# Tecnologías del lenguaje para Explainable-AI y su impacto en el soporte a la decisión Algunas aplicaciones a salud

## K. Gibert

*Knowledge Engineering&Machine Learning group at*

*Intelligent Data Science and Artificial Intelligence Res.C. (KEMLG-at-IDEAI).*

*Universitat Politècnica de Catalunya-BarcelonaTech, Spain*

*Vice dean of Big Data, Data Science and Artificial Intelligence*

*Official Chamber of Informatics Engineering of Catalonia*

*Interest Group IABiomed-Spanish Association of Artificial Intelligence (CAEPIA)*

*InfoDay sobre tecnologías del Lenguaje en sanidad y Biomedicina*
*BSC, Barcelona 2, diciembre 2019*

# Outline

- Introduction: Automatic interpretation of profiles

- Knowledge acquisition tools

  - Prior knowledge bases

  - Ontologies

  - Termometer

  - Super-concept based distance

- Explainability through embedded strategies in Data Science methods

  - Clustering based on rules and ontologies

- Profiles oriented Explainability tools

  - Visual: TLP, a-TLP

  - Conceptual: CCEC, CI-IMS

  - Dinamic: Trajectory map, Adherence map

- Knowldedge production tools

- Other cases: Topic modelling, Explainability in ANN

- Conclusions

# Gap Data Mining- Decision making

**Data**

**Lack of data-miners**
[Hal Varian 2008]
[AS Pentland 2013]
[Sooraij Shah 2013]

**Knowledge Expertise**

Data Mining

**Interpretation**

Data-driven results

Decision – making

*The Fact Gap: The Disconnect Between Data and Decisions* *[Hammond 2004]*

*No analysis* *No understandable* *No trust*

Explainability

*Needs to be general literacy about data interpretation* *[A "Sandy" Pentland]*
*keynote Campus Party Europa Sept 4th 2013 Head of MediaLab Enterpreneurship MIT*

© K. Gibert

# Data Science concept

- 2018: Gibert, Horsburg, Athanasiadis, Holmes *[ENVSOFT, 2018]*

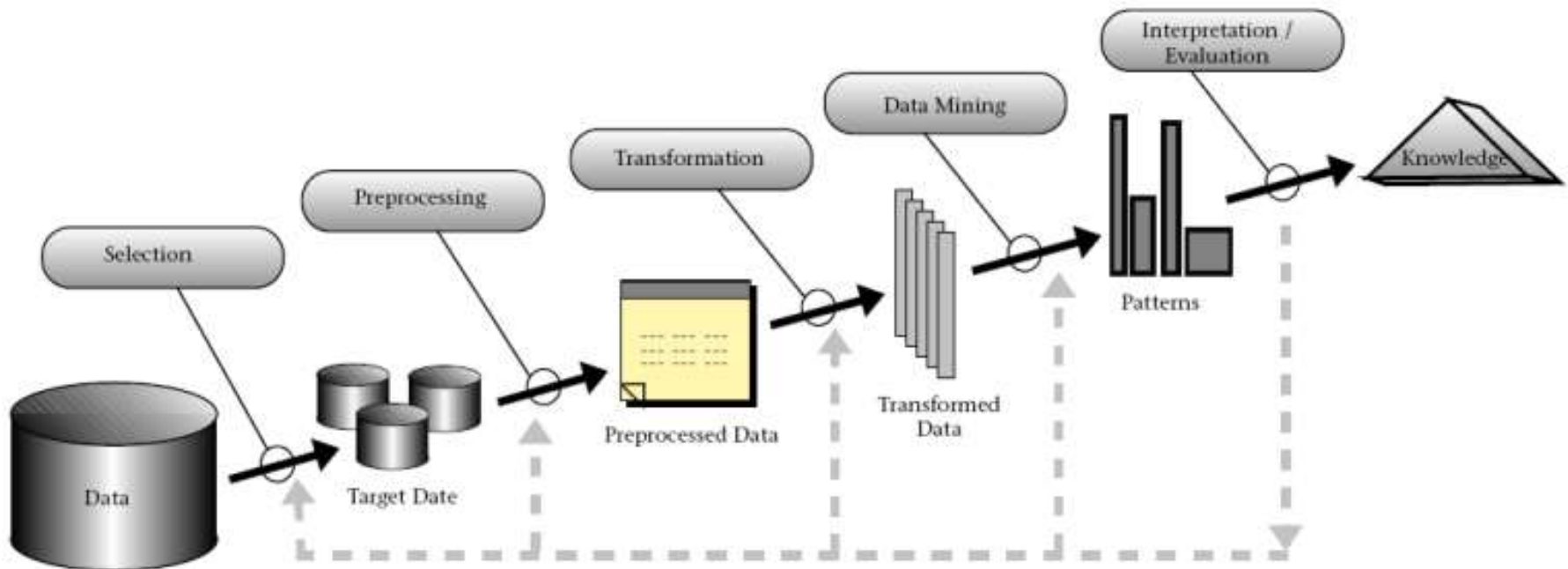*Data science : emergent multidisciplinary field combining*

