

EVERYTHING IS

X A N T L

THIS PRESENTATION IS ABOUT

- taxonomy
- model (un)certainty
- gradCAMs
- the Nigerian prince

THIS PRESENTATION IS ABOUT

- taxonomy
- model (un)certainty
- gradCAMs
- the Nigerian prince

thoughts about the presentation this morning?

TAXONOMY OF XAI

→ **by default**

1. linear models
2. tree-based models

→ **black boxes**

- global vs local methods
- model-specific vs model-agnostic

Interpretable Machine Learning

A Guide for Making
Black Box Models Explainable



Second
Edition

TAXONOMY OF XAI

→ **by default**

1. linear models
2. tree-based models

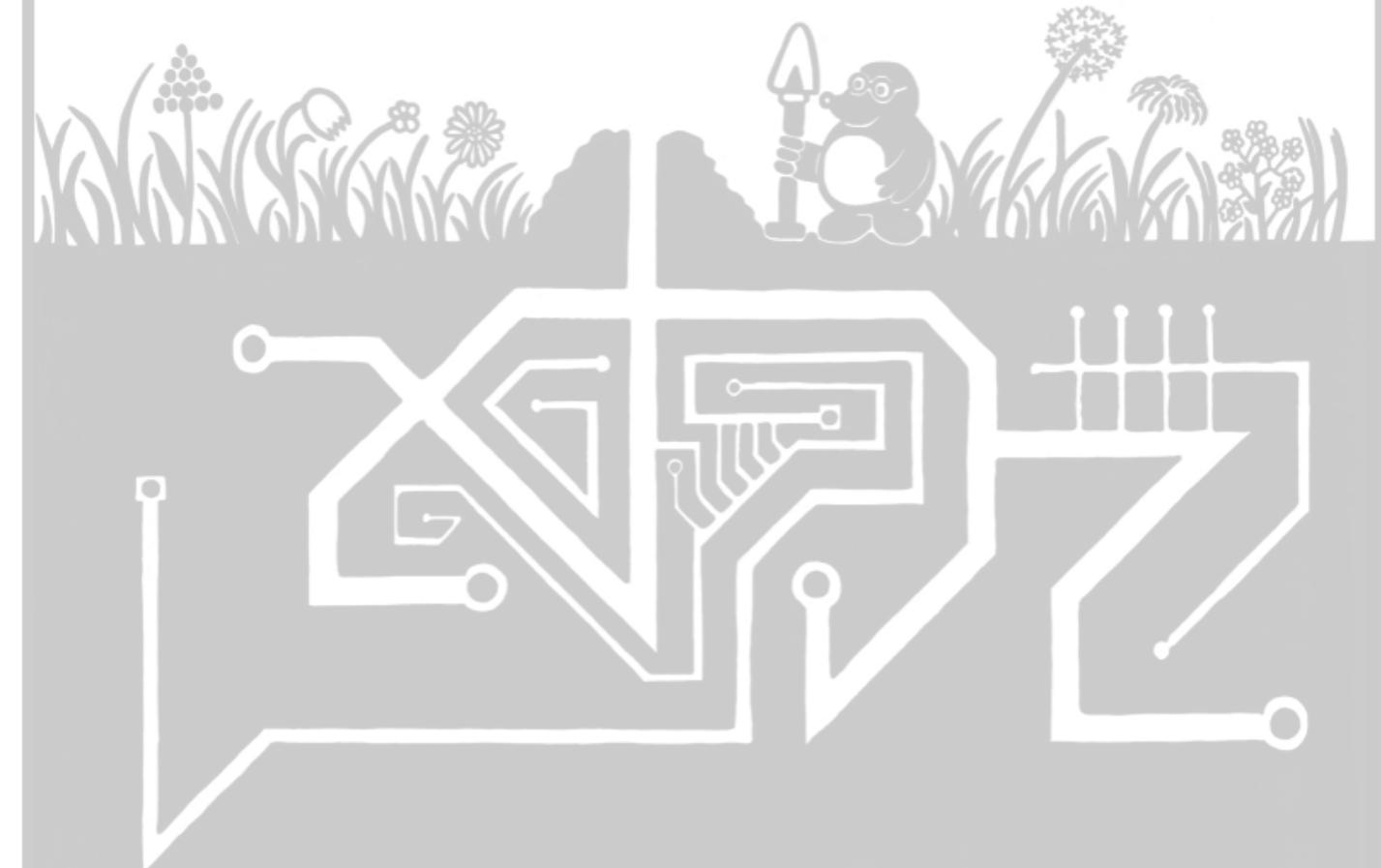
→ **black boxes**

- global vs local methods
- model-specific vs model-agnostic

classify gradCAM according to the taxonomy

Interpretable Machine Learning

A Guide for Making
Black Box Models Explainable



Christoph Molnar

Second
Edition

TAXONOMY OF XAI

→ **by default**

1. linear models
2. tree-based models

→ **black boxes**

- global vs **local** methods
- **model-specific** vs model-agnostic

gradCAM is a local, model-specific method to explain the predictions of a CNN

Interpretable Machine Learning

A Guide for Making
Black Box Models Explainable



Christoph Molnar

Second
Edition

MODEL (UN)CERTAINTY

“It is remarkable that science, which originated in the consideration of games of chance, should have become the most important object of human knowledge.”

— Pierre Simon Laplace

MODEL (UN)CERTAINTY I

aleatoric uncertainty

statistical in nature, refers to
random variation

$$\text{🔔 } P(Y|X; \theta)$$

MODEL (UN)CERTAINTY I

aleatoric uncertainty

statistical in nature, refers to random variation

$$\text{🔔 } P(Y|X; \theta)$$

more data is not a solution

data quality, model averaging, and prediction intervals help.

MODEL (UN)CERTAINTY II

epistemic
uncertainty

human in nature, refers to
ignorance

$$\text{🔔 } P(Y|X; \theta)$$

MODEL (UN)CERTAINTY II

**epistemic
uncertainty**

human in nature, refers to ignorance

$$\text{🔔 } P(Y|X; \theta)$$

more data is a solution

