

Efficient Development of Electronic Health Record Based Algorithms to Identify Rheumatoid Arthritis

Robert Carroll, Vanderbilt University

Rheumatoid Arthritis

- Prototypic chronic disease
- Inflammatory polyarthritis afflicting 1.3 million adults in the United States¹
- 50% increased risk of premature mortality²
- Life expectancy is decreased by 3 to 10 years compared with the general population²
- Autoimmune disorder
- Primarily affects joints

1. Helmick CG, Felson DT, Lawrence RC, Gabriel S, Hirsch R, Kwok CK, et al. Estimates of the prevalence of arthritis and other rheumatic conditions in the United States. Part I. Arthritis Rheum. 2008 Jan;58(1):15-25.

2. Myasoedova E, Davis JM 3rd, Crowson CS, Gabriel SE. Epidemiology of rheumatoid arthritis: rheumatoid arthritis and mortality. Curr Rheumatol Rep. 2010 Oct;12(5):379-385.

How to find RA in the EHR

- Structured data
 - Billing Codes
 - Lab results
- Medications
 - Codified from the order entry system
 - Free text from other sources
- Free text data
 - Physician notes
 - Radiology reports

Billing data

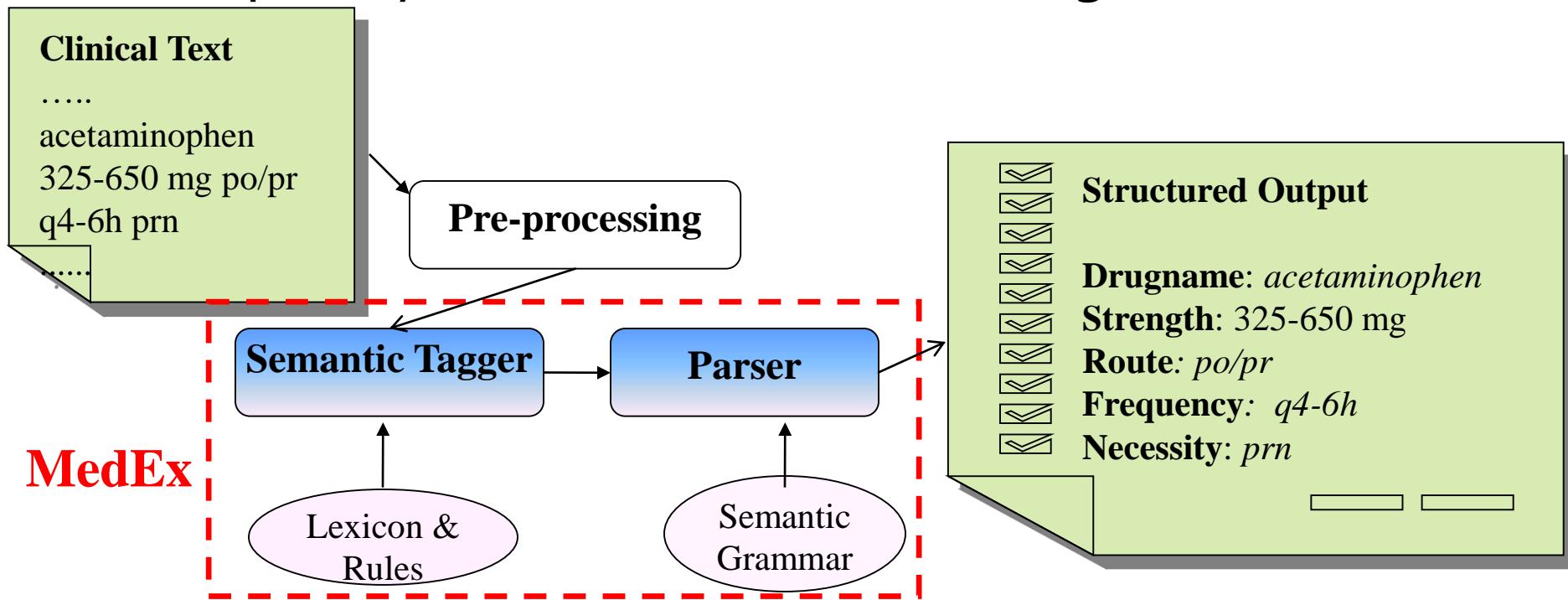
- Why not just look for billing codes?
 - A VA database study showed that only 57% of individuals with one RA ICD-9 billing code were RA positive.¹
 - Another study at Partners HealthCare showed that only 56% of individuals with three or more codes were RA positive.²

1. Singh JA, Holmgren AR, Noorbaloochi S. Accuracy of Veterans Administration databases for a diagnosis of rheumatoid arthritis. *Arthritis Rheum* 2004;51:952–7.

2. Liao KP, Cai T, Gainer V, et al. Electronic medical records for discovery research in rheumatoid arthritis. *Arthritis Care Res (Hoboken)* 2010;62:1120–7.

MedEx

- Medex identifies medication references from clinical text.
- It additionally records the dose, route, and frequency information for the drug reference.



KnowledgeMap Concept Identifier (KMCI)

- Identifies UMLS concepts from text
- It uses context to help disambiguate words and phrases.
- Incorporates:
 - NegEx - identifies negation
 - SecTag - identifies note sections, aiding in context determination. It also helps separate out Family History from the rest of the patient's information.

How to find RA in the EHR

- Codified data
 - Billing Codes
 - Lab results
- Medications
 - Codified from the order entry system
 - Free text from other sources
- Free text data
 - Physician notes
 - Radiology reports

Phenotype algorithms

- We now have access to multimodal data.
- How can we incorporate it all?
- Deterministic Algorithm

Apply a Deterministic Algorithm

- Development cycle
 - Create a set of rules
 - Iterate:
 - Generate a set of cases and controls
 - Judge the algorithm
 - If performance has not reached threshold, modify the rules
- This method did replicate a known genetic association for RA.

RA Deterministic Algorithm

ICD 9 codes (any of the below)

- 714 Rheumatoid arthritis and other inflammatory polyarthropathies
- 714.0 Rheumatoid arthritis
- 714.1 Felty's syndrome
- 714.2 Other rheumatoid arthritis with visceral or systemic involvement

AND

Medications (any of the below)

methotrexate [MTX][amethopterin] sulfasalazine [azulfidine]; Minocycline [minocin][solodyn]; hydroxychloroquine [Plaquenil]; adalimumab [Humira]; etanercept [Enbrel] infliximab [Remicade]; Gold [myochrysine]; azathioprine [Imuran]; rituximab [Rituxan] [MabThera]; anakinra [Kineret]; abatacept [Orencia]; leflunomide [Arava]

AND

Keywords (any of the below)

rheumatoid [rheum] [reumatoid] arthritis [arthritides] [arthriris] [arthristis] [arthritus] [arthrtis] [arthritis]

RA Deterministic Algorithm

Exclusions

AND NOT
ICD 9 codes (any of the below)

