

Introduction to Machine Learning algorithms

Dr. Fabien Plisson

Chemoinformatics in Drug Discovery

LANGEBIO, UGA CINVESTAV

October 15-18, 2019 - Irapuato, Mexico

Program

General Concept

Data Mining

Pre-processing

Exploratory Data Analysis

Predictive Modelling

Data Visualisation

Program

General Concept

Data Mining

Pre-processing

Exploratory Data Analysis

Predictive Modelling

Data Visualisation

What is Machine Learning?

Use and development of algorithms to **identify patterns in the data that allow computers to learn**, and thus predict future behaviours. Computers are not explicitly programmed for such a task.

Supervised ML: The result is known, the computer looks for the best way to predict this result.

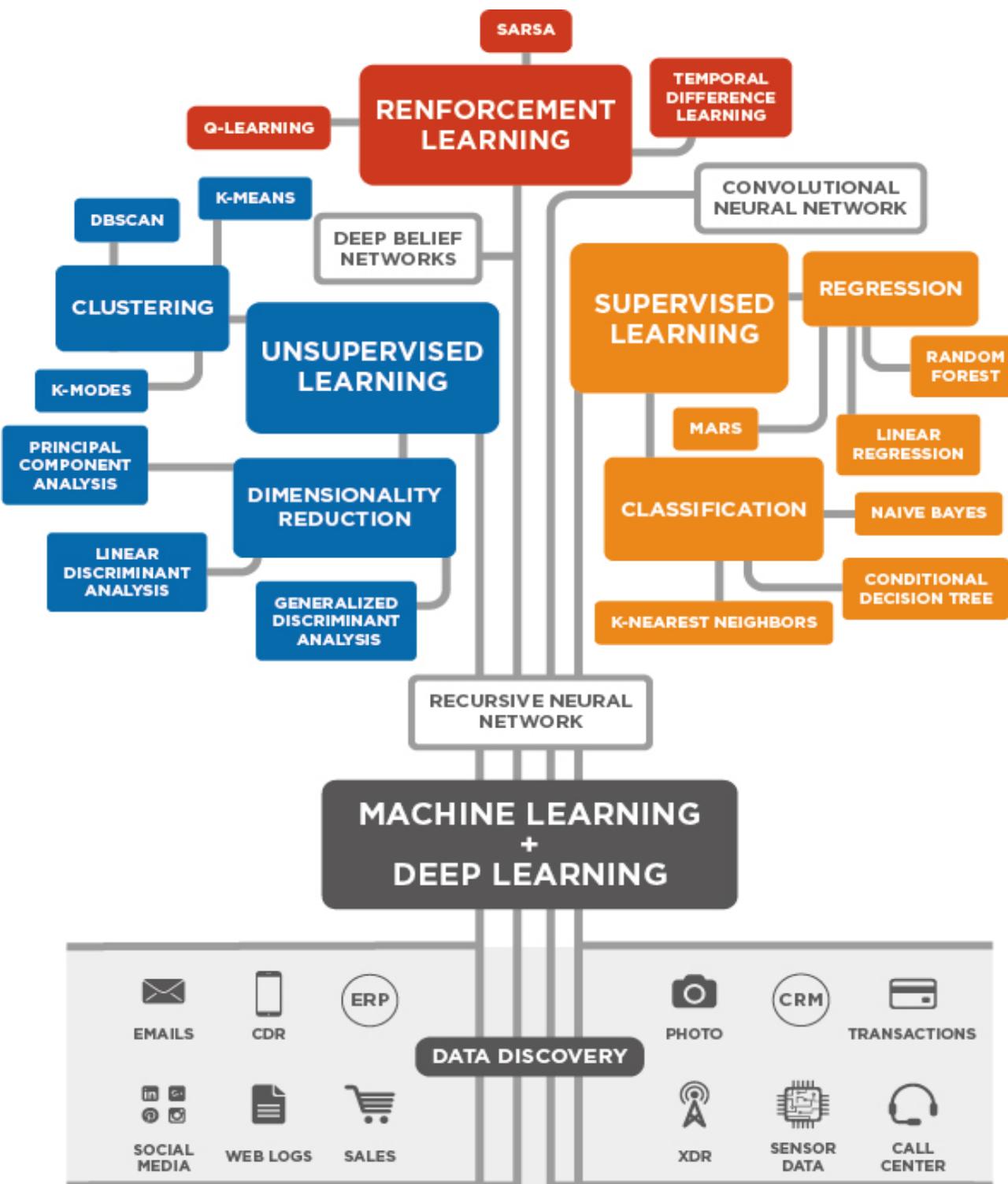
Unsupervised ML: The computer does not know the result, it groups objects based on similar patterns.

Reinforcement ML.



python



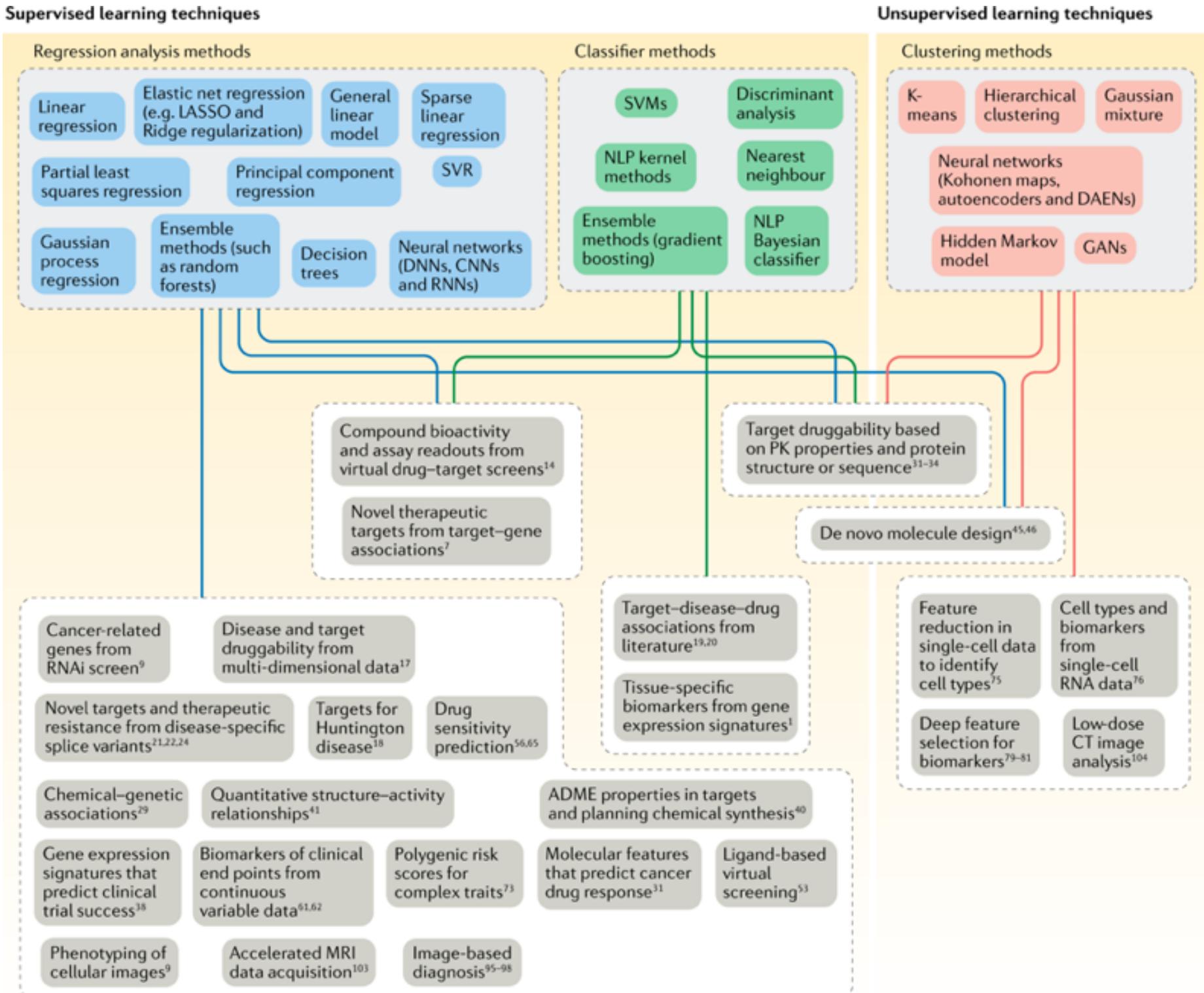


Applications

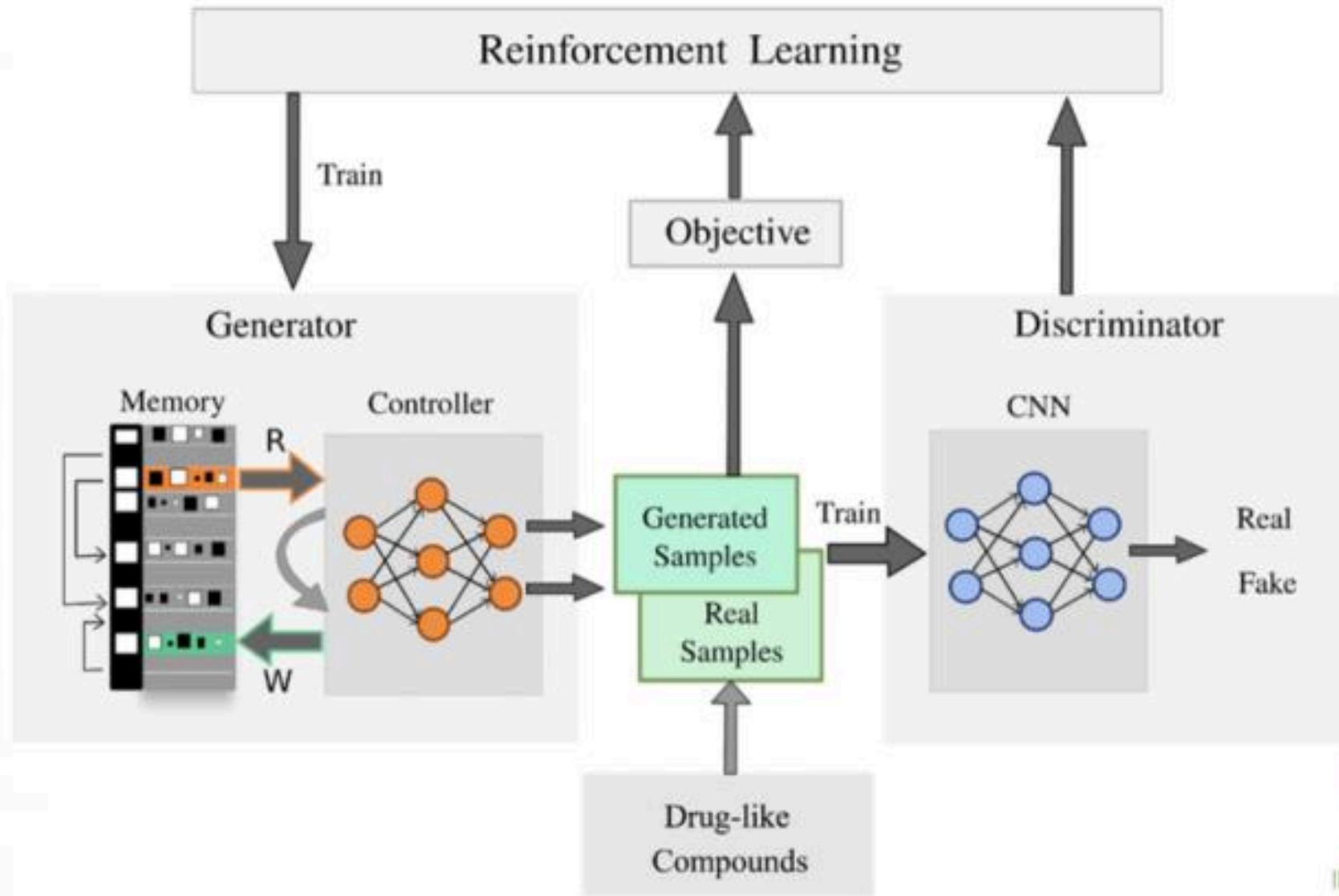
Database marketing
Handwriting recognition
Recommendation systems
Information retrieval
Object recognition in computer vision
Optical character recognition
Spam detection
Pattern recognition
Speech recognition
Bioinformatics
...



