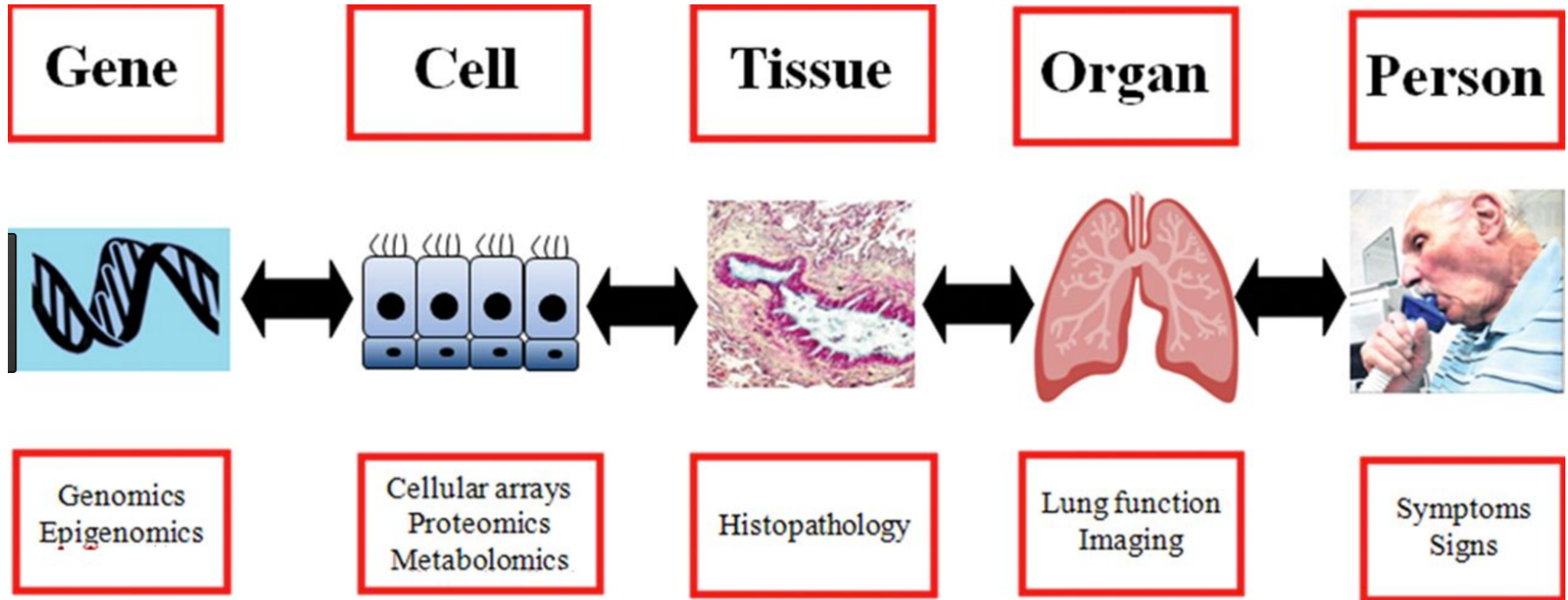


EHR-Based Phenotyping: Bulk Learning and Evaluation (with Infectious Diseases)

Po-Hsiang (Barnett) Chiu

Phenotypes and phenotyping



Physically observable traits of genotypes (and their interactions with environments)

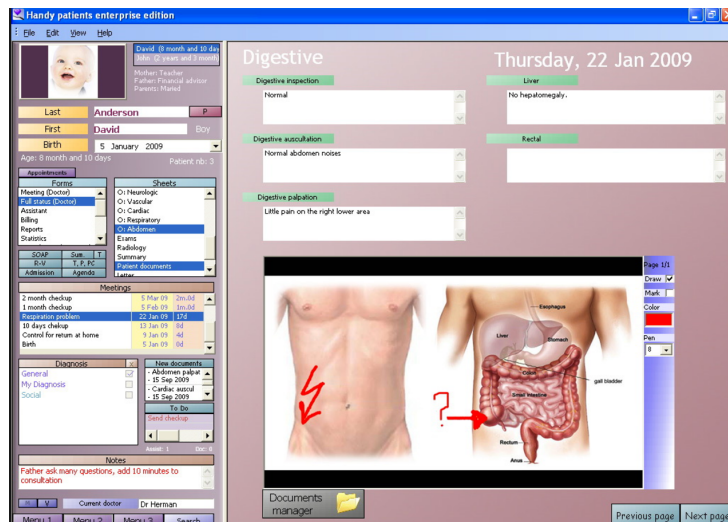
Biochemical or physiological properties, behavior, and products of behavior

Attributions of diseases (e.g. susceptibility)

Diseases (and disease subtypes)

Data-Driven Phenotyping

- Data-driven phenotyping
 - Two main methodologies
 - Rule-based approach (e.g. eMerge, <https://emerge.mc.vanderbilt.edu>)
 - **Predictive Analytics**
 - Data sources:
 - EHRs/EMRs: Medicinal treatments, diagnoses, lab measurements, etc.
 - Genomic data: SNP arrays, copy number variation (CNVs), etc.
 - Phenotypes
 - Diseases, subtypes, or variables attributed to disease predictions



Diagnostic Concept Units

- Various diseases sharing the same set of diagnostic concept units
- Infectious diseases
 - Lab tests
 - Microorganism, blood, urine, body tissues, stool
 - Medications
 - Antibiotic, antiviral, anthelmintic
- Build statistical models for each diagnostic component and combine them appropriately
 - Ensemble learning

Bulk Learning in a Nutshell ...

Bulk Learning is a **batch-phenotyping framework** that uses multiple diseases collectively (i.e. **bulk learning set**) as a substrate for model learning and evaluation wherein (a given) **medical ontology** is used to perform **feature selection** and **model stacking** is used to construct **abstract feature representation** of low sample complexity in order to **reduce training requirements**.

Key Concepts:

1. Build phenotyping models on top of multiple diseases
2. Automatic feature selection using an existing ontology
3. Models are combined via model stacking (a form of ensemble learning)
4. Abstract features

Dimensionality reduction

5. Less labeled data required for model evaluations

Phenotyping via Bulk Learning

- Under model stacking, we then arrive at the notion of “concept-driven phenotyping”
 - A subset or combinations of lab tests are more attributable to some diseases while the others are better explained by medications
- In this study, infectious diseases associated with 100 ICD-9 codes as the domain of study for **bulk learning**
 - For simplicity, consider different diagnostic codes as different diseases ...
 - Why 100 codes?
 - Code selection strategy?

