

# Machine Learning on Knowledge Bases

---

Marcus Pfeiffer

19th June 2018

Introduction

Data Preprocessing

Machine Learning Algorithms for Premise Selection

Premise Selectors in Detail

Comparison and Outlook

# Introduction

---

# Knowledge Bases and Provers

- Knowledge Bases

# Knowledge Bases and Provers

- Knowledge Bases
  - Mizar Mathematical Library (MML)

# Knowledge Bases and Provers

- Knowledge Bases
  - Mizar Mathematical Library (MML)
  - Library of Isabelle/HOL

# Knowledge Bases and Provers

- Knowledge Bases
  - Mizar Mathematical Library (MML)
  - Library of Isabelle/HOL
  - SUMO
  - Cyc

# Knowledge Bases and Provers

- Knowledge Bases
  - Mizar Mathematical Library (MML)
  - Library of Isabelle/HOL
  - SUMO
  - Cyc
- Automatic Theorem Provers



# Knowledge Bases and Provers

- Knowledge Bases
  - Mizar Mathematical Library (MML)
  - Library of Isabelle/HOL
  - SUMO
  - Cyc
- Automatic Theorem Provers
  - E
  - Vampire
  - CVC4
  - SPASS
  - Z3

# Knowledge Bases and Provers

- Knowledge Bases
  - Mizar Mathematical Library (MML)
  - Library of Isabelle/HOL
  - SUMO
  - Cyc
- Automatic Theorem Provers
  - E
  - Vampire
  - CVC4
  - SPASS
  - Z3
- Integration of Automatic Provers in Interactive Provers

## Definition (Premise Selection problem)

- ATP  $A$
- large number of premises  $\mathcal{P}$
- conjecture  $c$

## Definition (Premise Selection problem)

- ATP  $A$
- large number of premises  $\mathcal{P}$
- conjecture  $c$

Predict those premises from  $\mathcal{P}$  that are likely of use to  $A$  for constructing a proof for  $c$

## Introduction Example

$$\text{map } f \text{ } xs = ys \implies \text{zip}(\text{rev } xs)(\text{rev } ys) = \text{rev}(\text{zip } xs \text{ } ys)$$

## Introduction Example

$$\text{map } f \text{ } xs = ys \implies \text{zip}(\text{rev } xs)(\text{rev } ys) = \text{rev}(\text{zip } xs \text{ } ys)$$

A proof can be found including the following two lemmas:

$$\text{length } xs = \text{length } ys \implies \text{zip}(\text{rev } xs)(\text{rev } ys) = \text{rev}(\text{zip } xs \text{ } ys) \quad (1)$$

$$\text{length}(\text{map } f \text{ } xs) = \text{length } xs \quad (2)$$

## Introduction Example

`used[] = used evs`

## Introduction Example

$$\text{used}[\ ] = \text{used evs}$$

A straightforward proof uses the following lemmas:

$$\text{used}[\ ] = \bigcup_B \text{parts}(\text{initState } B) \quad (3)$$

$$X \in \text{parts}(\text{initState } B) \implies X \in \text{used evs} \quad (4)$$

