

LAS LIMITACIONES EXPERIMENTALES COMO FUENTE DE INFORMACION PARA REPLICACIONES DE EXPERIMENTOS EN EDUCACION DE REQUISITOS SOFTWARE

Dante Carrizo¹, Oscar Dieste², Marta López³

¹Dpto. de Ing. Informática y Cs. de la Computación, U. de Atacama, Copiapó, Chile

²Facultad de Informática, U. Politécnica de Madrid, Madrid, España

³Xunta de Galicia, España

dante.carrizo@uda.cl

RESUMEN

La entrevista es la técnica más extensamente usada en ingeniería de requisitos (IR). A pesar de su importancia, la investigación en entrevistas es bastante limitada, en particular de una perspectiva experimental. Hemos realizado una serie de experimentos que exploran la eficacia relativa de entrevistas estructuradas y no estructuradas. Esta línea de investigación ha sido activa en sistemas información en los años pasados, de modo que nuestros experimentos puedan ser agregados a ellos para obtener guías prácticas. La agregación experimental es una tarea exigente. No sólo requiere un número grande de experimentos, sino también considerar la influencia de las variables moderadoras existentes. Sin embargo, en el estado actual de la práctica en IR, aquellas moderadoras son desconocidas. Creemos que el análisis de las amenazas a la validez en experimentos sobre entrevistas puede dar luz sobre como mejorar nuevas replications y las agregaciones correspondientes. Es probable que esta estrategia pueda ser aplicada también en otras áreas de la Ingeniería de Software.

Palabras claves: Requisitos, educación de requisitos, experimentación, entrevistas

ABSTRACT

Interviews are the most widely used elicitation technique in Requirements Engineering (RE). Despite its importance, research in interviews is quite limited, in particular from an experimental perspective. We have performed a series of experiments exploring the relative effectiveness of structured and unstructured interviews. This line of research has been active in Information Systems in the past years, so that our experiments can be aggregated together with existing ones to obtain guidelines for practice. Experimental aggregation is a demanding task. It requires not only a large number of experiments, but also considering the influence of the existing moderators. However, in the current state of the practice in RE, those moderators are unknown. We believe that analyzing the threats to validity in interviewing experiments may give insight about how to improve further replications and the corresponding aggregations. It is likely that this strategy may be applied in other Software Engineering areas as well.

Keywords: Requirements, requirements elicitation, experimentation, Interviews

1. INTRODUCCIÓN

La IR, como una disciplina dentro del desarrollo de sistemas y software, ha sido ampliamente reconocida como crucial [1]. Requisitos inadecuados, incompletos, o ambiguos tienen un impacto crítico sobre la calidad del software y la cantidad de correcciones para desarrollar el producto final [2]. Varios autores destacan la necesidad de centrarse en el proceso de educación de requisitos para conseguir las especificaciones correctas [1] [2] [3]. Las entrevistas juegan un rol principal en la educación o captura de requisitos ya que ellas son las técnicas más extensamente aplicadas [4]. A pesar de su importancia, poca investigación ha sido realizada sobre como evaluar y mejorar la eficacia de las entrevistas [4]. Investigación empírica es particularmente escasa, existiendo solo unos pocos experimentos que consideran las entrevistas en IR. Esto contrasta con otras disciplinas, como Psicología o Finanzas, donde considerable investigación ha sido ejecutado para analizar empíricamente la eficacia, exactitud, influencia de roles, etc. de las entrevistas.

Nuestro trabajo sobre experimentación aplicada a la educación comienza con una revisión sistemática sobre técnicas de educación [5]. Nosotros detectamos que había algunos estudios experimentales comparativos que exploran la efectividad relativa de las entrevistas estructuradas y no estructuradas [6] [7] [8] [9]. En general, las entrevistas estructuradas se comportan mejor que las no estructuradas [5]. Sin embargo, creemos que se deben definir guías prácticas recomendadas por ingeniería de software basada en evidencia [10], y esto implica el uso de métodos rigurosos de agregación como meta-análisis.

