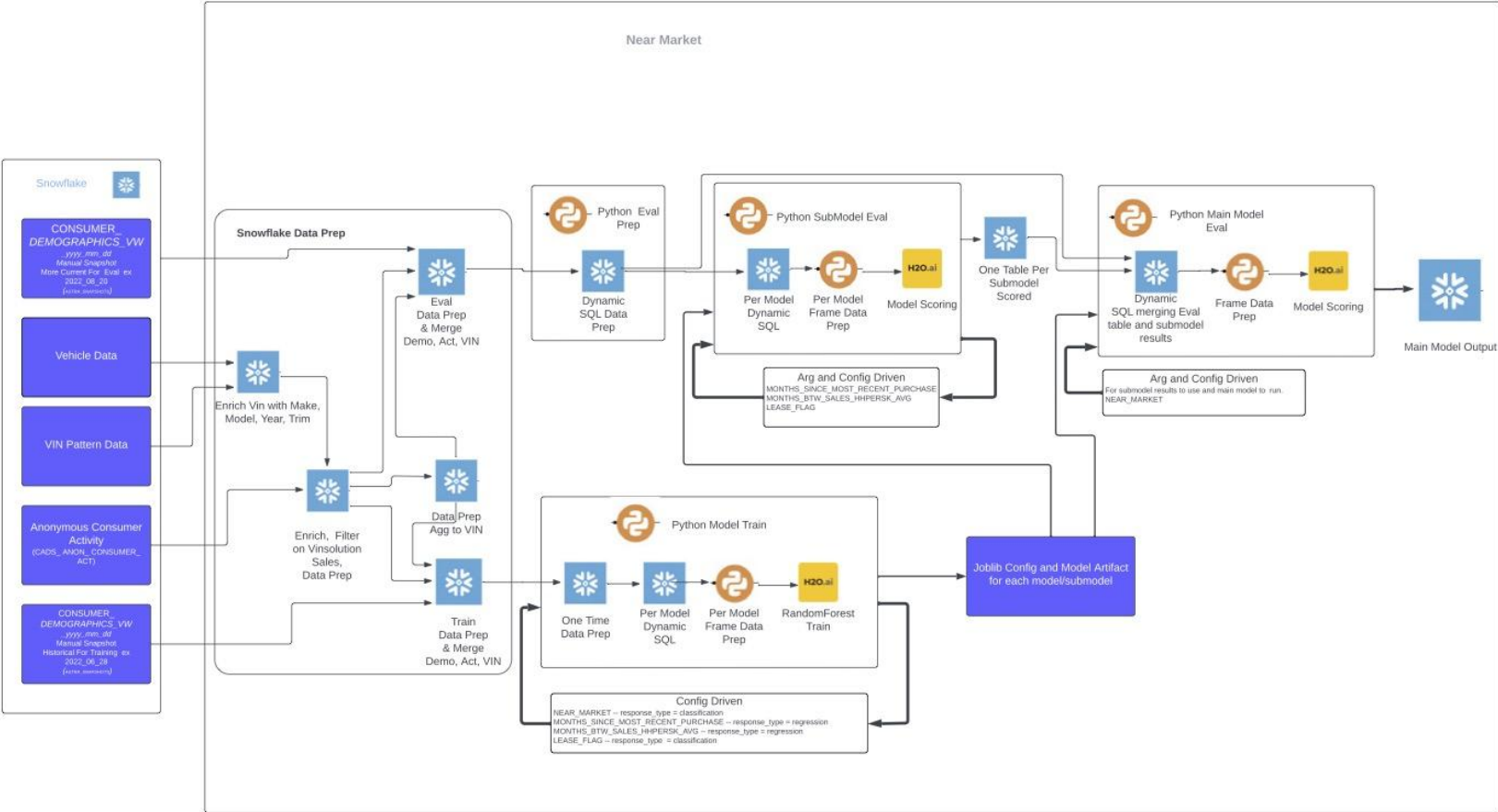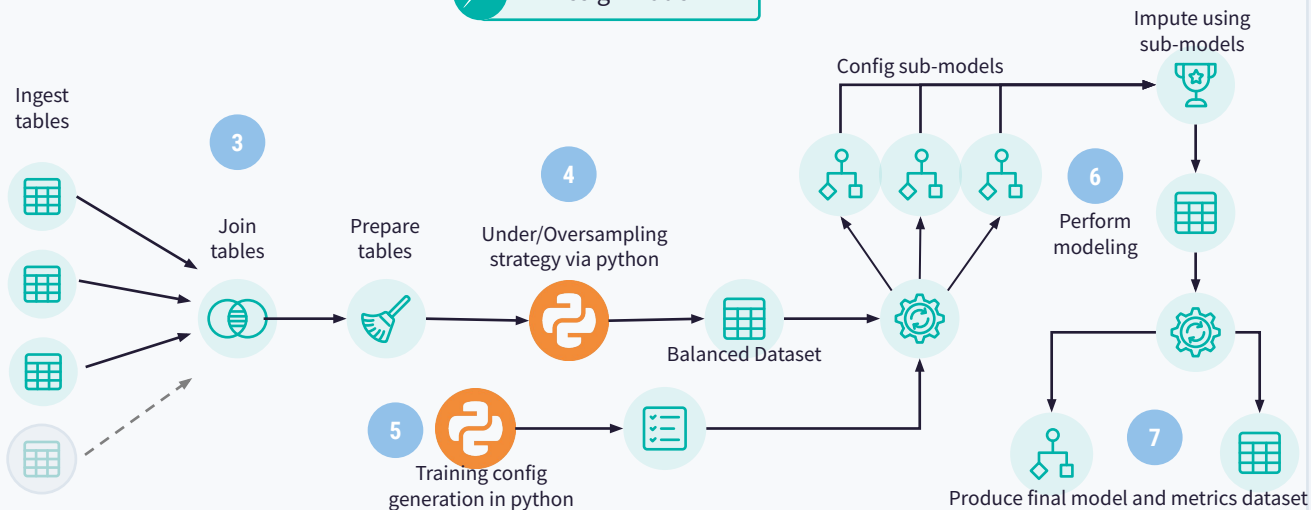**data iku**

