

This is a self-archived version of an original article. This version may differ from the original in pagination and typographic details.

Author(s): Kopra, Juho; Mäkelä, Pia; Tolonen, Hanna; Jousilahti, Pekka; Karvanen, Juha

Title: Follow-up data improve the estimation of the prevalence of heavy alcohol consumption

Year: 2018

Version: Accepted version (Final draft)

Copyright: © The Author(s) 2018. Medical Council on Alcohol and Oxford University Press.

Rights: In Copyright

Rights url: <http://rightsstatements.org/page/InC/1.0/?language=en>

Please cite the original version:

Kopra, J., Mäkelä, P., Tolonen, H., Jousilahti, P., & Karvanen, J. (2018). Follow-up data improve the estimation of the prevalence of heavy alcohol consumption. *Alcohol and Alcoholism*, 53(5), 586-596. <https://doi.org/10.1093/alcalc/agy019>

Follow-up data improve the estimation of the prevalence of heavy alcohol consumption

January 24, 2018

Juho Kopra¹, Pia Mäkelä², Hanna Tolonen², Pekka Jousilahti² and Juha Karvanen¹

¹Department of Mathematics and Statistics,
University of Jyväskylä, Finland

²Department of Public Health Solutions,
National Institute for Health and Welfare,
Helsinki, Finland

Corresponding author:

Juho Kopra
Department of Mathematics and Statistics
P.O. Box 35 (MaD)
FI-40014 University of Jyväskylä
Finland
Email: juho.j.kopra@jyu.fi
Telephone: +358408053455
Fax: –

Abstract

Aims: We aim to adjust for potential non-participation bias in the prevalence of heavy alcohol consumption.

Methods: Population survey data from Finnish health examination surveys conducted in 1987-2007 were linked to the administrative registers for mortality and morbidity follow-up until end of 2014. Utilising these data, available for both participants and non-participants, we model the association between heavy alcohol consumption and alcohol-related disease diagnoses.

Results: Our results show that the estimated prevalence of heavy alcohol consumption is on average 1.5 times higher for men and 1.8 times higher for women than what was obtained from participants only (complete case analysis). The magnitude of the difference in the mean estimates by year varies from 0 to 9 percentage points for men and from 0 to 2 percentage points for women.

Conclusion: The proposed approach improves the prevalence estimation but requires follow-up data on non-participants and Bayesian modelling.

Keywords: survey, Bayesian analysis, register linkage, non-response bias, data missing not at random

Introduction

Reliable information about the prevalence of heavy alcohol consumption is important because alcohol-related health problems and undesired social consequences (Klingemann and Gmel, 2001) cause significant costs in many countries (Rehm et al., 2009). Prevalence estimates can be obtained through health surveys, but the low participation rates (Galea and Tracy, 2007) imperil the reliability of the results. If the participation is selective with respect to alcohol consumption, the estimates of alcohol use suffer from non-participation bias, which hinders their usability for decision-making. If non-participants have worse health than participants, the bias usually leads to an overly positive image of the health of the population.

Empirical evidence suggests that participation is often selective concerning alcohol consumption. Studies from Canada (Zhao et al., 2009), England (Boniface et al., 2017), Finland (Karvanen et al., 2016; Kopra et al.,

2017b), Norway (Torvik et al., 2012), Scotland (Gorman et al., 2014) Sweden (Romelsjö, 1989), and the United States (Dawson et al., 2014) conclude that non-participants drink more alcohol than participants. Follow-up studies have shown that non-participants tend to have a higher risk of alcohol-related diseases (Romelsjö, 1989; Jousilahti et al., 2005; Gorman et al., 2014; Christensen et al., 2015; Karvanen et al., 2016), and increased risk of hospitalisations and death (Jousilahti et al., 2005; Christensen et al., 2015; Karvanen et al., 2016), which indicates that non-participants tend to use more alcohol than participants. A study from Netherlands (Lahaut et al., 2002) found that non-participants are more often abstainers than participants, which is not directly interpretable as a conflicting result because many social factors may be associated with abstaining. An older study from Sweden did not find an indication of selective participation (Halldin, 1985).

In addition to selective non-participation, bias may be introduced by imperfect coverage of the target population by the survey sampling frame, and by questionnaire design. First, if some individuals of target population cannot be invited to a survey, the sample does not represent the population of interest and the estimates will be biased. Mäkelä and Huhtanen (2010) observed that in Finland, persons who cannot be invited to a survey due to missing home address have about four times higher risk for alcohol-related deaths. This caused a small bias in the population estimates. Second, Livingston and Callinan (2015) claim that quantity-frequency design of the alcohol use questions underestimates alcohol consumption by one-third compared to asking about drinking with a within-location beverage-specific design. Gmel (2000) reported that alcohol as a subject of survey study does not have an impact on participation in comparison to similar questionnaire without alcohol-related questions.

Studies from the United States (Dawson et al., 2014) and Finland (Mäkelä, 2003) have shown that the non-participation bias cannot be adjusted using just weights depending on basic demographic variables. Some studies adjust for selective non-participation utilising continuum of resistance model (Zhao et al., 2009; Meiklejohn et al., 2012; Boniface et al., 2017) but there are also other methods (Karvanen et al., 2016; Kopra et al., 2017b).

In Finnish health surveys participation has been decreasing, while reported alcohol consumption has been mainly increasing from 1960s to 2007 (Mäkelä et al., 2012). Jousilahti et al. (2005) report that in Finland non-participants have a higher risk of alcohol-specific diseases and death (Jousilahti et al., 2005), which is why we expect the estimates of heavy alcohol

consumption to be biased. The difference in the disease risk could be explained by heavier alcohol consumption among non-participants. From previous studies (Harald et al., 2007; Hirvonen, 2017), we know that participation in the FINRISK Study is affected by age, gender, area, and education (Reinikainen et al., 2017).

We aim to adjust for selective non-participation for heavy alcohol consumption and to estimate the prevalence of heavy alcohol consumption with reduced bias. We present a Bayesian solution that is based on mortality and morbidity follow-up data.

Methods

Data

We used data from the National FINRISK Study, which is a series of cross-sectional health examination surveys (Borodulin et al. , 2017) conducted in Finland every fifth year since 1972. We analyzed data for the years 1987–2007. Years 1972–1982 as well as 2012 were excluded because the questions of alcohol consumption were too different from the questions in 1987–2007.

