# Machine Learning Analysis of Mass Spectrometry Data in the Life Sciences

William Stafford Noble, University of Washington, Seattle

Arzu Tugce Guler, Amsterdam University Medical Center, Amsterdam

Claire O'Donovan, European Bioinformatics Institute, Cambridge

## Agenda

- ❖ Overview of machine learning techniques in mass spectrometry data analysis
- ❖ Panel:
  - *Machine-learning for the proteomic masses: learning peptide properties and clustering spectra,* **Lukas Käll, Science for Life Laboratory, Stockholm**
  - *Revolutionizing MS-based proteomics using machine learning techniques: fragmentation prediction, relative peptide intensity prediction, missing value imputation.* **Lars Jensen, University of Copenhagen**
  - *Bayesian learning and MS big data,* **Sam Payne, Brigham Young University, Salt Lake City**
  - *Machine learning applications for real-time analysis,* **Devin Schweppe, Harvard University**
- ❖ Discussion

## What is a learning behaviour?

Given a task *T*, a performance criterion *C*, and experience *E*, a system **learns** from *E* if it becomes better at solving task *T*, as measured by criterion *C*, by exploiting the information in *E*. *

We need to know:

➢ which task is going to be performed
➢ how the performance on the task is measured
➢ what kind of information is used by the system

T. Mitchell. *Machine Learning.* McGraw-Hill 1996

## Machine Learning in Mass Spectrometry Analysis



Csordas et al. *Database* 2012

Perez-Riverol et al. *Nucleic Acids Research* 2012

## Machine Learning in Mass Spectrometry Analysis

Data is complex, noisy, non-trivial



1. Doana. *Metabolic Phenotyping in Personalized and Public Healthcare* 2016
2. https://phosphopedia.gs.washington.edu/PhosphoproteomicsAssay/index.xhtml
3. Saeed et. al *IEEE/ACM Trans Comput Biol Bionform* 2018

## Retention time prediction

Input: peptide sequence

Output: chromatographic retention time

Utility: boost statistical power to detect peptides

*Proc. Natl. Acad. Sci. USA*
Vol. 77, No. 3, pp. 1632–1636, March 1980
Medical Sciences

**Prediction of peptide retention times in high-pressure liquid chromatography on the basis of amino acid composition**
(lipophilicity/separation techniques)

JAMES L. MEEK

Laboratory of Preclinical Pharmacology, National Institute of Mental Health, Saint Elizabeth Hospital, Washington, D.C. 20032

*Communicated by Bruce Merrifield, December 17, 1979*
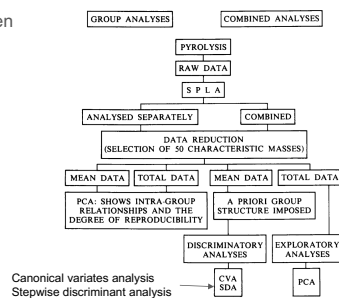
## Species identification

Input: MS1 or MS2 data from a given sample

Output: (list of) detected species

Utility:



Shute et al. *Microbiology* 1984

## Classifying peptide-spectrum matches

Input: a vector of features associated with a peptide-spectrum match

Output: Is this peptide responsible for generating this spectrum?

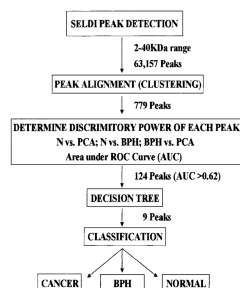Utility: boost statistical power to detect peptides / proteins



Keller et al. *Analytical Chemistry* 2002

## Phenotype / disease classification

Input: MS1 or MS2 data from an individual

Output: predicted phenotype

Utility: early diagnosis or disease prognosis



Adam et al. *Cancer Research* 2002

## Clustering of mass spectra

Input: large collection of mass spectra

Output: cluster assignments

Utility:

- boost statistical power to detect peptides
- speed up database search



Tabb et al. *Analytical Chemistry* 2003

Guthals et al. *Molecular BioSystems* 2012

## SciLifeLab

# MS and ML

Lukas Käll
Royal Institute of Technology - KTH
School of Biotechnology
Stockholm, Sweden

