Learning Structured Probability Matrices

Qingqing Huang

2016 February





Laboratory for Information & Decision Systems

Based on joint work with Sham Kakade, Weihao Kong and Greg Valiant.





Infer about the underlying rule θ
 (Estimation, approximation, property testing, optimization of f(θ))



- Infer about the underlying rule θ
 (Estimation, approximation, property testing, optimization of f(θ))
- Challenge:

Exploit our prior for structure of the underlying θ to design fast algorithm that uses as few as possible data X to achieve the target accuracy in learning θ



- Infer about the underlying rule θ
 (Estimation, approximation, property testing, optimization of f(θ))
- Challenge:

Exploit our prior for structure of the underlying θ to design fast algorithm that uses as few as possible data X to achieve the target accuracy in learning θ

Computation Complexity

Sample Complexity

$$\dim(\theta)$$



Discrete probability distribution (support over M outcomes)

N i.i.d. samples $\ N \sim {\rm Poisson}$ (frequency counts over M outcomes)

p

B = Poisson(Np)

.10



0





Goal: find \widehat{p} such that $\|\widehat{p}-p\|_1 \le \epsilon$ N (\uparrow M) **sample complexity:** upper / lower bound









?

Goal: find rank-2 \widehat{B} such that $\|\widehat{B} - \mathbb{B}\|_1 \le \epsilon$ N (\uparrow M) sample complexity: upper / lower bound

Connection to Learning Mixture Models

topic H = 1 or 2

Topic

model



size M vocabulary

 $\Pr(\text{word}_1, \text{word}_2 | \text{topic} = T_1) = pp^\top$ $\Pr(\text{word}_1, \text{word}_2 | \text{topic} = T_2) = qq^\top$

 ${\mathbb B}$ joint distribution over word pairs

Connection to Learning Mixture Models

topic H = 1 or 2





 $\Pr(\text{word}_1, \text{word}_2 | \text{topic} = T_1) = pp^\top$ $\Pr(\text{word}_1, \text{word}_2 | \text{topic} = T_2) = qq^\top$

 ${\mathbb B}\$ joint distribution over word pairs

 $\Pr(\text{output}_1, \text{output}_2 | \text{state} = S_i) = O_i (OQ_i)^\top$

 ${\mathbb B}\$ distribution of past and future outputs

HMM



size M output alphabet

Connection to Learning Mixture Models

topic H = 1 or 2



HMM



output, output_{t+1}

size M output alphabet

 $\Pr(\text{word}_1, \text{word}_2 | \text{topic} = T_1) = pp^\top$ $\Pr(\text{word}_1, \text{word}_2 | \text{topic} = T_2) = qq^\top$

 ${\mathbb B}\$ joint distribution over word pairs

 $\Pr(\text{output}_1, \text{output}_2 | \text{state} = S_i) = O_i (OQ_i)^\top$

 ${\mathbb B}$ distribution of past and future outputs

Spectral Algorithm: N data samples \downarrow Estimate model parameters \uparrow find low rank \hat{B} close to \mathbb{B}

Structured Distribution Learning



Sample complexity	Unstructured	Low rank structure
Estimation	Linear	?
Property Testing	Sub-Linear	?





MLE is non-convex optimization \otimes let's try something intuitive \odot



$$\frac{1}{N}B = \frac{1}{N}\operatorname{Poisson}(N\mathbb{B}) \to \mathbb{B}, \text{ as } N \to \infty$$

• Set \widehat{B} to be the rank 2 truncated SVD of $\frac{1}{N}B$



- + To achieve accuracy $\|\widehat{B} \mathbb{B}\|_1 \le \epsilon$ need $N = \Omega(M^2 \log M)$



• Set
$$\widehat{B}$$
 to be the rank 2 truncated SVD of $\frac{1}{N}B$

- + To achieve accuracy $\|\widehat{B} \mathbb{B}\|_1 \le \epsilon$ need $N = \Omega(M^2 \log M)$
- Not sample efficient! Hopefully $N = \Omega(M)$



$$\frac{1}{N}B = \frac{1}{N}$$
Poisson $(N\mathbb{B}) \to \mathbb{B}$, as $N \to \infty$

- Set \widehat{B} to be the rank 2 truncated SVD of $\frac{1}{N}B$
- + To achieve accuracy $\|\widehat{B} \mathbb{B}\|_1 \le \epsilon$ need $N = \Omega(M^2 \log M)$
- Not sample efficient! Hopefully $N = \Omega(M)$
- Small data in practice!

