

TADA!

Topics in Algorithmic Data Analysis

Pauli Miettinen & Jilles Vreeken



UNIVERSITÄT
DES
SAARLANDES



mpi max planck institut
informatik

Organization

- 5 credit points
- Lectures 2 h per week (now, here)
- Divided into four topics
- No weekly tutorials
- Four written assignments + final exam
 - First given next week
- What/when/where/how/why explained next week

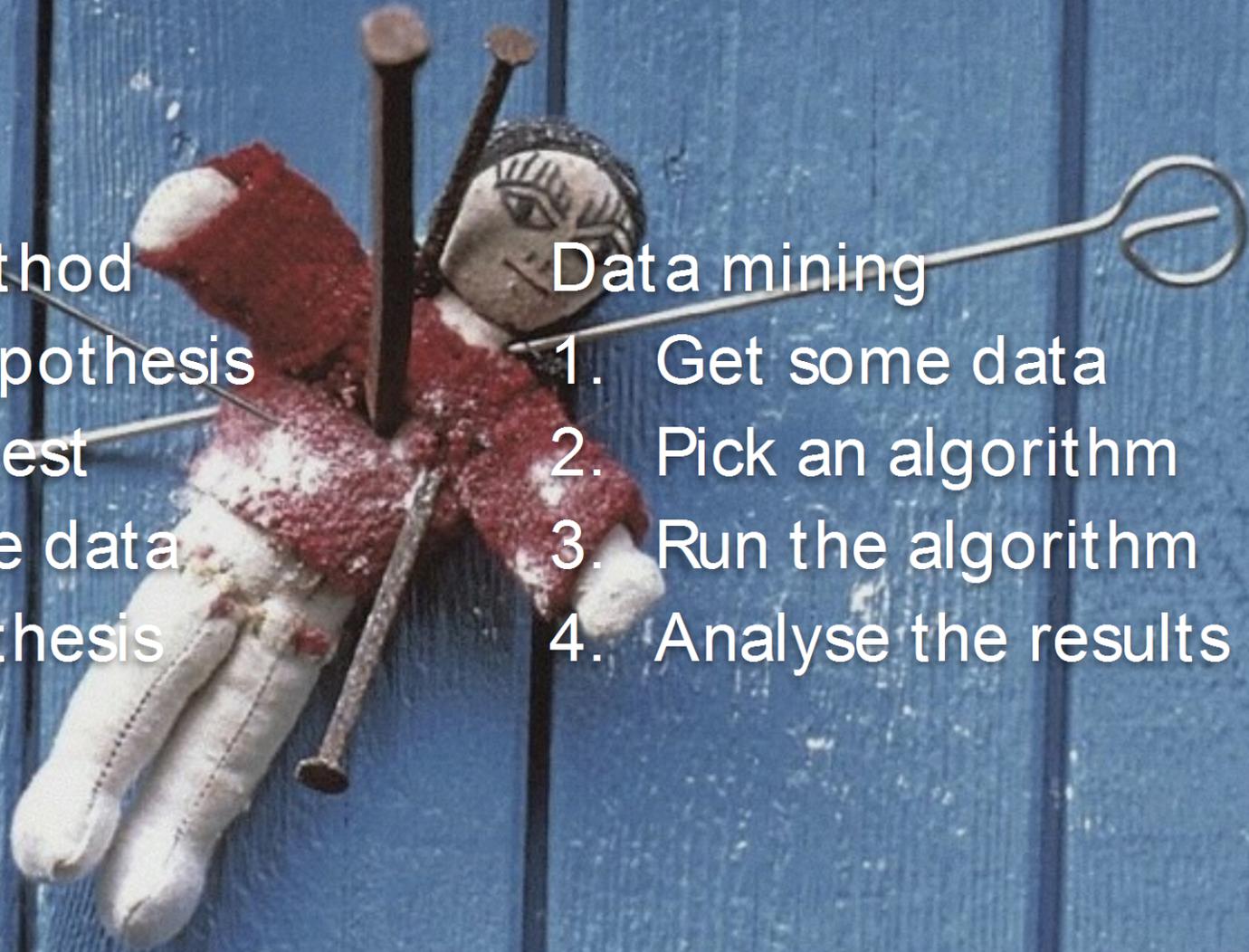
Data mining = voodoo science

Scientific method

1. Form a hypothesis
2. Design a test
3. Collect the data
4. Test hypothesis

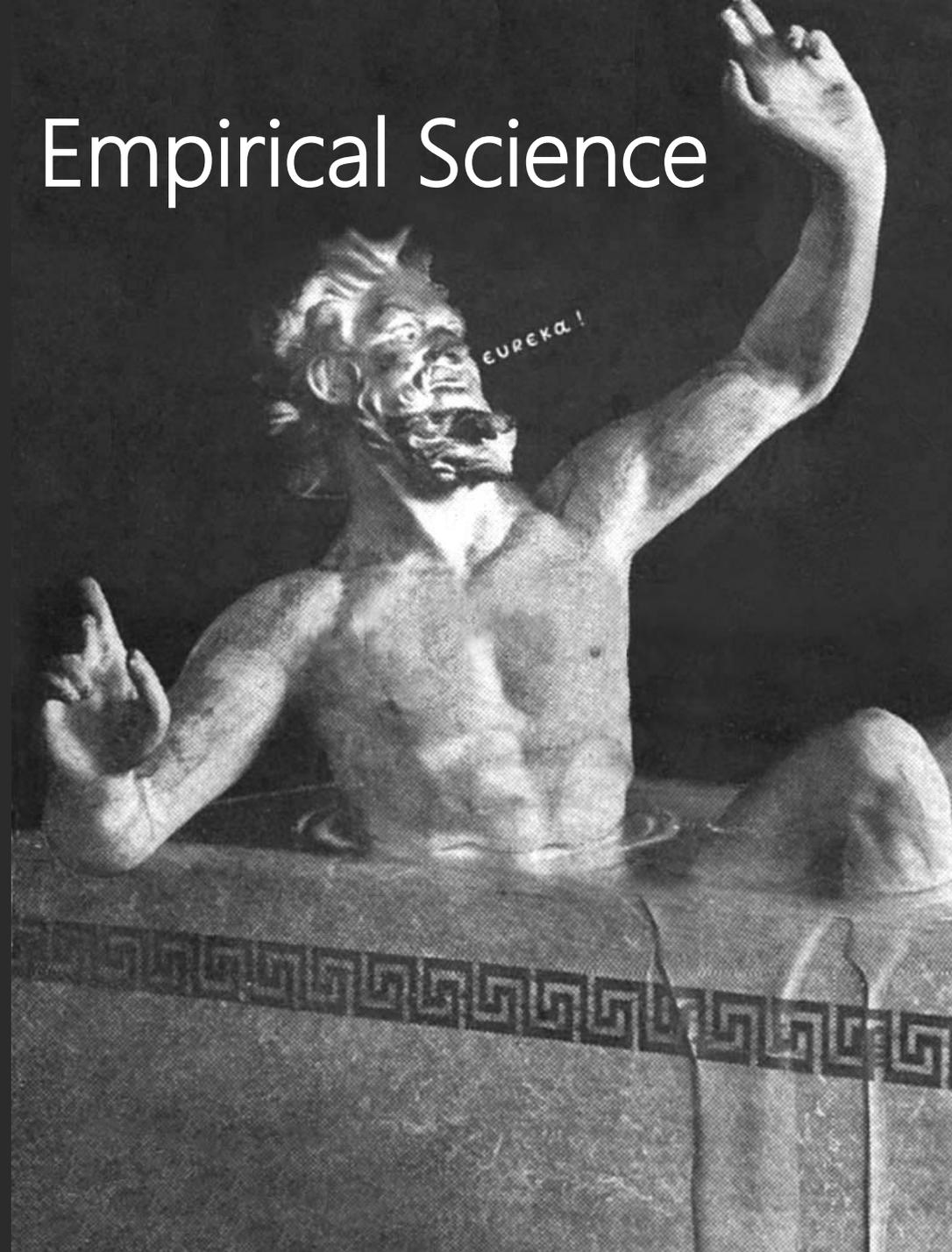
Data mining

1. Get some data
2. Pick an algorithm
3. Run the algorithm
4. Analyse the results



1st Paradigm: Empirical Science

For thousands of years,
science was empirical:
describing natural
phenomena



2nd Paradigm: Theoretical Science

The last few hundred years science was **theoretical**:
used models, generalizations, made **predictions**



