# Analysis and prediction of US flights delay
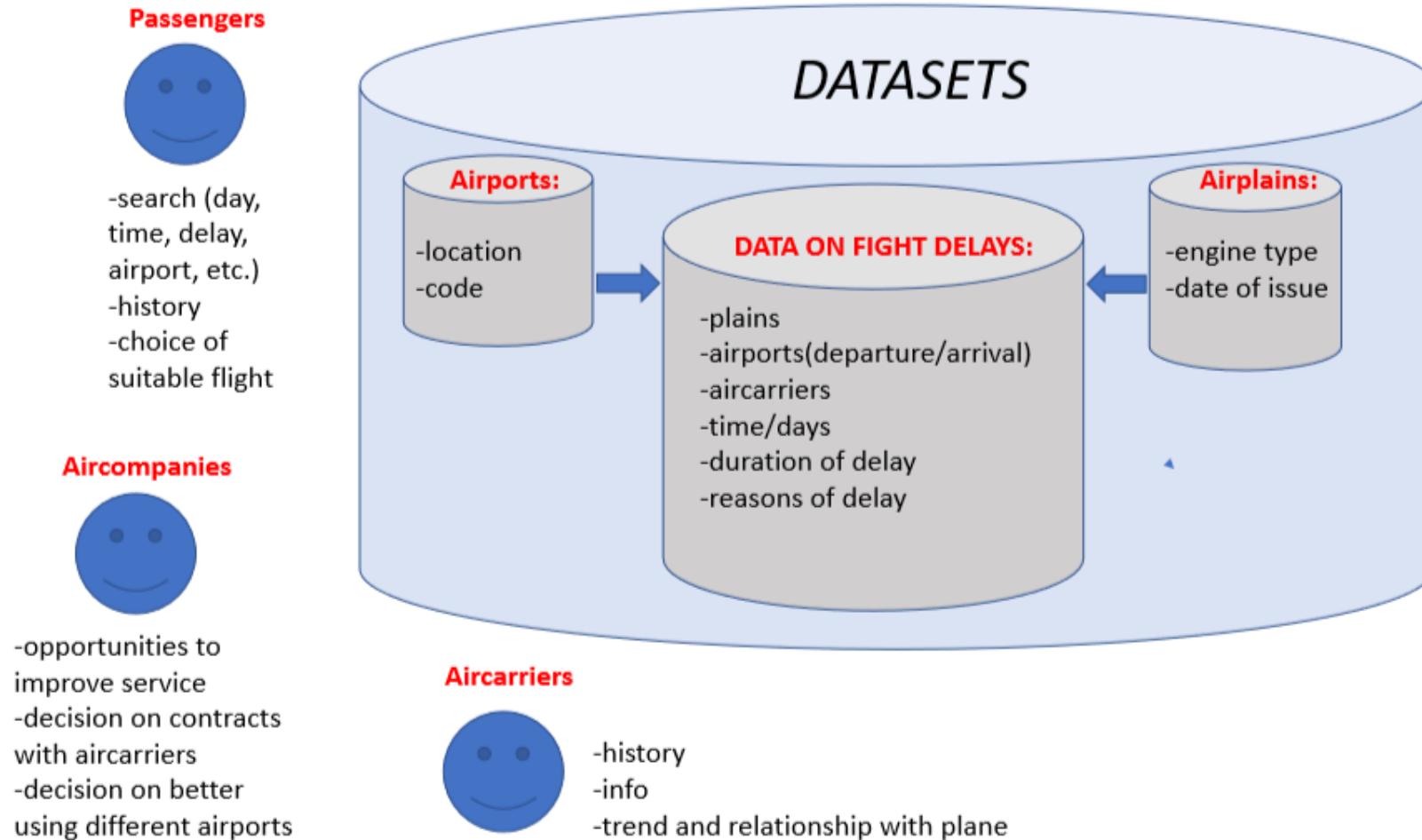
by Korchagina Evgeniya

(CEBD 1260)

