

Nonparametric Bayesian Multi-facet Clustering for Longitudinal Data

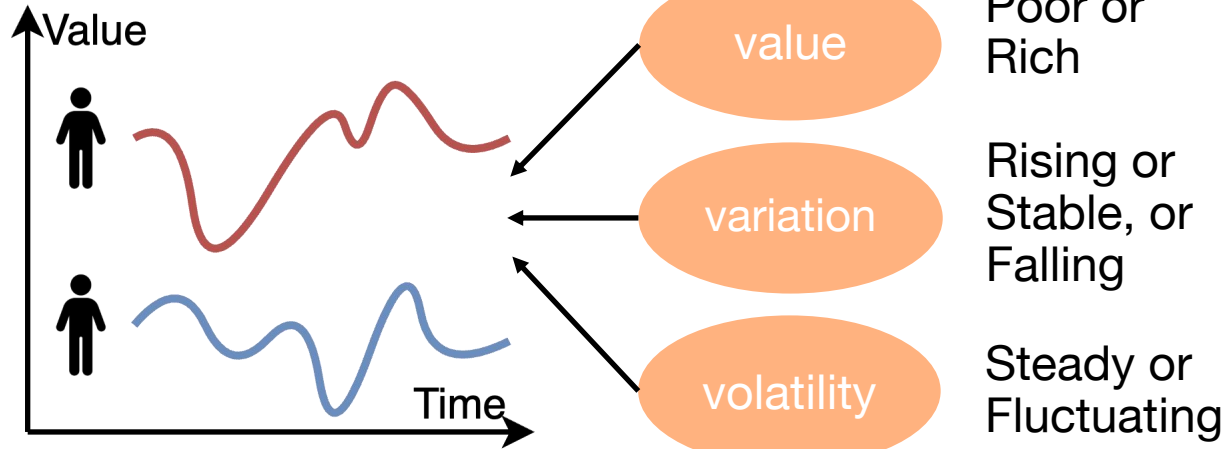
Luwei Wang, Kieran Richards, Sohan Seth

Data Science Unit, School of Informatics, University of Edinburgh

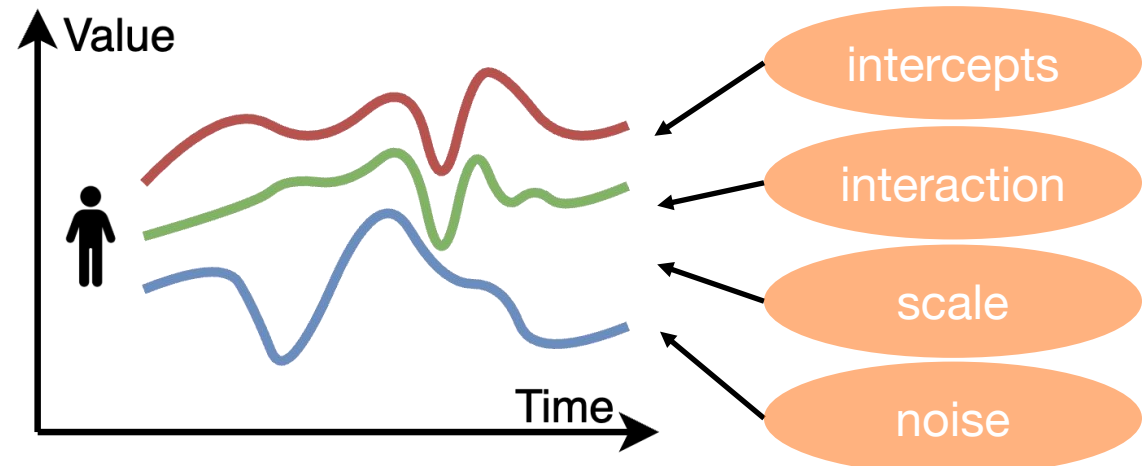


