

# UK Biobank Showcase User Guide: Getting Started

## 1 Introduction

UK Biobank holds an unprecedented amount of data on half a million participants aged 40-69 years (with a roughly even number of men and women) recruited between 2006 and 2010 throughout the UK. Showcase (available at <http://www.ukbiobank.ac.uk>) aims to present the data available for health-related research in a comprehensive and concise way, and to provide technical information for researchers considering applying to use the resource.

This user guide is designed to give you an overview of the data and provides some instructions on how to navigate your way through the system.

### Suggestions and information for new users:

- Have a printout of this user guide handy when you first use Showcase
- Read the background information about UK Biobank and details on access procedures at <http://www.ukbiobank.ac.uk/>
- Take time to familiarise yourself with the Showcase structure, the accompanying documentation and the descriptions provided for each data-field before completing a preliminary application to use the resource.
- While it is worth checking what variables are available in the Showcase Resource, there is no requirement to construct a Showcase basket until the Main Application stage.
- Note that Showcase is continually under development, as new data on exposure and health outcomes is incorporated into the database. More information on the timelines for future data is availability in the [Essential information](#) section of data Showcase.

If you encounter problems or faults, please email [showcase@ukbiobank.ac.uk](mailto:showcase@ukbiobank.ac.uk)



## 2 Data included in UK Biobank

### 2.1 Data collected at the Assessment Centre

All participants in UK Biobank were recruited through assessment centres, designed specifically for this purpose (a map of the 22 assessment centres is provided in the Essential Information section of the Showcase). Data collected at the assessment visit included information on a participant's health and lifestyle, hearing and cognitive function, collected through a touchscreen questionnaire and brief verbal interview. A range of physical measurements were also performed, which included: blood pressure; arterial stiffness; eye measures (visual acuity, refractometry, intraocular pressure, optical coherence tomography); body composition measures (including impedance); hand-grip strength; ultrasound bone densitometry; spirometry; and an exercise/fitness test with ECG. Samples of blood, urine and saliva were also collected. Some of these measures were incorporated into the Assessment visit towards the end of the recruitment period and are therefore not available for all 500,000 participants (see the 'Essential Information' section of the Showcase for more information on timelines).

During 2006, over 3,000 participants were included in the pilot phase of recruitment. Where possible, data collected from the pilot and the main recruitment phases have been combined. Where modifications to the protocol were made after the pilot study, the data-fields from the pilot and main recruitment phase are listed separately (e.g., touchscreen questions on medications, family history, qualifications and household income). Pilot data-fields can be identified easily; they include 'pilot' in the data-field name and have a field ID number in the 10000s). In addition, cognitive function tests that were felt to be too time-consuming and/or relatively uninformative were omitted from the main phase of recruitment (i.e. the light memory test' on the touchscreen questionnaire and the 'word test' that was performed during the verbal interview stage).

A web-based dietary questionnaire was included in the Assessment visit towards the end of the recruitment period, and participants were also invited via e-mail to complete the questionnaire on four further occasions (between 2011-2012).

#### **\*\*NEW\*\***

**REPEAT ASSESSMENT DATA:** A repeat assessment of 20,000 participants was carried out between August 2012 and June 2013 at the UK Biobank Co-ordinating Centre, Stockport, UK. Participants who lived within a 35 km radius of the assessment centre were invited to attend an appointment and undergo a repeat assessment of all the baseline measures. These data are now available in data Showcase. For further details; please consult the guide to [Repeat Assessment Data](#) (available in the Essential Information section of the Showcase).

## 2.2 Future data availability

Please see the 'Essential Information' section of the Showcase for an outline of which measures will be incorporated into the Resource over the next 12-18 months.

## 3 Finding data in Showcase

**\*\*NEW\*\***

**RECOMMENDED CATEGORIES:** To make it quicker and easier for you to build your dataset, we have created groups of commonly requested data-fields (i.e. variables) to act as a useful starting point so that you don't have to select individual data-fields from each category. These have been loosely grouped into categories such as demographics, cognitive function (that contain the main summary data-fields for all tests), physical measures (that contain the main summary data-fields from all measures), etc. You can then add or remove individual data-fields to or from your pre-populated basket to create a bespoke dataset for your research project. You can find these in the **CATEGORY** listing in the **CATALOGUE** section (<http://biobank.ctsu.ox.ac.uk/crystal/cats.cgi>); alternatively, you can link to them from the Essential Information section.

