ASPIRE: A System for Product Improvement, Review, and Evaluation

Paul Biemer RTI International and University of North Carolina at Chapel Hill



Co-developers

Heather Bergdahl, Lilli Japec, Åke Pettersson Statistics Sweden and Dennis Trewin Australian Statistician

What is ASPIRE?

- General system evaluating the quality of a statistical product
- Currently focuses on Accuracy, but has been also been applied to other quality dimensions
- Based upon a TSE decomposition of product quality
- Assigns quality ratings to products based upon
 - The intrinsic risk of each error source to quality
 - Risk mitigation activities by error source under five criteria
 - A weighted average of ratings over TSE components
- Presents visual summaries of the results accessible by upper management

TSE Decomposition

NOTATION:

- \hat{Y}_{P} = preliminary published estimate
- \hat{Y} = final published estimate
- Y = error-free parameter under the operational specification
- X = true parameter or preferred specification

$$\hat{Y}_{P} - X = (\hat{Y}_{P} - \hat{Y}) + (\hat{Y} - Y) + (Y - X)$$

total error = (revision error)
+ (sampling and other nonsampling errors)
+ specification error

Sampling and Other Nonsampling Errors

$$\hat{Y} - Y = \mathcal{E}_1 + \dots + \mathcal{E}_6$$

where the ε 's correspond to these six error sources:

- 1. Sampling error
- 2. Frame error
- 3. Nonresponse error
- 4. Measurement error
- 5. Data processing error
- 6. Model/estimation error

Products Reviewed

Survey Products

- Foreign Trade of Goods Survey (FTG)
- Labour Force Survey (LFS)
- Annual Municipal Accounts (RS)
- Structural Business Survey (SBS)
- Living Conditions Survey (LCS)
- Consumer Price Index (CPI)

Error Sources

- Revision error Specification error Sampling error Other nonsampling errors
- Frame error
- Nonresponse error
- Measurement error
- Data processing error
- Model/estimation error

Products to be Reviewed (cont'd)

Registers	Error Sources
Business Register (BR) Total Population Register (TPR)	Specification error Frame: Overcoverage Undercoverage Duplication Missing Data Content Error
Compilations	Error Sources
Gross Domestic Product (GDP)Annual GDP (production)Quarterly GDP (production)	Customized error profile

Process for Estimating GDP by Current and Constant Price Approaches



Published Constant GDP Estimate

8

is one such example.

*NOTE: Some items follow the deflation process in the opposite direction and are complied starting with information on volume change from the previous year. The volume estimate is then reflated with the price index in order to come to the current price estimate. Items within the Energy sector

GDP Error Sources

- Input data error (up to four sources)
 - Subject to all sampling error as well as the all relevant nonsampling error sources
- Compilation error
- Data Processing Error
- Model/Estimation Error
- Deflation/Reflation Error
- Balancing Error
- Revision Error

Effects of Input Data Errors on GDP Accuracy

The GDP input data sources give rise to a set of variables:

 x_1, x_2, \ldots, x_p

These input variables are used in the calculation of GDP and are subject to error; i.e.,

 $x_p = x'_p + \varepsilon_p$

where χ'_p is the true value and ${\cal E}_p$ is the error.

The true GDP has the form $GDP' = g(x'_1, x'_2, \dots, x'_P)$

Therefore, the estimate of GDP can be written as

$$GDP = f(x_1, x_2, \dots, x_P) = GDP' + e \leftarrow$$

The ASPIRE project is particularly interested in how the errors, $\mathcal{E}_1, \mathcal{E}_2, \dots, \mathcal{E}_P$ affect the magnitude of the error *e*.

Input 1

Input 2

Input K

Quality Criteria were Applied to Each Error Source

Intrinsic Risks to Data Quality by Error Source High, Medium, Low

Criteria by Error Source

- 1. Knowledge of risks
- 2. Communication with users
- Compliance with standards and best practices
- 4. Available expertise
- Achievement toward risks mitigation and/or improvement plans

Ratings by Criterion Poor () Fair () Good () Very Good () Excellent ()

