

A decorative grid pattern is located in the top left corner of the slide. It consists of a grid of squares in white, light blue, and red. Some squares are outlined in red, and some are filled with color. The pattern is arranged in a way that suggests a data table or a complex structure.

# > Improving Cyber-Security through Predictive Analysis

Peter Frometa  
SPSS

21 November 2005

# Applications of DM & TM in Security



Identify Cyber Threats

Identify Insider Threats

Detect Money Laundering Behavior

Detect Fraud



Detect Smuggling and Drug Trafficking Associations

Predict Non-Compliance of Foreign Visitors

Identify True Identity of Individuals

Better Allocate Security Resources



Ensure Maritime Domain Awareness

Discover Potentially Suspicious Inbound Cargo

Discover Terrorist Organization Affiliation

Predict Risk Potential for Threats to Safety



# Data Mining Defined

## *What is Data Mining?*

A hot buzzword for a class of techniques that find patterns or relationships that have not previously been detected

A user-centric, interactive process which leverages analysis technologies and computing power

Not reliant on an existing database

A relatively easy task that requires knowledge of the organization's problem and subject matter expertise

# Text Mining Defined

## *What is Text Mining?*

A user-centric process which leverages analysis technologies and computing power to access valuable information within unstructured text data sources

Driven by the use of Natural Language Processing and Linguistic based algorithms...not a search engine

No real value unless used in conjunction with Data Mining

Ultimately, eliminates need to manually read unstructured data sources

# Beyond Historical Reporting

## REACTIVE

### *Historical Reporting*

- What are the top source/destination countries?
- What are the top attacked products?
- What are the top offending ISP's
- What IP's have led to the most events?
- What ports are scanned the most?
- What are the IDS event counts for the past day, week, month?

## PROACTIVE

### *Predictive Analytics*

- What activity within the unidentified data of the top source/destination countries is likely to be malicious?
- Of those identified to be malicious, which is likely to be the worst?
- What are the characteristics/ vulnerabilities of products likely to be attacked?
- Which activity did not get flagged that is likely to be malicious?
- Which flagged items are most likely to be false positives?
- What can we expect the IDS event counts to be next week, month, etc.?
- What activity is likely to be associated with worm emergence?
- Which isolated events are likely to be associated with originating from the same source?

# NetFlow Data

	srcip	dstip	srcport	dstport	protocol	packets	bytes	flags	strttime	duration	endtime	sr
1	66.94.234.13	192.35.251.135	53	113	17	1	135	A	28-SEP-04	0	28-SEP-04	PM
2	216.239.57.99	206.33.100.4	53	2050	17	1	135	A	28-SEP-04	1	28-SEP-04	PM
3	17.254.3.183	192.35.251.74	53	30555	17	1	121	A	28-SEP-04	3	28-SEP-04	PM
4	209.202.248.202	192.35.251.74	53	30555	17	1	170	A	28-SEP-04	0	28-SEP-04	PM
5	208.45.133.23	200.150.13.120	53	1004	17	1	194	A	28-SEP-04	1	28-SEP-04	PM

sIP - source IP address (IP address of host that sent the IP packet)

dIP - destination IP address

sPort - port used by the source IP address

dPort - port used by the destination IP address

pro - IP protocol (primarily TCP/UDP in this data)

packets - packet count

bytes - byte count

flags - Transfer Control Protocol header flags

sTime - start time of flow (GMT)

dur - duration of flow (eTime - sTime)

eTime - end time of flow (GMT)

sr - source name or ID of sensor that picked up the data

# Where Does Predictive Analytics Fit?

## Level 1 Analysis

Stop an intrusion event

Log Data

Report

### Based on:

Net Flow data

(source and destination IP address,  
protocol, packets, bytes, flags...)

## Level 2 Analysis

Associate multiple similar intrusions

Classify organized intrusions

Identify state sponsored cyber terror

Net Flow data

Response data

Correlated Incidents, historical data

Keystroke logging

other...

valid web traffic v4.2\* - Clementine 11.0

File Edit Insert View Tools SuperNode Window Help

Streams Outputs Models

- Stream1
- Port Scanning v3.0
- valid web traffic v4.2

sp80-22jan\_04-08.txt Partition known or unknown known\_webservers.txt Table

Strip Blanks Type Merge

Table (15 fields, 501,227 records)

	sIP	dIP	sPort	dPort	pro	packets	bytes	flags	sTime	dur
99420	207.46.156.88	131.37.206.6	80	13674	6	14	16455	S P...	2005-12-01 04:19:51	0
99421	207.46.156.88	131.37.206.6	80	13818	6	1	52	A	2005-12-01 04:19:56	0
99422	207.46.156.88	131.37.206.6	80	13706	6	5	557	F5 P...	2005-12-01 04:19:52	0
99423	207.46.156.88	131.37.206.6	80	13785	6	3	988	F5 P...	2005-12-01 04:19:55	0
99424	207.46.156.88	144.183.224.2	80	24093	6	14	16455	S P...	2005-12-01 04:33:31	0
99425	207.46.156.88	144.147.1.66	80	44403	6	14	16455	S P...	2005-12-01 04:29:04	0
	131.37.206.6	80	13674	6	1	52	A	2005-12-01 04:19:51	0	
	144.147.1.66	80	11001	6	14	16455	S P...	2005-12-01 04:19:20	0	
	144.183.224.2	80	21803	6	14	16455	S P...	2005-12-01 04:01:17	0	
	138.27.1.2	80	34032	6	108	140...	S P...	2005-12-01 04:18:14	18	
	131.37.206.6	80	13658	6	3	988	F5 P...	2005-12-01 04:19:50	1	
	144.183.224.2	80	24021	6	14	16455	S P...	2005-12-01 04:32:19	0	
	131.37.206.6	80	13674	6	3	2459	F5 P...	2005-12-01 04:19:51	0	

