

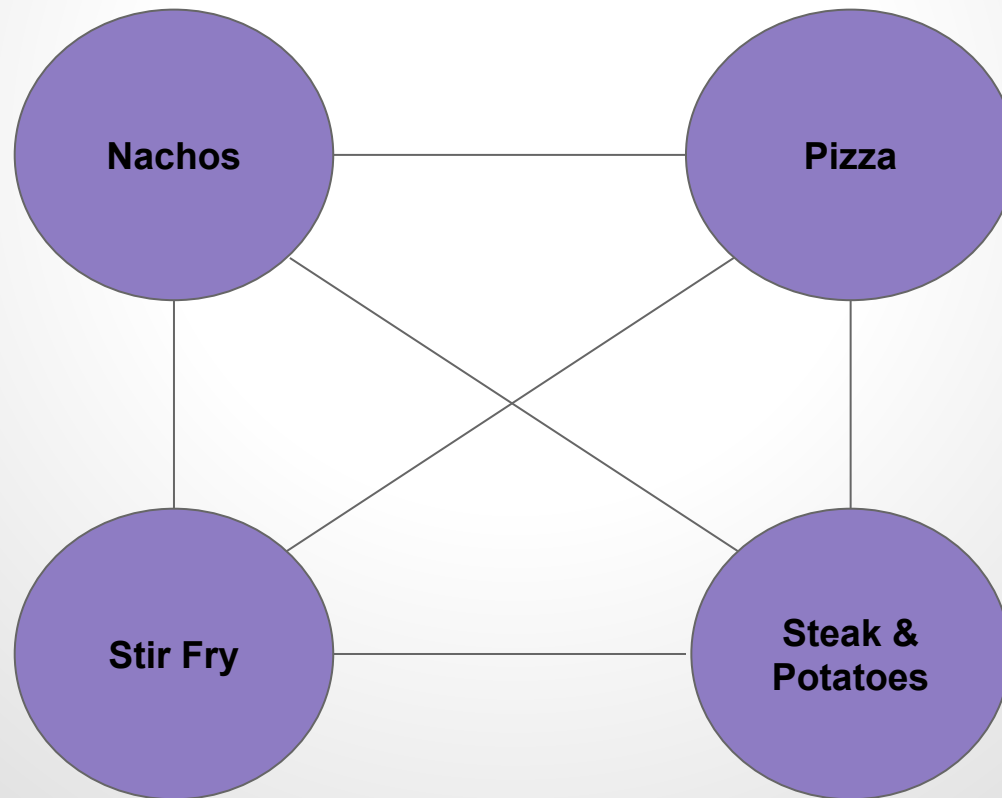
Higher Order Networks & BuildHon+

Steven Krieg

The Problem

How do we represent big data as a network, while accurately preserving dependencies?

Quite a problem, indeed...



First-order network

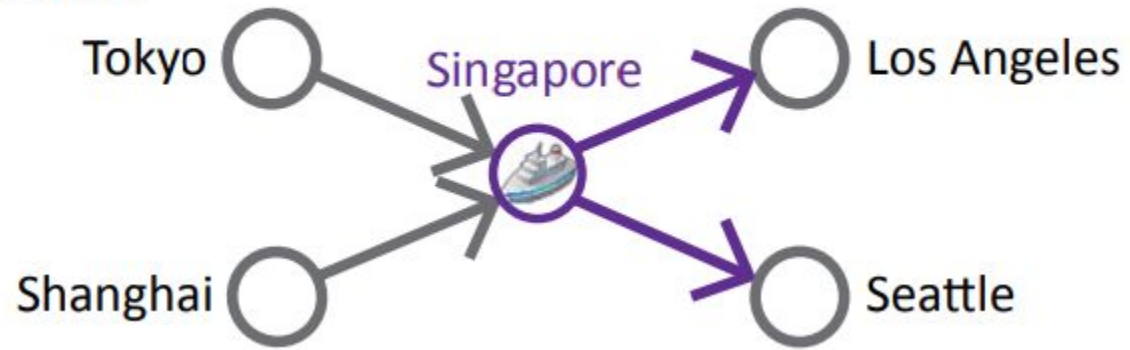
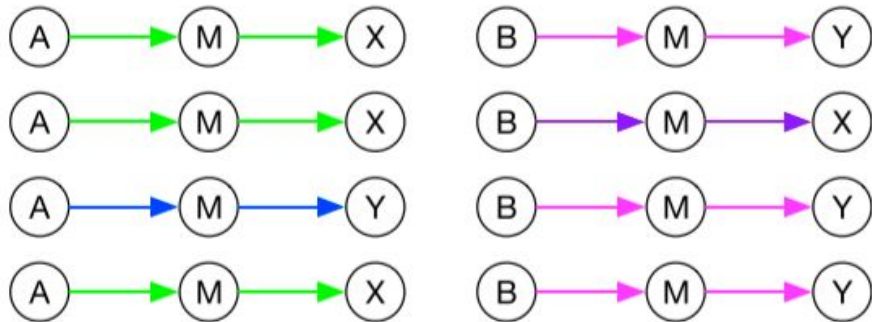