### Context

Cox Automotive wants to create a targeted outreach program aimed at certain customers they've classified as "near market", or looking to purchase a vehicle in the next 6-12 months. The project is a lift-and-shift from Domino Data Lab into Dataiku and will utilize Dataiku's more traceable project flow, metrics and checks, and suite of available algorithms. Cox Auto doesn't care about false positives but doesn't want to miss out on potential positives (false negatives), so we recommend using recall as a preferred performance metric when selecting a candidate model.

| | Current State | End State | Actions To Close Gap |
|---|---|---|---|
| **Objective** | The current implementation is extremely complicated and only comprehensible to experienced data scientists, consisting of pure python scripts only. | The new project will make use of Dataiku's more traceable assets like visual recipes whenever possible. The datasets are large so push-down compute to Snowflake will be a centerpiece. The project will be pushed to the automation node for scheduled execution. | 1. Confirm the functionality of Dataiku Design and Automation nodes and establish all required data connections<br>2. Create python code environment with all required package dependencies<br>3. Join and clean the various Snowflake tables as designed on the Lucidchart diagram |
| **Technical Solution** | Domino is currently housing, training, and deploying the H2O random forest models. The project consists of multiple sub-models that effectively impute missing values in the testing set. The "main" model is a binary classifier for whether a customer is near market. | Want to try a series of different algorithms, including importing H2O models into the training loop. Will employ MLFlow to import existing, trained H2O random forest model if necessary. Will explore using app-as-recipes to contain repetitive, customizable tasks such as over/undersampling. | 4. Create class balanced training data via under/oversampling with python and in the Lab<br>5. Use python to generate model training configurations and store as json in managed folder<br>6. Create config-driven models that includes H2O random forest as custom python and other algorithms<br>7. Store model performance metrics in an output dataset |
| **Data Details** | Snowflake is housing all of the input data. Currently relying on python-generated SQL queries to Snowflake tables for data prep tasks. The input data is also of dubious quality and a series of python transformations are used to control and modify the input as necessary. | Snowflake data cleaning and joining will be converted to visual recipes and run in-database whenever possible. Will implement a series of metrics and checks to ensure data quality on the input datasets. | 8. Lift and shift existing Streamlit webapp into Dataiku via code studios<br>9. Create email reporter scheme tailored to different audiences based on scenario outcomes<br>10. Output a scored dataset for API access<br>11. Run model retraining and back testing on a to-be-determined cadence |

