# Visualizing Airfare Trends ✈️

Ishan Patel, Sristi Karamchandani, Juntae Park, Royce Arockiasamy, Ved Sanjanwala, Luc Gau

## Motivation:

Airfare pricing has traditionally been driven by short-term factors like booking timelines and seat demand. However, these dynamic pricing models often fail to account for broader, long-term influences such as economic conditions and global disruptions. As a result, price trends appear unpredictable and reactive, limiting the ability of both consumers and industry stakeholders to make informed decisions.
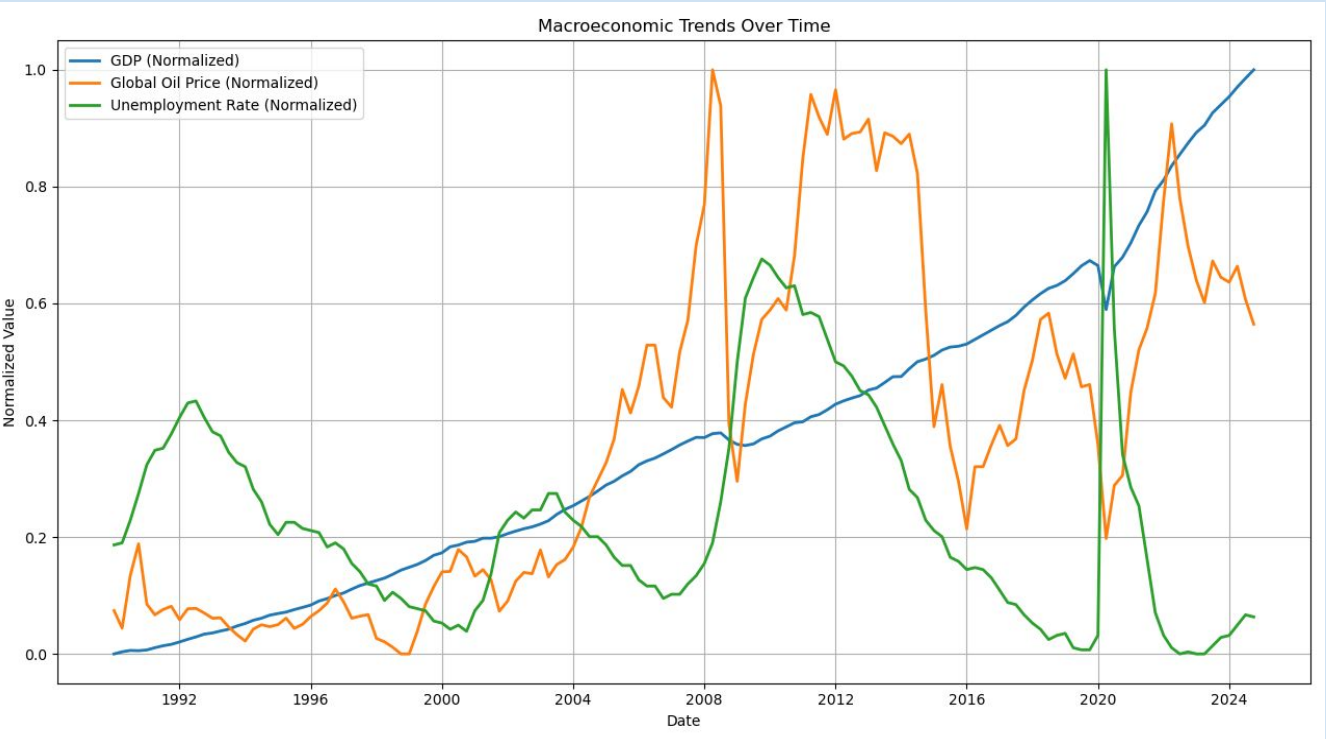
The airline industry is a critical driver of global connectivity and economic development. Understanding the forces behind airfare fluctuations is essential for travelers planning ahead, airlines optimizing revenue strategies, and policymakers ensuring market transparency. By identifying and visualizing how external factors influence airfare trends over time, our project offers a foundation for smarter forecasting and better decision-making for airfare pricing



