# semi-supervised data organization
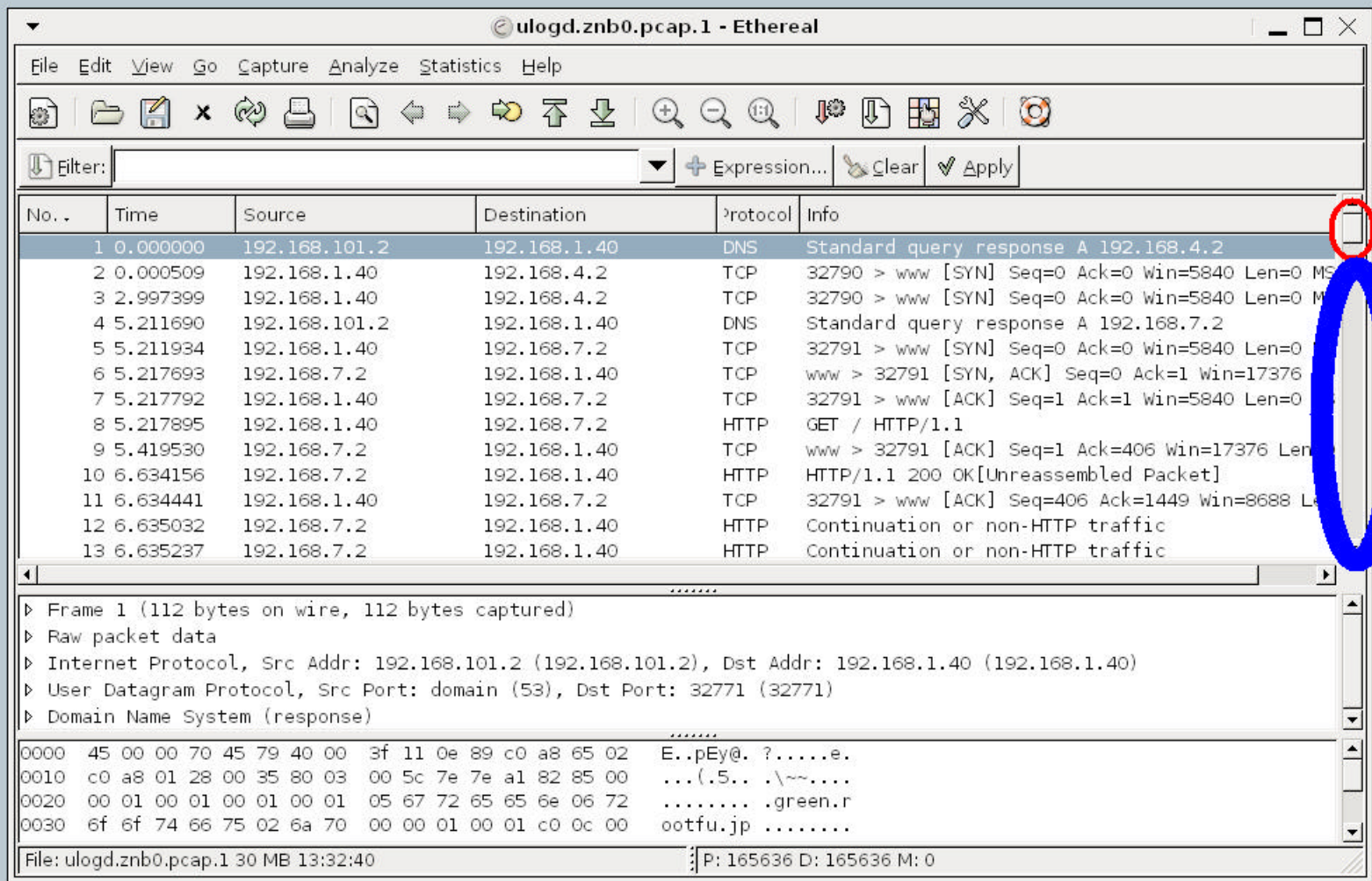
Javed Aslam

Sergey Bratus

► Virgil Pavlu

| Source | Destination | Protocol | |
|---|---|---|---|
| 192.168.1.60 | 192.168.7.2 | ICMP | Echo (ping) request |
| 192.168.1.60 | 192.168.7.2 | ICMP | Echo (ping) request |
| 192.168.1.60 | 192.168.7.2 | ICMP | Echo (ping) request |
| 192.168.1.60 | 192.168.7.2 | ICMP | Echo (ping) request |
| 192.168.1.60 | 192.168.7.2 | ICMP | Echo (ping) request |
| 192.168.1.120 | 192.168.5.2 | ICMP | Echo (ping) request |
| 192.168.1.60 | 192.168.7.2 | ICMP | Echo (ping) request |
| 192.168.1.120 | 192.168.5.2 | ICMP | Echo (ping) request |
| 192.168.1.60 | 192.168.7.2 | ICMP | Echo (ping) request |
| 192.168.1.120 | 192.168.5.2 | ICMP | Echo (ping) request |
| 192.168.1.60 | 192.168.7.2 | ICMP | Echo (ping) request |
| 192.168.1.120 | 192.168.5.2 | ICMP | Echo (ping) request |
| 192.168.1.60 | 192.168.7.2 | ICMP | Echo (ping) request |

84 bytes on wire, 84 bytes captured)

Src Addr: 192.168.1.60 (192.168.1.60), Dst Addr: 192.168.7.2

Control Message Protocol

# motivation : too much data

- log analysis :
  - lots of data
  - hard to search/visualize
  - very few labeled records

- but
  - easy to cluster/classify
  - interesting clusters have high density linkage
  - lots of similar records

# motivation: lots of log data

# Ethereal : Filter



- Filter in Ethereal :
  - if we would know what to look for....
  - similar records do not necessarily match boolean logic
  - filters get too long
  - one at a time, loses big picture
    - difficult browsing

