

## Week 4 - Example-Based Interpretability

Wednesday, February 5, 2025 1:11 PM

### Motivation

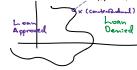
- Reasoning by analogy to examples can be powerful. We'd like ML systems that support this type of reasoning.
- These tools are especially useful when we want to avoid ML jargon. A general audience might not be able to make sense of an explanation of "CNN Filter 1432," but they "can" make sense of how the model behaves on examples.
- Unlike feature attribution, there doesn't seem to be a consensus around a few central approaches (like SHAP). So in these notes, we'll review a few representative methods (examples!) from the literature.

### Synthesizing Examples

#### Counterfactuals

- Surprisingly, some of the simplest approaches to example-based interpretability don't refer to examples in the training data. Instead, they "synthesize" examples that help make sense of the original predictions. We'll focus on Wachter et al. (2018)'s "counterfactual" and Hendricks et al. (2018)'s "visual" explanations.
- Their motivation comes from a close reading of GDPR. They realized that describing internal model mechanics wouldn't help average citizens with:
  - Understanding why a prediction was made.
  - Contesting the factors that led to a prediction.
  - Adopting to help achieve their goals ("algorithmic recourse")

To resolve this, they defined counterfactual explanations as examples that "are as close as possible" to the real example of interest, but which have a different (desired) model prediction.



- There are many ways to operationalize this definition. They suggest solving the optimization:

$$x^* = \arg \max_{x'} \lambda (f(x') - y)^2 + d(x, x')$$

For a target prediction  $y$  and a target sample  $x$ , solve:

$$d(x, x') = \sum_k \frac{1}{2k} \|x_k - x'_k\|^2$$

- Their illustrations are relatively simple (2-layer neural net, few features), but their main arguments are conceptual. In this example (Table 1 in the paper), they trained a model to predict GPA in the first year of law school based on characteristics from an admissions classifier. The left panel gives original examples, and those in the middle/right are counterfactuals with target scores  $y = 0$ . The change in score from the original example to the counterfactual is the far left column "score."

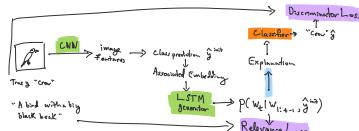
	Normalized L2		L1		L2		L1		L2		L1	
	GPA	LSAT	Race	GPA	LSAT	Race	GPA	LSAT	Race	GPA	LSAT	Race
0.17	3.1	39.0	0	3.0	37.0	0.2	3.0	34.0	0	3.0	37.0	0
0.41	3.3	38.0	1	3.5	38.0	0.4	3.4	33.0	0	3.4	38.0	0
-0.77	3.3	38.0	1	3.5	39.8	0.4	3.4	33.4	0	3.4	38.0	0
-0.83	2.4	38.5	1	2.7	37.4	0.2	2.6	35.7	0	2.7	37.4	0
-0.57	2.7	38.5	0	2.8	39.1	-0.4	2.9	34.1	0	2.7	38.5	0

#### Visual Explanations

- Counterfactual examples look like the original data. Hendricks et al. (2016) instead synthesize text explanations to accompany the predictions on image data. For being the oldest paper we review this week, it feels the most present, both in its use of multimodal data and on using language models to generate explanatory text.

- The key challenge they emphasize is that description ≠ explanation. A description only needs to accurately summarize what's in an image. An explanation needs to highlight those features that are relevant to the model's class prediction.

- Their architecture has this form:



The model parameters are highlighted green. The sentence classifier (orange) is pretrained in advance. The blue edge involves sampling a new sentence. This complicates optimization, but it is solved using the REINFORCE algorithm (a low sample-efficiency approach, but it seems to work).

- The two losses (purple) serve complementary purposes:

- Relevance Loss: Do the sampled sentences look like the original sentence description from the training data?
- Discriminator Loss: Do the sampled sentences help predict the correct class?

The second term is what makes this different from a caption model, say, it forces the generated sentences to mention features that help with classification.

This is a **Bronzed Cuckoo**.  
Definition: This bird is black with blue on its wings and has a long pointy beak.  
Description: This bird is mostly all black with a short pointy bill.  
Explanation-Lab: This is a black bird with a long beak.  
Explanation-Ds: This is a black bird with a red eye and a white beak.  
Explanation: This is a black bird with a red eye and a pointy black beak.

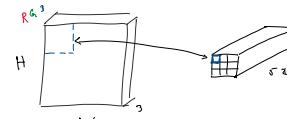
### Finding Examples

#### Parts-Based Explanation

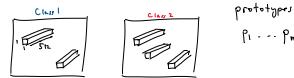
- Chen et al. (2019) try defining an intrinsically interpretable model that:
  - Defines representative image patches ("prototypes").
  - Uses similarity to those prototypes to guide classification.

Since the prototypes can be matched with training data examples, this helps users identify which parts of which images led to the final classification.