Machine Learning in Drug Discovery

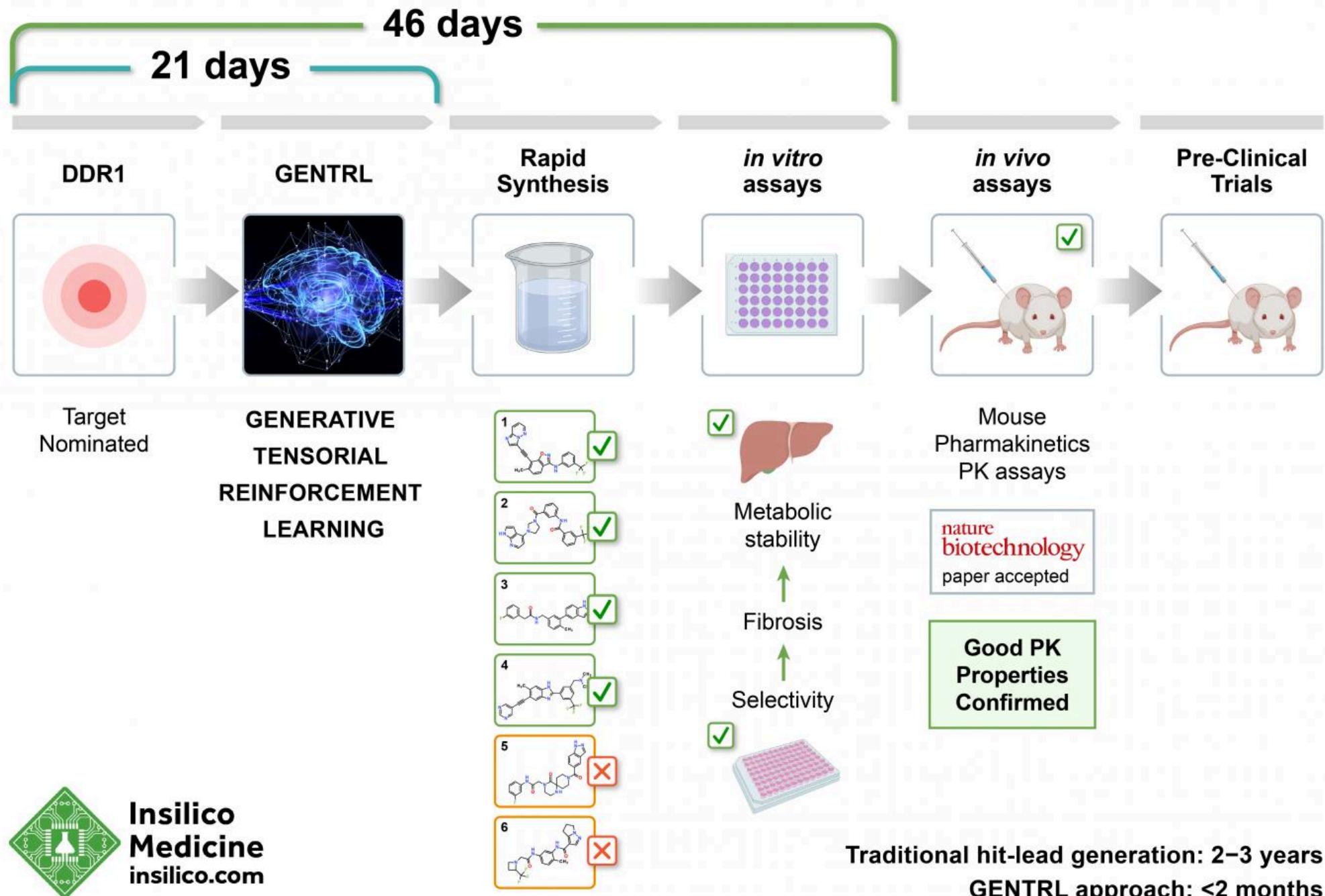


Machine Learning in Drug Discovery



INSILICO MEDICINE
insilico.com

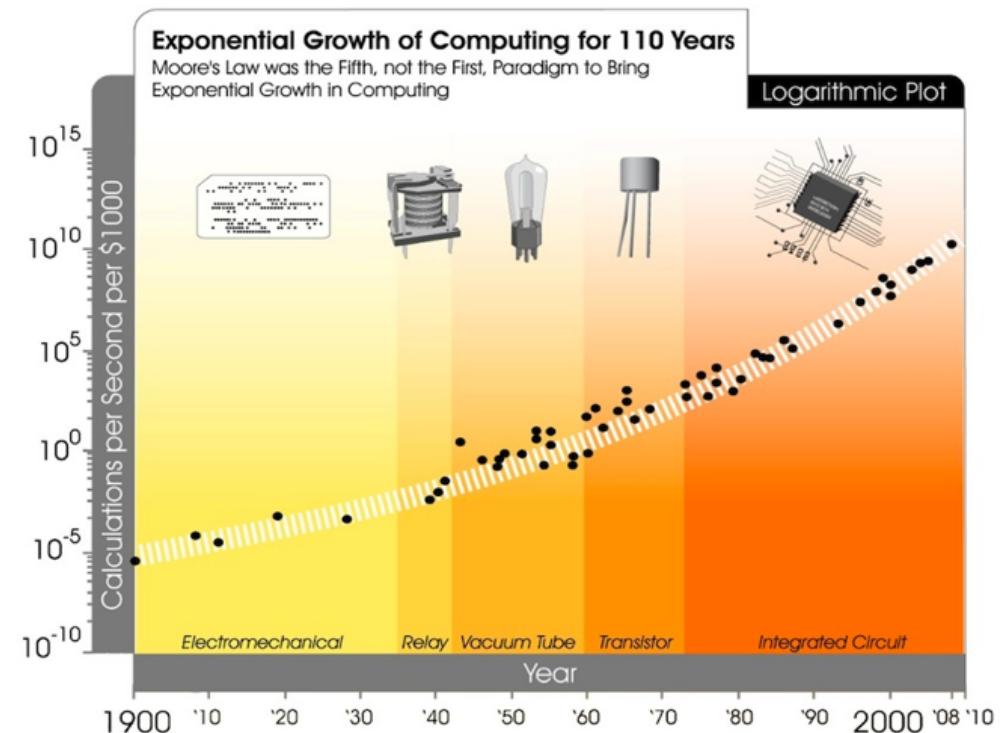
DEEP LEARNING ENABLES RAPID IDENTIFICATION OF POTENT DDR1 KINASE INHIBITORS



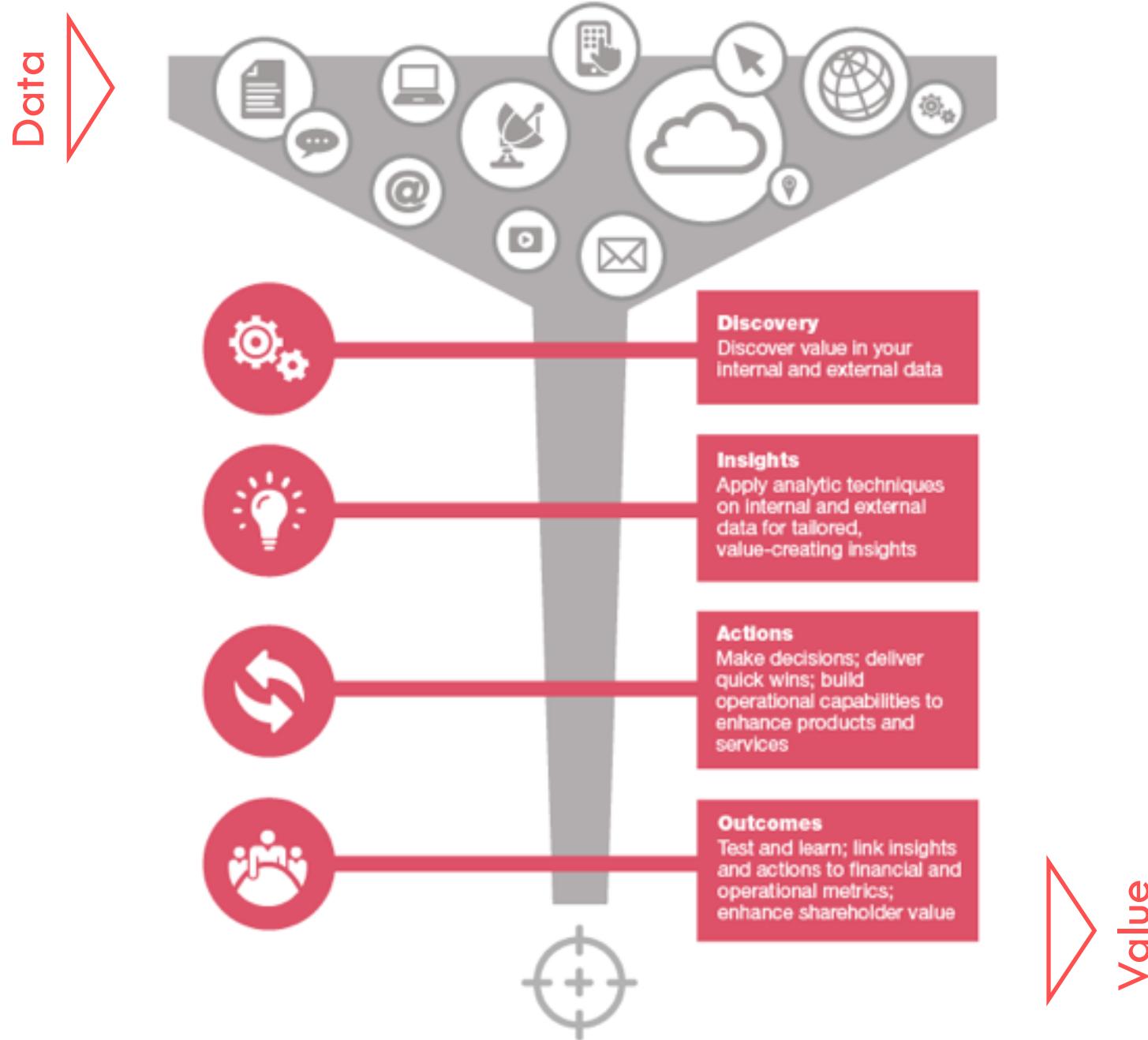
Insilico
Medicine
insilico.com

Driving Forces

Data Deluge, Faster Computing & Cheaper Storage



General Idea



Protocol

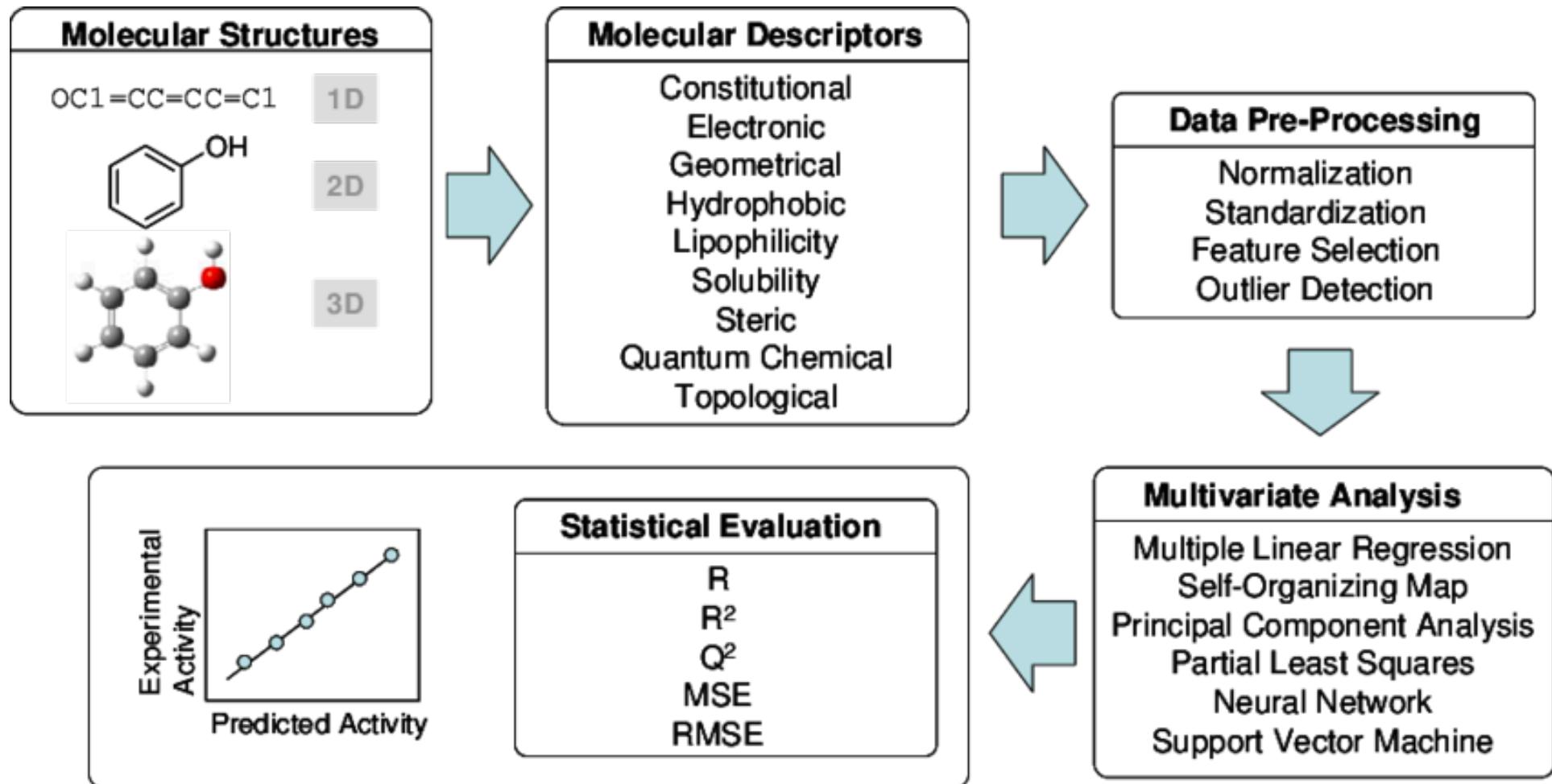


Step By Step

1. Identify Problem
2. Data Mining
3. Pre-processing
4. Exploratory Data Analysis
5. Feature Engineering
6. Predictive Modelling
7. Data Visualisation



Predicting biological activity based on Molecular Descriptors



Program

Machine Learning - General Concept

Data Mining

Pre-processing

Exploratory Data Analysis

Predictive Modelling

Data Visualisation

Data Mining

Collecting the relevant data to your problem.

(Chemical and Biological) Data can be:

- curated and stored in **databases / repositories**
- **manual curation**

Common programs retrieve data from databases:



Queries are usually the language used to access the relevant information from database(s)



Public / Private Repositories

ZINC Substances Catalogs Tranches Biological More

ZINC15

Welcome to ZINC, a free database of commercially-available compounds for virtual screening. ZINC contains over 230 million purchasable compounds in ready-to-dock, 3D formats. ZINC also contains over 750 million purchasable compounds you can search for analogs in under a minute.

DBAASP
v2.702

DATABASE OF ANTIMICROBIAL ACTIVITY
AND STRUCTURE OF PEPTIDES

NEWS

Monomer: 13398

Dimer: 171

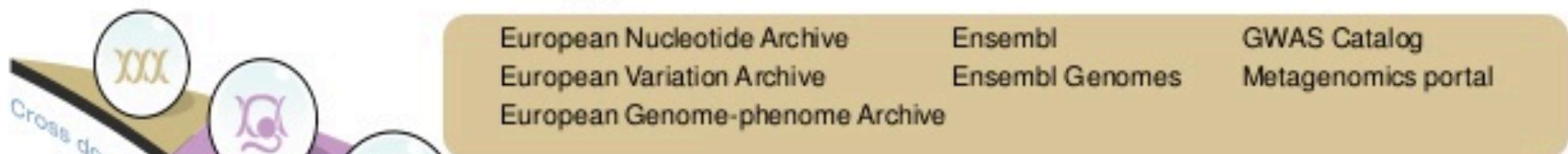
Multi Peptide: 175

The screenshot shows the UniProt homepage. At the top, there is a navigation bar with links for BLAST, Align, Retrieve/ID mapping, and Peptide search. The Peptide search link is circled in red. Below the navigation bar, there is a mission statement: "The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information". To the right of the mission statement, there are four main sections: UniProtKB (UniProt Knowledgebase), UniRef (Sequence clusters), UniParc (Sequence archive), and Proteomes. Each section has a small icon and a brief description. Below these sections, there are two rows of supporting data: Literature citations, Cross-ref. databases, Taxonomy, Diseases, and Subcellular locations, Keywords.

