# TIME-SERIES MODELING

# A. Project Setup
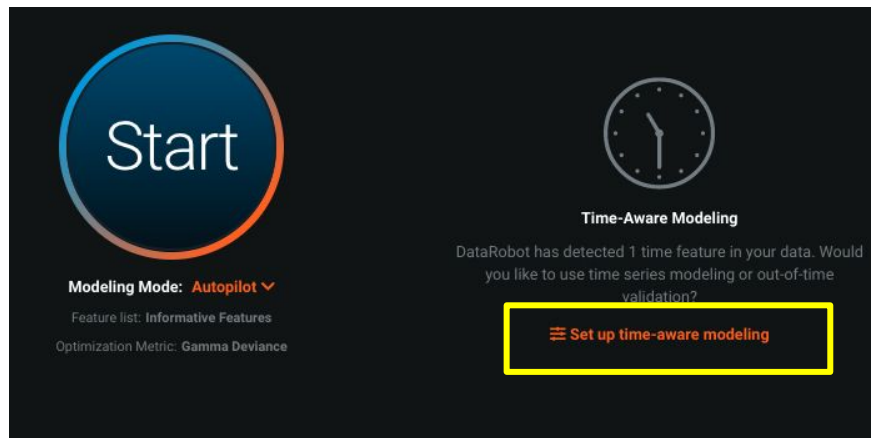
# Data Requirements for Time Series
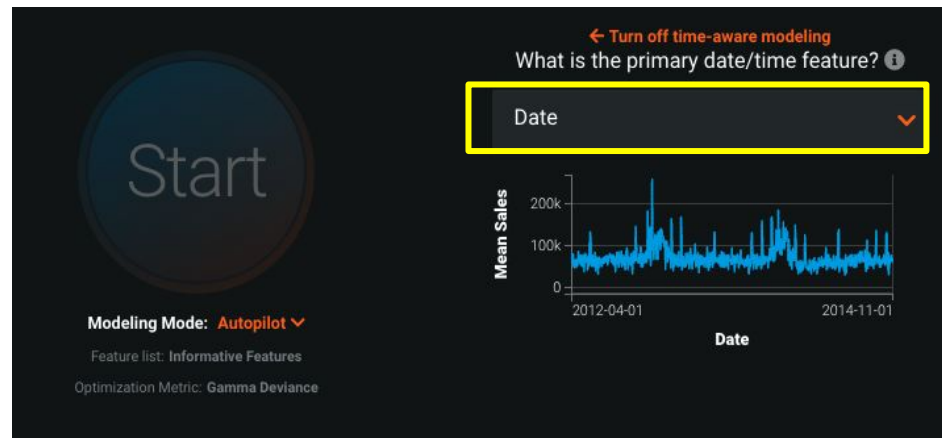
| Date | Hospital | Total Patients | Total Staff | Other Covariates |
|------|----------|----------------|-------------|------------------|
| 3/1/20 | A | 230 | 34 | ... |
| 3/2/20 | A | 219 | 28 | ... |
| 3/3/20 | A | 211 | 25 | ... |
| 3/4/20 | A | 245 | 32 | ... |
| 3/5/20 | A | 249 | 39 | ... |
| 3/6/20 | A | 230 | 41 | ... |
| 3/7/20 | A | 218 | 25 | ... |

- Each row must have unique timestamp; should be in regular intervals (days, hours, minutes, etc.).
- Aggregate rows to the level that makes sense for the use case.
- Multiple series are supported (hospital, department, clinic, etc.)
- Use time series when you want to predict a column into the future.
- Time Series supports covariates (like the AutoML product) - example: "total staff"

# Setting up a DataRobot Time Series Project



After you select your Target column,
select "Set up time-aware modeling"

Supply your date/time column

# Setup your Time Series project & specify series, if applicable

# Feature derivation window and forecast distance...



Feature derivation window tells DataRobot how far back to look to derive features. Short windows will recover to changes quickly, but longer windows will be more stable (and not overreact to spikes).

Forecast distance tells DataRobot how far into the future it should make predictions. This should be based on what provides business value.

# Validation in a time series project (i.e. "Backtests")



- Validation in time series (green or red bars in image) is based on a time range <u>after</u> the training period (blue bars) to prevent leakage and provide accurate performance metrics.
- Project configuration will affect the validation windows somewhat - move these around to cover areas of interest (e.g. before, during, after COVID)

# Additional Setup



1) Known-in-Advance features are columns in your data that the model can use to make the current prediction, such as new product rollout or planned outages, etc.
If you're not sure, leave it <u>off</u> this list.

2) Holidays are also known in advance but can be provided as a separate csv file.
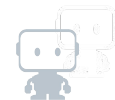
# Ready to start building models?



