

Validating and Extracting Information from Finnish National Identification Numbers with *hetu* Package

by Pyry Kantanen, Leo Lahti

Abstract National identification numbers (NIN) are a widely used method for uniquely identifying individuals in Finland and many other countries. Their widespread use in both public sector registers and private sector customer data systems prompts the need for openly available methods and tools for NIN analysis and validation. *hetu* R package provides functions for extracting information from Finnish NINs, personal identity codes, as well as generating valid and temporary codes for testing purposes. Personal identity codes contain information about person's birth date and sex as well as a control character used to Our package contains comprehensive documentation and examples to make usage of the package easy for domain experts from non-technical backgrounds. This work provides tools that are relevant mainly in the Finnish context, but also contribute to developing similar tools for other countries. *hetu* R package contributes to methods development in computational social science.

Introduction

Reliably keeping track of individuals in a given territory has long been seen as an important feature of modern governance. The scope of this tracking has spread from disciplinary institutions, such as prisons, military units and schools to society at large (Foucault, 2009). In everyday life settings, a combination of personal attributes, such as name, occupation and address may be sufficiently *unique* to identify individuals even in large groups of people (think of SQL's *composite key*). Such an approach has its limitations. The problem with such composite keys is, however, that parts that make the composite key might not be permanent, making it impossible to compare different data points in the same dataset or combining data from different datasets.

An ideal personal identification system should strive to be both unique and "self-same" over time. Self-sameness refers to a degree of immutability that allows organizations to identify and reidentify a person over time. A combination of attributes¹ such as name, occupation and address would probably form a unique identifier even in relatively large crowds, but such attributes might not stay the same over time. (Brensinger and Eyal, 2021).

According to Alterman (as cited in Ajana (2013, 8)) a distinction can be made between *biocentric data* and *indexical data*. The former is biometric data connected to the *body* of the individual whereas the latter has no distinguishable relation to the individual, physiologically, psychologically or otherwise. An example of biocentric data could be a fingerprint or an iris scan and an example of indexical data could be a randomly assigned number from which nothing could be deduced.²

In casual settings the combination of person's name and physiological characteristics (or *biometric identifiers*) are enough to identify friends, acquaintances and strangers from one another. Names can change over time and even fingerprints and iris pattern samples need to be refreshed from time to time, making them poor from self-sameness perspective. Interestingly, Lauer (2017) and Rule (1973) (as cited in Brensinger and Eyal (2021, 31–32)) noted that in the 1960's American credit bureaus and other organizations struggled with this same problem having used a combination of attributes to identify people in their registers and needed to find a better way to identify individuals. In a moment that may be described as *function creep* they settled for using the *social security number* (SSN) as the single identifying number³, horizontal and vertical information sharing possible.

For a number of reasons, many identification numbers are not just random strings. The American SSN originally contained information about the person's birth year and where the number was first registered (Brensinger and Eyal, 2021, 32) whereas Finnish and other Nordic national identification

¹Think of SQL's *composite key*

²Brensinger and Eyal (2021) discuss the concept of *dividuals*, manufactured objects that represent the living individual: address, fingerprints, name and so on. These dividuals need to go through the process of disembedding, standardization and reembedding to be useful. Disembedding means data gathering (taking a fingerprint sample), standardization means making the disembedded transcription into a standardized digital sample that can be easily compared with other similar samples and reembedding means linking these standardized records back to their actual flesh-and-blood counterparts. Without a way to reembed a huge and well standardized archive of fingerprints back to the population it is essentially useless. This is also a reason why biometric samples such as iris scans or fingerprints can never replace

³Think of SQL's *primary key*

numbers often contained information on the individual's birth date and sex (Watson, 2010; Salste, 2021). Therefore they may be seen as biocentric data rather than pure indexical data, making them more sensitive to handle.

