Rancho bio ciences



Aggregation and Integration of UK Biobank Phenotype Data for Downstream Analysis

Yulia Skovpen¹, Emily Wong², Christine Loreth¹, Alena Fedarovich¹, Tyler Kolisnik¹, Oksana Tyurina¹, Tania Khasanova¹, David Merberg³, Sándor Szalma², Julie Bryant¹

¹Rancho BioSciences, LLC, 16955 Via Del Campo #220, San Diego, CA 92127

^{2, 3}Takeda Pharmaceuticals, 9625 Towne Centre Drive, San Diego, CA 92101 and 35 Lansdowne Street, Cambridge, MA, 02139

Background

The UK Biobank is a rich dataset from 500,000 participants and contains phenotypic data types such as participant questionnaires and HER-derived data, genomic, biomarkers, and imaging data. It is a live dataset that is regularly updated in real-time throughout the year. As a result, data that is not readily amenable for analysis, or consumption accumulates in size. A curation step is required to harmonize incoming data to make it available for downstream data analysis.





- 502,600 subjects
- 3,390 variable fields
- Longitudinal data
 16 000 data fields
 - 16,000 data fields



Phenotype Harmonization

Around 100 harmonization rules have been proposed:

- Conversion of years (of events) into subject age, replacement of "ongoing data" by actual values of age or year
- replacement of categorical values in "Continuous" or "Integer" data fields by relevant numeric values
- harmonization of discrepancies in subjects' answers at multiple visits when an answer implies a single answer
- -define and harmonize data where subject's answer at the last visit

Mapping of Ontologies

Diagnoses and diseases – ICD10
 Operative procedures – SNOMED CT
 Drugs and medications – RXNORM, MeSH
 – Indication-based drugs classification - MeSH

Category	Count	
🗄 🧰 ear/nose/throat cancer	662	
🗄 🚞 gastrointestinal cancer	3913	
🗄 🚞 neurological system cancer	496	
🗄 🔄 urinary tract cancer	_	
kidney/renal cell cancer	713	
bladder cancer	1300	
other cancer of urinary tract	37	
breast cancer	12626	
🗄 🚞 genital tract cancer	11849	
🗄 🚞 haematological malignancy	2822	
🗄 🚞 skin cancer	13512	
🗄 🧰 other cancer	1223	
🗄 🚞 respiratory / intrathoracic cancer	898	
unclassifiable	_	

ICD10_meaning_coding19	node_id	parent_id	selectable
C03.9 Gum, unspecified	1229	1226	Y
C07 Malignant neoplasm of parotid gland	1247	44	Y
C08.9 Major salivary gland, unspecified	1252	1248	Y
C15.9 Oesophagus, unspecified	1292	1284	Υ
C16.9 Stomach, unspecified	1302	1293	Y
C17.9 Small intestine, unspecified	1309	1303	Υ
C18.9 Colon, unspecified	1320	1310	Υ
C21.0 Anus, unspecified	1324	1323	Υ
C18.7 Sigmoid colon	1318	1310	Y
C20 Malignant neoplasm of rectum	1322	45	Y
C22.9 Liver, unspecified	1335	1328	Y
C23 Malignant neoplasm of gallbladder	1336	45	Y
C25.9 Pancreas, unspecified	1350	1342	Υ
C34.9 Bronchus or lung, unspecified	1380	1374	Υ
C34.9 Bronchus or lung, unspecified	1380	1374	Y
C47.9 Peripheral nerves and autonomic n	en 1459	1450	Υ
C69.8 Overlapping lesion of eve and adne	exa 1561	1553	Y

- (instance) has most recent information
- identify data fields and subjects with values outside coding files, make appropriate curation decisions

fid	Field name	Data type	Code file	Value range	Code meaning/ Problem	Solution
20420	Longest period spent worried or anxious	Integer	517	-1, 1, 2, 3, , 839, - 999	-999 All my life/as long as I can remember	Replace -999 by Age of subject in months

• Solution example: value -999 "All my life/as long as I can remember" to convert into Age of subject at the time of answering the questionnaire

