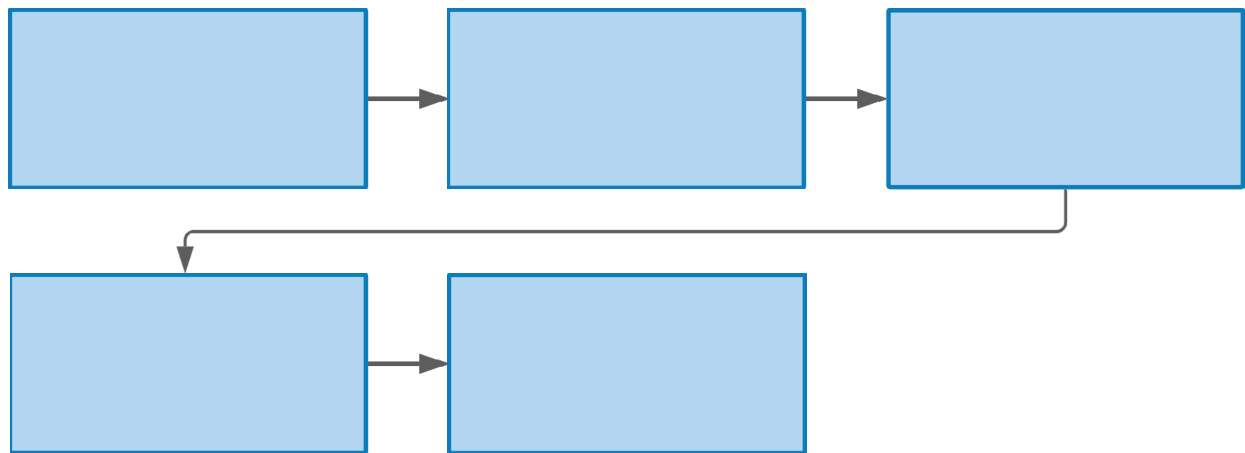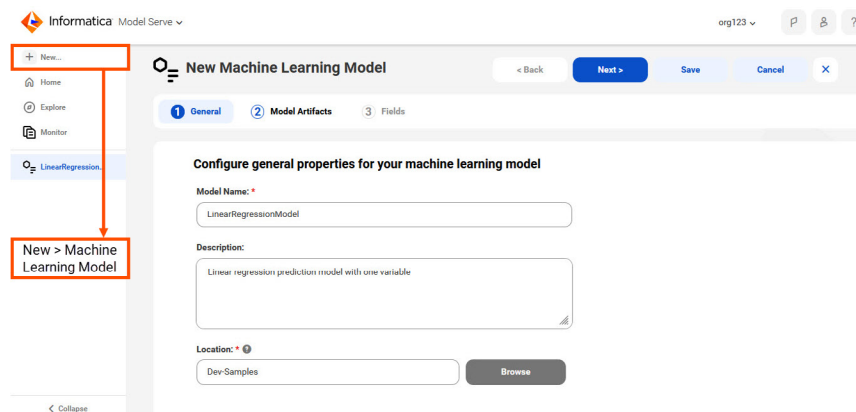# Model Serve Quick Start Guide

Model Serve provides a fully-managed platform to deploy and manage your data science and artificial intelligence models. You upload your machine learning model, and Model Serve deploys the model in a cloud environment. You can then make API requests to a URL endpoint to get predictions from the model.

Use this quick start guide to register and deploy a machine learning model:
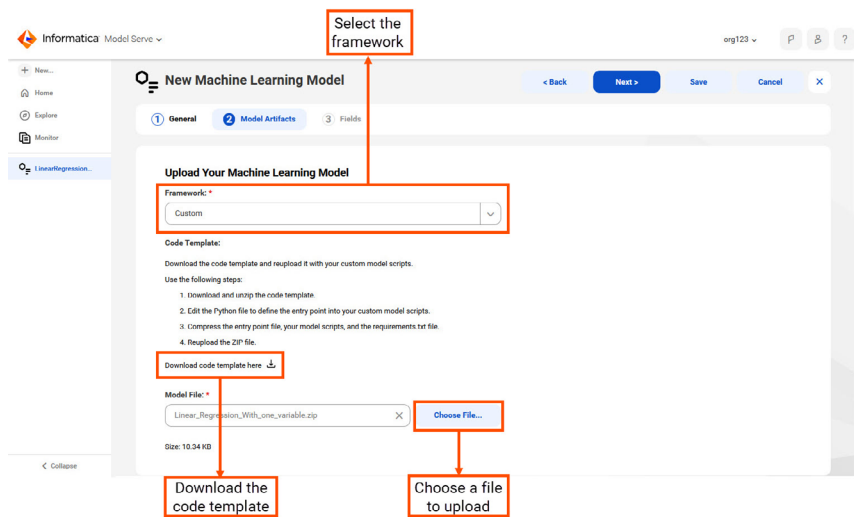