- *Data analysis*
- *Data processing*
- *Domain expertise*

*To transform data into understandable and actionable knowledge*

*Relevant for informed decision making (reduces the Fact Gap)*

- *involves intensive consumption of available and required data*
- *Copes with data heterogeneity*
- *BigData is a tool, not the focus, but domain complexity*
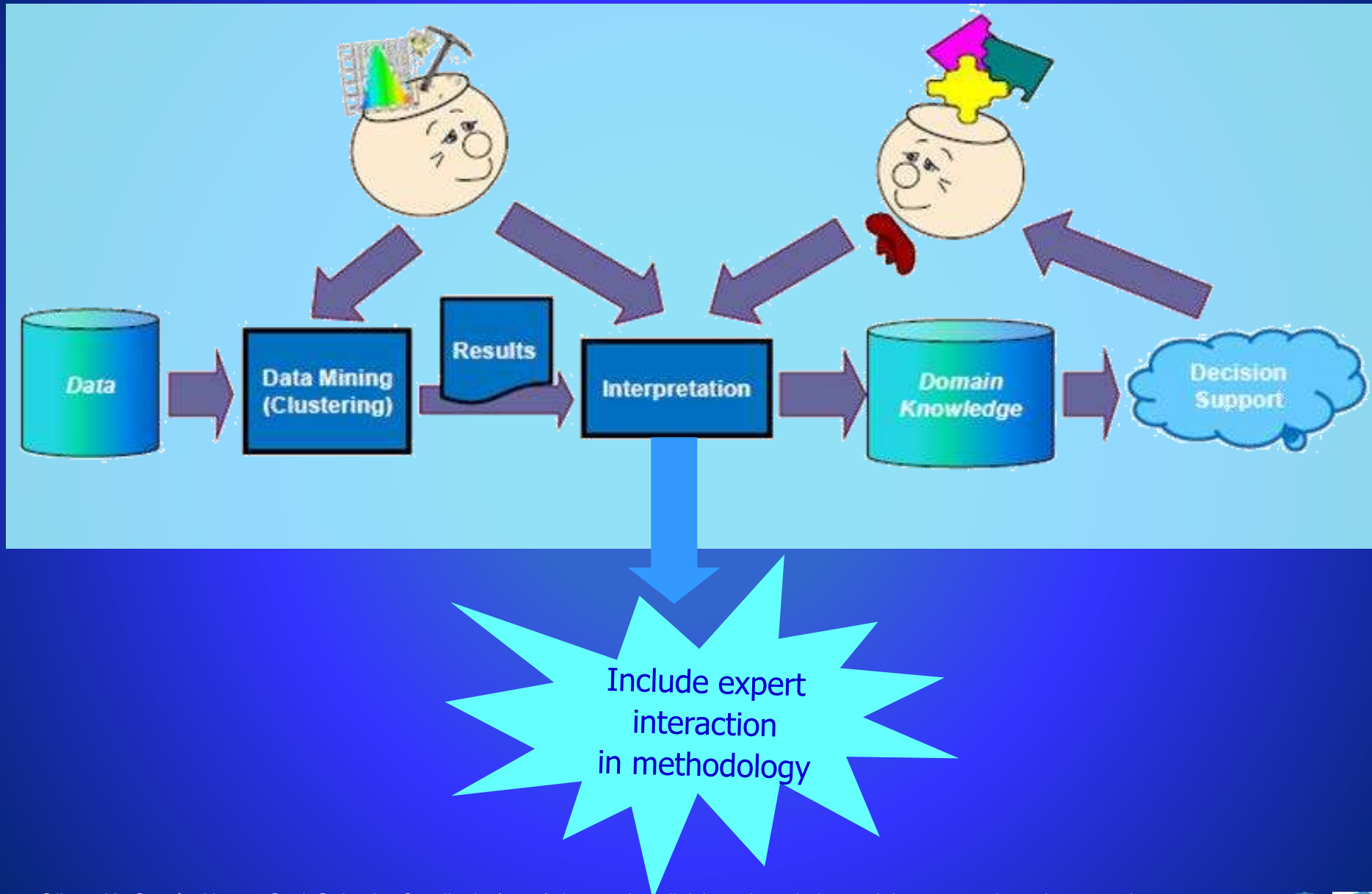
# Data Mining and Knowledge Discovery

- Knowledge Discovery System [Fayy96]:



*Focus: Clustering/Profiling*

# Expert-based collaborative Analysis (EbCA)



Include expert interaction in methodology

*Gibert, K., García-Alonso, C., & Salvador-Carulla, L. (2010). Integrating clinicians, knowledge and data: expert-based cooperative analysis in healthcare decision support. Health research policy and systems, 8(1), 28 DOI:10.1186/1478-4505-8-28*

# Profiling mental health systems in LAMIC



**Prior Expert Knowledge Acquisition**

**WHO-AIMS data**

Rows: 42 LAMIC
Columns:
 19 WHO-AIMS indicators
 + WHO-AIMS region
 + Income group
 (World Bank classification)

**KLASS**
<u>Clustering based on rules</u> *[Gib'96]*
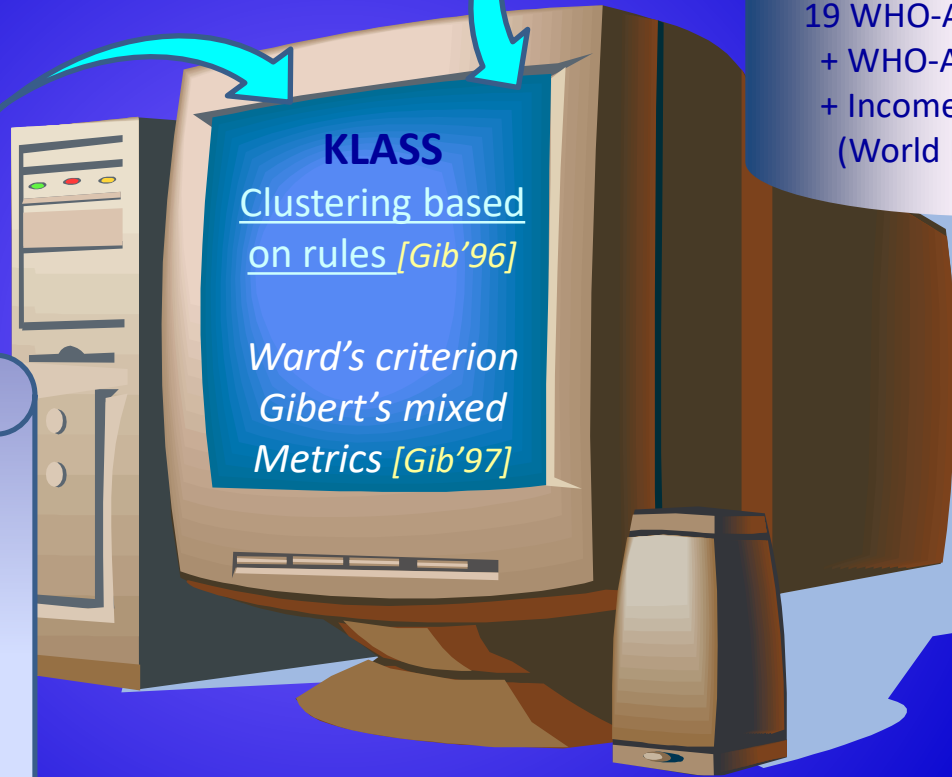
*Ward's criterion
Gibert's mixed
Metrics [Gib'97]*

**Prior Expert Knowledge**

*r0: Region=euro and
 income= lower*
 **-> PreSoviet**

*r1: Region=AMR and
 income = lower*
 **-> poorAmerica**

*r2: Region= SEAR and
 population < 10000000*
 **-> smallAsia**

Gibert, K., Izquierdo, J., Sànchez-Marrè, M., Hamilton, S. H., Rodríguez-Roda, I., & Holmes, G. (2018). Which method to use? An assessment of data mining methods in Environmental Data Science. Environmental modelling & software, 110, 3-27.

