

BETTERFLY

FINAL PROJECT – HIT
DOR SLAGTER



RESEARCH QUESTION

Is it possible to predict the price of a flight on a variety of sites like kayak using the characteristics of the flight and the seasons of the year?

RESEARCH QUESTION

Is it possible to predict the price of a flight on a variety of sites like kayak using the characteristics of the flight and the seasons of the year?





SCRAPING

SCRAPING TOOL

User Input

```
In [ ]: # get user input for routes
sources = []
destinations = []
print("Please enter -1 when done.")
print("-"*10)
while True:
    sources.append(input("From which city?\n"))
    if "-1" in sources:
        sources.pop(-1)
        break
    destinations.append(input("Where to?\n"))
    if "-1" in destinations:
        destinations.pop(-1)
        break
    print("-"*10)

print("\nRoutes:")
for i in range(len(sources)):
    print(f"{sources[i]} => {destinations[i]}")
```

```
In [ ]: # get user input for period (start and end date)
start_date = np.datetime64(input('Start Date, Please use YYYY-MM-DD format only '))
end_date = np.datetime64(input('End Date, Please use YYYY-MM-DD format only '))
days = end_date - start_date
num_days = days.item().days
```


SCRAPING

- Source
- Destination
- Date
- Price
- Duration
- Total stops
- Airline

The screenshot shows the Kayak flight search interface. At the top, the search parameters are: Round-trip, Tel Aviv (TLV) to Bangkok (BKK), Friday 14.07 to Friday 21.07, 1 adult, Economy class. The search results are sorted by 'Best' and show a price of ₪2,306 for a 1-stop flight. The flight details for the 'Best' option are: 22:00 - 15:25 (1 stop, 13h 25m) and 07:50 - 19:10 (1 stop, 15h 20m). The airlines listed are Turkish Airlines and Gulfair. The price is ₪2,306. The 'Cheapest' option is also shown with a price of ₪2,306 and a duration of 13h 47m. The 'Quickest' option is ₪4,676 with a duration of 12h 07m. The page also features a 'Track prices' toggle, 'Recommended filters', and a 'Fare Assistant' section. A Club Med advertisement is visible on the right side of the page.

Sort	Flight Details	Price
Cheapest	22:00 - 15:25 [†] (1 stop, 13h 47m)	₪2,306
Best	22:00 - 15:25 [†] (1 stop, 13h 25m) and 07:50 - 19:10 (1 stop, 15h 20m)	₪2,306
Quickest	22:00 - 15:25 [†] (1 stop, 12h 07m)	₪4,676

SCRAPED ROUTES



GRU



BKK



TLV



LUX



PAR



MEX



YYZ



OSL



FCO



AMS



NYC

ISSUES

Security Check:

Please confirm that you are a real KAYAK user.

- Closing and opening the driver every 10 days (pages).
- Alert the user to go solve the captcha and resume scraping



