

Big Data

Lista zadań

Jacek Cichoń, WPPT PWr, 2016/17

1 Wstęp

Zadanie 1 — Znajdź źródła swojej ulubionej książki. Zapisz je w formacie utf-8. W tm zadaniu zastosuj swój ulubiony język programowania (plik ma być takich rozmiarów aby w całości mieścił się w pamięci komputera).

1. Wczytaj książkę i podziel je na słowa. Usuń z tej listy stop-words (możesz ja znaleźć na stronie <https://pl.wikipedia.org/wiki/Wikipedia:Stopwords>
2. Wyznacz częstotliwości występowania wszystkich słów. Masz zbudować listę postaci $[[\text{słowo} \Rightarrow \text{liczba}], [\text{słowo} \Rightarrow \text{liczba}], \dots]$.
3. Posortuj otrzymaną listę według drugiego parametru.
4. Wyświetl kilkadziesiąt pierwszych elementów. Usuń z niej kilkanaście początkowych elementów i zapisz listę do pliku tekstowego.
5. Zbuduj chmurę wyrazów (word cloud) z otrzymanej listy. Możesz skorzystać np. z serwisu <http://www.wordclouds.com/>

Celem tego zadania jest wygenerowanie mniej więcej takiego obrazka (dla książki "Pan Tadeusz"):



Zadanie 2 — To jest kontynuacja poprzedniego zadania.

1. Podziel swoją książkę na rozdziały.
2. Każdy rozdział potraktuj jako dokumenty.
3. Podziel dokumenty na słowa. Wyznacz indeksy TF.IDF wszystkich słów we wszystkich rozdziałach
4. Zbuduj chmury wyrazów dla wszystkich rozdziałów i jedną chmurę dla całego dokument.

Zadanie 3 — Zainstaluj język Scala na swoim komputerze i pobaw się w konsoli REPL podstawowymi obiektami tego języka. Rozszyfruj i zapamiętaj skrót REPL

Zadanie 4 — Oprogramuj w języku Scala funkcje gcd (największy wspólny dzielnik) oraz lcm (najmniejsza wspólna wielokrotność). Ustalmy, że $\text{gcd}(0, 0) = \text{lcm}(0, 0) = 0$ oraz $\text{gcd}(a, b) = \text{gcd}(|a|, |b|)$ i $\text{lcm}(a, b) = \text{lcm}(|a|, |b|)$.

1. Oprogramuj następnie funkcję Eulera phi zdefiniowaną wzorem

$$\phi(n) = |\{k \in \{1, \dots, n\} : \text{gcd}(k, n) = 1\}|$$

Spróbuj zastosować (mocno nieefektywną w tym przypadku) metodę `count` do zakresu `Range(1,n+1)`. Sprawdź, czy na pewno otrzymasz $\phi(1) = 1$.

2. Sprawdź poprawność napisanych funkcji obliczając
`Range(1,101).filter(100%_==0).map(x=>Euler(x)).sum`
w konsoli REPL. Powinno wyjść 100
3. Poznaj prosty dowód tego, że $\sum_{k|n} \phi(k) = n$ dla każdej liczby naturalnej $n \geq 1$.

Zadanie 5 — Zrealizuj Zadanie 1 w języku Scala.

1. Zaimportuj bibliotekę `io.Source` (`import scala.io.Source`)
2. Skorzystaj z polecenia `Source.fromFile(source, "UTF-8")` do wczytania pliku, zamień go na łańcuch (`mkString`) i następnie podziel na wyrazy (`split("\\s+")`). Możesz to zrobić jednym poleceniem.
3. Wczytaj stop-słowa i podziel je na wyrazy.
4. Usuń z książki stop słowa (coś w stylu `Book.filterNot(Stop.contains(_))`)
5. Pogrupuj słowa książki (coś w stylu `Filtered.groupBy(x=>x)`)
6. Zredukuj (coś w rodzaju `Grouped.mapValues(x=>x.length)`)
7. Posortuj według drugiego parametru (`np. reduced.toSeq.sortWith((x,y)=>x._2>y._2)`)
8. Zapisz wynik do pliku. Uwaga: możesz skorzystać z obiektu `PrintWriter` z bibliotek `java.io`

