

# Forecasting Travel Times for Uber Riders in Los Angeles

May 24<sup>th</sup>, 2023



# Agenda

**1: Problem  
Statement &  
Approach**

**2: Data  
Exploration**

**3: Modeling**

**4: Takeaways  
Recommendations, &  
Future Work**

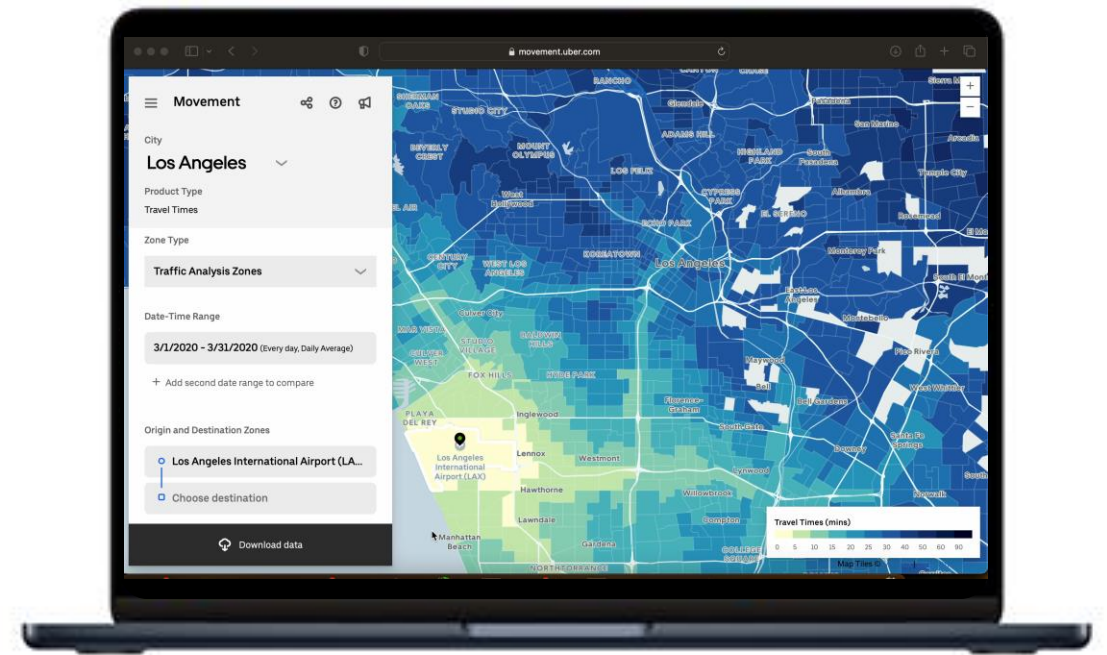


# Problem Statement

Uber Movement shares anonymized data aggregated from over ten billion trips to help urban planning around the world.

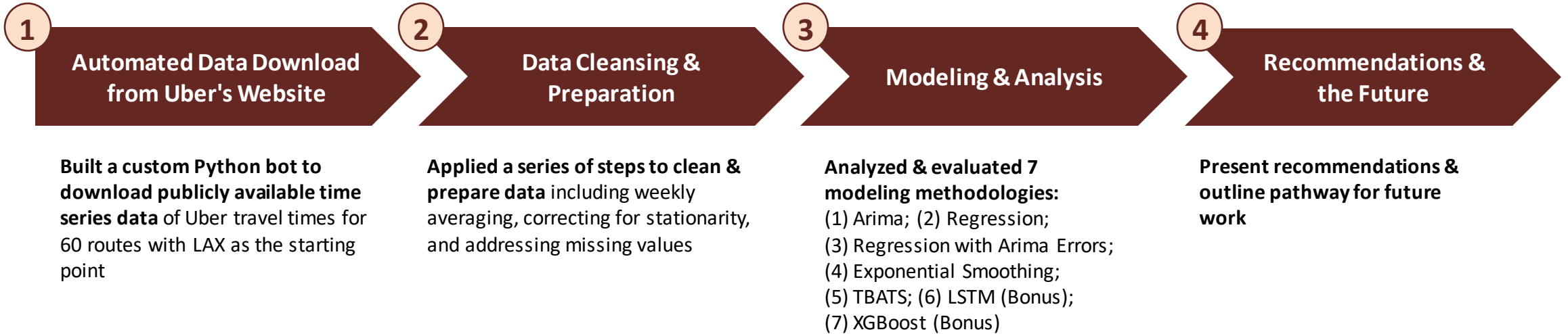
One of these data sets is based on travel times between two points.

**How can we build accurate and robust time series models to drive cost savings within Uber Operations while continuing to meet customer needs?**



# Our high-level analysis approach

We leveraged the power of our automation-fueled pipeline to address the problem statement



```
getUberMovementData2 - Jupyter
localhost:8888/notebooks/ChicagoSpring2023/TimeSeries%20MISCAN%201006/Project/getUberMovementData2.ipynb
File Edit View Insert Cell Kernel Widgets Help
Python 3 (Spyder)
getUberMovementData2 (unsaved changes)
In [18]: # Inputs
data_range = generate_date_ranges("2016-01-01", "2020-03-31")
dest_dct = pd.read_csv("LAX All Destinations.csv")
destinations = list(dest_dct['Destination Movement ID'].values)
origin = dest_dct['Origin Movement ID'].unique()[0]
destinations.remove(origin)
download_path = "C:/Users/DeepakVijani/Library/Mobile Documents/com~apple~CloudDocs/ChicagoSpring 2023/Time Series (MISCAN)
download_path = "C:/Users/DeepakVijani/Library/Application Support/Google/Chrome/Default"

In [53]: # Run Crawler
download_city = "Los Angeles", data_range = data_range,
origin_code = origin, destinations = destinations[1:],
download_path = download_path
```

```
df_list <- df_list_all
neighborhoods <- names(df_list)
models_df <- rbind()
for (route in neighborhoods){
  cat("*****\n")
  cat(paste0("Building models for LAX to ", route, "\n"))
  weekly_df <- get_weekly_avg(df_list[[route]])
  if(missing_brnach(weekly_df)){
    cat(paste0("Removing ", route, " from analysis, as too many missing values\n"))
    next
  }
  else {
    weekly_df <- clean_df(weekly_df)
    # plot_series(weekly_df)
    is_stationary <- check_stationarity(weekly_df)
    cat(paste0("Series is stationary:", is_stationary), "\n")
    temp <- cbind(Destination = route,
                  Stationary = is_stationary,
                  model(weekly_df, features = features, stationary = is_stationary))
    models_df <- rbind(models_df, temp)
  }
}
write.csv(models_df_2, "models_df_2.csv")
```



