

Exercise 10 - Solution

Network Interpretability and Visualization

Reliable and Interpretable Artificial Intelligence
ETH Zurich



Figure 1: Images strongly activating convolutional filters in three layers created by optimization. Order permuted.

Problem 1 (Feature Visualization). Figure 1 shows three images created by optimization (with Lucid¹). These images strongly activate convolutional filters in different layers in a convolutional neural network. Order the images according to their depth in the neural network. Provide an argument to justify your answer.

Solution 1. Figure 2 shows the solution. We observe that the first image is a very simple filter for edges in the image. The second image shows some red-blue texture diagonal lines in it. The last image shows an abstract concept like some from of abstract scaled animal. While the interpretations here might be subjective, the overall insight that earlier layers produce primitive filters and later layers high level concepts is clear.

¹Googles framework for neural network visualization. <https://github.com/tensorflow/lucid>



Figure 2: Images strongly activating convolutional filters in three layers created by optimization.. Order by their position in the network. Each shows the activation of channel zero in the post-ReLU layers conv2d0, mixed3a, mixed5a in GoogLeNet respectively.

Problem 2 (Shapley Values). Consider the classifier $g : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ and data points \mathbf{x} and \mathbf{x}' .

$$\mathbf{x} := \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \quad \mathbf{x}' := \begin{pmatrix} -1 \\ -2 \\ 3 \end{pmatrix}$$

$$g(\mathbf{x}) := \mathbf{B} \text{ReLU}(\mathbf{A}\mathbf{x})$$

$$\mathbf{A} := \begin{pmatrix} 3 & 0 & 0 \\ 0 & -2 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \mathbf{B} := \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix}$$

1. Compute the Shapley Values for all components of both inputs and both target classes. Set non-included features to zero. Discuss the results. (Consider implementing repetitive calculation in python.)
2. Argue why Shapley Values are impractical (applied directly; without approximations) for large images.

Solution 2.

1. Since we have a fixed input and relatively simple neural network we can simplify $g(\mathbf{x})$: For the first logit we obtain $g(\mathbf{x})_1 = \text{ReLU}(3\mathbf{x}_1) + \text{ReLU}(-2\mathbf{x}_2)$ and for the second $g(\mathbf{x})_2 = \text{ReLU}(-2\mathbf{x}_2) + \text{ReLU}(1\mathbf{x}_3)$. Plugging in the values for \mathbf{x} we see that the first logit is directly proportional to the first component of the input and the second logit to the third component of the input. To calculate the Shapley Value for the first input \mathbf{x}_0 we let $P := \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ and set $f(S) := g(\mathbf{x}_S)_1$ the value of the first logit, where \mathbf{x}_S denotes the vector \mathbf{x} with components not included in S

set to 0. Thus, for the first component \mathbf{x}_1 we calculate:

S	\emptyset	$\{\mathbf{x}_2\}$	$\{\mathbf{x}_3\}$	$\{\mathbf{x}_2, \mathbf{x}_3\}$
$\frac{ S !(P - S -1)!}{ P !}$	$\frac{2}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{2}{6}$
$f(S)$	0	0	0	0
$f(S \cup \{\mathbf{x}_1\})$	3	3	3	3

Plugging these values in the formula for Shapley Values, we obtain $C_{\mathbf{x}_1}^1 = 3$. Similarly we can obtain $C_{\mathbf{x}_2}^1 = C_{\mathbf{x}_3}^1 = 0$. We also provide code for the computation `solution10.py` (where indices are zero-based).

For the second logit we obtain:

$$C_{\mathbf{x}_1}^2 = 0, \quad C_{\mathbf{x}_2}^2 = 0, \quad C_{\mathbf{x}_3}^2 = 3.$$

A feature of Shapley values is that the attributions (C) add up to the original function value. For $g(\mathbf{x})$ both logits are 3 and they are only influenced by the first and last feature respectively thus we obtain $C_{\mathbf{x}_1}^1 = C_{\mathbf{x}_3}^2 = 3$ and zero for all other attributions. Nothing is attributed to \mathbf{x}_2 as the feature is discarded by the ReLU.

In the case of \mathbf{x}' , where both logits attain value 7, we get a similar picture, but now the second feature also influences the result:

$$\begin{array}{ll} C_{\mathbf{x}'_1}^1 = 3 & C_{\mathbf{x}'_1}^2 = 0 \\ C_{\mathbf{x}'_2}^1 = 0 & C_{\mathbf{x}'_2}^2 = 4 \\ C_{\mathbf{x}'_3}^1 = 4 & C_{\mathbf{x}'_3}^2 = 3 \end{array}$$

2. To compute the attribution we need to consider $2^{|P|-1}$ subsets per feature. And while a lot of the computations can be reused for different features this quickly becomes very expensive. Even for a single MNIST image ($28 \times 28 = 784$ pixels) this results in 784×2^{783} subsets to be considered, which without approximations is clearly intractable.

