# Machine Learning for Metabolite Identification

Huibin Shen
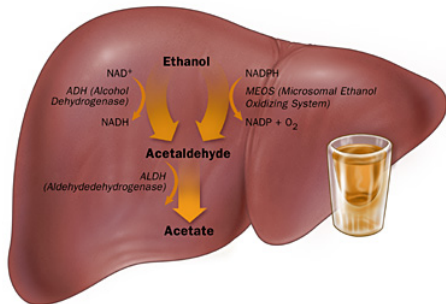
Department of Computer Science
School of Science, Aalto University

March 30, 2017

# Content

- Metabolite identification
- Machine learning
- Machine learning for metabolite identification
- Conclusion

**Aalto University**
School of Science
and Technology

**Huibin Shen**
March 30, 2017
2/21

# What is metabolite?

▶ Metabolism is the set of life-sustaining chemical transformations within the cells of of living organisms.



▶ Metabolites are the intermediates and products of metabolism.

Aalto University
School of Science
and Technology

Huibin Shen
March 30, 2017
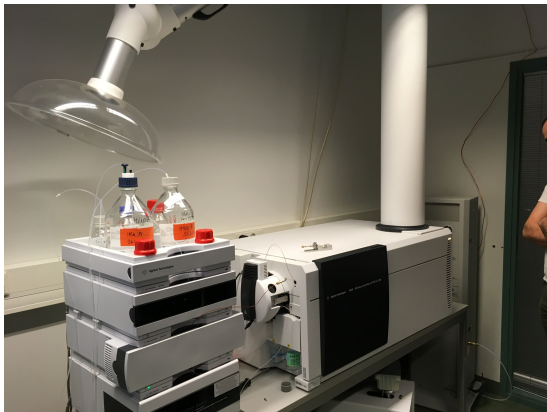3/21

# Why to identify the metabolites?

Because they relate to many things:

- ▶ health & nutrition
- ▶ pharmaceuticals
- ▶ biotechnology
- ▶ regulatory affairs (drug trafficking, anti-doping).

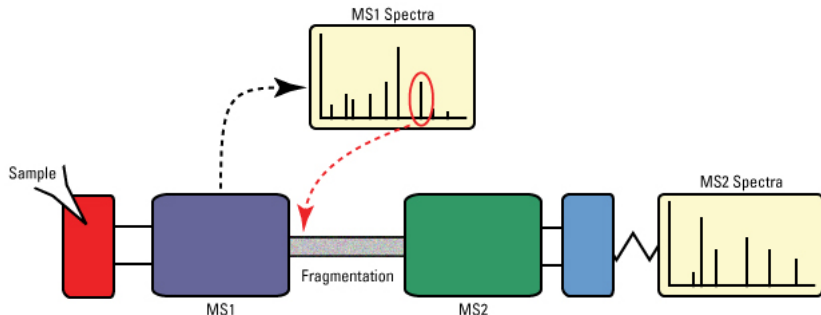It is a pre-requisite step for many subsequent analyses.

**Aalto University**
School of Science
and Technology

**Huibin Shen**
March 30, 2017
4/21

# How to identify the metabolites?

The main technology is the tandem mass spectrometry.

**Aalto University**
School of Science
and Technology

**Huibin Shen**
March 30, 2017
5/21

# How to identify the metabolites?
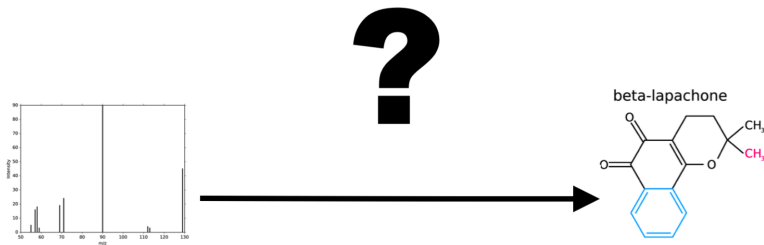
The main technology is the tandem mass spectrometry.