## Data:


Flight Data


Marco Data


Macroeconomic Trends Over Time

**U.S. Airline Fare Data (Kaggle):**
- Years Covered: 1993–2024
- Size: 63 MB, 2 million records
- Frequency: Quarterly
- Key Features:
  - Origin/destination cities and airport codes
  - Distance, number of passengers
  - Average fare and carrier-specific fare info
- Purpose: Core dataset used to model fare trends across domestic routes

**U.S. GDP Data (FRED):**
- Years Covered: 1947–2024
- Frequency: Quarterly
- Units: Billions of U.S. Dollars (Seasonally Adjusted Annual Rate)
- Purpose: Captures macroeconomic growth, used to understand long-term pricing trends

**Brent Crude Oil Prices (FRED):**
- Years Covered: 1990–2025
- Frequency: Quarterly
- Units: U.S. Dollars per Barrel
- Purpose: Included as a key cost driver of airline operations; captures fuel-related fare fluctuations
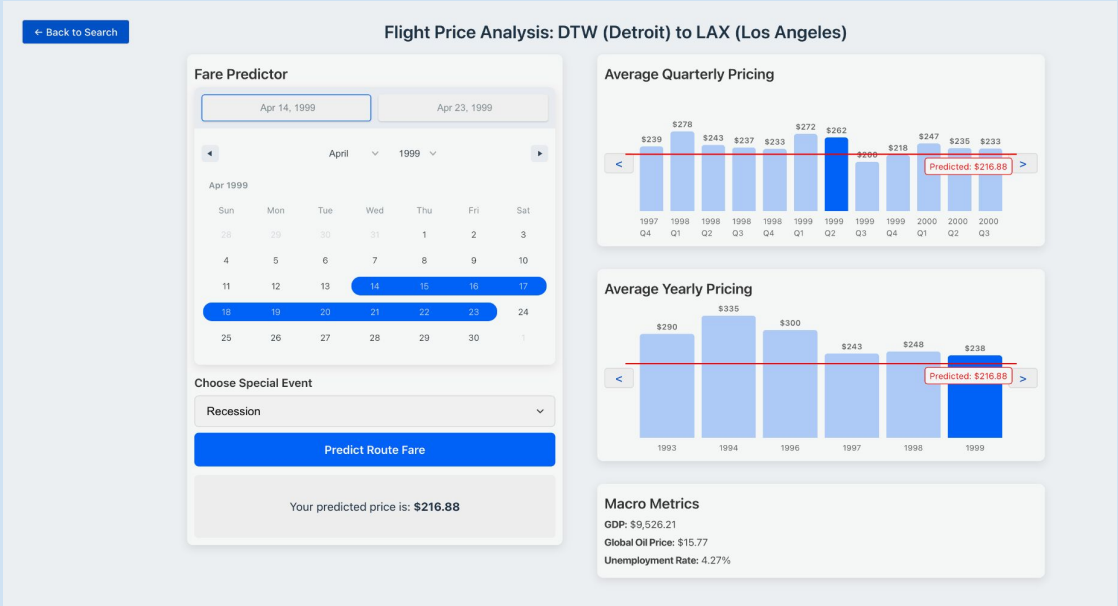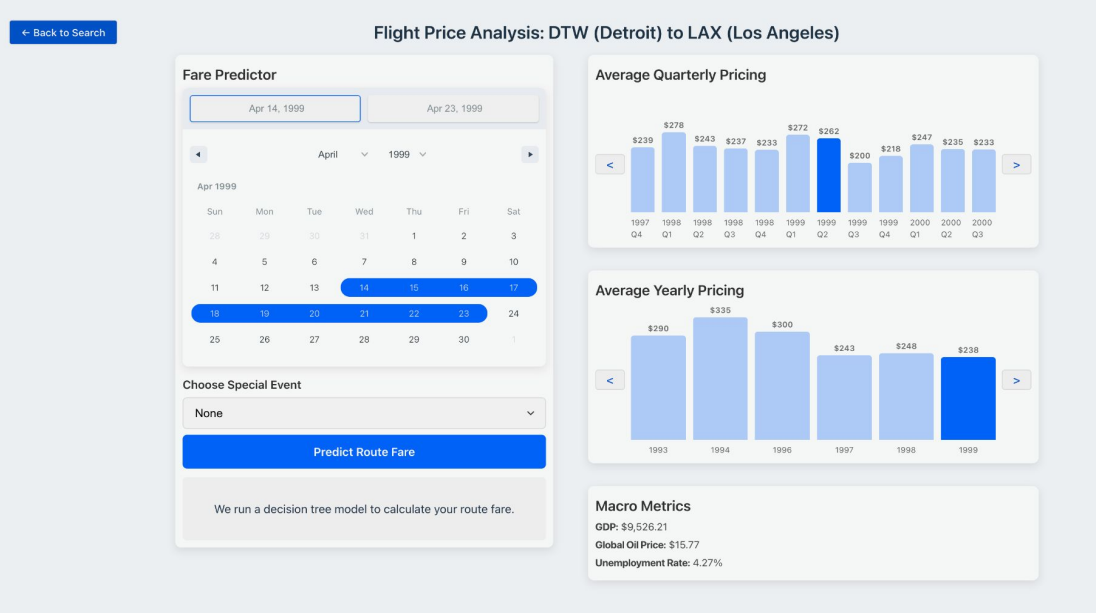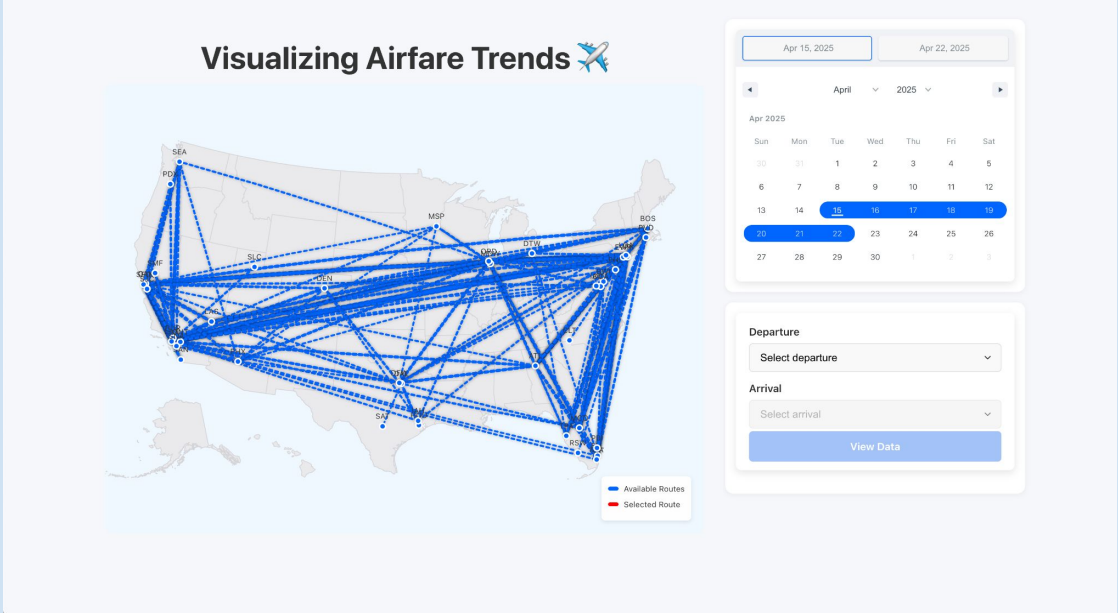
**U.S. Unemployment Rate (FRED):**
- Years Covered: 1948–2025
- Frequency: Monthly
- Units: Percentage (Seasonally Adjusted)
- Purpose: Indicator of consumer demand and travel behavior; aggregated to quarterly for model input

