

Unbiased Disease Stratification within the IMI U-BIOPRED Severe Asthma Programme Using **Topological Data Mining**

Dominic Burg^{1,2}, Doroteya Staykova¹, <u>Pek Lum³</u>, Xian Yang⁴, Yike Guo⁴, Anthony Rowe⁵, Ratko Djukanović², Paul Skipp¹ and the U-BIOPRED consortium 1. Centre for Proteomic Research, University of Southampton, UK. 2. NIHR Southampton University Hospital, UK. 3. Ayasdi Inc, USA. 4. Imperial College, London, UK. **5.** Janssen R&D, UK

(p = 0.0018)

Background and approach The U-BIOPRED consortium is an EU-wide collective of academics (20 institutions), biopharma industry (12), SME's (3), and patient organisations (6) working collaboratively to improve understanding of severe asthma. Representing the largest study cohort recruited for this disease, a variety of sample types are currently being analysed in parallel using a range of 'omics technologies to map molecular and clinical phenotypes of severe asthma in an unbiased manner. The heterogeneity of the disease, combined with the complexity of the study cohort (e.g. participants on a combination of medications, and varying co-morbidities), the range of biofluids and tissues analysed (each with corresponding challenges), and the multiple analytes being measured (e.g. lipids, proteins, mRNA) have necessitated the development of multiple data analysis pipelines to mine these complex datasets. One of the approaches used by the consortium is Topological Data Analysis (TDA), implemented via the Ayasdi software platform. TDA generates topological networks (Figures 1&2) that allow the scientist to explore, condense, visualise and extract useful information from these complex and multi-modal data. The subsequent sections show examples of how the Ayasdi platform is being used to combine and analyse proteomic data produced in the U-BIOPRED study, particularly highlighting the utility of the approach in exploring the biology of the data and as an unbiased feature selection tool. The datasets used to construct the example analyses represent only a small fraction of the final UBIOPRED cohort. As such, all results and interpretations must be viewed as exploratory and illustrative of approach and may not be representative of any final outcomes of the study. A) Original A) Original point cloud point cloud B) Gain 2.0 resolution 12 B) Colouring by filter value

Figure 1(Left) Summary of the Ayasdi approach¹ A) A 3D object (hand) represented as a point cloud. This point cloud can also be generated through e.g. correlation analysis of a complex dataset **B**) A filter value is applied to the point cloud and the object is now coloured by the values of the filter function. **C)** The data set is binned into overlapping groups. **D)** Each bin is clustered and a network is built. **Figure 2** (right) Resolution and gain of graphs from : A) analysis of a point cloud . (B&C) Resolution alters the number of data points in each node (D&E) Gain alters the bin overlap and the subsequent number of edges between nodes

C) Gain 2.0 resolution 30

D) Gain 3.0 resolution 30

E) Gain 4.0 resolution 30

....



C) Binning by filter value

D) Clustering and network construction



Innovative Medicines Initiative

Exploring features of complex data from multiple sources

Proteomic data, obtained from the analysis of sputum supernatant and top-12 depleted serum using the data independent acquisition approach MS^E, were searched against the UNIPROT database for identity assignment and quantification using the Hi-3 method². Each dataset was normalised and batch effect correction performed where necessary using modified ComBat³ scripts. Proteomic data were subsequently aligned and clinical information were added as metadata, resulting in a complete dataset for 80 asthmatic participants. Missing values in all data types were not imputed and left as 'null'. Data were loaded into the Ayasdi platform, including proteins with up to 60% null values, and analysed with a normalised correlation metric and multidimensional scaling (MDS) lenses. The resultant graph was explored using the proteomic and clinical data, and groups selected using the clinical meta-data (Figures 3 & 4).



Figure 3. Exploration of the combined graph with clinical and proteomic features. Panel A shows distribution of assigned cohorts . **Panels B-E** show various aspects of characterised by eosinophilic inflammation and atopy (sensitisation to common aero-allergens like house dust mites) with subtle differences in airway vs. systemic indicators. **Panels F-H** contrast eosinophilia from previous panels with neutrophilia. Interesting differences in sputum neutrophil counts vs systemic neutrophil counts can also be observed. Panels I & J contrast participant use of short acting beta agonists and oral corticosteroids (converted from categorical data). There appears to be correlation between short acting beta agonist use and a subset of participants with sputum neutrophilia



LICH

Aerocrine





Figure 4. Selection and analysis of subgroups using standard statistical tests (Non parametric hypergeometric) revealed features of specific groups and their participant membership. This information is being used as a basis for hypothesis generation

Feature and class selection in cryptic datasets

Traditional statistical analyses of data from high throughput proteomic analysis of serum from asthmatics, produces datasets with very few differentiating features (Table 1 and Figure 6). Using the Ayasdi software platform we were able to identify groups of participants based on these serum data, with specific group assignment guided by persistence of structure and contrasting clinical metadata (Figure 7). These topological groups were subsequently used as class labels in supervised machine learning approaches, implemented in InforSense⁴ (e.g. support vector machine, Figures 8 & 9), and resulted in an improved classification performance and predictive models over cohort information alone.



	Features: p < 0.05		Average
	KS test	FDR correction	Fold Change
ers	3	1	1.18
•	15	2	1.31
ate	9	1	1.26

