Amazon SageMaker Autopilot Candidate Definition Notebook This notebook was automatically generated by the AutoML job Assignment2-2. This notebook allows you to customize the candidate definitions and execute the SageMaker Autopilot workflow. The dataset has 9 columns and the column named Outcome is used as the target column. This is being treated as a BinaryClassification problem. The dataset also has 2 classes. This notebook will build a BinaryClassification model that maximizes the "F1" quality metric of the trained models. The "F1" metric applies for binary classification with a positive and negative class. It mixes between precision and recall, and is recommended in cases where there are more negative examples compared to positive examples. As part of the AutoML job, the input dataset has been randomly split into two pieces, one for training and one for validation. This notebook helps you inspect and modify the data transformation approaches proposed by Amazon SageMaker Autopilot. You can interactively train the data transformation models and use them to transform the data. Finally, you can execute a multiple algorithm hyperparameter optimization (multi-algo HPO) job that helps you find the best model for your dataset by jointly optimizing the data transformations and machine learning algorithms. Available Knobs Look for sections like this for recommended settings that you can change. **Contents** 1. Sagemaker Setup A. Downloading Generated Candidates B. SageMaker Autopilot Job and Amazon Simple Storage Service (Amazon S3) Configuration-Configuration) 2. Candidate Pipelines A. Generated Candidates **B.** Selected Candidates 3. Executing the Candidate Pipelines A. Run Data Transformation Steps B. Multi Algorithm Hyperparameter Tuning 4. Model Selection and Deployment A. Tuning Job Result Overview B. Model Deployment Sagemaker Setup Before you launch the SageMaker Autopilot jobs, we'll setup the environment for Amazon SageMaker • Check environment & dependencies. • Create a few helper objects/function to organize input/output data and SageMaker sessions. **Minimal Environment Requirements** Jupyter: Tested on JupyterLab 1.0.6, jupyter\_core 4.5.0 and IPython 6.4.0 Kernel: conda\_python3 • Dependencies required ■ sagemaker-python-sdk>=2.40.0 • Use !pip install sagemaker==2.40.0 to download this dependency. o Kernel may need to be restarted after download. Expected Execution Role/permission S3 access to the bucket that stores the notebook. **Downloading Generated Modules** Download the generated data transformation modules and an SageMaker Autopilot helper module used by this notebook. Those artifacts will be downloaded to **Assignment2-2-artifacts** folder. In [ ]: !mkdir -p Assignment2-2-artifacts !aws s3 sync s3://mohan-cogintern-bucket/Assignment2.1/OUTPUT/Assignment2-2/sagemaker-automl-candidates/Assignm !aws s3 sync s3://mohan-cogintern-bucket/Assignment2.1/OUTPUT/Assignment2-2/sagemaker-automl-candidates/Assignm import sys sys.path.append("Assignment2-2-artifacts") SageMaker Autopilot Job and Amazon Simple Storage Service (Amazon S3) Configuration The following configuration has been derived from the SageMaker Autopilot job. These items configure where this notebook will look for generated candidates, and where input and output data is stored on Amazon S3. In [ ]: from sagemaker\_automl import uid, AutoMLLocalRunConfig # Where the preprocessed data from the existing AutoML job is stored BASE\_AUTOML\_JOB\_NAME = 'Assignment2-2' BASE\_AUTOML\_JOB\_CONFIG = { 'automl\_job\_name': BASE AUTOML JOB NAME, 'automl\_output\_s3\_base\_path': 's3://mohan-cogintern-bucket/Assignment2.1/OUTPUT/Assignment2-2', 'data\_transformer\_image\_repo\_version': '2.5-1-cpu-py3', 'algo\_image\_repo\_versions': {'xgboost': '1.3-1-cpu-py3', 'linear-learner': 'training-cpu', 'mlp': 'training 'algo\_inference\_image\_repo\_versions': {'xgboost': '1.3-1-cpu-py3', 'linear-learner': 'inference-cpu', 'mlp # Path conventions of the output data storage path from the local AutoML job run of this notebook LOCAL AUTOML JOB NAME = 