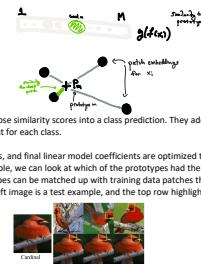
- CNNs transform images into "long" cubes. The depth dimension can be interpreted as different learned image features. The spatial position in that cube ( $W \times H$ ) corresponds to a spatial patch in the original image. E.g., the top left corner in the cube corresponds to the top left patch in the original image.



- The idea is to create class-specific prototypes with the same depth as these representations and which are just  $1 \times 1$  pixel tall/wide. You can imagine comparing each prototype to the features associated with a path.



- Specifically, the "long" summaries of each image patch can be compared to each prototype. Patches that are close to a prototype are more likely to come from the same class as that prototype. For prototype  $m$ , find the most similar patch in  $x[j]$  and store that similarity score in coordinate  $m$  of the vector  $g(f(x[j]))$ .



- A final linear model maps those similarity scores into a class prediction. They add a sparsity penalty so that only a few prototypes are relevant for each class.

- The prototypes, CNN weights, and final linear model coefficients are optimized to maximize classification accuracy. For any new example, we can look at which of the prototypes had the most influence on the final classification. These prototypes can be matched up with training data patches that give embeddings close to those prototypes. E.g., the left image is a test example, and the top row highlights relevant training patches.



- Can you come up with geometric interpretation of the "Cist" or "Sep" terms in objective on page 5? It should look similar to the drawing under point 15) above.

#### Prototypes and Criticism

- Kim et al. (2016) try to extract relevant examples without reference to any specific model architecture. They identify prototypes (now meaning "typical examples") and criticism (unusual examples) to help people understand a dataset. In principle they could help clarify differences in model behavior on these selected examples. However, their human subjects study simply considers whether participants can identify the correct class (not necessarily any model's prediction).

- We could use any clustering method to find prototypes and any anomaly detection method to find criticisms. They discuss how to use Maximum Mean Discrepancy (MMD) for both tasks.

$$\text{MMD}(\mathcal{F}, \mathcal{P}, \mathcal{Q}) = \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\mathcal{P}}[f(x)] - \mathbb{E}_{\mathcal{Q}}[f(x)] \right|$$

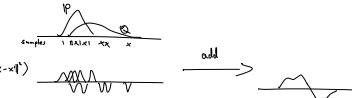
- If  $\mathcal{F} = \{x, x'\}$ , then for which pairs of  $\mathcal{P}$  and  $\mathcal{Q}$  is the MMD 0?

- To contrast two distributions, we can use a witness function,

$$J(x) = \left| \mathbb{E}_{\mathcal{P}}[K(x, x')] - \mathbb{E}_{\mathcal{Q}}[K(x, x')] \right|$$

Imagine drawing  $N$  samples from both  $\mathcal{P}$  and  $\mathcal{Q}$ . Then we can approximate the witness function.

$$J(x) = \frac{1}{N} \left| \sum_{x_i \sim \mathcal{P}} K(x, x_i) - \sum_{x_i \sim \mathcal{Q}} K(x, x_i) \right|$$



- Let  $\mathcal{P}$  be the empirical distribution over the entire training set. Let  $\mathcal{Q}$  be the empirical distribution over a selected subset  $S$ . To find prototypes, we look for a subset  $S$  so that  $\mathcal{P}$  and  $\mathcal{Q}$  have small MMD distance from one another. Note that the subset  $S$  that minimizes the MMD also maximizes

$$J(S) = \frac{1}{n_S} \sum_{i \in S} K(x, x_i) - \text{MMD}^2(\mathcal{F}, \mathcal{P}|_{S^c}, \mathcal{Q}|_S)$$

It turns out that this is a submodular function, which means greedy heuristics can be used to solve the problem efficiently.

- Prototypes are chosen to be representative of the overall training data distribution. However, these summaries are never complete, and criticism can clarify regions of the training data distribution that are not reflected in the prototypes. We can use the witness function between the overall distribution  $\mathcal{P}$  and the prototypical subset  $S$  to identify these criticism samples:

$$C^* = \arg \max_{\text{subset } S \subseteq \mathcal{C}} \sum_{x \in S} |f(x)|$$

- The idea is to create class-specific prototypes with the same depth as these representations and which are just  $1 \times 1$  pixel tall/wide. You can imagine comparing each prototype to the features associated with a path.

- They recommend also adding a regularization term to this objective function to help "spread out" the selected criticism points. The regularizer is based on the log determinant of the kernel matrix.

- The criticisms seem to capture examples from a class that are important for characterizing it, but which are not the typical cases. Here are some examples of prototypes and criticisms from a dog breed classification dataset:



Beyond the discussion in this paper, one reason we might be interested in criticisms is because there are some indications that some of the successes of AI/deep learning models comes from the fact that they are better than other models at memorizing the labels associated with these rare examples. For example, it seems that the performance improvements in larger models seems to come from their improved accuracy on clusters of outlying, non-representative samples; e.g., [1].