

## **Concept-based Explainable AI**

Explainable and Trustworthy AI

Gabriele Ciravegna



- 1. Motivation
- 2. Concept-based eXplainable AI (C-XAI)
- 3. Testing with Concept Activation Vectors (T-CAV)
- 4. Concept Bottleneck Models (CBM)
- 5. Concept Embedding Models (CEM)



# 3. Testing with Concept Activation Vectors (T-CAV)



#### Example: Post-training explanation

- To use machine learning responsibly, we need to ensure that
  - Our values are aligned
  - Our knowledge is reflected
- Standard XAI Solutions
  - Interpretable ML model (e.g. linear model)
    - Simple but we significantly lose the performance
  - Post-training explanation
    - E.g. Perturbation-based/sensitivity analysis-based methods
    - May be difficult to trust for standard users

#### Example: Post-training explanation



• Why was this a cash machine?

#### Problem Objective

Corresponding **Given Image** Saliency Map Hour Cas

**Prediction: Cash-machine** 

**Prediction:** 

- Did the **'human'** concept matters?
- Did the 'paper' concept matters?
- Did the 'ATM' or 'Cash' concept matters?

Sliding door

#### **TCAV objective**:

Quantitatively measure how

important are "user- chosen concepts"

#### TCAV: Overview



#### TCAV components

- a) A dataset annotated with both **examples of concepts** and **random images**
- b) The dataset with the **original classes**
- c) The **model** to explain
- d) The Concept Activation Vectors (CAV)
- e) The TCAV score showing the **influence** of a concept on a given class

#### TCAV: (1) How to define CAV?



#### Sorting Images with CAVs

- Given a set of images (e.g., belonging to the same class)
- Compute the cosine similarity between
  - the latent representation of an image  $f_l(x)$
  - the CAV  $v_C^l$  of the selected concept

#### CEO concept: most similar striped images



#### CEO concept: least similar striped images



#### Model Women concept: most similar necktie images



#### Model Women concept: least similar necktie images





#### TCAV: (2) How to compute TCAV scores?



$$S_{C,k,l}(\textcircled{\baselineskip})$$

$$S_{C,k,l}(\textcircled{\baselineskip})$$

$$S_{C,k,l}(\textcircled{\baselineskip})$$

$$S_{C,k,l}(\textcircled{\baselineskip})$$

$$\mathrm{TCAVQ}_{C,k,l} = \frac{|\{\boldsymbol{x} \in X_k : S_{C,k,l}(\boldsymbol{x}) > 0\}|}{|X_k|}$$

#### **Directional derivative with CAV:**

- $S_{C,k,l}(x) > 0$ : positive influence
- $S_{C,k,l}(x) < 0$ : negative influence
- The TCAV score is the number of class samples having a positive directional derivative w.r.t. the CAV

## TCAV score characteristcs

- $TCAV_{C,k,l} \in [0,1]$ 
  - $TCAV_{C,k,l} > 0.5$  : positive influence  $TCAV_{C,k,l} < 0.5$  : negative influence
  - Of concept *C*
  - Over class k
  - Computed in layer *l*

#### TCAV Example 1 (Zebra)



#### TCAV Example 2 (Doctor)



Was Woman concept important to this doctor image classifier?



TCAV tells that Woman has a negative importance for the classification of doctors

**BIAS IDENTIFICATION!** 

#### When and where can concept be learnt?

- Accuracy of the «linear probe»
  - High implies the network has automatically learnt a concept
  - Low implies the network does not use that concept for predicting the final class



mixed3a mixed3b mixed4a mixed4b mixed4c mixed4d mixed4e mixed5a mixed5b logit

- Simpler concepts have high accuracy throughout the NN
- High-level concepts can be detected better at higher layers

## 2. Concept Bottleneck Models (CBMs)



#### End-2-End models are difficult to interact with



## Ideal: Interact through high-level concepts



### CBMs Explicitly Represents Concepts



#### CBMs Allows Interactions!



#### CBMs Allows Interactions!





**CONCEPTS** 

### Importance of Concept Intervention



#### Concept bottleneck models architecture



### Different training strategy

• Indipendent: 
$$\hat{f} = \arg \min_{f} \sum_{i} L_{y}(f(c_{i}), y_{i})$$
 f is trained using the truth concepts  
 $\hat{g} = \arg \min_{g} \sum_{i} L_{c}(g(x_{i}), c_{i})$ 

