

Dept. of Statistics and Biostatistics

Personalized disease networks (PDN) for understanding and predicting cardiovascular diseases and other complex processes.

Javier Cabrera

Fei Wang

Nabil Adam, Jaideep Vaidya, Dhammika Amaratunga, William Kostis, John Kostis.

Research supported by a Chancellor seed grant from Rutgers Newark Rutgers, The State University of New Jersey

RUTGERS



Outline

- Big Data in medicine.
- MIDAS: cardiovascular diseases database.
- What are personalized disease networks?
- Analysis Plan.
- Enriched linear and non-linear PCA
- Proof of concept. Do PDN's help improve the prediction of cardiovascular outcomes.



Big Data in medicine

Electronic Medical Records



Big Data

- Databases of hospitalizations: ICD-9 or ICD-10 diagnosis and procedural codes with dates.
- Labs Data: Test results
- Health Insurance Providers: Costs of treatment, hospitalization, procedures and drugs.
- Medical Charts: Doctors handwriting, text on paper records.
 Patient data from devices: collected in real time by devices and uploaded into the cloud has many new applications.
- Objectives: Patient Summary Charts for doctor's uses.
 Medical Research Analysis of Big Data

The MIDAS dataset

MIDAS: All hospitalizations for cardiovascular diseases in New Jersey since 1980.

15 million hospitalizations

5 million patients with cardiovascular diseases.

MIDAS is housed at Rutgers Cardiovascular Institute

Plan: From MIDAS database PDN for cardiovascular diseases.

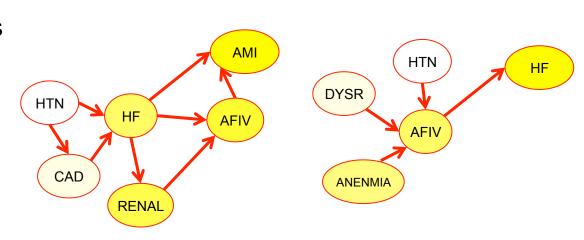
MIDAS database:

- Demographical information of the patients
- ICD9 codes for the diseases and the procedural codes and the corresponding date of the disease/procedure.
- Negotiating data sharing agreement with BlueCross/S

RUTGERS What are personalized disease networks?

- Personalized (patient) disease networks(PDN),
- Represent the evolution of a disease on a patient across subsequent diagnoses and medical interventions.
- Evolutionary steps of various cardiovascular diseases and procedures
- Network: represents a data based construction in a form of a connected graph to derive the disease evolution of each individual patient
- To build a PDN we use a list of events and their date
- Each event connected to another by a directed edge if the dates fulfill a certain criteria.

Example: PDN's of 2 patients that diagnosed with Atrium Fibrillation and Heart Failure. The arrows represent the order of events.





Example of data from MIDAS

| MR | CMTHY | CHD | HTN | DM | NEO | COPD | RENAL | STROKE | АМІ |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|---------|
| 8/23/06 | 10/4/11 | 8/29/95 | 8/30/95 | 2/27/98 | NA | 3/20/03 | 9/9/10 | NA | NA |
| 5/22/95 | NA | 5/22/95 | 1/13/95 | NA | NA | 5/14/95 | 11/30/10 | 4/27/07 | NA |
| NA | NA | 9/25/07 | 9/18/99 | 2/12/07 | NA | 9/18/99 | 4/21/09 | NA | NA |
| 12/27/11 | NA | 10/14/97 | 10/19/01 | 8/3/07 | 3/8/03 | 1/22/02 | 10/14/97 | 10/19/01 | 4/30/03 |
| NA | 8/14/08 | 12/19/00 | 10/18/01 | 12/19/00 | 11/18/10 | 12/15/01 | 11/5/08 | NA | NA |
| NA | NA | 10/27/97 | 10/21/04 | 10/31/97 | NA | NA | NA | NA | NA |
| NA | 4/3/95 | 4/3/95 | 4/3/95 | 10/16/95 | NA | 3/31/02 | 9/3/09 | NA | NA |
| NA | NA | 11/22/96 | 11/22/96 | NA | 8/29/05 | 3/9/11 | NA | NA | NA |
| NA | NA | 5/2/95 | 2/25/98 | NA | 1/15/08 | 2/25/98 | 9/8/10 | NA | NA |
| NA | 11/18/07 | 6/2/06 | 3/26/02 | 3/26/02 | NA | NA | 11/18/07 | NA | NA |

Column labels correspond to ICD9 codes and dates .

