

Machine Learning for Healthcare

6.871x

Learning with noisy labels

David Sontag



INSTITUTE FOR MEDICAL
ENGINEERING & SCIENCE



HEALTH SCIENCES
& TECHNOLOGY

Labels may be noisy

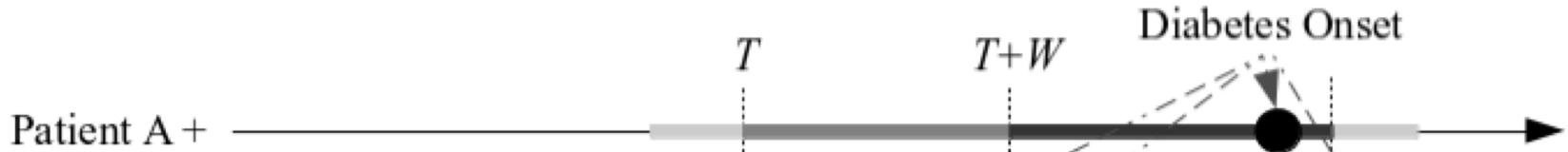
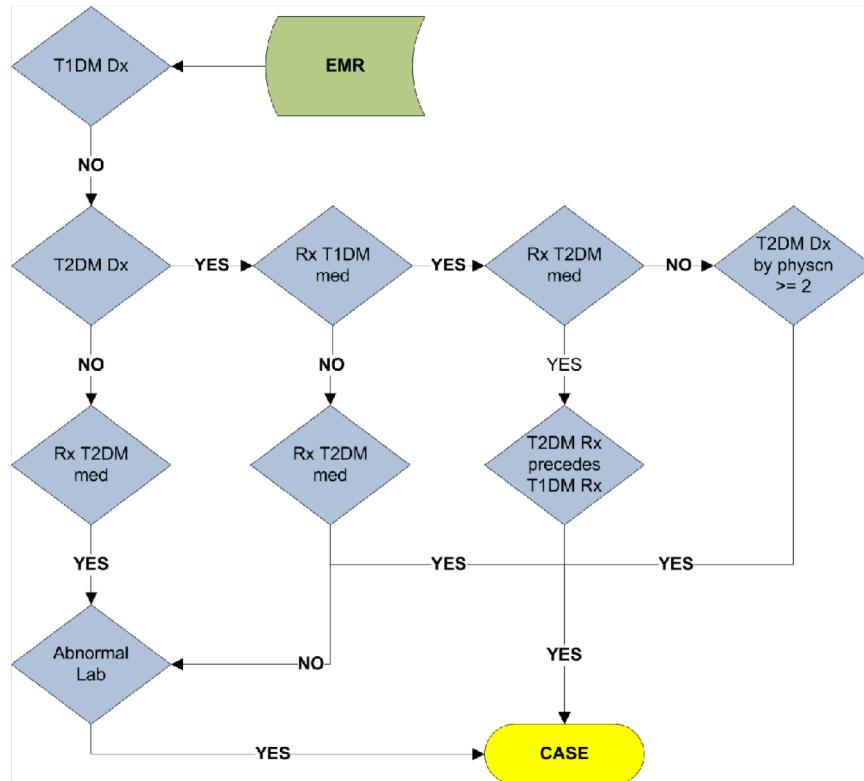


