

Univerzita Jana Evangelisty Purkyně v Ústí
nad Labem

Přírodovědecká fakulta

Katedra informatiky

OLAP a DuckDB

Seminární práce

Rok: 2025

Vypracoval: Valdemar Pospíšil

Obsah

1	Úvod do DuckDB	2
1.1	Instalace DuckDB na Linux	2
1.2	Vytvoření a použití databáze v DuckDB	2
2	Dataset	3
3	Vytvoření datového skladu	4
3.1	Dimenzní tabulky	4
3.2	Faktová tabulka	4
3.3	SQL kód pro vytvoření struktury datového skladu	4
3.4	Naplnění dimenzních a faktové tabulky	5
4	Analytické dotazy	6
4.1	Rozložení pozorování UFO podle států	6
4.2	Distribuce délky pozorování	7
4.3	Pozorování v průběhu dne a roku	8
4.4	Analýza popisů pozorování	9
5	Data mining	11
5.1	Shlukování (Clustering)	11
5.2	Asociační pravidla	13
6	Závěr	13

1 Úvod do DuckDB

DuckDB je moderní analytický databázový systém typu OLAP (Online Analytical Processing), který je navržen pro rychlé dotazování a analýzu velkých objemů dat. Na rozdíl od tradičních databázových systémů jako PostgreSQL nebo MySQL, které jsou primárně zaměřeny na OLTP (Online Transaction Processing), je DuckDB optimalizován pro analytické dotazy, které zpracovávají velké množství dat a provádějí agregace.

Mezi hlavní výhody DuckDB patří:

- **Jednoduchost** - lze používat jako vestavěnou databázi bez nutnosti spouštění samostatného serveru
- **Rychlost** - optimalizovaný pro analytické dotazy a operace OLAP
- **Sloupcově orientované úložiště** - efektivní pro analytické dotazy, které často pracují s omezeným počtem sloupců
- **Integrace s Pythonem** - snadné použití v datových analýzách a vědeckých výpočtech
- **SQL kompatibilita** - plná podpora standardních SQL příkazů

1.1 Instalace DuckDB na Linux

Pro instalaci DuckDB na operačním systému Void Linux jsem použil správce balíčků pip pro Python. Instalace je jednoduchá a přímočará:

```
1 # Aktualizace pip
2 pip install --upgrade pip
3
4 # Instalace DuckDB
5 pip install duckdb
```

Listing 1: Instalace DuckDB pomocí pip

Na Linux možné si nainstalovat duckdb jakožto samostatný databázový server pomocí jednoduchého příkazu:

```
1 # Instalace duckdb
2 curl https://install.duckdb.org | sh
```

Listing 2: Instalace duckdb na linux zařízeních pomocí příkazu

1.2 Vytvoření a použití databáze v DuckDB

Vytvoření databáze v DuckDB je velmi jednoduché. Databázi lze vytvořit a používat přímo z Pythonu bez nutnosti spouštět samostatný databázový server:

```
1 import duckdb
2
3 # Vytvoreni pripojeni k databazi (pokud soubor neexistuje, bude vytvoren)
4 con = duckdb.connect("ufo.db")
5
6 # Vytvoření tabulky
7 con.execute("""
```

```

8      CREATE TABLE example_table (
9          id INTEGER PRIMARY KEY,
10         name VARCHAR(100),
11         value DOUBLE
12     )
13     """)
14
15 # Vložení dat
16 con.execute("INSERT INTO example_table VALUES (1, 'Test', 10.5)")
17
18 # Dotaz
19 result = con.execute("SELECT * FROM example_table").fetchall()
20 print(result)
21
22 # Uzavření připojení
23 con.close()

```

Listing 3: Vytvoření a připojení k databázi v DuckDB

Databáze je uložena v souboru, který je specifikován při vytváření připojení. Pokud soubor neexistuje, DuckDB ho automaticky vytvoří. Tím se DuckDB liší od většiny databázových systémů, které vyžadují spuštění serveru.

2 Dataset

Pro projekt jsem si stáhl dataset **UFO Sightings** z Kaggle (<https://www.kaggle.com/datasets/sahityasetu/ufo-sightings>). Dataset obsahuje informace o pozorování UFO včetně času, místa, popisu a tvaru objektu.