The screenshot shows the MarinLit homepage. At the top, there is a header with the text "MarinLit" and "A database of the marine natural products literature". To the right of the header, there is a "Number of Articles" button with options 3, 4, 0, 6, 1 and a "Last Update : 10-10-2019" link. Below the header, there is a large image of a coral reef. To the right of the image, there is a text block about the database: "MarinLit is a database dedicated to marine natural products research. The database was established in the 1970s by Professors John Blunt and Murray Munro at the University of Canterbury, New Zealand. It was designed as an in-house system to fulfil the needs of the University of Canterbury Marine Group and has evolved to contain unique searchable features and powerful dereplication tools. The extremely comprehensive range of data contained along with these powerful features makes MarinLit the database of choice for marine natural products researchers." At the bottom left, there is a "Get started with MarinLit" button with links to "Introduction to the database", "How to access", and "Contact us". At the bottom right, there is a "Sign in" button with the text "with your subscriber username and password or via your home institution".

Data resources at EMBL-EBI

Genes, genomes & variation



Gene & protein expression

ArrayExpress
Expression Atlas **PRIDE**

Literature & ontologies

Europe PubMed Central
Gene Ontology
Experimental Factor
Ontology

Reactions, interactions & pathways

IntAct Reactome MetaboLights

Protein sequences, families & motifs

InterPro Pfam UniProt

Molecular structures

Protein Data Bank in Europe
Electron Microscopy Data Bank

Chemical biology

ChEMBL ChEBI

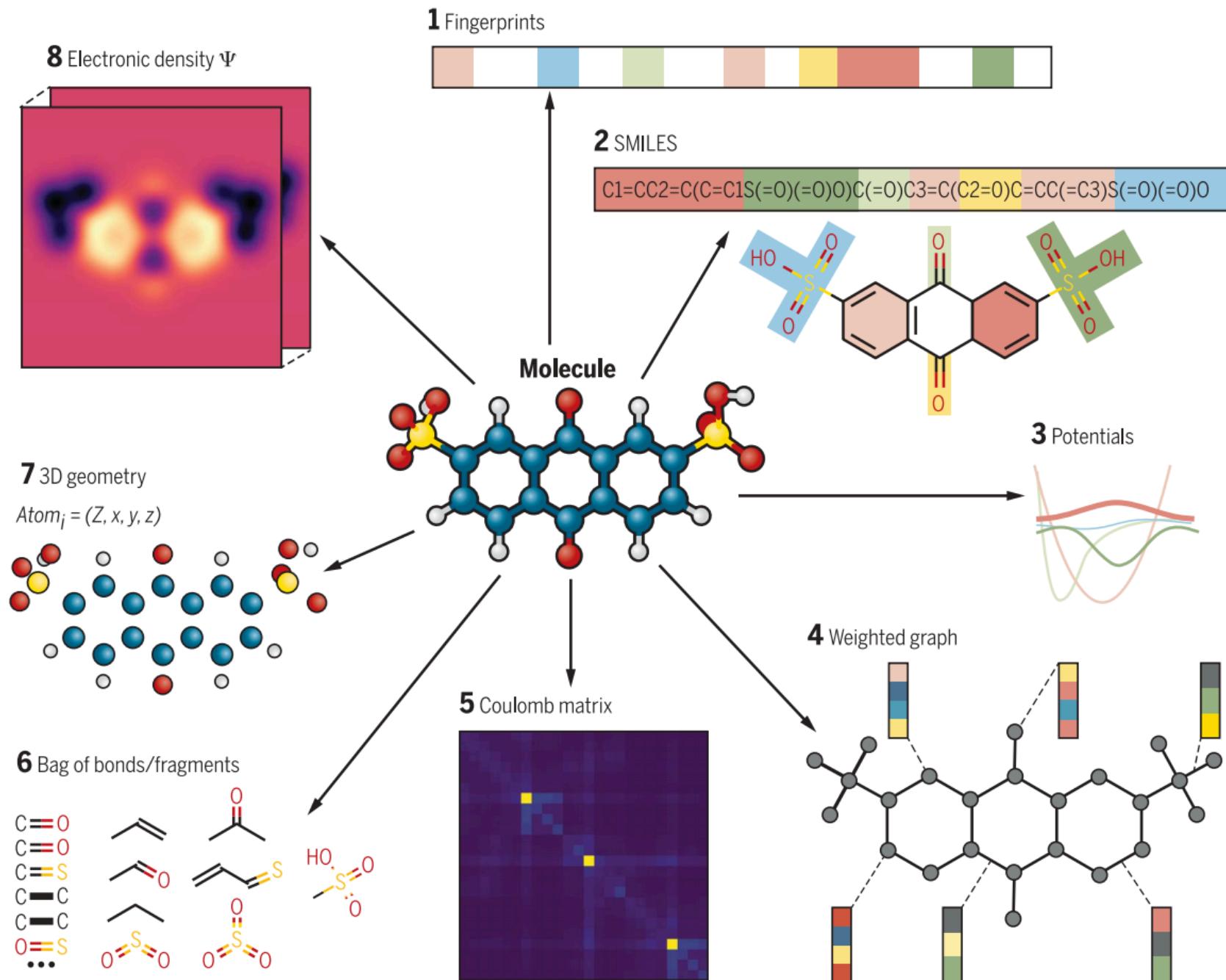
Systems

BioModels Enzyme Portal BioSamples

Digitalization



Encoding Molecules



Encoding Molecules

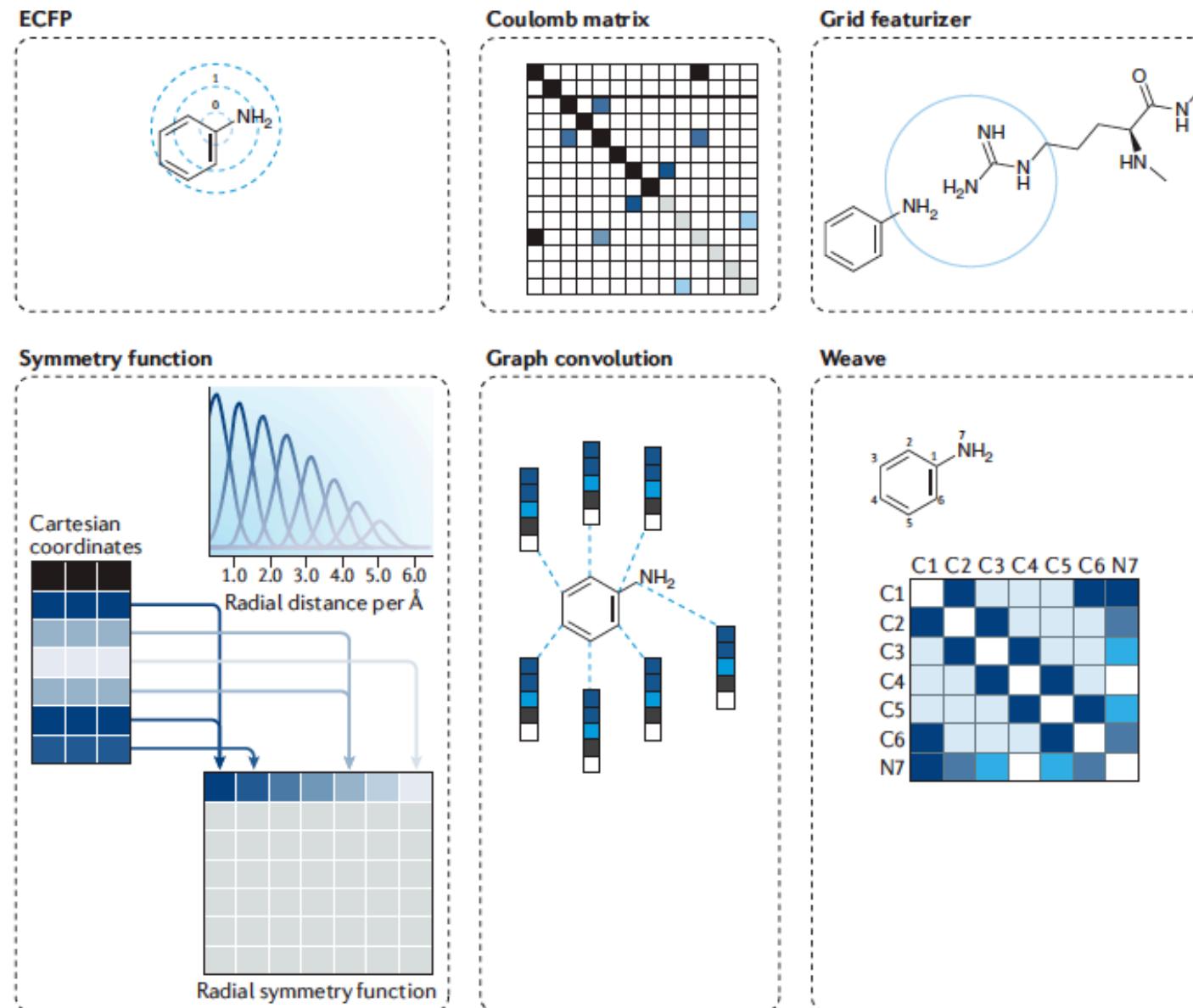


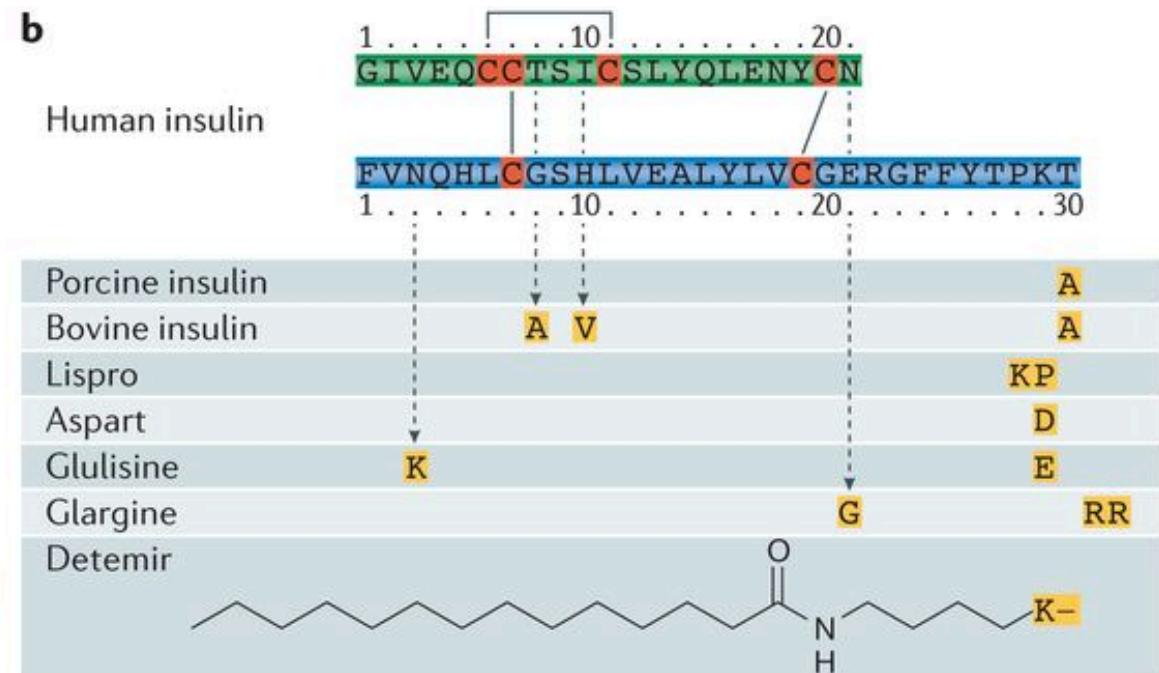
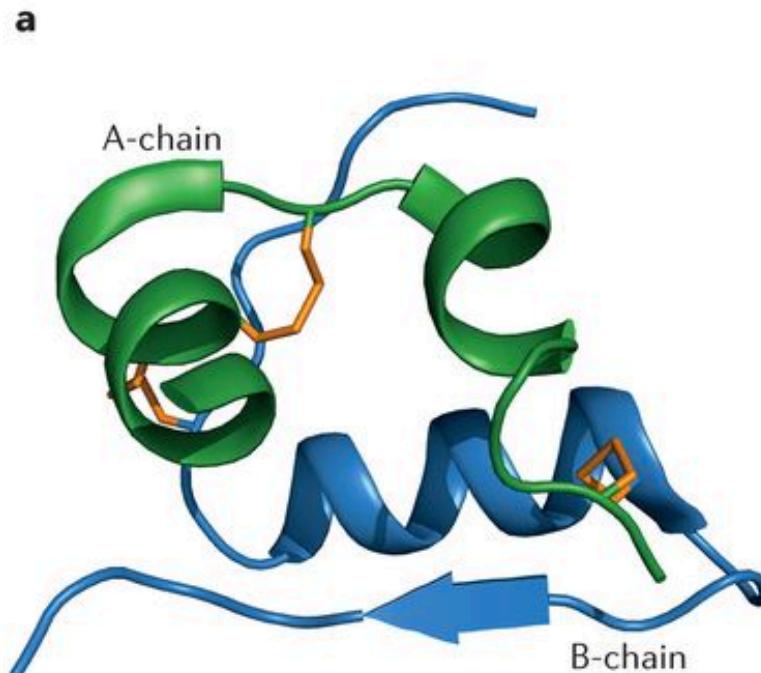
Fig. 3 | The challenges of compound structure representation in machine learning models.

Encoding Proteins

A protein P is represented by the linear chain given by its collapsed graph at amino acid level. A reduced molecular graph representation $G(V,E,C)$ known as a *string signature* where;

V: residues a , E: contiguous in sequence, C: amino acid type

$$^h\sigma(P) = \sum_{a \in A} ^h\sigma(a)$$



Encoding Images

0	0	0	0	0	0	0	...
0	156	155	156	158	158	158	...
0	153	154	157	159	159	159	...
0	149	151	155	158	159	159	...
0	146	146	149	153	158	158	...
0	145	143	143	148	158	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	0	...
0	167	166	167	169	169	169	...
0	164	165	168	170	170	170	...
0	160	162	166	169	170	170	...
0	156	156	159	163	168	168	...
0	155	153	153	158	168	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	0	...
0	163	162	163	165	165	165	...
0	160	161	164	166	166	166	...
0	156	158	162	165	166	166	...
0	155	155	158	162	167	167	...
0	154	152	152	157	167	167	...
...

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

+

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

+

-25				...
				...
				...
				...
...

