

Machine Learning for Healthcare

6.7930 [6.871], HST.956

Lecture 16:

Learning with Noisy Labels, Unsupervised Learning Applications, Weak Supervision

Prof. Manolis Kellis

Slides credit:
David Sontag



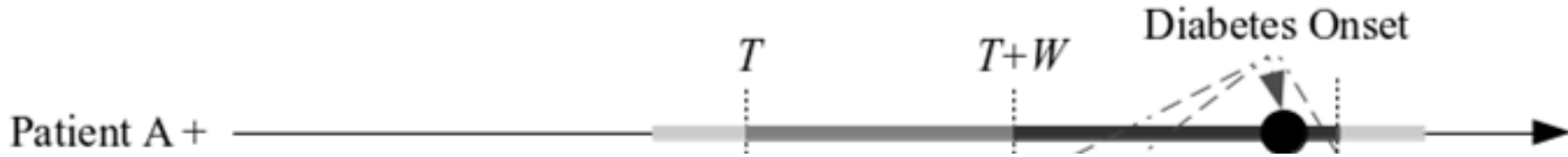
Outline for today's class

1. Learning with noisy labels

- Consistent estimation under class-conditional noise (Natarajan et al., NeurIPS '13)
- Application in health care (Halpern et al., JAMIA '16)

2. Weak supervision

Labels may be noisy



e.g. Predicting who will develop Diabetes.

- Who has diabetes at time T + Window W

- Decision tree based on:

- Diagnosis for T1D or T2D
- Medications for T1D or T2D
- Laboratory test result for T1D or T2D
- Lab test ordering for T1D or T2d

➔ Some 'hints' are informative, but noisy

➔ Some resulting labels will be wrong

➔ How does it affect learning?

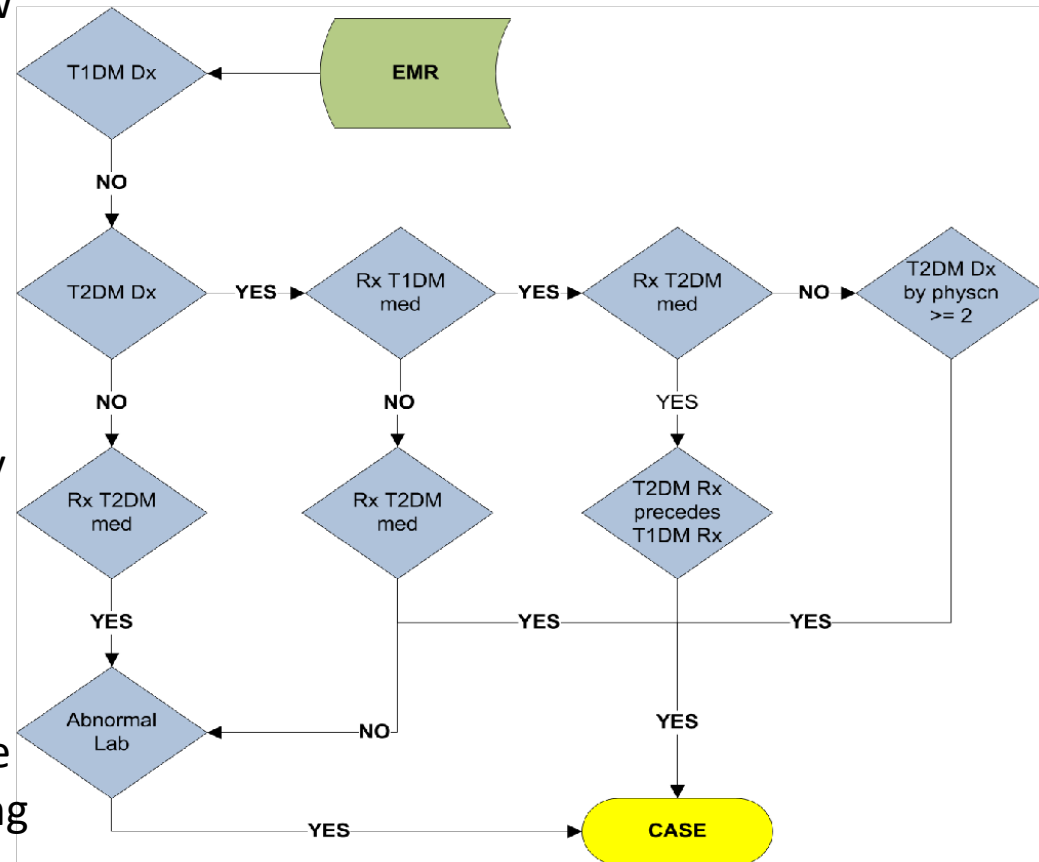
e.g. Billing codes, not diagnosis codes.

e.g. Lab test might actually come negative

Both false pos and false neg affect learning

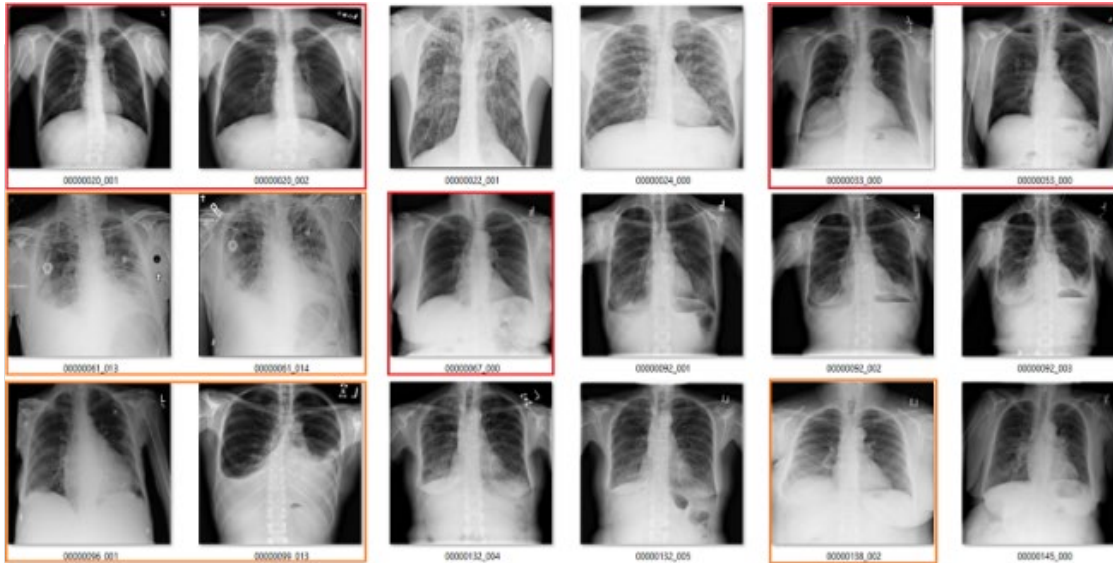
➔ How to deal with noisy labels?

Figure 1: Algorithm for identifying T2DM cases in the EMR.