hyperparameter tuning, regularization etc. help

HOW IS THIS USEFUL?

if we train a classifier on digits [0,8], what happens when we run..?



HOW IS THIS USEFUL?

if we train a classifier on digits [0,8], what happens when we run..?
→ model.predict([2])



HOW IS THIS USEFUL?

if we train a classifier on digits [0,8], what happens when we run..?

- model.predict([2])
- model.predict([7])



HOW IS THIS USEFUL?

if we train a classifier on digits [0,8], what happens when we run..?

- model.predict([2])
- model.predict([7])
- model.predict([9])



HOW IS THIS USEFUL?

if we train a classifier on digits [0,8], what happens when we run..?

- model.predict([2])
- model.predict([7])
- model.predict([9])

! never forget that the models you build are probabilistic in nature





Limitations

🎵 MIND YOUR STEP

**even the best models
can be wrong, and with
horrible consequences**

May occasionally generate
incorrect information

May occasionally produce
harmful instructions or biased
content

Limited knowledge of world and
events after 2021

🎵 MIND YOUR STEP

**even the best
models can be
wrong, and with
horrible
consequences**

remember the end user!



Limitations

May occasionally generate
incorrect information

May occasionally produce
harmful instructions or biased
content

Limited knowledge of world and
events after 2021

EMBRACE THE UNCERTAINTY

**Successful decisions are built
on**

- minimizing our ignorance
- accepting inherent randomness
- knowing the difference between the two.

think self-driving cars, medical diagnosis, or exploding rockets





may occasionally kill passengers

EXPLAINING DEEP NEURAL NETWORKS

EXPLAINING DEEP NEURAL NETWORKS

**feature
visualization**

*visualize the
activations of a
deep neural
network*

**feature
attribution**

**adversarial
examples**

EXPLAINING DEEP NEURAL NETWORKS

feature visualization

*visualize the
activations of a
deep neural
network*

feature attribution

*visualize the
features that
contribute strongly
to the activations
of CNN layers for a
given input image
and class label.*