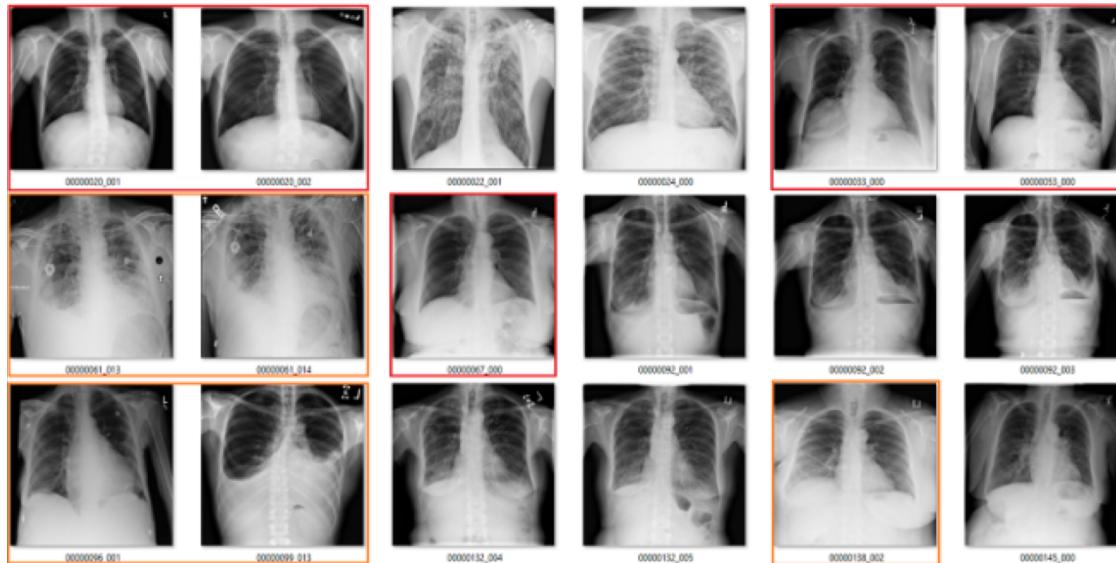
Figure 1: Algorithm for identifying T2DM cases in the EMR.



If the derived label
is noisy, how does
it affect learning?

Labels may be noisy

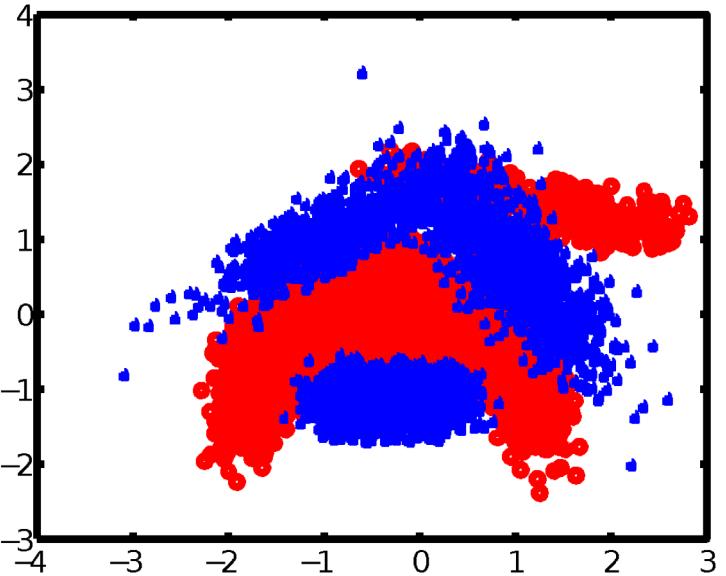
Fibrosis



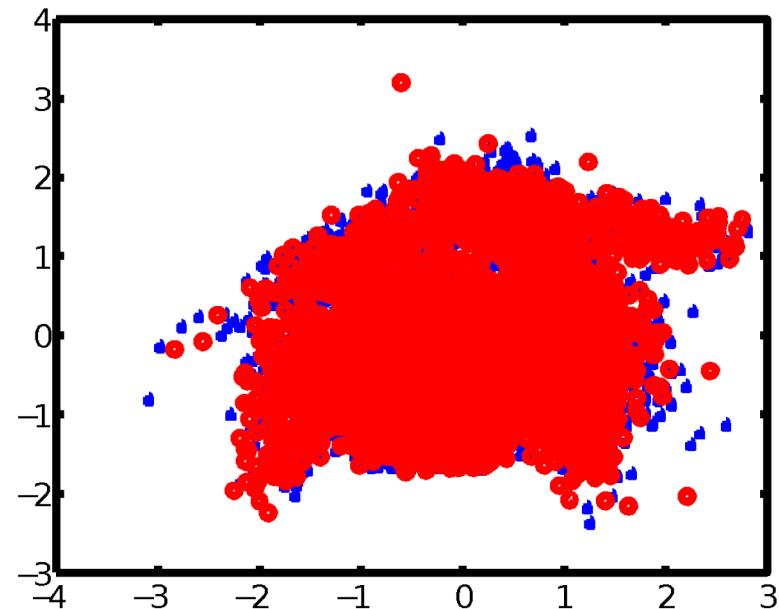
red = mislabeled
orange = maybe
mislabeled