3rd Paradigm: Computational Science

The last decades, science was **computational**:
complex models **simulating** complex **phenomena**



4th Paradigm: Data-Intensive Science

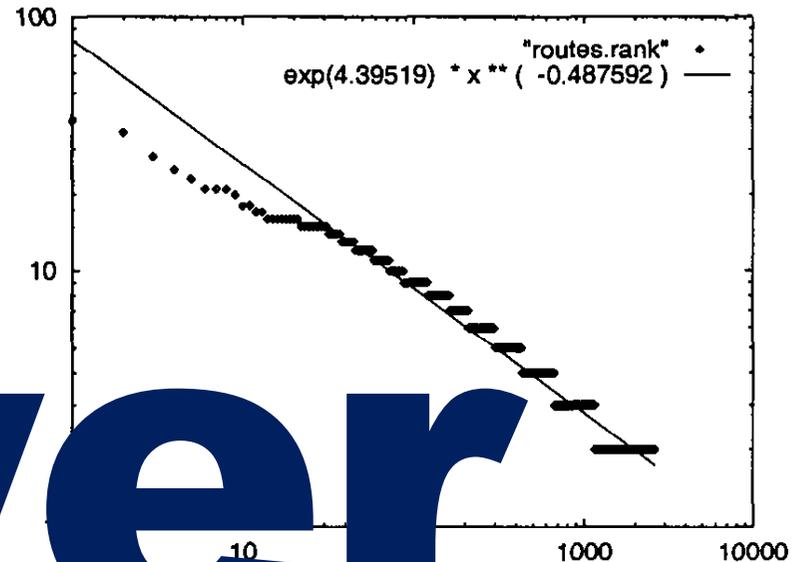
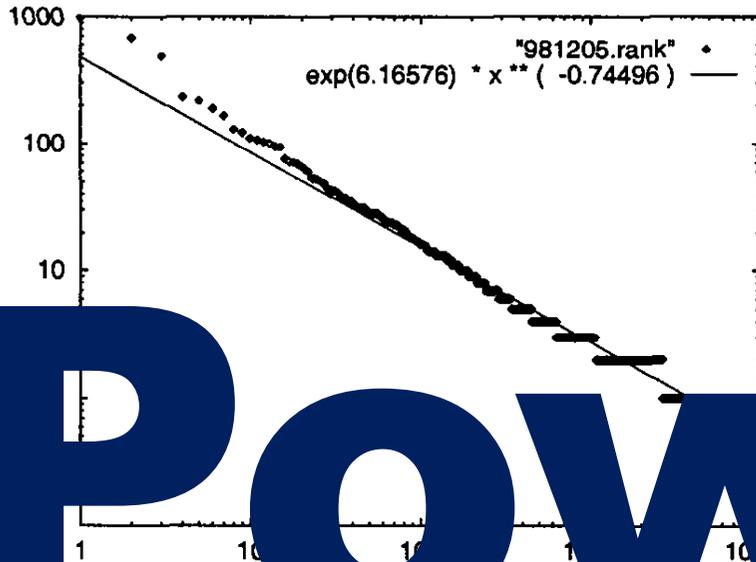
Interesting phenomena are **too complex** to come up with good hypotheses.

We need to
unify theory, experimentation, and simulation

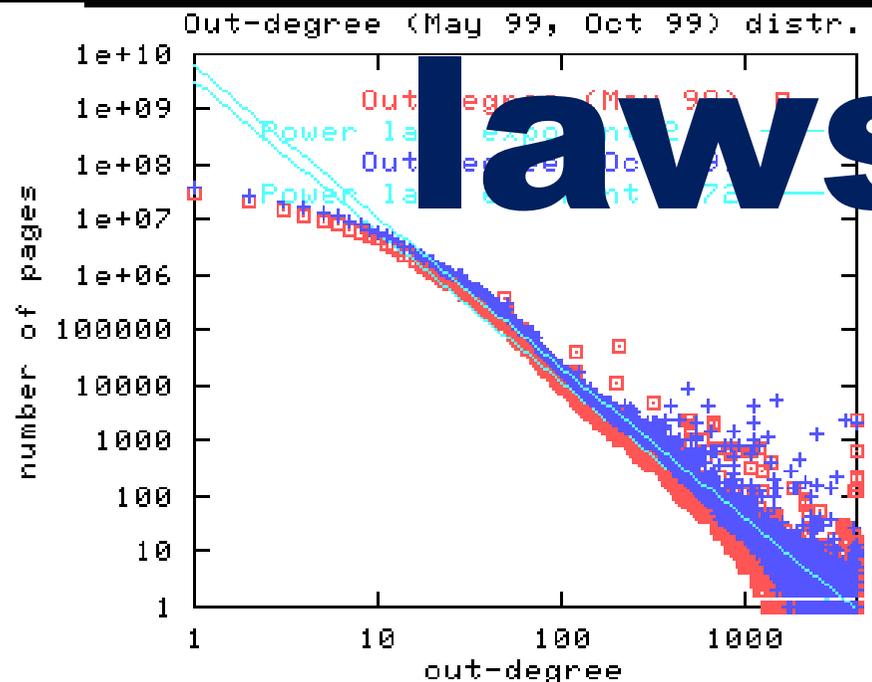
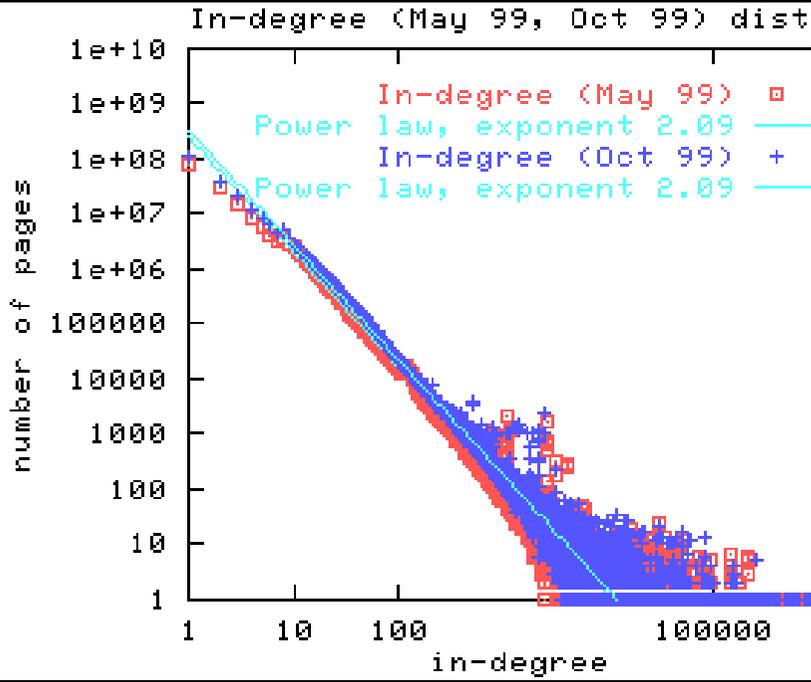
capture data, **mine** hypotheses, **inspect** and evaluate,
generate extra data to **select** the best ones, **iterate**

