פרויקט- אנליזה של ביג דאטה

https://www.kaggle.com/datasets/dilwong/flightprices/data :מאגר מידע

מאגר המידע מכיל 27 עמודות, כ82 מיליון שורות בנפח כולל של כ6B31. כל שורה במאגר מתאר חיפוש טיסה באתר Expedia, עבור כל טיסה מתועד ID יחודי לטיסה, תאריך חיפוש, תאריך הטיסה, מוצא, יעד, זמן טיסה, מרחק טיסה, ימי טיסה, מידע אודות תעריף ותוספות לכרטיס, האם הטיסה ישירה, תעריף בסיס, מחיר מלא, מס׳ המושבים שנותרו בטיסה זו, חברת תעופה, מחלקה, סוג מטוס ועוד.

בנוסף, לצורך מימוש אחת השאילתות הוספנו קובץ CSV שיצרנו, הקובץ מכיל 3 עמודות (CSV בנוסף, לצורך מימוש אחת השאילתות הוספנו קובץ

פירוט עמודות חשובות:

- .1 legId מזהה ייחודי לכל טיסה.משמש לזיהוי ולמעקב אחרי טיסות מסוימות.
- . Expedia תאריך החיפוש searchDate .2 מאפשר ניתוח קשר בין זמן החיפוש למחיר הכרטיס.
- .3 flightDate תאריך הטיסה.מאפשר לנתח כיצד מחירי הכרטיסים משתנים ביחס לתאריך הטיסה ותאריך ההזמנה.
- המוצא והיעד startingAirport / destinationAirport .4 (IATA).
 - מסייע בניתוח מסלולים בין שדות תעופה.
 - .5 totalTravelDistance המרחק הכולל של הטיסה במיילים. עשוי להיות קשור למחיר הטיסה ולמשך הטיסה.

מחירים ומחלקות:

- 6. totalFare המחיר הכולל של הכרטיס כולל מיסים ועמלות (בדולרים).
- (Boolean) האם מדובר בטיסה ישירה ללא עצירות ביניים. (isNonStop .7 עשוי להשפיע על נוחות הנוסע ועל המחיר.

מני טיסות:

- ימני המראה segmentsDepartureTimeRaw / segmentsArrivalTimeRaw .11 ונחיתה לכל קטע טיסה.
- מאפשר לנתח שעות פופולריות (ואת ההשפעה של השעה על המחיר) ולחשב זמני טיסות בפועל.
 - .12 travelDuration משך הטיסה בפורמט שעות ודקות.

מאפייני נוסעים ומטוסים:

.14 - seatsRemaining מספר המושבים שנותרו בטיסה. עשוי להשפיע על המחיר עקב ביקוש והיצע.

:DUCKD הכנסת הדאטה

על מנת להכניס את הדאטה לDUCKDB השתמשנו ביכולות של DUCKDB, מפאת גודל הדאטה, הקובץ נוצר רק אם הוא לא קיים, שכן מדובר בהמון זמן ריצה.

:הקוד להלן

```
if not file_exists("database.duckdb"):
    duckdb_conn = duckdb.connect("database.duckdb")
    duckdb_conn.execute("CREATE TABLE main AS SELECT * FROM
read_csv_auto('itineraries.csv')")
else:
    duckdb_conn = duckdb.connect("database.duckdb")
```

ביצענו: DUCKDB) ביצענו שדות התעופה) של מיקובץ השני (של מיקובץ השני (של מיקובץ השני (של מיקובץ השני

```
duckdb_conn.execute("CREATE TABLE airportsLocation AS SELECT * FROM
read csv auto('airportsLocation.csv')")
```

השאלות שנשאל על הדאטה והשאילתות SQL שנבצע עבורן:

1. כיצד משתנים מחירי הטיסות ביחס לזמן החיפוש מול מועד הטיסה?

מהטבלה שקיבלנו, ניתן לראות כי יש באופן כללי יש ירידת מחירים קלה במחיר הממוצע כאשר מזמינים יותר מ45 יום לפני הטיסה, לאחר מכן המחיר עולה באופן רציף עד למועד הטיסה.

