

# TOP-K THRESHOLD ALGORITHMUS

Vyhledávání top-k treshold výsledků v DB libovolných produktů s možností volby agregační funkce + GUI. Zdrojový kód: [HTTPS://GITHUB.COM/MSKL/FIT-BI-VMW-PROJECT](https://github.com/MSKL/FIT-BI-VMW-PROJECT)

## POPIS PROJEKTU

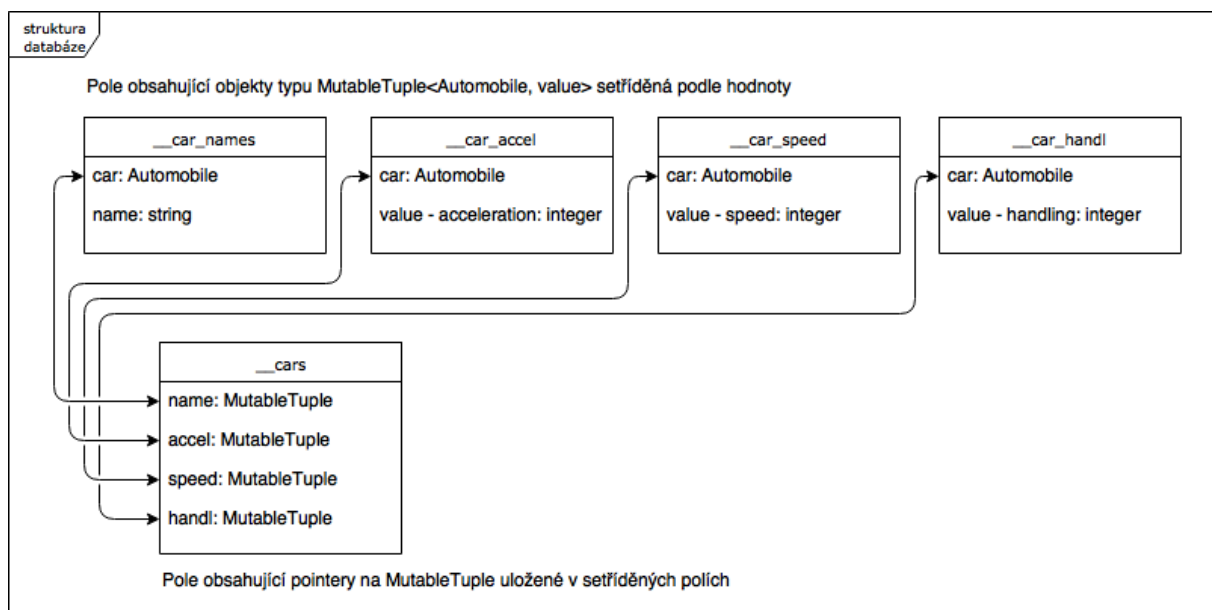
Implementoval jsem jednoduchou webovou aplikaci pro demonstraci použití top-k treshold algoritmu v jazyce Python. Mé řešení obsahuje navíc i naivní algoritmus a faginův algoritmus. Je tak možné přehledně vyzkoušet, jak jednotlivé algoritmy fungují a porovnat jejich efektivitu.

## ZPŮSOB ŘEŠENÍ

Jako databázi využívám čtyři pole obsahující KeyValuePair ve tvaru <pointer na třídu Automobile, numerická hodnota>. Všechna data jsou nahrána přímo v paměti počítače. Jednotlivá **pole jsou seříděna podle hodnoty při načtení dat do databáze**. Jednotlivé dotazy pořadí prvků v polích nemění.

Implementace algoritmů pro top-k query se nachází v souboru /source/classes/CarDatabase.py.

Následující obrázek popisuje hlavní strukturu uložení dat (jednotlivá seříděná pole):



Samotný top-k treshold algoritmus využívá **minimovou haldu** k uchovávání objektů a efektivní odstraňování prvků s nejnižším agregačním score. Mé řešení obsahuje **několik agregačních funkcí**.