Bias = 1

Output

Program

General Concept

Data Mining

Pre-processing

Exploratory Data Analysis

Predictive Modelling

Data Visualisation

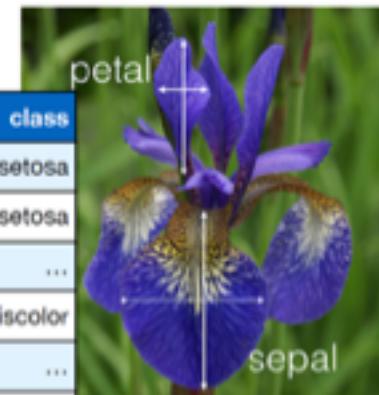
Nomenclature

IRIS

<https://archive.ics.uci.edu/ml/datasets/Iris>

Instances (samples, observations)

	sepal_length	sepal_width	petal_length	petal_width	class
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
...
50	6.4	3.2	4.5	1.5	versicolor
...
150	5.9	3.0	5.1	1.8	virginica



Features (attributes, dimensions)

Classes (targets)

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1D} \\ x_{21} & x_{22} & \cdots & x_{2D} \\ x_{31} & x_{32} & \cdots & x_{3D} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{ND} \end{bmatrix}$$

$$\mathbf{y} = [y_1, y_2, y_3, \dots, y_N]$$

Feature Engineering

Process of creating, modifying or maintaining features making them more appropriate for their use in modelling and maximum information can be extracted out of them.

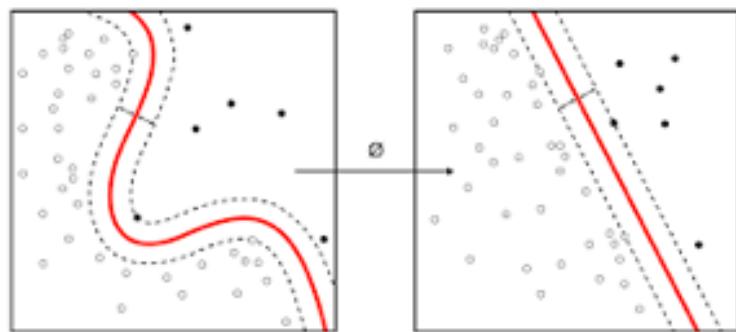


Pros

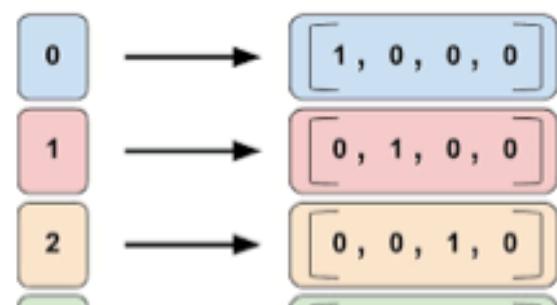
The better the variables that are analysed, the better and more reliable the results will be.
"Garbage In, Garbage Out"

Cons

Time-consuming process
In-depth knowledge of features
Compatibility features' appearance / model



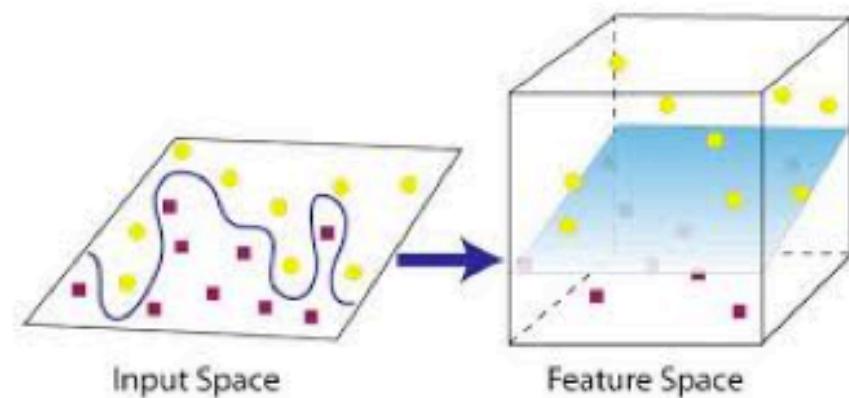
1. TRANSFORMATION



3. CONSTRUCTION

$$\frac{x - X_1}{\max(X) - \min(X)}$$

2. SCALING

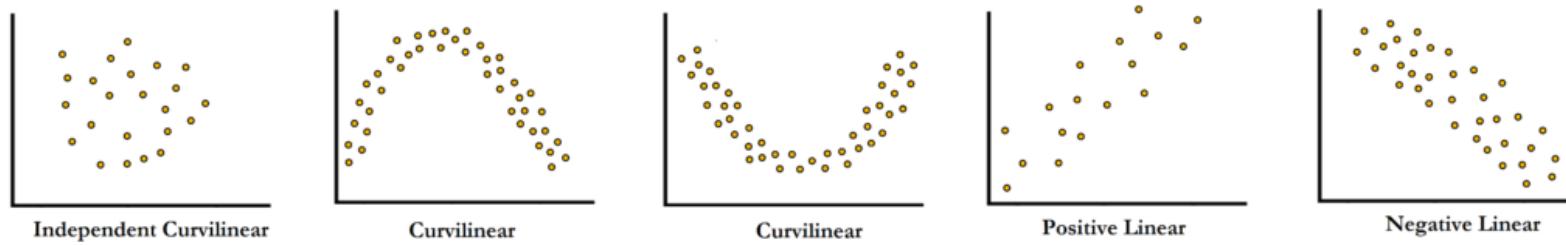


4. REDUCTION

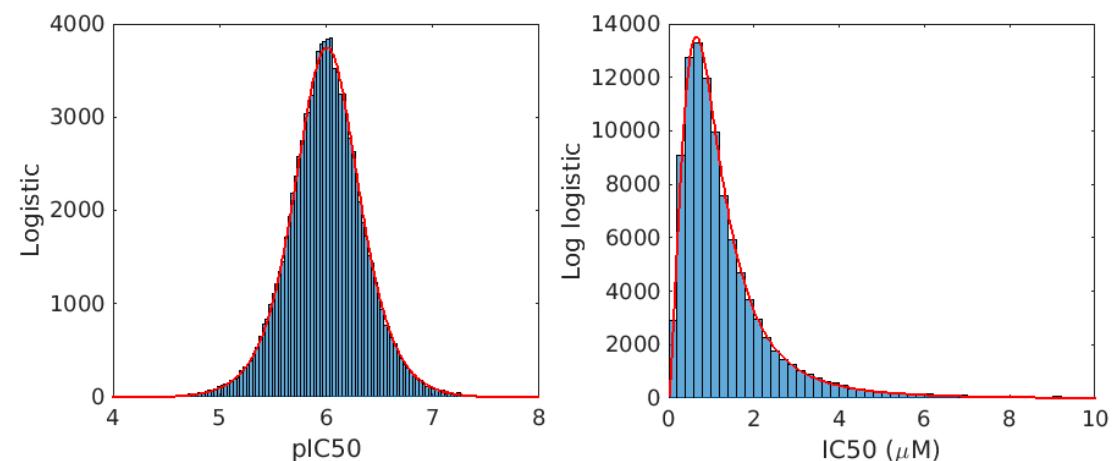
1. Transformation

Replacement of a variable by a function that helps (often) to transform:

- non-linear relationships into linear relationships
- non-normal (skewed) distributions into normal distributions



1. Logarithm
2. Square / Cube Root
3. Arcsine
4. Box-Cox



2. Scaling

Many ML algorithms require feature scaling to prevent the model from giving greater weightage to certain attributes (variables) as compared to others.

- **Models** e.g. K-means requires Euclidean distance matrix
- **Gradient descent** (optimisation algorithm) - scaling reduce computing time
- **Feature Extraction** e.g. scaling prevents emphasising variables in PCA or LDA

1. Min-Max Scaling (**Normalisation**)

Scales ranges: [0, 1] or [-1, 1]

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

$$x_{normalized}_{[-1,1]} = \frac{x_{current} - ((X_{max} + X_{min})/2)}{(X_{max} - X_{min})/2}$$

2. Z score Normalisation (**Standardisation**)

Normal distribution properties

with mean = 0, standard deviation = 1

$$z = \frac{x - \mu}{\sigma}$$

3. Construction

Generating new variables (features) based on existing variables to improve reliability of (linear) predictive models or to reduce noise (random data points)

1. **Binning** e.g. People age, IQ levels correlations
2. **Encoding** e.g. One-hot encoding countries
3. **Derived Variables** e.g. Key Performance Indicators (KPI)

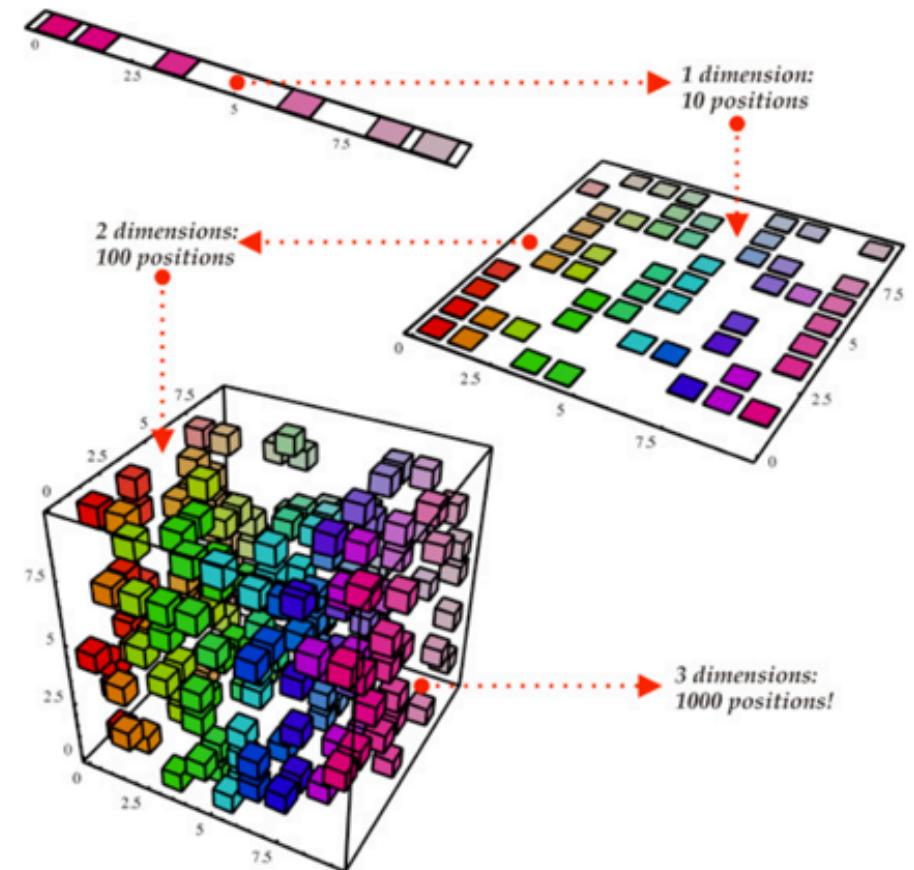
Derived Variables / Features



4. Reduction

Dimension(ality) reduction is the process of reducing the number of random variables under consideration, by obtaining a set of principal variables. It can be divided into **feature selection** and **feature extraction**.

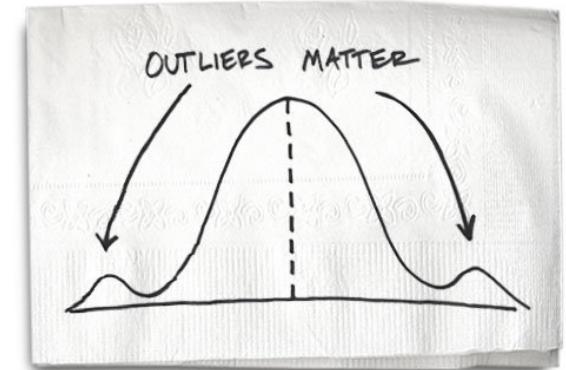
1. Unsupervised methods
2. Supervised methods



Outlier Detection

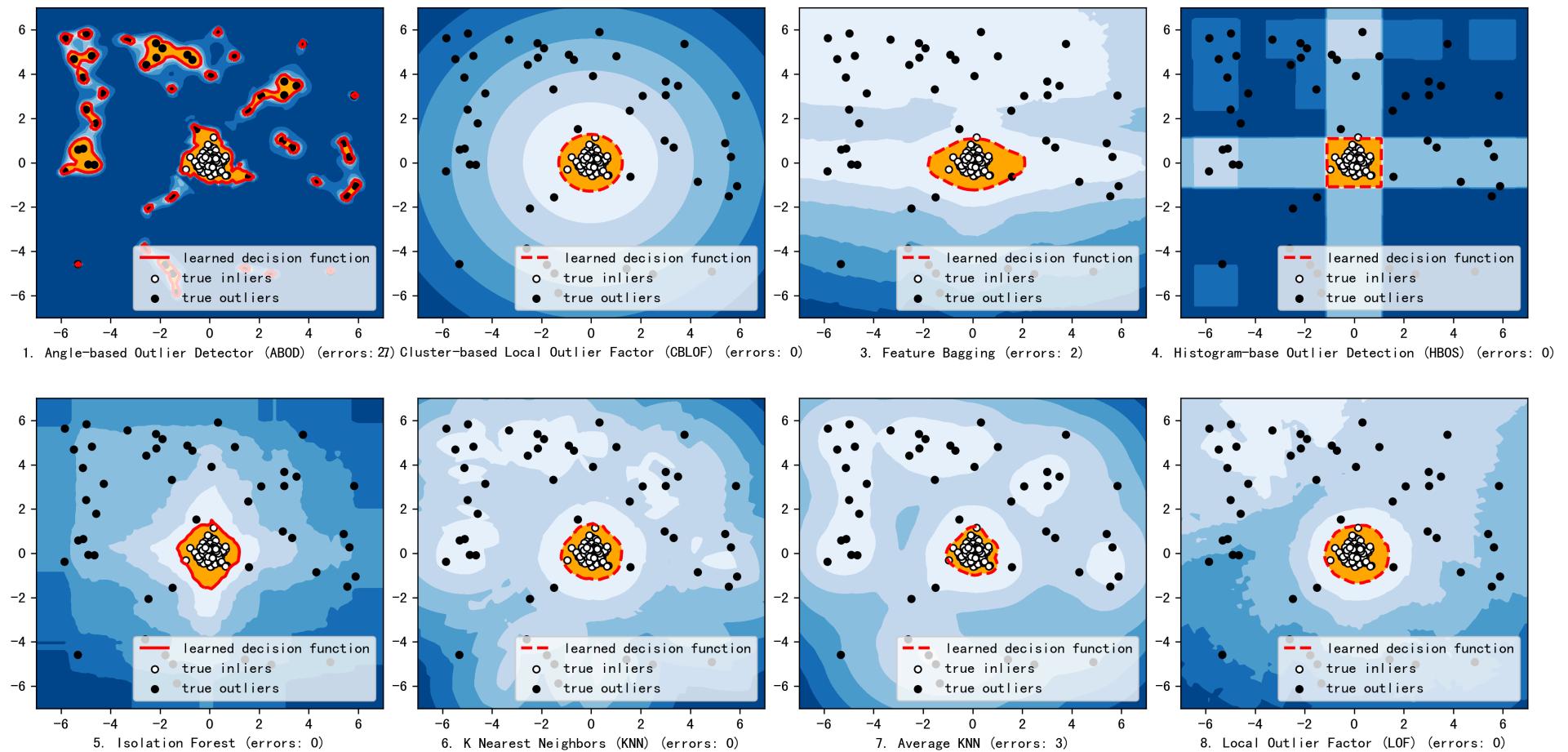
Outliers are extreme values that deviate from other observations on data , they may indicate a variability in a measurement, experimental errors or a novelty. In other words, an outlier is an observation that diverges from an overall pattern on a sample.

