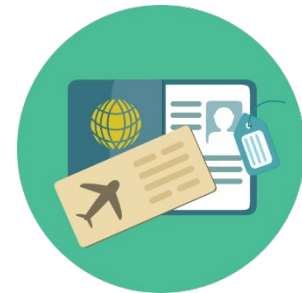


Introduction

Diana Shen

A Little about Myself



DataScience@SMU

On-Time Performance Forecast Project

Diana Shen

What Is On-Time Performance?

- On-time performance (OTP) is one of the key KPIs that are closely monitored in the airline industry. It measures the ability of an airline to be operated on time. Specifically, it is a metric calculating the percentage of on-time flights remaining on the published schedule.
- The OTP used in this project is a daily metric, determined by the number of arrival on-time flights vs. total number of scheduled flights per day, following the 15-minute on-time rule.
 - $OTP = (\text{total on-time flights} / \text{total scheduled flights}) * 100\%$
- 15-minute on-time rule: a flight that arrives within 15 minutes of its scheduled arrival time is considered on time

Project Background

- My team: a data science/operation research consultancy group
- My internal customers/stakeholder
- Analytics in airline operations: challenges and opportunities
- Opportunities and threats for our team
- Lineup for this project

Why Predicting OTP?

- POC: breaking the ground for predictive analytics in airline operation using weather
- Strategic reasons
- Providing decision-driven insights

DataScience@SMU

Data Cleaning and Wrangling

Diana Shen

Data Overview

- What factors could impact on-time performance?



Seasonality



Weather



Passenger



Crew

- Data considered in this project:
 - Weather forecast data
 - Passenger data: load factors
 - Schedule data: daily number of scheduled flights
 - Seasonality: day of week, month
 - Yesterday's OTP

Working with Weather Data

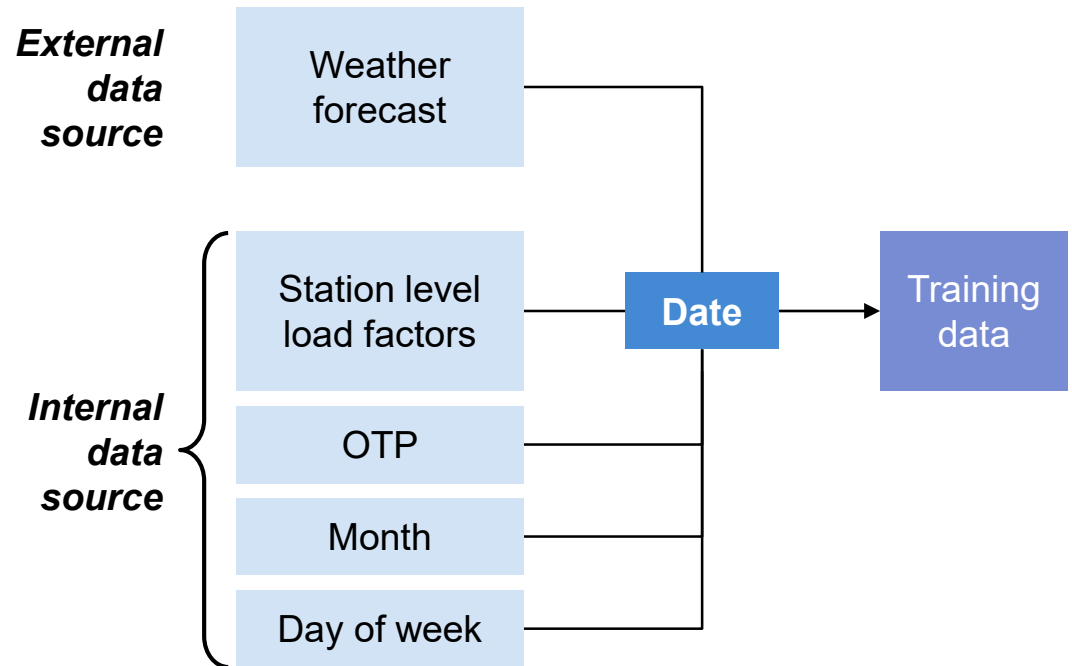
- Multiple weather forecast data sources to choose from:



- Things to consider when choosing the most suitable weather forecast source:
 - Data quality
 - Data availability
 - Data discrepancy
 - Easiness to obtain

Data Blending

- Querying external data: Python web scraping
- Querying Internal data: SQL
- Blending: Python
- Remember to do sanity check