Dataset obsahuje následující hlavní sloupce:

- **Date_time** - datum a čas pozorování
- **Date_documented** - datum, kdy bylo pozorování zaznamenáno
- **Year, Month, Hour** - extrahované časové údaje
- **Season** - roční období
- **Country_Code, Country, Region, Locale** - údaje o lokalitě
- **Latitude, Longitude** - geografické souřadnice
- **UFO_shape** - tvar pozorovaného objektu
- **Length_of_encounter_seconds** - doba trvání pozorování v sekundách
- **Encounter_Duration** - textový popis délky pozorování
- **Description** - podrobný popis pozorování

Pro načtení dat z CSV souboru do databáze jsem použil příkaz COPY:

```

1 -- Vytvoření základní tabulky z CSV
2 CREATE TABLE sightings (
3     Date_time TIMESTAMP,
4     Date_documented DATE,
5     Year INT,

```

```

6      Month INT,
7      Hour INT,
8      Season VARCHAR(20),
9      Country_Code VARCHAR(10),
10     Country VARCHAR(100),
11     Region VARCHAR(100),
12     Locale VARCHAR(100),
13     Latitude DOUBLE,
14     Longitude DOUBLE,
15     UFO_shape VARCHAR(50),
16     Length_of_encounter_seconds BIGINT,
17     Encounter_Duration VARCHAR(50),
18     Description TEXT
19 );
20
21 -- Načtení dat z CSV souboru s oddelovacem čarkou
22 COPY sightings FROM 'ufo_sightings.csv' (DELIMITER ',', HEADER);

```

Listing 4: Načtení dat z CSV souboru

3 Vytvoření datového skladu

Vytvořil jsem datovou strukturu ve tvaru **hvězdy (star schema)** s jednou faktovou tabulkou a třemi dimenzními tabulkami. Tato struktura je typická pro OLAP analýzy a umožňuje efektivní dotazování na různých úrovních granularity.

3.1 Dimenzní tabulky

Dimenzní tabulky obsahují popisné atributy, které jsou využívány pro filtrování a seskupování dat:

- **dim_ufo**: obsahuje různé tvary UFO.
- **dim_time**: obsahuje časové údaje (rok, měsíc, den, hodina).
- **dim_location**: obsahuje údaje o zemi, regionu a lokalitě.

3.2 Faktová tabulka

Faktová tabulka obsahuje měřené hodnoty (fakta) a cizí klíče, které odkazují na dimenzní tabulky:

- **fact_sightings**: obsahuje jednotlivá pozorování, která jsou propojena na dimenze přes cizí klíče.

3.3 SQL kód pro vytvoření struktury datového skladu

Nejprve jsem musel vytvořit sekvence pro generování primárních klíčů:

```

1 -- Vytvoření sekvencí pro primární klíče
2 CREATE SEQUENCE time_id_seq START 1;
3 CREATE SEQUENCE location_id_seq START 1;
4 CREATE SEQUENCE ufo_id_seq START 1;

```

```
5 CREATE SEQUENCE sighting_id_seq START 1;
```

Listing 5: Vytvoření sekvencí

Poté jsem vytvořil dimenzní tabulky:

```
1 -- Dimenzní tabulka pro čas
2 CREATE TABLE dim_time (
3     time_id INTEGER DEFAULT nextval('time_id_seq') PRIMARY KEY,
4     date_time TIMESTAMP,
5     date_documented DATE,
6     year INT,
7     month INT,
8     hour INT,
9     season VARCHAR(20)
10 );
11
12 -- Dimenzní tabulka pro lokality
13 CREATE TABLE dim_location (
14     location_id INTEGER DEFAULT nextval('location_id_seq') PRIMARY KEY,
15     country_code VARCHAR(10),
16     country VARCHAR(100),
17     region VARCHAR(100),
18     locale VARCHAR(100),
19     latitude DOUBLE,
20     longitude DOUBLE
21 );
22
23 -- Dimenzní tabulka pro tvary UFO
24 CREATE TABLE dim_ufo (
25     ufo_id INTEGER DEFAULT nextval('ufo_id_seq') PRIMARY KEY,
26     ufo_shape VARCHAR(50)
27 );
```

Listing 6: Vytvoření dimenzních tabulek

Následně jsem vytvořil faktovou tabulku, která obsahuje cizí klíče odkazující na dimenzní tabulky:

```
1 -- Faktová tabulka
2 CREATE TABLE fact_sightings (
3     sighting_id INTEGER DEFAULT nextval('sighting_id_seq') PRIMARY KEY,
4     time_id INT REFERENCES dim_time(time_id),
5     location_id INT REFERENCES dim_location(location_id),
6     ufo_id INT REFERENCES dim_ufo(ufo_id),
7     length_of_encounter_seconds BIGINT,
8     encounter_duration VARCHAR(50),
9     description TEXT
10 );
```

