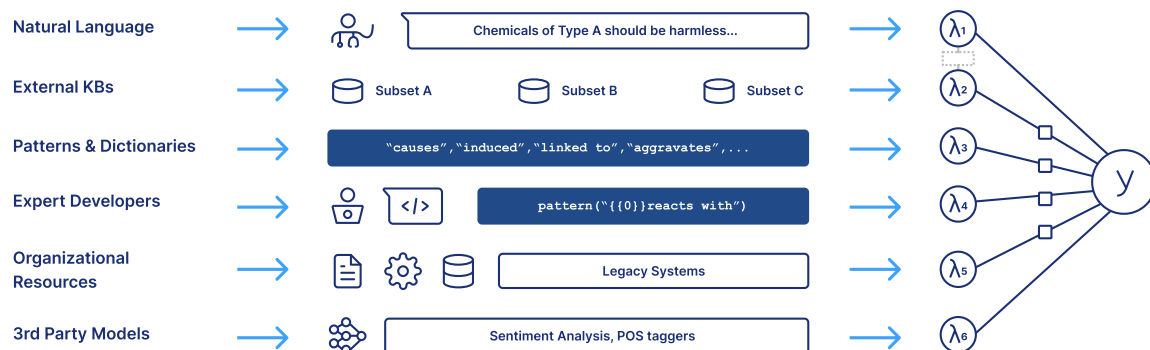


What is Weak Supervision and How Does Weak Supervision Work?

Weak supervision is an approach to machine learning in which high-level and often noisier sources of supervision are used to create much larger training sets much more quickly than could otherwise be produced by manual supervision (i.e. labeling examples manually, one by one).

If you have high-level, scalable, but potentially noisy sources of signal, you can combine them using multiple sources of supervision. At Snorkel AI, we use labeling functions to do this.

By observing when and where these different labeling functions agree or disagree with one another, you can automatically learn—in unsupervised ways—when, where, and how much to trust each of them. You can thus learn their areas of expertise, and the overall level of expertise, so that when you combine their votes you end up with the highest quality label possible for each data point.



Weak Supervision Interfaces with Snorkel

Weak Supervision Interfaces with Snorkel

Read this article [>](#)

When Should Weak Supervision Be Used?

Weak supervision enables the creation of very large training sets very quickly. If your particular problem would be better addressed with 100,000 “pretty good” labels, compared to 100 “perfect” labels, it may be worth looking at higher-level interfaces for gathering more data.

Additionally, weak supervision is great to use in any situation in which you need to adapt and iterate regularly and rapidly. If there are frequent shifts in the distribution of your data, such as in an adversarial setting (such as fraud detection) or just because your needs frequently change, weak supervision enables you to do anything from adding novel classes to incorporating and reflecting new realities about your problem.

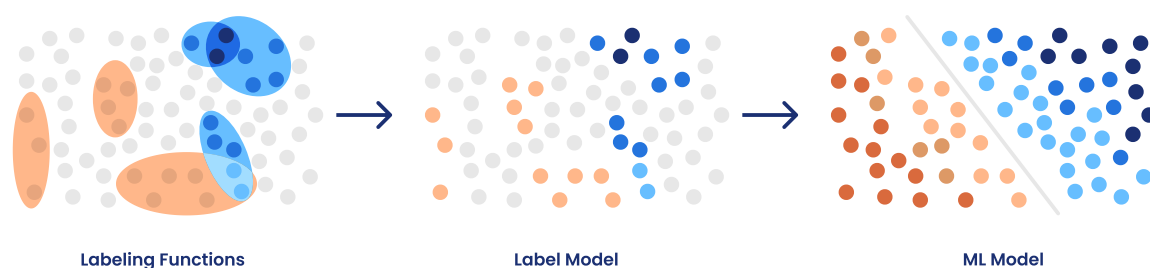
Weak Supervision vs. Rule-Based Classifiers

Weak supervision has some similarities—and some very important differences—to rule-based classifiers. The obvious similarity is that the inputs to each look like rules (i.e., simple functions that output labels or predictions). The important difference between them is that the rule-based classifier stops there—the rules are the classifier. Such systems are generally brittle because they do not generalize to other examples, even ones that are very similar to those that are labeled by one or more rules.

With weak supervision, on the other hand, the rules (or “labeling functions”) are used to create a training set for a machine-learning-based model. That model can be much more powerful, utilize a much richer feature set, and take advantage of other state-of-the-art techniques in machine learning, such as transfer learning from foundation models.

Each labeling function suggests training labels for multiple unlabeled data points, based on human-provided subject matter expertise. A label model (Snorkel Flow includes multiple variants optimized for different problem types) aggregates those weak labels into one training label per data point to create a training set. The ML model is trained on that training set and learns to generalize beyond just those data points that were labeled by labeling functions.

As a result, the model is generally much more robust than a corresponding rule-based classifier.



Weak Supervision using Label Model to train ML Model

How Snorkel Flow Makes Weak Supervision Practical

Snorkel AI has applied weak supervision to many problems over the years and we have learned a lot about which features and workflows make it most accessible and practical for users. We built Snorkel Flow specifically with that experience in mind.

Snorkel Flow is a data-centric platform for building AI applications powered by weak supervision and other modern machine learning techniques. In Snorkel Flow, users manage data throughout the full AI lifecycle by writing simple programs (labeling function) to label, manipulate, and monitor training data. These programmatic inputs are modeled and integrated using theoretically-grounded statistical techniques, made

Snorkel Flow provides you with the ability to easily express many different types of signal, whether that is importing existing labels or models that you already have and applying them, or allowing you to write new labeling functions that are rule- or heuristic-based. Snorkel Flow then gives you access to the label-model algorithms that we have developed. They automatically combine these different scalable (but potentially noisy) sources of supervision to create high-quality labels for each of your data points.

Lastly, Snorkel Flow equips you with ready-made infrastructure for the application of these functions. The platform guides you through the process and supplies integrated model training so

accessible to both developer and non-developer users alike via both a no-code UI and Python SDK.

that you can loop back and make adjustments yourself to your weak supervision sources as you go.

Use Cases for Weak Supervision

Weak supervision can be applied to many problems. The Snorkel AI team has applied it to text data (long and short), conversations, time series, PDFs, images, videos, and more. So long as domain-relevant resources exist or labeling heuristics can be described, weak supervision can be applied. Some of the use cases include:

Some of the use cases include:

- Text and document classification
- Information extraction from unstructured text, PDF, HTML and more
- Rich document processing
- Structured data classification
- Conversational AI and utterance classification
- Entity linking
- Image and cross-modal classification
- Time series analysis
- Video classification