

# VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

## FAKULTA INFORMAČNÍCH TECHNOLOGIÍ

### UPA – COVID-19

Filip Bali (xbalif00)  
Natália Holková (xholko02)  
Roland Žitný (xzitny01)

16. decembra 2021

# Obsah

<b>1</b>	<b>Prvá časť: návrh spracovania a uloženia dát</b>	<b>2</b>
1.1	Vybrané dotazy . . . . .	2
1.1.1	Dotazy skupiny A . . . . .	2
1.1.2	Dotaz skupiny B . . . . .	2
1.1.3	Dotaz skupiny C . . . . .	2
1.1.4	Vlastné dotazy . . . . .	2
1.1.5	Zdroje potrebných dát . . . . .	3
1.2	Schéma dát a ich spracovanie . . . . .	3
1.3	Indexovanie . . . . .	4
<b>2</b>	<b>Druhá časť: implementovaný systém pre získavanie, ukladanie a spracovanie</b>	<b>5</b>
2.1	Riešenie dotazu A1 . . . . .	5
2.2	Riešenie dotazu A3 . . . . .	5
2.3	Riešenie dotazu B1 . . . . .	6
2.4	Riešenie dolovacej úlohy C1 . . . . .	7
2.5	Riešenie 1. vlastného dotazu . . . . .	9
2.6	Riešenie 2. vlastného dotazu . . . . .	10
2.7	Grafické užívateľské rozhranie . . . . .	11
<b>3</b>	<b>Vybrané technológie a spustenie</b>	<b>13</b>
3.1	Technológie . . . . .	13
3.2	Inštalácia . . . . .	13
3.3	Spustenie 1. časti projektu . . . . .	13
3.4	Spustenie 2. časti projektu . . . . .	13

# 1 Prvá časť: návrh spracovania a uloženia dát

## 1.1 Vybrané dotazy

V prvej časti projektu sme si vybrali tieto dotazy:

### 1.1.1 Dotazy skupiny A

#### A1

Vytvorte čiarový (spojnicový) graf zobrazujúci vývoj covidovej situácie po mesiacoch pomocou nasledujúcich hodnôt: počet novo nakazených za mesiac, počet novo vyliečených za mesiac, počet novo hospitalizovaných osôb za mesiac, počet vykonaných testov za mesiac. Pokiaľ nebude výsledný graf dobre čitateľný, zvážte logaritmickú mierku alebo rozdeľte hodnoty do viacerých grafov.

Pre bod **A1** je nutné získať počet *nakazených*, *vyliečených* a *hospitalizovaných* osôb za mesiac a počet vykonaných *testov* za mesiac.

#### A3

Vytvorte seriu stĺpcových grafov, ktoré zobrazia:

1. graf: počty vykonaných očkovaní v jednotlivých krajoch (celkový počet od začiatku očkovania).
2. graf: počty vykonaných očkovaní ako v predchádzajúcom bode navyše rozdelené podľa pohlavia. Diagram môže mať napr. dve časti pre jednotlivé pohlavia.
3. graf: počty vykonaných očkovaní, ešte ďalej rozdelené podľa vekovej skupiny. Pre potreby tohto diagramu postacia 3 vekové skupiny (0-24 rokov, 25-59, nad 59).

V bode **A3** sa musíme oboznámiť s jednotlivými *krajmi*, ich priradenými *kódmi* a následne s počtami *očkovaní* v jednotlivých krajoch a ich závislosťami na *pohlavie* a *vekovej skupine*.

### 1.1.2 Dotaz skupiny B

#### B1

Zostavte 4 rebríčky krajov "best in covid" za posledné 4 štvrťroky (1 štvrťrok = 1 rebríček). Ako kritérium voľte počet novo nakazených prepočítaný na jedného obyvateľa kraja. Pre jeden štvrťrok zobrazte výsledky tiež graficky. Graf bude pre každý kraj zobrazovať celkový počet novo nakazených, celkový počet obyvateľov a počet nakazených na jedného obyvateľa. Graf môžete zhotoviť kombinácie dvomi grafovi do jedného (jeden stĺpcový graf zobrazí prvé dve hodnoty a druhý, čiarový graf, hodnotu tretiu).

V bode **B1** sa zameriame znova na *kraje*, počet ich *obyvateľov* a následne počty *nakazených* v prepočte na jedného obyvateľa daných krajov.

### 1.1.3 Dotaz skupiny C

#### C1

Hľadanie skupín podobných miest z hľadiska vývoja covidu a vekového rozloženia obyvateľov.

Atribúty: počet nakazených za posledné 4 štvrťroky, počet očkovaných za posledné 4 štvrťroky, počet obyvateľov vo vekovej skupine 0-14 rokov, počet obyvateľov vo vekovej skupine 15 - 59, počet obyvateľov nad 59 rokov. Pre potreby projektu vyberte ľubovoľných 50 miest, pre ktoré nájdete potrebné hodnoty (môžete napr. využiť nejaký rebríček 50 najľudnatejších miest v ČR).

Vo finálnom dotaze **C1** je nutné získať data o jednotlivých *mestách* a *veku* ich obyvateľov, kedy budeme sledovať počet *nakazených* a *očkovaných*.

### 1.1.4 Vlastné dotazy

#### Prvý vlastný dotaz

Vytvorte stĺpcový graf zobrazujúci vývoj covidovej situácie na základe počtov úmrtí v jednotlivých krajoch. Pre tento dotaz je nutné získať počet *úmrtí* na kraje.

## Druhý vlastný dotaz

Vytvorte stĺpcový graf popisujúci štatistiku o nasledujúcich dvoch hodnotách:

- Počet očkovaných osôb v jednotlivých krajocho.
- Počet úmrtí v jednotlivých krajocho.

Pre tento dotaz je nutné sledovať počty *očkování* a *úmrtí* v jednotlivých krajocho.

