



The Leverhulme Trust

# **Simulation of adaptive survey design in a longitudinal survey**

**Peter Lundquist and Anton Johansson (Statistics Sweden)**  
**Gabriele B. Durrant (University of Southampton)**

5<sup>th</sup> Workshop on Advances in Adaptive and Responsive Survey Design,  
University of Michigan, November 6<sup>th</sup>-7<sup>th</sup>, 2017

# Introduction and background

- Large survey resources are being spent on making unproductive calls e.g. contact attempts that do not lead to an interview.
- Unstable data collection and unclear collection strategy
- We have used work by Durrant et al. (2013, 2015) who assess the prediction of nonresponse models using paradata from previous and current wave
- The ambition is to find a new data collection strategy for the Swedish Labour Force Survey
- A simulation tool to use in the fieldwork planning

# Data used in Simulation

## The Swedish Labour Force Survey (LFS)

- Longitudinal survey with 8 waves
  - “A new sample” every week
  - One data collection mode: telephone
  - Two interviewer groups: field and call-center
- 
- Data used in modeling: **LFS in January 2016**
    - Initial sample size 5,164; Week 4 sample
  - Model used on next wave in **April 2016**:
    - 7/8 of the sample is the same

# Data collection in LFS

Approximately 100 Field + 100 call center interviewers sign up for the following shifts:

Shift/Day	Mon-Thu	Friday	Saturday	Sunday
08:00-11:59				
12:00-16:59				
17:00-18:59				
19:00-21:59		Contact by agreement		

- We assume that the resources could be better allocated to the different subsamples (e.g. Week 1-Week 4 in January 2016)