Zadanie 6 — Załóżmy, że mamy dostęp do bazy zakupów klientów w sieci hurtowni środków chemicznych z poprzedniego roku. W ciągu roku 10^7 klientów odwiedza ją 10 razy i za każdym razem kupuje średnio 10 różnego typu produktów z puli 200 dostępnych typów produktów. Załóżmy że znaleźliśmy w tej bazie danych dwóch klientów którzy zakupili choć raz ten sam koszyk produktów. Czy jest to czysty przypadek?

2 Funkcje haszujące

Zadanie 7 — Rozważmy funkcję haszującą zadaną wzorem $h(x) = x \bmod 21$. Stosujemy ją do liczb podzielnych przez pewną stałą c . Dla jakich stałych c jest to odpowiednia funkcja haszująca, czyli dla jakich stałych c można się spodziewać, że rozkład załadowania kubeków $\{0, \dots, 20\}$ będzie jednostajny?

Zadanie 8 — Znajdź wzór na rząd elementu $k \in \{0, \dots, n-1\}$ w grupie $C_n = (\{0, \dots, n-1\}, \oplus_n)$? Jaki jest związek tego zadania z poprzednim zadaniem?

Zadanie 9 — Mamy n kubeków. Rzucamy do nich k kul.

1. Oszacuj k taki aby z dużym prawdopodobieństwem doszło do 3-kolizji, czyli aby a jakimś kubku znalazły się 3 kulki.
2. Sprawdź eksperymentalnie otrzymany wynik
3. Uogólnij zadanie na a - kolizje

Zadanie 10 — Dwóch studentów ma dzban wypełniony 8 litrami napoju. Mają do dyspozycji dzbanek o pojemności 5 litrów oraz drugi dzbanek o pojemności 3 litrów. Chcą podzielić się równo napojem. Jak mogą to zrobić? Zagadanie to można potraktować jako system przepisujący o stanie początkowym $\{8, 0, 0\}$ Następujący kod (język Mathematica) opisuje pojedynczy, losowy krok transformacji stanu.

```
Move[C_] := Block[{x,y,V={8,5,3},Kopia,suma},
  Kopia = C;
  {x,y} = RandomSample[{1,2,3},2]; (*Chcę przelać z x do y*)
  suma = C[[x]]+C[[y]];
  If[suma<=V[[y]],
    Kopia[[x]]=0;Kopia[[y]]=suma,
    Kopia[[x]]=suma-V[[y]];Kopia[[y]]=V[[y]];
  ];
  Kopia
]
```

Można go uruchomić w pętli, czekając aż osiągniemy stan $\{4, 4, 0\}$. Jednak jest to kiepskie rozwiązanie - algorytm taki wpada bardzo często w pętle. Zastosuj tablicę mieszającą (`hashCode`) do kontroli historii przebiegu tego

algorytmu (ma ona służyć do unikania zapętleń).

Wskazówka: możesz użyć np. java: `java.util.Hashtable`; Scala: `scala.collection.mutable.Set.empty[List[Int]]`; Python: `np. set`; Wszystkie te klasy są oparte na HashTables.

3 Model MapReduce

Na razie zadania programistyczne realizujemy w standardowych językach programowania (Java, Python, Scala).

3.1 Działania

Zadanie 11 — Załóżmy, że \star jest działaniem łącznym.

1. Pokaż, że $(a \star b) \star (c \star d) = a \star (b \star (c \star d)) = ((a \star b) \star c) \star d$.
2. Uogólnij to zadanie na dowolną liczbę zmiennych.
3. Ile różnych wyrażeń możesz zbudować dla pięciu zmiennych? Wskazówka: Może przydać się zapisanie tych wyrażeń w postaci drzew.