In 1987 and 1992 studies the questions were essentially the same. In 1997 the study questions regarding the usage of cider or mild wine (alc. vol under 5%) were added. Otherwise the study remained the same as earlier. In 2002, question regarding consumption of red wine and other wines were separated from each other. Also in 2002, the total alcohol consumption was no more based on several alcohol beverage-specific questions but participants were advised to calculate their number of alcohol portions (standard drinks) consumed and sign asuitable class from quantity-frequency table. In 2007, the alcohol questions were the same as in 2002, but the instructions for the calculation of daily alcohol consumption were improved.

The surveys provide data on 25–74 -year-old adults from six regions of Finland. We restrict the data to people aged 25–65 -years since oldest age group 65–74 years old was not available in all areas until 2007. The survey consists of questionnaires, and a health examination carried out at a local study clinic. The sample was drawn from the National Population Register and was stratified by region, gender and 10-year age-group. In total, there are 44,317 invitees including 31,567 participants. The survey data contains self-reported alcohol consumption and background variables, age, gender,

region and study year for the whole sample. The survey utilised a beverage-specific quantity-frequency questions on alcohol use in the first three surveys and graduate frequency measure in the latter two surveys. The questions related to alcohol consumption are provided in Appendix B. The study questionnaires in 1987, 1992 and 1997 asked alcohol usage one type of alcohol beverages at a time: beer, spirits/vodka, long drink or cider (in 1997), wines and mild wines (in 1992 and 1997). In 2002 and 2007 the questionnaire was different, and individuals reported their alcohol consumption for all beverage types in one question.

From these questions, a number of standard drinks consumed per week during a previous 12 months was calculated. One standard drink equals 12 grams of pure alcohol which is equivalent to, e.g, one bottle of beer (33cl, 4.7 volume percent of alcohol). Based on the number of standard drinks, the (self-reported) total amount of 100% alcohol consumed in the previous 12 months was calculated.

Since our main interest was to estimate the prevalence of heavy alcohol users, we classified participants as heavy alcohol consumers and others (non-heavy alcohol consumers) as follows. The persons who reported consuming on average at least 24 standard drinks per week for men or at least 16 standard drinks per week for women during the one-year period before the examination were considered heavy alcohol consumers.

The survey data were linked to three registers: The Register of Completed Education and Degrees (Statistics Finland, 2016), Care Register for Health Care (National Institute for Health and Welfare , 2017) and Cause of Death Register (Statistics Finland, 2017) using personal identification code. The register-based data were available for both participants and non-participants. The level of education is categorised according to the International Standard Classification of Education (ISCED, 2011). We classified education into three levels: 1) high level (tertiary education, ISCED levels 5-8), 2) middle level (secondary education, ISCED levels 3-4) and 3) low level (primary education or less or unknown, ISCED levels 0-2). The Care Register gives data about the hospital visits with dates and ICD-codes for both participants and non-participants. From the Causes of Death Register, we obtain data about dates and ICD-codes of the cause of death.

Follow-up data contains the time-to-event (age) and ICD code of the first alcohol-related disease diagnosis or death. The ICD-codes we considered to be alcohol related are listed in Table 1. The follow-up begins from the survey and ends at the end of 2014. Persons who have neither alcohol-related disease

diagnosis nor alcohol-related cause of death are censored at the end of the follow-up. Deaths not related to alcohol are treated as censorings.

Table 1: The ICD-codes interpreted as alcohol-related events.

ICD-9:	
291	Alcohol-induced mental disorders
303	Alcohol dependence syndrome
357.5	Alcoholic polyneuropathy
425.5	Alcoholic cardiomyopathy
535.3	Alcoholic gastritis
571.0	Alcoholic fatty liver
571.1	Acute alcoholic liver disease
571.2	Alcoholic cirrhosis of liver
571.3	Alcoholic liver damage, unspecified
577.0D-F	Alcoholic disease of the pancreas, acute
577.1C-D	Alcoholic disease of the pancreas, chronic
980.0	Toxic effect ethyl alcohol
980.2	Toxic effect of isopropyl alcohol
980.8	Toxic effect of other specified alcohols
980.9	Toxic effect of other unspecified alcohol
E851	Accidental poisoning by alcohol
ICD-10:	
F10	Mental and behavioural disorders due to use of alcohol
G31.2	Degeneration of nervous system due to alcohol
G62.1	Alcoholic polyneuropathy
G72.1	Alcoholic myopathy
I42.6	Alcoholic cardiomyopathy
K29.2	Alcoholic gastritis
K70	Alcoholic liver disease
K85.2	Alcohol-induced acute pancreatitis
K86.0	Alcohol-induced chronic pancreatitis
T51	Toxic effect of alcohol
X45	Accidental poisoning or other exposure to alcohol
Y15	Poisoning by and exposure to alcohol, undetermined intent

Complete case analysis

The complete case analysis (e.g. mean estimate from the participants) assumes that participation is not selective concerning alcohol consumption. Violations of this assumption lead to bias. We compared the results of complete case analysis to a Bayesian approach, which relies on more realistic assumptions and allows for selective non-participation concerning heavy alcohol use.

Modelling approach

We applied a Bayesian approach introduced in (Kopra et al., 2017a) to estimate the prevalence of heavy alcohol consumption. The Bayesian model consists of three sub-models which are fitted simultaneously. The sub-models are:

1. Participation model,
2. Risk factor model, and
3. Survival model.

The mathematical formulas for the models are given in Appendix A.

The participation model describes which variables affect participation. Participation is defined as a binary indicator (0 or 1) for the availability information on alcohol consumption. This model is a logistic regression model with linear covariates for study year and age, and categorical variables for the region (4 levels), education (3 levels) and the alcohol consumption (binary). The model also takes into account the possible interactions of gender and study year, gender and alcohol consumption, and study year and alcohol consumption.

The risk factor model describes how alcohol consumption (heavy or non-heavy) varies by background variables. By background variables we mean age, gender, region, study year and education. The model is a logistic regression model with interactions for the year of birth with gender, region, study year and education.

The survival model describes the relationship between alcohol consumption and alcohol-related diseases. All disease events are combined and modelled as one survival outcome. The survival model is a piecewise constant

hazard model with one-year baseline hazard period terms. The model assumes monotonically increasing baseline hazard, which is accomplished using prior specification. In addition to baseline hazard, alcohol consumption is used as a regressor. Both baseline hazard terms and the regression coefficient are gender-specific. The model assumes that the disease risk of non-participants must be between the risks of heavy alcohol consumer participants and other participants. This follows from our reasoning that if the risk of non-participants were the same as the risk of heavy alcohol consumers, we would expect all of them to be heavy alcohol consumers. Similarly, if the risk of non-participants equaled to the risk of non-heavy alcohol consumers, we would expect none of them to be heavy alcohol consumers.