# Datasets

# Some use cases

**Passengers-search &Passengers-choice:**

-should I take flight at this time with this aircompany for urgent business trip
-should I take this flight for family trip (how reliable aircompany)
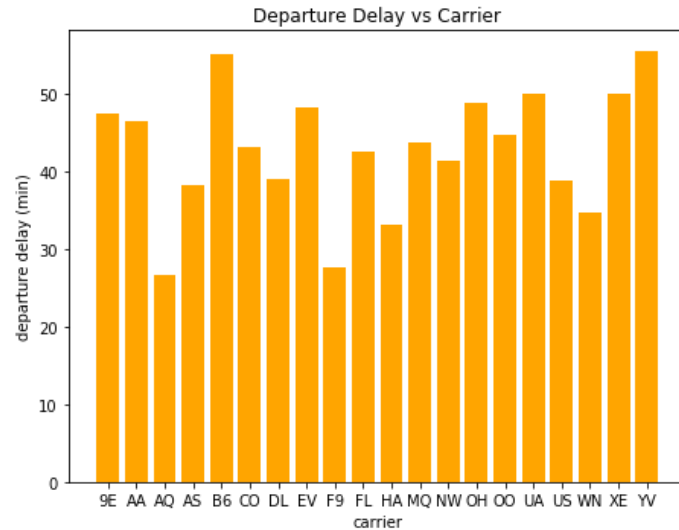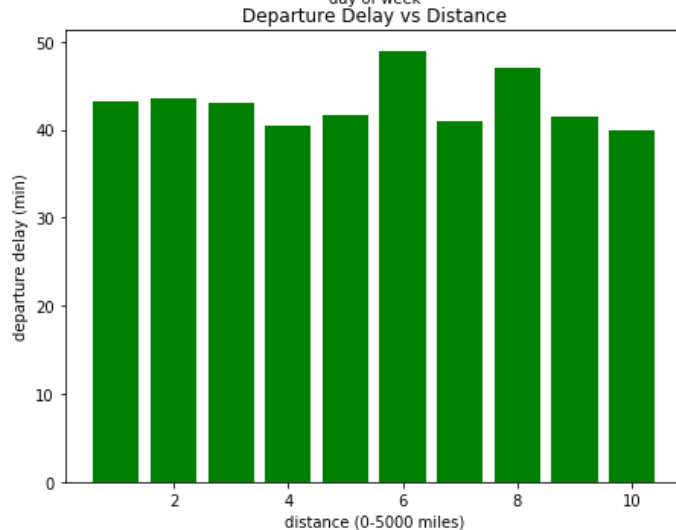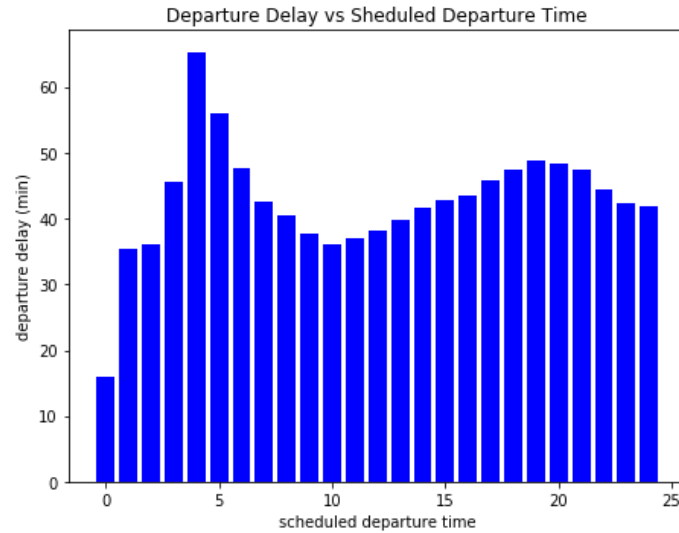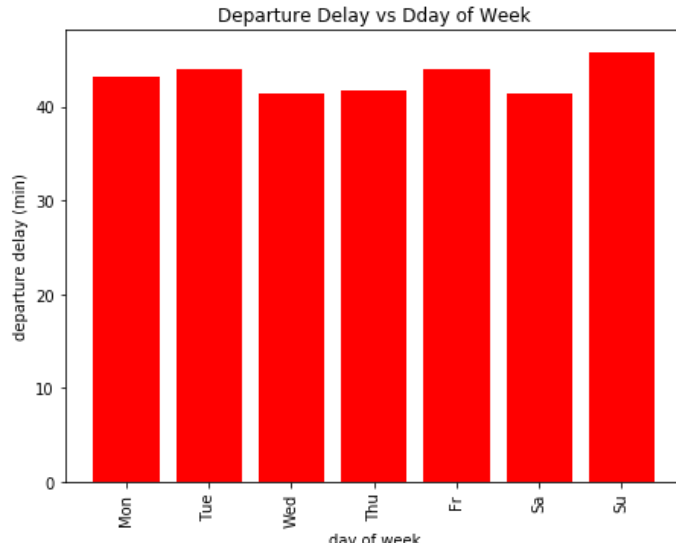-how high chances to have delay with this particular flight considering weather conditions
-etc.

