

Post Training Analysis

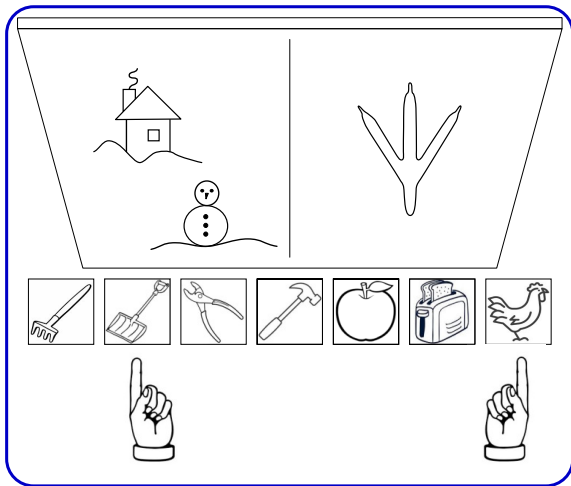
Deep Learning



- After training, analyze the network.
- One purpose of the analysis is to explain network decision making.
- In many applications (e.g., medicine) users need to understand the decisions in order to have confidence in them.
- The objective of **explainability** is to provide a narrative that gives the user confidence in the decisions.
- We will discuss techniques that **explain** how deep networks are making their decisions.



How good are human explanations?



Thoughts on explainability

- If deep networks take on complex tasks, their explanations may also be complex.
- We will discuss methods that provide insights into deep network operation.
- They don't provide complete [explanations](#).
- We use [explainability](#) in a technical way that does not include all meanings of its everyday usage.



- Merriam Webster's – to give the reason for or cause of.
- Related words – clarify, explicate, illustrate, interpret.
- How can we apply this concept to deep learning and find a quantitative measure?
- There will never be "ground truth", in that we will never be certain our explanation is correct.
- Explainability methods are mainly judged on their ability to provide a convincing narrative.



Weight of evidence

$$W(H : E) = \ln \frac{P(E|H)}{P(E|\bar{H})}$$

- H is an hypothesis, and E is evidence for that hypothesis.
- $P(E|H)$ is the conditional probability of E given H .
- \bar{H} is the complement of H .
- If E was more likely to occur when H was true, the weight of evidence would be positive.



One layer logsig network.

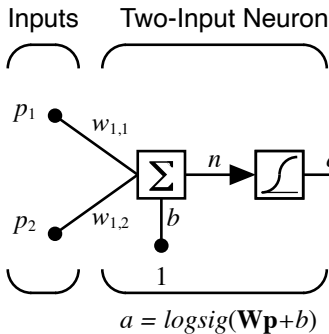
$$n_1 = w_{1,1}p_1 + w_{1,2}p_2$$

$$a_1 = \frac{1}{1 + e^{-n_1}}$$

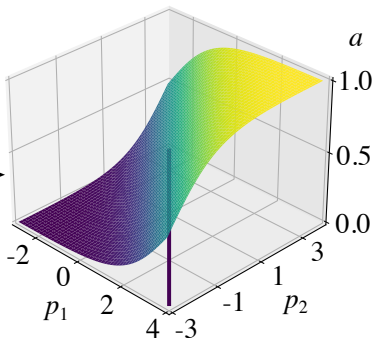
$$\begin{aligned} W(H : E) &= \ln \frac{P(E|H)}{P(E|\bar{H})} = \ln \frac{\frac{P(H|E)P(E)}{P(H)}}{\frac{P(\bar{H}|E)P(E)}{P(\bar{H})}} = \ln \frac{P(H|E)P(\bar{H})}{P(\bar{H}|E)P(H)} \\ &= \ln \frac{P(H|E)}{P(\bar{H}|E)} = \ln \frac{a_1}{1 - a_1} = \ln \frac{\frac{1}{1+e^{-n_1}}}{\frac{e^{-n_1}}{1+e^{-n_1}}} = \ln e^{n_1} = n_1 \\ &= \underbrace{w_{1,1}p_1}_{\text{first input contribution}} + w_{1,2}p_2 \end{aligned}$$



Illustrative example: single neuron



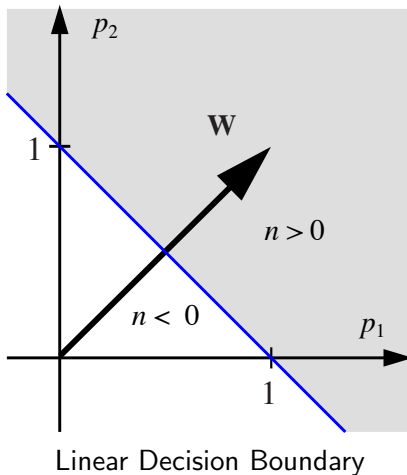
One Layer Network



Network Response

$$a = f(n) = f(\mathbf{W}\mathbf{p} + b) = \text{logsig}(1p_1 + 1p_2 - 1)$$





$$n = 1p_1 + 1p_2 - 1$$

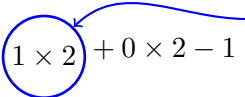
$$a = \text{logsig}(n)$$

- The global explanation is the weight (decision boundary).
- The weights have equal magnitude, so globally both inputs are equally important.
- For a particular example, the weights times the inputs are the local explanations.