Distribution of known or unknown

Value	Proportion	%	Count
unknown web traffic		89.36	447899
valid web traffic		10.64	53328

Table Graph Annotations

Favorites Sources Record Ops Field Ops

Type Filter Derive Filler Reclassify Binning Partition SetToFlag Restructure Transpose Time Intervals History Field Reorder SPSS Transform

Server: Local Server 77MB / 132MB

### Data Preparation/Manipulation

- Merge
- Compare
- Re-Classify
- Match cases to existing lists of valid IPs
- Calculate time between access attempts from same IP
- Calculate average bytes passed to same IP
- plus other data prep steps depending on desired analysis



**valid web traffic v4.2\* - Clementine 11.0**  
File Edit Insert View Tools SuperNode Window Help

Streams Outputs Models  
Stream1  
Port Scanning v3.0  
valid web traffic v4.2

CRISP-DM Classes  
Understanding  
Understanding  
Model  
Deployment

Output Export  
Histogram Neural Net Kohonen C5.0 C&RT K-Means Table Flat File Database

77MB / 132MB

The screenshot displays the Clementine 11.0 interface. At the top is the title bar with the file name 'valid web traffic v4.2\* - Clementine 11.0' and standard window controls. Below is a menu bar (File, Edit, Insert, View, Tools, SuperNode, Window, Help) and a toolbar with various icons for file operations and execution. The main workspace contains a workflow diagram with nodes: 'sp80-22jan\_04-08.txt', 'Partition', 'known or unknown', 'known\_webservers.txt', 'Table', 'Strip Blanks', 'Type', and 'Merge'. A 'known or unknown' node is highlighted with a yellow diamond icon. To the right, a sidebar shows a project tree with 'Streams', 'Outputs', and 'Models' tabs. Below that are 'CRISP-DM' and 'Classes' sections. At the bottom, an 'Output' section lists various model types like Histogram, Neural Net, Kohonen, C5.0, C&RT, K-Means, Table, Flat File, and Database. An overlay in the bottom right corner features a circular diagram with three colored nodes: 'Cluster' (orange), 'Data Mining' (yellow), and 'Predict' (blue), connected by green and blue arrows. A red box highlights the 'Predict' node. In the foreground, a 'known or unknown' dialog box is open, showing a list of rules for 'valid web traffic'. Rule 8, 9, and 10 are visible, each with conditions and a 'then' clause.

```
graph LR
  sp80[sp80-22jan_04-08.txt] --> Partition
  known[known or unknown] --> Partition
  known --> known_web[known_webservers.txt]
  known --> Table1[Table]
  known --> Table2[Table]
  known --> Database[Database]
  known --> Analysis[Analysis]
  known --> C50[C5.0]
  known --> C50_kou[known or unknown]
```

**known or unknown**

File Generate

1 2 3 All

Rule 8 for valid web traffic (39; 0.976)  
if bytes > 3691  
and bytes <= 3758  
and bytes per packet <= 364  
and dur <= 9  
and flags = FSPA  
and packets <= 12  
then valid web traffic

Rule 9 for valid web traffic (569; 0.972)  
if bytes <= 610  
and bytes per packet > 186  
and flags = A  
then valid web traffic

Rule 10 for valid web traffic (34; 0.972)  
if bytes per packet > 363  
and bytes per packet <= 364  
and dur <= 0

Model Summary Settings Annotations

OK Cancel Apply Reset

known or unknown

known or unknown

Cluster  
Data Mining  
Predict

Associate

Table (14 fields, 1,320 records) #1

	date	time	src	proto	s_port	num	orig	action	dst	octet 1	src octet 2	src oct
1	1-Apr-05	0:01:54	144.136.82.3	tcp	55241	2134	144.73.71.201	drop	144.73.205.185	144	136	
2	1-Apr-05	0:04:01	144.136.82.3	tcp	59867	4786	144.73.71.201	drop	144.73.169.165	144	136	
3	1-Apr-05	0:04:05	144.136.82.3	tcp	59936	4826	144.73.71.201	drop	144.73.224.130	144	136	
4	1-Apr-05	0:04:40	144.136.82.3	tcp	61245	5460	144.73.71.201	drop	144.73.230.132	144	136	
5	1-Apr-05	0:05:12	144.136.82.3	tcp	62378	6139	144.73.71.201	drop	144.73.35.187	144	136	
5	1-Apr-05	0:05:13	144.136.82.3	tcp	62419	6158	144.73.71.201	drop	144.73.251.250	144	136	
7	1-Apr-05	0:05:36	144.136.82.3	tcp	63243	6552	144.73.71.201	drop	144.73.194.174	144	136	
3	1-Apr-05	0:06:21	144.136.82.3	tcp	64826	7508	144.73.71.201	drop	144.73.88.20	144	136	
9	1-Apr-05	0:08:59	144.136.82.3	tcp	8297	10721	144.73.71.201	drop	144.73.201.46	144	136	
10	1-Apr-05	0:09:12	144.136.82.3	tcp	8766	11041	144.73.71.201	drop	144.73.171.35	144	136	