### 1.1.5 Zdroje potrebných dát

- **Nakazení:** dátum | vek | pohlavie | kraj\_kód | okres\_kód | zahraničie | zem\_zahranicie (link)
- **Vyliečení:** dátum | vek | pohlavie | kraj\_kód | okres\_kód (link)
- **Hospitalizovaní:** dátum | prvý\_záznam | kum\_prvý\_záznam | počet | bez\_priznakov | ... (link)
- **Testy:** dátum | pocet\_PCR | pocet\_AG | typologie | pozitivny (link)
- **Očkovanie(regiony):** dátum | vakcina | kraj\_kód | kraj\_názov | veková\_skupina | typy\_dávok | celkom\_dávok (link)
- **Očkovanie(ludia):** dátum | vakcina | kód | poradie | veková\_skupina | pohlavie | pocet\_dávok (link)
- **Očkovanie(geograficky):** dátum | názov\_vakcíny | kód\_vakcíny | poradie\_dávky | kraj\_kod | cznuts | okres\_kod | orp\_kod | počet\_dávok (link)
- **Mŕtvi:** dátum | vek | pohlavie | kraj\_kód | okres\_kód (link)
- **Prehľad vykázaných očkování pod profesí:** dátum | vakcina | kraj\_nuts\_kod | kraj\_nazev | zarizeni\_kod | zarizeni\_nazev | poradí\_davky | vekova\_skupina | orp\_bydliste | orp\_bydliste\_kod | pohlavi | vakcina\_kod | ukoncuji\_davka (link)
- **Obyvateľstvo podľa päť ročných vekových skupín a pohlavia v krajocho a okresoch:** idhod | hodnota | stapro\_kód | pohlavia\_číselník | pohlavie\_kód | vek\_číselník | vek\_kód | oblast\_číselník | oblast\_kód | dátum | pohlavie\_text | vek\_text (link)
- **Číselník krajocho:** hodnota | text | cznuts | ruian | kraj\_skratka (link)
- **Číselník okresov:** hodnota | text | cznuts | ruian | okres\_lau | okres\_skratka (link)

## 1.2 Schéma dát a ich spracovanie

**Nakazení:** výsledné data ohľadne nakazených osôb bolo nutné očistiť od nepotrebných informácií a previesť pohlavie podľa slovníka.

date	age	gender	region	district
------	-----	--------	--------	----------

**Vyliečení:** tieto data zostali v pôvodnej podobe s prevedením pohlavia podľa slovníka.

date	age	gender	region	district
------	-----	--------	--------	----------

**Hospitalizovaní:** v dátach o hospitalizovaných nám postačujú iba informácie o dátume a počte hospitalizovaných. Ostatné boli odstránené a hodnota dátumu upravená pre naše potreby. Počty hospitalizovaných osôb sú zoskupené podľa mesiacov.

month	patients
-------	----------

**Testy:** pri počte vykonaných testov sa ponechali iba dátum a počet PCR a AG testov, ktoré sa následne zoskupili do jednej hodnoty vyjadrujúcej počet celkových testov za jednotlivé mesiace.

month	tests
-------	-------

**Očkovanie(regióny):** pre získanie počtu očkování v jednotlivých krajocho sme museli vyčistiť data od nepotrebných informácií a následne ich zoskupiť podľa krajocho.

region	count
--------	-------

date	age_group	gender
------	-----------	--------

**Očkovanie(ludia):** pre pozorovanie počtu očkovaných ľudí vzhľadom na pohlavie a vekovú skupinu, bolo nutné získať data zo spomenutého druhého zdroja a jemne očistiť data.

**Očkovanie(geograficky):** data o úmrtiach budú ďalej vhodné pri odpovedaní na vlastné dotazy.

date	age	gender	region	district
------	-----	--------	--------	----------

**Mŕtvi:** data o úmrtiach budú ďalej vhodné pri odpovedaní na vlastné dotazy.

date	age	gender	region	district
------	-----	--------	--------	----------

**Obyvateľstvo podľa päť ročných vekových skupín a pohlavia v krajocho a okresoch:** data o úmrtiach budú ďalej vhodné pri odpovedaní na vlastné dotazy.

date	age	gender	region	district
------	-----	--------	--------	----------

**Číselník krajov:** data o úmrtiach budú ďalej vhodné pri odpovedaní na vlastné dotazy.

date	age	gender	region	district
------	-----	--------	--------	----------

**Číselník okresov:** data o úmrtiach budú ďalej vhodné pri odpovedaní na vlastné dotazy.

date	age	gender	region	district
------	-----	--------	--------	----------

### 1.3 Indexovanie

## 2 Druhá časť: implementovaný systém pre získavanie, ukladanie a spracovanie

Pre prehľadnosť boli jednotlivé dotazy implementované vo vlastných súboroch, ktoré sa nachádzajú v *Python* balíčku `queries`. V prípade potreby je tieto dotazy možné volať samostatne, všetky naraz pomocou súboru `vizualizer.py`, alebo interaktívne zobrazovať cez grafické užívateľské rozhranie. Bližšie informácie k spusteniu sa nachádzajú v sekcii 3.4.

Pre získavanie dát z NoSQL databáze *MongoDB* sa používa knižnica `pymongo`. Na extrakciu dát do `.csv` súboru, jeho opätovné načítanie a úpravu pred vytváraním grafov slúži knižnica `pandas`. Zobrazovanie grafov sa rieši cez knižnicu `seaborn`.

### 2.1 Riešenie dotazu A1

#### Extrakcia dát z NoSQL databáze

Extrakcia dát z *MongoDB* prebieha vo funkcii `A1.extract_csv`, ktorá sa nachádza v súbore `queriesA1.py`. `aggregate infected - group by month and year and count total number in that month` Najskôr agregujeme počty nakazených, ktorý sa nachádzajú v tabuľke *infected* (tabuľka 1.2), podľa rokov a mesiacov a výsledok prevedieme do `Pandas.DataFrame`. Rovnako postupujeme pri dátach o vyliečených, ktoré sú dostupné v tabuľke 1.2. Dáta o hospitalizovaných a počtoch vykonaných testov nie je potrebné agregovať, stačí ich iba previesť do `Pandas.DataFrame`. Jednotlivé `DataFrame` sú zlúčené podľa mesiacov a výsledok je zapísaný do `.csv` súboru.

