# THE UNBEARABLE FUZZINES OF NMR DATA ?!?

... a brief reflection by Stan Sykora (ebyte.it) ...

Given the limited time,
I will present only a few NMR-based illustrations of a broader
question that is presently becoming acute due to advances in
*Big Data* handling, *Artificial Intelligence*, and other areas:

**How *hard* are the *hard Sciences* data, really?**

Tentative (but qualified) answers:

- The simulated ones are cute, hard and as sharp as it goes ☺
- The 'real' rest is *disgustingly soft* and as fuzzy as it goes ☹

**NMR is no exception – at all levels!**

# My own involvement in these things (just to explain)



An Mnova engine room

MESTRELAB RESEARCH
Chemistry Software Solutions

Stan

In the last 10 years I was basically trying to teach a computer to do the work of a spectroscopist (which is what 'automatic' really means)

# Expectations and Reality (in general)
## I. Fuzziness in the rules of the game and in imperfect comprehension

Some fellow programmers tell me:

You write software for NMR spectra analysis? OK, so your life is EASY!
Just ask some chemist what the rules are, implement them, and you are done.

Wouldn't that be nice! Alas, they were *just* programmers, poor things.
Some selected problems with coding what spectroscopists do:

- *No spectroscopist can describe how s-he does it, not in a flowchart manner*
- *At best, they provide some (good and useful) examples, and some weak rules*
- *For each such rule, I easily find tens of exceptions!*
- *If I ask two of them why an exception occured, they rebuke me that, of course, exceptions are to be expected (afterwards, they quarrel among themselves)*
- *Hence, there is no 'NMR Spectroscopy Master Book' that would hold <u>always</u>*
- *About 10% of published assignments (for example) are known to be wrong*
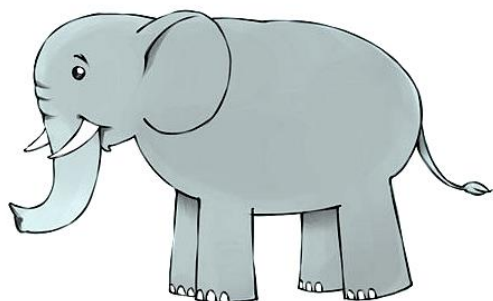
# … an interlude …

Yet, despite all the odds,
experienced spectroscopists are almost always correct !!!
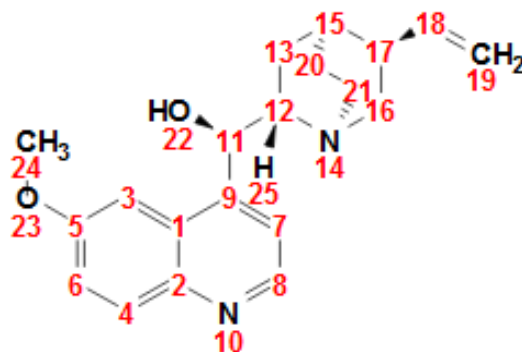
(That is really vexing!)

And I am to build software that, statistically, should beat them !??

(no worry, I am in sight of the target, but still a long way to go)
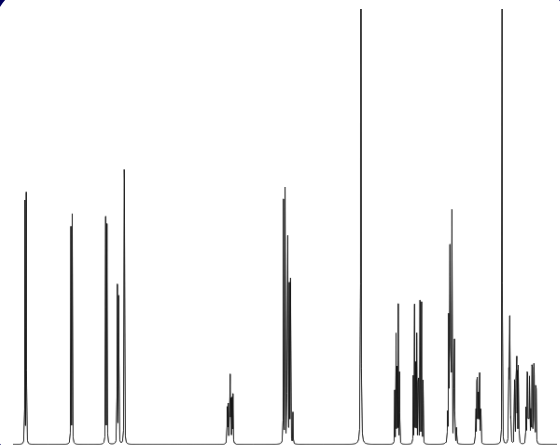
# Expectations and Reality (still in general)

## II. Reality is infinitely complex and never totally predictable. Hence, it is fuzzy!



This is NOT a **<u>real</u>** elephant! It's just a simple drawing of an elephant.



This is NOT a **<u>real</u>** molecule! It's just a structural sketch of a molecule.



This is NOT a **<u>real</u>** spectrum! It's a naive, simulated NMR spectrum.

Most of my code ( >80% ) is about whether what looks like XYZ, is really XYZ. On the opposite, *true Science* is the trivial part of the code, amounting to < 5%.