2. האם משתלם יותר לטוס בטיסה ישירה או בטיסה עם עצירות (קונקשן) בהתחשב במרחקי טיסה שונים?

```
duckdb_conn.execute("""
CREATE TABLE query2 AS
SELECT
    CASE
        WHEN totalTravelDistance < 500 THEN 'Short (<500 miles)'
        WHEN totalTravelDistance < 1000 THEN 'Medium (500-1000 miles)'
        ELSE 'Long (>1000 miles)'
        END as distance_category,
        ROUND (AVG (CASE WHEN isNonStop THEN totalFare END), 2) as
direct_avg_fare,
        ROUND (AVG (CASE WHEN NOT isNonStop THEN totalFare END), 2) as
connection_avg_fare,
        COUNT (CASE WHEN isNonStop THEN 1 END) as direct_flights_count,
        COUNT (CASE WHEN NOT isNonStop THEN 1 END) as
connection_flights_count
FROM main
WHERE totalTravelDistance IS NOT NULL
GROUP BY 1
ORDER BY distance_category;
""")
```

מהטבלה שקיבלנו, ניתן לראות באופן מפתיע שמחירן הממוצע של טיסות ישירות בכל הקטגוריות יהיה זול יותר מטיסות לא ישירות.

3. איך שעת ההמראה משפיעה על המחיר, הזמינות ואופי הטיסה (ישירה\לא ישירה)!

מהטבלה שקיבלנו, נראה כי אין דפוס אחיד לגמרי למחיר הממוצע לטיסה ביחס לשעה, אך ניתן לראות שבשעות 6-10 בבוקר וכן בשעות 22-00 בלילה מחירי הטיסות יקרים יותר. לאחר השעה 10 בבוקר המחירים יורדים (כמעט תמיד. במקומות בהן אין ירידה, העלייה היא לא משמעותית) באופן רציף עד לשעה 22 ומשם מתחילה עלייה משמעותית.

בנוסף נראה כי אחוז הטיסות ללא עצירה עולה באופן רציף (כמעט) עד לשעה 20 בערב שבה 50 אחוז מהטיסות הן ישירות. לאחר מכן אחוז הטיסות הישירות יורד משמעותית כך שבשעות הלילה יש הכי מעט טיסות ישירות.

אין הרבה פער במספר המושבים הזמינים לאורך שעות היום, מדובר בפער של פחות מכיסא בממוצע לאורך כל שעות היום. כמו כן אין חוקיות מסויימת.

4. מהם ימי השבוע היקרים והזולים ביותר לטיסה והאם יש הבדל בין טיסות ישירות לטיסות עם עצירות?

```
GROUP BY flight_day
ORDER BY CASE flight_day
WHEN 'Sunday' THEN 0
WHEN 'Monday' THEN 1
WHEN 'Tuesday' THEN 2
WHEN 'Wednesday' THEN 3
WHEN 'Thursday' THEN 4
WHEN 'Friday' THEN 5
WHEN 'Saturday' THEN 6
END
""")
```

מהטבלה שקיבלנו, ניתן לראות כי הימים שלישי ורביעי הם הזולים ביותר לטיסה, ולאחר מכן מחיר הטיסה הממוצע עולה לאורך סוף השבוע (עם ירידה קלה בשבת) עד ליום ראשון שהוא היום היקר ביותר לטיסה. לאחר יום ראשון המחירים יורדים.

קיימת קורולציה כמעט מלאה בין תנודות המחיר הממוצע של טיסות לא ישירות לטיסות ישירות ביחס ליום בשבוע אך טיסה לא ישירה תמיד תהיה יקרה יותר.

5. כמה טיסות מתקיימות לכל קו טיסה (ישירות ולא ישירות) והאם יש קשר בין מספר הטיסות למחיר הטיסה?

מטבלה זו נוכל לראות את מספר הטיסות לכל קו (מוצא ויעד), בהתייחס לטיסות ישירות ולא ישירות, בנוסף יוצגו בה הקאורדינטות של מיקום כל שדה תעופה וכן המחיר הממוצע לכל קו.

מורכב מאוד להסיק תובנות על סמך ההתבוננות בטבלה, בויזואליזציה שתוצג בהמשך יהיה קל יותר לראות את ההשפעות של מספר הטיסות על המחיר. באופן כללי ניתן לראות כי בד"כ, ככל שיש פחות היצע – המחיר עולה.

הסיפור שנרצה לספר:

בעולם התיירות והטיסות, קיימים גורמים שונים המשפיעים על מחירי הטיסה. נרצה לנתח אותם ולהבין מהם הגורמים לכך ולזהות תנודות במחירי הטיסות.

