

Understanding the Interaction between Diseases using Big Healthcare Data

Peter Lucas Martijn Lappenschaar Arjen Hommersom

Radboud University Nijmegen
Institute for Computing and Information Sciences

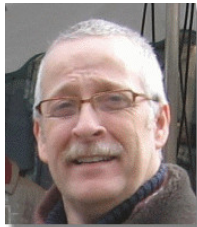
12th December, 2014

- Research focus: we try to **combine** knowledge and innovation from computing science with that of medicine
- Clinical knowledge is not **shallow** and therefore requires a decent **knowledge representation** method (difficult issue, much progress during last three decades, but still a long way to go)
- **Uncertainty** is an essential ingredient of any form of clinical decision making
- \Rightarrow **Machine learning** should take knowledge representation and uncertainty into account (compare medical statistics)

Pandora's box

Rob is 67 years old and has been ill for some time, in particular he is currently treated for:

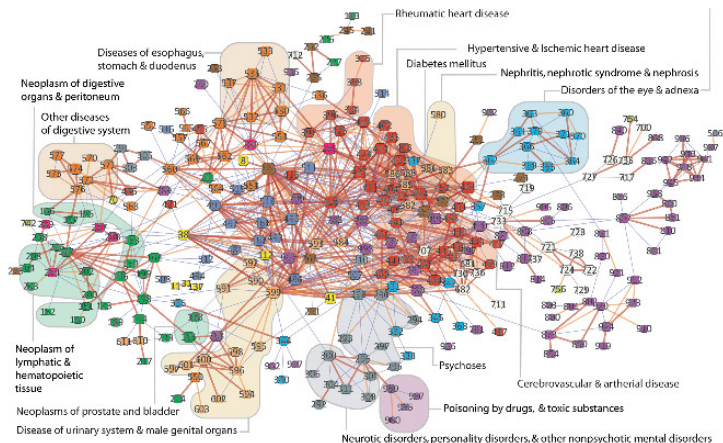
- diabetes mellitus type 2
- status after myocardial infarction
- chronic obstructive pulmonary disease



Rob is **not unique** ...

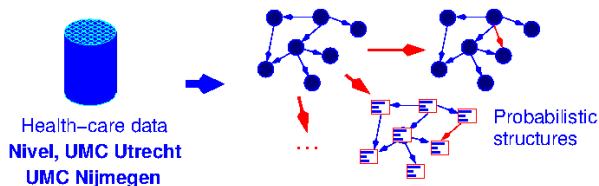
- **2/3rd** of patients older than 65 years have **2 or more disorders** at the same time
- However, medicine is organised around **single disorders!** = **problem of multimorbidity**

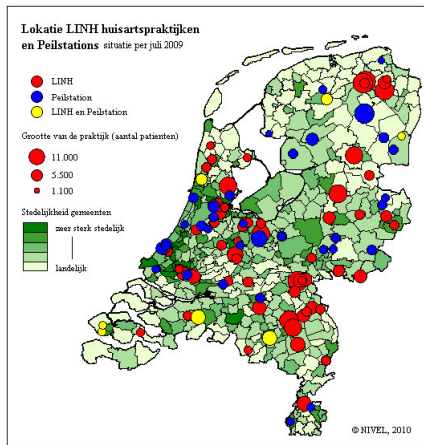
Challenge: dealing with diseases and their interactions



- **Complexity:** many individual diseases and classes of disease
- **Probabilistic relationships:** uncertain interactions between diseases, regional, social, and gender differences in prevalence

- **Knowledge representation and reasoning:**
 - how to exploit probabilistic graphical models to capture clinical knowledge
 - model-based diagnosis, prediction and decision-theoretic planning
- **Decision support and clinical guidelines:** how to integrate task execution with probabilistic reasoning
- **Learning probabilistic models** about disease interactions from large health-care databases:





- Multiple sources
 - practices
- Source characteristics
 - urbanicity
 - size
 - type
- Statistical methods
 - multi-level

- Paired comparison of frequency of occurrence of signs and symptoms given two disorders as **likelihood ratio** or **odds ratio**:

$$\frac{P(f | d_1)}{P(f | d_2)} \quad \text{or} \quad \frac{Odds(d_1 | f)}{Odds(d_2 | f)}$$

with f a feature, e.g. symptom, lab result, and d_1, d_2 two disorders

- Example: Odds Ratios derived from a clinical research:

	Diabetes Mellitus
Stroke	1.46
Heart Failure	1.76
Diabetes Mellitus	1.0
Hypertension	2.65

- Measures to compare two disorders are determined by:
 - **Linear regression** for continuous outcome O on explanatory variables (predictors) e :

$$P(O | e) \sim \mathcal{N}(\mu, \Sigma) \text{ with } \mu = E[O | e] = \beta^T e$$

- **Logistic regression** for dichotomous outcomes:

$$P(O | e) \sim \text{Bernoulli}(p) \text{ with } \text{logit}(E[O | e]) = \beta^T e$$

- Example logistic regression with an **interaction term**:
 - $\text{logit}(E[F | d_1, d_2]) = \beta_0 + \beta_1 d_1 + \beta_2 d_2 + \beta_{12} d_1 d_2$

Regression: interaction term β_{ij}

Disorders D_1 and D_2 and patient findings F :

$$\begin{aligned} \exp(\beta_{12}) &= \frac{\text{Odds}(f \mid d_1, d_2) \text{Odds}(f \mid \bar{d}_1, \bar{d}_2)}{\text{Odds}(f \mid \bar{d}_1, d_2) \text{Odds}(f \mid d_1, \bar{d}_2)} \\ &= \frac{\lambda(d_1, d_2 \mid f) \lambda(\bar{d}_1, \bar{d}_2 \mid f)}{\lambda(\bar{d}_1, d_2 \mid f) \lambda(d_1, \bar{d}_2 \mid f)} \\ &= \frac{\frac{P(d_1, d_2 \mid f)}{P(d_1, d_2 \mid \bar{f})} \frac{P(\bar{d}_1, \bar{d}_2 \mid f)}{P(\bar{d}_1, \bar{d}_2 \mid \bar{f})}}{\frac{P(\bar{d}_1, d_2 \mid f)}{P(\bar{d}_1, d_2 \mid \bar{f})} \frac{P(d_1, \bar{d}_2 \mid f)}{P(d_1, \bar{d}_2 \mid \bar{f})}} \\ &= \left\{ \frac{P(d_1, d_2 \mid f) P(\bar{d}_1, \bar{d}_2 \mid f)}{P(d_1, \bar{d}_2 \mid f) P(\bar{d}_1, d_2 \mid f)} \right\} \left\{ \frac{P(\bar{d}_1, d_2 \mid \bar{f}) P(d_1, \bar{d}_2 \mid \bar{f})}{P(d_1, d_2 \mid \bar{f}) P(\bar{d}_1, \bar{d}_2 \mid \bar{f})} \right\} \end{aligned}$$

Regression: other measures

Regression gives us outcomes like:

- $Odds(d_1 | d_2, F)$
- $Odds(d_2 | d_1, F)$

But with some calculation we can also obtain measures like:

- $Odds(d_1 | F)$
- $Odds(d_2 | F)$
- $Odds(d_1, d_2 | F)$

For example, using the odds derived in clinical research on one of the previous slides, we obtain:

$$Odds(hypertension, diabetes | heartfailure) = 1.88$$

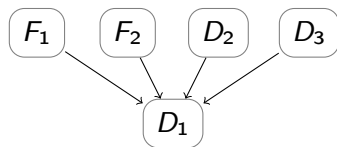
Capturing interaction by regression

Given a set of outcomes and observations, to obtain joint probabilities using regression, in order to investigate **interactions within multimorbidity**, we need:

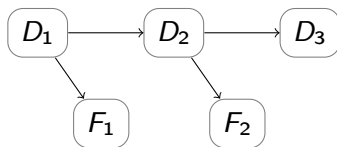
- a regression model **for each** outcome variable of interest
- within each regression model **all possible interaction terms**

Graphical representation

wrong, diagnostic model



right, causal model



The **diagnostic model** represents **regression analysis** of D_1 . It assumes all remaining variables are independent and certain, whereas in the **causal model** all true (possible) dependencies are modeled

Probabilistic graphical models, such as Bayesian networks, support explicit modelling by a graph (uncovered by **structure learning**)

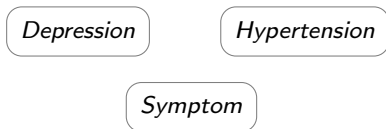
Concurrent multimorbidity

Independent diseases co-occur at the same time (unconditionally independent)

- $P(D_i, D_j) = P(D_i)P(D_j)$

No common signs and symptoms:

- $\forall F$: Conditional independence
 - $P(D_i, D_j | F) = P(D_i | F)P(D_j | F)$
 - Logistic regression: $\beta_{ij} = 0$
 - Structure learning: no edges (paths)



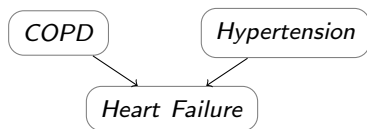
Concurrent multimorbidity

Independent diseases co-occur at the same time (unconditionally independent)

- $P(D_i, D_j) = P(D_i)P(D_j)$

Common signs and symptoms:

- $\exists F$: Conditional dependence
 - $P(D_i, D_j | F) \neq P(D_i | F)P(D_j | F)$
 - Logistic regression: $\beta_{ij} \neq 0$
 - Structure learning: $D_i \rightarrow F \leftarrow D_j$



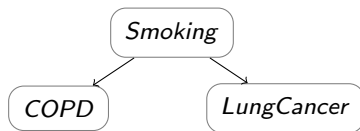
Dependent diseases:

- $P(D_i, D_j) \neq P(D_i)P(D_j)$

because of

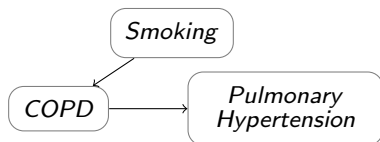
- $\exists F$: **Common cause** (conditional independence)

- $P(D_i, D_j | F) = P(D_i | F)P(D_j | F)$
- Logistic regression: $\beta_{ij} = 0$
- Structure learning: $D_i \leftarrow F \rightarrow D_j$

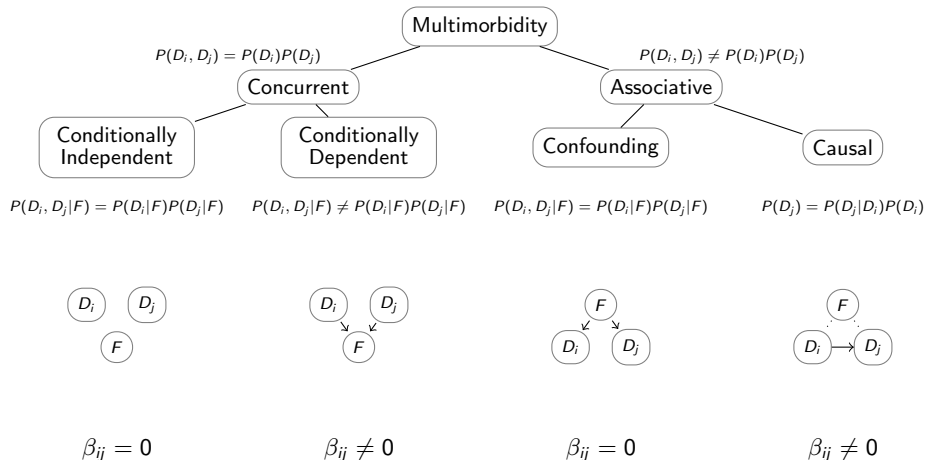


Dependent diseases:

- $P(D_i, D_j) \neq P(D_i)P(D_j)$
- D_j dependent of D_i
 - $P(D_i, D_j | F) = P(D_j | D_i)P(D_i | F)$
 - Logistic regression: $\beta_{ij} \neq 0$
 - Structure learning: $F \rightarrow D_i \rightarrow D_j$



Multimorbidity – types of correlation



■ Logistic Regression

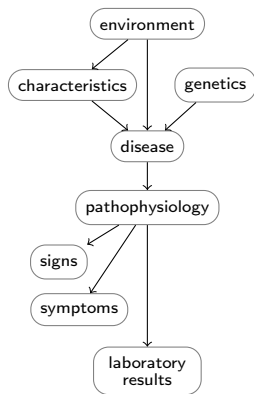
- Diseases often used as ...
 - outcome variable in one model (A)
 - explanatory variable in another model (B)
 - \Rightarrow multiple models
- Use of interaction terms:
 - $\beta_{ij} = 0 \rightarrow$ True Independence or Confounding?
 - $\beta_{ij} \neq 0 \rightarrow$ Conditional Dependence or Causality?