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

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



40% label noise



Machine learning

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

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	-1
64	1	1
75	1	-1

- Data sampled from $P(X, \tilde{Y}) = \sum_y P(X, Y = y, \tilde{Y})$

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 $\tilde{Y} \perp X|Y$:

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

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

- Since Y is binary, need two parameters to fully define $P(\tilde{Y}|Y)$: $\rho_+ = P(\tilde{Y} = -1|Y = 1)$ & $\rho_- = P(\tilde{Y} = 1|Y = -1)$
- Assume that $\rho_+ + \rho_- < 1$ and that ρ_+, ρ_- are known

Learning an unbiased estimator

- If we could learn $\eta(X) = P(Y|X)$, then we would be able to predict optimally.

Learning an unbiased estimator

- If we could learn $\eta(X) = P(Y|X)$, then we would be able to predict optimally.

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

Learning an unbiased estimator

- If we could learn $\eta(X) = P(Y|X)$, then we would be able to predict optimally.

$$\begin{aligned}\tilde{\eta}(X) &= P(\tilde{Y} = 1|X) \\ &= P(\tilde{Y} = 1, Y = 1|X) + P(\tilde{Y} = 1, Y = -1|X)\end{aligned}$$

Learning an unbiased estimator

- If we could learn $\eta(X) = P(Y|X)$, then we would be able to predict optimally.

$$\begin{aligned}\tilde{\eta}(X) &= P(\tilde{Y} = 1|X) \\ &= P(\tilde{Y} = 1, Y = 1|X) + P(\tilde{Y} = 1, Y = -1|X) \\ &= P(Y = 1|X)P(\tilde{Y} = 1|Y = 1) \\ &\quad + P(Y = -1|X)P(\tilde{Y} = 1|Y = -1)\end{aligned}$$

Learning an unbiased estimator

- If we could learn $\eta(X) = P(Y|X)$, then we would be able to predict optimally.

$$\begin{aligned}\tilde{\eta}(X) &= P(\tilde{Y} = 1|X) \\ &= P(\tilde{Y} = 1, Y = 1|X) + P(\tilde{Y} = 1, Y = -1|X) \\ &= P(Y = 1|X)P(\tilde{Y} = 1|Y = 1) \\ &\quad + P(Y = -1|X)P(\tilde{Y} = 1|Y = -1) \\ &= \eta(X)(1 - \rho_+) + (1 - \eta(X))\rho_-\end{aligned}$$

Learning an unbiased estimator

- If we could learn $\eta(X) = P(Y|X)$, then we would be able to predict optimally.

$$\begin{aligned}\tilde{\eta}(X) &= P(\tilde{Y} = 1|X) \\ &= P(\tilde{Y} = 1, Y = 1|X) + P(\tilde{Y} = 1, Y = -1|X) \\ &= P(Y = 1|X)P(\tilde{Y} = 1|Y = 1) \\ &\quad + P(Y = -1|X)P(\tilde{Y} = 1|Y = -1) \\ &= \eta(X)(1 - \rho_+) + (1 - \eta(X))\rho_- \\ &= \eta(X)(1 - \rho_+ - \rho_-) + \rho_-\end{aligned}$$