# LFS – fieldwork January 2016

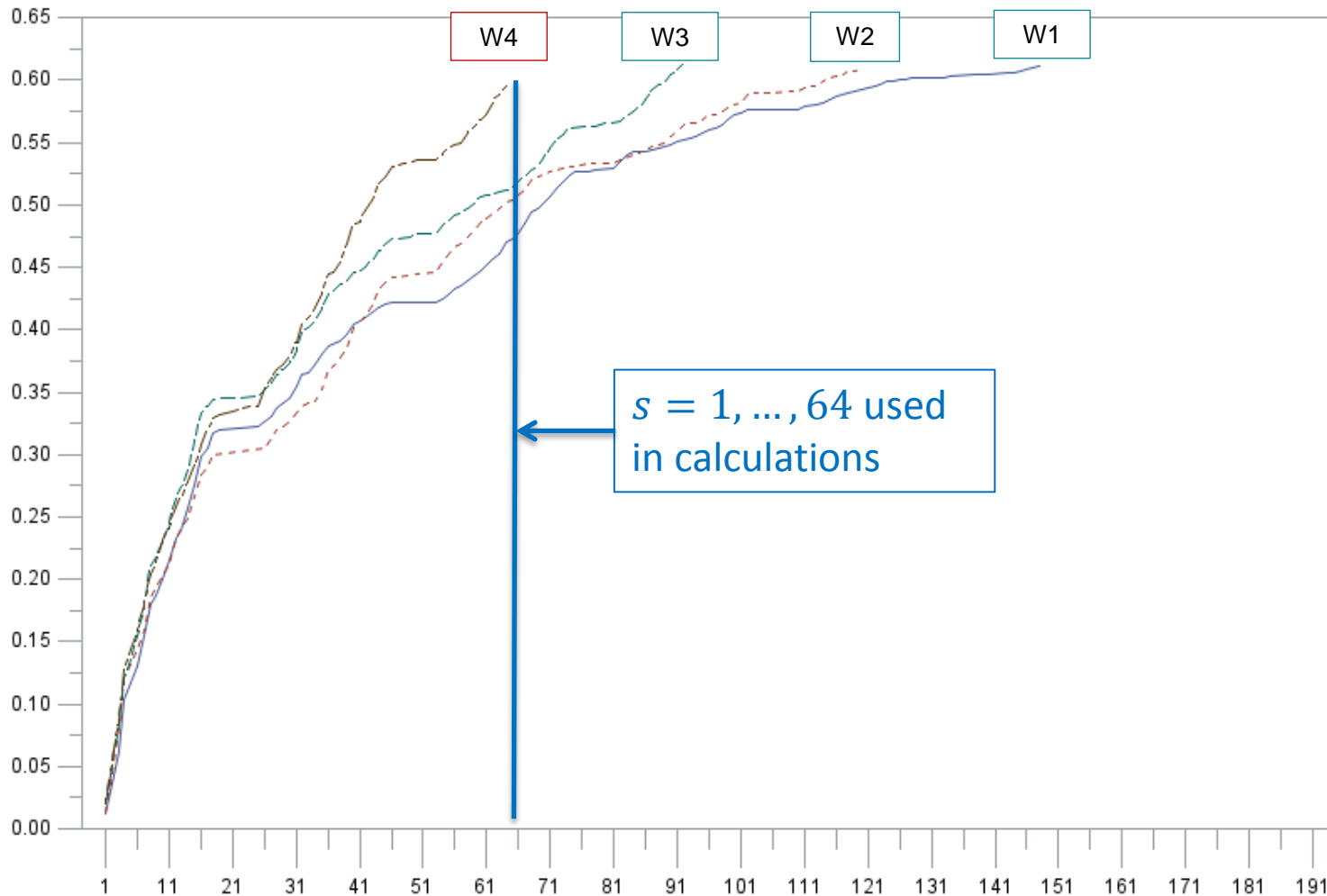
## Number of call attempts

Week	1	2	3	4	5	6 (1)+(2)
W1	(1): 17,896		(2): 12,412			30,308
W2		(1): 20,680		(2): 9,767		30,447
W3			(1): 20,813		(2): 8,219	29,032
W4				(1): 28,992		28,992
Fieldday	1	8	15	22	29	37

**(1) Primary fieldwork 16 days (equal all samples W1-W4)**

**(2) Extended fieldwork**

# Response progression over time-slots for 4 LFS-weeks



# How to decide the three phases?

Week 4 has a short data collection: 16 days

**Phase 1:** All units should have been contacted two times within the first four days.

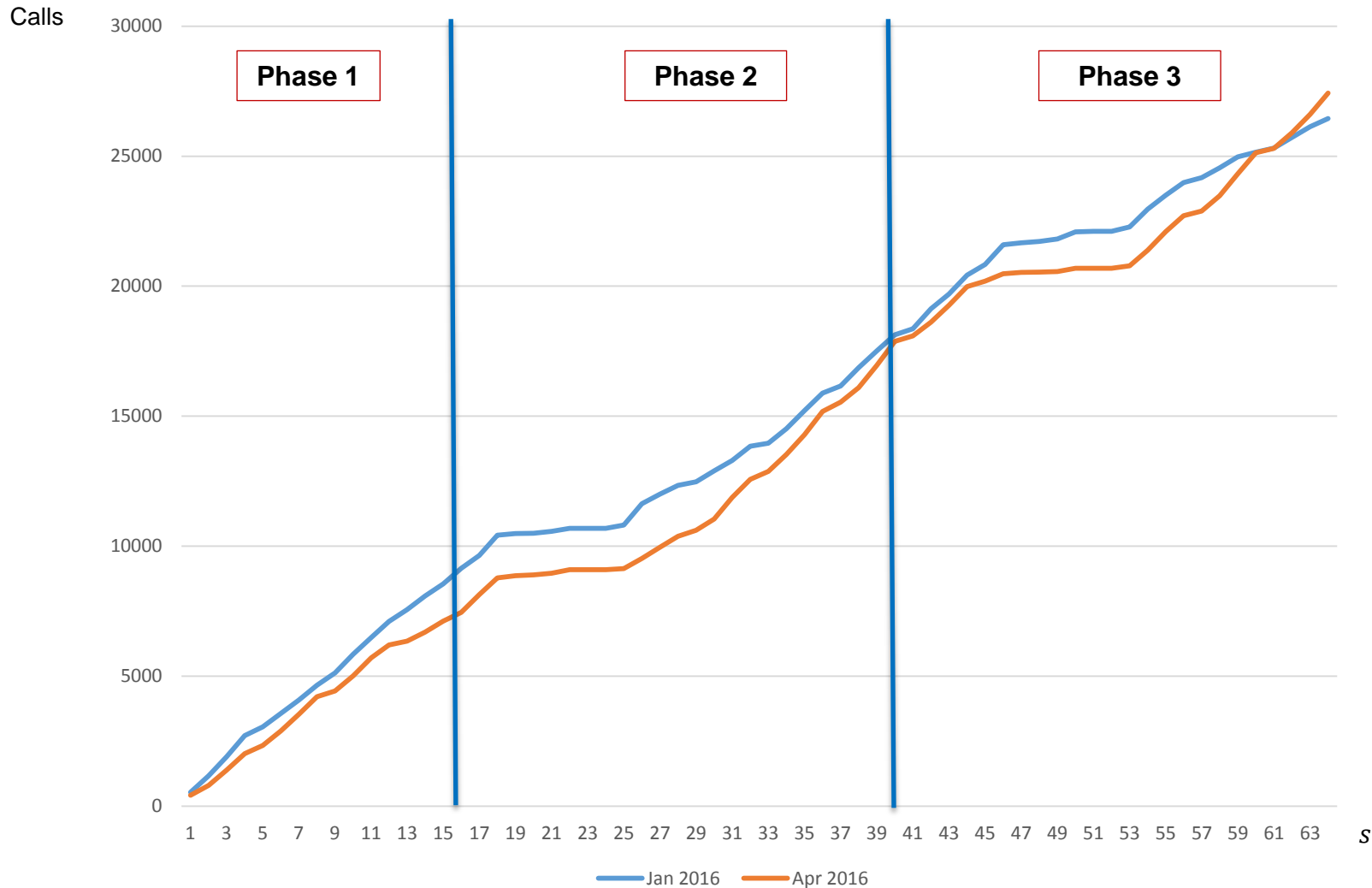
**Phase 2:** Active units should have seven calls (day 5 to 10).

