Liger: Fusing weak supervision with foundation model embeddings

Mayee F. Chen*, Daniel Y. Fu*, Dyah Adila, Michael Zhang, Fred Sala, Kayvon Fatahalian, Christopher Ré

Snorkel AI MLW, April 8th, 2022







Labeled data 🔊

Deep neural network trained from scratch





Labeled data 🔊

Deep neural network trained from scratch

Both: expensive and time-consuming!





Deep neural network trained from scratch

Unlabeled data







Weak Supervision!

=



Unlabeled data +

Noisier sources of signal: crowd workers, heuristics, external KBs, etc.



Weak Supervision!



Large, pre-trained models (BERT, GPT-3, DALL-E 2, CLIP)



Unlabeled data +

Noisier sources of signal: crowd workers, heuristics, external KBs, etc.



Weak Supervision!



Large, pre-trained models (BERT, GPT-3, DALL-E 2, CLIP)

Foundation models (FMs)*

=

*Bommasani et al. "On the opportunities and risks of foundation models." 2021.

Q: How do we combine Weak Supervision and Foundation Models?

Outline

- 1. Background
 - a. The weak supervision pipeline
 - b. Foundation models
- 2. Technical review of Weak Supervision
- 3. Liger 🦁 🐯
- 4. Why it works: embedding "smoothness"
- 5. Results
- 6. Summary and Future Directions



 $x\in \mathcal{X}$



Unlabeled dataset

Unknown labels $y \in \{-1,1\}$



def L_1: SPAM if "check out" def L_2: NOT SPAM if "love" def L_3: SPAM if "subscribe"

Unlabeled dataset

Unknown labels $y \in \{-1,1\}$

Users write labeling functions (LFs)/ Specify weak sources





def L_1:
 SPAM if "check out"

def L_2:
 NOT SPAM if "love"

def L_3:
 SPAM if "subscribe"



Unlabeled dataset

Unknown labels $y \in \{-1,1\}$

Users write labeling functions (LFs)/ Specify weak sources

Learn model parameters (latent variable estimation)



Unlabeled dataset

Unknown labels $y \in \{-1,1\}$

Users write labeling functions (LFs)/ Specify weak sources

Learn model parameters

Output probabilistic labels



Foundation Models (FMs)

- Off-the-shelf access, fine-tune the FM on your target task
- Limited interfaces:
 - Cannot access model weights/costly to fine-tune FM
 - Access FM embeddings of data as f(x), with fixed mapping f



Simple approaches?

Use weak supervision to make a dataset $(x, ilde{y})$

Then what?

Simple approaches?

Use weak supervision to make a dataset $(x, ilde{y})$

Then what?

- Fine-tuning? X
- KNN with f(x)
- Linear probes, adapters (simple MLPs on f(x))

Can we do better than sequential application of weak supervision and FM embeddings?













Weak Supervision 101

Method setup

Input:

- Unlabeled dataset $\mathcal{D}=\{x_i\}_{i=1}^n$ with $x\in\mathcal{X}$, unknown label $y\in\{-1,1\}$
- Weak sources $\lambda_1,\ldots,\lambda_m:\mathcal{X} o \{-1,0,1\}$
- FM embedding mapping $f:\mathcal{X} o \mathbb{R}^d$

Method setup

Input:

- Unlabeled dataset $\mathcal{D}=\{x_i\}_{i=1}^n$ with $x\in\mathcal{X}$, unknown label $y\in\{-1,1\}$
- Weak sources $\lambda_1,\ldots,\lambda_m:\mathcal{X} o \{-1,0,1\}$
- FM embedding mapping $f:\mathcal{X} o \mathbb{R}^d$

Desired output: $\Pr(y=1|\lambda_1,\ldots,\lambda_m,x)$

- Given a datapoint and a list of votes on its label, what is its true label?
- Intuitively: some weighted combination of votes

