

It may be possible to recognize a specific object class even without having seen an image of that object, if the text describes the object well enough. For example, having read that a cat has four legs and pointy ears, the learner might be able to guess that an image is a cat, without having seen a cat before.

Zero-data learning (Larochelle *et al.*, 2008) and zero-shot learning (Palatucci *et al.*, 2009; Socher *et al.*, 2013b) are only possible because additional information has been exploited during training. We can think of the zero-data learning scenario as including three random variables: the traditional inputs  $\mathbf{x}$ , the traditional outputs or targets  $\mathbf{y}$ , and an additional random variable describing the task,  $T$ . The model is trained to estimate the conditional distribution  $p(\mathbf{y} \mid \mathbf{x}, T)$  where  $T$  is a description of the task we wish the model to perform. In our example of recognizing cats after having read about cats, the output is a binary variable  $y$  with  $y = 1$  indicating “yes” and  $y = 0$  indicating “no.” The task variable  $T$  then represents questions to be answered such as “Is there a cat in this image?” If we have a training set containing unsupervised examples of objects that live in the same space as  $T$ , we may be able to infer the meaning of unseen instances of  $T$ . In our example of recognizing cats without having seen an image of the cat, it is important that we have had unlabeled text data containing sentences such as “cats have four legs” or “cats have pointy ears.”

Zero-shot learning requires  $T$  to be represented in a way that allows some sort of generalization. For example,  $T$  cannot be just a one-hot code indicating an object category. Socher *et al.* (2013b) provide instead a distributed representation of object categories by using a learned word embedding for the word associated with each category.

A similar phenomenon happens in machine translation (Klementiev *et al.*, 2012; Mikolov *et al.*, 2013b; Gouws *et al.*, 2014): we have words in one language, and the relationships between words can be learned from unilingual corpora; on the other hand, we have translated sentences which relate words in one language with words in the other. Even though we may not have labeled examples translating word  $A$  in language  $X$  to word  $B$  in language  $Y$ , we can generalize and guess a translation for word  $A$  because we have learned a distributed representation for words in language  $X$ , a distributed representation for words in language  $Y$ , and created a link (possibly two-way) relating the two spaces, via training examples consisting of matched pairs of sentences in both languages. This transfer will be most successful if all three ingredients (the two representations and the relations between them) are learned jointly.

Zero-shot learning is a particular form of transfer learning. The same principle explains how one can perform **multi-modal learning**, capturing a representation