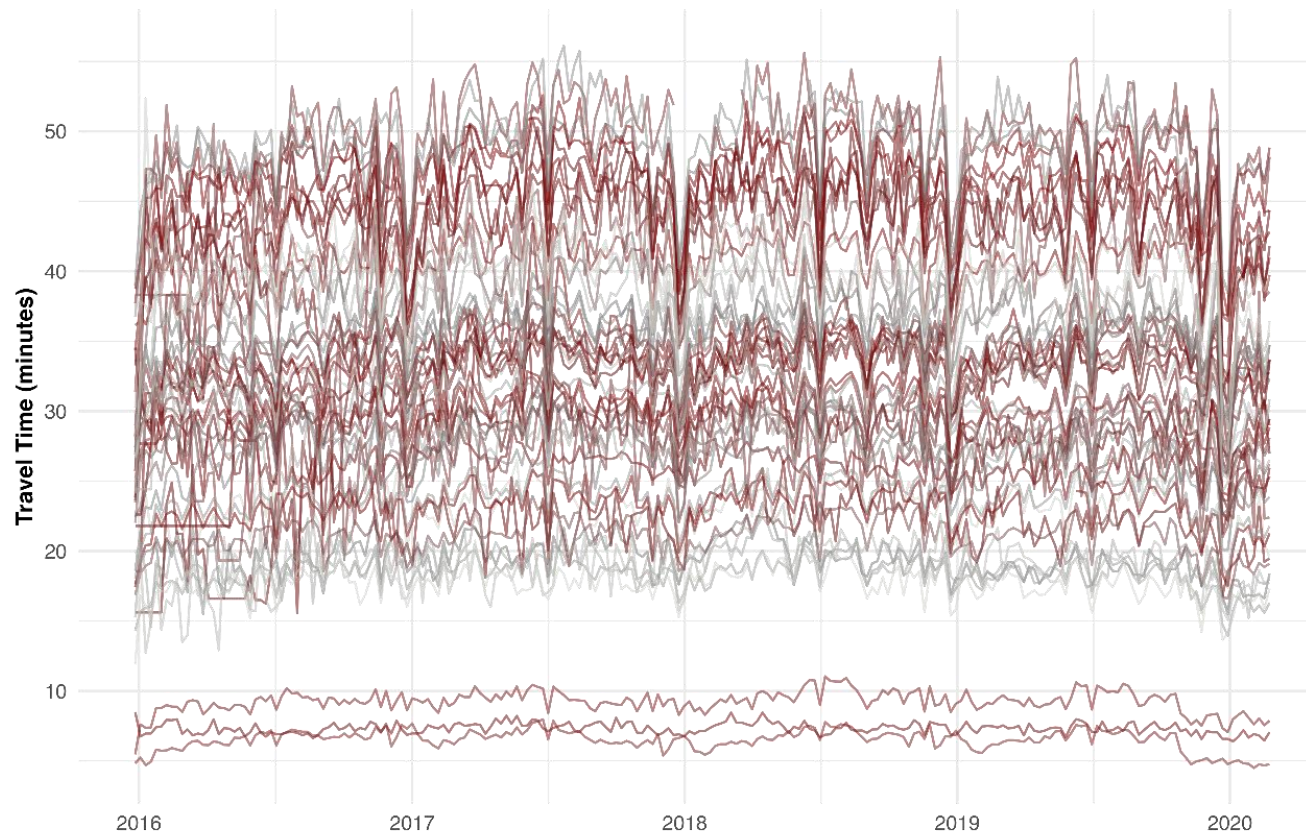
# Key features of our data

Data

Modeling

Conclusion

Average Weekly Travel Time between LAX and 50 different locations in LA (Jan 2016 - Feb 2020)



- **COVID Impact:** Data is available from January 2016 – March 2020, but we exclude March to avoid one-time COVID impact.
- **Missing Values:** 60 routes filtered down to 50 using a missing value threshold of 15%. Any remaining missing values were addressed using Last Observation Carried Forward imputation.
- **Seasonality & Trend:** Frequency of 52 weeks, with clear signs of annual seasonality. No clear trend present.
- **Stationarity:** All routes modeled were stationary using ADF test threshold of 0.05. Some routes required seasonal and/or non-seasonal differencing of order 1.
- **External Features:** Our analysis considered external variables including publicly available data for temperature in LA, relative humidity in LA, and a national holiday indicator.