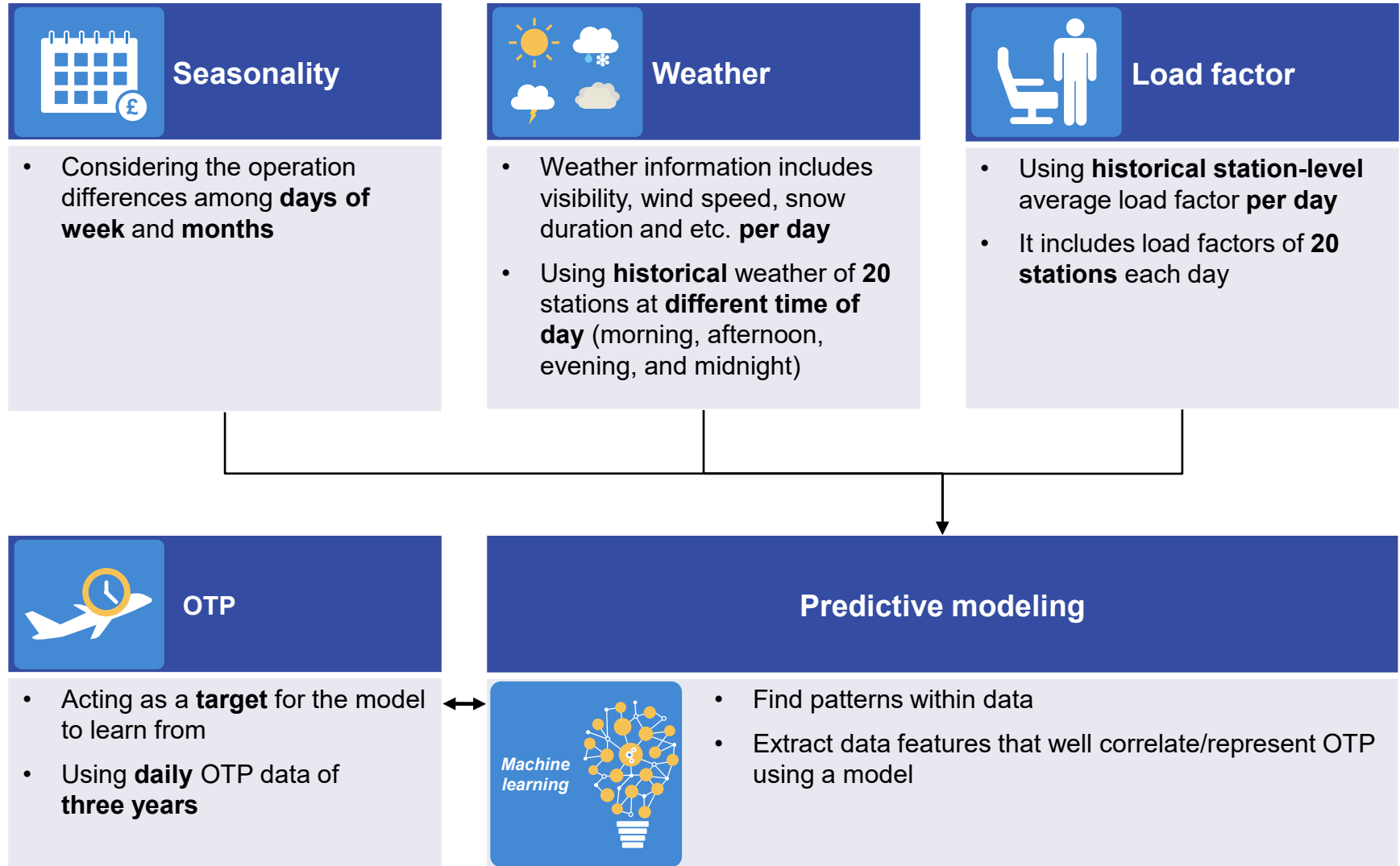


DataScience@SMU




Modeling

Diana Shen

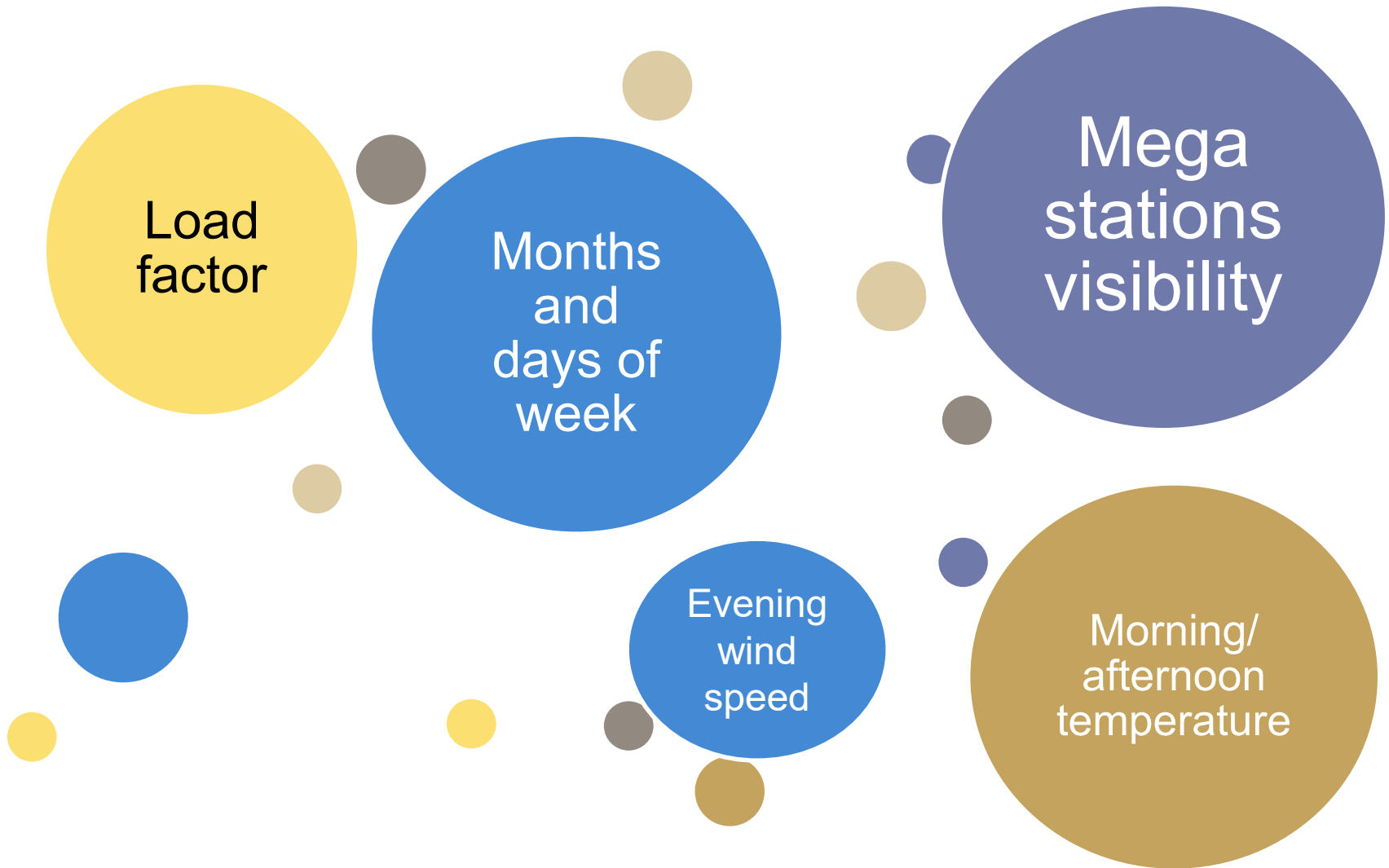
Overview of OTP Model



Machine Learning Models Investigated

 <p>Tree-based models</p>	Gradient Boosting Model
	Random Forest
 <p>Linear models</p>	Multivariate Linear Regression
	Lasso
	Ridge
Artificial intelligence	Neural Network 

Influential Variables in the Model



* Bubble size is proportional to influential level

Cross Validation

- Variables used vs. data points available for training
 - Why three years of training data?
- Benefits of using cross validation
- **RMSE**

Final Model

- Factors considered when choosing a model:
 - Performance
 - Model accuracy
 - Dimension reduction
 - Interpretability
 - Time needed to run
- Lasso is the final model implemented in this project

