

# **Izveštaj o realizaciji projekta: Analiza letova u MongoDB-u**

Oktober 2025.

## **Contents**

<b>1</b>	<b>Uvod</b>	<b>2</b>
<b>2</b>	<b>Inicijalna priprema skupa podataka</b>	<b>2</b>
<b>3</b>	<b>Analiza skupa podataka</b>	<b>3</b>
<b>4</b>	<b>Analiza performansi upita</b>	<b>3</b>
<b>5</b>	<b>Optimizacija performansi upita</b>	<b>4</b>
<b>6</b>	<b>Prerada upita</b>	<b>4</b>
<b>7</b>	<b>Uporedna analiza performansi</b>	<b>4</b>
<b>8</b>	<b>Vizualizacija rezultata</b>	<b>5</b>
<b>9</b>	<b>Zaključak</b>	<b>5</b>

# 1 Uvod

Ovaj izveštaj sažima ključne elemente za realizaciju projekta analize podataka o letovima u MongoDB-u, baziran na zahtevima za predispitnu obavezu. Projekt koristi podatke iz CSV fajlova (US\_flights\_2023.csv, Cancelled\_Diverted\_2023.csv, airports\_geolocation.csv, weather\_meteo\_by\_airport.csv, airports.csv, airport\_frequencies.csv, runways.csv) i uključuje pripremu podataka, MongoDB upite, optimizaciju performansi, i vizualizaciju rezultata. Izveštaj pokriva sedam faza projekta, sa fokusom na realno stanje podataka i praktičnu implementaciju.

## 2 Inicijalna priprema skupa podataka

Podaci su učitani iz CSV fajlova koristeći Python biblioteku pymongo. Strukturirani su prema sledećoj logičkoj šemi:

- **flights:** Podaci o letovima (FlightDate, Airline, Dep\_Airport, Dep\_Delay, Delay\_Weather, itd.).
- **airports:** Informacije o aerodromima (iata\_code, ident, type, elevation\_ft).
- **weather:** Vremenski podaci (airport\_id, time, prcp, wspd).
- **frequencies:** Frekvencije aerodroma (airport\_ident, type, frequency\_mhz).
- **runways:** Podaci o pistama (airport\_ident, length\_ft, surface).
- **geolocation:** Geolokacija aerodroma (IATA\_CODE, CITY, LATITUDE).

Python kod za učitavanje podataka:

```
1 import pandas as pd
2 import pymongo
3 import os
4 client = pymongo.MongoClient("mongodb://localhost:27017/")
5 db = client["flights_db"]
6 folder = r"C:\Users\tijan\Desktop\sbp"
7 files = {
8     "US_flights_2023": "US_flights_2023.csv",
9     "Cancelled_Diverted_2023": "Cancelled_Diverted_2023.csv",
10    "Airports_Geolocation": "airports_geolocation.csv",
11    "Weather_Meteo_by_Airport": "weather_meteo_by_airport.csv",
12    "Airports": "airports.csv",
13    "Airport_frequencies": "airport_frequencies.csv",
14    "Runways": "runways.csv"
15 }
16 for name, filename in files.items():
17     df = pd.read_csv(os.path.join(folder, filename))
18     df = df.where(pd.notnull(df), None)
19     collection = db[name.lower()]
20     collection.insert_many(df.to_dict("records"))
```

**Preporuke:** Testiraj učitavanje sa manjim skupom podataka (npr. nrows=10000) zbog velikih fajlova (6.7M redova u US\_flights\_2023.csv). Proveri nedostajuće vrednosti i osiguraj da MongoDB server radi.

### 3 Analiza skupa podataka

Deset kompleksnih agregacionih upita je razvijeno za analizu kašnjenja letova, povezujući kolekcije preko `$lookup`, `$unwind`, `$match`, `$group`, i `$sort`. Primer upita (aerodromi sa kašnjenjem > 15 minuta, padavine > 10 mm, piste > 10000 ft):

```
1 db.flights.aggregate([
2   {"$match": {"departure.delay": {"$gt": 15}}},
3   {"$lookup": {
4     "from": "weather", "localField": "departure.airport",
5     "foreignField": "airport_id", "as": "weather_info"
6   }},
7   {"$unwind": "$weather_info"},
8   {"$match": {"weather_info.weather.prcp": {"$gt": 10}}},
9   {"$lookup": {
10    "from": "airports", "localField": "departure.airport",
11    "foreignField": "iata_code", "as": "airport_info"
12  }},
13  {"$unwind": "$airport_info"},
14  {"$match": {"airport_info.type": "large_airport"}},
15  {"$lookup": {
16    "from": "runways", "localField": "airport_info.ident",
17    "foreignField": "airport_ident", "as": "runway_info"
18  }},
19  {"$unwind": "$runway_info"},
20  {"$match": {"runway_info.length_ft": {"$gt": 10000}}},
21  {"$group": {
22    "_id": "$departure.airport",
23    "avg_delay": {"$avg": "$departure.delay"},
24    "count": {"$sum": 1}
25  }},
26  {"$sort": {"avg_delay": -1}},
27  {"$limit": 5}
28 ])
```

**Preporuke:** Pokreni upite u MongoDB Compass-u ili Python-u. Proveri rezultate i osiguraj da sve kolekcije imaju podatke.