1. Learn relationship between y and $\lambda_1, \ldots, \lambda_m$ via graphical model

 $ext{Pr}_{ heta}(y,\lambda_1,\ldots,\lambda_m) \iff egin{array}{c} y \ heta_1 \ heta_2 \ heta_3 \ heta_1 \ heta_2 \ heta_3 \ heta_2 \ heta_3 \ heta_2 \ heta_3 \ heta_2 \ heta_3 \ heta_3 \ heta_2 \ heta_3 \ heta_3 \ heta_3 \ heta_2 \ heta_3 \ hea_3 \ heta_3 \ he$

1. Learn relationship between y and $\lambda_1, \ldots, \lambda_m$ via graphical model



- Learn accuracy parameters $heta_i = \mathbb{E}[\lambda_i y]$
- Under the hood: algorithms for *latent variable estimation*^{*} by computing covariances between LFs, e.g. estimating $\mathbb{E}[\lambda_i \lambda_j]$ on \mathcal{D} or via maximum likelihood estimation^{**}

*Fu et. al. "Fast and three-rious: Speeding up weak supervision with triplet methods." ICML 2020. *Ratner et. al. "Training Complex Models with Multi-Task Weak Supervision." AAAI 2019. **Ratner et. al. "Data Programming: Creating Large Training Sets, Quickly". NeurIPS 2016.

2. Inference

- Given *x*, output estimate of $\Pr_{\theta}(y=1|\lambda_1,\ldots,\lambda_m)$
- Under the hood: Bayes rule + weighing each λ_i by a function of θ_i

2. Inference

- Given x, output estimate of $\Pr_{ heta}(y=1|\lambda_1,\ldots,\lambda_m)$
- Under the hood: Bayes rule + weighing each λ_i by a function of θ_i

No x!

What's wrong with this?

• We don't use x directly in inference, only its votes and $\theta_1, \ldots, \theta_m$ \rightarrow coarse accuracies

2. Inference

- Given *x*, output estimate of $\Pr_{\theta}(y=1|\lambda_1,\ldots,\lambda_m)$
- Under the hood: Bayes rule + weighing each λ_i by a function of θ_i

No x!

What's wrong with this?

- We don't use x directly in inference, only its votes and $\theta_1, \ldots, \theta_m$ \rightarrow coarse accuracies
- When $\lambda_i = 0$ on x, the algorithm discards it and computes $\Pr_{\theta}(y = 1 | \lambda_i, \dots, \lambda_{i-1}, \lambda_{i+1}, \lambda_m)$ \rightarrow abstains result in worse probabilistic labels

Method

Two simple modifications using FM embeddings *f*:

Two simple modifications using FM embeddings *f*:





Two simple modifications using FM embeddings *f*:



Two simple modifications using FM embeddings *f*:



Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains

Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains



Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains



Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains



Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains



Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains



Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains



Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains



Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains



Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains



Two simple modifications using FM embeddings *f*:

- 1. **Partition** the FM embedding space and estimate a set of parameters per part \rightarrow finer grained accuracies, better approximation of $\Pr_{\theta}(y = 1 | \lambda_1, \dots, \lambda_m, x)$
- 2. **Extend** votes of LFs to points that are close by in FM embedding space to construct $\bar{\lambda}_1, \ldots \bar{\lambda}_m \rightarrow$ fewer abstains



Why does this work?















- A Partitioning into too many sets = good local estimates, high variance
- Extending too far in embedding space = LF votes become incorrect because true label changes value
- Need to control these depending on how smooth FM is!

- A Partitioning into too many sets = good local estimates, high variance
- Extending too far in embedding space = LF votes become incorrect because true label changes value
- Need to control these depending on how smooth FM is!

Theoretical results: improvement over standard WS depends on FM smoothness, amount of data, and above choices.

