (3)

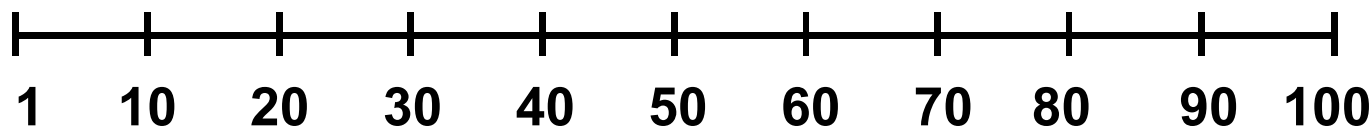
- **Statystyki mogą być gromadzone automatycznie (przez dedykowany proces SZBD) lub ręcznie (na żądanie użytkownika) przy użyciu pakietu DBMS_STATS.**
- **W przypadku braku statystyk dla obiektów używanych w zapytaniu przed wykonaniem zapytania optymalizator realizuje dynamiczne próbkowanie statystyk.**



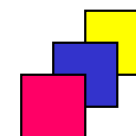
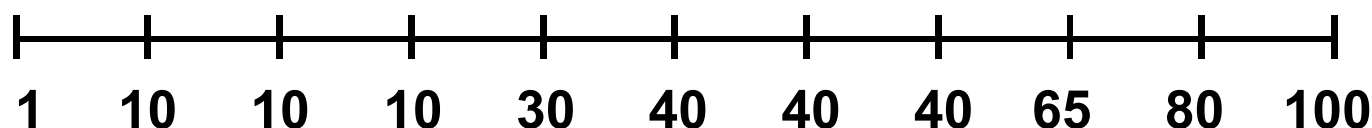
Histogramy (1)

- **Histogram** – szczegółowa statystyka opisująca rozkład wartości określonej kolumny relacji.
- **Rodzaje:**
 - histogram o zrównoważonej wysokości (ang. height balanced) – zbiór wartości kolumny dzielony jest na przedziały o tej samej (w przybliżeniu) liczbie rekordów; przykład (zakres wartości: $\langle 1, 100 \rangle$, liczba przedziałów: 10):

- **równomierny rozkład wartości atrybutu:**

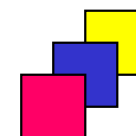


- **nierównomierny rozkład wartości atrybutu:**



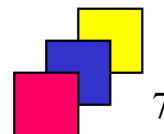
Histogramy (2)

- **Rodzaje (cd):**
 - **histogram częstotliwości (ang. frequency) – każda wartość kolumny odpowiada jednemu przedziałowi, każdy przedział zawiera liczbę wystąpień tej wartości; tworzony wtedy, gdy liczba wartości kolumny jest mniejsza bądź równa żądanej liczbie przedziałów histogramu.**
- **Histogramy należy tworzyć tylko dla kolumn z nierównomiernym rozkładem wartości (ang. skewed data), często używanych w warunkach zapytania.**
- **Gdy zmieni się rozkład danych kolumny, konieczne jest ponowne wygenerowanie histogramu,**



Ręczne zbieranie statystyk

- **Metody:**
 - na podstawie pełnych danych,
 - szacowanie na podstawie próbki, próbka określana w procentach liczby rekordów.
- **Procedury zbierające statystyki:**
 - `DBMS_STATS.GATHER_INDEX_STATS` – dla indeksu,
 - `DBMS_STATS.GATHER_TABLE_STATS` – dla relacji.
- **Procedury usuwające statystyki:**
 - `DBMS_STATS.DELETE_INDEX_STATS` – dla indeksu,
 - `DBMS_STATS.DELETE_TABLE_STATS` – dla relacji,
 - `DBMS_STATS.DELETE_COLUMN_STATS` – dla kolumny.



