

ML Pipeline with Kubeflow on Azure: A Comprehensive Guide

Comprehensive Guide: ML Pipeline with Kubeflow on Azure

Prerequisites

1. System Requirements

- Linux/macOS/Windows with WSL2
- Python 3.8+
- Docker Desktop
- Git

2. Cloud & Tools

```
1 # Install Azure CLI
2 curl -sL <https://aka.ms/InstallAzureCLIDeb> | sudo bash
3
4 # Install kubectl
5 curl -LO "<https://dl.k8s.io/release/$(curl -L -s
6 <https://dl.k8s.io/release/stable.txt>/bin/linux/amd64/kubectl>"
7 sudo install -o root -g root -m 0755 kubectl /usr/local/bin/kubectl
8
9 # Install Kustomize
10 curl -s "<https://raw.githubusercontent.com/kubernetes-sigs/kustomize/master/hack/install_kustomize.sh>" |
11 bash
12 sudo mv kustomize /usr/local/bin/
13
14 # Python packages
15 pip install azure-cli azure-ml-sdk kubernetes kfp azureml-sdk pandas scikit-learn tensorflow
```

3. Azure Resources

- Active Azure subscription
- Resource group
- Storage account
- Azure Container Registry (ACR)

Setup Infrastructure

1. Create AKS Cluster

```
1 # Create resource group
2 az group create --name ml-resources --location eastus
3
4 # Create AKS cluster
5 az aks create \
6   --resource-group ml-resources \
7   --name ml-cluster \
8   --node-count 3 \
9   --enable-addons monitoring \
10  --generate-ssh-keys \
```

```

11 --node-vm-size Standard_DS3_v2
12
13 # Get credentials
14 az aks get-credentials --resource-group ml-resources --name ml-cluster

```

2. Deploy Kubeflow

```

1 # Clone manifests
2 git clone <https://github.com/kubeflow/manifests.git>
3 cd manifests
4
5 # Install Kubeflow
6 while ! kustomize build example | kubectl apply -f -; do
7     echo "Retrying to apply resources";
8     sleep 30;
9 done
10
11 # Verify deployment
12 kubectl get pods -n kubeflow

```

Sample ML Pipeline Implementation

1. Dataset Preparation

```

1 # dataset_prep.py
2 import pandas as pd
3 from sklearn.datasets import load_diabetes
4 from sklearn.model_selection import train_test_split
5
6 def prepare_data():
7     # Load sample diabetes dataset
8     diabetes = load_diabetes()
9     df = pd.DataFrame(diabetes.data, columns=diabetes.feature_names)
10    df['target'] = diabetes.target
11
12    # Split data
13    train_df, test_df = train_test_split(df, test_size=0.2, random_state=42)
14
15    # Save to files
16    train_df.to_csv('train.csv', index=False)
17    test_df.to_csv('test.csv', index=False)
18
19    return 'train.csv', 'test.csv'
20
21 if __name__ == '__main__':
22    prepare_data()

```

2. Model Training

```

1 # train.py
2 import pandas as pd
3 import numpy as np
4 from sklearn.linear_model import Ridge
5 import pickle
6
7 def train_model(data_path):
8     # Load training data

```

```

9     train_data = pd.read_csv(data_path)
10    X = train_data.drop('target', axis=1)
11    y = train_data['target']
12
13    # Train model
14    model = Ridge(alpha=1.0)
15    model.fit(X, y)
16
17    # Save model
18    with open('model.pkl', 'wb') as f:
19        pickle.dump(model, f)
20
21    return 'model.pkl'
22
23 if __name__ == '__main__':
24     import argparse
25     parser = argparse.ArgumentParser()
26     parser.add_argument('--data', type=str)
27     args = parser.parse_args()
28     train_model(args.data)

```

3. Model Evaluation

```

1  # evaluate.py
2  import pandas as pd
3  import pickle
4  from sklearn.metrics import mean_squared_error, r2_score
5  import json
6
7  def evaluate_model(model_path, test_data_path):
8      # Load model and test data
9      with open(model_path, 'rb') as f:
10         model = pickle.load(f)
11
12     test_data = pd.read_csv(test_data_path)
13     X_test = test_data.drop('target', axis=1)
14     y_test = test_data['target']
15
16     # Make predictions
17     y_pred = model.predict(X_test)
18
19     # Calculate metrics
20     metrics = {
21         'mse': mean_squared_error(y_test, y_pred),
22         'r2': r2_score(y_test, y_pred)
23     }
24
25     # Save metrics
26     with open('metrics.json', 'w') as f:
27         json.dump(metrics, f)
28
29     return 'metrics.json'
30
31 if __name__ == '__main__':
32     import argparse
33     parser = argparse.ArgumentParser()
34     parser.add_argument('--model', type=str)
35     parser.add_argument('--test-data', type=str)

```

```
36     args = parser.parse_args()
37     evaluate_model(args.model, args.test_data)
```