הניתוח מתמקד בדפוסי המחירים של כרטיסי טיסה, תוך בחינה של קשרים בין גורמים שונים כמו הזמן בין ההזמנה עד הטיסה, מרחק הטיסה, זמני יציאה, היצע וביקוש, וסוגי טיסות (ישירה או עם עצירות).

השאילתות מספקות תובנות ברורות שמראות איך ומתי המחירים משתנים ומהם הגורמים שמניעים את השינויים הללו.

הכנסת הדאטה לSQLITE:

כדי ליצור טבלה לדוגמא בת 500 שורות, נשתמש בשאילתה הבאה שדוגמת את הטבלה הגדולה:

```
duckdb_conn.execute("CREATE TABLE sample AS SELECT * FROM main USING
SAMPLE 500;")
```

כדי להעלות את הדאטה הקטן לתוך SQLITE, נשתמש בפקודות הבאות:

נתקין את התוסף של SQLite, לאחר מכן נתחבר לדאטה בייס ונקרא לו sqliteDB. נבצע העתקה של הטבלאות בעזרת שאילתה שמתחילה ב- CREATE TABLE ולבסוף נקבל את הדאטה הקטן שלנו בSOLITE.

```
duckdb_conn.execute("INSTALL sqlite;")
duckdb_conn.execute("LOAD sqlite;")
duckdb_conn.execute("ATTACH 'database.sqlite' AS sqliteDB (TYPE
SQLITE);")

duckdb_conn.execute("CREATE TABLE sqliteDB.sample AS SELECT * FROM
sample")
duckdb_conn.execute("CREATE TABLE sqliteDB.airportsLocation AS SELECT
* FROM airportsLocation")
for i in range(1, 6):
    duckdb_conn.execute(f"CREATE TABLE sqliteDB.query{i} AS SELECT *
FROM query{i}")
```

שיטות הקטנת הנתונים

במהלך כל אחת מהשאילתות נעשה שימוש בטכניקות ספציפיות להקטנת הדאטה, תוך שמירה על מידע חיוני לתובנות וניתוחים. הנה פירוט רחב יותר:

- 1. קיבוץ:(GROUP BY) קיבוץ הנתונים לפי קטגוריות נבחרות, כמו:
- מספר ימים עד הטיסה (Query 1): מקבץ את כל הנתונים לפי ערכי days_until_flight, שומרים סיכומים עבור כל מספר ימים.
- **קטגוריות מרחק (Query 2):** שימוש בקטגוריות "Short/Medium/Long" כדי להפחית את המידע על טיסות ספציפיות ולהתמקד בתבניות כלליות.
- **שעות יציאה (Query 3):** מעבר משימוש בזמנים מדויקים (שעות ודקות) לשעות כלליות, כדי לצמצם את מספר הקטגוריות ולזהות מגמות רחבות יותר.

• ימים בשבוע (Query 4): קיבוץ לפי יום בשבוע מאפשר לזהות השפעות לפי ימי חול או סוף שבוע, במקום להתמקד בתאריכים ספציפיים.

רציונל :הקיבוץ מפחית את מספר השורות ומחליף את המידע הגולמי בערכים סטטיסטיים מרכזיים, כמו ממוצעים, מינימום ומקסימום, שחשובים לקבלת תובנות.

2. פונקציות אגרגציה (Aggregations):

בכל שאילתה מחושבים ערכים כמו:

- ממוצעים (AVG): למשל, מחיר ממוצע של כרטיסי טיסה כדי להבין את עלות הטיסות בכל קטגוריה.
 - **ספירות (COUNT):** למשל, מספר החיפושים או הטיסות בכל קטגוריה כדי למדוד פעילות.
- **ערכים קיצוניים (MIN/MAX):** כמו מחיר מינימלי ומקסימלי של טיסות, כדי לזהות שונות במחירים.
 - עמשל, אחוז הטיסות הישירות ב.(PERCENTAGE): למשל, אחוז הטיסות הישירות ב

רציונל: פונקציות האגרגריה מאפשרות ריכוז מידע משמעותי במקום שמירה על כל הנתונים הגולמיים, שמכבידים על תהליך הניתוח.

