

# Accelerating Machine Learning with Training Data Management

Alex Ratner

Stanford University

# Training data is the key ingredient in ML



# But it's created and managed in manual, ad hoc ways

#### **KEY RESEARCH QUESTION**

# Can we add mathematical & systems structure to the way people build & manage training sets today?



# Running Example: Chest X-Ray Triage



#### Motivation: Case prioritization for e.g. lowresource hospitals

[Dunnmon et. al., Radiology 2018; Dunnmon & Ratner et. al., 2019; Khandewala et. al., NeurIPS ML4H 2017]

# Running Example: Chest X-Ray Triage



± 1 point due to model choice

#### Model dev is often radically easier today!

[Dunnmon et. al., Radiology 2018; Dunnmon & Ratner et. al., 2019; Khandewala et. al., NeurIPS ML4H 2017]

(All scores: ROC AUC)



# Running Example: Chest X-Ray Triage



 $\pm$  9 points due to training set size

 $\pm$  8 points due to training set quality

± 1 point due to model choice

#### Training data is often the key differentiator

[Dunnmon et. al., Radiology 2018; Dunnmon & Ratner et. al., 2019; Khandewala et. al., NeurIPS ML4H 2017]

(All scores: ROC AUC)

# Challenges of Training Data Management

- *Volume* is critical
  - But training data is largely hand-labeled: slow & expensive
- Quality is critical
  - But this is challenging to assess
- *Flexibility* is critical
  - But training sets are completely static





#### **Our research: building systems that**

1 Let users specify training sets in higher-level, programmatic ways

2) Clean and integrate this input

+ + + + +

3) Use as training data for ML models

A new way to specify ML models---in hours rather than months



### This talk: Three systems that support and accelerate critical steps of training data creation & management



















#### Problem: Hand-labeling is slow, expensive, & static

Idea: Enable users to label training data programmatically



#### **KEY TECHNICAL IDEA:**

View training set labeling as a *noisy programmatic process* that we can model



# The Snorkel Pipeline

snorkel.stanford.edu







Users write *labeling functions* to heuristically label data



Snorkel *cleans and combines* the LF labels



The resulting training database used to train an ML model

#### **Note: No hand-labeled training data!**

# Snorkel: Real-World Deployments



#### In many cases: From person-months of handlabeling to hours





to heuristically label data



combines the LF labels

training database used to train an ML model



# (1) Writing Labeling Functions

	def LF_short_report(x): if len(X.words) < 15: return "NORMAL"	
	<pre>def LF_off_shelf_classifier(x):  if off_shelf_classifier(x) == 1:      return "NORMAL"</pre>	
	def LF_pneumo(x): if re.search(r'pneumo.*', X.text): return "ABNORMAL"	
	def LF_ontology(x): if DISEASES & X.words: return "ABNORMAL"	

LABELING FUNCTIONS

Labeling function:

 $\lambda: \mathcal{X} \mapsto \mathcal{Y} \cup \{0\}$ Data Labels Abstain

A simple abstraction for expressing domain heuristics or other noisy label sources



# Simple Example: Pattern Matching

"Indication: Chest pain. Findings: Focal consolidation and pneumothorax."



def LF\_pneumo(x):
 if re.search(r'pneumo.\*', X.text):
 return "ABNORMAL"



Labeling functions (LFs) are simple UDFs for expressing domain expertise



# Simple Example: Pattern Matching



# LFs can also be noisy---we can estimate their accuracies to handle this (next)



# A Simple Formalism for Weak Supervise Strategies

def LF pneumo(x):

- Pattern matching
- Distant supervision

Indication: Chest

pain. Findings: ocal consolidation

Functions of features



if re.search(r'pneumo.*', X.text): return "ABNORMAL"		
def LF_ontology(x): if DISEASES & X.words: return "ABNORMAL"		
def LF_short_report(x): if len(X.words) < 15: return "NORMAL"		

[e.g. Hearst 1992, Zhang 2017]

e.g.	Mintz	2009]
------	-------	-------

[e.g. Varma 2017]

And many others: crowdsourcing, other models, etc.

def LF circular mass(x):

return "ABNORMAL"

if c.radius > 1:

c = off\_shelf\_circle\_finder(x)[0]



# Result: Supervision as Code

	def LF_short_report(x): if len(X.words) < 15: return "NORMAL"
	def LF_off_shelf_classifier(x): if off_shelf_classifier(x) == 1: return "NORMAL"
	def LF_pneumo(x): if re.search(r'pneumo.*', X.text): return "ABNORMAL"
	def LF_ontology(x): if DISEASES & X.words: return "ABNORMAL"
	LABELING ELINCTIONS

But, very messy supervision...