Zadanie 12 — Niech $x \oplus y = x + y + 1$ oraz $x \odot y = xy + x + y$ dla $x, y \in \mathbb{R}$. Pokaż, że są to działania łączne i przemienne na \mathbb{R} . Wskazówka: Spróbuj to zrobić z minimalną liczbą rachunków.

Zadanie 13 — Podaj kilka przykładów działań nieprzemiennych. Podaj kilka przykładów działań które nie są łączne.

Zadanie 14 — Pokaż, że operacje $\min(x, y)$ i $\max(x, y)$ są przemienne i łączne. Czy operacja $s(x, y) = \frac{x+y}{2}$ jest łączna?

Zadanie 15 — Co robią następujące polecenia języka Python?

1. `list(filter(lambda x: x%2==0, range(1, 100)))`
2. `list(map(lambda x: x*x, range(1, 10)))`
3. `reduce(lambda x, y: x+y, [1, 2, 3, 4, 5])`
4. `reduce(lambda x, y: x*y, [1, 2, 3, 4, 5])`
5. `reduce(lambda x, y: x/y, [1, 2, 3, 4, 5])`

Uwaga: funkcję `reduce` zaimportuj z biblioteki **functions**.

3.2 Algorytmy MapReduce

Zadanie 16 — Wymień jakie aspekty działania systemu MapReduce są poza zasięgiem programisty. Które elementy kontroluje programista?

Zadanie 17 — Zaprojektuj algorytm MapReduce który dostaje bardzo duży zbiór liczb całkowitych i produkuje na wyjściu:

1. Największą liczbę.
2. Średnią wszystkich liczb.
3. Ten sam zbiór ale bez powtórzeń.
4. Liczbę różnych elementów bez powtórzeń.

Zadanie 18 — (Odwrocenie grafu) Dany jest graf w postaci listy sąsiadów: $[w, [w_i, w_{i1}, w_{i2}, \dots, w_{i,n_i}]]$ zapisany w zbiorze tekstowym, np

```
[
  [1, [3, 4, 5]],
  [2, [1, 3]],
  [3, [4, 5]],
  [4, [1, 2]],
  [5, [4, 5]]
]
```

Zastosuj technologię MapReduce do zbudowania grafu z odwróconymi linkami.

Wskazówka: Jeśli programujesz w języku Python, to możesz skorzystać z funkcji `groupby` z biblioteki `itertools`; pamiętaj, że lista par którą chce się grupować musi być posortowana. W języku Scala jest jeszcze łatwiej: przyjrzyj się metodzie `groupBy` stosowalnej do klasy `Traversable`.

Zadanie 19 — (Częste produkty) Mamy dany duży zbiór koszyków zakupowych z hipermarketu. Wyznacz zbiór wszystkich częstych par, czyli takich par produktów, które często występują w jednym koszyku. Załóżmy, że zbiór wszystkich możliwych par występujących w jednym koszyku jest tak duży, że nie jesteśmy w stanie ich wszystkich przetworzyć w realnym czasie.

Wskazówka: Jeśli para jest częsta to i każdy z jej składników jest częsty.

Zadanie 20 — (Odwrócony Indeks) Mając dany zbiór dokumentów zbuduj inverted index słów w nich występujących.

Zadanie 21 — Zaprojektuj algorytm MapReduce, który wyznacza złączenie dwóch relacji o schemacie $R(A,B,C)$ i $S(X,Y,Z)$ według połączenia $B=X$ oraz $C=Y$, czyli wyznacz tabelę

$$\{(A, Y) : (\exists B, C)(R(A, B, C) \wedge S(B, C, Y))\}.$$

Zadanie 22 — W pliku `TwoCollisions.csv`, do którego link znajduje się na stronie wykładu, w każdej linii znajduje się `(NumerHotelu, NumerDnia, NumerOsoby)`. Znajdź takie osoby, które w dwóch różnych dniach znajdowały się w tym samym hotelu.