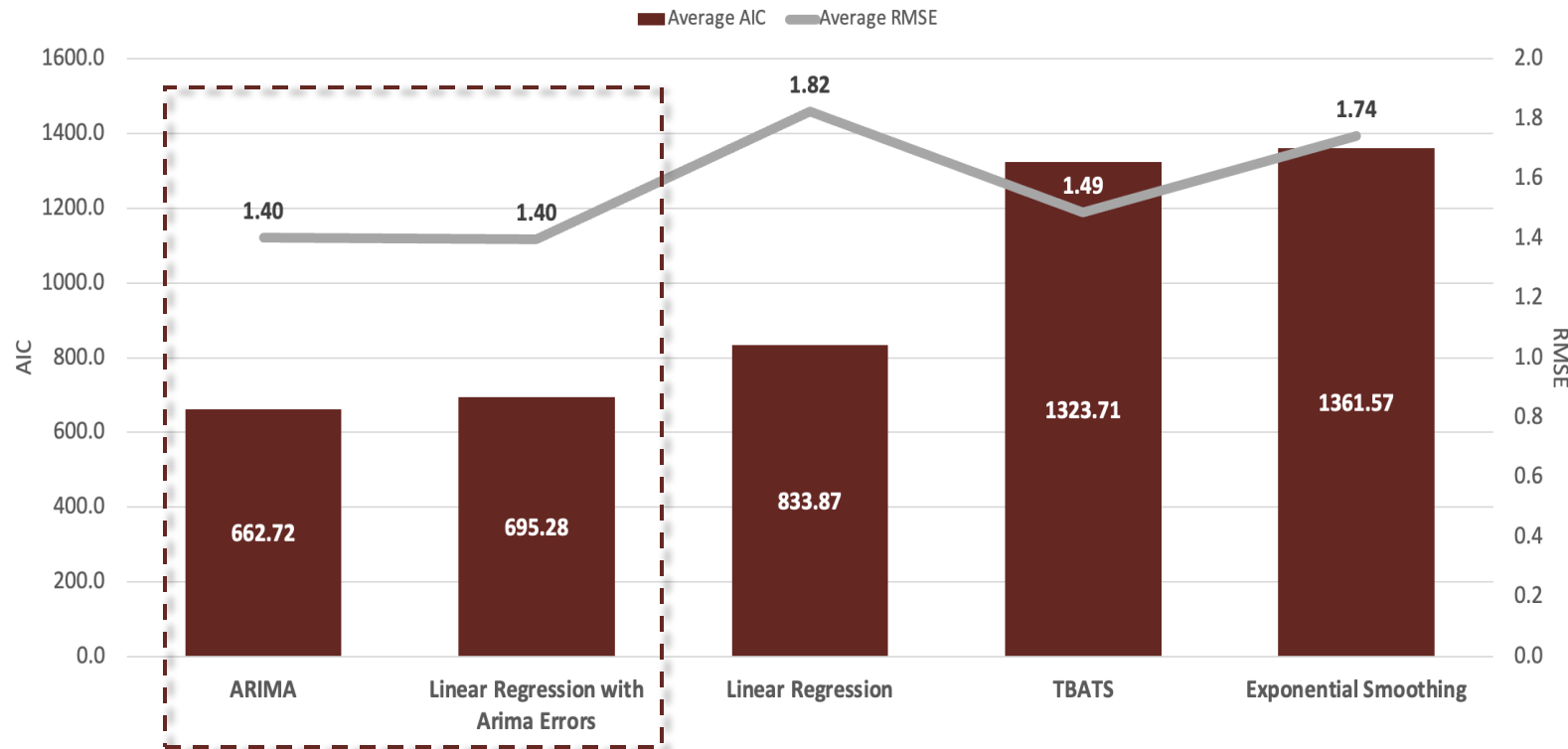
# Comparing model methodologies

Data

Modeling

Conclusion

Average AIC & RMSE by Modeling Methodology



**ARIMA's show the lowest average training RMSE & AIC across all routes considered, followed by Linear Regression with ARIMA errors. As such, we focus our analysis on these two methods.**



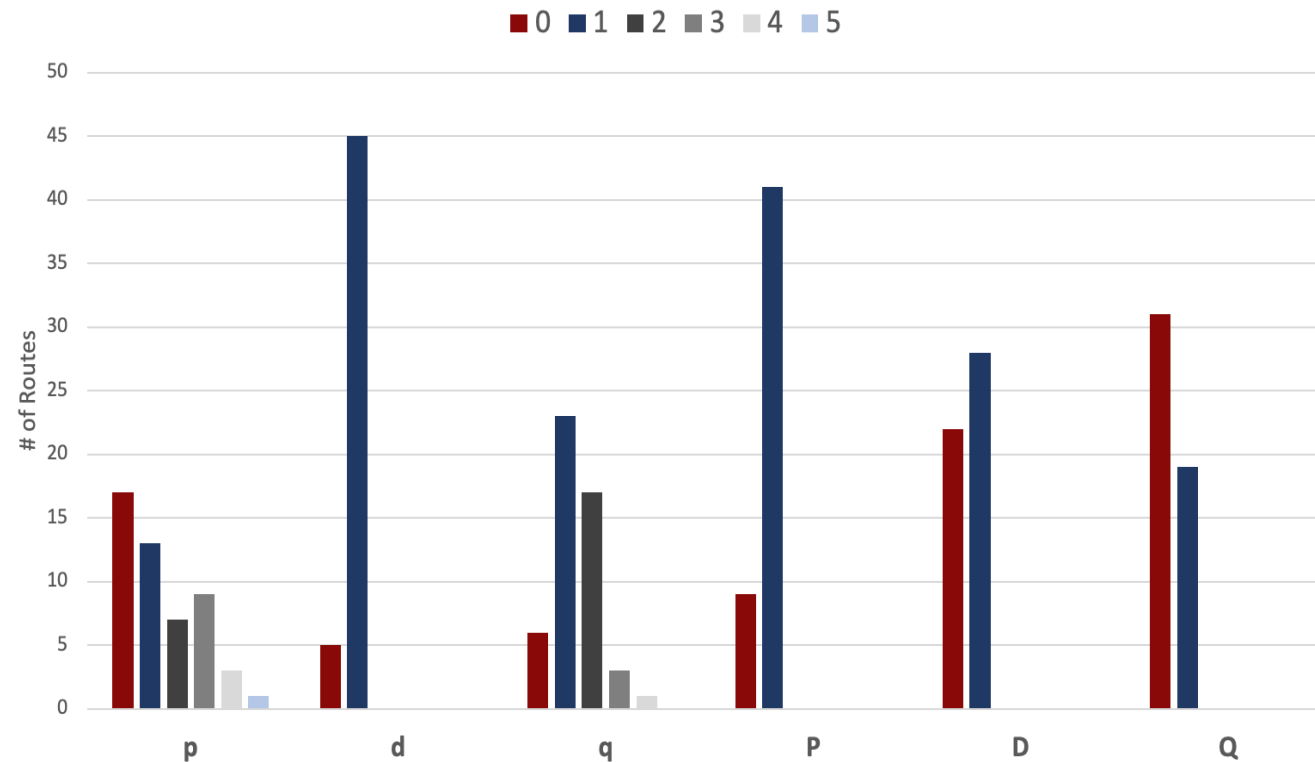
# Identifying Candidate Models (Arima)