Clementine

Streams Outputs Models

- Stream1
- Port Scanning v2.2
- valid web traffic v4.2
- valid dns behavior v4.6

Table Annotations

Table Table



Favorites Sources Record Ops Field

Table Matrix Analysis Data Audit Statistics Quality Report Set Globals Publisher



**Clementine**

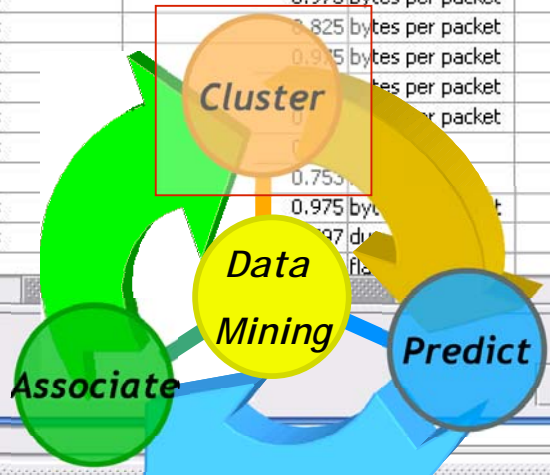
Streams Outputs Models

Anomaly

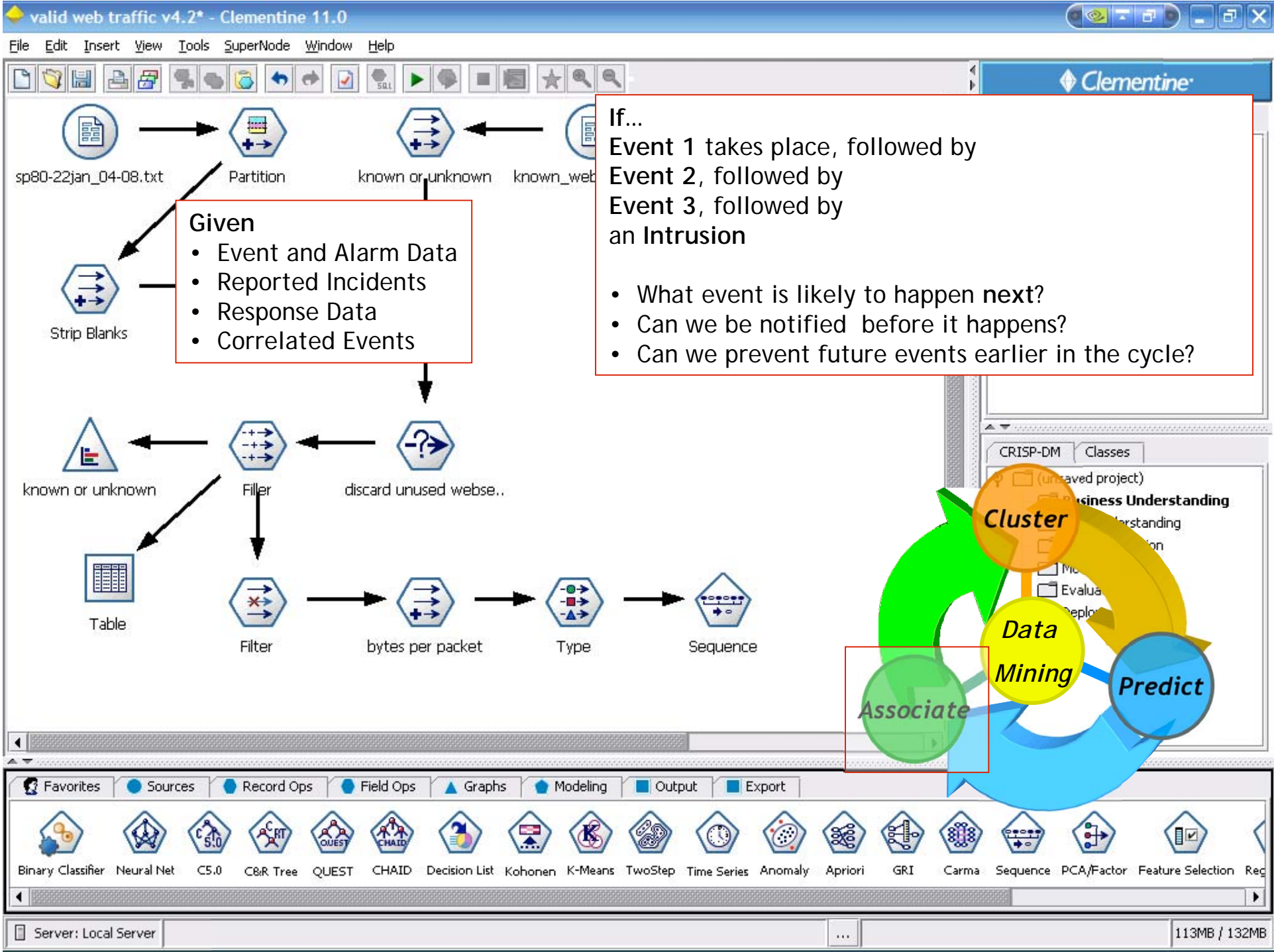


Table (23 fields, 4,883 records)