© K. Gibert

# Profiling mental health systems in LAMIC

**Prior Expert Knowledge Acquisition**

**WHO-AIMS data**

Rows: 42 LAMIC
Columns:
 19 WHO-AIMS indicators
 + WHO-AIMS region
 + Income group
 (World Bank classification)

**KLASS**
Clustering based on rules *[Gib'96]*

*Ward's criterion*
*Gibert's mixed*
*Metrics [Gib'97]*

**Prior Expert Knowledge**

*r0: Region=euro and*
*income= lower*
 **-> PreSoviet**
*r1: Region=AMR and*
*income = lower*
 **-> poorAmerica**
*R2: Region= SEAR and*
*population < 10000000*
 *-> smallAsia*

**Profiles**

**M H S in L A M I C**



IDEAI

# Profiling mental health systems in LAMIC countries for healthcare policy-making at WHO
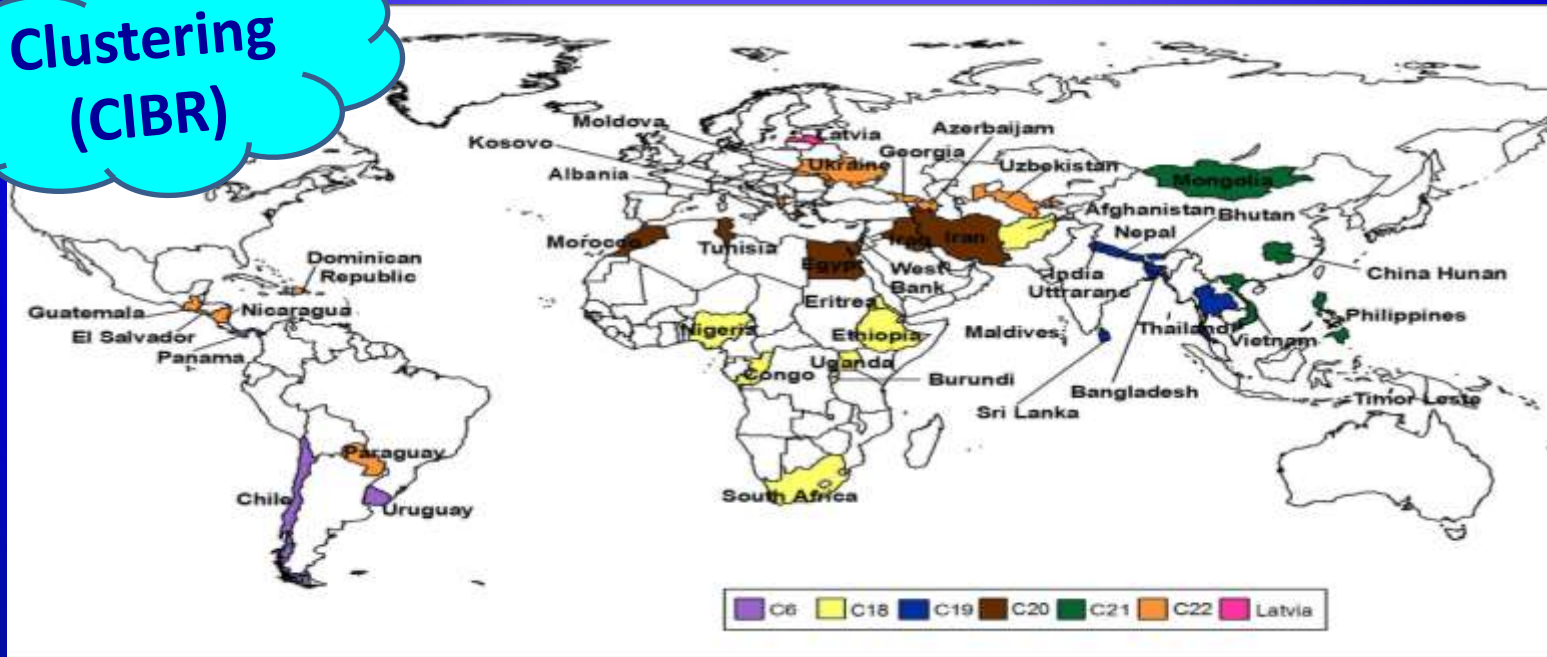
## Use WHO-AIMS DB to learn a **typology** of MHS in LAMIC

- Easy **understanding** of reality

- Assessment to countries

- Intervention design: guidelines, mental health policies....
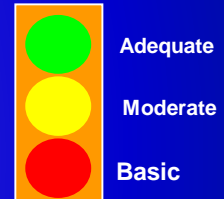
**Postprocessing CPG, TLP**

**Clustering (CIBR)**

Goal



| C6 | C18 | C19 | C20 | C21 | C22 | Latvia |

Gibert, K, L. Salvador-Carulla, J. Morris, A. Lora, S. Saxena (2017) The data mining approach as the starting point for Mental Health policy-making in low and middle income countries at World Health Organization. In Proc of ISI 2017. Belcasem Abdou et al, Marrakesh, Morocco, Jul 2017

© K. Gibert

# TLP elicits clustering criteria
# Conceptualization

## Induces categories of variables and classes

| BLOCK | | CLASS | CARE CAPACITY | | | | | | CARE ARRANGEMENT | | | | | POLICY FRAMEW | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Incom | HR | $MHe | Treat pre | Cap-ratio | close beds | %$m-hosp | LTC-pacs | comc arew | Lund | Manua | Legis | Pol-plan | Gov-Rep |
| I | *Upper-Moderate* | CI | UpMid | Highst | Highst | Highst | Highst | Lowest | High++ | High | High | Lowest | No | yes | yes | yes |
| | | C6 | UpMid | Mod | Mod | Mod | Low | Low | Mod | Highest | Low | High | Some | yes | yes | Some |
| | *Low-Mod* | C22 | LMid | Mod | Low | Mod | Mod | Mod | Highest | High | Mod | Mod | No | Some | Some | No |
| II | *Mid-Limited* | C21 | LMid | Low | Low | Low | Mod | Mod | High | High | Mod | Mod | yes | No | yes | Some |
| | | C20 | LMid | Low | Low | Low | Low | Mod | High | High | Highest | Highest | yes | Some | yes | Some |
| | | C19 | LMid | Low | Low | Lowest | Lowest | High | High | Low | Low | High | Some | No | Most | Most |
| | *Very Lim* | C18 | Low | Low | Low | Low | Low | Highst | High++ | Low | Low | Mod | No | Few | Most | No |

- 🟢 Adequate
- 🟡 Moderate
- 🔴 Basic

## Supports data-driven Ontologies

*Gibert K, D. Conti, D. Vrecko (2012) Assisting the end-user in the interpretation of profiles for decision support. An application to wastewater treatment plants. Environmental Engineering and Management Journal 11(5): 931-944*

# Knowledge Production

## MHS for LAMIC ontology



Intervention plans designed for each type
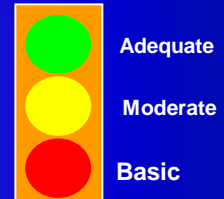
© K. Gibert