- 714.30 Polyarticular juvenile rheumatoid arthritis, chronic or unspecified
- 714.31 Polyarticular juvenile rheumatoid arthritis, acute
- 714.32 Pauciarticular juvenile rheumatoid arthritis
- 714.33 Monoarticular juvenile rheumatoid arthritis
- 695.4 Lupus erythematosus
- 710.0 Systemic lupus erythematosus
- 373.34 Discoid lupus erythematosus of eyelid
- 710.2 Sjogren's disease
- 710.3 Dermatomyositis
- 710.4 Polymyositis
- 555 Regional enteritis
- 555.0 Regional enteritis of small intestine
- 555.1 Regional enteritis of large intestine
- 555.2 Regional enteritis of small/large intestine
- 555.9 Regional enteritis of unspecified site
- 564.1 Irritable Bowel Syndrome
- 135 Sarcoidosis
- 696 Psoriasis and similar disorders
- 696.0 Psoriatic arthropathy
- 696.1 Other psoriasis and similar disorders excluding psoriatic arthropathy
- 696.8 Other psoriasis and similar disorders
- 099.3 Reiter's disease
- 716.8 Arthropathy, unspecified
- 274.0 Gouty arthropathy
- 358.0 myasthenia gravis
- 358.00 myasthenia gravis without acute exacerbation
- 358.01 myasthenia gravis with acute exacerbation
- 775.2 neonatal myasthenia gravis
- 719.3 Palindromic rheumatism
- 719.30 Palindromic rheumatism, site unspecified
- 719.31 Palindromic rheumatism involving shoulder region
- 719.32 Palindromic rheumatism involving upper arm
- 719.33 Palindromic rheumatism involving forearm
- 719.34 Palindromic rheumatism involving hand
- 719.35 Palindromic rheumatism involving pelvic region and thigh
- 719.36 Palindromic rheumatism involving lower leg
- 719.37 Palindromic rheumatism involving ankle and foot
- 719.38 Palindromic rheumatism involving other specified sites
- 719.39 Palindromic rheumatism involving multiple sites
- 720 Ankylosing spondylitis and other inflammatory spondylopathies
- 720.0 Ankylosing spondylitis
- 720.8 Other inflammatory spondylopathies
- 720.81 Inflammatory spondylopathies in diseases classified elsewhere
- 720.89 Other inflammatory spondylopathies
- 720.9 Unspecified inflammatory spondylopathy
- 721.2 Thoracic spondylosis without myelopathy
- 721.3 Lumbosacral spondylosis without myelopathy
- 729.0 Rheumatism, unspecified and fibrosis
- 340 Multiple sclerosis
- 341.9 Demyelinating disease of the central nervous system unspecified
- 323.9 transverse myelitis
- 710.1 Systemic sclerosis
- 245.2 Hashimoto's thyroiditis
- 242.0 Toxic diffuse goiter
- 443.0 Raynaud's syndrome

OR

Keywords (any of the below)

juvenile [juv] rheumatoid [rheum] [reumatoid] [rhumatoid] arthritis [arthritides] [arthriris] [arthristis] [arthritis] [arthrtis] [arthritis]
juvenile [juv] arthritis arthritis [arthritides] [arthriris] [arthristis] [arthritis] [arthrtis] [arthritis]
juvenile chronic arthritis [arthritides] [arthriris] [arthristis] [arthritis] [arthrtis] [arthritis]
juvenile [juv] RA; JRA
Inflammatory [inflammatory] [inflam] osteoarthritis [osteoarthrosis] [OA]
Reactive [psoriatic] arthritis [arthropathy] [arthritides] [arthriris] [arthristis] [arthritis] [arthrtis] [arthritis]

Phenotype algorithms

- Deterministic
 - Provides reasonable performance
 - Relies on the expert designer
- Machine Learning (ML)

Applying ML Algorithms

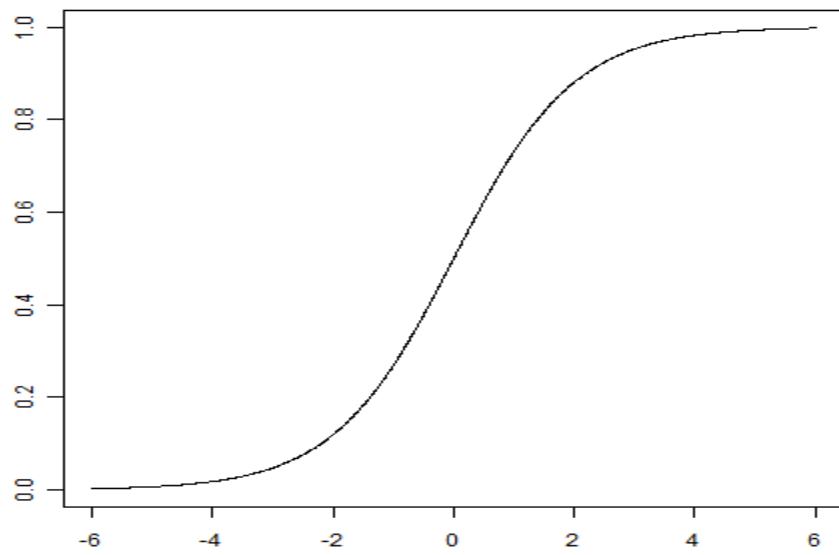
- Select attributes of interest, including related conditions or potential misdiagnoses
- Review subjects for training set
- Train a model and apply it to the record pool
- Select predictive positive patients and review
- Ensure performance meets standard
- This method has also been shown to be effective in genetic studies

Liao KP, Cai T, Gainer V, et al. Electronic medical records for discovery research in rheumatoid arthritis. *Arthritis Care Res (Hoboken)* 2010;62:1120–7

Kurreeman F, Liao K, Chibnik L, Hickey B, Stahl E, Gainer V, et al. Genetic basis of autoantibody positive and negative rheumatoid arthritis risk in a multi-ethnic cohort derived from electronic health records. *Am. J. Hum. Genet.* 2011 Jan 7;88(1):57–69.

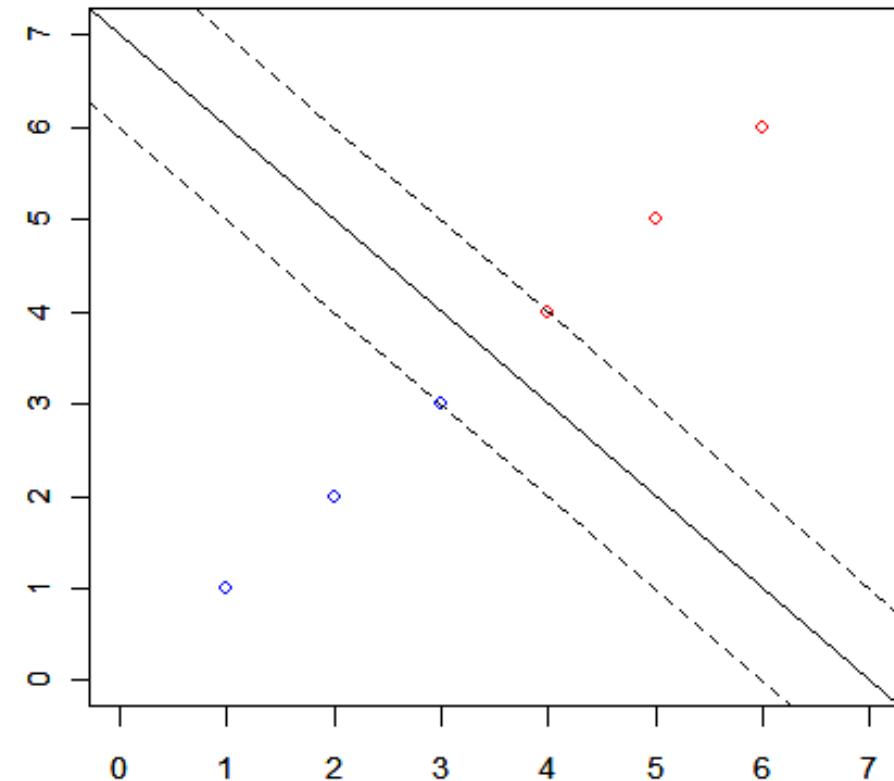
Logistic Regression

- Two-group classifier
- Finds the optimal weight for every attribute to best predict the class
 - $Z = \text{Intercept} + W_1 * \text{Attribute}_1 + W_2 * \text{Attribute}_2 + \dots$
 - Maps the Z score onto a logistic curve
 - $F(Z) = 1/(1+e^{-Z})$
 - Outputs a probability
 - Weights have some interpretation