'Assignment-notebook-run-{}'.format(uid()) LOCAL\_AUTOML\_JOB\_CONFIG = { 'local\_automl\_job\_name': LOCAL\_AUTOML\_JOB\_NAME, 'local\_automl\_job\_output\_s3\_base\_path': 's3://mohan-cogintern-bucket/Assignment2.1/OUTPUT/Assignment2-2/{} 'data\_processing\_model\_dir': 'data-processor-models', 'data\_processing\_transformed\_output\_dir': 'transformed-data', 'multi\_algo\_tuning\_output\_dir': 'multi-algo-tuning' AUTOML LOCAL RUN CONFIG = AutoMLLocalRunConfig( role='arn:aws:iam::313830654669:role/service-role/AmazonSageMaker-ExecutionRole-20220612T142446', base\_automl\_job\_config=BASE\_AUTOML\_JOB\_CONFIG, local\_automl\_job\_config=LOCAL AUTOML JOB CONFIG, security config={'EnableInterContainerTrafficEncryption': False, 'VpcConfig': {}}) AUTOML\_LOCAL\_RUN\_CONFIG.display() **Candidate Pipelines** The AutoMLLocalRunner keeps track of selected candidates and automates many of the steps needed to execute feature engineering and tuning steps. In [ ]: from sagemaker automl import AutoMLInteractiveRunner, AutoMLLocalCandidate automl interactive runner = AutoMLInteractiveRunner(AUTOML LOCAL RUN CONFIG) **Generated Candidates** The SageMaker Autopilot Job has analyzed the dataset and has generated 6 machine learning pipeline(s) that use 3 algorithm(s). Each pipeline contains a set of feature transformers and an algorithm. Available Knobs 1. The resource configuration: instance type & count 1. Select candidate pipeline definitions by cells 1. The linked data transformation script can be reviewed and updated. Please refer to the [README.md](./Assignment2-2artifacts/generated\_module/README.md) for detailed customization instructions. dpp0-xgboost: This data transformation strategy first transforms 'numeric' features using RobustImputer (converts missing values to nan). It merges all the generated features and applies RobustStandardScaler. The transformed data will be used to tune a xgboost model. Here is the definition: In [ ]: automl interactive runner.select candidate({ "data transformer": { "name": "dpp0", "training resource config": { "instance type": "ml.m5.12xlarge", "instance count": 1, "volume size in gb": "transform resource config": { "instance type": "ml.m5.4xlarge", "instance count": 1, "transforms label": True, "transformed data format": "text/csv", "sparse encoding": False "algorithm": { "name": "xgboost", "training resource config": { "instance type": "ml.m5.12xlarge", "instance count": 1, }) dpp1-xgboost: This data transformation strategy first transforms 'numeric' features using RobustImputer. It merges all the generated features and applies RobustPCA followed by RobustStandardScaler. The transformed data will be used to tune a xgboost model. Here is the definition: In [ ]: automl interactive runner.select candidate({ "data transformer": { "name": "dpp1", "training resource config": { "instance type": "ml.m5.12xlarge", "instance count": 1, "volume size in gb": 50 "transform resource config": { "instance type": "ml.m5.4xlarge", "instance count": 1, "transforms label": True, "transformed data format": "text/csv", "sparse encoding": False "algorithm": { "name": "xgboost", "training resource config": { "instance type": "ml.m5.12xlarge", "instance count": 1, } }) dpp2-linear-learner: This data transformation strategy first transforms 'numeric' features using combined RobustImputer and RobustMissingIndicator followed by QuantileExtremeValuesTransformer, 'categorical' features using ThresholdOneHotEncoder. It merges all the generated features and applies RobustPCA followed by RobustStandardScaler. The transformed data will be used to tune a linearlearner model. Here is the definition: In [ ]: automl interactive runner.select candidate({ "data transformer": { "name": "dpp2", "training resource config": { "instance type": "ml.m5.12xlarge", "instance count": 1, "volume size in gb": 50 "transform\_resource\_config": { "instance type": "ml.m5.4xlarge", "instance count": 1, "transforms label": True, "transformed data format": "application/x-recordio-protobuf", "sparse encoding": False "algorithm": { "name": "linear-learner", "training