

# High-Performing Machine Learning Models for Phenotype Development

Victor A. Rodriguez\*, MPhil, Tony Y. Sun\*, BS, Phyllis M. Thangaraj\*, MPhil, Karthik Natarajan, PhD, Patrick Ryan, PhD (\* Contributed equally)

Department of Biomedical Informatics, Columbia University

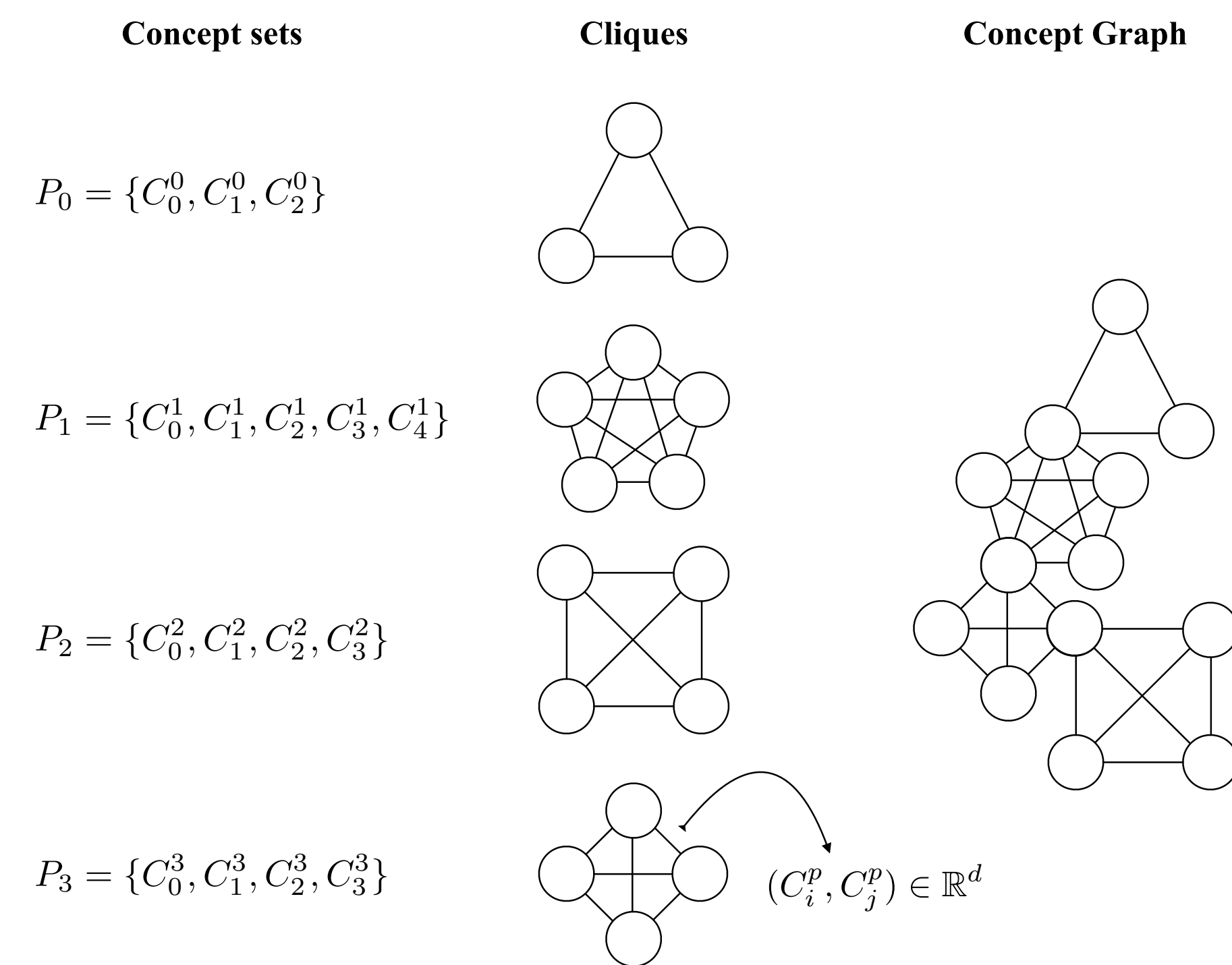


## 1 Introduction

- Phenotyping algorithms are essential tools for conducting clinical research on observational data, as high-throughput phenotype development remains an open challenge.
- eMERGE (electronic Medical Records and Genomics) phenotypes** are manually defined by the eMERGE Network of clinicians and informaticians in terms of multiple clinical concepts, or **concept sets**.
- We propose a framework for learning from the structure of eMERGE phenotype concept sets to aid construction of novel concept set phenotype definitions.

