

Open in app ↗



Search

Write



Databricks End-To-End Machine Learning - Create An Ingest-To-Serving MLOps Pipeline

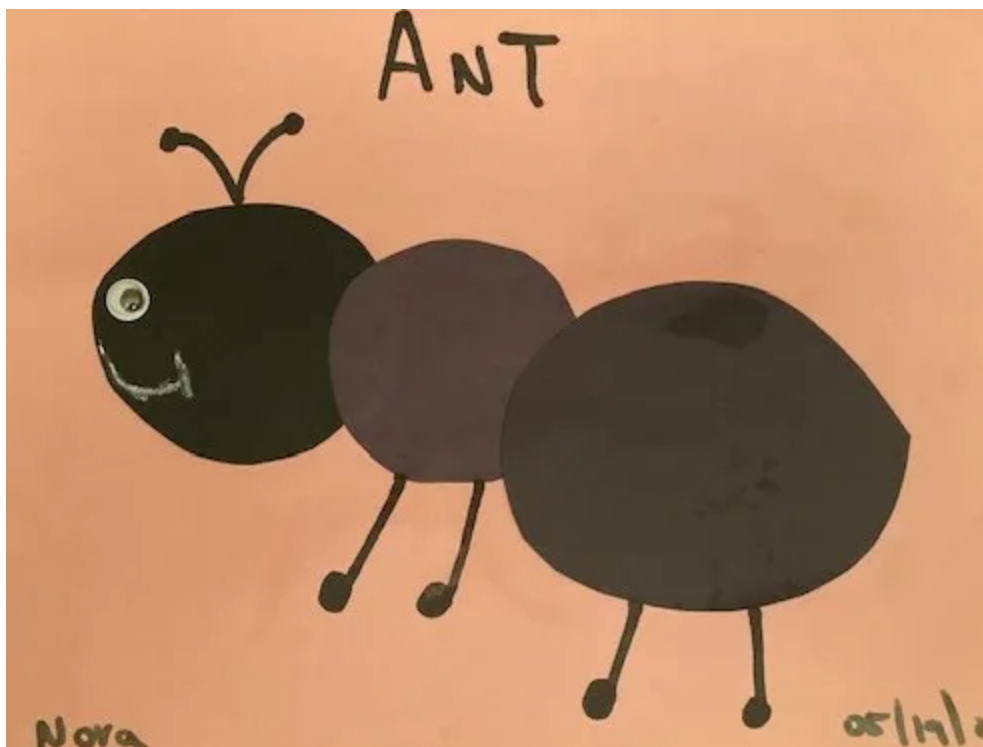


Thomas Jaensch

9 min read · Oct 12, 2022



13



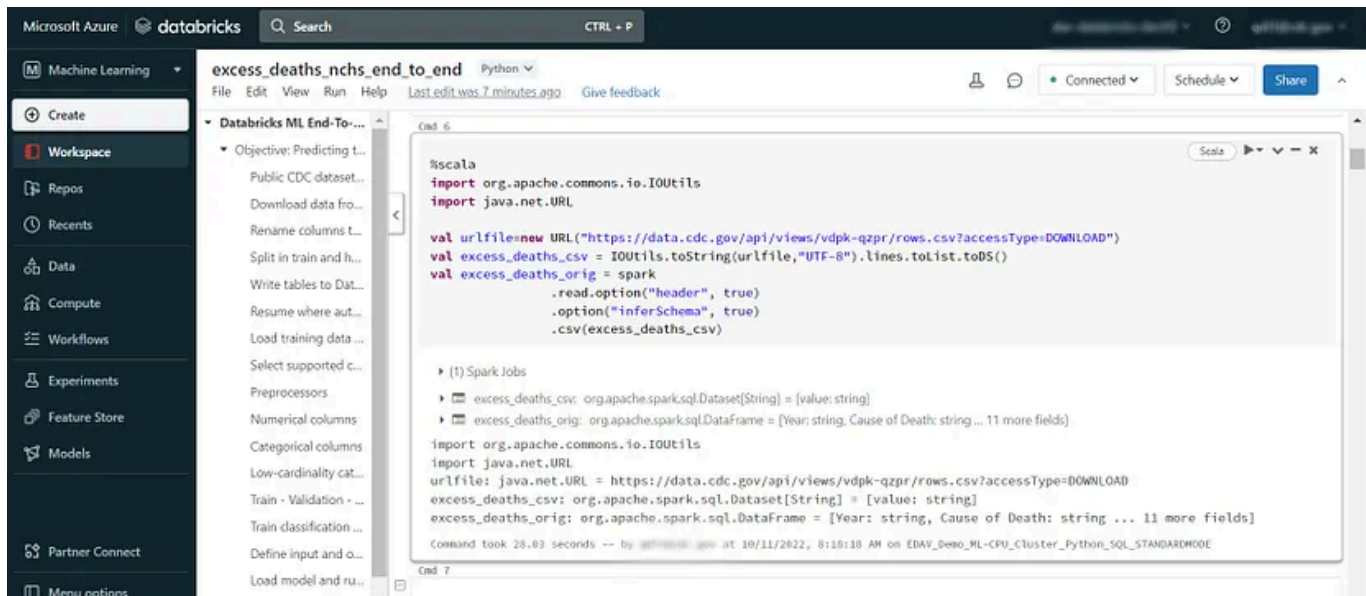
Art my daughter Nora produced at daycare

Hands on workflow how to build a Machine Learning XGBoost multi-class classification model all the way from real world data ingest (public CDC dataset) to model deployment and making real-time predictions. And once done, automate the whole enchilada MLOps-style. *

The complete code to recreate the following demo yourself can be found in this GitHub repo: https://github.com/tjaensch/excess_deaths_nchs *

Part 1: Ingest Data Into Databricks

Create a Python notebook in your Databricks workspace and attach it to a suitable Databricks ML cluster. I'm only using Scala in the following steps because it appeared to be the easiest way to get the data from the public CDC URL into a Spark dataframe without having to download files locally. This can be done in a Databricks Python notebook by using the Scala magic command `%scala` at the top of each cell with Scala code. All the actual Machine Learning code later on will be written in Python.



Download data from URL into a Spark dataframe:

```
%scala
```

```
import org.apache.commons.io.IOUtils
import java.net.URL
```

```
val urlfile=new URL("https://data.cdc.gov/api/views/vdpk-qzpr/rows.csv?accessType=DOWNLOAD")
val excess_deaths_csv = IOUtils.toString(urlfile,"UTF-8").lines.toList.toDS()
val excess_deaths_orig = spark
    .read.option("header", true)
    .option("inferSchema", true)
    .csv(excess_deaths_csv)
```

Rename columns to be able to save as Databricks table:

```
%scala
```

```
val excess_deaths = excess_deaths_orig.withColumnRenamed("Year",
"year")
    .withColumnRenamed("Cause of Death",
"cause_of_death")
    .withColumnRenamed("State", "state")
    .withColumnRenamed("State FIPS Code",
```

```
"state_fips_code")
    .withColumnRenamed("HHS Region", "hhs_region")
    .withColumnRenamed("Age Range", "age_range")
    .withColumnRenamed("Benchmark", "benchmark")
    .withColumnRenamed("Locality", "locality")
    .withColumnRenamed("Observed Deaths",
"observed_deaths")
    .withColumnRenamed("Population", "population")
    .withColumnRenamed("Expected Deaths",
"expected_deaths")
    .withColumnRenamed("Potentially Excess Deaths",
"potentially_excess_deaths")
    .withColumnRenamed("Percent Potentially Excess
Deaths", "percent_potentially_excess_deaths")
```

Split in training and holdout set:

```
%scala

val Array(training, holdout) = excess_deaths.randomSplit(Array(0.95,
0.05), seed = 12345)
```

I am going to use [Databricks AutoML](#) in the next step which does its own training/evaluation/test split so the above is mainly to have some data for testing (holdout) the best AutoML model after it has been created on data that the AutoML process has not seen at all yet.

