Grid Enabled Occupational Data Environment

GEODE project workshop: 'Handling Occupational Data'

# Practical session: Implementing GEODE

*Tuesday 16[th] January 2007, University of Stirling*

We suggest five activities to try out during the lab session, for which we have prepared some guidance notes and examples:

1) **Entering the GEODE portal and searching for occupational information resources**
2) **Comparator: accessing and using occupational information through conventional internet sources**
3) **Depositing new occupational information with GEODE**
4) **Matching occupational information to regular format survey data files through GEODE**
5) **Matching occupational information to a complex survey data file through GEODE**

**In addition, we would welcome the opportunity to spend some time during the lab with any workshop participants who would like to try linking their own data resources to the GEODE service**

*This handout includes descriptions of implementing each procedure. The exercises will often require accessing further files from other sources, as well as deploying other software packages. It will also be useful to have to hand the following documents:*

- **Instructions for Using the GEODE Portal** (GEODE project technical paper, ed. 0.3).

- **GEODE website and portal instruction pages**:
    o Matching occupational information:

       http://www.geode.stir.ac.uk/matching_occupational_data.html

    o Listing of occupational unit group schemes: http://www.geode.stir.ac.uk/ougs.html

    o Curating occupational data: http://www.geode.stir.ac.uk/geode_m_curation.html

    o Translating file formats: http://www.geode.stir.ac.uk/file_convert_info.html

*Members of the GEODE project will be on hand during the lab session to illustrate the exercises and to help users with their implementations.*

**HEALTH WARNING..!!** *An aim of GEODE is to provide a data service which allows users to search for and download occupational information resources, to deposit resources, and to match their own data with existing occupational information resources. At time writing [Jan 07],GEODE can do all of these functions, but **the GEODE portal is not yet particularly user friendly**.*

*Please be patient in this regard!! We continue to work on the accessibility of the portal and welcome further comments. For the purposes of the lab exercises, we will lead you through the current procedures for implementing the GEODE portal, several of which we envisage will shortly be updated to a more user-friendly format.*

In the exercises below we will use some pre-prepared data files and occupational information resources, as follows:

The data files [A-C] are available from: http://www.geode.stir.ac.uk/workshop/data/

| Micro-data files |
|---|
| **[A]** **bhps_w1_extract**.{sav/dta/dat}<br>Subset of British Household Panel Survey from 1991, 957 cases, 7 variables.<br>Key linking variables:<br>   - **ajbsoc** (SOC-90 for current job)<br>   - **ukempst** (employment status classification for current job) |
| **[B]** **lfs_2002extract**.{sav/dta/dat}<br>Subset of Labour Force Survey from 2002, 1527 cases; 55 variables.<br>Key linking variables:<br>   - **soc2km** (SOC-2000 for current job)<br>   - **ukempst** (employment status classification for current job) |
| **[C]** **ess_2001extract**.{sav/dta/dat}<br>Subset of European Social Survey from 2001, 9 countries, 2142 cases, 14 variables.<br>Key linking variables:<br>   - **iscoco** (ISCO-88 for current job)<br>   - **stdempst** (employment status for current job – standardised codes)<br>   - **empstat** (employment status for current job – dichotomy)<br>   - **supvn** (supervisory status in current job)<br>   - **cntry** (country from which data comes, string format description) |