KTH
VETENSKAP
OCH KONST

http://percolator.ms
http://kaell.org

---

# Background on machine learning for MS-based proteomics

---

## Identification by comparing observations to predictions



---

## MS-related separation techniques and their predictors



---

# Mining Repositories

- Clustering offers a nice way to condense information from prior experiments
  Frank et al 2008;Griss et al

- Yade yade

---

## 1 case

# Quandenser:
## Combining quantification and clustering

Input: Raw-files from LFQ experiments
Output: proteins with *Pr(diff. exp.)*
Binaries: https://github.com/statisticalbiotechnology/quandenser
Paper: The & Käll Nature Communications (in press)

Dr. Matthew The

MS1 features can often be retrieved across samples…



… but MS2 spectra are sparse



## Flipping the pipeline

*Identification-first*



## Flipping the pipeline

*Identification-first*        *Quantification-first*



## Benefits of quantification-first

- No need to rediscover the same peptides for each run
- Lowering the number of spectra
  - Enables more advanced identification strategies
  - Faster identification
  - Fewer hypothesis tests
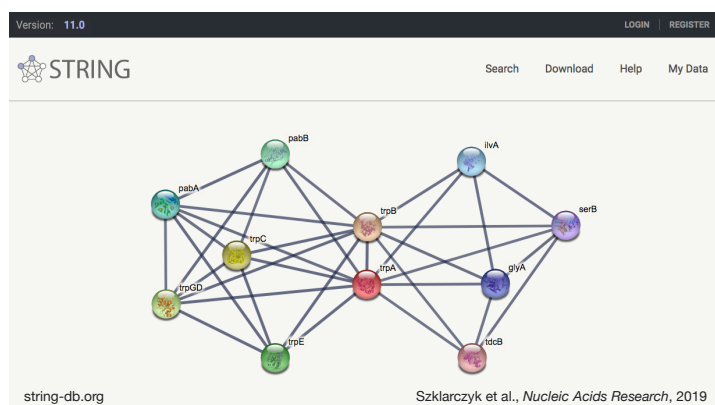
Focus on the spectra that matters!

# Dramatic increase in number of differentially quantified proteins at 5% FDR
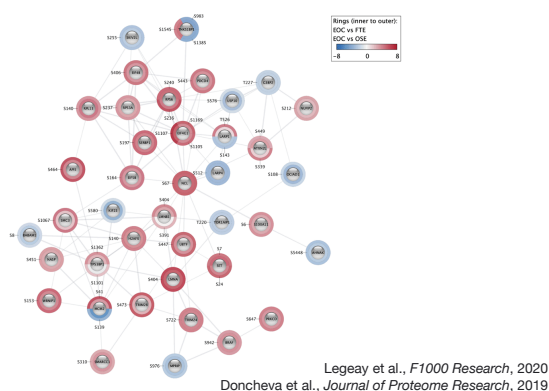
# Three problems to tackle with machine learning

Lars Juhl Jensen
jensenlab.org

network biology



data visualization

string-db.org
Szklarczyk et al., *Nucleic Acids Research*, 2019

**1**

**peptide fragmentation**

Legeay et al., *F1000 Research*, 2020
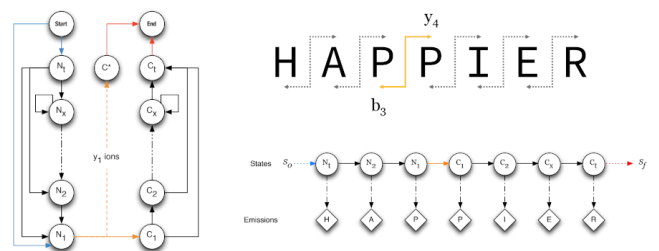Doncheva et al., *Journal of Proteome Research*, 2019