• Sequential:  $\hat{f} = \arg \min_f \sum_i L_y(f(g(x_i)), y_i)$  g is trained first as above, then freezed

• Joint:  $\hat{f}, \hat{g} = \arg \min_{f} \sum_{i} L_{y}(f(c_{i}), y_{i}) + \lambda \arg \min_{g} \sum_{i} L_{c}(g(x_{i}), c_{i})$  for some  $\lambda > 0$ 

• Standard:  $\hat{f}, \hat{g} = \arg \min_{f} \sum_{i} L_{y}(f(c_{i}), y_{i})$ 

It ignores the concepts loss

# Different interpretability/performance trade-offs



- Sequential and indipendent are the more «trustworthy» beacause they ensure no concept leakage
- Joint strategy provides better task accuracy
  - Different trade-offs according to the λ value
- **Standard** model still has higher accuracy on average

## Explicitly concept training ensure model learns the concepts



Standard E2E trained model

Method	X-Ray Concept Error (↓)
Independent	0.53
Sequential	0.53
Joint	0.54
TCAV [Probe]	0.68

In a trained model, identifying some concepts may not be possible, because it might not have learnt them automatically

 $\rightarrow$  Only by explicitly training a model we can ensure it represents all concepts!

#### CBM Drawbacks

#### Poor Trade-offs

Struggle to compromise between accuracy and explainability

Accuracy-Explainability Trade-Off 100 No concepts Task Accuracy (%) 90 80 70 60 CBM Fuzzy 50 (Koh et al.) 50 60 90 100 70 80 Concept Alignment (%)

#### Low Concept Efficiency

CBMs do not scale in real-world conditions



# 3. Concept Embedding Models (CEM)



## Concept Embedding Models: overview



## Concept Embedding workflow

- 1.  $h = \psi(x)$ : the latent space of the model
- 2.  $c_i^+ = \phi_i^+(x)$ : neural model dedicated to represent the i-th **positive** concept embedding
- 3.  $p_i = s([c_i^+, c_i^-])$ : the *concept score* (i.e., probability of presence of the ith concept) is a function shared among concepts working on the concatenations of the concept embeddings
- 4.  $\hat{c}_1 = p_i c_i^+ + (1 p_i) c_i^-$ : the *concept embedding* is represented by the weighted combination of the positive and negative concept embeddings according to its presence
- 5.  $f([\hat{c}_1, \dots, \hat{c}_k])$ : the task predictor works on the concatenation of all the concept embeddings

## CEM: A neural-symbolic approach

Neural

Symbolic (CBM)

Concepts are represented with: unsupervised **embeddings**  Concepts are represented with: **supervised** scalars

Neural Symbolic (CEM)

Concepts are represented with: pairs of **supervised embeddings** 

 $\mathbf{c}_{i} \in \mathbb{R}^{k}$ 

c<sub>i</sub> ∈ [0,1]

 $C_i \in \mathbb{R}^k$  $c_i = agg(c_i^+, c_i^-)$  CEM Advanatages







Beyond Trade-offs

CEMs overcome the current accuracy-explainability trade-off

#### **High Concept Efficiency**

CEMs scale to real-world conditions where concept supervisions are scarce

#### **Effective interventions**

CEMs are responsive to concept interventions

## CEM vs Hybrid approach

#### • PROS:

- Retain high accuracy
- Has high concept efficiency like CEM

- CONS:
  - Prevent any effect of concept intervention
    - Changing the predicted scores has no effect on the task prediction
  - All the information required to predict the task is encoded by the unsupervised neurons



#### Have we lost something?

#### Interpretability

CBM: Interpretable



#### CEM: NON-Interpretable



 $\hat{\mathbf{c}}_{\text{yellow}} = [2.3, 0.3, -3.5, \dots]^T$ 

## Can we create an Interpretable Model over Concept Embeddings?



## Come on Monday to the Project Presentation!

- You will form groups of about 4 people
- We will provide 8-10 different projects among which you will have to choose
- The remaining of the lecture we will do:
  - A laboratory on C-XAI
  - A guided laboratory on XAI for Text Data with Prof. Eliana Pastor