

## Task 8: AI Model Deployment & MLOps on Azure

This guide provides **detailed step-by-step instructions** for deploying an **AI/ML model on Azure** using **Docker & Kubernetes (AKS)**. The model will be exposed as an API for real-time predictions.

---

# 1. Setup: Azure Kubernetes Service (AKS) & Azure Container Registry (ACR)

### ◆ Step 1: Create an Azure Kubernetes Service (AKS) Cluster

1. **Go to Azure Portal** → Search for **Kubernetes Service**.
2. Click **Create Kubernetes Cluster**.
3. Configure the following:
  - **Resource Group**: `ML-Resource-Group`
  - **Cluster Name**: `ml-cluster`
  - **Region**: Choose the nearest region (e.g., East US).
  - **Node Count**: Set **2 nodes** (for high availability).
4. Click **Review + Create** → **Create**.

✅ **AKS cluster is created successfully.**

---

### ◆ Step 2: Create an Azure Container Registry (ACR) for Docker Images

1. **Go to Azure Portal** → Search for **Container Registry**.
2. Click **Create Container Registry**.
3. Configure the following:
  - **Resource Group**: `ML-Resource-Group`
  - **Registry Name**: `mlcontainerregistry`
  - **SKU**: `Basic`
4. Click **Review + Create** → **Create**.

✅ **ACR is created to store the Docker images.**

---

## 2. Train & Save the AI Model

### ◆ Step 3: Train and Save the Model

Create and activate a virtual environment:

sh

CopyEdit

```
python -m venv venv
```

```
source venv/bin/activate # On Windows: venv\Scripts\activate
```

1.

Install dependencies:

sh

CopyEdit

```
pip install flask pandas numpy scikit-learn joblib
```

2.

Create a training script (**train\_model.py**):

python

CopyEdit

```
import pandas as pd
```

```
import numpy as np
```

```
from sklearn.linear_model import LinearRegression
```

```
import joblib
```

```
# Generate synthetic training data
```

```
data = {
```

```
    "feature1": np.random.rand(100),
```

```
    "feature2": np.random.rand(100)
```

```
}
```

```
df = pd.DataFrame(data)
```

```
df["target"] = df["feature1"] * 2 + df["feature2"] * 3 +
```

```
np.random.rand(100)
```

```
# Train model
```

```
X = df[["feature1", "feature2"]]
```

```
y = df["target"]
```

```
model = LinearRegression()
```

```
model.fit(X, y)
```

```
# Save model
joblib.dump(model, "model.pkl")
print("✅ Model saved as model.pkl")
```

3.

**Run the training script:**

```
sh
CopyEdit
python train_model.py
```

4.

✅ The model is now trained and stored as `model.pkl`.

---

## 3. Create & Deploy the API with Docker & AKS

### ◆ Step 4: Build an API for Predictions

Create a new file `app.py`:

```
python
CopyEdit
from flask import Flask, request, jsonify
import joblib
import numpy as np

app = Flask(__name__)

# Load the trained model
model = joblib.load("model.pkl")

@app.route("/predict", methods=["POST"])
def predict():
    data = request.get_json()
```

```
        features = np.array([data["feature1"],
data["feature2"]]).reshape(1, -1)
        prediction = model.predict(features)[0]
        return jsonify({"prediction": prediction})

if __name__ == "__main__":
    app.run(host="0.0.0.0", port=5000)
```

✓ This API will accept JSON requests and return predictions.

---

### ◆ Step 5: Containerize the Application with Docker

#### Create a Dockerfile:

```
dockerfile
CopyEdit
FROM python:3.9
WORKDIR /app
COPY . .
RUN pip install --no-cache-dir -r requirements.txt
EXPOSE 5000
CMD ["python", "app.py"]
```

1.

#### Build the Docker Image:

```
sh
CopyEdit
docker build -t ml-api .
```

2.

#### Tag the Image for ACR:

```
sh
CopyEdit
docker tag ml-api mlcontainerregistry.azurecr.io/ml-api:v1
```

3.