Prior distributions

We used weakly informative prior distributions that reflect the existing knowledge but have variances large enough to allow for surprises. This approach is recommended in textbooks on Bayesian statistics (Gelman et al., 2014), and there exists guidelines for elicitation of prior distributions (O’Hagan et al., 2006). The participation model needed an informative prior for the effect of heavy alcohol consumption on the participation, i.e., for the strength of selectivity mechanism. Some degree of subjectivity cannot be avoided in the prior specification. To define a weakly informative prior, we took a 45-year-old non-heavy alcohol consumer who participates with probability 0.7 as a reference and considered the prior probability for a heavy alcohol consumer who is otherwise similar. We elicit that there is 25% chance that person participates with probability p lower than 0.5 ($P(p \leq 0.5) = 0.25$), 35% chance for $p \leq 0.6$, and 50% chance for $p \leq 0.7$. The functional form of the prior distribution was chosen to be logistic distribution. These elicitations lead to logistics prior distribution with expected value zero and variance 1/2.05.

In the survival model, we applied monotonically increasing baseline hazards separately for men and women. The prior distribution for the first hazard term (25–26-year-olds) was a uniform distribution with a range from 0 to 20. From second hazard term (26–27-year-olds) to the last hazard term (99–100-years-old), each term had a uniform distribution with the lower limit being the value of previous baseline hazard term and upper limit 20.

All the remaining model parameters had normally distributed priors with zero mean and variance 1000. The prior distributions are presented using mathematical notation in Table 4 of Appendix A.

Imputations and model fitting

Alcohol consumption was missing for the non-participants. These missing values (heavy or non-heavy) were imputed simultaneously with Bayesian model fitting using data augmentation (Tanner and Wong, 1987). The model was fitted using Markov chain Monte Carlo (MCMC) (Robert and Casella, 2004) and implemented with Just Another Gibbs Sampler (JAGS) -software (Plummer, 2003) and R software (R Core Team, 2017) with `rjags` -package (Plummer, 2015). The convergence of MCMC chains were investigated using Brooks-Gelman \hat{R} diagnostics (Brooks and Gelman, 1998), and all the \hat{R} s were below 1.01 which indicates convergence. The model fitting utilised computational resources of IT Center for Science Ltd (CSC).

Results

Descriptive statistics

In Table 2, we present the descriptive statistics on age, education and gender. These variables are examined by study year, and comparisons can be made between participants and non-participant as well as between heavy alcohol consumers and other alcohol consumers.

The average age of the non-participants was lower than the average age of the participants. Over the years, the average age appears not to have changed much for the non-participants, but it has slightly increased for the participants. Among participants, the average age of heavy alcohol consumers has increased more rapidly than for non-heavy alcohol consumers. The average age of heavy alcohol consumers was 41.7 in 1987 (44.4 for non-heavy) and it has increased between each study being 47.3 for heavy alcohol consumers and 45.5 for non-heavy alcohol consumers in 2007. The average age of non-heavy alcohol consumers has also increased between the studies, except between the 1997 and 2002 when it decreased by 0.2 years.

The level of education has increased for both participants and non-participants during the study period. The non-participants tend to have low education more often than participants, and participants tend to have high education more often than non-participants. In 1987, there were a higher proportion of highly educated participants among heavy alcohol consumers than among non-heavy alcohol consumers. In 2007, the situation was opposite; the proportion of highly educated persons is higher for non-heavy alcohol consumers than for the heavy alcohol consumers. The proportion of women among participants has slightly increased from 52.0% to 53.4% during 1987–2007. Among non-heavy alcohol consumers, the proportion is higher: 52.7%–55.3%. Women are a minority among heavy alcohol consumers. There were 15.9% women among heavy alcohol consumers in 1987, and the proportion has notably increased being 27.8% in 2007.

The proportion of women was higher among the participants than among non-participants. The proportion of women among participating heavy alcohol users has been rapidly increasing over the years, while the corresponding proportion had not increased by much among non-heavy alcohol consumers.

The number of invitees, the participation rate and the number of events for both participant and non-participant men and women are presented in Table 3. During the study period, the proportion of heavy alcohol consumers

Table 2: Description of background information by study year for non-participants, participants, and heavy and non-heavy alcohol consumers among participants.

Year	Non-participants	Participants		
		All	Heavy alcohol consumers	Moderate alcohol consumers
Average age:				
1987	42.5	44.4	41.7	44.4
1992	41.8	44.7	44.5	44.7
1997	42.8	45.0	45.0	45.0
2002	42.3	44.9	45.7	44.8
2007	41.9	45.6	47.3	45.5
High education (%):				
1987	13.6	18.5	23.0	18.4
1992	18.4	26.6	30.2	26.5
1997	24.7	29.9	31.4	29.8
2002	25.6	35.7	32.7	35.8
2007	27.7	38.6	34.3	38.8
Middle education (%):				
1987	30.6	31.8	31.0	31.8
1992	35.6	34.9	34.4	34.9
1997	37.7	38.1	37.6	38.1
2002	41.8	40.6	43.5	40.4
2007	45.3	44.2	45.5	44.1
Low education (%):				
1987	55.8	49.7	46.0	49.8
1992	46.1	38.5	35.4	38.6
1997	37.6	32.1	31.0	32.1
2002	32.6	23.8	23.9	23.8
2007	27.0	17.2	20.2	17.1
Women (%):				
1987	42.3	52.0	15.9	52.7
1992	40.5	53.0	21.2	54.1
1997	43.2	52.9	22.8	54.3
2002	41.6	53.6	29.0	54.9
2007	42.7	53.9	27.8	55.3

Table 3: Number of invitees, the participation rate, the prevalence of heavy alcohol consumption based on participants and Bayesian modelling (posterior mean), and the number of alcohol-related incident events (per 1000 follow-up years) for the non-participant and the participant men and women.

Year	Invited N	Participation rate	Prevalence for participants	Posterior mean	Alcohol-related incident events (per 1000 follow-up years)	
					Participants	Non-participants
Men						
1987	3910	79.5%	5.0%	9.6%	202 (2.8)	91 (5.2)
1992	3888	73.3%	7.7%	15.0%	168 (2.9)	128 (6.5)
1997	4034	70.0%	9.2%	9.7%	150 (3.2)	103 (5.4)
2002	3955	66.5%	14.4%	22.9%	118 (3.6)	62 (3.8)
2007	3202	61.8%	11.0%	12.8%	47 (3.1)	35 (3.7)
Women						
1987	3961	85.1%	0.7%	2.4%	52 (0.6)	29 (2.0)
1992	3951	81.0%	2.5%	4.4%	46 (0.7)	23 (1.5)
1997	4031	75.8%	3.7%	5.3%	25 (0.5)	37 (2.3)
2002	4019	75.4%	5.7%	5.9%	34 (0.9)	13 (1.0)
2007	3278	71.3%	4.0%	5.8%	12 (0.7)	10 (1.3)

has increased for both men and women among participants, and simultaneously the participation rate has decreased.