**Phase 3:** Active units should have 13 calls (day 11 to 16).

- Every day is divided in four shifts ( $s = 1, 2, \dots, 64$ )

# Comparison January and April 2016

## Cumulative number of calls for week 4





# Data collection strategies

Two different data collection strategies are studied in the simulation:

1. Maximize the response rate (comparison strategy)
2. Use assigned time slots in the first phase (desired strategy)

The data collection use the same three phases, i.e. Phase 1 (day 1-4), Phase 2 (day 5-10) and Phase 3 (day 11-16).

# Assigned time slots

**The LFS-interview (in January) ends with questions about a suitable time for the next interview (in April).**

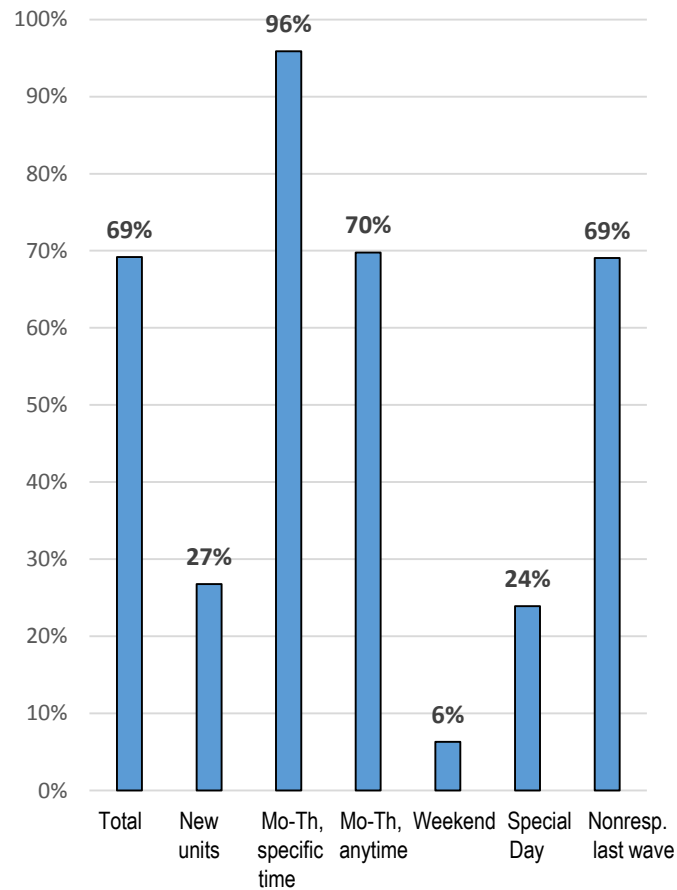
- For the respondents, it's then possible to place the first calls in appropriate time intervals (and days) in April.

The time slots are: 08:00-11:59, 12:00-16:59, 17:00-18:59 and 19:00-21:59.