$$(\bigwedge x. x \in A \implies x \in B) \implies A \subseteq B \quad (5)$$

$$b \in \bigcup_{x \in A} Bx \iff \exists x \in A. b \in Bx \quad (6)$$



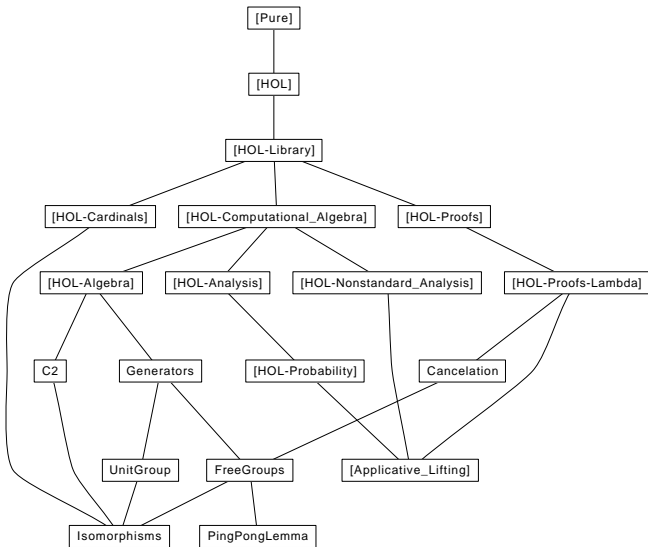
# Data Preprocessing

---

## **Definition (Dependency)**

A definition or theorem  $T$  depends on some definition, lemma or theorem  $T'$ , if  $T'$  is needed for the proof of  $T$

# Theory dependencies graph of Free Groups



## Example of feature extraction

For given depth 2,  $g(hxa)$  with  $x :: \tau$  we get the following structural features:

## Example of feature extraction

For given depth 2,  $g(hx a)$  with  $x :: \tau$  we get the following structural features:

$x$	$a$	$g \_$	$g(h \_ \_)$
$h \_ \_$	$hx \_$	$h \_ a$	$hx a$

## Example of feature extraction

For given depth 2,  $g(hx a)$  with  $x :: \tau$  we get the following structural features:

$x$	$a$	$g$	$g(h \_ \_)$
$h \_ \_$	$hx \_$	$h \_ a$	$hx a$

which are simplified to

$\tau$	$a$	$g$	$g(h)$
$h$	$h(\tau)$	$h(a)$	$h(\tau, a)$

## Another example of feature extraction

$$\text{transpose}(\text{map}(\text{map } f) \text{ xss}) = \text{map}(\text{map } f)(\text{transpose } \text{xss})$$

has the features:

map	map(list_list)	fun
map(fun)	map(map, list list)	list
map(map)	transpose	list_list
map(transpose)	transpose(map)	List
map(map, transpose)	transpose(list_list)	

$(c, \text{usedPremises}(c), F(c))$



# Machine Learning Algorithms for Premise Selection

---

## k-Nearest-Neighbors

"Nearness" of two formulas  $a$  and  $b$  in terms of shared features:

"Nearness" of two formulas  $a$  and  $b$  in terms of shared features:

$$n(a, b) = \sum_{t \in F(a) \cap F(b)} w(t)^{\tau_1} \quad (7)$$

## k-Nearest-Neighbors

"Nearness" of two formulas  $a$  and  $b$  in terms of shared features:

$$n(a, b) = \sum_{t \in F(a) \cap F(b)} w(t)^{\tau_1} \quad (7)$$

$N := k$  nearest neighbors of  $c$

## k-Nearest-Neighbors

"Nearness" of two formulas  $a$  and  $b$  in terms of shared features:

$$n(a, b) = \sum_{t \in F(a) \cap F(b)} w(t)^{\tau_1} \quad (7)$$

$N := k$  nearest neighbors of  $c$

Relevance $_c(p) =$

$$\left( \tau_2 \sum_{q \in N \mid p \in \text{usedPremises}(q)} \frac{n(q, c)}{|\text{usedPremises}(q)|} \right) + \begin{cases} n(p, c) & \text{if } p \in N \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

$$P(\text{p is used in the proof of c}) \quad (9)$$

$$P(\text{p is used in the proof of } c) \quad (9)$$

$$\approx P(\text{p is used to prove } c' \mid c' \text{ has features } F(c)) \quad (10)$$

$$P(\text{p is used in the proof of c}) \quad (9)$$

$$\approx P(\text{p is used to prove c}' | \text{c}' \text{ has features } F(c)) \quad (10)$$

$P(\text{p is used in the proof of c}' )$

- $\prod_{t \in F(c) \cap \bar{F}(p)} P(\text{c}' \text{ has feature } t | \text{p is used in the proof of c}' )$
- $\prod_{t \in F(c) - \bar{F}(p)} P(\text{p has feature } t | \text{p is not used in the proof of c}' )$
- $\prod_{t \in \bar{F}(p) - F(c)} P(\text{c}' \text{ does not have feature } t | t \text{ is used in the proof of c}' )$



- $r(q)$  = number of times a fact  $q$  occurs as a dependency

- $r(q)$  = number of times a fact  $q$  occurs as a dependency
- $s(q, t)$  = number of times a fact  $q$  occurs as a dependency of a fact described by feature  $t$ .