El meta-análisis es un exigente método en términos del número de experimentos requeridos, y por lo tanto la base experimental es insuficiente [6] [7] [8] [9]. Comenzamos a diseñar y ejecutar repeticiones de aquellos experimentos para expandir el conjunto de datos existente, adaptando los tipos particulares de entrevistas, variables respuesta y tareas experimentales a nuestro contexto (experimentos de laboratorio con estudiantes

de Informática). Ejecutamos un estudio piloto en 2006, y posteriormente experimentos anuales (excepto 2009 por razones logísticas): 2007 (analizado y publicado [11]), 2008 (a ser entregado), y 2010 (a ser analizado).

Diseñar repeticiones es una tarea desafiante, debido a que los experimentos originales no son normalmente descritos con los detalles necesarios para ser repetidos. Además, nosotros exploramos las variables moderadoras potenciales para ganar un profundo entendimiento acerca del fenómeno de interés (entrevistas en nuestro caso), pero estas variables son bastantes desconocidas. Según nuestra opinión, el único que se orienta al problema correcto es [16]. Sin embargo, requiere que las repeticiones sean similares para analizar la influencia posible de las variables moderadoras. Esto no es posible usando la base experimental existente.

Una fuente alternativa variables moderadoras potenciales puede provenir de las limitaciones y amenazas a la validez de cada experimento. Pero ya que las amenazas a la validez son restricciones metodológicas, ellas no pueden ser usadas para identificar variables moderadoras. Así, nosotros estamos interesados en las limitaciones. Cabe notar que a veces los autores identifican como limitaciones:

- 1) Aspectos que amenazan la validez externa de los experimentos (por ejemplo, muestras de conveniencias),
- 2) Características de una técnica aplicada (por ejemplo, variabilidad sobre la aplicación de una técnica de entrevista estructurada), muestra o proceso, y
- 3) Posibles razones que explican los resultados experimentales (por ejemplo, la influencia de los entrevistados).

De un conjunto inicial de limitaciones-amenazas, solamente este último subconjunto de aspectos apuntan a la existencia de variables moderadoras, y esa información puede ser usada para diseñar nuevas repeticiones. En este artículo, presentamos un ejemplo exploratorio del uso de limitaciones para encontrar variables

moderadoras potenciales, usando los experimentos sobre entrevistas como ejemplos. Las limitaciones necesitan, al menos, de una breve descripción de los experimentos. Sin embargo, por razones de espacio, esta información no es incluida aquí, aunque están disponibles en <http://www.grise.upm.es/sites/extras/1/>. La estructura de este artículo es como sigue: la sección 2 presenta el análisis de las limitaciones-amenazas experimentales; la sección 3 describe las moderadoras identificadas; finalmente, las conclusiones son presentadas en sección 4.

2. ANÁLISIS DE LIMITACIONES Y AMENAZAS EXPERIMENTALES

La Tabla 1 muestra las limitaciones-amenazas identificadas por Agarwal et al. [6], Browne et al. [7], Carrizo et al. [11], Marakas et al. [8], and Pitts et al. [9] (usamos solo el primer autor en adelante). En promedio, cada experimento identifica 5-6 limitaciones-amenazas, excepto Browne que identifica solo 2. Las limitaciones-amenazas han sido clasificadas de acuerdo al aspecto que hacen referencia. Por ejemplo, Agarwal declara que "los expertos constituyen una muestra de conveniencia, más que una aleatoria, incluso aunque su asignación a los grupos fue realizado aleatoriamente". Esta limitación-amenaza es referida a la muestra utilizada en el experimento.

Como muestra la Tabla 1, todas las limitaciones-amenazas han sido agrupadas de acuerdo a tres principales categorías finalmente obtenidas: Proceso seguido (cómo fue llevado a cabo el experimento), Muestra (las características de los entrevistadores y los entrevistados), y Técnicas (cómo las entrevistas fueron aplicadas). Estas categorías son un tanto generales y probablemente puedan ser utilizadas para clasificar limitaciones-amenazas en otras áreas (por ejemplo, pruebas de software).

Analizando las limitaciones-amenazas por categoría (filas), podemos obtener coincidencias, frecuencias, etc. Por ejemplo, puede notarse que la principal limitación-amenaza identificada para la categoría Proceso es la que indica la experimentación

en laboratorio, con todo lo que esto implica, como describe Marakas. En todos los otros casos, el foco del experimento se centra sobre aspectos como número de sesiones [6] o la complejidad de las tareas experimentales [7].

