

Mobile Phone User Polling for the German Federal Election 2017

Election Forecasting Project

Moritz Hemmerlein & Alexander Sacharow

12 April, 2017

Abstract

Using online surveys to forecast election outcomes imposes severe challenges to pollsters. Non-representative samples and likely-voter bias skew gathered information and require adequate statistical adjustment. This paper proposes a research design to employ different methods on data from mobile phone app users surveyed on their vote intention in the upcoming German federal election.

Contents

1	Introduction	2
2	Related Literature	3
3	Data and Potential Biases	4
3.1	Europulse Survey	4
3.2	Potential Sources of Bias	4
3.3	Data for Post-stratification and Weighting	7
4	Methodology	7
4.1	Approach number one	9
4.2	Approach number two	9
5	Data Overview	9
6	Results	10
7	Conclusion	10
8	References	10

1 Introduction

The digitalization is challenging the way polling was done for many decades. Calling people on their landline phones has become difficult as response rates dropped and households equipped with landline phones are getting less and less (Skibba 2016). In response, polling institutions have resorted to other methods for polling ranging from face-to-face interviews to mobile phone calling. However, these methods either face similar difficulties to ensure representativeness or are too expensive for regular polling. In order to tackle these obstacles, pollsters are increasingly using online polls, which are cheap and fast, but can be highly non-representative (e.g. W. Wang et al. 2015).

The eventual aim of pollsters in online as well as traditional polling is to collect sample data that reflects the view of a population of interest. The major difference between both methods is that for various reasons online polls cannot ensure representativeness before the actual poll takes place. For instance, respondents of online surveys are more likely to be from certain demographic groups or share a particular political background depending on the website or app where the survey is conducted. However, such non-representative polls can be statistically adjusted to match the demographic composition of the population.

Additionally, online election polls, like traditional ones, face another problem. Election forecasters are naturally not only interested in the population as such, but in the population of actual voters. By the time a poll is made representative in demographic terms, it is still in question whether it reflects the group of people who actually cast their ballots. This, however, is crucial in order to make an accurate prediction. Traditional polling tries to account for this using likely voter models and could perform fairly well (“Understanding Gallup’s Likely Voter Models” 2010, Keeter, Igielnik, and Weisel (2016)). Online surveys will also have to be adjusted to actual voting population in order to provide accurate predictions.

In this paper, we want to propose a research design to explore and analyse how different approaches to adjust online polls perform. Our first approach will be a two step procedure where the polling data is first made representative of the population and then likely voter methods are used to resemble the probable voters population. Our second approach attempts to combine both steps into one by adjusting the online polls directly to exit polls and voting statistics from previous elections. For this, we work with individual level data from mobile-phone app users who were surveyed on their vote intention in the German federal election 2017. As the ultimate election will only be after the end of research, we use other forecasts of the 2017 federal election as a benchmark.

The structure of the paper is as follows: Section two will survey the literature on non-representative polls and approaches to employ such data to forecast elections. Subsequently, we present our data and discuss possible problems with it. Afterwards, we discuss the methodology we want to use for adjusting our data at hand.

2 Related Literature

Traditional polling and in particular election polling has relied heavily on telephone surveys for the last decades. To ensure representativeness, the standard was randomized digit dialing (RDD). The selection of random respondents was intended to eliminate the sample bias of the survey. However, for several reasons this approach has become unreliable. First, response rates have declined heavily (Keeter et al. 2006; Holbrook, Krosnick, and Pfent 2007). The Pew Research Center (2012) reported that in the U.S. response rates dropped down to 9% in 2012, compared to 36% in 1997. This is fostered by technical changes that, for instance, make it possible to identify the caller before taking the call and increased the likelihood of people not answering survey requests. Second, more and more people do not get landlines telephones after moving to new places or just give them up as mobile phone and other means of communication have increasingly become popular. As a result, random polls are often exposed to non-response bias mitigating the probabilistic approach of RDD. Hence, classical representative polling is becoming less reliable, a trend that will rather continue than cease. Unsurprisingly, lacking representativeness of surveys has been identified as a core reason for recent election polling failures, e.g. in the UK General elections 2015 (Mellon and Prosser 2015).