The probabilities for not having alcohol-related disease diagnosis up to the given age for men and women are presented by Kaplan-Meier survival plots (Figure 1). The top row shows that the non-participants were more likely to have alcohol-related diagnoses than participants. The lower row shows that the risk for non-participants lies between the risks of heavy and non-heavy alcohol consumers, which is a requirement for the utilised Bayesian model. The number of persons with a disease diagnosed in each group is reported next to the survival curve in the Figure 1.

Adjusted prevalences of heavy alcohol consumption

Figure 2 presents the trends of the prevalence of heavy alcohol consumption, based on complete case analysis and the Bayesian modelling. It can be seen that the mean estimates of the Bayesian approach lie above the estimates of the complete case analysis. The numeric values are presented in Table 3.

To compare the prevalence estimates based on participants only, and the

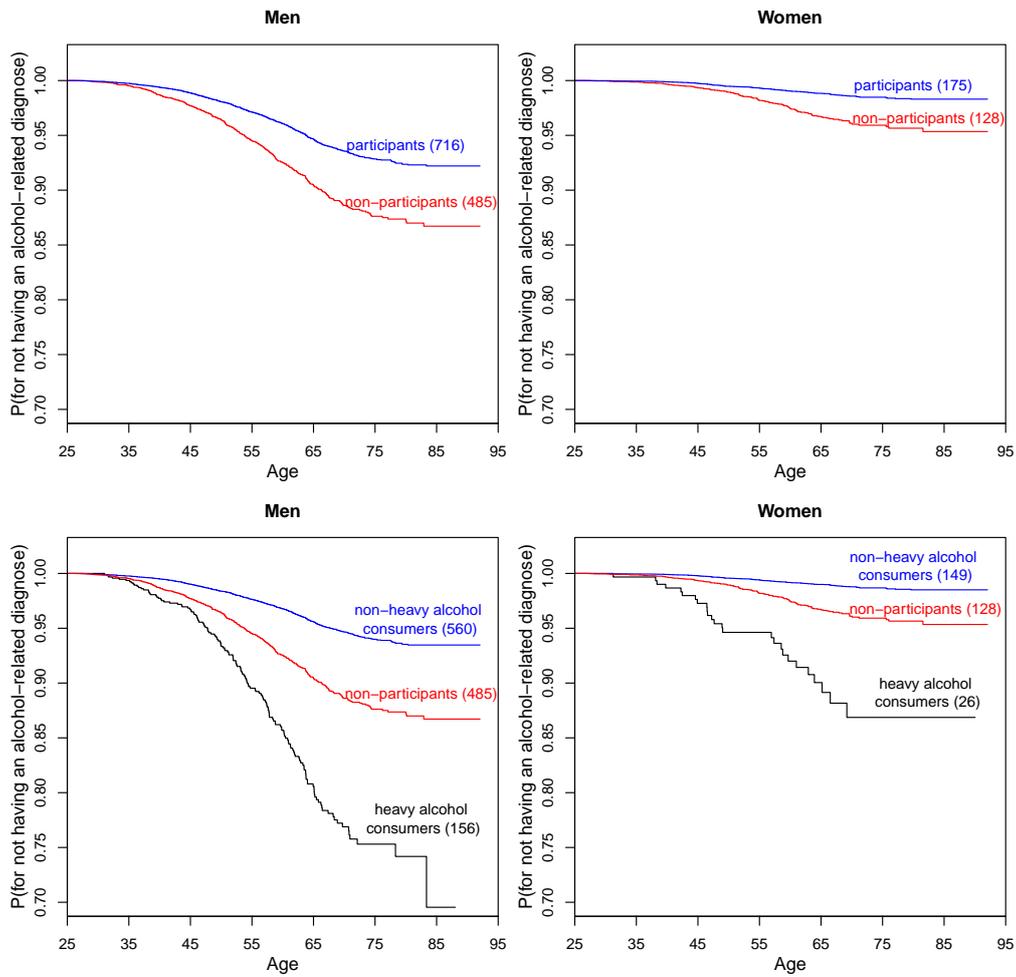


Figure 1: Kaplan-Meier survival plots for men and women comparing the probabilities of not having alcohol-related diagnoses among participants and non-participants (upper panels) and for heavy, non-heavy alcohol consumers and non-participants (lower panels). The number of persons with a disease diagnosed in each group is reported within parenthesis.

posterior estimate for the prevalence of entire survey, absolute and relative differences can be calculated. For men, the absolute difference of the yearly prevalence estimates for 1987–2007 are 4.6, 7.3, 0.5, 8.6, and 1.8 percentage points calculated from Table 3, respectively. Those lead to average difference of 4.6 percentage points. The corresponding relative differences for men are 1.93 (i.e. almost a two-fold difference), 1.95, 1.06, 1.6 and 1.17, respectively, and average relative difference is 1.5. For women, the corresponding values are yearly absolute differences; 1.7, 1.9, 1.6, 0.3 and 1.9, respectively, leading to average absolute difference of 1.5 percentage points. The yearly relative differences are 3.39, 1.77, 1.42, 1.04 and 1.47, respectively, leading to average relative difference of 1.8, see Table 3.

For men, the mean estimates based on Bayesian model vary year by year, but the credible intervals do not exclude the possibility of a monotonically increasing trend from 1987 to 2002. The complete case estimates are outside of the 90% credible interval of Bayesian trends in 1987, 1992, and 2002.

The credible intervals are narrower for women than for men. For women, the complete case prevalence estimates are outside of the 90% credible intervals of Bayesian trends in 1987 and 1992, and are within the credible interval in 1997, 2002 and 2007.