En los cinco experimentos la categoría más frecuentemente citada como limitación-amenaza es la Muestra, la cual incluye los entrevistadores, entrevistados, codificadores o cualquier otro rol necesario para ejecutar la experimentación. Algunos autores tratan específicos problemas de muestra, como su carácter conveniente [6] o su motivación [8]. Sin embargo, todos los experimentos convienen en reconocer la experiencia y juego de rol como limitaciones-amenazas. Cuatro de los trabajos señalan este aspecto centrándose en el entrevistador y sólo Carrizo en el entrevistado.

El criterio menos relevante. De acuerdo a la Tabla 1, es la técnica. Centrándose en las técnicas de educación aplicadas en las entrevistas, Marakas y Carrizo presentan perspectivas opuestas, si ellas aceptan o no la variabilidad en la aplicación de técnicas estructuradas. Esta diferencia solo muestra dos posibles enfoques de la experimentación, dependiendo sobre el tipo de control de estas técnicas de educación pero sin invalidar el experimento. Son sólo características de aquellos diseños experimentales. Otro tipo de técnicas son aquellas utilizadas para representar los datos educidos, como los DFD de Marakas, o aquellas aplicadas para codificar datos extraídos de la técnica educidos, de Browne.

El próximo paso fue analizar cada limitación-amenaza en Tabla 1 de acuerdo a la siguiente clasificación:

- 1) Amenazas reales a la validez externa de los experimentos; esto es, aspectos que evitan que los resultados experimentales sean extrapolados a poblaciones generales. Ya que son restricciones metodológicas, ellas no pueden ser usadas para identificar variables moderadoras.
- 2) Características, o cualquier aspecto que de hecho es una faceta del proceso,

muestra o técnica. Ellas no pueden ser usadas para identificar variables moderadoras tampoco.

3) Indicio o posibles razones hipotetizadas para explicar resultados experimentales

De acuerdo al diccionario Merriam-Webster, un indicio (*inkling*) es definido como “una idea vaga o noción; entendimiento leve”. Por esto, usamos el término para denotar una noción sugerida la cual ni es probada o juzgada ya que no era analizada en el experimento. Por ejemplo, en Salud podemos encontrar estudios que demuestran ciertas

conexiones lógicas no pretendidas en la literatura científica; estas conexiones potencialmente revelan nuevo conocimiento u ocultan hipótesis [13]. Por ejemplo, la conexión oculta entre la deficiencia de magnesio y las migrañas en revistas médicas la que es sólo detectada a través de minería de textos. Esto es un indicio; ni una noción experimentada, ni una certidumbre explícita, sino una noción conocida o sugerida que es mencionada en diferentes resultados experimentales. En esta línea, nosotros estamos buscando indicios en nuestra base experimental, de hecho, la fuente de potenciales variables moderadoras.

Tabla 1. Número y descripción de todas las limitaciones (o limitaciones-amenazas) por experimento

	Agarwalet al. [6]	Browne et al. [7]	Carrizo et al. [11]	Marakas et al. [8]	Pitts et al. [9]	Total							
Proceso	Solo un dominio de problema	-	1	Problemas no complejos	2	Laboratorio	1	Solamente una estrategia de determinación	7				
	Sin proyectos "oro"									Naturaleza exploratoria			
	Solamente una sesión												
Muestra	Muestra conveniente	1	Experiencia previa de entrevistadores	2	Entrevistado vs. problema	3	Tamaño muestra	3	Experiencia entrevistadores vs. Dominio del problema	11			
	Juego de Rol										Sin preferencias predefinidas sobre técnicas de los entrevistados	Compromiso de la muestra	Sólo un entrevistado
											Motivación de la muestra		
Técnicas	-	1	Precisión del esquema de codificación	1	Variabilidad de la técnica aplicada	2	Mínimo uso de 70% del tiempo	2	Sólo una medida de reglas de parada cognitiva	6			
											Sólo un método de representación	Esquema de codificación centrado sobre taxonomía predefinida de requisitos	
Total	5	2	4	7	6	24							

Por ejemplo, Marakas identifica como limitaciones-amenazas que su estudio es un experimento en laboratorio. Obviamente, los experimentos en laboratorios son limitados en el conocimiento que pueden obtener por muchas razones: entorno fuertemente controlado, configuraciones idealizadas, etc., pero esto no implica la existencia de algún moderador. Igualmente, Agarwal identifica como una limitación-amenaza el hecho que los sujetos utilizados en su experimento fueron una muestra de conveniencia. Similarmente, esto es una restricción metodológica. Una muestra sin irrestricta tendría que ser usada en una situación ideal. Sin embargo, esto no apunta a la existencia de alguna variable moderadora.