1. **Univariate outliers** can be found when looking at a distribution of values in a single feature space.
2. **Multivariate outliers** can be found in a n-dimensional space (of n-features).
3. **Point outliers** are single data points that lay far from the rest of the distribution.
4. **Contextual outliers** can be noise in data
 - e.g. *punctuation symbols / text analysis*
 - background noise signal / speech recognition*
5. **Collective outliers** can be subsets of novelties in data
 - e.g. *signal that may indicate the discovery of new phenomena*



Detection Methods

- Z-Score or Extreme Value Analysis (**parametric**)
- Probabilistic and Statistical Modelling (**parametric**)
- Linear Regression Models (PCA, LMS)
- Proximity Based Models (**non-parametric**)
- Information Theory Models
- High Dimensional Outlier Detection Methods (high dimensional sparse data)



Program

General Concept

Data Mining

Pre-processing

Exploratory Data Analysis

Predictive Modelling

Data Visualisation

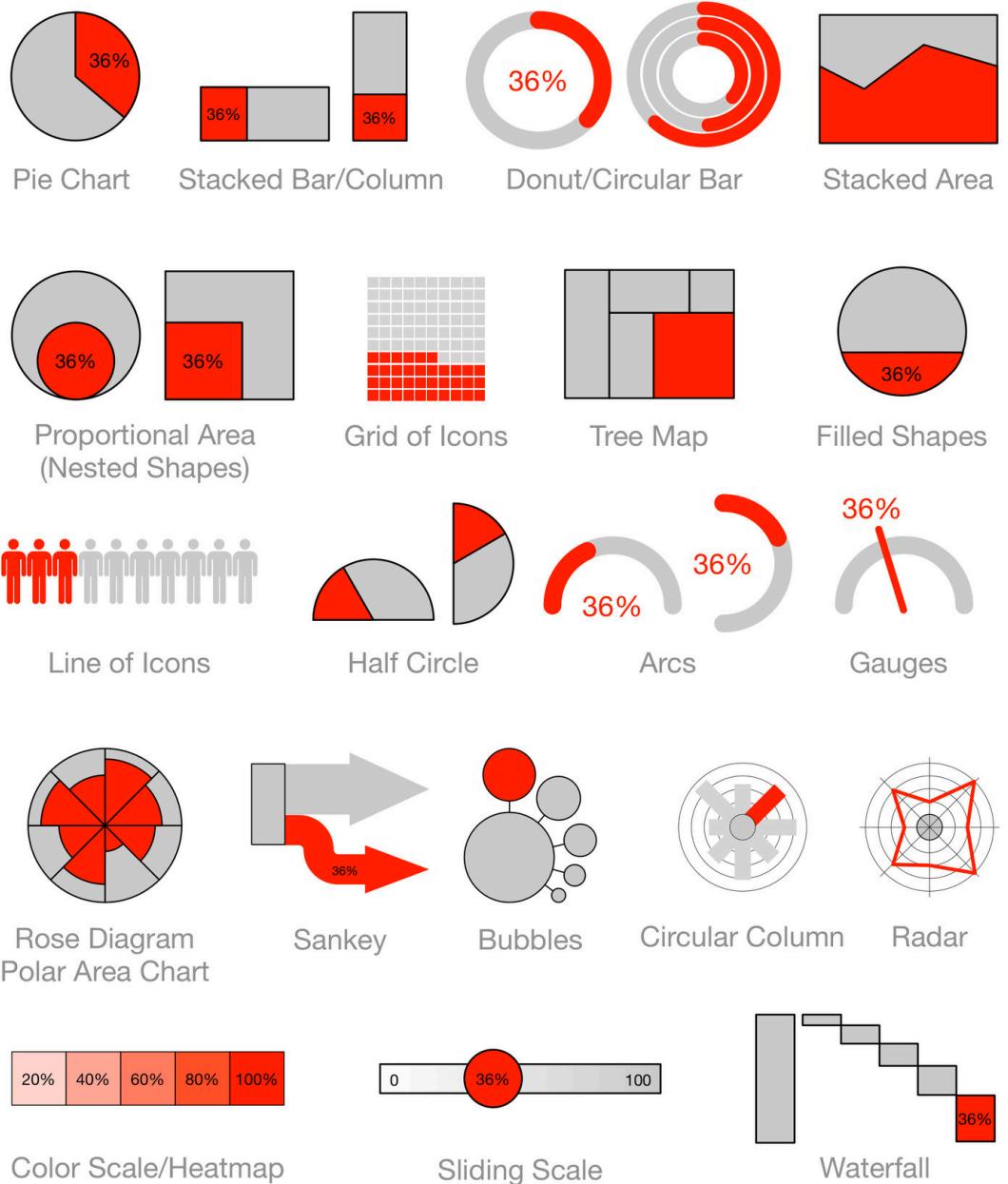
Telling a story

*En el OXXO siempre hay
dos weyes, uno que te
cobra, y otro que te dice
"de aquel lado le cobran"*



Exploratory Data Analysis

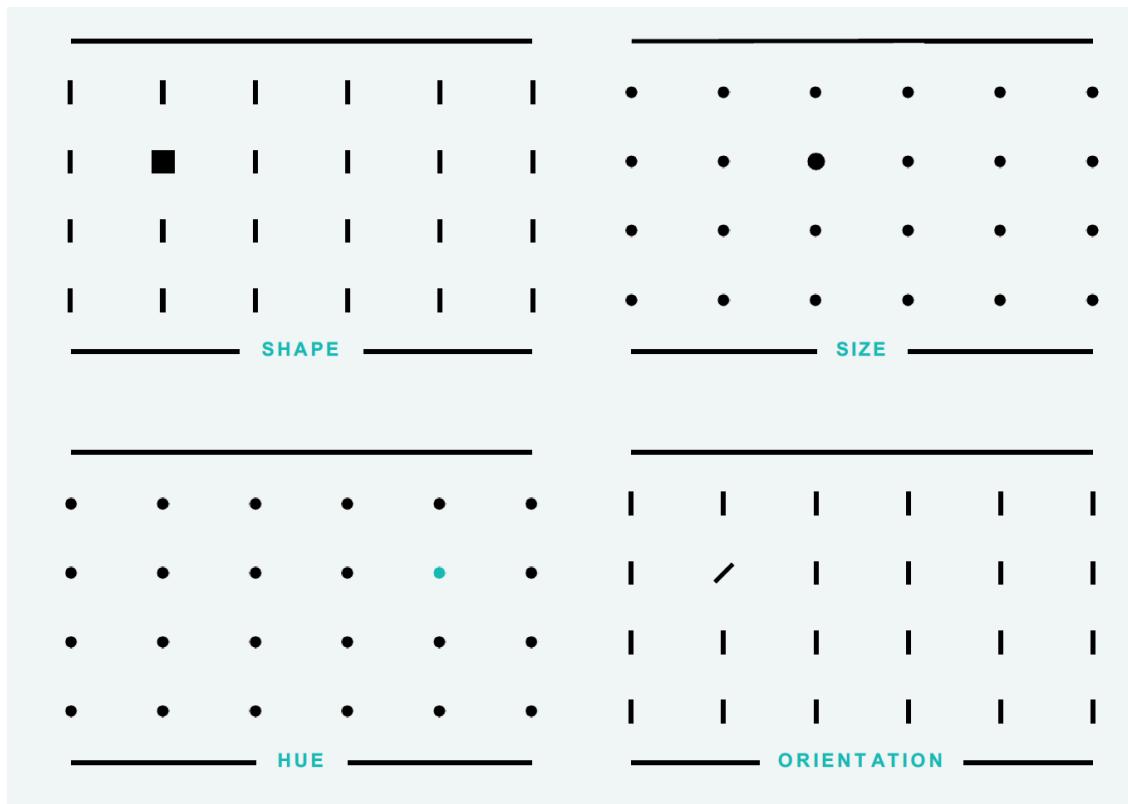
In statistics, **exploratory data analysis** (EDA) is an approach to analyzing **data sets** to summarize their main characteristics, often with visual methods



Our Brains Love Visuals

Through the visual system, the human brain quickly recognizes, stores and recalls images, seamlessly and subconsciously cementing ideas in long-term memory.

Our brain gathers information through pre-attentive processing of visual cues in our environment, which we unconsciously absorb and filter—within 250 milliseconds.



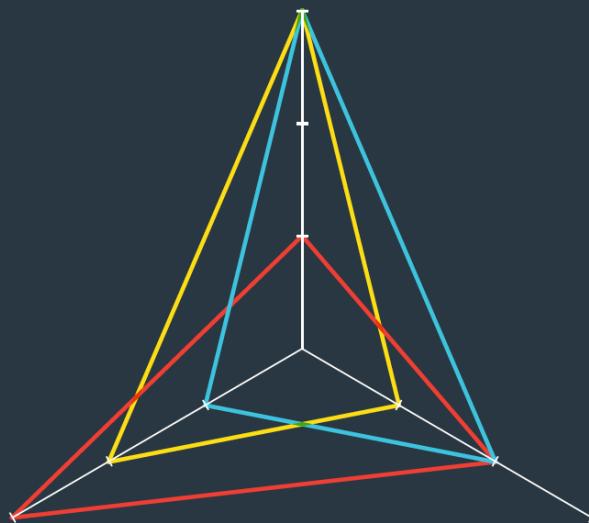
Notice that your eye is naturally drawn to these variations on the left.

VALUE OF VISUALIZATION

-  Academic/Scientific
-  Marketing
-  Editorial

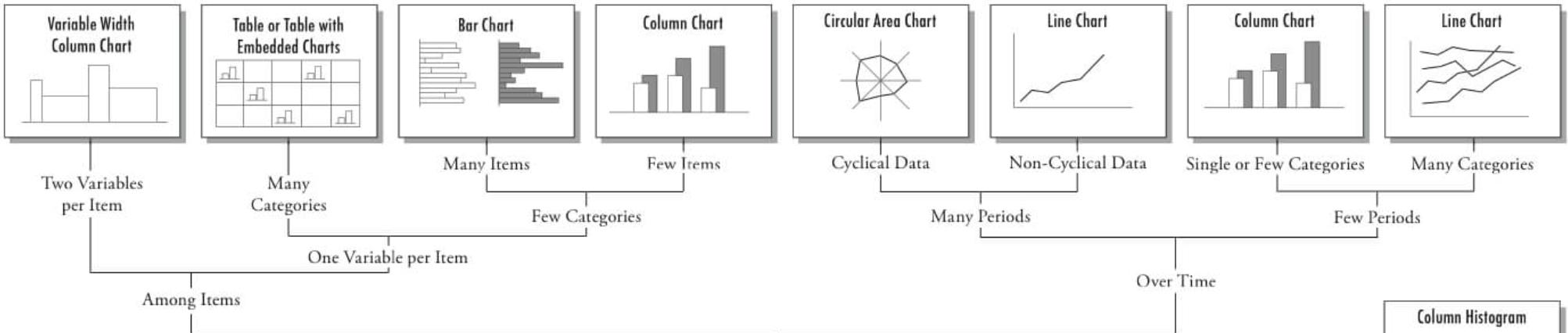
COMPREHENSION The brain is pre-wired to automatically interpret relationships between objects, allowing for instant comprehension with minimal effort. Representing these relationships visually, as opposed to merely describing them, means that your message is understood quickly, clearly and with significantly greater joy.

APPEAL Well-designed information is stimulating, attractive and engaging. These qualities pique interest even before information is processed. Aesthetics are not superficial; they are how you get people's attention.



RETENTION Visualizations trigger us to pull information from our long-term memory, allowing for rapid connections to already stored information, which help to cement the concept in the brain.

Chart Suggestions—A Thought-Starter



Comparison

Relationship

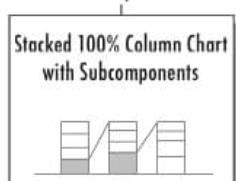
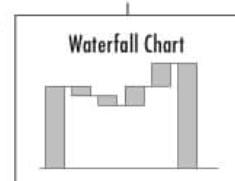
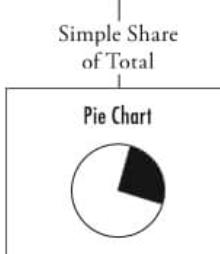
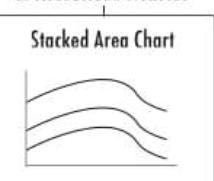
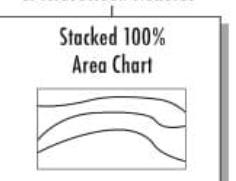
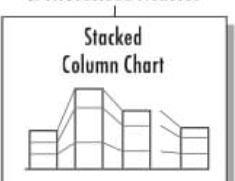
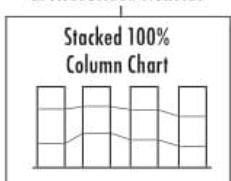
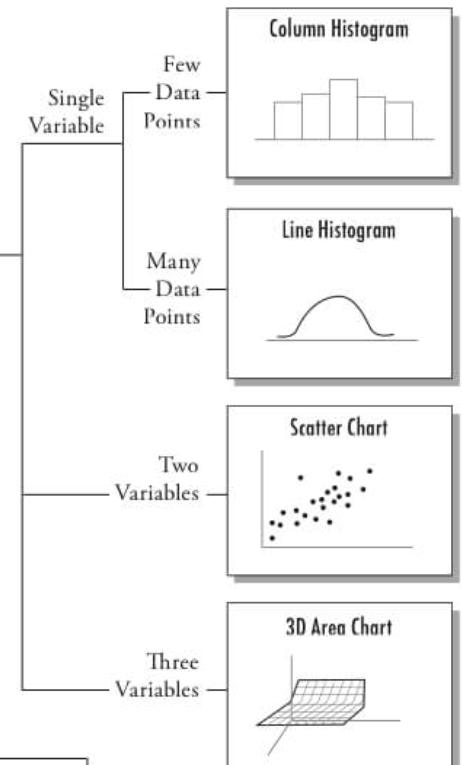
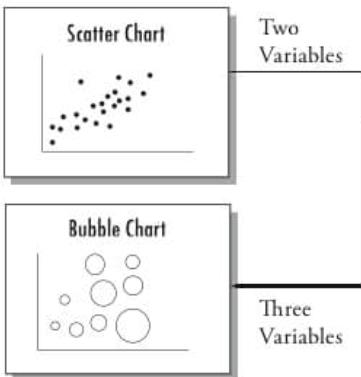
What would you like to show?