:(Filtering) סינון.3

- מחיקת נתונים לא רלוונטיים: לדוגמה, ב Query 5 מסננים רק מקרים שבהם השינוי היומי במחיר היה מעל 20% או מתחת ל--20%.
- **התמקדות בערכים לא ריקים:** ב 2 Query מתמקדים בטיסות עם מרחקים מוגדרים totalTravelDistance ומסננים שורות שבהן

רציונל :סינון מקטין את הדאטה בכך שהוא מוודא שהניתוח יתבצע רק על מידע משמעותי ורלוונטי, תוך התעלמות מפרטים שלא תורמים להבנה.

4. חלוקה לקטגוריות (Bucketing):

- ב-2 Query המרחקים מחולקים לטווחים קבועים מראש, כמו "Short/Medium/Long".
 - ב Query 4 ימים בשבוע מחולקים לסדר מוגדר כדי לזהות הבדלים בין ימי חול לסופי שבוע.

רציונל :חלוקה לקטגוריות מפשטת את ניתוח הנתונים ומאפשרת זיהוי מגמות ברורות בתוך קבוצות מוגדרות היטב.

5. פונקציות חלון:(Window Functions)

ב - Query 5 נעשה שימוש ב - LAG כדי לחשב את המחיר הממוצע של יום קודם ולהשוות אותו למחיר הממוצע הנוכחי. למחיר הממוצע הנוכחי.

רציונל :פונקציות חלון שומרות על הקשר בין שורות ומאפשרות ניתוח השינויים לאורך זמן, תוך שמירה על מבנה נתונים מסודר לצמצום מידע לא רלוונטי.

מה קיבלנו מכל זה?

- 1. צמצום נפח המידע: במקום לשמור מיליוני שורות של נתונים גולמיים, התוצאה היא טבלה קטנה יותר וממוקדת יותר.
- 2. **פשטות בניתוח** :במקום לעבד נתונים מורכבים, יש נתונים מסוכמים שקל יותר לעבוד איתם.
 - 3. זיהוי מגמות ודפוסים :המידע המעובד מאפשר לזהות מגמות ברורות ולבצע אופטימיזציות בקלות.
- 4. **שיפור ביצועים** :פחות נתונים לעיבוד מובילים לזמן תגובה קצר יותר וביצועים טובים יותר עבור שאילתות עתידיות.

הטבלאות שנמצאות SQLITE:

- 1. query1 הטבלה מכילה 5 עמודות ו60 שורות. מטרתה להציג את שינוי מחירי הטיסות ביחס לזמן החיפוש אל מול מועד הטיסה. כלומר, נרצה לראות אם יש קשר בין מועד ההזמנה למחירי הטיסה.
 - 2. query הטבלה מכילה 5 עמודות ו3 שורות, מטרתה לנתח כדאיות של טיסה ישירה auery 2 מול טיסה לא ישירה תוך התייחסות לטווח הטיסה.
- 3. query3 הטבלה מכילה 7 עמודות ו21 שורות, שורה עבור כל שעה עגולה שתועדה בה טיסה. מטרתה לנתח שעות שיותר יקר או יותר זול לטוס בהן, שעות בהן יש זמינות מושבים גבוהה יותר וכן את השעות בהן יש יותר טיסות ישירות.
- 4. **query4** הטבלה מכילה 5 עמודות ו7 שורות, שורה עבור כל יום בשבוע. מטרתה לנתח את הימים הזולים והיקרים ביותר באותו שבוע באופן ממוצע וכן בהתייחסות לטיסות ישירות או לא ישירות.
- 5. query5 הטבלה מכילה 11 עמודות 450 שורות. כל שורה תציג את מספר הטיסות לכל קובטיסה ישירה ולא ישירה וכן את המחיר הממוצע, המקסימלי והמינימלי לאותו קו. מטרתה לבדוק את הקשר בין כמות הטיסות למחיר הטיסה (היצע וביקוש).
 - של מטרתה להוות מעין יימדגםיי של Sample .6. הטבלה מכילה 27 עמודות 500 שורות. מטרתה להוות מעין יימדגםיי של הטבלה הגדולה.
- 2. airportsLocation הטבלה מכילה 3 עמודות 161 שורות. מכיל מידע גולמי על מיקום airportsLocation של כל שדה תעופה בדאטה.

