HYPERGRAPH-BASED INFERENCE ON NETWORK ROUTES

Rebecca Willett With Jorge Silva, Maxim Raginsky, and Svetlana Lazebnik (UNC)





Monitor p routers.

Each time a packet received, record co-occurrence of routers in packet route.





p routers being monitored $n \ \underline{\text{unlabeled}}$ training observations



- 1. How can we detect anomalous co -occurrences?
- 2. How can we determine which routers are most predictive?

One-class SVM



Scholkopf, Platt, and Shawne-Taylor, 2001

K-point minimum spanning trees



Hero, 2007





Observations

 $x \in \{0,1\}^p$ $x_j = 1$ if j^{th} entity participates in co-occurrence / hyperedge

Challenge: 2^p possible hyperedges (and p is really big)

Orthogonal series

 2^p possible hyperedges = 2^p hypercube vertices = 2^p basis functions = intractable...





Anomalous hyperedges (what we want to estimate):

$$\mathcal{A}^*=\{x:(1-\pi)f(x)<\pi\mu(x)\}$$

Variational approximation



Because $x_j \sim$ Bernoulli, $f_j(x_j) \equiv \theta_j^{x_j} (1-\theta_j)^{1-x_j}$

(Sufficiently rich for many high-dimensional problems)

Leads to O(np)Variational Expectation-Maximization algorithm and computation of annotations

> Beylkin, Garcke, and Mohlenkamp, 2007 McLachlan and Krishnan, 1996 Peterson and Anderson, 1987

Unimodal: variational approximation



Multimodal: mixture of variational approximations

$$f(x) = \sum_{k=1}^{K} \alpha_k \prod_{j=1}^{p} f_{k,j}(x_j)$$

Weight on k^{th}

Hidden data

Consider unobserved binary r.v.:

$$Y \equiv \begin{cases} 1, & \boldsymbol{X} \sim \mu \\ 0, & \boldsymbol{X} \sim f \end{cases}$$



Let
$$\eta(x)\equiv \mathbb{P}(Y=1|X=x,f,\pi)$$

Note \mathcal{A}^* $=$ $\{x:(1-\pi)f(x)<\pi\mu(x)\}$

$$= \{x: \eta(x) > 1/2\}$$

Good estimate of $\eta \Rightarrow$ good estimate of \mathcal{A}^*

Variational EM Algorithm

$$\widehat{\eta}^{(t+1)}(x_i) = \frac{\widehat{\pi}^{(t)}\mu(x_i)}{(1 - \widehat{\pi}^{(t)})\widehat{f}^{(t)}(x_i) + \widehat{\pi}^{(t)}\mu(x_i)}$$

M-step

$$\hat{\pi}^{(t+1)} = \frac{1}{n} \sum_{i=1}^{n} \hat{\eta}_i^{(t+1)}$$

$$\widehat{\theta}_{j}^{(t+1)} = \frac{\sum_{i=1}^{n} (1 - \widehat{\eta}_{i}^{(t+1)}) x_{i,j}}{\sum_{i=1}^{n} (1 - \widehat{\eta}_{i}^{(t+1)})}$$

$$\widehat{f}^{(t+1)}(x_i) = \prod_{j=1}^p (\widehat{\theta}_j^{(t+1)})^{x_{i,j}} (1 - \widehat{\theta}_j^{(t+1)})^{1 - x_{i,j}}$$

 Densities only need to be evaluated at training hyperedges (not all 2^p possible hyperedges)

•Since marginals Bernoulli, complexity linear (not quadratic) in *n*

Wolfe, 1970

Annotations

$$\mathcal{A}_i = \{ x : f(x) < f(x_i) \}$$

hyperedges more anomalous than x_i

smaller $\mathcal{A}_i \Leftrightarrow x_i$ more anomalous

$$p\mathsf{FDR}(\mathcal{A}_i) = \mathbb{P}(X \sim f | X \in \mathcal{A}_i)$$

$$\gamma_i \equiv 1 - p\mathsf{FDR}(\mathcal{A}_i)$$

$$= \frac{\pi \mathbb{U}(\mathcal{A}_i)}{(1 - \pi)\mathbb{F}(\mathcal{A}_i) + \pi \mathbb{U}(\mathcal{A}_i)}$$

Easy to estimate! (Monte Carlo efficient because of variational approx.)

Scott and Kolaczyk, 2007 Storey, 2003

NICO data

n = 250 router-level paths through p = 1105 nodes 3 sources, 83 destinations average path length of 15

Collected using traceroute, but *order* information ignored

Injected 25 random anomalies with average path length of 15

Anomalous distribution: $\mu(x)$ =

 θ_{j}



$$= \prod_{j=1}^{p} \theta_j^{x_j} (1 - \theta_j)^{1 - x_j}$$
$$= \frac{15}{p}$$

Rabbat, Figueiredo and Nowak, 2007

NICO results





NICO results

Mixture weights

w₁ •^w2 0.5 0.5 w₃ Mixture Weights 0.3 0.3 400 1000 200 600 800 π 0.5 800 200 400 600 1000 0.1 0.5 $\overset{0}{\overset{}_{0}}{}^{\scriptscriptstyle \mathsf{L}}$ 10 15 EM iteration 5 15 20 25 200 600 Host index 400 800 1000

Variational approx. coeffs.

