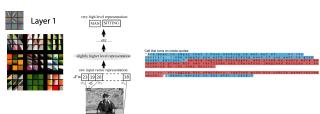


Week 7 - Concept-Based Explanations

Wednesday, March 5, 2020 9:02 PM

Activation Analysis

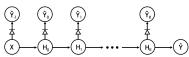
- Deep learning models are made from simple layers. When an input is passed through the network, each layer transforms the output of the previous layer. The intermediate outputs at a given layer are often called that layer's "activations." Each activation is associated with one neuron.
- There is a long, if informal, practice of relating model activations to properties of the data. For example, one of the earliest neural network explainability papers found that many neurons early on in the network activate when presented with specific orientations.



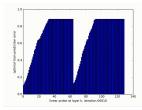
- One of the motivations behind building deep models was that the later layers in deep models learn meaningful abstractions. For example, neurons in early RNN language models were found to activate only on text that appear within a block of quotes. These explanations are arguably at a higher level than the orientation-based activation patterns seen in the first few layers of samples.
- The recent literature on concept-based explanations formalizes these early explorations. The appeal of the approach is that it helps bridge the vocabulary we already use to reason about a topic with the computational representations living within networks – it is an exercise in creating a shared language. If we let ourselves indulge in a little sci-fi, if we expect that one day we'll have machines capable of solving extremely complex problems, then the hope is that the methods descending from this literature will be to discover the underlying problem solving techniques.

Linear Probes

- A "linear probe" is just a linear model that predicts a feature of interest from intermediate model activations. If the model has high accuracy, then we can say this feature has been learned by the model in that layer.



- This can be useful for probing things about "How" models achieve their performance. In the toy example in their paper, they show that a skip connection across layers 1-14 prevents the model from learning anything useful within those layers. The implication is that while skip connections may be popular in many types of architectures, they might also lead to a lot of wasted computation and inadvertent shallowness if not attended to carefully (Click the image below to see how probe prediction accuracy only improves for the deep layers).



Network Dissection

- The network dissection paper advances over linear probes in two ways: (i) The authors realize that richly annotated datasets are a great source for probing what a model knows and (ii) they attempt to understand models at the level of individual neurons, not just layers.
- For this, the authors use the Brodatz "British" dataset of textures and objects. For each image in this dataset, each pixel is annotated with many complementary labels, some at several levels of resolution. There are 5 types of labels, which are easier to understand through an example.



They refer to the possible label values as concepts. E.g., "cottage", "sky", ... "lined" are concept

- The authors want to see whether a given neuron k is related to a concept c . To this end:

- For every input image x , compute the activations for the current neuron across all spatial patches of the image.
- Make a histogram of those activations across all inputs and patches. For each input, create a binary mask that selects patches whose activations are above a given quantile.
- Compute the interaction-over-units (IOU) of these binary masks with the concept's annotation mask, $\text{large}(b_k, v)$, which means that neuron k is highly active when concept c is present in a patch.



- The overall process makes it possible to create images like the one below. E.g., it seems that neuron 1410 in the ResNet model always activates when it sees a specific type of conical house.



- A skeptical reader would argue that these associations could be purely due to chance. We have searched so many neuron-concept pairs that some are bound to seem associated. If this is true, then the picture above would be overinterpreting something that has no inherent meaning (only large numbers of random associations). To address this concern, the authors propose to run an experiment where they apply a random rotation to all the activation vectors and rerun the dissection algorithm. They find that the IOUs in this case are much lower, which suggests that learned neurons really do have some meaning.
- One of the applications of this idea is to support interactive image generation. Network dissection can isolate concept-related neurons in generative image models, and by cleverly controlling the activations of those neurons, the user can interactively paint new objects in an existing generative image output. Click the link to see the demo.



Concept Activation Vectors

- Concept Activation Vectors have a similar motivation as Network Dissection – using human-annotated datasets to probe what a model knows. Their approach applies more generally (not just images) and requires less dense annotations. Their approach applies more generally (not just images) and requires less dense annotations. Their approach applies more generally (not just images) and requires less dense annotations. Their approach applies more generally (not just images) and requires less dense annotations.
- For this, the authors use the Brodatz "British" dataset of textures and objects. For each image in this dataset, each pixel is annotated with many complementary labels, some at several levels of resolution. There are 5 types of labels, which are easier to understand through an example.

- Step 1: Define a collection of images \mathcal{X} that represents a concept. Also construct a control pool of random images \mathcal{X}^* that represent other concepts. In this example, the concept is "stripes."

$$\mathcal{X}_c = \arg \min_{\mathcal{X}} \frac{1}{N} \sum_{n=1}^N L(x_n, x_n^*)^T c + b$$

$$x_n = \mathbb{1}_{\{x_n \in \text{True-Pos Images}\}}$$

$$h(x_n) / h(x_n^*)$$

$$L(x_n, x_n^*) = \left| \begin{array}{l} \{x_n : \text{class } k\} \\ \{x_n^* : \text{class } k'\} \end{array} \right|$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1 | \dots | c_m]^T f(x)$$

$$T(x^*) = [c_1^* | \dots | c_m^*]^T f(x^*)$$

$$T(x) = [c_1$$