

# Homework 1: Explainability

助教：莫易川

邮箱：mo666666@stu.pku.edu.cn





# **CONTENTS**

**01/ Brief Introduction**

**02/ Coding**

**03/ Paper reading**



# CONTENTS

---

## 01/ Brief Introduction

---

## Brief Introduction:

### Points

- Coding (20 points: 3+5+5+7)
- Paper reading (5 points: 3+1+1)

### Requirements

- Word/pdf is both ok.
- Write a report (at most **8** pages).
- Send your report and code to  
**trustworthy\_ai@163.com**  
Theme: Homework1-name-ID
- In Chinese/ English

**Due: 4/7 24:00** (+3 days delay is allowed)

### Language and wheel

- Python
- PyTorch

### Bonus

- The first five student who submits their homework will receive **1** extra bonus point



## Environment:



Not



and



## Reference:

Amaconda Installation:

[https://blog.csdn.net/qq\\_42257666/article/details/121383450](https://blog.csdn.net/qq_42257666/article/details/121383450)

The usage of Jupyter:

<https://zhuanlan.zhihu.com/p/33105153>

The documents of Pytorch:

<https://pytorch.org/docs/stable/index.html>

If you like, you could also use **Pycharm** and **Vscode** for coding



# CONTENTS

---

**01/ Brief Introduction**

**02/ Coding**

---

## 02

## Coding (Preparation)

## Package installation:

```

In 3 1 # Install other package if necessary.
      2 !pip install lime==0.1.1.37

```

Requirement already satisfied: scikit-learn>=0.18 in /home1/moyichuan/miniconda3/envs/skm/lib/python3.7/site-packages (from lime==0.1.1.37) (1.0.2)  
Requirement already satisfied: matplotlib in /home1/moyichuan/miniconda3/envs/skm/lib/python3.7/site-packages (from lime==0.1.1.37) (3.5.2)  
Requirement already satisfied: scipy in /home1/moyichuan/miniconda3/envs/skm/lib/python3.7/site-packages (from lime==0.1.1.37) (1.7.3)  
Requirement already satisfied: numpy in /home1/moyichuan/miniconda3/envs/skm/lib/python3.7/site-packages (from lime==0.1.1.37) (1.21.6)  
Requirement already satisfied: progressbar in /home1/moyichuan/miniconda3/envs/skm/lib/python3.7/site-packages (from lime==0.1.1.37) (2.5)

## Test whether jupyter works well:

```

1 # Test whether the jupyter works well.
2 np.random.seed(10)
3 img = Image.open('./test.JPEG')
4 img = img.convert('RGB')
5 transform = transforms.Compose([
6     transforms.Resize((224,224)),
7     transforms.ToTensor(),
8 ])
9 img = transform(img)
10 img_hwc = img.permute(1, 2, 0)
11 plt.imshow(img_hwc)
12 plt.axis("off")
13 plt.show()
14 img = img.unsqueeze(dim=0)

```

## Import package

```

In 4 1 # Import package
      2 import torch
      3 import json
      4 import cv2
      5 from PIL import Image
      6 import numpy as np
      7 import matplotlib.pyplot as plt
      8 from lime import lime_image
      9 from torchvision.datasets import ImageFolder
     10 import torch.nn as nn
     11 from torchvision.models.resnet import resnet50
     12 import torchvision.transforms as transforms
     13 from torch.utils.data import DataLoader
     14 from skimage.segmentation import slic
     15 from torch.autograd import Variable
     16 from typing import Callable, List, Tuple

```



## Load model

```
In 8 1  # Load the model
      2  class Normalize(nn.Module):
      3      def __init__(self, mean, std):
      4          super(Normalize, self).__init__()
      5          self.mean = torch.tensor(mean)
      6          self.std = torch.tensor(std)
      7
      8      def forward(self, x):
      9          return (x - self.mean.type_as(x)[None, :, None, None]) / self.std.type_as(x)[None, :, None, None]
10  imagenet_mean = (0.485, 0.456, 0.406)
11  imagenet_std = (0.229, 0.224, 0.225)
12  net = resnet50(num_classes=1000, pretrained=True)
13  model = nn.Sequential(Normalize(mean=imagenet_mean, std=imagenet_std), net)
14  model.eval()
```