Recently, non-representative polling has been proposed as a solution to these problems. Since the famously failure of the Literature Digest poll in the 1936 U.S. presidential election, pollsters have been sceptical of non-representative polling (Squire 1988, Goel, Obeng, and Rothschild (2017)) and this scepticism is still widespread. Yeager et al. (2011), for example, argue that phone surveys are still more accurate than online polls. However, they base their argument on simple correction approaches for non-probability sampled online polls.¹

W. Wang et al. (2015) in contrast are much more optimistic about the possibilities of non-representative polling. They used polling results from Xbox users, which were highly unrepresentative of the population, to forecast the 2012 U.S. presidential elections. By employing a sophisticated multi-stage approach to post-stratify and calibrate the data they were able to generate accurate forecasts of the elections. However, their methodology is mainly suited for large data sets which have sufficient observations for the combined sub-groups (strata) in the sample. Goel, Obeng, and Rothschild (2017) showed that smaller non-representative online surveys can also be accurate. They conducted polls on [Amazon Mechanical Turk](#) and a mobile phone app and achieved a level of accuracy sufficient for most practical applications. These works shows that online survey can be used in a meaningful way, if appropriate methods are used to stratify and calibrate the non-representative data. Still, how such methods perform for the case of election forecasting remains contested.

¹Goel, Obeng, and Rothschild (2017 Footnote 2) where able to get more accurate results using the basic approach, but in a different way.

3 Data and Potential Biases

3.1 Europulse Survey

In this paper, we are using data from [Dalia Research](#), an online polling firm which is conducting market and opinion research exclusively through smartphones. To ensure to collect data from a broad variety of target populations, Dalia is using a diverse set of app and website categories such as sports, news, entertainment or games. To control participants answer the survey seriously, an algorithm analyses the consistency and the response behaviour and computes a “trust score” to every respondent. Dalia praises its methodology as distinctively accounting for potential biases such as interviewer effect, social desirability bias or interviewer data entry errors. (Dalia Research 2016)

Our forecasting project utilizes data of Dalia Research’s Europulse Survey which is conducted quarterly in all EU countries. The survey consists of seven waves, but for this project we only use two waves of the survey from December 2016 and March 2017. The first wave is freely available on [Kaggle](#), the second wave was provided to us directly by Dalia Research.² Each wave consists of about 11000 individuals, of which roughly 1900 were from Germany which is the fraction of respondents we will focus on in the following analysis. The data is already pre-stratified by Dalia Research based on micro census data for age and gender.

The Europulse data is not particularly collected for election forecasting purposes but contains data on a variety of questions such as online behaviour, media consumption and personal views on political and societal development in the European Union and the respondent’s country of origin. Moreover, the survey contains information on the respondent’s personal background, demographic data and his or her financial situation. These data can be utilized in order to improve the representativeness of the survey through weighting. This will be explained more detailed in the methodology section.

The election-related variables collected in the Europulse survey are similar to traditional vote intention polling questions. First of all, the respondents are asked if and for which party they will vote in the upcoming election and for which party they voted in the previous election. Moreover, they are asked to rank political parties and describe the degree of certainty to cast a ballot for a particular party.

3.2 Potential Sources of Bias

As with all surveys, the methodology of Dalia Research has several sources of potential biases. In the case of the Europulse survey they are primarily from flawed measurement and representation

²The data for the first wave was collected between 5th and 15th of December 2016. The second wave was conducted between 13th to 27th of March 2017.

of the surveyed population, as listed in figure 1 (Groves et al. 2009).