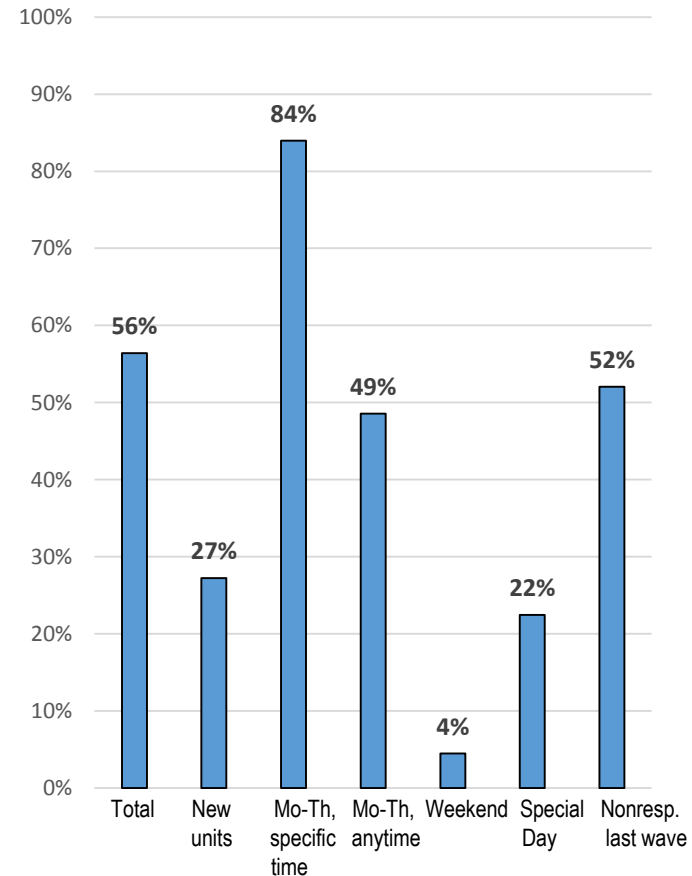
- For the nonrespondents register information, process data and information from earlier waves are used.
- For wave 1 only register information is used.

# The proportion in initially selected time slot

## 1<sup>st</sup> Call attempt



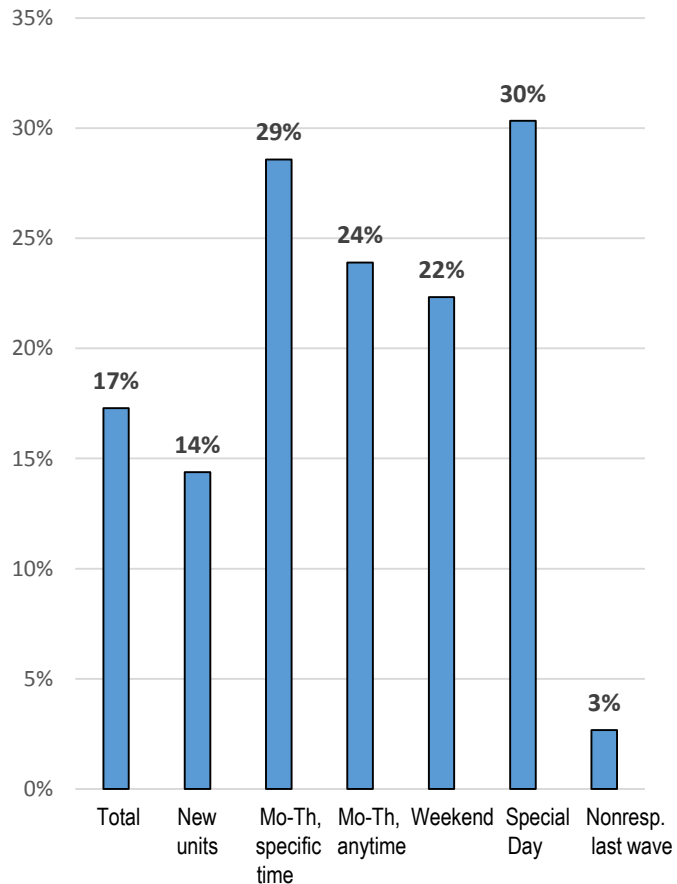
## 2<sup>nd</sup> Call attempt



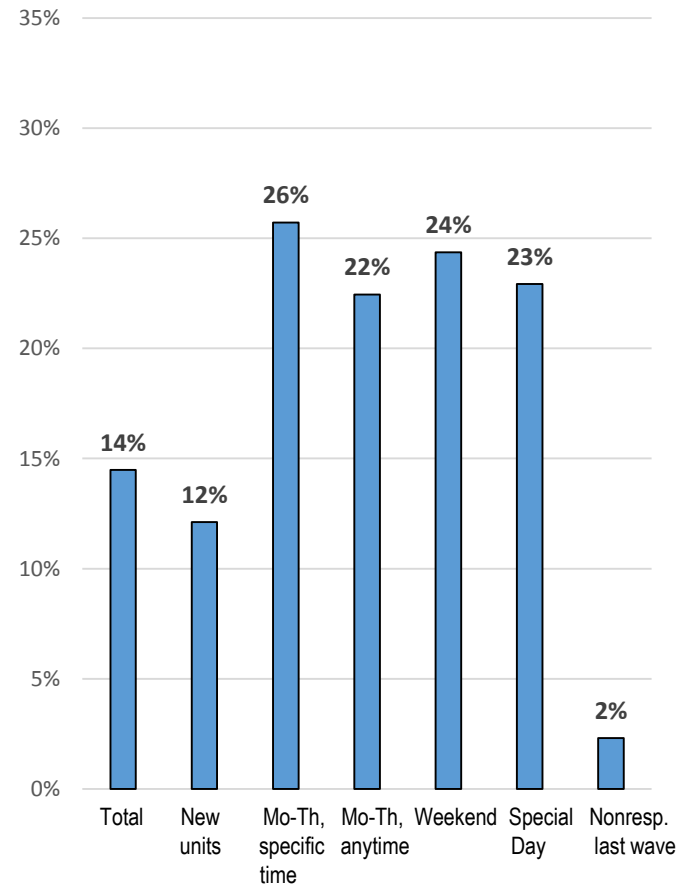
Adapted from Johansson and Olsson internal document 2017