Distribution

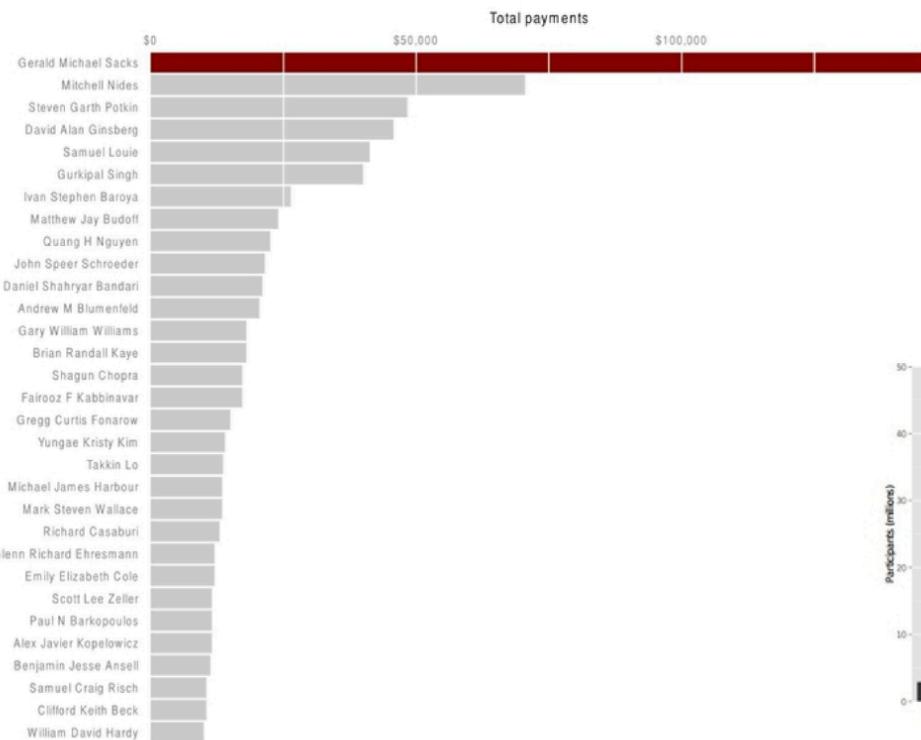
Composition

Changing Over Time

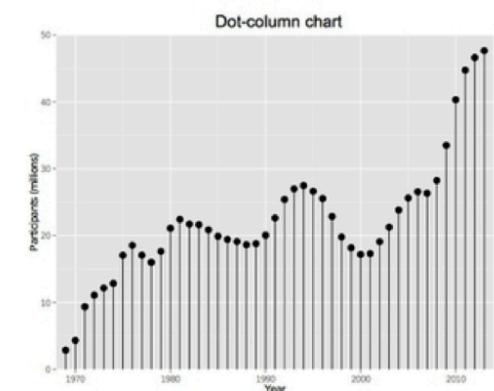
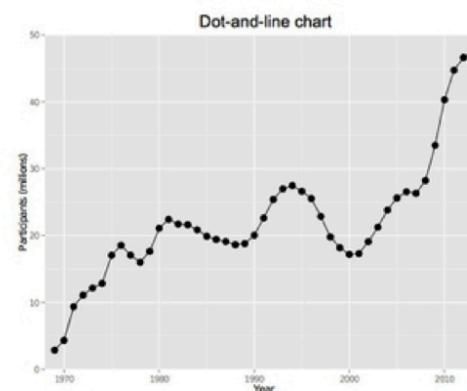
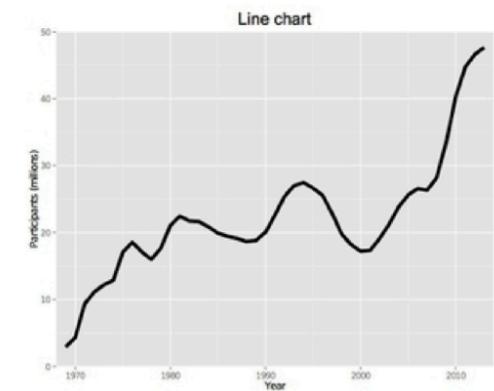
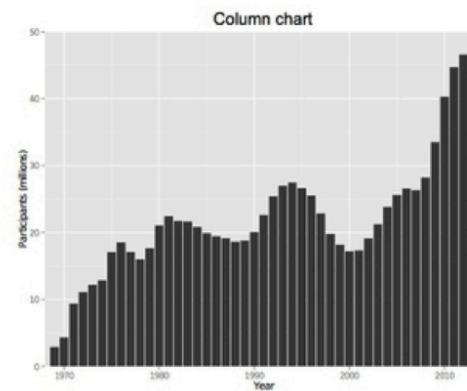
Static



Comparison

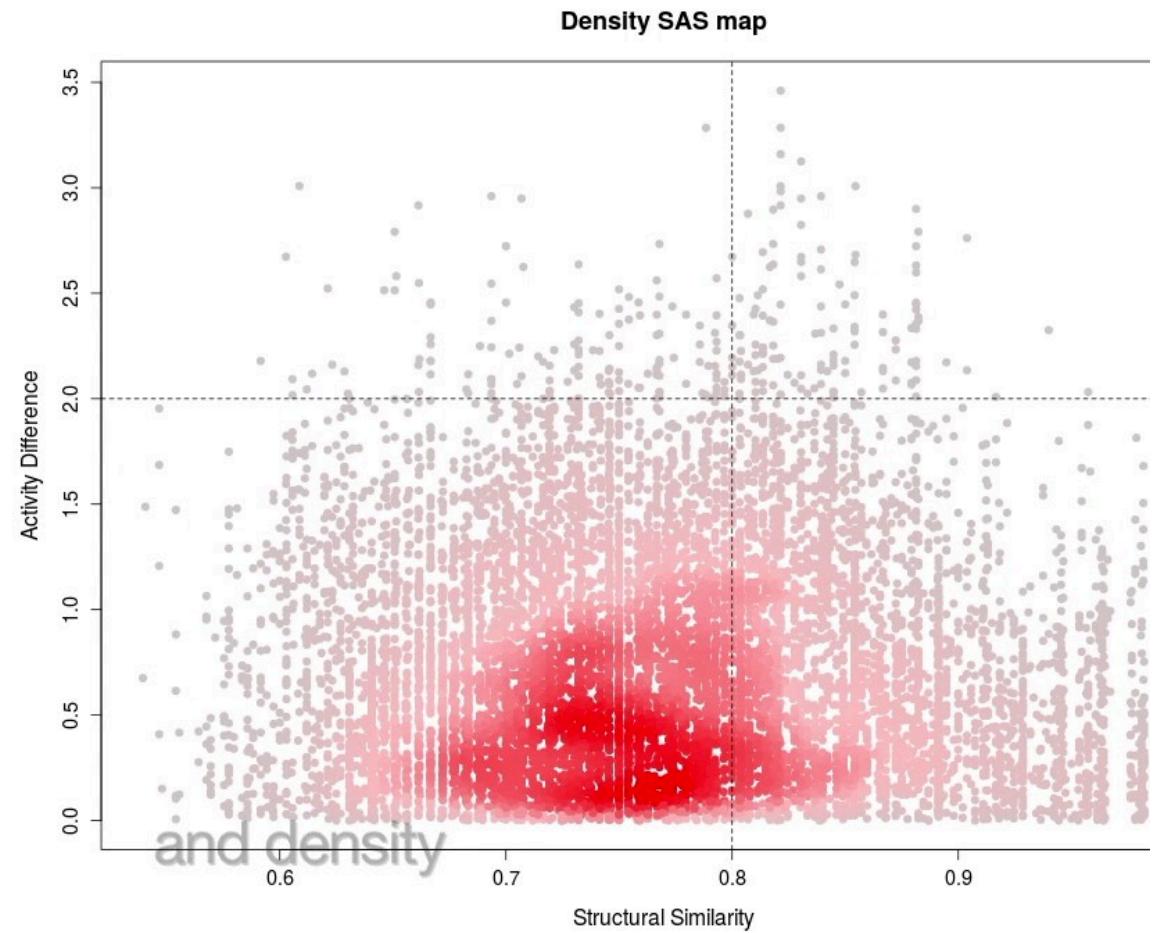
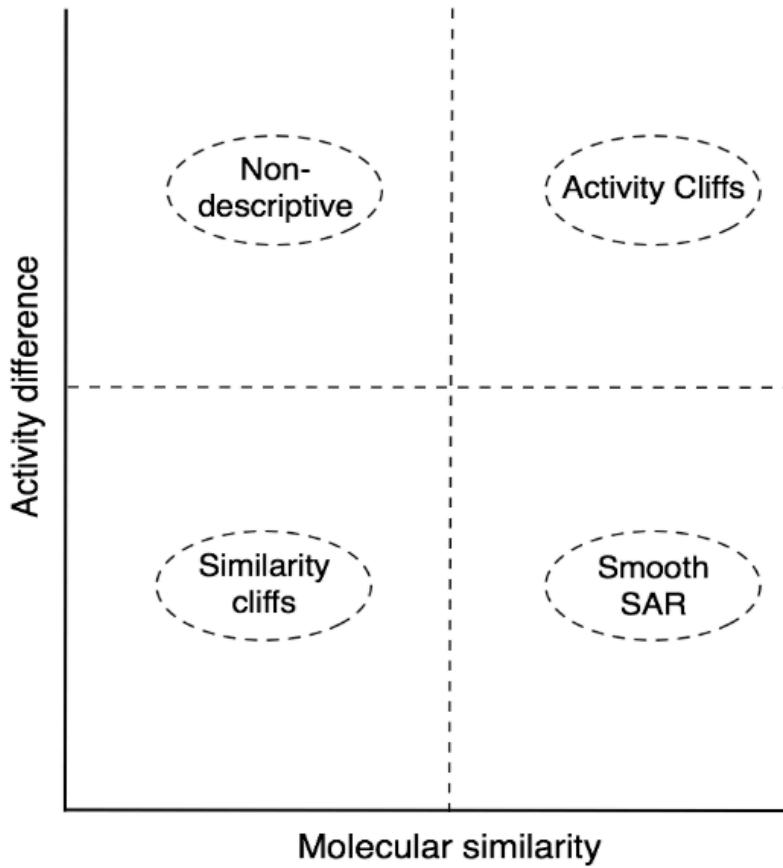