Data

Modeling

Conclusion

ARIMA(p,d,q)(P,D,Q) Specifications Across All Routes



Candidate Models

#	Models	% of routes covered
1	ARIMA(0,1,1)(1,1,1)[52]	10%
2	ARIMA(0,1,1)(1,1,0)[52]	8%
3	ARIMA(0,1,2)(1,0,0)[52]	6%
4	ARIMA(1,1,2)(1,1,0)[52]	6%

Across all routes, optimal ARIMA models contain 0-1 non-seasonal (p) and seasonal (P) autoregressive terms, 1 non-seasonal (d) and seasonal difference (D), and 0-2 lagged seasonal (q) and non-seasonal errors (Q).





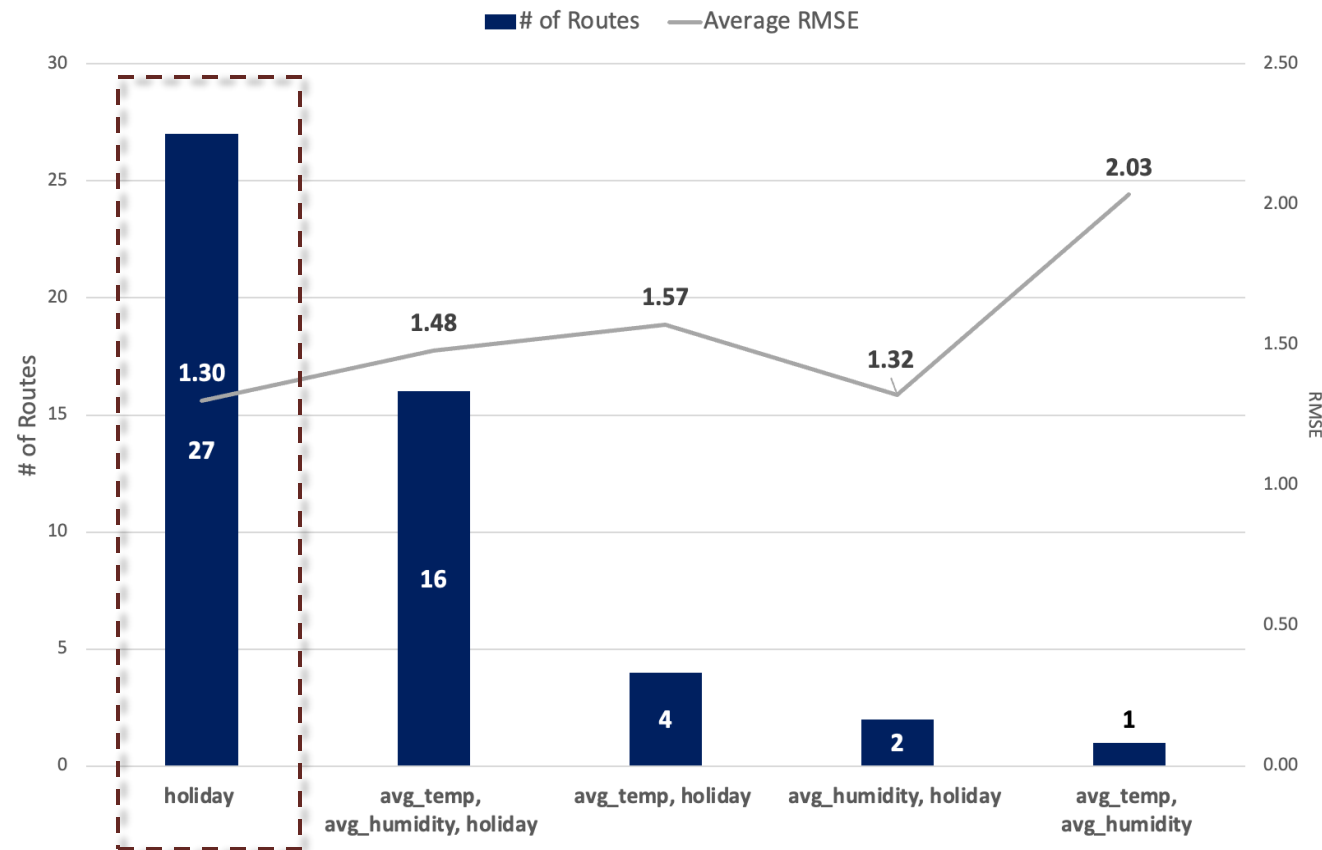
# Identifying Candidate Models (Regression w/errors)

Data

Modeling

Conclusion

Regressors present in Linear Regression (with Arima Errors) across all routes



Candidate Models

#	Models	% of routes covered
1	X = holiday (Errors: ARIMA(0,0,1)(1,0,0)[52])	4%
2	X = holiday (Errors: ARIMA(1,0,1)(1,0,0)[52])	4%
3	X = holiday (Errors: ARIMA(0,0,1)(1,1,1)[52])	4%
4	X = holiday (Errors: ARIMA(5,0,0)(1,0,0)[52])	4%

Linear regressions with a single holiday indicator variable outperforms models with additional variables.