At the time of writing there were 199317 examples in the training, and 10563 in the holdout dataset.

Split in train and holdout set

Cmd 12

```
%scala
val Array(training, holdout) = excess_deaths.randomSplit(Array(0.95, 0.05), seed = 12345)

▶ training: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [year: string, cause_of_death: string ... 11 more fields]
▶ holdout: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [year: string, cause_of_death: string ... 11 more fields]

training: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [year: string, cause_of_death: string ... 11 more fields]
holdout: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [year: string, cause_of_death: string ... 11 more fields]

Command took 0.79 seconds -- by qm... at 10/11/2022, 8:25:54 AM on DBRW_Demo_ML-CPU_Cluster_Python_SQL_STANDARDMODE
```

Cmd 13

```
%scala
println(training.count)
println(holdout.count)
```

▶ (4) Spark Jobs

```
199317
10563
```

Write the Spark dataframes as tables to the Databricks default database to be used in later steps for training and testing:

```
%scala
```

```
excess_deaths.write.mode("overwrite").format("parquet").saveAsTable("default.excess_deaths_nchs_all")
training.write.mode("overwrite").format("parquet").saveAsTable("default.excess_deaths_nchs_training")
test.write.mode("overwrite").format("parquet").saveAsTable("default.excess_deaths_nchs_holdout")
```

Part 2: Create AutoML Classification Experiment

Instead of writing all the necessary ML code from scratch in a notebook to build a classification model, I am going to leverage Databricks' AutoML capabilities. The AutoML Experiment (see below) is easy to configure and will run dozens of trials with different algorithms and in the end produce

one winning experiment that I am going to use to build the end-to-end MLOps pipeline.

The screenshot shows the Databricks AutoML Experiment Configuration page. The left sidebar contains navigation options: Machine Learning, Create, Workspace, Repos, Recents, Data, Compute, Workflows, Experiments (selected), Feature Store, Models, Partner Connect, and Menu options. The main panel is titled 'AutoML Experiment Configuration' and includes the following settings:

- Cluster:** Demo_ML-CPU_Cluster_Python_SQL_STANDAR...
- ML problem type:** Classification
- Dataset:** default.excess_deaths_nchs_training
- Prediction target:** cause_of_death
- Experiment name:** cause_of_death_excess_deaths_nchs_training-2022_10_05

Below these settings is the 'Advanced Configuration (optional)' section with a 'Start AutoML' button. To the right, the 'Schema' table for 'default.excess_deaths_nchs_training' is displayed:

Include	Column n...	Data type	Impute with
<input checked="" type="checkbox"/>	year	string	Auto
<input checked="" type="checkbox"/>	cause_of_de...	string	Auto
<input checked="" type="checkbox"/>	state	string	Auto
<input checked="" type="checkbox"/>	state_fips_co...	string	Auto
<input checked="" type="checkbox"/>	hhs_region	string	Auto
<input checked="" type="checkbox"/>	age_range	string	Auto
<input checked="" type="checkbox"/>	benchmark	string	Auto
<input checked="" type="checkbox"/>	locality	string	Auto
<input checked="" type="checkbox"/>	observed_de...	int	Auto

As the experiment is running, AutoML is discovering data issues that are automatically being addressed.

The screenshot shows the 'Warnings' tab in the Databricks AutoML interface. It displays a list of warnings identified by AutoML, with a link to 'View data exploration notebook'. The warnings are as follows:

Severity	Type	Affected Data	Action
Medium	More than 5% of values are null.	expected_deaths, observed_deaths, percent_potent..._excess_deaths, potentially_excess_deaths	AutoML imputed the null values.
Medium	Ratio of the least frequent label (DC) to the most frequent label (Heart Disease) is 0.096		AutoML will try to balance the target labels before training models
Low	Data exploration notebook ran on a subset of rows because dataset is too big		Modify the data exploration notebook and rerun it to profile the entire dataset.
Low	High correlation columns	hhs_region, cause_of_death, observed_deaths, expected_deaths, percent_potent..._excess_deaths, + 5 more	Correlations found. Refer to data exploration notebook for more details.

Once the training run has finished after one hour which was the default timeout at the time of writing, the Experiments page will display a long list of trials it ran during the experiment with the best run at the top.

The screenshot shows the Databricks AutoML Experiments page. The left sidebar contains navigation links: Machine Learning, Create, Workspace, Repos, Recents, Data, Compute, Workflows, Experiments (selected), Feature Store, and Models. The main panel displays the AutoML progress bar with 'Configure' and 'Train' steps completed, and 'Evaluate' in progress. Below the progress bar, it states 'AutoML Evaluation complete' and 'All runs have completed, and have been added to the table below.' It also highlights the 'Model with best val_f1_score' and provides instructions on how to view details or review the data exploration notebook. A table of runs is shown with columns for Created, Duration, Run Name, User, Source, Version, Models, and Metrics. The table lists three runs: xgboost (13.7min), xgboost (5.6min), and lightgbm (35.1min). The lightgbm run has the highest test_f1_score of 0.919.

Created	Duration	Run Name	User	Source	Version	Models	test_accuracy_sc	test_f1_score	test_log_loss
16 minutes ago	13.7min	xgboost	quintin@ucsf.edu	Notebook...	-	sklearn	0.92	0.92	0.189
17 minutes ago	5.6min	xgboost	quintin@ucsf.edu	Notebook...	-	sklearn	0.919	0.919	0.196
53 minutes ago	35.1min	lightgbm	quintin@ucsf.edu	Notebook...	-	sklearn	0.919	0.919	0.184

Now the resulting best model could be used as is for making predictions on new data, or the holdout data we saved previously, but the point of this post is to demonstrate how the code created by the best AutoML experiment can be reused to create a repeatable process that downloads fresh data each time from the CDC download URL, and runs the same code that worked best this time on this type of data as a scheduled process on new data in case the CDC updates the public dataset.

The neat thing about Databricks ML is that the AutoML experiments generate notebooks (under 'Source' in the image above) for each trial run that can be used to rerun the exact same trial with the same specs, inputs, etc., but that auto-generated code can also be used and enhanced which is what I am going to do. The auto-generated code for the best experiment I got at the time of writing can be found [here](#).

The winning algorithm out of the trials that AutoML produced was, rather unsurprisingly, XGBoost. The AutoML process automatically did a bunch of preprocessing, imputing of missing values, scaling, one-hot-encoding, etc. without me having to do a thing at all, and it also did an internal train-validation-test split on the training data we fed into it above. 60 percent of the training dataset was then used for training the model, and 20 percent each for validation and testing. I am going to change this distribution later in the end-to-end notebook to 90 percent for training and 5 percent each for validation and testing since the more data for training the better the model metrics I achieved.

