

Semester Thesis: Measuring leakage in Concept-based methods

Mikael Makonnen, Moritz Vandenhirtz

November 19, 2024

1 Background & Motivation

Concept Bottleneck Models Concept bottleneck models (CBM) [Koh et al. \(2020\)](#); [Lampert et al. \(2009\)](#); [Kumar et al. \(2009\)](#) are a simple class of interpretable neural networks typically trained on data points $(\mathbf{x}, \mathbf{c}, \mathbf{y})$, comprising the covariates $\mathbf{x} \in \mathcal{X}$ and targets $\mathbf{y} \in \mathcal{Y}$ additionally annotated by the concepts $\mathbf{c} \in \mathcal{C}$. Consider a neural network f_{θ} parameterised by θ and a slice $\langle g_{\psi}, h_{\phi} \rangle$ [Leino et al. \(2018\)](#) s.t.

$$f_{\theta}(\mathbf{x}) = g_{\psi}(h_{\phi}(\mathbf{x})) \quad (1)$$

for all $\mathbf{x} \in \mathcal{X}$, where $\hat{\mathbf{y}} := f_{\theta}(\mathbf{x}) = g_{\psi}(h_{\phi}(\mathbf{x}))$ denote the output of the network, *i.e.* the predicted targets. CBMs enforce a concept bottleneck $\hat{\mathbf{c}} := h_{\phi}(\mathbf{x})$: the model’s final output depends on the covariates \mathbf{x} solely through the predicted concepts $\hat{\mathbf{c}}$. Thus, in addition to the target prediction loss applied to the final output, $h_{\phi}(\cdot)$ is trained to predict the ground-truth concept values.

Interpretability The interpretability of CBMs is achieved by the set of high-level, human-understandable concepts. Often, these are C binary-valued attributes, *i.e.* $\mathcal{C} = \{0, 1\}^C$ that can be easily detected from the covariates \mathbf{x} and are predictive of the targets \mathbf{y} . Although CBMs make no assumptions on (anti)causal relationships among \mathbf{x} , \mathbf{c} , and \mathbf{y} , they implicitly assume that concepts \mathbf{c} are a sufficient statistic [Yeh et al. \(2020\)](#) for predicting \mathbf{y} based on \mathbf{x} [Havasi et al. \(2022\)](#); [Marcinkevičs et al. \(2024\)](#), *i.e.* $\mathbf{y} \perp\!\!\!\perp \mathbf{x} \mid \mathbf{c}$.

Leakage Leakage is an instance of shortcut learning ([Geirhos et al., 2020](#)). [Margeloiu et al. \(2021\)](#); [Mahinpei et al. \(2021\)](#); [Havasi et al. \(2022\)](#) show that leakage occurs in cases where the conditional independence assumption does not hold. The distribution of the predicted concept values encodes more information than solely the probability of concept presence. This additional information can then be exploited by the classifier $g_{\psi}(\cdot)$. This is an issue since the predicted concept values encode information different from the human-understandable concepts, thus, prohibiting the interpretation of the predicted probability as probability of concept presence. [Mahinpei et al. \(2021\)](#) show that even if the predicted concepts are not soft (*i.e.* $\mathbf{c} \in [0, 1]$) but hard (*i.e.* $\mathbf{c} \in \{0, 1\}$), leakage happens, albeit weaker. Therefore, any perception of interpretability for standard CBMs is void if $\mathbf{y} \perp\!\!\!\perp \mathbf{x} \mid \mathbf{c}$ is not fulfilled, which is often the case in real-world problems. Examples of works (unintentionally) committing this fallacy are [Espinosa Zarlenga et al. \(2022\)](#); [Marconato et al. \(2022\)](#); [Ismail et al. \(2023\)](#). To understand how strongly the interpretability of concept probabilities is restricted, we need a metric that is able to measure the leakage within these concept embeddings.

2 Related Work

To measure leakage, [Zarlenga et al. \(2023\)](#) propose metrics that estimate the degree of excessive information with respect to other concepts, which they call impurity. To resolve leakage, [Margeloiu et al. \(2021\)](#) recommend using the *independent* training procedure with hard concepts. However, this comes at the cost of decreasing performance since the encoder and predictor head can not communicate anymore. Thus, [Havasi et al. \(2022\)](#) propose to include a hard side-channel, in which the additional information can be learned explicitly, as well as an autoregressive structure over the hard concept predictions, such that their correlations can be captured. At intervention time, they use importance-weighted MCMC sampling to implicitly learn the effect of a concept intervention on the other concepts. It will be interesting to see whether their approach fully eradicates leakage.

3 Methods

Consider a neural network NN_{θ} parameterised by θ and a slice $\langle g_{\psi}, h_{\phi} \rangle$ [Leino et al. \(2018\)](#) s.t.

$$NN_{\theta}(\mathbf{x}) = g_{\psi}(h_{\phi}(\mathbf{x})) \quad (2)$$

For sake of intuition, think of it as a CBM, where $\mathbf{z} = h_{\phi}(\mathbf{x}) = \hat{\mathbf{c}}$ is trained via the prediction of concepts, but this formulation allows for a more general interpretation.

What we are interested in for leakage, is the information contained within \mathbf{z} , which is informative for the label \mathbf{y} but independent/non-informative of concepts \mathbf{c} :

$$I(\mathbf{z}; \mathbf{y} \mid \mathbf{c}) = H(\mathbf{y} \mid \mathbf{c}) - H(\mathbf{y} \mid \mathbf{z}, \mathbf{c})$$

Estimating $H(\mathbf{y} \mid \mathbf{c})$ and $H(\mathbf{y} \mid \mathbf{z}, \mathbf{c})$ is the goal of this thesis. A straightforward approximation is

$$H(\mathbf{y} \mid \mathbf{z}, \mathbf{c}) = \mathbb{E}[-\log p(\mathbf{y} \mid \mathbf{z}, \mathbf{c})] \approx -\frac{1}{N} \sum_{i=1}^N \log g_{a,\psi}(h_{\phi}(\mathbf{x}_i), \mathbf{c}_i)_{y_i}, \quad (3)$$

$$H(\mathbf{y} \mid \mathbf{c}) = \mathbb{E}[-\log p(\mathbf{y} \mid \mathbf{c})] \approx -\frac{1}{N} \sum_{i=1}^N \log g_{b,\psi}(\mathbf{c}_i)_{y_i}, \quad (4)$$

where $g_{a,\psi}$ and $g_{b,\psi}$ are two classifiers trained to predict \mathbf{y} from \mathbf{z}, \mathbf{c} and from \mathbf{c} , respectively.