$$\mathbf{p} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

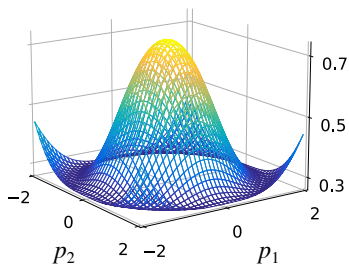
$$n = 1 \times 2 + 0 \times 2 - 1$$

main contribution

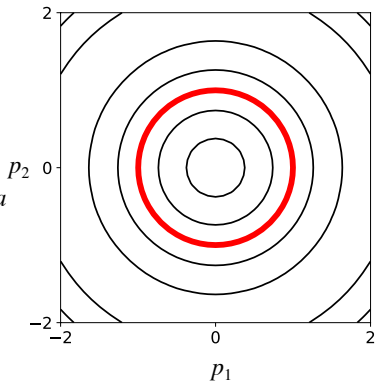




Global Explanation



Network Response Surface



Contour Plot and Decision Boundary



Taylor Series Expansion

$$a(\mathbf{p}) \cong a(\mathbf{p}^*) + \nabla a(\mathbf{p})^T|_{\mathbf{p}^*} (\mathbf{p} - \mathbf{p}^*)$$

- \mathbf{p}^* is the input where the expansion takes place.
- $\nabla a(\mathbf{p})$ is the gradient.

$$\nabla a(\mathbf{p}) = \left[\frac{\partial a(\mathbf{p})}{\partial p_1} \quad \frac{\partial a(\mathbf{p})}{\partial p_2} \quad \cdots \quad \frac{\partial a(\mathbf{p})}{\partial p_R} \right]^T$$

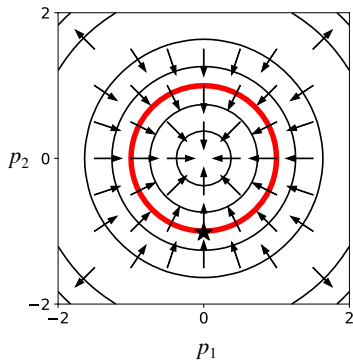
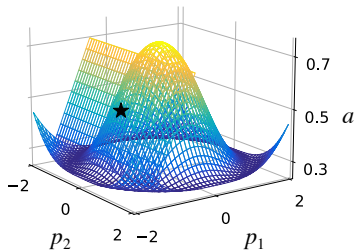
- The gradient is the weight in the approximate linear network.



Linearize the radial network

$$\mathbf{p}^* = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$a(\mathbf{p}) \cong 0.5 + \begin{bmatrix} 0 & 0.4 \end{bmatrix} \left(\mathbf{p} - \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right) = \begin{bmatrix} 0 & 0.4 \end{bmatrix} \mathbf{p} + 0.9$$



Notation and General Explainability Definitions

- We use Good/Poulin definition of **contribution**, but extended to multilayer networks.
- For a single-layer network, it is equal to the product of the input times the weight.
- We use two methods to extend this approach to multilayer networks.
- First, linearize the network about the relevant input vector.
- Second, backpropagate the contributions through the network.



Desired properties of contributions

- ① **Continuity**: If $a_c^M(\mathbf{p})$ is continuous, then $\mathbf{c}(\mathbf{p})$ is continuous.
- ② **Implementation Invariance**: If two networks are functionally equivalent (same input produces same output), then the computed contributions are the same.
- ③ **Sensitivity**
 - ① If an input and a baseline differ in one feature and have different network outputs, a nonzero contribution is assigned to that feature.
 - ② If the network output does not depend on a feature, the contribution of that feature is zero.
- ④ **Completeness**: The sum of all contributions is equal to the network output (or the difference between the output and a baseline output).