Rule: Arrow connecting a Diagnosis/Procedure A to another Diagnosis/Procedure B where A and B are observed at times T_A and T_B

1. Naive Rules

If $T_A < T_B$ then $A \rightarrow B$, If $T_A = T_B$ it will depend on A and B. If $0 < T_B - T_A < K$ then $A \rightarrow B$, If $T_A = T_B$ it will depend on A and B.

2. Flexible Rule:

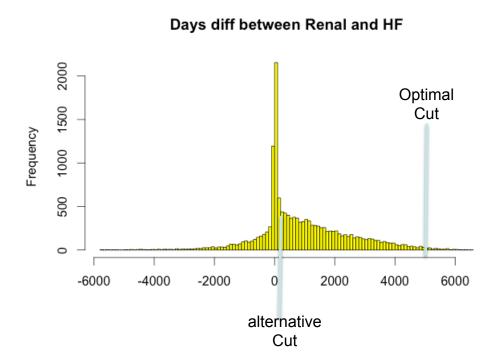
If $0 < T_B - T_A < K(A,B)$ then $A \rightarrow B$, If $T_A = T_B$ it will depend on A and B

To Calculate K(A,B) we may use several methods
(i) Expert Physicians: (It is been currently performed)
(ii) Find K(A,B) based on the prediction of some response.

3. Expert Physicians input. So far we got 3 physicians to give info.

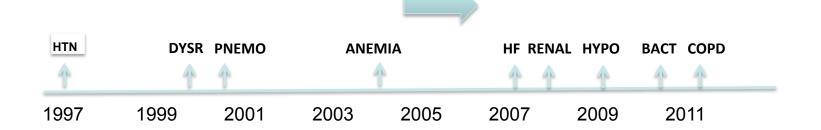
So far we got 4 physicians: Drs J Kostis, W Kostis, Dr Moreira, Dr Pantazopoulos from CVI/Cardiology at Rutgers

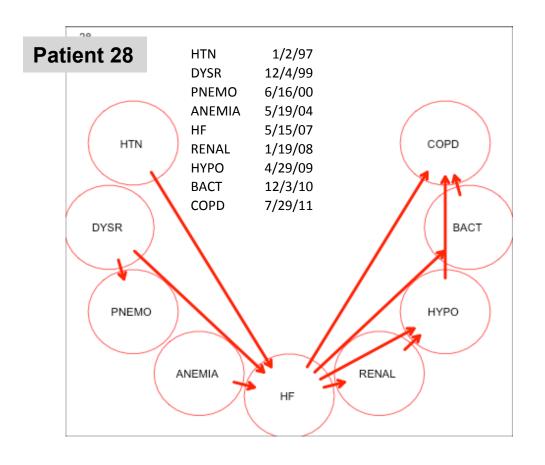
Optimal cuts presented to for input and a decision is reached Example: Renal dysfunction (kidney failure, dialysis) -> Heart Failure



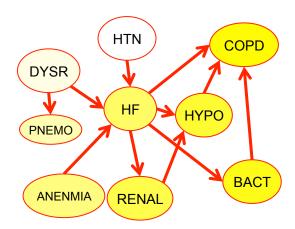


Personalized disease networks(Flexible rule . 2)

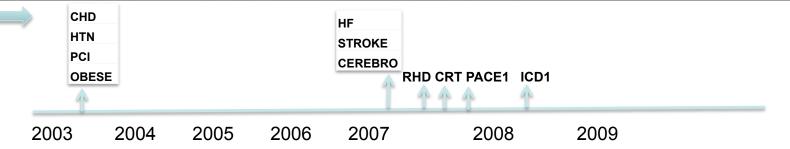




Patient 28



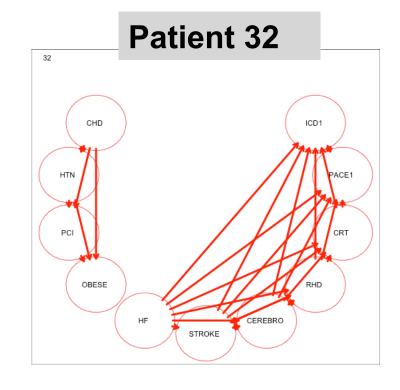
RUTGERS Naïve rule for building PDE's: K = 1000 days