Multiple facets of longitudinal data

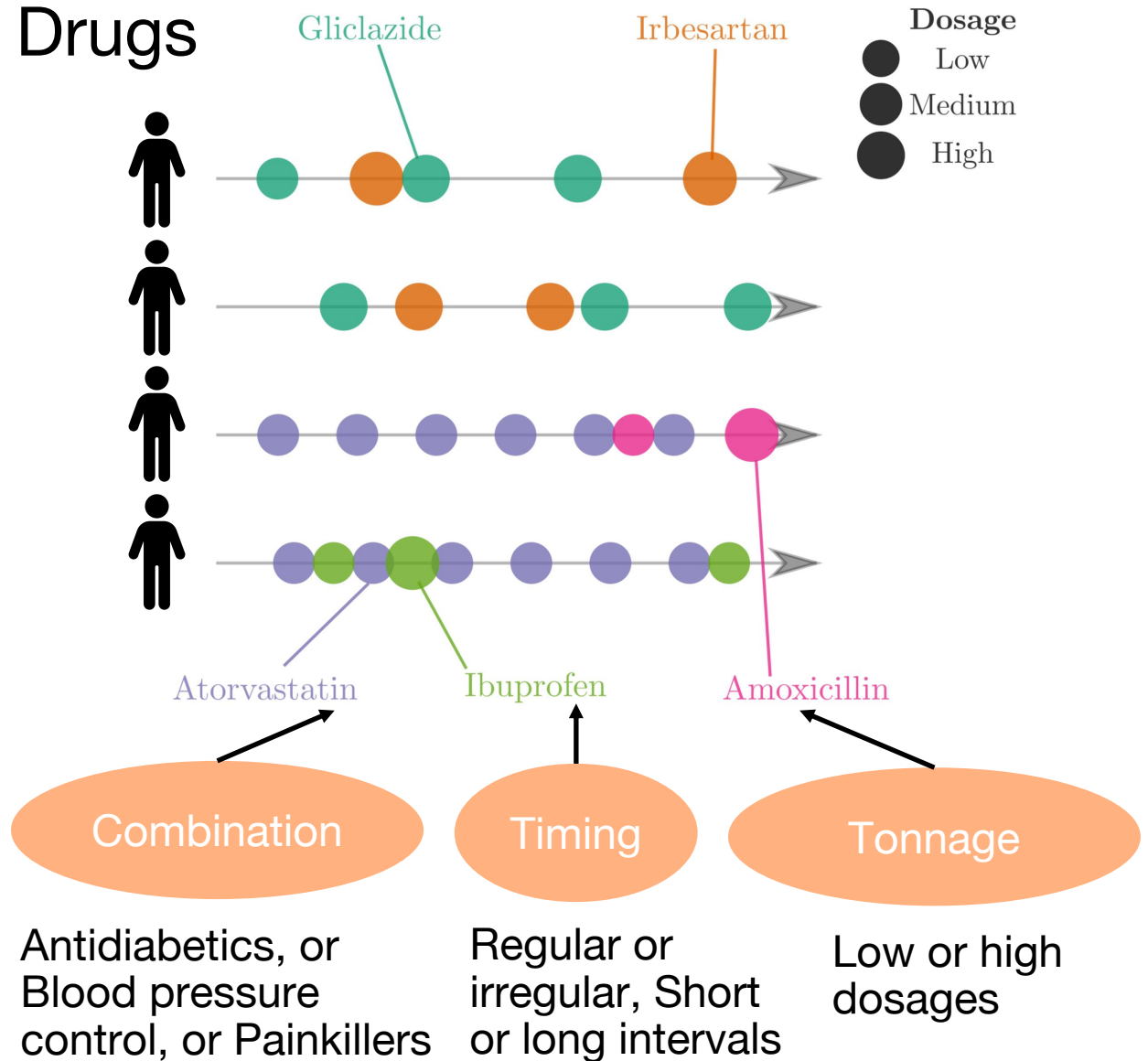
Income

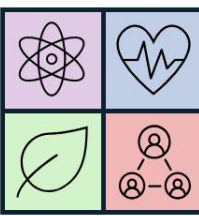


Frailty, wellbeing & social isolation



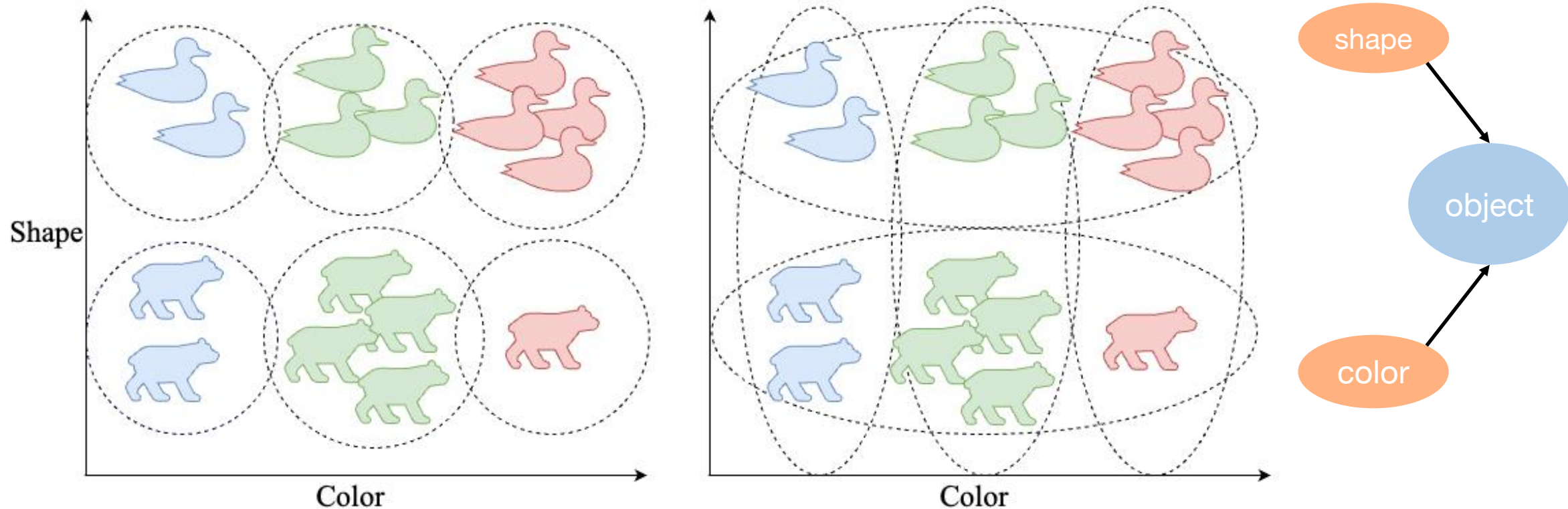
Drugs



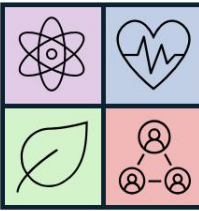


Disentangle facets for interpretability

- **Traditional clustering** aggregates all facets/characteristics
→ **poor interpretability**
- **Multi-facet clustering:** along different facets/characteristics



Single-facet versus multi-facet clustering.

