

A UNIVERSAL DATA MODEL (UDM) FOR LONGITUDINAL INTEGRATION OF DISPARATE BIOMARKER AND IN-LIFE PATIENT DATA AUGMENTED BY MACHINE LEARNING

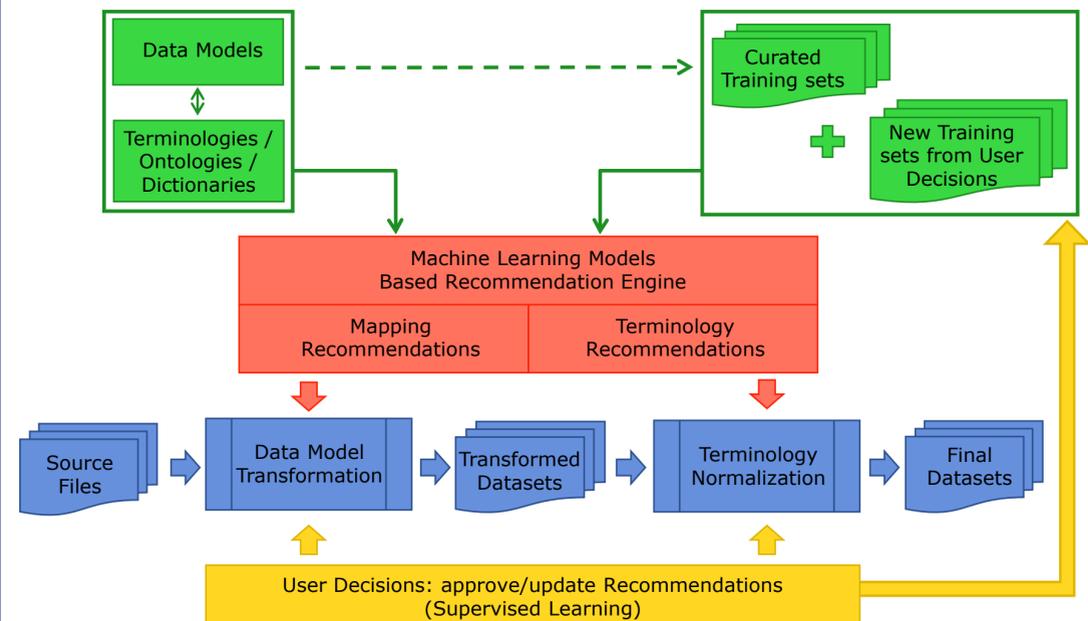
Introduction

Translational and precision medicine development in immuno and gene therapies rely on biomarker data from assays including genomics, proteomics, IHC, Flow cytometry and cell phenotyping data. Biomarkers are not only generated from the patient biosamples, but also from the biomanufacturing sites such as in adoptive immune cell biomanufacturing for novel immunotherapies. Extracting valuable insights from these disparate data sources and integrating it to clinical data to relate to patient outcome and/or discover and validate relevant biomarkers are met with challenges of ingestion, harmonization and integration of disparate data with the clinical data. Key decisions and ideas that impact study design of future clinical trials depend on gaining insights rapidly from such integrated data on patients or stratified cohorts. Systematic curation with self-validation for completeness and consistency is time consuming and difficult without technology. Xbiom, built on Machine Learning and Universal Data Modelling effectively solves these challenges of disparate, big and varied data sources. Xbiom's Smart Transformation platform, can make ingestion, curation and harmonization process automatic and is also capable of processing both streamed data as well as in batch mode.

Data Transformation Workflow

Xbiom's Machine Learning algorithms work on the principle of supervised learning. First, for the given data set Recommendation Engine recommends the Target model along with referencing to external registries for Data transformation. It is then followed by Mapping recommendations by machine learning models based on the trained datasets. The Xbiom platform is agile by having the ability to retrain the algorithm for the newer/unseen type of data types and structures but also robust in predicting the column for a recommendation in manual approval in case of unseen data types. Machine Recommendations can be further approved by the users which will be considered as learning in the next iteration.

Mapped data is further enriched by the Terminology Recommendation, wherein, for instance for protein names in Biomarker data will be enriched with species-specific Harmonization of proteins with reference to the Global registry UniProt. Similarly, Harmonization of Genes, cell-expressed proteins, cytokines etc., will be done with respect to the referenced registries. The curated and transformed output is quality checked and which can be viewed in the QC dashboard for the quality process.