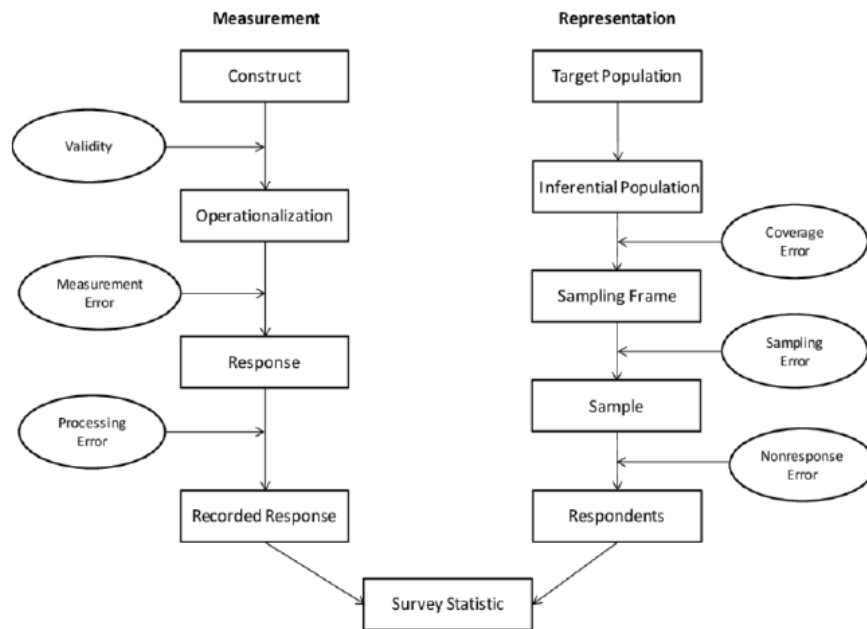
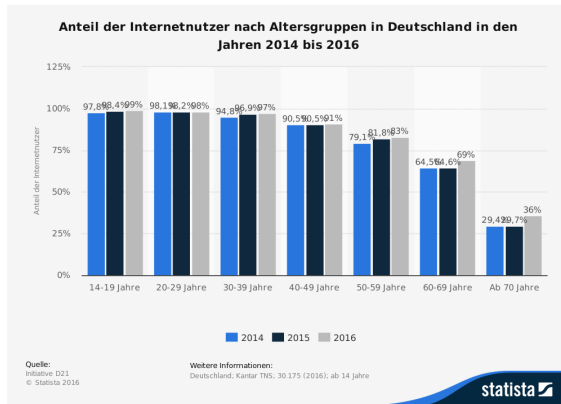


Figure 1: Potential sources of survey error

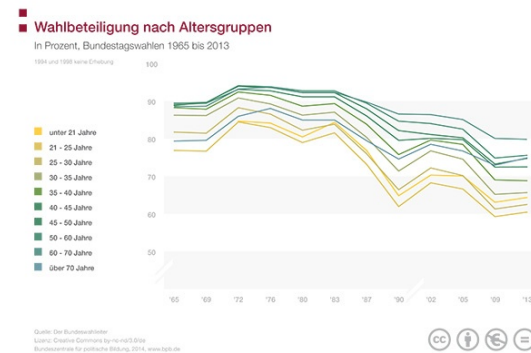
With regard to measuring vote intention to vote for a particular party, the Europulse survey uses a similar approach as traditional surveys: it asks respondents directly which party they intent to vote for at the upcoming election. Whether such questions measure correctly the actual voting behaviour at the election day is questionable but the approach does not differ from other polling methods with regards to validity. However, online surveys, such as Europulse, can reasonably claim to avoid some sources of measurement error such as social desirability bias or interviewer bias. Since such surveys are often anonymous, the social pressure on the respondent is presumably neglectable, and interaction with the interviewer does not bias the response.

Regarding representativeness, the Europulse survey has some limitations in comparison to face-to-face or RDD interviews. First of all, Europulse' framework does not select participants at random but offers visitors of certain websites or app-users the opportunity to participate in the survey. Hence, this approach carries the risk of self-selection. Moreover, representativeness would require that the Europulse survey in principle should be available to the entire population. However, not everyone is using smart phones and even if they are, even less individuals used the applications Dalia uses for its surveys. The users Dalia Research actually tends to reach are likely to be much younger and technology oriented than the general population would be (see figure). This is in particular a problem for election polling, as older voters are systematically underrepresented in such online methodologies.

Dalia Research tries to account for these problems. First, they try to increase the diversity of their



(a) Internet usage by age



(b) Voter turnout by age

Figure 2: Internet usage and turnout by age

respondents by presenting their surveys on various apps and platforms targeting different user groups. Second, they pre- and post-stratify their data. Pre-stratification is done on the basis of age and gender, using self-reported demographic information. For each age and gender strata they target a certain number of respondents so that their sample resembles roughly the German population. In a second step, they use data from the German Census and compute weights for combined cluster of age, gender, education and whether the respondent lives in an urban or rural area. These weights can finally be used to post-stratify the sample to match the demographic composition of the population more closely. Third, the selection bias is slightly decreased by the fact that respondents are not informed that they are answering an election related survey before they actually start it and in general Dalia Research has high completion rates, hence participants are only rarely dropping out after they started a survey.