Labels may be noisy

Fibrosis



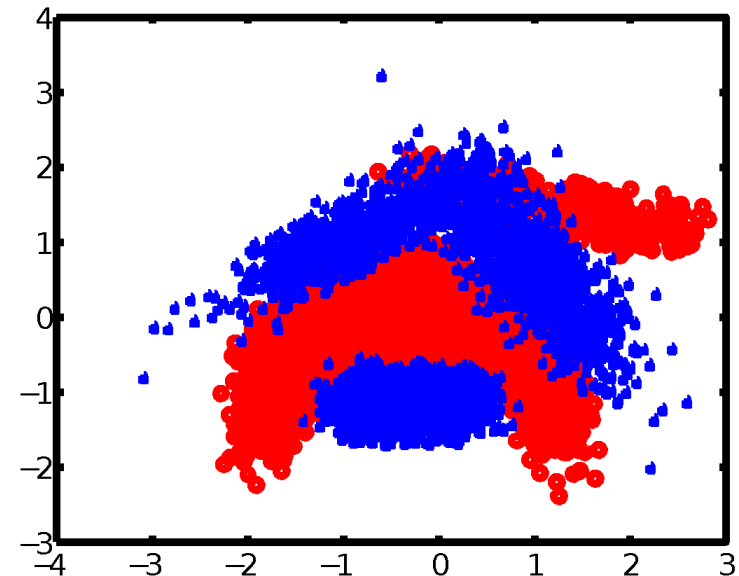
red = mislabeled
orange = maybe
mislabeled

Example:

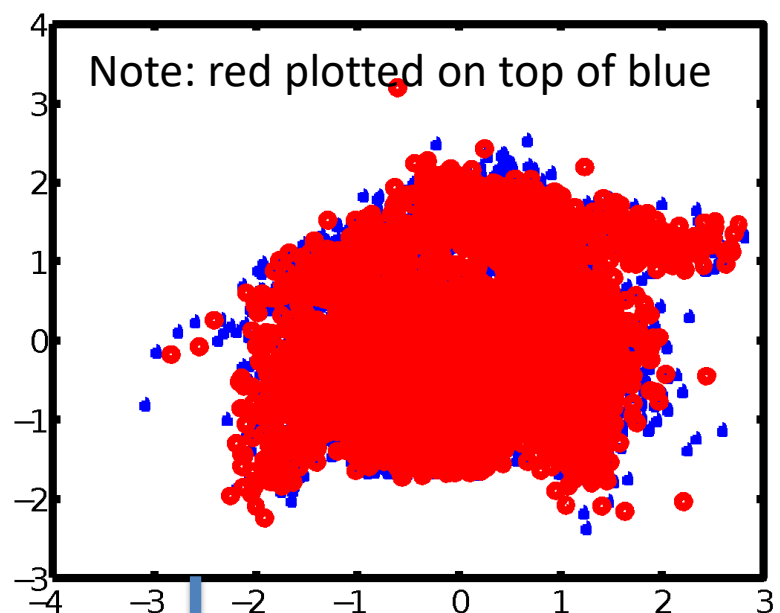
- Dataset released for training machine learning for chest x-ray
- Machine learning algorithms trained on it
- No-one had actually looked at the dataset closely
- Closer inspect revealed many errors and incorrectly-labeled images
- ➔ Can we ever build perfect ground-truth datasets? Perhaps not.
- ➔ How does that affect machine learning and the resulting models?
- ➔ How much noise can our ML algorithms tolerate and still learn correct features?
- ➔ Can we build machine learning methods that are robust to noisy labels?
- ➔ Can we 'bootstrap' from ML results to flag problematic labels in our dataset? =

[Wang et al., "Chest X-ray8"]

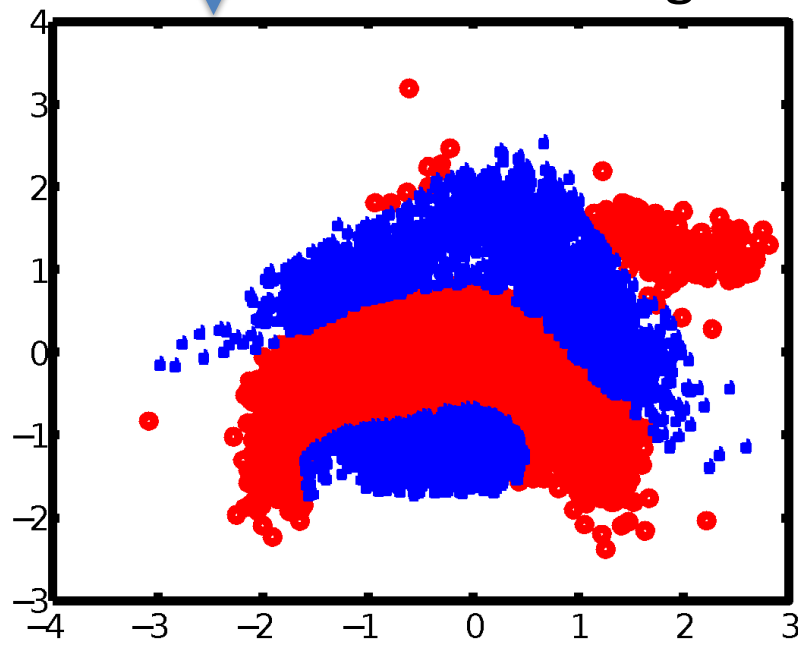
figure credit: <https://lukeoakdenrayner.wordpress.com/2017/12/18/the-chestxray14-dataset-problems/>



add 40%
label noise



Machine learning



Synthetic dataset

- Two features Red/Blue, +1/-1.

Add noise:

- Flip label for 40% of the dataset
- High noise rate: 0.4

Train classifier.

- Still learns correct function
- Can still accurately predict original 'true' label, even for artificially-flipped-label data points

Reflections:

- Makes you wonder if some original blue points in 'red regions' were incorrectly-labelled as well
- Sufficient generalization should allow you to 'question' some incorrectly labelled data points
- But if you only fix those, does that add bias, what about all those other incorrect ones you don't catch?

Learning with noisy labels

We will show that if we have

a) *class-conditional* label noise and

b) *lots* of training data,

learning as usual, substituting noisy labels, works!

This opens the door to using noisy labels for training, and coming up with clever ways of deriving these for free

- We won't know that we're doing really well.
- Our calculated "False positive rate" will be quite high (cuz some of our 'false' predictions will actually be true)
- But our model parameters will be quite good (good generalization power for future data, and in fact more accurate for some of our 'ground-truth' data that's incorrectly labelled)

Natarajan et al: Introduction

- Features X
- True unobserved labels $Y \in \{-1, 1\}$
- Noisy observed labels $\tilde{Y} \in \{-1, 1\}$
- True distribution $P(X, Y, \tilde{Y})$

X (age)	Y (diabetic)	\tilde{Y} (noisy version)
30	-1	
64	1	
75	1	

- Reason about joint distribution of X, Y, \tilde{Y}
(even though we do not observe true labels Y ,
still reason about joint distribution to prove guarantees about them)
- Assume process generating all three variables

Natarajan et al: Introduction

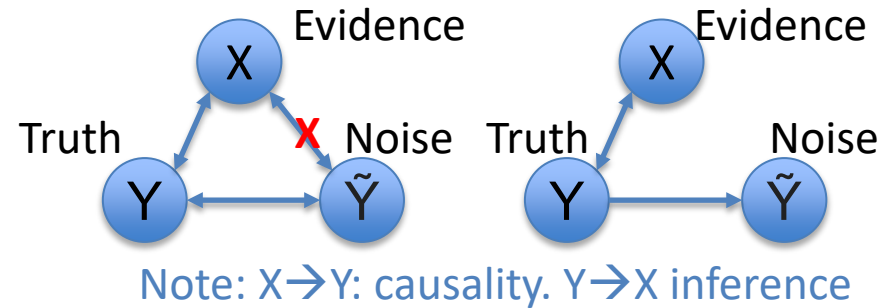
- Features X
- True unobserved labels $Y \in \{-1, 1\}$
- Noisy observed labels $\tilde{Y} \in \{-1, 1\}$ (Mnemonic: approx. $Y: \sim Y \rightarrow \tilde{Y}$)
- True distribution $P(X, Y, \tilde{Y})$