#### Štruktúra `.csv` súboru

Súbor vytvorený funkciou `A1.extract_csv` má štruktúru:

- *month* - mesiac v tvare YYYY-MM
- *infected* - počet nakazených za daný mesiac
- *cured* - počet vyliečených za daný mesiac
- *hospitalized* - počet hospitalizovaných za daný mesiac
- *tests* - počet vykonaných testov za daný mesiac

#### Grafické zobrazenie výsledkov z `.csv` súboru

Vytvorenie grafu z extrahovaných dát v `.csv` súbore prebieha vo funkcii `A1.plot_graph`, ktorá sa nachádza v súbore `queries/A1.py`. Bol vytvorený spojnicový graf, ktorý zobrazuje vývoj počtov nakazených, vyliečených, hospitalizovaných a vykonaných testov od začiatku pandémie. Počty sú agregované po mesiacoch. Údaje o vykonaných testoch začali byť vedené až neskoršie, čo možno vidieť na grafe. Pre prehľadnosť je použitá logaritmická mierka na osu y. Z grafu možno jasne vidieť nástup jednotlivých vln.

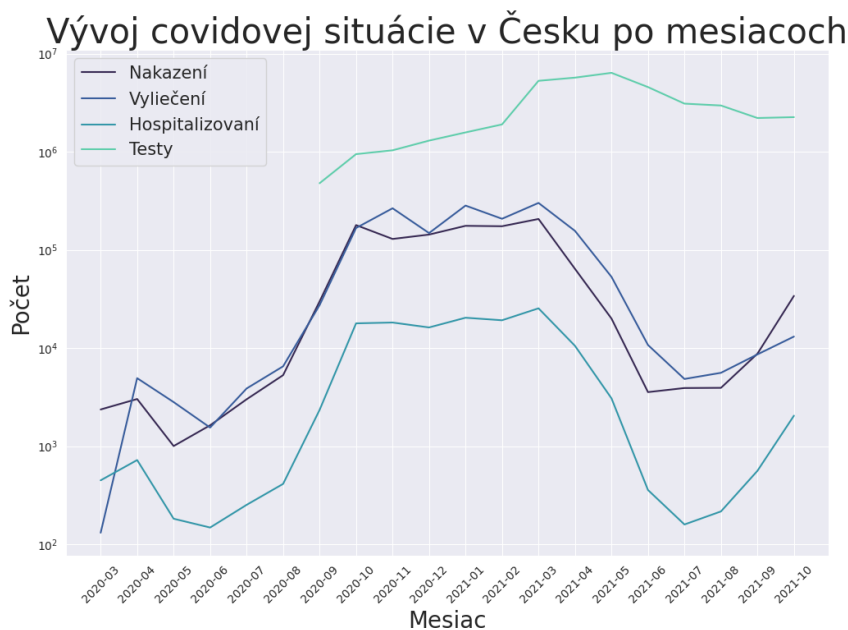
Výsledný graf je možné vidieť na obrázku 2.1.

### 2.2 Riešenie dotazu A3

#### Extrakcia dát z NoSQL databáze

Extrakcia dát z *MongoDB* prebieha vo funkcii `A3.extract_csv`, ktorá sa nachádza v súbore `queriesA3.py`. Proces agregácie je zreťazené spracovanie dokumentu obsahujúci očkovaných podľa profesií. Agregácia je rozdelená do nasledujúcich etáp:

- *match* - Z dokumentu vyberieme len záznamy, ktorých hodnota v atribúte *gender* je hodnota *M* alebo *Z*
- *project* - Následne vyberieme atribúty dokumentu *region*, *gender*, *age\_group*
- *group* - Zoskupíme záznamy dokumentu pod *\_id*, pričom je zložené z atribútov v slede: *region*, *gender*, *age\_group*. Ich následný počet je sčítaný pod atribút *count*.
- *sort* - Záznamy usporiadame podľa atribútu *region*.



Obr. 1: Dotaz A1 - graf

## Štruktúra .csv súboru

Súbor vytvorený funkciou `A3.extract_csv` má štruktúru:

- *region* - skratka kraja
- *count* - počet zaočkovaných ľudí kraja so zvoleným pohlavím a vekovou skupinou
- *gender* - pohlavie (M - muži, Ž - ženy)
- *age\_group* - vekové skupiny delené po 5 ročných intervaloch

## Grafické zobrazenie výsledkov z .csv súboru

Vytvorenie grafu z extrahovaných dát v .csv súbore prebieha vo funkcii `A3.plot_graph`, ktorá sa nachádza v súbore `queries/A3.py`.

Boli vytvorené celkovo tri stĺpcové grafy. Prvý graf zobrazuje celkový počet doterajších zaočkovaných na základe krajov. Druhý graf je podobný prvému, ale očkovaní sú tu ďalej rozdelení podľa pohlavia. Posledný graf ukazuje rozdelenie podľa vekovej skupiny, ktorou môžu byť *deti a mladí* (0-24), *dospelí* (25-59) a *starší* (60+). Kraje vo všetkých grafoch sú zoradené zostupne podľa počtu zaočkovaných.

Z grafov možno vidieť, že najviac zaočkovaných je v kraji *Hlavní město Praha*.

Výsledný graf je možné vidieť na obrázku 2.2.