# The proportion interviews

## 1<sup>st</sup> Call attempt



## 2<sup>nd</sup> Call attempt



# We now turn to the simulations...

- Logistic regression
- Register data and process data
- Number of calls in time slots
- Modelling the strategies
- Indicators

# Models for Phase 1, 2 and 3

## Logistic regression (binary), three different models

### Model 1

Phase 1 time slots:  $s = 1, 2, \dots, 16$  (day 1 to 4)

### Model 2

Phase 2 time slots:  $s = 17, \dots, 40$  (day 5 to 10)

### Model 2

Phase 3 time slots:  $s = 41, \dots, 64$  (day 11 to 16)

## Dependent variable:

response (interview or **not** interview)

- Objects are individuals (not households)
- Within each  $s$  is only one call attempt allowed

# Models for Phase 1, 2 and 3

## Logistic regression (binary), explanatory variables

Explanatory variables		Model 1	Model 2	Model 3
Register data	Age (16-54 or 55-74 yrs)	<b>X</b>	<b>X</b>	<b>X</b>
	Born in Sweden or not	<b>X</b>	<b>X</b>	
	Education (high) or not	<b>X</b>	<b>X</b>	<b>X</b>
	Married or not	<b>X</b>	<b>X</b>	<b>X</b>
	Employed or not	<b>X</b>		<b>X</b>
Process data	More than 2 calls	<b>X</b>		
	1 <sup>st</sup> wave	<b>X</b>	<b>X</b>	<b>X</b>
	2 <sup>nd</sup> -8 <sup>th</sup> wave and interview last wave	<b>X</b>	<b>X</b>	<b>X</b>
	2 <sup>nd</sup> -8 <sup>th</sup> wave and <b>no</b> interview last wave	<b>X</b>	<b>X</b>	<b>X</b>
	2 <sup>nd</sup> -8 <sup>th</sup> wave and LONG -1 (more than 6 call attempts or not <b>last wave</b> )		<b>X</b>	<b>X</b>
	LONG (more than 6 call attempts or not)		<b>X</b>	<b>X</b>
	Day shift (8-12, 12-17, 17-19, 19-22)	<b>X</b>	<b>X</b>	<b>X</b>
	Time slot (1, 2, ..., 64)	<b>X</b>	<b>X</b>	<b>X</b>

# Models January 2016

## Logistic regression, response as dependent variable

	Pseudo R <sup>2</sup>	ROC
Phase 1: Model 1	0.13	0.705
Phase 2: Model 2	0.15	0.750
Phase 3: Model 3	0.13	0.757

**Pseudo-R2 Statistic** : the proportion of variation in the dependent variable that is explained by the model.

**ROC (Receiver Operating Curve)**: the greater area under the curve (AUC) the greater predictive power. AUC values range from 0.5 (no discrimination) to 1 (perfect discrimination).



# Predicted response propensities

For unit  $k \in \text{sample}$  in phase  $j$  ( $j = 1, 2, 3$ ) and  $s = 1, 2, \dots, 64$