X (age)	Y (diabetic)	\tilde{Y} (noisy version)
30	-1	-1
64	1	1
75	1	-1

- Data sampled from $P(X, \tilde{Y}) = \sum_y P(X, Y = y, \tilde{Y})$
 - Marginalize over Y
 - Y is binary 1 or -1
 - Marginal for each of the two options
- Y exists, but it is hidden during training*

X (age)	\tilde{Y} (noisy version)
30	-1
64	1
75	-1

Assumption: class-conditional label noise



- Assume that $\tilde{Y} \perp\!\!\!\perp X|Y$

$$P(X, Y, \tilde{Y}) = P(X, Y) \underline{P(\tilde{Y}|Y)}$$

\tilde{Y} only depends on Y : label noise is independent of input features

- Class-Conditional random label Noise (CCN):
Misclassification noise only depends on Y , not X !
- Since Y is binary, need two parameters to fully define $P(\tilde{Y}|Y)$
- Misclassification rate ρ , separate for true pos (ρ_+), true neg (ρ_-)

flipping from +1 to -1 = ρ_+ $\rho_+ = P(\tilde{Y} = -1|Y = 1)$

flipping from -1 to +1: ρ_- $\rho_- = P(\tilde{Y} = 1|Y = -1)$

- Assume (for now) that these rates are known
 - However, the method works when weights are not known
- Only constraint: $\rho_+ + \rho_- < 1$, but either could be large

Learning with class-conditional noise

- If we could learn $\eta(X) = P(Y = 1|X)$, (prob of true 1 given data X), then we would be able to predict optimally. (but we've never observed Y, so how can we estimate it?)
- Instead, start by estimating prob of noisy-label Y=1 given data X

$$\tilde{\eta}(X) = P(\tilde{Y} = 1|X)$$

$$= P(\tilde{Y} = 1, Y = 1|X) + P(\tilde{Y} = 1, Y = -1|X)$$

Separately for True and for False, only two options → Marginalize

$$= P(Y = 1|X) P(\tilde{Y} = 1|Y = 1) + P(Y = -1|X) P(\tilde{Y} = 1|Y = -1)$$

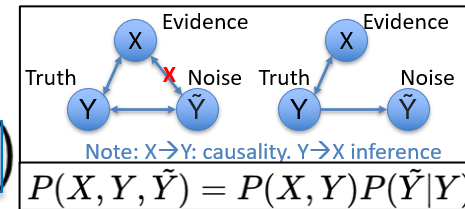
Independence assumption

$$= \eta(X) (1 - \rho_+) + (1 - \eta(X)) \rho_-$$

$$= \eta(X) (1 - \rho_+ - \rho_-) + \rho_-$$

$$\rightarrow \eta(X) = \frac{\tilde{\eta}(X) - \rho_-}{1 - \rho_+ - \rho_-}$$

(Natarajan et al., *Learning with Noisy Labels*. NeurIPS '13)



$$\rho_+ = P(\tilde{Y} = -1|Y = 1)$$

$$\rho_- = P(\tilde{Y} = 1|Y = -1)$$

Learning with class-conditional noise

$$\rightarrow \eta(X) = \frac{\tilde{\eta}(X) - \rho_-}{1 - \rho_+ - \rho_-} \quad \eta(X) \text{ is monotonically increasing in } \tilde{\eta}(X)$$

- Learn $\tilde{\eta}(X)$ using any ML algorithm which returns calibrated classifiers. Substitute $\tilde{\eta}(X)$ in the above equation to get $\eta(X)$!
- Thus: learn as if the labels were correct
 - ➔ Same ordering of true/false classification regardless.
 - ➔ Simple transformation to get actual true confidence
- Denominator must be >0 (why we needed $\rho_+ + \rho_- < 1$ constraint) [notice, if exactly 1, divide by zero, not good. If >1 , could flip labels]
- When might noise be helpful?
 - Privacy: e.g. sexually transmitted disease. Deniability
 - Force generalization, by avoiding overfitting
 - Could *predict* true classes, but only if perfect predictor (typically impossible)

Outline for today's class

1. Learning with noisy labels
 - Consistent estimation under class-conditional noise (Natarajan et al., NeurIPS '13)
 - **Application in health care (Halpern et al., JAMIA '16)**
2. Weak supervision

Application to electronic phenotyping

Hundreds of relevant clinical variables

Main Info

Ana Betz Female -> Male Age: 34y 4m 0d

Critical Information

Z88.0 : Personal history of allergy to penicillin Severe allergic reactions to β -lactams.
E10 : Type 1 diabetes mellitus
A90 : Dengue fever [classical dengue]
A40 : Streptococcal sepsis
Z31.2 : In vitro fertilization

Disabilities / Barriers:

General Info Functioning and Disability Surgeries Socioeconomics Lifestyle OB/GYN Genetics Medication

Main Misc

GP: Cordara, Cameron Family: Zenon-Betz Single Insurance: Insurator : 938291

Conditions

Condition	Status	Ac	In	Severity	All	Pr	Date of Diagno	Healed	Remarks
Z88.0 : Personal history of allergy to pen	unchanged	<input checked="" type="checkbox"/>	<input type="checkbox"/>	Severe	<input type="checkbox"/>	<input type="checkbox"/>	01/07/1991		
E10 : Type 1 diabetes mellitus	chronic	<input checked="" type="checkbox"/>	<input type="checkbox"/>	Moderate	<input type="checkbox"/>	<input type="checkbox"/>	11/10/1993		
A90 : Dengue fever [classical dengue]	acute	<input checked="" type="checkbox"/>	<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>			
A40 : Streptococcal sepsis		<input checked="" type="checkbox"/>	<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>			
Z31.2 : In vitro fertilization		<input checked="" type="checkbox"/>	<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>			

- Abdominal pain
- Active malignancy
- Altered mental status
- Cardiac etiology
- Renal failure
- Infection
- Urinary tract infection
- Shock
- Smoker
- Pregnant
- Lower back pain
- Motor Vehicle accident
- Psychosis
- Anticoagulated
- Type II diabetes

...

Figure source:

https://commons.wikimedia.org/wiki/File:GNU_Health_patient_main_screen.png

Simplest approach: rules


- We would like to estimate, for every patient, which phenotypes apply to them (at some point in time)
- Common practice is to derive manual rules:

		Nursing home?		physician response (gold standard)
		T	F	
text contains: "nursing home"	T	297	129	PPV 0.70
	F	1,319	34511	297/(297+129)

Sensitivity
0.18


Slow, expensive, poor sensitivity.

Often we can find noisy labels WITHIN the data!

Phenotype	Example of noisy label (“anchor”) 
Diabetic (type I)	gsn:016313 (insulin) in Medications
Strep Throat	Positive strep test in Lab results
Nursing home	“from nursing home” in Text
Pneumonia	“pna” in Text
Heart attack	ICD10 I21 in Billing codes

How can we use these for machine learning?

Often we can find noisy labels WITHIN the data!