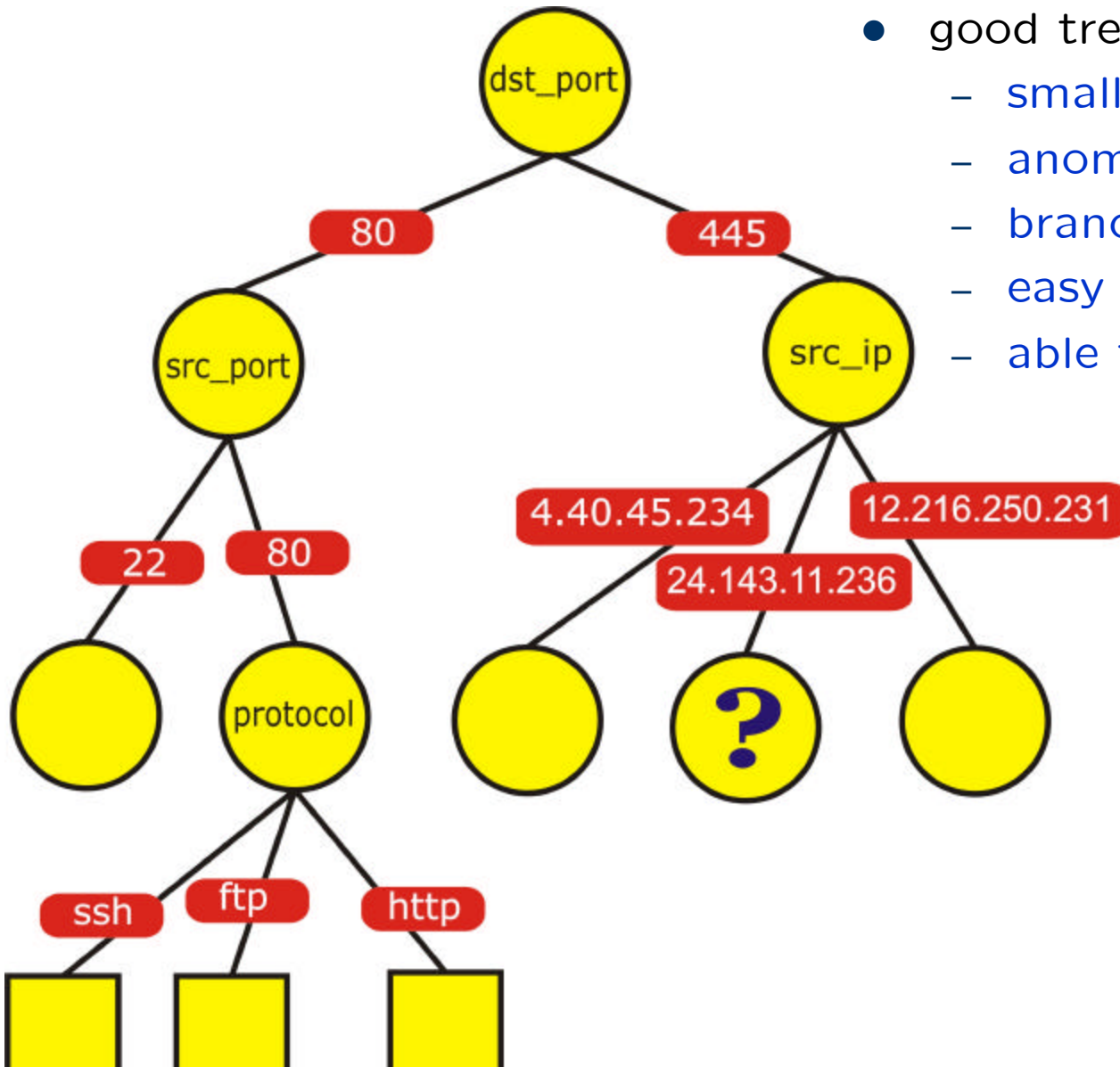# TreeView



TreeView

| File | Edit | Template |

**[1339]** Snort portscan alerts
- ⊞ **[1135]** dst_port: 445 src_ip: [55] dst_ip: [75]
- ⊞ **[70]** dst_port: 80 src_ip: [8] dst_ip: [30]
- ⊞ **[26]** dst_port: 21 src_ip: 80.141.141.173 dst_ip: [11]
- ⊞ **[22]** dst_port: 4899 src_ip: 218.103.195.242 dst_ip: [22]
- ⊞ **[20]** dst_port: 4000 src_ip: [2] dst_ip: [8]
- ⊞ **[15]** dst_port: 139 src_ip: 129.170.125.243 dst_ip: [8]
- ⊞ **[15]** dst_port: 443 src_ip: 211.5.239.5 dst_ip: [9]
- ⊞ **[12]** dst_port: 1524 src_ip: 192.139.15.34 dst_ip: [12]
- ⊞ **[9]** dst_port: 1 src_ip: 209.15.84.72 dst_ip: [9]
- ⊞ **[3]** dst_port: 8000 src_ip: 194.208.40.120 dst_ip: [2]
- ⊞ **[3]** dst_port: 1080 src_ip: 194.208.40.120 dst_ip: [2]
- ⊞ **[3]** dst_port: 3128 src_ip: 194.208.40.120 dst_ip: [2]
- ⊞ **[3]** dst_port: 8100 src_ip: 194.208.40.120 dst_ip: [2]
- ⊟ **[3]** dst_port: 8080 src_ip: 194.208.40.120 dst_ip: [2]
  - ⊟ **[3]** src_ip: 194.208.40.120 dst_ip: [2]
    - Apr 15 19:54:10 annon snort: 194.208.40.120 4743 -> 129