fragment ion spectrum

theoretical spectra
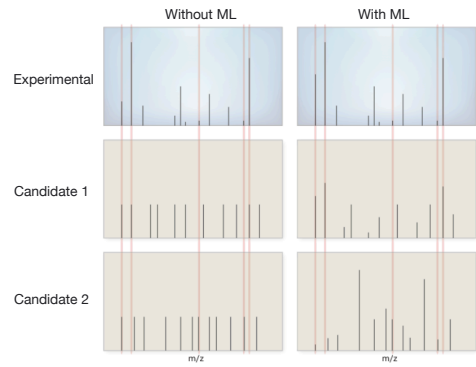
count matches

assume all are equal

fragmentation is predictable



Kirik, Refsgaard & Jensen, *Journal of Proteome Research*, 2019

better identification



Kirik, Refsgaard & Jensen, *Journal of Proteome Research*, 2019

better PTM localization

**2** **peptide abundance**

one protein

multiple unique peptides

equal abundance

different intensities

plenty of training data

peptides + MS parameters

relative intensities

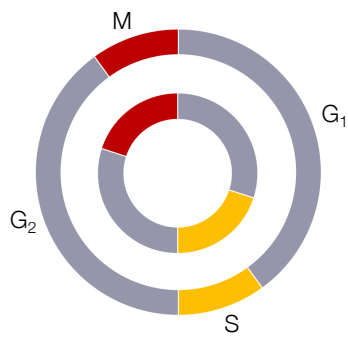better protein quantification

cross-sample comparison

phosphoproteomics

kinase motif analysis

regulation of CDK activity

different growth rates

a few common effects

affect numerous peptides

in a consistent manner

many unrelated experiments

dimensionality reduction

auto encoder

learn effect signatures
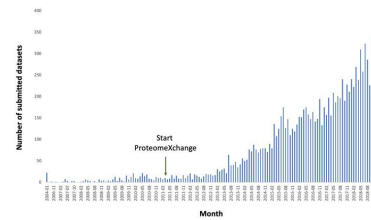
quantify their impact

residual signals

**questions**
?

# Bayesian Learning
# and MS Big Data

Samuel Payne
Brigham Young University
ASMS 2020 - Bioinformatics Interest Group

## Data, Data Everywhere. Are we learning?

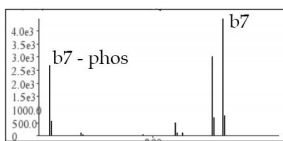How can we learn from public mass spectrometry data?



"We have seen this spectrum a lot of times. Why do we always pretend to not know anything about it?"

--- Mike MacCoss

## Learning from our data

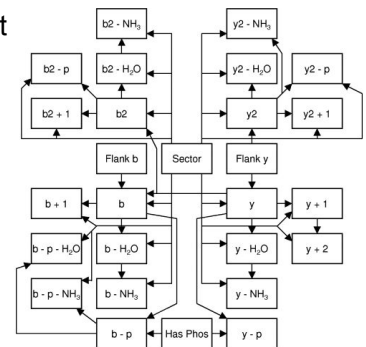What can we learn about expected intensity from data mining?



| B ion | B-PO$_4$ ion | |
|---|---|---|
| strong | strong | 22.5% |
| strong | medium | 20.4% |
| strong | weak | 3.8% |
| strong | absent | 53.3% |
| medium | strong | 8.1% |
| medium | medium | 25.1% |
| medium | weak | 10.8% |
| medium | absent | 55.9% |
| weak | strong | 3.1% |
| weak | medium | 12.9% |
| weak | weak | 14.5% |
| weak | absent | 69.5% |
| absent | strong | 5.4% |
| absent | medium | 15.9% |
| absent | weak | 9.3% |
| absent | absent | 69.4% |

Payne et al. 2008
https://doi.org/10.1021/pr800129m

## Building your Bayes Net

Identify mutual information

More nodes is not always better



## Training

**Inspect** (deprecated) trained on 170,000 phospho-peptide spectra (2008) to learn the probabilities (and joint probabilities) of fragment ion intensity

**MSGF+** trained on 2.8 million spectra (2014) to learn the probabilities of fragment ion intensities

## New Areas for Bayesian Scoring in MS id

"We have seen this spectrum a lot of times. Why do we always pretend to not know anything about it?"
  --- Mike MacCoss

Spectral Library Matching (DIA or metabolomics)
     - are my peak relative intensities as expected?
     - how often does this peak have interference?

