

General Latent Feature Modeling for Data Exploration Tasks

Isabel Valera

University of Cambridge



Joint work with M. F. Pradier & Z. Ghahramani

Data exploration



???

- Before being able to solve, or even define, a predictive task we need to *explore* the data
- Helps checking assumptions required for model fitting and hypothesis testing

Data exploration

Objects

Attributes

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|----|---|---------|--------|--------|---------|----|--------|--------|
| 1 | 1 | 14.2300 | 1.7100 | 2.4300 | 15.6000 | 58 | 2.8000 | 3.0600 |
| 2 | 1 | 13.2000 | 1.7800 | 2.1400 | 11.2000 | 31 | 2.6500 | 2.7600 |
| 3 | 1 | 13.1600 | 2.3600 | 2.6700 | 18.6000 | 32 | 2.8000 | 3.2400 |
| 4 | 1 | 14.3700 | 1.9500 | 2.5000 | 16.8000 | 44 | 3.8500 | 3.4900 |
| 5 | 1 | 13.2400 | 2.5900 | 2.8700 | 21 | 49 | 2.8000 | 2.6900 |
| 6 | 1 | 14.2000 | 1.7600 | 2.4500 | 15.2000 | 43 | 3.2700 | 3.3900 |
| 7 | 1 | 14.3900 | 1.8700 | 2.4500 | 14.6000 | 27 | 2.5000 | 2.5200 |
| 8 | 1 | 14.0600 | 2.1500 | 2.6100 | 17.6000 | 52 | 2.6000 | 2.5100 |
| 9 | 1 | 14.8300 | 1.6400 | 2.1700 | 14 | 28 | 2.8000 | 2.9800 |
| 10 | 1 | 13.8600 | 1.3500 | 2.2700 | 16 | 29 | 2.9800 | 3.1500 |
| 11 | 1 | 14.1000 | 2.1600 | 2.3000 | 18 | 36 | 2.9500 | 3.3200 |
| 12 | 1 | 14.1200 | 1.4800 | 2.3200 | 16.8000 | 26 | 2.2000 | 2.4300 |
| 13 | 1 | 13.7500 | 1.7300 | 2.4100 | 16 | 20 | 2.6000 | 2.7600 |
| 14 | 1 | 14.7500 | 1.7300 | 2.3900 | 11.4000 | 22 | 3.1000 | 3.6900 |
| 15 | 1 | 14.3800 | 1.8700 | 2.3800 | 12 | 33 | 3.3000 | 3.6400 |
| 16 | 1 | 13.6300 | 1.8100 | 2.7000 | 17.2000 | 43 | 2.8500 | 2.9100 |
| 17 | 1 | 14.3000 | 1.9200 | 2.7200 | 20 | 51 | 2.8000 | 3.1400 |
| 18 | 1 | 13.8300 | 1.5700 | 2.6200 | 20 | 46 | 2.9500 | 3.4000 |
| 19 | 1 | 14.1900 | 1.5900 | 2.4800 | 16.5000 | 39 | 3.3000 | 3.9300 |
| 20 | 1 | 13.6400 | 3.1000 | 2.5600 | 15.2000 | 47 | 2.7000 | 3.0300 |
| 21 | 1 | 14.0600 | 1.6300 | 2.2800 | 16 | 57 | 3 | 3.1700 |
| 22 | 1 | 12.9300 | 3.8000 | 2.6500 | 18.6000 | 33 | 2.4100 | 2.4100 |
| 23 | 1 | 13.7100 | 1.8600 | 2.3600 | 16.6000 | 32 | 2.6100 | 2.8800 |
| 24 | 1 | 12.8500 | 1.6000 | 2.5200 | 17.8000 | 26 | 2.4800 | 2.3700 |
| 25 | 1 | 13.5000 | 1.8100 | 2.6100 | 20 | 27 | 2.5300 | 2.6100 |

Standard techniques:

Visualization

+

Unsupervised Learning

- Dimensionality reductions: PCA, clustering, etc.
- Visualization of the lower-dimensional representation

Data exploration

Objects

Attributes

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|----|---|---------|--------|--------|---------|----|--------|--------|
| 1 | 1 | 14.2300 | 1.7100 | 2.4300 | 15.6000 | 58 | 2.8000 | 3.0600 |
| 2 | 1 | 13.2000 | 1.7800 | 2.1400 | 11.2000 | 31 | 2.6500 | 2.7600 |
| 3 | 1 | 13.1600 | 2.3600 | 2.6700 | 18.6000 | 32 | 2.8000 | 3.2400 |
| 4 | 1 | 14.3700 | 1.9500 | 2.5000 | 16.8000 | 44 | 3.8500 | 3.4900 |
| 5 | 1 | 13.2400 | 2.5900 | 2.8700 | 21 | 49 | 2.8000 | 2.6900 |
| 6 | 1 | 14.2000 | 1.7600 | 2.4500 | 15.2000 | 43 | 3.2700 | 3.3900 |
| 7 | 1 | 14.3900 | 1.8700 | 2.4500 | 14.6000 | 27 | 2.5000 | 2.5200 |
| 8 | 1 | 14.0600 | 2.1500 | 2.6100 | 17.6000 | 52 | 2.6000 | 2.5100 |
| 9 | 1 | 14.8300 | 1.6400 | 2.1700 | 14 | 28 | 2.8000 | 2.9800 |
| 10 | 1 | 13.8600 | 1.3500 | 2.2700 | 16 | 29 | 2.9800 | 3.1500 |
| 11 | 1 | 14.1000 | 2.1600 | 2.3000 | 18 | 36 | 2.9500 | 3.3200 |
| 12 | 1 | 14.1200 | 1.4800 | 2.3200 | 16.8000 | 26 | 2.2000 | 2.4300 |
| 13 | 1 | 13.7500 | 1.7300 | 2.4100 | 16 | 20 | 2.6000 | 2.7600 |
| 14 | 1 | 14.7500 | 1.7300 | 2.3900 | 11.4000 | 22 | 3.1000 | 3.6900 |
| 15 | 1 | 14.3800 | 1.8700 | 2.3800 | 12 | 33 | 3.3000 | 3.6400 |
| 16 | 1 | 13.6300 | 1.8100 | 2.7000 | 17.2000 | 43 | 2.8500 | 2.9100 |
| 17 | 1 | 14.3000 | 1.9200 | 2.7200 | 20 | 51 | 2.8000 | 3.1400 |
| 18 | 1 | 13.8300 | 1.5700 | 2.6200 | 20 | 46 | 2.9500 | 3.4000 |
| 19 | 1 | 14.1900 | 1.5900 | 2.4800 | 16.5000 | 39 | 3.3000 | 3.9300 |
| 20 | 1 | 13.6400 | 3.1000 | 2.5600 | 15.2000 | 47 | 2.7000 | 3.0300 |
| 21 | 1 | 14.0600 | 1.6300 | 2.2800 | 16 | 57 | 3 | 3.1700 |
| 22 | 1 | 12.9300 | 3.8000 | 2.6500 | 18.6000 | 33 | 2.4100 | 2.4100 |
| 23 | 1 | 13.7100 | 1.8600 | 2.3600 | 16.6000 | 32 | 2.6100 | 2.8800 |
| 24 | 1 | 12.8500 | 1.6000 | 2.5200 | 17.8000 | 26 | 2.4800 | 2.3700 |
| 25 | 1 | 13.5000 | 1.8100 | 2.6100 | 20 | 27 | 2.5300 | 2.6100 |