The hyperparameters of the best XGBoost model produced look as follows:

```
xgbc_classifier = TransformedTargetClassifier(  
    classifier=XGBClassifier(  
        colsample_bytree=0.6925102629069848,  
        learning_rate=0.16153054387601637,  
        max_depth=8,  
        min_child_weight=6,  
        n_estimators=1573,  
        n_jobs=100,  
        subsample=0.7774586806850589,  
        verbosity=0,  
        random_state=538597927,  
    ),  
    transformer=LabelEncoder() # XGBClassifier requires the target  
values to be integers between 0 and n_class-1  
)
```

AutoML also generates MLFlow code automatically that tracks all experiments and models in Databricks and can be reused as well going forward. I didn't have to write any of this myself:


```
# Enable automatic logging of input samples, metrics, parameters, and models
mlflow.sklearn.autolog(log_input_examples=True, silent=True)

with mlflow.start_run(experiment_id="1860757565312910", run_name="xgboost") as mlflow_run:
    # AutoML balanced the data internally and use _automl_sample_weight_4cd4 to calibrate the probability distribution
    xgbc_sample_weight = X_train.loc[:, "_automl_sample_weight_4cd4"].to_numpy()

    model.fit(X_train, y_train, classifier__early_stopping_rounds=5, classifier__verbose=False, classifier__eval_set=[(X_val_processed, y_val_processed)])

    # Log metrics for the training set
    xgbc_training_metrics = mlflow.sklearn.eval_and_log_metrics(model, X_train, y_train, prefix="training_", sample_weight=xgbc_sample_weight)

    # Log metrics for the validation set
    xgbc_val_metrics = mlflow.sklearn.eval_and_log_metrics(model, X_val, y_val, prefix="val_")

    # Log metrics for the test set
    xgbc_test_metrics = mlflow.sklearn.eval_and_log_metrics(model, X_test, y_test, prefix="test_")

    # Display the logged metrics
    xgbc_val_metrics = {k.replace("val_", ""): v for k, v in xgbc_val_metrics.items()}
    xgbc_test_metrics = {k.replace("test_", ""): v for k, v in xgbc_test_metrics.items()}
    display(pd.DataFrame([xgbc_val_metrics, xgbc_test_metrics], index=["validation", "test"]))
```

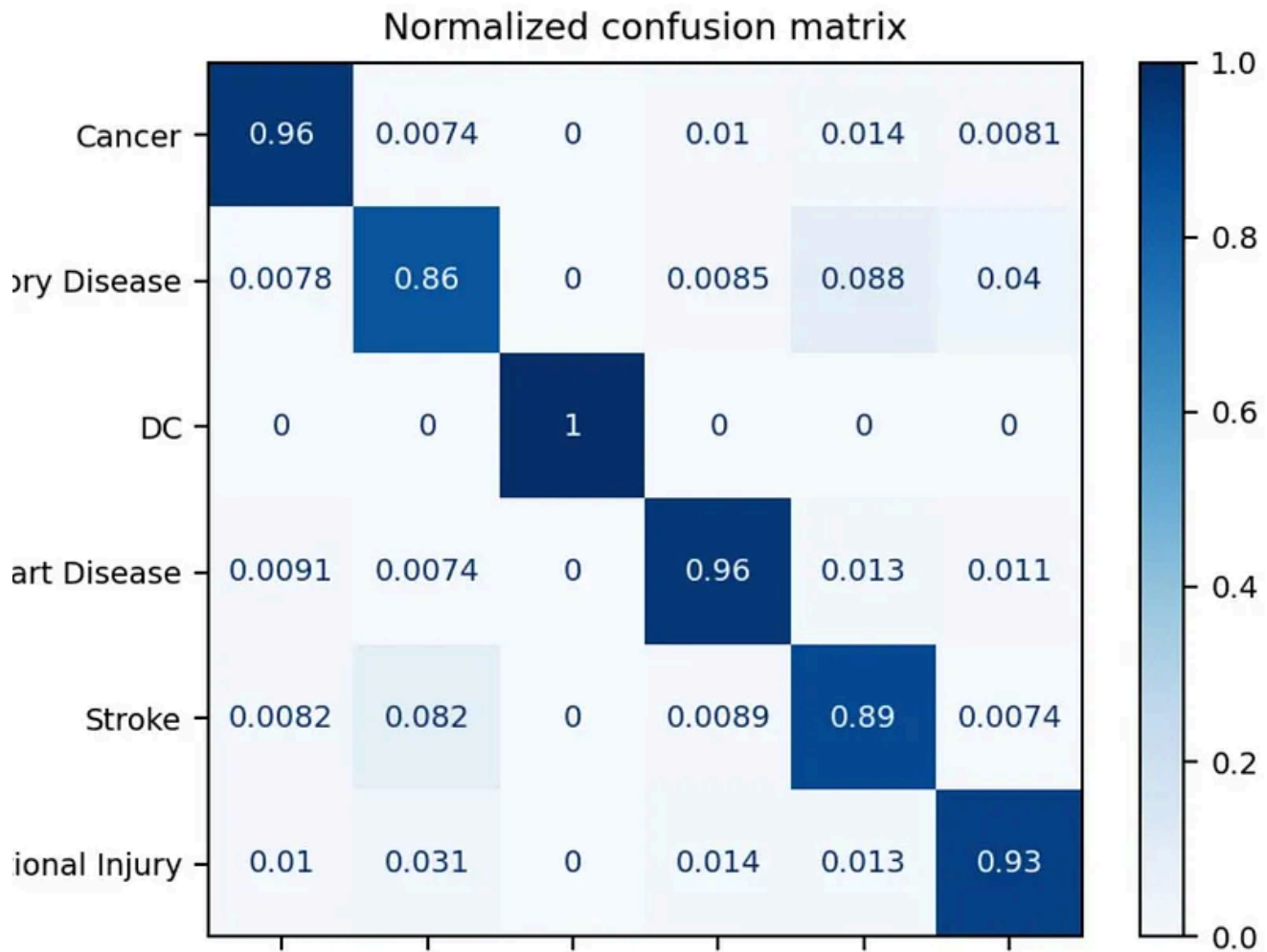
/databricks/python/lib/python3.9/site-packages/xgboost/sklearn.py:1224: UserWarning: The use of label encoder in XGBClassifier is deprecated and will be removed in a future release. To remove this warning, do the following: 1) Pass option use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

warnings.warn(label_encoder_deprecation_msg, UserWarning)

2022/10/05 20:20:26 WARNING mlflow.sklearn.utils: roc_auc_score failed. The metric training_roc_auc_score will not be recorded. Metric error: sample_weight is not supported for multiclass one-vs-one ROC AUC, 'sample_weight' must be None in this case.

	precision_score	recall_score	f1_score	accuracy_score	log_loss	roc_auc_score	score
validation	0.921369	0.921325	0.921315	0.921325	0.187503	0.995584	0.921325

In addition to that, the auto-generation also includes code to determine feature importance, stubs how to register and load the generated model, and even a confusion matrix generated on the validation data AutoML produced internally:



Part 3: Build End-To-End MLOps Workflow From Data Ingest To Model Deployment

What I did next is basically reuse the auto-generated model code from the best AutoML trial and enhance it with data ingest and model deployment code so the resulting notebook [here](#) can be run from start to finish and that way each time pull fresh data from the public CDC URL, save it to the Databricks default database, train an XGBoost model with the best hyperparameters from the top AutoML run, do batch inference on the holdout dataset, and in case the batch predictions are better than the 90 percent success threshold, automatically deploy the model into the [Databricks Model Registry](#) so it can be used either for batch predictions or real-time model serving via a REST API interface.

