

Projekt 1. část: návrh zpracování a uložení dat

Vysoké učení technické v Brně
Fakulta informačních technologií
UPA - Ukládání a příprava dat
25. října 2020

1 Zvolené téma

04: COVID-19 (dr. Rychlý)

2 Řešitelé

- Bc. Kapoun Petr – xkapou04,
- Bc. Nováček Pavel – xnovac16,
- Bc. Willaschek Tomáš – xwilla00.

3 Zvolené dotazy a formulace vlastního dotazu

- **Dotaz A:** vytvořte popisné charakteristiky pro alespoň 4 údaje (např. věk, pohlaví, okres, zdroj nákazy) z datové sady COVID-19: Přehled osob s prokázanou nákazou dle hlášení krajských hygienických stanic (využijte krabicové grafy, histogramy, atd.).
- **Dotaz B:** určete vliv počtu nemocných a jeho změny v čase na sousední okresy (aneb zjistěte jak se šíří nákaza přes hranice okresů).
- **Vlastní dotaz:** určete vliv věku na délku nemoci a úmrtnost.

4 Stručná charakteristika zvolené datové sady

Základním zdrojem dat pro všechny dotazy jsou **otevřené datové sady COVID-19 v ČR**[3] dostupné skrze veřejné API¹ ve formátech JSON a CSV.

Datová sada COVID-19: Přehled osob s prokázanou nákazou dle hlášení krajských hygienických stanic (v2) obsahuje následující charakteristiky o nakažených osobách:

- **datum** – datum, kdy byla osoba pozitivně testována,
- **vek** – věk nakažené osoby,
- **pohlavi** – pohlaví nakažené osoby,
- **kraj_nuts_kod** – identifikátor kraje podle klasifikace NUTS 3, ve kterém byla pozitivní nákaza hlášena krajskou hygienickou stanicí,
- **okres_lau_kod** – identifikátor okresu podle klasifikace LAU 1,
- **nakaza_v_zahranici** – příznak, zda došlo k nákaze mimo ČR,
- **nakaza_zeme_csu_kod** – identifikátor státu v zahraničí, kde došlo k nákaze (dvoumístný kód z číselníku zemí CZEM).

¹<https://onemocneni-aktualne.mzcr.cz/api/v2/covid-19>

Další dvě použité datové sady jsou COVID-19: Přehled vyléčených dle hlášení krajských hygienických stanic a COVID-19: Přehled úmrtí dle hlášení krajských hygienických stanic, které mají shodnou strukturu:

- `datum` – datum vyléčení nebo úmrtí osoby,
- `vek`,
- `pohlavi`,
- `kraj_nuts_kod`,
- `okres_lau_kod`.

Jelikož tyto datové sady používají identifikátory *NUTS 3*² pro kraje a *LAU 1*³ pro okresy, využijeme dále číselníky od Českého statistického úřadu, které obsahují mapování těchto identifikátorů na odpovídající kraje a okresy a informace o nich.

Číselník okresů⁴ a **číselník krajů**⁵ jsou oba ve formátu CSV a obsahují:

- `KODJAZ` – kód jazykové verze textů,
- `AKRCIS` – akronym číselníku/klasifikace,
- `KODCIS` – kód číselníku/klasifikace,
- `CHODNOTA` – kód položky,
- `ZKRTEXT` – zkrácený název položky,
- `TEXT` – plný název položky,
- `ADMPLD` – počátek administrativní platnosti,
- `ADMNEPO` – konec administrativní platnosti.

Zvoleným programovacím jazykem je Python[4]. K získání dat použijeme standardní moduly jazyka Python[4] `json`⁶ a `csv`⁷. K nahrání dat použijeme modul `elasticsearch`⁸.

5 Zvolený způsob uložení uložených surových dat

Pro uložení dat jsme zvolili NoSQL řešení `Elasticsearch`[2]⁹. `Elasticsearch` je dokumentově orientovaný vyhledávací engine naprogramovaný v jazyce Java s použitím `Lucene`[1], díky kterému je velmi efektivní při komplexním vyhledávání na velkých objemech dat. Schéma struktur kolekcí odpovídá struktuře uložených dat ve vstupních datových sadách. Index je typu klíč-hodnota a pro každý vstupní soubor bude existovat právě jeden index. Klíč nabývá hodnoty `CHODNOTA` u číselníku okresů a krajů a ve zbývajících indexech se generuje automaticky, protože data neobsahují žádnou unikátní hodnotu.

²Nomenklatura územních statistických jednotek; úroveň 3 odpovídá krajům

³Local Administrative Units; úroveň 1 odpovídá okresům

⁴Číselník okresů: http://apl.czso.cz/iSMS/cisexp.jsp?kodcis=109&typdat=0&cisvaz=80007_210&datapohl=20.10.2020&cisjaz=203&format=2&separator=,

⁵Číselník krajů: http://apl.czso.cz/iSMS/cisexp.jsp?kodcis=100&typdat=0&cisvaz=80007_885&datapohl=20.10.2020&cisjaz=203&format=2&separator=,

⁶Dokumentace modulu `json` jazyka Python[4]: <https://docs.python.org/3/library/json.html>

⁷Dokumentace modulu `csv` jazyka Python[4]: <https://docs.python.org/3/library/csv.html>

⁸Modul `elasticsearch` pro Python: <https://pypi.org/project/elasticsearch/>

⁹Dokumentace `Elasticsearch`[2]: <https://www.elastic.co/guide/en/elasticsearch/reference/current/index.html>

Reference

- [1] APACHE SOFTWARE FOUNDATION. *Apache Lucene* [online]. 1999 [cit. 2020-10-20]. Dostupné na: <<https://lucene.apache.org/>>.
- [2] ELASTIC NV. *Elasticsearch* [online]. 2010 [cit. 2020-10-19]. Dostupné na: <<https://www.elastic.co/>>.
- [3] KOMENDA M., KAROLYI M., BULHART V., ŽOFKA J., BRAUNER T., HAK J., JARKOVSKÝ J., MUŽÍK J., BLAHA M., KUBÁT J., KLIMEŠ D., LANGHAMMER P., DAŇKOVÁ Š., MÁJEK O., BARTŮŇKOVÁ M., DUŠEK L.. COVID-19: Přehled aktuální situace v ČR. Onemocnění aktuálně. In [online]. Praha: Ministerstvo zdravotnictví ČR, 2020 [cit. 2020-10-20]. Vývoj: společné pracoviště ÚZIS ČR a IBA LF MU. Dostupné na: <<https://onemocneni-aktualne.mzcr.cz/covid-19>>. ISSN 2694-9423.
- [4] PYTHON SOFTWARE FOUNDATION. *Python* [online]. 1991 [cit. 2020-10-20]. Dostupné na: <<https://www.python.org/>>.