# 02

## Coding (Preparation)

### Load dataset

```
In 9 1  # Load the dataset
      2  dataset = ImageFolder(root="./ImageNet_subset/", transform=transform)
      3  dataset_loader = DataLoader(dataset, batch_size=10, num_workers=6)
      4  for image, _ in dataset_loader:
      5      result = model(image)
      6      real_label = torch.argmax(result, dim = -1).numpy()
      7
      8
      9  # Print the label
     10  for j, label in enumerate(real_label):
     11      class_name = json.load(open("imagenet_class_index.json"))[str(label.item())]
     12      print("The label of image {:d} is :".format(j), class_name[1])
     13
     14
     15  # Show the image
     16  img_indices = [i for i in range(10)]
     17  all_image, _ = next(iter(dataset_loader))
     18  all_image = all_image.mul(255).add_(0.5).clamp_(0, 255).permute(0, 2, 3, 1).to('cpu', torch.uint8).numpy()
     19  fig, axs = plt.subplots(1, len(img_indices), figsize=(15, 8))
     20  for i, img in enumerate(all_image):
     21      axs[i].imshow(img)
     22      axs[i].set_xticks([])
     23      axs[i].set_yticks([])
```

## 02

## Coding (Line) 3point

```

In 10 1 def predict(input):
2     # input: numpy array, (batches, height, width, channels)
3     # output: the output of the model.
4     # pass
5     model.eval()
6     input = torch.FloatTensor(input).permute(0, 3, 1, 2)
7     output = model(input)
8     return output.detach().numpy()
9     #####
10    # write the code here
11    # return output.detach().numpy()
12    #####
13
14 def segmentation(input):
15     # split the image into 200 pieces with the help of segmentaion from skimage
16     # doc: https://scikit-image.org/docs/stable/api/skimage.segmentation.html#slic
17     return slic(input, n_segments=200, compactness=1, sigma=1, start_label=1)
18
19 img_indices = [i for i in range(10)]
20 fig, axs = plt.subplots(1, len(img_indices), figsize=(15, 8))
21
22
23 # fix the random seed to make it reproducible
24 for idx, (image, label) in enumerate(zip(all_image, real_label)):
25     x = (image/255).astype(np.double)
26

```

2 point

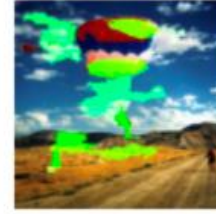
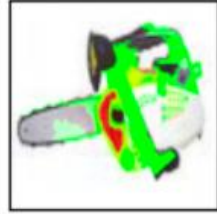
```

27 # Refer the doc: https://lime-ml.readthedocs.io/en/latest/lime.html?highlight=explain_instance#lime
28 .lime_image.LimeImageExplainer.explain_instance
29 #####
30 # write the code here
31 #####
32 explainer = lime_image.LimeImageExplainer()
33 explanation = explainer.explain_instance(image=x, classifier_fn=predict, segmentation_fn=segmentation)
34
35 # Turn the result from explainer to the image
36 # doc: https://lime-ml.readthedocs.io/en/latest/lime.html?highlight=get_image_and_mask#lime.lime_image
37 .ImageExplanation.get_image_and_mask
38 lime_img, mask = explanation.get_image_and_mask(label=label.item(), positive_only=False, hide_rest=False,
39 num_features=11, min_weight=0.05)
40 axs[idx].imshow(lime_img)
41 axs[idx].set_xticks([])
42 axs[idx].set_yticks([])
43
44 plt.xticks([])
45 plt.yticks([])
46 plt.axis('off')
47 plt.show()
48 plt.close()

```

1 point

Results:



## 02

## Coding (Saliency Map) 5 point

```
12 def normalize(image):
13     return (image - image.min()) / (image.max() - image.min())
14
15 def compute_saliency_maps(x, y, model):
16     # input: the input image, the ground truth label, the model
17     # output: the saliency maps of the images
18     # We need to normalize each image, because their gradients might vary i
19     pass
20     #####
21     # write the code here
22     # return saliencies
23     #####
24
25 image, _ = next(iter(dataset_loader))
26 result = model(image)
27 label = torch.argmax(result, dim = -1)
28
29
30
31 saliencies = compute_saliency_maps(image, label, model)
32 # visualize
33 fig, axs = plt.subplots(2, len(img_indices), figsize=(15, 8))
```