Since early stopping is enabled on the MLflow experiment, model training now is much quicker since it doesn't improve anymore on the previously determined best hyperparameters and just stops right there, and the model metrics look like this:

```
# Log metrics for the test set
xgbc_test_metrics = mlflow.sklearn.eval_and_log_metrics(model, X_test, y_test, prefix="test_")

# Display the logged metrics
xgbc_val_metrics = {k.replace("val_", ""): v for k, v in xgbc_val_metrics.items()}
xgbc_test_metrics = {k.replace("test_", ""): v for k, v in xgbc_test_metrics.items()}
display(pd.DataFrame([xgbc_val_metrics, xgbc_test_metrics], index=["validation", "test"]))
```

▼ (1) MLflow run

Logged 1 run to an [experiment](#) in MLflow. [Learn more](#)

/databricks/python/lib/python3.9/site-packages/xgboost/sklearn.py:1224: UserWarning: The use of label encoder in XGBClassifier is deprecated and will be removed in a future release. To remove this warning, do the following: 1) Pass option use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

warnings.warn(label_encoder_deprecation_msg, UserWarning)

Table Data Profile

	precision_score ▲	recall_score ▲	f1_score ▲	accuracy_score ▲	log_loss ▲	roc_auc_score ▲
1	0.9175218139817276	0.9177723269037796	0.9176010479449093	0.9177723269037796	0.20029566126041656	0.9955439093436397
2	0.9281755783860051	0.928284854563691	0.9280877351868554	0.928284854563691	0.18478242467713377	0.9964303373217934

Next I am defining input and output for batch inference on the holdout dataset:

excess_deaths_nchs_end_to_end Python ▾

File Edit View Run Help Last edit was 1 hour ago Give feedback

Objective: Predicting t...
Public CDC dataset...
Download data fro...
Rename columns t...
Split in train and h...
Write tables to Dat...
Resume where aut...
Load training data ...
Select supported c...
Preprocessors
Numerical columns
Categorical columns
Low-cardinality cat...
Train - Validation - ...
Train classification ...
Define input and o...
Load model and ru...

Command took 0.07 seconds -- by qd13@cdc.gov at 10/11/2022, 8:34:00 AM on EDAV_Demo_ML-CPU_Cluster_Python_SQL_STANDARDMODE

Cmd 36

Define input and output for batch inference

The table path assigned to `input_table_name` will be used for batch inference and the predictions will be saved to `output_table_path`.

Cmd 37

```
# redefining key variables here because %pip and %conda restarts the Python interpreter
input_table_name = "default.excess_deaths_nchs_holdout"
output_table_path = "/FileStore/batch-inference/excess_deaths_nchs"
```

Command took 0.03 seconds -- by qd13@cdc.gov at 10/11/2022, 8:35:33 AM on EDAV_Demo_ML-CPU_Cluster_Python_SQL_STANDARDMODE

Cmd 38

```
# load table as a Spark DataFrame
table = spark.table(input_table_name)
```

► table: pyspark.sql.dataframe.DataFrame = [year: string, cause_of_death: string ... 11 more fields]

Command took 0.73 seconds -- by qd13@cdc.gov at 10/11/2022, 8:35:40 AM on EDAV_Demo_ML-CPU_Cluster_Python_SQL_STANDARDMODE

Load the generated model and run batch inference on the holdout dataset:

excess_deaths_nchs_end_to_end Python ▾

File Edit View Run Help Last edit was 2 hours ago Give feedback

Objective: Predicting t...
Public CDC dataset...
Download data fro...
Rename columns t...
Split in train and h...
Write tables to Dat...
Resume where aut...
Load training data ...
Select supported c...
Preprocessors
Numerical columns
Categorical columns
Low-cardinality cat...
Train - Validation - ...
Train classification ...
Define input and o...
Load model and ru...
If percentage of co...

Cmd 40

Load model and run inference

```
import mlflow
from pyspark.sql.functions import struct

model_uri = f"runs://{mlflow_run.info.run_id}/model"

# create spark user-defined function for model prediction
predict = mlflow.pyfunc.spark_udf(spark, model_uri, result_type="string")

2022/10/11 12:35:50 WARNING mlflow.pyfunc: Calling `spark_udf()` with `env_manager="local"` does not recreate the same environment that was used during training, which may lead to errors or inaccurate predictions. We recommend specifying `env_manager="conda"`, which automatically recreates the environment that was used to train the model and performs inference in the recreated environment.
```

Command took 1.52 seconds -- by qd13@cdc.gov at 10/11/2022, 8:35:48 AM on EDAV_Demo_ML-CPU_Cluster_Python_SQL_STANDARDMODE

Cmd 41

```
output_df = table.withColumn("prediction", predict(struct(*table.columns)))
```

► output_df: pyspark.sql.dataframe.DataFrame = [year: string, cause_of_death: string ... 12 more fields]

Command took 0.26 seconds -- by qd13@cdc.gov at 10/11/2022, 8:36:01 AM on EDAV_Demo_ML-CPU_Cluster_Python_SQL_STANDARDMODE

At the time of writing this post, the holdout dataset contained 10563 examples and the generated model predicted 9708 causes of death correctly which comes down to about 91.9% correct predictions:

excess_deaths_nchs_end_to_end

Python

File Edit View Run Help

Last edit was 2 hours ago

Give feedback

Terminated

Schedule

Share

Objective: Predicting t...

Public CDC dataset...

Download data fro...

Rename columns t...

Split in train and h...

Write tables to Dat...

Resume where aut...

Load training data ...

Select supported c...

Preprocessors

Numerical columns

Categorical columns

Low-cardinality cat...

Train - Validation - ...

Train classification ...

Define input and o...

Load model and ru...

If percentage of co...

Update the model

all_pred = output_df.select('year', 'cause_of_death', 'state', 'prediction').count()

all_pred

(2) Spark Jobs

Out[17]: 10563

Command took 0.60 seconds -- by q@hadoop-gov at 10/11/2022, 8:36:41 AM on hdfs-Demo_ML-CPU_Cluster_Python_SQL_STANDARDMODE

Cmd 44

correct_pred = output_df.select('year', 'cause_of_death', 'state', 'prediction').filter('cause_of_death = prediction').count()

correct_pred

(2) Spark Jobs

Out[18]: 9708

Command took 10.81 seconds -- by q@hadoop-gov at 10/11/2022, 8:36:48 AM on hdfs-Demo_ML-CPU_Cluster_Python_SQL_STANDARDMODE

Cmd 45

percentage_correct_preds = correct_pred/all_pred*100

print("Percentage correct predictions: ", percentage_correct_preds)

Percentage correct predictions: 91.9057086055098

Considering that this is multi-class classification and there are 5 possible options for causes of death that is a pretty solid outcome, given the minimal amount of effort regarding data preparation and cleansing, feature selection, feature engineering, etc.

In the end-to-end notebook I then added code to register the model to the Databricks Model Registry in case the correct predictions on the holdout dataset are greater than 90 percent.

excess_deaths_nchs_end_to_endPython

FileEditViewRunHelpLast edit was 7 minutes agoGive feedback

Objective: Predicting t...Public CDC dataset...Download data fro...Rename columns t...Split in train and h...Write tables to Dat...Resume where aut...Load training data ...Select supported c...PreprocessorsNumerical columnsCategorical columnsLow-cardinality cat...Train - Validation - ...Train classification ...

Cmd 46

If percentage of correct predictions on external test set greater than 90%, register model to model registry

Cmd 47

```
success_threshold = 90

if percentage_correct_preds > success_threshold:

    model_name = "excess_deaths_nchs"
    model_uri = f"runs://{mlflow_run.info.run_id}/model"
    registered_model_version = mlflow.register_model(model_uri, model_name)
```

Registered model 'excess_deaths_nchs' already exists. Creating a new version of this model...
2022/10/11 15:18:25 INFO mlflow.tracking._model_registry.client: Waiting up to 300 seconds for model version to finish creation.
Model name: excess_deaths_nchs, version 12
Created version '12' of model 'excess_deaths_nchs'.