resource config": { "instance type": "ml.m5.12xlarge", "instance count": 1, }, } }) dpp3-xgboost: This data transformation strategy first transforms 'numeric' features using RobustImputer (converts missing values to nan), 'categorical' features using ThresholdOneHotEncoder. It merges all the generated features and applies RobustStandardScaler. The transformed data will be used to tune a xgboost model. Here is the definition: In [ ]: automl interactive runner.select candidate({ "data transformer": { "name": "dpp3", "training resource config": { "instance type": "ml.m5.12xlarge", "instance count": 1, "volume size in gb": "transform resource config": { "instance type": "ml.m5.4xlarge", "instance count": 1, "transforms label": True, "transformed data format": "application/x-recordio-protobuf", "sparse encoding": True "algorithm": { "name": "xgboost", "training resource config": { "instance type": "ml.m5.12xlarge", "instance count": 1, }, }) dpp4-xgboost: This data transformation strategy first transforms 'numeric' features using RobustImputer (converts missing values to nan), 'categorical' features using ThresholdOneHotEncoder. It merges all the generated features and applies RobustStandardScaler. The transformed data will be used to tune a xqboost model. Here is the definition: In [ ]: automl\_interactive\_runner.select\_candidate({ "data transformer": { "name": "dpp4", "training resource\_config": { "instance\_type": "ml.m5.12xlarge", "instance count": 1, "volume\_size\_in\_gb": "transform resource\_config": { "instance\_type": "ml.m5.4xlarge", "instance count": 1, "transforms label": True, "transformed data format": "application/x-recordio-protobuf", "sparse encoding": True "algorithm": { "name": "xgboost", "training\_resource\_config": { "instance\_type": "ml.m5.12xlarge", instance\_count": }) dpp5-mlp: This data transformation strategy transforms 'numeric' features using RobustImputer followed by RobustStandardScaler. The transformed data will be used to tune a mlp model. Here is the definition: In [ ]: automl interactive runner.select candidate({ "data transformer": { "name": "dpp5", "training resource config": { "instance type": "ml.m5.12xlarge", "instance count": 1, "volume size in gb": 50 "transform resource config": { "instance type": "ml.m5.4xlarge", "instance count": 1, "transforms label": True, "transformed data format": "text/csv", "sparse encoding": False "algorithm": { "name": "mlp", "training resource config": { "instance type": "ml.m5.12xlarge", "instance count": 1, "candidate specific static hyperparameters": { "num categorical features": '0', }) **Selected Candidates** You have selected the following candidates (please run the cell below and click on the feature transformer links for details): In [ ]: automl interactive runner.display candidates() The feature engineering pipeline consists of two SageMaker jobs: 1. Generated trainable data transformer Python modules like dpp0.py, which has been downloaded to the local file system 2. A **training** job to train the data transformers 3. A batch transform job to apply the trained transformation to the dataset to generate the algorithm compatible data The transformers and its training pipeline are built using open sourced sagemaker-scikit-learn-container and sagemaker-scikit-learnextension. **Executing the Candidate Pipelines** Each candidate pipeline consists of two steps, feature transformation and algorithm training. For efficiency first execute the feature transformation step which will generate a featurized dataset on S3 for each pipeline. After each featurized dataset is prepared, execute a multi-algorithm tuning job that will run tuning jobs in parallel for each pipeline. This tuning job will execute training jobs to find the best set of hyper-parameters for each pipeline, as well as finding the overall best performing pipeline. Run Data Transformation Steps Now you are ready to start execution all data transformation steps. The cell below may take some time to finish, feel free to go grab a cup of coffee. To expedite the process you can set the number of