Bulk Learning Basics I

- Addresses two central issues in **predictive analytical** approach to computational phenotyping
 - Feature engineering
 - **Medical ontology** for feature decomposition
 - Medical Entities Dict (<http://med.dmi.columbia.edu>)
 - Data annotation
 - **Ensemble learning** (e.g. **stacked generalization** [Wolpert 1992])
 - **Feature abstraction** for dimensionality reduction

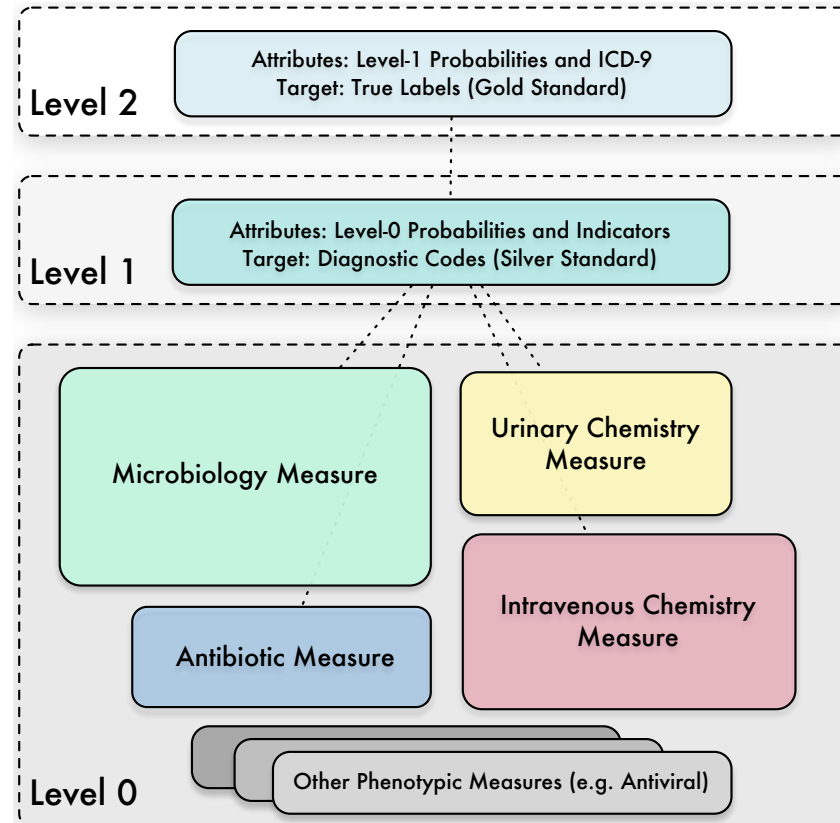
Medical Ontology for Grouping Features

- Snapshot of **Medical Entities Dictionary**
(<http://med.dmi.columbia.edu>)

Hierarchy		Slots
1 Parent		2235 - Microbiology Procedure
P a r e n t s	32458 - Organism Panels [2]	1-UMLS-CODE: C0085672 5-SYNONYMS: 7-HAS-PARTS * 8-PART-OF * 11-DEFINITION: 14-ASSESSES-SAMPLE * 16-ENTITY-MEASURED * 17-UNITS: 23-TEST-->RESULT-TYPE * 1067 - Smear Result 38-CPMC-NORMAL-VALUE: 39-CPMC-LOW-NORMAL-VALUE: 40-CPMC-HIGH-NORMAL-VALUE: 41-CPMC-MALE-LOW-NORMAL-VALUE: 42-CPMC-MALE-HIGH-NORMAL-VALUE: 43-CPMC-FEMALE-LOW-NORMAL-VALUE: 44-CPMC-FEMALE-HIGH-NORMAL-VALUE: 45-CPMC-NORMAL-RANGES-TEXT: 50-MAIN-MESH: 51-SUPPLEMENTARY-MESH: 95-ACTIVE-SYSTEM-ITEM-(MAPS-TO)->LEGACY-ITEM * 126-CPT4-CODE: 138-IS-DISPLAY-PARAMETER-OF * 139-HAS-TEST-DISPLAY-CLASS-NAME: 148-HAS-PROC-DISPLAY-CLASS-NAME:
2235 - Microbiology Procedure		
C h i l d r e n	32411 - Microbiology Blood Procedure [9] 33896 - Gonococcus Detection Procedures [2] 42238 - Microbiology Non-Sensitivity Procedures [72] 42247 - Microbiology Culture and Sensitivity Procedure [6] 49925 - New York Hospital (NYH) Microbiology Tests [3] 75025 - Microbiology Urine Procedure [28] 125810 - Millennium Microbiology Test [2] 157988 - Post Mortem Culture Procedure [4]	
8 Children		
Select new Medcode: <input type="text"/> * Submit Clear * Search the MED: <input type="text"/> * Submit Clear On Slot: All <input type="button" value="v"/>		

