# Deployment

## What is Model Deployment?

Deployment is the method by which you integrate a [machine learning](#) model into an existing production environment to make practical business decisions based on data. It is one of the last stages in the [machine learning life cycle](#) and can be one of the most cumbersome. Often, an organization's IT systems are incompatible with traditional model-building languages, forcing data scientists and programmers to spend valuable time and brainpower rewriting them.

## Why is Model Deployment Important?

In order to start using a [model](#) for practical decision-making, it needs to be effectively deployed into production. If you cannot reliably get practical [insights](#) from your model, then the impact of the model is severely limited.
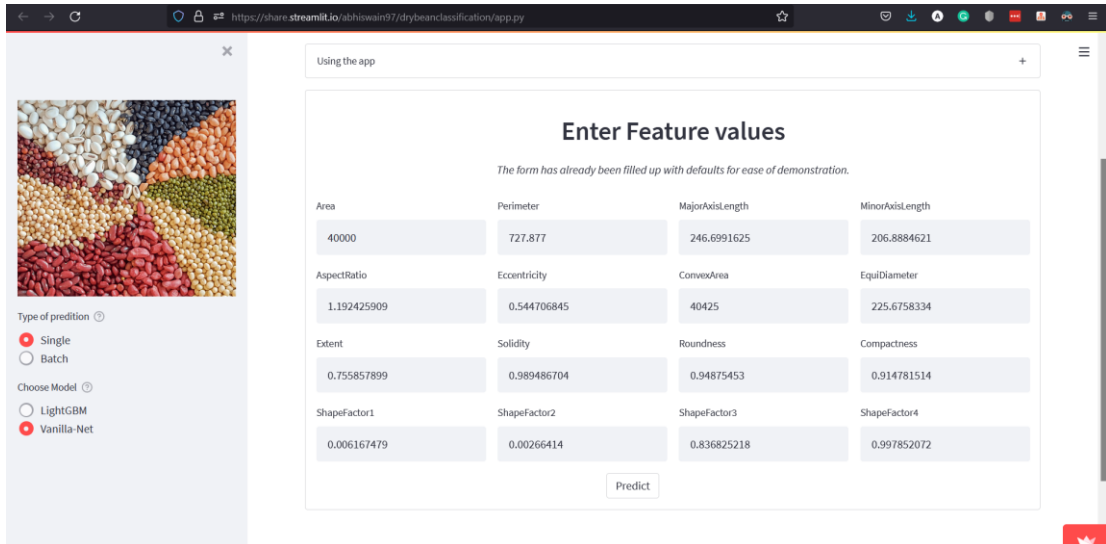
Model deployment is one of the most difficult processes of gaining value from machine learning. It requires coordination between data scientists, IT teams, software developers, and business professionals to ensure the model works reliably in the organization's production environment. This presents a major challenge because there is often a discrepancy between the programming language in which a machine learning model is written and the languages your production system can understand, and re-coding the model can extend the project timeline by weeks or months.

In order to get the most value out of machine learning models, it is important to seamlessly deploy them into production so a business can start using them to make practical decisions.

## Deploying our app

I have created a Streamlit app named "Dry bean classifier". It has two ways you can make classification.

1. Single prediction



The form has been pre-filled with default values for ease of demonstration.

2. Batch prediction using .csv file

You can also use a .csv file or a link to a csv file for doing prediction in bulk

Upload CSV directly

Paste a link to a CSV file



## Deploying the Tensorflow model

So, loading the TF model along with the app in Streamlit makes the app slower. So, to solve this I have the Vanilla Net model served at a docker container on Heroku using Tensorflow serving.

The logs of the deployed model look like this, you can see that version 5 of the model is deployed.



So, when you click on predict button the app is just making calls to this endpoint! The interesting thing about it is that all you need to do is just hit this endpoint for the prediction, so you can do it on your own without the app by a simple python script.