TOP-K THRESHOLD ALGORITHMUS FUNGUJE NÁSLEDUJÍCÍM ZPŮSOBEM:

1. Nastavím treshold  $t$  jako agregační hodnotu prvků zjištěných v přístupu do databáze.
  2. Náhodným přístupem vypočítám skóre nalezených objektů dle agregační funkce.
  3. V minimové haldě si uchovávám seznam top-k objektů na které jsem doposud narazil.
  4. Zastavím se, pokud jsou skóre top-k objektů v haldě větší nebo rovny tresholdu.
- Konec – vrátím top-k objektů které jsem prozkoumal.

## IMPLEMENTACE

Aplikace je napsaná v programovacím jazyce Python s využitím webového frameworku FLASK pro generování webové stránky. Samotná webová stránka je HTML dokument s několika funkcemi psanými v jazyce JavaScript (kontrola vstupu, animace,...).

### SPUŠTĚNÍ APLIKACE

Pro správný běh aplikace je potřeba mít nainstalovaný jazyk python3 a modul flask. Instalaci modulu je možné provést například pomocí příkazu:

```
pip3 install flask
```

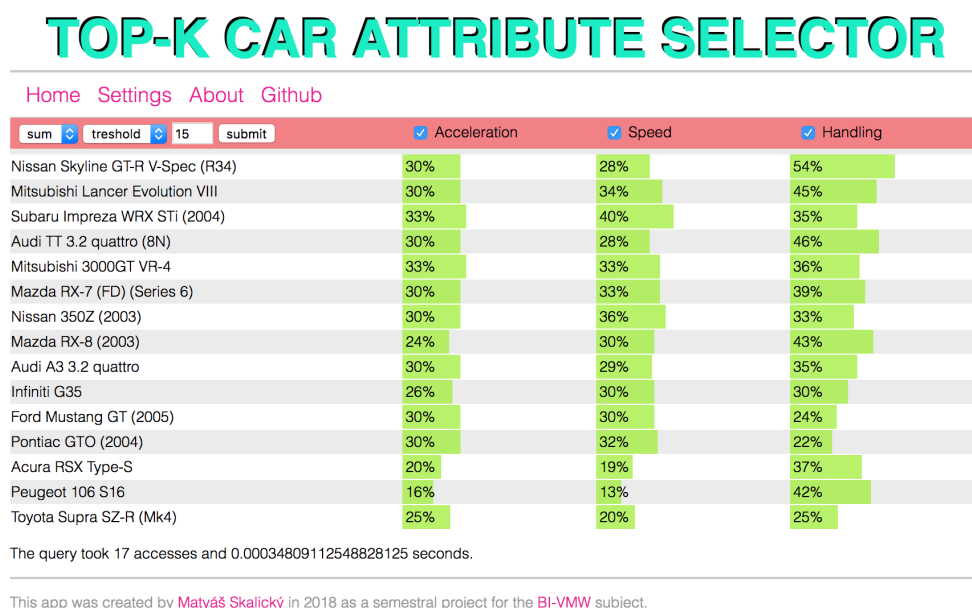
Samotná aplikace se spouští příkazem:

```
python3 run.py
```

Webový server běží ve výchozí konfiguraci na adrese <http://127.0.0.1:5000/>.

## PŘÍKLAD VÝSTUPU

Na screenshotu vidíme ukázkový dotaz nad databází aut ze hry Need for Speed: Most Wanted (pozn. data lze změnit v záložce Settings). Na query byla využita jako agregační funkce suma a byl aplikován top-k threshold algoritmus. Uživatel si v tomto případě přál získat 15 objektů a seřadit je podle všech 3 parametrů (akcelerace, rychlost, ovládání) s pomocí agregační funkce suma. Pod vypsanou tabulkou přehledně vidíme, že query trvala 0.0003s a využila celkem 17 přístupů do databáze.