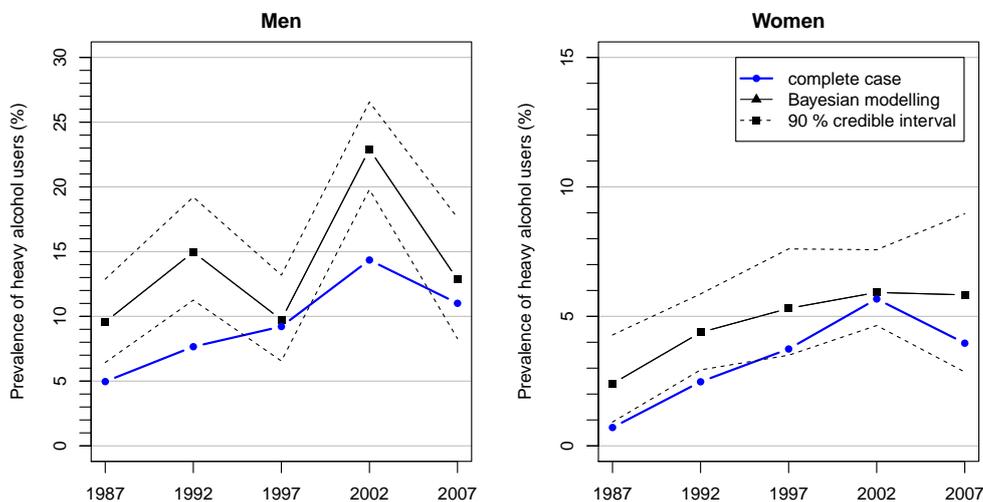


Figure 2: Comparison of prevalence estimates of complete case analysis and Bayesian multiple imputation adjusted for education. Note that the scales of the vertical axis for men and women are different from each other.

Discussion

There is evidence that non-participation in a survey asking about alcohol consumption is selective with respect to heavy alcohol consumption in Finland and in many other countries. We studied the prevalence of heavy alcohol consumption based on data from the National FINRISK Study, which suffer from selective non-participation. In FINRISK data, the average self-reported alcohol consumption for men was equal to 5.9 liters and for women 1.9 liters of pure 100% alcohol per year. For comparison, the national consumption statistics by National Institute for Health and Welfare (2016) show that the average yearly consumption of 100% alcohol for persons at least 15 years old was in the range of 10–13 liters per person during 1987–2007. Thus, in FINRISK data the self-reported consumption is about 60–70% lower what has been reported in the national consumption statistics (which were not used in our modelling in any way). Although many reasons can partly explain the differences between the consumption statistics and self-reported data, e.g. questionnaire design and imperfect matching of survey frame with the target population, the differences between non-participants and participants in the follow-up data summarized in Figure 1 suggest that selection bias is present.

We observed differences in alcohol-related events for participants and non-participants. Non-participants had significantly increased risk for alcohol-related disease or death compared to participants, and men had a higher risk than women. This phenomenon has also been observed for other data, see (Romelsjö, 1989), (Gorman et al., 2014) and (Christensen et al., 2015).

When participation is selective with respect to variables to be studied, which is the case for alcohol use, the estimates from complete case analysis are affected by non-participation bias and the real level of uncertainty is hidden, e.g. confidence intervals are not wide enough when complete case analysis is used. Mäkelä (2003) and Dawson et al. (2014) demonstrated that this kind of bias cannot be reduced for alcohol data with demographic information.

We compared the estimates obtained by a complete case analysis to estimates obtained by adjusting for non-participation with a full Bayesian modelling approach. The Bayesian approach gave a higher estimate of heavy alcohol consumption than the complete case analysis. Our approach reduced the bias and made the uncertainty visible. We estimated that the magnitude of bias is 0–9 percentage points for men and 0–2 percentage points for women in the FINRISK data. The Bayesian mean estimate was on average 1.5 times higher for men and 1.8 times higher for women compared to participants.

The use of our approach requires follow-up data and background variables for the entire invited sample (including non-participants), follow-up time long enough to observe alcohol-related disease events and Bayesian modelling. The first requirement cannot be fulfilled in many countries because of lack of register data or legal restrictions for data linkage. The second requirement means that the prevalence estimates will be available only several years after the survey. This requirement may be relaxed if there exist earlier surveys that can be assumed to share the same model parameters with the current survey. The third requirement is the easiest to fulfill because it only calls for statistical expertise that is widely available.

To conclude, the prevalence of heavy alcohol consumption based on survey participants only appears to be biased downward for both men and women. The magnitude of observed absolute bias was larger for men than women. The proposed non-participation adjustment approach is useful in context of alcohol research when follow-up data on non-participants are available, and the modelling requirements are met. The follow-up data can be used to improve the estimation of the prevalence of heavy alcohol consumption.

Declaration of conflicting interest

The Authors declare that there is no conflict of interest.

Funding

This work was supported by the Finnish Foundation for Alcohol Studies and Academy of Finland [grant numbers 266251 and 311877].

References

- Boniface S, Scholes S, Shelton N, Connor J. (2017) Assessment of non-response bias in estimates of alcohol consumption: applying the continuum of resistance model in a general population survey in England. *PloS one*, 12, e0170892.
- Borodulin K, Tolonen H, Jousilahti P et al. (2017) Cohort Profile: The National FINRISK Study. *Int J Epidemiol* doi: 10.1093/ije/dyx239

- <http://dx.doi.org/10.1093/ije/dyx239> (28 November 2017, date last accessed).
- Brooks SP, Gelman A. (1998) General methods for monitoring convergence of iterative simulations. *J Comput Graph Stat* 7: 434-455.
- Christensen AI, Ekholm O, Gray L, Glümer C, Juel K. (2015) What is wrong with non-participants? Alcohol-, drug- and smoking related mortality and morbidity in a 12-year follow up study of participants and non-participants in the Danish Health and Morbidity Survey. *Addiction* 110: 1505-12.
- Dawson DA, Goldstein RB, Pickering RP, Grant BF. (2014) Nonresponse bias in survey estimates of alcohol consumption and its association with harm. *J Stud Alcohol Drugs* 75: 695-703.
- Galea S, Tracy M. (2007) Participation rates in epidemiologic studies. *Ann Epidemiol* 17: 643-653.
- Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB. (2014) *Bayesian Data Analysis*, Third edition. Boca Raton, FL, USA: Chapman & Hall/CRC.
- Gmel G (2000) The effect of mode of data collection and of nonresponse on reported alcohol consumption: a splitsample study in Switzerland. *Addiction* 95: 123-134.
- Gorman E, Leyland AH, McCartney G et al. (2014) Assessing the representativeness of population-sampled health surveys through linkage to administrative data on alcohol-related outcomes. *Am J Epidemiol* 180: 941-948.
- Gray L, McCartney G, White IR, Katikireddi SV, Rutherford L, Gorman E, Leyland AH. (2013) Use of record-linkage to handle non-response and improve alcohol consumption estimates in health survey data: a study protocol. *BMJ Open* 3:e002647.
- Halldin J. (1985) Alcohol consumption and alcoholism in an urban population in central Sweden. *Acta Psychiat Scand* 71: 128-140.
- Harald K, Salomaa V, Jousilahti P, Koskinen S, Vartiainen E. (2007) Non-participation and mortality in different socioeconomic groups: the FIN-RISK population surveys in 1972-92. *J Epidemiol Commun H* 61: 449-454.