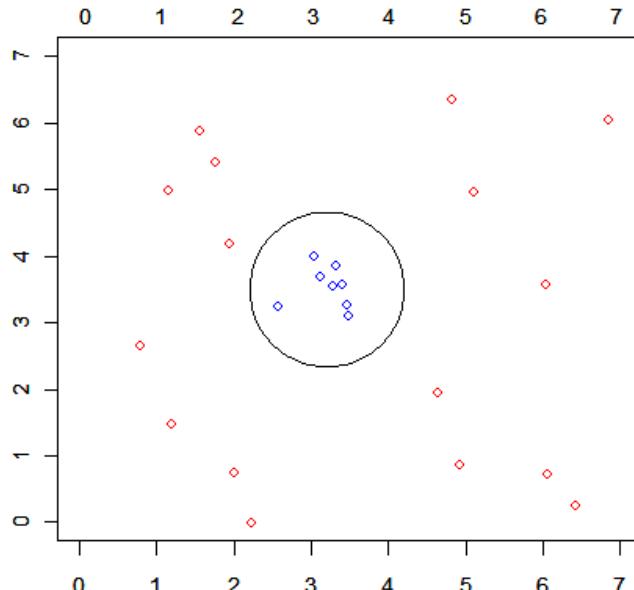
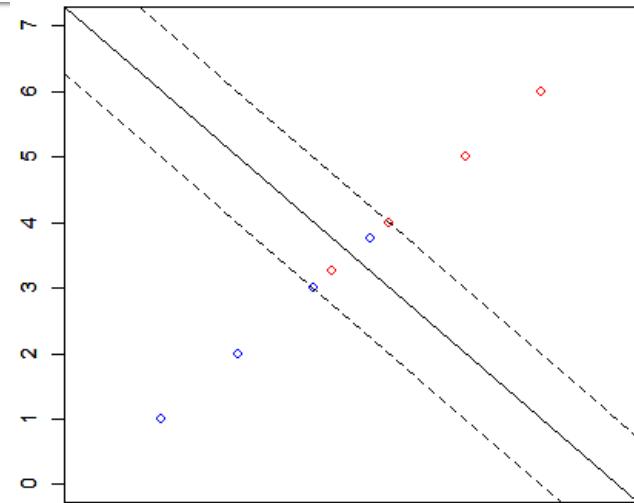
Support Vector Machine

- Two-group classifier
- Finds the best hyperplane to separate the two classes
 - Each attribute becomes a dimension
 - In the case of two dimensions, the hyperplane is a line.
 - Maximizes the distance from the points to the hyperplane.



SVM Modifications

- Soft margins: Allows points to be within the margin or across the hyperplane, but with a penalty.
- Kernel methods: Maps the attributes to another dimension, allowing for interactions between the attributes

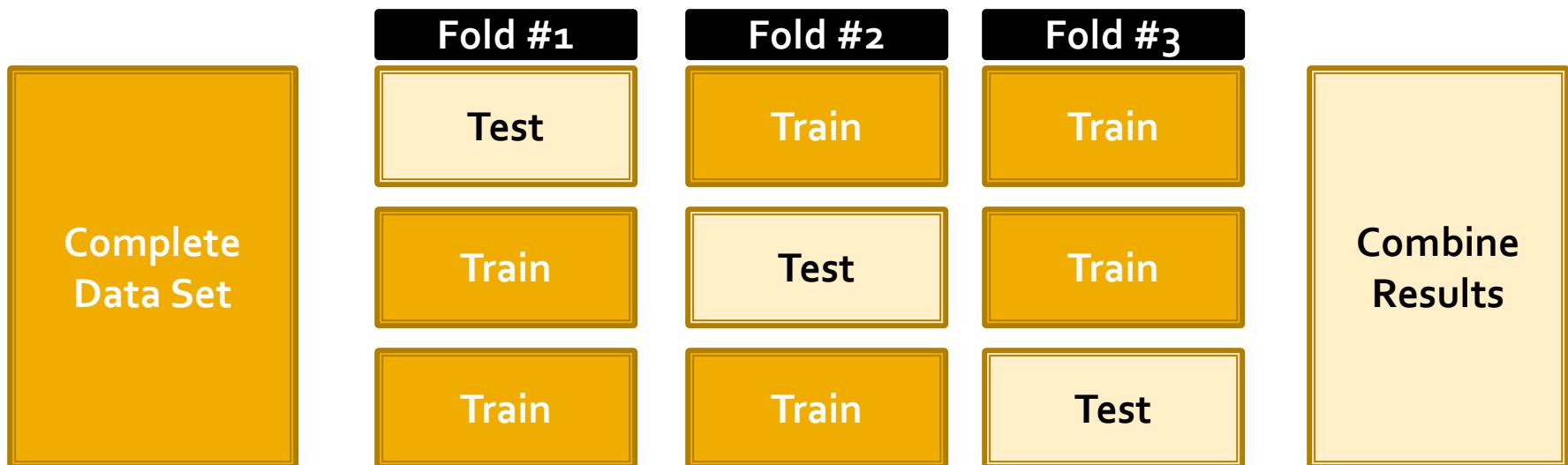


Overfitting

- Machine learning algorithms are designed to identify patterns.
- They can only learn what they are shown however.
 - Their behavior is not always well defined outside of the training space.
 - They can't tell the difference between erroneous patterns and clinically significant patterns.
- More attributes means more coincidental patterns, so choosing relevant attributes can help limit overfitting the training data.

Cross Validation

- Used to estimate more accurately the performance for algorithms being trained and tested on the same pool of data.
- Splits up the data set into groups.
 - 3-fold CV means three different groups
 - Train on two groups and test on the third
 - Repeat three times



Phenotype algorithms

- Deterministic
 - Provides reasonable performance
 - Relies on the expert designer
- Machine Learning (ML)
 - Allows for more complex relationships
 - Reliance on the perceived relevance of attributes
- The naïve approach

Naïve Electronic Health Record Phenotype Identification for Rheumatoid Arthritis

Hypothesis

A support vector machine (SVM) trained using a naïve data set can identify rheumatoid arthritis (RA) cases as well as one trained using an expert defined set of attributes.