Despite these efforts it is questionable whether the data can be regarded as representative. Printing a simple frequency table of the vote intent at the next general election and using self-reported voting intent to estimate the turnout shows results significantly diverging from other forecasts (see www.wahlumfragen.org).

The large divergence from can be due to a variety of flaws that bias Europulse. First of all, even if the sample represents the German population in terms of age groups and gender, it can be question whether the old female or male users are truly representative of their age group. Moreover, not revealing the survey content might sort out political motivated respondents, but interest in political matters is likely correlated to general willingness to respond to surveys. Hence, in order to utilize the data to election forecasting further weighting and post-stratification will be necessary.

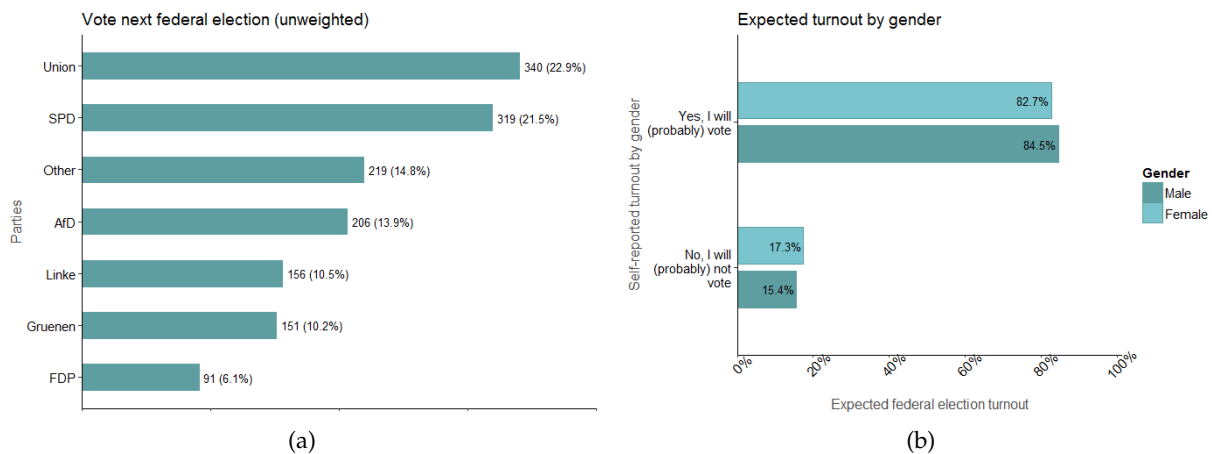


Figure 3: Vote intent and self-reported turnout (Europulse December 2016)

3.3 Data for Post-stratification and Weighting

The data used for post-stratification and weighting come from several sources. To adjust our sample to demographic composition of the German population we use data of the German [Zensus 2011](#). The data is freely available and we obtained combined frequencies for characteristics such as age, gender, education, employment status and confession. The Census claims to reflect the actual demographic distribution of the German population and hence is suitable in order to post-stratify our sample to demographic criteria.

Furthermore, we collected data from election polls conducted by Forschungsgruppe Wahlen e.V., Infratest dimap and from the official German election statistics (Neu 2013). While the latter contains only information on turnout and the demographic dimensions age and gender, Forschungsgruppe Wahlen e.V. has also issued election results and turnout along groups with different education, different employment status and confession. Such data is the closest estimate of the actual voting population and their voting behaviour across demographic groups and across the spectrum of political parties we can get.

Finally, in order to set our forecasts into context we will use polling data from other institutes. For this we plan to webscrape the respective informations from [wahlrecht.de](#).

4 Methodology

We want to employ two approaches generate election forecasts from the data we have. In the first approach the survey data is adjusted in a two step procedure: First, the representativeness of survey data is increased by finding appropriate weights. Then a likely voting model is used to get to the election forecast. Our second approach does this in one step: The survey data will be

directly adjusted to the likely voter population based on exit poll data from the last German federal election.