Workflow of Data Transformation

Smart Transformation and Harmonization

Raw Input File

LLOQ	D	M	Y	STUDY/ID	STAT	VENDOR	SUBJ/ID	TPT	TESTCD	ORRES	ORRESU	METHOD
1.1	23	4	2019	PC123		CROSSTICS	1001	C1D1	IL_6	6.6	pg/mL	multiplex immunoassay
1.1	23	4	2019	PC123		CROSSTICS	1001	C1D1	VEGF_A	68	pg/mL	multiplex immunoassay
1.1	23	4	2019	PC123		CROSSTICS	1001	C1D1	VEGF_C	1432	pg/mL	multiplex immunoassay
1.1	23	4	2019	PC123		CROSSTICS	1001	C1D1	VEGF_D	BLQ	pg/mL	multiplex immunoassay
1.1	23	4	2019	PC123		CROSSTICS	1001	C1D1	TMEM16B	12.5	pg/mL	multiplex immunoassay
1.1	23	4	2019	PC123		CROSSTICS	1001	C1D1	DPPI	23	pg/mL	multiplex immunoassay
1.1	23	4	2019	PC123		CROSSTICS	1001	C1D1	CYP11A6	15.7	pg/mL	multiplex immunoassay
1.1	23	4	2019	PC123		CROSSTICS	1001	C1D1	Connexin2	1234	pg/mL	multiplex immunoassay
1.1	23	4	2019	PC123		CROSSTICS	1001	C1D1	FMR2P	112.34	pg/mL	multiplex immunoassay
1.1	23	4	2019	PC123		CROSSTICS	1001	C1D1	AFAP-110	11.5	pg/mL	multiplex immunoassay

Target Model for Transformation

Batch Details

Summary

Batch Identifier: PC123

Data Source: Clinical Study Data Repository

Study List: PC123

Target Model Name: Protein Immunoassay

CT Version: RDIS Terminologies

Output Format: CSV

Data Transformations

Metadata Mappings

#	Source	Mapping	Target	Recommendation	No.
1	STUDYID		STUDYID	✓	
2	SUBJ_ID		SUBJ_ID	✓	
3	CONCAT_WS(';', STUDYID, SUBJ_ID)		USUBJID	✓	
4	TPT		XBPT	✓	
5	TESTCD		XBPTOT	✓	
6	ORRES		XBORRES	✓	
7	ORRESU		XBORRESU	✓	
8	LLOQ		XBLOQ	✓	
9	STAT		XBSTAT	✓	
10	REASND		XBREASND	✓	
11	D		XBREASND	✓	
12	M		XBREASND	✓	

Data Harmonization

IA.csv

USUBJID	XBSPC	XBPROT	XBORRES	XBORRESU	XBLOQ	XBCTC	XBPT	XBPT	XBPT
PC123-1001	blood plasma	Interleukin-6	6.6	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Vascular endothelial growth factor A	68	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Vascular endothelial growth factor C	1432	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Vascular endothelial growth factor D	BLQ	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	TMEM16B	12.5	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Dependent glycoprotein 1	23	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Cytochrome P450 2A6	15.7	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	BEN domain-containing protein 1B	1234	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	AFAP110 family member 2	112.34	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Actin filament-associated protein 1	11.5	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Chondroitin sulfate glucuronyltransferase	BLQ	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	SP9001 protein Chof116	BLQ	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Ceramidase kinase	113	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	CREB-regulated transcription coactivator 1	345	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Corticotropin-releasing factor receptor 1	239	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Beta-defensin 103	456	pg/mL	1.1	2019-04-23	C1D1	multipl	
PC123-1001	blood plasma	Gap junction alpha-3 protein	13	pg/mL	1.1	2019-04-23	C1D1	multipl	

Standardized Output File

Disparate Data - Sample formats for Immunoassay data received from Vendors

STUDYID	SUBJID	SPID	TPT	TESTCD	ORRES	ORRESU	LLOQ	STAT	REASND	D	M	Y
PC123	1001	501	C1D1	IL_6	6.6	pg/mL	1.1			23	4	2019
PC123	1001	501	C1D1	VEGF_A	68	pg/mL	1.1			23	4	2019
PC123	1001	501	C1D1	VEGF_C	1432	pg/mL	1.1			23	4	2019
PC123	1001	501	C1D1	VEGF_D	BLQ	pg/mL	1.1			23	4	2019

Machine Readable Columnar Format

Human Readable Pivoted Format

Source files transformation to a selected Target model along with reference to external registries for Data harmonization

