

Utilizing Machine Learning to Accelerate Automated Assignment of Backbone NMR Data

Background

NMR
Machine
Learning

Algorithm

Overview
Model Training
Preprocessing
The Search

Results

Outlook

Joel Venzke^{1,2}, David Mascharka¹, Paxten Johnson^{1,2},
Rachel Davis¹, Katherine Roth¹, Leah Robison¹, Timothy
Urness¹ and Adina Kilpatrick²

¹Department of Mathematics and Computer Science

²Department of Physics and Astronomy
Drake University

joel.venzke@drake.edu

April 16, 2015

Overview

① Background

NMR

Machine Learning

② Algorithm

Overview

Model Training

Preprocessing

The Search

③ Results

④ Outlook

Harvard University Conference

NMR

Assignment
with Machine
Learning

J. Venzke
D. Mascharka
P. Johnson
R. Davis
K. Roth
L. Robison
T. Urness
A. Kilpatrick

Background

NMR

Machine
Learning

- 1898

Algorithm

Overview
Model Training
Preprocessing
The Search

- Say some things

Results

Outlook

Harvard University Conference

NMR

Assignment
with Machine
Learning

J. Venzke
D. Mascharka
P. Johnson
R. Davis
K. Roth
L. Robison
T. Urness
A. Kilpatrick

Background

NMR
Machine
Learning

Algorithm

Overview
Model Training
Preprocessing
The Search

Results

Outlook

- 1898
- Say some things

Algorithmic Overview

Model Training

- Performed once during algorithm development
- Provides model used in Preprocessing

Preprocessing

- Imports NMR data set
- Filters NMR data using machine learning model

The Search

- Uses results from Preprocessing
- Assigns NMR data set
- Records results

Model Training

Training Data Set

Biological Magnetic Resonance Bank (BMRB)

- 9,736 datasets containing chemical shifts for the C_α and C_β resonances of 689,977 residues
- Removing outliers leaves 681,363 pairs of C_α and C_α
 - 3 standard deviations from the mean
 - Avoids over-fitting
 - Improves algorithmic performance

Training the Model

Preformed Once

- Time consuming task
- Trained once, used many times

Models Trained

- DecisionTable, j.48, LMT

Reading Data

Protein Sequence

- Read in as letters
- Converted to BMRB average values
- Used for comparison in the search

NMR Data Set

- Read in C_α , C_β for Residue i and $i - 1$
- Stored in Tile

Tile

Residue $i - 1$ C_α, C_β **Residue i** C_α, C_β **Confidence Levels** $P_1, P_2, \dots, P_{19}, P_{20}$

Confidence Level Calculation

Tile

Residue $i - 1$

C_α, C_β

Residue i

C_α, C_β

Confidence Levels

$P_1, P_2, \dots, P_{19}, P_{20}$

Machine Learning

Background

NMR

Machine
Learning

Algorithm

Overview

Model Training

Preprocessing

The Search

Results

Outlook

Missing Data

Blank Tile Creation

- Compare length of protein sequence to NMR Data set
- Blank tiles are created to make up the gap

Proline

- Lacks H-N spin system
- Does not produce C_{α} , C_{β} values
- Protein sequence is examined
- Special flags are set

Blank Tile

Residue $i - 1$

- , -

Residue i

- , -

Confidence Levels

1.0, 1.0, \dots , 1.0, 1.0

Proline

yes/no

Harvard University Conference

NMR

Assignment
with Machine
Learning

J. Venzke

D. Mascharka

P. Johnson

R. Davis

K. Roth

L. Robison

T. Urness

A. Kilpatrick

Background

NMR

Machine
Learning

- 1898

Algorithm

Overview

Model Training

Preprocessing

The Search

- Say some things

Results

Outlook

Machine Learning Algorithms

NMR

Assignment
with Machine
Learning

J. Venzke
D. Mascharka
P. Johnson
R. Davis
K. Roth
L. Robison
T. Urness
A. Kilpatrick

Background

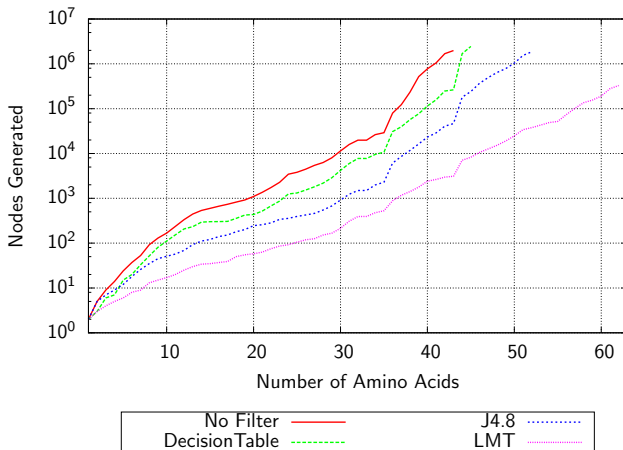
NMR
Machine
Learning

Algorithm

Overview
Model Training
Preprocessing
The Search

Results

Outlook



NMR

Assignment with Machine Learning

J. Venzke
D. Mascharka
P. Johnson
R. Davis
K. Roth
L. Robison
T. Urness
A. Kilpatrick

Background

NMR
Machine
Learning

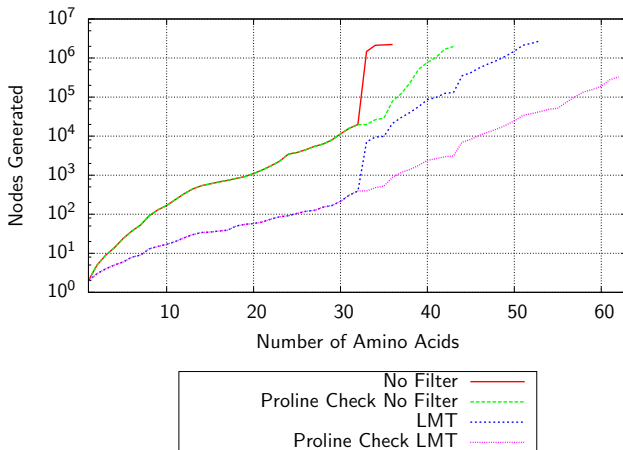
Algorithm

Overview
Model Training
Preprocessing
The Search

Results

Outlook

Proline Checking



Future Research

Extend the Proline checking to other amino acids

Include a hysteric for assignment cost prediction

Assign subsets and combine to generate full assignments

NMR

Assignment
with Machine
Learning

J. Venzke

D. Mascharka

P. Johnson

R. Davis

K. Roth

L. Robison

T. Urness

A. Kilpatrick

Thank You

Background

NMR

Machine
Learning

Algorithm

Overview

Model Training

Preprocessing

The Search

Results

Outlook