### 1.1 Concept Sets → Concept Graph

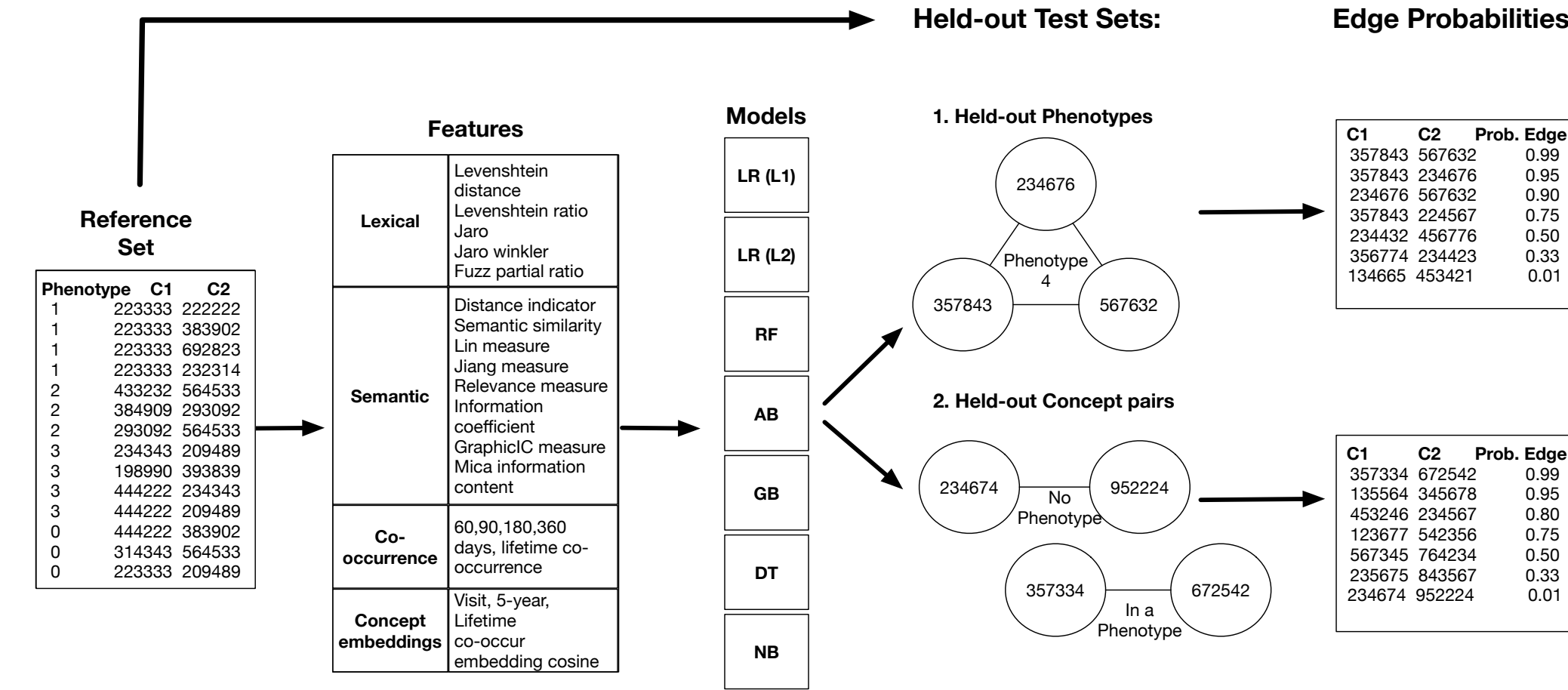
- We study the structure of eMERGE phenotype concept sets by considering **all possible concept-concept pairs within a concept set**.
- If we treat **concepts as nodes**, and **concept-concept pairs as edges**, then a **concept set is a fully connected subgraph, or clique**, within a larger **concept graph**.
- From this perspective, constructing a phenotype is equivalent to constructing a concept clique.
- We learn how to build concept set cliques by first learning how to predict concept-concept pair edges using rich edge features.