La limitación-amenaza relacionada con el número de sesiones, número de problemas, complejidad de los problemas, y número de técnicas aplicadas están relacionadas con el costo, esfuerzo, y disponibilidad de los individuos involucrados. Estas son claras restricciones que afectan la generalización de los resultados experimentales, pero ellas no implican la existencia de algún indicio que afecte la efectividad de la entrevista (esto es, un moderador). En esta misma línea, la limitación-amenaza técnica de Carrizo, Marakas y Pitts puede ser visualizada en la misma dirección, ya que ellas describen particularidades de aquellos experimentos que restringen la generalización de resultados como en el caso previo.

La Tabla 2 presenta los indicios provenientes del análisis de la Tabla 1. Estos pueden ser denotados como indicios debido a que apuntan a la falta de validez de los resultados experimentales en cada contexto experimental (esto es, validez interna). Sin embargo, ellos no son errores de diseño. Por ejemplo, la experiencia de los entrevistados puede ser uno de los indicios pues Browne señala que la experiencia puede afectar la efectividad de los entrevistadores. Por lo tanto, si los experimentos no toman medidas para controlar la experiencia de los sujetos, los resultados pueden ser inválidos. En otras palabras: la experiencia del entrevistador es un potencial moderador.

Otros ejemplos son la "precisión del esquema de codificación" de Browne, o el "esquema de

codificación centrado en taxonomía predefinida de requisitos" de Pitts. Estos pueden ser considerados como indicios en el sentido que nosotros lo proponemos ya que ellos pueden ser una fuente de un sesgo de medición, lo cual puede afectar al análisis de la hipótesis. También, los indicios agrupados bajo el criterio Muestra pueden ser considerados como fuente potencial de sesgo, excepto el aspecto del tamaño de la muestra, el cual no es considerado un indicio sino un factor que influye sobre el poder estadístico. Un tamaño de muestra más grande solo incrementa la confianza de un estimador. Todos estos indicios seleccionados son considerados en el análisis siguiente.

3. IDENTIFICACIÓN DE MODERADORES

En la literatura experimental de otros campos científicos (Finanzas, Salud, etc.) los indicios que identificamos en la Tabla 2 son usualmente relacionados a tipos específicos de sesgo. Sesgo, en este contexto, tiene el típico significado de Error Sistemático, señalando influencias no deseadas de origen diverso que necesita ser removido o minimizado para incrementar la precisión.

Aparentemente los autores de experimentos sobre entrevistas tuvieron un similar punto de vista acerca de los indicios de la Tabla 2 y esto es el por qué ellos los listaron bajo la sección *Amenazas a la Validez* de sus trabajos. En algunos casos, proceder de esta forma es completamente justificado, ya que el indicio es claramente un sesgo. Por ejemplo, los indicios listados bajo la categoría Técnicas en la Tabla 2 son instancias de sesgo de medición (riesgo para exacta determinación de los valores de las variables respuesta). Otro ejemplo, es la motivación de la muestra, la cual es un ejemplo de sesgo de motivación. Motivación es un prerrequisito para ejecutar adecuadamente una tarea sin consideración de su campo y no parece ser un objeto de investigación legítimo. En ambos casos, no encajan con nuestro propósito.

Sin embargo, en muchos otros casos esto no es verdad, particularmente en la categoría Muestra, la cual es nuevamente una de las más pobladas en la Tabla 2. Lo que puede

ser un sesgo para algunas disciplinas (como Finanzas), puede ser un legítimo objeto de investigación en IR. Es el caso de la experiencia de los sujetos, por citar un claro ejemplo. No es sorprendente que la mayoría de los indicios de este tipo provienen de la categoría Muestra. En IR, particularmente en entrevistas, estamos interesados en los *stakeholders*, sus particularidades y las relaciones que establecen con el problema bajo estudio. Por lo tanto, aquellos aspectos no son sesgos o riesgos, sino aspectos que tienen que ser considerados para entender apropiadamente cuándo y cómo las entrevistas trabajan. Esta es la razón por la que nos dimos cuenta que este tipo de indicios estaban realmente mostrándonos potenciales variables moderadoras.

