

# Before using cosine similarity between label representations consider:

- a) Which classifier your representations are coming from and
- b) Whether cosine similarity can answer your question



## What Cosine Similarity of Label Representations Can and Cannot Tell us

### Sigmoid classifier

$$p(y_v | \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{g}(v)^\top \mathbf{f}(\mathbf{x})}}$$

E.g. Word2Vec, multi-label classifier, labels can be present or not present.

Possible label combinations are captured by the **Sign Pattern**,  $\mathbf{W} \in \mathbb{R}^{k \times d}$  has  $\mathbf{g}(v)$  as rows

$$\mathcal{S}(\mathbf{W}) = \{\text{sign}(\mathbf{W}\mathbf{z}) : \mathbf{z} \in \mathbb{R}^d, \mathbf{w}_i^\top \mathbf{z} \neq 0, 1 \leq i \leq k\}$$

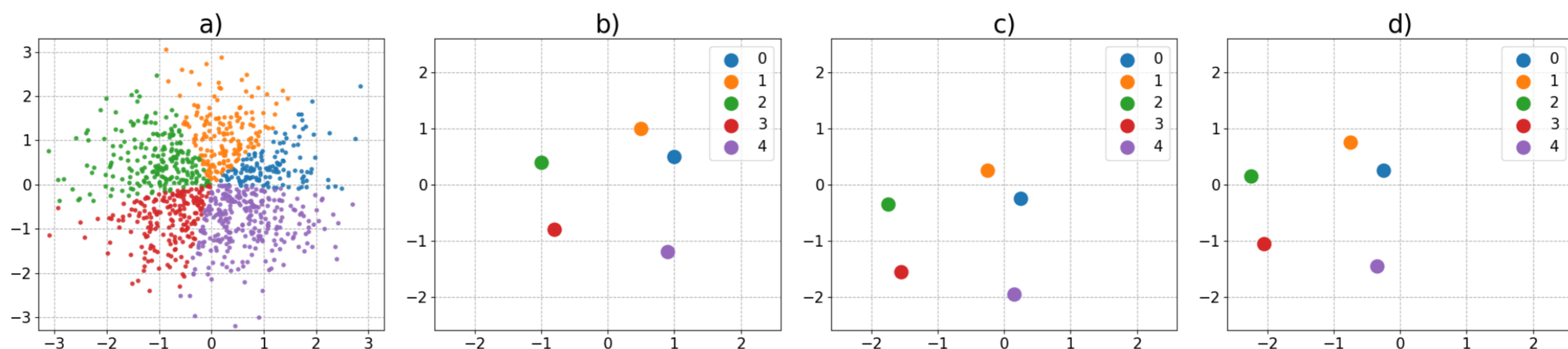
### Softmax classifier

$$p(y | \mathbf{x}) = \frac{\exp(\mathbf{f}(\mathbf{x})^\top \mathbf{g}(y))}{\sum_{y' \in \mathcal{Y}} \exp(\mathbf{f}(\mathbf{x})^\top \mathbf{g}(y'))}$$

E.g. Language model, image classifier, gives a ranking of labels.

Possible rankings are captured by the **Ranking Pattern**

$$\mathcal{R}(\mathbf{W}) = \left\{ \text{sign}(\mathbf{B}\mathbf{W}\mathbf{z}) : \begin{array}{l} \mathbf{z} \in \mathbb{R}^d, \\ (\mathbf{w}_i - \mathbf{w}_j)^\top \mathbf{z} \neq 0, \\ 1 \leq i < j \leq k \end{array} \right\}$$



Above, 4 softmax models with the same probabilities, but different cosine similarities between unembeddings.