Image from [1]

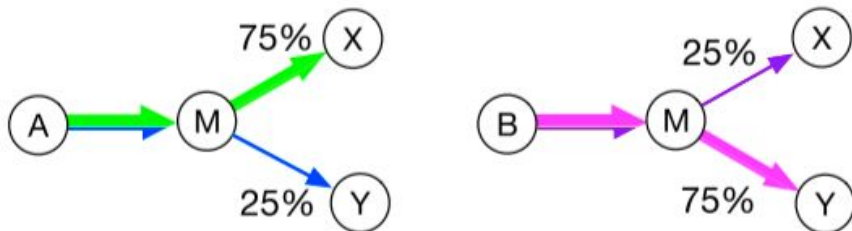
A Solution!

Raw event sequence data



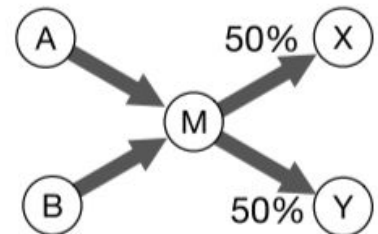
Extract higher-order dependencies from raw event sequences

Higher-order dependencies



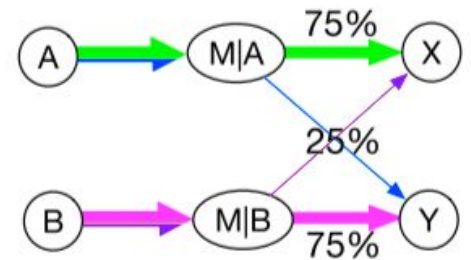
Count number of pairwise interactions as edge weights

First-order network

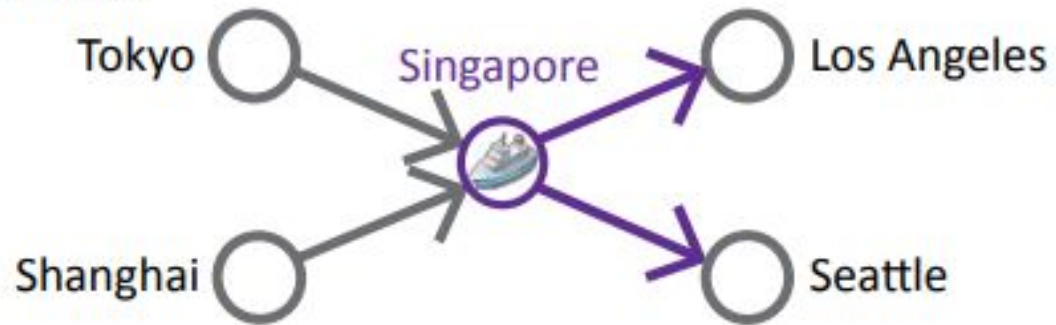


Construct HON based on the extracted rules

Higher-order network



First-order network



Higher-order network (HON)



Image from [1]

