# Algorithm Selection for Maximum Common Subgraph

Paulius Dilkas

School of Computing Science, University of Glasgow

16th January 2018

# 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 $m\colon \mathcal{I} \times \mathcal{A} \to \mathbb{R}$, the *algorithm selection problem* is to find a mapping $s\colon \mathcal{I} \to \mathcal{A}$ that optimises $\mathbb{E}[m(i, s(i))]$.
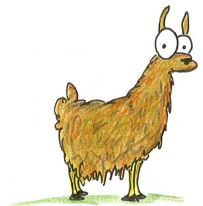
# Algorithm selection

### Definition (Bischl et al. 2016)

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

LLAMA (Kotthoff 2013)

# Algorithms

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

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
$\mu \colon V \to \{0, \ldots, N - 1\}$ 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
$\mu\colon V \to \{0, \dots, N-1\}$ 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)

Algorithm
Selection for
Maximum
Common
Subgraph
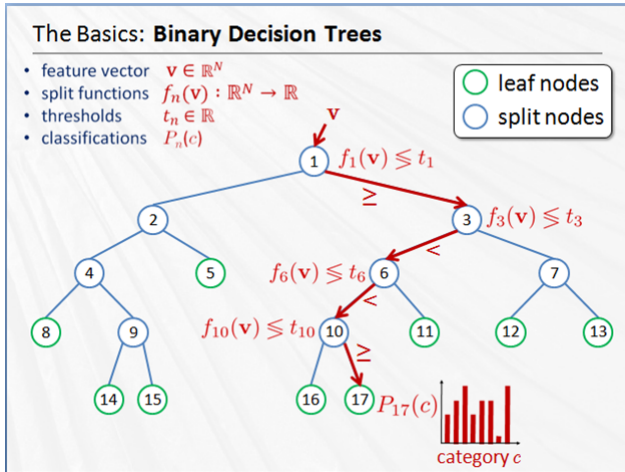
Paulius Dilkas

Algorithm
selection
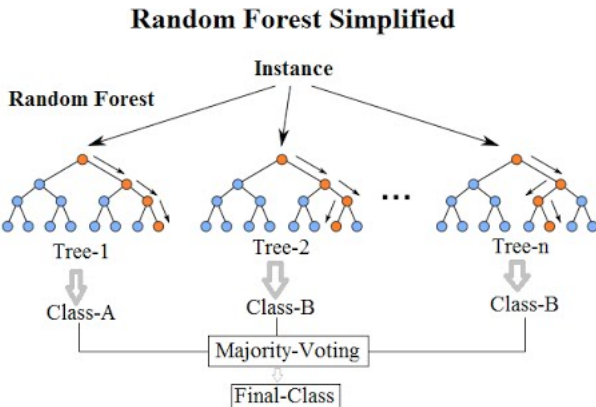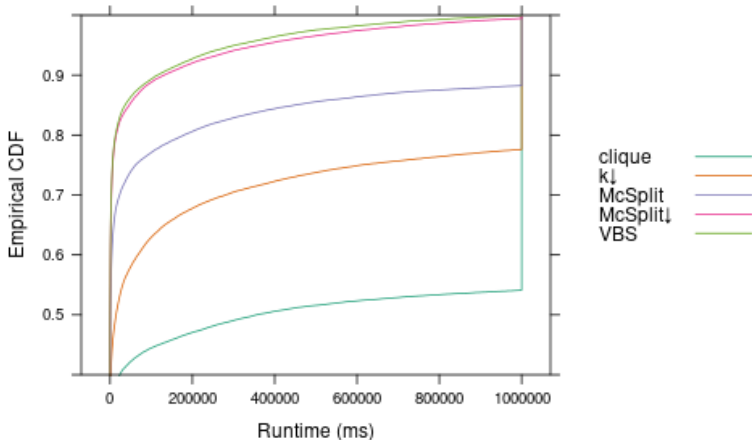
Algorithms