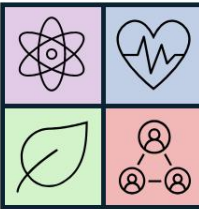


Method

- **Multi-Facet Mixture Model (MMM):** separate clustering per facet

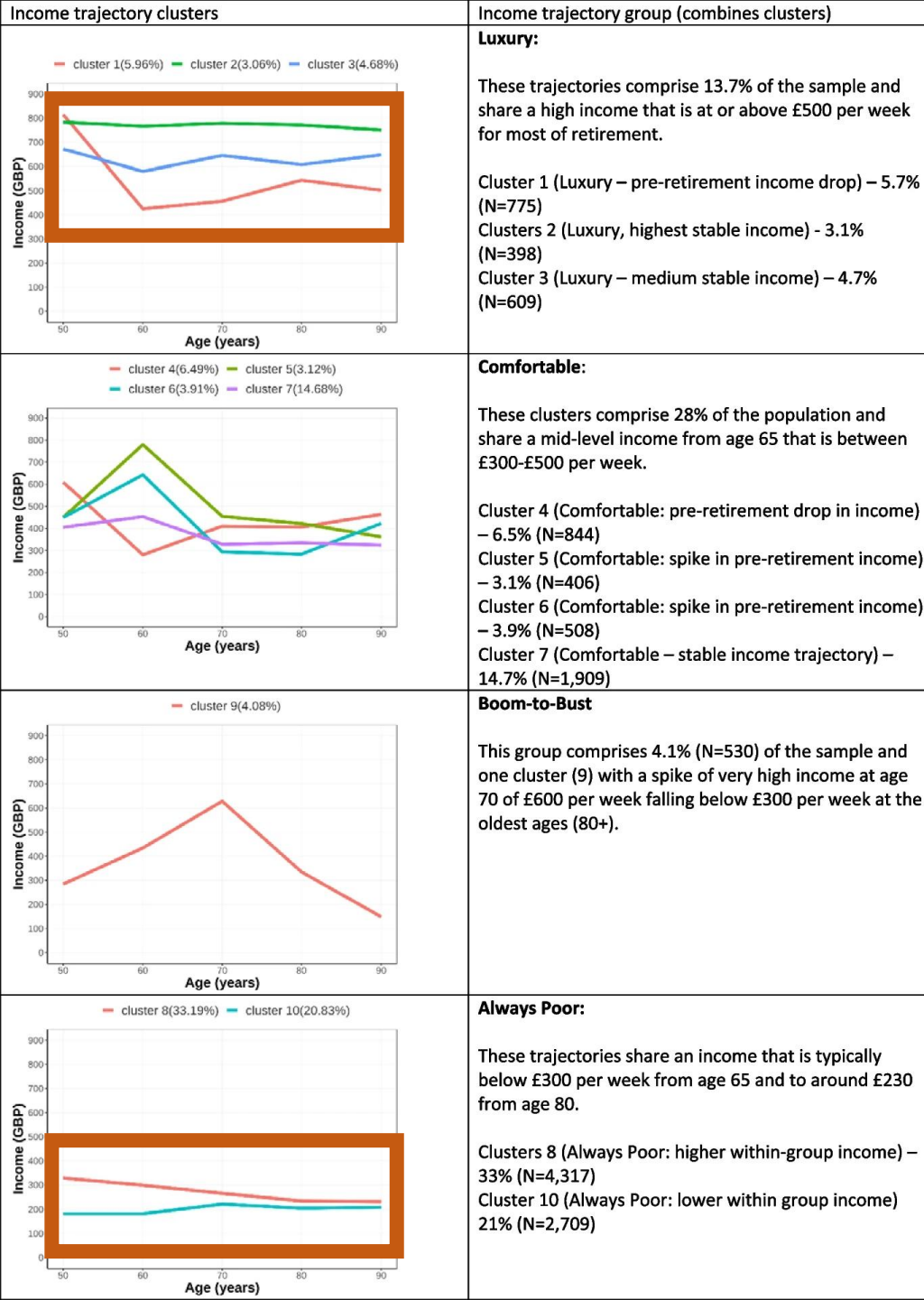
$$\sum_{k_1=1}^{K_1} \cdots \sum_{k_F=1}^{K_F} \pi_{k_1}^{(1)} \cdots \pi_{k_F}^{(F)} p(\mathbf{y} \mid \boldsymbol{\theta}_{k_1}^{(1)}, \dots, \boldsymbol{\theta}_{k_F}^{(F)})$$

- Independent facets apriori
- Time series data i.e. $y(t)$
- Facets defined by users
- *Dirichlet Process prior* for automatic cluster discovery
- *Variational Bayesian inference* for scalability vs MCMC
- Implemented for two time series models
 - Nonlinear Growth Model (NLG) & Vector Autoregressive Model (VAR)
 - More being added!