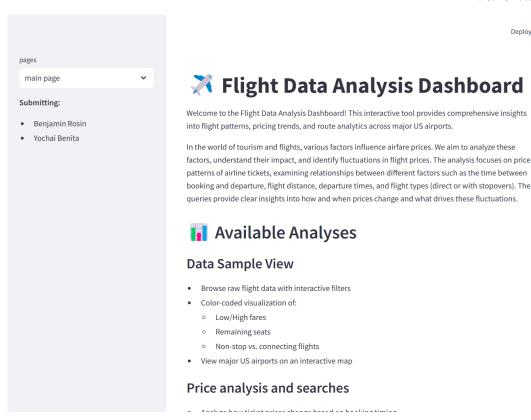
תיאור סכמתי של הדשבורד והסברים על הויזואליזציות:

הדשבורד מורכב מ-7 עמודים עם סרגל ניווט בצד:

Main Page .1

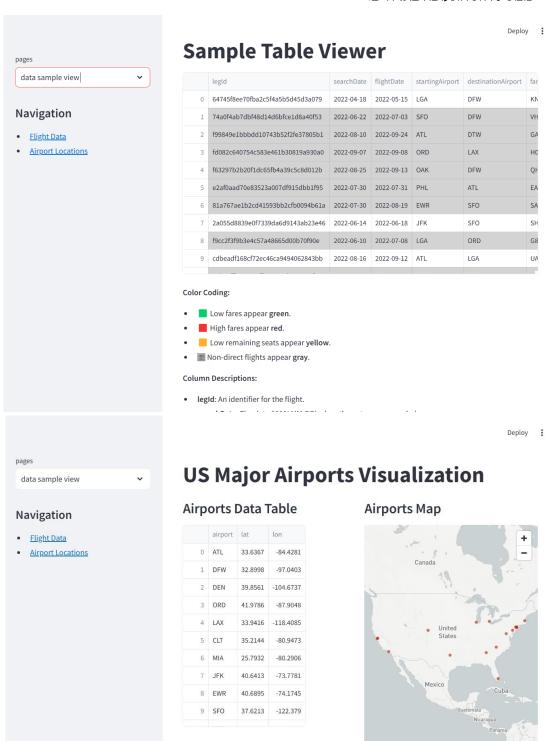
Deploy :

דף נחיתה עם תיאור הפרויקט והעבודה שביצענו, הסיפור שנרצה לספר והסבר על הדאשבורד.



Data Sample .2

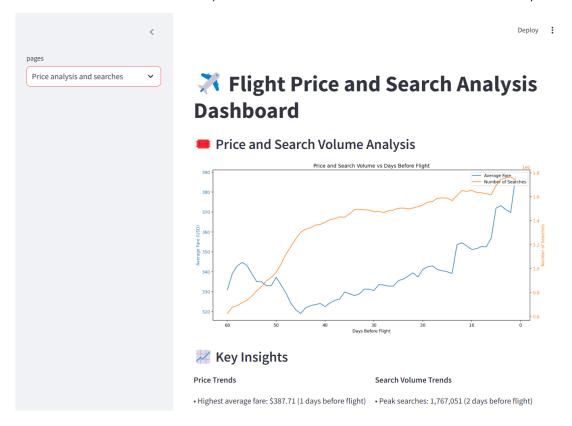
- טבלת נתוני הטיסות, בת 500 שורות כאשר היא מייצגת מדגם של הדאטה סט.
 - מפת שדות התעופה בארהייב.



(a) mapbox

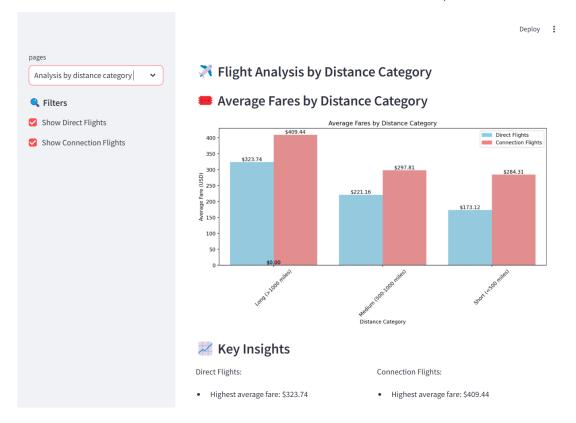
Price Analysis and Searches .3

• גרף המציג את מחיר הכרטיס ומספר החיפושים ביחס לזמן שלפני הטיסה.