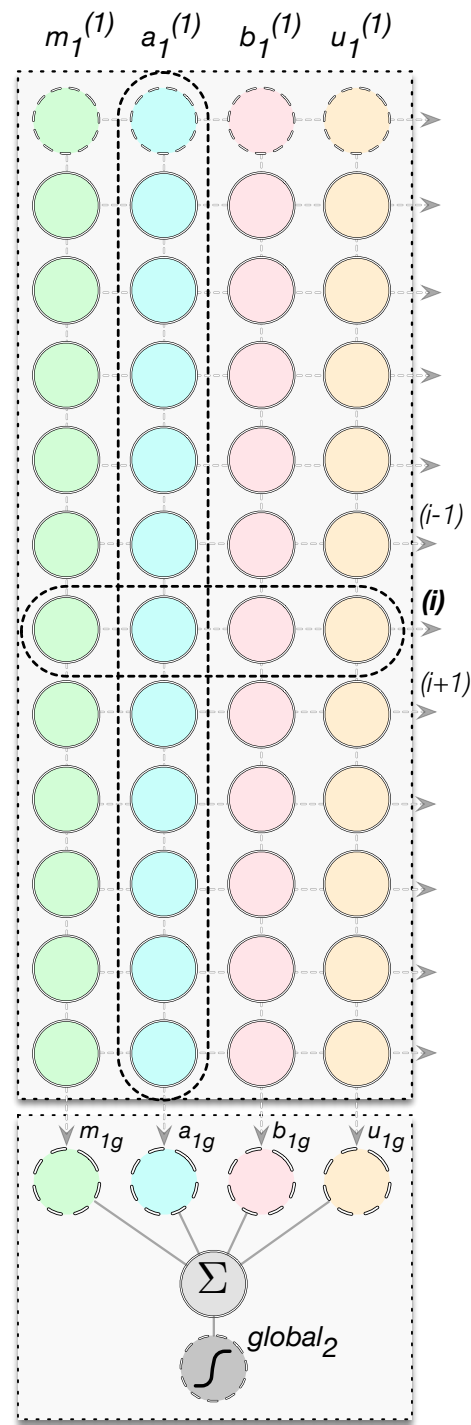
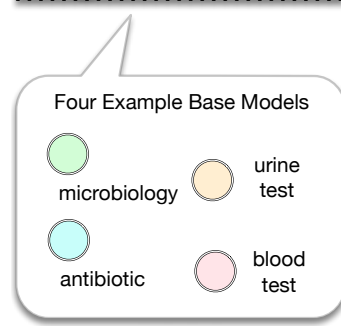
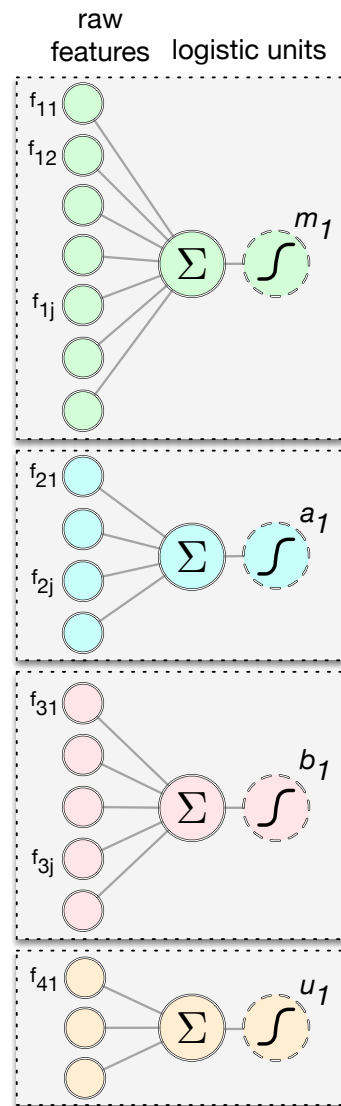
Model Stacking

- Why inspecting multiple (infectious) diseases?
 - Using **multiple diseases** as substrate and identify their common elements
 - Example stacking architecture (under stacked generalization method)



Surrogate Labels vs True Labels

- Model stacking is used to achieve:
 - Improve upon base model performances
 - Transform EHR data to a denser form
- Uses diagnostic codes (e.g. ICD-9) as surrogate labels to establish “approximate predictive models.”
- Why surrogate labels (e.g. ICD-9)?
 - Features extracted from EHR can be large
 - Used to derive compact representation of the training data
 - “Free” supervised signals that are sufficiently close but can be obtained without extra work
- Objective: Build statistical models in abstract feature space
 - Create a sparse annotation set (i.e. gold standard) that serves a proxy dataset for downstream model evaluations
 - 83 annotated cases

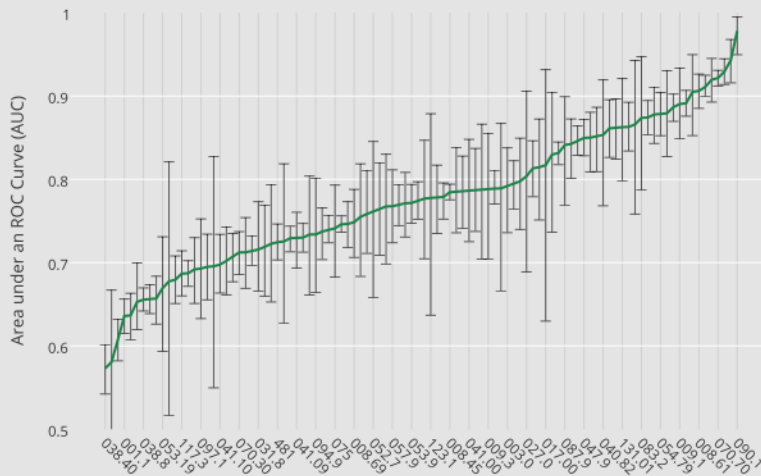


Performance Evaluations

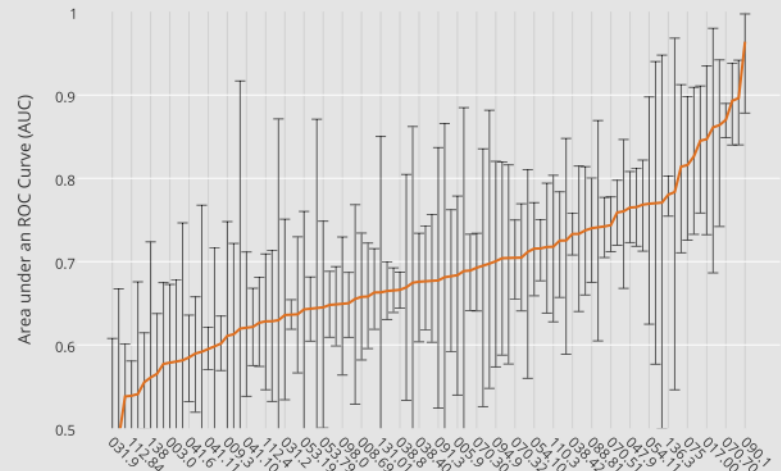
- How well does the model predict ICD-9s (using a separate test data)?
- How well does the model predict annotated data (assoc. with “true labels”)?
 - (Binarized) ICD-9 becomes a candidate feature among abstract features (e.g. probability scores, indicators)
 - Annotated sample consists of randomly selected cases in which errors of ICD-9 coding are corrected
 - Data annotations and coding procedures are two independent processes

