- snowflake-ml-python version: 1.2.2
- Feature Store PrPr Version: 0.5.1
- Updated date: 2/12/2024

# Basic Feature Demo

This notebook demonstrates feature store with simple features. It includes an end-2-end ML experiment cycle: feature creation, training and inference. It also demonstrate the interoperation between Feature Store and Model Registry.

```python
In [ ]:  from snowflake.snowpark import Session
         from snowflake.snowpark import functions as F
         from snowflake.ml.feature_store import (
             FeatureStore,
             FeatureView,
             Entity,
             CreationMode
         )
         from snowflake.ml.utils.connection_params import SnowflakeLoginOptions
```

# Setup Snowflake connection and database

For detailed session connection config, please follow this tutorial.

```python
In [ ]:  session = Session.builder.configs(SnowflakeLoginOptions()).create()
```

Below cell creates temporary database, schema and warehouse for this notebook. All temporary resources will be deleted at the end of this notebook. You can rename with your own name if needed.

```python
In [ ]:  # database name where test data and feature store lives.
         FS_DEMO_DB = f"FEATURE_STORE_BASIC_FEATURE_NOTEBOOK_DEMO"
         # schema where test data lives.
         TEST_DATASET_SCHEMA = 'TEST_DATASET'
         # feature store name.
         FS_DEMO_SCHEMA = "AWESOME_FS_BASIC_FEATURES"
         # Model registry database name.
         MR_DEMO_DB = f"FEATURE_STORE_BASIC_FEATURE_NOTEBOOK_MR_DEMO"
         # warehouse name used in this notebook.
         FS_DEMO_WH = "FEATURE_STORE_BASIC_FEATURE_NOTEBOOK_DEMO"

         session.sql(f"DROP DATABASE IF EXISTS {FS_DEMO_DB}").collect()
         session.sql(f"DROP DATABASE IF EXISTS {MR_DEMO_DB}").collect()
         session.sql(f"CREATE DATABASE IF NOT EXISTS {FS_DEMO_DB}").collect()
         session.sql(f"""
             CREATE SCHEMA IF NOT EXISTS
```

```
    {FS_DEMO_DB}.{TEST_DATASET_SCHEMA}""").collect()
session.sql(f"CREATE WAREHOUSE IF NOT EXISTS {FS_DEMO_WH}").collect()
```

# Create FeatureStore Client

Let's first create a feature store client.

We can pass in an existing database name, or a new database will be created upon the feature store initialization. Replace `DEMO_DB` and `DEMO_SCHEMA` with your database and schema.

In [ ]:
```python
fs = FeatureStore(
    session=session,
    database=FS_DEMO_DB,
    name=FS_DEMO_SCHEMA,
    default_warehouse=FS_DEMO_WH,
    creation_mode=CreationMode.CREATE_IF_NOT_EXIST,
)
```

# Prepare test data

We will use wine quality dataset for this demo. Download the public dataset from kaggle if you dont have it already: https://www.kaggle.com/datasets/uciml/red-wine-quality-cortez-et-al-2009. Replace `TEST_CSV_FILE_PATH` with your local file path.

In [ ]:
```python
TEST_CSV_FILE_PATH = 'winequality-red.csv'
session.file.put(
    f"file://{TEST_CSV_FILE_PATH}",session.get_session_stage())

from snowflake.snowpark.types import (
    StructType,
    StructField,
    IntegerType,
    FloatType
)
input_schema = StructType(
    [
        StructField("fixed_acidity", FloatType()),
        StructField("volatile_acidity", FloatType()),
        StructField("citric_acid", FloatType()),
        StructField("residual_sugar", FloatType()),
        StructField("chlorides", FloatType()),
        StructField("free_sulfur_dioxide", IntegerType()),
        StructField("total_sulfur_dioxide", IntegerType()),
        StructField("density", FloatType()),
        StructField("pH", FloatType()),
        StructField("sulphates", FloatType()),
        StructField("alcohol", FloatType()),
        StructField("quality", IntegerType())
    ]
)
```

```
df = session.read.options({"field_delimiter": ";", "skip_header": 1}) \
    .schema(input_schema) \
    .csv(f"{session.get_session_stage()}/winequality-red.csv")
full_table_name = f"{FS_DEMO_DB}.{TEST_DATASET_SCHEMA}.WINE_DATA"
df.write.mode("overwrite").save_as_table(full_table_name)
```

## Create and register a new Entity

We will create an Entity called *wine* and register it with the feature store.

You can retrieve the active Entities in the feature store with list_entities() API.

```
In [ ]:  entity = Entity(name="WINE", join_keys=["WINE_ID"])
         fs.register_entity(entity)
         fs.list_entities().show()
```

## Load source data and do some simple feature engineering

Then we will load from the source table and conduct some simple feature engineerings.

