

STOR566: Introduction to Deep Learning

Lecture 22: Explainability of ML Models

Yao Li
UNC Chapel Hill

Nov 15, 2022

Materials are from *Deep Learning (UCLA)*

Motivations

Why Explainability

- Understand predictions/decisions of machine learning models

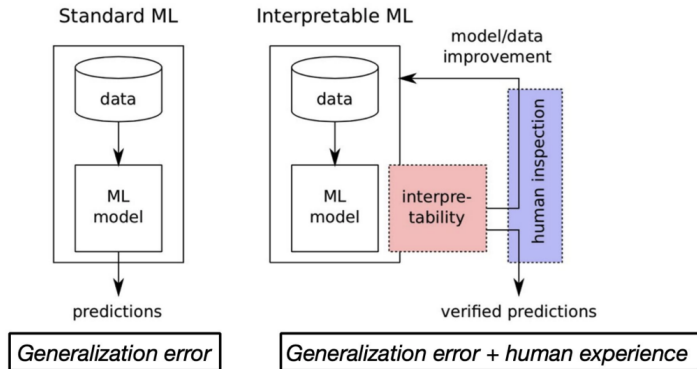


Top label: **“clog”**

Why did the network label this image as **“clog”**?

Why Explainability

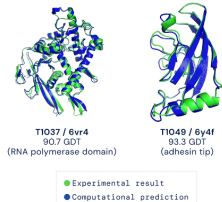
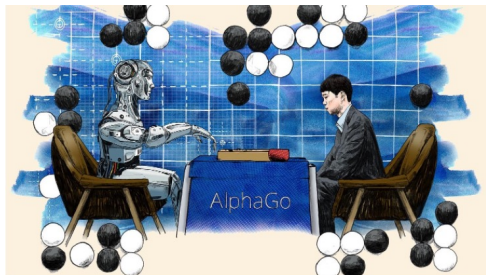
- Improve machine learning models



Credit: Samek, Binder, Tutorial on Interpretable ML, MICCAI'18

Why Explainability

- Learn new insights



Approaches

Model based

Build interpretable ML
models

Post-hoc

Generate explanations
for a given ML model

today's focus

Approaches

What is an “**explanation**”?

- Feature attribution

The model makes this prediction because of feature (pixel) X

- Data attribution

The model makes this prediction because of which training data

- Surrogate model

Approximate the complex model using a simple explainable surrogate model

Approaches

What is an “**explanation**”?

- Feature attribution

The model makes this prediction because of feature (pixel) X

- Data attribution

The model makes this prediction because of which training data

- Surrogate model

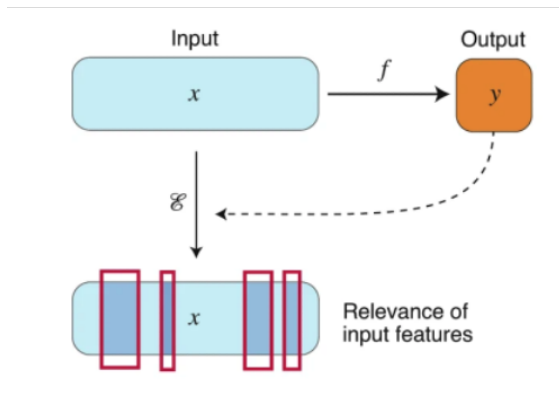
Approximate the complex model using a simple explainable surrogate model

We focus on the **first two types** in today's lecture.

Feature Attribution

What is Feature Attribution?

- Given a model f and input x
 - Assign a relevant score to each input feature
- R_i : how much does feature i contribute to the prediction



Perturbation-based Analysis

- Assumption: Feature i is important \rightarrow Perturbing \mathbf{x}_i (the i -th feature of input \mathbf{x}) will **significantly** change $f(\mathbf{x})$
- Therefore, assign relevant score for each feature by

$$R_i \leftarrow f(\mathbf{x}) - f(\mathbf{x} + \delta \mathbf{e}_i),$$

R_i : importance score of i -th input feature

\mathbf{e}_i : vector of zeros except that the i -th element is one. The dimension is the same as \mathbf{x} .