	\$O-Anomaly	\$O-AnomalyIndex	\$O-PeerGroup	\$O-Field-1	\$O-FieldImpact-1	\$O-Field-2	\$O-f
1	T	4.564	4	flags	0.975	bytes per packet	
2	T	5.000	4	dur	0.888	bytes per packet	
3	T	5.000	4	flags	0.368	bytes per packet	
4	T	5.000	4	bytes per packet	0.512	flags	
5	T	5.000	4	flags	0.981	bytes per packet	
5	T	4.568	4	flags	0.975	bytes per packet	
7	T	4.567	4	flags	0.975	bytes per packet	
8	T	4.771	4	flags	0.825	bytes per packet	
9	T	4.567	4	flags	0.975	bytes per packet	
10	T	4.966	4	flags		bytes per packet	
11	T	4.567	4	flags		bytes per packet	
12	T	5.000	4	flags		bytes per packet	
13	T	5.000	4	dur	0.753		
14	T	4.568	4	flags	0.975	bytes per packet	
15	T	5.000	4	flags	0.97	dur	
16	T	4.636	2	dur		flags	



Binary Classifier Neural Net C5.0 C&R Tree QUEST CHAID Decision List Kohonen K-Means TwoStep Time Series Anomaly Apriori GRI Carma Sequence PCA/Factor Feature Selection Reg



honeypot v2.1 - Clementine 9.0

File Edit Insert View Tools SuperNode Window Help

honeypot\_1224.txt

Table

Filler

Append

normal traffic.txt

port 53 October 8.tx...

traffic type

JTF-GNO Online Report - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites Media

JTF-GNO Predictive Analysis Report  
Joint Task Force - Global Network Operations

Excel Export DeepSight Reporting Save as Table

The following incidents matched characteristics of malicious behavior with a confidence level of greater than 90%..

Source IP 65.38.128.10  
Destination IP 199.252.155.234  
Start Time 10/01/2004 00:01:11  
Confidence Level 99.400002%

Source IP 129.72.231.228  
Destination IP 192.156.13.40  
Start Time 10/01/2004 00:01:23  
Confidence Level 99.400002%

Source IP 204.248.20.66  
Destination IP 204.208.20.3

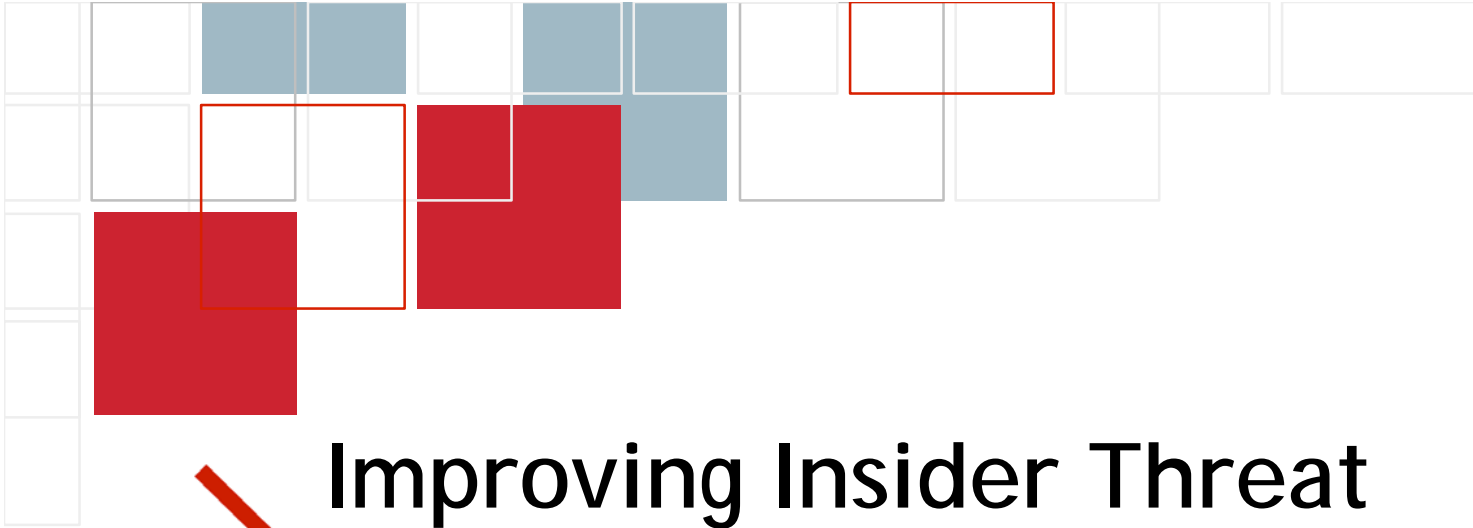
Done My Computer

Server: Local Server 117MB / 376MB

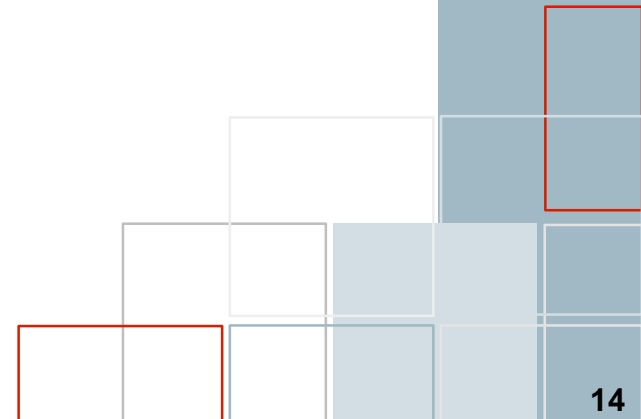
## Deployment

Predictive Analysis as:

- an ad-hoc analysis tool
- a background process



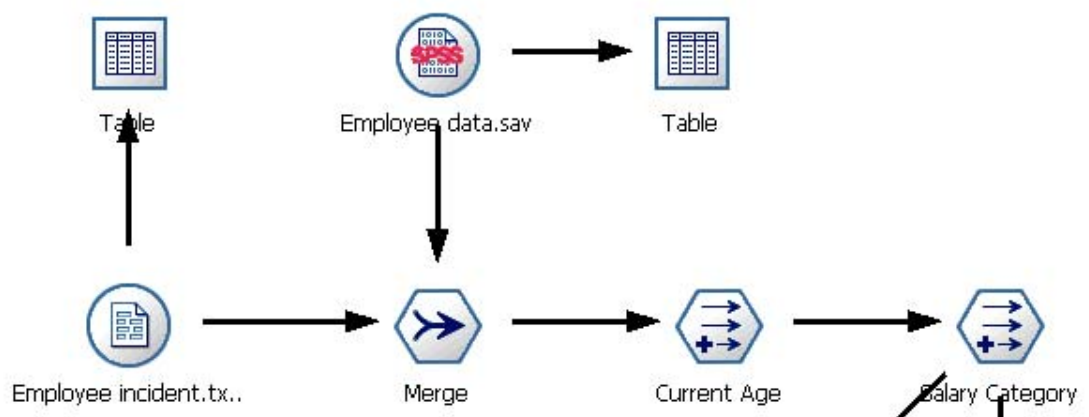
> **Improving Insider Threat  
Detection through  
Predictive Analysis**





Streams Outputs Models

- Stream1
- Port Scanning v3.0
- valid web traffic v4.2
- honeypot v3.0
- honeypot clustering v.1
- Employee Data



Data Audit of [12 fields] #1

Field	Graph	Type	Min	Max	Mean	Std. Dev	Skewness
Date of Birth		Range	1929-07-26	1970-03-14	--	--	--
Educational Level (years)		Range	8	21	13.345	2.920	-0.17
Employment Category		Set	--	--	--	--	--
Current Salary		Range	15750	135000	33247.636	17566.665	2.72
Beginning Salary		Range	9000	79980	16316.321	8440.891	3.89

Audit Quality Annotations

OK

113MB / 132MB

Favorites

Table Matrix

Server: Local Se

match keywords - Clementine 8.0

File Edit Insert View Tools SuperNode Window Help

Streams Outputs Models

- Stream1
- match keywords

Selected documents are run through Clementine's linguistics based text extraction engine. Linguistics-based text mining finds meaning in text much as people do—by recognizing a variety of word forms as having similar meanings and by analyzing sentence structure to provide a framework for understanding the text. This approach offers the speed and cost-effectiveness of statistics-based systems, but it offers a far higher degree of accuracy while requiring far less human intervention.

```

    graph TD
      In(( )) --> TextMining((Text Mining))
      TextMining --> Table1[Table]
      TextMining --> Uppertolower1[uppertolower]
      Uppertolower1 --> Table2[Table]
      Uppertolower1 --> SelectFromDistGraph[SelectFromDist.Graph]
      SelectFromDistGraph --> Table3[Table]
      SelectFromDistGraph --> Uppertolower2[uppertolower]
      Uppertolower2 --> Table4[Table]
      Uppertolower2 --> Knutse((knutse))
      Knutse --> Table5[Table]
      Table5 --> Concept((Concept))
  