# The KLASS thermometer-tool
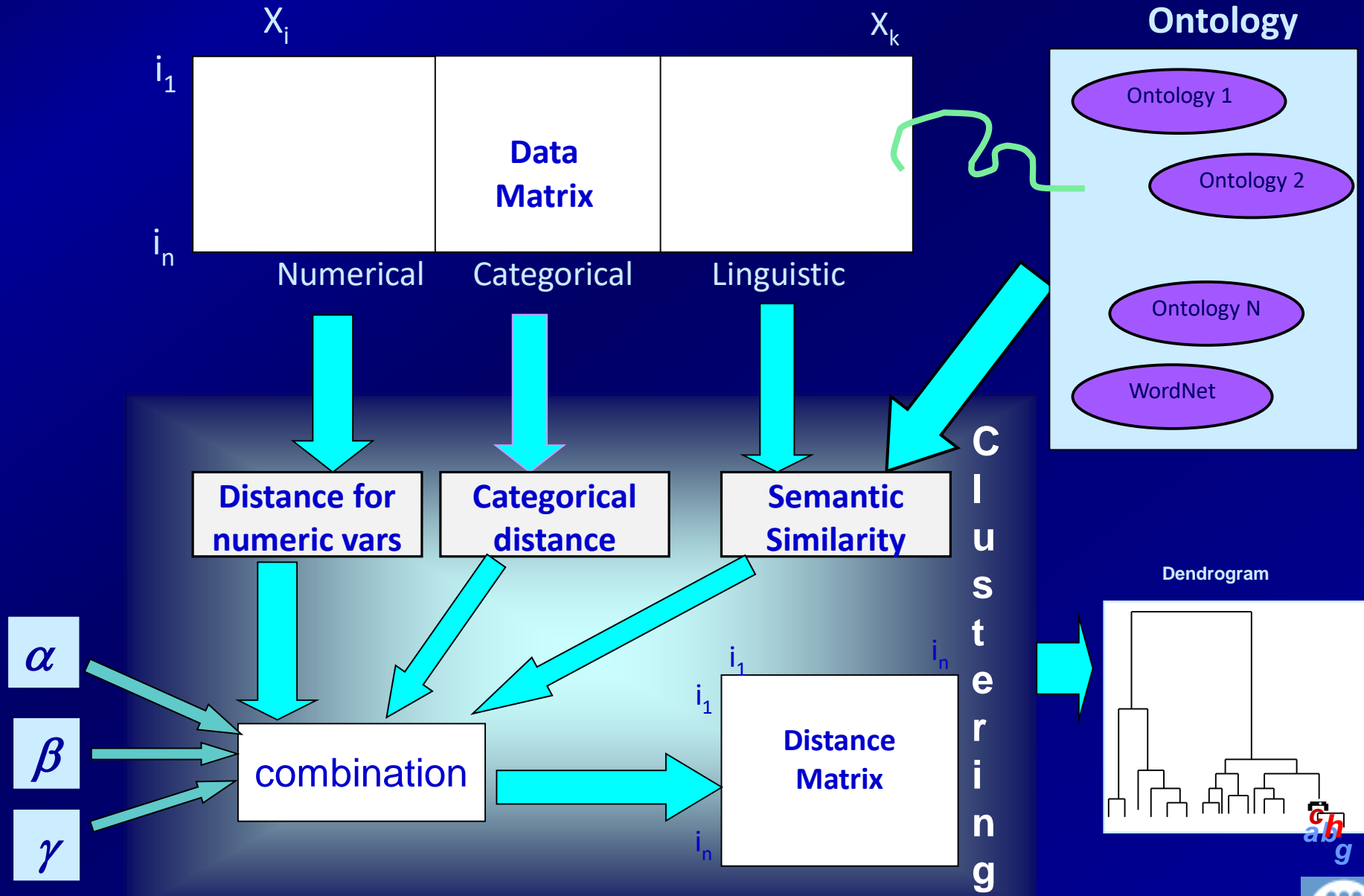
# TLP elicits clustering criteria
## Conceptualization

## Induces categories of variables and classes

| BLOCK | | CLASS | CARE CAPACITY | | | | | | CARE ARRANGEMENT | | | | | POLICY FRAMEW | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Incom | HR | $MHe | Treat pre | Cap-ratio | close beds | %$m-hosp | LTC-pacs | comc arew | Lund | Manua | Legis | Pol-plan | Gov-Rep |
| I | *Upper-Moderate* | CI | UpMid | Highst | Highst | Highst | Highst | Lowest | High++ | High | High | Lowest | No | yes | yes | yes |
| | | C6 | UpMid | Mod | Mod | Mod | Low | Low | Mod | Highest | Low | High | Some | yes | yes | Some |
| | *Low-Mod* | C22 | LMid | Mod | Low | Mod | Mod | Mod | Highest | High | Mod | Mod | No | Some | Some | No |
| II | *Mid-Limited* | C21 | LMid | Low | Low | Low | Mod | Mod | High | High | Mod | Mod | yes | No | yes | Some |
| | | C20 | LMid | Low | Low | Low | Low | Mod | High | High | Highest | Highest | yes | Some | yes | Some |
| | | C19 | LMid | Low | Low | Lowest | Lowest | High | High | Low | Low | High | Some | No | Most | Most |
| | *Very Lim* | C18 | Low | Low | Low | Low | Low | Highst | High++ | Low | Low | Mod | No | Few | Most | No |

- 🟢 Adequate
- 🟡 Moderate
- 🔴 Basic

## Supports data-driven Ontologies

*Gibert K, D. Conti, D. Vrecko (2012) Assisting the end-user in the interpretation of profiles for decision support. An application to wastewater treatment plants. Environmental Engineering and Management Journal 11(5): 931-944*

# Semantic Distances



**Ontology**

- Ontology 1
- Ontology 2
- Ontology N
- WordNet

$X_i$    $X_k$

$i_1$   **Data Matrix**   $i_n$

Numerical   Categorical   Linguistic

**Distance for numeric vars**   **Categorical distance**   **Semantic Similarity**

Clustering

$\alpha$   $\beta$   $\gamma$

combination   **Distance Matrix**

$i_1$   $i_n$

**Dendrogram**

© **K. Gibert**

UPC

# a-TLP: going further (WWTP case)

*©K. Gibert*

# CCEC: Conceptual Caracterization by Embedded Conditioning

## Exploits dendrogramm structure to induce classification rules



*r1.BC0.−r50r0: ((treatpre∈[18,57,172,77))∧ (comcarewor∈[0,0197,0,1098)))∧(Region∈ {AFR})−→(NovaClasseBLN7)C18*

*r2.BC1.−r2−r46−r50r0−r35−r37−r39 : (((treatpre∈[18,57,172,77))∧(comcarewor∈ [0,0197,0,1098))) ∧(((Region∈{SEAR})∨ (lundpararectrail∈[0,49,0,53)))∨(comcarewor∈ [0,0197,0,0255))))∧ ((((Region∈{SEAR})∧ (treatpre∈ [31,81, 87,59]))∧ (lundpararectrail = 0,49))∧(comcarewor = 0,0197))−→ (NovaClasseBLN7)C19*

*r4.BC3.−r46r53:(treatpre∈[18,57,172,77))∧(comca rewor∈(0,1313,0,624])−→(NovaClasseBLN7)C20*