    - Apr 15 19:55:00 annon snort: 194.208.40.120 4914 -> 129
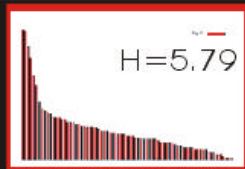    - Apr 15 19:55:06 annon snort: 194.208.40.120 4914 -> 129

**Attributes**

| Field | # | Value |
|---|---|---|
| _day | | 15 |
| _minute | | 54 |
| type | | SYN |
| _month | | Apr |
| loghost | | annon |
| _hour | | 19 |
| _line | | Apr 15 19:54:10 annon snort: 194.208.40.... |
| dst_port | | 8080 |
| src_ip | | 194.208.40.120 |
| src_port | | 4743 |
| flags | | *******S* |
| program | | snort |
| dst_ip | | 129.170.166.39 |
| _second | | 10 |

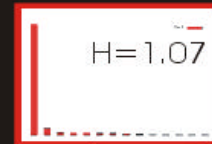# data organization : trees



- good tree :
  - small branching factor
  - anomalies grouped together
  - branches are different
  - easy browsing
  - able to use feedback

# from table to tree

# information theory



$$H(X,Y) = H(X) + H(Y) - I(X,Y) =$$
$$= H(X) + H(Y|X)$$

$$H(Y|X) = H(X,Y) - H(X) =$$
$$= H(Y) - I(X,Y)$$

$$I(X,Y) = H(X,Y) - H(X|Y) - H(Y|X) =$$
$$= H(X) - H(Y|X)$$

# information theory

# information theory

# information theory



H(Y|X)=0.76

H(Y|X)=2.216

H(Y|X)=0.39

Pick me!