Running times

	Variational EM	OCSVM	K-MST
Simulation,	31.7	6.3	3.7
n=100, p=2000			
Simulation,	32.7	57.2	49.5
n=300, p=2000			
NICO data,	5	-	30
n=277, p=1101	_		
Enron data,	504.0	1182.5	≫ 3600
n=800,			
<i>p</i> =75511			

Identify routers with strong predictive power

Our mission: Estimate 2^p Walsh basis coefficients with *fast* multiscale algorithm



Walsh functions in 1-d: $x \in \{0, 1\}$



1 x

Walsh, 1923 Talagrand, 1994 Dinur, Friedgut, Kindler and O'Donnell, 2007

Walsh basis in *p* dimensions

 $s \in \{0,1\}^p$ \checkmark Index of basis functions

$$\varphi_s(x) \stackrel{\triangle}{=} \prod_{j=1}^p \varphi_{s(j)}(x(j))$$

Similar to variational approx. from earlier. Each marginal = 1-d Walsh function.

Results in orthonormal basis of $\{0,1\}^p$

$$f = \sum_{s \in \{0,1\}^p} \theta_s \varphi_s$$

Liang and Krishnaiah, 1985 Ott and Kronmal, 1976

Orthogonal series estimation

1. Compute empirical basis coefficients:

$$\widetilde{\theta}_s = (1/n) \sum_{i=1}^n \varphi_s(X_i)$$

2. Threshold:

$$\widehat{\theta}_s = I_{\{T(\widetilde{\theta}_s) \ge \lambda\}}$$

3. Sum weighted basis functions:

$$\widehat{f} = \sum_{s \in \{0,1\}^p} \widehat{\theta}_s \varphi_s$$

Problem: 2^p basis coefficients!

Key insight



 $\varphi_s = \varphi_u \otimes \varphi_v$

whenever concatenation uv = s

Key insight



Determining if a prefix is "significant"

Main idea: Fix some $u \in \{0, 1\}^k$. If coefficient θ_s is "insignificant" (i.e. going to be thresholded) whenever s has prefix u, then don't compute any θ_s explicitly.

Consider energy of coefficients with index prefix \boldsymbol{u}

$$W_u \stackrel{\triangle}{=} \sum_{v \in \{0,1\}^{p-k}} \theta_{uv}^2$$

If $W_u < \lambda$ then all $\theta_{uv}^2 < \lambda$, so u is not significant.

 W_u can be estimated from data.

Hierarchical estimation



Raginsky, Lazebnik, Willett and Silva, 2008 Goldreich and Levin, 1989 Kushilevitz and Mansour, 1993

Performance

Theorem: If f is in a weak- ℓ_p ball in the Walsh basis, so that

$$|\theta_{(m)}| \le R \cdot m^{-r-1/2},$$

and thresholds are sufficiently large, then

$$\sup_{f \in \mathcal{F}_d(p)} \mathbb{E}_f \, \|f - f\|_2^2 \le C 2^{-d} \left(\frac{\log n}{n}\right)^{2r/(2r+1)}$$

Minimax rate:

$$\inf_{f} \sup_{f \in \mathcal{F}_{d}(p)} \mathbb{E}_{f} \| f - f \|_{2}^{2} \ge C 2^{-d/(2r+1)} n^{-2r/(2r+1)}$$

Raginsky, Lazebnik, Willett and Silva, 2008 Johnstone, 1994

Experimental validation

p = 15, n = 10000, true density = product of marginals



Theorem: For any $\delta \in (0,1)$, if thresholds depend on δ and are sufficiently large, then with probability at least $1-\delta$ the computational complexity of the algorithm is

$$O(n^2 p(n/2^p \log n)^{1/(2r+1)} K)),$$

where K = K(d) is a bound on the thresholds.

MSE vs. Complexity



NICO Results

- p = 20 routers,
- n = 250 routes



VARIATIONAL APPROXIMATION FOR HYPERGRAPH ANOMALY DETECTION

- O(np) computational complexity
- Implementable and effective even when p=10,000+ without any specialized data structures or hardware
- No parameter tuning or bandwidth selection
- False discovery rate can be quickly and accurately calculated

WALSH BASIS FOR DENSITY ESTIMATION

- Multiscale algorithm fast and accurate
- Allows predictive constellations of routers to be indentified

http://www.ee.duke.edu/nislab/