```

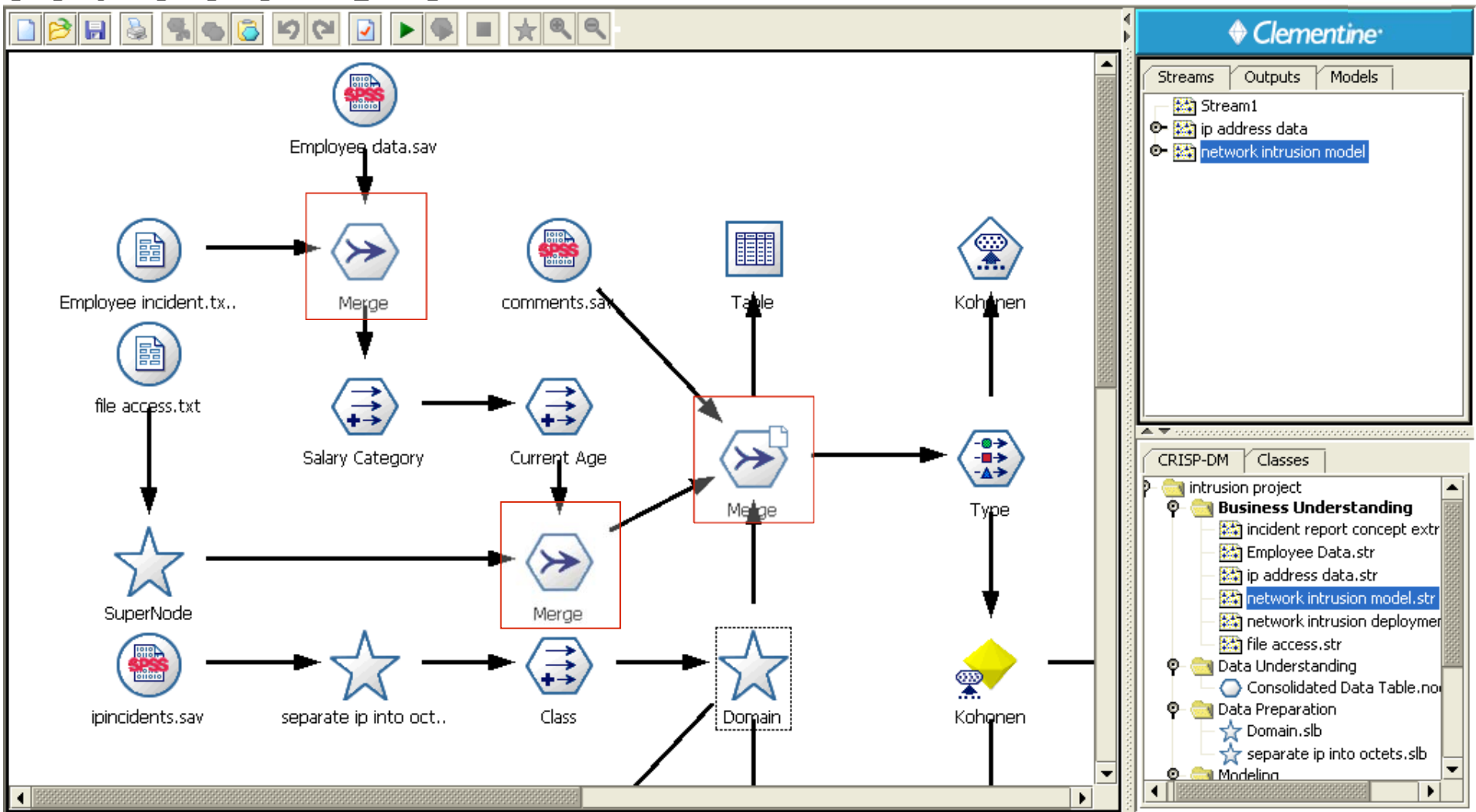
Table (5 fields, 74 records)

	Document	DocID	Concept	Data	Type
1	en\sales\ballbearings sales projections - IPD.doc	1	company_confidential_sales	1	U
2	en\sales\ballbearings sales projections - IPD.doc	1	internal_property	1	U
3	en\sales\ballbearings sales projections - IPD.doc	2	sales_report	1	U
4	en\sales\ballbearings sales projections - IPD.doc	4	\$3,000,000	1	c
5	en\sales\ballbearings sales projections - IPD.doc	4	ball_bearing	1	U
6	en\sales\ballbearings sales projections - IPD.doc	4	expenditures	1	U
7	en\sales\ballbearings sales projections - IPD.doc	6	stock_options_reserve	1	U
8	en\sales\ballbearings sales projections - IPD.doc	7	4.2_million	1	U
9	en\sales\ballbearings sales projections - IPD.doc	7	total_cost	1	U
10	en\sales\ballbearings sales projections - IPD.doc	8	new_investors	1	U
11	en\sales\ballbearings sales projections - IPD.doc	8	reserve_funding	1	U
12	en\sales\ballbearings sales projections - IPD.doc	8	stock_holders	1	U
13	en\eng\Engineering Confidential - Specs for Ba...	12	internal_confidential_memo	1	U
14	en\eng\Engineering Confidential - Specs for Ba...	13	engineering_report	1	U
15	en\eng\Engineering Confidential - Specs for Ba...	14	chief_engineer	1	U
16	en\eng\Engineering Confidential - Specs for Ba...	14	mary_olsen	1	P
17	en\eng\Engineering Confidential - Specs for Ba...	16	2.3mm	1	#
18	en\eng\Engineering Confidential - Specs for Ba...	16	current_tolerance	1	U
19	en\eng\Engineering Confidential - Specs for Ba...	16	new_bearings	1	U
20	en\eng\Engineering Confidential - Specs for Ba...	17	new_bearings	1	U
21	en\eng\Engineering Confidential - Specs for Ba...	18	current_bearings	1	U