H(Y|X)=3.35

# Jensen-Shannon divergence

$$JS(D_1, \ldots, D_n) = H\left(\sum_{k=1}^{n} p_k D_k\right) - \sum_{k=1}^{n} p_k H(D_k)$$

- measures (dis)similarity between several distributions
- almost a distance
- represents information reduction from encoding the distributions together rather than separately
- zero iff all distributions are identical
- better analytical properties than relative entropy

$$JS_{\pi_i}(\mathcal{R}_1, \ldots, \mathcal{R}_{|F_i|}) = H(\sum_{k=1}^{|F_i|} p_k^i \mathcal{R}_k) - \sum_{k=1}^{|F_i|} p_k^i H(\mathcal{R}_k)$$

- $\overline{\mathcal{R}} = \mathcal{R}_1 \cup \ldots \cup \mathcal{R}_{|F_i|}$

- measures dissimmilarities between tree branches

- bar on top of each node indicates the number of records and their class labels(unknown)

# information bottleneck

- [Tishby, Pereira, Bialek]
  - X is the set of objects to be clustered/compressed
  - Y = relevant feature(s)
  - find cluster C to achieve

$$\underset{C}{argmin}\ I(X;C) - \beta I(Y;C)$$

# information bottleneck and JS

- information bottleneck formula

$$\underset{C}{argmin}\, I(X;C) - \beta I(Y;C)$$

- applied with JS divergence

$$F_{(1)} = \underset{\{F_i | H(F_i) \neq 0\}}{argmin} H(F_i) - \beta \cdot JS_{\pi_i}(\mathcal{R}_1, \ldots, \mathcal{R}_{|F_i|}),$$

# semi-supervised

- L=set of labels provided by the user
  - Only a tiny percentage of records will be marked either way.
  - Not all copies of identical records (or very similar) records will be marked
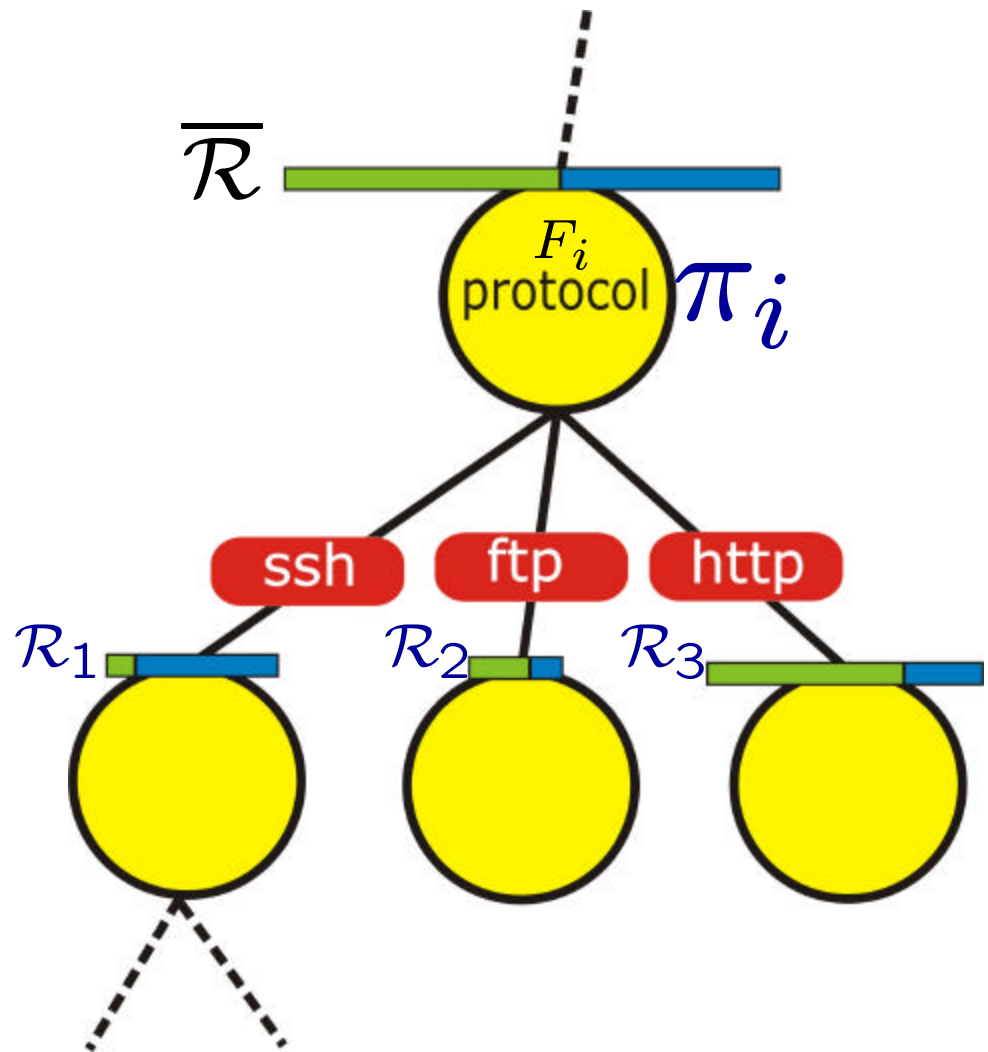
- all quantities of interest become conditionals of L:

$$H(F|L) \;=\; H(F) - I(F;L)$$

$$JS_{\pi_i}(\mathcal{R}_1, \dots, \mathcal{R}_{|F_i|}|L) \;=\; JS_{\pi_i}(\mathcal{R}_1, \dots, \mathcal{R}_{|F_i|})$$
$$- \; \left(I(\overline{\mathcal{R}}; L) - \sum_{k=1}^{|F_i|} p_k^i I(\mathcal{R}_k; L)\right)$$

# semi-supervised
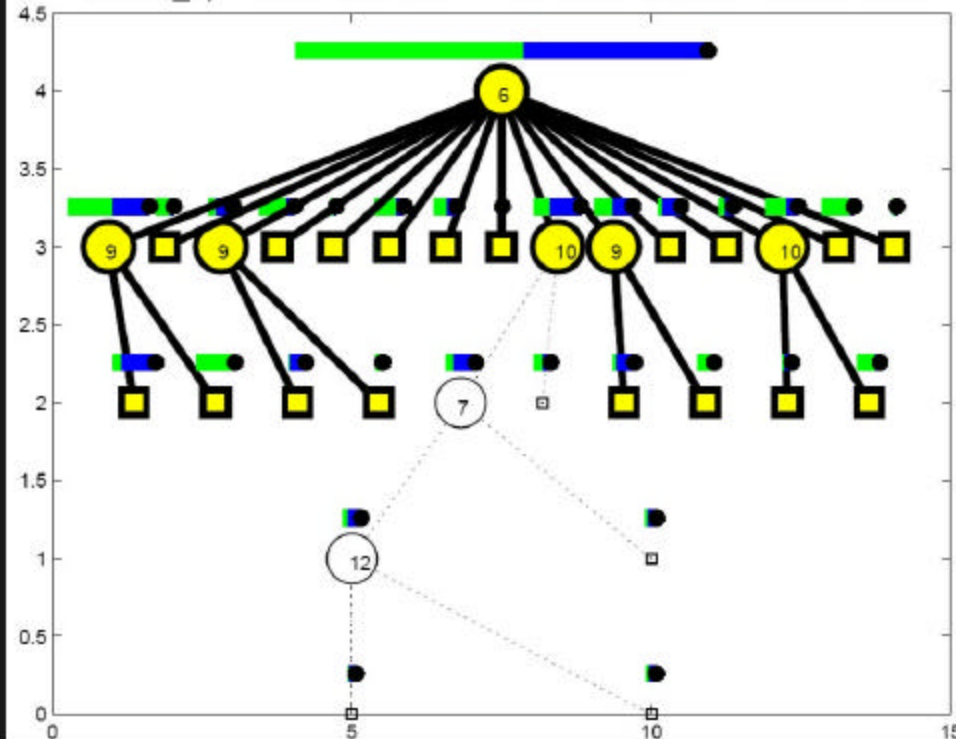
$$F^{(1)} = \underset{\{F_i | H(F_i) \neq 0\}}{argmin} \; H(F_i|L) - \beta \cdot JS_{\pi_i}(\mathcal{R}_1, \ldots, \mathcal{R}_{|F_i|}|L),$$