Random Walker Results*

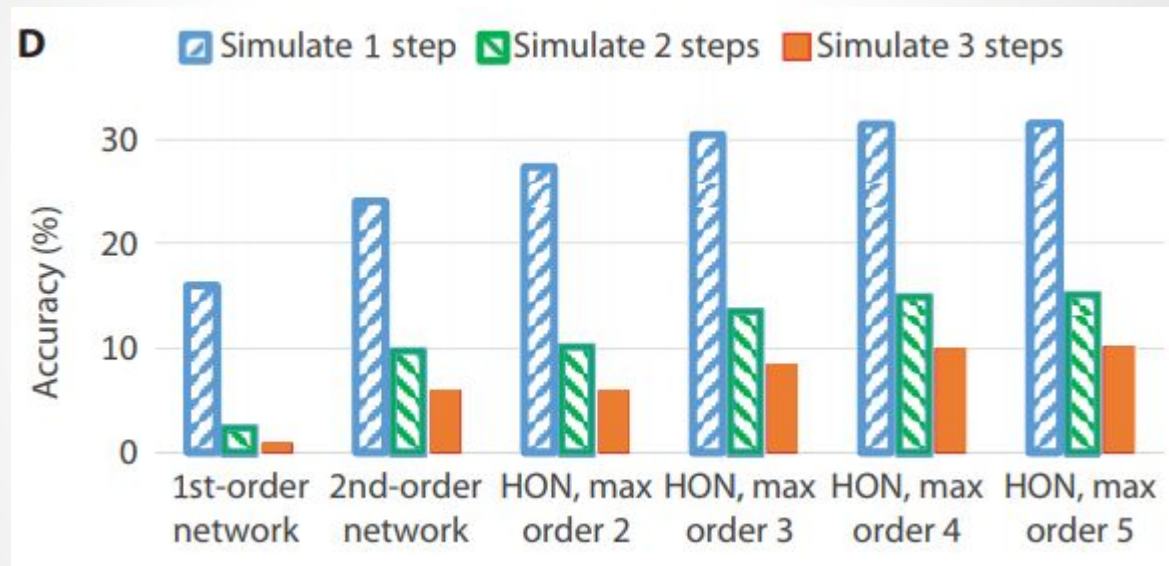


Image from [1]

*The goal of this experiment was not link prediction but to demonstrate the **improvement in representation** using HON.

HON Highlights

- is a network representation of a weighted digraph
- rewires the existing network so that nodes represent a series rather than a singular entity
- enables higher accuracy without needing new analysis tools/algorithms

HON Applications

- ranking (PageRank, etc.)
- anomaly detection
- NLP
- social
- biology
- etc...

The Kernel: BuildHon

Algorithm used to construct the HON

Has 2 main steps:

1. Rule extraction
2. Network rewiring

Step 1: Rule Extraction

1. Count the number of sequential node interactions at the first-order (basically the normal network)
2. Normalize the distributions for each pairwise interaction
3. For each fork node, add the preceding step and see how that changes the distribution of the sequence
4. If the change is “significant” (above a selected threshold), add a second-order dependency and repeat the process recursively to determine higher orders

Step 2: Network Rewiring

1. Construct a conventional first-order network
2. For every second order rule, add the corresponding node; Rewire the previous first-order link to connect to the new higher-order node; then repeat the process for third order rules and so on
3. Once we finish the highest order, rewire all out-edges from that order to connect to nodes with the highest orders possible.