# CIMS: Cluster Interpretation based on Integrated Marginal Significance

*Same differences with same conceptualizations in all classes*
*Consistency Inter Classes: Generalized Test –Value*

Numerical: $\tau_\nu = \dfrac{\bar{X}^C - \bar{X}}{\sqrt{\left(1 - \dfrac{n_c}{n}\right)\dfrac{s^2}{\nu}}} \sim t_{\nu-1}$

Qualitative: $\pi_\nu = \dfrac{p_{sc} - p_s}{\sqrt{\left(1 - \dfrac{n_c}{n}\right)\dfrac{p_s(1 - p_s)}{\nu}}} \sim z$

## Sensitivity Analysis

$\downarrow \nu \rightarrow \uparrow$ *p-value*

| | $\epsilon_2$ 0.2 | $\epsilon_1$ 0.3 | $n$ | $\epsilon_1$ 0.3 | $\epsilon_2$ 0.2 |
|---|---|---|---|---|---|
| | **0.5n** | **0.7n** | **n** | **1.3n** | **1.5n** |
| **Descripion-Power ($\Pi$)** | $\nu_1$ | $\nu_2$ | $\nu_3$ | $\nu_4$ | $\nu_5$ |
| **Robust Non-Significant ($\overline{R}$)** | ✘ | ✘ | ✘ | ✘ | ✘ |
| **Moderate Non-Significant ($\overline{M}$)** | ✘ | ✘ | ✘ | ✘ | ✔ |
| **Weak Non-Significant ($\overline{W}$)** | ✘ | ✘ | ✘ | ✔ | ✔ |
| **Weak Significant (W)** | ✘ | ✘ | ✔ | ✔ | ✔ |
| **Moderate Significant (M)** | ✘ | ✔ | ✔ | ✔ | ✔ |
| **Robust Significant (R)** | ✔ | ✔ | ✔ | ✔ | ✔ |
| **Basic Descriptor (B)** | B | B | B | B | B |

## Class Descriptor

**< W, C, description-power, sense>**

$W = \begin{cases} X & \textit{if X numerical} \\ <X, s> & \textit{if X categorical} \wedge \\ & \quad s \textit{ category} \in D_X \end{cases}$

$sense \in \{\uparrow, \downarrow\}$

## Regular Expressions

*Proportion of Smokers (Tobacco) is higher in C1*

*Weight is high in class C2*

*Age    is lowest in class C1*

*K. Gibert*
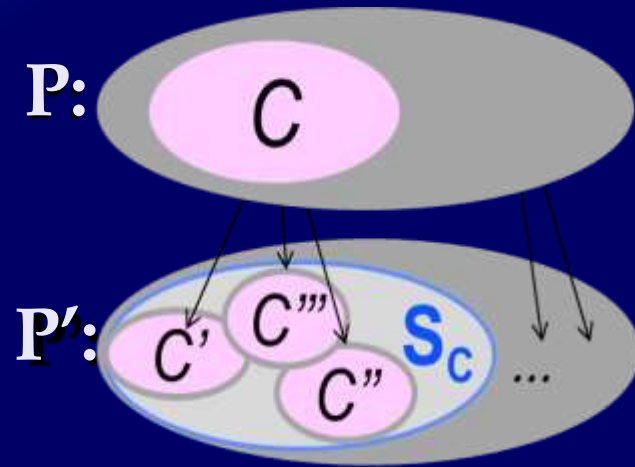
# Interpreting X in Nested Partitions

$$C = \bigcup_{C' \in S_C} C' \in P'$$

**P:** C

**P':** C' C''' $S_C$ C'' ...

Relationship between interpretation of X in C and $S_C$

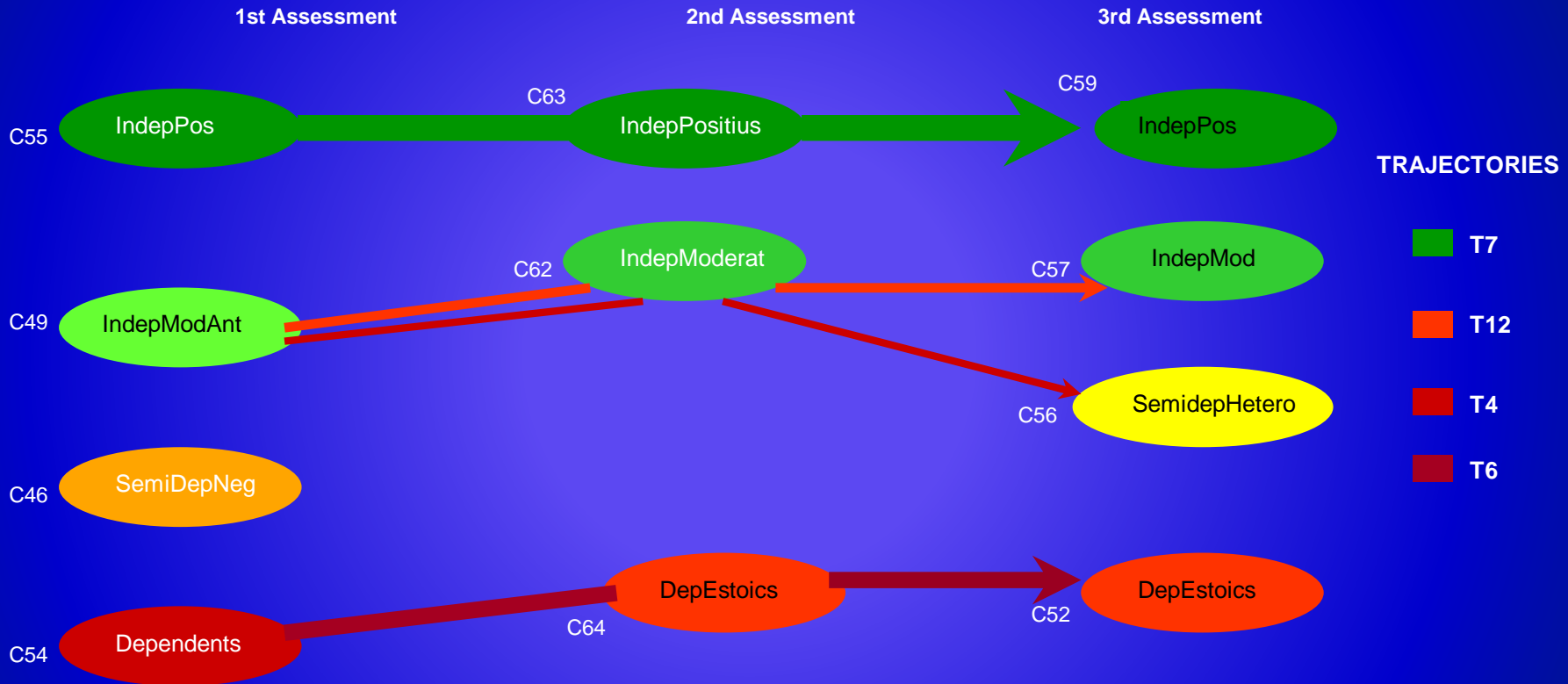| Super Class / Sub Classes | Non-Significant | Significant |
|---|---|---|
| **Non-Significant** | Irrelevance<br>$\forall C' \in S_C: R(C, C')$<br>$= Irrelevance$ | Inconsistency<br>$\forall C' \in S_C: R(C, C')$<br>$= Inconsistency$ |
| **Significant** | Specification<br>$\exists C' \in S_C: R(C, C')$<br>$= Specification$ | Inheritance<br>$\exists C' \in S_C: R(C, C')$<br>$= Inheritance$ |