Pairwise cosine similarity of label representations (unembeddings)

**cannot tell us** the probabilities (Thm 3.4)

**can tell us** the sign pattern (Lem 5.3.)

**cannot tell us** the ranking pattern (Thm 5.5)

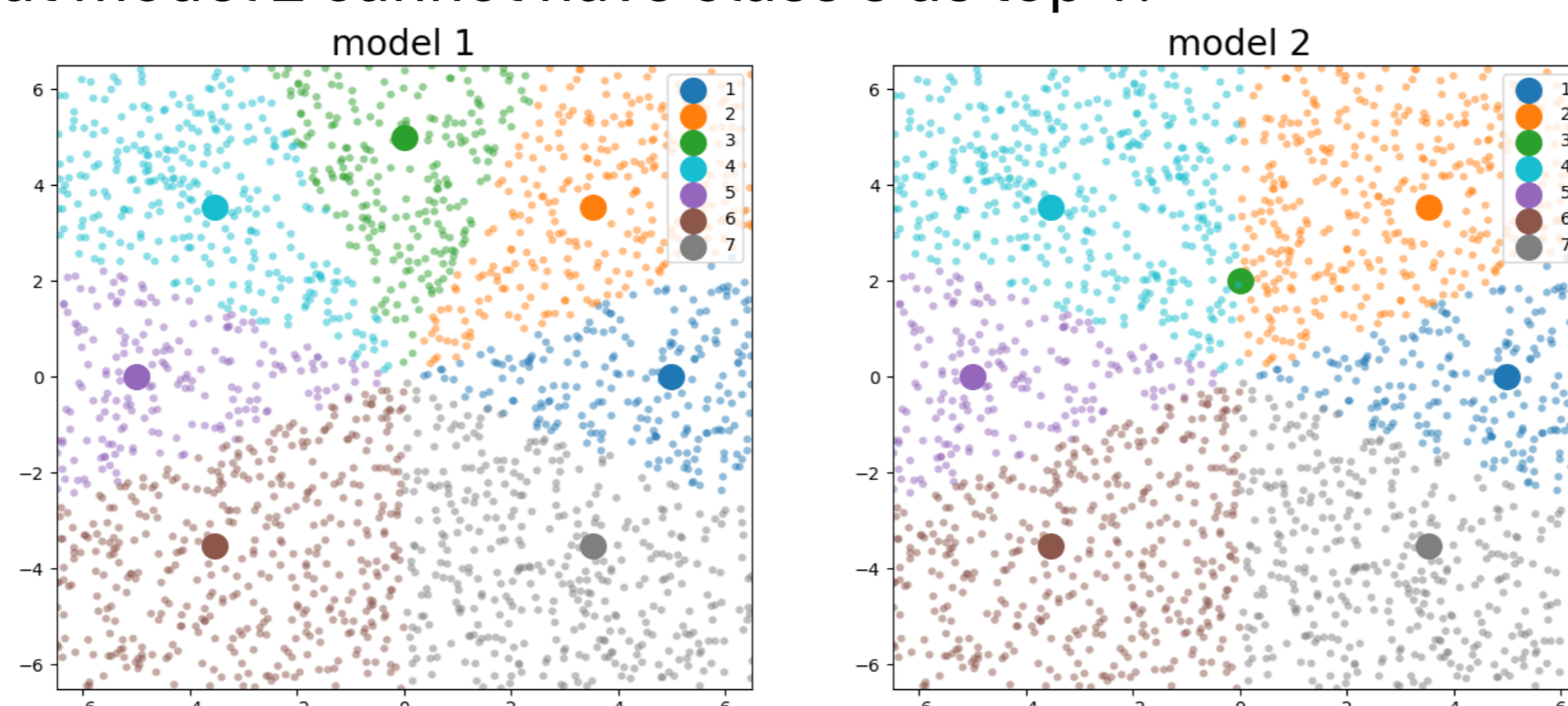
Cosine similarities of differences of unembeddings

**can tell us** the ranking pattern (Thm 5.6)

Below, Cosine similarities of unembeddings are the same, but model 2 cannot have class 3 as top 1.



Arxiv



Code