You can find the rest of the data that you need through two main routes:

**BROWSE:** Use this to navigate your way through hierarchical categories and subcategories of interest to data-fields (i.e. variables) of interest. **This will be the most appropriate tool for most researchers wishing to find and select data for their application to use the Resource.**

**SEARCH:** This is based on a text search of the data-field name and its notes, and uses the Boolean operators '&' and '|' to denote the 'AND' and 'OR' functions. By default, only whole word matches are returned, although you can use the asterisk '\*' as a wild-card character at the beginning or end of text to search for words containing that text. Please see the **HELP** page on 'Searching text' for more details. The **Full Search** facility allows you to conduct a search using specific criteria based on the type of data-field (see Section 5 for more details).

A full list of data-fields, categories and documents can be found in **CATALOGUES**.

## 4 Data categories and sub-categories

Data are organised in a tree structure, accessible via **BROWSE**, with the main categories based on the origin of data collection (Figure 1). These include:

- Base characteristics (some general characteristics of participants known before arrival)
- UK Biobank Assessment Centre (data obtained at the Assessment Centre)
- Laboratory biological samples (data on biological samples)
- Additional exposure data (data collected outside the Assessment Centre)
- Health-related outcomes (data from linkage of participants to health-related records).

Please see the **HELP** page on 'Browse' for more details.

The **Fields** column lists the number of data-fields in each category (and its sub-categories)

The **Help** button provides more information about items, as listed in the **Glossary** (at the bottom of the Help page)

Clicking on the **'Show Level'** button is an easy way to jump to a more detailed level

Category ID	Description	Fields
	Base characteristics	7
	UK Biobank Assessment Centre	850
	Laboratory biological samples	105
	Additional exposure data	0
	Health-related outcomes	0

Figure 1. Illustration of the tree structure via **BROWSE**

Clicking on the **Category ID** or the **Description** leads you to the subcategories and/ or data-fields contained within that category.

Category ID	Description	Fields
100021	Recruitment	8
100004	General assessment centre data	10
100025	Touchscreen	471
100071	Verbal interview	31
100006	Physical measures	329
100001	Biological sampling	10

Figure 2. Illustration of sub-categories within the 'Touchscreen' category

The tree structure assigns data-fields to one location only, and is not currently cross-referenced. It is therefore important to look in all parts of the tree that might contain data-fields relevant to your research question(s). In general, you should not rely on the **SEARCH** facility to find all fields of relevance for a particular topic.

## 5 Data-field information

The panel in the top-half of the data-field screen provides a brief description and category location of the data-field within the tree structure (Figure 3). It also includes more detailed technical information about each data-field. This includes information on:

- **Participants:** the number of participants that have the data item
- **Item Count:** the number of data items available
- **Stability:** whether the data-field is complete or changes over time

- **Value type:** the format and units of the data-field
- **Item type:** whether the data-field is a simple data point, relates to an inventory of biological samples, or is a large data object
- **Strata:** the likely relevance to researchers of the data-field
- **Sexed:** whether the data-field is available for both sexes
- **Instances:** how many occasions participants have this measurement performed
- **Array:** whether there are multiple data items for each instance. For example, Figure 3 shows that data on diastolic blood pressure is presented in an array with 2 values per measure (because the measurement was performed twice). Please see the **HELP** page for more details.

