

# Atrybuty

Dla każdych dwóch urządzeń porównujemy ich atrybuty i otrzymujemy ciąg podobieństwa, o wartościach z przedziału  $[0, 1]$ . Oznaczmy przez  $D_1, D_2$  porównywane urządzenia.

## 1 Czas korzystania

Reprezentujemy przez ciąg 24 - elementowy. Jeżeli w całym badanym okresie, w danej godzinie pojawiła się aktywność, to stawiamy 1, w p.p 0. Miarę ustalimy na podstawie analizy przykładów użytkowników z kilkoma urządzeniami i zależności występujących pomiędzy nimi. Czy podobieństwo czasu jest informacją pozytywną czy negatywną? W obu przypadkach można zastosować miarę  $\frac{\min\{c_1, c_2\}}{\max\{c_1, c_2\}}$ . (Chwilowo ignorujemy błędy wynikające z różnic stref czasowych do czasu uzyskania odpowiednich informacji.)

## 2 Dominujący kraj

Dla każdego urządzenia tworzymy ciąg  $(k_1, k_2, c)$ , gdzie  $k_1, k_2$  – to pierwszy i drugi najczęściej występujący kraj (przy czym kraj drugi wybieramy tylko, gdy występuje w co najmniej 30% zapytań, w przeciwnym przypadku  $k_2 = 0$ ),  $c$  - liczba wszystkich krajów. Za tę samą wartość na miejscach 1 i 2 przyznajemy odpowiednio wartości 0.5, 0.2. Liczbę krajów porównujemy następująco:  $0.3 \cdot \frac{\min\{c_1, c_2\}}{\max\{c_1, c_2\}}$ , gdzie  $c_1$  liczba krajów dla  $D_1$ ,  $c_2$  - dla  $D_2$ .

## 3 Dominujący region

Jak dla kraju.

## 4 Dominujące IP

Dla każdej godziny porównujemy dominujące IP. Liczbę takich samych wyników dzielimy przez liczbę godzin, w których oba urządzenia były aktywne.

## 5 Dominujący ISP

Jak przy IP.

## 6 Typ użytkownika określony względem URL

Z pliku requests usuwamy wszystkie informacje poza dev i url; grupujemy je po takim samym dev i klastrujemy devices za pomocą programu R, na podstawie podobieństwa zbiorów url przypisanych do jednego device (typ użytkownika). Informacje o przyporządkowaniu do danego klastra porównujemy dla dwóch urządzeń zero-jedynkowo.

## 7 Liczba odwiedzonych stron w ciągu godziny

Dla urządzenia tworzymy parę  $(Max, med)$ , gdzie  $Max$  - to maksimum, a  $med$  - mediana liczby stron otwartych w ciągu godziny. Dla  $D_1$  i  $D_2$  mamy odpowiednio pary  $(Max_1, med_1)$ ,  $(Max_2, med_2)$ . Niech  $M = \max\{Max_1, Max_2\}$ ,  $m = \min\{Max_1, Max_2\}$ . Miara:

$$\frac{1}{(1 + M - m)^2} \cdot$$

Podobnie porównujemy medianę.

## 8 Liczba stron startowych

Atrybut określony na podstawie pierwszego połączenia w ciągu dnia. Porównujemy podobnie jak w punkcie 7, z tym że ustalamy dodatkowy parametr *days* - liczba dni aktywności urządzenia. Jeżeli *days* < 10, to pomijamy porównywanie.

## 9 Typ połączenia (connection type)

Mamy 5 typów połączeń:

1. The database identifies dial-up (modem telefoniczny)
2. cellular (komórkowy)
3. cable
4. DSL
5. corporate connection speeds (wszyscy mają jedno IP)

Miara jak dla kraju.