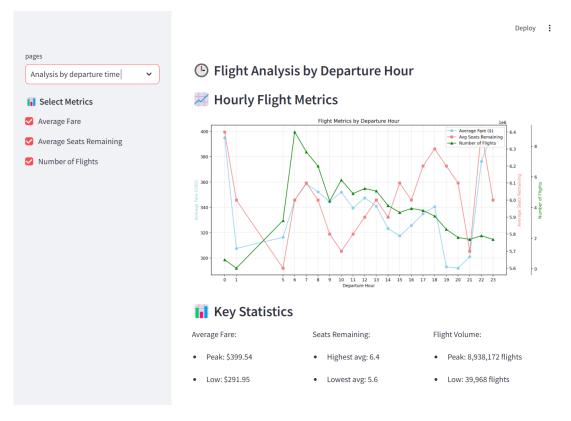
Analysis by Distance Category .4

- . גרף אינטראקטיבי שמציג השוואת מחירים ממוצעים לפי קטגוריות מרחק.
 - קיימת אפשרות לסנן בין טיסות ישירות וטיסות קונקשיין.
 - מוצגת סטטיסטיקה לכל סוג טיסה.



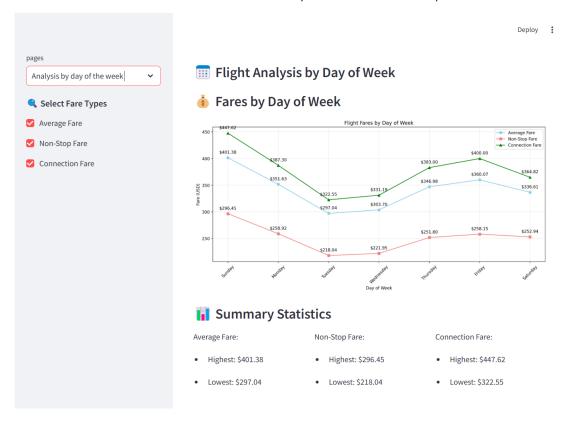
Analysis by Departure Time .5

- גרף אינטראקטיבי המציג נתונים שונים ביחס לשעת ההמראה, כגון מחיר ממוצע, מושבים פנויים בממוצע ומספר טיסות.
 - מוצגות סטטיסטיקות לכל מדד המוצג בגרף.



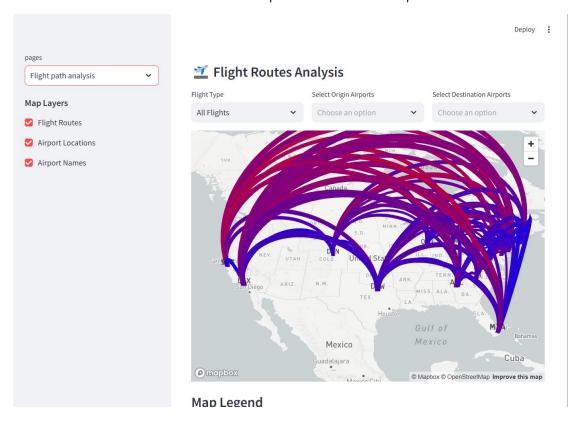
Analysis by Day of the Week .6

- גרף אינטראקטיבי המציג השוואת מחירים לפי ימי השבוע.
- מאפשר לסנן ולהציג מחיר ממוצע כללי, מחירים עבור טיסות ישירות וכן מחירים עבור טיסות לא ישירות.
 - מוצגות סטטיסטיקות לכל מדד המוצג בגרף.



Flight Path Analysis .7

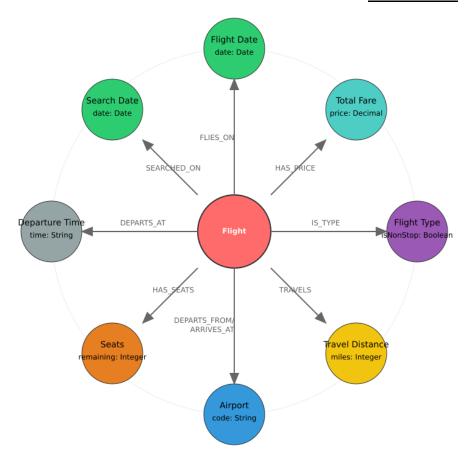
- מפה אינטראקטיבית של נתיבי טיסה, כאשר צבעי הקווים מייצגים מחירים (כחול זול, אדום יקר) ככל שהצבע יותר קרוב לאדום המחיר יותר יקר.
 - . ניתן לסנן את התוכן המוצג לפי סוג הטיסה, מוצא ויעד.
 - ניתן לבחור מה יוצג במפה: מיקומי נמלי התעופה, שמות הנמלים, וכן נתיבי הטיסה.
 - מוצגות סטטיסטיקות לכל נתיב המוצג בגרף.