- Hirvonen E. (2017) Puuttavuuden mallintaminen FINRISKI -tutkimuksessa (in Finnish). Masters thesis, University of Jyväskylä, <http://urn.fi/URN:NBN:fi:jyu-201706192945>. (19 October 2017, date last accessed).
- International Standard Classification of Education: ISCED 2011. UIS, Montreal, Quebec.
- Jousilahti P, Salomaa V, Kuulasmaa K, Niemelä M, Vartiainen E. (2005) Total and cause specific mortality among participants and non-participants of population based health surveys: a comprehensive follow up of 54 372 Finnish men and women. *J Epidemiol Commun H* 59: 310-315.
- Karvanen J, Tolonen H, Härkänen T, Jousilahti P, Kuulasmaa K. (2016) Selection bias was reduced by re-contacting nonparticipants. *J Clin Epidemiol* 76: 209-217.
- Klingemann, H, Gmel, G (Eds.) (2001) Mapping the social consequences of alcohol consumption. Dordrecht, The Netherlands, Kluwer Academic Publishers.
- Kopra J, Härkänen T, Tolonen H, Karvanen J. (2015) Correcting for non-ignorable missingness in smoking trends. *Stat*, 4: 1-14.
- Kopra J, Karvanen J, Härkänen T. (2017a) Bayesian models for data missing not at random in health examination surveys. *Stat Model*, Advance online publication, doi:10.1177/1471082X17722605.
- Kopra, J, Härkänen, T, Tolonen, H, Jousilahti, P, Kuulasmaa, K, Reinikainen, J, Karvanen, J. (2017b) Adjusting for selective non-participation with re-contact data in the FINRISK 2012 survey. *Scand J Public Health*, Advance online publication, doi:<https://doi.org/10.1177/1403494817734774>.
- Lahaut VM, Jansen HA, Van de Mheen D, Garretsen HF. (2002) Adjusting for selective nonparticipation with recontact and hospitalisation history data. *Alcohol Alcoholism* 37: 256-260.
- Livingston M, Callinan S. (2015) Underreporting in alcohol surveys: whose drinking is underestimated? *J Stud Alcohol Drugs* 76: 158-164.

- Meiklejohn J, Connor J, Kypri K (2012) The effect of low survey response rates on estimates of alcohol consumption in a general population survey. *PLoS One* 7: e35527.
- National Institute for Health and Welfare *Care Register for Health Care*. Helsinki, Finland. <http://www.thl.fi/en/web/thlfi-en/statistics/information-on-statistics/register-descriptions/care-register-for-health-care>. (26 October 2017, date last accessed).
- National Institute for Health and Welfare: *Alcoholic Beverage Consumption 2016 [e-publication]*. Helsinki, Finland. <https://www.thl.fi/fi/tilastot/tilastot-aiheittain/paihteet-ja-riippuvuudet/alkoholi/alkoholijuomien-kulutus>. (14 September 2017, date last accessed).
- Mäkelä P. (2003) Impact of correcting for nonresponse by weighting on estimates of alcohol consumption. *J Stud Alcohol* 64: 589-596.
- Mäkelä P, Huhtanen P. (2010) The effect of survey sampling frame on coverage: The level of and changes in alcoholrelated mortality in Finland as a test case. *Addiction* 105: 1935-1941.
- Mäkelä P, Tigerstedt C, Mustonen H. (2012) The Finnish drinking culture: change and continuity in the past 40 years. *Drug Alcohol Rev* 31: 831-840.
- O'Hagan A, Buck CE, Daneshkhah A, Eiser JR, Garthwaite PH, Jenkinson DJ, Oakley JE, Rakow T. (2006) *Uncertain Judgements: Eliciting Experts' Probabilities*. Chichester, England, John Wiley & Sons.
- Plummer M. (2003). *JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling*. In Proceedings of the 3rd international workshop on distributed statistical computing. **124**, p. 125. Wien, Austria: Technische Universit at Wien.
- Plummer M. (2015). *rjags: Bayesian Graphical Models using MCMC*. R package version 3-15. <https://CRAN.R-project.org/package=rjags>. (8 November 2017, date last accessed).
- R Foundation for Statistical Computing (2017) *R: A Language and Environment for Statistical Computing*. Vienna, Austria. <https://www.R-project.org/>. (8 November 2017, date last accessed).

- Robert CP, Casella G. (2004) Monte Carlo Statistical Methods. Vienna, Austria, Springer Texts in Statistics.
- Rehm J, Mathers C, Popova S, Thavorncharoensap M, Teerawattananon Y, Patra J. (2009) Global burden of disease and injury and economic cost attributable to alcohol use and alcohol-use disorders. *Lancet*, 373: 2223-2233.
- Reinikainen J, Tolonen H, Borodulin K et al. (2017) Participation rates by educational levels have diverged during 25 years in Finnish health examination surveys. *Eur J Public Health*, Advance online publication, doi:10.1093/eurpub/ckx151.
- Romelsjö A. (1989) The relationship between alcohol consumption and social status in Stockholm. Has the social pattern of alcohol consumption changed?. *Int J Epidemiol* 18: 842-851.
- Official Statistics of Finland (OSF): *Causes of death [e-publication]*. Helsinki, Finland. http://tilastokeskus.fi/til/ksyyt/index_en.html. (26 October 2017, date last accessed).
- Statistics Finland *The Register of Completed Education and Degrees*. Helsinki, Finland. https://www.stat.fi/til/kou_en.html. (8 November 2017, date last accessed).
- Tanner M, Wong W. (1987) *The calculation of posterior distributions by data augmentation*. *J Am Stat Assoc* 82: 528-540.
- Torvik FA, Rognmo K, Tambs K. (2012) Alcohol use and mental distress as predictors of non-response in a general population health survey: the HUNT study. *Soc Psych Psych Epid* 47: 805-816.
- Vartiainen E, Laatikainen T, Peltonen M et al. (2009) Thirty-five-year trends in cardiovascular risk factors in Finland. *Int J Epidemiol* 39: 504-518.
- Zhao J, Stockwell TIM, MacDonald S. (2009) Nonresponse bias in alcohol and drug population surveys. *Drug Alcohol Rev* 28: 648-657.

Appendix A

The statistical model is based on work presented in (Kopra et al., 2017a), which utilises similar model to estimate the prevalence of smoking.