**Near Market**

**Snowflake**

**CONSUMER_
DEMOGRAPHICS_VW**
_yyyy_mm_dd_
Manual Snapshot
More Current For Eval ex
2022_08_20
(xcrre_snapshots)

**Vehicle Data**

**VIN Pattern Data**

**Anonymous Consumer
Activity**
(CADS_ANON_CONSUMER_
ACT)

**CONSUMER_
DEMOGRAPHICS_VW**
_yyyy_mm_dd_
Manual Snapshot
Historical For Training ex
2022_06_28
(xcrre_snapshots)

**Snowflake Data Prep**

Enrich Vin with Make,
Model, Year, Trim

Enrich, Filter
on Vinsolution
Sales,
Data Prep

Eval
Data Prep
& Merge
Demo, Act, VIN

Data Prep
Agg to VIN

Train
Data Prep
& Merge
Demo, Act, VIN

**Python Eval Prep**

Dynamic
SQL Data
Prep

**Python SubModel Eval**

Per Model
Dynamic
SQL

Per Model
Frame Data
Prep

**H2O.ai**
Model Scoring

One Table Per
Submodel
Scored

**Python Main Model
Eval**

Dynamic
SQL merging Eval
table and submodel
results

Frame Data
Prep

**H2O.ai**
Model Scoring

Main Model Output

**Arg and Config Driven**
MONTHS_SINCE_MOST_RECENT_PURCHASE
MONTHS_BTW_SALES_HHPERSK_AVG
LEASE_FLAG

**Arg and Config Driven**
For submodel results to use and main model to run.
NEAR_MARKET

**Python Model Train**

One Time
Data Prep

Per Model
Dynamic
SQL

Per Model
Frame Data
Prep

**H2O.ai**
RandomForest
Train

**Joblib Config and Model Artifact
for each model/submodel**

**Config Driven**
NEAR_MARKET -- response_type = classification
MONTHS_SINCE_MOST_RECENT_PURCHASE -- response_type = regression
MONTHS_BTW_SALES_HHPERSK_AVG -- response_type = regression
LEASE_FLAG -- response_type = classification

# SOLUTION DIAGRAM

Design Node

Ingest tables

**3**

Join tables

Prepare tables

**4** Under/Oversampling strategy via python

Balanced Dataset

**5** Training config generation in python

Config sub-models

Impute using sub-models

**6** Perform modeling

**7** Produce final model and metrics dataset

**Pre-work (actions not shown):**

**1** Confirm the functionality of Dataiku Design and Automation nodes and establish all required data connections

**2** Create python code environment with all required package dependencies

**3** Join and clean the various Snowflake tables as designed on the Lucidchart diagram

**4** Create class balanced training data via under/oversampling with python and in the Lab

**5** Use python to generate model training configurations and store as json in managed folder
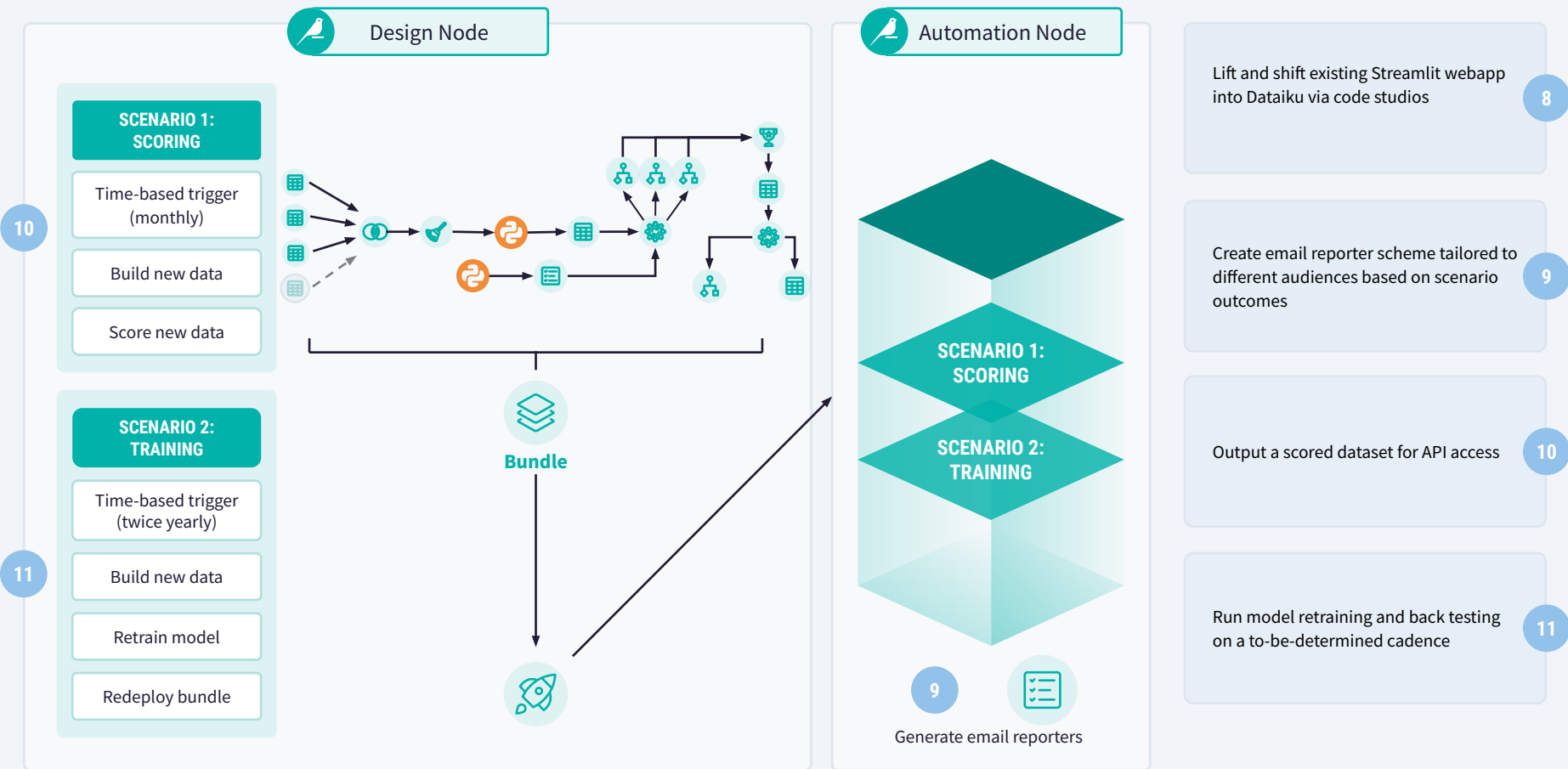
**6** Create config-driven models that includes H2O random forest as custom python and other algorithms