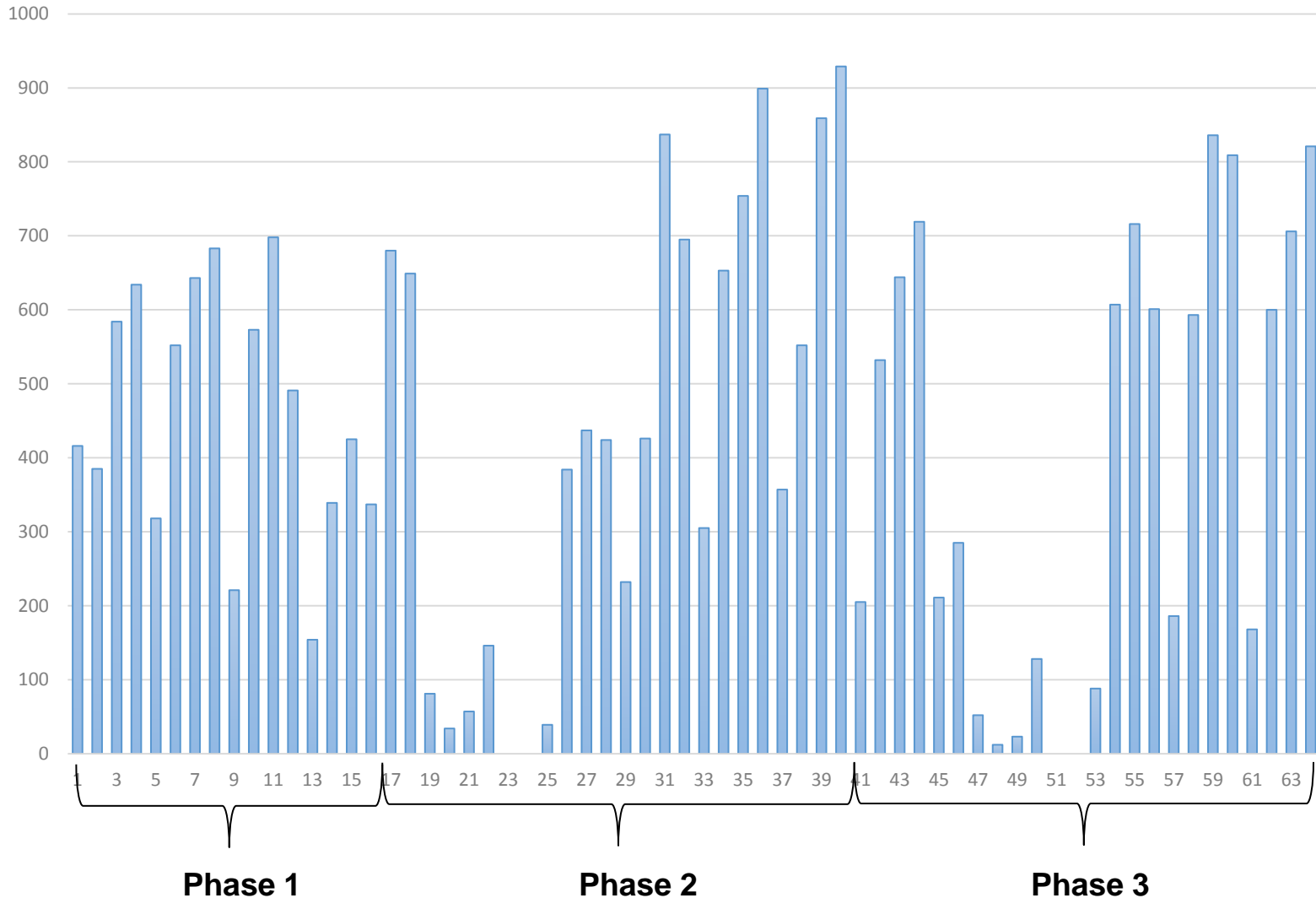
$$\hat{p}_k = \frac{\exp(\hat{\boldsymbol{\beta}}_j \mathbf{x}_{k,j(s)})}{1 + \exp(\hat{\boldsymbol{\beta}}_j \mathbf{x}_{k,j(s)})}$$

are used in the simulations.

# How to decide the workload?

- In each time slot is a fixed number of calls possible.
- We will in the modelling use the **workload** for week 4 in January, but in the simulation then **workload** for week 4 in April
- It is however clear that this will not always work.
- In Choudhry et al. (2011) the **workload** was optimized to minimize the data collection costs