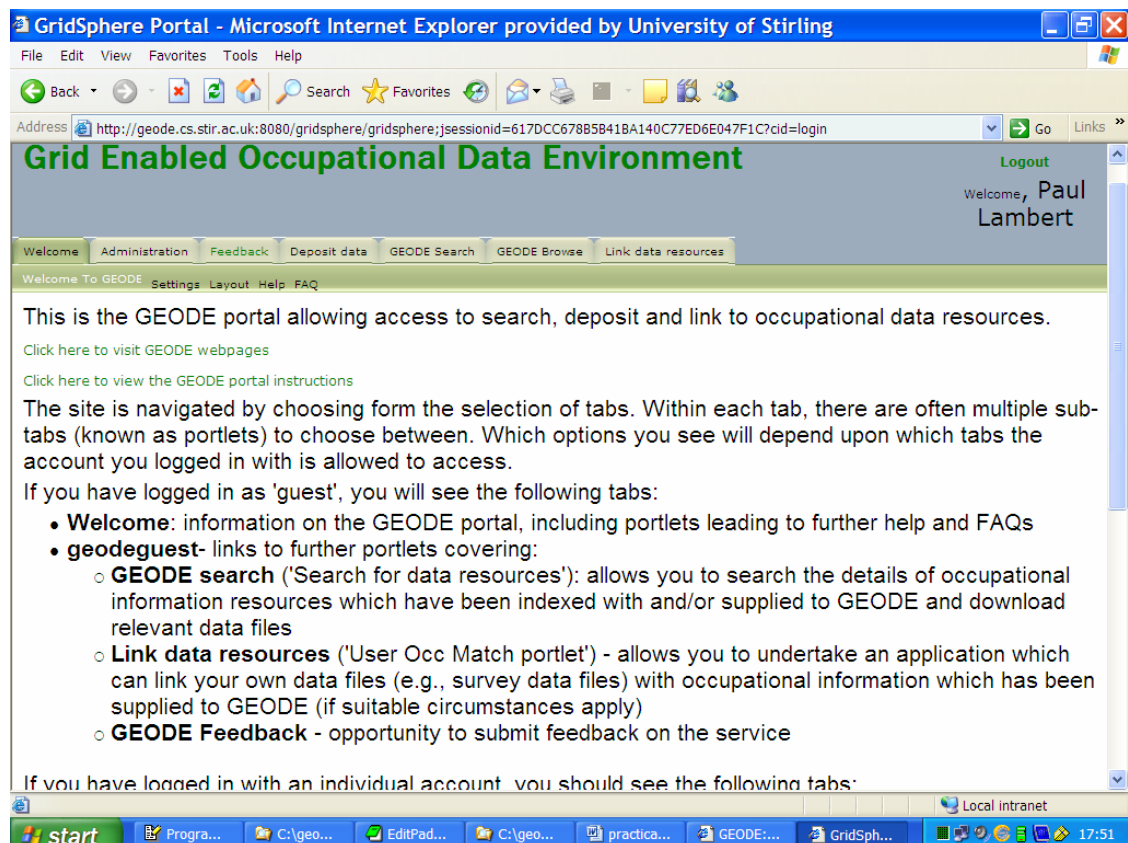| Occupational information resources *(details in text)* | |
|---|---|
| **(1)**-[A] | Hakim (1998) gender segregation statistics, for UK SOC-90 occupations |
| **(2)**-[A] | Index files for UK CAMSIS scale for SOC-90 occupations |
| **(3)**-[B] | Index files for UK CAMSIS scale for SOC-2000 occupations |
| **(4)**-[C] | Translation file for ISEI scale for ISCO-88 occupations across Europe |
| **(5)**-[C] | Translation file for EGP class for ISCO-88 occupations across Europe |
| **(6)**-[C] | Translation file for Skill levels for ISCO-88 occupations across Europe |
| **(7)**-[C] | Index file for CAMSIS scales for ISCO-88 occupations across Europe |

**Ex. 1) Entering the GEODE portal and searching for occupational information resources**

1.1)     **Login to the GEODE portal** by following the links from www.geode.stir.ac.uk , using your username and password

Your username: …………………………. Password: ………………………….

*(ask the instructors if you don't yet have a personalised GEODE account)*

You should see the 'Gridsphere' interface, which constitutes the GEODE portal, looking something like this:

## 1.2) Explore the GEODE portal a little.

The portal is navigated by following the tabs and sub-tabs (known as 'portlets'). The contents of the different tabs are descrbed briefly on the 'welcome' tab, and are also described in the technical paper 'Instructions on using the GEODE portal'. The 'FAQ' portlet under the 'welcome' tab may be especially helpful.