iterative procedure
between **world and model**,
scientist in the middle

Power laws



(b) Figure 95



laws

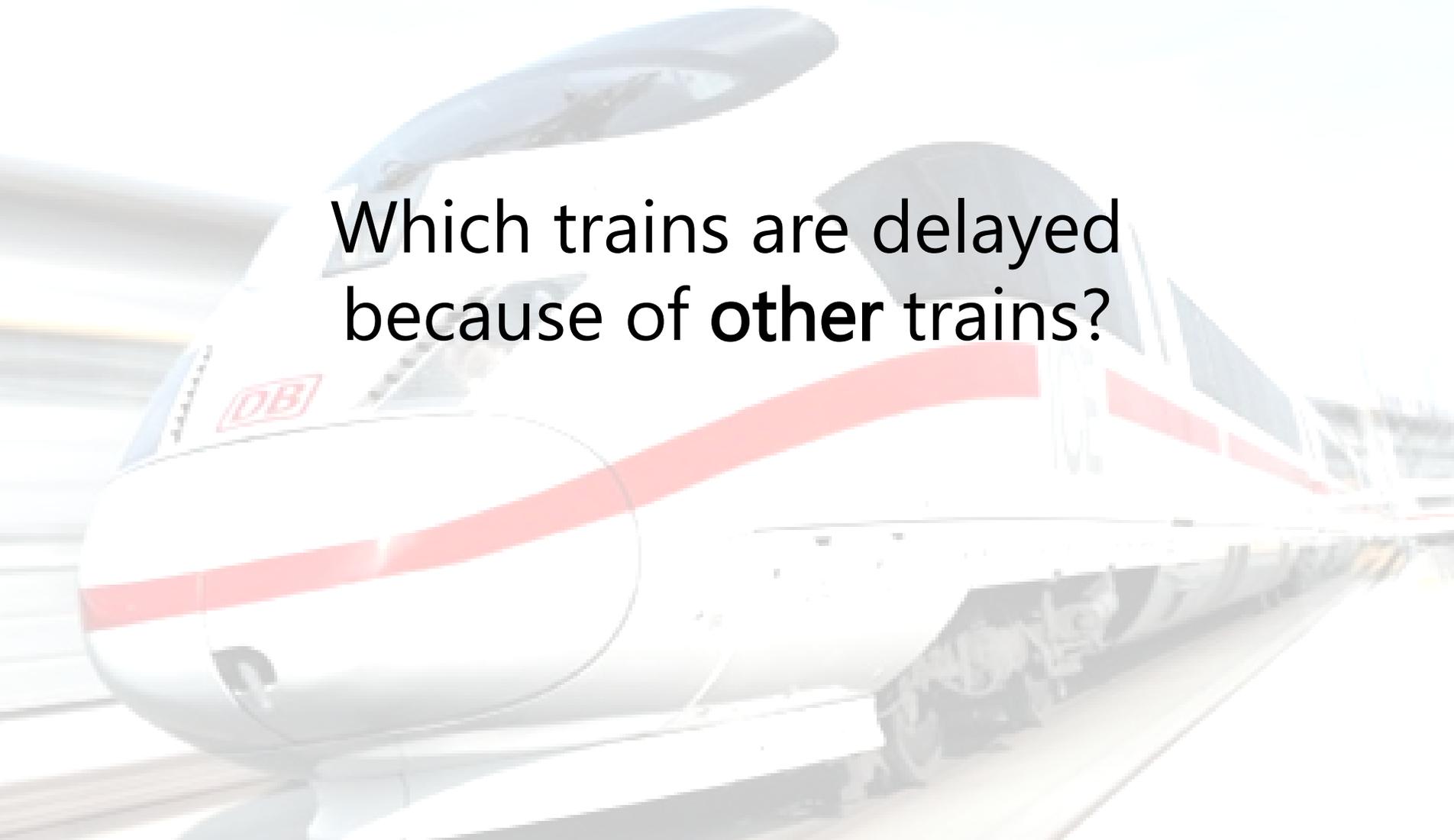
Shopping Data

A woman with long dark hair, wearing a black jacket and dark pants, is walking through a grocery store aisle. She is carrying a black shopping basket in her right hand and a black shoulder bag. The aisle is filled with shelves of various products, including boxes and bags of goods. The lighting is bright, and the overall scene is a typical grocery store environment.

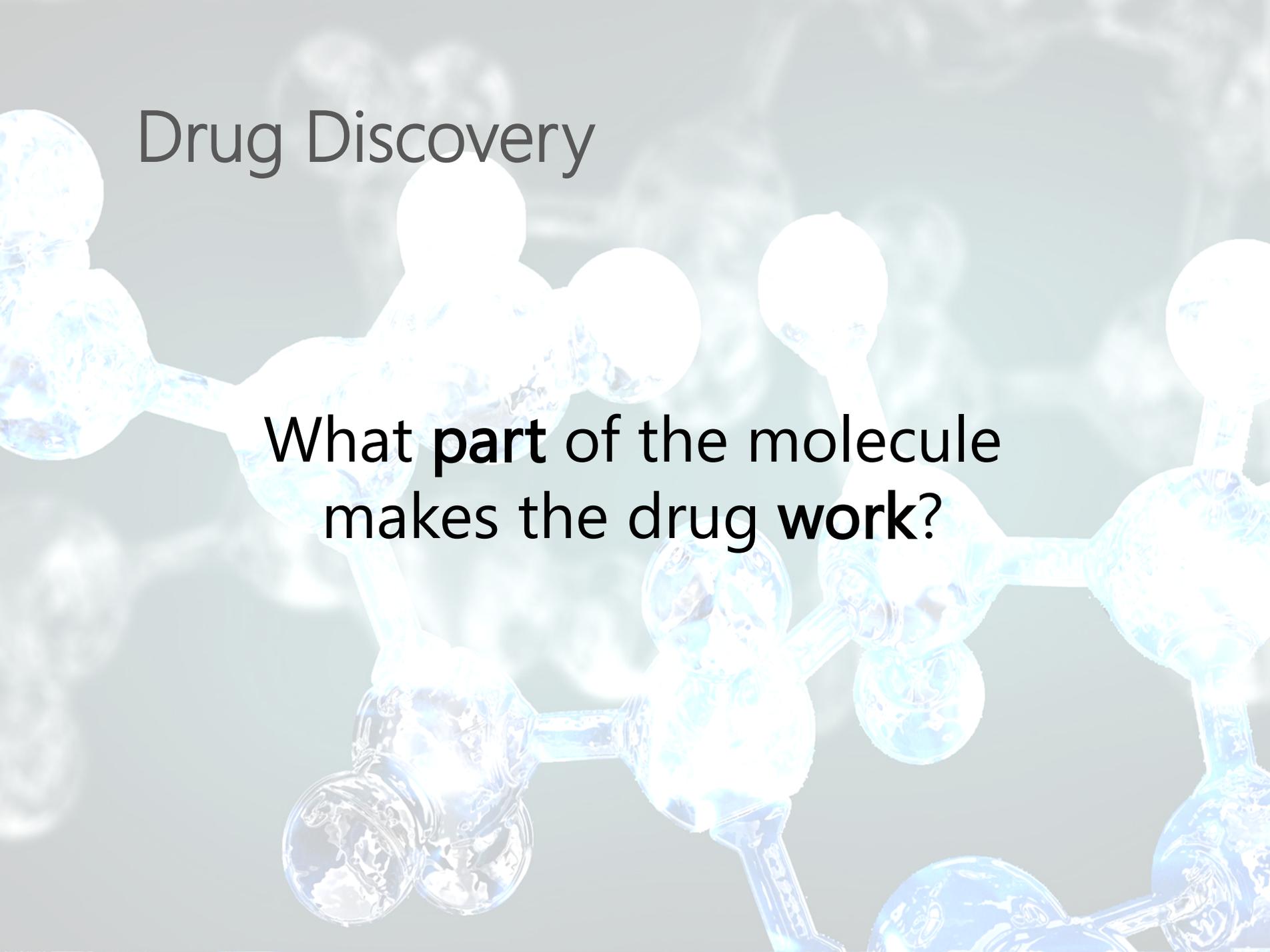
Which products are
often bought
together?