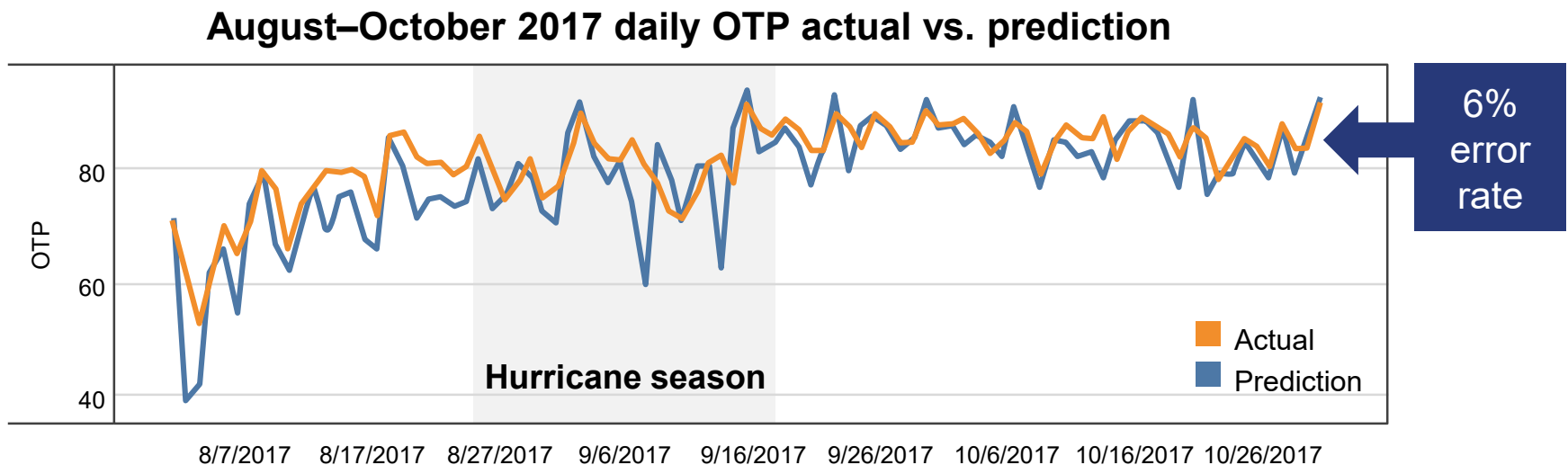
DataScience@SMU

Visualization

Diana Shen

Visualization on Model Performance

What plot could be used? Q-Q plot? Scatterplot?



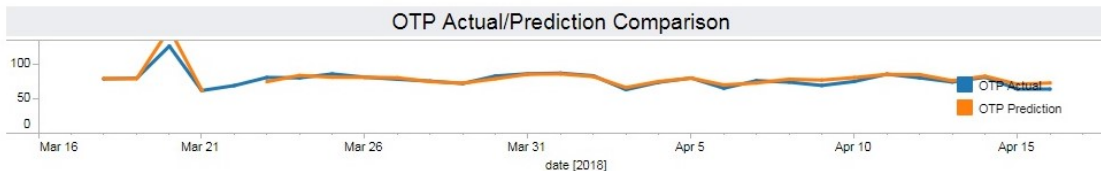
Dashboard Presented to Clients

OTP Prediction for Tuesday, April 17, 2018

Lower Bound
77.877

OTP Prediction
82.601

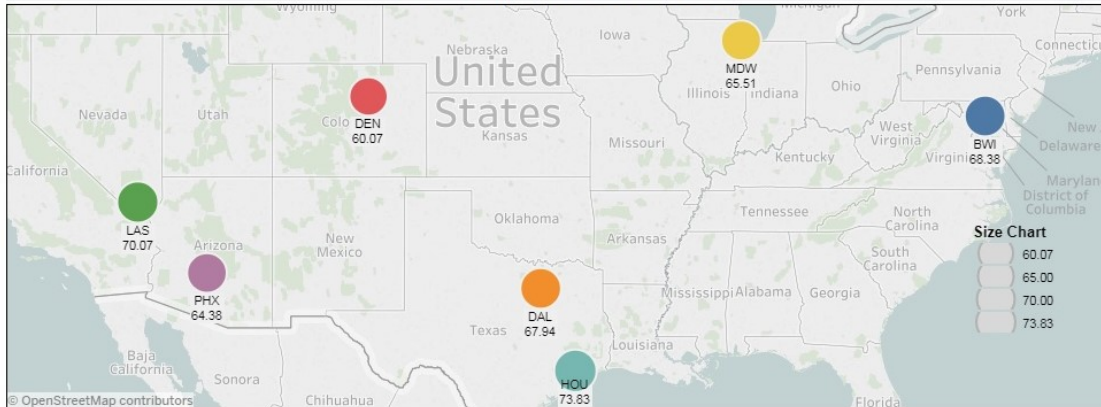
Upper Bound
87.006



Top 10 Factors Explaining Today's Prediction Results

-0.0070	-0.0087	-0.3525	-0.0165	-0.0167	-0.0180	-0.0091	-0.0284	-0.1437	-0.0421
LF_KORF	mon_11	OTP	P12_KDTW	P12_KPHL	SNW_KROC	TMP_KECP	TMP_KMSP	total_fit	WDS_KDEN

Station Turn Compliance Prediction



← A prediction with confidence intervals

← OTP comparison for the last 7 days

← Show what factors influence most on the prediction

← A map view showing station level predictions

DataScience@SMU

Deployment

Diana Shen

Predicting Tomorrow

- When to run?
- When will the data be ready?
- Channel for showing the prediction

Model Maintenance

- Model retraining frequency
- Process of maintaining the model/dashboard

DataScience@SMU

Business Challenges

Diana Shen

Customer Feedback

DataScience@SMU

Future Directions

Diana Shen

Further Expand the Project

- Network level OTP—station level OTP—flight level OTP (flight delays)
- Real decision-driven models
- More models informing customers of potential upcoming repercussion

DataScience@SMU