# NMR fuzziness at NMR ground level (the Reception Hall)

We would like to 'register' only clean **Peaks** and **Multiplets** (1D) or **Clusters** (2D)
But what enters through the front doors includes:

- Noise
  (from sample, probe, switching diodes, preamp, external, ...)
- Instrumental artifacts
  (dead time, imperfect shimming, rotation sidebands, spikes, ...)
- Sequence artifacts
  (ill defined expected intensities, missing and extra clusters, ...)
- Acquisition artifacts
  (fast repetition, FID truncation , limited digital resolution, ...)
- Evaluation artifacts
  (bad phase or baseline, referencing errors, DFT distorsions, ...)
- Expected impurities
  (solvent and its water, S- and Q-reference, residual solvents, ...)
- Unexpected impurities
  (various contaminants, fingerprints, ...)

It could definitely fill a full Book (I already promised to write one ☹)

# A few examples of undesirable spectral artifacts



Noise

Baseline roll

Imperfect phasing

Bruker / Jeol
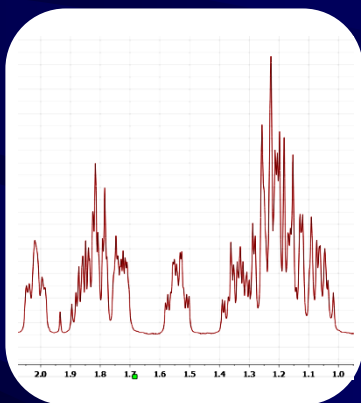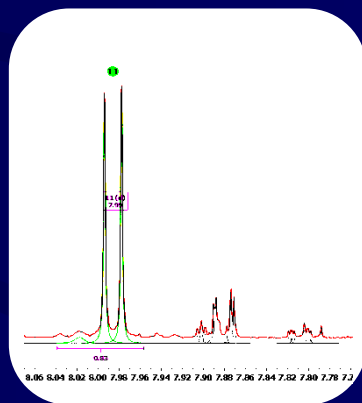smileys / brownies

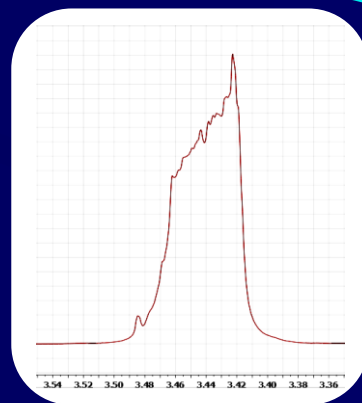FID truncation
effects

Underdigitization

Rotation sidebands

Peaks overlap

Impurities peaks

Sample temperature
drift effects

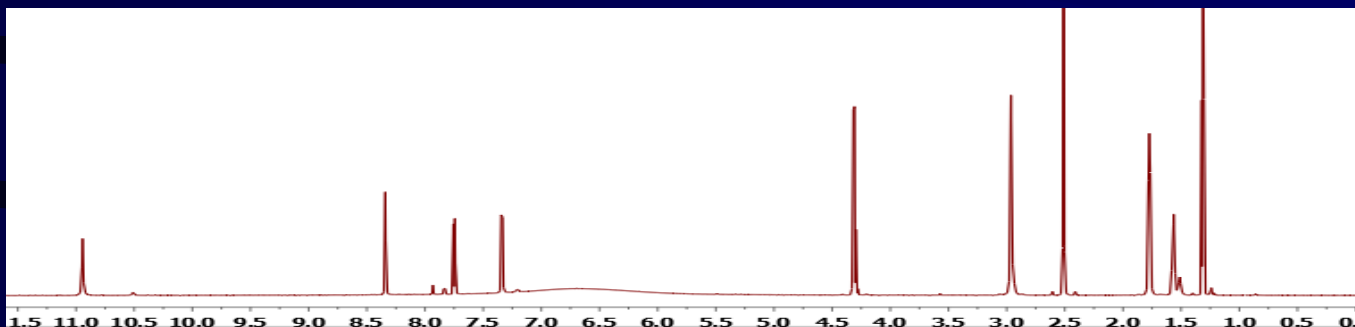=> There is so much to get rid of before doing anything serious with a spectrum !!!