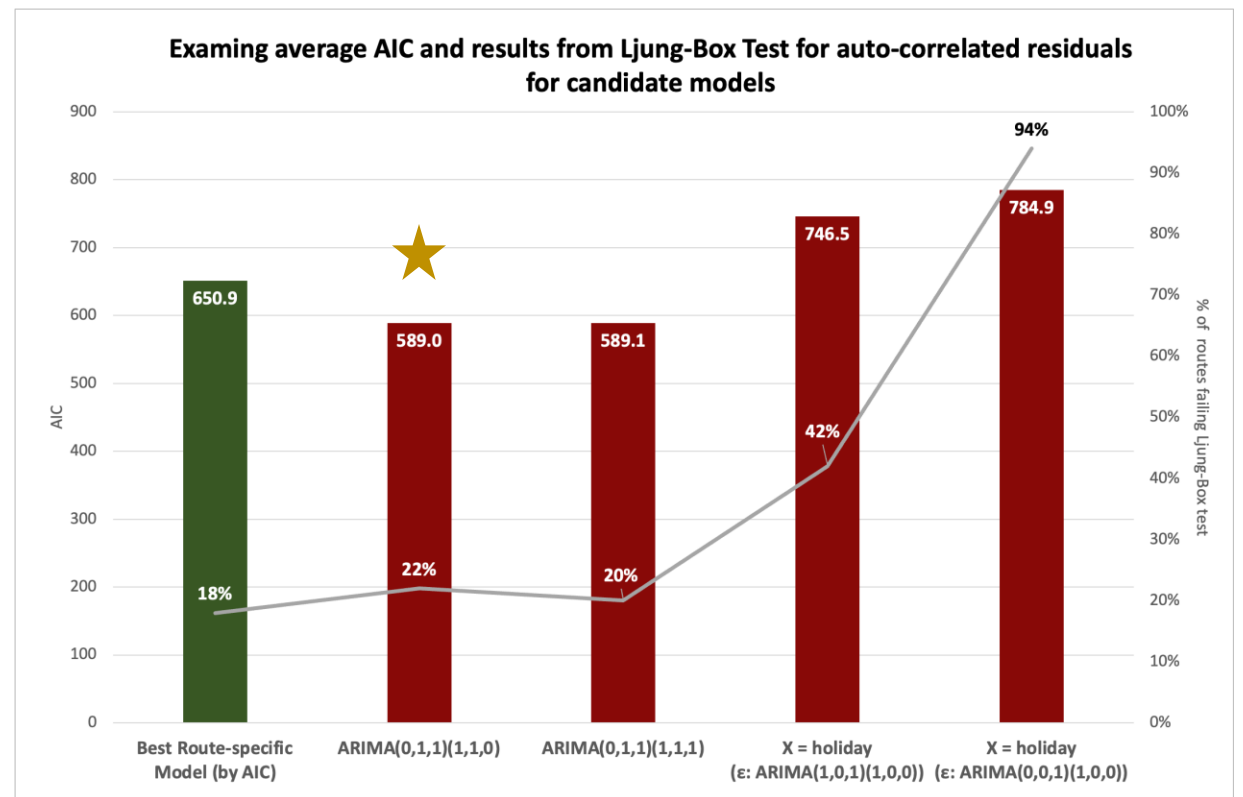
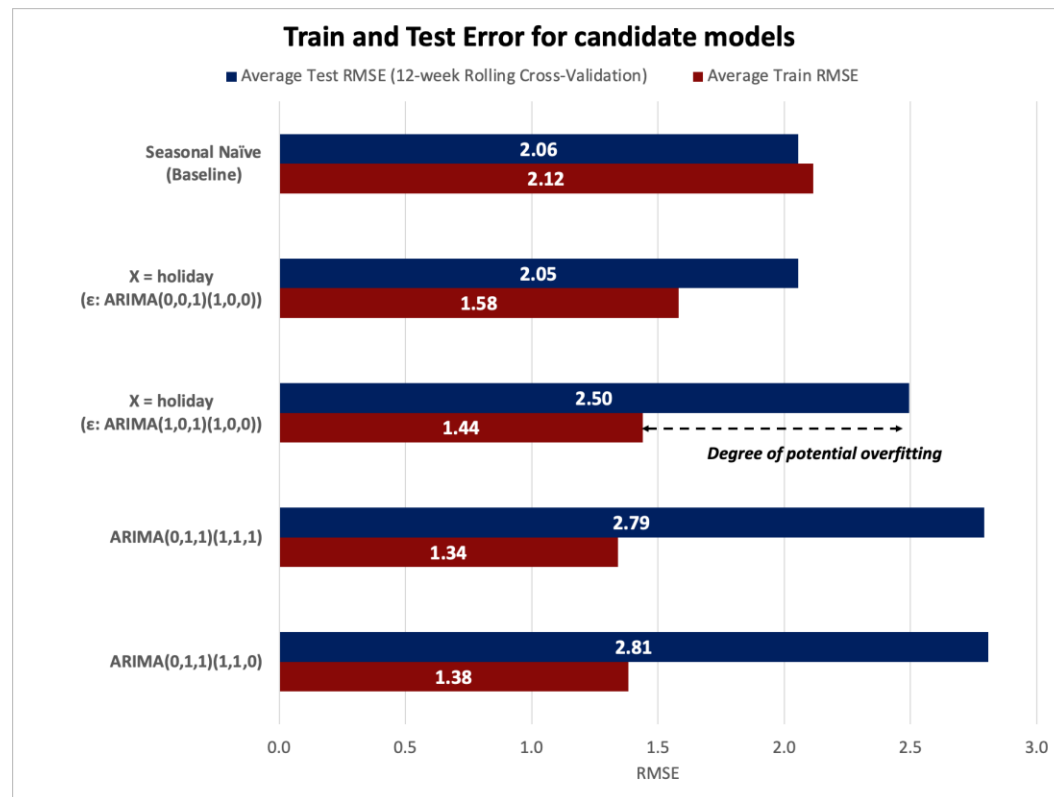
# Meeting our champion, Arima (0,1,1)(1,1,0)