Continuous

Discrete

- Real-valued
- Positive real
- Interval

- Categorical
- Ordinal
- Count

Challenge – Heterogeneous Attributes

Example – User survey contains information on gender, age, education level, income...

Data exploration

Objects

Attributes

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|----|---|---------|--------|--------|---------|----|--------|--------|
| 1 | 1 | 14.2300 | 1.7100 | 2.4300 | 15.6000 | 58 | 2.8000 | 3.0600 |
| 2 | 1 | 13.2000 | 1.7800 | 2.1400 | 11.2000 | 31 | 2.6500 | 2.7600 |
| 3 | 1 | 13.1600 | 2.3600 | 2.6700 | 18.6000 | 32 | 2.8000 | 3.2400 |
| 4 | 1 | 14.3700 | 1.9500 | 2.5000 | 16.8000 | 44 | 3.8500 | 3.4900 |
| 5 | 1 | 13.2400 | 2.5900 | 2.8700 | 21 | 49 | 2.8000 | 2.6900 |
| 6 | 1 | 14.2000 | 1.7600 | 2.4500 | 15.2000 | 43 | 3.2700 | 3.3900 |
| 7 | 1 | 14.3900 | 1.8700 | 2.4500 | 14.6000 | 27 | 2.5000 | 2.5200 |
| 8 | 1 | 14.0600 | 2.1500 | 2.6100 | 17.6000 | 52 | 2.6000 | 2.5100 |
| 9 | 1 | 14.8300 | 1.6400 | 2.1700 | 14 | 28 | 2.8000 | 2.9800 |
| 10 | 1 | 13.8600 | 1.3500 | 2.2700 | 16 | 29 | 2.9800 | 3.1500 |
| 11 | 1 | 14.1000 | 2.1600 | 2.3000 | 18 | 36 | 2.9500 | 3.3200 |
| 12 | 1 | 14.1200 | 1.4800 | 2.3200 | 16.8000 | 26 | 2.2000 | 2.4300 |
| 13 | 1 | 13.7500 | 1.7300 | 2.4100 | 16 | 20 | 2.6000 | 2.7600 |
| 14 | 1 | 14.7500 | 1.7300 | 2.3900 | 11.4000 | 22 | 3.1000 | 3.6900 |
| 15 | 1 | 14.3800 | 1.8700 | 2.3800 | 12 | 33 | 3.3000 | 3.6400 |
| 16 | 1 | 13.6300 | 1.8100 | 2.7000 | 17.2000 | 43 | 2.8500 | 2.9100 |
| 17 | 1 | 14.3000 | 1.9200 | 2.7200 | 20 | 51 | 2.8000 | 3.1400 |
| 18 | 1 | 13.8300 | 1.5700 | 2.6200 | 20 | 46 | 2.9500 | 3.4000 |
| 19 | 1 | 14.1900 | 1.5900 | 2.4800 | 16.5000 | 39 | 3.3000 | 3.9300 |
| 20 | 1 | 13.6400 | 3.1000 | 2.5600 | 15.2000 | 47 | 2.7000 | 3.0300 |
| 21 | 1 | 14.0600 | 1.6300 | 2.2800 | 16 | 57 | 3 | 3.1700 |
| 22 | 1 | 12.9300 | 3.8000 | 2.6500 | 18.6000 | 33 | 2.4100 | 2.4100 |
| 23 | 1 | 13.7100 | 1.8600 | 2.3600 | 16.6000 | 32 | 2.6100 | 2.8800 |
| 24 | 1 | 12.8500 | 1.6000 | 2.5200 | 17.8000 | 26 | 2.4800 | 2.3700 |
| 25 | 1 | 13.5000 | 1.8100 | 2.6100 | 20 | 27 | 2.5300 | 2.6100 |

Our approach:
Visualization
+
General Latent
Feature Modeling

Challenge – Heterogeneous Attributes

Example – User survey contains information on gender, age, education level, income...

General latent feature modeling

Suitable for data exploratory analysis?