$$\overline{\mathcal{R}} = \mathcal{R}_1 \cup \ldots \cup \mathcal{R}_{|F_i|}$$
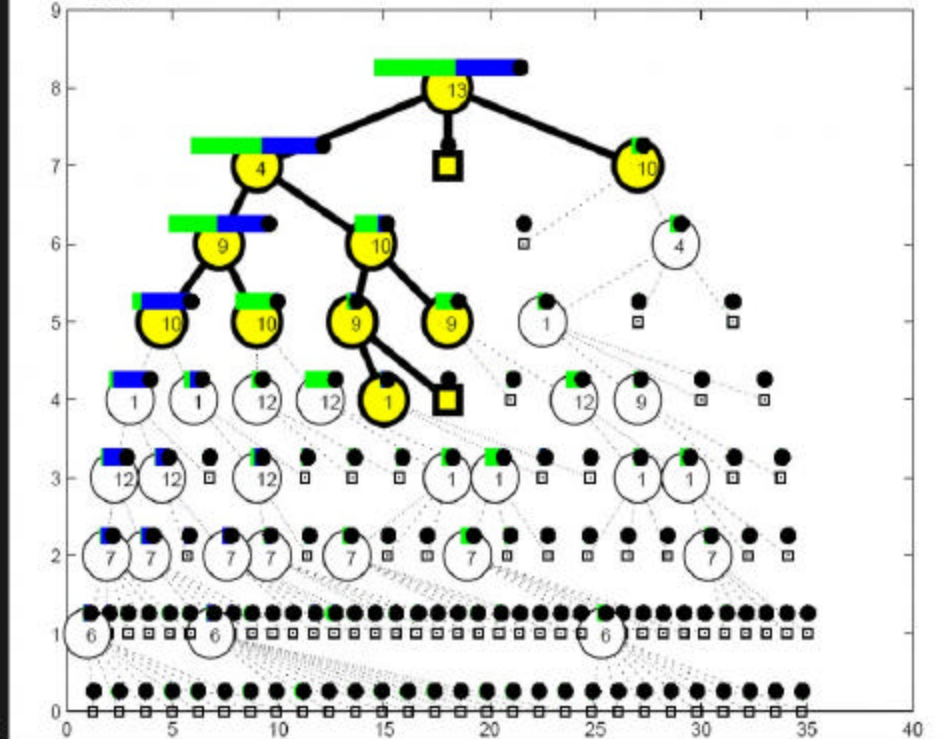