Phenotype	Example of noisy label (anchor) 
Heart attack	ICD10 I21 in Billing codes

- Suppose we want to know, was a patient admitted to the emergency department for a heart attack?
- Billing codes not available at prediction time (not useful for real-time use), but can be used for labels
- Reasonable to assume that $\rho_- = P(\tilde{Y} = 1|Y = -1) \approx 0$ (fraud), but because of noisy nature of billing codes, $\rho_+ = P(\tilde{Y} = -1|Y = 1)$ likely non-zero

Called “positive only” noise since it implies $P(Y = 1|\tilde{Y} = 1) = 1$

Anchor & Learn Algorithm

(special case for anchors derived from future data)

Training


1. Treat the anchors (noisy labels) as “true” labels then censor them from the dataset (*as they don't appear in practice for real-time prediction*)
2. Learn a classifier to predict presence/absence of anchor (*whether noisy-label (anchor) \tilde{Y} appears*)
3. Calibration step: divide by $\frac{1}{|P|} \sum_P P(\tilde{Y} = 1|X)$

P = data points with $\tilde{Y} = 1$

Test time

1. First check if anchor is present.
 - If yes: high value, cuz of high positive predictive value of anchor
 - If not: apply the learned classifier + multiply by calibration constant

Often we can find noisy labels WITHIN the data!

Phenotype	Example of noisy label (anchor)	
Nursing home	“from nursing home” in Triage note	

- We again assume that $\rho_- = P(\tilde{Y} = 1|Y = -1) \approx 0$, but because many ways to write “from nursing home” in text, we have $\rho_+ = P(\tilde{Y} = -1|Y = 1)$ likely non-zero
- If we simply learn to predict \tilde{Y} using the notes, we will learn a trivial classifier! It will simply extract mentions of this phrase!
- This is a clear violation of the assumption $\tilde{Y} \perp X|Y$, since \tilde{Y} is derived from X

*In this dataset, we have some natural candidates for anchors that we can take from different parts of the record. For example, we may use **medications** that are specific to a single disease. **Lab tests**, **phrases** in the patient notes, or ICD9 or 10 **billing codes**.*

Anchor & Learn Algorithm

Training

1. Treat the anchors as “true” labels
2. Learn a classifier to predict whether the **anchor** appears **based on all other features** (throw away all features used to predict true label)
3. Calibration step: divide by $\frac{1}{|P|} \sum_P P(\tilde{Y} = 1|X)$

Test time

1. If the anchor is present: Predict 1
2. Else: Predict using the learned classifier (with calibration)

P = data points with $\tilde{Y} = 1$

Evaluating phenotypes

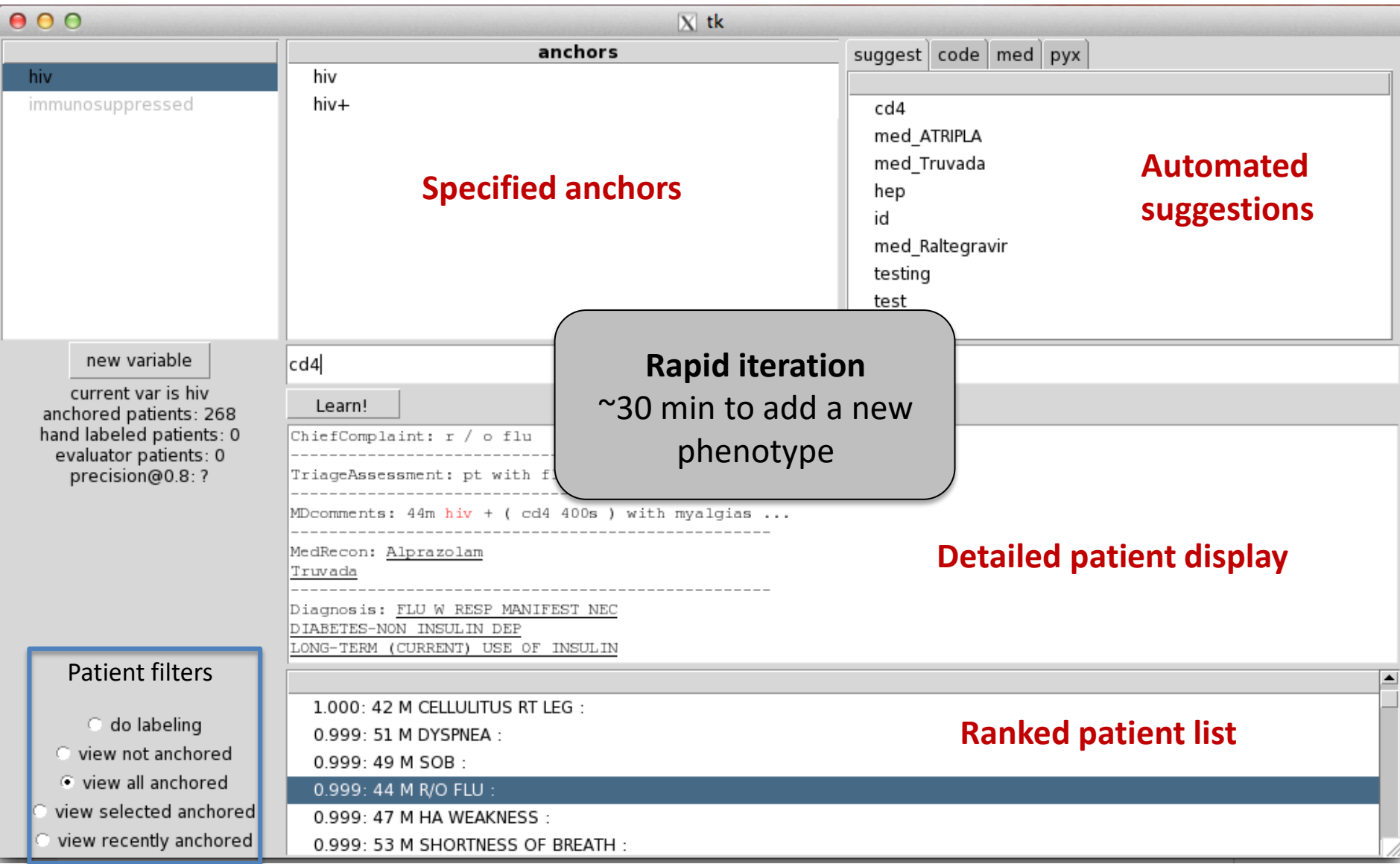
- Derived anchors and learned phenotypes using 270,000 patients' emergency department medical records

<u>History</u>	<u>Acute</u>		
Alcoholism	Abdominal pain	Deep vein thrombosis	Laceration
Anticoagulated	Allergic reaction	Employee exposure	Motor vehicle accident
Asthma/COPD	Ankle fracture	Epistaxis	Pancreatitis
Cancer	Back pain	Gastroenteritis	Pneumonia
Congestive heart failure	Bicycle accident	Gastrointestinal bleed	Psych
Diabetes	Cardiac etiology	Geriatric fall	Obstruction
HIV+	Cellulitis	Headache	Septic shock
Immunosuppressed	Chest pain	Hematuria	Severe sepsis
Liver malfunction	Cholecystitis	Intracerebral hemorrhage	Sexual assault
	Cerebrovascular accident	Infection	Suicidal ideation
		Kidney stone	Syncope
			Urinary tract infection



Then used in real-time to predict labels in a hospital setting
Evaluate using ground-truth data

