

Tensorflow Deployment.

Serving Tensorflow model with TF Serving.

① Prepare your model:

- ↳ Suppose you have an ArcFace model with the input tensors named
 - input_1. (None x 170 x 170 x 3). → Image dimensions.
 - input_2. (None x 93742). → Number of classes.

- ↳ The model needs to be saved in tf SavedModel format:

```
>>> arcface.load_weights("model.weights.hdf5")
```

```
>>> arcface.summary() # To see the input tensor names.
```

```
>>> tf.keras.models.save_model(arcface, "<abs_path>/model_name/<version#>")
```

↳ saves models
in tf SavedModel format.

① The model.

↳ ② Model path:

- + 1. Abs path to model directory
- + 2. Model dir name.
- + 3. Version number.

- ↳ Example:

```
* || -> tf.keras.models.save_model(arcface, "/home/hieu/arcface/1")
```

↳ /home/hieu

↳ /arcface

↳ /1

↳ /assets

↳ saved_model.pb

↳ /variables.

② Installing docker tf serving.

2.1. Installing Docker.

- ↳ Follow docs.docker.com/engine/install/ubuntu.

2.2. Installing Tf serving.

- ↳ * || \$ docker pull tensorflow/serving

③ Serving the model.

```
* || $ sudo docker run -t --rm -p {LOCAL PORT}:8501 \
-v "<abs_path>/model_name:/models/model_name" \
-e "MODEL_NAME=model_name" \
tensorflow/serving
```

- ↳ For example: model_name is arcface-keras.

↳ <abs_path>

↳ /arcface-keras

↳ /1

↳ /assets, /variables

↳ saved_model.pb

⇒ docker command:

```
$ docker -t --rm -p 8502:8501
      ↑
      run
      -v "<abs_path>/arcface_keras:/models/arcface_keras"
      -e "MODEL_NAME=arcface_keras"
      tensorflow/serving
```

↳ Now the model can be used for prediction at:

⊛ `http://localhost:8502/v1/models/arcface_keras:predict`

port # ver # model name.

④ Using the model.

↳ You can send data for prediction using python requests. Eg:

```
import cv2
import json
import time
import requests
import numpy as np
```

```
url = "http://localhost:8502/v1/models/arcface_keras:predict"
test_img = cv2.imread("2.jpg")
```

```
def make_request(model_url, img):
    img = cv2.resize(img, (170, 170))
    lbl = np.zeros(93742)
```

```
    data = json.dumps({
        "signature_name": "serving_default",
        ⊛ "instances": [
            ① Input tensors names. {
                "input-1": img.tolist(),
                "input-2": lbl.tolist()
            }
        ]
    })
```

② inputs have to be plain lists.

```
    3)
    ⊛ headers = {"content-type": "application/json"}
    ⊛ r = requests.post(model_url, data=data, headers=headers)
```

```
    predictions = np.array(json.loads(r.text)["predictions"])
    return predictions.
```