Train Delays

Which trains are delayed
because of **other** trains?

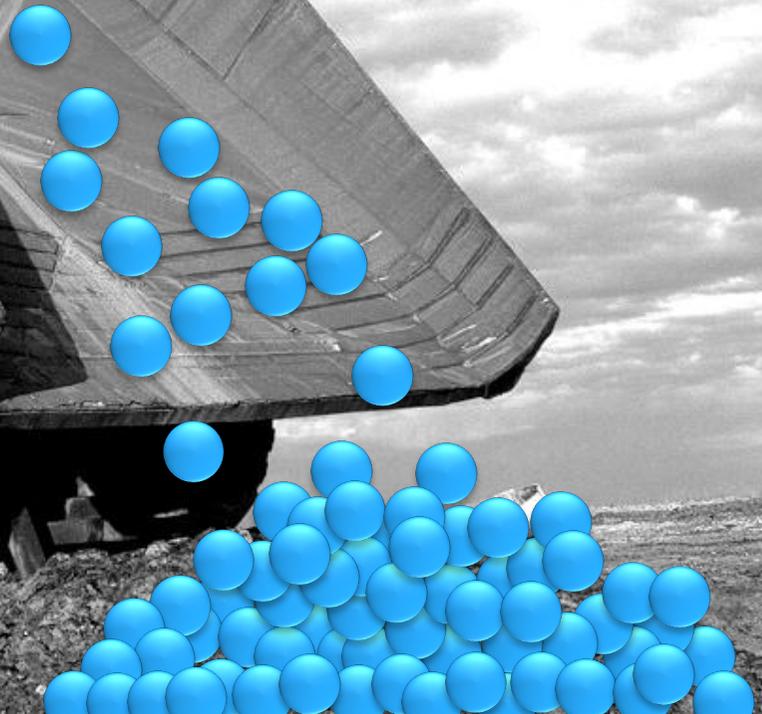


Drug Discovery

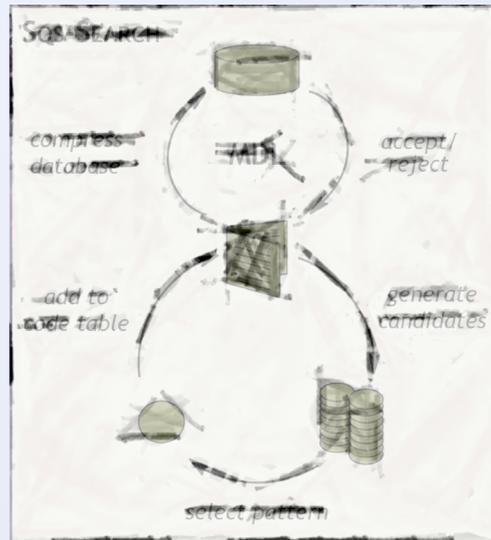


What **part** of the molecule
makes the drug **work**?

More patterns
than you
can shake
a stick
at

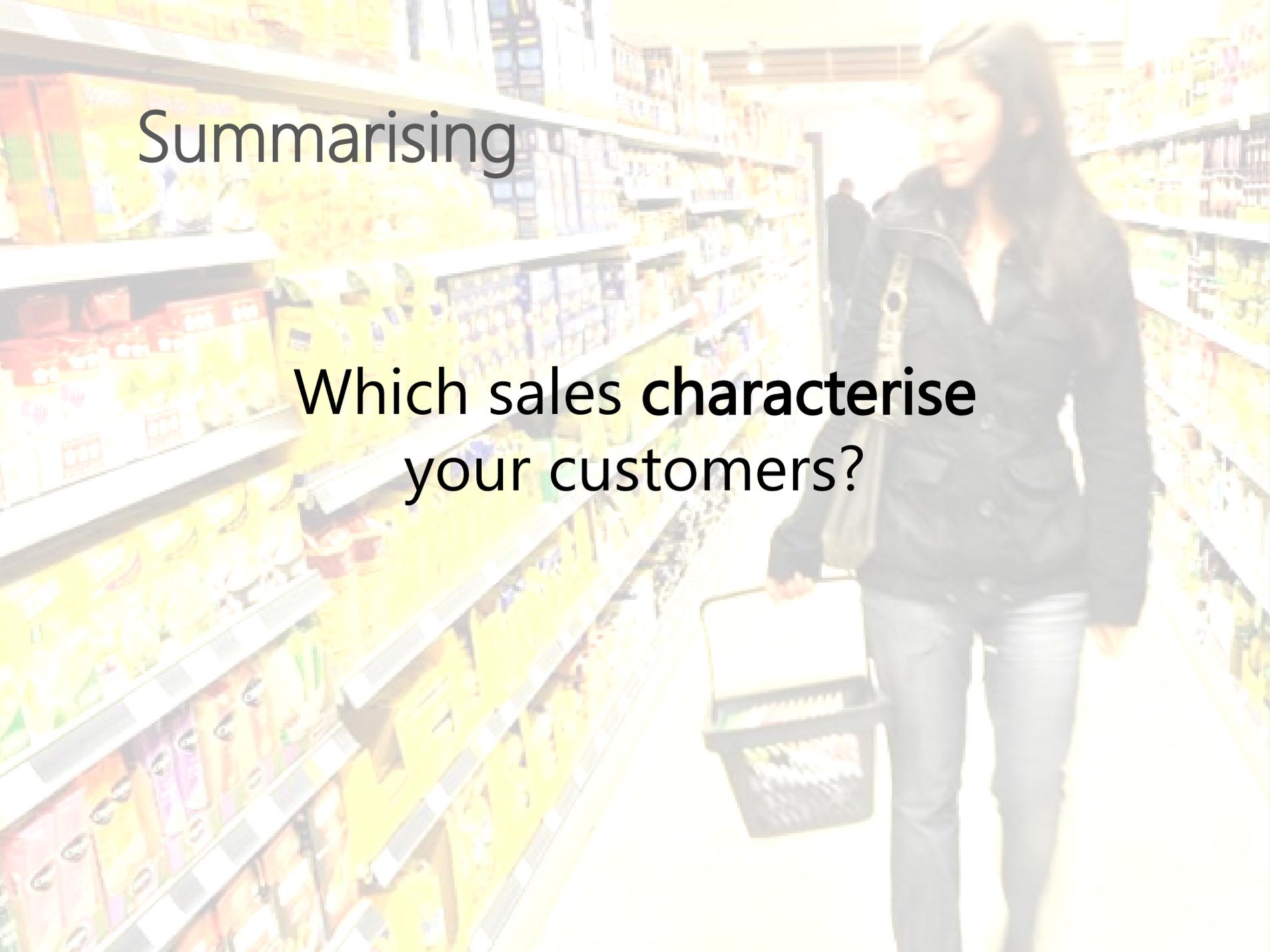


Pattern-based Modelling



support vector machin svm
associ rule mine
nearest neighbor
frequent itemset mine
naïv bay
linear discrimin analysi lda
cluster high dimension
state art
frequent pattern mine
algorithm
synthet real

