Machine learning



Jacob Kautzky Group Meeting 12-5-18

# Machine learning outline



What is machine learning?

**Common machine learning algorithms** 

Machine learning applied to chemistry

Retrosynthesis

Reaction discovery

**Reaction optimization** 

Drug optimization

Materials chemistry

# What is machine learning?

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with the experience E.



# Data gathering

#### Collect the data yourself

need to gather a sufficient quantity of data

High throughput

Automation

#### Get the data from other sources

In the past (and to lesser extents today), it was difficult to convert data from published papers and patents into data that could be utilized for machine learning

Datasets can have untrustworthy data and a way needs to be devised to prevent them from overly influencing the conclusions

Published papers generally omit data that does not work, which can cause problems in trying to build complete datasets

# Supervised vs unsupervised learning

Supervised learning — working with a set of labeled training data where each piece of data has an input and output object



Unsupervised learning (data mining) — the algorithm finds hidden patterns in a load of data

# Algorithm types: linear and logistic

Most basic forms for regression and classification respectively

# Algorithm types: linear and logistic

#### Most basic forms for regression and classification respectively

#### Linear - fitting a polynomial curve to data

- Easily interpretable results
- Fast to train a machine
- Low average prediction accuracy
- Easily skewed by outliers

# Algorithm types: linear and logistic

#### Most basic forms for regression and classification respectively

#### Linear - fitting a polynomial curve to data

- Easily interpretable results
- Fast to train a machine
- Low average prediction accuracy
- Easily skewed by outliers

Logistic - measures the relationship between the categorical dependent variable and one or more independent variables by using a logistic function



#### **Binary logistic curve**

Similar to a linear model, however the outputs are restricted between 0 and 1 and there is a different conditional basis set (Bernoulli vs Gaussian)

Similar properties to linear, but slighly more difficult for the user to interpret the results

# Algorithm types: k nearest neighbor

"Lazy learning" method that is dependent on the local structure of the data

For classification, the object is classified by a majority vote of its k closest neighbors

For regression, the value is the average of its k closest neighbors



### Algorithm types: naive bayes classifier



Despite their apparent oversimplification, naive Bayes classifiers have worked quite well in several real world applications

The advantage of a Bayes model is that it requires a small amount of training data necessary to attain the parameters required for classification, however it suffers from slightly lower average predictive accuracy in several instances

### Algorithm types: decision tree



# Algorithm types: random forests



The reasoning behind the algorithm is more difficult for someone to understand

Can be skewed toward certain descriptors if there are more variables associated with a given descriptor

Can be avoided by using the cforest algorithm which can avoid descriptor selection bias

# Algorithm types: artificial neural network

Computing system inspired by biological neural networks

A framework for many machine learning algorithms to work together and process coplex data imputs

Based on a collection of connected nodes called artificial neurons

An artificial neuron can transmit a signal to other neurons attached to it by an edge that in turn can signal other neurons



Examples of this include learning to identify images that contain cats by being shown images that have been manually labeled cat or no cat

This has been used by groups in the past to convert chemical structures in patents and old papers into SMILES and other readable formats for machine learning

Challenging to understand what the algorithm is doing

Deals with irrelevant features well

Slow training speed

High average predictive accuracy

# Machine learning outline



What is machine learning?

**Common machine learning algorithms** 

### Machine learning applied to chemistry

Retrosynthesis

Reaction discovery

**Reaction optimization** 

Drug optimization

Materials chemistry

## Machine learning

Retrosynthesis

- Reaction optimization
- Identifying ideal conditions
  - Predicting yields

Chemistry

- Medicinal Chemistry
- Structure generation
- Binding prediction
- Estimating ADMET properties

#### Materials chemistry

- Identifying ideal alloys
- Synthesis/ Crystalization

#### Reaction discovery

# History of machine learning and synthesis planning



Provided a list of heuristic conditions to guide the choice of synthetic disconnections

Short lived and eventually split in two

# Organic Chemical Simulation of Synthesis (OCSS)



Corey, E. J.; Wipke, W. T. Science 1969, 166, 178.