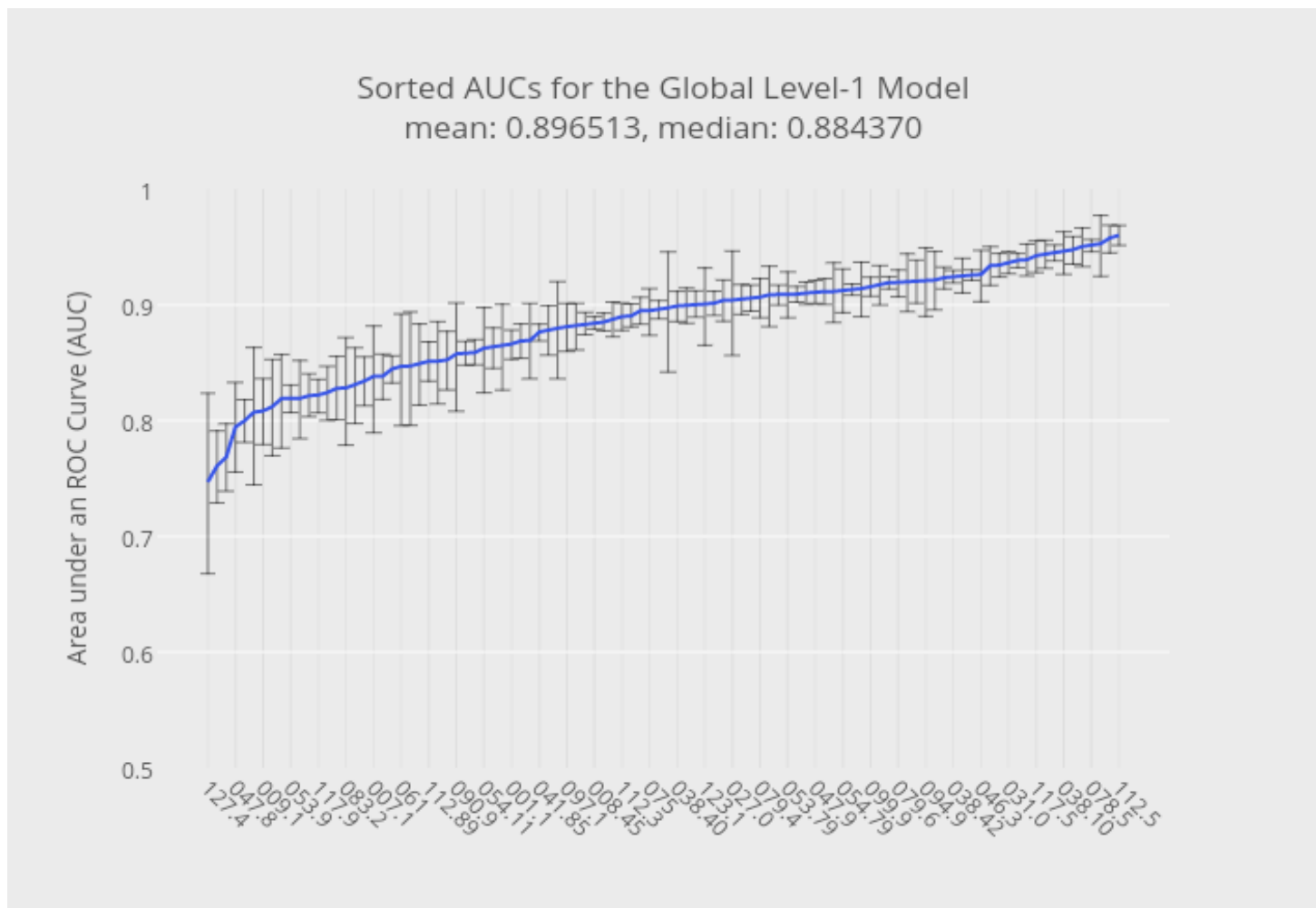
Base Level Performances

Sorted Performance of the Microbiology Model
grand mean: 0.775399, median: 0.776603



Sorted Performance of the Urine Test Model
grand mean: 0.676575, median: 0.685738





127.4 Enterobiasis

009.1 Gastroenteritis ...

117.9 Mycoses

047.8 (Other) viral meningitis

053.9 Herpez zoster

Settings	Sensitivity	Specificity	Mean AUC (Repeated 10-fold with 30 cycles)
Level 1 (L1)	1029/1170 (0.88)	212/1320 (0.16)	0.59 (0.51 ~ 0.66)
Level 2 (L2)	812/1170 (0.69)	456/1320 (0.35)	0.52 (0.45 ~ 0.60)
L1 + ICD9	1158/1170 (0.99)	771/1320 (0.58)	0.85 (0.80 ~ 0.89)
L2 + ICD9	910/1170 (0.78)	836/1320 (0.63)	0.74 (0.67 ~ 0.82)
Big Logistic	768/1170 (0.66)	866/1320 (0.66)	0.65 (0.59 ~ 0.72)
Big SVM	784/1170 (0.67)	862/1320 (0.65)	0.53 (0.51 ~ 0.56)

Table 7b. Comparison by annotation types among different meta-classifiers trained by mixing virtual annotations.

Settings	Type TP (39)	Type FP (15)	Type TN (29)	Type FN (0)
Level 1 (L1)	1029/1170 (0.88)	102/450 (0.23)	110/870 (0.13)	n/a
Level 2 (L2)	812/1170 (0.69)	158/450 (0.35)	298/870 (0.34)	n/a
L1 + ICD9	1158/1170 (0.99)	10/450 (0.02)	761/870 (0.87)	n/a
L2 + ICD9	910/1170 (0.78)	104/450 (0.23)	732/870 (0.84)	n/a
Big Logistic	768/1170 (0.66)	276/450 (0.61)	590/870 (0.68)	n/a
Big SVM	784/1170 (0.67)	291/450 (0.65)	571/870 (0.66)	n/a

Other Components

- Semi-supervised learning and virtual annotation set
- The 3rd tier in model stacking hierarchy
 - Trade-off between learned abstract features and the ICD-9 codes as surrogate labels.
 - Performance evaluation on predicting annotated labels
- Ontology-based feature engineering
- Proper design of treatment and control (training) data