An Example of the Rating Guidelines – Knowledge of Risks

Poor 🔴	Fair 🖰	Good 🔿	Very Good 🝚	Excellent 🔾
Internal	Internal	Some work has	Studies have	There is an ongoing program of
program	program	been done to	estimated relevant	research to evaluate all the
documentation	documentatio	assess the	bias and variance	relevant MSE components
does not	n	potential	components	associated with the error source
acknowledge	acknowledges	impact of the	associated with	and their implications for data
the source of	error source	error source on	the error source	analysis. The program is well-
error as a	as a potential	data quality.	and are well-	designed and appropriately
potential factor	factor in data		documented.	focused, and provides the
for product	quality.			information required to address
accuracy.				the risks from this error source.
The second second	But: No or	But:	But: Studies have	
	very little	Evaluations	not explored the	
	work has	have only	implications of the	
	been done to	considered	errors on various	
	assess these	proxy measures	types of data	
	risks	(example, error	analysis including	
		rates) of the	subgroup, trend,	
		impact with no	and multivariate	
		evaluations of	analyses	
		MSE		
		components		

Example: LFS Accuracy Ratings for 2011

	Quality criteria									
	Average	Knowledge	Commun-	Available	Compliance to	Plans to	Risk to			
	Score	of risks	ication to	expertise	standards& best	mitigate	data			
Error sources			users		practices	risk	quality			
Specification error	66%	-	-	-	-	Ο	L			
Frame error	58%	-	-	-	-	0	L			
Non-response error	66%	-	Ο	0	-	Ο	н			
Measurement error	50%	0	Ο	0		-	Н			
Dataprocessing error	54%	0		-	-	0	Μ			
Sampling error	70%	-	-	-	-	-	М			
Model error	46%	Ο	0		0	Ο	М			
Revision error							N/A			
Total score	58%									

The Evaluation Process

- Pre-interview activities
 - Background reading by the two evaluators
 - Self-assessments by each program area
- The Quality Interview
 - 1/2 day sessions involving 4-5 key product owners
 - Overview discussions of product processes
 - Detailed assessment of each of the 5 criteria
- Post-interview activities
 - Review of and comment on ratings by product owners
 - Ratings adjustments by evaluators to achieve equity

Example: LFS Accuracy Ratings for 2012

Error source	Score round 1	Score round 2	Knowledge of Risks	Communica tion to Users	Available Expertise	Compliance with standards & best practices	Plans towards mitigation of risks	Risk to data quality
Specification error	66	70	-	-	-	-	-	L
Frame error	58	58	-	-	-	-	0	L
Non-response error	56	52	0	0	0	0	0	Н
Measurement error	50	56	0	0	0	0	-	Н
Data processing error	54	62	0	0	-	-	-	М
Sampling error	70	78	-	0	-	-	0	М
Model/estimation error	50	60	0	0	0	-	-	М
Revision error	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
Total score	56,4	60,9						

Example: LFS Accuracy Ratings for 2012

	Score round 1	Score round 2	Knowledge of Risks	Communica tion to Users	Available Expertise	Compliance with standards & best	Plans towards mitigation of risks	Risk to data quality
Error source			/	Blue inc	licates	practices		
Specification error	66	70	-	deterior	ation	-	-	L
Frame error	58	58	-	in qualit	:y	Pink indicates		L
Non-response error	56	52	0	0	Ø	improve	ment	н
Measurement error	50	56	0	0	0	in quality	y	н
Data processing error	54	62	0	0	_	-		М
Sampling error	70	78	-	0	-	-	0	М
Model/estimation error	50	60	0	0	0	-	-	М
Revision error	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
Total score	56,4	60,9						