parallel\_jobs to be up to 10. Please check the account limits to increase the limits before increasing the number of jobs to run in parallel. automl interactive runner.fit data transformers(parallel jobs=7) Multi Algorithm Hyperparameter Tuning Now that the algorithm compatible transformed datasets are ready, you can start the multi-algorithm model tuning job to find the best predictive model. The following algorithm training job configuration for each algorithm is auto-generated by the AutoML Job as part of the recommendation. Available Knobs 1. Hyperparameter ranges 2. Objective metrics 3. Recommended static algorithm hyperparameters. Please refers to [Xgboost tuning](https://docs.aws.amazon.com/sagemaker/latest/dg/xgboost-tuning.html) and [Linear learner tuning] (https://docs.aws.amazon.com/sagemaker/latest/dg/linear-learner-tuning.html) for detailed explanations of the parameters. The AutoML recommendation job has recommended the following hyperparameters, objectives and accuracy metrics for the algorithm and problem type: In [ ]: ALGORITHM OBJECTIVE METRICS = { 'xgboost': 'validation:f1 binary', 'linear-learner': 'validation:binary f beta', 'mlp': 'validation:binary f beta', STATIC HYPERPARAMETERS = { 'xgboost': { 'objective': 'binary:logistic', 'eval metric': 'accuracy, f1 binary, auc, balanced accuracy, precision, recall, logloss', ' kfold': 5, ' num cv round': 3, 'linear-learner': { 'predictor type': 'binary classifier', 'ml application': 'linear learner', 'loss function': 'SoftmaxCrossEntropyLoss', 'reporting metrics': 'binary classification accuracy, binary f beta, roc auc score, balanced accuracy, pred 'eval metric': 'binary f beta', 'kfold': 5, 'num cv rounds': 3, 'prediction storage mode': 'store cv avg predictions', }, 'problem type': 'binary classification', 'positive example weight mult': 'auto', 'ml application': 'mlp', 'use batchnorm': 'true', 'activation': 'relu', 'warmup epochs': 10, 'reporting metrics': 'accuracy, binary f 1, roc auc, balanced accuracy, precision, recall, logloss', 'eval metric': 'binary f 1', 'kfold': 5, 'num cv rounds': 3, 'prediction storage mode': 'store cv avg predictions', }, The following tunable hyperparameters search ranges are recommended for the Multi-Algo tuning job: In [ ]: from sagemaker.parameter import CategoricalParameter, ContinuousParameter, IntegerParameter ALGORITHM TUNABLE HYPERPARAMETER RANGES = { 'xgboost': { 'num\_round': IntegerParameter(64, 1024, scaling\_type='Logarithmic'), 'max\_depth': IntegerParameter(2, 8, scaling\_type='Logarithmic'), 'eta': ContinuousParameter(1e-3, 1.0, scaling\_type='Logarithmic'), 'gamma': ContinuousParameter(1e-6, 64.0, scaling\_type='Logarithmic'), 'min\_child\_weight': ContinuousParameter(1e-6, 32.0, scaling\_type='Logarithmic'), 'subsample': ContinuousParameter(0.5, 1.0, scaling\_type='Linear'), 'colsample\_bytree': ContinuousParameter(0.3, 1.0, scaling\_type='Linear'), lambda': ContinuousParameter(1e-6, 2.0, scaling\_type='Logarithmic' 'alpha': ContinuousParameter(1e-6, 2.0, scaling\_type='Logarithmic'), 'linear-learner': { 'mini batch size': IntegerParameter(128, 512, scaling\_type='Linear'), 'wd': ContinuousParameter(1e-12, 1e-2, scaling\_type='Logarithmic'), 'learning rate': ContinuousParameter(1e-6, 1e-2, scaling type='Logarithmic'), 'mini\_batch\_size': IntegerParameter(128, 512, scaling\_type='Linear'), 'learning rate': ContinuousParameter(1e-6, 1e-2, scaling type='Logarithmic'), 'weight\_decay': ContinuousParameter(1e-12, 1e-2, scaling\_type='Logarithmic'), 'dropout\_prob': ContinuousParameter(0.25, 0.5, scaling\_type='Linear'), 'embedding\_size\_factor': ContinuousParameter(0.65, 0.95, scaling\_type='Linear'), 'network\_type': CategoricalParameter(['feedforward', 'widedeep']), 'layers': CategoricalParameter(['256', '50, 25', '100, 50', '200, 100', '256, 128', '300, 150', '200, }, **Prepare Multi-Algorithm Tuner Input** To use the multi-algorithm HPO tuner, prepare some inputs and parameters. Prepare a dictionary whose key is the name of the trained pipeline candidates and the values are respectively: 1. Estimators for