Alternative Specifications in Machine Learning

Percy Liang

Jacob Steinhardt, Fereshte Khani



FLI Workshop

January 5, 2016



Train $p_0(x)$



Train $p_0(x)$



Classic statistical learning theory:

training distribution = test distribution



Train $p_0(x)$



Classic statistical learning theory:

training distribution = test distribution

Relaxation: domain adaptation, mild adversaries

training distribution $\,\approx\,$ test distribution



Train $p_0(x)$



Classic statistical learning theory:

training distribution = test distribution

Relaxation: domain adaptation, mild adversaries

training distribution $\,\approx\,$ test distribution

ssue:

doesn't address large changes (disasters, adversaries)

Changes and changes



Changes and changes



Long-term risks of AI: unknown unknowns

What's the right specification?

Specification: standard machine learning-

Input: training data

Output: model that does obtains low expected test error

Is **expected** test error enough?

What's the right specification?

Specification: standard machine learning-

Input: training data

Output: model that does obtains low expected test error

Is **expected** test error enough?

Scenario:

- Err on 1% on instances
- Agents maximize, adversaries minimize, could drive us there!





[ACL 2016]

- Specification: selective prediction Input: training data Output: model that outputs correct answer or "don't know"

New specification 1/2

Previous work: Chow (1970); Tortorella (2000); El-Yaniv & Wiener (2010); Balsubramani (2016)



Assumption: exists mapping with zero error





Models consistent with training data:

$$\mathcal{C} = \{ M \ge 0 : SM = T \}$$



Models consistent with training data:

$$\mathcal{C} = \{ M \ge 0 : SM = T \}$$

Challenge:

Checking all consistent $M \in \mathcal{C}$ is slow...

Fast two point scheme



- Choose $M_1, M_2 \in \mathcal{C}$ randomly enough
- Return "don't know" iff M_1 and M_2 disagree

Experimental results

• GeoQuery semantic parsing dataset (800 train, 280 test)

What is the population of Texas?







Jacob Steinhardt

[NIPS 2016]



Previous work: Donmez et al. (2010); Dawid/Skene (1979); Zhang et al. (2014); Jaffe et al. (2015); Balasubramanian et al. (2011)

Is this possible?





Compute $\mathbb{E}[loss(x, y; \theta)]$

Assumptions

Conditional independence:



Assumptions

Conditional independence:



Loss function decomposes:

$$A(x;\theta) - f_1(x_1, y;\theta) - f_2(x_2, y;\theta) - f_3(x_3, y;\theta)$$

Assumptions

Conditional independence:



Loss function decomposes:

$$A(x;\theta) - f_1(x_1,y;\theta) - f_2(x_2,y;\theta) - f_3(x_3,y;\theta)$$

only conditional independence structure

Intuition



Intuition



Three views disagree \rightarrow high error



(k labels, views v = 1, 2, 3)

$$\begin{array}{c}
 f_v(x,1) \\
 \dots \\
 f_v(x,k)
 \end{array}$$



(k labels, views v = 1, 2, 3)

$$M_{v} = \begin{bmatrix} \mathbb{E}[f_{v}(x,1) \mid y=1] & \dots & \mathbb{E}[f_{v}(x,1) \mid y=k] \\ & \dots & & \dots \\ \mathbb{E}[f_{v}(x,k) \mid y=1] & \dots & \mathbb{E}[f_{v}(x,k) \mid y=k] \end{bmatrix}$$

13



(k labels, views v = 1, 2, 3)

• Observe $\mathbb{E}[f_1(x,a)f_2(x,b)]$



(k labels, views v = 1, 2, 3)

• Observe $\mathbb{E}[f_1(x,a)f_2(x,b)f_3(x,c)]$



(k labels, views v = 1, 2, 3)

- Observe $\mathbb{E}[f_1(x,a)f_2(x,b)f_3(x,c)]$
- Perform tensor factorization to obtain

$$M_{vba} = \mathbb{E}[f_v(x,b) \mid y = a]$$



(k labels, views v = 1, 2, 3)

- Observe $\mathbb{E}[f_1(x,a)f_2(x,b)f_3(x,c)]$
- Perform tensor factorization to obtain

$$M_{vba} = \mathbb{E}[f_v(x,b) \mid y = a]$$

• Use to compute risk (up to label permutation)

 $\mathbb{E}[A(x;\theta) - f_1(x_1, y;\theta) - f_2(x_2, y;\theta) - f_3(x_3, y;\theta)]$

Results



Discussion





- Maximize expected accuracy \Rightarrow selective prediction, unsupervised risk estimation
- Key question: Can we weaken the assumptions?





worksheets.codalab.org

Collaborators



Fereshte Khani



Jacob Steinhardt



Thank you!