Listing 7: Vytvoření faktové tabulky

3.4 Naplnění dimenzních a faktové tabulky

Po vytvoření tabulek jsem do nich naplnil data z původní tabulky sightings:

```
1 -- Naplnění dimenzní tabulky času
2 INSERT INTO dim_time (date_time, date_documented, year, month, hour,
3     season)
4 SELECT DISTINCT Date_time, date_documented, Year, Month, Hour, Season
```

```

4 FROM sightings;
5
6 -- Naplneni dimenzni tabulky lokality
7 INSERT INTO dim_location (country_code, country, region, locale,
8     latitude, longitude)
9 SELECT DISTINCT Country_Code, Country, Region, Locale, latitude,
10     longitude
11 FROM sightings;
12
13 -- Naplneni dimenzni tabulky UFO
14 INSERT INTO dim_ufo (ufo_shape)
15 SELECT DISTINCT UFO_shape
16 FROM sightings;

```

Listing 8: Naplnění dimenzních tabulek

Nakonec jsem naplnil faktovou tabulku, přičemž jsem vytvořil spojení s dimenzními tabulkami:

```

1 -- Naplneni faktove tabulky
2 INSERT INTO fact_sightings (time_id, location_id, ufo_id,
3     length_of_encounter_seconds, encounter_duration, description)
4 SELECT
5     t.time_id,
6     l.location_id,
7     u.ufo_id,
8     s.length_of_encounter_seconds,
9     s.Encounter_Duration,
10    s.Description
11 FROM sightings s
12 JOIN dim_time t ON s.Date_time = t.date_time
13                 AND s.date_documented = t.date_documented
14                 AND s.Year = t.year
15                 AND s.Month = t.month
16                 AND s.Hour = t.hour
17                 AND s.Season = t.season
18 JOIN dim_location l ON s.Country_Code = l.country_code
19                     AND s.Country = l.country
20                     AND s.Region = l.region
21                     AND s.Locale = l.locale
22                     AND s.latitude = l.latitude
23                     AND s.longitude = l.longitude
24 JOIN dim_ufo u ON s.UFO_shape = u.ufo_shape;

```

Listing 9: Naplnění faktové tabulky

4 Analytické dotazy

Po vytvoření datového skladu jsem provedl několik analytických dotazů k získání zajímavých informací o pozorováních UFO. Dotazy byly implementovány v Pythonu s využitím knihovny DuckDB a výsledky byly vizualizovány pomocí knihoven jako **Matplotlib**, **Seaborn**, **Tabulate**, **Plotly** a **Folium**.

4.1 Rozložení pozorování UFO podle států

První analýza se zaměřila na počet pozorování UFO v jednotlivých státech USA. Implementoval jsem ji pomocí následujícího SQL dotazu:

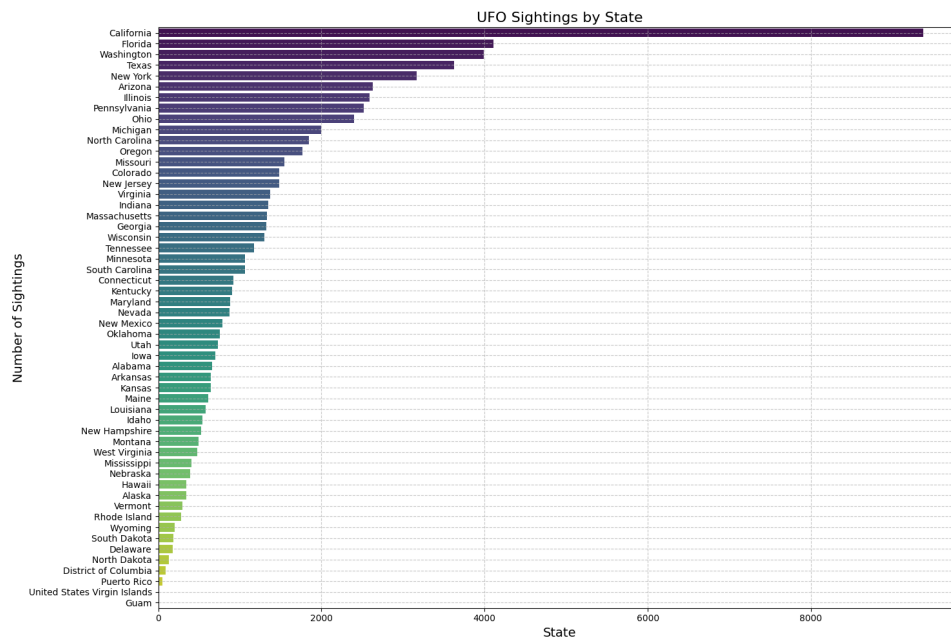
```

1 WITH state_sightings AS (
2     SELECT
3         l.region,
4         f.length_of_encounter_seconds
5     FROM fact_sightings f
6     JOIN dim_location l ON f.location_id = l.location_id
7     WHERE l.country_code = 'USA' AND l.region IS NOT NULL
8 )
9 SELECT
10     region AS state,
11     COUNT(*) AS sightings_count,
12     AVG(length_of_encounter_seconds) AS avg_encounter_seconds,
13     MEDIAN(length_of_encounter_seconds) AS median_encounter_seconds
14 FROM state_sightings
15 GROUP BY region
16 ORDER BY sightings_count DESC

```

Listing 10: SQL dotaz pro počet pozorování podle států

Výsledky byly vizualizovány pomocí horizontálního sloupčového grafu a interaktivní choroplethové mapy pomocí knihovny Plotly.