Theoretical results

Thm 1 (informal). Suppose we partition data into *s* sets, and *d* is the average diameter of a set in embedding space. *K* is a smoothness constant (lower K = more smooth). Then, Liger's error (no extensions) is:

$$Error \leq K \cdot d + \mathcal{O}igg(rac{ms}{n}igg) + H(y|\lambda_1, \dots \lambda_m, x)$$
Bias Variance Irreducible Error (conditional entropy)

Theoretical results

Thm 1 (informal). Suppose we partition data into *s* sets, and *d* is the average diameter of a set in embedding space. *K* is a smoothness constant (lower K = more smooth). Then, Liger's error (no extensions) is:

$$Error \leq K \cdot d + \mathcal{O}igg(rac{ms}{n}igg) + H(y|\lambda_1, \dots \lambda_m, x)$$
Bias Variance Irreducible Error (conditional entropy)

- Bias-variance tradeoff in size/number of partitions, depending on smoothness
- Irreducible error: amount of randomness in *y* after observing *x* and LF votes

Theoretical Results

Next, what does extending and using $ar{\lambda}$ do?

- Increases bias (larger diameter due to extending)
- Decreases variance (more coverage = more points to estimate on)
- Irreducible error: $H(y|\lambda,x)
 ightarrow H(y|ar{\lambda},x)$ unclear!

Theoretical Results

Thm 2 (informal). Suppose we extend one LF λ_i by r. Define p_i as the proportion of the region where λ_i is extended, and $p(\lambda_{-i})$ as the model accuracy when using all LFs but λ_i . Let M(r) be a function describing smoothness (lower = smoother). Then,

$$H(y|\lambda,x)-H(y|ar{\lambda},x)\geq 2p_i(1-p(\lambda_{-i}))^2(heta_i-M(r))^2$$

Theoretical Results

Thm 2 (informal). Suppose we extend one LF λ_i by r. Define p_i as the proportion of the region where λ_i is extended, and $p(\lambda_{-i})$ as the model accuracy when using all LFs but λ_i . Let M(r) be a function describing smoothness (lower = smoother). Then,

$$H(y|\lambda,x)-H(y|ar{\lambda},x)\geq 2p_i(1-p(\lambda_{-i}))^2(heta_i-M(r))^2$$

- Tradeoff in how much we extend $(p_i \vee M(r))$, depending on smoothness
- Improvement when $\bar{\lambda}_i$ has better-than-random accuracy
- Improvement is less when other LFs already are highly accurate



Empirical Results

Weak Supervision Datasets + GPT-3 embeddings (for text), CLIP embeddings (for video)

	Weak Sources Only					
	Task	WS-kNN	WS-Adapter	WS-LM	LIGER	Δ Coverage
NLP	Spam	72.8	92.3	83.6	95.0	+45.5
	Weather	62.0	86.0	78.0	98.0	+90.2
	Spouse	16.9	17.1	47.0	52.2	+12.1
Video	Basketball	33.3	48.9	27.9	69.6	+8.3
	Commercial	84.7	92.8	88.4	93.5	+18.8
	Tennis	83.0	83.8	82.0	83.3	+32.5

A closer look at smoothness



Matches theory that more smooth FM embeddings = better performance

A closer look at smoothness

What's the best way to embed sentences? Spouse Lipschitzness % Points with Changed Value 40 Smoother Prompting F1-score 20 No Prompt 48.5 **Prompt Beginning** 50.2 **Prompt End** 52.2 0.002 0.006 **Cosine Distance Radius** No Prompt — Prompt Beginning Prompt End

Matches theory that more smooth FM embeddings = better performance

Summary

- Unlabeled data + alternative forms of knowledge lying around: weak sources, foundation models
- Liger smartly combines weak supervision with foundation models with two simple steps that exploit FM embedding *smoothness*
 - **Partition:** Estimate finer-grained accuracy parameters
 - **Extend:** Improve coverage of labeling functions
- Improves over standard WS and simple baselines based on FM smoothness

Future Directions

View FMs as blackbox sources of information with easy (partial) access.

How do we understand and best utilize them in inexpensive ways?

- Simple training methods on their embeddings (e.g. adapters)
- How to choose a good FM for your task?
- Can we combine multiple FMs themselves a la Weak Supervision?



Contact: mfchen@stanford.edu, danfu@cs.stanford.edu

Shoring Up the Foundations: Fusing Model Embeddings and Weak Supervision https://arxiv.org/abs/2203.13270