**Figure 1:** eMERGE phenotypes form cliques within a concept graph. We study concept set construction by learning how to predict concept-concept pair edges using rich edge features. Here  $P_i$  indicates the  $i^{th}$  eMERGE phenotype.  $C_j^i$  indicates the  $j^{th}$  concept in the  $i^{th}$  concept set.  $(C_i^p, C_j^p)$  is a concept-concept pair represented by features in  $\mathbb{R}^d$ .

## 2 Methods

- We train several models to perform edge prediction within our concept graph.
- Our dataset is comprised of concept-concept pairs (edges), each described by a rich feature vector (See Figures 1 and 2)
- All edges appearing within at least one concept set are treated as **positive instances**; all other edges are considered **possible negative instances**.



**Figure 2:** Methods visual summary. C1 and C2 represent concept 1 and concept 2 pairs, Model abbreviations: LR (L1)- Logistic Regression with L1 penalty, LR (L2)- Logistic Regression with L2 penalty, RF- Random Forest, AB- Adaboost, GB- Gradient Boosting, DT- Decision Tree, and NB- Naive Bayes.

### 2.1 Concept-Concept Pair Features

- A set of rich features was developed for each possible concept-concept pair.
- These include **lexical**, **semantic**, **co-occurrence**, and **concept embedding** features.

### 2.2 Training & Test Set Splits

- We generate training and test partitions in two ways: **random** & **phenotype aware**.
  - Random:**
    - Randomly select 90% positive edges for training; remaining 10% used for testing.
    - Sample negative edges to equal number of positive edges in training and test sets.
  - Phenotype Aware:**
    - Select concept sets comprising  $\sim 10\%$  of total positive edges; use for testing; remaining positive edges used for training.
    - Sample negative edges to equal number of positive edges in training and test sets.

### 2.3 Models

- We train several simple and ensemble binary classifiers for the purpose of edge prediction (See Figure 2 for a full list; all models implemented in `scikit-learn`).

## 3 Results

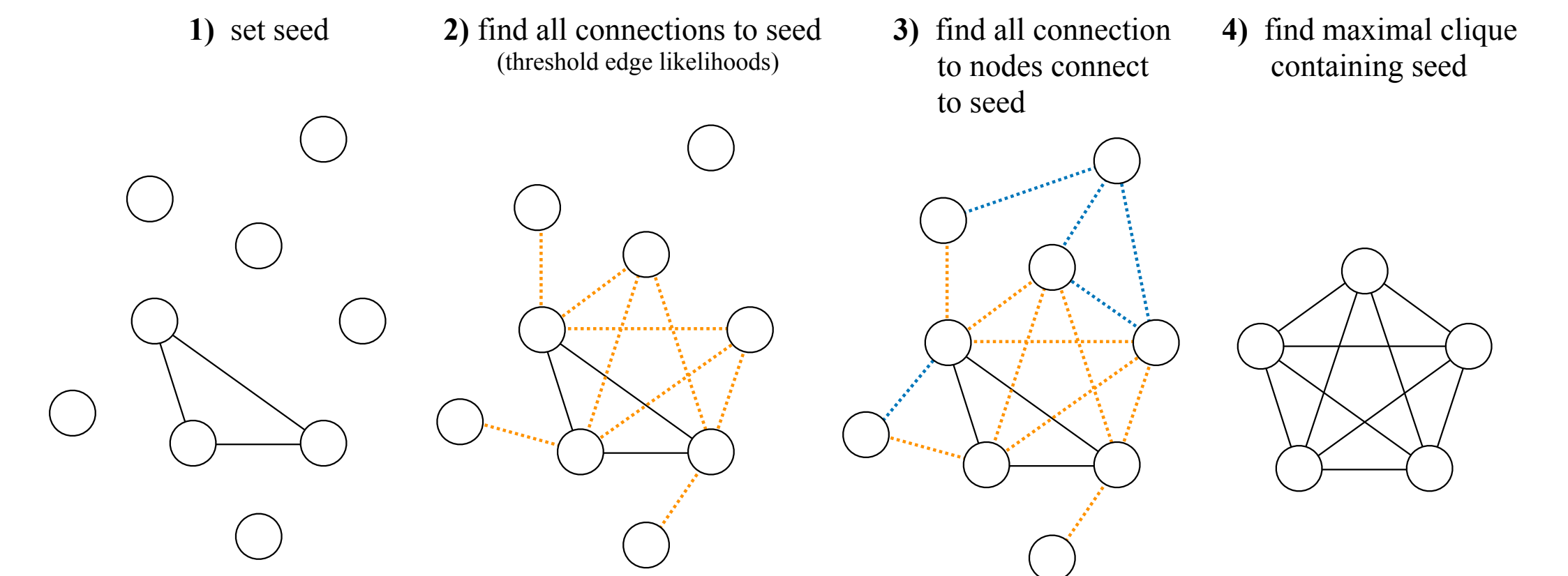
- Concept-concept pair prediction results for all models are show in Table 1
- For each model we evaluate the areas under the ROC & PR curves (AUROC, AUPRC), the maximum F1 (Max. F1) and the Precision at 50% of the test set (Prec.@50%).