- Correct any data issues that DataRobot highlighted during setup
- Verify your **feature derivation** windows & **forecast distance** ranges
- Adjust your **backtest** time ranges if needed (or just accept the defaults)
- Double-check your **known-in-advance** features + **holidays**
- Run Autopilot

# B.Evaluation

# Feature List Hierarchy

# Check out the Derived Features



Explore the different feature lists and their distributions.

You can also get derivation log...

# How do we assess quality? | Leaderboard

**Mean Absolute Scaled Error (MASE)** $=$ $\dfrac{\text{Mean Absolute Error (MAE) of } \textbf{Model of Interest}}{\text{Mean Absolute Error (MAE) of } \textbf{Baseline}}$

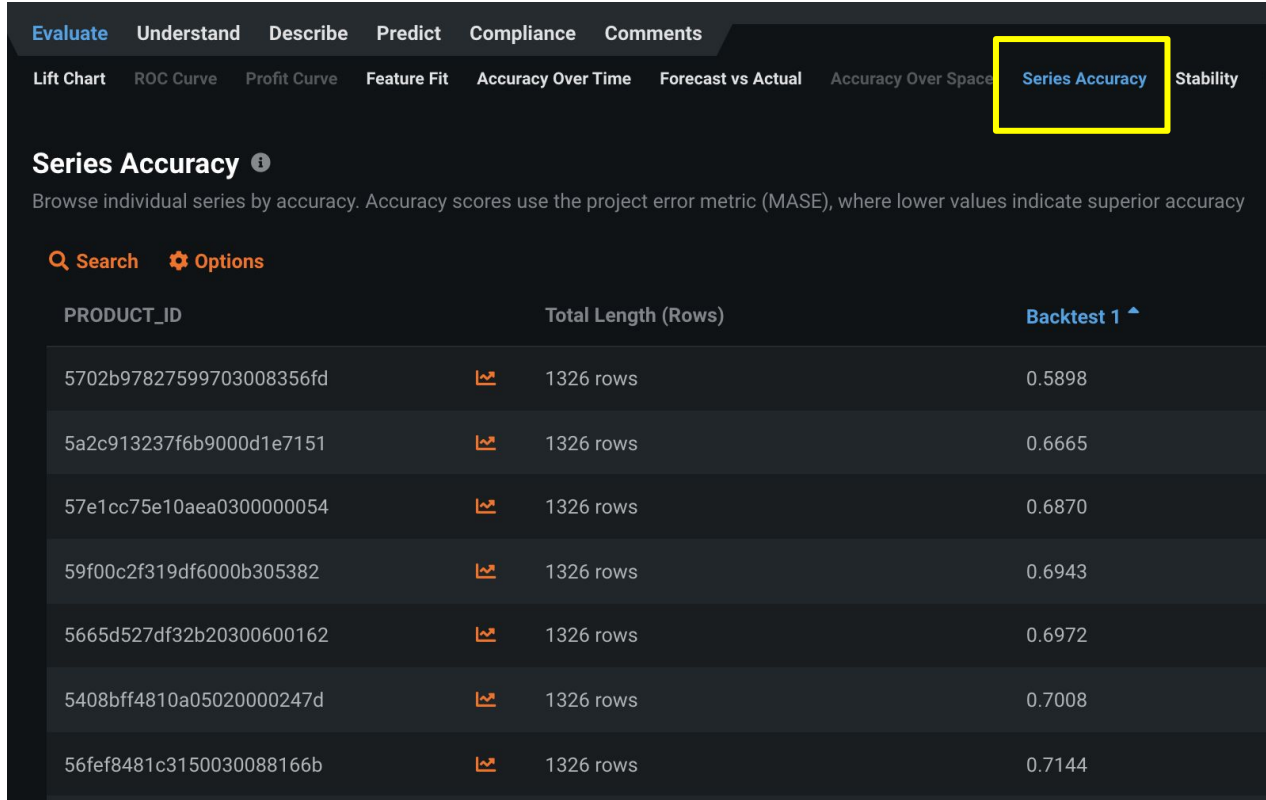# Find the Recommended Model on Leaderboard



You can adjust the Metrics to re-rank the models. There's no accuracy on the recommended model - it's trained to end of the data,

# Stability - How does model perform over time?



For RMSE, lower is better - accuracy is significantly better in backtest 3

# Series Accuracy - which series has best/worst accuracy?



Which series had the best accuracy? The worst?

From here you can identify potential areas for improvement or exclusion.

# Accuracy over Time - evaluate predictions



Can display for each backtest, series, and forecast distance.

# Making Predictions



You must prepare a file with a record for each series and each desired forecast date (i.e. next 7 days…) and any known-in-advance features.

# Prediction Previews & Intervals



The intervals* estimate the range of values DataRobot expects actual values of the target to fall within. They are based on the residual errors measured during the model's backtesting.

* to generate intervals, change and update preview, then download