# Lipids and Bayes Learning

Learn the fragment ions
Learn their intensity, relative intensity



Kyle et al. 2017
https://doi.org/10.1093/bioinformatics/btx046
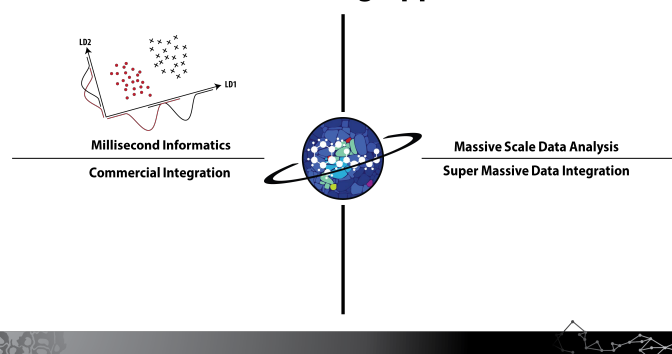
# New Areas for Bayesian Scoring

Protein Inference =? Bayesian Inference

- Do peptide intensities within a protein have a reliable relationship?

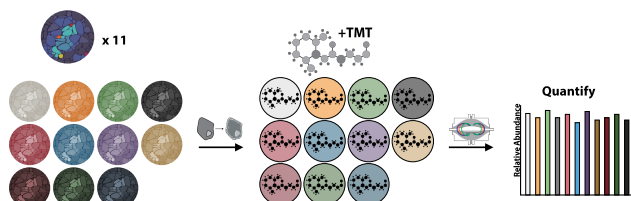- Do we need to learn for each tissue or cell line?

# Live Learning
## ML applications for real-time analysis

Devin K Schweppe, PhD
ML Analysis of MS Data in the Life Sciences Workshop
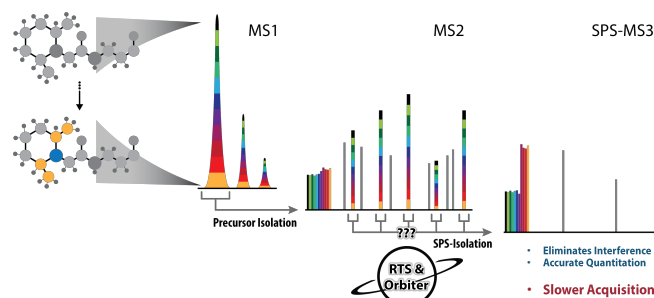ASMS 2020 Reboot
Thurs, June 4th

---

## Machine Learning Applications

LD2
LD1

Millisecond Informatics
Commercial Integration

Massive Scale Data Analysis
Super Massive Data Integration

---

## Basis: Multiplexed proteomics
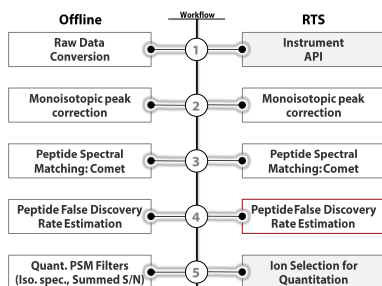
x 11

+TMT

Quantify

Relative Abundance

**Advantages**
- Throughput – one run, multiple samples
- Fewer missing values – quantitation for every sample
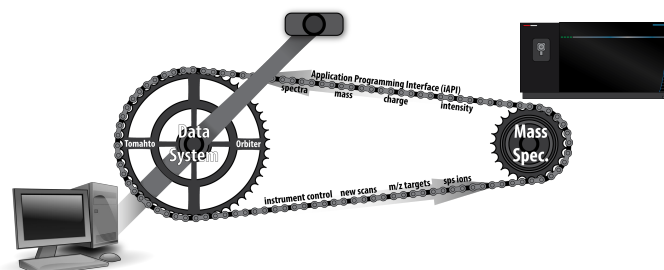- Complex experimental design – Doses, Mutants, Time

---

## Slowed down by SPS-MS³

MS1          MS2          SPS-MS3

Precursor Isolation

???

