



Log Grouping and Causality Analysis

Fateme Faraji Daneshgar
Research Associate

Polytechnique Montréal

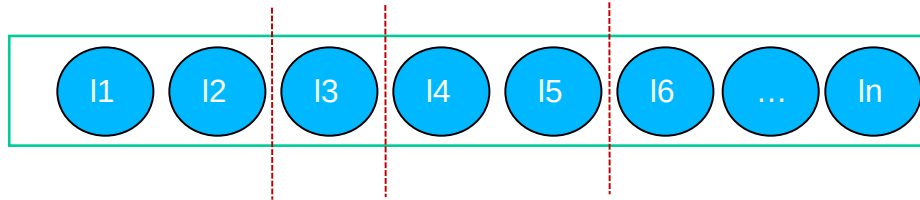
DORSAL Laboratory

Agenda

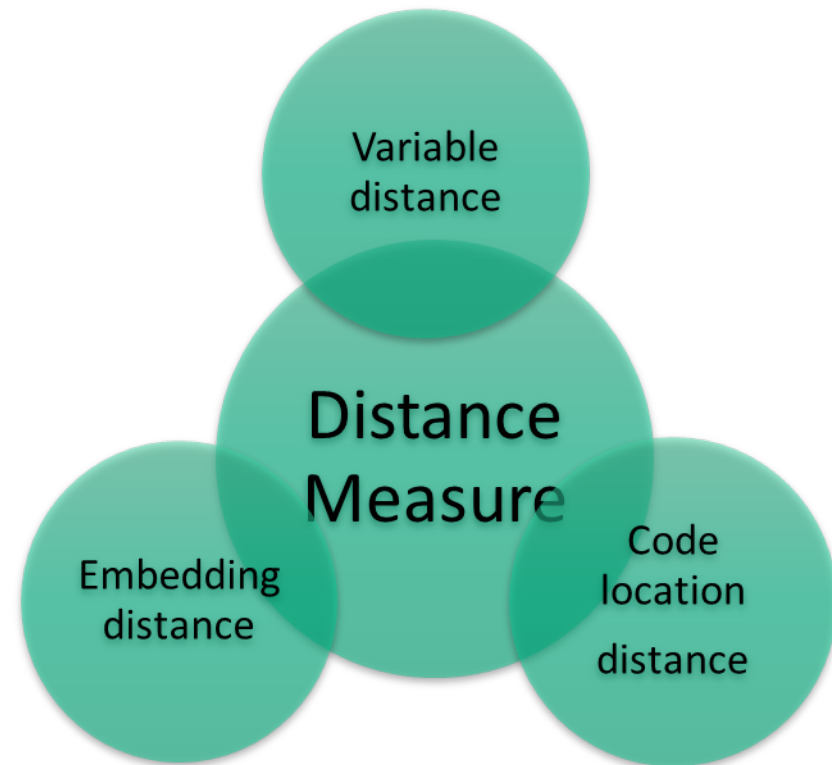
- Log grouping
- Incremental Prefix Tree
- Co-occurrence Probability
- Causality Analysis