Command took 6.60 seconds -- by user@mlflow on 10/11/2022, 11:13:04 AM on mlflow-ml-cpu-cluster-python-sql-standardmode

If a registered model of the same name already exists (as pictured above), a new version of the model will automatically be created and all previous model versions will still be in there as well.

Finally, promote the newly minted model to 'Production' stage and set the previous one to 'Archived'.



In Databricks Model Registry that will look something like this:

Microsoft Azure | databricks | Search | CTRL + P

Machine Learning

Create

Workspace

Repos

Recents

Data

Compute

Workflows

Experiments

Feature Store

Models

Partner Connect

Menu options

Registered Models > excess_deaths_nchs

excess_deaths_nchs

Permissions Use model for inference

Details Serving

Notify me about All new activity

Created Time: 2022-10-06 15:00:48 Last Modified: 2022-10-11 11:18:32 Creator: qdt@bndk.gov

Description Edit

Tags

Versions All Active 1 Compare

Version	Registered at	Created by	Stage	Pending Requests	Description
Version 12	2022-10-11 11:18:25	qdt@bndk.gov	Production	-	
Version 11	2022-10-11 11:01:59	qdt@bndk.gov	Archived	-	
Version 10	2022-10-11 08:37:38	qdt@bndk.gov	Archived	-	

Part 4: Schedule MLOps Pipeline And Make Real-Time Predictions

Now, to make this an automated, repeatable process that will build and deploy a new model on fresh data in case the CDC updates the source dataset and the new model surpasses the set threshold of 90 percent correct predictions on the new holdout dataset, all that needs to be done is to schedule a Databricks Workflow which can be achieved on the notebook page view by clicking on 'Schedule':

Microsoft Azure | databricks | Search | CTRL + P

Machine Learning

Create

Workspace

Repos

Recents

Data

Compute

Workflows

Experiments

Feature Store

Models

Partner Connect

excess_deaths_nchs_end_to_end Python

File Edit View Run Help Last edit was 1 hour ago Give feedback

Terminated Schedule Share

Databricks ML End-To-...

Objective: Predicting t...

Public CDC dataset...

Download data fro...

Rename columns t...

Split in train and h...

Write tables to Dat...

Resume where out...

Load training data ...

Select supported c...

Preprocessors

Numerical columns

Categorical columns

Low-cardinality cat...

Train - Validation - ...

Train classification ...

Define input and o...

Job name excess_deaths_nchs_end_to_end

Schedule Manual Scheduled

Every Week on Tuesday at 14 : 00 (UTC-04:00) Eastern Time...

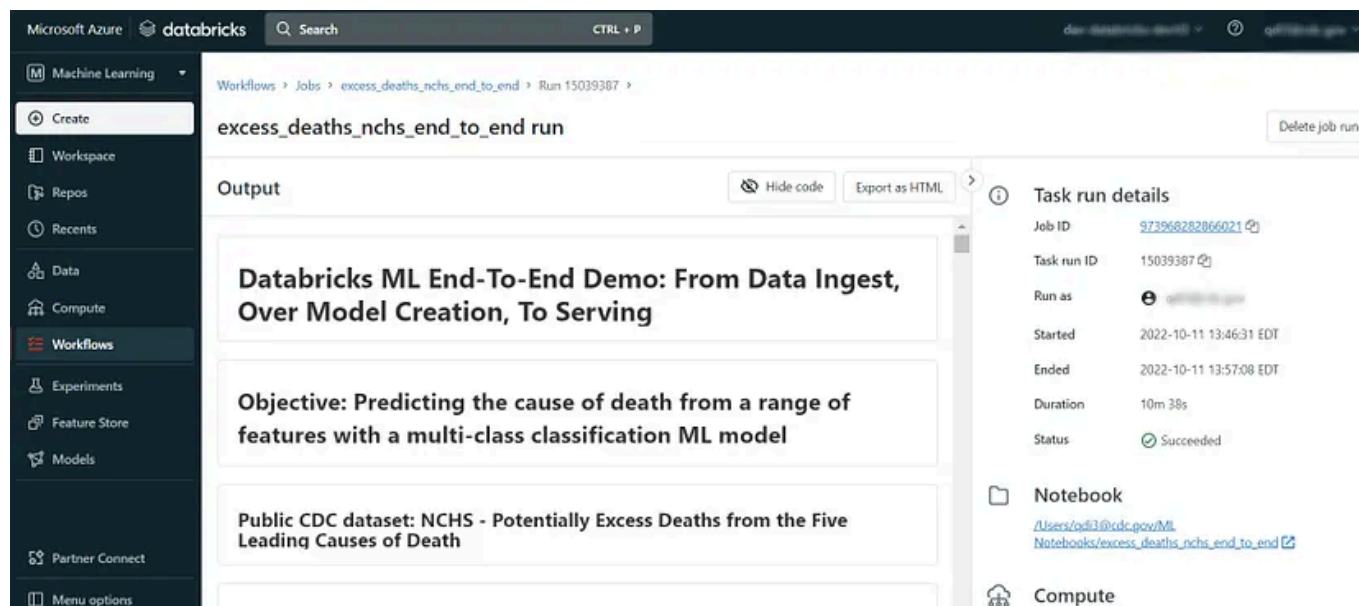
Cluster Demo_ML-CPU_Cluster_Python_SQL_STANDARDMODE D8R 11.2 ML - Spark 3.3.0 ...

Parameters Add

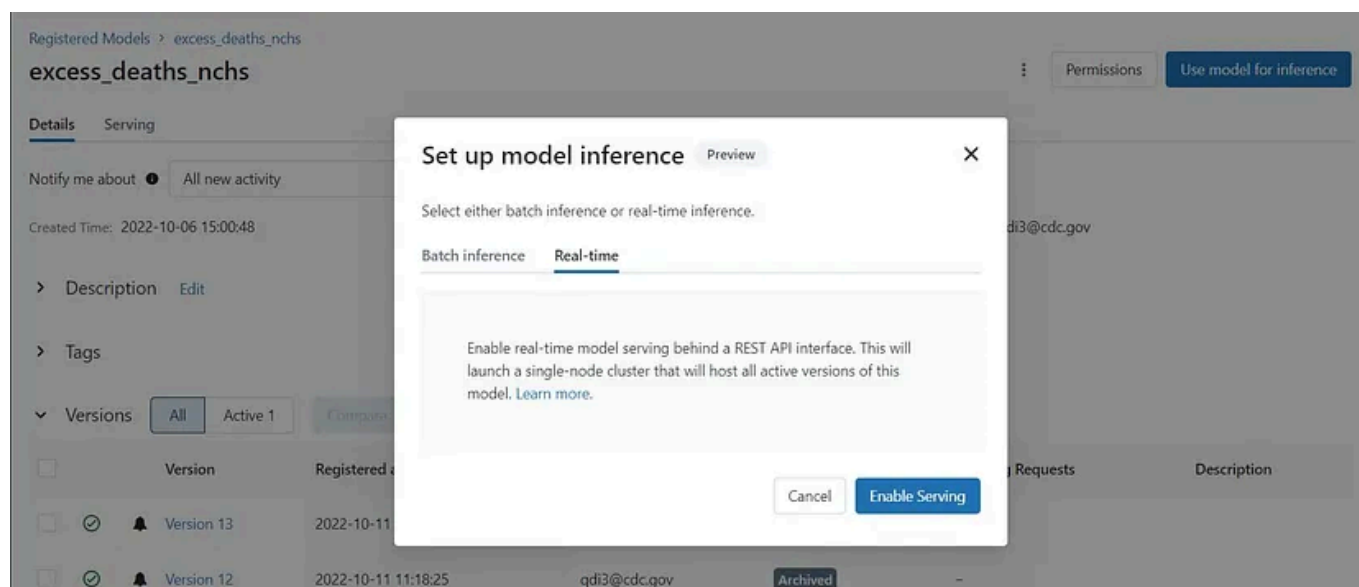
Alerts qdt@bndk.gov Start Success Failure Add