LFS Change Ratings between Round 1 and Round 2

	Score round 1	Score round 2	Knowledge of Risks	Communica tion to Users	Available Expertise	Compliance with standards & best	Plans towards mitigation of risks	Risk to data quality	Correction from 2011 rating Improvement from 2011 rating Deterioration from 2011 rating
						practices			Comments on changes
Error source									
Specification error	66	70	7	7	7	7	5 → 7 ¹	L	¹ Planning cognitive lab work to reduce specification error. Reint being planned will also help in this regard.
Frame error	58	58	7	7	7	3	5	L	
Non-response error	56	52	7→6 ¹	5	9 5 ²	7 6→5 ³	5→5⁴	н	 ¹Knowledge of the causes of nonresponse have deteriorated. Altheories, the true causes of the increases in both intrinsic and results be sorted out. ²Corrected due to level of expertise in data-collection ³This is both a correction to the Round 1 ratings and a deteriorate for telephone panels is to use face to face interviewing for Wave reasons but foremost is to reduce nonresponse bias. There are o practices as well. ⁴Despite the considerable planning effort, this rating stayed at "G mitigation activities have been slow to materialize while residuate a "critical" or "crisis level. This actually represents somewhat thus we note it even though there was no change in the rating.
Measurement error	50	56	5	5	5	3 <i>→</i> 5 ¹	7 <i>→</i> 8 ²	н	¹ Monitoring of TIs has commenced and further cognitive testing i questionnaire. However, to achieve compliance with best practice examination of measurement error is needed; for example, to bett causes and effects of rotation group bias, and removal of the fact a large extent aware of which calls are being monitored. ² Plans are in place to conduct reinterview survey; however, more measurement errors in the labour force estimates.
Data processing error	54	62	5	3→5 ¹	7	7	5→7 ²	М	¹ QD documents data editing and provide information on coding e planned in conjunction with ISO standards work. ² Plans to review the automated coding quality are in place.
Sampling error	70	78	7	$7 \rightarrow 9^1$	7	7 <i>→</i> 9 ²	7	М	¹ QD documents sample design and sampling error. ² Work on sampling error is well regarded and is consistent with
Model/estimation error	50	60	5	5	3 5 →6 ¹	3→7 ²	5→7 ³	М	¹ Error corrected in last year's evaluation of seasonal adjustment ² Work on time series adjustment regarded as state of the art. Also estimation is very good. ³ Plans in place to revise estimation approach have been approve implementation is underway.

Summary of Results from Round 1 – Dec. 2011

Error Source	RS	СРІ	FTG	LFS	NA	SBS	BR	TPR	Avg
Specification	74	68	62	66	56	46	62	44	60
Frame	36	42	62	58		62			
Overcov.							48	52	10
Undercov.							40	34	49
Duplication				· · · · · · · · · · · · · · · · · · ·			46	64	
NR/Miss. data	62	36	62	66	64	74	40	60	57
Meas/Content	52	40	54	50	58	50	42	50	50
Data proc.	46	70	46	54	44	52			52
Sampling		54		72	44	80			64
Model/est'n	54	64	66	46	44	60			56
Revision	74		62		62	58			64
Total	57	56	59	58	51	59	45	52	55

¹⁸Red Bold = High Risk, Black Bold = Medium Risk, No Bold = Low Risk

Summary of Results from Round 2 – Dec. 2012

Error Source	RS	СРІ	FTG	LFS	SBS	LCS	BR	TPR	Rating
Specification error	N/A	68	58	70	54	34	66	46	57
Frame error	60	62	58	58	64	42	55	62	58
Overcoverage							56	56	
Undercoverage							46	60	
Duplication							63	70	
Nonresponse error/Missing	52	55	66	52	70	40	48	66	56
Measurement error/Content	58	62	62	56	52	46	46	58	55
Data processing error	48	76	60	62	60	42	N/A	N/A	58
Sampling error	N/A	66	N/A	78	84	54	N/A	N/A	71
Model/estimation error	38	52	80	60	60	38	N/A	N/A	55
Revision error	58	N/A	76	N/A	56	N/A	N/A	N/A	63
Round 2 Mean Rating	49,6	63,9	<mark>65,8</mark>	60,9	61,4	42,1	52,2	58,0	57
Round 1 Mean Rating	46,7	60,3	57,3	56,4	59,6	N/A	47,2	52,2	54
Improvement	2.0	2.0	0.5	4.5	1.0	NI/A		ГО	2.5
Improvement	2,9	3,6	8,5	4,5	1,8	N/A	5,0	5,8	Z,5

RED BOLD = HIGH RISK BLACK BOLD = MEDIUM RISK REGULAR FONT = LOW RISK N/A = NOT APPLICABLE

Results for National Accounts from Round 2 (December 2012)

Error Source	GDP Quarterly	GDP Annual
Input data source (Average)	53	66
Structural Business Survey (SBS)	N/A	66
Index of Service Production (ISP)	58	N/A
Index of Industrial Production (IIP)	58	N/A
Merchanting Service of global enterprises	42	n.e.
Compilation error (modelling)	48	48
Compilation error (data processing)	40	35
Deflation error (including specification error)	48	48
Balancing error	56	50
Revision error	56	54
Round 2 Mean Rating	50,5	49,9

Future Work

- Work is continuing to address areas of high risk that have fair to poor ratings
- Criteria, checklists, etc for all quality dimensions will be revised and enhanced
- Criteria for National Accounts will be further revised and detailed quality criteria checklists will be developed
- Round 3 of the Accuracy assessment will begin in November, 2014
- User dimensions to be evaluated by internal quality units rather than external consultants