In order to follow both procedures, we need methods to compute weights for sub-groups of the survey sample aiming at increased representativeness. A classical way to get these weights is *raking*. With raking weights are assigned to each respondent in order to match the marginal distribution of characteristics in the ‘true’ population. For example, if we know the distribution of education level and employment status in the population, for each we can compute a weight so that the weighted survey sample has the same distribution of these characteristics as the ‘true’ population. Basically, raking computes a joint distributions for each combination of characteristics.

In formal terms, we can describe each combination of an individual and a characteristic with $x_{i,j} \in \{0, 1\}$ where i stands for the individual and j for the characteristic. To illustrate, if $x_{i,j} = 1$ an individual may be female and if $x_{i,j} = 0$ an individual is not female. Hence, each characteristic is modeled with a binary variable. c_j expresses the prevalence of a characteristic in the population and raking estimates the weights w_i such that:

$$c_j = \frac{\sum_{i=1}^n w_i x_{i,j}}{\sum_{i=1}^n w_i} \quad \forall j \quad (1)$$

The weights can then be used to compute the raking estimates for each observation in the survey y :

$$\hat{y}^{rake} = \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i} \quad (2)$$

However, raking is not ideal. It is often used if only marginal probabilities are available. This is also the case for some of our data. However, if the joint distribution of characteristics is known we can use *post-stratification* instead. In contrast to raking, post stratification takes into account the distribution of combination of characteristics in the population. In other words: It sub-divides the population into strata for each combination. Unfortunately, publically available data often only includes the distribution of pairs of characteristics or in rare cases triples (e.g. age, gender and past vote). Ideally, we would like to know the distribution of the population for each strata (combination of the characteristics).

Formally the post-stratification estimator can be expressed as follows (Goel, Obeng, and Rothschild 2017):

$$y^{post} = \frac{\sum_{j=1}^J N_j \hat{y}_j}{\sum_{j=1}^J N_j} \quad (3)$$

where \hat{y}^{post} is the estimator of y in the strata j and N_j the size of the j -th strata in the population. But even if the size of N_j is known for each strata, the number of strata’s grows exponentially with each

characteristic included. If we use for instance two gender categories and four age groups we have eight stratas, if we add past vote (7 categories) we have immediately 54. As a result we might only have a few individuals for each strata or even none. Hence, a few respondents in the strata female, old (60+) and FDP voter in 2013 will have an overproportional influence on the strata estimate.

If possible, we might use *model-based post-stratification* to counter this effect. In this approach the estimates for each strata are not based on the average in the strata, but the result of a multinomial logistic regression. In order to arrive at this regression results, demographic variables in the sample can be used. In order to execute this we will orient ourself at the work of Goel, Obeng, and Rothschild (2017).³ The larger question here is, whether it makes any sense at all to use model-based post-stratification if we only have the strata size for triples at best.

4.1 Approach number one

For stratifying the data we orient ourself at the work of XXX. The basic idea is to compute clusters of voters along several demographic categories and use their past votes to compute weights.

First, if possible we use post-stratification (see Lumpley, ch. 7?) to compute weights for subgroups in the sample. Post-stratification tries to make a sample representative of the actual population by ensuring the relative size of subgroup resembles the relative size of the same subgroup in the 'true' population. For forecasting the 'true' population is not known, as it is a question of who will actually turn out to vote.

How we plan to make poll representative

Compare different stratification approaches

Likely problems we will encounter: 1. empty clusters or clusters with low number of observations. Implications: If empty, there is a real problem. If the number of observation is low, e.g. below 20, the weights will amplify the impact of this small group in the total forecasting result.

Benchmark -> other publically available polls. This is straight-forward, but also problematic as it might induce a herding effect. The final evaluation is only possible after the election

4.2 Approach number two

5 Data Overview

how representative our data already is

³They are developing an r package for this (postr) and might be willing to share their r script (as they already announced to make it public).