Cancel Create

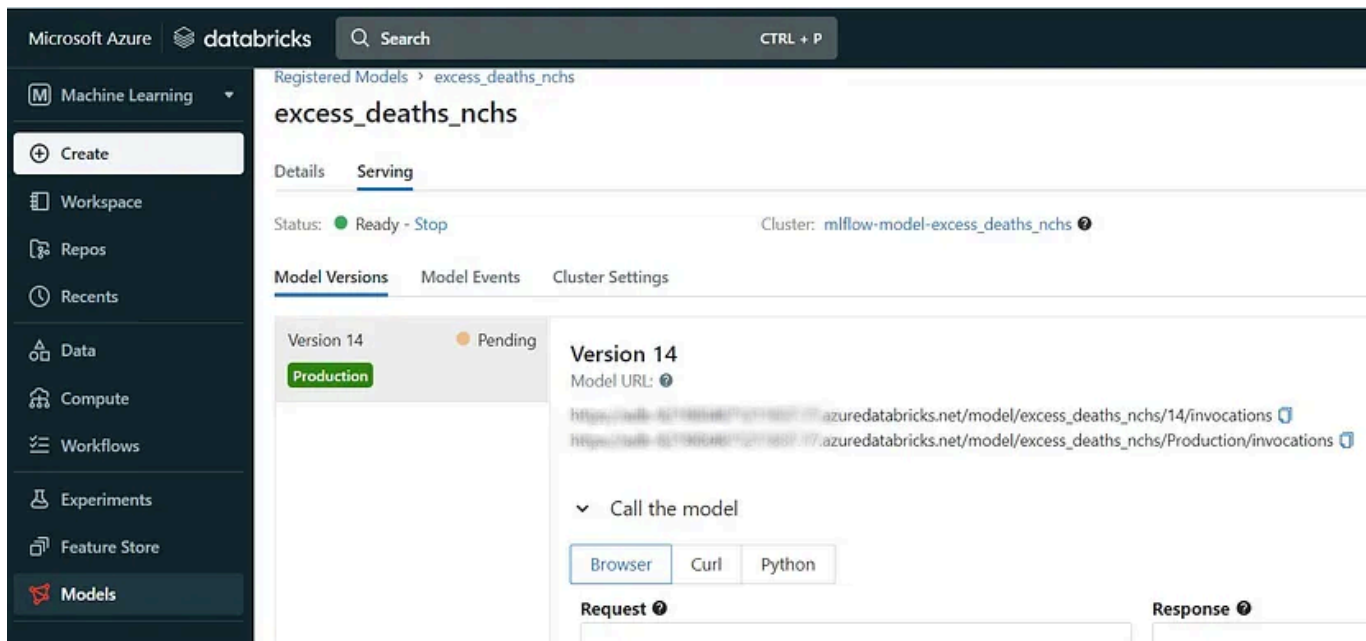
When a run has finished and succeeded, you will receive an email with a link to the task run details:



To use the 'Production' version of the model for real-time inference, serving needs to be enabled:



This will take a little while to deploy and once ready look as follows:



Microsoft Azure | databricks | Search | CTRL + P

Machine Learning

Create

Workspace

Repos

Recents

Data

Compute

Workflows

Experiments

Feature Store

Models

Registered Models > excess_deaths_nchs

excess_deaths_nchs

Details | Serving

Status: ● Ready - Stop Cluster: mlflow-model-excess_deaths_nchs

Model Versions | Model Events | Cluster Settings

Version 14 ● Pending

Production

Version 14

Model URL: https://adb-617411181219807.azuredatabricks.net/model/excess_deaths_nchs/14/invocations

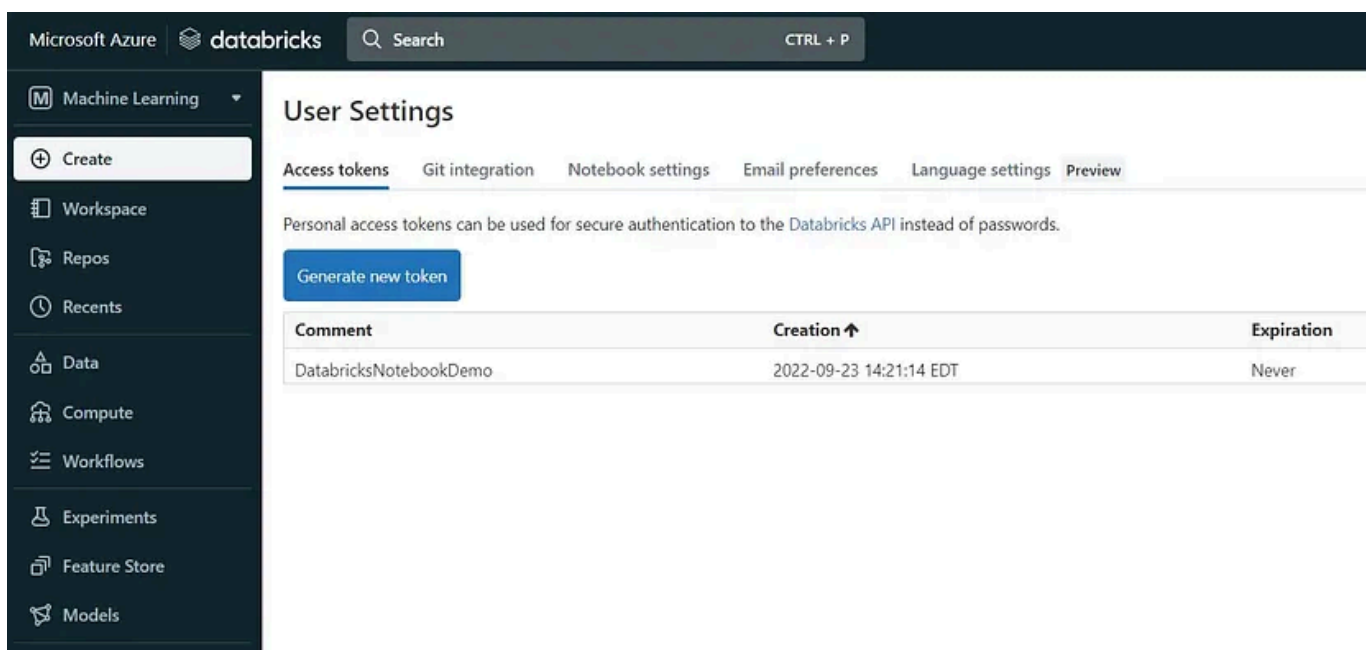
https://adb-617411181219807.azuredatabricks.net/model/excess_deaths_nchs/Production/invocations

Call the model

Browser | Curl | Python

Request [?](#) Response [?](#)

To make real-time predictions through the Databricks REST API using the deployed model whose version is going to update automatically each time a new model is created that surpasses the 90 percent success threshold, a Databricks token is required. The token can be created through the User Settings page:



Microsoft Azure | databricks | Search | CTRL + P

Machine Learning

Create

Workspace

Repos

Recents

Data

Compute

Workflows

Experiments

Feature Store

Models

User Settings

Access tokens | Git integration | Notebook settings | Email preferences | Language settings | Preview

Personal access tokens can be used for secure authentication to the Databricks API instead of passwords.