Obrázek 1: Počet pozorování UFO podle států

4.2 Distribuce délky pozorování

Další analýza se zaměřila na délku pozorování UFO podle tvaru objektu. Implementoval jsem ji pomocí následujícího SQL dotazu:

```

1 SELECT
2     d.ufo_shape,
3     MEDIAN(f.length_of_encounter_seconds)::INTEGER AS median_seconds,
4     COUNT(*) AS sightings_count,
5     AVG(f.length_of_encounter_seconds)::INTEGER AS avg_seconds
6 FROM fact_sightings f
7 JOIN dim_ufo d ON f.ufo_id = d.ufo_id

```



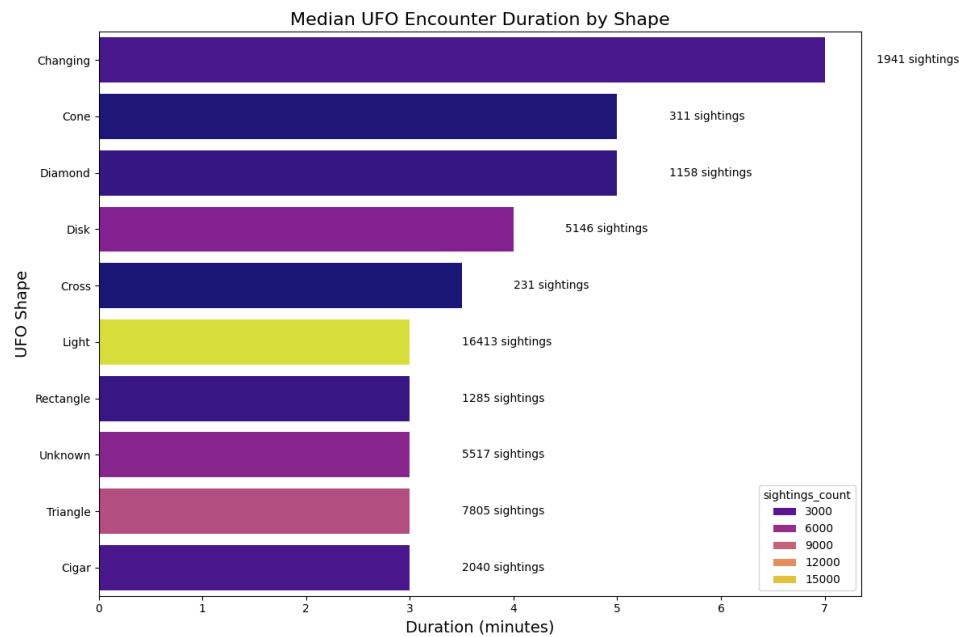
```

8 WHERE f.length_of_encounter_seconds IS NOT NULL AND f.
   length_of_encounter_seconds > 0
9 GROUP BY d.ufo_shape
10 HAVING COUNT(*) > 30
11 ORDER BY median_seconds DESC
12 LIMIT 10

```

Listing 11: SQL dotaz pro délku pozorování podle tvaru UFO

Výsledky byly vizualizovány pomocí horizontálního sloupcového grafu, který ukazuje medián délky pozorování pro různé tvary UFO:



Obrázek 2: Délka pozorování UFO podle tvaru objektu

4.3 Pozorování v průběhu dne a roku

Třetí analýza zkoumala, jak se počet pozorování UFO mění v průběhu dne a roku. Použil jsem následující SQL dotaz:

```

1 SELECT
2     CASE t.month
3         WHEN 1 THEN 'January'
4         WHEN 2 THEN 'February'
5         WHEN 3 THEN 'March'
6         WHEN 4 THEN 'April'
7         WHEN 5 THEN 'May'
8         WHEN 6 THEN 'June'
9         WHEN 7 THEN 'July'
10        WHEN 8 THEN 'August'
11        WHEN 9 THEN 'September'
12        WHEN 10 THEN 'October'
13        WHEN 11 THEN 'November'
14        WHEN 12 THEN 'December'
15    END AS month,
16    CASE
17        WHEN t.hour >= 5 AND t.hour < 12 THEN 'Morning'