Data

Modeling

Conclusion

This model shows the best balance across performance, simplicity, and robustness



# Net Takeaways & future considerations

Data

Modeling

Conclusion

One Generalized Model to  
"rule them all"

## MODEL RECOMMENDATION SPECTRUM

"We've got compute power  
for days"

Arima  
(0,1,1)(1,1,0)

Arima & Regression with Errors  
with fixed parameters

Arima & Regression with Errors with varying  
parameters

Best Model,  
per route

### Conclusions for Uber

**We built a fully automated modeling suite empowering Uber Operations** with the option to build route-specific models & more generalized models.

**We have model recommendations across it all.** Ultimately it boils down to a choice of simplicity & operational efficiency or large-scale improvement & better fit.

### Future considerations

**Investigation into what's important for Uber.** Did our problem statement address a pressing need?

**More routes, more locations, more test & learn.** We need to test our modeling assumptions more broadly;

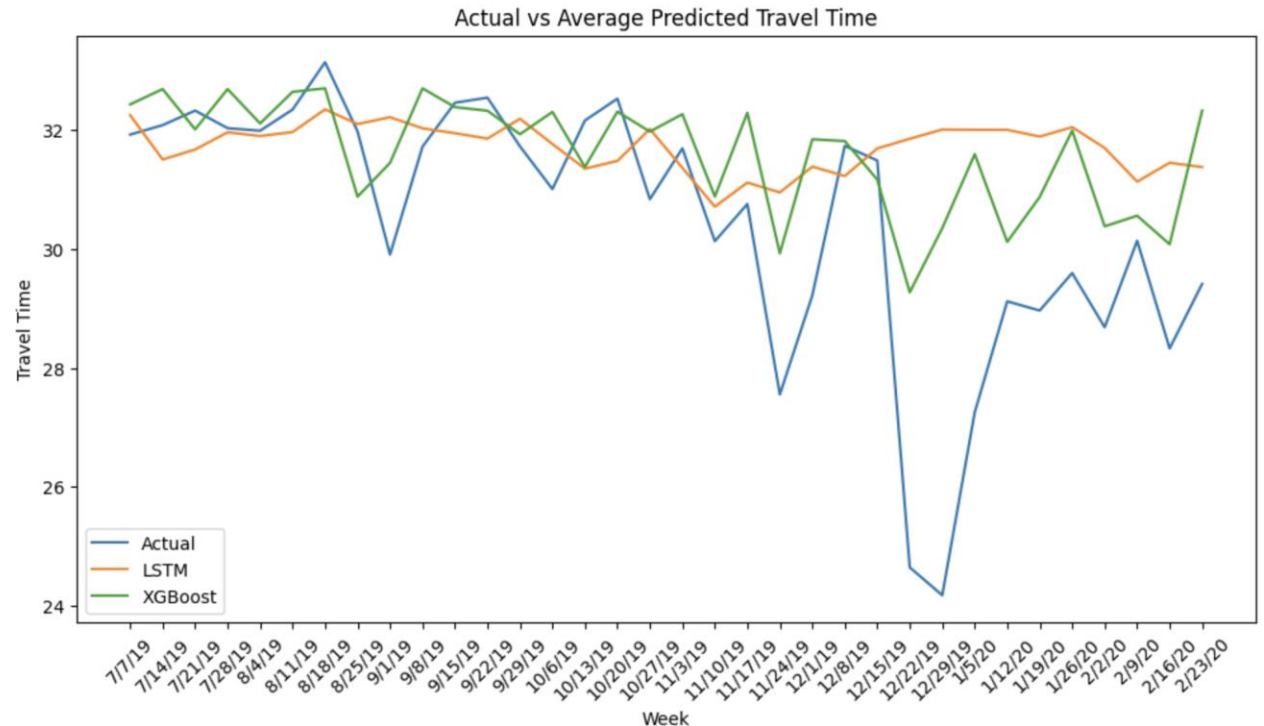
- New starting points, end points, across major metros, beyond urban, and extended globally
- Including new externalities. We leveraged weather & Holidays in our modeling. But what are Uber's most important additional features?



# BONUS: LSTM & XGBoosted Regression

## Methodology:

- Computational Efficiency → used average travel times of all locations, instead of average output of every model iteration through the locations.
- Experiment with more sophisticated models to assess their prowess to yield more performant results.
- Strong outcomes can be achieved with LSTM and XGBoost models.
- Results are close to the evaluation scores of the core 5 route-specific models → a model that incorporates all locations might also be viable.
- Details on each model's methodologies are in the appendix section.