Log Grouping

Log Sequence



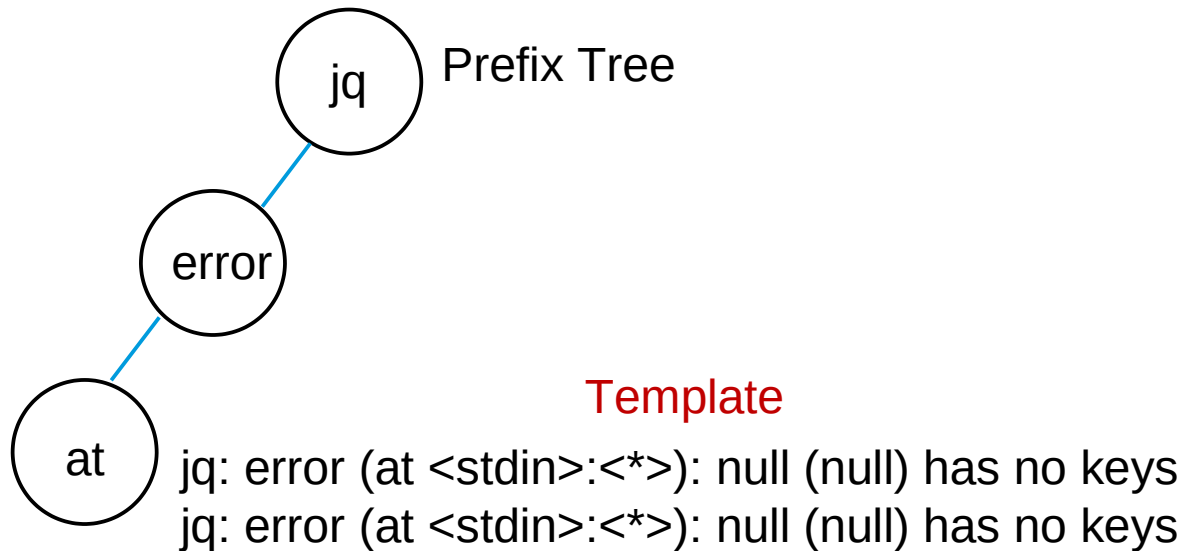
- Log
 - Template
 - Embedding distance
 - Variable
 - Cosine distance



Incremental Prefix Tree

- Parsing logs
 - Drain method
 - Prefix Tree

- jq: error (at <stdin>:45): null (null) has no keys
- jq: error (at <stdin>:21): null (null) has no keys



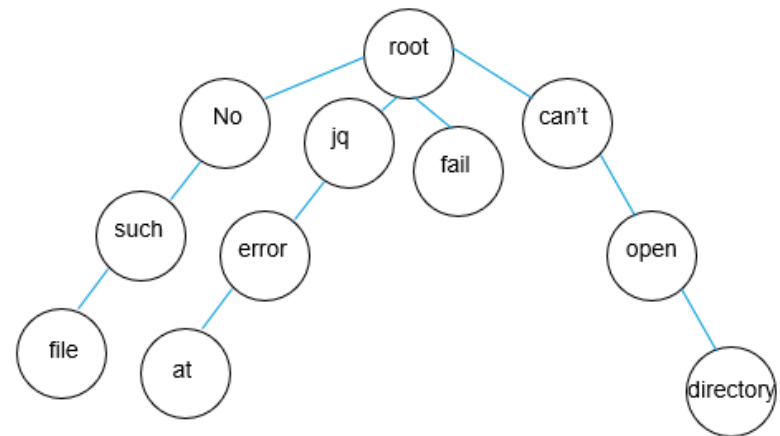
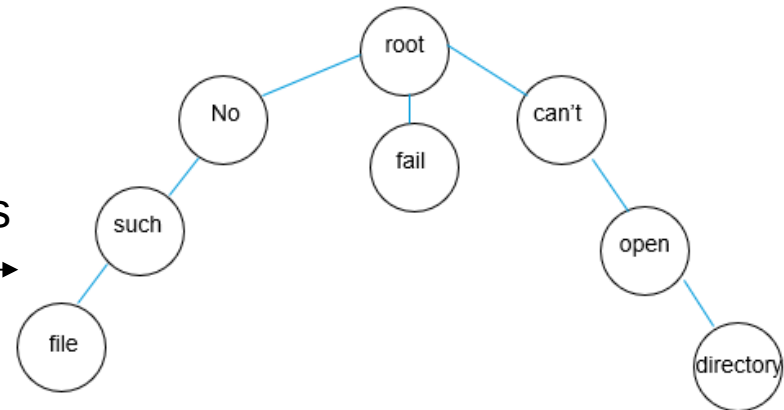
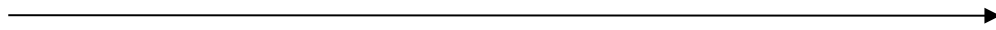
Template