# Naive Bayes

- $r(q)$  = number of times a fact  $q$  occurs as a dependency
- $s(q, t)$  = number of times a fact  $q$  occurs as a dependency of a fact described by feature  $t$ .
- $K$  = total number of known proofs.

# Naive Bayes

- $r(q)$  = number of times a fact  $q$  occurs as a dependency
- $s(q, t)$  = number of times a fact  $q$  occurs as a dependency of a fact described by feature  $t$ .
- $K$  = total number of known proofs.

## Naive Bayes

- $r(q)$  = number of times a fact  $q$  occurs as a dependency
- $s(q, t)$  = number of times a fact  $q$  occurs as a dependency of a fact described by feature  $t$ .
- $K$  = total number of known proofs.

$$P(p \text{ is used in the proof of } c') = \frac{r(p)}{K} \quad (12)$$

# Naive Bayes

- $r(q)$  = number of times a fact  $q$  occurs as a dependency
- $s(q, t)$  = number of times a fact  $q$  occurs as a dependency of a fact described by feature  $t$ .
- $K$  = total number of known proofs.

$$P(p \text{ is used in the proof of } c') = \frac{r(p)}{K} \quad (12)$$

$$P(c' \text{ has feature } t | p \text{ is used in the proof of } c') = \frac{s(p, t)}{r(p)} \quad (13)$$

# Naive Bayes

- $r(q)$  = number of times a fact  $q$  occurs as a dependency
- $s(q, t)$  = number of times a fact  $q$  occurs as a dependency of a fact described with feature  $t$ .
- $K$  = total number of known proofs.

$$\begin{aligned} \text{Relevance}_c(p) = & \sigma_1 \ln(r(p)) + \sum_{t \in F(c') \cap \bar{F}(p)} w(f) \ln \frac{\sigma_2 s(p, t)}{r(p)} \\ & + \sigma_3 \sum_{t \in \bar{F}(p) - F(c)} w(f) \ln \left( 1 - \frac{s(p, t)}{r(t)} \right) + \sigma_4 \sum_{t \in F(c) - \bar{F}(p)} w(f) \end{aligned}$$

## Premise Selectors in Detail

---



For a set  $K$  of known symbols let

$k$  = number of known symbols occurring in the fact

$u$  = number of unknown symbols occurring in the fact

For a set  $K$  of known symbols let

$k$  = number of known symbols occurring in the fact

$u$  = number of unknown symbols occurring in the fact

1. For each fact, compute as its score  $k/(k + u)$
2. Select the facts with the highest score;  
add the features of the selected facts to the set of known symbols  $K$ .

- MaSh: Machine Learning for Sledgehammer with k-NN and Naive Bayes

- MaSh: Machine Learning for Sledgehammer with k-NN and Naive Bayes
- MeSh: Combination of MaSh and a MePo like selector

- MaSh: Machine Learning for Sledgehammer with k-NN and Naive Bayes
- MeSh: Combination of MaSh and a MePo like selector
- In practice: Combined with a proximity selector

- MaSh: Machine Learning for Sledgehammer with k-NN and Naive Bayes
- MeSh: Combination of MaSh and a MePo like selector
- In practice: Combined with a proximity selector

