# step1:pytorch训练流程

安装库

```
1  pip install pytorch_quantization
```

train.py中导入库

```
1
2  from pytorch_quantization import nn as quant_nn
3  logging.set_verbosity(logging.FATAL)  # Disable logging as they are too noisy in
   notebook
4  from pytorch_quantization import quant_modules
5  quant_nn.TensorQuantizer.use_fb_fake_quant = True
6
7  quant_modules.initialize()
```

模型就正常训练

# step2:onnx导出

量化导出onnx文件，如下图



```python
def export_onnx(model, onnx_filename, batch_onnx):
    model.eval()
    opset_version = 13

    dummy_input = torch.randn(batch_onnx, 3, 192, 192, device='cuda') #TODO: switch
    input dims by model
    torch.onnx.export(model, dummy_input, onnx_filename, verbose=False,training=False,
    opset_version=opset_version, enable_onnx_checker=False, do_constant_folding=False)
```

```
7        return True
```

# step3:openvino推理调用

网络声明

```
def detect():
    # gpu = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
    torch.set_grad_enabled(False)
    palm_detector = BlazePalm().to(torch.device('cpu'))
    openvino_executor = openvino.runtime.Core()
    model = openvino_executor.compile_model(model = openvino_executor.read_model(model="quant_model.onnx"), device_name="CPU")
```

```
1    import openvino.runtime

2    openvino_executor = openvino.runtime.Core()

3    model = openvino_executor.compile_model(model =
     openvino_executor.read_model(model="quant_model.onnx"), device_name="CPU")
```

加入需要使用的工具函数

```
1  from typing import Union
2  def convert_any_to_numpy(
3      x: Union[torch.Tensor, np.ndarray, int, float, list, tuple],
4      accept_none: bool=True) -> np.ndarray:
5      if x is None and accept_none: return None
6      if x is None and not accept_none: raise ValueError('Trying to convert an empty
   value.')
7      if isinstance(x, np.ndarray): return x
8      elif isinstance(x, int) or isinstance(x, float): return np.array([x, ])
9      elif isinstance(x, torch.Tensor):
10         if x.numel() == 0 and accept_none: return None
11         if x.numel() == 0 and not accept_none: raise ValueError('Trying to convert an
   empty value.')
12         if x.numel() >= 1: return x.detach().cpu().numpy()
13     elif isinstance(x, list) or isinstance(x, tuple):
14         return np.array(x)
15     else:
16         raise TypeError(f'input value {x}({type(x)}) can not be converted as numpy
   type.')
```

替换原始的model，生成结果即可

```python
out_=model([convert_any_to_numpy(img)])
                for key in out_.items():
                    # print(key)
                    out_list.append(torch.tensor(key[-1]))
```

```python
img = img.to(torch.device('cpu'))
img = img.float() / 255.
out_list=[]
with torch.no_grad():
    # out = palm_detector(img)
    out_=model([convert_any_to_numpy(img)])
    for key in out_.items():
        # print(key)
        out_list.append(torch.tensor(key[-1]))
detections = _tensors_to_detections(out_list[1], out_list[0], anchors)
normalized_palm_detections = []
for i in range(len(detections)):
    faces = _weighted_non_max_suppression(detections[i])
```

用生成的onnx预测结果