Variables

45
21

Incremental Prefix Tree

jq: error (at <stdin>:<*>): null (null) has no keys



Template

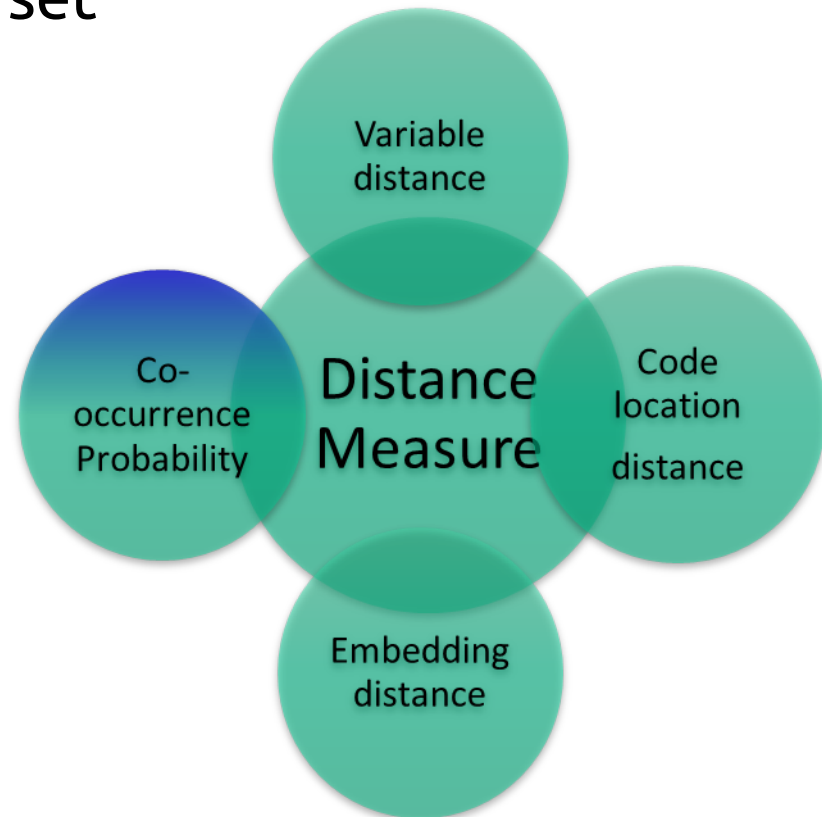
Variables

jq: error (at <stdin>:<*>): null (null) has no keys

--

Co-occurrence Probability

- A log file is an ordered sequence of logs
 - Sequential information
- Representing logs as sequence set
 - Each sequence is a log file
 - Each event is a log