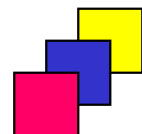
Zbieranie statystyk dla indeksu

```
exec DBMS_STATS.GATHER_INDEX_STATS(  
  ownname => <nazwa_schematu>, indname => <nazwa_indeksu>,  
  estimate_percent => <procentowa_wielkość_próbki>);
```

- jeśli wartość <procentowa_wielkość_próbki> określono jako:
 - null, wówczas statystyki zbierane na podstawie pełnych danych,
 - liczbę z przedziału <0,00001; 100>, wówczas szacowanie na podstawie próbki o zadanym rozmiarze,
 - DBMS_STATS.AUTO_SAMPLE_SIZE – rozmiar próbki dobiera system.

```
exec DBMS_STATS.GATHER_INDEX_STATS(  
  ownname => 'SCOTT', indname => 'PK_PRAC', estimate_percent => 20);
```

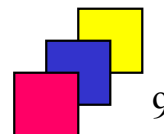
- **Uwaga!** Od Oracle10g statystyki dot. indeksów są gromadzone automatycznie podczas tworzenia lub przebudowy indeksu.



Zbieranie statystyk dla relacji (1)

```
exec DBMS_STATS.GATHER_TABLE_STATS(  
  ownname => <nazwa_schematu>, tablename => <nazwa_relacji>,  
  estimate_percent => <procentowa_wielkość_próbki>,  
  method_opt => <rodzaj_statystyk>,  
  cascade =><DBMS_STATS.AUTO_CASCADE | TRUE | FALSE> );
```

- **METHOD_OPT** – określa zakres zbieranych statystyk:
 - **FOR TABLE** – tylko statystyki dla tabeli bez statystyk dla kolumn,
 - **FOR ALL COLUMNS** [**<klauzula SIZE>**] – statystyki dla tabeli i statystyki dla wszystkich kolumn,
 - **FOR ALL INDEXED COLUMNS** [**<klauzula SIZE>**] – statystyki dla tabeli i statystyki dla poindeksowanych kolumn,
 - **FOR COLUMNS** [**<klauzula SIZE>**] kolumna1 [**<klauzula SIZE>**], kolumna2 [**<klauzula SIZE>**], ... – statystyki dla tabeli i statystyki dla wskazanych kolumn.



Zbieranie statystyk dla relacji (2)

- **<klauzula SIZE> – SIZE { liczba | REPEAT | AUTO | SKEWONLY }:**
 - liczba – liczba przedziałów w histogramie, zakres: <1, 254> ,
 - **REPEAT** – powtórzenie zbierania histogramów dla kolumn, które mają już histogramy,
 - **AUTO** – SZBD określi, dla których kolumn zbierać histogramy na podstawie obciążenia i rozkładu danych kolumny,
 - **SKEWONLY** – SZBD określi, dla których kolumn zbierać histogramy tylko na podstawie rozkładu danych kolumny (bez analizy obciążenia).
- **FOR ALL COLUMNS SIZE AUTO** – wartość domyślna dla par. **METHOD_OPT:**
 - statystyki tabeli,
 - podstawowe statystyki wszystkich kolumn tabeli,
 - histogramy dla kolumn wyznaczonych na podstawie wcześniejszych obserwacji dotyczących obciążenia i rozkładu wartości.

Zbieranie statystyk dla relacji (3)

```
exec DBMS_STATS.GATHER_TABLE_STATS(  
  ownname => 'SCOTT', tabname => 'PRACOWNICY',  
  estimate_percent => DBMS_STATS.AUTO_SAMPLE_SIZE,  
  method_opt => 'FOR COLUMNS placa_pod SIZE AUTO, nazwisko SIZE AUTO');
```

```
exec DBMS_STATS.GATHER_TABLE_STATS(  
  ownname => 'SCOTT', tabname => 'PRACOWNICY',  
  method_opt => 'FOR ALL INDEXED COLUMNS',  
  cascade => TRUE);
```

- Uwaga! Od Oracle12c statystyki dotyczące tabel zostają zebrane automatycznie w sytuacji, gdy tabela, do której ładowane są dane ścieżką bezpośrednią (polecenie INSERT /*+ APPEND */, dane umieszczane od razu w plikach bazy danych z pominięciem bufora bazy danych), była poprzednio pusta:
 - tabela została dopiero co utworzona i nie posiada jeszcze rekordów, lub
 - usunięto z tabeli wszystkie rekordy.