Labelling

Features

Random
forests

Results
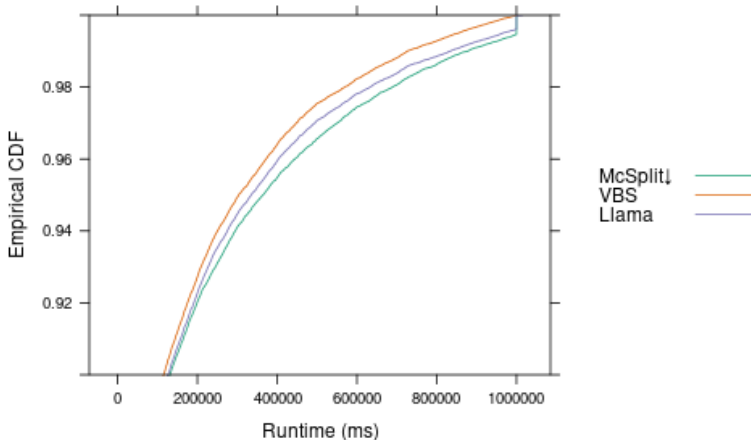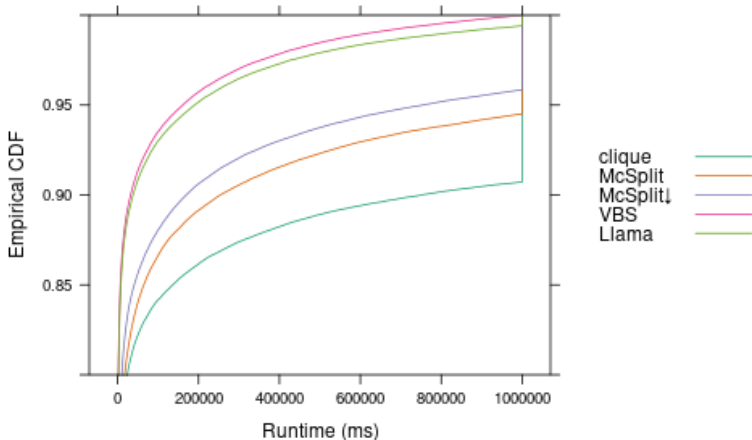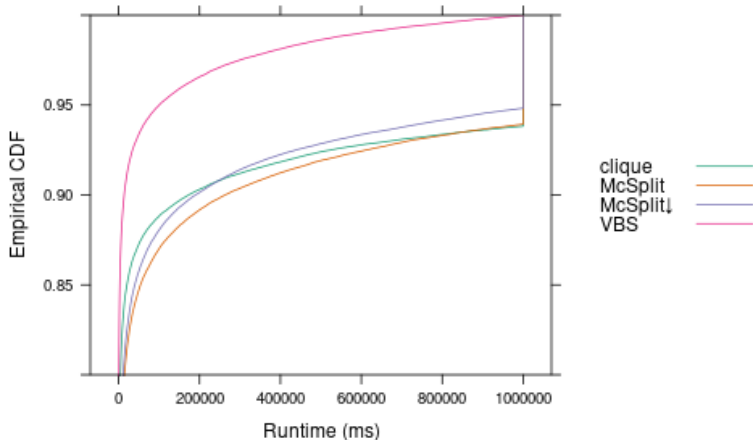
What happens
when labelling
changes?

Future work

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

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

**Vertex labels**

Vertex labels

# Results

**Both labels**

# Results (88%)

# Errors

- Out-of-bag error
- For each algorithm
  - $1 - recall$

### 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 |
| McSplit | 29 | 22 | 29 |
| McSplit $\downarrow$ | 11 | 11 | 11 |
| $k \downarrow$ | 80 | | |

# Convergence of errors for unlabelled graphs

# What happens when labelling changes?

Algorithm
Selection for
Maximum
Common
Subgraph
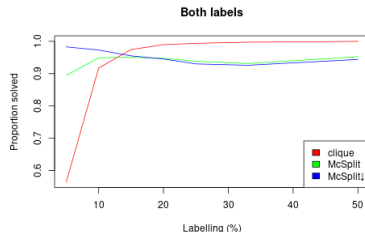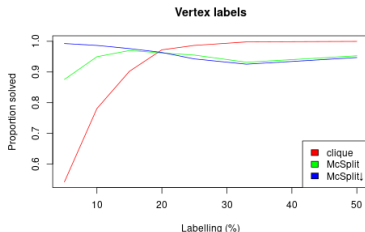
Paulius Dilkas
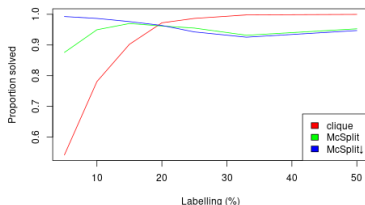
Algorithm
selection
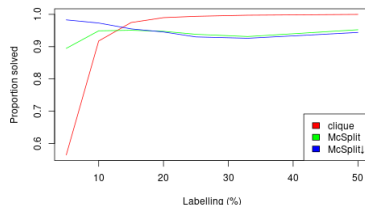
Algorithms

Labelling

Features

Random
forests

Results

What happens
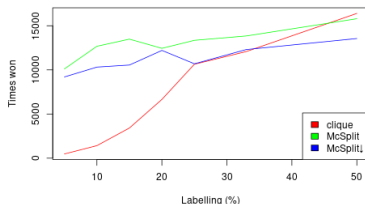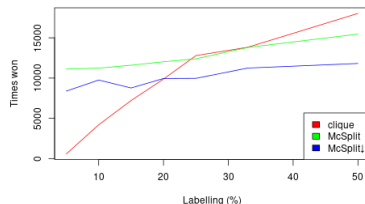when labelling
changes?

Future work

# What happens when labelling changes?

- 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?