Derivative of Network Output with Respect to Input

$$c^{sal}(p_i) = \frac{\partial n_c^M(\mathbf{p})}{\partial p_i} \times p_i$$

$$\mathbf{s}^m = \dot{\mathbf{F}}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T \mathbf{s}^{m+1}$$

$$\mathbf{s}^M = \frac{\partial n_c^M}{\partial \mathbf{n}^M} = \boldsymbol{\epsilon}_c$$

$$\mathbf{c}^{sal}(\mathbf{p}) = \frac{\partial n_c^M}{\partial \mathbf{p}} \circ \mathbf{p} = \left(\frac{\partial (\mathbf{n}^1)^T}{\partial \mathbf{p}} \frac{\partial n_c^M}{\partial \mathbf{n}^1} \right) \circ \mathbf{p} = ((\mathbf{W}^1)^T \mathbf{s}^1) \circ \mathbf{p}$$



- The gradient can be zero at \mathbf{p} , even when the network is confident.
- Average the gradient over a line between the current input \mathbf{p} and a baseline input $\bar{\mathbf{p}}$ to get a more consistent contribution.

$$\mathbf{c}^{ig}(\mathbf{p}) = \left[\frac{1}{T} \sum_{t=0}^{T-1} \frac{\partial n_c^M(\tilde{\mathbf{p}})}{\partial \tilde{\mathbf{p}}} \Big|_{\tilde{\mathbf{p}}=\mathbf{p}^t} \right] \circ (\mathbf{p} - \bar{\mathbf{p}})$$
$$\mathbf{p}^t = \bar{\mathbf{p}} + \frac{t}{T-1} (\mathbf{p} - \bar{\mathbf{p}})$$



Contribution propagation

Compute contribution at last layer and then backpropagate it.

$$\mathbf{c}(\mathbf{a}^M) = n_c^M \boldsymbol{\epsilon}_c$$

$$\mathbf{c}(\mathbf{a}^m) = \mathbf{C}(\mathbf{a}^m | \mathbf{a}^{m+1}) \mathbf{c}(\mathbf{a}^{m+1}), \quad m = M - 1, M - 2, \dots, 0$$

$$[\mathbf{C}(\mathbf{a}^m | \mathbf{a}^{m+1})]_{i,j} = c(a_i^m | a_j^{m+1})$$

To maintain total contribution.

$$\sum_{i=1}^{S^m} c(a_i^m) = n_c^M, \quad \text{if} \quad \sum_{i=1}^{S^m} c(a_i^m | a_j^{m+1}) = 1$$



Layerwise Relevance Propagation (LRP)

LRP-0

$$c^{lrp-0} \left(a_i^m | a_j^{m+1} \right) = \frac{w_{j,i}^{m+1} a_i^m}{\sum_{k=1}^{S^m} w_{j,k}^{m+1} a_k^m}$$

LRP- ϵ

$$c^{lrp-\epsilon} \left(a_i^m | a_j^{m+1} \right) = \frac{w_{j,i}^{m+1} a_i^m}{\epsilon + \sum_{k=1}^{S^m} w_{j,k}^{m+1} a_k^m}$$

LRP- γ

$$c^{lrp-\gamma} \left(a_i^m | a_j^{m+1} \right) = \frac{\left(w_{j,i}^{m+1} + \gamma \left(w_{j,i}^{m+1} \right)^+ \right) a_i^m}{\sum_{k=1}^{S^m} \left(w_{j,k}^{m+1} + \gamma \left(w_{j,k}^{m+1} \right)^+ \right) a_k^m}$$



DeepLift measures contributions from a baseline $\bar{\mathbf{p}}$.

$$\Delta n_c^M = n_c^M(\mathbf{p}) - n_c^M(\bar{\mathbf{p}}) = n_c^M - \bar{n}_c^M$$

Contributions are defined in terms of changes.

$$\begin{aligned} c^{dl}(\Delta a_i^m) &= c^{dl}(\Delta a_i^m | \Delta n_c^M) = c(a_i^m) \\ c^{dl}(\Delta a_i^m | \Delta a_j^{m+1}) &= c(a_i^m | a_j^{m+1}) \Delta a_j^{m+1} \end{aligned}$$

Modified multipliers are divided by change.

$$\begin{aligned} q(\Delta a_i^m | \Delta a_j^{m+1}) &= \frac{c^{dl}(\Delta a_i^m | \Delta a_j^{m+1})}{\Delta a_i^m} = c(a_i^m | a_j^{m+1}) \frac{\Delta a_j^{m+1}}{\Delta a_i^m} \\ q(\Delta a_i^m) &= \frac{c^{dl}(\Delta a_i^m)}{\Delta a_i^m} = c(a_i^m) \frac{1}{\Delta a_i^m} \end{aligned}$$



DeepLIFT propagation equations

$$q(\Delta a_c^M) = 1$$

$$q(\Delta a_i^m) = \sum_{j=1}^{S^{m+1}} q(\Delta a_i^m | \Delta n_j^{m+1}) q(\Delta n_j^{m+1} | \Delta a_j^{m+1}) q(\Delta a_j^{m+1})$$

Rescale Rule

$$q(\Delta n_j^{m+1} | \Delta a_j^{m+1}) = \frac{\Delta a_j^{m+1}}{\Delta n_j^{m+1}}$$

Linear Rule

$$q(\Delta a_i^m | \Delta n_j^{m+1}) = w_{j,i}^{m+1}$$

DeepLIFT Contributions

$$c^{dl}(\Delta p_i) = q(\Delta p_i) \Delta p_i$$