## 2.3 Riešenie dotazu B1

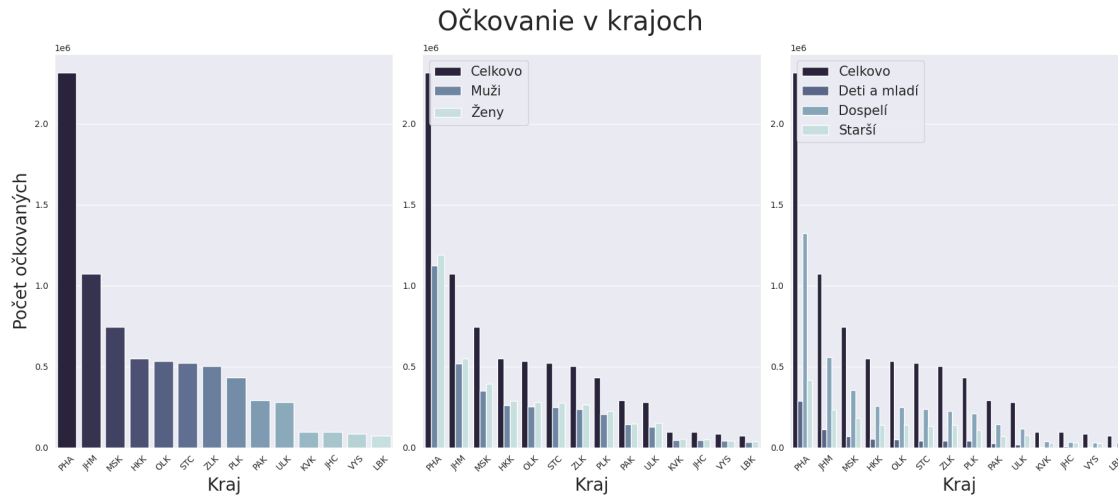
### Extrakcia dát z NoSQL databáze

Extrakcia dát z *MongoDB* prebieha vo funkcii `B1.extract_csv`, ktorá sa nachádza v súbore `queriesB1.py`.

Úloha vyžaduje agregáciu z viacerých dokumentov (infikovaný, demografické dáta, číselník krajov).

Dokument obsahujúci dáta infikovaných je v agregácii spracovaný v nasledujúcich zreťazených etápach:

- *match* - Z dokumentu vyberieme len záznamy, ktorých rok v atribúte *date* je 2021.
- *project* - Následne vyberieme atribúty dokumentu *date*, *region*, *quarter*, kde *quarter* je atribút odvodený od atribútu *date* a to tak, že špecifikuje v ktorom kvartáli roka 2021 sa záznam nachádza.
- *group* - Zoskupíme záznamy dokumentu pod *\_id*, pričom je zložené z atribútov v slede: *region*, *quarter*. Ich následný počet je sčítaný pod atribút *count*.



Obr. 2: Dotaz A3 - graf

- sort - Záznamy usporiadame podľa atribútov v slede: *region, quarter*.

Dokument demografických dát je spracovaný v jedinej etape *project*, ktorá vyberá atribúty *value, territory\_code, territory\_txt, valid\_date, gender\_code, gender\_txt, age\_code, age\_txt*.

Dokument obsahujúci číselník krajov je rovnako spracovaný v jedinej etape *project*, ktorá vyberá atribúty *cznuts, region\_shortcut, value*.

### Štruktúra .csv súboru

Funkcia *B1\_extract\_csv* vytvára 4 .csv súbory *B1\_Q1.csv, B1\_Q2.csv, B1\_Q3.csv* a *B1\_Q4.csv*, jeden pre každý štvrťrok. Každý súbor má nasledovnú štruktúru:

- *region\_shortcut* - skratka kraja
- *quarter* - číslo štvrťroka v poradí
- *population* - celková populácia kraja
- *infected* - počet nakazených v kraji za daný štvrťrok
- *infected\_normalized* - počet nakazených normalizovaný vzhľadom na populáciu kraja

### Grafické zobrazenie výsledkov z .csv súboru

Vytvorenie grafu z extrahovaných dát v .csv súbore prebieha vo funkcii *B1\_plot\_graph*, ktorá sa nachádza v súbore *queries/B1.py*.

Predvolene sme zvolili, že budeme vykresľovať rebríčky "best in covid" krajov z posledného štvrťroka roku 2021. Horný stĺpcový graf zobrazuje populáciu kraja a počet nakazených za daný štvrťrok. Dolný spojnicový graf predstavuje počet nakazených na 1 obyvateľa na základe kraja. Obe grafy majú zoradenú os x zostupne podľa tohto počtu nakazených na 1 obyvateľa.

Výsledný graf je možné vidieť na obrázku 2.3.