adversarial examples

EXPLAINING DEEP NEURAL NETWORKS

feature visualization

visualize the activations of a deep neural network

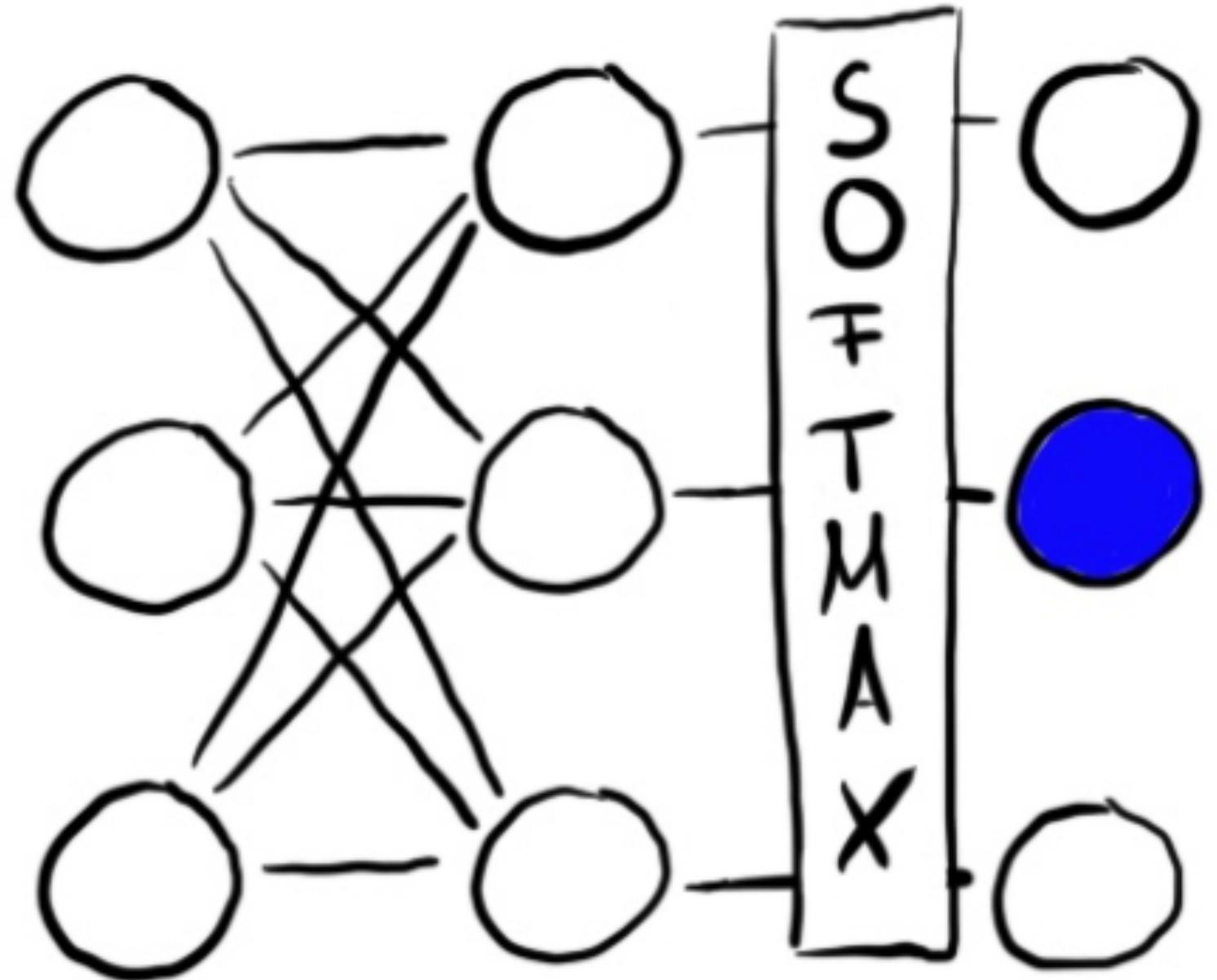
feature attribution

visualize the features that contribute strongly to the activations of CNN layers for a given input image and class label.