# History of machine learning and synthesis planning



Grzybowski, B. A. et al Angew. Cem. Int. Ed. 2016, 55, 5904.

# Six principles of Logic and Heuristics Applied to Synthetic Analysis (LHASA)

**Transformation-based strategy** – Identification of a powerful simplifying transformation (ex. Diels- Alder, Aldol cyclization, etc.)

**Mechanistic transforms** — The target is converted to a reactive intermediate from which other intermediates of synthetic value can be generated



Corey, E. J.; Long, A. K.; Rubenstein, S. D.. Science 1985, 228, 408.

# Six principles of Logic and Heuristics Applied to Synthetic Analysis (LHASA)

Structure-goal — Identification of a potential starting material, building block, retron-containing subunit or initiating chiral element



**Topological strategies** — Identification of one or more bonds that could lead to major simplifications

Stereochemical strategies - Stereoselective reactions, or steric based arguments used to reduce stereocomplexity

**Functional group-oriented strategies** — Functional group interconversions and determining logical disconnections based off of functional group arrangement(s)

Corey, E. J.; Long, A. K.; Rubenstein, S. D.. Science 1985, 228, 408.

### History of machine learning and synthesis planning



Grzybowski, B. A. et al Angew. Cem. Int. Ed. 2016, 55, 5904.

### Hendrickson's SYNGEN

#### Focus on skeletal construction

**Bondset** – a set of bonds  $\lambda$  which need to be constructed

The number of ways to break apart a skeleton is equal to the number of possible bondsets



*b*! / (*b*- $\lambda$ )! possible bond sets for *b* bonds

### Hendrickson's SYNGEN



Hendrickson, J. B. Angew. Chem. Int. Ed. 1990, 29, 1286.

# History of machine learning and synthesis planning



# History of machine learning and synthesis planning



Grzybowski, B. A. et al Angew. Cem. Int. Ed. 2016, 55, 5904.

The purpose of a Monte Carlo tree search is to given a game state, choose the most promising next move

The purpose of a Monte Carlo tree search is to given a game state, choose the most promising next move



The purpose of a Monte Carlo tree search is to given a game state, choose the most promising next move



The purpose of a Monte Carlo tree search is to given a game state, choose the most promising next move



Expand the tree by giving it further options off of its current leaf

Segler, M. H. S.; Preuss, M.; Waller, M. P. Nature 2018, 555, 604.

The purpose of a Monte Carlo tree search is to given a game state, choose the most promising next move



Segler, M. H. S.; Preuss, M.; Waller, M. P. Nature 2018, 555, 604.

Merged a Monte Carlo tree search (MCTS) with three separate neural networks to devise retrosyntheses

- Utilized all reactions on Reaxy's prior to 2015 to generate a training set
- Extracted 301,671 rules from this dataset and utilized them to train neural networks

Neural networks learn the context in which reactions can occur (functional group tolerance)

- Developed three neural networks
  - 1. Expansion policy guides the search in promising directions by proposing a restricted number of transformations
  - 2. Examine the feasibility of the proposed transformation
  - 3. Estimates the position value



Merged a Monte Carlo tree search (MCTS) with three separate neural networks to devise retrosyntheses



Segler, M. H. S.; Preuss, M.; Waller, M. P. Nature 2018, 555, 604.

#### **Evaluating Performance**

- The neural network reaction checker has a false positive rate of 1.5% and false negative of 14%
- The neural network predicted the correct solution 31% of the time and had the correct solution in the top 50 73% of the time



Segler, M. H. S.; Preuss, M.; Waller, M. P. Nature 2018, 555, 604.

### Chematica

Unites network theory, modern high-power computing, artificial intelligence, and expert chemical knowledge to rapidly design synthetic pathways

Thousands of reactions hand coded by experts



Grzybowski, B. A. et al. Chem 2018, 4, 522.