Mining
Algorithm

A woman with long dark hair, wearing a black jacket and dark pants, is walking through a supermarket aisle. She is holding a black shopping basket in her right hand. The aisle is filled with shelves of various products, including packaged goods and fresh produce. The lighting is bright, and the overall atmosphere is that of a typical grocery store.

Summarising

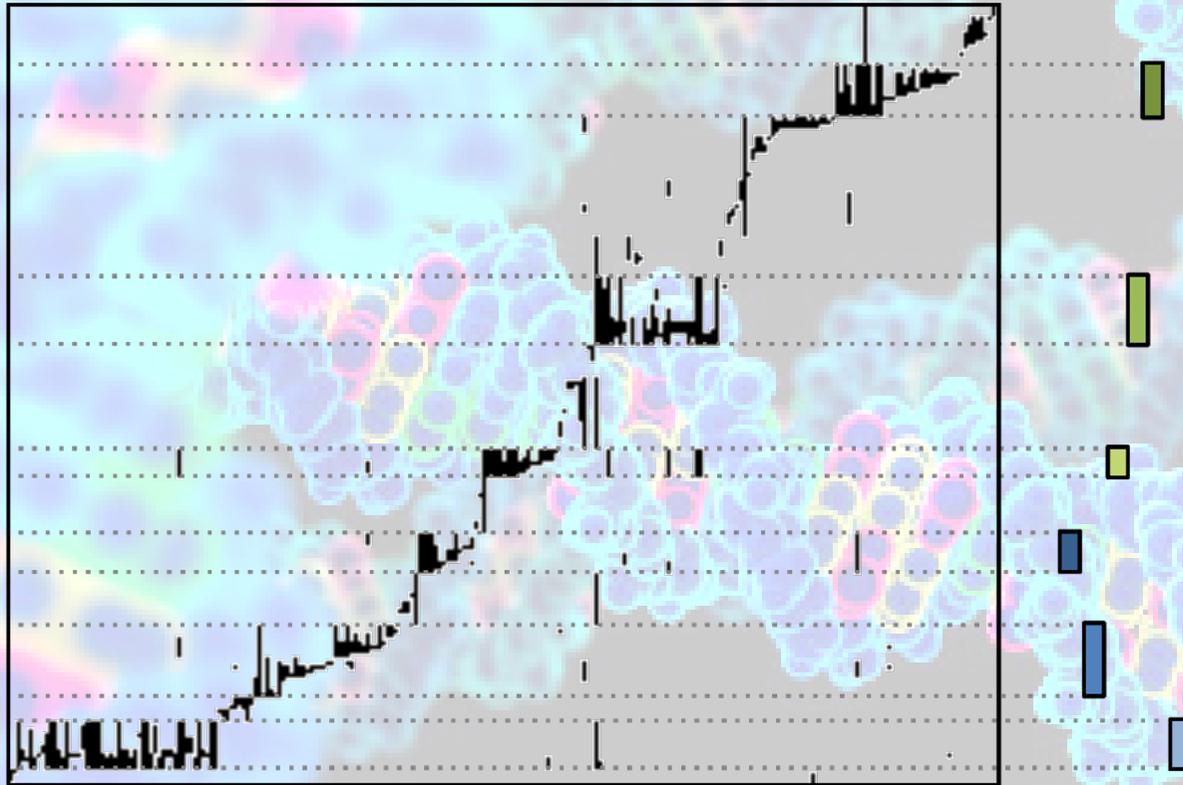
**Which sales characterise
your customers?**

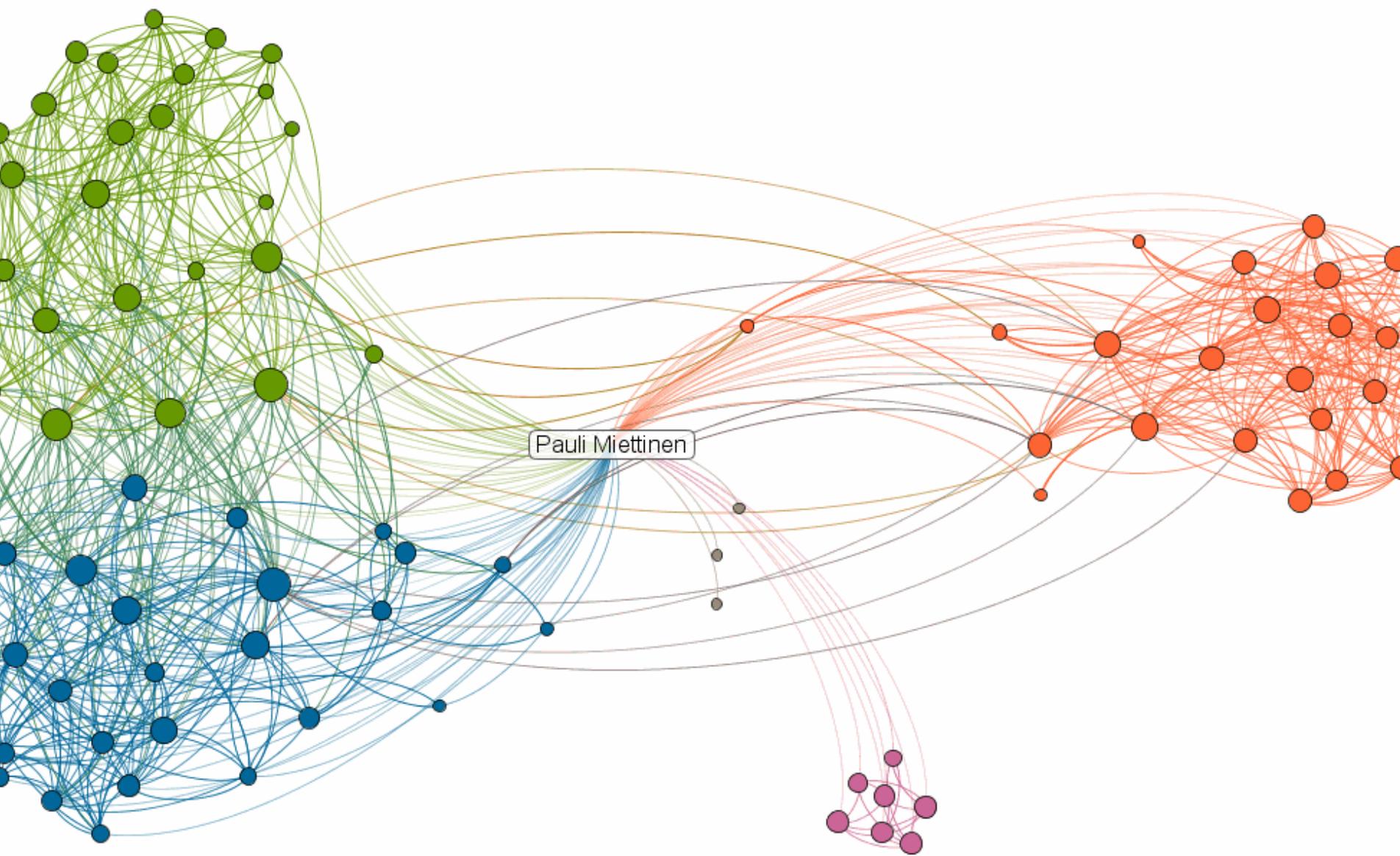
Summarising

Transactions (patients)