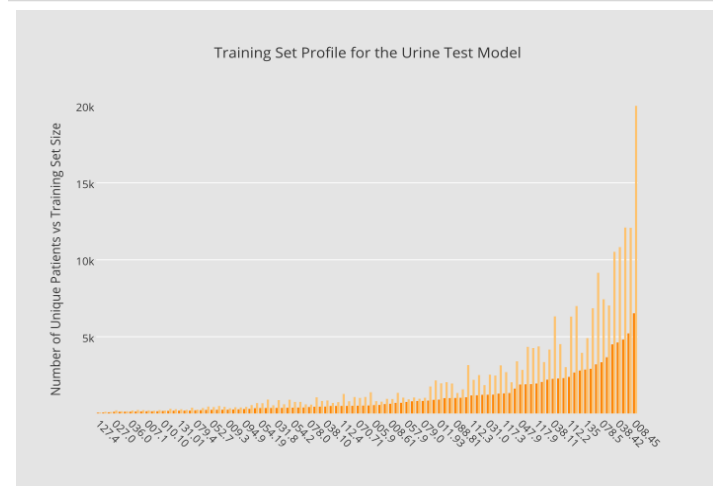
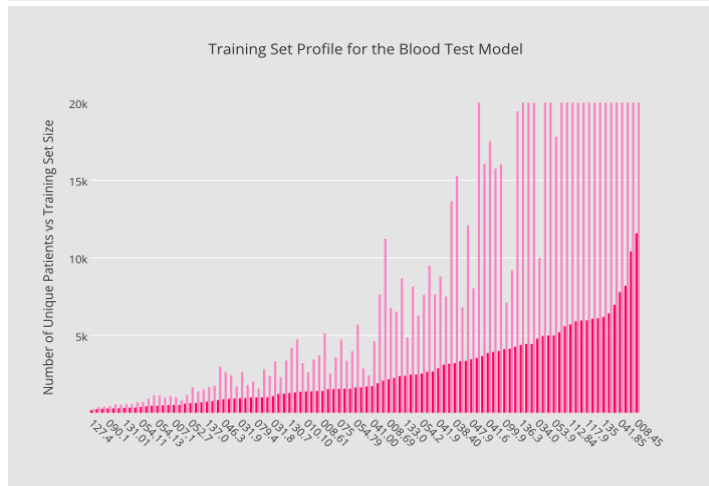
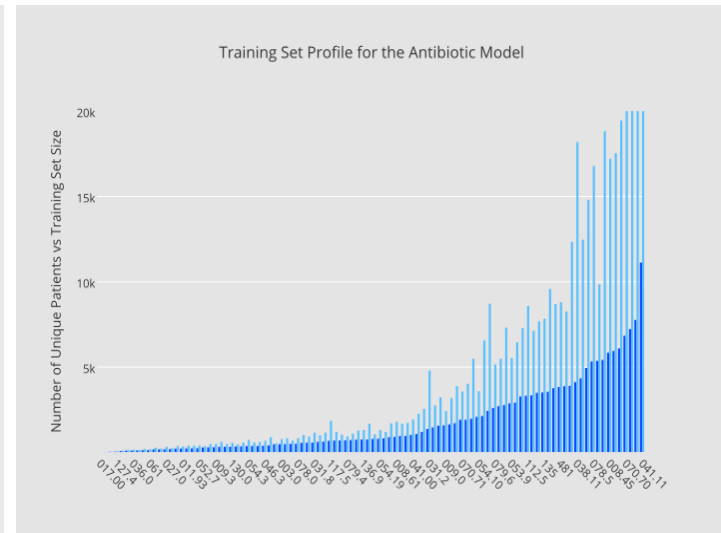
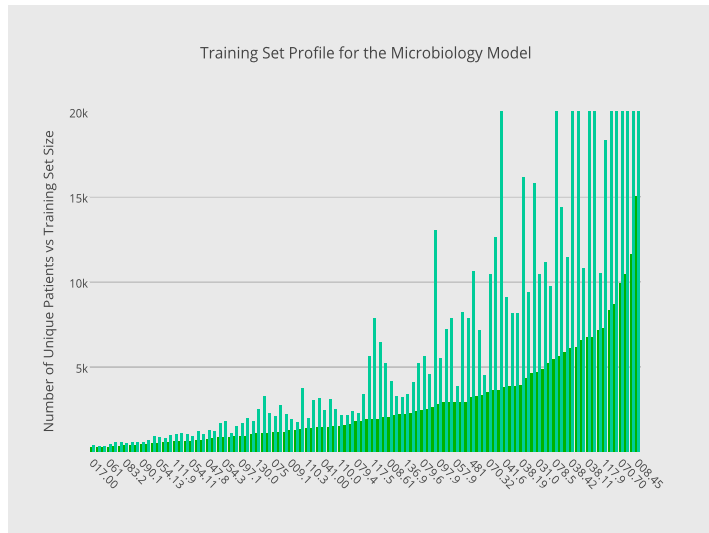
Modeling Perspective

- EHR data consist of **observations** and **latent variables**
 - Observations can be directly answered via simple queries
 - Did the patient have tests on E. Coli?
 - Did the patient take Ceftriaxon?
- Latent variables represent quantities that cannot be directly observed in EHR or computed via simple queries
 - Does the patient have an infection?
 - Diagnostic questions: specifically which infections do the patient have?
- Learn classifiers to predict latent variables (with only access to observations)

Medical Perspective

- Seemingly different infectious diseases may share similar sets of lab tests and medications
 - Staph. aureus
 - Skin infections, pneumonia, blood poisoning
 - Ceftriaxone
 - Meningitis
 - Infections at different sites of the body (e.g. bloodstream, lungs, urinary tracts)
- Multiple classifiers for the same disease
 - 4 classifiers per ICD-9 code, each of which is binary classifier
 - 400 classifiers at base level

Data Distribution Perspective



“Can we build a joint model applicable to all diseases?”

Abstract Feature Representation: Design Choices

- Related work in constructing high-level features
 - PCA, unsupervised feature learning, manifold learning, etc.
- Design choices
 - Data characteristics
 - Interpretability
- Deep Neural Network
 - Linear combination
 - Non-linear transformation (e.g. sigmoid, rectifier, etc.)
- Feature set: continuous, dense, and “homogeneous”
 - Image pixels
 - Times series of lab measurements
 - word2vec
- EHR data however are very different
 - sparse and incomplete
 - consist of many different types (binary, categorical, continuous, etc.)
 - Features associated with multiple concepts

Moving Forward ...