# Chematica

Unites network theory, modern high-power computing, artificial intelligence, and expert chemical knowledge to rapidly design synthetic pathways

Thousands of reactions hand coded by experts



Pruned down options by removing branches with unlikely structural motifs, that proceed through strained intermediates, etc.

A scoring algorithm evaluates the resultant substrates produced and the reactions required to make the set

Simplified graphic example of the synthetic possibilities

### Chematica

Unites network theory, modern high-power computing, artificial intelligence, and expert chemical knowledge to rapidly design synthetic pathways

Thousands of reactions hand coded by experts

The software was tested on several compounds of commercial interest to MilliporeSigma that currently had troublesome syntheses

The software designed routes that were then executed without changes except for straightforward adjustments to the reaction conditions (e.g. temperature, solvent, etc.)

In every instance improvements were made (shorter routes, fewer chromatographic steps, higher yields, more reproducible)



Chematica was also able to design a pathway that broke a patented route to a compound while more than doubling the yield!

Grzybowski, B. A. et al. Chem 2018, 4, 522.

# Machine learning outline



What is machine learning?

**Common machine learning algorithms** 

Machine learning applied to chemistry

Retrosynthesis

**Reaction optimization** 

Reaction discovery

Drug optimization

Materials chemistry



#### Why are isoxazoles challenging in the Buchwald Hartwig reaction?



Ahneman, D. T.; Estrada, J. G.; Lin, S.; Dreher, S. D.; Doyle, A. G. Science 2018, 360, 186.



Ahneman, D. T.; Estrada, J. G.; Lin, S.; Dreher, S. D.; Doyle, A. G. Science 2018, 360, 186.



Examining training set data size

Ahneman, D. T.; Estrada, J. G.; Lin, S.; Dreher, S. D.; Doyle, A. G. Science 2018, 360, 186.



The observation of oxidative addition appeared to track well with the important reaction parameters

Ahneman, D. T.; Estrada, J. G.; Lin, S.; Dreher, S. D.; Doyle, A. G. Science 2018, 360, 186.

Can machine learning explore chemical space and predict areas of reactivity?

Examined different combinations of two- and three- component reactions from a pool of 18 starting materials using an automated robot

**Automated Robot** 

- Executes six experiments in parallel
- Analyzes reactions by <sup>1</sup>H NMR, MS, and IR



The reactions are classified as reactive or unreactive by a supported vector machine (SVM) with a linear kernel

This algorithm compared the IR and NMR spectra of the product to that of the starting material and registered differences as reactivity hits

The SVM was trained on a set of 72 reactive and non reactive mixtures and could classify mixtures with an accuracy of 86%



SVM results were used to train a linear discriminant analysis (LDA) model that could construct a model of chemical space

Granda, J. M.; Donina, L.; Dragone V.; Long, D.-L.; Cronin, L. Nature 2018, 559, 377.



Granda, J. M.; Donina, L.; Dragone V.; Long, D.-L.; Cronin, L. Nature 2018, 559, 377.



By surveying 10% of the reaction space, the model could accurately predict the reactivity of 80% of the chemical space

Granda, J. M.; Donina, L.; Dragone V.; Long, D.-L.; Cronin, L. Nature 2018, 559, 377.

manual examination of the reactive combinations identified by the machine lead to the discovery of four novel transformations



Granda, J. M.; Donina, L.; Dragone V.; Long, D.-L.; Cronin, L. Nature 2018, 559, 377.

#### This system could also be utilized for the prediction of yields





Granda, J. M.; Donina, L.; Dragone V.; Long, D.-L.; Cronin, L. Nature 2018, 559, 377.





#### The SVM model was converted to a decision tree to make it chemically interpretable



This decision tree could be utilized to make chemical hypotheses



## Machine learning

Retrosynthesis

- Reaction optimization
- Identifying ideal conditions
  - Predicting yields

Chemistry

- Medicinal Chemistry
- Structure generation
- Binding prediction
- Estimating ADMET properties

#### Materials chemistry

- Identifying ideal alloys
- Synthesis/ Crystalization

#### Reaction discovery

# Questions?