Word distribution in language has fat tail. More sample documents ${\cal N}$, larger the vocabulary size ${\cal M}$

Main Results

- + Our upper bound algorithm:
 - ✓ Rank-2 estimate \widehat{B} with accuracy $\|\widehat{B} \mathbb{B}\|_1 \le \epsilon$ $\forall \epsilon > 0$
 - ✓ Using $N = O(M/\epsilon^2)$ number of samples
 - ✓ Runtime $O(M^3)$

Main Results

- Our upper bound algorithm:
 - ✓ Rank-2 estimate \widehat{B} with accuracy $\|\widehat{B} \mathbb{B}\|_1 \le \epsilon$ $\forall \epsilon > 0$
 - ✓ Using $N = O(M/\epsilon^2)$ number of samples
 - ✓ Runtime $O(M^3)$
- We prove (strong) lower bound:
 - ✓ Need a sequence of $N = \Omega(M)$ observations to **test** whether the sequence is i.i.d. of unif (M) or generated by a 2-state HMM

Main Results

- + Our upper bound algorithm:
 - ✓ Rank-2 estimate \widehat{B} with accuracy $\|\widehat{B} \mathbb{B}\|_1 \le \epsilon$ $\forall \epsilon > 0$
 - ✓ Using $N = O(M/\epsilon^2)$ number of samples
 - ✓ Runtime $O(M^3)$
- We prove (strong) lower bound:
 - ✓ Need a sequence of $N = \Omega(M)$ observations to **test** whether the sequence is i.i.d. of unif (M) or generated by a 2-state HMM

Sample complexity	Unstructured	Low rank structure
Estimation	Linear	Linear
Property Testing	Sub-Linear	Linear

We capitalize the idea of community detection in sparse random network, SBM is a special case of our problem formulation with homogeneous nodes

We capitalize the idea of community detection in sparse random network, SBM is a special case of our problem formulation with homogeneous nodes

M nodes 2 communities

Expected connection $\mathbb{B} = pp^{\top} + qq^{\top}$ Adjacency matrix $B = \text{Bernoulli}(N\mathbb{B})$



.09	.09	.09	.02	.02	.02
.09	.09	.09	.02	.02	.02
.09	.09	.09	.02	.02	.02
.02	.02	.02	.09	.09	.09
.02	.02	.02	.09	.09	.09
.02	.02	.02	.09	.09	.09

B



We capitalize the idea of community detection in sparse random network, SBM is a special case of our problem formulation with homogeneous nodes

M nodes 2 communities

Expected connection $\mathbb{B} = pp^{\top} + qq^{\top}$ Adjacency matrix $B = \text{Bernoulli}(N\mathbb{B})$



.09	.09	.09	.02	.02	.02
.09	.09	.09	.02	.02	.02
.09	.09	.09	.02	.02	.02
.02	.02	.02	.09	.09	.09
.02	.02	.02	.09	.09	.09
.02	.02	.02	.09	.09	.09

B

.30	.03		1	1	0	0	1	0
.30	.03		1	1	1	0	1	1
.30	.03	generate	0	1	1	0	1	0
.03	.30	estimate	0	0	0	0	1	1
.03	.30		1	1	1	1	1	1
.03	.30		0	1	0	0	1	1
p	\overline{q}				ŀ	3		

We capitalize the idea of community detection in sparse random network, SBM is a special case of our problem formulation with homogeneous nodes

M madea 2 composition	Expected connection	$\mathbb{B} = pp^\top + qq^\top$
M nodes Z communities	Adjacency matrix	$B = \operatorname{Bernoulli}(N\mathbb{B})$

Regularize Truncated SVD:

[Le, Levina, Vershynin]

remove heavy row/column from B, run rank-2 SVD on the remaining graph

.09	.09	.09	.02	.02	.02	
.09	.09	.09	.02	.02	.02	
.09	.09	.09	.02	.02	.02	
.02	.02	.02	.09	.09	.09	
.02	.02	.02	.09	.09	.09	
.02	.02	.02	.09	.09	.09	

B



We capitalize the idea of community detection in sparse random network, SBM is a special case of our problem formulation with homogeneous nodes

M nodes 2 communities	Expected connection	$\mathbb{B} = pp^\top + qq^\top$
	Adjacency matrix	$B = \operatorname{Bernoulli}(N\mathbb{B})$

Regularize Truncated SVD: [Le, Levina, Vershynin] remove heavy row/column from B, run rank-2 SVD on the remaining graph

$M \times M$ m	matrix	Probability matrix	$\mathbb{B} = \rho \rho^\top + \Delta \Delta^\top$
	matrix	Sample counts	$B = \text{Poisson}(N\mathbb{B})$

Key Challenge: In our general setup, we have heterogeneous nodes/ marginal probabilities

Algorithmic Idea 1, Binning

We group the M nodes according to the empirical marginal probability, divide the matrix into blocks, then apply regularized t-SVD to each block



Algorithmic Idea 1, Binning

We group the M nodes according to the empirical marginal probability, divide the matrix into blocks, then apply regularized t-SVD to each block



- ✓ Binning is not exact, we need to deal with spillover!
- ✓ We need to piece together estimates over bins!

Algorithmic Idea 2, Refinement

The coarse estimation from Phase 1 gives some global information. Make use of that to do local refinement for each row / column of B

Phase 2

- 1. Refine the estimate for each node use linear regression
- 2. Achieve sample complexity $N = O(M/\epsilon^2)$ minmax optimal



Take-Away Message



- We identify a problem that lies at the core of many learning problems
- Spectral algorithm is not solving for the non-convex MLE, we need carefully designed algorithm to improve its statistical efficiency
- Coming soon: estimation / approximation / property testing
 of structured probability distribution