*…Gridsphere users beware…:*
- *hitting the 'back' button on a web browser doesn't always lead to the expected locations in Gridsphere. It's usually better to click on the available tab or portlet links*
- *you can open new tabs in new windows by right clicking etc, but we don't recommend this*
- *the windowing in gridsphere can be confusing – within a portlet there is sometimes an icon in the top right which minimises or maximises subsections of the portlet (potentially hiding the subsections which you were expecting to see)*
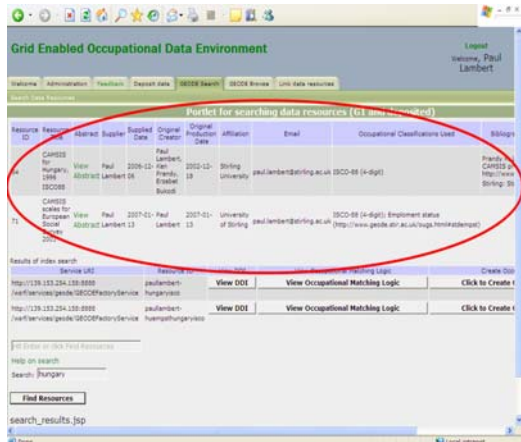- *you will be logged out of the portal after 30 minutes of inactivity*

## 1.3) Use the search tab - search for 'isco'

*The GEODE search tab allows you to search across all occupational information resources which have been indexed in the GEODE service. At present there are around 100 resources on the service. In time, there are likely to be far more.*
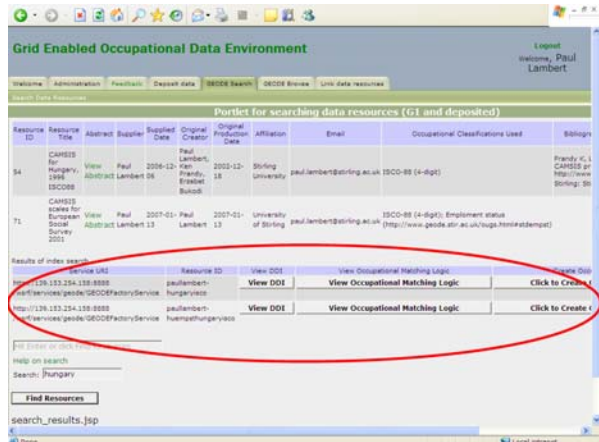
*There are a couple of important things to know about searching for occupational information resources through the GEODE portal:*

- **(A) There are two groups of resources which your search may uncover.** These are **(1) 'uncurated' resources**, which are data files, or links to data files, in any format, which have been notified to the GEODE portal, but haven't been subjected to any further treatment (curation). They are available for other users to access in their original format (e.g. to download), but are not fully integrated into GEODE. There are also **(2) 'curated' resources**, which have been integrated into the GEODE file matching process. These are the links that appear towards the end of the file searching linkage. These resources, which connect to data files available from corresponding 'uncurated' resources, are data files which have had sufficient metadata added to them that they can be used to run the GEODE portal file matching procedure on them.

| Uncurated resources: | Curated resources: |
|---|---|

- **(B) The rules which generate search results operate differently for the curated and uncurated resources** *[this situation is likely to be revised in the coming months]*. At present, uncurated resources are searched using an array of search logic terms which are described on the portal link 'help for search'. Curated resources however are searched only for exact matches of terms.

1.4)     **Use the 'GEODE browse' tab -** look for occupational information resources for the nineteenth century

The 'browse' function provides an alternative means to locate occupational information resources which have been indexed at GEODE.