```

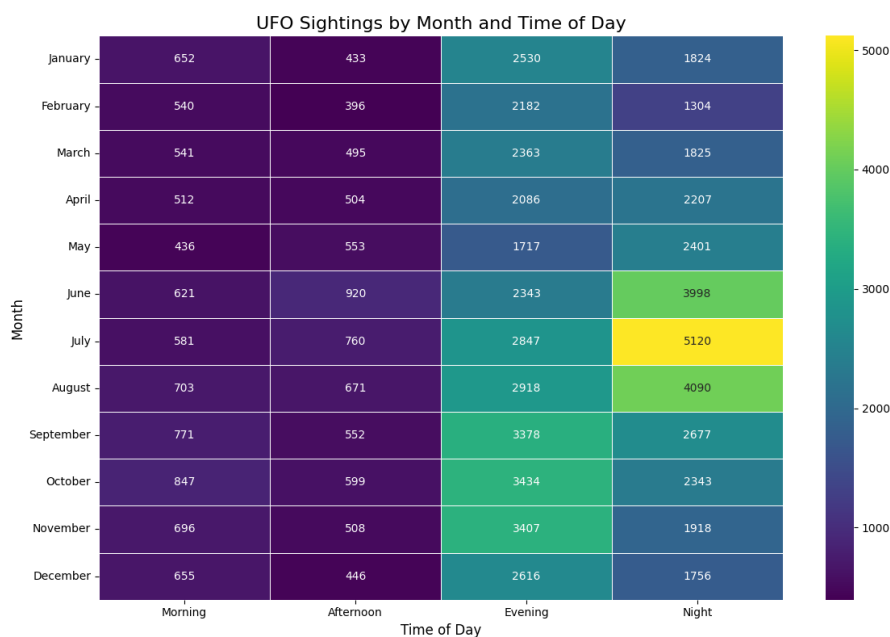
```

18     WHEN t.hour >= 12 AND t.hour < 17 THEN 'Afternoon'
19     WHEN t.hour >= 17 AND t.hour < 22 THEN 'Evening'
20     ELSE 'Night'
21   END AS time_of_day,
22   COUNT(*) AS sightings_count
23 FROM fact_sightings f
24 JOIN dim_time t ON f.time_id = t.time_id
25 GROUP BY t.month, time_of_day
26 ORDER BY
27   t.month,
28   CASE
29     WHEN time_of_day = 'Morning' THEN 1
30     WHEN time_of_day = 'Afternoon' THEN 2
31     WHEN time_of_day = 'Evening' THEN 3
32     WHEN time_of_day = 'Night' THEN 4
33   END

```

Listing 12: SQL dotaz pro počet pozorování podle měsíce a denní doby

Výsledky byly vizualizovány pomocí teplotní mapy (heatmap), která ukazuje počet pozorování v různých měsících a denních dobách:



Obrázek 3: Distribuce pozorování UFO podle měsíce a denní doby

4.4 Analýza popisů pozorování

Provedl jsem také analýzu textových popisů pozorování UFO, abych identifikoval nejčastější slova spojená s jednotlivými tvary objektů. Implementace zahrnovala předzpracování textu a využití knihovny Counter z modulu collections:

```

1 # SQL prikaz pro získání tvaru a popisku ufo
2 query = """
3 SELECT
4     d.ufo_shape,
5     f.description