- Summary
 - Bulk learning is a framework with at least the following system choices
 - The bulk learning set (of target conditions) => base models
 - Classification algorithms (guideline: probabilistic classifiers + well-calibrated)
 - Stacking architecture (multiple tiers => levels of abstractions)
 - Strategy for combining individual (local) disease models to a global model
 - Advantage: Can use a **small annotated sample** for model construction and evaluation within the abstract feature space (e.g. level-1 data)
 - 83 clinical cases were labeled in this study
 - Challenge: The model involving the interaction between abstract features and ICD-9 do not generalize well into the region of the data where the ICD-9 coding was incorrect
 - Multiple types of surrogate labels

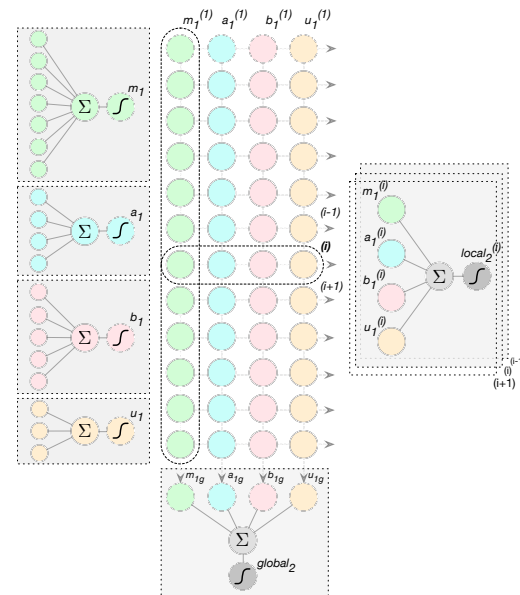
- Ongoing and future work

Complex decision boundary?

Other surrogate labels

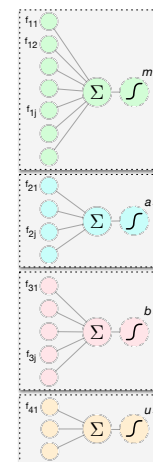
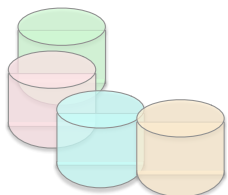
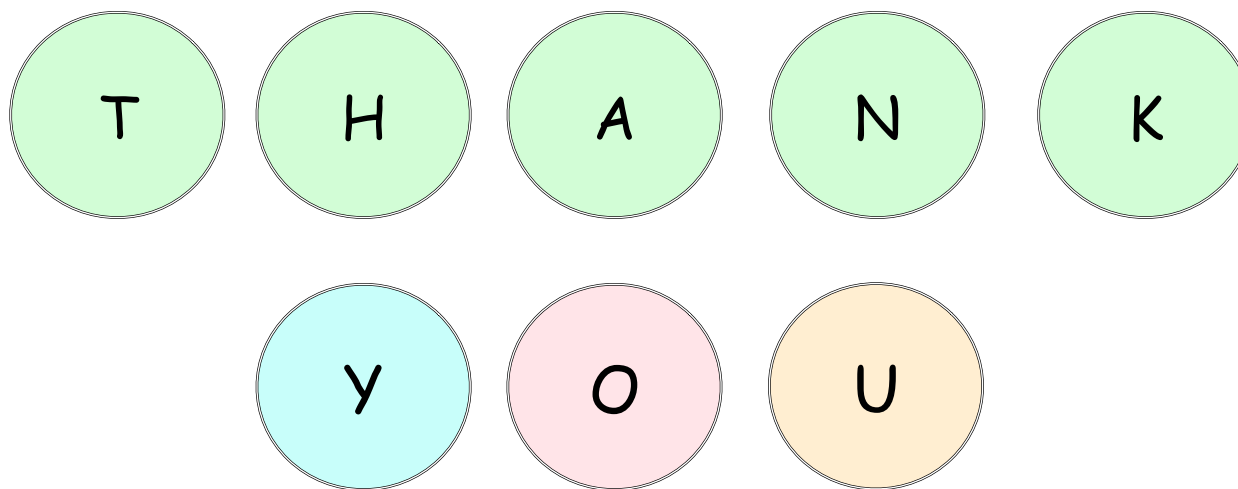
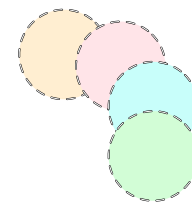
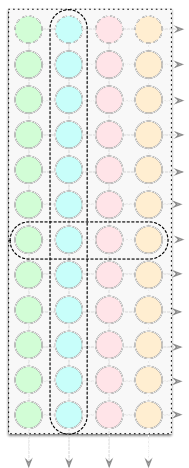
Semi-supervised learning