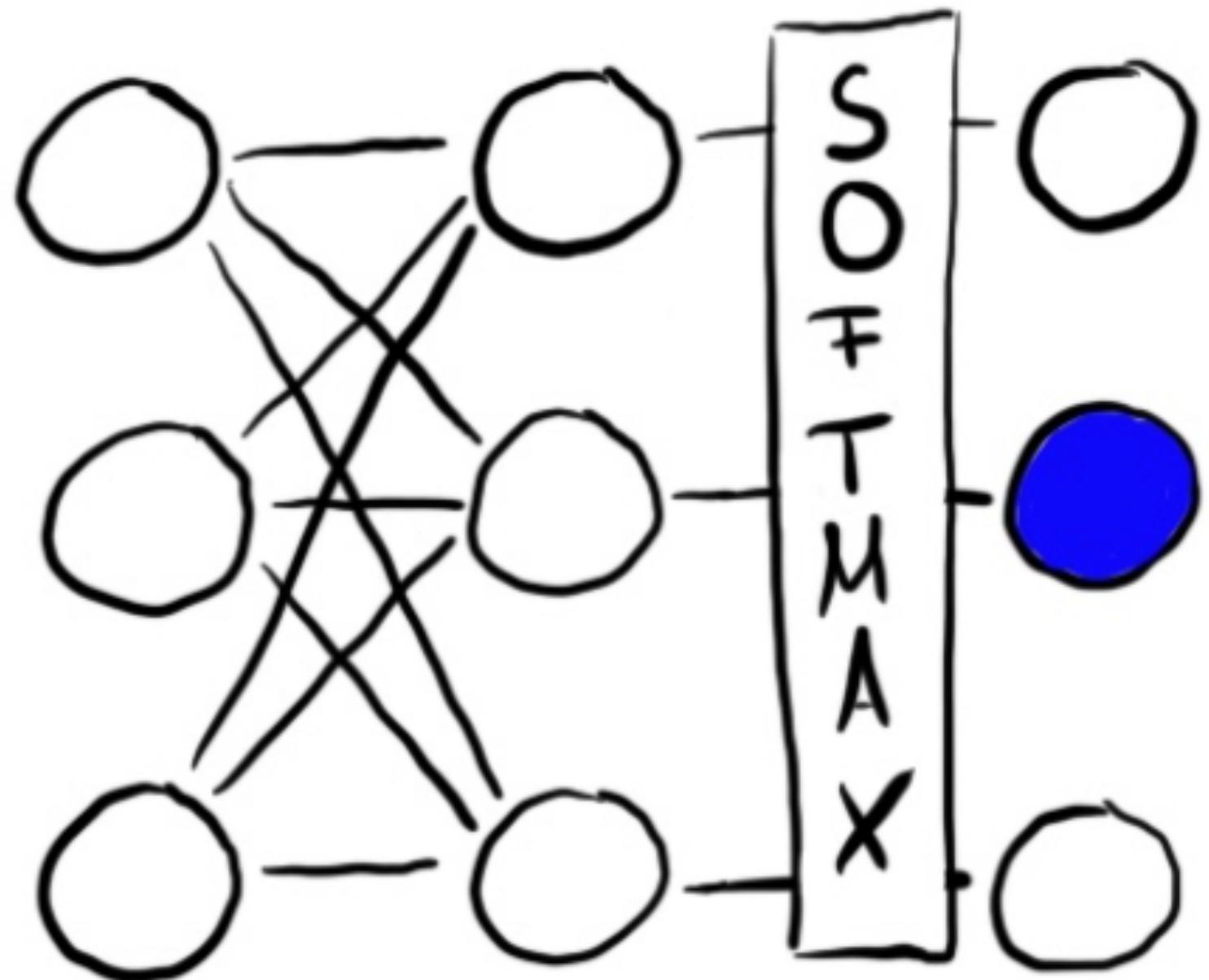
adversarial examples

understand the robustness of a deep neural network by generating adversarial examples giving insight into the decision boundaries of the model.