Algorithm 1 HON+ rule extraction algorithm. Given the raw sequential data T , extracts arbitrarily high orders of dependencies, and output the dependency rules R . Optional parameters include $MaxOrder$, $MinSupport$, and $ThresholdMultiplier$

```

1: define global  $C$  as nested counter
2: define global  $D, R$  as nested dictionary
3: define global  $SourceToExtSource$ ,  $StartingPoints$  as dictionary
4:
5: function EXTRACTRULES( $T$ , [ $MaxOrder$ ,  $MinSupport$ ,  $ThresholdMultiplier = 1$ ])
6:   global  $MaxOrder$ ,  $MinSupport$ ,  $Aggressiveness$ 
7:   BUILDFIRSTORDEROBSERVATIONS( $T$ )
8:   BUILDFIRSTORDERDISTRIBUTIONS( $T$ )
9:   GENERATEALLRULES( $MaxOrder$ ,  $T$ )
10:
11: function BUILDFIRSTORDEROBSERVATIONS( $T$ )
12:   for  $t$  in  $T$  do
13:     for ( $Source, Target$ ) in  $t$  do
14:        $C[Source][Target] += 1$ 
15:        $IC.add(Source)$ 
16:
17: function BUILDFIRSTORDERDISTRIBUTIONS( $T$ )
18:   for  $Source$  in  $C$  do
19:     for  $Target$  in  $C[Source]$  do
20:       if  $C[Source][Target] < MinSupport$  then
21:          $C[Source][Target] = 0$ 
22:       for  $Target$  in  $C[Source]$  do
23:         if  $then C[Source][Target] > 0$ 
24:            $C[Source][Target] / (\sum C[Source][*])$ 
25:
26: function GENERATEALLRULES( $MaxOrder$ ,  $T$ )
27:   for  $Source$  in  $D$  do
28:     ADDTORULES( $Source$ )
29:     EXTENDRULE( $Source$ ,  $Source$ , 1,  $T$ )
30:
31: function KLDTHRESHOLD( $NewOrder, ExtSource$ )
32:   return  $ThresholdMultiplier \times NewOrder / \log_2(1 + \sum C[ExtSource][*])$ 
33: function EXTENDRULE( $Valid, Curr$ ,  $order$ ,  $T$ )
34:   if  $Order \leq MaxOrder$  then
35:     ADDTORULES( $Source$ )
36:   else
37:      $Distr = D[Valid]$ 
38:     if  $-\log_2(\min(Distr[*].vals)) < KLDTHRESH-$ 
39:      $OLD(order + 1), Curr$  then
40:       ADDTORULES( $Valid$ )
41:     else
42:        $NewOrder = order + 1$ 
43:        $Extended = EXTENDSOURCE(Curr)$ 
44:       if  $Extended = \emptyset$  then
45:         ADDTORULES( $Valid$ )
46:       else
47:         for  $ExtSource$  in  $Extended$  do
48:            $ExtDistr = D[ExtSource]$ 
49:            $divergence = KLD(ExtDistr, Distr)$ 
50:           if  $divergence > KLDTHRESH-$ 
51:            $OLD(NewOrder, ExtSource)$  then
52:             EXTEN-
53:             DRULE( $ExtSource, ExtSource, NewOrder, T$ )
54:           else
55:             EXTEN-
56:             DRULE( $Valid, ExtSource, NewOrder, T$ )