La lista de abajo muestra una clasificación de los indicios de la Tabla 2, bajo la perspectiva de potencial sesgo al que pueden dar origen. No es una lista exhaustiva (nosotros incluimos sólo los ejemplos más claros) pero es útil para un análisis rápido:

- Sesgo de artefacto, relacionado los indicios 'Entrevistado vs. problema' y 'Entrevistador vs. Dominio del problema.'
- Sesgo de entrevistador, o cualquier error sistemático debido a la obtención consciente o inconsciente de datos por parte del entrevistador. Se relacionan con estos indicios Juego de rol y Experiencia de los entrevistadores.
- Sesgo de los entrevistados, relativos a juego de rol y, de acuerdo a la Tabla 1, con la posible identificación entrevistador-problema y potencial sesgo derivado de la preferencia de uso de una técnica sobre otra.

Como se explicó antes, los tres ítems de arriba no pueden ser considerados sesgos desde la perspectiva de entrevistas en IR debido a que son aspectos que necesitamos conocer para explicar las razones de la efectividad de las entrevistas y lo que necesitamos tomar en cuenta para realizar la educación en la práctica. Por lo tanto, de estos indicios (o potenciales sesgos) podemos identificar los moderadores siguientes:

Tabla 2. Número y descripción de indicios por experimento

	Agarwalet al. [6]	Browne et al. [7]	Carrizo et al. [11]	Marakas et al. [8]	Pitts et al. [9]	Total
Proceso	-	-	-	-	-	-
Muestra	1	1	2	2	3	9
	Juego de rol	Experiencia previa de los entrevistadores	Entrevistado vs. problema Sin preferencias predefinidas sobre técnicas de los entrevistados	Motivación de la muestra Compromiso de la muestra	Experiencia entrevistadores vs. Dominio del problema Experiencia cuantificada de los entrevistadores Sólo un entrevistado	
Técnicas	-	1	-	-	1	2
		Precisión del esquema de codificación			Esquema de codificación centrado sobre taxonomía predefinida de requisitos	
Total	1	2	2	2	4	11

- Problema
- Experiencia
- Características Personales (psicológicas)

Y de estos moderadores, las recomendaciones naturales son las siguientes:

- Realizar entrevistas acerca de diferentes tipos de problemas, de diferentes tamaños y complejidad, de preferencia, de diferentes dominios.
- Análisis de la experiencia de los sujetos. Mientras más detalle se obtenga, mejor control del experimento y más alta calidad de los datos obtenidos. Por ejemplo, y donde sea posible y apropiado, aparte de los años, consultar por el número y tamaño de los proyectos.
- Análisis del rol asignado a cada sujeto, basado en su experiencia, aptitudes, conocimiento del campo, etc. Quizás introducir tests psicológicos medidas relacionadas para estudiar la personalidad de los sujetos pueden ser útiles. No olvidar que hay más roles que entrevistados y entrevistadores y que ellos también pueden tener influencia sobre el resultado final de la experimentación.

Es, por supuesto, difícil aplicar todas estas recomendaciones en la práctica debido a las características específicas del experimento, métricas, falta de un apropiado grupo de sujetos, etc. Sin embargo, estos moderadores pueden tener una influencia y deben ser considerados en experimentos de entrevistas.

Para diferentes áreas de la ingeniería de software los moderadores seguramente diferirán (quizás no en el caso de la experiencia de los sujetos). Sin embargo, pensamos que un similar procedimiento aplicado a experimentos de entrevistas puede funcionar.