I'm not a robot



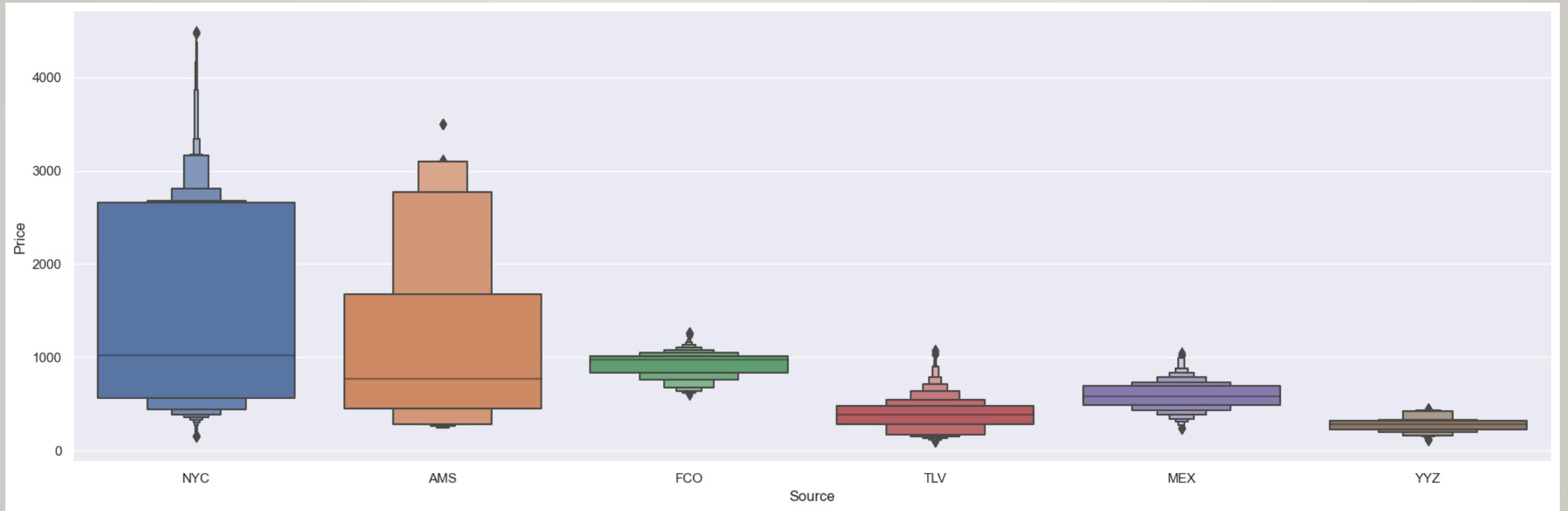
reCAPTCHA
Privacy - Terms

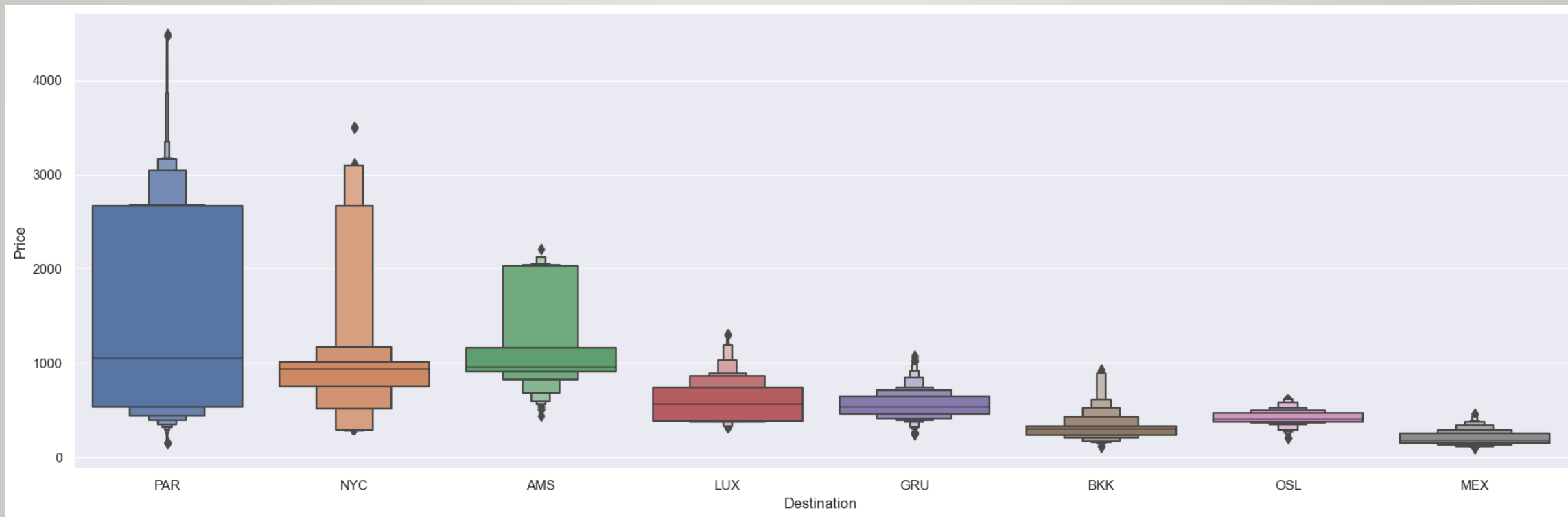
THE DATASET

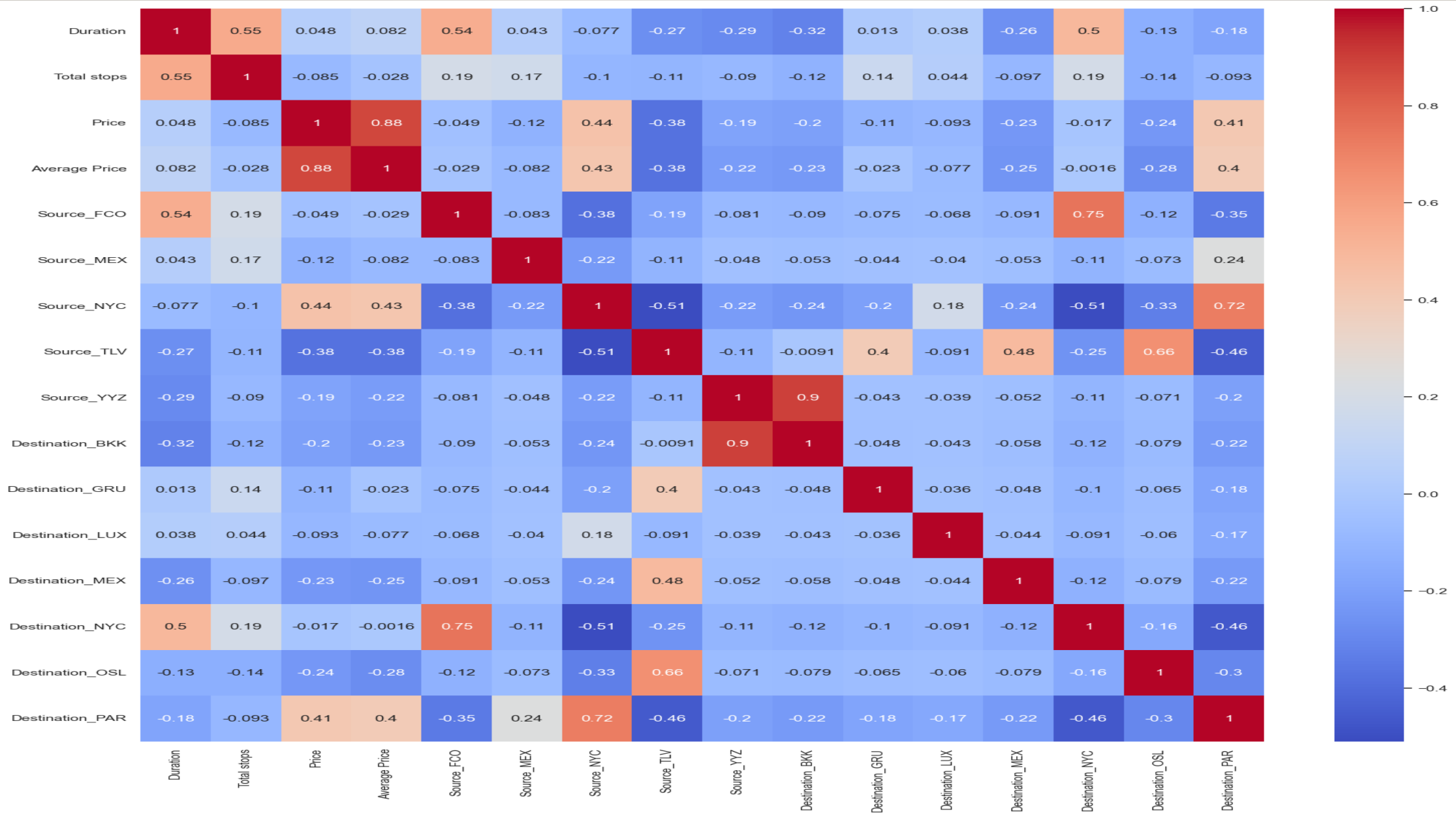
	Airline	Source	Destination	Duration	Total stops	Price	Date	Average Price
0	Aeroflot	AMS	NYC	615	nonstop	70.69	2023-03-01	67.915656
1	Aeroflot	AMS	NYC	615	nonstop	66.06	2023-03-01	67.915656
2	Aeroflot	AMS	NYC	615	nonstop	66.06	2023-03-01	67.915656
3	Aeroflot	AMS	NYC	615	nonstop	70.69	2023-03-02	67.915656
4	Aeroflot	AMS	NYC	615	nonstop	66.06	2023-03-02	67.915656
...
2174	Air France, airBaltic	YYZ	BKK	655	3 stops	295.73	2023-04-09	283.559226
2175	Air France, Pobeda	YYZ	BKK	400	1 stop	286.67	2023-04-03	273.818786
2176	Transavia France, Finnair	YYZ	BKK	695	2 stops	202.13	2023-04-05	215.548002
2177	Norwegian, LOT	YYZ	BKK	825	2 stops	192.53	2023-04-09	251.441247
2178	SWISS, Pobeda	YYZ	BKK	520	2 stops	285.87	2023-04-11	672.917667



EDA







RESULTS

RESULTS

model / metric	train(60%)	Val(20%)	MAE	MSE	RMSE
Linear regression	0.804	0.789	225.09	152995.68	391.15
Polynomial(5)	0.888	0.871	149.72	93366.04	305.56
Lasso	0.803	0.789	223.49	153097.06	391.28
Ridge	0.804	0.789	225.08	152995.61	391.15
Elastic Net	0.790	0.775	224.41	163034.63	403.78
Random Forest	0.965	0.945	61.72	40035.32	200.09

Actual VS Predicted