Example of personalized (patient 32) disease networks(PDN)

Patient 32

| CODE | DATE |
|---------|---------|
| CHD | 4/12/03 |
| HTN | 4/12/03 |
| PCI | 4/12/03 |
| OBESE | 4/12/03 |
| HF | 6/7/07 |
| STROKE | 6/7/07 |
| CEREBRO | 6/7/07 |
| RHD | 6/28/07 |
| CRT | 9/19/07 |
| PACE1 | 9/27/07 |
| ICD1 | 5/28/08 |

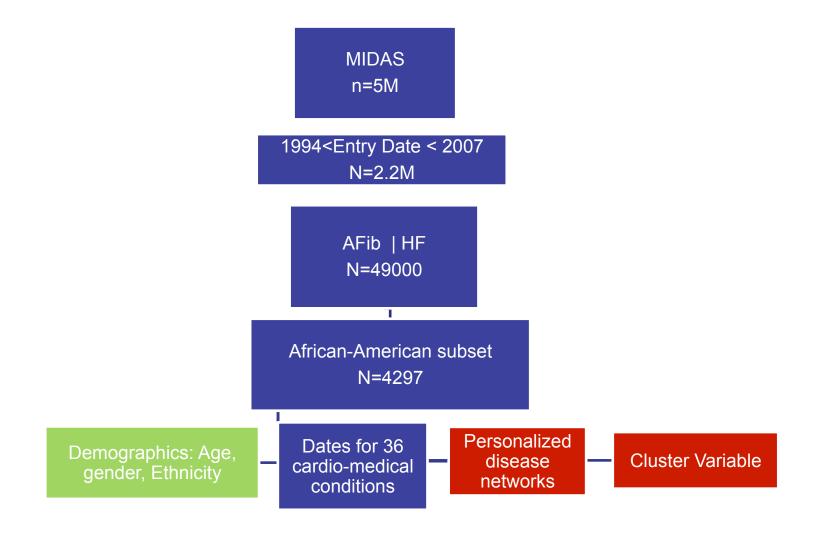




Data Matrix for personalized disease networks(PDN),

| AFIB to HF | HF to AFIB | AFIB to MR | MR to AFIB | AFIB to CMTHY | CMTHY to AFIB | AFIB to CHD | CHD to AFIB | AFIB to HTN | HTN to AFIB |
|---------------|---------------|---------------|---------------|---------------|---------------|-------------|----------------|----------------|----------------|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Data Subsets



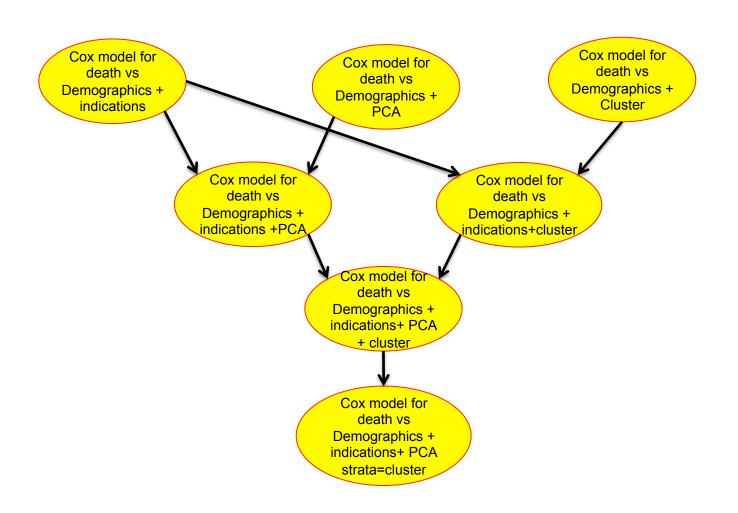
Demographics: age, gender, ethnicity

Dates for 36 cardio-medical indications

Personalized disease networks

→ Enriched PCA

Factor



Enriched PCA Scores

Enriched PCA scores. Amaratunga, Cabrera, Shkedy(2014)