Methods

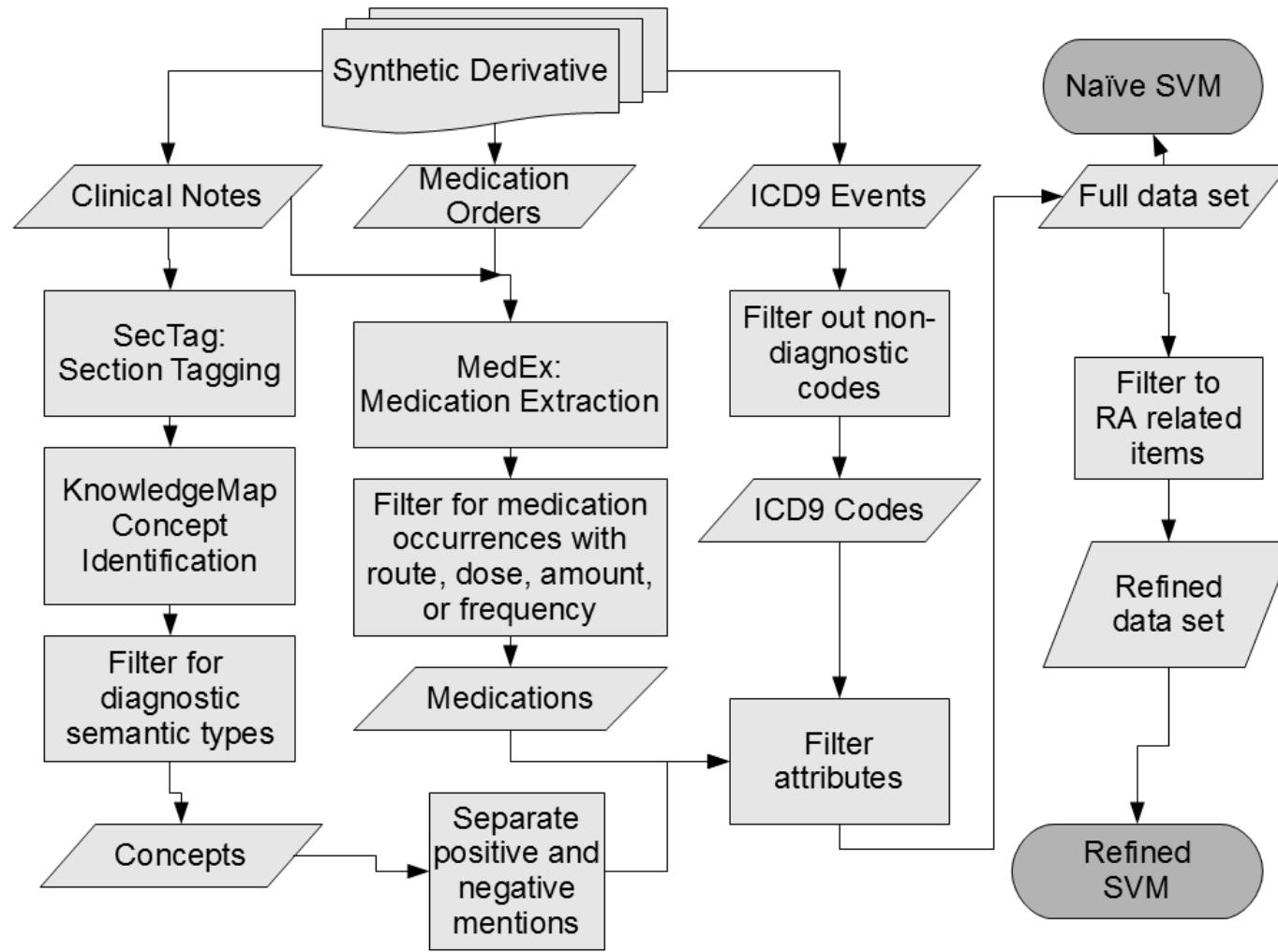
- Cohort
 - 376 records
 - Rheumatologist reviewed
 - Individuals with at least one RA ICD-9
- EHR-derived attributes
 - ICD-9 counts
 - Medication Entries
 - Unified Medical Language System (UMLS) Concepts

Our population

	Vanderbilt (n=376)	
	RA	Non-RA
Total	185 (49.2%)	191 (50.8%)
Age	52.9 ± 13.1	56.2 ± 16.5
Female	148 (80.0%)	141 (73.8%)
Ethnicity		
Caucasian	143 (77.3%)	155 (81.2%)
African American	14 (7.6%)	26 (13.6%)
Hispanic	1 (0.5%)	1 (0.5%)
Other	3 (1.6%)	2 (1.0%)
Unknown	24 (13.0%)	7 (3.7%)
Drugs		
Anti-TNF use	88 (47.6%)	26 (13.6%)
MTX	133 (71.9%)	63 (33.0%)
Codes		
RA	185 (100.0%)	191 (100.0%)
SLE	14 (7.6%)	32 (16.8%)
JRA	6 (3.2%)	8 (4.2%)
PsA	6 (3.2%)	14 (7.3%)
EHR Followup*	9.97 ± 4.06	9.06 ± 4.32

*Mean ± standard deviation in years

Attribute Creation



Methods: Evaluation

- Algorithms
 - Prior Published Deterministic¹
 - Support Vector Machine
 - Naïve and Gaussian kernel
 - Refined and Gaussian kernel
- Method
 - 10 fold cross validation
 - Nested CV for gamma and cost parameters
 - Subset evaluations (e.g., ICD9 only, Meds only)

1. Ritchie MD, et al. Am. J. Hum. Genet. 2010;86(4):560-572

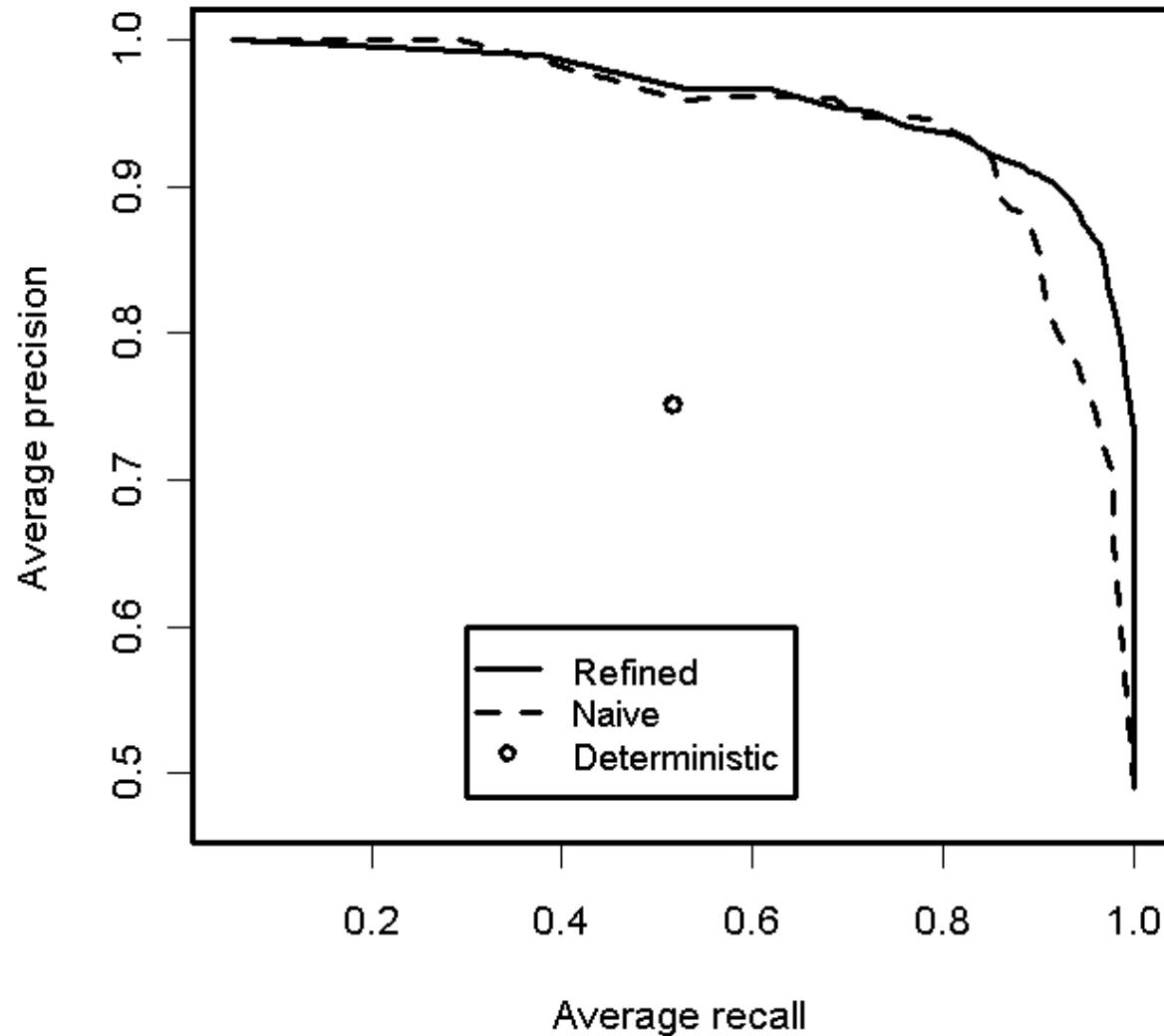
Measurements

- Recall (Sensitivity)
 - $tp/(tp+fn)$
- Precision (Positive Predictive Value):
 - $tp/(tp+fp)$
- Specificity
 - $tn/(tn+fp)$
- F-measure
 - $2*(precision*recall)/(precision+recall)$
- Area under the receiver operator characteristic curve (AUC)