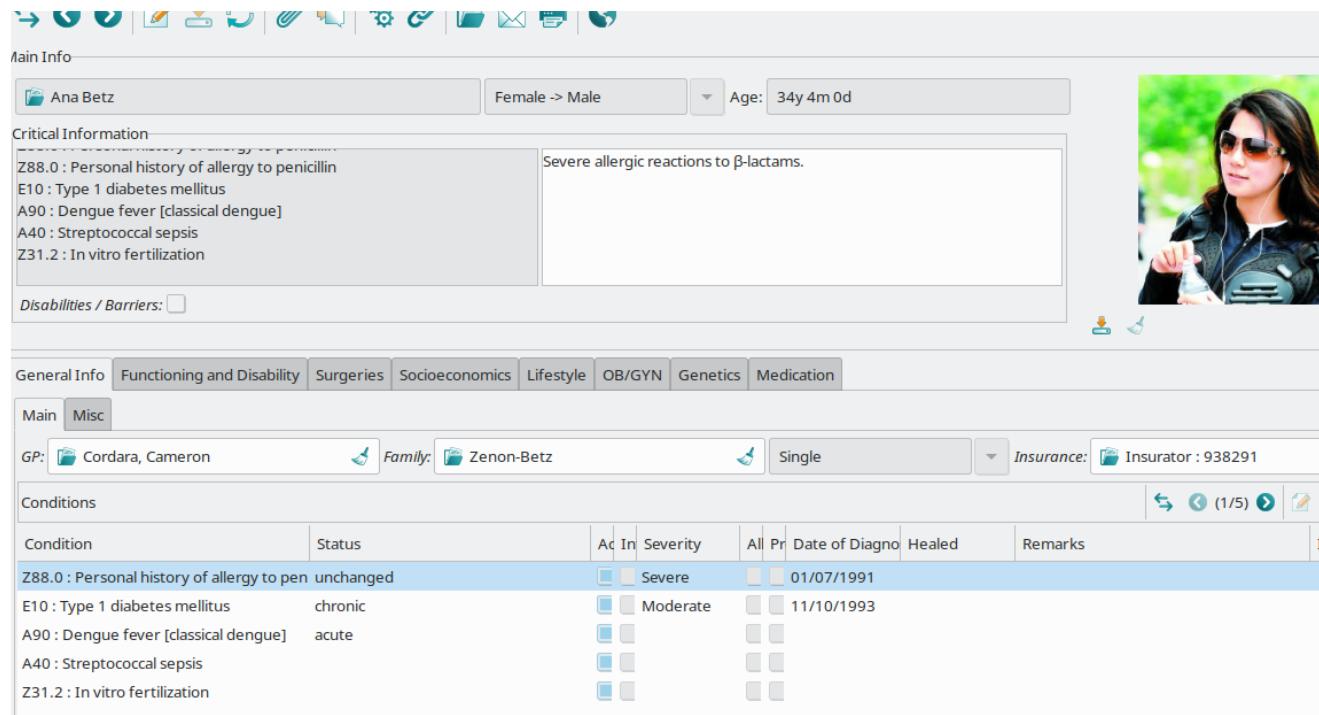
Learning an unbiased estimator

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

- Learn $\tilde{\eta}(X)$ using any ML algorithm which returns calibrated classifiers. Substitute $\tilde{\eta}(X)$ in the above equation to get an unbiased estimate of $\eta(X)$!

Application to electronic phenotyping

Hundreds of relevant clinical variables



The screenshot shows a patient record for Ana Betz. The top bar includes a toolbar with various icons and dropdown menus. The main info section shows the patient's name, gender (Female), age (34y 4m 0d), and a note about severe allergic reactions to β-lactams. A photograph of a woman wearing sunglasses and a dark jacket is displayed. Below this are tabs for General Info, Functioning and Disability, Surgeries, Socioeconomics, Lifestyle, OB/GYN, Genetics, and Medication. The General Info tab is selected, showing details like GP (Cordara, Cameron), Family (Zenon-Betz), marital status (Single), and insurance (Insurator: 938291). The Conditions section lists medical history items such as Z88.0 (Personal history of penicillin allergy), E10 (Type 1 diabetes mellitus), A90 (Dengue fever), A40 (Streptococcal sepsis), and Z31.2 (In vitro fertilization), each with status, severity (Severe or Moderate), and date of diagnosis.

Figure source:

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

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
...

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?		
	T	F
T	297	129
F	1,319	34511

physician response
(gold standard)

PPV
0.70

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, but can be used for labels
- Reasonable to assume that $\rho_- = P(\tilde{Y} = 1 | Y = -1) \approx 0$, 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 as “true” labels
2. Learn a classifier to predict whether the *anchor* \tilde{Y} appears
3. Calibration step: divide by $\frac{1}{|P|} \sum_P P(\tilde{Y} = 1 | X)$

Test time

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

1. Predict using the learned classifier (with calibration)

Often we can find noisy labels WITHIN the data!