# Strategy 1 and 2: Number of calls April 2016



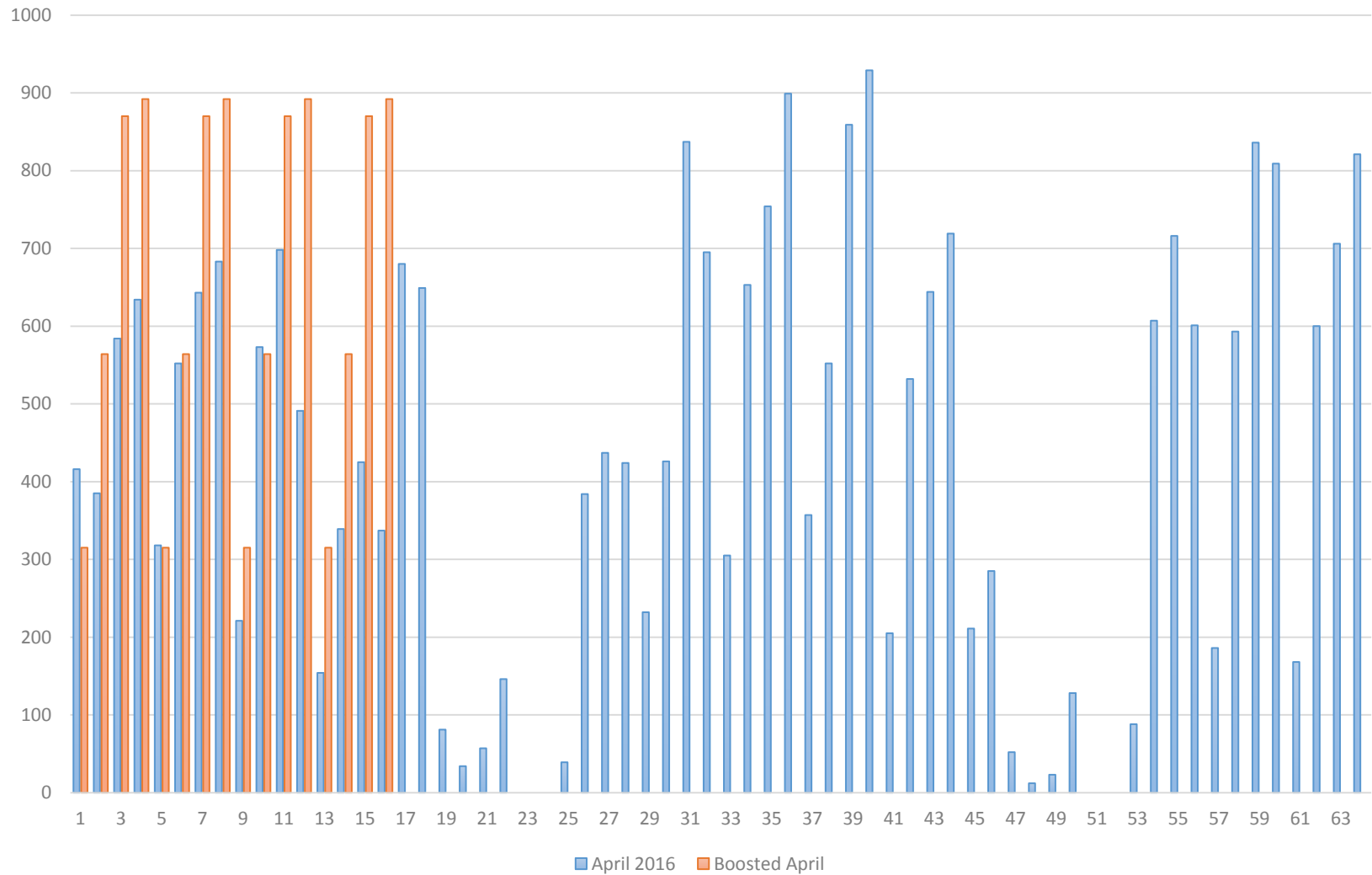
# Data collection: Strategy 1

Phase	Maximize the response rate (comparison strategy)
1) Day 1-4 s=1-16 Model 1	<ul style="list-style-type: none"><li>• For each <math>s</math> the <math>\hat{p}_k</math> objects in treatment are sorted (descending)</li><li>• A maximum of 3 calls</li></ul>
2) Day 5-10 s=17-40 Model 2	<ul style="list-style-type: none"><li>• For each <math>s</math> the <math>\hat{p}_k</math> objects in treatment are sorted (descending)</li><li>• A maximum of 7 calls</li><li>• After 4 calls: stop individuals in wave 2-8 that refused to participate last wave</li></ul>
3) Day 11-16 s =41-64 Model 3	<ul style="list-style-type: none"><li>• For each <math>s</math> the <math>\hat{p}_k</math> objects in treatment are sorted (descending)</li><li>• A maximum of 13 calls</li></ul>