## Define a Machine Learning Model



1. In the navigation bar, select **New > Machine Learning Model** to create a machine learning model.

2. Configure the name, description, and location of the model.

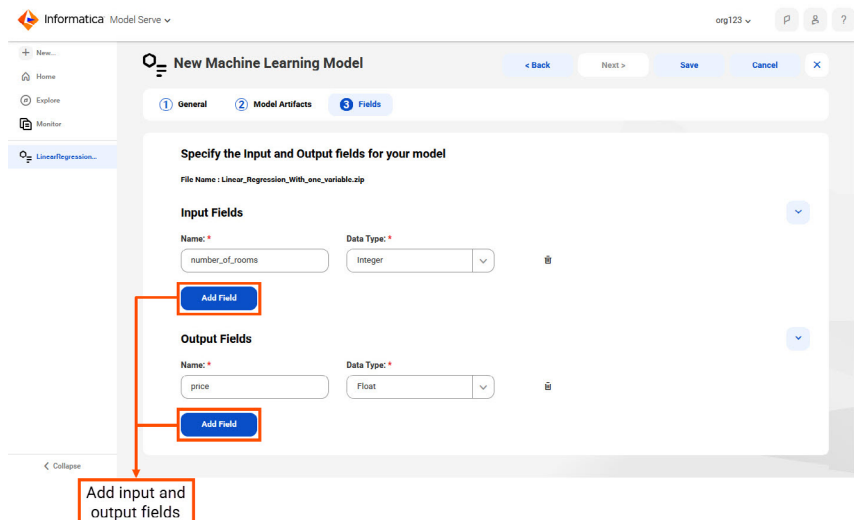Learn more about defining a machine learning model.

# Upload the Model Artifacts



1. Select the machine learning framework that you used to build the model.

2. If you're using a custom model, download and edit the code template to define the entry point to your custom files.

3. Upload a ZIP file containing the model artifacts, including the code template, if applicable.

Learn more about uploading the model artifacts.

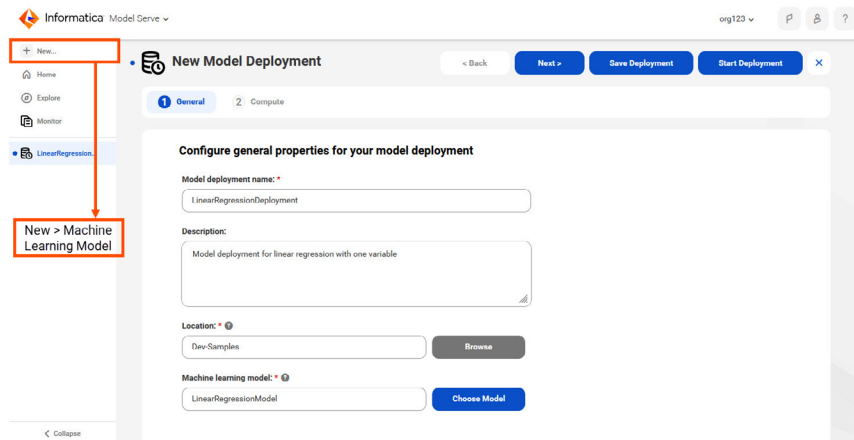# Define Input and Output Fields



Define input and output fields. Configure the name and data type of each input that the model expects and output that the model returns.

Then save the machine learning model.

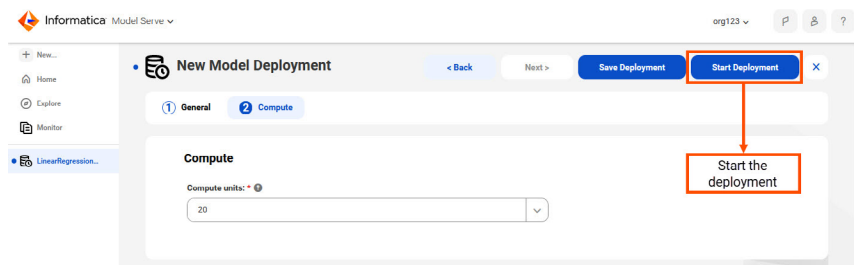Learn more about defining input and output fields.

# Configure a Model Deployment



1.  In the navigation bar, select **New > Model Deployment** to create a model deployment.

2.  On the General tab, configure the name, description, and location of the deployment. Also, select the model that you'll deploy.

3.  On the Compute tab, configure the maximum number of compute units that the model deployment can use.

Learn more about configuring a model deployment.

# Start the Deployment



After you save a valid model deployment, you can start the deployment. Model Serve provisions the cloud resources and deploys your model.

Learn more about starting a model deployment.

# Next Steps

1.  Send a sample request to test that the model is generating predictions as you expect.

2.  In your own application, add an API call that sends requests to the URL endpoint of the deployed model.

3.  Monitor the model deployment and stop it when it's not in use.