1. raw data forecast. Compared to other forecasts the data under represents the CDU as well as the SPD. (Verify)
2. Show distribution of respondents on different demographic clusters and compare to zensus / exit polls / election statistics
3. Raw (voted last election)

6 Results

What the result is of making it representative

1. Election forecast Weighted with exit polls
2. Election forecast weighted with election statics
3. Election forecast weighted with zensus

Compare the three different weighten approaches

Comment Moritz: I think we have to weight the data with Zensus data in any case at least for gender and age as long we don't want to use the Dalia weights; then we can either use election statistics or exit poll data (that we don't have)

My approach would be:

1. Weighting with Zensus (accounts for self-selection of the survey)
2. Different mixes of weighting with election statistics with education (accounts for likely voter bias)

7 Conclusion

Summary of the core finding

Further implications

8 References

Dalia Research. 2016. "Dalia Research Methodology." <https://daliaresearch.com/wp-content/uploads/2016/08/Methodology-PDF-1.pdf>.

Goel, Sharad, Adam Obeng, and David Rothschild. 2017. "Online, Opt-in Surveys: Fast and Cheap,

but Art They Accurate?" Working Paper.

Groves, Robert M., Floyd J. Fowler Jr., Mick P. Couper, James M. Lepkowski, Eleanor Singer, and Roger Tourangeau. 2009. *Survey Methodology*. 2nd ed. New Jersey: John Wiley & Sons Inc.

Holbrook, Allyson L., Jon A. Krosnick, and Alison Pfent. 2007. "The Causes and Consequences of Response Rates in Surveys by the News Media and Government Contractor Survey Research Firms." In *Advances in Telephone Survey Methodology*, edited by James M. Lepkowski, Clyde Tucker, J. Michael Brick, Edith D. de Leeuw, Lilli Japac, Paul J. Lavrakas, Michael W. Link, and Roberta L. Sangster, 499–528. Hoboken, NJ, USA: John Wiley & Sons, Inc. doi:[10.1002/9780470173404.ch23](https://doi.org/10.1002/9780470173404.ch23).

Keeter, Scott, Ruth Igielnik, and Rachel Weisel. 2016. "Can Likely Voter Models Be Improved? Evidence from the 2014 U.S. House Elections." Pew Research Center.

Keeter, Scott, Courtney Kennedy, Michael Dimock, Jonathan Best, and Peyton Craighill. 2006. "Gauging the Impact of Growing Nonresponse on Estimates from a National Rdd Telephone Survey." *The Public Opinion Quarterly* 70 (5): 759–79. <http://www.jstor.org/stable/4124225>.

Mellon, Jonathan, and Chris Prosser. 2015. "Investigating the Great British Polling Miss: Evidence from the British Election Study." *SSRN Electronic Journal*. doi:[10.2139/ssrn.2631165](https://doi.org/10.2139/ssrn.2631165).

Neu, Viola. 2013. "Bundestagswahl in Deutschland Am 22. September 2013." Konrad Adenauer Stiftung. http://www.kas.de/upload/dokumente/2013/09/Anhang_gesamt_neu.pdf.

Pew Research Center. 2012. "Assessing the Representativeness of Public Opinion Surveys." Pew Research Center. <http://www.people-press.org/files/legacy-pdf/Assessing%20the%20Representativeness%20of%20Public%20Opinion%20Surveys.pdf>.

Skibba, Ramin. 2016. "The Polling Crisis: How to Tell What People Really Think." *Nature* 538 (7625): 304–6. doi:[10.1038/538304a](https://doi.org/10.1038/538304a).

Squire, Peverill. 1988. "Why the 1936 Literary Digest Poll Failed." *The Public Opinion Quarterly* 52 (1): 125–33. <http://www.jstor.org/stable/2749114>.

"Understanding Gallup's Likely Voter Models." 2010. <http://www.gallup.com/poll/143372/understanding-gallup-likely-voter-models.aspx?version=print>.

Wang, Wei, David Rothschild, Sharad Goel, and Andrew Gelman. 2015. "Forecasting Elections with Non-Representative Polls." *International Journal of Forecasting* 31 (3): 980–91. doi:[10.1016/j.ijforecast.2014.06.001](https://doi.org/10.1016/j.ijforecast.2014.06.001).

Yeager, D. S., J. A. Krosnick, L. Chang, H. S. Javitz, M. S. Levendusky, A. Simpser, and R. Wang. 2011. "Comparing the Accuracy of Rdd Telephone Surveys and Internet Surveys Conducted with Probability and Non-Probability Samples." *Public Opinion Quarterly* 75 (4): 709–47. doi:[10.1093/poq/nfr020](https://doi.org/10.1093/poq/nfr020).