The peaks in the end characterize the structure of metabolite.

**Aalto University**
School of Science
and Technology

**Huibin Shen**
March 30, 2017
6/21

# Metabolite identification

Given a tandem MS/MS spectrum, what is the molecular structure?



This is the core problem this dissertation trying to improve.

Aalto University
School of Science
and Technology

Huibin Shen
March 30, 2017
7/21

# Content

- Metabolite identification
- Machine learning
- Machine learning for metabolite identification
- Conclusion

**Aalto University**
**School of Science**
**and Technology**

**Huibin Shen**
**March 30, 2017**
**8/21**

# Machine learning

**Aalto University**
School of Science
and Technology

**Huibin Shen**
March 30, 2017
9/21

# Machine learning



classification

kernel approximation

NOT WORKING

SVC Ensemble Classifiers

KNeighbors Classifier

SGD Classifier

Naïve Bayes

Text Data

Linear SVC

<100K samples

scikit-learn algorithm cheat-sheet

get more data

>50 samples

START

regression

SGD Regressor

ElasticNet Lasso

SVR(kernel='rbf') EnsembleRegressors

predicting a category

<100K samples

few features should be important

RidgeRegression SVR (kernel='linear')

labeled data

do you have

predicting a quantity

clustering

Spectral Clustering GMM

KMeans

number of categories known

<10K samples

MiniBatch KMeans

MeanShift VBGMM

<10K samples

just looking

Randomized PCA

Isomap Spectral Embedding

LLE

<10K samples

kernel approximation

dimensionality reduction

tough luck

predicting structure

# Machine learning

- ▶ We use supervised learning, where labels are available.
- ▶ We focus on a special family of the supervised learning, where non-tabular format of data can be handled.



- ▶ They are known as kernel methods.

**Aalto University**
School of Science
and Technology

**Huibin Shen**
March 30, 2017
11/21

# Kernel methods

- Data as $D = \{(x_i, y_i)\}_{i=1}^{n}$, $x \in \mathcal{X}$ and $y \in \{+1, -1\}$.
- A kernel function $k(x, x')$ measures similarity.
- For example, support vector machine (SVM) solves:

$$\max_{\alpha} \ \sum_{i=1}^{n} \alpha_i - \sum_{i,j=1}^{n} \alpha_i \alpha_j y_i y_j \boxed{k(x_i, x_j)}$$

$$s.t. \ \ 0 \leq \alpha_i \leq C, i = 1, \ldots, n,$$

$$\sum_{i=1}^{n} \alpha_i y_i = 0.$$

- We only need to plug in the similarities!

**Aalto University**
School of Science
and Technology

**Huibin Shen**
March 30, 2017
12/21

# Content

- Metabolite identification
- Machine learning
- Machine learning for metabolite identification
- Conclusion

**Aalto University**
**School of Science**
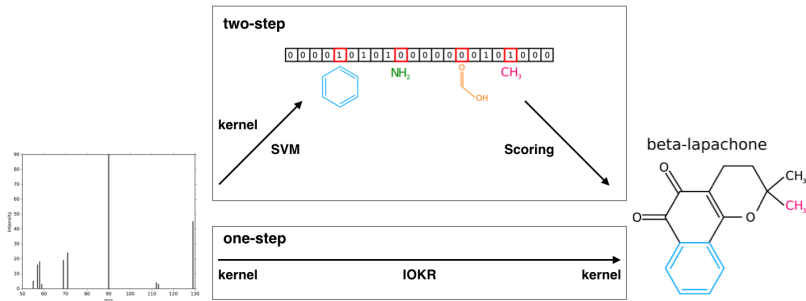**and Technology**

**Huibin Shen**
**March 30, 2017**
**13/21**

# Molecular fingerprint

Molecular fingerprints are a representation of the molecules, encoding structures or other properties.