הוראות הרצה לדאשבורד:

- : יש להתקין את הספריות הנדרשות להרצת הפרויקט עיי הפקודה הבאה .1 pip install -r requirements.txt
 - 2. יש להריץ בטרמינל: streamlit run ./dashboard.py

:GRAPH DB



סכמת הגרף מציגה את מבנה הדאטה בייס שעליו אנו עובדים.

במרכז הסכמה נמצא הצומת המרכזי Flight (טיסה) שמקושר ל-8 צמתים המקיפים אותו, כל אחד מייצג מאפיין שונה של הטיסה.

הצמתים המקושרים לטיסה הם:

- .FLIES ON מתי הטיסה אמורה להתקיים, מקושר דרך (Flight Date) מתי הטיסה. 1.
 - 2. מחיר כולל (Total Fare) עלות הטיסה, מקושר דרך HAS_PRICE.
- .IS TYPE האם זו טיסה ישירה או עם עצירות, מקושר דרך (Flight Type) . סוג טיסה
- .4 מרחק טיסה (Travel Distance) כמה מיילים הטיסה מכסה, מקושר דרך
 - 5. שדה תעופה (Airport) קוד שדה התעופה, מקושר דרך

DEPARTS_FROM/ARRIVES_AT כי כל טיסה יוצאת משדה תעופה אחד ונוחתת בשדה DEPARTS_FROM/ARRIVES_AT תעופה אחר.

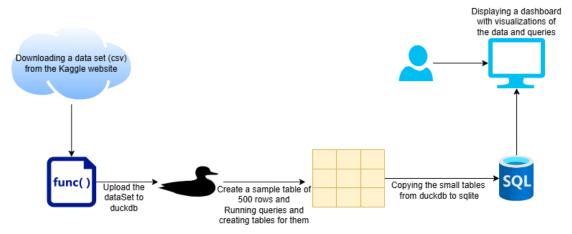
- 6. מושבים (Seats) כמה מושבים פנויים בטיסה, מקושר דרך Keats).
- .DEPARTS_AT מתי הטיסה ממריאה, מקושר דרך (Departure Time) מתי הטיסה
 - 8. תאריך חיפוש של הטיסה, מקושר דרך (Search Date) מתי בוצע החיפוש של הטיסה. SEARCHED ON

כל הקשרים בין הצומת המרכזי לשאר הצמתים (nodes) מסומנים בחיצים המראים את כיוון כל הקשרים בין הצומת המרכזי לשאר הצמתים (FLIES_ON, HAS_PRICE וכו').

הסכמה הזו מאפשרת לעקוב אחר מחירי טיסות, לנתח זמינות מושבים, להשוות בין סוגי טיסות שונים ולזהות מגמות במחירים לאורך זמן.

מבנה של המערכת:

את קבצי הCSV עם הדאטה נעלה בעזרת יכולות של הDUCKBD לקובץ בפורמט הנייל, על גבי הטבלאות הנייל נריץ שאילתות שיוצרות מסי טבלאות קטנות, חלק מהטבלאות מהוות יטעימהי מהדאטה הגדול וחלקן מהוות תשובה לשאילתה ספציפית. לאחר מכן, בעזרת כלים של BQLITE ניצור קובץ SQLITE ונעתיק אליו את כל הטבלאות הקטנות. בכך סיימנו את השלב הראשון. בשלב שני, בהרצת הקוד, על בסיס המידע שנמצא בקובץ הSQLITE נציג ויזואליזציות מרשימות שמנגישות את המסקנות מהמידע לצרכן.



העבודה על הפרויקט הנ״ל התחלקה באופן שווה בין שנינו לאורך הזמן, כל האלמנטים בפרויקט נכתבו בעבודה משותפת.