## Model Evaluation RMSE:

XGBoost= 1.9496

LSTM = 1.0841



THE UNIVERSITY OF  
**CHICAGO**

**THANK YOU**



# Contributions

**The entire team touched on nearly every aspect of this project.** We met weekly (9 hours in total), coordinated deliverables, and paced, as a group, to our finished presentation. **Below we have only highlighted the components each member over-indexed in.**

Deepak Vanjani = R & Python Scripts (Web Scraper, Automation, Aggregation, Reporting), Modeling, Project Brief

Danil Meresenschi = Model Recommendations, Data Strategy, Parallel Python work

Dharmi Gala = Model Recommendations, Data Collection, Outcome Strategy

Prayut Jain = Model Recommendations, Data Collection, Outcome Strategy

Christopher Marasco = Project Management, Data Collection, Narrative & Presentation

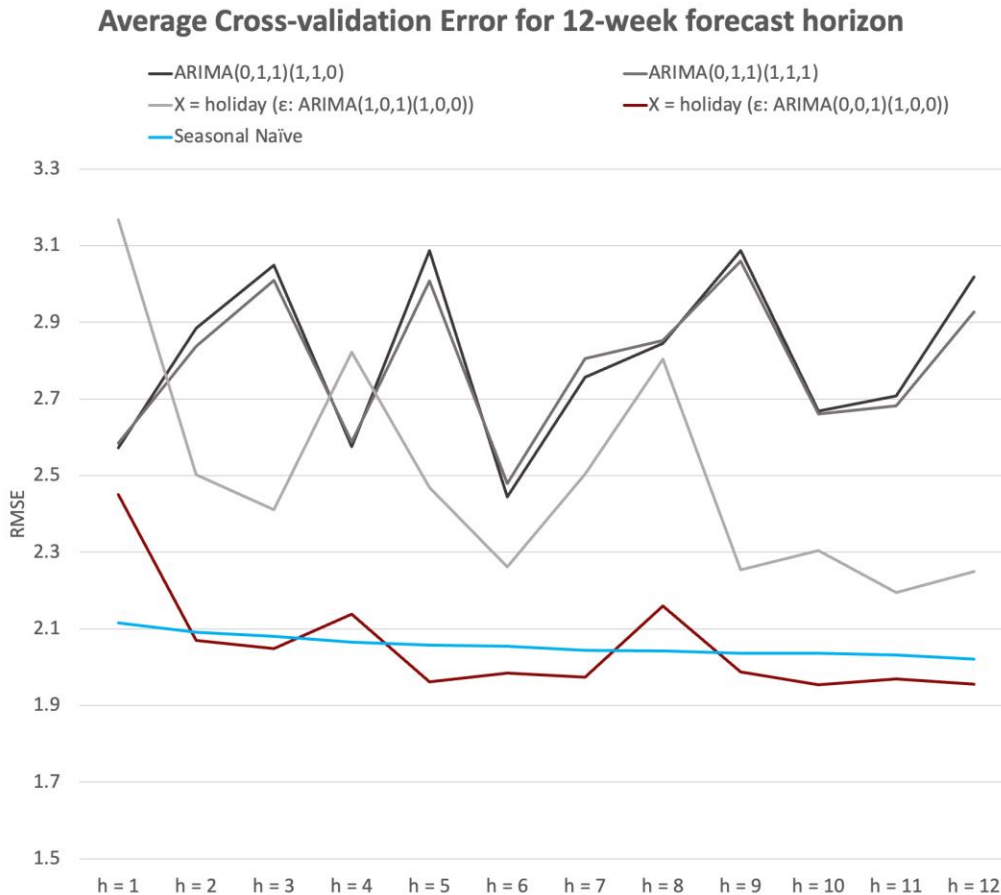


THE UNIVERSITY OF  
**CHICAGO**

# APPENDIX



# Appendix 1. Cross-validation results for candidate models



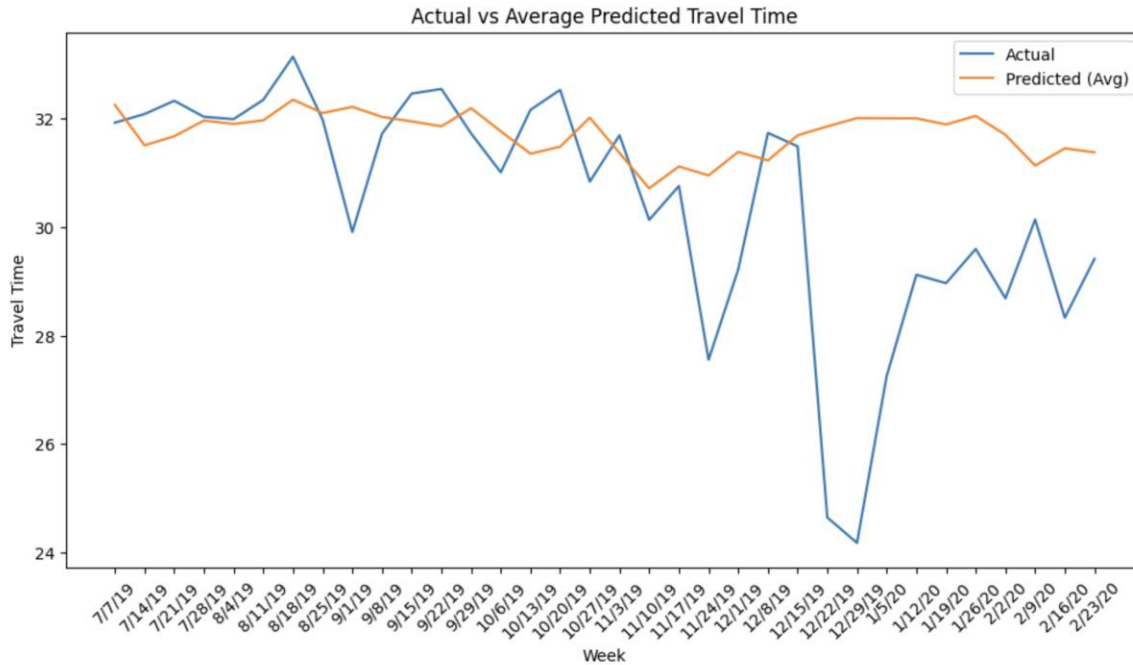
## Methodology:

- Applying cross-validation with rolling 12-week windows shows that X = holiday (Errors: ARIMA(0,0,1)(1,0,0)) performs the best in terms of average RMSE at each forecasting horizon





# Appendix 2. LSTM



## Model Evaluation:

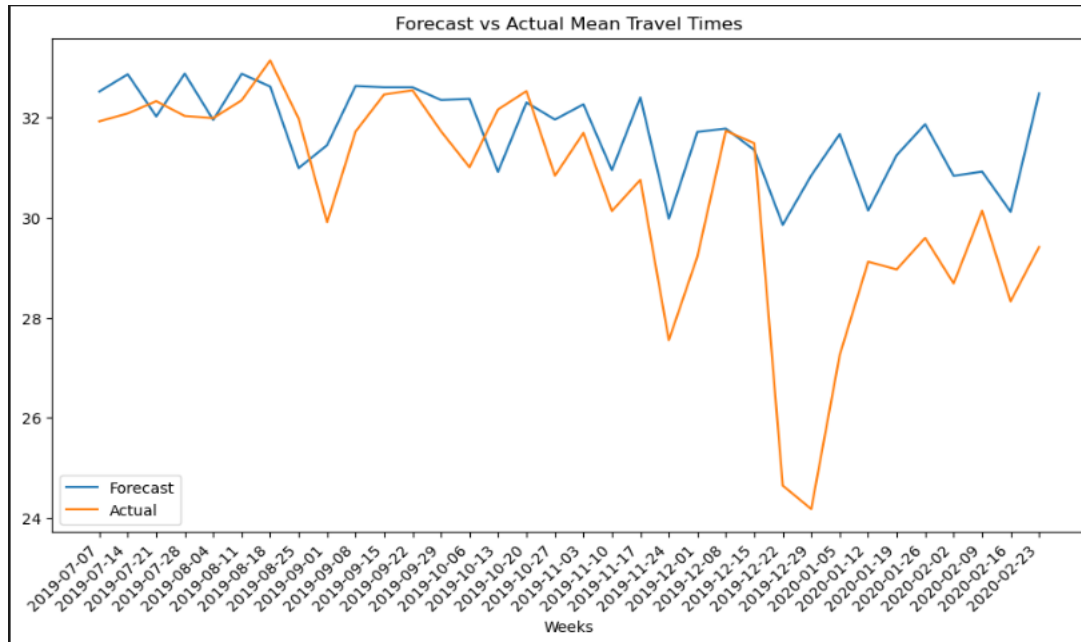
RMSE = 1.0841

## Methodology:

- Used normalized variables.
- Cross-validated using TimeSeriesSplit with 5 folds
- Implemented with 64 units and ReLU activation function.
- Compiled with the Adam optimizer and mean squared error loss function
- Trained using the training data for 100 epochs and a batch size of 32



# Appendix 3. XGBoost Regression with Randomized Search.



## Order of important features:

1. Holiday (True/False)
2. Average Temperature
3. Average Humidity

## Model Evaluation:

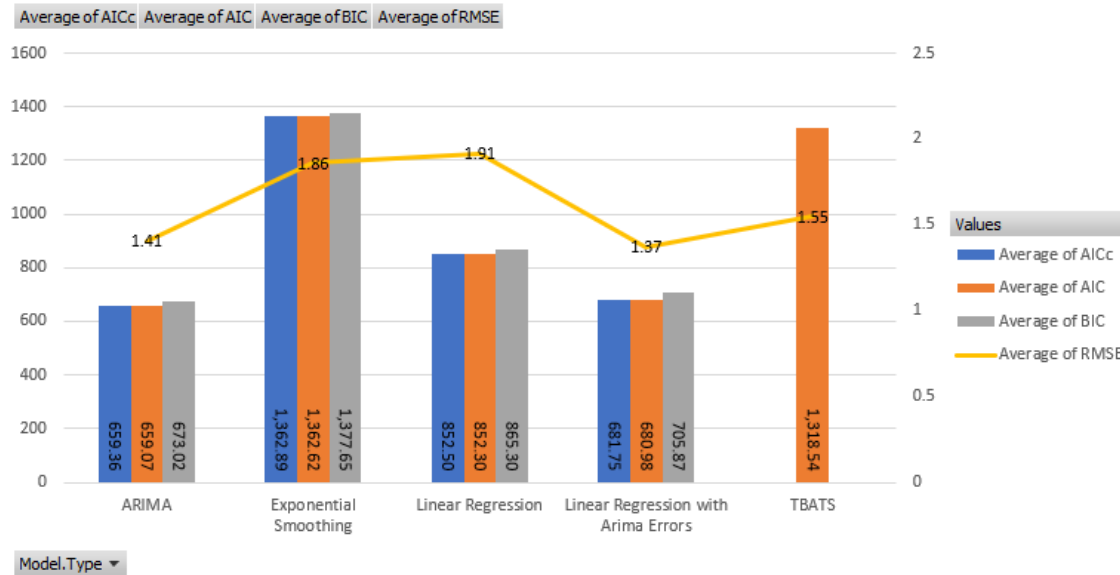
RMSE = 1.9496

## Methodology:

- Utilized Randomized Search to determine the best parameters.
- Implemented 5-Fold Cross-Validation in Randomized Search.
- The best XGBoost Regression model had the following parameters:
  - Learning Rate = 0.1
  - Max Depth = 3
  - Gamma = 1
  - Number of Estimators = 100



# Appendix 4. Tourist Locations



## Methodology:

- We included tourist-centric locations in our analysis to provide an opportunity to isolate the seasonality of travel times to potentially episodic locations.
- But from what you can see, the best options remain ARIMA and Linear Regression with ARIMA errors!