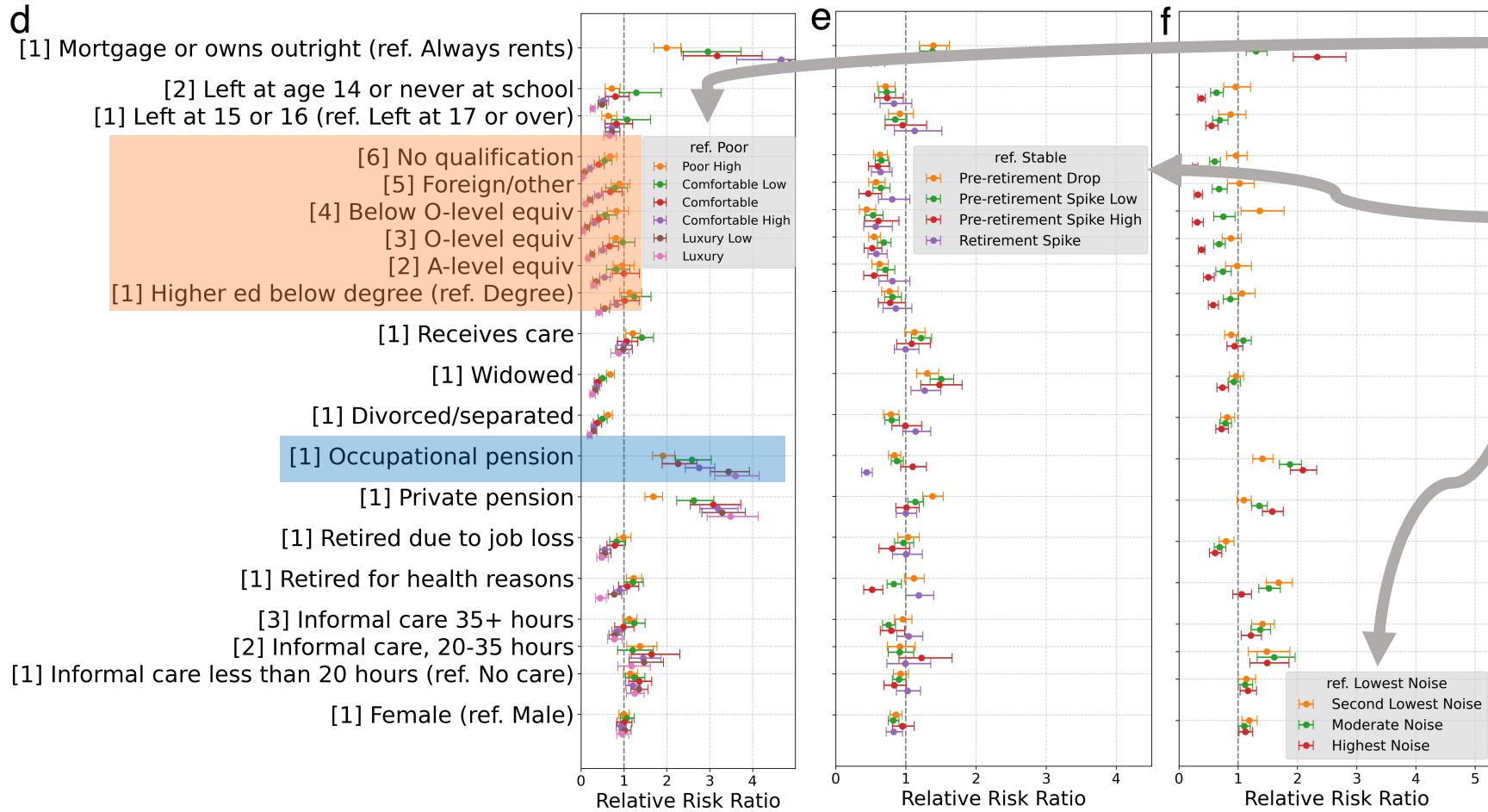
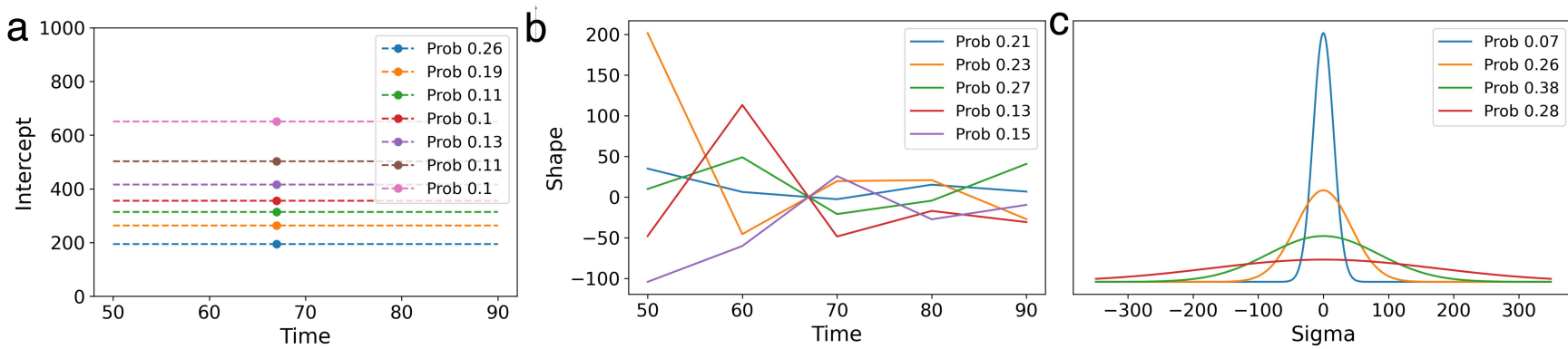
Single-facet clustering