Problem 3 (Robust Features vs Non-Robust Features). (Adapted from [1])

Consider a distribution \mathcal{D} of data points $(\mathbf{x}, y) \sim \mathcal{D}$ where

$$y \stackrel{\text{u.a.r.}}{\sim} \{-1, 1\}, \quad \mathbf{x} \in \mathbb{R}^{d+1}, \quad \mathbf{x}_1 = \begin{cases} +y & \text{w.p. } p \\ -y & \text{w.p. } 1-p \end{cases}, \quad \mathbf{x}_2, \dots, \mathbf{x}_{d+1} \sim \mathcal{N}(\eta y, 1).$$

Here, $\mathcal{N}(\mu, \sigma^2)$ denotes the Gaussian distribution with mean μ and variance σ^2 . Further consider the family of linear classifiers $f_i(\mathbf{x}) := \text{sign}((\mathbf{w}^i)^T \mathbf{x})$. We use two different instantiation with $\mathbf{w}^1 := (0, \frac{1}{d}, \dots, \frac{1}{d})^T \in \mathbb{R}^{d+1}$ and $\mathbf{w}^2 := (1, 0, \dots, 0)^T \in \mathbb{R}^{d+1}$. In this example we will use $\eta = \frac{3}{\sqrt{d}}$ and $p = 0.9975$. Thus, we refer to the first component of \mathbf{x}_1 as a robust feature (as it is strongly associated with y), and $\mathbf{x}_2, \dots, \mathbf{x}_{d+1}$ as weak or non-robust features as for large d (e.g., $d \geq 100$) there is only weak correlation y and the individual features.

1. Compute the expected accuracy of f_1, f_2 . That is, determine value of $\mathcal{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [f(\mathbf{x}) = y]$.
2. Consider the data point $y = 1, \mathbf{x} = (1, 0.1, \dots, 0.1)^T$. Attempt to find an adversarial example \mathbf{x}' with $\|\mathbf{x} - \mathbf{x}'\|_\infty \leq 0.15$, such that it gets misclassified ($f(\mathbf{x}') = -1$). For both, f_1 and f_2 , find such an example or describe why it does not exist.

Solution 3.

1. We start by considering f_1 .

$$\begin{aligned} \mathcal{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [f_1(\mathbf{x}) = y] &= \mathcal{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [\text{sign}((\mathbf{w}^1)^T \mathbf{x}) = y] = \mathcal{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [y(\mathbf{w}^1)^T \mathbf{x} > 0] \\ &= \mathcal{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[y \sum_{i=1}^{d+1} \mathbf{w}_i^1 \mathbf{x}_i > 0 \right] = \mathcal{P}_y \left[y \sum_{i=1}^{d+1} \mathbf{w}_i^1 \mathcal{N}(\eta y, 1) > 0 \right] = \mathcal{P}_y \left[y \sum_{i=2}^{d+1} \frac{1}{d} \mathcal{N}(\eta y, 1) > 0 \right] \\ &= \mathcal{P}_y \left[\frac{y}{d} \sum_{i=1}^d \mathcal{N}(\eta y, 1) > 0 \right] = \mathcal{P}_y \left[\frac{y}{d} \mathcal{N}(d\eta y, d) > 0 \right] = \mathcal{P}_y \left[\frac{1}{d} \mathcal{N}(d\eta y^2, dy^2) \right] \\ &\stackrel{(*)}{=} \mathcal{P} \left[\frac{1}{d} \mathcal{N}(d\eta, d) > 0 \right] = \mathcal{P} \left[\mathcal{N}\left(\frac{d\eta}{d}, \frac{d}{d^2}\right) > 0 \right] = \mathcal{P} \left[\mathcal{N}\left(\eta, \frac{1}{d}\right) > 0 \right] = 1 - \mathcal{P} \left[\mathcal{N}\left(-\eta, \frac{1}{d}\right) \leq 0 \right] \\ &= 1 - \Phi \left(\frac{0 - (-\eta)}{\sqrt{\frac{1}{d}}} \right) = 1 - \Phi(\sqrt{d}\eta) = 1 - \Phi(3) = 0.9987 \end{aligned}$$

Here Φ denotes the CDF of the standard normal distribution. The equality $(*)$ holds as $y^2 = 1$ for $y \stackrel{\text{u.a.r.}}{\sim} \{-1, 1\}$, as we use, indicated by \mathcal{P}_y . Thus for $\eta = \frac{3}{\sqrt{d}}$ the accuracy $\mathcal{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [f_1(\mathbf{x}) = y]$ attains 0.9987 independent of d .

For f_2 :

$$\mathcal{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [f_2(\mathbf{x}) = y] = \mathcal{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [\text{sign}((\mathbf{w}^2)^T \mathbf{x}) = y] = \mathcal{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [\text{sign}(x_1) = y] = p = 0.9975$$

Thus the classifier f_2 , relying just on the first feature (x_1) attains slightly lower accuracy (0.9975) than the classifier f_1 relying on all the “weak” features ($\mathbf{x}_2, \dots, \mathbf{x}_{d+1}$), which attains 0.9987 independently of the value of d . So we can see that depending on η and p the classifier relying on non-robust features $\mathbf{x}_2, \dots, \mathbf{x}_{d+1}$ can be more accurate than the classifier relying just on the robust features.

2. Here we will apply FGSM to attempt to compute the adversarial examples. To this end we consider the gradient, given by $\nabla_{\mathbf{x}} \mathbf{w}^T \mathbf{x} = \mathbf{w}$.

Thus for f_1 we find $\mathbf{x}' = \mathbf{x} - 0.15 \text{sign}(\mathbf{w}^1) = \mathbf{x} - (0, 0.15, \dots, 0.15) = (1, -0.05, \dots, -0.05)$, which gets classified to class -1.

However, for f_2 the same process results in $\mathbf{x}' = (0.85, 0.1, \dots, 0.1)$, which still classifies to call 1. In fact, f_2 just considered the first component of \mathbf{x} , which has value 1. Since we can only change this by 0.15, we can not flip its sign. Thus we can not find an adversarial example for $\epsilon = 0.15$ for f_2 .

So, while f_1 has a higher standard accuracy, f_2 is more robust (and has a higher robust accuracy).

References

- [1] Dimitris Tsipras et al. “Robustness May Be at Odds with Accuracy”. In: *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL: <https://openreview.net/forum?id=SyxAb30cY7>.