Variable Label	Value Coded By	Controlled Terms	Format	Core	Data Type	Notes
Study Identifier				Required	Text	Study Identifier
Domain				Required	Text	Domain Name
Subject Identifier for the Study				Required	Text	Subject Identifier
Unique Subject Identifier				Required	Text	Identifier used to uniquely identify a subject ac...
Secret Flag				Integer	Text	To mark Secret flag for Subjects (0-Restricted...
Specimen Type	Controlled Terms	SPECIYPE			Text	Sample Type or Matrix or Tissue type
Ontology ID for Specimen Type					Text	Ontology ID for Sample type
Category of Specimen Type					Text	Sample Category (Ex: BODY LIQUIDS)
Source Visit Name					Text	Visit Name
Visit Name					Text	Visit Name. Needs to map to visit name from ...
Visit Number					Float	Visit Number. Needs to map to visit name from ...
Protein Name	External Dictionary	UniProt		Required	Text	Source Protein Name
Source Protein ID					Text	Protein ID
Gene Symbol	External Dictionary	NCBI Gene Info			Text	Gene Symbol
Source Gene ID					Text	Source Gene Identifier NOT identical to NCBI ID
Result or Finding in Original Units					Text	Results or Findings
Original Units					Text	Units of the Result collected

Mapping of Source data to Target model variables, Mapping Recommendations by machine learning models based on the Training datasets & Users previous decision. Machine Recommendations can be further approved by the user which will be considered as Learning in the next iterations

Species specific Harmonization of Proteins with the Global registry UniProt. Recommendation of the Protein names based on look-up of source protein names in the registries.

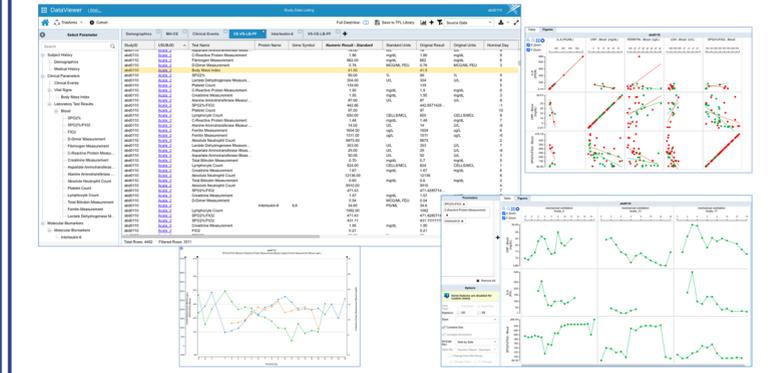
UniProt

Source Term	Protein Name	Protein ID	Recommen...
AFAP-110	Actin filament-associated protein 1	Q9N556	
CRFR-1	Corticotropin-releasing factor receptor 1	P34998	
CYP2A6	Cytochrome P450 2A6	P11509	
Ch3y-3	Chondroitin sulfate glucuronyltransferase	Q9P263	
Connexin2	BEN domain-containing protein 1B	Q6P351	
Cx46	Gap junction alpha-3 protein	Q9Y948	
DPPI	Dependent glycoprotein 1	P53534	
FMR2P	AFAP110 family member 2	P51616	
HD03	Beta-defensin 103	P61534	
IL6	Interleukin-6	P05231	
IL6A	Ceramidase kinase	Q8TCT0	
Pncp1	UPP0991 protein Chof116	Q5B946	
TMEM16B	TMEM16B		
TORC-1	CREB-regulated transcription coactivator 1	Q6U9V9	

Longitudinal Integration of Biomarker Data & patient data facilitated by UDM for visualization and analysis

Curated and harmonized data is stored in a Universal Data Model (UDM) which holds study and assay data of subjects or samples and their attributes in a simple, indexable form. The indexation facilitates instant search through query masks to yield highly stratified cohorts and data. Search and detection of signals is enabled by the integration of assay data into longitudinal patient data.

As shown in the immersive graphics below on the left sidebar listing the available data for the searched cohort enables handpicking of Biomarker which can be correlated across the cohort with another biomarker(IHC, TMB, FACS, Nanostring etc.) or with the clinical data (adverse event, clinical endpoints, Medical history etc.). UDM can be unique to the organization and is extendable for the future needs.



Conclusion

Challenges in translational medicine and other research in drug development involving different biomarkers are the difficulty in getting the cleaned biomarker data which is integrated with the clinical data to ask questions to gain insights into the study. Manual curation and harmonization can be time-consuming and would be challenging with steamed data in the ongoing trials.

Automated curation, standardization and harmonization by the Machine Learning platform in Xbiom is quality controlled and can process the streamed data. Data in the UDM is ready for the statistical analysis and enables the monitoring of ongoing trials for early insights and decision making. Harmonized data in the UDM can also act as a training data set in machine learning designed for biomarker discovery and other insights generation as in precision and translational medicine. The data governance system in Xbiom's allowing for the role-based access can be accessed safely in the browser as well as through the API. Quick and timely decision enabling process can cut down the time and cost in drug development.