*At time of writing [13.1.07], the 'browse' facility in GEODE is partially functional. It currently only refers to 'uncurated' occupational information resources, and only to those examples of uncurated information for which the appropriate 'browse' categorisations were specifically declared by the data depositor.*

**Ex. 2) Comparator: accessing and using occupational data through conventional internet sources**

We have alluded at several points to the difficulty of exploiting occupational information resources under current models of internet provision. In order to provide a comparison, we have prepared an illustrative example of a relatively complex exercise in matching occupational data, which corresponds to an exercise which we will shortly undertake by using the GEODE portal (Ex. 5). Exercise 2 involves obtaining appropriate data files from internet sources, then running either an SPSS or Stata demonstration command file to perform the linkages.

***\*\*We don't suggest trying this exercise unless you are fairly familiar with using SPSS and/or Stata to match data between different files\*\****

The exercise corresponds to datasets **(4)-[C]** and **(7)-[C]** noted on page 2. The starting point is a micro-data file from the European Social Survey (ESS) in 2001. Our aim is to match national specific CAMSIS scale scores, and internationally standarised ISEI scale scores, to the ISCO occupational units for each of the nine ESS countries for which such scale scores are available.

| SPSS instructions: | Stata instructions: |
|---|---|
| **2.1a) Open SPSS** | **2.1a) Open Stata** |
| **2.2a) You will need all of the following files in a suitable location on your machine:** | **2.2b) You will need all of the following files in a suitable location on your machine:** |
| - ess_spss_example.sps | - ess_stata_example.do |
| *(example SPSS command file)* | *(example Stata command file)* |
| - ess_extract2001.sav | - ess_extract2001.dta |
| *(original micro-data – extract from ESS)* | *(original micro-data – extract from ESS)* |
| - iskoisei.sps | - iskoisei.ado |
| *(SPSS format ISCO-ISEI linkage)* | *(Stata format ISCO-ISEI linkage)* |
| - gb91isco88.sav | - gb91isco88.dat |
| - ch90isco88.sav | - ch90isco88.dat |
| - cz94isco88.sav | - cz94isco88.dat |
| - hu96isco88.sav | - hu96isco88.dat |
| - ie96isco88.sav | - ie96isco88.dat |
| - plcherisco88.sav | - plcherisco88.dat |
| - ptcherisco88.sav | - ptcherisco88.dat |
| - sv94isco88.sav | - sv94isco88.dat |
| - se90isco88.sav | - se90isco88.dat |
| *(SPSS format CAMSIS-ISCO linkage files for each country)* | *(plain text format CAMSIS-ISCO linkage files for each country)* |
| **\*\*We have packaged up these two groups of files into the following zip archives\*\*:** | |
| www.geode.stir.ac.uk/workshop/data/ex2_SPSS.zip | |
| www.geode.stir.ac.uk/workshop/data/ex2_Stata.zip | |
| *Ordinarily this downloading process is a very long one. The micro-data would usually be obtained from www.europeansocialsurvey.org . The nine national specific CAMSIS files would normally be available from 9 different locations within www.camsis.stir.ac.uk/versions.html. The ISEI translation file in SPSS format could be downloaded from http://home.fsw.vu.nl/~ganzeboom/pisa/ . The ISEI translation file in Stata format could be downloaded from http://ideas.repec.org/c/boc/bocode/s425802.html .* | |

| | |
|---|---|
| **2.3a) Open the example Syntax file** (ess_spss_example.sps) in your SPSS session. | **2.3b) Open the example do file** (ess_stata_example.do) in the do-file editor of your Stata session. |
| **2.4a) Adjust the paths on the syntax file** to correspond to where you have stored the files mentioned above. | **2.4b) Adjust the paths on the do file** to correspond to where you have stored the files mentioned above. |
| **2.5a) Proceed through the syntax file, replicating the commands given** (which run processes on the files you have just downloaded). | **2.5b) Proceed through the do file, replicating the commands given** (which run processes on the files you have just downloaded). |
| *The end result is a micro-data file for the ESS which features ISEI and CAMSIS scale scores associated with the ISCO occupations supplied on the original survey.* | |
| **THAT'S ALL!!** | |