the recommended algorithm 2. Hyperparameters search ranges 3. Objective metrics In [ ]: multi\_algo\_tuning\_parameters = automl\_interactive\_runner.prepare\_multi\_algo\_parameters( objective metrics=ALGORITHM OBJECTIVE METRICS, static\_hyperparameters=STATIC\_HYPERPARAMETERS, hyperparameters search ranges=ALGORITHM TUNABLE HYPERPARAMETER RANGES) Below you prepare the inputs data to the multi-algo tuner: In [ ]: multi\_algo\_tuning\_inputs = automl\_interactive\_runner.prepare\_multi\_algo\_inputs() **Create Multi-Algorithm Tuner** With the recommended Hyperparameter ranges and the transformed dataset, create a multi-algorithm model tuning job that coordinates hyper parameter optimizations across the different possible algorithms and feature processing strategies. Available Knobs 1. Tuner strategy: [Bayesian](https://en.wikipedia.org/wiki/Hyperparameter\_optimization#Bayesian\_optimization), [Random Search](https://en.wikipedia.org/wiki/Hyperparameter\_optimization#Random\_search) 2. Objective type:  $M \in imize$ ,  $Ma\xi mize$ , see [optimization](https://en.wikipedia.org/wiki/Mathematical\_optimization) 3. Max Job size: the max number of training jobs HPO would be launching to run experiments. Note the default value is \*\*250\*\* which is the default of the managed flow. 4. Parallelism. Number of jobs that will be executed in parallel. Higher value will expedite the tuning process. Please check the account limits to increase the limits before increasing the number of jobs to run in parallel 5. Please use a different tuning job name if you rerun this cell after applied customizations. In [ ]: from sagemaker.tuner import HyperparameterTuner base tuning job name = "{}-tuning".format(AUTOML LOCAL RUN CONFIG.local automl job name) tuner = HyperparameterTuner.create( base\_tuning\_job\_name=base\_tuning\_job\_name, strategy='Bayesian', objective type='Maximize', max parallel\_jobs=7,  $max_jobs=250$ , \*\*multi\_algo\_tuning\_parameters, Run Multi-Algorithm Tuning Now you are ready to start running the **Multi-Algo Tuning** job. After the job is finished, store the tuning job name which you use to select models in the next section. The tuning process will take some time, please track the progress in the Amazon SageMaker Hyperparameter tuning jobs console. In [ ]: from IPython.display import display, Markdown # Run tuning tuner.fit(inputs=multi algo tuning inputs, include cls metadata=None) tuning job name = tuner.latest tuning job.name display( Markdown(f"Tuning Job {tuning job name} started, please track the progress from [here](https://{AUTOML LOCA # Wait for tuning job to finish tuner.wait() **Model Selection and Deployment** This section guides you through the model selection process. Afterward, you construct an inference pipeline on Amazon SageMaker to host the best candidate. Because you executed the feature transformation and algorithm training in two separate steps, you now need to manually link each trained model with the feature transformer that it is associated with. When running a regular Amazon SageMaker Autopilot job, this will automatically be done for you. **Tuning Job Result Overview** The performance of each candidate pipeline can be viewed as a Pandas dataframe. For more interactive usage please refers to model tuning monitor. In [ ]: from pprint import pprint from sagemaker.analytics import HyperparameterTuningJobAnalytics SAGEMAKER SESSION = AUTOML LOCAL RUN CONFIG.sagemaker session SAGEMAKER ROLE = AUTOML LOCAL RUN CONFIG.role tuner analytics = HyperparameterTuningJobAnalytics( tuner.latest\_tuning\_job.name, sagemaker\_session=SAGEMAKER\_SESSION) df\_tuning\_job\_analytics = tuner\_analytics.dataframe() # Sort the tuning job analytics by the final metrics value df tuning job analytics.sort values ( by=['FinalObjectiveValue'], inplace=True, ascending=False if tuner.objective type == "Maximize" else True) # Show detailed analytics for the top 20 models df\_tuning\_job\_analytics.head(20) The best training job can be selected as below:  $\P$  **Tips:** You could select alternative job by using the