SPS-Isolation

RTS & Orbiter

- Eliminates Interference
- Accurate Quantitation
- **Slower Acquisition**

---

## Millisecond Informatics: Orbiter

| Offline | Workflow | RTS |
|---|---|---|
| Raw Data Conversion | 1 | Instrument API |
| Monoisotopic peak correction | 2 | Monoisotopic peak correction |
| Peptide Spectral Matching: Comet | 3 | Peptide Spectral Matching: Comet |
| Peptide False Discovery Rate Estimation | 4 | Peptide False Discovery Rate Estimation |
| Quant. PSM Filters (Iso. spec., Summed S/N) | 5 | Ion Selection for Quantitation |

---

## iAPI: Cycling through Real-time data

Application Programming Interface (iAPI)
spectra    mass    charge    intensity

Data System
Tomahto    Orbiter

Mass Spec.

instrument control    new scans    m/z targets    sps ions

## Monoisotopic Peak Correction

Monocle (C#)

Relative AA Abundance

Averagine Distribution

Senko et al. *JASMS* 1995

Monocle

Instrument

$m_{obs} \pm X$ scans

Intensity weighted m/z dist.

Isolated m/z

Observed

Correction: Monoisotopic peak Probable charge = 3

Monocle m/z
Expect 13.756 ppm 1.14

Instrument Assigned m/z
Expect 4.421 ppm 22.60

GLVCGSELGNDELNPER
GLVCGSELGNDELNPER

Yeast LabelFree
HeLa LabelFree
HeLa TMT10

+43% +81%
+43% +82%
+52% +98%

Total Valid PSMs (x1000)

No Correction  Monocle  RawConverter

Rad et al *(in preparation)*

---

## Comet: Open source search engine

**Comet: An open-source MS/MS sequence database search tool**

Jimmy K. Eng[1], Tahmina A. Jahan[1] and Michael R. Hoopmann[2]

[1] Department of Genome Sciences, University of Washington, Seattle, WA, USA
[2] Institute for Systems Biology, Seattle, WA, USA

FASTA

Header
Mass Index
Digested peptides
Protein List

Comet Search → XCorr

Real-time Comet XCorr

Offline Comet XCorr

Yeast (F/R) | VM = Ox (M)
Human (F/R) | VM = Ox (M)

Relative Database size

Single Scan Search Time (ms)

Human
VM = Ox (M)

Yeast
VM = Ox (M)

Real-time Comet scoring with indexed databases shows consistent XCorr results compared to offline search.

---

## Millisecond Informatics: Orbiter

Workflow    RTS

1. Instrument API
2. Monoisotopic peak correction
3. Peptide Spectral Matching: Comet
4. Peptide False Discovery Rate Estimation
5. Ion Selection for Quantitation

Correction: Monoisotopic peak Probable charge = 1

Comet

---

## Towards real-time FDR filtering

Target | Decoy

Count

XCorr

deltaCorr

Target | Decoy

Combining peptide spectral match scores (e.g. XCorr and deltaCorr) aids separation of target and decoy populations.

---

## Linear Discriminant Analysis (LDA)

XCorr    charge    pep. length
deltaCorr    frxn ions matched
|ppm error|    missed cleavage