4. Kubeflow Pipeline Definition

```
1  # pipeline.py
2  from kfp import dsl
3  from kfp import components
4
5  def prepare_data_op():
6      return dsl.ContainerOp(
7          name='Prepare Data',
8          image='python:3.8',
9          command=['python', 'dataset_prep.py'],
10         file_outputs={
11             'train': '/train.csv',
12             'test': '/test.csv'
13         }
14     )
15
16  def train_model_op(train_data):
17      return dsl.ContainerOp(
18          name='Train Model',
19          image='python:3.8',
20          command=['python', 'train.py'],
21          arguments=['--data', train_data],
22          file_outputs={'model': '/model.pkl'}
23      )
24
25  def evaluate_model_op(model, test_data):
26      return dsl.ContainerOp(
27          name='Evaluate Model',
28          image='python:3.8',
29          command=['python', 'evaluate.py'],
30          arguments=[
31              '--model', model,
32              '--test-data', test_data
33          ],
34          file_outputs={'metrics': '/metrics.json'}
35      )
36
37  @dsl.pipeline(
38      name='Diabetes Prediction Pipeline',
39      description='Train and evaluate a diabetes prediction model'
40  )
41  def diabetes_pipeline():
42      data_prep = prepare_data_op()
43      train = train_model_op(data_prep.outputs['train'])
44      evaluate = evaluate_model_op(
45          train.outputs['model'],
46          data_prep.outputs['test']
47      )
48
49  # Compile pipeline
50  from kfp.compiler import Compiler
51  Compiler().compile(diabetes_pipeline, 'diabetes_pipeline.yaml')
```

5. Deploy Model to Azure ML

```

1 # deploy.py
2 from azureml.core import Workspace, Model, Environment
3 from azureml.core.webservice import AciWebservice
4 from azureml.core.model import InferenceConfig
5
6 def deploy_model(model_path):
7     # Initialize workspace
8     ws = Workspace.from_config()
9
10    # Register model
11    model = Model.register(
12        workspace=ws,
13        model_path=model_path,
14        model_name='diabetes-predictor'
15    )
16
17    # Create environment
18    env = Environment.from_conda_specification(
19        name='diabetes-env',
20        file_path='environment.yml'
21    )
22
23    # Define inference config
24    inference_config = InferenceConfig(
25        entry_script='score.py',
26        environment=env
27    )
28
29    # Define deployment config
30    deployment_config = AciWebservice.deploy_configuration(
31        cpu_cores=1,
32        memory_gb=1,
33        auth_enabled=True
34    )
35
36    # Deploy model
37    service = Model.deploy(
38        workspace=ws,
39        name='diabetes-service',
40        models=[model],
41        inference_config=inference_config,
42        deployment_config=deployment_config
43    )
44
45    service.wait_for_deployment(show_output=True)
46    return service.scoring_uri
47
48 if __name__ == '__main__':
49     import argparse
50     parser = argparse.ArgumentParser()
51     parser.add_argument('--model', type=str)
52     args = parser.parse_args()
53     deploy_model(args.model)

```

6. Scoring Script

```

1 # score.py
2 import json

```

```

3 import numpy as np
4 import pandas as pd
5 from azureml.core.model import Model
6 import pickle
7
8 def init():
9     global model
10    model_path = Model.get_model_path('diabetes-predictor')
11    with open(model_path, 'rb') as f:
12        model = pickle.load(f)
13
14 def run(raw_data):
15     try:
16         data = json.loads(raw_data)
17         data = pd.DataFrame(data)
18         prediction = model.predict(data)
19         return json.dumps({'prediction': prediction.tolist()})
20     except Exception as e:
21         return json.dumps({"error": str(e)})

```

Environment Configuration

Conda Environment

```

1 # environment.yml
2 name: diabetes-env
3 channels:
4   - conda-forge
5   - defaults
6 dependencies:
7   - python=3.8
8   - scikit-learn
9   - pandas
10  - numpy
11  - pip
12  - pip:
13    - azureml-defaults
14    - azureml-core

```

Usage Instructions

1. Set up infrastructure:

```
1 # Follow the setup steps above to create AKS cluster and deploy Kubeflow
```

2. Prepare environment:

```

1 # Create conda environment
2 conda env create -f environment.yml
3 conda activate diabetes-env

```

3. Run pipeline:

```

1 # Upload pipeline to Kubeflow
2 python pipeline.py
3
4 # Access Kubeflow UI
5 kubectl port-forward svc/istio-ingressgateway -n istio-system 8080:80

```

```
6 # Open browser at <http://localhost:8080>
```

4. Monitor deployment:

```
1 # Check pod status
2 kubectl get pods -n kubeflow
3
4 # Check service status
5 kubectl get services -n kubeflow
```

Monitoring

1. Azure Monitor Integration:

```
1 # Enable monitoring
2 from azureml.core import Workspace
3 from azureml.core.webservice import AciWebservice
4
5 ws = Workspace.from_config()
6 service = AciWebservice(ws, 'diabetes-service')
7
8 # Get logs
9 print(service.get_logs())
10
11 # Get metrics
12 print(service.get_metrics())
```

2. Kubeflow Dashboard:

- Access metrics and logs through the Kubeflow UI
- Monitor pipeline runs and component status
- Track model versions and deployments

This implementation provides a complete ML pipeline that:

- Uses the diabetes dataset as an example
- Implements data preparation, model training, and evaluation
- Deploys the model to Azure ML
- Includes monitoring capabilities

The code is production-ready but may need adjustments based on your specific requirements and scale.