Perturbation-based Analysis

- Assumption: Feature i is important \rightarrow Perturbing \mathbf{x}_i (the i -th feature of input \mathbf{x}) will **significantly** change $f(\mathbf{x})$
- Therefore, assign relevant score for each feature by

$$R_i \leftarrow f(\mathbf{x}) - f(\mathbf{x} + \delta \mathbf{e}_i),$$

R_i : importance score of i -th input feature

\mathbf{e}_i : vector of zeros except that the i -th element is one. The dimension is the same as \mathbf{x} .

- Questions:
 - How to choose perturbation δ ?
 - Efficiency: may need $O(d)$ function evaluations. d is the dimension of \mathbf{x} .
 - Can this capture the correlation between features?

Linear Model

- Linear prediction

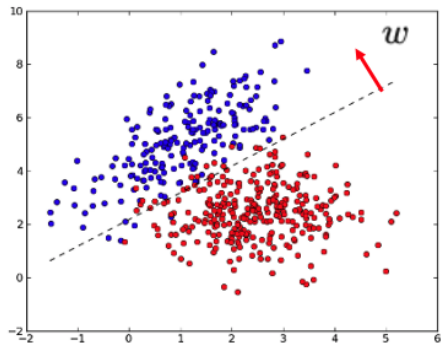
$$\begin{aligned}f(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} \\ &= w_1 x_1 + w_2 x_2 + \dots + w_d x_d\end{aligned}$$

- For any fixed perturbation

$$f(\mathbf{x}) - f(\mathbf{x} + \delta \mathbf{e}_i) = \delta w_i \propto w_i$$

- Feature importance:

$$R_i \leftarrow w_i$$



Gradient

- Consider the case when $\delta \rightarrow 0$

$$\lim_{\delta \rightarrow 0} f(\mathbf{x}) - f(\mathbf{x} - \delta \mathbf{e}_i) = \frac{\partial}{\partial x_i} f(\mathbf{x})$$

- Therefore, we can use gradient to measure importance of each feature

$$R_i \leftarrow \frac{\partial}{\partial x_i} f(\mathbf{x})$$

(Usually set f as the logit of the target model)

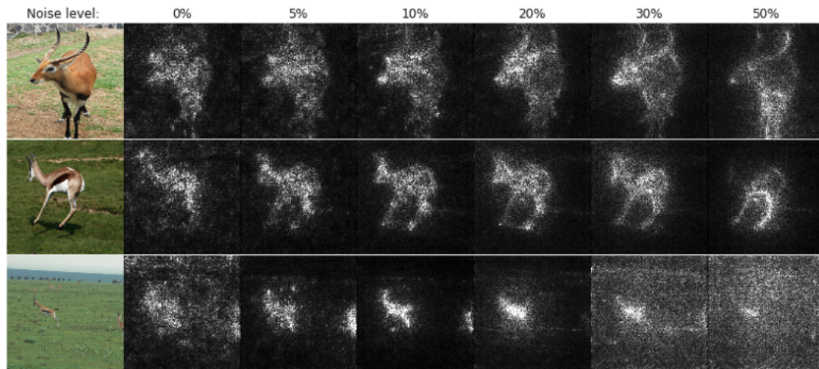
- Saliency map: visualize pixels with positive gradients



Smoothed Gradient

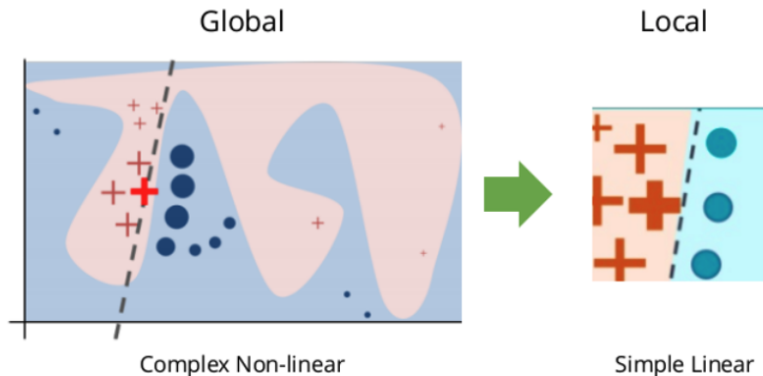
- Gradient maps are often noisy (visually)
- Smoothed gradient:

$$R(\mathbf{x}) = E_{z \sim N(0, \sigma)} \nabla f(\mathbf{x} + \mathbf{z})$$



LIME

- Build a local linear model, but not as local as gradient.



LIME

- Find a local linear model to mimic target (complex) model

$$\arg \min_{g \in G} L(f, g, \pi_x) + \Omega(g)$$

Family of simple
models

Target
model

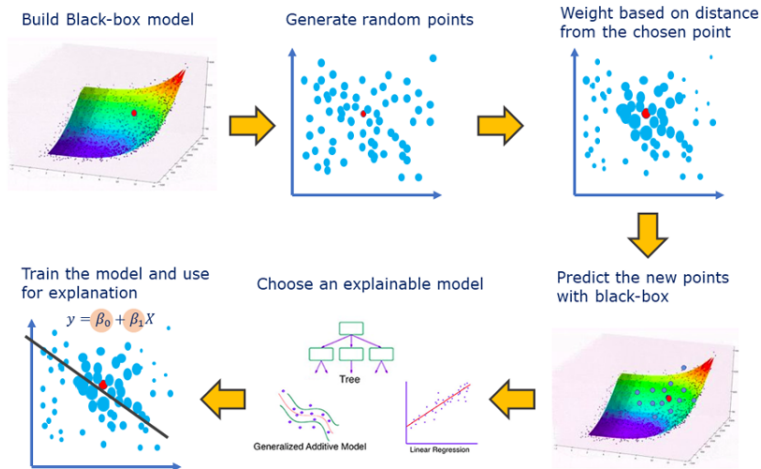
Simple
model

Sample
weights

Regularization

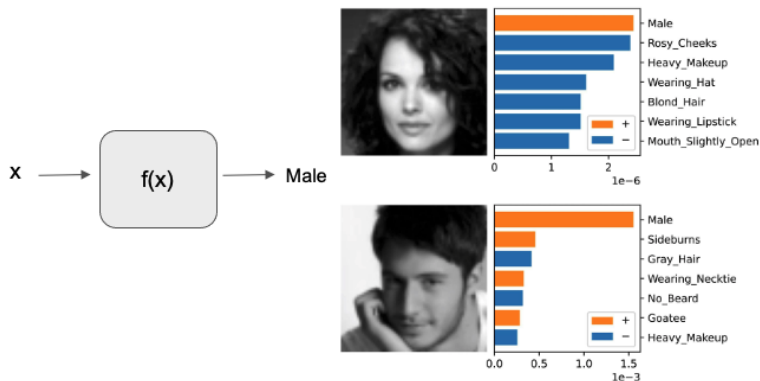
- Sample weights: more weights to local samples

- Overview



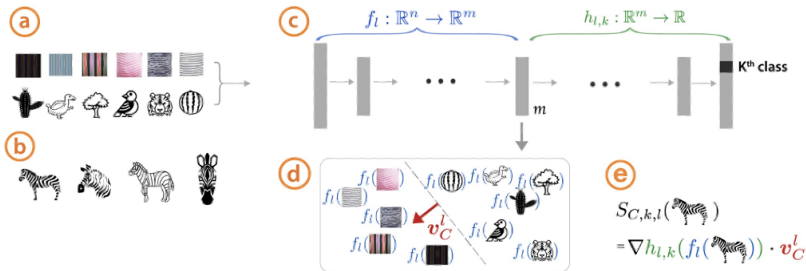
Concept-based Explanations

- Attribution to raw features may not be **human understandable**
- Can we attribute the prediction to **high-level concepts** instead of low-level features (pixels)?