■ Bayesian Networks

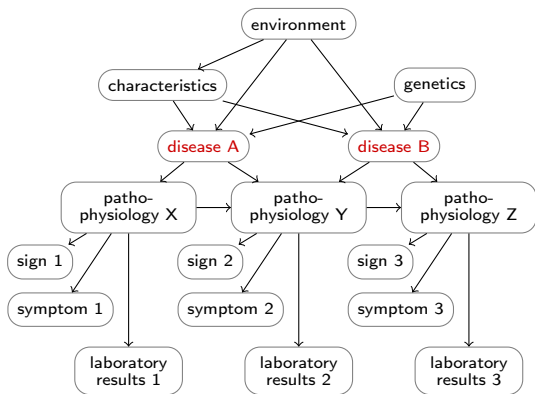
- All variables treated as uncertain
 - one model!
 - (possible) representation of underlying processes
- Interactions automatically incorporated
- Allows distinguishing between various forms of multimorbidity

Disease modelling by Bayesian networks

single disease



multiple diseases



Abstract model of a single disease (left) and multiple diseases (right)

Big data: multilevel regression

- To model variation of outcomes between various groups (e.g. different general practices), taking into account correlation within groups
- Formulation in terms of regression models (with l being a vector of higher level variables):

- **multilevel linear regression:**

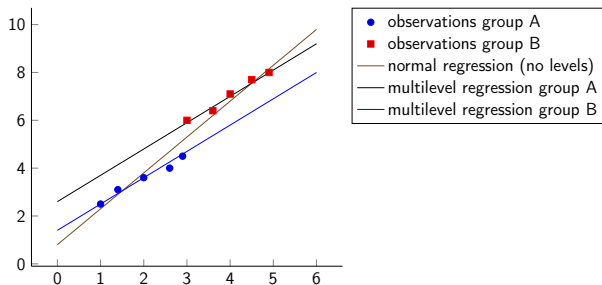
$$P(O_k | e, l) \sim \mathcal{N}(\mu, \Sigma) \text{ with } \mu = E[O | e, l] = \beta_k e = (\delta_k + \gamma_k l)^T e$$

- **multilevel logistic regression:**

$$P(O_k | e, l) \sim \text{Bernoulli}(p) \text{ with } \text{logit}(E[O_k | e, l]) = (\delta_k + \gamma_k l)^T e$$

with k the **group index**, γ_k the **level parameters**, and δ_k the **group variation**

Multilevel regression

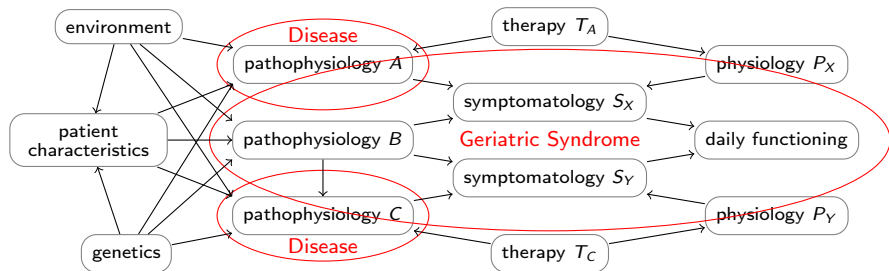


Limitations:

- Only comparison between one outcome variable and predictors
- Only predictions are treated as uncertain
- No explicit knowledge about relationships between predictors
- Within multimorbidity some variables are both outcome and predictor

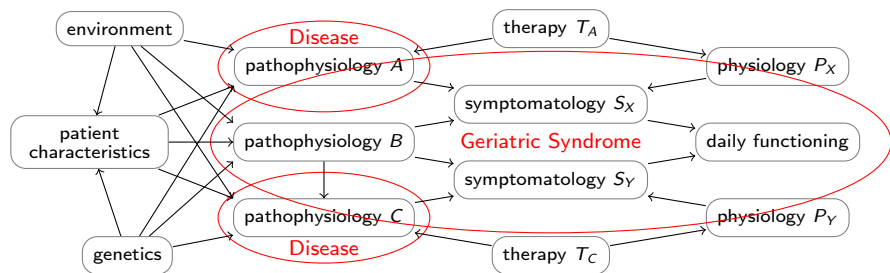
Disease modelling of multimorbidity

Graphical representation of risks, pathophysiology, and symptomatology:



Disease modelling of multimorbidity

Graphical representation of risks, pathophysiology, and symptomatology:



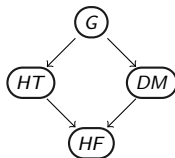
- $\text{logit}(E[\text{DiseaseA} \mid \text{Age}, \text{Gender}, \text{SymptomX}]) = \beta_{0A} + \beta_{1A}\text{Age} + \beta_{2A}\text{Gender} + \beta_{3A}\text{SymptomX}$
- $\text{logit}(E[\text{DiseaseB} \mid \text{Age}, \text{Gender}, \text{SymptomX}, \text{SymptomY}]) = \beta_{0B} + \beta_{1B}\text{Age} + \beta_{2B}\text{Gender} + \beta_{3B}\text{SymptomX} + \beta_{4B}\text{SymptomY}$
- $\text{logit}(E[\text{DiseaseC} \mid \text{Age}, \text{Gender}, \text{SymptomY}, \text{DiseaseB}]) = \beta_{0C} + \beta_{1C}\text{Age} + \beta_{2C}\text{Gender} + \beta_{3C}\text{SymptomX} + \beta_{4C}\text{DiseaseB}$

As a Bayesian network

A **Bayesian network** is a tuple $\mathcal{B} = (G, X, P)$, with $G = (V, E)$ a directed acyclic graph, $X = \{X_v \mid v \in V\}$ a set of random variables indexed by V , and P a joint probability distribution such that:

$$P(X_1 = x_1 \wedge \cdots \wedge X_n = x_n) = \prod_{v \in V} P(X_v = x_v \mid X_j = x_j \text{ for all } j \in \pi(v))$$

Simple example:



← explanatory variables

← outcome variables

$$P(V) = P(X_{HF} \mid X_{HT}, X_{DM})P(X_{HT} \mid X_G)P(X_{DM} \mid X_G)P(X_G)$$

Structure and **parameters** of a Bayesian network can be learned from data.

Combination of concepts

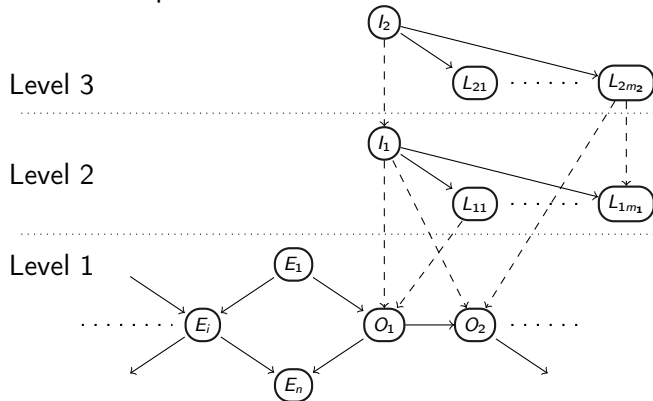
In summary, with patient data acquired from general practices and the aim of modelling multiple disease, we are facing:

- 1 hierarchical data structures
→ which can be analysed using multilevel regression
- 2 multiple diseases with multiple possible interactions
→ which can be modelled using probabilistic graphical methods
 - Bayesian networks
 - undirected graphs
 - hybrid graphs

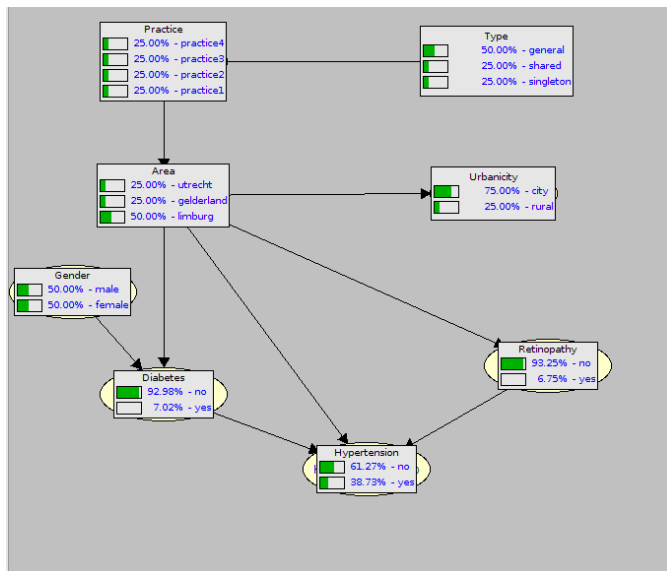
Our goal → adopting both concepts into **multilevel Bayesian networks**