Contradiction

# NCI-IMS: *Cluster Interpretation based on Integrated Marginal Significance for Nested partitions*

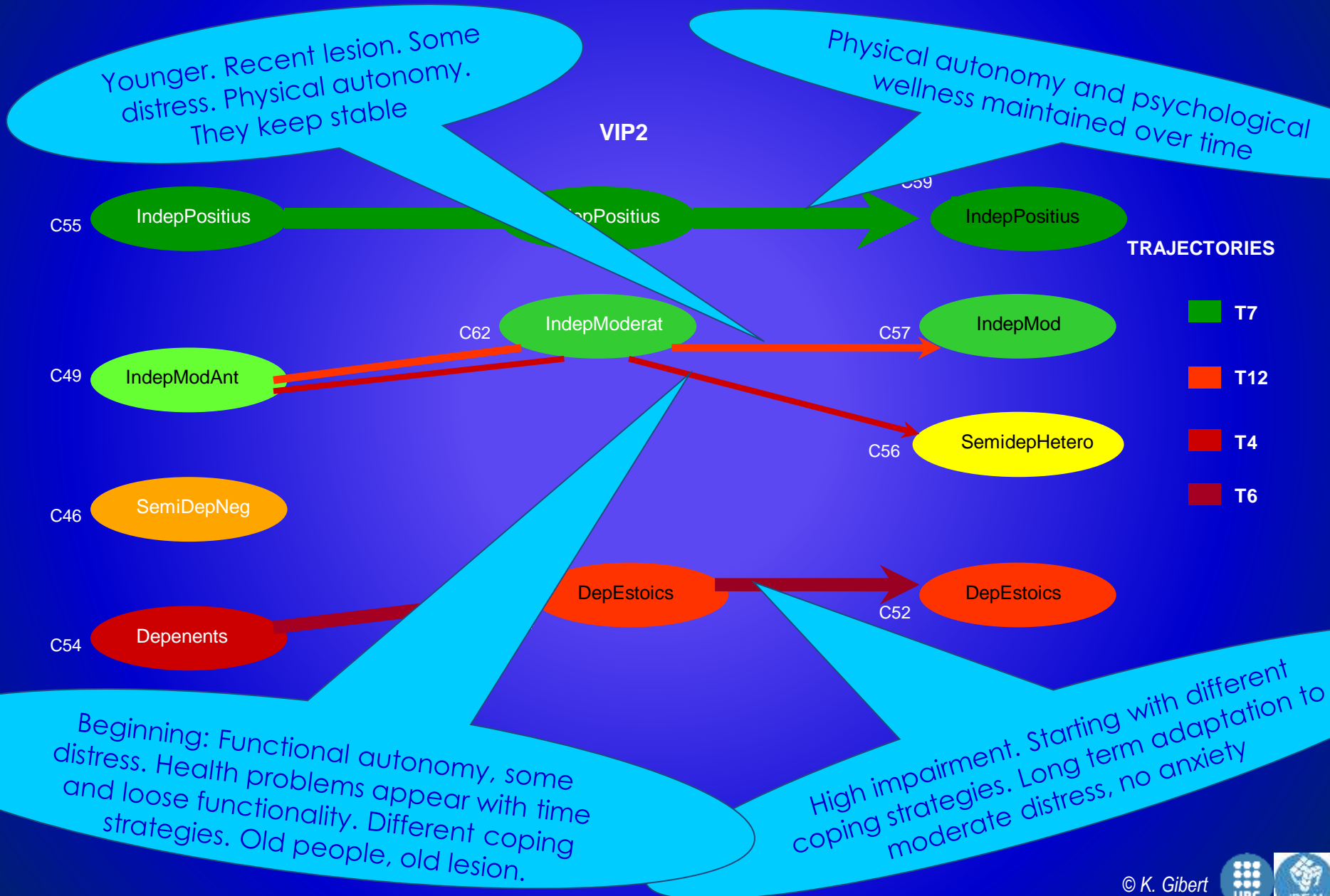Table $\mathcal{A}$ : Actions associated to Table $\mathcal{R}$

| Super Class / SubClass | $\overline{R}$ | $\overline{M}$ | $\overline{W}$ | W | M | R/B |
|---|---|---|---|---|---|---|
| $\overline{R}$ | | | | | W in description of C | |
| $\overline{M}$ | W ignored in description of C and C' $\in S_C$ | | | | | |
| $\overline{W}$ | | | | | W in description of C and C' $\in S_C$ | |
| W | | | | | | |
| M | W in description of C' $\in S_C$ | | | | | |
| R/B | | | | | | |

Gibert, K, B. Sevilla-Villanueva, M. Sànchez-Marrè (2016) *The role of significance tests in consistent interpretation of nested partitions.*
*Journal of Computational and Applied Mathematics, 292: 623-633, Elsevier, Amsterdam, NL (htpps://doi.org/10.1016/j.cam.2015.01.031)*