value from  $Tra \in gJobName$  column above and assign to  $best_tra\in g_job$  below attached\_tuner = HyperparameterTuner.attach(tuner.latest\_tuning\_job.name, sagemaker\_session=SAGEMAKER\_SESSION) best\_training\_job = attached\_tuner.best\_training\_job() print("Best Multi Algorithm HPO training job name is {}".format(best\_training\_job)) **Linking Best Training Job with Feature Pipelines** Finally, deploy the best training job to Amazon SageMaker along with its companion feature engineering models. At the end of the section, you get an endpoint that's ready to serve online inference or start batch transform jobs! Deploy a PipelineModel that has multiple containers of the following: 1. Data Transformation Container: a container built from the model we selected and trained during the data transformer sections 2. Algorithm Container: a container built from the trained model we selected above from the best HPO training job. 3. Inverse Label Transformer Container: a container that converts numerical intermediate prediction value back to non-numerical label Get both best data transformation model and algorithm model from best training job and create an pipeline model: from sagemaker.estimator import Estimator from sagemaker import PipelineModel from sagemaker\_automl import select\_inference\_output # Get a data transformation model from chosen candidate best candidate = automl interactive runner.choose candidate(df tuning\_job\_analytics, best\_training\_job) best\_data\_transformer\_model = best\_candidate.get\_data\_transformer\_model(role=SAGEMAKER\_ROLE, sagemaker\_session= # Our first data transformation container will always return recordio-protobuf format best\_data\_transformer\_model.env["SAGEMAKER\_DEFAULT\_INVOCATIONS\_ACCEPT"] = 'application/x-recordio-protobuf' # Add environment variable for sparse encoding if best candidate.data transformer\_step.sparse\_encoding: best data transformer model.env["AUTOML SPARSE ENCODE RECORDIO PROTOBUF"] = '1' # Get a algo model from chosen training job of the candidate algo estimator = Estimator.attach(best training job) best\_algo\_model = algo\_estimator.create\_model(\*\*best\_candidate.algo\_step.get\_inference\_container\_config()) # Final pipeline model is composed of data transformation models and algo model and an # inverse label transform model if we need to transform the intermediates back to non-numerical value model\_containers = [best\_data\_transformer model, best algo model] if best candidate.transforms\_label: model\_containers.append(best\_candidate.get\_data\_transformer\_model( transform mode="inverse-label-transform", role=SAGEMAKER ROLE, sagemaker session=SAGEMAKER SESSION)) # This model can emit response ['predicted\_label', 'probability', 'labels', 'probabilities']. To enable the model # of the response content, pass the keys to `output key` keyword argument in the select inference output method model containers = select inference output ("BinaryClassification", model containers, output keys=['predicted la pipeline model = PipelineModel( name="AutoML-{}".format(AUTOML LOCAL RUN CONFIG.local automl job name), role=SAGEMAKER ROLE, models=model containers, vpc\_config=AUTOML\_LOCAL\_RUN\_CONFIG.vpc\_config) **Deploying Best Pipeline** Available Knobs 1. You can customize the initial instance count and instance type used to deploy this model. 2. Endpoint name can be changed to avoid conflict with existing endpoints. Finally, deploy the model to SageMaker to make it functional. In [ ]: pipeline model.deploy(initial instance count=1, instance type='ml.m5.2xlarge', endpoint name=pipeline model.name, wait=True) Congratulations! Now you could visit the sagemaker endpoint console page to find the deployed endpoint (it'll take a few minutes to be in service). To rerun this notebook, delete or change the name of your endpoint! If you rerun this notebook, you'll run into an error on the last step because the endpoint already exists. You can either delete the endpoint from the endpoint console page or you can change the endpoint\_name in the previous code block.