## **Comparison and Outlook**

---

for  $p_1, \dots, p_n$  selected by selection algorithm

Full recall:

$$\min_{k \in N} \text{usedPremises}(c) \subset \{p_1, \dots, p_k\}$$



for  $p_1, \dots, p_n$  selected by selection algorithm

Full recall:

$$\min_{k \in N} \text{usedPremises}(c) \subset \{p_1, \dots, p_k\}$$

k-Coverage

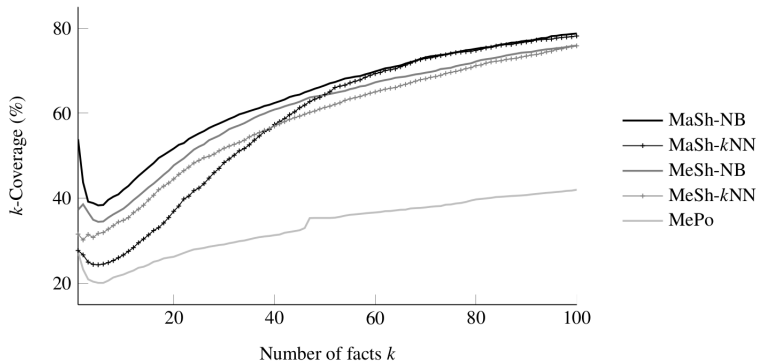
$$\frac{|\{p_1, \dots, p_k\} \cap \text{usedPremises}(c)|}{\min\{k, |\text{usedPremises}(c)|\}}$$

## Comparison in Metric

Formalization	MePo	MaSh		MeSh	
		NB	kNN	NB	kNN
Auth	647	104	143	<b>96</b>	112
IsaFoR	1332	<b>513</b>	604	517	570
Jinja	839	244	306	<b>229</b>	256
List	1083	<b>234</b>	263	259	271
Nominal2	1045	<b>220</b>	276	229	264
Probability	1324	424	422	<b>393</b>	395

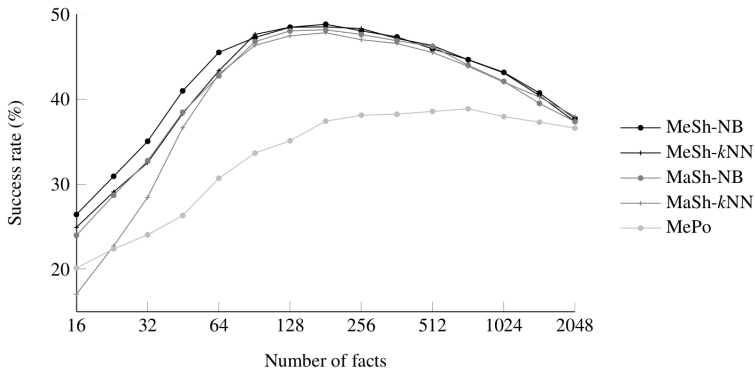
**Figure 1:** Average full recall

# Comparison in Metric



**Figure 2:**  $k$ -Coverage for IsaFoR

# Comparison in Performance



**Figure 3:** Success rates

## Summary and Outlook

- Fact Selection is an important problem in Automatic Theorem Proving

## Summary and Outlook

- Fact Selection is an important problem in Automatic Theorem Proving
- MePo performs well on problems with rare symbols

## Summary and Outlook

- Fact Selection is an important problem in Automatic Theorem Proving
- MePo performs well on problems with rare symbols
- Fine grained dependency analysis and good feature extraction enable the successful application of machine learning algorithms

## Summary and Outlook

- Fact Selection is an important problem in Automatic Theorem Proving
- MePo performs well on problems with rare symbols
- Fine grained dependency analysis and good feature extraction enable the successful application of machine learning algorithms
- MaSh and MeSh outperform MePo in terms of the chosen metrics



## Summary and Outlook

- Fact Selection is an important problem in Automatic Theorem Proving
- MePo performs well on problems with rare symbols
- Fine grained dependency analysis and good feature extraction enable the successful application of machine learning algorithms
- MaSh and MeSh outperform MePo in terms of the chosen metrics
- Similar results for SNoW and MOR compared to Vampire-SInE for MML

## Summary and Outlook

- Fact Selection is an important problem in Automatic Theorem Proving
- MePo performs well on problems with rare symbols
- Fine grained dependency analysis and good feature extraction enable the successful application of machine learning algorithms
- MaSh and MeSh outperform MePo in terms of the chosen metrics
- Similar results for SNoW and MOR compared to Vampire-SInE for MML
- Possibilities of applying other machine learning algorithms, e.g. Neural Networks