

Maximum Common Subgraph

Algorithms and Algorithm Portfolios

Paulius Dilkas

School of Computing Science
University of Glasgow

4th March 2018

Outline

- 1 Algorithms
- 2 Algorithm selection
- 3 Labelling
- 4 Features
- 5 Random forests
- 6 Results
- 7 What happens when labelling changes?
- 8 Future work

Maximum Common Subgraph

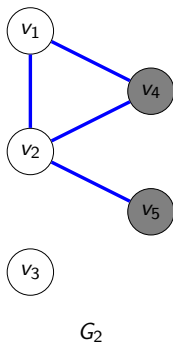
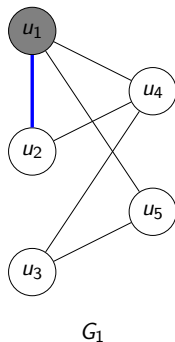
Definition

A *maximum common (induced) subgraph* between graphs G_1 and G_2 is a graph G_3 such that $G_3 = (V_3, E_3)$ is isomorphic to induced subgraphs of both G_1 and G_2 with $|V_3|$ maximised.

Algorithms

- MCSPLIT, MCSPLIT \downarrow
 - (McCreesh, Prosser and Trimble 2017)
- clique encoding
 - (McCreesh, Ndiaye et al. 2016)
- $k \downarrow$
 - (Hoffmann, McCreesh and Reilly 2017)

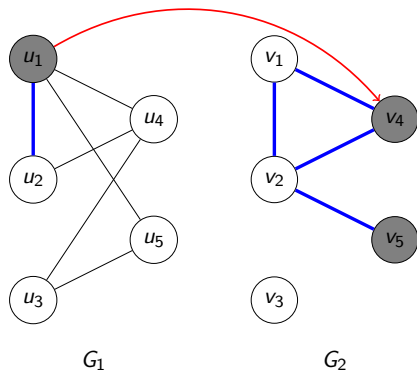
MCSPPLIT: a Branch and Bound Algorithm



Partial solution:
Upper bound: 4

Label	G_1	G_2
0	u_2, u_3, u_4, u_5	v_1, v_2, v_3
1	u_1	v_4, v_5

McSPIT: a Branch and Bound Algorithm



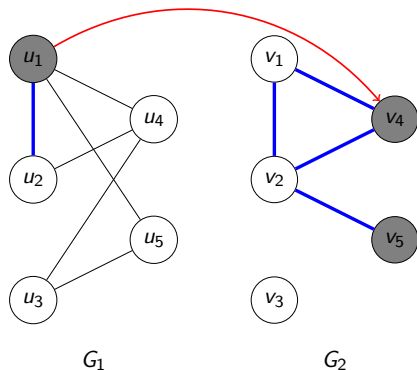
Partial solution:

Upper bound: 4

Label	G_1	G_2
0	u_2, u_3, u_4, u_5	v_1, v_2, v_3
1	u_1	v_4, v_5

Decision: $u_1 \mapsto v_4$

McSPIT: a Branch and Bound Algorithm

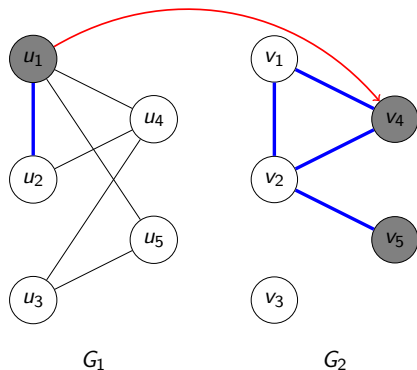


Partial solution:

Upper bound: 4

Label	G_1	G_2
00	u_3	v_3
01	u_4, u_5	\emptyset
02	u_2	v_1, v_2
10	\emptyset	v_5

McSPPLIT: a Branch and Bound Algorithm

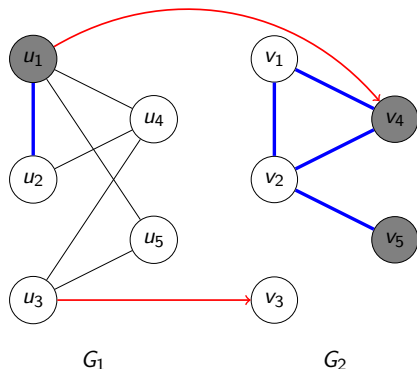


Partial solution: $u_1 \mapsto v_4$

Upper bound: $1 + 2$

Label	G_1	G_2
00	u_3	v_3
01	u_2	v_1, v_2

MCSP_{PLIT}: a Branch and Bound Algorithm



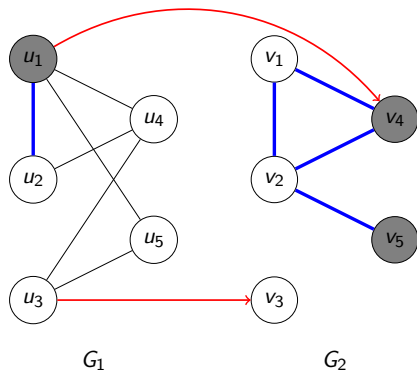
Partial solution: $u_1 \mapsto v_4$

Upper bound: $1 + 2$

Label	G_1	G_2
00	u_3	v_3
01	u_2	v_1, v_2

Decision: $u_3 \mapsto v_3$

MCSPPLIT: a Branch and Bound Algorithm

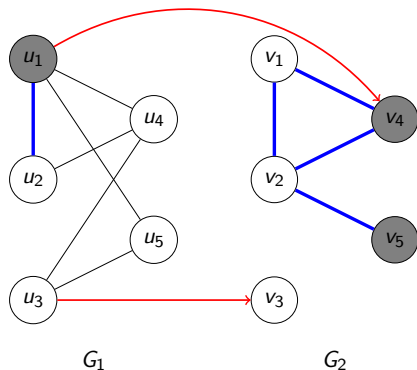


Partial solution: $u_1 \mapsto v_4$

Upper bound: $1 + 2$

Label	G_1	G_2
010	u_2	v_1, v_2
011	u_4, u_5	\emptyset

MCSPPLIT: a Branch and Bound Algorithm

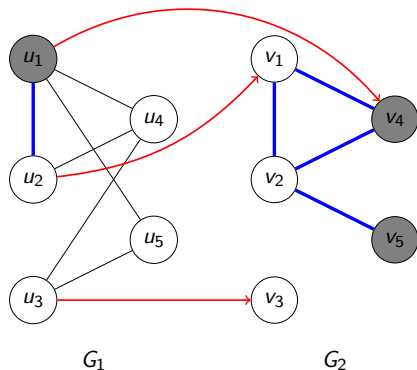


Partial solution: $u_1 \mapsto v_4, u_3 \mapsto v_3$

Upper bound: $2 + 1$

Label	G_1	G_2
010	u_2	v_1, v_2

McSPIT: a Branch and Bound Algorithm



Partial solution: $u_1 \mapsto v_4$, $u_3 \mapsto v_3$

Upper bound: $2 + 1$

Label	G_1	G_2
010	u_2	v_1, v_2

Decision: $u_2 \mapsto v_1$

Found a solution!

Backtrack to confirm optimality

Algorithm selection

Definition (Bischl et al. 2016)

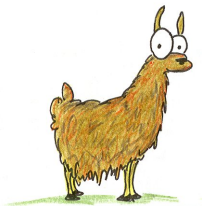
Given a set \mathcal{I} of problem instances, a space of algorithms \mathcal{A} , and a performance measure $m: \mathcal{I} \times \mathcal{A} \rightarrow \mathbb{R}$, the *algorithm selection problem* is to find a mapping $s: \mathcal{I} \rightarrow \mathcal{A}$ that optimises $\mathbb{E}[m(i, s(i))]$.

Algorithm selection

Definition (Bischl et al. 2016)

Given a set \mathcal{I} of problem instances, a space of algorithms \mathcal{A} , and a performance measure $m: \mathcal{I} \times \mathcal{A} \rightarrow \mathbb{R}$, the *algorithm selection problem* is to find a mapping $s: \mathcal{I} \rightarrow \mathcal{A}$ that optimises $\mathbb{E}[m(i, s(i))]$.