Model Performance

Naïve	Precision	Recall	F measure	AUC	Attributes
Full	93.3 ± 0.5	79.7 ± 5.2	85.1 ± 3.7	94.2 ± 1.3	17110
ICD-9	94.1 ± 0.2	87.1 ± 2.8	90.3 ± 1.6	95.6 ± 1.0	795
NLP	92.2 ± 0.6	68.2 ± 5.6	77.4 ± 4.1	90.4 ± 2.1	15171
Medication	88.9 ± 1.8	51.0 ± 5.4	63.5 ± 5.5	84.6 ± 2.6	1148
Refined	Precision	Recall	F measure	AUC	Attributes
Full	93.7 ± 0.6	85.8 ± 5.7	88.6 ± 4.0	96.6 ± 1.1	59
ICD-9	93.2 ± 0.5	78.1 ± 5.2	84.2 ± 3.5	95.5 ± 1.3	12
NLP	91.8 ± 1.0	68.8 ± 7.5	76.8 ± 5.7	89.5 ± 2.1	33
Medication	86.6 ± 1.6	40.5 ± 5.4	53.8 ± 5.2	83.3 ± 2.5	18
Deterministic	Precision	Recall	F measure	AUC	Attributes
Full	75.2 ± 2.5	51.6 ± 2.6	60.5 ± 2.6	N/A	N/A

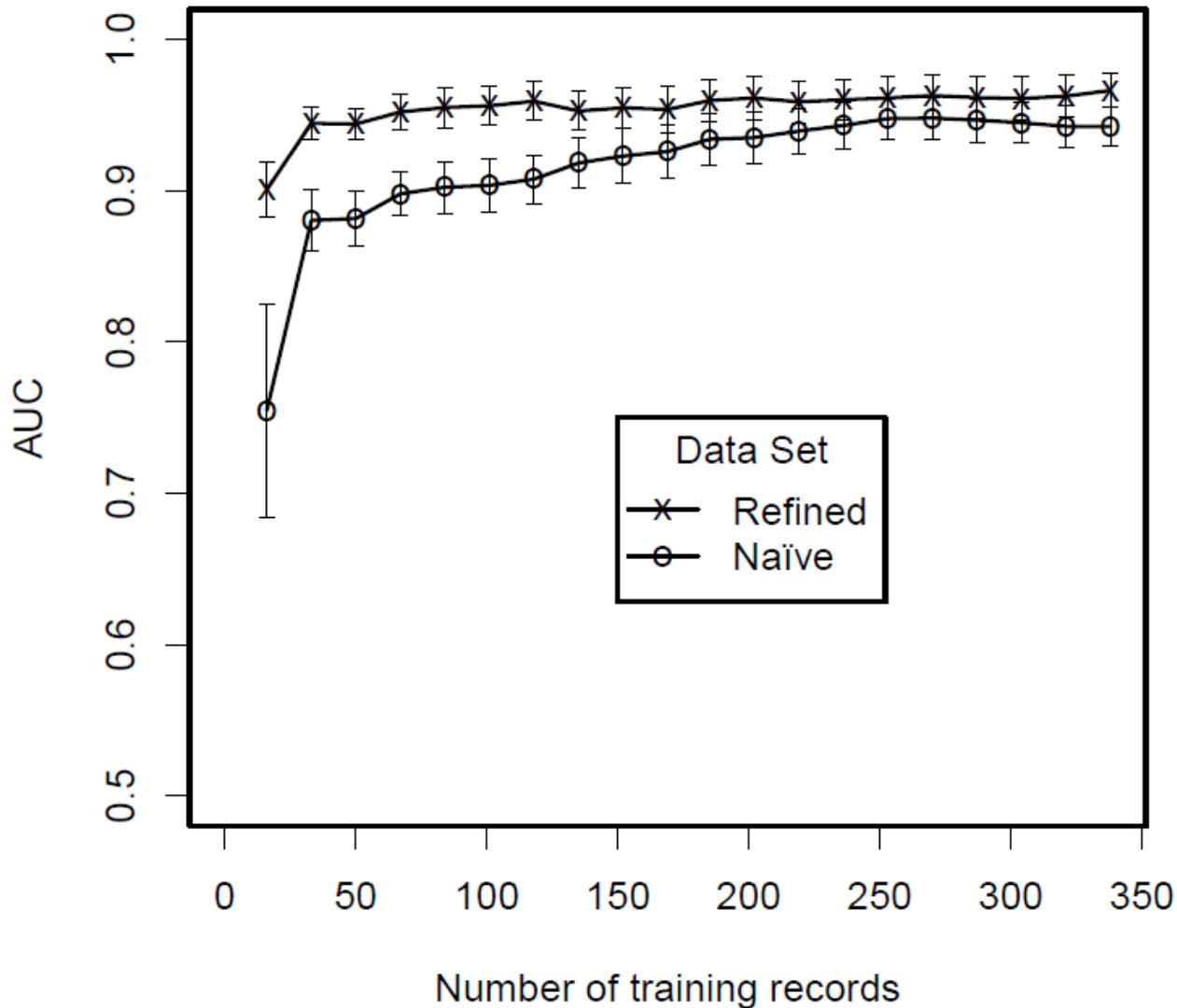
Precision Recall Curves



Evaluation: How many cases are needed to train?

- We used the same cross validation, but repeated it 20 times.
- Each time we increased the number of training samples available to the fold
- Approximated an iterative algorithm development approach

How many cases are required to train an accurate algorithm?



Discussion

- Both naïve and refined SVMs performed better than expert-designed algorithms
- Naïve SVMs can substitute for refined SVMs
- Naïve SVMs have some disadvantages
 - They require more training (100 cases vs. 20)
 - They have lower recall (6%)
- The NLP and medication methods had large standard error
- Medications alone had low recall and poor precision

Limitations

- RA is well represented by distinct ICD9 codes (though only 50% of cases were positive)
- RA is a chronic condition, yielding many relevant clinical encounters over time
- Only data from one institution was analyzed

Portability of an Algorithm to Identify Rheumatoid Arthritis in Electronic Health Records

Are phenotype algorithms portable to different EHRs?

- Previously published logistic regression model
 - Trained at Partners Healthcare
 - ICD9, Labs, Meds, NLP
 - Tested at Northwestern and Vanderbilt
- Is the disease signature the same across
 - Differing healthcare environments?
 - Differing EHR systems?
 - Differing information extraction techniques?