Zadanie 23 — Niech $F : ((\mathbb{N} \setminus \{0\}) \times \mathbb{R})^2 \rightarrow (\mathbb{N} \setminus \{0\}) \times \mathbb{R}$ będzie funkcją określoną wzorem

$$F([c_1, x_1], [c_2, x_2]) = [c_1 + c_2, \frac{c_1 x_1 + c_2 x_2}{c_1 + c_2}]$$

1. Pokaż, że F jest działaniem przemienne i łącznym.
2. Oznaczmy przez \odot działanie $x \odot y = F(x, y)$. Znajdź zwartą formułę dla

$$[c_1, x_1] \odot [c_2, x_2] \odot \dots \odot [c_n, x_n].$$

3. Zastosuj tę własność funkcji do zastosowania combainera dla problemu wyznaczania średniej i wariancji.

Zadanie 24 — Zastosuj metodę map-reduce do wyznaczenia średniej geometrycznej i harmonicznej.

Zadanie 25 — Zastosuj metodę map-reduce do wyznaczenia wszystkich anagramów występujących w zbiorze tekstowym.

Zadanie 26 — Multizbiorem o skończonym nośniku Ω nazywamy funkcję $F : \Omega \rightarrow \mathbb{N}$. Dla $F, G : \Omega \rightarrow \mathbb{N}$ określamy $(F \cup G)(\omega) = \max\{F(\omega), G(\omega)\}$, $(F \cap G)(\omega) = \min\{F(\omega), G(\omega)\}$, $(F \setminus G)(\omega) = \max\{F(\omega) - G(\omega), 0\}$. Zaprojektuj map-reduce algorytm do wyznaczania tych trzech operacji. Algorytm na wejściu dostaje listę elementów zbioru

$$\{(1, \omega, F(\omega)) : \omega \in \Omega \wedge F(\omega) > 0\} \cup \{(2, \omega, G(\omega)) : \omega \in \Omega \wedge G(\omega) > 0\}$$

Zadanie 27 — W pliku `word-count.scala` znajduje się skrypt symulujący pracę systemu MapReduce dla problemu word-count.

1. Przekształć ten skrypt w bardziej realistyczny model - zapisz wynik pośredni (zmienna `keyval` z funkcji `TextMapper`) do pliku roboczego. Funkcja `TextReducer` ma pobierać wyniki z tego pliku.
2. Skróć przekształcony skrypt. Na przykład, dwie linijki z pliku `word-count.scala`

```
val grouped = keyval.groupBy(_.1)
val reduced = grouped.mapValues(_.size)
```

mogą być skrócone do jednej linijki

```
val reduced = keyval.groupBy(_.1).mapValues(_.size)
```

4 Podobieństwo tekstów

Zadanie 28 — Pokaż, że funkcja $d(A, B) = |A \triangle B|$ jest metryką na przestrzeni niepustych skończonych podzbiorów ustalonego zbioru X .

Zadanie 29 — Niech $f : [0, \infty) \rightarrow [0, \infty)$ będzie funkcją rosnącą i wklęsłą.

1. Pokaż, że dla $a, b \geq 0$ mamy $f(a + b) \leq f(a) + f(b)$.

Wskazówka: Zauważ, że możemy założyć, że $a + b > 0$; następnie zauważ, że $a = (a + b) \frac{a}{a+b}$ oraz $b = (a + b) \frac{b}{a+b}$ i zastosuj nierówność Jensena dla funkcji wklęsłych.

2. Załóżmy dodatkowo, że $f(0) = 0$. Niech d będzie metryką na zbiorze X . Pokaż, że funkcja $\rho(x, y) = f(d(x, y))$ jest również metryką na zbiorze X .
3. Pokaż, że jeśli $\epsilon \in (0, 1)$ oraz d jest metryką na zbiorze X , to funkcja $\rho(x, y) = d(x, y)^\epsilon$ jest metryką na zbiorze X .
4. Pokaż, że jeśli d jest metryką na zbiorze X , to funkcja $\rho(x, y) = \frac{d(x, y)}{1 + d(x, y)}$ jest metryką na zbiorze X .