Itemsets

Items (oncogenes)





Google Flu

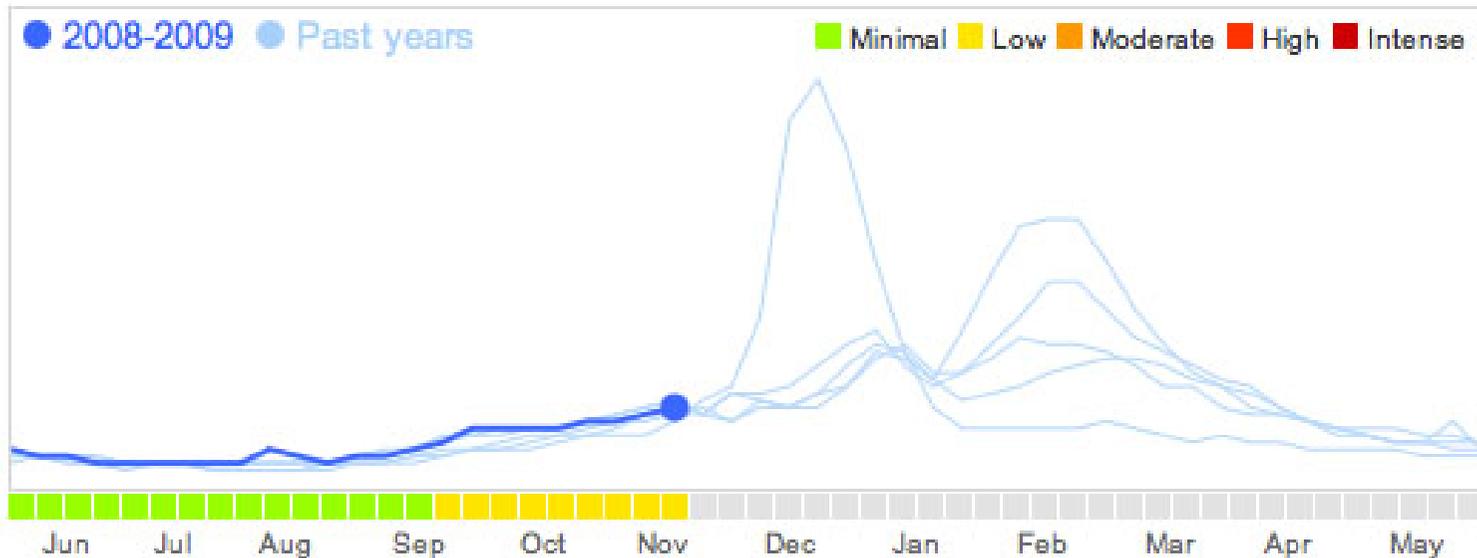
google.org Flu Trends

United States flu activity: ■ Low

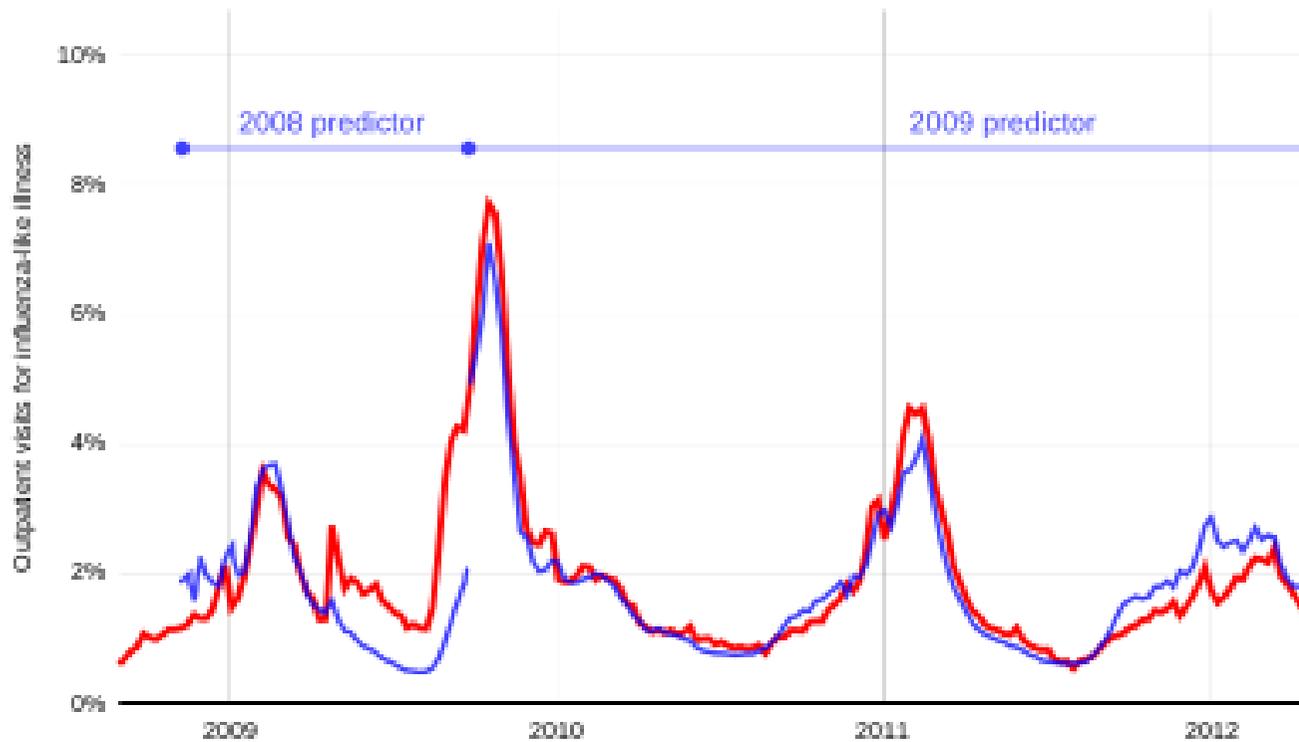
Entire United States

● 2008-2009 ● Past years

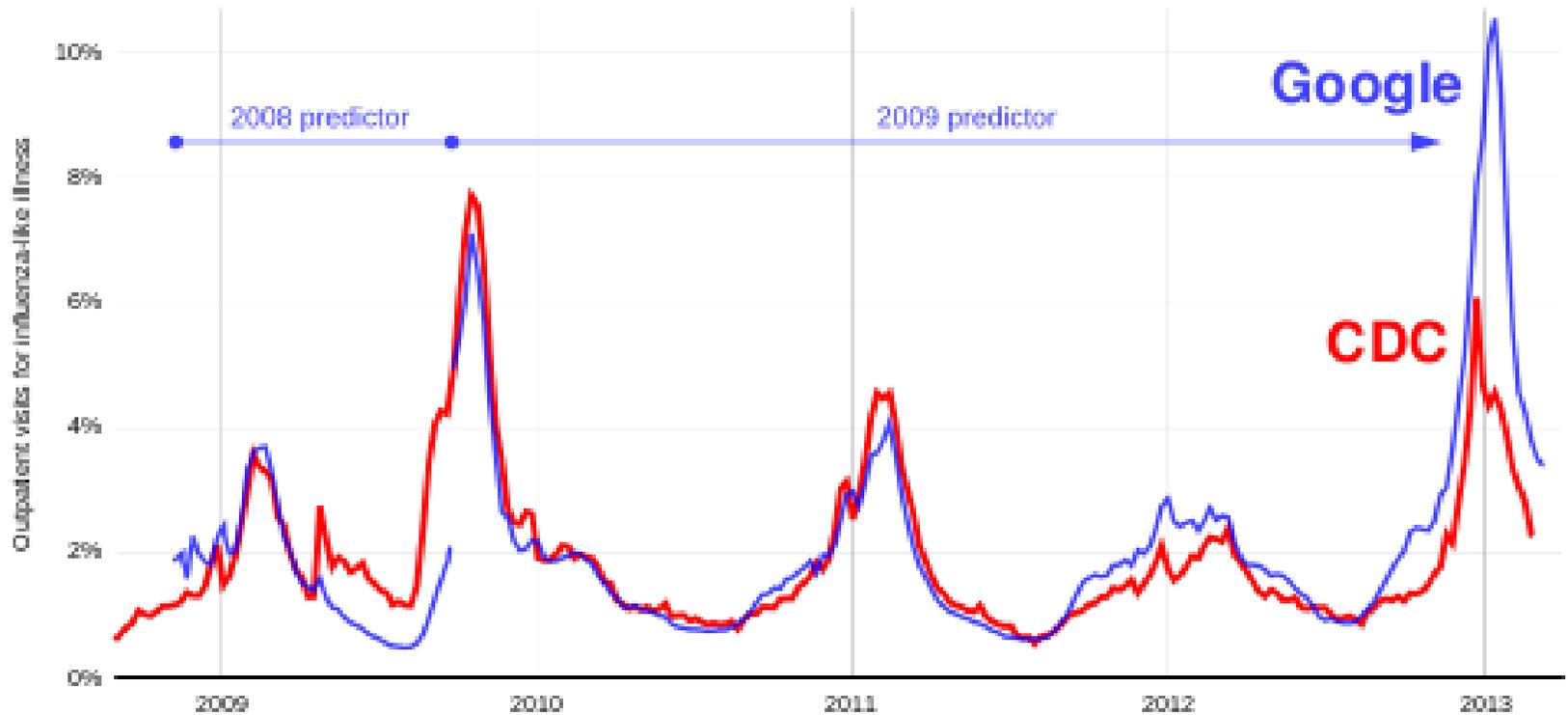
■ Minimal ■ Low ■ Moderate ■ High ■ Intense