Objects

Attributes

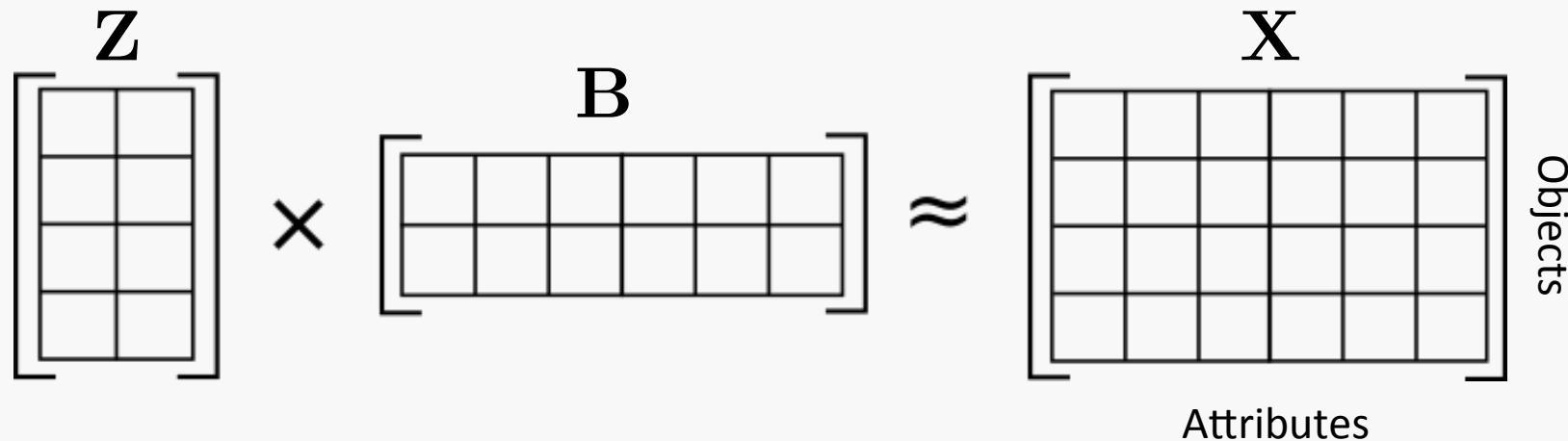
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|----|---|---------|--------|--------|---------|----|--------|--------|
| 1 | 1 | 14.2300 | 1.7100 | 2.4300 | 15.6000 | 58 | 2.8000 | 3.0600 |
| 2 | 1 | 13.2000 | 1.7800 | 2.1400 | 11.2000 | 31 | 2.6500 | 2.7600 |
| 3 | 1 | 13.1600 | 2.3600 | 2.6700 | 18.6000 | 32 | 2.8000 | 3.2400 |
| 4 | 1 | 14.3700 | 1.9500 | 2.5000 | 16.8000 | 44 | 3.8500 | 3.4900 |
| 5 | 1 | 13.2400 | 2.5900 | 2.8700 | 21 | 49 | 2.8000 | 2.6900 |
| 6 | 1 | 14.2000 | 1.7600 | 2.4500 | 15.2000 | 43 | 3.2700 | 3.3900 |
| 7 | 1 | 14.3900 | 1.8700 | 2.4500 | 14.6000 | 27 | 2.5000 | 2.5200 |
| 8 | 1 | 14.0600 | 2.1500 | 2.6100 | 17.6000 | 52 | 2.6000 | 2.5100 |
| 9 | 1 | 14.8300 | 1.6400 | 2.1700 | 14 | 28 | 2.8000 | 2.9800 |
| 10 | 1 | 13.8600 | 1.3500 | 2.2700 | 16 | 29 | 2.9800 | 3.1500 |
| 11 | 1 | 14.1000 | 2.1600 | 2.3000 | 18 | 36 | 2.9500 | 3.3200 |
| 12 | 1 | 14.1200 | 1.4800 | 2.3200 | 16.8000 | 26 | 2.2000 | 2.4300 |
| 13 | 1 | 13.7500 | 1.7300 | 2.4100 | 16 | 20 | 2.6000 | 2.7600 |
| 14 | 1 | 14.7500 | 1.7300 | 2.3900 | 11.4000 | 22 | 3.1000 | 3.6900 |
| 15 | 1 | 14.3800 | 1.8700 | 2.3800 | 12 | 33 | 3.3000 | 3.6400 |
| 16 | 1 | 13.6300 | 1.8100 | 2.7000 | 17.2000 | 43 | 2.8500 | 2.9100 |
| 17 | 1 | 14.3000 | 1.9200 | 2.7200 | 20 | 51 | 2.8000 | 3.1400 |
| 18 | 1 | 13.8300 | 1.5700 | 2.6200 | 20 | 46 | 2.9500 | 3.4000 |
| 19 | 1 | 14.1900 | 1.5900 | 2.4800 | 16.5000 | 39 | 3.3000 | 3.9300 |
| 20 | 1 | 13.6400 | 3.1000 | 2.5600 | 15.2000 | 47 | 2.7000 | 3.0300 |
| 21 | 1 | 14.0600 | 1.6300 | 2.2800 | 16 | 57 | 3 | 3.1700 |
| 22 | 1 | 12.9300 | 3.8000 | 2.6500 | 18.6000 | 33 | 2.4100 | 2.4100 |
| 23 | 1 | 13.7100 | 1.8600 | 2.3600 | 16.6000 | 32 | 2.6100 | 2.8800 |
| 24 | 1 | 12.8500 | 1.6000 | 2.5200 | 17.8000 | 26 | 2.4800 | 2.3700 |
| 25 | 1 | 13.5000 | 1.8100 | 2.6100 | 20 | 27 | 2.5300 | 2.6100 |

- Latent structure in the data
- Dependencies among objects and attributes
- Compact in a few features the immense redundant information

General latent feature modeling

How does it work?

Latent structure – Low-rank representation



$$p(\mathbf{X}|\mathbf{Z}, \mathbf{B}) = \prod_{d=1}^D p(\mathbf{x}_d|\mathbf{Z}, \mathbf{b}_d)$$

General latent feature modeling

How does it work?

Interpretability – Objects as binary latent features

$$\begin{array}{c} \mathbf{z}_n \\ \text{Latent features explaining each object} \\ \hline \begin{matrix} \text{Avatar 1} & [1 \ 0 \ 1 \ 1 \ \dots] \\ \text{Avatar 2} & [0 \ 0 \ 0 \ 1 \ \dots] \\ \text{Avatar 3} & [1 \ 0 \ 0 \ 1 \ \dots] \end{matrix} \\ \hline \end{array} \times \begin{bmatrix} 1.2 \\ 3.2 \\ -2.4 \\ 0.15 \\ \vdots \end{bmatrix} = \begin{array}{c} -1.05 \\ 0.15 \\ 1.35 \\ \hline \text{Observed values} \end{array}$$

The diagram illustrates the calculation of observed values from latent features and weights. It shows three objects (Avatars) with their corresponding latent feature vectors (\mathbf{z}_n). These vectors are multiplied by a weight vector (\mathbf{b}^d) to produce the observed values. Red brackets at the bottom group the latent features, the weight vector, and the observed values.

General latent feature modeling

How does it work?