Pfizer paid doctors

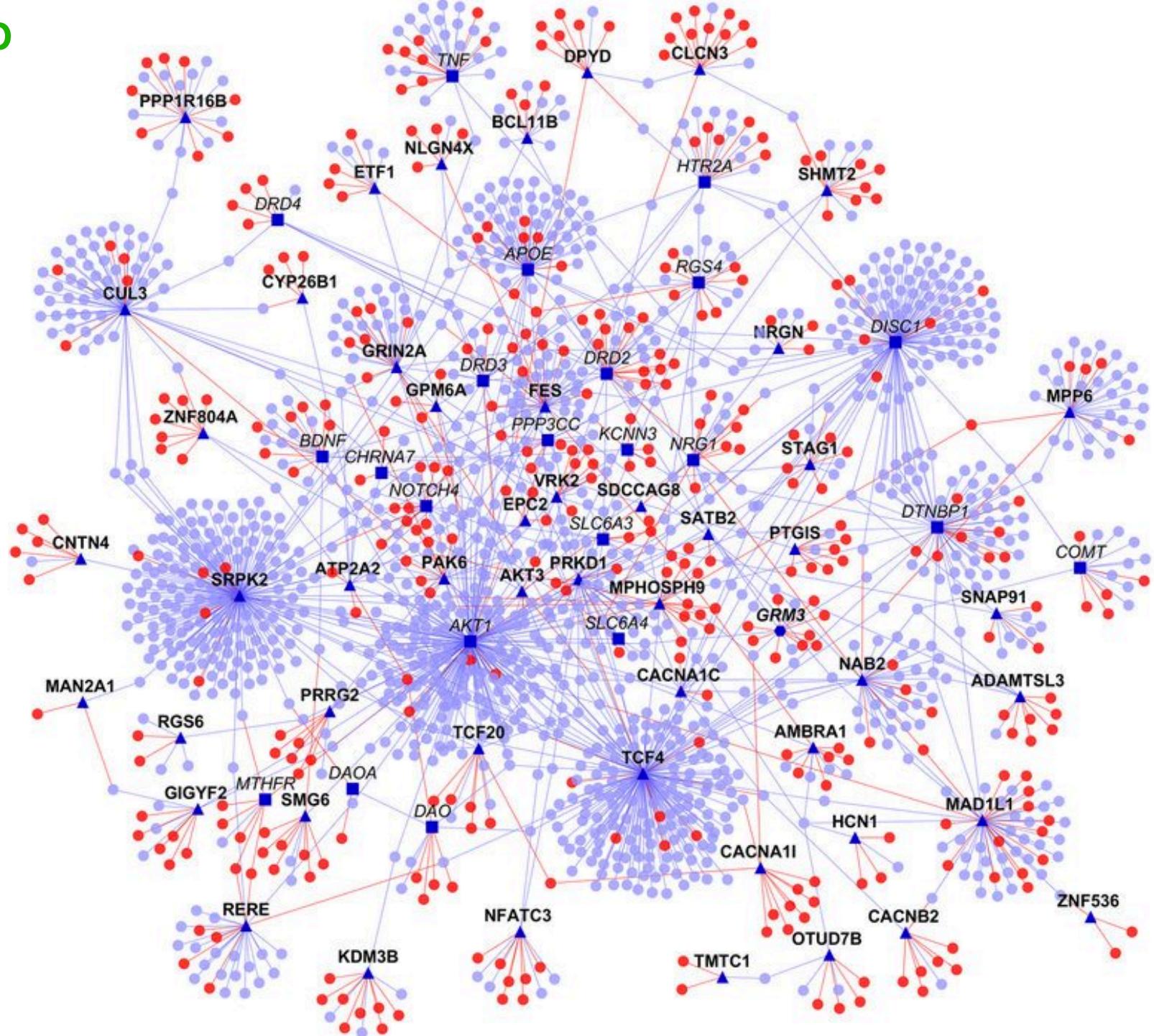


Gvt food stamps 1969-2014

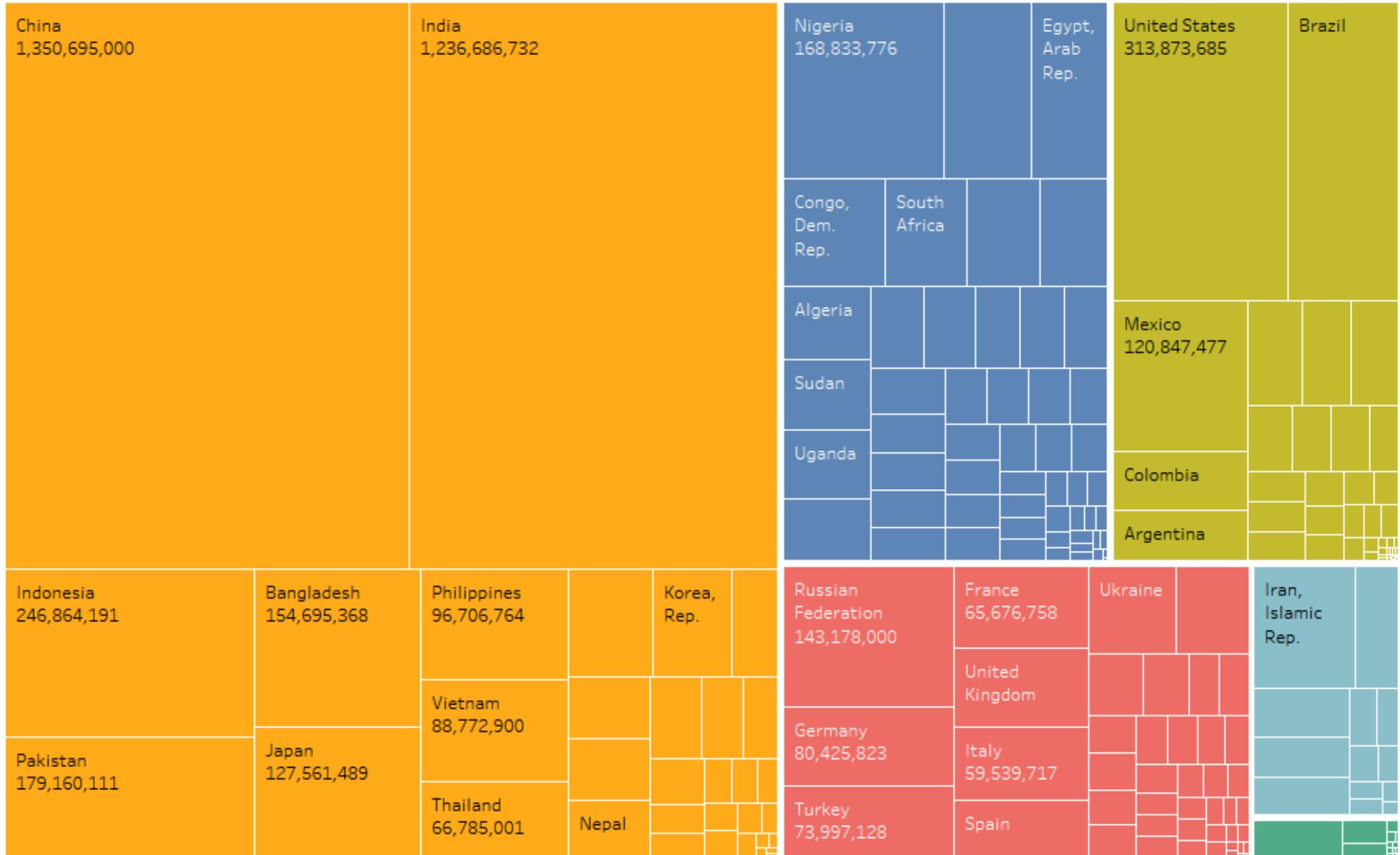
Relationship



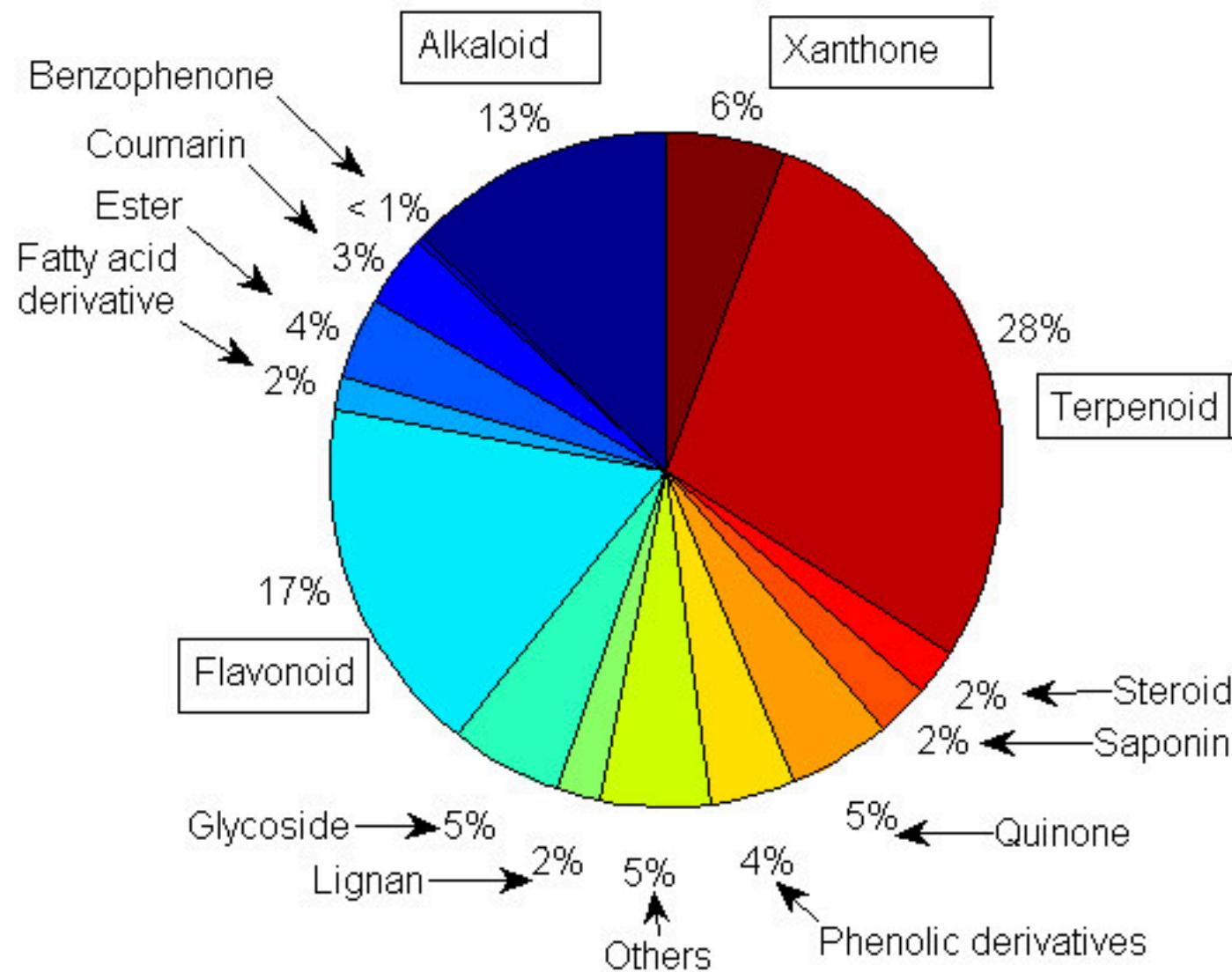
Relationship



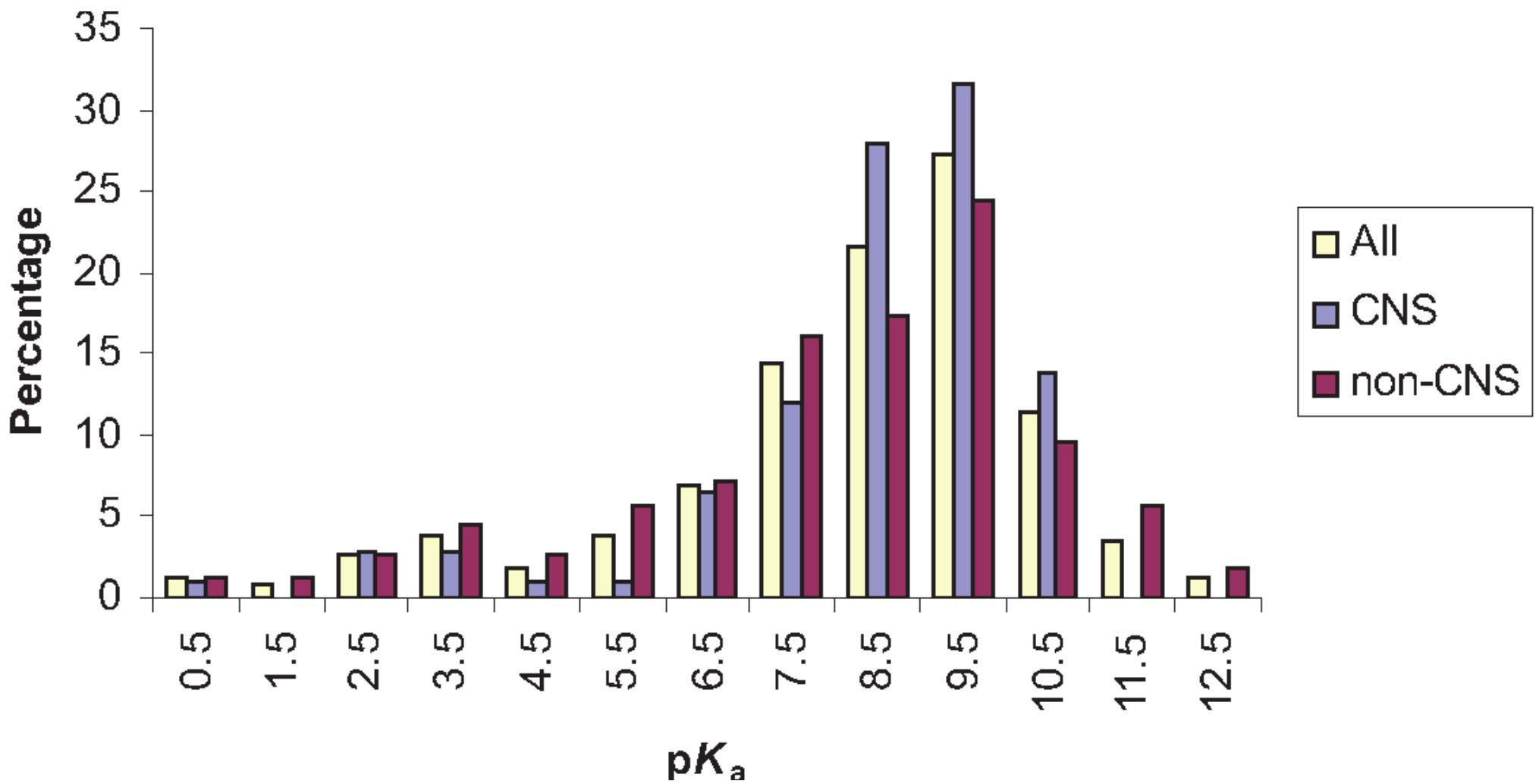
Composition



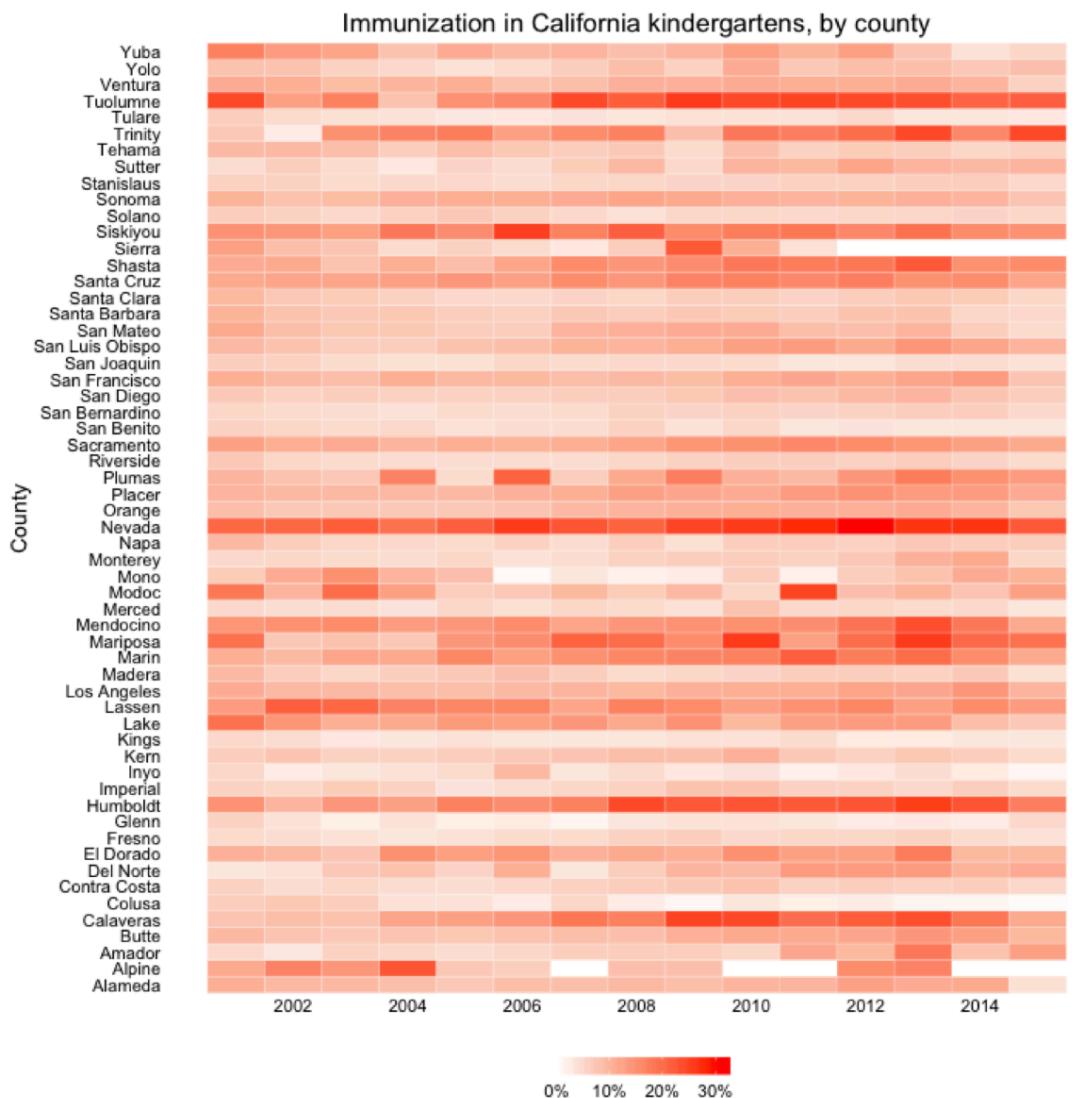
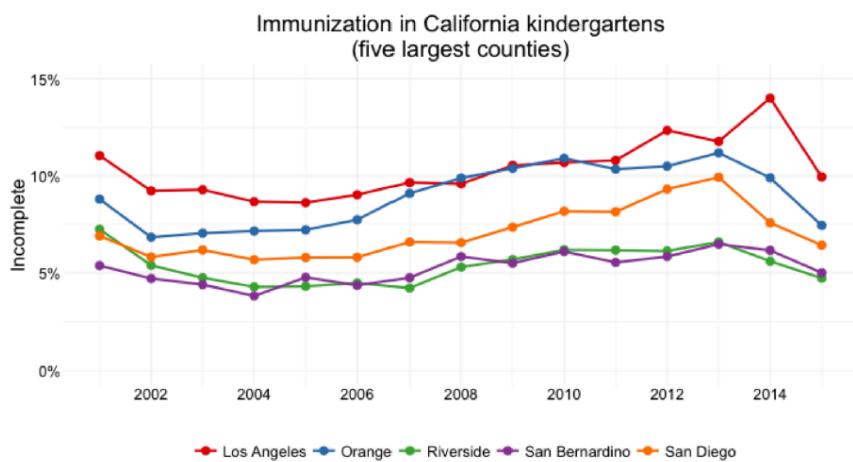
Composition