Co-occurrence Probability

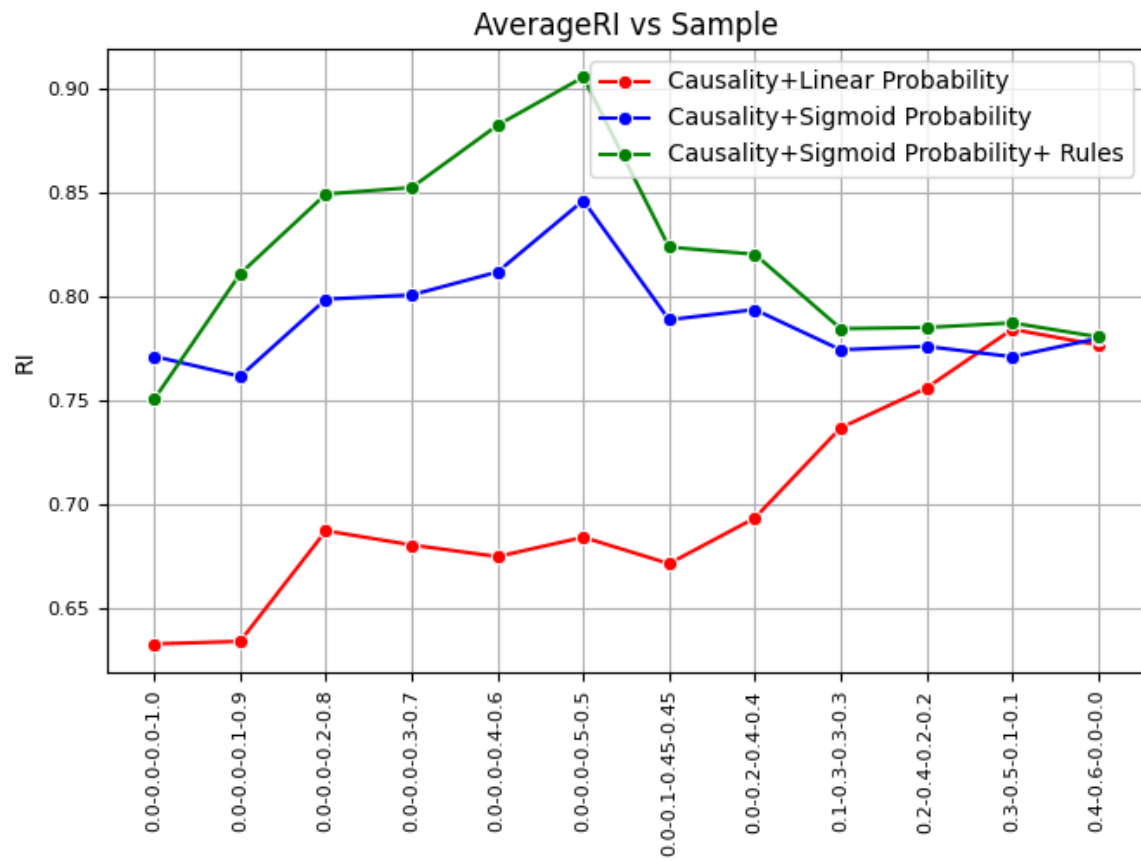
- Considering the time window T:
 - Linear probability
 - Non linear probability
- Suppose
 - M_{ij} : the number of co occurrence time intervals,
 - N_i : the number of time intervals with logj
 - $CP[i][i] = 1$
 - $CP[i][j] = CP[j][i] = \text{Max}(CP[i][j], CP[j][i])$
- Linear : $CP[i][j] = \frac{M_{ij}}{N_j}$,
- Non linear: $CP[i][j] = \frac{1}{1+e^{-x}}$, $x = k \left(M_{ij} - \frac{N_i}{2} \right)$, $k = 1$
-

Co-occurrence Probability Matrix

1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0.4	0.12	0.84	0	0	0	0	0	0	0	0.24	0	0	0	0	0	0	0
0	0.75	0	0.25	0.75	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0.0286	0	0.8	0.0571	0.4286	0	0.0286	0	0	0	0	0	0.2571	0.0286	0	0	0.2	0.1143	0.1429
0	0	0	0	0.8966	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0.0055	0	0.8866	0.0066	0.4262	0.0407	0.0352	0.0286	0.0011	0	0.1828	0.0441	0.011	0.0242	0	0	0
0	0	0	0	0	0	0.9	0.9	0.8	0.7	0.6	0.6	0.4	0	0	0	0	0	0	0
0	0	0	0.0032	0	0.0317	0.0032	0.8508	0.1032	0.0921	0.1016	0.0032	0	0.0063	0.0952	0.054	0.1746	0	0	0
0	0	0	0	0	0.1053	0.0263	0.6053	0.5789	0.9474	0.8947	0.0263	0	0	0.1053	0.0263	0	0	0	0
0	0	0	0	0	0.1111	0.0278	0.6667	0.6111	0.5833	0.8889	0.0278	0	0	0.1111	0.0278	0	0	0	0
0	0	0	0	0	0.0652	0.0217	0.6087	0.5217	0.4783	0.6957	0.0217	0	0	0.1957	0.0217	0	0	0	0
0	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0
0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0.0238	0	0.3095	0	0.4643	0.0417	0.0417	0.0357	0	0	0.869	0.1369	0.25	0.0417	0	0	0
0	0	0	0	0.0822	0.3014	0	0.0411	0	0	0	0	0	0.137	0.3014	0	0.1233	0	0	0
0	0	0	0	0	0	0	0.6364	0	0	0	0	0	0	0.0909	0.2273	0.2273	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.0909	0	0.8663	0	0	0
0	0	0	0.0294	0.1397	0	0	0	0	0	0	0	0	0	0.0699	0	0	0.8346	0.4007	0.2643
0	0	0	0.0122	0.0854	0	0	0	0	0	0	0	0	0	0.0244	0	0	0.939	0.5122	0.2073
0	0	0	0.0556	0.2778	0	0	0	0	0	0	0	0	0	0.1111	0	0	1	0	0
0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0.9091	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Experimental Results