Enriched PCA:
$$W_{G\times G} = Diag(W_i)$$

$$X^* = X W = U^* D^* V^{*'}$$

 $n \times G \quad n \times G \quad G \times G \quad G \times G$

Covariance or Correlation $R^* = V^* D^{*2} V^{*3}$

$$R^* = V^* D^{*2} V^{*'}$$

 $G \times G G \times G G \times G$

n<G

$$X^* = W \quad X = U^* \quad D^* \quad V^*$$
 $G_{\times n} \quad G_{\times G} \quad G_{\times n} \quad G_{\times n} \quad n_{\times n} \quad n_{\times n}$

Covariance or Correlation $R^* = U^* D^{*2} U^{*3}$

$$R^* = U^* D^{*2} U^{*'}$$
 $G \times G \qquad G \times n \qquad n \times n \qquad n \times G$

W matrix can be calculated as a measure of importance of a variable.

RUTGERS

Enriched PCA scores: Amaratunga, Cabrera, Shkedy(2014)

How to determine the weights? $W = Diag(W_i)$

The idea is to reduce the effect of the spurious signal that is create by *G* large number of variables,

Use any model $Y = f(X_i)$ and obtain the model *p-value* p_i .

If G is large convert p-values into q-values q_i .

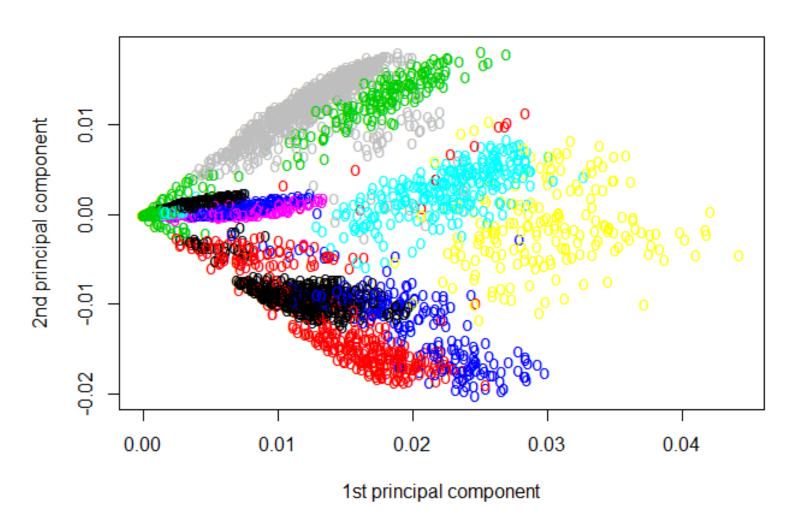
Define the weights:

$$W_i = -log(q_i)$$
$$W_i = 1/q_i$$

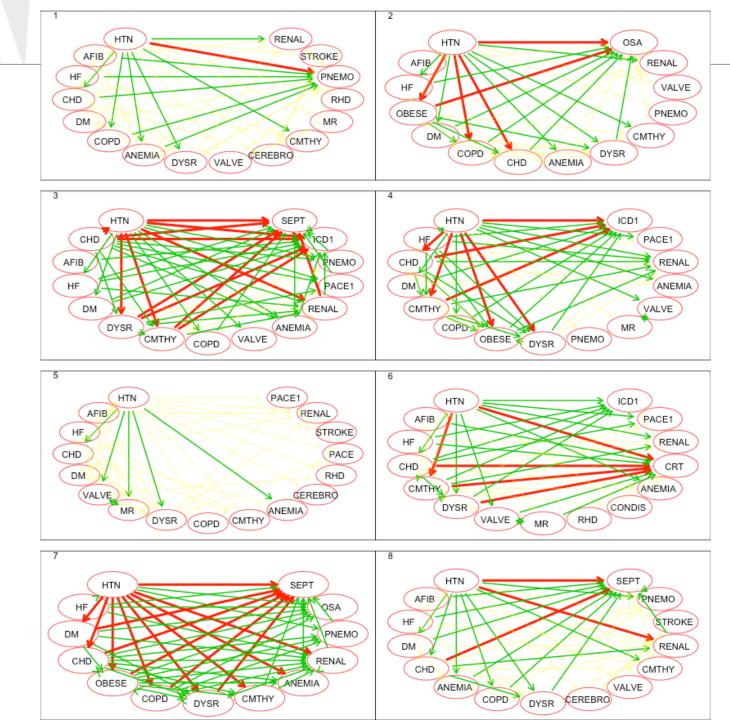
Enriched PCA and Clustering

Enriched PCA: 43 PCA

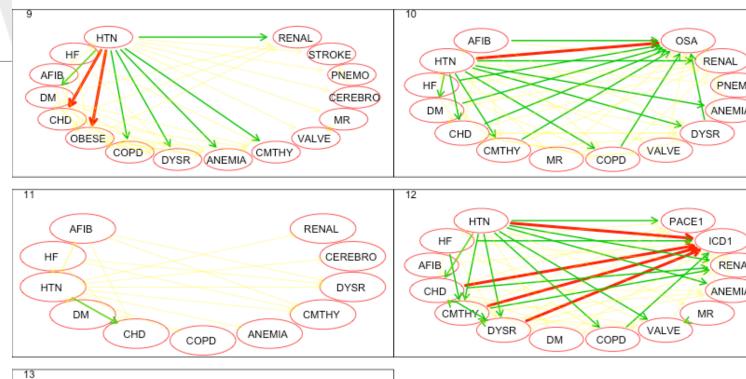
Hierarachical Clustering: 13 Clusters.