Train/Test Split	Model	AUROC	AUPRC	Max. F1	Prec.@50%
Random	LR (L1)	0.9417 $\pm$ 0.0008	0.9337 $\pm$ 0.0014	0.8701 $\pm$ 0.0011	0.9700 $\pm$ 0.0205
	LR (L2)	0.9402 $\pm$ 0.0009	0.9321 $\pm$ 0.0015	0.8700 $\pm$ 0.0010	0.9560 $\pm$ 0.0377
	Naive Bayes	0.9318 $\pm$ 0.0009	0.9202 $\pm$ 0.0014	0.8436 $\pm$ 0.0024	0.9620 $\pm$ 0.0316
	Decision Tree	0.9448 $\pm$ 0.0009	0.9224 $\pm$ 0.0030	0.8838 $\pm$ 0.0008	0.9700 $\pm$ 0.0241
	Random Forest	0.9507 $\pm$ 0.0007	0.9424 $\pm$ 0.001	0.8848 $\pm$ 0.0018	0.9720 $\pm$ 0.0204
	Gradient Boosting	<b>0.9540 <math>\pm</math> 0.0008</b>	0.9430 $\pm$ 0.0017	<b>0.8909 <math>\pm</math> 0.0013</b>	0.9780 $\pm$ 0.0166
	AdaBoost	0.9516 $\pm$ 0.0007	<b>0.9431 <math>\pm</math> 0.0013</b>	0.8840 $\pm$ 0.0010	<b>0.9800 <math>\pm</math> 0.0179</b>
Phen. Aware	LR (L1)	0.8635 $\pm$ 0.0293	0.8694 $\pm$ 0.0254	0.8030 $\pm$ 0.0223	0.9700 $\pm$ 0.0257
	LR (L2)	0.8632 $\pm$ 0.0331	0.8771 $\pm$ 0.0262	0.8005 $\pm$ 0.0252	0.9580 $\pm$ 0.0166
	Naive Bayes	0.8691 $\pm$ 0.0290	0.8731 $\pm$ 0.0232	0.7659 $\pm$ 0.0509	0.9460 $\pm$ 0.0156
	Decision Tree	0.8853 $\pm$ 0.0261	0.8698 $\pm$ 0.0257	0.8179 $\pm$ 0.0133	0.9340 $\pm$ 0.0559
	Random Forest	<b>0.8953 <math>\pm</math> 0.0234</b>	0.8833 $\pm$ 0.0213	<b>0.8314 <math>\pm</math> 0.0101</b>	0.8340 $\pm$ 0.2057
	Gradient Boosting	0.8924 $\pm$ 0.0241	<b>0.8856 <math>\pm</math> 0.0232</b>	0.8204 $\pm$ 0.0116	0.9440 $\pm$ 0.0280
	AdaBoost	0.8704 $\pm$ 0.0273	0.8651 $\pm$ 0.0259	0.8024 $\pm$ 0.0189	<b>0.9780 <math>\pm</math> 0.0140</b>