4. CONCLUSIONES

Cuando se reportan los experimentos, los investigadores tienden a mezclar amenazas a la validez y otras limitaciones. Sin embargo,

es interesante diferenciarlos debido a restricciones metodológicas, como el tipo de muestra o el número de sujetos, afectan la validez externa. Otras limitaciones o indicios, tales como la experticia de los entrevistadores (en el caso de experimentos de entrevista), no son amenazas a la validez. Aquellos indicios son, en realidad, partes del marco teórico de la correspondiente área científica. Por ejemplo, la experticia del entrevistador puede tener una influencia sobre la efectividad de las entrevistas, como el sentido común sugiere. Sin embargo, tratar con teorías es un asunto complicado en estos días en Ingeniería de Software Empírica. Es más fácil pensar de ellas como variables moderadoras posibles que puede influenciar los resultados de los experimentos.

Nosotros creemos que el análisis de los indicios identificados en experimentos puede ser una útil estrategia para encontrar moderadores. Los moderadores así encontrados pueden ser incluidos en el diseño de nuevas repeticiones. En este artículo, hemos aplicado estas ideas a un conjunto existente de experimentos acerca de técnicas de entrevistas. Nosotros somos conscientes que nuestra propuesta no tiene una formulación rigurosa y sistemática. Pensamos mejorarlo en el futuro.

5. REFERENCIAS

- [1] Abran, A., Moore, J.W., Bourque, P., Dupuis, R., and Tripp, L.L. Guide to the Software Engineering Body of Knowledge (SWEBOK). IEEE, 2004.
- [2] Bell, T.E. and Thayer, T.A., Software Requirements: Are they really a problem?. 2nd International Conference on Software Engineering (ICSE'76) (San Francisco, CA, 1976), 61-68.
- [3] Leuser, J., Porta, N., Bolz, A., and Raschke, A. Empirical Validation of a Requirements Engineering Process Guide. 13th Int .Conf. on Evaluation and Assessment in Software Engineering (EASE'09) (UK, April 20-21, 2009), 1-10.
- [4] Hickey, A., Davis, A., and Kaiser, D. Requirements Elicitation Techniques:

Analyzing the Gap Between Technology Availability and Technology Use. *Comparative Technology Transfer and Society* 1, 3 (Dec. 2003), 279-302.

[5] Dieste, O., and Juristo, N. Systematic Review and Aggregation of Empirical Studies on Elicitation Techniques. *IEEE Transactions on Software Engineering*, 37, 2, (March/April 2011), 283-304.

[6] Agarwal, R., and Tanniru, M. R. Knowledge Acquisition Using Structured Interviewing: An Empirical Investigation. *Journal of Mng. Inf. Systems*, 7, 1 (Summer 1990), 123-140.

[7] Browne, G. J., and Rogich, M. B.: An Empirical Investigation of User Requirements Elicitation: Comparing the Effectiveness of Prompting Techniques. *Journal of Manag. Inf. Systems*, 17, 4 (Spring 2001), 223-250.

[8] Marakas, G.M. and Elam, J.J. Semantic Structuring in Analyst Acquisition and Representation of Facts in Requirements Analysis. *Inform. Systems Research*, 9, 1 (March 1998), 37-63.

[9] Pitts, M.G., and Browne, G.J. Stopping Behavior of Systems Analysts During Information Requirements Elicitation. *Journal of Mng. Inf. Systems*, 21,1 (Summer 2004), 203-226.

[10] Kitchenham, B., Dybå, T., and Jørgensen, M. Evidence-based Software Engineering. 26th Int. Conference on Software Engineering (ICSE'04) (Edinburgh, UK, May 23-28, 2004). IEEE Computer Society, Washington DC, USA, 2004, 273-281.

[11] Carrizo, D., Dieste, O., Juristo, N., and Lopez, M. Estudio Experimental de la Efectividad de la Entrevista Abierta frente a la Entrevista Independiente de Contexto. 14th Workshop on Requirements Eng. (WER'11) (Brazil, April 27-29, 2011), 41.

[12] Juristo, N., and Vegas, S. Using Differences among Replications of Software Engineering Experiments to Gain Knowledge. 3rd Int. Symposium on Empirical Software Engineering and Measurement (ESEM'09) (Lake Buena Vista, Florida, USA, October 15-16, 2009). IEEE, 356-366.

[13] Swanson, D.R. Two Medical Literatures that are Logically but not Bibliographically Connected. *American Society for Information Science*, 38, 4 (July 1987), 228-233.