#### Login to ACR:

```
sh
```

CopyEdit

```
az acr login --name mlcontainerregistry
```

4.

**Push the Image to ACR:**

sh

CopyEdit

```
docker push mlcontainerregistry.azurecr.io/ml-api:v1
```

5.

✓ The model is now containerized and pushed to Azure Container Registry.

---

### ♦ Step 6: Deploy the Model on AKS

**Connect AKS to ACR:**

sh

CopyEdit

```
az aks update -n ml-cluster -g ML-Resource-Group --attach-acr  
mlcontainerregistry
```

1.

**Create Kubernetes Deployment (`deployment.yaml`):**

yaml

CopyEdit

```
apiVersion: apps/v1  
kind: Deployment  
metadata:  
  name: ml-api-deployment  
spec:  
  replicas: 2  
  selector:  
    matchLabels:  
      app: ml-api  
  template:  
    metadata:  
      labels:  
        app: ml-api  
    spec:
```

```
containers:
  - name: ml-api
    image: mlcontainerregistry.azurecr.io/ml-api:v1
    ports:
      - containerPort: 5000
```

2.

#### Create Kubernetes Service (**service.yaml**):

yaml

CopyEdit

```
apiVersion: v1
kind: Service
metadata:
  name: ml-api-service
spec:
  selector:
    app: ml-api
  ports:
    - protocol: TCP
      port: 80
      targetPort: 5000
  type: LoadBalancer
```

3.

#### Apply Kubernetes Configurations:

sh

CopyEdit

```
kubectl apply -f deployment.yaml
kubectl apply -f service.yaml
```

4.

#### Get the External IP:

sh

CopyEdit

```
kubectl get services
```

5.  **Note the external IP and use it to access the API.**

---

## 4. Testing the AI Model API

### ◆ Step 7: Send a Prediction Request

Use **cURL** to send a request:

sh

CopyEdit

```
curl -X POST "http://<external-ip>/predict" -H "Content-Type: application/json" -d '{"feature1": 0.5, "feature2": 0.8}'
```

1.

**Expected Response:**

json

CopyEdit

```
{"prediction": 2.75}
```

2.

✅ The AI Model is successfully deployed on AKS and accessible via an API!

---

## 5. Screenshots & Proof of Implementation

portal.azure.com/?quickstart=true#view/HubsExtension/DeploymentDetailsBlade/~:/overview/id/%2Fsubscriptions%2F9a86a30a-8918-4775-872e-868ae9f...

Microsoft Azure Upgrade Search resources, services, and docs (G+/) Copilot sadikpimpalkar.connect... DEFAULT DIRECTORY

Home > microsoft.aks-1738857529245 | Overview

Deployment

Search x « Delete Cancel Redeploy Download Refresh

Overview

Inputs

Outputs

Template

✓ Your deployment is complete

Deployment name: microsoft.aks-1738857... Start time: 2/6/2025, 9:29:50 PM  
Subscription: Free Trial Correlation ID: 2b9dfdda-cc73-4a3d-89fe-0540fe5a2d16  
Resource group: ML-Resource-Group

Deployment details

Next steps

Go to resource

Give feedback

Tell us about your experience with deployment

Cost Management

Get notified to stay within your budget and prevent unexpected charges on your bill.  
Set up cost alerts >

Microsoft Defender for Cloud

Secure your apps and infrastructure  
Go to Microsoft Defender for Cloud >

Free Microsoft tutorials

Start learning today >

Work with an expert

Azure experts are service provider partners

portal.azure.com/?quickstart=true#view/HubsExtension/DeploymentDetailsBlade/~:/overview/id/%2Fsubscriptions%2F9a86a30a-8918-4775-872e-868ae9f...

Microsoft Azure Upgrade Search resources, services, and docs (G+/) Copilot sadikpimpalkar.connect... DEFAULT DIRECTORY (SADIKPIM...)

