

神经元与概念

[Bau D, Zhou B, Khosla A, et al. Network Dissection: Quantifying Interpretability of Deep Visual Representations\[C\]// Computer Vision and Pattern Recognition. IEEE, 2017:3319-3327.](#)

量化神经元 k 与概念 c 之间的关系：

1. 输入 x ，前向传播至 k 时输出 $A_k(x)$ ，插值得原输入大小的 $S_k(x)$ ，以阈值 T_k 得掩码
 $M_k(x) = I(S_k(x) \geq T_k)$
2. 每个输入 x 人工标注概念 c 的掩码 $L_c(x)$
3. 用交并比 $IoU_{k,c} = \frac{\sum_x M_k(x) \cap L_c(x)}{\sum_x M_k(x) \cup L_c(x)}$ 反映神经元 k 与概念 c 的关系

- What is a disentangled representation, and how can its factors be quantified and detected?
- Do interpretable hidden units reflect a special alignment of feature space, or are interpretations a chimera?
- What conditions in state-of-the-art training lead to representations with greater or lesser entanglement?

emergent interpretability is an axis-aligned property of a representation that can be destroyed by rotation without affecting discriminative power.

1. Identify a broad set of human-labeled visual concepts.
2. Gather hidden variables' response to known concepts.
3. Quantify alignment of hidden variable–concept pairs.