[Halpern, Horng, Choi, Sontag, AMIA '14]
[Halpern, Horng, Choi, Sontag, JAMIA '16]



Specified anchors

Automated suggestions

**Rapid iteration
~30 min to add a new
phenotype**

Detailed patient display

Ranked patient list

Patient filters

- do labeling
- view not anchored
- view all anchored
- view selected anchored
- view recently anchored

Evaluating phenotypes

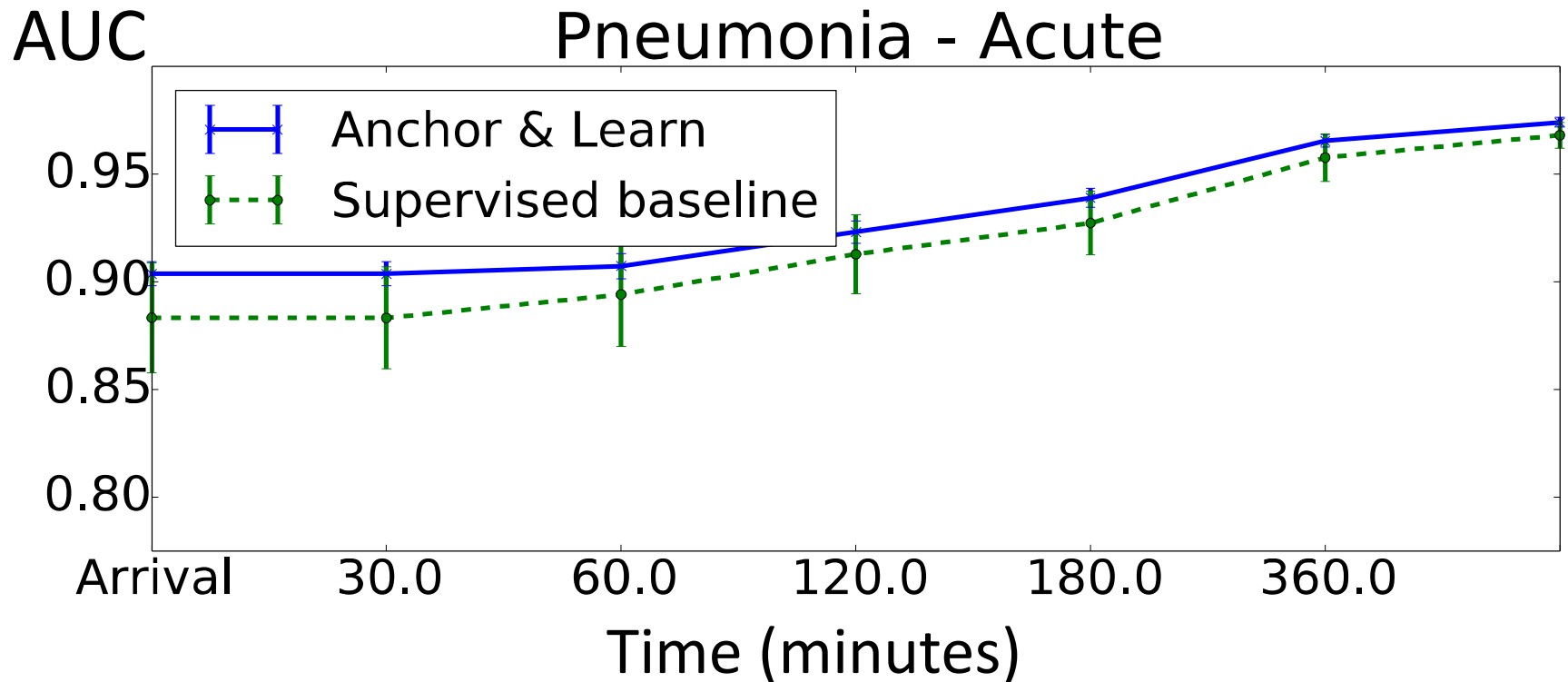
- Derived anchors and learned phenotypes using 270,000 patients' emergency department medical records
- To obtain ground truth, added a small number of questions to patient discharge procedure, rotated randomly

Does the patient have an active malignancy? ⁱ

Unlikely Unsure Likely



Evaluating phenotypes



Comparison to supervised learning using labels for 5000 patients
Gets better over time, as predictions become more confident
Performance similar to supervised baseline (possibly even better?)

Evaluating phenotypes – example model (cardiac etiology)

Anchors

ICD9 codes

410.* acute MI

411.* other acute ...

413.* angina pectoris

785.51 card. shock

Pyxis

coron. vasodilators

loop diuretic

Highly weighted terms

Ages

age=80-90

age=70-80

age=90+

Medications

lasix

furosemide

cp

chest pain

edema

cmed

chf exacerbation

sob

pedal edema

Sex=M

Pyxis

aspirin

clopidogrel

Heparin Sodium

Metoprolol

Tartrate

Morphine Sulfate

Integrilin

Labetalol

Unstructured text

Evaluating phenotypes – example model (cardiac etiology)

Anchors

ICD9 codes

410.* acute MI

411.* other acute ...