```

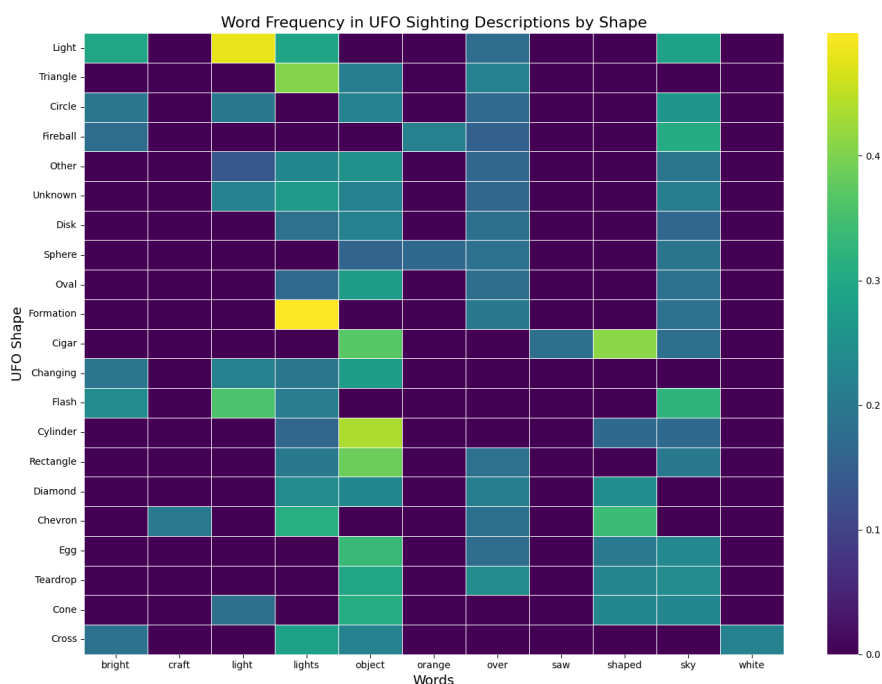
```

6 FROM fact_sightings f
7 JOIN dim_ufo d ON f.ufo_id = d.ufo_id
8 WHERE f.description IS NOT NULL
9       AND LENGTH(f.description) > 20
10      AND d.ufo_shape IS NOT NULL
11      AND d.ufo_shape != 'unknown'
12 LIMIT 10000
13 """
14
15 # Funkce pro cistení a tokenizaci textu
16 def clean_text(text):
17     if not isinstance(text, str):
18         return []
19     # Prevod na mala pismena
20     text = text.lower()
21     # Odstraneni specialnich znaku a cislic
22     text = re.sub(r'[^a-zA-Z\s]', '', text)
23     # Tokenizace
24     tokens = text.split()
25     # Odstraneni stop slov (bezn ch slov, kter neprinaseji velky
26     # význam)
27     stop_words = {'the', 'and', 'to', 'of', 'a', 'in', 'that', 'it', '...', 'is', 'was', 'are', 'were', 'do', 'does', 'did', 'have', 'has', 'had', 'be', 'been', 'am', 'are', 'was', 'were', 'do', 'does', 'did', 'have', 'has', 'had', 'be', 'been', 'am'}
28     tokens = [token for token in tokens if token not in stop_words and len(token) > 2]
29     return tokens

```

Listing 13: Kód pro analýzu textových popisů

Pro vizualizaci výsledků jsem vytvořil teplotní mapu (heatmap) a word cloudy pro nejčastější tvary UFO:



Obrázek 4: Analýza slov v popisech pozorování UFO podle tvaru objektu

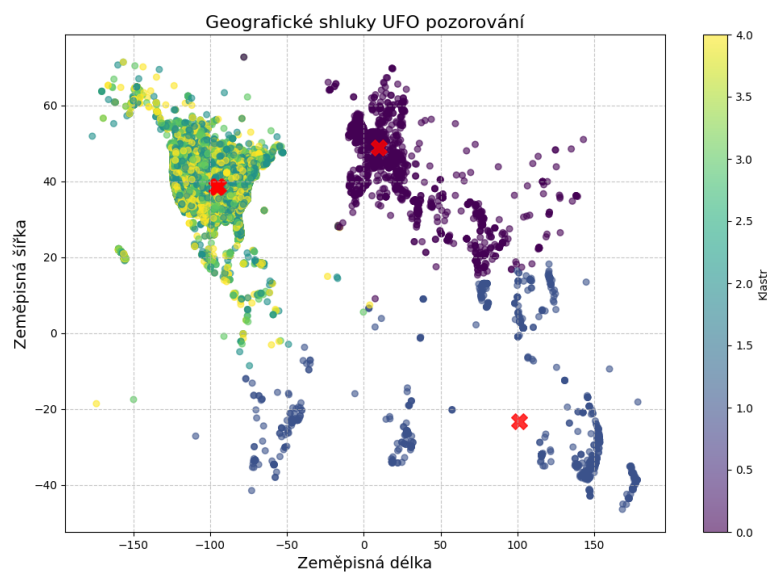

```

17     inertia.append(kmeans.inertia_)
18
19 # Zvolení konečného počtu shluku
20 k = 5
21 kmeans = KMeans(n_clusters=k, random_state=42, n_init=10)
22 clusters = kmeans.fit_predict(scaled_features)

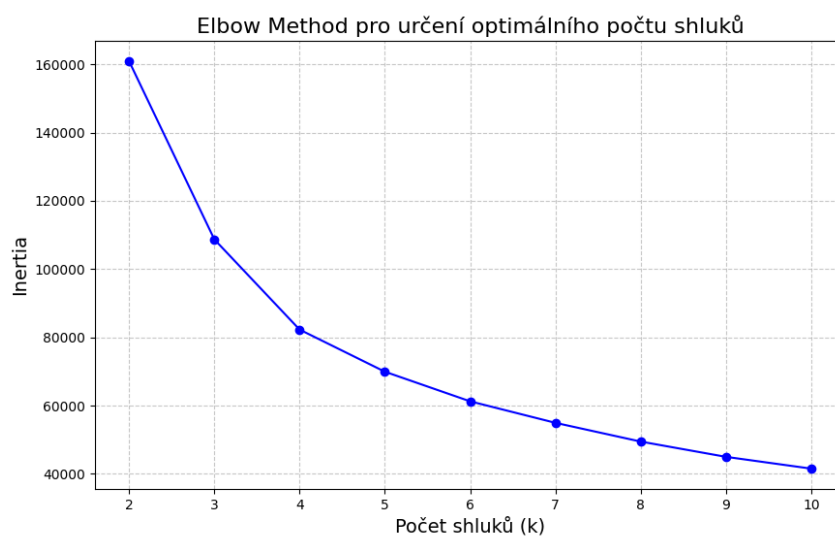
```