## 2.4 Riešenie dolovacej úlohy C1

### Extrakcia dát z NoSQL databáze

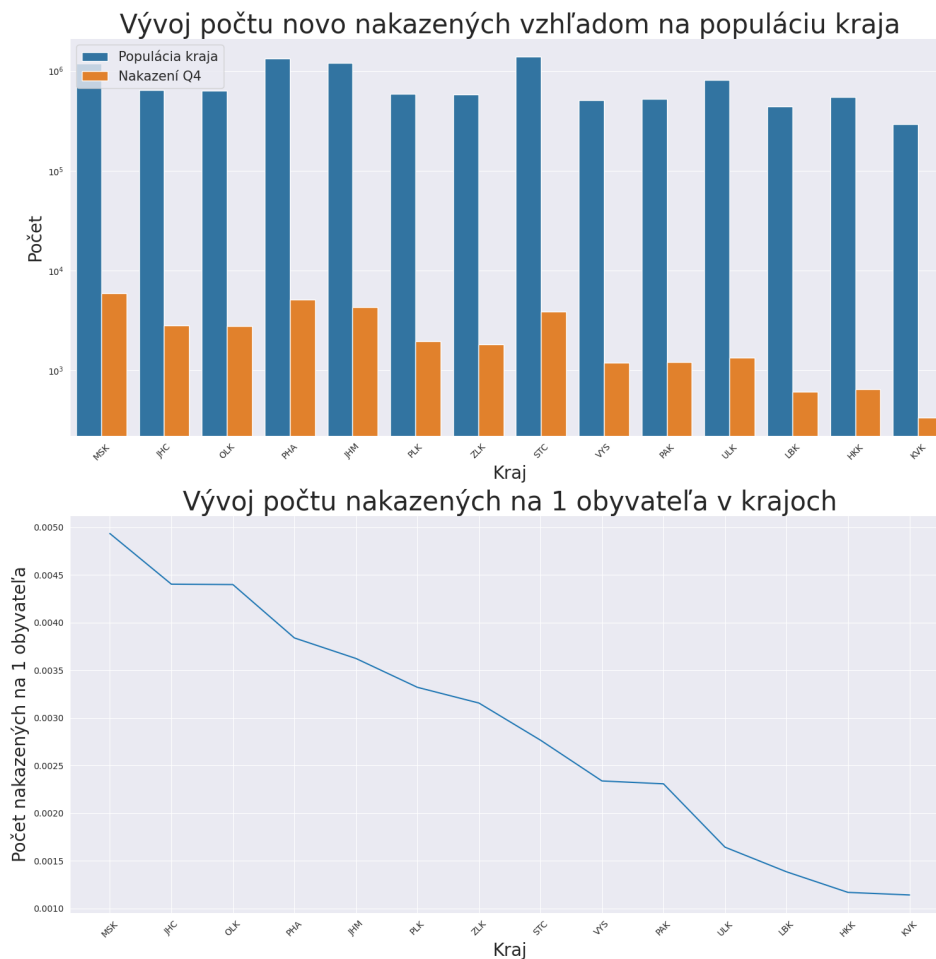
Úloha vyžaduje agregáciu z viacerých dokumentov (infikovaný, demografické dáta, číselník okresov, geografický prehľad vykázaných očkovaní).

Dokument obsahujúci dáta infikovaných je v agregácii spracovaný v nasledujúcich zreťazených etapách:

- match - Z dokumentu vyberieme len záznamy, ktorých rok v atribúte *date* je 2021.
- project - Následne vyberieme atribúty dokumentu *date, district*.



## "Best in covid" za štvrtrok Q4 2021



Obr. 3: Dotaz B1 - graf

- group - Zoskupíme záznamy dokumentu pod *\_id*, ktorý obsahuje atribút *region*. Ich následný počet je sčítaný pod atribút *count*.

Dokument obsahujúci demografické dáta je v agregácii spracovaný v nasledujúcich zreteľovaných etapách:

- match - Z dokumentu vyberieme len záznamy, ktorých dátum v atribúte *valid\_date* je 2020-12-31 (tento dátum predstavuje najnovšie dáta v datovej sade). Zároveň je prítomný atribút *territory\_txt* z ktorého je vybraný zoznam 50 miest. Zoznam miest je pevne daný.
- project - Následne vyberieme atribúty dokumentu *value*, *territory\_txt*, *territory\_code*, *age\_txt*, *age\_code*, *gender\_txt*, *gender\_code*.

Dokument obsahujúci geografický prehľad vykázaných očkovaní je v agregácii spracovaný v nasledujúcich zreteľovaných etapách:

- match - Z dokumentu vyberieme len záznamy, ktorých rok v atribúte *date* je 2021. Zároveň je prítomný operátor *or*, ktorý zaisťuje, že vybrané záznamy budú pochádzať len od plne zaočkovaných osôb (Minimálne dve dávky u vakcín s kódom: *CO01*, *CO02*, *CO03*. Minimálne jedna dávka u vakcíny s kódom *CO04*.)
- project - Následne vyberieme atribúty dokumentu *date*, *district\_name*, *orp*, *shot\_count*.

- `group` - Zoskupíme záznamy dokumentu pod `_id`, ktorý je zložený z atribútu `district_name`. Následné sú naprieč všetkými záznamami v jednotlivých krajoch sčítané hodnoty z atribútu `shot_count` a výsledná suma každý kraj je pod atribútom `vaccinated_count`.

Dokument obsahujúci číselník okresov je spracovaný v jedinej etape `project`, ktorá vyberá atribúty `text`, `cznuts`, `okres_lau`, `region_shortcut`.

### Štruktúra .csv súborov

Obe .csv súbory majú nasledovnú štruktúru:

- `district_name` - názov okresného mesta
- `infected_count` - počet nakazených v meste
- `vaccinated_count` - počet očkovaných v meste
- `population_count(age_0-14)` - veľkosť populácie vo vekovej skupine 0-14 rokov
- `population_count(age_15-59)` - veľkosť populácie vo vekovej skupine 15-59 rokov
- `population_count(age_59-above)` - veľkosť populácie vo vekovej skupine 60+ rokov

### Normalizácia, diskretizácia a spracovanie odláhlých hodnôt

Zvolili sme si normalizovať stĺpec `population_count(age_0-14)`. Normalizáciu sme vykonali podelením pôvodného stĺpca celkovou populáciou mesta (teda súčet všetkých troch vekových skupín).

Rozhodli sme sa odstraňovať odláhlé hodnoty zo stĺpca `vaccinated_count` a nahradiť ich mediánom. Pri každej hodnote počítame, či  $|\text{počet očkovaných} - \text{medián}| > \text{štandardná odchýlka}$ , v tom prípade nahradíme mediánom.

Diskretizácia bola aplikovaná na vekové rozloženie v jednotlivých krajoch. Výsledkom je rozdelenie dát vekového rozloženia do troch skupín:

- Vek v intervale 0 až 14
- Vek v intervale 15 až 59
- Vek v intervale 60 a viac

## 2.5 Riešenie 1. vlastného dotazu

### Extrakcia dát z NoSQL databáze

Extrakcia dát z *MongoDB* prebieha vo funkcii `VL1_extract_csv`, ktorá sa nachádza v súbore `queries/VL1.py`. Úloha vyžaduje agregáciu z viacerých dokumentov (mŕtvý, číselník krajov).