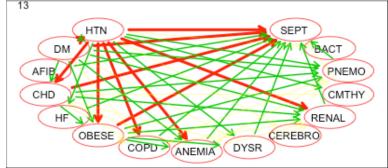


RUTGERS



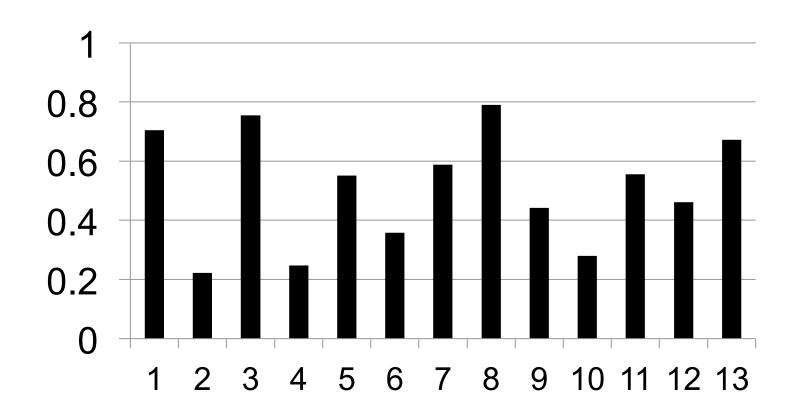
RUTGERS





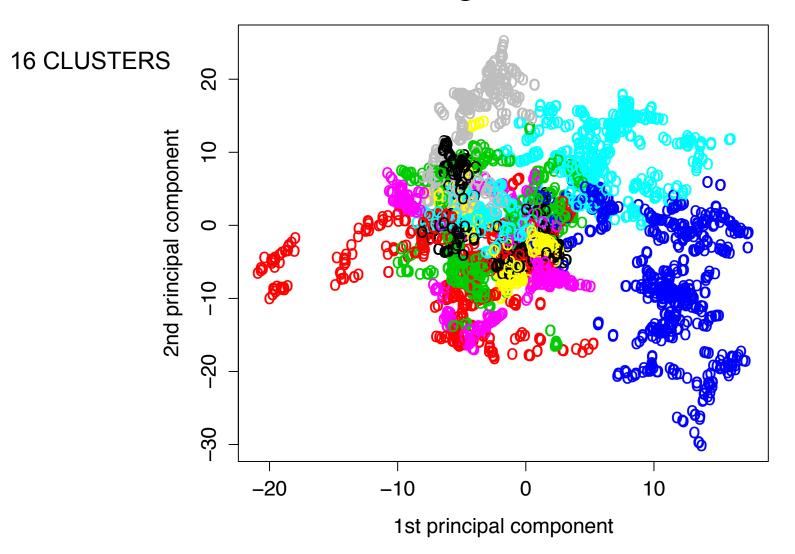


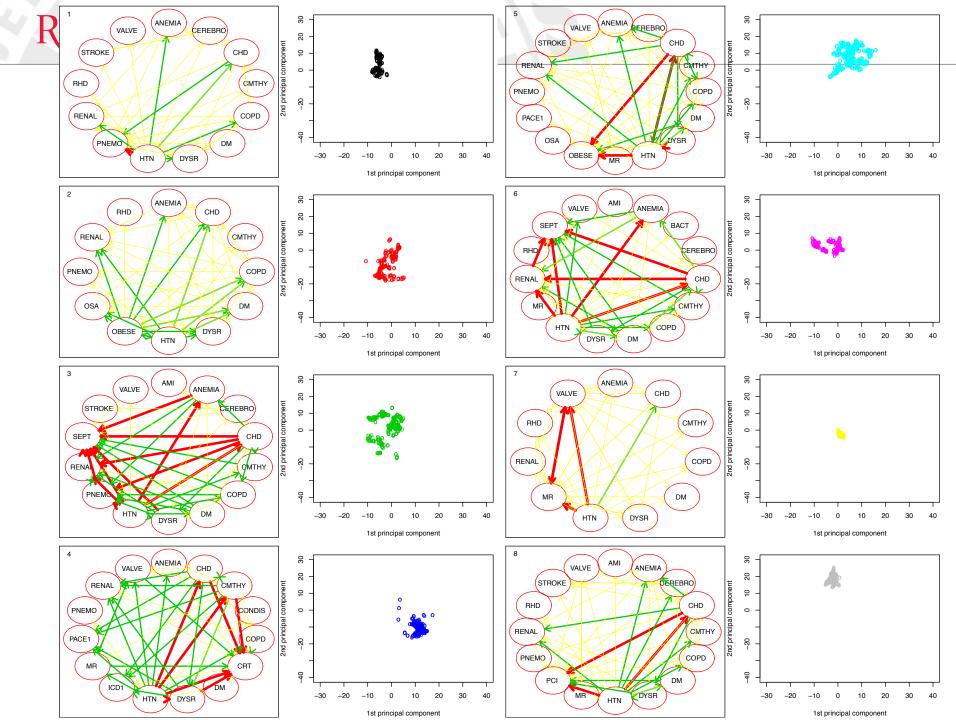
Mortality in each cluster



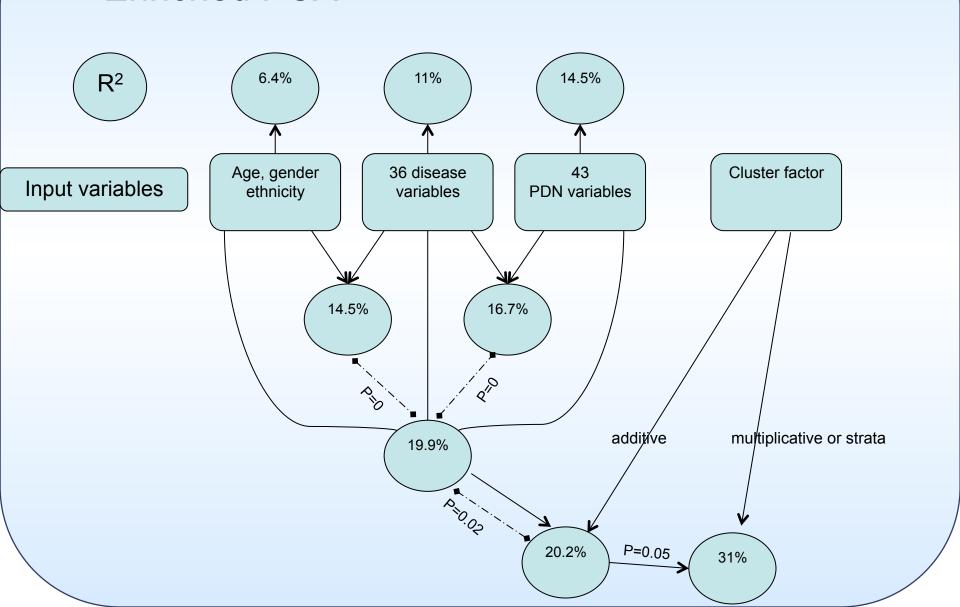
H Clustering with enriched nonlinear PCA.

tsna: R package L.J.P. van der Maaten e.a.(2008) Based on k-nearest neighbors.