# Ground floor data management (continued)

The main task of a Reception is to get rid of all the bullshit that enters.
In my world, this is one of the tasks of GSD (**Global Spectral Deconvolution**).
Why was GSD born? The basic idea, going back to 2006:

## What does a spectroscopist see?

Peaks, multiplets (singlets, doublets, AB quartets, triplets, quadruplets, …),
labiles, 13C satellite peaks, aromatic peaks, d-solvent and, reference peaks,
water peaks, impurities, reaction solvent residuals, spinning sidebands, …



1.5  11.0  10.5  10.0  9.5  9.0  8.5  8.0  7.5  7.0  6.5  6.0  5.5  5.0  4.5  4.0  3.5  3.0  2.5  2.0  1.5  1.0  0.5  0.0

## What does a programmer see?

Just an unexciting array of complex-valued data!
He can't understand what is the chemist talking about!

Implication: there is a big communication problem

# GSD: Historic notes

- I have started pushing the basic idea of Peaks List in 2006

- The rather complex algorithm was finalized in Summer 2008

- Since then, GSD was well tested and proved to be very robust

- Imitators and potential competitors appear around 2013 (CRAFT) ☺

- Today, GSD was applied in many different contexts, including structure verification, drug discovery, NMR quantitation, metabolomics, and others, and often turns out to be an enabling key to various novel avenues.

References:
  DOI:  10.3247/SL2Nmr08.011
  DOI:  10.3247/SL3Nmr09.003



F

# Ground floor data management (continued)

The GSD receptionist turns the input into a table of peaks that appear to be meaningful.
They get room keys and proceed to higher floors. All the rest is turned out!

## What does a spectroscopist see?

### Peaks,

(solvent, reference, impurity,…), multiplets, ...



GSD

| | ppm | Intensity | Width | Area | Type | Flags | Kurtosis |
|---|---|---|---|---|---|---|---|
| 44 | 2.612 | 0.474 | 1.224 | 8.289 | Compound | None | 0.000 |
| | | .437 | 11.624 | 239.109 | Compound | None | -0.019 |
| 46 | 2.955 | 150.983 | 4.326 | 9158.016 | Compound | None | 0.132 |
| 47 | 2.963 | 243.975 | 5.934 | 19624.604 | Compound | None | 0.376 |
| 48 | 2.970 | 132.463 | 3.810 | 6486.636 | Compound | None | 0.744 |
| 49 | 3.058 | 1.449 | 11.493 | 236.917 | Compound | None | 0.022 |
| | | .302 | 1.667 | 43.038 | Compound | None | 1.600 |
| 51 | 4.189 | 0.393 | 3.473 | 14.592 | Compound | None | 1.866 |
| 52 | 4.198 | 1.292 | 2.248 | 35.925 | Compound | None | 1.000 |
| 53 | 4.208 | 1.254 | 2.166 | 33.616 | Compound | None | 0.994 |
| 54 | 4.218 | 0.425 | 2.243 | 12.427 | Compound | None | 0.639 |
| | | 419 | 2.366 | 210.333 | Compound | None | 1.200 |
| 56 | 4.289 | 22.854 | 2.197 | 595.927 | Compound | None | 1.259 |

## What does a programmer see?

An array of … **peaks**, finally! That's **GREAT**

Implications: language synchronization => better communication

# GSD peaks-list example (a 400 MHz strychnine spectrum)

# Examples of peaks detection (WOW)



F

# Examples of peaks detection (GSD limits and artifacts)



- Possible missing peaks because of low S/N
- Possible missing peaks (marginal resolution)
- Possible extra peaks (marginal resolution)
- The resolution increase is not always welcome

**Fuzziness never goes [totally] away,
but it changes looks !!!**

# Automatic Peaks Editing

Having identified and tabulated all peaks, what more needs to be done ?

GSD by itself does not address issues like classifying each peak as:

- compound,
- primary or secondary solvent,
- potential labile,
- 13C satellite,
- valid member of a multiplet,
- impurity,
- S or Q reference,
- artifact,
- etc.

Nor does GSD group 1D peaks into multiplets and/or clusters and classify those.

All these are tasks which, unavoidably, bring-in fuzziness of their own.

# Automatic peaks editing example

## water peak recognition (primary and secondary solvent peaks are drawn in red)



Peaks editing uses greedily any available info (1H, HSQC, molecular predictions, ...).
It is the first plank in the evaluation hierarchy where specific NMR know-how is used.