5 point

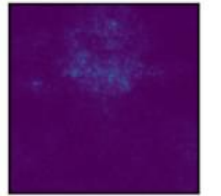
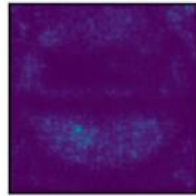
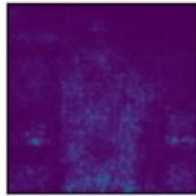
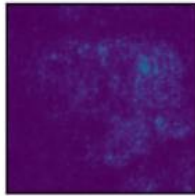
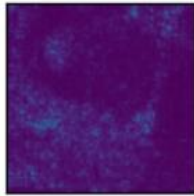
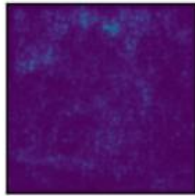
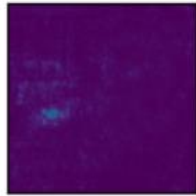
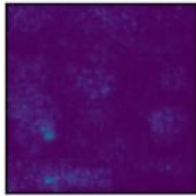
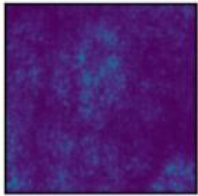
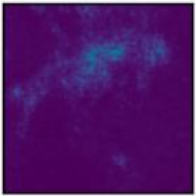
```
for row, target in enumerate([image, saliencies]):
    for column, img in enumerate(target):
        if row==0:
            axs[row][column].imshow(img.permute(1, 2, 0).detach().mul(255).add_(0.5).clamp_(0, 255).to('cpu',
            torch.uint8).numpy())
            axs[row][column].set_xticks([])
            axs[row][column].set_yticks([])
        else:
            axs[row][column].imshow(img.numpy())
            axs[row][column].set_xticks([])
            axs[row][column].set_yticks([])

plt.show()
plt.close()
```

# 02

## Coding (Saliency Map) 5 point

Results:



## 02

## Coding (Smooth-grad) 5 point

```

def normalize(image):
    return (image - image.min()) / (image.max() - image.min())

def smooth_grad(x, y, model, epoch, param_sigma_multiplier):
    # input: the input image, the ground truth label, the model,
    #        calculating std of noise
    # output: the saliency maps of the images
    pass
    #####
    # write the code here
    # return smooth
    #####

smooth = []
image, _ = next(iter(dataset_loader))
result = model(image)
label = torch.argmax(result, dim = -1)

for i, l in zip(image, label):
    smooth.append(smooth_grad(i, l, model, 10, 0.4))
smooth = np.stack(smooth)

```

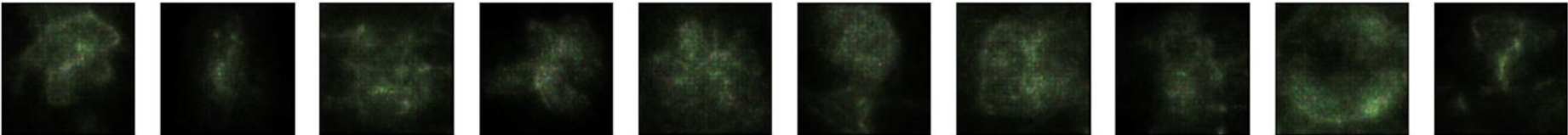
5 point

```

27 fig, axs = plt.subplots(2, len(img_indices), figsize=(15, 8))
28 for row, target in enumerate([image, smooth]):
29     for column, img in enumerate(target):
30         if row==0:
31             axs[row][column].imshow(img.permute(1, 2, 0).detach().mul(255).add_(0.5).clamp_(0, 255).to('cpu',
32                                     torch.uint8).numpy())
33             axs[row][column].set_xticks([])
34             axs[row][column].set_yticks([])
35         else:
36             axs[row][column].imshow(np.transpose(img.reshape(3,224,224), (1,2,0)))
37             axs[row][column].set_xticks([])
38             axs[row][column].set_yticks([])