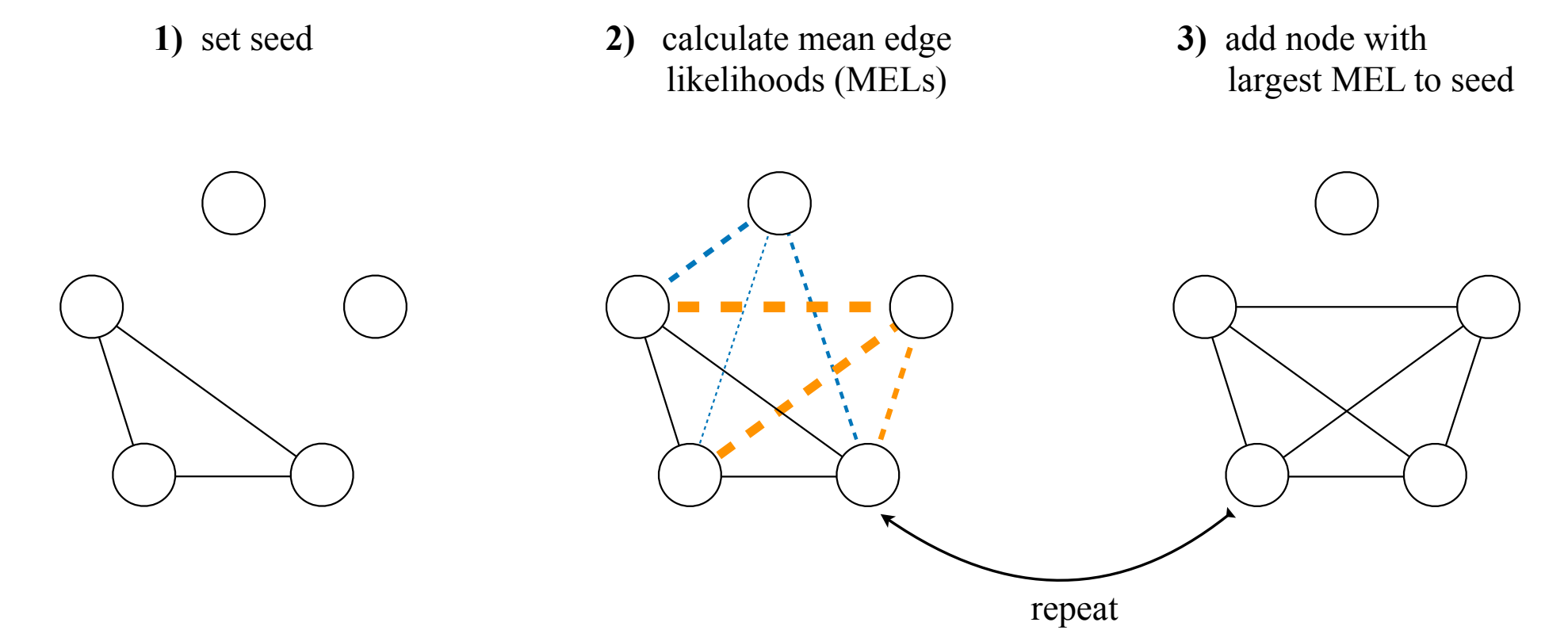
**Table 1:** Concept-concept pair prediction evaluation

## 4 Ongoing Work – Recovering Concept Sets

- Our overall goal is to leverage concept-concept pair prediction models to inform construction of phenotype concept sets.
- We experiment with two algorithms to grow a “seed” of concepts within our concept graph to recover held-out eMERGE phenotypes in our phenotype-aware test sets (See Figures 3 & 4 and Table 2).



**Figure 3:** Phenotype recovery algorithm 1.



**Figure 4:** Phenotype recovery algorithm 2.

Seed size	Precision @ 100%	Proportion @ 100%
10%	0.631 $\pm$ 0.291	15.82 $\pm$ 30.36
20%	0.664 $\pm$ 0.278	9.599 $\pm$ 22.46
30%	0.694 $\pm$ 0.269	7.117 $\pm$ 19.15
40%	0.712 $\pm$ 0.270	5.289 $\pm$ 17.89
50%	0.721 $\pm$ 0.276	4.784 $\pm$ 17.70
60%	0.745 $\pm$ 0.277	4.057 $\pm$ 14.81
70%	0.768 $\pm$ 0.268	3.661 $\pm$ 12.75
80%	0.806 $\pm$ 0.253	2.147 $\pm$ 8.179
90%	0.842 $\pm$ 0.282	1.056 $\pm$ 0.297

**Table 2:** Concept set recovery using algorithm 2 and LR (L1) edge likelihoods

## 5 Conclusion

- Simple and ensemble binary classifiers are capable of faithfully predicting held-out concept-concept pairs.
- This is true even when using phenotype aware train-test splits, suggesting utility in predicting concept-concept pairs for novel concept set construction.
- A simple, greedy algorithm for recovering held-out concept sets performs reasonably well, providing a path forward to high-throughput concept set construction.