In Nordic countries comprehensive national identification number systems were developed and implemented from the 1940's to 1960's (Watson, 2010). Finnish personal identity code has its roots in specialized employment pension number from 1962, which was gradually expanded to cover the whole population in the form of a social security number in 1964-1968. The structure of pension number was designed by mathematician Erkki Pale with punch card machines in mind (probably influenced by the Swedish födelsenummer or birth number). This design has proved to be resilient and with some tweaks it continues to be used, with the modern iteration being called a personal identity code (In Finnish: henkilötunnus, or *hetu* for short, hence the name of the package) (Salste, 2021).

Background

Similar R packages include **sweidnumbr** for Swedish personal numbers and organisation numbers (Magnusson and Bulow, 2020), **numbersBR** for Brazilian identity numbers for individuals, vehicles and organisations (Freitas, 2018), and **generator** for generating various types of Personally Identifiable Information (PII), such as fake e-mail addresses, names and United States Social Security Numbers (Hendricks, 2015).

The *hetu* R package is an rOpenGov project. rOpenGov is a community of R package developers interested in open government data analytics and related topics. Handling and analyzing personal identity codes is slightly different in scope than retrieving and analyzing European statistical data (Lahti et al., 2017) or Finnish administrative borders **geofi**. Another rOpenGov package, *sweidnumbr*, is especially important with regards to *hetu* package's development history as many function and parameter names were directly inspired by it and both packages share a key author, Måns Magnusson.

Finnish personal identity code generators and validators are nothing new. In fact, handling *hetu*-codes seems to be a common entry-level task to familiarize new computer science students with regular expressions, dates, string subsetting and similar concepts. The reasoning behind packaging these functions as a separate R package is the same as stated in Wickham and Bryan's R Packages Introduction: code sharing, transparency, user manuals, semantic versioning and documenting changes between versions (Wickham and Bryan, 2022). These are often lacking in hobbyist coding exercises and simple web applets.

Extracting information from Finnish personal identity codes

The method of validating and extracting information from identification numbers is manually doable and simple in principle but in practice becomes unfeasible with datasets larger than a few dozen observations. *Hetu*-package provides easy-to-use tools for programmatic handling of Finnish personal identity codes and Business ID codes (In Finnish: Yritys- ja Yhteisötunnus, or Y-tunnus for short). *Hetu*-package utilizes R's efficient vectorized operations and is able to generate and validate over 5 million Finnish personal identity codes or Business Identity Codes in under 10 minutes. This covers the practical upper limit set by the current population of Finland (5.5 million people), providing adequate headroom for handling of relatively large registry datasets.

Printing a table containing extracted information in a structured form:

```
x <- c("010101A0101", "111111-111C", "290201A010M")
hetu(x)
```

The results can be seen in Table 1. We can already see that there is a problem with the last personal identity code. 29th of February 2001 is not a valid date as 2001 is not a leap year.

	hetu	sex	p.num	ctrl.char	date	day	month	year	century	valid.pin
1	010101A0101	Female	010	1	2001-01-01	1	1	2001	A	TRUE
2	111111-111C	Male	111	C	1911-11-11	11	11	1911	-	TRUE
3	290201A010W	Female	010	M	<NA>	29	02	2001	A	FALSE

Table 1: Table printed with `hetu(pin = c("010101A0101", "111111-111C", "290201A010W"))`.

The generic way of outputting information found on individual columns is to use the standard `hetu()`-function with `extract`-parameter. For example:

```
hetu("010101A0101", extract = "sex")
[1] "Female"
hetu("010101A0101", extract = "date")
[1] "2001-01-01"
```

A valid extract parameter needs to be a column name in the table printed out by hetu-function. Most commonly used columns have their own function wrappers that are identical in output:

```
hetu_sex("010101A0101")
[1] "Female"
hetu_date("010101A0101")
[1] "2001-01-01"
```

By relying on [lubridate](#) we have also implemented a special function to calculate the age of individuals in years (default), months, weeks or days. The default option is to calculate the age at the current moment but it is possible to set a date at which the date is calculated.

```
hetu_age("010101A0101", date = "2004-02-01", timespan = "months")
The age in months has been calculated at 2004-02-01.
[1] 37
```