- **English Longitudinal Study of Ageing (ELSA)**
 - N = 13,002
 - Age: 50-90
 - Missingness: 86.8%
- **10** income trajectory clusters
- Manually grouped to **4 super-clusters** based on income at retirement
 - Luxury
 - Comfortable
 - Boom-to-bust
 - Always Poor
- **Similar variations** of trajectory
- **No volatility**





Income clusters & their drivers



7 value clusters
(Poor → Luxury)

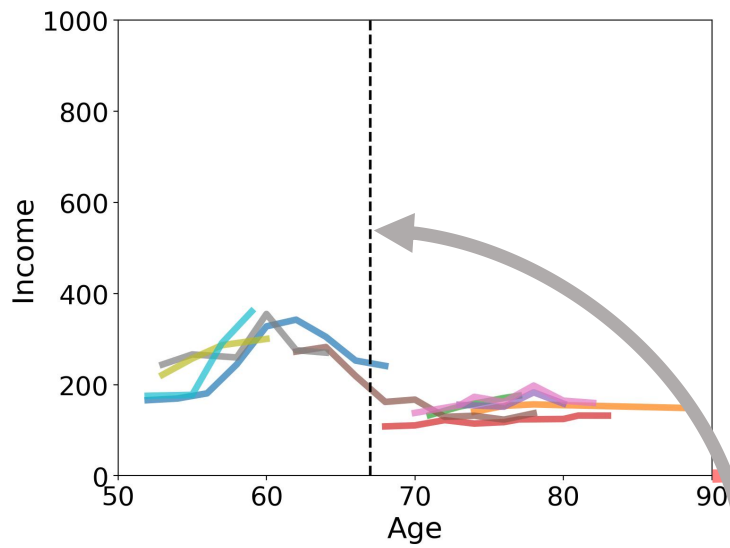
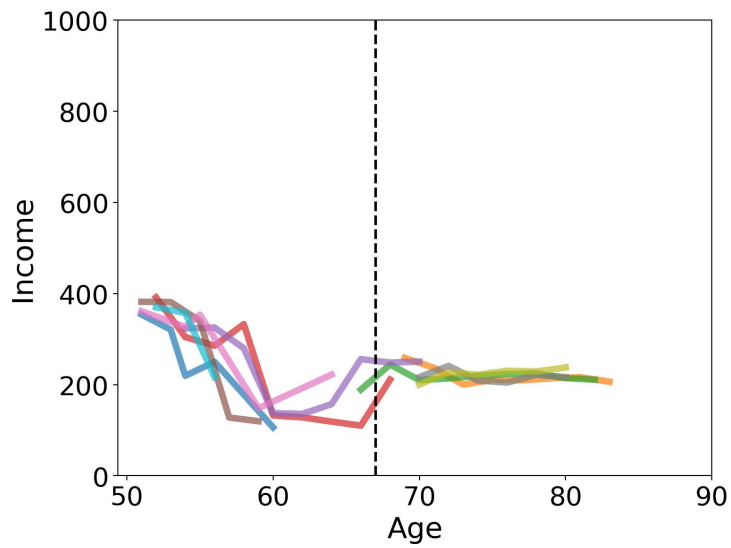
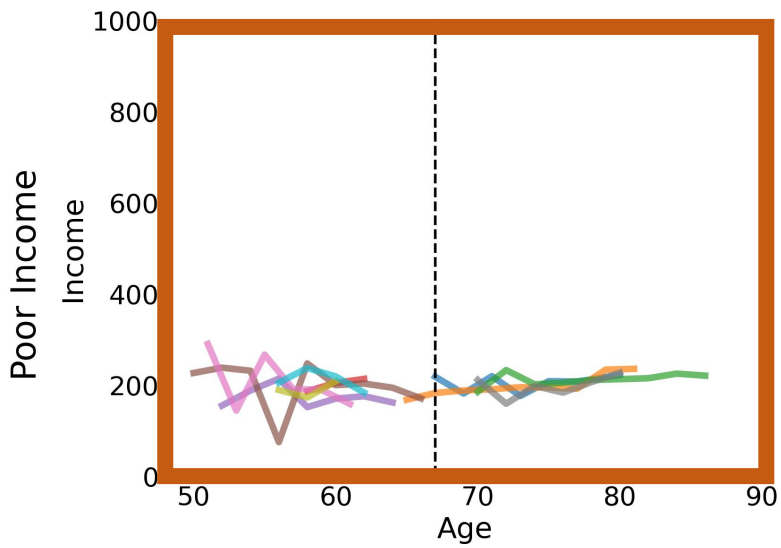
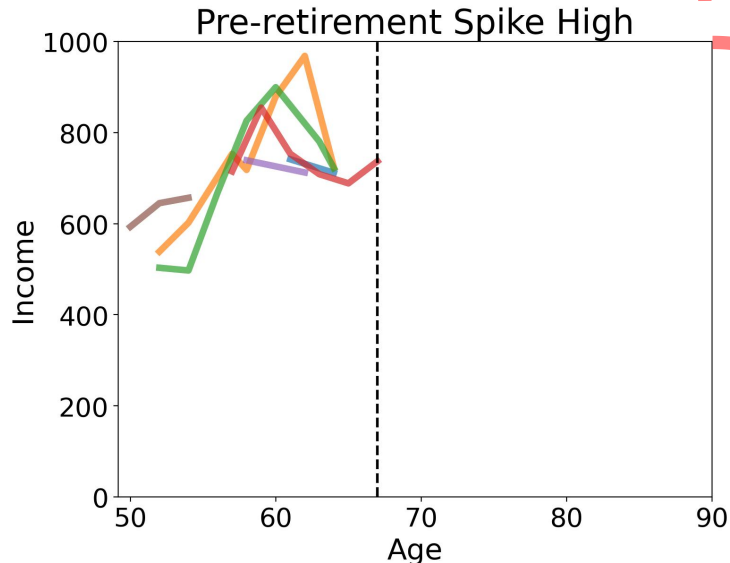
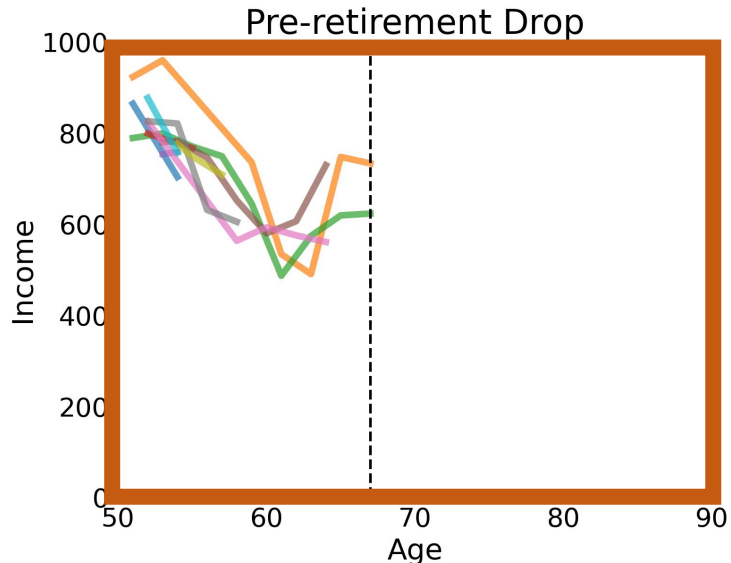
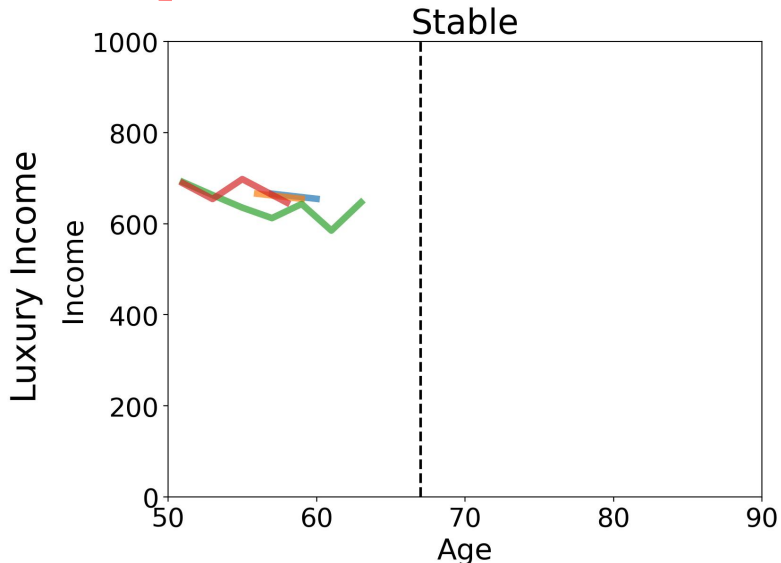
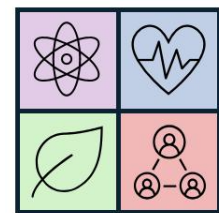
5 variation clusters

4 volatility clusters

Low education ↓ luxury

Pensions ↑ luxury

Facet 2: Income variation

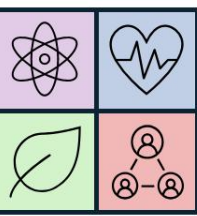


Facet 1:
Income
value

retirement age at 67

Summary

Real-world data are multi-faceted



1. **Interpretable** by disentangling multiple facets
2. Current implementation supports
 - a) Large datasets
 - b) Automatic number of clusters
3. Broader **applicability** to various trajectories
 - a) Socioeconomics
 - b) Health
 - c) and many more!

Thank you and get in touch!