Hypothesis

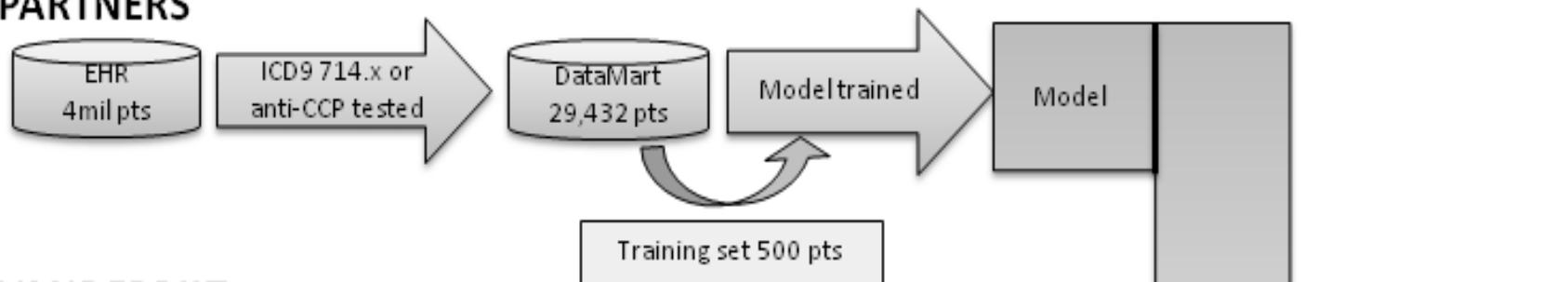
- In spite of differing healthcare environments, EHR systems, and information extraction techniques, a logistic regression model trained at one institution will predict RA status at other institutions.

Overview of underlying differences

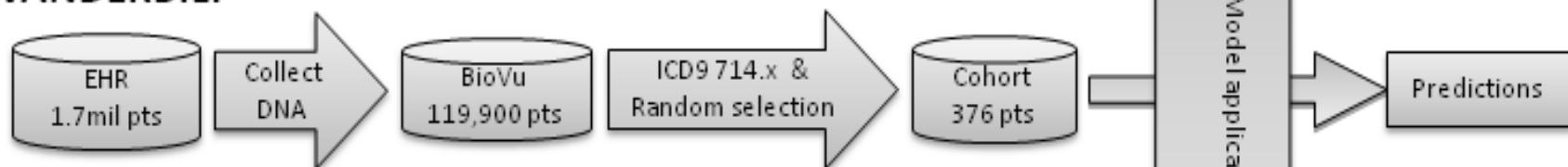
	Implementations by Institution		
	Partners Boston, MA	Northwestern Chicago, IL	Vanderbilt Nashville, TN
EHR system	Internally-developed	EpicCare (Outpatient) and Cerner PowerChart (Inpatient)	Internally-developed
Number of patients	4 million	2.2 million	1.7 million
Research EHR data	Enterprise Data Warehouse	Enterprise Data Warehouse	De-identified image of EHR (Synthetic Derivative)
Medication Source	Structured medication entries (inpatient and outpatient) and text queries	Structured outpatient medication entries and inpatient and outpatient text queries	NLP (MedEx) for outpatient medications and structured inpatient records
NLP system (disease concepts, lab results, medications, erosions)	HITEx	HITEx	KnowledgeMap Concept Identifier
NLP concept queries	Customized RegEx queries	Customized RegEx queries from Partners	Generic UMLS concepts, derived from KnowledgeMap web interface

Overview

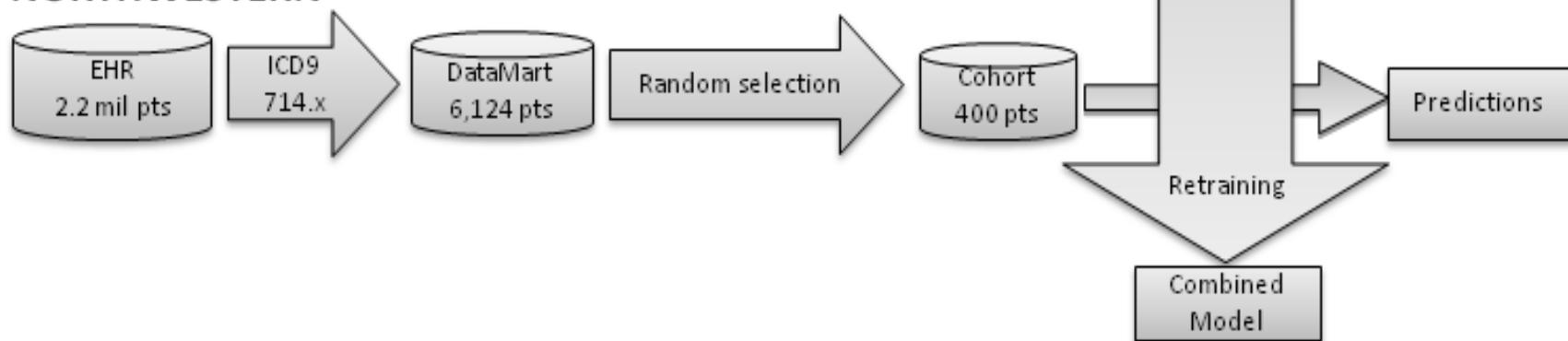
PARTNERS



VANDERBILT



NORTHWESTERN



Methods: Partners RA algorithm

- Selected and grouped ICD-9 codes for four diseases: RA, Juvenile Rheumatoid Arthritis, Psoriatic Arthritis, and Systemic Lupus Erythematosus
- Chose three categories of RA medications.
- Uses some NLP data for each of the four diseases as well.

Methods

- Physician reviewed charts are the gold standard.
- Parameter selection and weight minimization using lasso.
- Five-fold cross validation.
- We reworked the “fact” normalization to use the more ubiquitous ICD-9 code count.