LLAMA (Kotthoff 2013)



Labelling

Data from Foggia, Sansone and Vento 2001; Santo et al. 2003 (81400 pairs of graphs)

Labelling

Data from Foggia, Sansone and Vento 2001; Santo et al. 2003 (81400 pairs of graphs)

Definition

A *vertex-labelled graph* is a 3-tuple $G = (V, E, \mu)$, where $\mu: V \rightarrow \{0, \dots, N - 1\}$ is a vertex labelling function, for some $N \in \mathbb{N}$.

Labelling

Data from Foggia, Sansone and Vento 2001; Santo et al. 2003 (81400 pairs of graphs)

Definition

A *vertex-labelled graph* is a 3-tuple $G = (V, E, \mu)$, where $\mu: V \rightarrow \{0, \dots, N-1\}$ is a vertex labelling function, for some $N \in \mathbb{N}$.

Definition

A graph $G = (V, E, \mu)$ is said to have a $p\%$ (*vertex*) *labelling* if

$$N = \max \left\{ 2^n : n \in \mathbb{N}, 2^n < \left\lfloor \frac{p}{100\%} \times |V| \right\rfloor \right\}.$$

Labelling

Definition

A graph $G = (V, E, \mu)$ is said to have a $p\%$ (vertex) labelling if

$$N = \max \left\{ 2^n : n \in \mathbb{N}, 2^n < \left\lfloor \frac{p}{100\%} \times |V| \right\rfloor \right\}.$$

- 5% labelling - 20 vertices per label on average
- 50% labelling - 2 vertices per label on average

Labelling

Definition

A graph $G = (V, E, \mu)$ is said to have a $p\%$ (vertex) labelling if

$$N = \max \left\{ 2^n : n \in \mathbb{N}, 2^n < \left\lfloor \frac{p}{100\%} \times |V| \right\rfloor \right\}.$$

- 5% labelling - 20 vertices per label on average
- 50% labelling - 2 vertices per label on average
- Typical values explored: 33%, 50%, 75%

Labelling

Definition

A graph $G = (V, E, \mu)$ is said to have a $p\%$ (vertex) labelling if

$$N = \max \left\{ 2^n : n \in \mathbb{N}, 2^n < \left\lfloor \frac{p}{100\%} \times |V| \right\rfloor \right\}.$$

- 5% labelling - 20 vertices per label on average
- 50% labelling - 2 vertices per label on average
- Typical values explored: 33%, 50%, 75%
- In my data: 5%, 10%, 15%, 20%, 25%, 33%, 50%

Labelling

Definition

A graph $G = (V, E, \mu)$ is said to have a $p\%$ (*vertex*) *labelling* if

$$N = \max \left\{ 2^n : n \in \mathbb{N}, 2^n < \left\lfloor \frac{p}{100\%} \times |V| \right\rfloor \right\}.$$

- 5% labelling - 20 vertices per label on average
- 50% labelling - 2 vertices per label on average
- Typical values explored: 33%, 50%, 75%
- In my data: 5%, 10%, 15%, 20%, 25%, 33%, 50%
- 3 subproblems
 - no labels
 - vertex labels
 - vertex and edge labels

Features (34 in total)

1–8 are from Kotthoff, McCreesh and Solnon 2016

- ① number of vertices
- ② number of edges
- ③ mean/max degree
- ④ density
- ⑤ mean/max distance between pairs of vertices
- ⑥ number of loops
- ⑦ proportion of vertex pairs with distance $\geq 2, 3, 4$
- ⑧ connectedness

Features (34 in total)