**Comments:**

*We have shown this example to try to illustrate the amount of manual work which can be involved in merging micro-data with occupational information resources. The examples we've shown should have run relatively smoothly (since the pre-prepared command files have been tested previously..). However readers will perhaps be able to appreciate that in normal circumstances, numerous difficulties emerge when an analyst tries to implement such linkages on new data resources.*

*We will shortly illustrate undertaking the same exercises by using the GEODE portal (Ex. 5). This exercise should proceed very quickly. In the interests of fairness, we should point out that the exercise we present for the GEODE portal isn't exactly equivalent to that demonstrated above, since it exploits a single index data file for CAMSIS which was created as part of the GEODE project (rather than the nine different files used above).*

**Ex. 3) Depositing new occupational information with GEODE**

*One of the most important features of the GEODE portal is its ability to act as a centralised depository for electronic format occupational information resources, in contrast to the diverse internet distribution which is currently used.*

*As noted on the GEODE webpages (e.g. [http://www.geode.stir.ac.uk/geode_m_curation.html](http://www.geode.stir.ac.uk/geode_m_curation.html)), the supply of data to GEODE is essentially a two stage process:*

> *First, the data is sent to (or a URI link to the data is sent to) the GEODE portal, by means of the 'deposit data' portlet. This information becomes available immediately to other users of GEODE.*

> *Second, the data can have further information added to it (curated) in the form of xml format information files with standardised records. This process is nearly always undertaken by members of the GEODE project, but it is possible for other users to add this information.*

*In the following exercise we deposit a new data resource to GEODE in 'uncurated' form. It is possible to undertake further curation stages by following the instructions in the paper 'Instructions on using the GEODE portal'.*

**3.1) Creating an occupational information file**

In this example we suggest you use a simple occupational information file we have made available online:

[www.geode.stir.ac.uk/workshop/data/isco88_majorgroups_skill.xls](www.geode.stir.ac.uk/workshop/data/isco88_majorgroups_skill.xls)

This file is a listing of the four levels of skill associated with ISCO-88 jobs as advocated in an article by Elias (1997).

This is often used as a simple type of occupational class scheme. However, it might not seem ideal to many perspectives, for example it regards all occupations in major groups 0 and 1 as 'unable to classify'.

To **create a new occupational information resource: download the Excel file, and add in a further column giving your own views** on the appropriate skill categories of each major group.

Save this file to a new name on your machine.

**3.2) Upload the Excel file you have created to the GEODE portal:**

- Click on 'deposit data'
- Click on 'manage my deposited data resources'
- Fill out the online form giving details of the data resource you have created (some sections may be left blank).

*There is also a description of this process of depositing uncurated data in the paper 'Instructions on using the GEODE Portal'.*

**3.3) Use the search facility to locate the resource you have just deposited.** At this stage, your resource is visible to any GEODE user.

*Warning: It may take up to two minutes after depositing a resource before it is read by the search engine.*

*Tip: Search for terms that you know you included during the entry form process, e.g. a word you used in the abstract.*

**3.4) Use the GEODE portal to adjust the data description** that you have just supplied:

- Click on the 'Deposit data' tab and then 'Manage my deposited data resources'
- Check the box beside the resource you've just uploaded (at the bottom of the page underneath the form; most likely the only resource visible)
- Click on 'edit'
- Make a small change to the text of the abstract
- Input data in the three drop-down boxes located at the end of the listing (these are the categorisations used for GEODE's 'browse' function; note that specifying categories for browsing like this is only possible at this stage, and couldn't have been undertaken during the initial supply of the data). Once these categories are specified, your resource will be visible to other users who are looking through the 'Browse' facility.
- Click 'Save' to preserve your resources

**3.5) Extension data curation.** If you wish to look into the issues of further curating a data resource at GEODE, there are instructions, oriented around the example of this skill classification data file, available in the paper 'Instructions on using the GEODE Portal'.

