# Data programming for PubAnnotation via Snorkel

**JUAN M. BANDA**

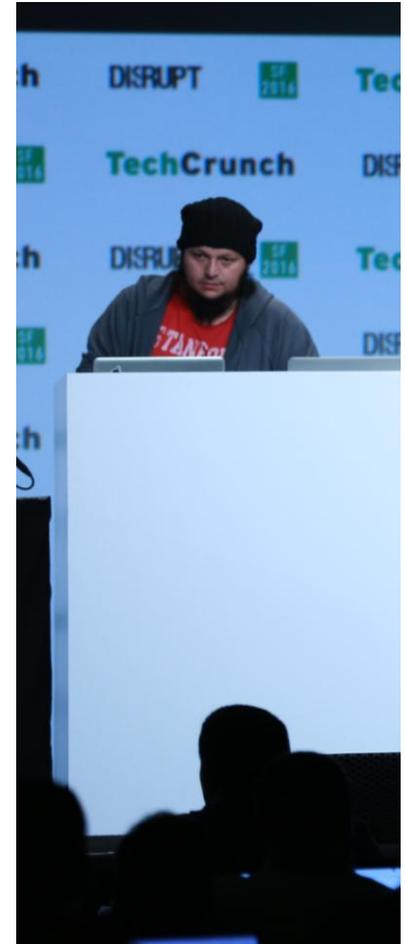**STANFORD CENTER FOR BIOMEDICAL INFORMATICS RESEARCH, STANFORD, CA**

Biomedical Linked Data Hackathon – BLAH3

Tokyo, Japan

January 16th

# About Me

- Research Scientist at Nigam Shah's lab at Stanford Center for Biomedical Informatics Research

- Hackathon enthusiast (attended 50+, 20 in 2016)

- Research Interests:
  - Electronic Phenotyping
  - Machine Learning:
    - Predictive Modeling
    - Noisy Labeling
  - Pharmacovigilance
  - Semantic Web
  - Image Retrieval/Classification (past)
  - CBIR systems (past)



**OHDSI**
OBSERVATIONAL HEALTH DATA SCIENCES AND INFORMATICS

**Stanford University**

# BLAH3 Project: Data programming for PubAnnotation via Snorkel

Why is this needed?

As the volumes of medical data being generated continues to increase with wide adoption of electronic health records, cheaper sequencing technology, etc. and their subsequent use in medical publications increases, we are faced with the challenge of producing high-quality labeled datasets for research
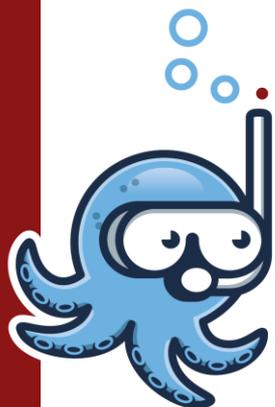
Data programming, have been proposed to rely instead on 'lower quality' auto generated noisy data to train weakly supervised models and then use generative models to denoise the noisy training data

**TRAINING DATA**

**The *New* New Oil**

# What is Snorkel? (1)

- Developed as a 'light' version of Deep Dive by Chris Ré's Hazy Research group at Stanford University

- Snorkel is framework for developing **structured information extraction applications** for domains in which large labeled training sets are not available or easy to obtain, using the *data programming* paradigm

- In **Data programming** approach to developing a machine learning system, the developer focuses on writing a set of *labeling functions*, which create a large but noisy training set.

- Snorkel then learns a generative model of this noise—learning, essentially, which labeling functions are more accurate than others—and uses this to train a discriminative classifier.
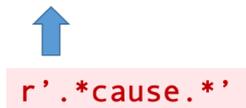
# What is Snorkel? (2)

- Labeling functions:

```
def lf1(x):
  cid = (x.chemical_id, x.disease_id)
  return 1 if cid in KB else 0
```
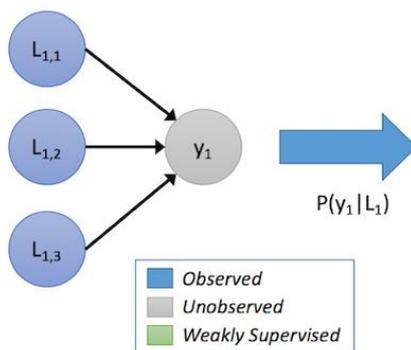
```
def lf2(x):
  m = re.search(r'.*cause.*', x.between)
  return 1 if m else 0
```

"Chemical A is found to
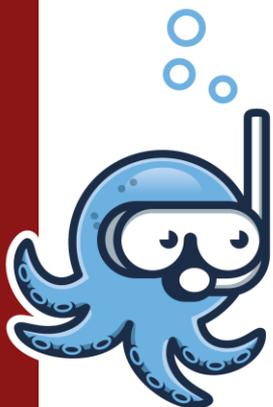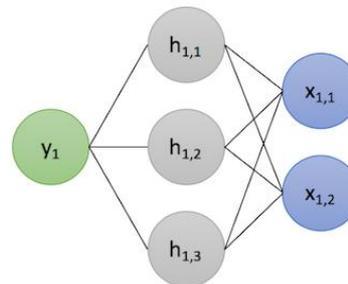cause disease B under
certain conditions…"     ➡ Label=TRUE

"Chemical A is found to
cause disease B under
certain conditions…"     ➡ Label=TRUE

Existing KB     Contains(A, B)

r'.*cause.*'

- Modeling steps:

Generative Model

$L_{1,1}$
$L_{1,2}$      $y_1$
$L_{1,3}$

$P(y_1 | L_1)$

Discriminative Model

$y_1$      $h_{1,1}$     $x_{1,1}$
           $h_{1,2}$     $x_{1,2}$
           $h_{1,3}$

Observed
Unobserved
Weakly Supervised

**snorkel**

**Stanford University**

# Snorkle annotation procedure



From: Ehrenberg, H.R., Shin, J., Ratner, A.J., Fries, J.A., R´e, C.:  Data programming with ddlite: Putting huma different part of the loop. In: Proceedings of the Workshop on Human-In-the-Loop Data Analytics. HILDA '16, pp. 13–1136. ACM, New York, NY, USA (2016). http://doi.acm.org/10.1145/2939502.2939515

Stanford University

# BLAH3 Goals

- Build extensions for Snorkle to produce all annotations in PubAnnotation format. This will allow us to take advantages of all of PubAnnotation's capabilities and cross functionality with other tools like Brat and Tagtog

- Build extensions for Snorkle to extract annotations from PubAnnotation for validation/testing of the labeling functions

**Maybe:**

- Find a collaborator that is providing an annotated corpus to BLAH3 and wants to build labeling functions to validate the power and usability of Snorkle

**Reach goals:**

- Build a bare-bones web-interface to interact with a running instance of Snorkle. Essentially build a REST-like tool that allows rules to be constructed and tested via a browser with minimal interaction with Python

# Acknowledgements

**DBCLS/BLAH organizers**

- Thanks for the travel support and the opportunity to work on this project

**Stanford**

- Nigam Shah's lab
- Hazy Research Lab

# QUESTIONS?

# THANK YOU