1–8 are from Kotthoff, McCreesh and Solnon 2016

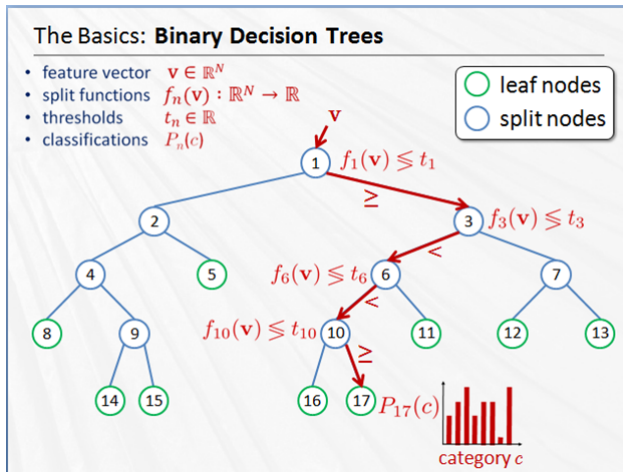
- ① number of vertices
- ② number of edges
- ③ mean/max degree
- ④ density
- ⑤ mean/max distance between pairs of vertices
- ⑥ number of loops
- ⑦ proportion of vertex pairs with distance $\geq 2, 3, 4$
- ⑧ connectedness
- ⑨ standard deviation of degrees
- ⑩ labelling percentage

Features (34 in total)

1–8 are from Kotthoff, McCreesh and Solnon 2016

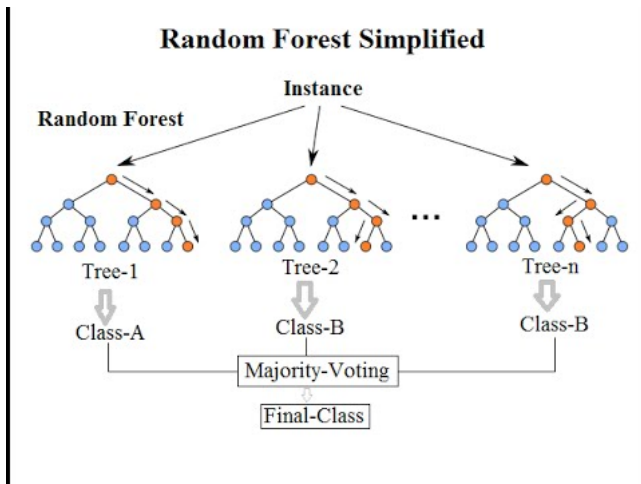
- ① number of vertices
- ② number of edges
- ③ mean/max degree
- ④ density
- ⑤ mean/max distance between pairs of vertices
- ⑥ number of loops
- ⑦ proportion of vertex pairs with distance $\geq 2, 3, 4$
- ⑧ connectedness
- ⑨ standard deviation of degrees
- ⑩ labelling percentage
- ⑪ ratios of features 1–5

Random forests (Breiman 2001)



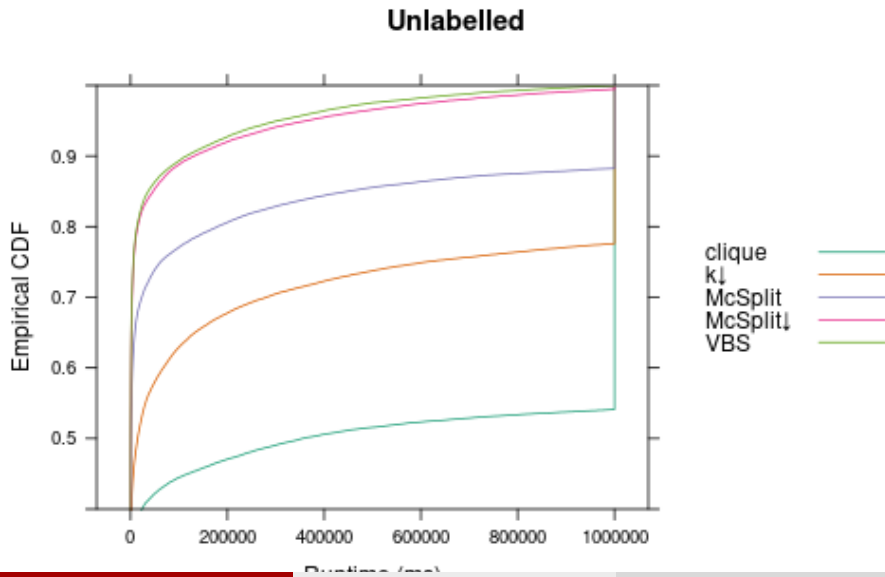
Source: Tae-Kyun Kim & Bjorn Stenger, Intelligent Systems and Networks (ISN) Research Group, Imperial College London

