Snorkel DryBell: A Case Study in Deploying Weak Supervision at Industrial Scale

Stephen H. Bach, Daniel Rodriguez, Yintao Liu, Chong Luo, Haidong Shao, Cassandra Xia, Souvik Sen, Alexander Ratner, Braden Hancock, Houman Alborzi, Rahul Kuchhal, Christopher Ré, Rob Malkin

(4)

(5)

Motivation



- Organizations using machine learning often manage many different training data sets
- These training sets are **costly to create** and need to be updated as strategy changes

How can we **replace hand-labeled training** data with existing organizational resources?

Case Studies at Google

Product Classification

- Existing classifier used to detect products in a category of interest
- Goal: move accessories to positive class
- Instant depreciation of investment in labels!

Topic Classification

Emerging topic of interest in Web content







Google

Stanford

Can we **transfer knowledge** from resources not servable in production environments?

- Goal: develop new classifier to identify topic
- Default procedure is to collect hundreds of thousands of labels for new topic!



www

Old:

Products

Architecture



Experiments

Ve first compare		Products		Topics		
Snorkel DryBell with training on the		Rel. F1	Lift	Rel. F1	Lift	
validation data (~10k examples)	Train on Val. Data	100%		100%		
and using the	Generative Model	103%	+3%	94%	-6%	
generative model to make predictions	Snorkel DryBell	105%	+5%	118%	+18%	



- Snorkel DryBell's generative model estimates the accuracies of the labeling functions and the true labels without access to ground truth labeled data
- 5. The estimated labels are used to train production machine learning systems

Example Code

In this example, we write a labeling function for classifying whether Web content is related to celebrity news, using a natural language processing (NLP) classifier as a knowledge resource.



Generative Modeling



Snorkel DryBell combines the votes of labeling functions using the modeling framework of Snorkel (Ratner et al., VLDB 2017).

accuracy) of the labeling functions by

maximizing the log likelihood of their

 $\operatorname{arg\,max} \log p_{\boldsymbol{\theta}}(\boldsymbol{\lambda})$

More info.: <u>http://snorkel.stanford.edu</u>

observed outputs:



Conclusions

- Useful resources for weak supervision are ulletabundant. Further, many open-source analogs exist, suggesting directions for future work.
- **Rethink systems for large-scale weak** supervision. Snorkel was originally designed for novice users, but we find that expert users need more flexibility and modularity.
- Cross-feature transfer is essential. Many resources that can be used for background knowledge are not servable. Using them as supervision enables transfer to classifiers.