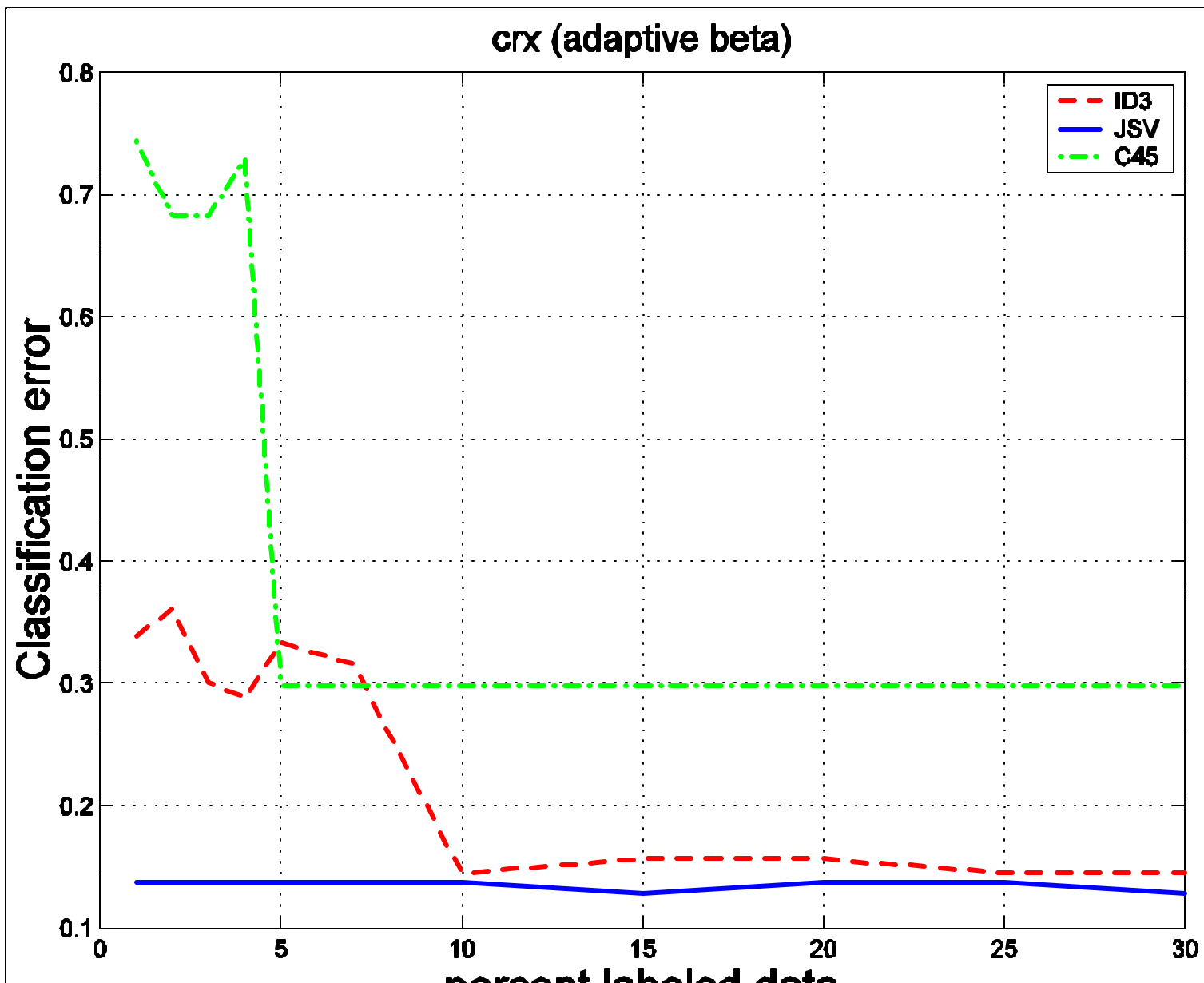
# results
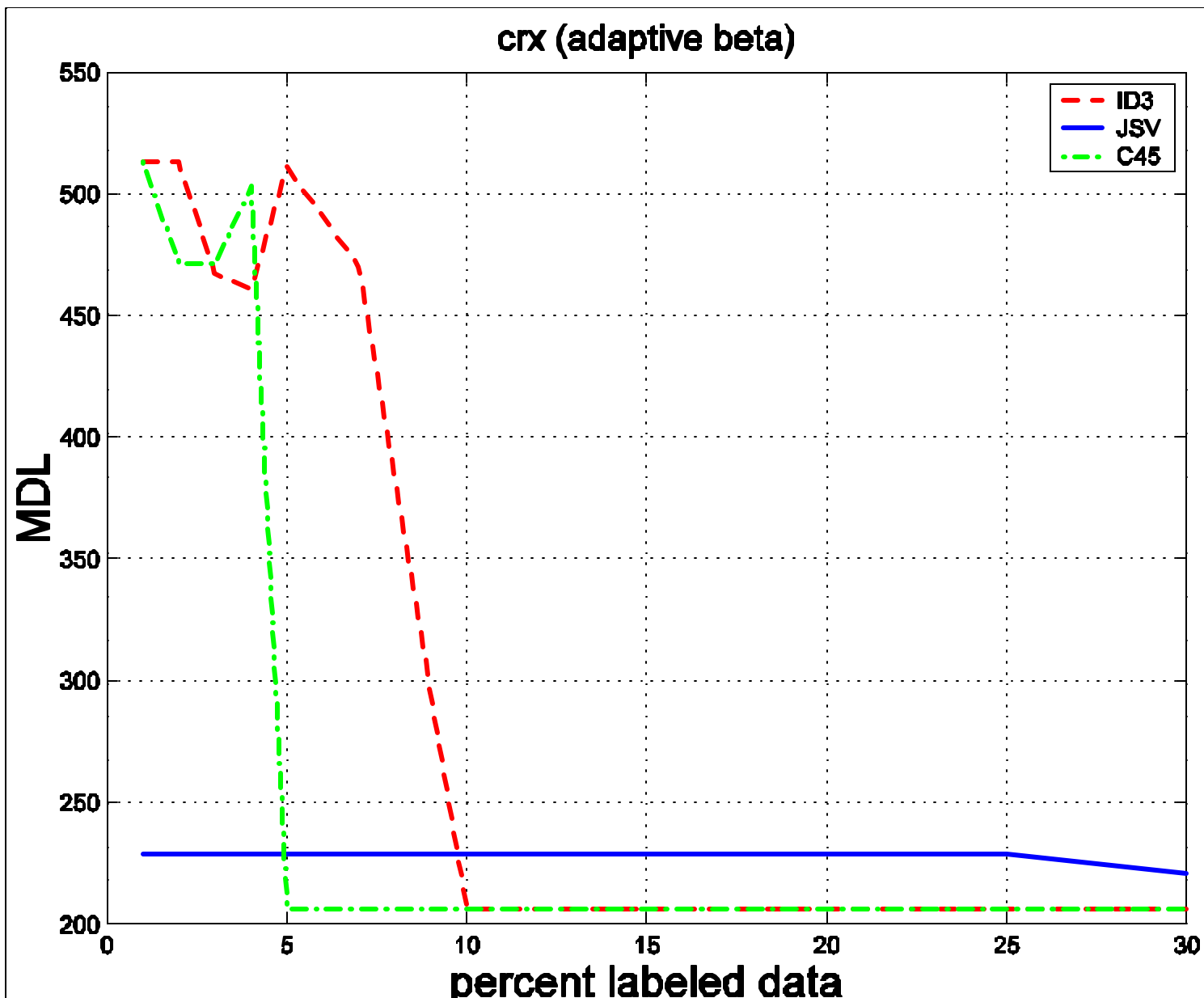