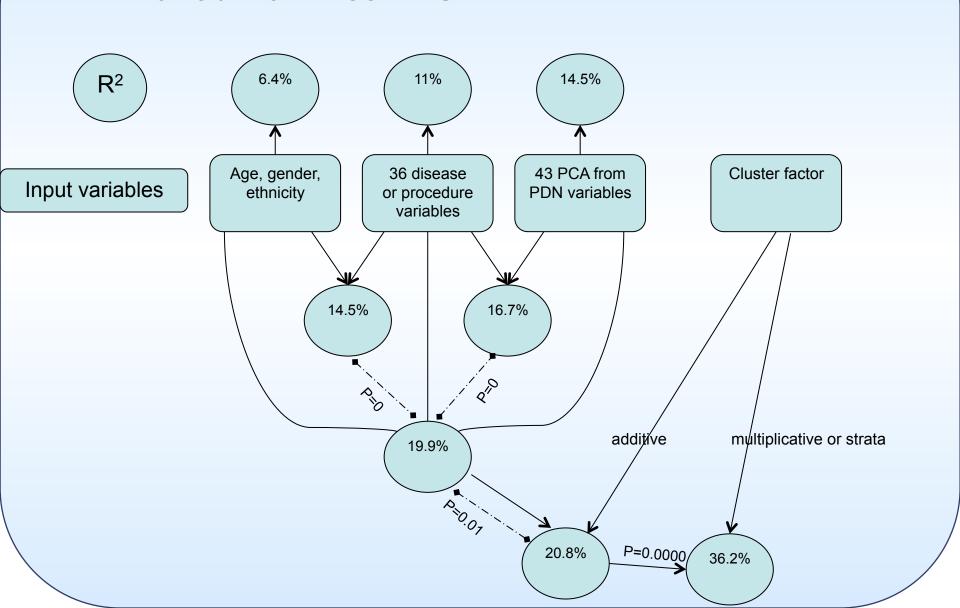




Modeling all Mortality by using Cox PH models Enriched PCA



Modeling all Mortality by using Cox PH models Enriched nonlinear PCA





Conclusions

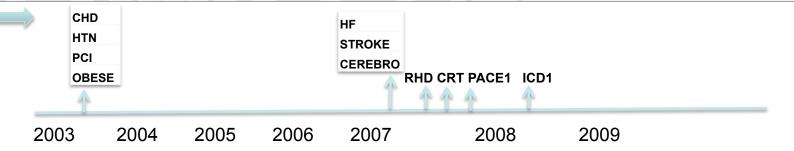
- We develop a novel algorithm for building patient level networks for prediction of medical outcomes
- We develop cluster methodologies (Hierarchical Cluster using enriched PCA scores) for individual network patient data and ways to summarize the clusters and study the within variability.
- In clinical outcome prediction, PDN's and cluster factors as input variables significantly improved goodness of fit of survival models in our dataset.

Future directions

- Optimize the network visualization. Time sequence, from horse shoes to more complex spatial network display.
- Network summaries. Median, mean, quantiles. Variability.
- Prediction of cardiovascular outcomes like cardiovascular death, Heart attack, Stroke, etc using the method. Training and testing sets. Cross validation.
- Apply the method for the larger subsets of MIDAS for other demographics and conditions.
- Simulations and
- Build an R package to implement our methodology.



RUTGERS Construction of Personalized disease networks



Example of personalized (patient 32) disease networks(PDN)

Patient 32

| CODE | DATE |
|---------|---------|
| CHD | 4/12/03 |
| HTN | 4/12/03 |
| PCI | 4/12/03 |
| OBESE | 4/12/03 |
| HF | 6/7/07 |
| STROKE | 6/7/07 |
| CEREBRO | 6/7/07 |
| RHD | 6/28/07 |
| CRT | 9/19/07 |
| PACE1 | 9/27/07 |
| ICD1 | 5/28/08 |