Random forests (Breiman 2001)

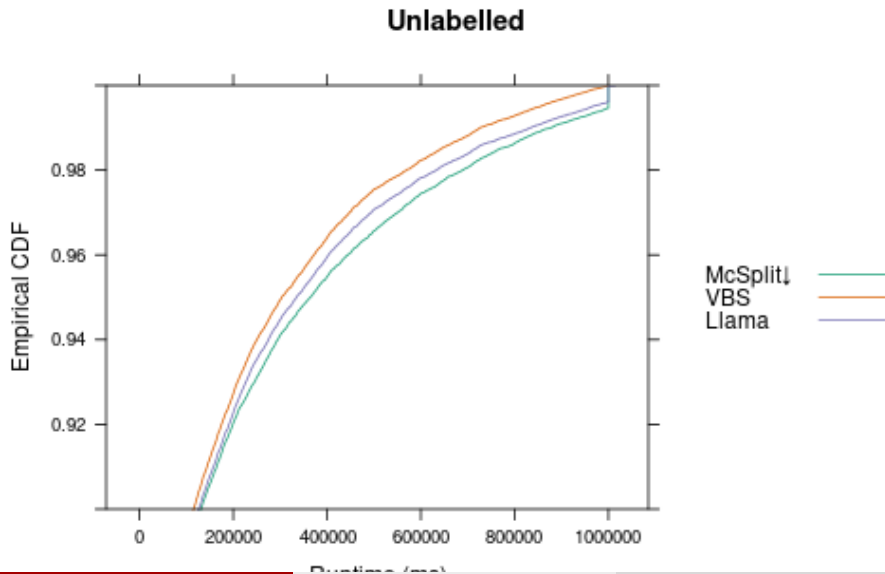


Source: Random Forests(r), Explained, Ilan Reinstein, KDnuggets

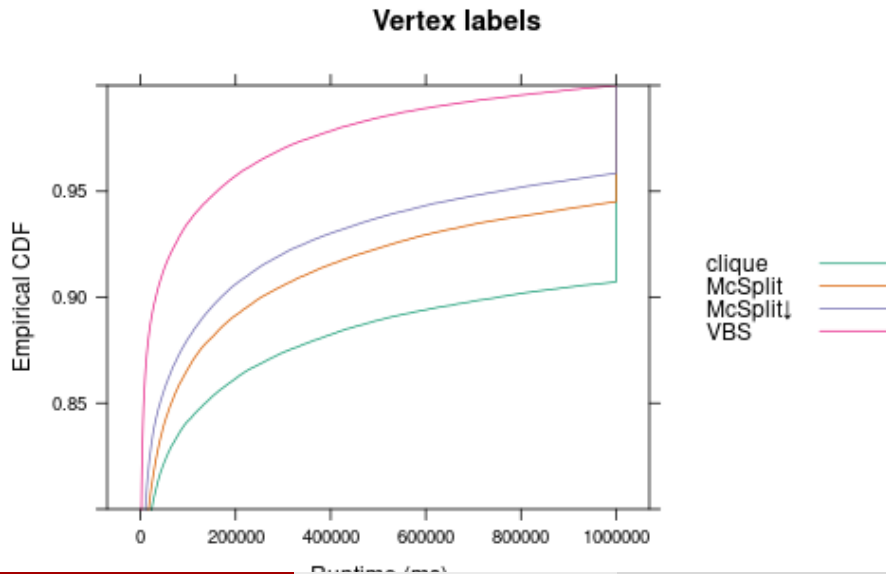