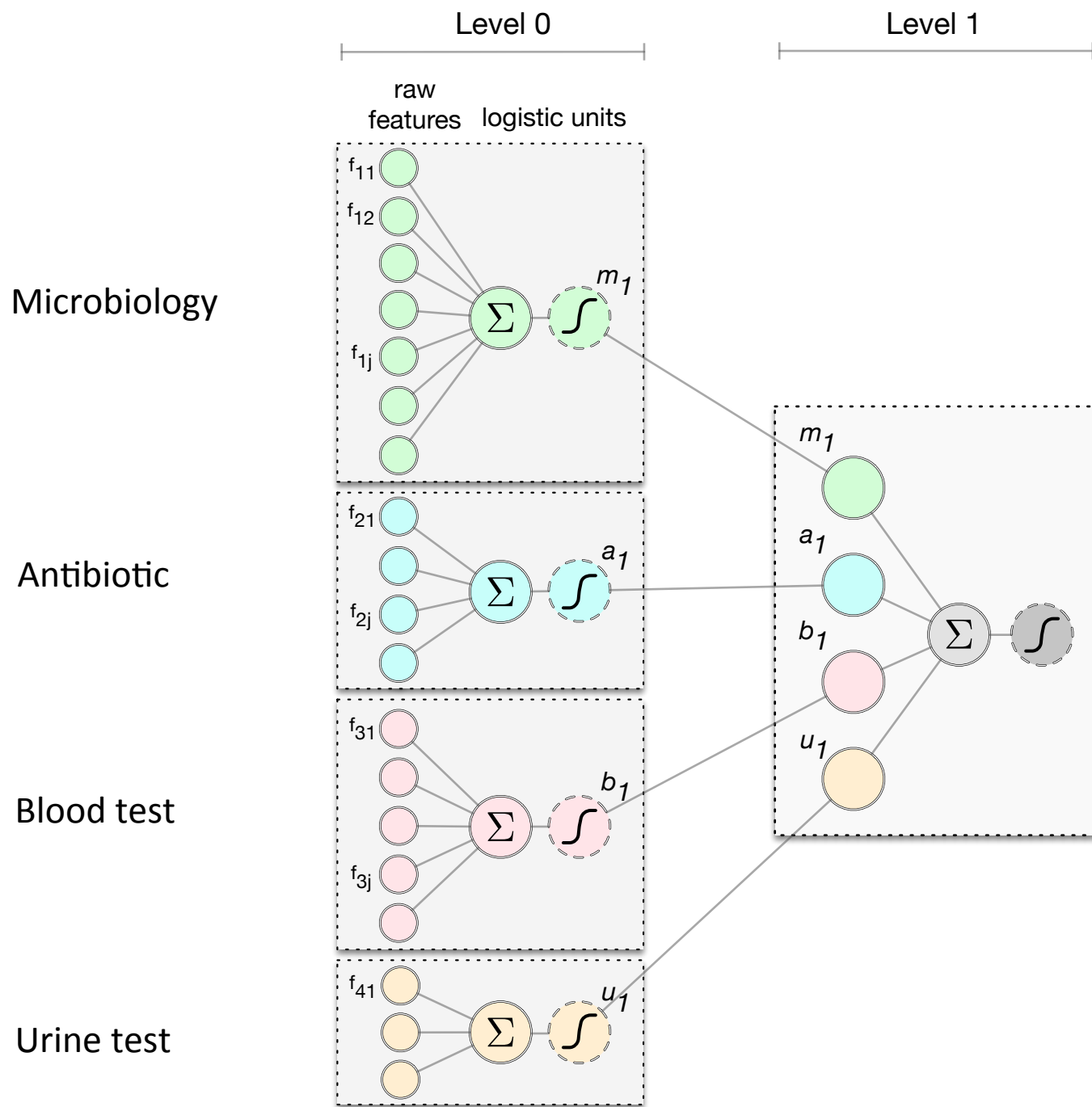
Active learning



Reference

- [1] D.H. Wolpert, Stacked generalization, Neural Networks. 5 (1992) 241–259.
- [2] K.M. Ting, I.H. Witten, Issues in stacked generalization, J. Artif. Intell. Res. 10 (1999) 271–289.
- [3] J. Jin Chen, C. Cheng Wang, R. Runsheng Wang, Using Stacked Generalization to Combine SVMs in Magnitude and Shape Feature Spaces for Classification of Hyperspectral Data, IEEE Trans. Geosci. Remote Sens. 47 (2009) 2193-2205.
- [4] David Baorto, James Cimino, et al.
Available: <http://med.dmi.columbia.edu>. Access date: Oct 20, 2016.
- [5] T.A. Lasko, J.C. Denny, M.A. Levy, Computational Phenotype Discovery Using Unsupervised Feature Learning over Noisy, Sparse, and Irregular Clinical Data, PLoS One. 8 (2013) e66341.





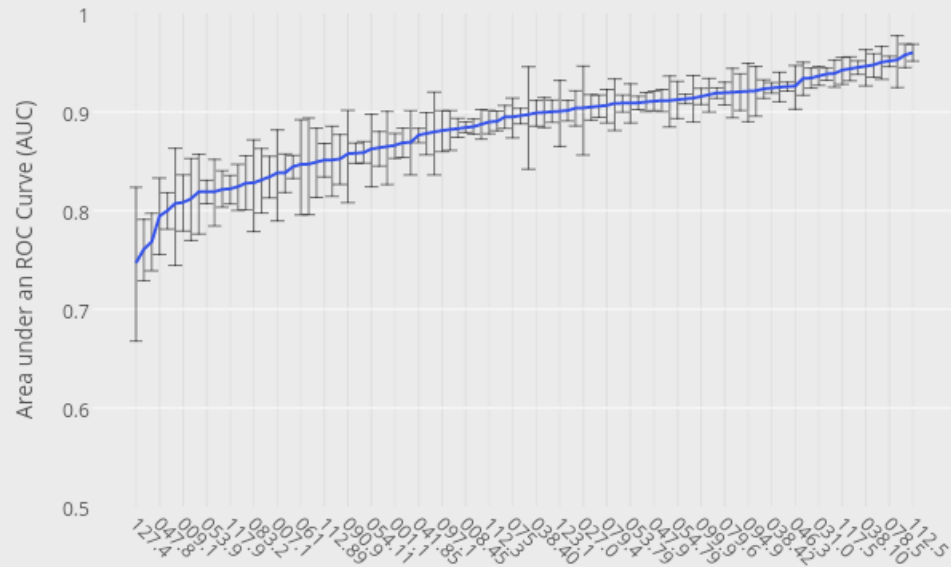
Example Features

Microorganism Lab Test (Microbiology)		Antibiotic Prescription (Antibiotic)	
MedCode	Description	MedCode	Description
935	Organism Result: Escherichia Coli	72900	Piperacillin/Tazobactam
799	Organism Result: Candida Albicans	72702	Vancomycin
774	Organism Result: Staphylococcus Aureus	100198	Ceftriaxone
910	Organism Result: Klebsiella Pneumoniae	66042	Levofloxacin
31826	Organism Result: Enterococcus Faecalis	61003	Tobramycin
59993	Negative for Clostridium Difficile Toxin A and Toxin B	60671	Azithromycin
39576	Rule Out Influenza Virus	62375	Meropenem
316	No Ova or Parasites Found	61461	Amoxicillin
994	Positive for Gram Negative Rods	60918	Dapsone
36453	Susceptibility Type: Microscan Mic	62879	Cephalexin

Intravenous Chemistry Test (Blood)		Urinary Chemistry Test (Urine)	
MedCode	Description	MedCode	Description
69494	Lab Test: Vitamin B12	36265	Lab Test: Ketone
35995	Lab Test: Lactate, Arterial	36267	Lab Test: Potassium, Random Urine
39564	Lab Test: Cyclosporine, Whole Blood	36260	Lab Test: Urine Glucose
65906	Lab Test: Hemoglobin A1c	36269	Lab Test: Urine Leukocyte Esterase
36300	Lab Test: Vancomycin	36286	Lab Test: Urine Protein
59415	Lab Test: Tacrolimus	1390	Urine Blood Test
46418	Blood Bank: ABO Antigen Determination	1395	Urine pH Measurement
46421	Blood Bank: Antierythrocyte Antibody Screen	1388	Urine Urobilinogen Test
59942	Lab Test: Glucose Wholeblood	1394	Urine Albumin Test
59047	Lab Test: Creatine Kinase	1392	Urine Acetone Test