```

## Results:





## 02

## Coding (grad-cam) 7 point

```

In _ 1 def scale_cam_image(cam, target_size=None):
2     result = []
3     for img in cam:
4         img = img - np.min(img)
5         img = img / (1e-7 + np.max(img))
6         if target_size is not None:
7             img = cv2.resize(img, target_size)
8         result.append(img)
9     result = np.float32(result)
10    return result
11
12 class ActivationsAndGradients:
13     """ Class for extracting activations and
14     registering gradients from targeted intermediate layers """
15
16     def __init__(self, model, target_layers):
17         self.model = model
18         self.gradients = []
19         self.activations = []
20         self.handles = []
21         for target_layer in target_layers:
22             self.handles.append(
23                 target_layer.register_forward_hook(self.save_activation))
24             self.handles.append(
25                 target_layer.register_forward_hook(self.save_gradient))

```

```

27     def save_activation(self, module, input, output):
28         activation = output
29         self.activations.append(activation.cpu().detach())
30
31     def save_gradient(self, module, input, output):
32         if not hasattr(output, "requires_grad") or not output.requires_grad:
33             # You can only register hooks on tensor requires grad.
34             return
35         # Gradients are computed in reverse order
36         def _store_grad(grad):
37             self.gradients = [grad.cpu().detach()] + self.gradients
38         output.register_hook(_store_grad)
39
40     def __call__(self, x):
41         self.gradients = []
42         self.activations = []
43         return self.model(x)
44
45     def release(self):
46         for handle in self.handles:
47             handle.remove()

```



## 02

## Coding (grad-cam) 7 point

```

def show_cam_on_image(img: np.ndarray,
                      mask: np.ndarray,
                      use_rgb: bool = False,
                      colormap: int = cv2.COLORMAP_JET,
                      image_weight: float = 0.5) -> np.ndarray:
    """ This function overlays the cam mask on the image as an heatmap.
    By default the heatmap is in BGR format.
    :param img: The base image in RGB or BGR format.
    :param mask: The cam mask.
    :param use_rgb: Whether to use an RGB or BGR heatmap, this should be set to True if 'img' is in RGB format.
    :param colormap: The OpenCV colormap to be used.
    :param image_weight: The final result is image_weight * img + (1-image_weight) * mask.
    :returns: The default image with the cam overlay.
    """
    heatmap = cv2.applyColorMap(np.uint8(255 * mask), colormap)
    if use_rgb:
        heatmap = cv2.cvtColor(heatmap, cv2.COLOR_BGR2RGB)
    heatmap = np.float32(heatmap) / 255

    if np.max(img) > 1:
        raise Exception(
            "The input image should np.float32 in the range [0, 1]")

    if image_weight < 0 or image_weight > 1:
        raise Exception(
            f"image_weight should be in the range [0, 1].\n"
            f"Got: {image_weight}")

    cam = (1 - image_weight) * heatmap + image_weight * img
    cam = cam / np.max(cam)
    return np.uint8(255 * cam)

```

```

85 class ClassifierOutputTarget:
86     def __init__(self, category):
87         self.category = category
88
89     def __call__(self, model_output):
90         if len(model_output.shape) == 1:
91             return model_output[self.category]
92         return model_output[:, self.category]
93
94
95 class GradCAM:
96     def __init__(self,
97                  model: torch.nn.Module,
98                  target_layers: List[torch.nn.Module]) -> None:
99         self.model = model.eval()
100        self.target_layers = target_layers
101        self.model = model
102        self.activations_and_grads = ActivationsAndGradients(
103            self.model, target_layers)
104
105        """ Get a vector of weights for every channel in the target layer.
106        Methods that return weights channels,
107        will typically need to only implement this function. """
108
109    def get_cam_weights(self,
110                       grads: torch.Tensor) -> np.ndarray:
111        return np.mean(grads, axis=(2, 3))
112