Listing 14: Implementace K-means shlukování

Výsledky shlukování byly vizualizovány na mapě, která ukazuje rozložení jednotlivých shluků a jejich centra:



Obrázek 6: Mapa shluků pozorování UFO



Obrázek 7: Graf metody elbow pro určení optimálního počtu shluků

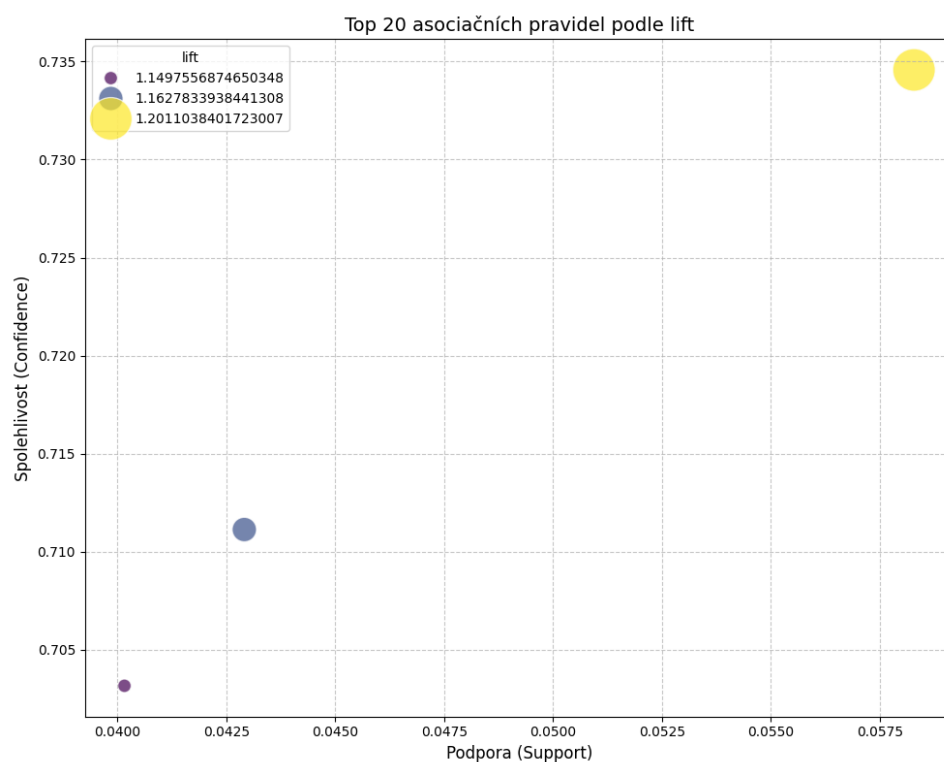
5.2 Asociační pravidla

Pro identifikaci zajímavých asociací v datech jsem použil algoritmus Apriori pro dolování asociačních pravidel. Tato metoda umožňuje objevit, které kombinace vlastností se často vyskytují společně:

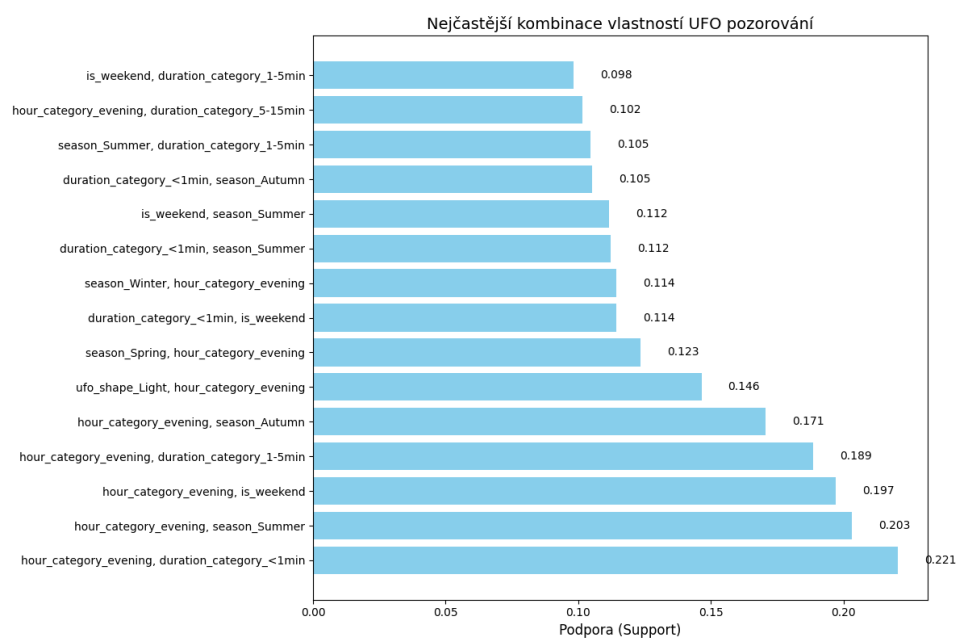
```
1 # Příprava kategorických příznaků
2 categorical_features = ['ufo_shape', 'season', 'is_weekend']
3
4 # Přidání diskretizovaných numerických příznaků
5 df['duration_category'] = pd.cut(
6     df['length_of_encounter_seconds'],
7     bins=[0, 60, 300, 900, 3600, 86400],
8     labels=['<1min', '1-5min', '5-15min', '15min-1hr', '>1hr']
9 )
10
11 df['hour_category'] = pd.cut(
12     df['hour'],
13     bins=[-0.1, 5.9, 11.9, 17.9, 23.9],
14     labels=['night', 'morning', 'afternoon', 'evening']
15 )
16
17 categorical_features.extend(['duration_category', 'hour_category'])
18
19 # One-hot kodování kategorických příznaků
20 encoded_df = pd.get_dummies(df[categorical_features])
21
22 # Nalezení častých itemsetů
23 min_support = 0.03 # Minimální podpora
24 frequent_itemsets = apriori(encoded_df, min_support=min_support,
25                             use_colnames=True)
26 frequent_itemsets['length'] = frequent_itemsets['itemsets'].apply(lambda
27     x: len(x))
28
29 # Generování asociacních pravidel
30 min_threshold = 0.7 # Minimální spolehlivost
31 rules = association_rules(frequent_itemsets, metric="confidence",
32                           min_threshold=min_threshold)
33
34 # Seřazení pravidel podle liftu
35 rules = rules.sort_values('lift', ascending=False)
```

Listing 15: Implementace dolování asociačních pravidel

Výsledky byly vizualizovány pomocí grafu scatter plot, který ukazuje vztah mezi podporou, spolehlivostí a liftem jednotlivých pravidel:



Obrázek 8: Asociační pravidla v datech o pozorování UFO



Obrázek 9: Nejčastější kombinace vlastností UFO pozorování

6 Závěr

V tomto projektu jsem úspěšně demonstroval použití DuckDB pro OLAP analýzy na reálném datasetu o pozorováních UFO. DuckDB se ukázal jako výkonný a snadno použitelný nástroj pro analytické zpracování dat, který kombinuje jednoduchost použití s vysokým výkonem.

Hlavní přínosy projektu:

- **Praktické použití DuckDB** - instalace, vytvoření databáze a provedení komplexních analytických dotazů
- **Vytvoření dimenzionálního modelu** - implementace struktury hvězdy (star schema) pro efektivní analytické dotazování
- **Implementace OLAP dotazů** - využití SQL a Python API pro analýzu dat z různých perspektiv
- **Aplikace metod data miningu** - využití shlukování a asociačních pravidel pro objevení skrytých vzorů v datech
- **Vizualizace výsledků** - vytvoření rozmanitých vizualizací pro lepší interpretaci získaných poznatků

Na základě analýz jsem zjistil několik zajímavých vzorů v datech o pozorováních UFO:

- Kalifornie, Florida a Washington patří mezi státy s nejvyšším počtem hlášených pozorování UFO
- Většina pozorování se odehrává v nočních hodinách, přičemž v letních měsících je počet pozorování obecně vyšší
- Existují určité geografické shluky, které mají specifické charakteristiky z hlediska délky trvání a typu objektu
- Z analýzy textových popisů vyplývá, že určité tvary UFO jsou častěji spojovány s konkrétními slovy a popisy

DuckDB se osvědčil jako ideální nástroj pro tento typ analytických úloh díky své jednoduchosti, výkonu a schopnosti efektivně zpracovávat analytické dotazy bez nutnosti složité konfigurace.

Projekt je kompletně připraven v repozitáři a doplněn o vizualizace výsledků. Zdrojový kód je organizován do samostatných skriptů pro jednotlivé analýzy, což usnadňuje údržbu a rozšíření projektu v budoucnu.