LDA

LD2

LD1

XCorr > 1
Decoy
Target

Frequency

LDA Score

LDA enables fast target-decoy classification and FDR estimation

---

## Accumulate Training Data

1000 - 2000 PSZ

LD1 Coefficients

Coefficient Value

# of PSMs for LDA Training

deltaXCorr
Ion Frxn

Charge

XCorr
Peptide Length

Missed Cleavages

|Adj| ppm error|

XCorr/deltaCorr/PPM

LDA + FDR

1000-2000 total PSMs collected
At least 10 decoy PSMs*

*Huttlin et al 2007

## Improved sensitivity & quantitative accuracy
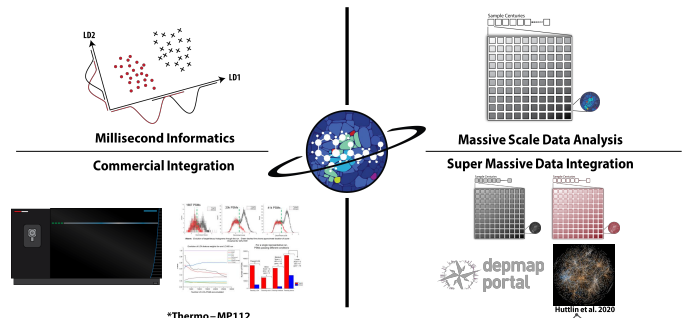


Real-time FDR filtering increases the number of quantified peptides and eliminated isobaric tag interference.
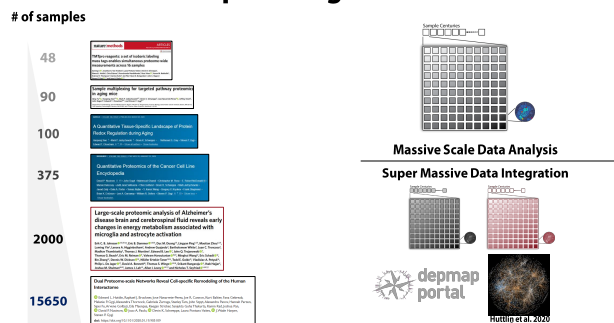
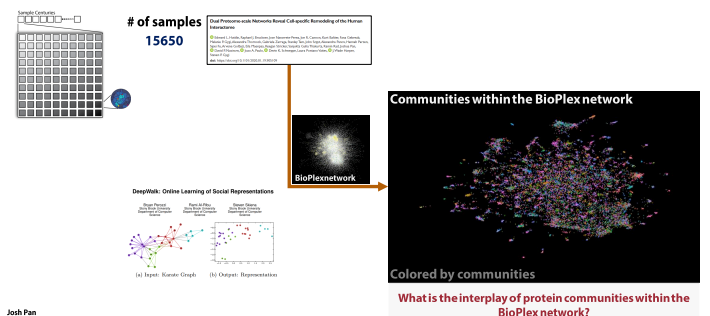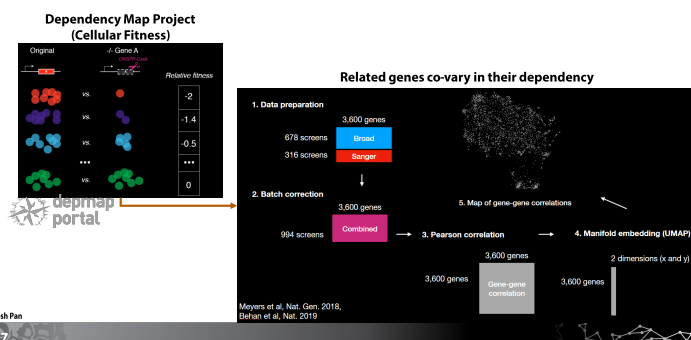Schweppe et al *UPR 2020*

13

## Moving Forward with ML



Millisecond Informatics
Commercial Integration

Massive Scale Data Analysis
Super Massive Data Integration

*Thermo – MP112

Huttlin et al. 2020

14

## As sample size grows…



# of samples

48
90
100
375
2000
15650

Massive Scale Data Analysis
Super Massive Data Integration

Huttlin et al. 2020

15

## Latent representations of graphs with DeepWalk



# of samples
15650

Communities within the BioPlex network

BioPlex network

DeepWalk: Online Learning of Social Representations

Colored by communities

What is the interplay of protein communities within the BioPlex network?

Josh Pan
Ed Huttlin

16

## Manifold embedding of cellular fitness profiles



Dependency Map Project
(Cellular Fitness)

Related genes co-vary in their dependency

Meyers et al, Nat. Gen. 2018,
Behan et al, Nat. 2019

Josh Pan

17

## Manifold Embedding to determine latent cell states



BioPlex network

Dependency Map

Incorporating multiple large-scale analyses to identify novel complex members, pathways or allied proteins, particularly for uncharacterized proteins.

Other datasets:
CCLE Proteome (Nusinow et al *Cell 2020*)
Cell perturbation
OxiMouse (Xiao et al *Cell 2020*)

Josh Pan
Ed Huttlin

18