```

Algorithm 1 (continued)

```

53: function ADDTORULES( $Source$ ):
54:   for  $order$  in  $[1..len(Source) + 1]$  do
55:      $s = Source[0 : order]$ 
56:     if  $not s$  in  $D$  or  $len(D[s]) == 0$  then
57:       EXTENDSOURCE( $s[1:]$ )
58:     for  $t$  in  $C[s]$  do
59:       if  $C[s][t] > 0$  then
60:          $R[s][t] = C[s][t]$ 
61:
62: function EXTENDSOURCE( $Curr$ )
63:   if  $Curr$  in  $SourceToExtSource$  then
64:     return  $SourceToExtSource[Curr]$ 
65:   else
66:     EXTENDOBSERVATION( $Curr$ )
67:     if  $Curr$  in  $SourceToExtSource$  then
68:       return  $SourceToExtSource[Curr]$ 
69:     else
70:       return  $\emptyset$ 
71:
72: function EXTENDOBSERVATION( $Source$ )
73:   if  $length(Source) > 1$  then
74:     if  $not Source[1:]$  in  $ExtC$  or  $ExtC[Source] = \emptyset$  then
75:       EXTENDOBSERVATION( $Source[1:]$ )
76:      $order = length(Source)$ 
77:     define  $ExtC$  as nested counter
78:     for  $Tindex, index$  in  $StartingPoints[Source]$  do
79:       if  $index - 1 \leq 0$  and  $index + order <$ 
80:        $length(T[Tindex])$  then
81:          $ExtSource = T[Tindex][index - 1 : index +$ 
82:          $order]$ 
83:          $ExtC[ExtSource][Target] += 1$ 
84:          $StartingPoints[ExtSource].add((Tindex, index -$ 
85:          $1))$ 
86:       if  $ExtC = \emptyset$  then
87:         return
88:       for  $S$  in  $ExtC$  do
89:         for  $t$  in  $ExtC[s]$  do
90:           if  $ExtC[s][t] < MinSupport$  then
91:              $ExtC[s][t] = 0$ 
92:            $C[s][t] += ExtC[s][t]$ 
93:            $CsSupport = \sum ExtC[s][*]$ 
94:         for  $t$  in  $ExtC[s]$  do
95:           if  $ExtC[s][t] > 0$  then
96:              $D[s][t] = ExtC[s][t] / CsSupport$ 
97:              $SourceToExtSource[s[1:]].add(s)$ 
98:
99: function BUILDSOURCETOEXTSOURCE( $order$ )
100:   for  $source$  in  $D$  do
101:     if  $len(source) = order$  then
102:       if  $len(source) > 1$  then
103:          $NewOrder = len(source)$ 
104:         for  $starting$  in  $[1..len(source)]$  do
105:            $curr = source[starting:]$ 
106:           if  $not curr$  in  $SourceToExtSource$  then
107:              $SourceToExtSource[curr] = \emptyset$ 
108:           if  $not NewOrder$  in
109:            $SourceToExtSource[curr]$  then
110:              $SourceToExtSource[curr][NewOrder] =$ 
111:              $\emptyset$ 
112:            $SourceToExtSource[curr][NewOrder].add(source)$ 

```

Scalability :S

Network representation (global shipping data)	Number of edges	Number of nodes	Network density	* Clustering time (mins)	** Ranking time (s)
Conventional first-order	31,028	2,675	4.3×10^{-3}	4	1.3
Fixed second-order	116,611	19,182	3.2×10^{-4}	73	7.7
HON, max order two	64,914	17,235	2.2×10^{-4}	45	4.8
HON, max order three	78,415	26,577	1.1×10^{-4}	63	6.2
HON, max order four	83,480	30,631	8.9×10^{-5}	67	7.0
HON, max order five	85,025	31,854	8.4×10^{-5}	68	7.6

* Using MapEquation with 1000 iterations

** Using PageRank

Time Complexity

$$L * N * \sum_{i=1}^k ((i + 1) R_i)$$

where L is the count of records in the raw data;
N is the number of unique nodes in the raw data;
k is the maximum order of dependency;
 R_i is the count of dependencies at order i

(*the theoretical upper bound is exponential but is not really helpful for real data sets, in which orders of dependency tend to follow an inverse power law)

Data Sets

- Synthetic web clickstreams (11 billion nodes)
- Global shipping data (3,415,577 voyages made by 65,591 ships between May 1st , 2012 and April 30th, 2013)

Related Work

- Going beyond Markov decisions to higher-order dependencies is not a new idea, but most work has focused on stochastic processes rather than on the representation problem.
- VOM (variable-order Markov) models for predictions are related, but occupy a different niche.

Further Resources

- [1] Xu, J., Wickramaratne, T., & Chawla, N. (2016). Representing higher-order dependencies in networks. *Science Advances*, 2(5), E1600028.
- [2] <http://www.higherordernetwork.com/>
- [3] <https://github.com/xyjprc/hon>
- [4] Xu, J., Saebi, M., Ribeiro, B., Kaplan, L., & Chawla, N. (2017). Detecting Anomalies in Sequential Data with Higher-order Networks.
- [5] Cui Jiao, Guo Jun, Zhang Cangsong, & Chang Xiaojun. (2012). Implementation of random walk algorithm by parallel computing. *Fuzzy Systems and Knowledge Discovery (FSKD)*, 2012 9th International Conference on, 2477-2481.