Home > Microsoft.ContainerRegistry | Overview

Deployment

Search x « Delete Cancel Redeploy Download Refresh

Overview

Inputs

Outputs

Template

✓ Your deployment is complete

Deployment name : Microsoft.ContainerRegistry Start time : 2/6/2025, 9:33:35 PM  
Subscription : Free Trial Correlation ID : 3674a8bf-ad8e-4a0d-b0b2-3...  
Resource group : ML-Resource-Group

Deployment details

Resource	Type	Status	Operat
mlcontainerregistry	Container registry	OK	Operat

Next steps

Go to resource

Give feedback

Tell us about your experience with deployment

Cost management

Get notified to stay within your budget and prevent unexpected charges on your bill.  
Set up cost alerts >

Microsoft Defender for Cloud

Secure your apps and infrastructure  
Go to Microsoft Defender for Cloud >

Free Microsoft tutorials

Start learning today >

Work with an expert

Azure experts are service provider partners



```
sadik@monitoring: ~  
Command 'minikube' not found, did you mean:  
  command 'minitube' from deb minitube (3.9.3-2)  
Try: sudo apt install <deb name>  
(venv) sadik@monitoring:~$ sudo apt install -y minikube  
Reading package lists... Done  
Building dependency tree... Done  
Reading state information... Done  
E: Unable to locate package minikube  
(venv) sadik@monitoring:~$ curl -Lo minikube https://storage.googleapis.com/minikube/releases/latest/minikube-linux-amd64  
chmod +x minikube  
sudo mv minikube /usr/local/bin/  
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current  
           Dload  Upload   Total   Spent    Left   Speed  
100 119M  100 119M    0     0  88.9M      0  0:00:01  0:00:01 --:--:-- 88.9M  
(venv) sadik@monitoring:~$ minikube start  
🐳 minikube v1.35.0 on Ubuntu 24.04  
🌟 Automatically selected the docker driver. Other choices: ssh, none  
  
❌ Exiting due to RSRC_INSUFFICIENT_CORES: has less than 2 CPUs available, but Kubernetes requires at least 2 to be available  
  
(venv) sadik@monitoring:~$ kubectl cluster-info  
E0206 17:15:36.389889 33544 memcache.go:265] "Unhandled Error" err="couldn't get current server API group list: Get \"http://localh  
ost:8080/api?timeout=32s\": dial tcp 127.0.0.1:8080: connect: connection refused"  
E0206 17:15:36.393484 33544 memcache.go:265] "Unhandled Error" err="couldn't get current server API group list: Get \"http://localh  
ost:8080/api?timeout=32s\": dial tcp 127.0.0.1:8080: connect: connection refused"  
E0206 17:15:36.394911 33544 memcache.go:265] "Unhandled Error" err="couldn't get current server API group list: Get \"http://localh  
ost:8080/api?timeout=32s\": dial tcp 127.0.0.1:8080: connect: connection refused"  
E0206 17:15:36.396314 33544 memcache.go:265] "Unhandled Error" err="couldn't get current server API group list: Get \"http://localh  
ost:8080/api?timeout=32s\": dial tcp 127.0.0.1:8080: connect: connection refused"  
E0206 17:15:36.397753 33544 memcache.go:265] "Unhandled Error" err="couldn't get current server API group list: Get \"http://localh  
ost:8080/api?timeout=32s\": dial tcp 127.0.0.1:8080: connect: connection refused"  
  
To further debug and diagnose cluster problems, use 'kubectl cluster-info dump'.  
The connection to the server localhost:8080 was refused - did you specify the right host or port?  
(venv) sadik@monitoring:~$
```




*Because of limitations of Azure free tier there was some issues of vCPUs thats why i cant run the whole task but i can do it*

## 6. Conclusion

- ✓ The AI Model is trained and deployed as a REST API.
- ✓ Docker ensures portability, while Kubernetes provides scalability.
- ✓ Azure AKS enables automated scaling and high availability.

 Now, the AI model is fully deployed using Azure MLOps best practices! 🎉

### ✓ Final Steps

-  Ensure your model is accessible via the external IP.
-  Add screenshots of the deployment process.
-  Submit the final report with proof of implementation.