GRAD-CAM

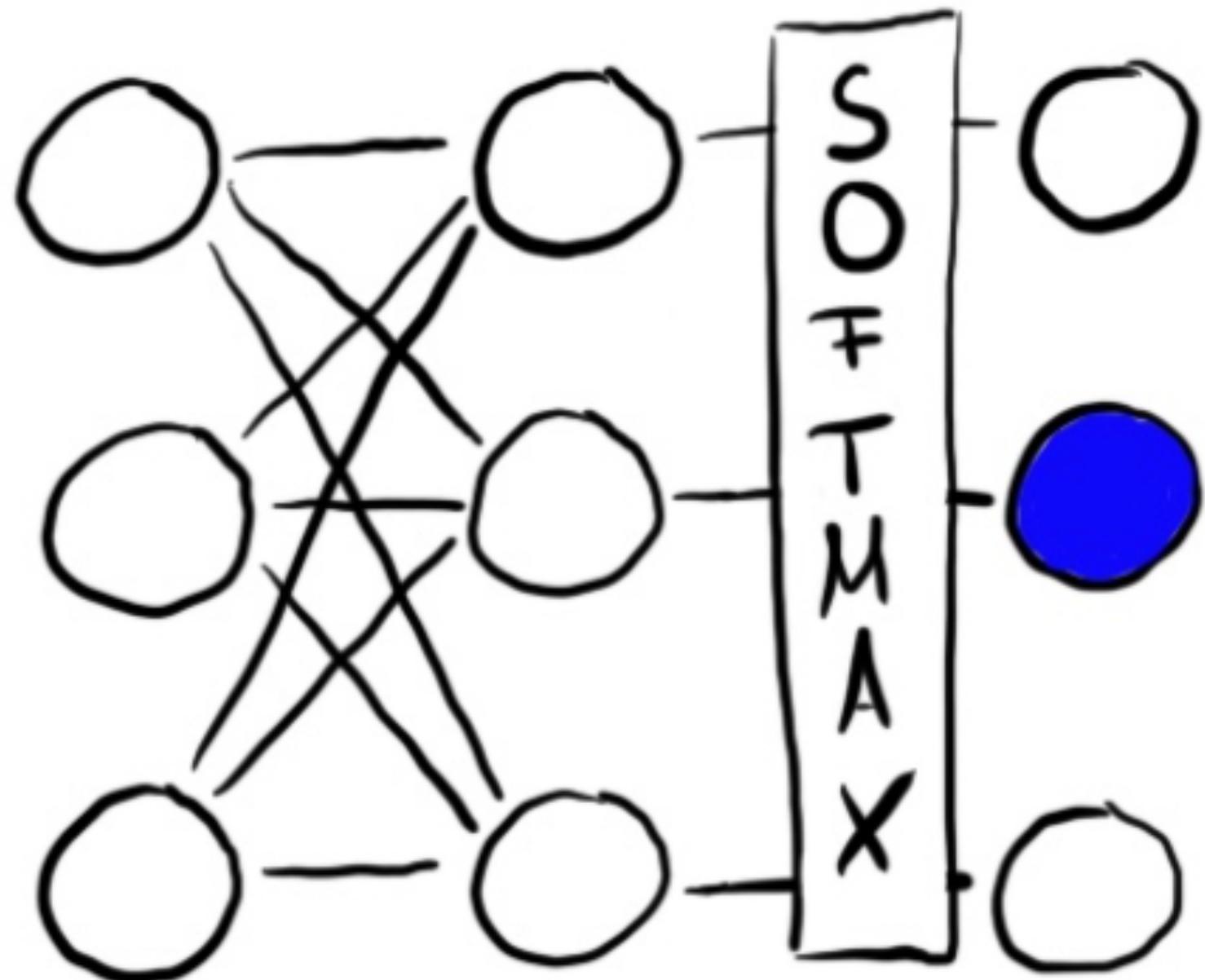


GRAD-CAM



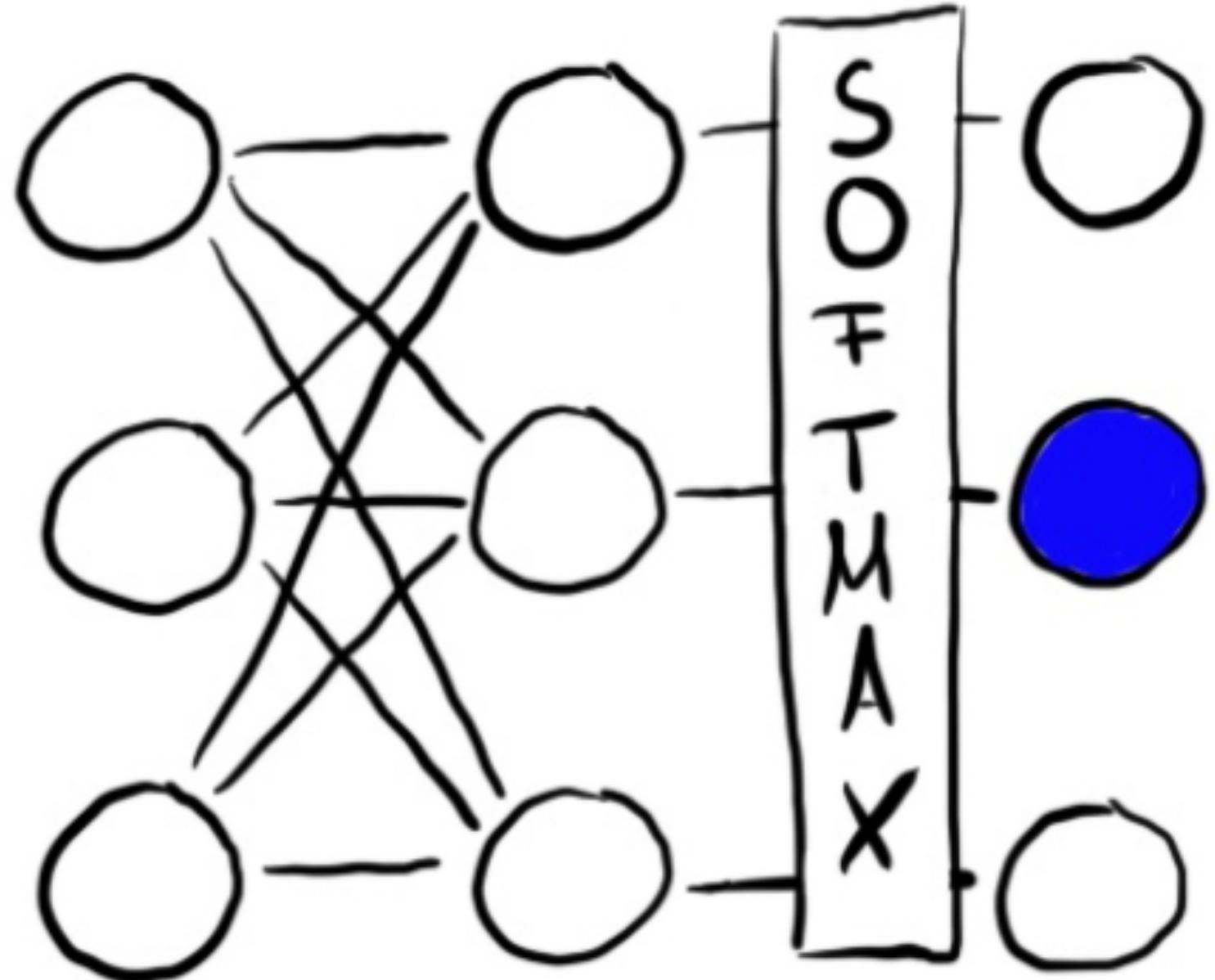
→ grad-CAM is a technique used in deep learning to highlight which parts of an image were important in predicting its classification.

GRAD-CAM



- grad-CAM is a technique used in deep learning to highlight which parts of an image were important in predicting its classification.
- grad-CAM works by computing the gradients of the output class score with respect to the feature maps of the last convolutional layer in the CNN.

GRAD-CAM



- grad-CAM is a technique used in deep learning to highlight which parts of an image were important in predicting its classification.
- grad-CAM works by computing the gradients of the output class score with respect to the feature maps of the last convolutional layer in the CNN.
- grad-CAM works by computing the gradients of the output class score with respect to the feature maps of the last convolutional layer in the CNN. These gradients are then used to compute a weight map for each feature map, which are then multiplied together and summed to produce the final heatmap¹.

¹Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization. In Proceedings of the IEEE international conference on computer vision (pp. 618-626).