### 4 Analiza performansi upita

Performanse upita su analizirane koristeći `explain()`. Uska grla uključuju:

- **Puni scanovi:** Bez indeksa na poljima poput `departure.airport`, upiti skeniraju celu kolekciju.
- **Spori \$lookup:** Velike kolekcije (runways: 41k redova) usporavaju povezivanje.
- **Veliki rezultati:** `$unwind` generiše mnogo dokumenata.

Primer merenja vremena u Python-u:

```
1 import time
2 start = time.time()
3 result = list(db.flights.aggregate(pipeline))
4 end = time.time()
5 print(f"Vreme šizvravanja: {end - start} sekundi")
```

**Preporuke:** Koristi `explain("executionStats")` za detalje o scan-ovima. Izmeri vreme za svaki upit (npr. 10 ponavljanja).

## 5 Optimizacija performansi upita

Optimizacija uključuje:

- **Indeksiranje:** Kreiranje indeksa na ključnim poljima.
- **Restrukturiranje šeme:** Denormalizacija (npr. ugnežđavanje `weather_summary` u `flights`).

Primer indeksiranja:

```
1 db.flights.create_index([("departure.airport", 1)])
2 db.airports.create_index([("iata_code", 1), ("ident", 1)])
3 db.weather.create_index([("airport_id", 1), ("time", 1)])
```

**Preporuke:** Testiraj performanse nakon indeksiranja. Denormalizuj često korišćene podatke (npr. `runway_length` u `airports`).

## 6 Prerada upita

Upiti su prilagođeni za novu šemu. Primer: Umesto `$lookup` za `weather`, koristi se `departure.weather.prcp`.

```
1 db.flights.aggregate([
2     {"$match": {
3         "departure.delay": {"$gt": 15},
4         "departure.weather.prcp": {"$gt": 10},
5         "departure.airport_type": "large_airport",
6         "departure.runway_length": {"$gt": 10000}
7     }},
8     {"$group": {
9         "_id": "$departure.airport",
10        "avg_delay": {"$avg": "$departure.delay"},
11        "count": {"$sum": 1}
12    }},
13     {"$sort": {"avg_delay": -1}},
14     {"$limit": 5}
15 ])
```

**Preporuke:** Proveri kompatibilnost upita sa novom šemom u Compass-u.

## 7 Uporedna analiza performansi

Performanse su upoređene pre i posle optimizacije:

- **Pre:** COLLSCAN, sporo (npr. 2.5s za upit 1).
- **Posle:** IXSCAN, brzo (npr. 0.5s za upit 1).

Primer za dijagram:

```
1 import matplotlib.pyplot as plt
2 queries = ["Upit 1", "Upit 2", "Upit 3"]
3 before = [2.5, 3.0, 2.8]
4 after = [0.5, 0.7, 0.6]
5 plt.bar(queries, before, width=0.4, label="Pre optimizacije")
```

```

6 plt.bar([q + 0.4 for q in range(len(queries))], after, width=0.4, label="
  Posle optimizacije")
7 plt.xlabel("Upiti")
8 plt.ylabel("Vreme (s)")
9 plt.legend()
10 plt.savefig("performance_comparison.png")

```

**Preporuke:** Uporedi `executionTimeMillis` iz `explain()`. Vizualizuj rezultate koristeći `matplotlib`.

## 8 Vizualizacija rezultata

Rezultati upita su vizualizovani u Metabase-u:

- **Podešavanje:** Poveži Metabase sa MongoDB (`mongodb://localhost:27017/flightsdb`). **Dashboard** *Barchart – ovizakanjenja, piechart – ovizaavio – kompanije*.  
 Primer: Bar chart prosečnih kašnjenja po aerodromima za upit 1.  
**Preporuke:** Eksportuj grafike iz Metabase-a kao PNG/JPG za prezentaciju.

## 9 Zaključak

Ovaj izveštaj pruža osnovu za realizaciju projekta, uključujući učitavanje podataka, kompleksne upite, optimizaciju, i vizualizaciju. Sledeći koraci:

- Testiraj učitavanje podataka i upite.
- Primeni indekse i denormalizaciju.
- Vizualizuj rezultate u Metabase-u.