With `hetu_diagnostic()` function we can take a closer look at diagnostics results of each personal identity code. The function prints information about 10 different checks done on the code. In this case we only want to check results related to validity of the personal number component, whether the control character is valid, whether the control character is correct and whether the date is valid.

```
diagnostics <- hetu_diagnostic("290201A010M")
diagnostics[,c("valid.p.num", "valid.checksum", "correct.checksum", "valid.date")]
```

The results of `hetu_diagnostic` can be seen in Table 2. We can see that the checksum character is valid, meaning that the character is either a number or an ASCII letter, excluding letters that could be mistaken for another character: G, I, O, Q or Z. Correctness of the checksum means whether the character is correct when date and personal number parts are concatenated together, divided by 31 and the remainder of this operator is used as a lookup value from a table containing valid control characters.

When data is inputted manually without validity checks, input errors can creep in. Control character in Finnish personal identity codes combined with validity checks in `hetu`-package can help catch most obvious errors. In our example of we see that the date is incorrect but the control character is also incorrect, meaning that there has been an error in date or personal number input. Because the date is obviously wrong, we try three different dates and see if any of these are correct: 28th of February 2001, 29th of March 2001 and 29th of February 2000 (which was a leap year). Table 2 shows this comparison.

	hetu	valid.p.num	valid.checksum	correct.checksum	valid.date
1	290201A010M	TRUE	TRUE	FALSE	FALSE
2	280201A010M	TRUE	TRUE	FALSE	TRUE
3	290301A010M	TRUE	TRUE	FALSE	TRUE
4	290200A010M	TRUE	TRUE	TRUE	TRUE

Table 2: A subset of diagnostics of invalid personal identity code "290201A010M", it's variations, and finally the correct personal identity code "290200A010M".

Hetu-package can generate a large number of personal identity codes with `rhetu` function. The date range of generated identity codes can be changed with parameters, but it has a hard coded lower limit in the year 1860 and upper limit in the current date. It has been theorized that the oldest individuals that received a personal identity code in 1960s were born in 1860s. Personal identity codes are never assigned beforehand and therefore it is impossible to have valid personal identity codes that have a future date.

The function can also be used to generate so called temporary identity codes. Temporary identity codes are never used as a persistent and unique identifier for a single individual but as a placeholder in institutions such as hospitals when a person does not have a Finnish personal identity code or it is not known. They can be identified by having a personal number in the range of 900-999.

Here is an example of generating 4 temporary PINs and checking their validity with `hetu_ctrl` function:

```
x <- rhetu(n = 4, p.male = 0.25, p.temp = 1.0)
x
[1] "160237-938R" "131166-950X" "151184-9241" "250104A954R"
hetu_ctrl(x, allow.temp = TRUE)
[1] TRUE TRUE TRUE TRUE
```

As additional features, our package also supports similarly generating and checking the validity of Finnish Business ID (BID) numbers (in Finnish: Yhteisötunnus or Y-tunnus for short) and Finnish Unique Identification (FINUID) numbers (in Finnish: Sähköinen asiointitunniste or SATU for short). As the name implies, BIDs are used as unique identifiers for companies, organizations and other legal persons. Unlike personal identity codes, BIDs do not contain any information about the company. BIDs consist of a random string of 7 numbers followed by a dash and 1 control character, a number between 0 and 9. FINUID numbers are similar in the sense that they do not contain any biocentric data on the individual and are mainly used by government authorities in IT systems. FINUID numbers consist of 8 numbers and 1 control character calculated in the same way as in personal identity codes. Both of these ID numbers are an example of indexical data, as mentioned earlier.

```
bid_ctrl(c("0000000-0", "0000001-9"))
[1] TRUE TRUE

satu_ctrl("10000001N")
[1] TRUE
```

Discussion

Hetu R package provides a free and open source methods for validating and extracting data from a large number of Finnish personal identity codes.