Distribution

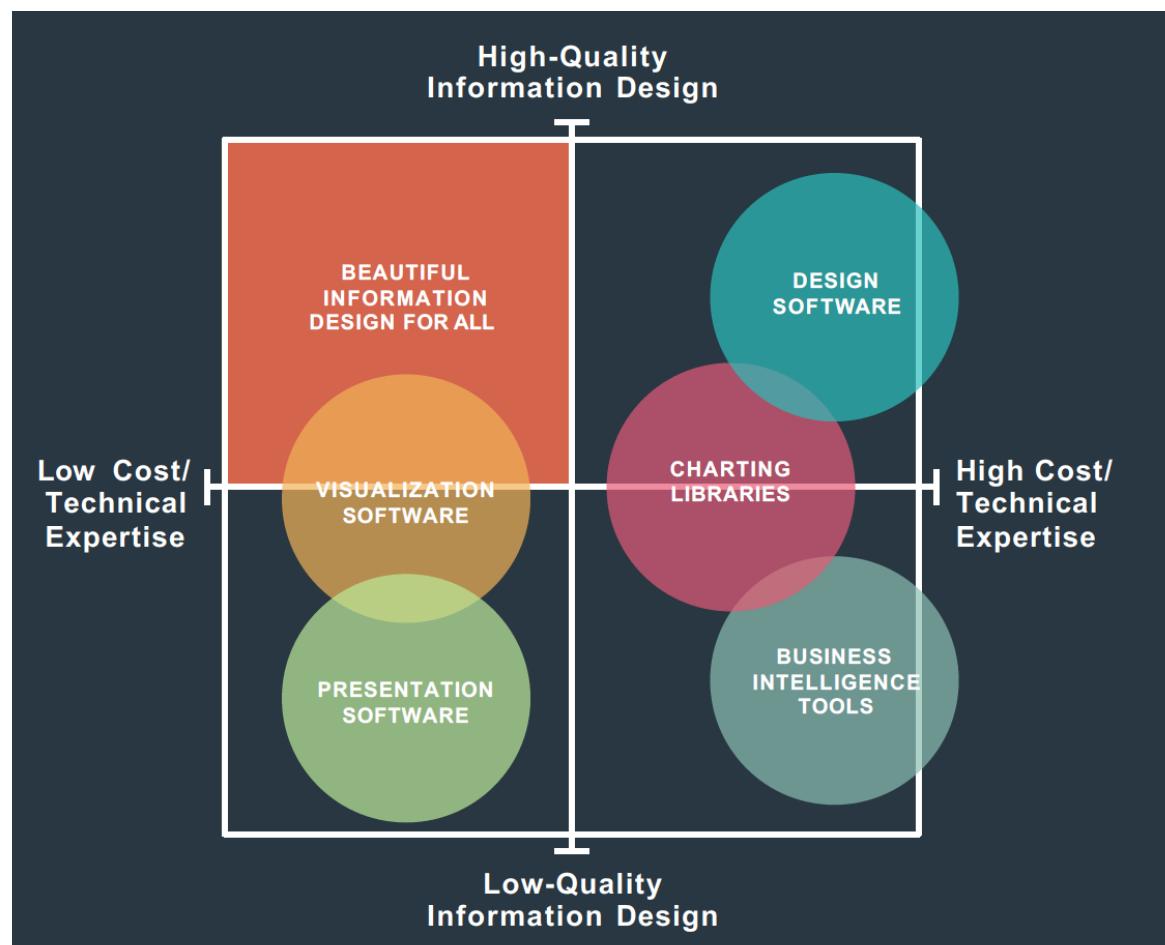


Distribution

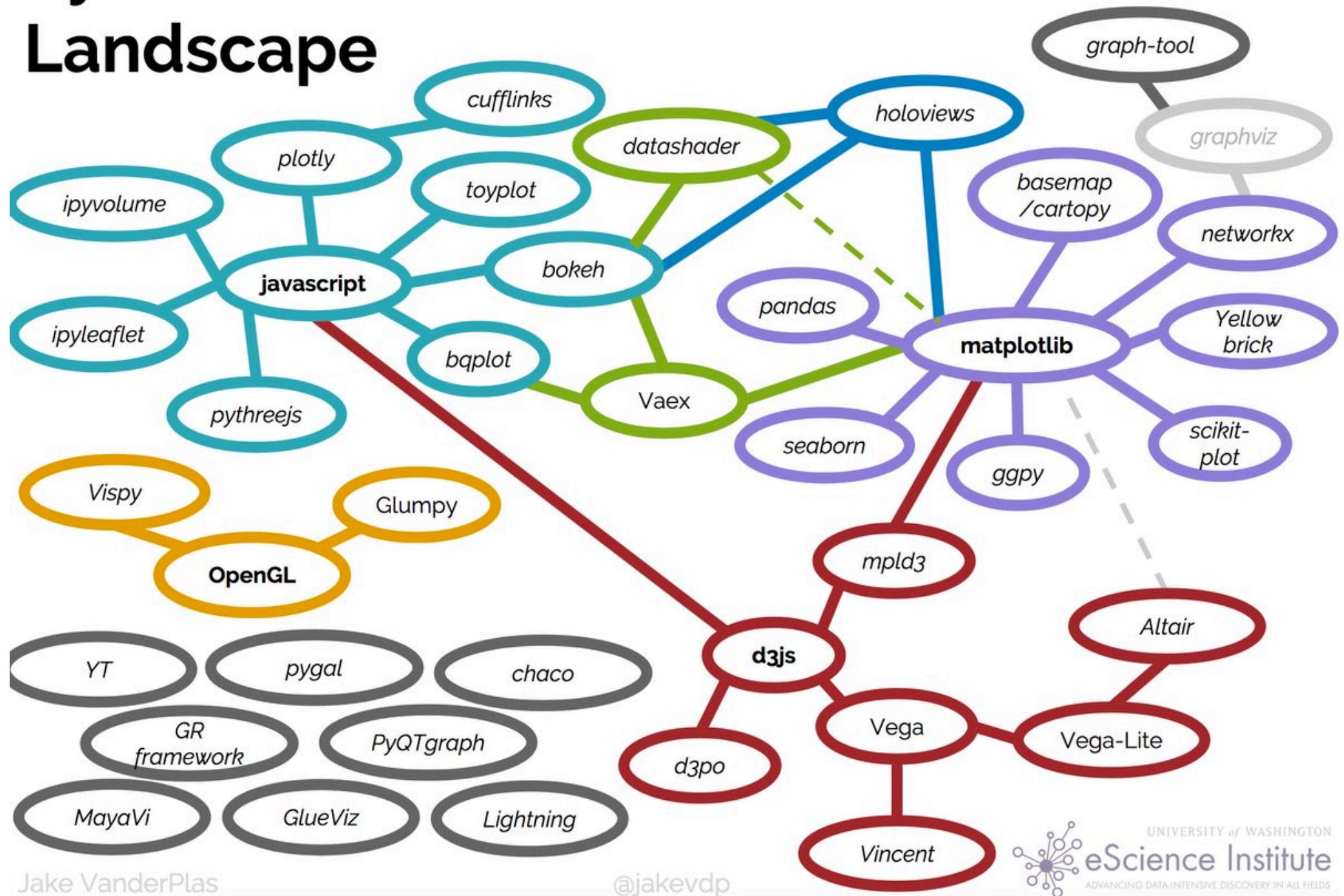


Design Tools

The traditional tools and design programs available are incomplete, challenging to master or limited in their aesthetic. Visualization software is an emerging field, providing better design tools that are easier to use.

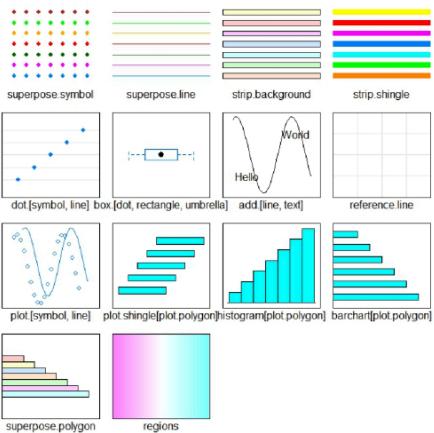


Python's Visualization Landscape

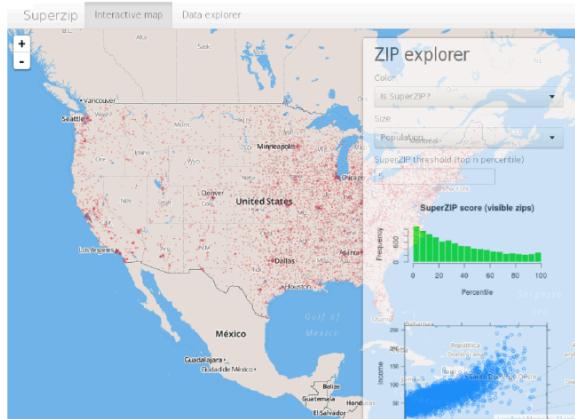


R Data Viz Packages

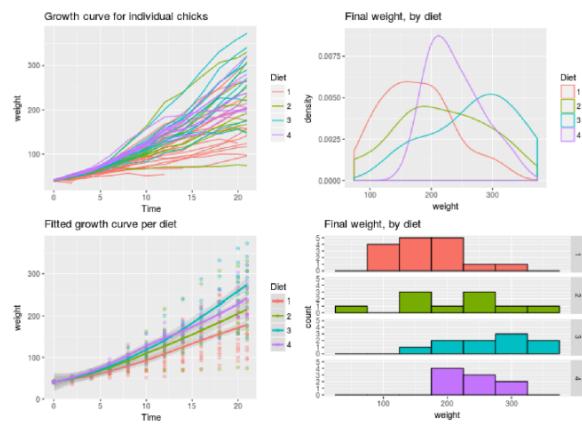
Lattice



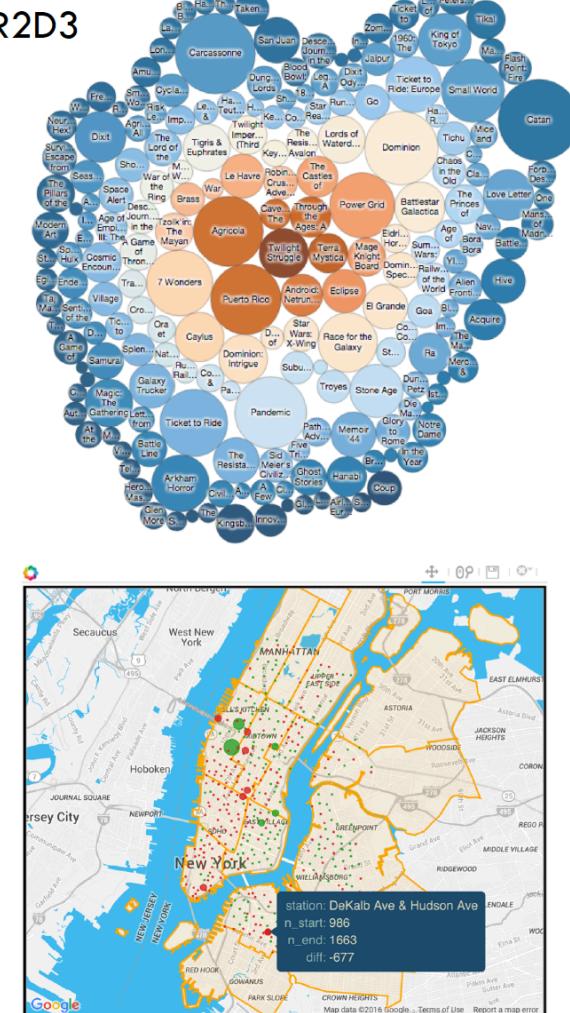
Rshiny



Ggplot2



Rbokeh



Program

General Concept

Data Mining

Pre-processing

Exploratory Data Analysis

Predictive Modelling

Data Visualisation

CHOOSING THE RIGHT ROBOT



STAY TUNED WITH DATA SCIENCE & DATA VIZ



Podcasts

Linear Machines 101
More or Less: Behind Stats
CodeBreaker
FT Tech Tonic
Data Skeptic
Linear Digressions
Partially Derivatives
The O'Reilly Data Show



Websites

Data Science Central
Open Data Science
Machine Learning Mastery
R Views
No Free Hunch
Simply Statistics
FlowingData
Medium



Online Videos

Siraj Raval / Udacity
Brandon Foltz / Stats 101
Derek Kane
Google Developers
Marginal Revolution Uni
ODSC DataCamp
DataSchool
Re.Work Events
Coursera / EdX

References

1. Dealing with missing data: Key assumptions and methods for applied analysis. Marina Soley
2. Data Scientist: The Sexiest Job of the 21st Century. Davenport T, Patil D. Harvard Business Review, 2012.
3. Deconstructing Data Science: Breaking The Complex Craft Into It's Simplest Part. Goldstein A. The Mission. 2017.
4. About Feature Scaling and Normalization. Sebastian Raschka; 2014.
5. Alberto Cairo, The Functional Art
6. Nathan Yau <https://flowingdata.com/>