Table Annotations

Server: Local Server





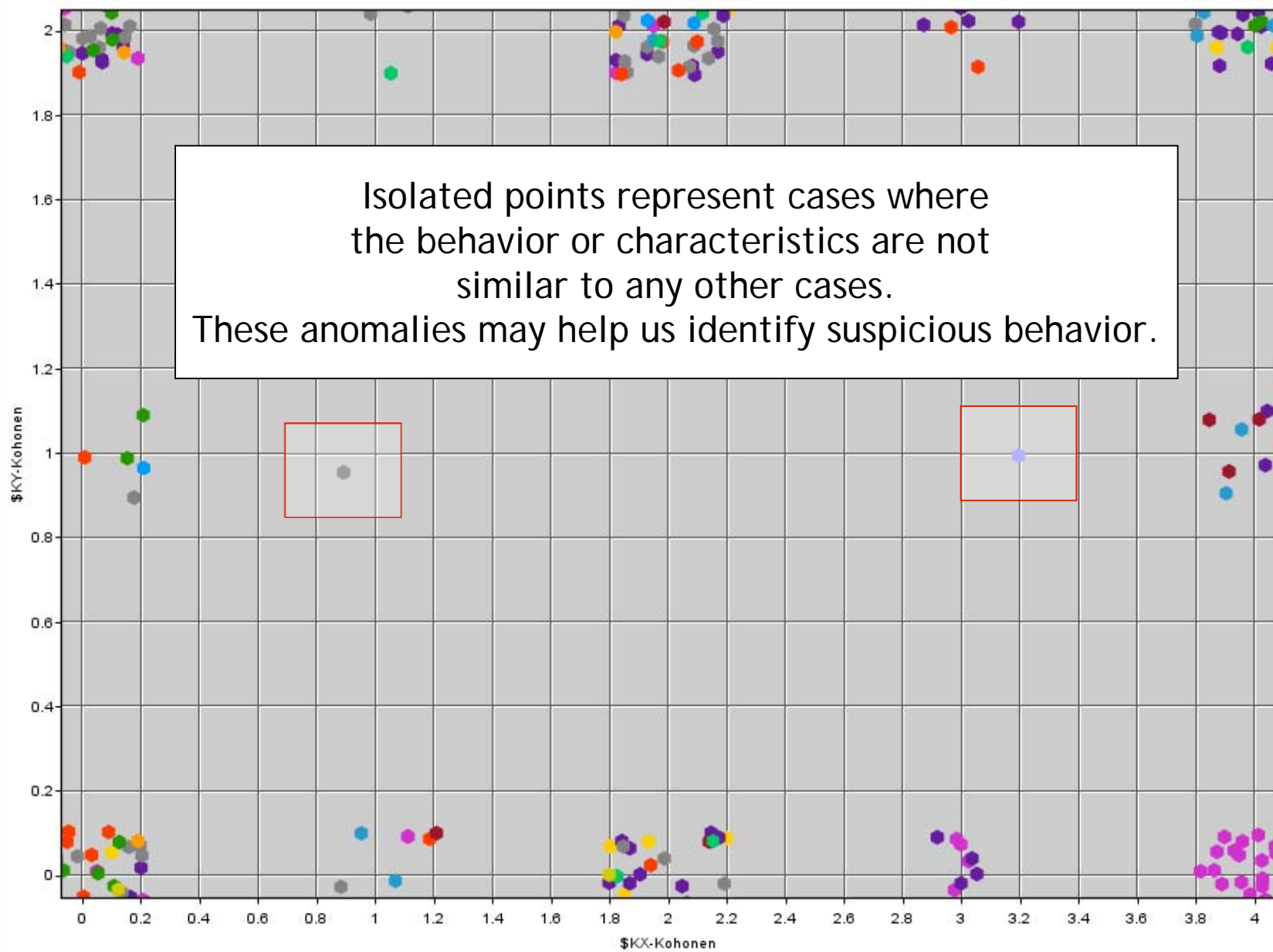
**Clementine**

Streams   Outputs   Models

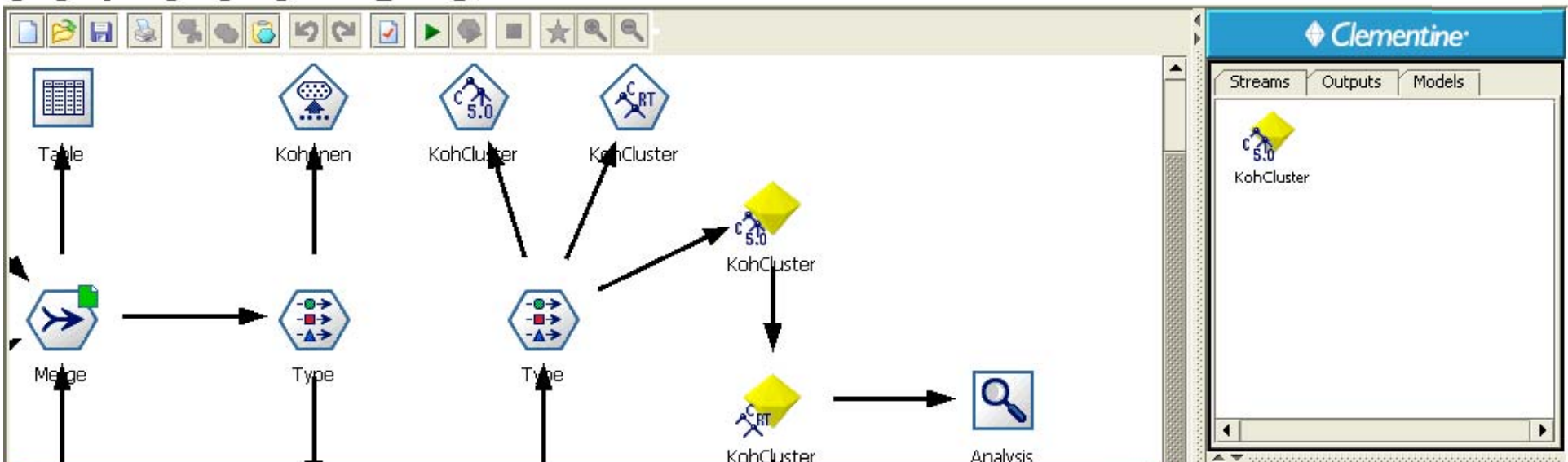
- Stream1
- ip address data
- network intrusion model