[Generate new token](#)

Comment	Creation ↑	Expiration
DatabricksNotebookDemo	2022-09-23 14:21:14 EDT	Never

In a terminal, add the token to your environment variables:

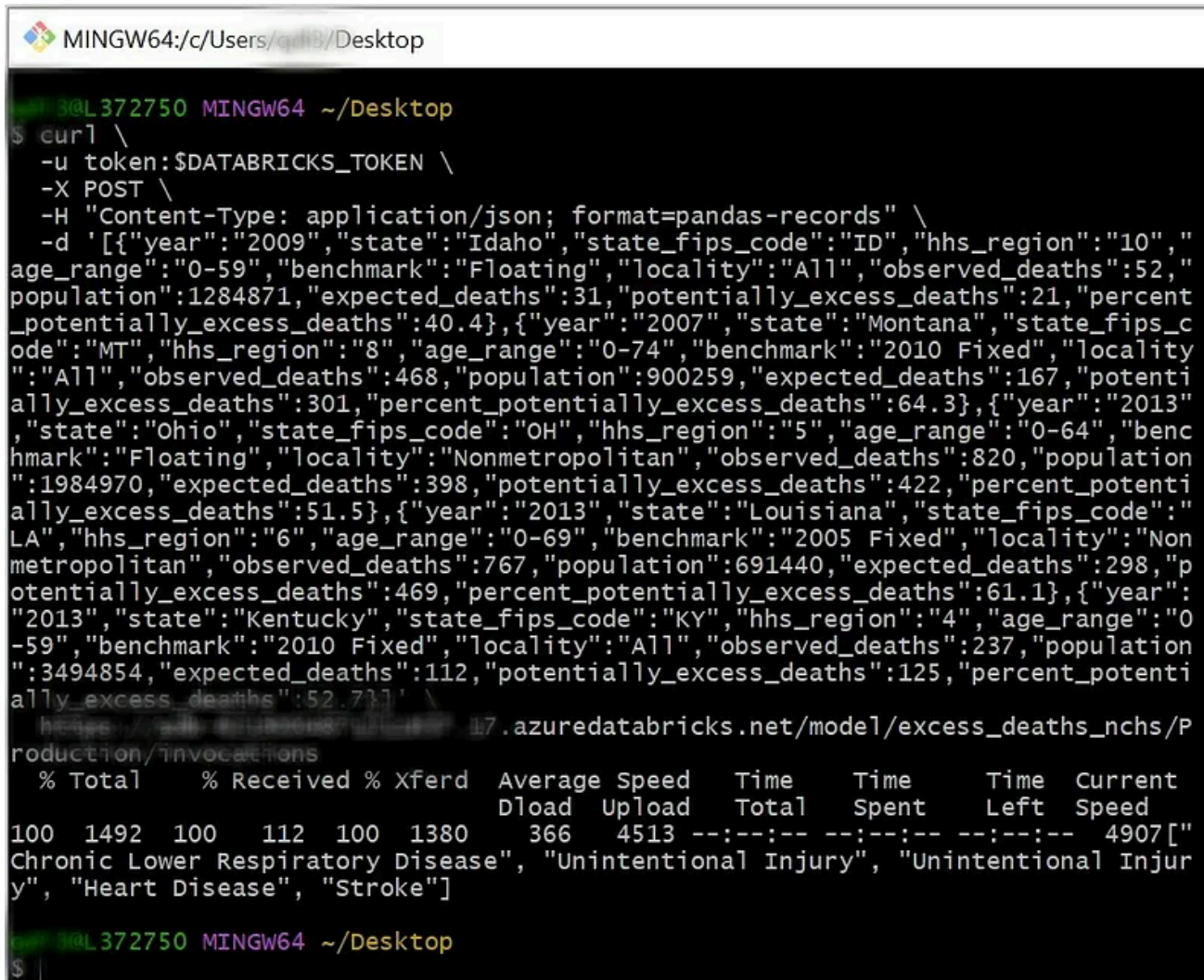


Make predictions using cURL and the following example command (you will have to update the model URL at the end to your own):

```
curl \
  -u token:$DATABRICKS_TOKEN \
  -X POST \
  -H "Content-Type: application/json; format=pandas-records" \
  -d
' [{"year": "2009", "state": "Idaho", "state_fips_code": "ID", "hhs_region":
"10", "age_range": "0-
59", "benchmark": "Floating", "locality": "All", "observed_deaths": 52, "pop
ulation": 1284871, "expected_deaths": 31, "potentially_excess_deaths": 21,
"percent_potentially_excess_deaths": 40.4},
{"year": "2007", "state": "Montana", "state_fips_code": "MT", "hhs_region":
"8", "age_range": "0-74", "benchmark": "2010
Fixed", "locality": "All", "observed_deaths": 468, "population": 900259, "ex
pected_deaths": 167, "potentially_excess_deaths": 301, "percent_potential
ly_excess_deaths": 64.3},
{"year": "2013", "state": "Ohio", "state_fips_code": "OH", "hhs_region": "5"
, "age_range": "0-
64", "benchmark": "Floating", "locality": "Nonmetropolitan", "observed_dea
ths": 820, "population": 1984970, "expected_deaths": 398, "potentially_exce
ss_deaths": 422, "percent_potentially_excess_deaths": 51.5},
{"year": "2013", "state": "Louisiana", "state_fips_code": "LA", "hhs_region
": "6", "age_range": "0-69", "benchmark": "2005
Fixed", "locality": "Nonmetropolitan", "observed_deaths": 767, "population
": 691440, "expected_deaths": 298, "potentially_excess_deaths": 469, "perce
nt_potentially_excess_deaths": 61.1},
{"year": "2013", "state": "Kentucky", "state_fips_code": "KY", "hhs_region"
: "4", "age_range": "0-59", "benchmark": "2010
Fixed", "locality": "All", "observed_deaths": 237, "population": 3494854, "e
xpected_deaths": 112, "potentially_excess_deaths": 125, "percent_potential
ly_excess_deaths": 52.7}] ]' \
  https://adb-
```

[XXX.azuredatabricks.net/model/excess_deaths_nchs/Production/invocations](https://xxx.azuredatabricks.net/model/excess_deaths_nchs/Production/invocations)