Quite Healthy



Patient Deceased



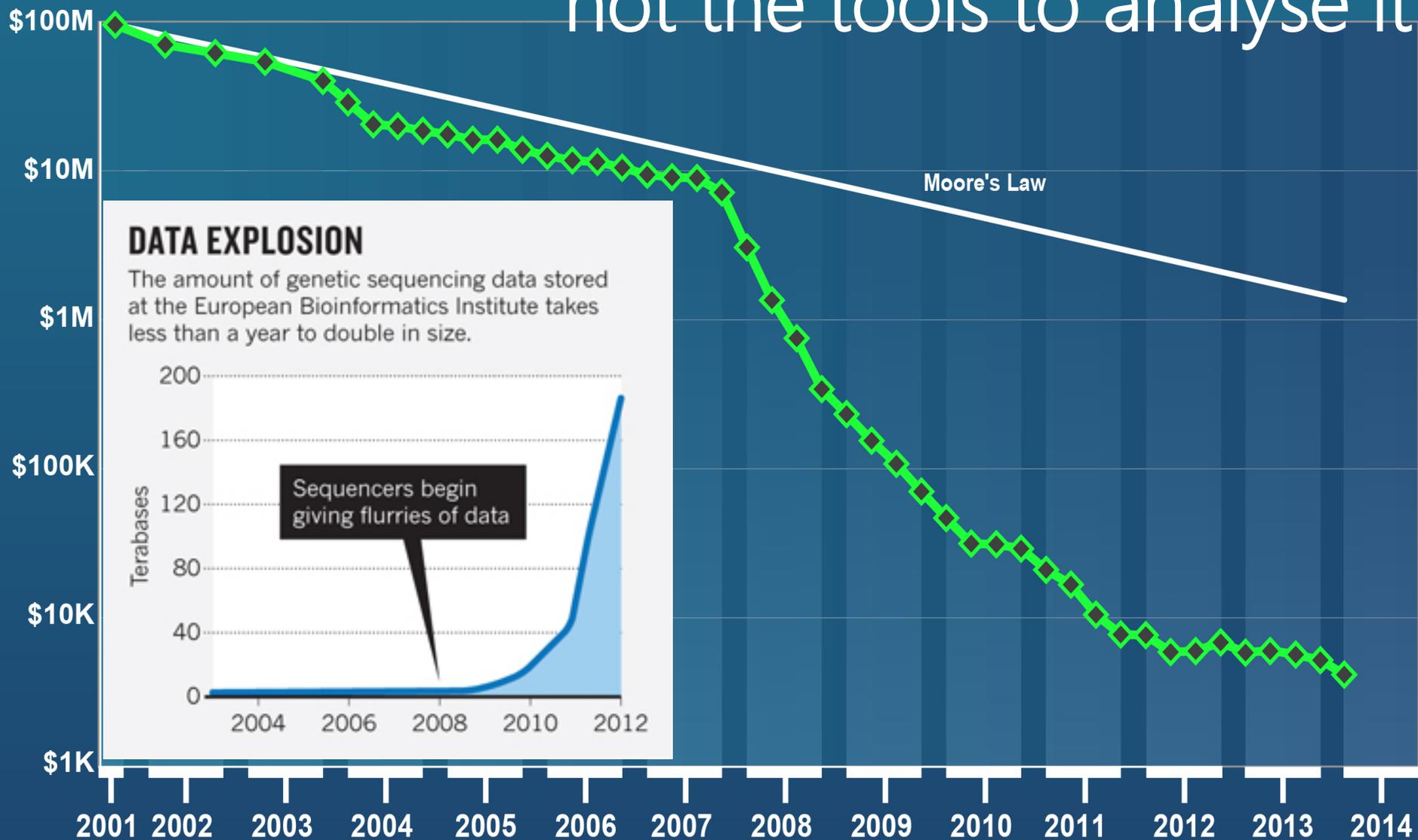
BIG DATA



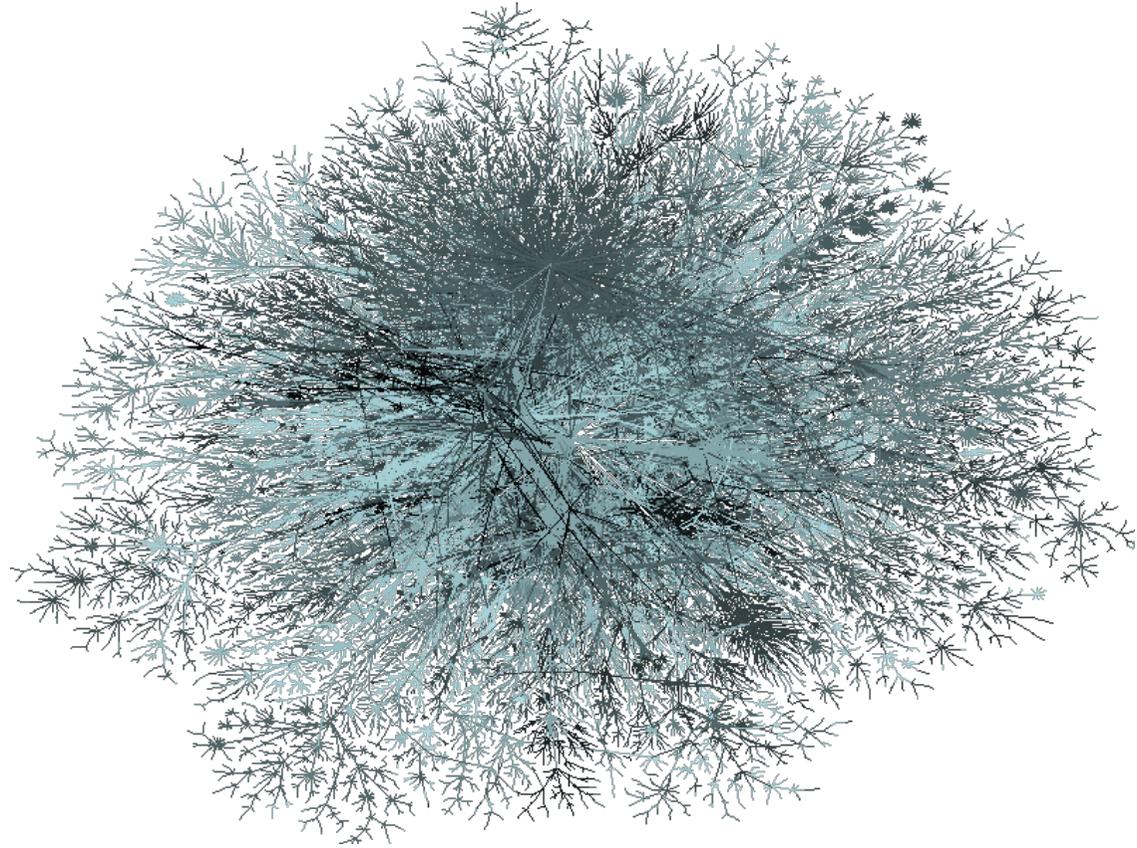
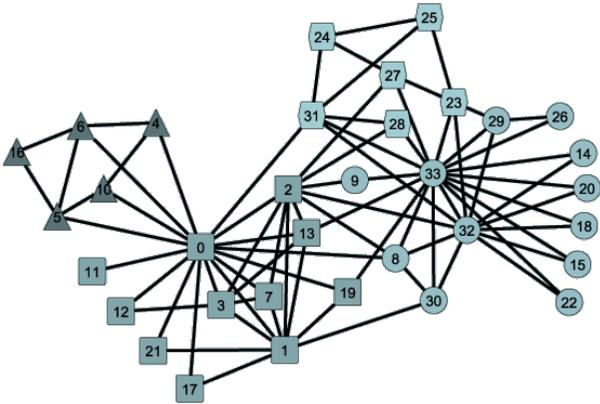
No model is perfect



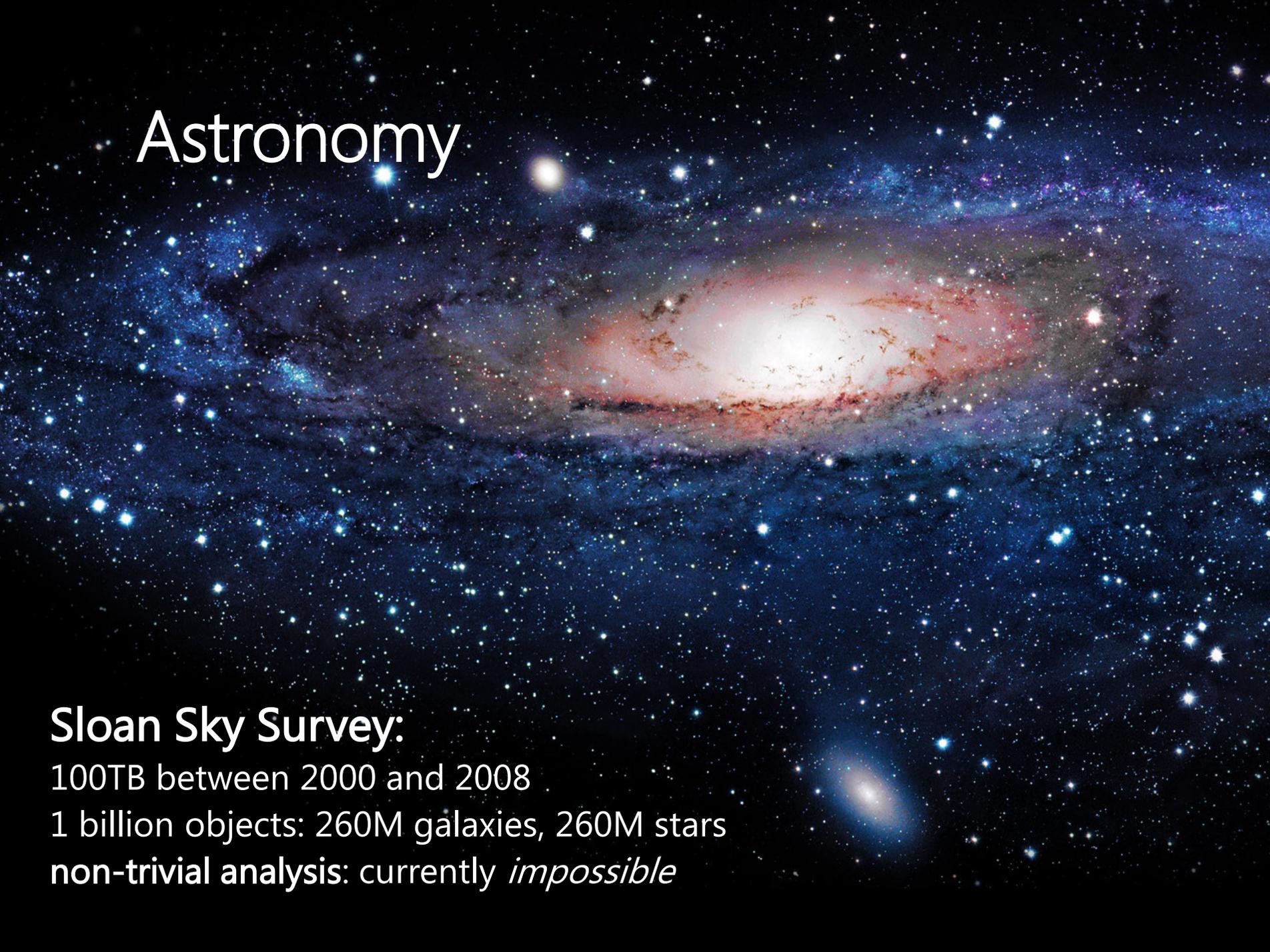
Science has lots of data, not the tools to analyse it



Social Science & the Web



Astronomy



Sloan Sky Survey:

100TB between 2000 and 2008

1 billion objects: 260M galaxies, 260M stars

non-trivial analysis: currently *impossible*

We Can Do It!