Enrichment of Cohorts



Example of Results

Treatment/medication code mapping to ontologies

RXNORM name	RXNORM ID	MESH name	MESH ID
vitamin e	RXNORM:11256	vitamin e	MESHF:D014810
chondroitin sulfates	RXNORM:2473	chondroitin	MESHF:D002807
coenzyme q10	RXNORM:21406	coenzyme q10	MESHF:C024989
unknown/not mapped	NA	unknown/not mapped	NA
fish oils	RXNORM:4419	fish oils	MESHF:D005395
ascorbic acid	RXNORM:1151	ascorbic acid	MESHF:D001205
evening primrose oil	RXNORM:203219	efamol	MESHF:C028498
bioflavonoids	RXNORM:1562	dietary supplements	MESHF:D019587
hypericum extract	RXNORM:1309228	hypericum	MESHF:D020902
aloe extract	RXNORM:91263	aloe	MESHF:D000504
serenoa preparation	RXNORM:350488	serenoa	MESHF:D028024
bioflavonoids	RXNORM:1562	dietary supplements	MESHF:D019587
unknown/not mapped	NA	unknown/not mapped	NA
aluminum hydroxide;magnesium hydrox	RXNORM:612;RXNORM:6581	aluminum hydroxide;mag	MESHF:D000536;MESHF:D008
aluminum hydroxide;magnesium hydrox	RXNORM:612;RXNORM:6581;RXNO	aluminum hydroxide;mag	MESHF:D000536;MESHF:D008
aluminum hydroxide;magnesium hydrox	RXNORM:612;RXNORM:6581	aluminum hydroxide;mag	MESHF:D000536;MESHF:D008
	RXNORM name ✓ vitamin e ✓ chondroitin sulfates ✓ coenzyme q10 ✓ unknown/not mapped ✓ fish oils ✓ ascorbic acid ✓ evening primrose oil ✓ bioflavonoids ✓ hypericum extract ✓ aloe extract ✓ serenoa preparation ✓ bioflavonoids ✓ unknown/not mapped ✓ aluminum hydroxide;magnesium hydrox ✓ aluminum hydroxide;magnesium hydrox ✓ aluminum hydroxide;magnesium hydrox ✓	RXNORM nameRXNORM IDvitamin eRXNORM:11256chondroitin sulfatesRXNORM:2473coenzyme q10RXNORM:21406unknown/not mappedNAfish oilsRXNORM:4419ascorbic acidRXNORM:1151evening primrose oilRXNORM:203219bioflavonoidsRXNORM:1309228aloe extractRXNORM:91263serenoa preparationRXNORM:350488bioflavonoidsRXNORM:1562unknown/not mappedNAaluminum hydroxide;magnesium hydroxRXNORM:612;RXNORM:6581;RXNOaluminum hydroxide;magnesium hydroxRXNORM:612;RXNORM:6581	RXNORM nameRXNORM IDMESH namevitamin eRXNORM:11256vitamin echondroitin sulfatesRXNORM:2473chondroitincoenzyme q10RXNORM:21406coenzyme q10unknown/not mappedNAunknown/not mappedfish oilsRXNORM:419fish oilsascorbic acidRXNORM:1151ascorbic acidevening primrose oilRXNORM:203219efamolbioflavonoidsRXNORM:1562dietary supplementshypericum extractRXNORM:1309228hypericumaloe extractRXNORM:1562dietary supplementsbioflavonoidsRXNORM:1562dietary supplementsunknown/not mappedNAunknown/not mappedaloe extractRXNORM:1562dietary supplementsunknown/not mappedNAunknown/not mappedaluminum hydroxide;magnesium hydroxRXNORM:612;RXNORM:6581aluminum hydroxide;magnesium hydroxaluminum hydroxide;magnesium hydroxRXNORM:612;RXNORM:6581aluminum hydroxide;magnesium hydrox

in size will add power to statistics for downstream analysis.

Cohorts enrichment

1	f.eid>disease_001>disease_002>disease_003>disease_004>disease_005>disease_006>disease_007>di	i
2	$\texttt{5727058} \times \texttt{M518} \longrightarrow \texttt{M9596} \longrightarrow \longrightarrow$	_
3	$\texttt{2339680} \rightarrow \texttt{M758} \longrightarrow \texttt{M430} \longrightarrow \texttt{I10} \rightarrow \texttt{E785} \longrightarrow \longrightarrow$	_
4	$\texttt{3186124} \times \texttt{N40} \times \texttt{I10} \times \texttt{K30} \times \texttt{M179} \longrightarrow \texttt{K409} \longrightarrow \texttt{E785} \longrightarrow \longrightarrow$	_
5	$4748085 \rightarrow D061 \longrightarrow D069 \longrightarrow R91 \rightarrow J304 \longrightarrow J189 \longrightarrow \longrightarrow$	_
6	$5046970 \times C504 \longrightarrow Z853 \longrightarrow V134 \longrightarrow R55 \times R11 \times I10 \times S809 \longrightarrow C509 \longrightarrow \longrightarrow$	_
7	$3884101 \rightarrow Z602 \longrightarrow J459 \longrightarrow Z301 \longrightarrow R634 \longrightarrow R55 \rightarrow R002 \longrightarrow N920 \longrightarrow M5446 \longrightarrow M512 \longrightarrow A099 \longrightarrow M518 \longrightarrow CONTRACTOR (CONTRACTOR (C$	_
8	$\texttt{5429609} \times \texttt{Z538} \longrightarrow \texttt{R05} \times \texttt{R042} \longrightarrow \texttt{N394} \longrightarrow \texttt{R490} \longrightarrow \texttt{R32} \times \texttt{R13} \times \texttt{K625} \longrightarrow \texttt{F458} \longrightarrow \texttt{D259} \longrightarrow \longrightarrow \longrightarrow \longrightarrow \longrightarrow \texttt{R490} \times \texttt{R490} \longrightarrow R490$	_
9	$\texttt{3070554} \texttt{R31} \texttt{M8999} \longrightarrow \longrightarrow$	-
10	$4403228 \times N950 \longrightarrow N952 \longrightarrow \longrightarrow$	-
11	$3678822 \rightarrow \longrightarrow \longrightarrow$	-
12	$5374646 \rightarrow J459 \longrightarrow G439 \longrightarrow \longrightarrow$	_
13	$1127007 \rightarrow \texttt{Z864} \longrightarrow \texttt{K660} \longrightarrow \texttt{K562} \longrightarrow \texttt{K529} \longrightarrow \texttt{K219} \longrightarrow \texttt{F412} \longrightarrow \texttt{Z871} \longrightarrow \texttt{R13} \rightarrow \texttt{Z922} \longrightarrow \texttt{Z720} \longrightarrow \texttt{I10} \rightarrow \texttt{Z824} \longrightarrow \texttt{L10} \rightarrow \texttt{L10} \rightarrow$	_