# Challenges of Supervision as Code



- Different unknown
  accuracies
- Different unknown
  correlations
- No ground truth

A new type of data cleaning and integration problem





*labeling functions* to heuristically label data



Snorkel cleans and combines the LF labels



The resulting training database used to train an ML model



#### How can we do this without ground-truth labels?



# Simple generative model of the labeling process

- Represent the LF outputs as RVs
- Model with a single parameter
- Can be extended

#### Include pairwise dependencies

- Assume edges are known
- We (provably) estimate from unlabeled data [ICML'17, Arxiv'19]

#### This talk: How to learn model without observing *Y*?



# Prior Work on Weak Supervision Modeling

#### Some prior work:

- 1. Crowdsourcing
  - 1. EM-based [e.g. Dawid & Skene, 1979]
  - 2. Spectral [e.g. Gosh, 2011; Anandkumar 2012]
- 2. Data fusion [e.g. Dong, 2015; Rekatsinas 2017]
- *3. Others (see snorkel.stanford.edu)*

Differences highlighted in this portion of the talk:

- 1. Complex dependencies between LFs
- 2. End-to-end theoretical guarantees







Key idea: Learn from the agreements & *disagreements* between the LFs

[Ratner et. al., AAAI '19] [Ratner et. al., NeurIPS '16]





 $\lambda_4$ 

 $\lambda_3$ 



This encodes the observed LF agreements / disagreements







Idea: Use graph-sparsity of  $\Sigma^{-1}$ 



We know the zeros of  $\Sigma^{-1}$  from our model [Loh & Wainwright 2013, Ratner 2019]

This encodes our knowledge of the dependency structure

Observed

Unobserved

snorke





Let  $\Omega$  be the set of 0 entry indices in  $\Sigma^{-1}$ 

Observed

snorke





 $(\Sigma_{0}^{-1} + zz^{T})_{0} = 0$ 

This is similar to a matrix completion problem!



# Result: Recovering the LF accuracies & correlations

Why is this nice?

- Simple to optimize: E.g. SGD
- Scalable: No dependence on n!
- Theoretical guarantees: We can leverage random matrix & perturbation tools to prove convergence



Let:

- *n* = number of *unlabeled* data points
- *d* = number of LF cliques

# Recovery Results (Informal)

Given:

- n *unlabeled* data points
- A set of LFs that are *on average* better than 50% accurate
- A sufficiently sparse dependency structure (per deterministic test)


## Recovery Results (Informal)

Given:

- n *unlabeled* data points
- A set of LFs that are *on average* better than 50% accurate
- A sufficiently sparse dependency structure (per deterministic test)

Then:  

$$E[\|\hat{z} - z^*\|] = O\left(\frac{1}{\sqrt{n}}\right)$$

Parameter (LF accuracy & correlation strength) estimation error

**Decreases with** *unlabeled* data!



# Result: Recovering the LF accuracies and correlations



This gives us a simple---and provably consistent---way to clean and integrate the LF outputs





Key question: How do we communicate the lineage (quality) of the training labels?



## Ex: Importance of Label Lineage



# Result: Average training label quality barely better than 60% accuracy



# Solution: Modify Loss to Use *Probabilistic Labels*

- Standard ERM:
  - Minimize  $\frac{1}{n} \sum_{i=1}^{n} l(x^{(i)}, y^{(i)})$
- We use a noise-aware loss:

• Minimize  $\frac{1}{n} \sum_{i=1}^{n} E_{\tilde{y} \sim p_{z}(\cdot|\lambda^{(i)})}[l(x^{(i)}, \tilde{y})]$ 



A simple tweak to the loss function to communicate lineage / quality!

### **Recap: The Snorkel Pipeline**





Users write *labeling functions* to heuristically label data



Snorkel cleans and combines the LF labels



The resulting training database used to train an ML model

## End-to-End Recovery Results (Informal)

Result: Given conditions from before, and some loose assumptions about the end model, generalization error decreases at the same rate

$$E[||l_{\widehat{W}} - l_{W^*}||] = O\left(\frac{1}{\sqrt{n}}\right)$$
  
Expected test error of the end model

**Decreases with unlabeled data!** 

### Same asymptotic rate as with labeled data!

[Ratner et. al., AAAI 2019; Ratner et. al., NeurIPS 2016]

## Question: Why train a final model at all?



## Question: Why train a final model at all?

### • (1) Generalization

- Often hard to write good, high-coverage LFs
- We can leverage commodity ML models and tools to do better! [VLDB 2018]
- (2) Cross-modal transfer
  - Write LFs over one feature set → train model over a different one





The resulting training database used to train an ML model

## Highlight: Cross-Modal Transfer





# Training data as a medium of transferring domain knowledge across modalities







def LF_ if "p r def LF_ if rej	pleural_effusion leural effusion eturn "ABNORMAL" normal_report(c len(NORMAL_TERMS ort_words)) > t	<pre>n(c): " in c.report.text: " , thresh=2):</pre>
def LF_ if rej	normal_report(c len(NORMAL_TERMS ort_words)) > t	, thresh=2):
-	eturn "NORMAL"	S.intersection(c. ;hresh:
India Media <b>norma</b> with conso <b>pleu</b> acute abnor	ation: Chest stinal conto limits. He n normal lin lidation, pr al effusion cardiopulmo mality.	t pain. Findings: ours are within eart size is nits. No focal neumothorax or . Impression: No onary

### 20 Labeling Functions



<pre>def LF_pneumothorax(c):     if re.search(r'pneumo.*', c.report.text):         return "ABNORMAL"</pre>
<pre>def LF_pleural_effusion(c):     if "pleural effusion" in c.report.text:         return "ABNORMAL"</pre>
<pre>def LF_normal_report(c, thresh=2):     if len(NORMAL_TERMS.intersection(c.       report.words)) &gt; thresh:       return "NORMAL"</pre>
Indication: Chest pain. Findings: Mediastinal contours are within normal limits. Heart size is within normal limits. No focal consolidation, pneumothorax or pleural effusion. Impression: No acute cardiopulmonary abnormality.