*K. Gibert*

# Trajectory maps
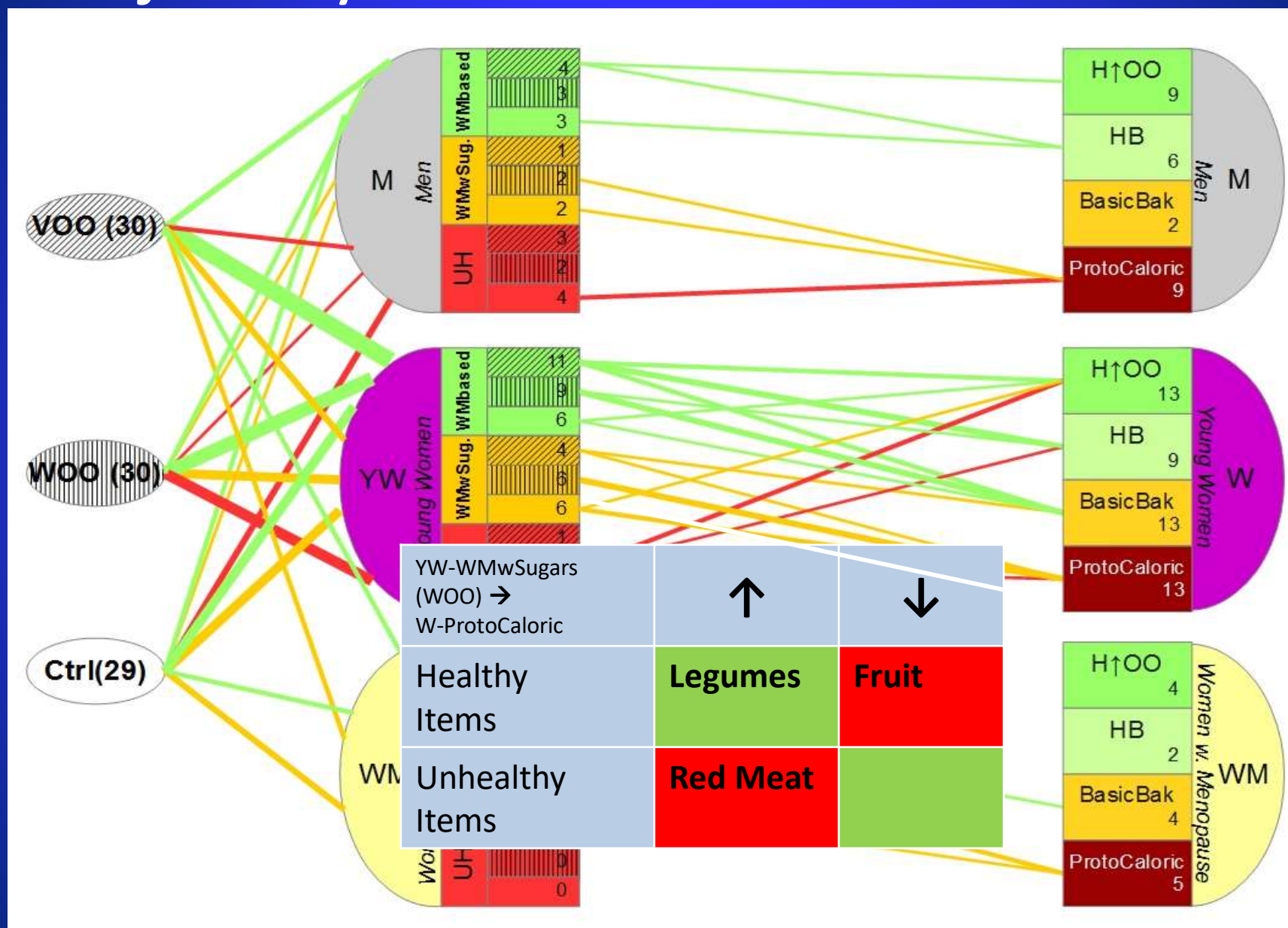
## More typical patterns ($\gamma \geq 0.05$)

# Expert's conceptualization of patterns

# Trajectory Characterization. Adherence



| YW-WMwSugars (WOO) → W-ProtoCaloric | ↑ | ↓ |
|---|---|---|
| Healthy Items | Legumes | Fruit |
| Unhealthy Items | Red Meat | |

© K. Gibert

# Assignment of the profile of a new patient

## Given a new patient:

Estimate $\pi_{High}$ by applying equation 1

If $\quad p_{High}$ is $>= \xi$ *then* assign patient to *High* profile

*Else,* Estimate $\pi_{IntII}$ by applying equation 2

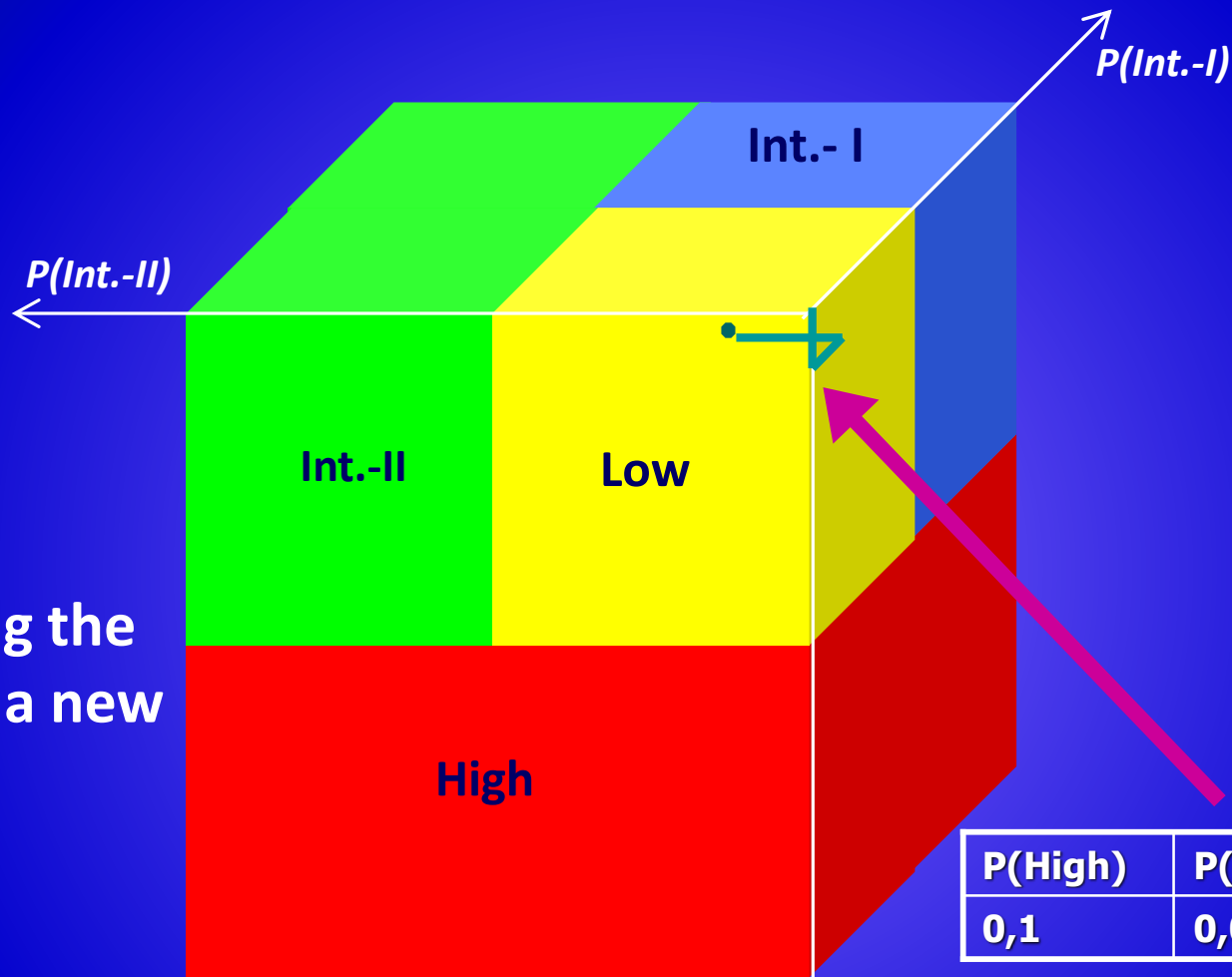*If* $p_{IntII}$ is $>= \xi$ *then* assign patient to *IntermediateII* profile.