Results



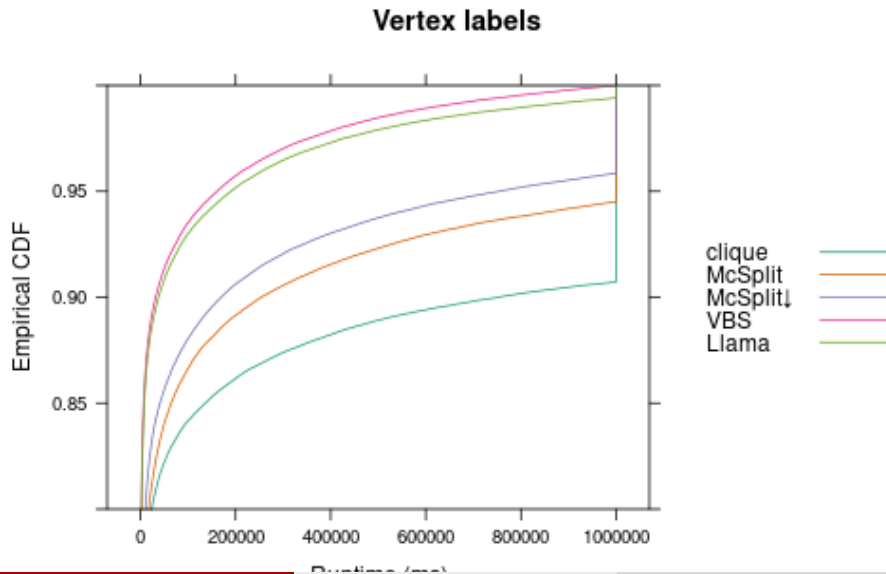
Results (27%)



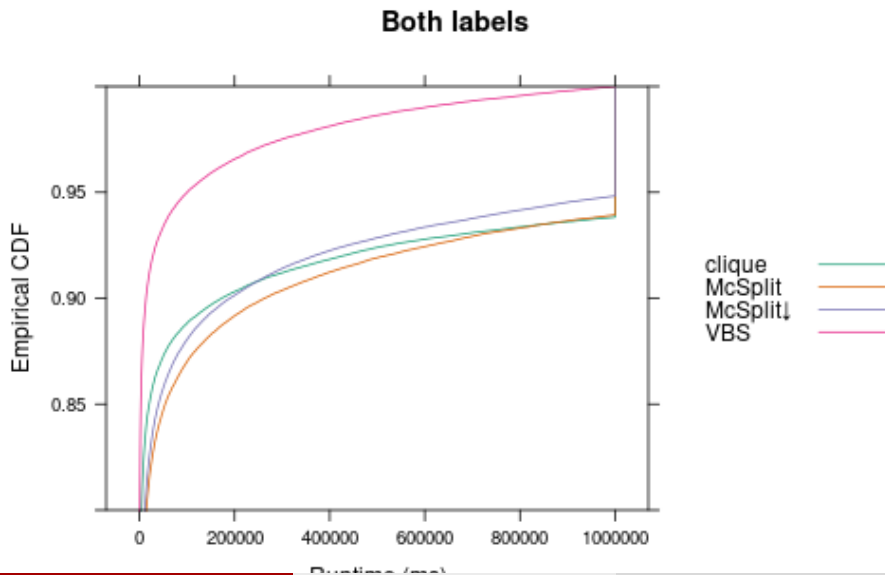
Results



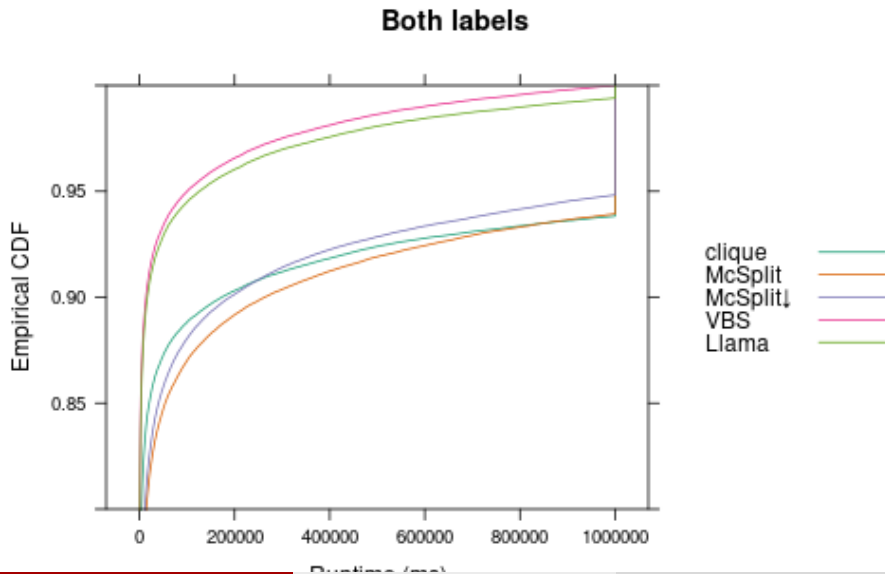
Results (86%)