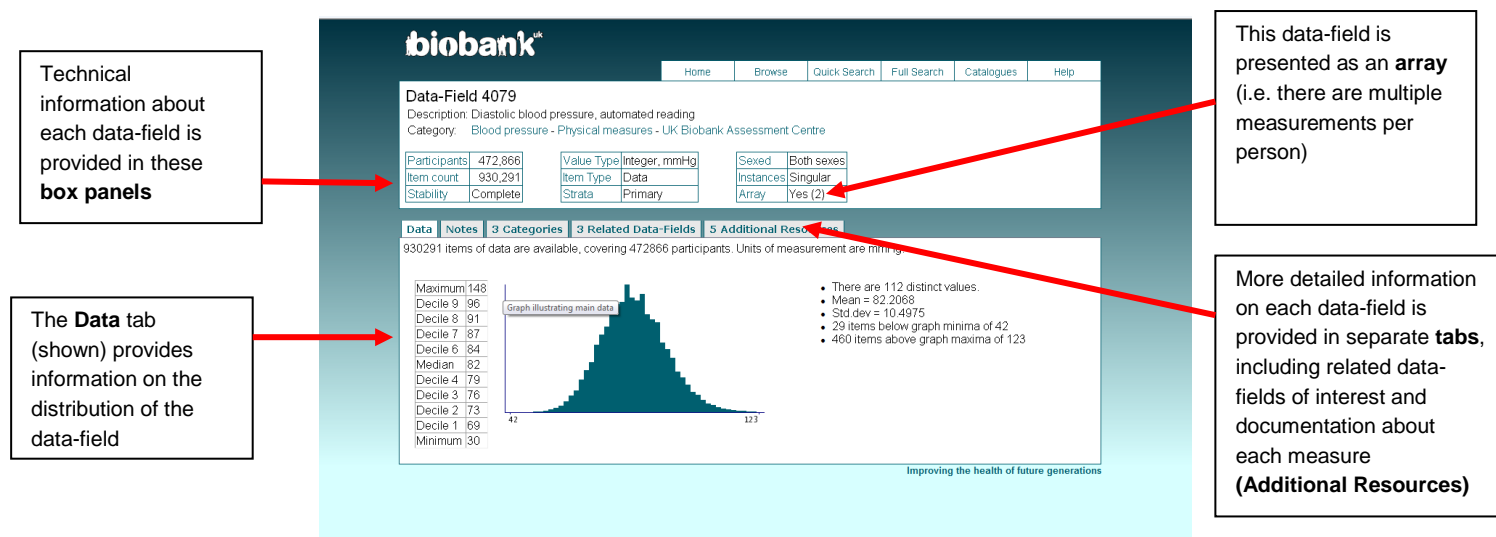


Figure 3. Illustration of a data-field

The univariate distribution of each data-field is presented in graphical or tabular format (or both) in the **Data** tab (Figure 3). Data are not presented if they are free-text, curve data (i.e. spirometry curves) or bulk data items (i.e. too large to be downloaded: eye images and exercise/ECG results). Distributions of data-fields that are of a sensitive nature (e.g., number of sexual partners) are not shown, although approved researchers can still request such data in their application.

The **Notes** tab includes the full description of the data-field, and for touchscreen questions, provides the exact text of the question that was asked, together with other details.

The **Categories** tab lists the categories and sub-categories of which the data-field is a member. This is also shown horizontally in the category tree, at the top of the page.

The **Related Fields** tab lists other data-fields to which the current data-field is related. For example, the data-field for 'diastolic blood pressure, automated reading' (ID: 4079) is related to 3 data-fields: one on diastolic blood pressure from a manual reading, one on systolic blood pressure, and one on pulse taken by the same device (see Figure 3).

The **Additional Resources** tab contains explanatory documentation related to each data-field. This may include screen-shots of the touchscreen questions, details of how each measurement was performed (in downloadable pdf format), photos and video-links.

Some data-fields that are not of primary interest to most researchers may nonetheless be of interest for some research purposes, and these have been classified as supporting or auxillary data-fields (in **Strata**). Examples include the keystroke history of a participant during a touchscreen question, and serial numbers of devices/equipment. Supporting and auxillary data-fields are not searched using **Quick Search**, although you can use the **Full Search** facility to specify the type of data-fields that are included in the search request (see Figure 4). You can also find these data fields using the **BROWSE** function.

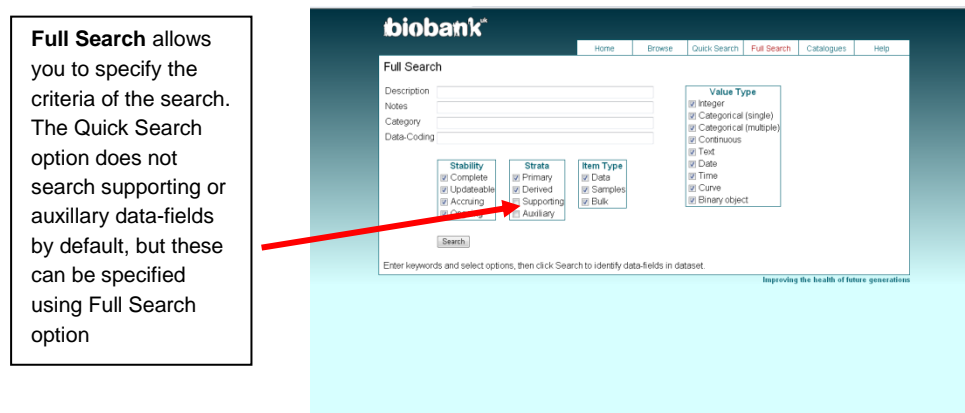
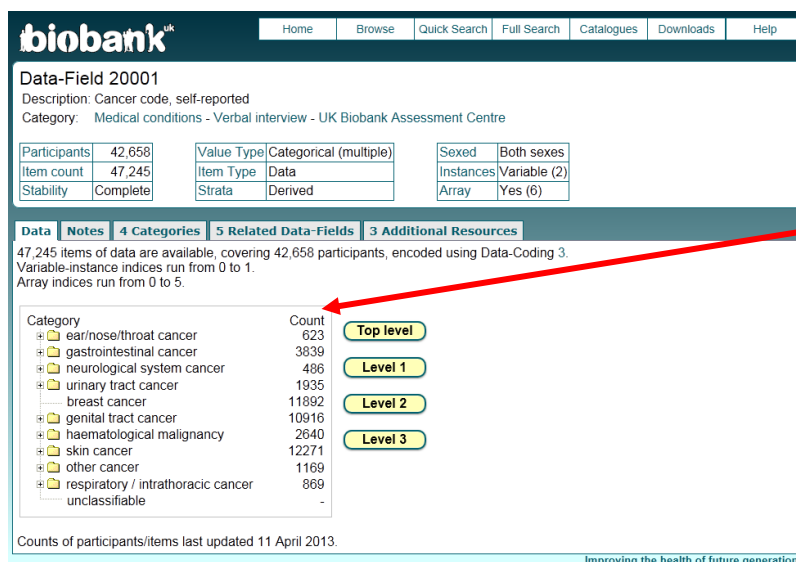