Here we are just doing two simple data manipulation (but more complex ones are carried out the same way):

1. Assign a WINE_ID column to the source
2. Derive a new column by multipying two existing feature columns

```
In [ ]:  source_df = session.table(full_table_name)
         source_df.show()
```

```
In [ ]:  def addIdColumn(df, id_column_name):
             # Add id column to dataframe
             columns = df.columns
             new_df = df.withColumn(id_column_name, F.monotonically_increasing_id())
             return new_df[[id_column_name] + columns]

         def generate_new_feature(df):
             # Derive a new feature column
             return df.withColumn(
                 "MY_NEW_FEATURE", df["FIXED_ACIDITY"] * df["CITRIC_ACID"])

         df = addIdColumn(source_df, "WINE_ID")
         feature_df = generate_new_feature(df)
         feature_df = feature_df.select(
             [
                 'WINE_ID',
                 'FIXED_ACIDITY',
                 'VOLATILE_ACIDITY',
                 'CITRIC_ACID',
                 'RESIDUAL_SUGAR',
```

```
        'CHLORIDES',
        'FREE_SULFUR_DIOXIDE',
        'TOTAL_SULFUR_DIOXIDE',
        'DENSITY',
        'PH',
        'MY_NEW_FEATURE',
    ]
)
feature_df.show()
```

## Create a new FeatureView and materialize the feature pipeline

Once the FeatureView construction is done, we can materialize the FeatureView to the Snowflake backend and incremental maintenance will start.

In [ ]:
```python
# NOTE:
# Due to a known issue on backend pipeline creation,
# if the source data is created right before the
# feature pipeline, there might be a chance for
# dataloss, so sleep for 60s for now.
# This issue will be fixed soon in upcoming patch.

import time
time.sleep(60)
```

In [ ]:
```python
fv = FeatureView(
    name="WINE_FEATURES",
    entities=[entity],
    feature_df=feature_df,
    refresh_freq="1 minute",
    desc="wine features"
)
fv = fs.register_feature_view(
    feature_view=fv,
    version="V1",
    block=True
)
```

In [ ]:
```python
# Examine the FeatureView content
fs.read_feature_view(fv).show()
```

## Explore additional features

Now I have my FeatureView created with a collection of features, but what if I want to explore additional features on top?

Since a materialized FeatureView is immutable (due to singe DDL for the backend dynamic table), we will need to create a new FeatureView for the additional features and then merge them.

```
In [ ]:  extra_feature_df = df.select(
             [
                 'WINE_ID',
                 'SULPHATES',
                 'ALCOHOL',
             ]
         )

         new_fv = FeatureView(
             name="EXTRA_WINE_FEATURES",
             entities=[entity],
             feature_df=extra_feature_df,
             refresh_freq="1 minute",
             desc="extra wine features"
         )
         new_fv = fs.register_feature_view(
             feature_view=new_fv,
             version="V1",
             block=True
         )
```

```
In [ ]:  # We can easily retrieve all FeatureViews for a given Entity.
         fs.list_feature_views(entity_name="WINE"). \
             select(["NAME", "ENTITIES", "FEATURE_DESC"]).show()
```

## Create new feature view with combined feature results [Optional]

Now we have two FeatureViews ready, we can choose to create a new one by merging the two (it's just like a join and we provide a handy function for that). The new FeatureView won't incur the cost of feature pipelines but only the table join cost.

Obviously we can also just work with two separate FeatureViews (most of our APIs support multiple FeatureViews), the capability of merging is just to make the features better organized and easier to share.

```
In [ ]:  full_fv = fs.merge_features(
             features=[fv, new_fv], name="FULL_WINE_FEATURES")
         full_fv = fs.register_feature_view(feature_view=full_fv, version="V1")
```

## Generate Training Data

After our feature pipelines are fully setup, we can start using them to generate training data and later do model prediction.

```
In [ ]:  spine_df = session.table(f"{FS_DEMO_DB}.{TEST_DATASET_SCHEMA}.WINE_DATA")
         spine_df = addIdColumn(source_df, "WINE_ID")
         spine_df = spine_df.select("WINE_ID", "QUALITY")
         spine_df.show()
```

Generate training data is easy since materialized FeatureViews already carry most of the metadata like join keys, timestamp for point-in-time lookup, etc. We just need to provide the spine data (it's called spine because we are essentially enriching the data by joining features with it). We can also generate dataset with a subset of features in the feature view by `slice`.

```python
In [ ]: training_dataset_full_path = \
            f"{FS_DEMO_DB}.{FS_DEMO_SCHEMA}.WINE_TRAINING_DATA_TABLE"
        session.sql(f"DROP TABLE IF EXISTS {training_dataset_full_path}") \
            .collect()
        training_data = fs.generate_dataset(
            spine_df=spine_df,
            features=[
                full_fv.slice([
                    "FIXED_ACIDITY",
                    "VOLATILE_ACIDITY",
                    "CITRIC_ACID"
                ])
            ],
            materialized_table="WINE_TRAINING_DATA_TABLE",
            spine_timestamp_col=None,
            spine_label_cols=["QUALITY"],
            save_mode="merge",
            exclude_columns=['WINE_ID']
        )

        training_data.df.show()
```