CRISP-DM   Classes

- intrusion project
  - Business Understanding
    - incident report concept extr
    - Employee Data.str
    - ip address data.str
    - network intrusion model.str
    - network intrusion deployer
    - file access.str
  - Data Understanding
    - Consolidated Data Table.no
  - Data Preparation
    - Domain.slb
    - separate ip into octets.slb
  - Modeling



- Employment Category
- Management - secure clearance
- Management - Level 5 Clearance
- Management - Level 4 Clearance
- Management - Level 3 Clearance
- Management - Level 2 Clearance
- Management - Level 1 Clearance
- IT staff - non secure
- IT staff - Level 5 Clearance
- IT staff - Level 4 Clearance
- IT staff - Level 3 Clearance
- IT staff - Level 2 Clearance
- Human Resources - non secure
- Accounting - Level 3 Clearance



### KohCluster

File Generate

1 2 3 4 All

```
and Concept_unauthorized_procedure = False
then 4 - 2
```

- Rules for High Risk - contains 3 rule(s)
  - Rule 1 for High Risk
    - if Employment Category = Management - Level 5 Clearance
    - then High Risk
  - Rule 2 for High Risk
    - if Class in [ "Class C" "loopback" ]
    - and Employment Category in [ "Human Resources - non secure" "IT staff - Level 4 Clearance" "Management - secure clearance" ]
    - and Salary Category in [ "Between \$30,000 and \$60,000" ]
    - then High Risk
  - Rule 3 for High Risk
    - if Class in [ "Class C" "Class E" ]
    - and Salary Category in [ "Between \$30,000 and \$60,000" "Between \$60,000 and \$100,000" ]
    - and Concept\_unusual\_clickstream\_activity = True
    - then High Risk
- Default: High Risk

Model Summary Annotations

### Clementine

Streams Outputs Models

KohCluster

CRISP-DM Classes

- Incision project
  - Business Understanding
    - incident report concept extraction.str
    - Employee Data.str
    - ip address data.str
    - network intrusion model.str
    - network intrusion deployment.str
    - file access.str
  - Data Understanding
    - Consolidated Data Table.nod
  - Data Preparation
    - Domain.slb
    - separate ip into octets.slb
  - Modeling

Tree K-Means Table Flat File Database

25Mb / 44Mb

network intrusion deployment - Clementine 8.0

File Edit Insert View Tools SuperNode Window Help

consolidated data.sa..

Kohonen

Model Vote

Select risk

Filter

risk

Table

Clementine

Streams Outputs Models

- Stream1
  - match keywords
  - network intrusion model
  - network intrusion deployment

CRISP-DM Classes

intrusion project

- Business Understanding
  - incident report concept extraction
  - Employee Data.str
  - ip address data.str
  - network intrusion model.str
  - network intrusion deployment.str
  - file access.str

Table (12 fields, 9 records)

File Edit Generate

	INCIDENT id	IPADDRESS	Class	domain	Gender	Employment Category	Last Overall Performance	Automobile	Models Vote	Employee Code	Salary Category
1	3	152.63.54.10	Class B	152.63	Female	IT staff - Level 2 Clearance		6 MidSize	High Risk	25	Between \$15,000 and \$30,000
2	29	146.188.162.165	Class B	146.188	Female	Management - Level 1 Clearance		6 MidSize	High Risk	473	Between \$15,000 and \$30,000
3	1755	157.130.33.58	Class B	157.130	Male	Management - Level 3 Clearance	\$null\$	Luxury - Any Size	High Risk	44	Between \$15,000 and \$30,000
4	1955	146.188.162.165	Class B	146.188	Female	Accounting - Level 3 Clearance		6 Sport Utility	High Risk	221	Between \$15,000 and \$30,000
5	1982	157.130.33.58	Class B	157.130	Female	Management - Level 1 Clearance		6 SubCompact	High Risk	82	Between \$15,000 and \$30,000
5	2942	146.188.136.165	Class B	146.188	Male	Human Resources - non secure		4 Sport Utility	High Risk	19	Between \$30,000 and \$60,000
7	3367	157.130.33.58	Class B	157.130	Female	Management - Level 3 Clearance		7 Truck	High Risk	224	Between \$15,000 and \$30,000
3	3371	157.130.33.58	Class B	157.130	Male	Accounting - Level 3 Clearance		8 SubCompact	High Risk	97	Between \$30,000 and \$60,000
3	3439	157.130.194.73	Class B	157.130	Male	Accounting - Level 3 Clearance		6 Sport Utility	High Risk	43	Between \$15,000 and \$30,000

Table Annotations

# Final Thoughts and Questions

- Visit [www.spss.com](http://www.spss.com)
- Future Questions:
  - Peter Frometa - [pfrometa@spss.com](mailto:pfrometa@spss.com)