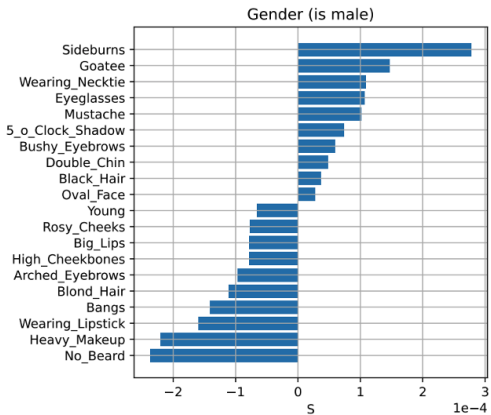
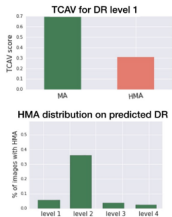
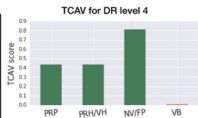


TCAV

- A concept can be given as positive/negative instances
- Assume layer l in DNN captures the concept $c = \langle v_C^l, x^l \rangle$
- Then attribution to concept is the product of gradient $\frac{\partial f}{\partial x^l}$ and v_C^l



Experimental results



Data Attribution

Explaining by Training Data

- Which training data causes the prediction?

test id5727
rhinoceros predicted as
rhinoceros



because?

Training data

train id29490
zebra predicted as
zebra



or

train id23304
rhinoceros predicted as
rhinoceros



Influence Function

- What's the relationship between training data and model?

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, \theta)$$

$\hat{\theta}$: model parameters

\mathbf{x}_i : training sample i

Influence Function

- What's the relationship between training data and model?

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, \theta)$$

$\hat{\theta}$: model parameters

\mathbf{x}_i : training sample i

- Each training sample contributes to the model equally ($\frac{1}{n}$ is the weight of each sample)
- Can we compute the influence when the weight of a training sample slightly increased or decreased?

Koh et al., Understanding Black-box Predictions via Influence Functions. ICML, 2017.

Influence Function

- Assume adding more weight (ϵ) to \mathbf{x}_j (the j -th training sample)
- The model will become

$$\hat{\theta}_{\epsilon, \mathbf{x}_j} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, \theta) + \epsilon L(\mathbf{x}_j, \theta)$$

- Gradient of loss w.r.t. ϵ :

$$\begin{aligned} l_{\text{up,loss}} &:= \left. \frac{dL(\mathbf{x}_{\text{test}}, \hat{\theta}_{\epsilon, \mathbf{x}_j})}{d\epsilon} \right|_{\epsilon=0} \\ &= -\nabla_{\theta} L(\mathbf{x}_{\text{test}}, \hat{\theta})^T H_{\hat{\theta}}^{-1} \nabla_{\theta} L(\mathbf{x}_j, \hat{\theta}) \end{aligned}$$

\mathbf{x}_{test} : a test sample

$$H_{\hat{\theta}} := \frac{1}{n} \nabla_{\theta}^2 L(\mathbf{x}_i, \hat{\theta})$$

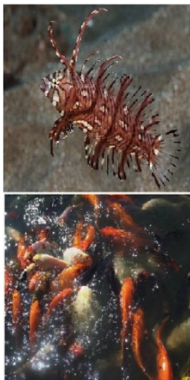
Influence Function

- Experimental results

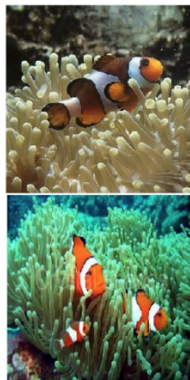
Test image



**RBF SVM
(raw pixels)**



**Logistic regression
(Inception features)**




Influence Function

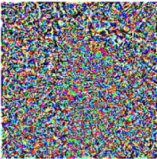
- Can be used for poisoning attack to identify important training samples.