```

## 02

## Coding (grad-cam) 7 point

```

def show_cam_on_image(img: np.ndarray,
                      mask: np.ndarray,
                      use_rgb: bool = False,
                      colormap: int = cv2.COLORMAP_JET,
                      image_weight: float = 0.5) -> np.ndarray:
    """ This function overlays the cam mask on the image as an heatmap.
    By default the heatmap is in BGR format.
    :param img: The base image in RGB or BGR format.
    :param mask: The cam mask.
    :param use_rgb: Whether to use an RGB or BGR heatmap, this should be set to True if 'img' is in RGB format.
    :param colormap: The OpenCV colormap to be used.
    :param image_weight: The final result is image_weight * img + (1-image_weight) * mask.
    :returns: The default image with the cam overlay.
    """
    heatmap = cv2.applyColorMap(np.uint8(255 * mask), colormap)
    if use_rgb:
        heatmap = cv2.cvtColor(heatmap, cv2.COLOR_BGR2RGB)
    heatmap = np.float32(heatmap) / 255

    if np.max(img) > 1:
        raise Exception(
            "The input image should np.float32 in the range [0, 1]")

    if image_weight < 0 or image_weight > 1:
        raise Exception(
            f"image_weight should be in the range [0, 1].\n"
            f"Got: {image_weight}")

    cam = (1 - image_weight) * heatmap + image_weight * img
    cam = cam / np.max(cam)
    return np.uint8(255 * cam)

```

```

85 class ClassifierOutputTarget:
86     def __init__(self, category):
87         self.category = category
88
89     def __call__(self, model_output):
90         if len(model_output.shape) == 1:
91             return model_output[self.category]
92         return model_output[:, self.category]
93
94
95 class GradCAM:
96     def __init__(self,
97                  model: torch.nn.Module,
98                  target_layers: List[torch.nn.Module]) -> None:
99         self.model = model.eval()
100        self.target_layers = target_layers
101        self.model = model
102        self.activations_and_grads = ActivationsAndGradients(
103            self.model, target_layers)
104
105        """ Get a vector of weights for every channel in the target layer.
106        Methods that return weights channels,
107        will typically need to only implement this function. """
108
109    def get_cam_weights(self,
110                       grads: torch.Tensor) -> np.ndarray:
111        return np.mean(grads, axis=(2, 3))
112

```

## 02

## Coding (grad-cam) 7 point

```

113 def get_cam_image(self,
114     activations: torch.Tensor,
115     grads: torch.Tensor) -> np.ndarray:
116     # input: the activation, the gradient(4D tensor)
117     # output: cam of a specific layer
118     pass
119     #####
120     # write the code here
121     # return cam
122     #####
123
124 def forward(self,
125     input_tensor: torch.Tensor,
126     targets: List[torch.nn.Module]) -> np.ndarray:
127     outputs = self.activations_and_grads(input_tensor)
128     self.model.zero_grad()
129     loss = sum([target(output)
130         for target, output in zip(targets, outputs)])
131     loss.backward(retain_graph=True)
132     # In most of the saliency attribution papers, the saliency is
133     # computed with a single target layer.
134     # Commonly it is the last convolutional layer.
135     # Here we support passing a list with multiple target layers.
136     # It will compute the saliency image for every image,
137     # and then aggregate them (with a default mean aggregation).
138     # This gives you more flexibility in case you just want to
139     # use all conv layers for example, all Batchnorm layers,
140     # or something else.
141     cam_per_layer = self.compute_cam_per_layer(input_tensor,
142         targets)
143     return self.aggregate_multi_layers(cam_per_layer)

```

3 point

```

146 def compute_cam_per_layer(
147     self,
148     input_tensor: torch.Tensor) -> np.ndarray:
149     activations_list = [a.cpu().data.numpy()
150         for a in self.activations_and_grads.activations]
151     grads_list = [g.cpu().data.numpy()
152         for g in self.activations_and_grads.gradients]
153     target_size = (input_tensor.size(-1), input_tensor.size(-2))
154     cam_per_target_layer = []
155     # Loop over the saliency image from every layer
156     for i in range(len(self.target_layers)):
157         target_layer = self.target_layers[i]
158         layer_activations = None
159         layer_grads = None
160         if i < len(activations_list):
161             layer_activations = activations_list[i]
162         if i < len(grads_list):
163             layer_grads = grads_list[i]
164         #####
165         # write the code here
166         #####
167         return cam_per_target_layer
168
169 def aggregate_multi_layers(
170     self,
171     cam_per_target_layer: np.ndarray) -> np.ndarray:
172     cam_per_target_layer = np.concatenate(cam_per_target_layer, axis=1)
173     cam_per_target_layer = np.maximum(cam_per_target_layer, 0)
174     result = np.mean(cam_per_target_layer, axis=1)
175     return scale_cam_image(result)