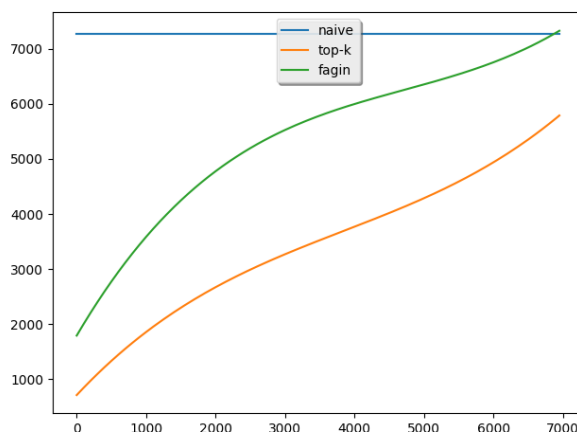


## EXPERIMENTÁLNÍ SEKCE

V následující sekci se pokusím porovnat Top-k, Faginův a naivní algoritmus na základě počtu přístupů do databáze, který je klíčový. Pro přehlednost dodávám porovnání i na reálné době běhu algoritmu. Pro experimentální účely jsem do aplikace přidal navíc databázi o 7 000 a 160 000 položkách s náhodně vygenerovanými hodnotami.

### POROVNÁNÍ POČTU PŘÍSTUPŮ DO DATABÁZE

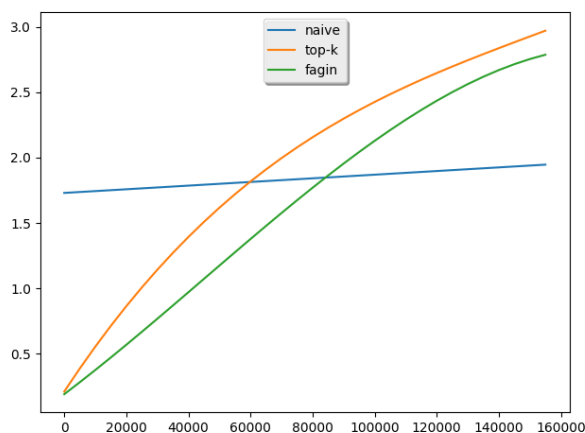
Je vidět, že algoritmus top-k využívá ve srovnání s Faginovým algoritmem značně menší počet přístupů do databáze. V prostředí s pomalým přístupem k databázi se proto i za cenu větší výpočetních nároků vyplatí a dotazy budou rychlejší.



OSA X – VELIKOST DOTAZU (K), OSA Y – POČET PŘÍSTUPŮ DO DATABÁZE

### POROVNÁNÍ DOBY BĚHU

Srovnání doby běhu algoritmů na databázi o velikosti 160 000 objektů. Je vidět, že ačkoliv použije algoritmus top-k nižší množství náhodných přístupů, vyžaduje udržování minimální haldy pro uchování top-k elementů delší výpočetní čas. Pokud mám rychlý přístup k databázi s náhodným přístupem, vyplatí se využít Faginův algoritmus.



OSA X – VELIKOST DOTAZU (K), OSA Y – ČASOVÁ DÉLKA DOTAZU

## DISKUZE

Projekt, který jsem vypracoval zahrnuje vlastní implementaci Top-K threshold algoritmu. Využíval jsem zjednodušený scénář, kdy všechna data nahraji přímo do paměti a nemusím proto řešit databázový software ani komplikované dotazování. V reálném případě použití by bylo potřeba zakomponovat podobný model do databázové vrstvy aplikace. V mém projektu jsem neměl prostor řešit optimalizaci jednotlivých algoritmů a datových struktur. Další práce v této oblasti by mohla značně zrychlit celé mé řešení. Python rovněž není nejrychlejším jazykem, tudíž aplikace například v C++ by byla bez pochyb rychlejší, ale náročnější na implementaci.

## ZÁVĚR

Mým zadáním bylo implementovat a ověřit funkci Top-K threshold algoritmu. Zadané téma mě začalo velmi bavit a užíval jsem si i grafickou stránku aplikace. Nakonec mě navíc zajímalo, jak je na tom top-k threshold algoritmus v porovnání s ostatními algoritmy. Do mé aplikace jsem tedy přidal i Faginův a naivní algoritmus pro porovnání funkcionality a rovněž několik agregačních funkcí. Celou aplikaci jsem vybudoval ve webovém frameworku Flask, což se ukázalo být nesmírně přínosné. Flask jsem následně využil i pro další semestrální práci.