## Train model with Snowpark ML

Now let's training a simple random forest model, and evaluate the prediction accuracy.

```python
In [ ]: from snowflake.ml.modeling.ensemble import RandomForestRegressor
        from snowflake.ml.modeling import metrics as snowml_metrics
        from snowflake.snowpark.functions import abs as sp_abs, mean, col

        def train_model_using_snowpark_ml(training_data):
            train, test = training_data.df.random_split([0.8, 0.2], seed=42)
            feature_columns = [col for col in training_data.df.columns if col != "QU
            label_column = "QUALITY"

            rf = RandomForestRegressor(
                input_cols=feature_columns, label_cols=[label_column],
                max_depth=3, n_estimators=20, random_state=42
            )

            rf.fit(train)
            predictions = rf.predict(test)

            mse = snowml_metrics.mean_squared_error(
                df=predictions,
                y_true_col_names=label_column,
```

```
        y_pred_col_names="OUTPUT_" + label_column)

    accuracy = 100 - snowml_metrics.mean_absolute_percentage_error(
        df=predictions,
        y_true_col_names=label_column,
        y_pred_col_names="OUTPUT_" + label_column
    )

    print(f"MSE: {mse}, Accuracy: {accuracy}")
    return rf

rf = train_model_using_snowpark_ml(training_data)
print(rf)
```

# [Predict Optional 1] With local model

Now let's predict with the model and the feature values retrieved from feature store.

```
In [ ]:  test_df = spine_df.limit(3).select("WINE_ID")

         enriched_df = fs.retrieve_feature_values(
             test_df, training_data.load_features())
         enriched_df = enriched_df.drop('WINE_ID')
```

```
In [ ]:  pred = rf.predict(enriched_df.to_pandas())

         print(pred)
```

# [Predict Option 2] Using Model Registry

## Step 1: Log the model along with its training dataset metadata into Model Registry

```
In [ ]:  from snowflake.ml.registry import model_registry

         registry = model_registry.ModelRegistry(
             session=session,
             database_name=MR_DEMO_DB,
             create_if_not_exists=True
         )
```

Register the dataset into model registry with `log_artifact`. Artifact is a generalized concept of ML pipeline outputs that are needed for subsequent execution. Refer to https://docs.snowflake.com/LIMITEDACCESS/snowflake-ml-model-registry for more details about the API.

```
In [ ]:  DATASET_NAME = "MY_DATASET"
         DATASET_VERSION = "V1"

         my_dataset = registry.log_artifact(
```

```
        artifact=training_data,
        name=DATASET_NAME,
        version=DATASET_VERSION,
    )
```

Now you can log the model together with the registered artifact (which is a dataset here).

```python
In [ ]:  model_name = "MY_MODEL"

         model_ref = registry.log_model(
             model_name=model_name,
             model_version="V2",
             model=rf,
             tags={"author": "my_rf_with_training_data"},
             artifacts=[my_dataset],
             options={"embed_local_ml_library": True},
         )
```

## Step 2 : Restore model and predict with features

We retrieve the training dataset from registry then construct dataframe of latest feature values. Then we restore the model from registry. At last, we can predict with latest feature values.

```python
In [ ]:  from snowflake.ml.dataset.dataset import Dataset

         registered_dataset = registry.get_artifact(
             DATASET_NAME,
             DATASET_VERSION)
         test_df = spine_df.limit(3).select("WINE_ID")

         enriched_df = fs.retrieve_feature_values(
             test_df, registered_dataset.load_features())
         enriched_df = enriched_df.drop('WINE_ID')
```

```python
In [ ]:  model_ref = model_registry.ModelReference(
             registry=registry,
             model_name=model_name,
             model_version="V2"
         )
         restored_model = model_ref.load_model()
         restored_prediction = restored_model.predict(enriched_df.to_pandas())

         print(restored_prediction)
```

## Cleanup notebook

Cleanup resources created in this notebook.

```python
In [ ]:  session.sql(f"DROP DATABASE IF EXISTS {FS_DEMO_DB}").collect()
```

```python
session.sql(f"DROP DATABASE IF EXISTS {MR_DEMO_DB}").collect()
session.sql(f"DROP WAREHOUSE IF EXISTS {FS_DEMO_WH}").collect()
```