Aalto University
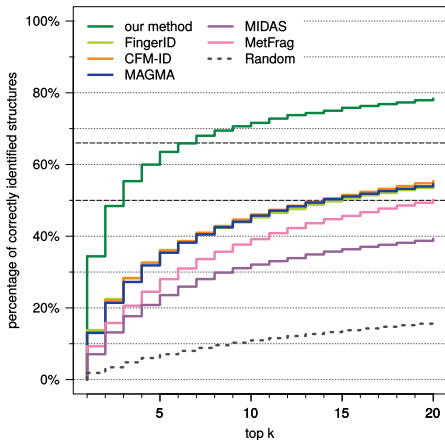School of Science
and Technology

# The proposed methods

Two step approach CSI-FingerID [Heinonen et al., 2012, Shen et al., 2013, Shen et al., 2014, Dührkop et al., 2015].



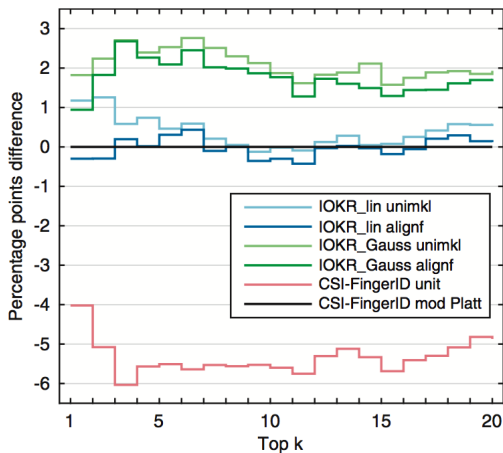One-step approach IOKR [Brouard et al., 2016].

Aalto University
School of Science
and Technology

Huibin Shen
March 30, 2017
15/21

# Results of the two-step approach CSI-FingerID

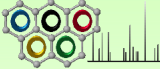Cross validation on 5923 mass spectra and compounds pairs.

Aalto University
School of Science
and Technology

Huibin Shen
March 30, 2017
16/21

# Results of the one-step apparoch IOKR

# Results of the one-step apparoch IOKR

Table: Running time evaluation

|  | Training time | Test time |
|---|---|---|
| CSI:FingerID | 82 h 28 min 23 s | 1 h 11 min 31 s |
| IOKR linear | **42 s** | **1 min 15 s** |
| IOKR polynomial | **38 s** | 21 min 58 s |
| IOKR Gaussian | **41 s** | 33 min 15 s |

**Aalto University**
School of Science
and Technology

**Huibin Shen**
March 30, 2017
18/21

# Results on the CASMI 2016 challenge

Aalto University
School of Science
and Technology

# Conclusion

▶ Metabolite identification is a major bottleneck in computational metabolomics.

▶ Kernel based machine learning for metabolite identification is the new state of the art.

▶ There are many possibilities to improve in the future.

**Aalto University**
School of Science
and Technology

**Huibin Shen**
March 30, 2017
20/21

# Reference

Brouard, C., Shen, H., Dührkop, K., d'Alché Buc, F., Böcker, S., and Rousu, J. (2016).
Fast metabolite identification with input output kernel regression.
*Bioinformatics*, 32(12):i28–i36.

Dührkop, K., Shen, H., Meusel, M., Rousu, J., and Böcker, S. (2015).
Searching molecular structure databases with tandem mass spectra using csi:
Fingerid.
*Proceedings of the National Academy of Sciences*, 112(41):12580–12585.

Heinonen, M., Shen, H., Zamboni, N., and Rousu, J. (2012).
Metabolite identification and molecular fingerprint prediction through machine
learning.
*Bioinformatics*, 28(18):2333–2341.

Shen, H., Dührkop, K., Böcker, S., and Rousu, J. (2014).
Metabolite identification through multiple kernel learning on fragmentation trees.
*Bioinformatics*, 30(12):i157–i164.

Shen, H., Zamboni, N., Heinonen, M., and Rousu, J. (2013).
Metabolite identification through machine learning—tackling casmi challenge using
fingerid.
*Metabolites*, 3(2):484–505.

**Aalto University**
School of Science
and Technology

**Huibin Shen**
March 30, 2017
21/21