Heterogeneous data – Data transformation

$$\begin{array}{llll} \mathbf{z}_n & \mathbf{B}^d & y_n^d & x_n^d \\ [1 \ 0 \ 1 \ 0 \ \dots] & \begin{bmatrix} 1.2 \\ 3.2 \\ -2.4 \\ 0.15 \\ \vdots \end{bmatrix} & -1.5 & 1 \ (\text{"low"}) \\ [0 \ 0 \ 0 \ 1 \ \dots] & \times & 0.15 & \xrightarrow{f_d(\cdot)} 2 \ (\text{"medium"}) \\ [1 \ 0 \ 0 \ 1 \ \dots] & & 1.35 & 2 \ (\text{"medium"}) \end{array}$$

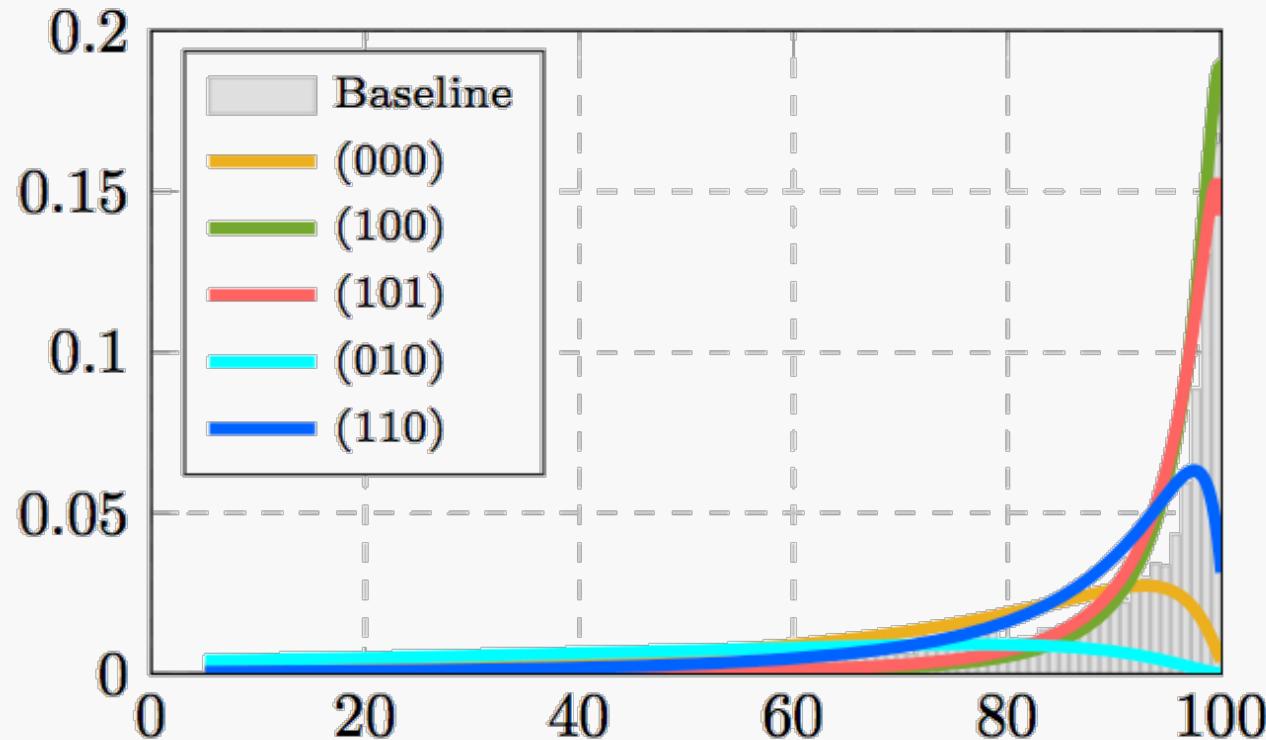
Low-rank representation Data transformation

The diagram illustrates the decomposition of a low-rank representation \mathbf{z}_n into a product of a row vector and a matrix \mathbf{B}^d , followed by a data transformation $f_d(\cdot)$ resulting in a scalar y_n^d and a categorical value x_n^d . The transformation $f_d(\cdot)$ maps the scalar y_n^d to a categorical value. Brackets at the bottom group the first two columns as 'Low-rank representation' and the last two columns as 'Data transformation'.

General latent feature modeling

How does it work?

Visualization – We can compute and visualize
 $p(\text{attribute} | \text{feature vector})$, i.e., $p(x_n^d | \mathbf{z}_n, \mathbf{B}^d)$



Application to Clinical Trials

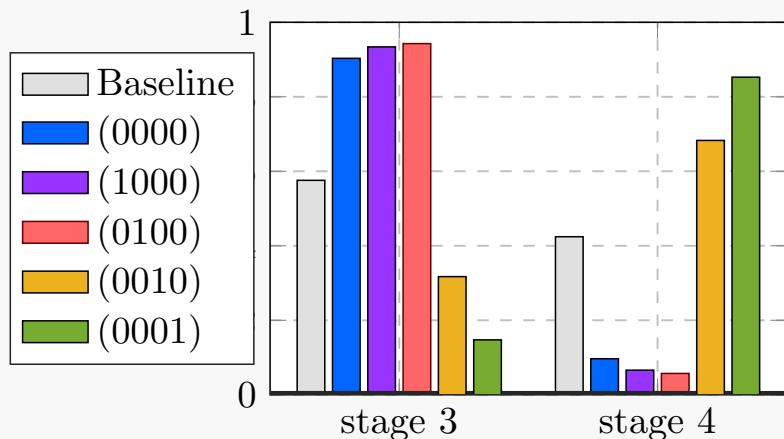
Goal: effects of a new drug for prostate cancer

Dataset contains 502 patients and 16 attributes, from which we make use of the following attributes:

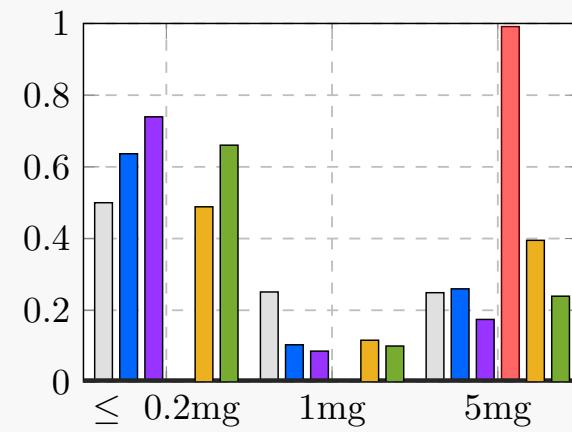
| Attribute description | Type of variable |
|---|-------------------------------|
| Stage of the cancer | Categorical with 2 categories |
| DES treatment level | Ordinal with 3 categories |
| Tumor size in cm ² | Count data |
| Serum Prostatic Acid Phosphatase (PAP) | Positive real-valued |
| Prognosis Status (outcome of the disease) | Categorical with 4 categories |