Notation for the data

For each individual $i = 1, \dots, N$ invited to a survey, we denote M_i being indicator of participation ($M_i = 1$ for participants and $M_i = 0$ for non-participants). Background information X_i consists of both survey frame and education information and is available for both participants and non-participants. The survey frame has variables gender g_i , region r_i , age a_i and study year s_i . The education is denoted by e_i , and values 1, 2 and 3 corresponds to high, middle and low education, respectively. Thus $X_i = (g_i, r_i, a_i, s_i, e_i)$. The heavy alcohol consumption Y_i is a binary variable such that heavy alcohol consumers have $Y_i = 1$ and non-heavy alcohol consumers have $Y_i = 0$. The variable T_i is the age at the first diagnosis of any of the alcohol-related diseases. The T_i is right censored and left-truncated at the age when the person entered the study.

Participation model

The participation model

$$\begin{aligned} \text{logit}(P(M_i = 1|X_i, Y_i)) = & \alpha_{0[g_i, s_i]} + \alpha_{1[g_i, s_i, e_i]} + \eta_{[g_i, s_i]} Y_i \\ & + \alpha_{2[g_i, Y_i]}(a_i - 45) + \alpha_{3[r_i]}, \end{aligned} \quad (1)$$

is a logistic regression model with following parameters. First, parameter $\alpha_{0[g_i, s_i]}$ is a constant where notation $[g_i, s_i]$ indicates that there are independent α_0 parameters for all levels of gender g_i and study year s_i . Second, parameter $\alpha_{1[g_i, s_i, e_i]}$ is the regression coefficient for education levels. For the lowest education level this parameter is forced to be 0. The parameter $\eta_{[g_i, s_i]}$ describes how heavy alcohol consumption affects participation. For this parameter, we need an informative prior. The parameter $\alpha_{2[g_i, Y_i]}$ describes how age at study affect participation. Finally, $\alpha_{3[r_i]}$ is a term for the region. For one of the regions, this parameter is forced to be 0. We selected a model that included important factors affecting participation while ensuring the convergence of the MCMC chains in Bayesian inference.

Risk factor model

The model for risk factor (heavy alcohol consumption) is

$$\text{logit}(P(Y_i = 1|X_i)) = \beta_{0[g_i, r_i, s_i, e_i]} + (s_i - a_i - 1938)\beta_{1[g_i, r_i, s_i, e_i]}. \quad (2)$$

The risk factor model is stratified by gender g_i , region r_i , study year s_i and education e_i using similar notation as in (1). The parameter $\beta_{0[g_i, r_i, s_i, e_i]}$ is constant for persons born in 1938. The parameter $\beta_{1[g_i, r_i, s_i, e_i]}$ determines how the heavy alcohol consumption prevalence changes with the year of birth.

Survival model

Let $dN_i(t)$ be the number of new events (increment) for the individual i at the time t . The increment follows a Poisson distribution with intensity parameters $\lambda_i(t)$. The intensity $\lambda_i(t)$ is modelled independently for both genders consisting of one-year period piecewise-constant baseline hazard terms $h_{0,0}(t)$ for men and $h_{0,1}(t)$ for women, and heavy alcohol consumption term $\exp(\gamma_1 Y_i)$ and $\exp(\gamma_2 Y_i)$ indicating the effect of heavy alcohol consumption for men and women, respectively

$$dN_i(t) \sim \text{Poisson}(\lambda_i(t))$$
$$\lambda_i(t) = \begin{cases} \exp(\gamma_1 Y_i) h_{0,0}(t), & \text{given that } T_i \geq t \text{ and } g = 0 \\ \exp(\gamma_2 Y_i) h_{0,1}(t), & \text{given that } T_i \geq t \text{ and } g = 1 \\ 0, & T_i < t. \end{cases}$$

Prior distributions

The prior distributions are specified in Table 4. The prior distributions for piecewise constant hazard terms $h_{0,0}(t)$ and $h_{0,1}(t)$ are specified such that the hazard becomes increasing function with respect to t .

Table 4: Prior distributions.

Notation	Distribution	Interpretation
Participation model		
η	Logistic(0, $\tau = 2.05$)	How heavy alcohol consumption affects the participation.
$\alpha_0, \alpha_1, \alpha_2, \alpha_3$	N(0, 1000^{-1})	Other parameters.
Risk factor model		
β_0, β_1	N(0, 1000^{-1})	Other parameters.
Survival model		
$h_{0,0}(25)$	Unif(0, 20)	Hazard for men at age 25–26.
$h_{0,1}(25)$	Unif(0, 20)	Hazard for women at age 25–26.
$h_{0,0}(t), t = 26, 27, \dots$	Unif($h_{0,0}(t-1)$, 20)	Hazard for men at age t .
$h_{0,1}(t), t = 26, 27, \dots$	Unif($h_{0,1}(t-1)$, 20)	Hazard for women at age t .
γ_1	N(0, 1000^{-1})	How heavy alcohol consumption affects hazard for men.
γ_2	N(0, 1000^{-1})	How heavy alcohol consumption affects hazard for women.

Appendix B

The study questions in 1987

CONSUMPTION OF ALCOHOL

1. Do you use any alcoholic drinks, even occasionally (f. ex. beer, wine or spirits)?

- 1 yes
- 2 no, but I have not quitted completely
- 3 no, because I quit using alcohol years ago
- 4 I have never used alcohol

If you have quitted alcohol use, please specify, why did you quit?

	no	yes
For health reasons	1	2
For economic reasons	1	2
For other reasons	1	2

2. Have you during the past year (last 12 months) had any alcohol (beer, wine or spirits)?

- 1 yes
- 2 no (for your part, the questions are completed)

3. How often do you usually drink beer (III or IV A)?

- 1 daily
- 2 a few times a week
- 3 about once a week
- 4 few times a month
- 5 about once a month
- 6 about once in a few months
- 7 3 - 4 times a year
- 8 twice a year
- 9 once a year or more seldom

0 never

4. How much do you usually drink beer at a time?

1 less than one bottle

2 1 bottle

3 2 bottles

4 3 bottles

5 4 - 5 bottles

6 6 - 9 bottles

7 10 - 14 bottles

8 15 bottles or more

9 I do not drink beer

5. How often do you usually drink wine (light or strong, also home made)?

1 daily

2 a few times a week

3 about once a week

4 a few times a month

5 about once a month

6 about once in a few months

7 3 - 4 times a year

8 twice a year

9 once a year or more seldom

0 never

6. How much do you usually drink wine at a time?

1 half a glass

2 one glass

3 two glasses

4 about half a big bottle

5 a little less than one big bottle

6 about one big bottle

7 from one to two big bottles

- 8 more than two big bottles
- 9 I do not drink wine

7. How often do you usually drink spirits?

- 1 daily
- 2 a few times a week
- 3 about once a week
- 4 a few times a month
- 5 about once a month
- 6 about once in a few months
- 7 3 - 4 times a year
- 8 twice a year
- 9 once a year or more seldom
- 0 never