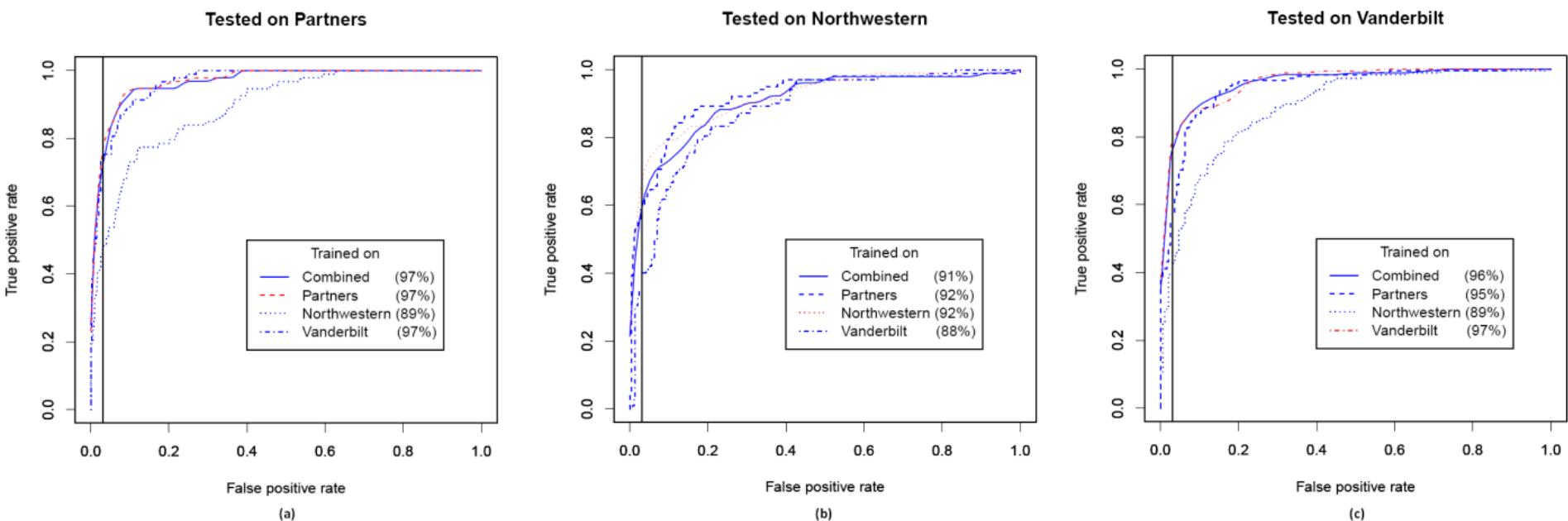
Patient Demographics

	Partners (n=500)		Northwestern (n=390)		Vanderbilt (n=376)	
	RA	Non-RA	RA	Non-RA	RA	Non-RA
Total	96 (19.2%)	404 (80.8%)	102 (26.2%)	288 (73.8%)	185 (49.2%)	191 (50.8%)
Age	60.7 ± 15.9	56.0 ± 18.6	54.3 ± 14.8	58.9 ± 16.8	52.9 ± 13.1	56.2 ± 16.5
Female	74 (77.1%)	303 (75.0%)	83 (81.4%)	209 (72.6%)	148 (80.0%)	141 (73.8%)
Ethnicity						
Caucasian	64 (66.7%)	286 (70.8%)	40 (39.2%)	120 (41.7%)	143 (77.3%)	155 (81.2%)
African American	3 (3.1%)	46 (11.4%)	18 (17.6%)	46 (16.0%)	14 (7.6%)	26 (13.6%)
Hispanic	2 (2.1%)	29 (7.2%)	6 (5.9%)	18 (6.3%)	1 (0.5%)	1 (0.5%)
Other	6 (6.3%)	7 (1.7%)	13 (12.7%)	44 (15.3%)	3 (1.6%)	2 (1.0%)
Unknown	21 (21.9%)	36 (8.9%)	25 (24.5%)	60 (20.8%)	24 (13.0%)	7 (3.7%)
Drugs						
Anti-TNF use	50 (52.1%)	50 (12.4%)	67 (65.7%)	37 (12.8%)	88 (47.6%)	26 (13.6%)
MTX	77 (80.2%)	105 (26.0%)	70 (68.6%)	61 (21.2%)	133 (71.9%)	63 (33.0%)
Codes						
RA	93 (96.9%)	329 (81.4%)	102 (100.0%)	283 (98.3%)	185 (100.0%)	191 (100.0%)
SLE	2 (2.1%)	37 (9.2%)	3 (2.9%)	22 (7.6%)	14 (7.6%)	32 (16.8%)
JRA	7 (7.3%)	28 (6.9%)	1 (1.0%)	18 (6.3%)	6 (3.2%)	8 (4.2%)
PsA	2 (2.1%)	21 (5.2%)	0 (0.0%)	12 (4.2%)	6 (3.2%)	14 (7.3%)
EHR Followup*	9.38 ± 6.77	10.14 ± 6.85	6.30 ± 4.69	6.05 ± 4.85	9.97 ± 4.06	9.06 ± 4.32

Results

Algorithm	Testing Set											
	Partners			Northwestern			Vanderbilt			Average		
	PPV	Sens	AUC	PPV	Sens	AUC	PPV	Sens	AUC	PPV	Sens	AUC
Published Algorithm	88%*	79%*	97%*	87%	60%	92%	95%	57%	95%	90%	65%	95%
Retrained with:												
Northwestern	79%	47%	89%	87%	73%	92%	93%	43%	89%	86%	54%	90%
Vanderbilt	85%	74%	97%	82%	40%	88%	97%	81%	97%	88%	65%	94%
Combined	86%	71%	97%	86%	65%	91%	97%	82%	96%	90%	72%	95%
ICD-9 Only: [‡]												
>1 RA code	22%	97%	N/A	26%	100%	N/A	49%	100%	N/A	33%	99%	N/A
>3 RA code	55%	81%	N/A	42%	87%	N/A	73%	98%	N/A	57%	89%	N/A
>Optimal	80%	49%	88%	80%	36%	84%	93%	43%	93%	84%	43%	88%
Optimal Code Count	53			29			48			43.3		

ROC Curves



Discussion

- These results show that a previously published logistic regression method developed at one institution is portable.
- The published logistic regression model improved sensitivity by 22% and PPV by 6% compared to the optimal ICD-9 count threshold, demonstrating the added value of the more complex phenotyping algorithm.
- Northwestern has a shorter follow-up time for their patients.

Discussion- Normalization and Portability

- The original algorithm normalized an attribute by the number of “facts”.
- Every piece of data available was a fact- from icd-9 codes, to lab results, and NLP results.
- This is hard to replicate, but normalizing by the record size is important.
- We selected a new proxy for record size, the number of ICD-9 codes.
- Using linear regression, we were able to approximate the number of facts and still apply the originally trained model.

Discussion- Betas

Attribute	Description	Original	Retrained		
			Combined	Vanderbilt	Northwestern
(Intercept)	Regression Intercept	-5.2088	-4.00186	-4.75161	-2.49521
age	Age of the patient	-0.00096	-0.00426	0.010474	-0.00769
sex	Binary: Female is 1	-0.10729	-0.15874	-0.10372	0
ICD-9 RA	Number of encounters with the specified billing code. Defined as sets of ICD-9s >7 days apart. Natural log transformed	0.639367	0.786732	1.036392	0.234517
ICD-9 PsA		0	-0.44454	0	-0.78851
ICD-9 SLE		-0.95937	-0.36747	-0.12847	-0.37016
ICD-9 JRA		-2.25118	-0.67657	-0.49168	0
Normalized ICD-9 RA	ICD-9 RA before transformation, normalized by the total ICD-9 counts.	66.02406*	91.33932	123.3725	18.72591
Methotrexate	Binary variable. Denotes exposure to this medication from codified sources.	0	0	0	0.541961
Anti-TNF		0.958811	0.745813	0	0.818504
Other Medications		0	0	0	0.147411
RF Negative	Binary variable for the Rheumatoid Factor test.	0.850944	-0.42402	-0.32315	-0.58007
RF Positive		0	0.748071	0.795551	0
NLP RA	Natural log transformed count of the number of notes with the specified concept.	0.969956	0.733087	0.645841	0.914909
NLP SLE		-0.52562	-0.21839	-0.09926	0
NLP JRA		0	-1.00356	-1.27468	0
NLP PsA		-0.85581	-0.05785	0	0
Methotrexate	Binary variable. Denotes exposure to this medication from narrative sources.	0.631764	0.442066	0	0
Anti-TNF		0.520743	0.321728	0	0.908849
Other Medications		0.298111	0.479028	0	0.020784
Cyclic citrullinated Peptide	Binary variable. Denotes positive mention of this test in narrative sources.	1.312583	0.701096	0	0
NLP RF	Binary variable. Denotes positive mention of this test in narrative sources.	0	-0.28693	0	1.279047
Seropositive	Binary variable. Denotes positive mention of this term in narrative sources.	2.773642	1.04373	0	0.16429
Erosions	Binary variable. Denotes positive mention of this finding in narrative sources.	1.259249	0.540583	0.464227	0

- The normalized RA count weight was very big, but the value was also very small
- Weights are very different, and sometimes opposite, but the model still predicted well.
- Vanderbilt has no medication attributes that were selected by the model.