# Data collection: Strategy 2

Phase	Use time slots (desired strategy)
1) Day 1-4 s=1-16 Model 1	<ul style="list-style-type: none"><li>• For each s the <math>\hat{p}_k</math> objects in treatment are sorted (randomly)</li><li>• Use time slots</li><li>• A maximum of 3 calls</li></ul>
2) Day 5-10 s=17-40 Model 2	<ul style="list-style-type: none"><li>• For each s the <math>\hat{p}_k</math> objects in treatment are sorted (descending)</li><li>• A maximum of 7 calls</li></ul>
3) Day 11-16 s =41-64 Model 3	<ul style="list-style-type: none"><li>• For each s the <math>\hat{p}_k</math> objects in treatment are sorted (descending)</li><li>• A maximum of 13 calls</li><li>• Select easy objects still in treatment and wave 1-2</li></ul>

# Strategy 2B: Boosted number of calls in phase 1



# Evaluation of the new strategy

$P$  = the weighted response rate in per cent

$IMB$  = the imbalance measure measures the difference between the response set  $r$  and the selected sample  $s$  for a chosen  $\mathbf{x}$ -vector.

It could be demonstrated\* that  $IMB$  is equal to the variance for the response propensities for the chosen  $\mathbf{x}$ -vector

$$CV_s = \frac{\sqrt{IMB}}{P}$$

**\*Särndal & Lundquist 2014**

**Note: auxiliary variables (register data) depends on available variables and the indicators depends on the sample.**

# Simulation results, averages and standard errors based on 1,000 simulations

LFS-April 2016	$P$	$IMB$	$CV_s$	Time (min)	Number of calls	
				<i>Per case</i>	<i>Total</i>	<i>Per case</i>
<b>Strategy 1</b>	64.8 (0.5)	1.67 (0.11)	19.9 (0.7)	17.8	27,424	5.19
<b>Strategy 2</b>	61.2 (0.5)	1.43 (0.10)	19.5 (0.7)	16.5	25,166	4.76
<b>Strategy 2B*</b>	62.6 (0.5)	1.37 (0.10)	18.7 (0.7)	17.2	26,573	5.03

The response rate  $P$  is weighted in percent,  $IMB$ ,  $CV_s$ , are multiplied with 100.

\*) Increased number of calls in phase 1

**x**-vector used in computations: Age, High Education, Origin, Civil, Gender  
(3) (2) (2) (2) (2)



# Discussion

- **Work with the logistic regression models:**
  - **Further elaboration with factors**
  - **Bayesian models** were not necessary in this situation, however in other kind of simulations this could be of interest (as in a real data collection)
  - **Development of the tool:** Indicators for the different phases and maximum interviewer resources
- **The simulation tool:**
  - Makes it possible to find better strategies
  - Generates new ideas how to develop and control the data collection planning
- **Experiments** -This possibility should be noted!
- **BADEN:** use tool by Schouten et al. 2017.

Thank you!

[peter.lundquist@scb.se](mailto:peter.lundquist@scb.se)  
[anton.johansson@scb.se](mailto:anton.johansson@scb.se)

# Simulation on April 2016

## EXAMPLE Strategy 1

- The response propensities are assumed to be  $Be(\hat{p}_k)$  distributed in the logistic regression models.
- For time slot 1: 5,282 random selections were made from Model 1. The randomization corresponds to the outcome if all the individuals in the sample were contacted for  $s = 1$ . The  $n_1$  highest response propensities are inspected and those who "respond" are set aside.
- The "nonresponders" continue to the 2<sup>nd</sup> time slot,  $s = 2$ . The  $n_2$  highest response propensities are inspected and those who "respond" are set aside...
- The procedure continues until time slot 64, where the "data collection" ends.

**The data collection is replicated 1,000 times for the described strategies.**



# References

Choudhry, Hidirolou & Laflamme (2011). "Optimizing CATI workload to minimize data collection cost." Proceedings of the Survey Research Methods Section, 1904-1913, ASA.

Durrant, D'Arrigo & Müller (2013). "Modeling Call Record Data: Examples from Cross-Sectional and Longitudinal Surveys," in Improving Surveys with Paradata: Analytic Use of Process Information, ed. Kreuter, F., 281–308, Hoboken, NJ: John Wiley and Sons.

Durrant, Maslovskaya & Smith (2015). "Modeling final outcome and length of call sequence to improve efficiency in interviewer call scheduling." Journal of Survey Statistics and Methodology, 3, 397–424.

Johansson, Lundquist & Durrant (2016). "Stopping rules in a longitudinal survey – impact on cost and survey quality." Presented at AAPOR.

Särndal & Lundquist (2014). "Accuracy in estimation with nonresponse: A function of degree of imbalance and degree of explanation." Journal of Survey Statistics and Methodology, 2, 361-387.

Schouten, Mushkudiani et al. (2017). "A Bayesian analysis of design parameters in survey data collection." Paper submitted.