8. How much do you usually drink spirits at a time?

- 1 less than one restaurant measure (less than 4 cl)
- 2 one restaurant measure (about 4 cl)
- 3 two restaurant measures (about 8 cl)
- 4 3 - 4 restaurant measures
- 5 5 - 6 restaurant measures (about quarter liter)
- 6 7 - 10 restaurant measures
- 7 about a half liter bottle
- 6 more than a half liter bottle
- 7 I do not drink spirits

9. How often have you during the last 12 months had so much beer, wine or spirits that you have felt intoxicated?

- 1 a few times a week or more often
- 2 about once a week
- 3 a few times a month
- 4 about once a month
- 5 about once in two months
- 6 4 - 5 times a year
- 7 2 - 3 times a year
- 8 once a year
- 9 not even once

The changes in questions from 1987 to 1992

The questions 1, 6 and 7 have with changes in text. We have highlighted the removed text with strikethrough font (e.g. ~~removed~~) and added text with italic font (e.g. *added*). The changes are in comparison with the previous survey.

1. Do you use any alcoholic drinks, even occasionally (f. ex. beer, wine or spirits)?

- 1 ~~yes~~ *yes, at least once a month*
- 2 ~~no, but I have not quitted completely~~ *yes, less than once a month*
- 3 no, because I quit using alcohol years ago
- 4 I have never used alcohol

~~If you have quitted alcohol use, please specify, why did you~~

		no	yes
quit?	For health reasons	1	2
	For economic reasons	1	2
	For other reasons	1	2

6. How much do you usually drink wine at a time?

- 1 half a glass
- 2 one glass
- 3 two glasses
- 4 ~~about one small bottle~~ *about half a big bottle*
- 5 a little less than one big bottle
- 6 about one big bottle
- 7 from one to two big bottles
- 8 more than two big bottles
- 9 I do not drink wine

7. How much do you usually drink spirits at a time?

- 1 less than one restaurant measure (less than 4 cl)
- 2 one restaurant measure (about 4 cl)
- 3 two restaurant measures (~~about 8 cl~~)
- 4 3 - 4 restaurant measures
- 5 5 - 6 restaurant measures (~~about quarter liter~~)

- 6 7 - 10 restaurant measures
- 7 about a half liter bottle
- 6 more than a half liter bottle
- 7 I do not drink spirits

The changes in questions from 1992 to 1997

The questions 4 and 6 have with changes in text. We have highlighted the removed text with strikethrough font (e.g. ~~removed~~) and added text with italic font (e.g. *added*). The changes are in comparison with the previous survey.

4. How much do you usually drink beer at a time? (*1 bottle = 1/3 liters.*)

- 1 less than one bottle
- 2 1 bottle
- 3 2 bottle
- 4 3 bottles
- 5 4 - 5 bottles
- 6 6 - 9 bottles
- 7 10 - 14 bottles
- 8 15 bottles or more
- 9 I do not drink beer

6. How much do you usually drink wine at a time?

- 1 half a glass
- 2 one glass (*1 glass = c. 12 cl*)
- 3 two glasses
- 4 about half a ~~big~~ bottle (*1 bottle = 0,75 l*)
- 5 a little less than one ~~big~~ bottle
- 6 about one ~~big~~ bottle
- 7 from one to two ~~big~~ bottles
- 8 more than two ~~big~~ bottles
- 9 I do not drink wine

The changes in questions from 1997 to 2002

In 2002 the questions 3–8 have been replaced with a new question number 3.

3. **How often did you drink the following amounts in one day during the last 12 months?** Instruction: Start answering from the first row. Mark (x) the most suitable 'How often?' alternative. Then continue row at a time down in the same manner. Please mark only one alternative per row.

1 dose =	bottle (1/3 liter) beer (class III)	
	<i>or</i> a glass (12 cl) of light wine	
	<i>or</i> a glass (8 cl) of strong wine	
	<i>or</i> a glass (4 cl) of spirits or other strong liquor	
Bottle (0.33 liter) beer (class IV), Gin Long Drink or strong cider	=	1.25 doses
Large bottle (0.5 liter) beer (class III)	=	1.5 doses
Large bottle (0.5 liter) beer (class IV)	=	2 doses
Bottle (0.75 liter) wine	=	7 doses
Bottle (0.75 liter) strong wine	=	10 doses
Bottle (0.5 liter) strong alcohol (e.g. Koskenkorva)	=	12 doses

Doses per day	Never	Once a month or more seldom	2-3 times a month	About once a week	2-3 times a week	4-5 times a week	6-7 times a week
15 or more	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
13-14	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
11-12	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
9-10	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
7-8	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
5-6	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
3-4	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
1-2	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

The changes in questions from 2002 to 2007

In 2007 the new question number 4 has been updated with a small change in the instructions and a change in the categories of consumed doses per day.

4. **How often did you drink the following amounts in one day during the last 12 months?**

Instruction: Start answering from the first row. Mark (~~x~~) the most suitable 'How often?' alternative. Then continue row at a time down

in the same manner. Please mark only one alternative per row.

1 dose =	bottle (1/3 liter) beer (class III)	
	<i>or</i> a glass (12 cl) of light wine	
	<i>or</i> a glass (8 cl) of strong wine	
	<i>or</i> a glass (4 cl) of spirits or other strong liquor	
Bottle (0.33 liter) beer (class IV), Gin Long Drink or strong cider	=	1.25 doses
Large bottle (0.5 liter) beer (class III)	=	1.5 doses
Large bottle (0.5 liter) beer (class IV)	=	2 doses
Bottle (0.75 liter) wine	=	7 doses
Bottle (0.75 liter) strong wine	=	10 doses
Bottle (0.5 liter) strong alcohol (e.g. Koskenkorva)	=	12 doses

Doses per day	At least 4 times a week	2-3 times a week	About once a week	1-2 times a month	3-10 times a year	1-2 times a year	Never
18 or more	1	2	3	4	5	6	7
13-17	1	2	3	4	5	6	7
8-12	1	2	3	4	5	6	7
5-7	1	2	3	4	5	6	7
3-4	1	2	3	4	5	6	7
1-2	1	2	3	4	5	6	7