Partners algorithm: Lots of regular expressions!

nlpRF	RFPositive	(?i)(?m)(?s)(borderline\s+ slightly\s+ strongly\s+ low\s+ weakly\s+ low\s*-\?\s*titer\s+ high\s*-\?\s*titer\s+)?(\+\s* pos(itive itivity)\?\s+ elevated\s+ increased\s+)((test\s+)?for\s+)?(IgG\s+)?((anti\s*-\?\s*CCP\s+antibod(ies y) anti\s*-\?\s*CCP CCP\s+antibod(ies y) (anti\s*-\?\s*)?cyclic\s+citrullinated\s+peptide(\s+antibod(ies y)) a(\s*-\s*)?CCP anti\s*-\?\s*citrullinated\s+protein\s*? \s*peptide\s+antibod(ies y) anti\s*-\?\s*citrullinated\s+(protein peptide)\s+antibod(ies y) ACPA ACP\s+antibod(ies y) CCP)\s+and\s+)?(\bRF rheumatoid\s+factor(s)?)
nlpRF	RFPositive	(?i)(?m)(?s)(\bRF rheumatoid\s+factor(s)?)\s+and\s+(anti\s*-\?\s*CCP\s+antibod(ies y) anti\s*-\?\s*CCP CCP\s+antibod(ies y) (anti\s*-\?\s*)?cyclic\s+citrullinated\s+peptide(\s+antibod(ies y)) a(\s*-\s*)?CCP anti\s*-\?\s*citrullinated\s+protein\s*? \s*peptide\s+antibod(ies y) anti\s*-\?\s*citrullinated\s+(protein peptide)\s+antibod(ies y) ACPA ACP\s+antibod(ies y) CCP)\s+is\s+ was\s+ are\s+ has\s+ also\s+)?been\s+)?(also\s+)?(found\s+)?(to\s+be\s+)?(also\s+)?(borderline\s+ slightly\s+ strongly\s+ low\s+ weakly\s+ low\s*-\?\s*titer\s+ high\s*-\?\s*titer\s+)?(present\b +)\s+positive\s+positivity\s+pos\b elevated\s+increased\s+greater\s+than\s+\d{2,})
nlpRF	RFPositive	(?i)(?m)(?s)(\bRF rheumatoid\s+factor(s)?)\s+and\s+(anti\s*-\?\s*CCP\s+antibod(ies y) anti\s*-\?\s*CCP CCP\s+antibod(ies y) (anti\s*-\?\s*)?cyclic\s+citrullinated\s+peptide(\s+antibod(ies y)) a(\s*-\s*)?CCP anti\s*-\?\s*citrullinated\s+protein\s*? \s*peptide\s+antibod(ies y) anti\s*-\?\s*citrullinated\s+(protein peptide)\s+antibod(ies y) ACPA ACP\s+antibod(ies y) CCP)\s*+
nlpRF	RFPositive	(?i)(?m)(?s)(?<!(no\s{1,100}))\s+of\s+circulating\s+ presence\s+of\s+circulating\s+ presence\s+of\s+circulating\s+)?(IgG\s+)?((anti\s*-\?\s*CCP\s+antibod(ies y) anti\s*-\?\s*CCP CCP\s+antibod(ies y) (anti\s*-\?\s*)?cyclic\s+citrullinated\s+peptide(\s+antibod(ies y)) a(\s*-\s*)?CCP anti\s*-\?\s*citrullinated\s+protein\s*? \s*peptide\s+antibod(ies y) anti\s*-\?\s*citrullinated\s+(protein peptide)\s+antibod(ies y) ACPA ACP\s+antibod(ies y) CCP)\s+and\s+)?(\bRF rheumatoid\s+factor(s)?)
nlpRF	RFStandalone	(?i)(?m)(?s)(\bRF\b \b\rheumatoid\s+factor\b)
nlpccp	CCPStandalone	(?i)(?m)(?s)\b(anti\s*-\?\s*CCP\s+antibod(ies y) anti\s*-\?\s*CCP CCP\s+antibod(ies y) (anti\s*-\?\s*)?cyclic\s+citrullinated\s+peptide(\s+antibod(ies y)) a(\s*-\?\s*)?CCP anti\s*-\?\s*citrullinated\s+protein\s*? \s*peptide\s+antibod(ies y) anti\s*-\?\s*citrullinated\s+(protein peptide)\s+antibod(ies y) ACPA ACP\s+antibod(ies y) CCP)\b
nlpccp	CCPPositive	(?i)(?m)(?s)(borderline\s+ slightly\s+ strongly\s+ low\s+ weakly\s+ low\s*-\?\s*titer\s+ high\s*-\?\s*titer\s+)?(\+\s* pos(itive itivity)\?\s+ elevated\s+ increased\s+)((test\s+)?for\s+)?(IgG\s+)?((\bRF rheumatoid\s+factor(s)?)\s+and\s+)?(anti\s*-\?\s*CCP\s+antibod(ies y) anti\s*-\?\s*CCP CCP\s+antibod(ies y) (anti\s*-\?\s*)?cyclic\s+citrullinated\s+peptide(\s+antibod(ies y)) a(\s*-\s*)?CCP anti\s*-\?\s*citrullinated\s+protein\s*? \s*peptide\s+antibod(ies y) anti\s*-\?\s*citrullinated\s+(protein peptide)\s+antibod(ies y) ACPA ACP\s+antibod(ies y) CCP)
nlpccp	CCPPositive	(?i)(?m)(?s)(anti\s*-\?\s*CCP\s+antibod(ies y) anti\s*-\?\s*CCP CCP\s+antibod(ies y) (anti\s*-\?\s*)?cyclic\s+citrullinated\s+peptide(\s+antibod(ies y)) a(\s*-\?\s*)?CCP anti\s*-\?\s*citrullinated\s+protein\s*? \s*peptide\s+antibod(ies y) anti\s*-\?\s*citrullinated\s+(protein peptide)\s+antibod(ies y) ACPA ACP\s+antibod(ies y) CCP)\s+and\s+ (\bRF rheumatoid\s+factor(s)?)\s+)?(is\s+ was\s+ are\s+ has\s+ also\s+)?been\s+)?(also\s+)?(found\s+)?(to\s+be\s+)?(also\s+)?(borderline\s+ slightly\s+ strongly\s+ low\s+ weakly\s+ low\s*-\?\s*titer\s+ high\s*-\?\s*titer\s+)?(present\b +)\s+positive\s+positivity\s+pos\b elevated\s+increased\s+greater\s+than\s+\d{2,})
nlpccp	CCPPositive	(?i)(?m)(?s)(anti\s*-\?\s*CCP\s+antibod(ies y) anti\s*-\?\s*CCP CCP\s+antibod(ies y) (anti\s*-\?\s*)?cyclic\s+citrullinated\s+peptide(\s+antibod(ies y)) a(\s*-\?\s*)?CCP anti\s*-\?\s*citrullinated\s+protein\s*? \s*peptide\s+antibod(ies y) anti\s*-\?\s*citrullinated\s+(protein peptide)\s+antibod(ies y) ACPA ACP\s+antibod(ies y) CCP)\s+and\s+ (\bRF rheumatoid\s+factor(s)?)\s+)?\s*+
nlpccp	CCPPositive	(?i)(?m)(?s)(?<!(no\s{1,100}))\s+of\s+circulating\s+ presence\s+of\s+circulating\s+)?(IgG\s+)?((anti\s*-\?\s*CCP\s+antibod(ies y) anti\s*-\?\s*CCP CCP\s+antibod(ies y) (anti\s*-\?\s*)?cyclic\s+citrullinated\s+peptide(\s+antibod(ies y)) a(\s*-\s*)?CCP anti\s*-\?\s*citrullinated\s+protein\s*? \s*peptide\s+antibod(ies y) anti\s*-\?\s*citrullinated\s+(protein peptide)\s+antibod(ies y) ACPA ACP\s+antibod(ies y) CCP)

Limitations

- Still only one chronic disease well represented in ICD-9s
- May not represent portability of other machine learning methods