Dokument obsahujúci dáta mŕtvých je v agregácii spracovaný v nasledujúcich zreteľovaných etapách:

- `match` - Z dokumentu vyberieme len záznamy, z atribútu `date` v ktorých je overený dátum a sú z desaťročia 2020.
- `project` - Následne vyberieme atribúty dokumentu `date`, `region`.
- `group` - Zoskupíme záznamy dokumentu pod `_id`, ktorý obsahuje atribút `region`. Následne sú pod atribút `dead_count` sčítané všetky záznamy patriace pod jednotlivé kraje.

Dokument obsahujúci číselník krajov je spracovaný v jedinej etape `project`, ktorá vyberá atribúty `region_shortcut`, `cznuts`.

### Štruktúra .csv súboru

Súbor vytvorený funkciou `VL1_extract_csv` má štruktúru:

- `region_shortcut` - skratka kraja
- `dead_count` - celkový počet mŕtvych v kraji

## Grafické zobrazenie výsledkov z .csv súboru

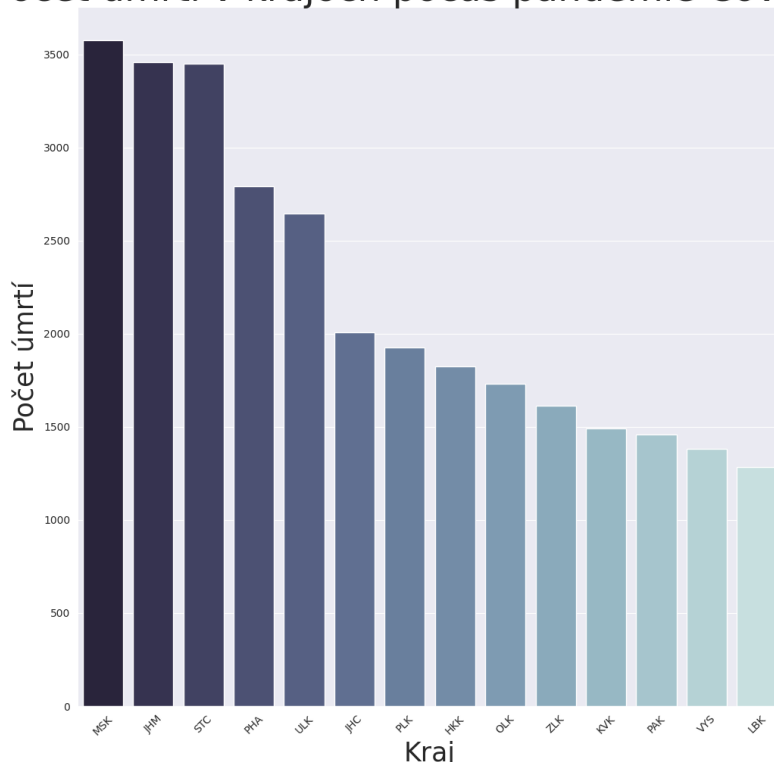
Vytvorenie grafu z extrahovaných dát v .csv súbore prebieha vo funkcii `VL1_plot_graph`, ktorá sa nachádza v súbore `queries/VL1.py`.

V tomto dotaze graficky zobrazujeme rozloženie ľudí, ktorí umreli na covid-19 počas celej pandémie, v závislosti na kraj ich pôvodu. Je to riešené spojnicovým grafom, ktorý je zostupne zoradený.

Z grafu je možné vidieť, že k najviac úmrtiam došlo v *Moravskoslezskom kraji*.

Výsledný graf je možné vidieť na obrázku 2.5.

Počet úmrtí v krajoch počas pandémie Covid-19



Obr. 4: Prvý vlastný dotaz - graf

## 2.6 Riešenie 2. vlastného dotazu

### Extrakcia dát z NoSQL databáze

Extrakcia dát z *MongoDB* prebieha vo funkcii `VL2_extract_csv`, ktorá sa nachádza v súbore `queries/VL2.py`. V tomto dotaze skúmame počty očkovaných a zomrelých v jednotlivých krajoch, preto využívame dáta z dvoch rôznych zdrojov.

Úloha vyžaduje agregáciu z viacerých dokumentov (úmrtia, demografické dáta, číselník krajov).

Dokument obsahujúci dáta mŕtvych je v agregácii spracovaný v nasledujúcich zreťazených etapách:

- `match` - Z dokumentu vyberieme len záznamy, z atribútu `date` v ktorých je overený dátum a sú z desaťročia 2020.
- `project` - Následne vyberieme atribúty dokumentu `date`, `region`.
- `group` - Zoskupíme záznamy dokumentu pod `_id`, ktorý obsahuje atribút `region`. Následne sú pod atribút `dead_count` sčítané všetky záznamy patriace pod jednotlivé kraje.

Dokument obsahujúci demografické dáta je spracovaný v jedinej etape *project*, ktorá vyberá atribúty *value*, *territory\_code*, *territory\_txt*, *valid\_date*, *gender\_code*, *gender\_txt*, *age\_code*, *age\_txt*.

Dokument obsahujúci číselník krajov je rovnako spracovaný v jedinej etape *project*, ktorá vyberá atribúty *text*, *cznuts*, *region\_shortcut*.