crx/id3_5p   REC=690 ATTR=15  MDL=373.00   cERR=0.22

crx/jsv_5p   REC=690 ATTR=15  BETA=1.0   MDL=213.00   cERR=0.1

# results



crx (adaptive beta)

# results



crx (adaptive beta)

# results

### Table 1. Performance with 5% of labeling

|      | MDL(ID3) | MDL(JSV) | Err(ID3) | Err(JVS) |
|------|----------|----------|----------|----------|
| BAL  | 436      | 422      | 0.30     | 0.28     |
| BAND | 359      | 343      | 0.29     | 0.39     |
| CAR  | 707      | 699      | 0.18     | 0.17     |
| CMC  | 1612     | 1648     | 0.68     | 0.52     |
| CRX  | 381      | 221      | 0.22     | 0.14     |
| MONK | 284      | 284      | 0.32     | 0.32     |
| TIC  | 663      | 595      | 0.33     | 0.25     |
| VOTE | 119      | 47       | 0.12     | 0.04     |

### Table 2. Performance with 10% of labeling

|      | MDL(ID3) | MDL(JSV) | Err(ID3) | Err(JVS) |
|------|----------|----------|----------|----------|
| BAL  | 431      | 422      | 0.30     | 0.29     |
| BAND | 346      | 348      | 0.28     | 0.35     |
| CAR  | 630      | 703      | 0.13     | 0.18     |
| CMC  | 1596     | 1638     | 0.60     | 0.53     |
| CRX  | 206      | 221      | 0.14     | 0.14     |
| MONK | 284      | 284      | 0.32     | 0.32     |
| TIC  | 521      | 589      | 0.29     | 0.24     |
| VOTE | 47       | 47       | 0.04     | 0.04     |

# application: ethereal plugin

- kerf.cs.dartmouth.edu

# Thank You