*Else* Estimate $\pi_{IntI}$ by applying equation 3.

*If* $\quad p_{IntI}$ is $>= \xi$

*then* assign patient to *IntermediateI* profile.
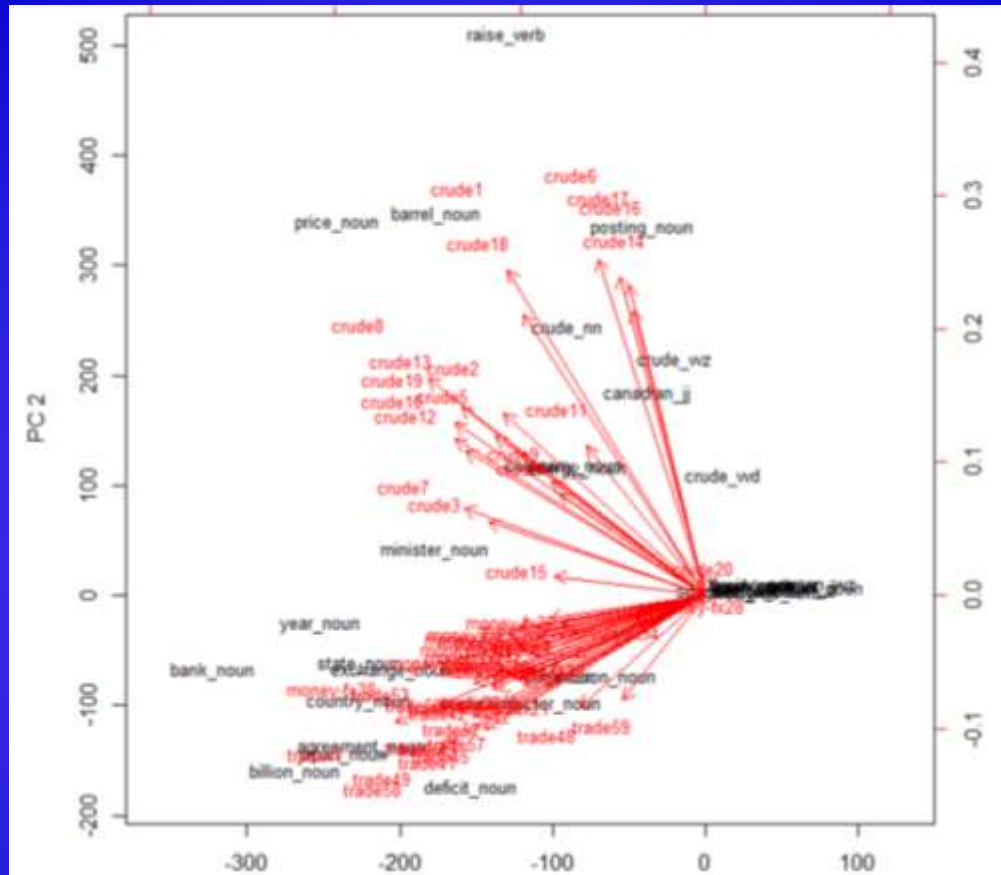
*Else* assign patient to Low profile.

Gibert K, Rodríguez-Silva G, Annicchiarico R (2013) Post-processing: bridging the gap between modelling and effective decision-support. The Profile Assessment Grid in Human Behaviour. Mathematical and Computer Modelling 57(7-8):1633-1639, Elsevier
https://doi.org/10.1016/j.envsoft.2009.11.004

# The Profile Assessment Grid



P(Int.-I)

P(Int.-II)

Int.- I

Int.-II    Low

High

**Identifying the profile of a new patient**

| | |
|---|---|
| B2 | 3 |
| B4 | 2 |
| B9 | 2 |
| S4 | 1 |
| S5 | 2 |
| S9 | 2 |

| P(High) | P(Int-I) | P(Int-II) |
|---|---|---|
| 0,1 | 0,07 | 0,20 |

**Misclassification Rate ($\varepsilon=0.5$): 8.3%**

$$P(High) = \frac{e^{-35.93+1.70*B2+3.35*B4+3.98*B9+2.20*S4}}{(1+ e^{-35.93+1.70*B2+3.35*B4+3.98*B9+2.20*S4})}$$

# PCA for topic modelling



Find terms with significant contributions to axes

Generalize in the reference  ontology (Wordnet by default)
Discover the latent variables (automatic interpretation of axes)

# Interpreting ANN



Visualization of Input Effect: VEC curve (Bank Marketing)

© K. Gibert

# Conclusions

- Explainable models required for trust and decisions
- Post-processing provide explainability
  - Visual tools: TLP/a-TLP (profiling), PAG (predictive models)
  - Conceptual: CCEC, CI-MIS (machine readable)
  - Dynamics: trajectory maps, adherence maps
- Prior knowledge transfer to models increase explainability
  - Termometers (semantics of variables, polarities)
  - Prior Knowledge Bases
  - Ontologies (semantic relations between terms)
- Language technologies play a relevant role in building these tools

*© K. Gibert*

# Tecnologías del lenguaje para Explainable-AI y su impacto en el soporte a la decisión
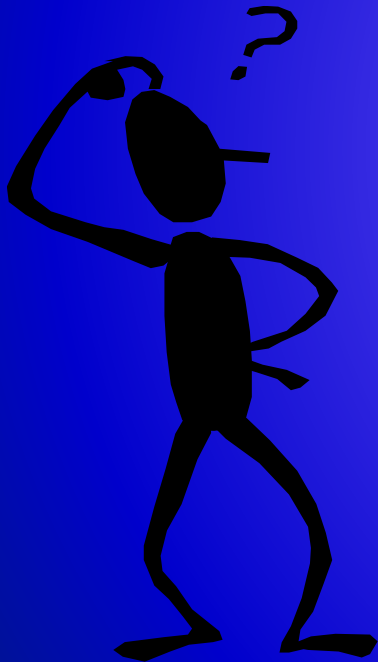# Algunas aplicaciones a salud

## K. Gibert

*karina.gibert@upc.edu*
*https://www.eio.upc.edu/en/homepages/karina*

*KEMLG-@-IDEAI: Knowledge Engineering and Machine Learning group at Intelligent Data Science and Artificial Intelligence Research Center*

*Universitat Politècnica de Catalunya, Barcelona*

### *Are there any questions?...*