### Štruktúra .csv súboru

Súbor vytvorený funkciou `VL2.extract_csv` má štruktúru:

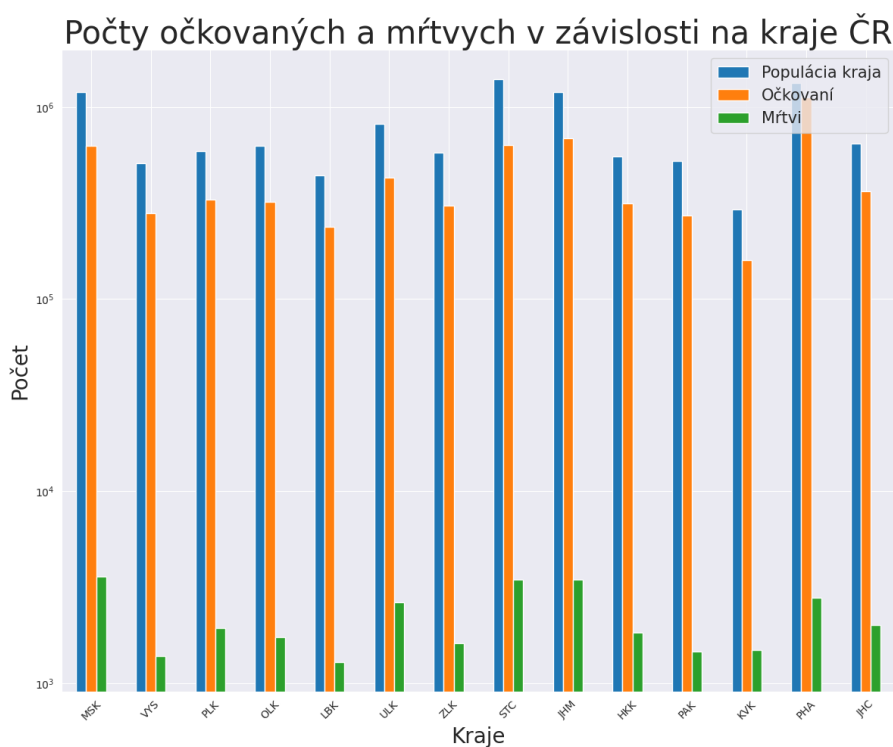
- *region\_shortcut* - skratka kraja
- *dead\_count* - celkový počet mŕtvych v kraji
- *vaccinated\_count* - celkový počet zaočkovaných v kraji
- *population* - celková populácia kraja

### Grafické zobrazenie výsledkov z .csv súboru

Vytvorenie grafu z extrahovaných dát v .csv súbore prebieha vo funkcii `VL2.plot_graph`, ktorá sa nachádza v súbore `queries/VL2.py`.

V tomto dotaze zobrazujeme celkovú populáciu kraja, počet očkovaných a počet zomrelých v kraji. Stĺpcový grafy používa logaritmickú mierku pre os y hlavne z dôvodu veľkého rozdielu medzi populáciou kraja a počtom zomrelých.

Výsledný graf je možné vidieť na obrázku 2.6.

