

# Semi-Supervised learning

*Mokshad Ketan Sankhe*

Semi-supervised learning is a class of machine learning techniques that make use of both labeled and unlabeled data during training. It falls between unsupervised learning, which does not require labeled data, and supervised learning, which relies solely on labeled data. Semi-supervised learning is useful in situations where labeled data is limited or expensive to obtain, but unlabeled data is readily available.

**Working:** Semi-supervised learning algorithms first use the unlabeled data to build a model that captures the basic structure of the data. This provides information about the shape and distribution of the input space. Next, the small labeled dataset is used to further improve the model. There are several common semi-supervised techniques:

- Self-training: The model makes predictions on unlabeled data, adds the most confident predictions to the labeled set, and retrains itself.
- Generative models: These estimate the joint probability distribution between inputs and outputs in order to generate more labeled data.
- Graph-based methods: These build a graph that connects similar data points, using both labeled and unlabeled data.
- Co-training: Uses two separate models to label unlabeled data for each other.

**Applications:** Semi-supervised learning is well-suited for classification and regression problems where labeled data is scarce. Key applications include:

1. Natural language processing: Text classification, sentiment analysis, etc.
2. Speech recognition: Transcribing audio when transcripts are limited.
3. Image classification: Identifying objects in images with few labeled examples.
4. Information retrieval: Document clustering, search results ranking.
5. Bioinformatics: Protein structure prediction with insufficient labeled data.

## **Advantages:**

- a) Increased model accuracy with small labeled dataset
- b) Avoids overfitting by using abundant unlabeled data
- c) Labeled data can be used more efficiently

## **Disadvantages:**

- a) Heavily depends on assumptions about data distribution
- b) Less robust compared to supervised learning
- c) Unlabeled data must be useful and representative

**Conclusion:** Semi-supervised techniques are becoming increasingly important with the rise of big data. As labeled data is scarce and expensive in many domains, semi-supervised learning provides a vital toolkit to train accurate models by fully exploiting both labeled and unlabeled data. With further research, semi-supervised methods could become as widely used as supervised learning.