**Data Collection and Cleaning:**
- Downloaded 20+ years of U.S. domestic airfare data from Kaggle
- Pulled macroeconomic indicators (GDP, oil prices, unemployment) from FRED
- Added event-based markers (e.g., 9/11, oil prices, 2008 recession) for contextual modeling
- Parsed GDP and oil prices directly from web tables (quarterly)
- Scraped unemployment data (monthly) and aggregated to quarterly for consistency
- Removed rows with NaNs, duplicates, and outliers (e.g., zero fares, extreme values)
- Merged macroeconomic indicators into airfare data by quarter and year
- Ensured alignment between route-level fare data and national-level economic trends

## Algorithm and Interactive Visualization:





### 🛠 Model Setup

We use a **Random Forest Regressor** to predict airfare trends. This model handles non-linearity and categorical variables well.



### 🧠 Fare Prediction

We predict airfare using the selected route and travel dates.



### 📊 Data & Modeling Libraries

We import necessary libraries for data manipulation, model training, and evaluation.



### 🔍 Small Subset Sampling

To speed up model training and testing, we use a small random sample of the training data.



### Our Approach:

➤ A hybrid machine learning model paired with an interactive web-based visualization tool that predicts and explains airline fare trends across U.S. routes

➤ Model: Learns from 20+ years of fare data, economic indicators (GDP, oil, unemployment), and special events to predict future prices

➤ Visualization: Lets users explore historical pricing, select routes/dates, and simulate the effects of disruptions like recessions in real time

➤ Traditional fare models react to short-term demand.

➤ Airfare pricing is driven by more than demand, it is influenced by seasonality, economic trends, and global events

➤ Ours model is proactive, capturing both linear patterns (e.g., inflation) and nonlinear shocks (e.g., pandemics)

➤ Our visualization makes these insights accessible and actionable

➤ Offers a user-facing platform that ties pricing, demand, and economic signals into one decision-support system

### Visualization Features:

➤ Route Map Visualization: Interactive U.S. map displays all available routes; users can explore pricing by selecting specific origin-destination pairs.

➤ Custom Date Selection: Allows users to view fare predictions for custom travel dates

➤ Event Simulation Tool: Users can simulate the impact of external events like recessions on fare predictions allowing for scenario-based analysis

➤ Quarterly and Yearly Pricing Trends: Bar charts show historical fare patterns over time, helping users detect seasonality, inflationary effects, and airline pricing cycles.

➤ Macro Metrics Dashboard: Displays GDP, oil prices, and unemployment rates relevant to the selected time period, helping connect fares to economic context

### ML Model:

➤ Random Forest: handles complex, nonlinear impacts, like oil prices and economic shocks

➤ Ridge Regression: captures smooth, long-term economic trends

## Experiments and Results:




Absolute Fare Prediction Error

| Statistic | Difference ($) |
| --- | --- |
| Mean | 8.46 |
| Standard Deviation | 7.92 |
| Minimum | 0 |
| 25th Percentile | 3.59 |
| 75th Percentile | 10.36 |
| Maximum | 77.92 |

**Model Performance Metrics**

| Metric | Value |
| --- | --- |
| Mean Absolute Error | 8.461986 |
| Root Mean Squared Error | 11.590751 |
| R² Score | 0.962821 |
| Mean Absolute Percentage Error | 4.402026 |

**Prediction vs Actual Fare Graphs:**

➤ Predictions closely track actual fare trends, reflecting real market behavior

➤ Seasonal spikes and long-term shifts are accurately captured

➤ Route-specific dynamics (e.g., busy vs. regional routes) are well preserved

**Absolute Fare Prediction Error:**

➤ 95.6% prediction accuracy

➤ Model fares align closely with real prices

➤ Most predictions are within a few dollars, making it practical for planning and pricing

➤ Performs well even in volatile periods, such as recessions and pandemics

➤ Tight error range shows consistent and trustworthy performance

➤ Model explains 96% of variation in airfare prices

➤ Benchmarked traditional models (e.g., linear regression) show lower performance of $R^2 \approx 0.61$

➤ Captures both macro and micro trends

➤ Provides a reliable, high-confidence tool for forecasting and analysis