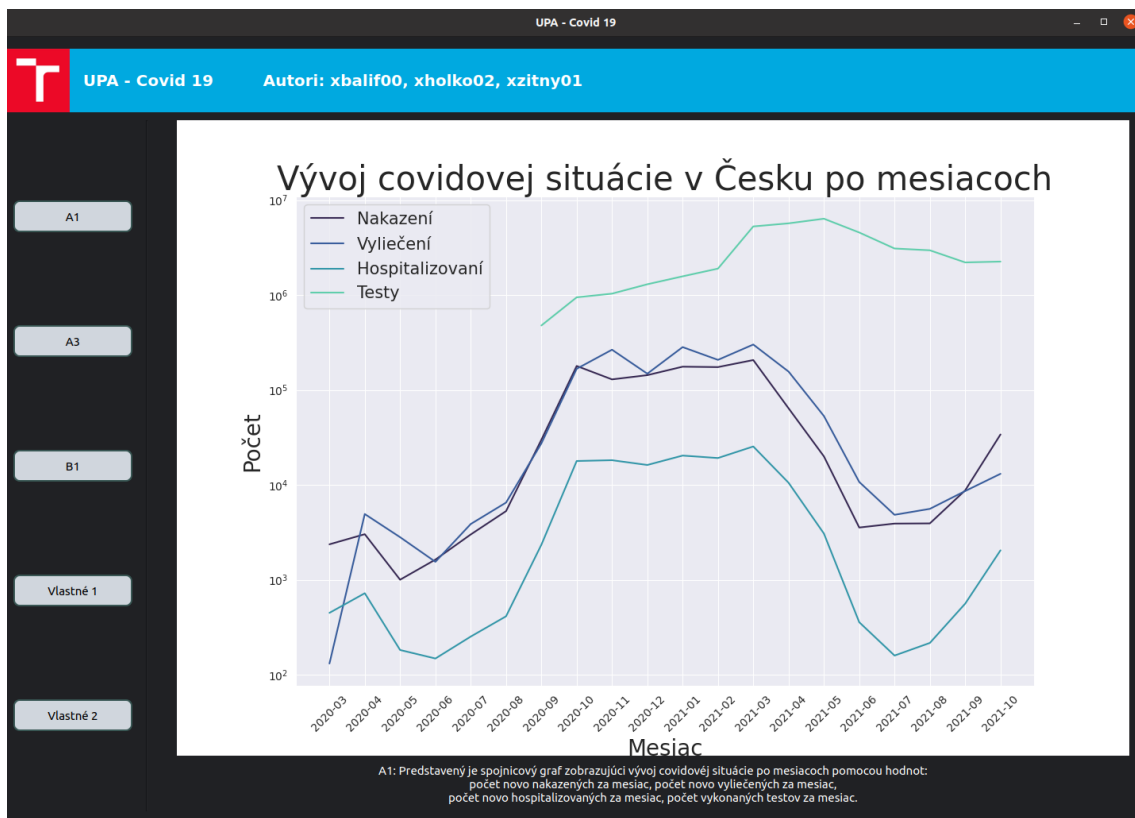


Obr. 5: Druhý vlastný dotaz - graf

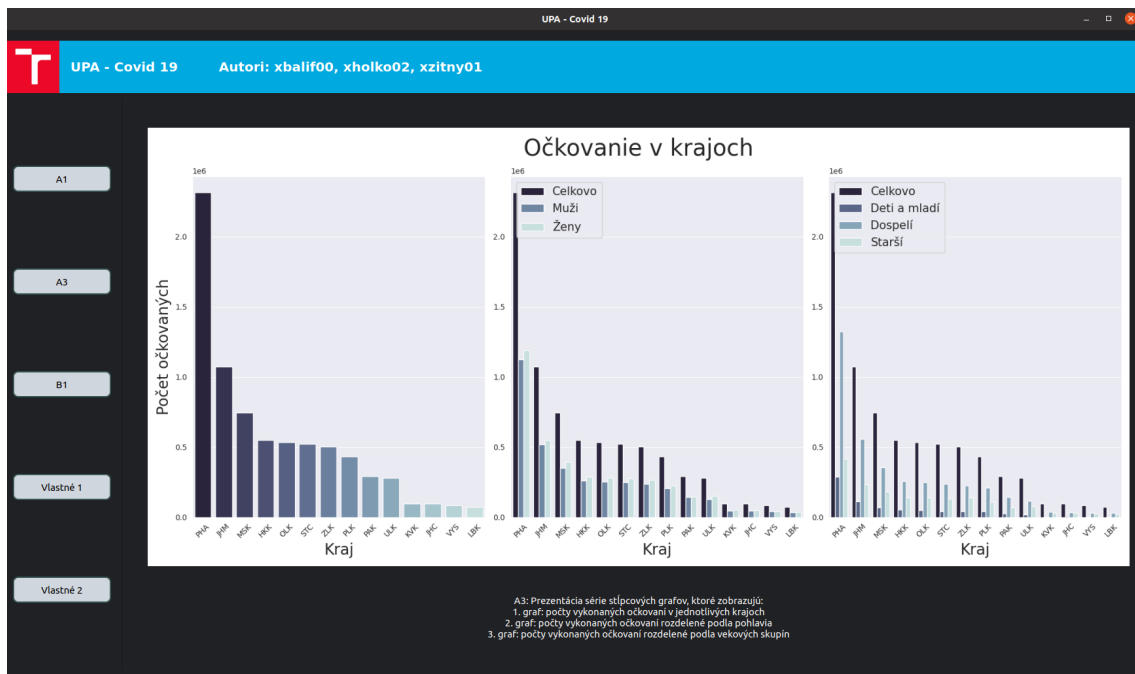
## 2.7 Grafické užívateľské rozhranie

Grafické užívateľské rozhranie bolo vytvorené pomocou multiplatformnej knižnici *PyQt*<sup>1</sup> a slúži pre jednoduché zobrazenie grafou vytvorených z dotazov preberaných v sekciách 2.1, 2.2, 2.3, 2.5 a 2.6. Toto rozhranie pozostáva z piatich tlačidiel, ktoré prezentujú výstup dotazov spolu s jednoduchým popisom, čo je možné vidieť na obrázkoch 2.7 a 2.7.

<sup>1</sup><https://wiki.python.org/moin/PyQt>



Obr. 6: GUI - dotaz A1



Obr. 7: GUI - dotaz A3

## 3 Vybrané technológie a spustenie

### 3.1 Technológie

- *Python3.8*
- *MongoDB* - NoSQL DB
- *Pandas* - manipulácia a analýza dát
- *Seaborn* - tvorba grafov
- *PyQt5* - tvorba GUI

### 3.2 Inštalácia

Inštalácia MongoDB - <https://docs.mongodb.com/manual/administration/install-community/>

```
python3 -m venv ./env-upa
source ./env-upa/bin/activate
pip3 install -r requirements.txt
```

### 3.3 Spustenie 1. časti projektu

Príklad spustenia 1. časti projektu:

```
python3 data_loader.py --mongo mongodb://localhost:27017/ -f download_data_folder -d UPA-db
```

### 3.4 Spustenie 2. časti projektu

Pred spustením 2. časti projektu je **nutné** znovu spustiť 1. časť podľa 3.3 aby bola zaručená správna štruktúra databáze. Bolo pripravených viacero možností, ako spúšťať dotazy a získavať výsledky. Pri všetkých možnostiach je potrebné používať rovnakú databázu, ktorá bola vytvorená v predchádzajúcom kroku.

#### Vygenerovanie všetkých .csv a grafov naraz

Pripravený skript `vizualizer.py` postupne extrahuje z databáze .csv pre všetky dotazy a zároveň vytvorí príslušné grafy.

```
python3 vizualizer.py --mongo mongodb://localhost:27017/ -d UPA-db
```

#### Spúšťanie úloh jednotlivo

V prípade nutnosti získať výsledky iba z jedinej úlohy, je možné toho docieľiť pomocou príkazov:

```
python3 queries/A1.py --mongo mongodb://localhost:27017/ -d UPA-db      # pre dotaz A1
python3 queries/A3.py --mongo mongodb://localhost:27017/ -d UPA-db      # pre dotaz A3
python3 queries/B1.py --mongo mongodb://localhost:27017/ -d UPA-db      # pre dotaz B1
python3 queries/C1.py --mongo mongodb://localhost:27017/ -d UPA-db      # pre dotaz C1
python3 queries/VL1.py --mongo mongodb://localhost:27017/ -d UPA-db     # pre dotaz VL1
python3 queries/VL2.py --mongo mongodb://localhost:27017/ -d UPA-db     # pre dotaz VL2
```

#### Spúšťanie cez grafické užívateľské rozhranie

Grafy z dotazov skupiny A, B a vlastných dotazov je navyše možné zobrazovať interaktívne cez grafické užívateľské rozhranie popísané v sekcii 2.7. Pred spustením je **nutné** spustiť vygenerovanie grafov cez `vizualizer.py`. Spustiť GUI je možné nasledovne:

```
python3 GUI/main.py
```