Statystyki w słowniku bazy danych

- Dla relacji:
 - **USER_TABLES, USER_TAB_STATISTICS**
- Dla kolumn:
 - **USER_TAB_COLUMNS, USER_TAB_COL_STATISTICS, USER_TAB_HISTOGRAMS**
- Dla indeksów:
 - **USER_INDEXES, USER_IND_STATISTICS**

```
SELECT num_rows, blocks, last_analyzed, sample_size  
FROM USER_TAB_STATISTICS  
WHERE table_name = 'PRACOWNICY';
```

```
SELECT num_distinct, low_value, high_value, num_buckets, histogram  
FROM USER_TAB_COL_STATISTICS  
WHERE table_name = 'PRACOWNICY'  
AND column_name = 'NAZWISKO';
```



Usuwanie statystyk

```
exec DBMS_STATS.DELETE_INDEX_STATS(  
  ownname => <nazwa_schematu>, indname => <nazwa_indeksu>);
```

```
exec DBMS_STATS.DELETE_TABLE_STATS(  
  ownname => <nazwa_schematu>, tabname => <nazwa_relacji>);
```

```
exec DBMS_STATS.DELETE_COLUMN_STATS(  
  ownname => <nazwa_schematu>, tabname => <nazwa_relacji>,  
  colname => <nazwa_kolumny>, col_stat_type => <rodzaj_usuwanych_statystyk>);
```

- **COL_STAT_TYPE:**
 - **HISTOGRAM** – usuwany jest histogram dla kolumny, podstawowe statystyki kolumny pozostają,
 - **ALL** – usuwane są wszystkie statystyki dla kolumny (wartość domyślna).

Wskaźnik zgrupowania indeksu (1)

- Minimalną jednostką operacji I/O jest blok dyskowy a nie rekord
- Statystyka, pozwalająca na porównanie kosztu operacji przeglądnięcia indeksu z kosztem pełnego przeglądnięcia tabeli
- Określa, jak mocno indeks jest "zsynchronizowany" z tabelą:
 - mała wartość – rekordy tabeli z tymi samymi (lub zbliżonymi) wartościami poindeksowanej kolumny są skupione w niewielkiej liczbie bloków
 - duża wartość – rekordy tabeli z tymi samymi (lub zbliżonymi) wartościami poindeksowanej kolumny są rozproszone w dużej liczbie bloków

Wskaźnik zgrupowania indeksu (2)

- Interpretacja:
 - mała wartość (równa lub bliska liczbie bloków tabeli) – dobrze, użycie indeksu jest korzystne w stosunku do pełnego przeglądnięcia tabeli z powodu konieczności wykonania mniejszej liczby operacji odczytu bloków tabeli (odczytu danych) po dostępie do indeksu (po odczycie adresów rekordów)
 - duża wartość (równa lub bliska liczbie rekordów tabeli) – źle, użycie indeksu jest niekorzystne w stosunku do pełnego przeglądnięcia tabeli z powodu konieczności wykonania większej liczby operacji odczytu bloków tabeli po dostępie do indeksu
- Słownik danych

```
SELECT clustering_factor FROM user_indexes  
WHERE index_name = 'PRAC_PK';
```

Wskaźnik zgrupowania indeksu (3)

- Przykład – tabela posiada 9 rekordów, poindeksowana kolumna K1 posiada trzy wartości A, B i C (po trzy rekordy), rekordy zajmują 3 bloki.
- Przypadek 1. Mała wartość wskaźnika. Niski koszt skanu indeksu – odczyt A wymaga dostępu do jednego bloku tabeli

BLOK 1			BLOK 2			BLOK 3		
A	A	A	B	B	B	C	C	C

- Przypadek 2. Duża wartość wskaźnika. Wyższy koszt skanu indeksu – odczyt A wymaga dostępu do wszystkich trzech bloków tabeli

BLOK 1			BLOK 2			BLOK 3		
A	B	C	A	B	C	A	B	C