The results will be returned as a list:



```

MINGW64/c/Users/qd18/Desktop
3@L372750 MINGW64 ~/Desktop
$ curl \
  -u token:$DATABRICKS_TOKEN \
  -X POST \
  -H "Content-Type: application/json; format=pandas-records" \
  -d '[{"year": "2009", "state": "Idaho", "state_fips_code": "ID", "hhs_region": "10", "age_range": "0-59", "benchmark": "Floating", "locality": "All", "observed_deaths": 52, "population": 1284871, "expected_deaths": 31, "potentially_excess_deaths": 21, "percent_potentially_excess_deaths": 40.4}, {"year": "2007", "state": "Montana", "state_fips_code": "MT", "hhs_region": "8", "age_range": "0-74", "benchmark": "2010 Fixed", "locality": "All", "observed_deaths": 468, "population": 900259, "expected_deaths": 167, "potentially_excess_deaths": 301, "percent_potentially_excess_deaths": 64.3}, {"year": "2013", "state": "Ohio", "state_fips_code": "OH", "hhs_region": "5", "age_range": "0-64", "benchmark": "Floating", "locality": "Nonmetropolitan", "observed_deaths": 820, "population": 1984970, "expected_deaths": 398, "potentially_excess_deaths": 422, "percent_potentially_excess_deaths": 51.5}, {"year": "2013", "state": "Louisiana", "state_fips_code": "LA", "hhs_region": "6", "age_range": "0-69", "benchmark": "2005 Fixed", "locality": "Nonmetropolitan", "observed_deaths": 767, "population": 691440, "expected_deaths": 298, "potentially_excess_deaths": 469, "percent_potentially_excess_deaths": 61.1}, {"year": "2013", "state": "Kentucky", "state_fips_code": "KY", "hhs_region": "4", "age_range": "0-59", "benchmark": "2010 Fixed", "locality": "All", "observed_deaths": 237, "population": 3494854, "expected_deaths": 112, "potentially_excess_deaths": 125, "percent_potentially_excess_deaths": 52.733}' \
  https://xxx.azuredatabricks.net/model/excess_deaths_nchs/Production/invocations
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
           Dload  Upload   Total     Spent    Left  Speed
100 1492  100   112  100   1380    366   4513  --:--:--  --:--:--  --:--:--  4907["Chronic Lower Respiratory Disease", "Unintentional Injury", "Unintentional Injury", "Heart Disease", "Stroke"]
3@L372750 MINGW64 ~/Desktop
$ |

```

Predictions can also be made directly in the browser or using Python.

Call the model

Browser Curl Python

Request

```
[
  {
    "year": "2010",
    "state": "Alabama",
    "state_fips_code": "AL",
    "hhs_region": "4",
    "age_range": "0-59",
    "benchmark": "Floating",
```

Response

```
[["Unintentional Injury", "Chronic Lower Respiratory Disease", "Heart Disease", "Chronic Lower Respiratory Disease", "Cancer"]]
```

Send Request Show Example

Call the model

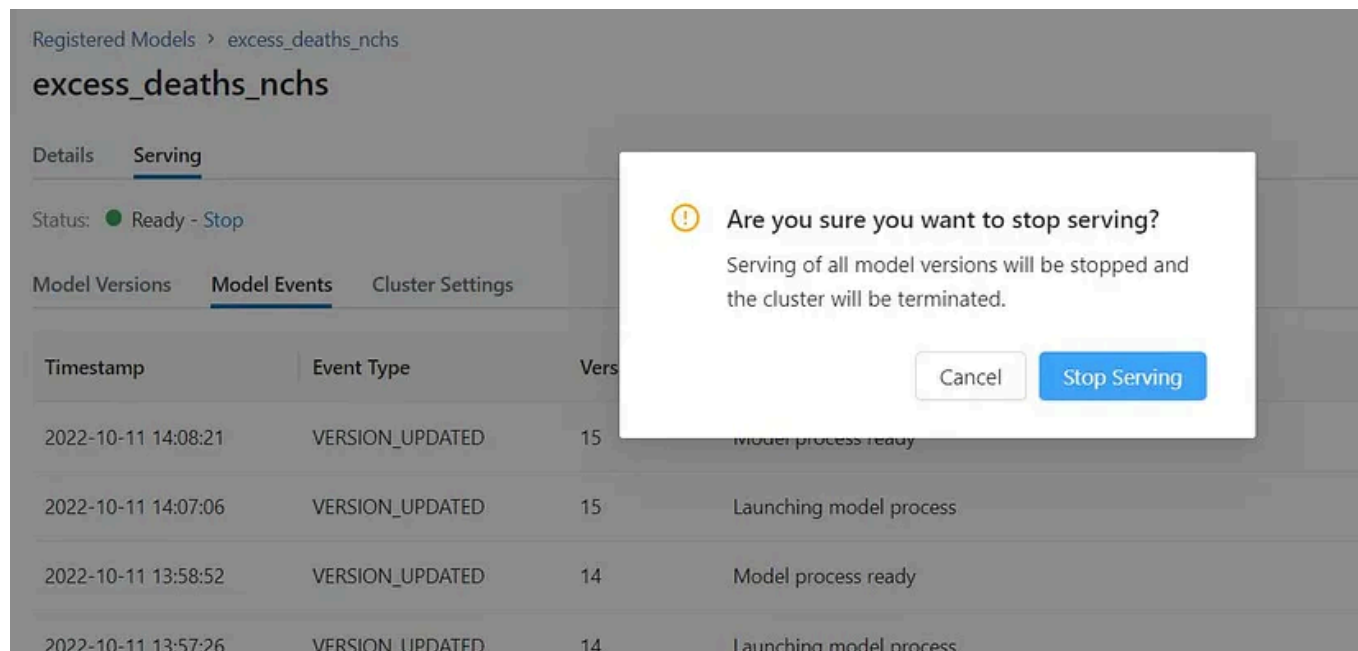
Browser Curl Python

```
import os
import requests
import numpy as np
import pandas as pd
import json

def create_tf_serving_json(data):
    return {'inputs': {name: data[name].tolist() for name in data.keys()} if isinstance(data, dict) else data.tolist()}

def score_model(dataset):
    url = 'https://adb-az1111111111111111.azuredatabricks.net/model/excess_deaths_nchs/15/invocations'
    headers = {'Authorization': f'Bearer {os.environ.get("DATABRICKS_TOKEN")}', 'Content-Type': 'application/json'}
    ds_dict = dataset.to_dict(orient='split') if isinstance(dataset, pd.DataFrame) else create_tf_serving_json(dataset)
    data_json = json.dumps(ds_dict, allow_nan=True)
    response = requests.request(method='POST', headers=headers, url=url, data=data_json)
    if response.status_code != 200:
        raise Exception(f'Request failed with status {response.status_code}, {response.text}')
    return response.json()
```

A version of a model with the status 'Production' cannot be deleted before the status has been changed to e.g. 'Archived', but the model server can be stopped anytime the model is not needed for real-time serving:



Registered Models > excess_deaths_nchs

excess_deaths_nchs

Details **Serving**

Status: ● Ready - Stop

Model Versions **Model Events** Cluster Settings

Timestamp	Event Type	Version	Description
2022-10-11 14:08:21	VERSION_UPDATED	15	Model process ready
2022-10-11 14:07:06	VERSION_UPDATED	15	Launching model process
2022-10-11 13:58:52	VERSION_UPDATED	14	Model process ready
2022-10-11 13:57:26	VERSION_UPDATED	14	Launching model process

If you need it again, just enable serving once again.

As I set it up above with the scheduled Databricks Workflow, the data-ingest-to-model-deployment notebook would now run on a weekly basis, download potentially updated source data from the CDC website, build a new XGBoost model with the new data, and test it on new data as well. Only if the predictions on the holdout data are better than 90 percent correct the new model will be deployed to production. In any case I would receive an email each time the process runs and if it succeeded or failed, and all the created model artifacts would be available for inspection and troubleshooting in case the predictive power of the model declines below the success threshold, or any other issues or problems arise.

Bottom Line

Databricks provides all the tools to develop an end-to-end Machine Learning solution with relative ease. There are certainly limitations (e.g. currently Databricks AutoML only works on tabular data, IDE support to be able to develop locally is limited and difficult to set up, one still has to create and

manage and wait on Spark clusters to become ready, just to name a few...), but overall it is a nice environment with a pleasant UI to build and manage Machine Learning code and artifacts.

Machine Learning

Databricks

Mlops

Python

Spark

**Written by Thomas Jaensch**[Edit profile](#)

21 Followers

WFH, staring at code

More from Thomas Jaensch