```

4 point

## 02

## Coding (grad-cam) 7 point

```

def __call__(self,
              input_tensor: torch.Tensor,
              targets: List[torch.nn.Module] = None) -> np.ndarray:
    return self.forward(input_tensor,
                        targets)

def __del__(self):
    self.activations_and_grads.release()

def __enter__(self):
    return self

def __exit__(self, exc_type, exc_value, exc_tb):
    self.activations_and_grads.release()
    if isinstance(exc_value, IndexError):
        # Handle IndexError here...
        print(
            f"An exception occurred in CAM with block: {exc_type}. Message: {exc_value}")
    return True

```

```

model = resnet50(pretrained=True)
label = torch.argmax(result, dim = -1).numpy()
target_layers = [model.layer4[-1]]
image, _ = next(iter(dataset_loader))
cam = GradCAM(model=model, target_layers=target_layers)
targets = [ClassifierOutputTarget(i) for i in label]
grayscale_cam = cam(input_tensor=image, targets=targets)
fig, axs = plt.subplots(2, len(img_indices), figsize=(15, 8))
for column, single_image in enumerate(image):
    axs[0][column].imshow(single_image.permute(1, 2, 0).detach().mul(255).add_(0.5).clamp_(0, 255).to('cpu',
        torch.uint8).numpy())
    axs[0][column].set_xticks([])
    axs[0][column].set_yticks([])
    axs[1][column].imshow(show_cam_on_image(single_image.permute(1, 2, 0).detach().to('cpu').numpy(),
        grayscale_cam[column], use_rgb=True))
    axs[1][column].set_xticks([])
    axs[1][column].set_yticks([])

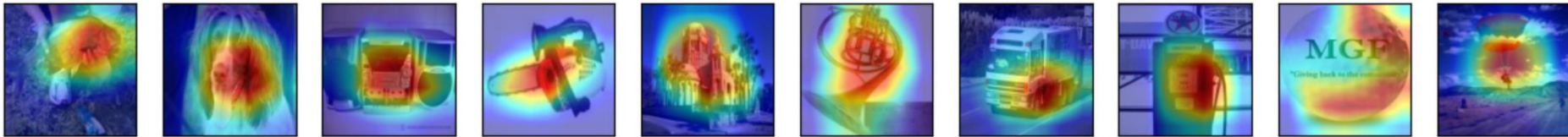
```

## 02

Coding (grad-cam) 5point

---

Results:





# CONTENTS

**01/ Brief Introduction**

**02/ Coding**

**03/ Paper reading**



# 03

## Paper reading

---

### Requirement:

- The brief summary(<120 words) **1point**
- The strength(<120 words) **1point**
- The weakness(<120 words) **1point**
- Quality **1point**
- Quantity ( >5 papers) **1point**

### Reference:

<https://openreview.net/>

You can get **4** points if you just read one paper!

# 论文列表

- [1] Learning Deep Features for Discriminative Localization
- [2] Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization
- [3] Grad-CAM++: Improved Visual Explanations for Deep Convolutional Networks
- [4] Smooth Grad-CAM++: An Enhanced Inference Level Visualization Technique for Deep Convolutional Neural Network Models
- [5] Score-CAM: Score-Weighted Visual Explanations for Convolutional Neural Networks
- [6] SS-CAM: Smoothed Score-CAM for Sharper Visual Feature Localization
- [7] IS-CAM: Integrated Score-CAM for axiomatic-based explanations
- [8] Axiom-based Grad-CAM: Towards Accurate Visualization and Explanation of CNNs



- [9] SmoothGrad: removing noise by adding noise.
- [10] Quantifying Attention Flow in Transformers
- [11] Transformer Interpretability Beyond Attention Visualization
- [12] Explaining deep neural networks and beyond: A review of methods and applications
- [13] Machine learning interpretability: A survey on methods and metrics
- [14] Towards explainable artificial intelligence
- [15] Explaining explanations: An overview of interpretability of machine learning

# Q&A

# Thanks