

# Przetwarzanie masywnych danych

Ćwiczenia przed kolokwium

Imię i nazwisko: \_\_\_\_\_

Numer indeksu: \_\_\_\_\_

Data: **15 stycznia 2019**

W odpowiedziach na pytania należy być **bardzo zwięzłym** oraz **pisać starannie**. Za każde pytanie można dostać do 10 punktów.

1. Sieć hoteli postanowiła założyć hurtownię danych. Na początek postanowiono stworzyć hurtownię danych dla procesu analizy zajętości pokoi hotelowych przy zabranych następujących informacjach. W ramach sieci rozróżnia się kategorię hotelu (standardowy, motel, zajazd, itp.). O każdym hotelu przechowuje się takie informacje jak jego kategorię, adres, kraj, region, położenie (centrum miasta, poza miastem, itp.). Pokoje w każdym hotelu opisane są rozmiarem, typem (apartament, pojedynczy, podwójny, itp.), liczbą łóżek, maksymalną liczbą osób, wyposażeniem dodatkowym (lodówka, kuchenka, itp.). Każdy hotel ma określoną liczbę pokoi danego typu.

Hurtownia danych ma umożliwić podanie sumarycznej liczby zajętych pokoi w różnych okresach czasu, z podziałem na kategorię hotelu, typ pokoju, dodatkowe wyposażenie, itp..

Zaproponuj schemat gwiazdy, wskaż tablicę faktów, tablice wymiarów, zaznacz podkreśleniem klucze główne oraz linią falowaną klucze obce, zaznacz miary.

**Odp.:** (10 pkt.)

3. Zaproponuj algorytm wykorzystujący indeks segmentowy (ang. *bit-sliced index*) dla warunków zakresowych postaci  $l < A < r$ , gdzie  $l$  i  $r$  definiują lewy i prawy koniec zakresu dla wartości atrybutu  $A$ . Niech wszystkie liczby będą zapisane na  $n$  bitach (najbardziej znaczący bit to  $n - 1$ , najmniej znaczący bit to 0). Do  $i$ -tego bitu liczb  $l$  i  $r$  odwołuj się jak do tablic, tzn. poprzez  $l[i]$  i  $r[i]$ . Indeks segmentowy składa się z  $n$  bitmap  $B[i]$ , gdzie  $i = 0, \dots, n - 1$ . Wynikiem jest bitmapa  $R$ . Pomocnicze bitmapy oznacz dużymi literami.  **Odp.:** (10 pkt.)

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

4. W systemie składowana jest następująca perspektywa zmateriałizowana (ang. *materialized view*) oznaczona jako  $V$ :

```
SELECT marka, model, rok, sum(cena) as cena
FROM Sprzedaz, Samochod
WHERE Sprzedaz.id_samochod = Samochod.id_samochod
GROUP BY marka, model, rok;
```

Do systemu zostało wysłane poniższe zapytanie:

```
SELECT model, sum(cena)
FROM Sprzedaz, Samochod
WHERE Sprzedaz.id_samochod = Samochod.id_samochod
WHERE Sprzedaz.marka = 'Ford'
GROUP BY model;
```

Pokaż jak system może wykonać przepisanie zapytania (ang. *query re-write*) do równoważnego zapytania, które wykorzystuje powyższą perspektywę  $V$  (10 pkt.):

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

2. Operatory agregacji można poklasyfikować ze względu na sposób ich obliczania. Rozróżnia się operatory rozproszone, algebraiczne oraz holistyczne. Przyporządkuj następujące operatory do odpowiednich grup: count(), max(), mediana(), min(), rank(), sum(), ave(), std(): (10 pkt.)

- operatory rozproszone: \_\_\_\_\_
- operatory algebraiczne: \_\_\_\_\_
- operatory holistyczne: \_\_\_\_\_



10. W problemie szukania najbliższych sąsiadów względem odległości Hamminga w binarnej przestrzeni 100-wymiarowej wykorzystano funkcje mieszające, które losowo wybierają element wielowymiarowego wektora opisującego dany obiekt. Jaką charakterystykę  $(d_1, d_2, p_1, p_2)$  ma taka funkcja mieszająca dla odległości  $d_1 = 5$  i  $d_2 = 95$ ? Co się zmieni jeżeli zastosujemy 50 takich funkcji w następującej konfiguracji: funkcje zostały podzielone na 5 grup, każda zawiera 10 elementów. Wewnątrz grup zastosowano operację OR, natomiast pomiędzy grupami operację AND.

Podaj charakterystykę  $(d_1, d_2, p_1, p_2)$  pojedynczej funkcji mieszającej dla  $d_1 = 5$  i  $d_2 = 95$ : (5 pkt.)

---

---

---

Podaj charakterystykę  $(d_1, d_2, p_1, p_2)$  dla 50 funkcji mieszających w konfiguracji podanej powyżej (również dla  $d_1 = 5$  i  $d_2 = 95$ ): (5 pkt.)