- Co-occurrence probability
 - Linear
 - Non linear
- removing noise

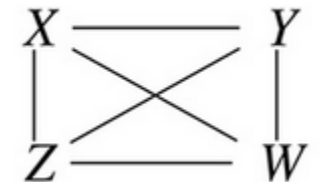
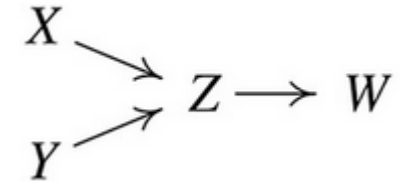


Conditional Independence

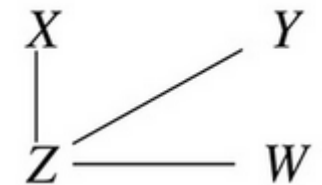
- Co-occurrence does not necessarily mean causality
 - Buying cold water
 - Buying ice cream
 - Hot weather
-
- Conditional Independency ($X \perp Y | Z$)
 - Statistical test, chi-square

PC Algorithm

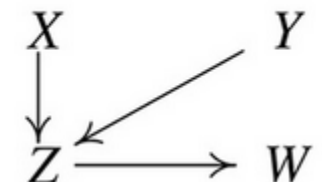
- Suppose we have 4 variable $\{X, Y, Z, W\}$
- Initialize a fully connected undirected graph



- Remove edges based on conditional independence tests
 - If $X \perp Y \mid Z$, remove the edge between X and Y



- Orient Edges to Form a Directed Graph



Chi-square Test

- Step1: create the contingency tables
- Suppose $A \perp B | C$

Group 1: $C = \text{Present}$

Event A	Event B	Frequency
Present	Present	1
Present	Absent	1
Absent	Present	1
Absent	Absent	1

Group 2: $C = \text{Absent}$

Event A	Event B	Frequency
Present	Present	1
Present	Absent	1
Absent	Present	1
Absent	Absent	2

Chi-square Test

- Calculate Expected Frequencies

$$\text{Expected Frequency} = \frac{(\text{Marginal Total of A}) \times (\text{Marginal Total of B})}{\text{Total Number of Observations}} = \frac{2 \times 2}{4} = 1$$

- Step 3: Compute Chi-Square Statistic

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

- If the chi-square statistic exceeds the critical value, we reject the null hypothesis and conclude that A and B are not independent given C.

Data Representation

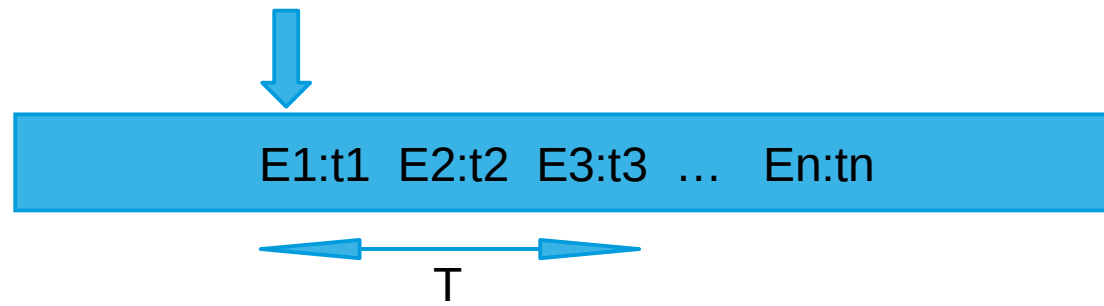
- Considering each log template as a variable to form a data frame

Seq	A	B	C	D
ABC	1	1	1	0
DC	0	0	1	1
ACD	1	0	1	1

- Time series:
 - For each log in each log file
 - The order and time of events are critical
 - Highly sparse
 - Statistical approaches rely on sufficient overlap and density in the data

Hybrid Representation

- The order is considered but not the absolute time



- Pair: $(E1, E2)$, $(E1, E3)$, $(E2, E3)$, $(E1 \perp E3 | E2)$
- Conditional independence test
- 92% accuracy with Ciena log files