Figure 4. Illustration of the Full Search facility

## 6 Self-reported medical conditions

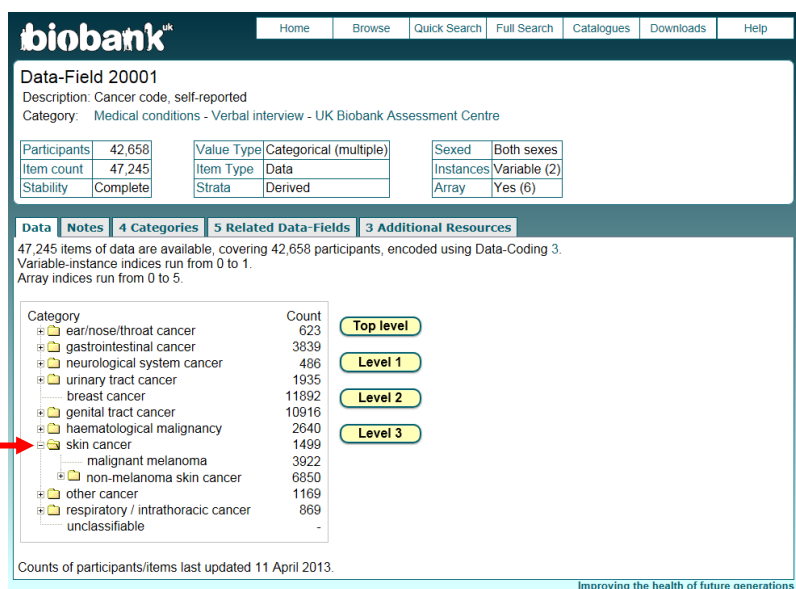
We advise that you use the **BROWSE** function (rather than **SEARCH**) to find data about a medical condition of interest. Self-reported medical conditions at assessment were indicated on the touchscreen questionnaire, and then confirmed through an interview with a trained member of staff (please see the category description of 'Medical conditions' for more details).

In Figure 5, the '**COUNT**' column shows the number of data items listed in the category - there are 12,271 data items for parent category 'skin cancer' for example. Clicking on + box reveals more detailed sub-classifications of the condition – e.g. there are 1499, 3922, and 6850 data items for 'skin cancer', 'malignant melanoma' and 'non-melanoma skin cancer' respectively (Figure 6).



The 'COUNT' column shows the number of data items listed in the category

Figure 5. Illustration of the coding for medical conditions in the verbal interview



The tree expands when you click on the + box, to give you sub-classifications of each condition.

Figure 6. Illustration of the coding for medical conditions in the verbal interview

## 7 Health-related Outcomes

UK Biobank acquires information on participants' health outcomes from a variety of different data sources. Death registrations, cancer registrations and hospital episode data are currently being obtained on a regular basis and are being made available to researchers via Data Showcase. These data can be found in the [Health-related outcomes](#) category (category ID: 100091).

UK Biobank holds both prospective and retrospective data. Health Outcomes Reports are published periodically, which aim to give researchers an indication of the number of

prevalent and incident cases for the most common conditions, by age group, sex and year. For the latest report; please see the additional resources tab of the Health-related outcomes category.

## **8 Data cleaning**

Data from the touchscreen questionnaire have been subject to data checks, as outlined in the explanatory documentation. Data from automatic devices were entered directly into the computer thereby minimising manual entry of data. Nonetheless, there may be occasions where wrong device numbers were entered, the date-time stamp was incorrect, or there were lapses in calibration. The majority of data that was entered as free-text (e.g., reason for skipping various measures) has been subsequently coded. However, some data-fields (such as serial device IDs for various physical measures) remain as free-text data items owing to the large number of devices that were used.