# Automatic peaks editing example
## CHCl$_3$ identification in an overlapping aromatic multiplet



It uses even the $^{13}C$ satellites (209.25 Hz apart) and their isotopic shift
(the satellite pair center is -2.67 ppb from the main peak)

F

Water:

Just one $^{13}C$ satellite peak on the left was detected. ☹
The other four 'contaminate' the compound multiplet.

DMSO:

F

Higher level NMR data evaluation tasks: Automatic Assignments (Thalidomid 600 MHz)

Notice the correctly identified labile peak.
Labiles are often a big pain in … They can be fuzzier than hell!

Etc.… **data evaluation tasks never finish** (they grow in a fractal way)

F

# … data evaluation tasks never finish …

OK, But !!!

At every new step new uncertainties arise and get added.
And the 'NMR master books' never tell the whole story!

Since, at each stage  even the input data are fuzzy,  what
about fuziness quantification and fuzziness propagation.

This is an a huge emerging area, needing more research.
NMR is just a little niche in it, but a nice example, too.

F

# The Scoring System

Among a dozen or so algorithms that I have developed (so far unpublished) a prominent role is held by the *mathematical concept of scoring system.*

## About votes and their Significances

A *scoring system entry* is a pair of real numbers associated with a query which was 'voted' upon by a 'voter' (or simply a test function):

{*vote*, *significance*}

where -1.0 (total rejection) ≤ vote ≤ +1.0 (total acceptance).

In general, the *significance* is an attribute of the 'voter' (0.0 ~ idiot, 2.0 ~ normal, 10.0 ~ expert).

The comparison of a scoring system to a committee set-up to vote on an issue is very appropriate (and even practically viable).

In such a context the *entries* are the votes cast by the voters, each vote accompanied by the particular voter's 'rating'.

# A graphical representation of a Scoring System

The graph below shows nine different entries into a scoring system, represented as peaks.

For each peak, its location corresponds to the *vote*, and its height to its *significance*.
The peak height is inversely proportional to its width since its integral must be 1.0
(*some* vote must be cast, anyway).



Positive votes      Undecided:      Negative votes

Insufficient data
and/or
Incapable voter

# Combining votes and a Scoring System output

A central idea is that a scoring system output is again a {*vote*, *significance*} pair.

This means that any number of such pairs can be combined into a single one,
(and N committees can form a super-committee).

To enable this game, we need a mathematical @ operator of the type
{$vote_1$, $significance_1$} @ {$vote_2$, $significance_2$} = {$vote$, $significance$}

The operator must give reasonable results and also satisfy several constraints.
The mathematical problem was solved [almost*] satisfactorily.
Examples:



*Note: A mathematically perfect, non-approximate solution does not exist

# Example: Using a Scoring System to locate a solvent

Each peak is scored on being the pivot peak of the primary solvent (DMSO):

# Usage of Scoring Systems: a summary

Mnova Data Processing package uses Scoring systems pervasively, starting from Peaks List Editing up to more complex tasks.

**Typically over 1000 scoring systems are set-up and over 10000 votes are cast in a single run of any advanced NMR data evaluation task.**

It often looks almost as a magic to witness how the scoring-system algorithm can arrive at a correct deduction from very fuzzy data.

This type of technologies will no doubt flourish in near future also also outside NMR, especially in the social area and big data areas. There are other such emerging ideas (e.g., artificial intuition).

# Side remarks



Spectroscopists beware !

Your job is being stolen right now

In a recent talk I have concluded with a slide like this
**The NMR Spectroscopist in XXII Century ==>**

1) In 10 years time, given a set of spectra (say $^1$H, $^{13}$C, HSQC, COSY, HMBC) and five possible structures, will anybody go through the tedious 'manual' analysis of the thing? Obviously not, if a click will be enough to do it.

**Hence, why should anybody be interested to laboriously study how to do it?**

2) And will anybody still study spin Hamiltonians and product operators if absolutely anything you can do with them is in the reach of a click?

**Few are versed in these things even today, and their numbers are dropping!**

# Thank You for Your Attention

## Coworkers:

**The whole of Mestrelab team**, in particular:
Carlos Cobas, Felipe Seoane, Esther Vaz,
Santiago Dominguez, Maruxa Sordo,
Cristina Geada, Pablo Monje,
...



**Getting on-board:**
Ester Maria Vasini
(Extra Byte might yet take a new turn ☺)

**External collaborators and/or 'collaborators':**
far too many to name (this is partially fuzzy, too)

Presented at XLVI Annual Congress of GIDRM, Fisciano, Salerno (Italy), 27-29 Sep 2017