MLBN with independence and intra-level structure

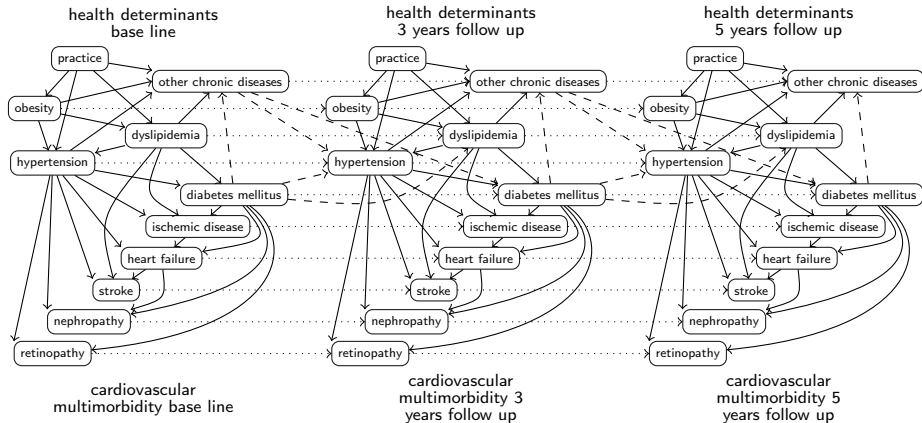
- Here all variables are uncertain (random) and expressed as such
- Representation of different levels of outcomes (and other variables)
- Inter-level dependence $--\rightarrow$
- Intra-level dependence \rightarrow



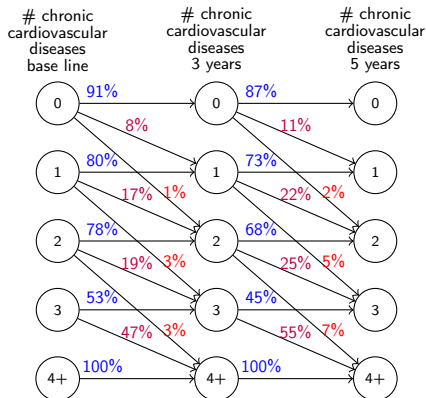
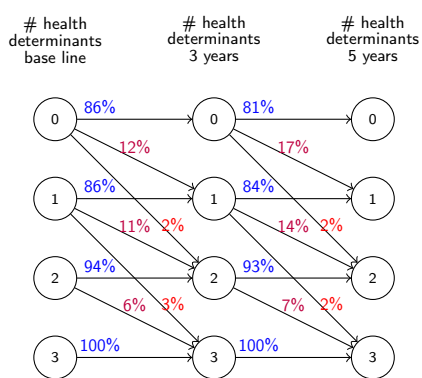
Toy example



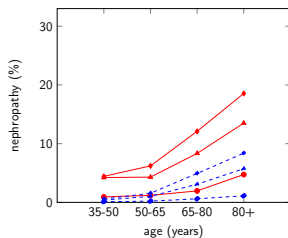
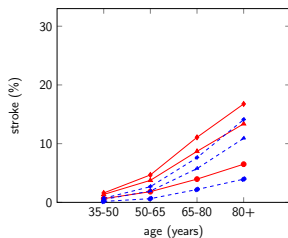
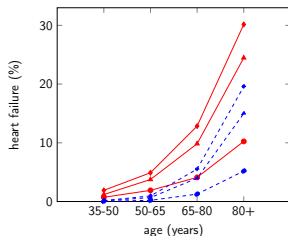
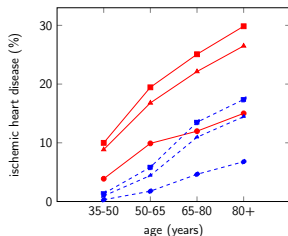
Cardiovascular model - MLBN at 3 time points



Transition probabilities



In context - diabetes mellitus



diabetics ● — base line ▲ — 3yr follow-up ■ — 5yr follow-up
non-diabetics ● — base line ▲ — 3yr follow-up ■ — 5yr follow-up

■ Machine learning in medicine

- Requires a combination of knowledge representation, reasoning and learning methods
- Big healthcare data: need for new methods

■ Methodology

- Integration of multilevel analysis and Bayesian networks
- Visualization of interactions between disease variables
- Personalization of patients (e.g., diabetics)
- Fundament towards clinical guidelines that deal with multimorbidity