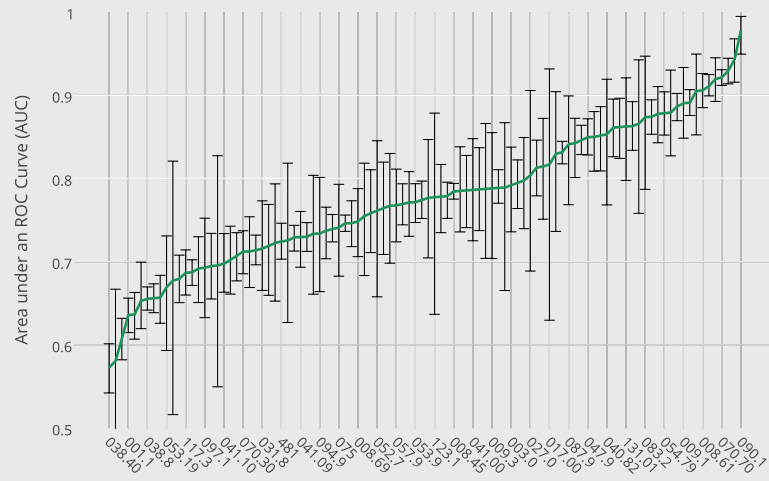
Sorted AUCs for the Global Level-1 Model

mean: 0.896513, median: 0.884370



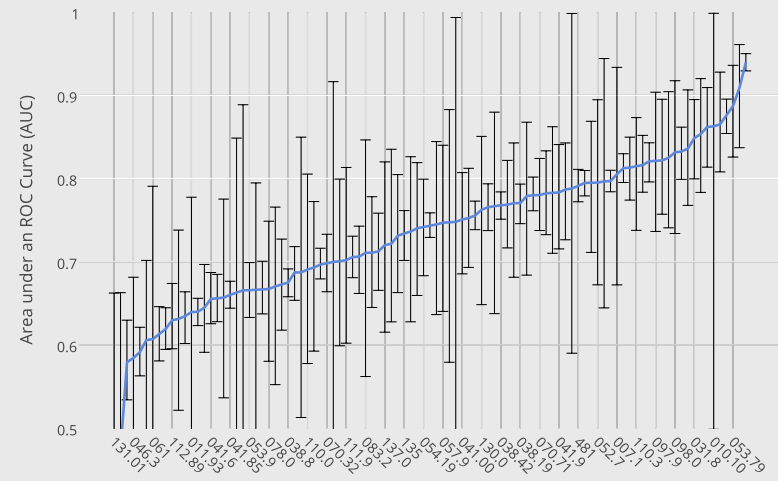
Sorted Performance of the Microbiology Model

grand mean: 0.775399, median: 0.776603



Sorted Performance of the Antibiotic Model

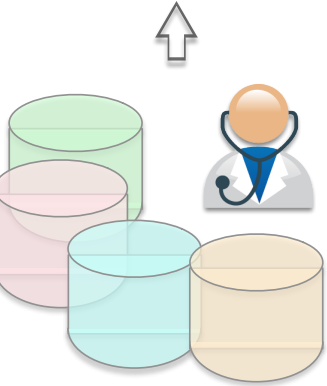
grand mean: 0.743400, median: 0.733976



1. Define Feature Groups Using Medical Ontology

1b. Use Medical Entities Dictionary to delineate feature scopes

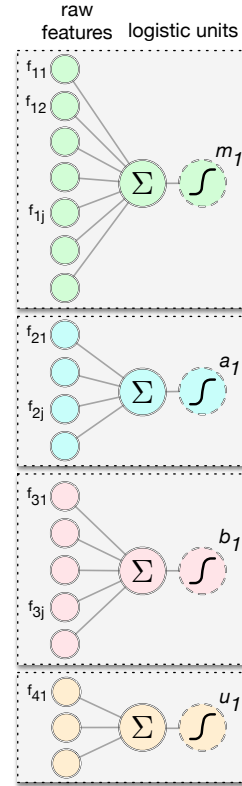
Hierarchy	Snits
1 Point	2225 - Microbiology Procedure
12458 - Organisms (2)	1: SNITS CODE CODES (2)
	2: SNITS CODE
	3: SNITS CODE
	4: SNITS CODE
	5: SNITS CODE
	6: SNITS CODE
	7: SNITS CODE
	8: SNITS CODE
	9: SNITS CODE
	10: SNITS CODE
	11: SNITS CODE
	12: SNITS CODE
	13: SNITS CODE
	14: SNITS CODE
	15: SNITS CODE
	16: SNITS CODE
	17: SNITS CODE
	18: SNITS CODE
	19: SNITS CODE
	20: SNITS CODE
	21: SNITS CODE
	22: SNITS CODE
	23: SNITS CODE
	24: SNITS CODE
	25: SNITS CODE
	26: SNITS CODE
	27: SNITS CODE
	28: SNITS CODE
	29: SNITS CODE
	30: SNITS CODE
	31: SNITS CODE
	32: SNITS CODE
	33: SNITS CODE
	34: SNITS CODE
	35: SNITS CODE
	36: SNITS CODE
	37: SNITS CODE
	38: SNITS CODE
	39: SNITS CODE
	40: SNITS CODE
	41: SNITS CODE
	42: SNITS CODE
	43: SNITS CODE
	44: SNITS CODE
	45: SNITS CODE
	46: SNITS CODE
	47: SNITS CODE
	48: SNITS CODE
	49: SNITS CODE
	50: SNITS CODE
	51: SNITS CODE
	52: SNITS CODE
	53: SNITS CODE
	54: SNITS CODE
	55: SNITS CODE
	56: SNITS CODE
	57: SNITS CODE
	58: SNITS CODE
	59: SNITS CODE
	60: SNITS CODE
	61: SNITS CODE
	62: SNITS CODE
	63: SNITS CODE
	64: SNITS CODE
	65: SNITS CODE
	66: SNITS CODE
	67: SNITS CODE
	68: SNITS CODE
	69: SNITS CODE
	70: SNITS CODE
	71: SNITS CODE
	72: SNITS CODE
	73: SNITS CODE
	74: SNITS CODE
	75: SNITS CODE
	76: SNITS CODE
	77: SNITS CODE
	78: SNITS CODE
	79: SNITS CODE
	80: SNITS CODE
	81: SNITS CODE
	82: SNITS CODE
	83: SNITS CODE
	84: SNITS CODE
	85: SNITS CODE
	86: SNITS CODE
	87: SNITS CODE
	88: SNITS CODE
	89: SNITS CODE
	90: SNITS CODE
	91: SNITS CODE
	92: SNITS CODE
	93: SNITS CODE
	94: SNITS CODE
	95: SNITS CODE
	96: SNITS CODE
	97: SNITS CODE
	98: SNITS CODE
	99: SNITS CODE
	100: SNITS CODE



1a. Gather EHR data according to medical concepts

1c. Apply feature selection within each concept group

2. Compute Base Models



Four Example Base Models

- microbiology
- antibiotic
- urine test
- blood test

3. Compute Meta Models (via Ensemble Learning)