The origins of this package can be traced to early 2010s when one curious individual wanted to analyze a large number of Finnish personal identity codes that were leaked to the internet by an anonymous hacker. The legality and morality of handling such dataset containing personal information was and is in a grey area at best. As developers of this package we cannot condone such activities, even if they are conducted out of curiosity and not ill intentions, but we acknowledge that we cannot prevent our users from doing that either.⁴

We have acknowledged beforehand that random personal identity codes generated with hetu-package could theoretically be used for purposes such as synthetic identity fraud.⁵ On the other hand it is important to note that such identity codes could also be created by hand as the same information that we have used in developing our algorithms is available at Finnish authorities web page (Digital and Population Data Services Agency, 2022). Our package does not make fraudulent activities significantly easier for malevolent individuals which is essential in judging the pros and cons of releasing this software to the public.

Similar data breaches have made people more wary of digital services. Privacy concerns can push Finland and other Nordic countries towards redesigning their national identification numbers to omit the embedded personal information sometime in the future. There is an ongoing government project led by Ministry of Finance to redesign the system of personal identity codes (Valtiovarainministeriö, 2022). These and other related policy and legislation changes will be closely monitored and, if necessary, the package functions will be adjusted accordingly.

The package is published under permissive GPL-2 license and we encourage users to study the code and modify it for their own use or submit improvements to our code repository on GitHub.

Bibliography

- B. Ajana. *Governing through Biometrics: The Biopolitics of Identity*. Palgrave Macmillan, New York, 2013. [p1]
- J. Brensinger and G. Eyal. The sociology of personal identification. *Sociological Theory*, 2021. doi: 10.1177/07352751211055771. URL <https://doi.org/10.1177/07352751211055771>. OnlineFirst. [p1, 4]

⁴The use of "Good, not evil" license texts is thought to be problematic and against the ethos of open source software

⁵see Brensinger and Eyal (2021, 32) for a short description of synthetic fraud related to American SSNs

- Digital and Population Data Services Agency. The personal identity code, 2022. URL <https://dvv.fi/en/personal-identity-code>. Accessed: 2022-01-17. [p4]
- M. Foucault. *Security, territory, population: lectures at the Collège de France, 1977-1978*. Palgrave Macmillan, New York, 2009. Editors: Michel Senellart, François Ewald, Alessandro Fontana, Arnold I. Davidson. [p1]
- W. Freitas. numbersbr: Validate, compare and format identification numbers from brazil, 2018. URL <https://CRAN.R-project.org/package=numbersBR>. R package version 0.0.2. [p2]
- P. Hendricks. generator: Generate data containing fake personally identifiable information, 2015. URL <https://CRAN.R-project.org/package=generator>. R package version 0.1.0. [p2]
- L. Lahti, J. Huovari, M. Kainu, and P. Biecek. Retrieval and Analysis of Eurostat Open Data with the eurostat Package. *The R Journal*, 9(1):385–392, 2017. doi: 10.32614/RJ-2017-019. URL <https://doi.org/10.32614/RJ-2017-019>. [p2]
- M. Magnusson and E. Bulow. sweidnumbr: R tools to handle of swedish identity numbers, 2020. URL <http://github.com/rOpenGov/sweidnumbr>. R package version 1.4.3. [p2]
- T. Salste. Henkilötunnus – ihmisten koodaaja. <https://www.tuomas.salste.net/doc/tunnus/henkilotunnus.html>, 2021. Accessed: 2021-12-13. [p2]
- Valtiovarainministeriö. Project on redesigning the system of personal identity codes, 2022. URL <https://vm.fi/en/project-on-redesigning-the-system-of-personal-identity-codes>. Accessed: 2022-01-17. [p4]
- I. Watson. A short history of national identification numbering in iceland. *Bifröst Journal of Social Science / Tímarit um félagsvísindi*, 1:51–89, 2010. ISSN 1670-7796. [p2]
- H. Wickham and J. Bryan. R packages, 2022. URL <https://r-pkgs.org/intro.html>. The book is mentioned to be a work-in-progress. Accessed: 2022-01-17. [p2]

Pyry Kantanen
University of Turku
Turku Data Science Group, Agora, 20014 University of Turku
Finland
pyry.kantanen@utu.fi

Leo Lahti
University of Turku
Turku Data Science Group, Agora, 20014 University of Turku
Finland
leo.lahti@utu.fi