Results



Results (88%)



Errors

- Out-of-bag error
- For each algorithm
 - $1 - \text{recall}$

Definition

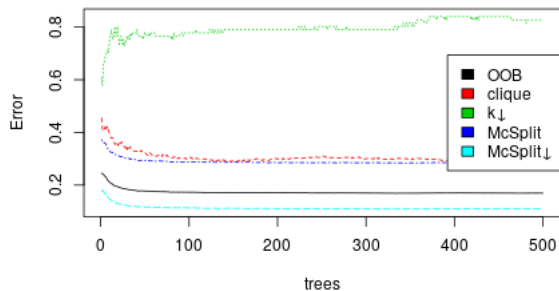
For an algorithm A , *recall* (sensitivity) is

$$\frac{\text{the number of instances that were correctly predicted as } A}{\text{the number of instances where } A \text{ is the correct prediction}}.$$

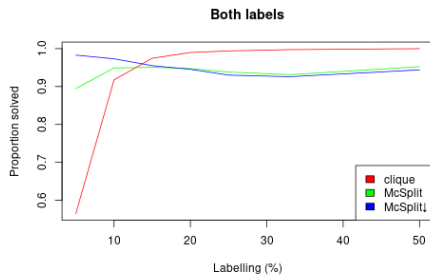
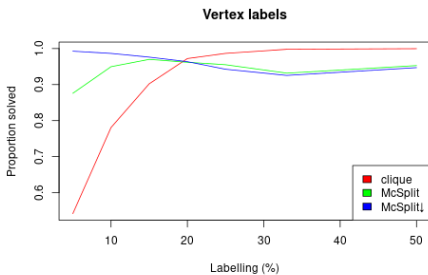
Errors (%)

Error	Labelling		
	no	vertex	both
out-of-bag	17	13	14
clique	30	8	7
McSP _{LIT}	29	22	29
McSP _{LIT} ↓	11	11	11
k ↓	80		

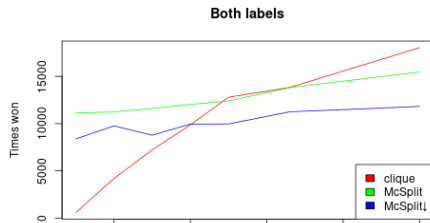
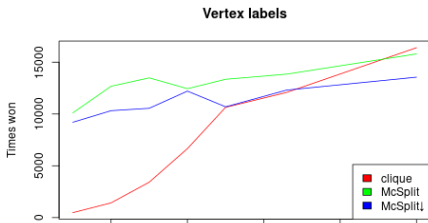
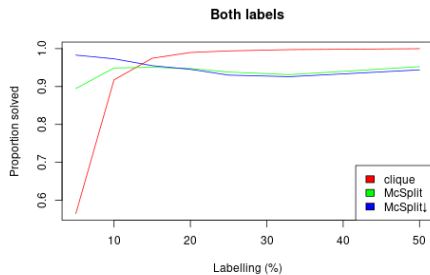
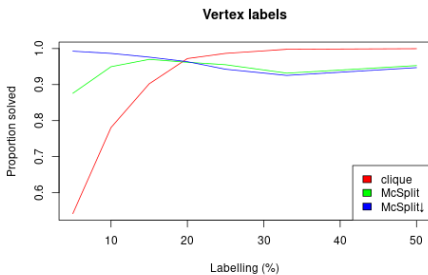
Convergence of errors for unlabelled graphs



What happens when labelling changes?



What happens when labelling changes?



Future work

- Relationships between clique algorithm's performance and properties of the association graph
- How the association graph changes after making a decision
- Can $k \downarrow$ and clique work together?