Application to Clinical Trials

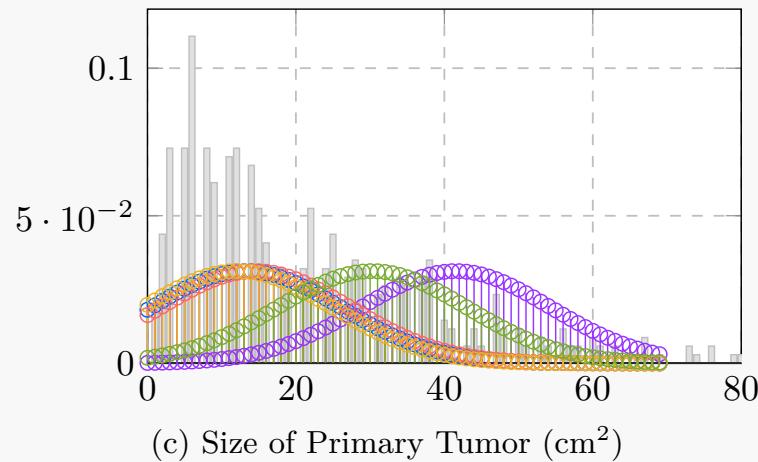
Feature patterns:



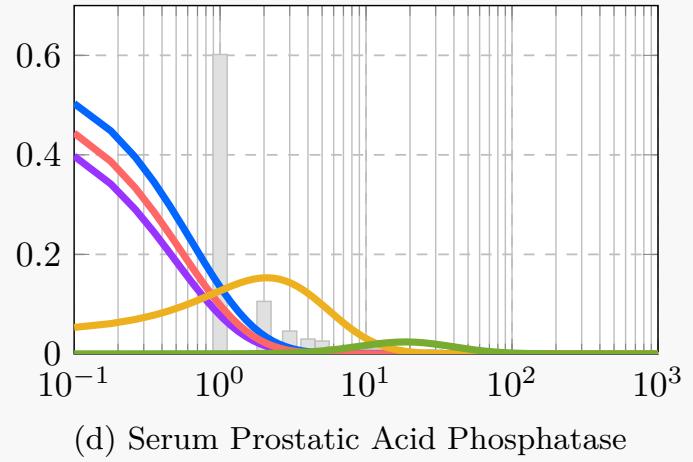
(a) Type of Cancer



(b) Drug Level



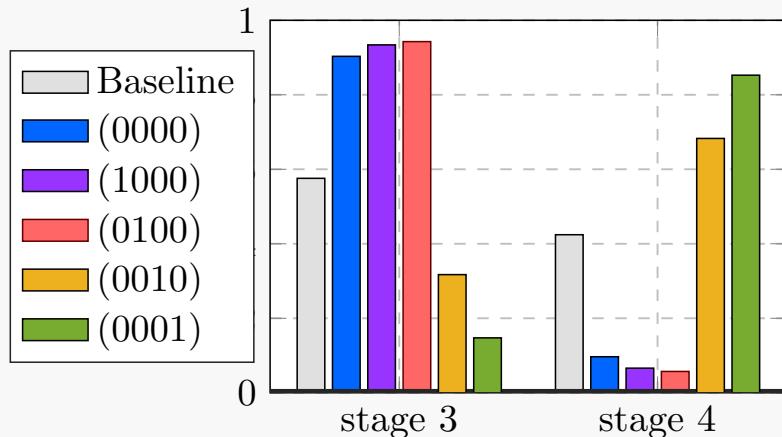
(c) Size of Primary Tumor (cm²)



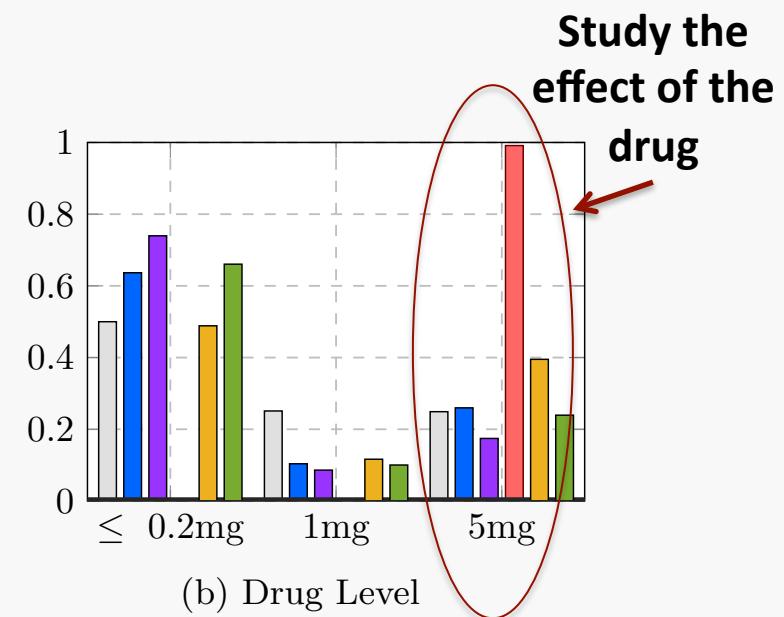
(d) Serum Prostatic Acid Phosphatase

Application to Clinical Trials

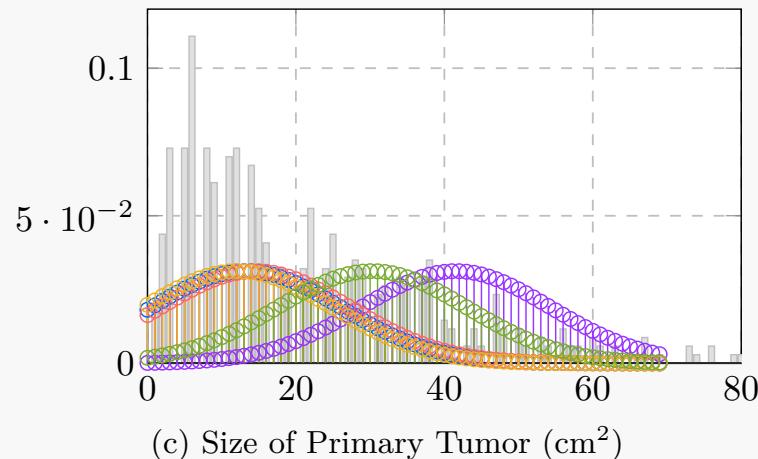
Feature patterns:



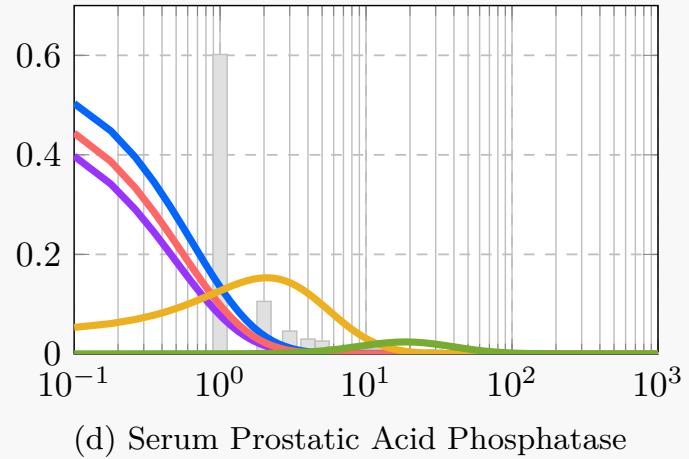
(a) Type of Cancer



(b) Drug Level



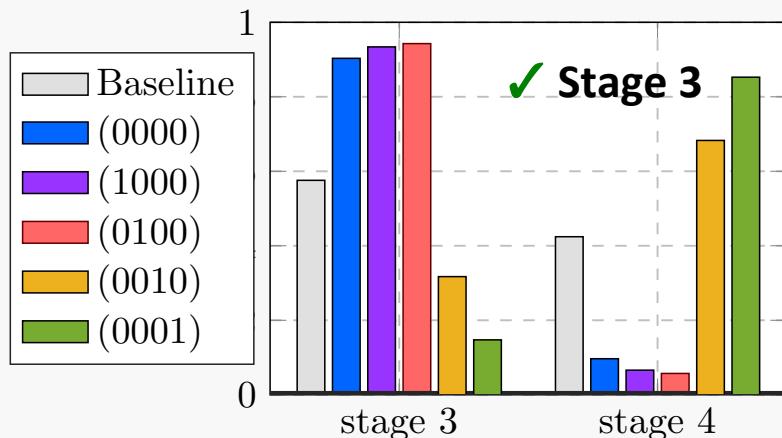
(c) Size of Primary Tumor (cm²)



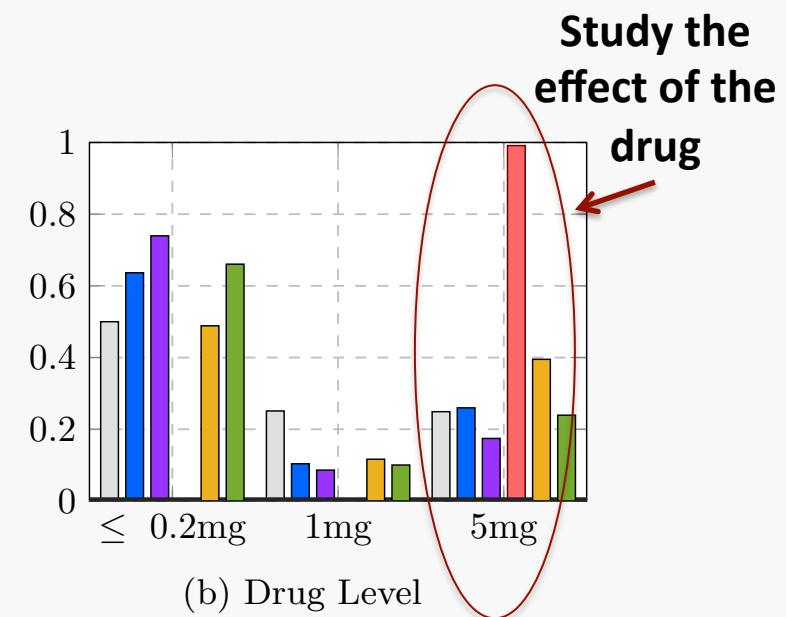
(d) Serum Prostatic Acid Phosphatase

Application to Clinical Trials

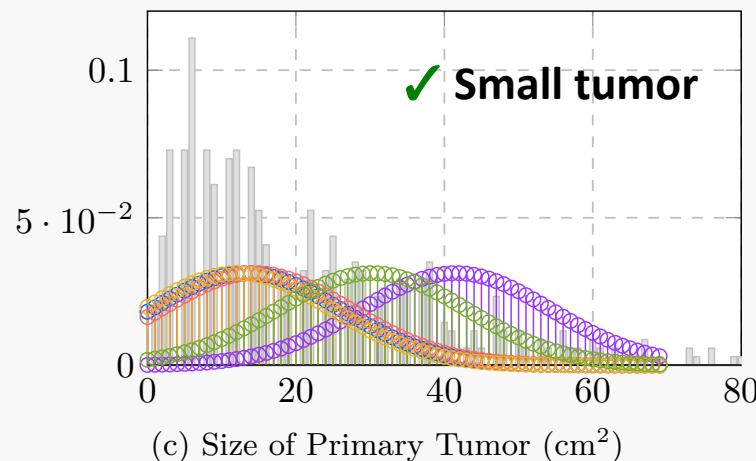
Feature patterns:



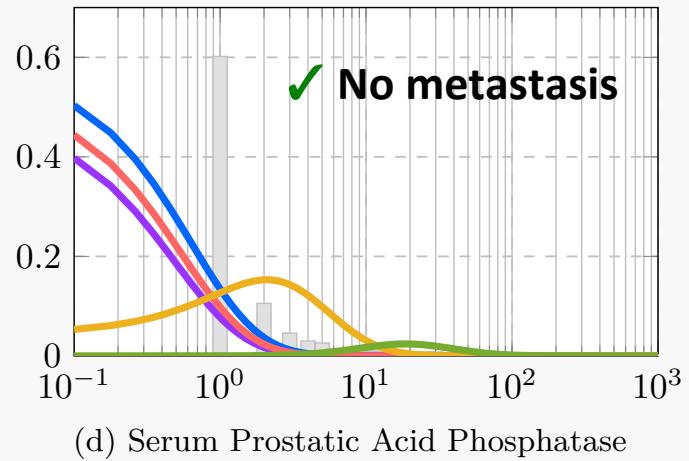
(a) Type of Cancer



(b) Drug Level



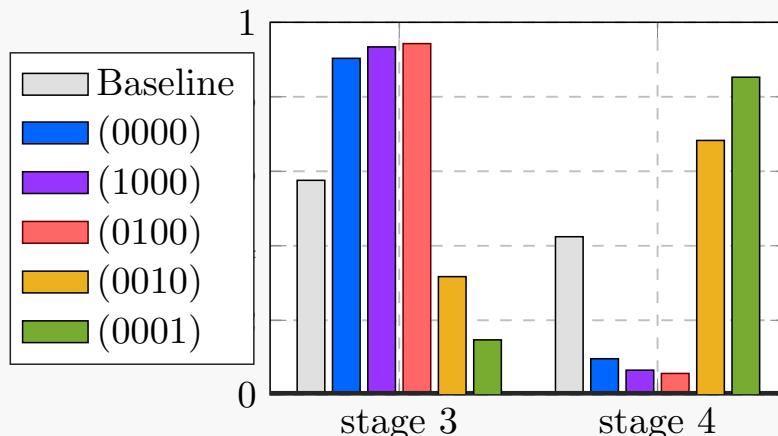
(c) Size of Primary Tumor (cm^2)



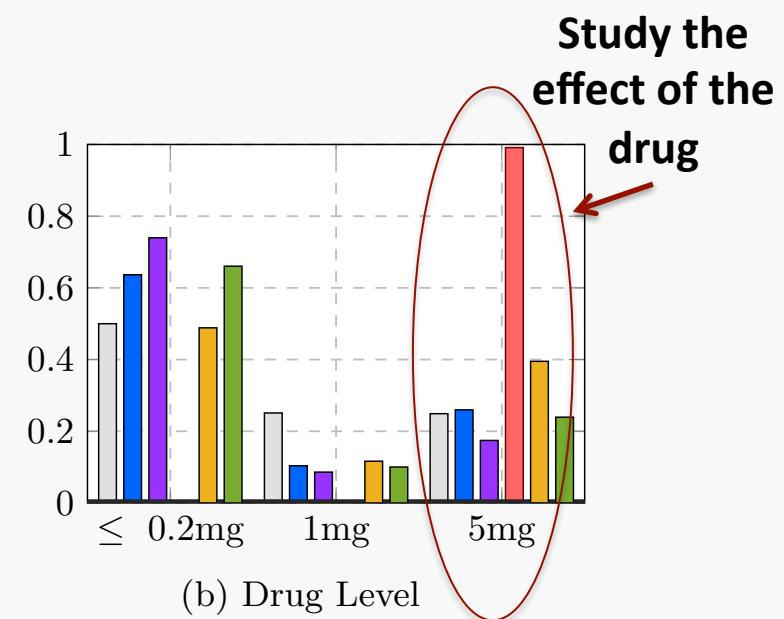
(d) Serum Prostatic Acid Phosphatase

Application to Clinical Trials

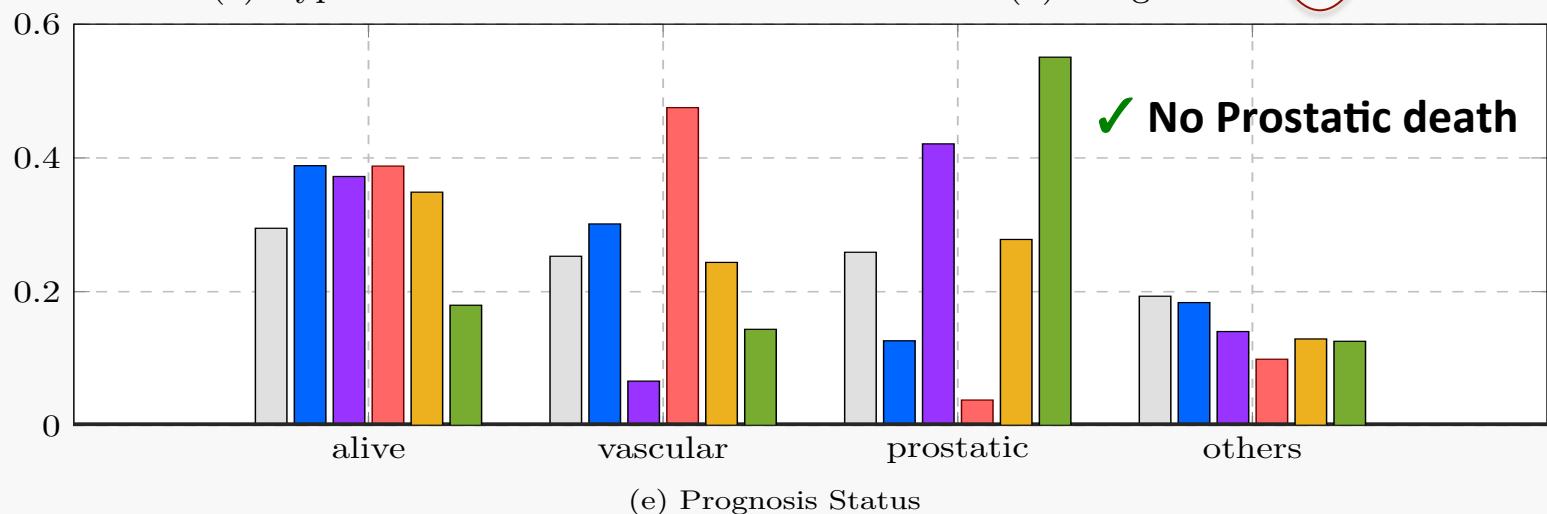
Feature patterns:



(a) Type of Cancer



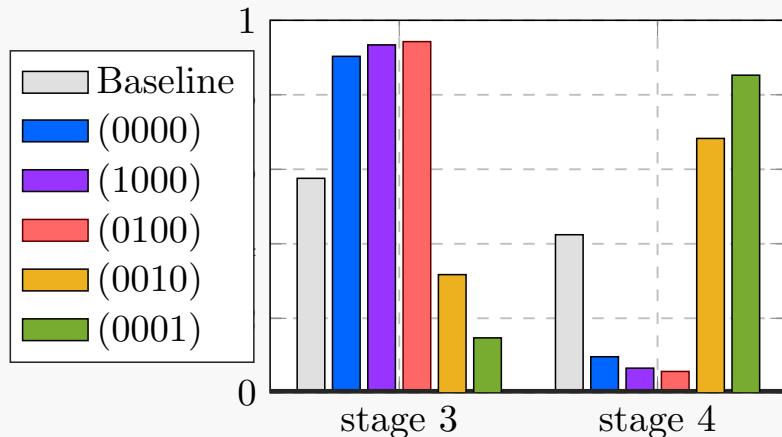
(b) Drug Level



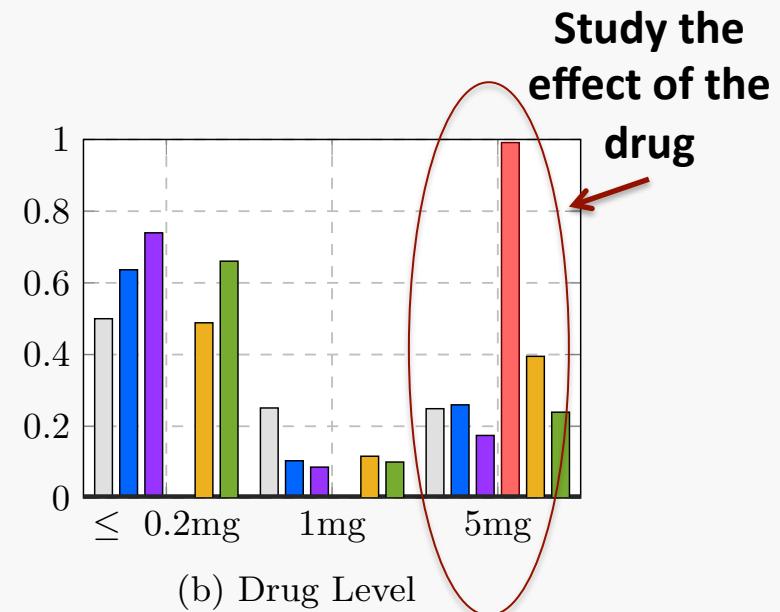
(e) Prognosis Status

Application to Clinical Trials

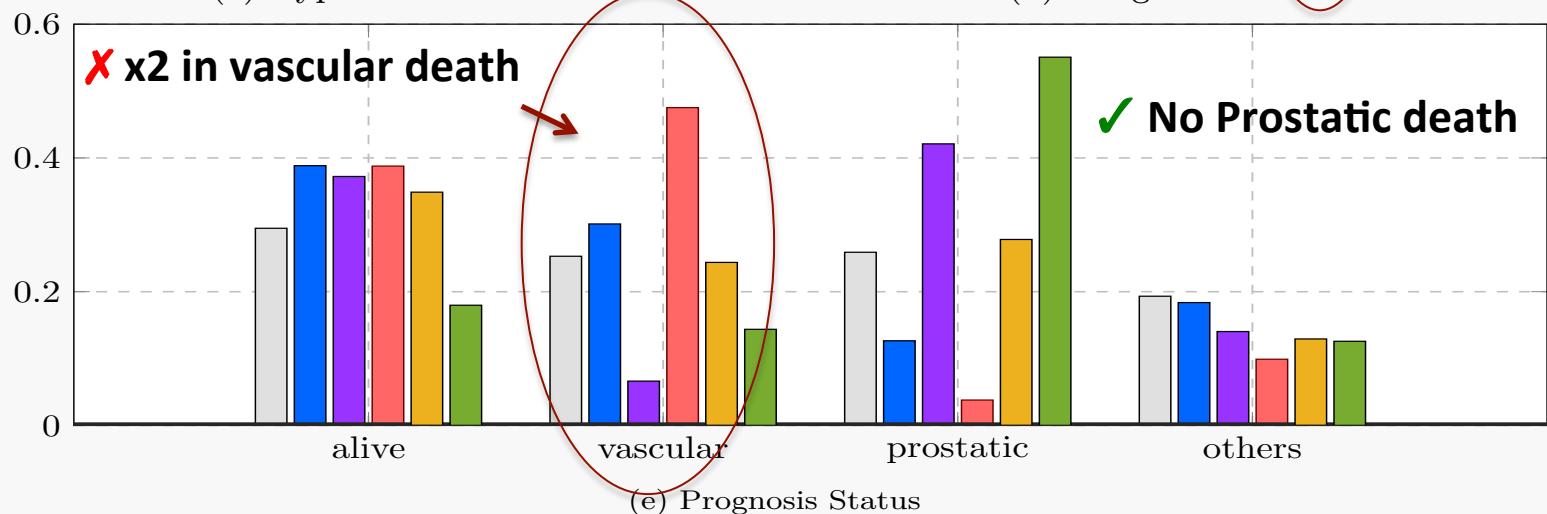
Feature patterns:



(a) Type of Cancer



(b) Drug Level



(e) Prognosis Status

Further results

- 1) Further details: <https://arxiv.org/abs/1706.03779>
 - i. Other datasets: Psychiatry, electoral elections...
 - ii. Other applications: Missing data estimation
- 2) Available Software: <https://ivaleram.github.io/GLFM/>
 - i. User interfaces in Matlab, Python and R.
 - ii. Available functions for:
 - Inference
 - Predictions
 - Visualization

Further results

- 1) Further details: <https://arxiv.org/abs/1706.03779>
 - i. Other datasets: Psychiatry, electoral elections...
 - ii. Other applications: Missing data estimation
- 2) Available Software: <https://ivaleram.github.io/GLFM/>
 - i. User interfaces in Matlab, Python and R.
 - ii. Available functions for:
 - Inference
 - Predictions
 - Visualization

Thanks!