A small perturbation to one **training** example:

Label: Fish




+ $\epsilon \cdot$




→

Label: Fish



Can change multiple **test** predictions:



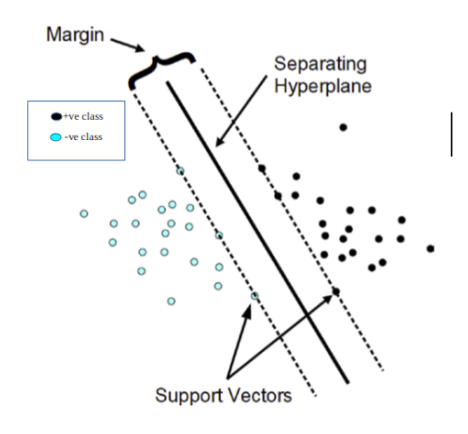
Orig (confidence):	Dog (97%)	Dog (98%)	Dog (98%)	Dog (99%)	Dog (98%)
New (confidence):	Fish (97%)	Fish (93%)	Fish (87%)	Fish (63%)	Fish (52%)

Representer Theorem in Linear Models

- Linear model: $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Representer theorem: model can be decomposed with training samples

$$\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{x}_i$$

- SVM: $\alpha_i \neq 0$ support vectors
- General: represent the importance of each sample



Representer Theorem in Linear Models

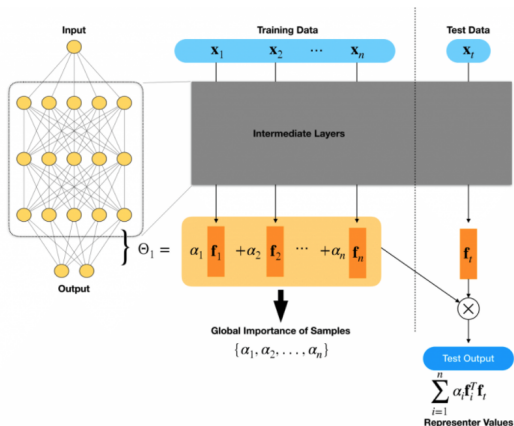
- For a test sample: \mathbf{x}_{test} :

$$f(\mathbf{x}_{test}) = \mathbf{w}^T \mathbf{x}_{test} = \sum_{i=1}^n \alpha_i \mathbf{x}_i^T \mathbf{x}_{test}$$

- $\alpha_i \mathbf{x}_i^T \mathbf{x}_{test}$: importance of sample \mathbf{x}_i in the final prediction based on \mathbf{x}_{test}
- A natural way to attribute prediction to each training sample

Representer Points in DNN

- Consider the final hidden layer output for each training sample:
 $f_i := f(\mathbf{x}_i), i = 1, \dots, n$
- Applying representer theorem to the final linear layer:
 $\Theta_1 = \sum_{i=1}^n \alpha_i f_i$
- Attribute the prediction by $F(\mathbf{x}_t) = \sum_{i=1}^n \alpha_i f_i^T f_t$
 f_t : final layer output based on test sample \mathbf{x}_t



Yeh et al., Representer Point Selection for Explaining Deep Neural Networks.

NeurIPS, 2018.

Representer points in DNN

test id3092
grizzly bear predicted as
grizzly bear



train id13033
grizzly bear predicted as
grizzly bear



POSITIVE Example

train id12728
grizzly bear predicted as
grizzly bear



POSITIVE Example

train id12742
grizzly bear predicted as
grizzly bear



POSITIVE Example

train id21249
polar bear predicted as
polar bear



NEGATIVE Example

train id1228
beaver predicted as
beaver



NEGATIVE Example

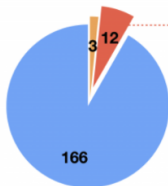
train id20730
pig predicted as
pig



NEGATIVE Example

Representer points in DNN

Test Points with Labels
Antelope



- Misclassified as Other
- Misclassified as Deer
- Correctly Classified



train id29372
predicted as zebra
true label is zebra



train id688
predicted as deer
true label is antelope



train id8090
predicted as elephant
true label is elephant



train id29208
predicted as zebra
true label is zebra



Conclusions

- Introduction to explainable ML
- Feature attribution
- Data attribution

Questions?