

# AE 10: Model comparison

Add your name here

## Packages

```
library(tidyverse)
library(tidymodels)
library(knitr)
```

## Data

For this application exercise we will work with a dataset of 25,000 randomly sampled flights that departed one of three NYC airports (JFK, LGA, EWR) in 2013.

```
flight_data <- read_csv("data/flight-data.csv")
```

Rows: 25000 Columns: 10

-- Column specification -----

Delimiter: ","

chr (4): origin, dest, carrier, arr\_delay

dbl (4): dep\_time, flight, air\_time, distance

dtm (1): time\_hour

date (1): date

i Use `spec()` to retrieve the full column specification for this data.

i Specify the column types or set `show\_col\_types = FALSE` to quiet this message.

1. Let's get started with some data prep: Convert all variables that are character strings to factors.

## Modeling prep

2. Split the data into testing (75%) and training (25%), and save each subset.

```
set.seed(222)
```

3. Specify a logistic regression model that uses the "glm" engine.

Next, we'll create two recipes and workflows and compare them to each other.

## Model 1: Everything and the kitchen sink

4. Define a recipe that predicts `arr_delay` using all variables except for `flight` and `time_hour`, which, in combination, can be used to identify a flight. Also make sure this recipe handles dummy coding as well as issues that can arise due to having categorical variables with some levels apparent in the training set but not in the testing set. Call this recipe `flights_rec1`.
5. Create a workflow that uses `flights_rec1` and the model you specified.
6. Fit this model to the training data using your workflow and display a tidy summary of the model fit.
7. Predict `arr_delay` for the testing data using this model.
8. Plot the ROC curve and find the area under the curve. Comment on how well you think this model has done for predicting arrival delay.

## Model 2: Let's be a bit more thoughtful

9. Define a new recipe, `flights_rec2`, that, in addition to what was done in `flights_rec1`, adds features for day of week and month based on `date` and also adds indicators for all US holidays (also based on `date`). A list of these holidays can be found in `timeDate::listHolidays("US")`. Once these features are added, `date` should be removed from the data. Then, create a new workflow, fit the same model (logistic regression) to the training data, and do predictions on the testing data. Finally, draw another ROC curve and find the area under the curve. Compare the predictive performance of this new model to the previous one. Based on the area under the curve statistic, which model does better?

## Putting it altogether

10. Create an ROC curve that plots both models, in different colors, and adds a legend indicating which model is which.

## Acknowledgement

This exercise was inspired by <https://www.tidymodels.org/start/recipes/>.