> **Keep GEODE tidy.** *Sometime after undertaking this exercise, we would recommend you use the data management facilities described in (3.4) to delete the resource you deposited here, if it serves no other function.*

**Ex.4) Matching occupational information to a regular format survey data file through GEODE**

> Exercises 4 and 5 now turn to using the GEODE portal in order to match occupational data resources to micro-data files.
>
> Instructions on this process are given at:
>
> - http://www.geode.stir.ac.uk/matching_occupational_data.html

*In Exercise 4 we use two relatively straightforward examples, referring to the data files named on page 2:*

Linking an extract from the British Household Panel Survey of 1991 [A] to:
⇕
Statistics on gender segregation for SOC-90 units (1)
CAMSIS scale scores for SOC-90 units (2).

Linking an extract from the Labour Force Survey of 2002 [B] to:
⇕
CAMSIS scale scores for SOC-2000 units (3)

This exercise allows for handling these files in SPSS and Stata formats. To begin the exercise you will need to download the datafiles [A] and [B] in your preferred format, from:

http://www.geode.stir.ac.uk/workshop/data/bhps_w1_extract.{sav/dta}
http://www.geode.stir.ac.uk/workshop/data/lfs_2002extract.{sav/dta}

**4.1) SOC-90 to gender segregation statistics**.

- Open the data file 'bhps_w1_extract' in your preferred format in SPSS or Stata
- Save it out to a plain text tab-delimited format (see instructions at http://www.geode.stir.ac.uk/file_convert_info.html)
- Enter the GEODE portal and search for 'soc90'
- Identify the 'G1' ('curated') resource which is called 'paullambert-hakimsoc' and click on 'Click to create occupational Resource'

- Switch to the 'Link data resources tab', and, beside the resource for 'hakimsoc', click on 'Match to my local data'
- This launches the JAVA application for matching the data files together. Fill out the form and run the matching process, following the instructions given on http://www.geode.stir.ac.uk/matching_occupational_data.html
- As your input data file, use the plain text extract you have just made from the SPSS/Stata file
- As your output data file, choose any name and location to export the plain text output data
- Note that the matching takes place 100 cases at a time.
- When the matching is finished, note that you've added a variable called 'genseg' to your data file (these are gender segregation statistics – the proportion of women in the SOC-90 occupation from the 1991 census).
- Use your preferred package SPSS or Stata to read the plain text output file (see the instructions on http://www.geode.stir.ac.uk/file_convert_info.html)

## 4.2) SOC-90 to CAMSIS scale scores.

- Enter the GEODE portal and search for 'soc90'
- Identify the 'G1' ('curated') resource which is called 'paullambert-gbsocukempst' and click on 'Click to create occupational Resource'
- Switch to the 'Link data resources tab', and, beside the resource for 'gbsocukempst', click on 'Match to my local data'
- This launches the JAVA application for matching the data files together. Fill out the form and run the matching process, following the instructions given on http://www.geode.stir.ac.uk/matching_occupational_data.html
- As your input data file, use the plain text data file that was created as output in (4.1)
- As your output data file, choose any name and location to export the plain text output data
- When the matching is finished, note that you've added several variables to the end of your data file (these are social classifications which can be identified by soc90 and employment status combinations, they include male and female CAMSIS scale scores, and the Registrar Generals Social Class scheme and the CASMIN social class scheme).
- Use your preferred package SPSS or Stata to read the plain text output file (see the instructions on http://www.geode.stir.ac.uk/file_convert_info.html)

> *Comment on converting micro-data files to plain text formats.*
>
> *We usually recommend that you extract only the minimal subset of your data in plain text format, before performing the occupational matching, then linking back to your original data (in order not to loose information in the formatting stages). Information on how to do this in SPSS and Stata is shown on the webpage http://www.geode.stir.ac.uk/file_convert_info.html .*
>
> *We hope in due course to initiate an automated file format transfer which will remove the need for users to perform these translations manually.*