413.* angina pectoris

785.51 card. shock

Pyxis

coron. vasodilators

cardiac medicine

BIDMC shortform

Highly weighted terms

Ages

age=80-90

age=70-80

age=90+

Medications

lasix

furosemide

cp

chest pain

edema

cmed

chf exacerbation

sob

pedal edema

Sex=M

Pyxis

aspirin

clopidogrel

Heparin Sodium

Metoprolol

Tartrate

Morphine Sulfate

Integrilin

Labetalol

Unstructured text

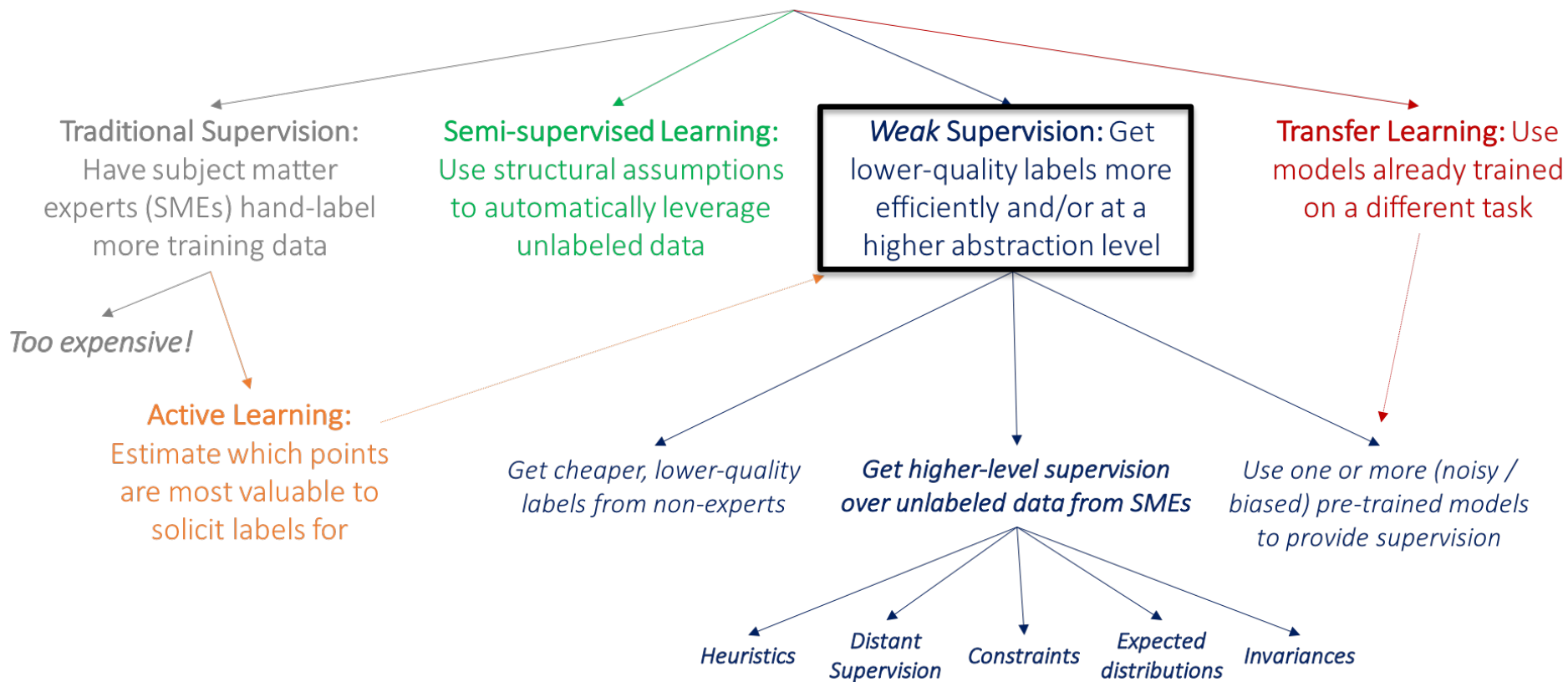
Outline for today's class

1. Learning with noisy labels

- Consistent estimation under class-conditional noise (Natarajan et al., NeurIPS '13)
- Application in health care (Halpern et al., JAMIA '16)

2. **Weak supervision**

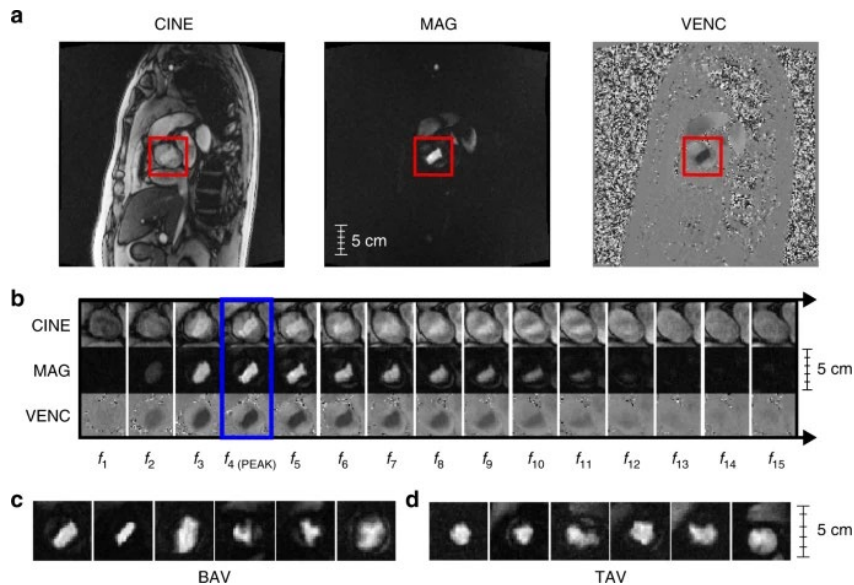
How to get more labeled training data?



Weak supervision

- Define one or more labeling functions $l(x)$ that outputs a label (or no label) for each example
- E.g., for sentiment analysis
 - “good” $\rightarrow +1$
 - “bad” $\rightarrow -1$
- Reconcile conflicting labels; ignore data points that are unlabeled
- Learn a model on the labeled data points

Classifying Aortic Valve Malformations

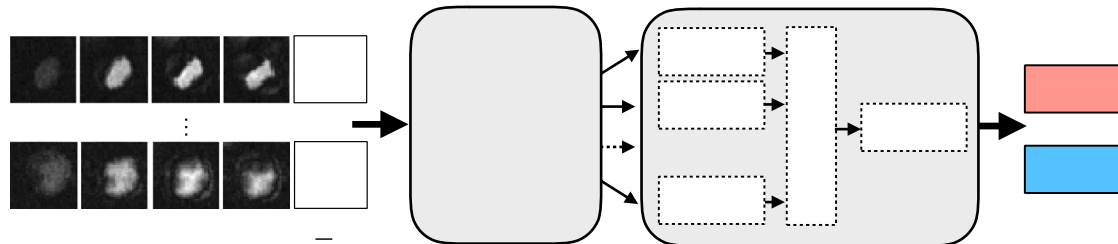


- Data: MRI sequences for 14,328 subjects from the UK Biobank
- True gold standard labels for *aortic valve malformations* (BAV) derived for 412 subjects
- Goal: Train a model which can classify BAV (positive or negative) when given a new MRI sequence

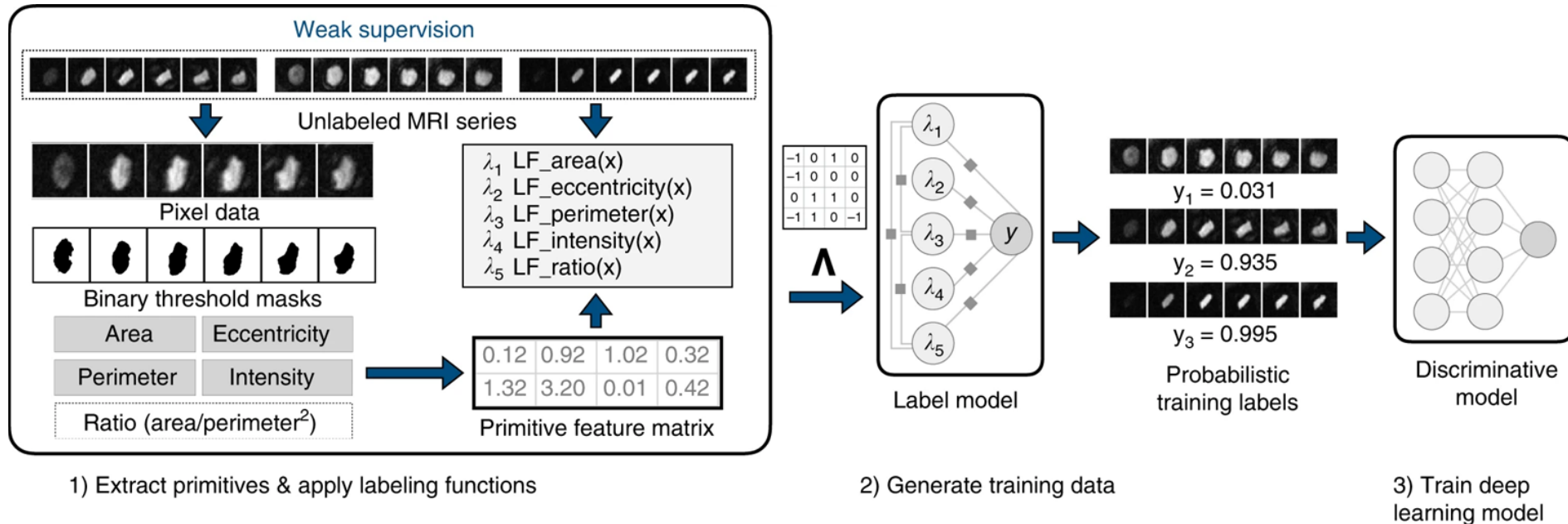
Classifying Aortic Valve Malformations

Methodology:

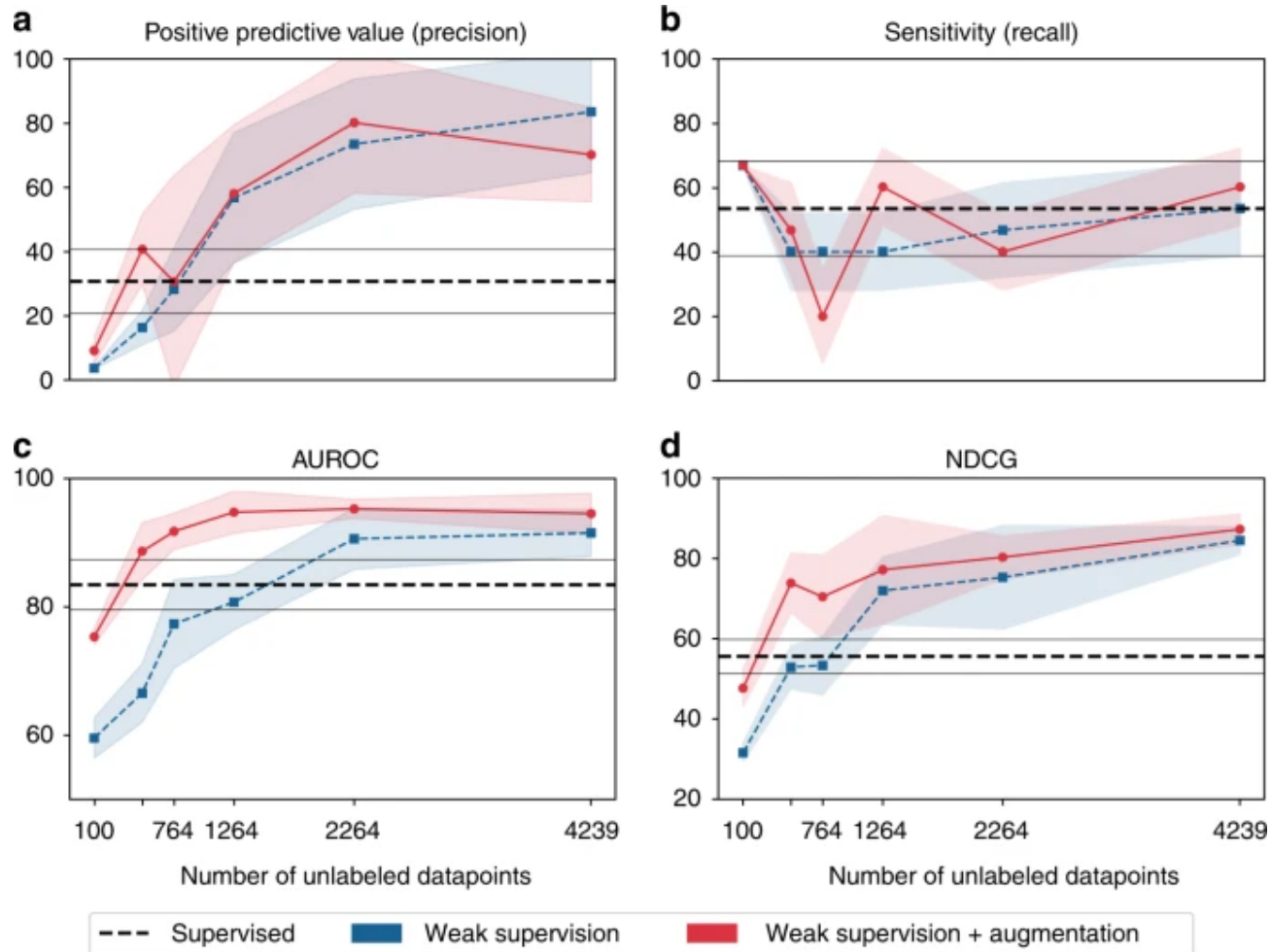
1. Train a factor graph-based model to predict noisy labels for all unlabeled examples
2. Train a hybrid convolutional NN / LSTM using the derived noisy labels



Classifying Aortic Valve Malformations



Classifying Aortic Valve Malformations



Weak supervision for text classification

- Example labeling functions:

```
3 # Setting LF output values
4 ABSTAIN_VAL = 0
5 HEMORRHAGE_VAL = 1
6 NO_HEMORRHAGE_VAL = -1
-
63 def LF_positive_hematoma(report):
64     """
65     Checking for words indicating hematoma
66     """
67     r1 = re.compile('(No|without|resolution|scalp|subgaleal)\\s{[\S]*\\s}{0,10}(hematoma)', re.IGNORECASE)
68     r = re.compile('hematoma', re.IGNORECASE)
69     for s in report.report.sentences:
70         if r.search(s.text) and (not r1.search(s.text)):
71             return HEMORRHAGE_VAL
72     return ABSTAIN_VAL
73
74 def LF_hemorrhage_hi_cover(report):
75     """
76     Checking for both hemorrhage and hematoma
77     """
78     if LF_positive_hemorrhage(report) == 0 and LF_positive_hematoma(report) == 0:
79         return NO_HEMORRHAGE_VAL
80     return HEMORRHAGE_VAL
```

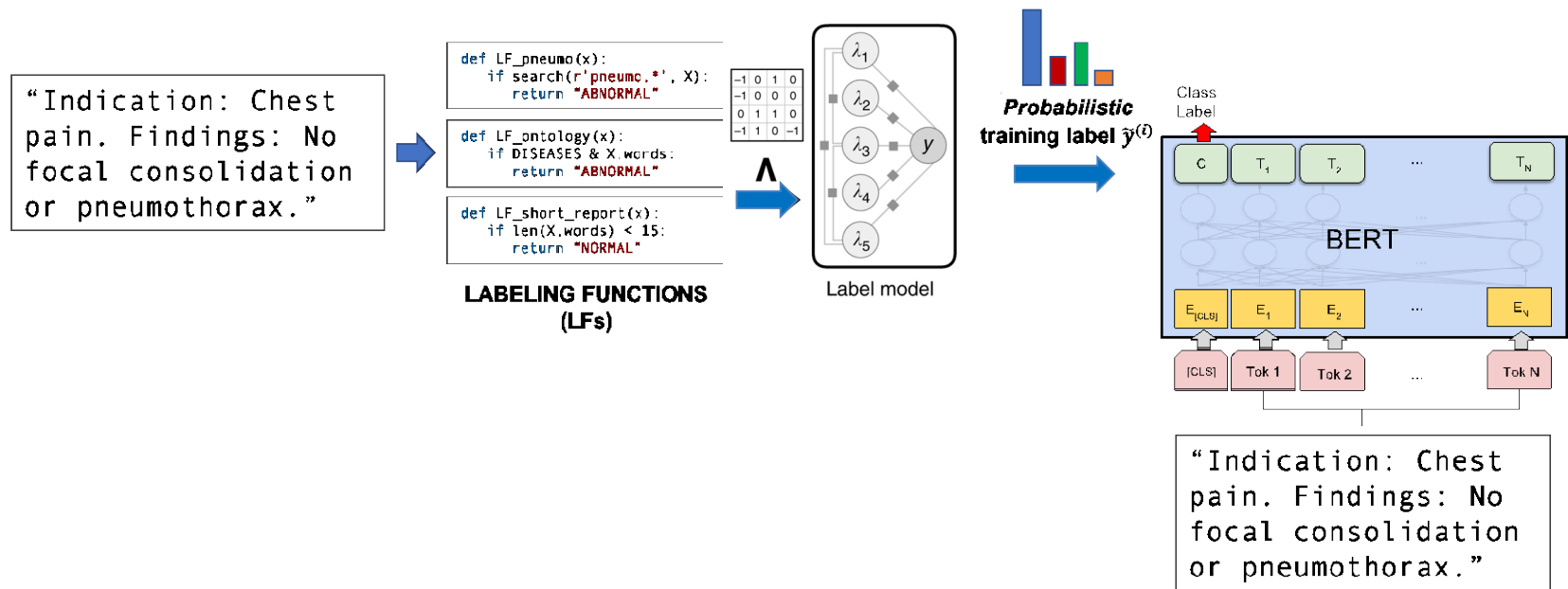
Weak supervision for text classification