TF EXPLAIN

```
!pip install tf_explain
!pip install opencv-python

#load libraries
import numpy as np
import tensorflow as tf
import PIL

#load GradCAM
from tf_explain.core.grad_cam import GradCAM

IMAGE_PATH = "./assets/images/cat.jpg"
class_index = 281

img = tf.keras.preprocessing.image.load_img(IMAGE_PATH, target_size=(224, 224))
img = tf.keras.preprocessing.image.img_to_array(img)

model = tf.keras.applications.VGG16(weights="imagenet", include_top=True)

#get model summary
model.summary()

#first create the input in a format that the explainer expects (a tuple)
input_img = (np.array([img]), None)

#initialize the explainer as an instance of the GradCAM object
explainer = GradCAM()

# Obtain explanations for your image using VGG 16 and GradCAM
grid = explainer.explain(input_img,
                          model,
                          class_index=class_index
                          )

#save the resulting image
explainer.save(grid, "./outputs/explain/", "grad_cam_cat.png")
```



TF EXPLAIN

```
!pip install tf_explain
!pip install opencv-python

#load libraries
import numpy as np
import tensorflow as tf
import PIL

#load GradCAM
from tf_explain.core.grad_cam import GradCAM

IMAGE_PATH = "./assets/images/cat.jpg"
class_index = 281

img = tf.keras.preprocessing.image.load_img(IMAGE_PATH, target_size=(224, 224))
img = tf.keras.preprocessing.image.img_to_array(img)

model = tf.keras.applications.VGG16(weights="imagenet", include_top=True)

#get model summary
model.summary()

#first create the input in a format that the explainer expects (a tuple)
input_img = (np.array([img]), None)

#initialize the explainer as an instance of the GradCAM object
explainer = GradCAM()

# Obtain explanations for your image using VGG 16 and GradCAM
grid = explainer.explain(input_img,
                          model,
                          class_index=class_index
                          )

#save the resulting image
explainer.save(grid, "./outputs/explain/", "grad_cam_cat.png")
```



```
e Module for Grad CAM Algorithm
```

```
import numpy as np
import tensorflow as tf
import cv2
from tf_explain.utils.display import grid_display, heatmap_display
from tf_explain.utils.saver import save_rgb
```

```
class GradCAM:
```

```
    """
    Perform Grad CAM algorithm for a given input
```

```
Paper: [Grad-CAM: Visual Explanations from Deep Networks  
via Gradient-based Localization](https://arxiv.org/abs/1610.02391).
```

```
"""
```

```
def explain(  
    self,  
    validation_data,  
    model,  
    class_index,  
    layer_name=None,  
    ...)
```

SOURCE CODE



<https://github.com/sicara/tf-explain>

- [x] gradCAM
- [x] smoothGrad
- [x] guided smoothgrad
- [x] integrated gradients
- [x] occlusion sensitivity
- [x] activations visualization

e Module for Grad CAM Algorithm

```
import numpy as np
import tensorflow as tf
import cv2
from tf_explain.utils.display import grid_display, heatmap_display
from tf_explain.utils.saver import save_rgb
```

class GradCAM:

"""

Perform Grad CAM algorithm for a given input

Paper: [Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization](<https://arxiv.org/abs/1610.02357>)

"""

```
def explain(
    self,
    validation_data,
    model,
    class_index,
    layer_name=None,
```

SOURCE CODE



<https://github.com/sicara/tf-explain>

- [x] gradCAM
- [x] smoothGrad
- [x] guided smoothgrad
- [x] integrated gradients
- [x] occlusion sensitivity
- [x] activations visualization

this is an active area of research, so expect more methods to be added in the future

THAT NIGERIAN PRINCE NEVER
EMAILLED BACK

THE NIGERIAN PRINCE

I HOPE HE'S OKAY

AT NIGERIAN PRINCE NEVER EMAILLED BACK

OPE HE'S OK

MOVING BEYOND FEATURE ATTRIBUTION

- deep neural networks are notoriously sensitive to small perturbations in the input
- by intentionally introducing small perturbations in the input, we can fool the model into making a different prediction

AT NIGERIAN PRINCE NEVER EMAILLED BACK



MOVING BEYOND FEATURE ATTRIBUTION

- deep neural networks are notoriously sensitive to small perturbations in the input
- by intentionally introducing small perturbations in the input, we can fool the model into making a different prediction

*it's a cat and mouse game
between the attacker and the
defender*

OPE HE'S OK

ADVERSARIAL TRAINING vs ADVERSARIAL ATTACKS



- Using **complete knowledge** of the model architecture, its sources of uncertainty, its parameters, and the training data, we can intentionally introduce adversarial examples to the model to test the robustness and reliability of the model. **adversarial training**. Used to improve the network's ability to make accurate predictions in real-world scenarios.