## 4.3) SOC-2000 to CAMSIS scale scores and social classifications

- Open the data file 'lfs_2002extract' in your preferred format in SPSS or Stata
- Save it out to a plain text tab-delimited format (see instructions at http://www.geode.stir.ac.uk/file_convert_info.html)
- Enter the GEODE portal and search for 'soc2000'
- Identify the 'G1' ('curated') resource which is called 'paullambert-gbsockkukempst' and click on 'Click to Create Occupational Resource'
- Switch to the 'Link data resources tab', and, beside the resource for 'gbsockkukempst', click on 'Match to my local data'
- This launches the JAVA application for matching the data files together. Fill out the form and run the matching process, following the instructions given on http://www.geode.stir.ac.uk/matching_occupational_data.html
- As your input data file, use the plain text data file that you have just created from the lfs data file
- As your output data file, choose any name and location to export the plain text output data
- When the matching is finished, note that you've added several variables to the end of your data file (these are social classifications which can be identified by soc2000 and employment status combinations, they include male and female CAMSIS scale scores, and the Registrar Generals, NS-SEC, and CASMIN social class schemes).
- Use your preferred package, SPSS or Stata, to read the plain text output file (see the instructions on http://www.geode.stir.ac.uk/file_convert_info.html)

*Comment: Strengths and weaknesses of the GEODE matching facility*

At present the matching facility described above in not especially user friendly, though we are working on improving it. We welcome feedback.

What we see as weaknesses include:
- *Currently there is difficulty in users observing the connection between 'uncurated' and 'curated' resources, and consequently in understanding what exactly the curated resources are doing (i.e. what is 'gbsockkukempst'??!!).*
- *Currently there is no control over the names assigned to the output variable(s) obtained during matching, meaning a danger of duplicate or otherwise problematic variable names*
- *Currently there is no facility to match multiple occupational variables at the same time (e.g. own occupation and spouses occupation) – the processor needs to be run once over for each additional matching variable*
- *The need to manually transform data into plain text format in order to use the processor*

On the other hand, we think that the leading strengths of the processor involve:
- *Its ability to match occupational data independently of data analysis software requirements*
- *Its ability to communicate with a wide range of data resources, including resources which are stored online outwith the GEODE server*

**Ex. 5) Matching occupational information to a complex survey data file through GEODE**

Our last example concerns matching a wider array of occupational information to a more complex micro-data file. Exercise 5 involves:

Linking an extract from the European Social Survey of 2001 [C] to:
⇕
ISEI scale scores for ISCO-88 units (4)
EGP class positions for ISCO-88 units (5)
Skill level positions for ISCO-88 units (6)
National specific CAMSIS scale scores for ISCO-88 units (7).

We assume users will be handling the micro-data file in SPSS or Stata formats. To begin the exercise you will need to download the data file [C] in your preferred format, from:

http://www.geode.stir.ac.uk/workshop/data/ess_2001extract.{sav/dta}

The curated data resources ('G1' resources) on GEODE which can be used to achieve these matches are found at:

(4)     'paullambert-iscoiseib' *(search for 'isco')*
(5)     'paullambert-iscoegp' *(search for 'egp')*
(6)     'paullambert-iscomajgps' *(search for 'isco')*
(7)     'paullambert-esscamsis' *(search for 'isco')*

To perform the matches, follow the steps as in Ex.4 (also see http://www.geode.stir.ac.uk/matching_occupational_data.html).

Note the following issues:
- (5) and (7) both require employment status information in different classifications – use the 'standard category uri' to get information on these
- (7) involves linking by different countries, but this is relatively straightforward assuming a consistent country identifier (ISO country codes)

END OF PRACTICAL SESSION