- Example labeling functions:

```
3 # Setting LF output values
4 ABSTAIN_VAL = 0
5 HEMORRHAGE_VAL = 1
6 NO_HEMORRHAGE_VAL = -1
12 def LF_normal_V01(report):
13     """
14     Checking for specific normal phrase
15     """
16     r = re.compile('Normal CT of the Head',re.IGNORECASE)
17     for s in report.report.sentences:
18         if r.search(s.text):
19             return NO_HEMORRHAGE_VAL
20     return ABSTAIN_VAL
21
22 def LF_normal_V02(report):
23     """
24     Checking for specific normal phrase
25     """
26     r = re.compile('No acute intracranial abnormality',re.IGNORECASE)
27     for s in report.report.sentences:
28         if r.search(s.text):
29             return NO_HEMORRHAGE_VAL
30     return ABSTAIN_VAL
```

Weak supervision for text classification

- Use BERT as “end model”



- Why does this not simply learn to reproduce the labeling functions?

Weak supervision for text classification

Table 1: Statistics of all the tasks, domains and datasets included in WRENCH.

Task (#)	Domain (#)	Dataset (#)	#L _{label}	#L _F	Train	Dev	Test
					#Data	#Data	#Data
Sentiment Class.	Movie	IMDb [61, 79]	2	5	20,000	2,500	2,500
	Review	Yelp [107, 79]	2	8	30,400	3,800	3,800
Spam Class.	Review	Youtube [1]	2	10	1,586	120	250
	Text Message	SMS [2, 3]	2	73	4,571	500	500
Topic Class.	News	AGNews [107, 79]	4	9	96,000	12,000	12,000
Question Class.	Web Query	TREC [49, 3]	6	68	4,965	500	500
Relation Class.	News	Spouse [11, 77]	2	9	22,254	2,811	2,701
	Biomedical	CDR [13, 77]	2	33	8,430	920	4,673
	Web Text	SemEval [31, 109]	9	164	1,749	200	692
	Chemical	ChemProt [41, 102]	10	26	12,861	1,607	1,607

Weak supervision for text classification

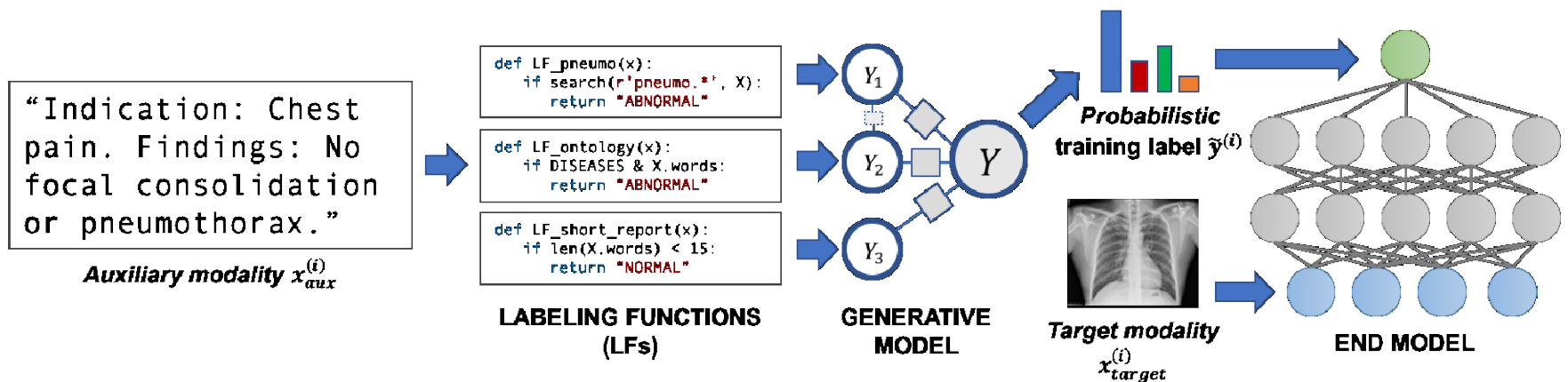
EM: end model (R=RoBERTa, RC=COSINE-RoBERTa, BC=COSINE-BERT)

LM: label model (MV="majority vote", WMC="weighted majority vote")

Dataset	Metric	Best Gold		Top 1			Top 2			Top 3		
		EM	Value	EM	LM	Value	EM	LM	Value	EM	LM	Value
IMDb	Acc.	R	93.25	RC	<u>MeTaL</u>	88.86	RC	<u>FS</u>	88.48	RC	<u>MV</u>	88.48
Yelp	Acc.	R	97.13	RC	FS	95.45	RC	<u>FS</u>	95.33	RC	<u>DS</u>	95.01
Youtube	Acc.	B	97.52	BC	<u>MV</u>	98.00	RC	MV	97.60	RC	<u>MV</u>	97.60
SMS	F1	B	96.96	RC	<u>WMV</u>	98.02	RC	MeTaL	97.71	RC	<u>WMV</u>	97.27
AGNews	Acc.	R	91.39	RC	DS	88.20	RC	<u>MV</u>	88.15	RC	<u>WMV</u>	88.11
TREC	Acc.	R	96.68	RC	DP	82.36	RC	<u>MeTaL</u>	79.84	BC	DP	78.72
Spouse	F1	–	–	BC	FS	56.52	–	MeTaL	46.62	RC	<u>MV</u>	46.28
CDR	F1	R	65.86	–	MeTaL	69.61	–	DP	63.51	RC	DP	61.40
SemEval	Acc.	B	95.43	BC	<u>DP</u>	88.77	BC	<u>MV</u>	86.80	RC	<u>DP</u>	86.73
ChemProt	Acc.	B	89.76	BC	<u>DP</u>	61.56	RC	<u>MV</u>	59.43	RC	<u>MV</u>	59.32

Weak supervision with multiple views

- Alternatively, one could just use the noisy labels from the label model to directly train the downstream model:



- Co-training (Blum & Mitchell, '98) can be used to improve performance further

Conclusion

- Can be difficult to get labeled data for machine learning in health care
- Often possible to quickly derive *noisy* labels (i.e., anchors or labeling functions)
- With conditionally independent noise, ML as usual can be used (with recalibration)
 - $x \perp \tilde{Y} | Y$ (noise rate constant for all examples)
 - Can sometimes *censor* the features to make this assumption more realistic (the anchor & learn method)
 - Alternatively, use pretrained representations