**Aircompanies-opportunities:**

-can we change boarding time for long-distance trips to avoid delays?
-can we change time of flights considering weather stats on time of day?
-should we use higher technical standards for planes (in case if the cause of delay is technical issue)
-should we consider some compensation to passengers to avoid their frustration?

**Aircarriers:**

-Do we need improve quality of aircraft to avoid technical issues (weather resistance, engine quality and time lasting, etc.)
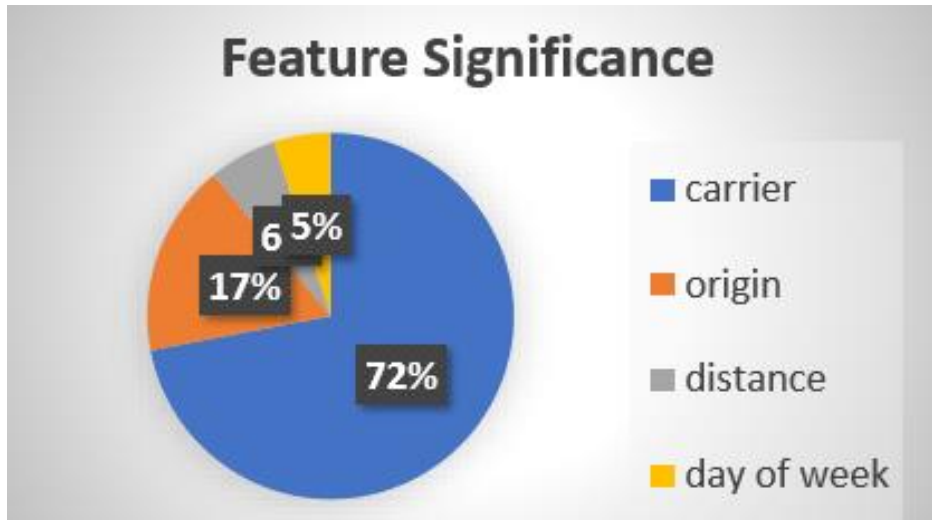-shall we do more often check up of planes travelling to cold countries/long distances
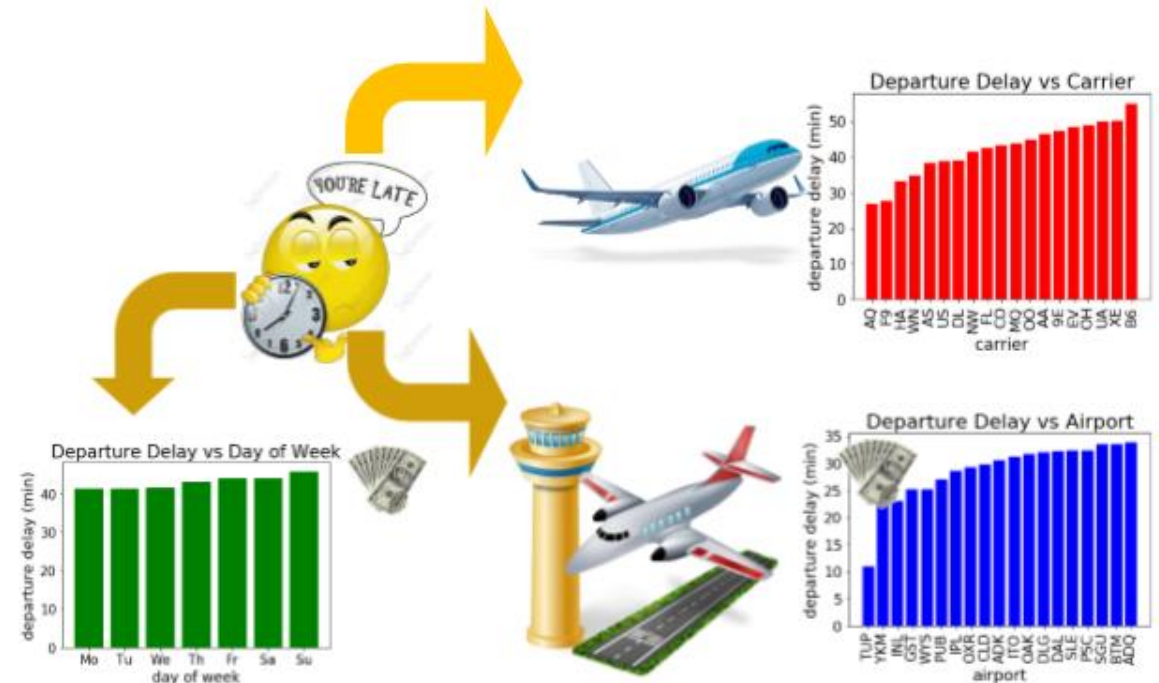
# Preprocessing and EDA