Zadanie 30 — Wybierzmy dwa losowe m -elementowe podzbiory A, B n -elementowego zbioru X . Jaka jest wartość oczekiwana podobieństwa Jaccarda $J(A, B)$?

Zadanie 31 — Korzystając z Twierdzenia o Gęstości Liczb Pierwszych (Prime Numbers Theorem) oszacuj liczbę liczb pierwszych z przedziału $[2^{64}, 2^{64} + 1000]$ i następnie wyznacz te liczby.

Zadanie 32 — (**Twierdzenie Steinhausa**) Niech d będzie metryką na zbiorze X . Ustalmy element $a \in X$ i zdefiniujmy funkcję

$$\rho(x, y) = \frac{2d(x, y)}{d(x, a) + d(y, a) + d(x, y)}$$

Celem tego zadania jest pokazanie, że ρ jest metryką na zbiorze X .

1. Pokaż najpierw, że jeśli $0 < p \leq q$ oraz $r \geq 0$ to $\frac{p}{q} \leq \frac{p+r}{q+r}$.
2. Wprowadź oznaczenia $p = d(x, y)$, $q = d(x, y) + d(x, a) + d(y, a)$ oraz $r = d(x, z) + d(y, z) - d(x, y)$ i zastosuj obserwację z poprzedniego punktu do pokazania nierówności trójkąta dla funkcji ρ .

Zadanie 33 — Zastosuj Twierdzenie Steinhausa do metryki $d(X, Y) = |X \triangle Y|$ na zbiorze skończonych podzbiorów zbioru Ω do pokazania, że funkcja $d(X, Y) = 1 - S(X, Y)$ (odległość Jaccarda) jest metryką.

Zadanie 34 — Załóżmy, że S jest takim podobieństwem obiektów przestrzeni Ω , że istnieje rodzina funkcji haszujących \mathcal{H} oraz prawdopodobieństwo na rodzinie \mathcal{H} takie, że dla dowolnych dwóch obiektów $A, B \in \Omega$ mamy

$$P_h[h(A) = h(B)] = S(A, B)$$

Pokaż, że wtedy funkcja $d(A, B) = 1 - S(A, B)$ jest metryką na zbiorze Ω .

Zadanie 35 — Uzupełnij szczegóły dowodu tego, że jeśli $\Omega = \{\omega_i : 1 \leq i \leq N\}$, π jest losową permutacją zbioru $\{1, \dots, N\}$ (wybraną zgodnie z rozkładem jednostajnym), oraz $h_\pi(X) = \min\{k : \omega_{\pi(k)} \in X\}$ dla $X \subseteq \Omega$ to

$$P_\pi[h_\pi(A) = h_\pi(B)] = S(A, B).$$

Zadanie 36 — Napisz procedurę o specyfikacji `jaccard(f1:String, f2:String, k:Integer):Double`, która dla plików o nazwach `f1, f2` wyznacza ich k -gramy i następnie wylicza ich odległość Jaccarda. Przed wyznaczeniem k -gramów pliki powinny być oczyszczone (minimum to usunięcie znaków nowej linii, tabulatorów oraz podwójnych spacji)

1. Zastosuj tę procedurę do kilku wariantów swoich plików z algorytmami (zastosuj 4-gramy)
2. Zastosuj tę procedurę do porównania kolejnych rozdziałów analizowanej w Zadaniu 2 książki.

Zadanie 37 — Zastosuj metodę minhash do poprzedniego zadania. Twoja procedura powinna zależeć od parametru H który określa liczbę funkcji haszujących stosowanych do budowania sygnatury tekstu.

1. Przetestuj tę procedurę na danych z poprzedniego zadania dla $H \in \{50, 100, 250\}$ - porównaj aproksymację odległości Jaccarda z jej dokładnymi wartościami.

Pamiętaj o wygenerowaniu wspólnej rodziny funkcji haszujących dla wszystkich analizowanych tekstów.

c.d.n.

Powodzenia,
Jacek Cichoń