

Amazon SageMaker Autopilot Candidate Definition Notebook

This notebook was automatically generated by the AutoML job `automl-dm-1624877542`. This notebook allows you to customize the candidate definitions and execute the SageMaker Autopilot workflow.

The dataset has 2 columns and the column named `sentiment` is used as the target column. This is being treated as a **MulticlassClassification** problem. The dataset also has 3 classes. This notebook will build a [MulticlassClassification](#) model that maximizes the "ACCURACY" quality metric of the trained models. The "ACCURACY" metric provides the percentage of times the model predicted the correct class.

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.
 - 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 **automl-dm-1624877542-artifacts** folder.

In []:

```
!mkdir -p automl-dm-1624877542-artifacts
!aws s3 sync s3://sagemaker-us-east-1-985209124379/autopilot/automl-dm-1624877542/sagemaker-automl-candidates/automl-dm-1624877542-pr-1-c888d863c77c4afba59434523d82fa5deaedc/generated_module automl-dm-1624877542-artifacts/generated_module --only-show-errors
!aws s3 sync s3://sagemaker-us-east-1-985209124379/autopilot/automl-dm-1624877542/sagemaker-automl-candidates/automl-dm-1624877542-pr-1-c888d863c77c4afba59434523d82fa5deaedc/notebooks/sagemaker_automl automl-dm-1624877542-artifacts/sagemaker_automl --only-show-errors

import sys
sys.path.append("automl-dm-1624877542-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 = 'automl-dm-1624877542'
BASE_AUTOML_JOB_CONFIG = {
    'automl_job_name': BASE_AUTOML_JOB_NAME,
    'automl_output_s3_base_path': 's3://sagemaker-us-east-1-985209124379/autopilot/automl-dm-1624877542',
    'data_transformer_image_repo_version': '0.2-1-cpu-py3',
    'algo_image_repo_versions': {'xgboost': '1.2-2-cpu-py3'},
    'algo_inference_image_repo_versions': {'xgboost': '1.2-2-cpu-py3'}
}

# Path conventions of the output data storage path from the local AutoML job run of this notebook
LOCAL_AUTOML_JOB_NAME = 'automl-dm--notebook-run-{}'.format(uid())
LOCAL_AUTOML_JOB_CONFIG = {
    'local_automl_job_name': LOCAL_AUTOML_JOB_NAME,
    'local_automl_job_output_s3_base_path': 's3://sagemaker-us-east-1-985209124379/autopilot/automl-dm-1624877542/{}'.format(LOCAL_AUTOML_JOB_NAME),
    '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::985209124379:role/c21581a4068141803767t1w9852-SageMakerExecutionRole-13EE24MKUSOXF',
    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 **3** machine learning pipeline(s) that use **1** 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 [\[README.md\]\(./automl-dm-1624877542-artifacts/generated_module/README.md\)](#) for detailed customization instructions.

dpp0-xgboost: This data transformation strategy first transforms 'text' features using [MultiColumnTfidfVectorizer](#). 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.m5.4xlarge",
            "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": True
    },
    "algorithm": {
        "name": "xgboost",
        "training_resource_config": {
            "instance_type": "ml.m5.4xlarge",
            "instance_count": 1,
        },
    },
})
```

dpp1-xgboost: This data transformation strategy first transforms 'text' features using [MultiColumnTfidfVectorizer](#). 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.m5.4xlarge",
            "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.m5.4xlarge",
        "instance_count": 1,
    },
}
})

```

dpp2-xgboost: This data transformation strategy first transforms 'text' features using [MultiColumnTfidfVectorizer](#). 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": "dpp2",
        "training_resource_config": {
            "instance_type": "ml.m5.4xlarge",
            "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": True
    },
    "algorithm": {
        "name": "xgboost",
        "training_resource_config": {
            "instance_type": "ml.m5.4xlarge",
            "instance_count": 1,
        },
    },
})

```

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-learn-extension](#).

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.

In []:

```
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:accuracy',
}

STATIC_HYPERPARAMETERS = {
    'xgboost': {
        'objective': 'multi:softprob',
        'num_class': 3,
        '_kfold': 5,
    },
}
```

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(2, 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'),
    },
}

```

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: `Minimize`, `Maximize`, 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 re-run 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-Algorithm 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_LOCAL_RUN_CONFIG.region}.console.aws.amazon.com/sagemaker/home?region={AUTOML_LOCAL_RUN_CONFIG.region}#/hyper-tuning-jobs/{tuning_job_name})")
)

# 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 refer 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:

Tip: You could select alternative job by using the value from `TrainingJobName` column above and assign to `best_training_job` below

In []:

```
attached_tuner = HyperparameterTuner.attach(tuner.latest_tuning_job.name, sagemaker_session=SAGEMAKER_SESSION)
```

```

on=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 value.

Get both best data transformation model and algorithm model from best training job and create an pipeline model:

In []:

```

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=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 to emit one or more
# of the response content, pass the keys to `output_key` keyword argument in the select_inference_output method.

model_containers = select_inference_output("MulticlassClassification", model_containers, output_keys=['predicted_label'])

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.