#### 20 Labeling Functions

### Applications: Diversity and Real-World Impact



Stanford MEDICINE





### **Medical Monitoring**

### (image, video, time series)

[Dunnmon & Ratner, 2019] [Khandwala, NeurIPS ML4H 2017] [Fries, 2018]

### **Knowledge Base Construction**

### (text, tables, PDFs, HTML)

[Ratner, VLDB 2018] [Wu, SIGMOD 2018] [Kuleshov, NeurIPS ML4H 2016]

#### Industry

#### (web, text, other)

[Bach, SIGMOD Industry 2019] [Mallinar, AAAI 2019] [Bringer, 2019]

### In many cases: exceeds the efficacy of personmonths of labeled data in hours

## Highlight: Scaling with unlabeled data



Takeaway: Add more *unlabeled* data---without changing the LFs---and get better end performance!

### NIH Snorkel Workshop and User Study





# Novice users: 45.5% better on average using Snorkel vs. hand-labeling

### Snorkel: A System for Rapidly Creating Training Sets



# Snorkel can enable a more *accessible*, *faster*, and *powerful* way of building ML applications

### Our Research: Training Data Management Systems



### Our Research: Training Data Management Systems



### One Critical Tool: Data Augmentation



Ex: 13.4 pt. avg. accuracy gain from data augmentation across top ten CIFAR-100 models

### Problem: Data augmentation is *critical*, but hard to hand-tune

Idea: Users provide transformations which we automatically tune and compose



### Automatic Data Augmentation from User-**Specified Invariances**





Jsers write transformation functions (TFs)

### TFs can express a diverse range of invariances



### How do we tune & compose these?

### Automating Tuning & Composing

### Idea #1: Treat this as a sequence modeling problem



# How do we generate diverse but realistic transformed images?

## Automating Tuning & Composing



Idea #2: Use adversarial approach to learn to generate realistic images from unlabeled data

## Automatic Data Augmentation from User-**Specified Invariances**



transformation functions (TFs)



model to tune & compose the TFs

64

end model

### Empirical Results: Gains Across Domains

- Gains over heuristic data augmentation approaches:
  - 4 pts. in accuracy on CIFAR-10
  - 1.4 F1 score pts. on a text relation extraction problem
  - 3.8 pts. in accuracy on a clinical mammography classification task
- Our core ideas have since been adopted in industry:
  - Google's AutoAugment, yielding new SOTA on Imagenet



### Our Research: Training Data Management Systems



### Our Research: Training Data Management Systems



# Most real-world settings: multiple related modeling tasks



# How do we approach these multiple related modeling tasks?

### Basic approach: T pipelines for T tasks



### For ex: three separate Snorkel pipelines

### Problem: We have to provide supervision (write LFs) for multiple tasks

Idea: Jointly model across multiple related tasks to do better with less



### Basic approach: T pipelines for T tasks



### Snorkel MeTaL: A System for Multi-Task Supervision




Use cross-task agreements / disagreements in an extended version of the matrix completion approach

### Snorkel MeTaL: A System for Multi-Task Supervision



Empirical results: Strong gains over single-task approach (avg. 4 F1 points) and easier interface

#### Our Research: Training Data Management Systems



### Research Agenda

# Make real-world ML applications radically easier and faster to build



## with data management systems that support critical steps outside of the model

# Building up the code-as-supervision





From supervision as labels to *supervision as code* 

78

# Building up the code-as-supervision



Goal: Move up the stack- make ML radically easier to program

Models

## Formalize and support other critical "preprocessing" steps of ML



Goal: Build data management systems that accelerate where ML developers actually spend their time



The Massively Multi-Task Ecosystem

**Today:** 



As it becomes faster to build training sets there will be *tens to hundreds of interacting models* 

New challenges:

- Incremental maintenance
- Handling complex data dependencies
- New formalisms

Goal: Support new model ecosystems at massive scale

### Our Research: Training Data Management Systems



## Key Idea: Add mathematical & systems structure to training data creation & management

### Our Research: Training Data Management Systems



**Thank you to:** Chris Ré, Daniel Rubin, Kunle Olukotun, John Duchi, Chris De Sa, Sen Wu, Daniel Selsam, Henry Ehrenberg, Jason Fries, Bryan He, Braden Hancock, Theo Rekatsinas, Paroma Varma, Fred Sala, Jared Dunnmon, the Stanford Bio-X SIG Fellowship, and the many other contributors, users, and sponsors of Snorkel