**Based on EDA:**

Current project will be focused on business use case concerning passengers search of flights covering some specific criteria. Result of data analysis suppousably will lead to an application that will predict delay time based on carrier, distance of flight, time of departure, airport of origin, tail number and day of week. Modeling will be based on regression analysis of these variables.

# Technical and business perspective



*DecisionTree Regressor and RandomForestRegressor (used for categorical data) showed the same absolute errors and standard deviations. Whiskers plot were visually undistinguishable.*
*>>any of the methods can be applied as a regression model*



**Scheme1.** *Scheme representing the main contributing factors to flight delay. The lower plots indicate money, as a factor affecting bar plot tendency.*

*Here we are focusing on prediction of possible flight delay based on the flight features (from passenger' perspective).*
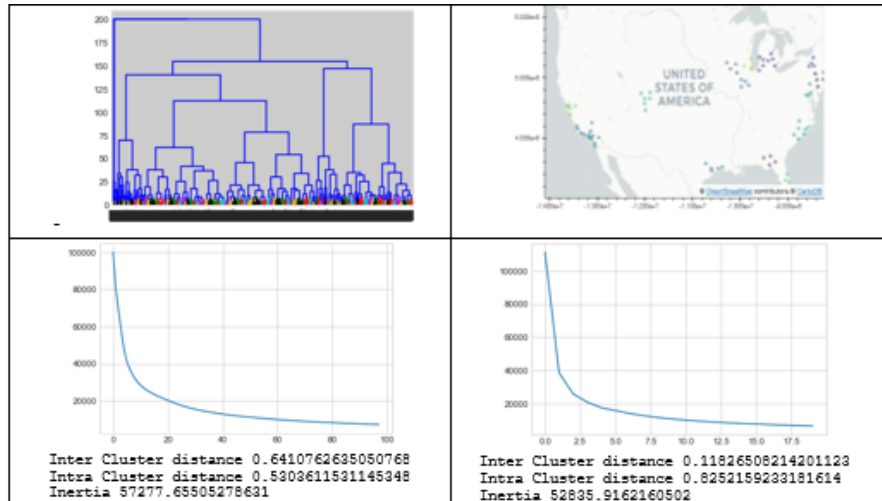
# Technical challenges



Fig.1. Top left – dendrograms of hierarchical clustering; top right – map of US airport created via DBSCAN; bottom left – inertia of hierarchical clustering; bottom right – inertia of K-means clustering mechanism.

APP:

- *Big size of dataset*

- *Parsing categorical value to dummy numerical ones*

- *time consuming regression processing on big set of data (impossible to apply models with categorical data of high variety)*

- *Hierarchical clustering and K-means do not for chosen data*

- *Due to the lack of time DBSCAN was used only for simple practice*

- *APP to be finished*