**7** Store model performance metrics in an output dataset

**Design Node**

**Automation Node**

**SCENARIO 1: SCORING**
- Time-based trigger (monthly)
- Build new data
- Score new data

**SCENARIO 2: TRAINING**
- Time-based trigger (twice yearly)
- Build new data
- Retrain model
- Redeploy bundle

**Bundle**

**SCENARIO 1: SCORING**

**SCENARIO 2: TRAINING**

9 — Generate email reporters

8 — Lift and shift existing Streamlit webapp into Dataiku via code studios

9 — Create email reporter scheme tailored to different audiences based on scenario outcomes

10 — Output a scored dataset for API access

11 — Run model retraining and back testing on a to-be-determined cadence

# TASKS & DELIVERABLES

| | **Phase 1**<br>**Project Setup & Discovery** | **Phase 2**<br>**Data Preparation** | **Phase 3**<br>**ML Modelling & Analysis** | **Phase 4**<br>**Automation & Reporting** |
|---|---|---|---|---|
| **Activities** | • Confirm the functionality of Dataiku Design and Automation nodes and establish all required data connections<br>• Create python code environment with all required package dependencies | • Join and clean the various Snowflake tables as designed on the Lucidchart diagram<br>• Create class balanced training data via under/oversampling with python and in the Lab<br>• Use python to generate model training configurations and store as json in managed folder | • Create config-driven models that includes H2O random forest as custom python and other algorithms<br>• Store model performance metrics in an output dataset | • Lift and shift existing Streamlit webapp into Dataiku via code studios<br>• Create email reporter scheme tailored to different audiences based on scenario outcomes<br>• Output a scored dataset for API access<br>• Run model retraining and back testing on a to-be-determined cadence |
| **Milestones** | • Dataiku project with input data pipeline | • Functional data prep pipeline<br>• Prepared dataset for model training | • First iteration of baseline models | • Dataiku project automated via a scenario to send out a report |
| **Assigned to** | Cox Auto | Cox Auto | Cox Auto | Cox Auto |
| **Level of effort** | **0-1 weeks**<br>(0-3 hours)* | **2-3 weeks**<br>(6-9 hours)* | **2-3 weeks**<br>(6-9 hours)* | **1-2 weeks**<br>(3-6 hours)* |

*Hours provided are projected **meeting hours** with your Dataiku data scientist. The customer team is expected to do work outside of coaching sessions. Additionally, the Dataiku data scientist may be required to spend additional hours outside of coaching sessions to prep or respond to follow-up questions.

| PHASE | December | | January | | | | February |
|---|---|---|---|---|---|---|---|
| | Dec 19-23 | Dec 26-30 | Jan 2-6 | Jan 9-13 | Jan 16-20 | Jan 23-27 | Jan 30-Feb 3 |
| **Project Setup & Discovery** | | | | | | | |
| **Data Preparation** | | | | | | | |
| **ML Modelling & Analysis** | | | | | | | |
| **Automation & Reporting** | | | | | | | |