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

- 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

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***
3. Calibration step: divide by $\frac{1}{|P|} \sum_P P(\tilde{Y} = 1 | X)$

Test time

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

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

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



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

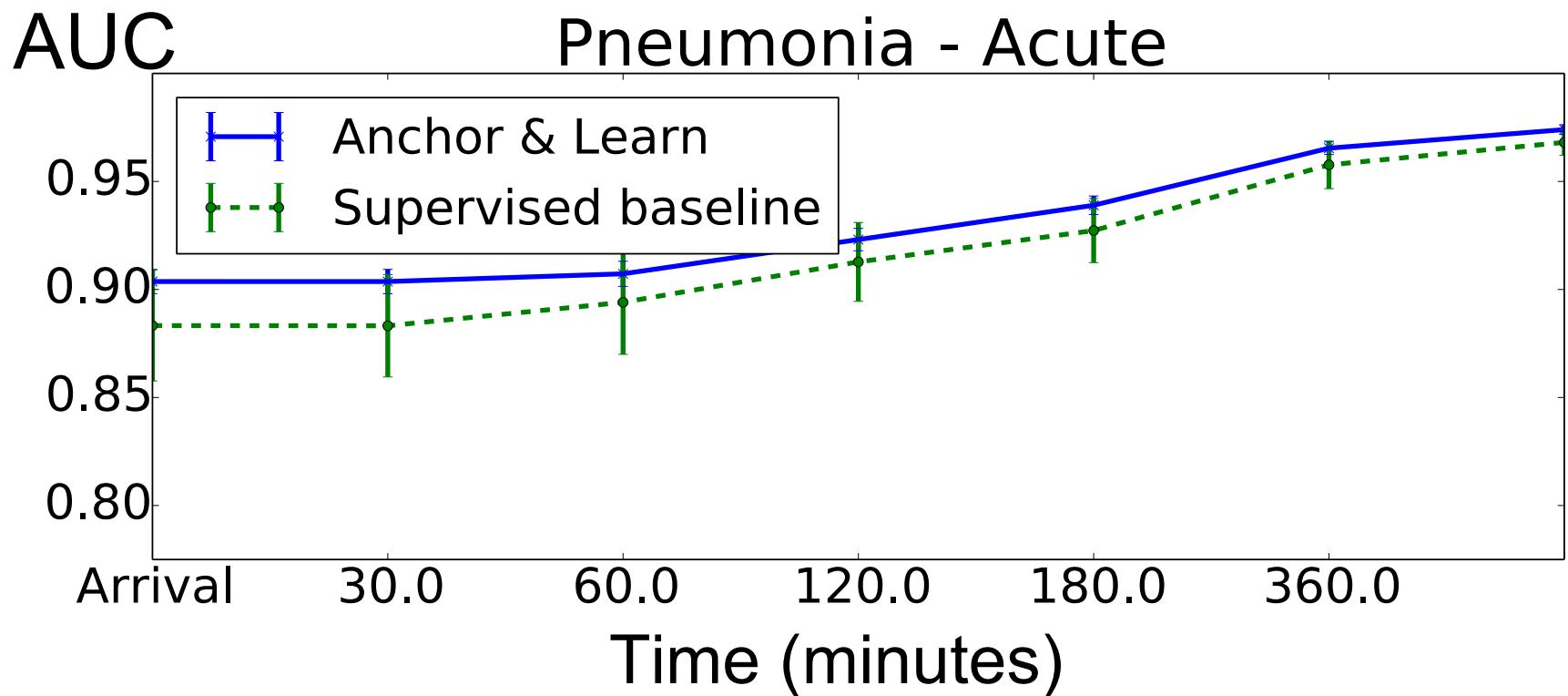
Does the patient have an active malignancy?ⁱ

Unlikely Unsure Likely

<- Previous Abort Next ->



Evaluating phenotypes



Comparison to supervised learning using labels for
5000 patients

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+

nstemi

stemi

ntg

lasix

nitro

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

beta blockers

cardiac medicine

BIDMC shortform

Highly weighted terms

Ages

age=80-90

age=70-80

age=90+

nstemi

stemi

ntg

lasix

nitro

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