*An example of adversarial training
is training a neural network to
recognize handwritten digits by
incorporating adversarial
examples of slightly modified
digits to the training data.*

ADVERSARIAL TRAINING vs ADVERSARIAL ATTACKS



→ Using **complete knowledge** of the model architecture, its sources of uncertainty, its parameters, and the training data, we can intentionally introduce adversarial examples to the model to test the robustness and reliability of the model. **adversarial training**. Used to improve the network's ability to make accurate predictions in real-world scenarios.

An example of adversarial training is training a neural network to recognize handwritten digits by incorporating adversarial examples of slightly modified digits to the training data.

→ Using **incomplete knowledge** of the model architecture, its sources of uncertainty, its parameters, and the training data, we can intentionally try to manipulate the input data to cause it to make incorrect predictions. **adversarial attacks**. Used to identify vulnerabilities in the network and to improve security.

An example of an adversarial attack is adding a small amount of noise to an image of a stop sign in order to make it appear as a go sign to an autonomous vehicle's image recognition system.



ADVERSARIAL ATTACKS



BY KÁROLY ZSOLNAI-FÉHER
PAPERS

WITH KÁROLY ZSOLNAI-FÉHER.

99.9% CAT

SUMMARY

SUMMARY

→ **XAI** is a field of research that aims to make machine learning models more interpretable and explainable

SUMMARY

- **XAI** is a field of research that aims to make machine learning models more interpretable and explainable
- **linear models** are intrinsically interpretable because they are linear functions of the input features

SUMMARY

- **XAI** is a field of research that aims to make machine learning models more interpretable and explainable
- **linear models** are intrinsically interpretable because they are linear functions of the input features
- **decision trees** are intrinsically interpretable because they are a series of if-then-else statements

SUMMARY

- **XAI** is a field of research that aims to make machine learning models more interpretable and explainable
- **linear models** are intrinsically interpretable because they are linear functions of the input features
- **decision trees** are intrinsically interpretable because they are a series of if-then-else statements
- **neural networks** are not intrinsically interpretable because they are non-linear functions of the input features. They need to be **interpreted** for high-stakes decision making applications.

SUMMARY

- **XAI** is a field of research that aims to make machine learning models more interpretable and explainable
- **linear models** are intrinsically interpretable because they are linear functions of the input features
- **decision trees** are intrinsically interpretable because they are a series of if-then-else statements
- **neural networks** are not intrinsically interpretable because they are non-linear functions of the input features. They need to be **interpreted** for high-stakes decision making applications.
- **quantifying uncertainty** is a key component of XAI regardless of the model architecture

SUMMARY

- **XAI** is a field of research that aims to make machine learning models more interpretable and explainable
- **linear models** are intrinsically interpretable because they are linear functions of the input features
- **decision trees** are intrinsically interpretable because they are a series of if-then-else statements
- **neural networks** are not intrinsically interpretable because they are non-linear functions of the input features. They need to be **interpreted** for high-stakes decision making applications.
- **quantifying uncertainty** is a key component of XAI regardless of the model architecture
- methods like **grad-CAM** help us understand which parts of the input image are most important for the model to make a prediction

SUMMARY

- **XAI** is a field of research that aims to make machine learning models more interpretable and explainable
- **linear models** are intrinsically interpretable because they are linear functions of the input features
- **decision trees** are intrinsically interpretable because they are a series of if-then-else statements
- **neural networks** are not intrinsically interpretable because they are non-linear functions of the input features. They need to be **interpreted** for high-stakes decision making applications.
- **quantifying uncertainty** is a key component of XAI regardless of the model architecture
- methods like **grad-CAM** help us understand which parts of the input image are most important for the model to make a prediction
- **adversarial attacks** and **adversarial training** are two sides of the same coin to testing the reliability and robustness of machine learning models

EXPLODING ROCKETS

wat

#\$@&%*?!**& ^%\$!

THANK YOU



[@nbhushan](#)